

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



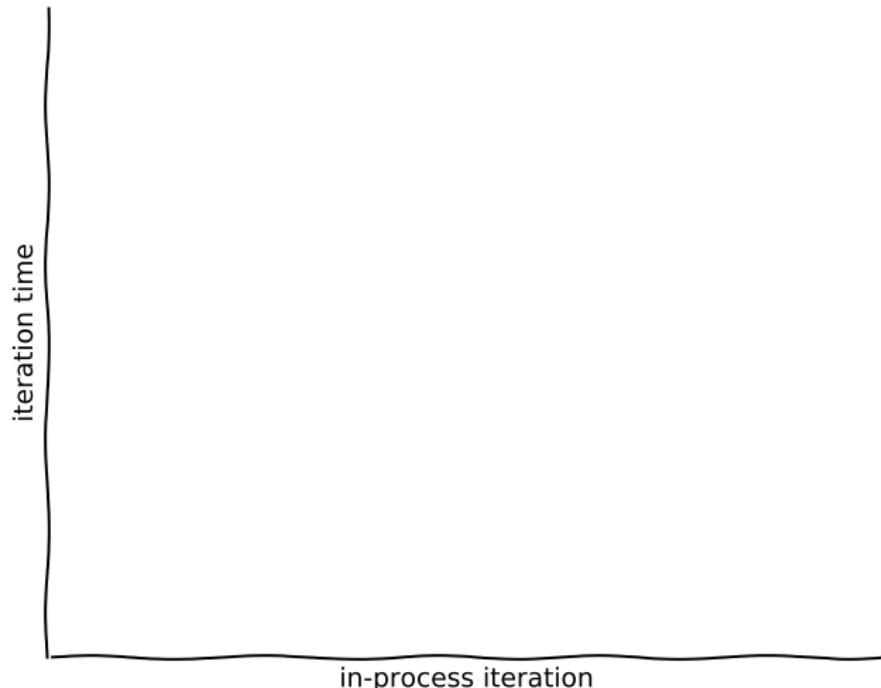
Software Development Team
2018-05-17

What to expect from this talk

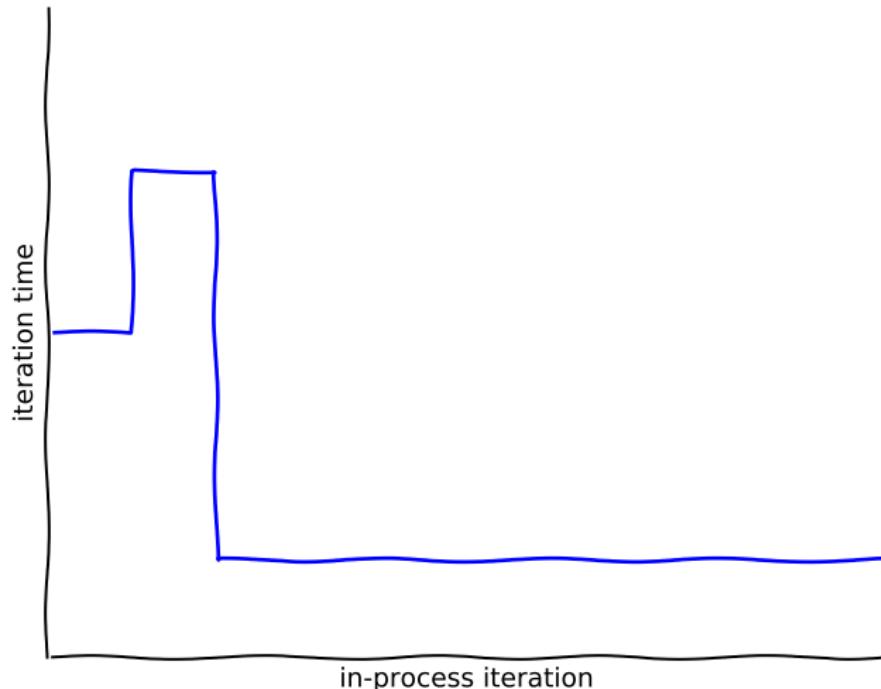
JVMs bring "gcc -O2"
to the masses

-Cliff Click: A JVM does that?

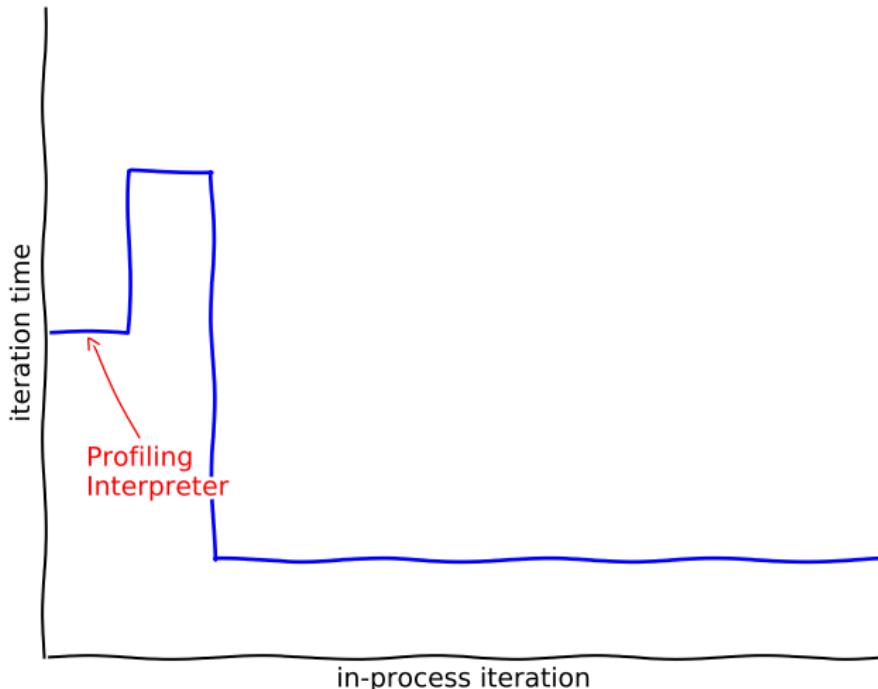
What do VM claims pertain to?



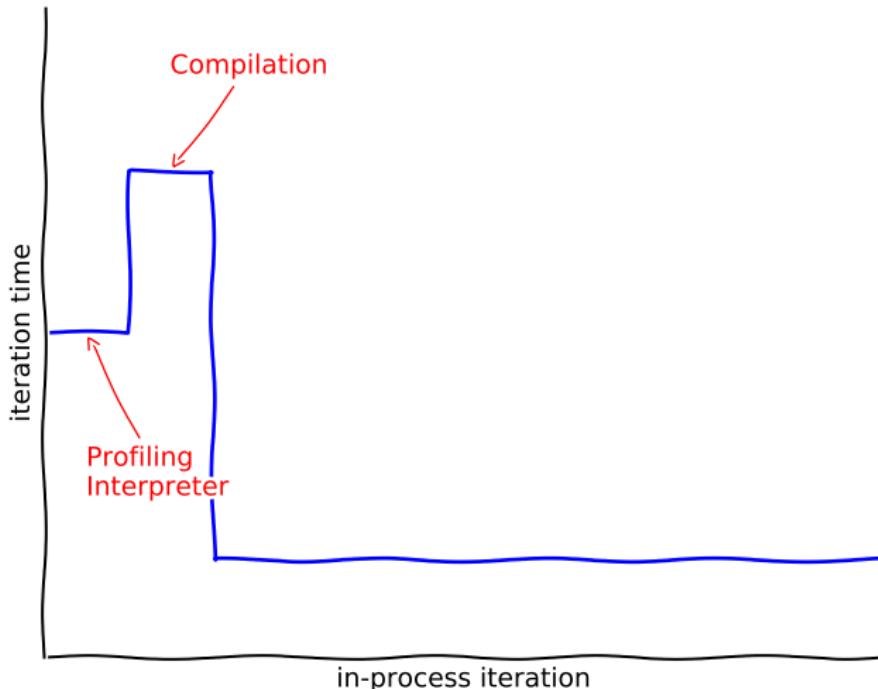
What do VM claims pertain to?



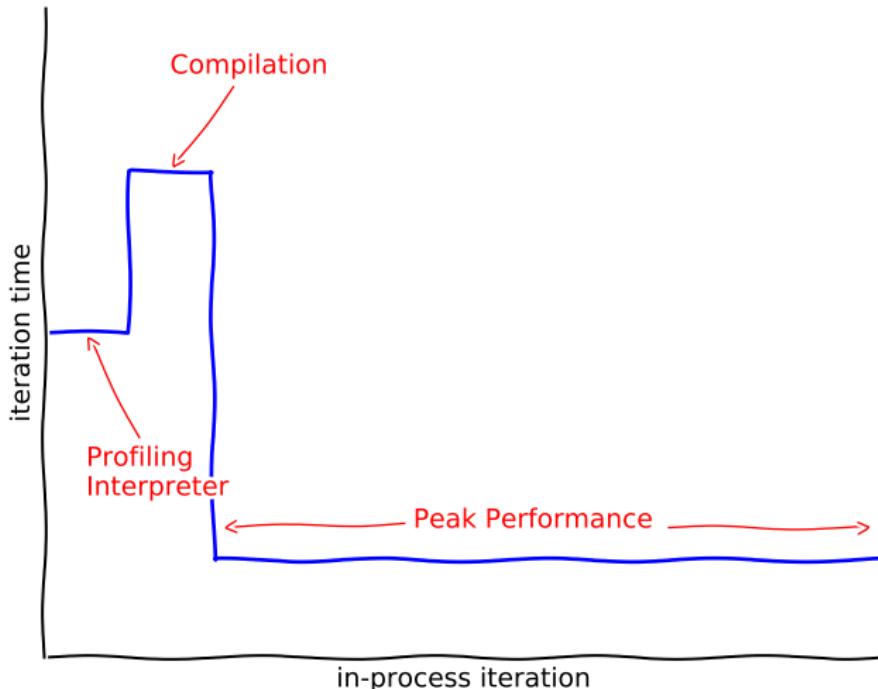
What do VM claims pertain to?



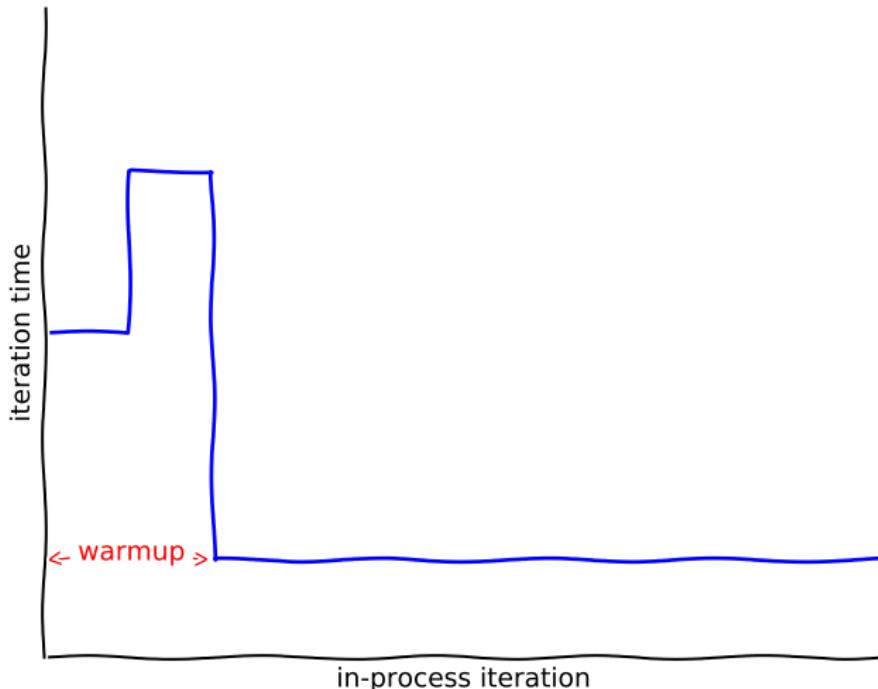
What do VM claims pertain to?



What do VM claims pertain to?



What do VM 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
- TruffleRuby 20170502
- 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
- Reruns any proc. exec where the CPU was throttled
- Checks `dmesg` for changes after each proc. exec
- Checks system at (roughly) same temperature for proc. execs
- Enforces kernel settings (tickless mode, CPU governors, ...)

And now for data...

