# ON THE COST OF TYPE-TAG SOUNDNESS

Ben Greenman    Zeina Migeed

Northeastern University

# ON THE COST OF TYPE-TAG SOUNDNESS

# TYPE-TAG SOUNDNESS

# Type Soundness

If $\vdash e : \tau$ then either:

- $e \longrightarrow^* v$ and $\vdash v : \tau$

- $e$ diverges

- $e \longrightarrow^* \textbf{Error}$ (division by zero, etc.)

**No undefined behavior**

**Type-based reasoning**

# Type Soundness

If $\vdash e : \tau$ then either:

- $e \longrightarrow^* v$ and $\vdash v : \tau$

- $e$ diverges

- $e \longrightarrow^* \textbf{Error}$  (division by zero, etc.)

# Tag Soundness

If $\vdash e : \tau$ then either:

- $e \longrightarrow^* v$ and $\vdash v : \lfloor \tau \rfloor$

- $e$ diverges

- $e \longrightarrow^* \textbf{Error}$ (division by zero, etc.)

# Tag Soundness

If $\vdash e : \tau$ then either:

- $e \longrightarrow^* v$ and $\vdash v : \lfloor \tau \rfloor$

- $e$ diverges

- $e \longrightarrow^* $ **Error**  (division by zero, etc.)

$$\boxed{\lfloor \tau \rfloor = K}$$

$$\lfloor \text{Int} \rfloor = \text{Int}$$

$$\lfloor \tau \times \tau' \rfloor = \text{Pair}$$

$$\lfloor \tau \rightarrow \tau' \rfloor = \text{Fun}$$

$$\cdots$$

# Tag Soundness

If $\vdash \mathbf{e} : \mathbf{\tau}$ then either:

- $\mathbf{e} \longrightarrow^* \mathbf{v}$ and $\vdash \mathbf{v} : \lfloor \mathbf{\tau} \rfloor$

- $\mathbf{e}$ diverges

- $\mathbf{e} \longrightarrow^* \mathbf{Error}$

$$\boxed{\lfloor \mathbf{\tau} \rfloor = \mathbf{K}}$$

$$\lfloor \mathbf{Int} \rfloor = \mathbf{Int}$$

$$\lfloor \mathbf{\tau} \times \mathbf{\tau'} \rfloor = \mathbf{Pair}$$

$$\lfloor \mathbf{\tau} \rightarrow \mathbf{\tau'} \rfloor = \mathbf{Fun}$$

$$\ldots$$

**No undefined behavior**

**Tag-based reasoning**

# Types vs. Tags

If $\vdash$ **e** : **Int × Int** and **e** $\longrightarrow$* **v** then **v** might be:

| Type Soundness | Tag Soundness |
|---|---|
| (0, 0) | ("A", 0) |
| (3, 2) | (0, 0) |
| (-7, 9) | (3, 2) |
| | (-7, 9) |
| | (0, (1, 2)) |

# Types vs. Tags

If $\vdash$ **e : Int × Int** and **e $\longrightarrow^*$ v** then **v** might be:

| Type Soundness | Tag Soundness |
|---|---|
| (0, 0)  (3, 2) | (0, 0)  (3, 2)  ("A", 0) |
| (-7, 9) | (-7, 9)  (0, (1, 2)) |

$\longrightarrow^{*}$ **fast** $\qquad$ $\longrightarrow^{*}$ **slow**

|  | $\xrightarrow{}^{*}$ **fast** | $\xrightarrow{}^{*}$ **slow** |
| --- | --- | --- |
| Type Sound? | ✘ | ✔ |
| Tag Sound? | ✔ | ✔ |

# PERFORMANCE COST OF SOUNDNESS

# Problem: Safe Interaction

$$e^\tau \longrightarrow e^\tau \longrightarrow e^\tau \longrightarrow \tau$$

**?**

# Gradual Typing

# User Input

$$e^\tau \longrightarrow^* \text{read()}^\tau \longrightarrow \tau$$

Enter a value:
>

# Deserialization

$$e^\tau \longrightarrow^* \mathbf{unzip()}^\tau \longrightarrow \tau$$

```
0110
1110
1011
```

# Primitive Operations ( δ )

$$e^{Int} \longrightarrow^* \quad v + v^{Int} \longrightarrow \quad Int$$

$$E[v,v] \longrightarrow \ldots \longrightarrow v'$$

# Unreliable Source

# Option 1: Trust

$$e^\tau \longrightarrow \nu^\tau$$

?

