

Why Aren't More Users More Happy With Our VMs?



Laurence
Tratt

Warmup work in collaboration with:
Edd Barrett, Carl Friedrich Bolz, Rebecca Killick, and Sarah Mount



Software Development Team
2017-11-21

What to expect from this talk

What to expect from this talk

The gap

What to expect from this talk

What we tell users to expect

The gap

What to expect from this talk

What we tell users to expect

The gap

What users experience

What to expect from this talk

The gap

What to expect from this talk

The gap
is bigger than we think

A stroppy user? Or rightfully disappointed?

A stroppy user? Or rightfully disappointed?

*"You told me I'd get
a 10x speed-up,
but I only saw 1.2x"*

Background

Background

Dynamic ('JIT') compilation utilises information not known statically

Is this just about 'dynamic typing'?

Is this statically or dynamically typed?

```
fn f(a: Option<i64>) {  
    match a {  
        Some(i) => ...,  
        None     => ...  
    };  
}
```

Is this just about 'dynamic typing'?

Is this statically or dynamically typed? Both!

```
fn f(a: Option<i64>) {  
    match a {  
        Some(i) => ...,  
        None     => ...  
    };  
}
```

VM claims

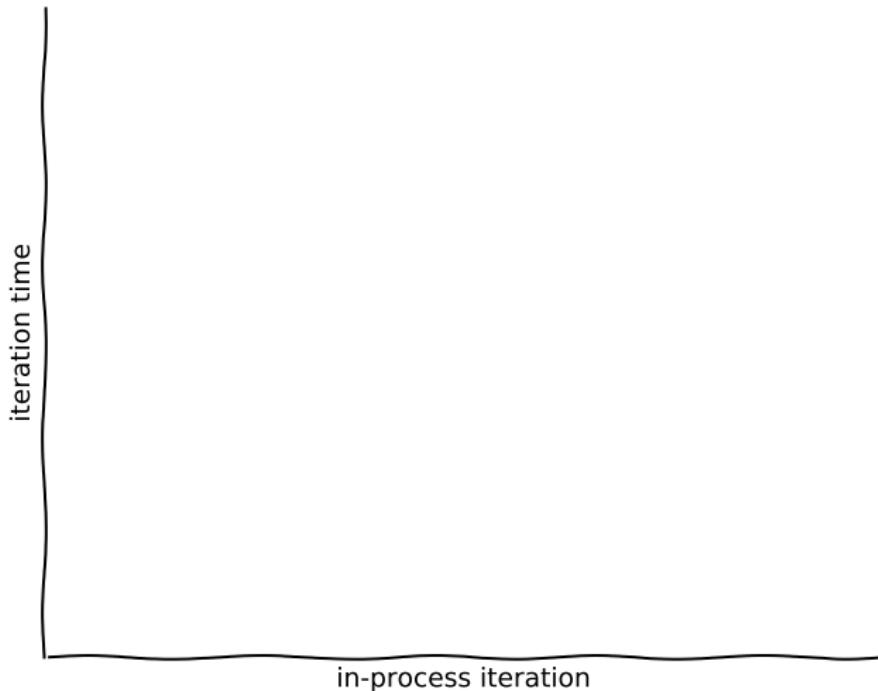
The best VMs are close in performance to,
and sometimes faster than,
`gcc -O2`

VM claims

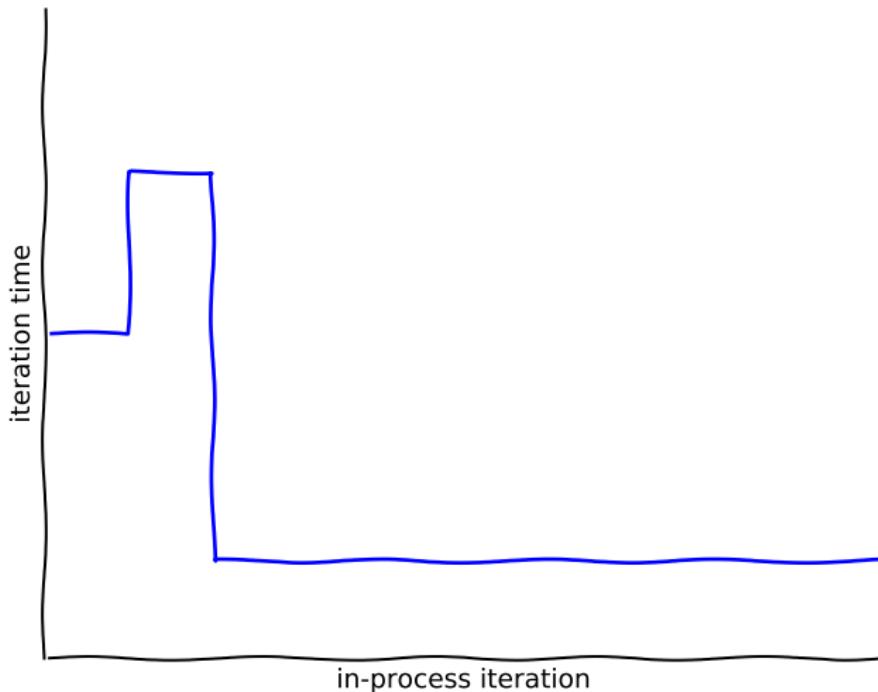
The best VMs are close in performance to,
and sometimes faster than,
`gcc -O2`

What's being measured?