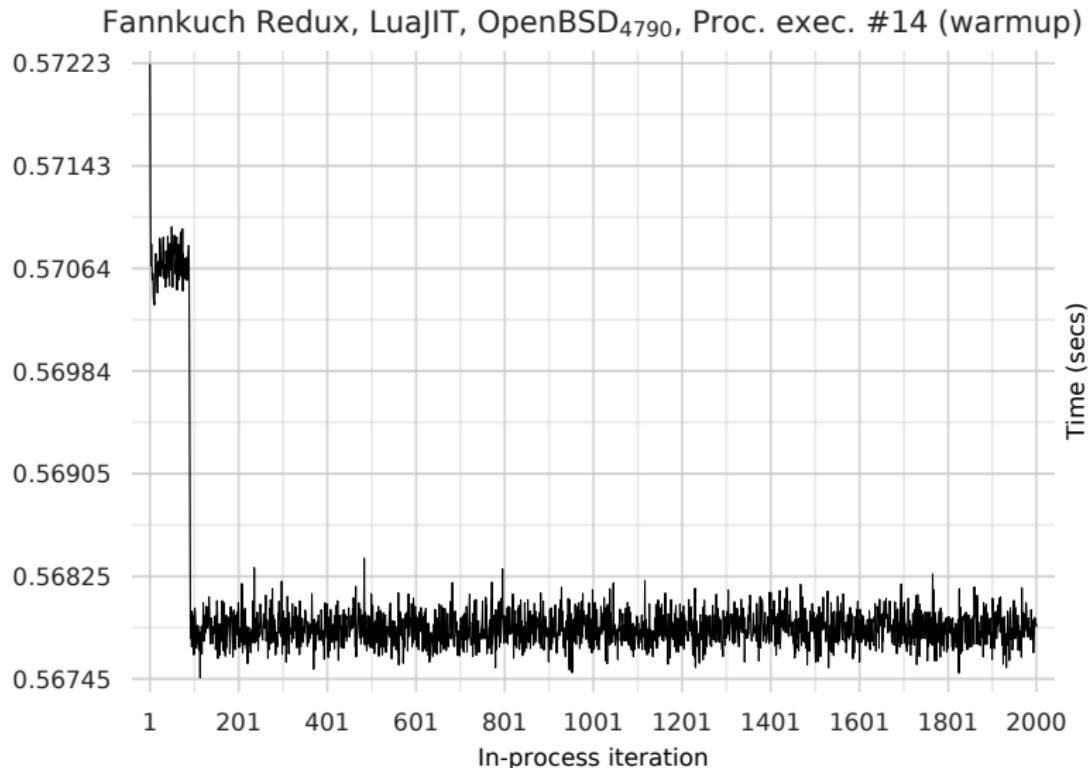
The experiment has gone through many versions.

And now for data...

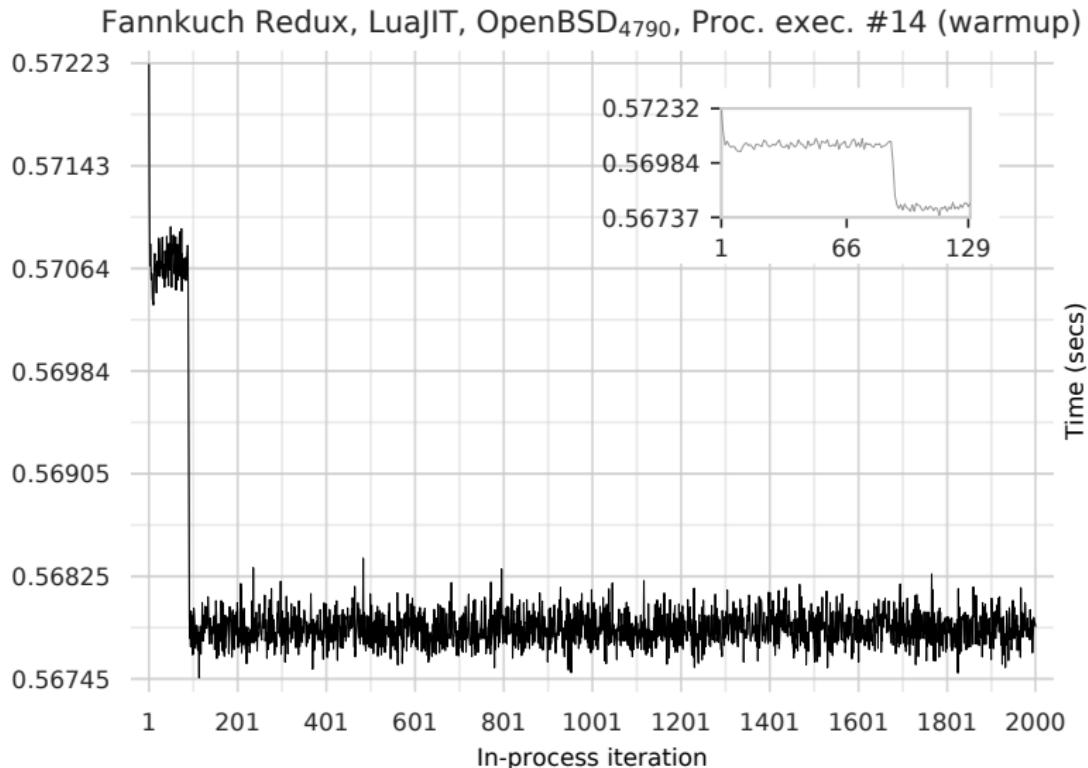
The experiment has gone through many versions.

The following data is from the 1.5 run.

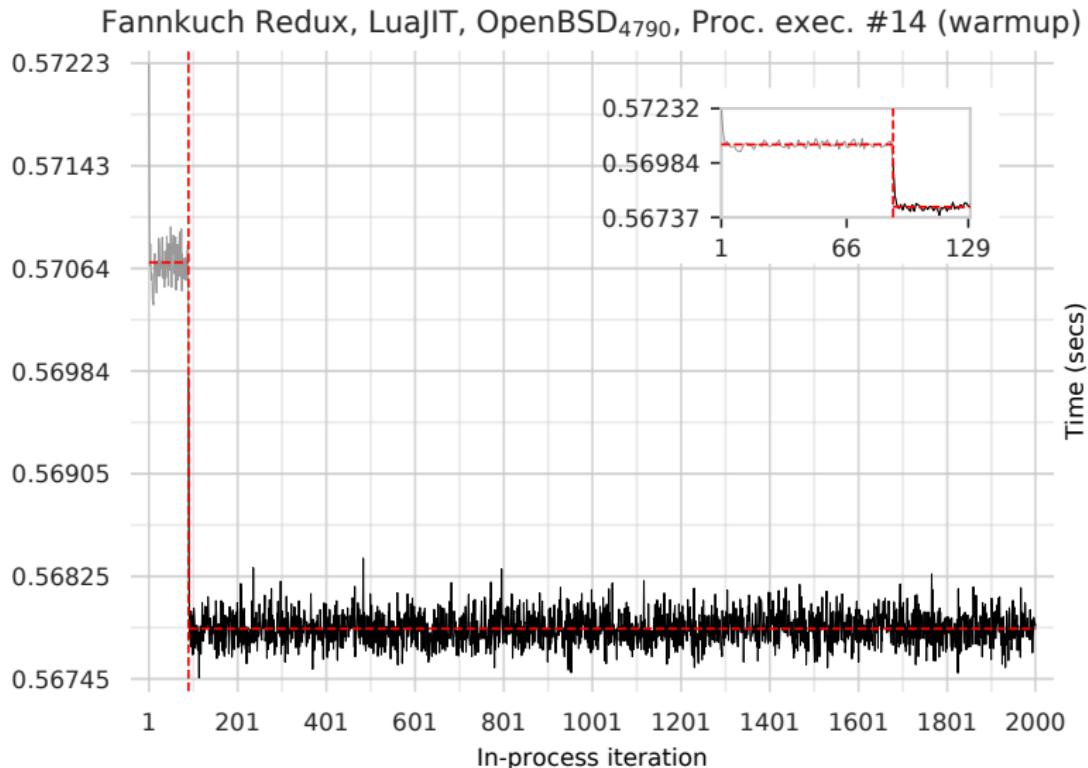
Warmup & flat (1)



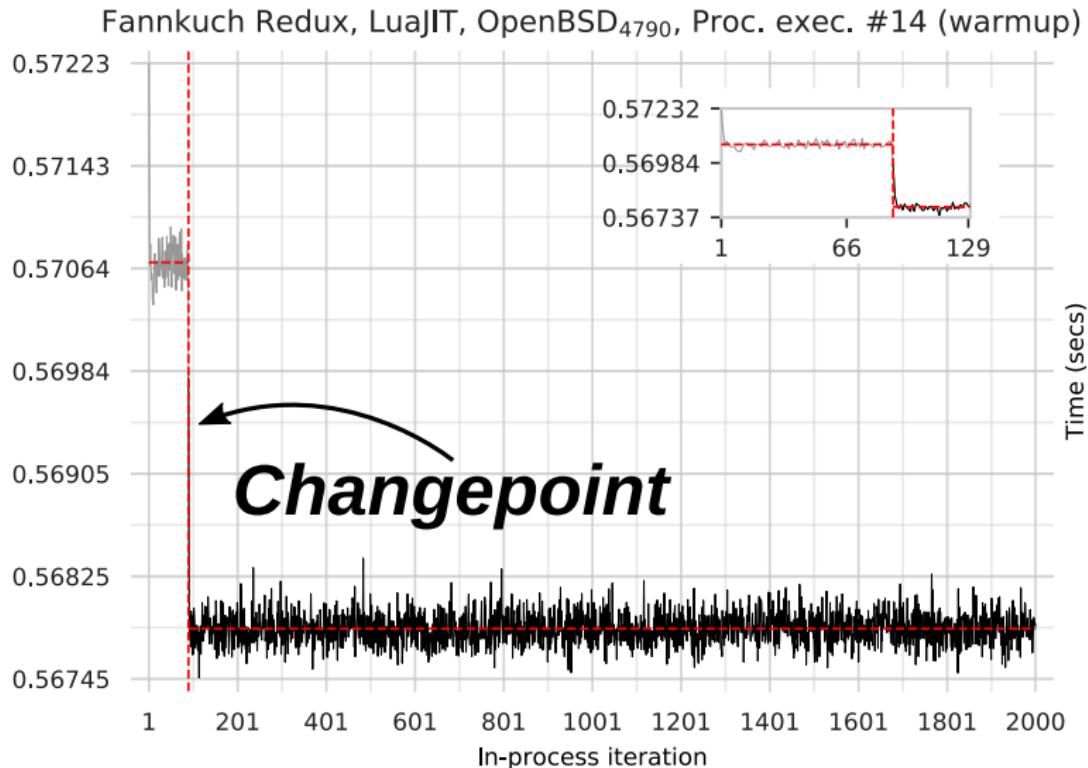
Warmup & flat (1)



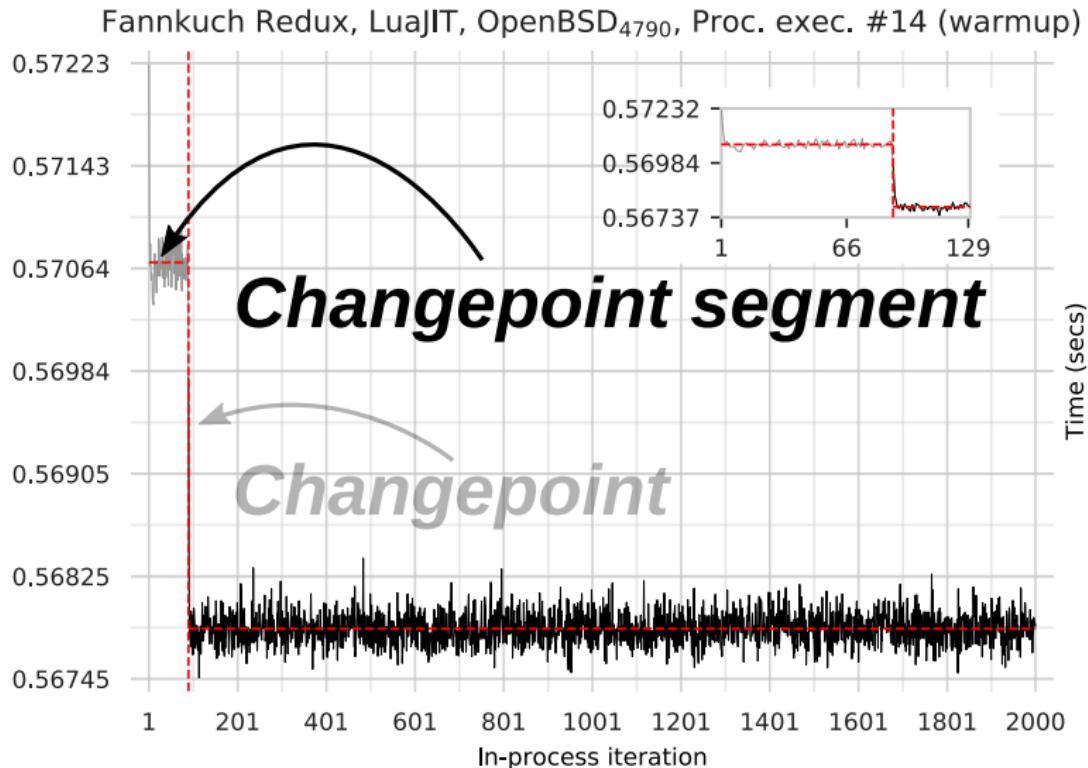
Warmup & flat (1)



