

Rapid and Robust Impact Assessment of Software Changes in Large Internet-based Services

Shenglin Zhang, Ying Liu, Dan Pei
Yu Chen, Xianping Qu, Shimin Tao, Zhi Zang



清華大學
Tsinghua University



Internet-based Services

- Search



- Shopping



- Social



- Portal



- Video



Software Change: Software Upgrade or Configuration Change

- Software upgrade



Introduce
new feature



Improve
performance



Fix bugs

Software Change: Software Upgrade or Configuration Change

- Software upgrade



Introduce
new feature



Improve
performance



Fix bugs

- Configuration change

- e.g., traffic switching for load balancing reasons

Software Change: Software Upgrade or Configuration Change

- Software upgrade



Introduce
new feature



Improve
performance



Fix bugs

- Configuration change

- e.g., traffic switching for load balancing reasons

- Occurs frequently

- 10K+ per day in Baidu

Impact of Erroneous Software Upgrades

2012.10, Google

Google Apps Incident Report

Gmail Partial Outage - December 10, 2012

Prepared for Google Apps customers

- An update to Google's load balancing software
- Poor performance to Gmail for 18 minutes

The following is the incident report for the Gmail partial outage on December 10, 2012. We understand the impact of this outage and apologize to our customers.

Issues

For a period of time, some users experienced poor performance when accessing Gmail. The number of users affected was approximately 10% of users. The root cause of the issue was a software update to Google's load balancing software.

Actions and Root Cause

Background: The load balancing software routes users' requests to Google data centers around the world for processing and serving content, such as search results and email.

Between 8:45 AM PT and 9:13 AM PT, a routine update to Google's load balancing software was rolled out to production. A bug in the software update caused it to incorrectly interpret a portion of Google data centers as being unavailable. The Google load balancers have a failsafe mechanism to prevent this type of failure from causing Google-wide service degradation, and they continued to route user traffic. As a result, most Google services, such as Google Search, Maps, and AdWords, were unaffected. However, some services, including Gmail, that require specific data center information to efficiently route users' requests, experienced a partial outage.

Impact of Erroneous Software Upgrades

2012.10, Google

Google Apps Incident Report

Gmail Partial Outage - December 10, 2012

Prepared for Google Apps customers

- An update to Google's load balancing software
- Poor performance to Gmail for 18 minutes

2014.11, Microsoft Azure

Update on Azure Storage Service Interruption

WEDNESDAY, NOVEMBER 19, 2014



JASON ZANDER
CVP, Microsoft Azure Team

- A performance update to Azure Storage
- Reduced capacity across services utilizing Azure Storage

Wednesday, November 19,

As part of a performance update to Azure Storage, an issue was discovered that resulted in reduced capacity across services utilizing Azure Storage, including Virtual Machines, Visual Studio Online, Websites, Search and other Microsoft services. Prior to applying the performance update, it had been tested over several weeks in a subset of our customer-facing storage service for Azure Tables. We typically call this "flighting," as we work to identify issues before we broadly deploy any updates. The flighting test demonstrated a notable performance improvement and we proceeded to deploy the update across the storage service. During the rollout we discovered an issue that resulted in storage blob front ends going into an infinite loop, which had gone undetected during flighting. The net result was an inability for the front ends to take on further traffic, which in turn caused other services built on top to experience issues.

Impact of Erroneous Configuration Changes

2014.1, Dropbox

Outage post-mortem

Akhil Gupta | January 13, 2014

0 [in](#) 21 [8+](#)

On Friday, we had a planned maintenance scheduled to upgrade the OS on some of our machines. During this process, the upgrade script checks to make sure there is no active data on the machine before installing the new OS.

When the script ran, it found that some machines had active data and therefore did not upgrade them. This resulted in a small number of machines being left in an inconsistent state, which caused the site to go down.

When the site went down, we realized that the upgrade script was not working as expected. We quickly rolled back the changes and restored the site to its previous state.

What happened? The upgrade script was not working as expected. It was supposed to check for active data on the machine before installing the new OS, but it failed to do so.

We use thousands of machines for redundancy. Each machine has one master and two replica machines for redundancy. We use a distributed database to store data, and incremental data backups and store them in a separate environment.

On Friday at 5:30 PM PT, we had a planned maintenance scheduled to upgrade the OS on some of our machines. During this process, the upgrade script checks to make sure there is no active data on the machine before installing the new OS.

A subtle bug in the script caused the command to reinstall a small number of active machines. Unfortunately, some master-replica pairs were impacted which resulted in the site going down.

Impact of Erroneous Configuration Changes

2014.1, Dropbox

Outage post-mortem

Akhil Gupta | January 13, 2014

On Friday
and run

back up
day.

on as

Wh

We use thousands of machines. Each machine has one master and two replica machines for redundancy. We use a distributed database and incremental data backups and store them in a separate environment.

On Friday at 5:30 PM PT, we had a planned maintenance scheduled to upgrade the OS on some of our machines. During this process, the upgrade script checks to make sure there is no active data on the machine before installing the new OS.

A subtle bug in the script caused the command to reinstall a small number of active machines. Unfortunately, some master-replica pairs were impacted which resulted in the site going down.

- Planned maintenance to upgrade the OS on some machines
- Dropbox service been down for three hours

2014.6, Facebook

Facebook outage caused by software system update

20 June 2014 | By Hollie Luxford

Print Email + Share

in Share

Tweet 0

Like 0

Social networking site Facebook suffered a worldwide outage yesterday after an issue while updating the configuration of one of its software systems.

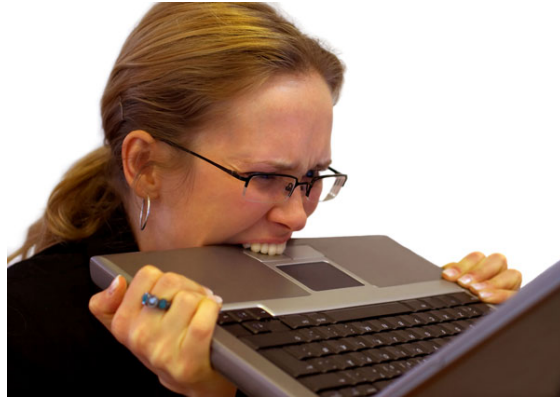
The worldwide outage lasted for 31 minutes.

Facebook

- Update the configuration of the software systems
- Failed Facebook for 31 minutes

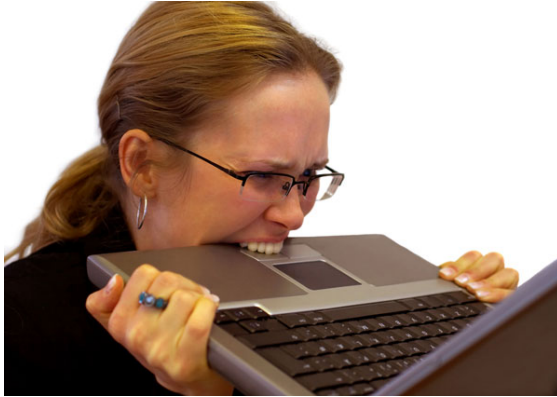
Impact of Erroneous Software Changes

- Poor user experience

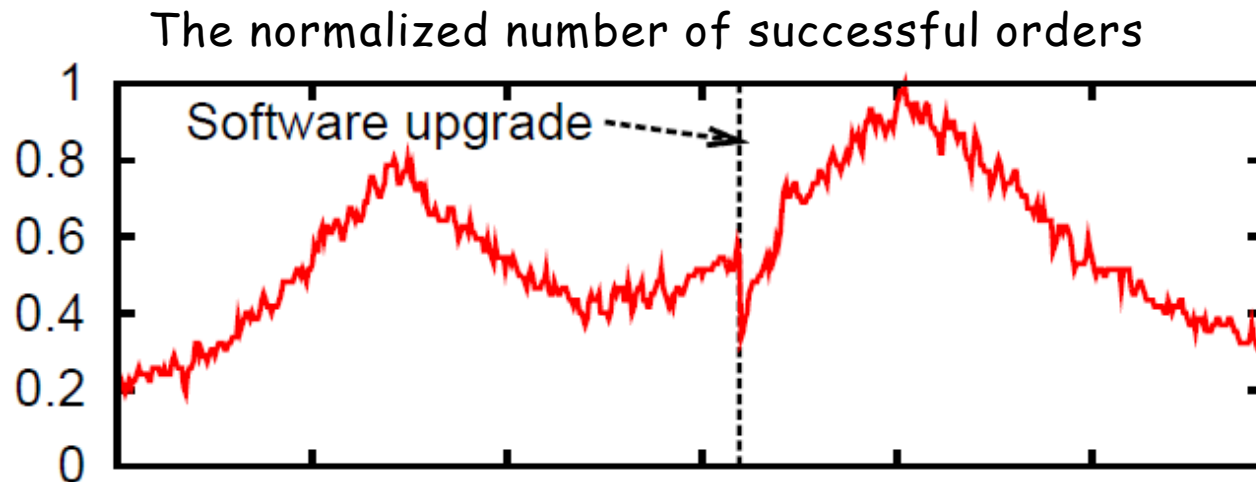


Impact of Erroneous Software Changes

- Poor user experience



- A drop in revenue



A real-world example

Manual Software Change Impact Assessment

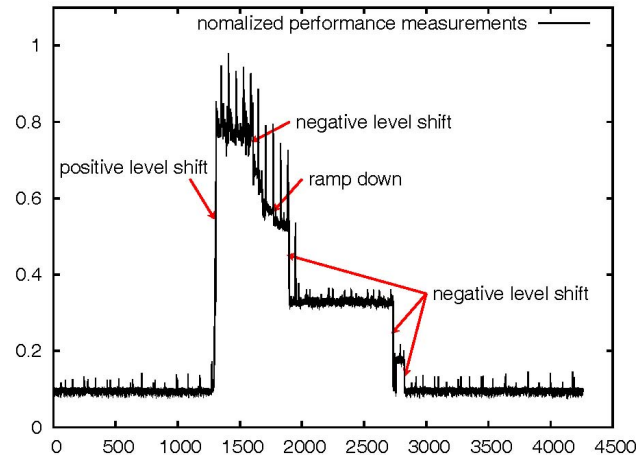


Select a **subset** of KPIs
that maybe impacted

Manual Software Change Impact Assessment

Select a **subset** of KPIs
that maybe impacted

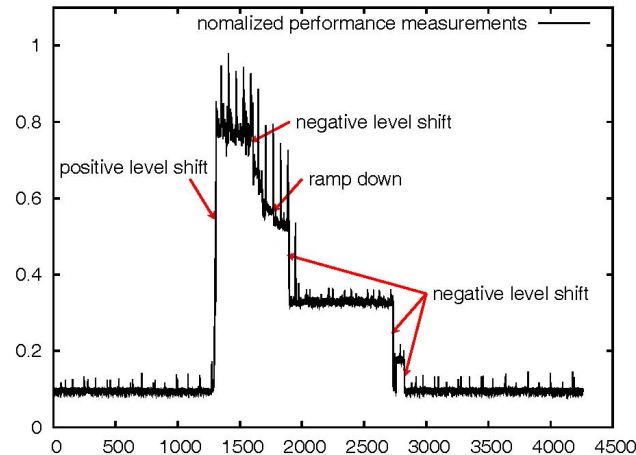
Inspect KPI changes



Manual Software Change Impact Assessment

Select a **subset** of KPIs that maybe impacted

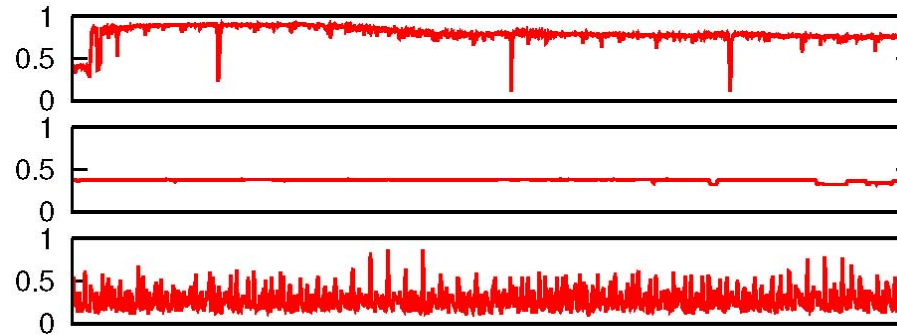
Inspect KPI changes



Decide whether to roll back

KPI (Key Performance Indicator) in Software Change

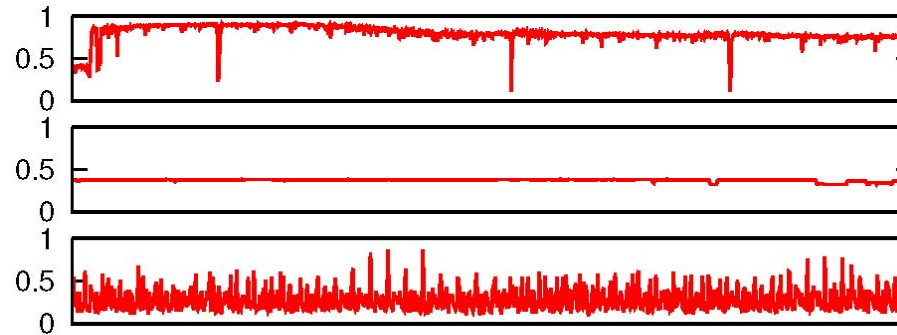
- KPIs of servers
 - CPU utilization
 - Memory utilization
 - NIC throughput
 - ...



KPI (Key Performance Indicator) in Software Change

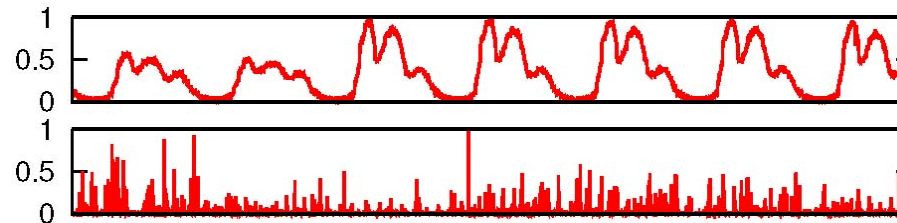
- KPIs of servers

- CPU utilization
- Memory utilization
- NIC throughput
- ...



- KPIs of modules/processes

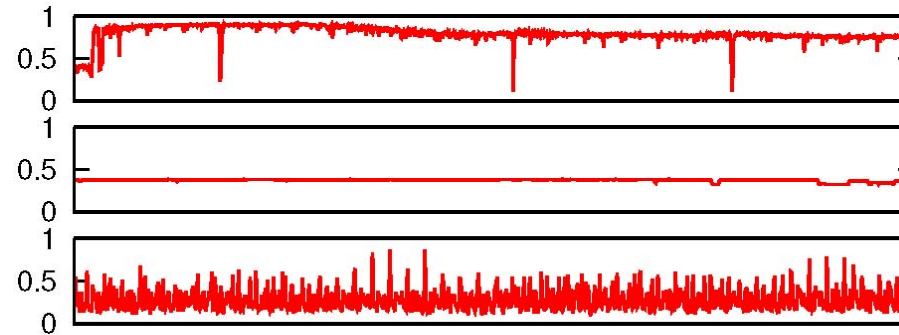
- Web page view count
- Web page view delay
- ...



KPI (Key Performance Indicator) in Software Change

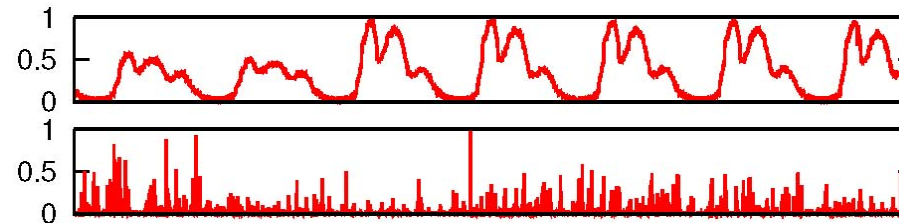
- KPIs of servers

- CPU utilization
- Memory utilization
- NIC throughput
- ...



