

Part Number 83-0902024A001 Rev A
Version Date 4 December 1996

DTOS FORMAL TOP-LEVEL SPECIFICATION (FTLS)

CONTRACT NO. MDA904-93-C-4209
CDRL SEQUENCE NO. A005

Prepared for:
Maryland Procurement Office

Prepared by:



Secure Computing Corporation
2675 Long Lake Road
Roseville, Minnesota 55113

Authenticated by _____ Approved by _____
(Contracting Agency) (Contractor)

Date _____ Date _____

Distribution limited to U.S. Government Agencies Only. This document contains NSA information (4 December 1996). Request for the document must be referred to the Director, NSA.

Not releasable to the Defense Technical Information Center per DOD Instruction 3200.12.

© Copyright, 1994, Secure Computing Corporation. All Rights Reserved. This material may be reproduced by or for the U.S. Government pursuant to the copyright license under the clause at DFARS 252.227-7013 (OCT.88).



Formal Top Level Specification

**DTOS FORMAL TOP-LEVEL SPECIFICATION
(FTLS)**

Secure Computing Corporation

Abstract

This report formally describes a portion of the DTOS kernel.

Part Number 83-0902024A001 Rev A
Created 2 December 1994
Revised 4 December 1996
Done for Maryland Procurement Office
Distribution Secure Computing and U.S. Government
CM /home/cmt/rev/dtos/docs/ftls/RCS/ftls.vdd,v 1.21 4 December 1996

This document was produced using the T_EX document formatting system and the L^AT_EX style macros.

LOCKserverTM, LOCKstationTM, NETCourierTM, Security That Strikes BackTM, SidewinderTM, and Type EnforcementTM are trademarks of Secure Computing Corporation.

LOCK[®], LOCKguard[®], LOCKix[®], LOCKout[®], and the padlock logo are registered trademarks of Secure Computing Corporation.

All other trademarks, trade names, service marks, service names, product names and images mentioned and/or used herein belong to their respective owners.

© Copyright, 1994, Secure Computing Corporation. All Rights Reserved. This material may be reproduced by or for the U.S. Government pursuant to the copyright license under the clause at DFARS 252.227-7013 (OCT.88).

Contents

1	Scope	1
1.1	Identification	1
1.2	System Overview	1
1.3	Document Overview	1
2	Applicable Documents	3
3	FTLS Overview	4
3.1	Internal Consistency Within the FTLS	4
3.2	Comments on Request Specifications	5
3.3	Typographic Conventions	6
4	Basic Kernel State Definition	8
4.1	Primitive Entities	8
4.2	Process Management	10
4.3	Port Name Space	19
4.4	Ports	24
4.5	Notifications	26
4.6	Special Ports	27
4.7	Total Send Rights	33
4.8	Registered Rights	34
4.9	Memory System	35
4.10	Messages	42
4.11	Processors and Processor Sets	53
4.12	Time	55
4.13	Devices	55
4.14	Summary	57
5	DTOS State Extensions	59
5.1	Subject Security Information	59
5.2	Object Security Information	60
5.3	Security Identifiers for Access Computations	62
5.4	Permissions	64
5.5	Access Vector Cache	71
5.6	Message Security Information	73
5.7	Task Creation Information	74
5.8	Server Ports	76
5.9	Memory Region Protections	76
5.10	Summary of DTOS Kernel State	77
6	Kernel Execution Model	78
6.1	Execution Summary	78
6.2	Utility Transitions	80
6.3	Trap Invocation	83
6.4	Initial mach_msg processing	84
6.5	Service Checks for IPC Based Kernel Requests	86

6.6	Request Validation	89
6.7	Definitions	91
7	System Trap Requests	97
7.1	Introduction to System Trap Requests	97
7.2	mach_thread_self	98
8	Port Requests	105
8.1	Introduction to Port Requests	105
8.2	mach_port_allocate	113
8.3	mach_port_get_receive_status	119
8.4	mach_port_get_refs	124
8.5	mach_port_get_set_status	130
8.6	mach_port_names	134
8.7	mach_port_rename	139
8.8	mach_port_request_notification	144
8.9	mach_port_set_mscount	158
8.10	mach_port_set_qlimit	162
8.11	mach_port_set_seqno	167
9	Thread Requests	172
9.1	Introduction to Thread Requests	172
9.2	thread_abort	183
9.3	thread_create and thread_create_secure	188
9.4	thread_depress_abort	196
9.5	thread_disable_pc_sampling	200
9.6	thread_enable_pc_sampling	203
9.7	thread_get_assignment	208
9.8	thread_get_sampled_pcs	210
9.9	thread_get_special_port	215
9.10	thread_get_state	222
9.11	thread_info	227
9.12	thread_max_priority	234
9.13	thread_policy	239
9.14	thread_priority	243
9.15	thread_resume and thread_resume_secure	248
9.16	thread_set_special_port	253
9.17	thread_set_state and thread_set_state_secure	260
9.18	thread_suspend	267
9.19	thread_terminate	270
10	Virtual Memory Requests	277
10.1	Introduction to Virtual Memory Requests	277
10.2	vm_allocate and vm_allocate_secure	286
10.3	vm_deallocate	296
10.4	vm_inherit	299
10.5	vm_protect	303
10.6	vm_write	308
11	Notes	314
11.1	Acronyms	314
11.2	Glossary	314

A	Bibliography	315
<hr/>		
B	Z Extensions	316
B.1	Disjointness and Partitions	316
B.2	Partial Orders	317
B.3	Sequences	318
<hr/>		
C	IPC	319
C.1	IPC Requests	319
C.2	mach_msg	320
<hr/>		
D	Refinements of the State Model	348
D.1	Additional Z Extensions	348
D.2	Refinement of IPC Name Spaces	348
D.3	Refinement of Pending Receives	354
D.4	Refinement of Virtual Memory	356
D.5	Miscellaneous Refinements	360

List of Figures

1	Utility Transitions	80
2	mach_msg Trap Invocation	84
3	Message Transmission	86
4	Request Validation	89
5	mach_thread_self Processing	103

List of Tables

1	Return Values for mach_thread_self	100
2	Return Values for mach_port_allocate	117
3	State Change Cases for mach_port_allocate	117
4	Return Values for mach_port_get_receive_status	122
5	Return Values for mach_port_get_refs	128
6	Return Values for mach_port_get_set_status	133
7	Return Values for mach_port_names	137
8	Return Values for mach_port_rename	142
9	Return Values for mach_port_request_notification	150
10	Return Values for mach_port_request_notification , port-destroyed notification	151
11	Return Values for mach_port_request_notification , no-senders notification	152
12	Return Values for mach_port_request_notification , dead-name notification	152
13	State Change Cases for mach_port_request_notification	154
14	Return Values for mach_port_set_mscount	161
15	Return Values for mach_port_set_qlimit	165
16	Return Values for mach_port_set_seqno	169
17	Return Values for thread_abort	185
18	Return Values for thread_create	191
19	Return Values for thread_create_secure	191
20	Return Values for thread_depress_abort	198
21	Return Values for thread_disable_pc_sampling	202
22	Return Values for thread_enable_pc_sampling	206
23	Return Values for thread_get_assignment	210
24	Return Values for thread_get_sampled_pcs	213
25	Return Values for thread_get_sampled_pcs	213
26	Return Values for thread_get_sampled_pcs	213
27	Return Values for thread_get_sampled_pcs	214
28	Return Values for thread_get_special_port	220
29	Return Values for thread_get_special_port	220
30	Return Values for thread_get_state	225
31	Return Values for thread_get_state	225
32	Return Values for thread_get_state	225
33	Return Values for thread_info	231
34	Return Values for thread_info	232
35	Return Values for thread_info	232
36	Return Values for thread_max_priority	237
37	Return Values for thread_policy	241
38	Return Values for thread_priority	246
39	Return Values for thread_resume	250
40	Return Values for thread_resume_secure	251
41	Return Values for thread_set_special_port	257
42	Return Values for thread_set_state	263
43	Return Values for thread_set_state_secure	264
44	Return Values for thread_suspend	268
45	Return Values for thread_terminate	271
46	Return Values for vm_allocate and vm_allocate_secure	292

47	Return Values for vm_deallocate	298
48	Return Values for vm_inherit	302
49	Return Values for vm_protect	306
50	Return Values for vm_write	311

Section **1**
Scope

1.1 Identification

This Formal Top Level Specification (FTLS) document presents a formal specification of a portion¹ of the prototype kernel developed on the Distributed Trusted Operating System (DTOS) program, contract MDA904-93-C-4209.

1.2 System Overview

The DTOS prototype is an enhanced version of the CMU Mach 3.0 kernel that provides support for a wide variety of security policies by enforcing access decisions provided to it by a *security server*. By developing different security servers, a wide range of policies can be supported by the same DTOS kernel. By developing a security server that allows all accesses, the DTOS kernel behaves essentially the same as the CMU Mach 3.0 kernel. Although this is uninteresting from a security standpoint, it demonstrates the compatibility of DTOS with Mach 3.0.

By using appropriately developed security servers, the DTOS kernel can support interesting security policies such as MLS (multi-level security) and type enforcement. The first security server planned for development is one that enforces a combination of MLS and type enforcement.

1.3 Document Overview

The report is structured as follows:

- Section 1, **Scope**, defines the scope and this overview of the document.
- Section 2, **Applicable Documents**, describes other documents that are relevant to this document.
- Section 3, **FTLS Overview**, provides an introduction to this document.
- Section 4, **Basic Kernel State Definition**, describes the data structures contained in the Mach kernel state.
- Section 5, **DTOS State Extensions**, describes extensions to the base Mach microkernel state that are needed to support the DTOS kernel.
- Section 6, **Kernel Execution Model**, describes the computational model used to model the DTOS kernel requests and the processing that is common to multiple DTOS requests.
- Section 7, **System Trap Requests**, describes a single request (**swtch**) that is invoked as a system trap.
- Section 8, **Port Requests**, describes a selection of the port manipulation requests.

¹ See Section 3 for a description the coverage.

- Section 9, **Thread Requests**, describes a selection of the thread manipulation requests.
- Section 10, **Virtual Memory Requests**, describes a selection of the virtual memory manipulation requests.
- Section 11, **Notes**, contains a partial list of acronyms and a small glossary for this document.
- Appendix A, **Bibliography**, provides the bibliographical information for the documents referenced in this document.
- Appendix B, **Z Extensions**, describes “extensions” to the Z specification language that are used in the DTOS FTLS.
- Appendix C, **IPC Requests**, describes the **mach_msg** request. This section has not yet been updated for DTOS. Currently, this section is a direct copy of the corresponding DTMach FTLS [5] section with minor changes required for DTOS sections that depend on this section and has been included only for easy reference.
- Appendix D, **Refinements of the State Model**, refines portions of the state model to a lower level of detail to model some of the data types and relationships that are used to implement the high-level abstract model described in the Basic Kernel State Definition and DTOS State Extensions chapters.

Section **2**
Applicable Documents

The following document provides a high level description of the Mach microkernel:

- OSF Mach Kernel Principles [7]

The following documents provide a detailed description of the Mach and DTOS microkernels:

- OSF Mach 3 Kernel Interface [6]
- DTOS Kernel Interface Document (KID) [8]

The DTOS security policy model is described in

- DTOS Formal Security Policy Model (FSPM) [9]

Much of this document was derived from the following document:

- Formal Top Level Specification for Distributed Trusted Mach [5]

The following documents were used as additional sources of information on Mach:

- A Mathematical Model of the Mach Kernel: Entities and Relations (Draft) [2]
- A Mathematical Description of the Mach Kernel: Virtual Memory Services (Draft) [1]
- A Mathematical Model of the Mach Kernel: Port Services (Draft) [3]
- A Mathematical Model of the Mach Kernel: Task and Thread Services (Draft) [4]

Section **3**
FTLS Overview

This document provides a *partial* formal top-level specification (FTLS) of the DTOS micro-kernel. We have made no attempt at complete coverage of the kernel interface. The FTLS includes

- A specification of the DTOS system state,
- The general properties of request execution common to most requests,
- One system trap request specification,
- 10 port request specifications,
- 21 thread request specifications, and
- Six virtual memory request specifications.

There are roughly 150 requests in DTOS, so this document covers approximately 25% of the DTOS kernel requests.

This document describes the system behavior both in English and in the Z formal specification language. Thus, readers who are unfamiliar with Z can simply ignore the formal Z specifications and read the English text.

Writing an FTLS is valuable because many behaviors of the system that have an impact on security can easily be overlooked in a less formal description. This is particularly true of behaviors which might be considered side effects of operations that have some other primary purpose. We have found this to be especially relevant for Mach since, by design, objects in the Mach microkernel have complex interactions. The process of formally specifying the behavior of a system frequently brings these behaviors to the surface.

3.1 Internal Consistency Within the FTLS

During the initial phases of the DTOS program, the primary goals of the FTLS were clarity, accuracy and completeness. As the program progressed, completeness in coverage was no longer a goal and the importance of accuracy was also diminished, leaving room for experimentation in the presentation of the FTLS. The result of this was that the FTLS is no longer internally consistent.²

There are two causes of inconsistency:

- Updates to the state model were not carried through to the request chapters. This only affects the virtual memory requests.
- The execution model was completely rewritten, and other than `mach_thread_self`, none of the request specifications were updated to conform with the new model. The specific effects of the new model on the request specifications is described in the following section.

²For further discussion of the evolution of the FTLS, the reasons for this evolution and general lessons learned while developing the FTLS, see the DTOS Lessons Learned Report[10].

3.2 Comments on Request Specifications

We have taken steps in writing this document to make the specifications easier to follow and to aid readers in locating desired information. First, requests are generally grouped based upon the type of kernel object to which the name given as the first parameter of the request is resolved (see the discussion of the client and kernel interfaces later in this section). Thus, requests whose first parameter is resolved to a thread (e.g., **thread_abort**), are in the Thread Requests chapter. Requests whose first parameter is resolved to a name space are in the Port Requests chapter, and requests whose first parameter is resolved to a memory map are in the VM Requests chapter.³ There are, however, some exceptions to this rule. For example, the first parameter of a **thread_create** request is resolved to a task, but the request specification has been placed in the Thread Requests chapter since it is so intimately linked to the rest of the material in that chapter. This guideline does not apply at all to system trap requests many of which do not even have any parameters. So, these requests are in a separate chapter.

Second, we have attempted to make the structure of each chapter and each request specification as consistent as possible. Each chapter begins with an introduction that describes processing that is common or similar for multiple requests in the chapter.

The remainder of each chapter specifies the behavior of individual requests. Each of these specifications has the following structure:

1. Client Interface - the interface visible to the thread sending the request message. This includes a C Synopsis similar to that given in the Kernel Interface Document.
 - (a) Input Parameters - the parameters included in the request message and the way in which the request is invoked.
 - (b) Output Parameters - the parameters included in the reply message generated by the request (note that some parameters may occur in both lists).

Editorial Note:

The client interface is not consistent with the new execution model because the new execution model makes no attempt to model the client interface. The old model included hooks for such modeling though they were not integrated at all with the rest of the specification.

2. Kernel Interface - the interface visible to the kernel when it calls the kernel service routine.
 - (a) Input Parameters - the input parameters included in the call.
 - (b) Output Parameters - the output parameters included in the call (note that some parameters may occur in both lists).

Editorial Note:

The kernel interface is not consistent with the new execution model. The old execution model was based upon an abstraction which considered kernel requests to be received off of a message queue rather than "directly" from a trap as in the updated execution model.

In the new execution model, the specification of the input interface is directly related to the extraction of parameters from a request (*ExtractRequest*) while specification of the output interface leads into the transition (*Return*) that describes the return from a kernel trap into user space.

³This draft of the FTLS does not explicitly represent name spaces and memory maps. It associates the names in a space and the regions in a map with the task that contains the name space and the memory map. However, the distinction between the task and the name space or memory map is still reflected in the division of the FTLS into chapters.

3. Request Criteria - the conditions that determine the behavior of the request (i.e., return values and state changes) in a given situation.
4. Return Values - a description of the value(s) returned in each possible situation as determined by the Request Criteria.
5. State Changes - a description of how the system state is changed by the request as determined by the Request Criteria.
6. Complete Request - ties together the multiple pieces of the specification including the common processing behavior specified in Section 6, **Kernel Execution Model**, and in the appropriate chapter introduction plus the return value and state change behavior described earlier in the given specification. This section is primarily of importance for the formal specification in Z. Readers who are ignoring the Z can skip the Complete Request sections entirely.

Editorial Note:

It is this section that is most directly influenced by the new execution model. With the new model, it should be possible to describe the total processing of a request more coherently in english as well as formally. See the specification of **mach_thread_self** for an example.

The following points should be made. First, some of the specifications contain a description of the parts of the system state that are invariant in the request. Other specifications do not contain such a description. In the cases where the invariants are specified, they have largely been inherited from the DTMach FTLS with only minor editing. They should not be considered to be highly reliable, particularly since the model of the system state presented in the DTOS FTLS is changed significantly from the model in the DTMach FTLS.

Second, the specification for each of the IPC-based kernel requests describes both the client and kernel interface. As stated above, the former describes the input and output parameters included in the invocation message and the reply message while the latter describes the input and output parameters included in the function call to the kernel routine. Typically, these parameters differ only in that some of the Mach names have been resolved to the object named (e.g., a thread, task, or port). We have not formally specified the relationship between these two interfaces.

3.3 Typographic Conventions

Finally, we have established typographic conventions for the identifiers used in the FTLS. In general, global objects in the specification contain capital letters while local objects are all in lower case. More specifically, the identifiers have the following forms:

- The names of Z schemas consist of capitalized words with no underscores between words (e.g., *SpecialThreadPorts*).
- Both schema components and variables consist of lower case words separated by underscores (e.g., *task_self*).
- Global constants defined through axiomatic, generic and free type definitions have the first word capitalized, the remaining words in lower case and all the words separated by underscores (e.g., *Values_partition*).

- All other identifiers (i.e., given types, free types, abbreviations and generic parameters) are printed in upper case with underscores to separate words (e.g., *PORT_CLASS*).

A request name appearing in the text is set in bold face with words separated by underscores as they would be in a call within a program (e.g., **thread_create_secure**). In the formalization, each request has an identifier which is declared as an axiomatic global constant and is therefore typeset with the first word capitalized, the remaining words in lower case and all the words separated by underscores (e.g., *Thread_create_secure_id*). One other typographic convention followed in this document is that components of the system state that are considered primitive (as opposed to being derived from some other piece or pieces of the system state) have their first character underlined (e.g., task_self_rel).

Section 4

Basic Kernel State Definition

The following describes the data structures contained in the Mach kernel state. The organization of this section is as follows:

- Section 4.1, **Primitive Entities**, describes the primitive entities in Mach. Mach is an object-based system having these primitive entities as the defined objects.
- Section 4.2, **Process Management**, describes data structures associated with process management.
- Section 4.3, **Port Name Space**, describes data structures associated with task port name spaces.
- Section 4.4, **Ports**, describes data structures associated with ports.
- Section 4.5, **Notifications**, describes data structures associated with registered notifications.
- Section 4.6, **Special Ports**, describes the various classes of ports associated with the primitive entities.
- Section 4.7, **Total Send Rights**, describes the way in which send rights are counted in the kernel.
- Section 4.8, **Registered Rights**, describes the data structures used to record the set of port rights registered for a task.
- Section 4.9, **Memory System**, describes the data structures associated with the virtual memory system.
- Section 4.10, **Messages**, describes the data structures associated with messages.
- Section 4.11, **Processors and Processor Sets**, describes the data structures associated with processors and processor sets.
- Section 4.12, **Time**, describes the data structures associated with clocks.
- Section 4.13, **Devices**, describes the data structures associated with devices.

The model of Mach presented in this section consists of both primitive and derived notions. The derived notions provide no additional information about the Mach state beyond that embodied in the primitive notions. In the following sections, derived notions are noted as being conveniences. For example, Section 4.2.1 introduces the derived notion embodied by the function *threads* to provide a more convenient representation for the primitive notion embodied by the relation *task_thread_rel*. Although any statement about *threads* can be reworded as a statement about *task_thread_rel*, it is often more desirable to write the statement in terms of *threads*. In many cases, the choice of whether to view a structure as being primitive or derived is subjective. For example, others might prefer to view *task_thread_rel* as being derived from *threads* instead of *threads* being derived from *task_thread_rel*.

As a convention, we underline the first letter in the identifier for each primitive structure in the Mach state. This is most useful when identifying which primitive structures are affected by DTOS services.

4.1 Primitive Entities

The primitive entities in Mach are:

Tasks — environments in which threads execute; a task consists of an address space, a port name space, and a set of threads

Threads — active entities comprised of an instruction pointer and a local register state

Ports — unidirectional communication channels between tasks

Messages — entities transmitted through ports

Memories — memory object representing a shared memory

Pages — logical units of memory; either a unit of physical memory or provided by a memory

Hosts — instances of the Mach kernel

Processors — devices capable of executing threads

Processor Sets — groups of processors, each belonging to a host, to which threads are assigned for scheduling

Devices — resources such as terminals and printers that can be used to transmit information between the system and its environment

Each of these primitive entities can be viewed as an abstract data type.

Mach Definition 1

$$[\textit{TASK}, \textit{THREAD}, \textit{PORT}, \textit{MESSAGE}, \textit{MEMORY}, \textit{PAGE}, \textit{HOST}, \textit{PROCESSOR}, \textit{PROCESSOR_SET}, \textit{DEVICE}]$$

At any given time, only certain primitive entities are present in the system. The sets $\underline{\textit{task_exists}}$, $\underline{\textit{thread_exists}}$, $\underline{\textit{port_exists}}$, $\underline{\textit{message_exists}}$, $\underline{\textit{memory_exists}}$, $\underline{\textit{page_exists}}$, $\underline{\textit{proc_exists}}$, $\underline{\textit{procset_exists}}$, and $\underline{\textit{device_exists}}$ denote the entities of each class that are present in the current system state.

Mach Definition 2

$$\begin{aligned} \textit{TaskExist} &\hat{=} [\underline{\textit{task_exists}} : \mathbb{P} \textit{TASK}] \\ \textit{ThreadExist} &\hat{=} [\underline{\textit{thread_exists}} : \mathbb{P} \textit{THREAD}] \\ \textit{MessageExist} &\hat{=} [\underline{\textit{message_exists}} : \mathbb{P} \textit{MESSAGE}] \\ \textit{MemoryExist} &\hat{=} [\underline{\textit{memory_exists}} : \mathbb{P} \textit{MEMORY}] \\ \textit{PageExist} &\hat{=} [\underline{\textit{page_exists}} : \mathbb{P} \textit{PAGE}] \\ \textit{ProcessorExist} &\hat{=} [\underline{\textit{proc_exists}} : \mathbb{P} \textit{PROCESSOR}] \\ \textit{ProcessorSetExist} &\hat{=} [\underline{\textit{procset_exists}} : \mathbb{P} \textit{PROCESSOR_SET}] \\ \textit{DeviceExist} &\hat{=} [\underline{\textit{device_exists}} : \mathbb{P} \textit{DEVICE}] \end{aligned}$$

$\underline{\textit{Ip_null}}$ and $\underline{\textit{Ip_dead}}$ are two special values in \textit{PORT} which are never in the set of existing ports. $\underline{\textit{port_pointer}}$ consists of $\underline{\textit{port_exists}}$ plus the special values $\underline{\textit{Ip_null}}$ and $\underline{\textit{Ip_dead}}$.

Mach Definition 3

$$\left| \begin{array}{l} \underline{\textit{Ip_null}}, \underline{\textit{Ip_dead}} : \textit{PORT} \\ \underline{\textit{Ip_null}} \neq \underline{\textit{Ip_dead}} \end{array} \right.$$

Mach Definition 4

$\begin{aligned} & \underline{PortExist} \\ & \underline{port_exists} : \mathbb{P} \text{ PORT} \\ & \underline{port_pointer} : \mathbb{P} \text{ PORT} \\ & \underline{Ip_null} \notin \underline{port_exists} \\ & \underline{Ip_dead} \notin \underline{port_exists} \\ & \underline{port_pointer} = \underline{port_exists} \cup \{ \underline{Ip_null}, \underline{Ip_dead} \} \end{aligned}$
--

Mach Definition 5

$\begin{aligned} & \underline{Exist} \\ & \underline{TaskExist} \\ & \underline{ThreadExist} \\ & \underline{PortExist} \\ & \underline{MessageExist} \\ & \underline{MemoryExist} \\ & \underline{PageExist} \\ & \underline{ProcessorExist} \\ & \underline{ProcessorSetExist} \\ & \underline{DeviceExist} \end{aligned}$
--

Note that in the model, the kernel itself is viewed as an existing task and is denoted by $\underline{k}ernel$.

Mach Definition 6

$\begin{aligned} & \underline{Kernel} \\ & \underline{k}ernel : \text{TASK} \\ & \underline{TaskExist} \\ & \underline{k}ernel \in \underline{task_exists} \end{aligned}$
--

4.2 Process Management

This section describes the data structures associated with process management. Multi-threaded processes are supported by allowing tasks to contain multiple threads.

4.2.1 Thread to Task Relationship

The relation $\underline{task_thread_rel}$ denotes the relationship between threads and tasks; a pair $(task, thread)$ is an element of $\underline{task_thread_rel}$ exactly when $thread$ is one of the threads contained in $task$. Each thread belongs to exactly one task. For convenience, the following additional notation is introduced:

- $\underline{owning_task}(thread)$ — the task to which $thread$ belongs
- $\underline{threads}(task)$ — the set of threads belonging to $task$

Mach Definition 7

$\begin{aligned} & \text{TasksAndThreads} \\ & \text{TaskExist} \\ & \text{ThreadExist} \\ & \underline{task_thread_rel} : TASK \leftrightarrow THREAD \\ & \text{owning_task} : THREAD \mapsto TASK \\ & \text{threads} : TASK \mapsto \mathbb{P} \text{THREAD} \\ \\ & \text{dom } \underline{task_thread_rel} \subseteq \underline{task_exists} \\ & \text{ran } \underline{task_thread_rel} = \underline{thread_exists} \\ & \text{owning_task} = \underline{task_thread_rel} \sim \\ & \text{threads} \\ & = (\lambda \text{task} : TASK \\ & \quad \text{task} \in \underline{task_exists} \\ & \quad \bullet \underline{task_thread_rel}(\{\text{task}\})) \end{aligned}$

4.2.2 Execution Status

The execution status of a thread identifies whether a thread is running, waiting on an event, waiting uninterruptibly, and/or halted. A thread holds some subset of these characteristics at any point in time. The type *RUN_STATES* defines the possible thread characteristics. *RUN_STATES* has possible values *Running*, *Stopped*, *Waiting*, *Uninterruptible* and *Halted*.

Mach Definition 8

$$RUN_STATES ::= Running \mid Stopped \mid Waiting \mid Uninterruptible \mid Halted$$

The values of this type have the following meanings:

- *Running* — The thread is either executing on a processor or is in a run queue waiting to execute.
- *Stopped* — The thread has been asked to stop (and might have done so). A stopped thread does not execute any instructions.
- *Waiting* — The thread is waiting for an event.
- *Uninterruptible* — The thread is waiting uninterruptibly.
- *Halted* — The thread is halted at what the kernel considers to be a “clean” point (i.e., it can be resumed properly).

The state *Uninterruptible* does not imply the state *Waiting*. A *run_state* that includes the former but not the latter can result when the procedure `clear_wait` is called on a thread that is both *Uninterruptible* and *Waiting*. The expression `run_state(thread)` indicates which of the above characteristics are held by an existing thread.

Each thread has an associated suspend count that determines whether the thread may execute user level instructions. This count is denoted by `thread_suspend_count(thread)`. A thread may execute such instructions only if the value of its suspend count is zero. It is a consequence of the operation of the system (and therefore is not stated as an axiom here) that only stopped threads have a suspend count greater than zero.

A thread may be swapped out. A thread that is swapped out has no kernel stack. The set of such threads is indicated by `swapped_threads`. Some threads may be wired into the system. A

wired thread may not be swapped out. The set $\underline{threads_wired}$ denotes the set of wired threads. Certain threads are called idle threads. An idle thread is one that runs on a processor that has no user threads to run. (That is, the thread keeps the processor “idling”.) User threads will not be marked as idle. We use $\underline{idle_threads}$ to denote the set of idle threads.

Mach Definition 9

$\underline{ThreadExecStatus}$ $\underline{ThreadExist}$ $\underline{run_state} : THREAD \mapsto \mathbb{P} RUN_STATES$ $\underline{thread_suspend_count} : THREAD \mapsto \mathbb{N}$ $\underline{swapped_threads} : \mathbb{P} THREAD$ $\underline{threads_wired} : \mathbb{P} THREAD$ $\underline{idle_threads} : \mathbb{P} THREAD$
$\text{dom } \underline{run_state} = \underline{thread_exists}$ $\text{dom } \underline{thread_suspend_count} = \underline{thread_exists}$ $\underline{swapped_threads} \subseteq \underline{thread_exists}$ $\underline{threads_wired} \subseteq \underline{thread_exists}$ $\underline{idle_threads} \subseteq \underline{thread_exists}$ $\underline{threads_wired} \cap \underline{swapped_threads} = \emptyset$

Each task also has a suspend count. The expression $\underline{task_suspend_count}(task)$ denotes the count associated with $task$. If this value is non-zero, then none of the threads in $task$ may execute regardless of their individual suspend counts.

Mach Definition 10

$\underline{TaskSuspendCount}$ $\underline{task_suspend_count} : TASK \mapsto \mathbb{N}$
--

Review Note:

We should probably specify the relationships between $\underline{task_suspend_count}$, $\underline{thread_suspend_count}$ and $\underline{run_state}$ here.

4.2.3 Priority Levels

Thread priority levels are used to determine thread execution scheduling priorities. Priority levels are represented as a subset of the integers (in particular by the numbers between 0 and 31 inclusive in current implementations). The set $\underline{Priority_levels}$ denotes the allowable priority levels. The relation $\underline{Lower_priority}$ indicates when a priority is lower than a second priority; in particular, (x, y) is an element of $\underline{Lower_priority}$ exactly when x is a lower priority than y . Since the implementation uses higher numbers to indicate lower priorities, x is lower than y when $x > y$. The relation $\underline{Higher_priority}$ is the inverse ordering indicating when a priority is higher than a second priority. The constants $\underline{Lowest_possible_priority}$ and $\underline{Highest_possible_priority}$ denote the maximum and minimum integers, respectively, in $\underline{Priority_levels}$.

Mach Definition 11

$Priority_levels : \mathbb{P} \mathbb{Z}$ $Lower_priority, Higher_priority : \mathbb{Z} \leftrightarrow \mathbb{Z}$ $Lowest_possible_priority, Highest_possible_priority : \mathbb{Z}$
$Lower_priority \subset Priority_levels \times Priority_levels$ $\forall x, y : Priority_levels \bullet (x, y) \in Lower_priority \Leftrightarrow x > y$ $Higher_priority = Lower_priority \sim$ $Lowest_possible_priority = \max Priority_levels$ $Highest_possible_priority = \min Priority_levels$

Using these relations, the minimum and maximum priorities in a set of priorities can be defined. These are denoted by $Lowest_priority(priority_set)$ and $Highest_priority(priority_set)$, respectively.

Mach Definition 12

$Lowest_priority, Highest_priority : \mathbb{P} \mathbb{Z} \mapsto \mathbb{Z}$
$\text{dom } Lowest_priority = \mathbb{P}_1 Priority_levels$ $\text{ran } Lowest_priority = Priority_levels$ $\text{dom } Highest_priority = \mathbb{P}_1 Priority_levels$ $\text{ran } Highest_priority = Priority_levels$ $\forall priority_set : \mathbb{P}_1 \mathbb{Z} \bullet Lowest_priority(priority_set) = \max priority_set$ $\forall priority_set : \mathbb{P}_1 \mathbb{Z} \bullet Highest_priority(priority_set) = \min priority_set$

There is a highest priority (equal to 12 in current implementations) normally granted to ordinary user threads. This priority is denoted by $Base_user_priority$.

Mach Definition 13

$Base_user_priority : \mathbb{Z}$
$Base_user_priority \in Priority_levels$

Three different types of priority values are associated with each thread.

- The expression $\underline{thread_priority}(thread)$ represents a base user-setable priority for $thread$.
- The expression $\underline{thread_max_priority}(thread)$ represents the maximum value to which $\underline{thread_priority}(thread)$ can be set.
- The expression $\underline{thread_sched_priority}(thread)$ represents the priority that the system uses to make scheduling decisions. This value is determined based upon $\underline{thread_priority}$ and the thread scheduling policy (discussed in Section 4.2.4), and is not directly set by the user. This value cannot exceed $\underline{thread_priority}(thread)$.

The priority level of a thread can temporarily be depressed by the request **switch_pri** or **thread_switch** to allow other threads to run. When a thread is depressed, its priority is set to the lowest possible priority.⁴ The set $\underline{depressed_threads}$ denotes those threads whose priority is currently depressed. The expression $\underline{priority_before_depression}(thread)$ denotes the priority level $thread$ had before depression if $thread$'s priority level has been depressed and $\underline{thread_priority}(thread)$ otherwise.

Mach Definition 14

⁴Note, however, that not all threads having the lowest possible priority are depressed.

$\underline{ThreadPri}$ $\underline{ThreadExist}$ $\underline{thread_priority} : THREAD \leftrightarrow \mathbb{Z}$ $\underline{thread_max_priority} : THREAD \leftrightarrow \mathbb{Z}$ $\underline{thread_sched_priority} : THREAD \leftrightarrow \mathbb{Z}$ $\underline{depressed_threads} : \mathbb{P} THREAD$ $\underline{priority_before_depression} : THREAD \leftrightarrow \mathbb{Z}$
$\text{ran } \underline{thread_priority} \subseteq \text{Priority_levels}$ $\text{ran } \underline{thread_max_priority} \subseteq \text{Priority_levels}$ $\text{ran } \underline{thread_sched_priority} \subseteq \text{Priority_levels}$ $\text{ran } \underline{priority_before_depression} \subseteq \text{Priority_levels}$ $\underline{depressed_threads} \subseteq \underline{thread_exists}$ $\text{dom } \underline{thread_priority} = \text{dom } \underline{thread_max_priority} = \text{dom } \underline{thread_sched_priority}$ $\quad = \text{dom } \underline{priority_before_depression} = \underline{thread_exists}$ $\forall \text{ thread} : THREAD \mid \text{thread} \in \text{dom } \underline{thread_priority}$ <ul style="list-style-type: none"> • $(\underline{thread_priority}(\text{thread}), \underline{thread_max_priority}(\text{thread})) \notin \text{Higher_priority}$ $\wedge (\underline{thread_sched_priority}(\text{thread}), \underline{thread_priority}(\text{thread})) \notin \text{Higher_priority}$ $\forall \text{ thread} : THREAD \mid \text{thread} \in \text{dom } \underline{thread_priority} \setminus \underline{depressed_threads}$ <ul style="list-style-type: none"> • $\underline{priority_before_depression}(\text{thread}) = \underline{thread_priority}(\text{thread})$ $\forall \text{ thread} : THREAD \mid \text{thread} \in \underline{depressed_threads}$ <ul style="list-style-type: none"> • $\underline{thread_priority}(\text{thread}) = \text{Lowest_possible_priority}$

Each existing task has an associated priority level, denoted by $\underline{task_priority}(\text{task})$, that is used to assign the initial priority for any thread created within the task.

Mach Definition 15

$\underline{TaskPriority}$ $\underline{TaskExist}$ $\underline{task_priority} : TASK \leftrightarrow \mathbb{Z}$
$\text{dom } \underline{task_priority} = \underline{task_exists}$ $\text{ran } \underline{task_priority} \subseteq \text{Priority_levels}$

4.2.4 Scheduling Policies

Each thread has an associated scheduling policy, represented by $\underline{thread_sched_policy}(\text{thread})$. The type $SCHED_POLICY$ represents the set of available scheduling policies. Examples of supported policies are Timesharing (*Timeshare*) and Fixed Priority (*Fixedpri*). Some scheduling policies have associated policy specific data that must be associated with each thread. For example, threads scheduled under the Fixed Priority policy must have an associated scheduling quantum. The type $SCHED_POLICY_DATA$ denotes policy specific scheduling data. The expression $\underline{thread_sched_policy_data}(\text{thread})$ denotes any such policy specific data associated with thread . The set $\underline{supported_sp}$ indicates which scheduling policies are actually supported by a given Mach system. All Mach systems are required to support *Timeshare* and each thread in a Mach system must be assigned one of the scheduling policies supported by the system.

Mach Definition 16

[$SCHED_POLICY, SCHED_POLICY_DATA$]

$Timeshare, Fixedpri : SCHED_POLICY$
$Timeshare \neq Fixedpri$

Mach Definition 17

$ThreadSchedPolicy$
$ThreadExist$
$\underline{thread_sched_policy} : THREAD \leftrightarrow SCHED_POLICY$
$\underline{thread_sched_policy_data} : THREAD \leftrightarrow SCHED_POLICY_DATA$
$\underline{supported_sp} : \mathbb{P} SCHED_POLICY$
$dom \underline{thread_sched_policy_data} \subseteq dom \underline{thread_sched_policy} = \underline{thread_exists}$
$Timeshare \in \underline{supported_sp}$
$ran \underline{thread_sched_policy} \subseteq \underline{supported_sp}$

4.2.5 Instruction Pointer

The set *VIRTUAL_ADDRESS* is used to denote the set of virtual addresses. These addresses are assumed to be ordered in some manner with *Vm_start* and *Vm_end* denoting, respectively, the smallest and largest addresses.

Mach Definition 18

$[VIRTUAL_ADDRESS]$
$Vm_start, Vm_end : VIRTUAL_ADDRESS$

Each thread has an associated instruction pointer indicating the address at which the thread is currently executing. The expression $\underline{instruction_pointer}(thread)$ denotes *thread's* current instruction pointer.

Mach Definition 19

$ThreadInstruction$
$\underline{instruction_pointer} : THREAD \leftrightarrow VIRTUAL_ADDRESS$

4.2.6 Emulation Environment

Mach supports binary compatibility by allowing tasks to establish user-level handlers for system calls. This is accomplished by associating an *emulation vector* with each task. Each entry in an emulation vector specifies a system call and a virtual address. Whenever the task executes a system call that has an entry in the emulation vector, the code at the specified virtual address for the system call is executed rather than the system call. The expression $\underline{emulation_vector}(task)$ denotes *task's* emulation vector.

Mach Definition 20

$EmulationVector$
$TaskExist$
$\underline{emulation_vector} : TASK \leftrightarrow \mathbb{N} \leftrightarrow VIRTUAL_ADDRESS$
$dom \underline{emulation_vector} = \underline{task_exists}$

4.2.7 Sampling

Any thread or task may be sampled. This causes the instruction pointer to be recorded in a buffer during clock interrupts or page faults if the thread or task is currently executing. The type *SAMPLE* represents the sampling information that is collected, and type *SAMPLE_TYPES* represents information that determines at which times during execution samples are collected for a given thread or task.

There are six recognized sample types. They are:

- *Sample_periodic* — each clock interrupt
- *Sample_vm_zfill_faults* — zero-filling a virtual memory page
- *Sample_vm_reactivation_faults* — reactivating a virtual memory page
- *Sample_vm_pagein_faults* — bringing a virtual memory page in
- *Sample_vm_cow_faults* — virtual memory copy-on-write faults
- *Sample_vm_faults_any* — all virtual memory page faults. This includes miscellaneous faults beyond the above mentioned four types of virtual memory faults.

These values comprise the elements of the set *Recognized_sample_types*.

Mach Definition 21

[*SAMPLE*, *SAMPLE_TYPES*]

$ \begin{aligned} & \textit{Sample_periodic}, \textit{Sample_vm_zfill_faults}, \\ & \textit{Sample_vm_reactivation_faults}, \textit{Sample_vm_pagein_faults}, \\ & \textit{Sample_vm_cow_faults}, \textit{Sample_vm_faults_any} : \textit{SAMPLE_TYPES} \\ & \textit{Recognized_sample_types} : \mathbb{P} \textit{SAMPLE_TYPES} \end{aligned} $
$ \begin{aligned} & (\textit{Sample_periodic}, \textit{Sample_vm_zfill_faults}, \\ & \textit{Sample_vm_reactivation_faults}, \textit{Sample_vm_pagein_faults}, \\ & \textit{Sample_vm_cow_faults}, \textit{Sample_vm_faults_any}) \\ & \textit{Values_partition} \textit{Recognized_sample_types} \end{aligned} $

For convenience, *SAMPLE_VM_FAULTS* is used as the combination of the sample types *Sample_vm_zfill_faults*, *Sample_vm_reactivation_faults*, *Sample_vm_pagein_faults* and *Sample_vm_cow_faults*.

There is a maximum number of samples (determined by the buffer size) that can be kept for any thread or task. This maximum is represented by *Max_samples*.

Mach Definition 22

$$\textit{SAMPLE_VM_FAULTS} == \{\textit{Sample_vm_zfill_faults}, \textit{Sample_vm_reactivation_faults}, \textit{Sample_vm_pagein_faults}, \textit{Sample_vm_cow_faults}\}$$

|
$$\textit{Max_samples} : \mathbb{N}_1$$

The set *sampled_threads* denotes the set of threads that are currently being sampled. For each sampled thread there is a set of sample types, denoted by *thread_sample_types(thread)*, indicating when a sample should be taken for the thread. Each sample taken for a thread is assigned a unique sequence number. The expression *thread_sample_sequence_number(thread)* denotes the sequence number of the most recent sample for a thread (or zero if no samples have

been collected). The expression $\underline{thread_samples}(thread)$ denotes the currently stored samples for $thread$. Each sample is stored with an associated sample number. Only the $Max_samples$ most recent samples are retained.

Mach Definition 23

<p><i>ThreadSampling</i></p> <p><i>ThreadExist</i></p> <p>$\underline{sampler_threads} : \mathbb{P} \text{ THREAD}$</p> <p>$\underline{thread_sample_types} : \text{ THREAD} \mapsto \mathbb{P} \text{ SAMPLE_TYPES}$</p> <p>$\underline{thread_sample_sequence_number} : \text{ THREAD} \mapsto \mathbb{N}$</p> <p>$\underline{thread_samples} : \text{ THREAD} \mapsto (\mathbb{N} \mapsto \text{ SAMPLE})$</p> <hr/> <p>$\underline{sampler_threads} \subset \underline{thread_exists}$</p> <p>$\text{dom } \underline{thread_sample_types} = \underline{sampler_threads}$</p> <p>$\text{dom } \underline{thread_sample_sequence_number} = \underline{sampler_threads}$</p> <p>$\text{dom } \underline{thread_samples} = \underline{sampler_threads}$</p> <p>$\forall \text{ smpls} : \mathbb{N} \mapsto \text{ SAMPLE}; \text{ thread} : \text{ THREAD};$ $\text{ num}, \text{ high} : \mathbb{N}$</p> <p>$(\text{ thread}, \text{ smpls}) \in \underline{thread_samples}$ $\wedge \text{ high} = \underline{thread_sample_sequence_number}(\text{ thread})$ $\wedge \text{ num} = \min \{ \text{ high}, \text{ Max_samples} \}$</p> <p>• $\text{ dom smpls} = \text{ high} - \text{ num} + 1 .. \text{ high}$</p>

The same sampling information is kept for tasks.

Mach Definition 24

<p><i>TaskSampling</i></p> <p><i>TaskExist</i></p> <p>$\underline{sampler_tasks} : \mathbb{P} \text{ TASK}$</p> <p>$\underline{task_sample_types} : \text{ TASK} \mapsto \mathbb{P} \text{ SAMPLE_TYPES}$</p> <p>$\underline{task_sample_sequence_number} : \text{ TASK} \mapsto \mathbb{N}$</p> <p>$\underline{task_samples} : \text{ TASK} \mapsto (\mathbb{N} \mapsto \text{ SAMPLE})$</p> <hr/> <p>$\underline{sampler_tasks} \subset \underline{task_exists}$</p> <p>$\text{dom } \underline{task_sample_types} = \underline{sampler_tasks}$</p> <p>$\text{dom } \underline{task_sample_sequence_number} = \underline{sampler_tasks}$</p> <p>$\text{dom } \underline{task_samples} = \underline{sampler_tasks}$</p> <p>$\forall \text{ smpls} : \mathbb{N} \mapsto \text{ SAMPLE}; \text{ task} : \text{ TASK};$ $\text{ num}, \text{ high} : \mathbb{N}$</p> <p>$(\text{ task}, \text{ smpls}) \in \underline{task_samples}$ $\wedge \text{ high} = \underline{task_sample_sequence_number}(\text{ task})$ $\wedge \text{ num} = \min \{ \text{ high}, \text{ Max_samples} \}$</p> <p>• $\text{ dom smpls} = \text{ high} - \text{ num} + 1 .. \text{ high}$</p>

4.2.8 Thread Time Statistics

The system records time statistics for each thread. The following information is recorded:

- $\underline{user_time}(thread)$ — the total user run time for $thread$
- $\underline{system_time}(thread)$ — the total system run time for $thread$

- $\underline{cpu_time}(thread)$ — *thread's scaled CPU usage*
- $\underline{sleep_time}(thread)$ — *the amount of time for which $thread$ has been sleeping*

Mach Definition 25

$\underline{ThreadStatistics}$ <hr/> $\underline{ThreadExist}$ $\underline{user_time} : THREAD \rightarrow \mathbb{N}$ $\underline{system_time} : THREAD \rightarrow \mathbb{N}$ $\underline{cpu_time} : THREAD \rightarrow \mathbb{N}$ $\underline{sleep_time} : THREAD \rightarrow \mathbb{N}$ <hr/> $\text{dom } \underline{user_time} = \text{dom } \underline{system_time} = \text{dom } \underline{cpu_time} = \text{dom } \underline{sleep_time}$ $= \underline{thread_exists}$

Review Note:

Should the domain of $\underline{sleep_time}$ be all threads or only those with a particular run state?

4.2.9 Machine State

The system records the machine state of each thread. Typically, the structure of the machine state varies depending upon the architecture of the machine to which the thread is assigned. The type $SUPP_MACHINE_ARCH$ represents the set of supported machine architectures. The set $THREAD_STATE_INFO_TYPES$ denotes the names of the various structures that are associated with the supported architectures. The type $THREAD_STATE_INFO$ denotes the possible values of the state information recorded for a thread.

The expression $\underline{State_info_avail}(arch)$ denotes the types of state information which the architecture supports.

Mach Definition 26

[$SUPP_MACHINE_ARCH$]
[$THREAD_STATE_INFO_TYPES, THREAD_STATE_INFO$]

$\underline{State_info_avail} : SUPP_MACHINE_ARCH$ $\rightarrow \mathbb{P} \text{ } THREAD_STATE_INFO_TYPES$

The expression $\underline{thread_state}(thread, info_type)$ returns the indicated type of state information recorded for $thread$.

Mach Definition 27

$\underline{ThreadMachineState}$ <hr/> $\underline{ThreadExist}$ $\underline{thread_state} : THREAD \times THREAD_STATE_INFO_TYPES$ $\rightarrow THREAD_STATE_INFO$ <hr/> $\text{dom } \underline{thread_state} = \underline{thread_exists} \times THREAD_STATE_INFO_TYPES$
--

Review Note:

Actually, the current instruction pointer is part of the machine state rather than being a separate state component.

Mach Definition 28

<p><i>Threads</i></p> <hr/> <p><i>TasksAndThreads</i></p> <p><i>ThreadPri</i></p> <p><i>ThreadSchedPolicy</i></p> <p><i>ThreadInstruction</i></p> <p><i>ThreadExecStatus</i></p> <p><i>ThreadStatistics</i></p> <p><i>ThreadMachineState</i></p> <p><i>ThreadSampling</i></p> <p><i>TaskSampling</i></p>
--

4.3 Port Name Space

Each task uses its own (local) set of names to refer to ports. The set *NAME* is used to name ports in a task's name space.

Mach Definition 29

[*NAME*]

The names *Mach_port_null* and *Mach_port_dead* are reserved. They will never be used as an index in a task's port name space. The remainder of this section discusses the three types of entities that can be in name spaces: port rights, port sets, and dead names.

Mach Definition 30

<p><i>Mach_port_dead</i> : <i>NAME</i></p> <p><i>Mach_port_null</i> : <i>NAME</i></p>

4.3.1 Port Rights

A port is only of use to a task if the task holds some kind of right to the port. The types of available rights are defined via the type *RIGHT*. A right for a port allows a task to either send or receive messages via that port. The task may have either a general right to send messages via a port or a one-time right to do so. Thus, the elements of type *RIGHT* are: *Send*, *Receive*, and *Send_once*.

A *Capability* is the combination of a port and a right to do something with that port.

Strictly speaking, a task associates a name with a particular right to a port, not simply with the port. The set *port_right_rel* relates the ports to which a task has rights with their right types and their local names. More specifically, each element of *port_right_rel* is a tuple of the form (*task*, *port*, *name*, *right*, *i*). Such a tuple is an element of *port_right_rel* only when *name* denotes in *task*'s name space a right of type *right* to *port*. The *i*-value is used to allow a task to accumulate multiple send rights under the same name. For send-once or receive rights, the

value of i is always equal to 1. For convenience, the expression $named_port(task, name)$ denotes the port associated with $name$ in $task$'s name space.

At most one task can receive messages from a port at any given time. The expression $receiver(port)$ denotes the task (if any) that is currently permitted to receive messages from $port$, and $receiver_name(port)$ denotes the receiver task's name for the port.

Many tasks may have $Send$ or $Send_once$ rights to a port. The relation $sender$ indicates the tasks currently permitted to send messages to a port; an element $(port, task)$ is in $sender$ exactly when $task$ has a send right to $port$.

Mach Definition 31

$$RIGHT ::= Send \mid Receive \mid Send_once$$

Capability

$$port : PORT$$

$$right : RIGHT$$

TasksAndPorts

TaskExist

PortExist

$$\underline{port_right_rel} : \mathbb{P}(TASK \times PORT \times NAME \times RIGHT \times \mathbb{N}_1)$$

$$named_port : TASK \times NAME \leftrightarrow PORT$$

$$receiver : PORT \leftrightarrow TASK$$

$$receiver_name : PORT \leftrightarrow NAME$$

$$sender : PORT \leftrightarrow TASK$$

$$\underline{port_right_rel} \subseteq \underline{task_exists} \times \underline{port_exists} \times NAME \times RIGHT \times \mathbb{N}_1$$

$$\forall task : TASK; port : PORT; right : RIGHT; i : \mathbb{N}_1$$

$$\bullet (task, port, Mach_port_null, right, i) \notin \underline{port_right_rel}$$

$$\wedge (task, port, Mach_port_dead, right, i) \notin \underline{port_right_rel}$$

$$named_port = \{ task : TASK; port : PORT; name : NAME; right : RIGHT; i : \mathbb{N}_1$$

$$\mid (task, port, name, right, i) \in \underline{port_right_rel} \bullet ((task, name), port) \}$$

$$receiver = \{ task : TASK; port : PORT; name : NAME$$

$$\mid (task, port, name, Receive, 1) \in \underline{port_right_rel} \bullet (port, task) \}$$

$$receiver_name = \{ task : TASK; port : PORT; name : NAME$$

$$\mid (task, port, name, Receive, 1) \in \underline{port_right_rel} \bullet (port, name) \}$$

$$sender = \{ task : TASK; port : PORT; name : NAME; right : RIGHT; i : \mathbb{N}_1$$

$$\mid (((task, port, name, right, i) \in \underline{port_right_rel}) \wedge right \in \{Send, Send_once\})$$

$$\bullet (port, task) \}$$

The i -value is called the user reference count. As noted above, it is equal to 1 for receive and send-once rights, but is of interest for send rights. The expression $s_right_ref_count(task, name)$ returns the user reference count for $name$ in $task$'s name space (when it is a send right). There is a system-wide maximum number of references to a given send right which a task may accumulate, represented by Max_right_refs .

Mach Definition 32

$$\mid Max_right_refs : \mathbb{N}_1$$

Mach Definition 33

<p><i>UserReferenceCount</i></p> <hr/> <p><i>TasksAndPorts</i></p> <p>$s_right_ref_count : TASK \times NAME \leftrightarrow \mathbb{N}_1$</p> <p>$\forall task : TASK; port : PORT; name : NAME; right : \{Receive, Send_once\}; i : \mathbb{N}_1$ $\bullet (task, port, name, right, i) \in \underline{port_right_rel} \Rightarrow i = 1$ $s_right_ref_count = \{ task : TASK; port : PORT; name : NAME; i : \mathbb{N}_1$ $\mid (task, port, name, Send, i) \in \underline{port_right_rel} \bullet ((task, name), i) \}$ $\forall task : TASK; name : NAME \bullet s_right_ref_count(task, name) \leq Max_right_refs$</p>

For convenience:

- The relations s_right , r_right , and so_right are used to identify the names of each of the types of rights which are associated with a given task. For example, $(task, name)$ is an element of s_right exactly when $name$ is a send right in $task$'s name space.
- The relation s_r_right is used to identify names that are either a receive or a send right.
- The relation $port_right_namep$ identifies names that are either receive, send, or send-once rights.

The semantics of Mach are such that send and receive rights within a task coalesce into a single name. In other words:

- If $name$ is a receive right for $port$ in $task$'s name space, then no other name in $task$'s name space may be a send right for $port$; the send rights must be associated with $name$, too.
- If $name$ is a send right for $port$ in $task$'s name space, then all of the send rights for $port$ in $task$'s name space are associated with $name$.

Note, however, that the same task can have multiple names associated with send-once rights for the same port. Mach prohibits a name that is a send or a receive right from also being a send-once right.

A message may be forcibly enqueued using a send right. In this case it will be added to the message queue of the named port even if the queue has reached its designated size limit. At most one message may be forcibly enqueued at a time using any given send right. After that message is removed from the queue, a message-accepted notification is sent and the send right can again be used to forcibly enqueue a message. The component $forcibly_queued(task, name)$ denotes the message, if any, forcibly enqueued using a send right $name$ in $task$'s ipc name space.

Mach Definition 34

<p><i>TasksAndRights</i></p> <p><i>MessageExist</i></p> <p><i>TasksAndPorts</i></p> <p>$s_right : TASK \leftrightarrow NAME$</p> <p>$r_right : TASK \leftrightarrow NAME$</p> <p>$so_right : TASK \leftrightarrow NAME$</p> <p>$s_r_right : TASK \leftrightarrow NAME$</p> <p>$port_right_namep : TASK \leftrightarrow NAME$</p> <p>$f_orcibly_queued : (TASK \times NAME) \leftrightarrow MESSAGE$</p> <hr/> <p>$s_right = \{ task : TASK; port : PORT; name : NAME; i : \mathbb{N}_1$ $\mid (task, port, name, Send, i) \in \underline{port_right_rel} \bullet (task, name) \}$</p> <p>$r_right = \{ task : TASK; port : PORT; name : NAME$ $\mid (task, port, name, Receive, 1) \in \underline{port_right_rel} \bullet (task, name) \}$</p> <p>$so_right = \{ task : TASK; port : PORT; name : NAME$ $\mid (task, port, name, Send_once, 1) \in \underline{port_right_rel} \bullet (task, name) \}$</p> <p>$s_r_right = s_right \cup r_right$</p> <p>$port_right_namep = s_r_right \cup so_right$</p> <p>$dom f_orcibly_queued \subseteq s_right$</p> <p>$ran f_orcibly_queued \subseteq \underline{message_exists}$</p> <p>$disjoint \{so_right, s_r_right\}$</p> <p>$\forall task : TASK; name_1, name_2 : NAME$</p> <p>$\bullet (task, name_1) \in s_r_right \wedge (task, name_2) \in s_r_right$ $\wedge \underline{named_port}(task, name_1) = \underline{named_port}(task, name_2)$ $\Rightarrow name_1 = name_2$</p>

Review Note:

I'd like to tie the message indicated by *f_orcibly_queued* back to the port indicated by the send right, but I'm not sure this will be accurate.

4.3.2 Port Sets

A port set is a set of ports associated with a particular task and name. A port set is used to allow the receiving of a message via any member of the port set. Given a task and a port set name, the expression $port_set(task, name)$ denotes the port set. The relation $port_set_namep$ identifies the port set names associated with each task. $containing_set(port)$ denotes the name of the port set containing *port*, if any. Note that a port can be in at most one port set.

Mach prohibits the reserved names *Mach_port_null* and *Mach_port_dead* from being port set names or the inclusion of the same receive right in two different port sets.

Mach Definition 35

<p><i>PortSets</i></p> <p><i>TaskExist</i></p> <p><i>TasksAndRights</i></p> <p>$\underline{port_set_rel} : \mathbb{P}(TASK \times NAME \times \mathbb{P} PORT)$</p> <p>$port_set : (TASK \times NAME) \leftrightarrow \mathbb{P} PORT$</p> <p>$port_set_namep : TASK \leftrightarrow NAME$</p> <p>$containing_set : PORT \leftrightarrow NAME$</p> <hr/> <p>$port_set = \{task : TASK; name : NAME; set_of_ports : \mathbb{P} PORT$ $\mid (task, name, set_of_ports) \in \underline{port_set_rel} \bullet ((task, name), set_of_ports)\}$</p> <p>$port_set_namep = \text{dom } port_set$</p> <p>$containing_set = \{task : TASK; name : NAME; port : PORT$ $\mid (task, name) \in port_set_namep \wedge port \in port_set(task, name)$ $\bullet (port, name)\}$</p> <p>$\text{dom } port_set_namep \subseteq \underline{task_exists}$</p> <p>$\forall task : TASK; name : NAME; port : PORT \mid (task, name) \in \text{dom } port_set$ $\bullet port \in port_set(task, name) \Rightarrow task = receiver(port)$</p> <p>$\forall task : TASK; set_of_ports : \mathbb{P} PORT$ $\bullet ((task, Mach_port_null), set_of_ports) \notin port_set$ $\wedge ((task, Mach_port_dead), set_of_ports) \notin port_set$</p> <p>$\forall task : TASK; name_1, name_2 : NAME$ $\mid (task, name_1) \in \text{dom } port_set \wedge (task, name_2) \in \text{dom } port_set$ $\bullet name_1 \neq name_2 \Rightarrow \text{disjoint } \langle port_set(task, name_1), port_set(task, name_2) \rangle$</p>
--

4.3.3 Dead Rights

A dead name is a name which previously named a send, receive, or send-once right for a task, but no longer does.⁵ Each dead name in a task can have an associated count that is analogous to the reference count associated with send rights. This count is initially set based on the user reference counts for the right previously bearing the name. The count may be modified by subsequent actions of the kernel. The relation $\underline{dead_right_rel}$ identifies the dead names and their associated counts for each task; an element $(task, name, i)$ is an element of $\underline{dead_right_rel}$ if $name$ is a dead name in $task$ with associated count i . The previously defined constant, Max_right_refs , is a system-wide maximum for the reference count of a given dead right. For convenience:

- The relation $dead_namep$ identifies the dead names associated with each task.
- The expression $dead_right_ref_count(task, name)$ denotes the count associated with $name$ in $task$ (when $name$ is a dead name).

Mach prohibits $Mach_port_null$ and $Mach_port_dead$ from being dead names.

Mach Definition 36

⁵A dead name may also be specified in the body of a message in place of an actual port right.

<p><i>DeadRights</i></p> <p>$\underline{dead_right_rel} : \mathbb{P}(TASK \times NAME \times \mathbb{N}_1)$ $\underline{dead_right_ref_count} : TASK \times NAME \mapsto \mathbb{N}_1$ $\underline{dead_namep} : TASK \leftrightarrow NAME$</p> <p>$\underline{dead_right_ref_count} = \{task : TASK; name : NAME; i : \mathbb{N}_1 \mid (task, name, i) \in \underline{dead_right_rel} \bullet ((task, name), i)\}$ $\underline{dead_namep} = \text{dom } \underline{dead_right_ref_count}$ $\forall task : TASK; name : NAME$ $\bullet \underline{dead_right_ref_count}(task, name) \leq \text{Max_right_refs}$ $\forall task : TASK$ $\bullet (task, \text{Mach_port_null}) \notin \underline{dead_namep}$ $\wedge (task, \text{Mach_port_dead}) \notin \underline{dead_namep}$</p>
--

4.3.4 Summary

A task's port right names (send, receive, and send-once), port set names, and dead names are mutually disjoint. The union of *port_right_namep*, *port_set_namep*, and *dead_namep* identifies the names in each task's name space. For convenience:

- The relation *local_namep* is used to denote this union.
- The expression *number_of_rights(task)* is used to denote the number of names that *local_namep* associates with *task*. This is the current size of *task*'s name space.

Mach Definition 37

<p><i>PortNameSpace</i></p> <p><i>TaskExist</i> <i>TasksAndPorts</i> <i>TasksAndRights</i> <i>UserReferenceCount</i> <i>PortSets</i> <i>DeadRights</i> $\underline{local_namep} : TASK \leftrightarrow NAME$ $\underline{number_of_rights} : TASK \mapsto \mathbb{N}$</p> <p>$\text{disjoint } \langle \underline{port_right_namep}, \underline{port_set_namep}, \underline{dead_namep} \rangle$ $\underline{local_namep} = \underline{port_right_namep} \cup \underline{port_set_namep} \cup \underline{dead_namep}$ $\text{dom } \underline{number_of_rights} = \underline{task_exists}$ $\forall task : TASK \mid task \in \underline{task_exists}$ $\bullet \underline{number_of_rights}(task) = \#(\underline{local_namep}(\{ task \}))$</p>
--

4.4 Ports

This section describes data structures associated with ports.

4.4.1 Make Send Count

Each time the receiver for a port creates a new send right for the port, the system increments a counter associated with the port. The expression *make_send_count(port)* denotes the value

of the counter associated with *port*. Note that this count does not necessarily represent the current number of send rights for the port since tasks other than the receiver can create send rights. Furthermore, the count does not necessarily represent the number of send rights the receiver has created because the count can directly be set to arbitrary values by user threads.

Mach Definition 38

$\underline{SendRightsCount}$ $\underline{PortExist}$ $\underline{make_send_count} : PORT \leftrightarrow \mathbb{N}$ $\text{dom } \underline{make_send_count} = \underline{port_exists}$

4.4.2 Message Queues

Each port has an associated message queue. A message queue can be thought of as a sequence of messages. In Mach, a task may set a limit on the number of messages that are permitted in a given message queue. The value *Mach_port_q_limit_default* represents the default limit the kernel uses for newly allocated ports. The value *Mach_port_q_limit_max* represents a system-imposed limit on the value a task may specify as the limit for a message queue.

Mach Definition 39

$\underline{Mach_port_q_limit_max} : \mathbb{N}$ $\underline{Mach_port_q_limit_default} : \mathbb{N}$
--

For each port, *q_limit(port)* indicates the current limit set for the port. This denotes an intended bound on the number of messages in the associated message queue. The expression *port_size(port)* indicates the number of messages that are actually present in *port*'s message queue. Although it is intended that *port_size(port)* is always less than or equal to *q_limit(port)*, the kernel does not actually guarantee that this property always holds. Examples of ways in which the property may be violated include:

- The intended bound on the number of messages in a queue can be decreased below the number of messages already in the queue.
- Messages sent with a send-once right are delivered regardless of whether the destination port's queue is already full.
- Each name for a send right to a port may be used to forcibly enqueue one message at a time to the named full port.

The expression *message_in_port_rel(port)* denotes the sequence of messages in the queue associated with *port*. Each message is contained in at most one message queue. For convenience, the expression *containing_port(message)* is used to indicate the port associated with the message queue to which *message* belongs.

Each port has an associated sequence number that is used to properly sequence messages received through the port. The expression *sequence_no(port)* indicates *port*'s current sequence number.

Mach Definition 40

<p><i>MessageQueues</i></p> <p><i>PortExist</i></p> <p>$\underline{q_limit} : PORT \leftrightarrow \mathbb{N}$</p> <p>$\underline{message_in_port_rel} : PORT \rightarrow \text{iseq } MESSAGE$</p> <p>$\underline{port_size} : PORT \leftrightarrow \mathbb{N}$</p> <p>$\underline{containing_port} : MESSAGE \leftrightarrow PORT$</p> <p>$\underline{sequence_no} : PORT \leftrightarrow \mathbb{Z}$</p> <p>$\underline{containing_port} = \{ message : MESSAGE; port : PORT$ $\quad message \in \text{ran}(\underline{message_in_port_rel}(port)) \bullet message \mapsto port \}$</p> <p>$(\forall port : \underline{port_exists}$ $\bullet \underline{port_size}(port) = \#(\underline{message_in_port_rel}(port))$ $\quad \wedge \underline{q_limit}(port) \leq Mach_port_q_limit_max)$</p> <p>$\text{dom } \underline{q_limit} = \underline{port_exists}$</p> <p>$\text{dom } \underline{message_in_port_rel} = \underline{port_exists}$</p> <p>$\text{dom } \underline{sequence_no} = \underline{port_exists}$</p>
--

4.4.3 Summary

The data structures defined in this section consist of make-send counts, message queues, and sequence numbers associated with ports.

Mach Definition 41

<p><i>PortSummary</i></p> <p><i>SendRightsCount</i></p> <p><i>MessageQueues</i></p>

4.5 Notifications

A task may request that a notification message be sent when one of the following changes occurs in the status of a port:

- The port is destroyed.
- The last send right for the port is deallocated.

A task may also request a notification message be sent when a send right becomes a dead name. In each case, the task requesting the notification must register a port to which the notification should be sent.

The relation $\underline{port_notify_destroyed_rel}$ identifies the ports for which a destroyed notification has been requested and the associated notification ports. For convenience, $\underline{port_notify_destroyed}(port)$ is used to denote the notification port registered for a destroyed notification on $port$.

The relation $\underline{port_notify_no_more_senders_rel}$ identifies the ports for which a no-more-senders notification has been requested and the associated notification ports. For convenience, $\underline{port_notify_no_more_senders}(port)$ is used to denote the notification port registered for a no-more-senders notification on $port$.

The relation $\underline{port_notify_dead_rel}$ identifies the task-name pairs for which a dead-name notification has been requested and the associated notification ports. For convenience,

$port_notify_dead(task, name)$ is used to denote the notification port registered for a dead-name notification on $name$ in $task$'s name space.

The registered notification ports remain in force as long as both the port in question and the registered port exist regardless of whether the same tasks remain related to these ports.

Mach Definition 42

<p><i>Notifications</i></p> <p><i>PortExist</i></p> <p><i>TasksAndPorts</i></p> <p>$\underline{port_notify_destroyed_rel} : PORT \leftrightarrow PORT$</p> <p>$\underline{port_notify_no_more_senders_rel} : PORT \leftrightarrow PORT$</p> <p>$\underline{port_notify_dead_rel} : \mathbb{P}(PORT \times TASK \times NAME)$</p> <p>$port_notify_destroyed : PORT \mapsto PORT$</p> <p>$port_notify_no_more_senders : PORT \mapsto PORT$</p> <p>$port_notify_dead : TASK \times NAME \mapsto PORT$</p> <hr/> <p>$port_notify_destroyed = \underline{port_notify_destroyed_rel}$</p> <p>$port_notify_no_more_senders = \underline{port_notify_no_more_senders_rel}$</p> <p>$\forall task : TASK; port : PORT; name : NAME$</p> <ul style="list-style-type: none"> • $((port, task, name) \in \underline{port_notify_dead_rel} \Leftrightarrow ((task, name), port) \in port_notify_dead)$ <p>$dom\ port_notify_destroyed = \underline{port_exists}$</p> <p>$dom\ port_notify_no_more_senders = \underline{port_exists}$</p> <p>$dom\ port_notify_dead = dom\ named_port$</p> <p>$ran\ port_notify_destroyed \subseteq \underline{port_exists} \cup \{Ip_null\}$</p> <p>$ran\ port_notify_dead \subseteq \underline{port_exists} \cup \{Ip_null\}$</p> <p>$ran\ port_notify_no_more_senders \subseteq \underline{port_exists} \cup \{Ip_null\}$</p>

Review Note:
Should the range of these functions also include Ip_dead ? It seems that it should because the port could die. Should look at the code to see what happens if we try to send a notification in this situation.

4.6 Special Ports

This section describes the special ports known to the kernel. Each of the special ports is associated with some kernel entity.

4.6.1 Task Ports

In addition to the ports referenced in its port name space, each task has four special ports. The self port is used to request the kernel to perform actions upon the task. Any task holding a send right to a second task may use that right to request operations on the second task. The kernel is always the receiver for a task's self port. A task's sself port is normally equal to its self port, but may refer to a different port and have a task other than the kernel, such as a debugger, as its receiver. The relations $\underline{task_self_rel}$ and $\underline{task_sself_rel}$ identify the self and sself ports associated with each task.

The other two special ports are the exception port and the bootstrap port. A task receives exception messages from the kernel via its exception port. A task's bootstrap port is provided as a start-up means for a task to obtain a send right to a service port for a server that can provide the task start-up information. The relations $\underline{task_eport_rel}$ and $\underline{task_bport_rel}$ identify the exception port and bootstrap port associated with each task. The $sself$, exception and bootstrap ports may be modified. Unlike the self port, they may become Ip_null or Ip_dead .

For convenience:

- The expression $task_self(task)$ denotes $task$'s self port.
- The expression $task_sself(task)$ denotes $task$'s sself port.
- The expression $task_eport(task)$ denotes $task$'s exception port.
- The expression $task_bport(task)$ denotes $task$'s bootstrap port.
- The expression $self_task(port)$ denotes the task (if any) having $port$ as its self port.

Mach Definition 43

<p><i>SpecialTaskPorts</i></p> <p><i>TaskExist</i></p> <p><i>PortExist</i></p> <p><i>Kernel</i></p> <p><i>TasksAndPorts</i></p> <p>$\underline{task_self_rel} : TASK \leftrightarrow PORT$</p> <p>$\underline{task_sself_rel} : TASK \leftrightarrow PORT$</p> <p>$\underline{task_eport_rel} : TASK \leftrightarrow PORT$</p> <p>$\underline{task_bport_rel} : TASK \leftrightarrow PORT$</p> <p>$task_self : TASK \rightsquigarrow PORT$</p> <p>$task_sself : TASK \leftrightarrow PORT$</p> <p>$task_eport : TASK \leftrightarrow PORT$</p> <p>$task_bport : TASK \leftrightarrow PORT$</p> <p>$self_task : PORT \rightsquigarrow TASK$</p> <hr/> <p>$task_self = \underline{task_self_rel}$</p> <p>$task_eport = \underline{task_eport_rel}$</p> <p>$task_bport = \underline{task_bport_rel}$</p> <p>$task_sself = \underline{task_sself_rel}$</p> <p>$dom\ task_self = dom\ task_sself = dom\ task_eport = dom\ task_bport = \underline{task_exists}$</p> <p>$ran\ task_self \subset \underline{port_exists}$</p> <p>$ran\ task_sself \subset \underline{port_pointer}$</p> <p>$ran\ task_eport \subset \underline{port_pointer}$</p> <p>$ran\ task_bport \subset \underline{port_pointer}$</p> <p>$self_task = \underline{port_exists} \triangleleft (task_self \sim)$</p> <p>$\forall task : TASK \mid task \in \underline{task_exists} \bullet receiver(task_self(task)) = \underline{kernel}$</p>

4.6.2 Thread Ports

Each thread has a self port, sself port, and an exception port with purposes parallel to the corresponding special ports for tasks. The relations and functions $\underline{thread_self_rel}$, $\underline{thread_sself_rel}$, $\underline{thread_eport_rel}$, $\underline{thread_self}$, $\underline{thread_sself}$, $\underline{thread_eport}$, and $\underline{self_thread}$ are used to denote these state components.

Mach Definition 44

<p><i>SpecialThreadPorts</i></p> <p><i>ThreadExist</i></p> <p><i>PortExist</i></p> <p><i>TasksAndPorts</i></p> <p><i>Kernel</i></p> <p>$\underline{t}hread_self_rel : THREAD \leftrightarrow PORT$</p> <p>$\underline{t}hread_sself_rel : THREAD \leftrightarrow PORT$</p> <p>$\underline{t}hread_eport_rel : THREAD \leftrightarrow PORT$</p> <p>$thread_self : THREAD \rightsquigarrow PORT$</p> <p>$thread_sself : THREAD \rightsquigarrow PORT$</p> <p>$thread_eport : THREAD \rightsquigarrow PORT$</p> <p>$self_thread : PORT \rightsquigarrow THREAD$</p> <hr/> <p>$\underline{t}hread_self_rel = thread_self$</p> <p>$\underline{t}hread_sself_rel = thread_sself$</p> <p>$\underline{t}hread_eport_rel = thread_eport$</p> <p>$dom\ thread_self = \underline{t}hread_exists$</p> <p>$dom\ thread_sself = \underline{t}hread_exists$</p> <p>$dom\ thread_eport = \underline{t}hread_exists$</p> <p>$ran\ thread_self \subset \underline{p}ort_exists$</p> <p>$ran\ thread_sself \subset \underline{p}ort_pointer$</p> <p>$ran\ thread_eport \subset \underline{p}ort_pointer$</p> <p>$self_thread = \underline{p}ort_exists \triangleleft (thread_self \sim)$</p> <p>$\forall thread : THREAD \mid thread \in \underline{t}hread_exists \bullet receiver(thread_self(thread)) = \underline{k}ernel$</p>

4.6.3 Memory Ports

A kernel and a memory object interact by engaging in a dialogue. The kernel sends messages to an object port and the object manager sends messages to a control port. There is also a name port used to identify the object in **vm_region** requests. The relations *object_port_rel*, *control_port_rel*, and *name_port_rel* are used to represent the binding between a memory and its associated ports. For a particular Mach host kernel, there is at most one of each type of port associated with a given memory. Furthermore, no object port is associated with more than one memory object. For convenience:

- The expressions *object_port(memory)*, *control_port(memory)*, and *name_port(memory)* are used to denote, respectively, the object, control, and name port for *memory*.
- The expression *object_memory(port)* denotes the memory object (if any) for which *port* is the object port.
- The expression *control_memory(port)* denotes the memory object (if any) for which *port* is the control port.

Memory objects are given a name port immediately upon allocation. However, they need not necessarily have object and control ports until a page that they back needs to be paged out.

Mach Definition 45

<p><i>MemoriesAndPorts</i></p> <hr/> <p><i>Kernel</i></p> <p><i>MemoryExist</i></p> <p><i>TasksAndPorts</i></p> <p>$\underline{o}bject_port_rel : MEMORY \leftrightarrow PORT$</p> <p>$\underline{c}ontrol_port_rel : MEMORY \leftrightarrow PORT$</p> <p>$\underline{n}ame_port_rel : MEMORY \leftrightarrow PORT$</p> <p>$object_port : MEMORY \rightsquigarrow PORT$</p> <p>$control_port : MEMORY \rightsquigarrow PORT$</p> <p>$name_port : MEMORY \rightsquigarrow PORT$</p> <p>$object_memory : PORT \rightsquigarrow MEMORY$</p> <p>$control_memory : PORT \rightsquigarrow MEMORY$</p> <hr/> <p>$\underline{o}bject_port_rel = object_port$</p> <p>$\underline{c}ontrol_port_rel = control_port$</p> <p>$\underline{n}ame_port_rel = name_port$</p> <p>$object_port^{\sim} = object_memory$</p> <p>$control_port^{\sim} = control_memory$</p> <p>$dom \underline{c}ontrol_port_rel = dom \underline{o}bject_port_rel \subseteq dom \underline{n}ame_port_rel$</p> <p>$dom \underline{n}ame_port_rel = \underline{m}emory_exists$</p> <p>$\forall port : PORT \mid port \in ran \underline{c}ontrol_port_rel$</p> <ul style="list-style-type: none"> • $port \in dom receiver \wedge receiver(port) = \underline{k}ernel$ <p>$\forall port : PORT \mid port \in ran \underline{n}ame_port_rel$</p> <ul style="list-style-type: none"> • $port \in dom receiver \wedge receiver(port) = \underline{k}ernel$
--

4.6.4 Host Ports

Each host has two associated ports: the control port and the name port. These ports are denoted by $\underline{h}ost_control_port$ and $\underline{h}ost_name_port$. The kernel is the receiver for each of these ports. The name port is used to service “unprivileged” requests while the control port is used to service “privileged” requests.

Mach Definition 46

<p><i>HostsAndPorts</i></p> <hr/> <p><i>Kernel</i></p> <p><i>TasksAndPorts</i></p> <p>$\underline{h}ost_control_port : PORT$</p> <p>$\underline{h}ost_name_port : PORT$</p> <hr/> <p>$(\underline{h}ost_name_port, \underline{k}ernel) \in receiver$</p> <p>$(\underline{h}ost_control_port, \underline{k}ernel) \in receiver$</p>
--

4.6.5 Processor Ports

Each processor has a port that is used to name it. The relation $\underline{p}rocessor_port_rel$ indicates the association between processors and their name ports. There is exactly one port associated with each processor. For convenience, $\underline{p}roc_self(proc)$ and $\underline{t}he_processor(port)$ are used to denote, respectively, the port associated with a given processor and the processor associated with a given port.

Each processor set has two associated ports: the control port and the name port. The relations $\underline{ps_control_port_rel}$ and $\underline{ps_name_port_rel}$ are used to represent the binding between a processor set and its associated ports. In Mach, there is exactly one of each type of port associated with each existing processor set. For convenience:

- The expression $\text{controlled_proc_set}(port)$ is used to indicate the processor set (if any) having $port$ as its control port.
- The expression $\text{procset_self}(procset)$ is used to indicate $procset$'s control port.
- The expression $\text{named_proc_set}(port)$ is used to indicate the processor set (if any) having $port$ as its name port.
- The expression $\text{procset_name_port}(procset)$ is used to indicate $procset$'s name port.

Mach Definition 47

<p style="margin: 0;"><i>ProcessorsAndPorts</i></p> <hr style="border: 0.5px solid black;"/> <p style="margin: 0;"><i>Kernel</i></p> <p style="margin: 0;"><i>TasksAndPorts</i></p> <p style="margin: 0;">$\underline{processor_port_rel} : PROCESSOR \leftrightarrow PORT$</p> <p style="margin: 0;">$\underline{ps_control_port_rel} : PROCESSOR_SET \leftrightarrow PORT$</p> <p style="margin: 0;">$\underline{ps_name_port_rel} : PROCESSOR_SET \leftrightarrow PORT$</p> <p style="margin: 0;">$\underline{proc_self} : PROCESSOR \rightarrow PORT$</p> <p style="margin: 0;">$\underline{the_processor} : PORT \rightarrow PROCESSOR$</p> <p style="margin: 0;">$\underline{controlled_proc_set} : PORT \rightarrow PROCESSOR_SET$</p> <p style="margin: 0;">$\underline{procset_self} : PROCESSOR_SET \rightarrow PORT$</p> <p style="margin: 0;">$\underline{named_proc_set} : PORT \rightarrow PROCESSOR_SET$</p> <p style="margin: 0;">$\underline{procset_name_port} : PROCESSOR_SET \rightarrow PORT$</p> <hr style="border: 0.5px solid black;"/> <p style="margin: 0;">$\text{dom } \underline{ps_control_port_rel} = \text{dom } \underline{ps_name_port_rel}$</p> <p style="margin: 0;">$\underline{processor_port_rel} \sim = \underline{the_processor}$</p> <p style="margin: 0;">$\underline{processor_port_rel} = \underline{proc_self}$</p> <p style="margin: 0;">$\underline{ps_control_port_rel} \sim = \underline{controlled_proc_set}$</p> <p style="margin: 0;">$\underline{ps_control_port_rel} = \underline{procset_self}$</p> <p style="margin: 0;">$\underline{ps_name_port_rel} \sim = \underline{named_proc_set}$</p> <p style="margin: 0;">$\underline{ps_name_port_rel} = \underline{procset_name_port}$</p> <p style="margin: 0;">$\forall port : PORT \mid port \in \text{ran } \underline{ps_control_port_rel}$</p> <ul style="list-style-type: none"> • $port \in \text{dom } receiver \wedge receiver(port) = \underline{k}ernel$ <p style="margin: 0;">$\forall port : PORT \mid port \in \text{ran } \underline{ps_name_port_rel}$</p> <ul style="list-style-type: none"> • $port \in \text{dom } receiver \wedge receiver(port) = \underline{k}ernel$ <p style="margin: 0;">$\forall port : PORT \mid port \in \text{ran } \underline{processor_port_rel}$</p> <ul style="list-style-type: none"> • $port \in \text{dom } receiver \wedge receiver(port) = \underline{k}ernel$
--

4.6.6 Device Ports

Each device is represented by a unique port. The relation $\underline{device_port_rel}$ identifies the device port representing each device. The kernel is the receiver for a device port. For convenience:

- The expression $\text{device_port}(dev)$ is used to denote dev 's device port.
- The expression $\text{port_device}(port)$ is used to denote the device (if any) having $port$ as its device port.

Mach Definition 48

<i>DevicesAndPorts</i> <i>TasksAndPorts</i> <i>Kernel</i> $\underline{d}evice_port_rel : DEVICE \leftrightarrow PORT$ $device_port : DEVICE \leftrightarrow PORT$ $port_device : PORT \leftrightarrow DEVICE$ <hr/> $device_port = \underline{d}evice_port_rel$ $port_device = \underline{d}evice_port_rel \sim$ $\forall port : PORT \mid port \in \text{ran } \underline{d}evice_port_rel$ <ul style="list-style-type: none"> • $port \in \text{dom } receiver \wedge receiver(port) = \underline{k}ernel$

4.6.7 Device Master Port

Tasks gain access to devices through the device master port which is denoted by *master_device_port*. The kernel is the receiver for this port.

Mach Definition 49

<i>MasterDevicePort</i> <i>TasksAndPorts</i> <i>Kernel</i> $\underline{m}aster_device_port : PORT$ <hr/> $(\underline{m}aster_device_port, \underline{k}ernel) \in receiver$
--

4.6.8 Summary

Each special port for which the kernel is always the receiver must be distinct from all of the other special ports for which the kernel is always the receiver. For example, no two tasks may have the same self port, and no port may be both a task self port and a thread self port. Note, however, that the kernel does not prohibit overlaps between the special ports for which the kernel is always the receiver and the other special ports. For example, a task's bootstrap port might be set to some others task's self port (even though this would probably not serve any useful purpose).

Mach Definition 50

<i>SpecialPurposePorts</i> <i>SpecialTaskPorts</i> <i>SpecialThreadPorts</i> <i>MemoriesAndPorts</i> <i>HostsAndPorts</i> <i>ProcessorsAndPorts</i> <i>DevicesAndPorts</i> <i>MasterDevicePort</i> <hr/> $\text{disjoint } \{ \text{ran } \underline{t}ask_self, \text{ran } \underline{t}hread_self, \text{ran } \underline{c}ontrol_port, \text{ran } \underline{o}bject_port, \\ \text{ran } \underline{n}ame_port, \{ \underline{h}ost_control_port \}, \{ \underline{h}ost_name_port \}, \\ \text{ran } \underline{p}s_control_port_rel, \text{ran } \underline{p}s_name_port_rel, \text{ran } \underline{p}rocessor_port_rel, \\ \text{ran } \underline{d}evice_port_rel, \{ \underline{m}aster_device_port \} \}$

Editorial Note:

The following needs some revision:

- Add port classes for pager name ports and pager (object) ports.
 - Correct the misunderstanding that a port in a port class must have the kernel as the receiver. While this is true for most classes, memory object (pager) ports are a notable exception.
-

The type *PORT_CLASS* denotes the classes of ports for which the kernel is the receiver. These are *Pc_task*, *Pc_thread*, *Pc_host_control*, *Pc_host_name*, *Pc_ps_control*, *Pc_ps_name*, *Pc_processor*, *Pc_memory*, and *Pc_device*.

If the kernel is the receiver for *port*, then the expression $\underline{port_class}(port)$ denotes *port*'s class.

Mach Definition 51

$$\begin{aligned} \text{PORT_CLASS} ::= & Pc_task \mid Pc_thread \mid Pc_host_control \mid Pc_host_name \\ & \mid Pc_ps_control \mid Pc_ps_name \mid Pc_processor \mid Pc_memory \\ & \mid Pc_device \end{aligned}$$

PortClasses

SpecialPurposePorts

$\underline{port_class} : \text{PORT} \mapsto \text{PORT_CLASS}$

$\forall port : \text{PORT}$

- $(port \in \text{ran } \underline{task_self} \Rightarrow (port, Pc_task) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{thread_self} \Rightarrow (port, Pc_thread) \in \underline{port_class})$
 - $\wedge (port = \underline{host_control_port} \Rightarrow (port, Pc_host_control) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{device_port_rel} \Rightarrow (port, Pc_device) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{control_port_rel} \Rightarrow (port, Pc_memory) \in \underline{port_class})$
 - $\wedge (port = \underline{host_name_port} \Rightarrow (port, Pc_host_name) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{ps_control_port_rel} \Rightarrow (port, Pc_ps_control) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{ps_name_port_rel} \Rightarrow (port, Pc_ps_name) \in \underline{port_class})$
 - $\wedge (port \in \text{ran } \underline{processor_port_rel} \Rightarrow (port, Pc_processor) \in \underline{port_class})$
-

4.7 Total Send Rights

In addition to the send rights contained in the port name spaces associated with the tasks, the kernel maintains so-called naked send rights to the special ports. We occasionally need to know the total number of send rights to a given port including both those recorded in a name space and the naked rights. Naked rights are associated with the following ports: *task_sself*, *task_eport*, *task_bport*, *thread_sself* and *thread_eport*. We define *port_right_seq* to be any sequence of the elements of the set *port_right_rel* (the precise ordering of elements is not important for our purposes). The expression *total_name_space_srights(port)* denotes the number of send rights to *port* in all name spaces, and *total_naked_srights(port)* denotes the total number of send rights to *port* that are not stored in any name space. The expression *total_srights(port)* is the sum of these two numbers.

Review Note:

Need to determine if naked send rights are implied by any other special port relationships. Note that a naked send right is *not* created for the self port relationships (e.g., *thread_self*).

Need to determine whether rights in messages count as naked send rights too.

Mach Definition 52

$ \begin{aligned} & \text{TotalSendRights} \\ & \text{PortExist} \\ & \text{TasksAndPorts} \\ & \text{SpecialPurposePorts} \\ & \text{port_right_seq} : \text{seq}(\text{TASK} \times \text{PORT} \times \text{NAME} \times \text{RIGHT} \times \mathbb{N}_1) \\ & \text{total_name_space_srights} : \text{PORT} \leftrightarrow \mathbb{N} \\ & \text{total_naked_srights} : \text{PORT} \leftrightarrow \mathbb{N} \\ & \text{total_srights} : \text{PORT} \leftrightarrow \mathbb{N} \\ \\ & \text{ran } \text{port_right_seq} = \text{port_right_rel} \\ & \# \text{port_right_seq} = \# \text{port_right_rel} \\ \\ & (\forall \text{port} : \text{PORT} \mid \text{port} \in \text{port_exists} \\ & \bullet \text{total_name_space_srights}(\text{port}) \\ & \quad = \text{Seq_plus}(\text{squash} \{ \text{task} : \text{TASK}; \text{name} : \text{NAME}; i, n : \mathbb{N}_1 \\ & \quad \mid (n, (\text{task}, \text{port}, \text{name}, \text{Send}, i)) \in \text{port_right_seq} \\ & \quad \bullet (n, i) \}) \\ & \wedge \text{total_naked_srights}(\text{port}) = \#(\text{task_self} \triangleright \{ \text{port} \}) \\ & \quad + \#(\text{task_eport} \triangleright \{ \text{port} \}) \\ & \quad + \#(\text{task_bport} \triangleright \{ \text{port} \}) \\ & \quad + \#(\text{thread_self} \triangleright \{ \text{port} \}) \\ & \quad + \#(\text{thread_eport} \triangleright \{ \text{port} \}) \\ & \wedge \text{total_srights}(\text{port}) = \text{total_name_space_srights}(\text{port}) + \text{total_naked_srights}(\text{port}) \end{aligned} $
--

4.8 Registered Rights

Each task has a finite array of send rights, intended to use for access to the Network Name Server, the Environment Manager, and the Service server (although they may have any use). These rights are called “registered,” to denote the fact that the kernel knows their identity. The expression *registered_rights(task)* denotes the set of names of rights registered for *task*. There may be more than three registered rights, in fact their number need only be less than or equal to the system constant *Task_port_register_max*. The kernel has three constants *Name_server_slot*, *Environment_slot*, and *Service_slot* which tell it which element of the array refers to each of these servers.

Mach Definition 53

$ \begin{aligned} & \text{Task_port_register_max} : \mathbb{N} \\ & \text{Name_server_slot} : \mathbb{N} \\ & \text{Environment_slot} : \mathbb{N} \\ & \text{Service_slot} : \mathbb{N} \end{aligned} $
--

$\begin{aligned} & \text{RegisteredRights} \\ & \text{TaskExist} \\ & \underline{\text{registered_rights}} : \text{TASK} \leftrightarrow \text{seq PORT} \\ & \text{dom } \underline{\text{registered_rights}} = \underline{\text{task_exists}} \\ & \forall \text{task} : \text{TASK} \mid \text{task} \in \underline{\text{task_exists}} \\ & \bullet \#(\underline{\text{registered_rights}}(\text{task})) \leq \text{Task_port_register_max} \end{aligned}$

4.9 Memory System

This section describes the components of the Mach system that are used to provide tasks with address spaces.

4.9.1 Memory

Each memory can be viewed as mapping a memory offset to a value. Essentially, a memory can be viewed as an array of values indexed by offsets; the only difference is that a memory may have holes in the sense that some offsets do not map to any value. The mapping from offsets to values is defined by the memory's manager. As described in Section 4.9.2, the kernel becomes aware of pieces of this mapping as data is cached in resident pages. The types *OFFSET* and *WORD* denote, respectively, the sets of memory offsets and memory values.

The kernel maintains a copy strategy for each memory object. This strategy is one of the following:

- *Memory_copy_none* —

Review Note:
We need to figure out the meaning of each strategy.

- *Memory_copy_call* —
- *Memory_copy_delay* —
- *Memory_copy_temporary* —

These values comprise the elements of the type *MEMORY_COPY_STRATEGY*. The expression *copy_strategy(memory)* denotes the copy strategy recorded for *memory*.

The kernel cannot request access permissions and data from a memory object until it has received a **memory_object_ready** command (normally in reply to a **memory_object_init** request). The set *initialized* denotes the set of memory objects for which this has occurred.

The kernel records which memory objects may be cached; the set *may_cache* denotes the set of such memory objects. The memory performance for a memory object is influenced by its copy strategy and whether it can be cached.

A memory can be either managed or unmanaged. The set *managed* denotes the set of memories that are managed. Corresponding to each such memory there is a task acting as the memory's manager. The manager for *memory* is denoted by *manager(memory)*. Each memory having an object port is managed.

Similarly, memories can be temporary or non-temporary. The set *temporary_rel* denotes the set of memories that are temporary.

If the page of data corresponding to a given memory-offset pair is not resident when a thread attempts access, then the thread is blocked on a page fault. The expression $\underline{memory_fault}(memory, offset)$ indicates the set of threads that are currently blocked on a page fault generated by access to a given memory-offset pair.

Temporary memory is backed by the default memory manager. The kernel records a port identifying the current default memory manager. This port is denoted by $\underline{default_mem_manager}$.

A null value is used to indicate the lack of a memory filling a particular function in a virtual memory map entry.

Review Note:

Need to figure out how $\underline{default_mem_manager}$ relates to $\underline{managed}$ and $\underline{manager}$.

Mach Definition 54

[*WORD, OFFSET*]

$$\underline{MEMORY_COPY_STRATEGY} ::= \underline{Memory_copy_none} \mid \underline{Memory_copy_call} \\ \mid \underline{Memory_copy_delay} \mid \underline{Memory_copy_temporary}$$

Memory

MemoriesAndPorts

PortExist

$\underline{copy_strategy} : \underline{MEMORY} \leftrightarrow \underline{MEMORY_COPY_STRATEGY}$

$\underline{initialized} : \mathbb{P} \underline{MEMORY}$

$\underline{may_cache} : \mathbb{P} \underline{MEMORY}$

$\underline{managed} : \mathbb{P} \underline{MEMORY}$

$\underline{manager} : \underline{MEMORY} \leftrightarrow \underline{TASK}$

$\underline{temporary_rel} : \mathbb{P} \underline{MEMORY}$

$\underline{memory_fault} : \underline{MEMORY} \times \underline{OFFSET} \leftrightarrow \mathbb{P} \underline{THREAD}$

$\underline{default_mem_manager} : \underline{PORT}$

$\underline{default_mem_manager} \in \underline{port_exists}$

$\underline{managed} = \text{dom } \underline{object_port}$

$\text{dom } \underline{object_port} \subseteq \text{dom } \underline{manager}$

$\underline{initialized} \subseteq \text{dom } \underline{object_port}$

$\underline{may_cache} \subseteq \underline{initialized}$

$\underline{initialized} = \text{dom } \underline{copy_strategy}$

$\forall \text{memory} : \underline{MEMORY}; \text{offset} : \underline{OFFSET}$

$\mid (\text{memory}, \text{offset}) \in \text{dom } \underline{memory_fault}$

$\wedge \underline{memory_fault}(\text{memory}, \text{offset}) \neq \emptyset$

• $\text{memory} \in \underline{managed}$

4.9.2 Pages

At the physical level, pages relate page offsets and values in much the same way as memories relate memory offsets and values. The relation $\underline{page_word_rel}$ identifies the binding between page-offset pairs and words of data. Since at most one value can be stored at a given page offset, $\underline{page_word_rel}$ is actually a function mapping page-offset pairs to values. For convenience,

$page_word_fun(page)(page_offset)$ is used to denote the word of data at offset $page_offset$ of page $page$.

Each page represents some area of memory. The relation $represents_rel$ indicates the binding between pages and memory-offset pairs. This relation should be interpreted as indicating the memory and offset within that memory of the beginning of the data that a page represents. Since each area of memory is represented by at most one page, the function $representing_page$ denoting the page representing an area of memory can be defined. Each page in the range of this function represents some area of memory. For convenience:

- The set $represents_memory$ is used to denote the set of pages that represent some area of memory.
- The set $represented$ is used to denote the set of memory-offset pairs that are represented by some page.
- The expressions $represented_memory(page)$ and $represented_offset(page)$ denote, respectively, the memory and offset that $page$ represents.

When a page is modified, it becomes dirty. The set $dirty_rel$ denotes the set of dirty pages. Upon evicting a page, the kernel checks whether the page is dirty. If it is, then the contents of the page are sent to the appropriate memory manager for it to record the updates. A memory manager may instruct the kernel that it will not retain a copy of a page that it has provided to the kernel by indicating that the page is precious. Whenever the kernel evicts a precious page, it sends the contents of the page to the appropriate memory manager regardless of whether the page is dirty. By instructing the kernel that a page is precious, a memory manager can relieve itself of the responsibility of retaining a copy of a page while the page is resident; the memory manager can rely on the kernel to inform it of the page's current contents whenever the page is evicted. The set $precious$ is used to denote the set of precious pages.

Mach Definition 55

[PAGE_OFFSET]

<p><i>PageAndMemory</i></p> <p>$\underline{page_word_rel} : \mathbb{P}((PAGE \times PAGE_OFFSET) \times WORD)$ $\underline{page_word_fun} : PAGE \mapsto PAGE_OFFSET \rightarrow WORD$ $\underline{represents_rel} : PAGE \leftrightarrow (MEMORY \times OFFSET)$ $\underline{representing_page} : MEMORY \times OFFSET \mapsto PAGE$ $\underline{represents_memory} : \mathbb{P} PAGE$ $\underline{represented} : \mathbb{P}(MEMORY \times OFFSET)$ $\underline{represented_memory} : PAGE \mapsto MEMORY$ $\underline{represented_offset} : PAGE \mapsto OFFSET$ $\underline{dirty_rel} : \mathbb{P} PAGE$ $\underline{precious} : \mathbb{P} PAGE$</p> <p>$(\forall page : PAGE; page_offset : PAGE_OFFSET; word : WORD$ $\bullet \underline{page_word_fun}(page)(page_offset) = word$ $\Leftrightarrow ((page, page_offset), word) \in \underline{page_word_rel})$ $\underline{represents_memory} \subseteq \text{dom } \underline{page_word_fun}$ $\underline{representing_page} = \underline{represents_rel} \sim$ $\underline{dirty_rel} \subseteq \text{ran } \underline{representing_page}$ $\underline{represented} = \text{dom } \underline{representing_page}$ $\underline{represented_memory} = \{memory : MEMORY; offset : OFFSET; page : PAGE$ $\quad (page, (memory, offset)) \in \underline{represents_rel} \bullet (page, memory)\}$ $\underline{represented_offset} = \{memory : MEMORY; offset : OFFSET; page : PAGE$ $\quad (page, (memory, offset)) \in \underline{represents_rel} \bullet (page, offset)\}$ $\underline{precious} \subseteq \underline{represents_memory}$</p>
--

Mach allows pages to be locked against particular types of accesses. This is represented by associating a set of protections with each page. The protections are of type *PROTECTION* which is comprised of the elements *Read*, *Write*, and *Execute*. The relation $\underline{page_lock_rel}$ indicates the access modes against which a page is locked. For convenience $\underline{page_locks}(page)$ is defined to be the set of access modes against which *page* is locked.

Mach Definition 56

$PROTECTION ::= Read \mid Write \mid Execute$

<p><i>Lock</i></p> <p>$\underline{page_lock_rel} : PAGE \leftrightarrow \mathbb{P} PROTECTION$ $\underline{page_locks} : PAGE \mapsto \mathbb{P} PROTECTION$</p> <p>$\underline{page_lock_rel} = \underline{page_locks}$</p>
--

4.9.3 Address Space

The set $\underline{allocated}$ is used to denote the set of *TASK-PAGE_INDEX* pairs that have been allocated in a task's address space. A task-index pair may be mapped to a memory area. Using the previously defined state components, these memory areas can be related to the physical pages used to contain the data when it is paged out. Thus, a task's address space completes the picture of mapping virtual addresses to physical pages and values. Note, however, that not all allocated addresses need be mapped to memory. The relation $\underline{map_rel}$ associates task-index pairs with memory-offset pairs. There is at most one memory-offset pair associated with each task-index pair. For convenience:

- The expressions $mapped_memory(task, index)$ and $mapped_offset(task, index)$ are used to denote the memory and offset corresponding to a given task-index pair.
- The set $mapped$ is used to denote the set of memories to which some task-index pair maps.

Mach Definition 57

[PAGE_INDEX]

<p><i>AddressSpace</i></p> $\underline{map_rel} : (TASK \times PAGE_INDEX) \leftrightarrow (MEMORY \times OFFSET)$ $mapped_memory : TASK \times PAGE_INDEX \mapsto MEMORY$ $mapped_offset : TASK \times PAGE_INDEX \mapsto OFFSET$ $\underline{allocated} : \mathbb{P}(TASK \times PAGE_INDEX)$ $mapped : \mathbb{P} MEMORY$
$\text{dom } \underline{map_rel} = \text{dom } mapped_memory = \text{dom } mapped_offset$ $\text{dom } \underline{map_rel} \subseteq \underline{allocated}$ $mapped = \text{ran } mapped_memory$ $\forall task_va_pair : TASK \times PAGE_INDEX ; memory : MEMORY ; offset : OFFSET$ <ul style="list-style-type: none"> • $(task_va_pair, (memory, offset)) \in \underline{map_rel}$ $\Leftrightarrow (mapped_memory(task_va_pair) = memory$ <li style="padding-left: 2em;">$\wedge mapped_offset(task_va_pair) = offset)$

4.9.4 Memory Protection

Mach protects memory objects by assigning protections to each page in a task’s address space. Three sets of protections are associated with each page in a task’s address space. The Mach protection holds currently applicable protection limits as indicated by users. The maximum protection limits the allowable values for the Mach protection. The third set, the current protections, is what actually limits a task’s access to a page. This is a DTOS addition and will be further defined in Section 5.9.⁶

We use $\underline{mach_protection}$ to denote the relation between tasks, pages, and Mach protection sets. The pair $((task, page_index), protection_set)$ is an element of $\underline{mach_protection}$ if $protection_set$ is the set of protections most recently established by a user request to set the Mach protections for $page_index$. We model maximum protections similarly by defining $\underline{max_protection}(task, page_index)$ to denote the maximum protection that $task$ is permitted to the memory it has mapped at $page_index$.

Mach Definition 58

<p><i>MachProtection</i></p> $\underline{mach_protection} : (TASK \times PAGE_INDEX) \mapsto \mathbb{P} PROTECTION$ $\underline{max_protection} : (TASK \times PAGE_INDEX) \mapsto \mathbb{P} PROTECTION$
$\text{dom } \underline{mach_protection} = \text{dom } \underline{max_protection}$ $\forall task_page_index : TASK \times PAGE_INDEX$ $ task_page_index \in \text{dom } \underline{mach_protection}$ <ul style="list-style-type: none"> • $\underline{mach_protection}(task_page_index) \subseteq \underline{max_protection}(task_page_index)$

⁶The Mach protection in DTOS is called the current protection in Mach and is used in Mach to control a task’s access of pages. The terminology has been changed here to remain consistent with the prototype which must take into account the decisions of the security server when determining the current protections.

4.9.5 Memory Inheritance

For each memory region within a task's address space, Mach records an inheritance attribute that indicates the manner in which child tasks inherit the memory. The possible options are:

- *Inheritance_option_share* — indicates the region should be shared by the parent and child
- *Inheritance_option_copy* — indicates the region should be shared by the parent and child until one of them writes to the region; once a modification occurs, a copy-on-write is performed
- *Inheritance_option_none* — indicates the region should not be made accessible to the child

These values comprise the elements of the type *INHERITANCE_OPTION*.

The expression *inheritance(task, page_index)* indicate the inheritance option associated with the region indicated by *page_index* in *task*'s address space.

Mach Definition 59

$$INHERITANCE_OPTION ::= Inheritance_option_share \mid Inheritance_option_copy \\ \mid Inheritance_option_none$$

$\underline{inheritance} : TASK \times PAGE_INDEX \mapsto INHERITANCE_OPTION$

4.9.6 Shadow Memories

A memory, *memory₁*, is said to back a second memory, *memory₂*, if *memory₁*'s manager takes responsibility for pages of *memory₂* that are not handled by *memory₂*'s manager. The relation *backing_rel* indicates when *memory₁* backs *memory₂* at a given offset within *memory₁*. Each memory is backed by at most one memory-offset pair. Furthermore, a memory may back at most one other memory. For convenience, *backing_memory(memory)* and *backing_offset(memory)* are used to denote, respectively, the memory and offset backing *memory*.

Whenever *memory₁* backs *memory₂*, *memory₂* is said to shadow *memory₁*. For convenience:

- The expression *shadow_memories(memory)* indicates the singleton set of memories backed by *memory*. *shadow_memories* is defined only for those memories that back another memory.
- The expression *backing_chain(memory)* indicates the sequence of memories backing *memory*.

If a memory is not backed by any memories, then its backing chain is empty. If *memory₁* is backed by *memory₂* then the backing chain for *memory₁* consists of *memory₂* followed by the backing chain for *memory₂*. For example, suppose *memory₂* backs *memory₁*, *memory₃* backs *memory₂*, and no memory backs *memory₃*. Then, the backing chains for *memory₃*, *memory₂*, and *memory₁* are, respectively, $\langle \rangle$, $\langle memory_3 \rangle$, and $\langle memory_2, memory_3 \rangle$. Mach does not permit cycles to occur in the sequence of memories backing a memory. Thus, we require that no memory be present in its backing chain.

Mach Definition 60

<p><i>ShadowMemories</i></p> <p>$\underline{backing_rel} : \mathbb{P}(MEMORY \times MEMORY \times OFFSET)$ $\underline{backing_memory} : MEMORY \leftrightarrow MEMORY$ $\underline{backing_offset} : MEMORY \leftrightarrow OFFSET$ $\underline{shadow_memories} : MEMORY \rightarrow \mathbb{P} MEMORY$ $\underline{backing_chain} : MEMORY \rightarrow \text{seq } MEMORY$</p> <p>$\forall memory_1, memory_2 : MEMORY ; offset : OFFSET$ $\bullet (memory_1, memory_2, offset) \in \underline{backing_rel}$ $\Leftrightarrow ((memory_2, memory_1) \in \underline{backing_memory}$ $\quad \wedge (memory_2, offset) \in \underline{backing_offset})$ $\text{dom } \underline{shadow_memories} = \text{ran } \underline{backing_memory}$ $\forall memory_1 : MEMORY \mid memory_1 \in \text{dom } \underline{shadow_memories}$ $\bullet \underline{shadow_memories}(memory_1)$ $= \{ memory_2 : MEMORY$ $\quad \mid (\exists offset : OFFSET \bullet (memory_1, memory_2, offset) \in \underline{backing_rel}) \}$ $\forall memory_1 : MEMORY \mid memory_1 \in \text{dom } \underline{shadow_memories}$ $\bullet \#(\underline{shadow_memories}(memory_1)) = 1$ $\forall memory : MEMORY$ $\bullet memory \notin \text{dom } \underline{backing_memory} \Rightarrow \#(\underline{backing_chain}(memory)) = 0$ $\wedge (memory \in \text{dom } \underline{backing_memory}$ $\quad \Rightarrow \underline{backing_chain}(memory)$ $\quad = \langle \underline{backing_memory}(memory) \rangle$ $\quad \quad \wedge \underline{backing_chain}(\underline{backing_memory}(memory)) \rangle$ $\forall memory : MEMORY \bullet memory \notin \text{ran}(\underline{backing_chain}(memory))$</p>

4.9.7 Page Wiring

To prevent critical pages from being evicted, Mach allows tasks to wire pages. For each page allocated in a task, a count is maintained of the number of times that the task has wired the page. The expression $\underline{wire_count}(task, page_index)$ denotes the number of times that $task$ has wired the page indicated by $page_index$ in its address space. As long as a task's count for $page_index$ remains nonzero, the physical page associated with $page_index$ must be retained in memory. In other words, a physical page may only be evicted when no task has the page wired. The set \underline{wired} denotes the set of physical pages that are wired by some task.