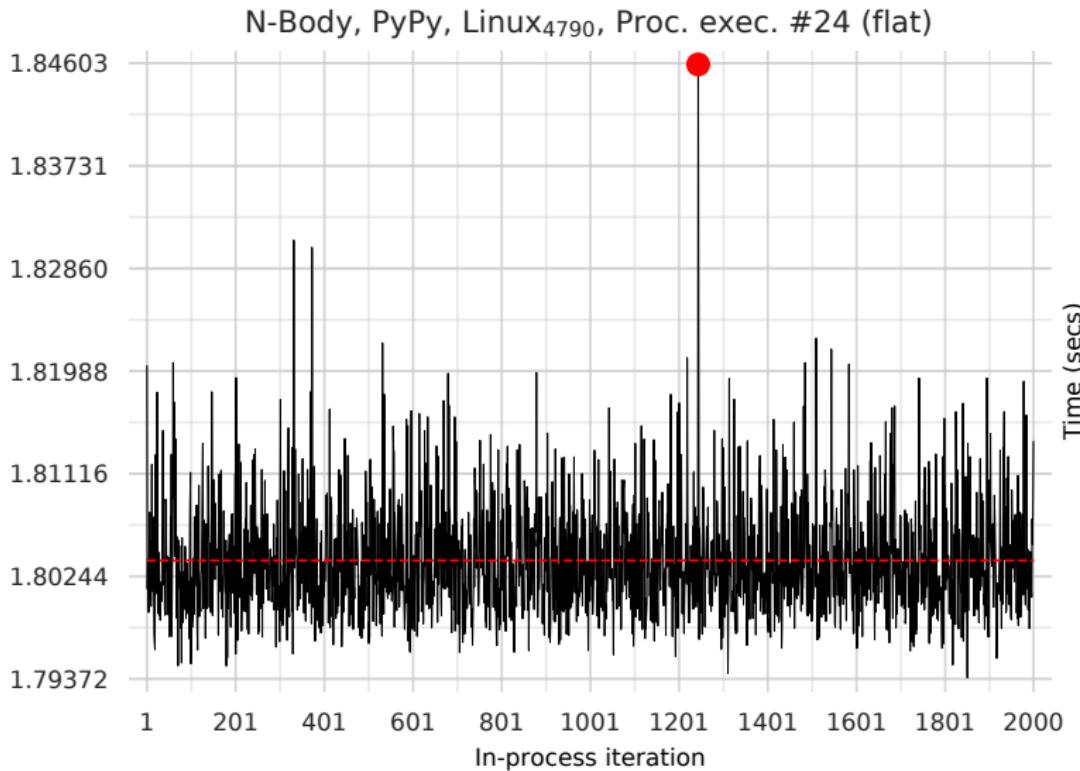
Warmup & flat (1)



Warmup & flat (1)



Warmup & flat (1)



Method 7: Classification

Classification algorithm (steps in order):

All segs are equivalent: *flat*

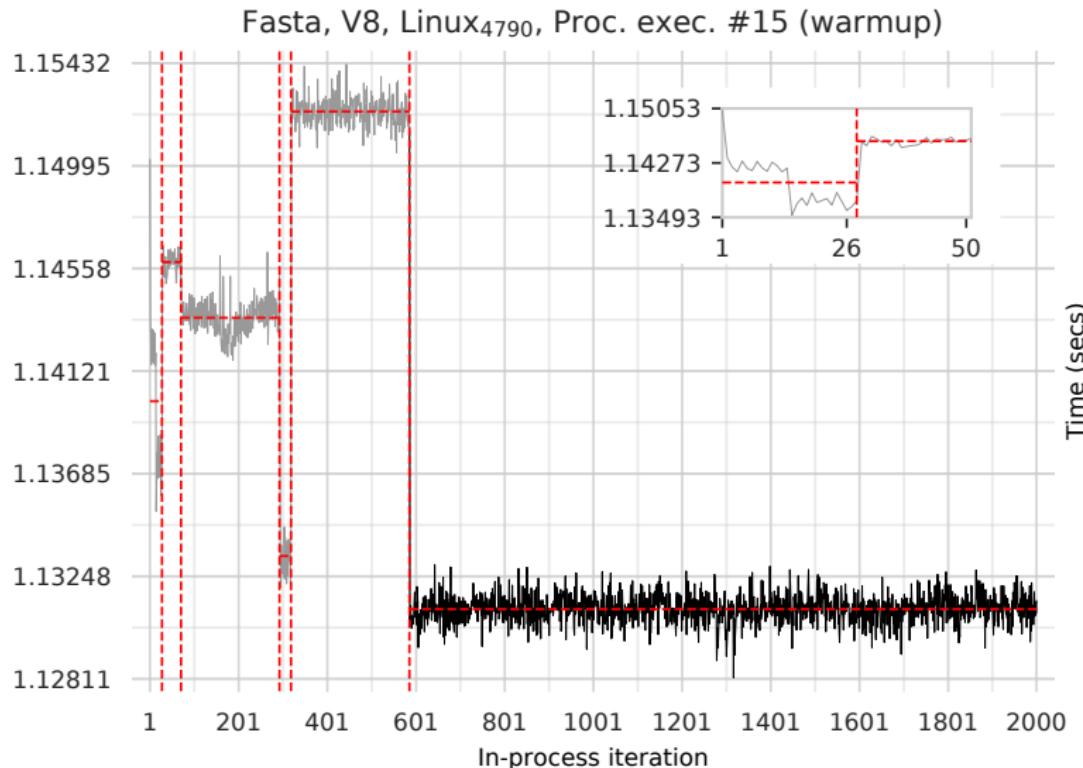
Method 7: Classification

Classification algorithm (steps in order):

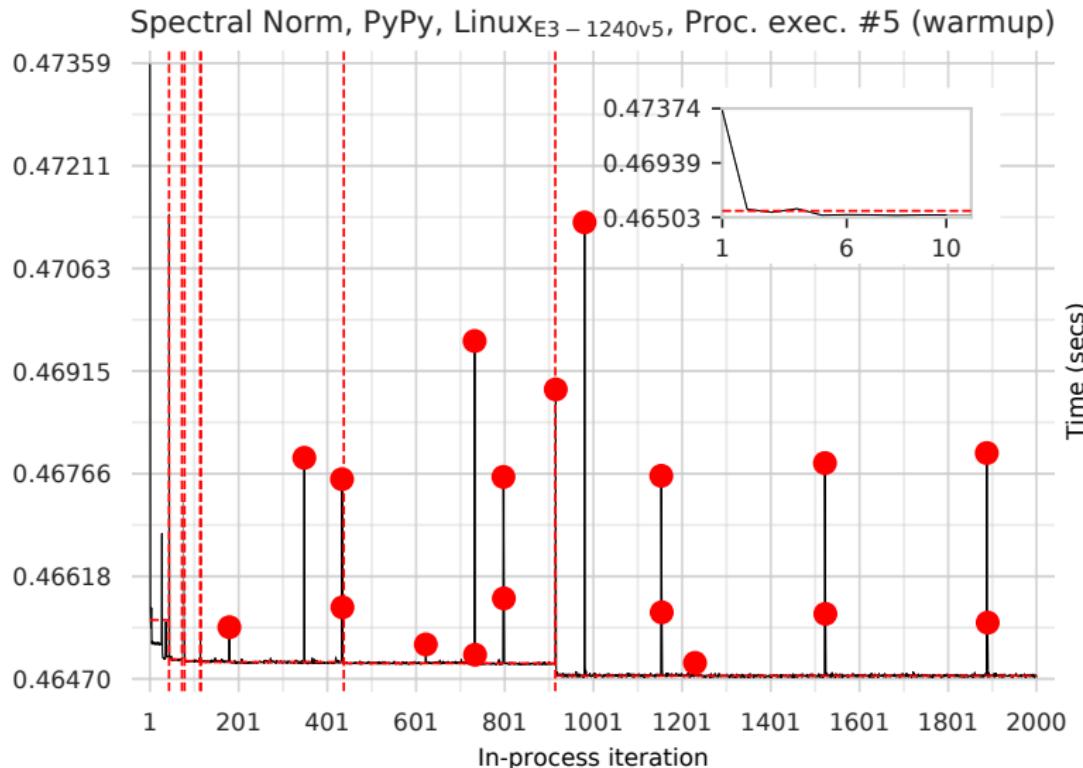
All segs are equivalent: *flat*

Final seg is in fastest set: *warmup*

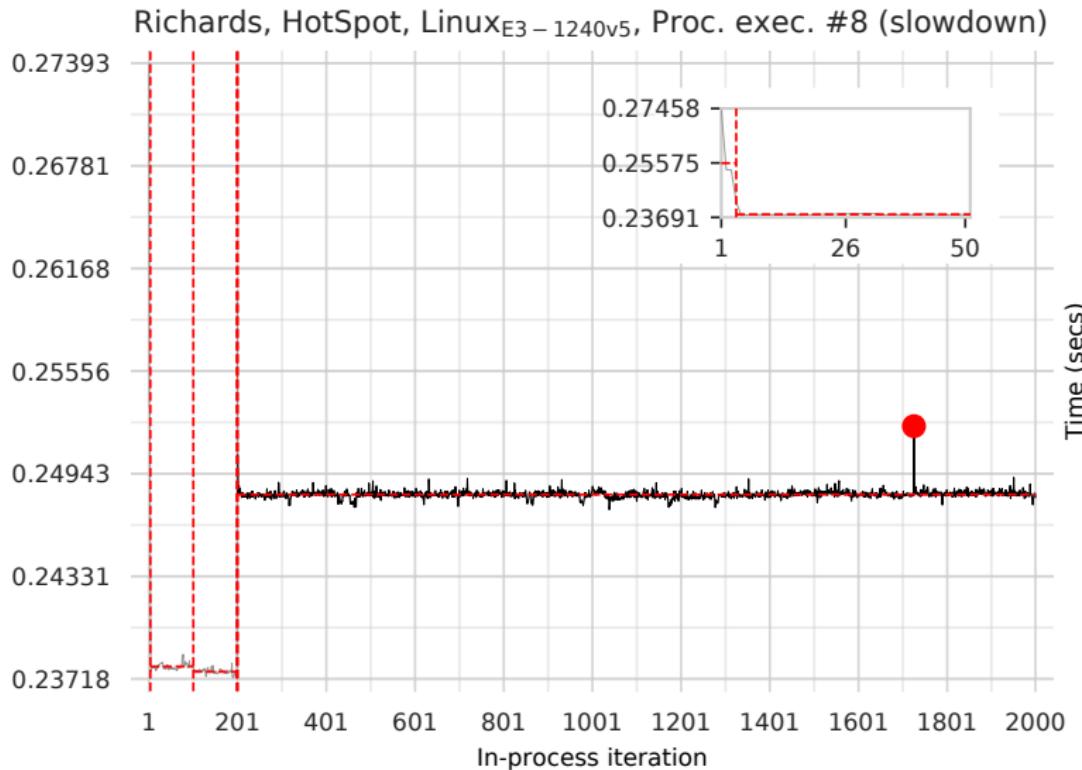
Warmup & flat (2)



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

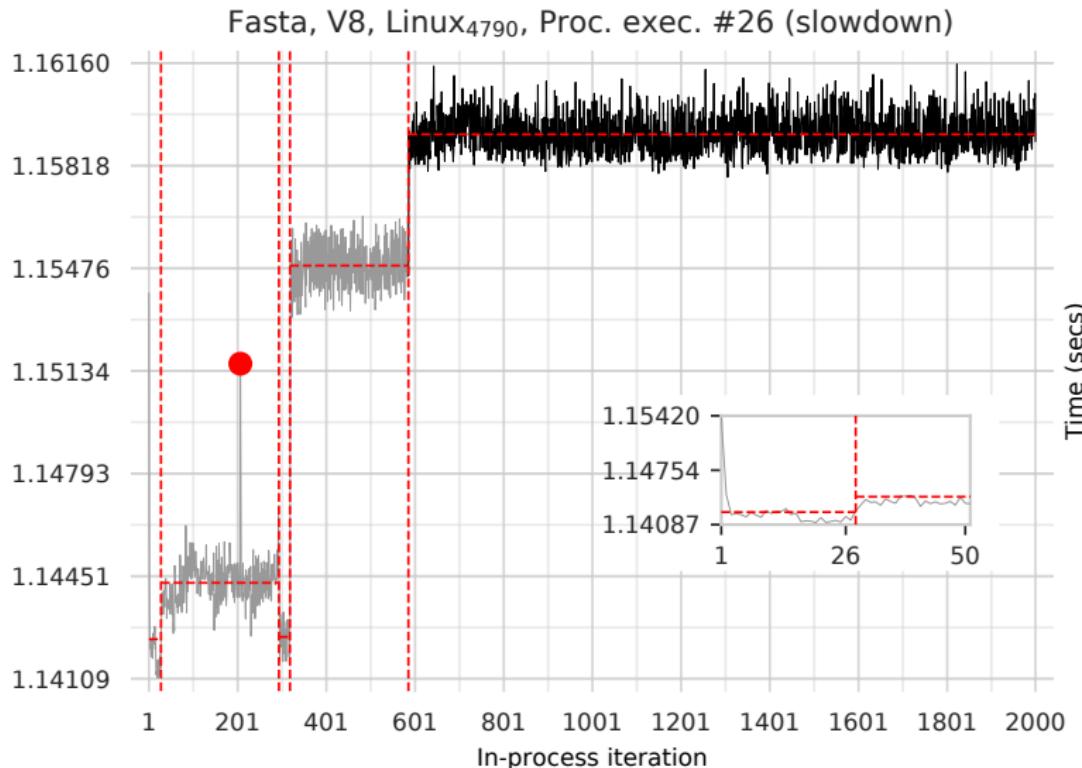
Classification algorithm (steps in order):

