

# Benchmarking: Are We Doing it Wrong?



Edd Barrett



Carl  
Friedrich  
Bolz-Tereick



Rebecca  
Killick  
(Lancaster)



Sarah Mount



Laurence  
Tratt

KING'S  
*College*  
LONDON

Software Development Team  
2017-05-27

# Introduction

To appear: OOPSLA 2017, Vancouver.

<https://arxiv.org/abs/1602.00602>

1

## Virtual Machine Warmup Blows Hot and Cold

EDD BARRETT, King's College London

CARL FRIEDRICH BOLZ-TEREICK, King's College London

REBECCA KILLICK, Lancaster University

SARAH MOUNT, King's College London

LAURENCE TRATT, King's College London

---

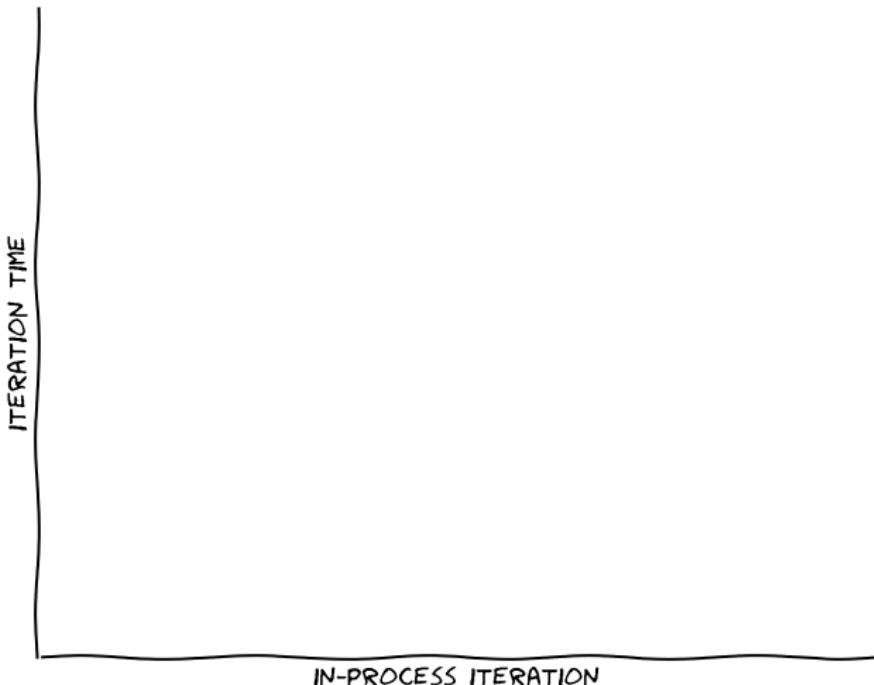
Virtual Machines (VMs) with Just-In-Time (JIT) compilers are traditionally thought to execute programs in two phases: the initial warmup phase determines which parts of a program would most benefit from dynamic compilation, before JIT compiling them into machine code; subsequently the program is said to be at a steady state of peak performance. Measurement methodologies almost always discard data collected during the warmup phase such that reported measurements focus entirely on peak performance. We introduce a fully automated statistical approach, based on changepoint analysis, which allows us to determine if a program has reached a steady state and, if so, whether that represents peak performance or not. Using this, we show that even when run in the most controlled of circumstances, small, deterministic, widely studied microbenchmarks often fail to reach a steady state of peak performance on a variety of common VMs. Repeating our experiment on 3 different machines, we found that at most 43.5% of (VM, benchmark) pairs consistently reach a steady state of peak performance.

CCS Concepts: •Software and its engineering → Software performance; Just-in-time compilers; *Inference engines*

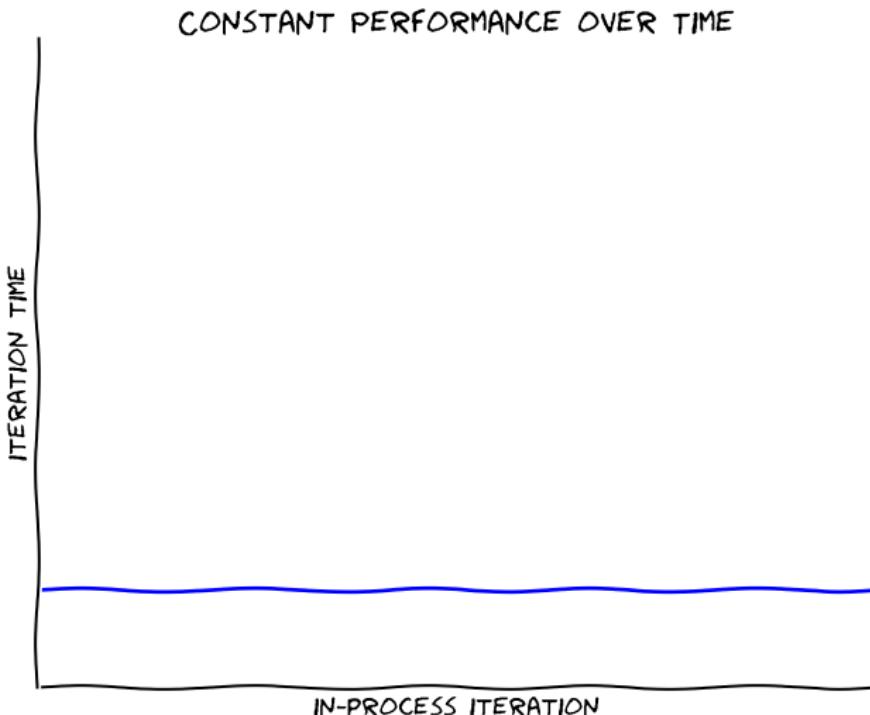
# Programming Language Benchmarking

- ▶ Used to decide if a language is performing well.
- ▶ Used by pretty much all language designers:
  - ▶ Google
  - ▶ Oracle
  - ▶ Facebook
  - ▶ Mozilla
  - ▶ ...
- ▶ Long established methods.

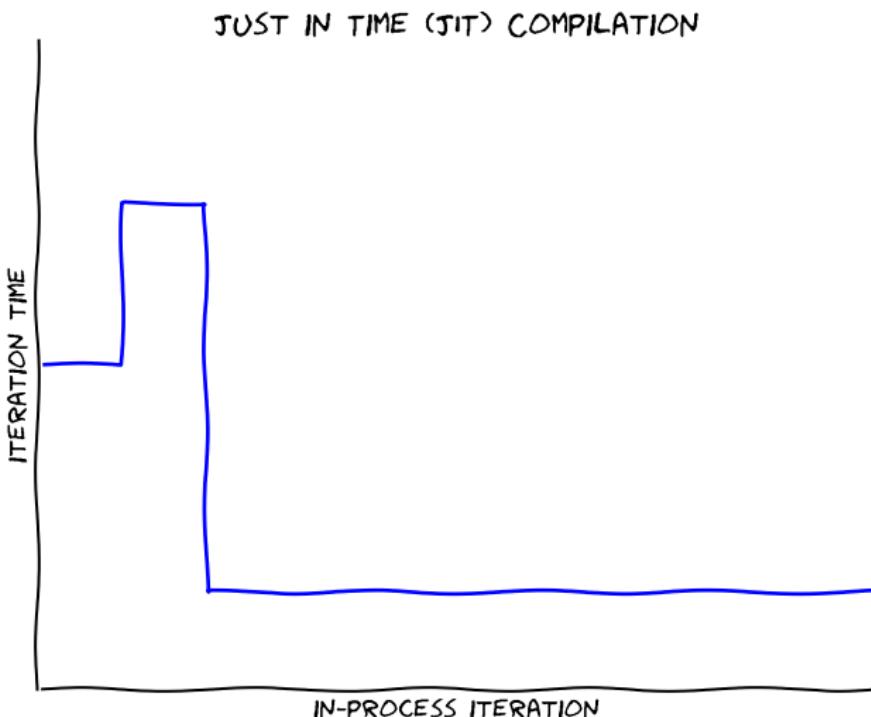
# The Current State of the Art of Benchmarking



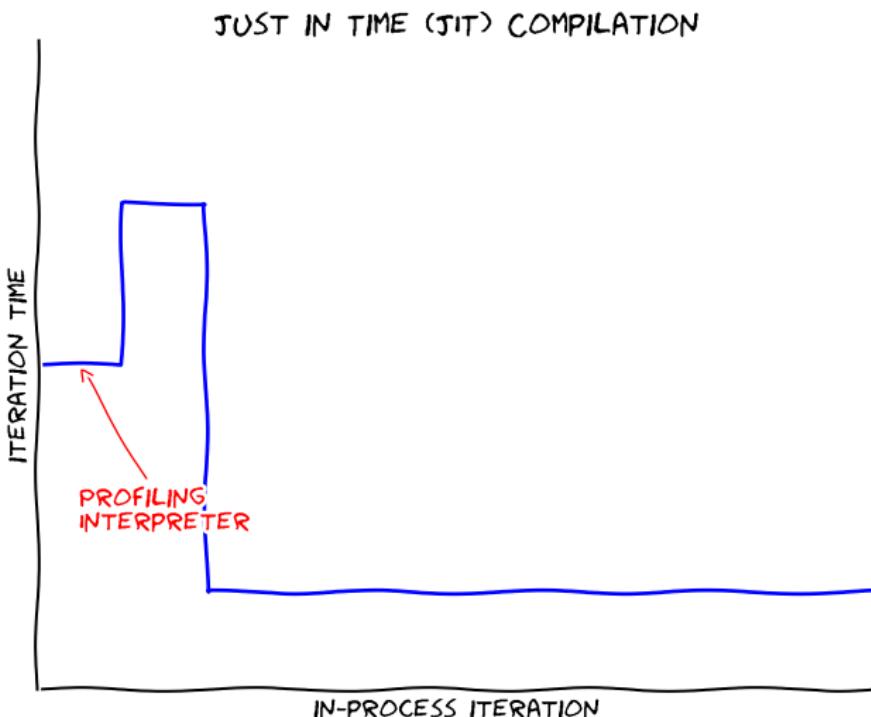
# The Current State of the Art of Benchmarking



