# The Alpha 21264 Out-of-Order Implementation



Reorder Buffer (ROB)

| Branch prediction and instr fetch | | Register File P1-P64 |

Committed Reg Map
R1→P1
R2→P2

Instr 1
Instr 2
Instr 3
Instr 4
Instr 5
Instr 6
Instr 7

R1 ← R1+R2
R2 ← R1+R3
BEQZ R2
R3 ← R1+R2
R1 ← R3+R2
LD  R4 ← 8[R3]
ST R4 → 8[R1]

Instr Fetch Queue

Decode & Rename

Speculative Reg Map
R1→P36
R2→P34

P33 ← P1+P2
P34 ← P33+P3
BEQZ P34
P35 ← P33+P34
P36 ← P35+P34
P37 ← 8[P35]
P37 → 8[P36]

Issue Queue (IQ)

ALU   ALU   ALU

Results written to regfile and tags broadcast to IQ

ALU

P37 ← [P35 + 8]
P37 → [P36 + 8]

LSQ

D-Cache

1 core

CMP Multi-core
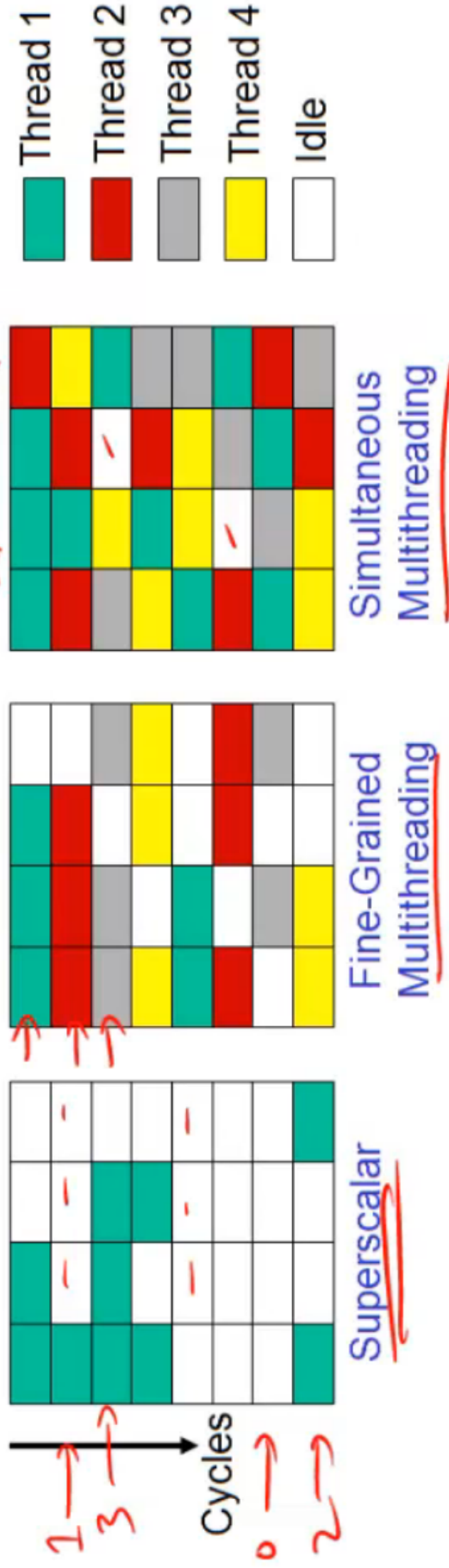
caches

9

# Thread-Level Parallelism

$iW > 4$

Avg IPC = 1.5 $\begin{cases} \\ \end{cases}$ 3.0

IPC = 1.5

- Motivation:
  - ➤ a single thread leaves a processor under-utilized for most of the time
  - ➤ by doubling processor area, single thread performance barely improves

