Benchmarking Computer Systems

Tomas Kalibera



My benchmarking background

Systems

CORBA, Mono, RT Java (Ovm, RTS, WebSphere), Java (OpenJDK, Ovm, JikesRVM), R (FastR, GNU-R)

Benchmarks

SciMark, CORBA micro, CSiBE, FFT, SPEC JVM, SPEC JBB, SPEC CPU, DaCapo, Shootout, AT&T, Collision Det.

Methodology

Planning multi-level experiments with non-determinism, impact of code-layout, ratio estimation, concurrency and memory sharing metrics

Teaching, PhD

Most PL/systems papers include execution time measurements

Papers from 2011

	Papers	Measured time
ISMM	13	12
PLDI	55	42
ASPLOS	32	25
TOPLAS	13	5
TACO	9	6
Total	122	90

Kalibera, Jones: Quantifying performance changes with effect size confidence intervals. http://www.cs.kent.ac.uk/pubs/2012/3233/

Practices (not?) to follow

In 2006, Potti and Nevins at Duke reported they could predict lung cancer using expression arrays, and started a company

By 2010, 3 major papers were retracted, the company was gone, Potti resigned, and a major investigation was underway

Due to a combination of bad science ranging from fraud, unsound methods, to off-by-one errors in Excel spreadsheets

Uncovered by a repetition study conducted by statisticians Baggerly and Coombes with access to raw data and 2,000 hours of effort



Experimental results in PL/systems papers ignore elementary statistics

Papers from 2011

	Papers	Measured time	No error estir (ignored unce	nate ertainty)
ISMM	13	12	5	42%
PLDI	55	42	39	93%
ASPLOS	32	25	18	72%
TOPLAS	13	5	4	80%
TACO	9	6	5	83%
Total	122	90	71	79 %

Kalibera, Jones: Quantifying performance changes with effect size confidence intervals. http://www.cs.kent.ac.uk/pubs/2012/3233/





Ħ	HTTP Ping - MONO Regression Benchmarking - Mozilla										
<u> </u>	<u>E</u> dit	<u>V</u> iew	<u>G</u> o	<u>B</u> ook	kmarks	<u>T</u> ools	<u>W</u> indov	v <u>H</u> elp			
Ì 🔶	•	}	- 🥹	1		🗼 htt	n://nenva	ms.mff.cuni.cz/projects/mono/benchmarks/bttp_ping.phtml?	⇒.	រារ	ก
Back	Fo	orward	Rel	oad	Stop		p.,, nenya		Print	00	<u>u</u>
Patrice		5		Nα	2 2 2 2	2.55	<u>a</u> <u>b</u> <u>a</u> <u>b</u>	8724848788628862487499588			•
		60	900	- - - - - - - - - - - - - - - - - - -	565 565	444 000	244				
		005-	005-	005-	002	005-	005	222222222222222222222222222222222222			
		й Х	ನ ನ	а К	រ៍ស៊ស៊	ងីស័	มีผีผีผี	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~			
								Marcian			
								Version			
						_					
Newer	r Vers	sion O	lder V	/ersi	on Cha	ange In	npact[%]			
2005	5-10-0)4	2005-0	09-12	2	21.70	5 %	diff used-class-libraries-diff			
2005	5-10-1	1	2005-	10-10	D	-30.2	8 %	diff used-class-libraries-diff			
2005	5-10-1	7	2005-	10-14	4	-5.82	. %	diff used-class-libraries-diff			
2005	5-10-2	28	2005-	10-25	5	3.01	%	diff used-class-libraries-diff			
2005	5-12-0)9	2005-	12-08	8	1.93	%	diff used-class-libraries-diff			
2005	5-12-1	9	2005-	12-17	7	-1.8	%	diff used-class-libraries-diff			
2000	6-02-0)4	2006-0	02-03	3	-2.68	8 %	diff used-class-libraries-diff			
2000	5-04-2	27	2006-0	02-07	7	-11.1	2 %	diff used-class-libraries-diff			
2000	6-04-3	30	2006-0	04-27	7	1.3	%	diff used-class-libraries-diff			
2000	6-06-0)3	2006-0	05-24	4	-1.58	8 %	diff used-class-libraries-diff			
2000	6-06-1	4	2006-0	06-13	3	1.43	%	diff used-class-libraries-diff			
2000	6-08-1	6	2006-(08-15	5	18.39	9 %	diff used-class-libraries-diff			
. 384		http:/	Inores	-	off	an Innai-	ete les eve	(diffe2old=2005_10_148 now=2005_10_178 house http://www.aine			-
	∞	nttp://	menya.	.ms.n	nit.cuni.	cz/proje	cts/mono	yanis;oia=2005-10-14&new=2005-10-17&bench=http_ping		-	3

