

An Empirical Study on the Discrepancy between Performance Testing Results from Virtual and Physical Environments



Muhammad Moiz Arif



Weiye Shang




Emad Shihab

Department of Engineering & Computer Science
Concordia University, Montreal, Canada.

What are performance tests?

- Performance assurance activities: ensuring software meets performance requirements.
- Mimicking user behaviour.

Why are performance tests important?

The background of the slide features a dark, textured surface with the Amazon logo and the words 'amazon web services' projected onto it in a bright, glowing light. The logo consists of a yellow cube-like shape above the word 'amazon' in a white, sans-serif font, with 'web services' in a smaller font below it.

Amazon estimates that a one-second page-load slowdown can cost up to \$1.6 billion!

Why are performance tests important?

HealthCare.gov

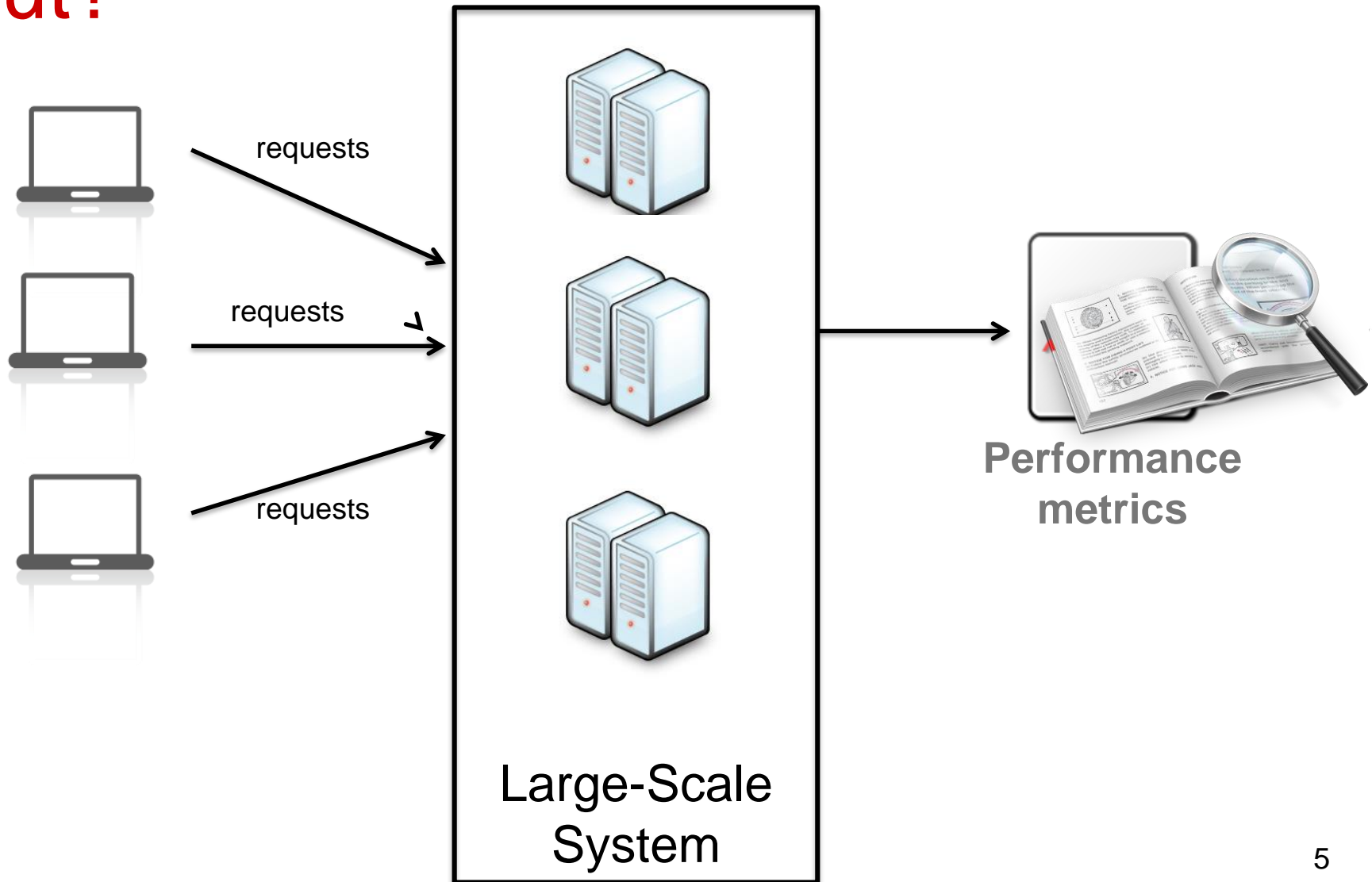
Individuals & Families

Small Businesses

The Health Insurance Marketplace online application isn't available from app
we make improvements. Additional down times may be possible as we work to
and the Marketplace call center remain available during these hours.

Healthcare.gov failed to rollout successfully
because of lack of performance testing!

How is a performance test carried out?

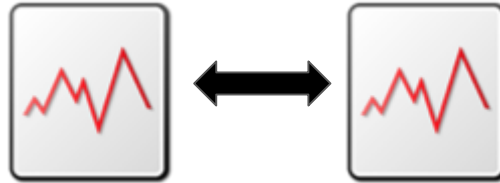


Literature review: Types of performance analysis



- Single performance metric

[Nguyen, T.H. et al, ICPE, 2012]



- Relationship between performance metrics

[Cohen, I. et al, OSDI, 2004]



- Statistical modeling on performance metrics

[Shang, W. et al, ICPE, 2015]

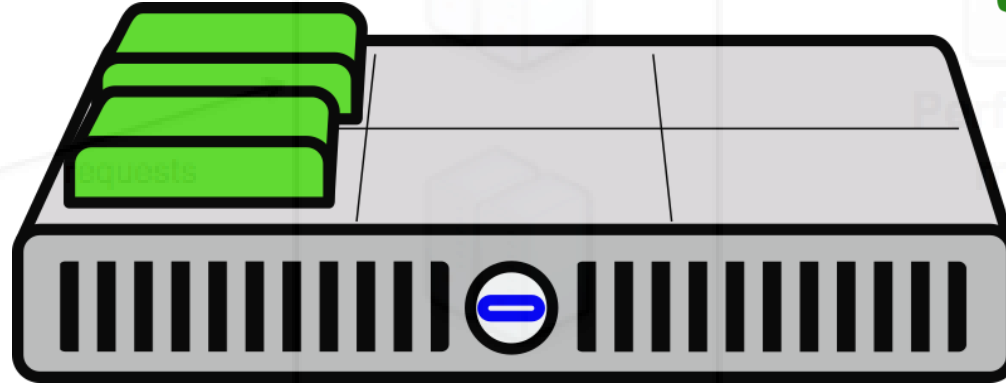
How is a performance test carried out?

What if you need to modify your environment?



How is a performance test carried out?

What if you need to modify your environment?



VIRTUAL MACHINES!

Examples from the industry


Load Testing SugarCRM in a Virtual Machine

Determining the CPU cost of virtualization with VMware ESX

Christopher L Merrill

©2008 Web Performance, Inc.

v1.1

 [Read and post comments](#)

Summary

The performance of our reference application under load (a default SugarCRM installation) on a virtualized server showed a 14% decrease, measured by total system capacity, compared to the same system running natively on equivalent hardware.

Overview

In a typical virtualization deployment scenario, a virtual machine replaces an existing physical machine. Multiple VMs will usually be deployed on a single host machine, sharing the resources of the host. For instance, VMs with a total of 16 cores might be deployed on a 4 core host machine. As long as utilization remains low, the processor resources can be shared while maintaining application performance goals. However, when the utilization of a single VM becomes very high, performance of one or all the VMs will suffer if the resource sharing continues. With the right hardware and software

Examples from the industry

Load Testing SugarCRM in a Virtual Machine

Determining the CPU cost of virtualization with VMware ESX

*“The **performance of our reference application** under load (a default SugarCRM installation) on a **virtualized server** showed a 14% decrease, measured by total system capacity, **compared to the same system running natively on equivalent hardware.**”*

Overview

In a typical virtualization deployment scenario, a single physical machine. Multiple VMs will usually be deployed on a single host machine, sharing the resources of the host. For instance, VMs with a total of 16 cores might be deployed on a 4 core host machine. As long as utilization remains low, the processor resources can be shared while maintaining application performance goals. However, when the utilization of a single VM becomes very high, performance of one or all the VMs will suffer if the resource sharing continues. With the right hardware and software

Examples from the industry

performance-testing systems on virtual machines that normally run on physical machines



7



1

My employer runs some of our systems on physical machines with attached hard drives. I am charged with performance-testing those systems. For cost reasons, I've been asked to test those systems running on virtual machines (using Xen) attached to a SAN. This is clearly not an apples-to-apples comparison. Some systems use a lot of disk I/O, and so the SAN issue is especially worrisome. Rather than responding with "can't be done" or "not reliable", I want to recommend what *is* possible.

Here are some things that come to mind or that I've found with Google searches:

- Measure SAN speed vs. hard drive and calculate a ratio
- Borrow a physical machine long enough to run a benchmark, do the same with a virtual machine and calculate a ratio
- Even if you can't predict absolute performance on physical machines, you may be able to predict relative performance (i.e. whether the candidate release will be faster or slower than what's currently in production)
- Measure multiple times at different times of day to mitigate resource contention issues, i.e. conflicts with other virtual machines running on the same physical machine or with other clients using the SAN.

Are there other things you can do to mitigate differences between physical machine performance and virtual machine performance in an environment similar to mine? I am particularly interested in actual experiences rather than educated guesses.

performance virtualization xen

asked

view

active



Ge

-
-
-

Rela

1

3

Examples from the industry

performance-testing systems on virtual machines that normally run on physical machines



7

My employer runs some of our systems on physical machines with attached hard drives. I am charged with performance-testing those systems. For cost reasons, I've been asked to test those systems running on virtual machines (using Xen) attached to a SAN. This is clearly not an apples-to-apples comparison. Some systems use a lot of disk I/O, and so the SAN issue is especially worrisome. Rather than responding with "can't be done" or "not reliable", I want to recommend what

...*“For **cost reasons**, I've been asked to test those systems running on **virtual machines** (using Xen) attached to a SAN. **This is clearly not an apples-to-apples comparison.**”*...

Here are some things that come to mind or that I've found with Google searches:

- Borrow a physical machine long enough to run a benchmark, do the same with a virtual
- Even if you can't predict absolute performance on physical machines, you may be able to
- Measure multiple times at different times of day to mitigate resource contention issues, i.e. conflicts with other virtual machines running on the same physical machine or with other clients using the SAN.

Are there other things you can do to mitigate differences between physical machine performance and virtual machine performance in an environment similar to mine? I am particularly interested in actual experiences rather than educated guesses.

performance

virtualization

xen

asked

view

activ

Go

Rela

1

3

Research hypothesis

- *For software testing activities there exists a discrepancy between **physical** and **virtual** environments.*
- *We believe that the approaches used so far do not take into account the heterogeneous environments.*

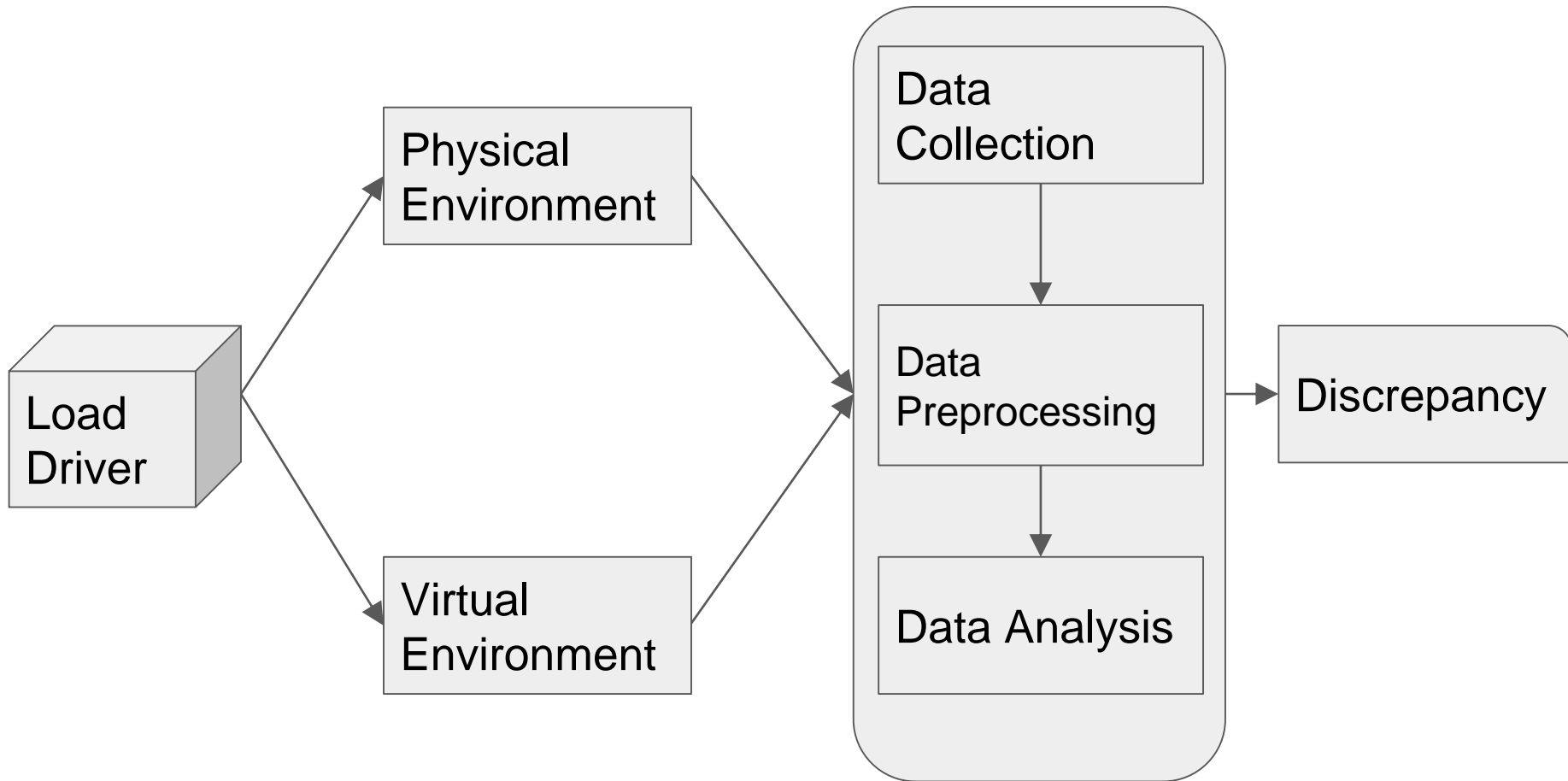
Research Questions

RQ1: Are the trend and distribution of a single performance metric similar across environments?

