

Examination of Performance Changes (at Code Level)

David Georg Reichelt^{1,2}

¹School of Computing and Communications, Lancaster University Leipzig,

²Universitätsrechenzentrum, Universität Leipzig

BEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung

15.4.2023

Why measure performance? [SSP19]

Which is faster?

```
1  StringBuilder buf
2    = new StringBuilder(16);
3  buf.append("Hello_ World");
4
5
6  return buf.toString();
```

```
1  StringBuilder buf
2    = new StringBuilder(16);
3  buf.append("Hello")
4    .append("_")
5    .append("World");
6  return buf.toString();
```

Why measure performance? [SSP19]

Which is faster?

```
1  StringBuilder buf
2    = new StringBuilder(16);
3  buf.append("Hello_ World");
4
5
6  return buf.toString();
```

```
1  StringBuilder buf
2    = new StringBuilder(16);
3  buf.append("Hello")
4    .append("_")
5    .append("World");
6  return buf.toString();
```

Why measure performance? [SSP19]

Which is faster?

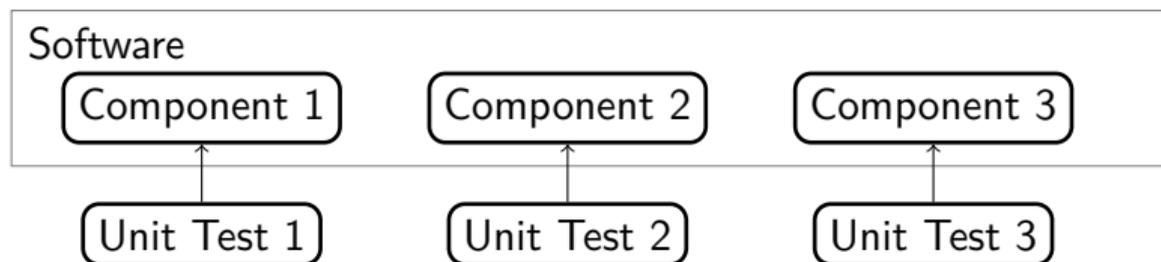
```
1 StringBuilder buf
2   = new StringBu
3 buf.append("Hell
4
5
6 return buf.toStr
```

Complex behaviour
⇒ impact of changes
needs to be measured

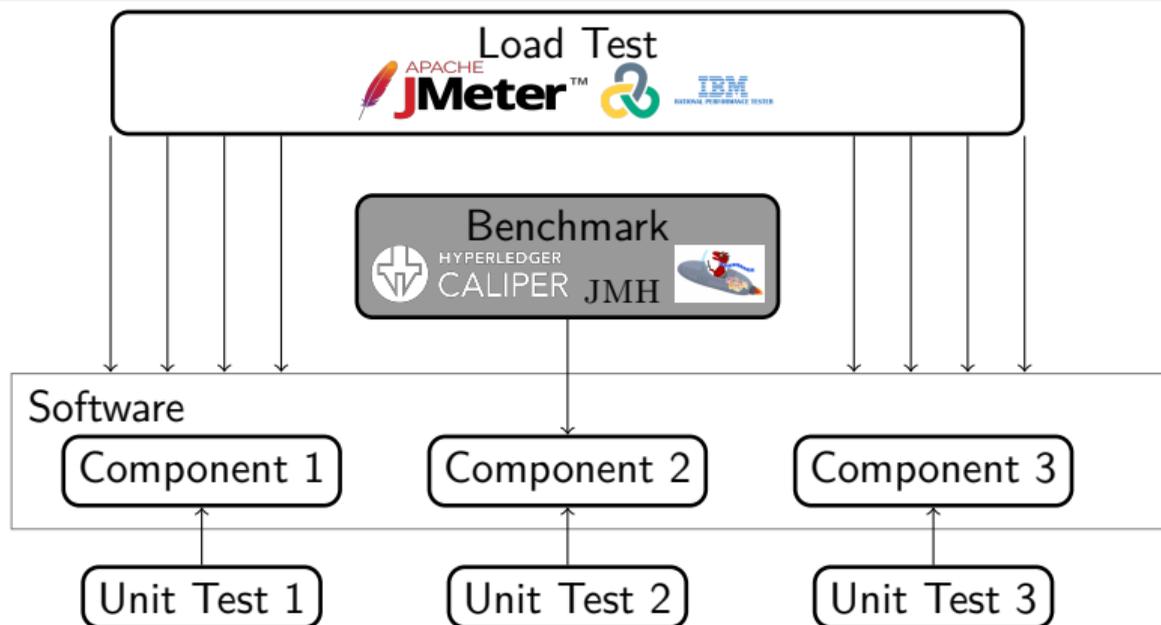
```
er buf
ingBuilder(16);
"Hello")
")
World");
toString();
```

Peass Method

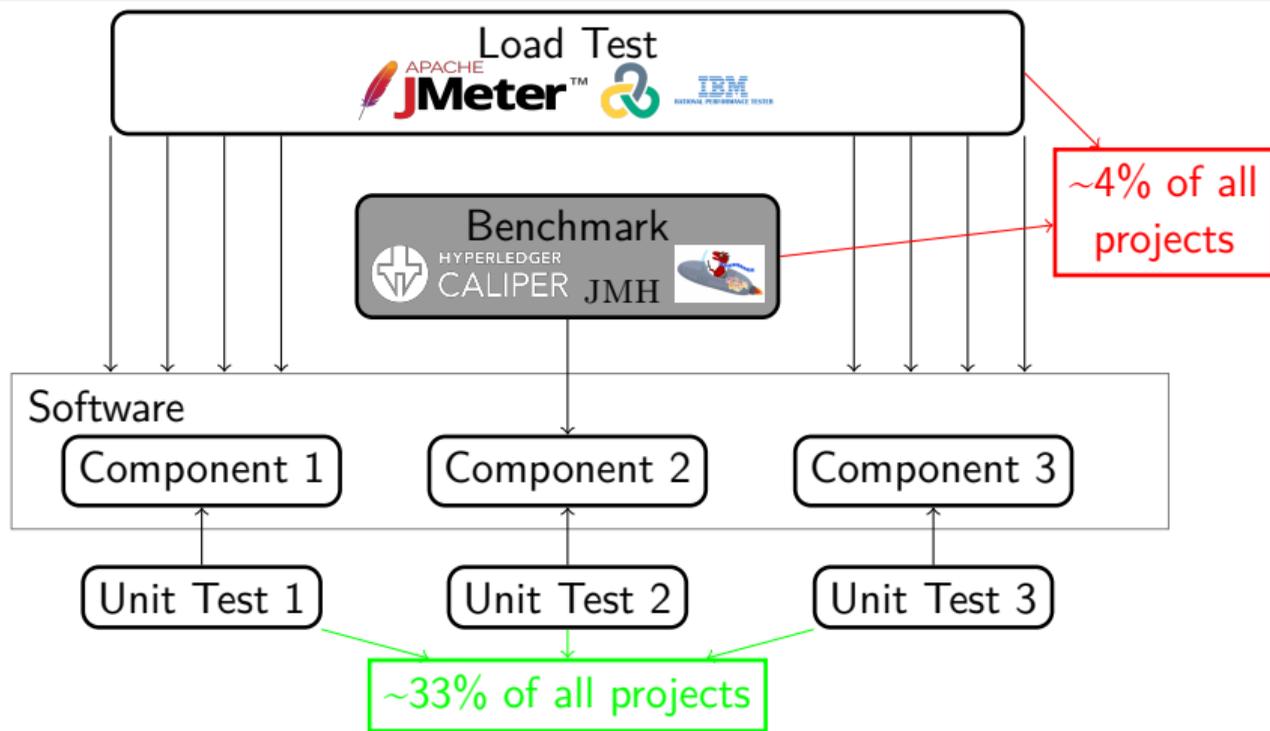
Why Unit Test Measurement?



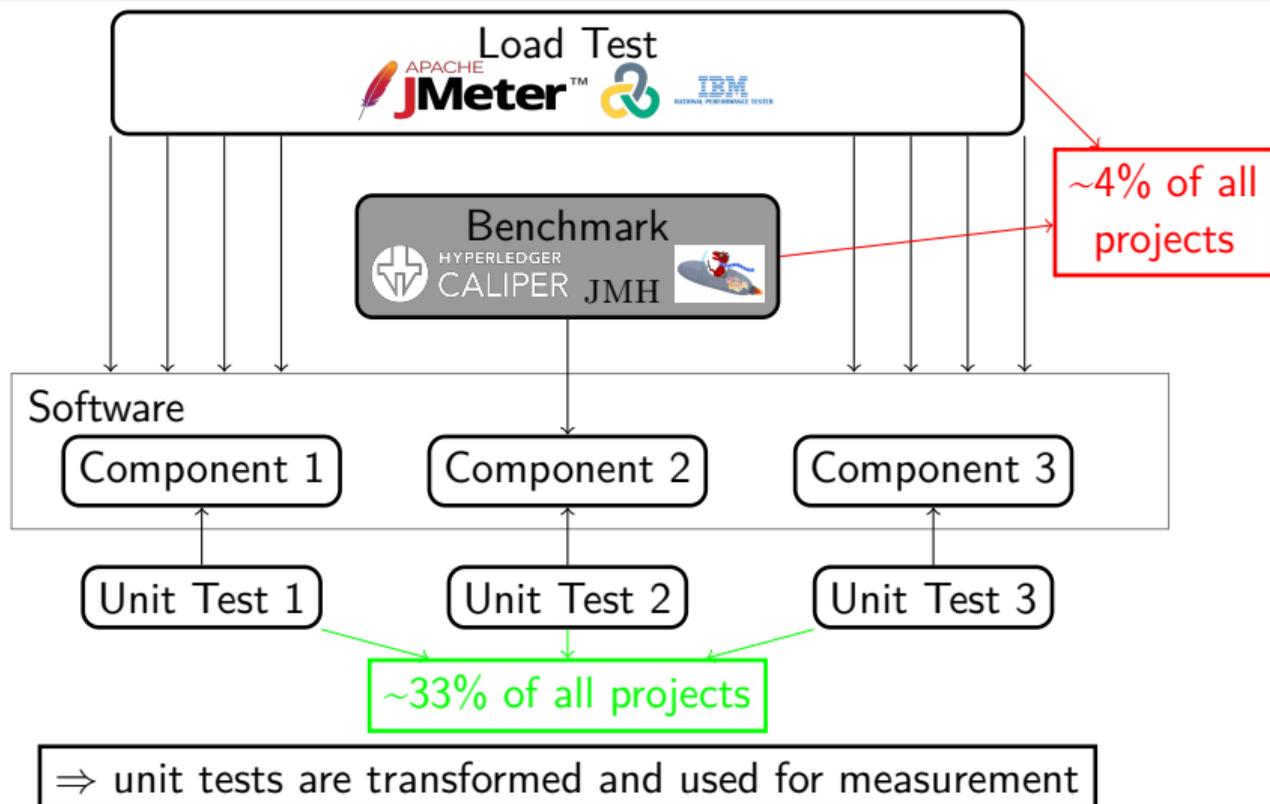
Why Unit Test Measurement?



Why Unit Test Measurement?



Why Unit Test Measurement?

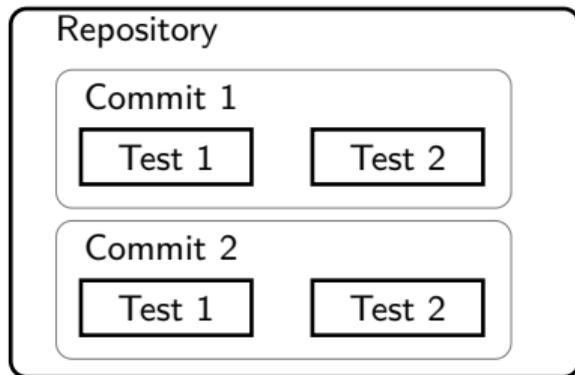


Peass

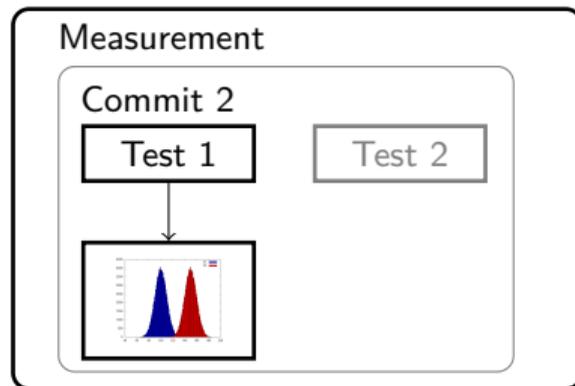
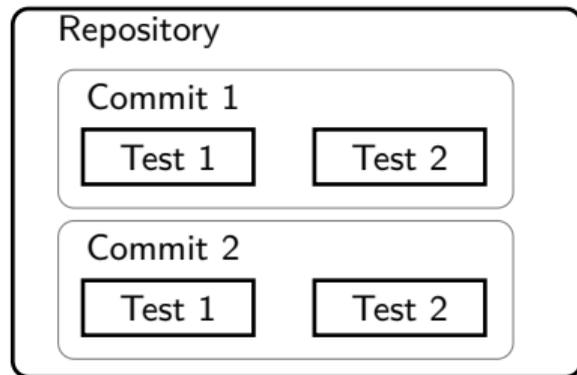


