# Flexible Tracing and Analysis of Applications' I/O Behavior

Tânia Esteves

Under the supervision of

**Prof. João Paulo**

**Prof. Rui Oliveira**

INESCTEC

HASLab
HIGH-ASSURANCE
SOFTWARE LABORATORY

Universidade do Minho

# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

▸ Large codebases

- Fluent Bit: ≈1M LoC, 5K files, 4 languages

- TensorFlow: >4M LoC, 20K files, 3K contributors

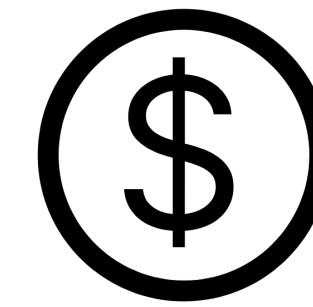▸ Several components

▸ Complex interactions (e.g., replication)

# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

  ▸ Large codebases

    - Fluent Bit: ≈1M LoC, 5K files, 4 languages

    - TensorFlow: >4M LoC, 20K files, 3K contributors

  ▸ Several components
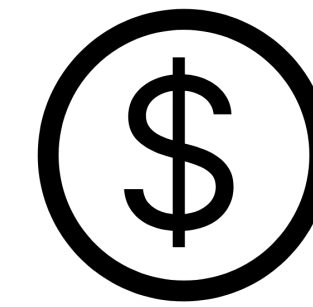
  ▸ Complex interactions (e.g., replication)

Healthcare

# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

  ‣ Large codebases

   - Fluent Bit: ≈1M LoC, 5K files, 4 languages

   - TensorFlow: >4M LoC, 20K files, 3K contributors

  ‣ Several components

  ‣ Complex interactions (e.g., replication)

