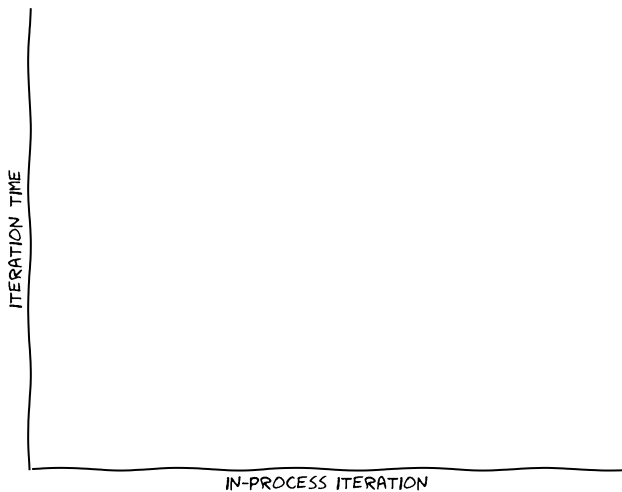# VM Warmup Blows Hot and Cold



Edd Barrett, Carl Friedrich Bolz, Rebecca Killick (Lancaster),
Vincent Knight (Cardiff), Sarah Mount, Laurence Tratt
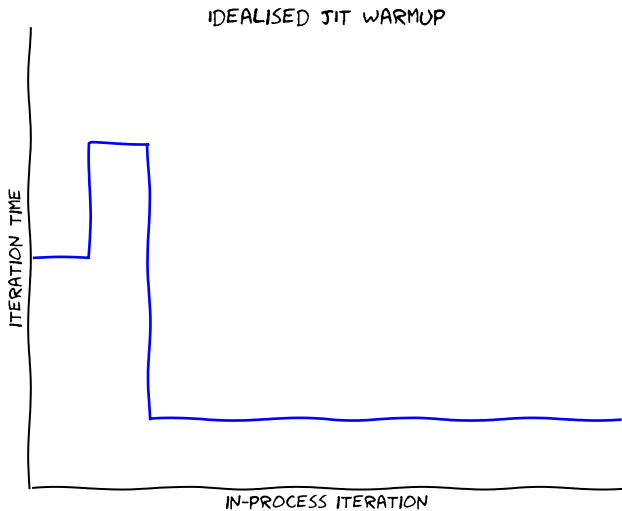

KING'S College LONDON

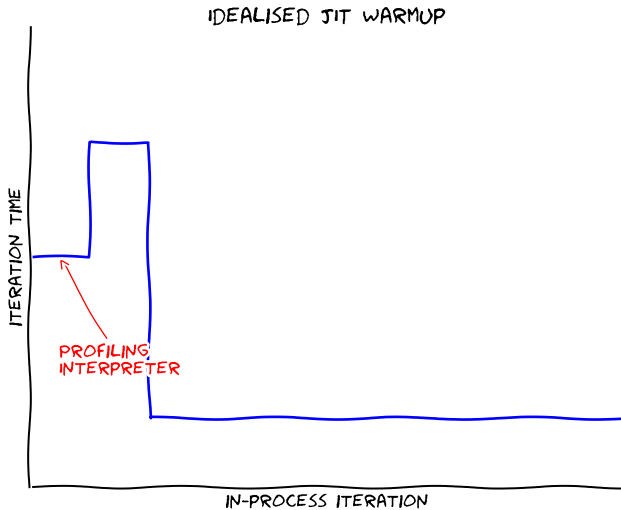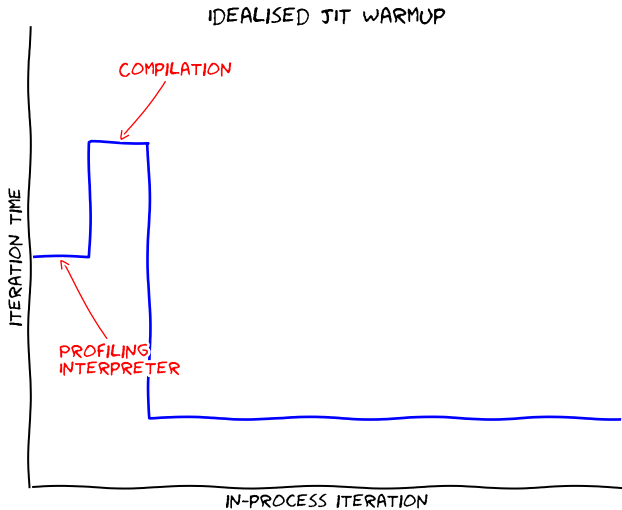Software Development Team
2016-09-27