- ▶ **Performance analysis of software systems**
- ▶ CLI-Tool: <https://github.com/DaGeRe/peass>
- ▶ Jenkins-Plugin: <https://github.com/jenkinsci/peass-ci-plugin>

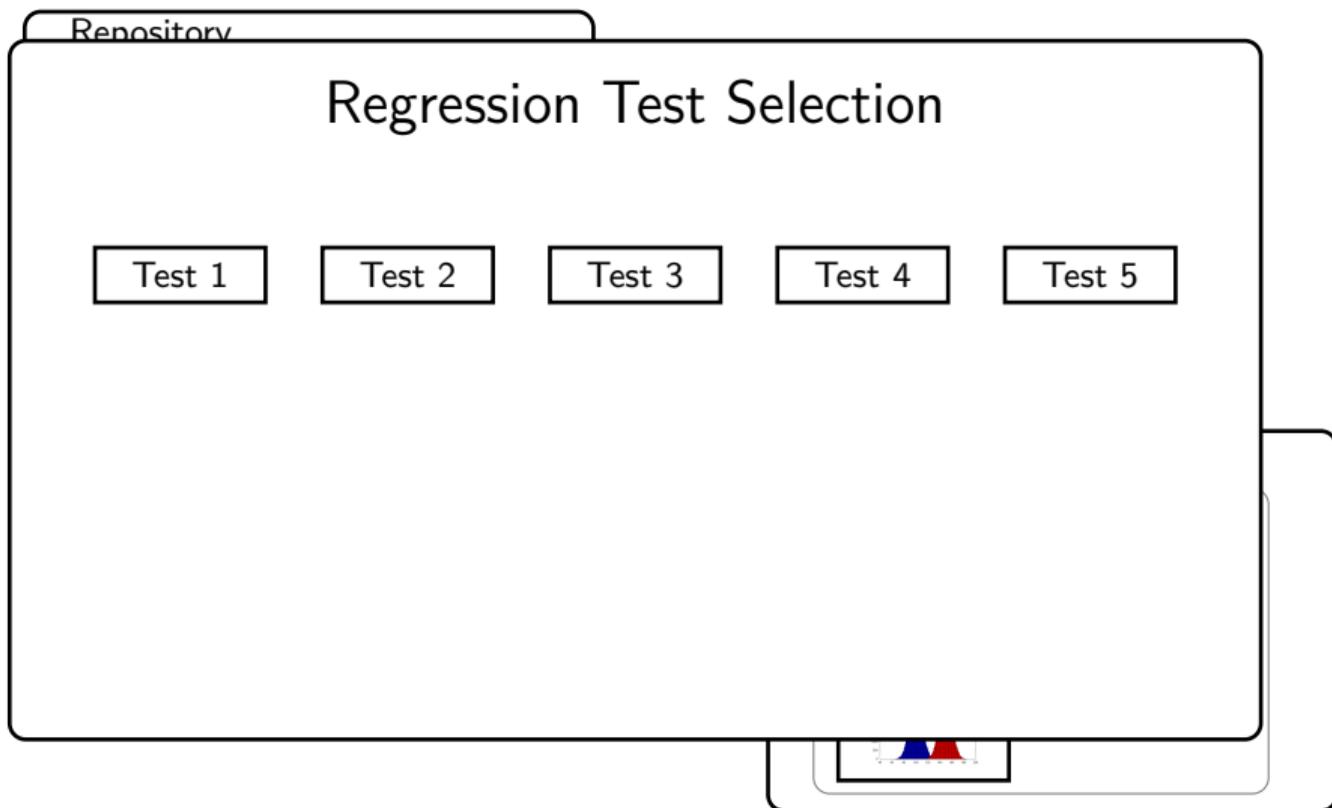
Process of Peass



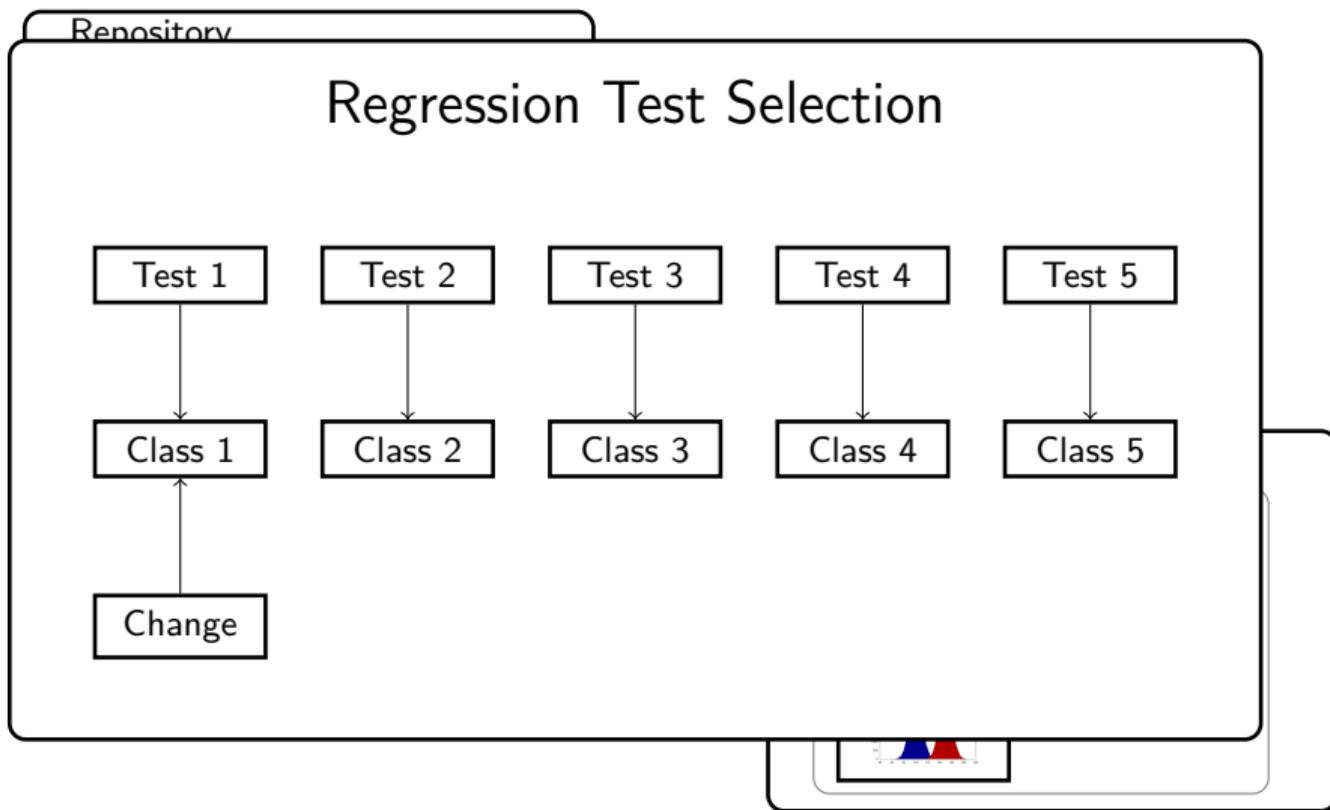
Process of Peass



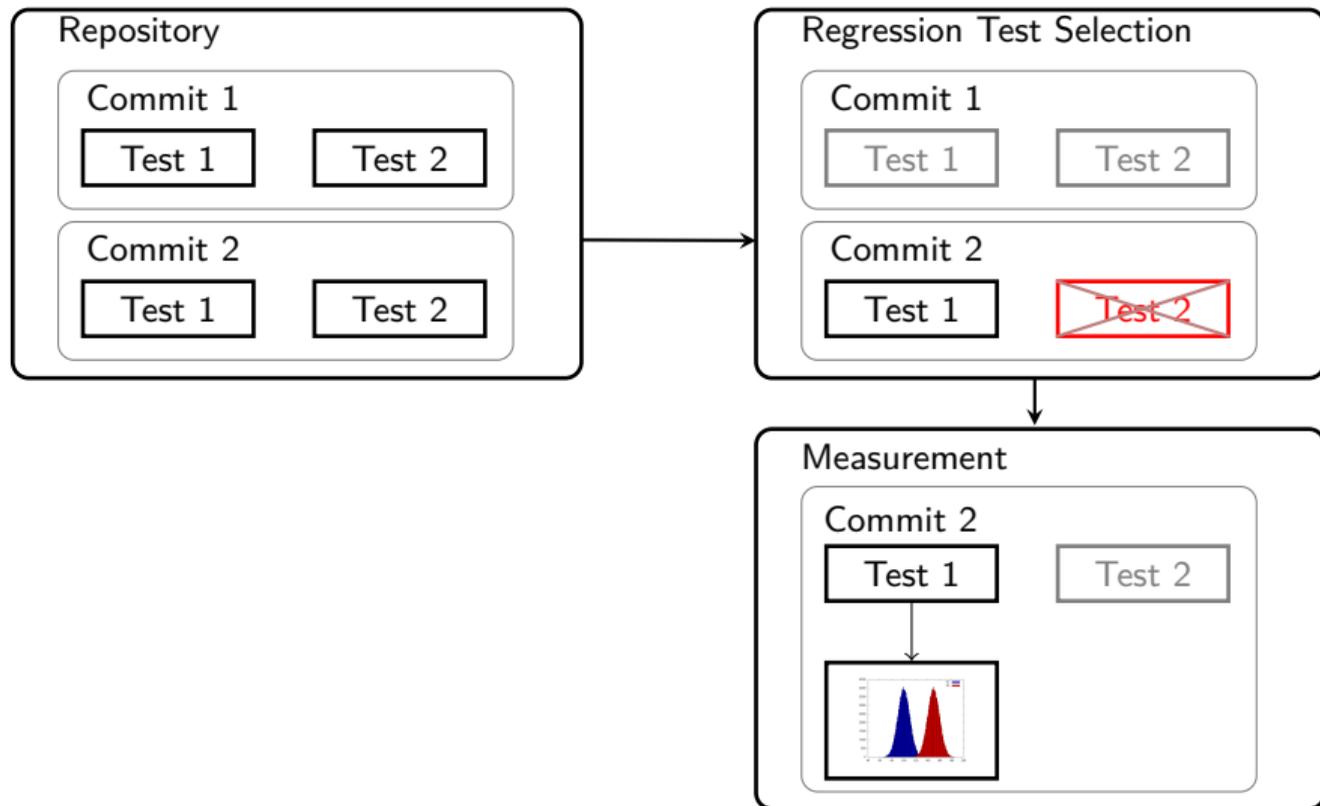
Process of Peass



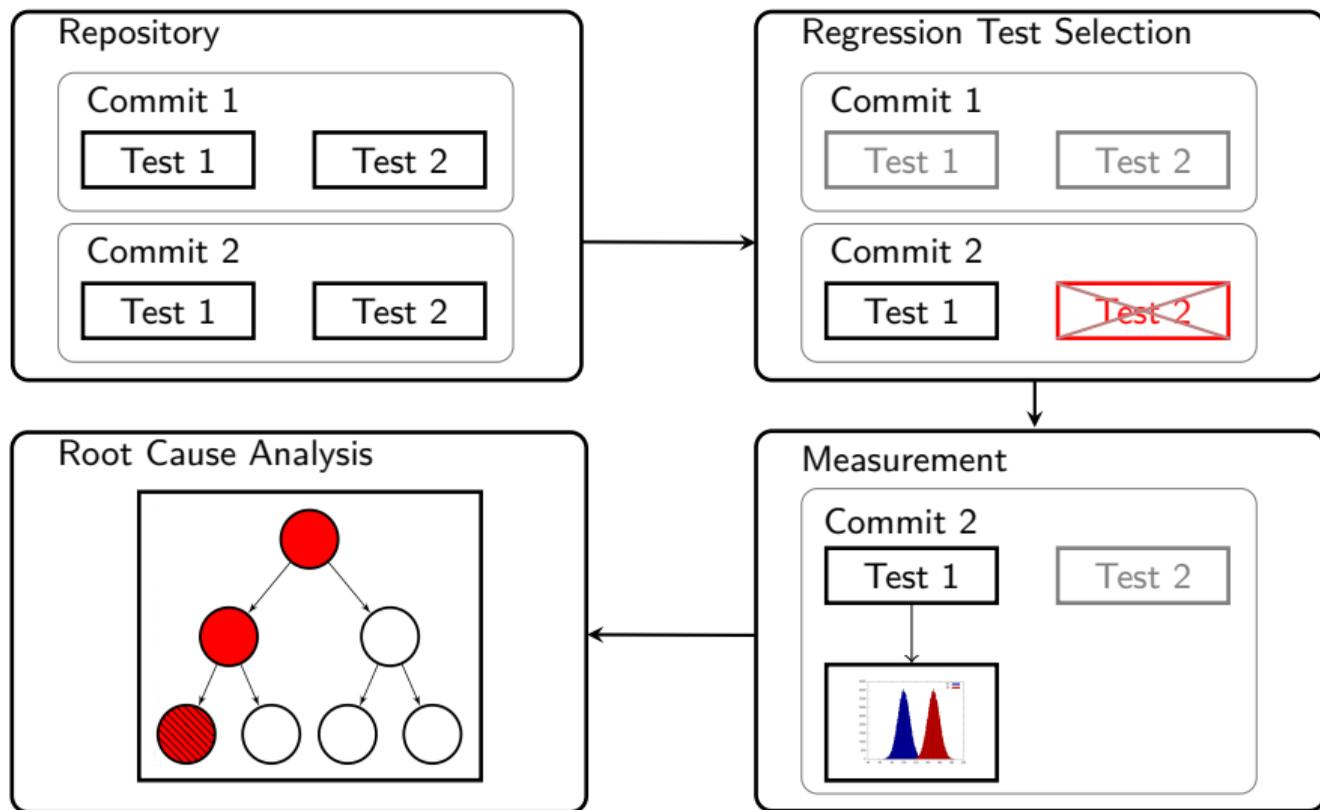
Process of Peass



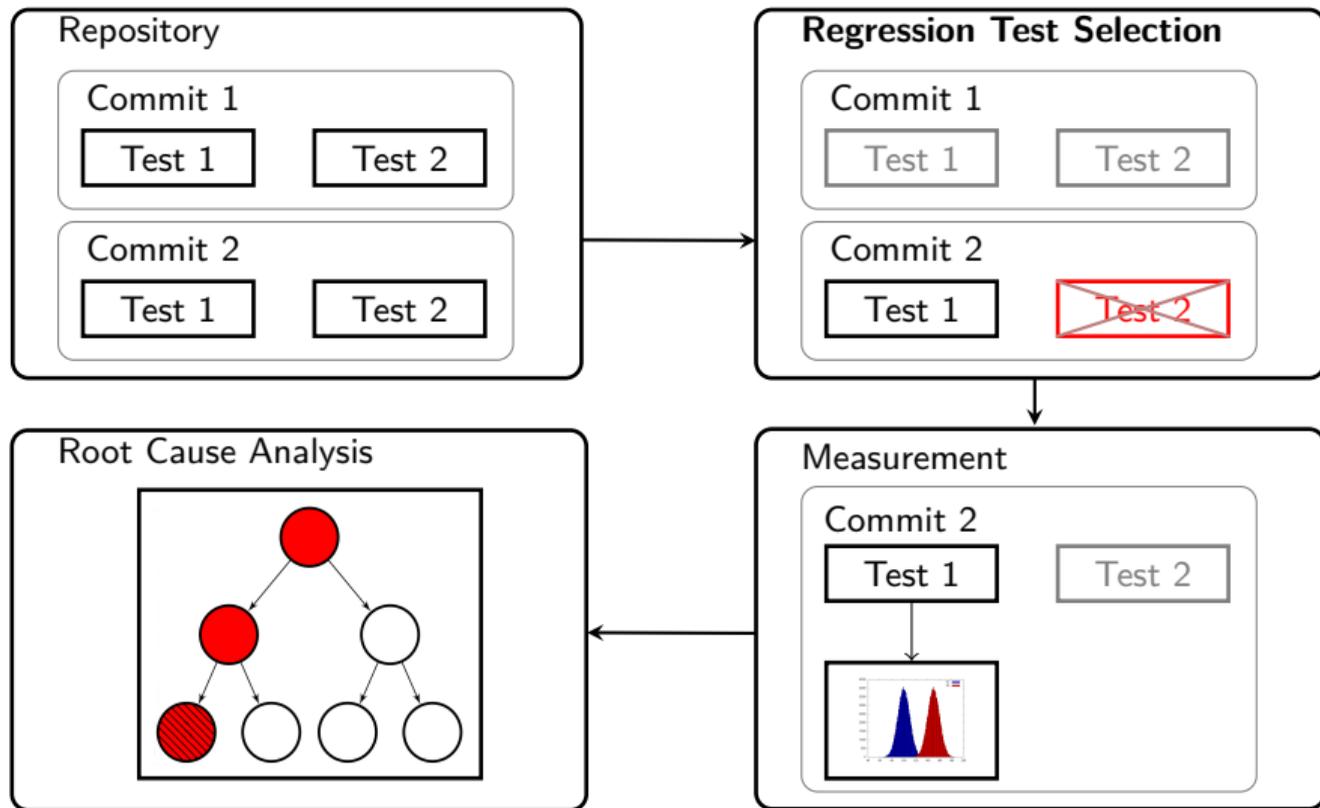
Process of Peass



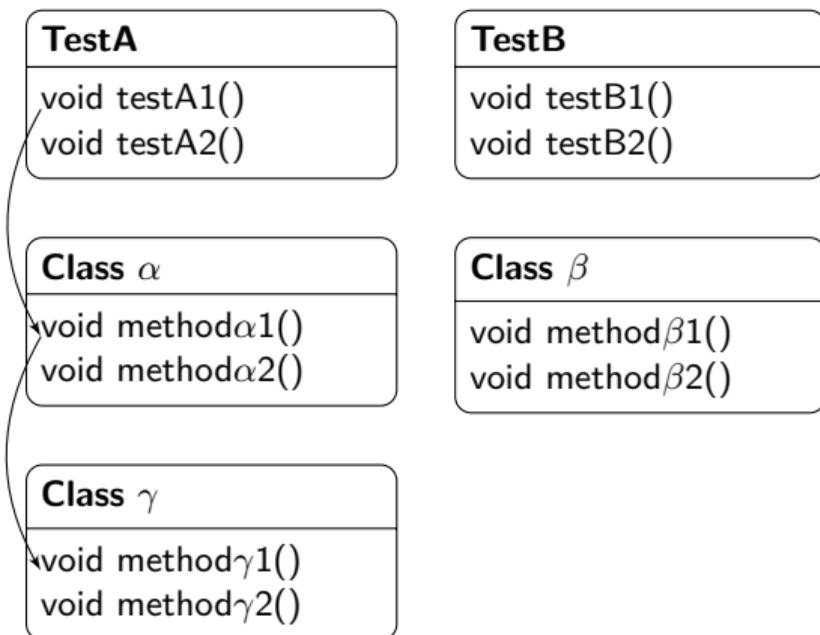
Process of Peass



Process of Peass



Regression Test Selection - Approach [LTB18]



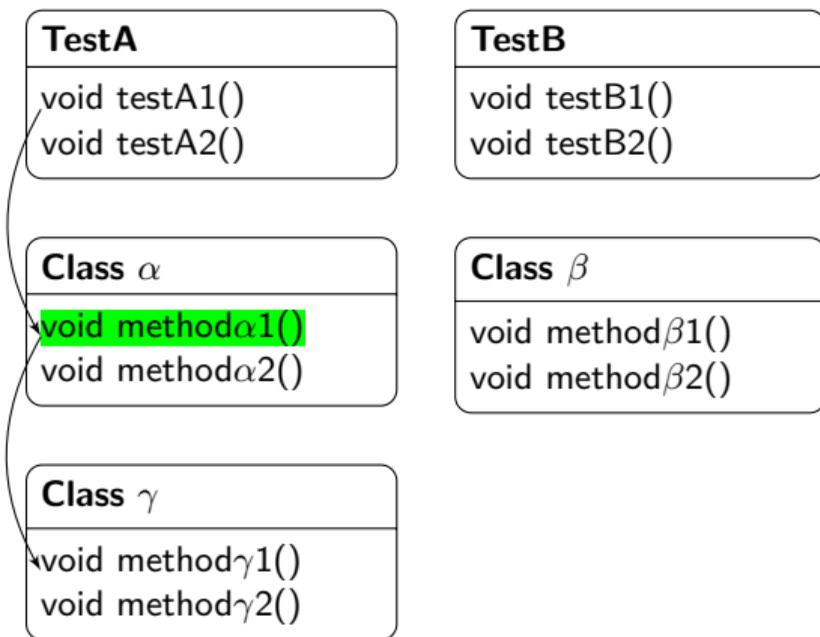
Trace:

TestA.testA1
Class α .method α 1
Class γ .method γ 1
Class α .method α 1
Class γ .method γ 1