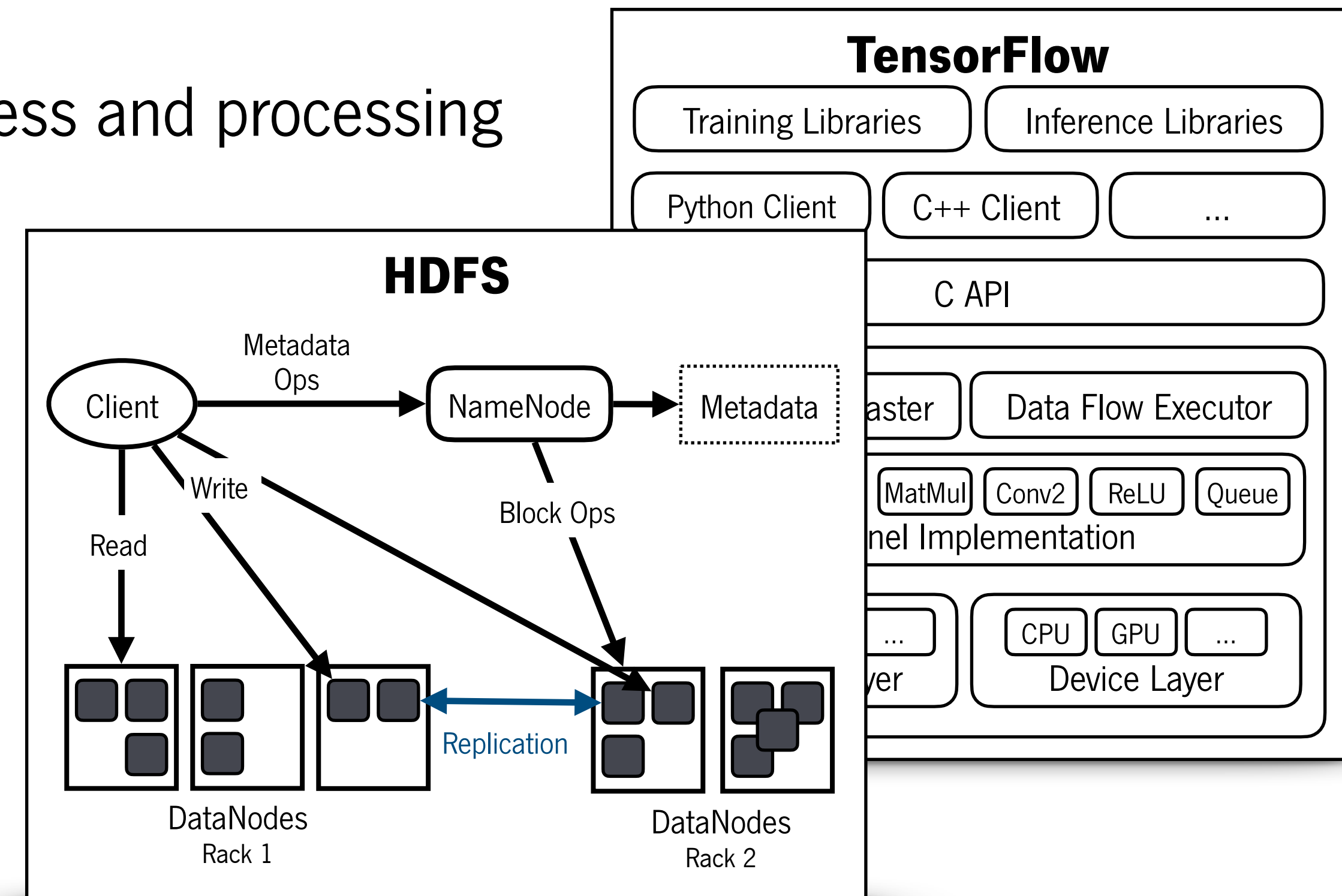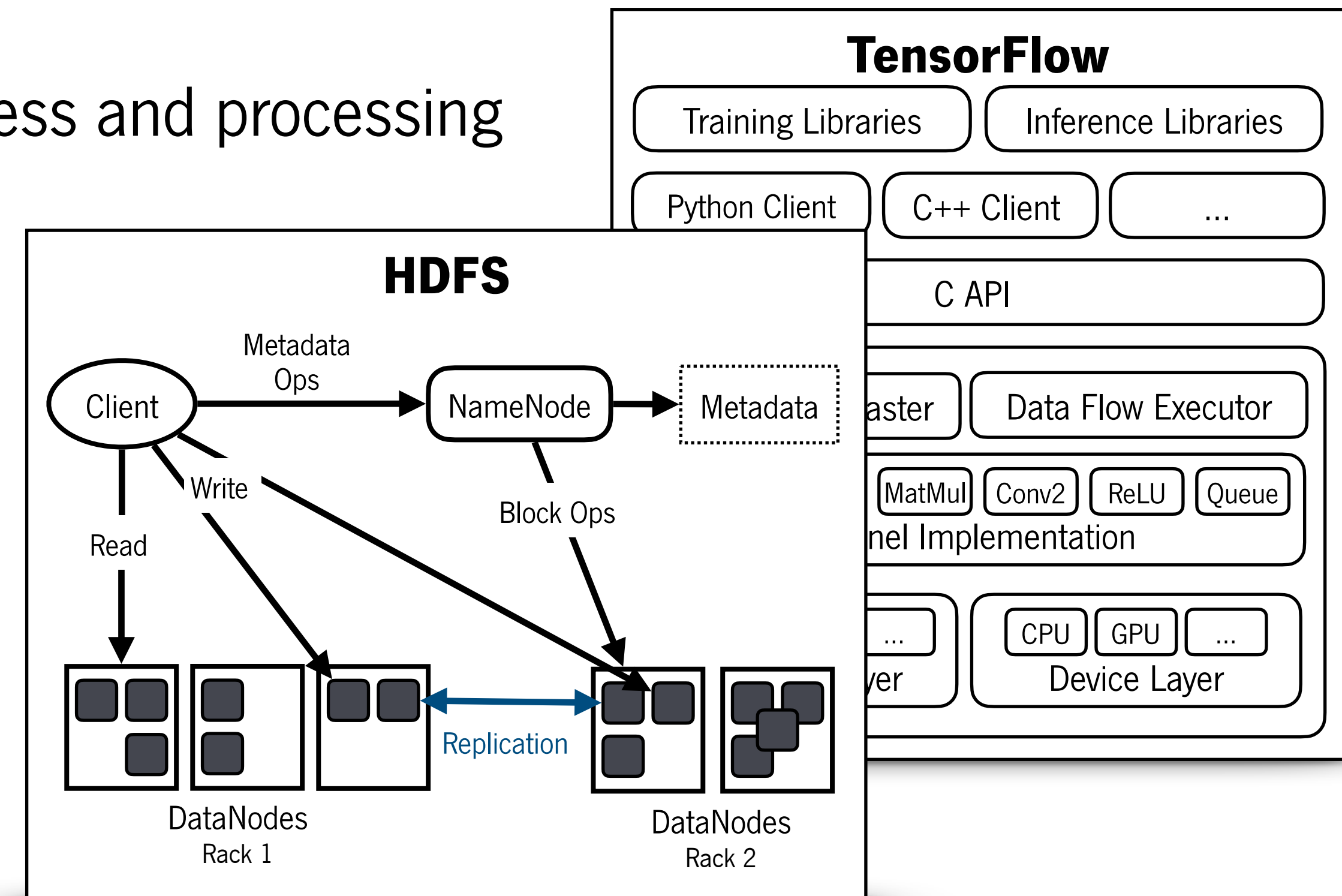Healthcare        Financial
                  Services

# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

▸ Large codebases

- Fluent Bit: ≈1M LoC, 5K files, 4 languages
- TensorFlow: >4M LoC, 20K files, 3K contributors

▸ Several components

▸ Complex interactions (e.g., replication)

Healthcare          Financial          Retail
                    Services

# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

▸ Large codebases
- Fluent Bit: ≈1M LoC, 5K files, 4 languages
- TensorFlow: >4M LoC, 20K files, 3K contributors

▸ Several components

▸ Complex interactions (e.g., replication)

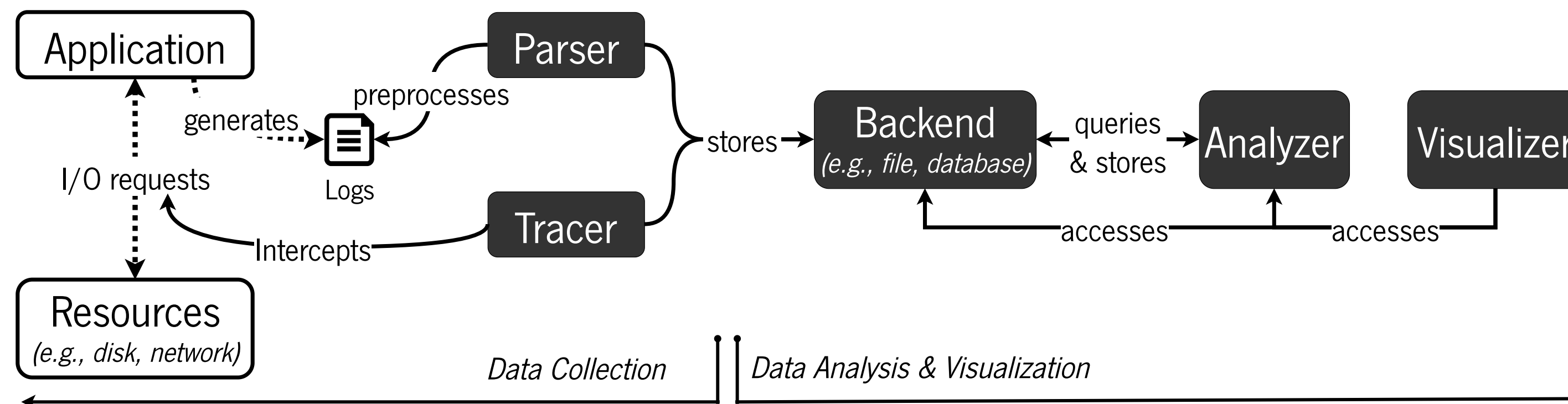# Data-centric and Distributed Systems

◉ Critical services increasingly rely on efficient data access and processing

◉ Complex architectures

▸ Large codebases
- Fluent Bit: ≈1M LoC, 5K files, 4 languages
- TensorFlow: >4M LoC, 20K files, 3K contributors