All segs are equivalent: *flat*

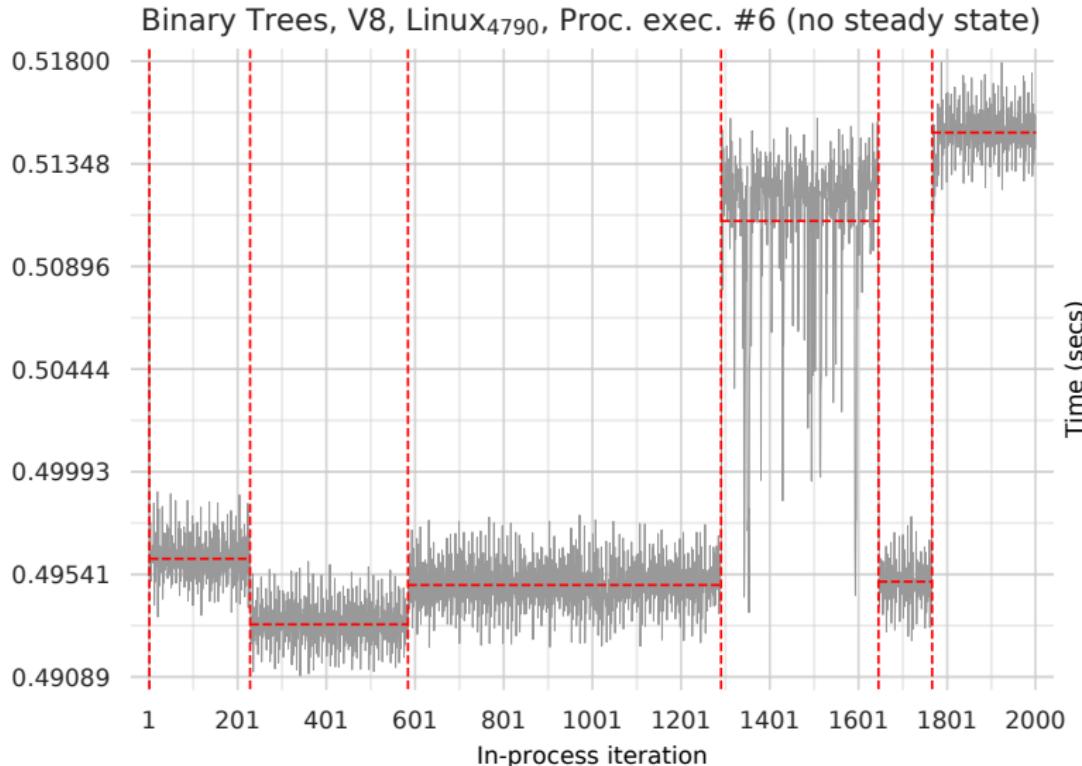
Final seg is in fastest set: *warmup*

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*

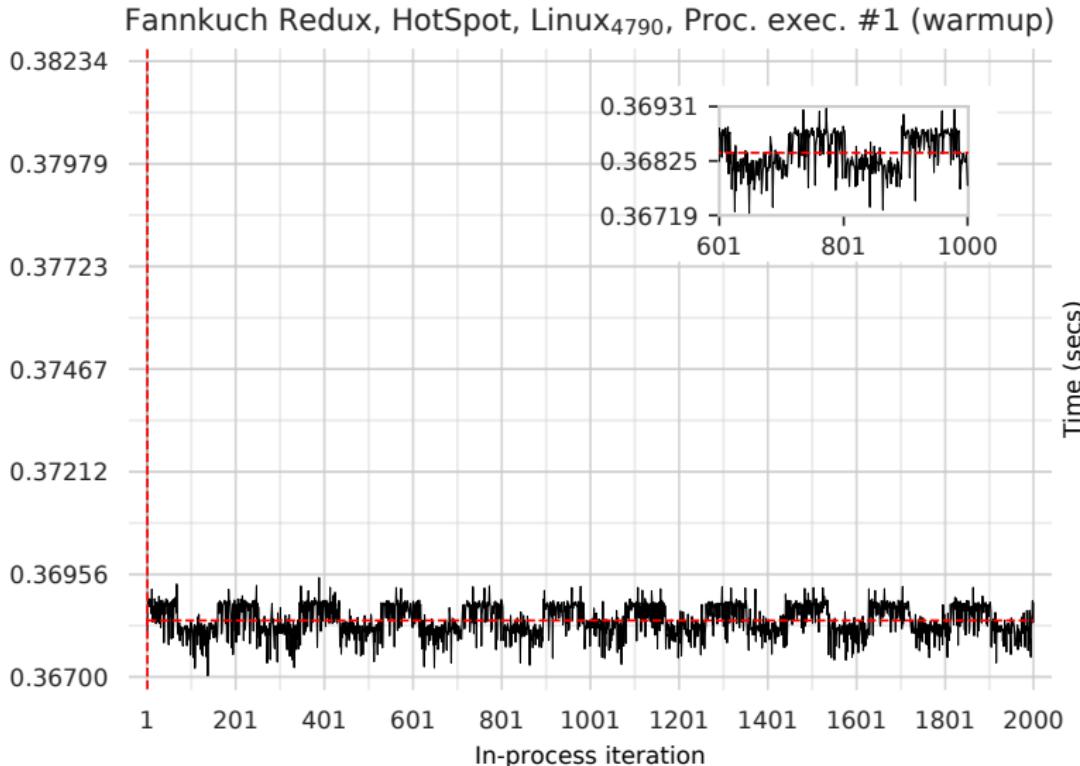
Final seg is in fastest set: *warmup*

Final seg is not in fastest set: *slowdown*

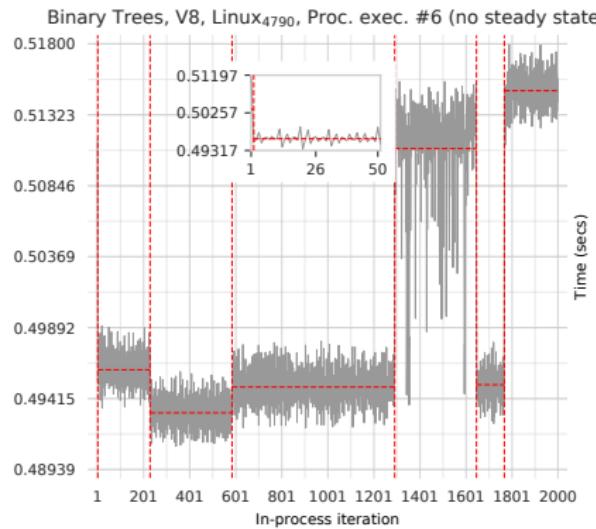
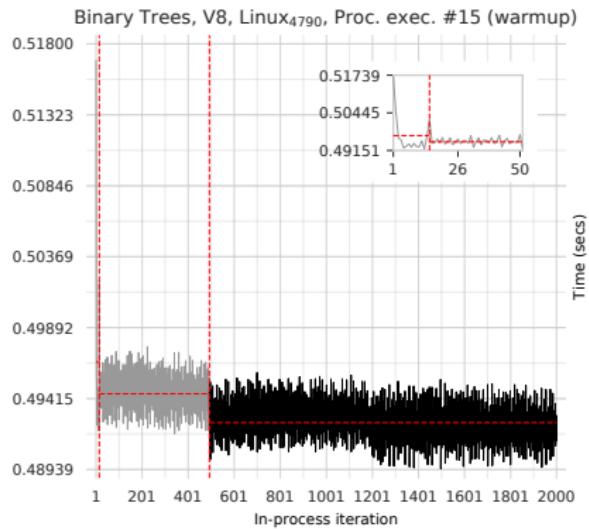
Else: *no steady state*

Bad

Warmup or no steady state?

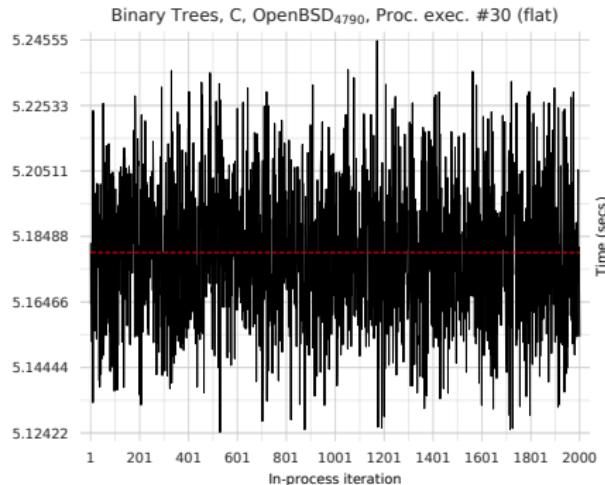
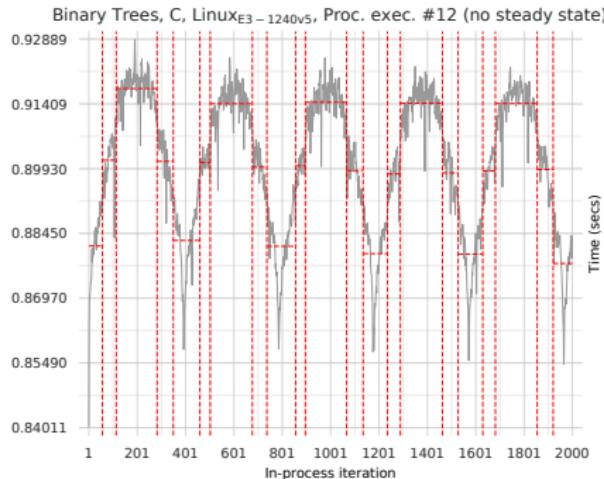


Inconsistent Process-executions



(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	w					-			0.40557 ±0.000185
Graal	■ (27L, 2f, 1w)				■	6.0 (5.0, 7.0)	0.86 0.00	0.13677 ±0.000035	
HHVM	■ (19L, 7f, 4w)				■	2.0 (2.0, 2.0)	0.14 0.00	0.13699 ±0.000029	
HotSpot	■ (15L, 13f, 2w)				■	1.0 (1.0, 1.0)	0.00 0.00	0.25401 ±0.002040	
LuaJIT	■ (25L, 3w, 2f)				■	— (0.000, 0.000)	— 0.00	0.18376 ±0.000463	
PyPy	■ (27w, 3f)				■	— (0.000, 0.000)	— 0.00	0.46998 ±0.001516	
TruffleRuby	■ (19f, 11L)	1003.0 (283.8, 1009.2)	2621.96 (749.093, 2639.554)	2.61136 ±0.036511	■	28L, 2f (68.0, 1003.0)	68.0 0.00	40.45 0.00	
V8	■ (13-, 9f, 8L)	228.0 (1.0, 588.0)	106.34 (0.000, 275.806)	0.46927 ±0.03135	■	= (20-, 10L) (1.0, 363.0)	1.0 0.00	0.00 0.00	
<i>binary trees</i>									
C	■ (25-, 5f)	1.0 (1.0, 538.0)	0.00 (0.000, 224.243)	0.41742 ±0.01343	■	18-, 8f, 4w (2.0, 35.3)	— 2.0	0.95 0.00	0.26465 ±0.007761
Graal	■ (26L, 4w)				■	— (0.879, 9.745)	— —	— —	— —
HHVM	■ (28L, 2w)								
HotSpot	■ (17f, 10L, 3w)				■	28f, 1w, 1L (1.023, 1.102)	— —	— —	— —
LuaJIT	—				■	18w, 7f, 5L (3.0, 189.8)	— 3.0	— 0.04	— 0.00
PyPy	■ (19-, 9f, 2f)	1.0 (1.0, 1120.3)	0.00 (0.000, 1802.797)	0.56287 ±0.00113	■	Richards (2.0, 2.0)	2.0 1021.0 (1.023, 1.102)	1.04 1521.33 0.00	0.88940 ±0.03495
TruffleRuby	■ (26f, 3w, 1L)				■	— (499.847, 1585.296)	— —	— —	— —
V8	= (28L, 2-)	2.0 (1.5, 25.0)	0.32 (0.142, 7.585)	0.30634 ±0.00160	■	— (3.0, 189.8)	— (0.961, 97.303)	0.47681 0.00	0.00 0.00
<i>funckuch reduce</i>									
C	—				■	775.0 (1.5, 780.0)	425.16 0.00	0.54581 ±0.03116	
Graal	■	4.0 (3.0, 34.4)	0.75 (0.535, 5.873)	0.16188 ±0.000738	■	— (2.0, 94.6)	14.0 13.60 (0.830, 98.737)	1.05685 0.00	
HHVM	■ (13f, 10L, 7w)				■	— (29L, 1w)	— —	— —	— —
HotSpot	■ (27L, 3f)	6.0 (5.0, 595.0)	0.62 (0.500, 70.062)	0.11691 ±0.001949	■	— (7.0, 7.5)	— 7.0 (1.902, 3.645)	0.31472 0.00	0.169143 ±0.002181
LuaJIT	■ (19f, 9-, 1w, 1L)				■	— (27-, 3L)	— 1.0 (1.0, 45.2)	— 0.00 (0.000, 20.597)	— 0.46480 0.00
PyPy	w				■	— (25L, 5w)	— 3.0 (3.0, 3.0)	— 0.52 (0.523, 0.526)	— 0.25362 0.00
TruffleRuby									
V8	■ (27L, 3L)	586.0 (318.0, 587.6)	612.12 (330.270, 622.328)	1.04904 ±0.020905	■				
<i>fasta</i>									
C	—				■	— (27L, 2-, 1L)	— 775.0 (1.5, 780.0)	— 425.16 0.00	0.54581 ±0.03116
Graal	■	4.0 (3.0, 34.4)	0.75 (0.535, 5.873)	0.16188 ±0.000738	■	— (2.0, 94.6)	14.0 13.60 (0.830, 98.737)	1.05685 0.00	
HHVM	■ (13f, 10L, 7w)				■	— (29L, 1w)	— —	— —	— —
HotSpot	■ (27L, 3f)	6.0 (5.0, 595.0)	0.62 (0.500, 70.062)	0.11691 ±0.001949	■	— (7.0, 7.5)	— 7.0 (1.902, 3.645)	0.31472 0.00	0.169143 ±0.002181
LuaJIT	■ (19f, 9-, 1w, 1L)				■	— (27-, 3L)	— 1.0 (1.0, 45.2)	— 0.00 (0.000, 20.597)	— 0.46480 0.00
PyPy	w				■	— (25L, 5w)	— 3.0 (3.0, 3.0)	— 0.52 (0.523, 0.526)	— 0.25362 0.00
TruffleRuby									
V8	■ (27L, 3L)	586.0 (318.0, 587.6)	612.12 (330.270, 622.328)	1.04904 ±0.020905	■				
<i>spectrenorm</i>									
C	—				■	— (27L, 2-, 1L)	— 775.0 (1.5, 780.0)	— 425.16 0.00	0.54581 ±0.03116
Graal	■	4.0 (3.0, 34.4)	0.75 (0.535, 5.873)	0.16188 ±0.000738	■	— (2.0, 94.6)	14.0 13.60 (0.830, 98.737)	1.05685 0.00	
HHVM	■ (13f, 10L, 7w)				■	— (29L, 1w)	— —	— —	— —
HotSpot	■ (27L, 3f)	6.0 (5.0, 595.0)	0.62 (0.500, 70.062)	0.11691 ±0.001949	■	— (7.0, 7.5)	— 7.0 (1.902, 3.645)	0.31472 0.00	0.169143 ±0.002181
LuaJIT	■ (19f, 9-, 1w, 1L)				■	— (27-, 3L)	— 1.0 (1.0, 45.2)	— 0.00 (0.000, 20.597)	— 0.46480 0.00
PyPy	w				■	— (25L, 5w)	— 3.0 (3.0, 3.0)	— 0.52 (0.523, 0.526)	— 0.25362 0.00
TruffleRuby									
V8	■ (27L, 3L)	586.0 (318.0, 587.6)	612.12 (330.270, 622.328)	1.04904 ±0.020905	■				

