



Cloud Benchmarking

Estimating Cloud Application Performance Based on Micro Benchmark Profiling

Joel Scheuner





Problem



→ Unpractical to Test all Instance Types





Motivation







Research Questions

RQ1 – Performance Variability *within* **Instance Types**

Does the performance of equally configured cloud instances vary relevantly?

RQ2 – Application Performance Estimation *across* Instance Types

Can a set of micro benchmarks estimate application performance for cloud instances of different configurations?



RQ2.1 – Estimation Accuracy

How accurate can a set of micro benchmarks estimate application performance?



RQ2.2 – Micro Benchmark Selection

Which subset of micro benchmarks estimates application performance most accurately?





Methodology







Performance Data Set

Webservices	Instance Type	vCPU	ECU*	RAM [GiB]	Virtualization	Network Performance		
Γ	m1.small	1	1	1.7	PV	Low	eu + us 🏾	
	m1.medium	1	2	3.75	PV	Moderate		
	m3.medium	1	3	3.75	PV /HVM	Moderate	eu + us FRQ	1
	m1.large	2	4	7.5	PV	Moderate		
	m3.large	2	6.5	7.5	HVM	Moderate	eu	
RQ2 -{	m4.large	2	6.5	8.0	HVM	Moderate		
	c3.large	2	7	3.75	HVM	Moderate		
	c4.large	2	8	3.75	HVM	Moderate		
	c3.xlarge	4	14	7.5	HVM	Moderate		
	c4.xlarge	4	16	7.5	HVM	High		
	c1.xlarge	8	20	7	PV	High		

* ECU := Elastic Compute Unit (i.e., Amazon's metric for CPU performance)



>240 Virtual Machines (VMs) à 3 Iterations \rightarrow ~750 VM hours

>60'000 Measurements





RQ1 – Approach







RQ1 – Results







RQ1 – Implications



Hardware heterogeneity exploiting approaches are not worthwhile anymore [OZL+13, OZN+12, FJV+12]



Smaller sample size required to confidently assess instance type performance



Fair offer

[OZL+13] Z. Ou, H. Zhuang, A. Lukyanenko, J. K. Nurminen, P. Hui, V. Mazalov, and A. Ylä- Jääski. Is the same instance type created equal? exploiting heterogeneity of public clouds. *IEEE Transactions on Cloud Computing*, 1(2):201–214, 2013
 [OZN+12] Zhonghong Ou, Hao Zhuang, Jukka K. Nurminen, Antti Ylä-Jääski, and Pan Hui. Exploiting hardware heterogeneity within the same instance type of amazon ec2. In *Proceedings of the 4th USENIX Conference on Hot Topics in Cloud Computing* (*HotCloud'12*), 2012
 [FJV+12] Benjamin Farley, Ari Juels, Venkatanathan Varadarajan, Thomas Ristenpart, Kevin D. Bowers, and Michael M. Swift. More for your money: Exploiting performance heterogeneity in public clouds. In *Proceedings of the 3rd ACM Symposium on Cloud Computing* (SoCC '12), pages 20:1–20:14, 2012





RQ2 – Approach

RQ2 – Application Performance Estimation *across* **Instance Types**

Can a set of micro benchmarks estimate application performance for cloud instances of different configurations?







RQ2.1 – Results



How accurate can a set of micro benchmarks estimate application performance?







RQ2.2 – Results



RQ2.2 – Micro Benchmark Selection

Which subset of micro benchmarks estimates application performance most accurately?

Estimation Results for WPBench Read – Response Time

	Relative Error [%]	R ² [%]
Benchmark		
Sysbench – CPU Multi Thread	12.5	99.2
Sysbench – CPU Single Thread	454.0	85.1
Baseline		
vCPUs	616.0	68.0
ECU	359.0	64.6





RQ2 – Implications



Suitability of **selected** micro benchmarks to estimate application performance



Benchmarks cannot be used interchangeable → Configuration is important



Baseline metrics vCPU and ECU are insufficient





Related Work



Application Performance Profiling

- System-level resource monitoring [ECA+16, CBMG16]
- Compiler-level program similarity [HPE+06]

[ECA+16] Athanasia Evangelinou, Michele Ciavotta, Danilo Ardagna, Aliki Kopaneli, George Kousiouris, and Theodora Varvarigou.