IDEALISED JIT WARMUP

IDEALISED JIT WARMUP

ITERATION TIME

PROFILING
INTERPRETER

IN-PROCESS ITERATION

IDEALISED JIT WARMUP

COMPILATION

ITERATION TIME

PROFILING INTERPRETER

IN-PROCESS ITERATION

IDEALISED JIT WARMUP

COMPILATION

ITERATION TIME

PROFILING INTERPRETER

PEAK PERFORMANCE

IN-PROCESS ITERATION

IDEALISED JIT WARMUP

MORE REALISTIC VM WARMUP

ITERATION TIME

IN-PROCESS ITERATION

MORE REALISTIC VM WARMUP

ITERATION TIME

SOME NOISE

IN-PROCESS ITERATION

MORE REALISTIC VM WARMUP

GC SPIKES

COMPILATION TIERS

SOME NOISE

ITERATION TIME

IN-PROCESS ITERATION

# Why care?

USERS HATE NOTICABLE WARMUP

ITERATION TIME

FRUSTRATING

HAPPY DAYS!

IN-PROCESS ITERATION

VM AUTHORS HATE ALL WARMUP

Warmup is bad for everyone.

Measure warmup of modern language
implementations

Measure warmup of modern language
implementations

*Hypothesis*: Small, deterministic programs exhibit
classical warmup behaviour

The language benchmark games are perfect for us
(unusually)

The language benchmark games are perfect for us
(unusually)

We removed any CFG non-determinism

The language benchmark games are perfect for us (unusually)

We removed any CFG non-determinism

We added checksums to all benchmarks

2000 *in-process iterations*

2000 *in-process iterations*

10 *process executions*

- Graal-0.13
- HHVM-3.12.0
- JRuby/Truffle (git #f82ac771)
- Hotspot-8u72b15
- LuaJit-2.0.4
- PyPy-4.0.1
- V8-4.9.385.21
- GCC-4.9.3

Note: same GCC (4.9) used for all compilation

- Linux-Debian8/i4790K, 24GiB RAM
- Linux-Debian8/i4790, 32GiB RAM
- OpenBSD-5.8/i4790, 32GiB RAM

- Linux-Debian8/i4790K, 24GiB RAM
- Linux-Debian8/i4790, 32GiB RAM
- OpenBSD-5.8/i4790, 32GiB RAM


- Turbo boost and hyper-threading disabled
- SSH blocked from non-local machines
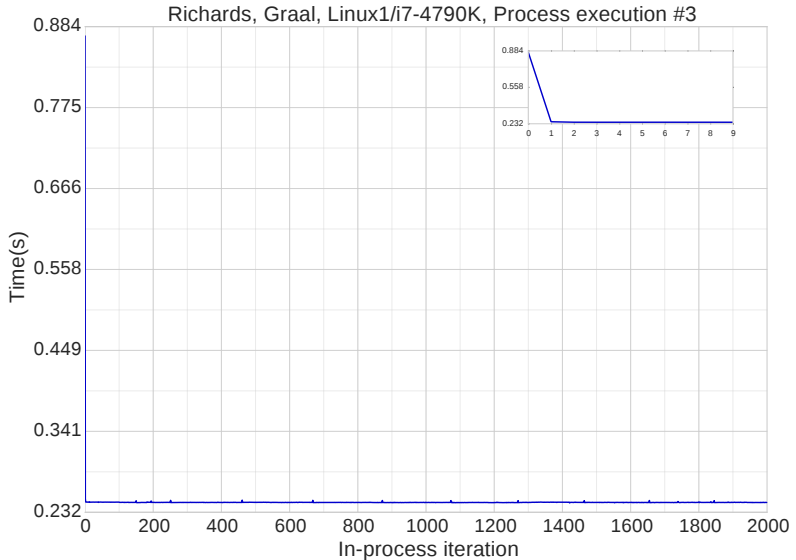- Daemons disabled (cron, smtpd)

