

How (NOT) To Benchmark

0. "I don't need to benchmark,
my algorithm is $O(n)$."
1. Use unrepresentative benchmarks.
(Blackburn et al. 2006)
2. Ignore basic statistics.
(Georges et al. 2007)
3. Measure the wrong thing.
(Mytkowicz et al. 2009)
4. Make faulty assumptions.
(Barrett et al. 2017)

O. Evaluating performance

2

defn. A function $f(n)$ is $O(g(n))$ if there exists constants C, n_0 such that for all $n \geq n_0$ then $f(n) \leq C \cdot g(n)$.

- From Week 5 of the second CS course. A full 15 weeks of this in "Baby algorithms."
- Sufficient? No way!
Necessary? Debatable...
- Architecture independence is seen as an advantage, But it can also lead to practical irrelevance.
- Constant factors matter.
Architecture matters.

e.g. Persistent vectors

2-3 finger Trees

$O(1)$	$O(\log n)$	$O(n)$
append	insert	concat
prepend	nth	

Bitmapped Vector Trie

$O(1)$	$O(\log n)$	$O(n)$
	append	concat
	nth	insert
		prepend

- In practice Bitmapped vector trie is faster. Why?
 - $O(\log_{32} n)$ effectively constant
 - Locality of Reference

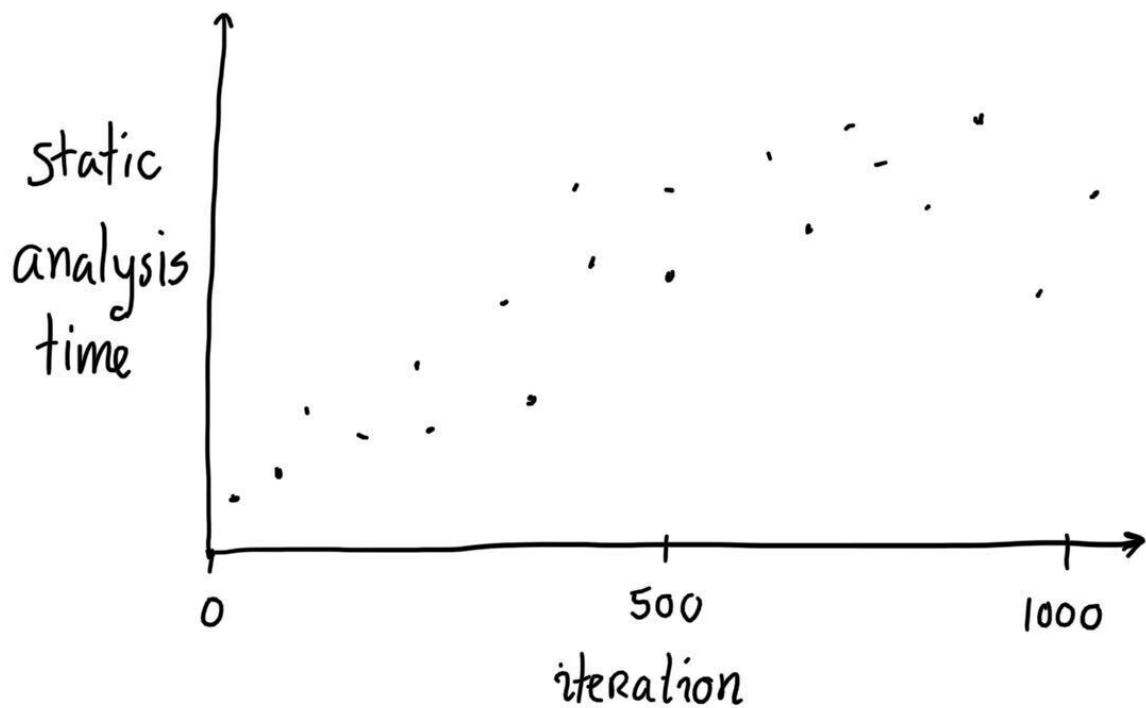
e.g. GTA Online

4

- Recently (Feb. 28, 2021) someone looked into why GTA Online is so slow.
- $\geq 81\%$ of players had a loading time of ≥ 3 minutes (for a 7 year old game on today's hardware)
- PROBLEM: Using `sscanf` to read tokens when parsing JSON (it calls `strlens`).
Result: 6min load \rightsquigarrow 2 min.
- Any developer who knows how to use a profiler could have found this in an afternoon.
- They didn't. I'm not surprised.