Enterprise applications cloud rightsizing through a joint benchmarking and optimization approach.

Future Generation Computer Systems, 2016

[CBMG16] Mauro Canuto, Raimon Bosch, Mario Macias, and Jordi Guitart. A methodology for full-system power modeling in heterogeneous data centers.

In Proceedings of the 9th International Conference on Utility and Cloud Computing (UCC '16), pages 20–29, 2016

[HPE+06] Kenneth Hoste, Aashish Phansalkar, Lieven Eeckhout, Andy Georges, Lizy K. John, and Koen De Bosschere.

Performance prediction based on inherent program similarity. In Proceedings of the 15th International Conference on Parallel Architectures and Compilation Techniques (PACT '06), pages 114–122, 2006



- Trace and reply with Cloud-Prophet [LZZ+11, LZK+11]
- Bayesian cloud configuration refinement for big data analytics [ALC+17]

[LZZ+11] Ang Li, Xuanran Zong, Ming Zhang, Srikanth Kandula, and Xiaowei Yang. Cloud-prophet: predicting web application performance in the cloud. *ACM SIGCOMM Poster*, 2011

[LZK+11] Ang Li, Xuanran Zong, Srikanth Kandula, Xiaowei Yang, and Ming Zhang. Cloud-prophet: Towards application performance prediction in cloud. In *Proceedings of the ACM SIGCOMM 2011 Conference (SIGCOMM '11)*, pages 426–427, 2011

[ALC+17] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang.

Cherrypick: Adaptively unearthing the best cloud configurations for big data analytics.

In 14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17), 2017





Conclusion

RQ1 – Performance Variability within Instance Types

Does the performance of equally configured cloud instances vary relevantly?

Outcome: No. Performance does *not* vary relevantly for most benchmarks in Amazon's EC2 cloud for all intensively tested configurations in two different regions.

RQ2 – Application Performance Estimation *across* **Instance Types**

Can a set of micro benchmarks estimate application performance for cloud instances of different configurations?

Outcome: Yes. Selective micro benchmarks are able to estimate certain application performance metrics with acceptable accuracy.



RQ2.1 – Estimation Accuracy

How accurate can a set of micro benchmarks estimate application performance?

Outcome: Scientific computing application relative error below 10% Web serving application relative error between 10% and 20%



RQ2.2 – Micro Benchmark Selection

Which subset of micro benchmarks estimates application performance most accurately?

Outcome: Multi Thread CPU benchmark





APPENDIX





Motivation



Related Work (1)

[OZL+13] Z. Ou, H. Zhuang, A. Lukyanenko, J. K. Nurminen, P. Hui, V. Mazalov, and A. Ylä- Jääski. Is the same instance type created equal? exploiting heterogeneity of public clouds. *IEEE Transactions on Cloud Computing*, 1(2):201–214, 2013

[OZN+12] Zhonghong Ou, Hao Zhuang, Jukka K. Nurminen, Antti Ylä-Jääski, and Pan Hui. **Exploiting hardware heterogeneity within the same instance type of amazon ec2**. In *Proceedings of the 4th USENIX Conference on Hot Topics in Cloud Computing (HotCloud'12)*, 2012

[FJV+12] Benjamin Farley, Ari Juels, Venkatanathan Varadarajan, Thomas Ristenpart, Kevin D. Bowers, and Michael M. Swift. **More for your money: Exploiting performance heterogeneity in public clouds**. In *Proceedings of the 3rd ACM Symposium on Cloud Computing (SoCC '12)*, pages 20:1–20:14, 2012 [DPC10] Jiang Dejun, Guillaume Pierre, and Chi-Hung Chi. **EC2 Performance Analysis for Resource Provisioning of Service-Oriented Applications**, pages 197–207. Springer, 2010