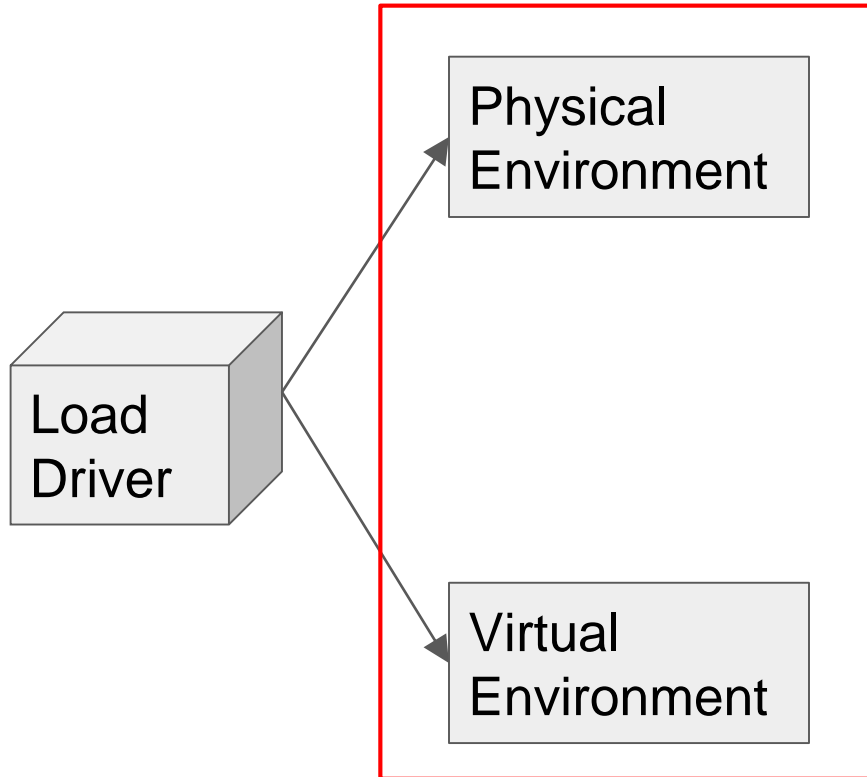
RQ2: To what extent does the relationship between the performance metrics change across environments?

RQ3: Can statistical performance models be applied across virtual and physical environments?

Approach



Approach



1. Dell DVD Store (DS2)
2. Cloudstore

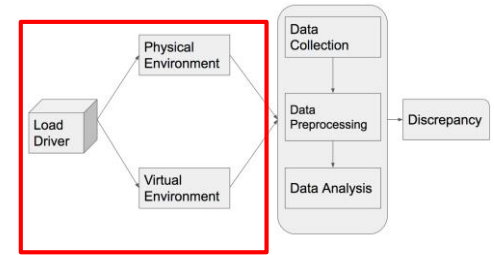
Approach

- Both the system are used in prior studies.
- DS2 comes with a load driver.
- Cloudstore, we used JMeter to replicate the load.



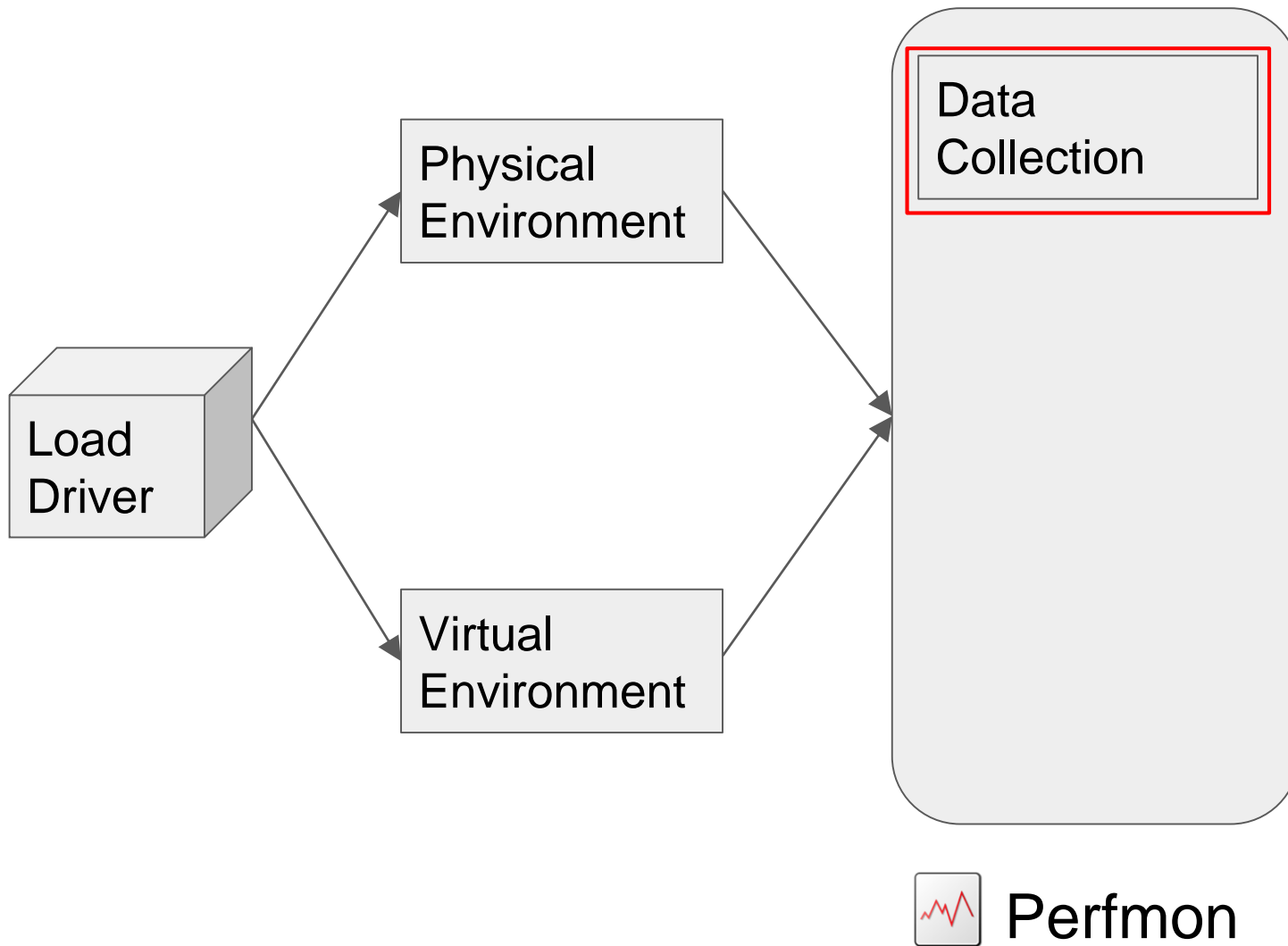
1. Dell DVD Store (DS2)
2. Cloudstore

Approach

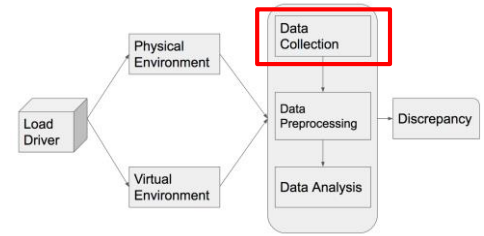


- We set up 3 nodes: Application, database and load driver.
- One virtual environment set up on each node: *single tenancy*.
- We set up identical config. between the two environments. i.e. CPU 2 cores, Memory 3GB.

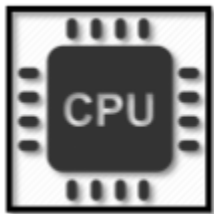
Approach



Approach



- We used ***PERFMON*** to monitor application and database server.
- We recorded *all* the performance metrics available.



CPU

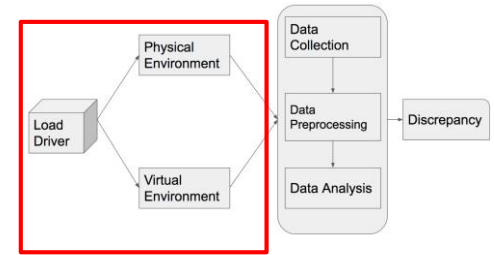


I/O Ops



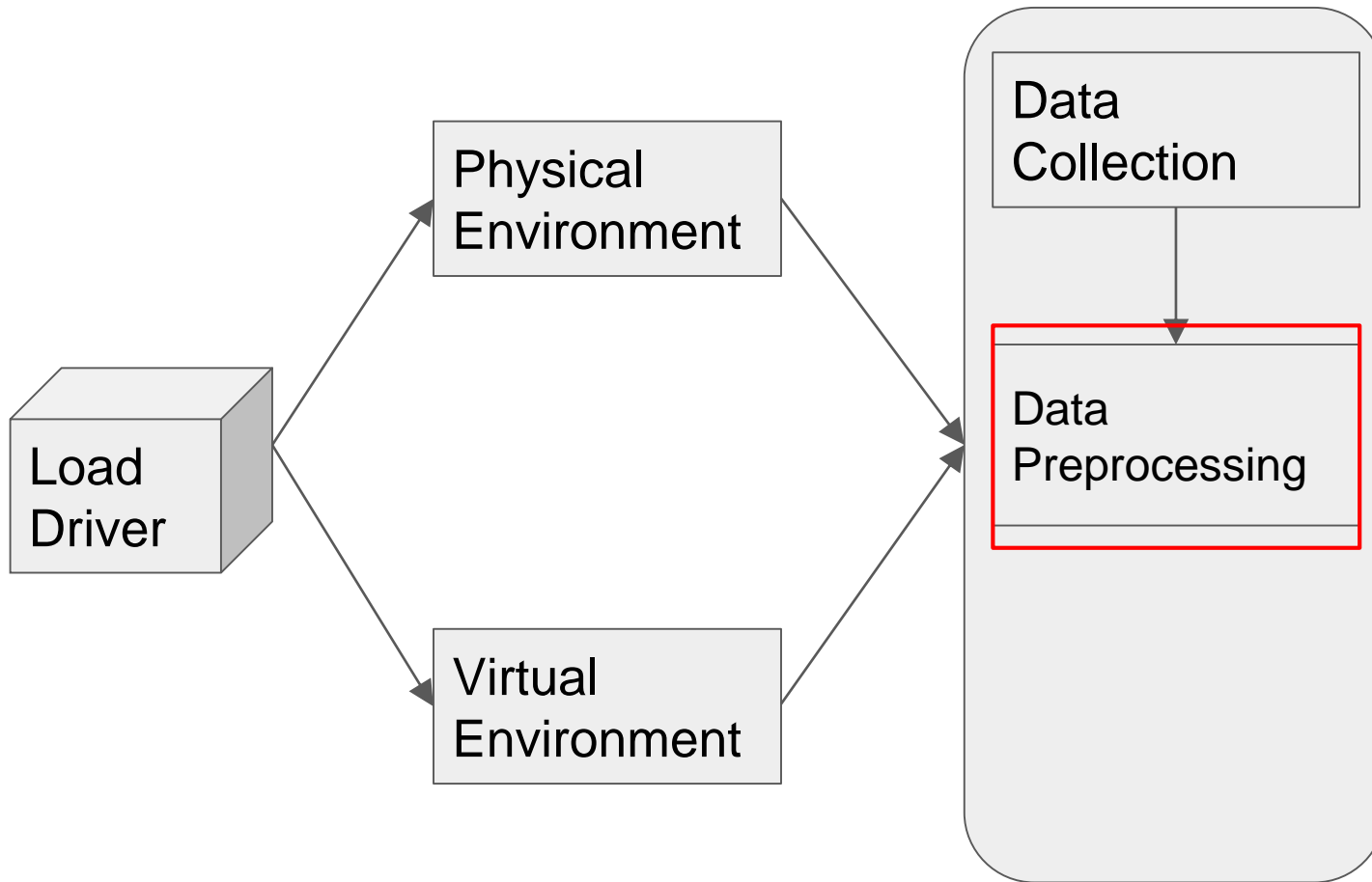
Memory

Approach

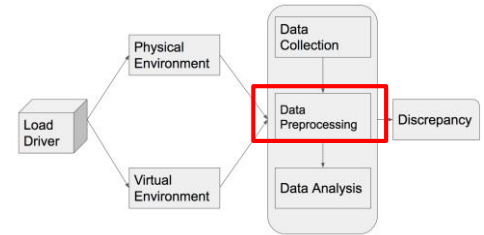


- *Workload* variation introduced by number of threads (threads=users) .
- The variation of *# of threads* was **identical**, **periodic** and **random** across the environments.
- Runtime for the test: 9+ hours.