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

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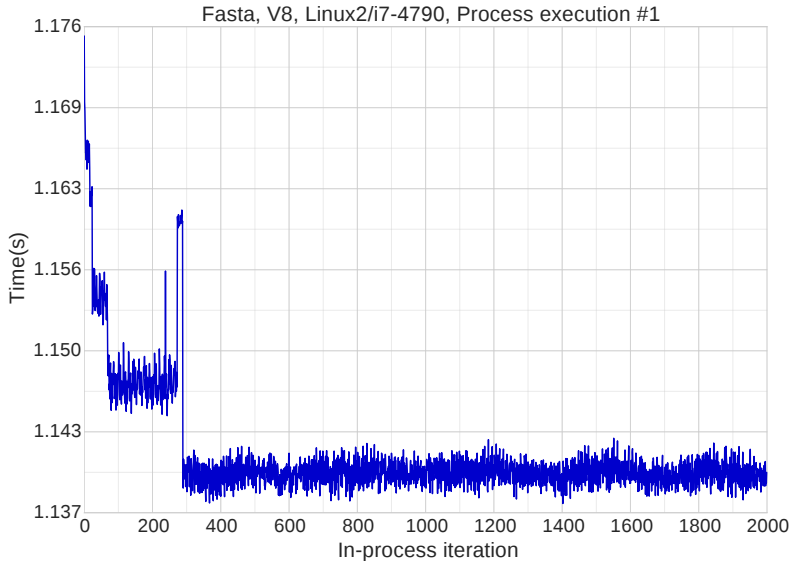confounding variables as possible e.g.:

- 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
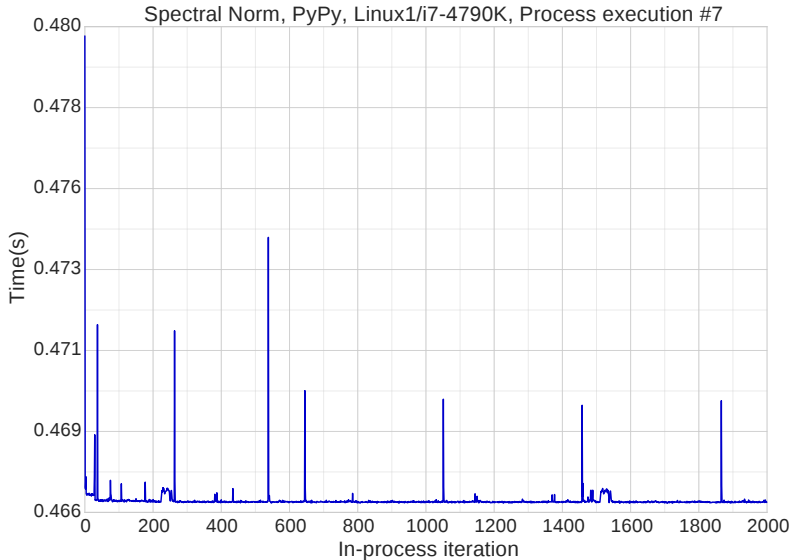- Enforces kernel settings (tickless mode, CPU governors, ...)
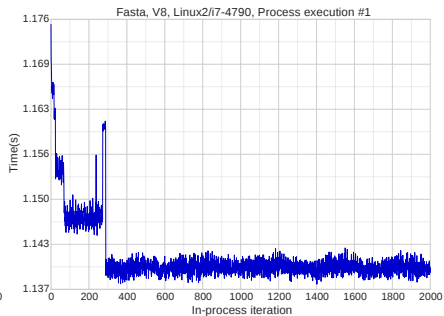
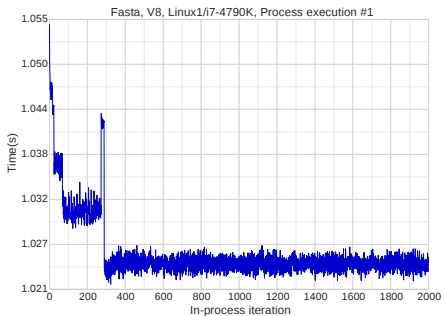# Preliminary results

# Classical Warmup



Richards, Graal, Linux1/i7-4790K, Process execution #3

Fasta, V8, Linux2/i7-4790, Process execution #1

# Classical Warmup



Spectral Norm, PyPy, Linux1/i7-4790K, Process execution #7

(Different machines)

Fannkuch Redux, LuaJIT, OpenBSD/i7-4790, Process execution #10

Richards, Hotspot, Linux2/i7-4790, Process execution #2

Fannkuch Redux, Hotspot, Linux1/i7-4790K, Process execution #1

Fannkuch Redux, Hotspot, OpenBSD/i7-4790, Process execution #4

# Cycles

# Never-ending Phase Changes



Fasta, LuaJIT, OpenBSD/i7-4790, Process execution #5

(Note: same machine)

(Note: different machines. Bouncing ball pattern Linux-specific)

Classical warmup occurs for only:

Classical warmup occurs for only:

50% of process executions

Classical warmup occurs for only:

50% of process executions

25% of (VM, benchmark) pairs

Classical warmup occurs for only:

50% of process executions

25% of (VM, benchmark) pairs

0% of benchmarks for all VMs, machines &
proc execs.