Mach Definition 61

<p><i>Wired</i></p> <p><i>AddressSpace</i> <i>PageAndMemory</i> $\underline{wire_count} : (TASK \times PAGE_INDEX) \leftrightarrow \mathbb{N}$ $\underline{wired_locations} : \mathbb{P}(TASK \times PAGE_INDEX)$ $\underline{wired} : \mathbb{P} PAGE$</p> <p>$\text{dom } \underline{wire_count} = \underline{allocated}$ $\underline{wired_locations} = \{ task : TASK ; page_index : PAGE_INDEX$ $\quad \mid \underline{wire_count}(task, page_index) > 0 \}$ $\underline{wired_locations} \subseteq \text{dom}(\underline{representing_page} \circ \underline{map_rel})$ $\underline{wired} = (\underline{representing_page} \circ \underline{map_rel}) \downarrow (\underline{wired_locations})$</p>

Review Note:

The *wire-count* component corresponds to the VM entry wire count. A page is wired if any VM entry that is mapped to it is wired. For efficiency the prototype maintains two wire counts, one on VM entries and another on pages. The latter denotes the number of VM entries that have the page wired ignoring multiple wirings by a single VM entry. We do not model the page wire count.

4.9.8 Summary

The memory system is comprised of memory objects, address spaces, pages, and backing chains.

Mach Definition 62

<i>MemorySystem</i> <i>Memory</i> <i>AddressSpace</i> <i>PageAndMemory</i> <i>MachProtection</i> <i>Lock</i> <i>ShadowMemories</i> <i>Inheritance</i> <i>Wired</i>
$\underline{allocated} = \text{dom } \underline{mach_protection}$ $\text{ran } \underline{represented_memory} \subseteq \text{dom } \underline{object_port}$ $\forall \text{task_va_pair} : \text{TASK} \times \text{PAGE_INDEX}$ $ \text{task_va_pair} \in \text{dom } \underline{map_rel}$ $\wedge \underline{map_rel}(\text{task_va_pair}) \in \text{dom } \underline{representing_page}$ $\bullet \underline{mach_protection}(\text{task_va_pair})$ $\subseteq \text{PROTECTION} \setminus \text{page_locks}(\underline{representing_page}(\underline{map_rel}(\text{task_va_pair})))$ $\text{dom } \underline{inheritance} = \underline{allocated}$ $\underline{mapped} \subseteq \text{dom } \underline{object_port}$

4.10 Messages

This section discusses the structure of messages.

4.10.1 Message Options

The type *MACH_MSG_OPTION* denotes the base values of the *options* parameter of **mach_msg**. The recognized values of this type are *Mach_send_msg*, *Mach_rcv_msg*, *Mach_send_cancel*, *Mach_send_notify*, *Mach_rcv_notify*, *Mach_rcv_large*, *Mach_send_timeout*, and *Mach_rcv_timeout*. The *options* parameter is set to some set of the base values.

Mach Definition 63

[*MACH_MSG_OPTION*]

<i>Mach_send_msg</i> : <i>MACH_MSG_OPTION</i> <i>Mach_rcv_msg</i> : <i>MACH_MSG_OPTION</i> <i>Mach_send_cancel</i> : <i>MACH_MSG_OPTION</i> <i>Mach_send_notify</i> : <i>MACH_MSG_OPTION</i> <i>Mach_rcv_notify</i> : <i>MACH_MSG_OPTION</i> <i>Mach_rcv_large</i> : <i>MACH_MSG_OPTION</i> <i>Mach_send_timeout</i> : <i>MACH_MSG_OPTION</i> <i>Mach_rcv_timeout</i> : <i>MACH_MSG_OPTION</i>
$\text{disjoint} \{ \{ \textit{Mach_send_msg} \}, \{ \textit{Mach_rcv_msg} \},$ $\{ \textit{Mach_send_cancel} \}, \{ \textit{Mach_rcv_large} \}, \{ \textit{Mach_send_notify} \},$ $\{ \textit{Mach_rcv_notify} \}, \{ \textit{Mach_send_timeout} \}, \{ \textit{Mach_rcv_timeout} \} \}$

4.10.2 Complex Messages

In addition to simply carrying data, a message can also carry port rights and memory regions. A message carrying port rights or memory regions is called a *complex* message. Each message carries a flag indicating whether the message contains port rights or memory regions. The type *COMPLEX_OPTION* consists of the elements *Co_carries_rights* and *Co_carries_memory*; the flag carried in each message is a set of these values. Note that a flag containing both elements indicates that the message contains both port rights and memory regions.

Mach Definition 64

[*COMPLEX_OPTION*]

<i>Co_carries_rights</i> : <i>COMPLEX_OPTION</i> <i>Co_carries_memory</i> : <i>COMPLEX_OPTION</i>
$\text{disjoint} \{ \{ \textit{Co_carries_rights} \}, \{ \textit{Co_carries_memory} \} \}$

4.10.3 Data Types

Each element in the body of a message is typed. The set *MACH_MSG_TYPE* denotes the set of data types recognized by the system.

Mach Definition 65

[*MACH_MSG_TYPE*]

Whenever a port right is sent in a message, the client indicates a transfer option for the port right. The collection of acceptable transfer options is denoted by *Recognized_transfer_options* and contain the values *Mmt_make_send*, *Mmt_copy_send*, *Mmt_move_send*, *Mmt_make_send_once*, *Mmt_move_send_once*, and *Mmt_move_receive*.

An element of type *Mmt_make_send* indicates a receive right held by the sender from which a send right is to be created for the receiver. Similarly, an element of type *Mmt_make_send_once* indicates a receive right held by the sender from which a send-once right is to be created for the receiver.

An element of type *Mmt_copy_send* indicates a send right that should be copied from the sender's port name space into the receiver's port name space. In other words, the sender retains the existing port right while passing the right to the receiver.

An element of type *Mmt_move_send* indicates a send right that should be moved from the sender's port name space into the receiver's port name space. In other words, the sender's reference count is decremented by one and the receiver's reference count is incremented by one. If the sender's reference count was one, then the sender loses the capability associated with the right. If the receiver's reference count was zero, then the receiver gains the capability associated with the right. Similarly, *Mmt_move_send_once* and *Mmt_move_receive* allow send-once and receive rights to be moved from the sender's name space to the receiver's name space.

Mach Definition 66

<pre> Mmt_make_send : MACH_MSG_TYPE Mmt_copy_send : MACH_MSG_TYPE Mmt_move_send : MACH_MSG_TYPE Mmt_make_send_once : MACH_MSG_TYPE Mmt_move_send_once : MACH_MSG_TYPE Mmt_move_receive : MACH_MSG_TYPE Recognized_transfer_options : P MACH_MSG_TYPE </pre>
<pre> ({ Mmt_make_send }, { Mmt_copy_send }, { Mmt_move_send }, { Mmt_make_send_once }, { Mmt_move_send_once }, { Mmt_move_receive }) partition Recognized_transfer_options </pre>

After the kernel translates the port rights to an internal representation, it is no longer relevant whether the right was moved, copied or made and the kernel simply records the type of right, *Mach_msg_type_port_receive*, *Mach_msg_type_port_send*, or *Mach_msg_type_port_send_once*. These values of *MACH_MSG_TYPE* comprise the set *Mach_msg_type_port_rights*.

Mach Definition 67

<pre> Mach_msg_type_port_receive : MACH_MSG_TYPE Mach_msg_type_port_send : MACH_MSG_TYPE Mach_msg_type_port_send_once : MACH_MSG_TYPE Mach_msg_type_port_rights : P MACH_MSG_TYPE </pre>
<pre> ({ Mach_msg_type_port_receive }, { Mach_msg_type_port_send }, { Mach_msg_type_port_send_once }) partition Mach_msg_type_port_rights </pre>

4.10.4 Message Headers

The header for a message residing in user-space memory or kernel-space memory contains the following data:

- *local_port* — specifies the reply port when sending a message (*Mach_port_null* indicates no reply port is specified)
- *local_rights* — the port rights for the local port (if one is specified)
- *remote_port* — specifies the destination port when sending a message
- *remote_rights* — the port rights for the remote port
- *size* — specifies the size, in bytes, of a message when receiving
- *msg_sequence_no* — specifies the sequence number when receiving a message
- *operation* — operation or function id set by message sender

In addition, a message header in kernel space contains a value *complex* which indicates whether the message carries port rights or memory regions or both. This

value is a set of elements of type *COMPLEX_OPTION*. In place of *complex*, a message header in user space contains a single value *complex_boolean* indicating whether the message carries port rights and/or memory regions. The possible values are *Co_carries_rights_and_or_memory* and *Co_carries_neither_rights_nor_memory*. If *complex_boolean* has value *Co_carries_neither_rights_nor_memory*, then the message contains no port rights nor memory regions regardless of what is indicated by the individual data elements of the message.

Mach Definition 68

[*OPERATION*]

COMPLEX_OPTION_BOOLEAN
 $::=$ *Co_carries_rights_and_or_memory*
 \mid *Co_carries_neither_rights_nor_memory*

<p><i>MachMsgHeader</i></p> <p><i>local_port</i> : <i>NAME</i> <i>local_rights</i> : \mathbb{P} <i>MACH_MSG_TYPE</i> <i>remote_port</i> : <i>NAME</i> <i>remote_rights</i> : <i>MACH_MSG_TYPE</i> <i>size</i> : \mathbb{N} <i>msg_sequence_no</i> : \mathbb{N} <i>operation</i> : <i>OPERATION</i> <i>complex_boolean</i> : <i>COMPLEX_OPTION_BOOLEAN</i></p> <p>$\#local_rights \leq 1$</p>
--

Messages residing in kernel space contain ports rather than names. Thus, the *remote_port* and *local_port* fields contain ports instead of names when a message is in transit. If *Mach_port_null* was specified as the name of the local port in the *MachMsgHeader*, then *local_port* is empty in the corresponding *MachInternalHeader*.

Mach Definition 69

<p><i>MachInternalHeader</i></p> <p><i>local_port</i> : \mathbb{P} <i>PORT</i> <i>local_rights</i> : \mathbb{P} <i>MACH_MSG_TYPE</i> <i>remote_port</i> : <i>PORT</i> <i>remote_rights</i> : <i>MACH_MSG_TYPE</i> <i>size</i> : \mathbb{N} <i>msg_sequence_no</i> : \mathbb{N} <i>operation</i> : <i>OPERATION</i> <i>complex</i> : \mathbb{P} <i>COMPLEX_OPTION</i></p> <p>$\#local_rights = \#local_port \leq 1$</p>

4.10.5 Outcall Operations

There are several sets of operation identifiers used in messages to external servers (e.g., the security server) and user tasks. Some of these identifiers are used by the kernel when sending outcalls. We use

- *Exception_ids* to denote the set of operations used by the kernel when sending an exception message, The only element of this set is *Mach_exception_id*.
- *Kernel_service_reply_ids* to denote the set of operations used by the kernel in reply messages to kernel service requests,
- *Security_server_ids* to denote the set of security server operations,
- *Audit_ids* to denote the set of audit operations,
- *Mem_obj_confirmation_ids* to denote the set of operations used by the kernel when sending confirmations of memory operations to a pager,
- *Pager_request_ids* to denote the set of pager operations,
- *Mach_notify_ids* to denote the set of operations used by the kernel in notification messages, and
- *Network_packet_ids* to denote the set of operations used by the kernel when forwarding network packets.

We give a partial description of the identifiers in these sets.

Mach Definition 70

$$\begin{array}{|l} \hline \textit{Exception_ids} : \mathbb{P} \textit{ OPERATION} \\ \textit{Mach_exception_id} : \textit{ OPERATION} \\ \hline \textit{Exception_ids} = \{ \textit{Mach_exception_id} \} \end{array}$$

Mach Definition 71

$$\begin{array}{|l} \hline \textit{Kernel_service_reply_ids} : \mathbb{P} \textit{ OPERATION} \\ \hline \end{array}$$

Mach Definition 72

$$\begin{array}{|l} \hline \textit{Security_server_ids} : \mathbb{P} \textit{ OPERATION} \\ \textit{SSI_compute_av_id} : \\ \textit{ OPERATION} \\ \hline \{ \textit{SSI_compute_av_id} \} \\ \subseteq \textit{Security_server_ids} \end{array}$$

Mach Definition 73

$$\begin{array}{|l} \hline \textit{Audit_ids} : \mathbb{P} \textit{ OPERATION} \\ \textit{Audit_batch_id}, \textit{Audit_id} : \\ \textit{ OPERATION} \\ \hline \{ \textit{Audit_batch_id}, \textit{Audit_id} \} \\ \subseteq \textit{Audit_ids} \end{array}$$

Mach Definition 74

$$\begin{array}{|l} \hline \textit{Mem_obj_confirmation_ids} : \mathbb{P} \textit{ OPERATION} \\ \textit{Memory_object_change_completed_id}, \textit{Memory_object_lock_completed_id}, \\ \textit{Memory_object_supply_completed_id} : \\ \textit{ OPERATION} \\ \hline \{ \textit{Memory_object_change_completed_id}, \textit{Memory_object_lock_completed_id}, \\ \textit{Memory_object_supply_completed_id} \} \\ \subseteq \textit{Mem_obj_confirmation_ids} \end{array}$$

Mach Definition 75

<p><i>Pager_request_ids</i> : \mathbb{P} OPERATION <i>Memory_object_copy_id</i>, <i>Memory_object_create_id</i>, <i>Memory_object_data_initialize_id</i>, <i>Memory_object_data_request_id</i>, <i>Memory_object_data_return_id</i>, <i>Memory_object_data_unlock_id</i>, <i>Memory_object_data_write_id</i>, <i>Memory_object_init_id</i>, <i>Memory_object_terminate_id</i> :</p> <p style="text-align: center;">OPERATION</p>
<p>{ <i>Memory_object_copy_id</i>, <i>Memory_object_create_id</i>, <i>Memory_object_data_initialize_id</i>, <i>Memory_object_data_request_id</i>, <i>Memory_object_data_return_id</i>, <i>Memory_object_data_unlock_id</i>, <i>Memory_object_data_write_id</i>, <i>Memory_object_init_id</i>, <i>Memory_object_terminate_id</i> }</p> <p style="text-align: center;">\subseteq <i>Pager_request_ids</i></p>

Mach Definition 76

<p><i>Mach_notify_ids</i> : \mathbb{P} OPERATION <i>Ipc_notify_dead_name_id</i>, <i>Ipc_notify_msg_accepted_id</i>, <i>Ipc_notify_no_senders_id</i>, <i>Ipc_notify_port_deleted_id</i>, <i>Ipc_notify_port_destroyed_id</i>, <i>Ipc_notify_send_once_id</i> :</p> <p style="text-align: center;">OPERATION</p>
<p>{ <i>Ipc_notify_dead_name_id</i>, <i>Ipc_notify_msg_accepted_id</i>, <i>Ipc_notify_no_senders_id</i>, <i>Ipc_notify_port_deleted_id</i>, <i>Ipc_notify_port_destroyed_id</i>, <i>Ipc_notify_send_once_id</i> }</p> <p style="text-align: center;">\subseteq <i>Mach_notify_ids</i></p>

Mach Definition 77

<p><i>Network_packet_ids</i> : \mathbb{P} OPERATION <i>Forward_net_packet_id</i> :</p> <p style="text-align: center;">OPERATION</p>
<p>{ <i>Forward_net_packet_id</i> }</p> <p style="text-align: center;">\subseteq <i>Network_packet_ids</i></p>

4.10.6 Message Bodies

The body of a message consists of a sequence of message elements. Each element contains the following:

- the number of data elements contained in the message element
- a data type
- a collection of data elements or a single address

A triple that contains a collection of data elements represents in-line data. The number of data elements in the collection is the same as the specified number of data elements, and each such element is of the specified type. A triple that contains a single address represents out-of-line data. The address specifies the start of the area of memory containing the data. The data in that area is interpreted as being a collection of the specified number of data elements of the specified data type. Each out-of-line element contains a flag indicating whether the memory should be deallocated from the sender's address space. The possible values of this flag are *Msg_deallocate* and *Msg_dont_deallocate*.

Mach Definition 78

[MSG_DATA]

$$\begin{aligned}
 OLSD &::= \text{Msg_deallocate} \mid \text{Msg_dont_deallocate} \\
 \text{BASE_MSG_ELEMENT} \\
 &::= \text{In_line} \langle \mathbb{N} \times \text{MACH_MSG_TYPE} \times \text{seq MSG_DATA} \rangle \\
 &\quad \mid \text{Out_of_line} \langle \mathbb{N} \times \text{MACH_MSG_TYPE} \times \text{VIRTUAL_ADDRESS} \times OLSD \rangle
 \end{aligned}$$

Thus, an in-line message element is denoted by:

$$\text{In_line}(n, \text{mach_msg_type}, \text{data_seq})$$

and an out-of-line message element is denoted by:

$$\text{Out_of_line}(n, \text{mach_msg_type}, \text{va}, \text{olsd})$$

The number of entries specified in a triple representing in-line data must be the same as the number of entries in the specified sequence of data elements. The set *Msg_element* denotes the set of valid message elements, and the set *MESSAGE_BODY* denotes the set of sequences of valid message elements. In other words, *MESSAGE_BODY* denotes the set of valid message bodies.

Mach Definition 79

$ \begin{aligned} &\text{Msg_element} : \mathbb{P} \text{BASE_MSG_ELEMENT} \\ \hline &\text{Msg_element} \\ &= \{ \text{msg_element} : \text{BASE_MSG_ELEMENT} \\ &\quad \mid (\exists n : \mathbb{N}; \text{mach_msg_type} : \text{MACH_MSG_TYPE}; \text{data_seq} : \text{seq MSG_DATA}; \\ &\quad \text{va} : \text{VIRTUAL_ADDRESS}; \text{olsd} : OLSD \\ &\quad \bullet (\text{msg_element} = \text{In_line}(n, \text{mach_msg_type}, \text{data_seq}) \\ &\quad \quad \wedge \# \text{data_seq} = n) \\ &\quad \vee \text{msg_element} = \text{Out_of_line}(n, \text{mach_msg_type}, \text{va}, \text{olsd}) \} \end{aligned} $

Mach Definition 80

$$\text{MESSAGE_BODY} == \text{seq Msg_element}$$

When a message is moved into kernel space, the port names appearing in the message are transformed into port identifiers and the virtual addresses indicating out-of-line data are transformed into memory-offset pairs. In other words, the client specific names for kernel entities are transformed into the appropriate global names used internal to the kernel. Thus, an element in a message body in kernel space is of one of the following forms:

- *Msg_value*(*n*, *mach_msg_type*, (*task*, *value_seq*)) — an in-line element; if *mach_msg_type* is an element of *Recognized_transfer_options* and some elements of *value_seq* have not yet been resolved to ports then further processing is required to transform the sequence of data into a sequence of ports.

Note that there are two forms for elements of *value_seq*. An entry of the form *V_data*(*msg_data*, *v_data_l*) denotes the data *msg_data* while an entry of the form

$V_port(port, v_data_l)$ denotes a port name that has been resolved into a port. In either case, v_data_l indicates whether the element came from an in-line data element or an out-of-line data element. The only time v_data_l will indicate an out-of-line data element is when the element is a port name from an out-of-line data element that has been resolved into a port.

- $Transit_right(n, mach_msg_type, (task, port_seq, v_data_l))$ — a sequence of port rights in transit; $task$ indicates the task that sent the message and v_data_l indicates whether the port right was sent in-line or out-of-line
- $Msg_region(n, mach_msg_type, (task, va, olsd))$ — an out-of-line element that requires further processing to transform the task-address pair into a memory-offset pair; $task$ indicates the task that sent the message and $olsd$ indicates whether the region should be deallocated from $task$'s address space
- $Transit_memory(n, mach_msg_type, (task, memory, offset))$ — an out-of-line element that has been transformed from a task-address pair to a memory-offset pair; $task$ indicates the task that sent the message

The number of entries specified in a triple representing in-line data must be the same as the number of entries in the specified sequence of data elements. The type $Internal_element$ denotes the set of valid message elements internal to the kernel, and the type $INTERNAL_BODY$ denotes the set of sequences of these elements. Thus, $INTERNAL_BODY$ denotes the set of message bodies that can be stored in the kernel.

Mach Definition 81

$$\begin{aligned}
 V_DATA_LOCATION &::= V_data_in \mid V_data_out \\
 MSG_VALUE &::= V_data\langle\langle MSG_DATA \times V_DATA_LOCATION \rangle\rangle \\
 &\quad \mid V_port\langle\langle PORT \times V_DATA_LOCATION \rangle\rangle \\
 BASE_INTERNAL_ELEMENT \\
 &::= Msg_value\langle\langle \mathbb{N} \times MACH_MSG_TYPE \times (TASK \times seq\ MSG_VALUE) \rangle\rangle \\
 &\quad \mid Transit_right\langle\langle \mathbb{N} \times MACH_MSG_TYPE \\
 &\quad \quad \times (TASK \times seq\ PORT \times V_DATA_LOCATION) \rangle\rangle \\
 &\quad \mid Msg_region\langle\langle \mathbb{N} \times MACH_MSG_TYPE \times (TASK \times VIRTUAL_ADDRESS \times OLSD) \rangle\rangle \\
 &\quad \mid Transit_memory\langle\langle \mathbb{N} \times MACH_MSG_TYPE \times (TASK \times MEMORY \times OFFSET) \rangle\rangle
 \end{aligned}$$

Editorial Note:

$Transit_right$ probably needs to be considered in the following.

$$Internal_element : \mathbb{P}\ BASE_INTERNAL_ELEMENT$$

$$Internal_element$$

$$\begin{aligned}
 = \{ &msg_element : BASE_INTERNAL_ELEMENT \\
 &\mid (\exists n : \mathbb{N}; mach_msg_type : MACH_MSG_TYPE; task : TASK; \\
 &\quad value_seq : seq\ MSG_VALUE; port_seq : seq\ PORT; \\
 &\quad memory : MEMORY; offset : OFFSET; va : VIRTUAL_ADDRESS; \\
 &\quad olsd : OLSD; v_data_l : V_DATA_LOCATION \\
 &\quad \bullet (msg_element = Msg_value(n, mach_msg_type, (task, value_seq)) \\
 &\quad \quad \wedge \#value_seq = n) \\
 &\quad \vee msg_element = Msg_region(n, mach_msg_type, (task, va, olsd)) \\
 &\quad \vee msg_element \\
 &\quad = Transit_memory(n, mach_msg_type, (task, memory, offset))\}
 \end{aligned}$$

```

INTERNAL_BODY == {body : seq Internal_element
  | (∃ task : TASK
    • (∀ n : ℕ; mach_msg_type : MACH_MSG_TYPE;
      value_seq : seq MSG_VALUE;
      old : OLSD; task1 : TASK; va : VIRTUAL_ADDRESS
      | Msg_value(n, mach_msg_type, (task1, value_seq)) ∈ ran body
      ∨ Msg_region(n, mach_msg_type, (task1, va, old)) ∈ ran body
      • task = task1))}

```

Review Note:
Should *Transit_memory* be added to the above?

Note that all of the elements in a single message body must contain the same task identifier. It is intended that this task identifier unambiguously defines the identity of the task that sent the message.

4.10.7 Message Status

Once a message enters the kernel, it can be in one of three states:

- *Msg_stat_send* — indicates that the kernel is performing processing to send the message
- *Msg_stat_pseudo* — indicates that the kernel is performing processing to return the message to the message sender as part of a failed send request
- *Msg_stat_rcv* — indicates that the kernel is performing processing to receive the message

These elements comprise the values of the type *MSG_STATUS*.

The following error conditions can arise during the processing of a message: *Msg_error_invalid_memory*, *Msg_error_invalid_right*, *Msg_error_invalid_type*, *Msg_error_msg_too_small*, *Msg_error_notify_in_progress*, and *Msg_error_timed_out*. These values comprise the set *MSG_ERROR*.

Mach Definition 82

$$MSG_STATUS ::= Msg_stat_send \mid Msg_stat_pseudo \mid Msg_stat_rcv$$

$$MSG_ERROR ::= Msg_error_invalid_memory \mid Msg_error_invalid_right \\ \mid Msg_error_invalid_type \mid Msg_error_msg_too_small \\ \mid Msg_error_notify_in_progress \mid Msg_error_timed_out$$

4.10.8 Message Structure

Each message is modeled as containing fields *header* and *body*. The type *Message* denotes the set of user space messages.

Mach Definition 83

$Message$ <hr style="border: 0; border-top: 1px solid black; margin: 5px 0;"/> <i>header</i> : MachMsgHeader <i>body</i> : MESSAGE_BODY
--

In addition to the header and body, messages in transit also contain the following fields:

- *option* — indicates the options specified by the client
- *time_out_at* — indicates when a given send or receive request will time out
If the set contained in this field is empty, then the message will not time out. Otherwise, the set contains exactly one value and this value defines the earliest time at which the associated send or receive request can time out.
- *status* — indicates future processing the kernel must perform on the message
- *error* — indicates the first error (if any) that occurred during the processing of the message.

Editorial Note:

status and *error* should be removed as the purpose they were intended to serve can now be accomplished more generally using the tools of the execution model.

The type *InternalMessage* denotes the possible values of messages in transit.

Mach Definition 84

<pre> InternalMessage header : MachInternalHeader body : INTERNAL_BODY option : P MACH_MSG_OPTION time_out_at : P N status : MSG_STATUS error : P MSG_ERROR #time_out_at ≤ 1 #error ≤ 1 </pre>

4.10.9 Pending Receives

Each port can have clients blocked on message receive requests waiting for messages to arrive at the port. Each pending receive request has the following associated information:

- *notify* — the notify port name specified by the receiving task
- *option* — the options specified by the receiving task
- *rcv_size* — the receive size specified by the receiving task
- *time_out_at* — the time at which the request will time out; this has the same format as the *time_out_at* component of *InternalMessage*.

Mach Definition 85

<pre> PendingReceive notify : NAME option : P MACH_MSG_OPTION rcv_size : N time_out_at : P N #time_out_at ≤ 1 </pre>

4.10.10 Reply Ports

The sender of a message can specify a reply port for the receiver to use to reply to the message. The sender does so by setting the *local_port* field to its name for the port. For convenience, the relation *reply_port_rel* is used to denote the reply port and transferred right in a message specifying a reply port. The interpretation of:

$$(message, (port, right))$$

being an element of *reply_port_rel* is that *message* transfers the type of right specified by *right* (send or send-once) for *port* to the receiver of *message*. The intent is that the receiver use the transferred right to send a reply message to *port*. Each message contains at most one reply port and right for that port. For convenience, the expressions *reply_port(message)* and *reply_port_right(message)* are used to denote the reply port and transferred right contained in a given message.

Mach Definition 86

<i>ReplyPorts</i>
$\underline{reply_port_rel} : MESSAGE \leftrightarrow (PORT \times \{Send, Send_once\})$
$reply_port : MESSAGE \mapsto PORT$
$reply_port_right : MESSAGE \mapsto \{Send, Send_once\}$
$reply_port = \{ message : MESSAGE; port : PORT; right : RIGHT$ $\quad (message, (port, right)) \in \underline{reply_port_rel} \bullet (message, port) \}$
$reply_port_right = \{ message : MESSAGE; port : PORT; right : RIGHT$ $\quad (message, (port, right)) \in \underline{reply_port_rel} \bullet (message, right) \}$

4.10.11 Summary

This section has defined the data structures used to model messages. The expression *msg_contents(message)* is used to denote the internal message structure associated with each message identifier, and the expression *pending_receives(task, name)* indicates the receive requests currently pending for threads in *task* that attempted to receive through the port named by *name*. The expression *task_received_msgs(task)* denotes the set of user-space messages that have been received by *task*.

For convenience, the expression *msg_operation(message)* is used to denote the type of operation requested by *message*. In other words, the returned value is the *operation* field of the message identified by *message*.

Mach Definition 87

<i>Operations</i>
$msg_operation : MESSAGE \mapsto OPERATION$

Mach Definition 88

<p><i>Messages</i></p> <p><i>TaskExist</i></p> <p><i>MessageExist</i></p> <p><i>Operations</i></p> <p><i>ReplyPorts</i></p> <p>$\underline{msg_contents} : MESSAGE \leftrightarrow InternalMessage$</p> <p>$\underline{pending_receives} : TASK \times NAME \leftrightarrow seq PendingReceive$</p> <p>$\underline{task_received_msgs} : TASK \rightarrow \mathbb{P} MESSAGE$</p> <hr/> <p>$dom \underline{reply_port} \subseteq \underline{message_exists}$</p> <p>$dom \underline{msg_operation} = dom \underline{msg_contents} = \underline{message_exists}$</p> <p>$\forall message : MESSAGE \mid message \in \underline{message_exists}$</p> <ul style="list-style-type: none"> $\bullet \underline{msg_operation}(message) = (\underline{msg_contents}(message)).header.operation$ <p>$\forall message : MESSAGE; port : PORT \mid message \in \underline{message_exists}$</p> <ul style="list-style-type: none"> $\bullet (message, (port, Send)) \in \underline{reply_port_rel}$ <p>$\Leftrightarrow ((\underline{msg_contents}(message)).header.local_port = \{port\})$ $\quad \wedge (\underline{msg_contents}(message)).header.local_rights$ $\quad \cap \{ Mmt_make_send, Mmt_move_send, Mmt_copy_send \} \neq \emptyset)$</p> <p>$\forall message : MESSAGE; port : PORT \mid message \in \underline{message_exists}$</p> <ul style="list-style-type: none"> $\bullet (message, (port, Send_once)) \in \underline{reply_port_rel}$ <p>$\Leftrightarrow ((\underline{msg_contents}(message)).header.local_port = \{port\})$ $\quad \wedge (\underline{msg_contents}(message)).header.local_rights$ $\quad \cap \{ Mmt_make_send_once, Mmt_move_send_once \} \neq \emptyset)$</p> <p>$\forall task : TASK$</p> <ul style="list-style-type: none"> $\mid task \notin \underline{task_exists}$ $\bullet \underline{task_received_msgs}(task) = \emptyset$
--

Review Note:
Must figure out what the axioms are on $\underline{pending_receives}$.

4.11 Processors and Processor Sets

Each host has a default processor set denoted by $\underline{default}$. Furthermore, each host has a master processor denoted by $\underline{master_proc}$.

Mach Definition 89

<p><i>HostsAndProcessors</i></p> <p><i>ProcessorsAndPorts</i></p> <p>$\underline{default} : PROCESSOR_SET$</p> <p>$\underline{master_proc} : PROCESSOR$</p> <hr/> <p>$\underline{default} \in dom \underline{ps_control_port_rel}$</p> <p>$\underline{master_proc} \in dom \underline{processor_port_rel}$</p>
--

Each processor is a member of a single processor set. The relation $\underline{member_rel}$ indicates which processors belong to each processor set. For convenience, the expressions $\underline{processors}(procset)$ and $\underline{proc_assigned_procset}(proc)$ are used to denote, respectively, the set of processors that belong to $procset$ and the processor set to which $proc$ belongs.

Mach Definition 90

<i>ProcessorAndProcessorSet</i> <i>ProcessorsAndPorts</i> $\underline{m}em_rel : PROCESSOR \leftrightarrow PROCESSOR_SET$ $\underline{p}rocessors : PROCESSOR_SET \rightarrow \mathbb{P} PROCESSOR$ $\underline{p}roc_assigned_procset : PROCESSOR \leftrightarrow PROCESSOR_SET$
$\text{dom } \underline{m}em_rel \subseteq \text{dom } \underline{p}rocessor_port_rel$ $\text{ran } \underline{m}em_rel \subseteq \text{dom } \underline{p}s_control_port_rel$ $\underline{p}roc_assigned_procset = \underline{m}em_rel$ $\underline{p}rocessors = (\lambda procset : PROCESSOR_SET \bullet \underline{m}em_rel \sim (\{\{procset\}\}))$

Each task is assigned to a single processor set. The relation $\underline{t}ask_assignment_rel$ indicates the association between tasks and processor sets. For convenience, the expressions $\underline{h}ave_assigned_tasks(procset)$ and $\underline{t}ask_assigned_to(task)$ are used to denote, respectively, the set of tasks assigned to $procset$ and processor set to which $task$ is assigned.

Mach Definition 91

<i>TaskAndProcessorSet</i> <i>SpecialTaskPorts</i> <i>ProcessorsAndPorts</i> $\underline{t}ask_assignment_rel : TASK \leftrightarrow PROCESSOR_SET$ $\underline{h}ave_assigned_tasks : PROCESSOR_SET \rightarrow \mathbb{P} TASK$ $\underline{t}ask_assigned_to : TASK \leftrightarrow PROCESSOR_SET$
$\text{dom } \underline{t}ask_assignment_rel = \text{ran } \underline{self_task}$ $\text{ran } \underline{t}ask_assignment_rel \subseteq \text{dom } \underline{p}s_control_port_rel$ $\underline{t}ask_assignment_rel = \underline{t}ask_assigned_to$ $\underline{h}ave_assigned_tasks = (\lambda procset : PROCESSOR_SET$ $\bullet \underline{t}ask_assignment_rel \sim (\{\{procset\}\}))$

Similarly, Each thread is assigned to a single processor set. The relation $\underline{t}hread_assignment_rel$ associates threads with processor sets. For convenience, the expressions $\underline{h}ave_assigned_threads(procset)$ and $\underline{t}hread_assigned_to(thread)$ are used to denote, respectively, the set of threads assigned to $procset$ and processor set to which $thread$ is assigned.

Each processor set has a set of enabled scheduling policies, denoted by $\underline{e}nabled_sp(procset)$ and a maximum priority for assigned threads, denoted by $\underline{p}s_max_priority(procset)$. The set of enabled scheduling policies for a thread's processor set is used to constrain the policies that can be assigned to that thread. The maximum scheduling priority for a processor set constrains the priorities that can be assigned to a newly created thread associated with that processor set.

Mach Definition 92

<p><i>ThreadAndProcessorSet</i></p> <p><i>ProcessorSetExist</i></p> <p><i>ProcessorsAndPorts</i></p> <p><i>SpecialThreadPorts</i></p> <p><i>ThreadSchedPolicy</i></p> <p>$\underline{t}hread_assignment_rel : THREAD \leftrightarrow PROCESSOR_SET$</p> <p>$have_assigned_threads : PROCESSOR_SET \rightarrow \mathbb{P} \text{ THREAD}$</p> <p>$\underline{t}hread_assigned_to : THREAD \leftrightarrow PROCESSOR_SET$</p> <p>$\underline{e}nabled_sp : PROCESSOR_SET \leftrightarrow \mathbb{P} \text{ SCHED_POLICY}$</p> <p>$\underline{p}s_max_priority : PROCESSOR_SET \leftrightarrow \mathbb{Z}$</p> <hr/> <p>$\underline{t}hread_assignment_rel = \underline{t}hread_assigned_to$</p> <p>$have_assigned_threads = (\lambda \text{procset} : PROCESSOR_SET$</p> <ul style="list-style-type: none"> • $\underline{t}hread_assignment_rel \sim (\{\{\text{procset}\}\})$) <p>$\text{dom } \underline{t}hread_assignment_rel = \text{dom } \underline{t}hread_self$</p> <p>$\text{ran } \underline{t}hread_assignment_rel \subseteq \text{dom } \underline{p}s_control_port_rel$</p> <p>$\text{dom } \underline{e}nabled_sp = \text{dom } \underline{p}s_max_priority = \underline{p}rocset_exists$</p> <p>$\bigcup (\text{ran } \underline{e}nabled_sp) \subseteq \underline{s}upported_sp$</p> <p>$\text{ran } \underline{p}s_max_priority \subseteq \underline{P}riority_levels$</p>
--

Each processor may have an active thread. The expression $\underline{a}ctive_thread(proc)$ indicates the thread (if any) that is active on $proc$.

Mach Definition 93

<p><i>ThreadsAndProcessors</i></p> <p><i>ThreadExist</i></p> <p><i>Exist</i></p> <p>$\underline{a}ctive_thread : PROCESSOR \rightarrow \text{ THREAD}$</p> <hr/> <p>$\text{dom } \underline{a}ctive_thread \subseteq \underline{p}roc_exists$</p> <p>$\text{ran } \underline{a}ctive_thread \subseteq \underline{t}hread_exists$</p>
--

4.12 Time

Each host provides a system clock. The current system time is denoted by $\underline{h}ost_time$.

Mach Definition 94

<p><i>HostTime</i></p> <p>$\underline{h}ost_time : \mathbb{N}$</p>
--

4.13 Devices

Each device has an associated count indicating how many times the device has been opened and not closed. We use $\underline{d}evice_open_count(dev)$ to indicate the count associated with dev . This count is incremented each time dev is opened and decremented each time dev is closed. Each device with a positive creation count has an associated device port that represents the device.

Mach Definition 95

<i>DeviceOpenCount</i> <i>DevicesAndPorts</i> $\underline{device_open_count} : DEVICE \rightarrow \mathbb{N}$ $\text{dom } \underline{device_port} = \{ dev : DEVICE \mid \underline{device_open_count}(dev) > 0 \}$
--

A kernel-space device driver may supply event counters for use by user-space device drivers. An event counter is used as a semaphore for events produced by kernel-space drivers. The counter is incremented when a relevant event occurs and decremented when a thread (e.g., a user-space device driver) indicates via the **evc_wait** trap that it wishes to process an event. Each task refers to an event by referencing its event counter. The appropriate event counter is communicated to a thread in a driver-specific way.⁷ The expression *EVENT_COUNTER* denotes the set of all event counters.

Mach Definition 96

[*EVENT_COUNTER*]

Each event counter may have at most one thread, denoted by $\underline{thread_waiting}(evc)$, waiting for it. Furthermore, each thread may be waiting on at most one event counter. The number of event that are queued and waiting to be processed by a thread is denoted by $\underline{event_count}(evc)$. The expression $\underline{supplying_device}$ denotes the kernel-space device driver that supplied the event counter.

Mach Definition 97

<i>Events</i> <i>ThreadExecStatus</i> $\underline{thread_waiting} : EVENT_COUNTER \rightarrow THREAD$ $\underline{event_count} : EVENT_COUNTER \rightarrow \mathbb{Z}$ $\underline{supplying_device} : EVENT_COUNTER \rightarrow DEVICE$ $\text{dom } \underline{event_count} = \text{dom } \underline{supplying_device}$ $\text{dom } \underline{thread_waiting} \subseteq \text{dom } \underline{event_count}$ $\text{ran } \underline{thread_waiting} \subseteq$ $\{ thread : THREAD \mid thread \in \underline{thread_exists} \wedge \text{Waiting} \in \underline{run_state}(thread) \}$
--

Devices can be associated with memory objects that can then be mapped into address spaces. We use $\underline{mapped_devices}$ to denote the set of devices that have been associated with memory objects.

Mach Definition 98

<i>MappedDevices</i> $\underline{mapped_devices} : \mathbb{P} DEVICE$

Each device has two associated queues of data records. We use $\underline{device_in}(dev)$ and $\underline{device_out}(dev)$ to denote, respectively, data input and output through the device. Data read from *dev* is dequeued from $\underline{device_in}(dev)$, and data written to *dev* is enqueued to $\underline{device_out}(dev)$.

⁷Threads may also wait for events that occur while the system is operating in kernel space (e.g., another thread becomes suspended). This is handled through a separate waiting mechanism that is not modeled in the FTLS.

Mach Definition 99

[*DEVICE_RECORD*]

Mach Definition 100

<i>DeviceData</i> <u><i>device_in</i></u> : <i>DEVICE</i> → seq <i>DEVICE_RECORD</i> <u><i>device_out</i></u> : <i>DEVICE</i> → seq <i>DEVICE_RECORD</i>
--

Each device can have associated filters that are used to route data received through the device. Each filter has an associated port to which data accepted by the filter is delivered. Furthermore, a priority can be associated with each port to indicate the ordering when there are multiple ports associated with the filter. We use *device_filter_info(dev)* to indicate the set of (*device_filter*, *port*, *filter_priority*) triples associated with *dev*.

Mach Definition 101

[*DEVICE_FILTER*, *FILTER_PRIORITY*]
DEVICE_FILTER_INFO == *DEVICE_FILTER* × *PORT* × *FILTER_PRIORITY*

Mach Definition 102

<i>DeviceFilterInfo</i> <u><i>device_filter_info</i></u> : <i>DEVICE</i> → P <i>DEVICE_FILTER_INFO</i>

Each device has an associated status. We use *device_status(dev)* to denote *dev*'s status.

Mach Definition 103

[*DEVICE_STATUS*]

Mach Definition 104

<i>DeviceStatus</i> <u><i>device_status</i></u> : <i>DEVICE</i> → <i>DEVICE_STATUS</i>

Mach Definition 105

<i>Devices</i> <i>DeviceOpenCount</i> <i>Events</i> <i>MappedDevices</i> <i>DeviceData</i> <i>DeviceFilterInfo</i> <i>DeviceStatus</i>
--

4.14 Summary

The data structures defined in the previous sections comprise the Mach system state. The type *Mach* is used to denote the set of Mach system states.

Mach Definition 106

$$\begin{aligned}
 \text{Process} &\hat{=} \text{Threads} \wedge \text{TaskPriority} \wedge \text{TaskSuspendCount} \\
 &\quad \wedge \text{EmulationVector} \\
 \text{Ipc} &\hat{=} \text{PortNameSpace} \wedge \text{RegisteredRights} \wedge \text{Notifications} \\
 &\quad \wedge \text{PortSummary} \wedge \text{PortClasses} \wedge \text{Messages} \\
 \text{Processor} &\hat{=} \text{HostsAndProcessors} \wedge \text{ProcessorAndProcessorSet} \wedge \text{TaskAndProcessorSet} \\
 &\quad \wedge \text{ThreadAndProcessorSet} \wedge \text{ThreadsAndProcessors}
 \end{aligned}$$
Mach Definition 107

<i>Mach</i> <i>Exist</i> <i>Process</i> <i>Ipc</i> <i>Processor</i> <i>MemorySystem</i> <i>HostTime</i> <i>Devices</i>
<u><i>manager</i></u> = <i>receiver</i> o <i>object_port</i>

Section 5

DTOS State Extensions

This section describes extensions to the base Mach microkernel state that are needed to support the DTOS kernel. The DTOS kernel is intended to support a wide range of policies. Thus, the state components described in this section are independent of any specific access control policy.

In general, an access control policy consists of three components. First, security attributes must be associated with the subjects accessing entities in the system. Second, security attributes must be associated with the entities in the system that subjects access. Finally, a rule must be defined that indicates the set of accesses that a subject with a given attribute can make to an entity with a given attribute. To provide policy flexibility, the DTOS kernel abstracts the security attributes associated with specific policies into sets of *security identifiers*. Although the kernel relies upon a security server to define the policy to be enforced, the kernel maintains a cache of accesses previously authorized by the security server.

In addition to providing a framework for access control policies, the DTOS kernel also enhances the security of the Mach IPC mechanism.

The organization of this section is as follows:

- Section 5.1, **Subject Security Information**, describes the security information recorded for subjects.
- Section 5.2, **Object Security Information**, describes the security information recorded for objects.
- Section 5.3, **Security Identifiers for Access Computations**, describes some security identifiers used only in access computations.
- Section 5.4, **Permissions**, describes the permissions enforced in DTOS.
- Section 5.5, **Access Vector Cache**, describes the DTOS kernel's access vector cache.
- Section 5.6, **Message Security Information**, describes the security information associated with messages to enhance the security of the Mach IPC mechanism.
- Section 5.7, **Task Creation Information**, describes information associated with tasks to enhance the security of the Mach approach for process initiation.
- Section 5.8, **Server Ports**, describes ports used by the kernel for communication with other servers.
- Section 5.9, **Memory Region Protections**, describes information associated with regions to allow the DTOS kernel to enforce access.

5.1 Subject Security Information

Subjects in DTOS are threads executing within tasks. Each task has a *subject security identifier* (SSI). The set *SSI* denotes the set of all SSIs.

We will occasionally need to identify two distinct components of each SID, a *mandatory* security identifier (MID) and an *authentication* identifier (AID). We use the types *MID* and *AID* to denote, respectively, MIDs and AIDs. The functions *Ssi_to_mid* and *Ssi_to_aid* are used to map SSIs to MIDs and AIDs.

DTOS Kernel Definition 1

[SSI]

[MID, AID]

$$\left| \begin{array}{l} Ssi_to_mid : SSI \rightarrow MID \\ Ssi_to_aid : SSI \rightarrow AID \end{array} \right.$$

The expressions $\underline{task_sid}(task)$, $task_mid(task)$ and $task_aid(task)$ are used to denote the SSI, MID and AID associated with a task. The expression $thread_sid(thread)$ denotes the SSI associated with a thread. It is defined to be the SSI of its parent task.

DTOS Kernel Definition 2

$\begin{array}{l} \underline{SubjectSid} \\ TaskExist \\ ThreadExist \\ TasksAndThreads \\ \underline{task_sid} : TASK \leftrightarrow SSI \\ task_mid : TASK \leftrightarrow MID \\ task_aid : TASK \leftrightarrow AID \\ thread_sid : THREAD \leftrightarrow SSI \\ \hline \text{dom } \underline{task_sid} = \text{dom } task_mid = \text{dom } task_aid = \underline{task_exists} \\ \text{dom } thread_sid = \underline{thread_exists} \\ task_mid = Ssi_to_mid \circ \underline{task_sid} \\ task_aid = Ssi_to_aid \circ \underline{task_sid} \\ thread_sid = \underline{task_sid} \circ \text{owning_task} \end{array}$
--

5.2 Object Security Information

Each port has an associated *object security identifier* (OSI) that represents the security attributes associated with the port. Similarly, each memory region has an associated OSI. The set *OSI* denotes the set of all OSIs.

The functions Osi_to_mid and Osi_to_aid are used to map OSIs to MIDs and AIDs.

DTOS Kernel Definition 3

[OSI]

$$\left| \begin{array}{l} Osi_to_mid : OSI \rightarrow MID \\ Osi_to_aid : OSI \rightarrow AID \end{array} \right.$$

The expressions $port_sid(port)$, $port_mid(port)$ and $port_aid(port)$ are used to denote the OSI, MID and AID associated with a port.

DTOS Kernel Definition 4

$\begin{aligned} & \underline{PortSid} \\ & \underline{PortExist} \\ & \underline{port_sid} : PORT \leftrightarrow OSI \\ & \underline{port_mid} : PORT \leftrightarrow MID \\ & \underline{port_aid} : PORT \leftrightarrow AID \end{aligned}$
$\begin{aligned} \text{dom } \underline{port_sid} &= \text{dom } \underline{port_mid} = \text{dom } \underline{port_aid} = \underline{port_exists} \\ \underline{port_mid} &= \underline{Osi_to_mid} \circ \underline{port_sid} \\ \underline{port_aid} &= \underline{Osi_to_aid} \circ \underline{port_sid} \end{aligned}$

Each task and thread has a self port on which the kernel receives requests to perform an action on the task or thread. The OSI of the self ports is derived from the SSI of the corresponding task. The expressions $\underline{Task_port_sid}(ssi)$ and $\underline{Thread_port_sid}(ssi)$ indicate the corresponding OSIs. When memory is allocated, it is labeled with an OSI that is derived from the SSI of the owning task. The expression $\underline{Default_vm_port_sid}(ssi)$ indicates the derived OSI. Similarly, when a port is created, it is labeled with an OSI derived from the SSI of the task in whose IPC name space it is allocated. The expression $\underline{Default_port_sid}(ssi)$ indicates the derived OSI.

DTOS Kernel Definition 5

$\begin{aligned} & \underline{Task_port_sid} : SSI \rightarrow OSI \\ & \underline{Thread_port_sid} : SSI \rightarrow OSI \\ & \underline{Default_vm_port_sid} : SSI \rightarrow OSI \\ & \underline{Default_port_sid} : SSI \rightarrow OSI \end{aligned}$
$\text{disjoint } \langle \text{ran } \underline{Task_port_sid}, \text{ran } \underline{Thread_port_sid}, \text{ran } \underline{Default_vm_port_sid}, \text{ran } \underline{Default_port_sid} \rangle$

$\begin{aligned} & \underline{KernelPortSid} \\ & \underline{TasksAndThreads} \\ & \underline{SpecialPurposePorts} \\ & \underline{SubjectSid} \\ & \underline{PortSid} \end{aligned}$
$\begin{aligned} \forall \text{ task} : \underline{task_exists} \\ & \bullet \underline{port_sid}(\underline{task_self}(\text{task})) = \underline{Task_port_sid}(\underline{task_sid}(\text{task})) \\ \forall \text{ thread} : \underline{thread_exists} \\ & \bullet \underline{port_sid}(\underline{thread_self}(\text{thread})) = \underline{Thread_port_sid}(\underline{task_sid}(\underline{owning_task}(\text{thread}))) \end{aligned}$

The expressions $\underline{page_sid}(\text{task}, \text{page_index})$, $\underline{page_mid}(\text{task}, \text{page_index})$ and $\underline{page_aid}(\text{task}, \text{page_index})$ are used to denote the OSI, MID and AID associated with a page. Note that $\underline{page_sid}$ effectively associates an OSI with each allocated address in a task's address space. If a page is managed and the manager is not the default memory manager, then the SID of the page is derived from the SID of the pager port of the object containing the page. The derivation of page SIDs from pager port SIDs is modeled by the function $\underline{Pp_to_page_sid}$.

DTOS Kernel Definition 6

$\underline{Pp_to_page_sid} : OSI \leftrightarrow OSI$

DTOS Kernel Definition 7

$PageSid$ $AddressSpace$ $Memory$ $TasksAndPorts$ $PortSid$ $\underline{page_sid} : TASK \times PAGE_INDEX \leftrightarrow OSI$ $\underline{page_mid} : TASK \times PAGE_INDEX \leftrightarrow MID$ $\underline{page_aid} : TASK \times PAGE_INDEX \leftrightarrow AID$
$dom \underline{page_sid} = dom \underline{page_mid} = dom \underline{page_aid} = \underline{allocated}$ $\underline{page_mid} = \underline{Osi_to_mid} \circ \underline{page_sid}$ $\underline{page_aid} = \underline{Osi_to_aid} \circ \underline{page_sid}$
$(\forall \underline{task_va_pair} : TASK \times PAGE_INDEX; \underline{memory} : MEMORY;$ $\underline{port} : PORT$ $ (\underline{task_va_pair}, \underline{memory}) \in \underline{mapped_memory}$ $\wedge (\underline{memory}, \underline{port}) \in \underline{object_port}$ $\wedge \underline{receiver}(\underline{port}) \neq \underline{receiver}(\underline{default_mem_manager})$ $\bullet \underline{page_sid}(\underline{task_va_pair}) = \underline{Pp_to_page_sid}(\underline{port_sid}(\underline{port}))$

Editorial Note:

Need to figure out if there is a better way to check that the memory is not being paged by the default memory manager.

DTOS Kernel Definition 8

$ObjectSid$ $PortSid$ $KernelPortSid$ $PageSid$
$dom \underline{Pp_to_page_sid} \subseteq \text{ran } \underline{port_sid}$ $\text{ran } \underline{Pp_to_page_sid} \subseteq \text{ran } \underline{page_sid}$

5.3 Security Identifiers for Access Computations

Access computations in the DTOS kernel are generally made based upon the SSI of the task accessing an object and the OSI of the accessed object. This section discusses a few special cases in which other security identifiers are used.

Sometimes kernel requests can have side effects resulting in outcalls from the kernel, for instance, to deliver dead name notifications. For fine grained control over such operations it is desirable to distinguish between the kernel sending such a message to a port as a side effect of another request and the client directly sending a message to the port. To provide for this, such side effects are sometimes controlled based not upon the SSI of the client but upon an SSI derived from the client's SSI and indicating that it is the kernel acting on behalf of a client with the given SSI. The function $\underline{Derive_kernel_as}$ maps an SSI s_1 to the derived SSI s_2 representing the kernel acting on behalf of a task with SSI s_1 . We use $\underline{kernel_as}(task)$ to denote the derived SSI indicating the kernel acting on behalf of a task $task$.

DTOS Kernel Definition 9

| $Derive_kernel_as : SSI \rightarrow SSI$

$KernelAs$ $SubjectSid$ $kernel_as : TASK \rightarrow SSI$
$kernel_as = \underline{task_sid} ; Derive_kernel_as$

One of the features of Mach is that it allows tasks to perform operations on other tasks that have not traditionally been provided by operating systems. For example, Mach allows tasks to access memory regions in other tasks while one of the features of traditional operating systems is the separation of address spaces. To provide finer control over task accesses, we define $Task_self_sid$ to be a value to be used in access computations governing accesses a task makes to itself. Similarly, we use $Thread_self_sid$ to be a value to be used in access computations governing accesses a task makes to threads that it owns. The security policy should normally be defined in such a way as to prevent any kernel entities from being assigned $Task_self_sid$ or $Thread_self_sid$ as their SID.⁸ Instead, these SIDs indicate to security servers that the kernel requires an access computation to be performed between a task and the task itself or between a task and one of the task's threads. One potential use of this finer control would be to contain a faulty task by preventing it from corrupting other tasks having the same SID.

We define $task_target(task_1, task_2)$ to be the OSI of $task_2$'s self port if $task_1$ and $task_2$ are different and $Task_self_sid$, otherwise. Analogously, we define $thread_target(task, thread)$ to be the OSI of $thread$'s self port if $thread$ does not belong to $task$ and $Thread_self_sid$, otherwise. When $task_1$ attempts to operate on $task_2$, the kernel enforces accesses on the pair $(\underline{task_sid}(task_1), task_target(task_1, task_2))$. Analogously, operations that $task$ performs on $thread$ are governed by the accesses recorded for $(\underline{task_sid}(task), thread_target(task, thread))$. This allows separate permissions sets to be applied when a task operates on itself versus operating on another process with the same SSI.

DTOS Kernel Definition 10

$Task_self_sid : OSI$ $Thread_self_sid : OSI$
$Task_self_sid \neq Thread_self_sid$

DTOS Kernel Definition 11

⁸This property is not guaranteed by the kernel. For example, a **mach_port_allocate_secure** request may specify a self SID as the SID for the newly created port. If the security server allows the client to add a name to the target task and allows the target task to hold a receive right for a port with the specified SID, the request will succeed and the port will be labeled with a self SID.

<p><i>TargetSids</i></p> <p><i>PortSid</i></p> <p><i>TasksAndThreads</i></p> <p><i>SpecialPurposePorts</i></p> <p>$task_target : TASK \times TASK \leftrightarrow OSI$</p> <p>$thread_target : TASK \times THREAD \leftrightarrow OSI$</p> <p>$\{Task_self_sid, Thread_self_sid\} \cap \text{ran } port_sid = \emptyset$</p> <p>$\text{dom } task_target = TASK \times \underline{task_exists}$</p> <p>$\text{dom } thread_target = TASK \times \underline{thread_exists}$</p> <p>$\forall task_1, task_2 : TASK$</p> <ul style="list-style-type: none"> • $task_target(task_1, task_2)$ = if $task_1 = task_2$ then $Task_self_sid$ else $port_sid(task_self(task_2))$ <p>$\forall task : TASK; thread : THREAD$</p> <ul style="list-style-type: none"> • $thread_target(task, thread)$ = if $task = \text{owning_task}(thread)$ then $Thread_self_sid$ else $port_sid(thread_self(thread))$
--

Editorial Note:

In the prototype $Task_self_sid$ and $Thread_self_sid$ are not implemented as constants. Rather, they are derived from the corresponding subject SID in the same way as the derived SIDs $Task_port_sid$, $Thread_port_sid$, $Default_vm_port_sid$ and $Default_port_sid$ which are described above. Given the way the self SIDs are used the two approaches are equivalent.

5.4 Permissions

The DTOS security policy constrains when clients may obtain *services*. The security policy is enforced by:

- associating a set of allowed permissions⁹ with each SSI-OSI pair,
- associating a set of required permissions with each service, and
- granting service only when the required permissions are contained in the allowed permissions for the client to the target for the operation.

The set *PERMISSION* denotes the set of all permissions. This set contains permissions governing kernel services as well as permissions governing services provided by user space servers.

The set *Kernel_permission* is used to denote the subset of *PERMISSION* that governs kernel services.

DTOS Kernel Definition 12

[*PERMISSION*]

| $Kernel_permission : \mathbb{P} PERMISSION$

⁹Note that the terms *access vector*, *service vector*, and *permission set* are used somewhat interchangeably.

The elements of *Kernel_permission* are enumerated in subsections 5.4.1-5.4.14. The operator *Values_partition* is formally defined in Appendix B. Informally, the expression $\langle val_1, \dots, val_n \rangle$ *Values_partition* *S* denotes that the values val_1, \dots, val_n are unique values that together comprise the set *val_set*.

5.4.1 IPC Permissions

The DTOS kernel enforces the following “IPC” permissions: *Can_receive*, *Can_send*, *Hold_receive*, *Hold_send*, *Hold_send_once*, *Interpose*, *Map_vm_region*, *Set_reply*, *Specify*, *Transfer_ool*, *Transfer_receive*, *Transfer_rights*, *Transfer_send*, *Transfer_send_once*. We use *Ip_permissions* to denote this set of permissions.

DTOS Kernel Definition 13

<p><i>Ip_permissions</i> : \mathbb{P} <i>PERMISSION</i> <i>Can_receive</i>, <i>Can_send</i>, <i>Hold_receive</i>, <i>Hold_send</i>, <i>Hold_send_once</i>, <i>Interpose</i>, <i>Map_vm_region</i>, <i>Set_reply</i>, <i>Specify</i>, <i>Transfer_ool</i>, <i>Transfer_receive</i>, <i>Transfer_rights</i>, <i>Transfer_send</i>, <i>Transfer_send_once</i> :</p>
<p><i>PERMISSION</i></p>
<p>\langle <i>Can_receive</i>, <i>Can_send</i>, <i>Hold_receive</i>, <i>Hold_send</i>, <i>Hold_send_once</i>, <i>Interpose</i>, <i>Map_vm_region</i>, <i>Set_reply</i>, <i>Specify</i>, <i>Transfer_ool</i>, <i>Transfer_receive</i>, <i>Transfer_rights</i>, <i>Transfer_send</i>, <i>Transfer_send_once</i> \rangle <i>Values_partition</i> <i>Ip_permissions</i></p>

5.4.2 Port Permissions

The DTOS kernel enforces the following permissions on port requests: *Add_name*, *Alter_pns_info*, *Extract_right*, *Lookup_ports*, *Manipulate_port_set*, *Observe_pns_info*, *Port_rename*, *Register_notification*, *Register_ports*, *Remove_name*. We use *Port_permissions* to denote this set of permissions.

DTOS Kernel Definition 14

<p><i>Port_permissions</i> : \mathbb{P} <i>PERMISSION</i> <i>Add_name</i>, <i>Alter_pns_info</i>, <i>Extract_right</i>, <i>Lookup_ports</i>, <i>Manipulate_port_set</i>, <i>Observe_pns_info</i>, <i>Port_rename</i>, <i>Register_notification</i>, <i>Register_ports</i>, <i>Remove_name</i> :</p>
<p><i>PERMISSION</i></p>
<p>\langle <i>Add_name</i>, <i>Alter_pns_info</i>, <i>Extract_right</i>, <i>Lookup_ports</i>, <i>Manipulate_port_set</i>, <i>Observe_pns_info</i>, <i>Port_rename</i>, <i>Register_notification</i>, <i>Register_ports</i>, <i>Remove_name</i> \rangle <i>Values_partition</i> <i>Port_permissions</i></p>

5.4.3 VM Permissions

The DTOS kernel enforces the following permissions on VM requests:

Access_machine_attribute, Allocate_vm_region, Chg_vm_region_prot, Copy_vm, Deallocate_vm_region, Get_vm_region_info, Get_vm_statistics, Read_vm_region, Set_vm_region_inherit, Wire_vm_for_task, Write_vm_region. We use *Vm_permissions* to denote this set of permissions.

DTOS Kernel Definition 15

<i>Vm_permissions</i> : \mathbb{P} PERMISSION <i>Access_machine_attribute, Allocate_vm_region, Chg_vm_region_prot, Copy_vm, Deallocate_vm_region, Get_vm_region_info, Get_vm_statistics, Read_vm_region, Set_vm_region_inherit, Wire_vm_for_task, Write_vm_region</i> : PERMISSION
\langle <i>Access_machine_attribute, Allocate_vm_region, Chg_vm_region_prot, Copy_vm, Deallocate_vm_region, Get_vm_region_info, Get_vm_statistics, Read_vm_region, Set_vm_region_inherit, Wire_vm_for_task, Write_vm_region</i> \rangle Values-partition <i>Vm_permissions</i>

5.4.4 Memory Object Permissions

The DTOS kernel enforces the following permissions on memory requests: *Have_execute, Have_read, Have_write, Page_vm_region.* We use *Memory_object_permissions* to denote this set of permissions.

DTOS Kernel Definition 16

<i>Memory_object_permissions</i> : \mathbb{P} PERMISSION <i>Have_execute, Have_read, Have_write, Page_vm_region</i> : PERMISSION
\langle <i>Have_execute, Have_read, Have_write, Page_vm_region</i> \rangle Values-partition <i>Memory_object_permissions</i>

5.4.5 Pager Permissions

The DTOS kernel enforces the following permissions on pager requests: *Change_page_locks, Destroy_object, Get_attributes, Invoke_lock_request, Make_page_precious, Provide_data, Remove_page, Revoke_ibac, Save_page, Set_attributes, Set_ibac_port, Supply_ibac.* We use *Pager_permissions* to denote this set of permissions.

DTOS Kernel Definition 17

<p><i>Pager_permissions</i> : P PERMISSION <i>Change_page_locks</i>, <i>Destroy_object</i>, <i>Get_attributes</i>, <i>Invoke_lock_request</i>, <i>Make_page_precious</i>, <i>Provide_data</i>, <i>Remove_page</i>, <i>Revoke_ibac</i>, <i>Save_page</i>, <i>Set_attributes</i>, <i>Set_ibac_port</i>, <i>Supply_ibac</i> : PERMISSION</p>
<p>(<i>Change_page_locks</i>, <i>Destroy_object</i>, <i>Get_attributes</i>, <i>Invoke_lock_request</i>, <i>Make_page_precious</i>, <i>Provide_data</i>, <i>Remove_page</i>, <i>Revoke_ibac</i>, <i>Save_page</i>, <i>Set_attributes</i>, <i>Set_ibac_port</i>, <i>Supply_ibac</i>) <i>Values_partition Pager_permissions</i></p>

5.4.6 Thread Permissions

The DTOS kernel enforces the following permissions on thread requests: *Abort_thread*, *Abort_thread_depress*, *Assign_thread_to_pset*, *Can_swch*, *Can_swch_pri*, *Depress_pri*, *Get_thread_assignment*, *Get_thread_exception_port*, *Get_thread_info*, *Get_thread_kernel_port*, *Get_thread_state*, *Initiate_secure*, *Raise_exception*, *Resume_thread*, *Sample_thread*, *Set_max_thread_priority*, *Set_thread_exception_port*, *Set_thread_kernel_port*, *Set_thread_policy*, *Set_thread_priority*, *Set_thread_state*, *Suspend_thread*, *Switch_thread*, *Terminate_thread*, *Wait_evt*, *Wire_thread_into_memory*. We use *Thread_permissions* to denote this set of permissions.

DTOS Kernel Definition 18

<p><i>Thread_permissions</i> : P PERMISSION <i>Abort_thread</i>, <i>Abort_thread_depress</i>, <i>Assign_thread_to_pset</i>, <i>Can_swch</i>, <i>Can_swch_pri</i>, <i>Depress_pri</i>, <i>Get_thread_assignment</i>, <i>Get_thread_exception_port</i>, <i>Get_thread_info</i>, <i>Get_thread_kernel_port</i>, <i>Get_thread_state</i>, <i>Initiate_secure</i>, <i>Raise_exception</i>, <i>Resume_thread</i>, <i>Sample_thread</i>, <i>Set_max_thread_priority</i>, <i>Set_thread_exception_port</i>, <i>Set_thread_kernel_port</i>, <i>Set_thread_policy</i>, <i>Set_thread_priority</i>, <i>Set_thread_state</i>, <i>Suspend_thread</i>, <i>Switch_thread</i>, <i>Terminate_thread</i>, <i>Wait_evt</i>, <i>Wire_thread_into_memory</i> : PERMISSION</p>
<p>(<i>Abort_thread</i>, <i>Abort_thread_depress</i>, <i>Assign_thread_to_pset</i>, <i>Can_swch</i>, <i>Can_swch_pri</i>, <i>Depress_pri</i>, <i>Get_thread_assignment</i>, <i>Get_thread_exception_port</i>, <i>Get_thread_info</i>, <i>Get_thread_kernel_port</i>, <i>Get_thread_state</i>, <i>Initiate_secure</i>, <i>Raise_exception</i>, <i>Resume_thread</i>, <i>Sample_thread</i>, <i>Set_max_thread_priority</i>, <i>Set_thread_exception_port</i>, <i>Set_thread_kernel_port</i>, <i>Set_thread_policy</i>, <i>Set_thread_priority</i>, <i>Set_thread_state</i>, <i>Suspend_thread</i>, <i>Switch_thread</i>, <i>Terminate_thread</i>, <i>Wait_evt</i>, <i>Wire_thread_into_memory</i>) <i>Values_partition Thread_permissions</i></p>

5.4.7 Task Permissions

The DTOS kernel enforces the following permissions on task requests: *Add_thread*, *Add_thread_secure*, *Assign_task_to_pset*, *Change_sid*, *Chg_task_priority*, *Create_task*, *Create_task_secure*, *Cross_context_create*, *Cross_context_inherit*, *Get_emulation*, *Get_task_assignment*, *Get_task_boot_port*, *Get_task_exception_port*, *Get_task_info*, *Get_task_kernel_port*, *Get_task_threads*, *Make_sid*, *Resume_task*, *Sample_task*, *Set_emulation*, *Set_ras*, *Set_task_boot_port*, *Set_task_exception_port*, *Set_task_kernel_port*, *Suspend_task*, *Terminate_task*, *Transition_sid*. We use *Task_task_permissions* to denote this set of permissions.

DTOS Kernel Definition 19

```

Task_task_permissions : P PERMISSION
Add_thread, Add_thread_secure, Assign_task_to_pset,
Change_sid, Chg_task_priority, Create_task,
Create_task_secure, Cross_context_create, Cross_context_inherit,
Get_emulation, Get_task_assignment, Get_task_boot_port,
Get_task_exception_port, Get_task_info, Get_task_kernel_port,
Get_task_threads, Make_sid, Resume_task,
Sample_task, Set_emulation, Set_ras,
Set_task_boot_port, Set_task_exception_port, Set_task_kernel_port,
Suspend_task, Terminate_task, Transition_sid :
PERMISSION

```

```

{Add_thread, Add_thread_secure, Assign_task_to_pset, Change_sid,
  Chg_task_priority, Create_task, Create_task_secure,
  Cross_context_create, Cross_context_inherit, Get_emulation,
  Get_task_assignment, Get_task_boot_port, Get_task_exception_port,
  Get_task_info, Get_task_kernel_port, Get_task_threads, Make_sid,
  Resume_task, Sample_task, Set_emulation, Set_ras, Set_task_boot_port,
  Set_task_exception_port, Set_task_kernel_port, Suspend_task,
  Terminate_task, Transition_sid}
Values_partition Task_task_permissions

```

We use *Task_permissions* to denote the union of *Task_task_permissions*, *Port_permissions*, and *Vm_permissions*.

DTOS Kernel Definition 20

```

Task_permissions : P PERMISSION
{Port_permissions, Vm_permissions, Task_task_permissions} partition Task_permissions

```

5.4.8 Host Name Port Permissions

The DTOS kernel enforces the following permissions on host name port requests: *Create_pset*, *Flush_permission*, *Get_audit_port*, *Get_authentication_port*, *Get_crypto_port*, *Get_default_pset_name*, *Get_host_control_port*, *Get_host_info*, *Get_host_name*, *Get_host_version*, *Get_negotiation_port*, *Get_network_ss_port*, *Get_security_master_port*, *Get_security_client_port*, *Get_special_port*, *Get_time*, *Pset_names*, *Set_audit_port*, *Set_authentication_port*, *Set_crypto_port*, *Set_negotiation_port*, *Set_network_ss_port*, *Set_security_master_port*, *Set_security_client_port*, *Set_special_port*. We use *Host_name_port_permissions* to denote this set of permissions.

DTOS Kernel Definition 21

<p><i>Host_name_port_permissions</i> : P PERMISSION <i>Create_pset, Flush_permission, Get_audit_port,</i> <i>Get_authentication_port, Get_crypto_port, Get_default_pset_name,</i> <i>Get_host_control_port, Get_host_info, Get_host_name,</i> <i>Get_host_version, Get_negotiation_port, Get_network_ss_port,</i> <i>Get_security_master_port, Get_security_client_port, Get_special_port,</i> <i>Get_time, Pset_names, Set_audit_port,</i> <i>Set_authentication_port, Set_crypto_port, Set_negotiation_port,</i> <i>Set_network_ss_port, Set_security_master_port, Set_security_client_port,</i> <i>Set_special_port :</i> PERMISSION</p>
<p>(<i>Create_pset, Flush_permission, Get_audit_port,</i> <i>Get_authentication_port, Get_crypto_port, Get_default_pset_name,</i> <i>Get_host_control_port, Get_host_info, Get_host_name, Get_host_version,</i> <i>Get_negotiation_port, Get_network_ss_port, Get_security_master_port,</i> <i>Get_security_client_port, Get_special_port, Get_time,</i> <i>Pset_names, Set_audit_port, Set_authentication_port,</i> <i>Set_crypto_port, Set_negotiation_port, Set_network_ss_port,</i> <i>Set_security_master_port, Set_security_client_port, Set_special_port)</i> <i>Values_partition Host_name_port_permissions</i></p>

5.4.9 Host Control Port Permissions

The DTOS kernel enforces the following permissions on host control port requests:

Get_boot_info, Get_host_processors, Pset_ctrl_port, Reboot_host, Set_default_memory_mgr,
Set_time, Wire_thread, Wire_vm. We use *Host_control_port_permissions* to denote this set of permissions.

DTOS Kernel Definition 22

<p><i>Host_control_port_permissions</i> : P PERMISSION <i>Get_boot_info, Get_host_processors, Pset_ctrl_port,</i> <i>Reboot_host, Set_default_memory_mgr, Set_time,</i> <i>Wire_thread, Wire_vm :</i> PERMISSION</p>
<p>(<i>Get_boot_info, Get_host_processors, Pset_ctrl_port, Reboot_host,</i> <i>Set_default_memory_mgr, Set_time, Wire_thread, Wire_vm)</i> <i>Values_partition Host_control_port_permissions</i></p>

5.4.10 Processor Permissions

The DTOS kernel enforces the following permissions on processor requests:

Assign_processor_to_set, Get_processor_assignment, Get_processor_info, May_control_processor.
We use *Processor_permissions* to denote this set of permissions.

DTOS Kernel Definition 23

<i>Processor_permissions</i> : \mathbb{P} <i>PERMISSION</i> <i>Assign_processor_to_set</i> , <i>Get_processor_assignment</i> , <i>Get_processor_info</i> , <i>May_control_processor</i> : <i>PERMISSION</i>
<i>(Assign_processor_to_set, Get_processor_assignment, Get_processor_info,</i> <i>May_control_processor)</i> <i>Values_partition Processor_permissions</i>

5.4.11 Processor Set Name Port Permissions

The DTOS kernel enforces the following permissions on processor set name port requests: *Get_pset_info*. We use *Procset_name_port_permissions* to denote this set of permissions.

DTOS Kernel Definition 24

<i>Procset_name_port_permissions</i> : \mathbb{P} <i>PERMISSION</i> <i>Get_pset_info</i> : <i>PERMISSION</i>
<i>(Get_pset_info)</i> <i>Values_partition Procset_name_port_permissions</i>

5.4.12 Processor Set Control Port Permissions

The DTOS kernel enforces the following permissions on processor set control port requests: *Assign_processor*, *Assign_task*, *Assign_thread*, *Chg_pset_max_pri*, *Define_new_scheduling_policy*, *Destroy_pset*, *Invalidate_scheduling_policy*, *Observe_pset_processes*. We use *Procset_control_port_permissions* to denote this set of permissions.

DTOS Kernel Definition 25

<i>Procset_control_port_permissions</i> : \mathbb{P} <i>PERMISSION</i> <i>Assign_processor</i> , <i>Assign_task</i> , <i>Assign_thread</i> , <i>Chg_pset_max_pri</i> , <i>Define_new_scheduling_policy</i> , <i>Destroy_pset</i> , <i>Invalidate_scheduling_policy</i> , <i>Observe_pset_processes</i> : <i>PERMISSION</i>
<i>(Assign_processor, Assign_task, Assign_thread, Chg_pset_max_pri,</i> <i>Define_new_scheduling_policy, Destroy_pset,</i> <i>Invalidate_scheduling_policy, Observe_pset_processes)</i> <i>Values_partition Procset_control_port_permissions</i>

We use *Procset_permissions* to denote the union of *Procset_name_port_permissions* and *Procset_control_port_permissions*.

DTOS Kernel Definition 26

<i>Procset_permissions</i> : \mathbb{P} <i>PERMISSION</i>
<i>(Procset_name_port_permissions, Procset_control_port_permissions)</i> <i>partition Procset_permissions</i>

5.4.13 Device Permissions

The DTOS kernel enforces the following permissions on device requests: *Close_device*, *Control_pager*, *Get_device_status*, *Map_device*, *Open_device*, *Read_device*, *Set_device_filter*, *Set_device_status*, *Write_device*. We use *Device_permissions* to denote this set of permissions.

DTOS Kernel Definition 27

$$\frac{\begin{array}{l} \textit{Device_permissions} : \mathbb{P} \textit{PERMISSION} \\ \textit{Close_device}, \textit{Control_pager}, \textit{Get_device_status}, \\ \textit{Map_device}, \textit{Open_device}, \textit{Read_device}, \\ \textit{Set_device_filter}, \textit{Set_device_status}, \textit{Write_device} : \\ \textit{PERMISSION} \end{array}}{\langle \textit{Close_device}, \textit{Control_pager}, \textit{Get_device_status}, \textit{Map_device}, \textit{Open_device}, \\ \textit{Read_device}, \textit{Set_device_filter}, \textit{Set_device_status}, \textit{Write_device} \rangle} \\ \textit{Values_partition Device_permissions}$$

5.4.14 Kernel Reply Port Permissions

The DTOS kernel enforces the following permissions on requests sent to kernel reply ports: *Provide_permission*. We use *Kernel_reply_permissions* to denote this set of permissions.

DTOS Kernel Definition 28

$$\frac{\begin{array}{l} \textit{Kernel_reply_permissions} : \mathbb{P} \textit{PERMISSION} \\ \textit{Provide_permission} : \\ \textit{PERMISSION} \end{array}}{\langle \textit{Provide_permission} \rangle} \\ \textit{Values_partition Kernel_reply_permissions}$$

We do not require that all of the above sets of permissions be non-overlapping. The only such requirement is that the *Ipc_permissions* do not overlap with any of the other sets. This is consistent with the current prototype in which permissions are simply integers specifying positions in access vectors. Because there are different types of access vector depending upon the type of target object, multiple permissions may specify the same access vector position. Every vector contains the IPC permissions stored at the same positions.

DTOS Kernel Definition 29

$$\begin{array}{l} \textit{Ipc_permissions} \\ \cap (\textit{Memory_object_permissions} \cup \textit{Pager_permissions} \\ \cup \textit{Thread_permissions} \cup \textit{Task_permissions} \\ \cup \textit{Host_name_port_permissions} \cup \textit{Host_control_port_permissions} \\ \cup \textit{Processor_permissions} \cup \textit{Procset_permissions} \\ \cup \textit{Device_permissions} \cup \textit{Kernel_reply_permissions}) \\ = \emptyset \end{array}$$

5.5 Access Vector Cache

The kernel receives an access decision from the security server as a *Ruling*. Each ruling consists of:

- *ssi* — a subject security identifier
- *osi* — an object security identifier
- *access_vector* — a set of granted permissions between the *ssi* and *osi*
- *control_vector* — the set of granted permissions which are allowed to be cached in the kernel for later access
- *expiration_value* — the time at which the cached permissions expire

DTOS Kernel Definition 30

Ruling

ssi : *SSI*
osi : *OSI*
access_vector : \mathbb{P} *PERMISSION*
control_vector : \mathbb{P} *PERMISSION*
expiration_value : \mathbb{N}

Review Note:

We need to be careful not to get bit by using *ssi* and *osi* in *Ruling*, since they are often used as “variables” also. Or else we could rename them here.

A ruling is usable for a given *ssi* and *osi* if the *ssi* and *osi* match those in the ruling and the ruling has not expired. The expression $Usable_ruling(ssi, osi, time)$ denotes the set of all such rulings with respect to *ssi*, *osi* and *time*, the time at which the ruling is consulted. When a ruling is initially received by the kernel, the kernel need only check the access vector and expiration time to see if a permission is granted. This is reflected by the function $Ruling_allows(ruling, ssi, osi)$ which returns the set of permissions in the access vector of *ruling* if *ssi* and *osi* are the same as in *ruling*.

Editorial Note:

The prototype does not currently check the expiration time in these cases, but we plan to correct this.

DTOS Kernel Definition 31

$Usable_ruling : SSI \times OSI \times \mathbb{N} \rightarrow \mathbb{P} Ruling$
 $Ruling_allows : Ruling \times SSI \times OSI \times \mathbb{N} \rightarrow \mathbb{P} PERMISSION$

$\forall ruling : Ruling; ssi : SSI; osi : OSI; time : \mathbb{N}; permission : PERMISSION$
 • $(ruling \in Usable_ruling(ssi, osi, time))$
 $\Leftrightarrow (ssi = ruling.ssi$
 $\quad \wedge osi = ruling.osi$
 $\quad \wedge time < ruling.expiration_value)$
 $\wedge permission \in Ruling_allows(ruling, ssi, osi, time)$
 $\Leftrightarrow (ruling \in Usable_ruling(ssi, osi, time)$
 $\quad \wedge permission \in ruling.access_vector)$

To enhance performance, the kernel is permitted to cache the rulings provided by security servers. A cached ruling is usable for a given *ssi*, *osi* and permission if the *ssi* and *osi* match those in the ruling, the permission is in the *control_vector* and the ruling has not expired. The expression $Usable_cached_ruling(ssi, osi, permission, time)$ denotes the set of all such rulings. Once cached, a ruling grants a particular *permission* from *ssi* to *osi* if the ruling is

usable and the permission is included in the *access_vector*. This is reflected by the function *Cached_ruling_allows*(*ruling*, *ssi*, *osi*, *time*), where *time* is the time at which the ruling is consulted.

DTOS Kernel Definition 32

$$\begin{array}{l}
 \hline
 Usable_cached_ruling : SSI \times OSI \times PERMISSION \times \mathbb{N} \rightarrow \mathbb{P} Ruling \\
 Cached_ruling_allows : Ruling \times SSI \times OSI \times \mathbb{N} \rightarrow \mathbb{P} PERMISSION \\
 \hline
 \forall ruling : Ruling; ssi : SSI; osi : OSI; time : \mathbb{N}; permission : PERMISSION \\
 \bullet (ruling \in Usable_cached_ruling(ssi, osi, permission, time) \\
 \Leftrightarrow (ruling \in Usable_ruling(ssi, osi, time) \\
 \wedge permission \in ruling.control_vector)) \\
 \wedge (permission \in Cached_ruling_allows(ruling, ssi, osi, time) \\
 \Leftrightarrow (ruling \in Usable_cached_ruling(ssi, osi, permission, time) \\
 \wedge permission \in ruling.access_vector))
 \end{array}$$

The kernel cache is a set of rulings, represented by *cache*. There may only be one unexpired ruling in the cache for each (*ssi*, *osi*) pair. The function *cache_allows*(*ssi*, *osi*) returns the set of permissions granted to the (*ssi*, *osi*) pair by the rulings in the cache according to the function *Cached_ruling_allows*. The quadruple (*ssi*, *osi*, *permission*, *ruling*) is in *cached_ruling_avail* if and only if *ruling* is in the cache and it is usable for *ssi*, *osi* and *permission* at the current time.

DTOS Kernel Definition 33

$$\begin{array}{l}
 \hline
 KernelCache \\
 \underline{cache} : \mathbb{P} Ruling \\
 cache_allows : SSI \times OSI \rightarrow \mathbb{P} PERMISSION \\
 cached_ruling_avail : \mathbb{P}(SSI \times OSI \times PERMISSION \times Ruling) \\
 HostTime \\
 \hline
 \forall ruling_1, ruling_2 : Ruling \\
 | \{ ruling_1, ruling_2 \} \subseteq \underline{cache} \\
 \wedge ruling_1.ssi = ruling_2.ssi \\
 \wedge ruling_1.osi = ruling_2.osi \\
 \wedge ruling_1.expiration_value > \underline{host_time} \\
 \wedge ruling_2.expiration_value > \underline{host_time} \\
 \bullet ruling_1 = ruling_2 \\
 \\
 \forall ssi : SSI; osi : OSI \\
 \bullet cache_allows(ssi, osi) = \bigcup \{ ruling : Ruling \mid ruling \in \underline{cache} \\
 \bullet Cached_ruling_allows(ruling, ssi, osi, \underline{host_time}) \} \\
 \\
 \forall ssi : SSI; osi : OSI; permission : PERMISSION; ruling : Ruling \\
 \bullet (ssi, osi, permission, ruling) \in cached_ruling_avail \\
 \Leftrightarrow (ruling \in \underline{cache} \\
 \cap Usable_cached_ruling(ssi, osi, permission, \underline{host_time}))
 \end{array}$$

5.6 Message Security Information

Each existing message has an SSI associated with it that indicates the SSI of the task that sent the message. The expression *msg_sending_sid*(*message*) indicates the SSI of the task that

sent *message*. In addition, certain messages have an associated SSI that indicates which tasks may receive the message. The set *msg_receiver_specified* indicates the set of messages that have a receiving SID specified, and *msg_receiving_sid(message)* indicates the receiving SSI for each message in this set. As part of the processing of a message, the sender's permissions to the destination port are computed and attached to the message. The set *msg_ruling_computed* denotes the set of messages for which the permissions have already been computed, and *msg_ruling(message)* indicates the associated set of permissions for each such message. A ruling must be computed for each message before the message can be enqueued at a port. An "effective" sending SID and access vector may optionally be specified by the sender of a message. The expressions *msg_specified_sid(message)* and *msg_specified_vector(message)* indicate, respectively, the "effective" SID and access vector specified by the sender.

Editorial Note:

Need to think about how to model the specified vectors. The current specification ignores the cache control and notification vectors. The prototype currently has all three vectors represented explicitly. It has been implemented to allow the number of vectors to be easily changed.

DTOS Kernel Definition 34

<p><i>DtosMessages</i></p> <p><i>MessageExist</i></p> <p><i>MessageQueues</i></p> <p><i>msg_sending_sid</i> : MESSAGE \leftrightarrow SSI</p> <p><i>msg_receiver_specified</i> : \mathbb{P} MESSAGE</p> <p><i>msg_receiving_sid</i> : MESSAGE \leftrightarrow SSI</p> <p><i>msg_ruling_computed</i> : \mathbb{P} MESSAGE</p> <p><i>msg_ruling</i> : MESSAGE \leftrightarrow Ruling</p> <p><i>msg_specified_sid</i> : MESSAGE \leftrightarrow SSI</p> <p><i>msg_specified_vector</i> : MESSAGE \leftrightarrow \mathbb{P} PERMISSION</p> <hr/> <p>dom <i>msg_sending_sid</i> = <i>message_exists</i></p> <p>dom <i>msg_receiving_sid</i> = <i>msg_receiver_specified</i> \subseteq <i>message_exists</i></p> <p>dom <i>msg_ruling</i> = <i>msg_ruling_computed</i> \subseteq <i>message_exists</i></p> <p>dom <i>containing_port</i> \subseteq <i>msg_ruling_computed</i></p> <p>dom <i>msg_specified_sid</i> \subseteq <i>message_exists</i></p> <p>dom <i>msg_specified_vector</i> \subseteq <i>message_exists</i></p>
--

5.7 Task Creation Information

Each task has a state used in controlling the secure initiation of threads within that task. The type *TASK_CREATION_STATE* is comprised of the possible values of this state. The recognized values of this type are:

- *Tcs_task_empty* — indicates a task that was created using **task_create_secure** and does not yet have any threads.
- *Tcs_thread_created* — indicates a task created using **task_create_secure** for which a thread has been created using **thread_create_secure** but has not had its initial state set.

- *Tcs_thread_state_set* — indicates a task created using **task_create_secure** for which a thread has been created using **thread_create_secure** that has had its initial state set using **thread_set_state_secure** but has not been resumed (i.e., started).
- *Tcs_task_ready* — indicates either a task that was not created using **task_create_secure** or a task that was created using **task_create_secure** and which has a thread that was created using **thread_create_secure**, has had its state set using **thread_set_state_secure**, and has been resumed using **thread_resume_secure**.

These states are used to ensure that processes initiated using **task_create_secure** follow the normal process initiation sequence of:

1. Create the task.
2. Create a thread within the task.
3. Set the state of the thread.
4. Resume the thread.

Review Note:

The above, particularly the description of *Tcs_task_ready*, must be checked against the prototype

This allows an untrusted process to create a trusted process using **task_create_secure** while prohibiting the untrusted process from (for example) changing the state of threads in the trusted process after the trusted process has started execution.

The expression $\underline{task_creation_state}(task)$ denotes the creation state of *task*.

DTOS Kernel Definition 35

$$TASK_CREATION_STATE ::= Tcs_task_empty \mid Tcs_thread_created \\ \mid Tcs_thread_state_set \mid Tcs_task_ready$$

$\underline{TaskCreationState}$ $\underline{TaskExist}$ $\underline{task_creation_state} : TASK \rightarrow TASK_CREATION_STATE$ $\text{dom } \underline{task_creation_state} = \underline{task_exists}$

The Mach model of process creation uses an existing task to serve as a “template” for each new task. This task is the *parent_task* parameter to **task_create**. A newly created task inherits parts of its environment, such as portions of its address space, from the “parent” task. To simplify the statement of the security requirements on task creation, we introduce $\underline{parent_task}(task)$ to denote *task*’s parent.¹⁰

DTOS Kernel Definition 36

$\underline{ParentTask}$ $\underline{parent_task} : TASK \rightarrow TASK$

¹⁰Note that this information is not actually recorded in the current design. Since we only use this information for stating requirements on task creation and this information is available at this point in the processing in the implementation, this deviation between the model and the implementation is tolerable.

5.8 Server Ports

The kernel records the ports to be used for communications with certain servers:

- *security_server_master_port* denotes the port used by the kernel to make requests of the security server.
- *security_server_client_port* denotes the port used by non-kernel clients to make requests of the security server.
- *authentication_server_port* denotes the port used to make requests of the authentication server.
- *audit_server_port* denotes the port used to make requests of the audit server.
- *crypto_server_port* denotes the port used to make requests of the crypto server.
- *negotiation_server_port* denotes the port used to make requests of the negotiation server.
- *network_ss_port* denotes the port used to make security requests over the network.

DTOS Kernel Definition 37

<pre> ServerPorts security_server_master_port : PORT security_server_client_port : PORT authentication_server_port : PORT audit_server_port : PORT crypto_server_port : PORT negotiation_server_port : PORT network_ss_port : PORT </pre>

When the kernel requests an access computation from the Security Server, it specifies a reply port to which the computed accesses should be sent. We use *kernel_reply_ports* to denote the set of ports that the kernel has specified as reply ports for requests to the Security Server.

DTOS Kernel Definition 38

<pre> KernelReplyPorts PortExist kernel_reply_ports : P PORT kernel_reply_ports ⊆ port_exists </pre>
--

5.9 Memory Region Protections

The current protection of a region limits a task's access to that region. It is calculated as the intersection of the Mach protection together with the accesses allowed for a task to a memory region by the relevant access vector. We use *protection(task, index)* to denote current protections of the region denoted by a given task-index pair.¹¹

Mach Definition 108

¹¹The prototype does not currently implement the enforcement of read-only access. The low-level memory routines in the prototype treat read and execute interchangeably.

Protection

MachProtection
 $protection : (TASK \times PAGE_INDEX) \leftrightarrow \mathbb{P} PROTECTION$

$\text{dom } protection = \text{dom } \underline{mach_protection}$
 $\forall task_page_index : TASK \times PAGE_INDEX$
 $| task_page_index \in \text{dom } protection$

- $protection(task_page_index) \subseteq \underline{mach_protection}(task_page_index)$

5.10 Summary of DTOS Kernel State

The DTOS kernel state is the Mach kernel state augmented with the access vector cache and the security information associated with subjects, objects, and messages.

DTOS Kernel Definition 39

DtosAdditions

SubjectSid
ObjectSid
TargetSids
KernelAs
KernelCache
DtosMessages
TaskCreationState
ParentTask
ServerPorts
KernelReplyPorts
Protection

DTOS Kernel Definition 40

Dtos

Mach
DtosAdditions

Section 6

Kernel Execution Model

This section describes the computational model used to represent the DTOS kernel requests and the additional data structures that are required to support this computational model. The organization of this section is as follows:

- Section 6.1, **Execution Summary**, gives a high level overview of the execution model and its data structures. The following sections give detailed descriptions of transitions which occur in the processing of all requests.
- Section 6.2, **Utility Transitions**, discusses several utility transitions that are used in various specifications.
- Section 6.3, **Trap Invocation**, discusses the transitions which occur at the start of any request.
- Sections 6.4 through 6.6 describe the initial processing common to all kernel requests which are made through the **mach_msg** trap.
- Section 6.7, **Definitions**, defines the data structures used to implement the transitions discussed in the previous sections.

6.1 Execution Summary

The DTOS execution model centers on the selection of a set of common “break points” in the processing of a kernel request. The break points are chosen to highlight significant processing events such as request invocation and service checking. We address the issue of atomicity by selecting an appropriate number and type of break points. One advantage of this approach is flexibility with respect to the level of detail in the model; we can easily change the amount of concurrency and level of detail by redefining the break points and the transitions which govern them.

Given a set of break points, every kernel request can be viewed as a sequence of transitions describing how processing moves from one break point to the next. In order to specify these transitions we need to know the current execution state for each thread. In the model, we maintain the execution status of each existing thread by setting the values of a function called *breaks*. The domain of this function is the set of existing threads and the values indicate what type of transitions have occurred as well as what information the thread needs to resume processing. In a sense the *breaks* function is analogous to a processor’s stack where information is stored between context switches, although the particular break points modeled by *breaks* do not in general coincide with actual context switches.

Every thread executing in user space maintains a value of *Bk_user_space*. To enter kernel space a thread issues an instruction to trap into kernel code. We model such a transition by changing the relevant value of the function *breaks* from *Bk_user_space* to *Bk_new_trap*. Similarly, we model a transition where a thread starts at a break point labeled *Bk_point_A* and ends at a break point labeled *Bk_point_B* by changing the relevant value of the function *breaks* from *Bk_point_A* to *Bk_point_B*. The following sections describe the specific break points, their interpretations in the execution model, the information needed to resume processing, and the flow of processing from one break point to the next.

In the Mach kernel, many different requests share common features of processing. Another advantage of our model of execution is that it is easy to specify transitions which are common to many requests in a reusable manner. In general we begin by discussing transitions which are common to all requests and then discuss specialized transitions and finally discuss the transitions which are specific to a particular request. This section details the common transitions, transitions common to a class of requests are discussed in request chapter introductions, while the request specific transitions are specified in individual request sections.

Now we describe the data structures constituting the execution model. The values of *breaks* are elements of the free type *BREAK_STATUS*. These values indicate the current processing status of a thread together with the environment needed to resume processing. The elements of *BREAK_STATUS* are discussed in the following sections as preconditions and postconditions to transitions. The formal definition of *BREAK_STATUS* is given in Section 6.7.5. We use the schema *Breaks* to define the *breaks* function.

<p><i>Breaks</i></p> <hr/> <p><i>ThreadExist</i></p> <p>$breaks : THREAD \mapsto BREAK_STATUS$</p> <hr/> <p>$dom\ breaks = \underline{thread_exists}$</p>

The state for the DTOS execution model consists of the components present in the DTOS kernel state together with the function *breaks*.

<p><i>DtosExec</i></p> <hr/> <p><i>Dtos</i></p> <p><i>Breaks</i></p>
--

We introduce a special schema *Transition* which serves as the signature for every main transition. This schema introduces the DTOS kernel state and declares four variables. The first, *cpu??*, is the processor on which the transition is occurring. The other three are derived from *cpu??* and are included as aliases to commonly used state elements. The variable *curr_th??* is the thread which is currently executing on *cpu??*, *curr_task??* is *curr_th??*'s parent task, and *curr_bk??* is the execution status of *curr_th??*.¹²

<p><i>Transition</i></p> <hr/> <p>$\Delta DtosExec$</p> <p>$cpu?? : PROCESSOR$</p> <p>$curr_th?? : THREAD$</p> <p>$curr_task?? : TASK$</p> <p>$curr_bk?? : BREAK_STATUS$</p> <hr/> <p>$cpu?? \in dom\ \underline{active_thread}$</p> <p>$curr_th?? = \underline{active_thread}(cpu??)$</p> <p>$curr_task?? = \underline{owning_task}(curr_th??)$</p> <p>$curr_bk?? = \underline{breaks}(curr_th??)$</p>
--

¹²The double questionmark decoration is used to provide a distinct look to these four components, since they have an interpretation distinct from either elements of the system state or inputs or outputs.

6.2 Utility Transitions

We begin our discussion of break points by specifying several transitions, shown in Figure 1, which are used as utilities. Each box in the diagram represents a complete transition; the first line in each box gives the name of the corresponding transition schema while the next two lines describe the break type of the precondition and postcondition, respectively. For example, the right-most box in Figure 1 describes the transition *RulingInCache* which has as precondition the existence of a break of type *Bk_check_pending*, and produces a break of type *Bk_have_ruling* as a postcondition. In these diagrams, a solid arrow from *TransitionOne* to *TransitionTwo* indicates that *TransitionOne* precedes *TransitionTwo* **and** no other transitions from this request intervene (of course concurrency allows transitions from other requests to occur.) By contrast, a dashed line indicates that intervening requests may occur. For example, the line from *RulingNotInCache* to *RulingObtained* is dashed to reflect the fact that when the kernel waits for a ruling from the security server the most general interaction could involve repeated failures and retries.

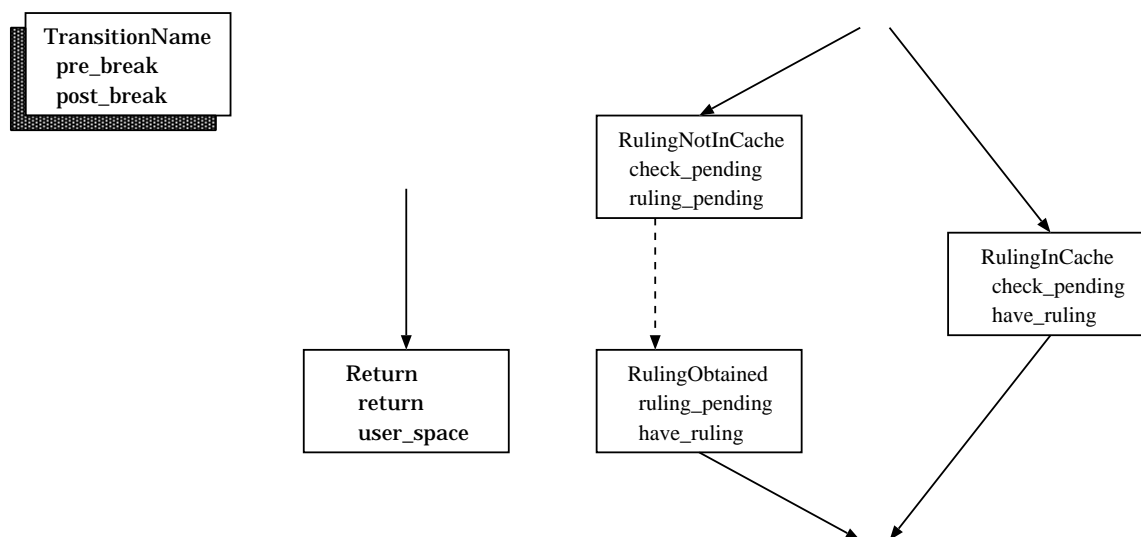


Figure 1: Utility Transitions

6.2.1 The Return Utility

Here we discuss the transition associated with request termination. The final transition in the processing of an IPC based request occurs when the kernel builds a return message containing status information and a specific kernel reply. We model such a transition with the schema *Return*. The precondition for this transition is the existence of a break of type *Bk_return* and the post condition is the creation of a break of type *Bk_user_space* signaling the fact that the thread has left kernel processing. The break *Bk_user_space* does not store any environment parameters (since no further processing is necessary) while *Bk_return* maintains the following information:

- *reply_to_port* — the port where the return message should be enqueued (if dead or null no message is sent),
- *operation* — the operation identifier for the terminating request,

- *reply* — an element of the type *KERNEL_REPLY* representing request specific output parameters supplied by the kernel,
- *return* — an element of the type *KERNEL_RETURN* describing the error status of request processing.

The kernel uses these parameters to build the return message. We model message construction with a set of functions: *Outputs_to_body*, *Reply_size*, *Reply_complex*, and *Reply_op*. These functions and the types *KERNEL_REPLY* and *KERNEL_RETURN* are discussed in detail in Section 6.7.1. It is worth noting that this reply message is not sent *viamach_msg*; rather the kernel builds the message and directly enqueues it at *reply_to_port*.

Editorial Note:
We do not currently specify enqueueing of the message.

<i>Return</i>
<i>Transition</i>
$\begin{aligned} &\exists \text{ message} : \text{MESSAGE}; \text{int_msg} : \text{InternalMessage}; \\ &\quad \text{reply_to_port} : \text{PORT}; \text{operation} : \text{OPERATION}; \\ &\quad \text{reply} : \text{KERNEL_REPLY}; \text{return} : \text{KERNEL_RETURN} \\ &\bullet \text{ curr_bk??} = \text{Bk_return}(\text{reply_to_port}, \text{operation}, \text{reply}, \text{return}) \\ &\quad \wedge \text{message} \notin \underline{\text{message_exists}} \\ &\quad \wedge \underline{\text{message_exists}}' = \underline{\text{message_exists}} \cup \{ \text{message} \} \\ &\quad \wedge \underline{\text{msg_contents}}(\text{message}) = \text{int_msg} \\ &\quad \wedge \text{int_msg.header.local_rights} = \emptyset \\ &\quad \wedge \text{int_msg.header.complex} \\ &\quad \quad = \text{Reply_complex}(\text{operation}, \text{Outputs_to_body}(\text{reply}, \text{return})) \\ &\quad \wedge \text{int_msg.header.size} \\ &\quad \quad = \text{Reply_size}(\text{operation}, \text{Outputs_to_body}(\text{reply}, \text{return})) \\ &\quad \wedge \text{int_msg.header.remote_port} = \text{reply_to_port} \\ &\quad \wedge \text{int_msg.header.local_port} = \emptyset \\ &\quad \wedge \text{int_msg.header.operation} = \text{Reply_op}(\text{operation}) \\ &\quad \wedge \text{int_msg.body} = \text{Outputs_to_body}(\text{reply}, \text{return}) \\ &\quad \wedge \text{int_msg.option} = \{ \text{Mach_send_msg} \} \\ &\quad \wedge \text{breaks}' = \text{breaks} \oplus \{ \text{curr_th??} \mapsto \text{Bk_user_space} \} \end{aligned}$

6.2.2 Permission Checking

Next we define the set of transitions involved in specifying a permission check. There are two possible transitions at the start of a permission check: *RulingInCache* or *RulingNotInCache*. The precondition for each of these transitions is the existence of a break of type *Bk_check_pending*. The transition *RulingInCache* examines the cache, determines that a cached ruling is available for the permission check, and creates a new break of type *Bk_have_ruling*. The transition *RulingNotInCache* examines the cache, determines that a cached ruling is not available, and creates a new break of type *Bk_ruling_pending*. In this case the kernel continues processing by consulting the Security Server. We model this transition by the schema *RulingObtained* which has as precondition the existence of a break of type *Bk_ruling_pending* and which creates a new break of type *Bk_have_ruling*.

The permission checking transitions need to maintain several environment parameters. These are interpreted as:

- *ssi* — the subject SID of the check,
- *osi* — the object SID of the check,
- *perm* — the required permission,
- *env* — the stored environment needed to resume processing,
- *op_allowed* — the boolean flag determining permission.

The first four of these are used in several places so we combine them in a structure called *CheckPending*:

```

CheckPending
-----
ssi : SSI
osi : OSI
perm : PERMISSION
env : ENVIRONMENT

```

There are three distinct contexts in which a permission check may be required: at the beginning of a system trap (e.g. **mach_thread_self**), at the beginning of an IPC based request (e.g. the service check for **thread_get_state**), or later in the processing of an IPC based request (e.g. the deferred check in **thread_get_special_port**.) As such, the parameter *env* needs to store one of three different types of data. To handle these three cases we define a free type called *ENVIRONMENT*, which is described in Section 6.7.4.

The first permission checking utility transition is *RulingInCache*. When a permission check is initiated, the kernel consults the cache to determine if there is an applicable ruling. The schema *RulingInCache* models the case where a permission check has been requested, and the kernel verifies that the cache contains an applicable ruling. The precondition of *RulingInCache* is the existence of a break of type *Bk_check_pending*, reflecting the condition that the processing of some request is waiting for a permission check. The postcondition, *Bk_have_ruling*, reflects the fact that an available ruling was found; in this case the result of the permission check is stored in the parameter *perm*. The parameter *env* is passed along unchanged.

```

RulingInCache
-----
Transition
-----
∃ CheckPending; ruling : Ruling; op_allowed : BOOLEAN
• curr_bk?? = Bk_check_pending(ssi, osi, perm, env)
  ∧ (ssi, osi, perm, ruling) ∈ cached_ruling_avail
  ∧ op_allowed
  = if perm ∈ Cached_ruling_allows(ruling, ssi, osi, host_time)
    then True
    else False
  ∧ breaks' = breaks ⊕ { curr_th?? ↦ Bk_have_ruling(perm, op_allowed, env) }

```

The schema *RulingNotInCache* models the case where a permission check has been requested and the kernel has determined that the cache does not contain an applicable ruling. Again the precondition is the existence of a break of type *Bk_check_pending*, but in this case the postcondition is a break of type *Bk_ruling_pending* reflecting the fact that the kernel is waiting for a ruling from the Security Server. As before, the parameter *env* is passed along unchanged.

<p><i>RulingNotInCache</i></p> <hr/> <p><i>Transition</i></p> <hr/> <p>\exists <i>CheckPending</i></p> <ul style="list-style-type: none"> • $curr_bk?? = Bk_check_pending(ssi, osi, perm, env)$ $\wedge (\forall ruling : Ruling \mid ruling \in \underline{cache}$ <ul style="list-style-type: none"> • $ruling \notin Usable_cached_ruling(ssi, osi, perm, \underline{host_time})$ $\wedge breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_ruling_pending(ssi, osi, perm, env) \}$

The schema *RulingObtained* models a transition where the kernel receives a valid ruling from the Security Server. The precondition is the existence of a break of type *Bk_have_ruling* and the postcondition is the creation of a new break of type *Bk_have_ruling*. The result of the permission check is stored in the parameter *perm*. As before, the parameter *env* is passed along unchanged.

Editorial Note:
The ruling obtained from the Security Server is modeled as a kernel input, but we do not specify how *ruling?* gets added to the cache.

<p><i>RulingObtained</i></p> <hr/> <p><i>Transition</i></p> <p><i>ruling?</i> : <i>Ruling</i></p> <hr/> <p>\exists <i>CheckPending</i>; <i>op_allowed</i> : <i>BOOLEAN</i></p> <ul style="list-style-type: none"> • $curr_bk?? = Bk_ruling_pending(ssi, osi, perm, env)$ $\wedge ruling? \in Usable_ruling(ssi, osi, \underline{host_time})$ $\wedge op_allowed$ $= \mathbf{if} \ perm \in Ruling_allows(ruling?, ssi, osi, \underline{host_time})$ $\mathbf{then} \ True$ $\mathbf{else} \ False$ $\wedge breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_have_ruling(perm, op_allowed, env) \}$
--

It is important to note that there are in fact many transitions between *RulingNotInCache* and *RulingObtained* including the sending of a message to the security server and receipt of a response.