[SSS+08] Will Sobel, Shanti Subramanyam, Akara Sucharitakul, Jimmy Nguyen, Hubert Wong, Arthur Klepchukov, Sheetal Patil, Armando Fox, and David Patterson. Cloudstone: Multi-platform, multi-language benchmark and measurement tools for web 2.0, 2008

[PSF16] Tapti Palit, Yongming Shen, and Michael Ferdman. **Demystifying cloud benchmarking**. In 2016 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), pages 122–132, 2016

Related Work (2)

Application Performance Profiling

- System-level resource monitoring and micro benchmarks [ECA+16, CBMG16]
- Compiler-level program similarity [HPE+06]

- Trace and reply with Cloud-Prophet [LZZ+11, LZK+11]
- Bayesian cloud configuration refinement for big data analytics [ALC+17]
- Multi-component queuing models [SS05]

[ECA+16] Athanasia Evangelinou, Michele Ciavotta, Danilo Ardagna, Aliki Kopaneli, George Kousiouris, and Theodora Varvarigou. **Enterprise applications cloud** rightsizing through a joint benchmarking and optimization approach. *Future Generation Computer Systems*, 2016

[CBMG16] Mauro Canuto, Raimon Bosch, Mario Macias, and Jordi Guitart. A methodology for full-system power modeling in heterogeneous data centers. In Proceedings of the 9th International Conference on Utility and Cloud Computing (UCC '16), pages 20–29, 2016

[HPE+06] Kenneth Hoste, Aashish Phansalkar, Lieven Eeckhout, Andy Georges, Lizy K. John, and Koen De Bosschere. **Performance prediction based on inherent program similarity**. In *Proceedings of the 15th International Conference on Parallel Architectures and Compilation Techniques (PACT '06)*, pages 114– 122, 2006

[LZZ+11] Ang Li, Xuanran Zong, Ming Zhang, Srikanth Kandula, and Xiaowei Yang. Cloud-prophet: predicting web application performance in the cloud. ACM SIGCOMM Poster, 2011

[LZK+11] Ang Li, Xuanran Zong, Srikanth Kandula, Xiaowei Yang, and Ming Zhang. **Cloud-prophet: Towards application performance prediction in cloud**. In *Proceedings of the ACM SIGCOMM 2011 Conference (SIGCOMM '11)*, pages 426–427, 2011

[ALC+17] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. Cherrypick: Adaptively unearthing the best cloud configurations for big data analytics. In 14th USENIX Symposium on Networked Systems Design and Implementation (NSDI 17), 2017

[SS05] Christopher Stewart and Kai Shen. **Performance modeling and system management for multi-component online services**. In *Proceedings of the 2nd Conference on Symposium on Networked Systems Design & Implementation - Volume 2*, NSDI'05, pages 71– 84, Berkeley, 2005

Molecular Dynamics

Benchmark Design

Benchmark Execution

Benchmark Execution – Data Set

Webservices	Instance Type	vCPU	ECU*	RAM [GiB]	Virtualization	Network Performance	-
Γ	m1.small	1	1	1.7	PV	Low	eu + us 🔵
	m1.medium	1	2	3.75	PV	Moderate	
	m3.medium	1	3	3.75	PV /HVM	Moderate	eu + us - RQ1
	m1.large	2	4	7.5	PV	Moderate	
	m3.large	2	6.5	7.5	HVM	Moderate	eu
RQ2 -{	m4.large	2	6.5	8.0	HVM	Moderate	
	c3.large	2	7	3.75	HVM	Moderate	
	c4.large	2	8	3.75	HVM	Moderate	
	c3.xlarge	4	14	7.5	HVM	Moderate	
	c4.xlarge	4	16	7.5	HVM	High	
	c1.xlarge	8	20	7	PV	High	_

* ECU := Elastic Compute Unit (i.e., Amazon's metric for CPU performance)

Benchmark

Execution

>240 Virtual Machines (VMs) à 3 Iterations \rightarrow ~750 VM hours

Data Pre-

Processing

Data

Analyses

>60'000 Measurements

Benchmark

Design

Data Pre-Processing

S. e. a. . software evolution & architecture lab

Department of Informatics – s.e.a.l.