*Hypothesis*: ~~Small, deterministic programs exhibit classical warmup behaviour~~

How can we measure anything any more?

How can we measure anything any more?

For how long has this been going on?

How can we measure anything any more?

For how long has this been going on?

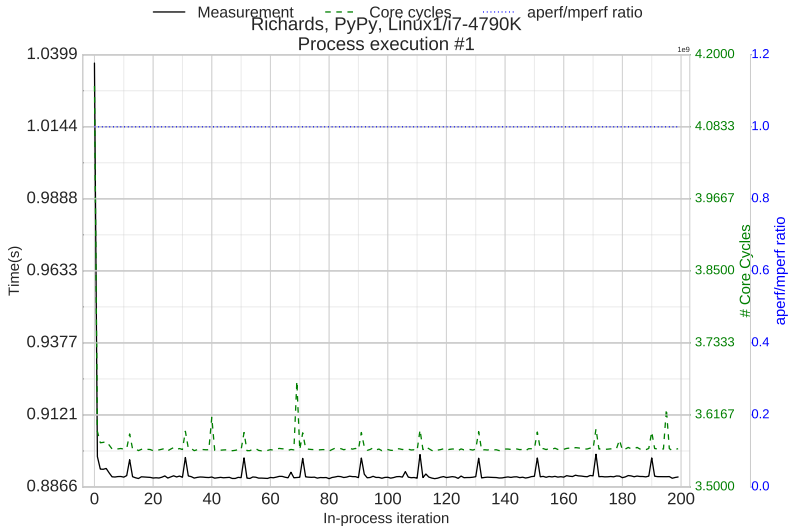Is this really the fault of the VMs?

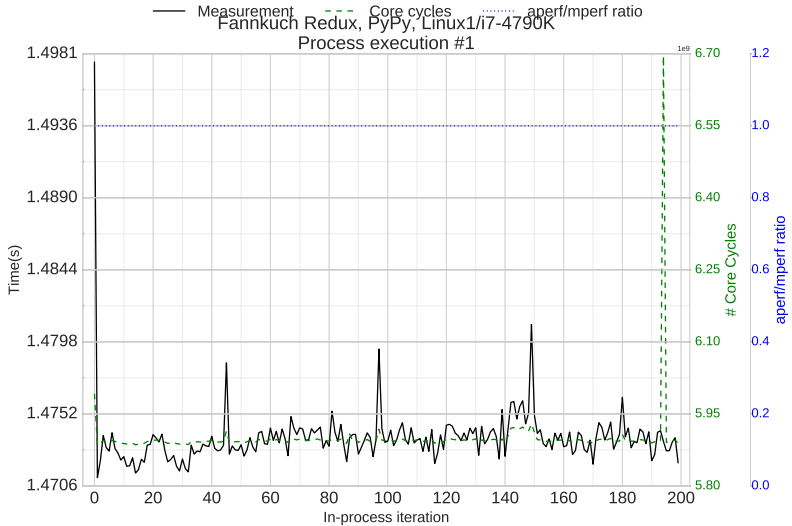# Ongoing/Future Work

(a.k.a. "making sense of our results")

- `CPU_CLK_UNHALTED.CORE`
  - Counts the number of "core cycles" executed per-core.

- `IA32_APERF` / `IA32_MPERF` ratio.
  - `IA32_MPERF` increments at a fixed reference frequency.
  - `IA32_APERF` increments proportional to "actual performance".
  - The ratio of two delta's indicates if clock speed changed.