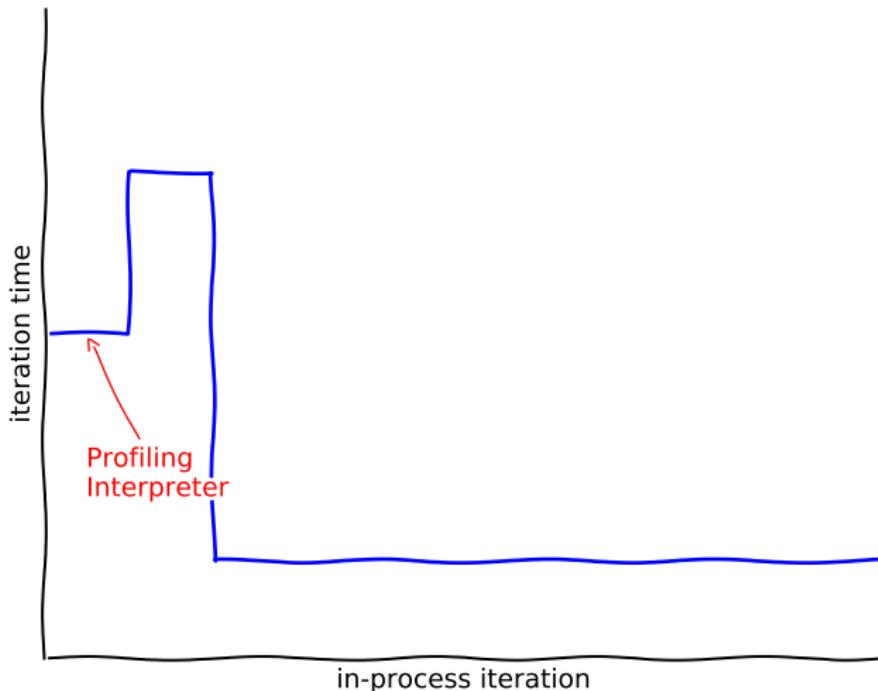
What our claims pertain to



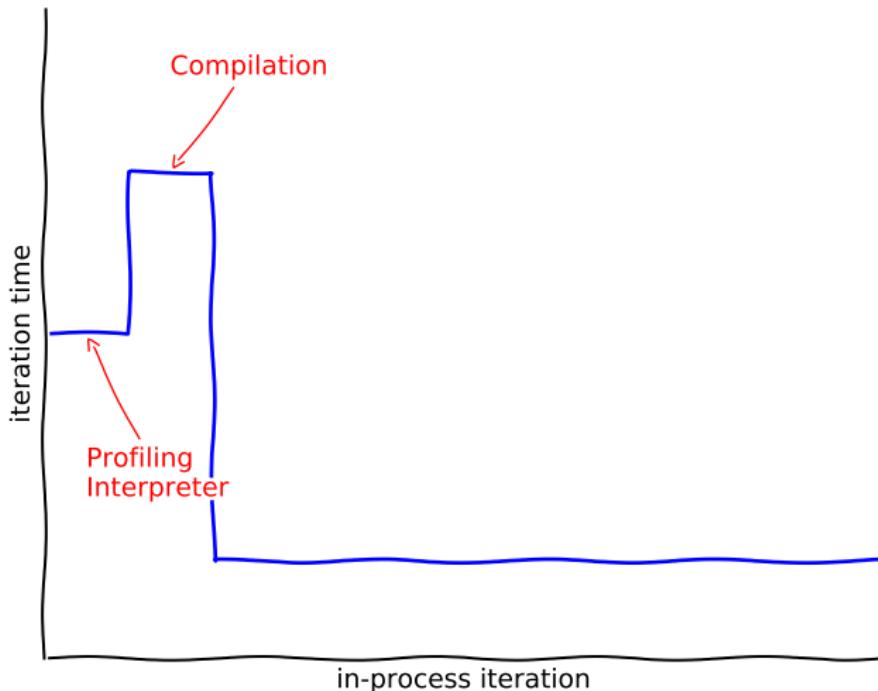
What our claims pertain to



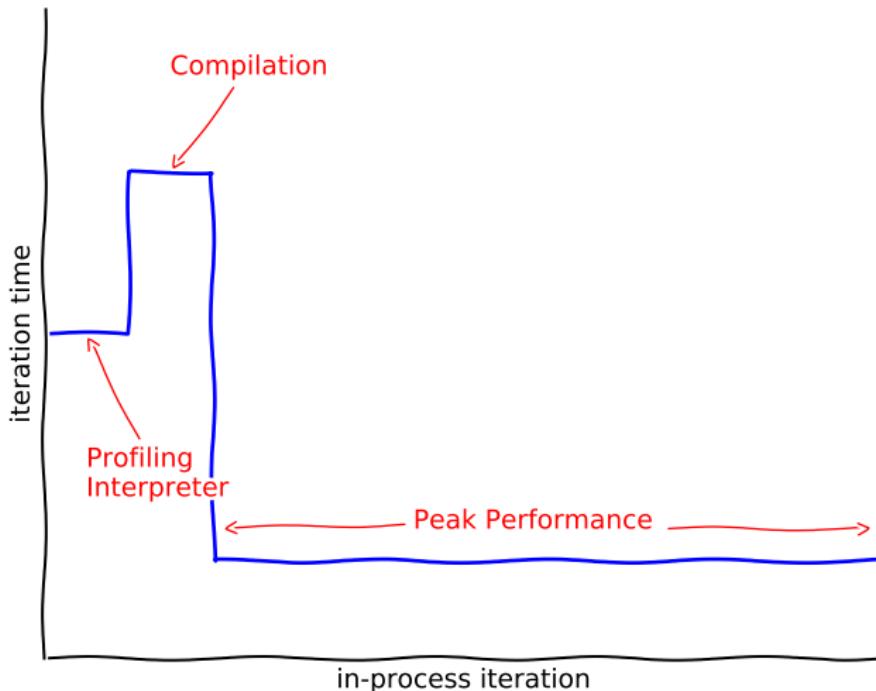
What our claims pertain to



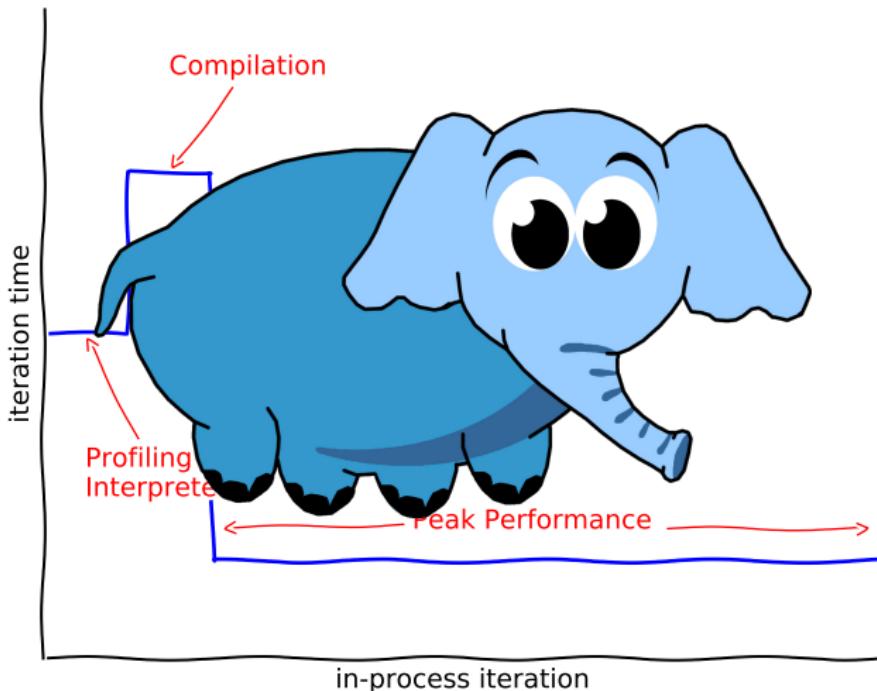
What our claims pertain to



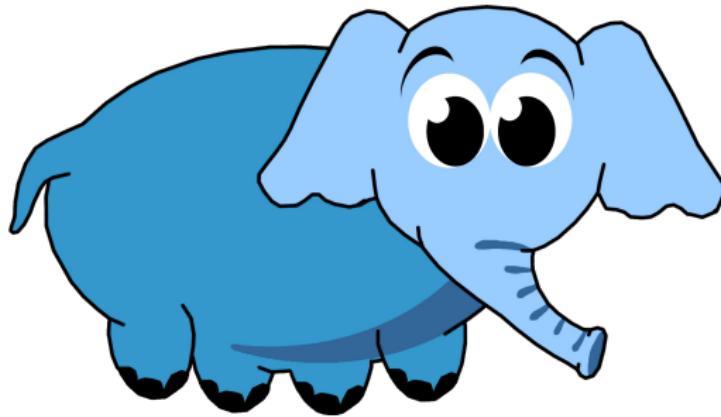
What our claims pertain to



What our claims pertain to

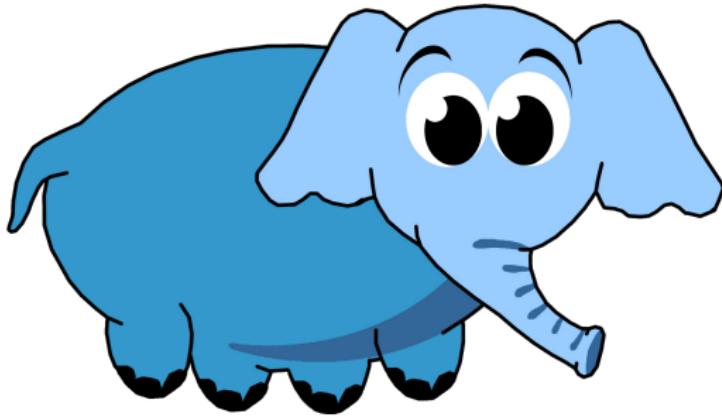


What our claims pertain to



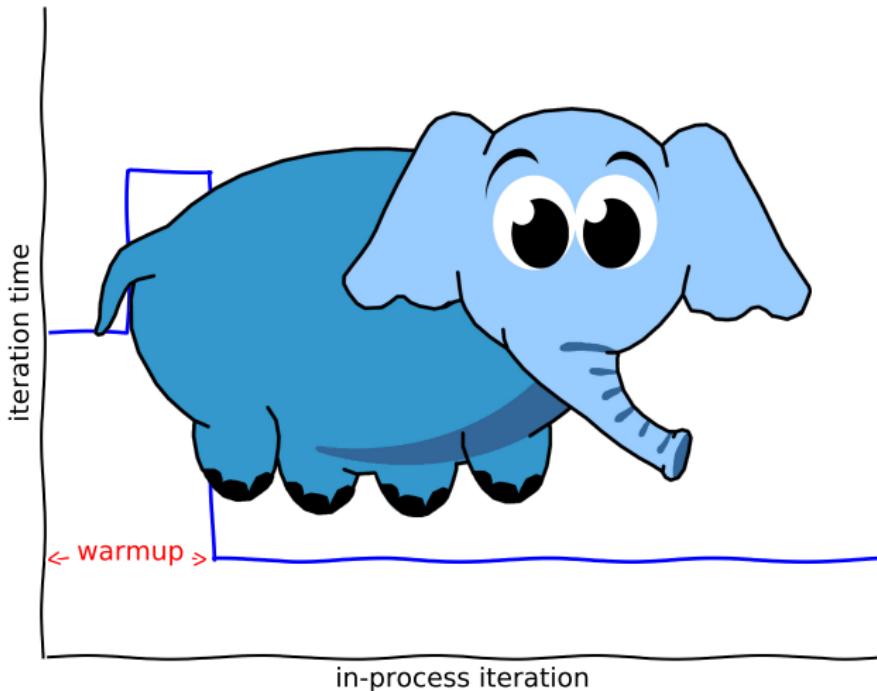
This is Barry

What our claims pertain to



This is Barry: the benchmarking elephant in the room

What our claims pertain to



Warmup

Users *always* perceive warmup

Warmup

Users *always* perceive warmup

Maybe we should know how long it is?

The Warmup Experiment

Measure warmup of modern language implementations

The Warmup Experiment

Measure warmup of modern language implementations

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

Method 1: Which benchmarks?

The language benchmark games are perfect for us
(unusually)

Method 1: Which benchmarks?

The language benchmark games are perfect for us
(unusually)

We removed any CFG non-determinism

Method 1: Which benchmarks?

The language benchmark games are perfect for us
(unusually)

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

- Linux₄₇₉₀, Debian 8, 24GiB RAM
- Linux_{E3-1240v5}, Debian 8, 32GiB RAM
- OpenBSD₄₇₉₀, OpenBSD 6.0, 32GiB RAM

Method 4: Machines

- Linux₄₇₉₀, Debian 8, 24GiB RAM
 - Linux_{E3-1240v5}, Debian 8, 32GiB RAM
 - OpenBSD₄₇₉₀, OpenBSD 6.0, 32GiB RAM
-
- Turbo boost and hyper-threading disabled
 - Network card turned off.
 - Daemons disabled (cron, smtpd)

Method 5: Krun

Benchmark runner: tries to control as many confounding variables as possible

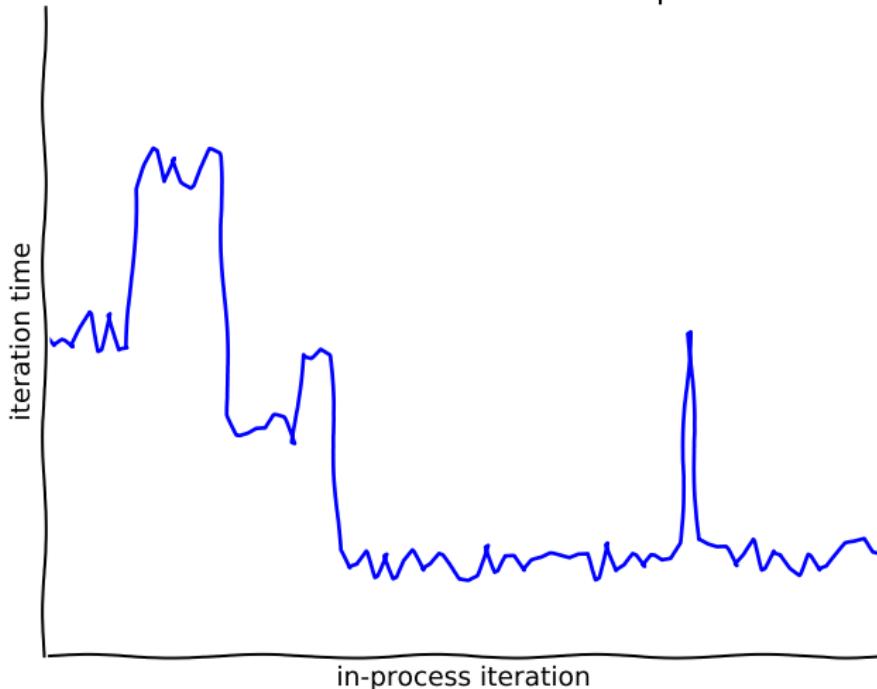
Method 5: Krun

Benchmark runner: tries to control as many confounding variables as possible e.g.:

- 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, ...)