Change:

...

Regression Test Selection - Approach [LTB18]



Trace:

TestA.testA1

Class α .method α 1

Class γ .method γ 1

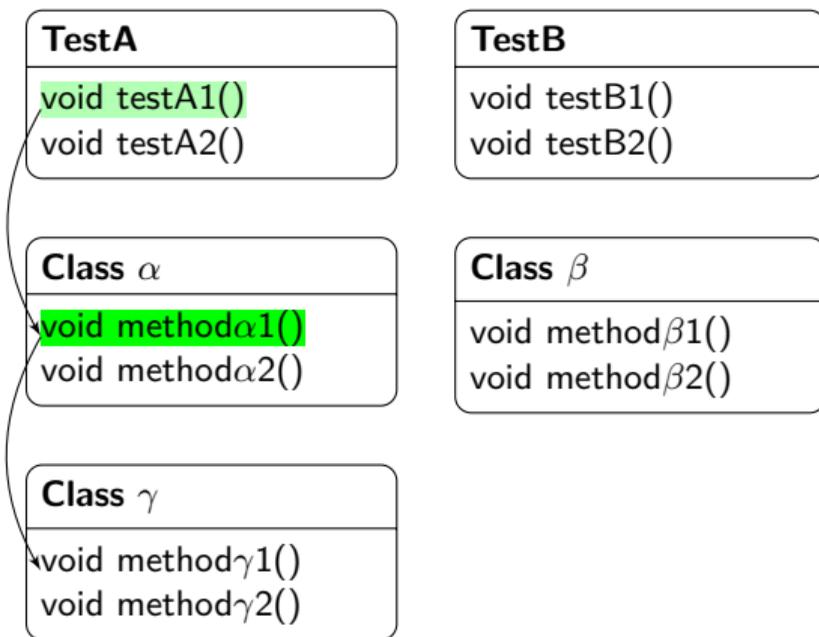
Class α .method α 1

Class γ .method γ 1

Change:

Class α .method α 1

Regression Test Selection - Approach [LTB18]



Trace:

TestA.testA1

Class α .method α 1

Class γ .method γ 1

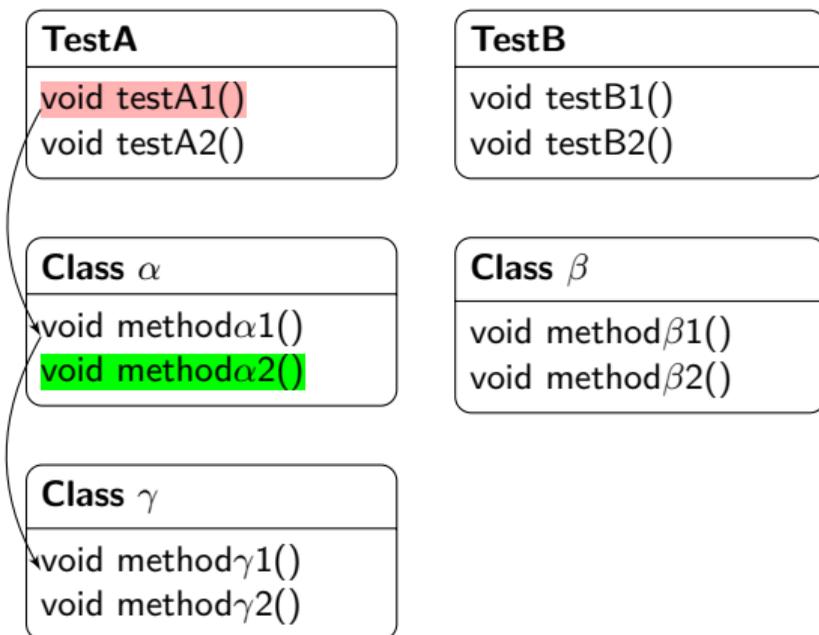
Class α .method α 1

Class γ .method γ 1

Change:

Class α .method α 1

Regression Test Selection - Approach [LTB18]



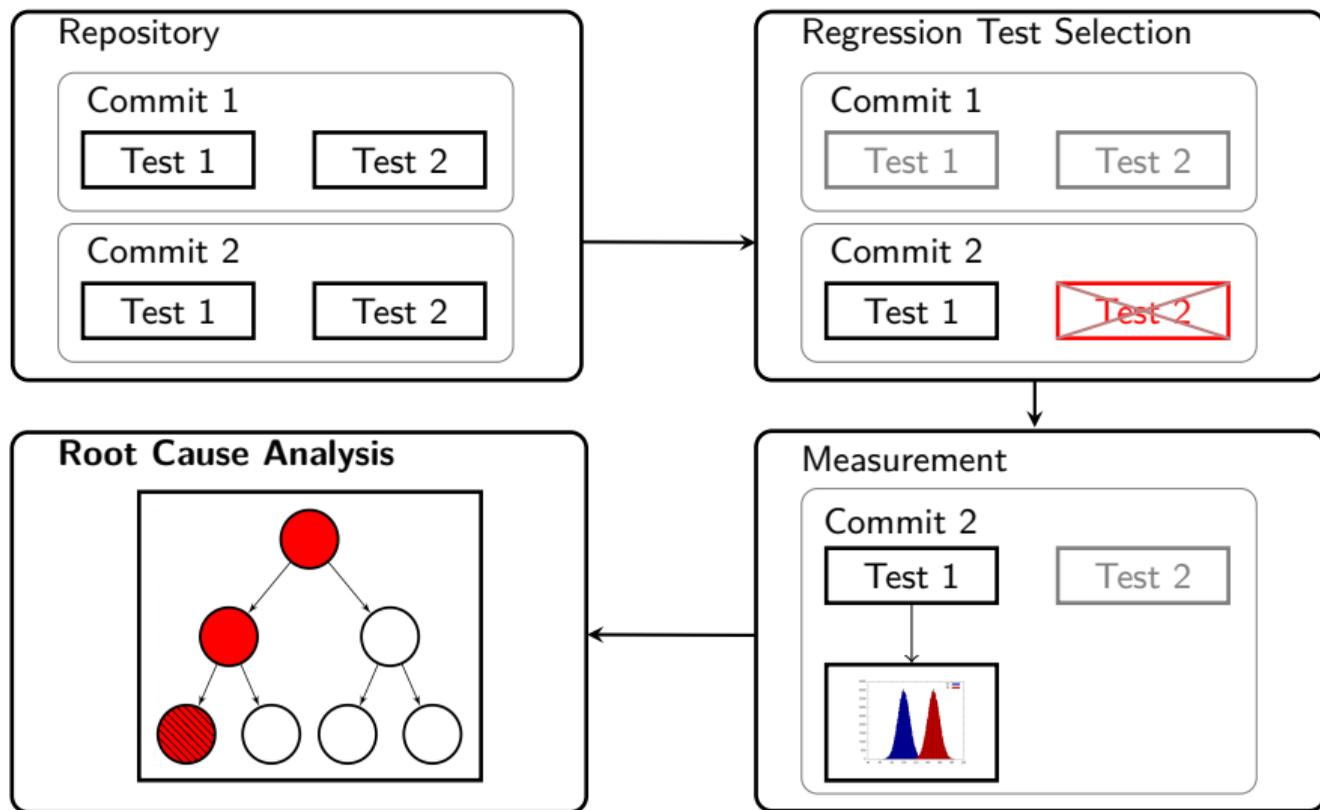
Trace:

```
TestA.testA1  
Class $\alpha$ .method $\alpha$ 1  
Class $\gamma$ .method $\gamma$ 1  
Class $\alpha$ .method $\alpha$ 1  
Class $\gamma$ .method $\gamma$ 1
```