Data Analyses – Implementation

	Views: Design Results	🚱 Need help? 💌
Repository ×	Process × Context ×	Parameters ×
🔾 Add Data 🛛 = 👻	🕢 Process + Forward Selection + 100% 🔑 🔑 📮 🧟 💣 🔛	% Performance (Performance (
Od jen Od jen Odo-menn-agregatid-selecci-f Odo-menn-agregatid-selecci-f Odo-menn-agregatid selecci-f Odo-menn-selecci laster selecci Odo-menn-selecci laster selecci laster selecci Odo-menn-selecci laster selecci Odo-menn-selecci laster selecci Odo-menn-selecci laster selecci Odo-menn-selecci laster selecci laster selecci Odo-menn-selecci laster selecci laster selecci laster selecci Odo-menn-selecci laster selecci laster seleci Odo-menn-selecci laster selecci laster selecci la	Average Contracts	main criterion relation_ent ♥ 0 ↑ ♥ root mean squared error ♥ ♥ abusize error ♥ 0 ♥ relative error ♥ 0 ♥ relative error solution ♥ relative error solution ♥ relative error solution ♥ relative error solution
Saldrift (1997) Data Access (47) A Bending (77) Clanning (76) Modeling (129) ▼ Predictive (51) b Clanning (7)		restricted abstance error 0 v restricted abstanced error 0 v This e absanced error 0 v Help ×
Bayesian (2) Trees (9) Get more operators from the Marketplace	Raccommended Operators 0 • Ty Spit Data g# 51% 1 Select Antibutes g# 33% 1 Select Antibutes g# 27%	Linear Regression RapidMiner Studio Core Tegr: Supervised. Classification.

🛞 Project: (Nor
-
E List +
Q,
1181 141 1181 -
15" "30" "7" "1.
-
🐵 Publish 👻
1
-
A 6.83
 (4)
 mülaros (su)

Benchmark

Design

joe4dev/cwb-analysis

This repository Sea	Pull requests issues Ma	ketplace Gist	▲ +- 願-
은 joe4dev / cwb-anal	ysis Private	O Unwatch -	1 ★Star 0 ÿFork 0
O Code ③ Issues ④	1) Pull requests (0) III Projects (0) III Wiki	© Settings Insights	Ŧ
CWB analysis scripts and cloud-workbench Manage to	d data pics		Edit
@ 73 commits	ÿ 3 branches	♥ 0 releases	AL 1 contributor
Branch: master * New put	ll request	Create new file Upload fil	es Find file Clone or download +
ioe4dev committed on G	itHub Merge pull request #1 from joe4dev/feature/rq2 📖		Latest commit ea8d578 7 days ago
🖿 data_interim	Update interim dataset readme		16 days ago
🖿 data_raw	Update data with additional instances for RQ2		16 days ago
M db_dumps	Update data with additional instances for RQ2		16 days ago
illi rq1	Use violin plot instead of boxplot for rq1		14 days ago
III rq2	Add instance type filter boundary variable		14 days ago
E README.md	Refine docs		26 days ago
III README.md			

CWB Data

This repository stores the data (raw and interim) and analysis scripts for the CWB cloud benchmarking operiments combining micro and application benchmarks. Every directory contains the data (incl. docs) and scipts (incl. docs) to produce or update this data.

Repository Structure

File	Explanation
lb_dumps	PostgreSQL database dumps from CWB
loud_workbench_production.sql	Exported from CWB via dump-cwb-db.sh
- dump-cwb-db.md	Dump CWB DB docs
ump-cwb-db.sh	Dumps CWB DB to cloud_workbench_production.sql
lata_interim	Pre-processed intermediate data
- cwb-data-interim.csv	Pre-processed raw data via pre-process.sh
- cwb-data-interim.md	Describes the interim (i.e., pre-processed) data set
pre-process.md	Pre-process docs
pre-process.rmp	RapidMiner Studioprocess
pre-process.sh	Pre-process runner script
lata_raw	Unaltered raw data as empirically observed
cwb-data-raw.csv	Exported from CWB via export-cwb-data.sh OF export-cwb- data_local.sh
wb-data-raw.md	Describes the raw data set
export-cwb-data.md	CWB exporter docs
export-cwb-data.rb	CWB raw data exporter script
export-cwb-data.sh	Exports CWB data to cwb-data-raw.csv
export-cwb-data-local.sh	Exports local CWB data to cwb-data-raw.csv

2017-06-15

Benchmark Execution Data Pre-Processing Data Analyses

Page 24

Data Analyses – Results

Guided by the Research Questions ...