- KPIs of modules/processes

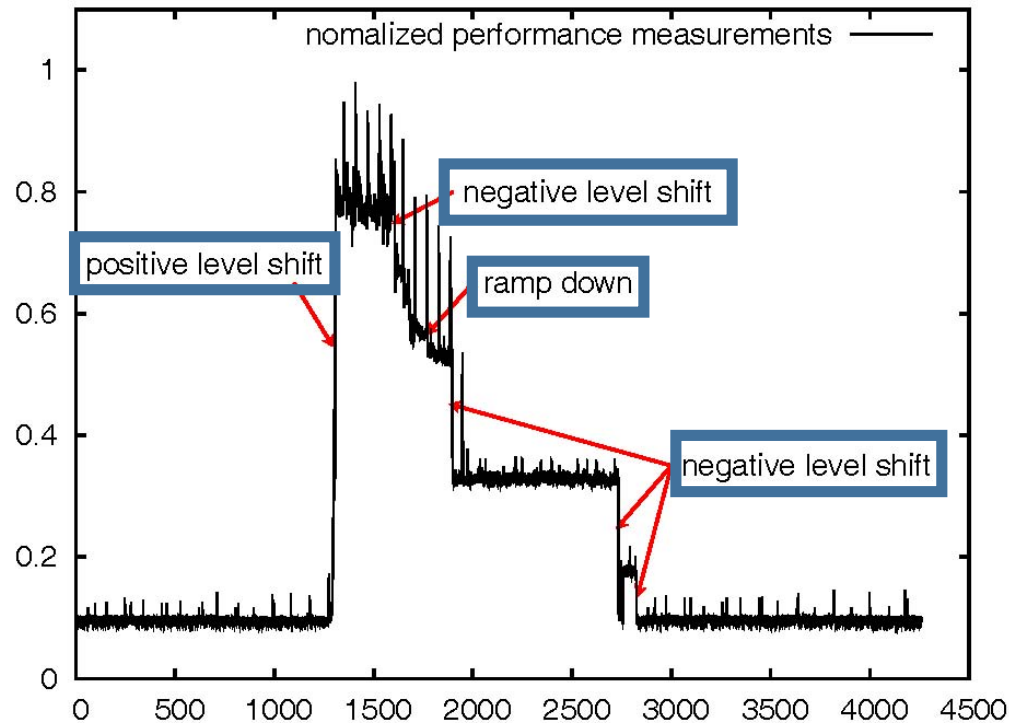
- Web page view count
- Web page view delay
- ...



- Up to hundreds of KPIs for a single software change

Definition of KPI Change: Level Shift or Ramp up/down

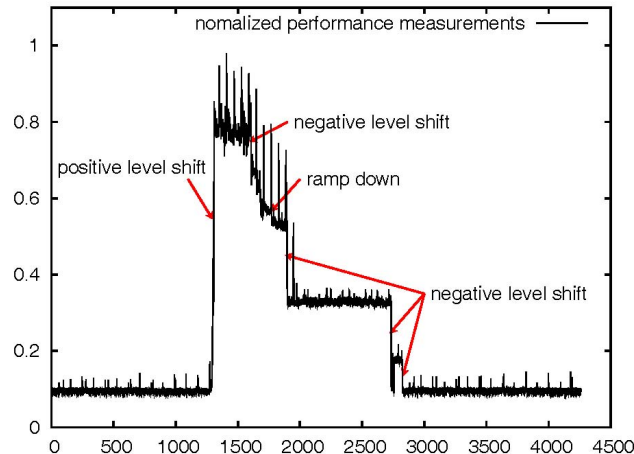
- KPI change
 - Indicative of performance increase/degradation
 - Hard to simulate in testbeds
 - Not reproducible



Manual Software Change Impact Assessment

Select a **subset** of KPIs that maybe impacted

Inspect KPI changes



Decide whether to roll back

- Labor-intensive
- Prone to error
- Not scalable

Design Goal

Software Change Impact Assessment System

Decide
whether to roll
back

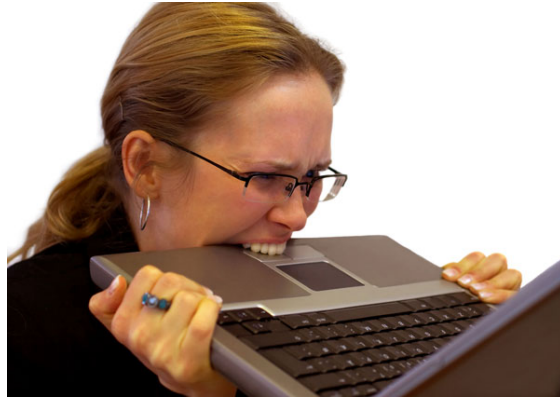
- Automatic
- Scalable
- Robust to various software changes and KPIs

Outline

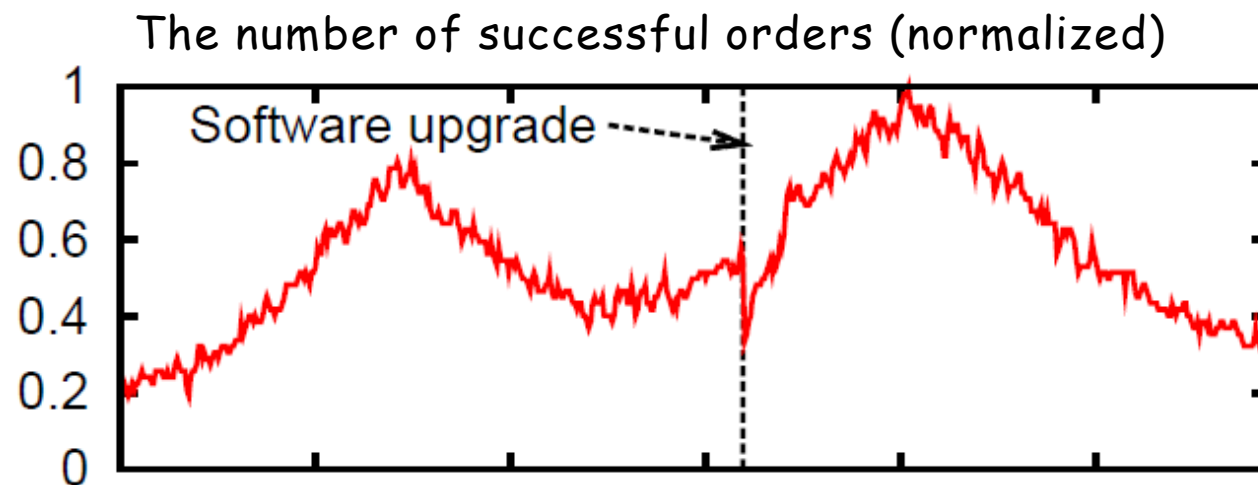
- Background and Motivation
- Challenges
- Key Ideas
- Results
- Conclusion

Challenge 1: Short Detection Delay Requirement Against Robustness

- Poor user experience



- A drop in revenue



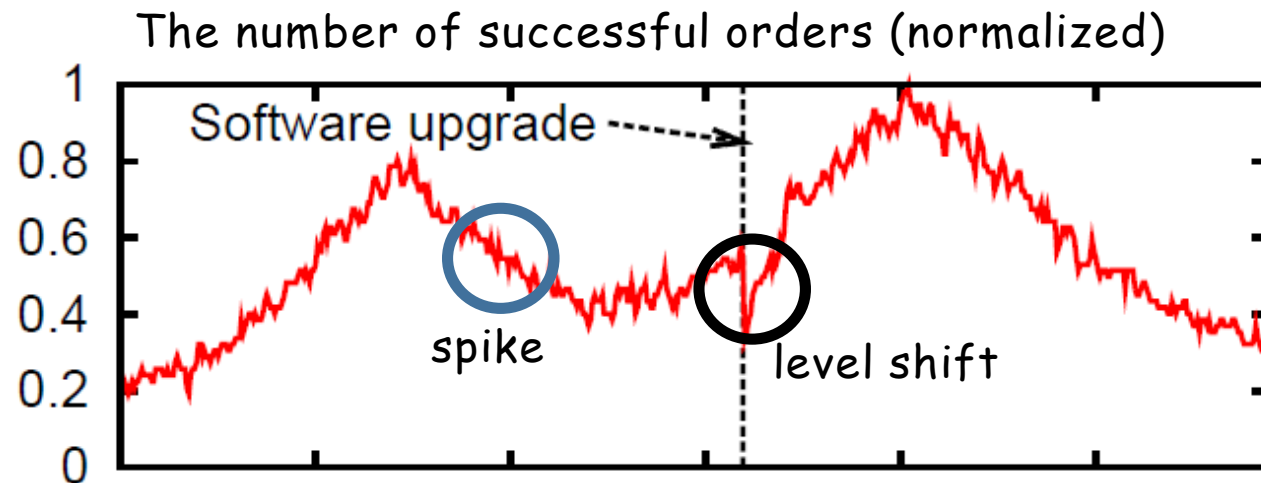
A real-world example

Challenge 1: Short Detection Delay Requirement Against Robustness

- Poor user experience



- A drop in revenue



A real-world example

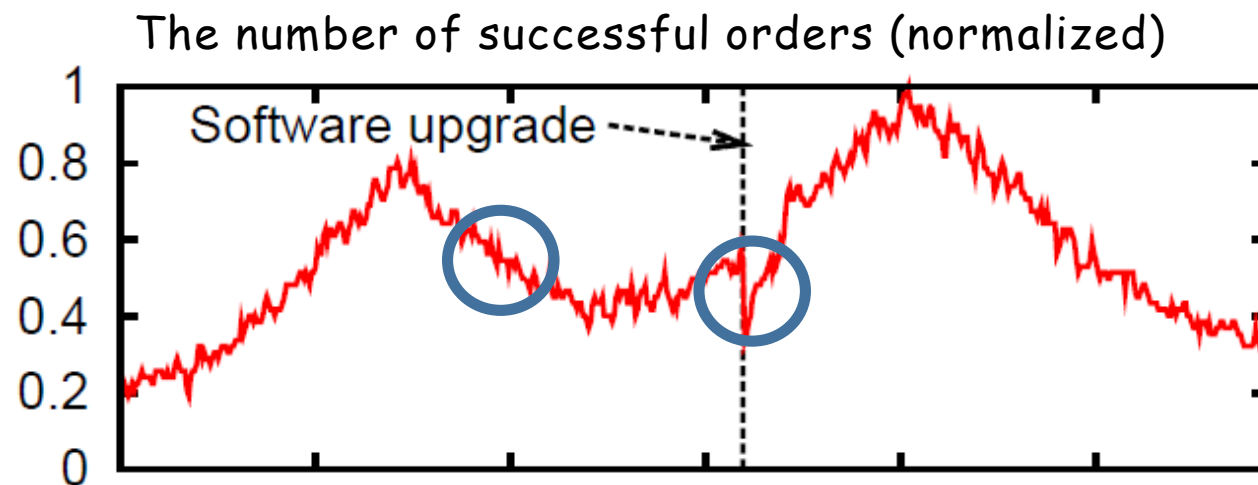
Challenge 1: Short Detection Delay Requirement Against Robustness