▸ Several components

▸ Complex interactions (e.g., replication)



**Question:** How can we ensure the correctness and good performance of these systems?

# Diagnosis Pipelines

⊙ Provide the <u>collection</u>, <u>analysis</u> and <u>visualization</u> of I/O requests made by applications

⊙ Useful for:

▸ <u>Debugging</u> - uncover the root cause of errors, inefficiencies and unattained performance

▸ <u>Validation</u> - validate applications' expected behaviors and the corrections of errors

▸ <u>Exploration</u> - understand how applications and storage systems handle data requests

# Challenges

# Challenges

**◎ #1 - Transparency**

▸ Provide <u>transparent</u> solutions that do not require modifications to applications' source code and are <u>generally applicable</u>

# Challenges

◉ **#1 - Transparency**

▸ Provide <u>transparent</u> solutions that do not require modifications to applications' source code and are <u>generally applicable</u>

◉ **#2 - Accuracy and Overhead**

▸ Balance the <u>amount and detail</u> of <u>collected data</u> with the <u>overhead</u> imposed on the targeted system

# Challenges

◎ **#1 - Transparency**

‣ Provide <u>transparent</u> solutions that do not require modifications to applications' source code and are <u>generally applicable</u>

◎ **#2 - Accuracy and Overhead**

‣ Balance the <u>amount and detail</u> of <u>collected data</u> with the <u>overhead</u> imposed on the targeted system

◎ **#3 - Analysis and Visualization**

‣ Offer <u>practical</u>, <u>integrated</u>, and <u>automated</u> components for analyzing, correlating, and visualizing I/O requests

# Challenges

## ◉ #1 - Transparency

▸ Provide <u>transparent</u> solutions that do not require modifications to applications' source code and are <u>generally applicable</u>

## ◉ #2 - Accuracy and Overhead

▸ Balance the <u>amount and detail</u> of <u>collected data</u> with the <u>overhead</u> imposed on the targeted system

## ◉ #3 - Analysis and Visualization

▸ Offer <u>practical</u>, <u>integrated</u>, and <u>automated</u> components for analyzing, correlating, and visualizing I/O requests

## ◉ #4 - Scope

▸ Design <u>comprehensive</u> solutions for diagnosing <u>different kinds</u> of I/O behaviors

# Contributions

◎ **Content-aware Diagnosis with CaT**

‣ Enables the collection and analysis of <u>distributed systems</u>' I/O requests

◎ **Comprehensive and Flexible Diagnosis with DIO**

‣ Provides customizable and insightful diagnosis of <u>data-centric applications</u>' storage I/O

◎ **Custom and Improved Analysis with CRIBA**

‣ Offers specialized and automated analysis of <u>cryptographic ransomware</u> I/O behavior

# Content-Aware Diagnosis

CAT, a framework for diagnosing
I/O flow of distributed systems

- Collects requests' context and content

- Combines causality inference with data similarity techniques

- Pinpoints data flow throughout the components of distributed systems

# CᴀT Design

Falcon's original component  Falcon's modified component  New components

# CᴀT Design



CᴀTʀᴀᴄᴇʀ

1 Collector

2 Handler

3 SigComp

CatLog

*<type, timestamp, pid, **content**, ...>*

☐ Falcon's original component     ⬚ Falcon's modified component     ▨ New components

◉ **CᴀTʀᴀᴄᴇʀ:** Kernel-level tracer that collects information about I/O requests (events), including their content

# CᴀT Design



Falcon's original component     Falcon's modified component     New components

◉**CᴀTʀᴀᴄᴇʀ:** Kernel-level tracer that collects information about I/O requests (events), including their content

◉**Tʀᴀᴄᴇ Pʀᴏᴄᴇssᴏʀ**: Parses and organizes events into different data structures

# CₐT Design



$<type, timestamp, pid, content,$ **$id, order, dependencies$**$, …>$

☐ Falcon's original component    ⬚ Falcon's modified component    ⬛ New components