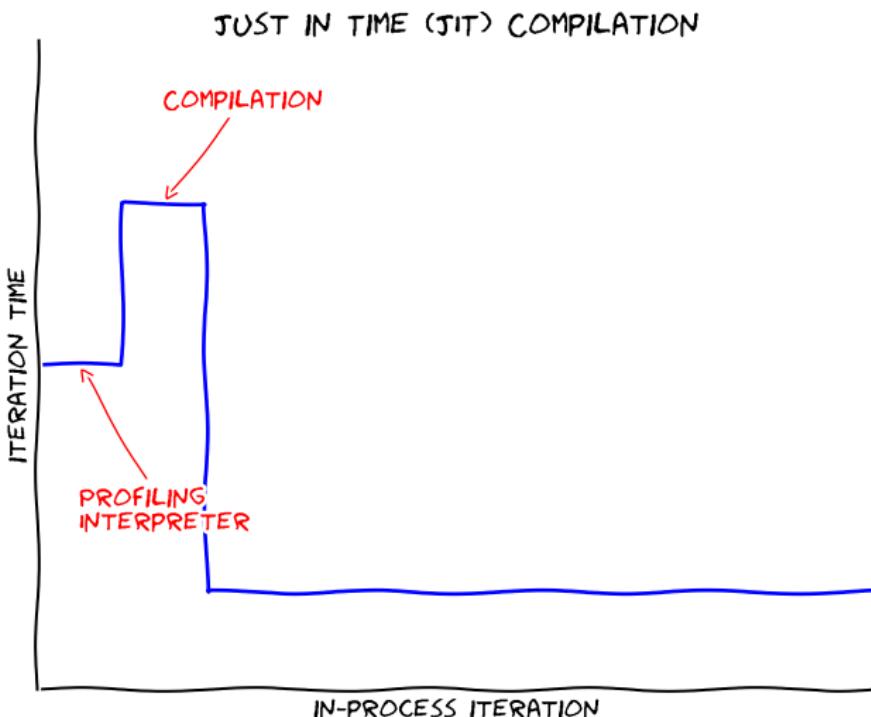
# The Current State of the Art of Benchmarking



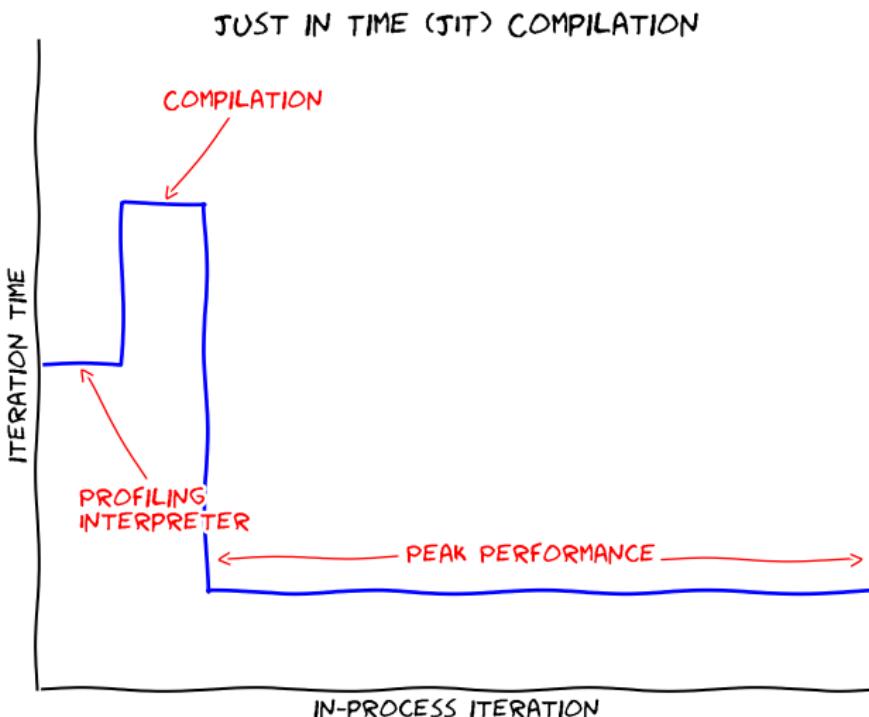
# The Current State of the Art of Benchmarking



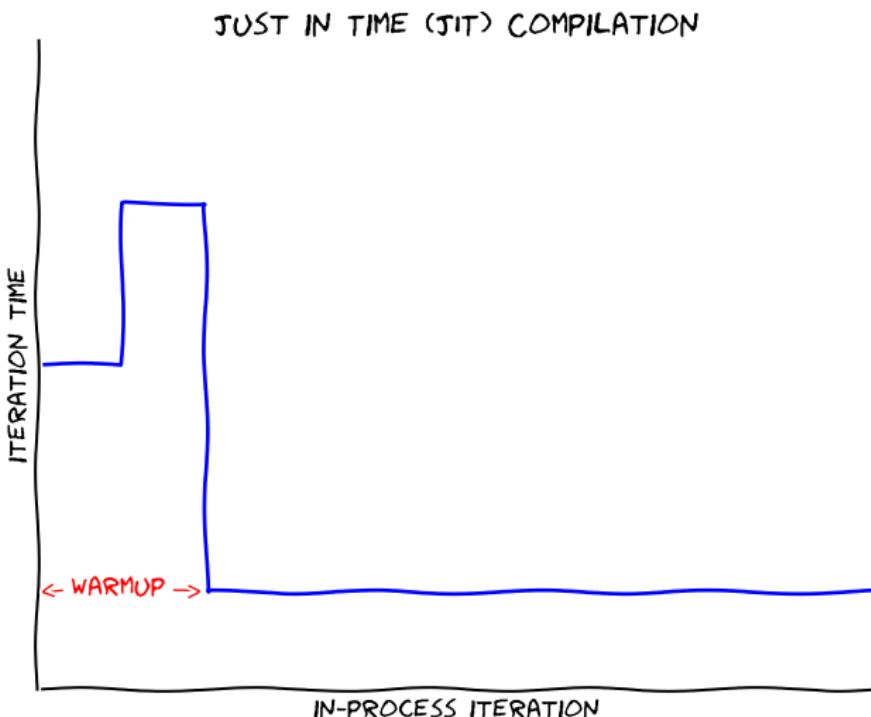
# The Current State of the Art of Benchmarking



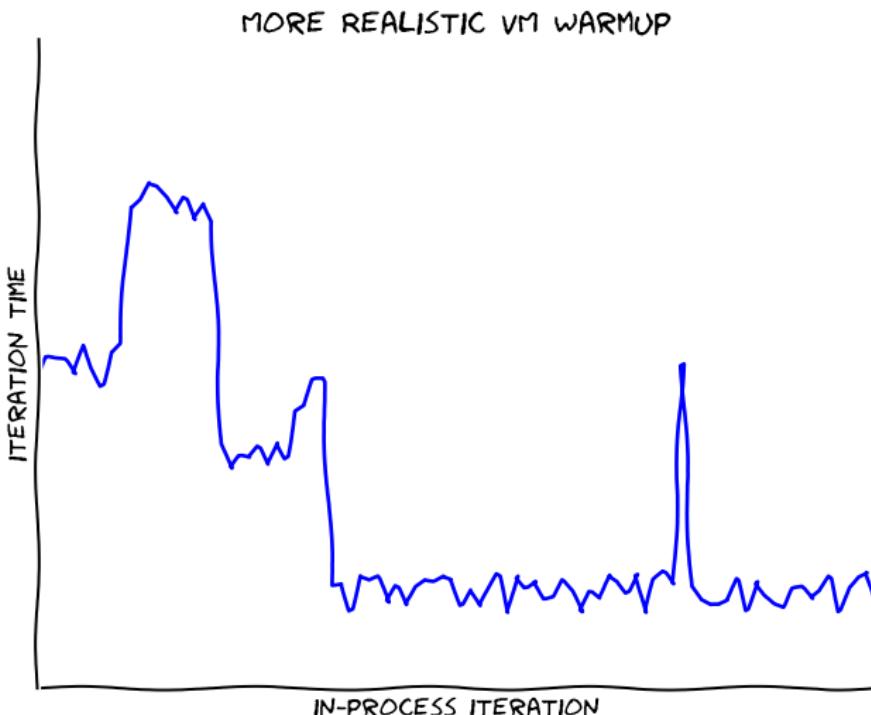
# The Current State of the Art of Benchmarking



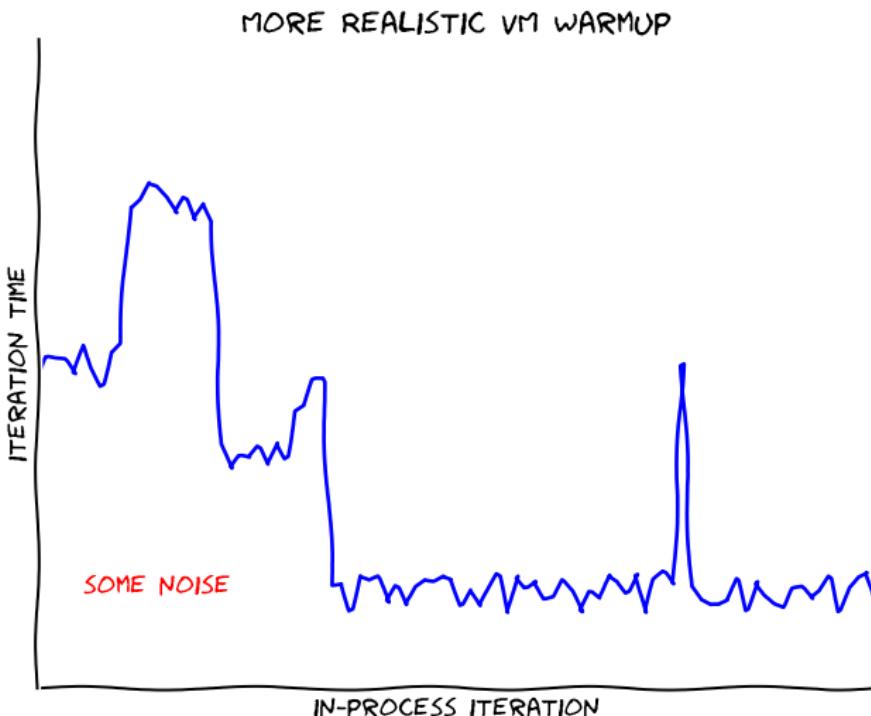
# The Current State of the Art of Benchmarking