- Poor user experience



Detect KPI changes rapidly and accurately

- A drop in revenue



A real-world example

Challenge 2: Large Number of KPIs



Challenge 2: Large Number of KPIs



100+ Internet-based services

20+ Internet-based services has 100+ million users

10k+ modules

500+ thousand servers



Challenge 2: Large Number of KPIs

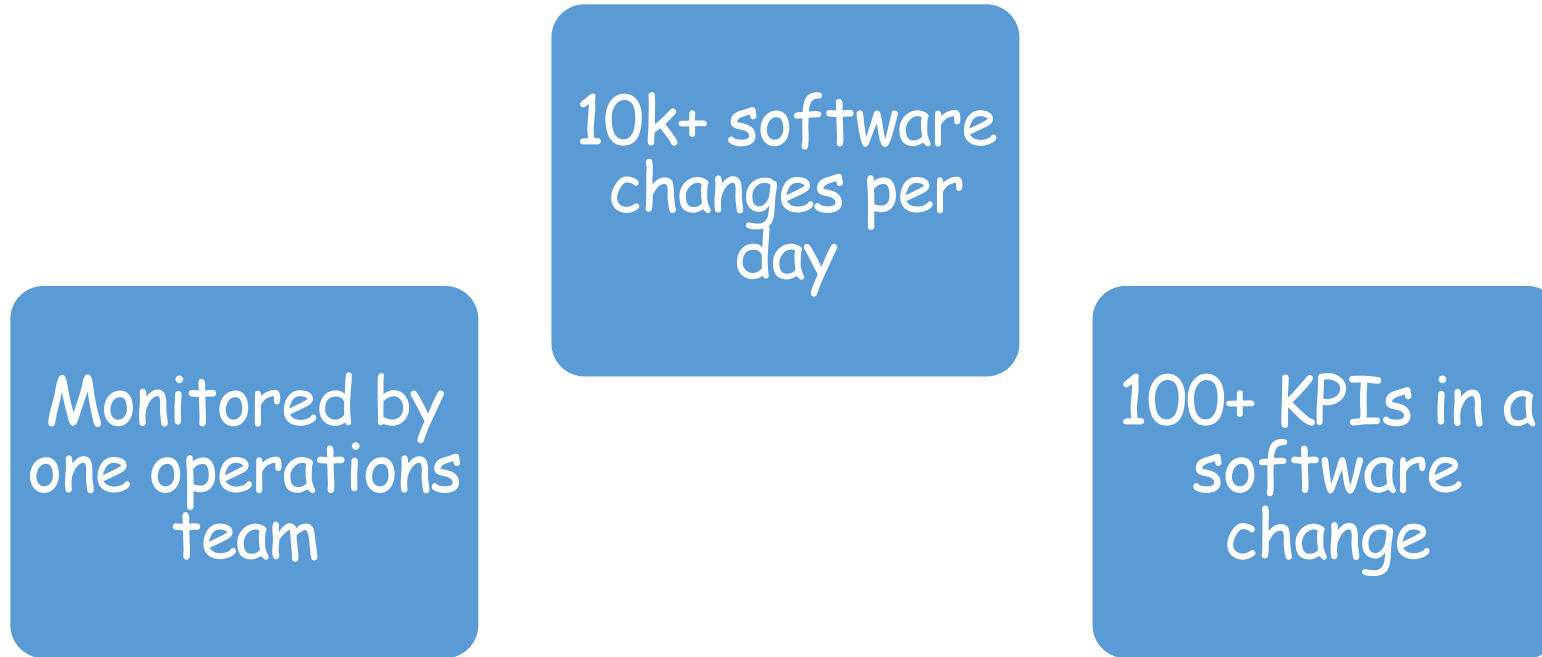
Monitored by
one operations
team

Challenge 2: Large Number of KPIs

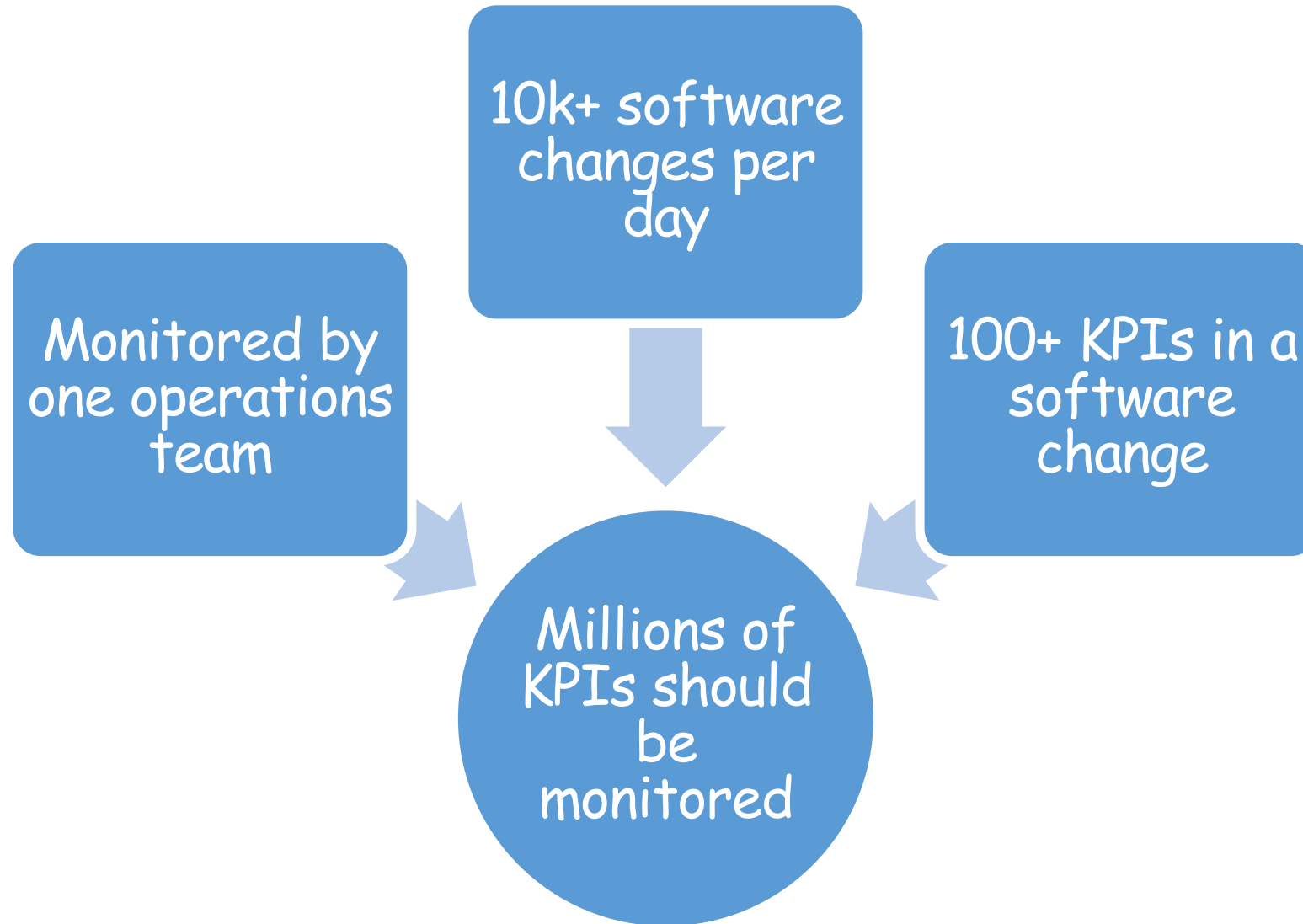
10k+ software
changes per
day

Monitored by
one operations
team

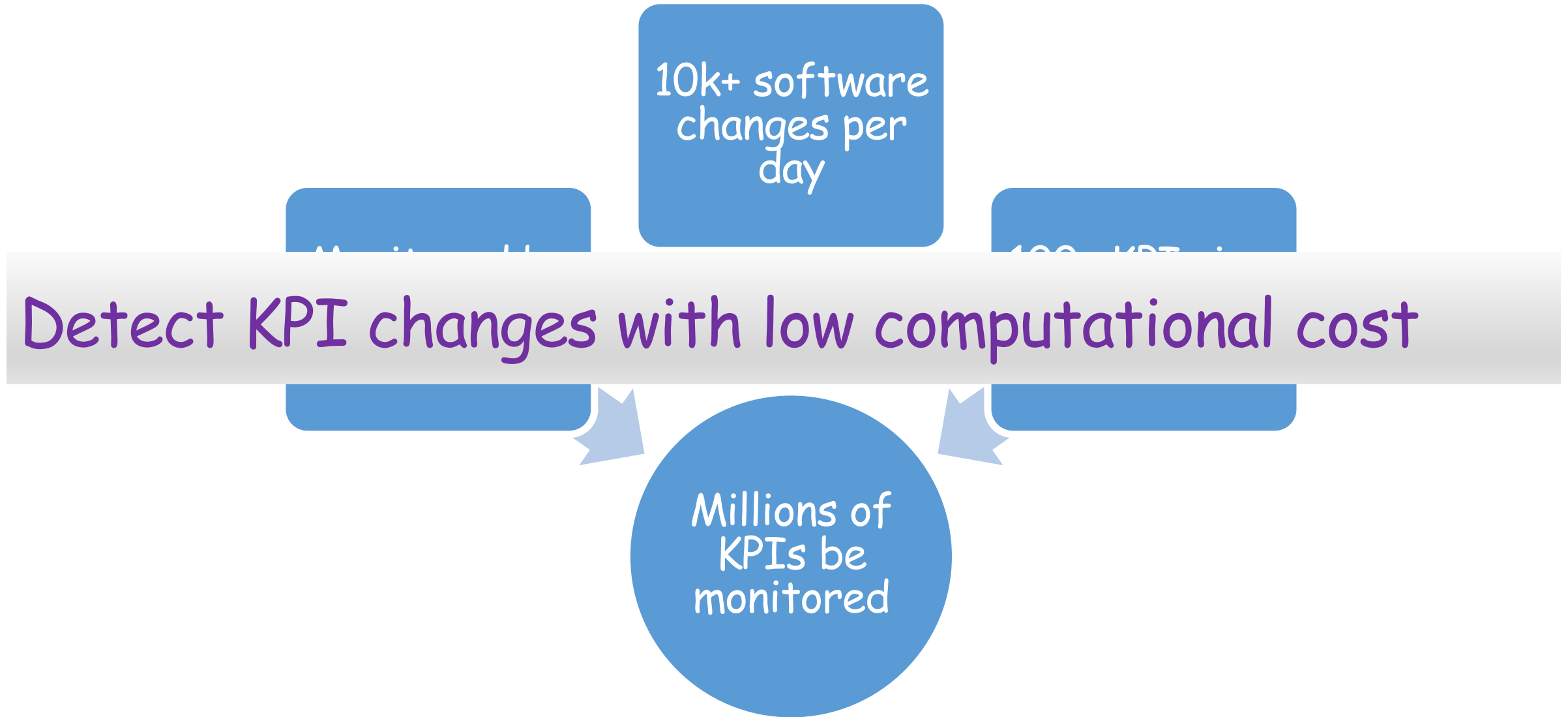
Challenge 2: Large Number of KPIs



Challenge 2: Large Number of KPIs



Challenge 2: Large Number of KPIs



Challenge 3: Diverse Types of Data

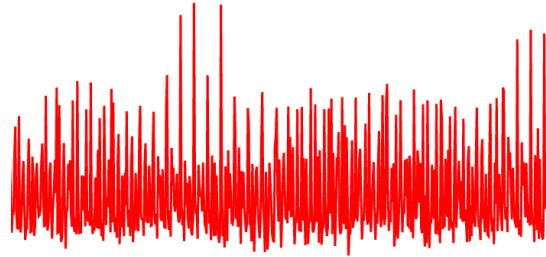
- Diverse types of KPI data

Seasonal



Page view count

Variable



NIC throughput

Stationary



Memory utilization

Challenge 3: Diverse Types of Data

- Diverse types of KPI data

Seasonal

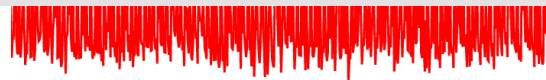
Variable

Stationary

Robust to various KPIs



Page view count



NIC throughput

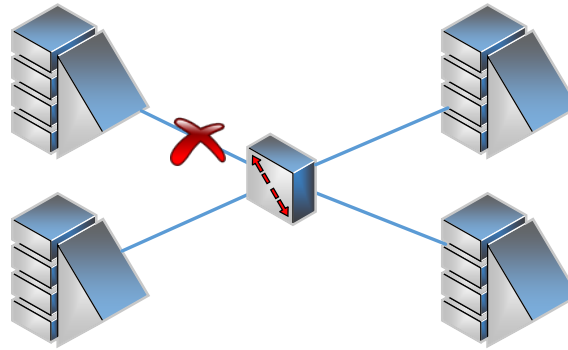
Memory utilization

Challenge 4: KPI Changes Maybe Caused by Other Factors

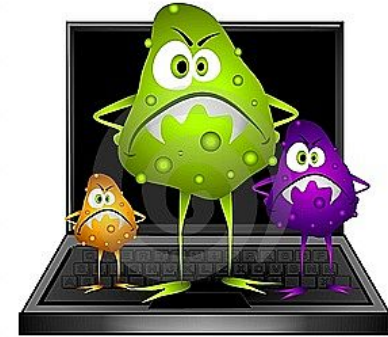
Seasonality



Network
breakdowns



Malicious
attacks



Challenge 4: KPI Changes Maybe Caused by Other Factors

Seasonality

Network
breakdowns

Malicious
attacks

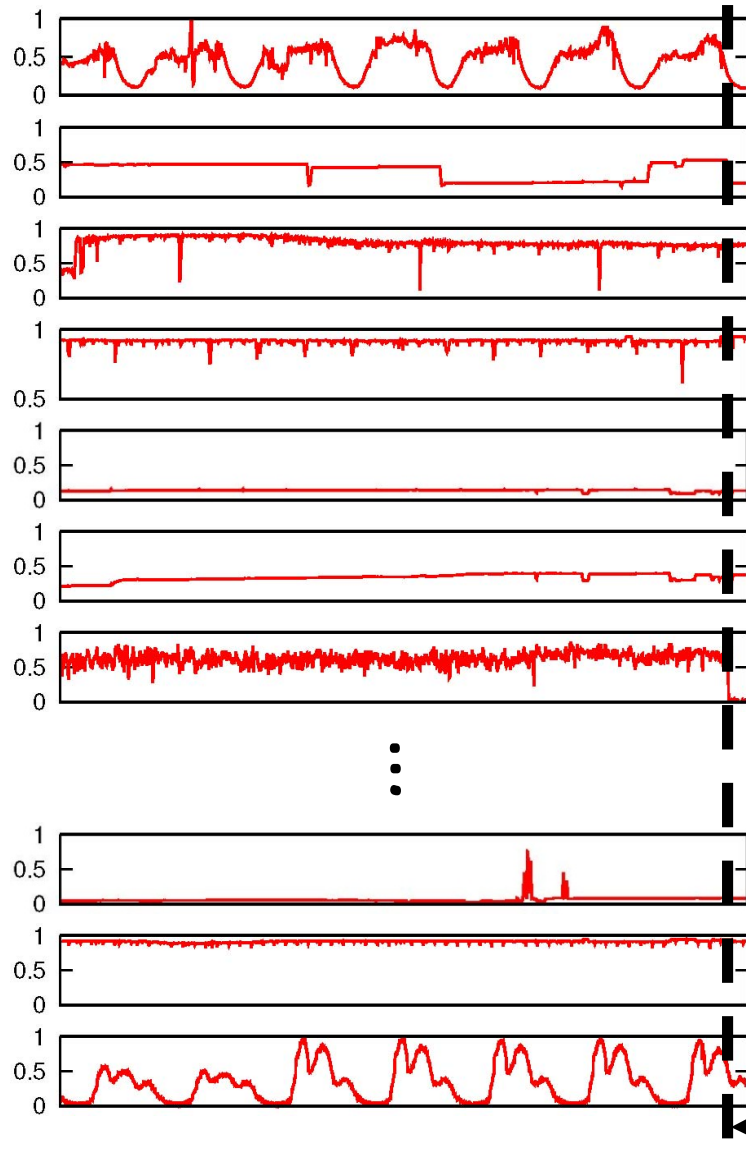
Eliminate KPI changes induced by other factors