◉**CₐTRACER:** Kernel-level tracer that collects information about I/O requests (events), including their content

◉**TRACE PROCESSOR**: Parses and organizes events into different data structures

◉**HB MODEL GENERATOR**: Infers the causality between events

# CₐT Design



**Falcon's original component**    ⌁ **Falcon's modified component**    ▮ **New components**

◉ **CₐTRACER:** Kernel-level tracer that collects information about I/O requests (events), including their content

◉ **TRACE PROCESSOR**: Parses and organizes events into different data structures

◉ **HB MODEL GENERATOR**: Infers the causality between events

◉ **CₐSOLVER**: Finds events with a high probability of operating over the same data flow

# CaT Design



$<type, timestamp, pid, content, id, order, dependencies, **similarities**,...>$

◻ Falcon's original component     ⦙ Falcon's modified component     ◼ New components

◉**CaTracer:** Kernel-level tracer that collects information about I/O requests (events), including their content

◉**Trace Processor**: Parses and organizes events into different data structures

◉**HB Model Generator**: Infers the causality between events

◉**CaSolver**: Finds events with a high probability of operating over the same data flow

◉**Visualizer**: Builds space-time diagrams representing the execution, the events' causal relationship and their data flow

# CₐT in Action
## Storage and Replication of a file in HDFS

◉ **HDFS:** Hadoop distributed file system composed of several DataNodes

◉ **Replication Process:**

▸ Clients send file(s) to one DataNode

▸ DataNodes forward data to other nodes and then persist it on disk

▸ The process is repeated until all DataNodes have the clients' data

◉ **3 Test Scenarios:**

▸ <u>Normal execution</u>

▸ <u>Storage corruption</u>: data modified before being persisted

▸ <u>Network corruption</u>: data modified before being transmitted

**DFS Client**

1. Write data          3. Ack

**DataNode 1**

2. Write data          3. Ack

**DataNode 1**

2. Write data          3. Ack

**DataNode 1**

# CᴀT in Action

**Storage and Replication of a file in HDFS**

# CAT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist client's data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
    client's data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

**b) Storage corruption**

Data persisted by DataNode 2
differs from clients' data

# CᴀT in Action
## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

**b) Storage corruption**

Data persisted by DataNode 2
differs from clients' data

# CAT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

**b) Storage corruption**

Data persisted by DataNode 2
differs from clients' data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist client's data

**b) Storage corruption**

Data persisted by DataNode 2 differs from clients' data

# CᴀT in Action

## Storage and Replication of a file in HDFS



a) **Normal execution**

All DataNodes persist
client's data

b) **Storage corruption**

Data persisted by DataNode 2
differs from clients' data

c) **Network corruption**

Data transmitted by DataNode
2 differs from clients' data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

**b) Storage corruption**

Data persisted by DataNode 2
differs from clients' data

**c) Network corruption**

Data transmitted by DataNode
2 differs from clients' data

# CᴀT in Action

## Storage and Replication of a file in HDFS



**a) Normal execution**

All DataNodes persist
client's data

**b) Storage corruption**

Data persisted by DataNode 2
differs from clients' data

**c) Network corruption**

Data transmitted by DataNode
2 differs from clients' data

# C<small>A</small>T Summary

◉ C<small>A</small>T's content-aware approach enables the detection of **data adulteration, corruption and leakage patterns** that would go **unnoticed with state-of-the-art context-based** solutions

◉ <u>**Open challenges**</u>:

▸ Comprehensive diagnosis of applications

▸ Practical and efficient analysis pipeline

# Comprehensive and Flexible Diagnosis

DIO, a generic tool for diagnosing applications' storage I/O

◉ Supports 42 storage-related system calls

◉ Collects their type, arguments, return value and extra context from the kernel

◉ Provides different strategies to customize the amount and detail of collected data