6.3 Trap Invocation

All kernel requests are initiated through invocation of a trap while a thread is executing in user space. The thread must specify the particular trap identifier as well as some collection of parameters.

The precondition for this transition is the existence of a break of type *Bk_user_space* indicating that the thread is currently executing in user space. The postcondition is the creation of a new break of type *Bk_new_trap*. A break of type *Bk_new_trap* maintains two parameters: *trap_id?* identifies the type of trap being invoked and *user_spec?* contains components for the user supplied parameters. These types are described in Section 6.7.2 and Section 6.7.3.1.

<i>Invoke</i> <i>Transition</i> $trap_id? : TRAP_ID$ $user_spec? : UserSpecified$
$curr_bk?? = Bk_user_space$ $trap_id? \in Trap_ids$ $breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_new_trap(trap_id?, user_spec?) \}$

6.4 Initial `mach_msg` processing

In this section we discuss the transitions shown in Figure 2 which specify the early processing associated with the invocation of a `mach_msg` trap.

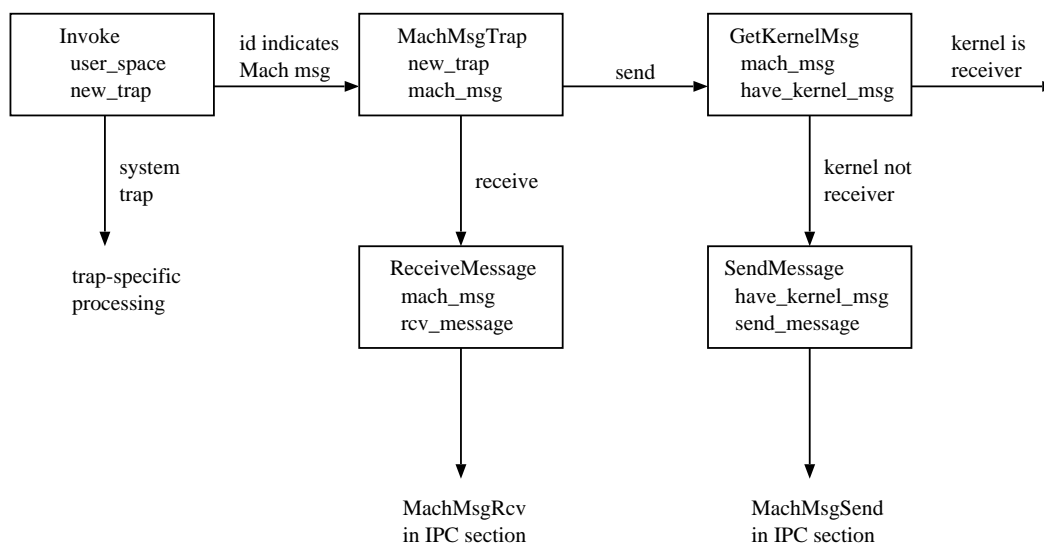


Figure 2: `mach_msg` Trap Invocation

If $trap_id?$ indicates that the new break is `Mach_msg_trap` processing continues with `MachMsgTrap`; the precondition is the existence of a break of type `Bk_new_trap` and the postcondition is the creation of a new break of type `Bk_mach_msg`. The break `Bk_new_trap` carries the trap identifier, `Mach_msg_trap`, as well as the user parameters `user_spec` which in this case contains the user space message.

<i>MachMsg</i> <i>Transition</i>
$\exists user_spec : UserSpecified$ <ul style="list-style-type: none"> $curr_bk?? = Bk_new_trap(Mach_msg_trap, user_spec)$ $\wedge breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_mach_msg(user_spec) \}$

Now we have two cases: this is a request to receive a message or a request to send a message.

Processing to receive a message is initiated in the transition `ReceiveMessage`. A client requests to receive a message by including `Mach_rcv_msg` in the set of options. The precondition also

includes the existence of a break of type Bk_mach_msg . The postcondition is the creation of a new break of type $Bk_rcv_message$. Subsequent processing of a receive request is described in Section C.1.

<i>ReceiveMessage</i>
<i>Transition</i>
$\begin{aligned} &\exists user_spec : UserSpecified \\ &\bullet curr_bk?? = Bk_mach_msg(user_spec) \\ &\quad \wedge Mach_rcv_msg \in user_spec.options \\ &\quad \wedge Mach_send_msg \notin user_spec.options \\ &\quad \wedge breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_rcv_message(user_spec) \} \end{aligned}$

Otherwise this is a request to send a message, in which case processing continues with the conversion of the user space message into an internal representation. The transition *GetKernelMsg* models a transition where the kernel resolves local name references to port references and virtual memory addresses to physical addresses. Here the precondition is the existence of a break of type Bk_mach_msg and the postcondition is a new break of type $Bk_have_kernel_msg$. The kernel message is modeled by an element of type *InternalMessage*; this is discussed in Section 6.7.3.2.

Editorial Note:
Currently this is modeled as a “black box” transition, but the utilities exist (in the IPC section) to specify the conversion.

<i>GetKernelMsg</i>
<i>Transition</i>
$\begin{aligned} &\exists user_spec : UserSpecified; int_msg : InternalMessage \\ &\bullet curr_bk?? = Bk_mach_msg(user_spec) \\ &\quad \wedge Mach_send_msg \in int_msg.option \\ &\quad \wedge breaks' = breaks \\ &\quad \oplus \{ curr_th?? \mapsto Bk_have_kernel_msg(int_msg) \} \end{aligned}$

Now there are two cases to consider depending on whether or not the kernel is the receiver for the message.

If the kernel is not the receiver, this is a request to send a message and we model continued processing with the transition *SendMessage*. The precondition is the existence of a break of type $Bk_have_kernel_msg$ and the postcondition is the creation of a new break of type $Bk_send_message$. Subsequent processing of a send request are described in Section C.1.

<i>SendMessage</i>
<i>Transition</i>
$\begin{aligned} &\exists int_msg : InternalMessage \\ &\bullet curr_bk?? = Bk_have_kernel_msg(int_msg) \\ &\quad \wedge \underline{k}ernel \neq receiver(int_msg.header.remote_port) \\ &\quad \wedge breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_send_message(int_msg) \} \end{aligned}$

The case in which the kernel is the receiver is considered further in the next section.

6.5 Service Checks for IPC Based Kernel Requests

Kernel requests which are received through the `mach_msg` trap generally must pass through an initial service check to determine if the client has permission to make the request. This is performed whenever the permission required by the request is dependent only upon the client, the port provided as the “target” port in the request and the operation identifier. For a few requests, this information is not sufficient and the permission check is deferred.

Figure 3 shows the transitions described within this section.

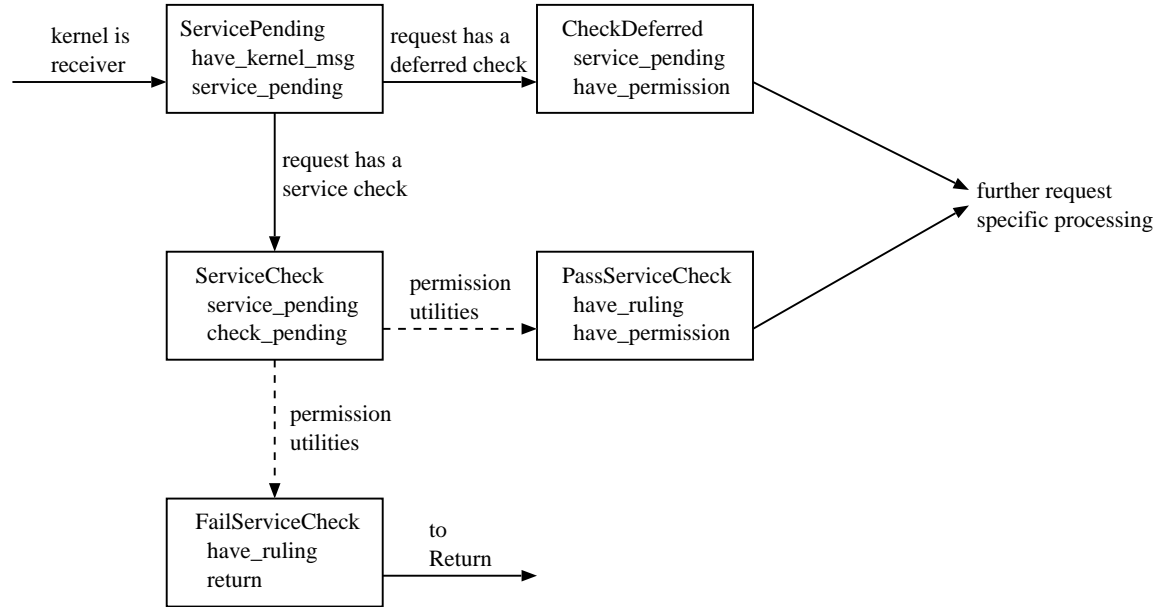
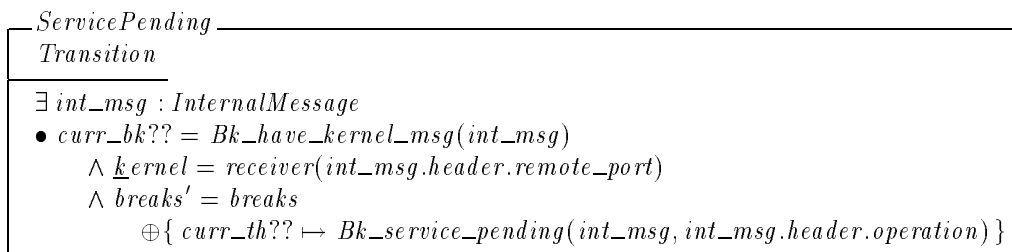


Figure 3: Message Transmission

The kernel prepares for the permission check by determining the operation that is being requested. We model this processing with the transition *ServicePending*. The precondition is the existence of a break of type *Bk_have_kernel_msg* and the postcondition is the creation of a new break of type *Bk_service_pending*. As before, the parameter *int_msg* is carried along for future use. In addition, this transition determines the operation and stores the value in the *operation* parameter of *Bk_service_pending*.



Next the kernel determines whether or not the client has permission to request the operation. Each operation typically has an associated *primary permission* that a client must have in order to successfully call the operation. Checking that the client has this primary permission is referred to as making the *service check*. For example, the primary permission

associated with the **thread_create** request is *Add_thread*. The **thread_create** request does not have any other permissions associated with it. As another example, the primary permission associated with the **mach_port_allocate** request is *Add_name*. However, there are other permissions such as *Hold_receive* that are also relevant to the **mach_port_allocate** request. The expression $Required_permission(operation)$ denotes the primary permission, if any, associated with operation *operation*. For certain operations the service check is deferred because the required permission depends on a parameter that must first be processed. Once the parameter has been extracted from the message the appropriate permission check is performed. The set *Service_check_deferred* is the set of all such operations. No operation in the domain of *Required_permission* can be in the set *Service_check_deferred*.

DTOS Kernel Definition 41

$$\begin{array}{|l} Required_permission : OPERATION \mapsto PERMISSION \\ Service_check_deferred : \mathbb{P} OPERATION \\ \hline Service_check_deferred \cap \text{dom } Required_permission = \emptyset \end{array}$$

The schema *CheckDeferred* models a transition in which the service check for a request is deferred until further processing can extract the appropriate parameters for the check. The precondition is the existence of a break of type *Bk_service_pending* and the resulting break, *Bk_have_permission*, takes an element of type *PERMISSION* as one of its arguments. We introduce a dummy permission called *permission_deferred* to act as the desired permission in a deferred service check.

$$| \quad permission_deferred : PERMISSION$$

$$\begin{array}{|l} CheckDeferred \\ Transition \\ \hline \exists int_msg : InternalMessage; operation : OPERATION \\ \bullet curr_bk?? = Bk_service_pending(int_msg, operation) \\ \quad \wedge operation \in Service_check_deferred \\ \quad \wedge breaks' = breaks \\ \quad \oplus \{ curr_th?? \mapsto Bk_have_permission(int_msg, permission_deferred) \} \end{array}$$

If the service check is not deferred, the schema *ServiceCheck* models a transition where the current thread waits for the kernel to obtain a ruling from the cache or from the security server. The precondition is the existence of a break of type *Bk_service_pending* and the resulting break, *Bk_check_pending*, contains the parameters for the permission check. The break *Bk_check_pending* requires an element of type *ENVIRONMENT*; the expression $E_kern(int_msg)$ packages the internal message into an element of this type. The function *E_kern* is discussed in Section 6.7.4.

<i>ServiceCheck</i>
<i>Transition</i>
$\exists \text{int_msg} : \text{InternalMessage}; \text{operation} : \text{OPERATION};$ $\text{ssi} : \text{SSI}; \text{osi} : \text{OSI}; \text{perm} : \text{PERMISSION}$ <ul style="list-style-type: none"> • $\text{curr_bk}?? = \text{Bk_service_pending}(\text{int_msg}, \text{operation})$ $\wedge \text{operation} \notin \text{Service_check_deferred}$ $\wedge \text{ssi} = \text{task_sid}(\text{curr_task}??)$ $\wedge \text{osi} = \text{port_sid}(\text{int_msg.header.remote_port})$ $\wedge \text{perm} = \text{Required_permission}(\text{operation})$ $\wedge \text{breaks}' = \text{breaks}$ $\oplus \{ \text{curr_th}?? \mapsto \text{Bk_check_pending}(\text{ssi}, \text{osi}, \text{perm}, \text{E_kern}(\text{int_msg})) \}$

Editorial Note:

The OSI in the previous might need more elaboration depending upon the use of the “self” SIDs. Also, if we ever began to correctly consider specified sender SIDs, that would need to be taken care of around this point in the processing.

After obtaining a ruling, the kernel examines the access vector to determine if the operation is allowed. If the check fails, the kernel builds a return message by extracting information from the internal message. The *local_port* component of the message header specifies the reply port (*Ip_null* indicates that no reply should be sent.) In the case of a failed permission check we use the element *Null_reply* of type *KERNEL_REPLY* to represent an empty kernel reply and the kernel returns the special value *Kern_insufficient_permission*. We model this processing with the transition *FailServiceCheck*. The precondition is the existence of a break of type *Bk_have_ruling* indicating that the permission checking transition(s) have occurred. The postcondition is either *Bk_return* in the case where a return message has been requested or *Bk_user_space* for an immediate return to user space.

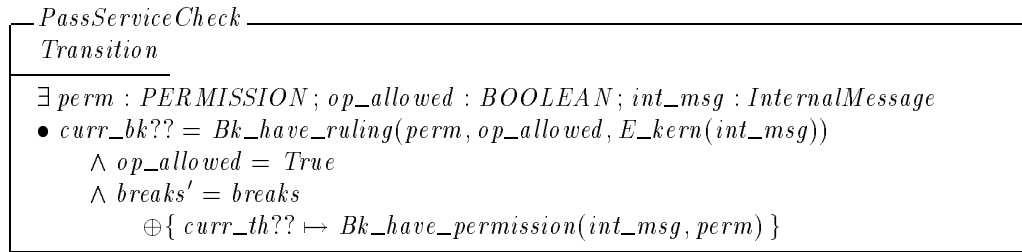
| *Null_reply* : *KERNEL_REPLY*

<i>FailServiceCheck</i>
<i>Transition</i>
$\exists \text{perm} : \text{PERMISSION}; \text{op_allowed} : \text{BOOLEAN}; \text{int_msg} : \text{InternalMessage};$ $\text{reply_to_port} : \text{PORT}; \text{operation} : \text{OPERATION};$ $\text{reply} : \text{KERNEL_REPLY}; \text{return} : \text{KERNEL_RETURN}$ <ul style="list-style-type: none"> • $\text{curr_bk}?? = \text{Bk_have_ruling}(\text{perm}, \text{op_allowed}, \text{E_kern}(\text{int_msg}))$ $\wedge \text{op_allowed} = \text{False}$ $\wedge \text{reply_to_port} = \text{int_msg.header.local_port}$ $\wedge \text{operation} = \text{int_msg.header.operation}$ $\wedge \text{reply} = \text{Null_reply}$ $\wedge \text{return} = \text{Kern_insufficient_permission}$ $\wedge \text{reply_to_port} = \text{Ip_null} \Rightarrow \text{breaks}' = \text{breaks}$ $\quad \oplus \{ \text{curr_th}?? \mapsto \text{Bk_user_space} \}$ $\wedge \text{reply_to_port} \neq \text{Ip_null} \Rightarrow \text{breaks}' = \text{breaks}$ $\quad \oplus \{ \text{curr_th}?? \mapsto \text{Bk_return}(\text{reply_to_port}, \text{operation}, \text{reply}, \text{return}) \}$

Editorial Note:

The schema is not exactly coherent with the rest of the model. This is because the state model consider *local_port* to be a set of ports (zero or one element) rather than allowing it to take null values as it should and as is assumed here.

The transition *PassServiceCheck* models the case were a ruling has been obtained and the operation is allowed. The precondition is the existence of a break of type *Bk_have_ruling*; such a break can only be produced be the permission checking utilities so this ensures that a permission check has occurred. The postcondition is the creation of a new break of type *Bk_have_permission*. The internal message is carried along as a parameter throughout the permission check as an element of type *ENVIRONMENT* (in this case *E_kern(int_msg)*); in this transition we convert back to an element of type *InternalMessage*.



In summary, if the service check passes or has been deferred there will be a break of type *Bk_have_permission*. Further processing is described in the next section.

6.6 Request Validation

The final request processing steps which are generally common to all IPC based kernel requests is validation of the request and extraction of the request parameters from the message body. These transitions are shown in Figure 4.

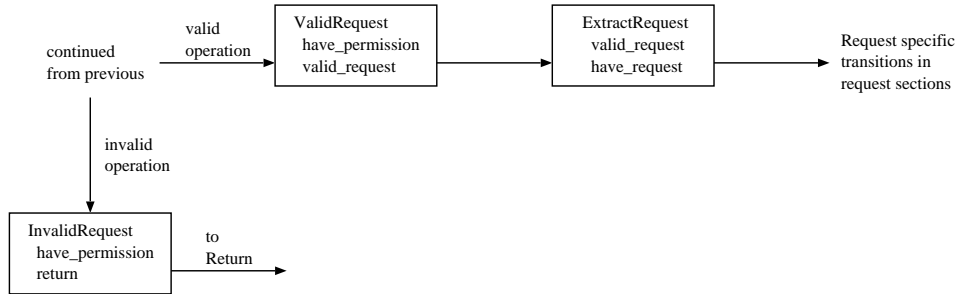


Figure 4: Request Validation

If the service check passes or is deferred the kernel next verifies that the specified operation is an allowed Mach operation. The set *Allowed_kernel_ops* denotes the set of recognized Mach operations.

$$| Allowed_kernel_ops : P OPERATION$$

If the kernel determines that *operation* is not an allowed kernel operation, an error message is generated and sent to the reply port. Again the kernel sends an empty reply, and if the reply port

is null no message is sent. The kernel reply is the special value *MIG_BAD_ID*. This is modeled by the transition *InvalidRequest*. The precondition is a break of type *Bk_have_permission* and the two possible postconditions are *Bk_user_space* for a return without message or *Bk_return* for a return with message.

$ \text{MIG_BAD_ID} : \text{KERNEL_REPLY}$
<hr/> <p style="text-align: center;"><i>InvalidRequest</i> Transition</p> <hr/> $\exists \text{int_msg} : \text{InternalMessage}; \text{perm} : \text{PERMISSION};$ $\text{reply_to_port} : \text{PORT}; \text{operation} : \text{OPERATION};$ $\text{reply} : \text{KERNEL_REPLY}; \text{return} : \text{KERNEL_RETURN}$ <ul style="list-style-type: none"> • $\text{curr_bk??} = \text{Bk_have_permission}(\text{int_msg}, \text{perm})$ <ul style="list-style-type: none"> $\wedge \text{operation} = \text{int_msg.header.operation}$ $\wedge \text{operation} \notin \text{Allowed_kernel_ops}$ $\wedge \text{reply_to_port} = \text{int_msg.header.local_port}$ $\wedge \text{reply} = \text{MIG_BAD_ID}$ $\wedge \text{return} = \text{Kern_invalid_value}$ $\wedge \text{reply_to_port} = \text{Ip_null} \Rightarrow \text{breaks}' = \text{breaks}$ <ul style="list-style-type: none"> $\oplus \{ \text{curr_th??} \mapsto \text{Bk_user_space} \}$ $\wedge \text{reply_to_port} \neq \text{Ip_null} \Rightarrow \text{breaks}' = \text{breaks}$ <ul style="list-style-type: none"> $\oplus \{ \text{curr_th??} \mapsto \text{Bk_return}(\text{reply_to_port}, \text{operation}, \text{reply}, \text{return}) \}$ <hr/>

Editorial Note:

The schema is not exactly coherent with the rest of the model. This is because the state model consider *local_port* to be a set of ports (zero or one element) rather than allowing it to take null values as it should and as is assumed here.

If the request is valid, we model processing with *ValidRequest*. The precondition is the existence of a break of type *Bk_have_permission* and the postcondition *Bk_valid_request* indicates that the operation is allowed.

<hr/> <p style="text-align: center;"><i>ValidRequest</i> Transition</p> <hr/> $\exists \text{int_msg} : \text{InternalMessage}; \text{perm} : \text{PERMISSION}$ <ul style="list-style-type: none"> • $\text{curr_bk??} = \text{Bk_have_permission}(\text{int_msg}, \text{perm})$ <ul style="list-style-type: none"> $\wedge \text{int_msg.header.operation} \in \text{Allowed_kernel_ops}$ $\wedge \text{breaks}' = \text{breaks} \oplus \{ \text{curr_th??} \mapsto \text{Bk_valid_request}(\text{int_msg}) \}$ <hr/>
--

Finally, if the operation is valid, the kernel extracts the request parameters. We model this with the transition *ExtractRequest*. The precondition is the existence of a break of type *Bk_valid_request* and the postcondition *Bk_have_request*, which maintains an element of type *Request*, indicates that a request has been extracted. The components of *Request* are discussed in Section 6.7.3.3.

Editorial Note:

This is modeled as a “black box” conversion. Potentially there are many different extraction transitions, depending on the types of the parameters. At some point it is also important to deal with the possibility that the extraction code (MIG) gets too confused about the types.

ExtractRequest <hr/> Transition <hr/> $\exists \text{int_msg} : \text{InternalMessage}; \text{request} : \text{Request}$ <ul style="list-style-type: none"> • $\text{curr_bk}?? = \text{Bk_valid_request}(\text{int_msg})$ $\wedge \text{breaks}' = \text{breaks} \oplus \{ \text{curr_th}?? \mapsto \text{Bk_have_request}(\text{request}) \}$

Further request processing is discussed in the chapter introductions and request sections.

6.7 Definitions

In this section we define the types and constructors used to describe the break points.

6.7.1 Reply Messages

First we discuss the functions and types connected with kernel reply messages. We have two types to represent the information returned by the kernel in reply messages. The elements of *KERNEL_REPLY* represent the various types of output that the kernel can supply to a client through a reply message. Elements of type *KERNEL_REPLY* are request dependent, so here we define *KERNEL_REPLY* as an abstract set; particular elements are discussed in the request specifications.

[*KERNEL_REPLY*]

The set *KERNEL_RETURN* is an enumerated type representing the possible return statuses that a request can generate. The set of statuses in DTOS consists of:

- *Kern_success* — the request was successful,
- *Kern_failure* — an implementation dependent failure occurred,
- *Kern_invalid_argument* — an attempt was made to perform an operation on the wrong type of entity; for example, an attempt was made to perform a task operation on a thread,
- *Kern_protection_failure* — an attempt was made to access memory in violation of the protections in force,
- *Kern_invalid_address* — an invalid address was specified,
- *Kern_no_space* — an attempt was made to allocate space in a task whose address space or name space was full,
- *Kern_invalid_host* — an attempt was made to perform a host operation on an entity other than a host,
- *Kern_resource_shortage* — insufficient resources were available for service to be provided,
- *Kern_invalid_right* — the wrong type of port right was provided,
- *Kern_invalid_value* — a parameter value that was out of range was provided,
- *Kern_name_exists* — an attempt was made to reuse a name that was already used in the target task's address space,
- *Kern_invalid_name* — a name provided as a port right was not currently in use,
- *Kern_not_in_set* — a name provided as the element of a port set was not in any port set,
- *Kern_urefs_overflow* — an operation was attempted that would cause a user reference count to overflow,
- *Kern_memory_present* —

*Review Note:*Need to determine what this is used for if it is really used

- *Kern_invalid_task* — an attempt was made to perform a task operation on an entity other than a task,
- *Kern_eml_bad_cnt* — an invalid syscall number was specified for an emulation vector entry,
- *Kern_invalid_capability* — a provided name is not a right of the appropriate type,
- *Kern_insufficient_permission* — a security checked failed in the processing of the request.

[*KERNEL_RETURN*]

```

Kern_success : KERNEL_RETURN
Kern_failure : KERNEL_RETURN
Kern_invalid_argument : KERNEL_RETURN
Kern_protection_failure : KERNEL_RETURN
Kern_invalid_address : KERNEL_RETURN
Kern_no_space : KERNEL_RETURN
Kern_invalid_host : KERNEL_RETURN
Kern_resource_shortage : KERNEL_RETURN
Kern_invalid_right : KERNEL_RETURN
Kern_invalid_value : KERNEL_RETURN
Kern_name_exists : KERNEL_RETURN
Kern_invalid_name : KERNEL_RETURN
Kern_not_in_set : KERNEL_RETURN
Kern_urefs_overflow : KERNEL_RETURN
Kern_true : KERNEL_RETURN
Kern_false : KERNEL_RETURN
Kern_memory_present : KERNEL_RETURN
Kern_invalid_task : KERNEL_RETURN
Kern_insufficient_permission : KERNEL_RETURN
Kern_eml_bad_cnt : KERNEL_RETURN
Kern_invalid_capability : KERNEL_RETURN

```

```

Values_disjoint{Kern_success, Kern_failure, Kern_invalid_argument,
  Kern_protection_failure, Kern_invalid_address, Kern_no_space, Kern_invalid_host,
  Kern_resource_shortage, Kern_invalid_right, Kern_invalid_value, Kern_name_exists,
  Kern_invalid_name, Kern_not_in_set, Kern_urefs_overflow,
  Kern_true, Kern_false, Kern_memory_present,
  Kern_eml_bad_cnt, Kern_invalid_task,
  Kern_insufficient_permission, Kern_invalid_capability}

```

Note that all but *Kern_insufficient_permission* are Mach status codes, while *Kern_insufficient_permission* is a DTOS addition.

Next we define several functions which package the return parameters into a message. The expression *Outputs_to_body(reply, return)* converts a group of output parameters *reply* and a status *return* to the message body structure, and the expression *Reply_size(operation, Outputs_to_body(reply, return))* denotes the size of a message carrying *reply* and *return* as output from a request of type *operation*.

$$\left| \begin{array}{l} \text{Outputs_to_body} : \text{KERNEL_REPLY} \times \text{KERNEL_RETURN} \longrightarrow \text{INTERNAL_BODY} \\ \text{Reply_size} : \text{OPERATION} \times \text{INTERNAL_BODY} \longrightarrow \mathbb{N} \end{array} \right.$$

As with general messages in Mach, a reply message can be either simple or complex as specified by the *complex* field of the message header. The expression $\text{Reply_complex}(\text{operation}, \text{Outputs_to_body}(\text{reply}, \text{return}))$ denotes the value that should be assigned to this field when returning *reply* and *return* as output from a request of type *operation*.

$$\left| \text{Reply_complex} : \text{OPERATION} \times \text{INTERNAL_BODY} \longrightarrow \mathbb{P} \text{ COMPLEX_OPTION} \right.$$

The kernel also needs to assign a value to the *operation* field of the reply message. The expression $\text{Reply_op}(\text{operation})$ denotes the value that is used to indicate a reply message for a request of type *operation*.¹³

$$\left| \text{Reply_op} : \text{OPERATION} \twoheadrightarrow \text{OPERATION} \right.$$

6.7.2 Trap Identifiers

We define the set *TRAP_ID* which represents the set of all trap operations. The set of identifiers used to represent traps is *Trap_ids*.

[*TRAP_ID*]

$$\left| \begin{array}{l} \text{Evc_wait_trap}, \text{Mach_thread_self_trap}, \text{Swch_trap}, \text{Swch_pri_trap}, \\ \text{Thread_switch_trap}, \text{Mach_msg_trap} : \text{TRAP_ID} \\ \text{Trap_ids} : \mathbb{P} \text{ TRAP_ID} \\ \hline \langle \text{Evc_wait_trap}, \text{Mach_thread_self_trap}, \text{Swch_trap}, \text{Swch_pri_trap}, \\ \text{Thread_switch_trap}, \text{Mach_msg_trap} \rangle \text{Values_partition } \text{Trap_ids} \end{array} \right.$$

6.7.3 Environment Parameters

As processing of a request progresses, the parameters involved are subject to several transformations. To handle various data contexts we model three types of parameters: user specified parameters, kernel parameters, and abstract request parameters. In this section we define structures which represent these types of data.

As an example, there are three distinct contexts in which a permission check may be required: at the beginning of a system trap (e.g. **mach_thread_self**), at the beginning of an IPC based request (e.g. the service check for **thread_get_state**), or later in the processing of an IPC based request (e.g. the deferred check in **thread_get_special_port**.) As such, the parameter *env* supplied to the permission checking utilities needs to store one of three different types of data. To handle these three cases we define a free type called *ENVIRONMENT*, which uses three different constructor functions to store the three types of parameters.

6.7.3.1 User Parameters In the case of a system trap, we use the structure *UserSpecified*. The components consist of all possible user inputs for the various traps. Most of these inputs come from *Mach_msg_trap*. The components have the following interpretations:

¹³The current implementation defines $\text{Reply_op}(\text{operation})$ to be $\text{operation} + 100$.

- *trap_id* — the identifier of the originating trap,
- *priority* — the priority argument to the **swtch_pri** trap,
- *thread_switch_name* — the name argument to the **thread_switch** trap,
- *thread_switch_option* — the option to the **thread_switch** trap,
- *timeout* — a timeout parameter, used by the **thread_switch** and **mach_msg** traps,
- *message* — the user message being sent via **mach_msg**,
- *options* — the send/receive options specified in the **mach_msg** trap,
- *send_size* — the size of a message being sent,
- *receive_size* — the maximum size message that can be received,
- *receiver* — where return messages will be received,
- *notify* — where to send notifications.

The type *THREAD_SWITCH_OPTION* consists of the following three values:

```
THREAD_SWITCH_OPTION ::= Thread_switch_none | Thread_switch_depress
                        | Thread_switch_wait
```

UserSpecified

```
trap_id : TRAP_ID
priority : Z
thread_switch_name : NAME
thread_switch_option : THREAD_SWITCH_OPTION
timeout : N
message : Message
options : P MACH_MSG_OPTION
send_size : N
receive_size : N
receiver : NAME
notify : NAME
```

6.7.3.2 Kernel Parameters In the case of a service check, the user space parameters contained in an IPC message have been converted into internal representations—names have become ports and virtual memory references have become physical addresses. We model this by storing the relevant processing information in a structure of type *InternalMessage*, which is described in Section 4.10.8.

6.7.3.3 Request Parameters In the case of a deferred check, the kernel has performed additional processing on the message parameters to extract the request parameters. We use an element of type *Request* to represent the parameters of the abstract service being requested. These values are obtained from the contents of the message and the port through which the message was received. The components have the following interpretations:

- *operation* — the type of operation specified in the message,
- *service_port* — the port through which the message was received,
- *pc* — *service_port*'s class (task, thread, memory control, ...),
- *reply_to_port* — a set containing the port to which the reply message should be sent, the empty set indicates no reply should be sent,
- *message* — the received message,
- *sending_sid* — the SID of the subject that sent the message ,

- *receiver_specified* — a Boolean indicating whether an intended receiving SID is specified,
- *receiving_sid* — the intended receiving SID (if any).

BOOLEAN ::= True | False

<i>Request</i> <i>operation</i> : <i>OPERATION</i> <i>service_port</i> : <i>PORT</i> <i>pc</i> : <i>PORT_CLASS</i> <i>reply_to_port</i> : \mathbb{P} <i>PORT</i> <i>message</i> : <i>MESSAGE</i> <i>sending_sid</i> : <i>SSI</i> <i>receiver_specified</i> : <i>BOOLEAN</i> <i>receiving_sid</i> : <i>SSI</i>

6.7.4 Environment

The three preceding data types are used to build the free type *ENVIRONMENT*, which is used to store parameters between transition breaks:

ENVIRONMENT ::= E_user $\langle\langle$ *UserSpecified* $\rangle\rangle$
 | *E_kern* $\langle\langle$ *InternalMessage* $\rangle\rangle$
 | *E_req* $\langle\langle$ *Request* $\rangle\rangle$

6.7.5 Break Status

Finally, we give the formal definition of the free type *BREAK_STATUS*. All of the constructor functions defining *BREAK_STATUS* have been discussed in the preceding sections. They consist of:

- *Bk_return* indicates that a return message is being built and processing is terminating (Section 6.2),
- *Bk_check_pending* indicates that a thread is waiting to check whether a usable ruling for a given permission is available in the cache (Section 6.2),
- *Bk_ruling_pending* indicates that a thread is waiting for a ruling to be received from the Security Server (Section 6.2),
- *Bk_have_ruling* indicates that a thread has obtained a ruling and permission (Section 6.2),
- *Bk_user_space* indicates that a thread is executing in user space (Section 6.3),
- *Bk_new_trap* indicates that a thread has issued a trap into kernel space (Section 6.3),
- *Bk_mach_msg* indicates that a thread has issued a trap of type *Mach_msg_trap* (Section 6.4),
- *Bk_rcv_message* indicates that a thread is waiting to receive a message (Section 6.4),
- *Bk_have_kernel_msg* indicates that a user space message has been converted to an internal kernel message (Section 6.4),
- *Bk_send_message* indicates that a thread is waiting to send a message and the kernel is not the receiver (Section 6.4),
- *Bk_service_pending* indicates that a thread has sent a message to the kernel and is waiting for an IPC based request to be performed (Section 6.5),

- *Bk_have_permission* indicates that a permission check has terminated and the operation is allowed (Section 6.5),
- *Bk_valid_request* indicates that the indicated operation is a recognized Mach operation (Section 6.6),
- *Bk_have_request* indicates that the MIG interface has extracted the internal request parameters (Section 6.6).

The following abbreviations are used for some of the parameters in the constructor functions. Values of type *BK_PERM_REQUEST* indicate the permission check for which a given request is waiting. Values of type *BK_PERM_RESULT* indicate the permission check and its results for a given request. Values of type *BK_RETURN* contain the information necessary to build a return message.

```

BK_PERM_REQUEST == SSI × OSI × PERMISSION × ENVIRONMENT
BK_PERM_RESULT == PERMISSION × BOOLEAN × ENVIRONMENT
BK_RETURN == PORT × OPERATION × KERNEL_REPLY × KERNEL_RETURN

```

The type *BREAK_STATUS* is defined by:

```

BREAK_STATUS ::= Bk_return⟨BK_RETURN⟩
| Bk_check_pending⟨BK_PERM_REQUEST⟩
| Bk_ruling_pending⟨BK_PERM_REQUEST⟩
| Bk_have_ruling⟨BK_PERM_RESULT⟩
| Bk_user_space
| Bk_new_trap⟨Trap_ids × UserSpecified⟩
| Bk_mach_msg⟨UserSpecified⟩
| Bk_rcv_message⟨UserSpecified⟩
| Bk_have_kernel_msg⟨InternalMessage⟩
| Bk_send_message⟨InternalMessage⟩
| Bk_service_pending⟨InternalMessage × OPERATION⟩
| Bk_have_permission⟨InternalMessage × PERMISSION⟩
| Bk_valid_request⟨InternalMessage⟩
| Bk_have_request⟨Request⟩

```

Section 7 System Trap Requests

7.1 Introduction to System Trap Requests

This chapter describes the system trap requests in DTOS.

7.1.1 Kernel Processing

Every trap begins with the transition *Invoke* which produces a new break of type *Bk_new_trap(trap_id?, user_spec?)*. Typically there are no user inputs, but several traps (e.g. **thread_switch**) require one or more inputs. These parameters are stored in components of the *user_spec?* structure. Permission checks are handled by the utilities discussed in Section 6.2.

When a trap produces an error, it does not affect the system state. We define the following transition to describe the situation where trap processing encounters an error and execution leaves kernel space without changing the state.

$$\begin{array}{l}
 \textit{TrapOnlyObserves} \\
 \exists \textit{Dtos} \\
 \textit{Transition} \\
 \hline
 \textit{breaks}' = \textit{breaks} \oplus \{ \textit{curr_th??} \mapsto \textit{Bk_user_space} \}
 \end{array}$$

We now describe the individual system trap requests.

7.1.2 Adding a Send Right

Editorial Note:

The following schema is used both in this chapter and the thread requests chapter. As currently used, it forces us to violate the goal of having a common signature for all state transition schemas. It is not clear that there is a simple way to avoid this, or whether this indicates a weakness in Z or in our use of Z.

The following schema describes a successful addition of a send right with name equal to *name* for the port *port* to the port name space for task *task*. This action is one result of several thread requests.

The name for the new right cannot already name a send-once right, dead name or port set. If it does name either a send or receive right, this right must be for *port*. If it does not already name a send right, then a send right for port *port* with name *name* is added to the port name space for *task* with a user reference count of 1. If it already names a send right belonging to *task*, the user reference count for this send right is incremented by 1.

<p><i>AddSendRight</i></p> <p>Δ <i>PortNameSpace</i></p> <p>Δ <i>SendRightsCount</i></p> <p><i>name</i> : <i>NAME</i></p> <p><i>task</i> : <i>TASK</i></p> <p><i>port</i> : <i>PORT</i></p> <p>$(task, name) \notin so_right \cup dead_namep \cup port_set_namep$</p> <p>$(task, name) \in s_r_right \Rightarrow named_port(task, name) = port$</p> <p>$(task, name) \notin s_right$ $\Rightarrow port_right_rel' = port_right_rel \cup \{(task, port, name, Send, 1)\}$</p> <p>$(task, name) \in s_right$ $\Rightarrow port_right_rel' = port_right_rel$ $\quad \setminus \{(task, port, name, Send, s_right_ref_count(task, name))\}$ $\quad \cup \{(task, port, name, Send, s_right_ref_count(task, name) + 1)\}$</p>

7.2 mach_thread_self

The request **mach_thread_self** places a send right for a thread's kernel port in the name space of the owning task of the thread. It is a system trap.

7.2.1 Kernel Interface

```
mach_port_t mach_thread_self
    0;
```

7.2.1.1 Input Parameters No input parameters are provided to the kernel for a **mach_thread_self** request.

7.2.1.2 Output Parameters The following output parameters are returned by the kernel for a **mach_thread_self** request:

- *kernel_port_name!* — the name for a send right to the thread's kernel port in the port name space of the thread's owning task

<p><i>MachThreadSelf Outputs</i></p> <p><i>kernel_port_name!</i> : <i>NAME</i></p>
--

7.2.2 Request Criteria

The following criteria are defined for the **mach_thread_self** request.

- **C1** — Permission *Get_thread_kernel_port* has been obtained.

<i>C1Mach ThreadSelfGetThreadKernelPort</i>
<i>Transition</i>
$\exists env : ENVIRONMENT$ $\bullet curr_bk?? = Bk_have_ruling(Get_thread_kernel_port, True, env)$

<i>NotC1Mach ThreadSelfGetThreadKernelPort</i>
<i>Transition</i>
$\exists env : ENVIRONMENT$ $\bullet curr_bk?? = Bk_have_ruling(Get_thread_kernel_port, False, env)$

- **C2** — The sself port for the current thread is not *Ip_dead*.

<i>C2Mach ThreadSelfKernelPortNotDead</i>
<i>Transition</i>
$thread_sself(curr_th??) \neq Ip_dead$

NotC2Mach ThreadSelfKernelPortNotDead
 $\hat{=} Transition \wedge \neg C2Mach ThreadSelfKernelPortNotDead$

- **C3** — The sself port for the current thread is not *Ip_null*.

<i>C3Mach ThreadSelfKernelPortNotNull</i>
<i>Transition</i>
$thread_sself(curr_th??) \neq Ip_null$

NotC3Mach ThreadSelfKernelPortNotNull
 $\hat{=} Transition \wedge \neg C3Mach ThreadSelfKernelPortNotNull$

- **C4** — Either the current task already holds a send or receive right to the thread's sself port (so that a new IPC entry need not be created), or the kernel has the available resources to create an IPC entry in the current task's name space. We do not actually model the consumption of resources by the kernel. So, we will use the set *Resources_available_to_create_ipc_entry* to indicate the set of states where there are sufficient resources to create an IPC entry.

| *Resources_available_to_create_ipc_entry* : $\mathbb{P} DtosExec$

<i>C4Mach ThreadSelfResourcesAvailable</i>
<i>Transition</i>
$((\exists name : NAME; i : \mathbb{N}_1; right : \{Send, Receive\})$ $\bullet (curr_task??, thread_sself(curr_th??), name, right, i)$ $\in \underline{port_right_rel})$ $\vee \theta DtosExec \in Resources_available_to_create_ipc_entry)$

NotC4Mach ThreadSelfResourcesAvailable
 $\hat{=} Transition \wedge \neg C4Mach ThreadSelfResourcesAvailable$

■ **C5** — Permission *Hold_send* has been obtained.

<i>C5MachThreadSelfHoldSend</i> <i>Transition</i>
$\exists env : ENVIRONMENT$ <ul style="list-style-type: none"> • $curr_bk?? = Bk_have_ruling(Hold_send, True, env)$

<i>NotC5MachThreadSelfHoldSend</i> <i>Transition</i>
$\exists env : ENVIRONMENT$ <ul style="list-style-type: none"> • $curr_bk?? = Bk_have_ruling(Hold_send, False, env)$

7.2.3 Return Values

Table 1 describes the values returned at the completion of the request and the conditions under which each value is returned. We note that **C2** and **C3** may not both be false simultaneously.

<i>kernel_port_name!</i>	C1	C2	C3	C4	C5
<i>name</i>	T	T	T	T	T
<i>Mach_port_null</i>	T	T	T	T	F
<i>Mach_port_null</i>	T	T	T	F	-
<i>Mach_port_null</i>	T	T	F	-	-
<i>Mach_port_dead</i>	T	F	T	-	-
<i>Return_status_to_name(Kern_insufficient_permission)</i>	F	-	-	-	-

Table 1: Return Values for **mach_thread_self**

We call attention here to the fact that **C1** – **C4** are checked in the the transition *MachThreadSelfMiddle* (see below) while only **C5** is checked in transition *MachThreadSelfEnd*. These are distinct transitions between which an arbitrary number of other transitions may occur. Furthermore, **C5** is checked only if in some earlier transition for that trap request **C1** – **C4** are all true. We also note that *name* is constrained by *AddSendRight* when *RVMachThreadSelfGood* is combined with *MachThreadSelfState*.

<i>RVMachThreadSelfGood</i> <i>C5MachThreadSelfHoldSend</i> <i>MachThreadSelf Outputs</i> <i>Transition</i> <i>name : NAME</i>
$kernel_port_name! = name$ $breaks' = breaks \oplus \{ curr_th?? \mapsto Bk_user_space \}$

<i>RVMachThreadSelfNoHoldSend</i> <i>NotC5MachThreadSelfHoldSend</i> <i>MachThreadSelf Outputs</i>
$kernel_port_name! = Mach_port_null$

<i>RV</i> MachThreadSelfResourceShortage <i>C1</i> MachThreadSelfGetThreadKernelPort <i>C2</i> MachThreadSelfKernelPortNotDead <i>C3</i> MachThreadSelfKernelPortNotNull <i>NotC4</i> MachThreadSelfResourcesAvailable MachThreadSelf Outputs
<i>kernel_port_name!</i> = <i>Mach_port_null</i>

<i>RV</i> MachThreadSelfKernelPortNull <i>C1</i> MachThreadSelfGetThreadKernelPort <i>NotC3</i> MachThreadSelfKernelPortNotNull MachThreadSelf Outputs
<i>kernel_port_name!</i> = <i>Mach_port_null</i>

<i>RV</i> MachThreadSelfKernelPortDead <i>C1</i> MachThreadSelfGetThreadKernelPort <i>NotC2</i> MachThreadSelfKernelPortNotDead MachThreadSelf Outputs
<i>kernel_port_name!</i> = <i>Mach_port_dead</i>

| *Return_status_to_name* : *KERNEL_RETURN* \rightarrow *NAME*

<i>RV</i> MachThreadSelfNoGetThreadKernelPort <i>NotC1</i> MachThreadSelfGetThreadKernelPort MachThreadSelf Outputs
<i>kernel_port_name!</i> = <i>Return_status_to_name</i> (<i>Kern_insufficient_permission</i>)

7.2.4 State Changes

A successful ***mach_thread_self*** request results in the addition of a send right for the kernel port, *port*, of *curr_th??* into the port name space of the task, *task*, containing *curr_th??*. This is accomplished by schema *AddSendRight* which describes the relationships between *name*, *task*, and *port* when a send right is successfully added to a name space.

<i>MachThreadSelfState</i> <hr/> <i>Transition</i> <i>ThreadInvariants</i> \exists <i>ThreadExist</i> \exists <i>PortExist</i> \exists <i>TaskExist</i> \exists <i>DeadRights</i> \exists <i>Threads</i> \exists <i>SpecialThreadPorts</i> \exists <i>ThreadAndProcessorSet</i> <i>AddSendRight</i> <hr/> <i>task</i> = <i>curr_task??</i> <i>port</i> = <i>thread_sself(curr_th??)</i> <hr/> $\underline{port_set_rel}' = \underline{port_set_rel}$

7.2.5 Complete Request

Here we discuss the transitions shown in Figure 5 which describe the general form of the ***mach_thread_self*** request.

1. A ***mach_thread_self*** request is invoked through a system trap that has the *trap_id?* field set to *Mach_thread_self_trap*.

<i>InvokeMachThreadSelf</i> <hr/> <i>Invoke</i> <hr/> <i>trap_id?</i> = <i>Mach_thread_self_trap</i>
--

2. *MachThreadSelfPermCheckGTKP* suspends processing to wait for the availability of a ruling on the *Get_thread_kernel_port* permission. The permission check is handled by the utilities described in Section 6.2 and consists of one or two transitions depending upon the availability of a ruling in the cache.

Editorial Note:

This should eventually be covered by generic trap processing if possible.

<i>MachThreadSelfPermCheckGTKP</i> <hr/> <i>Transition</i> <hr/> \exists <i>CheckPending</i> ; <i>user_spec</i> : <i>UserSpecified</i> <ul style="list-style-type: none"> • <i>curr_bk??</i> = <i>Bk_new_trap(Mach_thread_self_trap, user_spec)</i> \wedge <i>ssi</i> = <i>thread_sid(curr_th??)</i> \wedge <i>osi</i> = <i>Thread_self_sid</i> \wedge <i>breaks'</i> = <i>breaks</i> $\oplus \{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Get_thread_kernel_port, E_user(user_spec)) \}$

3. Either *MachThreadSelfMiddleBad* or *MachThreadSelfMiddleGood* occurs. This models a transition where a ruling on the *Get_thread_kernel_port* permission has been obtained

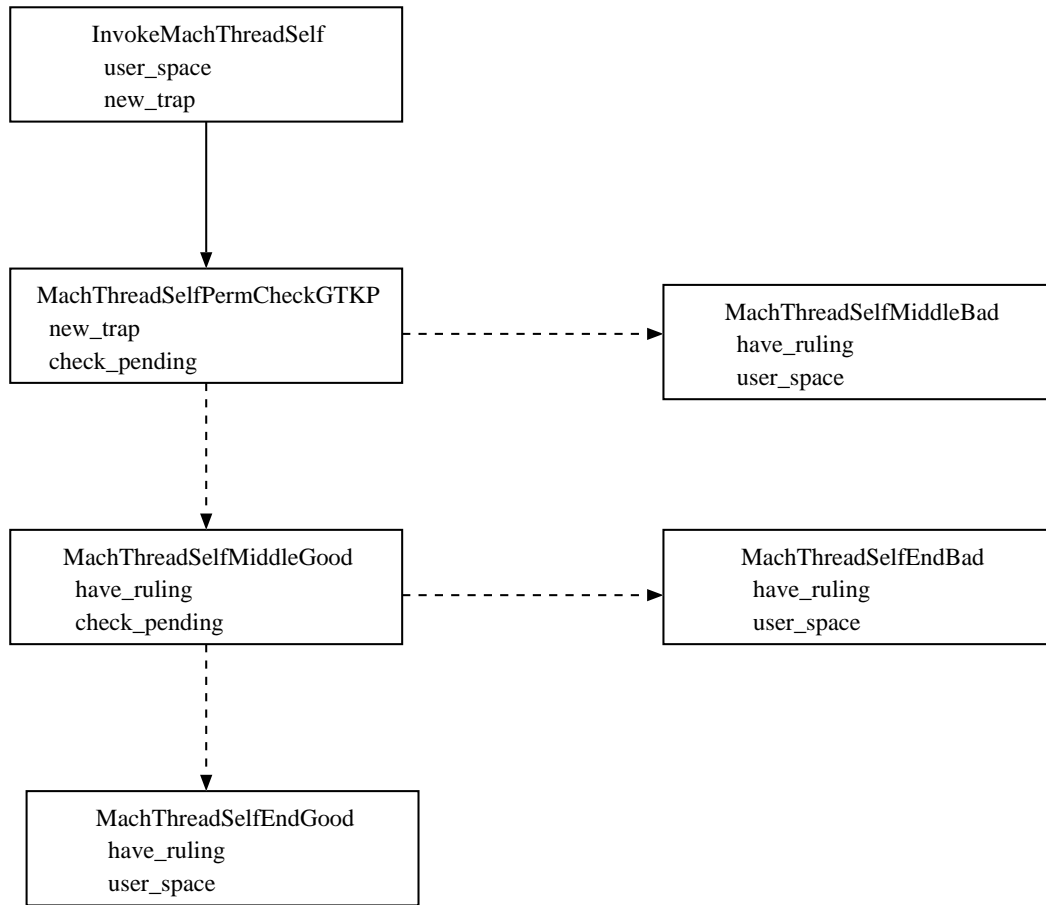


Figure 5: **mach_thread_self** Processing

from the cache or Security Server. Conditions **C1–C4** are checked. If any of them is false, no state changes are made, an error value is returned and processing terminates with this step.

$$\begin{aligned}
 & MachThreadSelfMiddleBad \\
 & \hat{=} (RVMachThreadSelfResourceShortage \\
 & \quad \vee RVMachThreadSelfKernelPortNull \\
 & \quad \vee RVMachThreadSelfKernelPortDead \\
 & \quad \vee RVMachThreadSelfNoGetThreadKernelPort) \\
 & \wedge TrapOnlyObserves
 \end{aligned}$$

If the four conditions are all true, then trap processing is suspended (via *MachThreadSelfPermCheckHS*) while waiting to check on the availability of a ruling for the *Hold_send* permission.

Editorial Note:

The value used for *osi* in the following schema is not in agreement with the requirements in the FSPM which state that the *osi* should be *port_sid(port)*. However, this is the value that is being used in the prototype as of 17 May 1996. CAR# 5024 describes this discrepancy.

Also, env should not be a free variable but should be explicitly specified, in this case most likely as the “empty” environment.

<p style="margin: 0;"><i>MachThreadSelfPermCheckHS</i></p> <hr/> <p style="margin: 0;"><i>Transition</i></p> <hr/> <p style="margin: 0;">\exists <i>CheckPending</i>; $env : ENVIRONMENT$; $port : PORT$</p> <ul style="list-style-type: none"> • $ssi = thread_sid(curr_th??)$ <li style="padding-left: 20px;">$\wedge port = thread_self(curr_th??)$ <li style="padding-left: 20px;">$\wedge osi = \mathbf{if} (curr_th??, port) \in thread_self$ <li style="padding-left: 40px;">then <i>Thread_self_sid</i> <li style="padding-left: 20px;">else if $(curr_task??, port) \in task_self$ <li style="padding-left: 40px;">then <i>Task_self_sid</i> <li style="padding-left: 20px;">else $port_sid(port)$ <p style="margin: 0;">$\wedge breaks' = breaks$</p> <p style="margin: 0;">$\oplus \{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Hold_send, env) \}$</p>

MachThreadSelfMiddleGood
 $\hat{=} C1MachThreadSelfGetThreadKernelPort$
 $\wedge C2MachThreadSelfKernelPortNotDead$
 $\wedge C3MachThreadSelfKernelPortNotNull$
 $\wedge C4MachThreadSelfResourcesAvailable$
 $\wedge MachThreadSelfPermCheckHS$

4. Either *MachThreadSelfEndGood* or *MachThreadSelfEndBad* occurs. This models a transition where a ruling on the *Hold_send* permission has been obtained from the cache or the Security Server. The *Hold_send* permission is checked. If it is granted, the state changes in *MachThreadSelfState* occur, and the name of the new send right is returned.

$MachThreadSelfEndGood \hat{=} RVMachThreadSelfGood \wedge MachThreadSelfState$

Otherwise, *Mach_port_null* is returned and no state changes occur.

$MachThreadSelfEndBad \hat{=} RVMachThreadSelfNoHoldSend \wedge TrapOnlyObserves$

Section 8

Port Requests

8.1 Introduction to Port Requests

This chapter describes the port kernel requests in DTOS.

8.1.1 Request Identifiers

We first define the identifier that is used to represent each port request. The kernel accepts all port requests through task self ports.

```

Mach_port_allocate_id, Mach_port_allocate_secure_id, Mach_port_allocate_name_id,
Mach_port_allocate_name_secure_id, Mach_port_deallocate_id,
Mach_port_destroy_id, Mach_port_extract_right_id,
Mach_port_get_receive_status_id, Mach_port_get_refs_id,
Mach_port_get_set_status_id, Mach_port_insert_right_id,
Mach_port_mod_refs_id, Mach_port_move_member_id, Mach_port_names_id,
Mach_port_rename_id, Mach_port_request_notification_id,
Mach_port_set_mscount_id, Mach_port_set_qlimit_id, Mach_port_set_seqno_id,
Mach_port_type_id, Mach_port_type_secure_id : OPERATION

```

```

Port_operations : P OPERATION

```

```

{Mach_port_allocate_id, Mach_port_allocate_secure_id, Mach_port_allocate_name_id,
Mach_port_allocate_name_secure_id, Mach_port_deallocate_id,
Mach_port_destroy_id, Mach_port_extract_right_id,
Mach_port_get_receive_status_id, Mach_port_get_refs_id,
Mach_port_get_set_status_id, Mach_port_insert_right_id,
Mach_port_mod_refs_id, Mach_port_move_member_id, Mach_port_names_id,
Mach_port_rename_id, Mach_port_request_notification_id,
Mach_port_set_mscount_id, Mach_port_set_qlimit_id, Mach_port_set_seqno_id,
Mach_port_type_id, Mach_port_type_secure_id}

```

```

Values_partition Port_operations

```

```

Port_operations ⊆ Allowed_mach_services(Pc_task)

```

8.1.2 Required Permissions

For each operation there is a primary permission that is required to perform the operation. We define here the portion of the *Required_permission* function that pertains to port requests.

```
{(Mach_port_allocate_id, Add_name),
 (Mach_port_get_receive_status_id, Observe_pns_info),
 (Mach_port_get_refs_id, Observe_pns_info),
 (Mach_port_get_set_status_id, Observe_pns_info),
 (Mach_port_names_id, Observe_pns_info),
 (Mach_port_rename_id, Port_rename),
 (Mach_port_set_mscount_id, Alter_pns_info),
 (Mach_port_set_qlimit_id, Alter_pns_info),
 (Mach_port_set_seqno_id, Alter_pns_info)}
⊂ Required_permission
```

8.1.3 Invariant Information

Review Note:

This section will be completed in a future draft when all of the port requests are completed.

8.1.4 General Information

This section contains bits of information which are common to several port requests.

Review Note:

In a future draft, it may be appropriate to move some of this information to the state chapter.

The names *Mach_port_dead* and *Mach_port_null* are referred to in this chapter as *Reserved_names*.

<i>Reserved_names</i> : \mathbb{P} NAME <i>Reserved_names</i> = { <i>Mach_port_dead</i> , <i>Mach_port_null</i> } <i>Mach_port_null</i> \neq <i>Mach_port_dead</i>
--

Parameters to several of the port requests include a type of right which includes not only send, send-once and receive rights but also dead rights and port sets. These parameters are given the type *RIGHT_TYPE*.

[*RIGHT_TYPE*]

<i>Mach_port_right_send</i> , <i>Mach_port_right_receive</i> , <i>Mach_port_right_send_once</i> , <i>Mach_port_right_port_set</i> , <i>Mach_port_right_dead_name</i> : <i>RIGHT_TYPE</i> <i>Values_disjoint</i> (<i>Mach_port_right_send</i> , <i>Mach_port_right_receive</i> , <i>Mach_port_right_send_once</i> , <i>Mach_port_right_port_set</i> , <i>Mach_port_right_dead_name</i>)
--

The **mach_port_request_notification** request has a parameter of type *MACH_MSG_ID* which indicates which of the three types of notification is being requested.

[*MACH_MSG_ID*]

<p><i>Mach_notify_port_destroyed</i>, <i>Mach_notify_no_senders</i>, <i>Mach_notify_dead_name</i> : <i>MACH_MSG_ID</i></p> <p><i>Values_disjoint</i>(<i>Mach_notify_port_destroyed</i>, <i>Mach_notify_no_senders</i>, <i>Mach_notify_dead_name</i>)</p>
--

The **mach_port_get_receive_status** request returns a record *PortStatus* consisting of the following information:

- *port_set_name* — if the receive right is a member of a port set, the name of this port set; otherwise, the name *Mach_port_null*
- *make_send_count_value* — the make-send count for the port
- *port_destroyed_notification_requested* — takes on the value *True* if a port-destroyed notification request is currently active for the port; otherwise *False*
- *no_more_senders_notification_requested* — takes on the value *True* if a no-more-senders notification request is currently active for the port; otherwise *False*
- *msg_count_value* — the number of messages queued at the port
- *qlimit_value* — the limit on the number of messages which can be queued to the port
- *sequence_no_value* — the current sequence number for the port
- *number_of_send_once_rights* — the number of send-once rights which exist to the port
- *any_send_rights* — takes on the value *True* if there exist send rights to the port; otherwise *False*

<p><i>PortStatus</i></p> <p><i>port_set_name</i> : <i>NAME</i></p> <p><i>make_send_count_value</i> : \mathbb{N}</p> <p><i>port_destroyed_notification_requested</i> : <i>BOOLEAN</i></p> <p><i>no_more_senders_notification_requested</i> : <i>BOOLEAN</i></p> <p><i>msg_count_value</i> : \mathbb{N}</p> <p><i>qlimit_value</i> : \mathbb{N}</p> <p><i>sequence_no_value</i> : \mathbb{Z}</p> <p><i>number_of_send_once_rights</i> : \mathbb{N}</p> <p><i>any_send_rights</i> : <i>BOOLEAN</i></p>
--

The **mach_port_names**, **mach_port_type**, and **mach_port_type_secure** requests return a mask *PortTypeMask* consisting of the following flags:

- *mach_port_type_send* — equal to *True* if and only if the name is a send right
- *mach_port_type_receive* — equal to *True* if and only if the name is a receive right
- *mach_port_type_send_once* — equal to *True* if and only if the name is a send-once right
- *mach_port_type_port_set* — equal to *True* if and only if the name is a port set
- *mach_port_type_dead_name* — equal to *True* if and only if the name is a dead name
- *mach_port_type_dead_name_request* — equal to *True* if and only if there is an outstanding dead name notification request for the name
- *mach_port_type_msg_accepted_request* — equal to *True* if and only if the name is not available for use as a send right since it has been used to force a message on a message queue

```

PortTypeMask
mach_port_type_send : BOOLEAN
mach_port_type_receive : BOOLEAN
mach_port_type_send_once : BOOLEAN
mach_port_type_port_set : BOOLEAN
mach_port_type_dead_name : BOOLEAN
mach_port_type_dead_name_request : BOOLEAN
mach_port_type_msg_accepted_request : BOOLEAN

```

8.1.5 Parameter Packaging Functions

When invoking a kernel request, the following functions package the input parameters into a message body:

```

Mach_port_allocate_inputs_to_body : RIGHT_TYPE → MESSAGE_BODY
Mach_port_extract_right_inputs_to_body
    : NAME × MACH_MSG_TYPE → MESSAGE_BODY
Mach_port_get_receive_status_inputs_to_body : NAME → MESSAGE_BODY
Mach_port_get_refs_inputs_to_body : NAME × RIGHT_TYPE → MESSAGE_BODY
Mach_port_get_set_status_inputs_to_body : NAME → MESSAGE_BODY
Mach_port_insert_right_inputs_to_body
    : NAME × NAME × MACH_MSG_TYPE → MESSAGE_BODY
Mach_port_move_member_inputs_to_body : NAME × NAME → MESSAGE_BODY
Mach_port_rename_inputs_to_body : NAME × NAME → MESSAGE_BODY
Mach_port_request_notification_inputs_to_body
    : NAME × MACH_MSG_ID × ℕ × NAME × MACH_MSG_TYPE → MESSAGE_BODY
Mach_port_set_mscount_inputs_to_body : NAME × ℕ → MESSAGE_BODY
Mach_port_set_qlimit_inputs_to_body : NAME × ℕ → MESSAGE_BODY
Mach_port_set_seqno_inputs_to_body : NAME × ℕ → MESSAGE_BODY

```

When creating a reply message from a request, the following functions package the output parameters into a kernel reply:

```

Mach_port_allocate_outputs_to_reply : NAME → KERNEL_REPLY
Mach_port_extract_right_outputs_to_reply
    : PORT × MACH_MSG_TYPE → KERNEL_REPLY
Mach_port_get_receive_status_outputs_to_reply : PortStatus → KERNEL_REPLY
Mach_port_get_refs_outputs_to_reply : ℕ → KERNEL_REPLY
Mach_port_get_set_status_outputs_to_reply : ℙ NAME × ℕ → KERNEL_REPLY
Mach_port_names_outputs_to_reply
    : seq NAME × ℕ × seq PortTypeMask × ℕ → KERNEL_REPLY
Mach_port_request_notification_outputs_to_reply : PORT → KERNEL_REPLY

```

When receiving a reply message from the kernel the following functions unpack the message body to obtain the output parameters (including the return status):

$Body_to_mach_port_allocate_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN \times NAME$
$Body_to_mach_port_extract_right_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN \times NAME \times MACH_MSG_TYPE$
$Body_to_mach_port_get_receive_status_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN \times PortStatus$
$Body_to_mach_port_get_refs_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN \times \mathbb{N}$
$Body_to_mach_port_get_set_status_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN \times \mathbb{P} NAME \times \mathbb{N}$
$Body_to_mach_port_insert_right_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN$
$Body_to_mach_port_move_member_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN$
$Body_to_mach_port_names_outputs$ $: MESSAGE_BODY$ $\rightarrow KERNEL_RETURN \times seq NAME \times \mathbb{N} \times seq PortTypeMask \times \mathbb{N}$
$Body_to_mach_port_rename_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN$
$Body_to_mach_port_request_notification_outputs$ $: MESSAGE_BODY \rightarrow KERNEL_RETURN \times NAME$
$Body_to_mach_port_set_mscount_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN$
$Body_to_mach_port_set_qlimit_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN$
$Body_to_mach_port_set_seqno_outputs : MESSAGE_BODY \rightarrow KERNEL_RETURN$

8.1.6 Kernel Processing

The initial kernel processing of any request, when removing the request from the bag of validated requests, is described by the *ProcessRequest* schema in Section 6. In this section, we consider additional initial processing which is shared by all port requests.

The service port through which a port request is received must be the self port for some task. If it is not, then the request immediately terminates and returns *Kern_invalid_task*.

$ProcessPortRequestBad$
$ProcessRequest$ $SpecialTaskPorts$ $return! : KERNEL_RETURN$
$operation? \in Port_operations$ $service_port? \notin dom self_task$ $return! = Kern_invalid_task$

$$PortRequestBad \hat{=} ProcessPortRequestBad \gg RequestNoOp$$

Otherwise, the task owning the self port is identified and the kernel continues to process the request.

<i>ProcessPortRequestGood</i> <i>ProcessRequest</i> <i>SpecialTaskPorts</i> <i>task? : TASK</i>
<i>operation? ∈ Port_operations</i> <i>service_port? ∈ dom self_task</i> <i>task? = self_task(service_port?)</i>

Review Note:

Several of the port requests have the potential of returning *Kern_resource_shortage*. The model cannot accurately reflect when this is returned, so the original plans were to include another version of the "bad" schema here which dealt with those cases. However, this has not been done.

Review notes have been added to each of the individual requests for which *Kern_resource_shortage* is a possible return value.

8.1.7 Notifications

Several of the port requests can result in the sending of notifications, messages which inform about some change to the state of a port or port right. There are six kinds of notifications:

dead-name A dead-name notification for a port right is registered by ***mach_msg*** or ***mach_port_request_notification***. The notification is sent if the port right becomes dead due to the destruction of the port.

port-deleted A port-deleted notification is registered whenever a dead-name notification is registered. The notification is sent if the port right becomes unusable due to the right itself being destroyed or moved.

msg-accepted A msg-accepted notification for a port right is registered when the port right is used to forcibly enqueue a message on a port using ***mach_msg***. The notification is sent when the message is removed from the queue (either to be received or destroyed).

no-senders A no-senders notification for a port is registered using ***mach_port_request_notification***. The notification is sent when the last send or send-once right for the port is destroyed (with one exception depending upon the parameter *sync* to ***mach_port_request_notification***; see Section 8.8).

port-destroyed A port-destroyed notification for a port is registered using ***mach_port_request_notification***. The notification is sent when the port would otherwise be destroyed, and contains a receive right for the port, thus saving it from destruction.

send-once A send-once notification for a port is sent whenever a send-once right is destroyed without being used to send a message.

Note that the kernel cannot guarantee to send notifications. It can fail to send a notification, for instance if it determines that an attempt to send the notification could result in an infinite loop.

Editorial Note:

Port-destroyed notifications could interact in interesting ways with the kernel's ability to identify "circularities" of receive rights in transit.

A circularity in receive rights exists when the receive right for port 1 is contained in a message destined for port 2, the receive right for port 2 is contained in a message destined for port 3, ... and the receive right for port N is contained in a message for port 1. If this occurs, then none of the messages can be received and the kernel sets out to clean them all up.

However, if the kernel cleans out a message and finds a receive right with a registered port-destroyed notification port, it instead sends the receive right to that port, saving the right from destruction. This breaks the circle, so other messages in the circle may once again be reachable.

We have not spent any time trying to determine how this interesting case is handled. It should be considered when port-destroyed notifications are modeled.

Sending of each kind of notification is modeled in the following sections.

Review Note:

For now, only dead-name notifications will be considered. When other types of notifications are considered, processing which is common to multiple notifications should be moved to this introductory section.

8.1.7.1 Dead Name Notifications The internal message header for a dead-name notification message is filled in as follows:

- The *local_port* and *local_rights* fields are empty since there is no expected reply from a dead-name notification.
- The *remote_port* and *remote_rights* fields contain a send-once right to the notification port.
- The *size* field contains the size of a dead-name notification, which is a constant.
- The *msg_sequence_no* field is initialized to zero. This field is ignored until the message is received.
- The *operation* field contains the identifier *Ip_notify_dead_name_id*.
- The *complex* field is empty since the body of a dead-name notification contains no port rights or out-of-line memory.

$$\begin{aligned} \text{Dead_name_notification_header} &: \text{PORT} \rightarrow \text{MachInternalHeader} \\ \text{Dead_name_notification_size} &: \mathbb{N} \end{aligned}$$

$$\forall \text{port} : \text{PORT}$$

- (let $\text{header} == \text{Dead_name_notification_header}(\text{port})$
 - $\text{header.local_port} = \emptyset$
 - $\wedge \text{header.local_rights} = \emptyset$
 - $\wedge \text{header.remote_port} = \text{port}$
 - $\wedge \text{header.remote_rights} = \text{Mach_msg_type_port_send_once}$
 - $\wedge \text{header.size} = \text{Dead_name_notification_size}$
 - $\wedge \text{header.msg_sequence_no} = 0$
 - $\wedge \text{header.operation} = \text{Ip_notify_dead_name_id}$
 - $\wedge \text{header.complex} = \emptyset$)

The body of a dead-name notification consists of a single element containing the name of the now dead port right.

Review Note:

I tried to specify the body of the message explicitly in terms of the model, but was unable to do so. The lack of a direct map between the model and the code caused some difficulty, but the final nail was the fact that the model requires a task to be associated with each message element, and this made no sense in the current situation.

| $Dead_name_notification_body : NAME \rightarrow INTERNAL_BODY$

The entire dead-name notification message contains a header and body as just described.

Review Note:

Once again, I would like to fill in the other values from the schema *InternalMessage*, but these values do not correspond to anything in the code and it is difficult to determine what values to choose.

| $Dead_name_notification_message : PORT \times NAME \rightarrow InternalMessage$
 | $\forall port : PORT; name : NAME$
 | $\bullet (Dead_name_notification_message(port, name)).header = Dead_name_notification_header(port)$
 | $\wedge (Dead_name_notification_message(port, name)).body = Dead_name_notification_body(name)$

The following occurs when a dead name notification is sent:

- A new message (*new_message*) is created, whose contents are give by $Dead_name_notification_message(notify_port, dead_right_name)$.
- *new_message* is added to the queue for *notify_port*.
- The sending SID for *new_message* is set to reflect the kernel acting as *current_task*.
- No receiving, sending SID or access vector is specified for the message.

Review Note:

Failure can occur here if the kernel is unable to allocate memory for the message. In this case the notification message is not sent and the kernel simply continues processing as if the message were successfully sent.

Review Note:

There is a permission check here which has not been specified. This check occurs after the message is allocated, however, in the current prototype the message is not deallocated when the permission check fails. The kernel simply continues processing as if the message were successfully sent.

This permission check will need to be added later when the execution model is updated. The *Ruling* associated with the message will also need to be added at that time.

There is also a much bigger list of invariants that could be added to the following schema.

Review Note:

The code also has to make sure that *notify_port* is not *Ip_dead*. This should come for free whenever the schema is used in this chapter, but that should be doublechecked. If *notify_port* is *Ip_dead*, then the message is destroyed.

QueueDeadNameNotification

Δ *Messages*
 Δ *DtosMessages*
KernelAs
current_task : *TASK*
notify_port : *PORT*
dead_right_name : *NAME*

\exists *new_message* : *MESSAGE*

- (*new_message* \notin *message_exists*
 - \wedge *message_exists'* = *message_exists* \cup {*new_message*}
 - \wedge *msg_contents'* = *msg_contents*
 - \oplus {*new_message* \mapsto *Dead_name_notification_message*(*notify_port*, *dead_right_name*)}
 - \wedge *message_in_port_rel'* = *message_in_port_rel*
 - \oplus {*notify_port* \mapsto *message_in_port_rel*(*notify_port*) \cap {*new_message*}
 - \wedge *msg_sending_sid'* = *msg_sending_sid*
 - \oplus {*new_message* \mapsto *kernel_as*(*current_task*)}

msg_receiving_sid' = *msg_receiving_sid*
msg_specified_sid' = *msg_specified_sid*
msg_specified_vector' = *msg_specified_vector*

Review Note:

It's worth noting that the recipient of a dead-name notification only receives the name of the dead right, and no indication of which ipc name space contains the dead right.

8.2 mach_port_allocate

A **mach_port_allocate** request creates a new receive right, port set, or dead name in a task's name space.

8.2.1 Client Interface

kern_return_t	mach_port_allocate	
	(mach_port_t	<i>task_name</i> ,
	mach_port_right_t	<i>right_type</i> ,
	mach_port_t*	<i>new_name</i>);
kern_return_t	mach_port_allocate_secure	
	(mach_port_t	<i>task_name</i> ,
	mach_port_right_t	<i>right_type</i> ,
	mach_port_t*	<i>new_name</i> ,
	security_id_t	<i>obj_sid</i>);

Review Note:

No attempt has been made to integrate **mach_port_allocate_secure** into this specification.

8.2.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_allocate** request:

- *task_name?* — the client's name for the task in whose name space the new receive right, port set or dead name is created
- *right_type?* — the type of right to be created, either *Mach_port_right_receive*, *Mach_port_right_port_set*, or *Mach_port_right_dead_name*.

```

MachPortAllocate ClientInputs
-----
task_name? : NAME
right_type? : RIGHT_TYPE

```

A **mach_port_allocate** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_allocate_id* and has a body consisting of *right_type?*.

```

Invoke MachPortAllocate
-----
Invoke MachMsg
MachPortAllocate ClientInputs
-----
name? = task_name?
operation? = Mach_port_allocate_id
msg_body = Mach_port_allocate_inputs_to_body(right_type?)

```

8.2.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_allocate** request:

- *return!* — the status of the request
- *new_name!* — the name of the new right

```

MachPortAllocate ClientOutputs
-----
return! : KERNEL_RETURN
new_name! : NAME

```

```

MachPortAllocate ReceiveReply
-----
Invoke MachMsgRcv
MachPortAllocate ClientOutputs
-----
(return!, new_name!) = Body_to_mach_port_allocate_outputs(msg_body)

```

8.2.2 Kernel Interface

8.2.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_allocate** request:

- *task?* — the task known to the client by *task_name?*

- *right_type?* — provided by the client

<i>MachPortAllocateInputs</i> <i>task?</i> : <i>TASK</i> <i>right_type?</i> : <i>RIGHT_TYPE</i>

8.2.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_allocate** request:

- *return!* — the status of the request
- *new_name!* — the name of the new right

<i>MachPortAllocateOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i> <i>new_name!</i> : <i>NAME</i>

Upon completion of the processing of a **mach_port_allocate** request, a reply message is built from the output parameters.

<i>MachPortAllocateReply</i> <i>RequestReturn</i> <i>new_name?</i> : <i>NAME</i>
<i>reply?</i> = <i>Mach_port_allocate_outputs_to_reply(new_name?)</i>

8.2.3 Request Criteria

The following criteria are defined for the **mach_port_allocate** request:

- **C1** — *right_type?* is equal to one of the three values *Mach_port_right_receive*, *Mach_port_right_port_set*, or *Mach_port_right_dead_name*.

<i>C1MachPortAllocateRightIsValid</i> <i>right_type?</i> : <i>RIGHT_TYPE</i> <i>right_type?</i> ∈ { <i>Mach_port_right_receive</i> , <i>Mach_port_right_port_set</i> , <i>Mach_port_right_dead_name</i> }
--

NotC1MachPortAllocateRightIsValid
 $\hat{=} \neg C1MachPortAllocateRightIsValid$

- **C2** — *right_type?* is equal to *Mach_port_right_receive*.

<i>C2MachPortAllocateReceive</i> <i>right_type?</i> : <i>RIGHT_TYPE</i> <i>right_type?</i> = <i>Mach_port_right_receive</i>

- **C3** — *right_type?* is equal to *Mach_port_right_port_set*.

<i>C3MachPortAllocatePortSet</i>
<i>right_type?</i> : <i>RIGHT_TYPE</i>
<i>right_type?</i> = <i>Mach_port_right_port_set</i>

- **C4** — *right_type?* is equal to *Mach_port_right_dead_name*.

<i>C4MachPortAllocateDeadName</i>
<i>right_type?</i> : <i>RIGHT_TYPE</i>
<i>right_type?</i> = <i>Mach_port_right_dead_name</i>

- **C5** — The number of rights belonging to *task?* is less than the maximum.

Review Note:

The maximum number of rights per task used to be modeled by a constant value. This was not correct. However, the correct model has not been determined either. Therefore there is no predicate in the following schema. Note that this means that the schema *NotC5MachPortAllocateRoomInNameSpace* is empty.

<i>C5MachPortAllocateRoomInNameSpace</i>
<i>PortNameSpace</i>
<i>task?</i> : <i>TASK</i>

NotC5MachPortAllocateRoomInNameSpace
 $\hat{=} \text{PortNameSpace} \wedge \neg \text{C5MachPortAllocateRoomInNameSpace}$

8.2.4 Return Values

Table 2 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

The table has been checked to agree with the code in CM as of 14Sep94.

Review Note:

This request can also return *Kern_resource_shortage*. If there were an attempt to model this, there would be three separate criteria for each of the three types of possible rights. Moreover, in the case when there is no room in the name space and no memory available for allocation, it appears that details in the current state of the name space may determine whether *Kern_resource_shortage* or *Kern_no_space* is returned.

Review Note:

According to the KID, there should also be a check that *task?* has *Hold_receive* privilege to the new port whenever *right_type?=mach_port_type_receive*. This is not currently in the model nor was it in the code at the time the model was written.

<i>return!</i>	<i>new_name!</i>	C1	C5
<i>Kern_invalid_value</i>	-	F	-
<i>Kern_no_space</i>	-	T	F
<i>Kern_success</i>	<i>new_name!</i>	T	T

Table 2: Return Values for `mach_port_allocate`

<i>RVMachPortAllocateInvalid Value</i>
<i>MachPortAllocate Outputs</i>
<i>NotC1MachPortAllocateRightIs Valid</i>
<i>return!</i> = <i>Kern_invalid_value</i>

<i>RVMachPortAllocateNoSpace</i>
<i>MachPortAllocate Outputs</i>
<i>C1MachPortAllocateRightIs Valid</i>
<i>NotC5MachPortAllocateRoomInNameSpace</i>
<i>return!</i> = <i>Kern_no_space</i>

Review Note:
The way in which *new_name!* is actually determined is dependent upon the “next” available index in the name space and the generation number of that index. Since these details do not appear in the model, we cannot accurately model how *new_name!* is determined. Therefore all that is done here is to state properties on *new_name!*.

<i>RVMachPortAllocateSuccess</i>
<i>MachPortAllocate Outputs</i>
<i>C1MachPortAllocateRightIs Valid</i>
<i>C5MachPortAllocateRoomInNameSpace</i>
<i>return!</i> = <i>Kern_success</i>
<i>new_name!</i> \notin <i>Reserved_names</i>
<i>(task?, new_name!)</i> \notin <i>local_namep</i>

8.2.5 State Changes

Table 3 lists the possible successful executions of a `mach_port_allocate` request. Note that C2, C3 and C4 are mutually exclusive.

Case	C2	C3	C4	C5
<i>MachPortAllocateStateReceive</i>	T	-	-	T
<i>MachPortAllocateStatePortSet</i>	-	T	-	T
<i>MachPortAllocateStateDeadName</i>	-	-	T	T

Table 3: State Change Cases for `mach_port_allocate`

If the request is to allocate a new receive right, then

- a new port, *port!*, is created,
- *task?* is given a receive right to *port!* with the name *new_name!*,
- the make-send count for *port!* is initialized to 0,
- the queue size limit for *port!* is initialized to the default value,
- the message queue for *port!* is initialized to be empty,
- the sequence number for *port!*'s message queue is initialized to 0, and
- the SID for *port!* is initialized to the default for *task?*.

MachPortAllocateStateReceive

Δ *PortNameSpace*
 Δ *PortSummary*
 Δ *ObjectSid*
PortSid
C2MachPortAllocateReceive
C5MachPortAllocateRoomInNameSpace
port! : *PORT*
new_name! : *NAME*

$\{port!\} = \underline{port_exists}' \setminus \underline{port_exists}$
 $\underline{port_right_rel}' = \underline{port_right_rel} \cup \{(task?, port!, new_name!, Receive, 1)\}$
 $\underline{make_send_count}' = \underline{make_send_count} \oplus \{port! \mapsto 0\}$
 $\underline{q_limit}' = \underline{q_limit} \oplus \{port! \mapsto Mach_port_q_limit_default\}$
 $\underline{message_in_port_rel}' = \underline{message_in_port_rel} \oplus \{port! \mapsto \langle \rangle\}$
 $\underline{sequence_no}' = \underline{sequence_no} \oplus \{port! \mapsto 0\}$
 $\underline{port_sid}' = \underline{port_sid} \oplus \{port! \mapsto Default_port_sid(\underline{task_sid}(task?))\}$

If the request is to allocate a new port set, then an empty port set with the name *new_name!* is added to *task?*'s name space.

MachPortAllocateStatePortSet

Δ *PortNameSpace*
C3MachPortAllocatePortSet
C5MachPortAllocateRoomInNameSpace
new_name! : *NAME*

$\underline{port_set_rel}' = \underline{port_set_rel} \cup \{(task?, new_name!, \emptyset)\}$

If the request is to allocate a new dead right, then a dead right with the name *new_name!* and a user reference count of 1 is added to *task?*'s name space.

MachPortAllocateStateDeadName

Δ *PortNameSpace*
C4MachPortAllocateDeadName
C5MachPortAllocateRoomInNameSpace
new_name! : *NAME*

$\underline{dead_right_rel}' = \underline{dead_right_rel} \cup \{(task?, new_name!, 1)\}$

Review Note:
Invariants should be stated here as well.

8.2.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_allocate** request is described in Section 8.1.

$\begin{aligned} & \textit{ProcessingMachPortAllocate} \\ & \textit{ProcessPortRequestGood} \\ & \textit{operation?} = \textit{Mach_port_allocate_id} \end{aligned}$

An unsuccessful **mach_port_allocate** request results in no changes to the Mach state and returns only the appropriate error status.

$$\begin{aligned} & \textit{MachPortAllocateBad} \\ & \hat{=} (\textit{RVMachPortAllocateInvalidValue} \\ & \quad \vee \textit{RVMachPortAllocateNoSpace}) \\ & \gg \textit{RequestNoOp} \end{aligned}$$

A successful **mach_port_allocate** request alters the Mach state as described in Section 8.2.5 and returns a reply message.

$$\begin{aligned} & \textit{MachPortAllocateGood} \\ & \hat{=} ((\textit{MachPortAllocateStateReceive} \\ & \quad \vee \textit{MachPortAllocateStatePortSet} \\ & \quad \vee \textit{MachPortAllocateStateDeadName}) \\ & \wedge \textit{RVMachPortAllocateSuccess}) \\ & \gg \textit{MachPortAllocateReply} \end{aligned}$$

The complete specification of kernel processing of a **mach_port_allocate** request consists of the initial processing followed by an unsuccessful or successful execution.

$$\begin{aligned} & \textit{MachPortAllocate} \\ & \hat{=} \textit{ProcessingMachPortAllocate} \\ & \textcircled{\&} (\textit{MachPortAllocateBad} \\ & \quad \vee \textit{MachPortAllocateGood}) \end{aligned}$$

8.3 mach_port_get_receive_status

A **mach_port_get_receive_status** request returns the current status of the port associated with a receive right.

8.3.1 Client Interface

$\text{kern_return_t } \mathbf{mach_port_get_receive_status}$	$(\text{mach_port_t}$ mach_port_t $\text{mach_port_status_t}^*$	$\textit{task_name},$ $\textit{right_name},$ $\textit{port_status});$
---	---	--

8.3.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_get_receive_status** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of a receive right for the port whose status is requested

```

MachPortGetReceiveStatus ClientInputs _____
task_name? : NAME
right_name? : NAME

```

A **mach_port_get_receive_status** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_get_receive_status_id* and has a body consisting of *right_name?*.

```

InvokeMachPortGetReceiveStatus _____
InvokeMachMsg
MachPortGetReceiveStatus ClientInputs
name? = task_name?
operation? = Mach_port_get_receive_status_id
msg_body = Mach_port_get_receive_status_inputs_to_body(right_name?)

```

8.3.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_get_receive_status** request:

- *return!* — the status of the request
- *port_status!* — the status information for *right_name?*, as described in Section 8.1.4.

```

MachPortGetReceiveStatus ClientOutputs _____
return! : KERNEL_RETURN
port_status! : PortStatus

```

```

MachPortGetReceiveStatus ReceiveReply _____
InvokeMachMsgRcv
MachPortGetReceiveStatus ClientOutputs
(return!, port_status!) = Body_to_mach_port_get_receive_status_outputs(msg_body)

```

8.3.2 Kernel Interface

8.3.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_get_receive_status** request:

- *task?* — the task known to the client by *task_name?*
- *right_name?* — provided by the client

```

MachPortGetReceiveStatusInputs _____
task? : TASK
right_name? : NAME

```

8.3.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_get_receive_status** request:

- *return!* — the status of the request
- *port_status!* — the status information for *right_name?*, as described above

<i>MachPortGetReceiveStatus Outputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i> <i>port_status!</i> : <i>PortStatus</i>

Upon completion of the processing of a **mach_port_get_receive_status** request, a reply message is built from the output parameters.

<i>MachPortGetReceiveStatus Reply</i>
<i>RequestOnlyObserves</i> <i>port_status?</i> : <i>PortStatus</i>
<i>reply?</i> = <i>Mach_port_get_receive_status_outputs_to_reply(port_status?)</i>

8.3.3 Request Criteria

The following criteria are defined for the **mach_port_get_receive_status** request:

- **C1** — *right_name?* represents a right in *task?*'s name space.

<i>C1MachPortGetReceiveStatusNameIsARight</i>
<i>PortNameSpace</i> <i>task?</i> : <i>TASK</i> <i>right_name?</i> : <i>NAME</i>
$(task?, right_name?) \in local_namep$

NotC1MachPortGetReceiveStatusNameIsARight
 $\hat{=} PortNameSpace \wedge \neg C1MachPortGetReceiveStatusNameIsARight$

- **C2** — *right_name?* represents a receive right in *task?*'s name space.

<i>C2MachPortGetReceiveStatusNameIsAReceiveRight</i>
<i>PortNameSpace</i> <i>task?</i> : <i>TASK</i> <i>right_name?</i> : <i>NAME</i>
$(task?, right_name?) \in r_right$

NotC2MachPortGetReceiveStatusNameIsAReceiveRight
 $\hat{=} PortNameSpace \wedge \neg C2MachPortGetReceiveStatusNameIsAReceiveRight$

8.3.4 Return Values

Table 4 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

The order of these checks agrees with the code in CM on 14Sep94.

<i>return!</i>	<i>port_status!</i>	C1	C2
<i>Kern_invalid_name</i>	-	F	-
<i>Kern_invalid_right</i>	-	T	F
<i>Kern_success</i>	as described below	T	T

Table 4: Return Values for *mach_port_get_receive_status*

<i>RVMachPortGetReceiveStatusInvalidName</i>
<i>MachPortGetReceiveStatusOutputs</i>
<i>NotC1MachPortGetReceiveStatusNameIsARight</i>
<i>return!</i> = <i>Kern_invalid_name</i>

<i>RVMachPortGetReceiveStatusInvalidRight</i>
<i>MachPortGetReceiveStatusOutputs</i>
<i>C1MachPortGetReceiveStatusNameIsARight</i>
<i>NotC2MachPortGetReceiveStatusNameIsARight</i>
<i>return!</i> = <i>Kern_invalid_right</i>

In the successful case when *right_name?* refers to a receive right in *task?*'s name space, then the request returns a record with the following fields. Here *port?* refers to the port named by *right_name?*.

- *port_set_name* — if *right_name?* is a member of a port set, the name of this port set; otherwise, the name *Mach_port_null*
- *make_send_count_value* — the make-send count for *port?*
- *port_destroyed_notification_requested* — a boolean value indicating if a port-destroyed notification request is currently active for *port?*
- *no_more_senders_notification_requested* — a boolean value indicating if a no-more-senders notification request is currently active for *port?*
- *msg_count_value* — the number of messages queued at *port?*
- *qlimit_value* — the limit on the number of messages which can be queued at *port?*
- *sequence_no_value* — the current sequence number for the message queue at *port?*
- *number_of_send_once_rights* — the number of send-once rights which exist to *port?*
- *any_send_rights* — a boolean value indicating if there are existing send rights to *port?*

*RV*MachPortGetReceiveStatusSuccess

MachPortGetReceiveStatus Outputs
*C1*MachPortGetReceiveStatusNameIsARight
*C2*MachPortGetReceiveStatusNameIsAReceiveRight
PortSummary
Notifications

```

return! = Kern_success
(let port? == named_port(task?, right_name?)
  • port_status!.port_set_name = if port? ∈ ∪(ran port_set)
    then containing_set(port?)
    else Mach_port_null
  ∧ port_status!.make_send_count_value = make_send_count(port?)
  ∧ (port_status!.port_destroyed_notification_requested = True
    ⇔ port? ∈ dom port_notify_destroyed)
  ∧ (port_status!.no_more_senders_notification_requested = True
    ⇔ port? ∈ dom port_notify_no_more_senders)
  ∧ port_status!.msg_count_value = port_size(port?)
  ∧ port_status!.qlimit_value = q_limit(port?)
  ∧ port_status!.sequence_no_value = sequence_no(port?)
  ∧ (port_status!.number_of_send_once_rights
    = #{task : TASK; name : NAME; i : ℕ1
      | (task, port?, name, Send_once, i) ∈ port_right_rel
      • name})
  ∧ (port_status!.any_send_rights = True
    ⇔ (∃ task : TASK; name : NAME; i : ℕ1
      • (task, port?, name, Send, i) ∈ port_right_rel)))

```

Review Note:
The values of *port_status!.number_of_send_once_rights* and *port_status!.any_send_rights* may be incorrect in the specification, since they only count rights which currently exist in some name space. It is unclear whether the code also counts rights in transit.

8.3.5 State Changes

A **mach_port_get_receive_status** request does not make any state changes since it only observes the system state.

8.3.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_get_receive_status** request is described in Section 8.1.

Processing MachPortGetReceiveStatus

ProcessPortRequestGood

operation? = *Mach_port_get_receive_status_id*

An unsuccessful **mach_port_get_receive_status** request results in no changes to the Mach state and returns only the appropriate error status.

```

MachPortGetReceiveStatus Bad
  ≐ (RVMachPortGetReceiveStatusInvalidName
     ∨ RVMachPortGetReceiveStatusInvalidRight)
  >> RequestNoOp

```

A successful **mach_port_get_receive_status** request results in no changes to the Mach state and returns a reply message.

```

MachPortGetReceiveStatus Good
  ≐ RVMachPortGetReceiveStatusSuccess
  >> MachPortGetReceiveStatusReply

```

The complete specification of kernel processing of a **mach_port_get_receive_status** request consists of the initial processing followed by an unsuccessful or successful execution.

```

MachPortGetReceiveStatus
  ≐ ProcessingMachPortGetReceiveStatus
  ; (MachPortGetReceiveStatus Bad
     ∨ MachPortGetReceiveStatus Good)

```

8.4 mach_port_get_refs

A **mach_port_get_refs** request returns the number of user references a task has for a right.

8.4.1 Client Interface

```

kern_return_t mach_port_get_refs(
    mach_port_t      port,
    mach_port_t      right,
    mach_port_right_t right_type,
    mach_port_urefs_t* urefs,
    task_name_t      task_name,
    right_name_t     right_name,
    right_type_t     right_type,
    refs_t           refs);

```

8.4.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_get_refs** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of the right whose reference count is desired
- *right_type?* — the type of right for which the reference count is requested, either *Mach_port_right_send*, *Mach_port_right_receive*, *Mach_port_right_send_once*, *Mach_port_right_port_set*, or *Mach_port_right_dead_name*

```

MachPortGetRefsClientInputs
  task_name? : NAME
  right_name? : NAME
  right_type? : RIGHT_TYPE

```

A **mach_port_get_refs** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_get_refs_id* and has a body consisting of *right_name?* and *right_type?*.

<i>InvokeMachPortGetRefs</i> <i>InvokeMachMsg</i> <i>MachPortGetRefs ClientInputs</i> <hr/> <i>name?</i> = <i>task_name?</i> <i>operation?</i> = <i>Mach_port_get_refs_id</i> <i>msg_body</i> = <i>Mach_port_get_refs_inputs_to_body(right_name?, right_type?)</i>

8.4.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_get_refs** request:

- *return!* — the status of the request
- *refs!* — the number of user references to the type of right indicated by *right_type?* and associated with *right_name?*

<i>MachPortGetRefs ClientOutputs</i> <hr/> <i>return!</i> : <i>KERNEL_RETURN</i> <i>refs!</i> : \mathbb{N}
--

<i>MachPortGetRefs ReceiveReply</i> <hr/> <i>InvokeMachMsgRcv</i> <i>MachPortGetRefs ClientOutputs</i> <hr/> (<i>return!</i> , <i>refs!</i>) = <i>Body_to_mach_port_get_refs_outputs(msg_body)</i>

8.4.2 Kernel Interface

8.4.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_get_refs** request:

- *task?* — the task known to the client by *task_name?*
- *right_name?* — provided by the client
- *right_type?* — provided by the client

<i>MachPortGetRefs Inputs</i> <hr/> <i>task?</i> : <i>TASK</i> <i>right_name?</i> : <i>NAME</i> <i>right_type?</i> : <i>RIGHT_TYPE</i>

8.4.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_get_refs** request:

- *return!* — the status of the request
- *refs!* — the number of user references for *right_type?* associated with *right_name?*

<i>MachPortGetRefsOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i> <i>refs!</i> : \mathbb{N}

Upon completion of the processing of a **mach_port_get_refs** request, a reply message is built from the output parameters.

<i>MachPortGetRefsReply</i> <i>RequestOnlyObserves</i> <i>refs?</i> : \mathbb{N}
<i>reply?</i> = <i>Mach_port_get_refs_outputs_to_reply(refs?)</i>

8.4.3 Request Criteria

The following criteria are defined for the **mach_port_get_refs** request:

- **C1** — The value of *right_type?* is one of the five values *Mach_port_right_send*, *Mach_port_right_receive*, *Mach_port_right_send_once*, *Mach_port_right_port_set*, or *Mach_port_right_dead_name*.

<i>C1MachPortGetRefsRightIsRecognized</i> <i>right_type?</i> : <i>RIGHT_TYPE</i>
$right_type? \in \{Mach_port_right_send, Mach_port_right_receive, Mach_port_right_send_once, Mach_port_right_port_set, Mach_port_right_dead_name\}$

NotC1MachPortGetRefsRightIsRecognized
 $\hat{=} \neg C1MachPortGetRefsRightIsRecognized$

- **C2** — *right_name?* represents a right in *task?*'s name space.

<i>C2MachPortGetRefsNameIsARight</i> <i>PortNameSpace</i> <i>task?</i> : <i>TASK</i> <i>right_name?</i> : <i>NAME</i>
$(task?, right_name?) \in local_namep$

NotC2MachPortGetRefsNameIsARight
 $\hat{=} PortNameSpace \wedge \neg C2MachPortGetRefsNameIsARight$

- **C3** — *right_type?* and *right_name?* both refer to a send right.

$C3MachPortGetRefsSend$ $PortNameSpace$ $task? : TASK$ $right_name? : NAME$ $right_type? : RIGHT_TYPE$ <hr/> $right_type? = Mach_port_right_send$ $(task?, right_name?) \in s_right$

$$NotC3MachPortGetRefsSend$$

$$\hat{=} PortNameSpace \wedge \neg C3MachPortGetRefsSend$$

- **C4** — *right_type?* and *right_name?* both refer to a dead name.

$C4MachPortGetRefsDeadName$ $PortNameSpace$ $task? : TASK$ $right_name? : NAME$ $right_type? : RIGHT_TYPE$ <hr/> $right_type? = Mach_port_right_dead_name$ $(task?, right_name?) \in dead_namep$
--

$$NotC4MachPortGetRefsDeadName$$

$$\hat{=} PortNameSpace \wedge \neg C4MachPortGetRefsDeadName$$

- **C5** — *right_type?* and *right_name?* both refer to a receive right, send-once right, or port set.

$C5MachPortGetRefsOther$ $PortNameSpace$ $task? : TASK$ $right_name? : NAME$ $right_type? : RIGHT_TYPE$ <hr/> $(right_type? = Mach_port_right_receive$ $\wedge (task?, right_name?) \in r_right)$ $\vee (right_type? = Mach_port_right_send_once$ $\wedge (task?, right_name?) \in so_right)$ $\vee (right_type? = Mach_port_right_port_set$ $\wedge (task?, right_name?) \in port_set_namep)$
--

$$NotC5MachPortGetRefsOther$$

$$\hat{=} PortNameSpace \wedge \neg C5MachPortGetRefsOther$$

8.4.4 Return Values

Table 5 describes the values returned at the completion of the request and the conditions under which each value is returned. Note that criteria C3 through C5 are mutually exclusive by definition.

Review Note:

The order of the checks agree with the code in CM as of 14Sep94.

<i>return!</i>	<i>refs!</i>	C1	C2	C3	C4	C5
<i>Kern_invalid_value</i>	-	F	-	-	-	-
<i>Kern_invalid_name</i>	-	T	F	-	-	-
<i>Kern_success</i>	0	T	T	F	F	F
<i>Kern_success</i>	<i>s_right_ref_count(task?,right_name?)</i>	T	T	T	-	-
<i>Kern_success</i>	<i>dead_right_ref_count(task?,right_name?)</i>	T	T	-	T	-
<i>Kern_success</i>	1	T	T	-	-	T

Table 5: Return Values for **mach_port_get_refs**

<i>RVMachPortGetRefsInvalidValue</i>
<i>MachPortGetRefsOutputs</i>
<i>NotC1MachPortGetRefsRightIsRecognized</i>
<i>return! = Kern_invalid_value</i>

<i>RVMachPortGetRefsInvalidName</i>
<i>MachPortGetRefsOutputs</i>
<i>C1MachPortGetRefsRightIsRecognized</i>
<i>NotC2MachPortGetRefsNameIsARight</i>
<i>return! = Kern_invalid_name</i>

If *right_name?* does not represent a right of type *right_type?*, the value 0 is returned.

<i>RVMachPortGetRefsWrongRight</i>
<i>MachPortGetRefsOutputs</i>
<i>C1MachPortGetRefsRightIsRecognized</i>
<i>C2MachPortGetRefsNameIsARight</i>
<i>NotC3MachPortGetRefsSend</i>
<i>NotC4MachPortGetRefsDeadName</i>
<i>NotC5MachPortGetRefsOther</i>
<i>return! = Kern_success</i>
<i>refs! = 0</i>

If *right_type? = Mach_port_right_send* and *right_name?* is a send right, then the number of user references to the send right is returned.

<i>RVMachPortGetRefsSend</i> <i>MachPortGetRefsOutputs</i> <i>C1MachPortGetRefsRightIsRecognized</i> <i>C2MachPortGetRefsNameIsARight</i> <i>C3MachPortGetRefsSend</i>
<i>return!</i> = <i>Kern_success</i> <i>refs!</i> = <i>s_right_ref_count(task?, right_name?)</i>

If *right_type?* = *Mach_port_right_dead_name* and *right_name?* is a dead right, then the number of user references to the dead right is returned.

<i>RVMachPortGetRefsDeadName</i> <i>MachPortGetRefsOutputs</i> <i>C1MachPortGetRefsRightIsRecognized</i> <i>C2MachPortGetRefsNameIsARight</i> <i>C4MachPortGetRefsDeadName</i>
<i>return!</i> = <i>Kern_success</i> <i>refs!</i> = <i>dead_right_ref_count(task?, right_name?)</i>

If *right_type?* and *right_name?* both refer to a receive right, send-once right, or port set, then the value 1 is returned since there is only one receive right, send-once right or port set associated with any name.

<i>RVMachPortGetRefsOther</i> <i>MachPortGetRefsOutputs</i> <i>C1MachPortGetRefsRightIsRecognized</i> <i>C2MachPortGetRefsNameIsARight</i> <i>C5MachPortGetRefsOther</i>
<i>return!</i> = <i>Kern_success</i> <i>refs!</i> = 1

8.4.5 State Changes

A ***mach_port_get_refs*** request does not make any state changes since it only observes the system state.

8.4.6 Complete Request

The initial processing by the kernel upon receipt of the ***mach_port_get_refs*** request is described in Section 8.1.

<i>ProcessingMachPortGetRefs</i> <i>ProcessPortRequestGood</i>
<i>operation?</i> = <i>Mach_port_get_refs_id</i>

An unsuccessful ***mach_port_get_refs*** request results in no changes to the Mach state and returns only the appropriate error status.

```

MachPortGetRefs Bad
  ≐ (RVMachPortGetRefsInvalidValue
     ∨ RVMachPortGetRefsInvalidValue)
  >> RequestNoOp

```

A successful **mach_port_get_refs** request results in no changes to the Mach state and returns a reply message.

```

MachPortGetRefs Good
  ≐ (RVMachPortGetRefsWrongRight
     ∨ RVMachPortGetRefsSend
     ∨ RVMachPortGetRefsDeadName
     ∨ RVMachPortGetRefsOther)
  >> MachPortGetRefs Reply

```

The complete specification of kernel processing of a **mach_port_get_refs** request consists of the initial processing followed by an unsuccessful or successful execution.

```

MachPortGetRefs
  ≐ ProcessingMachPortGetRefs
  ; (MachPortGetRefs Bad
     ∨ MachPortGetRefs Good)

```

8.5 mach_port_get_set_status

A **mach_port_get_set_status** request returns the names of the members of a given port set.

8.5.1 Client Interface

```

kern_return_t mach_port_get_set_status
    (mach_port_t
     mach_port_t
     mach_port_array_t*
     mach_msg_type_number_t*
     task_name,
     right_name,
     member_names,
     count);

```

8.5.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_get_set_status** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of the port set whose members are returned

```

MachPortGetSetStatus ClientInputs
task_name? : NAME
right_name? : NAME

```

A **mach_port_get_set_status** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_get_set_status_id* and has a body consisting of *right_name?*.

```

InvokeMachPortGetSetStatus _____
InvokeMachMsg
MachPortGetSetStatus ClientInputs
_____
name? = task_name?
operation? = Mach_port_get_set_status_id
msg_body = Mach_port_get_set_status_inputs_to_body(right_name?)

```

8.5.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_get_set_status** request:

- *return!* — the status of the request
- *member_names!* — the names of the members of the port set *right_name?*
- *count!* — the number of members of the port set *right_name?*

```

MachPortGetSetStatus ClientOutputs _____
return! : KERNEL_RETURN
member_names! : P NAME
count! : N

```

```

MachPortGetSetStatus ReceiveReply _____
InvokeMachMsgRcv
MachPortGetSetStatus ClientOutputs
_____
(return!, member_names!, count!)
= Body_to_mach_port_get_set_status_outputs(msg_body)

```

8.5.2 Kernel Interface

8.5.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_get_set_status** request:

- *task?* — the task known to the client by *task_name?*
- *right_name?* — provided by the client

```

MachPortGetSetStatus Inputs _____
task? : TASK
right_name? : NAME

```

8.5.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_get_set_status** request:

- *return!* — the status of the request
- *member_names!* — the names of the members of the port set *right_name?*
- *count!* — the number of members of the port set *right_name?*

MachPortGetSetStatus Outputs

return! : *KERNEL_RETURN*
member_names! : \mathbb{P} *NAME*
count! : \mathbb{N}

Upon completion of the processing of a **mach_port_get_set_status** request, a reply message is built from the output parameters.

MachPortGetSetStatus Reply

RequestOnlyObserves
member_names? : \mathbb{P} *NAME*
count? : \mathbb{N}

reply? = *Mach_port_get_set_status_outputs_to_reply(member_names?, count?)*

8.5.3 Request Criteria

The following criteria are defined for the **mach_port_get_set_status** request:

- **C1** — *right_name?* represents a right in *task?*'s name space.

C1MachPortGetSetStatusNameIsARight

PortNameSpace
task? : *TASK*
right_name? : *NAME*

$(task?, right_name?) \in local_namep$

NotC1MachPortGetSetStatusNameIsARight

$\hat{=} PortNameSpace \wedge \neg C1MachPortGetSetStatusNameIsARight$

- **C2** — *right_name?* represents a port set in *task?*'s name space.

C2MachPortGetSetStatusNameIsAPortSet

PortNameSpace
task? : *TASK*
right_name? : *NAME*

$(task?, right_name?) \in port_set_namep$

NotC2MachPortGetSetStatusNameIsAPortSet

$\hat{=} PortNameSpace \wedge \neg C2MachPortGetSetStatusNameIsAPortSet$

8.5.4 Return Values

Table 6 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

The order of the checks agrees with the code in CM on 14Sep94.

Review Note:

This request can also return *Kern_resource_shortage*, when allocating memory for *member_names*.

In the code from which this model was produced, one page of memory is allocated for *member_names* before checking either C1 or C2. If C1 and C2 are both true, then an exhaustive search of *task?*'s name space is performed to search for receive rights to any of the ports in the port set named by *right_name?*. If more rights are found than will fit on one page of memory, then additional memory is allocated at this time.

Thus *Kern_resource_shortage* may be returned before or after checking C1 and C2. If it is returned after the other checks, *member_names* will contain a partial list of names.