Richards, PyPy, Linux1/i7-4790K
Process execution #1

Fannkuch Redux, PyPy, Linux1/i7-4790K
Process execution #1
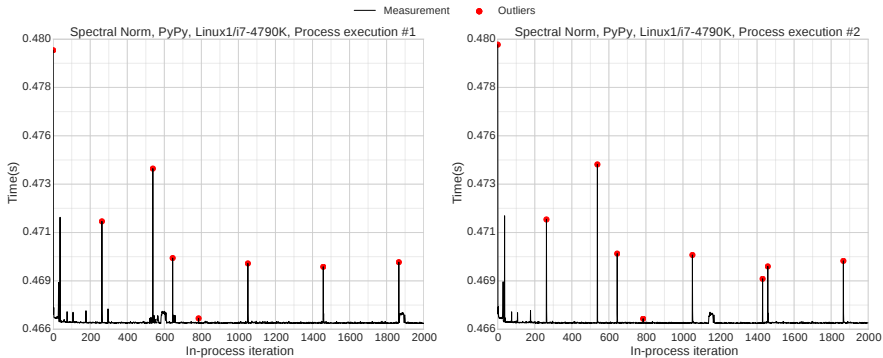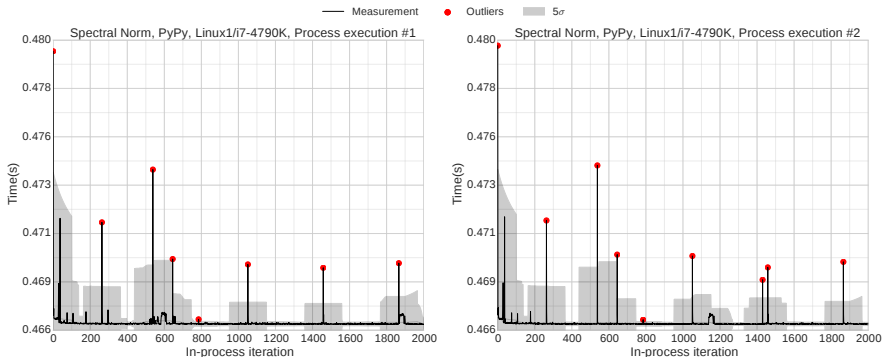
For PyPy and Hotspot, record:

- Time spent in GC.
- Time spent in Compilation.

# Outlier Detection

# Outlier Detection



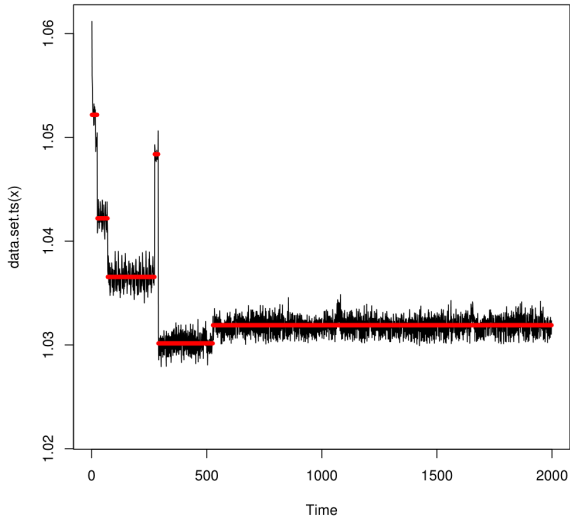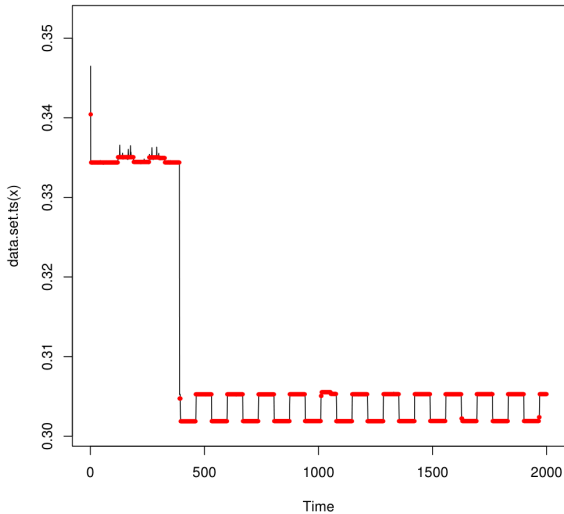Outliers outside $5\sigma$ of rolling average

Recurring outliers

fasta:V8:default−javascript , run: 5

fannkuch_redux:Hotspot:default–java , run: 1

binarytrees:PyPy:default–python , run: 1

```
https://archive.org/download/softdev_warmup_
          experiment_artefacts/v0.2/
```

- `all_graphs.pdf` All plots in one huge PDF.
- `warmup_results*.json.bz2` Raw results.

(Note: newer results available)

**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



Fannkuch Redux, Hotspot, Linux1/i7-4790K, Process execution #1



Fasta, LuaJIT, OpenBSD/i7-4790, Process execution #5