



BENCHMARKING
THE SECURITY OF SOFTWARE SYSTEMS OR
TO BENCHMARK OR NOT TO BENCHMARK

LASIGE Workshop'2018
 Lisbon, Portugal
 June 07th, 2018

Marco Vieira
mvieira@dei.uc.pt
 Department of Informatics Engineering
 University of Coimbra - Portugal



1



BENCHMARKING

**Assessing and comparing
 computer systems and/or components
 according to specific quality attributes**

- Performance benchmarking
 - Well established both in terms of research and application
 - Supported by organizations like TPC and SPEC
 - Mostly for marketing
- Dependability benchmarking
 - Well established from a research perspective
 - No endorsement from the industry

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 2

2

BENCHMARKING

**Assessing and comparing
computer systems and/or components
according to specific quality attributes**

- Security benchmarking
 - Several works can be found
 - No common approach available yet

Performance benchmarks *Dependability benchmarks* *Security benchmarks*

1972 1983 1985 1988 1987 1999 2000 2017

Whetstone Wisconsin Bench TP1 DebitCredit Orange Book TPC & SPEC EMBC SIGDeB CIS Common Criteria

Release of commercial performance benchmarks ... Research projects on dependability & security benchmarks

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 3

3


OUTLINE

The past: Performance & Dependability Benchmarking

- The present: Security Benchmarking
- Benchmarking the **Security of Systems**
 - Approach: Qualification + Trustworthiness Assessment
 - Example: Benchmarking Web Service Frameworks
- Benchmarking **Security Tools**
 - Approach: Vulnerability and Attack Injection
 - Example: Benchmarking Intrusion Detection Systems
- Challenges and Conclusions

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 4

4




PERFORMANCE BENCHMARKING


Assessing and comparing
computer systems and/or components
in terms of performance

Marco Vieira *LASIGE'2018*, Lisbon, Portugal, June 07th, 2018 5

5



PERFORMANCE BENCHMARKING




```
graph LR; W[Workload] --> SUB[SUB]; SUB --> M[Metrics]
```


- Workload:
 - Set of representative operations
- Metrics:
 - Throughput
 - Response time
 - Latency
 - ...

Marco Vieira *LASIGE'2018*, Lisbon, Portugal, June 07th, 2018 6

6



TPC-C (1992)




```

graph LR
    W[Workload] --> DBMS[DBMS]
    DBMS --> M[Metrics]
  
```

- Workload:
 - Database transactions
- *Although some integrity tests are performed, it assumes that nothing fails*
 - Transaction rate (tpmC)
 - Price per transaction (\$/tpmC)

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 7

7



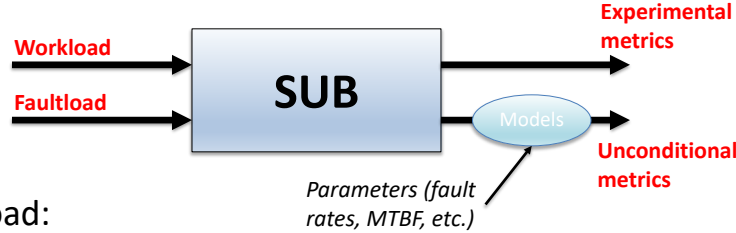
DEPENDABILITY BENCHMARKING

**Assessing and comparing
computer systems and/or components
considering dependability attributes**

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 8

8

DEPENDABILITY BENCHMARKING



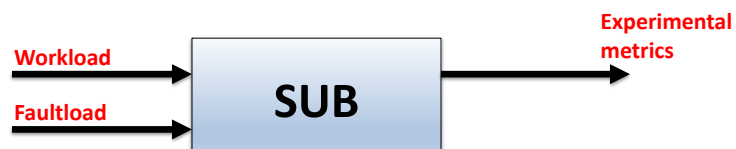
The diagram shows a central box labeled 'SUB'. Two arrows labeled 'Workload' and 'Faultload' point into the box from the left. Two arrows point out of the box to the right: one labeled 'Experimental metrics' and another labeled 'Unconditional metrics'. A blue oval labeled 'Models' is connected to the 'SUB' box by a double-headed arrow. An arrow points from the text 'Parameters (fault rates, MTBF, etc.)' to the 'Models' oval.