<i>return!</i>	<i>member_names!</i>	<i>count!</i>	C1	C2
<i>Kern_invalid_name</i>	-	-	F	-
<i>Kern_invalid_right</i>	-	-	T	F
<i>Kern_success</i>	as described below	as described below	T	T

Table 6: Return Values for **mach_port_get_set_status**

<i>RVMachPortGetSetStatusInvalidName</i> <i>MachPortGetSetStatus Outputs</i> <i>NotC1MachPortGetSetStatusNameIsARight</i> <hr/> <i>return!</i> = <i>Kern_invalid_name</i>
--

<i>RVMachPortGetSetStatusInvalidRight</i> <i>MachPortGetSetStatus Outputs</i> <i>C1MachPortGetSetStatusNameIsARight</i> <i>NotC2MachPortGetSetStatusNameIsAPortSet</i> <hr/> <i>return!</i> = <i>Kern_invalid_right</i>

<i>RVMachPortGetSetStatusSuccess</i> <i>MachPortGetSetStatus Outputs</i> <i>C1MachPortGetSetStatusNameIsARight</i> <i>C2MachPortGetSetStatusNameIsAPortSet</i> <hr/> <i>return!</i> = <i>Kern_success</i> <i>member_names!</i> = { <i>port</i> : <i>PORT</i> <i>port</i> ∈ <i>port_set(task?, right_name?)</i> • <i>receiver_name(port)</i> } <i>count!</i> = # <i>member_names!</i>

8.5.5 State Changes

A **mach_port_get_set_status** request does not make any state changes since it only observes the system state.

8.5.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_get_set_status** request is described in Section 8.1.

<i>ProcessingMachPortGetSetStatus</i> _____ <i>ProcessPortRequestGood</i> <hr/> <i>operation? = Mach_port_get_set_status_id</i>

An unsuccessful **mach_port_get_set_status** request results in no changes to the Mach state and returns only the appropriate error status.

MachPortGetSetStatusBad
 $\hat{=}$ (*RVMachPortGetSetStatusInvalidName*
 \vee *RVMachPortGetSetStatusInvalidRight*)
 \gg *RequestNoOp*

A successful **mach_port_get_set_status** request results in no changes to the Mach state and returns a reply message.

MachPortGetSetStatusGood
 $\hat{=}$ *RVMachPortGetSetStatusSuccess*
 \gg *MachPortGetSetStatusReply*

The complete specification of kernel processing of a **mach_port_get_set_status** request consists of the initial processing followed by an unsuccessful or successful execution.

MachPortGetSetStatus
 $\hat{=}$ *ProcessingMachPortGetSetStatus*
 \ddagger (*MachPortGetSetStatusBad*
 \vee *MachPortGetSetStatusGood*)

8.6 mach_port_names

A **mach_port_names** request returns information about all of the rights in a task's port name space. The same information for a single right can be retrieved using **mach_port_type**.

8.6.1 Client Interface

kern_return_t **mach_port_names**
 (mach_port_t *task_name*,

<i>mach_port_array_t*</i>	<i>right_names,</i>
<i>mach_msg_type_number_t*</i>	<i>ncount,</i>
<i>mach_port_type_array_t*</i>	<i>type_masks,</i>
<i>mach_msg_type_number_t*</i>	<i>tcount);</i>

8.6.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_names** request:

- *task_name?* — the client's name for the task whose port name space is returned

```

MachPortNamesClientInputs
task_name? : NAME

```

A **mach_port_names** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_names_id* and an empty body.

```

InvokeMachPortNames
InvokeMachMsg
MachPortNamesClientInputs
name? = task_name?
operation? = Mach_port_names_id

```

8.6.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_names** request:

- *return!* — the status of the request
- *right_names!* — a sequence consisting of all names in the port name space of *task_name?*
- *ncount!* — the number of elements in the sequence *right_names!*
- *type_masks!* — a sequence of port type masks, as described in Section 8.1.4, corresponding to each element of the sequence *right_names!*
- *tcount!* — the number of elements in the sequence *type_masks!* (which is the same as *ncount!*)

```

MachPortNamesClientOutputs
return! : KERNEL_RETURN
right_names! : seq NAME
ncount! : N
type_masks! : seq PortTypeMask
tcount! : N

```

```

MachPortNamesReceiveReply
InvokeMachMsgRcv
MachPortNamesClientOutputs
(return!, right_names!, ncount!, type_masks!, tcount!)
= Body_to_mach_port_names_outputs(msg_body)

```

8.6.2 Kernel Interface

8.6.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_names** request:

- *task?* — the task known to the client by *task_name?*

```

MachPortNamesInputs
task? : TASK

```

8.6.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_names** request:

- *return!* — the status of the request
- *right_names!* — a sequence consisting of all names in the port name space of *task_name?*
- *ncount!* — the number of elements in the sequence *right_names!*
- *type_masks!* — a sequence of port type masks, as described in Section 8.1.4, corresponding to each element of the sequence *right_names!*
- *tcount!* — the number of elements in the sequence *type_masks!* (which is the same as *ncount!*)

```

MachPortNamesOutputs
return! : KERNEL_RETURN
right_names! : seq NAME
ncount! : N
type_masks! : seq PortTypeMask
tcount! : N

```

Upon completion of the processing of a **mach_port_names** request, a reply message is built from the output parameters.

```

MachPortNamesReply
RequestOnlyObserves
right_names? : seq NAME
ncount? : N
type_masks? : seq PortTypeMask
tcount? : N

reply?
= Mach_port_names_outputs_to_reply(right_names?, ncount?, type_masks?, tcount?)

```

8.6.3 Request Criteria

There are no criteria for the **mach_port_names** request.

8.6.4 Return Values

Table 7 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>	<i>right_names!</i>	<i>ncount!</i>	<i>type_masks!</i>	<i>tcount!</i>
<i>Kern_success</i>	described below	# <i>right_names!</i>	described below	# <i>type_masks!</i>

Table 7: Return Values for **mach_port_names**

Review Note:

There is also the possibility that *Kern_resource_shortage* can be returned, though this has not been modeled.

Review Note:

No specification is given for the *mach_port_type_msg_accepted_request* field, since the model currently does not include the necessary information.

Review Note:

There is actually also a COMPAT field in the return mask. However, in the prototype it will always return false so it need not be modeled.

right_names! is a sequence consisting of all names in the name space for *task?*. No element of the name space is repeated, so the number of elements in *right_names!* is *number_of_rights(task?)*.

ncount! is the number of elements in *right_names!*.

type_masks! is a sequence, whose length *tcount!* is the same as *ncount!*, consisting of a *PortTypeMask* for each corresponding name in *right_names!*. Each mask contains boolean flags indicating the following:

- *mach_port_type_send*— if the name refers to a send right
- *mach_port_type_receive*— if the name refers to a receive right
- *mach_port_type_send_once*— if the name refers to a send-once right
- *mach_port_type_port_set*— if the name refers to a port set
- *mach_port_type_dead_name*— if the name refers to a dead right
- *mach_port_type_dead_name_request*— if there is an outstanding dead-name notification request for the name

<i>RVMachPortNamesSuccess</i> <i>MachPortNames Outputs</i> <i>PortNameSpace</i> <i>Notifications</i> <i>task? : TASK</i>
<i>return!</i> = <i>Kern_success</i> <i>ran right_names!</i> = <i>local_namep</i> ({ <i>task?</i> }) <i>ncount!</i> = # <i>right_names!</i> = <i>number_of_rights</i> (<i>task?</i>) <i>tcount!</i> = # <i>type_masks!</i> = <i>ncount!</i> $\forall i : 1 \dots ncount!$ <ul style="list-style-type: none"> • ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_send</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>s_right</i>) \wedge ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_receive</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>r_right</i>) \wedge ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_send_once</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>so_right</i>) \wedge ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_port_set</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>port_set_namep</i>) \wedge ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_dead_name</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>dead_namep</i>) \wedge ((<i>type_masks!</i>(<i>i</i>)).<i>mach_port_type_dead_name_request</i> = <i>True</i> \Leftrightarrow (<i>task?</i>, <i>right_names!</i>(<i>i</i>) \in <i>dom port_notify_dead</i>)

8.6.5 State Changes

A **mach_port_names** request does not make any state changes since it only observes the system state.

8.6.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_names** request is described in Section 8.1.

<i>ProcessingMachPortNames</i> <i>ProcessPortRequestGood</i>
<i>operation?</i> = <i>Mach_port_names_id</i>

A **mach_port_names** request results in no changes to the Mach state and returns a reply message.

MachPortNames Good
 $\hat{=}$ *RVMachPortNamesSuccess*
 \gg *MachPortNames Reply*

The complete specification of kernel processing of a **mach_port_names** request consists of the initial processing followed by execution.

MachPortNames
 $\hat{=}$ *ProcessingMachPortNames*
 \S *MachPortNames Good*

8.7 mach_port_rename

A **mach_port_rename** request allows a client to change the name by which a task knows a port, port set or dead name.

8.7.1 Client Interface

```
kern_return_t mach_port_rename
    (mach_port_t          task_name,
     mach_port_t          old_name,
     mach_port_t          new_name);
```

8.7.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_rename** request:

- *task_name?* — the client's name for the task whose port name space is to be changed
- *old_name?* — the current name of the right to be renamed
- *new_name?* — the new name for the right

```
MachPortRename ClientInputs
task_name? : NAME
old_name?  : NAME
new_name?  : NAME
```

A **mach_port_rename** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_rename_id* and has a body consisting of *old_name?* and *new_name?*.

```
InvokeMachPortRename
InvokeMachMsg
MachPortRename ClientInputs
name? = task_name?
operation? = Mach_port_rename_id
msg_body = Mach_port_rename_inputs_to_body(old_name?, new_name?)
```

8.7.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_rename** request:

- *return!* — the status of the request

```

MachPortRenameClientOutputs _____
return! : KERNEL_RETURN

```

```

MachPortRenameReceiveReply _____
InvokeMachMsgRcv
MachPortRenameClientOutputs
return! = Body_to_mach_port_rename_outputs(msg_body)

```

8.7.2 Kernel Interface

8.7.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_rename** request:

- *task?* — the task known to the client by *task_name?*
- *old_name?* — the current name of the right to be renamed
- *new_name?* — the new name for the right

```

MachPortRenameInputs _____
task? : TASK
old_name? : NAME
new_name? : NAME

```

8.7.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_rename** request:

- *return!* — the status of the request

```

MachPortRenameOutputs _____
return! : KERNEL_RETURN

```

8.7.3 Request Criteria

The following criteria are defined for the **mach_port_rename** request:

- **C1** — *new_name?* is not a reserved name.

```

C1MachPortRenameNewNameNotReserved _____
new_name? : NAME
new_name? ∉ Reserved_names

```

$$\begin{aligned} & \text{NotC1MachPortRenameNewNameNotReserved} \\ & \hat{=} \neg \text{C1MachPortRenameNewNameNotReserved} \end{aligned}$$

- **C2** — *new_name?* is not currently in the name space of *task?*.

$\begin{aligned} & \text{C2MachPortRenameNewNameNotInUse} \\ & \text{PortNameSpace} \\ & \text{task? : TASK} \\ & \text{new_name? : NAME} \\ & \text{(task?, new_name?)} \notin \text{local_namep} \end{aligned}$
--

$$\begin{aligned} & \text{NotC2MachPortRenameNewNameNotInUse} \\ & \hat{=} \text{PortNameSpace} \wedge \neg \text{C2MachPortRenameNewNameNotInUse} \end{aligned}$$

- **C3** — *old_name?* is currently in the name space of *task?*.

$\begin{aligned} & \text{C3MachPortRenameOldNameInUse} \\ & \text{PortNameSpace} \\ & \text{task? : TASK} \\ & \text{old_name? : NAME} \\ & \text{(task?, old_name?)} \in \text{local_namep} \end{aligned}$
--

$$\begin{aligned} & \text{NotC3MachPortRenameOldNameInUse} \\ & \hat{=} \text{PortNameSpace} \wedge \neg \text{C3MachPortRenameOldNameInUse} \end{aligned}$$

8.7.4 Return Values

Table 8 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:
The criteria are in the correct order according to the code in CM on 16Sep94.

Review Note:
This request can also return *Kern_resource_shortage*, because it may need to create a new entry. This check, if it occurs, comes between C1 and C2.

$\begin{aligned} & \text{RVMachPortRenameInvalidValue} \\ & \text{MachPortRenameOutputs} \\ & \text{NotC1MachPortRenameNewNameNotReserved} \\ & \text{return!} = \text{Kern_invalid_value} \end{aligned}$

<i>return!</i>	C1	C2	C3
<i>Kern_invalid_value</i>	F	-	-
<i>Kern_name_exists</i>	T	F	-
<i>Kern_invalid_name</i>	T	T	F
<i>Kern_success</i>	T	T	T

Table 8: Return Values for **mach_port_rename**

<i>RV MachPortRenameNameExists</i>
<i>MachPortRename Outputs</i>
<i>C1 MachPortRenameNewNameNotReserved</i>
<i>Not C2 MachPortRenameNewNameNotInUse</i>
<i>return! = Kern_name_exists</i>

<i>RV MachPortRenameInvalidName</i>
<i>MachPortRename Outputs</i>
<i>C1 MachPortRenameNewNameNotReserved</i>
<i>C2 MachPortRenameNewNameNotInUse</i>
<i>Not C3 MachPortRenameOldNameInUse</i>
<i>return! = Kern_invalid_name</i>

<i>RV MachPortRenameSuccess</i>
<i>MachPortRename Outputs</i>
<i>C1 MachPortRenameNewNameNotReserved</i>
<i>C2 MachPortRenameNewNameNotInUse</i>
<i>C3 MachPortRenameOldNameInUse</i>
<i>return! = Kern_success</i>

8.7.5 State Changes

If all of the criteria are satisfied, then the name space for *task?* is changed so that the name of any right currently with the name *old_name?* is changed to *new_name?*. Here right refers to a send, receive, or send once right, port set or dead name.

In addition, if there is an outstanding dead-name request for the right *old_name?*, then the request must be renamed.

Review Note:

In addition, the name of a message accepted request may need to be changed. However, that is not currently in the model.

There are several things in the model that should be changed by this request, but are not because the model is incorrect.

- A name in a port set. The model is incorrect, since the kernel actually considers port sets as a set of ports, not a set of names.
- A name in the set of registered rights (*registered_rights*) for a task. Again, the model is incorrect since this is really a set of ports, not a set of names. (This is modeled correctly in the FSPM.)

- The data structures associated with threads blocked for a pending receive (*PendingReceive* and *pending_receives*). Again, the model is incorrect in considering a port to be a name.

MachPortRenameState

Δ *PortNameSpace*
 Δ *Notifications*
C1MachPortRenameNewNameNotReserved
C2MachPortRenameNewNameNotInUse
C3MachPortRenameOldNameInUse

$\underline{port_right_rel}' = \underline{port_right_rel}$
 $\setminus \{ port : PORT; right : RIGHT; i : \mathbb{N}_1 \bullet (task?, port, old_name?, right, i) \}$
 $\cup \{ port : PORT; right : RIGHT; i : \mathbb{N}_1$
 $\quad | (task?, port, old_name?, right, i) \in \underline{port_right_rel}$
 $\quad \bullet (task?, port, new_name?, right, i) \}$

$\underline{port_set_rel}' = \underline{port_set_rel}$
 $\setminus \{ set_of_ports : \mathbb{P} PORT \bullet (task?, old_name?, set_of_ports) \}$
 $\cup \{ set_of_ports : \mathbb{P} PORT$
 $\quad | (task?, old_name?, set_of_ports) \in \underline{port_set_rel}$
 $\quad \bullet (task?, new_name?, set_of_ports) \}$

$\underline{dead_right_rel}' = \underline{dead_right_rel}$
 $\setminus \{ i : \mathbb{N}_1 \bullet (task?, old_name?, i) \}$
 $\cup \{ i : \mathbb{N}_1 \mid (task?, old_name?, i) \in \underline{dead_right_rel}$
 $\quad \bullet (task?, new_name?, i) \}$

$\underline{port_notify_dead_rel}' = \underline{port_notify_dead_rel}$
 $\setminus \{ port : PORT \bullet (port, task?, old_name?) \}$
 $\cup \{ port : PORT \mid (port, task?, old_name?) \in \underline{port_notify_dead_rel}$
 $\quad \bullet (port, task?, new_name?) \}$

Review Note:
 Invariants should be stated here as well.

8.7.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_rename** request is described in Section 8.1.

ProcessingMachPortRename

ProcessPortRequestGood

$operation? = Mach_port_rename_id$

An unsuccessful **mach_port_rename** request results in no changes to the Mach state and returns only the appropriate error status.

$$\begin{aligned}
 & MachPortRenameBad \\
 & \hat{=} (RVMachPortRenameNameExists \\
 & \quad \vee RVMachPortRenameInvalidName \\
 & \quad \vee RVMachPortRenameInvalidValue) \\
 & \gg RequestNoOp
 \end{aligned}$$

A successful **mach_port_rename** request alters the Mach state as described in Section 8.7.5 and returns a reply message.

$$\begin{aligned}
 & MachPortRenameGood \\
 & \hat{=} (MachPortRenameState \\
 & \quad \wedge RVMachPortRenameSuccess) \\
 & \gg RequestReturnOnlyStatus
 \end{aligned}$$

The complete specification of kernel processing of a **mach_port_rename** request consists of the initial processing followed by an unsuccessful or successful execution.

$$\begin{aligned}
 & MachPortRename \\
 & \hat{=} ProcessingMachPortRename \\
 & \ ; (MachPortRenameBad \\
 & \quad \vee MachPortRenameGood)
 \end{aligned}$$

8.8 mach_port_request_notification

A **mach_port_request_notification** request registers a notification message to be sent when a particular port event occurs. If a notification has already been requested, it returns the send-once right associated with the existing notification request.

Notifications are described in Section 8.1.7.

8.8.1 Client Interface

kern_return_t	mach_port_request_notification	
	(mach_port_t	<i>task_name,</i>
	mach_port_t	<i>right_name,</i>
	mach_msg_id_t	<i>variant,</i>
	mach_port_mscount_t	<i>sync,</i>
	mach_port_t	<i>notify_name,</i>
	mach_msg_type_name_t	<i>notify_type,</i>
	mach_port_t*	<i>previous_name);</i>

8.8.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_request_notification** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located

- *right_name?* — the name of a right, in *task_name?*'s name space, for which the notification is requested. If *variant?* is set to *Mach_notify_port_destroyed* or *Mach_notify_no_senders*, this must be a receive right.
- *variant?* — the type of event for which notification is requested, either *Mach_notify_port_destroyed*, *Mach_notify_no_senders*, or *Mach_notify_dead_name*
- *sync?* — When *variant?* is set to *Mach_notify_dead_name*, this must be set to zero. When *variant?* is set to *Mach_notify_no_senders*, this value is used to overcome race conditions
- *notify_name?* — the name of a right, in the client's name space, for the port to which the notification should be send
- *notify_type?* — the manner in which a send-once right should be extracted from *notify_name?*, either *Mmt_make_send_once* or *Mmt_move_send_once*

```

MachPortRequestNotificationClientInputs _____
task_name? : NAME
right_name? : NAME
variant? : MACH_MSG_ID
sync? : N
notify_name? : NAME
notify_type? : MACH_MSG_TYPE

```

A **mach_port_request_notification** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_request_notification_id* and has a body consisting of *right_name?*, *variant?*, *sync?*, *notify_name?*, and *notify_type?*.

```

InvokeMachPortRequestNotification _____
InvokeMachMsg
MachPortRequestNotificationClientInputs
name? = task_name?
operation? = Mach_port_request_notification_id
msg_body = Mach_port_request_notification_inputs_to_body(right_name?, variant?,
sync?, notify_name?, notify_type?)

```

8.8.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_request_notification** request:

- *return!* — the status of the request
- *previous_name!* — if the notification has already been requested, the previously registered send-once right; otherwise *Mach_port_null*

```

MachPortRequestNotificationClientOutputs _____
return! : KERNEL_RETURN
previous_name! : NAME

```

```

MachPortRequestNotificationReceiveReply _____
InvokeMachMsgRcv
MachPortRequestNotificationClientOutputs
(return!, previous_name!) = Body_to_mach_port_request_notification_outputs(msg_body)

```

8.8.2 Kernel Interface

8.8.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_request_notification** request:

- *task?* — the task known to the client by *task_name?*
- *right_name?* — provided by the client
- *variant?* — provided by the client
- *sync?* — provided by the client
- *notify?* — the port known to the client by *notify_name?*

```

MachPortRequestNotificationInputs _____
task? : TASK
right_name? : NAME
variant? : MACH_MSG_ID
sync? : N
notify? : PORT

```

8.8.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_request_notification** request:

- *return!* — the status of the request
- *previous!* — if the notification has already been requested, the previously registered notification port; otherwise *Ip_null*

```

MachPortRequestNotificationOutputs _____
return! : KERNEL_RETURN
previous! : PORT

```

Upon completion of the processing of a **mach_port_request_notification** request, a reply message is built from the output parameters.

```

MachPortRequestNotificationReply _____
RequestReturn
previous? : PORT
reply? = Mach_port_request_notification_outputs_to_reply(previous?)

```

8.8.3 Request Criteria

The following criteria are defined for the **mach_port_request_notification** request:

Review Note:

There is a number missing in the list because of improvements to the model. It doesn't seem to be worthwhile to renumber the criteria.

- **C1** — *notify?* is not *Ip_dead*.

<i>C1MachPortRequestNotificationNotifyNotDead</i> _____
<i>notify?</i> : <i>PORT</i>
<i>notify?</i> ≠ <i>Ip_dead</i>

NotC1MachPortRequestNotificationNotifyNotDead
 $\hat{=} \neg C1MachPortRequestNotificationNotifyNotDead$

- **C2** — *variant?* is set to *Mach_notify_port_destroyed*.

<i>C2MachPortRequestNotificationPortDestroyed</i> _____
<i>variant?</i> : <i>MACH_MSG_ID</i>
<i>variant?</i> = <i>Mach_notify_port_destroyed</i>

NotC2MachPortRequestNotificationPortDestroyed
 $\hat{=} \neg C2MachPortRequestNotificationPortDestroyed$

- **C3** — *variant?* is set to *Mach_notify_no_senders*.

<i>C3MachPortRequestNotificationNoSenders</i> _____
<i>variant?</i> : <i>MACH_MSG_ID</i>
<i>variant?</i> = <i>Mach_notify_no_senders</i>

NotC3MachPortRequestNotificationNoSenders
 $\hat{=} \neg C3MachPortRequestNotificationNoSenders$

- **C4** — *variant?* is set to *Mach_notify_dead_name*.

<i>C4MachPortRequestNotificationDeadName</i> _____
<i>variant?</i> : <i>MACH_MSG_ID</i>
<i>variant?</i> = <i>Mach_notify_dead_name</i>

NotC4MachPortRequestNotificationDeadName
 $\hat{=} \neg C4MachPortRequestNotificationDeadName$

■ **C5** — *sync?* is equal to zero.

$C5MachPortRequestNotificationSyncIsZero$
$sync? : \mathbb{N}$
$sync? = 0$

$$\begin{aligned} &NotC5MachPortRequestNotificationSyncIsZero \\ &\hat{=} \neg C5MachPortRequestNotificationSyncIsZero \end{aligned}$$

■ **C6** — *right_name?* represents a right in *task?*'s name space.

$C6MachPortRequestNotificationNameIsARight$
$PortNameSpace$
$task? : TASK$
$right_name? : NAME$
$(task?, right_name?) \in local_namep$

$$\begin{aligned} &NotC6MachPortRequestNotificationNameIsARight \\ &\hat{=} PortNameSpace \wedge \neg C6MachPortRequestNotificationNameIsARight \end{aligned}$$

■ **C7** — *right_name?* represents a receive right in *task?*'s name space.

$C7MachPortRequestNotificationNameIsAReceiveRight$
$PortNameSpace$
$task? : TASK$
$right_name? : NAME$
$(task?, right_name?) \in r_right$

$$\begin{aligned} &NotC7MachPortRequestNotificationNameIsAReceiveRight \\ &\hat{=} PortNameSpace \wedge \neg C7MachPortRequestNotificationNameIsAReceiveRight \end{aligned}$$

■ **C8** — This criteria only applies if *right_name?* indicates a receive right in *task?*'s name space, in which case it is true whenever:

- There are no send rights for the port indicated by *right_name?*,
- *sync?* is less than the make-send count value for that port, and
- *notify?* is not *Ip_null*.

Review Note:

This criteria is not completely defined in the schema because the state model does not capture the total number of send rights associated with a port.

$C8MachPortRequestNoSendersNotificationSendNow$
$C7MachPortRequestNotificationNameIsAReceiveRight$
$SendRightsCount$
$sync? : \mathbb{N}$
$notify? : PORT$
$sync? \leq \underline{make_send_count}(named_port(task?, right_name?))$
$notify? \neq Ip_null$

$$\begin{aligned} & \text{Not}C8\text{MachPortRequestNoSendersNotificationSendNow} \\ & \hat{=} C7\text{MachPortRequestNotificationNameIsAReceiveRight} \wedge \text{SendRightsCount} \\ & \quad \wedge \neg C8\text{MachPortRequestNoSendersNotificationSendNow} \end{aligned}$$

- **C10** — *right_name?* represents a send, send-once or receive right in *task?*'s name space.

$\begin{aligned} & C10\text{MachPortRequestNotificationNameIsAPortRight} \\ & \text{PortNameSpace} \\ & \text{task?} : \text{TASK} \\ & \text{right_name?} : \text{NAME} \end{aligned}$
$(\text{task?}, \text{right_name?}) \in \text{port_right_namep}$

$$\begin{aligned} & \text{Not}C10\text{MachPortRequestNotificationNameIsAPortRight} \\ & \hat{=} \text{PortNameSpace} \wedge \neg C10\text{MachPortRequestNotificationNameIsAPortRight} \end{aligned}$$

- **C11** — *right_name?* represents a deadname in *task?*'s name space.

$\begin{aligned} & C11\text{MachPortRequestNotificationNameIsADeadName} \\ & \text{PortNameSpace} \\ & \text{task?} : \text{TASK} \\ & \text{right_name?} : \text{NAME} \end{aligned}$
$(\text{task?}, \text{right_name?}) \in \text{dead_namep}$

$$\begin{aligned} & \text{Not}C11\text{MachPortRequestNotificationNameIsADeadName} \\ & \hat{=} \text{PortNameSpace} \wedge \neg C11\text{MachPortRequestNotificationNameIsADeadName} \end{aligned}$$

- **C12** — *notify?* is not *Ip_null* and *sync?* is non-zero.

$\begin{aligned} & C12\text{MachPortRequestNotificationNotifySyncNotZero} \\ & \text{notify?} : \text{PORT} \\ & \text{sync?} : \mathbb{N} \end{aligned}$
$\begin{aligned} & \text{notify?} \neq \text{Ip_null} \\ & \text{sync?} \neq 0 \end{aligned}$

$$\begin{aligned} & \text{Not}C12\text{MachPortRequestNotificationNotifyAndSyncNonZero} \\ & \hat{=} \neg C12\text{MachPortRequestNotificationNotifySyncNotZero} \end{aligned}$$

- **C13** — The number of user references for *right_name?* is less than the maximum allowed.

$\begin{aligned} & C13\text{MachPortRequestNotificationURefsNotAtMax} \\ & \text{PortNameSpace} \\ & \text{task?} : \text{TASK} \\ & \text{right_name?} : \text{NAME} \end{aligned}$
$\text{dead_right_ref_count}(\text{task?}, \text{right_name?}) < \text{Max_right_refs}$

$$\begin{aligned} & \text{Not}C13\text{MachPortRequestNotificationURefsNotAtMax} \\ & \hat{=} \text{PortNameSpace} \wedge \neg C13\text{MachPortRequestNotificationURefsNotAtMax} \end{aligned}$$

8.8.4 Return Values

Table 9 describes those values returned at the completion of the request which are common to all three values of *variant?*. Tables 10, 11 and 12 describe the return values specific to the cases in which *variant?* is set to *Mach_notify_port_destroyed*, *Mach_notify_no_senders*, and *Mach_notify_dead_name*, respectively.

Review Note:

The order of the checks in this section is accurate based upon the code in CM on 19Sep94.

Review Note:

In Table 9, note that C2, C3, and C4 are mutually exclusive.

<i>return!</i>	<i>previous!</i>	C1	C2	C3	C4
<i>Kern_invalid_capability</i>	-	F	-	-	-
<i>Kern_invalid_value</i>	-	T	F	F	F
See Table 10		T	T	-	-
See Table 11		T	-	T	-
See Table 12		T	-	-	T

Table 9: Return Values for *mach_port_request_notification*

<i>RV MachPortRequestNotificationInvalid Capability</i>
<i>MachPortRequestNotification Outputs</i>
<i>Not C1 MachPortRequestNotificationNotifyNotDead</i>
<i>return!</i> = <i>Kern_invalid_capability</i>

<i>RV MachPortRequestNotificationInvalid Value</i>
<i>MachPortRequestNotification Outputs</i>
<i>C1 MachPortRequestNotificationNotifyNotDead</i>
<i>Not C2 MachPortRequestNotificationPortDestroyed</i>
<i>Not C3 MachPortRequestNotificationNoSenders</i>
<i>Not C4 MachPortRequestNotificationDeadName</i>
<i>return!</i> = <i>Kern_invalid_value</i>

8.8.4.1 Port-Destroyed Notification Request Table 10 describes the return values for the case in which *variant?* is set to *Mach_notify_port_destroyed*.

<i>RV MachPortRequestPortDestroyedNotificationInvalid Value</i>
<i>MachPortRequestNotification Outputs</i>
<i>C1 MachPortRequestNotificationNotifyNotDead</i>
<i>C2 MachPortRequestNotificationPortDestroyed</i>
<i>Not C5 MachPortRequestNotificationSyncIsZero</i>
<i>return!</i> = <i>Kern_invalid_value</i>

<i>return!</i>	<i>previous!</i>	C1	C2	C5	C6	C7
<i>Kern_invalid_value</i>	-	T	T	F	-	-
<i>Kern_invalid_name</i>	-	T	T	T	F	-
<i>Kern_invalid_right</i>	-	T	T	T	T	F
<i>Kern_success</i>	<i>port_notify_destroyed(</i> <i>named_port(task?, right_name?))</i>	T	T	T	T	T

Table 10: Return Values for *mach_port_request_notification*, port-destroyed notification

RVMachPortRequestPortDestroyedNotificationInvalidName _____
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C2 MachPortRequestNotificationPortDestroyed
C5 MachPortRequestNotificationSyncIsZero
NotC6 MachPortRequestNotificationNameIsARight

return! = Kern_invalid_name

RVMachPortRequestPortDestroyedNotificationInvalidRight _____
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C2 MachPortRequestNotificationPortDestroyed
C5 MachPortRequestNotificationSyncIsZero
C6 MachPortRequestNotificationNameIsARight
NotC7 MachPortRequestNotificationNameIsAReceiveRight

return! = Kern_invalid_right

RVMachPortRequestPortDestroyedNotificationSuccess _____
Notifications
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C2 MachPortRequestNotificationPortDestroyed
C5 MachPortRequestNotificationSyncIsZero
C6 MachPortRequestNotificationNameIsARight
C7 MachPortRequestNotificationNameIsAReceiveRight

return! = Kern_success
previous! = port_notify_destroyed(named_port(task?, right_name?))

8.8.4.2 No-Senders Notification Request Table 11 describes the return values for the case in which *variant?* is set to *Mach_notify_no_senders*.

RVMachPortRequestNoSendersNotificationInvalidName _____
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C3 MachPortRequestNotificationNoSenders
NotC6 MachPortRequestNotificationNameIsARight

return! = Kern_invalid_name

<i>return!</i>	<i>previous!</i>	C1	C3	C6	C7
<i>Kern_invalid_name</i>	-	T	T	F	-
<i>Kern_invalid_right</i>	-	T	T	T	F
<i>Kern_success</i>	<i>port_notify_no_more_senders(</i> <i>named_port(task?, right_name?))</i>	T	T	T	T

Table 11: Return Values for **mach_port_request_notification**, no-senders notification

<i>RVMachPortRequestNoSendersNotificationInvalidRight</i>
<i>MachPortRequestNotification Outputs</i>
<i>C1 MachPortRequestNotificationNotifyNotDead</i>
<i>C3 MachPortRequestNotificationNoSenders</i>
<i>C6 MachPortRequestNotificationNameIsARight</i>
<i>NotC7 MachPortRequestNotificationNameIsAReceiveRight</i>
<i>return! = Kern_invalid_right</i>

<i>RVMachPortRequestNoSendersNotificationSuccess</i>
<i>Notifications</i>
<i>MachPortRequestNotification Outputs</i>
<i>C1 MachPortRequestNotificationNotifyNotDead</i>
<i>C3 MachPortRequestNotificationNoSenders</i>
<i>C6 MachPortRequestNotificationNameIsARight</i>
<i>C7 MachPortRequestNotificationNameIsAReceiveRight</i>
<i>return! = Kern_success</i>
<i>previous! = port_notify_no_more_senders(named_port(task?, right_name?))</i>

8.8.4.3 Dead-Name Notification Request Table 12 describes the return values for the case in which *variant?* is set to *Mach_notify_dead_name*.

Review Note:

In Table 12, note that C10 and C11 are mutually exclusive.

<i>return!</i>	<i>previous!</i>	C1	C4	C6	C10	C11	C12	C13
<i>Kern_invalid_name</i>	-	T	T	F	-	-	-	-
<i>Kern_invalid_right</i>	-	T	T	T	F	F	-	-
<i>Kern_success</i>	<i>port_notify_dead(</i> <i>task?, right_name?)</i>	T	T	T	T	-	-	-
<i>Kern_invalid_argument</i>	-	T	T	T	-	T	F	-
<i>Kern_urefs_overflow</i>	-	T	T	T	-	T	T	F
<i>Kern_success</i>	<i>Ip_null</i>	T	T	T	-	T	T	T

Table 12: Return Values for **mach_port_request_notification**, dead-name notification

RVMachPortRequestDeadNameNotificationInvalidName _____
MachPortRequestNotification Outputs
C1MachPortRequestNotificationNotifyNotDead
C4MachPortRequestNotificationDeadName
NotC6MachPortRequestNotificationNameIsARight

return! = *Kern_invalid_name*

RVMachPortRequestDeadNameNotificationInvalidRight _____
MachPortRequestNotification Outputs
C1MachPortRequestNotificationNotifyNotDead
C4MachPortRequestNotificationDeadName
C6MachPortRequestNotificationNameIsARight
NotC10MachPortRequestNotificationNameIsAPortRight
NotC11MachPortRequestNotificationNameIsADeadName

return! = *Kern_invalid_right*

RVMachPortRequestDeadNameNotificationSuccessOne _____
Notifications
MachPortRequestNotification Outputs
C1MachPortRequestNotificationNotifyNotDead
C4MachPortRequestNotificationDeadName
C6MachPortRequestNotificationNameIsARight
C10MachPortRequestNotificationNameIsAPortRight

return! = *Kern_success*
previous! = *port_notify_dead(task?, right_name?)*

Review Note:

Note that this case could also result in a resource shortage. This is due to the fact that many dead name notifications can be active for the same port, so memory is allocated for each additional notification request.

RVMachPortRequestDeadNameNotificationInvalidArgument _____
MachPortRequestNotification Outputs
C1MachPortRequestNotificationNotifyNotDead
C4MachPortRequestNotificationDeadName
C6MachPortRequestNotificationNameIsARight
C11MachPortRequestNotificationNameIsADeadName
NotC12MachPortRequestNotificationNotifyAndSyncNonZero

return! = *Kern_invalid_argument*

Editorial Note:

It seems rather surprising that the request fails when *name?* is a dead name and *notify?* is *Ip_null?* That seems like a case that Mach would usually allow, though the request in this case would do nothing.

```

RV MachPortRequestDeadNameNotificationUrefsOverflow _____
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C4 MachPortRequestNotificationDeadName
C6 MachPortRequestNotificationNameIsARight
C11 MachPortRequestNotificationNameIsADeadName
C12 MachPortRequestNotificationNotifySyncNotZero
NotC13 MachPortRequestNotificationURefsNotAtMax

return! = Kern_urefs_overflow

```

```

RV MachPortRequestDeadNameNotificationSuccessTwo _____
MachPortRequestNotification Outputs
C1 MachPortRequestNotificationNotifyNotDead
C4 MachPortRequestNotificationDeadName
C6 MachPortRequestNotificationNameIsARight
C11 MachPortRequestNotificationNameIsADeadName
C12 MachPortRequestNotificationNotifySyncNotZero
C13 MachPortRequestNotificationURefsNotAtMax

return! = Kern_success
previous! = Ip_null

```

In this last case, *previous!* returns *Ip_null* since *right_name?* identifies a dead right.

8.8.5 State Changes

Table 13 lists the possible successful executions of a **mach_port_request_notification** request.

Case	C1	C2	C3	C4	C5	C6	C7	C8	C10	C11	C12	C13
<i>MachPortRequestPortDestroyedNotificationStateChanges</i>	T	T	-	-	T	T	T	-	-	-	-	-
<i>MachPortRequestNoSendersNotificationStateChangesOne</i>	T	-	T	-	-	T	T	F	-	-	-	-
<i>MachPortRequestNoSendersNotificationStateChangesTwo</i>	T	-	T	-	-	T	T	T	-	-	-	-
<i>MachPortRequestDeadNameNotificationStateChangesOne</i>	T	-	-	T	-	T	-	-	T	-	-	-
<i>MachPortRequestDeadNameNotificationStateChangesTwo</i>	T	-	-	T	-	T	-	-	-	T	T	T

Table 13: State Change Cases for **mach_port_request_notification**

8.8.5.1 Port-Destroyed Notification Request A request for a port-destroyed notification is successful if the following are true:

- *notify?* is not *Ip_dead*
- *variant?* is set to *Mach_notify_port_destroyed*
- *sync?* is zero
- *right_name?* is a receive right in *task?*'s name space

If the request is successful, the kernel replaces the existing port-destroyed notification request port for the port named by *right_name?* with *notify?*. Note that either the existing notification request port or *notify?* could be *Ip_null*.

<p style="text-align: right; margin: 0;"><i>MachPortRequestPortDestroyedNotificationStateChanges</i> _____</p> <p>Δ Notifications</p> <p><i>C1 MachPortRequestNotificationNotifyNotDead</i></p> <p><i>C2 MachPortRequestNotificationPortDestroyed</i></p> <p><i>C5 MachPortRequestNotificationSyncIsZero</i></p> <p><i>C6 MachPortRequestNotificationNameIsARight</i></p> <p><i>C7 MachPortRequestNotificationNameIsAReceiveRight</i></p> <hr style="border: 0.5px solid black;"/> <p><i>port_notify_destroyed'</i> = <i>port_notify_destroyed</i> $\oplus \{ \text{named_port}(task?, right_name?) \mapsto notify? \}$</p>

8.8.5.2 No-Senders Notification Request A request for a no-senders notification is successful if the following are true:

- *notify?* is not *Ip_dead*
- *variant?* is set to *Mach_notify_no_senders*
- *right_name?* is a receive right in *task?*'s name space

If the request is successful, then kernel further determines whether a no-senders notification should immediately be sent, by checking if all of the following are true:

- There are no send rights for the port indicated by *right_name?*
- *sync?* is less than the make-send count value for that port
- *notify?* is not *Ip_null*

If any of these are not true, the kernel replaces the existing no-senders notification request port for the port named by *right_name?* with *notify?*. Note that either the existing notification request port or *notify?* could be *Ip_null*.

<p style="text-align: right; margin: 0;"><i>MachPortRequestNoSendersNotificationStateChangesOne</i> _____</p> <p>Δ Notifications</p> <p><i>C1 MachPortRequestNotificationNotifyNotDead</i></p> <p><i>C3 MachPortRequestNotificationNoSenders</i></p> <p><i>C6 MachPortRequestNotificationNameIsARight</i></p> <p><i>C7 MachPortRequestNotificationNameIsAReceiveRight</i></p> <p><i>NotC8 MachPortRequestNoSendersNotificationSendNow</i></p> <hr style="border: 0.5px solid black;"/> <p><i>port_notify_no_more_senders'</i> = <i>port_notify_no_more_senders</i> $\oplus \{ \text{named_port}(task?, right_name?) \mapsto notify? \}$</p>

If the kernel determines that it must immediately send a no-senders notification, it first removes any existing notification request port and then attempts to send the notification.

MachPortRequestNoSendersNotificationStateChangesTwo _____

Δ Notifications

C1 *MachPortRequestNotificationNotifyNotDead*

C3 *MachPortRequestNotificationNoSenders*

C6 *MachPortRequestNotificationNameIsARight*

C7 *MachPortRequestNotificationNameIsAReceiveRight*

C8 *MachPortRequestNoSendersNotificationSendNow*

$port_notify_no_more_senders' = port_notify_no_more_senders$
 $\oplus \{named_port(task?, right_name?) \mapsto Ip_null\}$

Review Note:

This second case still needs to be completed. Four things are missing, all of which should eventually be handled in the port chapter introduction:

- Make sure that a notification message can be allocated. If not, no notification is sent.
 - Build the notification message.
 - Check for send permission
 - Queue the message
-

Review Note:

It's interesting to note that the *sync?* value is used to determine whether a no-senders notification is immediately sent, but it has no relevance to notifications sent after the notification is registered. In that respect it seems that the race condition that *sync?* is apparently intended to prevent is not really prevented. The only way to truly avoid it is for the recipient of the notification to check the make send count returned with the notification.

MachPortRequestNoSendersNotificationStateChanges

$\hat{=} MachPortRequestNoSendersNotificationStateChangesOne$

$\vee MachPortRequestNoSendersNotificationStateChangesTwo$

8.8.5.3 Dead-Name Notification Request A request for a dead-name notification can be successful in two distinct cases.

The first successful case occurs if the following are true:

- *notify?* is not *Ip_dead*
- *variant?* is set to *Mach_notify_dead_name*
- *right_name?* is a send, send-once or receive right in *task?*'s name space

In this case, the kernel replaces the existing dead-name notification request port for *right_name?* with *notify?*. Note that either the existing notification request port or *notify?* could be *Ip_null*.

MachPortRequestDeadNameNotificationStateChangesOne _____

Δ Notifications

C1 *MachPortRequestNotificationNotifyNotDead*

C4 *MachPortRequestNotificationDeadName*

C6 *MachPortRequestNotificationNameIsARight*

C10 *MachPortRequestNotificationNameIsAPortRight*

$port_notify_dead' = port_notify_dead \oplus \{(task?, right_name?) \mapsto notify?\}$

The second successful case occurs if the following are true:

- *notify?* is not *Ip_dead*
- *variant?* is set to *Mach_notify_dead_name*
- *right_name?* is a dead right in *task?*'s name space
- *sync?* is non-zero
- *dead_right_ref_count(task?,right_name?)* is not at the maximum value (*Max_right_refs*)

In this case, the kernel immediately attempts to send a dead-name notification. Prior to doing this however, it increments the *dead_right_ref_count* for the dead right.

MachPortRequestDeadNameNotification.StateChangesTwo

Δ *DeadRights*

C1MachPortRequestNotificationNotifyNotDead

C4MachPortRequestNotificationDeadName

C6MachPortRequestNotificationNameIsARight

C11MachPortRequestNotificationNameIsADeadName

C12MachPortRequestNotificationNotifySyncNotZero

C13MachPortRequestNotificationURefsNotAtMax

$$dead_right_ref_count' = dead_right_ref_count$$

$$\oplus \{(task?, right_name?) \mapsto dead_right_ref_count(task?, right_name?) + 1\}$$

Editorial Note:

It is unclear why *dead_right_ref_count* is being incremented in this case. That does not agree with the understanding of this field presented in the state model. It might also be interesting to look and see whether it is decremented in the case that the send of the notification message fails.

Review Note:

This second case still needs to be completed. Four things are missing, all of which should eventually be handled in the port chapter introduction:

- Make sure that a notification message can be allocated. If not, no notification is sent.
 - Build the notification message.
 - Check for send permission
 - Queue the message
-

MachPortRequestDeadNameNotification.StateChanges

$$\hat{=} MachPortRequestDeadNameNotification.StateChangesOne$$

$$\vee MachPortRequestDeadNameNotification.StateChangesTwo$$

8.8.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_request_notification** request is described in Section 8.1.

ProcessingMachPortRequestNotification

ProcessPortRequestGood

operation? = *Mach_port_request_notification_id*

An unsuccessful **mach_port_request_notification** request results in no changes to the Mach state and returns only the appropriate error status.

```
MachPortRequestNotification Bad
  ≐ (RVMachPortRequestNotificationInvalidCapability
    ∨ RVMachPortRequestNotificationInvalidValue
    ∨ RVMachPortRequestPortDestroyedNotificationInvalidValue
    ∨ RVMachPortRequestPortDestroyedNotificationInvalidName
    ∨ RVMachPortRequestPortDestroyedNotificationInvalidRight
    ∨ RVMachPortRequestNoSendersNotificationInvalidName
    ∨ RVMachPortRequestNoSendersNotificationInvalidRight
    ∨ RVMachPortRequestDeadNameNotificationInvalidName
    ∨ RVMachPortRequestDeadNameNotificationInvalidRight
    ∨ RVMachPortRequestDeadNameNotificationInvalidArgument
    ∨ RVMachPortRequestDeadNameNotificationUrefsOverflow)
  >> RequestNoOp
```

A successful **mach_port_request_notification** request alters the Mach state as described in Section 8.8.5 and returns a reply message.

```
MachPortRequestNotification Good
  ≐ ((MachPortRequestPortDestroyedNotificationStateChanges
    ∨ MachPortRequestNoSendersNotificationStateChanges
    ∨ MachPortRequestDeadNameNotificationStateChanges)
  ∧ (RVMachPortRequestPortDestroyedNotificationSuccess
    ∨ RVMachPortRequestNoSendersNotificationSuccess
    ∨ RVMachPortRequestDeadNameNotificationSuccessOne
    ∨ RVMachPortRequestDeadNameNotificationSuccessTwo))
  >> MachPortRequestNotificationReply
```

The complete specification of kernel processing of a **mach_port_request_notification** request consists of the initial processing followed by an unsuccessful or successful execution.

```
MachPortRequestNotification
  ≐ ProcessingMachPortRequestNotification
  ; (MachPortRequestNotificationBad
    ∨ MachPortRequestNotificationGood)
```

8.9 mach_port_set_mscount

A **mach_port_set_mscount** request changes the make-send count for the port associated with a specified receive right in a task's port name space.

8.9.1 Client Interface

```
kern_return_t mach_port_set_mscount
    (mach_port_t
    mach_port_t
    mach_port_mscount_t
    task_name,
    right_name,
    mscount);
```


8.9.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_set_mscount** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of a receive right for the port whose make-send count is to be changed
- *mscount?* — the value to be assigned to the make-send count for the port associated with *right_name?*

```

MachPortSetMscount ClientInputs
task_name? : NAME
right_name? : NAME
mscount? : N

```

A **mach_port_set_mscount** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_set_mscount_id* and has a body consisting of *right_name?* and *mscount?*.

```

Invoke MachPortSetMscount
Invoke MachMsg
MachPortSetMscount ClientInputs
name? = task_name?
operation? = Mach_port_set_mscount_id
msg_body = Mach_port_set_mscount_inputs_to_body(right_name?, mscount?)

```

8.9.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_set_mscount** request:

- *return!* — the status of the request

```

MachPortSetMscount ClientOutputs
return! : KERNEL_RETURN

```

```

MachPortSetMscount ReceiveReply
Invoke MachMsgRcv
MachPortSetMscount ClientOutputs
return! = Body_to_mach_port_set_mscount_outputs(msg_body)

```

8.9.2 Kernel Interface

8.9.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_set_mscount** request:

- *task?* — the task known to the client by *task_name?*

- *right_name?* — provided by the client
- *mscount?* — provided by the client

<i>MachPortSetMscountInputs</i>
<i>task?</i> : <i>TASK</i>
<i>right_name?</i> : <i>NAME</i>
<i>mscount?</i> : <i>N</i>

8.9.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_set_mscount** request:

- *return!* — the status of the request

<i>MachPortSetMscountOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

8.9.3 Request Criteria

The following criteria are defined for the **mach_port_set_mscount** request:

- **C1** — *right_name?* represents a right in *task?*'s name space.

<i>C1MachPortSetMscountNameIsARight</i>
<i>PortNameSpace</i>
<i>task?</i> : <i>TASK</i>
<i>right_name?</i> : <i>NAME</i>
$(task?, right_name?) \in local_namep$

NotC1MachPortSetMscountNameIsARight
 $\hat{=} PortNameSpace \wedge \neg C1MachPortSetMscountNameIsARight$

- **C2** — *right_name?* represents a receive right in *task?*'s name space.

<i>C2MachPortSetMscountNameIsAReceiveRight</i>
<i>PortNameSpace</i>
<i>task?</i> : <i>TASK</i>
<i>right_name?</i> : <i>NAME</i>
$(task?, right_name?) \in r_right$

NotC2MachPortSetMscountNameIsAReceiveRight
 $\hat{=} PortNameSpace \wedge \neg C2MachPortSetMscountNameIsAReceiveRight$

8.9.4 Return Values

Table 14 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

The order of the checks agrees with the code in CM as of 20Sep94.

<i>return!</i>	C1	C2
<i>Kern_invalid_name</i>	F	-
<i>Kern_invalid_right</i>	T	F
<i>Kern_success</i>	T	T

Table 14: Return Values for `mach_port_set_mscount`

<i>RV</i> <code>MachPortSetMscountInvalidName</code>
<i>MachPortSetMscount</i> Outputs
<i>NotC1MachPortSetMscountNameIsARight</i>
<i>return!</i> = <i>Kern_invalid_name</i>

<i>RV</i> <code>MachPortSetMscountInvalidRight</code>
<i>MachPortSetMscount</i> Outputs
<i>C1MachPortSetMscountNameIsARight</i>
<i>NotC2MachPortSetMscountNameIsAReceiveRight</i>
<i>return!</i> = <i>Kern_invalid_right</i>

<i>RV</i> <code>MachPortSetMscountSuccess</code>
<i>MachPortSetMscount</i> Outputs
<i>C1MachPortSetMscountNameIsARight</i>
<i>C2MachPortSetMscountNameIsAReceiveRight</i>
<i>return!</i> = <i>Kern_success</i>

8.9.5 State Changes

If all of the criteria are satisfied, then the make-send count for the port associated with `right_name?` in `task?`'s name space is given the value `mscount?`.

<i>MachPortSetMscountState</i>
Δ <i>PortSummary</i>
<i>MachPortSetMscount</i> Inputs
<i>C1MachPortSetMscountNameIsARight</i>
<i>C2MachPortSetMscountNameIsAReceiveRight</i>
<i>make_send_count'</i>
= <i>make_send_count</i> \oplus { <i>named_port(task?, right_name?)</i> \mapsto <i>mscount?</i> }

*Review Note:*Invariants should be stated here as well.

8.9.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_set_mscount** request is described in Section 8.1.

<i>ProcessingMachPortSetMscount</i>
<i>ProcessPortRequestGood</i>
<i>operation? = Mach_port_set_mscount_id</i>

An unsuccessful **mach_port_set_mscount** request results in no changes to the Mach state and returns only the appropriate error status.

$$\begin{aligned}
 & MachPortSetMscountBad \\
 & \hat{=} (RV\text{MachPortSetMscountInvalidName} \\
 & \quad \vee RV\text{MachPortSetMscountInvalidRight}) \\
 & \gg RequestNoOp
 \end{aligned}$$

A successful **mach_port_set_mscount** request alters the Mach state as described in Section 8.9.5 and returns a reply message.

$$\begin{aligned}
 & MachPortSetMscountGood \\
 & \hat{=} (MachPortSetMscountState \\
 & \quad \wedge RV\text{MachPortSetMscountSuccess}) \\
 & \gg RequestReturnOnlyStatus
 \end{aligned}$$

The complete specification of kernel processing of a **mach_port_set_mscount** request consists of the initial processing followed by an unsuccessful or successful execution.

$$\begin{aligned}
 & MachPortSetMscount \\
 & \hat{=} ProcessingMachPortSetMscount \\
 & \ddagger (MachPortSetMscountBad \\
 & \quad \vee MachPortSetMscountGood)
 \end{aligned}$$

8.10 mach_port_set_qlimit

A **mach_port_set_qlimit** request changes the message queue limit for the port associated with a specified receive right in a task's port name space.

8.10.1 Client Interface

kern_return_t	mach_port_set_qlimit	
	(mach_port_t	<i>task_name,</i>
	mach_port_t	<i>rightname,</i>
	mach_port_mscount_t	<i>qlimit_value</i>);

8.10.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_set_qlimit** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of a receive right for the port whose message queue limit is to be changed
- *qlimit_value?* — the value to be assigned to the message queue limit for the port associated with *right_name?*

```

MachPortSetQlimitClientInputs
task_name? : NAME
right_name? : NAME
qlimit_value? : N

```

A **mach_port_set_qlimit** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_set_qlimit_id* and has a body consisting of *right_name?* and *qlimit_value?*.

```

InvokeMachPortSetQlimit
InvokeMachMsg
MachPortSetQlimitClientInputs
name? = task_name?
operation? = Mach_port_set_qlimit_id
msg_body = Mach_port_set_qlimit_inputs_to_body(right_name?, qlimit_value?)

```

8.10.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_set_qlimit** request:

- *return!* — the status of the request

```

MachPortSetQlimitClientOutputs
return! : KERNEL_RETURN

```

```

MachPortSetQlimitReceiveReply
InvokeMachMsgRcv
MachPortSetQlimitClientOutputs
return! = Body_to_mach_port_set_qlimit_outputs(msg_body)

```

8.10.2 Kernel Interface

8.10.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_set_qlimit** request:

- *task?* — the task known to the client by *task_name?*

- *right_name?* — provided by the client
- *qlimit_value?* — provided by the client

MachPortSetQlimitInputs

task? : *TASK*
right_name? : *NAME*
qlimit_value? : \mathbb{N}

8.10.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_set_qlimit** request:

- *return!* — the status of the request

MachPortSetQlimitOutputs

return! : *KERNEL_RETURN*

8.10.3 Request Criteria

The following criteria are defined for the **mach_port_set_qlimit** request:

- **C1** — *qlimit_value?* is no larger than the specified maximum, *Mach_port_q_limit_max*.

C1MachPortSetQlimitNameIsARight

qlimit_value? : \mathbb{N}
qlimit_value? \leq *Mach_port_q_limit_max*

NotC1MachPortSetQlimitNameIsARight

$\hat{=} \neg$ *C1MachPortSetQlimitNameIsARight*

- **C2** — *right_name?* represents a right in *task?*'s name space.

C2MachPortSetQlimitNameIsAReceiveRight

PortNameSpace
task? : *TASK*
right_name? : *NAME*
(task?, right_name?) \in *local_namep*

NotC2MachPortSetQlimitNameIsAReceiveRight

$\hat{=}$ *PortNameSpace* \wedge \neg *C2MachPortSetQlimitNameIsAReceiveRight*

- **C3** — *right_name?* represents a receive right in *task?*'s name space.

C3MachPortSetQlimitValueIsValid

PortNameSpace
task? : *TASK*
right_name? : *NAME*

$(task?, right_name?) \in r_right$

NotC3MachPortSetQlimitValueIsValid
 $\hat{=} PortNameSpace \wedge \neg C3MachPortSetQlimitValueIsValid$

8.10.4 Return Values

Table 15 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:
The order of the checks agrees with the code in CM as of 20Sep94.

<i>return!</i>	C1	C2	C3
<i>Kern_invalid_value</i>	F	-	-
<i>Kern_invalid_name</i>	T	F	-
<i>Kern_invalid_right</i>	T	T	F
<i>Kern_success</i>	T	T	T

Table 15: Return Values for *mach_port_set_qlimit*

RVMachPortSetQlimitInvalidValue

MachPortSetQlimitOutputs
NotC1MachPortSetQlimitNameIsARight

return! = *Kern_invalid_value*

RVMachPortSetQlimitInvalidName

MachPortSetQlimitOutputs
C1MachPortSetQlimitNameIsARight
NotC2MachPortSetQlimitNameIsAReceiveRight

return! = *Kern_invalid_name*

RVMachPortSetQlimitInvalidRight

MachPortSetQlimitOutputs
C1MachPortSetQlimitNameIsARight
C2MachPortSetQlimitNameIsAReceiveRight
NotC3MachPortSetQlimitValueIsValid

return! = *Kern_invalid_right*

```

RVMachPortSetQlimitSuccess
MachPortSetQlimitOutputs
C1MachPortSetQlimitNameIsARight
C2MachPortSetQlimitNameIsAReceiveRight
C3MachPortSetQlimitValueIsValid
return! = Kern_success

```

8.10.5 State Changes

If all of the criteria are satisfied, then the message queue limit for the port associated with *right_name?* in *task?*'s name space is given the value *qlimit_value?*.

Review Note:
This request may wake up threads which are blocked trying to send to the port, if the queue limit is increased. This does not currently fit into the model.

```

MachPortSetQlimitState
Δ PortSummary
C1MachPortSetQlimitNameIsARight
C2MachPortSetQlimitNameIsAReceiveRight
C3MachPortSetQlimitValueIsValid
q_limit' = q_limit ⊕ {named_port(task?, right_name?) ↦ qlimit_value?}

```

Review Note:
Invariants should be stated here as well.

8.10.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_set_qlimit** request is described in Section 8.1.

```

ProcessingMachPortSetQlimit
ProcessPortRequestGood
operation? = Mach_port_set_qlimit_id

```

An unsuccessful **mach_port_set_qlimit** request results in no changes to the Mach state and returns only the appropriate error status.

```

MachPortSetQlimitBad
≅ (RVMachPortSetQlimitInvalidName
  ∨ RVMachPortSetQlimitInvalidRight
  ∨ RVMachPortSetQlimitInvalidValue)
>>> RequestNoOp

```


A successful **mach_port_set_qlimit** request alters the Mach state as described in Section 8.10.5 and returns a reply message.

$$\begin{aligned} & MachPortSetQlimitGood \\ & \hat{=} (MachPortSetQlimitState \\ & \wedge RVMachPortSetQlimitSuccess) \\ & \gg RequestReturnOnlyStatus \end{aligned}$$

The complete specification of kernel processing of a **mach_port_set_qlimit** request consists of the initial processing followed by an unsuccessful or successful execution.

$$\begin{aligned} & MachPortSetQlimit \\ & \hat{=} ProcessingMachPortSetQlimit \\ & \ddagger (MachPortSetQlimitBad \\ & \vee MachPortSetQlimitGood) \end{aligned}$$

8.11 mach_port_set_seqno

A **mach_port_set_seqno** request changes the current sequence number for the port associated with a specified receive right in a task's port name space.

8.11.1 Client Interface

```
kern_return_t mach_port_set_seqno
    (mach_port_t
     mach_port_t
     mach_port_seqno_t
     task_name,
     right_name,
     seqno);
```

8.11.1.1 Input Parameters The following input parameters are provided by the client of a **mach_port_set_seqno** request:

- *task_name?* — the client's name for the task in whose name space *right_name?* is located
- *right_name?* — the name of a receive right for the port whose current sequence number is to be changed
- *seqno?* — the value to be assigned to the current sequence number for the port associated with *right_name?*

<p><i>MachPortSetSeqno ClientInputs</i></p> <p><i>task_name?</i> : NAME</p> <p><i>right_name?</i> : NAME</p> <p><i>seqno?</i> : ℤ</p>

A **mach_port_set_seqno** request is invoked by sending a message to the port indicated by *task_name?* that has the operation field set to *Mach_port_set_seqno_id* and has a body consisting of *right_name?* and *seqno?*.

```

InvokeMachPortSetSeqno _____
InvokeMachMsg
MachPortSetSeqno ClientInputs
_____
name? = task_name?
operation? = Mach_port_set_seqno_id
msg_body = Mach_port_set_seqno_inputs_to_body(right_name?, seqno?)

```

8.11.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **mach_port_set_seqno** request:

- *return!* — the status of the request

```

MachPortSetSeqno ClientOutputs _____
return! : KERNEL_RETURN
_____

```

```

MachPortSetSeqno ReceiveReply _____
InvokeMachMsgRcv
MachPortSetSeqno ClientOutputs
_____
return! = Body_to_mach_port_set_seqno_outputs(msg_body)
_____

```

8.11.2 Kernel Interface

8.11.2.1 Input Parameters The following input parameters are provided to the kernel for a **mach_port_set_seqno** request:

- *task?* — the task known to the client by *task_name?*
- *right_name?* — provided by the client
- *seqno?* — provided by the client

```

MachPortSetSeqno Inputs _____
task? : TASK
right_name? : NAME
seqno? : Z
_____

```

8.11.2.2 Output Parameters The following output parameters are returned by the kernel for a **mach_port_set_seqno** request:

- *return!* — the status of the request

```

MachPortSetSeqno Outputs _____
return! : KERNEL_RETURN
_____

```

8.11.3 Request Criteria

The following criteria are defined for the **mach_port_set_seqno** request:

- **C1** — *right_name?* represents a right in *task?*'s name space.

$\text{C1MachPortSetSeqnoNameIsARight} \text{ —————}$ PortNameSpace task? : TASK $\text{right_name? : NAME}$ <hr style="border: 0.5px solid black;"/> $(\text{task?}, \text{right_name?}) \in \text{local_namep}$
--

$$\text{NotC1MachPortSetSeqnoNameIsARight}$$

$$\hat{=} \text{PortNameSpace} \wedge \neg \text{C1MachPortSetSeqnoNameIsARight}$$

- **C2** — *right_name?* represents a receive right in *task?*'s name space.

$\text{C2MachPortSetSeqnoNameIsAReceiveRight} \text{ —————}$ PortNameSpace task? : TASK $\text{right_name? : NAME}$ <hr style="border: 0.5px solid black;"/> $(\text{task?}, \text{right_name?}) \in \text{r_right}$

$$\text{NotC2MachPortSetSeqnoNameIsAReceiveRight}$$

$$\hat{=} \text{PortNameSpace} \wedge \neg \text{C2MachPortSetSeqnoNameIsAReceiveRight}$$

8.11.4 Return Values

Table 16 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:
The order of the checks agrees with the code in CM as of 20Sep94.

<i>return!</i>	C1	C2
<i>Kern_invalid_name</i>	F	-
<i>Kern_invalid_right</i>	T	F
<i>Kern_success</i>	T	T

Table 16: Return Values for **mach_port_set_seqno**

$\text{RVMachPortSetSeqnoInvalidName} \text{ —————}$ $\text{MachPortSetSeqno Outputs}$ $\text{NotC1MachPortSetSeqnoNameIsARight}$ <hr style="border: 0.5px solid black;"/> $\text{return!} = \text{Kern_invalid_name}$
--

RV MachPortSetSeqnoInvalidRight _____
MachPortSetSeqno Outputs
C1 MachPortSetSeqnoNameIsA Right
Not C2 MachPortSetSeqnoNameIsA Receive Right

return! = *Kern_invalid_right*

RV MachPortSetSeqnoSuccess _____
MachPortSetSeqno Outputs
C1 MachPortSetSeqnoNameIsA Right
C2 MachPortSetSeqnoNameIsA Receive Right

return! = *Kern_success*

8.11.5 State Changes

If all of the criteria are satisfied, then the current sequence number for the port associated with *right_name?* in *task?*'s name space is given the value *seqno?*.

MachPortSetSeqnoState _____
 Δ *PortSummary*
MachPortSetSeqno Inputs
C1 MachPortSetSeqnoNameIsA Right
C2 MachPortSetSeqnoNameIsA Receive Right

$\underline{sequence_no}' = \underline{sequence_no} \oplus \{ \text{named_port}(\text{task?}, \text{right_name?}) \mapsto \text{seqno?} \}$

Review Note:
 Invariants should be stated here as well.

8.11.6 Complete Request

The initial processing by the kernel upon receipt of the **mach_port_set_seqno** request is described in Section 8.1.

Processing MachPortSetSeqno _____
ProcessPortRequestGood

operation? = *Mach_port_set_seqno_id*

An unsuccessful **mach_port_set_seqno** request results in no changes to the Mach state and returns only the appropriate error status.

MachPortSetSeqno Bad
 $\hat{=}$ (*RV MachPortSetSeqnoInvalidName*
 \vee *RV MachPortSetSeqnoInvalidRight*)
 \gg *RequestNoOp*

A successful **mach_port_set_seqno** request alters the Mach state as described in Section 8.11.5 and returns a reply message.

$$\begin{aligned} & \text{MachPortSetSeqno Good} \\ & \hat{=} (\text{MachPortSetSeqnoState} \\ & \wedge \text{RV MachPortSetSeqnoSuccess}) \\ & \gg \text{RequestReturnOnlyStatus} \end{aligned}$$

The complete specification of kernel processing of a **mach_port_set_seqno** request consists of the initial processing followed by an unsuccessful or successful execution.

$$\begin{aligned} & \text{MachPortSetSeqno} \\ & \hat{=} \text{ProcessingMachPortSetSeqno} \\ & \ddagger (\text{MachPortSetSeqno Bad} \\ & \quad \vee \text{MachPortSetSeqno Good}) \end{aligned}$$

Section 9

Thread Requests

9.1 Introduction to Thread Requests

This chapter describes the thread kernel requests in DTOS.

9.1.1 Constants and Types

We first define the identifier that is used to represent each thread request. The kernel accepts two thread requests through task kernel ports (*Thread_task_port_ops*) and most of the others through thread kernel ports (*Thread_thread_port_ops*).

```
Thread_abort_id, Thread_assign_id, Thread_assign_default_id,
Thread_depress_abort_id, Thread_disable_pc_sampling_id,
Thread_enable_pc_sampling_id, Thread_get_assignment_id,
Thread_get_sampled_pcs_id, Thread_get_special_port_id,
Thread_get_state_id, Thread_info_id, Thread_max_priority_id,
Thread_policy_id, Thread_priority_id, Thread_resume_id,
Thread_resume_secure_id, Thread_set_special_port_id,
Thread_set_state_id, Thread_set_state_secure_id,
Thread_suspend_id, Thread_terminate_id : OPERATION
Thread_thread_port_ops : P OPERATION
```

```
{Thread_abort_id, Thread_assign_id, Thread_assign_default_id,
  Thread_depress_abort_id, Thread_disable_pc_sampling_id,
  Thread_enable_pc_sampling_id, Thread_get_assignment_id,
  Thread_get_sampled_pcs_id, Thread_get_special_port_id,
  Thread_get_state_id, Thread_info_id, Thread_max_priority_id,
  Thread_policy_id, Thread_priority_id, Thread_resume_id,
  Thread_resume_secure_id, Thread_set_special_port_id,
  Thread_set_state_id, Thread_set_state_secure_id,
  Thread_suspend_id, Thread_terminate_id}
  Values_partition Thread_thread_port_ops
```

```
Thread_create_id, Thread_create_secure_id : OPERATION
Thread_task_port_ops : P OPERATION
```

```
{Thread_create_id, Thread_create_secure_id}
  Values_partition Thread_task_port_ops
```

Together these two disjoint sets of operations form the set *Thread_operations* denoting all thread operations. Each thread request must be received through a port of the appropriate port class.

<p><i>Thread_operations</i> : P OPERATION</p> <p>$\langle \textit{Thread_thread_port_ops}, \textit{Thread_task_port_ops} \rangle$ partition <i>Thread_operations</i></p> <p>$\textit{Thread_thread_port_ops} \subseteq \textit{Allowed_mach_services}(\textit{Pc_thread})$</p> <p>$\textit{Thread_task_port_ops} \subseteq \textit{Allowed_mach_services}(\textit{Pc_task})$</p>

9.1.2 Required Permissions

For each operation there is a primary permission that is required to perform the operation. We define here the portion of the *Required_permission* function that pertains to thread requests. The *Abort_thread* implementation service permission implies the *Set_thread_priority* and *Abort_thread_depress* permissions are automatically granted since the **thread_abort** request can set priorities and abort priority depression. The **thread_priority** request requires *Set_thread_priority* permission, but *Set_max_thread_priority* permission is also needed if the *set_max* parameter has value *True*. We also assume that *Initiate_secure* permission is granted whenever *Resume_thread* or *Set_thread_state* permission is granted.

Review Note:

Here are the full sets of permissions that are currently needed for each request (except the special port ones).

```
{(Thread_abort_id, Abort_thread, Set_thread_priority, Abort_thread_depress),
 (Thread_assign_id, Assign_thread_to_pset, Set_max_thread_priority,
  Set_thread_priority, Set_thread_policy, Assign_thread),
 (Thread_assign_default_id, Assign_thread_to_pset, Set_max_thread_priority,
  Set_thread_priority, Set_thread_policy, Assign_thread),
 (Thread_create_id, Add_thread),
 (Thread_depress_abort_id, Abort_thread_depress, Set_thread_priority),
 (Thread_disable_pc_sampling_id, Sample_thread),
 (Thread_enable_pc_sampling_id, Sample_thread),
 (Thread_get_assignment_id, Get_thread_assignment),
 (Thread_get_sampled_pcs_id, Sample_thread),
 (Thread_get_state_id, Get_thread_state),
 (Thread_info_id, Get_thread_info),
 (Thread_max_priority_id, Set_max_thread_priority, Set_thread_priority),
 (Thread_policy_id, Set_thread_policy),
 (Thread_priority_id, Set_thread_priority, Set_max_thread_priority),
 (Thread_resume_id, Resume_thread, Initiate_secure),
 (Thread_set_state_id, Set_thread_state, Initiate_secure),
 (Thread_suspend_id, Suspend_thread),
 (Thread_terminate_id, Terminate_thread, Sample_thread)}
```

This will be simplified when the FSPM is modified so services do not overlap so often.

Review Note:

Does **thread_assign_default** also require *Assign_thread* permission? I suspect so.

```

{( Thread_abort_id, Abort_thread),
  ( Thread_assign_id, Assign_thread_to_pset),
  ( Thread_assign_default_id, Assign_thread_to_pset),
  ( Thread_create_id, Add_thread),
  ( Thread_create_secure_id, Add_thread_secure),
  ( Thread_depress_abort_id, Abort_thread_depress),
  ( Thread_disable_pc_sampling_id, Sample_thread),
  ( Thread_enable_pc_sampling_id, Sample_thread),
  ( Thread_get_assignment_id, Get_thread_assignment),
  ( Thread_get_sampled_pcs_id, Sample_thread),
  ( Thread_get_state_id, Get_thread_state),
  ( Thread_info_id, Get_thread_info),
  ( Thread_max_priority_id, Set_max_thread_priority),
  ( Thread_policy_id, Set_thread_policy),
  ( Thread_priority_id, Set_thread_priority),
  ( Thread_resume_id, Resume_thread),
  ( Thread_resume_secure_id, Initiate_secure),
  ( Thread_set_state_id, Set_thread_state),
  ( Thread_set_state_secure_id, Initiate_secure),
  ( Thread_suspend_id, Suspend_thread),
  ( Thread_terminate_id, Terminate_thread)}
⊆ Required_permission

```

The permission required for a *Thread_get_special_port_id* or *Thread_set_special_port_id* operation depends upon the value of the *which_port?* parameter. Therefore the permission cannot be checked in the common processing, and the two operations are in the *set.Service_check_deferred*.

$$\{Thread_get_special_port_id, Thread_set_special_port_id\} \subseteq Service_check_deferred$$

9.1.3 Invariant Information

The thread requests operate on only certain components of the state. We use the following schema to provide a general framework for describing thread requests.

Review Note:

This list has problems in that some schemas are indirectly pulled in where they should not be. For example, \exists *SpecialTaskPorts* includes \exists *PortExist* which we do not want. Might be possible to get a better handle on this problem by doing fuzz -t and comparing Δ *DtosExec* to a schema with all of the *ThreadInvariants* except Δ *DtosExec*.

<i>ThreadInvariants</i>
Δ <i>DtosExec</i>
Ξ <i>TaskExist</i>
Ξ <i>MessageExist</i>
Ξ <i>MemoryExist</i>
Ξ <i>PageExist</i>
Ξ <i>ProcessorExist</i>
Ξ <i>ProcessorSetExist</i>
Ξ <i>DeviceExist</i>
Ξ <i>TaskSuspendCount</i>
Ξ <i>Kernel</i>
Ξ <i>RegisteredRights</i>
Ξ <i>MemoriesAndPorts</i>
Ξ <i>HostsAndPorts</i>
Ξ <i>ProcessorsAndPorts</i>
Ξ <i>SpecialTaskPorts</i>
Ξ <i>DevicesAndPorts</i>
Ξ <i>Notifications</i>
Ξ <i>MessageQueues</i>
Ξ <i>MemorySystem</i>
Ξ <i>Messages</i>
Ξ <i>HostsAndProcessors</i>
Ξ <i>ProcessorAndProcessorSet</i>
Ξ <i>TaskAndProcessorSet</i>
Ξ <i>PortClasses</i>
Ξ <i>TaskPriority</i>
Ξ <i>EmulationVector</i>
Ξ <i>MasterDevicePort</i>
Ξ <i>HostTime</i>

9.1.4 General Information

9.1.4.1 Special Ports The requests **thread_get_special_port** and **thread_set_special_port** each have an input parameter specifying the type of special port to be processed. The following type is used for these input parameters:

[THREAD_SPECIAL_PORTS]

There are two recognized values of this type. They are:

- *Thread_exception_port* — indicates the exception port
- *Thread_kernel_port* — indicates the sself port

We require these two values to be disjoint, but place no restrictions on other values of type *THREAD_SPECIAL_PORTS*.

$\begin{array}{l} \text{Thread_exception_port} : \text{THREAD_SPECIAL_PORTS} \\ \text{Thread_kernel_port} : \text{THREAD_SPECIAL_PORTS} \\ \text{Recognized_thread_special_ports} : \mathbb{P} \text{ THREAD_SPECIAL_PORTS} \end{array}$
$\begin{array}{l} \langle \text{Thread_exception_port}, \text{Thread_kernel_port} \rangle \\ \text{Values_partition Recognized_thread_special_ports} \end{array}$

9.1.4.2 Thread Information Types The request **thread_info** returns an array of information describing a thread. The array used to hold the information is of type *THREAD_INFO*.

[*THREAD_INFO*]

There are two recognized types of thread information.

- *Thread_basic_info* — information on execution statistics, execution status and priority
- *Thread_sched_info* — information on scheduling priorities and policies

These two types of information are in the set *THREAD_INFO_TYPE*.

[*THREAD_INFO_TYPE*]

We require the two values *Thread_basic_info* and *Thread_sched_info* of this type to be disjoint, but place no restrictions on other values of *THREAD_INFO_TYPE*.

$\begin{array}{l} \text{Thread_basic_info} : \text{THREAD_INFO_TYPE} \\ \text{Thread_sched_info} : \text{THREAD_INFO_TYPE} \\ \text{Recognized_thread_info_types} : \mathbb{P} \text{ THREAD_INFO_TYPE} \end{array}$
$\begin{array}{l} \langle \text{Thread_basic_info}, \text{Thread_sched_info} \rangle \\ \text{Values_partition Recognized_thread_info_types} \end{array}$

9.1.4.3 Execution Status Changes Several requests (e.g., **thread_suspend** and **switch**) can cause the execution of the current thread to be blocked. We describe here the changes that take place when a thread is blocked.

The blocking of a thread results in the thread being swapped out, and another thread moving onto the processor, unless there is nothing else for the processor to swap in. The run states will change for the thread moved off the processor and for the thread moved onto the processor. The thread moved onto the processor is determined by the scheduling algorithm. The algorithm may select the blocking thread in which case the thread remains on the processor. We model this selection of a new thread by the function *Select_next_thread*.

Review Note:

This function should be related to the *RunQueue* which is currently in the specification of the **switch** request but which should probably be moved to the state chapter. Also, what is the relationship between *Select_next_thread* and *thread_sched_priority*? Would it be useful to model this relationship?

| *Select_next_thread* : *PROCESSOR_SET* \mapsto *THREAD*

If the scheduling algorithm selects the blocking thread then that thread is marked as not swapped out and not uninterruptibly waiting, and the blocking operation is completed. Otherwise, the new thread selected by the scheduling algorithm receives these markings and, unless the blocking thread is being terminated, the following changes are made to the blocking thread:

1. It is added to \underline{s} wapped_threads.
2. It is marked as not *Running* if it is not in \underline{i} dle_threads, and it was either marked as
 - (a) *Stopped*, but not *Uninterruptible*, or
 - (b) *Waiting*.

The schema *ThreadBlock* describes these changes. The component *blocking_thread* is the thread that is being blocked, and *init_run_state* is the run state in effect when the thread blocking occurs. This may differ from the *run_state* function depending upon the context in which the blocking occurs. For example, when blocking occurs as part of a **thread_suspend** request the *Stopped* state will have been added to *run_state*(*blocking_thread*) to obtain *init_run_state*. In the case where the blocking thread is being terminated *blocking_thread* is not in the domain of *init_run_state*.

Review Note:

The need for *init_run_state* originates in our level of granularity in the specification. There are changes that various requests make to the run state of a thread in preparation for blocking the thread. Since *ThreadBlock* models only a portion of this processing, we need a way to specify what changes have been made to the run state in the request processing that precedes the blocking.

ThreadBlock

Δ *ThreadExecStatus*
ProcessorAndProcessorSet
blocking_thread : *THREAD*
cpu?? : *PROCESSOR*
init_run_state : *THREAD* \leftrightarrow \mathbb{P} *RUN_STATES*

cpu?? \in dom *proc_assigned_procset*
let *new_thread* == *Select_next_thread*(*proc_assigned_procset*(*cpu??*))
 • *new_thread* \notin \underline{s} wapped_threads'
 \wedge *run_state'*(*new_thread*) = *init_run_state*(*new_thread*) \setminus {*Uninterruptible*}
 \wedge ((*new_thread* \neq *blocking_thread* \wedge *blocking_thread* \notin dom *init_run_state*)
 \Rightarrow (*blocking_thread* \in \underline{s} wapped_threads'
 \wedge *run_state'*(*blocking_thread*)
 = **if** *blocking_thread* \notin \underline{i} dle_threads
 \wedge (*init_run_state*(*blocking_thread*)
 \cap {*Stopped*, *Uninterruptible*} = {*Stopped*}
 \vee *Waiting* \in *init_run_state*(*blocking_thread*))
then *init_run_state*(*blocking_thread*) \setminus {*Running*}
else *init_run_state*(*blocking_thread*)))
 \wedge (\forall *thread* : *THREAD* | *thread* \notin {*new_thread*, *blocking_thread*}
 • *run_state'*(*thread*) = *init_run_state*(*thread*)
 \wedge *thread* \in \underline{s} wapped_threads' \Leftrightarrow *thread* \in \underline{s} wapped_threads
 \wedge *thread* \in \underline{i} dle_threads' \Leftrightarrow *thread* \in \underline{i} dle_threads)

A request may also wait for a given thread to stop running. The component *stopping_thread* is the thread being stopped, and *init_run_state* is defined as for *ThreadBlock*.

ThreadDoWait $\Delta \text{ThreadExecStatus}$ $\text{stopping_thread} : \text{THREAD}$ $\text{init_run_state} : \text{THREAD} \mapsto \mathbb{P} \text{RUN_STATES}$ <hr/> $\text{run_state}' = \text{init_run_state}$ $\oplus \{ \text{stopping_thread} \mapsto \text{init_run_state}(\text{stopping_thread}) \setminus \{ \text{Running} \} \}$ $\text{swapped_threads}' = \text{swapped_threads}$ $\text{idle_threads}' = \text{idle_threads}$
--

Some requests (e.g., **thread_set_state** and **thread_get_state**) must wait for a thread to stop before they can do their work. When they are done modifying or observing the characteristics of the stopped thread they allow the thread to start again. For example, **thread_get_state** stops the thread, examines its machine state (e.g., machine registers) and then allows the thread to run again. The cumulative effect of this sequence of operations might include the side-effect of altering the run state. The run state will contain *Running* when it previously contained neither *Waiting* nor *Stopped*. It will contain *Halted* when it previously contained both *Halted* and *Stopped*. The *Stopped*, *Waiting* and *Uninterruptible* characteristics are unchanged.

Review Note:

We believe that $\text{Halted} \Rightarrow \text{Stopped}$ at the termination of a request. If this is true then the thread will be halted if and only if it was halted before the request. We also believe that at least one of the states *Running*, *Stopped* and *Waiting* must be contained in the run state. This means that *Running* could be removed from the run state by this operation, but never added.

$\text{ThreadDoWaitThenRelease}$ $\Delta \text{ThreadExecStatus}$ $\text{stopping_thread} : \text{THREAD}$ <hr/> $\forall \text{thread} : \text{THREAD} \mid \text{thread} \in \text{dom } \text{run_state} \wedge \text{thread} \neq \text{stopping_thread}$ <ul style="list-style-type: none"> • $\text{run_state}'(\text{thread}) = \text{run_state}(\text{thread})$ $\text{run_state}'(\text{stopping_thread}) \cap \{ \text{Stopped}, \text{Waiting}, \text{Uninterruptible} \}$ $= \text{run_state}(\text{stopping_thread}) \cap \{ \text{Stopped}, \text{Waiting}, \text{Uninterruptible} \}$ $\text{Halted} \in \text{run_state}'(\text{stopping_thread})$ $\Leftrightarrow \{ \text{Halted}, \text{Stopped} \} \subseteq \text{run_state}(\text{stopping_thread})$ $\text{Running} \in \text{run_state}'(\text{stopping_thread})$ $\Leftrightarrow \text{run_state}(\text{stopping_thread}) \cap \{ \text{Waiting}, \text{Stopped} \} = \emptyset$
--

9.1.4.4 Parameter Packaging Functions When invoking a kernel request, the following functions package the input parameters into a message body:

$$\begin{aligned}
& \text{Name_and_number_to_text} : \text{NAME} \times \mathbb{Z} \rightarrow \text{MESSAGE_BODY} \\
& \text{Name_to_text} : \text{NAME} \rightarrow \text{MESSAGE_BODY} \\
& \text{Number_and_boolean_to_text} : \mathbb{Z} \times \text{BOOLEAN} \rightarrow \text{MESSAGE_BODY} \\
& \text{Policy_and_data_to_text} : \text{SCHED_POLICY} \times \text{SCHED_POLICY_DATA} \\
& \quad \rightarrow \text{MESSAGE_BODY} \\
& \text{Sample_type_set_to_text} : \mathbb{P} \text{SAMPLE_TYPES} \rightarrow \text{MESSAGE_BODY} \\
& \text{Sequence_number_to_text} : \mathbb{N} \rightarrow \text{MESSAGE_BODY} \\
& \text{Thread_info_type_and_count_to_text} : \text{THREAD_INFO_TYPE} \times \mathbb{N} \\
& \quad \rightarrow \text{MESSAGE_BODY} \\
& \text{Thread_set_state_params_to_text} : \\
& \quad \text{THREAD_STATE_INFO_TYPES} \times \text{THREAD_STATE_INFO} \times \mathbb{N} \\
& \quad \rightarrow \text{MESSAGE_BODY} \\
& \text{Thread_special_port_and_name_to_text} : \text{THREAD_SPECIAL_PORTS} \times \text{NAME} \\
& \quad \rightarrow \text{MESSAGE_BODY} \\
& \text{Thread_special_ports_to_text} : \text{THREAD_SPECIAL_PORTS} \rightarrow \text{MESSAGE_BODY} \\
& \text{Thread_state_info_type_and_number_to_text} : \text{THREAD_STATE_INFO_TYPES} \times \mathbb{N} \\
& \quad \rightarrow \text{MESSAGE_BODY}
\end{aligned}$$

When creating a reply message from a request, the following functions package the output parameters into a kernel reply:

$$\begin{aligned}
& \text{Return_capability} : \text{Capability} \rightarrow \text{KERNEL_REPLY} \\
& \text{Return_sample_cnt} : \mathbb{N} \rightarrow \text{KERNEL_REPLY} \\
& \text{Return_samples} : (\mathbb{N} \times (\text{seq SAMPLE}) \times \mathbb{Z}) \rightarrow \text{KERNEL_REPLY} \\
& \text{Return_thread_info} : \text{THREAD_INFO} \times \mathbb{N} \rightarrow \text{KERNEL_REPLY} \\
& \text{Return_thread_state_info} : \text{THREAD_STATE_INFO} \times \mathbb{N} \rightarrow \text{KERNEL_REPLY} \\
& \text{Return_ticks} : \mathbb{N}_1 \rightarrow \text{KERNEL_REPLY}
\end{aligned}$$

When receiving a reply message from the kernel the following functions unpack the message body to obtain the output parameters (including the return status):

$$\begin{aligned}
& \text{Text_to_count_and_status} : \text{MESSAGE_BODY} \rightarrow (\mathbb{N} \times \text{KERNEL_RETURN}) \\
& \text{Text_to_info_and_count_and_status} : \text{MESSAGE_BODY} \\
& \quad \rightarrow (\text{THREAD_INFO} \times \mathbb{N} \times \text{KERNEL_RETURN}) \\
& \text{Text_to_name_and_status} : \text{MESSAGE_BODY} \rightarrow (\text{NAME} \times \text{KERNEL_RETURN}) \\
& \text{Text_to_seqno_and_PCs_and_count_and_status} : \text{MESSAGE_BODY} \\
& \quad \rightarrow (\mathbb{N} \times \text{seq SAMPLE} \times \mathbb{Z} \times \text{KERNEL_RETURN}) \\
& \text{Text_to_state_and_count_and_status} : \text{MESSAGE_BODY} \\
& \quad \rightarrow (\text{THREAD_STATE_INFO} \times \mathbb{N} \times \text{KERNEL_RETURN}) \\
& \text{Text_to_status} : \text{MESSAGE_BODY} \rightarrow \text{KERNEL_RETURN} \\
& \text{Text_to_ticks_and_status} : \text{MESSAGE_BODY} \rightarrow (\mathbb{N}_1 \times \text{KERNEL_RETURN})
\end{aligned}$$

9.1.4.5 Destroying a Port

Review Note:

The following may fit more naturally in the port requests chapter.

The following schema describes the destruction of a port. This is required to describe the **thread_terminate** request. The port destroyed is removed from the set of existing ports, and

all send, receive, and send-once rights to this port are removed from all port name spaces. New dead names are created for all previous send or send-once rights to this port. Note that the creation of notifications when these names turn into dead names should be added to this schema.

Editorial Note:

Destruction of a port representing a message queue for IPC can cause a chain reaction not represented in this schema. Whether the same chain reaction is possible for ports representing kernel objects (as the schema is used in this chapter) is unclear.

$\begin{array}{l} \textit{PortDestroy} \\ \Delta \textit{Ipc} \\ \textit{port} : \textit{PORT} \\ \\ \underline{\textit{port_exists}}' = \underline{\textit{port_exists}} \setminus \{\textit{port}\} \\ \underline{\textit{port_right_rel}}' = \underline{\textit{port_right_rel}} \\ \quad \setminus \{ \textit{task} : \textit{TASK}; \textit{name} : \textit{NAME}; \textit{right} : \textit{RIGHT}; \textit{i} : \mathbb{N}_1 \\ \quad \bullet (\textit{task}, \textit{port}, \textit{name}, \textit{right}, \textit{i}) \} \\ \underline{\textit{make_send_count}}' = \{\textit{port}\} \triangleleft \underline{\textit{make_send_count}} \\ \underline{\textit{dead_right_rel}}' = \underline{\textit{dead_right_rel}} \\ \quad \cup \{ \textit{task} : \textit{TASK}; \textit{name} : \textit{NAME} \\ \quad \quad \textit{named_port}(\textit{task}, \textit{name}) = \textit{port} \\ \quad \bullet (\textit{task}, \textit{name}, 1) \} \end{array}$

9.1.4.6 **Miscellaneous** The function *Thread_port_to_s_right* takes a port and returns a send right to the port.

$\begin{array}{l} \textit{Thread_port_to_s_right} : \textit{PORT} \leftrightarrow \textit{Capability} \\ \\ \forall \textit{port} : \textit{PORT} \\ \bullet (\textit{Thread_port_to_s_right}(\textit{port})).\textit{right} = \textit{Send} \\ \wedge (\textit{Thread_port_to_s_right}(\textit{port})).\textit{port} = \textit{port} \end{array}$
--

The function *Thread_state_count* returns the size of a given type of thread state information.

$\textit{Thread_state_count} : \textit{THREAD_STATE_INFO_TYPES} \longrightarrow \mathbb{N}$
--

9.1.5 Kernel Processing

The kernel performs processing for a thread request only when it detects a break indicating that a request has been received through a port of the appropriate class, *Pc_task* or *Pc_thread*. If the specified service port no longer exists, then a *Kern_invalid_argument* status code is returned.

ProcessThreadRequestBadAux

ProcessRequest
 $\exists Mach$
reply_to_port! : $\mathbb{P} PORT$
reply! : *KERNEL_REPLY*
return! : *KERNEL_RETURN*

$((pc? = Pc_thread \wedge operation? \in Thread_thread_port_ops$
 $\wedge service_port? \notin \text{dom } self_thread)$
 $\vee (pc? = Pc_task \wedge operation? \in Thread_task_port_ops$
 $\wedge service_port? \notin \text{dom } self_task))$

reply_to_port! = *reply_to_port?*
return! = *Kern_invalid_argument*

ProcessThreadRequestBad $\hat{=}$ *ProcessThreadRequestBadAux* \gg *RequestNoOp*

Otherwise, the kernel processes the request. In this case, we use the following schema to represent the parameters to thread requests which are processed via thread ports:

Editorial Note:
flavor? is omitted because it is used with two different types in different requests.

ThreadParameters

data? : *SCHED_POLICY_DATA*
new_state? : *THREAD_STATE_INFO*
new_state_cnt? : \mathbb{N}
old_state_cnt? : \mathbb{N}
policy? : *SCHED_POLICY*
priority? : \mathbb{Z}
procset? : *PROCESSOR_SET*
seqno? : \mathbb{N}
set_max? : *BOOLEAN*
special_port? : *PORT*
task? : *TASK*
thread? : *THREAD*
target_thread? : *THREAD*
thread_infoCnt? : \mathbb{N}
which_port? : *THREAD_SPECIAL_PORTS*

The interpretations of the components of these schemas are:

- *data?* — policy specific data used with the scheduling policy to determine the scheduling priority of a thread (**thread_policy**)
- *flavor?* — specific type of information (**thread_info**) or state information (**thread_↔get_state**, **thread_set_state** and **thread_set_state_secure**)
- *new_state?* — state information to be assigned to a thread (**thread_set_state** and **thread_set_state_secure**)

- *new_state_cnt?* — amount of storage provided by a client to hold state information to be assigned to a thread (**thread_set_state** and **thread_set_state_secure**)
- *old_state_cnt?* — amount of storage provided by a client to hold state information (**thread_get_state**)
- *policy?* — desired scheduling policy (**thread_policy**)
- *priority?* — desired priority for a thread (**thread_priority** and **thread_max_priority**)
- *procset?* — desired processor set for a thread (**thread_assign**); the control port for the processor set to which a thread is currently assigned (**thread_max_priority**)
- *seqno?* — the sequence number of the first sample that should be returned (**thread_↔get_sampled_pcs**)
- *set_max?* — a flag indicating whether the maximum priority should be reset when the *priority* is changed (**thread_priority**)
- *special_port?* — a port specified by the client to become the special port for the target thread (for **thread_set_special_port**)
- *task?* — the target task for the request (**thread_create** and **thread_create_secure**)
- *target_thread?* — the target thread for the request (an alternative name for *thread?* that is used in some requests)
- *thread?* — the target thread for the request
- *thread_infoCnt?* — amount of storage provided by a client to hold information on a thread (**thread_info**)
- *which_port?* — the type of special port specified by the client (**thread_get_special_port** and **thread_set_special_port**)

The following schema maps a message sent to a thread port to a value of type *ThreadParameters*:

<p><i>ThreadMessageToThreadParameters</i> _____</p> <p><i>Request?</i> <i>SpecialThreadPorts</i> <i>ThreadParameters</i></p> <hr style="width: 20%; margin-left: 0;"/> <p><i>pc?</i> = <i>Pc_thread</i> <i>operation?</i> ∈ <i>Thread_thread_port_ops</i> <i>service_port?</i> ∈ dom <i>self_thread</i> <i>thread?</i> = <i>self_thread(service_port?)</i> <i>target_thread?</i> = <i>thread?</i></p>

Review Note:

ThreadInvariants really belongs in the state changes schemas rather than here in the message processing. What I want here is \exists almosteverything. The same goes for the use a couple schemas further down.

<i>ProcessThreadViaThreadPortRequestGood</i> <i>ProcessRequest</i> <i>ThreadInvariants</i> <i>ThreadMessageToThreadParameters</i>
--

Similarly, we use the following function to map a message sent to a task port to a value of type *ThreadParameters*:

<i>TaskMessageToThreadParameters</i> <i>Request?</i> <i>SpecialTaskPorts</i> <i>ThreadParameters</i>
$pc? = Pc_task$ $operation? \in Thread_task_port_ops$ $service_port? \in \text{dom } self_task$ $task? = self_task(service_port?)$

<i>ProcessThreadViaTaskPortRequestGood</i> <i>ProcessRequest</i> <i>ThreadInvariants</i> <i>TaskMessageToThreadParameters</i>
--

We now describe the individual thread requests.

9.2 thread_abort

The request **thread_abort** helps to cleanly stop a thread by interrupting page faults and any **mach_msg** calls in progress by the thread. It causes an interrupt return code to be returned from any system trap in progress on behalf of the thread (even though the execution of the trap may finish). It also aborts any priority depressions. Note that **thread_abort** does not suspend a thread. If the thread did not already have the *Stopped* state, then at the conclusion of a **thread_abort** request it is neither *Stopped* nor *Halted*.

9.2.1 Client Interface

```
kern_return_t thread_abort
(mach_port_t target_thread_name);
```

9.2.1.1 Input Parameters The following input parameters are provided by the client of a **thread_abort** request:

- *target_thread_name?* — the client's name for the thread to which the abort will be applied

<i>ThreadAbortClientInputs</i> <i>target_thread_name? : NAME</i>

A **thread_abort** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_abort_id* and has no body.

<i>Invoke ThreadAbort</i> <i>Invoke MachMsg</i> <i>ThreadAbortClientInputs</i>
<i>name?</i> = <i>target_thread_name?</i> <i>operation?</i> = <i>Thread_abort_id</i>

9.2.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_abort** request:

- *return!* — the status of the request

<i>ThreadAbortClientOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>
--

<i>ThreadAbortReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadAbortClientOutputs</i>
<i>return!</i> = <i>Text_to_status(msg_body)</i>

9.2.2 Kernel Interface

9.2.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_abort** request:

- *target_thread?* — the thread to which the abort will be applied

<i>ThreadAbortInputs</i> <i>target_thread?</i> : <i>THREAD</i>

9.2.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_abort** request:

- *return!* — the status of the request

<i>ThreadAbortOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>
--

9.2.3 Request Criteria

The following criteria are defined for the **thread_abort** request.

- **C1** — The parameter *target_thread?* is the client thread (i.e., the thread currently active on the CPU).

$C1ThreadAbortClientThread$ $ThreadsAndProcessors$ $cpu?? : PROCESSOR$ $target_thread? : THREAD$
$cpu?? \in \text{dom } \underline{active_thread}$
$target_thread? = \underline{active_thread}(cpu??)$

$$NotC1ThreadAbortClientThread \\ \hat{=} ThreadsAndProcessors \wedge \neg C1ThreadAbortClientThread$$

9.2.4 Return Values

Table 17 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>	C1
<i>Kern_invalid_argument</i>	T
<i>Kern_success</i>	F

Table 17: Return Values for **thread_abort**

Review Note:

thread_abort can return *Kern_aborted* when

- there is a cycle of halt operations, or
- the client thread is interrupted while waiting for the target thread to halt.

IPC will convert *Kern_aborted* to an IPC interrupted error code. This behavior is not modeled.

$RVThreadAbortInvalidArgument$ $C1ThreadAbortClientThread$ $ThreadAbortOutputs$
$return! = Kern_invalid_argument$

$RVThreadAbortGood$ $NotC1ThreadAbortClientThread$ $ThreadAbortOutputs$
$return! = Kern_success$

9.2.5 State Changes

A successful **thread_abort** request will interrupt page faults and message primitive calls in use by the thread. The thread will resume execution at the point of return from the interrupted system call. This will occur upon return from this request unless the thread is in a *Stopped* state, in which case it will occur when the thread is resumed via **thread_resume**.

Review Note:

The granularity of the FTLS is not fine enough to model the interruption of page faults and message primitive calls in use by the thread.

As stated above, **thread_abort** does not suspend a thread. If the thread does not already have the *Stopped* state, then at the conclusion of a **thread_abort** request it is neither *Stopped* nor *Halted*. If the thread has been previously *Stopped*, the thread will remain *Stopped* upon completion of the **thread_abort** request until it is resumed. The thread will also be *Halted* at this point since **thread_abort** has insured that it is stopped at a clean point. A thread that is not already *Stopped* and is not *Waiting* will have *Running* added to its run state by **thread_abort** (assuming it is not already there).

ThreadAbortExecStatus
 Δ *ThreadExecStatus*
 $target_thread? : THREAD$
 $\forall thread : THREAD \mid thread \in \text{dom } \underline{run_state} \wedge thread \neq target_thread?$

- $\underline{run_state}'(thread) = \underline{run_state}(thread)$

 $\underline{run_state}'(target_thread?) \cap \{Stopped, Waiting, Uninterruptible\}$
 $= \underline{run_state}(target_thread?) \cap \{Stopped, Waiting, Uninterruptible\}$
 $Halted \in \underline{run_state}'(target_thread?) \Leftrightarrow Stopped \in \underline{run_state}(target_thread?)$
 $Running \in \underline{run_state}'(target_thread?)$
 $\Leftrightarrow \underline{run_state}(target_thread?) \cap \{Waiting, Stopped\} = \emptyset$
 $\underline{swapped_threads}' = \underline{swapped_threads}$
 $\underline{idle_threads}' = \underline{idle_threads}$
 $\underline{thread_suspend_count}' = \underline{thread_suspend_count}$
 $\underline{threads_wired}' = \underline{threads_wired}$

The sending of the return from this request means that the thread has received an interrupt return code from the program it was executing. It follows from the specification for the processing of invocable requests and **mach_msg** that the receipt of the return of *Kern_success* from this request will therefore occur only if the thread is not in *Stopped* state. If the target thread is in a stopped state, it will receive the return value when (and if) it is resumed (via **thread_resume**).

Any priority depression is also aborted. This returns the priority of the thread to its value before the depression. Note that the scheduling priority may also change, but since we do not have enough detail in our model to compute its value we will leave it unspecified.

$\begin{aligned} & \text{ThreadAbortPriority} \\ & \Delta \text{ThreadPri} \\ & \text{target_thread?} : \text{THREAD} \\ & \text{target_thread?} \in \text{dom } \underline{\text{priority_before_depression}} \\ & \underline{\text{thread_priority}}' = \underline{\text{thread_priority}} \\ & \quad \oplus \{ \text{target_thread?} \mapsto \underline{\text{priority_before_depression}}(\text{target_thread?}) \} \\ & \underline{\text{thread_max_priority}}' = \underline{\text{thread_max_priority}} \\ & \underline{\text{depressed_threads}}' = \underline{\text{depressed_threads}} \setminus \{ \text{target_thread?} \} \\ & \underline{\text{priority_before_depression}}' = \underline{\text{priority_before_depression}} \end{aligned}$

When the target thread resumes execution it will be at the return point from any interrupt of trap it might have been executing. The component *at_call_return* represents the address at which execution will resume if the thread is resumed.

Editorial Note:
Our model is not detailed enough to formally describe the value of *at_call_return*.

$\begin{aligned} & \text{ThreadAbortState} \\ & \text{ThreadInvariants} \\ & \text{ThreadAbortExecStatus} \\ & \text{ThreadAbortPriority} \\ & \exists \text{ThreadExist} \\ & \exists \text{PortExist} \\ & \Delta \text{Threads} \\ & \exists \text{TasksAndThreads} \\ & \exists \text{ThreadSchedPolicy} \\ & \Delta \text{ThreadInstruction} \\ & \Delta \text{Events} \\ & \exists \text{PortNameSpace} \\ & \exists \text{SpecialPurposePorts} \\ & \text{target_thread?} : \text{THREAD} \\ & \text{at_call_return} : \text{VIRTUAL_ADDRESS} \\ & \underline{\text{instruction_pointer}}' = \underline{\text{instruction_pointer}} \oplus \{ \text{target_thread?} \mapsto \text{at_call_return} \} \end{aligned}$

Review Note:
How can we represent here the fact that the execution of messages and traps might not be completed?
There might be some delay in halting the thread. Is this important?

9.2.6 Complete Request

The following schema defines the general form of a **thread_abort** request.

$\begin{aligned} & \text{Processing ThreadAbort} \\ & \text{Process Thread Via ThreadPortRequestGood} \\ & \text{operation?} = \text{Thread_abort_id} \end{aligned}$
--

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} \text{ThreadAbortGood} &\hat{=} (\text{RVThreadAbortGood} \wedge \text{ThreadAbortState}) \\ &\gg \text{RequestReturnOnlyStatus} \end{aligned}$$

An unsuccessful request returns an error status.

$$\text{ThreadAbortBad} \hat{=} \text{RVThreadAbortInvalidArgument} \gg \text{RequestNoOp}$$

Execution of the request consists of a good execution or an error execution.

$$\text{ExecuteThreadAbort} \hat{=} (\text{ThreadAbortGood} \vee \text{ThreadAbortBad}) \setminus (\text{at_call_return})$$

The full specification for kernel processing of a validated **thread_abort** request consists of processing the request followed by its execution.

$$\text{ThreadAbort} \hat{=} \text{ProcessingThreadAbort} \wp \text{ExecuteThreadAbort}$$

9.3 thread_create and thread_create_secure

The requests **thread_create** and **thread_create_secure** create a new thread within an existing task. The name of a send right to the kernel port of the new thread is returned. The **thread_create_secure** request (which is used in the secure initiation of threads within a task) expects the parent task to have task creation state *Tcs_task_empty* (see Section 5.7). It modifies the state to *Tcs_thread_created*.

9.3.1 Client Interface

```
kern_return_t thread_create
    (mach_port_t
     mach_port_t*
     parent_task_name,
     child_thread_name);
```

```
kern_return_t thread_create_secure
    (mach_port_t
     mach_port_t*
     parent_task_name,
     child_thread_name);
```

9.3.1.1 Input Parameters The following input parameters are provided by the client of a **thread_create** or **thread_create_secure** request:

- *parent_task_name?* — the client's name for the task that will be the parent for the newly created thread

<p><i>ThreadCreateClientInputs</i></p> <p><i>parent_task_name?</i> : NAME</p>

A **thread_create** request is invoked by sending a message to the port indicated by *parent_task_name?* that has the operation field set to *Thread_create_id* and has no body.

<i>Invoke ThreadCreate</i> <i>Invoke MachMsg</i> <i>ThreadCreate ClientInputs</i>
<i>name?</i> = <i>parent_task_name?</i> <i>operation?</i> = <i>Thread_create_id</i>

A **thread_create_secure** request is invoked by sending a message to the port indicated by *parent_task_name?* that has the operation field set to *Thread_create_secure_id* and has no body.

<i>Invoke ThreadCreateSecure</i> <i>Invoke MachMsg</i> <i>ThreadCreate ClientInputs</i>
<i>name?</i> = <i>parent_task_name?</i> <i>operation?</i> = <i>Thread_create_secure_id</i>

9.3.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_create** or **thread_create_secure** request:

- *child_thread_name!* — the name of a send right to the kernel port of the new thread
- *return!* — the status of the request

<i>ThreadCreate ClientOutputs</i> <i>child_thread_name!</i> : <i>NAME</i> <i>return!</i> : <i>KERNEL_RETURN</i>

<i>ThreadCreate ReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadCreate ClientOutputs</i> <i>(child_thread_name!, return!) = Text_to_name_and_status(msg_body)</i>

9.3.2 Kernel Interface

9.3.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_create** or **thread_create_secure** request:

- *parent_task?* — the task that will be the parent for the newly created thread

<i>ThreadCreate Inputs</i> <i>parent_task?</i> : <i>TASK</i>

9.3.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_create** or **thread_create_secure** request:

- *child_thread!* — the new thread
- *return!* — the status of the request

<i>ThreadCreateOutputs</i> <i>child_thread!</i> : <i>THREAD</i> <i>return!</i> : <i>KERNEL_RETURN</i>

Upon completion of the processing of a **thread_create** or **thread_create_secure** request a reply message is built from the output parameters. The reply message will contain a send right for the created thread's kernel port.

<i>ThreadCreateReply</i> <i>RequestReturn</i> <i>child_thread?</i> : <i>THREAD</i> <i>reply?</i> = <i>Return_capability(Thread_port_to_s_right(thread_sself(child_thread?)))</i>

9.3.3 Request Criteria

The following criteria are defined for the **thread_create** and **thread_create_secure** requests.