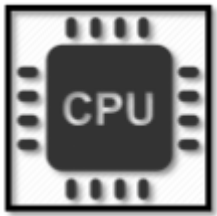
Approach



Approach



System Level



CPU



I/O Ops



Memory

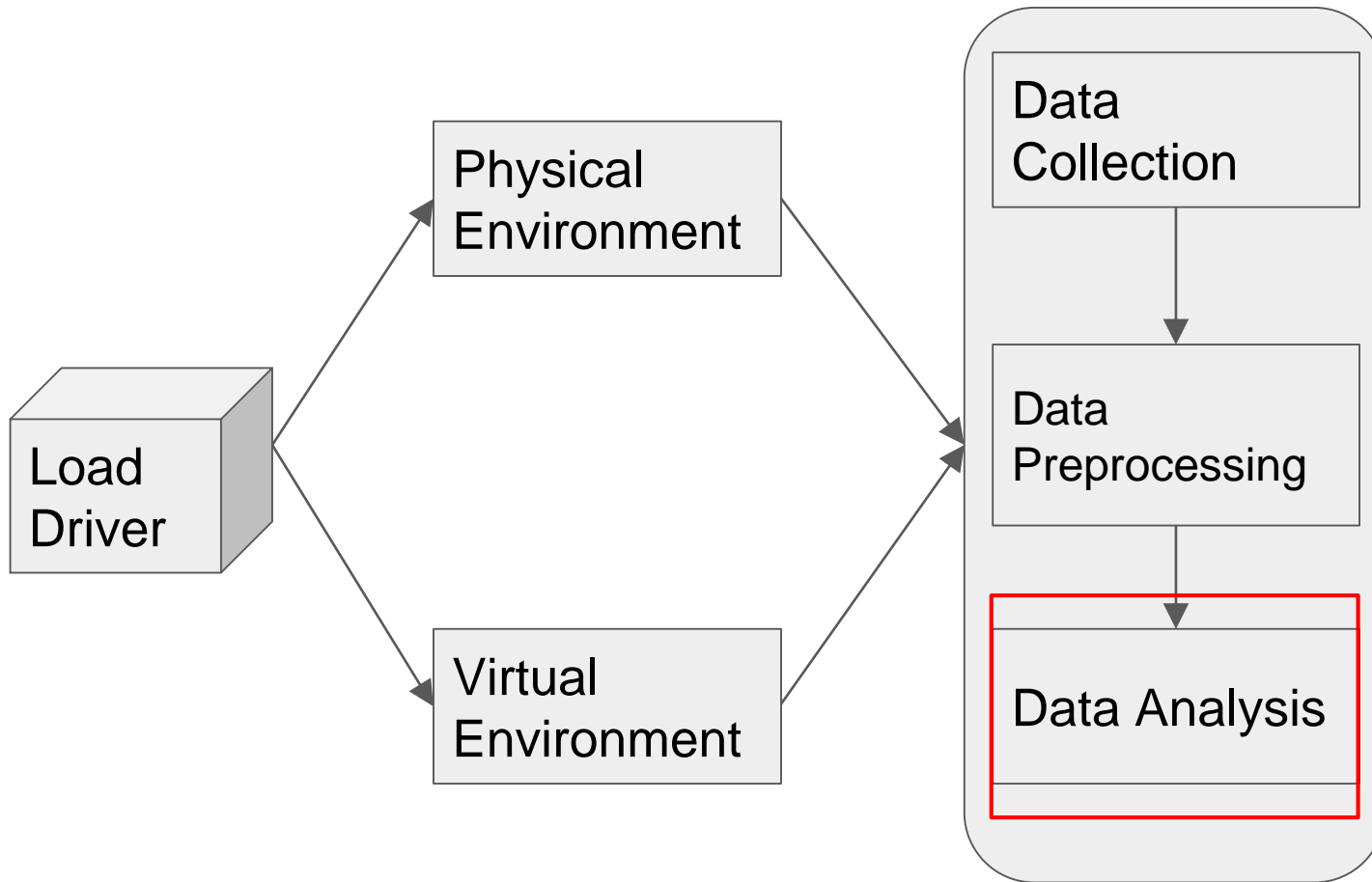
Application Level



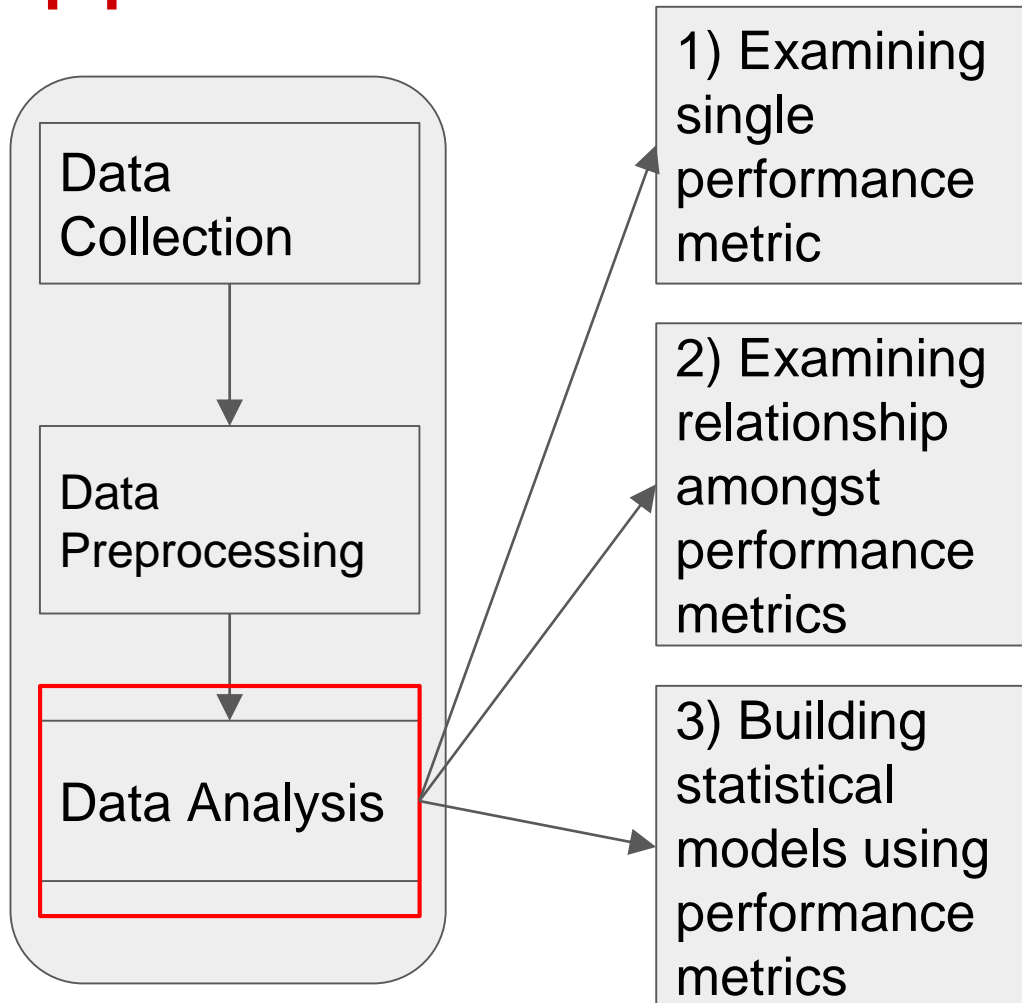
Throughput

Timestamps on logs used to calculate # of request/minute.

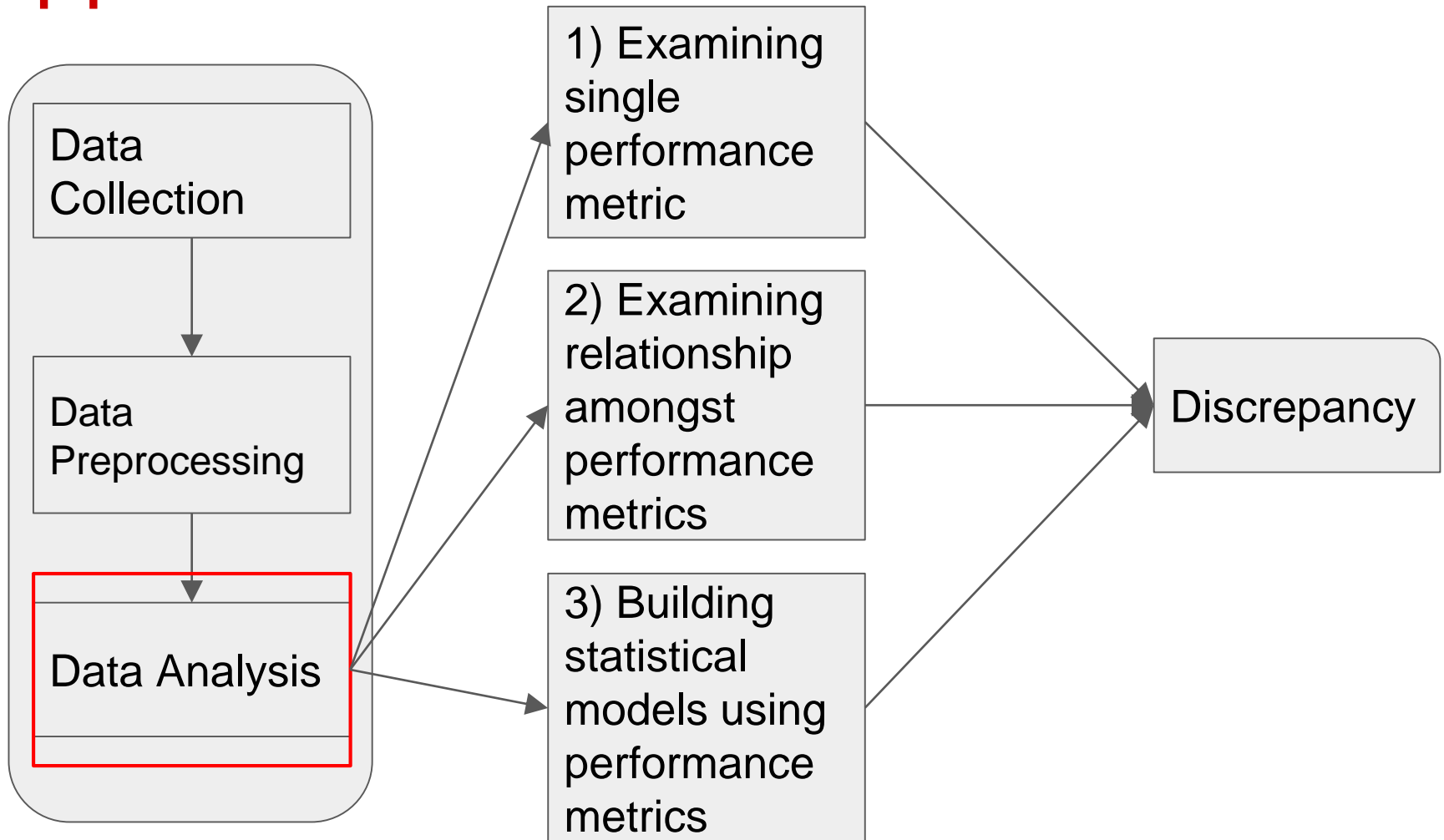
Approach



Approach



Approach



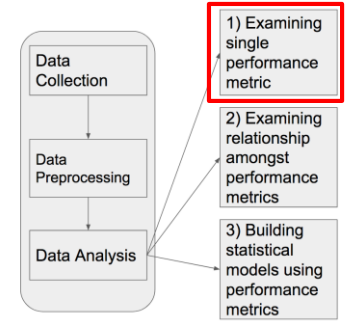
Research Questions

RQ1: Are the **trend and distribution** of a **single performance metric** similar across environments?

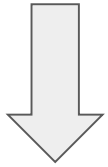
RQ2: To what extent does the relationship between the performance metrics change across environments?

RQ3: Can statistical performance models be applied across virtual and physical environments?

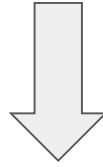
RQ1: Approach



Shape of the distribution

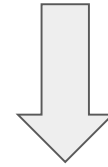


Quantile-
Quantile
(QQ)
Plots



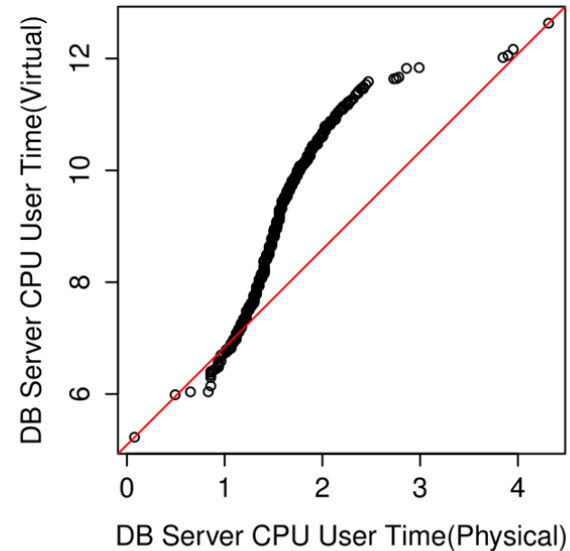
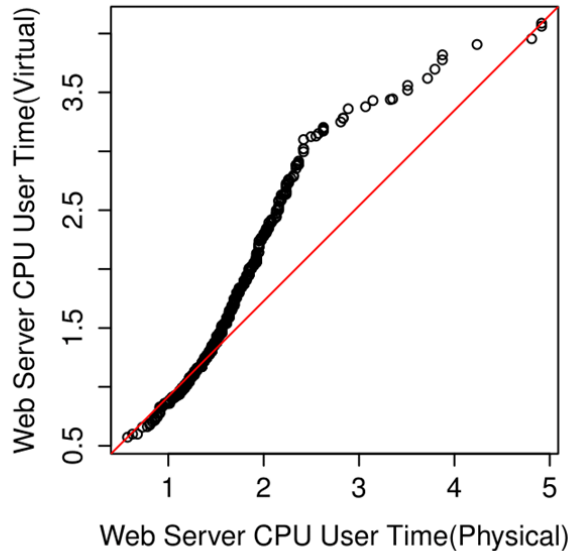
Normalized
KS-Test

Trend

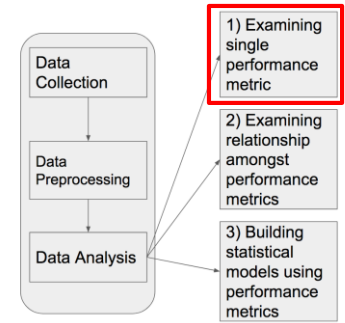


Spearman's
rank
correlation
coefficient

DS2: Most performance metrics do not follow the same shape of the distribution in virtual and physical environments.



KS-Test: Performance metrics that do not follow the same distribution.