Change:

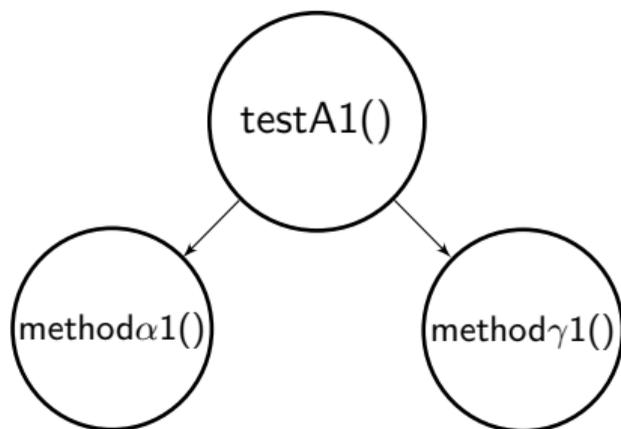
```
Class $\alpha$ .method $\alpha$ 2
```

Process of Peass



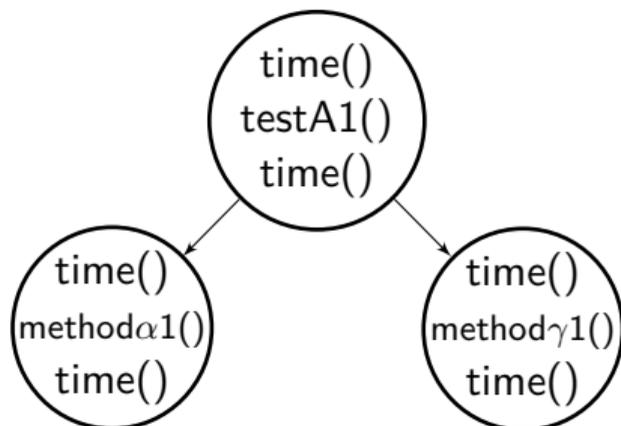
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree



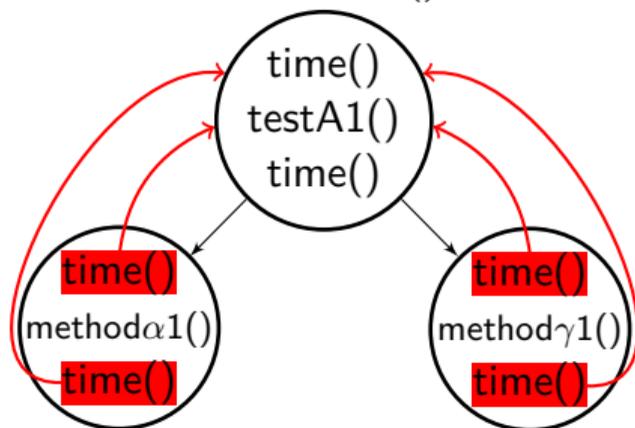
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes



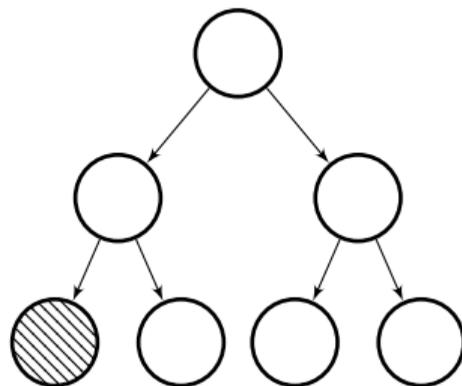
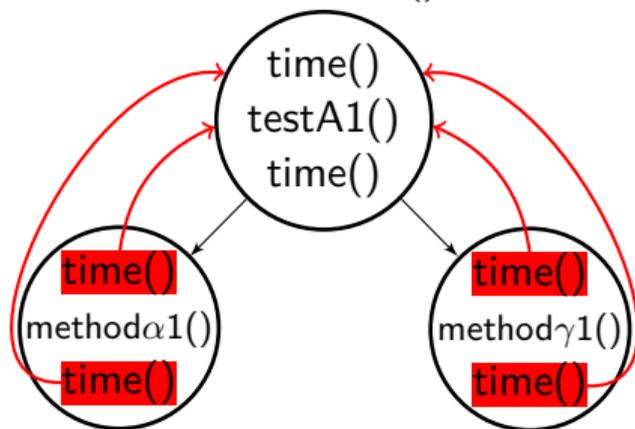
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



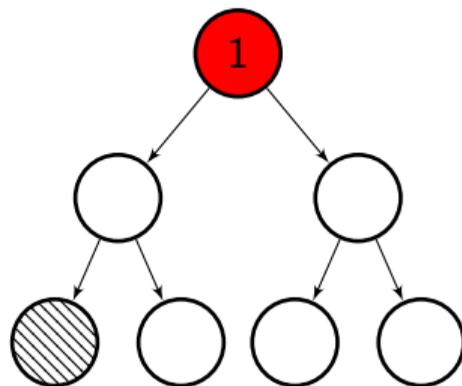
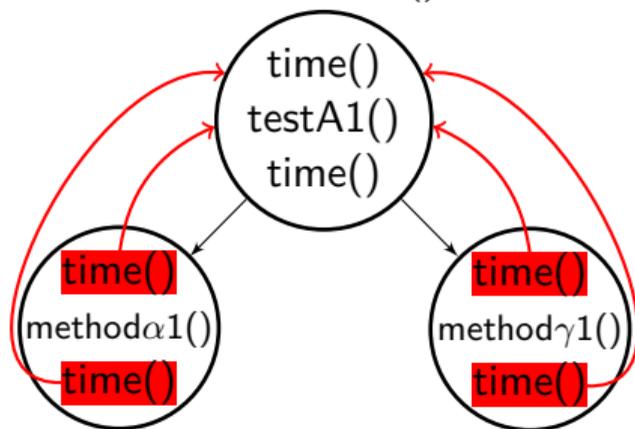
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



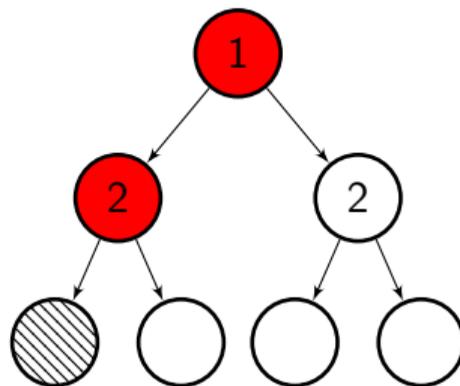
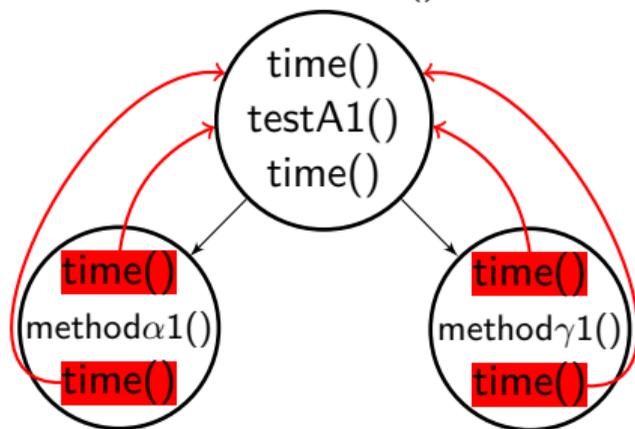
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



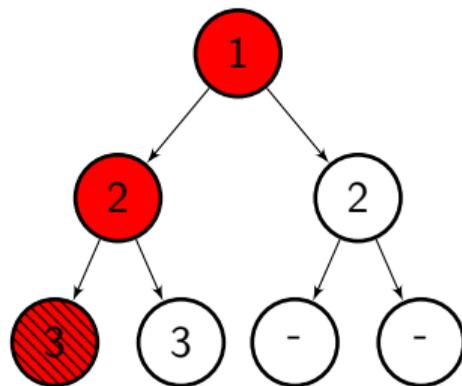
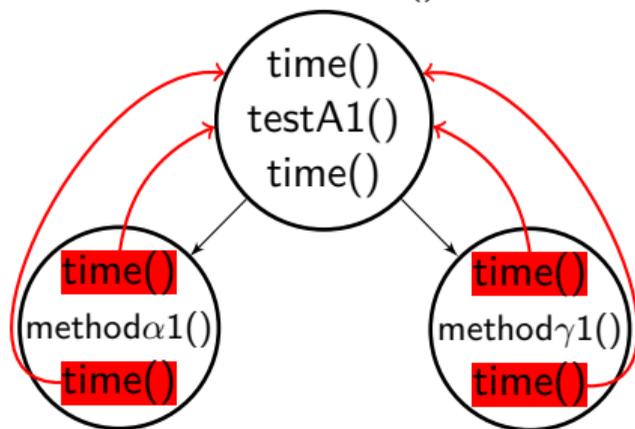
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



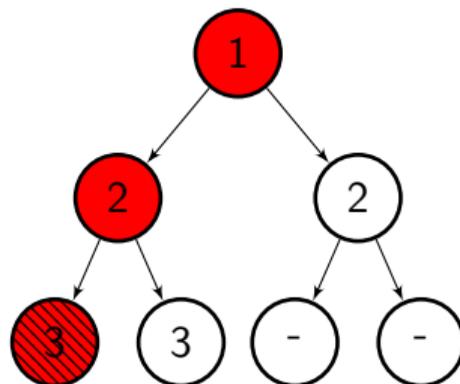
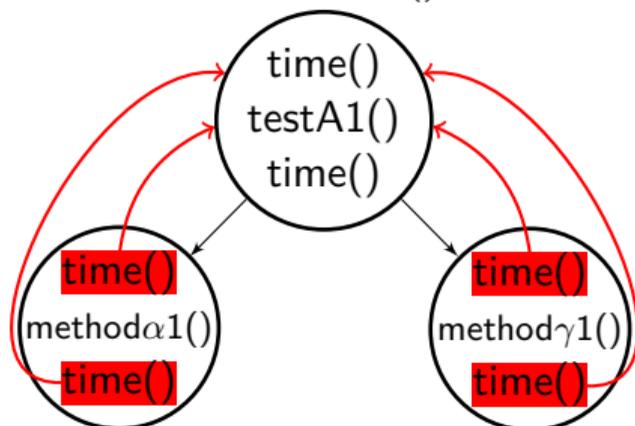
Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



Root Cause Analysis

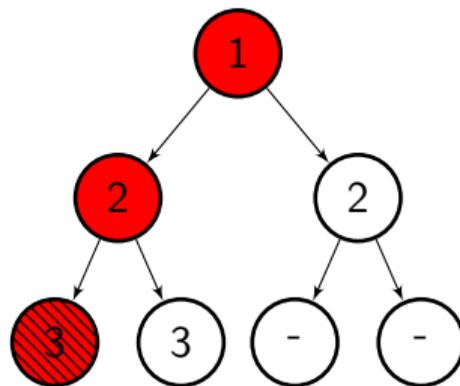
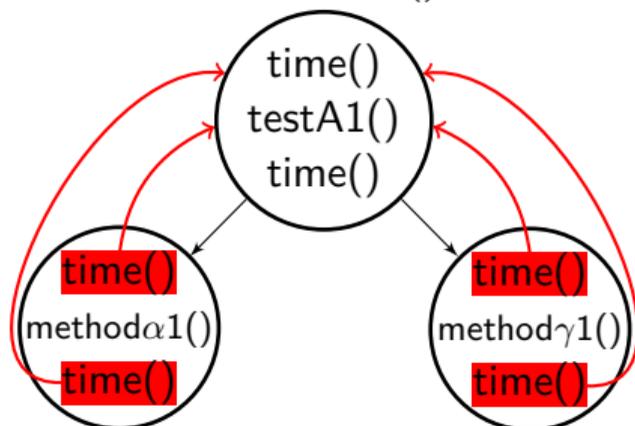
- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