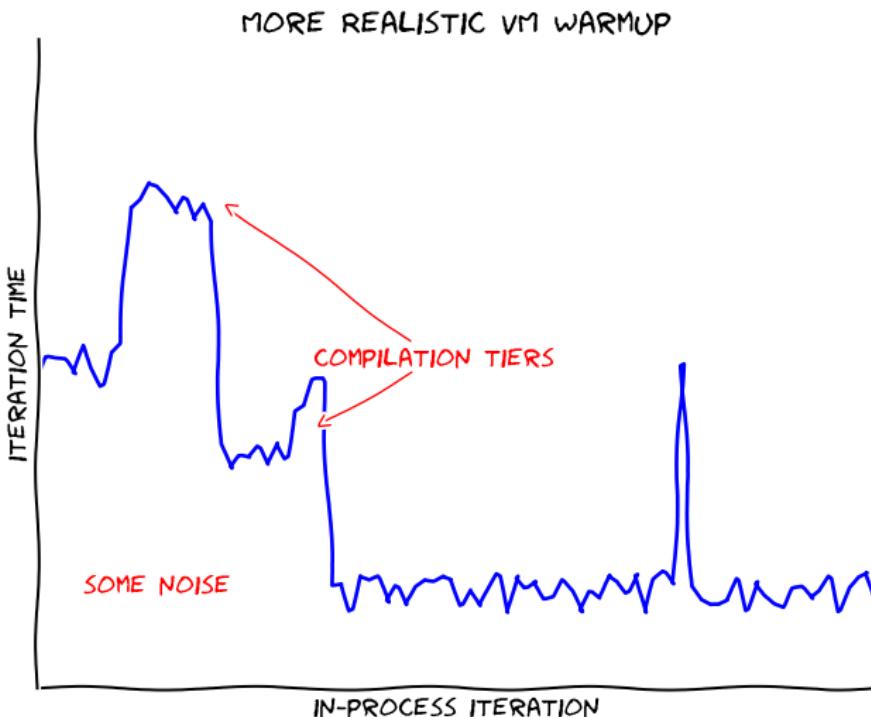
# The Current State of the Art of Benchmarking



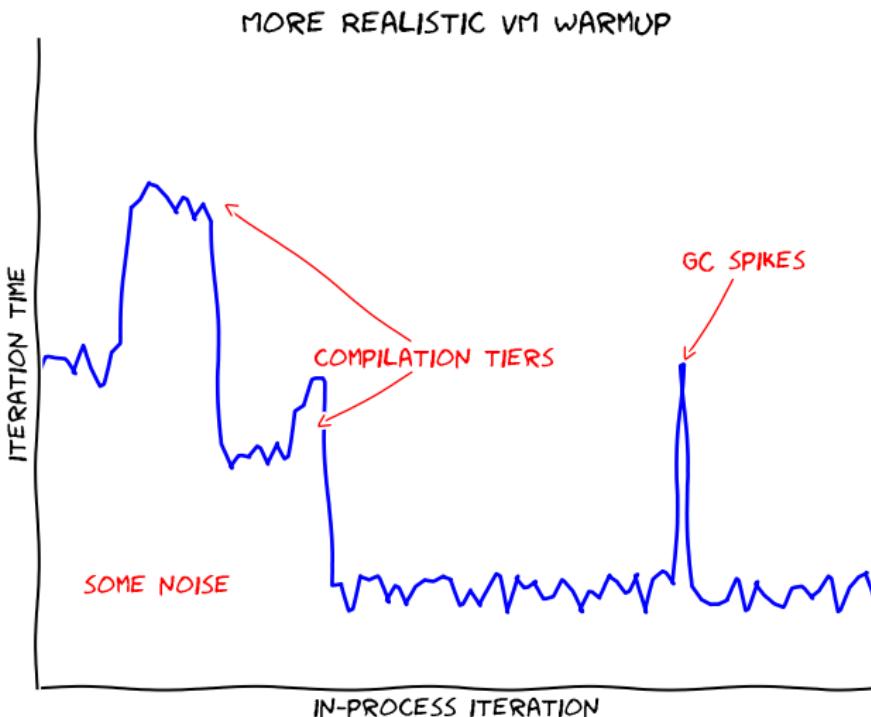
# The Current State of the Art of Benchmarking



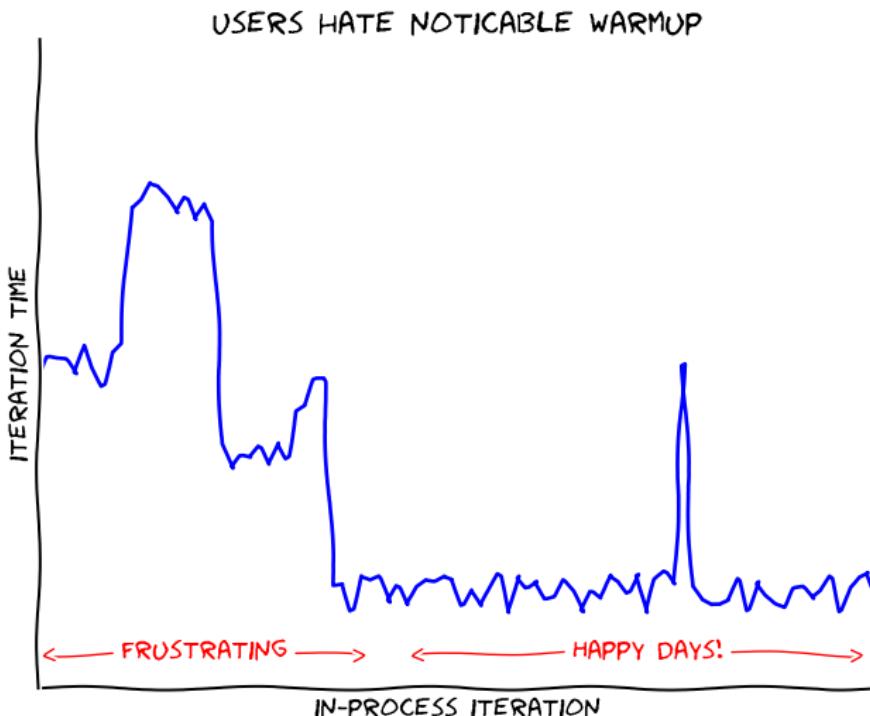
# The Current State of the Art of Benchmarking



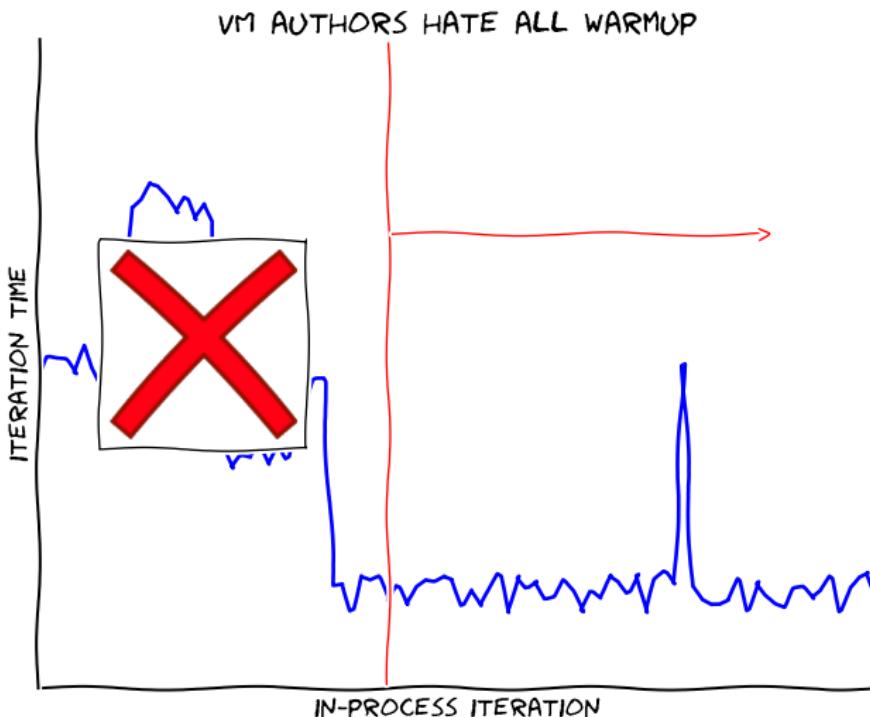
# The Current State of the Art of Benchmarking



# Warmup Matters



# Warmup Matters



# The Warmup Experiment

We should measure the warmup of modern language implementations

# The Warmup Experiment

We should measure the warmup of modern language implementations

*Hypothesis:* Small, deterministic programs reach a steady state of peak performance.

## Method 1: Which benchmarks?

The CLBG benchmarks are perfect for us (unusually)

<http://benchmarks.game.alioth.debian.org/>

## Method 1: Which benchmarks?

The CLBG benchmarks are perfect for us (unusually)

<http://benchmarks.game.alioth.debian.org/>

We removed any CFG non-determinism

## Method 1: Which benchmarks?

The CLBG benchmarks are perfect for us (unusually)

<http://benchmarks.game.alioth.debian.org/>

We removed any CFG non-determinism

We added checksums to all benchmarks

## Method 2: How long to run?

2000 *in-process iterations*

## Method 2: How long to run?

2000 *in-process iterations*

30 *process executions*

## Method 3: VMs

- Graal-0.22
- HHVM-3.19.1
- JRuby/Truffle (git #6e9d5d381777)
- Hotspot-8u121b13
- LuaJit-2.0.4
- PyPy-5.7.1
- V8-5.8.283.32
- GCC-4.9.4

Note: same GCC (4.9.4) used for all compilation

## Method 4: Machines

- Linux<sub>4790</sub>, Debian 8, 24GiB RAM
- Linux<sub>E3-1240v5</sub>, Debian 8, 32GiB RAM
- OpenBSD<sub>4790</sub>, OpenBSD 6.0, 32GiB RAM

## Method 4: Machines

- Linux<sub>4790</sub>, Debian 8, 24GiB RAM
  - Linux<sub>E3-1240v5</sub>, Debian 8, 32GiB RAM
  - OpenBSD<sub>4790</sub>, OpenBSD 6.0, 32GiB RAM
- 
- Turbo boost and hyper-threading disabled
  - Network card turned off.
  - Daemons disabled (crond, smtpd, sshd, atd, ...)