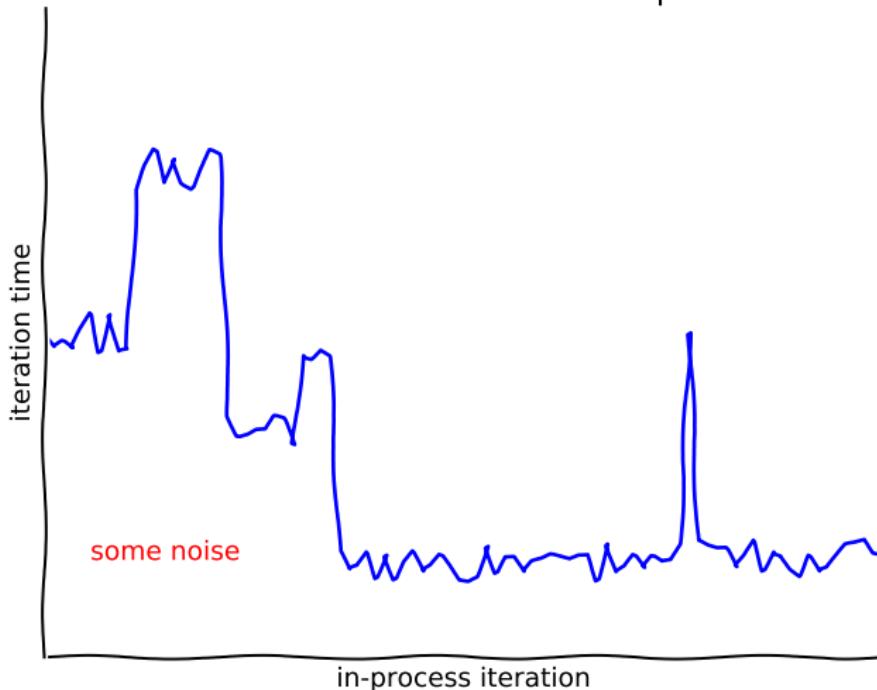
Method 6: Expectations

More Realistic VM Warmup

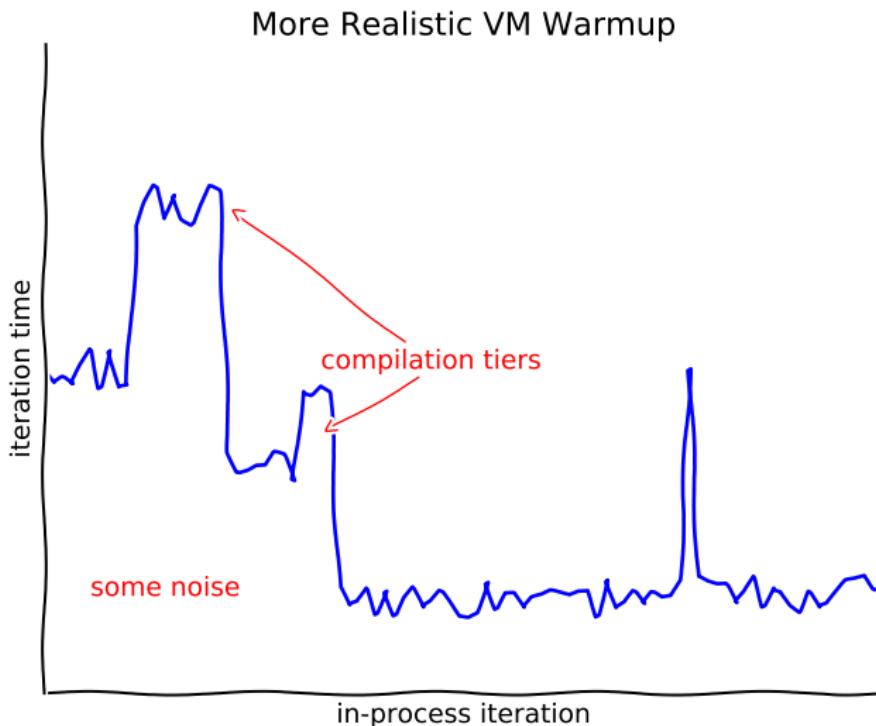


Method 6: Expectations

More Realistic VM Warmup

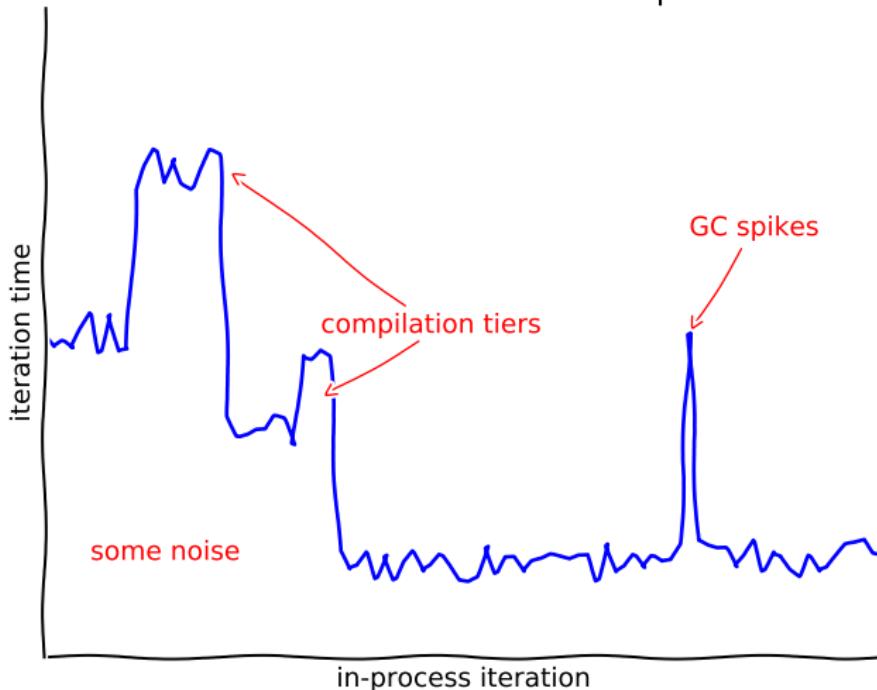


Method 6: Expectations

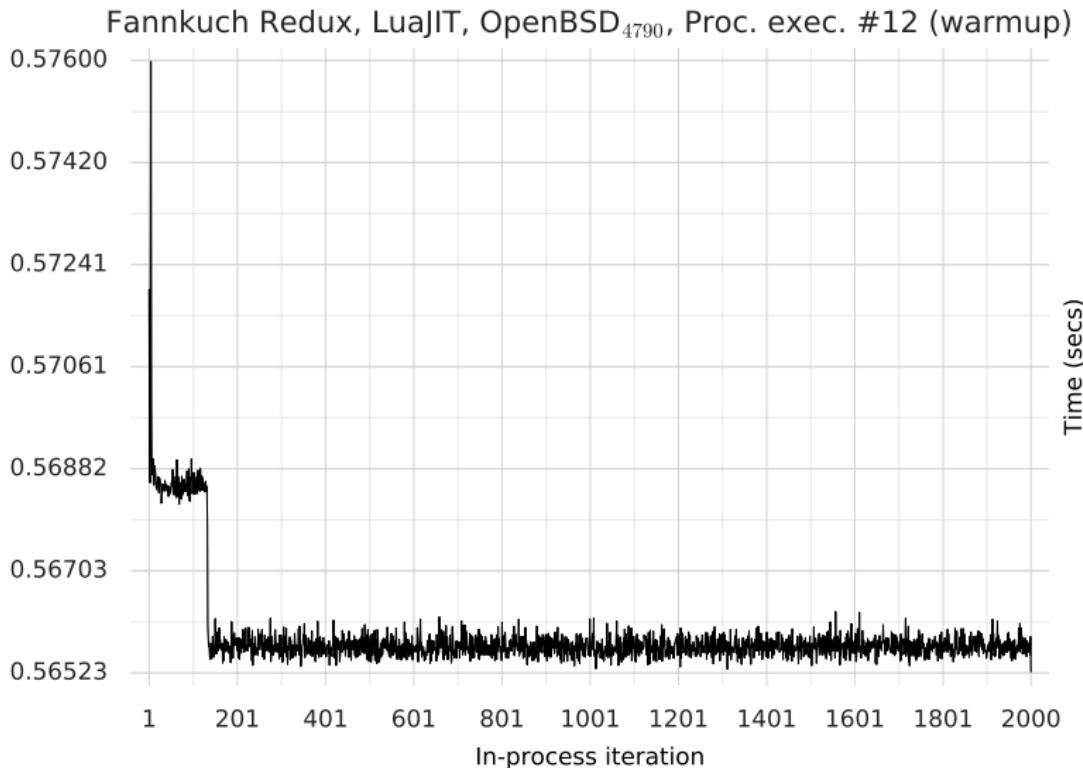


Method 6: Expectations

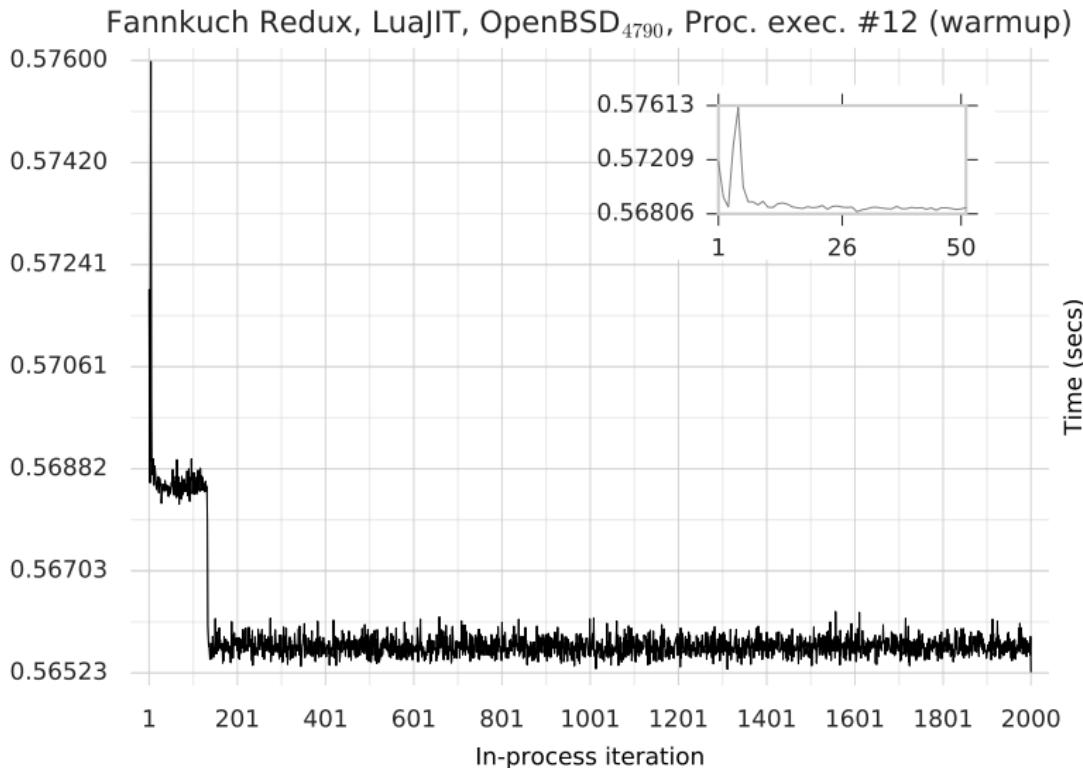
More Realistic VM Warmup



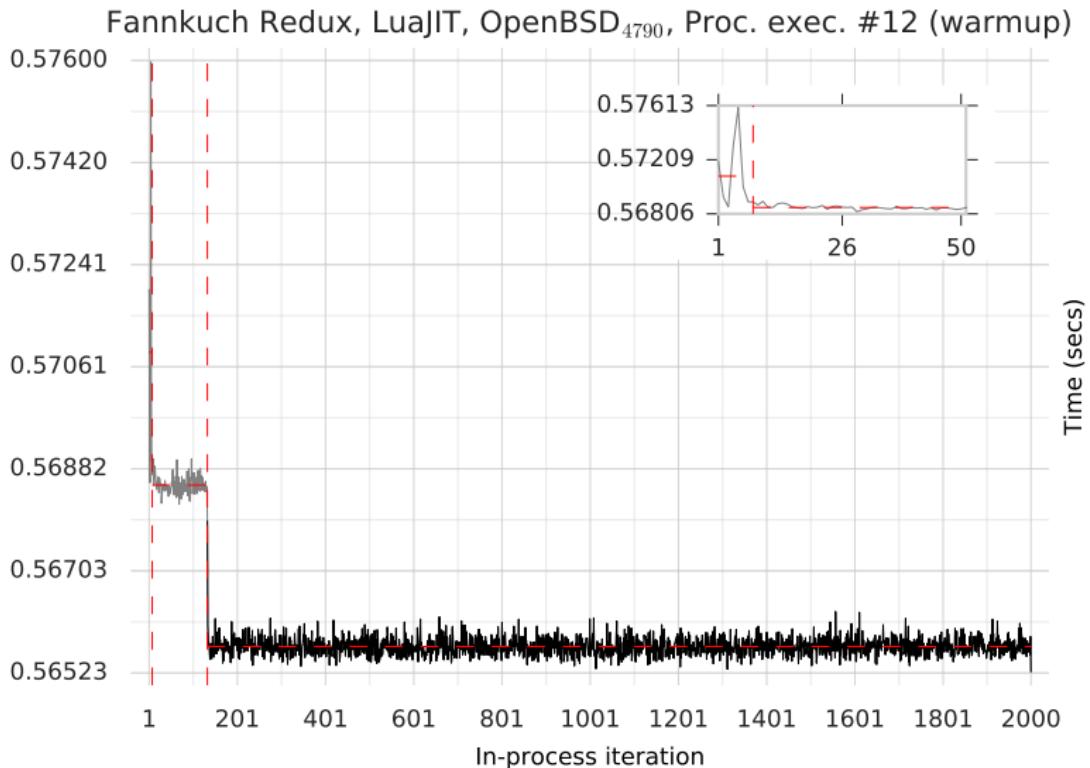
Warmup & flat (1)