KS-Test (P-value > 0.05)

DS2

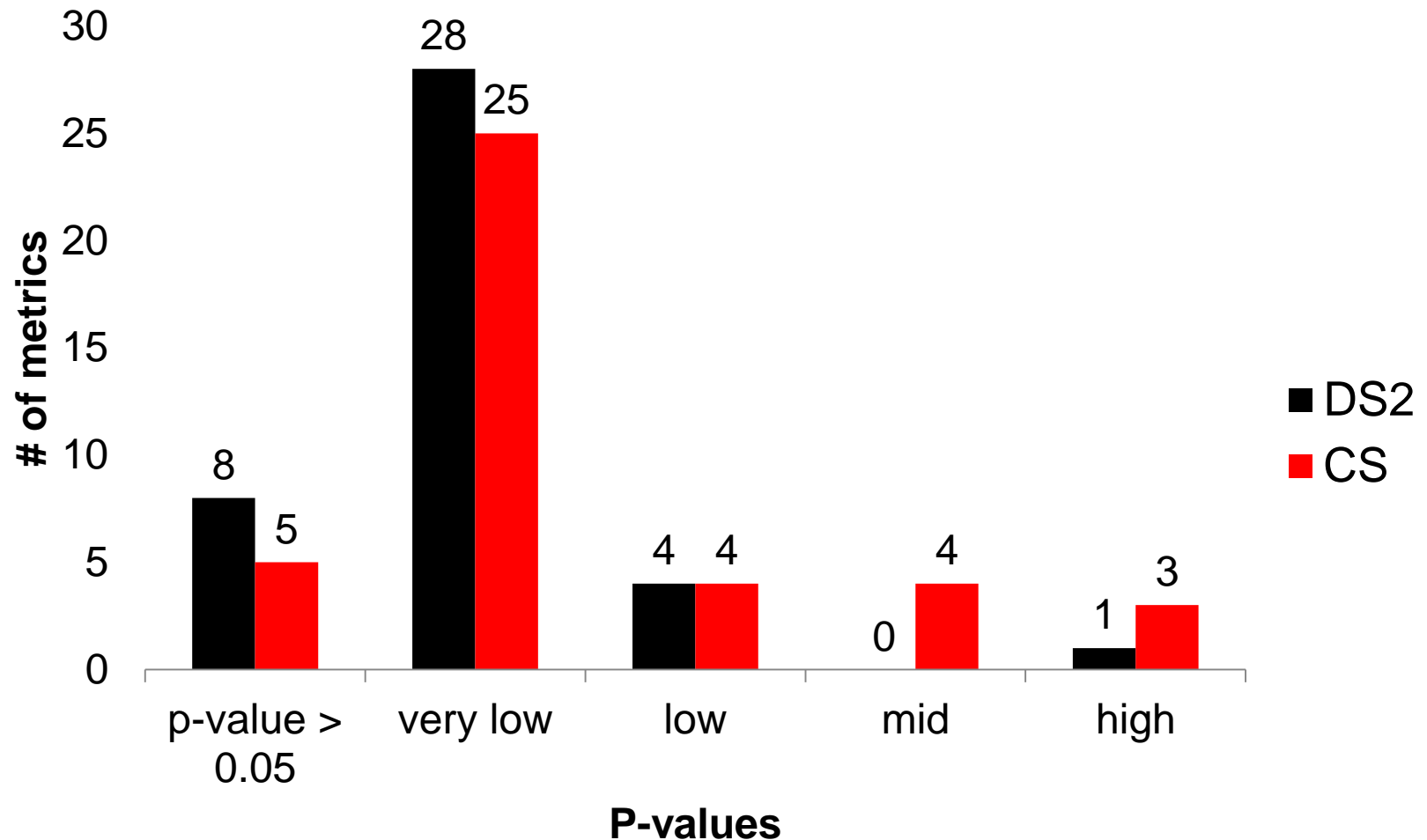
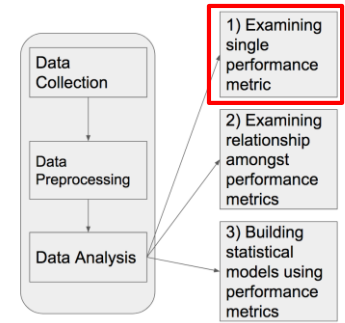
CloudStore

13

12

Spearman:

Most performance metrics do not have the same trend in virtual and physical environments.



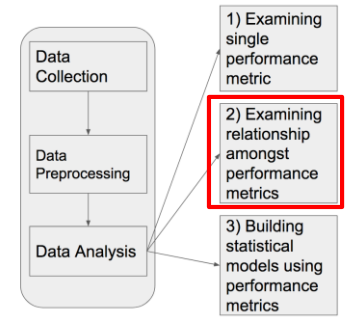
Research Questions

RQ1: Are the trend and distribution of a single performance metric similar across environments?

RQ2: To what extent does the **relationship between the performance metrics** change across environments?

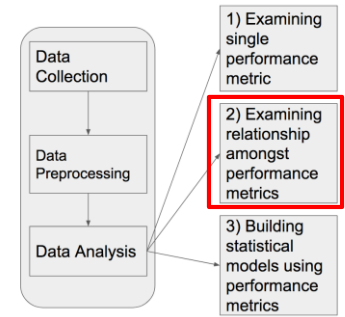
RQ3: Can statistical performance models be applied across virtual and physical environments?

RQ2: Approach

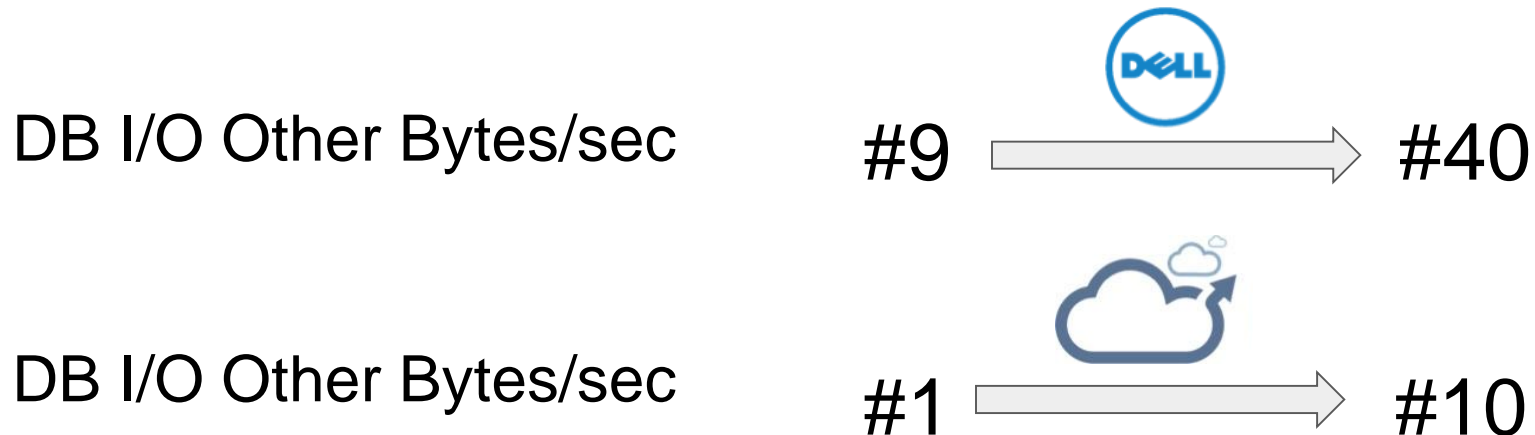


- Spearman's rank correlation coefficient:
 - a. We calculate against throughput.
 - b. We calculate the absolute difference for each pair of metrics. (represented by heatmaps)

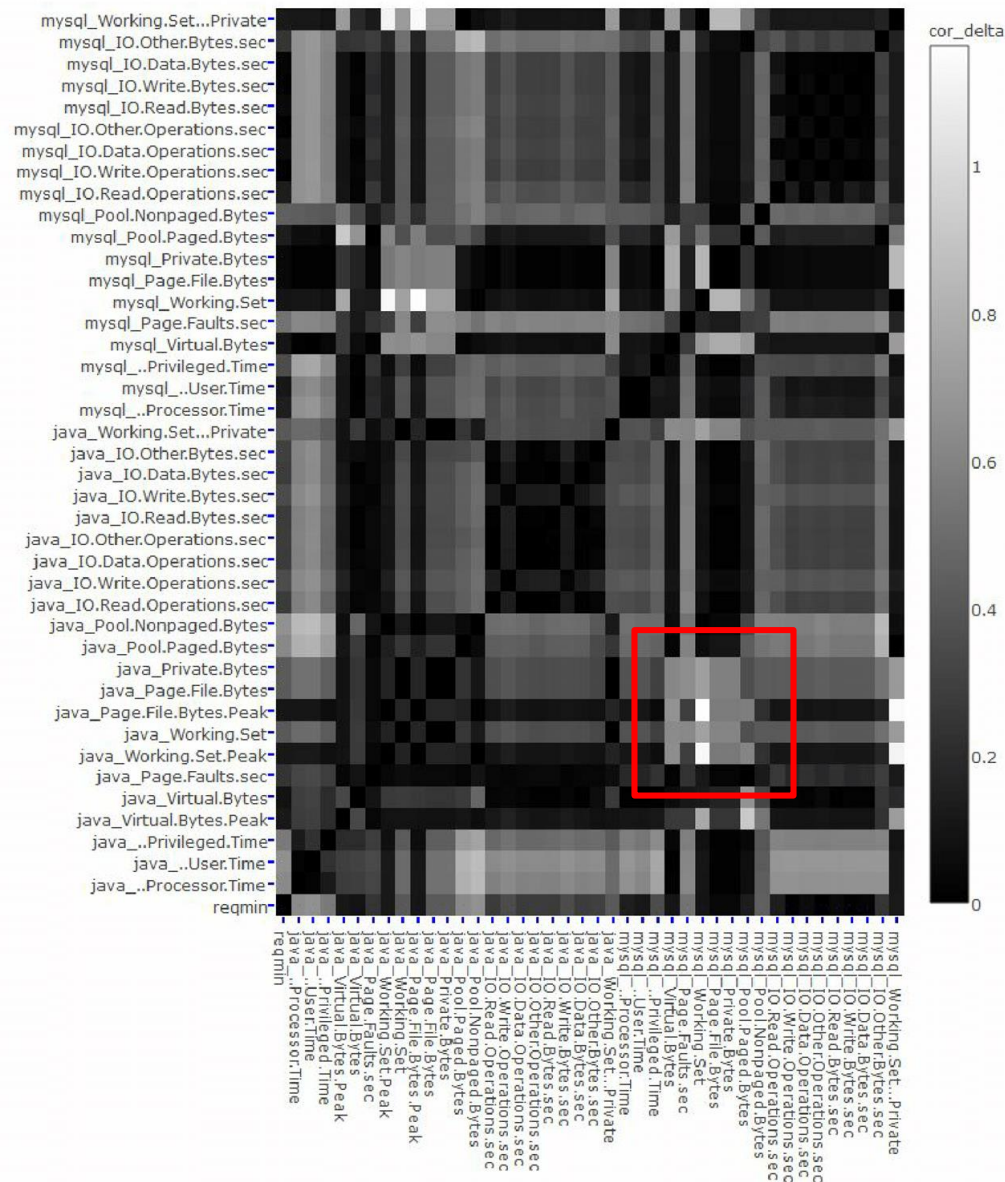
Spearman: The rank of the metrics changes across the environments.



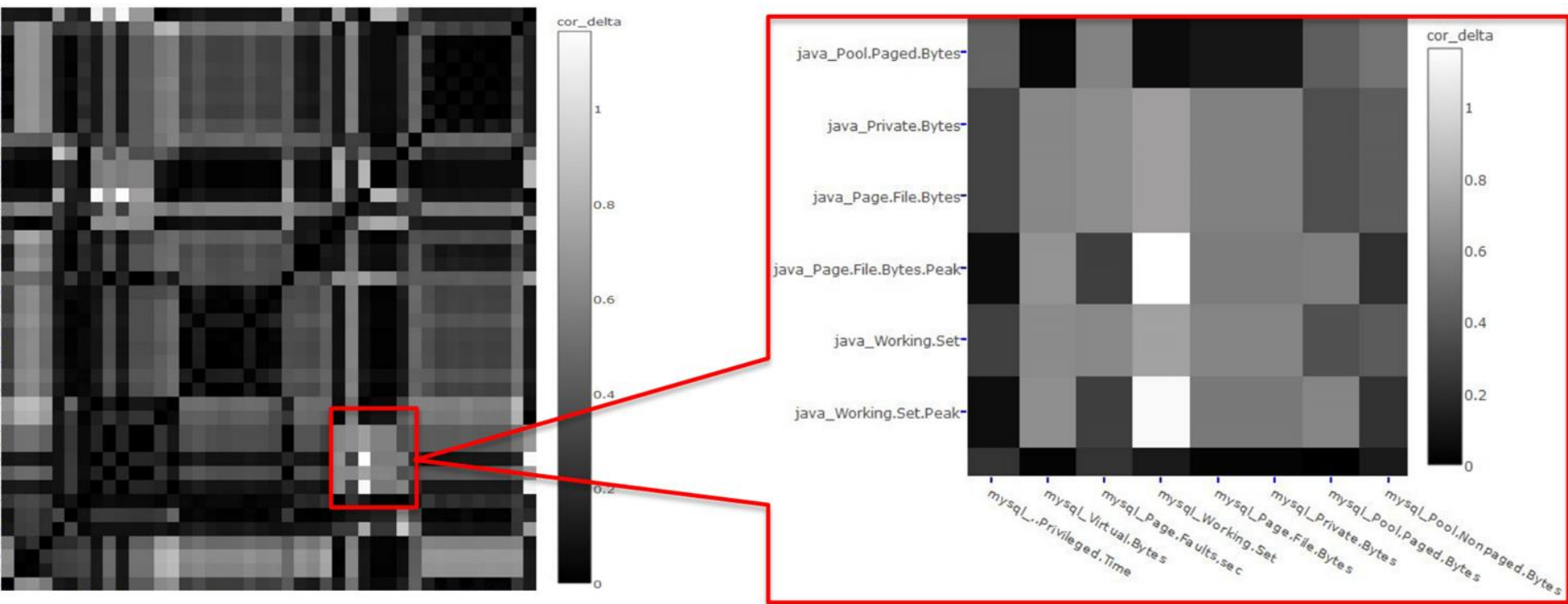
The *rank (vs. throughput)* of the metrics is **not the same** in the two environments.



CS: Heatmap



CS: There exist differences in correlation among the performance metrics from virtual and physical environments.