## Method 5: Benchmark Harness

KRUN

Control as many confounding variables as possible

## Method 5: Benchmark Harness

### KRUN

Control as many confounding variables as possible

- Minimises I/O
- Sets fixed heap and stack ulimits
- Drops privileges to a 'clean' user account
- Automatically reboots the system prior to each proc. exec
- Checks `dmesg` for changes after each proc. exec
- Checks system at (roughly) same temperature for proc. execs
- Enforces kernel settings (tickless mode, CPU governors, ...)

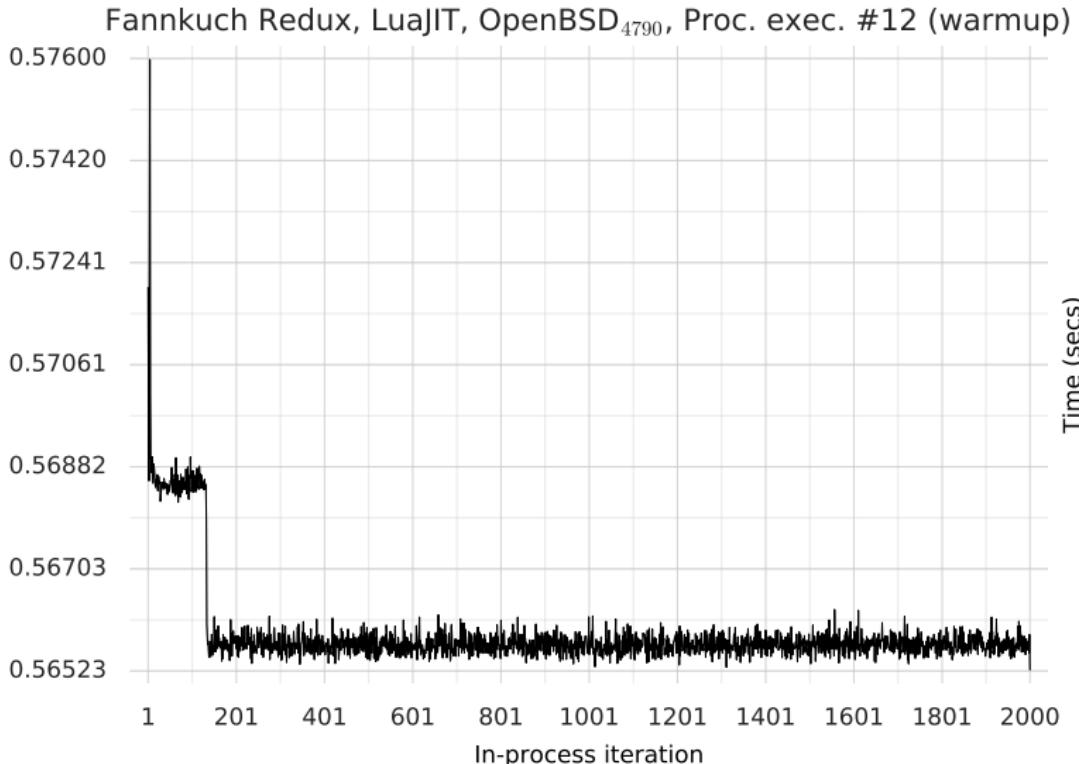
## Results and Classification

Let's look at some plots for our results.

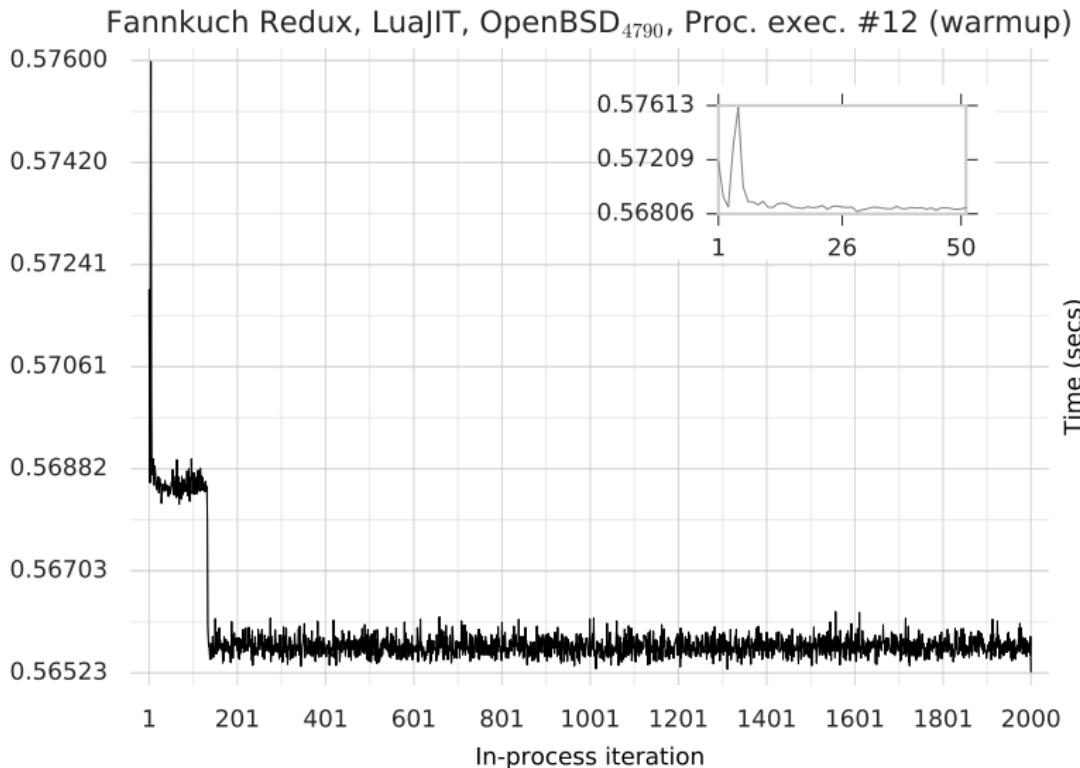
We also classify each process execution.

Classification uses *changepoint analysis*.