RQ1 – Performance Variability within Instance Types

Does the performance of equally configured cloud instances vary relevantly?

RQ2 – Application Performance Estimation across Instance Types

Can a set of micro benchmarks estimate application performance for cloud instances of different configurations?

े 0 0

RQ2.1 – Estimation Accuracy

How accurate can a set of micro benchmarks estimate application performance?

RQ2.2 – Micro Benchmark Selection

Which subset of micro benchmarks estimates application performance most accurately?

2017-06-15 Benchmark Benchmark Data Pre-Design Execution Processing Analyses

WPBench Write – Root Cause Analysis

Contributions

- 1. Extension of Cloud WorkBench (CWB) [SLCG14, SCLG15] with modular plugin system
- 2. Newly crafted Web serving application benchmark WPBench with three different load scenarios
- 3. Automated CWB benchmark that combines single-instance and multiinstance micro and application benchmarks
- 4. Raw and cleaned data set with performance metrics from Amazon EC2
- 5. Evaluation of an estimation model for application performance based on micro benchmark profiling

[SLCG14] Joel Scheuner, Philipp Leitner, Jürgen Cito, and Harald Gall. **Cloud WorkBench - Infrastructure-as-Code Based Cloud Benchmarking**. In *Proceedings of the 6th IEEE International Conference on Cloud Computing Technology and Science (CloudCom'14)*, 2014

Threats to Validity

Construct Validity

Almost 100% of benchmarking reports are wrong because benchmarking is "very very error-prone"¹ [senior performance architect @Netflix]

 \rightarrow Guidelines, rationalization, open source

External Validity (Generalizability)

Other cloud providers? Larger instance types? Other application domains? → Future work

Internal Validity

the extent to which cloud environmental factors, such as multi-tenancy, evolving infrastructure, or dynamic resource limits, affect the performance level of a VM instance

 \rightarrow Variability RQ1, stop interfering process

Reproducibility

the extent to which the methodology and analysis is repeatable at any time for anyone and thereby leads to the same conclusions

A dynamic cloud environment

 \rightarrow Fully automated execution, open source

Open Challenges

• How to identify a suitable micro benchmark estimator?

Future Work

More IaaS providers \rightarrow Custom instance types

Runtime performance data Dedicated performance testing \rightarrow Instance type selection as integral part of (vertical) scaling strategies

Multi-instance application architectures

Conclusion (1)

RQ1 – Performance Variability *within* Instance Types

Does the performance of equally configured cloud instances vary relevantly?

Outcome: No. Performance does *not* vary relevantly for most benchmarks in Amazon's EC2 cloud for all intensively tested configurations in two different regions.

RQ2 – Application Performance Estimation *across* **Instance Types**

Can a set of micro benchmarks estimate application performance for cloud instances of different configurations?

Outcome: Yes. Selective micro benchmarks are able to estimate certain application performance metrics with acceptable accuracy.

Conclusion (2)

RQ2.1 – Estimation Accuracy

How accurate can a set of micro benchmarks estimate application performance?

Outcome: A scientific computing application achieves relative error rates below 10% and the response time of a Web serving application is estimated with a relative error between 10% and 20%.

RQ2.2 – Micro Benchmark Selection

Which subset of micro benchmarks estimates application performance most accurately?

Outcome: A single CPU benchmark was able to estimate the duration of a scientific computing application and the response time of a Web serving application most accurately.

Summary

Research Questions

performance most accurately?