[Kraft, S. et al,
SIGSOFT, 2011.
Menon, A. et al.
ACM ICVEE, 2005]

Research Questions

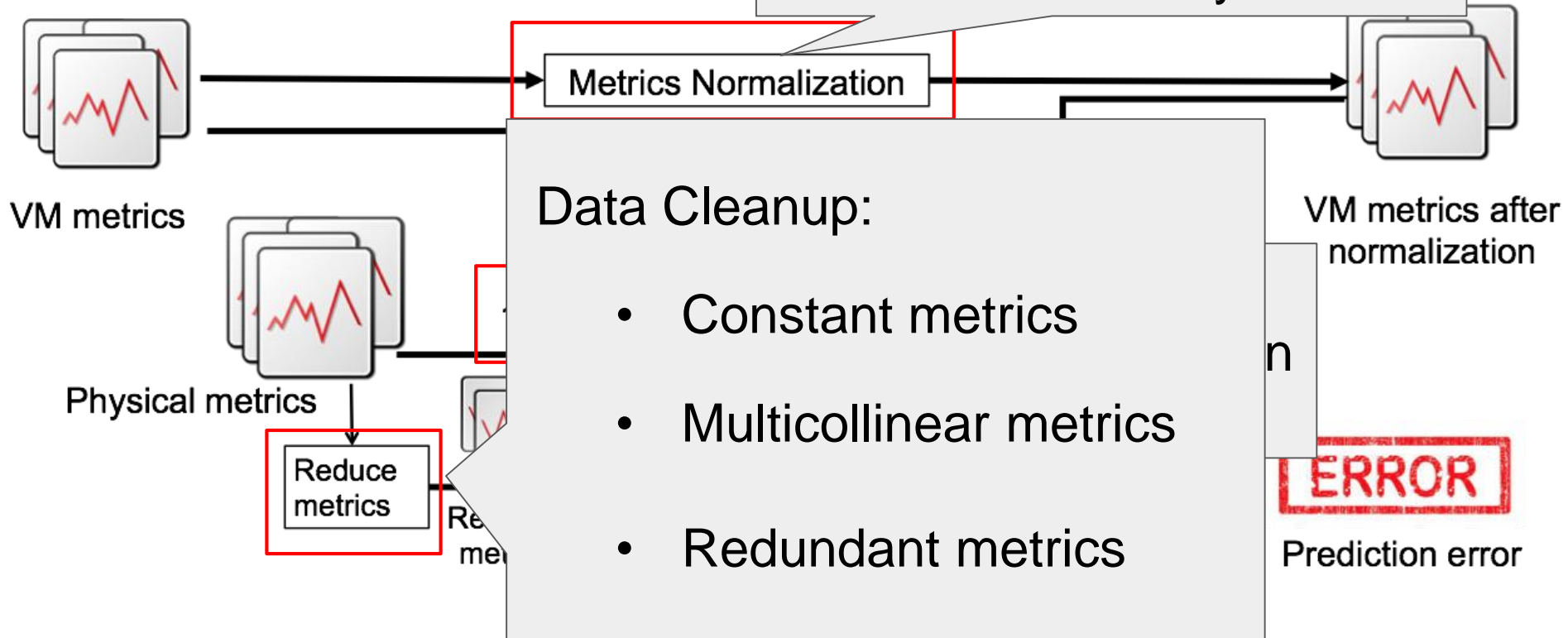
RQ1: Are the trend and distribution of a single performance metric similar across environments?

RQ2: To what extent does the relationship between the performance metrics change across environments?

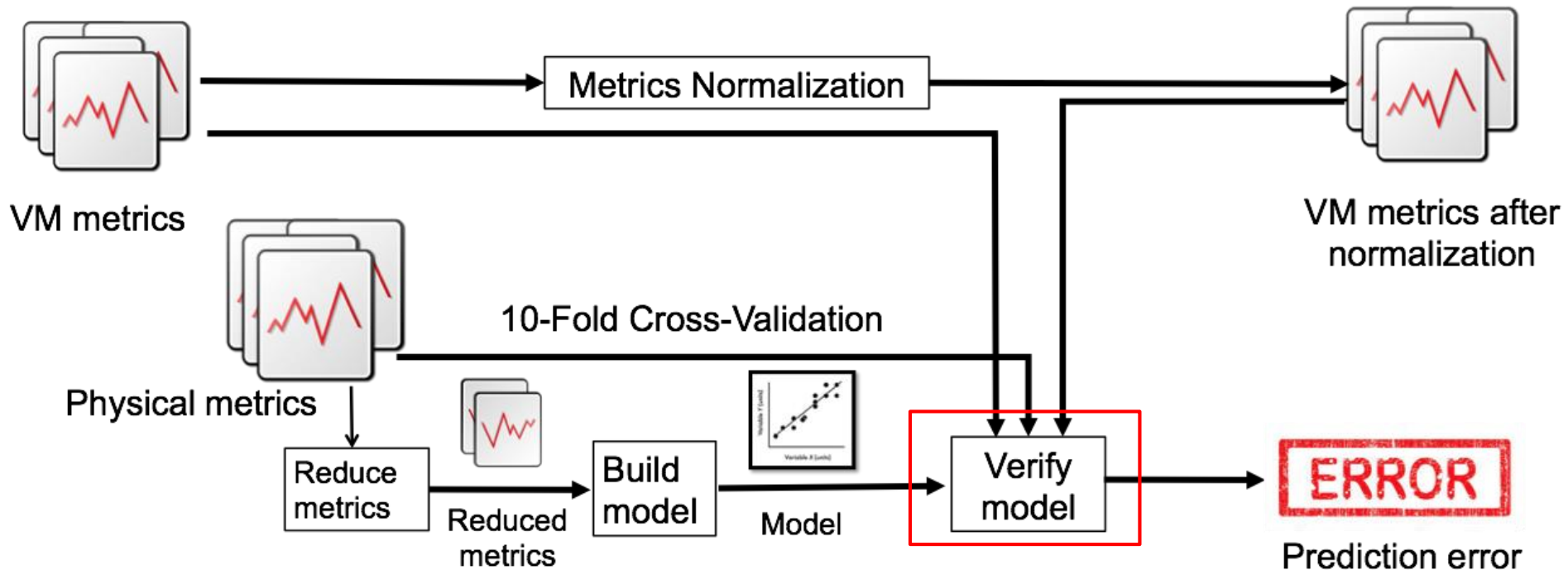
RQ3: Can **statistical performance models** be applied **across virtual and physical environments**?

RQ3: Approach -

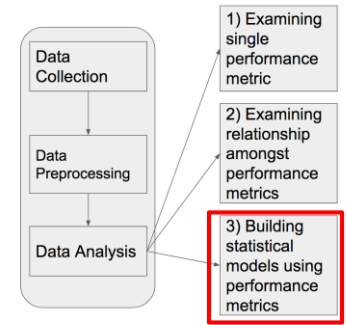
- Normalization by deviance
- Normalization by load



RQ3: Approach - Statistical modelling



RQ3: Model Verification



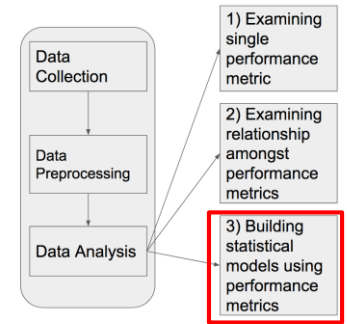
Two types of model verification:

1. Internal Validation

2. External Validation after:

- We normalize by deviance
- We normalize by load

RQ3: Model Verification

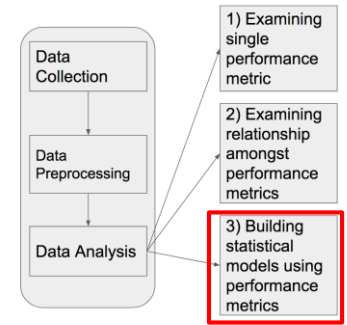


External Validation

- Normalization by deviance

$$M_{normalized} = \frac{M - \tilde{M}}{MAD(M)}$$

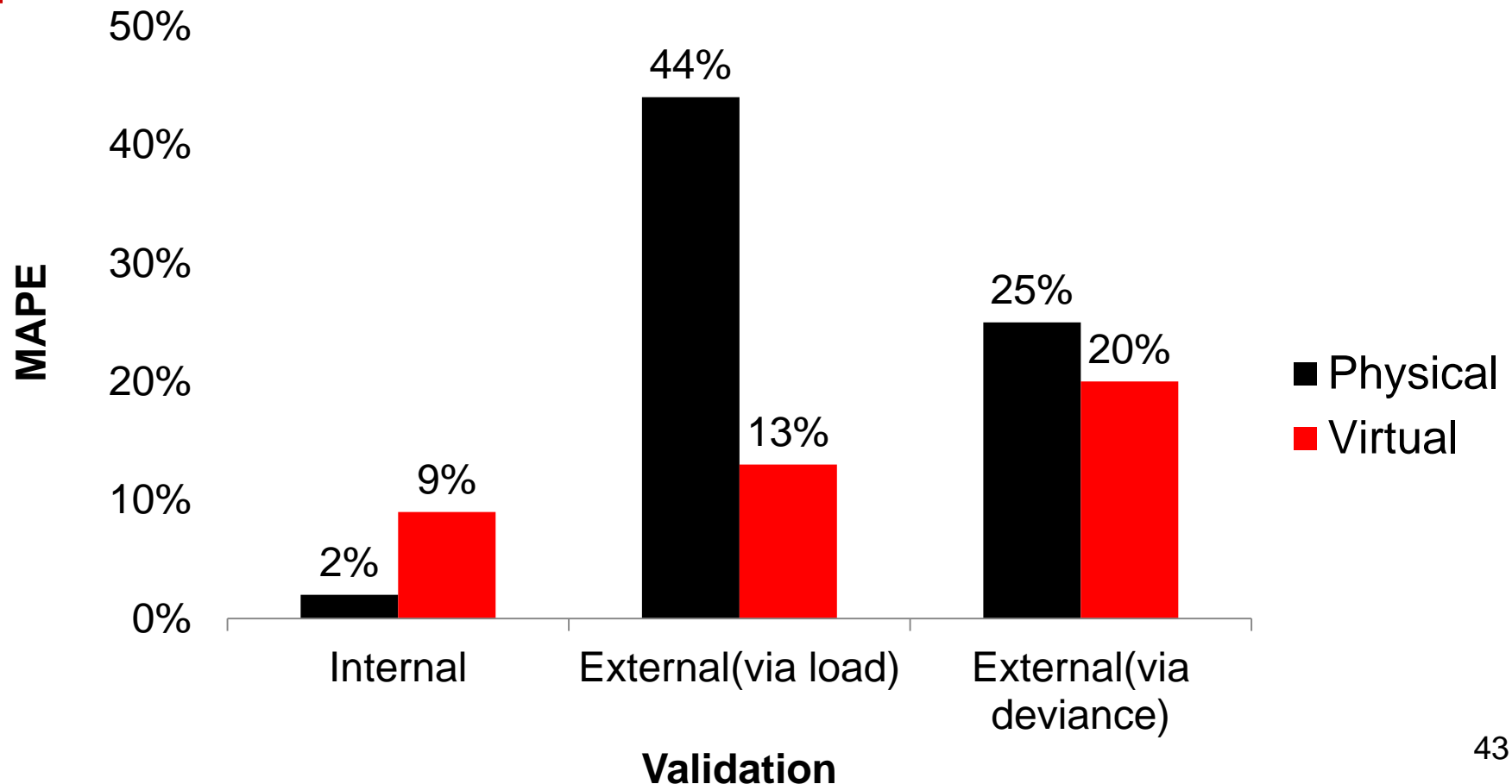
RQ3: Model Verification



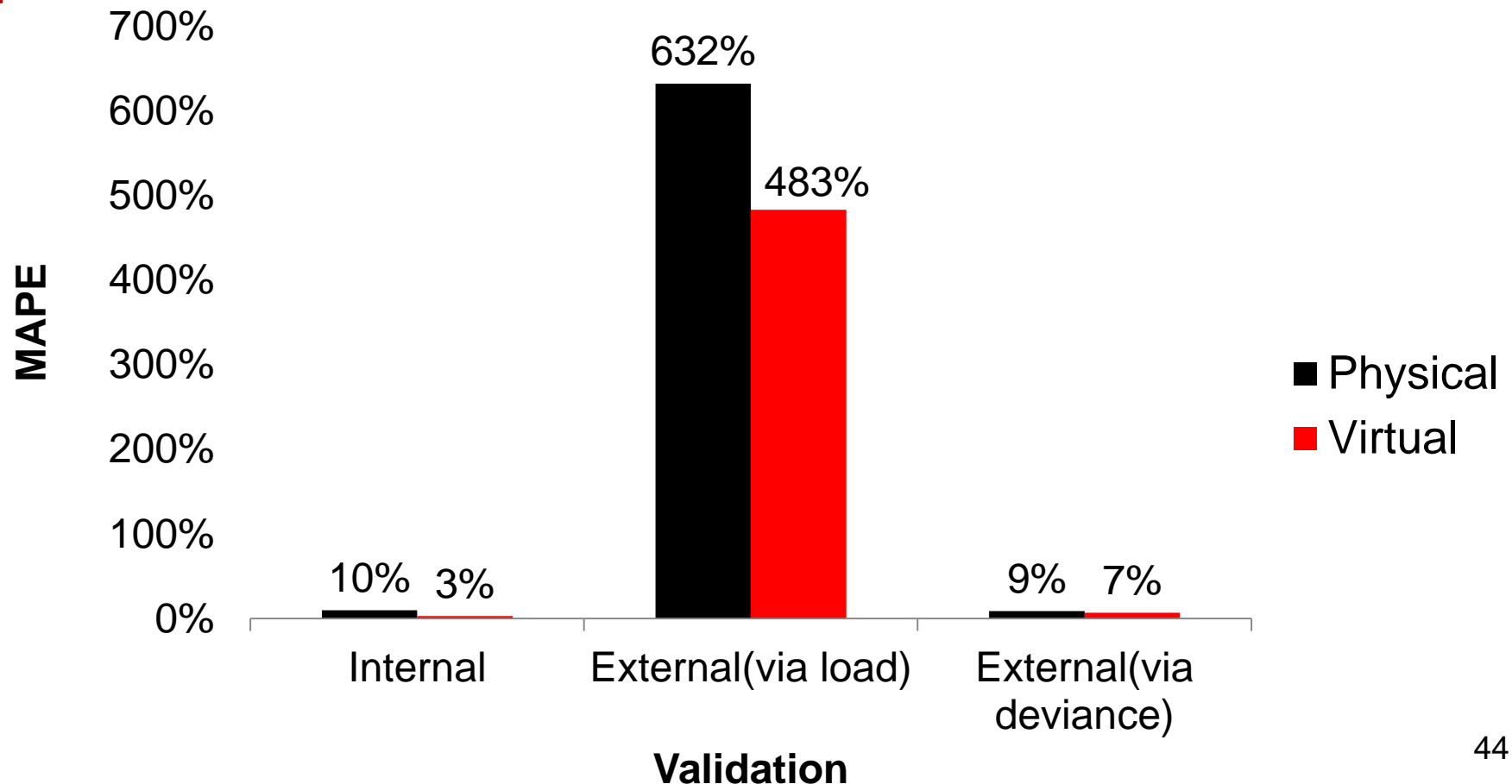
External Validation

- Normalization by deviance
- Followed by a min-max normalization to avoid calculation errors for negative values.