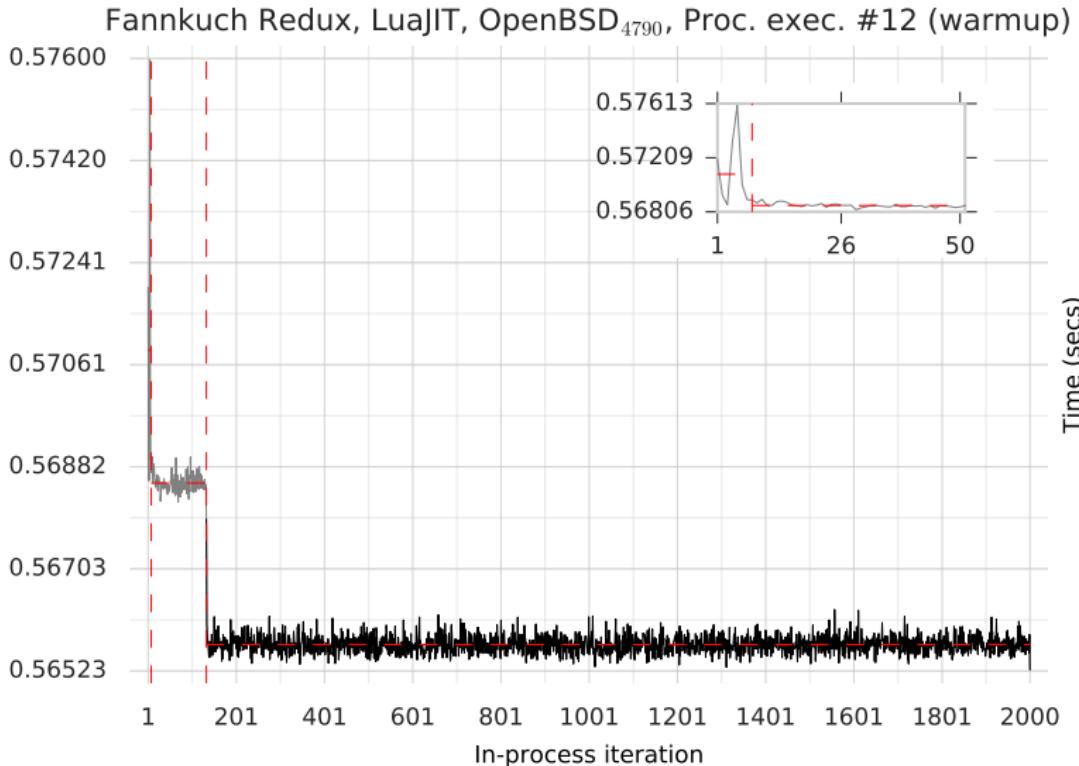
# Changepoint Analysis 101



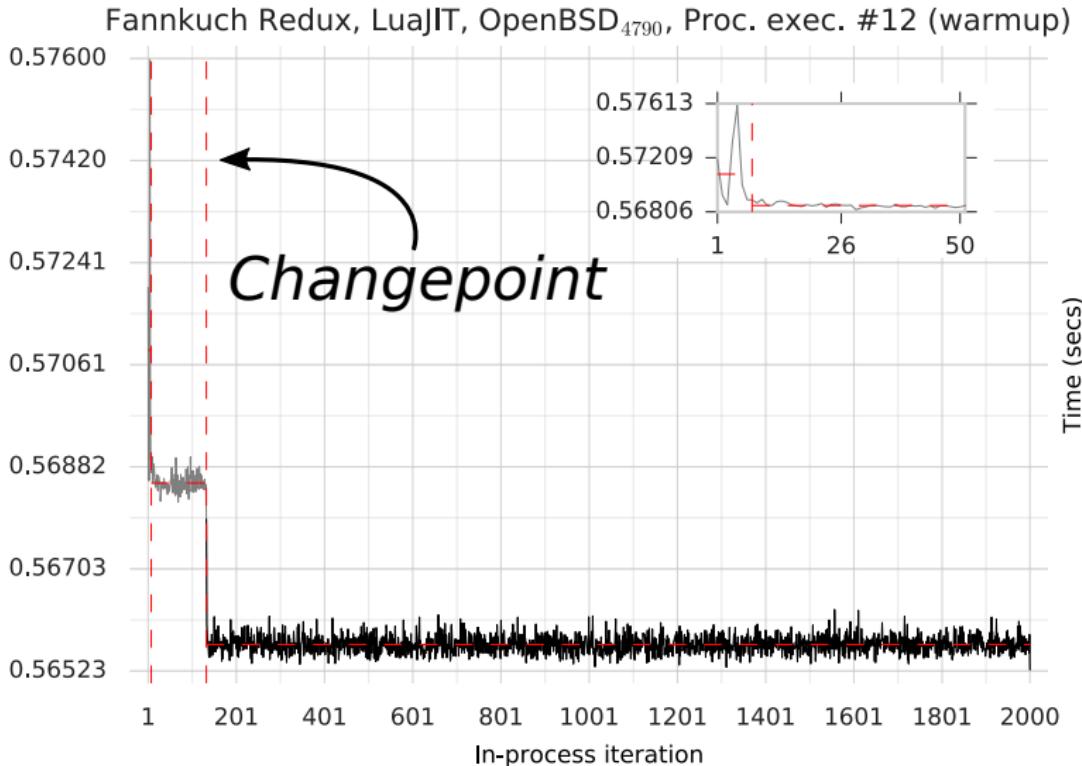
# Changepoint Analysis 101



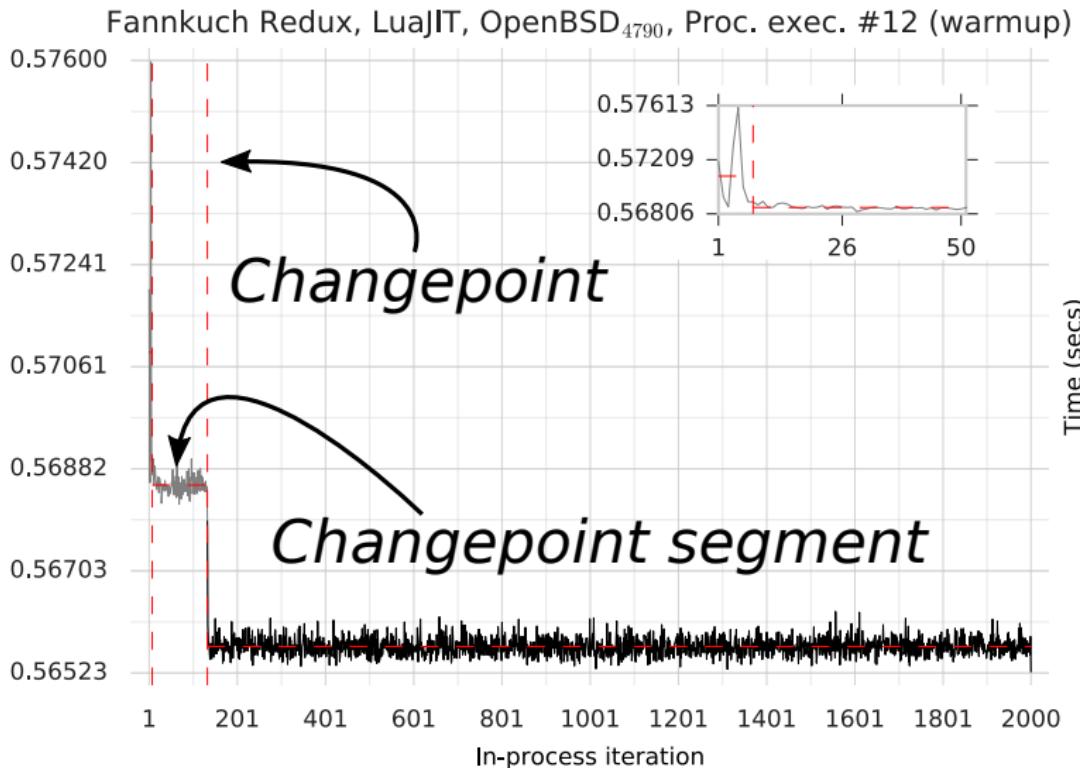
# Changepoint Analysis 101



# Changepoint Analysis 101



# Changepoint Analysis 101

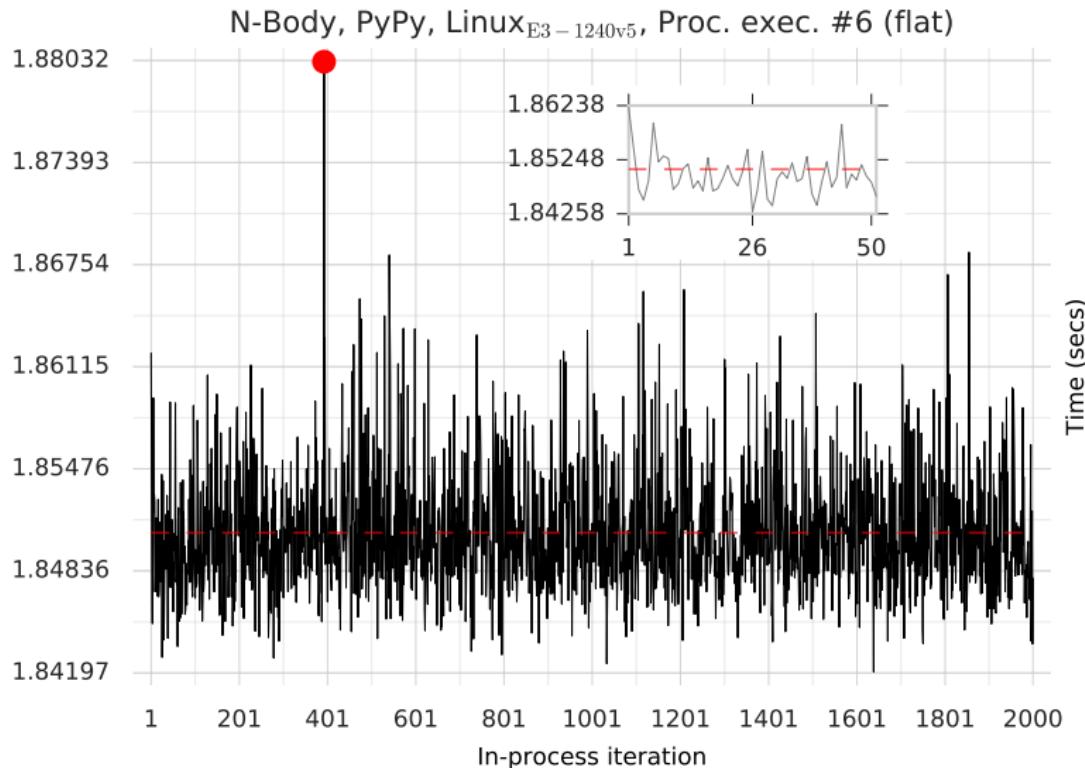


# Results and Classification

- ▶ All segs are equivalent:

FLAT

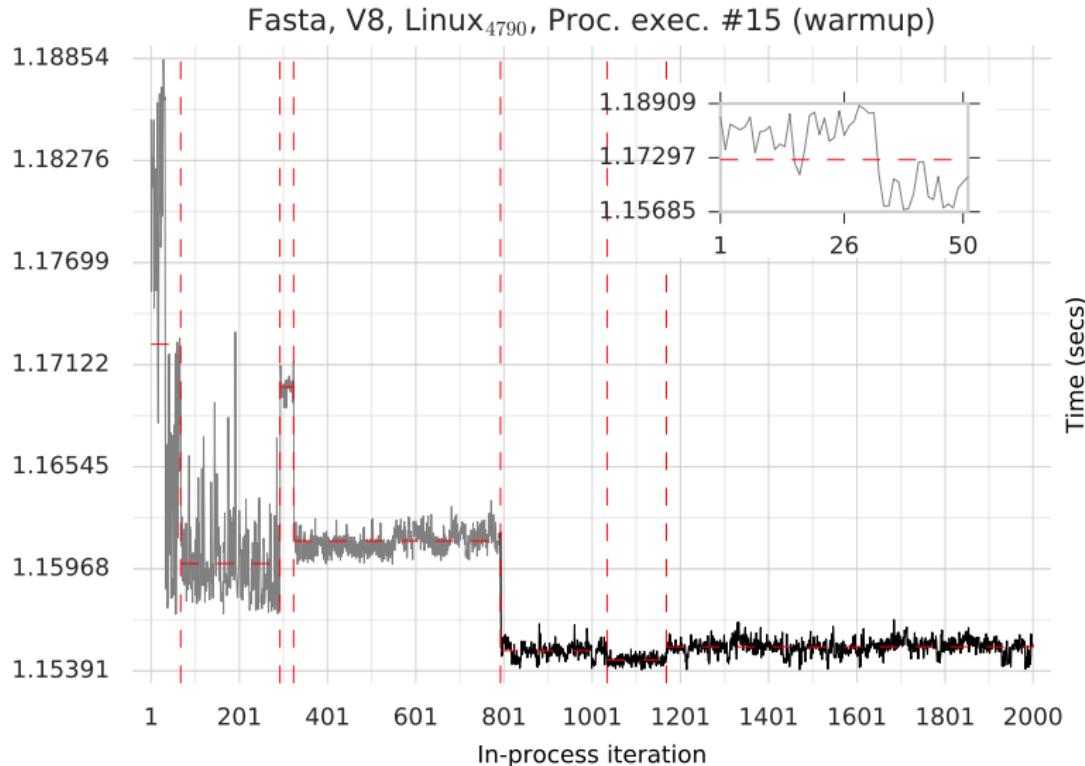
# Results and Classification: Flat



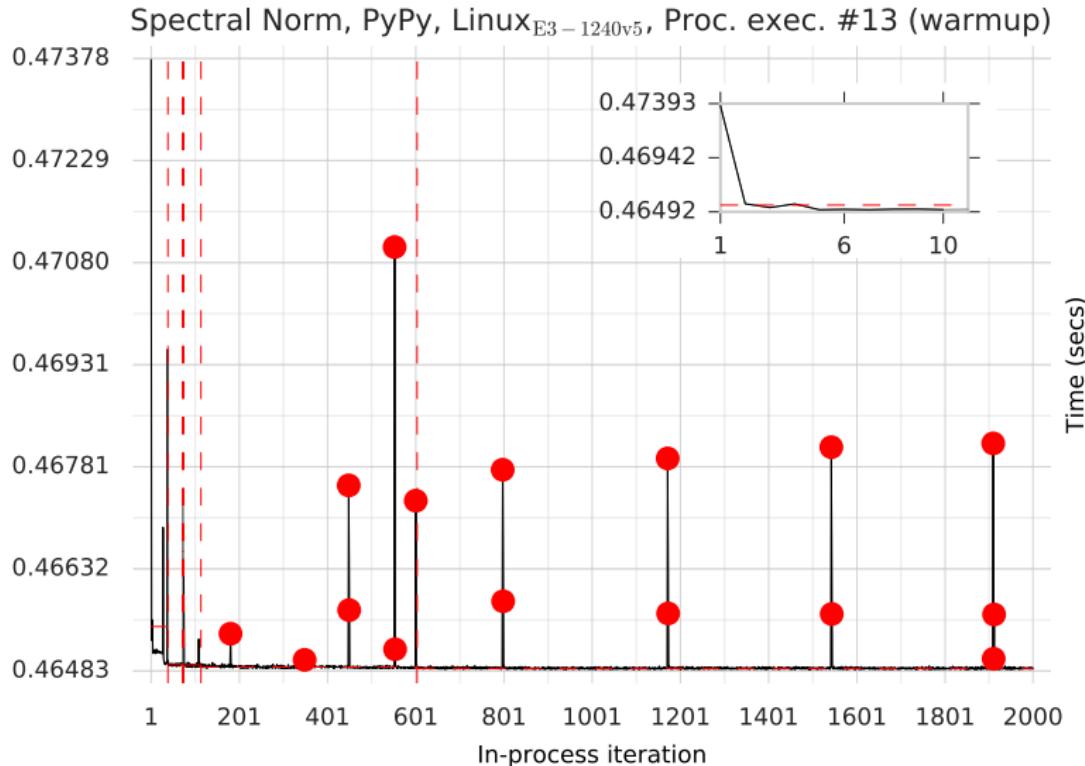
## Results and Classification

- ▶ All segs are equivalent:  
**FLAT**
- ▶ Final equivalent segs  $\geq$  500 iters:
  - ▶ Final seg is equivalent to the fastest seg:  
**WARMUP**

# Results and Classification: Warmup



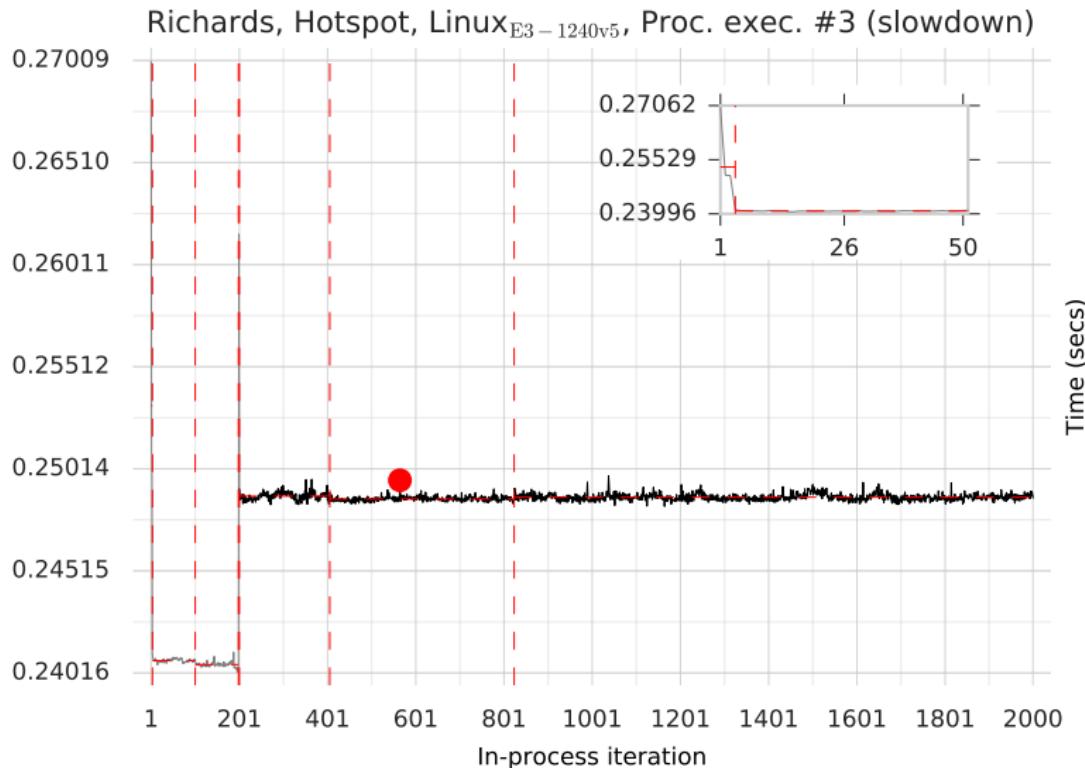
# Results and Classification: Warmup



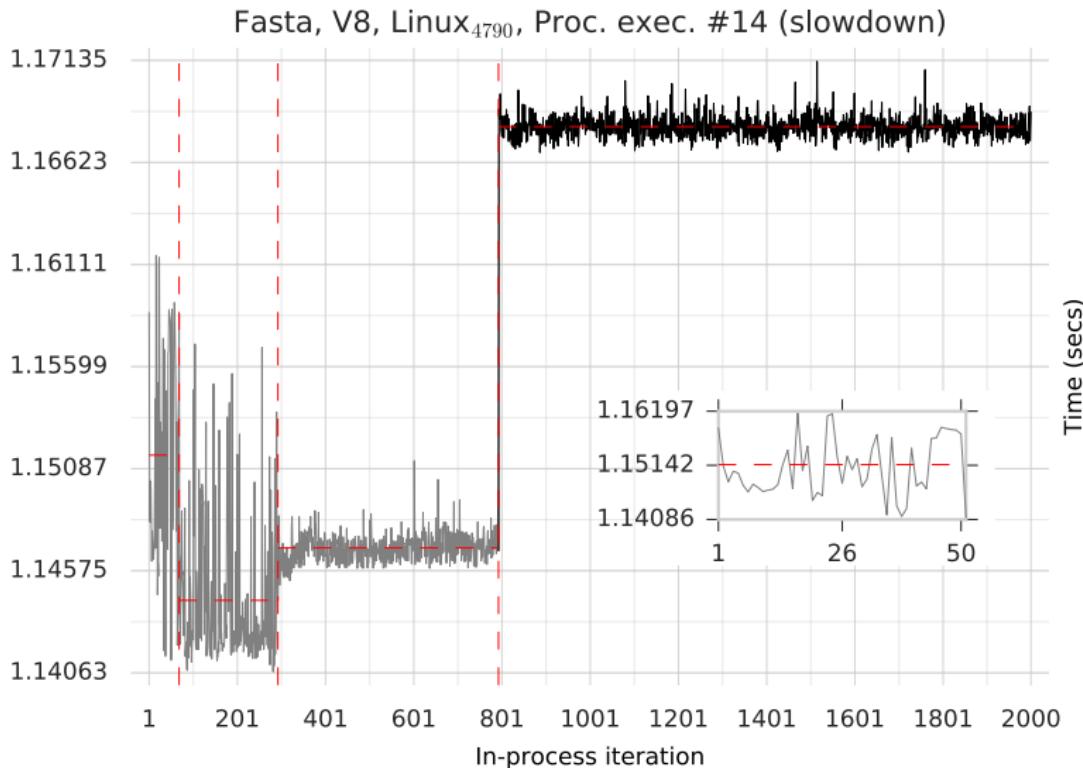
## Results and Classification

- ▶ All segs are equivalent:  
**FLAT**
- ▶ Final equivalent segs  $\geq$  500 iters:
  - ▶ Final seg is equivalent to the fastest seg:  
**WARMUP**
  - ▶ Final seg is not equivalent to the fastest seg:  
**SLOWDOWN**