Outline

- Background and Motivation
- Challenges
- Key Ideas
- Results
- Conclusion

Design Overview

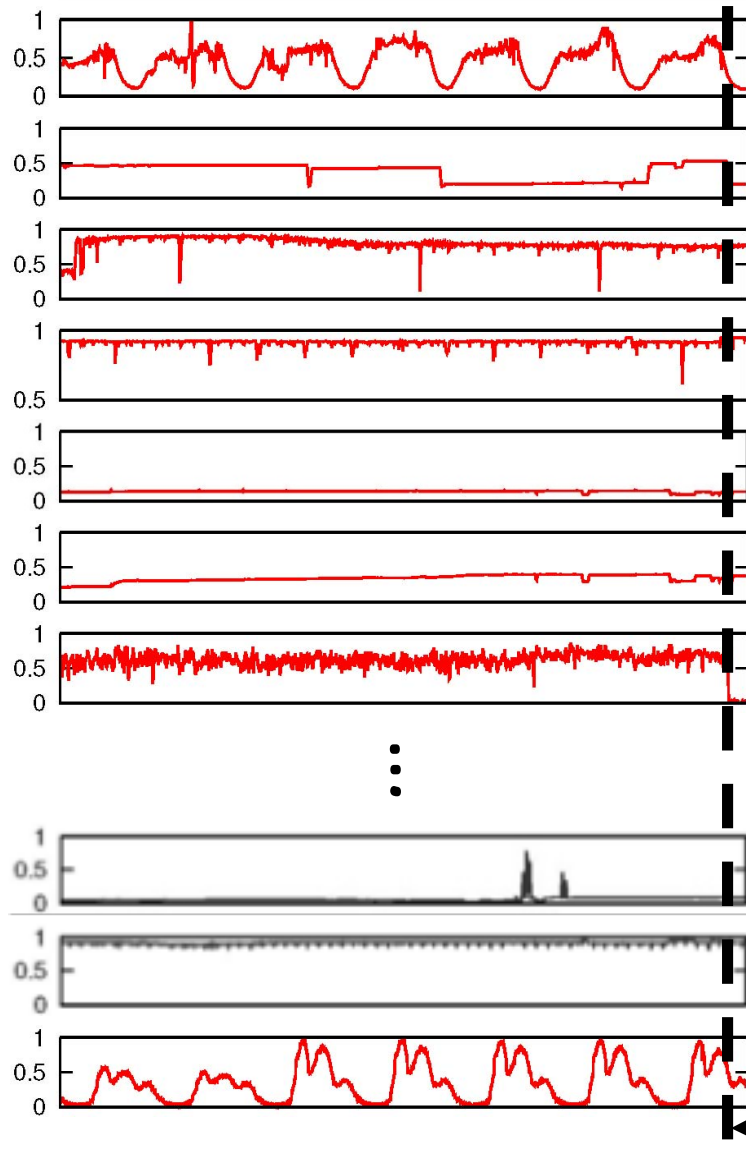


Step 1 - Identify the impact set

Step 1

Software change in module A

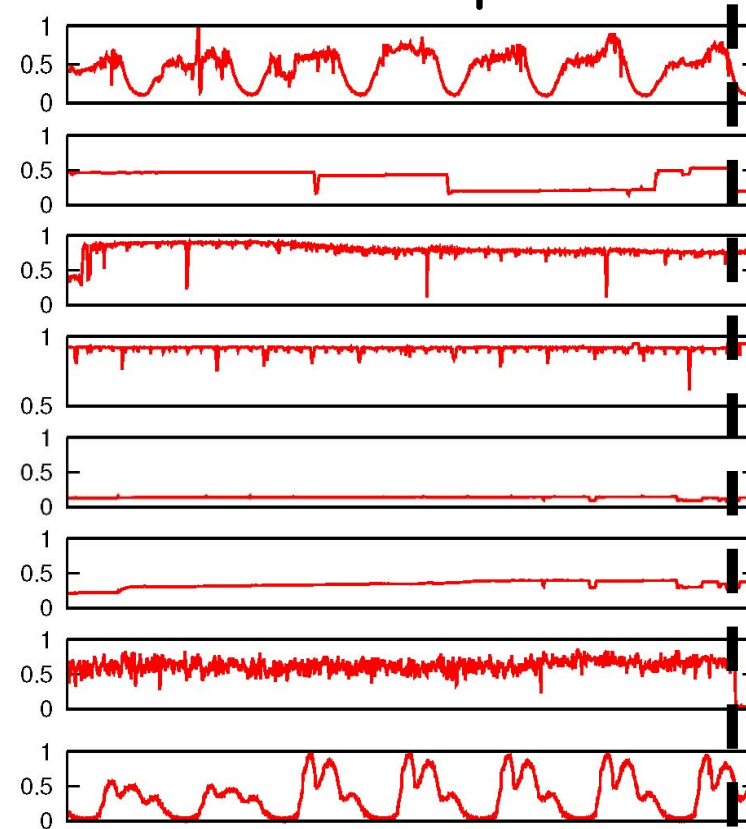
Design Overview



Step 1 - Identify the impact set

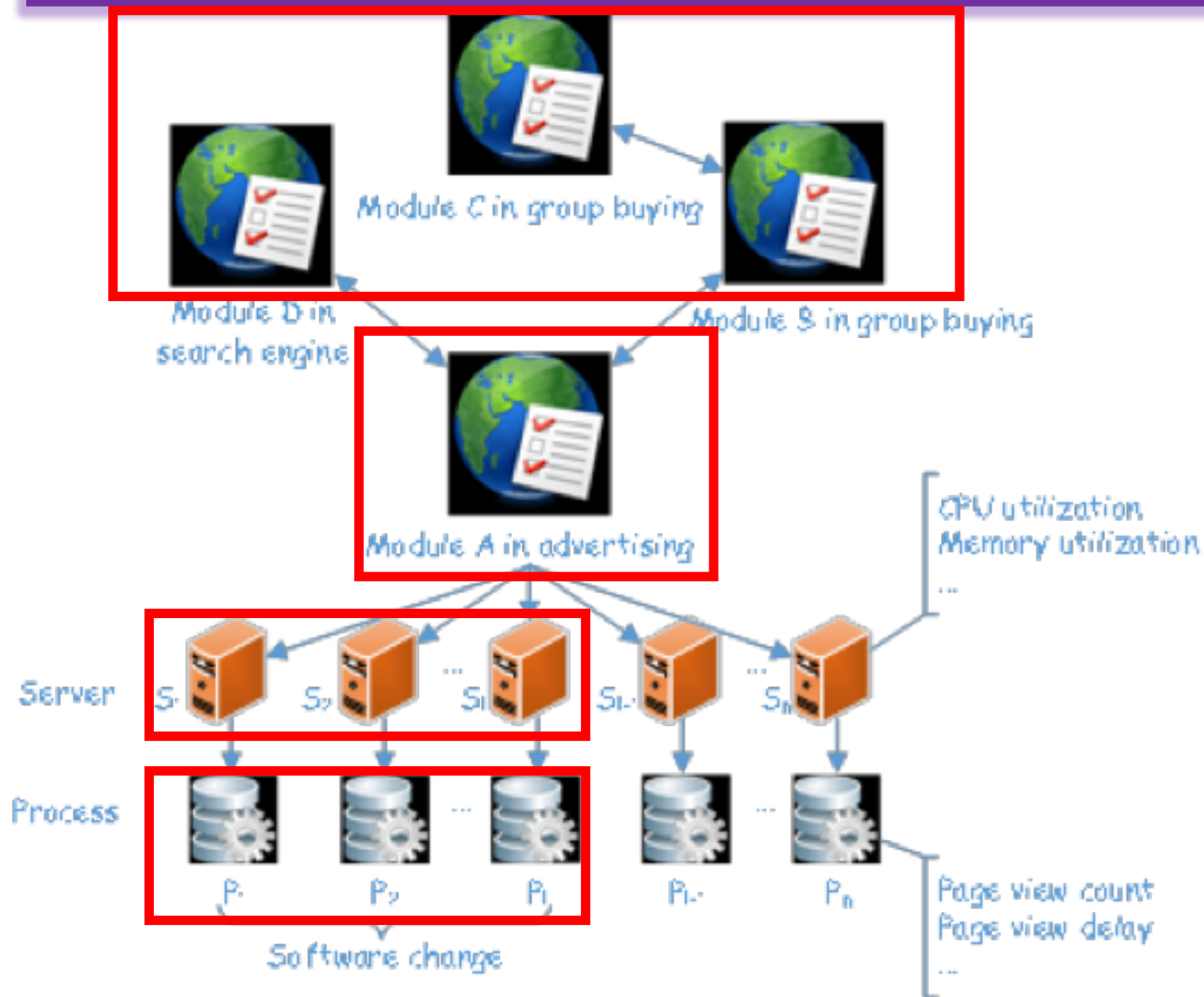
Step 1

KPIs in the impact set



Software change in module A

Identify the Impact Set: Automatically Retrieve the Relevant KPIs



Identify the Impact Set: Automatically Retrieve the Relevant KPIs

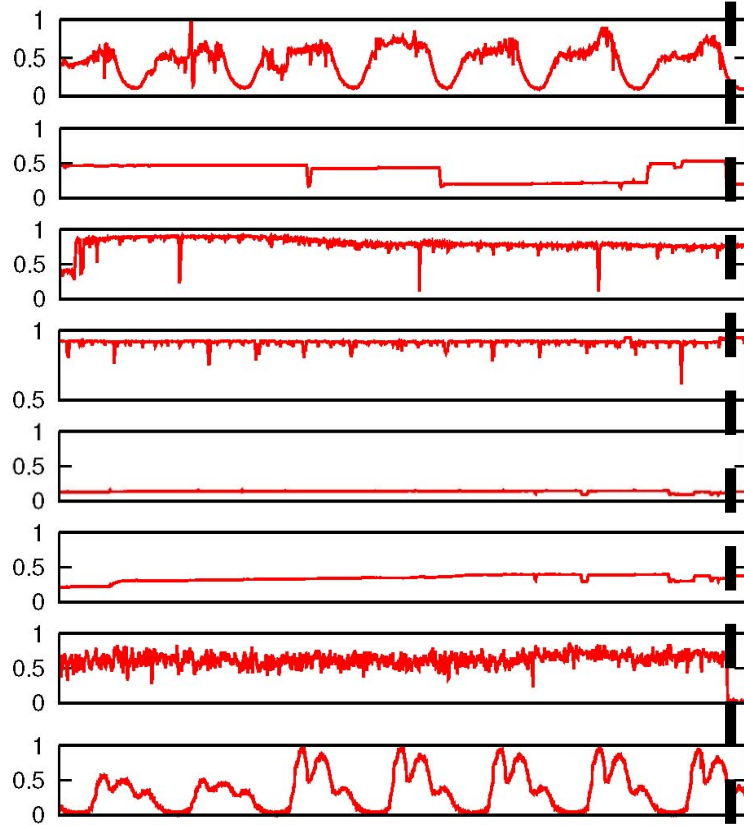


Input from operators

- Modules related module A: module B, C, D
- Servers/processes where the software change is deployed.

Design Overview

KPIs in the impact set



Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

Step 1

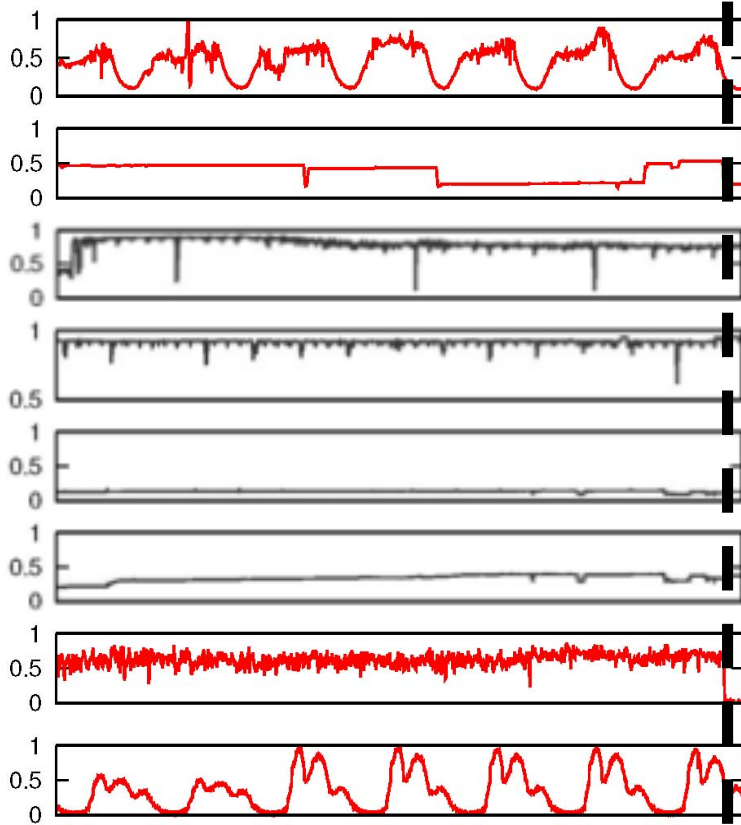
Step 2

Software change in module A

Design Overview

Step 1

KPIs in the impact set

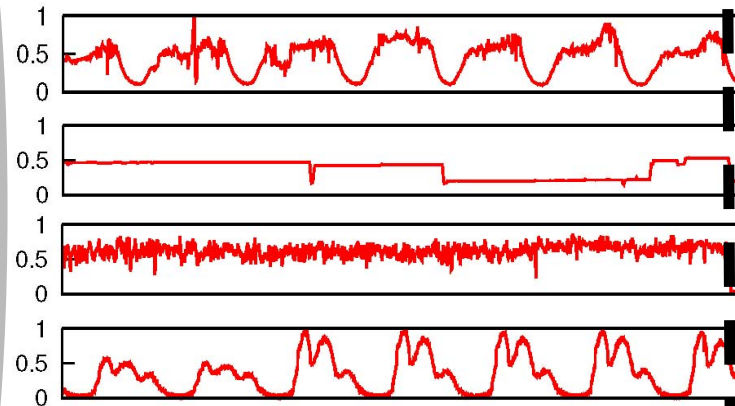


Step 2

Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

KPIs with behavior changes

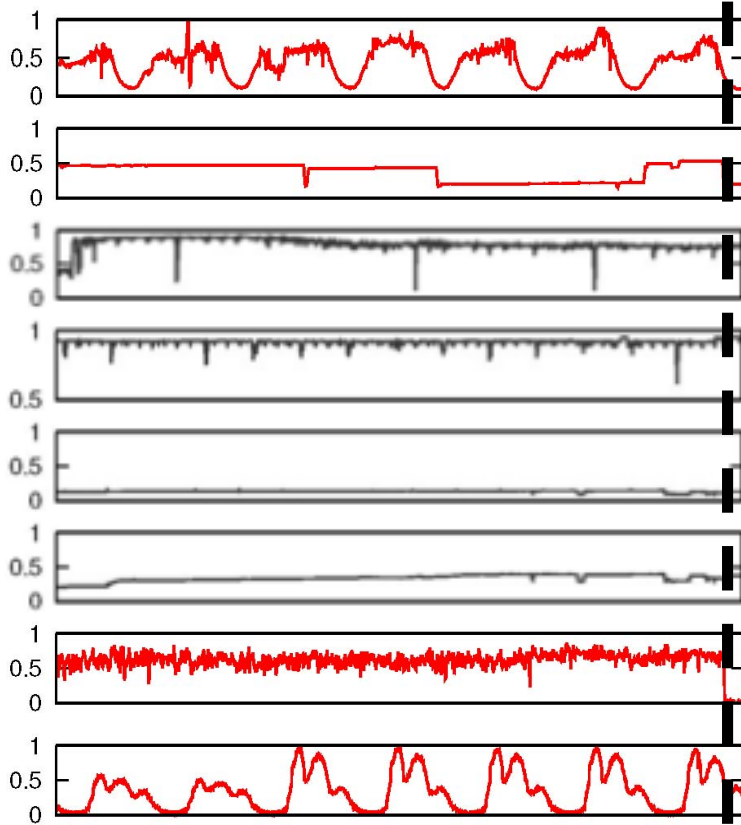


Software change in module A

Design Overview

Step 1

KPIs in the impact set

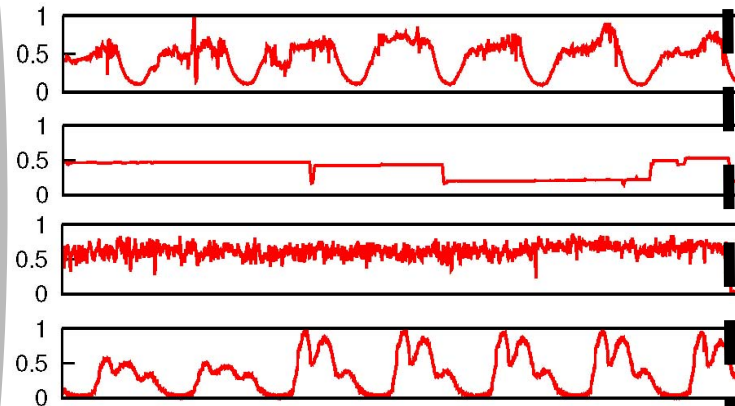


Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

Step 2

KPIs with behavior changes



Short detection delay
requirement against robustness

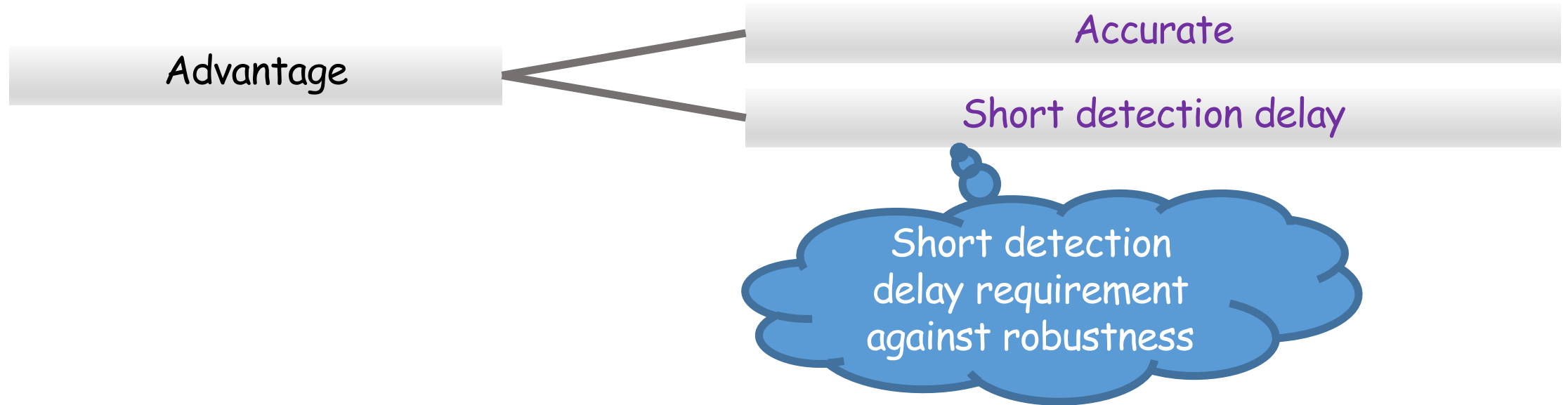
Diverse types of data

Large number of KPIs

Software change in module A

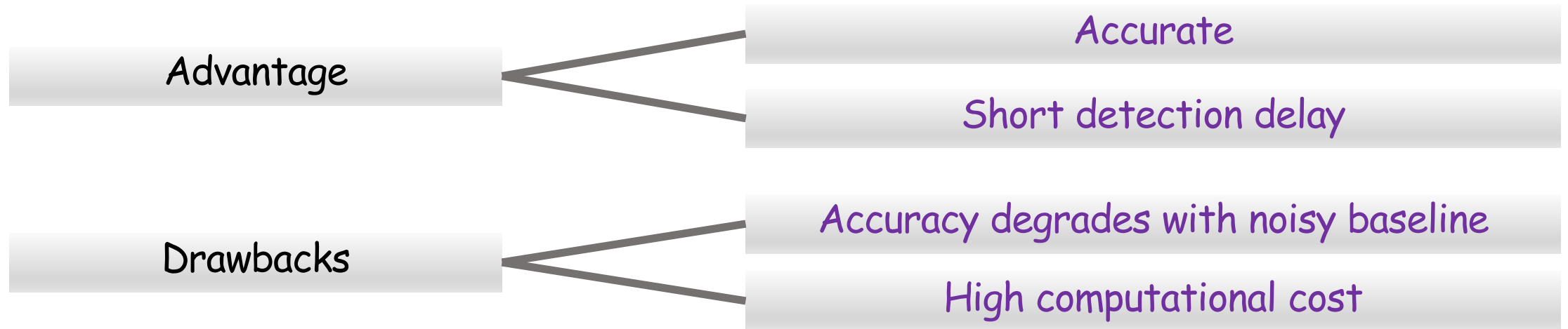
Improved Singular Spectrum Transform (SST)

- Improved singular spectrum transform (SST) $x_s(t) = 1 - \alpha(t)^T \beta(t)$



Improved Singular Spectrum Transform (SST)

- Improved singular spectrum transform (SST) $x_s(t) = 1 - \alpha(t)^T \beta(t)$



T. Idé and K. Tsuda, SDM 2007

Improved Singular Spectrum Transform (SST)

- Improved singular spectrum transform (SST)

$$\hat{x}(t) = \frac{\sum_{i=1}^{\eta} \lambda_i \times \varphi_i(t)}{\sum_{i=1}^{\eta} \lambda_i}$$

Advantage

Accurate

Short detection delay

Drawbacks

Accuracy degrades with noisy baseline

High computational cost

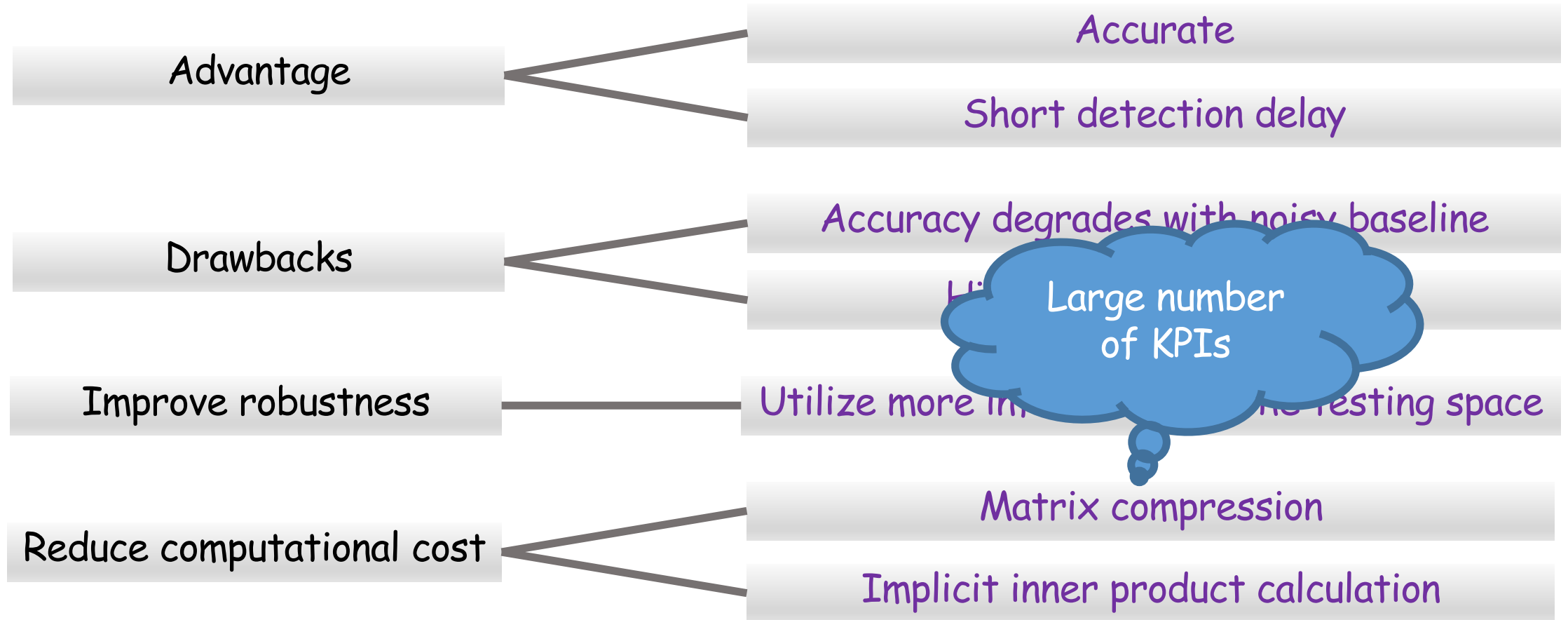
Improve robustness

Utilize more information in the testing space

Diverse types
of data

Improved Singular Spectrum Transform (SST)