- **C1** — The kernel has the necessary resources available to create the thread. We do not actually model the consumption of resources by the kernel. So, we will use the set *Resources_available_to_create_thread* to indicate the set of states where there are sufficient resources to create a thread.

| *Resources_available_to_create_thread* : \mathbb{P} *DtosExec*

<i>C1ThreadCreateResourcesAvailable</i> <i>DtosExec</i> <hr/> $\theta DtosExec \in Resources_available_to_create_thread$
--

NotC1ThreadCreateResourcesAvailable $\hat{=}$
 $DtosExec \wedge \neg C1ThreadCreateResourcesAvailable$

- **C2** — The task creation state of the parent task must be *Tcs_task_empty*. This criterion applies only to the **thread_create_secure** request.

<i>C2ThreadCreateSecureTaskEmpty</i> <i>TaskCreationState</i> <u><i>parent_task?</i></u> : <i>TASK</i> <hr/> <u><i>parent_task?</i></u> \in $\text{dom } \underline{task_creation_state}$ <u><i>task_creation_state</i></u> (<u><i>parent_task?</i></u>) = <i>Tcs_task_empty</i>
--

NotC2ThreadCreateSecureTaskEmpty
 $\hat{=}$ *TaskCreationState* $\wedge \neg C2ThreadCreateSecureTaskEmpty$

9.3.4 Return Values

Table 18 describes the values returned at the completion of the **thread_create** request and the conditions under which each value is returned. The value *thread* is the newly created thread (see the State Changes section). The design does not specify the value of *child_thread!* when an error occurs. It depends on the implementation, and we leave it unspecified.

<i>child_thread!</i>	<i>return!</i>	C1
<i>thread</i>	<i>Kern_success</i>	T
—	<i>Kern_resource_shortage</i>	F

Table 18: Return Values for **thread_create**

<i>RVThreadCreateGood</i>
<i>C1ThreadCreateResourcesAvailable</i>
<i>ThreadCreateOutputs</i>
<i>thread</i> : <i>THREAD</i>
<i>child_thread!</i> = <i>thread</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadCreateResourceShortage</i>
<i>NotC1ThreadCreateResourcesAvailable</i>
<i>ThreadCreateOutputs</i>
<i>return!</i> = <i>Kern_resource_shortage</i>

Table 19 describes the values returned at the completion of the **thread_create_secure** request and the conditions under which each value is returned. In the case where both C1 and C2 are false we assume that *Kern_insufficient_permission* is returned.

Review Note:
In the prototype the conditions are checked in the order C2, C1.

<i>child_thread!</i>	<i>return!</i>	C1	C2
<i>thread</i>	<i>Kern_success</i>	T	T
—	<i>Kern_resource_shortage</i>	F	T
—	<i>Kern_insufficient_permission</i>	-	F

Table 19: Return Values for **thread_create_secure**

<i>RVThreadCreateSecureGood</i>
<i>C1ThreadCreateResourcesAvailable</i>
<i>C2ThreadCreateSecureTaskEmpty</i>
<i>ThreadCreateOutputs</i>
<i>thread</i> : <i>THREAD</i>
<i>child_thread!</i> = <i>thread</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadCreateSecureResourceShortage</i> <i>NotC1ThreadCreateResourcesAvailable</i> <i>C2ThreadCreateSecureTaskEmpty</i> <i>ThreadCreateOutputs</i>
<i>return!</i> = <i>Kern_resource_shortage</i>

<i>RVThreadCreateSecureInsufficientPermission</i> <i>NotC2ThreadCreateSecureTaskEmpty</i> <i>ThreadCreateOutputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

9.3.5 State Changes

A successful **thread_create** or **thread_create_secure** request creates a new thread. The OSF documentation for this request states that, in addition, a send right to the thread's kernel port is given to the containing task. This is not shown explicitly here. We believe that the existence of a new thread self port is an "implicit" send right, not in the port name space (and not usable) for the containing task until the thread executes **amach_thread_self** request.

The creation of a new thread affects much of the state information associated with threads. We will consider each type of state information individually. We first define the things that do not change in a successful **thread_create** or **thread_create_secure** request. Note that the port name space of the receiving task on the reply port for this request will change after the invocable request created by the schema *Return* is processed, and not immediately upon completion of this request.

<i>ThreadCreateInvariants</i> <i>ThreadInvariants</i> \exists <i>PortNameSpace</i>
--

A new thread is created and added to the list of threads associated with the parent task.

<i>ThreadCreateTasksAndThreads</i> Δ <i>ThreadExist</i> Δ <i>TasksAndThreads</i> $\underline{parent_task?} : TASK$ $\underline{thread} : THREAD$
$\underline{thread} \notin \underline{thread_exists}$ $\underline{thread_exists}' = \underline{thread_exists} \cup \{\underline{thread}\}$ $\underline{task_thread_rel}' = \underline{task_thread_rel} \cup \{(\underline{parent_task?}, \underline{thread})\}$

A newly created thread takes its maximum priority to be the lower of the following two priorities:

- the maximum priority of the processor set to which it is assigned, or

- the *Base_user_priority* constant.

It takes its priority to be the lower of its parent task's priority and its own maximum priority. No other thread's priorities change due solely to the creation of this thread.

$\begin{array}{l} \text{ThreadCreatePriority} \\ \Delta \text{ThreadPri} \\ \text{TaskPriority} \\ \Delta \text{ThreadAndProcessorSet} \\ \underline{\text{parent_task?}} : \text{TASK} \\ \underline{\text{thread}} : \text{THREAD} \\ \\ \underline{\text{thread}} \in \text{dom } \underline{\text{thread_assigned_to'}} \\ \underline{\text{thread_assigned_to'}}(\underline{\text{thread}}) \in \text{dom } \underline{\text{ps_max_priority}} \\ \underline{\text{parent_task?}} \in \text{dom } \underline{\text{task_priority}} \\ \\ \underline{\text{thread_max_priority'}} = \underline{\text{thread_max_priority}} \\ \cup \{ \text{thread} \mapsto \text{Lowest_priority}(\{ \underline{\text{ps_max_priority}}(\underline{\text{thread_assigned_to'}}(\text{thread})), \\ \text{Base_user_priority} \}) \} \\ \underline{\text{thread_priority'}} = \underline{\text{thread_priority}} \\ \cup \{ \text{thread} \mapsto \text{Lowest_priority}(\{ \underline{\text{task_priority}}(\underline{\text{parent_task?}}), \\ \underline{\text{thread_max_priority'}}(\text{thread}) \}) \} \\ \text{dom } \underline{\text{thread_sched_priority'}} = \text{dom } \underline{\text{thread_sched_priority}} \cup \{ \text{thread} \} \\ \underline{\text{thread_sched_priority}} \subset \underline{\text{thread_sched_priority'}} \\ \underline{\text{depressed_threads'}} = \underline{\text{depressed_threads}} \\ \underline{\text{priority_before_depression'}} \\ = \underline{\text{priority_before_depression}} \cup \{ \text{thread} \mapsto \underline{\text{thread_priority'}}(\text{thread}) \} \end{array}$
--

The new thread's scheduling policy is *Timeshare*. Since the *Timeshare* policy does not require any scheduling policy data, there is no change to *thread_sched_policy_data*.

$\begin{array}{l} \text{ThreadCreateSchedPolicy} \\ \Delta \text{ThreadSchedPolicy} \\ \underline{\text{thread}} : \text{THREAD} \\ \\ \underline{\text{thread_sched_policy'}} = \underline{\text{thread_sched_policy}} \cup \{ \text{thread} \mapsto \text{Timeshare} \} \\ \underline{\text{thread_sched_policy_data'}} = \underline{\text{thread_sched_policy_data}} \end{array}$
--

The thread is created in a *Stopped* run state, and it is swapped out. Its suspend count is one larger than the suspend count of its parent task.

<p><i>ThreadCreateExecStatus</i></p> <p>Δ <i>ThreadExecStatus</i></p> <p><i>TaskSuspendCount</i></p> <p><i>thread</i> : <i>THREAD</i></p> <p><i>parent_task?</i> : <i>TASK</i></p> <hr/> <p><i>parent_task?</i> \in dom <i>task_suspend_count</i></p> <p><i>run_state'</i> = <i>run_state</i> \cup { <i>thread</i> \mapsto { <i>Stopped</i> } }</p> <p><i>swapped_threads'</i> = <i>swapped_threads</i> \cup { <i>thread</i> }</p> <p><i>thread_suspend_count'</i> = <i>thread_suspend_count</i> \cup { <i>thread</i> \mapsto <i>task_suspend_count</i>(<i>parent_task?</i>) + 1 }</p> <p><i>threads_wired'</i> = <i>threads_wired</i></p>

All of the thread's timing statistics are set to zero.

<p><i>ThreadCreateStatistics</i></p> <p>Δ <i>ThreadStatistics</i></p> <p><i>thread</i> : <i>THREAD</i></p> <hr/> <p><i>user_time'</i> = <i>user_time</i> \cup { <i>thread</i> \mapsto 0 }</p> <p><i>system_time'</i> = <i>system_time</i> \cup { <i>thread</i> \mapsto 0 }</p> <p><i>cpu_time'</i> = <i>cpu_time</i> \cup { <i>thread</i> \mapsto 0 }</p> <p><i>sleep_time'</i> = <i>sleep_time</i> \cup { <i>thread</i> \mapsto 0 }</p>
--

A new self port is created for the thread. This port is assigned to be the kernel (sself) port as well. There is no exception port assigned to the thread.

<p><i>ThreadCreateSpecialPorts</i></p> <p>Δ <i>PortExist</i></p> <p>Δ <i>SpecialThreadPorts</i></p> <p><i>thread</i> : <i>THREAD</i></p> <p><i>port</i> : <i>PORT</i></p> <hr/> <p><i>port</i> \notin <i>port_exists</i></p> <p><i>port_exists'</i> = <i>port_exists</i> \cup { <i>port</i> }</p> <p><i>thread_self'</i> = <i>thread_self</i> \cup { <i>thread</i> \mapsto <i>port</i> }</p> <p><i>thread_sself'</i> = <i>thread_sself</i> \cup { <i>thread</i> \mapsto <i>port</i> }</p> <p><i>thread_eport'</i> = <i>thread_eport</i></p>

The thread will be assigned to the processor set to which its parent task is assigned.

ThreadCreateThreadAndProcessorSet

Δ *ThreadAndProcessorSet*
TaskAndProcessorSet
 $\underline{parent_task?} : TASK$
 $\underline{thread} : THREAD$

$\underline{parent_task?} \in \text{dom } \underline{task_assigned_to}$

$\underline{thread_assignment_rel}' = \underline{thread_assignment_rel}$
 $\cup \{(thread, \underline{task_assigned_to}(\underline{parent_task?}))\}$

$\underline{enabled_sp}' = \underline{enabled_sp}$
 $\underline{ps_max_priority}' = \underline{ps_max_priority}$

ThreadCreateState

Δ *Threads*
ThreadCreateInvariants
ThreadCreateTasksAndThreads
ThreadCreatePriority
ThreadCreateSchedPolicy
ThreadCreateExecStatus
ThreadCreateStatistics
ThreadCreateSpecialPorts
ThreadCreateThreadAndProcessorSet
 Δ *SpecialPurposePorts*

For the **thread_create_secure** request the task creation state of the parent task is changed to *Tcs_thread_created*. There is no change to the task creation state of the parent task for a **thread_create** request.

ThreadCreateSecureState

Δ *TaskCreationState*
 $\underline{parent_task?} : TASK$
 $\underline{operation?} : OPERATION$

$(\underline{operation?} = \underline{Thread_create_secure_id}$
 $\wedge \underline{task_creation_state}' = \underline{task_creation_state}$
 $\oplus \{\underline{parent_task?} \mapsto Tcs_thread_created\})$
 $\vee (\underline{operation?} = \underline{Thread_create_id}$
 $\wedge \underline{task_creation_state}' = \underline{task_creation_state})$

The new port gets a SID based upon the parent task.

ThreadCreateDtosState

Δ *PortSid*
SubjectSid
 $\underline{port} : PORT$
 $\underline{parent_task?} : TASK$

$\underline{port_sid}' = \underline{port_sid} \cup \{\underline{port} \mapsto \underline{Thread_port_sid}(\underline{task_sid}(\underline{parent_task?}))\}$

9.3.6 Complete Request

The following schemas define the general form of the **thread_create** and **thread_create_↔secure** requests.

$\text{Processing Thread Create} \xrightarrow{\text{Process Thread Via Task Port Request Good}}$ $\text{operation?} = \text{Thread_create_id}$
--

$\text{Processing Thread Create Secure} \xrightarrow{\text{Process Thread Via Task Port Request Good}}$ $\text{operation?} = \text{Thread_create_secure_id}$

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} \text{Thread Create Good} &\hat{=} \\ &(\text{RVThread Create Good} \wedge \text{Thread Create State} \\ &\quad \wedge \text{Thread Create Secure State} \wedge \text{Thread Create Dtos State}) \\ &\gg \text{Thread Create Reply} \\ \text{Thread Create Secure Good} &\hat{=} \\ &(\text{RVThread Create Secure Good} \wedge \text{Thread Create State} \\ &\quad \wedge \text{Thread Create Secure State} \wedge \text{Thread Create Dtos State}) \\ &\gg \text{Thread Create Reply} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} \text{Thread Create Bad} &\hat{=} \text{RVThread Create Resource Shortage} \gg \text{Request No Op} \\ \text{Thread Create Secure Bad} &\hat{=} (\text{RVThread Create Secure Resource Shortage} \\ &\quad \vee \text{RVThread Create Secure Insufficient Permission}) \\ &\gg \text{Request No Op} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} \text{Execute Thread Create} &\hat{=} (\text{Thread Create Good} \vee \text{Thread Create Bad}) \setminus (\text{port}, \text{thread}) \\ \text{Execute Thread Create Secure} &\hat{=} (\text{Thread Create Secure Good} \vee \text{Thread Create Secure Bad}) \setminus (\text{port}, \text{thread}) \end{aligned}$$

The full specification for kernel processing of a validated **thread_create** or **thread_create_↔secure** request consists of processing the request followed by its execution.

$$\begin{aligned} \text{Thread Create} &\hat{=} \text{Processing Thread Create} \ ; \ \text{Execute Thread Create} \\ \text{Thread Create Secure} &\hat{=} \text{Processing Thread Create Secure} \ ; \ \text{Execute Thread Create Secure} \end{aligned}$$

9.4 thread_depress_abort

The request **thread_depress_abort** restores the original scheduling priority to a thread whose priority has been set to the lowest possible value by a **swtch_pri** or **thread_switch** request.

9.4.1 Client Interface

```
kern_return_t thread_depress_abort
                (mach_port_t thread_name);
```

9.4.1.1 Input Parameters The following input parameters are provided by the client of a **thread_depress_abort** request:

- *thread_name?* — the client's name for the thread whose priority depression will be canceled

```
ThreadDepressAbortClientInputs
thread_name? : NAME
```

A **thread_depress_abort** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_depress_abort_id* and has no body.

```
Invoke ThreadDepressAbort
InvokeMachMsg
ThreadDepressAbortClientInputs
name? = thread_name?
operation? = Thread_depress_abort_id
```

9.4.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_depress_abort** request:

- *return!* — the status of the request

```
ThreadDepressAbortClientOutputs
return! : KERNEL_RETURN
```

```
ThreadDepressAbortReceiveReply
InvokeMachMsgRcv
ThreadDepressAbortClientOutputs
return! = Text_to_status(msg_body)
```

9.4.2 Kernel Interface

9.4.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_depress_abort** request:

- *thread?* — the thread whose priority depression will be canceled

```
ThreadDepressAbortInputs
thread? : THREAD
```

9.4.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_depress_abort** request:

- *return!* — the status of the request

<i>ThreadDepressAbortOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

9.4.3 Request Criteria

No criteria are defined for the **thread_depress_abort** request.

9.4.4 Return Values

Table 20 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>
<i>Kern_success</i>

Table 20: Return Values for **thread_depress_abort**

<i>RVThreadDepressAbortGood</i>
<i>ThreadDepressAbortOutputs</i>
<i>return!</i> = <i>Kern_success</i>

9.4.5 State Changes

A successful **thread_depress_abort** request returns the priority of the thread to its value before the depression. If the priority of the thread is not currently depressed, no changes occur. Note that the scheduling priority may also change, but since we do not have enough detail in our model to compute its value we will leave it unspecified.

<p><i>ThreadDepressAbortState</i></p> <p>Δ <i>ThreadPri</i></p> <p>Δ <i>Threads</i></p> <p>Ξ <i>ThreadExist</i></p> <p>Ξ <i>PortExist</i></p> <p>Ξ <i>TasksAndThreads</i></p> <p>Ξ <i>ThreadSchedPolicy</i></p> <p>Ξ <i>ThreadInstruction</i></p> <p>Ξ <i>ThreadMachine.State</i></p> <p>Δ <i>ThreadStatistics</i></p> <p>Δ <i>Events</i></p> <p>Ξ <i>PortName.Space</i></p> <p>Ξ <i>SpecialPurposePorts</i></p> <p>Ξ <i>ThreadAndProcessorSet</i></p> <p><i>ThreadInvariants</i></p> <p><i>thread?</i> : <i>THREAD</i></p> <hr/> <p><i>thread?</i> \in dom <i>priority_before_depression</i></p> <p>$\underline{thread_priority}' = \underline{thread_priority}$ $\oplus \{ \text{thread?} \mapsto \text{priority_before_depression}(\text{thread?}) \}$</p> <p>$\underline{thread_max_priority}' = \underline{thread_max_priority}$</p> <p>$\underline{depressed_threads}' = \underline{depressed_threads} \setminus \{ \text{thread?} \}$</p> <p>$\underline{priority_before_depression}' = \underline{priority_before_depression}$</p>
--

9.4.6 Complete Request

The following schemas define the general form of a **thread_depress_abort** request.

<p><i>Processing ThreadDepressAbort</i></p> <p><i>Process Thread Via ThreadPortRequestGood</i></p> <hr/> <p><i>operation?</i> = <i>Thread_depress_abort_id</i></p>
--

A request makes the state changes described in the previous section.

$$\text{ThreadDepressAbortGood} \hat{=} (\text{RVThreadDepressAbortGood} \wedge \text{ThreadDepressAbortState}) \\ \gg \text{RequestReturnOnlyStatus}$$

Review Note:

This definition is included only for consistency with other request specifications.

Execution of the request consists of a good execution.

$$\text{Execute ThreadDepressAbort} \hat{=} \text{ThreadDepressAbortGood}$$

The full specification for kernel processing of a validated **thread_depress_abort** request consists of processing the request followed by its execution.

$$\text{ThreadDepressAbort} \hat{=} \text{Processing ThreadDepressAbort} \ ; \ \text{Execute ThreadDepressAbort}$$

9.5 thread_disable_pc_sampling

The request **thread_disable_pc_sampling** turns off all sampling for a thread.

9.5.1 Client Interface

```
kern_return_t thread_disable_pc_sampling
                (mach_port_t
                 int
                 thread_name,
                 *sample_cnt);
```

Review Note:

The DTOS KID incorrectly includes *flavor* as an input parameter of **thread_disable_pc_sampling**. This parameter is not present in the prototype. The request disables sampling of all types, not for just a particular type, and therefore there is no need for a *flavor* parameter.

9.5.1.1 Input Parameters The following input parameters are provided by the client of a **thread_disable_pc_sampling** request:

- *thread_name?* — the client's name for the thread for which sampling will be turned off

```
ThreadDisablePCSamplingClientInputs _____
thread_name? : NAME
```

A **thread_disable_pc_sampling** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_disable_pc_sampling_id* and has no body.

```
Invoke ThreadDisablePCSampling _____
Invoke MachMsg
ThreadDisablePCSamplingClientInputs
name? = thread_name?
operation? = Thread_disable_pc_sampling_id
```

9.5.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_disable_pc_sampling** request:

- *return!* — the status of the request
- *sample_cnt!* — the number of sample elements in the kernel for the thread

Editorial Note:

In the prototype this parameter is present, but unused. Thus, its output value will be whatever the input value is. The parameter should probably not be present at all since with the current semantics of thread sampling all samples are discarded when sampling is disabled for a thread. To reflect this we will define this value to be zero.

<p><i>ThreadDisablePCSamplingClientOutputs</i></p> <p><i>return!</i> : <i>KERNEL_RETURN</i></p> <p><i>sample_cnt!</i> : \mathbb{N}</p>

<p><i>ThreadDisablePCSamplingReceiveReply</i></p> <p><i>InvokeMachMsgRcv</i></p> <p><i>ThreadDisablePCSamplingClientOutputs</i></p>
<p>$(sample_cnt!, return!) = Text_to_count_and_status(msg_body)$</p>

9.5.2 Kernel Interface

9.5.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_disable_pc_sampling** request:

- *thread?* — the thread for which sampling will be turned off

<p><i>ThreadDisablePCSamplingInputs</i></p> <p><i>thread?</i> : <i>THREAD</i></p>

9.5.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_disable_pc_sampling** request:

- *return!* — the status of the request
- *sample_cnt!* — the number of sample elements in the kernel for the thread

<p><i>ThreadDisablePCSamplingOutputs</i></p> <p><i>return!</i> : <i>KERNEL_RETURN</i></p> <p><i>sample_cnt!</i> : \mathbb{N}</p>

Upon completion of the processing of a **thread_disable_pc_sampling** request a reply message is built from the output parameters.

<p><i>ThreadDisablePCSamplingReply</i></p> <p><i>RequestReturn</i></p> <p><i>sample_cnt?</i> : \mathbb{N}</p>
<p>$reply? = Return_sample_cnt(sample_cnt?)$</p>

9.5.3 Request Criteria

No criteria are defined for the **thread_disable_pc_sampling** request.

<i>sample_cnt!</i>	<i>return!</i>
0	<i>Kern_success</i>

Table 21: Return Values for **thread_disable_pc_sampling**

9.5.4 Return Values

Table 21 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>RVThreadDisablePCSamplingGood</i>
<i>ThreadDisablePCSamplingOutputs</i>
<i>sample_cnt! = 0</i>
<i>return! = Kern_success</i>

9.5.5 State Changes

A successful **thread_disable_pc_sampling** request removes the thread from the set of sampled threads and from the domains of the functions describing sampling. All samples are discarded.

<i>ThreadDisablePCSamplingState</i>
Δ <i>Threads</i>
Δ <i>ThreadSampling</i>
Ξ <i>ThreadPri</i>
Ξ <i>TasksAndThreads</i>
Ξ <i>ThreadSchedPolicy</i>
Ξ <i>ThreadInstruction</i>
Ξ <i>ThreadMachineState</i>
Ξ <i>ThreadExecStatus</i>
Ξ <i>Events</i>
Ξ <i>ThreadExist</i>
Ξ <i>ThreadAndProcessorSet</i>
Ξ <i>PortExist</i>
Ξ <i>PortNameSpace</i>
Ξ <i>SpecialPurposePorts</i>
<i>ThreadInvariants</i>
<i>thread? : THREAD</i>
$\underline{s}ampled_threads' = \underline{s}ampled_threads \setminus \{thread?\}$
$\underline{t}hread_sample_types' = \{thread?\} \triangleleft \underline{t}hread_sample_types$
$\underline{t}hread_sample_sequence_number' = \{thread?\} \triangleleft \underline{t}hread_sample_sequence_number$
$\underline{t}hread_samples' = \{thread?\} \triangleleft \underline{t}hread_samples$

9.5.6 Complete Request

The following schemas define the general form of a **thread_disable_pc_sampling** request.

$\text{Processing ThreadDisablePCSampling} \text{ -----}$ $\text{Process Thread Via ThreadPortRequestGood}$ $\text{operation?} = \text{Thread_disable_pc_sampling_id}$
--

A request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} & \text{ThreadDisablePCSamplingGood} \\ & \hat{=} (\text{RVThreadDisablePCSamplingGood} \wedge \text{ThreadDisablePCSamplingState}) \\ & \gg \text{ThreadDisablePCSamplingReply} \end{aligned}$$

Execution of the request consists of a good execution.

$$\text{Execute ThreadDisablePCSampling} \hat{=} \text{ThreadDisablePCSampling Good}$$

The full specification for kernel processing of a validated **thread_disable_pc_sampling** request consists of processing the request followed by its execution.

$$\begin{aligned} & \text{ThreadDisablePCSampling} \hat{=} \text{Processing ThreadDisablePCSampling} \\ & \quad ; \text{Execute ThreadDisablePCSampling} \end{aligned}$$

9.6 thread_enable_pc_sampling

The request **thread_enable_pc_sampling** turns on a given type of sampling for a thread.

9.6.1 Client Interface

kern_return_t	thread_enable_pc_sampling	
	(mach_port_t	<i>thread_name,</i>
	int	<i>*ticks,</i>
	sampled_pc_flavor_t	<i>flavor);</i>

9.6.1.1 Input Parameters The following input parameters are provided by the client of a **thread_enable_pc_sampling** request:

- *thread_name?* — the client's name for the thread for which sampling will be turned on
- *flavor?* — the type of samples to collect

$\text{ThreadEnablePCSamplingClientInputs} \text{ -----}$ $\text{thread_name?} : \text{NAME}$ $\text{flavor?} : \mathbb{P} \text{ SAMPLE_TYPES}$
--

A **thread_enable_pc_sampling** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_enable_pc_sampling_id* and has a body consisting of *flavor?*.

<i>Invoke ThreadEnablePCSampling</i> <i>Invoke MachMsg</i> <i>ThreadEnablePCSamplingClientInputs</i>
<i>name?</i> = <i>thread_name?</i> <i>operation?</i> = <i>Thread_enable_pc_sampling_id</i> <i>msg_body</i> = <i>Sample_type_set_to_text(flavor?)</i>

9.6.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_enable_pc_sampling** request:

- *return!* — the status of the request
- *ticks!* — the clock granularity (ticks per second) according to the kernel

<i>ThreadEnablePCSamplingClientOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i> <i>ticks!</i> : \mathbb{N}_1

<i>ThreadEnablePCSamplingReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadEnablePCSamplingClientOutputs</i>
<i>(ticks!, return!) = Text_to_ticks_and_status(msg_body)</i>

9.6.2 Kernel Interface

9.6.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_enable_pc_sampling** request:

- *thread?* — the thread for which sampling will be turned on
- *flavor?* — the type of samples to collect

<i>ThreadEnablePCSamplingInputs</i>
<i>thread?</i> : <i>THREAD</i> <i>flavor?</i> : \mathbb{P} <i>SAMPLE_TYPES</i>

9.6.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_enable_pc_sampling** request:

- *return!* — the status of the request
- *ticks!* — the clock granularity (ticks per second) according to the kernel

<i>ThreadEnablePCSamplingOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i> <i>ticks!</i> : \mathbb{N}_1

Upon completion of the processing of a **thread_enable_pc_sampling** request a reply message is built from the output parameters.

<i>ThreadEnablePCSamplingReply</i> <i>RequestReturn</i> <i>ticks?</i> : \mathbb{N}_1
<i>reply?</i> = <i>Return_ticks(ticks?)</i>

9.6.3 Request Criteria

The following criteria are defined for the **thread_enable_pc_sampling** request.

- **C1** — There are sufficient resources to create a sampling buffer. We do not actually model the consumption of resources by the kernel. So, we will use the set *Resources_available_to_create_sampling_buffer* to indicate the set of states where there are sufficient resources to create a sampling buffer.

| *Resources_available_to_create_sampling_buffer* : \mathbb{P} *DtosExec*

<i>C1ThreadEnablePCSamplingResourcesAvailable</i> <i>DtosExec</i>
θ <i>DtosExec</i> \in <i>Resources_available_to_create_sampling_buffer</i>

NotC1ThreadEnablePCSamplingResourcesAvailable $\hat{=}$
DtosExec \wedge \neg *C1ThreadEnablePCSamplingResourcesAvailable*

Note that no criterion is defined to check that *flavor?* is a set of recognized sample types. If an unrecognized type is included in *flavor?*, no error will occur. Unrecognized sample types will simply be ignored and produce no samples.

9.6.4 Return Values

Table 22 describes the values returned at the completion of the request and the conditions under which each value is returned. The design does not specify the value of *ticks!* when an error occurs. It depends on the implementation, and we leave it unspecified.

| *Ticks_per_second* : \mathbb{N}_1

Editorial Note:

Even though C1 examines resource availability, the kernel returns *Kern_invalid_argument* when C1 is false. In addition the following message is printed to standard output: "thread_enable_pc_sampling: kalloc failed".

<i>ticks!</i>	<i>return!</i>	C1
<i>Ticks_per_second</i>	<i>Kern_success</i>	T
—	<i>Kern_invalid_argument</i>	F

Table 22: Return Values for **thread_enable_pc_sampling**

<i>RVThreadEnablePCSamplingGood</i>
<i>C1ThreadEnablePCSamplingResourcesAvailable</i>
<i>ThreadEnablePCSamplingOutputs</i>
<i>ticks!</i> = <i>Ticks_per_second</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadEnablePCSamplingResourceShortage</i>
<i>NotC1ThreadEnablePCSamplingResourcesAvailable</i>
<i>ThreadEnablePCSamplingOutputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

9.6.5 State Changes

A successful **thread_enable_pc_sampling** request adds the thread to the set of sampled threads and records the type of samples to be collected. It also sets the sample sequence number for the thread to zero. If the thread was already being sampled, the flavor is reset, but the sequence number is unchanged. In this case, any samples currently in the buffer remain there.

<p><i>ThreadEnablePCSamplingState</i></p> <p>Δ <i>Threads</i></p> <p>Δ <i>ThreadSampling</i></p> <p>Ξ <i>ThreadPri</i></p> <p>Ξ <i>TasksAndThreads</i></p> <p>Ξ <i>ThreadSchedPolicy</i></p> <p>Ξ <i>ThreadInstruction</i></p> <p>Ξ <i>ThreadMachine.State</i></p> <p>Ξ <i>ThreadExecStatus</i></p> <p>Ξ <i>Events</i></p> <p>Ξ <i>ThreadExist</i></p> <p>Ξ <i>ThreadAndProcessorSet</i></p> <p>Ξ <i>PortExist</i></p> <p>Ξ <i>PortName.Space</i></p> <p>Ξ <i>SpecialPurposePorts</i></p> <p><i>ThreadInvariants</i></p> <p><i>thread?</i> : <i>THREAD</i></p> <p><i>flavor?</i> : \mathbb{P} <i>SAMPLE_TYPES</i></p> <hr/> <p>$\underline{\text{sampled_threads}}' = \underline{\text{sampled_threads}} \cup \{\text{thread?}\}$</p> <p>$\underline{\text{thread_sample_types}}' = \underline{\text{thread_sample_types}} \oplus \{\text{thread?} \mapsto \text{flavor?}\}$</p> <p>$\underline{\text{thread_sample_sequence_number}}' = \{\text{thread?} \mapsto 0\}$ $\oplus \underline{\text{thread_sample_sequence_number}}$</p> <p>$\underline{\text{thread_samples}}' = \{\text{thread?} \mapsto \langle \rangle\} \oplus \underline{\text{thread_samples}}$</p>

9.6.6 Complete Request

The following schemas define the general form of a **thread_enable_pc_sampling** request.

<p><i>Processing ThreadEnablePCSampling</i></p> <p><i>Process Thread Via ThreadPortRequestGood</i></p> <hr/> <p><i>operation?</i> = <i>Thread_enable_pc_sampling_id</i></p>

A request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} & \textit{ThreadEnablePCSampling Good} \\ & \hat{=} (\textit{RVThreadEnablePCSamplingGood} \wedge \textit{ThreadEnablePCSamplingState}) \\ & \gg \textit{ThreadEnablePCSampling Reply} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} & \textit{ThreadEnablePCSampling Bad} \\ & \hat{=} \textit{RVThreadEnablePCSamplingResourceShortage} \gg \textit{RequestNoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} & \textit{Execute ThreadEnablePCSampling} \hat{=} \textit{ThreadEnablePCSampling Good} \\ & \wedge \textit{ThreadEnablePCSampling Bad} \end{aligned}$$

The full specification for kernel processing of a validated **thread_enable_pc_sampling** request consists of processing the request followed by its execution.

$$\text{ThreadEnablePCSampling} \hat{=} \text{ProcessingThreadEnablePCSampling} \\ \text{;} \text{ExecuteThreadEnablePCSampling}$$

9.7 thread_get_assignment

The request **thread_get_assignment** returns a send right to the name port of the processor set to which a thread is assigned. This port can only be used to obtain information about the processor set.

9.7.1 Client Interface

```
kern_return_t thread_get_assignment
    (mach_port_t
     mach_port_t*
     thread_name,
     processor_set_name);
```

9.7.1.1 Input Parameters The following input parameters are provided by the client of a **thread_get_assignment** request:

- *thread_name?* — the client's name for the thread whose processor set name port is requested

```
ThreadGetAssignmentClientInputs _____
thread_name? : NAME
```

A **thread_get_assignment** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_get_assignment_id* and has no body.

```
Invoke ThreadGetAssignment _____
Invoke MachMsg
ThreadGetAssignmentClientInputs
name? = thread_name?
operation? = Thread_get_assignment_id
```

9.7.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_get_assignment** request:

- *processor_set_name!* — a send right to the name port of the desired processor set
- *return!* — the status of the request

```
ThreadGetAssignmentClientOutputs _____
processor_set_name! : NAME
return! : KERNEL_RETURN
```

<i>ThreadGetAssignmentReceiveReply</i> <i>InvokeMachMsgRcv</i> <i>ThreadGetAssignmentClientOutputs</i> <i>(processor_set_name!, return!) = Text_to_name_and_status(msg_body)</i>

9.7.2 Kernel Interface

9.7.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_get_assignment** request:

- *thread?* — the thread whose processor set name port is requested

<i>ThreadGetAssignmentInputs</i> <i>thread? : THREAD</i>

9.7.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_get_assignment** request:

- *processor_set!* — the desired processor set
- *return!* — the status of the request

<i>ThreadGetAssignmentOutputs</i> <i>processor_set! : PROCESSOR_SET</i> <i>return! : KERNEL_RETURN</i>
--

Upon completion of the processing of a **thread_get_assignment** request a reply message is built from the output parameters.

<i>ThreadGetAssignmentReply</i> <i>RequestOnlyObserves</i> <i>processor_set? : PROCESSOR_SET</i> let <i>port</i> == <i>ps_name_port_rel(processor_set?)</i> <ul style="list-style-type: none"> ● <i>reply?</i> = <i>Return_capability(Thread_port_to_s_right(port))</i>

9.7.3 Request Criteria

No criteria are defined for the **thread_get_assignment** request.

9.7.4 Return Values

Table 23 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>processor_set!</i>	<i>return!</i>
<i>thread_assigned_to(thread?)</i>	<i>Kern_success</i>

Table 23: Return Values for **thread_get_assignment**

<i>RVThreadGetAssignmentGood</i> <i>ThreadAndProcessorSet</i> <i>ProcessorsAndPorts</i> <i>ThreadGetAssignmentOutputs</i> <i>thread? : THREAD</i>
<i>thread? ∈ dom thread_assigned_to</i>
<i>return! = Kern_success</i> <i>processor_set! = thread_assigned_to(thread?)</i>

9.7.5 State Changes

A **thread_get_assignment** request does not make any state changes since it only observes the system state.

9.7.6 Complete Request

The **thread_get_assignment** request has the following general form.

<i>Processing ThreadGetAssignment</i> <i>Process Thread Via ThreadPortRequestGood</i>
<i>operation? = Thread_get_assignment_id</i>

The full specification for kernel processing of a validated **thread_get_assignment** request consists of processing the request, execution of the request, and the creation of a kernel reply.

$$\text{ThreadGetAssignment} \hat{=} \text{Processing ThreadGetAssignment} \\ \quad ; (\text{RVThreadGetAssignmentGood} \gg \text{ThreadGetAssignmentReply})$$

9.8 thread_get_sampled_pcs

The request **thread_get_sampled_pcs** returns the samples collected for a given thread.

9.8.1 Client Interface

```
kern_return_t thread_get_sampled_pcs
    (mach_port_t
     unsigned
     sampled_pc_t
     int
     thread_name,
     *seqno,
     sampled_pcs[ ],
     *sample_cnt);
```

9.8.1.1 Input Parameters The following input parameters are provided by the client of a **thread_get_sampled_pcs** request:

- *thread_name?* — the client's name for the thread whose samples will be returned
- *seqno?* — the sequence number of the first sample that should be returned. If this sample is no longer available due to insufficient space in the sampling buffer, the earliest available sample will be the starting point.

Review Note:

The DTOS KID reports *seqno* as an output parameter only. In the prototype it is used for both input and output.

ThreadGetSampledPCs ClientInputs

thread_name? : NAME
seqno? : N

A **thread_get_sampled_pcs** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_get_sampled_pcs_id* and has a body consisting of *seqno?*.

Invoke ThreadGetSampledPCs

Invoke MachMsg
ThreadGetSampledPCs ClientInputs

name? = *thread_name?*
operation? = *Thread_get_sampled_pcs_id*
msg_body = *Sequence_number_to_text(seqno?)*

9.8.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_get_sampled_pcs** request:

- *return!* — the status of the request
- *seqno!* — the sequence number of the most recently collected sample
- *sampled_pcs!* — the samples returned
- *sample_cnt!* — the number of samples returned

ThreadGetSampledPCs ClientOutputs

return! : KERNEL_RETURN
seqno! : N
sampled_pcs! : seq SAMPLE
sample_cnt! : Z

ThreadGetSampledPCs ReceiveReply

Invoke MachMsgRcv
ThreadGetSampledPCs ClientOutputs

(*seqno!*, *sampled_pcs!*, *sample_cnt!*, *return!*)
= *Text_to_seqno_and_PCs_and_count_and_status(msg_body)*

9.8.2 Kernel Interface

9.8.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_get_sampled_pcs** request:

- *thread?* — the thread whose samples will be returned
- *seqno?* — the sequence number of the first sample that should be returned. If this sample is no longer available due to insufficient space in the sampling buffer, the earliest available sample will be the starting point.

Review Note:

The DTOS KID reports *seqno* as an output parameter only. In the prototype it is used for both input and output.

ThreadGetSampledPCsInputs _____*thread?* : *THREAD**seqno?* : \mathbb{N}

9.8.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_get_sampled_pcs** request:

- *return!* — the status of the request
- *seqno!* — the sequence number of the most recently collected sample
- *sampled_pcs!* — the samples returned
- *sample_cnt!* — the number of samples returned

ThreadGetSampledPCsOutputs _____*return!* : *KERNEL_RETURN**seqno!* : \mathbb{N} *sampled_pcs!* : seq *SAMPLE**sample_cnt!* : \mathbb{Z}

Upon completion of the processing of a **thread_get_sampled_pcs** request a reply message is built from the output parameters.

ThreadGetSampledPCsReply _____*RequestOnlyObserves**seqno?* : \mathbb{N} *sampled_pcs?* : seq *SAMPLE**sample_cnt?* : \mathbb{Z}

reply? = *Return_samples(seqno?, sampled_pcs?, sample_cnt?)*

9.8.3 Request Criteria

The following criteria are defined for the **thread_get_sampled_pcs** request.

- **C1** — Sampling is currently enabled for the thread.

<i>C1ThreadSamplingEnabled</i>
<i>ThreadSampling</i>
<i>thread? : THREAD</i>
<i>thread? ∈ <u>s</u>ampled_threads</i>

$$\begin{aligned} & \text{NotC1ThreadSamplingEnabled} \\ & \hat{=} \text{ThreadSampling} \wedge \neg \text{C1ThreadSamplingEnabled} \end{aligned}$$

9.8.4 Return Values

Tables 24–27 describe the values returned at the completion of the request and the conditions under which each value is returned. The specification does not state what should be returned in *seqno!*, *sampled_pcs!* and *sample_cnt!* when the thread is not being sampled. So, these values depend on the implementation and we leave them unspecified.

<i>seqno!</i>	C1
<i>thread_sample_sequence_number(thread?)</i>	T
—	F

Table 24: Return Values for **thread_get_sampled_pcs**

<i>sampled_pcs!</i>	C1
<i>Samples_returned(thread_samples(thread?), seqno?, thread_sample_sequence_number(thread?))</i>	T
—	F

Table 25: Return Values for **thread_get_sampled_pcs**

<i>sample_cnt!</i>	C1
<i># sampled_pcs!</i>	T
—	F

Table 26: Return Values for **thread_get_sampled_pcs**

The function *Samples_returned_count(from, to)* returns the number of samples in the range *from* to *to* or the buffer size *Max_samples*, whichever is smaller. A negative return value is interpreted as zero samples. The function *Samples_returned(sample_sequence, from, to)* returns a number of samples as indicated by *Samples_returned_count(from, to)* from *sample_sequence* ending with the sample with sequence number *to*.

<i>return!</i>	C1
<i>Kern_success</i>	T
<i>Kern_failure</i>	F

Table 27: Return Values for **thread_get_sampled_pcs**

$\begin{aligned} & \text{Samples_returned} : (\text{seq } SAMPLE) \times \mathbb{N} \times \mathbb{N} \leftrightarrow \text{seq } SAMPLE \\ & \text{Samples_returned_count} : (\mathbb{N} \times \mathbb{N}) \leftrightarrow \mathbb{Z} \end{aligned}$
$\begin{aligned} & \forall \text{sample_sequence} : \text{seq } SAMPLE; \text{from}, \text{to} : \mathbb{N} \\ & \bullet \text{Samples_returned_count}(\text{from}, \text{to}) = \min \{ \text{Max_samples}, \text{to} - \text{from} + 1 \} \\ & \wedge \text{Samples_returned}(\text{sample_sequence}, \text{from}, \text{to}) \\ & \quad = \{ j : \mathbb{N} \mid \text{to} - \text{Samples_returned_count}(\text{from}, \text{to}) < j \leq \text{to} \} \upharpoonright \text{sample_sequence} \end{aligned}$
$\begin{aligned} & \text{RVThreadGetSampledPCsGood} \\ & \text{C1ThreadSamplingEnabled} \\ & \text{ThreadGetSampledPCsOutputs} \\ & \text{ThreadGetSampledPCsInputs} \end{aligned}$
$\begin{aligned} & \text{seqno!} = \text{thread_sample_sequence_number}(\text{thread?}) \\ & \text{sampled_pcs!} = \text{Samples_returned}(\text{thread_samples}(\text{thread?}), \\ & \quad \text{seqno?}, \text{thread_sample_sequence_number}(\text{thread?})) \\ & \text{sample_cnt!} = \# \text{sampled_pcs!} \\ & \text{return!} = \text{Kern_success} \end{aligned}$
$\begin{aligned} & \text{RVThreadGetSampledPCsBad} \\ & \text{NotC1ThreadSamplingEnabled} \\ & \text{ThreadGetSampledPCsOutputs} \\ & \text{ThreadGetSampledPCsInputs} \end{aligned}$
$\text{return!} = \text{Kern_failure}$

9.8.5 State Changes

A **thread_get_sampled_pcs** request does not make any state changes since it only observes the system state.

9.8.6 Complete Request

The following schemas define the general form of a **thread_get_sampled_pcs** request.

$\begin{aligned} & \text{Processing ThreadGetSampledPCs} \\ & \text{Process Thread Via ThreadPortRequestGood} \end{aligned}$
$\text{operation?} = \text{Thread_get_sampled_pcs_id}$

A successful request creates a kernel reply.

$$\begin{aligned} & \text{ThreadGetSampledPCsGood} \hat{=} \text{RVThreadGetSampledPCsGood} \\ & \gg \text{ThreadGetSampledPCsReply} \end{aligned}$$

An unsuccessful request returns an error status.

$$\text{ThreadGetSampledPCsBad} \hat{=} \text{RVThreadGetSampledPCsBad} \gg \text{RequestNoOp}$$

Execution of the request consists of a good execution or an error execution.

$$\text{Execute ThreadGetSampledPCs} \hat{=} \text{ThreadGetSampledPCsGood} \vee \text{ThreadGetSampledPCsBad}$$

The full specification for kernel processing of a validated **thread_get_sampled_pcs** request consists of processing the request followed by its execution.

$$\begin{aligned} \text{ThreadGetSampledPCs} &\hat{=} \text{Processing ThreadGetSampledPCs} \\ &\quad \ddagger \text{Execute ThreadGetSampledPCs} \end{aligned}$$

9.9 thread_get_special_port

The request **thread_get_special_port** allows a task to obtain a send right to a specified special port for a specified thread.

9.9.1 Client Interface

```
kern_return_t thread_get_special_port
    (mach_port_t thread_name,
     int which_port,
     mach_port_t* special_port_name);
```

thread_get_exception_port

Macro form

```
kern_return_t thread_get_exception_port
    (mach_port_t thread_name,
     mach_port_t* special_port_name);
⇒ thread_get_special_port (thread_name, THREAD_EXCEPTION_PORT,
    special_port_name)
```

thread_get_kernel_port

Macro form

```
kern_return_t thread_get_kernel_port
    (mach_port_t thread_name,
     mach_port_t* special_port_name);
⇒ thread_get_special_port (thread_name, THREAD_KERNEL_PORT,
    special_port_name)
```

9.9.1.1 Input Parameters The following input parameters are provided by the client of a **thread_get_special_port** request:

- *thread_name?* — the client's name for the thread whose special port is to be returned
- *which_port?* — the type of special port that is to be returned

```
ThreadGetSpecialPortClientInputs _____
thread_name? : NAME
which_port? : THREAD_SPECIAL_PORTS
```

A **thread_get_special_port** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_get_special_port_id* and has a body consisting of *which_port?*.

```
Invoke ThreadGetSpecialPort _____
Invoke
ThreadGetSpecialPortClientInputs
trap_id? = Mach_msg_trap
user_spec?.message.header.operation = Thread_get_special_port_id
user_spec?.message.header.remote_port = thread_name?
user_spec?.message.body
    = Thread_special_ports_to_text(which_port?)
```

9.9.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_get_special_port** request:

- *special_port_name!* — the name of a send right (capability) for the requested special port
- *return!* — the status of the request

```
ThreadGetSpecialPortClientOutputs _____
special_port_name! : NAME
return! : KERNEL_RETURN
```

```
ThreadGetSpecialPortReceiveReply _____
Invoke MachMsgRcv
ThreadGetSpecialPortClientOutputs
(special_port_name!, return!) = Text_to_name_and_status(msg_body)
```

9.9.2 Kernel Interface

9.9.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_get_special_port** request:

- *thread?* — the thread whose special port is to be returned

- *which_port?* — the type of special port that is to be returned

<p><i>ThreadGetSpecialPortInputs</i></p> <p><i>thread?</i> : <i>THREAD</i></p> <p><i>which_port?</i> : <i>THREAD_SPECIAL_PORTS</i></p>
--

9.9.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_get_special_port** request:

- *special_port!* — the requested special port
- *return!* — the status of the request

<p><i>ThreadGetSpecialPortOutputs</i></p> <p><i>special_port!</i> : <i>PORT</i></p> <p><i>return!</i> : <i>KERNEL_RETURN</i></p>
--

Upon completion of the processing of a **thread_get_special_port** request a reply message is built from the output parameters. The reply message will contain a send right to the requested special port.

<p><i>ThreadGetSpecialPortReply</i></p> <p><i>OnlyObserves</i></p> <p><i>special_port?</i> : <i>PORT</i></p> <p><i>reply?</i> = <i>Return_capability(Thread_port_to_s_right(special_port?))</i></p>

9.9.3 Request Criteria

The following criteria are defined for the **thread_get_special_port** request.

- **C1** — An exception port request is made.

<p><i>C1ThreadGetExceptionPort</i></p> <p><i>which_port?</i> : <i>THREAD_SPECIAL_PORTS</i></p> <p><i>which_port?</i> = <i>Thread_exception_port</i></p>

$NotC1ThreadGetExceptionPort \hat{=} \neg C1ThreadGetExceptionPort$

- **C2** — A kernel port request is made.

<p><i>C2ThreadGetKernelPort</i></p> <p><i>which_port?</i> : <i>THREAD_SPECIAL_PORTS</i></p> <p><i>which_port?</i> = <i>Thread_kernel_port</i></p>

$NotC2ThreadGetKernelPort \hat{=} \neg C2ThreadGetKernelPort$

■ **C3** — The client has *Get_thread_exception_port* permission to the target thread.

Review Note:

In C3 and C4, we've begun the process of dealing with deferred permission checks under the new execution model. The first schema is used to initiate the permission checking routine and the criteria schemas will be used after the permission has been retrieved.

ThreadGetSpecialPortPermCheckGTEP

Transition

$\exists request : Request; CheckPending; ThreadGetSpecialPortInputs$

- $curr_bk?? = Bk_have_request(request)$
 - $\wedge request.operation = Thread_get_special_port_id$
 - $\wedge thread_self(thread?) = request.service_port$
 - $\wedge ssi = \underline{task_sid}(curr_task??)$
 - $\wedge osi = thread_target(curr_task??, thread?)$
 - $\wedge breaks' = breaks$
- $\oplus\{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Get_thread_exception_port, env)\}$

C3ThreadCanGetExceptionPort

Transition

env : ENVIRONMENT

$curr_bk?? = Bk_have_ruling(Get_thread_exception_port, True, env)$

NotC3ThreadCanGetExceptionPort

Transition

env : ENVIRONMENT

$curr_bk?? = Bk_have_ruling(Get_thread_exception_port, False, env)$

■ **C4** — The client has *Get_thread_kernel_port* permission to the target thread.

ThreadGetSpecialPortPermCheckGTKP

Transition

$\exists request : Request; CheckPending; ThreadGetSpecialPortInputs$

- $curr_bk?? = Bk_have_request(request)$
 - $\wedge request.operation = Thread_get_special_port_id$
 - $\wedge thread_self(thread?) = request.service_port$
 - $\wedge ssi = \underline{task_sid}(curr_task??)$
 - $\wedge osi = thread_target(curr_task??, thread?)$
 - $\wedge breaks' = breaks$
- $\oplus\{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Get_thread_kernel_port, env)\}$

C4ThreadCanGetKernelPort

Transition

env : ENVIRONMENT

$curr_bk?? = Bk_have_ruling(Get_thread_kernel_port, True, env)$

$NotC4ThreadCanGetKernelPort$ <i>Transition</i> $env : ENVIRONMENT$ $curr_bk?? = Bk_have_ruling(Get_thread_kernel_port, False, env)$

- **C5** — The exception port of the thread is defined.

$C5ThreadExceptionPortDefined$ <i>SpecialThreadPorts</i> $thread? : THREAD$ $thread? \in \text{dom } thread_eport$
--

$NotC5ThreadExceptionPortDefined$
 $\hat{=} SpecialThreadPorts \wedge \neg C5ThreadExceptionPortDefined$

- **C6** — The kernel port of the thread is defined.

$C6ThreadKernelPortDefined$ <i>SpecialThreadPorts</i> $thread? : THREAD$ $thread? \in \text{dom } thread_self$
--

$NotC6ThreadKernelPortDefined$
 $\hat{=} SpecialThreadPorts \wedge \neg C6ThreadKernelPortDefined$

9.9.4 Return Values

Tables 28 and 29 describe the values returned at the completion of the request and the conditions under which each value is returned. Note that C1 and C2 are mutually exclusive. It is possible that a thread has no exception or kernel (self) port since the port may have been deleted. The design does not specify the value of *special_port!* in this case. We assume that the null port is returned by the kernel routine, and that IPC will convert this into the name *Mach_port_null*. We leave unspecified the value returned in *special_port!* when the client does not have permission to get the requested special port or when the client does not ask for a valid type of special port. Note that C5 and C6 do not affect the return status.

Review Note:

We assume that the prototype will check the conditions in the order {C1, C2 }, {C3 or C4}, {C5 or C6}. However, the prototype is currently not checking C3 and C4.

Review Note:

It might make more sense to permit *Null_port* in the range of *thread_eport* and *thread_self*. (Note that the exception port is actually initialized to *Null_port* by the **thread_create** request in the prototype and that **thread_set_special_port** can set an exception or kernel port to *Null_port*.) This would remove the need for criteria C5 and C6.

<i>return!</i>	C1	C2	C3	C4
<i>Kern_success</i>	T	F	T	-
<i>Kern_success</i>	F	T	-	T
<i>Kern_insufficient_permission</i>	T	F	F	-
<i>Kern_insufficient_permission</i>	F	T	-	F
<i>Kern_invalid_argument</i>	F	F	-	-

Table 28: Return Values for **thread_get_special_port**

<i>special_port!</i>	C1	C2	C3	C4	C5	C6
<i>thread_eport(thread?)</i>	T	F	T	-	T	-
<i>Null_port</i>	T	F	T	-	F	-
<i>thread_sself(thread?)</i>	F	T	-	T	-	T
<i>Null_port</i>	F	T	-	T	-	F
—	otherwise					

Table 29: Return Values for **thread_get_special_port**

<i>RVThreadGetExceptionPortGood</i>
<i>C1 ThreadGetExceptionPort</i>
<i>NotC2 ThreadGetKernelPort</i>
<i>C3 ThreadCanGetExceptionPort</i>
<i>C5 ThreadExceptionPortDefined</i>
<i>ThreadGetSpecialPort Outputs</i>
<i>thread? : THREAD</i>
<i>return! = Kern_success</i>
<i>special_port! = thread_eport(thread?)</i>

<i>RVThreadGetExceptionPortNull</i>
<i>C1 ThreadGetExceptionPort</i>
<i>NotC2 ThreadGetKernelPort</i>
<i>C3 ThreadCanGetExceptionPort</i>
<i>NotC5 ThreadExceptionPortDefined</i>
<i>ThreadGetSpecialPort Outputs</i>
<i>return! = Kern_success</i>
<i>special_port! = Null_port</i>

<i>RVThreadGetKernelPortGood</i>
<i>NotC1 ThreadGetExceptionPort</i>
<i>C2 ThreadGetKernelPort</i>
<i>C4 ThreadCanGetKernelPort</i>
<i>C6 ThreadKernelPortDefined</i>
<i>ThreadGetSpecialPort Outputs</i>
<i>thread? : THREAD</i>
<i>return! = Kern_success</i>
<i>special_port! = thread_sself(thread?)</i>

<i>RVThreadGetKernelPortNull</i> <i>NotC1ThreadGetExceptionPort</i> <i>C2ThreadGetKernelPort</i> <i>C4ThreadCanGetKernelPort</i> <i>NotC6ThreadKernelPortDefined</i> <i>ThreadGetSpecialPortOutputs</i>
<i>return!</i> = <i>Kern_success</i> <i>special_port!</i> = <i>Null_port</i>

<i>RVThreadCannotGetExceptionPort</i> <i>C1ThreadGetExceptionPort</i> <i>NotC2ThreadGetKernelPort</i> <i>NotC3ThreadCanGetExceptionPort</i> <i>ThreadGetSpecialPortOutputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

<i>RVThreadCannotGetKernelPort</i> <i>NotC1ThreadGetExceptionPort</i> <i>C2ThreadGetKernelPort</i> <i>NotC4ThreadCanGetKernelPort</i> <i>ThreadGetSpecialPortOutputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

<i>RVThreadGetSpecialPortInvalidArgument</i> <i>NotC1ThreadGetExceptionPort</i> <i>NotC2ThreadGetKernelPort</i> <i>ThreadGetSpecialPortOutputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

9.9.5 State Changes

A **thread_get_special_port** request does not make any state changes since it only observes the system state.

9.9.6 Complete Request

The **thread_get_special_port** request has the following general form.

<i>Processing ThreadGetSpecialPort</i> <i>ProcessThreadViaThreadPortRequestGood</i>
<i>operation?</i> = <i>Thread_get_special_port_id</i>

A successful **thread_get_special_port** request causes the creation of a kernel reply.

$$\begin{aligned} & \textit{ThreadGetSpecialPortGood} \\ & \hat{=} (\textit{RVThreadGetExceptionPortGood} \vee \textit{RVThreadGetExceptionPortNull} \\ & \quad \vee \textit{RVThreadGetKernelPortGood} \vee \textit{RVThreadGetKernelPortNull}) \\ & \gg \textit{ThreadGetSpecialPortReply} \end{aligned}$$

$$\begin{aligned} & \textit{ThreadGetSpecialPortBad} \\ & \hat{=} (\textit{RVThreadCannotGetExceptionPort} \vee \textit{RVThreadCannotGetKernelPort} \\ & \quad \vee \textit{RVThreadGetSpecialPortInvalidArgument}) \\ & \gg \textit{NoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\textit{ExecuteThreadGetSpecialPort} \hat{=} (\textit{ThreadGetSpecialPortGood} \vee \textit{ThreadGetSpecialPortBad})$$

The full specification for kernel processing of a validated **thread_get_special_port** request consists of processing the request followed by its execution.

$$\textit{ThreadGetSpecialPort} \hat{=} \textit{ProcessingThreadGetSpecialPort} ; \textit{ExecuteThreadGetSpecialPort}$$

9.10 thread_get_state

The request **thread_get_state** returns an array containing state information about a specified thread (other than the client thread).

9.10.1 Client Interface

kern_return_t	thread_get_state	
	(mach_port_t	<i>target_thread_name,</i>
	int	<i>flavor;</i>
	thread_state_t	<i>old_state,</i>
	mach_msg_type_number_t*	<i>old_state_cnt</i>);

9.10.1.1 Input Parameters The following input parameters are provided by the client of a **thread_get_state** request:

- *target_thread_name?* — the client's name for the thread whose state information is to be returned
- *flavor?* — the type of state information that is to be returned
- *old_state_cnt?* — the maximum size of the array to be returned

$\begin{aligned} & \textit{ThreadGetStateClientInputs} \\ & \textit{target_thread_name?} : \textit{NAME} \\ & \textit{flavor?} : \textit{THREAD_STATE_INFO_TYPES} \\ & \textit{old_state_cnt?} : \mathbb{N} \end{aligned}$

A **thread_get_state** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_get_state_id* and has a body consisting of *flavor?* and *old_state_cnt?*.

<i>Invoke ThreadGetState</i> <i>Invoke MachMsg</i> <i>ThreadGetState ClientInputs</i>
<i>name?</i> = <i>target_thread_name?</i> <i>operation?</i> = <i>Thread_get_state_id</i> <i>msg_body</i> = <i>Thread_state_info_type_and_number_to_text(flavor?, old_state_cnt?)</i>

9.10.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_get_state** request:

- *old_state!* — state information of the specified type about the target thread
- *old_state_cnt!* — the size of the full array of available state information of the type specified
- *return!* — the status of the request

<i>ThreadGetState ClientOutputs</i>
<i>old_state!</i> : <i>THREAD_STATE_INFO</i> <i>old_state_cnt!</i> : \mathbb{N} <i>return!</i> : <i>KERNEL_RETURN</i>

<i>ThreadGetState ReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadGetState ClientOutputs</i>
(<i>old_state!</i> , <i>old_state_cnt!</i> , <i>return!</i>) = <i>Text_to_state_and_count_and_status(msg_body)</i>

9.10.2 Kernel Interface

9.10.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_get_state** request:

- *target_thread?* — the thread whose state information is to be returned
- *flavor?* — the type of state information that is to be returned
- *old_state_cnt?* — the maximum size of the array to be returned

<i>ThreadGetStateInputs</i>
<i>target_thread?</i> : <i>THREAD</i> <i>flavor?</i> : <i>THREAD_STATE_INFO_TYPES</i> <i>old_state_cnt?</i> : \mathbb{N}

9.10.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_get_state** request:

- *old_state!* — state information of the specified type about the target thread
- *old_state_cnt!* — the size of the full array of available state information of the type specified
- *return!* — the status of the request

<i>ThreadGetState Outputs</i> <i>old_state!</i> : <i>THREAD_STATE_INFO</i> <i>old_state_cnt!</i> : \mathbb{N} <i>return!</i> : <i>KERNEL_RETURN</i>
--

Upon completion of the processing of a **thread_get_state** request a reply message is built from the output parameters.

<i>ThreadGetState Reply</i> <i>RequestReturn</i> <i>old_state?</i> : <i>THREAD_STATE_INFO</i> <i>old_state_cnt?</i> : \mathbb{N} <i>reply?</i> = <i>Return_thread_state_info(old_state?, old_state_cnt?)</i>
--

9.10.3 Request Criteria

The following criteria are defined for the **thread_get_state** request.

- **C1** — The parameter *thread?* is not equal to the client thread, *flavor?* is a valid type of state information, and *old_state_cnt?* is large enough for the requested state information. (The function *Thread_state_count* returns the size required for the given type of state information.)

Editorial Note:

Nothing in the design states the reason that the client thread may not get its own state information. We believe that it is merely an implementation difficulty in that in order to get state information the thread must be temporarily stopped. If the client thread stopped itself, it could not collect the state information.

$C1ThreadGetStateGoodArgs$ $ThreadMachineState$ $ThreadsAndProcessors$ $cpu?? : PROCESSOR$ $thread? : THREAD$ $flavor? : THREAD_STATE_INFO_TYPES$ $old_state_cnt? : \mathbb{N}$
$cpu?? \in \text{dom } \underline{active_thread}$ $flavor? \in \text{dom } Thread_state_count$
$thread? \neq \underline{active_thread}(cpu??)$ $(thread?, flavor?) \in \text{dom } \underline{thread_state}$ $old_state_cnt? \geq Thread_state_count(flavor?)$

$$NotC1ThreadGetStateGoodArgs \hat{=} ThreadMachineState \wedge ThreadsAndProcessors \wedge \neg C1ThreadGetStateGoodArgs$$

9.10.4 Return Values

Tables 30–32 describe the values returned at the completion of the request and the conditions under which each value is returned. The design does not specify the values of $old_state!$ and $old_state_cnt!$ when an error occurs. These values therefore depend on the implementation algorithm and we leave them unspecified.

$return$	C1
$Kern_success$	T
$Kern_invalid_argument$	F

Table 30: Return Values for **thread_get_state**

$old_state!$	C1
$\underline{thread_state}(thread?, flavor?)$	T
—	F

Table 31: Return Values for **thread_get_state**

$old_state_cnt!$	C1
$Thread_state_count(flavor?)$	T
—	F

Table 32: Return Values for **thread_get_state**

<i>RVThreadGetStateGood</i> <i>C1ThreadGetStateGoodArgs</i> <i>ThreadMachineState</i> <i>ThreadGetStateOutputs</i>
<i>old_state_cnt!</i> = <i>Thread_state_count</i> (<i>flavor?</i>) <i>old_state!</i> = <i>thread_state</i> (<i>thread?</i> , <i>flavor?</i>) <i>return!</i> = <i>Kern_success</i>
<i>RVThreadGetStateInvalidArgument</i> <i>NotC1ThreadGetStateGoodArgs</i> <i>ThreadGetStateOutputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

9.10.5 State Changes

A successful **thread_get_state** request gets the state of the thread to the supplied state information. The run state of the thread may also change since the request must ensure that the thread is temporarily suspended and then perhaps restart it.

<i>ThreadGetStateState</i> Δ <i>Threads</i> Ξ <i>TasksAndThreads</i> Ξ <i>ThreadPri</i> Ξ <i>ThreadSchedPolicy</i> Ξ <i>ThreadInstruction</i> Δ <i>ThreadExecStatus</i> Ξ <i>ThreadStatistics</i> Ξ <i>ThreadMachineState</i> Ξ <i>Exist</i> Ξ <i>SpecialPurposePorts</i> Ξ <i>ThreadAndProcessorSet</i> <i>ThreadInvariants</i> <i>ThreadGetStateInputs</i> <i>target_thread?</i> : <i>THREAD</i>
<i>ThreadDoWaitThenRelease</i> [<i>target_thread?</i> / <i>stopping_thread</i>] <i>swapped_threads'</i> = <i>swapped_threads</i> <i>idle_threads'</i> = <i>idle_threads</i> <i>thread_suspend_count'</i> = <i>thread_suspend_count</i> <i>threads_wired'</i> = <i>threads_wired</i>

9.10.6 Complete Request

The following schemas define the general form of a **thread_get_state** request.

<i>ProcessingThreadGetState</i> <i>ProcessThreadViaThreadPortRequestGood</i>
<i>operation?</i> = <i>Thread_get_state_id</i>

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} \text{ThreadGetStateGood} &\hat{=} (\text{RVThreadGetStateGood} \wedge \text{ThreadGetStateState}) \\ &\gg \text{ThreadGetStateReply} \end{aligned}$$

An unsuccessful request returns an error status.

$$\text{ThreadGetStateBad} \hat{=} \text{RVThreadGetStateInvalidArgument} \gg \text{RequestNoOp}$$

Execution of the request consists of a good execution or an error execution.

$$\text{ExecuteThreadGetState} \hat{=} \text{ThreadGetStateGood} \vee \text{ThreadGetStateBad}$$

The full specification for kernel processing of a validated **thread_get_state** request consists of processing the request followed by its execution.

$$\text{ThreadGetState} \hat{=} \text{ProcessingThreadGetState} \ ; \ \text{ExecuteThreadGetState}$$

9.11 thread_info

The request **thread_info** returns a specified type of information about a thread. The two valid choices for information types are the thread's execution status and statistics, or its scheduling parameters.

9.11.1 Client Interface

kern_return_t	thread_info	
	(mach_port_t	<i>target_thread_name,</i>
	int	<i>flavor;</i>
	thread_info_t	<i>thread_info,</i>
	mach_msg_type_number_t*	<i>thread_infoCnt);</i>

9.11.1.1 Input Parameters The following input parameters are provided by the client of a **thread_info** request:

- *target_thread_name?* — the client's name for the thread whose information is to be returned
- *flavor?* — the type of information that is to be returned. The recognized information types are *Thread_basic_info* and *Thread_sched_info*.
- *thread_infoCnt?* — the maximum amount of return information that can be handled by the client

```

ThreadInfo ClientInputs
target_thread_name? : NAME
flavor? : THREAD_INFO_TYPE
thread_infoCnt? : N

```

A **thread_info** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_info_id* and has a body consisting of *flavor?* and *thread_infoCnt?*.

```

Invoke ThreadInfo
Invoke MachMsg
ThreadInfo ClientInputs
name? = target_thread_name?
operation? = Thread_info_id
msg_body = Thread_info_type_and_count_to_text(flavor?, thread_infoCnt?)

```

9.11.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_info** request:

- *return!* — the status of the request
- *thread_info!* — the information about the target thread
- *thread_infoCnt!* — the size of the information returned in *thread_info!*

```

ThreadInfo ClientOutputs
return! : KERNEL_RETURN
thread_info! : THREAD_INFO
thread_infoCnt! : N

```

```

ThreadInfo ReceiveReply
Invoke MachMsgRcv
ThreadInfo ClientOutputs
(thread_info!, thread_infoCnt!, return!)
= Text_to_info_and_count_and_status(msg_body)

```

9.11.2 Kernel Interface

9.11.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_info** request:

- *target_thread?* — the thread whose information is to be returned
- *flavor?* — the type of information that is to be returned. The recognized information types are *Thread_basic_info* and *Thread_sched_info*.

- *thread_infoCnt?* — the maximum amount of return information that can be handled by the client

<i>ThreadInfoInputs</i>
<i>target_thread?</i> : <i>THREAD</i>
<i>flavor?</i> : <i>THREAD_INFO_TYPE</i>
<i>thread_infoCnt?</i> : \mathbb{N}

9.11.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_info** request:

- *return!* — the status of the request
- *thread_info!* — the information about the target thread
- *thread_infoCnt!* — the size of the information returned in *thread_info!*

<i>ThreadInfoOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>
<i>thread_info!</i> : <i>THREAD_INFO</i>
<i>thread_infoCnt!</i> : \mathbb{N}

Upon completion of the processing of a **thread_info** request a reply message is built from the output parameters.

<i>ThreadInfoReply</i>
<i>RequestOnlyObserves</i>
<i>thread_info?</i> : <i>THREAD_INFO</i>
<i>thread_infoCnt?</i> : \mathbb{N}
<i>reply?</i> = <i>Return_thread_info(thread_info?, thread_infoCnt?)</i>

The information returned for *Thread_basic_info* is comprised of the following items:

- *user_time_value* — the total user run time for the thread
- *system_time_value* — the total system run time for the thread
- *cpu_time_value* — the cpu time used for the thread
- *thread_base_priority_value* — the base user-setable priority for the thread
- *thread_sched_priority_value* — the priority value used by the system to make scheduling decisions. This is calculated by the system based on the *thread_base_priority_value*, the scheduling policy for the thread, and other system conditions.
- *run_state_value* — a set which either contains one of the values *Running*, *Stopped*, *Waiting*, *Uninterruptible*, and *Halted*, or is empty
- *flags* — a set which either contains one of the values *Thread_flags_swapped* or *Thread_flags_idle*, or is empty

- *thread_suspend_count_value* — a thread may execute user level instructions only if this value is zero
- *sleep_time_value* — the amount of time for which a thread has been sleeping

$THREAD_FLAGS ::= Thread_flags_swapped \mid Thread_flags_idle$

<i>ThreadBasicInfo</i>
<i>user_time_value</i> : \mathbb{N}
<i>system_time_value</i> : \mathbb{N}
<i>cpu_time_value</i> : \mathbb{N}
<i>thread_base_priority_value</i> : \mathbb{Z}
<i>thread_sched_priority_value</i> : \mathbb{Z}
<i>run_state_value</i> : $\mathbb{P} RUN_STATES$
<i>flags</i> : $\mathbb{P} THREAD_FLAGS$
<i>thread_suspend_count_value</i> : \mathbb{N}
<i>sleep_time_value</i> : \mathbb{N}
<i>run_state_value</i> ≤ 1
<i>flags</i> ≤ 1

The information returned for *Thread_sched_info* is comprised of the following items:

- *thread_policy_value* — the scheduling policy in force for the thread
- *thread_sched_policy_data_value* — policy-specific data that may influence the functioning of the policy in force
- *thread_base_priority_value* — see above
- *thread_max_priority_value* — the highest priority to which *thread_base_priority_value* can be set
- *thread_sched_priority_value* — see above
- *depressed_indicator_value* — equal to *True* if the thread's scheduling priority is currently depressed to the lowest possible value, and equal to *False* otherwise. Priority depression is accomplished via **thread_switch** or **swtch_pri**.
- *priority_before_depression_value* — if the thread's scheduling priority is currently depressed, the scheduling priority of the thread before it was depressed; otherwise, equal to *thread_base_priority_value*

<i>ThreadSchedInfo</i>
<i>thread_policy_value</i> : <i>SCHED_POLICY</i>
<i>thread_sched_policy_data_value</i> : <i>SCHED_POLICY_DATA</i>
<i>thread_base_priority_value</i> : \mathbb{Z}
<i>thread_max_priority_value</i> : \mathbb{Z}
<i>thread_sched_priority_value</i> : \mathbb{Z}
<i>depressed_indicator_value</i> : <i>BOOLEAN</i>
<i>priority_before_depression_value</i> : \mathbb{Z}

The actual space required for each type of thread information is represented by the constants *Thread_basic_info_count* and *Thread_sched_info_count*.

$$\left| \begin{array}{l} \textit{Thread_basic_info_count} : \mathbb{N} \\ \textit{Thread_sched_info_count} : \mathbb{N} \end{array} \right.$$

9.11.3 Request Criteria

The following criteria are defined for the **thread_info** request.

- **C1** — Basic information (i.e., execution statistics, status and priorities) is requested, and the client has provided enough space to hold the information.

$$\begin{array}{|l} \hline \textit{C1BasicInfo} \\ \hline \textit{flavor?} : \textit{THREAD_INFO_TYPE} \\ \textit{thread_infoCnt?} : \mathbb{N} \\ \hline \textit{flavor?} = \textit{Thread_basic_info} \\ \textit{thread_infoCnt?} \geq \textit{Thread_basic_info_count} \\ \hline \end{array}$$

$$\textit{NotC1BasicInfo} \hat{=} \neg \textit{C1BasicInfo}$$

- **C2** — Scheduling information (i.e., priorities and policies) is requested, and the client has provided enough space to hold the information.

$$\begin{array}{|l} \hline \textit{C2SchedInfo} \\ \hline \textit{flavor?} : \textit{THREAD_INFO_TYPE} \\ \textit{thread_infoCnt?} : \mathbb{N} \\ \hline \textit{flavor?} = \textit{Thread_sched_info} \\ \textit{thread_infoCnt?} \geq \textit{Thread_sched_info_count} \\ \hline \end{array}$$

$$\textit{NotC2SchedInfo} \hat{=} \neg \textit{C2SchedInfo}$$

9.11.4 Return Values

Tables 33–35 describe the values returned at the completion of the request and the conditions under which each value is returned. The specification does not state what should be returned in *thread_infoCnt!* and *thread_info!* when there is an error. These values depend in the implementation, and we leave them unspecified. Note that C1 and C2 cannot simultaneously be true.

<i>thread_infoCnt!</i>	C1	C2
<i>Thread_basic_info_count</i>	T	F
<i>Thread_sched_info_count</i>	F	T
—	F	F

Table 33: Return Values for **thread_info**

<i>thread_info!</i>	C1	C2
<i>Format_basic_info(basic_info)</i>	T	F
<i>Format_sched_info(sched_info)</i>	F	T
—	F	F

Table 34: Return Values for **thread_info**

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	F
<i>Kern_success</i>	F	T
<i>Kern_invalid_argument</i>	F	F

Table 35: Return Values for **thread_info**

Each of the two types of information must be reformatted into *thread_info!* to be returned. This formatting is represented by the functions *Format_basic_info* and *Format_sched_info*.

$$\left| \begin{array}{l} \textit{Format_basic_info} : \textit{ThreadBasicInfo} \rightarrow \textit{THREAD_INFO} \\ \textit{Format_sched_info} : \textit{ThreadSchedInfo} \rightarrow \textit{THREAD_INFO} \end{array} \right.$$

The kernel maintains for each thread a set of run states drawn from the values *Halted*, *Running*, *Uninterruptible*, *Stopped* and *Waiting*. At most one of these values is returned with the basic information. The value selected is the first one in the above list that is in the run state of the thread. In the event the run state set is empty, an empty set is returned. The function *Primary_run_state* represents this mapping from internal run state to basic information.

Similarly, the kernel will only return one of the flags *Thread_flags_swapped* or *Thread_flags_idle*. If *thread?* is in *swapped_threads*, then *Thread_flags_swapped* is always returned regardless of the value of *idle_threads*.

$$\left| \begin{array}{l} \textit{Primary_run_state} : \mathbb{P} \textit{RUN_STATES} \rightarrow \mathbb{P} \textit{RUN_STATES} \\ \textit{Run_state_order} : \textit{seq} \textit{RUN_STATES} \\ \textit{Run_state_order} = \langle \textit{Halted}, \textit{Running}, \textit{Uninterruptible}, \textit{Stopped}, \textit{Waiting} \rangle \\ \forall r : \mathbb{P} \textit{RUN_STATES} \\ \bullet \textit{Primary_run_state}(r) = (\textit{Run_state_order} \upharpoonright r)\{\{1\}\} \end{array} \right.$$

Review Note:

It is important to take the relational image under $\{1\}$ rather than taking the head of the sequence since the sequence might be empty.

RVThreadBasicInfo

ThreadStatistics
ThreadPri
ThreadExecStatus
C1BasicInfo
NotC2SchedInfo
ThreadInfo Outputs
thread? : THREAD

(\exists *basic_info* : *ThreadBasicInfo*)

- *basic_info.user_time_value* = *user_time*(*thread?*)
- \wedge *basic_info.system_time_value* = *system_time*(*thread?*)
- \wedge *basic_info.cpu_time_value* = *cpu_time*(*thread?*)
- \wedge *basic_info.thread_base_priority_value* = *thread_priority*(*thread?*)
- \wedge *basic_info.thread_sched_priority_value* = *thread_sched_priority*(*thread?*)
- \wedge *basic_info.run_state_value* = *Primary_run_state*(*run_state*(*thread?*))
- \wedge *basic_info.flags* = **if** *thread?* \in *swapped_threads*
 - then** { *Thread_flags_swapped* }
 - else** (**if** *thread?* \in *idle_threads*
 - then** { *Thread_flags_idle* }
 - else** \emptyset)
- \wedge *basic_info.thread_suspend_count_value* = *thread_suspend_count*(*thread?*)
- \wedge *basic_info.sleep_time_value* = *sleep_time*(*thread?*)
- \wedge *thread_info!* = *Format_basic_info*(*basic_info*)

thread_infoCnt! = *Thread_basic_info_count*
return! = *Kern_success*

RVThreadSchedInfo

ThreadPri
ThreadSchedPolicy
NotC1BasicInfo
C2SchedInfo
ThreadInfo Outputs
thread? : THREAD

(\exists *sched_info* : *ThreadSchedInfo*)

- *sched_info.thread_policy_value* = *thread_sched_policy*(*thread?*)
- \wedge *sched_info.thread_sched_policy_data_value* = *thread_sched_policy_data*(*thread?*)
- \wedge *sched_info.thread_base_priority_value* = *thread_priority*(*thread?*)
- \wedge *sched_info.thread_max_priority_value* = *thread_max_priority*(*thread?*)
- \wedge *sched_info.thread_sched_priority_value* = *thread_sched_priority*(*thread?*)
- \wedge *sched_info.depressed_indicator_value* = **if** *thread?* \in *depressed_threads*
 - then** *True*
 - else** *False*
- \wedge *sched_info.priority_before_depression_value* = *priority_before_depression*(*thread?*)
- \wedge *thread_info!* = *Format_sched_info*(*sched_info*)

thread_infoCnt! = *Thread_sched_info_count*
return! = *Kern_success*

$RVThreadInfoBad$ $NotC1BasicInfo$ $NotC2SchedInfo$ $ThreadInfoOutputs$
$return! = Kern_invalid_argument$

9.11.5 State Changes

A **thread_info** request does not make any state changes since it only observes the system state.

9.11.6 Complete Request

The following schemas define the general form of a **thread_info** request.

$ProcessingThreadInfo$ $ProcessThreadViaThreadPortRequestGood$
$operation? = Thread_info_id$

A successful request creates a kernel reply.

$$ThreadInfoGood \hat{=} (RVThreadBasicInfo \vee RVThreadSchedInfo) \\ \gg ThreadInfoReply$$

An unsuccessful request returns an error status.

$$ThreadInfoBad \hat{=} RVThreadInfoBad \gg RequestNoOp$$

Execution of the request consists of a good execution or an error execution.

$$ExecuteThreadInfo \hat{=} ThreadInfoGood \vee ThreadInfoBad$$

The full specification for kernel processing of a validated **thread_info** request consists of processing the request followed by its execution.

$$ThreadInfo \hat{=} ProcessingThreadInfo ; ExecuteThreadInfo$$

9.12 thread_max_priority

The request **thread_max_priority** sets the maximum scheduling priority of a specified thread. This value limits the value to which the user can set the priority of the thread (using **thread_↔priority**). Since the client thread must have access to the control port of the processor set to which the thread is assigned, the request may set the maximum priority to any legal value including one that is higher than the current value. This contrasts with **thread_priority** which does not require control port access and can only lower the maximum priority.

9.12.1 Client Interface

```

kern_return_t thread_max_priority
    (mach_port_t
     mach_port_t
     int
     thread_name,
     processor_set_name,
     priority);

```

9.12.1.1 Input Parameters The following input parameters are provided by the client of a **thread_max_priority** request:

- *thread_name?* — the client's name for the the thread whose maximum priority is to be set
- *processor_set_name?* — the client's name for the control port of the processor set to which the target thread is currently assigned
- *priority?* — the desired maximum priority

```

ThreadMaxPriorityClientInputs
thread_name? : NAME
processor_set_name? : NAME
priority? : Z

```

A **thread_max_priority** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_max_priority_id* and has a body consisting of *processor_set_name?* and *priority?*.

```

Invoke ThreadMaxPriority
Invoke MachMsg
ThreadMaxPriorityClientInputs
name? = thread_name?
operation? = Thread_max_priority_id
msg_body = Name_and_number_to_text(processor_set_name?, priority?)

```

9.12.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_max_priority** request:

- *return!* — the status of the request

```

ThreadMaxPriorityClientOutputs
return! : KERNEL_RETURN

```

```

ThreadMaxPriorityReceiveReply
Invoke MachMsgRcv
ThreadMaxPriorityClientOutputs
return! = Text_to_status(msg_body)

```

9.12.2 Kernel Interface

9.12.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_max_priority** request:

- *thread?* — the thread whose maximum priority is to be set
- *processor_set?* — the processor set to which the target thread is currently assigned
- *priority?* — the desired maximum priority

<i>ThreadMaxPriorityInputs</i>
<i>thread?</i> : <i>THREAD</i>
<i>processor_set?</i> : <i>PROCESSOR_SET</i>
<i>priority?</i> : <i>Z</i>

9.12.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_max_priority** request:

- *return!* — the status of the request

<i>ThreadMaxPriorityOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

9.12.3 Request Criteria

The following criteria are defined for the **thread_max_priority** request.

Review Note:

It is assumed here that the existence of *processor_set?* has been verified by the IPC processing of the request. However, since an arbitrary time delay might occur between the IPC processing and the kernel request processing, we should probably have an additional check here that the *processor_set?* still exists.

- **C1** — The parameter *priority?* is a valid priority level.

<i>C1 ThreadMaxPriorityValidArguments</i>
<i>priority?</i> : <i>Z</i>
<i>priority?</i> ∈ <i>Priority_levels</i>

NotC1 ThreadMaxPriorityValidArguments
 $\hat{=} \neg C1 ThreadMaxPriorityValidArguments$

- **C2** — The parameter *processor_set?* is the processor set to which *thread?* is assigned.

$\begin{aligned} & \text{C2 ThreadMaxPriorityAssignedProcessorSet} \\ & \text{ThreadAndProcessorSet} \\ & \text{thread?} : \text{THREAD} \\ & \text{processor_set?} : \text{PROCESSOR_SET} \\ & \text{(thread?, processor_set?)} \in \text{thread_assigned_to} \end{aligned}$
--

$$\begin{aligned} & \text{NotC2 ThreadMaxPriorityAssignedProcessorSet} \\ & \hat{=} \text{ThreadAndProcessorSet} \wedge \neg \text{C2 ThreadMaxPriorityAssignedProcessorSet} \end{aligned}$$

9.12.4 Return Values

Table 36 describes the values returned at the completion of the request and the conditions under which each value is returned. In the case where both C1 and C2 are false, we assume *Kern_invalid_argument* is returned.

Review Note:
The prototype checks the conditions in the order C1, C2.

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	T
<i>Kern_failure</i>	T	F
<i>Kern_invalid_argument</i>	F	-

Table 36: Return Values for `thread_max_priority`

$\begin{aligned} & \text{RVThreadMaxPriorityGood} \\ & \text{C1 ThreadMaxPriorityValidArguments} \\ & \text{C2 ThreadMaxPriorityAssignedProcessorSet} \\ & \text{ThreadMaxPriorityOutputs} \end{aligned}$
$\text{return!} = \text{Kern_success}$

$\begin{aligned} & \text{RVThreadMaxPriorityWrongProcessorSet} \\ & \text{C1 ThreadMaxPriorityValidArguments} \\ & \text{NotC2 ThreadMaxPriorityAssignedProcessorSet} \\ & \text{ThreadMaxPriorityOutputs} \end{aligned}$
$\text{return!} = \text{Kern_failure}$

$\begin{aligned} & \text{RVThreadMaxPriorityInvalidArgument} \\ & \text{NotC1 ThreadMaxPriorityValidArguments} \\ & \text{ThreadMaxPriorityOutputs} \end{aligned}$
$\text{return!} = \text{Kern_invalid_argument}$

9.12.5 State Changes

A successful **thread_max_priority** request sets the maximum priority of the thread as requested. If the thread's priority is higher than this new maximum value, it is reset to the new maximum value. If the thread is currently depressed and if its priority before depression is higher than the new maximum, then the priority before depression is reset to the new maximum value. In this way the modification of $\underline{thread_max_priority}$ will be reflected in the priority of the thread when its priority depression is removed. Note that if the thread is currently depressed $\underline{thread_priority}(thread?)$ will be equal to the lowest possible priority. Thus, $Lowest_priority\{\underline{thread_priority}(thread?), priority?\} = \underline{thread_priority}(thread?)$ and no change is made to $\underline{thread_priority}$. The thread's current scheduling priority may or may not change as a result of this request, so we state no constraint on the value of $\underline{thread_sched_priority}$.

$\underline{ThreadMaxPriorityState}$ Δ <i>Threads</i> Ξ <i>TasksAndThreads</i> Δ <i>ThreadPri</i> Ξ <i>ThreadSchedPolicy</i> Ξ <i>ThreadInstruction</i> Ξ <i>ThreadExecStatus</i> Ξ <i>ThreadStatistics</i> Ξ <i>ThreadMachineState</i> Ξ <i>Exist</i> Ξ <i>SpecialPurposePorts</i> Ξ <i>ThreadAndProcessorSet</i> <i>ThreadInvariants</i> $thread? : THREAD$ $priority? : \mathbb{Z}$
$\underline{depressed_threads}' = \underline{depressed_threads}$ $\underline{thread_max_priority}' = \underline{thread_max_priority} \oplus \{thread? \mapsto priority?\}$ $\underline{priority_before_depression}' = \mathbf{if} \ thread? \in \underline{depressed_threads}$ $\quad \mathbf{then} \ \underline{priority_before_depression} \oplus \{thread?$ $\quad \quad \mapsto Lowest_priority\{\underline{priority_before_depression}(thread?), priority?\}$ $\quad \mathbf{else} \ \underline{priority_before_depression} \oplus \{thread? \mapsto \underline{thread_priority}'(thread?)\}$ $\underline{thread_priority}' = \underline{thread_priority}$ $\quad \oplus \{thread? \mapsto Lowest_priority\{\underline{thread_priority}(thread?), priority?\}$

9.12.6 Complete Request

The following schemas define the general form of a **thread_max_priority** request.

$\underline{ProcessingThreadMaxPriority}$ $\underline{ProcessThreadViaThreadPortRequestGood}$ $operation? = Thread_max_priority_id$
--

A successful request makes the state changes described in the previous section.

$$\underline{ThreadMaxPriorityGood} \hat{=} (RVThreadMaxPriorityGood \wedge ThreadMaxPriorityState) \\ \gg RequestReturnOnlyStatus$$

An unsuccessful request returns an error status.

$$\begin{aligned} & \textit{ThreadMaxPriorityBad} \\ & \hat{=} (RV\textit{ThreadMaxPriorityWrongProcessorSet} \vee RV\textit{ThreadMaxPriorityInvalidArgument}) \\ & \gg \textit{RequestNoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} & \textit{ExecuteThreadMaxPriority} \\ & \hat{=} \textit{ThreadMaxPriorityGood} \vee \textit{ThreadMaxPriorityBad} \end{aligned}$$

The full specification for kernel processing of a validated **thread_max_priority** request consists of processing the request followed by its execution.

$$\textit{ThreadMaxPriority} \hat{=} \textit{ProcessingThreadMaxPriority} \ ; \ \textit{ExecuteThreadMaxPriority}$$

9.13 thread_policy

The request **thread_policy** sets the scheduling policy for a specified thread. This value is used by the system (together with the thread priority and current conditions) to determine the current scheduling priority of the thread.

9.13.1 Client Interface

```
kern_return_t thread_policy
                (mach_port_t
                 int
                 int
                 thread_name,
                 policy,
                 data);
```

9.13.1.1 Input Parameters The following input parameters are provided by the client of a **thread_policy** request:

- *thread_name?* — the client's name for the thread whose scheduling policy is to be set
- *policy?* — the desired scheduling policy
- *data?* — policy specific data which may influence the operation of the scheduling policy

<pre><i>ThreadPolicyClientInputs</i> <i>thread_name?</i> : NAME <i>policy?</i> : SCHED_POLICY <i>data?</i> : SCHED_POLICY_DATA</pre>
--

A **thread_policy** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_policy_id* and has a body consisting of *policy?* and *data?*.

<i>Invoke ThreadPolicy</i> <i>Invoke MachMsg</i> <i>ThreadPolicyClientInputs</i>
<i>name?</i> = <i>thread_name?</i> <i>operation?</i> = <i>Thread_policy_id</i> <i>msg_body</i> = <i>Policy_and_data_to_text(policy?, data?)</i>

9.13.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_policy** request:

- *return!* — the status of the request

<i>ThreadPolicyClientOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>

<i>ThreadPolicyReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadPolicyClientOutputs</i>
<i>return!</i> = <i>Text_to_status(msg_body)</i>

9.13.2 Kernel Interface

9.13.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_policy** request:

- *thread?* — the thread whose scheduling policy is to be set
- *policy?* — the desired scheduling policy
- *data?* — policy specific data which may influence the operation of the scheduling policy

<i>ThreadPolicyInputs</i> <i>thread?</i> : <i>THREAD</i> <i>policy?</i> : <i>SCHED_POLICY</i> <i>data?</i> : <i>SCHED_POLICY_DATA</i>
--

9.13.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_policy** request:

- *return!* — the status of the request

<i>ThreadPolicyOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>

9.13.3 Request Criteria

The following criteria are defined for the **thread_policy** request.

- **C1** — The parameter *policy?* is an existing policy.

$C1ThreadPolicyValidPolicy$ $ThreadSchedPolicy$ $policy? : SCHED_POLICY$
$policy? \in \underline{supported_sp}$

$$NotC1ThreadPolicyValidPolicy \hat{=} ThreadSchedPolicy \wedge \neg C1ThreadPolicyValidPolicy$$

- **C2** — The parameter *policy?* is a permitted scheduling policy for the processor set to which *thread?* is assigned.

$C2ThreadPolicyProcessorSetPermits$ $ThreadAndProcessorSet$ $thread? : THREAD$ $policy? : SCHED_POLICY$
$thread? \in \text{dom } thread_assigned_to$ $thread_assigned_to(thread?) \in \text{dom } \underline{enabled_sp}$
$policy? \in \underline{enabled_sp}(thread_assigned_to(thread?))$

$$NotC2ThreadPolicyProcessorSetPermits \hat{=} ThreadAndProcessorSet \wedge \neg C2ThreadPolicyProcessorSetPermits$$

9.13.4 Return Values

Table 37 describes the values returned at the completion of the request and the conditions under which each value is returned. In the case where both C1 and C2 are false, we assume *Kern_invalid_argument* is returned.

Review Note:
The prototype checks the conditions in the order C1, C2.

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	T
<i>Kern_failure</i>	T	F
<i>Kern_invalid_argument</i>	F	-

Table 37: Return Values for **thread_policy**

<i>RVThreadPolicyGood</i> <i>C1ThreadPolicyValidPolicy</i> <i>C2ThreadPolicyProcessorSetPermits</i> <i>ThreadPolicyOutputs</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadPolicyNotPermitted</i> <i>C1ThreadPolicyValidPolicy</i> <i>NotC2ThreadPolicyProcessorSetPermits</i> <i>ThreadPolicyOutputs</i>
<i>return!</i> = <i>Kern_failure</i>

<i>RVThreadPolicyInvalidPolicy</i> <i>NotC1ThreadPolicyValidPolicy</i> <i>ThreadPolicyOutputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

9.13.5 State Changes

A successful **thread_policy** request sets the target thread's policy and policy specific data as requested.

<i>ThreadPolicyNewPolicy</i> Δ <i>ThreadSchedPolicy</i> <i>thread?</i> : <i>THREAD</i> <i>policy?</i> : <i>SCHED_POLICY</i> <i>data?</i> : <i>SCHED_POLICY_DATA</i>
$\underline{thread_sched_policy}' = \underline{thread_sched_policy} \oplus \{ \text{thread?} \mapsto \text{policy?} \}$ $\underline{thread_sched_policy_data}' = \underline{thread_sched_policy_data} \oplus \{ \text{thread?} \mapsto \text{data?} \}$ $\underline{supported_sp}' = \underline{supported_sp}$

The priority and maximum priority of the thread are not modified. However, the thread's current scheduling priority may or may not change as a result of this request, so we state no constraint on the value of *thread_sched_priority*.

<i>ThreadPolicyPriority</i> Δ <i>ThreadPri</i>
$\underline{thread_priority}' = \underline{thread_priority}$ $\underline{thread_max_priority}' = \underline{thread_max_priority}$ $\underline{depressed_threads}' = \underline{depressed_threads}$ $\underline{priority_before_depression}' = \underline{priority_before_depression}$

<i>ThreadPolicyState</i>
<i>ThreadPolicyNewPolicy</i>
<i>ThreadPolicyPriority</i>
Δ <i>Threads</i>
Ξ <i>TasksAndThreads</i>
Ξ <i>ThreadInstruction</i>
Ξ <i>ThreadExecStatus</i>
Ξ <i>ThreadStatistics</i>
Ξ <i>ThreadMachineState</i>
Ξ <i>Exist</i>
Ξ <i>SpecialPurposePorts</i>
Ξ <i>ThreadAndProcessorSet</i>
<i>ThreadInvariants</i>

9.13.6 Complete Request

The following schema defines the general form of **thread_policy**.

<i>Processing ThreadPolicy</i>
<i>ProcessThreadViaThreadPortRequestGood</i>
<i>operation?</i> = <i>Thread_policy_id</i>

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} \text{ThreadPolicyGood} &\hat{=} (RV\text{ThreadPolicyGood} \wedge \text{ThreadPolicyState}) \\ &\gg \text{RequestReturnOnlyStatus} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} \text{ThreadPolicyBad} &\hat{=} (RV\text{ThreadPolicyNotPermitted} \vee RV\text{ThreadPolicyInvalidPolicy}) \\ &\gg \text{RequestNoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} \text{Execute ThreadPolicy} &\hat{=} \text{ThreadPolicyGood} \vee \text{ThreadPolicyBad} \end{aligned}$$

The full specification for kernel processing of a validated **thread_policy** request consists of processing the request followed by its execution.

$$\text{ThreadPolicy} \hat{=} \text{Processing ThreadPolicy} \ ; \ \text{Execute ThreadPolicy}$$

9.14 thread_priority

The request **thread_priority** sets the priority of a specified thread. This priority is used by the system (together with the thread scheduling policy and current conditions) to determine the

current scheduling priority of the thread used when scheduling threads to run. If the priority of the thread is currently depressed, the priority before depression will be reset instead so that this request will take effect when the priority depression is aborted. Optionally, the request may also *lower* the thread maximum priority, which limits the value of the priority. The request **thread_max_priority** also changes the priority and maximum priority. However, **thread_max_priority** requires access to the control port of the processor set to which the thread is assigned, and it may *raise* the maximum priority.

9.14.1 Client Interface

```
kern_return_t thread_priority
                (mach_port_t      thread_name,
                 int               priority,
                 boolean_t        set_max);
```

9.14.1.1 Input Parameters The following input parameters are provided by the client of a **thread_priority** request:

- *thread_name?* — the client's name for the thread whose priority is to be set
- *priority?* — the desired priority
- *set_max?* — a boolean parameter equal to *True* if the thread's maximum priority value should also be set (lowered) to *priority?*

```
ThreadPriorityClientInputs
thread_name? : NAME
priority?   : Z
set_max?   : BOOLEAN
```

A **thread_priority** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_priority_id* and has a body consisting of *priority?* and *set_max?*.

```
Invoke ThreadPriority
Invoke MachMsg
ThreadPriorityClientInputs
name? = thread_name?
operation? = Thread_priority_id
msg_body = Number_and_boolean_to_text(priority?, set_max?)
```

9.14.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_priority** request:

- *return!* — the status of the request

<i>ThreadPriorityClientOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>

<i>ThreadPriorityReceiveReply</i> <i>InvokeMachMsgRcv</i> <i>ThreadPriorityClientOutputs</i> <i>return!</i> = <i>Text_to_status(msg_body)</i>
--

9.14.2 Kernel Interface

9.14.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_priority** request:

- *thread?* — the thread whose priority is to be set
- *priority?* — the desired priority
- *set_max?* — a boolean parameter equal to *True* if the thread's maximum priority value should also be set to *priority?*

<i>ThreadPriorityInputs</i> <i>thread?</i> : <i>THREAD</i> <i>priority?</i> : \mathbb{Z} <i>set_max?</i> : <i>BOOLEAN</i>
--

9.14.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_priority** request:

- *return!* — the status of the request

<i>ThreadPriorityOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>

9.14.3 Request Criteria

The following criteria are defined for the **thread_priority** request.

- **C1** — The parameter *priority?* is a valid priority level.

<i>C1ThreadPriorityValidPriority</i> <i>priority?</i> : \mathbb{Z} <i>priority?</i> \in <i>Priority_levels</i>
--

NotC1ThreadPriorityValidPriority
 $\hat{=}$ \neg *C1ThreadPriorityValidPriority*

- **C2** — The new priority is no higher than the maximum priority for *thread?*.

<i>C2 ThreadPriorityAllowedPriority</i> <i>ThreadPri</i> <i>thread? : THREAD</i> <i>priority? : Z</i> <i>(priority?, thread_max_priority(thread?)) ∉ Higher_priority</i>
--

NotC2 ThreadPriorityAllowedPriority
 $\hat{=} ThreadPri \wedge \neg C2 ThreadPriorityAllowedPriority$

9.14.4 Return Values

Table 38 describes the values returned at the completion of the request and the conditions under which each value is returned. In the case where both C1 and C2 are false we assume *Kern_invalid_argument* is returned.

Review Note:

The prototype checks the conditions in the order C1, C2.

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	T
<i>Kern_failure</i>	T	F
<i>Kern_invalid_argument</i>	F	-

Table 38: Return Values for *thread_priority*

<i>RVThreadPriorityGood</i> <i>C1 ThreadPriorityValidPriority</i> <i>C2 ThreadPriorityAllowedPriority</i> <i>ThreadPriorityOutputs</i> <i>return! = Kern_success</i>
--

<i>RVThreadPriorityPriorityTooHigh</i> <i>C1 ThreadPriorityValidPriority</i> <i>NotC2 ThreadPriorityAllowedPriority</i> <i>ThreadPriorityOutputs</i> <i>return! = Kern_failure</i>
--

<i>RVThreadPriorityInvalidPriority</i> <i>NotC1 ThreadPriorityValidPriority</i> <i>ThreadPriorityOutputs</i> <i>return! = Kern_invalid_argument</i>
--

9.14.5 State Changes

A successful **thread_priority** request makes the following changes to the system state. The *priority_before_depression* for the thread is reset to *priority?*. If the thread priority is not currently depressed, *thread_priority* for the thread is also reset to *priority?*. Otherwise, *thread_priority* does not change. In addition, if *set_max?* is *True*, the *thread_max_priority* for the thread will also be reset to *priority?*. Note that the thread's current scheduling priority may or may not change as a result of this request, so we state no constraint on the value of *thread_sched_priority*.

<p><i>ThreadPriorityInvariants</i></p> <p><i>ThreadInvariants</i></p> <ul style="list-style-type: none"> \exists <i>TasksAndThreads</i> \exists <i>ThreadSchedPolicy</i> \exists <i>ThreadInstruction</i> \exists <i>ThreadExecStatus</i> \exists <i>ThreadStatistics</i> \exists <i>ThreadMachineState</i> \exists <i>Exist</i> \exists <i>SpecialPurposePorts</i> \exists <i>ThreadAndProcessorSet</i>
--

<p><i>ThreadPriorityState</i></p> <p>Δ <i>Threads</i></p> <p>Δ <i>ThreadPri</i></p> <p><i>ThreadPriorityInvariants</i></p> <p><i>thread?</i> : <i>THREAD</i></p> <p><i>priority?</i> : \mathbb{Z}</p> <p><i>set_max?</i> : <i>BOOLEAN</i></p> <hr/> <p>$\underline{depressed_threads}' = \underline{depressed_threads}$</p> <p>$\underline{thread_priority}' = \text{if } thread? \in \underline{depressed_threads}$ then <i>thread_priority</i> else $\underline{thread_priority} \oplus \{thread? \mapsto priority?\}$</p> <p>$\underline{priority_before_depression}' = \underline{priority_before_depression}$ $\oplus \{thread? \mapsto priority?\}$</p> <p>$\underline{thread_max_priority}' = \text{if } set_max? = True$ then $\underline{thread_max_priority} \oplus \{thread? \mapsto priority?\}$ else <i>thread_max_priority</i></p>

9.14.6 Complete Request

The following schemas define the general form of a **thread_priority** request.

<p><i>Processing ThreadPriority</i></p> <p><i>Process Thread Via ThreadPortRequestGood</i></p> <hr/> <p><i>operation?</i> = <i>Thread_priority_id</i></p>

A successful request makes the state changes described in the previous section, and creates a kernel reply.

$$\begin{aligned} \textit{ThreadPriorityGood} &\hat{=} (\textit{RVThreadPriorityGood} \wedge \textit{ThreadPriorityState}) \\ &\gg \textit{RequestReturnOnlyStatus} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} \textit{ThreadPriorityBad} \\ &\hat{=} (\textit{RVThreadPriorityPriorityTooHigh} \vee \textit{RVThreadPriorityInvalidPriority}) \\ &\gg \textit{RequestNoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} \textit{ExecuteThreadPriority} \\ &\hat{=} \textit{ThreadPriorityGood} \vee \textit{ThreadPriorityBad} \end{aligned}$$

The full specification for kernel processing of a validated **thread_priority** request consists of processing the request followed by its execution.

$$\textit{ThreadPriority} \hat{=} \textit{ProcessingThreadPriority} \ ; \ \textit{ExecuteThreadPriority}$$

9.15 thread_resume and thread_resume_secure

The requests **thread_resume** and **thread_resume_secure** decrement the suspend count of a thread by 1. They may impact the thread's run states as a result. The **thread_resume_secure** request (which is used in the secure initiation of threads within a task) expects the parent task to have task creation state *Tcs_thread_state_set* (see Section 5.7). It modifies the state to *Tcs_task_ready*.

9.15.1 Client Interface

```
kern_return_t thread_resume
                (mach_port_t                                target_thread_name);

kern_return_t thread_resume_secure
                (mach_port_t                                target_thread_name);
```

9.15.1.1 Input Parameters The following input parameters are provided by the client of a **thread_resume** or **thread_resume_secure** request:

- *target_thread_name?* — the client's name for the thread that is to be resumed

<p><i>ThreadResumeClientInputs</i></p> <p><i>target_thread_name?</i> : NAME</p>

A **thread_resume** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_resume_id* and has no body.

<i>Invoke ThreadResume</i> <i>Invoke MachMsg</i> <i>ThreadResume ClientInputs</i>
<i>name?</i> = <i>target_thread_name?</i> <i>operation?</i> = <i>Thread_resume_id</i>

A **thread_resume_secure** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_resume_secure_id* and has no body.

<i>Invoke ThreadResumeSecure</i> <i>Invoke MachMsg</i> <i>ThreadResume ClientInputs</i>
<i>name?</i> = <i>target_thread_name?</i> <i>operation?</i> = <i>Thread_resume_secure_id</i>

9.15.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_resume** or **thread_resume_secure** request:

- *return!* — the status of the request

<i>ThreadResume ClientOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>
--

<i>ThreadResume ReceiveReply</i> <i>Invoke MachMsgRcv</i> <i>ThreadResume ClientOutputs</i>
<i>return!</i> = <i>Text_to_status(msg_body)</i>

9.15.2 Kernel Interface

9.15.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_resume** or **thread_resume_secure** request:

- *target_thread?* — the thread that is to be resumed

<i>ThreadResume Inputs</i> <i>target_thread?</i> : <i>THREAD</i>

9.15.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_resume** or **thread_resume_secure** request:

- *return!* — the status of the request

<i>ThreadResume Outputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>
--

9.15.3 Request Criteria

The following criteria are defined for the **thread_resume** and **thread_resume_secure** requests.

- **C1** — The suspend count of the target thread is positive.

$C1ThreadResumeSuspendCountPositive$ $ThreadExecStatus$ $target_thread? : THREAD$
$target_thread? \in \text{dom } \underline{thread_suspend_count}$ $\underline{thread_suspend_count}(target_thread?) > 0$

$$\text{Not}C1ThreadResumeSuspendCountPositive \hat{=} \\ ThreadExecStatus \wedge \neg C1ThreadResumeSuspendCountPositive$$

- **C2** — The task creation state of the target thread's owning task must be $Tcs_thread_state_set$. This criterion applies only to the **thread_resume_secure** request.

$C2ThreadResumeThreadStateSet$ $TasksAndThreads$ $TaskCreationState$ $target_thread? : THREAD$
$target_thread? \in \text{dom } \underline{owning_task}$ $\underline{owning_task}(target_thread?) \in \text{dom } \underline{task_creation_state}$ $\underline{task_creation_state}(\underline{owning_task}(target_thread?)) = Tcs_thread_state_set$

$$\text{Not}C2ThreadResumeThreadStateSet \\ \hat{=} TasksAndThreads \wedge TaskCreationState \wedge \neg C2ThreadResumeThreadStateSet$$

9.15.4 Return Values

Table 39 describes the values returned at the completion of the **thread_resume** request and the conditions under which each value is returned.

$return!$	C1
$Kern_success$	T
$Kern_failure$	F

Table 39: Return Values for **thread_resume**

$RVThreadResumeGood$ $C1ThreadResumeSuspendCountPositive$ $ThreadResumeOutputs$
$return! = Kern_success$

<i>RVThreadResumeFailure</i>
<i>NotC1ThreadResumeSuspendCountPositive</i>
<i>ThreadResumeOutputs</i>
<i>return!</i> = <i>Kern_failure</i>

Table 40 describes the values returned at the completion of the **thread_resume_secure** request and the conditions under which each value is returned. In the case where C1 and C2 are both false we assume *Kern_insufficient_permission* is returned.