- ▶ Challenges for unit test measurement
 - ▶ Overhead influences measurement quality and duration
⇒ Low overhead necessary
 - ▶ Configuration and node selection strategy a priori unknown

Root Cause Analysis

- ▶ Root cause analysis for performance changes (Heger et al., 2013)
 - ▶ Getting the call tree
 - ▶ Levelwise measurement of call tree nodes
 - ▶ `time()`-calls if `method α 1()` and `method γ 1()` influence measurement of `testA1()`



- ▶ Challenges for unit test measurement
 - ▶ **Overhead influences measurement quality and duration**
⇒ **Low overhead necessary**
 - ▶ Configuration and node selection strategy a priori unknown

Kieker Monitoring Process [WOSPC23]

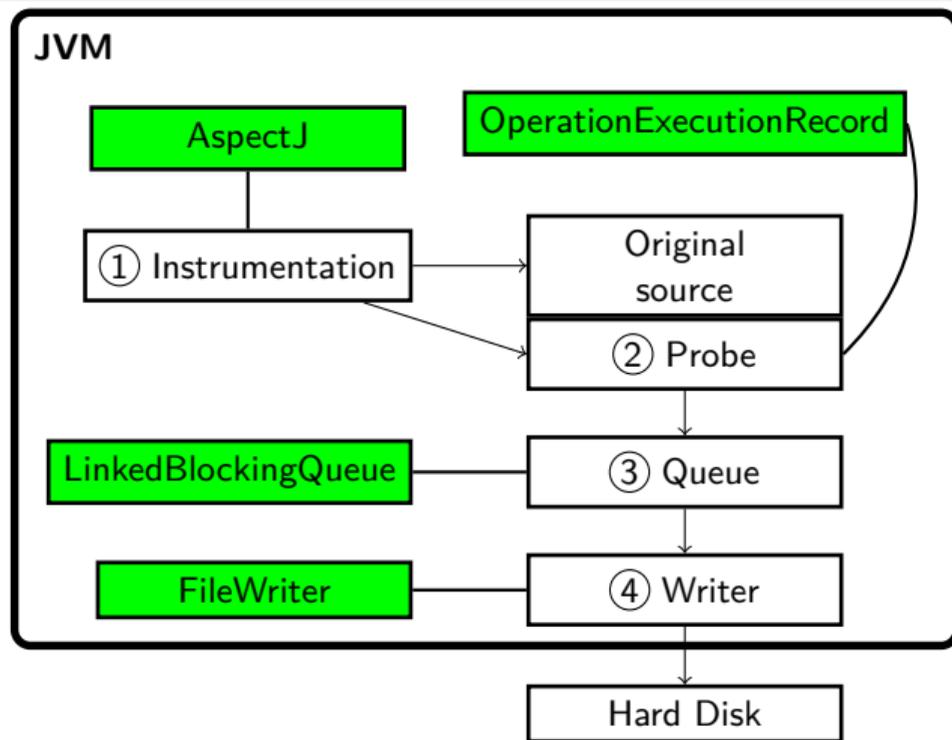


Figure: Monitoring Process in Kieker

Kieker Monitoring Process [WOSPC23]

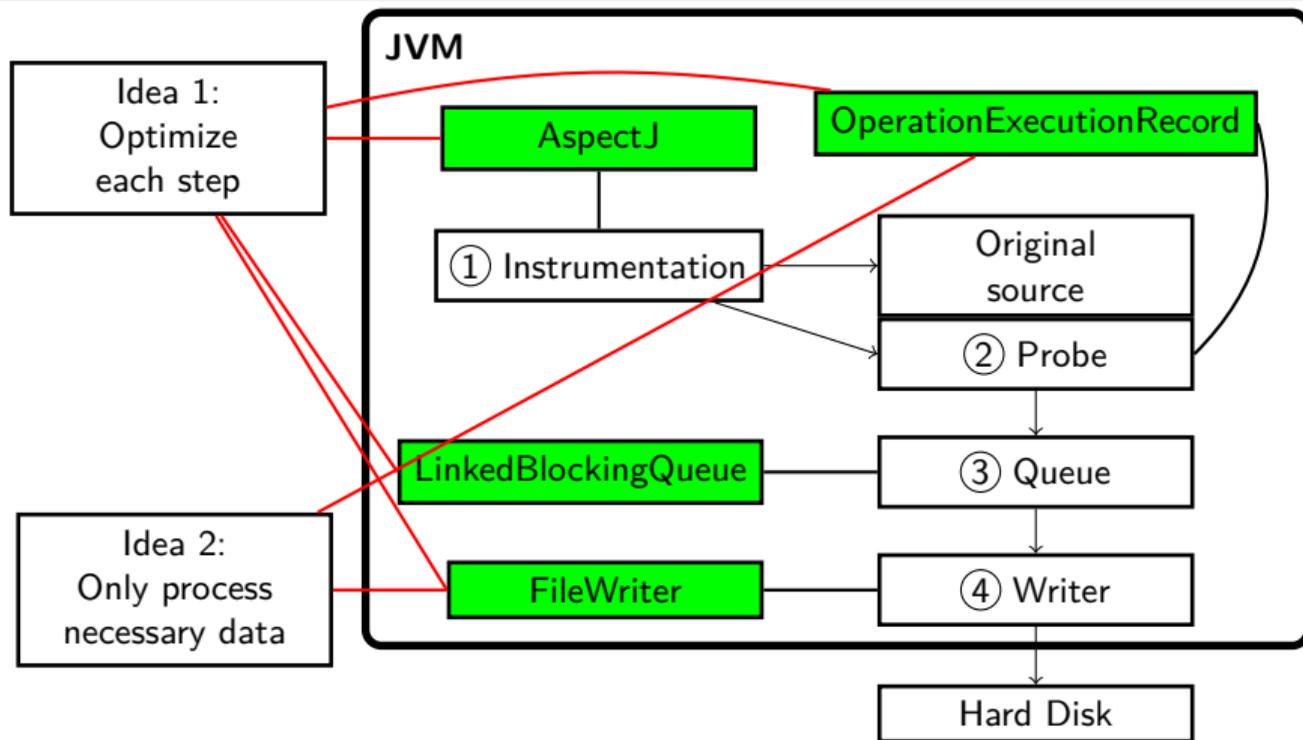


Figure: Monitoring Process in Kieker

Optimization Options [WOSPC23]

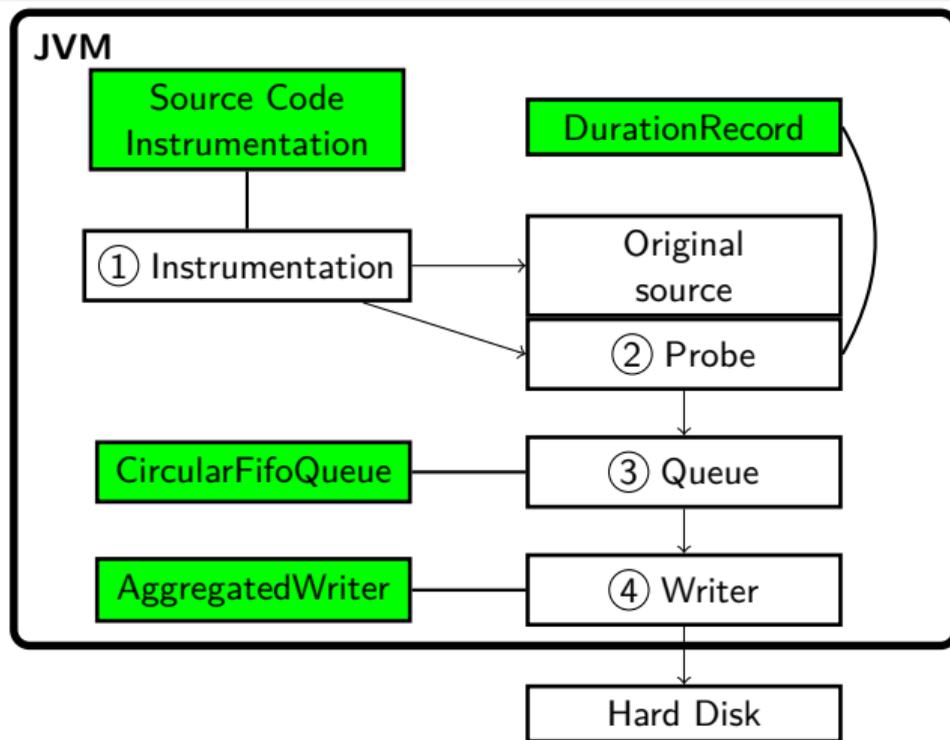


Figure: Possible Monitoring Optimizations

Optimization Options [WOSPC23]

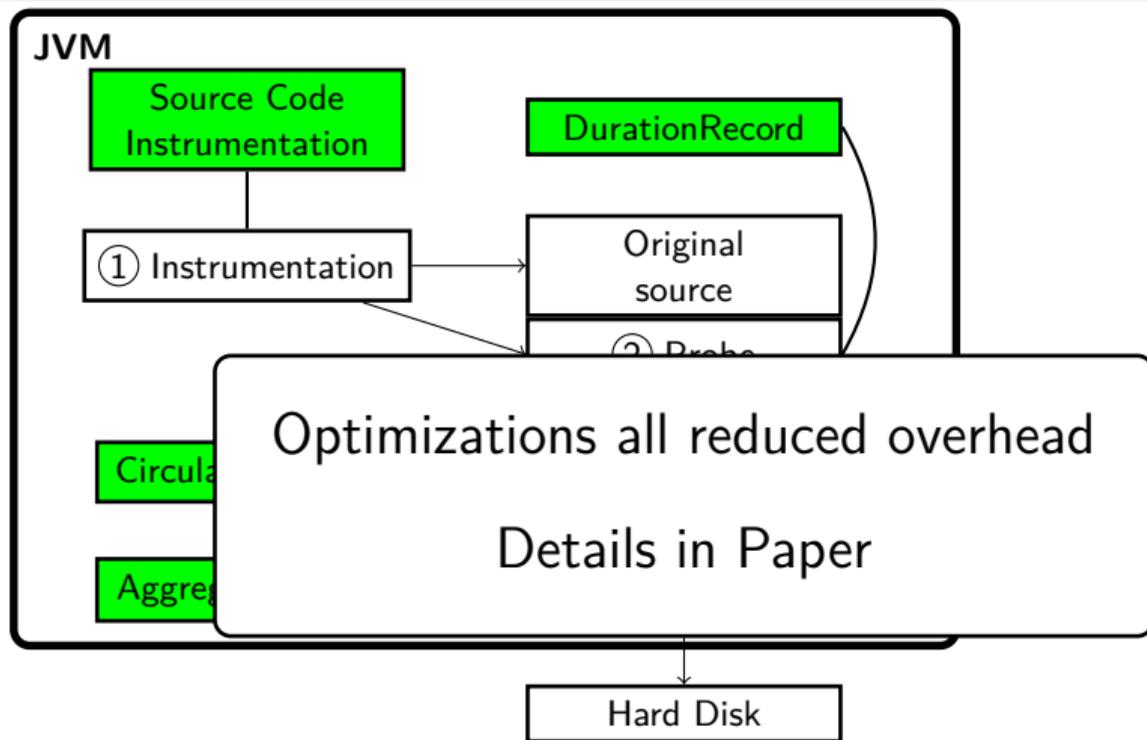
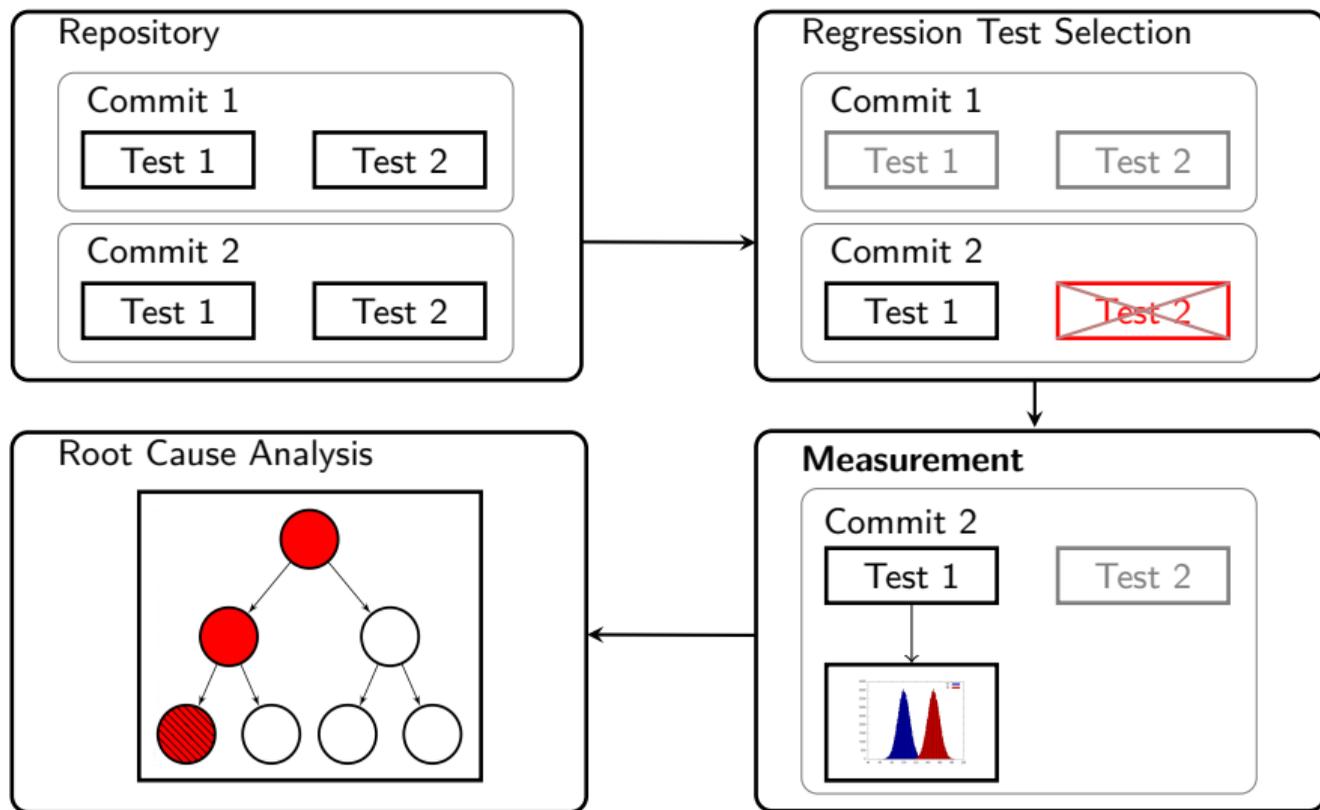