# Option 2: Check

$$e^\tau \longrightarrow v^\tau \longrightarrow \ldots \longrightarrow v^\tau$$

**?**

# Option 2: Check

# Cost of Types $\left(\xrightarrow{*}_{\text{slow}}\right)$

$$\overset{\text{Int} \times \text{Int}}{(6,1)} \longrightarrow \overset{\text{Int} \times \text{Int}}{(6,1)} \longrightarrow \overset{\text{Int} \times \text{Int}}{(6,1)}$$

$$\longrightarrow \overset{\text{Int} \times \text{Int}}{(6,1)}$$

( ? )

# Cost of Tags $\left( \xrightarrow{\quad}^{*}_{\text{fast}} \right)$

$$\text{Int} \times \text{Int} \qquad \text{Int} \times \text{Int}$$
$$(6,1) \longrightarrow (6,1)$$

?

# COST OF SOUNDNESS IN RETICULATED

# Retic vs. Python

$$e^\tau \longrightarrow e^\tau \longrightarrow e^\tau \longrightarrow \tau$$

# Reticulated

```
def dist(pt : Tuple(Int,Int)) -> Int:
    x = pt[0]
    y = pt[1]
    return abs(x + y)
```

# Reticulated

```
def dist(pt : Tuple(Int,Int)) -> Int:
    x = pt[0]
    y = pt[1]
    return abs(x + y)
```

$$dist((0, 0)) \longrightarrow^* 0$$

# Reticulated

```
def dist(pt : Tuple(Int,Int)) -> Int:
    x = pt[0]
    y = pt[1]
    return abs(x + y)
```

dist("NaN") ⟶* Expected **Tuple**

# Reticulated

```
def dist(pt : Tuple(Int,Int)) -> Int:
    x = pt[0]
    y = pt[1]
    return abs(x + y)
```

dist((0, "NaN"))  ⟶*  Expected Int

# Evaluation Method

# 1. Fully-Typed

# 2. Configurations

# 3. Measure

|  | 11s |  |
|---|---|---|
| 7s | 9s | 2s |
| 5s | 24s | 9s |
| 14s | 5s | 21s |
| 9s | 6s | 9s |
| 8s | 4s | 5s |

**What % have at most Dx overhead?**

D = 4, vs.

|  | 11s |  |
|-----|------|-----|
| 7s | 9s | 2s |
| 5s | 24s | 9s |
| 14s | 5s | 21s |
| 9s | 6s | 9s |
| 8s | 4s | 5s |

# Evaluation Method

# EXPERIMENT & RESULTS

# Benchmarks

| DLS 2014 | POPL 2017 | PEPM 2018 |
|---|---|---|
| futen | call_method | espionage |
| http2 | call_simple | pythonflow |
| slowSHA | chaos | take5 |
| aespython | fannkuch | sample_fsm |
| stats | go | |
| | meteor | |
| | nbody | |
| | nqueens | |
| | pidigits | |
| | pystone | |
| | spectralnorm | |

# # Typed Components

| DLS 2014 | POPL 2017 | PEPM 2018 |
|---|---|---|
| 15 | 7 | 12 |
| 4 | 6 | 12 |
| 17 | 15 | 16 |
| 34 * | 1 | 19 * |
| 79 * | 7 | |
| | 8 | |
| | 5 | |
| | 2 | |
| | 5 | |
| | 14 | |
| | 5 | |

# Exhaustive Results

What % of configurations have at most **4x** overhead?

# Exhaustive Results

What % of configurations have at most **Dx** overhead?

# Approximate Results

What % of configurations have at most **4x** overhead, based on **R** samples of **S** configurations each?

# Approximate Results

What % of configurations have at most **Dx** overhead, based on **R** samples of **S** configurations each?