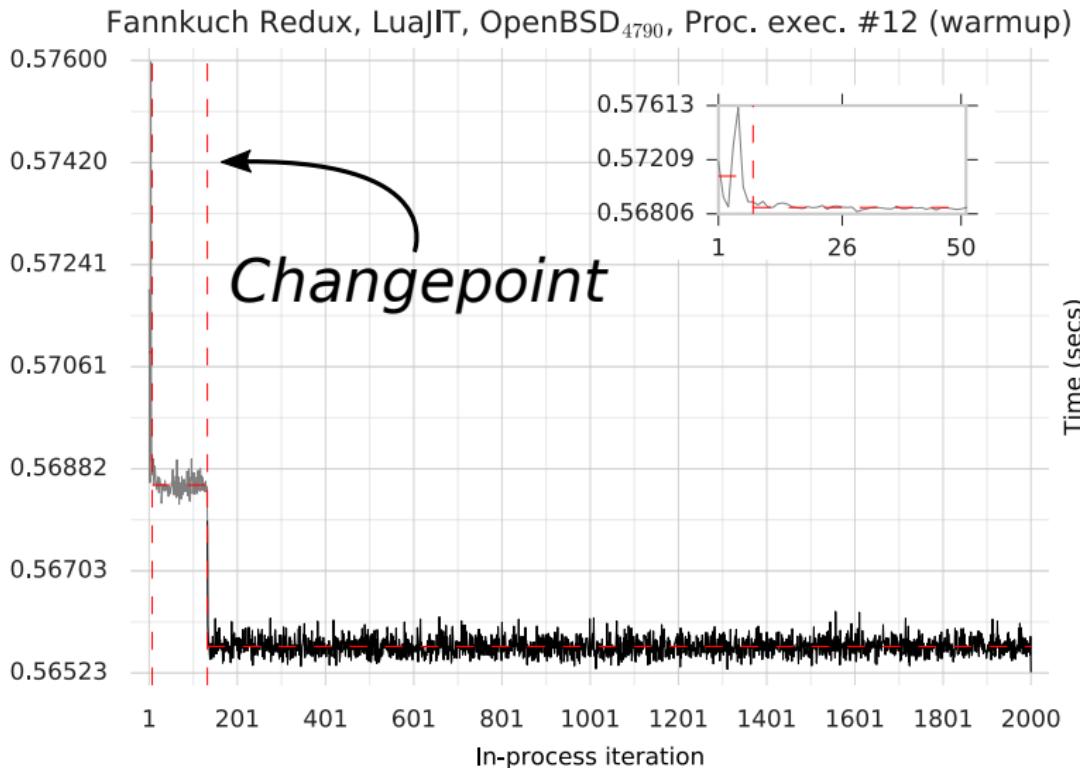
Warmup & flat (1)



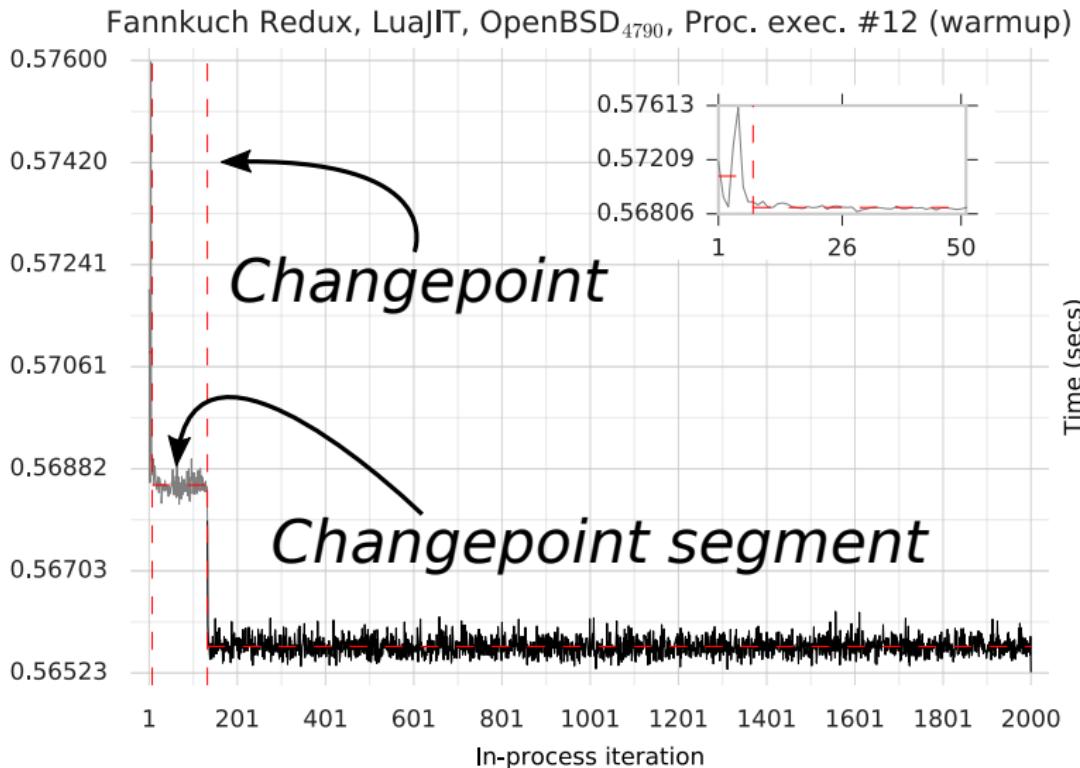
Warmup & flat (1)



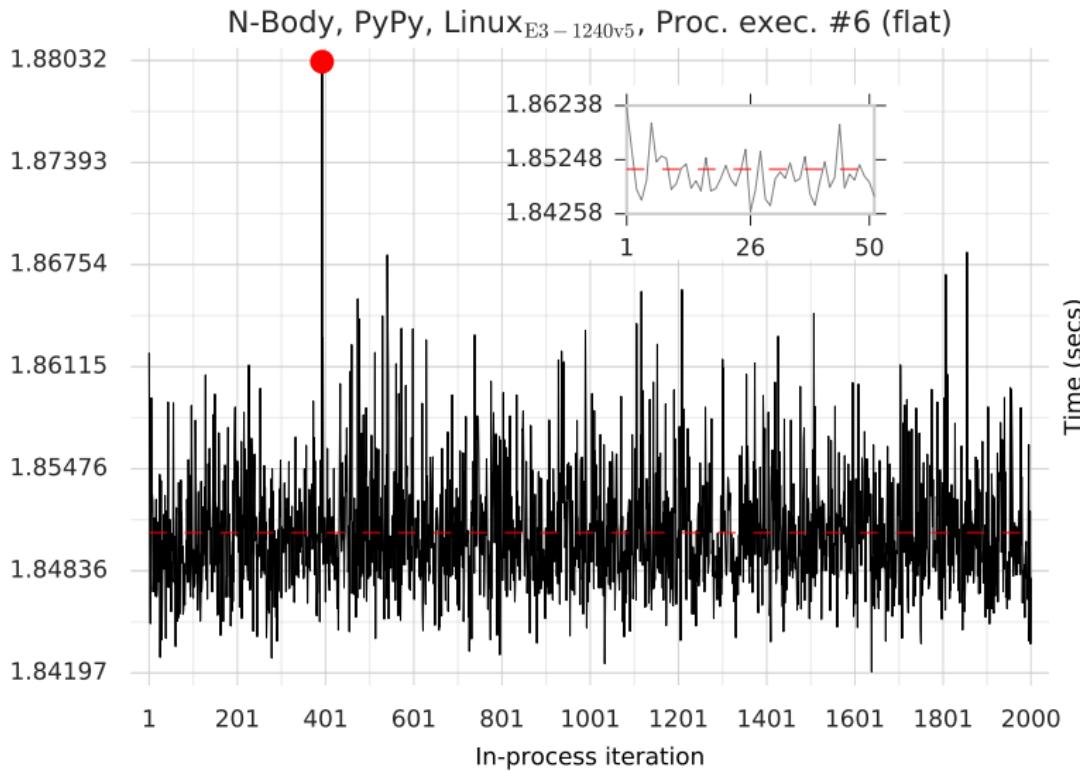
Warmup & flat (1)



Warmup & flat (1)



Warmup & flat (1)



Method 7: Classification

Classification algorithm (steps in order):

All segs are equivalent: *flat*

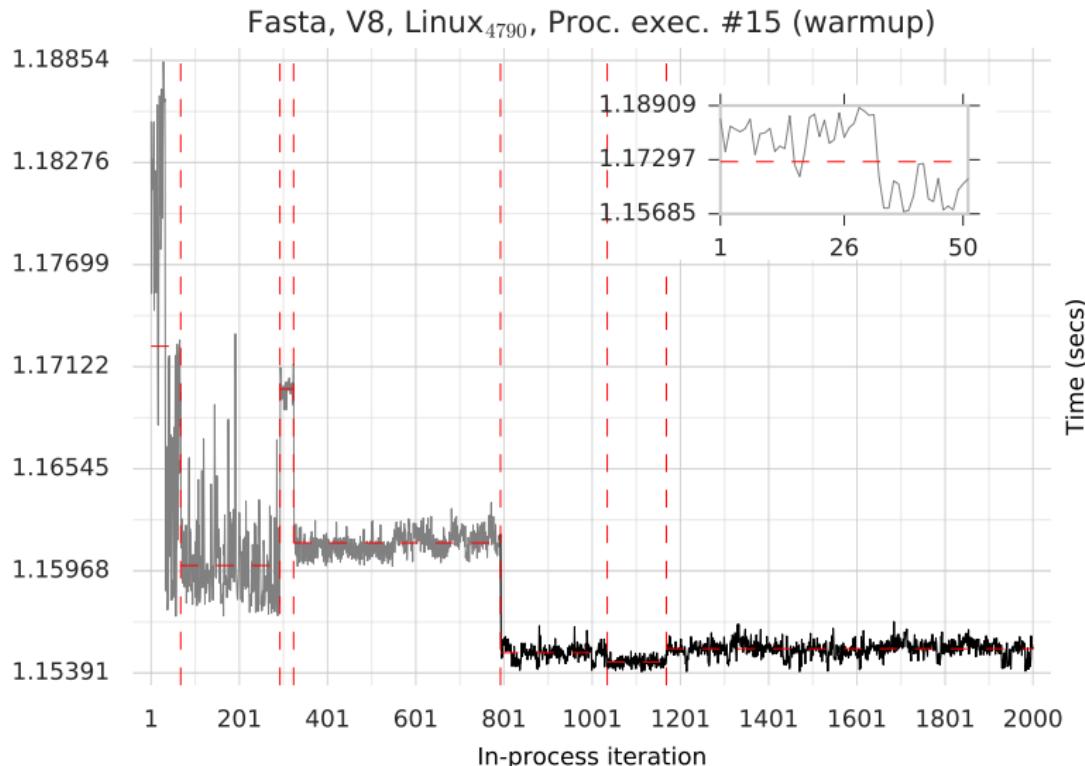
Method 7: Classification

Classification algorithm (steps in order):