Discrepancy - DS2: We find that the statistical models built by performance testing results in an environment cannot advocate for the other environment due to discrepancies present.



Discrepancy - CS: We find that the statistical models built by performance testing results in an environment cannot advocate for the other environment due to discrepancies present.



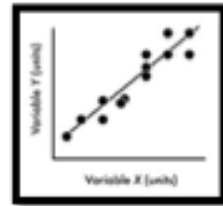
Examining the impact from other factors on our results

1. Instability of the virtual environment
2. Impact of the specific virtual machine software
3. Impact of allocated resources

1. Instability of the virtual environment



Repeat tests

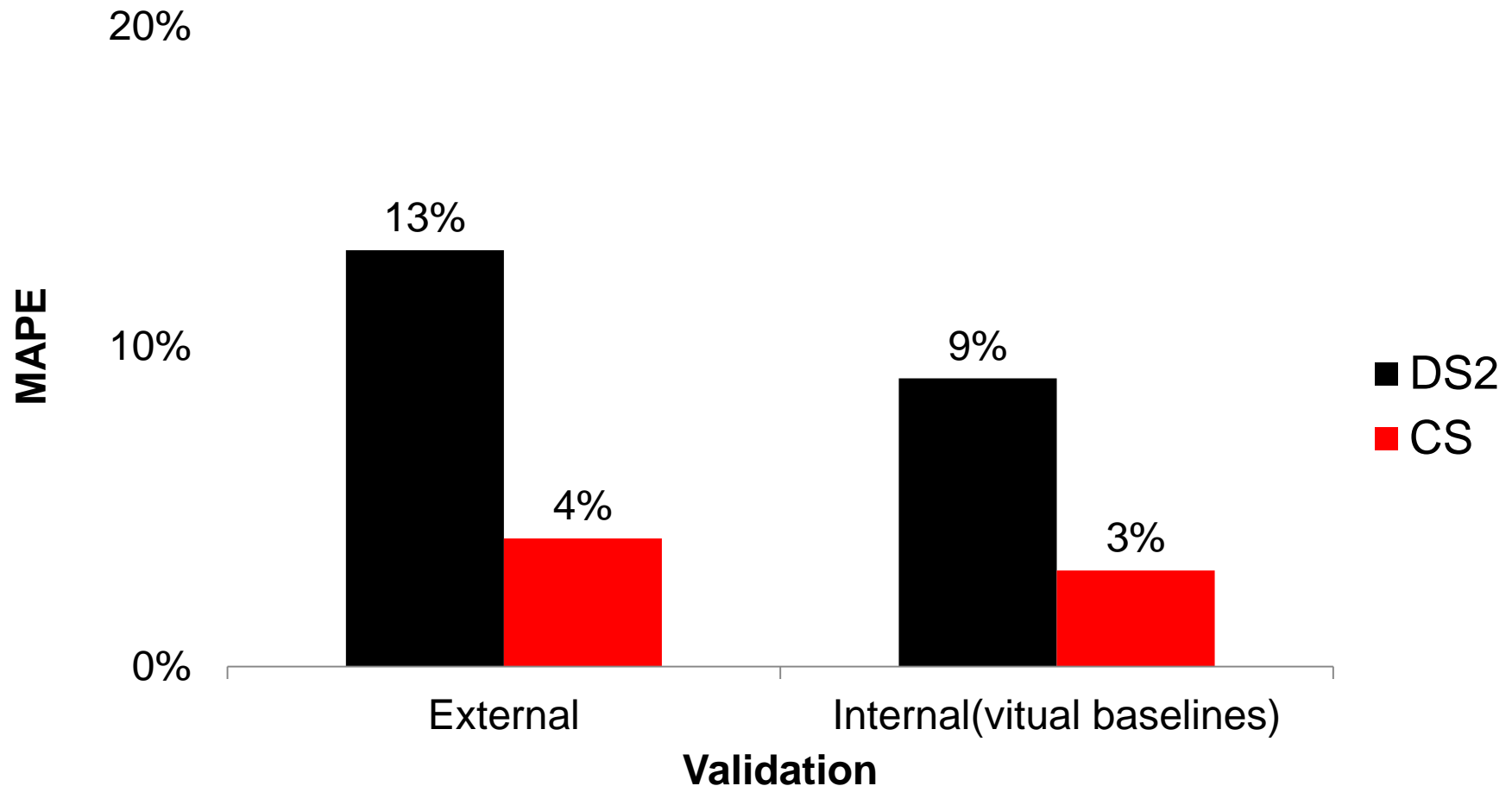


Build models

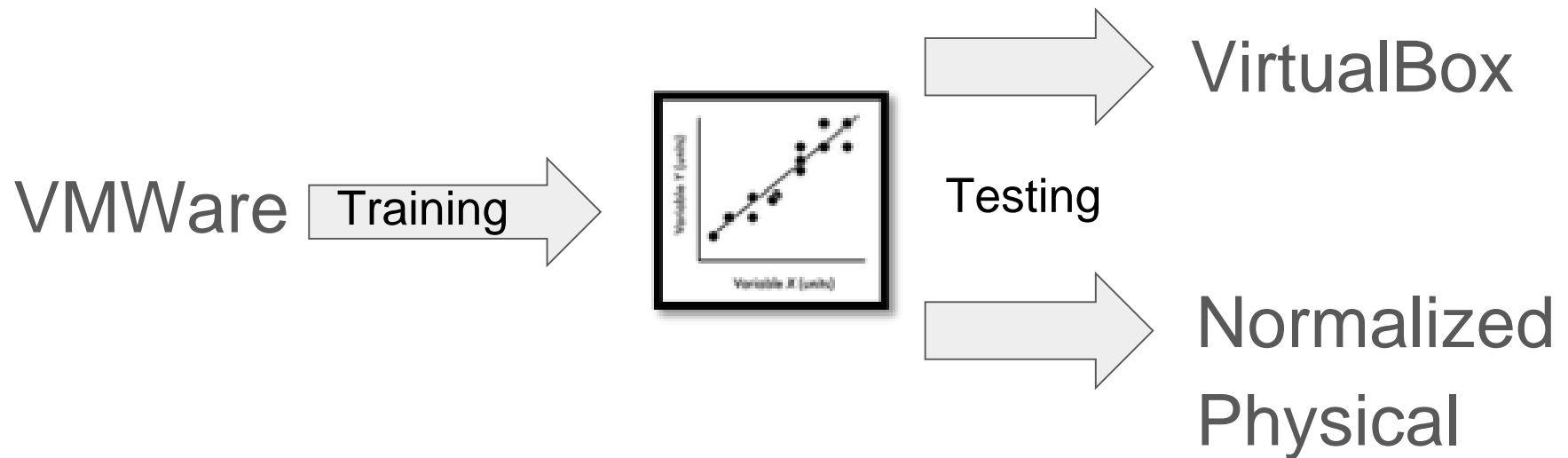


Compare
externally

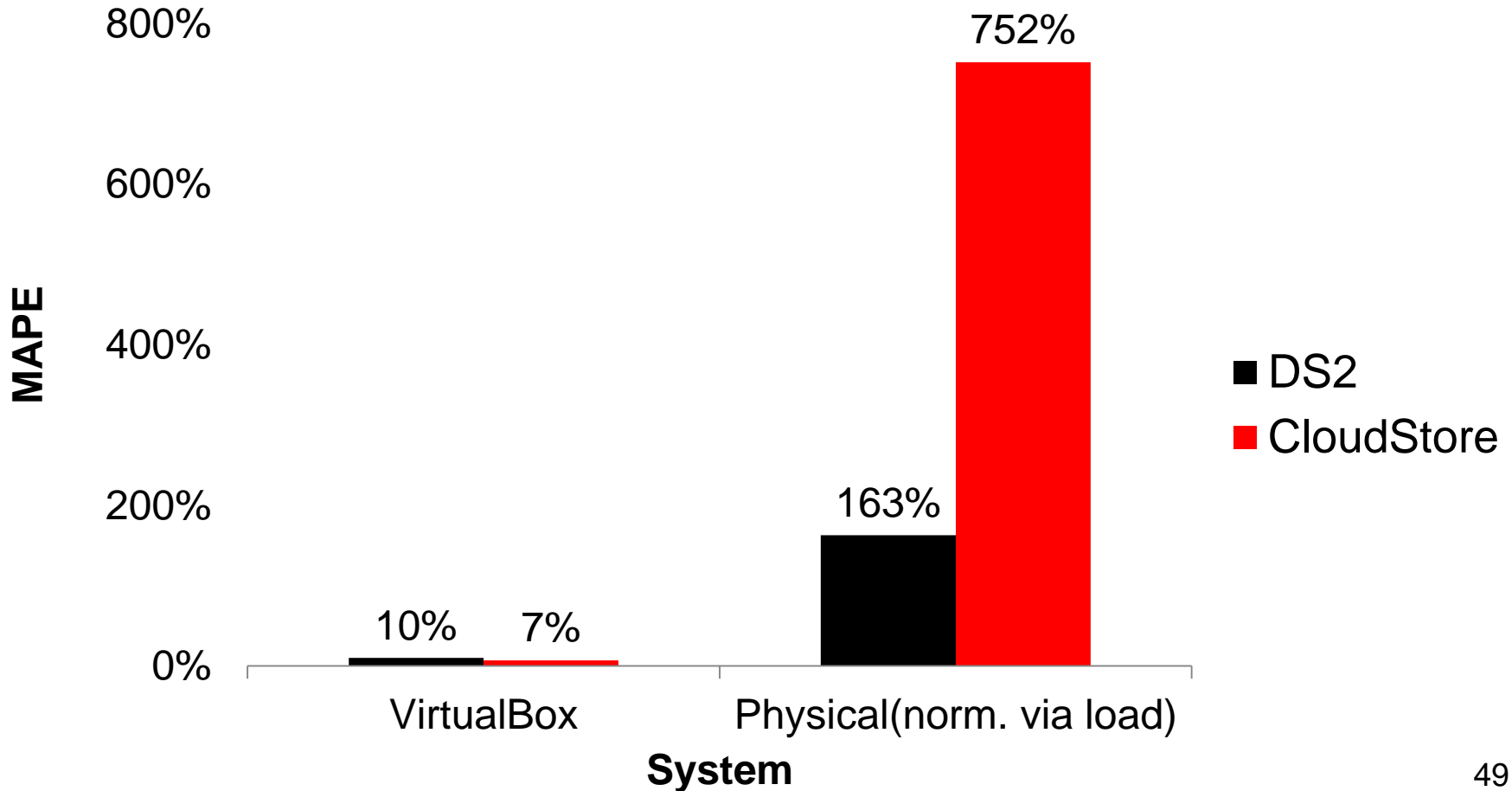
1. Performance testing results from the virtual environments are stable.



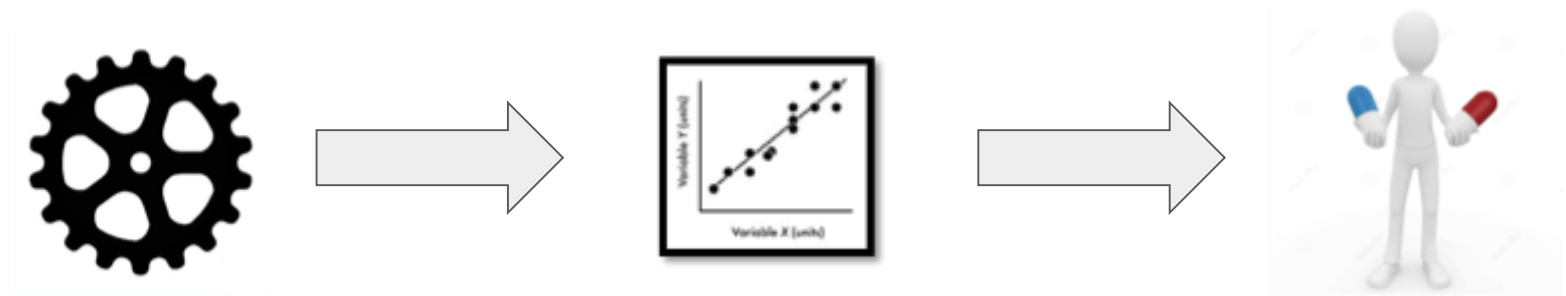
2. Impact of the specific virtual machine software