Individual benchmark stats

		Steady iter (#)	Steady iter (s)	Steady perf (s)
C	<i>spectralnorm</i>	✗ (27L, 2-, 1L) 14.0 ✗ (29L, 1w) 7.0 — = (27-, 3L) ✗ (25L, 5w) 3.0	775.0 (1.5,780.0) 14.0 (2.0,94.6) 7.0 (7.0,7.5) — 1.0 (1.0,45.2) 0.00 (0.000,20.597) 0.52 (0.523,0.526)	425.16 (0.246,426.809) 13.60 (0.830,98.737) 1.91 (1.902,3.645) 0.22181 ±0.000039 0.46480 ±0.000085 0.25362 ±0.000034
Graal				1.05685 ±0.000126
HHVM				
HotSpot				0.31472 ±0.169143
LuaJIT				0.22181 ±0.000039
PyPy				
TruffleRuby				
V8				

Overall benchmark stats

Class.	Linux ₄₇₉₀	Linux _{1240v5}	OpenBSD ₄₇₉₀ [†]
⟨VM, benchmark⟩ pairs			
—	8.9%	11.1%	13.3%
⊓	20.0%	17.8%	20.0%
⊓	4.4%	4.4%	3.3%
⊜	4.4%	4.4%	0.0%
=	11.1%	8.9%	13.3%
✗	51.1%	53.3%	50.0%
Process executions			
—	22.0%	23.3%	37.7%
⊓	48.3%	43.9%	35.2%
⊓	20.1%	22.1%	12.1%
⊜	9.6%	10.8%	15.0%

Overall benchmark stats

Class.	Linux ₄₇₉₀	Linux _{1240v5}	OpenBSD ₄₇₉₀ [†]
⟨VM, benchmark⟩ pairs			
—	8.9%	11.1%	13.3%
⊓	20.0%	17.8%	20.0%
⊓	4.4%	4.4%	3.3%
⊜	4.4%	4.4%	0.0%
=	11.1%	8.9%	13.3%
✗	51.1%	53.3%	50.0%
Process executions			
—	22.0%	23.3%	37.7%
⊓	48.3%	43.9%	35.2%
⊓	20.1%	22.1%	12.1%
⊜	9.6%	10.8%	15.0%

Summary

Classical warmup occurs for only:

Summary

Classical warmup occurs for only:

67.2%-70.3% of process executions

Summary

Classical warmup occurs for only:

67.2%-70.3% of process executions

37.8%-40.0% of (VM, benchmark) pairs

Summary

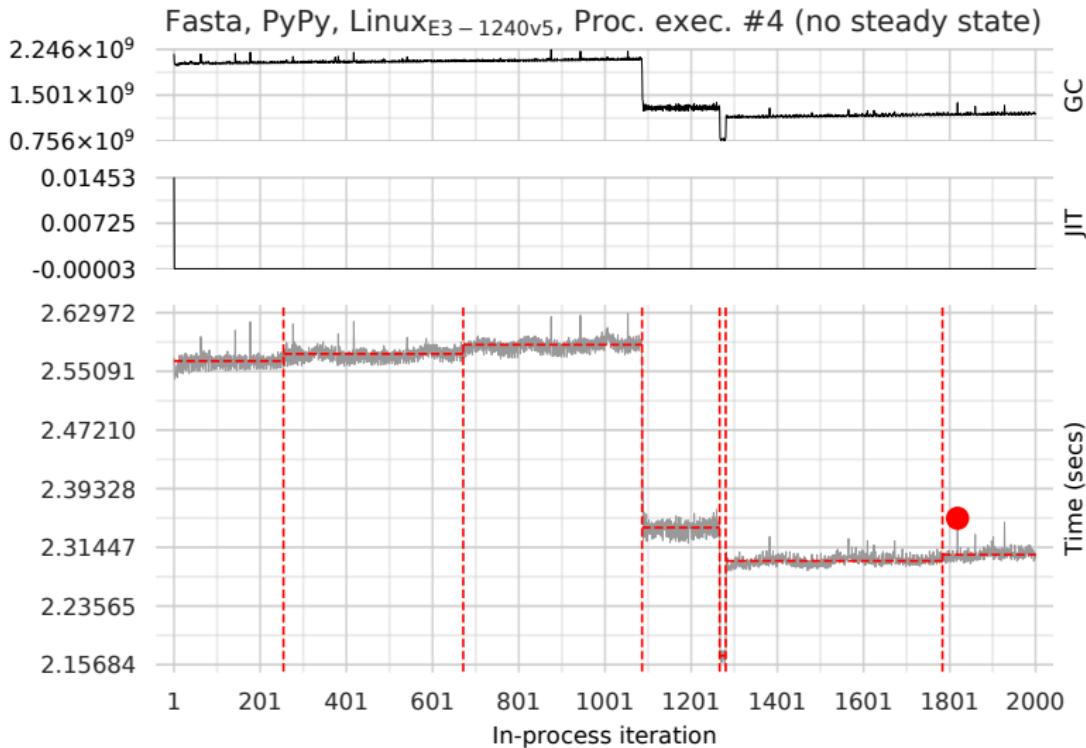
Classical warmup occurs for only:

67.2%-70.3% of process executions

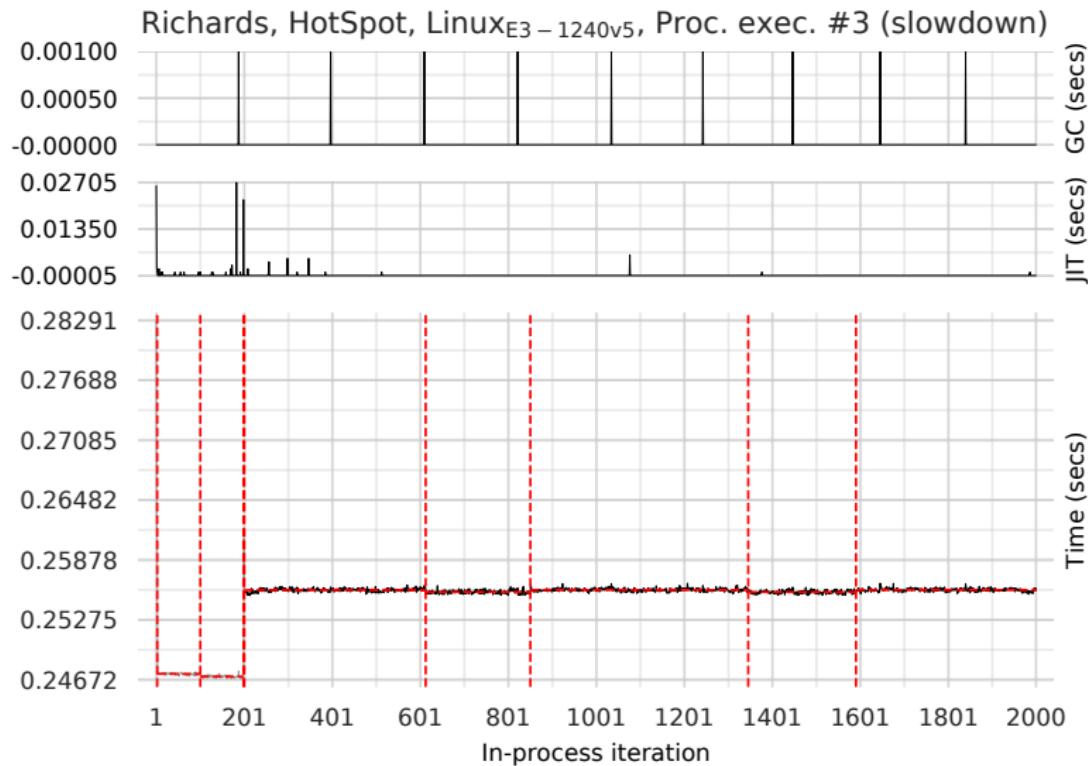
37.8%-40.0% of (VM, benchmark) pairs

12.5% of benchmarks for (VM, benchmark, machine) triples

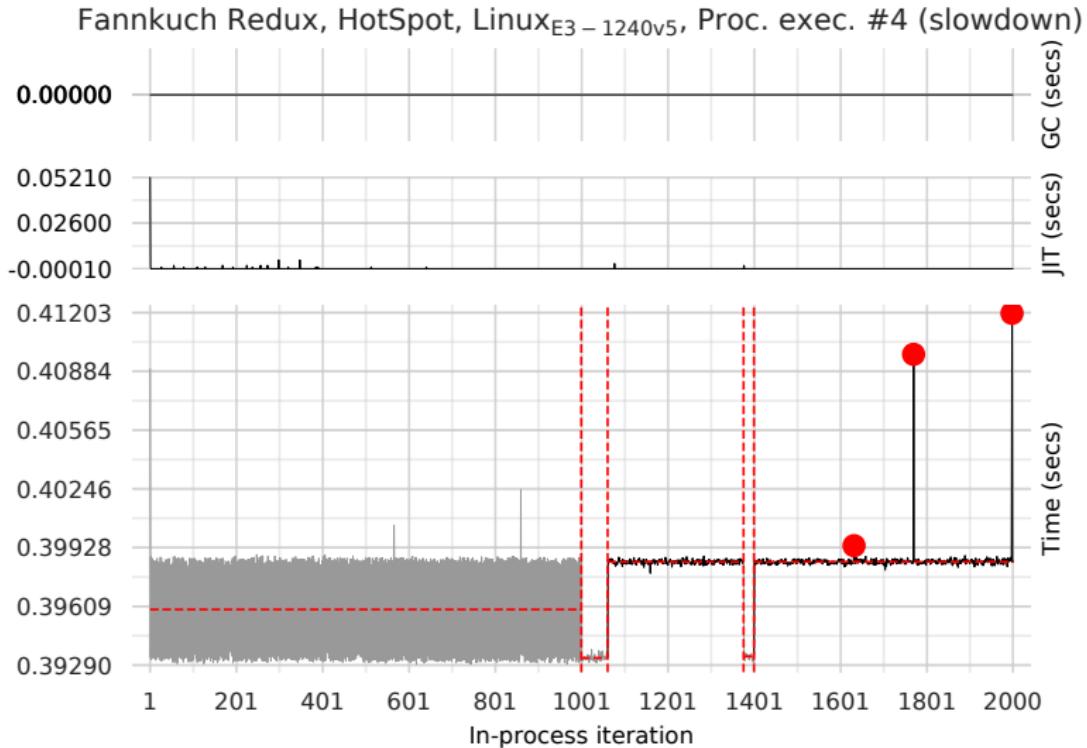
Are odd effects correlated with compilation and GC?