e.g. Long-running experiments. 5

- I'm not surprised because I've made this mistake too.
- Had experiment which took > 1 month to run



- Memory leak. Analysis retained information (and worked with it) on all subsequent iterations
- ~4 line fix \Rightarrow finishes in 2 weeks
- "anchoring effects"

O. Evaluating Performance.

6

- Language and library implementers also need to care about performance

Step 1: Invent new _____.

(data structure, algorithm,
compiler opt., runtime opt, etc.)

Step 2: Create or acquire some benchmark programs.

Step 3: Run them with and without improvement and measure _____.

(execution time, memory use, etc.)

Step 4: new < old? Declare victory!
Write a PLDI paper.

- What could possibly go wrong?

1. Bad benchmarks.

7

- In 2003, the DaCapo research group was unsatisfied with existing Java Benchmarks
- Standard suites: SPECjvm98, SPECjbb2000
 - Don't adequately make use of Java's features
 - Low memory footprint and overprovisioned heap (doesn't stress GC)
- New goals:
 - Diverse and relevant workloads.
 - Easy to measure and use for research.
- Developed a new suite meeting goals

1a. Benchmark selection

8

- Eleven reasonably large real-world actively maintained Java projects (e.g. antlr, eclipse)
- Evaluation on many qualitative and quantitative metrics
 - CK metrics for OO code complexity (methods per class, inheritance hierarchy depth)
 - Heap metrics (maximum objects allocated, nursery survival rate)
 - Architecture metrics (mix of different instructions)
- Use principle component analysis (PCA) to assess diversity.

16. Benchmark packaging

9

- Benchmarks should be easy to use and come with infrastructure
- Fixed, manageable workloads (feasible to run many iterations and setups)
- Harness to automatically launch, time, and test correctness
- Configuration options for input size, iterations, warmup, hardware tuning
- Remove need for multiple hosts and third-party components (e.g. databases)
- This is important but thankless work.
- “The [NSF] panel suggested we solve our own problems ... but didn’t provide additional funds.”

e.g. GTP LNM.

10

untyped lmm-plot.rkt.

```
(plot-pict
  (function
    (λ (N)
      ... )))
```

typed lmm-plot.rkt.

```
(plot-pict
  (function
    (λ ([N-Raw : Real])
      (: N Nonnegative-Real)
      (define N
        (if (≥ N 0) N ...)))
      ... )))
```

this is never run!

↑
↑
typo (N-Raw)

- Got lucky, didn't impact performance.