- **Faultload:**
 - Set of representative faults, injected into the system
- **Metrics:**
 - Performance and/or dependability
 - Both baseline and in the presence of faults
 - Unconditional and/or direct

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 9

9

DBENCH-OLTP (2005)



The diagram shows a central box labeled 'SUB'. Two arrows labeled 'Workload' and 'Faultload' point into the box from the left. One arrow labeled 'Experimental metrics' points out of the box to the right.

- **Workload:**
 - TPC-C transactions
- **Faultload:**
 - Operator faults + Software faults + HW component failures
- **Metrics:**
 - Performance: tpmC, \$/tpmC, Tf, \$/Tf
 - Dependability: Ne, AvtS, AvtC

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 10

10

DBENCH-OLTP (2005)

| System | Operating System | DBMS | DBMS Config. | Hardware |
|--------|-----------------------|----------------------|--------------|---|
| A | Windows 2K Prof. SP 3 | Oracle 8i R2 (8.1.7) | Config. A | <i>Processor:</i> Intel Pentium III 800 MHz <i>Memory:</i> 256MB <i>Hard Disks:</i> Four 20GB/7200 rpm <i>Network:</i> Fast Ethernet |
| B | Windows 2K Prof. SP 3 | Oracle 9i R2 (9.0.2) | Config. A | |
| C | Windows Xp Prof. SP 1 | Oracle 8i R2 (8.1.7) | Config. A | |
| D | Windows Xp Prof. SP 1 | Oracle 9i R2 (9.0.2) | Config. A | |
| E | Windows 2K Prof. SP 3 | Oracle 8i R2 (8.1.7) | Config. B | |
| F | Windows 2K Prof. SP 3 | Oracle 9i R2 (9.0.2) | Config. B | |
| G | SuSE Linux 7.3 | Oracle 8i R2 (8.1.7) | Config. A | |
| H | SuSE Linux 7.3 | Oracle 9i R2 (9.0.2) | Config. A | |
| I | SuSE Linux 7.3 | PostgreSQL 7.3 | - | |
| J | Windows 2K Prof. SP 3 | Oracle 8i R2 (8.1.7) | Config. A | <i>Processor:</i> Intel Pentium IV 2GHz <i>Memory:</i> 512MB <i>Hard Disks:</i> Four 20GB/7200 rpm <i>Network:</i> Fast Ethernet |
| K | Windows 2K Prof. SP 3 | Oracle 9i R2 (9.0.2) | Config. A | |

Faultload: Operator faults

Marco Vieira
LASIGE'2018, Lisbon, Portugal, June 07th, 2018
11

11

DBENCH-OLTP (2005)

Baseline Performance

| System | tpmC | \$/tpmC |
|--------|------|---------|
| A | 2200 | 15 |
| B | 2500 | 15 |
| C | 2300 | 15 |
| D | 2500 | 15 |
| E | 1500 | 15 |
| F | 1500 | 15 |
| G | 1800 | 15 |
| H | 1800 | 15 |
| I | 800 | 15 |
| J | 3500 | 15 |
| K | 3500 | 15 |


Performance With Faults

| System | Tf | \$/Tf |
|--------|------|-------|
| A | 1500 | 15 |
| B | 1800 | 15 |
| C | 1800 | 15 |
| D | 1800 | 15 |
| E | 1000 | 15 |
| F | 1000 | 15 |
| G | 1500 | 15 |
| H | 1500 | 15 |
| I | 1000 | 15 |
| J | 3000 | 15 |
| K | 3000 | 15 |

Does not take into account malicious behaviors (faults = vulnerability + attack)

Marco Vieira
LASIGE'2018, Lisbon, Portugal, June 07th, 2018
12

12




SECURITY BENCHMARKING

**Assessing and comparing
computer systems and/or components
considering security aspects**


- Benchmarking the Security of **Systems / Components**
 - Systems that should implement security requirements
 - OS, middleware, server software, etc.
- Benchmarking **Security Tools**
 - Tools used to improve the security of systems
 - Penetration testers, static analyzers, IDS, etc.