# Results and Slowdown



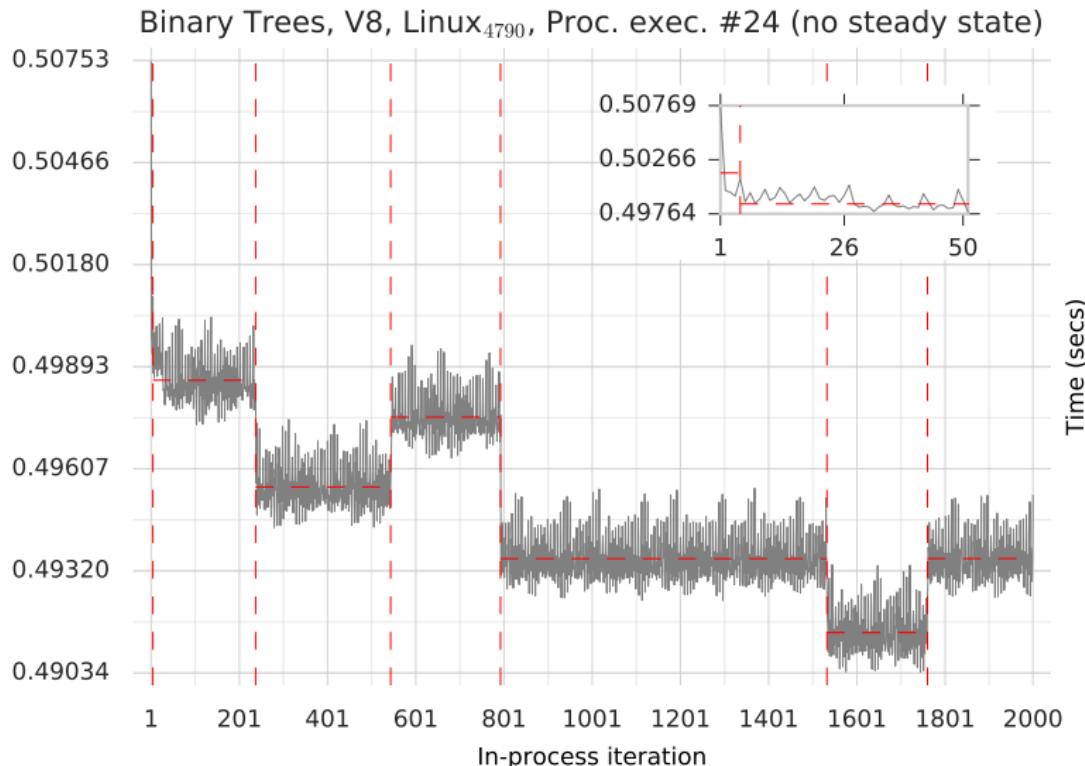
# Results and Slowdown



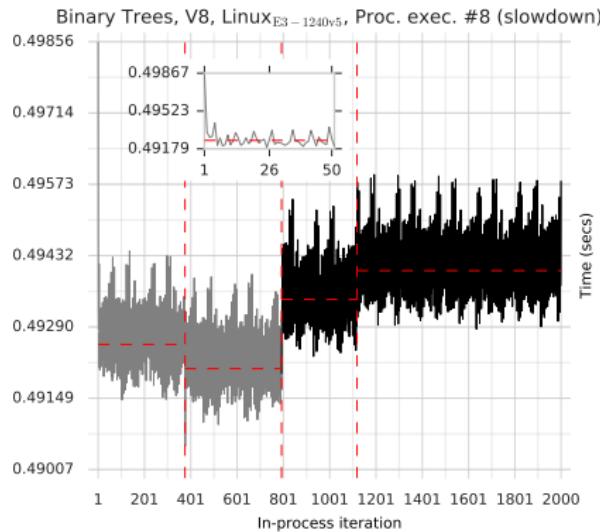
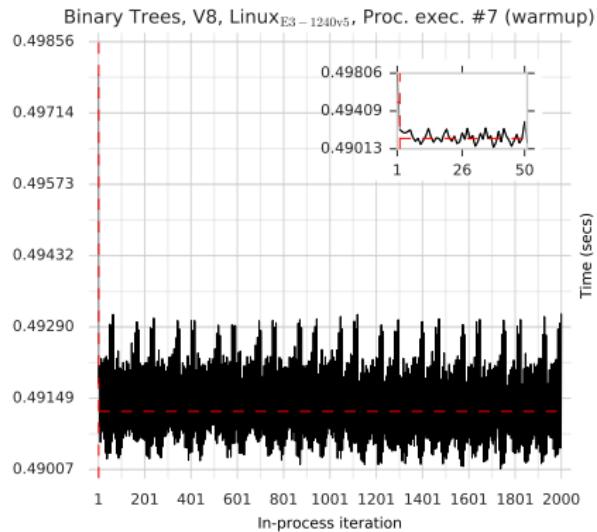
## Results and Classification

- ▶ All segs are equivalent:  
**FLAT**
- ▶ Final equivalent segs  $\geq$  500 iters:
  - ▶ Final seg is equivalent to the fastest seg:  
**WARMUP**
  - ▶ Final seg is not equivalent to the fastest seg:  
**SLOWDOWN**
- ▶ Otherwise:  
**No STEADY STATE**

# Results and Classification: No steady state

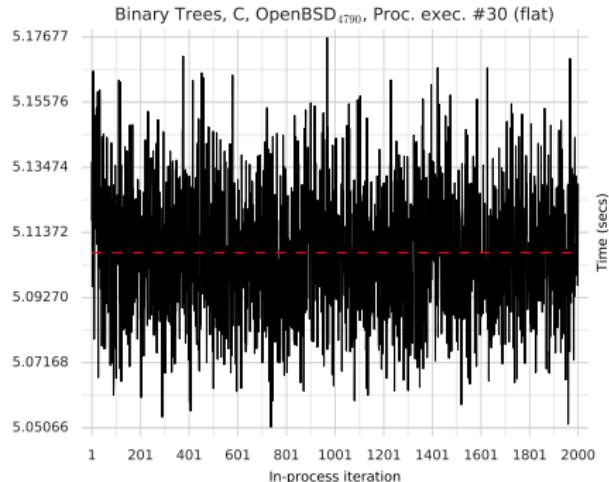
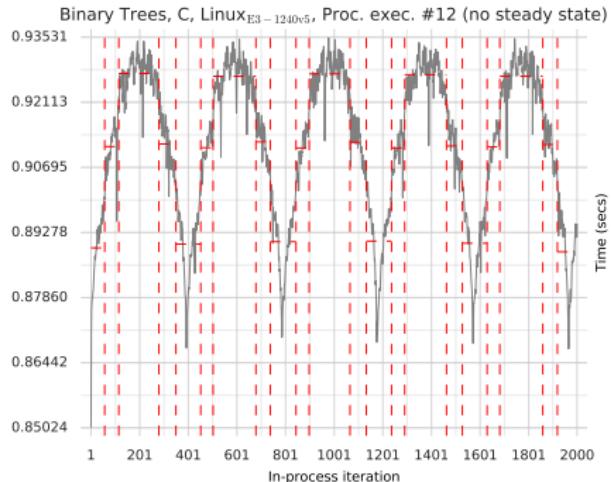


# Inconsistent Process-executions



(Same machine)

# Inconsistent Process-executions



(Different machines. Bouncing ball Linux-specific)

# Quantitative Results

	Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)		Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)
binary trees	C	~							0.40555 ±0.00510
	Graal	✗ (27L, 3L)	32.0 (17.6, 193.8)	6.60 —	0.18594 ±0.000315	L	8.0 (7.5, 8.5)	1.22 (1.126, 1.358)	0.13334 ±0.00045
	HHVM	✗ (24L, 4L, 2w)				✗ (16L, 11L, 3w)			
	HotSpot	✗ (25L, 5L)	7.0 (7.0, 53.5)	1.19 (1.182, 7.03)	0.18279 ±0.000116	L	2.0 (2.0, 2.0)	0.14 (0.141, 0.143)	0.13699 ±0.00032
	JRuby+Truffle	✗ (99L, 123L, 5)	1082.0 —	2219.59 (2039, 304, 2516, 021)	2.05150 ±0.017738	L	69.0 (69, 70.0)	17.95 (17, 716, 18, 127)	0.20644 ±0.01568
	LuaJIT	✗ (23L, 4L, 2-, 1w)				—			0.25399 ±0.00471
	PyPy	✗ (27L, 3w)				—			1.85835 ±0.012893
	V8	✗ (15-, 9L, 6L)	1.5 (1.0, 794.0)	0.25 (0.000, 391.026)	0.49237 ±0.003198	= (25-, 5L)	1.0 (1.0, 361.6)	0.00 (0.000, 87.367)	0.24138 ±0.00389
	C	✗ (21-, 6L, 2L, 1w)				✗ (19-, 5L, 4w, 2L)			
	Graal	✗ (28L, 1w, 1L)				✗ (28L, 1w, 1L)			
fannkuchredux	HHVM	L	10.0 (10.0, 10.0)	52.66 (52,660,52,708)	1.35779 ±0.011948	L			
	HotSpot	L	390.0 (2.0, 390.0)	153.70 (0.407, 155.254)	0.36202 ±0.02767	L	26.0 (26, 2w, 2L)		
	JRuby+Truffle	✗ (99L, 1023.1)	1016.5 —	1039.04 (1014,290,1059,967)	1.08833 ±0.008580	L	1021.0 (1014,9,1027.0)	917.30 (901,708,946,483)	0.89509 ±0.027590
	LuaJIT	—				✗ (21w, 7L, 2L)			
	PyPy	✗ (15L, 13-, 2L)	2.0 (1.0, 28.9)	1.57 (0.000, 43.483)	1.55442 ±0.020549	L	2.0 (2.0, 3.0)	1.12 (1.114, 2.054)	0.96809 ±0.010950
	V8	= (19L, 11-)	2.0 (1.0, 25.9)	0.31 (0.000, 7.525)	0.30401 ±0.001514	= (29L, 1-)	4.0 (4.0, 16.0)	1.44 (1.43, 7.135)	0.47421 ±0.001218
	C	—		0.07048 ±0.000210		✗ (28L, 1-, 1L)	997.0 (2.0, 1001.0)	546.38 ((0.546, 547.717)	0.54547 ±0.026562
	Graal	✗ (29L, 1w)				✗ (29L, 1-)	15.0 (2.0, 21.0)	12.40 ((0.812, 17.804)	0.89293 ±0.008887
	HHVM	✗ (27L, 2w, 1L)				L	35.0 (34.0, 41.0)	139.18 ((136, 909, 147, 626)	1.40690 ±0.011333
	HotSpot	✗ (18L, 12L)	261.0 (6.0, 595.0)	30.73 (0.614, 70.021)	0.11744 ±0.001725	L	7.0 (7.0, 8.6)	1.90 ((1.901, 2.425)	0.31470 ±0.000929
fasta	JRuby+Truffle	—				L	1011.0 (1007.5, 1014.0)	893.38 ((888, 985, 896, 562)	0.836333 ±0.014043
	LuaJIT	~				—			0.22435 ±0.000620
	PyPy	~				L	75.0 (75.0, 75.0)	34.43 ((34, 429, 34, 457)	0.46489 ±0.000646
	V8	✗ (19L, 10L, 1w)				L	3.0 (3.0, 3.0)	0.55 ((0.554, 0.354))	0.24963 ±0.00039
	Richards	—				specnrm			