# Approximate Results

What % of configurations have at most **Dx** overhead, based on **R** samples of **S** configurations each?
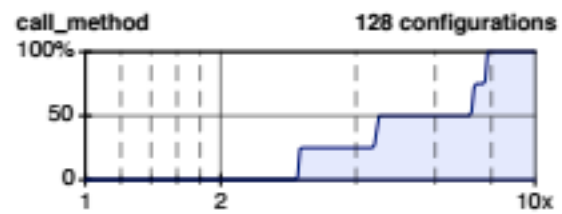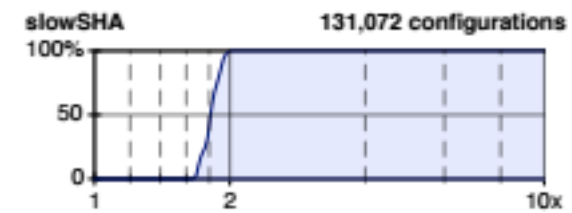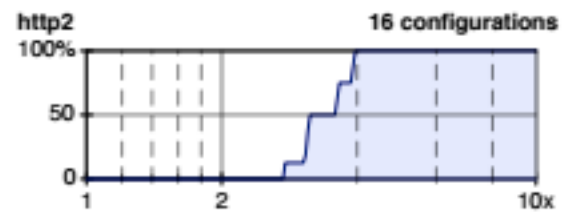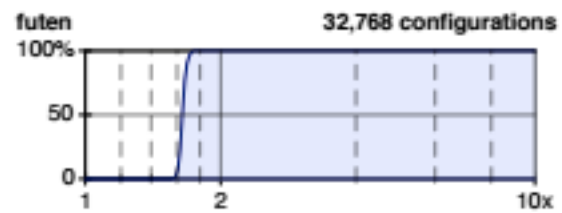
espionage — 4,096 configurations

aespython — 10 samples of 340 configurations

# Cost of Tag Soundness

- Worst-case overhead: under 10x

This is an APPLES to ORANGES comparison!

Type Soundness

Tag Soundness

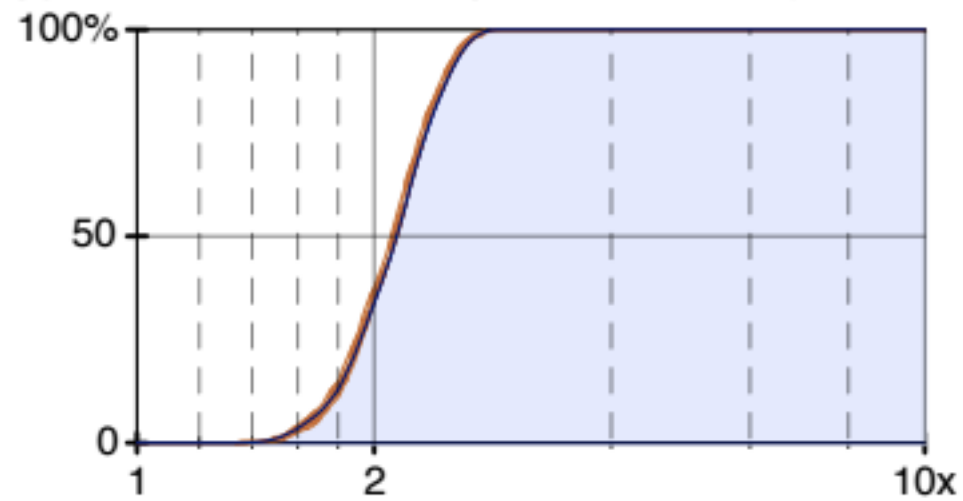1x 1x 2x 4x 1x 1x 5x 1x 5x 27x 29x 32x 10x 34x 43x 47x 292x 233x 139x 1527x