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 13

13



BENCHMARKING SECURITY OF SYSTEMS

Workload →  → Experimental metrics

*Attacking what? Do we know the vulnerabilities?
What are representative attacks?*

- *Does not work if one wants to benchmark how secure different systems are!*
- *e.g. does the number of vulnerabilities of a system represent anything?*
 - Performance + dependability
 - Security (e.g., number vulnerabilities, attack detection)

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 14

14

A DIFFERENT APPROACH...

```

graph TD
    SUBs --> SQ[Security Qualification]
    SQ -- Unacceptable --> S0[Security = 0]
  
```

- **Security Qualification:**
 - Apply state-of-the-art techniques and tools to detect vulnerabilities
 - SUBs with vulnerabilities are:
 - Disqualified!
 - Or vulnerabilities are fixed...

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 15

15

A DIFFERENT APPROACH...

```

graph LR
    SUBs --> SQ[Security Qualification]
    SQ -- Unacceptable --> S0[Security = 0]
    SQ -- Acceptable --> TWA[Trustworthiness Assessment]
    TWA -- Metrics --> MetricsOut[ ]
  
```

- **Trustworthiness Assessment:**
 - Gather evidences on how much one can trust
 - e.g., best coding practices, development process, bad smells

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 16

16

A DIFFERENT APPROACH...

```

graph LR
    SUBs[SUBs] --> SQ[Security Qualification]
    SQ -- Acceptable --> TA[Trustworthiness Assessment]
    SQ -- Unacceptable --> S0[Security = 0]
    TA -- Metrics --> Metrics[Metrics]
  
```

- Metrics:
 - Portray trust from a user perspective
 - Dynamic: may change over time
 - Depend on the type of evidences gathered
 - Different metrics for different attack vectors

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 17

17

EXAMPLE: WEB SERVICE FRAMEWORKS

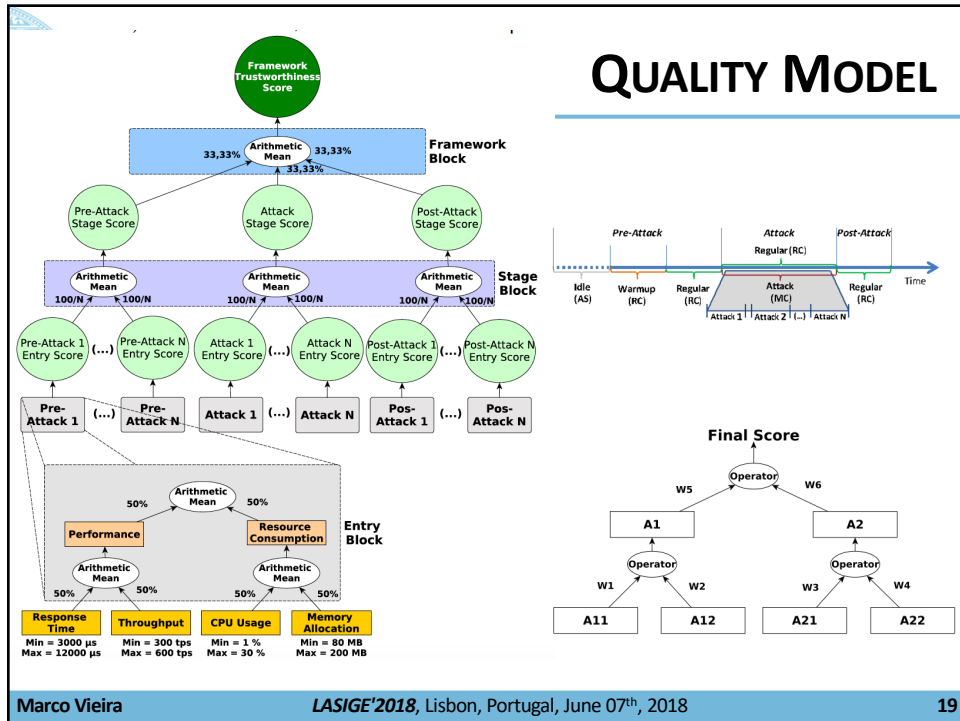
```

graph LR
    WSFs[WSFs] --> Q[Qualification (testing)]
    Q -- Acceptable --> A[Assessment (CPU + mem.)]
    Q -- Unacceptable --> S0[Security = 0]
    A -- Trust. Score --> TS[Trust. Score]
  
```

- Qualification
 - DoS Attacks
 - *Coercive Parsing, Malformed XML, Malicious Attachment, etc.*
- Trustworthiness Assessment:
 - Quality model to compute a score