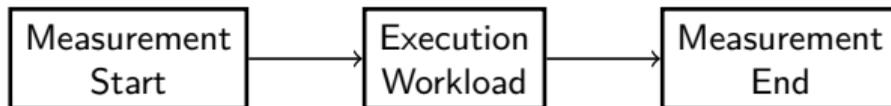
Figure: Possible Monitoring Optimizations

Performance Change Detection

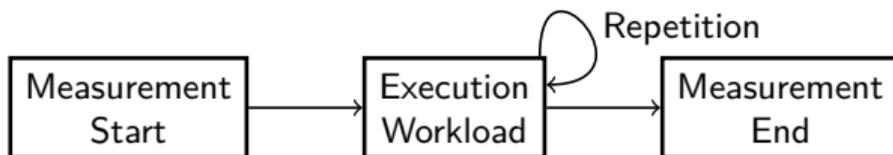
Process of Peass



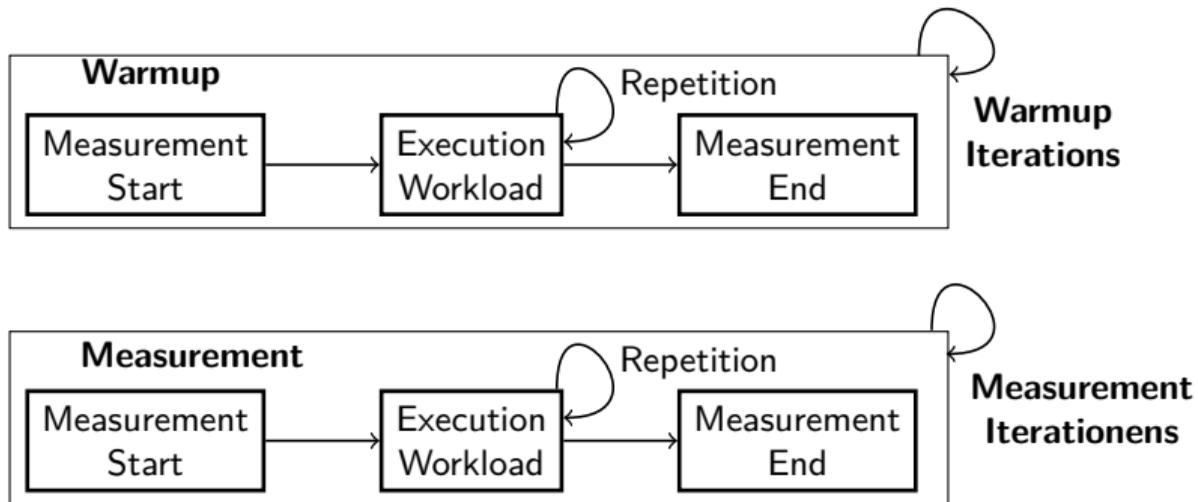
General Measurement Process (Adaption of Georges et al., 2007)



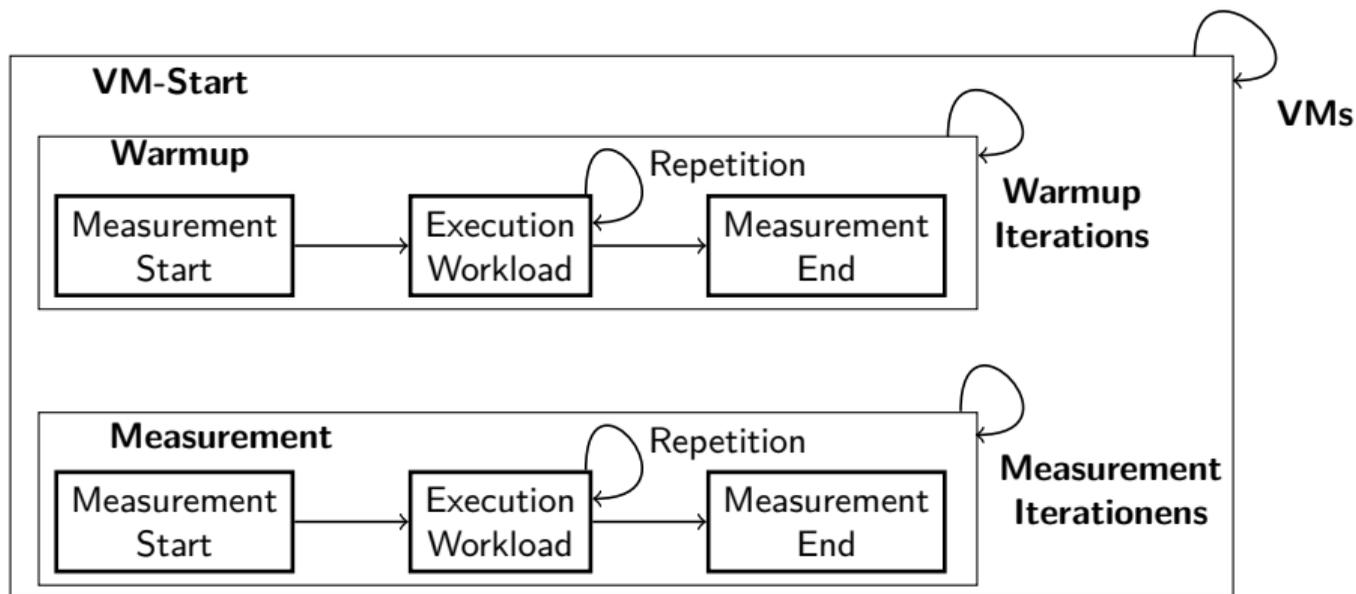
General Measurement Process (Adaption of Georges et al., 2007)



General Measurement Process (Adaption of Georges et al., 2007)



General Measurement Process (Adaption of Georges et al., 2007)



Example Workloads [QRS22]

- ▶ Goal: Determine statistical boundaries of performance change detection
- ▶ Statistical boundaries depend on
 - ▶ Type I error
 - ▶ Type II error
 - ▶ Standard deviation
 - ▶ Effect size (change size in relation to standard deviation)
 - ▶ VM count

Example Workloads [QRS22]

- ▶ Goal: Determine statistical boundaries of performance change detection
- ▶ Statistical boundaries depend on
 - ▶ Type I error
 - ▶ Type II error
 - ▶ Standard deviation
 - ▶ Effect size (change size in relation to standard deviation)
 - ▶ VM count
- ▶ For examination of performance measurements: Usage of example workloads with workload size s
- ▶ Used workloads
 - ▶ Add-workload: Adding s random numbers
 - ▶ RAM-workload: Reservation of s arrays consisting of 3 ints
 - ▶ Sysout-workload: Generating s random numbers and printing them to System.out

Example Workloads [QRS22]

- ▶ Goal: Determine statistical boundaries of performance change detection
- ▶ Statistical boundaries depend on
 - ▶ Type I error Decided: 1 %
 - ▶ Type II error
 - ▶ Standard deviation Empirical: 1 %
 - ▶ Effect size (change size in relation to standard deviation)
 - ▶ VM count
- ▶ For examination of performance measurements: Usage of example workloads with workload size s
- ▶ Used workloads
 - ▶ Add-workload: Adding s random numbers
 - ▶ RAM-workload: Reservation of s arrays consisting of 3 ints
 - ▶ Sysout-workload: Generating s random numbers and printing them to System.out

Statistical Boundaries of Detecting Changes [QRS22]

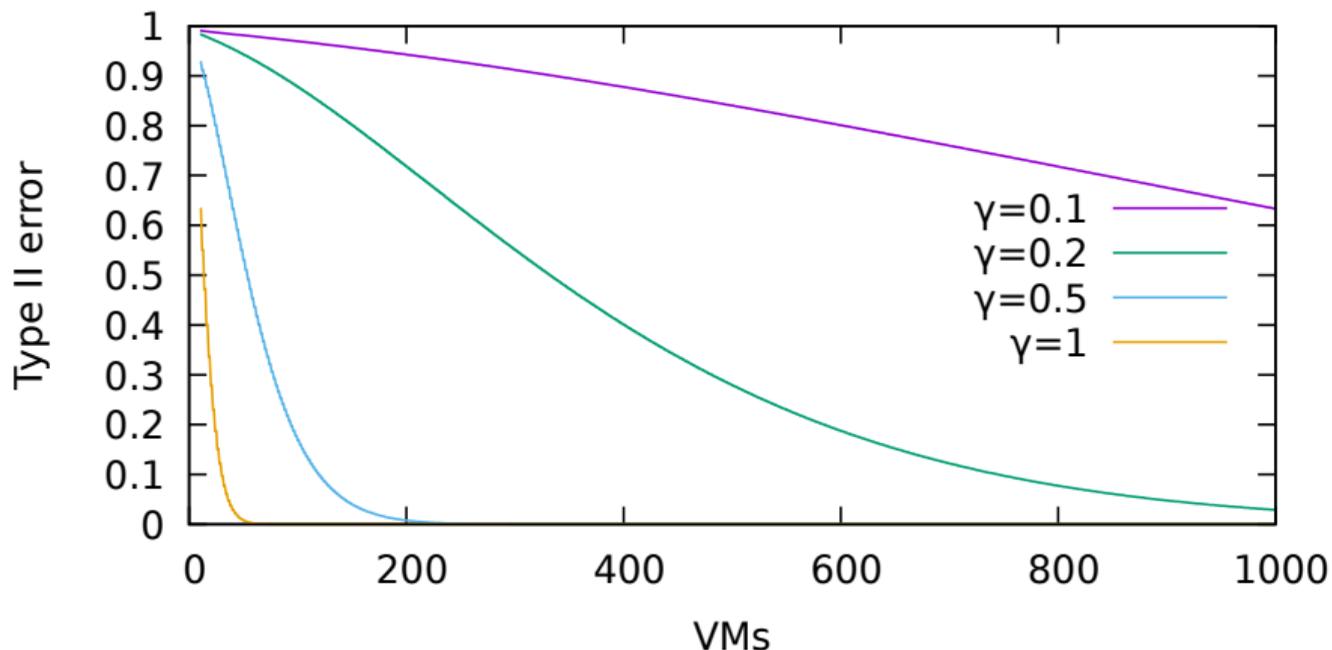


Figure: Type II error in relation to VM count and effect size γ

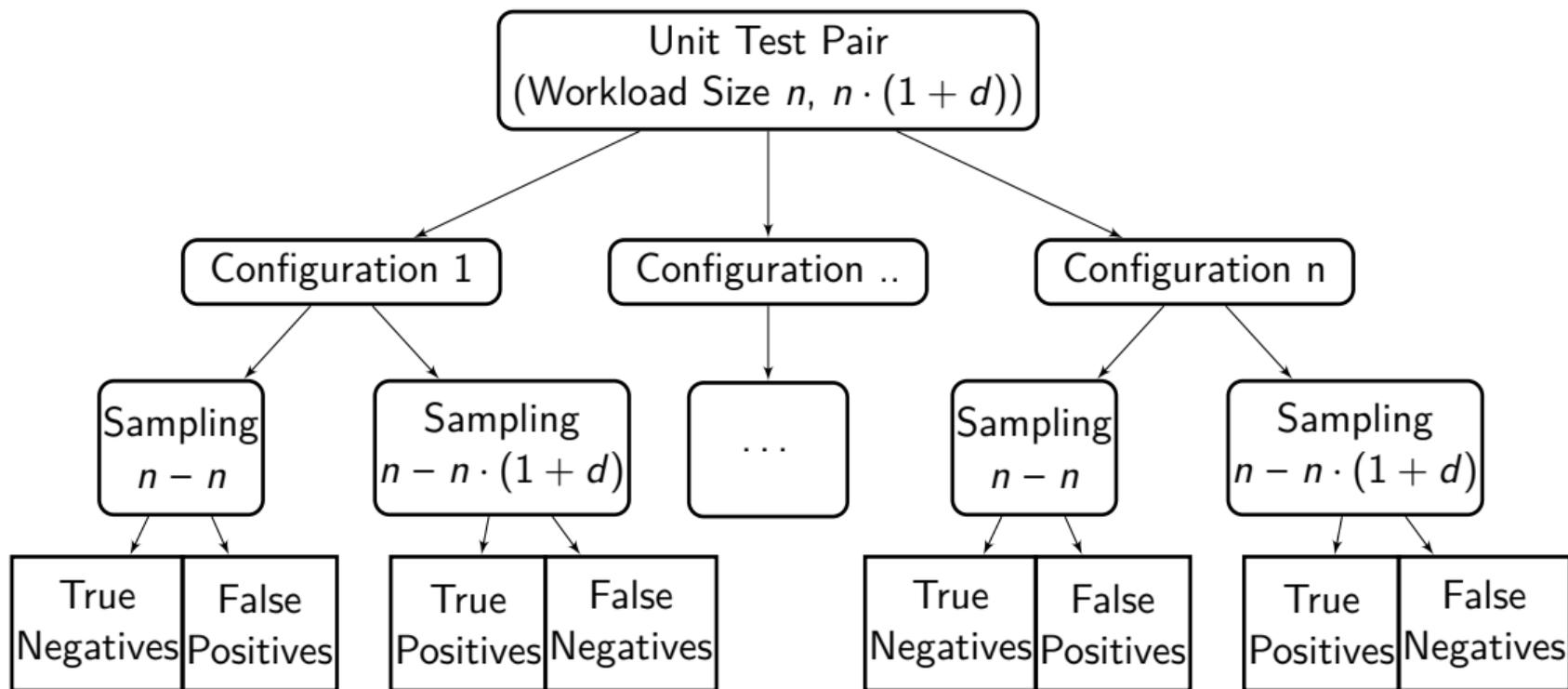
Statistical Boundaries of Detecting Changes [QRS22]