Experimental evaluation ingredients

Know the goal Know the system Know the real applications, benchmarks Run initial experiments, explore data Choose appropriate statistics Plan, dimension experiments Automate experiments execution Validate findings

Know the goal

"Our speed-up is 5%"

 $1 - \frac{\text{time on new system}}{\text{time on old system}}$

Papers from ISMM, PLDI, ASPLOS, TOPLAS, TACO (2011)

	Papers	Measured time	Reported speedup
Total	122	90	65

Over 70% of papers that measured time reported speed-up

Kalibera, Jones: Quantifying performance changes with effect size confidence intervals. http://www.cs.kent.ac.uk/pubs/2012/3233/

Know the goal

"Our speed-up is 5% ± 1.5% with 95% confidence"

 $1 - \frac{\text{time on new system}}{\text{time on old system}}$

Typically we want to compare two systems, estimating relative speed-up

Typically we focus on steady-state, mean performance

May have to re-visit as we progress

Know the system

Hardware

Frequency scaling, cores, caches, memory, performance counters

Operating system and system software

Scheduler, swapper, compiler, linker, page allocator Virtual machine

Just-in-time compiler, garbage collector

Know the system (FFT runs)



Individual samples, vertical lines denote new runs

Know the system (Mono builds)



Know real applications, benchmarks

Thread density

	Any Ops	Ali Madifs	Spok shared	Shared	Ap Os	At Modifs	Palate Not	Shared
			Density			Peri	odic density	
avrora ₉	25	22	25	25	22	22	22	22
$xalan_6$	8	8	8	8	8	8	8	8
tomcat ₉	9	15	12	14	9	8	7	8
$tradesoap_9$	13	44	32	28	7	13	12	11
$h2_9$	5	4	4	5	4	4	4	4
$tradebeans_9$	4	4	4	4	4	4	4	4
xalan ₉	4	4	4	4	4	4	4	4
pmd_9	5	5	5	5	4	4	3	3
hsqldb ₆	201	385	382	393	1	34	15	15
$lusearch_6$	63	59	66	64	44	0	57	59
lusearch ₉	4	5	5	5	4	0	4	4
$sunflow_9$	9	9	9	9	4	0	4	4
$antlr_6$	1	1	1	2	1	0	0	1
batik ₉	3	3	5	6	1	0	0	1
$bloat_6$	1	1	1	2	1	0	0	0
$chart_6$	1	2	1	3	1	0	0	0
eclipse ₆	2	6	5	5	1	0	0	1
eclipse ₉	28	197	191	72	1	1	0	1
fop ₆	1	1	1	2	1	0	0	0
fop ₉	1	1	1	2	1	0	0	1
jython ₆	1	1	1	2	1	0	0	0
jython ₉	1	1	2	2	1	0	0	1
$luindex_6$	1	2	2	2	1	0	0	1
luindex ₉	2	3	3	3	1	0	0	1
pmd ₆	1	1	1	2	1	0	0	1

Know real applications, benchmarks



Kalibera, Mole, Jones, Vitek: A black-box approach to understanding concurrency in DaCapo. OOPSLA'12

Know real applications, benchmarks



Kalibera, Mole, Jones, Vitek: A black-box approach to understanding concurrency in DaCapo. OOPSLA'12

FFT SciMark (Mono kernel benchmark)



🛄 FFT SciMark – run sequence plot



🛄 FFT SciMark – run sequence plot













FFT SciMark Execution Time [ms] **Original Randomly reordered**

lag 1

FFT SciMark – lag plot (lags 1,2,3,4)



Choose appropriate statistics

Interpret goal in precise/statistical terms, e.g.