1x 1x 1x 1x 3x 1x 2x 2x 2x 3x 3x 2x 5x 3x 7x 6x 6x 7x 7x 8x
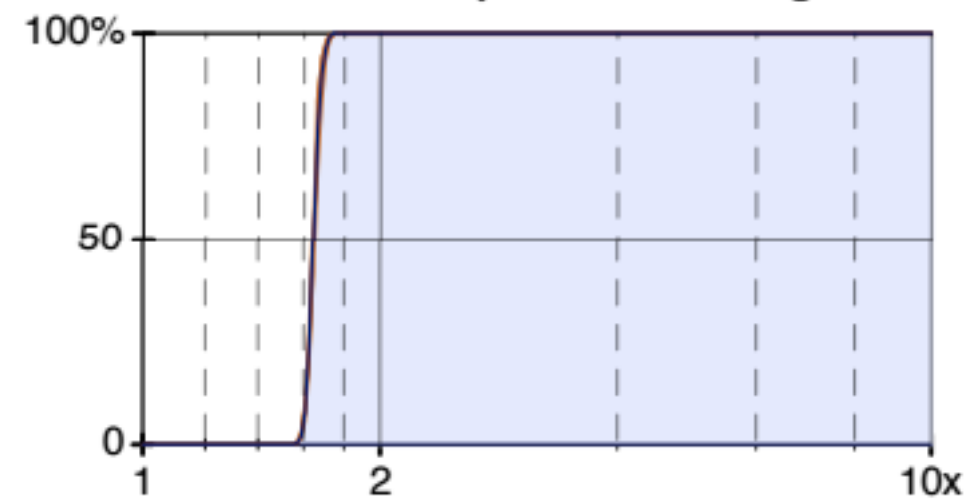
# Cost of Tag Soundness

- Worst-case overhead: under 10x

- Best-case overhead: 1x -- 4x

  - adding types never* improves performance

- Slowest configuration: fully-typed
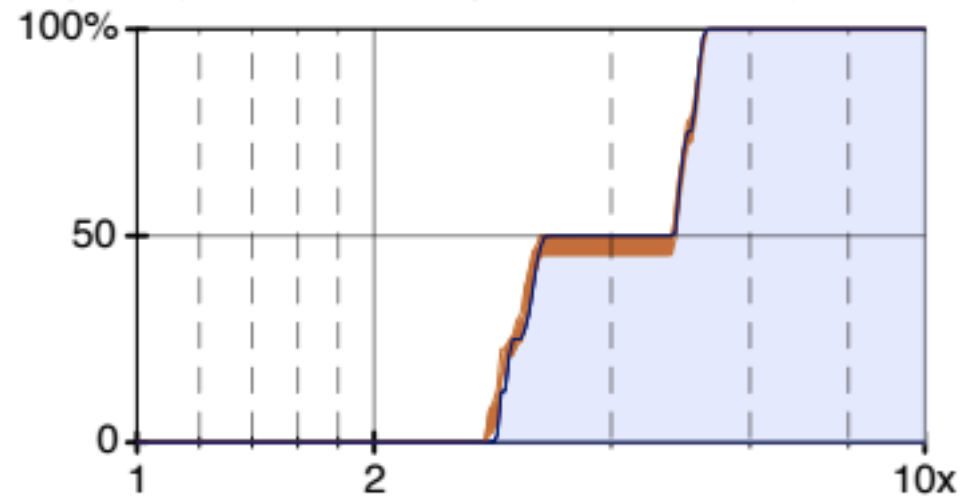
  - Overhead $\propto$ number of type annotations

# Runtime vs. # Types



spectralnorm

# Runtime vs. # Types



call_method

# Speedup?



- Unsound optimization for read-only values (tuples)

- Double-checks method calls

# Runtime vs. # Types

# Experiment

- granularity: functions & class-fields

- 10 samples of [10 * (F + C)] configurations

- Karst at Indiana University cluster (32GB RAM, 250GB other)

- Reticulated, master branch, commit e478343

- Python 3.4.3

- 40 iterations per configuration, report average

- 200 values of D on x-axis

# POPL 2017



**Figure 11.** Runtime comparison of Reticulated Python to standard Python 3.4. Experiments were performed on an Ubuntu 14.04 laptop with a 2.8GHz Intel i7-3840QM CPU and 16GB memory.

# Module Marshal

`module Marshal: sig .. end`

Marshaling of data structures.