◉ Includes an integrated pipeline for near real-time analysis and visualization

# DIO Design

DIO's components      ⟶ DIO main flow      ⇢ App flow

# DIO Design



Server1

User-space

Kernel-space

Application

Syscalls

Storage Device

DIO's components     →  DIO main flow     ⇢  App flow

# DIO Design



DIO's tracer runs along the targeted application, intercepting its syscalls

Server1 — User-space: Application, Tracer — Kernel-space: Syscalls, Storage Device

DIO's components — DIO main flow — App flow

# DIO Design



Server1

User-space

Kernel-space

Application

Tracer

Syscalls

P    P

P

Storage Device

DIO's tracer runs along the targeted application, intercepting its syscalls

DIO's components   DIO main flow   App flow

# DIO Design



Server1

User-space

Application

Tracer

1 attach

Syscalls

P

P

P

Kernel-space

Storage Device

DIO's tracer runs along the targeted application, intercepting its syscalls

■ DIO's components ——→ DIO main flow ⤍ App flow

# DIO Design



DIO's tracer runs along the targeted application, intercepting its syscalls

# DIO Design



DIO's tracer runs along the targeted application, intercepting its syscalls

# DIO Design



DIO's tracer runs along the targeted application, intercepting its syscalls

# DIO Design



Server1

User-space

Application

Tracer

write()   read()

1 attach

3 collect

Syscalls

P   P

P

Ring
Buffer

Kernel-space

2 intercepts

Storage Device

DIO's tracer runs along the targeted
application, intercepting its syscalls

■ DIO's components   → DIO main flow   ---→ App flow

# DIO Design



Collected information is sent directly to the B*ackend* component, which is responsible for indexing and persisting it

# DIO Design



**Server1**

*User-space*

Application

*write()*    *read()*

**1** attach

Tracer

**4** send → Backend

**Server2**

**3** collect

Syscalls

P    P

P

Ring Buffer

*Kernel-space*

**2** intercepts

Storage Device

Collected information is sent directly to the B*ackend* component, which is responsible for indexing and persisting it

■ DIO's components    →  DIO main flow    ⇢  App flow

# DIO Design



Collected information is sent directly to the B*ackend* component, which is responsible for indexing and persisting it

# DIO Design



As soon as the data reaches the *Backend*, it becomes available for visualization at the Visualizer

# DIO Design

# DIO Design



Users can query directly the backend and build correlation algorithms, or visually explore the data at the Visualizer

DIO's components — DIO main flow - - - ▶ App flow

# DIO in Action
## Finding the root cause of RocksDB's performance anomalies

⊙ **RocksDB:** An embedded key-value store

⊙ **Problem:** RocksDB clients observe high tail latency (1 & 3)

  ▸ Reproducible with db_bench benchmark



99th percentile latency for RocksDB client operations.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies

◉ **RocksDB:** An embedded key-value store

◉ **Problem:** RocksDB clients observe high tail latency (1 & 3)

▸ Reproducible with db_bench benchmark



99th percentile latency for RocksDB client operations.

# DIO in Action

**Finding the root cause of RocksDB's performance anomalies**

◉ **RocksDB:** An embedded key-value store

◉ **Problem:** RocksDB clients observe high tail latency (1 & 3)

‣ Reproducible with db_bench benchmark



99th percentile latency for RocksDB client operations.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

▸ (1&3) Multiple background threads perform I/O simultaneously, db_bench performance decreases

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

▸ (1&3) Multiple background threads perform I/O simultaneously, db_bench performance decreases

▸ (2&4) Few background threads perform I/O simultaneously, db_bench performance improves

# DIO in Action

## Finding the root cause of RocksDB's performance anomalies



Syscalls issued by RocksDB over time, aggregated by thread name.