# Results

		Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)
C		~			
Graal		⌘ (27L, 3J)	32.0 (17.0, 193.8)	6.60 (3.729, 36.608)	0.18594 ±0.000316
HHVM		⌘ (24L, 4J, 2w)			
HotSpot	binary trees	⌘ (25L, 5J)	7.0 (7.0, 53.5)	1.19 (1.182, 9.703)	0.18279 ±0.000116
JRuby+Truffle		J	1082.0 (999.0, 1232.5)	2219.59 (2039.304, 2516.021)	2.05150 ±0.017737
LuaJIT		⌘ (23L, 4J, 2-, 1w)			
PyPy		⌘ (27J, 3w)			
V8		⌘ (15-, 9L, 6J)	1.5 (1.0, 794.0)	0.25 (0.000, 391.026)	0.49237 ±0.003198

# Results: Summary

Class.	Linux <sub>4790</sub>	Linux <sub>1240v5</sub>	OpenBSD <sub>4790</sub>
⟨VM, benchmark⟩ pairs			
—	8.7%	13.0%	6.7%
⊓	28.3%	23.9%	10.0%
⊓	6.5%	6.5%	0.0%
≈	4.3%	6.5%	0.0%
=	6.5%	6.5%	13.3%
≵	45.7%	43.5%	70.0%
Process executions			
—	26.4%	20.9%	34.0%
⊓	48.3%	51.5%	52.1%
⊓	16.7%	17.9%	11.1%
≈	8.7%	9.6%	2.8%

# Results: Summary

Class.	Linux <sub>4790</sub>	Linux <sub>1240v5</sub>	OpenBSD <sub>4790</sub>
⟨VM, benchmark⟩ pairs			
—	8.7%	13.0%	6.7%
⊓	28.3%	23.9%	10.0%
⊓	6.5%	6.5%	0.0%
≈	4.3%	6.5%	0.0%
=	6.5%	6.5%	13.3%
≵	45.7%	43.5%	70.0%
Process executions			
—	26.4%	20.9%	34.0%
⊓	48.3%	51.5%	52.1%
⊓	16.7%	17.9%	11.1%
≈	8.7%	9.6%	2.8%

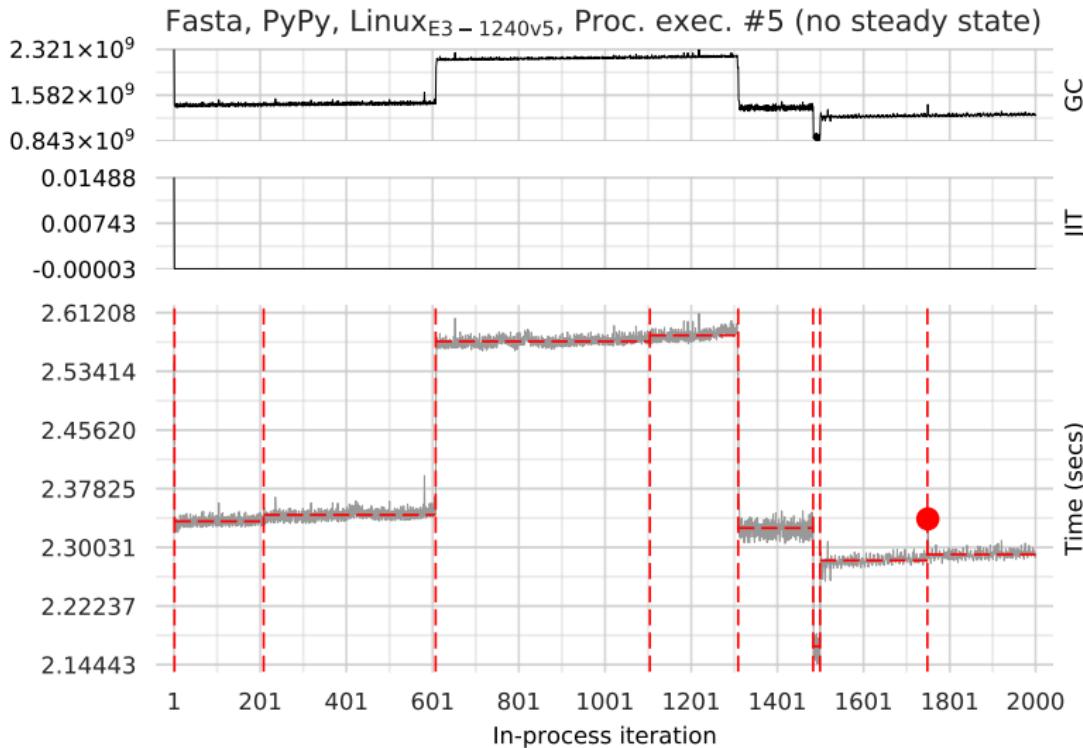
Best case no. good process executions: 86.1%

# Results: Summary

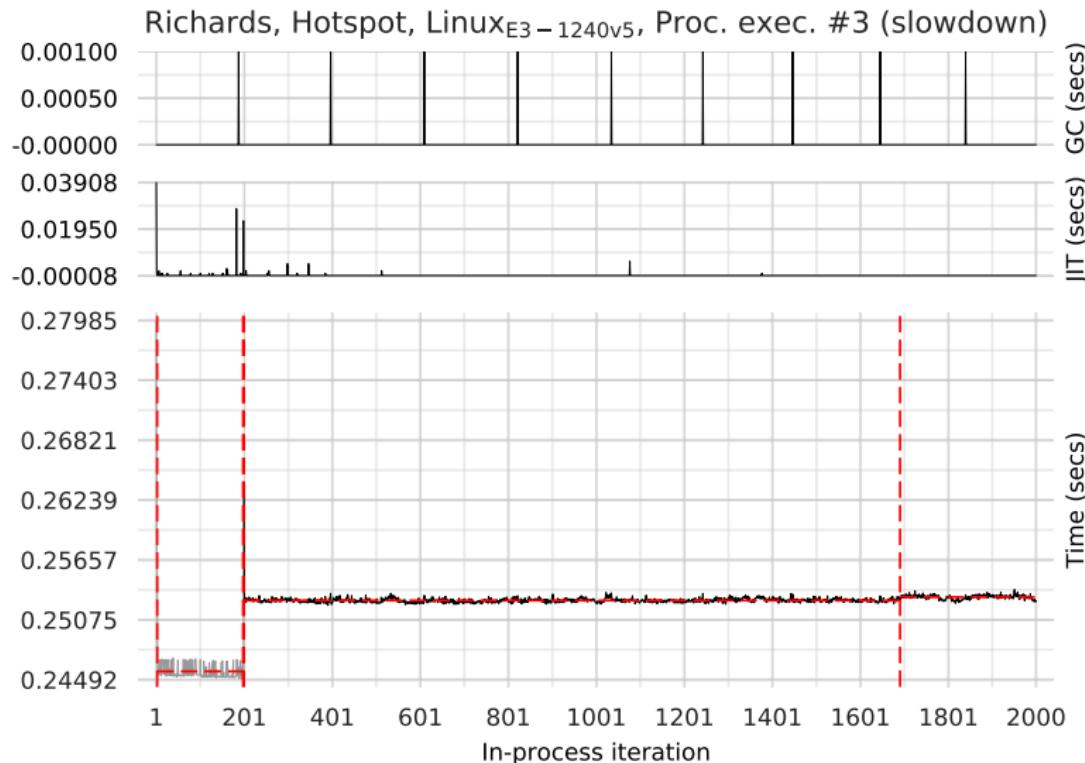
Class.	Linux <sub>4790</sub>	Linux <sub>1240v5</sub>	OpenBSD <sub>4790</sub>
⟨VM, benchmark⟩ pairs			
—	8.7%	13.0%	6.7%
⊓	28.3%	23.9%	10.0%
⊓	6.5%	6.5%	0.0%
≈	4.3%	6.5%	0.0%
=	6.5%	6.5%	13.3%
≈	45.7%	43.5%	70.0%
Process executions			
—	26.4%	20.9%	34.0%
⊓	48.3%	51.5%	52.1%
⊓	16.7%	17.9%	11.1%
≈	8.7%	9.6%	2.8%

Best case good ⟨VM, benchmark⟩ pairings: 43.5%

# Are the effects due to JIT and GC?



# Are the effects due to JIT and GC?



Are the effects due to JIT and GC?

However

In many cases, the JIT/GC can't explain oddness

# What Have We Learned?

- ▶ Benchmarks often don't warmup as we expect.
- ▶ Repeating a benchmark often gives a different warmup characteristic.
- ▶ Invalid benchmarking assumptions may have misled us!
  - ▶ Ineffectual or bad optimisations?

# What Can We Do?

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that peak performance may not occur.

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that peak performance may not occur.
- 3 Always report warmup time.

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that peak performance may not occur.
- 3 Always report warmup time.
- 4 Stop over-training on small benchmark suites.

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that peak performance may not occur.
- 3 Always report warmup time.
- 4 Stop over-training on small benchmark suites.
- 5 The more benchmarks, the better.

## What Can We Do?

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that peak performance may not occur.
- 3 Always report warmup time.
- 4 Stop over-training on small benchmark suites.
- 5 The more benchmarks, the better.
- 6 Focus on predictable performance.

# References

## **VM Warmup Blows Hot and Cold**

*E. Barrett, C. F. Bolz, R. Killick, V. Knight, S. Mount and L. Tratt.*

## **Rigorous Benchmarking in Reasonable Time**

*T. Kalibera and R. Jones*

## **Specialising Dynamic Techniques for Implementing the Ruby Programming Language**

*C. Seaton (Chapter 4)*

## **Quantifying performance changes with effect size confidence intervals**

*T. Kalibera and R. Jones*

# Thanks for Listening