1. JavaScript too

11

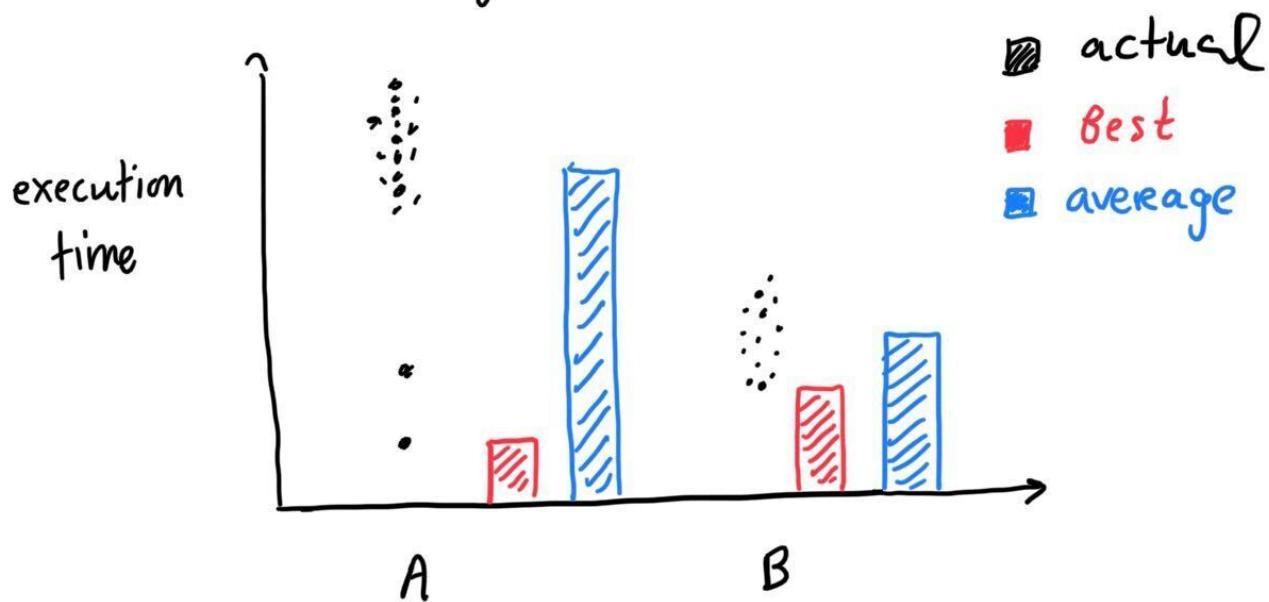
- Familiar story from Richards et al. in 2011, five years after DaCapo
- SunSpider suite has small benchmarks. V8 too (600–5000 LOC). Unrepresentative of real app
- Doesn't make use of dynamism found in wild. Totally different memory usage patterns. (Familiar?)
- Developed new benchmarks and show Firefox's 13x improvement on SunSpider was really 3x.
- Summary: Poor benchmarks can yield wasted effort chasing fruitless optimizations
- Where did new benchmarks come from?

1. Generating Benchmarks.

12

- A realistic corpus is difficult for JavaScript
 - Ecosystem rapidly evolving
 - Applications are interactive
- Wrote a tool for creating a Benchmark from any application via "record-and-replay"
 - Deterministic Replay
 - Browser agnostic
 - Accurately mimic behavior and performance characteristics of the original
- Instrument JavaScript and output trace. Reconstruct program and memoize non-deterministic events.

- Georges et al. (2007) argue for statistically rigorous data analysis
SPEC recommends
- What do you report? Best? Second Best? Average? Worst?



- Lots of sources of non-determinism (JIT, thread scheduling, GC, OS, ...). Need to account for Random variations
- Solution: Use statistics. 95% CI of mean. T-test. ANOVA. etc...
Many standard techniques.

2. How big is the problem?

14

- Q: Is runtime execution variability really that big?
- A: Startup: 8%, "Steady-State": 20% difference between max and min
- Q: Do nonrigorous data analyses lead to incorrect conclusions?
- A: Startup: 16%, "Steady-State": 20% prevalent data analysis methodology comes to "misleading" conclusion
- Summary: Yes it's a problem, a statistical rigor is the solution. Only one issue...
- STATISTICS IS HARD!

2. Easy you say? Ha.

15

From Georges et al.

"A 90% confidence interval... means that there is a 90% probability that the actual distribution mean of the underlying population, μ , is within the confidence interval."

No. See "Robust misinterpretation of confidence intervals." Hoekstra et al. 2014

"If they do overlap then we cannot conclude that the differences seen in the mean values are not due to random fluctuations in the measurements."

No. See "A Brief note on overlapping confidence intervals." Austin et al. 2002

2. Critiquing the critique.

16

- Those numbers I just presented could be totally wrong (it's unclear from their description.)

"To consult a statistician after an experiment is finished is often merely to ask him to conduct a post mortem examination."