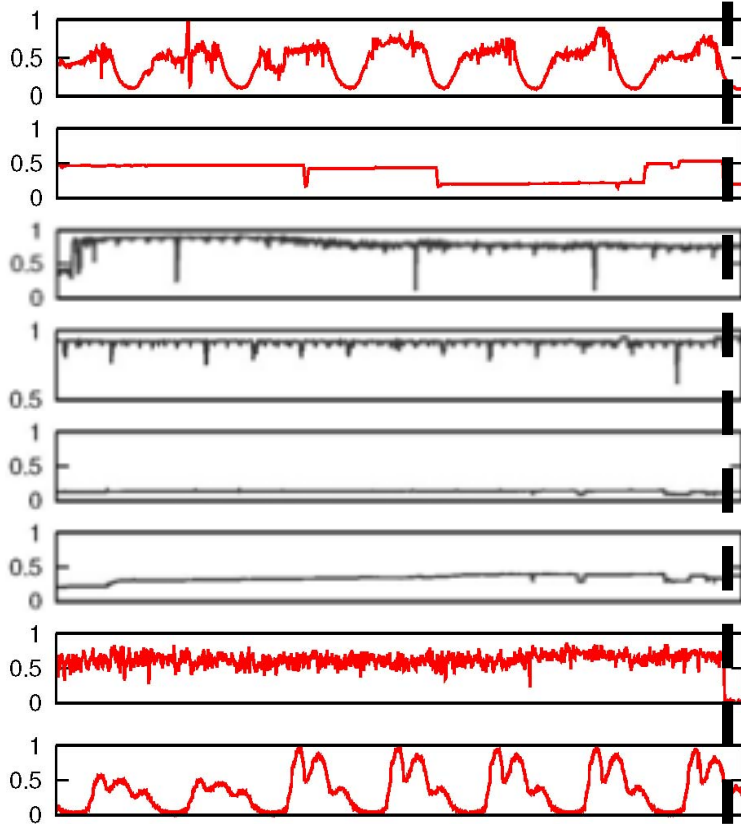
- Improved singular spectrum transform (SST) $\varphi_i(t) \simeq 1 - \sum_{j=1}^{\eta} x_j^2$



Design Overview

Step 1

KPIs in the impact set



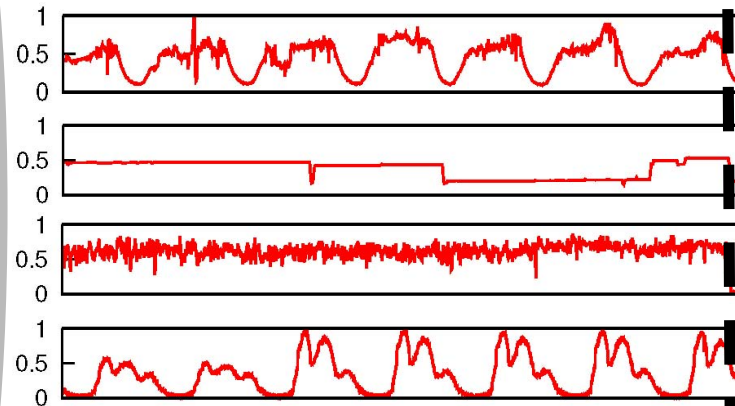
Step 2

Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

Step 3 - Eliminate KPI changes induced by other factors

KPIs with behavior changes

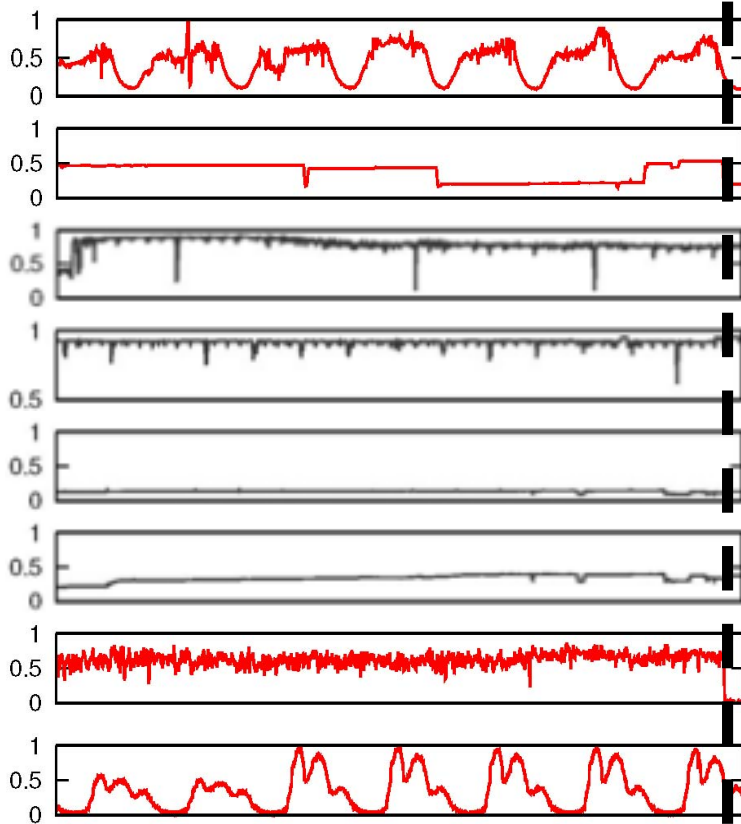


Software change in module A

Design Overview

Step 1

KPIs in the impact set



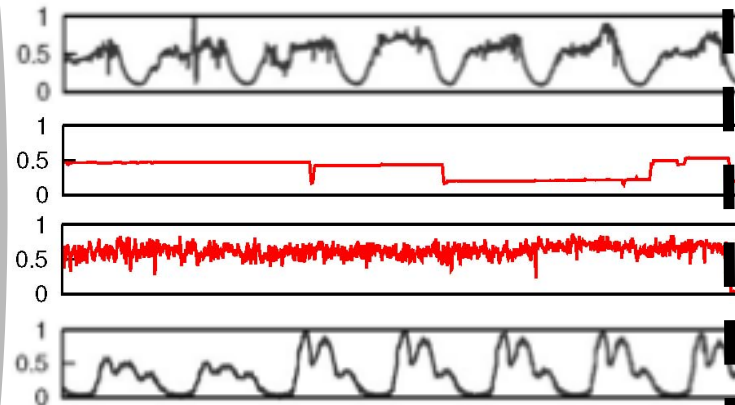
Step 2

Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

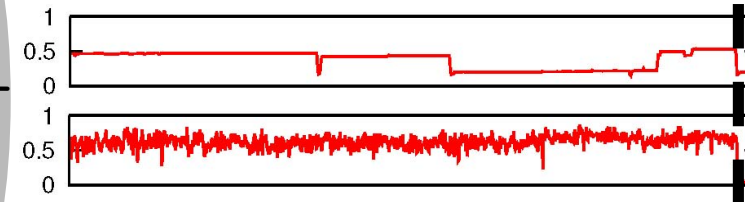
Step 3 - Eliminate KPI changes induced by other factors

KPIs with behavior changes



Step 3

KPIs with behavior changes induced by software change

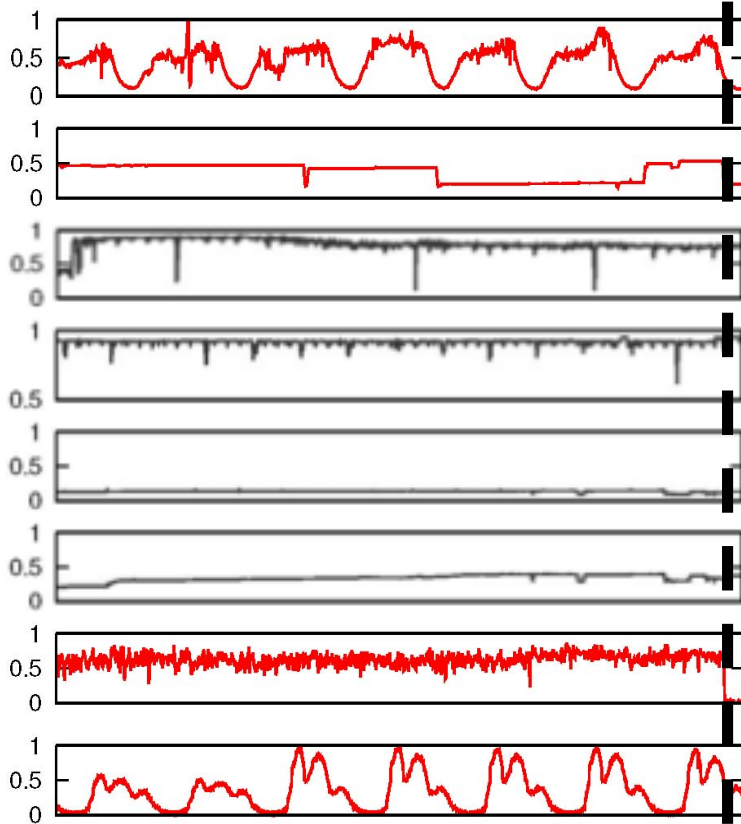


Software change in module A

Design Overview

Step 1

KPIs in the impact set



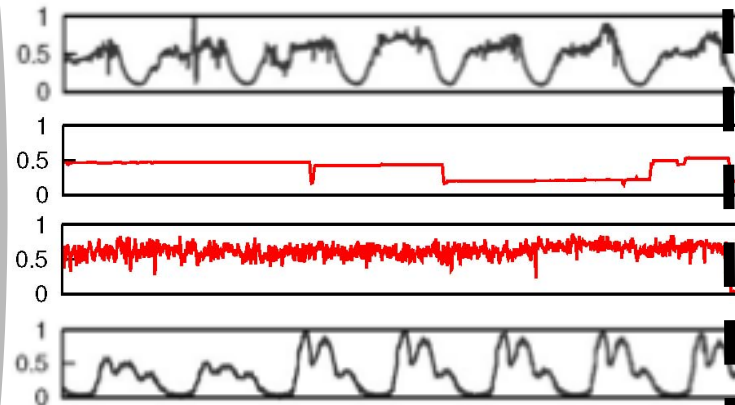
Step 2

Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

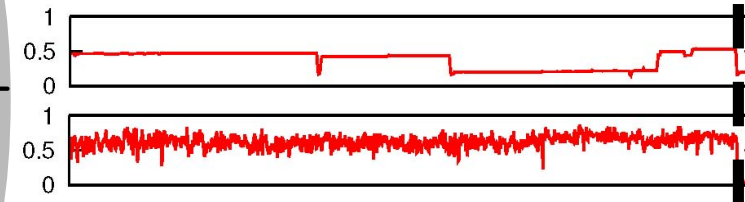
Step 3 - Eliminate KPI changes induced by other factors

KPIs with behavior changes



Step 3

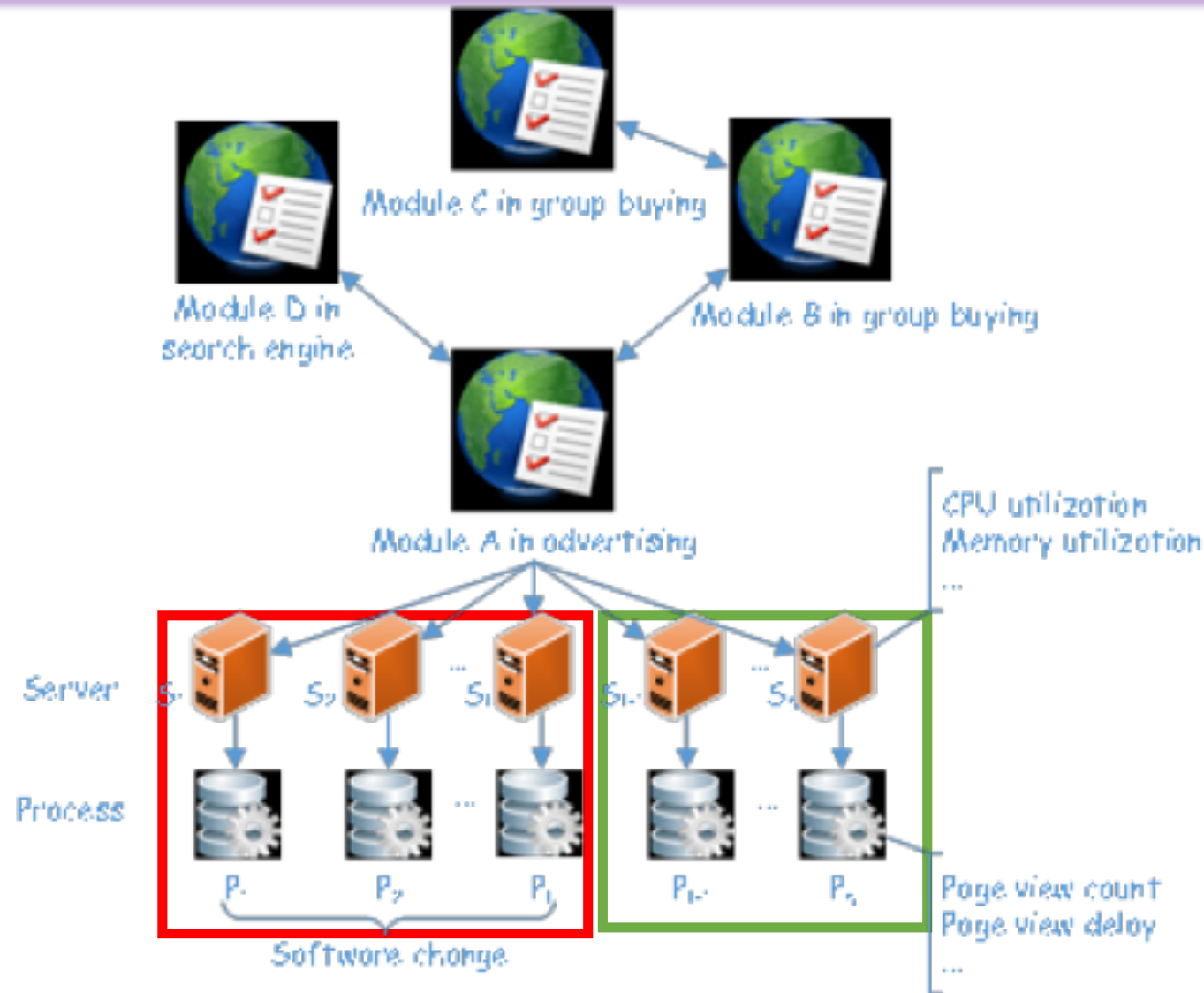
KPIs with behavior changes induced by software change



KPI changes maybe caused by other factors

Software change in module A

Eliminate KPI Changes Induced by Other Factors

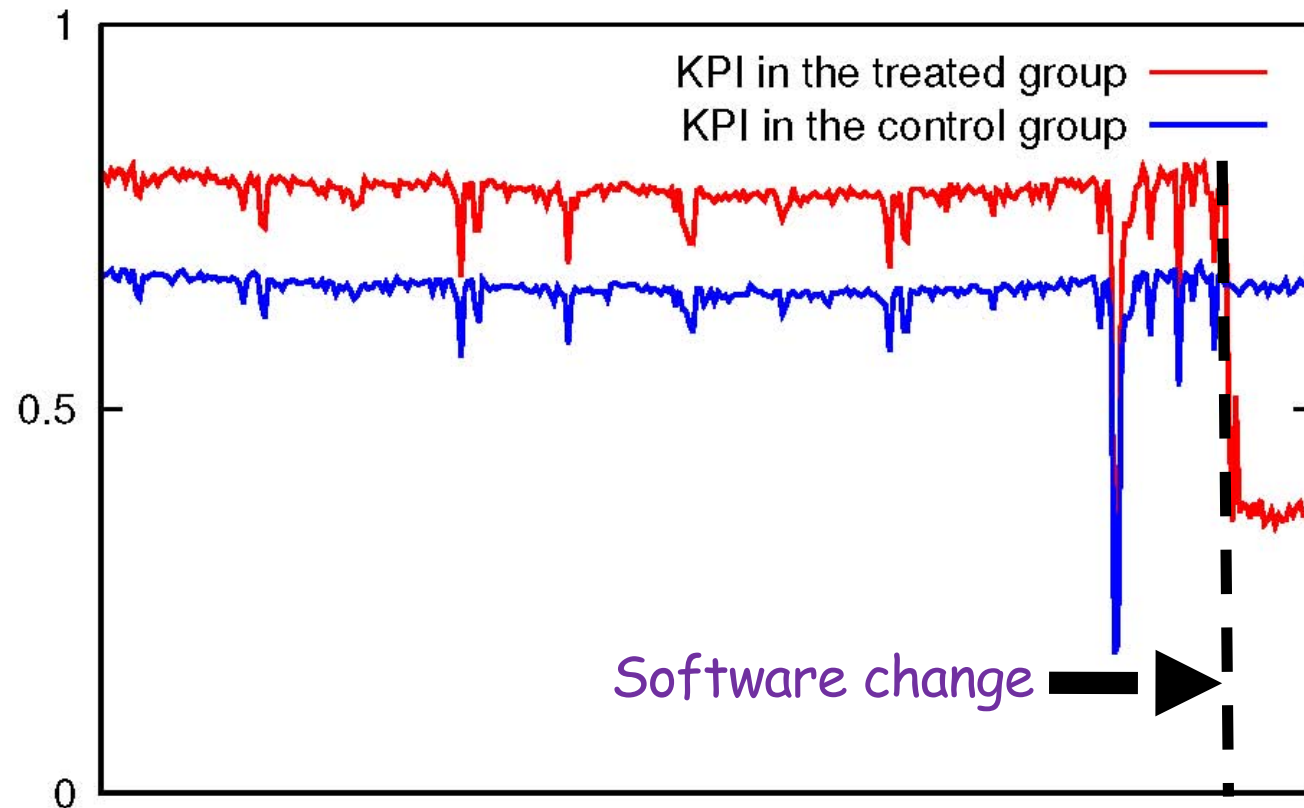


Eliminate KPI Changes Induced by Other Factors

- Split testing
 - Evaluation of interventions instituted at a specific time
 - Control group & treated group