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 18

18




19

SYSTEMS UNDER BENCHMARKING

| Framework | Version | Security Qualification |
|---------------|---------|------------------------|
| Apache Axis 1 | 1.4.1 | ✗ |
| Apache Axis 2 | 1.6.1 | ✓ |
| | 1.6.2 | ✗ |
| Apache CXF | 2.5.1 | ✓ |
| | 3.0.3 | ✓ |
| Oracle Metro | 2.1.1 | ✗ |
| | 2.3.1 | ✓ |
| XINS | 3.1 | ✗ |
| Spring JAX-WS | 1.9 | ✗ |
| Spring WS | 2.2.0 | ✗ |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 20

20

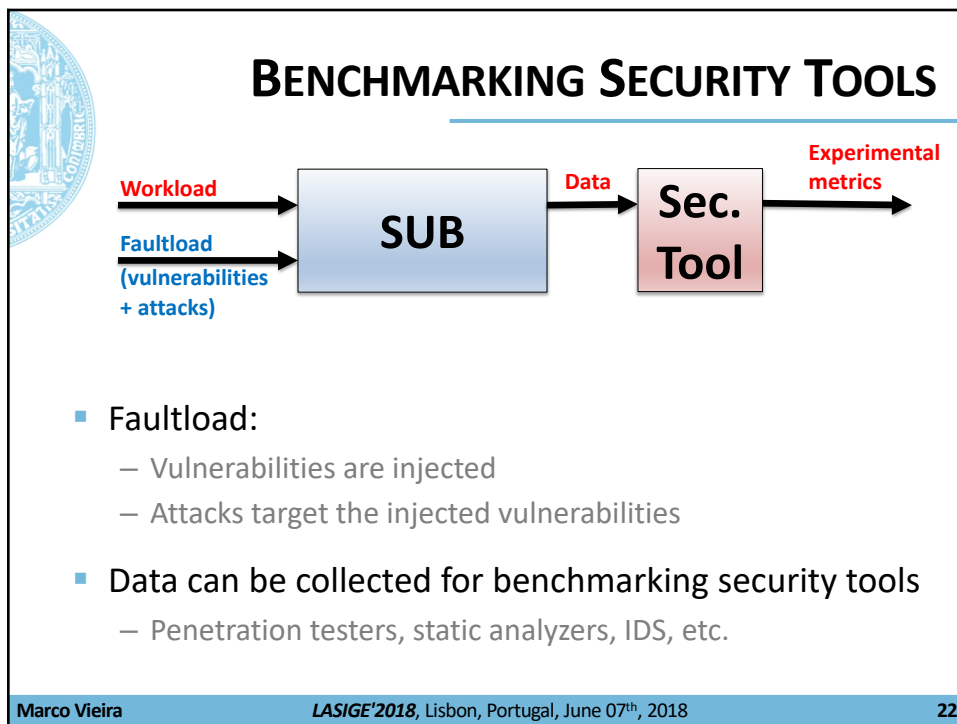


TRUSTWORTHINESS RESULTS

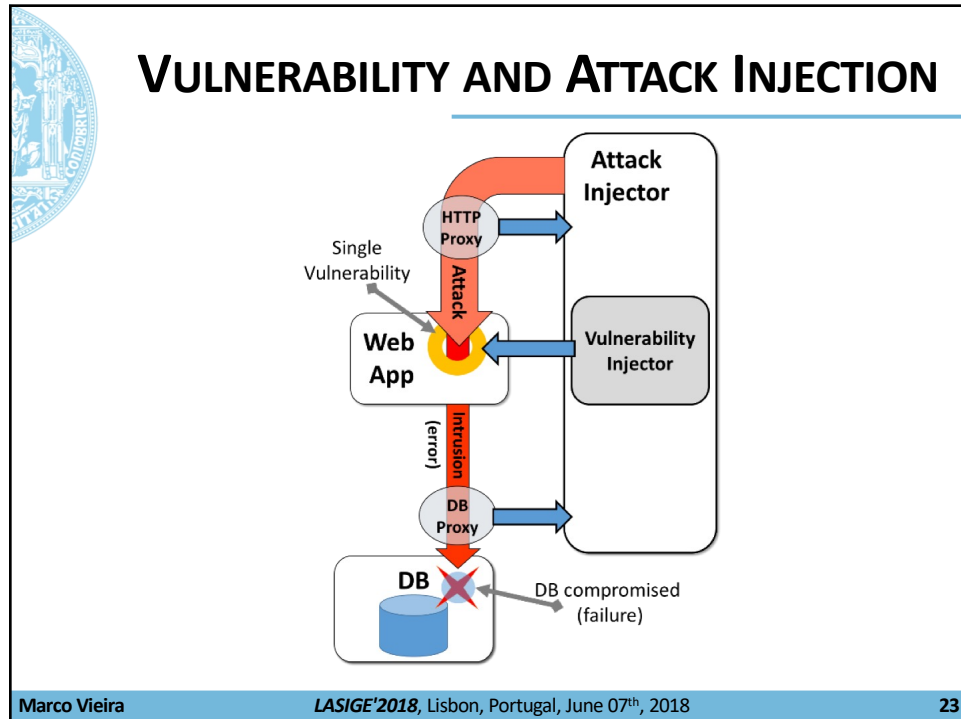