- Ratio of mean execution times of new / old system
- Maximum pause time
- Maximum live size
- Decide on summarization method
 - Confidence interval, which kind
 - How to summarize over benchmarks, platforms
 - Report raw results "we have seen max live size"...
- Validate assumptions

Choose appropriate statistics

Confidence interval for the ratio of means Bootstrap method – statistical simulation Fieller's method Delta method

Kalibera, Jones: Quantifying performance changes with effect size confidence intervals. http://www.cs.kent.ac.uk/pubs/2012/3233/

Kalibera, Jones: Rigorous benchmarking in reasonable time. ISMM'13

Choose appropriate statistics

confidence interval for the ratio ex / ey

```
cfratio <- function(x, y, R = 10000, conf = 0.95) {
```

```
means <- sapply(1:R, function(i) {
    xs <- sample(x, replace = TRUE)
    ys <- sample(y, replace = TRUE)
    mean(xs) / mean(ys)
})</pre>
```

```
alpha <- (1 - conf) / 2
cfi <- quantile( means, probs=c( alpha, 1 - alpha ) )
names(cfi) <- NULL
```

```
c(lower = cfi[1], upper = cfi[2], meansRatio = mean(x) / mean(y))
Percentile bootstrap method.
```

Plan, dimension experiments

Identify variables impacting performance

Controlled, fixed – record and later report

Random (non-determinism) – **plan for repetition**

Uncontrolled fixed – turn into either of above

Typical, known sources of non-determinism include

- Code layout (linker, JIT)
- Memory layout (memory allocator in OS/VM/app)

Scheduler, GC, system services, network

Randomized algorithms, random (unique) names

Self-optimization, profiling

Plan repetition

Must repeat at highest "level" of experiment with non-determinism

If compilation is non-deterministic or randomized due to layout dependency, must run benchmarks for multiple builds

Can also repeat at lower "levels" of experiment

This can save experimentation time, if repeating at levels with high contribution to result variance

Repetition count to be chosen via initial experiments

Decide on warm-up iterations

For confidence interval, inference, need independent identically distributed values

If benchmark iterations are **dependent**, there is no point repeating them (must repeat at higher level)

For steady-state measurements

Drop initial values that are prone to obvious initialization noise

Dimensioning, warm-up is platform dependent.

Kalibera, Jones: Rigorous benchmarking in reasonable time. ISMM'13

Decide on warm-up iterations

Warm-up in Dacapo6 (platform, VM dependent!)

	Initialized	Independent
bloat	2	4
chart	3	∞
eclipse	5	7
fop	10	180
hsqldb	6	6
jython	3	∞
luindex	13	∞
lusearch	10	85
pmd	7	∞
xalan	6	13



Kalibera, Jones: Rigorous benchmarking in reasonable time. ISMM'13

Warm-up: automated heuristics fail Warm-up in Dacapo6 (platform, VM dependent!)

	Initialized	Independent	Harness	Georges
bloat	2	4	8	∞
chart	3	∞	4	1
eclipse	5	7	7	4
fop	10	180	7	8
hsqldb	6	6	8	15
jython	3	∞	5	2
luindex	13	∞	4	8
lusearch	10	85	7	8
pmd	7	∞	4	1
xalan	6	13	15	139

Kalibera, Jones: Rigorous benchmarking in reasonable time. ISMM'13



Ping benchmark, Linux, TAO/C++, Dual Pentium 3 SMP



FFT-SciMark benchmark, Linux/Mono, C#, Pentium 4

More tips and advice

- Automate experiment execution
- Archive scripts, data for reference
- Obtain and archive machines config
- Keep a log of decisions
- Archive graphs used for decisions
- Archive code used for evaluation

Use R

All results must be subject to critical (self)-review. Attempts for automated or mechanical-checklist evaluation produces garbage.

References

Kalibera, Jones: Rigorous benchmarking in reasonable time. ISMM'13.

Kalibera, Jones: Quantifying performance changes with effect size confidence intervals. http://www.cs.kent.ac.uk/pubs/2012/3233/

Kalibera, Mole, Jones, Vitek: A black-box approach to understanding concurrency in DaCapo. OOPSLA'12.

Vitek, Kalibera: R3: repeatability, reproducibility and rigor. ACM Sigplan. Not. 2012

My favorite sources

NIST Engineering Handbook: Engineering Statistics, http://www.itl.nist.gov/div898/handbook/

Jain: The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation, and Modeling. Wiley'91

Lilja: Measuring Computer Performance: A Practitioner's Guide. Cambridge. U. P.'05

Wasserman: All of Statistics, A Concise Course in Statistical Inference. Springer'04

(and many more as cited in the Uni of Kent technical report)