— R.A. Fisher

- Shouldn't this be someone's job?

- Takeaway:

If you are performing a nontrivial empirical evaluation, ask someone with statistical experience for help.

3. Measuring the wrong thing.

17

- A basic experiment might compare the execution of a system without optimization (S) and with it (S').
- Each execution runs in some experimental context: OS, processor, compiler flags, ...
- Q: How sensitive are measurements to experimental setup?
- A: Often, very sensitive.
- Consequently, a seemingly innocuous change in experimental setup can yield a different result
- These are confounding variables and can cause omitted-variable bias. Systematic error, not random. Statistics can't help you.

3. Subtle confounders.

18

- Mytkowicz (2009) demonstrates two subtle confounders in experimental setup.
 1. Unix environment size
 2. Linking order
- They are commonplace (i.e. affect different platforms) and significant (i.e. can lead to incorrect conclusion).
- Q: How effective is O3 optimization varying these two factors?
- A: Link order: 0.92 - 1.10
Env size: 0.91 - 1.07
- Consequently, conclusion about whether O3 optimizes can flip depending on either of these.

3. Solutions.

- Obvious solution: Don't change anything between S and S'
- Can be easier said than done if you're unaware of possible confounders.
- Many possible confounders even beyond experimental setup
- Benchmark diversity can help cancel some confounding effects
- SPEC CPU 2006 suite was NOT diverse enough to mitigate effects due to experimental setup
- Result is you may end up measuring the wrong thing.

e.g. GTP ZondoZ.

20

untyped zo-find.rkt.

```
(define-values (title children)
  (parse-zo z))
```

...

typed zo-find.rkt.

```
(define-values (title children)
  (parse-zo z))
```

```
(define zstr
  (format "%a" z))
```

...
↑ from untyped
code

} unused with
no equivalent
in untyped

- Were NOT lucky. Worst case overhead goes from 14x to <4x slowdown.
- One line change can cause an order of magnitude difference

e.g. Branch ordering -

slow.rkt.

(match

$$\left[\begin{array}{l} ((R(B) (= v) (B)) \dots) \\ ((B(B) (= v) (B)) \dots) \\ \dots \end{array} \right] \quad \left. \begin{array}{l} \text{base} \\ \text{cases} \end{array} \right\}$$

$[(N a x b)]$

(cond

$$\left[\begin{array}{l} (< x v) \dots \\ (= x v) \dots \\ (> x v) \dots \end{array} \right])])$$

↑
↑
↑

hot
branches

better.rkt.

(match

$[(N a x b)]$

(cond

$$\left[\begin{array}{l} (< x v) \dots \\ (= x v) \quad \text{— Base cases —} \\ (> x v) \dots \end{array} \right]))$$

- A 2x slowdown disappears !

3. Is that good enough?

22

- Perfectly controlled confounders.
- Is that enough? Nope!
- Internal validity: Does my experiment establish a causal link?
- External validity: Does my conclusion generalize to a broad population?
- Controlling confounders improves internal validity ONLY.
- Additionally, randomize experimental setup. Can improve both.
- Also easier said than done. Requires more time and resources. Produces more data to analyze.

4. Making false assumptions.

23

"Later iterations usually have less compilation... and may reach a steady state."

— DaCapo team in CACM

"Since most of the JIT compilation is performed during startup, steady state performance suffers from less variability."