| Scenario | Axis 2 | CXF v2 | Metro | CXF v3 |
|------------------|----------|----------|----------|----------|
| <i>Neutral</i> | 72.3 (1) | 70.7 (2) | 58.1 (3) | 57.9 (4) |
| <i>Scenario1</i> | 73.4 (2) | 77.1 (1) | 66.5 (4) | 70.0 (3) |
| <i>Scenario2</i> | 67.4 (3) | 73.1 (1) | 66.6 (4) | 68.7 (2) |
| <i>Scenario3</i> | 61.8 (4) | 70.3 (1) | 63.6 (3) | 67.0 (2) |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 21

21



22



23

EXAMPLE: BENCHMARKING IDS

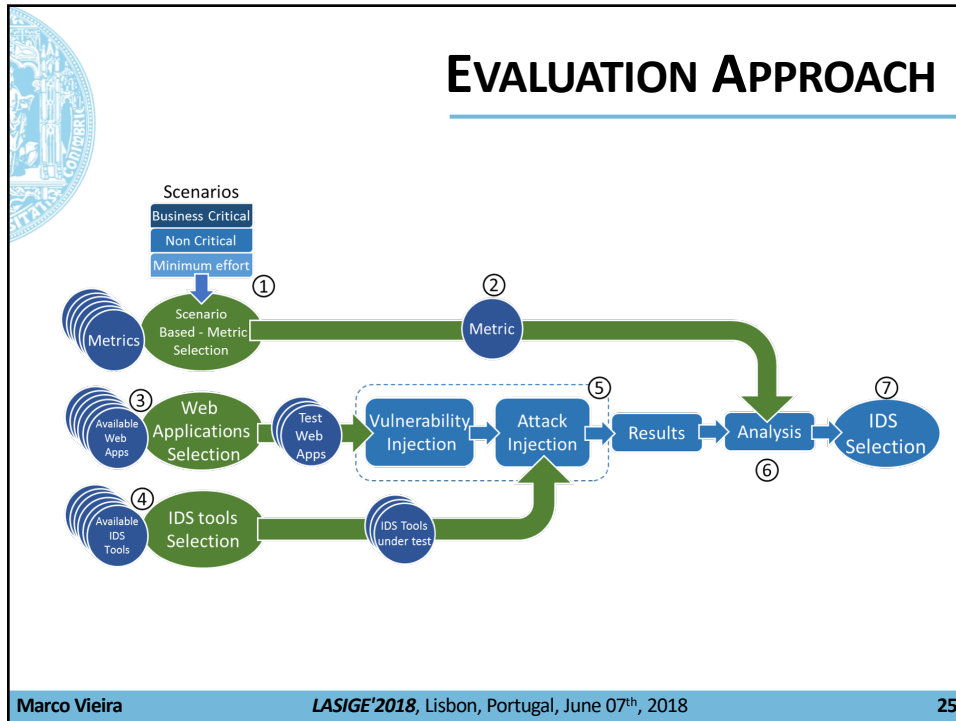
Security requires a defense in depth approach

- Coding best practices
- Testing
- Static analysis
- ...