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 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

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 =====
```

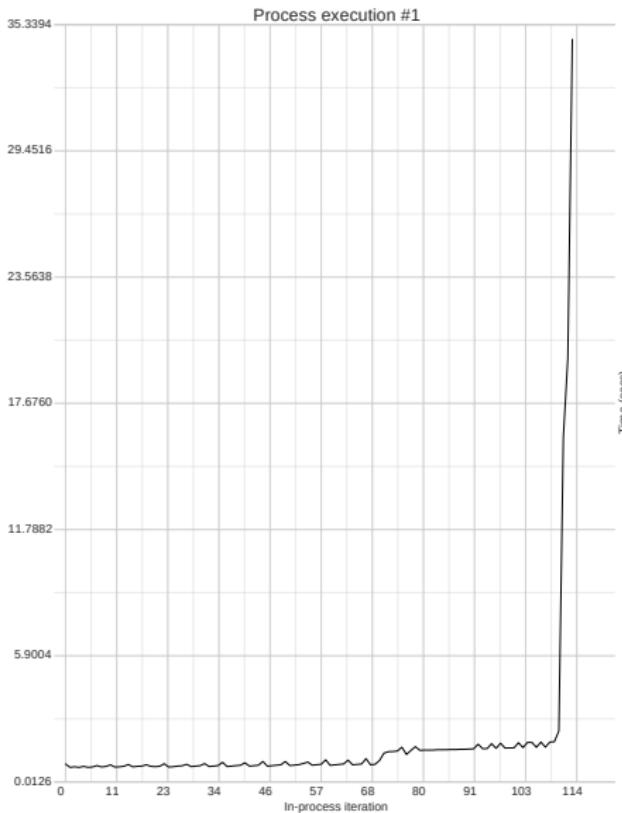
```
Security context: 0x20d333ad3ba9 <JS Object>
```

```
 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>,
```

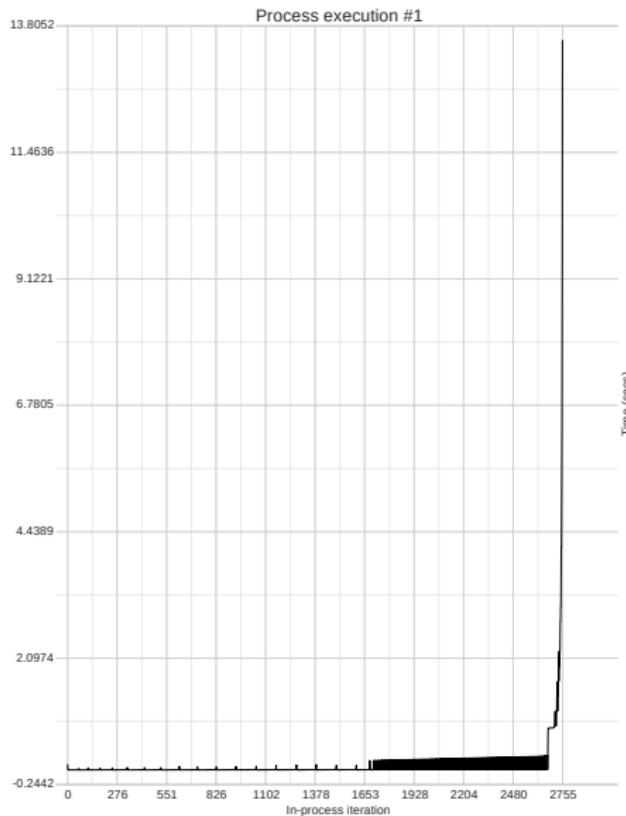
```
#  
# Fatal error in CALL_AND_RETRY_LAST  
# Allocation failed - process out of memory  
#
```

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

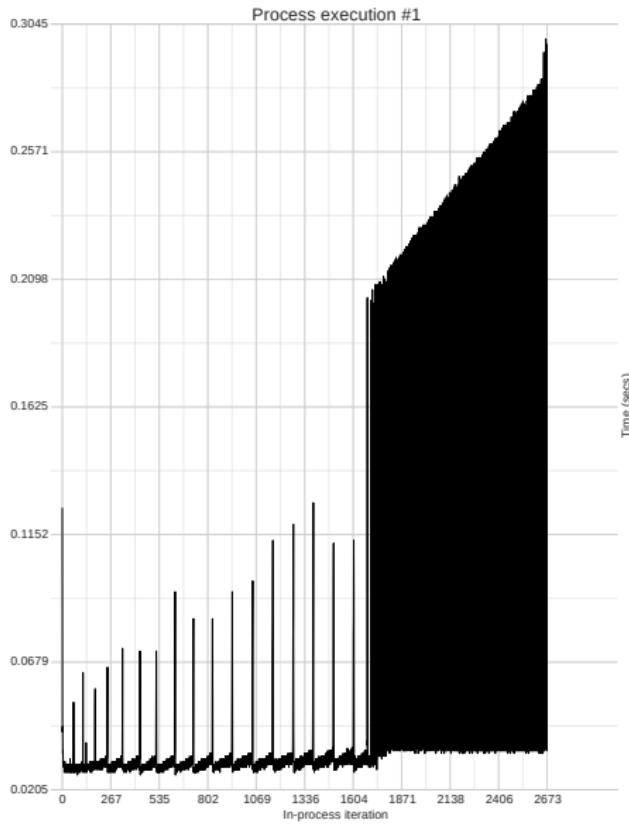
Octane: analysing pdf.js



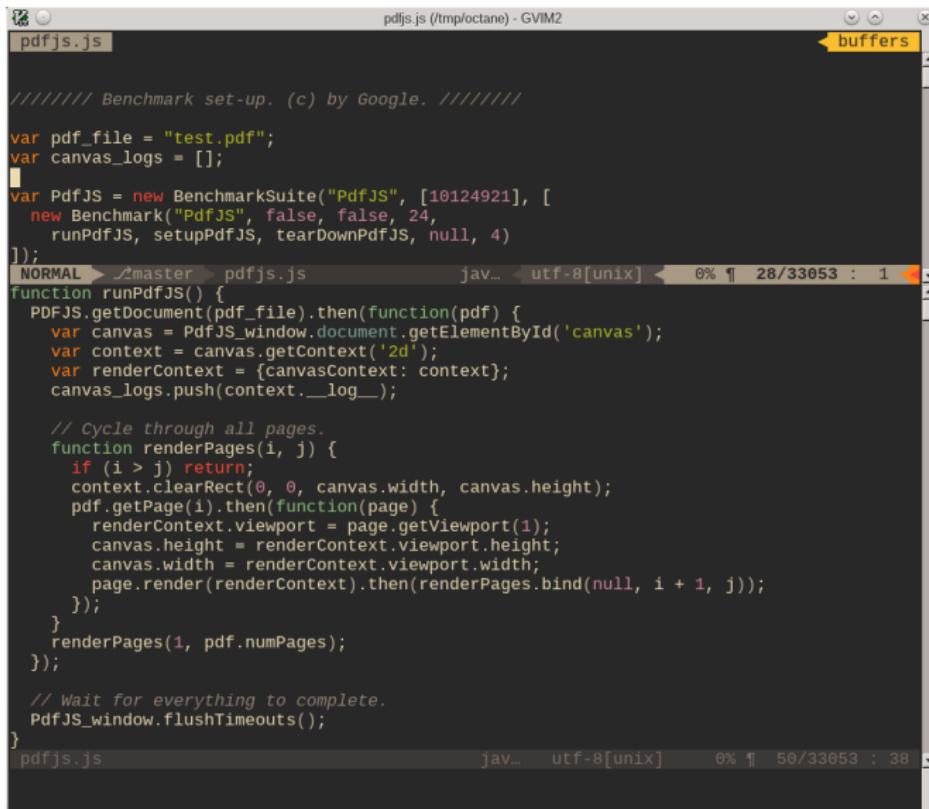
Octane: analysing pdf.js



Octane: analysing pdf.js



Octane: debugging



The screenshot shows a GVIM window with two buffers. The left buffer contains the following JavaScript code for PdfJS benchmark setup and rendering:

```
///////// 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 "Open" button indicates the merge status. The commit message is "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 a diff for "pdfjs.js". The diff highlights a new line of code: "+ canvas_logs.length = 0;". The code snippet is as follows:

```
@@ -43,6 +43,7 @@ function setupPdfJS() {  
 43   43 }  
 44   44  
 45   45 function runPdfJS() {  
+ 46   46   canvas_logs.length = 0;  
 47   47   PDFJS.getDocument(pdf_file).then(function(pdf) {  
 48   48     var canvas = PdfJS_window.document.getElementById('canvas');  
 49   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: benchmarking *and* benchmarks are performance destiny.

Summary

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

My thesis: benchmarking *and* benchmarks are performance destiny.

Ours 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 neither peak performance or a steady state may occur.

How to benchmark a bit better

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that neither peak performance or a steady state may occur.
- 3 Always report warmup time.

How to benchmark a bit better

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that neither peak performance or a steady state may 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 neither peak performance or a steady state may occur.
- 3 Always report warmup time.
- 4 Stop over-training on small benchmark suites.
- 5 Collect more benchmarks.

How to benchmark a bit better

- 1 Run benchmarks for longer to uncover issues.
- 2 Accept that neither peak performance or a steady state may occur.
- 3 Always report warmup time.
- 4 Stop over-training on small benchmark suites.
- 5 Collect more benchmarks.
- 6 Focus on predictable performance.

The big question

The big question

Can we fix existing VMs?

The big question

Can we fix existing VMs?

At least a bit... but a lot? Unclear.

The big question

Can we fix existing VMs?

At least a bit... but a lot? Unclear.

In case we can't, I have an idea...

Meta-tracing JITs

FL Interpreter

```
program_counter = 0; stack = []
vars = {...}
while True:
    jit_merge_point(program_counter)
    instr = load_instruction(program_counter)
    if instr == INSTR_VAR_GET:
        stack.push(
            vars[read_var_name_from_instruction()])
        program_counter += 1
    elif instr == INSTR_VAR_SET:
        vars[read_var_name_from_instruction()] =
            stack.pop()
        program_counter += 1
    elif instr == INSTR_INT:
        stack.push(read_int_from_instruction())
        program_counter += 1
    elif instr == INSTR_LESS_THAN:
        rhs = stack.pop()
        lhs = stack.pop()
        if isinstance(lhs, int) and isinstance(rhs, int):
            if lhs < rhs:
                stack.push(True)
            else:
                stack.push(False)
        else: ...
        program_counter += 1
```

```
elif instr == INSTR_IF:
    result = stack.pop()
    if result == True:
        program_counter += 1
    else:
        program_counter +=
            read_jump_if_instruction()
elif instr == INSTR_ADD:
    lhs = stack.pop()
    rhs = stack.pop()
    if isinstance(lhs, int) and isinstance(rhs, int):
        stack.push(lhs + rhs)
    else: ...
    program_counter += 1
```

Meta-tracing JITs

FL Interpreter

```
program_counter = 0; stack = []
vars = {...}
while True:
    jit_merge_point(program_counter)
    instr = load_instruction(program_counter)
    if instr == INSTR_VAR_GET:
        stack.push(
            vars[read_var_name_from_instruction()])
        program_counter += 1
    elif instr == INSTR_VAR_SET:
        vars[read_var_name_from_instruction()] =
            stack.pop()
        program_counter += 1
    elif instr == INSTR_INT:
        stack.push(read_int_from_instruction())
        program_counter += 1
    elif instr == INSTR_LESS_THAN:
        rhs = stack.pop()
        lhs = stack.pop()
        if isinstance(lhs, int) and isinstance(rhs, int):
            if lhs < rhs:
                stack.push(True)
            else:
                stack.push(False)
        else: ...
        program_counter += 1
```

Meta-tracing JITs

FL Interpreter

```
program_counter = 0; stack = []
vars = {...}
while True:
    jit_merge_point(program_counter)
    instr = load_instruction(program_counter)
    if instr == INSTR_VAR_GET:
        stack.push(
            vars[read_var_name_from_instruction()])
        program_counter += 1
    elif instr == INSTR_VAR_SET:
        vars[read_var_name_from_instruction()] =
            stack.pop()
        program_counter += 1
    elif instr == INSTR_INT:
        stack.push(read_int_from_instruction())
        program_counter += 1
    elif instr == INSTR_LESS_THAN:
        rhs = stack.pop()
        lhs = stack.pop()
        if isinstance(lhs, int) and isinstance(rhs, int):
            if lhs < rhs:
                stack.push(True)
            else:
                stack.push(False)
        else: ...
        program_counter += 1
```

User program (lang FL)

```
if x < 0:
    x = x + 1
else:
    x = x + 2
x = x + 3
```

Meta-tracing JITs

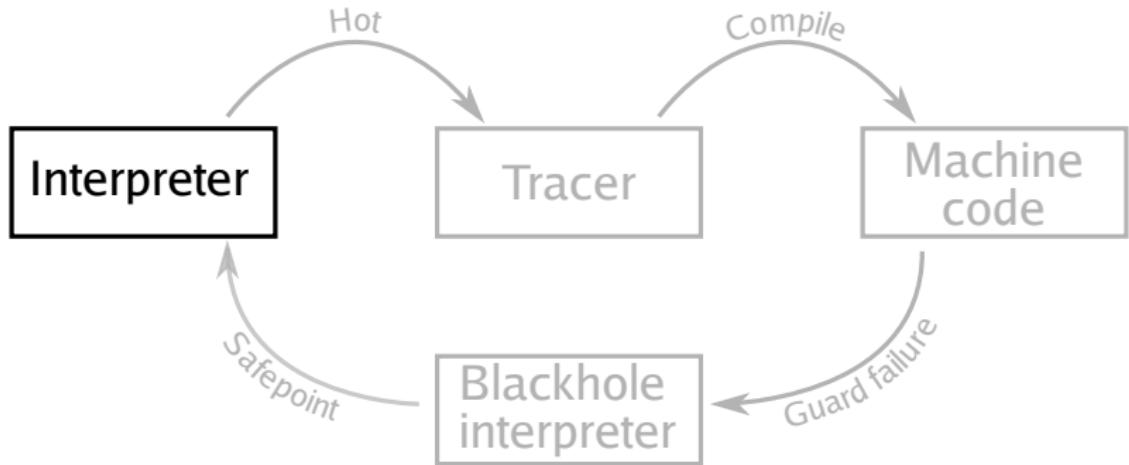
FL Interpreter