- ▶ Example: Workload size 1,000 will take 97 seconds for 10,000,000 repetitions
⇒ at maximum 891 VMs are possible when using 12 hours measurement time
- ▶ Boundaries depend on measurement duration; $\gamma = 0.1$ would require 4.808 VMs and will hardly be measurable

Statistical Boundaries of Detecting Changes [QRS22]

- ▶ Example: Workload size 1,000 will take 97 seconds for 10,000,000 repetitions
⇒ at maximum 891 VMs are possible when using 12 hours measurement time
- ▶ Boundaries depend on measurement duration; $\gamma = 0.1$ would require 4.808 VMs and will hardly be measurable
- ▶ Measurements via monitoring (e.g. with OpenTelemetry, Kieker) or load tests have higher deviation and will only detect more coarse-grained performance changes

Approach [QRS22]



Workload Size [QRS22]

- ▶ Examination of unit test sized workloads
 - ⇒ What is a unit test size?
(What is s ?)

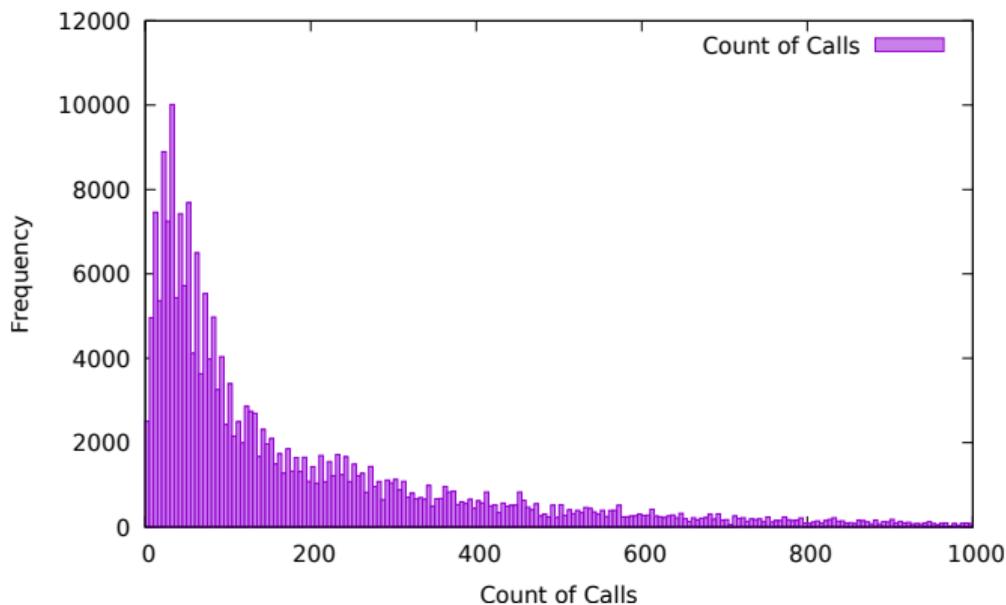
Workload Size [QRS22]

- ▶ Examination of unit test sized workloads
 - ⇒ What is a unit test size?
(What is s ?)
- ▶ Approach: Analysis of call counts of open source projects

Workload Size [QRS22]

- ▶ Examination of unit test sized workloads
⇒ What is a unit test size?
(What is s ?)

- ▶ Approach: Analysis of call counts of open source projects



Example: Comparison of Statistical Tests [QRS22]

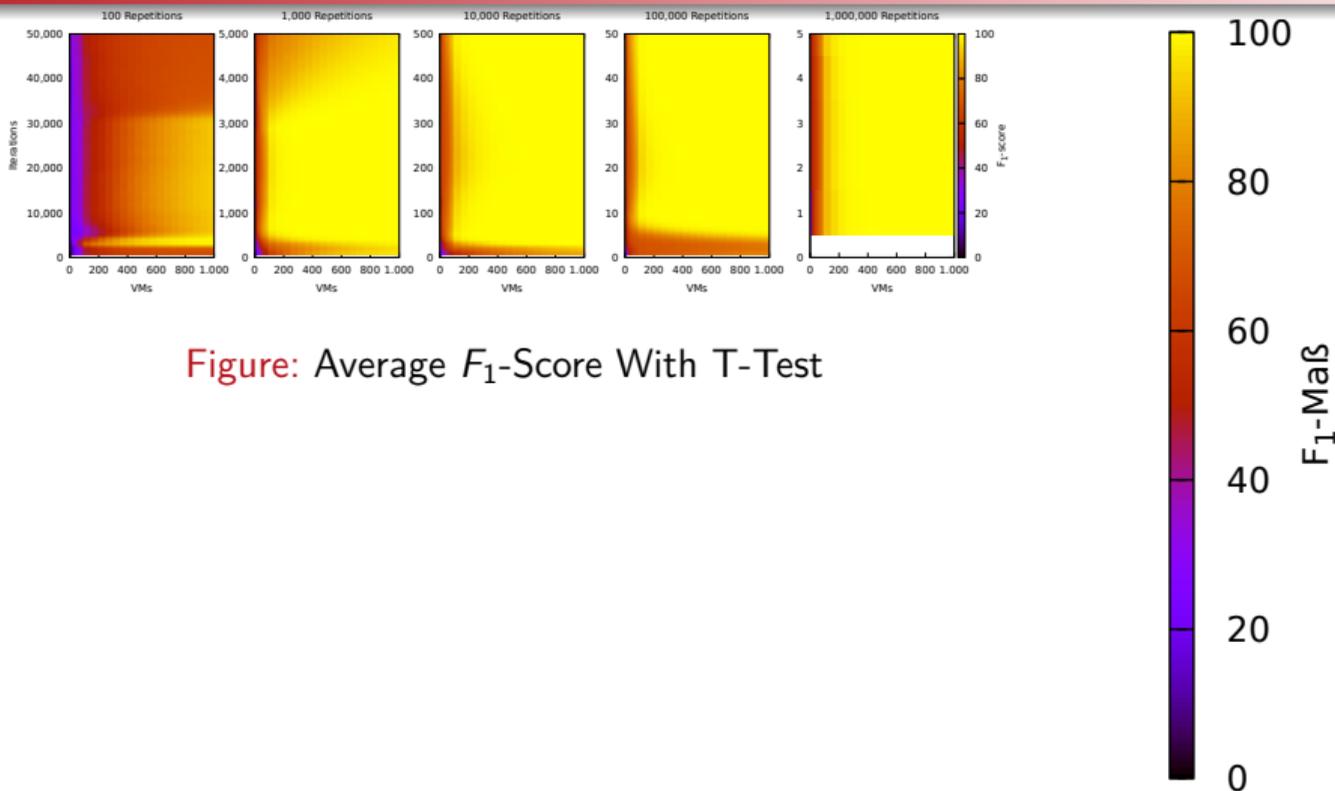


Figure: Average F_1 -Score With T-Test

Example: Comparison of Statistical Tests [QRS22]

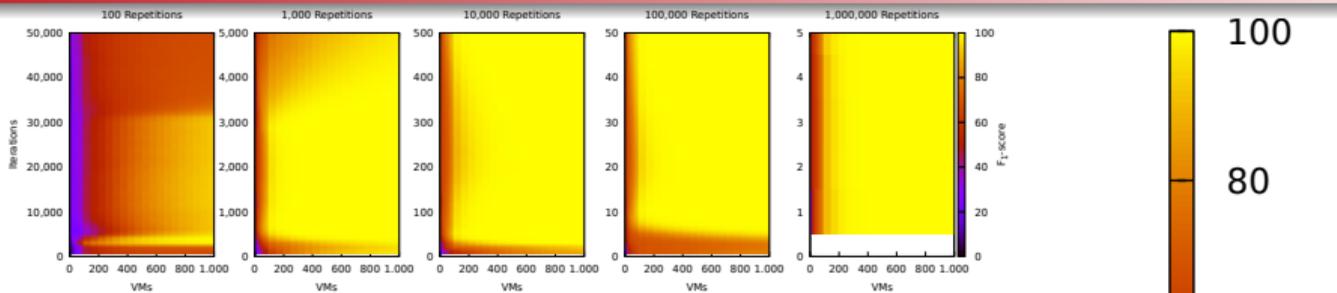


Figure: Average F_1 -Score With T-Test

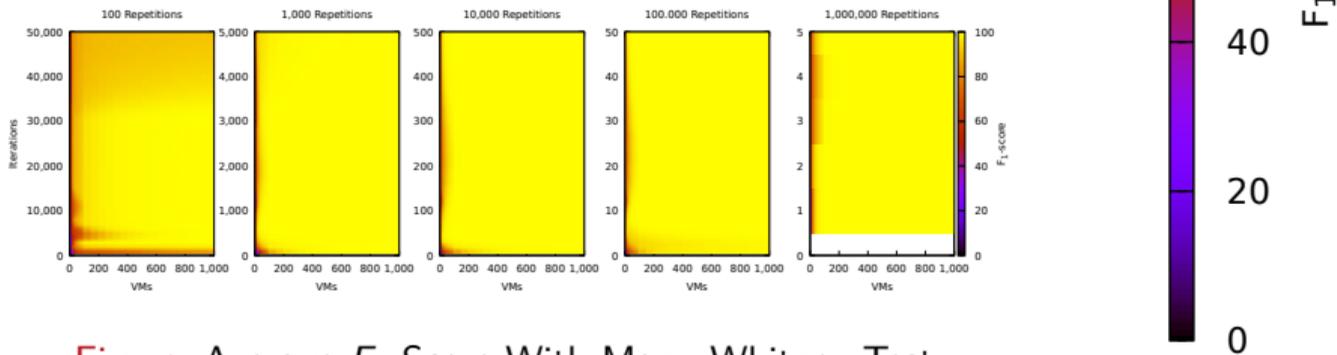
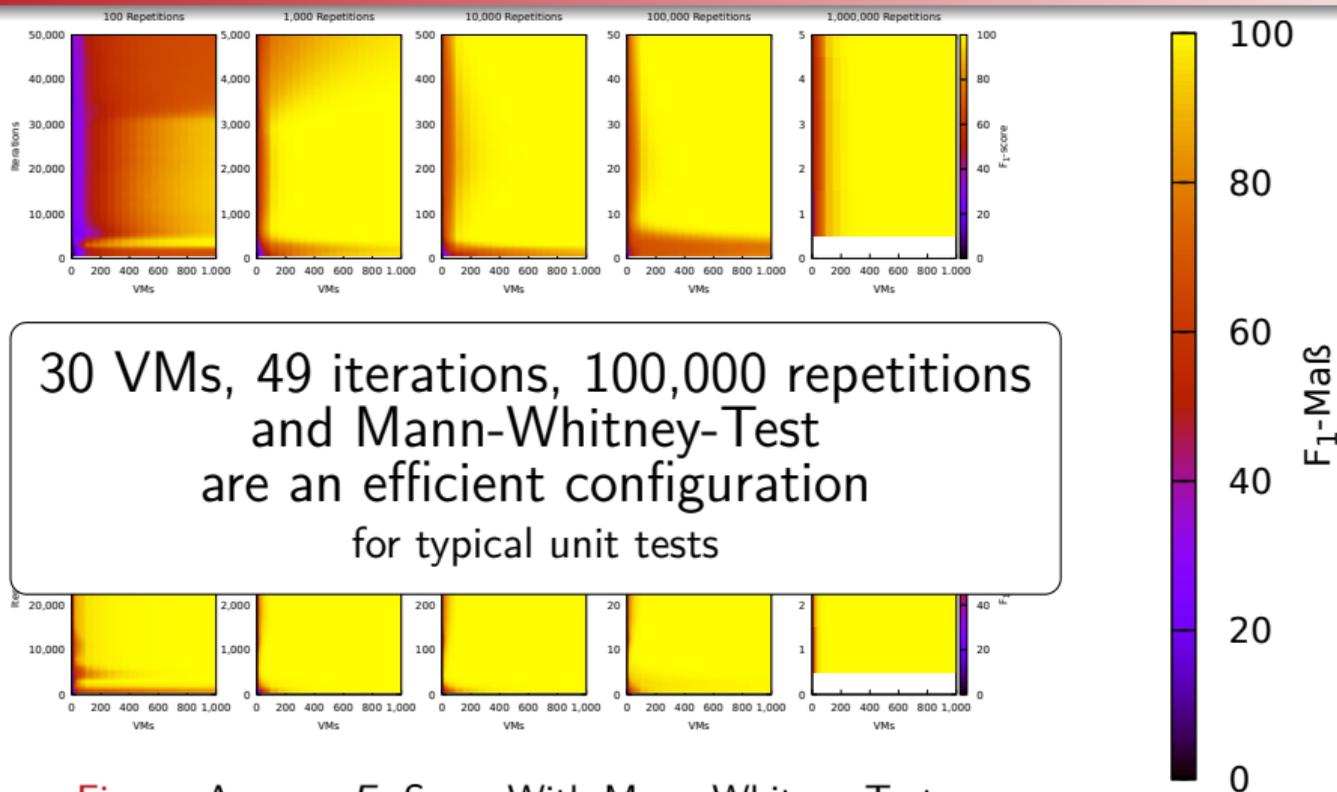