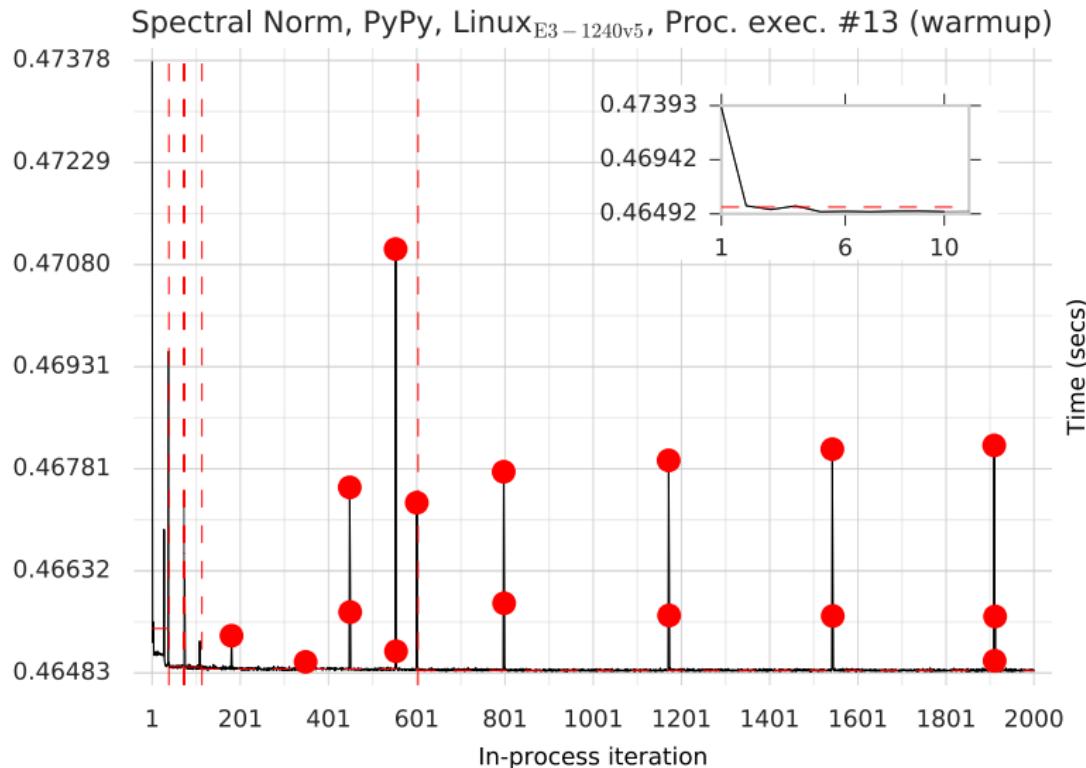
All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

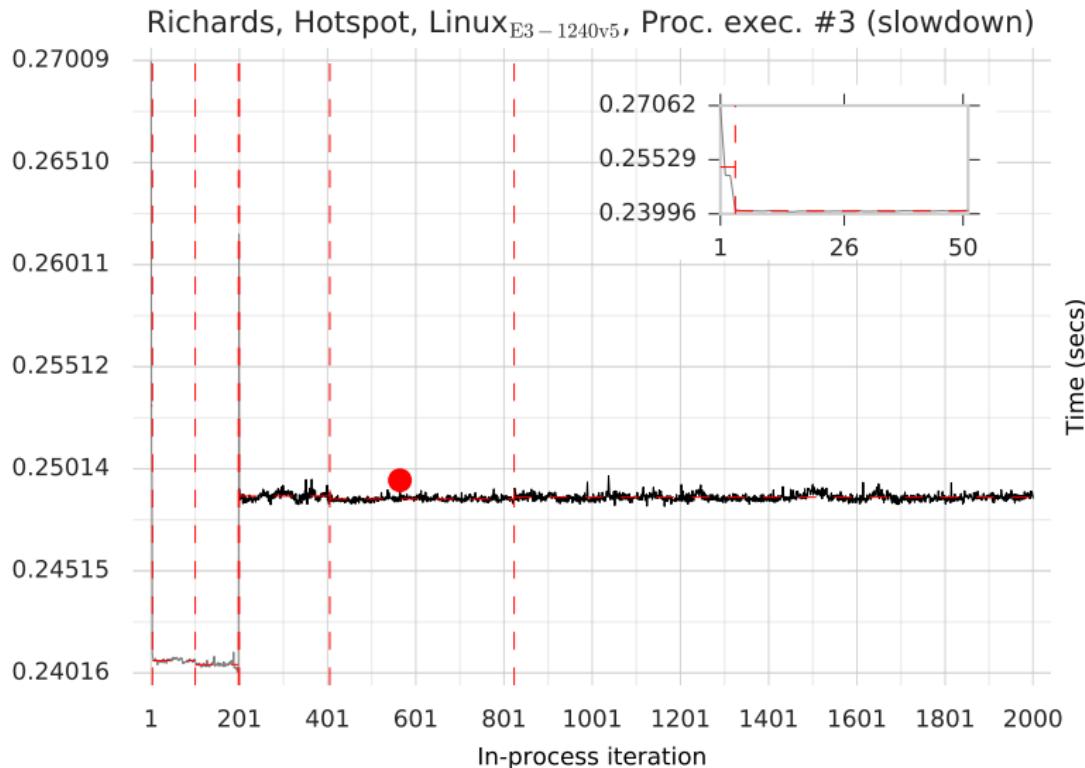
Warmup & flat (2)



Warmup & flat (2)



Slowdown (1)



Method 7: Classification

Classification algorithm (steps in order):

All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

Method 7: Classification

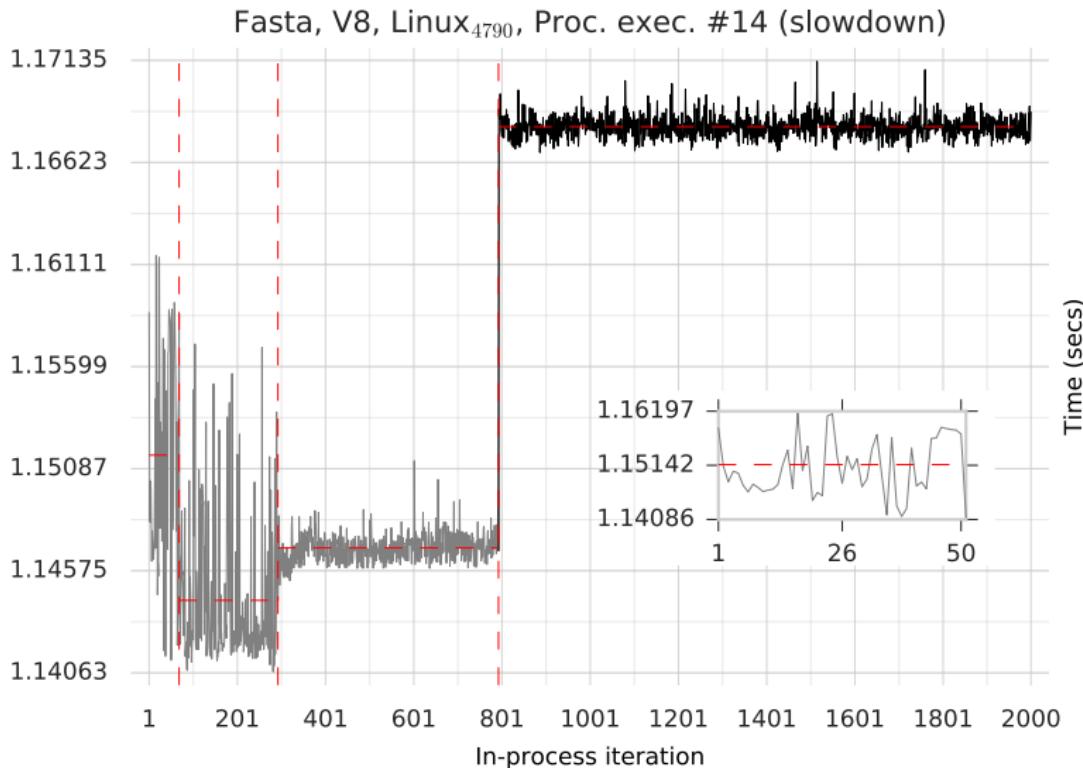
Classification algorithm (steps in order):

All segs are equivalent: *flat*

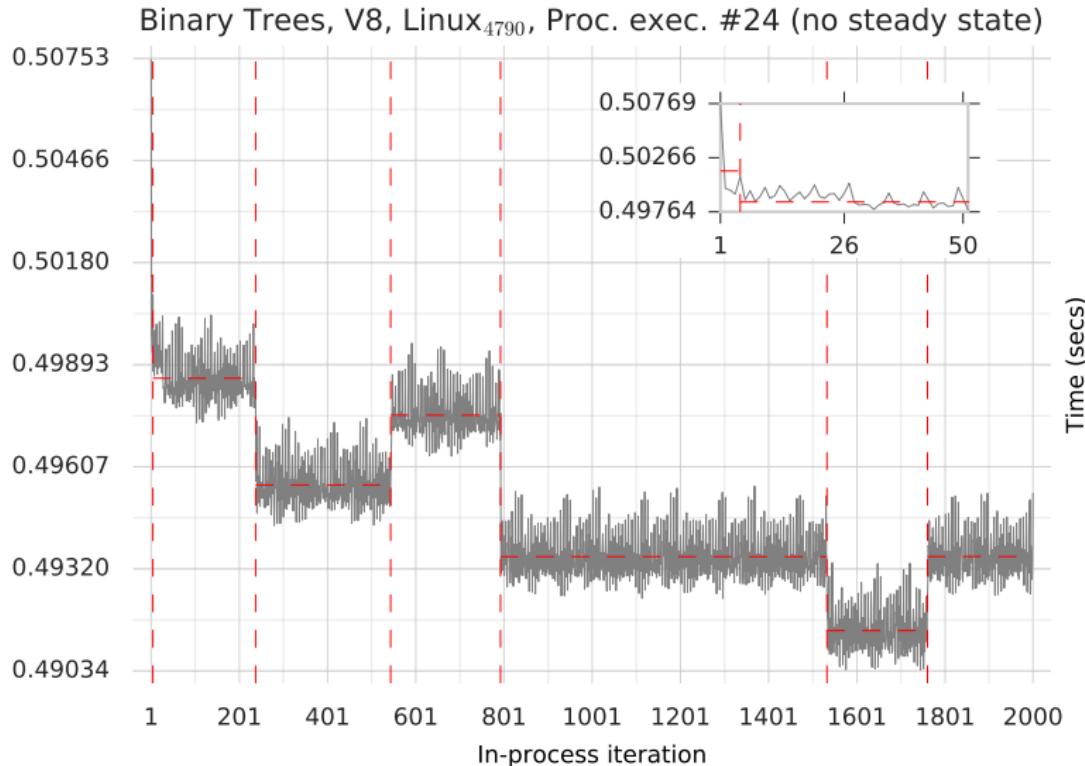
Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

Slowdown (2)



No steady state (1)



Classification (3)

Classification algorithm (steps in order):

All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

Classification (3)

Classification algorithm (steps in order):

All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

Else: *no steady state*

Classification (3)

Classification algorithm, in order:

All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

Else: *no steady state*

Good

Classification (3)

Classification algorithm, in order:

All segs are equivalent: *flat*

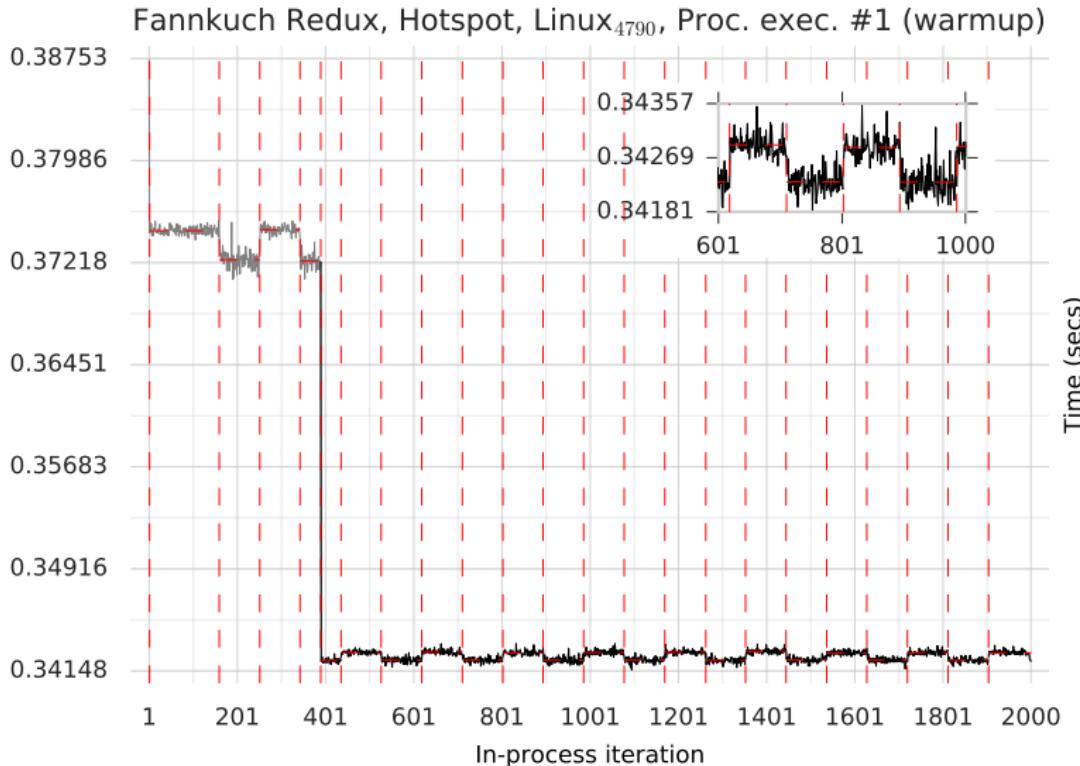
Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

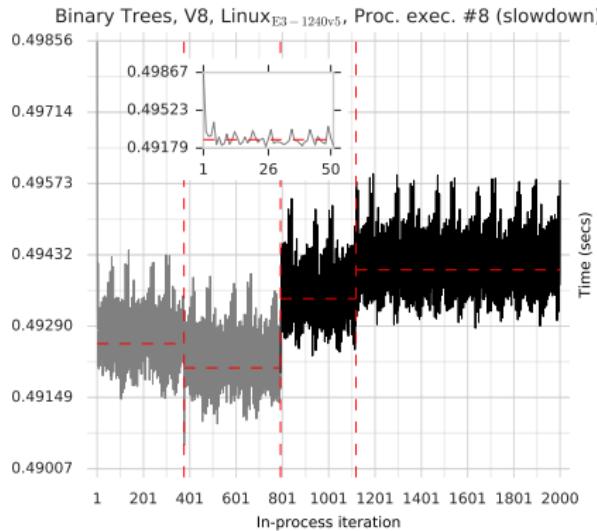
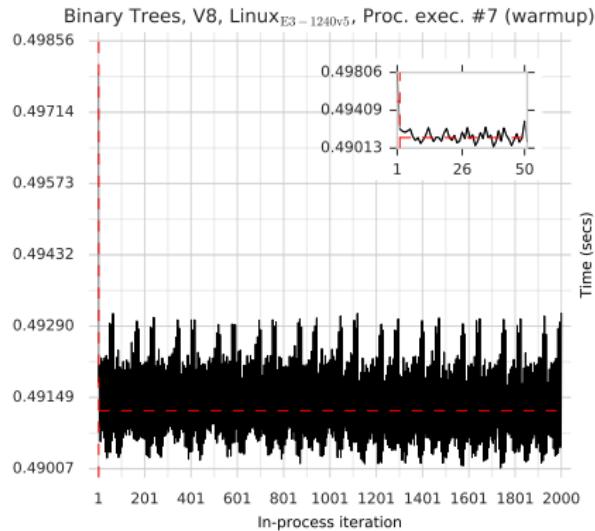
Else: *no steady state*

Bad

Warmup or no steady state?

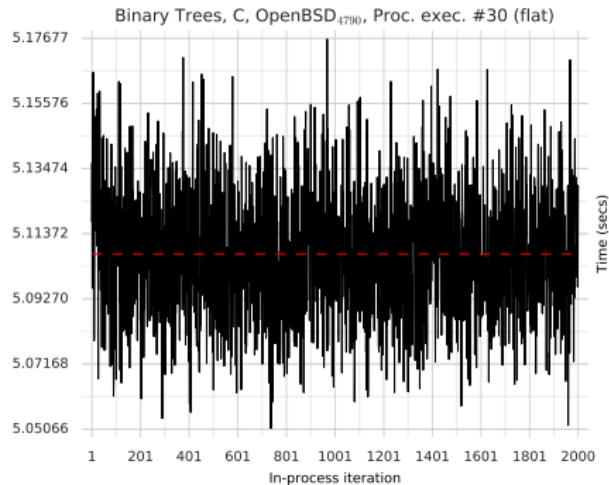
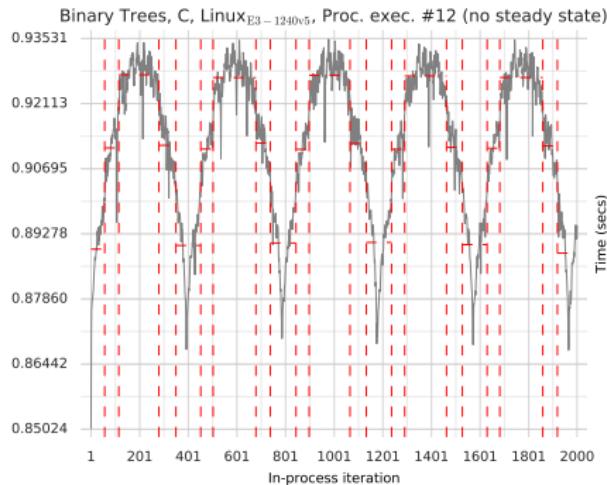


Inconsistent Process-executions



(Same machine)

Inconsistent Process-executions



(Different machines. Bouncing ball Linux-specific)

Individual benchmark stats

	Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)		Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)
C	**	32.0	6.60	0.18594 ±0.00935		-			0.40555 ±0.00510
Graal	✗ (27L, 3f)	(17.0, 193.8)		(3.729, 36.698)		L	8.0	(7.5, 8.5)	1.22 ±0.00045
HHVM	✗ (24L, 4f, 2w)					✗ (16L, 11f, 3w)			0.13334 ±0.00045
HotSpot	✗ (25L, 5f)	7.0	1.19	0.18279 ±0.00116		L	2.0	(2.0, 2.0)	0.13699 ±0.00032
JRuby+Truffle	L	1082.0	2219.59	2.05150 ±0.01738		L	69.0	(69.0, 70.0)	17.95 ±0.00568
LuaJIT	✗ (23L, 4f, 2w, 1w)	(999.0, 1232.5)	(2039.304, 2156.021)			n-body	-	(17.716, 18.127)	0.20644 ±0.00568
PyPy	✗ (27L, 3w)						-		0.25399 ±0.00447
V8	✗ (15L, 9L, 6f)	1.5	0.25	0.49237 ±0.003198		= (25L, 5L)	1.0	(1.0, 361.6)	0.00 ±0.00389
C	✗ (21L, 6L, 2f, 1w)					✗ (19L, 5f, 4w, 2L)			
Graal	✗ (28L, 1w, 1f)					✗ (28L, 1w, 1f)			
HHVM	L	10.0 (10.0, 10.0)	52.66 (52.660, 52.708)	1.35779 ±0.01948					
HotSpot	L	390.0 (2.0, 390.0)	153.70 (0.407, 155.254)	0.36202 ±0.02767					
JRuby+Truffle	L	1016.5 (999.0, 1923.1)	1039.04 (1014.290, 1059.967)	1.08833 ±0.038580					
LuaJIT	-			0.56285 ±0.00097					
PyPy	✗ (15L, 13L, 2f)	2.0 (1.0, 28.9)	1.57 (0.000, 43.483)	1.55442 ±0.026549					
V8	= (19L, 11L)	2.0 (1.0, 25.5)	0.31 (0.008, 7.525)	0.30401 ±0.000154					
C	-			0.07048 ±0.00020		Richards	1021.0 (1014.9, 1027.0)	917.30 (901.708, 946.683)	0.89509 ±0.02790
Graal	✗ (29L, 1w)								
HHVM	✗ (27L, 2w, 1f)								
HotSpot	✗ (18L, 12f)	261.0 (6.0, 595.0)	30.73 (0.614, 70.021)	0.11744 ±0.001723					
JRuby+Truffle	fasta								
LuaJIT	**								
PyPy	**								
V8	✗ (19L, 10f, 1w)								
						spectralnorm			
C	-					✗ (28L, 1w, 1f)	997.0 (2.0, 1001.6)	546.38 (0.546, 547.317)	0.54547 ±0.00562
Graal	✗ (29L, 1w)					✗ (29L, 1-)	15.0 (2.0, 21.0)	12.40 (0.812, 17.404)	0.89293 ±0.00087
HHVM	✗ (27L, 2w, 1f)						35.0 (34.0, 41.0)	139.18 (136.909, 147.826)	1.40690 ±0.01113
HotSpot	✗ (18L, 12f)	261.0 (6.0, 595.0)	30.73 (0.614, 70.021)	0.11744 ±0.001723			7.0 (7.0, 8.4)	(1.901, 2.025)	0.31470 ±0.00092
JRuby+Truffle							1011.0 (1007.5, 1014.0)	893.38 (888.985, 896.562)	0.83633 ±0.01403
LuaJIT	-								0.22435 ±0.00062
PyPy	-								0.46489 ±0.00096
V8	✗ (19L, 10f, 1w)								0.24963 ±0.00039

Individual benchmark stats

		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	binary trees	⌘ (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

Overall benchmark stats

Class.	Linux ₄₇₉₀	Linux _{1240v5}	OpenBSD ₄₇₉₀ [†]
⟨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%

Overall benchmark stats

Class.	Linux ₄₇₉₀	Linux _{1240v5}	OpenBSD ₄₇₉₀ [†]
⟨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%

Summary

Classical warmup occurs for only:

Summary

Classical warmup occurs for only:

72.4%-74.7% of process executions

Summary

Classical warmup occurs for only:

72.4%-74.7% of process executions

43.4%-43.5% of (VM, benchmark) pairs

Summary

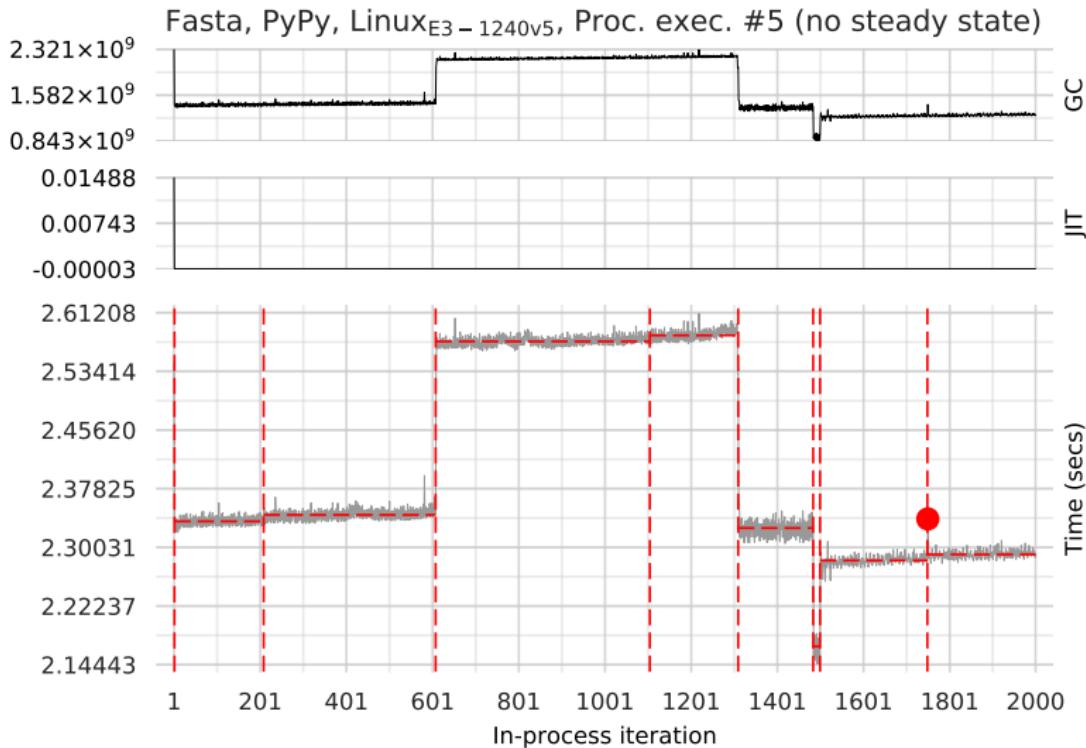
Classical warmup occurs for only:

72.4%-74.7% of process executions

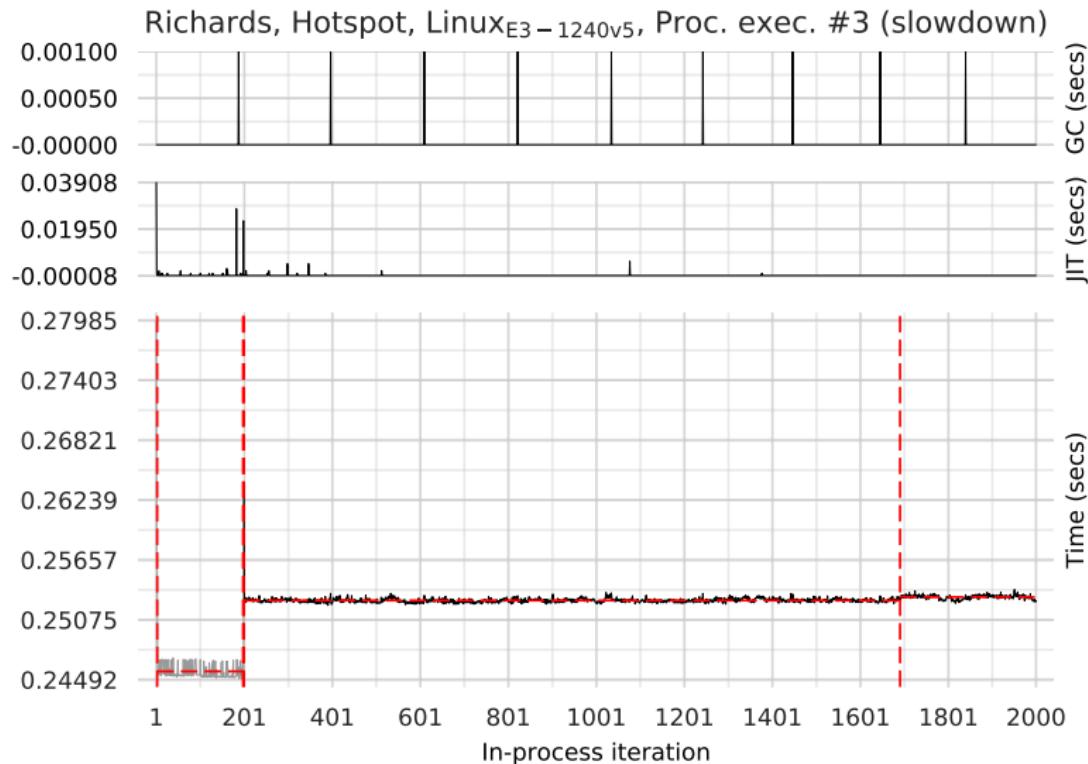
43.4%-43.5% of (VM, benchmark) pairs

0% of benchmarks for (VM, benchmark, machine) triples

Are odd effects correlated with compilation and GC?



Are odd effects correlated with compilation and GC?



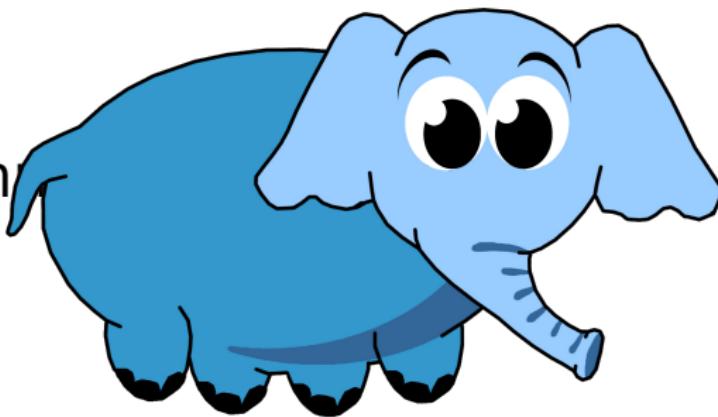
Benchmark suites

Benchmark suites

Benchmarks guide our optimisations

Benchmark suites

Benchmarks



Benchmark suites

Benchmarks guide our optimisations

Are they complete guides?

A war story

A war story

Symptom: poor performance of a Pyston
benchmark on PyPy

A war story

Symptom: poor performance of a Pyston
benchmark on PyPy

Cause: RPython traces recursion

A war story

Symptom: poor performance of a Pyston
benchmark on PyPy

Cause: RPython traces recursion

Fix: Check for recursion before tracing

A war story: the basis of a fix

```
diff --git a/rpython/jit/metainterp/pyjitpl.py b/rpython/jit/metainterp/pyjitpl.py
--- a/rpython/jit/metainterp/pyjitpl.py
+++ b/rpython/jit/metainterp/pyjitpl.py
@@ -951,9 +951,31 @@ 
     if warmrunnerstate.inlining:
         if warmrunnerstate.can_inline_callable(greenboxes):
+            # We've found a potentially inlinable function; now we need to
+            # see if it's already on the stack. In other words: are we about
+            # to enter recursion? If so, we don't want to inline the
+            # recursion, which would be equivalent to unrolling a while
+            # loop.
             portal_code = targetjitdriver_sd.mainjitcode
-            return self.metainterp.perform_call(portal_code, allboxes,
-                                                greenkey=greenboxes)
+            inline = True
+            if self.metainterp.is_main_jitcode(portal_code):
+                for gk, _ in self.metainterp.portal_trace_positions:
+                    if gk is None:
+                        continue
+                    assert len(gk) == len(greenboxes)
+                    i = 0
+                    for i in range(len(gk)):
+                        if not gk[i].same_constant(greenboxes[i]):
+                            break
+                    else:
+                        # The greenkey of a trace position on the stack
+                        # matches what we have, which means we're definitely
+                        # about to recurse.
+                        inline = False
+                        break
+            if inline:
+                return self.metainterp.perform_call(portal_code, allboxes,
+                                                    greenkey=greenboxes)
```

A war story: mixed fortunes

Success: slow benchmark now 13.5x faster

A war story: mixed fortunes

Success: slow benchmark now 13.5x faster

Failure: some PyPy benchmarks slow down

A war story: mixed fortunes

Success: slow benchmark now 13.5x faster

Failure: some PyPy benchmarks slow down

Solution: allow *some* tracing into recursion

A war story: data

#unrollings	1	2	3	5	7	10	
hexiom2	1.3	1.4	1.1	1.0	1.0	1.0	
raytrace-simple	3.3	3.1	2.8	1.4	1.0	1.0	
spectral-norm	3.3	1.0	1.0	1.0	1.0	1.0	
sympy_str	1.5	1.0	1.0	1.0	1.0	1.0	
telco	4	2.5	2.0	1.0	1.0	1.0	
polymorphism	0.07	0.07	0.07	0.07	0.08	0.09	

<http://marc.info/?l=pypy-dev&m=141587744128967&w=2>

A war story: conclusion

The benchmark suite said 7 levels, so that's what I suggested

A war story: conclusion

The benchmark suite said 7 levels, so that's what I suggested

Even though I doubted it was the right global value

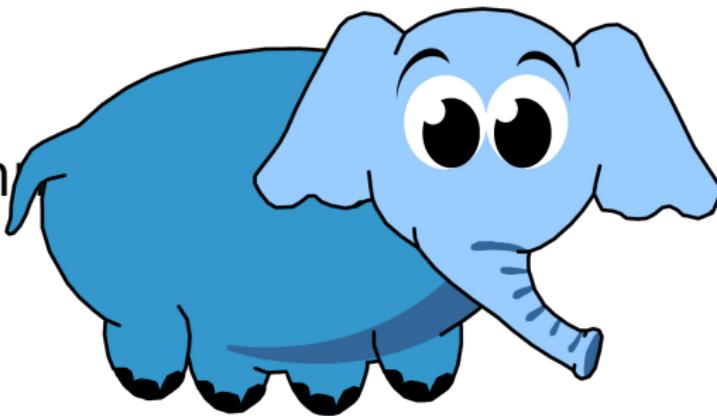
Benchmark suites (2)

Benchmark suites (2)

Benchmarks guide our optimisations

Benchmark suites (2)

Benchmarks



Benchmarks guide our optimisations

Are they correct guides?

17 JavaScript benchmarks from V8

17 JavaScript benchmarks from V8

Let's make each benchmark run for 2000 iterations

Octane: pdf.js explodes

```
$ d8 run.js
Richards
DeltaBlue
Encrypt
Decrypt
RayTrace
Earley
Boyer
RegExp
Splay
NavierStokes
PdfJS
```

```
<--- Last few GCs --->
```

```
14907865 ms: Mark-sweep 1093.9 (1434.4) -> 1093.4 (1434.4) MB, 274.8 / 0.0 ms [allocation failure] [GC in old space
14908140 ms: Mark-sweep 1093.4 (1434.4) -> 1093.3 (1434.4) MB, 274.4 / 0.0 ms [allocation failure] [GC in old space
14908421 ms: Mark-sweep 1093.3 (1434.4) -> 1100.5 (1418.4) MB, 280.9 / 0.0 ms [last resort gc].
14908703 ms: Mark-sweep 1100.5 (1418.4) -> 1107.8 (1418.4) MB, 282.1 / 0.0 ms [last resort gc].
```

```
<--- JS stacktrace --->
```

```
===== JS stack trace =====
```