Review Note:
C2 is checked first in the prototype.

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	T
<i>Kern_failure</i>	F	T
<i>Kern_insufficient_permission</i>	-	F

Table 40: Return Values for *thread_resume_secure*

<i>RVThreadResumeSecureGood</i>
<i>C1ThreadResumeSuspendCountPositive</i>
<i>C2ThreadResumeThreadStateSet</i>
<i>ThreadResumeOutputs</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadResumeSecureFailure</i>
<i>NotC1ThreadResumeSuspendCountPositive</i>
<i>C2ThreadResumeThreadStateSet</i>
<i>ThreadResumeOutputs</i>
<i>return!</i> = <i>Kern_failure</i>

<i>RVThreadResumeSecureInsufficientPermission</i>
<i>NotC2ThreadResumeThreadStateSet</i>
<i>ThreadResumeOutputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

9.15.5 State Changes

A successful **thread_resume** or **thread_resume_secure** request decrements the thread's suspend count. If, as a result, the thread's suspend count becomes zero, the run state of the thread will be modified as follows. First, the thread will be taken out of the *Stopped* and *Halted* states. In addition, if the thread is not in the *Waiting* state, it will be placed in the *Running*

state. The state *Uninterruptible* is not affected by this request. (Note that a thread may have the state *Uninterruptible* without having the state *Waiting*.) If the suspend count is not zero after it is decremented, the run state does not change. Nothing else changes due to the request.

$\text{ThreadResumeInvariants}$ ThreadInvariants $\exists \text{ThreadExist}$ $\exists \text{TasksAndThreads}$ $\exists \text{ThreadPri}$ $\exists \text{ThreadSchedPolicy}$ $\exists \text{SpecialPurposePorts}$ $\exists \text{ThreadAndProcessorSet}$
--

ThreadResumeState $\text{ThreadResumeInvariants}$ $\Delta \text{ThreadExecStatus}$ $\Delta \text{Threads}$ $\text{target_thread?} : \text{THREAD}$
$\underline{\text{thread_suspend_count}}' = \underline{\text{thread_suspend_count}}$ $\oplus \{ \text{target_thread?} \mapsto \underline{\text{thread_suspend_count}}(\text{target_thread?}) - 1 \}$ $(\underline{\text{thread_suspend_count}}'(\text{target_thread?}) = 0 \wedge \text{Waiting} \notin \underline{\text{run_state}}(\text{target_thread?}))$ $\Rightarrow \underline{\text{run_state}}' = \underline{\text{run_state}}$ $\oplus \{ \text{target_thread?} \mapsto ((\underline{\text{run_state}}(\text{target_thread?}) \setminus \{ \text{Stopped}, \text{Halted} \})$ $\cup \{ \text{Running} \}) \}$ $(\underline{\text{thread_suspend_count}}'(\text{target_thread?}) = 0 \wedge \text{Waiting} \in \underline{\text{run_state}}(\text{target_thread?}))$ $\Rightarrow \underline{\text{run_state}}' = \underline{\text{run_state}}$ $\oplus \{ \text{target_thread?} \mapsto (\underline{\text{run_state}}(\text{target_thread?}) \setminus \{ \text{Stopped}, \text{Halted} \}) \}$ $\underline{\text{thread_suspend_count}}'(\text{target_thread?}) \neq 0 \Rightarrow \underline{\text{run_state}}' = \underline{\text{run_state}}$ $\underline{\text{swapped_threads}}' = \underline{\text{swapped_threads}}$ $\underline{\text{idle_threads}}' = \underline{\text{idle_threads}}$ $\underline{\text{threads_wired}}' = \underline{\text{threads_wired}}$

For the **thread_resume_secure** request the task creation state of the parent task is changed to *Tcs_task_ready*. There is no change to the task creation state of the parent task for a **thread_resume** request.

$\text{ThreadResumeSecureState}$ $\Delta \text{TaskCreationState}$ $\exists \text{TasksAndThreads}$ $\text{target_thread?} : \text{THREAD}$ $\text{operation?} : \text{OPERATION}$
$(\text{operation?} = \text{Thread_resume_secure_id}$ $\wedge \text{target_thread?} \in \text{dom owning_task}$ $\wedge \underline{\text{task_creation_state}}' = \underline{\text{task_creation_state}}$ $\oplus \{ \text{owning_task}(\text{target_thread?}) \mapsto \text{Tcs_task_ready} \})$ $\vee (\text{operation?} = \text{Thread_resume_id}$ $\wedge \underline{\text{task_creation_state}}' = \underline{\text{task_creation_state}})$

9.15.6 Complete Request

The following schemas define the general forms of the **thread_resume** and **thread_resume_secure** requests.

<i>Processing Thread Resume</i>
<i>Process Thread Via Thread Port Request Good</i>
<i>operation?</i> = <i>Thread_resume_id</i>

<i>Processing Thread Resume Secure</i>
<i>Process Thread Via Thread Port Request Good</i>
<i>operation?</i> = <i>Thread_resume_secure_id</i>

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned}
 & \textit{ThreadResume Good} \\
 & \quad \hat{=} (RV\textit{ThreadResumeGood} \wedge \textit{ThreadResumeState} \wedge \textit{ThreadResumeSecureState}) \\
 & \quad \gg \textit{RequestReturnOnlyStatus} \\
 & \textit{ThreadResumeSecure Good} \\
 & \quad \hat{=} (RV\textit{ThreadResumeSecureGood} \wedge \textit{ThreadResumeState} \wedge \textit{ThreadResumeSecureState}) \\
 & \quad \gg \textit{RequestReturnOnlyStatus}
 \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned}
 & \textit{ThreadResumeBad} \hat{=} RV\textit{ThreadResumeFailure} \gg \textit{RequestNoOp} \\
 & \textit{ThreadResumeSecureBad} \\
 & \quad \hat{=} (RV\textit{ThreadResumeSecureFailure} \vee RV\textit{ThreadResumeSecureInsufficientPermission}) \\
 & \quad \gg \textit{RequestNoOp}
 \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned}
 & \textit{Execute ThreadResume} \hat{=} \textit{ThreadResume Good} \vee \textit{ThreadResumeBad} \\
 & \textit{Execute ThreadResumeSecure} \hat{=} \textit{ThreadResumeSecure Good} \vee \textit{ThreadResumeSecureBad}
 \end{aligned}$$

The full specification for kernel processing of a validated **thread_resume** or **thread_resume_secure** request consists of processing the request followed by its execution.

$$\begin{aligned}
 & \textit{ThreadResume} \hat{=} \textit{Processing ThreadResume} ; \textit{Execute ThreadResume} \\
 & \textit{ThreadResumeSecure} \hat{=} \textit{Processing ThreadResumeSecure} ; \textit{Execute ThreadResumeSecure}
 \end{aligned}$$

9.16 thread_set_special_port

The **thread_set_special_port** request allows a task to set a specified special port for a specified thread to be the port associated with one of the task's send rights.

9.16.1 Client Interface

```
kern_return_t thread_set_special_port
    (mach_port_t thread_name,
     int which_port,
     mach_port_t special_port_name);
```

thread_set_exception_port

Macro form

```
kern_return_t thread_set_exception_port
    (mach_port_t thread_name,
     mach_port_t special_port_name);
⇒ thread_set_special_port (thread_name, THREAD_EXCEPTION_PORT,
    special_port_name)
```

thread_set_kernel_port

Macro form

```
kern_return_t thread_set_kernel_port
    (mach_port_t thread_name,
     mach_port_t special_port_name);
⇒ thread_set_special_port (thread_name, THREAD_KERNEL_PORT,
    special_port_name)
```

9.16.1.1 Input Parameters The following input parameters are provided by the client of a **thread_set_special_port** request:

- *thread_name?* — the client's name for the thread whose special port is to be set
- *which_port?* — the type of special port that is to be set
- *special_port_name?* — the client's name for the port to which the target thread's specified special port should be set

```
ThreadSetSpecialPort ClientInputs
thread_name? : NAME
which_port? : THREAD_SPECIAL_PORTS
special_port_name? : NAME
```

A **thread_set_special_port** request is invoked by sending a message to the port indicated by *thread_name?* that has the operation field set to *Thread_set_special_port_id* and has a body consisting of *which_port?* and *special_port_name?*.


```

Invoke ThreadSetSpecialPort
Invoke
ThreadSetSpecialPort ClientInputs



```

9.16.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_set_special_port** request:

- *return!* — the status of the request

```

ThreadSetSpecialPort ClientOutputs
return! : KERNEL_RETURN
    
```

```

ThreadSetSpecialPort ReceiveReply
Invoke MachMsgRcv
ThreadSetSpecialPort ClientOutputs

return! = Text_to_status(msg_body)
    
```

9.16.2 Kernel Interface

9.16.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_set_special_port** request:

- *thread?* — the thread whose special port is to be set
- *which_port?* — the type of special port that is to be set
- *special_port?* — the port to which the target thread's specified special port should be set

```

ThreadSetSpecialPort Inputs
thread? : THREAD
which_port? : THREAD_SPECIAL_PORTS
special_port? : PORT
    
```

9.16.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_set_special_port** request:

- *return!* — the status of the request

```

ThreadSetSpecialPort Outputs
return! : KERNEL_RETURN
    
```

9.16.3 Request Criteria

The following criteria are defined for the **thread_set_special_port** request.

- **C1** — An exception port request is made.

<i>C1ThreadSetExceptionPort</i>
<i>which_port?</i> : <i>THREAD_SPECIAL_PORTS</i>
<i>which_port?</i> = <i>Thread_exception_port</i>

$$\text{Not}C1\text{ThreadSetExceptionPort} \hat{=} \neg C1\text{ThreadSetExceptionPort}$$

- **C2** — A kernel port request is made.

<i>C2ThreadSetKernelPort</i>
<i>which_port?</i> : <i>THREAD_SPECIAL_PORTS</i>
<i>which_port?</i> = <i>Thread_kernel_port</i>

$$\text{Not}C2\text{ThreadSetKernelPort} \hat{=} \neg C2\text{ThreadSetKernelPort}$$

- **C3** — The client has *Set_thread_exception_port* permission to the target thread.

Review Note:

In C3 and C4, we've begun the process of dealing with deferred permission checks under the new execution model. The first schema is used to initiate the permission checking routine and the criteria schemas will be used after the permission has been retrieved.

<i>ThreadSetSpecialPortPermCheckSTEP</i>
<i>Transition</i>
$\exists request : Request; CheckPending; ThreadSetSpecialPortInputs$ <ul style="list-style-type: none"> • $curr_bk?? = Bk_have_request(request)$ <ul style="list-style-type: none"> $\wedge request.operation = Thread_set_special_port_id$ $\wedge thread_self(thread?) = request.service_port$ $\wedge ssi = task_sid(curr_task??)$ $\wedge osi = thread_target(curr_task??, thread?)$ $\wedge breaks' = breaks$ $\oplus \{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Set_thread_exception_port, env) \}$

<i>C3ThreadCanSetExceptionPort</i>
<i>Transition</i>
<i>env</i> : <i>ENVIRONMENT</i>
<i>curr_bk??</i> = <i>Bk_have_ruling(Set_thread_exception_port, True, env)</i>

<i>NotC3ThreadCanSetExceptionPort</i>
<i>Transition</i>
<i>env</i> : <i>ENVIRONMENT</i>
<i>curr_bk??</i> = <i>Bk_have_ruling(Set_thread_exception_port, False, env)</i>

- **C4** — The client has *Set_thread_kernel_port* permission to the target thread.

<p><i>ThreadSetSpecialPortPermCheckSTKP</i></p> <p>Transition</p> <p>$\exists request : Request; CheckPending; ThreadSetSpecialPortInputs$</p> <ul style="list-style-type: none"> • $curr_bk?? = Bk_have_request(request)$ <ul style="list-style-type: none"> $\wedge request.operation = Thread_set_special_port_id$ $\wedge thread_self(thread?) = request.service_port$ $\wedge ssi = \underline{task_sid}(curr_task??)$ $\wedge osi = thread_target(curr_task??, thread?)$ $\wedge breaks' = breaks$ <p>$\oplus \{ curr_th?? \mapsto Bk_check_pending(ssi, osi, Set_thread_kernel_port, env) \}$</p>
--

<p><i>C4ThreadCanSetKernelPort</i></p> <p>Transition</p> <p>$env : ENVIRONMENT$</p> <p>$curr_bk?? = Bk_have_ruling(Set_thread_kernel_port, True, env)$</p>
--

<p><i>NotC4ThreadCanSetKernelPort</i></p> <p>Transition</p> <p>$env : ENVIRONMENT$</p> <p>$curr_bk?? = Bk_have_ruling(Set_thread_kernel_port, False, env)$</p>
--

9.16.4 Return Values

Table 41 describes the values returned at the completion of the request and the conditions under which each value is returned. Note that C1 and C2 are mutually exclusive.

Review Note:

We assume that the prototype will check the conditions in the order {C1, C2}, {C3 or C4}. However, the prototype is currently not checking C3 and C4.

<i>return</i>	C1	C2	C3	C4
<i>Kern_success</i>	T	F	T	-
<i>Kern_success</i>	F	T	-	T
<i>Kern_insufficient_permission</i>	T	F	F	-
<i>Kern_insufficient_permission</i>	F	T	-	F
<i>Kern_invalid_argument</i>	F	F	-	-

Table 41: Return Values for *thread_set_special_port*

<i>RVThreadSetExceptionPort</i> <i>C1ThreadSetExceptionPort</i> <i>NotC2ThreadSetKernelPort</i> <i>C3ThreadCanSetExceptionPort</i> <i>ThreadSetSpecialPort Outputs</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadSetKernelPort</i> <i>NotC1ThreadSetExceptionPort</i> <i>C2ThreadSetKernelPort</i> <i>C4ThreadCanSetKernelPort</i> <i>ThreadSetSpecialPort Outputs</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadCannotSetExceptionPort</i> <i>C1ThreadSetExceptionPort</i> <i>NotC2ThreadSetKernelPort</i> <i>NotC3ThreadCanSetExceptionPort</i> <i>ThreadSetSpecialPort Outputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

<i>RVThreadCannotSetKernelPort</i> <i>NotC1ThreadSetExceptionPort</i> <i>C2ThreadSetKernelPort</i> <i>NotC4ThreadCanSetKernelPort</i> <i>ThreadSetSpecialPort Outputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

<i>RVThreadSetSpecialPortInvalidPortType</i> <i>NotC1ThreadSetExceptionPort</i> <i>NotC2ThreadSetKernelPort</i> <i>ThreadSetSpecialPort Outputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

9.16.5 State Changes

A successful **thread_set_special_port** request sets the exception or kernel port of the thread to the given port.

ThreadSetExceptionPortState

ThreadInvariants
 \exists *Exist*
 \exists *Threads*
 \exists *PortNameSpace*
 Δ *SpecialPurposePorts*
 Δ *SpecialThreadPorts*
 \exists *ThreadAndProcessorSet*
thread? : *THREAD*
special_port? : *PORT*

thread_self' = *thread_self*
thread_sself' = *thread_sself*
thread_eport' = *thread_eport* \oplus {*thread?* \mapsto *special_port?*}

ThreadSetKernelPortState

ThreadInvariants
 \exists *Exist*
 \exists *Threads*
 \exists *PortNameSpace*
 Δ *SpecialPurposePorts*
 Δ *SpecialThreadPorts*
 \exists *ThreadAndProcessorSet*
thread? : *THREAD*
special_port? : *PORT*

thread_self' = *thread_self*
thread_sself' = *thread_sself* \oplus {*thread?* \mapsto *special_port?*}

Review Note:

The prototype also releases a send right on the former kernel or exception port. This can cause no-sender notifications to be sent if the number of send rights becomes zero. We have attempted to model the total number of send rights in *TotalSendRights*. However, we have not yet modeled the sending of notifications.

9.16.6 Complete Request

The general form of a **thread_set_special_port** request is

Processing ThreadSetSpecialPort

Process Thread Via ThreadPortRequestGood

operation? = *Thread_set_special_port_id*

A successful request makes the state changes described in the previous section and creates a

kernel reply.

$$\begin{aligned} & \textit{ThreadSetSpecialPortGood} \\ & \hat{=} ((RV\textit{ThreadSetExceptionPort} \wedge \textit{ThreadSetExceptionPortState}) \\ & \quad \vee (RV\textit{ThreadSetKernelPort} \wedge \textit{ThreadSetKernelPortState})) \\ & \gg \textit{ReturnOnlyStatus} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} & \textit{ThreadSetSpecialPortBad} \\ & \hat{=} (RV\textit{ThreadCannotSetExceptionPort} \vee RV\textit{ThreadCannotSetKernelPort} \\ & \quad \vee RV\textit{ThreadSetSpecialPortInvalidPortType}) \\ & \gg \textit{NoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\textit{ExecuteThreadSetSpecialPort} \hat{=} (\textit{ThreadSetSpecialPortGood} \vee \textit{ThreadSetSpecialPortBad})$$

The full specification for the kernel processing of a validated **thread_set_special_port** request consists of processing the request followed its execution.

$$\textit{ThreadSetSpecialPort} \hat{=} \textit{ProcessingThreadSetSpecialPort} \ ; \ \textit{ExecuteThreadSetSpecialPort}$$

9.17 thread_set_state and thread_set_state_secure

The requests **thread_set_state** and **thread_set_state_secure** set the machine state of a specified thread. **thread_set_state_secure** can be used only if the thread was created using **thread_create_secure**.

9.17.1 Client Interface

kern_return_t	thread_set_state	
	(mach_port_t	<i>target_thread_name,</i>
	int	<i>flavor;</i>
	thread_state_t	<i>new_state,</i>
	mach_msg_type_number_t	<i>new_state_cnt</i>);

kern_return_t	thread_set_state_secure	
	(mach_port_t	<i>target_thread_name,</i>
	int	<i>flavor;</i>
	thread_state_t	<i>new_state,</i>
	mach_msg_type_number_t	<i>new_state_cnt</i>);

9.17.1.1 Input Parameters The following input parameters are provided by the client of a **thread_set_state** or **thread_set_state_secure** request:

- *target_thread_name?* — the client's name for the thread whose state information is to be set

- *flavor?* — the type of state information that is to be set
- *new_state?* — state information of the specified type for the target thread
- *new_state_cnt?* — the maximum size that should be assumed for the state information supplied

```

ThreadSetState ClientInputs
-----
target_thread_name? : NAME
flavor? : THREAD_STATE_INFO_TYPES
new_state? : THREAD_STATE_INFO
new_state_cnt? : N

```

A **thread_set_state** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_set_state_id* and has a body consisting of *flavor?*, *new_state?* and *new_state_cnt?*.

```

Invoke ThreadSetState
-----
Invoke MachMsg
ThreadSetState ClientInputs
-----
name? = target_thread_name?
operation? = Thread_set_state_id
msg_body = Thread_set_state_params_to_text(flavor?, new_state?, new_state_cnt?)

```

A **thread_set_state_secure** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_set_state_secure_id* and has a body consisting of *flavor?*, *new_state?* and *new_state_cnt?*

```

Invoke ThreadSetStateSecure
-----
Invoke MachMsg
ThreadSetState ClientInputs
-----
name? = target_thread_name?
operation? = Thread_set_state_secure_id
msg_body = Thread_set_state_params_to_text(flavor?, new_state?, new_state_cnt?)

```

9.17.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_set_state** or **thread_set_state_secure** request:

- *return!* — the status of the request

```

ThreadSetState ClientOutputs
-----
return! : KERNEL_RETURN

```

```

ThreadSetState ReceiveReply
-----
Invoke MachMsgRcv
ThreadSetState ClientOutputs
-----
return! = Text_to_status(msg_body)

```

9.17.2 Kernel Interface

9.17.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_set_state** or **thread_set_state_secure** request:

- *target_thread?* — the thread whose state information is to be set
- *flavor?* — the type of state information that is to be set
- *new_state?* — state information of the specified type for the target thread
- *new_state_cnt?* — the maximum size that should be assumed for the state information supplied

<i>ThreadSetStateInputs</i>
<i>target_thread?</i> : <i>THREAD</i>
<i>flavor?</i> : <i>THREAD_STATE_INFO_TYPES</i>
<i>new_state?</i> : <i>THREAD_STATE_INFO</i>
<i>new_state_cnt?</i> : \mathbb{N}

9.17.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_set_state** or **thread_set_state_secure** request:

- *return!* — the status of the request

<i>ThreadSetStateOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

9.17.3 Request Criteria

The following criteria are defined for the **thread_set_state** and **thread_set_state_secure** requests.

- **C1** — The parameter *target_thread?* is not equal to the client thread.

Editorial Note:

Nothing in the design states the reason that the client thread may not set its own state information. We believe that it is merely an implementation difficulty in that in order to set state information the thread must be temporarily stopped. If the client thread stopped itself, it could not set the state information.

<i>C1ThreadSetStateNotClientThread</i>
<i>ThreadsAndProcessors</i>
<i>cpu??</i> : <i>PROCESSOR</i>
<i>target_thread?</i> : <i>THREAD</i>
$(cpu??, target_thread?) \notin \underline{active_thread}$

NotC1ThreadSetStateNotClientThread

$\hat{=} \text{ThreadsAndProcessors} \wedge \neg \text{C1ThreadSetStateNotClientThread}$

- **C2** — The parameter *flavor?* is a valid type of state information, and *new_state_cnt?* is large enough for the requested state information.

$C2ThreadSetStateGoodFlavorAndCount$ $ThreadMachineState$ $target_thread? : THREAD$ $flavor? : THREAD_STATE_INFO_TYPES$ $new_state_cnt? : \mathbb{N}$
$flavor? \in \text{dom } Thread_state_count$ $(target_thread?, flavor?) \in \text{dom } thread_state$ $new_state_cnt? \geq Thread_state_count(flavor?)$

$$NotC2ThreadSetStateGoodFlavorAndCount \hat{=} ThreadMachineState \wedge \neg C2ThreadSetStateGoodFlavorAndCount$$

- **C3** — The task creation state of the target thread's owning task must be *Tcs_thread_created*. This criterion applies only to the **thread_set_state_secure** request.

$C3ThreadSetStateSecureThreadCreated$ $TasksAndThreads$ $TaskCreationState$ $target_thread? : THREAD$
$target_thread? \in \text{dom } owning_task$ $owning_task(target_thread?) \in \text{dom } task_creation_state$ $task_creation_state(owning_task(target_thread?)) = Tcs_thread_created$

$$NotC3ThreadSetStateSecureThreadCreated \hat{=} TasksAndThreads \wedge TaskCreationState \wedge \neg C3ThreadSetStateSecureThreadCreated$$

9.17.4 Return Values

Table 42 describes the values returned at the completion of the **thread_set_state** request and the conditions under which each value is returned.

<i>return!</i>	C1	C2
<i>Kern_success</i>	T	T
<i>Kern_invalid_argument</i>	otherwise	

Table 42: Return Values for **thread_set_state**

$RVThreadSetStateGood$ $C1ThreadSetStateNotClientThread$ $C2ThreadSetStateGoodFlavorAndCount$ $ThreadSetStateOutputs$
$return! = Kern_success$

<i>RVThreadSetStateInvalidFlavorOrCount</i> <i>C1 ThreadSetStateNotClientThread</i> <i>NotC2 ThreadSetStateGoodFlavorAndCount</i> <i>ThreadSetState Outputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

<i>RVThreadSetStateInvalidThread</i> <i>NotC1 ThreadSetStateNotClientThread</i> <i>ThreadSetState Outputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

Table 43 describes the values returned at the completion of the **thread_set_state_secure** request and the conditions under which each value is returned. In all cases where C1 is false we assume *Kern_invalid_argument* is returned. In the case where C1 is true and both C2 and C3 are false we assume *Kern_insufficient_permission* is returned.

Review Note:

In the prototype the order in which conditions are checked is C1, C3, C2.

<i>return!</i>	C1	C2	C3
<i>Kern_success</i>	T	T	T
<i>Kern_insufficient_permission</i>	T	-	F
<i>Kern_invalid_argument</i>	otherwise		

Table 43: Return Values for `thread_set_state_secure`

<i>RVThreadSetStateSecureGood</i> <i>C1 ThreadSetStateNotClientThread</i> <i>C2 ThreadSetStateGoodFlavorAndCount</i> <i>C3 ThreadSetStateSecureThreadCreated</i> <i>ThreadSetState Outputs</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVThreadSetStateSecureInsufficientPermission</i> <i>C1 ThreadSetStateNotClientThread</i> <i>NotC3 ThreadSetStateSecureThreadCreated</i> <i>ThreadSetState Outputs</i>
<i>return!</i> = <i>Kern_insufficient_permission</i>

<i>RVThreadSetStateSecureInvalidFlavorOrCount</i> <i>C1 ThreadSetStateNotClientThread</i> <i>NotC2 ThreadSetStateGoodFlavorAndCount</i> <i>C3 ThreadSetStateSecureThreadCreated</i> <i>ThreadSetState Outputs</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

$\begin{aligned} &RVThreadSetStateSecureInvalidThread \\ &NotC1ThreadSetStateNotClientThread \\ &ThreadSetStateOutputs \\ \hline &return! = Kern_invalid_argument \end{aligned}$
--

9.17.5 State Changes

A successful **thread_set_state** or **thread_set_state_secure** request sets the state of the thread to the supplied state information. Any information in *new_state?* in addition to that expected for *flavor?* is ignored. The run state of the thread may also change since the request must ensure that the thread is temporarily suspended and then perhaps restart it.

$\begin{aligned} &ThreadSetStateState \\ \hline &\Delta Threads \\ &\Xi TasksAndThreads \\ &\Xi ThreadPri \\ &\Xi ThreadSchedPolicy \\ &\Xi ThreadInstruction \\ &\Delta ThreadExecStatus \\ &\Xi ThreadStatistics \\ &\Delta ThreadMachineState \\ &\Xi Exist \\ &\Xi SpecialPurposePorts \\ &\Xi ThreadAndProcessorSet \\ &ThreadInvariants \\ &ThreadSetStateInputs \\ &target_thread? : THREAD \\ \hline &\underline{thread_state}' = \underline{thread_state} \oplus \{(target_thread?, flavor?) \mapsto new_state?\} \\ &ThreadDoWaitThenRelease[target_thread?/stopping_thread] \\ &\underline{swapped_threads}' = \underline{swapped_threads} \\ &\underline{idle_threads}' = \underline{idle_threads} \\ &\underline{thread_suspend_count}' = \underline{thread_suspend_count} \\ &\underline{threads_wired}' = \underline{threads_wired} \end{aligned}$
--

$\begin{aligned} &ThreadSetStateSecureState \\ \hline &\Delta TaskCreationState \\ &\Xi TasksAndThreads \\ &target_thread? : THREAD \\ &operation? : OPERATION \\ \hline &(operation? = Thread_set_state_secure_id \\ &\quad \wedge target_thread? \in \text{dom } owning_task \\ &\quad \wedge \underline{task_creation_state}' = \underline{task_creation_state} \\ &\quad \oplus \{owning_task(target_thread?) \mapsto Tcs_thread_state_set\}) \\ &\vee (operation? = Thread_set_state_id \\ &\quad \wedge \underline{task_creation_state}' = \underline{task_creation_state}) \end{aligned}$

9.17.6 Complete Request

The following schemas define the general forms of the **thread_set_state** and **thread_set_state_secure** requests.

$\begin{aligned} & \textit{Processing ThreadSetState} \\ & \textit{Process Thread Via ThreadPortRequestGood} \\ & \textit{operation?} = \textit{Thread_set_state_id} \end{aligned}$
--

$\begin{aligned} & \textit{Processing ThreadSetStateSecure} \\ & \textit{Process Thread Via ThreadPortRequestGood} \\ & \textit{operation?} = \textit{Thread_set_state_secure_id} \end{aligned}$
--

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} \textit{ThreadSetState Good} & \hat{=} (\textit{RVThreadSetStateGood} \wedge \textit{ThreadSetStateState} \\ & \quad \wedge \textit{ThreadSetStateSecureState}) \\ & \gg \textit{RequestReturnOnlyStatus} \\ \textit{ThreadSetStateSecure Good} & \hat{=} (\textit{RVThreadSetStateSecureGood} \wedge \textit{ThreadSetStateState} \\ & \quad \wedge \textit{ThreadSetStateSecureState}) \\ & \gg \textit{RequestReturnOnlyStatus} \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} \textit{ThreadSetState Bad} & \hat{=} (\textit{RVThreadSetStateInvalidFlavorOrCount} \vee \textit{RVThreadSetStateInvalidThread}) \\ & \gg \textit{RequestNoOp} \\ \textit{ThreadSetStateSecure Bad} & \hat{=} (\textit{RVThreadSetStateSecureInvalidFlavorOrCount} \\ & \quad \vee \textit{RVThreadSetStateSecureInsufficientPermission} \\ & \quad \vee \textit{RVThreadSetStateSecureInvalidThread}) \\ & \gg \textit{RequestNoOp} \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$\begin{aligned} \textit{Execute ThreadSetState} & \hat{=} (\textit{ThreadSetState Good} \vee \textit{ThreadSetState Bad}) \\ \textit{Execute ThreadSetStateSecure} & \hat{=} (\textit{ThreadSetStateSecure Good} \vee \textit{ThreadSetStateSecure Bad}) \end{aligned}$$

The full specification for kernel processing of a validated **thread_set_state** or **thread_set_state_secure** request consists of processing the request followed by its execution.

$$\begin{aligned} \textit{ThreadSetState} & \hat{=} \textit{Processing ThreadSetState} \ ; \ \textit{Execute ThreadSetState} \\ \textit{ThreadSetStateSecure} & \hat{=} \textit{Processing ThreadSetStateSecure} \ ; \ \textit{Execute ThreadSetState} \end{aligned}$$

9.18 thread_suspend

The request **thread_suspend** increments the suspend count of a thread by 1. If the thread was not already stopped, it will cause the thread to be stopped.

9.18.1 Client Interface

```
kern_return_t thread_suspend
(mach_port_t target_thread_name);
```

9.18.1.1 Input Parameters The following input parameters are provided by the client of a **thread_suspend** request:

- *target_thread_name?* — the client's name for the thread that is to be suspended

```
ThreadSuspend ClientInputs
target_thread_name? : NAME
```

A **thread_suspend** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_suspend_id* and has no body.

```
Invoke ThreadSuspend
Invoke MachMsg
ThreadSuspend ClientInputs
name? = target_thread_name?
operation? = Thread_suspend_id
```

9.18.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_suspend** request:

- *return!* — the status of the request

```
ThreadSuspend ClientOutputs
return! : KERNEL_RETURN
```

```
ThreadSuspend Receive Reply
Invoke MachMsgRcv
ThreadSuspend ClientOutputs
return! = Text_to_status(msg_body)
```

9.18.2 Kernel Interface

9.18.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_suspend** request:

- *target_thread?* — the thread that is to be suspended

<i>ThreadSuspend Inputs</i> <i>target_thread? : THREAD</i>

9.18.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_suspend** request:

- *return!* — the status of the request

<i>ThreadSuspend Outputs</i> <i>return! : KERNEL_RETURN</i>
--

9.18.3 Request Criteria

No criteria are defined for the **thread_suspend** request.

9.18.4 Return Values

Table 44 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>
<i>Kern_success</i>

Table 44: Return Values for **thread_suspend**

<i>RVThreadSuspendGood</i> <i>ThreadSuspend Outputs</i> <i>return! = Kern_success</i>

9.18.5 State Changes

A successful **thread_suspend** request increments the thread's suspend count. The thread will obtain the run state of *Stopped*. (Note it is possible that the thread already has this state.) A thread in *Stopped* status cannot execute any user level instructions or system traps. If a thread is suspending itself, then it will block (see Section 9.1.4.3). Otherwise, the run state *Running* will be removed by *ThreadDoWait* (see Section 9.1.4.3). The OSF documentation states that any system traps which are in progress when a thread is suspended will return after the thread resumes (via **thread_resume**).

<p><i>ThreadSuspendState</i></p> <hr/> <p><i>ThreadInvariants</i></p> <p>\exists <i>ThreadExist</i></p> <p>Δ <i>Threads</i></p> <p>\exists <i>TasksAndThreads</i></p> <p>\exists <i>ThreadPri</i></p> <p>\exists <i>ThreadSchedPolicy</i></p> <p>Δ <i>ThreadStatistics</i></p> <p>Δ <i>ThreadExecStatus</i></p> <p>\exists <i>Events</i></p> <p>\exists <i>ThreadSampling</i></p> <p>\exists <i>SpecialPurposePorts</i></p> <p>\exists <i>ThreadAndProcessorSet</i></p> <p><i>ThreadsAndProcessors</i></p> <p><i>target_thread?</i> : <i>THREAD</i></p> <p><i>cpu??</i> : <i>PROCESSOR</i></p> <hr/> <p>$cpu?? \in \text{dom } \underline{active_thread}$</p> <p>$\underline{threads_wired}' = \underline{threads_wired}$</p> <p>$\underline{thread_suspend_count}' = \underline{thread_suspend_count}$</p> <p>$\oplus \{target_thread? \mapsto \underline{thread_suspend_count}(target_thread?) + 1\}$</p> <p>let <i>rs</i> == $\underline{run_state} \oplus \{target_thread? \mapsto \underline{run_state}(target_thread?) \cup \{Stopped\}\}$</p> <p>• ($(\underline{thread_suspend_count}(target_thread?) = 0$ $\wedge target_thread? = \underline{active_thread}(cpu??)$ $\Rightarrow ThreadBlock[target_thread?/blocking_thread, rs/init_run_state])$</p> <p>$\wedge ((\underline{thread_suspend_count}(target_thread?) = 0$ $\wedge target_thread? \neq \underline{active_thread}(cpu??)$ $\Rightarrow ThreadDoWait[target_thread?/stopping_thread, rs/init_run_state])$</p> <p>$\wedge (\underline{thread_suspend_count}(target_thread?) \neq 0$ $\Rightarrow \underline{run_state}' = \underline{run_state}$ $\wedge \underline{swapped_threads}' = \underline{swapped_threads}$ $\wedge \underline{idle_threads}' = \underline{idle_threads})$</p>

Review Note:

The DTOS KID states that unpredictable results may occur if a program suspends a thread and alters its user state so that its direction is changed upon resuming.

9.18.6 Complete Request

The following schema defines the general form of **thread_suspend**.

<p><i>Processing ThreadSuspend</i></p> <hr/> <p><i>Process Thread Via ThreadPortRequestGood</i></p> <hr/> <p><i>operation?</i> = <i>Thread_suspend_id</i></p>

A successful request makes the state changes described in the previous section and creates a

kernel reply.

$$\begin{aligned} \text{ThreadSuspend Good} &\hat{=} (\text{RVThreadSuspendGood} \wedge \text{ThreadSuspendState}) \\ &\gg \text{RequestReturnOnlyStatus} \end{aligned}$$

Execution of the request consists of a good execution.

$$\text{Execute ThreadSuspend} \hat{=} \text{ThreadSuspend Good}$$

The full specification for kernel processing of a validated **thread_suspend** request consists of processing the request followed by its execution.

$$\text{ThreadSuspend} \hat{=} \text{Processing ThreadSuspend} \ ; \ \text{Execute ThreadSuspend}$$

9.19 thread_terminate

The request **thread_terminate** permanently stops execution of a thread.

9.19.1 Client Interface

```
kern_return_t thread_terminate
(mach_port_t target_thread_name);
```

9.19.1.1 Input Parameters The following input parameters are provided by the client of a **thread_terminate** request:

- *target_thread_name?* — the client's name for the thread to be destroyed

```
Thread Terminate ClientInputs
target_thread_name? : NAME
```

A **thread_terminate** request is invoked by sending a message to the port indicated by *target_thread_name?* that has the operation field set to *Thread_terminate_id* and has no body.

```
Invoke Thread Terminate
Invoke MachMsg
Thread Terminate ClientInputs
name? = target_thread_name?
operation? = Thread_terminate_id
```

9.19.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **thread_terminate** request:

- *return!* — the status of the request

<i>Thread Terminate Client Outputs</i>
<i>return! : KERNEL_RETURN</i>

<i>Thread Terminate Receive Reply</i>
<i>Invoke MachMsgRcv</i>
<i>Thread Terminate Client Outputs</i>
<i>return! = Text_to_status(msg_body)</i>

9.19.2 Kernel Interface

9.19.2.1 Input Parameters The following input parameters are provided to the kernel for a **thread_terminate** request:

- *target_thread?* — the thread to be destroyed

<i>Thread Terminate Inputs</i>
<i>target_thread? : THREAD</i>

9.19.2.2 Output Parameters The following output parameters are returned by the kernel for a **thread_terminate** request:

- *return!* — the status of the request

<i>Thread Terminate Outputs</i>
<i>return! : KERNEL_RETURN</i>

9.19.3 Request Criteria

No criteria are defined for the **thread_terminate** request.

9.19.4 Return Values

Table 45 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>
<i>Kern_success</i>

Table 45: Return Values for **thread_terminate**

Review Note:

It is actually possible for the prototype to return *Kern_failure*. There are two cases where this appears to happen.

1. Someone else has already initiated a **thread_terminate** request on *target_thread_name?*. Although the current request returns *Kern_failure* and does not itself destroy the target thread, the thread is still destroyed by the other request which is in progress.
2. The client thread is currently being terminated itself. In this case, the client thread seems to hasten its own termination rather than finishing the current request. Thus, unless there were additional termination requests in progress for the target thread, it is *not* terminated.

It doesn't appear that our model is deep enough to handle either of these cases.

<i>RVThreadTerminateGood</i>
<i>ThreadTerminateOutputs</i>
<i>return!</i> = <i>Kern_success</i>

9.19.5 State Changes

A successful **thread_terminate** request destroys the thread. The terminated thread is removed from the set of existing threads, and from its relationship with its parent task.

<i>ThreadTerminateStateExist</i>
Δ <i>ThreadExist</i>
Δ <i>TasksAndThreads</i>
<i>target_thread?</i> : <i>THREAD</i>
$\underline{t}hread_exists' = \underline{t}hread_exists \setminus \{target_thread?\}$
$\underline{t}ask_thread_rel' = \underline{t}ask_thread_rel \triangleright \{target_thread?\}$

The processor assignment of the thread is also removed.

<i>ThreadTerminateStateThreadAndProcessorSet</i>
Δ <i>ThreadAndProcessorSet</i>
<i>target_thread?</i> : <i>THREAD</i>
$\underline{t}hread_assignment_rel' = \{target_thread?\} \triangleleft \underline{t}hread_assignment_rel$
$\underline{e}nabled_sp' = \underline{e}nabled_sp$
$\underline{p}s_max_priority' = \underline{p}s_max_priority$

The thread no longer has information associated with it regarding priorities, scheduling policies, statistics and sampling.

<i>ThreadTerminateStatePriority</i>
Δ <i>ThreadPri</i>
<i>target_thread?</i> : <i>THREAD</i>
$\underline{t}hread_priority' = \{target_thread?\} \triangleleft \underline{t}hread_priority$
$\underline{t}hread_max_priority' = \{target_thread?\} \triangleleft \underline{t}hread_max_priority$
$\underline{t}hread_sched_priority' = \{target_thread?\} \triangleleft \underline{t}hread_sched_priority$
$\underline{d}epressed_threads' = \underline{d}epressed_threads \setminus \{target_thread?\}$
$\underline{p}riority_before_depression' = \{target_thread?\} \triangleleft \underline{p}riority_before_depression$

ThreadTerminateStateSchedPolicy

Δ *ThreadSchedPolicy*
target_thread? : *THREAD*

$\underline{t}hread_sched_policy' = \{target_thread?\} \triangleleft \underline{t}hread_sched_policy$
 $\underline{t}hread_sched_policy_data = \{target_thread?\} \triangleleft \underline{t}hread_sched_policy_data$
 $\underline{s}upported_sp' = \underline{s}upported_sp$

ThreadTerminateStateStatistics

Δ *ThreadStatistics*
target_thread? : *THREAD*

$\underline{u}ser_time' = \{target_thread?\} \triangleleft \underline{u}ser_time$
 $\underline{s}ystem_time' = \{target_thread?\} \triangleleft \underline{s}ystem_time$
 $\underline{c}pu_time' = \{target_thread?\} \triangleleft \underline{c}pu_time$
 $\underline{s}leep_time' = \{target_thread?\} \triangleleft \underline{s}leep_time$

ThreadTerminateStateSampling

Δ *ThreadSampling*
target_thread? : *THREAD*

$\underline{s}ampled_threads' = \underline{s}ampled_threads \setminus \{target_thread?\}$
 $\underline{t}hread_sample_types' = \{target_thread?\} \triangleleft \underline{t}hread_sample_types$
 $\underline{t}hread_sample_sequence_number'$
 $\quad = \{target_thread?\} \triangleleft \underline{t}hread_sample_sequence_number$
 $\underline{t}hread_samples' = \{target_thread?\} \triangleleft \underline{t}hread_samples$

The thread no longer has information associated with it regarding execution status. If the current thread is terminating itself, then it uses thread blocking to start another thread running.

Review Note:

The active thread on the CPU might change too. We have not modeled this change at all in the FTLS and say nothing about it here either.

ThreadTerminateStateExecStatus

Δ *ThreadExecStatus*
ThreadsAndProcessors
ProcessorAndProcessorSet
target_thread? : *THREAD*
cpu?? : *PROCESSOR*

$((cpu??, target_thread?) \in \underline{a}ctive_thread$
 $\Rightarrow (\mathbf{let} \underline{i}nit_run_state == \{target_thread?\} \triangleleft \underline{r}un_state$
 $\quad \bullet \text{ThreadBlock}[target_thread?/blocking_thread]))$
 $((cpu??, target_thread?) \notin \underline{a}ctive_thread$
 $\Rightarrow \underline{r}un_state' = \{target_thread?\} \triangleleft \underline{r}un_state$
 $\quad \wedge \underline{s}wapped_threads' = \underline{s}wapped_threads \setminus \{target_thread?\})$
 $\underline{t}hread_suspend_count' = \{target_thread?\} \triangleleft \underline{t}hread_suspend_count$
 $\underline{t}hreads_wired' = \underline{t}hreads_wired \setminus \{target_thread?\}$

All special ports are removed from the thread.

Review Note:

The prototype releases send rights on the sself port and exception port. This can cause no-sender notifications to be sent if the number of send rights becomes zero. We have attempted to model the total number of send rights in *TotalSendRights*. However, we have not yet modeled the sending of notifications.

Thread TerminateStateSpecialPorts

Δ *SpecialPurposePorts*
 Δ *SpecialThreadPorts*
target_thread? : *THREAD*

$\underline{thread_self}' = \{target_thread?\} \triangleleft \underline{thread_self}$
 $\underline{thread_sself}' = \{target_thread?\} \triangleleft \underline{thread_sself}$
 $\underline{thread_cport}' = \{target_thread?\} \triangleleft \underline{thread_cport}$

The self port of the thread is destroyed.

Thread TerminateStateSelfPort

Δ *Ipc*
 Δ *SpecialThreadPorts*
target_thread? : *THREAD*

let *port* == *thread_self*(*target_thread?*)
 • *PortDestroy*

Any event for which the thread was waiting is disassociated from the thread. The corresponding event count is also incremented by 1.

Review Note:

I purposely violated the indentation conventions in this schema to show the nesting of the logical formulas.

Thread TerminateStateEvent

Δ *Events*
target_thread? : *THREAD*

$\underline{thread_waiting}' = \underline{thread_waiting} \triangleright \{target_thread?\}$
 $\underline{event_count}'$
 = { *event* : *EVENT_COUNTER*; *count* : \mathbb{N}
 | *event* \in dom $\underline{event_count}$
 \wedge (((*event*, *target_thread?*) \notin $\underline{thread_waiting}$
 \wedge *count* = $\underline{event_count}$ (*event*))
 \vee ((*event*, *target_thread?*) \in $\underline{thread_waiting}$
 \wedge *count* = $\underline{event_count}$ (*event*) + 1))
 • (*event*, *count*) }

State information is no longer available for the thread.

$\text{Thread TerminateStateMachineState}$ $\Delta \text{ ThreadMachineState}$ $\text{target_thread?} : \text{THREAD}$
$\underline{\text{thread_state}}$ $= \{ \text{info} : \text{THREAD_STATE_INFO_TYPES}$ $\quad \bullet (\text{target_thread?}, \text{info}) \}$ $\triangleleft \underline{\text{thread_state}}$

The thread no longer has an instruction pointer.

$\text{Thread TerminateStateInstr}$ $\Delta \text{ ThreadInstruction}$ $\text{target_thread?} : \text{THREAD}$
$\underline{\text{instruction_pointer}}$ = $\{ \text{target_thread?} \} \triangleleft \underline{\text{instruction_pointer}}$

No other changes occur in the system state.

$\text{Thread TerminateState}$ ThreadInvariants $\Delta \text{ Threads}$ $\exists \text{ SpecialTaskPorts}$ $\text{Thread TerminateStateExist}$ $\text{Thread TerminateStatePriority}$ $\text{Thread TerminateStateSchedPolicy}$ $\text{Thread TerminateStateExecStatus}$ $\text{Thread TerminateStateStatistics}$ $\text{Thread TerminateStateSampling}$ $\text{Thread TerminateStateSpecialPorts}$ $\text{Thread TerminateStateThreadAndProcessorSet}$ $\text{Thread TerminateStateEvent}$ $\text{Thread TerminateStateMachineState}$ $\text{Thread TerminateStateInstr}$
--

9.19.6 Complete Request

The following schema defines the general form of **thread_terminate**.

$\text{Processing Thread Terminate}$ $\text{Process Thread Via ThreadPortRequestGood}$
$\text{operation?} = \text{Thread_terminate_id}$

A request makes the state changes described in the previous section and creates a kernel reply.

$$\text{Thread Terminate Good} \hat{=} (\text{RVThread Terminate Good} \wedge \text{Thread Terminate State})$$

$$\gg \text{RequestReturnOnlyStatus}$$

Execution of the request consists of a good execution.

$$\text{Execute Thread Terminate} \hat{=} \text{Thread Terminate Good}$$

The full specification for kernel processing of a validated **thread_terminate** request consists of processing the request followed by its execution.

ThreadTerminate $\hat{=}$ *Processing ThreadTerminate* ; *Execute ThreadTerminate*

Section 10 Virtual Memory Requests

10.1 Introduction to Virtual Memory Requests

This chapter describes the virtual memory kernel requests in DTOS.

10.1.1 Constants and Types

The following defines identifiers that are used to represent each of the requests. They are partitioned into *Vm_task_ops* and *Vm_wire_id*:

```

Vm_allocate_id, Vm_allocate_secure_id, Vm_copy_id, Vm_deallocate_id,
Vm_inherit_id, Vm_machine_attribute_id, Vm_map_id, Vm_protect_id,
Vm_read_id, Vm_region_id, Vm_region_secure_id, Vm_statistics_id,
Vm_write_id : OPERATION
Vm_wire_id : OPERATION
Vm_task_ops : P OPERATION

```

```

{Vm_allocate_id, Vm_allocate_secure_id, Vm_copy_id, Vm_deallocate_id,
Vm_inherit_id, Vm_machine_attribute_id, Vm_map_id, Vm_protect_id,
Vm_read_id, Vm_region_id, Vm_region_secure_id, Vm_statistics_id, Vm_write_id}
Values_partition Vm_task_ops
Vm_task_ops ⊆ Allowed_mach_services(Pc_task)
Vm_wire_id ∈ Allowed_mach_services(Pc_host_control)

```

10.1.2 Required Permissions

For each operation there is a primary permission that is required to perform the operation. We define here the portion of the *Required_permission* function that pertains to vm requests.

```

{(Vm_allocate_id, Allocate_vm_region),
 (Vm_allocate_secure_id, Allocate_vm_region),
 (Vm_copy_id, Copy_vm),
 (Vm_deallocate_id, Deallocate_vm_region),
 (Vm_inherit_id, Set_vm_region_inherit),
 (Vm_machine_attribute_id, Access_machine_attribute),
 (Vm_map_id, Map_vm_region),
 (Vm_protect_id, Chg_vm_region_prot),
 (Vm_read_id, Read_vm_region),
 (Vm_region_id, Get_vm_region_info),
 (Vm_region_secure_id, Get_vm_region_info),
 (Vm_statistics_id, Get_vm_statistics),
 (Vm_write_id, Write_vm_region)}
⊂ Required_permission

```

10.1.3 Invariant Information

No invariants are stated in this version of the VM Requests chapter.

10.1.4 General Information

10.1.4.1 Regions The following functions are needed to determine the pages specified by a request.

- *Get_page(va)* — determines the page index for the page of a virtual address *va*.
- *Get_offset(va)* — determines the offset on the page of a virtual address *va*.
- *Page_start(va)* — maps a virtual address *va* to the virtual address at the beginning of its page.
- *Address_num(va)* — maps a virtual address *va* to a number on which calculations can be performed.
- *Relative_addr(addr, n)* — calculates the address *n* bytes past the address *addr* if such an address exists.
- *Page_aligned* — denotes the set of virtual addresses that are the beginning of a virtual page.

We assume that *Vm_start* and *Relative_addr(Vm_end, 1)* are page aligned.

Review Note:

It might make sense to move these axioms and the *VAMWord* schema to the state chapter.

Editorial Note:

The definition of these functions as globals implies that there is a single global page size. This may not be true in a distributed environment with multiple processors of different types. The prototype uses a single global page size.

$Get_page : VIRTUAL_ADDRESS \twoheadrightarrow PAGE_INDEX$ $Get_offset : VIRTUAL_ADDRESS \twoheadrightarrow PAGE_OFFSET$ $Page_start : VIRTUAL_ADDRESS \rightarrow VIRTUAL_ADDRESS$ $Page_aligned : \mathbb{P} VIRTUAL_ADDRESS$ $Address_num : VIRTUAL_ADDRESS \twoheadrightarrow \mathbb{N}$ $Relative_addr : VIRTUAL_ADDRESS \times \mathbb{N} \rightarrow VIRTUAL_ADDRESS$
$Page_start \circ Page_start = Page_start$ $\forall va_1, va_2 : VIRTUAL_ADDRESS$ <ul style="list-style-type: none"> • $Get_page(va_1) = Get_page(va_2) \Leftrightarrow Page_start(va_1) = Page_start(va_2)$ <li style="padding-left: 2em;">$\wedge va_1 \in Page_aligned \Leftrightarrow va_1 = Page_start(va_1)$ $\forall va_1, va_2 : VIRTUAL_ADDRESS$ <ul style="list-style-type: none"> $Get_page(va_1) = Get_page(va_2) \wedge Get_offset(va_1) = Get_offset(va_2)$ • $va_1 = va_2$ $\text{dom } Relative_addr = \{ addr : VIRTUAL_ADDRESS; n : \mathbb{N} \mid Address_num(addr) + n \in \text{ran } Address_num \}$ $\forall addr : VIRTUAL_ADDRESS; n : \mathbb{N}$ <ul style="list-style-type: none"> $(addr, n) \in \text{dom } Relative_addr$ • $Relative_addr(addr, n) = Address_num \sim (Address_num(addr) + n)$ $Vm_start \in Page_aligned$ $Relative_addr(Vm_end, 1) \in Page_aligned$

The contents of a task's address space at a particular virtual address is denoted by the function va_word .

$VAWord$ $PageAndMemory$ $AddressSpace$ $va_word : TASK \times VIRTUAL_ADDRESS \rightarrow WORD$
$\forall task : TASK; va : VIRTUAL_ADDRESS$ <ul style="list-style-type: none"> $(task, Get_page(va)) \in \text{dom } \underline{map_rel}$ <li style="padding-left: 2em;">$\wedge \underline{map_rel}(task, Get_page(va)) \in \text{dom } representing_page$ <li style="padding-left: 2em;">$\wedge representing_page(\underline{map_rel}(task, Get_page(va))) \in \text{dom } page_word_fun$ <ul style="list-style-type: none"> • $va_word(task, va)$ <li style="padding-left: 2em;">$= (page_word_fun(representing_page(\underline{map_rel}(task, Get_page(va)))))(Get_offset(va))$

We use $Region_of(va, size)$ to denote the region of $size$ bytes starting at the page containing va in some task's address space. Since a region consists of a sequence of pages, the return from this function is the set of page indices denoting pages containing an address between va and $va + size - 1$. Because of this rounding to virtual page boundaries, the amount of memory in a region may be greater than $size$.

$Region_of : VIRTUAL_ADDRESS \times \mathbb{N} \rightarrow \mathbb{P} PAGE_INDEX$
$\forall va : VIRTUAL_ADDRESS; size : \mathbb{N}$ <ul style="list-style-type: none"> • $Region_of(va, size) = \{ va_1 : VIRTUAL_ADDRESS \mid Address_num(va_1) \in Address_num(va) .. (Address_num(va) + size - 1) \}$ • $Get_page(va_1)$

$VmRegionInUse[task, address, size]$ denotes that $Region_of(address, size)$ contains at least one page that is allocated in $task$'s address space, and $VmRegionNotInUse$ denotes that none of the

pages are allocated.

$VmRegionInUse$ <hr/> $AddressSpace$ $task : TASK$ $address : VIRTUAL_ADDRESS$ $size : \mathbb{N}$ <hr/> $Region_of(address, size) \cap \underline{allocated}(\{ task \}) \neq \emptyset$

$$VmRegionNotInUse \hat{=} AddressSpace \wedge \neg VmRegionInUse$$

All addresses within a valid region must lie in the range $Vm_start..Vm_end$. We use $VmGoodRegion[address, size]$ to denote that the region of length $size$ starting at $address$ is valid.

Review Note:

Since we are assuming Vm_start and $Vm_end + 1$ are page aligned we do not need to round $address$ and $size$.

$VmGoodRegion$ <hr/> $address : VIRTUAL_ADDRESS$ $size : \mathbb{N}$ <hr/> $Address_num(address) .. Address_num(address) + size - 1$ $\subseteq Address_num(Vm_start) .. Address_num(Vm_end)$
--

Set_region_attr defines a function that maps all of the pages in a virtual memory region to a particular attribute.

$[R]$ <hr/> $Set_region_attr : (\mathbb{P}(TASK \times PAGE_INDEX) \times R)$ $\rightarrow ((TASK \times PAGE_INDEX) \leftrightarrow R)$ <hr/> $\forall region : \mathbb{P}(TASK \times PAGE_INDEX); x : R$ <ul style="list-style-type: none"> • $Set_region_attr(region, x) =$ $\{ task_va_pair : TASK \times PAGE_INDEX \mid task_va_pair \in region$ <ul style="list-style-type: none"> • $task_va_pair \mapsto x \}$

The **vm-write** request takes a vm map copy parameter that describes a region of virtual memory including the offset, the size and the task from whose address space the memory was copied. We model this with *MapCopy*.

Review Note:

In the prototype, a map copy does not contain a direct reference to the task. Although we are uncertain, it is even possible that the task from whose address space the map copy was produced no longer exists. It is conceivable that the task was destroyed after the map copy was created, and the map entries are still present since the map copy holds a reference to them. The correct solution here would be to model maps as entities in their own right independent of tasks. This would require significant changes to the state description.

<p><i>MapCopy</i></p> <p><i>task</i> : <i>TASK</i></p> <p><i>offset</i> : <i>VIRTUAL_ADDRESS</i></p> <p><i>size</i> : \mathbb{N}</p>

10.1.4.2 Parameter Packaging Functions When invoking a kernel request, the following functions package the parameters into a message body:

<p><i>Address_to_body</i> : <i>VIRTUAL_ADDRESS</i> \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Region_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times \mathbb{N}) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Region_bool_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ <i>BOOLEAN</i>) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Region_inheritance_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ <i>INHERITANCE_OPTION</i>) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Region_bool_sid_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ <i>BOOLEAN</i> \times <i>OSI</i>) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Region_bool_prot_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ <i>BOOLEAN</i> \times \mathbb{P} <i>PROTECTION</i>) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Address_data_to_body</i> : (<i>VIRTUAL_ADDRESS</i> \times <i>MapCopy</i> \times \mathbb{N}) \rightarrow <i>MESSAGE_BODY</i></p> <p><i>Name_region_prot_to_body</i> : (<i>NAME</i> \times <i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ \mathbb{P} <i>PROTECTION</i>) \rightarrow <i>MESSAGE_BODY</i></p>
--

When creating a reply message from a request, the following functions package the output parameters into a kernel reply:

<p><i>Address_to_reply</i> : <i>VIRTUAL_ADDRESS</i> \rightarrow <i>KERNEL_REPLY</i></p> <p><i>Attributes_to_reply</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ \mathbb{P} <i>PROTECTION</i> \times \mathbb{P} <i>PROTECTION</i> \times <i>INHERITANCE_OPTION</i> \times <i>BOOLEAN</i> \times <i>Capability</i> \times <i>OFFSET</i>) \rightarrow <i>KERNEL_REPLY</i></p> <p><i>Secure_attributes_to_reply</i> : (<i>VIRTUAL_ADDRESS</i> \times $\mathbb{N} \times$ \mathbb{P} <i>PROTECTION</i> \times <i>INHERITANCE_OPTION</i> \times <i>BOOLEAN</i> \times <i>Capability</i> \times \mathbb{P} <i>PROTECTION</i> \times <i>OFFSET</i> \times <i>OSI</i> \times \mathbb{P} <i>Kernel_permission</i>) \rightarrow <i>KERNEL_REPLY</i></p>

When receiving a reply message from the kernel the following functions unpack the message body to obtain the output parameters (including the return status):

```

Text_to_address_and_status : MESSAGE_BODY
  →(VIRTUAL_ADDRESS × KERNEL_RETURN)
Text_to_region_info_and_status : MESSAGE_BODY
  →(VIRTUAL_ADDRESS × N × P PROTECTION
     × P PROTECTION × INHERITANCE_OPTION × BOOLEAN
     × Capability × OFFSET × KERNEL_RETURN)
Text_to_region_secure_info_and_status : MESSAGE_BODY
  →(VIRTUAL_ADDRESS × N × P PROTECTION
     × P PROTECTION × INHERITANCE_OPTION × BOOLEAN
     × Capability × OFFSET × P PROTECTION
     × OSI × P Kernel_permission × KERNEL_RETURN)

```

Review Note:

The command *Text_to_status* is also used in this chapter. It is declared in the Thread Request chapter introduction.

10.1.5 Kernel Processing

The kernel performs processing for a VM request only when it detects a break indicating that a request has been received through a port of the appropriate class, *Pc_task* or *Pc_host_control*.

For a request sent to a task port, if the specified service port no longer exists, then a *Kern_invalid_argument* status code is returned.

```

NotTaskPort
ProcessRequest
≡ Mach
reply_to_port! : P PORT
reply! : KERNEL_REPLY
return! : KERNEL_RETURN

pc? = Pc_task
operation? ∈ Vm_task_ops
service_port? ∉ dom self_task

reply_to_port! = reply_to_port?
return! = Kern_invalid_argument

```

For a **vm_wire** request, which must be sent to a host control port, if the service port no longer exists, then a *Kern_invalid_host* status code is returned.

<p><i>NotHostTaskPort</i></p> <hr/> <p><i>ProcessRequest</i> $\exists Mach$ <i>reply_to_port!</i> : $\mathbb{P} PORT$ <i>reply!</i> : <i>KERNEL_REPLY</i> <i>return!</i> : <i>KERNEL_RETURN</i></p> <hr/> <p><i>pc?</i> = <i>Pc_host_control</i> <i>operation?</i> = <i>Vm_wire_id</i> <i>service_port?</i> \neq <i>host_control_port</i></p> <hr/> <p><i>reply_to_port!</i> = <i>reply_to_port?</i> <i>return!</i> = <i>Kern_invalid_host</i></p>
--

$ProcessVMRequestBad \hat{=} (NotTaskPort \vee NotHostTaskPort) \gg RequestNoOp$

Otherwise, the kernel processes the request. In this case, we use the following schema to represent the parameters to the requests:

<p><i>VMParameters</i></p> <hr/> <p><i>address?</i> : <i>VIRTUAL_ADDRESS</i> <i>anywhere?</i> : <i>BOOLEAN</i> <i>copy?</i> : <i>BOOLEAN</i> <i>count?</i> : \mathbb{N} <i>cur_protection?</i> : $\mathbb{P} PROTECTION$ <i>data?</i> : <i>MapCopy</i> <i>data_count?</i> : \mathbb{N} <i>dest_address?</i> : <i>VIRTUAL_ADDRESS</i> <i>host_priv?</i> : <i>HOST</i> <i>inheritance?</i> : <i>INHERITANCE_OPTION</i> <i>mask?</i> : <i>VIRTUAL_ADDRESS</i> <i>max_protection?</i> : $\mathbb{P} PROTECTION$ <i>memory_object?</i> : <i>Capability</i> <i>new_inheritance?</i> : <i>INHERITANCE_OPTION</i> <i>new_protection?</i> : $\mathbb{P} PROTECTION$ <i>protection?</i> : $\mathbb{P} PROTECTION$ <i>obj_sid?</i> : <i>OSI</i> <i>offset?</i> : <i>OFFSET</i> <i>set_maximum?</i> : <i>BOOLEAN</i> <i>shared?</i> : <i>BOOLEAN</i> <i>size?</i> : \mathbb{N} <i>source_address?</i> : <i>VIRTUAL_ADDRESS</i> <i>target_task?</i> : <i>TASK</i> <i>wired_access?</i> : $\mathbb{P} PROTECTION$</p>
--

The interpretation of the components of this schema are:

address? — starting address for a region.

anywhere? — a Boolean indicating whether the region can be anywhere in the target task's address space.

copy? — a Boolean indicating whether a copy is made of an area of a memory object.

count? — the number of bytes in a region.

cur_protection? — the initial current protection for a region.

data? — a copy of a portion of a memory map.

data_count? — the number of bytes in a data array (ignored).

dest_address? — starting address for the destination region.

host_priv? — the host on which the target task executes.

inheritance? — the inheritance attribute for the region.

mask? — alignment restrictions for the starting address of a region.

max_protection? — the maximum protection for a region.

memory_object? — the port naming a memory object.

new_inheritance? — the new inheritance attribute for the region.

new_protection? — the new protection for the region.

protection? — the current protection for a region including those protections

obj_sid? — the security identifier for a region.

offset? — an offset within a memory object, in bytes.

set_maximum? — a Boolean indicating whether the maximum protection or the current protection should be set.

shared? — a Boolean indicating whether the region is shared with another task.

size? — the number of bytes in a region.

source_address? — starting address for the source region.

target_task? — the task to whose address space the command applies.

wired_access? — the pageability of a region.

The following schema determines the target task based upon the task service port to which a task operation request has been sent.

<i>MessageToVMParameters</i> <i>ProcessRequest</i> <i>SpecialTaskPorts</i> <i>VMParameters</i>
<i>pc?</i> = <i>Pc_task</i> <i>operation?</i> ∈ <i>Vm_task_ops</i> <i>service_port?</i> ∈ dom <i>self_task</i> <i>target_task?</i> = <i>self_task(service_port?)</i>

The following schema verifies that a **vm_wire** request has been sent to the host control port.

<p><i>MessageToHostParameters</i></p> <hr/> <p><i>ProcessRequest</i> <i>HostsAndPorts</i> <i>VMParameters</i></p> <hr/> <p><i>pc?</i> = <i>Pc_host_control</i> <i>operation?</i> = <i>Vm_wire_id</i> <i>service_port?</i> = <i>host_control_port</i></p>
--

10.1.6 Security Server Request

For some requests (e.g., **vm_allocate**) a second security check is needed. In this case the access vector cache will be checked for the needed information. If the information is not present (or not valid for the client thread) the security server is queried, and the kernel must wait for the response before continuing the execution of the request. We represent this waiting time by adding an element to the set of pending requests that contains the current request, the client thread, and the OSI associated with the security server request. The schema *VmSecurityRequest* checks the cache for permission *perm* from the subject *ssi* to the object *osi*. If it is not found (i.e., *Cache_undefined*), the request is placed in *PENDINGREQUEST*.

$$\left| \begin{array}{l} Vm_request_to_pending_request : Request \times THREAD \times OSI \\ \rightarrow PENDINGREQUEST \end{array} \right.$$

<p><i>VmSecurityRequest</i></p> <hr/> <p><i>Transition</i> <i>KernelAllows</i> Δ <i>PendingRequests</i> <i>ThreadsAndProcessors</i> <i>Request?</i> <i>perm</i> : <i>PERMISSION</i> <i>ssi</i> : <i>SSI</i> <i>osi</i> : <i>OSI</i></p> <hr/> <p>let <i>thread</i> == <i>active_thread(cpu??)</i> <ul style="list-style-type: none"> • <i>cache_allows(thread, ssi, osi, perm)</i> = <i>Cache_undefined</i> \wedge <i>PENDINGREQUEST'</i> = <i>PENDINGREQUEST</i> \uplus [<i>Vm_request_to_pending_request</i>(θ <i>Request?</i>, <i>thread</i>, <i>osi</i>)] </p>
--

After the security request has been processed, the kernel request is removed from the set of pending requests by the schema *VmContinue*. The client thread and the OSI supplied in the security server request are also retrieved.

VmContinue

Δ *PendingRequests*
SpecialTaskPorts
VMPParameters
Request?
thread' : THREAD
obj_sid' : OSI

\exists *pending_request : PENDREQUEST*

- *pending_request* \in *PENDINGREQUEST*
 - \wedge *pending_request*
 - $=$ *Vm_request_to_pending_request*(θ *Request?*, *thread'*, *obj_sid'*)
 - \wedge *PENDINGREQUEST'* $=$ *PENDINGREQUEST* \cup [*pending_request*]
 - \wedge *pc?* $=$ *Pc_task*
 - \wedge *operation?* \in *Vm_task_ops*

If the required permission is already in the access vector cache, the security server request will not be necessary. *VmNoSecurityRequest* describes this case.

VmNoSecurityRequest

Transition
KernelAllows
ThreadsAndProcessors
Request?
perm : PERMISSION
ssi : SSI
osi : OSI

let *thread* $==$ *active_thread*(*cpu??*)

- *cache_allows*(*thread*, *ssi*, *osi*, *perm*) \neq *Cache_undefined*

We now describe the individual virtual memory requests.

10.2 **vm_allocate** and **vm_allocate_secure**

The **vm_allocate** and **vm_allocate_secure** task requests allocate a zero-filled region of memory in the target task's address space. The physical memory is not allocated until an executing thread references the new virtual memory, and a memory object managed by the default manager is not created until the region must be swapped out. **vm_allocate_secure** allows the client to specify a security identifier for the allocated region, while **vm_allocate** uses a default security identifier.

10.2.1 Client Interface

kern_return_t vm_allocate (mach_port_t vm_address_t* vm_size_t boolean_t	<i>target_task_name,</i> <i>address,</i> <i>size,</i> <i>anywhere);</i>
---	--

kern_return_t vm_allocate_secure	
(mach_port_t	<i>target_task_name,</i>
vm_address_t*	<i>address,</i>
vm_size_t	<i>size,</i>
boolean_t	<i>anywhere,</i>
security_id_t	<i>obj_sid);</i>

10.2.1.1 Input Parameters The following input parameters are provided by the client of a **vm_allocate** request:

- *target_task_name?* — the client's name for the task in whose virtual address space the region is to be allocated
- *address?* — the requested starting address for the region. This parameter is ignored if *anywhere?* is *True*. Otherwise, it is rounded down to the start of a page boundary.
- *size?* — the number of bytes to allocate. It is rounded up to an integer number of pages. (This differs from the interpretation of *size?* used in the other VM requests.)
- *anywhere?* — a Boolean indicating whether the allocated region can be placed anywhere in the target task's address space or must be placed at *address?*

```

VmAllocateClientInputs
target_task_name? : NAME
address? : VIRTUAL_ADDRESS
size? : N
anywhere? : BOOLEAN

```

The following additional parameter must be provided by the client of a **vm_allocate_secure** request:

- *obj_sid?* — security identifier that will be attached to the newly allocated region.

```

VmAllocateSecureClientInputs
VmAllocateClientInputs
obj_sid? : OSI

```

A **vm_allocate** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_allocate_id* and has a body consisting of *address?*, *size?*, and *anywhere?*.

```

InvokeVmAllocate
InvokeMachMsg
VmAllocateClientInputs
name? = target_task_name?
operation? = Vm_allocate_id
msg_body = Region_bool_to_body (address?, size?, anywhere?)

```

A **vm_allocate_secure** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_allocate_secure_id* and has a body consisting of *address?*, *size?*, *anywhere?*, and *obj_sid?*.

```

Invoke VmAllocateSecure
Invoke MachMsg
VmAllocateSecureClientInputs
-----
name? = target_task_name?
operation? = Vm_allocate_secure_id
msg_body = Region_bool_sid_to_body(address?, size?, anywhere?, obj_sid?)

```

10.2.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **vm_allocate** or **vm_allocate_secure** request:

- *address!* — the actual starting address for the memory object
- *return!* — the status of the request

```

VmAllocateClientOutputs
-----
address! : VIRTUAL_ADDRESS
return! : KERNEL_RETURN

```

```

VmAllocateReceiveReply
Invoke MachMsgRcv
VmAllocateClientOutputs
-----
(address!, return!) = Text_to_address_and_status(msg_body)

```

10.2.2 Kernel Interface

10.2.2.1 Input Parameters The following input parameters are provided to the kernel for a **vm_allocate** request:

- *target_task?* — the task in whose virtual address space the region is to be allocated
- *address?* — the requested starting address for the region. This parameter is ignored if *anywhere?* is *True*. Otherwise, it is rounded down to the start of a page boundary.
- *size?* — the number of bytes to allocate. It is rounded up to an integer number of pages. (This differs from the interpretation of *size?* used in the other VM requests.)
- *anywhere?* — a Boolean indicating whether the allocated region can be placed anywhere in the target task's address space or must be placed at *address?*

```

VmAllocateInputs
-----
target_task? : TASK
address? : VIRTUAL_ADDRESS
size? : N
anywhere? : BOOLEAN

```

The following additional parameter must be provided by the client of a **vm_allocate_secure** request:

- *obj_sid?* — security identifier that will be attached to the newly allocated region.

```

VmAllocateSecureInputs _____
VmAllocateInputs
obj_sid? : OSI

```

10.2.2.2 Output Parameters The following output parameters are returned by the kernel for a **vm_allocate** or **vm_allocate_secure** request:

- *address!* — the actual starting address for the memory object
- *return!* — the status of the request

```

VmAllocateOutputs _____
address! : VIRTUAL_ADDRESS
return! : KERNEL_RETURN

```

Upon completion of the processing of either a **vm_allocate** or a **vm_allocate_secure** request, a reply message is built from the output parameters.

```

VmAllocateReply _____
RequestReturn
address? : VIRTUAL_ADDRESS
reply? = Address_to_reply( address? )

```

10.2.3 Request Criteria

The following criteria are defined for the **vm_allocate** and **vm_allocate_secure** requests.

- **C1** — The security identifier for the new region can be specified by the client thread, *active_thread(cpu??)*, as determined from the result of a security policy query. The value of *obj_sid* is either *vm_port_sid(target_task?)* for a **vm_allocate** request or the input parameter *obj_sid?* for a **vm_allocate_secure** request. The binding of *obj_sid* is determined by the appropriate processing schema from Section 10.2.6.

```

C1VmAllocateGoodSecurityId _____
SubjectSid
KernelAllows
ThreadsAndProcessors
thread : THREAD
obj_sid : OSI
cpu?? : PROCESSOR
thread = active_thread(cpu??)
thread ∈ dom thread_sid
cache_allows(thread, thread_sid(thread), obj_sid, Map_vm_region)
= Cache_allowed

```

$$\text{Not}C1VmAllocateGoodSecurityId \hat{=} SubjectSid \wedge KernelAllows \\ \wedge ThreadsAndProcessors \wedge \neg C1VmAllocateGoodSecurityId$$

- **C2** — The task remains after a possible second security server query has been made. The port *service_port?* is the port through which the request was received.

$C2VmAllocateTaskRemains$
$SpecialTaskPorts$
$target_task? : TASK$
$service_port? : PORT$
$(service_port?, target_task?) \in self_task$

$$\text{Not}C2VmAllocateTaskRemains \hat{=} SpecialTaskPorts \wedge \neg C2VmAllocateTaskRemains$$

- **C3** — The parameter *size?* is greater than zero.

$C3VmAllocatePositiveSize$
$size? : \mathbb{N}$
$size? > 0$

$$\text{Not}C3VmAllocatePositiveSize \hat{=} \neg C3VmAllocatePositiveSize$$

- **C4** — The parameter *anywhere? = True*, or the addresses specified for the region (when rounded) are valid.

$C4VmAllocateGoodAddress$
$address? : VIRTUAL_ADDRESS$
$size? : \mathbb{N}$
$anywhere? : BOOLEAN$
$anywhere? = True$
$\vee (\text{let } address == Page_start(address?); size == size?$
<ul style="list-style-type: none"> • <i>VmGoodRegion</i>)

$$\text{Not}C4VmAllocateGoodAddress \hat{=} \neg C4VmAllocateGoodAddress$$

- **C5** — There is room in *target_task?*'s address space to allocate a region of length *size?* starting at a page boundary. If *anywhere? = False*, there is room starting at the beginning of the page containing *address?*.

<p><i>C5 VmAllocateRoomToAllocate</i></p> <p><i>AddressSpace</i></p> <p><i>target_task?</i> : <i>TASK</i></p> <p><i>address?</i> : <i>VIRTUAL_ADDRESS</i></p> <p><i>size?</i> : \mathbb{N}</p> <p><i>anywhere?</i> : <i>BOOLEAN</i></p> <p><i>address</i> : <i>VIRTUAL_ADDRESS</i></p>
<p>\exists <i>size</i> : \mathbb{N}</p> <ul style="list-style-type: none"> • (<i>anywhere?</i> = <i>True</i> \vee <i>address</i> = <i>Page_start</i>(<i>address?</i>)) <ul style="list-style-type: none"> \wedge <i>address</i> \in <i>Page_aligned</i> \wedge <i>size</i> = <i>size?</i> \wedge <i>VmGoodRegion</i> \wedge <i>VmRegionNotInUse</i>[<i>target_task?</i>/<i>task</i>]

$$\text{Not } C5 \text{ VmAllocateRoomToAllocate} \hat{=} \text{AddressSpace} \wedge \neg C5 \text{ VmAllocateRoomToAllocate}$$

Review Note:

Do we also require *Have_execute*, *Have_read* and *Have_write* permissions for the target task?

10.2.4 Return Values

Table 46 describes the values returned at the completion of the request and the conditions under which each value is returned. The value *address* is any address that satisfied the criterion C5. When *anywhere?* is *False*, this is the address at the start of the page containing *address?*. When *anywhere?* is *True*, the starting address of the allocated region depends upon *address?*, *size?* and the allocated pages of *target_task?*. The relationship between these three items and the address returned depends upon the implementation algorithm. In the prototype if C3 is false and C1 and C2 are true, the zero address, *Address_num*[~](0), is returned. We leave unspecified the precise address returned in cases where *return!* \neq *Kern_success*.

Editorial Note:

The algorithm currently used in the prototype will never yield a page that starts earlier in the memory than the beginning of the page containing *address?*. Thus, if the client specifies the last page and it is already allocated, the return value will be *Kern_no_space* even if there are pages available earlier in the address space.

The value of *address!* when an error occurs is undefined in the design and therefore also depends on the implementation algorithm and is left unspecified. In the case where more than one error occurs we assume that the first applicable return status from the following list is returned: *Kern_insufficient_permission*, *Kern_invalid_argument*, *Kern_invalid_address* and *Kern_no_space*.

Review Note:

The prototype checks the conditions in the order C1, C2, C3, C4 and C5.

<i>address!</i>	<i>return!</i>	C1	C2	C3	C4	C5
<i>address</i>	<i>Kern_success</i>	T	T	T	T	T
—	<i>Kern_no_space</i>	T	T	T	T	F
—	<i>Kern_invalid_address</i>	T	T	T	F	-
<i>Address_num</i> ~(0)	<i>Kern_success</i>	T	T	F	-	-
—	<i>Kern_invalid_argument</i>	T	F	-	-	-
—	<i>Kern_insufficient_permission</i>	F	-	-	-	-

Table 46: Return Values for **vm_allocate** and **vm_allocate_secure**

<i>RVVmAllocateSuccessful</i>
<i>AddressSpace</i>
<i>VmAllocate Outputs</i>
<i>C1 VmAllocateGoodSecurityId</i>
<i>C2 VmAllocateTaskRemains</i>
<i>C3 VmAllocatePositiveSize</i>
<i>C4 VmAllocateGoodAddress</i>
<i>C5 VmAllocateRoomToAllocate</i>
<i>address!</i> = <i>address</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVVmAllocateNoSpace</i>
<i>AddressSpace</i>
<i>VmAllocate Outputs</i>
<i>C1 VmAllocateGoodSecurityId</i>
<i>C2 VmAllocateTaskRemains</i>
<i>C3 VmAllocatePositiveSize</i>
<i>C4 VmAllocateGoodAddress</i>
<i>NotC5 VmAllocateRoomToAllocate</i>
<i>return!</i> = <i>Kern_no_space</i>

<i>RVVmAllocateBadAddress</i>
<i>AddressSpace</i>
<i>VmAllocate Outputs</i>
<i>C1 VmAllocateGoodSecurityId</i>
<i>C2 VmAllocateTaskRemains</i>
<i>C3 VmAllocatePositiveSize</i>
<i>NotC4 VmAllocateGoodAddress</i>
<i>return!</i> = <i>Kern_invalid_address</i>

RVVmAllocateVacuous

AddressSpace
VmAllocateOutputs
C1VmAllocateGoodSecurityId
C2VmAllocateTaskRemains
NotC3VmAllocatePositiveSize

address! = *Address_num*~(0)
return! = *Kern_success*

RVVmAllocateBadArgument

AddressSpace
VmAllocateOutputs
C1VmAllocateGoodSecurityId
NotC2VmAllocateTaskRemains

return! = *Kern_invalid_argument*

RVVmAllocateBadSecurityId

AddressSpace
VmAllocateOutputs
NotC1VmAllocateGoodSecurityId

return! = *Kern_insufficient_permission*

10.2.5 State Changes

When the request is successful, a new region *size?* in length is added to the mapped address space for *target_task?* starting at *address!* (one of the outputs calculated above). This region is initially mapped to the null memory object. The maximum protections for the new region are set so that they allow all accesses, and the current protections allow reading and writing. The *inheritance* for the region is initialized to *Inheritance_option_copy*. The initial value of 0 will be set later when the region is first accessed.

Review Note:

The *c_protection*' should take into account the access vector contents. It should be the intersection of read and write with the permissions allowed from the target task to its vm port sid.

VmAllocateState

Δ *AddressSpace*
 Δ *Protection*
 Δ *Inheritance*
Memory
address! : *VIRTUAL_ADDRESS*
size? : \mathbb{N}
target_task? : *TASK*

let *region* == { *target_task?* } \times *Region_of*(*address!*, *size?*)

- *allocated'* = *allocated* \cup *region*
 - \wedge *map_rel'*(*region*) \subseteq ({ *Null_memory* } \times *OFFSET*)
 - \wedge *m_protection'* = *m_protection*
 \oplus *Set_region_attr*(*region*, { *Read*, *Write*, *Execute* })
 - \wedge *c_protection'* = *c_protection* \oplus *Set_region_attr*(*region*, { *Read*, *Write* })
 - \wedge *i_inheritance'* = *i_inheritance* \oplus *Set_region_attr*(*region*, *Inheritance_option_copy*)

VmAllocateSecureState

Δ *PageSid*
address! : *VIRTUAL_ADDRESS*
size? : \mathbb{N}
target_task? : *TASK*
obj_sid : *OSI*

page_sid' = *page_sid*
 \oplus *Set_region_attr*(({ *target_task?* } \times *Region_of*(*address!*, *size?*)), *obj_sid*)

10.2.6 Complete Request

The general form of a **vm_allocate** request received through a task port has the following form. If a security server request is needed, then after the processing is begun the security request is made and the kernel request is marked as pending. It will later be continued by *VmContinue*. Note that *obj_sid'* is set to the default virtual memory security identifier for the target task.

ProcessingVmAllocateSignature

PortSid
MessageToVMParameters
 Δ *DtosExec*
 \exists *Mach*
 \exists *DtosAdditions*
 \exists *ValidatedRequests*
obj_sid' : *OSI*
thread' : *THREAD*

ProcessingVmAllocateNoRequest

Transition
ProcessingVmAllocateSignature

operation? = *Vm_allocate_id*
thread' = *active_thread(cpu??)*
obj_sid' = *vm_port_sid(target_task?)*
let *subject* == *thread_sid(thread')*

- *VmNoSecurityRequest[subject/ssi, obj_sid'/osi, Map_vm_region/perm]*

ProcessingVmAllocateWithRequest

Transition
ProcessingVmAllocateSignature

operation? = *Vm_allocate_id*
thread' = *active_thread(cpu??)*
obj_sid' = *vm_port_sid(target_task?)*
let *subject* == *thread_sid(thread')*

- *VmSecurityRequest[subject/ssi, obj_sid'/osi, Map_vm_region/perm]*

The general form of a **vm_allocate_secure** request received through a task port has the following form. If a security server request is needed, then after the processing is begun the security request is made and the kernel request is marked as pending. It will later be continued by *VmContinue*. Note that *obj_sid'* is set to the security identifier specified by *obj_sid?*.

ProcessingVmAllocateSecureNoRequest

Transition
ProcessingVmAllocateSignature

operation? = *Vm_allocate_secure_id*
thread' = *active_thread(cpu??)*
obj_sid' = *obj_sid?*
let *subject* == *thread_sid(thread')*

- *VmNoSecurityRequest[subject/ssi, obj_sid'/osi, Map_vm_region/perm]*

ProcessingVmAllocateSecureWithRequest

Transition
ProcessingVmAllocateSignature

operation? = *Vm_allocate_secure_id*
thread' = *active_thread(cpu??)*
obj_sid' = *obj_sid?*
let *subject* == *thread_sid(thread')*

- *VmSecurityRequest[subject/ssi, obj_sid'/osi, Map_vm_region/perm]*

A successful request makes the state changes described in the previous section and creates a kernel reply.

VmAllocateGood
 $\hat{=}$ $((RVVmAllocateSuccessful \vee RVVmAllocateVacuous)$
 $\wedge VmAllocateState \wedge VmAllocateSecureState)$
 $\gg VmAllocateReply$

An unsuccessful request returns an error status.

$$\begin{aligned} VmAllocateBad & \\ & \hat{=} (RVVmAllocateBadSecurityId \vee RVVmAllocateBadArgument \\ & \quad \vee RVVmAllocateBadAddress \vee RVVmAllocateNoSpace) \\ & \gg RequestNoOp \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

Review Note:

The component *address* is hidden so that *ExecuteVmAllocate* has a signature consistent with other requests.

$$ExecuteVmAllocate \hat{=} (VmAllocateGood \vee VmAllocateBad) \setminus (address)$$

The full specification for kernel processing of a validated **vm_allocate** or **vm_allocate_secure** request consists of processing the request, waiting until the correct information is in the access vector cache (if necessary), and then executing the request.

$$\begin{aligned} VmAllocate & \\ & \hat{=} ([VmContinue \mid operation? = Vm_allocate_id] \vee ProcessingVmAllocateNoRequest) \\ & \quad ; ExecuteVmAllocate \\ VmAllocateSecure & \\ & \hat{=} ([VmContinue \mid operation? = Vm_allocate_secure_id] \\ & \quad \vee ProcessingVmAllocateSecureNoRequest) \\ & \quad ; ExecuteVmAllocate \end{aligned}$$

10.3 vm_deallocate

The **vm_deallocate** task request deallocates a region of memory in the target task's address space.

10.3.1 Client Interface

```
kern_return_t vm_deallocate
    (mach_port_t
     vm_address_t
     vm_size_t
     target_task_name,
     address,
     size);
```

10.3.1.1 Input Parameters The following input parameters are provided by the client of a **vm_deallocate** request:

- *target_task_name?* — the client's name for the task in whose virtual address space the region is to be deallocated
- *address?* — starting address for the region

- *size?* — the number of bytes to deallocate. Any page that contains an address in the range *address? .. (address? + size? - 1)* will be deallocated.

```

VmDeallocate ClientInputs
target_task_name? : NAME
address? : VIRTUAL_ADDRESS
size? : N

```

A **vm_deallocate** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_deallocate_id* and has a body consisting of *address?* and *size?*.

```

Invoke VmDeallocate
Invoke MachMsg
VmDeallocate ClientInputs
name? = target_task_name?
operation? = Vm_deallocate_id
msg_body = Region_to_body (address?, size?)

```

10.3.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **vm_deallocate** request:

- *return!* — the status of the request

```

VmDeallocate ClientOutputs
return! : KERNEL_RETURN

```

```

VmDeallocate ReceiveReply
Invoke MachMsgRcv
VmDeallocate ClientOutputs
return! = Text_to_status (msg_body)

```

10.3.2 Kernel Interface

10.3.2.1 Input Parameters The following input parameters are provided to the kernel for a **vm_deallocate** request:

- *target_task?* — the task in whose virtual address space the region is to be deallocated
- *address?* — starting address for the region
- *size?* — the number of bytes to deallocate. Any page that contains an address in the range *address? .. (address? + size? - 1)* will be deallocated.

```

VmDeallocateInputs
target_task? : TASK
address? : VIRTUAL_ADDRESS
size? : N

```

10.3.2.2 Output Parameters The following output parameters are returned by the kernel for a **vm_deallocate** request:

- *return!* — the status of the request

<i>VmDeallocate Outputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

10.3.3 Request Criteria

No criteria are defined for the **vm_deallocate** request.

10.3.4 Return Values

Table 47 describes the values returned at the completion of the request and the conditions under which each value is returned.

Editorial Note:

As noted by CLI, the OSF KID states that *Kern_invalid_address* is returned if there are any unallocated pages in the region to be deallocated. The prototype always returns *Kern_success*.

<i>return!</i>
<i>Kern_success</i>

Table 47: Return Values for **vm_deallocate**

<i>RVVmDeallocateSuccessful</i>
<i>VmDeallocate Outputs</i>
<i>return!</i> = <i>Kern_success</i>

10.3.5 State Changes

A successful **vm_deallocate** request deallocates virtual memory. It also deletes any system attributes that are only defined for allocated memory (protections, inheritance, security identifier).

<p><i>VmDeallocateState</i></p> <p>Δ <i>AddressSpace</i></p> <p>Δ <i>Protection</i></p> <p>Δ <i>Inheritance</i></p> <p><i>address?</i> : <i>VIRTUAL_ADDRESS</i></p> <p><i>size?</i> : \mathbb{N}</p> <p><i>target_task?</i> : <i>TASK</i></p> <p>let <i>region</i> == { <i>target_task?</i> } \times <i>Region_of</i>(<i>address?</i>, <i>size?</i>)</p> <ul style="list-style-type: none"> • <i>allocated'</i> = <i>allocated</i> \ <i>region</i> \wedge <i>map_rel'</i> = <i>region</i> \triangleleft <i>map_rel</i> \wedge <i>m_protection'</i> = <i>region</i> \triangleleft <i>m_protection</i> \wedge <i>c_protection'</i> = <i>region</i> \triangleleft <i>c_protection</i> \wedge <i>i_inheritance'</i> = <i>region</i> \triangleleft <i>i_inheritance</i>
<p><i>VmDeallocateSecureState</i></p> <p>Δ <i>PageSid</i></p> <p><i>address?</i> : <i>VIRTUAL_ADDRESS</i></p> <p><i>size?</i> : \mathbb{N}</p> <p><i>target_task?</i> : <i>TASK</i></p> <p><i>page_sid'</i> = ({ <i>target_task?</i> } \times <i>Region_of</i>(<i>address?</i>, <i>size?</i>)) \triangleleft <i>page_sid</i></p>

10.3.6 Complete Request

The general form of a **vm_deallocate** request received through a task port has the following form.

<p><i>Processing VmDeallocate</i></p> <p><i>Message To VMParameters</i></p> <p><i>operation?</i> = <i>Vm_deallocate_id</i></p>
--

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$VmDeallocate\ Good \hat{=} (RVVmDeallocateSuccessful \wedge VmDeallocateState) \gg RequestReturnOnlyStatus$$

Execution of the request consists of a good execution.

$$Execute\ VmDeallocate \hat{=} VmDeallocate\ Good$$

The full specification for kernel processing of a validated **vm_deallocate** request consists of processing the request followed by its execution.

$$VmDeallocate \hat{=} Processing\ VmDeallocate \ ; \ Execute\ VmDeallocate$$

10.4 vm_inherit

The **vm_inherit** task request sets the inheritance attribute for a region within a specified task's address space.

10.4.1 Client Interface

kern_return_t vm_inherit	
(mach_port_t	<i>target_task_name,</i>
vm_address_t	<i>address,</i>
vm_size_t	<i>size,</i>
vm_inherit_t	<i>new_inheritance);</i>

10.4.1.1 Input Parameters The following input parameters are provided by the client of a **vm_inherit** request:

- *target_task_name?* — the client's name for the task in whose virtual address space the region is contained
- *address?* — starting address for the region
- *size?* — the number of bytes in the region. The inheritance attributes will be modified for any page that contains an address in the range *address? .. (address? + size? - 1)*.
- *new_inheritance?* — the new inheritance attribute for the region

```

_VmInheritClientInputs_____
target_task_name? : NAME
address? : VIRTUAL_ADDRESS
size? : N
new_inheritance? : INHERITANCE_OPTION

```

A **vm_inherit** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_inherit_id* and has a body consisting of *address?*, *size?*, and *new_inheritance?*.

```

_InvokeVmInherit_____
_InvokeMachMsg
_VmInheritClientInputs_____
name? = target_task_name?
operation? = Vm_inherit_id
msg_body = Region_inheritance_to_body(address?, size?, new_inheritance?)

```

10.4.1.2 Output Parameters The following output parameters are received through the reply provided by the client of a **vm_inherit** request:

- *return!* — the status of the request

```

_VmInheritClientOutputs_____
return! : KERNEL_RETURN

```

<i>VmInheritReceiveReply</i> <i>InvokeMachMsgRcv</i> <i>VmInheritClientOutputs</i> <i>return!</i> = <i>Text_to_status(msg_body)</i>
--

10.4.2 Kernel Interface

10.4.2.1 Input Parameters The following input parameters are provided to the kernel for a **vm_inherit** request:

- *target_task?* — the task in whose virtual address space the region is contained
- *address?* — starting address for the region
- *size?* — the number of bytes in the region. The inheritance attributes will be modified for any page that contains an address in the range *address? .. (address? + size? - 1)*.
- *new_inheritance?* — the new inheritance attribute for the region

<i>VmInheritInputs</i> <i>target_task?</i> : <i>TASK</i> <i>address?</i> : <i>VIRTUAL_ADDRESS</i> <i>size?</i> : \mathbb{N} <i>new_inheritance?</i> : <i>INHERITANCE_OPTION</i>

10.4.2.2 Output Parameters The following output parameters are returned by the kernel for a **vm_inherit** request:

- *return!* — the status of the request

<i>VmInheritOutputs</i> <i>return!</i> : <i>KERNEL_RETURN</i>
--

10.4.3 Request Criteria

The following criteria are defined for the **vm_inherit** request.

- **C1** — The value of *new_inheritance?* is valid.

<i>C1VmInheritGoodInheritance</i> <i>new_inheritance?</i> : <i>INHERITANCE_OPTION</i> <i>new_inheritance?</i> $\in \{ Inheritance_option_share, Inheritance_option_copy, Inheritance_option_none \}$
--

$NotC1VmInheritGoodInheritance \hat{=} \neg C1VmInheritGoodInheritance$

10.4.4 Return Values

Table 48 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

Although the OSF KID states that *Kern_invalid_address* is returned if the address is illegal or specifies a non-allocated region, in the prototype, *Kern_invalid_address* is never returned for this request. It appears that *Kern_success* is returned in the case of a bad address. CLI has also noted this discrepancy.

<i>return!</i>	C1
<i>Kern_success</i>	T
<i>Kern_invalid_argument</i>	F

Table 48: Return Values for **vm_inherit**

<i>RVVmInheritSuccessful</i>
<i>VmInheritOutputs</i>
<i>C1VmInheritGoodInheritance</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVVmInheritBadInheritance</i>
<i>VmInheritOutputs</i>
<i>NotC1VmInheritGoodInheritance</i>
<i>return!</i> = <i>Kern_invalid_argument</i>

10.4.5 State Changes

A successful **vm_inherit** sets the inheritance attribute for the region defined by *address?* and *size?* to the value specified by *new_inheritance?*.

<i>VmInheritState</i>
Δ <i>Inheritance</i>
Ξ <i>AddressSpace</i>
<i>target_task?</i> : <i>TASK</i>
<i>address?</i> : <i>VIRTUAL_ADDRESS</i>
<i>size?</i> : \mathbb{N}
<i>new_inheritance?</i> : <i>INHERITANCE_OPTION</i>
$\text{let } \textit{region} == \{ \textit{page_index} : \textit{PAGE_INDEX}$ $\quad \textit{page_index} \in \textit{Region_of}(\textit{address?}, \textit{size?})$ $\quad \wedge (\textit{target_task?}, \textit{page_index}) \in \underline{\textit{allocated}}$ $\quad \bullet (\textit{target_task?}, \textit{page_index}) \}$ $\bullet \underline{\textit{inheritance}}' = \underline{\textit{inheritance}} \oplus \textit{Set_region_attr}(\textit{region}, \textit{new_inheritance?})$

10.4.6 Complete Request

The general form of a **vm_inherit** request received through a task port has the following form.

<i>Processing VmInherit</i>
<i>Message To VM Parameters</i>
<i>operation? = Vm_inherit_id</i>

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$VmInheritGood \hat{=} (RVVmInheritSuccessful \wedge VmInheritState) \\ \gg RequestReturnOnlyStatus$$

An unsuccessful request returns an error status.

$$VmInheritBad \\ \hat{=} RVVmInheritBadInheritance \gg RequestNoOp$$

Execution of the request consists of a good execution or an error execution.

$$Execute VmInherit \hat{=} VmInheritGood \vee VmInheritBad$$

The full specification for kernel processing of a validated **vm_inherit** request consists of processing the request followed by its execution.

$$VmInherit \hat{=} Processing VmInherit ; Execute VmInherit$$

10.5 vm_protect

The **vm_protect** task request sets the current and/or maximum protections for a region within a specified task's address space. If the parameter *set_maximum?* is *False*, only the current protections are set. If *set_maximum?* is *True*, the maximum protections are set, and the current protections are also set so that they do not exceed the new maximum. Note that this request cannot be used to increase the maximum protections but only to decrease them.

10.5.1 Client Interface

```
kern_return_t vm_protect
(mach_port_t
vm_address_t
vm_size_t
boolean_t
vm_prot_t
target_task_name,
address,
size,
set_maximum,
new_protection);
```

10.5.1.1 Input Parameters The following input parameters are provided by the client of a **vm_protect** request:

- *target_task_name?* — the client's name for the task in whose virtual address space the region is contained
- *address?* — starting address for the region
- *size?* — the number of bytes in the region. The protections will be modified for any page that contains an address in the range $address? \dots (address? + size? - 1)$.
- *set_maximum?* — a Boolean indicating whether the maximum protection should be set. A value of *True* indicates the maximum protection should be set. (The current protection is also set if it violates the new maximum.) A value of *False* indicates only the current protection is set.
- *new_protection?* — the new protection for the region

```

VmProtectClientInputs
target_task_name? : NAME
address? : VIRTUAL_ADDRESS
size? : N
set_maximum? : BOOLEAN
new_protection? : P PROTECTION

```

A **vm_protect** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_protect_id* and has a body consisting of *address?*, *size?*, *set_maximum?*, and *new_protection?*.

```

Invoke VmProtect
Invoke MachMsg
VmProtectClientInputs
name? = target_task_name?
operation? = Vm_protect_id
msg_body
    = Region_bool_prot_to_body( address?, size?, set_maximum?, new_protection? )

```

10.5.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **vm_protect** request:

- *return!* — the status of the request

```

VmProtectClientOutputs
return! : KERNEL_RETURN

```

```

VmProtectReceiveReply
Invoke MachMsgRcv
VmProtectClientOutputs
return! = Text_to_status(msg_body)

```

10.5.2 Kernel Interface

10.5.2.1 Input Parameters The following input parameters are provided to the kernel for a **vm_protect** request:

- *target_task?* — the task in whose virtual address space the region is contained
- *address?* — starting address for the region
- *size?* — the number of bytes in the region. The protections will be modified for any page that contains an address in the range $address? \dots (address? + size? - 1)$.
- *set_maximum?* — a Boolean indicating whether the maximum protection should be set. A value of *True* indicates the maximum protection should be set. (The current protection is also set if it violates the new maximum.) A value of *False* indicates only the current protection is set.
- *new_protection?* — the new protection for the region

<i>VmProtectInputs</i>
<i>target_task?</i> : <i>TASK</i>
<i>address?</i> : <i>VIRTUAL_ADDRESS</i>
<i>size?</i> : \mathbb{N}
<i>set_maximum?</i> : <i>BOOLEAN</i>
<i>new_protection?</i> : \mathbb{P} <i>PROTECTION</i>

10.5.2.2 Output Parameters The following output parameters are returned by the kernel for a **vm_protect** request:

- *return!* — the status of the request

<i>VmProtectOutputs</i>
<i>return!</i> : <i>KERNEL_RETURN</i>

10.5.3 Request Criteria

The following criteria are defined for the **vm_protect** request.

- **C1** — The new protection is less than the existing maximum protection.

<i>C1VmProtectGoodProtection</i>
<i>Protection</i>
<i>target_task?</i> : <i>TASK</i>
<i>address?</i> : <i>VIRTUAL_ADDRESS</i>
<i>size?</i> : \mathbb{N}
<i>new_protection?</i> : \mathbb{P} <i>PROTECTION</i>
\forall <i>page_index</i> : <i>PAGE_INDEX</i>
<i>page_index</i> \in <i>Region_of</i> (<i>address?</i> , <i>size?</i>)
• (<i>target_task?</i> , <i>page_index</i>) \in dom <i>m_protection</i>
\wedge <i>new_protection?</i> \subseteq <i>m_protection</i> (<i>target_task?</i> , <i>page_index</i>)

$$\text{NotC1VmProtectGoodProtection} \hat{=} \text{Protection} \wedge \neg \text{C1VmProtectGoodProtection}$$

10.5.4 Return Values

Table 49 describes the values returned at the completion of the request and the conditions under which each value is returned.

Review Note:

Although the OSF KID states that *Kern_invalid_address* is returned if the address is illegal or specifies a non-allocated region, in the prototype, *Kern_invalid_address* is never returned for this request. It appears that *Kern_success* is returned in the case of an unallocated page. CLI has also noted this discrepancy.

<i>return!</i>	C1
<i>Kern_success</i>	T
<i>Kern_protection_failure</i>	F

Table 49: Return Values for **vm_protect**

<i>RVVmProtectSuccessful</i>
<i>VmProtectOutputs</i>
<i>C1VmProtectGoodProtection</i>
<i>return!</i> = <i>Kern_success</i>

<i>RVVmProtectBadProtection</i>
<i>VmProtectOutputs</i>
<i>NotC1VmProtectGoodProtection</i>
<i>return!</i> = <i>Kern_protection_failure</i>

10.5.5 State Changes

A successful **vm_protect** sets either the maximum or the current memory protection (read, write, and/or execute) allowed for the region, depending on whether *set_maximum?* is *True* or *False*. If the maximum is set below the current protection, the current protection must also be adjusted to remove any permissions that are not within the new maximum.

VmProtectState

Δ *Protection*
 Ξ *AddressSpace*
target_task? : *TASK*
address? : *VIRTUAL_ADDRESS*
size? : \mathbb{N}
set_maximum? : *BOOLEAN*
new_protection? : \mathbb{P} *PROTECTION*

let *region* == { *page_index* : *PAGE_INDEX*
| *page_index* \in *Region_of*(*address?*, *size?*)
 \wedge (*target_task?*, *page_index*) \in *allocated*
• (*target_task?*, *page_index*) }
• (*set_maximum?* = *True*
 \wedge *m_protection'* = *m_protection* \oplus *Set_region_attr*(*region*, *new_protection?*)
 \wedge *c_protection'* = *c_protection*
 \oplus { *task_va_pair* : *TASK* \times *PAGE_INDEX* | *task_va_pair* \in *region*
• *task_va_pair* \mapsto *c_protection*(*task_va_pair*) \cap *new_protection?* })
 \vee (*set_maximum?* = *False*
 \wedge *c_protection'* = *c_protection* \oplus *Set_region_attr*(*region*, *new_protection?*))

Review Note:

This ignores the protections coming from the security server. Right now these protections are not in the model of the state.

10.5.6 Complete Request

The general form of a **vm_protect** request received through a task port has the following form.

Processing VmProtect

Message To VM Parameters

operation? = *Vm_protect_id*

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$VmProtectGood \hat{=} (RVVmProtectSuccessful \wedge VmProtectState) \\ \gg RequestReturnOnlyStatus$$

An unsuccessful request returns an error status.

$$VmProtectBad \\ \hat{=} RVVmProtectBadProtection \gg RequestNoOp$$

Execution of the request consists of a good execution or an error execution.

$$Execute VmProtect \hat{=} VmProtectGood \vee VmProtectBad$$

The full specification for kernel processing of a validated **vm_protect** request consists of processing the request followed by its execution.

$$VmProtect \hat{=} Processing VmProtect \ ; \ Execute VmProtect$$

10.6 vm_write

The **vm_write** task request writes an allocated region in the target task's address space.

10.6.1 Client Interface

```
kern_return_t vm_write
    (mach_port_t          target_task_name,
     vm_address_t        address,
     vm_offset_t         data,
     mach_msg_type_number_t data_count);
```

10.6.1.1 Input Parameters The following input parameters are provided by the client of a **vm_write** request:

- *target_task_name?* — the client's name for the task in whose virtual address space the region is to be written
- *address?* — starting address for the destination region, which must be the start of a page boundary
- *data?* — the data to be written

Editorial Note:

In the DTOS KID, this parameter is described as a page-aligned array of data. However, in the prototype the *data?* parameter is a pointer to a `vm_map_copy` structure which encodes information about the source region to copy including its offset, size, a type and a memory map. The type describes how this structure represents the data. The three possibilities are an entry list, an object and a page list (only entry lists are currently supported by the prototype). This structure is returned by a **vm_read** request. We will model this structure as a *MapCopy* containing the offset, the size and the task from whose address space the copy was made.

- *data_count?* — ignored

Editorial Note:

In the DTOS KID this parameter denotes the number of bytes in the array pointed to by the *data?* parameter. However, the number of bytes is included in the `vm_map_copy` structure, and the *data_count?* parameter is ignored in the prototype.

Vm Write ClientInputs

```
target_task_name? : NAME
address? : VIRTUAL_ADDRESS
data? : MapCopy
data_count? : N
```

A **vm_write** request is invoked by sending a message to the port indicated by *target_task_name?* that has the operation field set to *Vm_write_id* and has a body consisting of *address?*, *data?*, and *data_count?*.

Invoke Vm Write
Invoke MachMsg
Vm Write ClientInputs

name? = *target_task_name?*
operation? = *Vm_write_id*
msg_body = *Address_data_to_body*(*address?*, *data?*, *data_count?*)

10.6.1.2 Output Parameters The following output parameters are received through the reply port provided by the client of a **vm_write** request:

- *return!* — the status of the request

Vm Write ClientOutputs
return!: *KERNEL_RETURN*

Vm Write Receive Reply
Invoke MachMsgRcv
Vm Write ClientOutputs

return! = *Text_to_status*(*msg_body*)

10.6.2 Kernel Interface

10.6.2.1 Input Parameters The following input parameters are provided to the kernel for a **vm_write** request:

- *target_task?* — the task in whose virtual address space the region is to be written
- *address?* — starting address for the destination region, which must be the start of a page boundary
- *data?* — the data to be written

Editorial Note:

In the DTOS KID, this parameter is described as a page-aligned array of data. However, in the prototype the *data?* parameter is a pointer to a *vm_map_copy* structure which encodes information about the source region to copy including its offset, size, a type and a memory map. The type describes how this structure represents the data. The three possibilities are an entry list, an object and a page list (only entry lists are currently supported by the prototype). This structure is returned by a **vm_read** request. We will model this structure as a *MapCopy* containing the offset, the size and the task from whose address space the copy was made.

- *data_count?* — ignored

Editorial Note:

In the DTOS KID this parameter denotes the number of bytes in the array pointed to by the *data?* parameter. However, the number of bytes is included in the *vm_map_copy* structure, and the *data_count?* parameter is ignored in the prototype.

VmWriteInputs

```
target_task? : TASK
address? : VIRTUAL_ADDRESS
data? : MapCopy
data_count? : N
```

10.6.2.2 Output Parameters The following output parameters are returned by the kernel for a **vm_write** request:

- *return!* — the status of the request

VmWriteOutputs

```
return! : KERNEL_RETURN
```

10.6.3 Request Criteria

The following criteria are defined for the **vm_write** request.