2. Discrepancy observed during our experiment also exists with the virtual environments that are set up with VMWare.



3. Impact of allocated resources



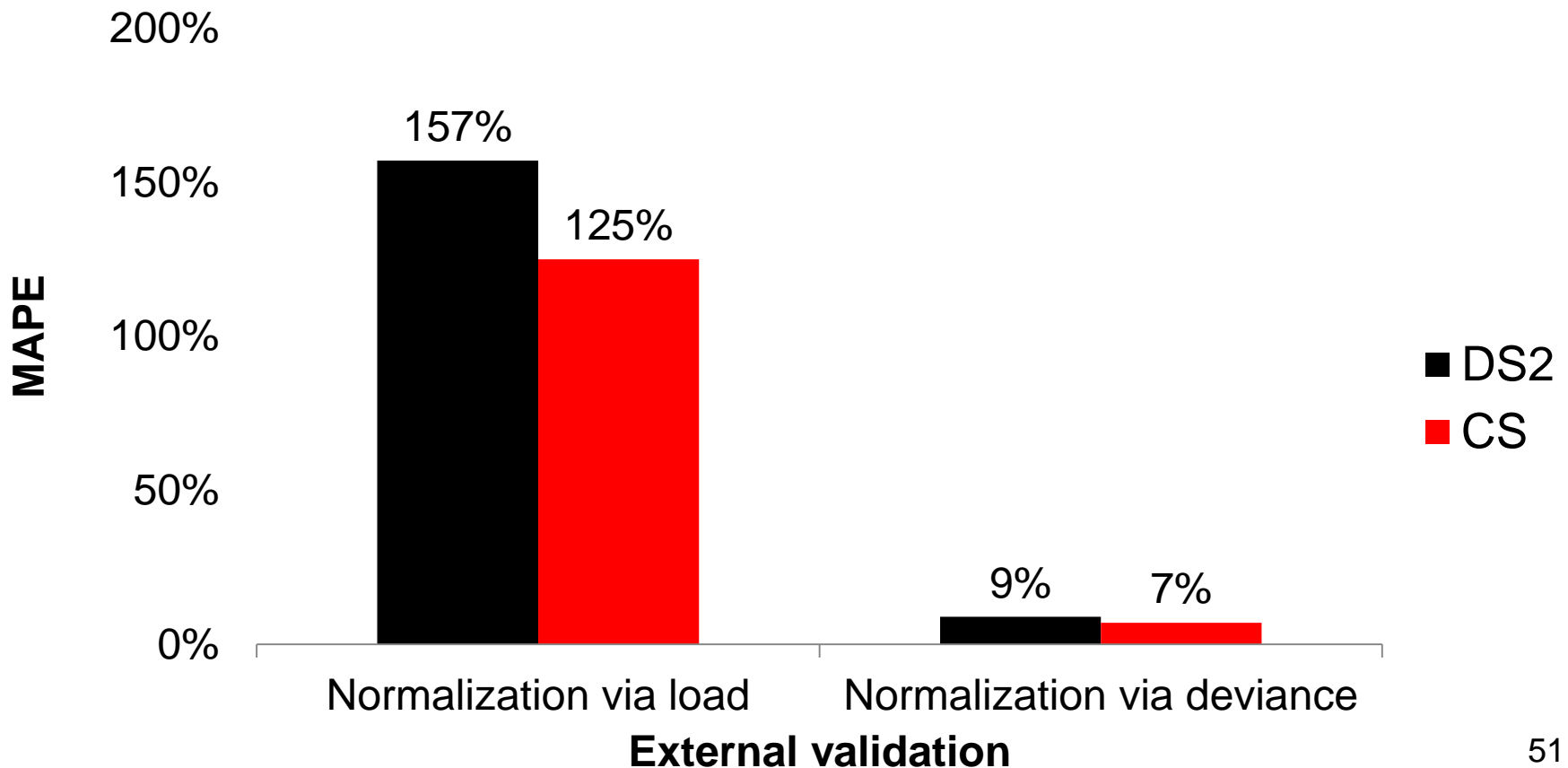
Modify
configuration

Build models

Compare
externally

We increase the resources to
3 cores and 5GB memory.

3. Our findings still hold when the allocated resources are changed and this change has minimal impact on the results of our case studies.



Contributions

This is the ***first research attempt*** to evaluate the ***discrepancy*** between performance testing results in virtual and physical environments.

We identified the ***performance metrics that contribute the most to the discrepancies.***

We find that **normalizing performance metrics based on deviance** may ***reduce*** the discrepancy

Implications

Practitioners ***cannot assume a straight forward overhead.***

Practitioners should always verify whether the ***inconsistency of correlations between performance metrics*** are due to virtual environments.

Normalization by deviance may minimize such discrepancy.

Summary

How is a performance regression detected?

What if you need to modify your environment?

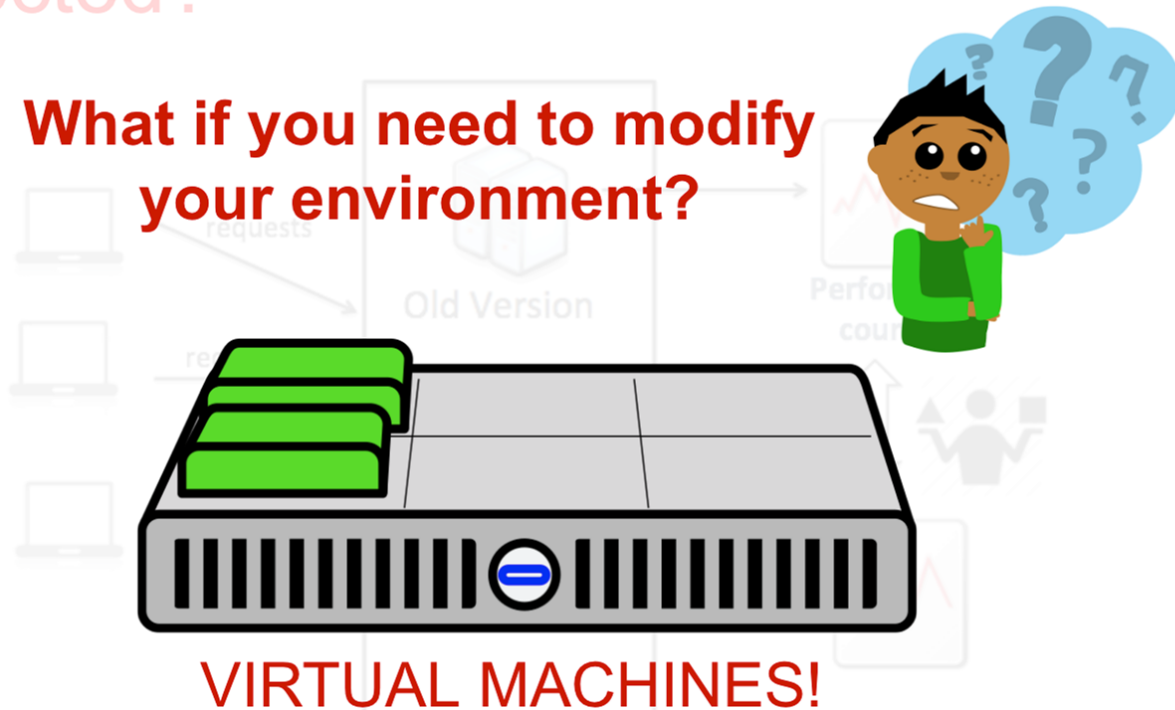
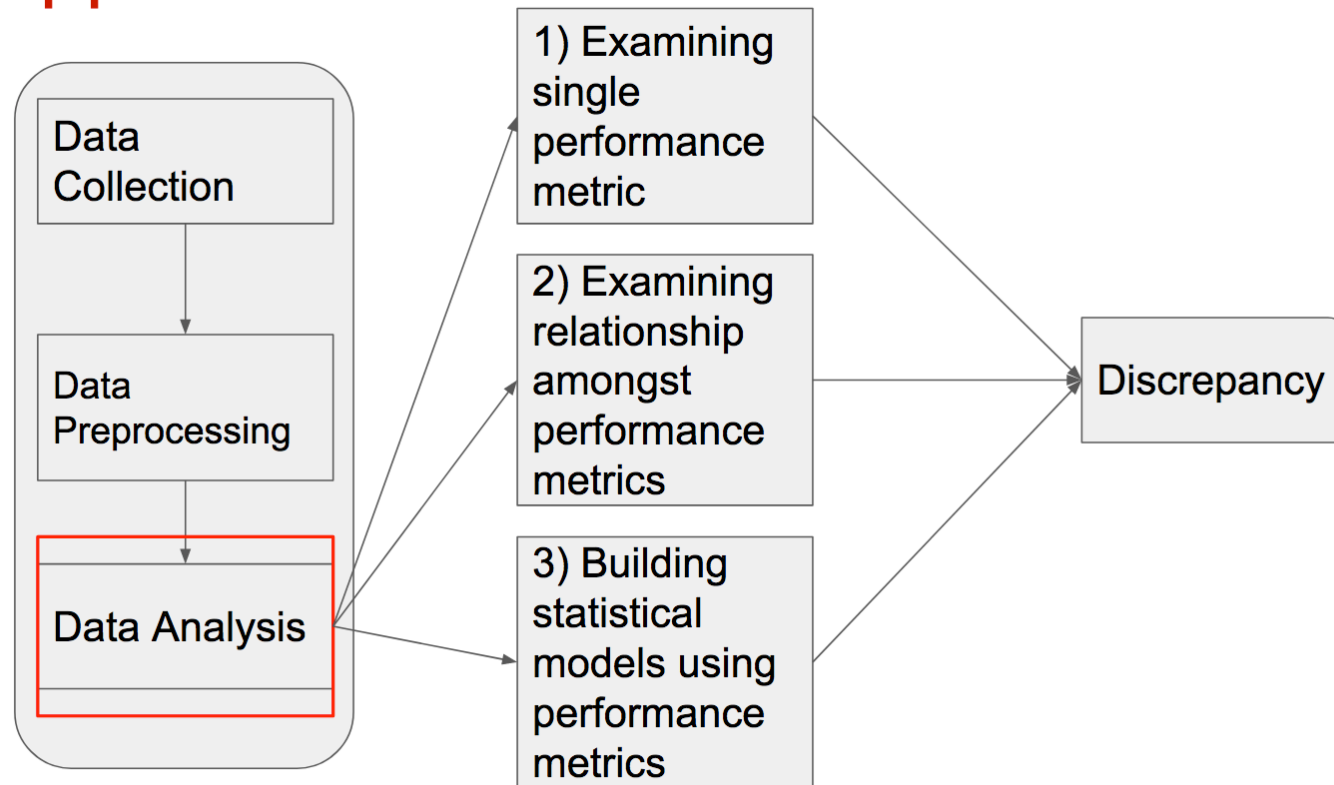


Image:
<https://openclipart.org/detail/191766/question-guy>
<https://openclipart.org/detail/167198/two-virtual-machines>

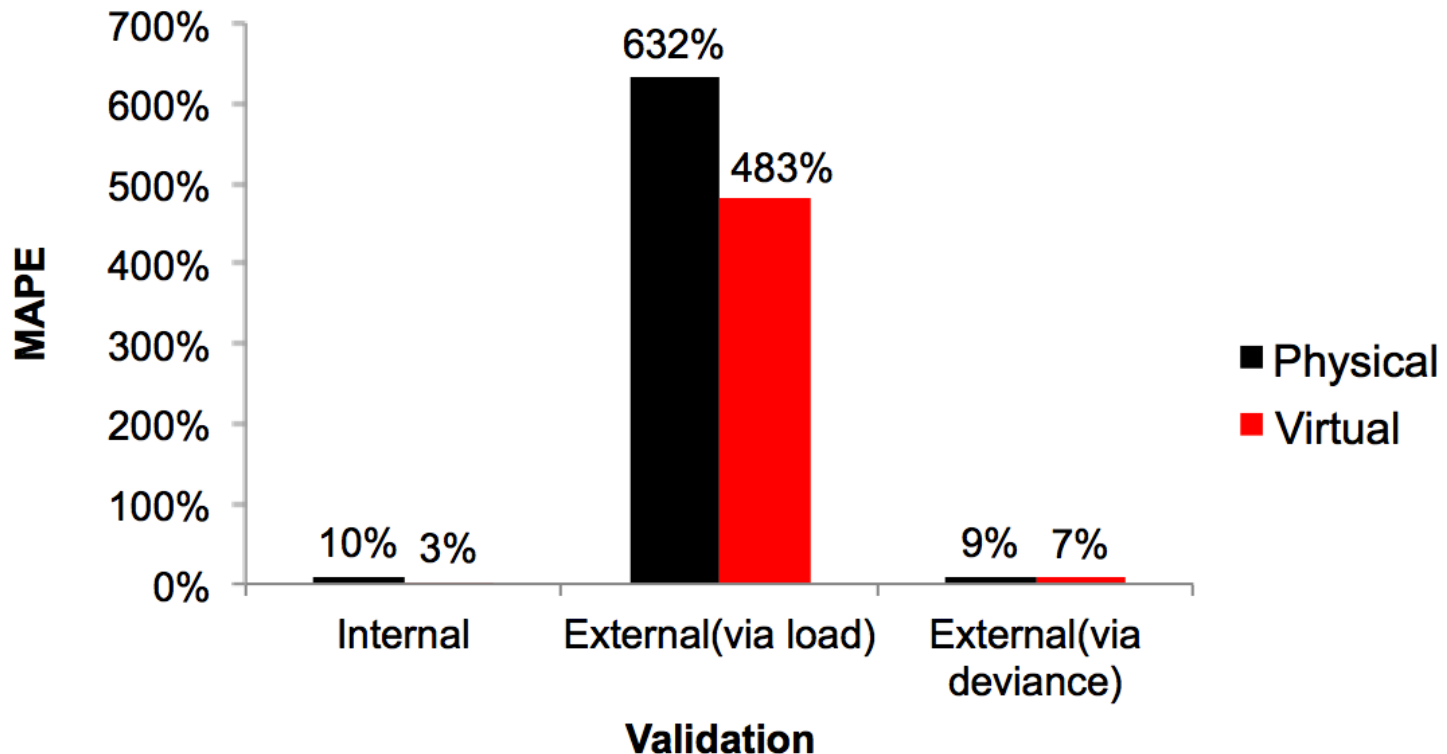
Summary

Approach



Summary

Discrepancy - CS: We find that the statistical models built by performance testing results in an environment cannot advocate for the other environment due to discrepancies present.



Summary

Contributions

- This is the ***first research attempt*** to evaluate the ***discrepancy*** between performance testing results in virtual and physical environments.
- We identified the ***performance metrics that contribute the most to the discrepancies.***
- We find that **normalizing performance metrics based on deviance** may ***reduce*** the discrepancy

Summary

How is a performance regression detected?

What if you need to modify your environment?