```
program_counter = 0; stack = []
vars = {...}
while True:
    jit_merge_point(program_counter)
    instr = load_instruction(program_counter)
    if instr == INSTR_VAR_GET:
        stack.push(
            vars[read_var_name_from_instruction()])
        program_counter += 1
    elif instr == INSTR_VAR_SET:
        vars[read_var_name_from_instruction()] =
            stack.pop()
        program_counter += 1
    elif instr == INSTR_INT:
        stack.push(read_int_from_instruction())
        program_counter += 1
    elif instr == INSTR_LESS_THAN:
        rhs = stack.pop()
        lhs = stack.pop()
        if isinstance(lhs, int) and isinstance(rhs, int):
            if lhs < rhs:
                stack.push(True)
            else:
                stack.push(False)
        else: ...
        program_counter += 1
```

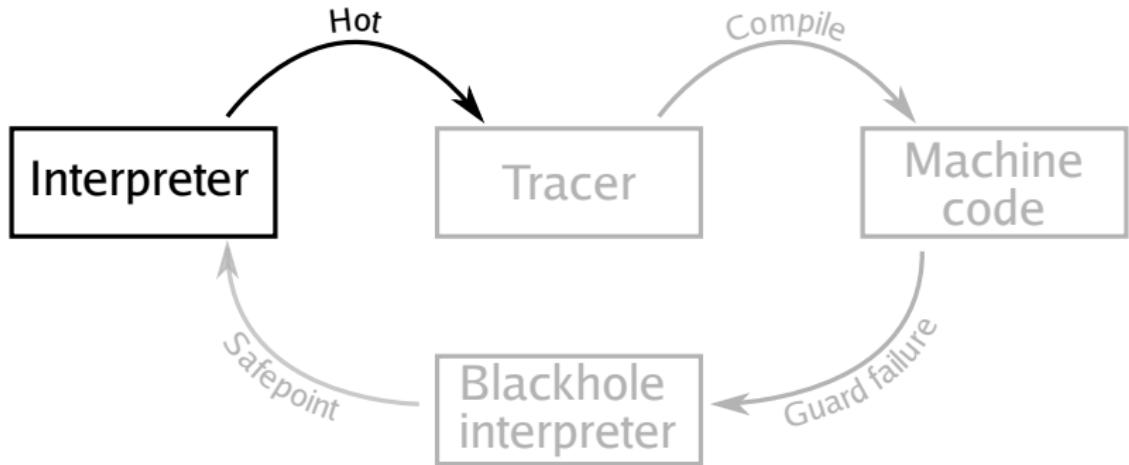
Initial trace