- **C1** — The parameter *address?* and the offset included in the parameter *data?* are on a page boundary. Also, the size included in *data?* is an integer number of pages.

C1VmWritePageAligned

```
address? : VIRTUAL_ADDRESS
data? : MapCopy
```

$$\{address?, data?.offset, Address_num \sim (data?.size)\} \subseteq Page_aligned$$

$NotC1VmWritePageAligned \hat{=} \neg C1VmWritePageAligned$

- **C2** — The addresses specified for the destination region are valid and are allocated.

C2VmWriteGoodAddress

```
AddressSpace
target_task? : TASK
address? : VIRTUAL_ADDRESS
data? : MapCopy
```

```
let data_size == data?.size
• VmGoodRegion[address?/address, data_size/size]
  ∧ Region_of(address?, data_size) ⊆ allocated({ target_task? })
```

$NotC2VmWriteGoodAddress \hat{=} AddressSpace \wedge \neg C2VmWriteGoodAddress$

- **C3** — The target task has permission to write to the region.

Review Note:
I believe the security server should be queried to make sure the target task still has write permission. However, I don't think the prototype currently makes this check (11/15/94).

<p><i>C3 Vm Write Writable</i></p> <p><i>Protection</i></p> <p><i>target_task?</i> : <i>TASK</i></p> <p><i>address?</i> : <i>VIRTUAL_ADDRESS</i></p> <p><i>data?</i> : <i>MapCopy</i></p>
<p>\forall <i>page</i> : <i>PAGE_INDEX</i></p> <p> <i>page</i> \in <i>Region_of</i>(<i>address?</i>, <i>data?.size</i>)</p> <p> \wedge (<i>target_task?</i>, <i>page</i>) \in <i>dom c_protection</i></p> <p>• <i>Write</i> \in <i>c_protection</i>(<i>target_task?</i>, <i>page</i>)</p>

$$\text{Not}C3\ Vm\ Write\ Writable \hat{=} Protection \wedge \neg C3\ Vm\ Write\ Writable$$

Review Note:

Should we state that read permission is required on the source region? The prototype does require this, but I am not certain what protections there will be on the map copy object.

10.6.4 Return Values

Table 50 describes the values returned at the completion of the request and the conditions under which each value is returned.

<i>return!</i>	C1	C2	C3
<i>Kern_success</i>	T	T	T
<i>Kern_protection_failure</i>	T	T	F
<i>Kern_invalid_address</i>	T	F	-
<i>Kern_invalid_argument</i>	F	-	-

Table 50: Return Values for *vm_write*

<p><i>RV Vm Write Successful</i></p> <p><i>Vm Write Outputs</i></p> <p><i>C1 Vm Write Page Aligned</i></p> <p><i>C2 Vm Write Good Address</i></p> <p><i>C3 Vm Write Writable</i></p>
<p><i>return!</i> = <i>Kern_success</i></p>

<p><i>RV Vm Write Protect Fail</i></p> <p><i>Vm Write Outputs</i></p> <p><i>C1 Vm Write Page Aligned</i></p> <p><i>C2 Vm Write Good Address</i></p> <p><i>Not C3 Vm Write Writable</i></p>
<p><i>return!</i> = <i>Kern_protection_failure</i></p>

RVVmWriteBadAddress

VmWriteOutputs
C1 VmWritePageAligned
NotC2 VmWriteGoodAddress

return! = *Kern_invalid_address*

RVVmWriteInvalidArg

VmWriteOutputs
NotC1 VmWritePageAligned

return! = *Kern_invalid_argument*

Review Note:

The prototype also checks (after C1) whether the size of the *MapCopy* is zero. If so, it returns *Kern_success* without checking conditions C2 and C3. No changes are made to the state. Since, when the size is zero, conditions C2 and C3 are automatically true and no words are changed below, this circumstance is covered by the *Kern_success* case, and we do not state the extra criterion.

10.6.5 State Changes

A successful **vm_write** request writes the data to the memory pages associated with the specified area of virtual memory. The data written into the address space of the target task originated from the address space of some task (e.g., it was read from that address space using **vm_read**).

VmWriteState

Δ *VaWord*
target_task? : *TASK*
address? : *VIRTUAL_ADDRESS*
data? : *MapCopy*

$\forall x : 0 \dots data?.size - 1$
 $| (address?, x) \in \text{dom } Relative_addr$
 $\bullet va_word'(target_task?, Relative_addr(address?, x))$
 $= va_word(data?.task, Relative_addr(data?.offset, x))$

Review Note:

It would be better to model memory maps explicitly (independent of tasks) in the state description instead of just associating virtual addresses with tasks. This would allow map copies to be modeled as a special memory map that has no directly associated task.

10.6.6 Complete Request

The general form of a **vm_write** request received through a task port has the following form.

Processing VmWrite

MessageToVMParameters

operation? = *Vm_write_id*

A successful request makes the state changes described in the previous section and creates a kernel reply.

$$\begin{aligned} VmWriteGood &\hat{=} (RVVmWriteSuccessful \wedge VmWriteState) \\ &\gg RequestReturnOnlyStatus \end{aligned}$$

An unsuccessful request returns an error status.

$$\begin{aligned} VmWriteBad & \\ &\hat{=} (RVVmWriteInvalidArg \vee RVVmWriteBadAddress \vee RVVmWriteProtectFail) \\ &\gg RequestNoOp \end{aligned}$$

Execution of the request consists of a good execution or an error execution.

$$ExecuteVmWrite \hat{=} VmWriteGood \vee VmWriteBad$$

The full specification for kernel processing of a validated **vm_write** request consists of processing the request followed by its execution.

$$VmWrite \hat{=} ProcessingVmWrite ; ExecuteVmWrite$$

Review Note:

No interaction with memory managers for the region being written is specified. Any pages not backed by Null_memory must not be locked against writing, but we only have locking information for cached segments of the memory objects.

Section **11**
Notes

11.1 Acronyms

CCA Covert Channel Analysis

CMU Carnegie Mellon University

DTOS Distributed Trusted Operating System

FSPM Formal Security Policy Model

IPC Interprocess Communication

KID Kernel Interface Document

MLS Multi-Level Secure

OSC Object Security Context

OSF Open Software Foundation

OSI Object Security Identifier

SID Security Identifier

SSC Subject Security Context

SSI Subject Security Identifier

VM Virtual Memory

11.2 Glossary

dirty page A page in kernel memory is dirty if the pager associated with the page has not yet been made aware of modifications that have been made to the page.

permission A permission is an access mode enforced by the kernel. The kernel ensures that a service is provided only when the client of the service has the appropriate permission.

precious page A page in kernel memory is precious if the pager associated with the page has indicated that it is not maintaining a copy of the page. Regardless of whether the page is dirty, the kernel must send the contents of the page to the pager before removing the page from memory.

security server A security server is a user space task that provides access computations to the kernel.

Appendix **A**
Bibliography

- [1] William R. Bevier and Lawrence M. Smith. A Mathematical Description of the Mach Kernel: Virtual Memory Services (Draft). Technical report, Computational Logic, Incorporated, August 1993.
- [2] William R. Bevier and Lawrence M. Smith. A Mathematical Model of the Mach Kernel: Entities and Relations (Draft). Technical report, Computational Logic, Incorporated, April 1993.
- [3] William R. Bevier and Lawrence M. Smith. A Mathematical Model of the Mach Kernel: Port Services (Draft). Technical report, Computational Logic, Incorporated, August 1993.
- [4] William R. Bevier and Lawrence M. Smith. A Mathematical Model of the Mach Kernel: Task and Thread Services (Draft). Technical report, Computational Logic, Incorporated, August 1993.
- [5] Todd Fine, Carol Muehrcke, and Edward A. Schneider. Formal Top Level Specification for Distributed Trusted Mach. Technical report, Secure Computing Corporation, 2675 Long Lake Road, Roseville, Minnesota 55113-2536, April 1993. DTMach CDRL A012.
- [6] Keith Loepere. *Mach 3 Kernel Interfaces* Open Software Foundation and Carnegie Mellon University, November 1992.
- [7] Keith Loepere. *OSF Mach Kernel Principles* Open Software Foundation and Carnegie Mellon University, final draft edition, May 1993.
- [8] Secure Computing Corporation. DTOS Kernel Interfaces Document. Technical report, Secure Computing Corporation, 2675 Long Lake Road, Roseville, Minnesota 55113-2536, April 1995. DTOS CDRL A003.
- [9] Secure Computing Corporation. DTOS Formal Security Policy Model (FSPM). Technical report, Secure Computing Corporation, 2675 Long Lake Road, Roseville, Minnesota 55113-2536, September 1996. DTOS CDRL A004.
- [10] Secure Computing Corporation. DTOS Lessons Learned Report. Technical report, Secure Computing Corporation, 2675 Long Lake Road, Roseville, Minnesota 55113-2536, October 1996. DTOS CDRL A008.

Appendix *B*

Z Extensions

This section describes “extensions” to the Z specification language that are used in the DTOS FTLS. All of these extensions are defined in terms of constructs in the Z specification language, so they are not technically extensions to the language.

B.1 Disjointness and Partitions

It is often necessary to indicate that each element of a collection of values is unique. For example, consider specifying that val_1, \dots, val_n are unique values. Since n might be relatively large, it is undesirable to enumerate each pair:

$$val_1 \neq val_2 \wedge val_1 \neq val_3 \wedge val_1 \neq val_4 \dots$$

Although `disjoint` is part of the Z mathematical toolkit, it addresses disjointness of sets instead of disjointness of values. While we could convert values to singleton sets of values as follows:

$$\text{disjoint } \langle \{ val_1 \}, \dots, \{ val_n \} \rangle$$

this is somewhat inconvenient. Another possibility would be to specify that:

$$\langle val_1, \dots, val_n \rangle$$

is, when viewed as a function, injective. However, the expression:

$$\langle val_1, \dots, val_n \rangle \in \mathbb{N} \mapsto X$$

is a rather unintuitive way to express disjointness.

Instead, the generic predicate *Values_disjoint* is defined to state such disjointness properties. The expression *Values_disjoint* $\langle val_1, \dots, val_n \rangle$ denotes that val_1, \dots, val_n are unique values.

Mach Definition 109

$\begin{aligned} & \text{Values_disjoint } _ : \mathbb{P}(\text{seq } X) \\ & \forall val_seq : \text{seq } X \\ & \bullet \text{Values_disjoint } val_seq \\ & \Leftrightarrow (\forall i_1, i_2 : \mathbb{N} \mid i_1 \in \text{dom } val_seq \wedge i_2 \in \text{dom } val_seq \wedge i_1 \neq i_2 \\ & \quad \bullet val_seq(i_1) \neq val_seq(i_2)) \end{aligned}$

Similarly, the expression $\langle val_1, \dots, val_n \rangle$ *Values_partition* S denotes that the values val_1, \dots, val_n are unique values that together comprise the set val_set .

Mach Definition 110

$\begin{aligned} & \text{-- Values_partition_} : (\text{seq } X) \leftrightarrow \mathbb{P} X \\ & \forall \text{val_seq} : \text{seq } X; \text{val_set} : \mathbb{P} X \\ & \bullet \text{val_seq Values_partition val_set} \\ & \Leftrightarrow (\text{Values_disjoint val_seq} \wedge \text{val_set} = \text{ran val_seq}) \end{aligned}$
--

B.2 Partial Orders

A partial ordering is a relation that is *reflexive*, *antisymmetric*, and *transitive*.

A reflexive relation is one that relates each element to itself; in other words, the identity relation is contained in every reflexive relation.

An antisymmetric relation is a relation containing no cycles of the form $(val_1, val_2) \in R \wedge (val_2, val_1) \in R$ for distinct val_1 and val_2 . Since $(val_2, val_1) \in R$ is equivalent to $(val_1, val_2) \in R^\sim$, a relation is antisymmetric exactly when $(val_1, val_2) \in R \wedge (val_1, val_2) \in R^\sim$ only holds for $val_1 = val_2$. In other words, a relation is antisymmetric when its intersection with its inverse is contained in *id*.

A relation is transitive when:

$$(val_1, val_2) \in R \wedge (val_2, val_3) \in R \Rightarrow (val_1, val_3) \in R$$

In other words, whenever it is possible to get from val_1 to val_3 through repeated iteration of R , R relates val_1 to val_3 directly. This is equivalent to R^2 being contained in R . For each type X , the following sets of relations are defined:

- *Reflexive*[X] — the set of all reflexive relations on X
- *Anti_symmetric*[X] — the set of all antisymmetric relations on X
- *Transitive*[X] — the set of all transitive relations on X
- *Poset*[X] — the set of all relations on X that are posets; this is simply the intersection of *Reflexive*[X], *Anti_symmetric*[X], and *Transitive*[X]

Mach Definition 111

$\begin{aligned} & \text{Poset} : \mathbb{P}(X \leftrightarrow X) \\ & \text{Reflexive} : \mathbb{P}(X \leftrightarrow X) \\ & \text{Anti_symmetric} : \mathbb{P}(X \leftrightarrow X) \\ & \text{Transitive} : \mathbb{P}(X \leftrightarrow X) \end{aligned}$
$\begin{aligned} & \text{Poset} = \text{Reflexive} \cap \text{Anti_symmetric} \cap \text{Transitive} \\ & \text{Reflexive} = \{R : X \leftrightarrow X \mid \text{id } X \subseteq R\} \\ & \text{Anti_symmetric} = \{R : X \leftrightarrow X \mid R \cap R^\sim \subseteq \text{id } X\} \\ & \text{Transitive} = \{R : X \leftrightarrow X \mid R^2 \subseteq R\} \end{aligned}$

B.3 Sequences

The expression $val_seq\ Add_value\ val$ is used to denote the sequence resulting from adding the element val to the end of the sequence val_seq . The expression $s\ Wrap_value\ val$ is used to denote the sequence resulting from replacing the first element of val_seq with val .

Mach Definition 112

$\begin{array}{l} \text{---}[X] \text{---} \\ \text{---} \\ _Add_value\ _ : (\text{seq } X) \times X \rightarrow (\text{seq } X) \\ _Wrap_value\ _ : (\text{seq } X) \times X \rightarrow (\text{seq } X) \\ \text{---} \\ \forall val_seq : \text{seq } X; val : X \\ \bullet\ val_seq\ Add_value\ val = val_seq \hat{\ } \{1 \mapsto val\} \\ \wedge (\#val_seq > 0 \Rightarrow val_seq\ Wrap_value\ val = val_seq \oplus \{1 \mapsto val\}) \end{array}$

The expression $Seq_plus(S)$ where S is a sequence of numbers returns the sum of the numbers in S .

Mach Definition 113

$\begin{array}{l} Seq_plus : \text{seq } \mathbb{Z} \rightarrow \mathbb{Z} \\ \text{---} \\ Seq_plus(\langle \rangle) = 0 \\ \forall S : \text{seq}_1 \mathbb{Z} \\ \bullet\ Seq_plus(S) = head(S) + Seq_plus(tail(S)) \end{array}$

Appendix **C** **IPC**

C.1 IPC Requests

This section describes the **mach_msg** request.

Review Note:

This section has not yet been updated for DTOS. Currently, this section is a direct copy of the corresponding DTMach section with minor changes required for DTOS sections that depend on this section.

C.1.1 Constants and Types

We use the following type to denote **mach_msg** return codes:

[*MACH_MSG_RETURN*]

The return values defined in Mach are:

```
Mm_no_op : MACH_MSG_RETURN
Mm_send_msg_too_small : MACH_MSG_RETURN
Mm_send_no_buffer : MACH_MSG_RETURN
Mm_send_invalid_header : MACH_MSG_RETURN
Mm_send_invalid_dest : MACH_MSG_RETURN
Mm_send_invalid_reply : MACH_MSG_RETURN
Mm_send_invalid_notify : MACH_MSG_RETURN
Mm_rcv_invalid_notify : MACH_MSG_RETURN
Mm_rcv_invalid_name : MACH_MSG_RETURN
Mm_rcv_in_set : MACH_MSG_RETURN
Mm_rcv_timed_out : MACH_MSG_RETURN
Mm_rcv_too_large : MACH_MSG_RETURN
Mm_send_will_notify : MACH_MSG_RETURN
Mm_success : MACH_MSG_RETURN
Mm_send_invalid_right : MACH_MSG_RETURN
Mm_send_invalid_memory : MACH_MSG_RETURN
Mm_send_invalid_type : MACH_MSG_RETURN
Mm_rcv_port_died : MACH_MSG_RETURN
Mm_rcv_port_changed : MACH_MSG_RETURN

Values_disjoint(Mm_no_op, Mm_send_msg_too_small, Mm_send_invalid_header,
Mm_send_invalid_dest, Mm_send_invalid_reply, Mm_send_invalid_notify,
Mm_rcv_invalid_name, Mm_rcv_in_set, Mm_rcv_timed_out,
Mm_rcv_too_large, Mm_send_will_notify, Mm_success,
Mm_send_invalid_right, Mm_send_invalid_memory, Mm_send_invalid_type,
Mm_rcv_port_died, Mm_rcv_port_changed)
```

C.2 mach_msg

Review Note:

This section has not yet been updated for DTOS. Currently, this section is a direct copy of the corresponding DTMach section with minor changes required for DTOS sections that depend on this section.

The request **mach_msg** allows a thread to send and receive messages.¹⁴

The request has the following input parameters:

- *client?* — the thread sending or receiving a message
- *msg?* — the message header; note that this is only relevant when a message is being sent
- *option?* — message options
- *send_size?* — specifies the size of *msg?* when a message is being sent
- *rcv_size?* — specifies the size of *msg?* when a message is being received
- *rcv_name?* — specifies the port or port set from which to receive a message when a message is being received
- *time_out?* — specifies the amount of time to wait for the operation to complete before giving up
- *notify?* — specifies the notification port to use in the case in which notifications are requested
- *msg_body?* — the message body; note that this is only relevant when a message is being sent

The request has the following output parameters:

- *msg!* — the message buffer; note that this is only relevant when a message is being received
- *rcv_size!* — specifies the size of the message when an attempt is made to receive a message that is too large
- *msg_body!* — the in-line data portion of the message; note that this is only an output in the case when a message is being received
- *msg_return!* — the status of the request

The request is initiated by a schema of the following form.

¹⁴The specification of this request is incomplete. See Section C.2.3 for a description of the work that remains to be done.

<p><i>MachMsgSignature</i></p> <p>Δ <i>DtosExec</i></p> <p><i>client?</i> : <i>THREAD</i></p> <p><i>msg</i> : <i>MachMsgHeader</i></p> <p><i>option?</i> : \mathbb{P} <i>MACH_MSG_OPTION</i></p> <p><i>send_size?</i> : \mathbb{N}</p> <p><i>rcv_size?</i> : \mathbb{N}</p> <p><i>rcv_size!</i> : \mathbb{N}</p> <p><i>rcv_name?</i> : <i>NAME</i></p> <p><i>time_out?</i> : \mathbb{N}</p> <p><i>notify?</i> : <i>NAME</i></p> <p><i>msg_body</i> : <i>MESSAGE_BODY</i></p> <p><i>msg_return!</i> : <i>MACH_MSG_RETURN</i></p> <hr/> <p><i>client?</i> \in <i>dom_owning_task</i></p> <p><i>owning_task client?</i> \in <i>task_exists</i></p>
--

where the meaning of the parameters is as described earlier.¹⁵

If *option?* includes neither *Mach_send_msg* nor *Mach_rcv_msg*, then no processing occurs.¹⁶

<p><i>MachMsgNoOp</i></p> <p>\exists <i>DtosExec</i></p> <p><i>MachMsgSignature</i></p> <hr/> <p>$\{ \textit{Mach_send_msg}, \textit{Mach_rcv_msg} \} \cap \textit{option?} = \emptyset$</p> <p><i>msg_return!</i> = <i>Mm_no_op</i></p>
--

C.2.1 Message Send

The **mach_msg** request can be used to send a message by including *Mach_send_msg* in *option?* and not including *Mach_rcv_msg*.

<p><i>MachMsgSend</i></p> <p><i>MachMsgSignature</i></p> <hr/> <p><i>Mach_send_msg</i> \in <i>option?</i></p> <p><i>Mach_rcv_msg</i> \notin <i>option?</i></p>
--

There are four general cases to consider:

- An error condition occurs during the initial processing, and the request is a no-op.
- A subsequent error condition occurs and the message is returned through a pseudo-receive operation.
- A subsequent error condition occurs and some of the message is lost during delivery.

¹⁵Note that the *msg?* and *msg!* parameters are both represented by *msg*. Similarly, *msg_body?* and *msg_body!* are both represented by *msg_body*.

¹⁶The DTMach Kernel Interface does not define a return status for this case. We have introduced the return status *Mm_no_op* to denote the return status in this case.

- The message is successfully delivered.

The first case is discussed in Section C.2.1.1. The remaining cases are described in Section C.2.1.2.

C.2.1.1 Initial Processing We use the following schema to describe send operations that are processed as no-ops due to error conditions that arise during the initial processing of the request:

$\begin{array}{l} \text{MachMsgSendNoOp} \\ \exists Dtos \\ \text{MachMsgSend} \end{array}$

If *send_size?* is too small, then an error message is returned and no further processing occurs. We define the following constant to denote the minimum send size.

$\text{Min_send_size} : \mathbb{N}$

The case in which the message is too small is specified as follows:

$\begin{array}{l} \text{MachMsgSendMsgTooSmall} \\ \text{MachMsgSendNoOp} \\ \text{send_size?} < \text{Min_send_size} \\ \text{msg_return!} = \text{Mm_send_msg_too_small} \end{array}$

If there is not enough memory available for the kernel to process the request, then an error message is returned and no further processing occurs. We use the following predicate to indicate when there is insufficient memory available:

$\text{cannot_allocate_send_buffer} : \mathbb{P} \text{ Mach}$

The specification of the processing is as follows:¹⁷

$\begin{array}{l} \text{MachMsgSendSizeOk} \\ \text{send_size?} : \mathbb{N} \\ \text{send_size?} \geq \text{Min_send_size} \end{array}$
--

$\begin{array}{l} \text{MachMsgSendNoBuffer} \\ \text{MachMsgSendNoOp} \\ \text{MachMsgSendSizeOk} \\ \text{cannot_allocate_send_buffer}(\theta \text{ Mach}) \\ \text{msg_return!} = \text{Mm_send_no_buffer} \end{array}$
--

¹⁷For convenience, we define a schema representing the negation of the earlier tests before defining the schema representing the processing for a given case. For example, the schema *MachMsgSendSizeOk* is the negation of the previously described test of whether *send_size?* is too small. The schema *MachMsgSendNoBuffer* uses *MachMsgSendSizeOk* to define the processing for the case in which there is insufficient memory available to process the request.

If the rights specified for the local or remote port are invalid, then an error message is returned and no further processing occurs. The rights specified for the remote port are valid only if they provide the receiver with either a *Send* or *Send_once* right. Thus, the rights are invalid if they do not include any of *Mmt_make_send*, *Mmt_copy_send*, *Mmt_move_send*, *Mmt_make_send_once*, and *Mmt_move_send_once*. We define the set *TRANSFER_SEND_RIGHTS* to denote this set of values of type *MACH_MSG_TYPE*.

$$\text{TRANSFER_SEND_RIGHTS} == \{ \text{Mmt_make_send}, \text{Mmt_copy_send}, \text{Mmt_move_send}, \\ \text{Mmt_make_send_once}, \text{Mmt_move_send_once} \}$$

The remote port rights are valid exactly when they contain an element of this set. Similarly, the local port rights are valid when they contain an element of *TRANSFER_SEND_RIGHTS*. In addition, the local port rights are also valid when they are empty and the local port is null.

The specification for the case in which either the remote or local port rights are invalid is as follows:

<p><i>MachMsgSendCanAllocateBuffer</i> _____</p> <p><i>Mach</i></p> <p><i>MachMsgSendSizeOk</i></p> <p>$\neg \text{cannot_allocate_send_buffer}(\theta \text{Mach})$</p>
<p><i>MachMsgSendInvalidHeader</i> _____</p> <p><i>MachMsgSendCanAllocateBuffer</i></p> <p><i>MachMsgSendNoOp</i></p> <p>$(\text{msg}h.\text{remote_rights} \notin \text{TRANSFER_SEND_RIGHTS} \vee \\ \text{TRANSFER_SEND_RIGHTS} \cap \text{msg}h.\text{local_rights} = \emptyset \wedge \\ (\text{msg}h.\text{local_rights} \neq \emptyset \vee \\ \text{msg}h.\text{local_port} \neq \text{Mach_port_null}) \vee \\ \neg \text{msg}h.\text{remote_rights} \in \text{Recognized_transfer_options} \vee \\ \neg \text{msg}h.\text{local_rights} \subseteq \text{Recognized_transfer_options}) \\ \text{msg_return!} = \text{Mm_send_invalid_header}$</p>

Otherwise, if the client task does not have the right required by *msg.h.remote_rights*, then an error message is returned and no further processing occurs. We use the following function to denote the right required for each type of transfer:

<p><i>Required_right</i> : <i>Recognized_transfer_options</i> \rightarrow <i>RIGHT</i></p> <p><i>Required_right</i> =</p> <p>{ <i>Mmt_make_send</i> \mapsto <i>Receive</i>, <i>Mmt_move_send</i> \mapsto <i>Send</i>, <i>Mmt_copy_send</i> \mapsto <i>Send</i>, <i>Mmt_make_send_once</i> \mapsto <i>Receive</i>, <i>Mmt_move_send_once</i> \mapsto <i>Send_once</i>, <i>Mmt_move_receive</i> \mapsto <i>Receive</i> }</p>
--

This function captures the following semantics of port right transfers in Mach:

- A receive right can be moved or used to create a send or send-once right.

- A send right can be moved or copied.
- A send-once right can be moved.

Using this function, we specify the case in which the destination port is valid as follows:

<p><i>MachMsgSendValidHeader</i></p> <hr/> <p><i>MachMsgSendCanAllocateBuffer</i> <i>msg</i> : <i>MachMsgHeader</i></p> <hr/> <p>$msg.remote_rights \in TRANSFER_SEND_RIGHTS$ $(TRANSFER_SEND_RIGHTS \cap msg.local_rights \neq \emptyset \vee$ $(msg.local_rights = \emptyset \wedge$ $msg.local_port = Mach_port_null))$ $msg.remote_rights \in Recognized_transfer_options$ $msg.local_rights \subseteq Recognized_transfer_options$</p>
<p><i>MachMsgSendInvalidDest</i></p> <hr/> <p><i>MachMsgSendValidHeader</i> <i>MachMsgSendNoOp</i></p> <hr/> <p>(let $needed_rights == Required_right(msg.remote_rights);$ $port == named_port(owning_task(client?), msg.remote_port)$ • $\forall i : \mathbb{N}$ • $(owning_task(client?), port, msg.remote_port, needed_rights, i)$ $\notin \underline{port_right_rel}$) $msg_return! = Mm_send_invalid_dest$</p>

Otherwise, if the client task does not have the right specified in *msg.local_rights*, then an error message is returned and no further processing occurs. Note that if the client task specifies *Mmt_move_send* or *Mmt_move_send_once* in *msg.remote_rights*, then it loses a reference to *msg.remote_port*. This change in the number of references must be accounted for when determining whether the client has sufficient rights for *msg.local_port*.

Before defining the functions for manipulating the port name space, we first define the following schema to denote that the previous checks were successful:

<p><i>MachMsgSendValidDest</i></p> <hr/> <p><i>MachMsgSendValidHeader</i> $client? : THREAD$</p> <hr/> <p>(let $needed_rights == Required_right(msg.remote_rights);$ $port == named_port(owning_task(client?), msg.remote_port)$ • $\exists i : \mathbb{N}$ • $(owning_task(client?), port, msg.remote_port, needed_rights, i)$ $\in \underline{port_right_rel}$)</p>

We use the following type to denote the kernel data structure defining the port name spaces for each task. This structure has the same format and meaning as the relation *port_right_rel* in the definition of the Mach system state.

$$PORT_NAME_SPACE == \mathbb{P}(TASK \times PORT \times NAME \times RIGHT \times \mathbb{N}_1)$$

The function *Change_ref_count* is used to change the reference count associated with a name and a right in a task's port name space by a specified amount. If subtracting the specified amount from the current count results in a positive value, then the new count is that positive value. Otherwise, the name and right are not present in the port name space returned by this function.

$$\begin{array}{|l}
 \hline
 \text{Change_ref_count} : \mathbb{Z} \times \text{TASK} \times \text{NAME} \times \text{RIGHT} \times \text{PORT_NAME_SPACE} \longrightarrow \\
 \text{PORT_NAME_SPACE} \\
 \hline
 \forall \text{task, task}_1 : \text{TASK}; \text{port} : \text{PORT}; \text{name, name}_1 : \text{NAME}; \text{right, right}_1 : \text{RIGHT}; \\
 i : \mathbb{N}_1; n : \mathbb{Z}; \text{pns} : \text{PORT_NAME_SPACE} \bullet \\
 (\text{task, port, name, right, i}) \in \\
 \text{Change_ref_count}(n, \text{task}_1, \text{name}_1, \text{right}_1, \text{pns}) \Leftrightarrow \\
 ((\text{task, port, name, right, i}) \in \text{pns} \wedge \\
 (\text{task, name, right}) \neq (\text{task}_1, \text{name}_1, \text{right}_1)) \vee \\
 ((\text{task, port, name, right, i} + n) \in \text{pns} \wedge \\
 (\text{task, name, right}) = (\text{task}_1, \text{name}_1, \text{right}_1))
 \end{array}$$

The functions *Change_receive_count*, *Change_send_count*, and *Change_send_once_count* use *Change_ref_count* to change the count associated with a receive, send, or send-once right.

$$\begin{array}{|l}
 \hline
 \text{Change_receive_count} : \mathbb{Z} \times \text{TASK} \times \text{NAME} \times \text{PORT_NAME_SPACE} \longrightarrow \\
 \text{PORT_NAME_SPACE} \\
 \hline
 \forall n : \mathbb{Z}; \text{task} : \text{TASK}; \text{name} : \text{NAME}; \text{pns} : \text{PORT_NAME_SPACE} \bullet \\
 \text{Change_receive_count}(n, \text{task}, \text{name}, \text{pns}) = \\
 \text{Change_ref_count}(n, \text{task}, \text{name}, \text{Receive}, \text{pns})
 \end{array}$$

$$\begin{array}{|l}
 \hline
 \text{Change_send_count} : \mathbb{Z} \times \text{TASK} \times \text{NAME} \times \text{PORT_NAME_SPACE} \\
 \longrightarrow \text{PORT_NAME_SPACE} \\
 \hline
 \forall n : \mathbb{Z}; \text{task} : \text{TASK}; \text{name} : \text{NAME}; \text{pns} : \text{PORT_NAME_SPACE} \bullet \\
 \text{Change_send_count}(n, \text{task}, \text{name}, \text{pns}) = \\
 \text{Change_ref_count}(n, \text{task}, \text{name}, \text{Send}, \text{pns})
 \end{array}$$

$$\begin{array}{|l}
 \hline
 \text{Change_send_once_count} : \mathbb{Z} \times \text{TASK} \times \text{NAME} \times \text{PORT_NAME_SPACE} \longrightarrow \\
 \text{PORT_NAME_SPACE} \\
 \hline
 \forall n : \mathbb{Z}; \text{task} : \text{TASK}; \text{name} : \text{NAME}; \text{pns} : \text{PORT_NAME_SPACE} \bullet \\
 \text{Change_send_once_count}(n, \text{task}, \text{name}, \text{pns}) = \\
 \text{Change_ref_count}(n, \text{task}, \text{name}, \text{Send_once}, \text{pns})
 \end{array}$$

The function *Process_right* computes a new port name space from an old port name space, a task, a name, and a set of transfer options. If none of the transfer options requires moving a right, then the resulting name space is the same as the input name space. Otherwise, the count for each type of right that is moved is decremented.

$$\begin{array}{l}
\text{Process_right} : \mathbb{P} \text{ Recognized_transfer_options} \times \text{TASK} \times \text{NAME} \times \\
\text{PORT_NAME_SPACE} \rightarrow \text{PORT_NAME_SPACE} \\
\hline
\forall \text{mmt_set} : \mathbb{P} \text{ Recognized_transfer_options}; \text{task} : \text{TASK}; \text{name} : \text{NAME}; \\
\text{pns} : \text{PORT_NAME_SPACE} \bullet \\
(\{ \text{Mmt_move_receive}, \text{Mmt_move_send}, \text{Mmt_move_send_once} \} \cap \text{mmt_set} = \emptyset \Rightarrow \\
\text{Process_right}(\text{mmt_set}, \text{task}, \text{name}, \text{pns}) = \text{pns}) \wedge \\
(\text{Mmt_move_receive} \in \text{mmt_set} \Rightarrow \\
(\text{let } \text{mmt_set}_1 == \text{mmt_set} \setminus \{ \text{Mmt_move_receive} \}; \\
\text{pns}_1 == \text{Change_receive_count}(1, \text{task}, \text{name}, \text{pns}) \bullet \\
\text{Process_right}(\text{mmt_set}, \text{task}, \text{name}, \text{pns}) = \\
\text{Process_right}(\text{mmt_set}_1, \text{task}, \text{name}, \text{pns}_1))) \wedge \\
(\text{Mmt_move_send} \in \text{mmt_set} \Rightarrow \\
(\text{let } \text{mmt_set}_1 == \text{mmt_set} \setminus \{ \text{Mmt_move_send} \}; \\
\text{pns}_1 == \text{Change_send_count}(1, \text{task}, \text{name}, \text{pns}) \bullet \\
\text{Process_right}(\text{mmt_set}, \text{task}, \text{name}, \text{pns}) = \\
\text{Process_right}(\text{mmt_set}_1, \text{task}, \text{name}, \text{pns}_1))) \wedge \\
(\text{Mmt_move_send_once} \in \text{mmt_set} \Rightarrow \\
(\text{let } \text{mmt_set}_1 == \text{mmt_set} \setminus \{ \text{Mmt_move_send_once} \}; \\
\text{pns}_1 == \text{Change_send_once_count}(1, \text{task}, \text{name}, \text{pns}) \bullet \\
\text{Process_right}(\text{mmt_set}, \text{task}, \text{name}, \text{pns}) = \\
\text{Process_right}(\text{mmt_set}_1, \text{task}, \text{name}, \text{pns}_1)))
\end{array}$$

Using these functions, the case in which the reply port is invalid can be specified as follows.

$$\begin{array}{l}
\text{MachMsgSendInvalidReply} \\
\text{MachMsgSendValidHeader} \\
\text{MachMsgSendNoOp} \\
\hline
\text{let } \text{needed_rights} == \text{Required_right}(\text{msg}h.\text{local_rights}); \\
\text{new_port_right_rel} == \\
\text{Process_right}(\{ \text{msg}h.\text{remote_rights} \}, \text{owning_task client?}, \\
\text{msg}h.\text{remote_port}, \text{port_right_rel}); \\
\text{port} == \text{named_port}(\text{owning_task client?}, \text{msg}h.\text{local_port}) \bullet \\
\exists \text{right} : \text{RIGHT} \mid \\
\text{right} \in \text{needed_rights} \bullet \\
\forall i : \mathbb{N} \bullet \\
(\text{owning_task client?}, \text{port}, \text{msg}h.\text{local_port}, \text{right}, i) \\
\notin \text{new_port_right_rel} \\
\text{msg_return!} = \text{Mm_send_invalid_reply}
\end{array}$$

This is analogous to the case in which the destination port is invalid. The main difference is that the reference counts for the destination port are decremented, if necessary, before testing the reply port. When the reply and destination ports are the same, this decrementing can influence whether the reply port is valid.

If the client specifies the *Mach_send_cancel* option and *msg.h.notify* does not denote a receive right, then an error message is returned and no further processing takes place. As before, it is necessary to decrement reference counts associated with earlier right transfers before performing this test.

MachMsgSendValidReply
MachMsgSendValidHeader
MachMsgSendNoOp

let *needed_rights* == *Required_right*(*msg*.*local_rights*);
new_port_right_rel ==
Process_right({ *msg*.*remote_rights*}, *owning_task_client?*,
msg.*remote_port*, *port_right_rel*);
port == *named_port*(*owning_task_client?*, *msg*.*local_port*) •
 \forall *right* : *RIGHT* |
right \in *needed_rights* •
 \exists *i* : \mathbb{N} •
(*owning_task_client?*, *port*, *msg*.*local_port*, *right*, *i*)
 \in *new_port_right_rel*

MachMsgSendInvalidNotify
MachMsgSendValidReply
MachMsgSendNoOp

Mach_send_cancel \in *option?*
let *new_port_right_rel* ==
Process_right({ *msg*.*remote_rights*}, *owning_task_client?*,
msg.*remote_port*, *port_right_rel*)
• let *new_port_right_rel*₁ ==
Process_right(*msg*.*local_rights*, *owning_task_client?*,
msg.*local_port*, *new_port_right_rel*) •
(\forall *port* : *PORT*; *i* : \mathbb{N} •
(*owning_task_client?*, *port*, *notify?*, *Receive*, *i*) \notin *new_port_right_rel*₁)
msg_return! = *Mm_send_invalid_notify*

The following schema denotes the case in which the kernel can continue processing the message.

MachMsgSendValid
MachMsgSignature
MachMsgSendValidReply
option? : \mathbb{P} *MACH_MSG_OPTION*
notify? : *NAME*

Mach_send_cancel \notin *option?* \vee
(**let** *new_port_right_rel* ==
Process_right({ *msg*.*remote_rights*}, *owning_task_client?*,
msg.*remote_port*, *port_right_rel*)
• **let** *new_port_right_rel*₁ ==
Process_right(*msg*.*local_rights*, *owning_task_client?*,
msg.*local_port*, *new_port_right_rel*) •
(\exists *port* : *PORT*; *i* : \mathbb{N} •
(*owning_task_client?*, *port*, *notify?*, *Receive*, *i*) \in *new_port_right_rel*₁))

In this case, the message specified by the client is added to the set of messages in user space. Before describing this processing, we first describe functions that convert a message from its format in user space to its format in kernel space.

The function *Msgh_to_internal_msgh* converts the message header. The *remote_port* and *local_port* fields are filled in as specified by the two port parameters to the function. The remaining fields are copied without change.

$$\begin{array}{l} \text{Msgh_to_internal_msgh} : \text{PORT} \times \mathbb{P} \text{PORT} \times \text{MachMsgHeader} \longrightarrow \text{MachInternalHeader} \\ \hline \forall \text{msgh} : \text{MachMsgHeader}; \text{int_msgh} : \text{MachInternalHeader}; \text{port}_1 : \text{PORT}; \text{port}_2 : \mathbb{P} \text{PORT} \mid \\ \quad \text{Msgh_to_internal_msgh}(\text{port}_1, \text{port}_2, \text{msgh}) = \text{int_msgh} \bullet \\ \quad \text{msgh.local_rights} = \text{int_msgh.local_rights} \wedge \\ \quad \text{msgh.remote_rights} = \text{int_msgh.remote_rights} \wedge \\ \quad \text{msgh.size} = \text{int_msgh.size} \wedge \\ \quad \text{msgh.operation} = \text{int_msgh.operation} \wedge \\ \quad \text{port}_1 = \text{int_msgh.remote_port} \wedge \\ \quad \text{port}_2 = \text{int_msgh.local_port} \end{array}$$

Editorial Note:

The previous definition used to state that the *complex* field of the internal message header was copied from the user space message header. It appears that it is really generated by the kernel parsing the message.

The function *Msg_data_to_msg_value* converts an element of type *MSG_DATA* to an element of type *MSG_VALUE*.

$$\begin{array}{l} \text{Msg_data_to_msg_value} : \text{MSG_DATA} \longrightarrow \text{MSG_VALUE} \\ \hline \forall \text{msg_data} : \text{MSG_DATA} \bullet \\ \quad \text{Msg_data_to_msg_value} \text{ msg_data} = \text{V_data}(\text{msg_data}, \text{V_data_in}) \end{array}$$

The function *Msg_data_seq_to_msg_value_seq* converts an element of type *seq MSG_DATA* to an element of type *seq MSG_VALUE*.

$$\begin{array}{l} \text{Msg_data_seq_to_msg_value_seq} : \text{seq MSG_DATA} \longrightarrow \text{seq MSG_VALUE} \\ \hline \forall \text{data_seq} : \text{seq MSG_DATA} \bullet \\ \quad \text{Msg_data_seq_to_msg_value_seq} \text{ data_seq} = \text{Msg_data_to_msg_value} \circ \text{data_seq} \end{array}$$

The function *Msge_to_internal_msge* converts a single element of a message body. Elements in a message in user space are either *In_line* or *Out_of_line*. The former are converted to *Msg_value* entries, and the latter are converted to *Msg_region* entries. The function's task parameter is associated with the element to record the task in whose space out-of-line data and port rights should later be resolved.

$$\begin{array}{l} \text{Msge_to_internal_msge} : \text{TASK} \times \text{Msg_element} \longrightarrow \text{Internal_element} \\ \hline \forall n : \mathbb{N}; \text{mach_msg_type} : \text{MACH_MSG_TYPE}; \text{data_seq} : \text{seq MSG_DATA}; \\ \text{va} : \text{VIRTUAL_ADDRESS}; \text{int_msge} : \text{Internal_element}; \text{task} : \text{TASK}; \text{olzd} : \text{OLSD} \bullet \\ \quad (\text{let value_seq} == \text{Msg_data_seq_to_msg_value_seq} \text{ data_seq} \bullet \\ \quad \quad (\text{Msge_to_internal_msge}(\text{task}, \text{In_line}(n, \text{mach_msg_type}, \text{data_seq})) \\ \quad \quad \quad = \text{int_msge} \\ \quad \quad \quad \Rightarrow \text{int_msge} = \text{Msg_value}(n, \text{mach_msg_type}, (\text{task}, \text{value_seq}))) \\ \quad \quad \wedge (\text{Msge_to_internal_msge}(\text{task}, \text{Out_of_line}(n, \text{mach_msg_type}, \text{va}, \text{olzd})) \\ \quad \quad \quad = \text{int_msge} \\ \quad \quad \quad \Rightarrow \text{int_msge} = \text{Msg_region}(n, \text{mach_msg_type}, (\text{task}, \text{va}, \text{olzd})))))) \end{array}$$

The function $Msgb_to_internal_msgb$ converts a message body by applying $Msge_to_internal_msge$ to each element in the body.

$$\left| \begin{array}{l} Msgb_to_internal_msgb : TASK \times MESSAGE_BODY \rightarrow INTERNAL_BODY \\ \hline \forall msgb : MESSAGE_BODY; int_msgb : INTERNAL_BODY; task : TASK \mid \\ \quad Msgb_to_internal_msgb(task, msgb) = int_msgb \bullet \\ \quad \#msgb = \#int_msgb \wedge \\ \quad (\forall i : \mathbb{N} \mid i \in \text{dom } msgb \bullet \\ \quad \quad int_msgb(i) = Msge_to_internal_msge(task, msgb(i))) \end{array} \right.$$

Finally, a message in user space is converted to a message in user space by using $Msgb_to_internal_msgb$ to convert the header and using $Msgb_to_internal_msgb$ to convert the body.

$$\left| \begin{array}{l} Msg_to_internal_msg : \mathbb{N} \times \mathbb{N} \times \mathbb{P} MACH_MSG_OPTION \times PORT \times PORT \times TASK \times \\ \quad Message \rightarrow InternalMessage \\ \hline \forall task : TASK; port_1, port_2 : PORT; msg : Message; int_msg : InternalMessage; \\ \quad current_time, time_out : \mathbb{N}; option? : \mathbb{P} MACH_MSG_OPTION \\ \mid Msg_to_internal_msg(current_time, time_out, option?, port_1, port_2, task, msg) \\ \quad = int_msg \\ \bullet Msgb_to_internal_msgb(port_1, \{port_2\}, msg.header) = int_msg.header \wedge \\ \quad Msgb_to_internal_msgb(task, msg.body) = int_msg.body \wedge \\ \quad int_msg.option = option? \wedge \\ \quad int_msg.time_out_at = \mathbf{if} Mach_send_timeout \in option? \\ \quad \quad \mathbf{then} \{ current_time + time_out \} \\ \quad \quad \mathbf{else} \emptyset \wedge \\ \quad int_msg.status = Msg_stat_send \wedge \\ \quad int_msg.error = \emptyset \end{array} \right.$$

Note that:

- The $time_out_at$ field is set to indicate the earliest time at which the send request can time out.
If the client specified a time out was desired, then this field is set to the current time plus the specified time out duration. Otherwise, the $time_out_at$ field is set to \emptyset to denote that the send request should block rather than time out.
- The $status$ field is set to indicate that the message should be processed as part of a send request.
- The $error$ field is initialized to \emptyset .

The function $Msgb_to_internal_msgb$ requires inputs indicating the remote and local ports. The remote port can be determined by using $named_port$ to resolve the remote port name in the task's name space. When the local port name is not null, the same approach can be used to determine the local port. In the cases in which the local port name is null, we use $Null_port$ to denote the local port.

$$\mid Null_port : PORT$$

Before the message is moved into kernel space, the appropriate reference counts are decremented and the make send count is incremented for any port for which a send right was made.

The function *Update_ms_count* defines the changes that need to be made to *make_send_count*. The count for the remote port must be incremented by 1 if a send right was made. Similarly, the count for the local port must be incremented by 1 if it exists and a send right was made for it.

$$\begin{array}{l}
 \text{Update_ms_count} : (\text{PORT} \times \mathbb{P} \text{ Recognized_transfer_options}) \times \\
 (\text{PORT} \times \mathbb{P} \text{ Recognized_transfer_options}) \times \mathbb{P} \text{ PORT} \times \\
 (\text{PORT} \mapsto \mathbb{N}) \rightarrow (\text{PORT} \mapsto \mathbb{N}) \\
 \hline
 \forall \text{port}_1, \text{port}_2 : \text{PORT}; \text{old_ms_count} : \text{PORT} \mapsto \mathbb{N}; \\
 \text{mmt_set}_1, \text{mmt_set}_2 : \mathbb{P} \text{ MACH_MSG_TYPE}; \text{port_set} : \mathbb{P} \text{ PORT} \bullet \\
 \text{Update_ms_count}((\text{port}_1, \text{mmt_set}_1), (\text{port}_2, \text{mmt_set}_2), \text{port_set}, \text{old_ms_count}) = \\
 \text{if } \text{port}_2 \notin \text{port_set} \\
 \quad \text{then } \text{old_ms_count} \oplus \{ \text{port}_1 \mapsto \\
 \quad \quad \text{old_ms_count } \text{port}_1 + \#(\{ \text{Mmt_make_send} \} \cap \text{mmt_set}_1) \} \\
 \text{else if } \text{port}_1 = \text{port}_2 \\
 \quad \text{then } \text{old_ms_count} \oplus \{ \text{port}_1 \mapsto \\
 \quad \quad \text{old_ms_count } \text{port}_1 + \#(\{ \text{Mmt_make_send} \} \cap \text{mmt_set}_1) + \\
 \quad \quad \#(\{ \text{Mmt_make_send} \} \cap \text{mmt_set}_2) \} \\
 \text{else} \\
 \quad \text{old_ms_count} \oplus \{ \text{port}_1 \mapsto \\
 \quad \quad \text{old_ms_count } \text{port}_1 + \#(\{ \text{Mmt_make_send} \} \cap \text{mmt_set}_1), \\
 \quad \quad \text{port}_2 \mapsto \text{old_ms_count } \text{port}_2 + \#(\{ \text{Mmt_make_send} \} \cap \text{mmt_set}_2) \}
 \end{array}$$

Note that the above definition accounts for the possibility that the local and remote ports are the same by counting the send rights made for either port against the common port.

The function *Update_name_space* performs any necessary decrementing of the reference counts for the local and remote ports. It does so by first using *Process_right* to address the remote port and then using *Process_right* on the result to address the local port.

$$\begin{array}{l}
 \text{Update_name_space} : (\mathbb{P} \text{ MACH_MSG_TYPE} \times \text{NAME}) \times \\
 (\mathbb{P} \text{ MACH_MSG_TYPE} \times \text{NAME}) \times \\
 \text{TASK} \times \text{PORT_NAME_SPACE} \rightarrow \text{PORT_NAME_SPACE} \\
 \hline
 \forall \text{mmt_set}_1, \text{mmt_set}_2 : \mathbb{P} \text{ MACH_MSG_TYPE}; \text{task} : \text{TASK}; \text{name}_1, \text{name}_2 : \text{NAME}; \\
 \text{pns} : \text{PORT_NAME_SPACE} \bullet \\
 \text{Update_name_space}((\text{mmt_set}_1, \text{name}_1), (\text{mmt_set}_2, \text{name}_2), \text{task}, \text{pns}) = \\
 (\text{let } \text{pns}_1 == \text{Process_right}(\text{mmt_set}_1, \text{task}, \text{name}_1, \text{pns}) \bullet \\
 \quad \text{Process_right}(\text{mmt_set}_2, \text{task}, \text{name}_2, \text{pns}_1))
 \end{array}$$

Using the previously defined functions, the entering of a send message request into kernel space can be specified as follows:

$\frac{\text{MachMsgSendStart}}{\text{MachMsgSendValid}}$ $\exists \text{message} : \text{MESSAGE}; \text{msg} : \text{Message};$ $\text{port}_1, \text{port}_2 : \text{PORT}; \text{task} : \text{TASK} \bullet$ $\text{message} \in \underline{\text{message_exists}}' \setminus \underline{\text{message_exists}} \wedge$ $\text{msg.header} = \text{msgh} \wedge$ $\text{msg.body} = \text{msg_body} \wedge$ $\text{task} = \text{owning_task client?} \wedge$ $\text{port}_1 = \text{named_port}(\text{task}, \text{msgh.remote_port}) \wedge$ $(\text{port}_2 = \text{if} (\text{task}, \text{msgh.local_port}) \in \text{dom named_port}$ $\quad \text{then named_port}(\text{task}, \text{msgh.local_port})$ $\quad \text{else Null_port}) \wedge$ $\underline{\text{msg_contents}}' = \underline{\text{msg_contents}} \cup \{ \text{message} \mapsto$ $\quad \text{Msg_to_internal_msg}(\underline{\text{host_time}}, \text{time_out?}, \text{option?}, \text{port}_1, \text{port}_2,$ $\quad \text{task}, \text{msg}) \}$ $\wedge \underline{\text{make_send_count}}' =$ $\quad (\text{let pair}_1 == (\text{port}_1, \{ \text{msgh.remote_rights} \});$ $\quad \text{pair}_2 == (\text{port}_2, \text{msgh.local_rights}) \bullet$ $\quad \text{Update_ms_count}(\text{pair}_1, \text{pair}_2, \underline{\text{port_exists}}, \underline{\text{make_send_count}}))$ $(\text{let pair}_1 == (\{ \text{msgh.remote_rights} \}, \text{msgh.remote_port});$ $\text{pair}_2 == (\text{msgh.local_rights}, \text{msgh.local_port}) \bullet$ $\quad \underline{\text{port_right_rel}}' = \text{Update_name_space}(\text{pair}_1, \text{pair}_2, \text{owning_task client?},$ $\quad \underline{\text{port_right_rel}}))$

C.2.1.2 Kernel Processing In this section, we describe the processing of messages in kernel space that are not yet queued at a port.

The function *Unprocessed_rights* returns the set of port rights in transit that must be processed before a message can be enqueued. An element $(\text{message}, i, j)$ belongs to the resulting set exactly when the i^{th} element of *message*'s body is a data element whose j^{th} entry is an unresolved port right. Note that regardless of the types of the data elements in a message body, no rights are transferred unless the *complex* field of the message header indicates rights are being transferred.

$\frac{\text{Unprocessed_rights} : \text{Mach} \rightarrow \mathbb{P}(\text{MESSAGE} \times \mathbb{N} \times \mathbb{N})}{\forall \text{mach_st} : \text{Mach} \bullet}$ $\text{Unprocessed_rights mach_st} = \{ \text{message} : \text{MESSAGE}; i, j : \mathbb{N} \mid$ $\text{message} \in \text{mach_st.}\underline{\text{message_exists}} \wedge$ $\text{Co_carries_rights} \in (\text{mach_st.}\underline{\text{msg_contents}} \text{ message}).\text{header.complex} \wedge$ $(\text{mach_st.}\underline{\text{msg_contents}} \text{ message}).\text{status} = \text{Msg_stat_send} \wedge$ $(\text{let int_msgb} == (\text{mach_st.}\underline{\text{msg_contents}} \text{ message}).\text{body} \bullet$ $\quad i \in \text{dom int_msgb} \wedge$ $(\exists n : \mathbb{N}; \text{mach_msg_type} : \text{MACH_MSG_TYPE};$ $\quad \text{value_seq} : \text{seq MSG_VALUE}; \text{task} : \text{TASK};$ $\quad \text{msg_data} : \text{MSG_DATA}; \text{v_data_l} : \text{V_DATA_LOCATION}$ $\quad \bullet \text{int_msgb}(i) = \text{Msg_value}(n, \text{mach_msg_type}, (\text{task}, \text{value_seq})) \wedge$ $\quad \text{mach_msg_type} \in \text{Recognized_transfer_options} \wedge$ $\quad j \in \text{dom value_seq} \wedge$ $\quad \text{value_seq}(j) = \text{V_data}(\text{msg_data}, \text{v_data_l}))) \}$

The function *Unprocessed_memories* returns the set of memory objects in transit that must be processed before a message can be enqueued. An element $(message, i)$ belongs to the resulting set exactly when the i^{th} element of *message*'s body is an unprocessed out-of-line data element. Note that regardless of the types of the data elements in a message body, no memories are transferred unless the *complex* field of the message header indicates memories are being transferred.

$$\begin{array}{l}
 \underline{Unprocessed_memories : Mach \rightarrow \mathbb{P}(MESSAGE \times \mathbb{N})} \\
 \forall mach_st : Mach \bullet \\
 \quad Unprocessed_memories\ mach_st = \{ message : MESSAGE; i : \mathbb{N} \mid \\
 \quad \quad message \in mach_st.\underline{message_exists} \wedge \\
 \quad \quad Co_carries_memory \in (mach_st.\underline{msg_contents}\ message).header.complex \wedge \\
 \quad \quad (mach_st.\underline{msg_contents}\ message).status = Msg_stat_send \wedge \\
 \quad \quad (\text{let } int_msgb == (mach_st.\underline{msg_contents}\ message).body \bullet \\
 \quad \quad \quad i \in \text{dom } int_msgb \wedge \\
 \quad \quad \quad (\exists n : \mathbb{N}; mach_msg_type : MACH_MSG_TYPE; oldsd : OLSD; \\
 \quad \quad \quad \quad task : TASK; va : VIRTUAL_ADDRESS \bullet \\
 \quad \quad \quad \quad \quad int_msgb(i) = Msg_region(n, mach_msg_type, (task, va, oldsd)))) \}
 \end{array}$$

The function *Element_type* returns the type of an element in a message body.

$$\begin{array}{l}
 \underline{Element_type : Internal_element \rightarrow MACH_MSG_TYPE} \\
 \forall inte : Internal_element; n : \mathbb{N}; mach_msg_type : MACH_MSG_TYPE; \\
 \quad value_seq : \text{seq } MSG_VALUE; task : TASK; \\
 \quad va : VIRTUAL_ADDRESS; oldsd : OLSD; \\
 \quad memory : MEMORY; offset : OFFSET \mid \\
 \quad \quad inte \in \\
 \quad \quad \quad \{ Msg_value(n, mach_msg_type, (task, value_seq)), \\
 \quad \quad \quad \quad Msg_region(n, mach_msg_type, (task, va, oldsd)), \\
 \quad \quad \quad \quad \quad Transit_memory(n, mach_msg_type, (task, memory, offset)) \} \bullet \\
 \quad \quad \quad \quad Element_type\ inte = mach_msg_type
 \end{array}$$

The set *Invalid_msg_types* indicates the set of message elements having invalid data types. An element $(message, i)$ belongs to the resulting set exactly when the type specified for the i^{th} element of *message*'s body is invalid. The set *Valid_data_types* defines the set of valid data types.

$$\begin{array}{l}
 \underline{Valid_data_types : \mathbb{P} MACH_MSG_TYPE} \\
 \underline{Invalid_msg_types : Mach \rightarrow \mathbb{P}(MESSAGE \times \mathbb{N})} \\
 \quad Recognized_transfer_options \subseteq Valid_data_types \\
 \forall mach_st : Mach \bullet \\
 \quad \quad Invalid_msg_types\ mach_st = \{ message : MESSAGE; i : \mathbb{N} \mid \\
 \quad \quad \quad message \in mach_st.\underline{message_exists} \wedge \\
 \quad \quad \quad i \in \text{dom}(mach_st.\underline{msg_contents}\ message).body \wedge \\
 \quad \quad \quad \quad Element_type((mach_st.\underline{msg_contents}\ message).body(i)) \notin Valid_data_types \}
 \end{array}$$

The set *Processed_messages* indicates the set of messages that are not yet enqueued but require no further processing. In other words, these are messages that have no elements with invalid data types or unprocessed rights or memories and that are not present in any message queue.

$$\begin{array}{l}
 \underline{Processed_messages} : Mach \rightarrow \mathbb{P} MESSAGE \\
 \forall mach_st : Mach \bullet \\
 \quad \underline{Processed_messages} mach_st = \\
 \quad \quad mach_st.\underline{message_exists} \setminus \\
 \quad \quad \quad (\{ message : MESSAGE; i, j : \mathbb{N} \mid \\
 \quad \quad \quad \quad (message, i, j) \in \underline{Unprocessed_rights} mach_st \bullet message \} \cup \\
 \quad \quad \quad \{ message : MESSAGE; i : \mathbb{N} \mid \\
 \quad \quad \quad \quad (message, i) \in \underline{Invalid_msg_types} mach_st \bullet message \} \cup \\
 \quad \quad \quad \{ message : MESSAGE; i : \mathbb{N} \mid \\
 \quad \quad \quad \quad (message, i) \in \underline{Unprocessed_memories} mach_st \bullet message \} \cup \\
 \quad \quad \quad \{ message : MESSAGE \mid (\exists port : PORT \bullet \\
 \quad \quad \quad \quad message \in \text{ran}(mach_st.\underline{message_in_port_rel} port)) \} \cup \\
 \quad \quad \quad \{ message : MESSAGE \mid (mach_st.\underline{msg_contents} message).status \\
 \quad \quad \quad \quad \neq Msg_stat_send \})
 \end{array}$$

The function *Address_to_index* is used to convert a virtual address into a page index.

$$\underline{Address_to_index} : VIRTUAL_ADDRESS \rightarrow PAGE_INDEX$$

Before describing the processing of message elements, we define the following schema to represent parts of the processing that are common to the various cases to be considered:

$$\begin{array}{l}
 \underline{GeneralSendProcessing} \\
 \Delta DtosExec \\
 \quad message : MESSAGE \\
 \quad i, n : \mathbb{N} \\
 \quad int_msg_1, int_msg_2 : InternalMessage \\
 \quad v_data_l : V_DATA_LOCATION \\
 \quad task : TASK \\
 \quad value_seq_1, value_seq_2 : seq MSG_VALUE \\
 \quad mach_msg_type : MACH_MSG_TYPE \\
 \quad va : VIRTUAL_ADDRESS \\
 \quad olsd : OLSD \\
 \quad memory : MEMORY \\
 \quad offset : OFFSET \\
 \quad error : MSG_ERROR \\
 \quad page_index : PAGE_INDEX \\
 \quad \underline{message} \in \underline{message_exists} \\
 \quad int_msg_1 = \underline{msg_contents} message \\
 \quad i \in \text{dom}(int_msg_1.body) \\
 \quad int_msg_1.body(i) \in \\
 \quad \quad \{ Msg_value(n, mach_msg_type, (task, value_seq_1)), \\
 \quad \quad \quad Msg_region(n, mach_msg_type, (task, va, olsd)), \\
 \quad \quad \quad Transit_memory(n, mach_msg_type, (task, memory, offset)) \} \\
 \quad \underline{msg_contents}' = \underline{msg_contents} \oplus \{ message \mapsto int_msg_2 \} \\
 \quad int_msg_2.header = int_msg_1.header \\
 \quad int_msg_2.option = int_msg_1.option \\
 \quad int_msg_2.time_out_at = int_msg_1.time_out_at \\
 \quad int_msg_1.error \neq \emptyset \Rightarrow int_msg_2.error = int_msg_1.error \\
 \quad page_index = \underline{Address_to_index} va
 \end{array}$$

This schema requires that *message* is an existing message and *i* is a valid index for the body of the message associated with *message*. The processing of the element is accomplished by modifying the body of the message. The components *int_msg₁* and *int_msg₂* are introduced to denote the initial and final values for the message. It is required that the *header*, *option*, and *time_out_at* fields of the message are not altered. Furthermore, it is required that the *error* field cannot be altered if it is nonempty. The remaining components of the schema are introduced to define the general form of the *ith* element of the message body.

The function *replace_entry* replaces a specified entry in a sequence with a specified value.

$$\begin{array}{l} \text{---}[X] \text{---} \\ \text{---} \\ \text{replace_entry} : \mathbb{N}_1 \times X \times \text{seq } X \longrightarrow \text{seq } X \\ \text{---} \\ \forall i : \mathbb{N}_1; x : X; x_seq : \text{seq } X \mid i \in \text{dom } x_seq \bullet \\ \text{replace_entry}(i, x, x_seq) = x_seq \oplus \{ i \mapsto x \} \\ \text{---} \end{array}$$

The function *Data_to_name* converts an element of type *MSG_DATA* to an element of type *NAME*. It is assumed that this function is an injection.

$$\mid \text{Data_to_name} : \text{MSG_DATA} \rightsquigarrow \text{NAME}$$

From the schema *GeneralSendProcessing*, we build the following schema for processing port rights:

$$\begin{array}{l} \text{---GeneralSendProcessing2 \text{---} \\ GeneralSendProcessing \\ j : \mathbb{N} \\ value_seq_2 : \text{seq } \text{MSG_VALUE} \\ port : \text{PORT} \\ msg_data : \text{MSG_DATA} \\ name : \text{NAME} \\ \text{---} \\ int_msg_1.body(i) = \text{Msg_value}(n, mach_msg_type, (task, value_seq_1)) \\ mach_msg_type \in \text{Recognized_transfer_options} \\ Co_carries_rights \in int_msg_1.header.complex \\ j \in \text{dom } value_seq_1 \\ value_seq_1(j) = V_data(msg_data, v_data_l) \\ name = \text{Data_to_name } msg_data \\ port = \text{if } (task, name) \in \text{dom } \text{named_port} \wedge \\ (\exists k : \mathbb{N} \bullet \\ (task, \text{named_port}(task, name), name, \text{Required_right}(mach_msg_type), k) \in \\ \text{port_right_rel}) \\ \text{then } \text{named_port}(task, name) \\ \text{else } \text{Null_port} \\ value_seq_2 = \text{replace_entry}(j, V_port(port, v_data_l), value_seq_1) \\ int_msg_2.body = \\ \text{replace_entry}(i, \text{Msg_value}(n, mach_msg_type, (task, value_seq_2)), \\ int_msg_1.body) \\ \text{---} \end{array}$$

This schema requires that the message element being processing is of the *Msg_value* form and the type of the message element indicates that a port right is being transferred. The component *j* indicates the index into *value_seq₁* denoting the right to be processed. The component

$value_seq_2$ is introduced to denote the new sequence of values to be stored as the i^{th} element in the body. The new sequence is obtained by replacing the j^{th} entry of the original sequence with an entry indicating the port to which the transferred right resolves. The component $port$ is introduced to denote this port. The component $name$ is introduced to represent the $name$ of the transferred right. The name is defined by the data in the j^{th} element of the sequence. If the name is in the task's name space and the task has the appropriate rights to transfer the port, then $port$ is defined to be the port associated with the name in the task's name space. Otherwise, $port$ is defined to be the $Null_port$.

If a message contains an element with an invalid data type, then progress can be made in processing the message by processing that element. The message element is processed by removing it from the body and recording the error condition if no error condition has previously been recorded.

The function $remove_entry$ removes a specified entry from a sequence.

$$\begin{array}{l} \text{---}[X] \text{---} \\ \text{---} \\ remove_entry : \mathbb{N} \times seq\ X \rightarrow seq\ X \\ \text{---} \\ \forall i : \mathbb{N}; x_seq : seq\ X \bullet \\ \quad remove_entry(i, x_seq) = ((1 \dots i) \upharpoonright x_seq) \hat{\ } (((i + 1) \dots \#x_seq) \upharpoonright x_seq) \\ \text{---} \end{array}$$

Using this function, the processing of a data element with an invalid type is as follows:

$$\begin{array}{l} \text{---}ProcessInvalidType \text{---} \\ \text{---} \\ GeneralSendProcessing \\ \text{---} \\ (message, i) \in Invalid_msg_types(\theta Mach) \\ int_msg_2.body = remove_entry(i, int_msg_1.body) \\ int_msg_2.status = int_msg_1.status \\ int_msg_1.error = \emptyset \Rightarrow \\ \quad int_msg_2.error = \{ Msg_error_invalid_type \} \\ \text{---} \end{array}$$

If a message contains port rights that have not yet been processed, then progress can be made in processing the message by processing one of the port rights. The first case to consider is that in which the name being processed does not denote a right appropriate for the type of transfer requested. In this case, $GeneralSendProcessing2$ resolves the name to $Null_port$. Thus, the processing for this case can be specified as follows:

$$\begin{array}{l} \text{---}ProcessRightBad \text{---} \\ \text{---} \\ GeneralSendProcessing2 \\ \text{---} \\ (message, i, j) \in Unprocessed_rights(\theta Mach) \\ port = Null_port \\ int_msg_2.status = int_msg_1.status \\ int_msg_1.error = \emptyset \Rightarrow \\ \quad int_msg_2.error = \{ Msg_error_invalid_right \} \\ \text{---} \end{array}$$

In other words, the invalid right is replaced by a right for $Null_port$. The conversion from an entry of type V_data to V_port makes progress towards completion of the request since there is one less unprocessed right in the resulting state.

The only differences between this and the processing of a valid port right are:

- it is not necessary to record an error for a valid port right
- the task's port name space and port's make-send count must be updated

ProcessRightGood
GeneralSendProcessing2

```

(message, i, j) ∈ Unprocessed_rights(θMach)
port ≠ Null_port
port_right_rel' = Process_right({ mach_msg_type }, task, name, port_right_rel)
make_send_count' = make_send_count ⊕ { port ↦
    make_send_count port + #({ mach_msg_type } ∩ { Mmt_make_send }) }
int_msg2.status = int_msg1.status ∧
int_msg2.error = int_msg1.error

```

The case in which an out-of-line memory region is inaccessible to the sending task is specified as follows:

ProcessMemoryBad
GeneralSendProcessing

```

(message, i) ∈ Unprocessed_memories(θMach)
int_msg1.body(i) = Msg_region(n, mach_msg_type, (task, va, old))
((task, page_index) ∉ allocated ∨
    Read ∉ protection(task, page_index))
int_msg2.body = remove_entry(i, int_msg1.body)
int_msg2.status = int_msg1.status ∧
int_msg1.error = ∅ ⇒
    int_msg2.error = { Msg_error_invalid_memory }

```

If an out-of-line memory region is accessible and does not carry any port rights, then the element of form *Msg_region* can be converted into an element of form *Transit_memory*.

ProcessMemoryGood
GeneralSendProcessing

```

(message, i) ∈ Unprocessed_memories(θMach)
int_msg1.body(i) = Msg_region(n, mach_msg_type, (task, va, old))
(task, page_index) ∈ allocated
Read ∉ protection(task, page_index)
(mach_msg_type ∉ Recognized_transfer_options ∨
    Co_carries_rights ∉ int_msg1.header.complex)
((task, page_index), (memory, offset)) ∈ map_rel
(let inte == Transit_memory(n, mach_msg_type, (task, memory, offset)) •
    int_msg2.body = replace_entry(i, inte, int_msg1.body))
map_rel' = if old = Msg_deallocate
    then { (task, page_index) } ≺ map_rel else map_rel
int_msg2.status = int_msg1.status ∧
int_msg1.error = int_msg2.error

```

Note that the transferred memory is deallocated from the address space of the sending task if *old* indicates that it should be deallocated.

If an out-of-line memory region is accessible, carries port rights, and is currently in memory, then the element of form *Msg_region* can be converted into an element of form *Msg_value*. The resulting element will subsequently be processed by *ProcessRightGood* or *ProcessRightBad*.

The function *Va_offset* is used to add an integer to a virtual address.

$$| \quad Va_offset : VIRTUAL_ADDRESS \times \mathbb{N} \rightsquigarrow VIRTUAL_ADDRESS$$

The function *Index_to_offset* is used to convert a page index to a page offset.

$$| \quad Index_to_offset : PAGE_INDEX \rightsquigarrow PAGE_OFFSET$$

The function *Word_to_data* is used to convert a word on a page to a data item.

$$| \quad Word_to_data : WORD \rightsquigarrow MSG_DATA$$

ProcessAvailableOutOfLineRights

GeneralSendProcessing

```
(message, i) ∈ Unprocessed_memories(θMach)
int_msg1.body(i) = Msg_region(n, mach_msg_type, (task, va, old))
mach_msg_type ∈ Recognized_transfer_options
Co_carries_rights ∈ int_msg1.header.complex
∃ page_reference_set : P(N × VIRTUAL_ADDRESS × PAGE_INDEX) •
  page_reference_set = { m : N; va1 : VIRTUAL_ADDRESS;
    page_index1 : PAGE_INDEX |
    m ∈ 1..n ∧
    va1 = Va_offset(va, m) ∧
    page_index1 = Address_to_index va1 } ∧
value_seq2 = { k : N; va2 : VIRTUAL_ADDRESS;
  page : PAGE; page_offset : PAGE_OFFSET; word : WORD;
  page_index2 : PAGE_INDEX; memory2 : MEMORY; offset2 : OFFSET;
  msg_data : MSG_DATA |
  (k, va2, page_index2) ∈ page_reference_set ∧
  page_index2 = Address_to_index va2 ∧
  (task, page_index2) ∈ allocated ∧
  Read ∈ protection(task, page_index2) ∧
  ((task, page_index2), (memory2, offset2)) ∈ map_rel ∧
  (page, (memory2, offset2)) ∈ represents_rel ∧
  page_offset = Index_to_offset page_index2 ∧
  ((page, page_offset), word) ∈ page_word_rel ∧
  msg_data = Word_to_data word •
  (k, V_data(msg_data, V_data_out)) } ∧
(let inte == Msg_value(n, mach_msg_type, (task, value_seq2)) •
  int_msg2.body = replace_entry(i, inte, int_msg1.body)) ∧
map_rel' = if old = Msg_deallocate
  then { r : N; va3 : VIRTUAL_ADDRESS; page_index3 : PAGE_INDEX |
    (r, va3, page_index3) ∈ page_reference_set •
    (task, page_index3) } ≺ map_rel
  else map_rel
int_msg2.status = int_msg1.status ∧
int_msg1.error = int_msg2.error
```

The set *page_reference_set* denotes the pages that are referenced by the message element. If each page is in memory, the necessary data can be read from the pages and stored in *value_seq₂*. Note that if *oldd* indicates that the region should be deallocated, then each of the referenced pages is removed from the address space of the sending task.

If a page referenced by the message element is not accessible, then the processing is analogous to that described by *ProcessMemoryBad*.

<i>ProcessRightsMemoryBad</i>
<p><i>GeneralSendProcessing</i></p> <p>$(message, i) \in Unprocessed_memories(\theta Mach)$ $int_msg_1.body(i) = Msg_region(n, mach_msg_type, (task, va, oldd))$ $mach_msg_type \in Recognized_transfer_options$ $Co_carries_rights \in int_msg_1.header.complex$ $\exists page_reference_set : \mathbb{P}(\mathbb{N} \times VIRTUAL_ADDRESS \times PAGE_INDEX) \bullet$ $page_reference_set = \{ m : \mathbb{N}; va_1 : VIRTUAL_ADDRESS;$ $page_index_1 : PAGE_INDEX \mid$ $m \in 1..n \wedge$ $va_1 = Va_offset(va, m) \wedge$ $page_index_1 = Address_to_index\ va_1 \} \wedge$ $(\exists k : \mathbb{N}; va_2 : VIRTUAL_ADDRESS;$ $page_index_2 : PAGE_INDEX \bullet$ $(k, va_2, page_index_2) \in page_reference_set \wedge$ $page_index_2 = Address_to_index\ va_2 \wedge$ $((task, page_index_2) \notin \underline{allocated} \vee$ $Read \notin protection(task, page_index_2))) \wedge$ $\underline{map_rel}' = \mathbf{if}\ oldd = Msg_deallocate$ $\mathbf{then} \{ r : \mathbb{N}; va_3 : VIRTUAL_ADDRESS; page_index_3 : PAGE_INDEX \mid$ $(r, va_3, page_index_3) \in page_reference_set \bullet$ $(task, page_index_3) \} \triangleleft \underline{map_rel}$ $\mathbf{else} \underline{map_rel}$ $int_msg_2.body = remove_entry(i, int_msg_1.body)$ $int_msg_2.status = int_msg_1.status$ $int_msg_1.error = \emptyset \Rightarrow$ $int_msg_2.error = \{ Msg_error_invalid_memory \}$</p>

If a page referenced by the message element is accessible but is not in memory, then the kernel must request the page's data from the page's memory manager.

The following function is used to build the header for the request sent to the memory manager.

<p><i>Mach_object_data_request</i> : OPERATION <i>Build_data_request_header</i> : PORT × PORT → MachInternalHeader</p> <p>$\forall port_1, port_2 : PORT; int_msg_h : MachInternalHeader \mid$ $Build_data_request_header(port_1, port_2) = int_msg_h \bullet$ $int_msg_h.local_rights = \{ Mmt_copy_send \} \wedge$ $int_msg_h.remote_rights = Mmt_make_send_once \wedge$ $int_msg_h.complex = \emptyset \wedge$ $int_msg_h.remote_port = port_1 \wedge$ $int_msg_h.local_port = \{ port_2 \} \wedge$ $int_msg_h.operation = Mach_object_data_request$</p>
--

The first port parameter is the remote port and the second one is the local port. The operation is specified as being *Mach_object_data_request*.

The following functions are used to build the body for the request sent to the memory manager.

$ \begin{aligned} &Mmt_integer : MACH_MSG_TYPE \\ &Mmt_protection : MACH_MSG_TYPE \\ &Integer_to_data : \mathbb{N} \rightarrow MSG_DATA \\ &Protection_to_data : \mathbb{P} PROTECTION \rightarrow MSG_DATA \\ &Build_data_request_body : TASK \times \mathbb{N} \times \mathbb{N} \times \mathbb{P} PROTECTION \rightarrow INTERNAL_BODY \end{aligned} $
$ \begin{aligned} &\forall task : TASK; i, j : \mathbb{N}; prot_set : \mathbb{P} PROTECTION \bullet \\ &Build_data_request_body(task, i, j, prot_set) = \\ &\quad (\text{let } value_seq_1 == \\ &\quad \langle V_data(Integer_to_data\ i, V_data_in), V_data(Integer_to_data\ j, V_data_in) \rangle; \\ &\quad \quad value_seq_2 == \langle V_data(Protection_to_data\ prot_set, V_data_in) \rangle \bullet \\ &\quad \quad \langle Msg_value(2, Mmt_integer, (task, value_seq_1)), \\ &\quad \quad \quad Msg_value(1, Mmt_protection, (task, value_seq_2)) \rangle) \end{aligned} $

The integers indicate, respectively, the desired offset in the memory object and length of the data. The set of protections specify the access modes desired for the object.

The following function is used to build the request sent to the memory manager.

$ \begin{aligned} &Build_data_request : TASK \times PORT \times PORT \times \mathbb{N} \times \mathbb{N} \times \mathbb{P} PROTECTION \rightarrow \\ &\quad InternalMessage \end{aligned} $
$ \begin{aligned} &\forall task : TASK; port_1, port_2 : PORT; i, j : \mathbb{N}; prot_set : \mathbb{P} PROTECTION; \\ &int_msg : InternalMessage \mid \\ &\quad int_msg = Build_data_request(task, port_1, port_2, i, j, prot_set) \bullet \\ &\quad \quad int_msg.header = Build_data_request_header(port_1, port_2) \wedge \\ &\quad \quad int_msg.body = Build_data_request_body(task, i, j, prot_set) \wedge \\ &\quad \quad int_msg.option = \{ Mach_send_msg \} \wedge \\ &\quad \quad int_msg.time_out_at = \emptyset \wedge \\ &\quad \quad int_msg.status = Msg_stat_send \wedge \\ &\quad \quad int_msg.error = \emptyset \end{aligned} $

The function *Index_to_nat* is used to convert a page index to an integer.

$Index_to_nat : PAGE_INDEX \rightarrow \mathbb{N}$

The constant *Page_size* denotes the number of words on a page.

$Page_size : \mathbb{N}$

Using these definitions, the sending of a request to the memory manager can be specified as follows:

*RequestRightsData**GeneralSendProcessing**int_msg!* : *InternalMessage*

```

(message, i) ∈ Unprocessed_memories( $\theta$  Mach)
int_msg1.body(i) = Msg_region(n, mach_msg_type, (task, va, oldd))
mach_msg_type ∈ Recognized_transfer_options
Co_carries_rights ∈ int_msg1.header.complex
∃ page_reference_set :  $\mathbb{P}(\mathbb{N} \times \text{VIRTUAL\_ADDRESS} \times \text{PAGE\_INDEX})$  •
  page_reference_set = { m :  $\mathbb{N}$ ; va1 : VIRTUAL_ADDRESS;
    page_index1 : PAGE_INDEX |
    m ∈ 1 .. n ∧
    va1 = Va_offset(va, m) ∧
    page_index1 = Address_to_index va1 } ∧
  (∃ k :  $\mathbb{N}$ ; page : PAGE; va2 : VIRTUAL_ADDRESS;
    page_index2 : PAGE_INDEX; memory2 : MEMORY; offset2 : OFFSET •
    (k, va2, page_index2) ∈ page_reference_set ∧
    page_index2 = Address_to_index va2 ∧
    (task, page_index2) ∈ allocated ∧
    Read ∈ protection(task, page_index2) ∧
    ((task, page_index2), (memory2, offset2)) ∈ map_rel ∧
    (page, (memory2, offset2)) ∉ represents_rel ∧
    int_msg! =
      (let port1 == object_port memory2;
        port2 == control_port memory2;
        prot_set == { Read, Write, Execute };
        r == Index_to_nat page_index2;
        s == max{ Page_size, n - (k - 1) * Page_size } •
          Build_data_request(k kernel, port1, port2, r, s, prot_set)))
int_msg2.body = remove_entry(i, int_msg1.body)
int_msg2.status = int_msg1.status
int_msg1.error = int_msg2.error

```

Note that *int_msg!* denotes the message that should be sent to the memory manager. The “sending” of this message would be represented by adding it to the range of *msg_contents*. For simplicity, we do not address that processing here.

If the destination port for a message that has been processed does not exist, then the message can be discarded.

Editorial Note:

The model of the state component *f_orcibly_queued* was previously as a function from a port to a message. This model was based upon the Kernel Principles document which states “*mach_msg* provides an option allowing one message to be left waiting to be queued.” However, this is not one message per port, but one message per port right. The model has now been fixed, but its ramifications on this section have not been determined. Therefore all mention of *f_orcibly_queued* within the Z has been commented out in this section (though none of the text has been changed).

CAR 4041 has been filed to address this issue.

PortDied $\Delta \text{DtosExec}$ $\text{message} : \text{MESSAGE}$ <hr/> $\text{message} \in \text{Processed_messages}(\theta \text{Mach})$ $(\underline{\text{msg_contents}} \text{message}).\text{header}.\text{remote_port} \notin \underline{\text{port_exists}}$ $\underline{\text{message_exists}}' = \underline{\text{message_exists}} \setminus \{ \text{message} \}$

A message that has been processed can be queued at its destination port if that port exists and there is room in the message queue associated with the port or if the message was sent using a send-once right. The return status is as defined by the *error* component of the message if it is nonempty. Otherwise, the status is *Mm_success*.

The function *Error_to_status* converts an element of type *MSG_ERROR* to an element of type *MACH_MSG_RETURN*.

$\text{Error_to_status} : \mathbb{P} \text{MSG_ERROR} \rightarrow \text{MACH_MSG_RETURN}$ $\text{Error_to_status} = \{ \{ \text{Msg_error_invalid_memory} \} \mapsto \text{Mm_send_invalid_memory},$ $\{ \text{Msg_error_invalid_right} \} \mapsto \text{Mm_send_invalid_right},$ $\{ \text{Msg_error_invalid_type} \} \mapsto \text{Mm_send_invalid_type},$ $\{ \text{Msg_error_msg_too_small} \} \mapsto \text{Mm_send_msg_too_small} \}$
--

EnqueueMsg $\Delta \text{DtosExec}$ $\text{message} : \text{MESSAGE}$ $\text{msg_return!} : \text{MACH_MSG_RETURN}$ <hr/> $\text{message} \in \text{Processed_messages}(\theta \text{Mach})$ $(\text{let } \text{port} == (\underline{\text{msg_contents}} \text{message}).\text{header}.\text{remote_port} \bullet$ $\text{port} \in \underline{\text{port_exists}} \wedge$ $(\underline{\text{q_limit}} \text{port} > \text{port_size } \text{port} \vee$ $(\underline{\text{msg_contents}} \text{message}).\text{header}.\text{remote_rights} \in$ $\{ \text{Mmt_move_send_once}, \text{Mmt_make_send_once} \}) \wedge$ $\underline{\text{message_in_port_rel}}' = \underline{\text{message_in_port_rel}} \oplus \{ \text{port} \mapsto$ $\underline{\text{message_in_port_rel}} \text{port} \hat{\ } (\text{message}) \})$ $\text{msg_return!} =$ $(\text{let } \text{msg_error} == (\underline{\text{msg_contents}} \text{message}).\text{error} \bullet$ $\text{if } \text{msg_error} \neq \emptyset \text{ then } \text{Error_to_status } \text{msg_error}$ $\text{else } \text{Mm_success})$

If the following conditions hold:

- The message was sent using a send right rather than a send-once right.
- The client specified the *Mach_send_notify* option and the message either does not have a time out specified or the time out period has passed.
- The destination port exists and has a full message queue.
- No message is currently forcibly enqueued at the port.

then the message can be forcibly enqueued at the port.

If there already is a message forcibly queued at the port, then an error message is returned and a pseudo-receive is initiated. We represent that a pseudo receive has been initiated by changing the message status from Msg_stat_send to Msg_stat_pseudo .

If a time out was specified and the time out period has passed, then the message can time out with a pseudo receive operation being generated.

<p><i>MsgSendTimeOut</i></p> <p>Δ <i>DtosExec</i></p> <p><i>message</i> : MESSAGE</p> <p><i>msg_return!</i> : MACH_MSG_RETURN</p> <hr/> <p><i>message</i> \in <i>Processed_messages</i>(θMach)</p> <p>$\exists i : \mathbb{N} \mid i \in (\underline{msg_contents} \text{ message}).time_out_at \bullet$</p> <p>$i \leq \underline{host_time}$</p> <p><i>Mach_send_notify</i> \in (<i>msg_contents</i> message).option</p> <p>(<i>msg_contents'</i> message).option = (<i>msg_contents</i> message).option</p> <p>(<i>msg_contents'</i> message).time_out_at = (<i>msg_contents</i> message).time_out_at</p> <p>(<i>msg_contents'</i> message).status = <i>Msg_stat_pseudo</i></p> <p>(<i>msg_contents'</i> message).error =</p> <p style="padding-left: 40px;">if (<i>msg_contents</i> message).error = \emptyset then { <i>Msg_error_timed_out</i> }</p> <p style="padding-left: 40px;">else (<i>msg_contents</i> message).error</p>

C.2.2 Message Receive

The **mach_msg** request can be used to receive a message by including *Mach_rcv_msg* in *option?* and not including *Mach_send_msg*.

<p><i>MachMsgRcv</i></p> <p><i>MachMsgSignature</i></p> <hr/> <p><i>Mach_rcv_msg</i> \in <i>option?</i></p> <p><i>Mach_send_msg</i> \notin <i>option?</i></p>

C.2.2.1 Initial Processing We use the following schema to describe receive operations that are processed as no-ops due to error conditions that arise during the initial processing of the request:

<p><i>MachMsgRcvNoOp</i></p> <p>\exists <i>Dtos</i></p> <p><i>MachMsgRcv</i></p>

If *rcv_name?* does not denote a receive right or a port set for the client task, then an error message is returned and no further processing occurs.

<p><i>MachMsgRcvInvalidName</i></p> <p><i>MachMsgNoOp</i></p> <hr/> <p>(<i>owning_task client?</i>, <i>rcv_name?</i>) \notin (<i>r_right</i> \cup <i>port_set_namep</i>)</p> <p><i>msg_return!</i> = <i>Mm_rcv_invalid_name</i></p>

Otherwise, if $rcv_name?$ is a member of a port set, then an error message is returned and no further processing occurs.

<p><i>MachMsgRcvValidName</i></p> <p><i>Mach</i> <i>client?</i> : <i>THREAD</i> <i>rcv_name?</i> : <i>NAME</i></p> <p>$(owning_task\ client?, rcv_name?) \in r_right \cup port_set_namep$</p>

<p><i>MachMsgRcvInSet</i></p> <p><i>MachMsgRcvValidName</i> <i>MachMsgRcvNoOp</i></p> <p>$(owning_task\ client?, rcv_name?) \in port_set_namep$ $msg_return! = Mm_rcv_in_set$</p>

Otherwise, the request is queued at the end of the list of pending receives.

<p><i>MachMsgRcvNotInSet</i></p> <p><i>MachMsgRcvValidName</i></p> <p>$(owning_task\ client?, rcv_name?) \notin port_set_namep$</p>
--

<p><i>MachMsgRcvMakePending</i></p> <p><i>MachMsgRcv</i> <i>MachMsgRcvNotInSet</i></p> <p>$\forall p_rcv : PendingReceive \mid$ $p_rcv.notify = notify? \wedge$ $p_rcv.option = option? \wedge$ $p_rcv.rcv_size = rcv_size? \wedge$ $p_rcv.time_out_at = \mathbf{if} Mach_rcv_timeout \in option?$ $\quad \mathbf{then} \emptyset \mathbf{else} \{ time_out? + \underline{host_time} \} \bullet$ $\quad \underline{pending_receives}' = \underline{pending_receives} \oplus$ $\quad \{ (owning_task\ client?, rcv_name?) \mapsto$ $\quad \quad \underline{pending_receives}(owning_task\ client?, rcv_name?) \wedge \langle p_rcv \rangle \}$</p>

C.2.2.2 Kernel Processing Only the first request in the sequence associated with a port can be processed when a message is detected at the port. We introduce the following schema to denote processing of the first request.

<p><i>GeneralRcvProcessing</i></p> <p>$\Delta DtosExec$ $p_rcv : PendingReceive$ $task : TASK$ $name : NAME$</p> <p>$(task, name) \in local_namep \cap \text{dom } \underline{pending_receives}$ $\#(\underline{pending_receives}(task, name)) \neq 0$ $(\underline{pending_receives}(task, name))(1) = p_rcv$</p>

The components $task$ and $name$ denote the task that initiated the receive operation and the $name$ that $task$ specified as $rcv_name?$. We require that $task$ is an existing task, $name$ is a name that is in use in $task$'s name space, there is a sequence of pending receive requests associated with $(task, name)$, and the sequence of requests is nonempty. We introduce the component p_rcv to denote the first request in the sequence.

For a receive operation to be successful, the name specified by the client must either be a receive right or the name of a port set. The following schema defines the processing for the case in which the name is neither a receive right nor a port set:

$\frac{RcvPortDied}{GeneralRcvProcessing}$ $msg_return! : MACH_MSG_RETURN$
$(task, name) \notin (r_right \cup port_set_namep)$ $\underline{pending_receives}' = \underline{pending_receives} \oplus \{(task, name) \mapsto tail(\underline{pending_receives}(task, name))\}$ $msg_return! = Mm_rcv_port_died$

The request also fails if the specified name is a receive right that belongs to a port set.

$\frac{RcvPortChanged}{GeneralRcvProcessing}$ $msg_return! : MACH_MSG_RETURN$
$(task, name) \in r_right$ $\exists name_1 : NAME \bullet$ $(task, name_1) \in port_set_namep \wedge$ $named_port(task, name) \in port_set(task, name_1)$ $\underline{pending_receives}' = \underline{pending_receives} \oplus \{(task, name) \mapsto tail(\underline{pending_receives}(task, name))\}$ $msg_return! = Mm_rcv_port_changed$

The following schemas define the negation of the previous checks.

$\frac{GeneralRcvProcessing2}{GeneralRcvProcessing}$ $name_1 : NAME$ $port : PORT$ $message : MESSAGE$ $int_msg_1 : InternalMessage$
$(((\exists i : \mathbb{N} \bullet$ $(task, port, name, Receive, i) \in \underline{port_right_rel}) \wedge$ $name_1 = name) \vee$ $((task, name) \in port_set_namep \wedge$ $port \in port_set(task, name) \wedge$ $named_port(task, name_1) = port))$ $(Mach_rcv_large \notin p_rcv.option \vee$ $(\underline{msg_contents}(\underline{message_in_port_rel\ port\ 1}).header.size \leq p_rcv.rcv_size)$ $\#(\underline{message_in_port_rel\ port}) \neq 0$ $(\underline{message_in_port_rel\ port})(1) = message$ $int_msg_1 = \underline{msg_contents\ message}$

If *Mach_rcv_large* is specified in *option?* and the message to be received is larger than *rcv_size*, then an error message is returned, *msg!.rcv_size* is set to the size of the message, and no further processing occurs.

<p><i>MachMsgRcvTooLarge</i></p> <p><i>GeneralRcvProcessing2</i></p> <p><i>rcv_size!</i> : \mathbb{N}</p> <p><i>msg_return!</i> : <i>MACH_MSG_RETURN</i></p>
<p><i>Mach_rcv_large</i> \in <i>p_rcv.option</i></p> <p>$(\underline{msg_contents}(\underline{message_in_port_rel\ port\ 1}).header.size > p_rcv.rcv_size$</p> <p>$\underline{pending_receives}' = \underline{pending_receives} \oplus \{ (task, name) \mapsto$</p> <p>$\quad tail(\underline{pending_receives}(task, name)) \}$</p> <p>$rcv_size! = (\underline{msg_contents}(\underline{message_in_port_rel\ port\ 1}).header.size$</p> <p>$msg_return! = Mm_rcv_too_large$</p>

<p><i>GeneralRcvProcessing3</i></p> <p><i>GeneralRcvProcessing2</i></p> <p><i>int_msg₂</i> : <i>InternalMessage</i></p>
<p>$\underline{msg_contents}' = \underline{msg_contents} \oplus \{ message \mapsto int_msg_2 \}$</p> <p>$int_msg_2.header = int_msg_1.header$</p> <p>$int_msg_2.time_out_at = int_msg_1.time_out_at$</p>

Using this schema, we represent the initiation of the processing by setting the status of the message to *Msg_stat_rcv*.

<p><i>InitiateMsgRcv</i></p> <p><i>GeneralRcvProcessing3</i></p>
<p>$int_msg_1.body = int_msg_2.body$</p> <p>$int_msg_2.status = Msg_stat_rcv$</p> <p>$int_msg_2.error = \emptyset$</p>

Subsequent processing occurs only on messages having a status of *Msg_stat_rcv*. We use the following schema to represent processing of that form.

<p><i>GeneralRcvProcessing4</i></p> <p><i>GeneralRcvProcessing2</i></p>
<p>$int_msg_1.status = Msg_stat_rcv$</p>

<p><i>GeneralRcvProcessing5</i></p> <p><i>GeneralRcvProcessing3</i></p>
<p>$int_msg_1.status = Msg_stat_rcv$</p>

If the client did not specify *Mach_rcv_large* and the message is larger than the specified receive size, then the message is dequeued and destroyed.

The following schema denotes the dequeuing and destruction of a message:

$\begin{array}{l} \textit{DestroyMessage} \\ \Delta \textit{DtosExec} \\ \textit{message} : \textit{MESSAGE} \\ \textit{port} : \textit{PORT} \\ \textit{msg_return!} : \textit{MACH_MSG_RETURN} \\ \\ \textit{port} \in \text{dom } \underline{\textit{message_in_port_rel}} \\ \#(\underline{\textit{message_in_port_rel}} \textit{port}) \neq 0 \\ (\underline{\textit{message_in_port_rel}} \textit{port})(1) = \textit{message} \\ \underline{\textit{message_exists}}' = \underline{\textit{message_exists}} \setminus \{ \textit{message} \} \\ \underline{\textit{msg_contents}}' = \{ \textit{message} \} \triangleleft \underline{\textit{msg_contents}} \\ \underline{\textit{message_in_port_rel}}' = \underline{\textit{message_in_port_rel}} \oplus \{ \textit{port} \mapsto \\ \quad \textit{tail}(\underline{\textit{message_in_port_rel}} \textit{port}) \} \end{array}$
--

Using this schema, the processing of a message that is too large can be specified as follows:

$\begin{array}{l} \textit{MachMsgRcvTooLarge2} \\ \textit{GeneralRcvProcessing4} \\ \textit{DestroyMessage} \\ \\ \textit{Mach_rcv_large} \notin \textit{p_rcv.option} \\ \textit{int_msg}_1.\textit{header.size} > \textit{p_rcv.rcv_size} \\ \textit{msg_return!} = \textit{Mm_rcv_too_large} \end{array}$
--

If the client specified the *Mach_rcv_notify* option and the notify argument does not denote a valid receive right, then the processing is similar.

$\begin{array}{l} \textit{MachMsgReceiveInvalidNotify} \\ \textit{GeneralRcvProcessing4} \\ \textit{DestroyMessage} \\ \\ (\textit{Mach_rcv_large} \in \textit{p_rcv.option} \vee \\ \quad \textit{int_msg}_1.\textit{header.size} \leq \textit{p_rcv.rcv_size}) \\ \textit{Mach_rcv_notify} \in \textit{p_rcv.option} \\ (\textit{task}, \textit{p_rcv.notify}) \notin \textit{r_right} \\ \textit{msg_return!} = \textit{Mm_rcv_invalid_notify} \end{array}$

C.2.3 Notes

In this section we describe aspects of the **mach_msg** processing that are not addressed in the preceding section and issues concerning the correctness of the specification. The main gaps in the current specification are the kernel processing of receive and pseudo-receive requests.

The majority of this processing is concerned with transforming a message from type *InternalMessage* to type *Message*. To a large extent, this processing is simply the reverse of the processing described for the send request to transform a message from type *Message* to type *InternalMessage*. There do not appear to be any major obstacles to defining this “reverse” processing.

A more serious problem with completing the specification of receive requests is the concurrency present in the system. For example, it is not clear from the available documentation what happens if a task loses the receive right for a port while the kernel is in the middle of dequeuing a message from that port. The kernel interface document states that a status of *Mm_rcv_port_died* is returned to the client, but is unclear about other aspects of the processing. In particular, it is not clear from the documentation whether transferred port rights become visible only after the kernel commits to dequeuing the message.

A related problem is addressing side-effects of send and receive operations. For example, when a receive operation dequeues and destroys a message, receive rights for ports can be destroyed. This requires the modeling of the destruction of the ports and the generation of notification messages that must be sent.

Another problem that is common to all of the specifications in the FTLS is that input and output parameters are represented by value while they are actually implemented as references. In reality, the client specifies a virtual address rather than specifying a message header and body. The kernel assumes that the message header starts at the specified address and that the message body starts directly after the message header. One example of the ramifications of this simplification in the specification is that the specification does not address the case in which the sender of a message does not have access to the memory containing the header or body. In this case, the implementation treats the request as a no-op and returns a status of *mm_send_invalid_data*. A more complicated example is that the specification does not address the case in which the memory indicated by the virtual address is not resident. In this case, the kernel must enter a dialogue with a memory manager to determine the message header and body to use for the request.

Several subtle aspects of Mach are unclear from the available documentation. Examples include:

- It is not clear what the types of the *remote_rights* and *local_rights* fields of the message header are. The specifications models them both as sets of *MACH_MSG_TYPE*. This means that a send or send-once right must be transferred to the receiving task. Without examining the Mach source code, it is not possible to tell whether this is really how Mach works. A related question is whether more than one type of right can be passed at a time. If not, then the sets of *MACH_MSG_TYPE* should be constrained to having at most one element.
- When determining whether an area of out-of-line memory is accessible by a client, it is not clear whether the client's access to all of the pages comprising the region must be checked or whether it suffices to simply check access to the first page. The latter would be more efficient, but it requires that the kernel ensure that a client have the same access to all pages comprising a memory object. Although this property is desirable, it is not yet captured in the FTLS state.
- For simplicity, the current specification assumes that all of the pages containing port rights passed in out-of-line memory must be resident before the kernel can process the rights. It is possible that the implementation allows the kernel to process the rights contained on resident pages while waiting for the data containing the other rights. Modeling this capability would slightly complicate the model since provision would have to be made for data elements that are partially in-line and partially out-of-line.

Appendix *D*

Refinements of the State Model

In this appendix we refine portions of the state model to a lower level of detail. This models some of the data types and relationships that are used to implement the high-level abstract model described in the Basic Kernel State Definition and DTOS State Extensions chapters.

D.1 Additional Z Extensions

We define a function Gen_set to model generic queues. This function will be used in refining many of the components of the state model. A generic queue has a head element that points to the first element of a linked list of queue elements. $HEAD$ is the generic type of the head element of the queue, and $ELEM$ is the generic type of the elements of the queue. If $head_fnc$ maps the head of a queue to the first element of the queue and $next_fnc$ maps a queue element to its successor, then the expression $Gen_set(head_fnc, next_fnc)$ denotes a function mapping each element of $HEAD$ to the set of elements in its queue. We define a function Gen_seq to model generic sequences. This function will also be used in refining components of the state model. A generic sequence has a head element that points to the first element of a sequence of elements. The expression $Gen_seq(head_fnc, next_fnc)$ denotes a function mapping each element of $HEAD$ to its sequence of elements.¹⁸ Note that for certain values of $next_fnc$ $Gen_seq(head_fnc, next_fnc)$ may be infinite and therefore not of the type $seq\ ELEM$.

$$\begin{array}{l}
 \boxed{[HEAD, ELEM]} \\
 \hline
 Gen_set : (HEAD \leftrightarrow ELEM) \times (ELEM \leftrightarrow ELEM) \\
 \quad \rightarrow (HEAD \rightarrow \mathbb{P}\ ELEM) \\
 Gen_seq : (HEAD \leftrightarrow ELEM) \times (ELEM \leftrightarrow ELEM) \\
 \quad \rightarrow (HEAD \rightarrow (\mathbb{N}_1 \leftrightarrow ELEM)) \\
 \hline
 \forall head_fnc : HEAD \leftrightarrow ELEM; \\
 \quad next_fnc : ELEM \leftrightarrow ELEM; \\
 \quad head : HEAD \\
 \bullet Gen_seq(head_fnc, next_fnc)(head) \\
 \quad = \{ i : \mathbb{N}_1; e : ELEM \mid (head, e) \in head_fnc \ ; \ (next_fnc^{i-1}) \} \\
 \wedge Gen_set(head_fnc, next_fnc)(head) \\
 \quad = \text{ran}(Gen_seq(head_fnc, next_fnc)(head))
 \end{array}$$

D.2 Refinement of IPC Name Spaces

In refining the specification of IPC name spaces we introduce the following additional types:

$$\begin{array}{l}
 [IPC_SPACE, IPC_ENTRY, IPC_OBJECT, PORT_SET, IPC_SPLAY_TREE, \\
 \quad IPC_TREE_NODE] \\
 IPC_TABLE == \mathbb{N}_1 \leftrightarrow IPC_ENTRY
 \end{array}$$

¹⁸The Z expression $Q \ ; \ R$ denotes the composition of two relations with Q applied first followed by R . The expression R^k denotes the relation resulting from k applications of relation R . If $k = 0$, R^k denotes the identity relation. Thus, $Q \ ; \ (R^k)$ denotes one application of Q followed by k applications of R .

IPC_SPACE is the representation of a name space. Each name space consists of a set of elements of type IPC_ENTRY . Some of these entries point to an IPC_OBJECT which may be either a port or a port set. Note that we are introducing an explicit given type for port sets rather than representing them merely as a set of ports. This agrees with the prototype and makes it easier to model properties of port sets such as the message queue of a port set. The entries in a space are organized into two data structures, an IPC_TABLE and an IPC_SPLAY_TREE . An IPC_TABLE is simply a sequence of IPC_ENTRY that may have gaps in it. A splay tree is a search tree containing nodes of type IPC_TREE_NODE . Each IPC_TREE_NODE points to an IPC_ENTRY .

The expression $\underline{task_space}(tk)$ denotes the IPC_SPACE associated with task tk . No two tasks have the same value for this expression. The set \underline{spacep} denotes the existing IPC name spaces. The expression $\underline{space_table}(sp)$ denotes the IPC_TABLE associated with space sp , and $\underline{space_table_size}(sp)$ denotes the current maximum size of this table. Note that this value may change dynamically to improve performance and memory utilization. The expression $\underline{space_tree}(sp)$ denotes the IPC_SPLAY_TREE associated with space sp . Every space has both a table and a splay tree although one or both of these could be empty.

$\underline{IpcSpace}$ $\underline{TaskExist}$ $\underline{task_space} : TASK \mapsto IPC_SPACE$ $\underline{spacep} : \mathbb{P} IPC_SPACE$ $\underline{space_table} : IPC_SPACE \mapsto IPC_TABLE$ $\underline{space_table_size} : IPC_SPACE \mapsto \mathbb{N}$ $\underline{space_tree} : IPC_SPACE \mapsto IPC_SPLAY_TREE$
$\text{dom } \underline{task_space} = \underline{task_exists}$ $\text{ran } \underline{task_space} = \underline{spacep}$ $\text{dom } \underline{space_table} = \text{dom } \underline{space_table_size} = \text{dom } \underline{space_tree} = \underline{spacep}$ $\forall sp : \underline{spacep} \bullet \max(\text{dom}(\underline{space_table} sp)) < \underline{space_table_size}(sp)$

We augment the set of rights to contain $Port_set_right$ and $Dead_name_right$. The former is the right associated with an entry for a port set, and the latter is a right that may be associated with a dead name. Because the marking of dead rights in Mach is performed lazily, a dead right need not be marked $Dead_name_right$. It is also recognized as dead if it points to an inactive IPC_OBJECT .

$$ALL_RIGHTS ::= Right_for_port\langle\langle RIGHT \rangle\rangle \mid Port_set_right \mid Dead_name_right$$

$\underline{Receive_right}, \underline{Send_right}, \underline{Send_once_right} : ALL_RIGHTS$
$\underline{Receive_right} = Right_for_port(Receive)$ $\underline{Send_right} = Right_for_port(Send)$ $\underline{Send_once_right} = Right_for_port(Send_once)$

The set \underline{entryp} denotes the existing IPC_ENTRY elements. An entry is marked with a generation that is used in determining whether it is out of date. The expression $\underline{entry_gen}(entry)$ denotes the generation of $entry$. The expression $\underline{entry_object}(entry)$ denotes the IPC_OBJECT associated with $entry$. The expression $\underline{entry_rights}(entry)$ denotes the set of ALL_RIGHTS associated with $entry$. Finally, the expression $\underline{entry_count}(entry)$ denotes the number of send rights denoted by $entry$ when a name denotes multiple send rights for a task. If $\underline{entry_count}(entry)$ is positive, then $entry$ must denote a $Send_right$, $Send_once_right$ or $Dead_name_right$.

Editorial Note:

Should probably add dead name notification requests to this.

<p><i>IpEntry</i></p> <p>$\underline{entryp} : \mathbb{P} IPC_ENTRY$ $\underline{entry_gen} : IPC_ENTRY \leftrightarrow \mathbb{N}$ $\underline{entry_object} : IPC_ENTRY \leftrightarrow IPC_OBJECT$ $\underline{entry_rights} : IPC_ENTRY \leftrightarrow \mathbb{P} ALL_RIGHTS$ $\underline{entry_count} : IPC_ENTRY \leftrightarrow \mathbb{N}$</p> <p>$\text{dom } \underline{entry_gen} = \text{dom } \underline{entry_object} = \text{dom } \underline{entry_rights} = \text{dom } \underline{entry_count}$ $= \underline{entryp}$ $\forall \text{entry} : IPC_ENTRY$ $\underline{entry_count}(\text{entry}) > 0$ $\bullet \underline{entry_rights}(\text{entry}) \cap \{Send_right, Send_once_right, Dead_name_right\} \neq \emptyset$</p>
--

Every entry in the table associated with a space must denote some right.

<p><i>IpTableEntry</i></p> <p><i>IpSpace</i></p> <p><i>IpEntry</i></p> <p>$\forall \text{entry} : IPC_ENTRY; \text{table} : IPC_TABLE$ $\text{table} \in \text{ran } \text{space_table} \wedge \text{entry} \in \text{ran } \text{table}$ $\bullet \underline{entry_rights}(\text{entry}) \neq \emptyset$</p>
--

The set $\underline{objectp}$ denotes the existing IPC_OBJECT elements. An existing object may be inactive. The active objects are denoted by $\underline{active_objects}$. An IPC_OBJECT may be either a port or a port set. The expressions $\underline{object_as_port}(\text{obj})$ and $\underline{object_as_port_set}(\text{obj})$ denote the associated port or port set. The domains of these two functions partition the set of existing objects.