Eliminate KPI Changes Induced by Other Factors

- Split testing
 - Evaluation of interventions instituted at a specific time
 - Control group & treated group



Eliminate KPI Changes Induced by Other Factors

Treated group

- Servers/processes in the impact set



Eliminate KPI Changes Induced by Other Factors

Treated group

- Servers/processes in the impact set

Control group

- Servers/processes in the same module
- Without software change



Eliminate KPI Changes Induced by Other Factors

Treated group

- Servers/processes in the impact set

Control group

- Servers/processes in the same module
- Without software change

DiD method



Eliminate KPI Changes Induced by Other Factors

Treated group

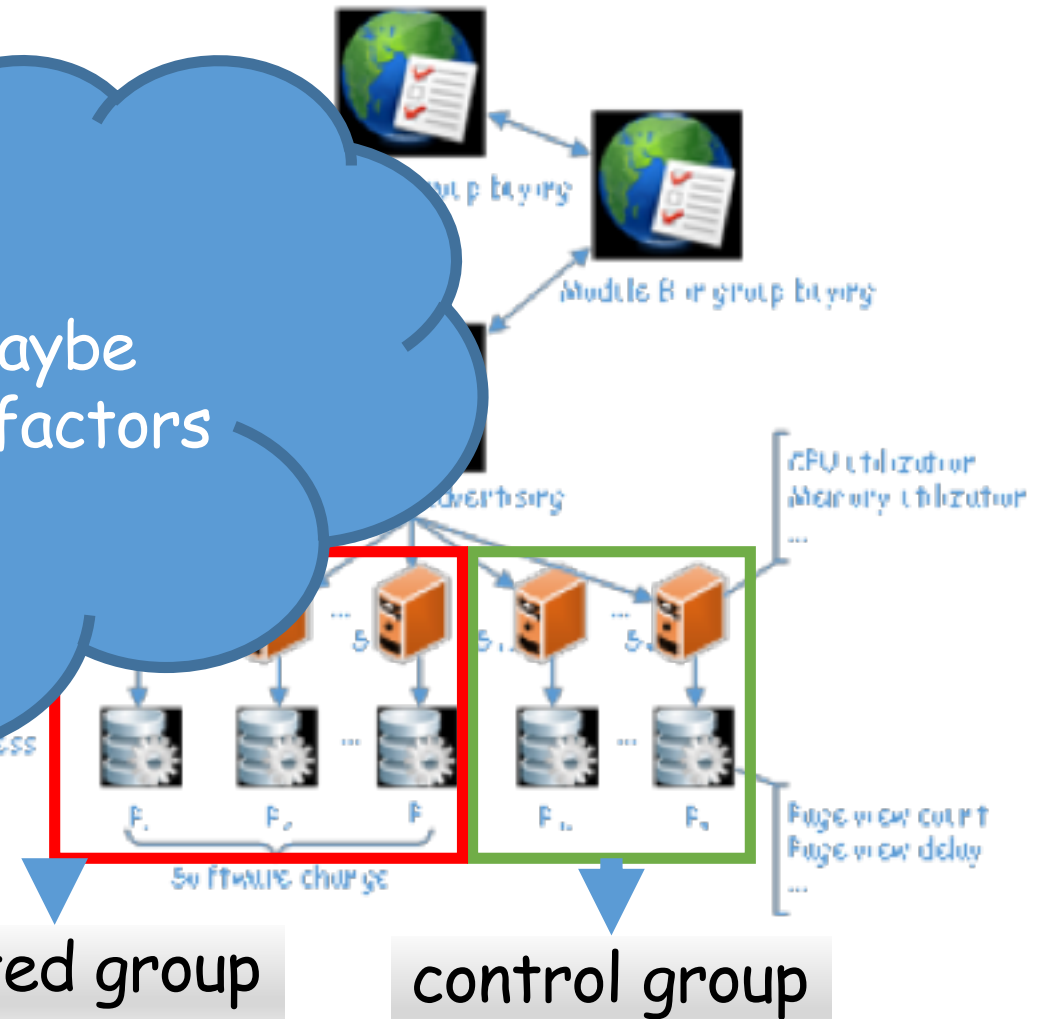
- Servers/processes in t

Control group

- Servers/processes
- Without software chan

DiD method

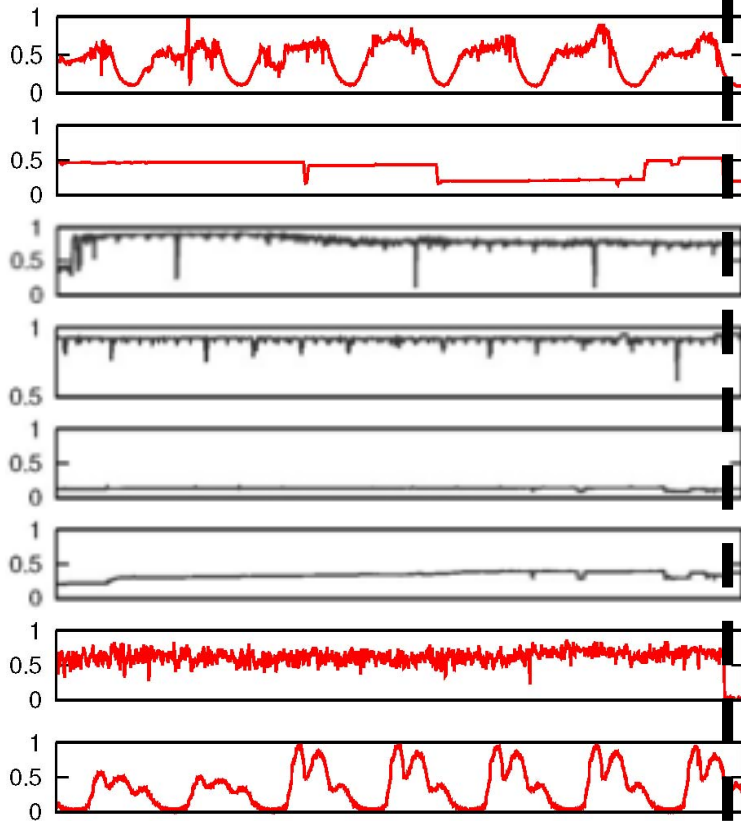
KPI changes maybe caused by other factors



Design Overview

Step 1

KPIs in the impact set



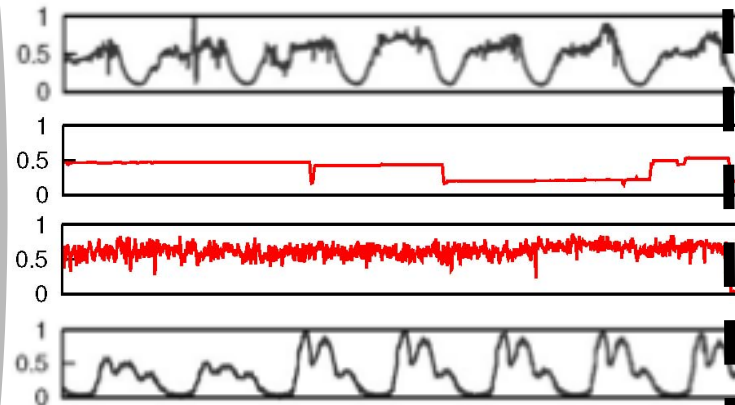
Step 1 - Identify the impact set

Step 2 - Detect behavior changes in KPIs

Step 3 - Eliminate KPI changes induced by other factors

Step 2

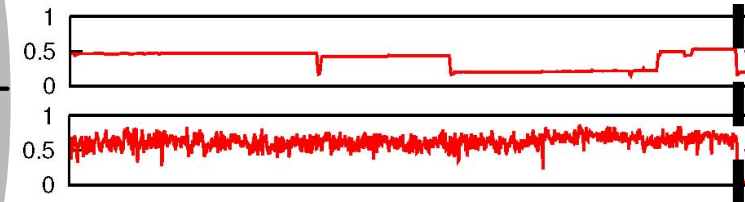
KPIs with behavior changes



improved SST

Step 3

KPIs with behavior changes induced by software change



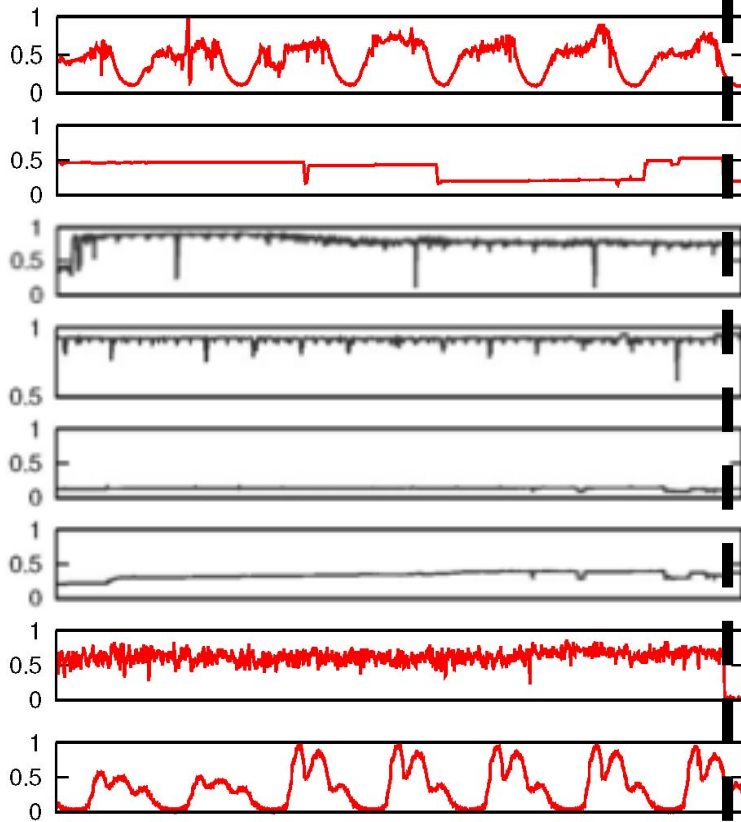
split testing

Software change in module A

Design Overview

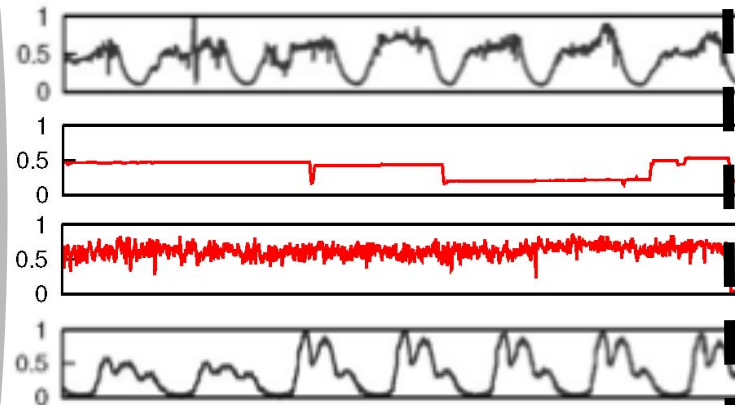
Step 1

KPIs in the impact set



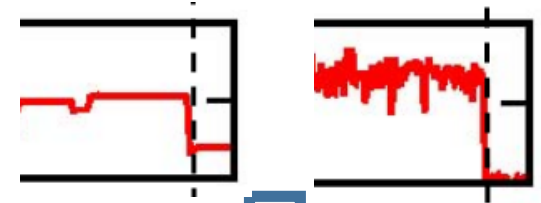
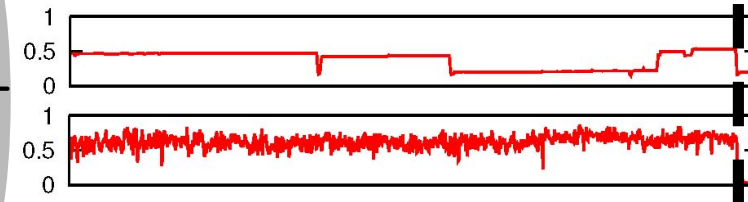
Step 2

KPIs with behavior changes



Step 3

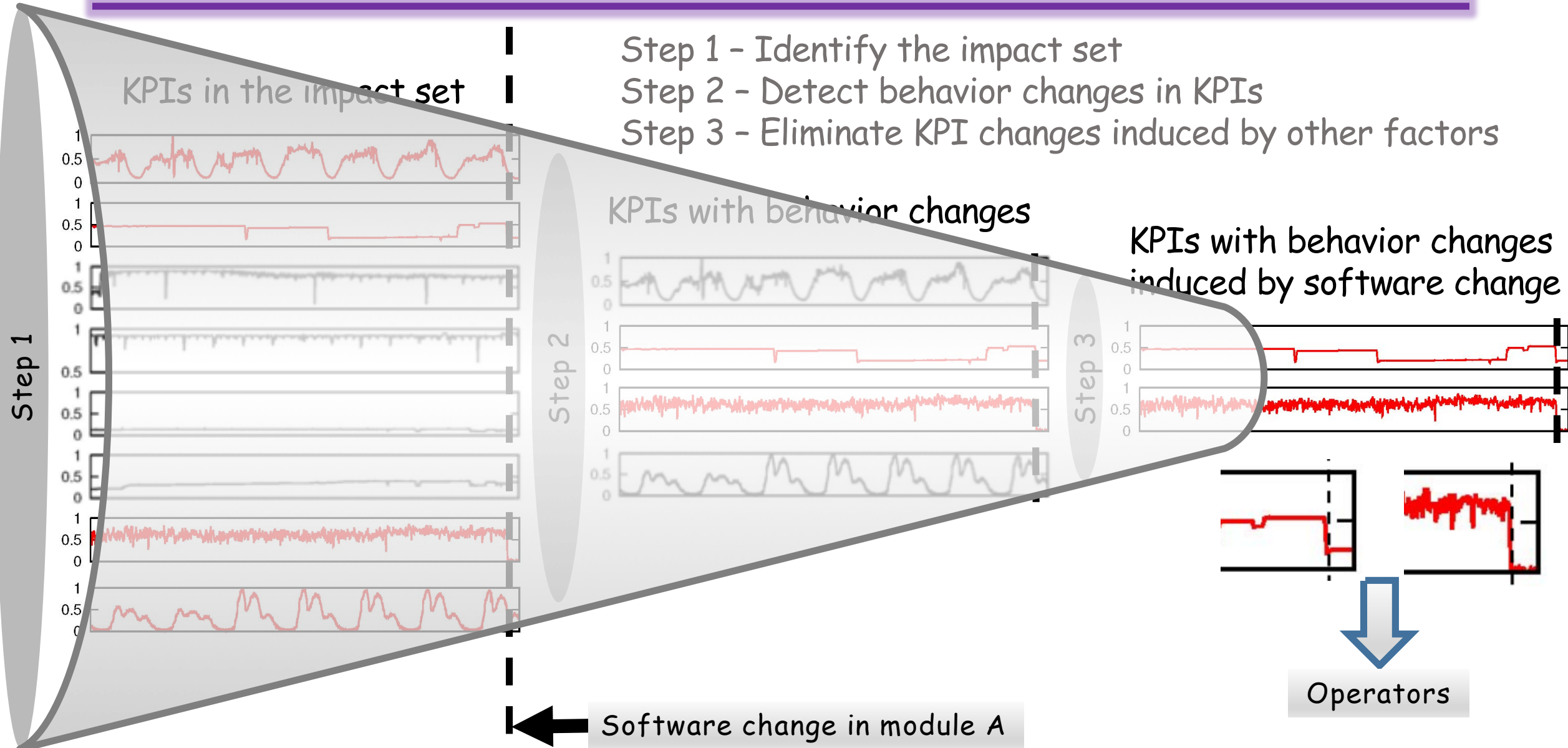
KPIs with behavior changes induced by software change