This module provides functions to encode arbitrary data structures as sequences of bytes, which can then be written on a file or sent over a pipe or network connection. The bytes can then be read back later, possibly in another process, and decoded back into a data structure. The format for the byte sequences is compatible across all machines for a given version of OCaml.
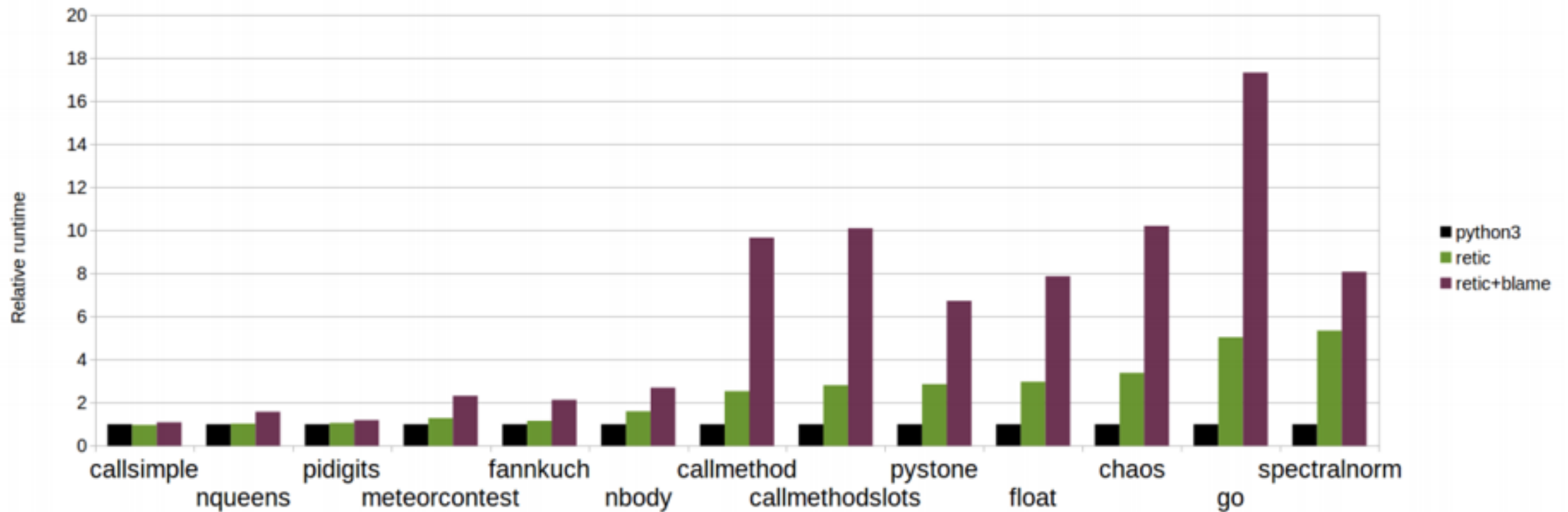
Warning: marshaling is currently not type-safe. The type of marshaled data is not transmitted along the value of the data, making it impossible to check that the data read back possesses the type expected by the context. In particular, the result type of the `Marshal.from_*` functions is given as `'a`, but this is misleading: the returned OCaml value does not possess type `'a` for all `'a`; it has one, unique type which cannot be determined at compile-time. The programmer should explicitly give the expected type of the returned value, using the following syntax:

- `(Marshal.from_channel chan : type)`. Anything can happen at run-time if the object in the file does not belong to the given type.

# References

- Vitousek, Swords, Siek. *Big Types in Little Runtime: Open-World Soundness and Collaborative Blame for Gradual Type Systems.* POPL 2017

- Takikawa, Feltey, Greenman, New, Vitek, Felleisen. *Is Sound Gradual Typing Dead?* POPL 2016.

# UNUSED SLIDES

$$\frac{\overline{\phantom{xxxxxxxxxxxxxxxx}}}{\vdash e' : \tau'} \quad ???$$

$$\vdots \qquad \vdots$$

$$\frac{\phantom{xxxxxxxxxxxx}}{\vdash e : \tau}$$

# Takikawa Method

- granularity

- experimental modules, fixed modules

- configurations

- baseline

- performance ratio