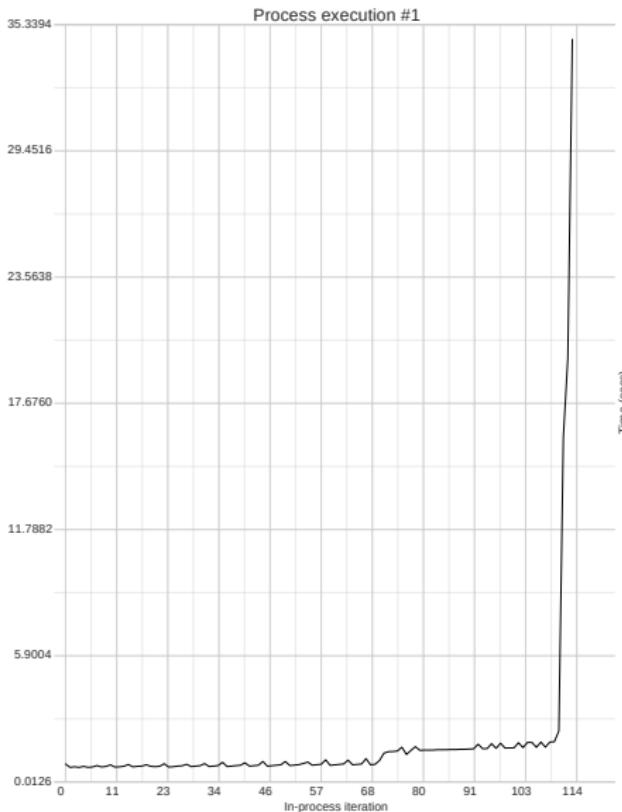
```
Security context: 0x20d333ad3ba9 <JS Object>
```

```
 2: extractFontProgram(aka Type1Parser_extractFontProgram) [pdfjs.js:17004] [pc=0x3a13b275421b] (this=0x3de358283
 3: new Type1Font [pdfjs.js:17216] [pc=0x3a13b2752078] (this=0x4603fbdae9 <a Type1Font with map 0x1f822134f7e1>,
```

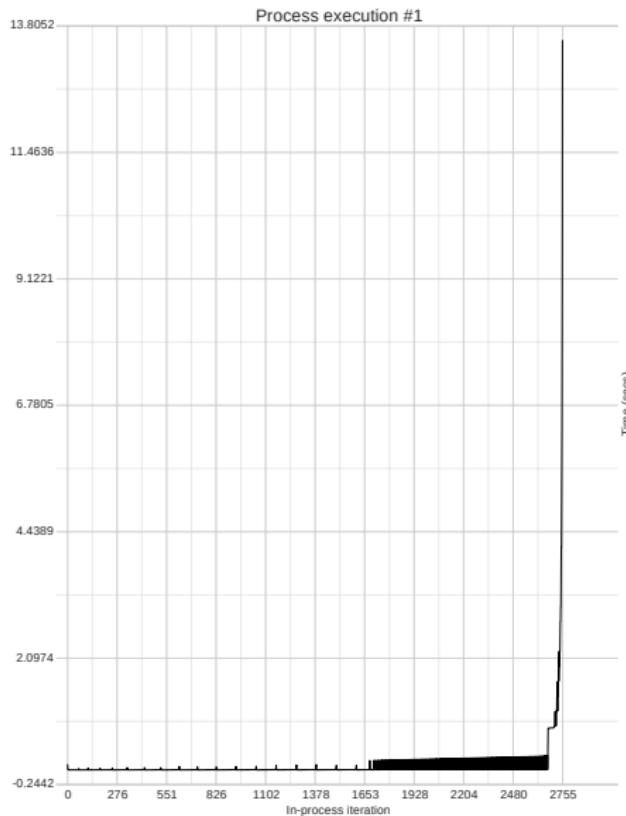
```
#  
# Fatal error in CALL_AND_RETRY_LAST  
# Allocation failed - process out of memory  
#
```

```
zsh: illegal hardware instruction  d8 run.js
```

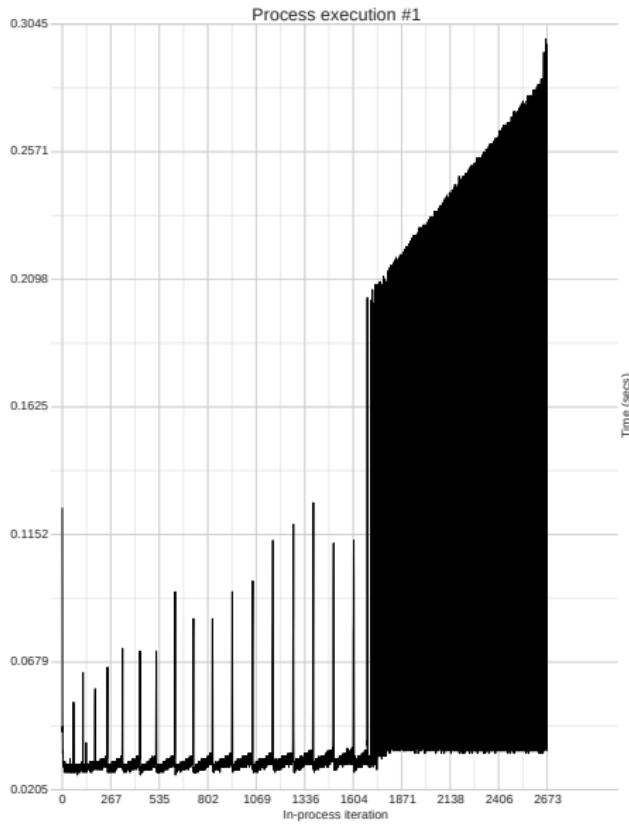
Octane: analysing pdf.js



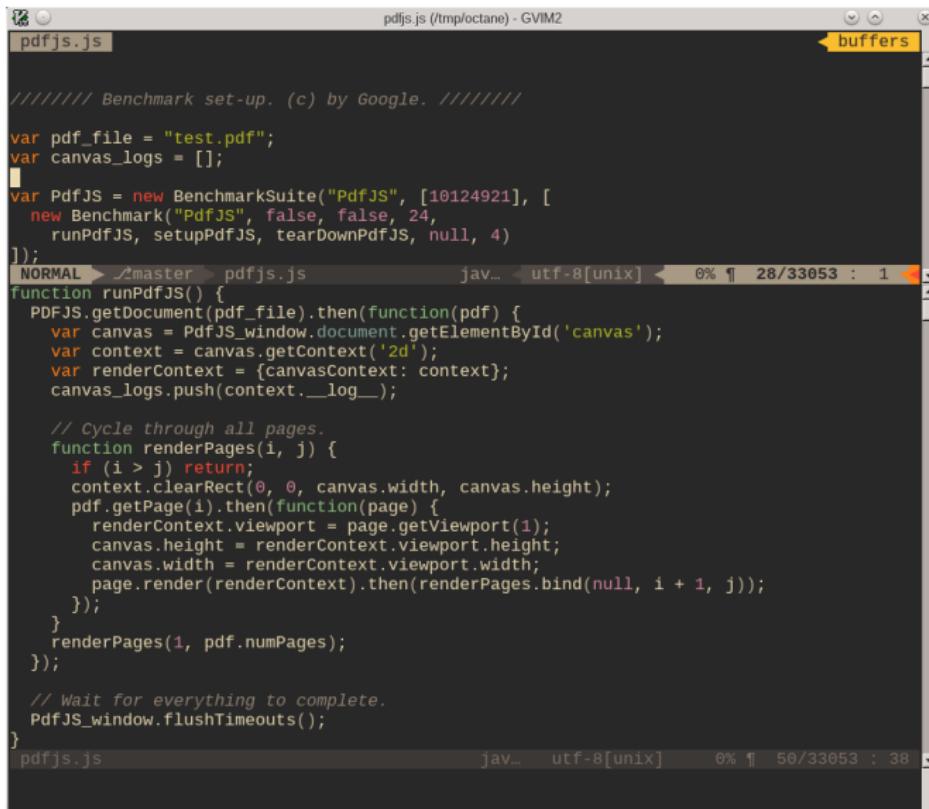
Octane: analysing pdf.js



Octane: analysing pdf.js



Octane: debugging



The screenshot shows a GVIM window with two buffers. The left buffer contains the following JavaScript code:

```
//////// Benchmark set-up. (c) by Google. /////////
var pdf_file = "test.pdf";
var canvas_logs = [];

var PdfJS = new BenchmarkSuite("PdfJS", [10124921], [
  new Benchmark("PdfJS", false, false, 24,
    runPdfJS, setupPdfJS, tearDownPdfJS, null, 4)
]);
function runPdfJS() {
  PDFJS.getDocument(pdf_file).then(function(pdf) {
    var canvas = PdfJS_window.document.getElementById('canvas');
    var context = canvas.getContext('2d');
    var renderContext = {canvasContext: context};
    canvas_logs.push(context.__log__);

    // Cycle through all pages.
    function renderPages(i, j) {
      if (i > j) return;
      context.clearRect(0, 0, canvas.width, canvas.height);
      pdf.getPage(i).then(function(page) {
        renderContext.viewport = page.getViewport(1);
        canvas.height = renderContext.viewport.height;
        canvas.width = renderContext.viewport.width;
        page.render(renderContext).then(renderPages.bind(null, i + 1, j));
      });
    }
    renderPages(1, pdf.numPages);
  });

  // Wait for everything to complete.
  PdfJS_window.flushTimeouts();
}
pdfjs.js
```

The right buffer is empty and labeled "buffers". The status bar at the bottom indicates the file is "pdfjs.js" and the line count is 38.

Octane: fixing

The screenshot shows a GitHub pull request page for issue #42. The title is "Fix memory leak in pdfjs.js. #42". A green button says "Open" and indicates "ltratt wants to merge 2 commits into chromium:master from ltratt:master". Below the title, there are three tabs: "Conversation 5", "Commits 2", and "Files changed 1". The "Files changed" tab is selected, showing changes from "all commits" to "1 file". The diff shows a single commit adding a line of code to initialize an array: "46 + canvas_logs.length = 0;".

```
1 pdfjs.js
@@ -43,6 +43,7 @@ function setupPdfJS() {
 43     }
 44
 45     function runPdfJS() {
 46         canvas_logs.length = 0;
 47         PDFJS.getDocument(pdf_file).then(function(pdf) {
 48             var canvas = PdfJS_window.document.getElementById('canvas');
 49             var context = canvas.getContext('2d');
```

Octane: other issues

pdfjs isn't the only problem

Octane: other issues

pdfjs isn't the only problem

CodeLoadClosure also has a memory leak

Octane: other issues

pdfjs isn't the only problem

CodeLoadClosure also has a memory leak

zlib complains that Cannot enlarge memory arrays in asm.js (a memory leak? I don't know)

Octane: other issues

pdfjs isn't the only problem

CodeLoadClosure also has a memory leak

zlib complains that Cannot enlarge memory arrays in asm.js (a memory leak? I don't know)

Timings are made with a non-monotonic microsecond timer

Summary

Summary

Why aren't more users more happy with
our VMs?

Summary

Why aren't more users more happy with our VMs?

My thesis: our benchmarking *and* our benchmarks have misled us

How to benchmark a bit better

How to benchmark a bit better

- 1 Run benchmarks for longer to uncover issues.

How to benchmark a bit better

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

How to benchmark a bit better

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

How to benchmark a bit better

- 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.

How to benchmark a bit better

- 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.

How to benchmark a bit better

- 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.

What we're doing next

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

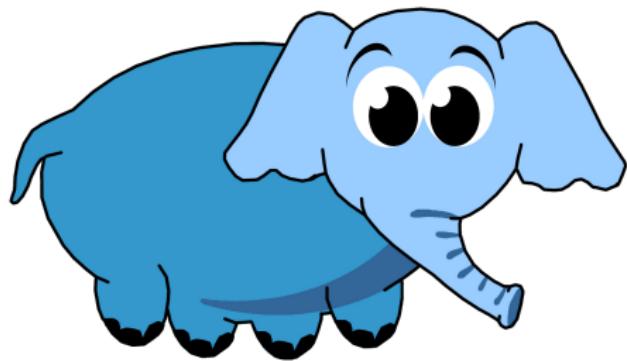
Thanks

- EPSRC: *COOLER* and *Lecture*.
- Oracle: various.

K Whiteford for Barry, the Benchmarking Elephant
in the Room

Thanks for listening

Thanks for listening



And don't forget Barry