<p><i>IpObject</i></p> <p>$\underline{objectp} : \mathbb{P} IPC_OBJECT$ $\underline{active_objects} : \mathbb{P} IPC_OBJECT$ $\underline{object_as_port} : IPC_OBJECT \rightsquigarrow PORT$ $\underline{object_as_port_set} : IPC_OBJECT \rightsquigarrow PORT_SET$</p> <p>$\underline{active_objects} \subseteq \underline{objectp}$ $(\text{dom } \underline{object_as_port}, \text{dom } \underline{object_as_port_set})$ $\text{partition } \underline{objectp}$</p>
--

A port port is in a port set P if and only if $(\text{port}, P) \in \underline{port_in_set}$.

<p><i>PortInSet</i></p> <p><i>PortExist</i></p> <p>$\underline{port_setp} : \mathbb{P} PORT_SET$ $\underline{port_in_set} : PORT \leftrightarrow PORT_SET$</p> <p>$\text{dom } \underline{port_in_set} \subseteq \underline{port_exists}$ $\text{ran } \underline{port_in_set} \subseteq \underline{port_setp}$</p>
--

Each *NAME* n encodes an index denoted $\underline{n_name_index}(n)$ and a generation denoted $\underline{n_name_gen}(n)$. If names n_1 and n_2 have the same index and generation, then they are the same name.

<i>Name</i>
$\underline{n_name_index} : NAME \rightarrow \mathbb{N}$ $\underline{n_name_gen} : NAME \rightarrow \mathbb{N}$
$\forall n_1, n_2 : NAME$ $ \underline{n_name_index}(n_1) = \underline{n_name_index}(n_2)$ $\quad \wedge \underline{n_name_gen}(n_1) = \underline{n_name_gen}(n_2)$ $\bullet n_1 = n_2$

The set $\underline{splay_treep}$ denotes the existing splay trees. For efficiency of lookup, a splay tree is represented internally by three (possibly empty) tree structures, a left, a right and a middle tree. The root *IPC_TREE_NODE* of each of these trees (when the tree is nonempty) is denoted, respectively, by $\underline{tree_left}(splay)$, $\underline{tree_right}(splay)$ and $\underline{tree_middle}(splay)$. A pair $(splay, node)$ is in $\underline{tree_trees}$ if and only if $node$ is the root of one of the three trees associated with $splay$.

<i>IpcSplayTree</i>
$\underline{splay_treep} : \mathbb{P} IPC_SPLAY_TREE$ $\underline{tree_middle} : IPC_SPLAY_TREE \leftrightarrow IPC_TREE_NODE$ $\underline{tree_left} : IPC_SPLAY_TREE \leftrightarrow IPC_TREE_NODE$ $\underline{tree_right} : IPC_SPLAY_TREE \leftrightarrow IPC_TREE_NODE$ $\underline{tree_trees} : IPC_SPLAY_TREE \leftrightarrow IPC_TREE_NODE$
$\text{dom } \underline{tree_trees} \subseteq \underline{splay_treep}$ $\underline{tree_trees} = \underline{tree_middle} \cup \underline{tree_left} \cup \underline{tree_right}$

The set $\underline{tree_nodep}$ denotes the set of existing *IPC_TREE_NODE* elements. Each tree node $node$ points to an *IPC_ENTRY* which is denoted $\underline{tree_node_entry}(node)$. The expression $\underline{tree_node_name}(node)$ denotes a *NAME* associated with $node$. Each tree node may have a left and a right child tree node. These are denoted by the expressions $\underline{tree_node_lchild}(node)$ and $\underline{tree_node_rchild}(node)$. A pair $(node_1, node_2)$ is in $\underline{tree_node_children}$ if and only if $node_2$ is either the left or right child of $node_1$.

<i>IpcTreeNode</i>
$\underline{tree_nodep} : \mathbb{P} IPC_TREE_NODE$ $\underline{tree_node_entry} : IPC_TREE_NODE \leftrightarrow IPC_ENTRY$ $\underline{tree_node_name} : IPC_TREE_NODE \leftrightarrow NAME$ $\underline{tree_node_lchild} : IPC_TREE_NODE \leftrightarrow IPC_TREE_NODE$ $\underline{tree_node_rchild} : IPC_TREE_NODE \leftrightarrow IPC_TREE_NODE$ $\underline{tree_node_children} : IPC_TREE_NODE \leftrightarrow IPC_TREE_NODE$
$\text{dom } \underline{tree_node_entry} = \text{dom } \underline{tree_node_name} = \underline{tree_nodep}$ $\text{dom } \underline{tree_node_lchild} \subseteq \underline{tree_nodep}$ $\text{dom } \underline{tree_node_rchild} \subseteq \underline{tree_nodep}$ $\underline{tree_node_children} = \underline{tree_node_lchild} \cup \underline{tree_node_rchild}$

We are now ready to define the relations $\underline{port_right_rel}$, $\underline{port_set_rel}$, and $\underline{dead_right_rel}$. All three share the requirement that the space of the task must contain an entry with the appropriate name. If the entry is in the table, then the index of the name must match the position

of the name in the table, and the generations of the name and entry must be identical. If the entry is in the splay tree, then the name in the tree node must equal the given name. This requirement is abstracted by the relation $entry_in_space$. A triple $(task, name, entry)$ is in this relation if and only if it satisfies the above requirement.

$EntryInSpace$ $TaskExist$ $IpC_TableEntry$ $Name$ $IpC_SplayTree$ $IpC_TreeNode$ $entry_in_space : \mathbb{P}(TASK \times NAME \times IPC_ENTRY)$
$\forall tk : TASK; n : NAME; entry : IPC_ENTRY$ <ul style="list-style-type: none"> • $(tk, n, entry) \in entry_in_space$ $\Leftrightarrow (tk \in \underline{task_exists}$ $\wedge entry \in \underline{entryp}$ $\wedge (((\underline{name_index}(n), entry) \in \underline{space_table}(\underline{task_space}(tk))$ $\wedge \underline{name_gen}(n) = \underline{entry_gen}(entry))$ $\vee (\exists splay : IPC_SPLAY_TREE;$ $\quad \underline{tree_nodes} : seq_1 IPC_TREE_NODE$ <ul style="list-style-type: none"> • $(\underline{task_space}(tk), splay) \in \underline{space_tree}$ $\wedge (splay, \underline{tree_nodes}(1)) \in \underline{tree_trees}$ $\wedge (\underline{last\ tree_nodes}, entry) \in \underline{tree_node_entry}$ $\wedge (\underline{last\ tree_nodes}, n) \in \underline{tree_node_name}$ $\wedge (\forall i : 2 \dots \#\underline{tree_nodes}$ <ul style="list-style-type: none"> • $(\underline{tree_nodes}(i - 1), \underline{tree_nodes}(i)) \in \underline{tree_node_children}))$

A 5-tuple $(tk, p, n, r, count)$ is in $\underline{port_right_rel}$ if and only if there exists an IPC_ENTRY , $entry$ such that

- $(tk, n, entry)$ is in $entry_in_space$,
- r is one of the rights associated with the entry all of which are rights to use a port (i.e., not port set rights nor dead rights),
- an active object is associated with the entry,
- the object is a port, and
- either
 - r is a receive right and $count$ is 1, or
 - r is not a receive right and $count$ is the right count of the entry.

<p><i>PortRightRefinement</i></p> <p><i>TasksAndPorts</i></p> <p><i>EntryInSpace</i></p> <p><i>IpcObject</i></p>
$\forall tk : TASK; p : PORT; n : NAME; r : RIGHT; count : \mathbb{N}_1$ <ul style="list-style-type: none"> • $(tk, p, n, r, count) \in \underline{port_right_rel}$ $\Leftrightarrow (\exists entry : IPC_ENTRY$ <ul style="list-style-type: none"> • $(tk, n, entry) \in \underline{entry_in_space}$ $\wedge \underline{Right_for_port}(r) \in \underline{entry_rights}(entry) \subseteq \text{ran } \underline{Right_for_port}$ $\wedge \underline{entry_object}(entry) \in \underline{active_objects}$ $\wedge (\underline{entry_object}(entry), p) \in \underline{object_as_port}$ $\wedge ((r = \text{Receive} \wedge count = 1)$ <li style="padding-left: 2em;">$\vee (r \neq \text{Receive} \wedge count = \underline{entry_count}(entry)))$

Editorial Note:

This is to cover up the name vs. port discrepancy in version 1.13 of the FTLS. When we incorporate this refinement is a new version of the FTLS in which this discrepancy is fixed, we must remove this schema and replace $\underline{new_port_set}$ below with $\underline{port_set}$.

<p><i>NewPortSets</i></p> <p><i>TasksAndRights</i></p> <p>$\underline{new_port_set} : (TASK \times NAME) \leftrightarrow \mathbb{P} PORT$</p>
--

We define $\underline{port_set_rel}$ indirectly by defining $\underline{new_port_set}$. A port $port$ is in $\underline{new_port_set}(tk, n)$ if and only if there exists an IPC_ENTRY , $entry$ and a port set PS such that

- $(tk, n, entry)$ is in $\underline{entry_in_space}$,
- the right associated with the entry is $\underline{Port_set_right}$,
- an active object is associated with the entry,
- the object is PS and
- $port$ is an element of PS .

<p><i>PortSetRefinement</i></p> <p><i>NewPortSets</i></p> <p><i>EntryInSpace</i></p> <p><i>IpcObject</i></p> <p><i>PortInSet</i></p>
$\forall tk : TASK; port : PORT; n : NAME$ <ul style="list-style-type: none"> • $port \in \underline{new_port_set}(tk, n)$ $\Leftrightarrow (\exists entry : IPC_ENTRY; PS : PORT_SET$ <ul style="list-style-type: none"> • $(tk, n, entry) \in \underline{entry_in_space}$ $\wedge \underline{entry_rights}(entry) = \{\underline{Port_set_right}\}$ $\wedge \underline{entry_object}(entry) \in \underline{active_objects}$ $\wedge (\underline{entry_object}(entry), PS) \in \underline{object_as_port_set}$ $\wedge (port, PS) \in \underline{port_in_set}$

A triple $(tk, n, count)$ is in $\underline{dead_right_rel}$ if and only if there exists an IPC_ENTRY , $entry$ such that

- $(tk, n, entry)$ is in $entry_in_space$,
- either the rights associated with the entry include $Dead_name_right$, or the entry points to an inactive object,
- $count$ is the right count of the entry,

$\underline{DeadRightRefinement}$ $DeadRights$ $EntryInSpace$ $IpObject$
$\forall tk : TASK; n : NAME; count : \mathbb{N}_1$ • $(tk, n, count) \in \underline{dead_right_rel}$ $\Leftrightarrow (\exists entry : IPC_ENTRY$ • $(tk, n, entry) \in entry_in_space$ $\wedge (Dead_name_right \in \underline{entry_rights}(entry)$ $\vee \underline{entry_object}(entry) \notin \underline{active_objects})$ $\wedge count = \underline{entry_count}(entry))$

The schema $IpRefinement$ defines the refinements for IPC name spaces.

$\underline{IpRefinement}$ $PortRightRefinement$ $PortSetRefinement$ $DeadRightRefinement$

D.3 Refinement of Pending Receives

The expression $\underline{port_waiting_threads_head}(port)$ denotes the first $THREAD$, if one exists, in the sequence of threads waiting to receive a message from $port$. The expression $\underline{port_set_waiting_threads_head}(pset)$ denotes the first $THREAD$, if one exists, in the sequence of threads waiting to receive a message from port set $pset$. The expression $\underline{next_waiting_thread}(th)$ denotes the successor $THREAD$ of th , if one exists, in the sequence of threads waiting to receive a message from the port or port set from which th is waiting to receive a message.

$\underline{WaitingThreads}$ $ThreadExist$ $PortExist$ $PortInSet$ $\underline{port_waiting_threads_head} : PORT \leftrightarrow THREAD$ $\underline{port_set_waiting_threads_head} : PORT_SET \leftrightarrow THREAD$ $\underline{next_waiting_thread} : THREAD \leftrightarrow THREAD$
$\text{dom } \underline{port_waiting_threads_head} \subseteq \underline{port_exists}$ $\text{dom } \underline{port_set_waiting_threads_head} \subseteq \underline{port_setp}$ $\text{dom } \underline{next_waiting_thread} \subseteq \underline{thread_exists}$

The expression $\underline{thread_pending_receive}(th)$ denotes the *PendingReceive* data stored in a thread th that is waiting to receive a message.

<p><i>StoredReceiveState</i></p> <p><i>ThreadExist</i></p> <p>$\underline{thread_pending_receive} : THREAD \leftrightarrow PendingReceive$</p> <p>$\text{dom } \underline{thread_pending_receive} \subseteq \underline{thread_exists}$</p>

The expression $\underline{named_port_set}(tk, nm)$ denotes the port set named by nm in the IPC name space of task tk .

Editorial Note:
It might make more sense to have this in the regular state, not the refinements. Since it is here, we will refine it right away.

<p><i>TasksAndPortSets</i></p> <p><i>EntryInSpace</i></p> <p><i>IpcObject</i></p> <p>$\underline{named_port_set} : (TASK \times NAME) \leftrightarrow PORT_SET$</p> <p>$\forall tk : TASK; n : NAME; ps : PORT_SET$</p> <ul style="list-style-type: none"> • $ps = \underline{named_port_set}(tk, n)$ <p>$\Leftrightarrow (\exists entry : IPC_ENTRY$</p> <ul style="list-style-type: none"> • $(tk, n, entry) \in \underline{entry_in_space}$ $\wedge \underline{entry_rights}(entry) = \{Port_set_right\}$ $\wedge \underline{entry_object}(entry) \in \underline{active_objects}$ $\wedge (\underline{entry_object}(entry), ps) \in \underline{object_as_port_set}$

The expression $\underline{waiting_for_port}(tk, nm)$ denotes the sequence of threads that are waiting for a message on the port named by nm in the IPC name space of tk . Note that $\underline{Gen_seq}(\underline{port_waiting_threads_head}, \underline{next_waiting_thread})$ denotes a function of type $PORT \leftrightarrow \text{seq } THREAD$ where a thread th is in the sequence associated with a port p if and only if th is waiting to receive a message from p .¹⁹ The expression $\underline{waiting_for_port_set}(tk, nm)$ denotes the sequence of threads that are waiting for a message on the port set named by nm in the IPC name space of tk . A name may not name both a port and a port set for the same task. Thus, the domains of $\underline{waiting_for_port}$ and $\underline{waiting_for_port_set}$ are disjoint. For convenience, we define $\underline{waiting_for_message}$ to be the union of the functions $\underline{waiting_for_port}$ and $\underline{waiting_for_port_set}$. Because the two domains are disjoint, $\underline{waiting_for_message}$ is necessarily a function. Every thread that is waiting for a message holds *PendingReceive* information.

¹⁹The Z expression $Q ; R$ denotes the composition of two relations with Q applied first followed by R .

<p><i>TasksAndWaiting</i></p> <p><i>WaitingThreads</i></p> <p><i>TasksAndPorts</i></p> <p><i>TasksAndPortSets</i></p> <p><i>StoredReceiveState</i></p> <p>$\underline{waiting_for_port} : TASK \times NAME \leftrightarrow \text{seq } THREAD$</p> <p>$\underline{waiting_for_port_set} : TASK \times NAME \leftrightarrow \text{seq } THREAD$</p> <p>$waiting_for_message : TASK \times NAME \leftrightarrow \text{seq } THREAD$</p>
<p>$\underline{waiting_for_port} = \text{named_port}$</p> <p> $\S Gen_seq(\underline{port_waiting_threads_head}, \underline{next_waiting_thread})$</p> <p>$\underline{waiting_for_port_set} = \text{named_port_set}$</p> <p> $\S Gen_seq(\underline{port_set_waiting_threads_head}, \underline{next_waiting_thread})$</p> <p>$waiting_for_message = \underline{waiting_for_port} \cup \underline{waiting_for_port_set}$</p> <p>$\forall tk : TASK; n : NAME$</p> <p>• $\text{ran}(waiting_for_message(tk, n)) \subseteq \text{dom } \underline{thread_pending_receive}$</p>

We now refine the definition of *pending_receives*. For any (tk, n) pair the sequence of *PendingReceive* values associated with name n for task tk is found by extracting (via *thread_pending_receive*) the *PendingReceive* data from each thread in the sequence *waiting_for_message*(tk, n).

<p><i>PendingReceiveRefinement</i></p> <p><i>TasksAndWaiting</i></p> <p><i>Messages</i></p>
<p>$\forall tk : TASK; n : NAME$</p> <p>• $\underline{pending_receives}(tk, n)$</p> <p> $= \text{waiting_for_message}(tk, n) \S \underline{thread_pending_receive}$</p>

D.4 Refinement of Virtual Memory

In refining the specification of virtual memory we introduce the following additional types:

[*VM_MAP*, *VM_ENTRY*, *VM_MAP_OBJECT*]

VM_MAP is the representation of a virtual address space. Each map consists of a sequence of elements of type *VM_ENTRY*. A *VM_ENTRY* denotes a contiguous range of virtual addresses that share the same properties (e.g., protections and inheritance options). Some of these entries point to a *VM_MAP_OBJECT* which may be either a memory object or a another memory map called a submap.

The expression $\underline{task_map}(tk)$ denotes the *VM_MAP* associated with task tk . Tasks running in kernel space may have the same map (the kernel map). No two kernel-external tasks have the same value for $\underline{task_map}$. The set $\underline{map_exists}$ denotes the existing VM maps. The expression $\underline{vm_entries_head}(map)$ denotes the first *VM_ENTRY*, if one exists, in the sequence of entries associated with map .

$\begin{array}{l} \text{VmMapStructure} \\ \text{TaskExist} \\ \underline{task_map} : TASK \leftrightarrow VM_MAP \\ \underline{map_exists} : \mathbb{P} VM_MAP \\ \underline{vm_entries_head} : VM_MAP \leftrightarrow VM_ENTRY \\ \hline \text{dom } \underline{task_map} = \underline{task_exists} \\ \text{dom } \underline{vm_entries_head} \subseteq \underline{map_exists} \end{array}$

The set $\underline{vm_entry_exists}$ denotes the existing VM_ENTRY elements, and the set $\underline{vm_entry_submap_p}$ denotes the entries that are submaps.²⁰ The following functions are defined on VM_ENTRY :

- $\underline{next_vm_entry}(e)$ — denotes the next entry after e , if there is one, in the sequence of VM entries associated with some VM map,
- $\underline{vm_entry_start}(e)$ — denotes the starting address of e ,
- $\underline{vm_entry_end}(e)$ — denotes the first address after the end of e ,
- $\underline{vm_entry_map_object}(e)$ — denotes the VM_MAP_OBJECT associated with e ,
- $\underline{vm_entry_offset}(e)$ — denotes the offset at which e is mapped into a memory object,
- $\underline{vm_entry_prot}(e)$ — denotes the current protections associated with e ,
- $\underline{vm_entry_max_prot}(e)$ — denotes the maximum protections that e may take,
- $\underline{vm_entry_inh}(e)$ — denotes the inheritance option in effect for e ,
- $\underline{vm_entry_wire_count}(e)$ — denotes the number of times that e has been wired, and
- $\underline{vm_entry_sid}(e)$ — denotes the OSI associated with e .

Every existing entry has a start and an end address, an offset, a protection, a maximum protection, an inheritance option and a wire count. Every entry that is a submap has an associated map object while an entry that is not a submap might not have any associated VM_MAP_OBJECT . For convenience we define $\underline{map_entries}(map)$ to denote the set of VM_ENTRY contained in map .

²⁰The only map known to have submaps is `kernel_map`. It has the following submaps:

- `device_io_map`,
- `ipc_kernel_map`,
- `kalloc_map`,
- `zone_map`, and
- the map of any task running in kernel space that does not use the entire kernel map.

$\begin{aligned} & \text{VmEntry} \\ & \text{VmMapStructure} \\ & \underline{vm_entry_exists} : \mathbb{P} \text{ VM_ENTRY} \\ & \underline{vm_entry_submap_p} : \mathbb{P} \text{ VM_ENTRY} \\ & \underline{next_vm_entry} : \text{ VM_ENTRY} \leftrightarrow \text{ VM_ENTRY} \\ & \underline{vm_entry_start} : \text{ VM_ENTRY} \leftrightarrow \text{ PAGE_INDEX} \\ & \underline{vm_entry_end} : \text{ VM_ENTRY} \leftrightarrow \text{ PAGE_INDEX} \\ & \underline{vm_entry_map_object} : \text{ VM_ENTRY} \leftrightarrow \text{ VM_MAP_OBJECT} \\ & \underline{vm_entry_offset} : \text{ VM_ENTRY} \leftrightarrow \text{ OFFSET} \\ & \underline{vm_entry_prot} : \text{ VM_ENTRY} \leftrightarrow \mathbb{P} \text{ PROTECTION} \\ & \underline{vm_entry_max_prot} : \text{ VM_ENTRY} \leftrightarrow \mathbb{P} \text{ PROTECTION} \\ & \underline{vm_entry_inh} : \text{ VM_ENTRY} \leftrightarrow \text{ INHERITANCE_OPTION} \\ & \underline{vm_entry_wire_count} : \text{ VM_ENTRY} \leftrightarrow \mathbb{N} \\ & \underline{vm_entry_sid} : \text{ VM_ENTRY} \leftrightarrow \text{ OSI} \\ & \underline{map_entries} : \text{ VM_MAP} \rightarrow \mathbb{P} \text{ VM_ENTRY} \\ \\ & \text{dom } \underline{vm_entry_start} = \text{dom } \underline{vm_entry_end} = \text{dom } \underline{vm_entry_offset} \\ & \quad = \text{dom } \underline{vm_entry_prot} = \text{dom } \underline{vm_entry_max_prot} = \text{dom } \underline{vm_entry_inh} \\ & \quad = \text{dom } \underline{vm_entry_wire_count} = \text{dom } \underline{vm_entry_sid} = \underline{vm_entry_exists} \\ & \text{dom } \underline{next_vm_entry} \subseteq \underline{vm_entry_exists} \\ & \underline{vm_entry_submap_p} \subseteq \text{dom } \underline{vm_entry_map_object} \subseteq \underline{vm_entry_exists} \\ & \underline{map_entries} = \text{Gen_set}(\underline{vm_entries_head}, \underline{next_vm_entry}) \end{aligned}$
--

The set $\underline{vm_map_object_exists}$ denotes the existing VM_MAP_OBJECT elements. A VM_MAP_OBJECT may be either a memory object or another VM map. The expressions $\underline{map_object_as_memory}(obj)$ and $\underline{map_object_as_submap}(obj)$ denote the associated memory object or VM map. The domains of these two functions partition the set of existing map objects. The function $\underline{vm_entry_map_object}$ maps a submap entry to an element of the domain of $\underline{map_object_as_submap}$, and it maps other entries to elements of the domain of $\underline{map_object_as_memory}$. For convenience, we define the functions $\underline{vm_entry_memory}$ and $\underline{vm_entry_submap}$ as the compositions of $\underline{vm_entry_map_object}$ with $\underline{map_object_as_memory}$ and $\underline{map_object_as_submap}$, respectively.

$\begin{aligned} & \text{VmMapObject} \\ & \text{VmEntry} \\ & \underline{vm_map_object_exists} : \mathbb{P} \text{ VM_MAP_OBJECT} \\ & \underline{map_object_as_memory} : \text{ VM_MAP_OBJECT} \rightsquigarrow \text{ MEMORY} \\ & \underline{map_object_as_submap} : \text{ VM_MAP_OBJECT} \rightsquigarrow \text{ VM_MAP} \\ & \underline{vm_entry_memory} : \text{ VM_ENTRY} \leftrightarrow \text{ MEMORY} \\ & \underline{vm_entry_submap} : \text{ VM_ENTRY} \leftrightarrow \text{ VM_MAP} \\ \\ & (\text{dom } \underline{map_object_as_memory}, \text{dom } \underline{map_object_as_submap}) \\ & \quad \text{partition } \underline{vm_map_object_exists} \\ & \text{dom } \underline{map_object_as_memory} = \text{ran}(\underline{vm_entry_submap_p} \triangleleft \underline{vm_entry_map_object}) \\ & \text{dom } \underline{map_object_as_submap} = \text{ran}(\underline{vm_entry_submap_p} \triangleleft \underline{vm_entry_map_object}) \\ & \underline{vm_entry_memory} = \underline{vm_entry_map_object} \ ; \ \underline{map_object_as_memory} \\ & \underline{vm_entry_submap} = \underline{vm_entry_map_object} \ ; \ \underline{map_object_as_submap} \end{aligned}$
--

We define a global function $\underline{Page_index_int}$ that maps a PAGE_INDEX to a non-negative integer. This allows the numeric comparison of page addresses.

| $\underline{Page_index_int} : \text{ PAGE_INDEX} \rightsquigarrow \mathbb{N}$

For use in refining the model of the VM system, we define the functions vm_map_lookup and $vm_map_lookup_entry$ which each map a $(tk, pindex)$ pair to a VM_ENTRY . For vm_map_lookup , the pair is mapped to $entry$ if and only if $entry$ is not a submap, and there is a non-empty sequence $lookup_seq$ of $VM_MAP \times VM_ENTRY$ pairs such that

- $task_map(tk)$ is the first component of the first element of $lookup_seq$,
- $entry$ is the second component of the last element of $lookup_seq$,
- for each element (m, e) of $lookup_seq$:
 - e is in the set of entries for m ,
 - $pindex$ is in the address range defined by the start and end addresses of e ,
 - if (m, e) is not the last element of $lookup_seq$, then e is a submap entry with the first component of the next pair in the sequence as its submap.

For $vm_map_lookup_entry$, a $(tk, pindex)$ pair is mapped to $entry$ if and only if $entry$ is in the set of entries for $task_map(tk)$, and $pindex$ is in the address range defined by the start and end addresses of $entry$,

$VmLookup$ $VmMapStructure$ $VmEntry$ $VmMapObject$ $vm_map_lookup_entry : TASK \times PAGE_INDEX \mapsto VM_ENTRY$ $vm_map_lookup : TASK \times PAGE_INDEX \mapsto VM_ENTRY$
$\forall entry : VM_ENTRY; tk : TASK; pindex : PAGE_INDEX$ <ul style="list-style-type: none"> ● $entry = vm_map_lookup_entry(tk, pindex)$ $\Leftrightarrow (\exists map : VM_MAP$ <ul style="list-style-type: none"> ● $task_map(tk) = map$ $\wedge entry \in map_entries(map)$ $\wedge Page_index_int(vm_entry_start(entry)) \leq Page_index_int(pindex)$ $< Page_index_int(vm_entry_end(entry))$)
$\forall entry : VM_ENTRY; tk : TASK; pindex : PAGE_INDEX$ <ul style="list-style-type: none"> ● $entry = vm_map_lookup(tk, pindex)$ $\Leftrightarrow (entry \notin vm_entry_submap_p$ <ul style="list-style-type: none"> $\wedge (\exists lookup_seq : seq_1(VM_MAP \times VM_ENTRY)$ <ul style="list-style-type: none"> ● $task_map(tk) = first(head\ lookup_seq)$ $\wedge entry = second(last\ lookup_seq)$ $\wedge (\forall i : 1.. \#lookup_seq; e : VM_ENTRY$ <ul style="list-style-type: none"> $e = second(lookup_seq(i))$ ● $e \in map_entries(first(lookup_seq(i)))$ $\wedge Page_index_int(vm_entry_start(e)) \leq Page_index_int(pindex)$ $< Page_index_int(vm_entry_end(e))$ $\wedge (i < \#lookup_seq$ <ul style="list-style-type: none"> $\Rightarrow (e, first(lookup_seq(i + 1))) \in vm_entry_submap))))$

Now, we define $mapped_memory$ and $mapped_offset$ by composing vm_map_lookup with vm_entry_memory and vm_entry_offset , respectively. We define $mapped_offset$, $protection$, $max_protection$, $inheritance$ and $wire_count$ by composing $vm_map_lookup_entry$ with the appropriate VM entry functions.

Editorial Note:

It is unclear whether *mapped_memory* and *mapped_offset* are best defined as below or whether it would be better to use *vm_map_lookup_entry* for them as well. The question is whether, when *index* denotes a submap for task *tk*, the pair $(tk, index)$ should be in the domains of *mapped_memory* and *mapped_offset*. The prototype appears to follow the submap link when dealing with page faults. However, when accessing and returning state information associated with a region, it does not look at the submap. Furthermore, *vm_region* returns a null name for the memory object when the address leads to a submap.

*VmRefinement**VmLookup**AddressSpace**Protection**Inheritance**Wired**PageSid**mapped_memory* = *vm_map_lookup* ; *vm_entry_memory**mapped_offset* = *vm_map_lookup* ; *vm_entry_offset**protection* = *vm_map_lookup_entry* ; *vm_entry_prot**max_protection* = *vm_map_lookup_entry* ; *vm_entry_max_prot**inheritance* = *vm_map_lookup_entry* ; *vm_entry_inh**wire_count* = *vm_map_lookup_entry* ; *vm_entry_wire_count**page_sid* = *vm_map_lookup_entry* ; *vm_entry_sid***D.5 Miscellaneous Refinements**

The expression *threads_head(tk)* denotes the first *THREAD*, if one exists, in the sequence of threads belonging to task *tk*. The expression *next_thread(th)* denotes the successor *THREAD* of *th*, if one exists, in the sequence of threads belonging to the owning task of thread *th*.

*ThreadList**TasksAndThreads**threads_head* : *TASK* \leftrightarrow *THREAD**next_thread* : *THREAD* \leftrightarrow *THREAD**threads* = *Gen_set*(*threads_head*, *next_thread*)

The expression *processors_head(pset)* denotes the first *PROCESSOR*, if one exists, in the sequence of processors belonging to processor set *pset*. The expression *next_processor(proc)* denotes the successor *PROCESSOR* of *proc*, if one exists, in the sequence of processors belonging to the processor set of which *proc* is a member.

*ProcessorList**ProcessorAndProcessorSet**processors_head* : *PROCESSOR_SET* \leftrightarrow *PROCESSOR**next_processor* : *PROCESSOR* \leftrightarrow *PROCESSOR**processors* = *Gen_set*(*processors_head*, *next_processor*)

The expression *assigned_tasks_head(pset)* denotes the first *TASK*, if one exists, in the sequence of tasks belonging to processor set *pset*. The expression *next_assigned_task(tk)* denotes the

successor *TASK* of *tk*, if one exists, in the sequence of tasks belonging to the processor set to which *tk* is assigned.

<i>AssignedTaskList</i> <i>TaskAndProcessorSet</i> $\underline{assigned_tasks_head} : PROCESSOR_SET \leftrightarrow TASK$ $\underline{next_assigned_task} : TASK \leftrightarrow TASK$ $\underline{have_assigned_tasks} = Gen_set(\underline{assigned_tasks_head}, \underline{next_assigned_task})$

The expression $\underline{assigned_threads_head}(pset)$ denotes the first *THREAD*, if one exists, in the sequence of threads belonging to processor set *pset*. The expression $\underline{next_assigned_thread}(th)$ denotes the successor *THREAD* of *th*, if one exists, in the sequence of threads belonging to the processor set to which *th* is assigned.

<i>AssignedThreadList</i> <i>ThreadAndProcessorSet</i> $\underline{assigned_threads_head} : PROCESSOR_SET \leftrightarrow THREAD$ $\underline{next_assigned_thread} : THREAD \leftrightarrow THREAD$ $\underline{have_assigned_threads} = Gen_set(\underline{assigned_threads_head}, \underline{next_assigned_thread})$

The expression $\underline{messages_head}(port)$ denotes the first *MESSAGE*, if one exists, in the sequence of messages waiting in *port*. The expression $\underline{next_message}(msg)$ denotes the successor *MESSAGE* of *msg*, if one exists, in the sequence of messages waiting in the port in which *msg* is waiting.

Editorial Note:

This says nothing about the messages in a queue of a port set. This queue is not currently modeled, so there is nothing to refine. If we add port set message queues, the refinement would appear nearly identical to the following refinement.

<i>MessageInPortList</i> <i>MessageQueues</i> $\underline{messages_head} : PORT \leftrightarrow MESSAGE$ $\underline{next_message} : MESSAGE \leftrightarrow MESSAGE$ $\underline{message_in_port_rel} = Gen_seq(\underline{messages_head}, \underline{next_message})$
--

Index

The italic numbers denote the pages where the corresponding entry is described, numbers underlined point to the definition, all others indicate the places where it is used.

Symbols	
<i>Abort_thread</i>	<u>67</u>
<i>Abort_thread_depress</i>	<u>67</u>
<i>Access_machine_attribute</i>	<u>66</u>
<i>a_active_thread</i>	<u>55</u>
<i>Add_name</i>	<u>65</u>
<i>AddressSpace</i>	<u>39</u>
<i>Add_thread</i>	<u>68</u>
<i>Add_thread_secure</i>	<u>68</u>
<i>Add_value</i>	<u>318</u>
<i>AID</i>	<u>59</u>
<i>Allocate_vm_region</i>	<u>66</u>
<i>allocated</i>	<u>38</u>
<i>Alter_pns_info</i>	<u>65</u>
<i>Anti_symmetric</i>	<u>317</u>
<i>Assign_processor</i>	<u>70</u>
<i>Assign_processor_to_set</i>	<u>69</u>
<i>Assign_task</i>	<u>70</u>
<i>Assign_task_to_pset</i>	<u>68</u>
<i>Assign_thread</i>	<u>70</u>
<i>Assign_thread_to_pset</i>	<u>67</u>
<i>Audit_ids</i>	<u>46</u>
<i>a_audit_server_port</i>	<u>76</u>
<i>authentication_server_port</i>	<u>76</u>
<i>backing_chain</i>	<u>40</u>
<i>backing_memory</i>	<u>40</u>
<i>backing_offset</i>	<u>40</u>
<i>backing_rel</i>	<u>40</u>
<i>BASE_MSG_ELEMENT</i>	<u>48</u>
<i>Base_user_priority</i>	<u>13</u>
<i>cache_allows</i>	<u>73</u>
<i>Cached_ruling_allows</i>	<u>73</u>
<i>Cached_ruling_allows</i>	<u>73</u>
<i>cached_ruling_avail</i>	<u>73</u>
<i>Can_receive</i>	<u>65</u>
<i>Can_send</i>	<u>65</u>
<i>Can_switc</i>	<u>67</u>
<i>Can_switc_pri</i>	<u>67</u>
<i>Capability</i>	<u>20</u>
<i>Change_page_locks</i>	<u>66</u>
<i>Chg_pset_max_pri</i>	<u>70</u>
<i>Chg_vm_region_prot</i>	<u>66</u>
<i>Change_sid</i>	<u>68</u>
<i>Chg_task_priority</i>	<u>68</u>
<i>Close_device</i>	<u>71</u>
<i>Co_carries_memory</i>	<u>43</u>
<i>Co_carries_rights</i>	<u>43</u>
<i>COMPLEX_OPTION_BOOLEAN</i>	<u>45</u>
<i>COMPLEX_OPTION</i>	<u>43</u>
<i>containing_port</i>	<u>25</u>
<i>containing_set</i>	<u>22</u>
<i>control_memory</i>	<u>29</u>
<i>Control_pager</i>	<u>71</u>
<i>controlled_proc_set</i>	<u>31</u>
<i>c_copy_strategy</i>	<u>35</u>
<i>Copy_vm</i>	<u>66</u>
<i>cpu_time</i>	<u>18</u>
<i>Create_pset</i>	<u>68</u>
<i>Create_task</i>	<u>68</u>
<i>Create_task_secure</i>	<u>68</u>
<i>Cross_context_create</i>	<u>68</u>
<i>Cross_context_inherit</i>	<u>68</u>
<i>crypto_server_port</i>	<u>76</u>
<i>dead_namep</i>	<u>23</u>
<i>dead_right_ref_count</i>	<u>23</u>
<i>dead_right_rel</i>	<u>23</u>
<i>DeadRights</i>	<u>24</u>
<i>Deallocate_vm_region</i>	<u>66</u>
<i>default_mem_manager</i>	<u>36</u>
<i>Default_port_sid</i>	<u>61</u>
<i>Default_vm_port_sid</i>	<u>61</u>
<i>Define_new_scheduling_policy</i>	<u>70</u>
<i>depressed_threads</i>	<u>13</u>
<i>Depress_pri</i>	<u>67</u>
<i>priority_before_depression</i>	<u>13</u>
<i>Derive_kernel_as</i>	<u>62</u>
<i>Destroy_object</i>	<u>66</u>
<i>Destroy_pset</i>	<u>70</u>
<i>DeviceData</i>	<u>57</u>
<i>DeviceExist</i>	<u>9</u>
<i>device_exists</i>	<u>9</u>
<i>device_filter_info</i>	<u>57</u>
<i>DeviceFilterInfo</i>	<u>57</u>
<i>DEVICE_FILTER_INFO</i>	<u>57</u>
<i>DEVICE_FILTER</i>	<u>57</u>
<i>device_in</i>	<u>56</u>
<i>device_open_count</i>	<u>55</u>
<i>DeviceOpenCount</i>	<u>55</u>
<i>device_out</i>	<u>56</u>
<i>Device_permissions</i>	<u>71</u>
<i>device_port</i>	<u>31</u>
<i>device_port_rel</i>	<u>31</u>

<i>DEVICE_RECORD</i>	<u>57</u>	<i>Get_thread_kernel_port</i>	<u>67</u>
<i>DEVICE</i>	<u>9</u>	<i>Get_thread_state</i>	<u>67</u>
<i>Devices</i>	<u>57</u>	<i>Get_time</i>	<u>68</u>
<i>DevicesAndPorts</i>	<u>32</u>	<i>Get_vm_region_info</i>	<u>66</u>
<i>device_status</i>	<u>57</u>	<i>Get_vm_statistics</i>	<u>66</u>
<i>DeviceStatus</i>	<u>57</u>	<i>Halted</i>	<u>11</u>
<i>DEVICE_STATUS</i>	<u>57</u>	<i>have_assigned_tasks</i>	<u>54</u>
<i>dirty_rel</i>	<u>37</u>	<i>have_assigned_threads</i>	<u>54</u>
<i>Dtos</i>	<u>77</u>	<i>Have_execute</i>	<u>66</u>
<i>DtosAdditions</i>	<u>77</u>	<i>Have_read</i>	<u>66</u>
<i>DtosMessages</i>	<u>74</u>	<i>Have_write</i>	<u>66</u>
<i>emulation_vector</i>	<u>15</u>	<i>Higher_priority</i>	<u>12</u>
<i>EmulationVector</i>	<u>15</u>	<i>Highest_possible_priority</i>	<u>12</u>
<i>enabled_sp</i>	<u>54</u>	<i>Hold_receive</i>	<u>65</u>
<i>Environment_slot</i>	<u>34</u>	<i>Hold_send</i>	<u>65</u>
<i>event_count</i>	<u>56</u>	<i>Hold_send_once</i>	<u>65</u>
<i>EVENT_COUNTER</i>	<u>56</u>	<i>host_control_port</i>	<u>30</u>
<i>Events</i>	<u>56</u>	<i>host_name_port</i>	<u>30</u>
<i>Exception_ids</i>	<u>46</u>	<i>Host_control_port_permissions</i>	<u>69</u>
<i>Exist</i>	<u>10</u>	<i>HOST</i>	<u>9</u>
<i>Extract_right</i>	<u>65</u>	<i>HostsAndPorts</i>	<u>30</u>
<i>FILTER_PRIORITY</i>	<u>57</u>	<i>HostsAndProcessors</i>	<u>53</u>
<i>Fixedpri</i>	<u>14</u>	<i>host_time</i>	<u>55</u>
<i>Flush_permission</i>	<u>68</u>	<i>HostTime</i>	<u>55</u>
<i>forcibly_queued</i>	<u>21</u>	<i>Host_name_port_permissions</i>	<u>68</u>
<i>Get_attributes</i>	<u>66</u>	<i>idle_threads</i>	<u>12</u>
<i>Get_audit_port</i>	<u>68</u>	<i>inheritance</i>	<u>40</u>
<i>Get_authentication_port</i>	<u>68</u>	<i>Inheritance</i>	<u>40</u>
<i>Get_boot_info</i>	<u>69</u>	<i>Inheritance_option_copy</i>	<u>40</u>
<i>Get_crypto_port</i>	<u>68</u>	<i>Inheritance_option_none</i>	<u>40</u>
<i>Get_default_pset_name</i>	<u>68</u>	<i>INHERITANCE_OPTION</i>	<u>40</u>
<i>Get_device_status</i>	<u>71</u>	<i>Inheritance_option_share</i>	<u>40</u>
<i>Get_emulation</i>	<u>68</u>	<i>initialized</i>	<u>35</u>
<i>Get_host_control_port</i>	<u>68</u>	<i>Initiate_secure</i>	<u>67</u>
<i>Get_host_info</i>	<u>68</u>	<i>In_line</i>	<u>47</u>
<i>Get_host_name</i>	<u>68</u>	<i>instruction_pointer</i>	<u>15</u>
<i>Get_host_processors</i>	<u>69</u>	<i>INTERNAL_BODY</i>	<u>49</u>
<i>Get_host_version</i>	<u>68</u>	<i>Internal_element</i>	<u>49</u>
<i>Get_negotiation_port</i>	<u>68</u>	<i>InternalMessage</i>	<u>51</u>
<i>Get_network_ss_port</i>	<u>68</u>	<i>Interpose</i>	<u>65</u>
<i>Get_processor_assignment</i>	<u>69</u>	<i>Invalidate_scheduling_policy</i>	<u>70</u>
<i>Get_processor_info</i>	<u>69</u>	<i>Invoke_lock_request</i>	<u>66</u>
<i>Get_pset_info</i>	<u>70</u>	<i>Ipc_permissions</i>	<u>65</u>
<i>Get_security_master_port</i>	<u>68</u>	<i>Ip_dead</i>	<u>9</u>
<i>Get_security_client_port</i>	<u>68</u>	<i>Ip_null</i>	<u>9</u>
<i>Get_special_port</i>	<u>68</u>	<i>kernel</i>	<u>10</u>
<i>Get_task_assignment</i>	<u>68</u>	<i>Kernel</i>	<u>10</u>
<i>Get_task_boot_port</i>	<u>68</u>	<i>kernel_as</i>	<u>62</u>
<i>Get_task_exception_port</i>	<u>68</u>	<i>KernelAs</i>	<u>63</u>
<i>Get_task_info</i>	<u>68</u>	<i>KernelCache</i>	<u>73</u>
<i>Get_task_kernel_port</i>	<u>68</u>	<i>Kernel_permission</i>	<u>64</u>
<i>Get_task_threads</i>	<u>68</u>	<i>KernelPortSid</i>	<u>61</u>
<i>Get_thread_assignment</i>	<u>67</u>	<i>Kernel_reply_permissions</i>	<u>71</u>
<i>Get_thread_exception_port</i>	<u>67</u>	<i>kernel_reply_ports</i>	<u>76</u>
<i>Get_thread_info</i>	<u>67</u>	<i>KernelReplyPorts</i>	<u>76</u>

<i>Kernel_service_reply_ids</i>	<u>46</u>	<i>Memory_copy_none</i>	<u>35</u>
<i>local_namep</i>	<u>24</u>	<i>MEMORY_COPY_STRATEGY</i>	<u>35</u>
<i>Lock</i>	<u>38</u>	<i>Memory_copy_temporary</i>	<u>35</u>
<i>Lookup_ports</i>	<u>65</u>	<i>name_port</i>	<u>29</u>
<i>Lower_priority</i>	<u>12</u>	<i>name_port_rel</i>	<u>29</u>
<i>Lowest_possible_priority</i>	<u>12</u>	<i>MEMORY</i>	<u>9</u>
<i>Mach</i>	<u>58</u>	<i>MemoriesAndPorts</i>	<u>30</u>
<i>Mach_exception_id</i>	<u>46</u>	<i>Memory</i>	<u>36</u>
<i>MachInternalHeader</i>	<u>45</u>	<i>MemoryExist</i>	<u>9</u>
<i>MachMsgHeader</i>	<u>45</u>	<i>memory_exists</i>	<u>9</u>
<i>MACH_MSG_OPTION</i>	<u>42</u>	<i>Mem_obj_confirmation_ids</i>	<u>46</u>
<i>MACH_MSG_TYPE</i>	<u>43</u>	<i>Memory_object_permissions</i>	<u>66</u>
<i>Mach_notify_ids</i>	<u>46</u>	<i>MemorySystem</i>	<u>42</u>
<i>Mach_port_dead</i>	<u>19</u>	<i>Message</i>	<u>50</u>
<i>Mach_port_null</i>	<u>19</u>	<i>Msg_element</i>	<u>48</u>
<i>Mach_port_q_limit_default</i>	<u>25</u>	<i>MessageExist</i>	<u>9</u>
<i>Mach_port_q_limit_max</i>	<u>25</u>	<i>message_exists</i>	<u>9</u>
<i>mach_protection</i>	<u>39</u>	<i>message_in_port_rel</i>	<u>25</u>
<i>MachProtection</i>	<u>39</u>	<i>MessageQueues</i>	<u>26</u>
<i>Mach_rcv_large</i>	<u>42</u>	<i>MESSAGE</i>	<u>9</u>
<i>Mach_rcv_msg</i>	<u>42</u>	<i>Messages</i>	<u>53</u>
<i>Mach_rcv_notify</i>	<u>42</u>	<i>MID</i>	<u>59</u>
<i>Mach_rcv_timeout</i>	<u>42</u>	<i>Lowest_priority</i>	<u>13</u>
<i>Mach_send_cancel</i>	<u>42</u>	<i>Mmt_copy_send</i>	<u>43</u>
<i>Mach_send_msg</i>	<u>42</u>	<i>Mmt_make_send</i>	<u>43</u>
<i>Mach_send_notify</i>	<u>42</u>	<i>Mmt_make_send_once</i>	<u>43</u>
<i>Mach_send_timeout</i>	<u>42</u>	<i>Mmt_move_receive</i>	<u>43</u>
<i>Make_page_precious</i>	<u>66</u>	<i>Mmt_move_send</i>	<u>43</u>
<i>make_send_count</i>	<u>24</u>	<i>Mmt_move_send_once</i>	<u>43</u>
<i>Make_sid</i>	<u>68</u>	<i>Mach_msg_type_port_receive</i>	<u>44</u>
<i>managed</i>	<u>35</u>	<i>Mach_msg_type_port_rights</i>	<u>44</u>
<i>manager</i>	<u>35</u>	<i>Mach_msg_type_port_send</i>	<u>44</u>
<i>Manipulate_port_set</i>	<u>65</u>	<i>Mach_msg_type_port_send_once</i>	<u>44</u>
<i>map_rel</i>	<u>38</u>	<i>MESSAGE_BODY</i>	<u>48</u>
<i>Map_device</i>	<u>71</u>	<i>msg_contents</i>	<u>52</u>
<i>mapped</i>	<u>39</u>	<i>MSG_DATA</i>	<u>48</u>
<i>mapped_devices</i>	<u>56</u>	<i>Msg_deallocate</i>	<u>47</u>
<i>MappedDevices</i>	<u>56</u>	<i>Msg_dont_deallocate</i>	<u>47</u>
<i>mapped_memory</i>	<u>39</u>	<i>Msg_error_invalid_memory</i>	<u>50</u>
<i>mapped_offset</i>	<u>39</u>	<i>Msg_error_invalid_right</i>	<u>50</u>
<i>Map_vm_region</i>	<u>65</u>	<i>Msg_error_invalid_type</i>	<u>50</u>
<i>master_device_port</i>	<u>32</u>	<i>Msg_error_msg_too_small</i>	<u>50</u>
<i>MasterDevicePort</i>	<u>32</u>	<i>Msg_error_notify_in_progress</i>	<u>50</u>
<i>master_proc</i>	<u>53</u>	<i>MSG_ERROR</i>	<u>50</u>
<i>Highest_priority</i>	<u>13</u>	<i>Msg_error_timed_out</i>	<u>50</u>
<i>max_protection</i>	<u>39</u>	<i>msg_operation</i>	<u>52</u>
<i>Max_right_refs</i>	<u>20</u>	<i>Operations</i>	<u>52</u>
<i>Max_samples</i>	<u>16</u>	<i>msg_receiving_sid</i>	<u>74</u>
<i>may_cache</i>	<u>35</u>	<i>Msg_region</i>	<u>49</u>
<i>May_control_processor</i>	<u>69</u>	<i>msg_ruling</i>	<u>74</u>
<i>member_rel</i>	<u>53</u>	<i>msg_sending_sid</i>	<u>73</u>
<i>control_port</i>	<u>29</u>	<i>msg_specified_sid</i>	<u>74</u>
<i>control_port_rel</i>	<u>29</u>	<i>msg_specified_vector</i>	<u>74</u>
<i>Memory_copy_call</i>	<u>35</u>	<i>Msg_stat_pseudo</i>	<u>50</u>
<i>Memory_copy_delay</i>	<u>35</u>	<i>Msg_stat_rcv</i>	<u>50</u>

<i>Msg_stat_send</i>	<u>50</u>	<i>Pc_processor</i>	<u>33</u>
<i>MSG_STATUS</i>	<u>50</u>	<i>Pc_ps_control</i>	<u>33</u>
<i>Msg_value</i>	<u>48</u>	<i>Pc_ps_name</i>	<u>33</u>
<i>MSG_VALUE</i>	<u>49</u>	<i>Pc_task</i>	<u>33</u>
<i>named_port</i>	<u>20</u>	<i>Pc_thread</i>	<u>33</u>
<i>named_proc_set</i>	<u>31</u>	<i>PendingReceive</i>	<u>51</u>
<i>NAME</i>	<u>19</u>	<i>pending_receives</i>	<u>52</u>
<i>Name_server_slot</i>	<u>34</u>	<i>PERMISSION</i>	<u>64</u>
<i>negotiation_server_port</i>	<u>76</u>	<i>port_aid</i>	<u>60</u>
<i>Network_packet_ids</i>	<u>46</u>	<i>port_class</i>	<u>33</u>
<i>network_ss_port</i>	<u>76</u>	<i>PORT_CLASS</i>	<u>33</u>
<i>Notifications</i>	<u>27</u>	<i>PortClasses</i>	<u>33</u>
<i>number_of_rights</i>	<u>24</u>	<i>port_device</i>	<u>31</u>
<i>object_memory</i>	<u>29</u>	<i>PortExist</i>	<u>10</u>
<i>object_port</i>	<u>29</u>	<i>port_exists</i>	<u>9</u>
<i>object_port_rel</i>	<u>29</u>	<i>port_mid</i>	<u>60</u>
<i>ObjectSid</i>	<u>62</u>	<i>PortNameSpace</i>	<u>24</u>
<i>Observe_pns_info</i>	<u>65</u>	<i>port_notify_dead</i>	<u>27</u>
<i>Observe_pset_processes</i>	<u>70</u>	<i>port_notify_dead_rel</i>	<u>26</u>
<i>OFFSET</i>	<u>35</u>	<i>port_notify_destroyed</i>	<u>26</u>
<i>OLSD</i>	<u>48</u>	<i>port_notify_destroyed_rel</i>	<u>26</u>
<i>Open_device</i>	<u>71</u>	<i>port_notify_no_more_senders</i>	<u>26</u>
<i>OPERATION</i>	<u>45</u>	<i>port_notify_no_more_senders_rel</i>	<u>26</u>
<i>OSI</i>	<u>60</u>	<i>Port_permissions</i>	<u>65</u>
<i>Osi_to_aid</i>	<u>60</u>	<i>port_pointer</i>	<u>9</u>
<i>Osi_to_mid</i>	<u>60</u>	<i>Port_rename</i>	<u>65</u>
<i>Out_of_line</i>	<u>47</u>	<i>port_right_rel</i>	<u>19</u>
<i>WORD</i>	<u>35</u>	<i>port_right_namep</i>	<u>21</u>
<i>threads</i>	<u>10</u>	<i>port_right_seq</i>	<u>33</u>
<i>owning_task</i>	<u>10</u>	<i>PORT</i>	<u>9</u>
<i>PageAndMemory</i>	<u>38</u>	<i>port_set</i>	<u>22</u>
<i>page_aid</i>	<u>61</u>	<i>port_set_namep</i>	<u>22</u>
<i>PageExist</i>	<u>9</u>	<i>port_set_rel</i>	<u>22</u>
<i>page_exists</i>	<u>9</u>	<i>PortSets</i>	<u>23</u>
<i>memory_fault</i>	<u>36</u>	<i>port_sid</i>	<u>60</u>
<i>PAGE_INDEX</i>	<u>39</u>	<i>PortSid</i>	<u>61</u>
<i>page_lock_rel</i>	<u>38</u>	<i>port_size</i>	<u>25</u>
<i>page_locks</i>	<u>38</u>	<i>PortSummary</i>	<u>26</u>
<i>page_mid</i>	<u>61</u>	<i>Poset</i>	<u>317</u>
<i>PAGE_OFFSET</i>	<u>37</u>	<i>Pp_to_page_sid</i>	<u>61</u>
<i>Pager_permissions</i>	<u>66</u>	<i>precious</i>	<u>37</u>
<i>Pager_request_ids</i>	<u>46</u>	<i>Priority_levels</i>	<u>12</u>
<i>PAGE</i>	<u>9</u>	<i>proc_assigned_procset</i>	<u>53</u>
<i>page_sid</i>	<u>61</u>	<i>Process</i>	<u>58</u>
<i>PageSid</i>	<u>62</u>	<i>ProcessorExist</i>	<u>9</u>
<i>Page_vm_region</i>	<u>66</u>	<i>proc_exists</i>	<u>9</u>
<i>page_word_rel</i>	<u>36</u>	<i>Processor_permissions</i>	<u>69</u>
<i>page_word_fun</i>	<u>37</u>	<i>processor_port_rel</i>	<u>30</u>
<i>SCHED_POLICY_DATA</i>	<u>14</u>	<i>PROCESSOR</i>	<u>9</u>
<i>parent_task</i>	<u>75</u>	<i>Pset_ctrl_port</i>	<u>69</u>
<i>ParentTask</i>	<u>75</u>	<i>Pset_names</i>	<u>68</u>
<i>Pc_device</i>	<u>33</u>	<i>processors</i>	<u>53</u>
<i>Pc_host_control</i>	<u>33</u>	<i>ProcessorsAndPorts</i>	<u>31</u>
<i>Pc_host_name</i>	<u>33</u>	<i>ProcessorAndProcessorSet</i>	<u>54</u>
<i>Pc_memory</i>	<u>33</u>		

<i>proc_self</i>	<u>30</u>	<i>run_state</i>	<u>11</u>
<i>ProcessorSetExist</i>	<u>9</u>	<i>RUN_STATES</i>	<u>11</u>
<i>procset_exists</i>	<u>9</u>	<i>s_ampled_tasks</i>	<u>17</u>
<i>procset_name_port</i>	<u>31</u>	<i>s_ampled_threads</i>	<u>16</u>
<i>Procset_control_port_permissions</i>	<u>70</u>	<i>Sample_periodic</i>	<u>16</u>
<i>PROCESSOR_SET</i>	<u>9</u>	<i>SAMPLE</i>	<u>16</u>
<i>procset_self</i>	<u>31</u>	<i>Sample_task</i>	<u>68</u>
<i>Procset_name_port_permissions</i>	<u>70</u>	<i>Sample_thread</i>	<u>67</u>
<i>Protection</i>	<u>76</u>	<i>SAMPLE_TYPES</i>	<u>16</u>
<i>Execute</i>	<u>38</u>	<i>Sample_vm_cow_faults</i>	<u>16</u>
<i>Read</i>	<u>38</u>	<i>SAMPLE_VM_FAULTS</i>	<u>16</u>
<i>PROTECTION</i>	<u>38</u>	<i>Sample_vm_faults_any</i>	<u>16</u>
<i>Write</i>	<u>38</u>	<i>Sample_vm_pagein_faults</i>	<u>16</u>
<i>Provide_data</i>	<u>66</u>	<i>Sample_vm_reactivation_faults</i>	<u>16</u>
<i>Provide_permission</i>	<u>71</u>	<i>Sample_vm_zfill_faults</i>	<u>16</u>
<i>ps_control_port_rel</i>	<u>31</u>	<i>Save_page</i>	<u>66</u>
<i>ps_max_priority</i>	<u>54</u>	<i>SCHED_POLICY</i>	<u>14</u>
<i>ps_name_port_rel</i>	<u>31</u>	<i>s_ecurity_server_client_port</i>	<u>76</u>
<i>q_limit</i>	<u>25</u>	<i>Security_server_ids</i>	<u>46</u>
<i>Raise_exception</i>	<u>67</u>	<i>s_ecurity_server_master_port</i>	<u>76</u>
<i>Read_device</i>	<u>71</u>	<i>self_task</i>	<u>28</u>
<i>Read_vm_region</i>	<u>66</u>	<i>self_thread</i>	<u>28</u>
<i>Reboot_host</i>	<u>69</u>	<i>Send</i>	<u>19</u>
<i>Receive</i>	<u>19</u>	<i>sender</i>	<u>20</u>
<i>receiver</i>	<u>20</u>	<i>Send_once</i>	<u>19</u>
<i>receiver_name</i>	<u>20</u>	<i>SendRightsCount</i>	<u>25</u>
<i>Recognized_sample_types</i>	<u>16</u>	<i>s_quence_no</i>	<u>25</u>
<i>Recognized_transfer_options</i>	<u>43</u>	<i>Seq_plus</i>	<u>318</u>
<i>Reflexive</i>	<u>317</u>	<i>ServerPorts</i>	<u>76</u>
<i>r_egistered_rights</i>	<u>34</u>	<i>Service_check_deferred</i>	<u>87</u>
<i>RegisteredRights</i>	<u>35</u>	<i>Service_slot</i>	<u>34</u>
<i>Register_notification</i>	<u>65</u>	<i>Set_attributes</i>	<u>66</u>
<i>Register_ports</i>	<u>65</u>	<i>Set_audit_port</i>	<u>68</u>
<i>Remove_name</i>	<u>65</u>	<i>Set_authentication_port</i>	<u>68</u>
<i>Remove_page</i>	<u>66</u>	<i>Set_crypto_port</i>	<u>68</u>
<i>reply_port</i>	<u>52</u>	<i>Set_ibac_port</i>	<u>66</u>
<i>r_eplay_port_rel</i>	<u>52</u>	<i>Set_default_memory_mgr</i>	<u>69</u>
<i>reply_port_right</i>	<u>52</u>	<i>Set_device_filter</i>	<u>71</u>
<i>ReplyPorts</i>	<u>52</u>	<i>Set_device_status</i>	<u>71</u>
<i>represented</i>	<u>37</u>	<i>Set_emulation</i>	<u>68</u>
<i>represented_memory</i>	<u>37</u>	<i>Set_vm_region_inherit</i>	<u>66</u>
<i>represented_offset</i>	<u>37</u>	<i>Set_max_thread_priority</i>	<u>67</u>
<i>representing_page</i>	<u>37</u>	<i>Set_negotiation_port</i>	<u>68</u>
<i>r_represents_rel</i>	<u>37</u>	<i>Set_network_ss_port</i>	<u>68</u>
<i>represents_memory</i>	<u>37</u>	<i>Set_ras</i>	<u>68</u>
<i>Required_permission</i>	<u>87</u>	<i>Set_reply</i>	<u>65</u>
<i>Resume_task</i>	<u>68</u>	<i>Set_security_master_port</i>	<u>68</u>
<i>Resume_thread</i>	<u>67</u>	<i>Set_security_client_port</i>	<u>68</u>
<i>Revoke_ibac</i>	<u>66</u>	<i>Set_special_port</i>	<u>68</u>
<i>RIGHT</i>	<u>19</u>	<i>Set_task_boot_port</i>	<u>68</u>
<i>r_right</i>	<u>21</u>	<i>Set_task_exception_port</i>	<u>68</u>
<i>Ruling</i>	<u>72</u>	<i>Set_task_kernel_port</i>	<u>68</u>
<i>Ruling_allows</i>	<u>72</u>	<i>Set_thread_exception_port</i>	<u>67</u>
<i>Ruling_allows</i>	<u>72</u>	<i>Set_thread_kernel_port</i>	<u>67</u>
<i>Running</i>	<u>11</u>	<i>Set_thread_policy</i>	<u>67</u>

<i>Set_thread_priority</i>	<u>67</u>	<i>TasksAndRights</i>	<u>22</u>
<i>Set_thread_state</i>	<u>67</u>	<i>TasksAndThreads</i>	<u>11</u>
<i>Set_time</i>	<u>69</u>	<i>task_self</i>	<u>28</u>
<i>shadow_memories</i>	<u>40</u>	<i>task_self_rel</i>	<u>27</u>
<i>ShadowMemories</i>	<u>41</u>	<i>Task_self_sid</i>	<u>63</u>
<i>sleep_time</i>	<u>18</u>	<i>task_sid</i>	<u>60</u>
<i>so_right</i>	<u>21</u>	<i>task_sself</i>	<u>28</u>
<i>SpecialPurposePorts</i>	<u>32</u>	<i>task_sself_rel</i>	<u>27</u>
<i>SpecialTaskPorts</i>	<u>28</u>	<i>task_suspend_count</i>	<u>12</u>
<i>SpecialThreadPorts</i>	<u>29</u>	<i>TaskSuspendCount</i>	<u>12</u>
<i>Specify</i>	<u>65</u>	<i>task_target</i>	<u>63</u>
<i>s_right</i>	<u>21</u>	<i>Task_task_permissions</i>	<u>68</u>
<i>s_right_ref_count</i>	<u>20</u>	<i>task_thread_rel</i>	<u>10</u>
<i>s_r_right</i>	<u>21</u>	<i>Tcs_task_empty</i>	<u>74</u>
<i>SSI</i>	<u>59</u>	<i>Tcs_task_ready</i>	<u>75</u>
<i>Ssi_to_aid</i>	<u>59</u>	<i>Tcs_thread_created</i>	<u>74</u>
<i>Ssi_to_mid</i>	<u>59</u>	<i>Tcs_thread_state_set</i>	<u>75</u>
<i>State_info_avail</i>	<u>18</u>	<i>temporary_rel</i>	<u>35</u>
<i>Stopped</i>	<u>11</u>	<i>Terminate_task</i>	<u>68</u>
<i>SubjectSid</i>	<u>60</u>	<i>Terminate_thread</i>	<u>67</u>
<i>Supply_ibac</i>	<u>66</u>	<i>the_processor</i>	<u>30</u>
<i>supplying_device</i>	<u>56</u>	<i>ThreadAndProcessorSet</i>	<u>55</u>
<i>SUPP_MACHINE_ARCH</i>	<u>18</u>	<i>thread_assigned_to</i>	<u>54</u>
<i>supported_sp</i>	<u>14</u>	<i>thread_assignment_rel</i>	<u>54</u>
<i>Suspend_task</i>	<u>68</u>	<i>thread_eport</i>	<u>28</u>
<i>Suspend_thread</i>	<u>67</u>	<i>thread_eport_rel</i>	<u>28</u>
<i>swapped_threads</i>	<u>11</u>	<i>ThreadExecStatus</i>	<u>12</u>
<i>Switch_thread</i>	<u>67</u>	<i>ThreadExist</i>	<u>9</u>
<i>system_time</i>	<u>17</u>	<i>thread_exists</i>	<u>9</u>
<i>TargetSids</i>	<u>63</u>	<i>ThreadInstruction</i>	<u>15</u>
<i>task_aid</i>	<u>60</u>	<i>ThreadMachineState</i>	<u>18</u>
<i>TaskAndProcessorSet</i>	<u>54</u>	<i>thread_max_priority</i>	<u>13</u>
<i>task_assigned_to</i>	<u>54</u>	<i>Thread_permissions</i>	<u>67</u>
<i>task_assignment_rel</i>	<u>54</u>	<i>Thread_port_sid</i>	<u>61</u>
<i>task_bport</i>	<u>28</u>	<i>ThreadPri</i>	<u>14</u>
<i>task_bport_rel</i>	<u>28</u>	<i>thread_priority</i>	<u>13</u>
<i>task_creation_state</i>	<u>75</u>	<i>THREAD</i>	<u>9</u>
<i>TaskCreationState</i>	<u>75</u>	<i>Threads</i>	<u>19</u>
<i>TASK_CREATION_STATE</i>	<u>74</u>	<i>thread_samples</i>	<u>17</u>
<i>task_eport</i>	<u>28</u>	<i>thread_sample_sequence_number</i>	<u>16</u>
<i>task_eport_rel</i>	<u>28</u>	<i>thread_sample_types</i>	<u>16</u>
<i>TaskExist</i>	<u>9</u>	<i>ThreadSampling</i>	<u>17</u>
<i>task_exists</i>	<u>9</u>	<i>ThreadsAndProcessors</i>	<u>55</u>
<i>task_mid</i>	<u>60</u>	<i>thread_sched_policy</i>	<u>14</u>
<i>Task_port_register_max</i>	<u>34</u>	<i>ThreadSchedPolicy</i>	<u>15</u>
<i>Task_port_sid</i>	<u>61</u>	<i>thread_sched_policy_data</i>	<u>14</u>
<i>task_priority</i>	<u>14</u>	<i>thread_sched_priority</i>	<u>13</u>
<i>TaskPriority</i>	<u>14</u>	<i>thread_self</i>	<u>28</u>
<i>task_received_msgs</i>	<u>52</u>	<i>thread_self_rel</i>	<u>28</u>
<i>TASK</i>	<u>9</u>	<i>Thread_self_sid</i>	<u>63</u>
<i>task_samples</i>	<u>17</u>	<i>thread_sid</i>	<u>60</u>
<i>task_sample_sequence_number</i>	<u>17</u>	<i>thread_sself</i>	<u>28</u>
<i>task_sample_types</i>	<u>17</u>	<i>thread_sself_rel</i>	<u>28</u>
<i>TaskSampling</i>	<u>17</u>	<i>thread_state</i>	<u>18</u>
<i>TasksAndPorts</i>	<u>20</u>	<i>THREAD_STATE_INFO</i>	<u>18</u>

<i>THREAD_STATE_INFO_TYPES</i>	<u>18</u>	<i>Cached_ruling_allows</i>	<u>73</u>
<i>ThreadStatistics</i>	<u>18</u>	<i>cached_ruling_avail</i>	<u>73</u>
<i>thread_suspend_count</i>	<u>11</u>	<i>crypto_server_port</i>	<u>76</u>
<i>threads_wired</i>	<u>12</u>	<i>Default_port_sid</i>	<u>61</u>
<i>thread_target</i>	<u>63</u>	<i>Default_vm_port_sid</i>	<u>61</u>
<i>thread_waiting</i>	<u>56</u>	<i>kernel_as</i>	<u>62</u>
<i>Timeshare</i>	<u>14</u>	<i>kernel_reply_ports</i>	<u>76</u>
<i>total_naked_srights</i>	<u>33</u>	<i>msg_receiving_sid</i>	<u>74</u>
<i>total_name_space_srights</i>	<u>33</u>	<i>msg_ruling</i>	<u>74</u>
<i>total_srights</i>	<u>33</u>	<i>msg_sending_sid</i>	<u>73</u>
<i>TotalSendRights</i>	<u>34</u>	<i>msg_specified_sid</i>	<u>74</u>
<i>Transfer_ool</i>	<u>65</u>	<i>msg_specified_vector</i>	<u>74</u>
<i>Transfer_receive</i>	<u>65</u>	<i>negotiation_server_port</i>	<u>76</u>
<i>Transfer_rights</i>	<u>65</u>	<i>network_ss_port</i>	<u>76</u>
<i>Transfer_send</i>	<u>65</u>	<i>Osi_to_aid</i>	<u>60</u>
<i>Transfer_send_once</i>	<u>65</u>	<i>Osi_to_mid</i>	<u>60</u>
<i>Transition_sid</i>	<u>68</u>	<i>page_aid</i>	<u>61</u>
<i>Transitive</i>	<u>317</u>	<i>page_mid</i>	<u>61</u>
<i>Transit_memory</i>	<u>49</u>	<i>page_sid</i>	<u>61</u>
<i>Transit_right</i>	<u>49</u>	<i>parent_task</i>	<u>75</u>
<i>Uninterruptible</i>	<u>11</u>	<i>port_aid</i>	<u>60</u>
<i>Usable_cached_ruling</i>	<u>72</u>	<i>port_mid</i>	<u>60</u>
<i>Usable_ruling</i>	<u>72</u>	<i>port_sid</i>	<u>60</u>
<i>UserReferenceCount</i>	<u>21</u>	<i>Pp_to_page_sid</i>	<u>61</u>
<i>user_time</i>	<u>17</u>	<i>Ruling_allows</i>	<u>72</u>
<i>Values_disjoint</i>	<u>316</u>	<i>security_server_client_port</i>	<u>76</u>
<i>Values_partition</i>	<u>316</u>	<i>security_server_master_port</i>	<u>76</u>
<i>V_data</i>	<u>48</u>	<i>Ssi_to_aid</i>	<u>59</u>
<i>V_data_in</i>	<u>49</u>	<i>Ssi_to_mid</i>	<u>59</u>
<i>V_DATA_LOCATION</i>	<u>49</u>	<i>task_aid</i>	<u>60</u>
<i>V_data_out</i>	<u>49</u>	<i>task_creation_state</i>	<u>75</u>
<i>VIRTUAL_ADDRESS</i>	<u>15</u>	<i>task_mid</i>	<u>60</u>
<i>Vm_end</i>	<u>15</u>	<i>Task_port_sid</i>	<u>61</u>
<i>Vm_permissions</i>	<u>66</u>	<i>Task_self_sid</i>	<u>63</u>
<i>Vm_start</i>	<u>15</u>	<i>task_sid</i>	<u>60</u>
<i>V_port</i>	<u>49</u>	<i>task_target</i>	<u>63</u>
<i>Wait_evt</i>	<u>67</u>	<i>Thread_port_sid</i>	<u>61</u>
<i>Waiting</i>	<u>11</u>	<i>Thread_self_sid</i>	<u>63</u>
<i>wire_count</i>	<u>41</u>	<i>thread_sid</i>	<u>60</u>
<i>wired</i>	<u>41</u>	<i>thread_target</i>	<u>63</u>
<i>Wired</i>	<u>41</u>	<i>Usable_cached_ruling</i>	<u>72</u>
<i>Wire_thread</i>	<u>69</u>	<i>Usable_ruling</i>	<u>72</u>
<i>Wire_thread_into_memory</i>	<u>67</u>		
<i>Wire_vm</i>	<u>69</u>	DTOS Types:	
<i>Wire_vm_for_task</i>	<u>66</u>	<i>AID</i>	<u>59</u>
<i>Wrap_value</i>	<u>318</u>	<i>MID</i>	<u>59</u>
<i>Write_device</i>	<u>71</u>	<i>OSI</i>	<u>60</u>
<i>Write_vm_region</i>	<u>66</u>	<i>PERMISSION</i>	<u>64</u>
<i>default</i>	<u>53</u>	<i>SSI</i>	<u>59</u>
<i>protection</i>	<u>76</u>	<i>TASK_CREATION_STATE</i>	<u>74</u>

D**DTOS Structures:**

<i>audit_server_port</i>	<u>76</u>
<i>authentication_server_port</i>	<u>76</u>
<i>cache_allows</i>	<u>73</u>

G**Global Identifiers:**

<i>Abort_thread</i>	<u>67</u>
<i>Abort_thread_depress</i>	<u>67</u>
<i>Access_machine_attribute</i>	<u>66</u>

<i>Add_name</i>	<u>65</u>	<i>Get_host_control_port</i>	<u>68</u>
<i>Add_thread</i>	<u>68</u>	<i>Get_host_info</i>	<u>68</u>
<i>Add_thread_secure</i>	<u>68</u>	<i>Get_host_name</i>	<u>68</u>
<i>Add_value</i>	<u>318</u>	<i>Get_host_processors</i>	<u>69</u>
<i>Allocate_vm_region</i>	<u>66</u>	<i>Get_host_version</i>	<u>68</u>
<i>Alter_pns_info</i>	<u>65</u>	<i>Get_negotiation_port</i>	<u>68</u>
<i>Anti_symmetric</i>	<u>317</u>	<i>Get_network_ss_port</i>	<u>68</u>
<i>Assign_processor</i>	<u>70</u>	<i>Get_processor_assignment</i>	<u>69</u>
<i>Assign_processor_to_set</i>	<u>69</u>	<i>Get_processor_info</i>	<u>69</u>
<i>Assign_task</i>	<u>70</u>	<i>Get_pset_info</i>	<u>70</u>
<i>Assign_task_to_pset</i>	<u>68</u>	<i>Get_security_master_port</i>	<u>68</u>
<i>Assign_thread</i>	<u>70</u>	<i>Get_security_client_port</i>	<u>68</u>
<i>Assign_thread_to_pset</i>	<u>67</u>	<i>Get_special_port</i>	<u>68</u>
<i>Audit_ids</i>	<u>46</u>	<i>Get_task_assignment</i>	<u>68</u>
<i>Base_user_priority</i>	<u>13</u>	<i>Get_task_boot_port</i>	<u>68</u>
<i>Cached_ruling_allows</i>	<u>73</u>	<i>Get_task_exception_port</i>	<u>68</u>
<i>Can_receive</i>	<u>65</u>	<i>Get_task_info</i>	<u>68</u>
<i>Can_send</i>	<u>65</u>	<i>Get_task_kernel_port</i>	<u>68</u>
<i>Can_swch</i>	<u>67</u>	<i>Get_task_threads</i>	<u>68</u>
<i>Can_swch_pri</i>	<u>67</u>	<i>Get_thread_assignment</i>	<u>67</u>
<i>Change_page_locks</i>	<u>66</u>	<i>Get_thread_exception_port</i>	<u>67</u>
<i>Chg_pset_max_pri</i>	<u>70</u>	<i>Get_thread_info</i>	<u>67</u>
<i>Chg_vm_region_prot</i>	<u>66</u>	<i>Get_thread_kernel_port</i>	<u>67</u>
<i>Change_sid</i>	<u>68</u>	<i>Get_thread_state</i>	<u>67</u>
<i>Chg_task_priority</i>	<u>68</u>	<i>Get_time</i>	<u>68</u>
<i>Close_device</i>	<u>71</u>	<i>Get_vm_region_info</i>	<u>66</u>
<i>Co_carries_memory</i>	<u>43</u>	<i>Get_vm_statistics</i>	<u>66</u>
<i>Co_carries_rights</i>	<u>43</u>	<i>Halted</i>	<u>11</u>
<i>Control_pager</i>	<u>71</u>	<i>Have_execute</i>	<u>66</u>
<i>Copy_vm</i>	<u>66</u>	<i>Have_read</i>	<u>66</u>
<i>Create_pset</i>	<u>68</u>	<i>Have_write</i>	<u>66</u>
<i>Create_task</i>	<u>68</u>	<i>Higher_priority</i>	<u>12</u>
<i>Create_task_secure</i>	<u>68</u>	<i>Highest_possible_priority</i>	<u>12</u>
<i>Cross_context_create</i>	<u>68</u>	<i>Hold_receive</i>	<u>65</u>
<i>Cross_context_inherit</i>	<u>68</u>	<i>Hold_send</i>	<u>65</u>
<i>Deallocate_vm_region</i>	<u>66</u>	<i>Hold_send_once</i>	<u>65</u>
<i>Define_new_scheduling_policy</i>	<u>70</u>	<i>Host_control_port_permissions</i>	<u>69</u>
<i>Depress_pri</i>	<u>67</u>	<i>Host_name_port_permissions</i>	<u>68</u>
<i>Derive_kernel_as</i>	<u>62</u>	<i>Inheritance_option_copy</i>	<u>40</u>
<i>Destroy_object</i>	<u>66</u>	<i>Inheritance_option_none</i>	<u>40</u>
<i>Destroy_pset</i>	<u>70</u>	<i>Inheritance_option_share</i>	<u>40</u>
<i>Device_permissions</i>	<u>71</u>	<i>Initiate_secure</i>	<u>67</u>
<i>Environment_slot</i>	<u>34</u>	<i>In_line</i>	<u>47</u>
<i>Exception_ids</i>	<u>46</u>	<i>Interpose</i>	<u>65</u>
<i>Extract_right</i>	<u>65</u>	<i>Invalidate_scheduling_policy</i>	<u>70</u>
<i>Fixedpri</i>	<u>14</u>	<i>Invoke_lock_request</i>	<u>66</u>
<i>Flush_permission</i>	<u>68</u>	<i>Ipc_permissions</i>	<u>65</u>
<i>Get_attributes</i>	<u>66</u>	<i>Ip_dead</i>	<u>9</u>
<i>Get_audit_port</i>	<u>68</u>	<i>Ip_null</i>	<u>9</u>
<i>Get_authentication_port</i>	<u>68</u>	<i>Kernel_permission</i>	<u>64</u>
<i>Get_boot_info</i>	<u>69</u>	<i>Kernel_reply_permissions</i>	<u>71</u>
<i>Get_crypto_port</i>	<u>68</u>	<i>Kernel_service_reply_ids</i>	<u>46</u>
<i>Get_default_pset_name</i>	<u>68</u>	<i>Lookup_ports</i>	<u>65</u>
<i>Get_device_status</i>	<u>71</u>	<i>Lower_priority</i>	<u>12</u>
<i>Get_emulation</i>	<u>68</u>	<i>Lowest_possible_priority</i>	<u>12</u>

<i>Mach_exception_id</i>	<u>46</u>	<i>Network_packet_ids</i>	<u>46</u>
<i>Mach_notify_ids</i>	<u>46</u>	<i>Observe_pns_info</i>	<u>65</u>
<i>Mach_port_dead</i>	<u>19</u>	<i>Observe_pset_processes</i>	<u>70</u>
<i>Mach_port_null</i>	<u>19</u>	<i>Open_device</i>	<u>71</u>
<i>Mach_port_q_limit_default</i>	<u>25</u>	<i>Out_of_line</i>	<u>47</u>
<i>Mach_port_q_limit_max</i>	<u>25</u>	<i>Pager_permissions</i>	<u>66</u>
<i>Mach_rcv_large</i>	<u>42</u>	<i>Pager_request_ids</i>	<u>46</u>
<i>Mach_rcv_msg</i>	<u>42</u>	<i>Page_vm_region</i>	<u>66</u>
<i>Mach_rcv_notify</i>	<u>42</u>	<i>Pc_device</i>	<u>33</u>
<i>Mach_rcv_timeout</i>	<u>42</u>	<i>Pc_host_control</i>	<u>33</u>
<i>Mach_send_cancel</i>	<u>42</u>	<i>Pc_host_name</i>	<u>33</u>
<i>Mach_send_msg</i>	<u>42</u>	<i>Pc_memory</i>	<u>33</u>
<i>Mach_send_notify</i>	<u>42</u>	<i>Pc_processor</i>	<u>33</u>
<i>Mach_send_timeout</i>	<u>42</u>	<i>Pc_ps_control</i>	<u>33</u>
<i>Make_page_precious</i>	<u>66</u>	<i>Pc_ps_name</i>	<u>33</u>
<i>Make_sid</i>	<u>68</u>	<i>Pc_task</i>	<u>33</u>
<i>Manipulate_port_set</i>	<u>65</u>	<i>Pc_thread</i>	<u>33</u>
<i>Map_device</i>	<u>71</u>	<i>Port_permissions</i>	<u>65</u>
<i>Map_vm_region</i>	<u>65</u>	<i>Port_rename</i>	<u>65</u>
<i>Highest_priority</i>	<u>13</u>	<i>Poset</i>	<u>317</u>
<i>Max_right_refs</i>	<u>20</u>	<i>Priority_levels</i>	<u>12</u>
<i>Max_samples</i>	<u>16</u>	<i>Processor_permissions</i>	<u>69</u>
<i>May_control_processor</i>	<u>69</u>	<i>Pset_ctrl_port</i>	<u>69</u>
<i>Memory_copy_call</i>	<u>35</u>	<i>Pset_names</i>	<u>68</u>
<i>Memory_copy_delay</i>	<u>35</u>	<i>Procset_control_port_permissions</i>	<u>70</u>
<i>Memory_copy_none</i>	<u>35</u>	<i>Procset_name_port_permissions</i>	<u>70</u>
<i>Memory_copy_temporary</i>	<u>35</u>	<i>Execute</i>	<u>38</u>
<i>Mem_obj_confirmation_ids</i>	<u>46</u>	<i>Read</i>	<u>38</u>
<i>Memory_object_permissions</i>	<u>66</u>	<i>Write</i>	<u>38</u>
<i>Msg_element</i>	<u>48</u>	<i>Provide_data</i>	<u>66</u>
<i>Lowest_priority</i>	<u>13</u>	<i>Provide_permission</i>	<u>71</u>
<i>Mmt_copy_send</i>	<u>43</u>	<i>Raise_exception</i>	<u>67</u>
<i>Mmt_make_send</i>	<u>43</u>	<i>Read_device</i>	<u>71</u>
<i>Mmt_make_send_once</i>	<u>43</u>	<i>Read_vm_region</i>	<u>66</u>
<i>Mmt_move_receive</i>	<u>43</u>	<i>Reboot_host</i>	<u>69</u>
<i>Mmt_move_send</i>	<u>43</u>	<i>Receive</i>	<u>19</u>
<i>Mmt_move_send_once</i>	<u>43</u>	<i>Recognized_sample_types</i>	<u>16</u>
<i>Mach_msg_type_port_receive</i>	<u>44</u>	<i>Recognized_transfer_options</i>	<u>43</u>
<i>Mach_msg_type_port_rights</i>	<u>44</u>	<i>Reflexive</i>	<u>317</u>
<i>Mach_msg_type_port_send</i>	<u>44</u>	<i>Register_notification</i>	<u>65</u>
<i>Mach_msg_type_port_send_once</i>	<u>44</u>	<i>Register_ports</i>	<u>65</u>
<i>Msg_deallocate</i>	<u>47</u>	<i>Remove_name</i>	<u>65</u>
<i>Msg_dont_deallocate</i>	<u>47</u>	<i>Remove_page</i>	<u>66</u>
<i>Msg_error_invalid_memory</i>	<u>50</u>	<i>Required_permission</i>	<u>87</u>
<i>Msg_error_invalid_right</i>	<u>50</u>	<i>Resume_task</i>	<u>68</u>
<i>Msg_error_invalid_type</i>	<u>50</u>	<i>Resume_thread</i>	<u>67</u>
<i>Msg_error_msg_too_small</i>	<u>50</u>	<i>Revoke_ibac</i>	<u>66</u>
<i>Msg_error_notify_in_progress</i>	<u>50</u>	<i>Ruling_allows</i>	<u>72</u>
<i>Msg_error_timed_out</i>	<u>50</u>	<i>Running</i>	<u>11</u>
<i>Msg_region</i>	<u>49</u>	<i>Sample_periodic</i>	<u>16</u>
<i>Msg_stat_pseudo</i>	<u>50</u>	<i>Sample_task</i>	<u>68</u>
<i>Msg_stat_rcv</i>	<u>50</u>	<i>Sample_thread</i>	<u>67</u>
<i>Msg_stat_send</i>	<u>50</u>	<i>Sample_vm_cow_faults</i>	<u>16</u>
<i>Msg_value</i>	<u>48</u>	<i>SAMPLE_VM_FAULTS</i>	<u>16</u>
<i>Name_server_slot</i>	<u>34</u>	<i>Sample_vm_faults_any</i>	<u>16</u>

<i>Sample_vm_pagein_faults</i>	<u>16</u>	<i>Transfer_receive</i>	<u>65</u>
<i>Sample_vm_reactivation_faults</i>	<u>16</u>	<i>Transfer_rights</i>	<u>65</u>
<i>Sample_vm_zfill_faults</i>	<u>16</u>	<i>Transfer_send</i>	<u>65</u>
<i>Save_page</i>	<u>66</u>	<i>Transfer_send_once</i>	<u>65</u>
<i>Security_server_ids</i>	<u>46</u>	<i>Transition_sid</i>	<u>68</u>
<i>Send</i>	<u>19</u>	<i>Transitive</i>	<u>317</u>
<i>Send_once</i>	<u>19</u>	<i>Transit_memory</i>	<u>49</u>
<i>Seq_plus</i>	<u>318</u>	<i>Transit_right</i>	<u>49</u>
<i>Service_check_deferred</i>	<u>87</u>	<i>Uninterruptible</i>	<u>11</u>
<i>Service_slot</i>	<u>34</u>	<i>Values_disjoint</i>	<u>316</u>
<i>Set_attributes</i>	<u>66</u>	<i>Values_partition</i>	<u>316</u>
<i>Set_audit_port</i>	<u>68</u>	<i>V_data</i>	<u>48</u>
<i>Set_authentication_port</i>	<u>68</u>	<i>V_data_in</i>	<u>49</u>
<i>Set_crypto_port</i>	<u>68</u>	<i>V_data_out</i>	<u>49</u>
<i>Set_ibac_port</i>	<u>66</u>	<i>Vm_end</i>	<u>15</u>
<i>Set_default_memory_mgr</i>	<u>69</u>	<i>Vm_permissions</i>	<u>66</u>
<i>Set_device_filter</i>	<u>71</u>	<i>Vm_start</i>	<u>15</u>
<i>Set_device_status</i>	<u>71</u>	<i>V_port</i>	<u>49</u>
<i>Set_emulation</i>	<u>68</u>	<i>Wait_evt</i>	<u>67</u>
<i>Set_vm_region_inherit</i>	<u>66</u>	<i>Waiting</i>	<u>11</u>
<i>Set_max_thread_priority</i>	<u>67</u>	<i>Wire_thread</i>	<u>69</u>
<i>Set_negotiation_port</i>	<u>68</u>	<i>Wire_thread_into_memory</i>	<u>67</u>
<i>Set_network_ss_port</i>	<u>68</u>	<i>Wire_vm</i>	<u>69</u>
<i>Set_ras</i>	<u>68</u>	<i>Wire_vm_for_task</i>	<u>66</u>
<i>Set_reply</i>	<u>65</u>	<i>Wrap_value</i>	<u>318</u>
<i>Set_security_master_port</i>	<u>68</u>	<i>Write_device</i>	<u>71</u>
<i>Set_security_client_port</i>	<u>68</u>	<i>Write_vm_region</i>	<u>66</u>
<i>Set_special_port</i>	<u>68</u>		
<i>Set_task_boot_port</i>	<u>68</u>		
<i>Set_task_exception_port</i>	<u>68</u>		
<i>Set_task_kernel_port</i>	<u>68</u>		
<i>Set_thread_exception_port</i>	<u>67</u>		
<i>Set_thread_kernel_port</i>	<u>67</u>		
<i>Set_thread_policy</i>	<u>67</u>		
<i>Set_thread_priority</i>	<u>67</u>		
<i>Set_thread_state</i>	<u>67</u>		
<i>Set_time</i>	<u>69</u>		
<i>Specify</i>	<u>65</u>		
<i>State_info_avail</i>	<u>18</u>		
<i>Stopped</i>	<u>11</u>		
<i>Supply_ibac</i>	<u>66</u>		
<i>Suspend_task</i>	<u>68</u>		
<i>Suspend_thread</i>	<u>67</u>		
<i>Switch_thread</i>	<u>67</u>		
<i>Task_port_register_max</i>	<u>34</u>		
<i>Task_task_permissions</i>	<u>68</u>		
<i>Tcs_task_empty</i>	<u>74</u>		
<i>Tcs_task_ready</i>	<u>75</u>		
<i>Tcs_thread_created</i>	<u>74</u>		
<i>Tcs_thread_state_set</i>	<u>75</u>		
<i>Terminate_task</i>	<u>68</u>		
<i>Terminate_thread</i>	<u>67</u>		
<i>Thread_permissions</i>	<u>67</u>		
<i>Timeshare</i>	<u>14</u>		
<i>Transfer_ool</i>	<u>65</u>		

M

Mach Structures:

<i>active_thread</i>	<u>55</u>
<i>allocated</i>	<u>38</u>
<i>backing_chain</i>	<u>40</u>
<i>backing_memory</i>	<u>40</u>
<i>backing_offset</i>	<u>40</u>
<i>backing_rel</i>	<u>40</u>
<i>containing_port</i>	<u>25</u>
<i>containing_set</i>	<u>22</u>
<i>control_memory</i>	<u>29</u>
<i>controlled_proc_set</i>	<u>31</u>
<i>copy_strategy</i>	<u>35</u>
<i>cpu_time</i>	<u>18</u>
<i>dead_namep</i>	<u>23</u>
<i>dead_right_ref_count</i>	<u>23</u>
<i>dead_right_rel</i>	<u>23</u>
<i>default_mem_manager</i>	<u>36</u>
<i>depressed_threads</i>	<u>13</u>
<i>priority_before_depression</i>	<u>13</u>
<i>device_exists</i>	<u>9</u>
<i>device_filter_info</i>	<u>57</u>
<i>device_in</i>	<u>56</u>
<i>device_open_count</i>	<u>55</u>
<i>device_out</i>	<u>56</u>
<i>device_port</i>	<u>31</u>
<i>device_port_rel</i>	<u>31</u>
<i>device_status</i>	<u>57</u>

<i>dirty_rel</i>	<u>37</u>	<i>port_device</i>	<u>31</u>
<i>emulation_vector</i>	<u>15</u>	<i>port_exists</i>	<u>9</u>
<i>enabled_sp</i>	<u>54</u>	<i>port_notify_dead</i>	<u>27</u>
<i>event_count</i>	<u>56</u>	<i>port_notify_dead_rel</i>	<u>26</u>
<i>forcibly_queued</i>	<u>21</u>	<i>port_notify_destroyed</i>	<u>26</u>
<i>have_assigned_tasks</i>	<u>54</u>	<i>port_notify_destroyed_rel</i>	<u>26</u>
<i>have_assigned_threads</i>	<u>54</u>	<i>port_notify_no_more_senders</i>	<u>26</u>
<i>host_control_port</i>	<u>30</u>	<i>port_notify_no_more_senders_rel</i>	<u>26</u>
<i>host_name_port</i>	<u>30</u>	<i>port_pointer</i>	<u>9</u>
<i>host_time</i>	<u>55</u>	<i>port_right_rel</i>	<u>19</u>
<i>idle_threads</i>	<u>12</u>	<i>port_right_namep</i>	<u>21</u>
<i>inheritance</i>	<u>40</u>	<i>port_right_seq</i>	<u>33</u>
<i>initialized</i>	<u>35</u>	<i>port_set</i>	<u>22</u>
<i>instruction_pointer</i>	<u>15</u>	<i>port_set_namep</i>	<u>22</u>
<i>kernel</i>	<u>10</u>	<i>port_set_rel</i>	<u>22</u>
<i>local_namep</i>	<u>24</u>	<i>port_size</i>	<u>25</u>
<i>mach_protection</i>	<u>39</u>	<i>precious</i>	<u>37</u>
<i>make_send_count</i>	<u>24</u>	<i>proc_assigned_procset</i>	<u>53</u>
<i>managed</i>	<u>35</u>	<i>proc_exists</i>	<u>9</u>
<i>manager</i>	<u>35</u>	<i>processor_port_rel</i>	<u>30</u>
<i>map_rel</i>	<u>38</u>	<i>processors</i>	<u>53</u>
<i>mapped</i>	<u>39</u>	<i>proc_self</i>	<u>30</u>
<i>mapped_devices</i>	<u>56</u>	<i>procset_exists</i>	<u>9</u>
<i>mapped_memory</i>	<u>39</u>	<i>procset_name_port</i>	<u>31</u>
<i>mapped_offset</i>	<u>39</u>	<i>procset_self</i>	<u>31</u>
<i>master_device_port</i>	<u>32</u>	<i>ps_control_port_rel</i>	<u>31</u>
<i>master_proc</i>	<u>53</u>	<i>ps_max_priority</i>	<u>54</u>
<i>max_protection</i>	<u>39</u>	<i>ps_name_port_rel</i>	<u>31</u>
<i>may_cache</i>	<u>35</u>	<i>q_limit</i>	<u>25</u>
<i>member_rel</i>	<u>53</u>	<i>receiver</i>	<u>20</u>
<i>control_port</i>	<u>29</u>	<i>receiver_name</i>	<u>20</u>
<i>control_port_rel</i>	<u>29</u>	<i>registered_rights</i>	<u>34</u>
<i>name_port</i>	<u>29</u>	<i>reply_port</i>	<u>52</u>
<i>name_port_rel</i>	<u>29</u>	<i>reply_port_rel</i>	<u>52</u>
<i>memory_exists</i>	<u>9</u>	<i>reply_port_right</i>	<u>52</u>
<i>message_exists</i>	<u>9</u>	<i>represented</i>	<u>37</u>
<i>message_in_port_rel</i>	<u>25</u>	<i>represented_memory</i>	<u>37</u>
<i>msg_contents</i>	<u>52</u>	<i>represented_offset</i>	<u>37</u>
<i>msg_operation</i>	<u>52</u>	<i>representing_page</i>	<u>37</u>
<i>named_port</i>	<u>20</u>	<i>represents_rel</i>	<u>37</u>
<i>named_proc_set</i>	<u>31</u>	<i>represents_memory</i>	<u>37</u>
<i>number_of_rights</i>	<u>24</u>	<i>r_right</i>	<u>21</u>
<i>object_memory</i>	<u>29</u>	<i>run_state</i>	<u>11</u>
<i>object_port</i>	<u>29</u>	<i>sampled_tasks</i>	<u>17</u>
<i>object_port_rel</i>	<u>29</u>	<i>sampled_threads</i>	<u>16</u>
<i>threads</i>	<u>10</u>	<i>self_task</i>	<u>28</u>
<i>owning_task</i>	<u>10</u>	<i>self_thread</i>	<u>28</u>
<i>page_exists</i>	<u>9</u>	<i>sender</i>	<u>20</u>
<i>memory_fault</i>	<u>36</u>	<i>sequence_no</i>	<u>25</u>
<i>page_lock_rel</i>	<u>38</u>	<i>shadow_memories</i>	<u>40</u>
<i>page_locks</i>	<u>38</u>	<i>sleep_time</i>	<u>18</u>
<i>page_word_rel</i>	<u>36</u>	<i>so_right</i>	<u>21</u>
<i>page_word_fun</i>	<u>37</u>	<i>s_right</i>	<u>21</u>
<i>pending_receives</i>	<u>52</u>	<i>s_right_ref_count</i>	<u>20</u>
<i>port_class</i>	<u>33</u>		

<i>s_r_right</i>	<u>21</u>	Mach Types:	
<i>supplying_device</i>	<u>56</u>	<i>BASE_MSG_ELEMENT</i>	<u>48</u>
<i>supported_sp</i>	<u>14</u>	<i>COMPLEX_OPTION_BOOLEAN</i>	<u>45</u>
<i>swapped_threads</i>	<u>11</u>	<i>COMPLEX_OPTION</i>	<u>43</u>
<i>system_time</i>	<u>17</u>	<i>DEVICE_FILTER_INFO</i>	<u>57</u>
<i>task_assigned_to</i>	<u>54</u>	<i>DEVICE_FILTER</i>	<u>57</u>
<i>task_assignment_rel</i>	<u>54</u>	<i>DEVICE_RECORD</i>	<u>57</u>
<i>task_bport</i>	<u>28</u>	<i>DEVICE</i>	<u>9</u>
<i>task_bport_rel</i>	<u>28</u>	<i>DEVICE_STATUS</i>	<u>57</u>
<i>task_eport</i>	<u>28</u>	<i>EVENT_COUNTER</i>	<u>56</u>
<i>task_eport_rel</i>	<u>28</u>	<i>FILTER_PRIORITY</i>	<u>57</u>
<i>task_exists</i>	<u>9</u>	<i>HOST</i>	<u>9</u>
<i>task_priority</i>	<u>14</u>	<i>INHERITANCE_OPTION</i>	<u>40</u>
<i>task_received_msgs</i>	<u>52</u>	<i>INTERNAL_BODY</i>	<u>49</u>
<i>task_samples</i>	<u>17</u>	<i>Internal_element</i>	<u>49</u>
<i>task_sample_sequence_number</i>	<u>17</u>	<i>MACH_MSG_OPTION</i>	<u>42</u>
<i>task_sample_types</i>	<u>17</u>	<i>MACH_MSG_TYPE</i>	<u>43</u>
<i>task_self</i>	<u>28</u>	<i>MEMORY_COPY_STRATEGY</i>	<u>35</u>
<i>task_self_rel</i>	<u>27</u>	<i>MEMORY</i>	<u>9</u>
<i>task_sself</i>	<u>28</u>	<i>MESSAGE</i>	<u>9</u>
<i>task_sself_rel</i>	<u>27</u>	<i>MESSAGE_BODY</i>	<u>48</u>
<i>task_suspend_count</i>	<u>12</u>	<i>MSG_DATA</i>	<u>48</u>
<i>task_thread_rel</i>	<u>10</u>	<i>MSG_ERROR</i>	<u>50</u>
<i>temporary_rel</i>	<u>35</u>	<i>MSG_STATUS</i>	<u>50</u>
<i>the_processor</i>	<u>30</u>	<i>MSG_VALUE</i>	<u>49</u>
<i>thread_assigned_to</i>	<u>54</u>	<i>NAME</i>	<u>19</u>
<i>thread_assignment_rel</i>	<u>54</u>	<i>OFFSET</i>	<u>35</u>
<i>thread_eport</i>	<u>28</u>	<i>OLSD</i>	<u>48</u>
<i>thread_eport_rel</i>	<u>28</u>	<i>OPERATION</i>	<u>45</u>
<i>thread_exists</i>	<u>9</u>	<i>WORD</i>	<u>35</u>
<i>thread_max_priority</i>	<u>13</u>	<i>PAGE_INDEX</i>	<u>39</u>
<i>thread_priority</i>	<u>13</u>	<i>PAGE_OFFSET</i>	<u>37</u>
<i>thread_samples</i>	<u>17</u>	<i>PAGE</i>	<u>9</u>
<i>thread_sample_sequence_number</i>	<u>16</u>	<i>SCHED_POLICY_DATA</i>	<u>14</u>
<i>thread_sample_types</i>	<u>16</u>	<i>PORT_CLASS</i>	<u>33</u>
<i>thread_sched_policy</i>	<u>14</u>	<i>PORT</i>	<u>9</u>
<i>thread_sched_policy_data</i>	<u>14</u>	<i>PROCESSOR</i>	<u>9</u>
<i>thread_sched_priority</i>	<u>13</u>	<i>PROCESSOR_SET</i>	<u>9</u>
<i>thread_self</i>	<u>28</u>	<i>PROTECTION</i>	<u>38</u>
<i>thread_self_rel</i>	<u>28</u>	<i>RIGHT</i>	<u>19</u>
<i>thread_sself</i>	<u>28</u>	<i>RUN_STATES</i>	<u>11</u>
<i>thread_sself_rel</i>	<u>28</u>	<i>SAMPLE</i>	<u>16</u>
<i>thread_state</i>	<u>18</u>	<i>SAMPLE_TYPES</i>	<u>16</u>
<i>thread_suspend_count</i>	<u>11</u>	<i>SCHED_POLICY</i>	<u>14</u>
<i>threads_wired</i>	<u>12</u>	<i>SUPP_MACHINE_ARCH</i>	<u>18</u>
<i>thread_waiting</i>	<u>56</u>	<i>TASK</i>	<u>9</u>
<i>total_naked_rights</i>	<u>33</u>	<i>THREAD</i>	<u>9</u>
<i>total_name_space_rights</i>	<u>33</u>	<i>THREAD_STATE_INFO</i>	<u>18</u>
<i>total_rights</i>	<u>33</u>	<i>THREAD_STATE_INFO_TYPES</i>	<u>18</u>
<i>user_time</i>	<u>17</u>	<i>V_DATA_LOCATION</i>	<u>49</u>
<i>wire_count</i>	<u>41</u>	<i>VIRTUAL_ADDRESS</i>	<u>15</u>
<i>wired</i>	<u>41</u>		
<i>default</i>	<u>53</u>	S	
<i>protection</i>	<u>76</u>	Schemas:	
		<i>AddressSpace</i>	<u>39</u>
		<i>Capability</i>	<u>20</u>

<i>DeadRights</i>	<u>24</u>	<i>PortClasses</i>	<u>33</u>
<i>DeviceData</i>	<u>57</u>	<i>PortExist</i>	<u>10</u>
<i>DeviceExist</i>	<u>9</u>	<i>PortNameSpace</i>	<u>24</u>
<i>DeviceFilterInfo</i>	<u>57</u>	<i>PortSets</i>	<u>23</u>
<i>DeviceOpenCount</i>	<u>55</u>	<i>PortSid</i>	<u>61</u>
<i>Devices</i>	<u>57</u>	<i>PortSummary</i>	<u>26</u>
<i>DevicesAndPorts</i>	<u>32</u>	<i>Process</i>	<u>58</u>
<i>DeviceStatus</i>	<u>57</u>	<i>ProcessorExist</i>	<u>9</u>
<i>Dtos</i>	<u>77</u>	<i>ProcessorsAndPorts</i>	<u>31</u>
<i>DtosAdditions</i>	<u>77</u>	<i>ProcessorAndProcessorSet</i>	<u>54</u>
<i>DtosMessages</i>	<u>74</u>	<i>ProcessorSetExist</i>	<u>9</u>
<i>EmulationVector</i>	<u>15</u>	<i>Protection</i>	<u>76</u>
<i>Events</i>	<u>56</u>	<i>RegisteredRights</i>	<u>35</u>
<i>Exist</i>	<u>10</u>	<i>ReplyPorts</i>	<u>52</u>
<i>HostsAndPorts</i>	<u>30</u>	<i>Ruling</i>	<u>72</u>
<i>HostsAndProcessors</i>	<u>53</u>	<i>SendRightsCount</i>	<u>25</u>
<i>HostTime</i>	<u>55</u>	<i>ServerPorts</i>	<u>76</u>
<i>Inheritance</i>	<u>40</u>	<i>ShadowMemories</i>	<u>41</u>
<i>InternalMessage</i>	<u>51</u>	<i>SpecialPurposePorts</i>	<u>32</u>
<i>Kernel</i>	<u>10</u>	<i>SpecialTaskPorts</i>	<u>28</u>
<i>KernelAs</i>	<u>63</u>	<i>SpecialThreadPorts</i>	<u>29</u>
<i>KernelCache</i>	<u>73</u>	<i>SubjectSid</i>	<u>60</u>
<i>KernelPortSid</i>	<u>61</u>	<i>TargetSids</i>	<u>63</u>
<i>KernelReplyPorts</i>	<u>76</u>	<i>TaskAndProcessorSet</i>	<u>54</u>
<i>Lock</i>	<u>38</u>	<i>TaskCreationState</i>	<u>75</u>
<i>Mach</i>	<u>58</u>	<i>TaskExist</i>	<u>9</u>
<i>MachInternalHeader</i>	<u>45</u>	<i>TaskPriority</i>	<u>14</u>
<i>MachMsgHeader</i>	<u>45</u>	<i>TaskSampling</i>	<u>17</u>
<i>MachProtection</i>	<u>39</u>	<i>TasksAndPorts</i>	<u>20</u>
<i>MappedDevices</i>	<u>56</u>	<i>TasksAndRights</i>	<u>22</u>
<i>MasterDevicePort</i>	<u>32</u>	<i>TasksAndThreads</i>	<u>11</u>
<i>MemoriesAndPorts</i>	<u>30</u>	<i>TaskSuspendCount</i>	<u>12</u>
<i>Memory</i>	<u>36</u>	<i>ThreadAndProcessorSet</i>	<u>55</u>
<i>MemoryExist</i>	<u>9</u>	<i>ThreadExecStatus</i>	<u>12</u>
<i>MemorySystem</i>	<u>42</u>	<i>ThreadExist</i>	<u>9</u>
<i>Message</i>	<u>50</u>	<i>ThreadInstruction</i>	<u>15</u>
<i>MessageExist</i>	<u>9</u>	<i>ThreadMachineState</i>	<u>18</u>
<i>MessageQueues</i>	<u>26</u>	<i>ThreadPri</i>	<u>14</u>
<i>Messages</i>	<u>53</u>	<i>Threads</i>	<u>19</u>
<i>Operations</i>	<u>52</u>	<i>ThreadSampling</i>	<u>17</u>
<i>Notifications</i>	<u>27</u>	<i>ThreadsAndProcessors</i>	<u>55</u>
<i>ObjectSid</i>	<u>62</u>	<i>ThreadSchedPolicy</i>	<u>15</u>
<i>PageAndMemory</i>	<u>38</u>	<i>ThreadStatistics</i>	<u>18</u>
<i>PageExist</i>	<u>9</u>	<i>TotalSendRights</i>	<u>34</u>
<i>PageSid</i>	<u>62</u>	<i>UserReferenceCount</i>	<u>21</u>
<i>ParentTask</i>	<u>75</u>	<i>Wired</i>	<u>41</u>
<i>PendingReceive</i>	<u>51</u>		