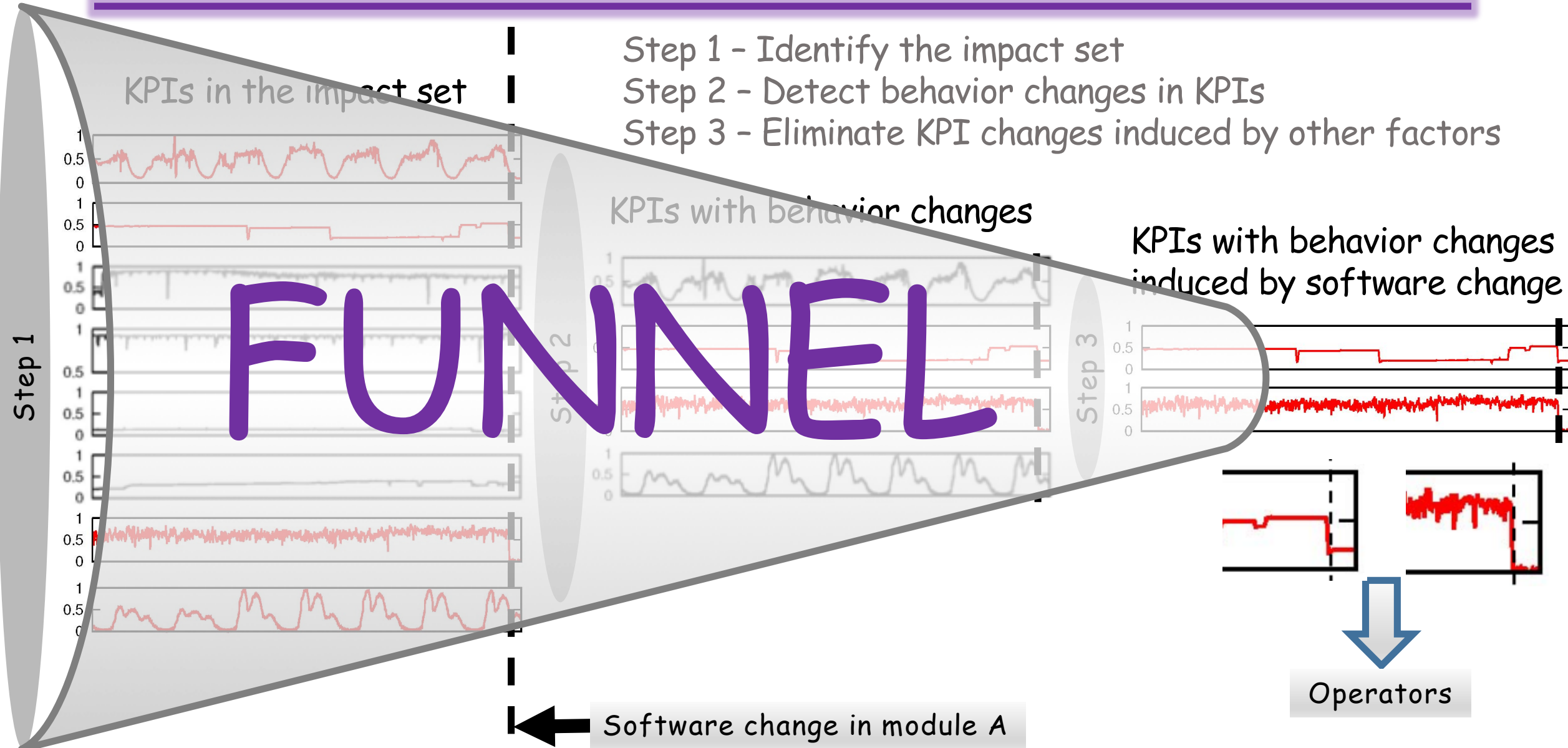
Operators

Software change in module A

Design Overview



Design Overview



Outline

- Background and Motivation
- Challenges
- Key Ideas
- Results
- Conclusion

Datasets of Evaluation

144 software changes of Baidu

72 introduced KPI changes

72 introduced no KPI changes

Datasets of Evaluation

144 software changes of Baidu

72 introduced KPI changes

72 introduced no KPI changes

Large amount of labelling work

9982 (software change,
server/module/process, KPI)s

Manually labelled by operators

Datasets of Evaluation

144 software changes of Baidu

72 introduced KPI changes

72 introduced no KPI changes

Large amount of labelling work

9982 (software change,
server/module/process, KPI)s

Manually labelled by operators

Diverse KPIs

Seasonal

Variable

Stationary

Datasets of Evaluation

144 software changes of Baidu

72 introduced KPI changes

72 introduced no KPI changes

Large amount of labelling work

9982 (software change,
server/module/process, KPI)s

Manually labelled by operators

Diverse KPIs

Seasonal

Variable

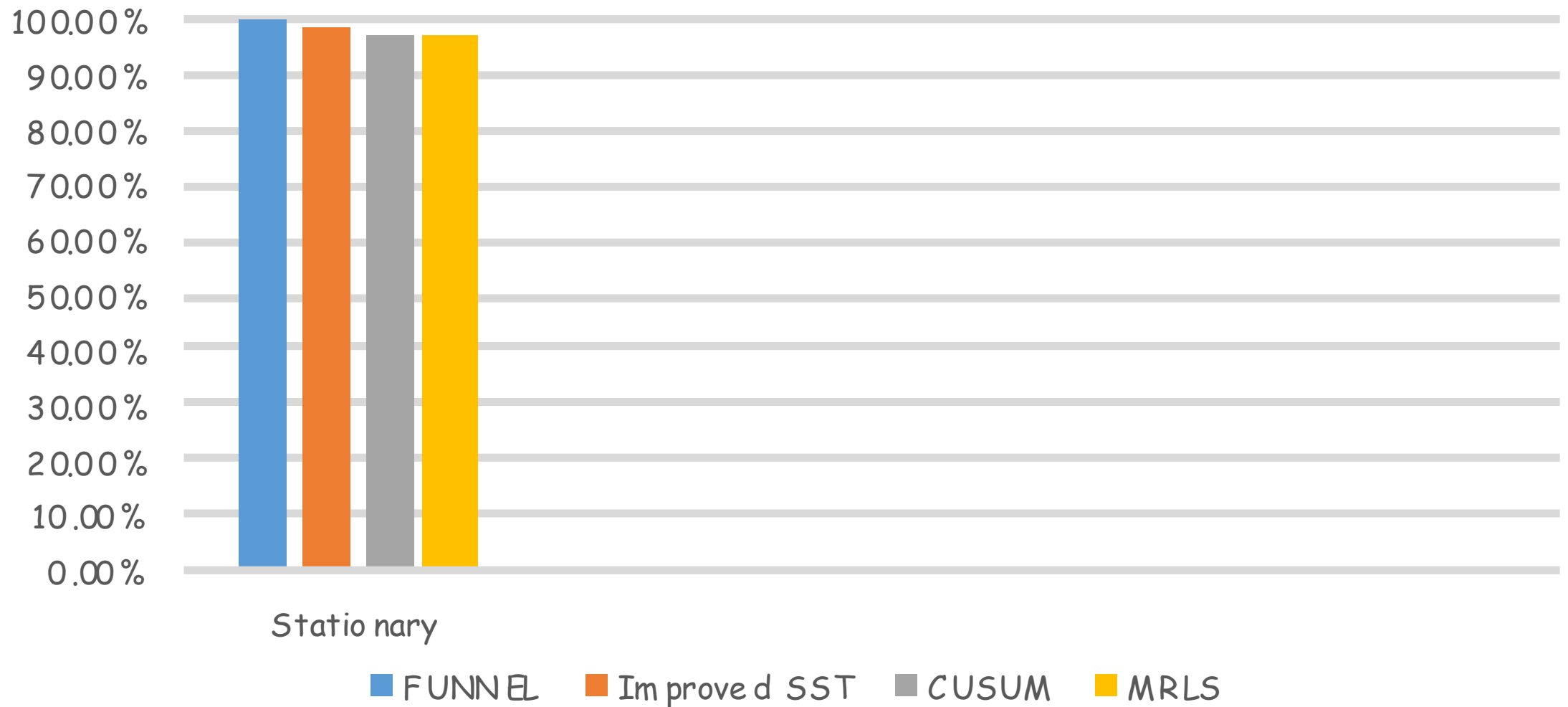
Stationary

Comparison baseline

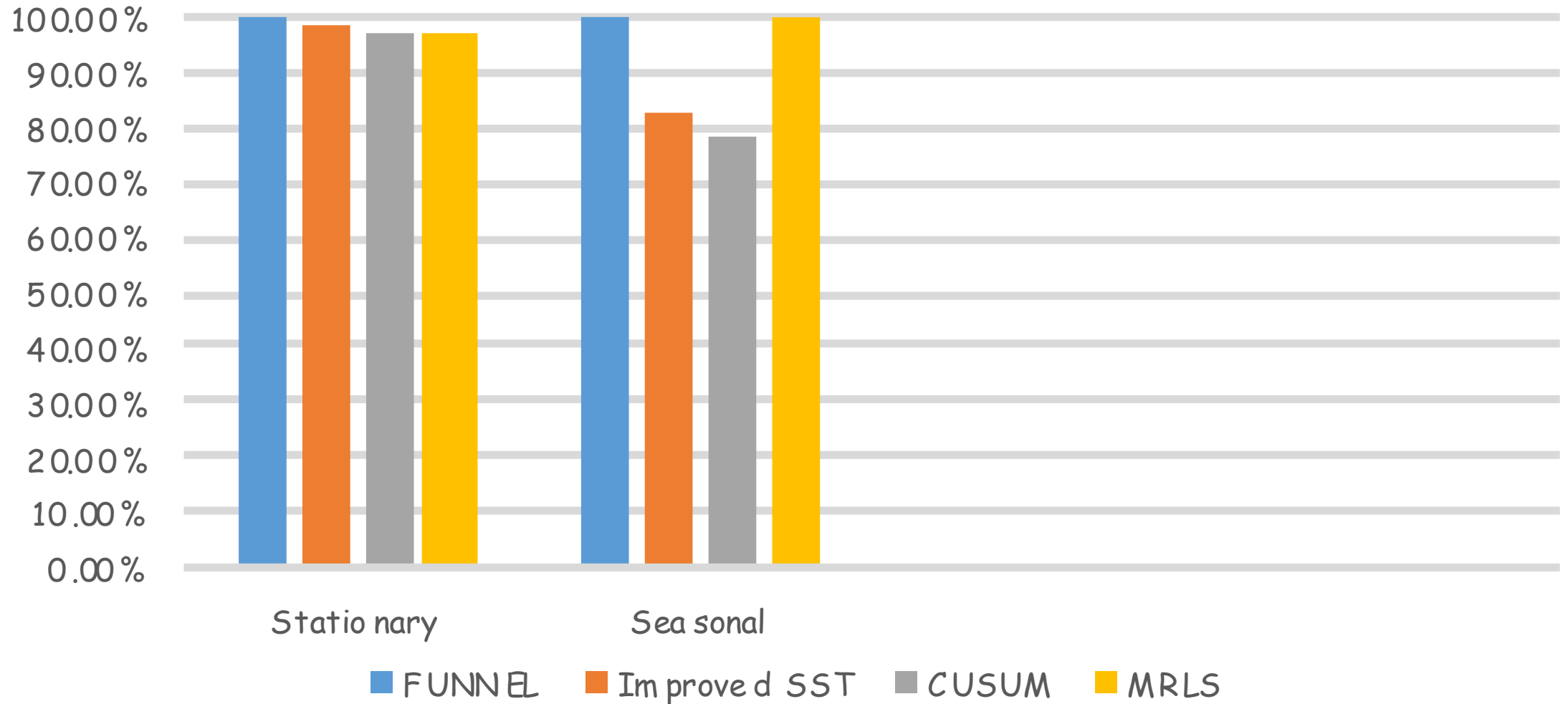
CUSUM (SIGCOMM 10)

Multiscale Robust Local Subspace
(CoNEXT 11)

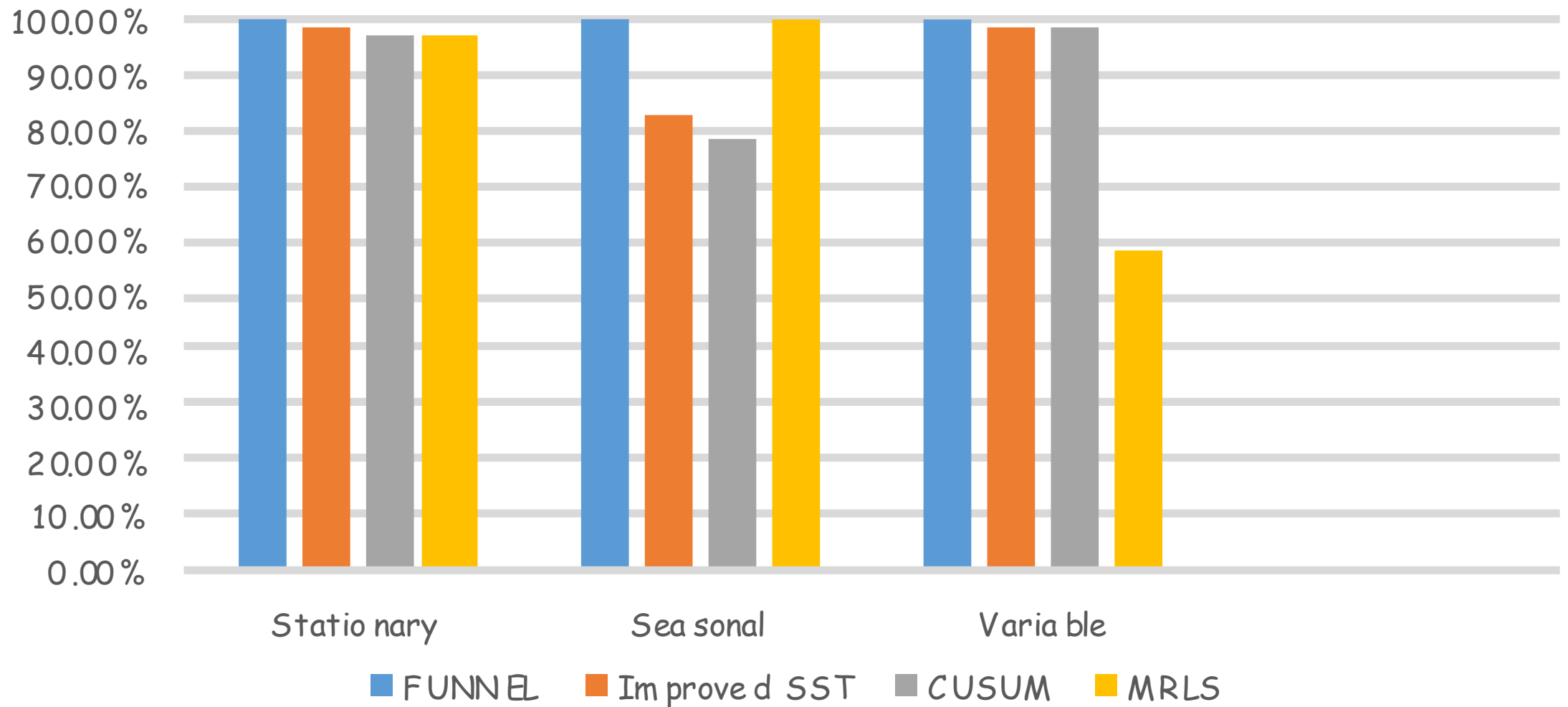
Comparison of Accuracy



Comparison of Accuracy



Comparison of Accuracy



Comparison of Computational Cost

- Real-world scenario
 - At least 1 million KPIs need to be monitored
 - The detection interval for each KPI is 1 minute
 - Runs on the same kinds of CPU as testing

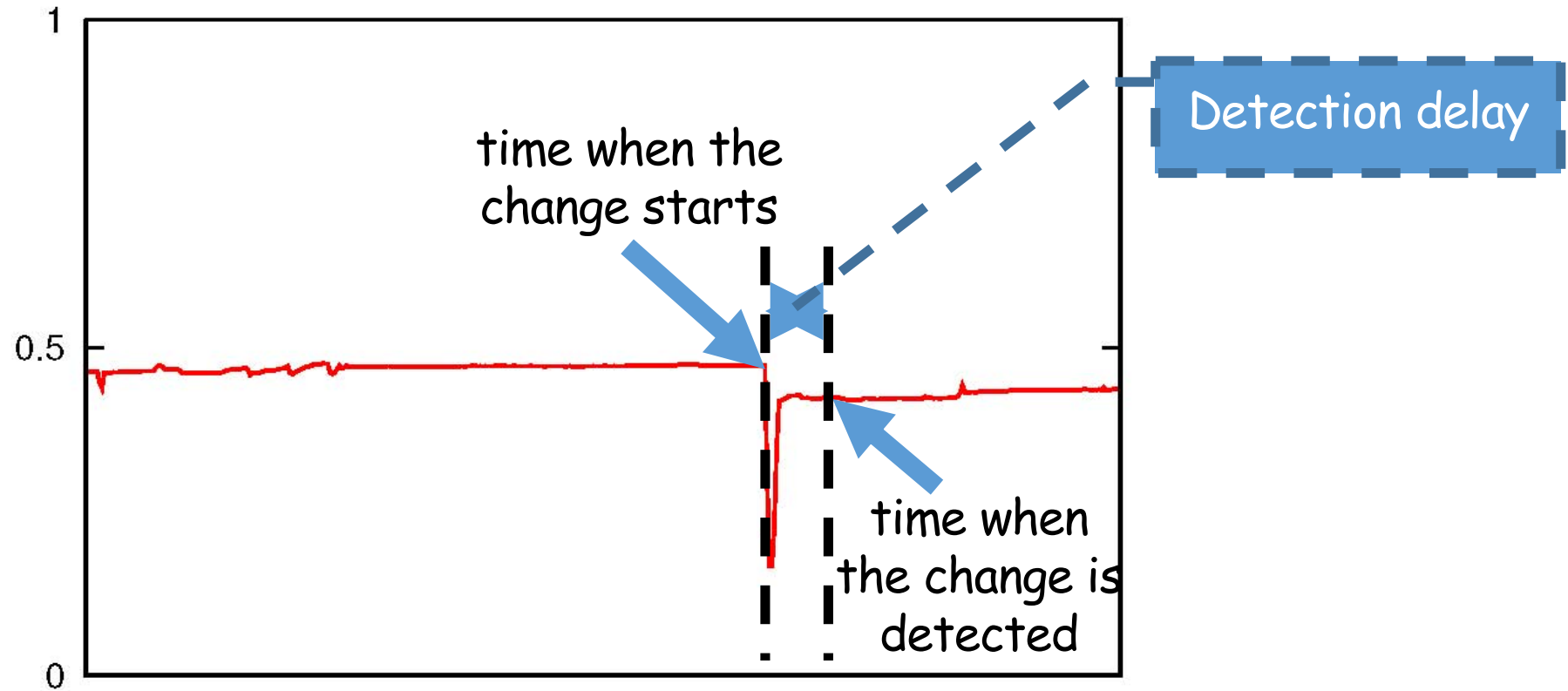
Comparison of Computational Cost

- Real-world scenario
 - At least 1 million KPIs need to be monitored
 - Each KPI is detected every 1 minute
 - Runs on the same kinds of CPU as testing
- Comparison results