- Vulnerability-free code is hard (or even impossible) to achieve...
- Intrusion detection tools support a post-deployment approach
 - For protecting against known and unknown attacks

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 24

24




25

EXAMPLES OF VULNERABILITIES INJECTED

| Original PHP code | Code with injected vulnerability | Operation performed |
|--|----------------------------------|---|
| <code>\$id=intval(\$_GET['id']);</code> | <code>\$id=\$_GET['id'];</code> | Removed the “intval” function allowing also non numeric values (i.e. SQL commands) in the “\$id” variable |
| <code>\$page = urlencode(\$page);</code> | <code>\$page = \$page;</code> | Removed the “urlencode” function allowing also alphanumeric values (i.e. SQL commands) in the “\$page” variable |
| ... | ... | ... |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 26

26




EXAMPLES OF ATTACKS

| Attack payloads | Expected result |
|----------------------------------|--|
| ' | Modifies the structure of the query; usually results in an error |
| or 1=1 | Modifies the structure of the query. Overrides the query restrictions by adding a statement that is always true. |
| ' or 'a'='a | Modifies the structure of the query. Overrides the query restrictions by adding a statement that is always true. |
| +connection_id()-connection_id() | Modifies the query result to 0 |
| +1-1 | Modifies the query result to 0 |
| +67-ASCII('A') | Modifies the query result to 0 |
| +51-ASCII(1) | Modifies the query result to 0 |
| ... | ... |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 27

27

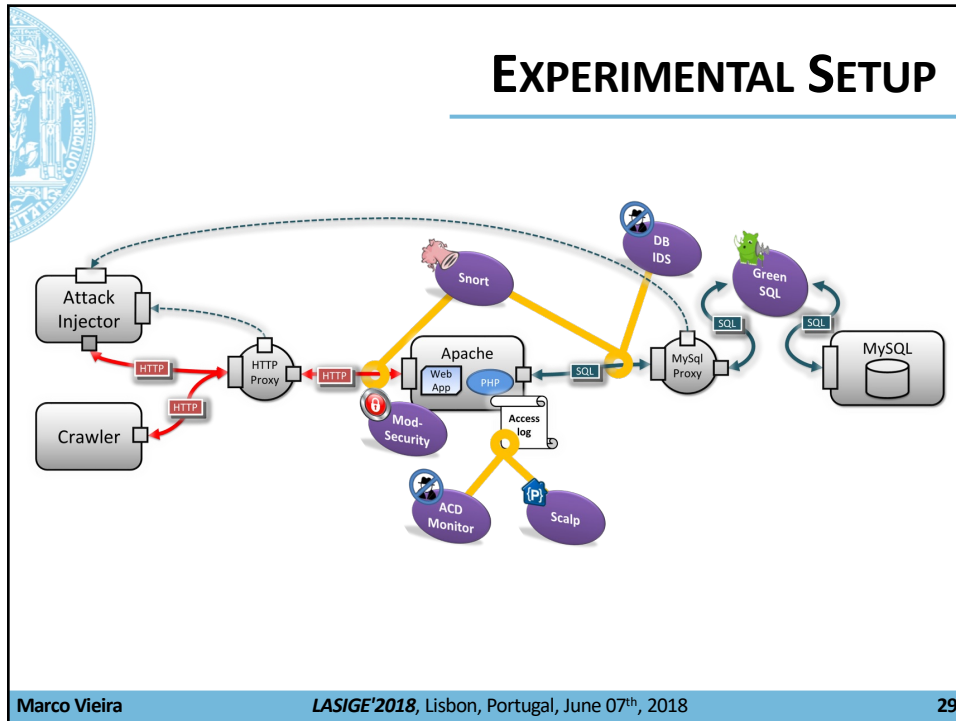


SYSTEMS UNDER BENCHMARKING

| Tool | Architectural Level monitored | Detection Approach | Data Source | Known Technology Limitations |
|------------------------------|-------------------------------|--------------------|---------------------|------------------------------|
| <i>ACD</i> | Application | Anomaly Based | Apache Log | Only GET method |
| <i>Apache Scalp</i> | Application | Signature Based | Apache Log | Only GET method |
| <i>ModSecurity</i> | Application | Signature Based | HTTP traffic | - |
| <i>Snort (v2.8 and v2.9)</i> | Network | Signature Based | Network Traffic | - |
| <i>GreenSQL</i> | Database | Signature Based | SQL Proxy Traffic | MySQL data |
| <i>DB IDS</i> | Database | Anomaly Based | SQL Sniffer Traffic | MySQL and Oracle data |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 28