Figure: Average F_1 -Score With Mann-Whitney-Test

Example: Comparison of Statistical Tests [QRS22]



Example Usage: Peass-CI

The screenshot shows the Jenkins web interface for a build named 'ServletFileUploadTest_testFoldedHeaders'. The left sidebar contains navigation options like 'Back to Project', 'Status', 'Changes', 'Console Output', etc. The main area displays performance metrics:

| Property | Predecessor | Current |
|------------------|--------------|--------------|
| Mean | 66890.885 µs | 77788.354 µs |
| Deviation | 7046.532 | 5002.706 |
| In-VM-Deviation | 0 | 0 |
| VMs: 100 T=12.61 | | |

Additional information includes the version '4ed6e923cb2033272fcb993978d69e325910a5aa' and the test case 'org.apache.commons.fileupload.ServletFileUploadTest#testFoldedHeaders'. A search bar at the top right contains the text 'Suchen'.

Overlaid text box:

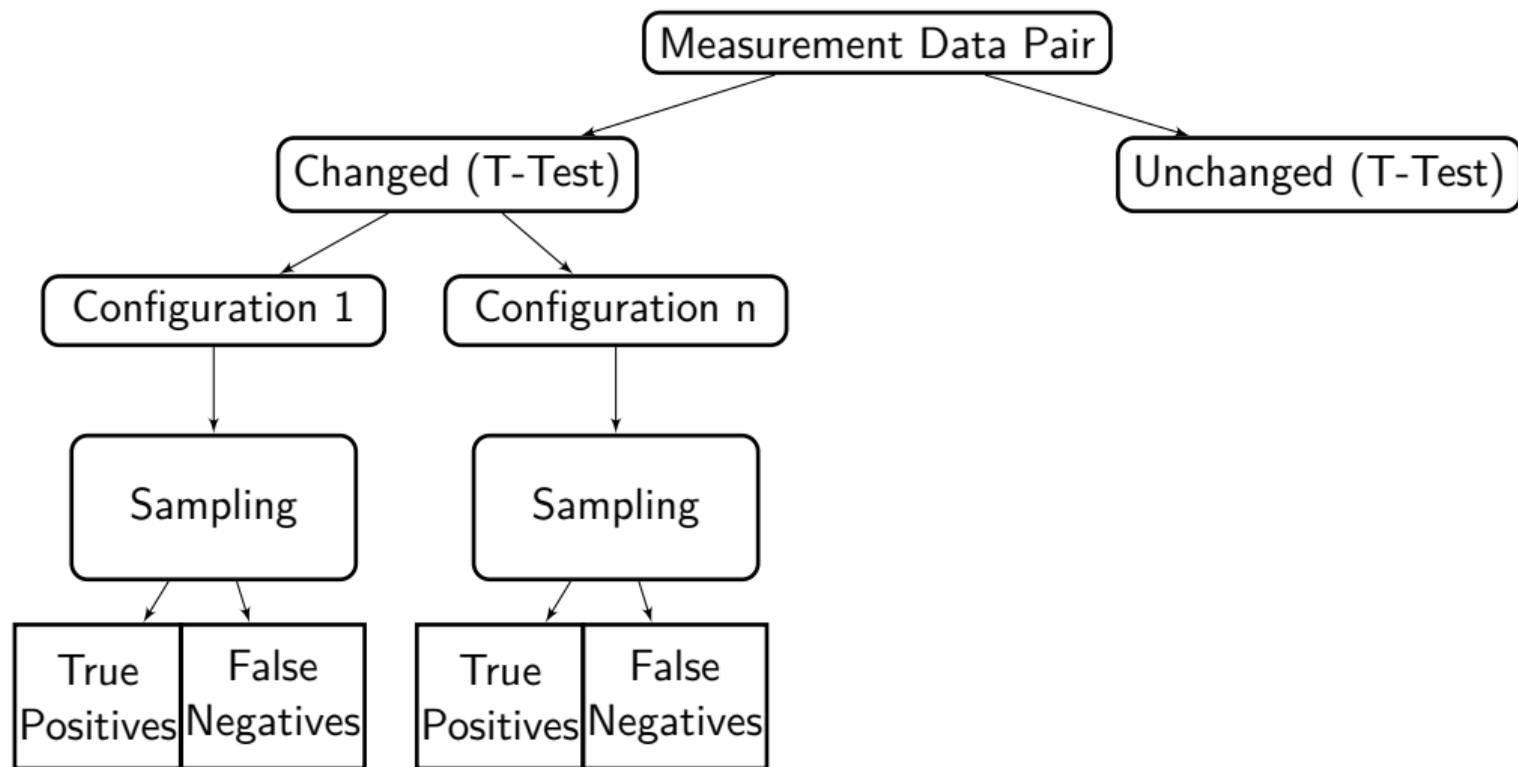
Detection of performance changes in CI:
<https://plugins.jenkins.io/peass-ci/>
 Steps of the approach
 are all implemented

Application to GraalVM Data

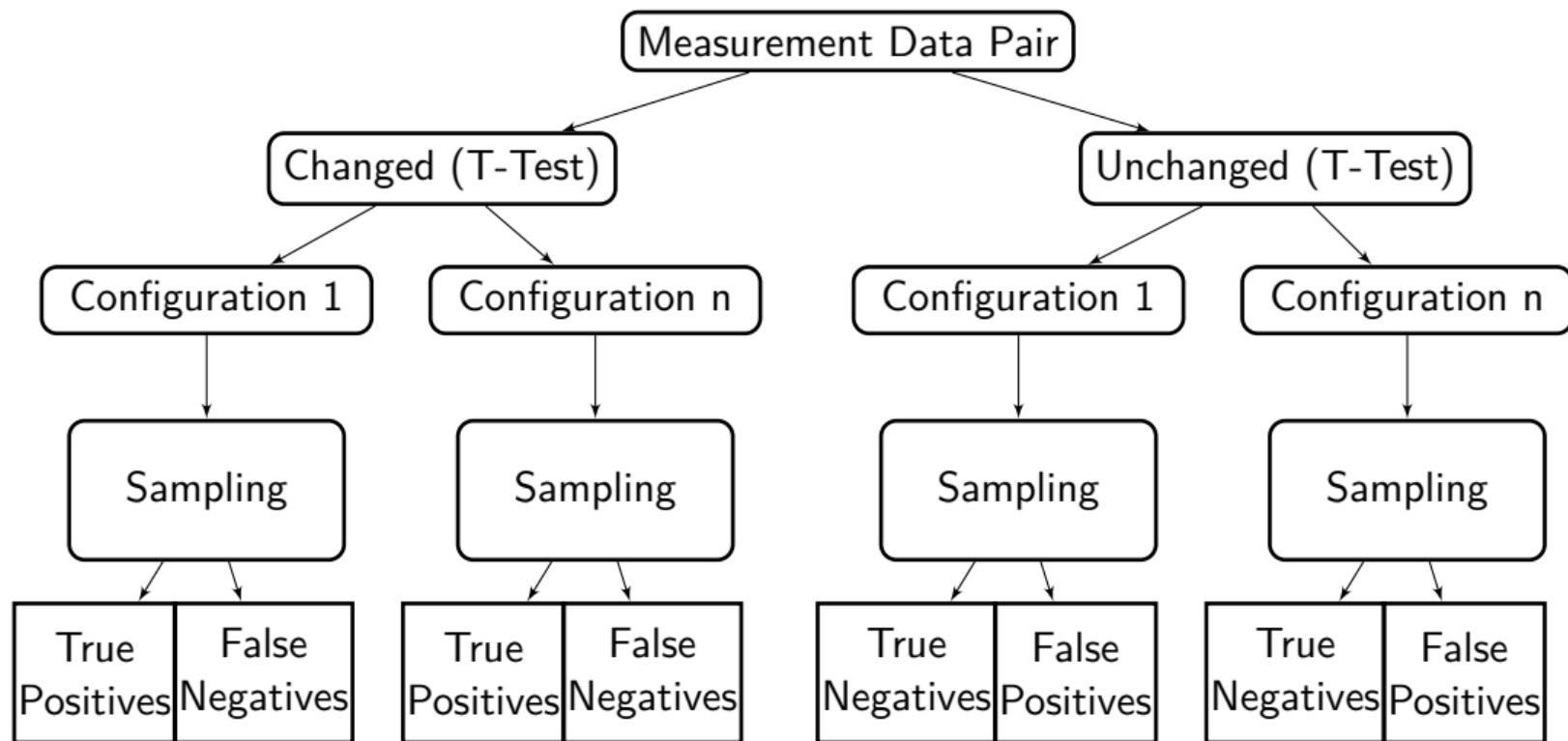
Application to GraalVM Data – Idea

- ▶ Problem: GraalVM benchmarks might be executed too often
- ▶ Goal: Determination of suitable VM count and iteration / run count
- ▶ Problem: Artificial benchmark pairs and option to repeat them more often than necessary not available ⇒ Approach: Sample from existing data

Application to GraalVM Data – Approach



Application to GraalVM Data – Approach



Application to GraalVM Data – Preliminary Results

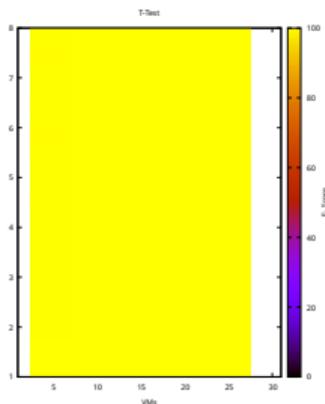


Figure: Difference
48994

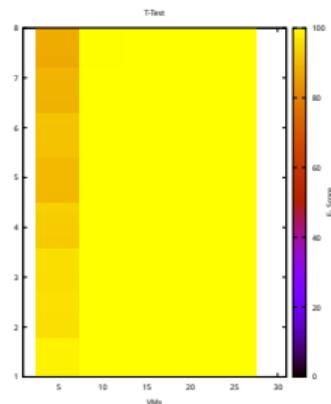


Figure: Difference
49001

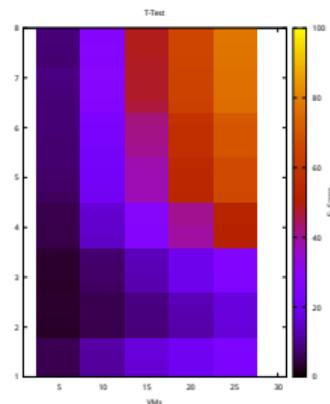


Figure: Difference
50180

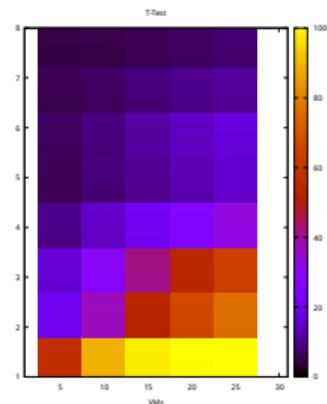


Figure: Difference
50405

Summary and Outlook

Summary and Outlook

- ▶ Previous research
 - ▶ Definition of a method for performance change detection using unit tests:
<https://github.com/DaGeRe/peass>
 - ▶ Regression test selection by trace analysis
 - ▶ Root cause analysis implementation using Kieker
- ▶ Performance change detection
 - ▶ Basic idea: Sample from artificial / known performance changes and derive measurement configuration
 - ▶ Current goal: Check configuration generation on GraalVM data

Thanks for your attention!

- ▶ David Georg Reichelt
- ▶ Assistant Professor
- ▶ Lancaster University Leipzig
- ▶ Contact: [d.g.reichelt \[a t \] lancaster.ac.uk](mailto:d.g.reichelt@lancaster.ac.uk)