```
v0 = <program_counter>
v1 = <stack>
v2 = <vars>
v3 = load_instruction(v0)
guard_eq(v3, INSTR_VAR_GET)
v4 = dict_get(v2, "x")
list_append(v1, v4)
v5 = add(v0, 1)
v6 = load_instruction(v5)
guard_eq(v6, INSTR_INT)
list_append(v1, 0)
v7 = add(v5, 1)
v8 = load_instruction(v7)
guard_eq(v8, INSTR_LESS_THAN)
v9 = list_pop(v1)
v10 = list_pop(v1)
guard_type(v9, int)
guard_type(v10, int)
guard_not_less_than(v9, v10)
list_append(v1, False)
v11 = add(v7, 1)
v12 = load_instruction(v11)
guard_eq(v12, INSTR_IF)
v13 = list_pop(v1)
guard_false(v13)
...
```

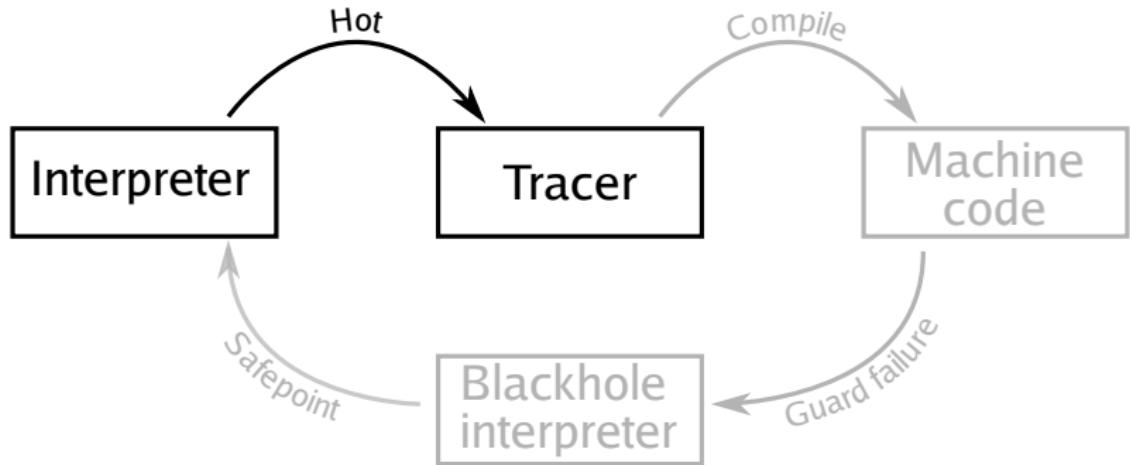
Meta-tracer states



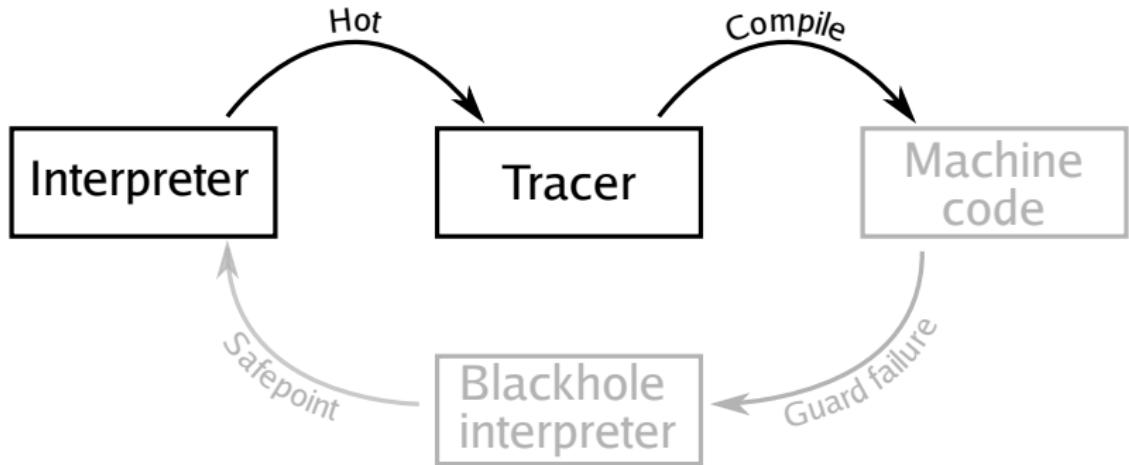
Meta-tracer states



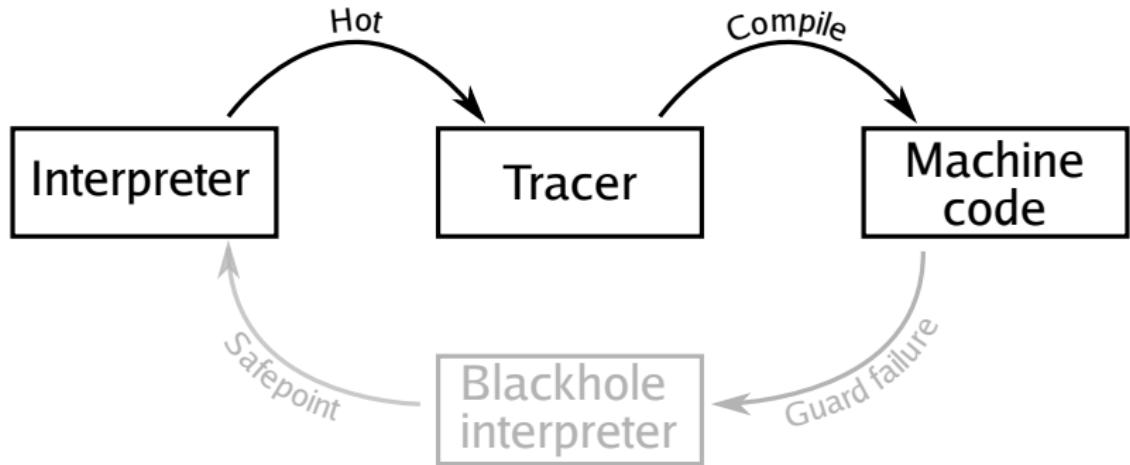
Meta-tracer states



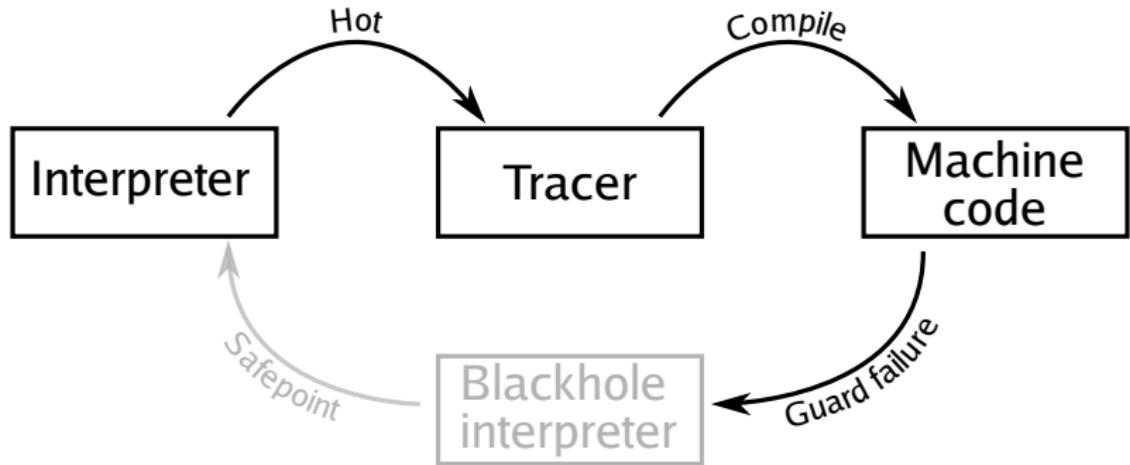
Meta-tracer states



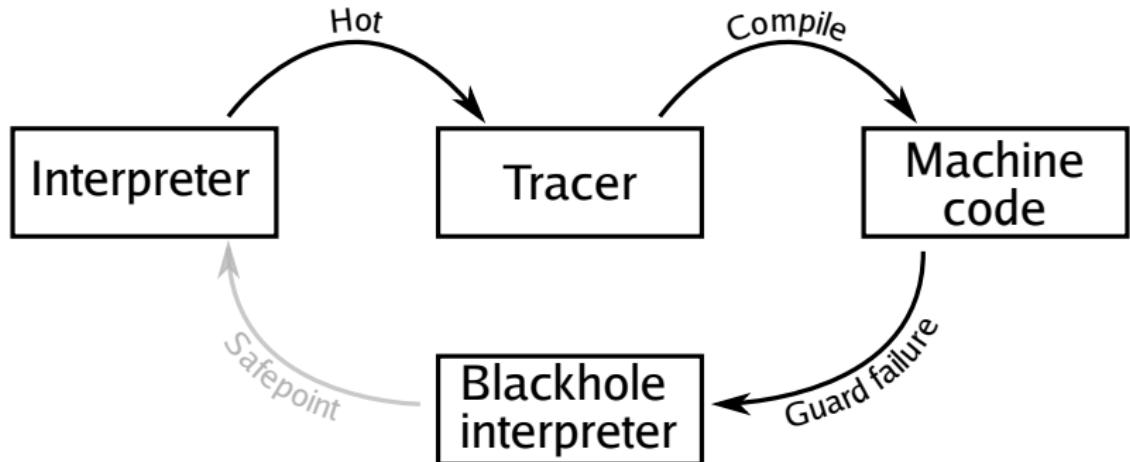
Meta-tracer states



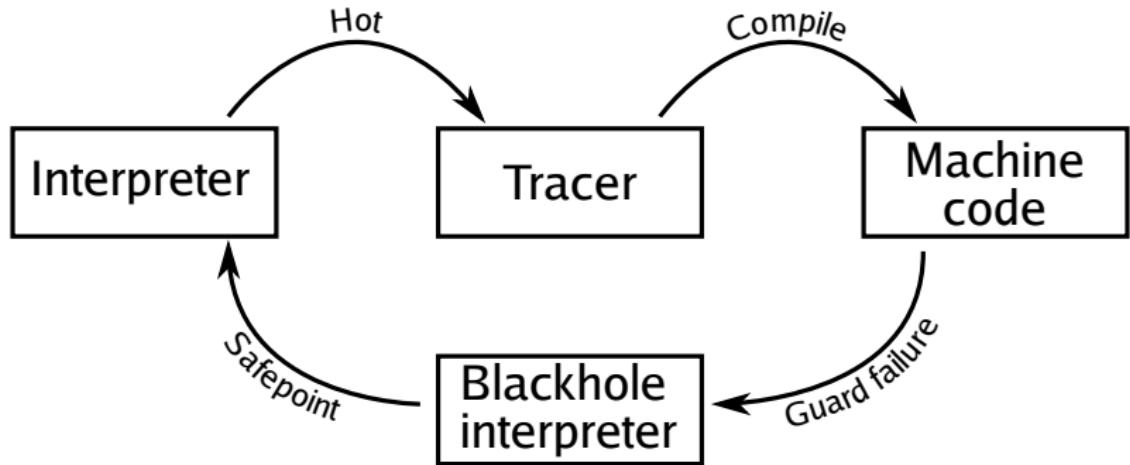
Meta-tracer states



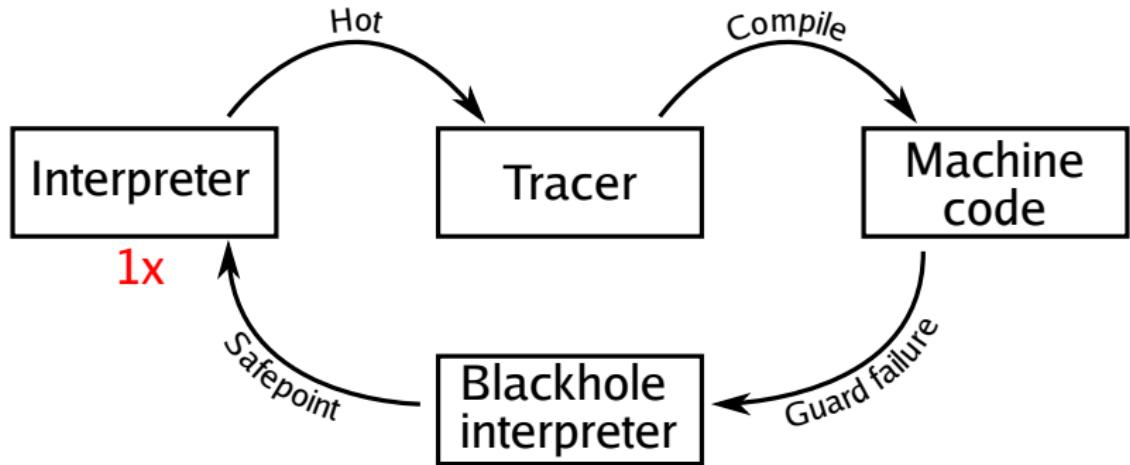
Meta-tracer states



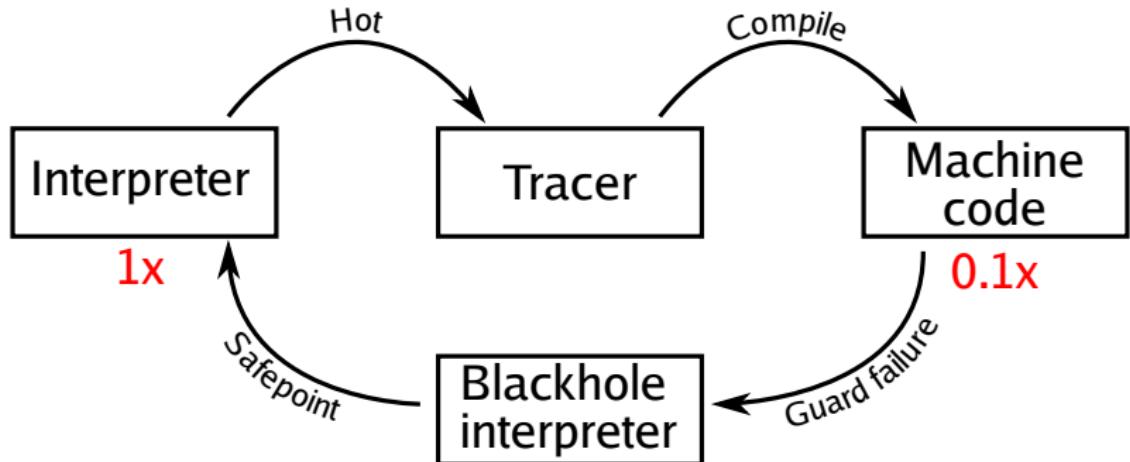
Meta-tracer states



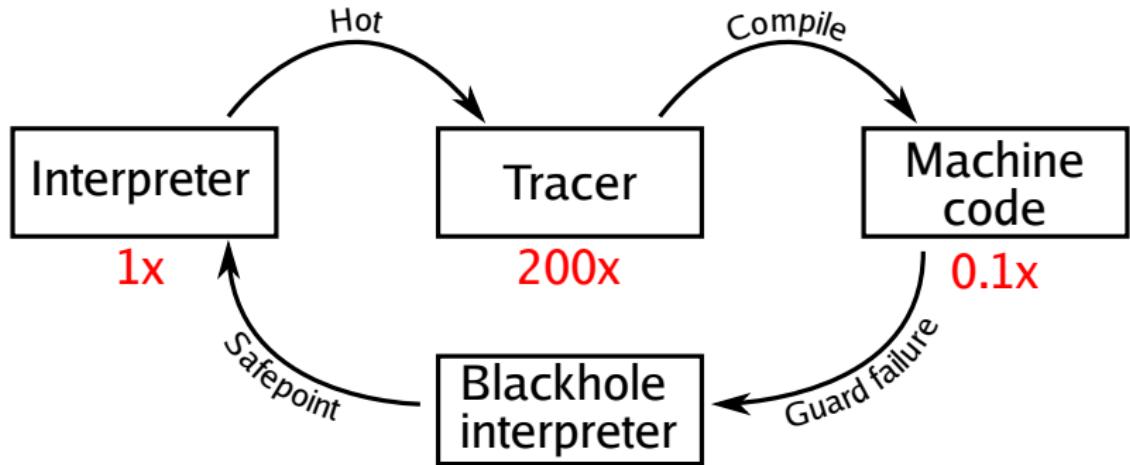
Meta-tracer performance (now)



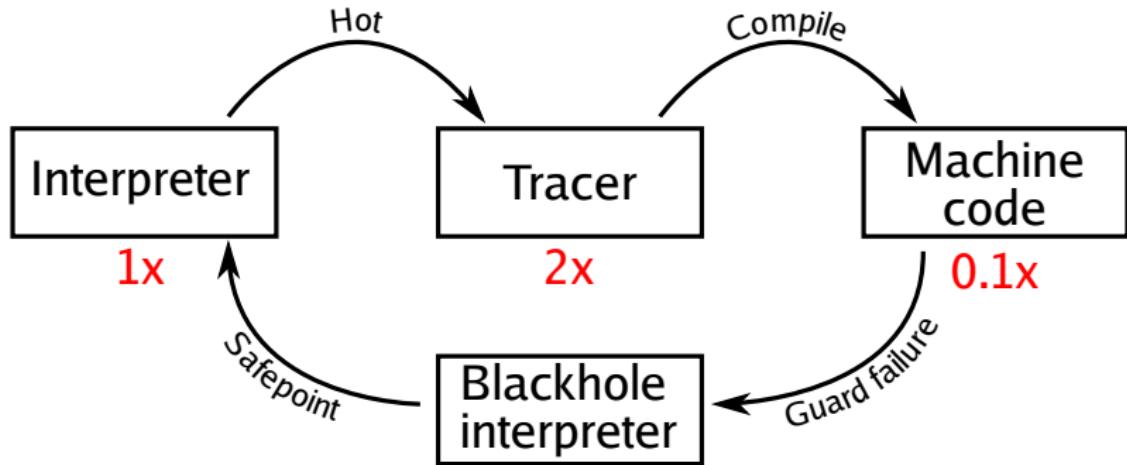
Meta-tracer performance (now)



Meta-tracer performance (now)



Meta-tracer performance (our aim)



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

Software

warmup_stats Use our statistical method on your VMs

http://soft-dev.org/src/warmup_stats/

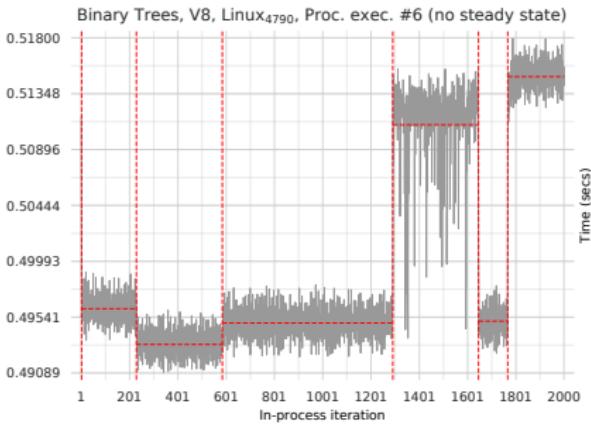
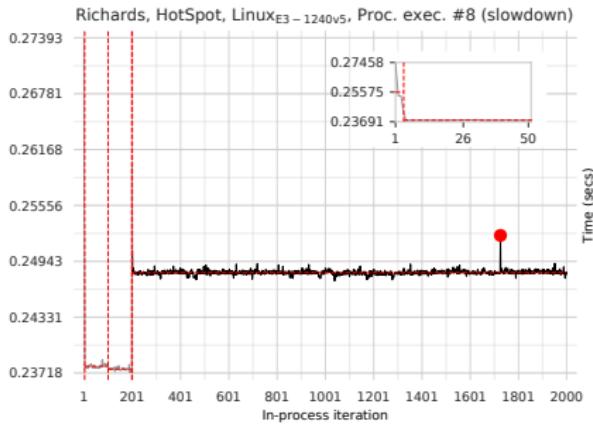
Krun Run experiments in a controlled manner

<http://soft-dev.org/src/krun/>

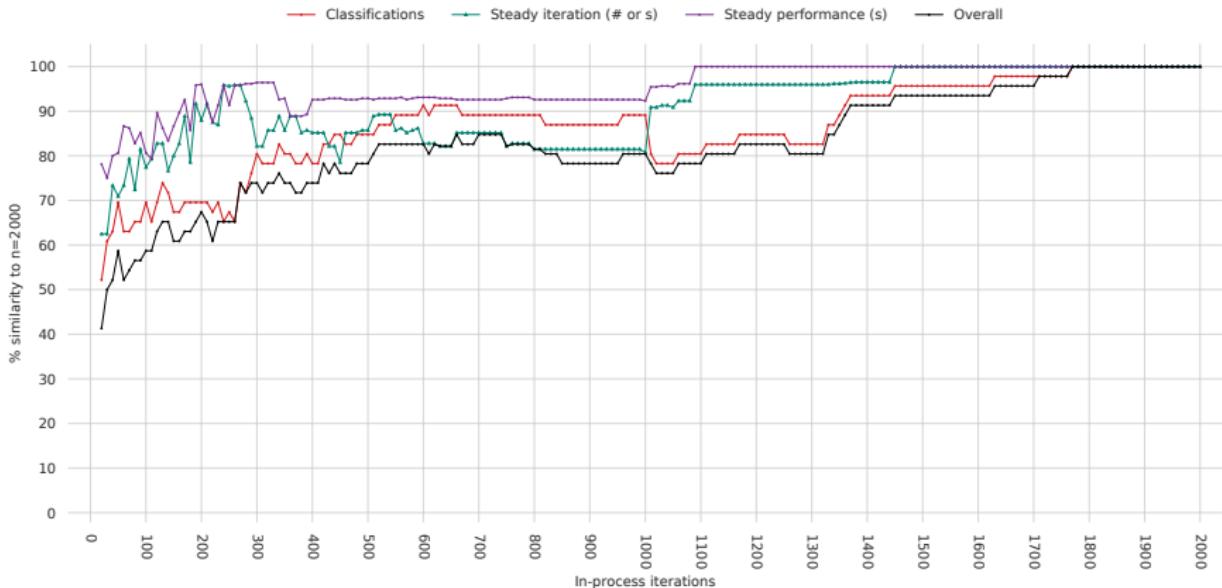
Thanks

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

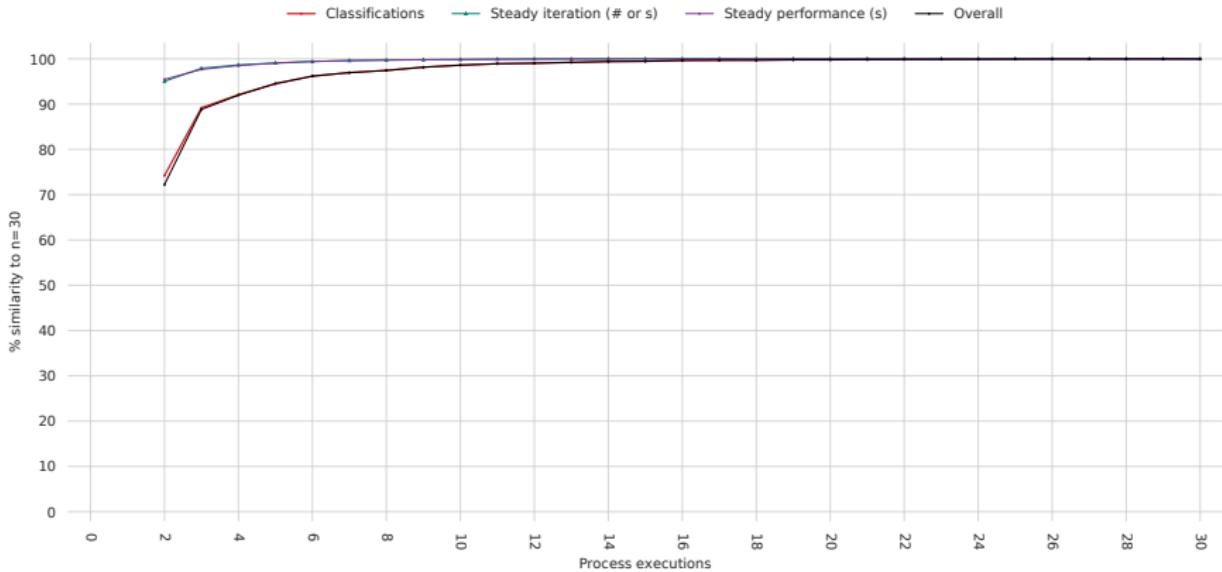
Thanks for listening



How long to run things for (0.8)



How long to run things for (0.8)



Diffing results ($0.8 \rightarrow 1.5$)

	Class.	Steady iter (#)	Steady iter (s)	Steady perf (s)
binarytrees		✗ (27L, 2F, 1w)		
fannkuch_redux		✗ (26L, 4w)		
fasta	L	4.0 (3.0,34.4)	0.75 (0.535,5.873)	0.16188 ± 0.000738
	L	6.0 $\delta = -2.0$ (5.0,7.0)	0.86 $\delta = -0.352$ (0.704,1.090)	0.13677 $\delta = +0.00343$ ± 0.000035
nbody	Graal			
richards	L	2.0 (2.0,35.3)	0.95 (0.879,9.745)	0.26465 ± 0.007761
spectralnorm	F	14.0 (2.0,94.6)	13.60 (0.830,98.737)	1.05685 $\delta = +0.16392$ ± 0.000126