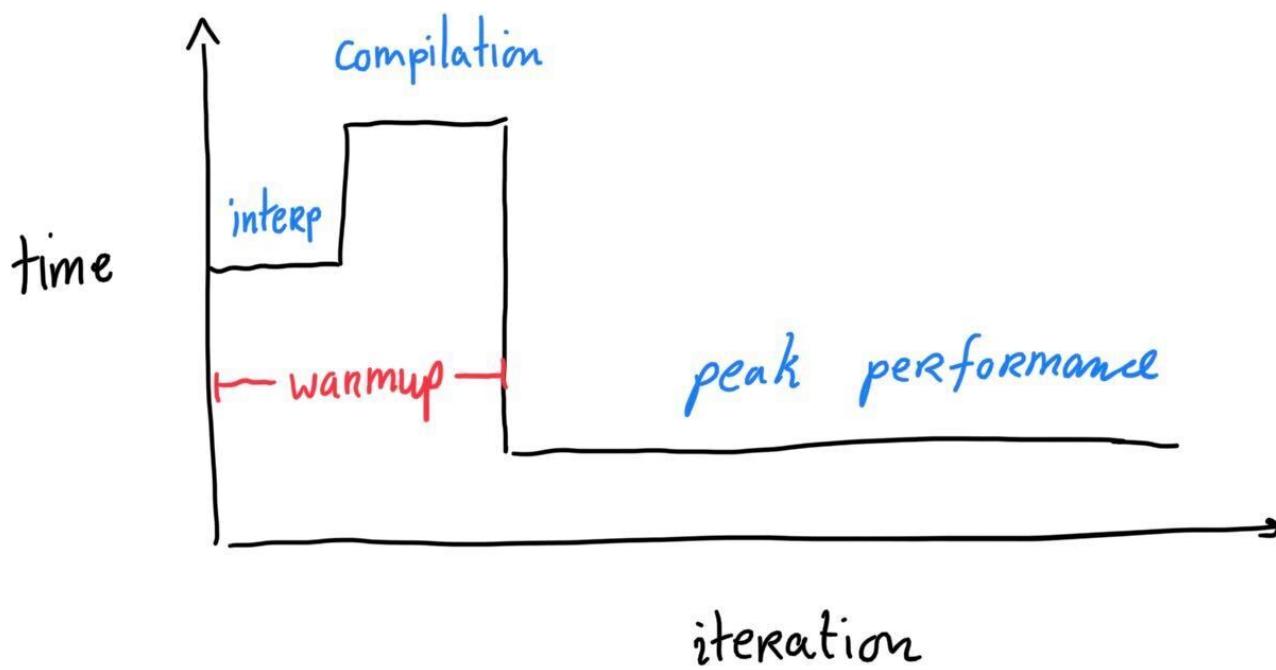
— Georges et al. (2007)

- Ad hoc approach: Drop the first 2 or 3 Benchmark iterations.
- More Rigorous: Fix desired iterations to keep K . Once coefficient of variation of last K iterations is below a threshold α , then take those.
Still somewhat ad hoc

4. How long is warmup?

24

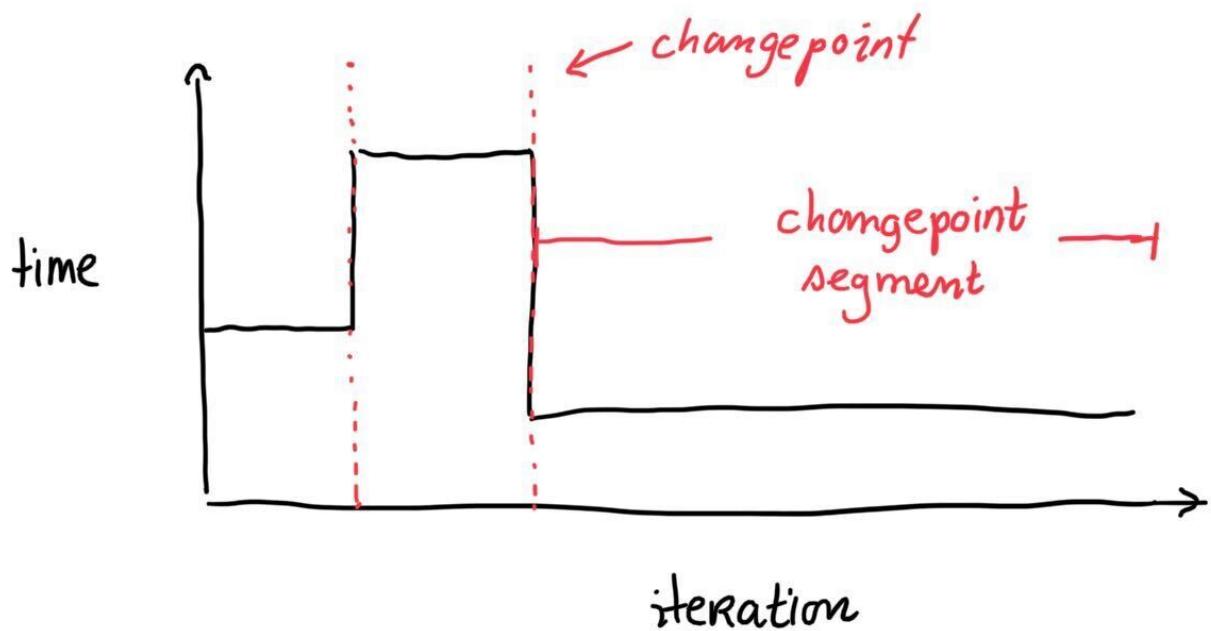
- Expected performance given knowledge of JIT compilation:



- Q: How long is the **warmup** period?
- Methodology: Take small benchmarks (more predictable), remove as many confounders as possible, several VMs, and measure performance across 2000 iterations

4. Data analysis.

25



Automatic performance classifier

- Segments equivalent : FLAT } good
- Final fastest : WARMUP }
- Final not fastest : SLOWDOWN }
- Never settles: NO STEADY STATE } bad

Good corresponds to classical expectations of JIT performance.

Bad corresponds to abnormal performance characteristics.

- For each VM - Benchmark pair only 40% consistently gave expected performance characteristics when repeating experiment

4.5% consistent slowdown

4.5% consistent no steady state

51% mix containing either slowdown or no steady state

- Takeaways:

- Your methodology cannot assume peak performance OR even a steady state will occur

- Even in carefully controlled experiments, performance can be unpredictable

5. Stepping Back

27

- Why is performance evaluation so difficult?
- Many answers to this. A few:
 - Satisfies 2/3 criteria for catastrophic failure potential
(Perrow 1984): complexity,
tight coupling.
 - Mistakes are silent, no typechecker for performance
 - Little academic incentive for going above and beyond in an evaluation
 - Too much folklore, not enough fact: "GC is bloated and slow"
(Zorn 1992)
 - Proper evaluation skills are not taught.

5. Broader context

28

- This isn't a problem that only affects systems or PL
- Same thing in other fields with a similar timeline
- Replication crisis: Starting with "Why Most Published Research Findings are Wrong" (2005)
 - Standard methodological practices have lead to incorrect results becoming commonplace
 - p hacking, publication bias, underpowered experiments
- Especially as PL adopts more empirical methods, this will be more prevalent

5. A note of optimism.

29

- Not too long ago papers often didn't come with associated data or software
- Or maybe the link to it died. Or it wasn't portable. Or it just didn't work.
- Now, peer-reviewed artifacts available at permanent storage links is normal
- Other disciplines are beginning to make progress too, e.g. preregistration of clinical trial to prevent p-hacking
- As we've seen, more work to be done though

Tools

Krum and ReBench for eliminating confounders due to operating system and hardware parameters.

Articles

- R3 - Repeatability, Reproducibility, Rigor (Vitek, and Kalibera 2011)
- Slides Empirical Evaluation Checklist

Books

- Experimental Methods for Science and Engineering Students (Kirkup)
- Experimental Design and Analysis (Seltman, free online!)
- Statistics Done Wrong (Reinhart)

CONCLUSION

31

DON'T Use small contrived
Benchmark programs

DO Use large real-world
programs solving a diverse
set of problems

DON'T Just report best or worst
performance

DO Use statistical techniques
as appropriate

DON'T Introduce experimental confounders

DO Eliminate as many as possible

DON'T Assume classical warmup
always happens

DO Visualize long-running Benchmarks
to double-check your assumptions

Most importantly DO evaluate
performance. We DON'T need more
folklore!