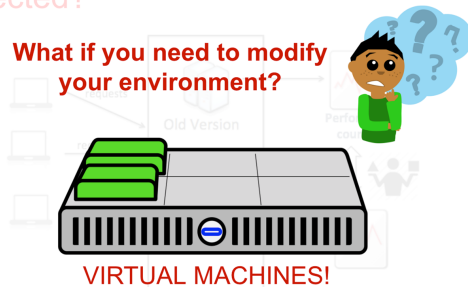
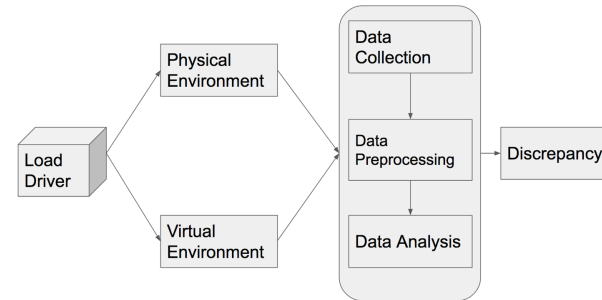


Image:
https://openstax.org/detail/181796/question-guy
https://openstax.org/detail/181796/how-virtual-machines

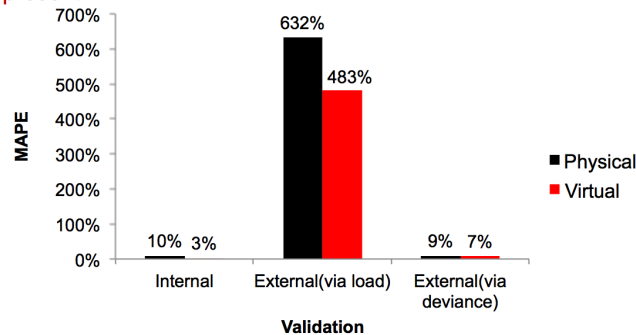
6

Approach



14

Discrepancy - CS: We find that the statistical models built by performance testing results in an environment cannot advocate for the other environment due to discrepancies present.



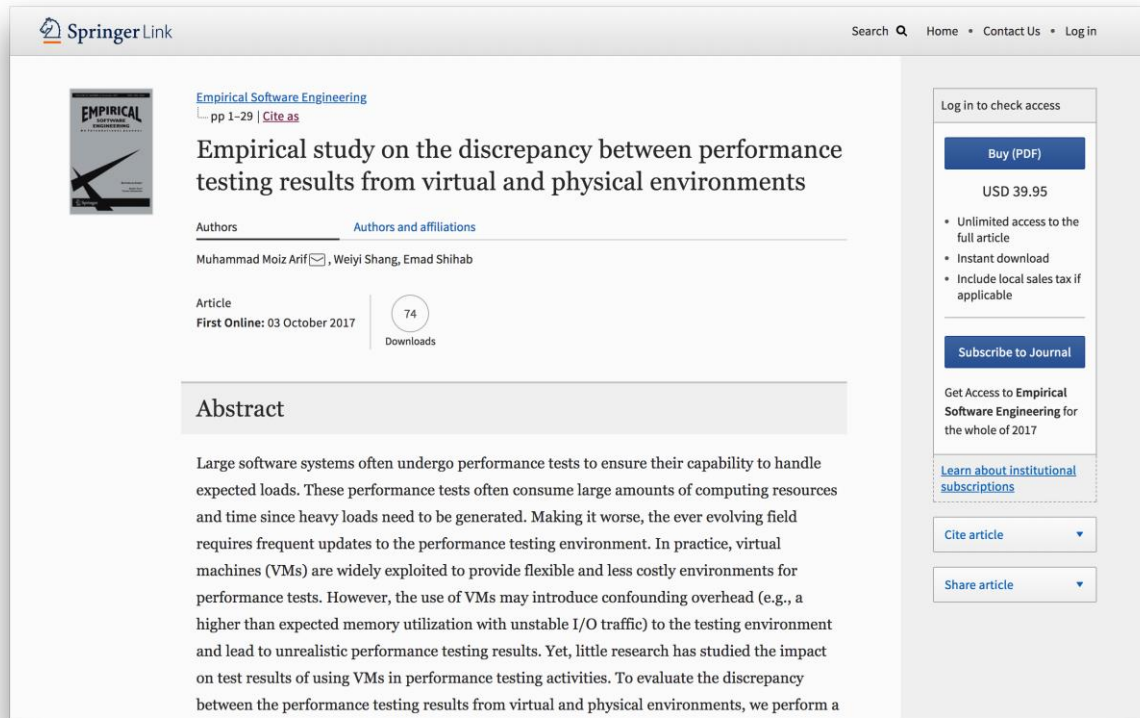
Contributions

- This is the **first research attempt** to evaluate the **discrepancy** between performance testing results in virtual and physical environments.
- We identified the **performance metrics that contribute the most to the discrepancies**.
- We find that **normalizing performance metrics based on deviance** may **reduce** the discrepancy

51

Publications

Arif, M.M., Shang, W. & Shihab, E. Empir
Software Eng (2017)



The screenshot shows the SpringerLink interface for the article "Empirical study on the discrepancy between performance testing results from virtual and physical environments". The page includes the journal cover, title, authors (Muhammad Moiz Arif, Weiyl Shang, Emad Shihab), and an abstract. The abstract discusses the challenges of performance testing in virtual environments and the need for research to evaluate the discrepancy between virtual and physical environments. The right sidebar contains options to buy the PDF for USD 39.95, subscribe to the journal, and links to cite or share the article.

SpringerLink

Search Home Contact Us Log in

Empirical Software Engineering
pp 1–29 | Cite as

Empirical study on the discrepancy between performance testing results from virtual and physical environments

Authors Authors and affiliations
Muhammad Moiz Arif✉, Weiyl Shang, Emad Shihab

Article
First Online: 03 October 2017

74 Downloads

Abstract

Large software systems often undergo performance tests to ensure their capability to handle expected loads. These performance tests often consume large amounts of computing resources and time since heavy loads need to be generated. Making it worse, the ever evolving field requires frequent updates to the performance testing environment. In practice, virtual machines (VMs) are widely exploited to provide flexible and less costly environments for performance tests. However, the use of VMs may introduce confounding overhead (e.g., a higher than expected memory utilization with unstable I/O traffic) to the testing environment and lead to unrealistic performance testing results. Yet, little research has studied the impact on test results of using VMs in performance testing activities. To evaluate the discrepancy between the performance testing results from virtual and physical environments, we perform a

Log in to check access

Buy (PDF)

USD 39.95

- Unlimited access to the full article
- Instant download
- Include local sales tax if applicable

Subscribe to Journal

Get Access to Empirical Software Engineering for the whole of 2017

Learn about institutional subscriptions

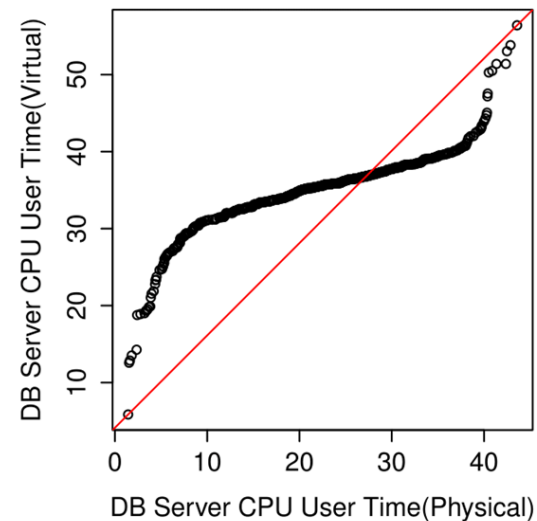
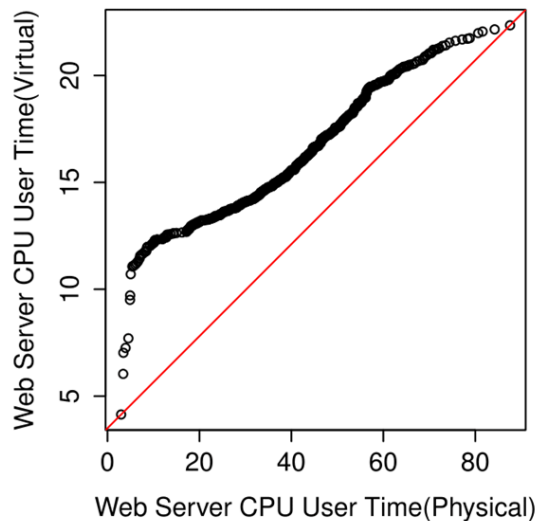
Cite article

Share article

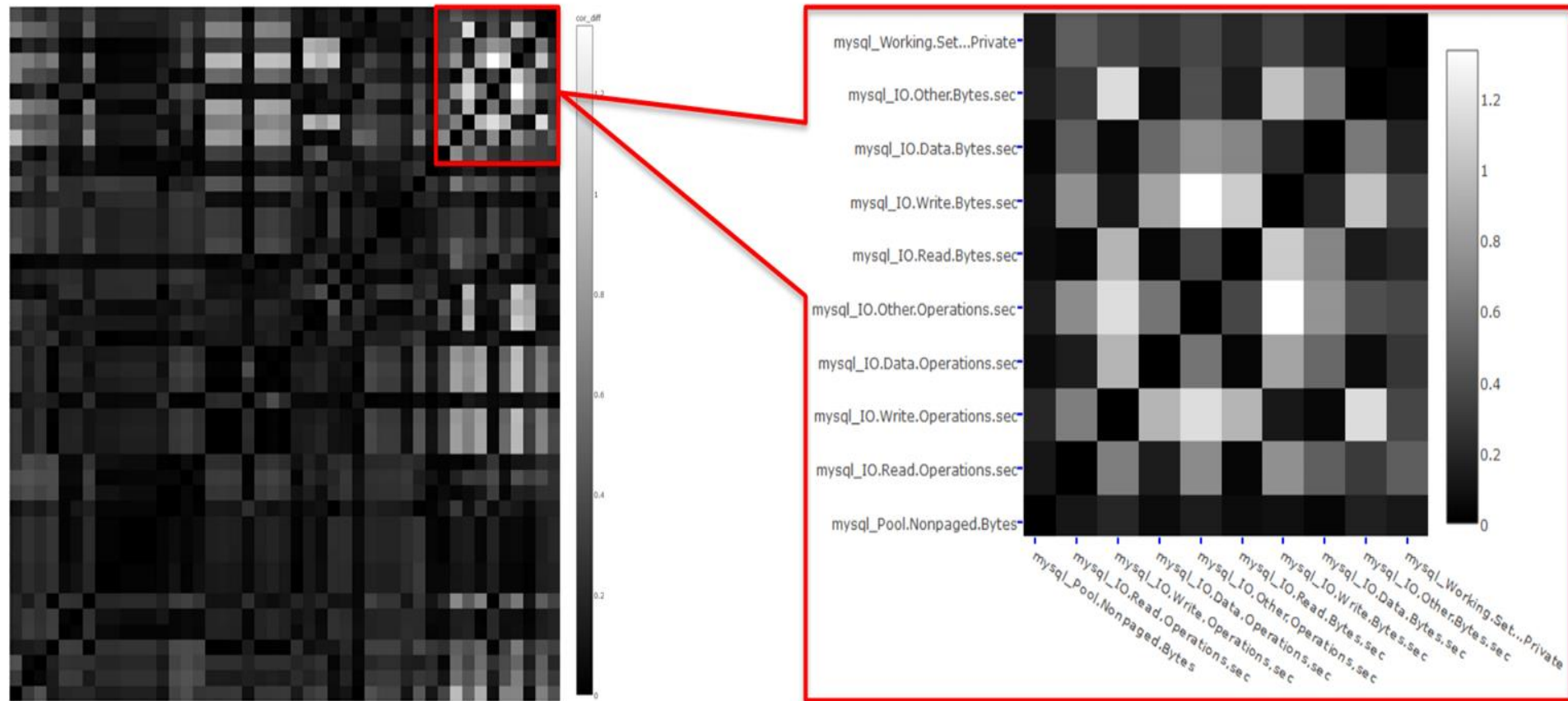
https://users.encs.concordia.ca/~shang/pubs/emse_moiz.pdf

m.moizarif@gmail.com

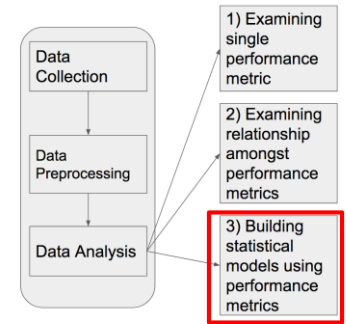
CS: Most performance metrics do not follow the same shape of the distribution in virtual and physical environments.



DS2: Heatmap



Model Verification

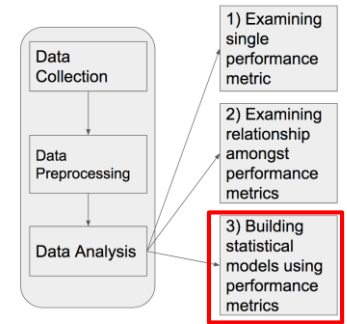


Internal Validation

Ten-fold Cross validation

k=1	Test	Train	Train	Train	Train	Train	Train	Train	Train	Train
k=2	Train	Test	Train	Train	Train	Train	Train	Train	Train	Train
k=3	Train	Train	Test	Train	Train	Train	Train	Train	Train	Train
...
k=10	Train	Train	Train	Train	Train	Train	Train	Train	Train	Test

Model Verification



External Validation

Normalization by load

$$throughput_p = \alpha_p \times M_p + \beta_p$$

$$throughput_v = \alpha_v \times M_v + \beta_v$$

$$M_{normalized} = \frac{(\alpha_v \times M_v) + \beta_v - \beta_p}{\alpha_p}$$

Discrepancy: Absolute % Error

Median Absolute Percentage Error (*MAPE*)

Throughput: 100 requests/minute

Predicted: 110 requests/minute

$$MAPE = \left(\frac{|110 - 100|}{100} \right) = 0.1$$

Limitations

More case studies

Quality of performance tests and recorded values.

Other types of workload variation.

Hint of overfitting of models.

Choice of intervals (for e.g. 9 hours, 10 seconds)

The complete system: user's point of view.

Future Work

Reproducing known performance regressions in heterogeneous environments.

Replicating our experiments in cloud environments.

Designing automating techniques.