Prepare

eanup

Background – Micro Benchmarks

File I/O: 4k random read I/O ubuntu@ip-172-31-4-0:/tmp\$ sysbench --test=fileio --file-total-size=40 --file-block-size=4K --file-test-mode=rndrd prepare sysbench 0.4.12: multi-threaded system evaluation benchmark 128 files, 32768Kb each, 4096Mb total ubuntu@ip-172-31-4-0:/tmp\$ sysbench --test=fileio --file-total-size= --file-block-size=4K --file-test-mode=rndrd run multi-innequeu system evaluation bench Running the test with following options: Number of threads: 1 Extra file open flags: 0 128 files, 32Mb each 4Gb total file size Block size 4Kb Number of random requests for random IO: 10000 Read/Write ratio for combined random IO test: 1.50 Periodic FSYNC enabled, calling fsync() each 100 requests. Calling fsync() at the end of test, Enabled. Using synchronous I/O mode Doing random read test Threads started! Done. Operations performed: 10000 Read, 0 Write, 0 Other Read 39.062Mb Written 0b Total transferred 39.062Mb (3.5793Mb/sec) 916.31 Requests/sec executed B Test execution summary: total time: 10.9133s total number of events: 10000 total time taken by event execution: 10.9092 per-request statistics: 0.00ms 1.09ms max: 596.05ms approx. 95 percentile: 3.27ms Threads fairness: events (ava/stddev): 10000.0000/0.00 execution time (avg/stddev): 10.9092/0.00 ubuntu@ip-172-31-4-0:/tmp\$ sysbench --test=fileio --file-total-size=4 --file-block-size=4K --file-test-mode=rndrd cleanup sysbench 0.4.12: multi-threaded system evaluation benchmark

emoving test files..

, 1 Bandwidth
Network
Server
ubuntu@ir <mark>-172-31-4-0:</mark> 'tmp\$ iperfserverlen 128k
Server listening on TCP port 5001 TCP window size: 85.3 KByte (default)
[4] local 172.31.4.0 port 5001 connected with 172.31.10.209 port 46432 [ID] Interval Transfel Bandwidth [4] 0.0-30.0 sec 3.39 GBytes 971 Mbits/sec
Client
ubuntu@ip-172-31-10-209:/tmp\$ iperf - 172.31.4.0 ·l 128k -t 30
Client connecting to 172.31.4.0, TCP port 5001 TCP window size: 325 KByte (default)
[3] local 172.31.10.209 port 46432 connected with 172.31.4.0 port 5001 [ID] Interval Transfer Pandwidth [3] 0.0-30.0 sec 3.39 GBytes 972 Mbits/sec
972 Mbits/sec

Background – Application Benchmarks

Test plan (JMeter)

Webapp (Wordpress)

Related Work

Combining Micro and Application Benchmarks

- B. Varghese and O. Akgun and I. Miguel and L. Thai and A. Barker,
 "Cloud Benchmarking For Maximising Performance of Scientific Applications" in IEEE TRANSACTIONS ON CLOUD COMPUTING (2016)
- A. Evangelinou and M. Ciavotta and D. Ardagna and A. Kopaneli and G. Kousiouris and T. Varvarigou, "Enterprise applications cloud rightsizing through a joint benchmarking and optimization approach" *Future Generation Computer Systems* - (2016)

Hardware / Performance Heterogeneity

- Farley, Benjamin and Juels, Ari and Varadarajan, Venkatanathan and Ristenpart, Thomas and Bowers, Kevin D. and Swift, Michael M., "More for Your Money: Exploiting Performance Heterogeneity in Public Clouds" Proceedings of the Third ACM Symposium on Cloud Computing (SoCC '12)
- Z. Ou and H. Zhuang and A. Lukyanenko and J. K. Nurminen and P. Hui and V. Mazalov and A. Ylä-Jääski, "Is the Same Instance Type Created Equal? Exploiting Heterogeneity of Public Clouds" IEEE Transactions on Cloud Computing (2013)

Typical Performance Data

- Instance metadata
 - CPU model
 - # CPU cores
- Benchmark metadata
 - Version number (including compiler)
- Benchmark
 - Execution time
 - Latency / response time
 - Throughput / Bandwidth
 - Operations / sec

Execution Methodology

Ali Abedi and Tim Brecht, Conducting Repeatable Experiments in Highly Variable Cloud Computing Environments (ICPE'17)