28



29

MAIN RESULTS

All

| Ivl | Tool | Review | | | Reported | | | | Prec. | Recall | Mark. | Infor. |
|-----|-------------|--------|-----|------|----------|-----|-----|-----|-------|--------|-------|--------|
| | | P | N | Pop | TP | TN | FN | FP | | | | |
| App | ACD | 1051 | 224 | 1275 | 376 | 174 | 675 | 50 | 0.883 | 0.358 | 0.088 | 0.135 |
| | Scalp | 1051 | 224 | 1275 | 206 | 224 | 845 | 0 | 1.000 | 0.196 | 0.210 | 0.196 |
| | ModSecurity | 826 | 225 | 1051 | 236 | 225 | 590 | 0 | 1.000 | 0.286 | 0.276 | 0.286 |
| Net | Snort 2.8 | 458 | 817 | 1275 | 0 | 817 | 458 | 0 | - | 0.000 | - | 0.000 |
| DB | GreenSQL | 458 | 817 | 1275 | 244 | 813 | 214 | 4 | 0.984 | 0.533 | 0.775 | 0.528 |
| | DB IDS | 458 | 817 | 1275 | 451 | 384 | 7 | 433 | 0.510 | 0.985 | 0.492 | 0.455 |
| Net | Snort 2.9 | 173 | 878 | 1051 | 0 | 878 | 173 | 0 | - | 0.000 | - | 0.000 |

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 30

30

WHAT IS WRONG?

Established benchmarks are mostly for marketing!

- **Strict benchmarking conditions**
 - Fixed workload & faultload + Small set of metrics
- **Workload & faultload:**
 - May not be representative of the user scenario
- **Metrics:**
 - Fixed! May not satisfy the user needs
 - Decision based on several metrics is difficult!

No security benchmark endorsed by any organization or industry

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 31

31

FIXED!

```

graph LR
    Activation --> SUB
    SUB --> Metrics
    Activation -.-> Fixed
    Metrics -.-> Fixed
  
```

- **Example:**
 - Benchmarking vulnerability detection tools
 - Typical metric: F-Measure
 - Is this good in all scenarios?
 - Business critical: recall
 - Best effort: F-Measure
 - Minimum effort: Markedness

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 32

32

A POTENTIAL APPROACH...

Benchmarking conditions adaptable to the user needs

- Include multiple usage scenarios:
 - Metrics depend on the scenario
 - Adaptable workload and faultload
- Use quality models instead of independent metrics
 - Quality models should also adapt to the scenario

```

graph TD
    A11[A11] -- W1 --> Op1((Operator))
    A12[A12] -- W2 --> Op1
    Op1 -- W5 --> A1[A1]
    A21[A21] -- W3 --> Op2((Operator))
    A22[A22] -- W4 --> Op2
    Op2 -- W6 --> A2[A2]
    A1 -- W5 --> Op3((Operator))
    A2 -- W6 --> Op3
    Op3 --> FS[Final Score]
      
```

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 33


33

SCENARIOS AND QUALITY MODELS

How to define scenarios? How to define quality models? How to adapt workloads and faultloads to the scenarios?

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 34

34




CHALLENGES

- Satisfy industry requirements
 - Representativeness, portability, scalability, non-intrusiveness, low cost, ...
 - Prevent “gaming”
- Satisfy user requirements
 - Representativeness, usefulness, simplicity of use...
 - Adaptable – allow “gaming”
- Endorsement by TPC, SPEC, ...
 - **How to?**

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 35

35




IS THERE A FUTURE?

- Resilience Benchmarking
 - Assess and compare the behavior of components and computer systems when subjected to changes
 - Which resilience metrics?
 - Comparable, consistent, understandable, meaningful, ...
 - Changeloads:
 - Representative, practical, portable, ...
- Trustworthiness Benchmarking
 - What evidences to collect?
 - What metrics?
 - Dynamicity of perception... social trust...

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 36

36




CONCLUSIONS

The benchmarking concept is well established!

- Acceptance by “big” industry depends on perceived utility for marketing
- Acceptance by users requires “adaptability”
- From a research perspective, performance and dependability benchmarking are well known
- Security benchmarking approaches are weak
- New types of benchmarks will bring additional challenges!

Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 37

37



QUESTIONS?



Marco Vieira
 Department of Informatics Engineering
 University of Coimbra
mvieira@dei.uc.pt
<http://eden.dei.uc.pt/~mvieira>



Marco Vieira LASIGE'2018, Lisbon, Portugal, June 07th, 2018 38

38