- Strategies for thread-level parallelism:
  - ➤ multiple threads share the same large processor → reduces under-utilization, efficient resource allocation Simultaneous Multi-Threading (SMT)
  - ➤ each thread executes on its own mini processor → simple design, low interference between threads Chip Multi-Processing (CMP) or multi-core
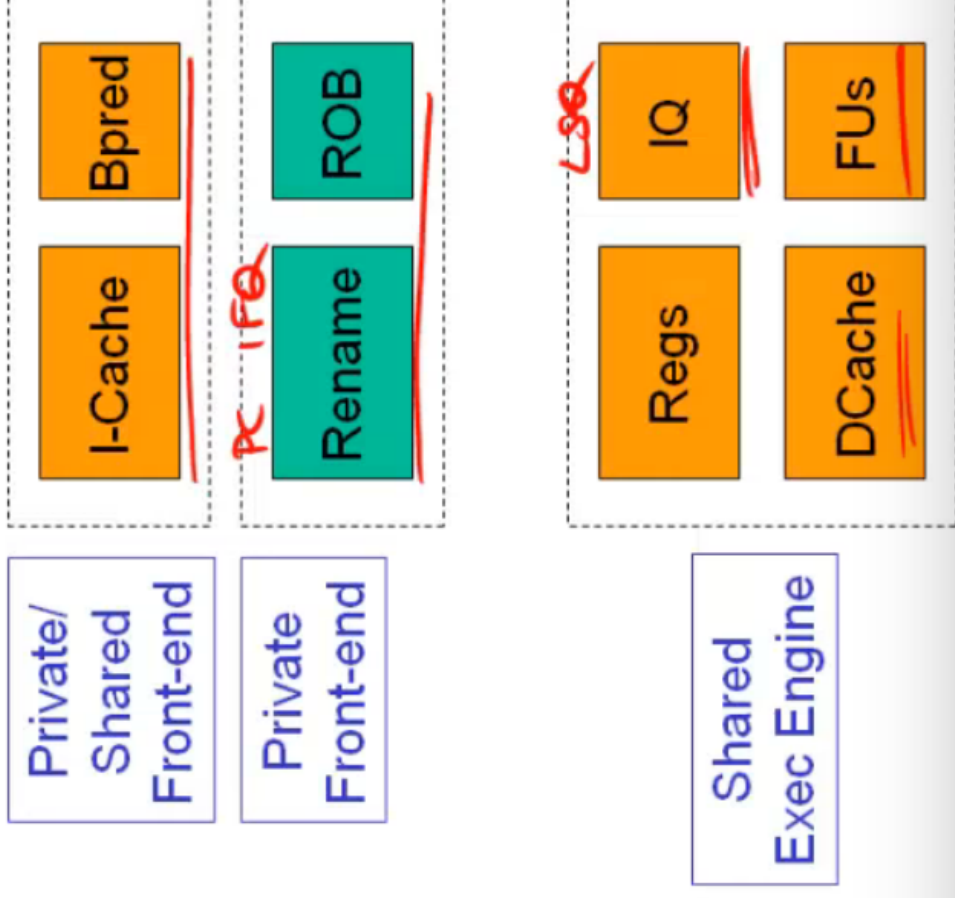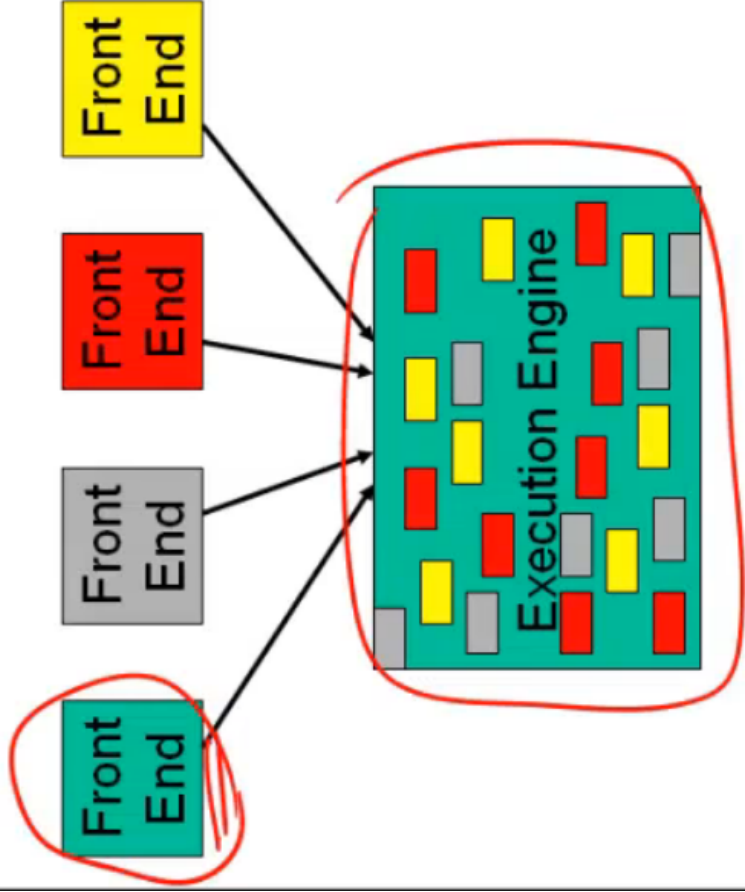
# How are Resources Shared?

Each box represents an issue slot for a functional unit. Peak thruput is 4 IPC.

Thread 1
Thread 2
Thread 3
Thread 4
Idle

Superscalar

Fine-Grained Multithreading

Simultaneous Multithreading

Cycles

- Superscalar processor has high under-utilization – not enough work every cycle, especially when there is a cache miss
- Fine-grained multithreading can only issue instructions from a single thread in a cycle – can not find max work every cycle, but cache misses can be tolerated
- Simultaneous multithreading can issue instructions from any thread every cycle – has the highest probability of finding work for every issue slot

11

# Pipeline Structure



Front End   Front End   Front End   Front End

Execution Engine

I-Cache   Bpred
Private/ Shared Front-end

Rename   ROB
PC   IFQ
Private Front-end

Regs   IQ   LSQ
DCache   FUs
Shared Exec Engine

9:13 / 11:07