Method	FUNNEL	CUSUM	MRLS
Number of cores for one million KPIs	7	31	47526

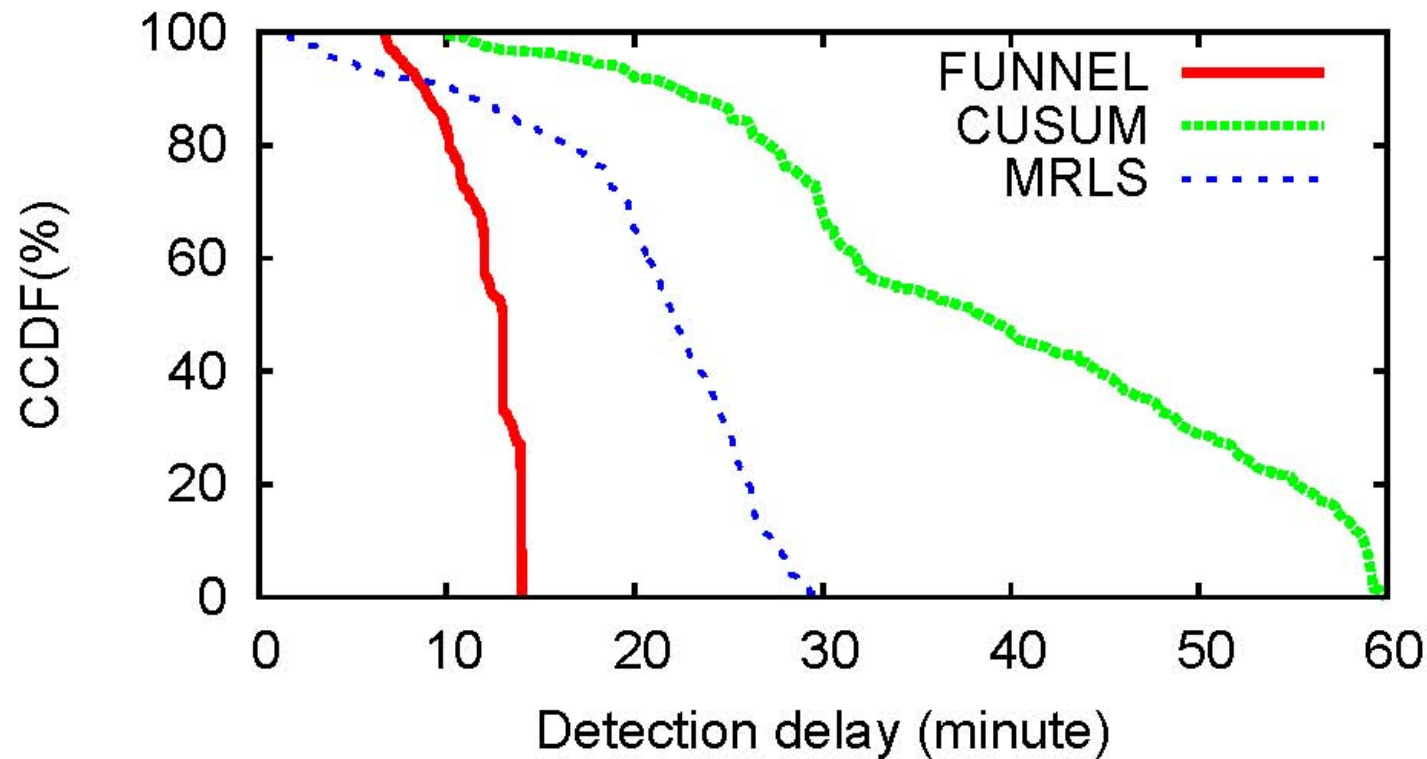
Comparison of Detection Delay

- Detection delay
 - time when a KPI change is detected - time when a KPI change starts



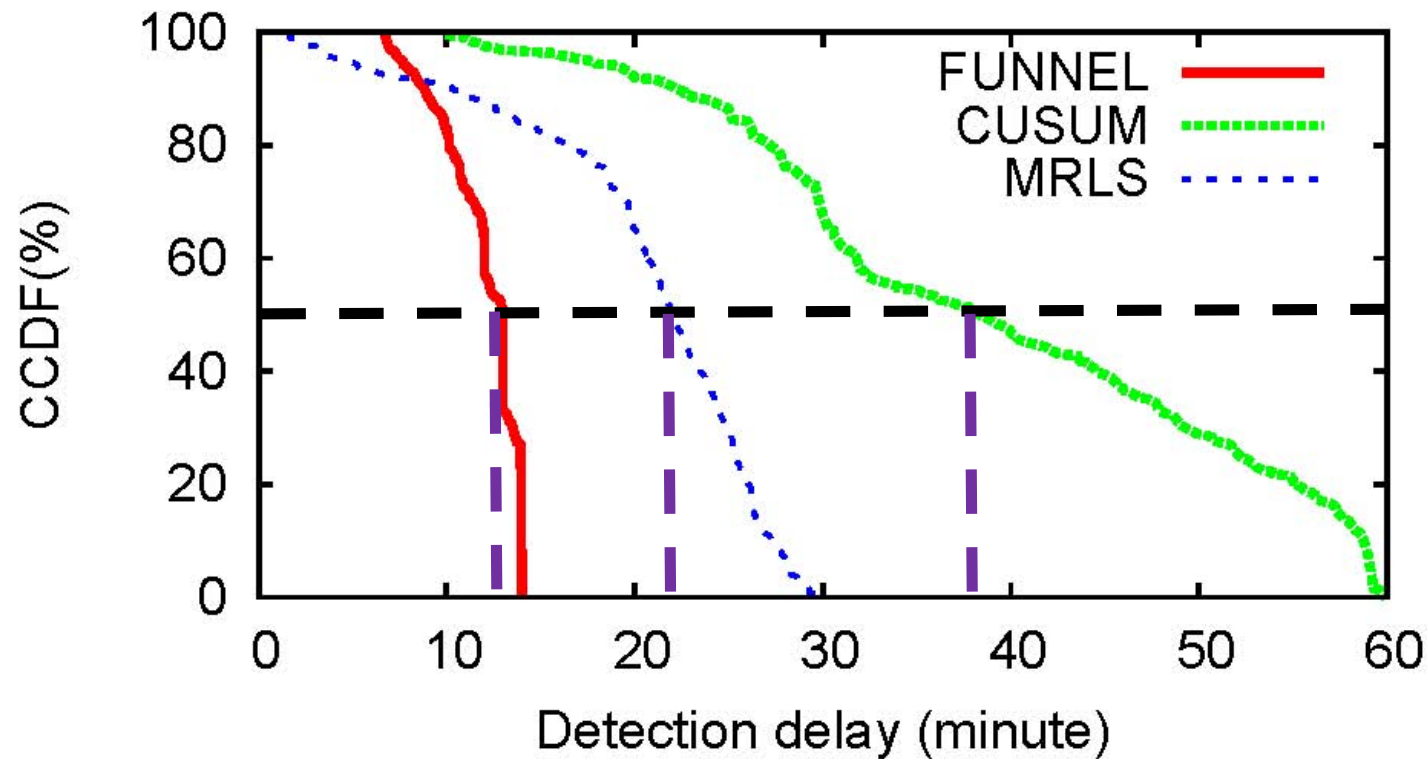
Comparison of Detection Delay

- Comparison results



Comparison of Detection Delay

- Comparison results

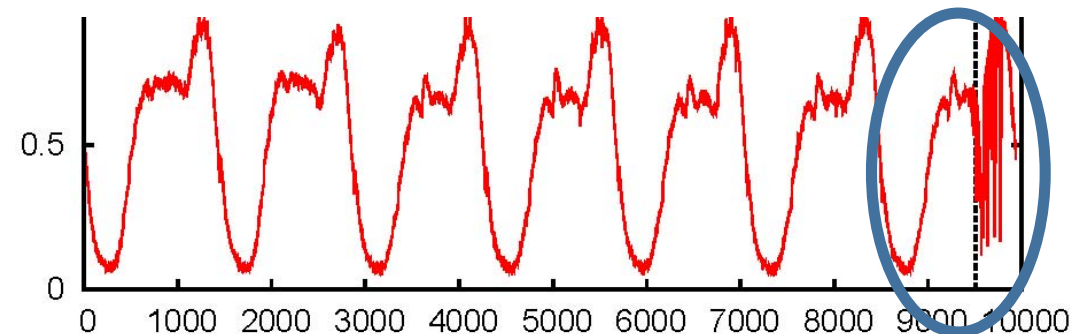
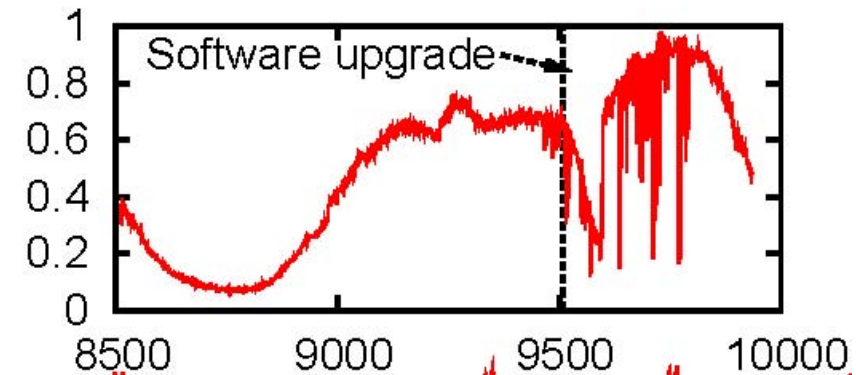


Case Study: An Erroneous Software Upgrade in Advertising

- Methodology
 - A fraction of software changes
 - Not deliver the results to the operators
 - The operators assessed the software changes independently

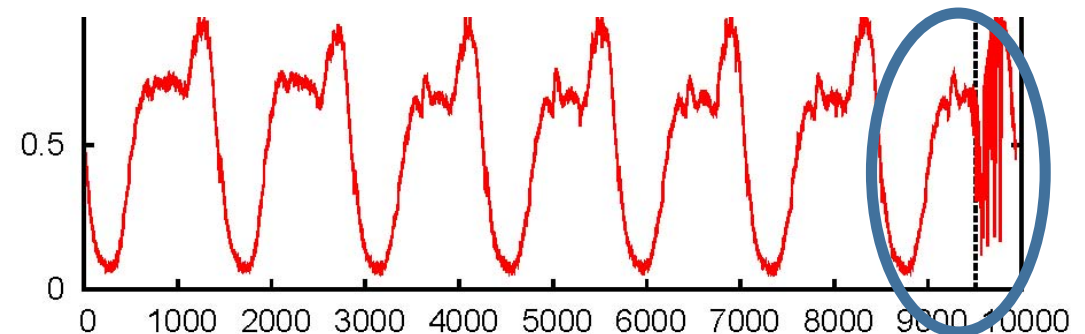
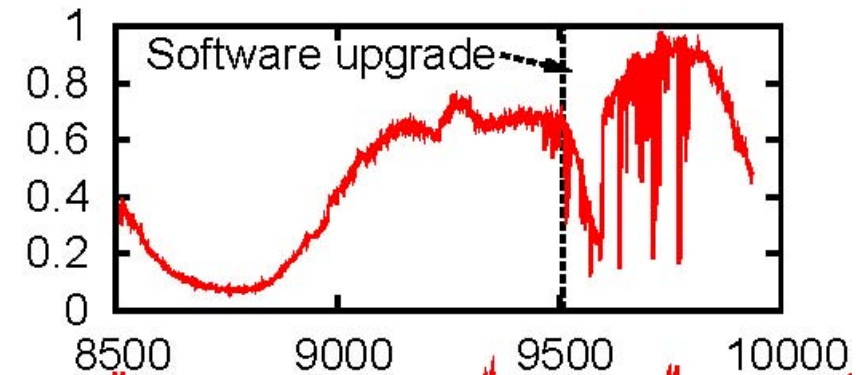
Case Study: An Erroneous Software Upgrade in Advertising

- Methodology
 - A fraction of software changes
 - Not deliver the results to the operators
 - The operators assess software changes independently
- FUNNEL
 - 10 minutes
 - Seasonal KPIs



Case Study: An Erroneous Software Upgrade in Advertising

- Methodology
 - A fraction of software changes
 - Not deliver the results to the operators
 - The operators assess software changes independently
- FUNNEL
 - 10 minutes
 - Seasonal KPIs
- The operators
 - 1.5 hours



Customer
complaints

Inspecting
KPIs

Troubleshooting

Outline

- Background and Motivation
- Challenges
- Key Ideas
- Results
- Conclusion

Conclusion

Challenges of automatic software change impact assessment

- Short detection delay requirement against robustness
- Large number of KPIs
- Diverse types of data
- KPI changes maybe caused by other factors

FUNNEL

- Improved SST - main algorithm contribution of the paper.
- Split testing

Evaluation

- Real-world software changes

Thank you!

zhangsl12@mails.tsinghua.edu.cn

Q&A

Why 144 Software Changes

- Evaluation needs ground truth
 - FUNNEL
 - detect KPI changes
 - determine whether KPIs changes are induced by software change
 - Operators
 - Label whether there is behavior change in KPI
 - Label whether a KPI changes is caused by software change
- 9982 (software change, server/module/process, KPI)s
- A huge amount of work
- Labelling for much more software changes is prohibitive

Why Using Cores

- The CPU utilization is 100% in testing
- Assume the CPU utilization is also 100% in deployment
- The operators care about how many servers/cores the system needs

Why just a single team

- For the efficiency purpose
- Build a single database to monitor all KPIs
- By natural

Unbalanced hotspot

- Split testing
- The number of hotspots is very small (3% in Microsoft)
- Compare the treated group and the control group
- The large number of KPIs in the control group makes the determination robust even in the face of hotspots.

The parameters of FUNNEL, CUSUM and MRLS

- Two parameters
 - a in DiD method
 - w in Improved SST
- Best for accuracy
- Operators care most about the accuracy
- Fair for the four methods

About the detection delay comparison

- Set a threshold for FUNNEL
- MRLS can detect behavior changes with smaller detection delay than FUNNEL at sometimes
- Sacrificing the accuracy

Why not Just Split Testing?

- Set threshold small
 - Sensitive to spikes
 - Many false positives
- Set threshold large
 - The detection delay is large
- Almost impossible to find a balance in our scenario
- The improved SST
 - Robust
 - Short detection delay

Obtain the Relationship of Modules

- The operators name the modules based on the module hierarchy
- The operators know the relationship of modules

Why not decide to roll out/back by FUNNEL?

- The KPI changes & the decision
 - Hard to learn
 - Few cases for a specific combination of KPI change and software change
- Rolling back a software change is a big thing
 - The operators would like to decide themselves.
- FUNNEL is helpful for the operators to make decision
 - The number of KPIs with behavior changes induced by software changes is small
 - The work of the operators is small.

About the Deployment

- Assess the software changes of a few dozens of Internet-based services

Number of software changes	Number of changes that have impact	Number of KPIs	Number of KPI changes	Precision
24119	268	2256390	10249	98.21%

If A Software Change is Deployed to All Servers ...

- Treated group
 - Measurements of KPIs in the impact set around the software change
- Control group
 - Measurements of KPIs in the impact set in the same period but on historical days

About the Number of Software Changes

- If a software change is deployed on a subset of servers firstly, and then on another subset of servers
- From the operators' perspective
- They are two software changes