**Root Cause:** Interference between client writes, flushes and compactions

First observed in SILK (ATC'19) by <u>instrumenting</u> and <u>manually inspecting</u> more than 440K LoC

# DIO Summary

◉ DIO enables the diagnosis of storage **correctness**, **dependability** and **performance** issues and avoids the need for combining multiple tools and running the application multiple times



**Fluent Bit**
Identification of erroneous actions
that lead to data loss



**Elasticsearch**
Top-down exploration
and I/O diagnosis



**Redis**
Debugging and validation of
inefficient I/O patterns

# DIO Summary

◉ DIO enables the diagnosis of storage **correctness**, **dependability** and **performance** issues and avoids the need for combining multiple tools and running the application multiple times



**Fluent Bit**
Identification of erroneous actions
that lead to data loss



**Elasticsearch**
Top-down exploration
and I/O diagnosis



**Redis**
Debugging and validation of
inefficient I/O patterns

# DIO Summary

⦿ DIO enables the diagnosis of storage **correctness**, **dependability** and **performance** issues and avoids the need for combining multiple tools and running the application multiple times



**Fluent Bit**
Identification of erroneous actions
that lead to data loss



**Elasticsearch**
Top-down exploration
and I/O diagnosis



**Redis**
Debugging and validation of
inefficient I/O patterns

# Custom and Improved Analysis

CRIBA, a tool for diagnosing the I/O behavior of Linux cryptographic ransomware

◉ Supports the collection of 13 network-related system calls and system metrics

◉ Enhances the analysis process with 6 correlation algorithms and 8 dashboards tailored for ransomware characterization

# Cryptographic Ransomware

◉ **Malicious software** that **encrypts** victim's data and demands a **ransom**

◉ New ransomware families are constantly appearing

◉ CRIBA allows the observation of **characteristic** ransomware behavior:

▸ Traverses all victims' directories
▸ Rewrites victims' files with encrypted data
▸ Adds a new file extension to encrypted files
▸ Leaves ransom notes to inform the victim
▸ Has high CPU consumption due to encryption algorithms

# CRIBA in Action

⦿ Study with 5 different Linux cryptographic ransomware families

# CRIBA in Action

⊙ Study with 5 different Linux cryptographic ransomware families

⊙ Total of 26 different observations regarding generic statistics, ransom notes creation, data access and encryption patterns, and evasion techniques

# CRIBA in Action

◉ Study with 5 different Linux cryptographic ransomware families

◉ Total of 26 different observations regarding generic statistics, ransom notes creation, data access and encryption patterns, and evasion techniques

   ‣ **Metadata-related** operations are the most predominant (lseek, stat, fstat)

# CRIBA in Action



◉ Study with 5 different Linux cryptog

◉ Total of 26 different observations regarding generic statistics, ransom notes creation, data access and encryption patterns, and evasion techniques

‣ **Metadata-related** operations are the most predominant (lseek, stat, fstat)

‣ **Concurrent encryption actions** in RANSOMEXX that may lead to **data corruption**

# CRIBA in Action

◉ Study with 5 different Linux cryptographic ransomware families

◉ Total of 26 different observations regarding generic statistics, ransom notes creation, data access and encryption patterns, and evasion techniques

▸ **Metadata-related** operations are the most predominant (lseek, stat, fstat)

▸ **Concurrent encryption actions** in RANSOMEXX that may lead to **data corruption**

▸ **Different system call sequences** for file access (based on the targeted file)

# CRIBA in Action

◉ Study with 5 different Linux cryptographic ransomv

◉ Total of 26 different observations regarding generic
access and encryption patterns, and evasion techn

▸ **Metadata-related** operations are the most pre

▸ **Concurrent encryption actions** in RANSOMEX

▸ **Different system call sequences** for file access (based on the targeted file)

▸ Some families process only **partial content** of files or target **specific file extensions**

# CRIBA Summary

◉ CRIBA highlights **DIO's usefulness** by extending it with tailored analysis and visualization

◉ Our study shows that **different features** must be considered for a **clear understanding** of ransomware's **intrinsic** behavior

◉ This knowledge is key for improving detection tools for Linux cryptographic ransomware

# Conclusion

◉ **CAT**, a framework for diagnosing storage and network I/O requests of **distributed systems**

   ‣ Follows a <u>content-aware</u> approach that allows observing how <u>data flows</u> across components

   ‣ Useful for uncovering <u>data corruption</u> and <u>adulteration issues</u>

◉ **DIO**, a generic tool for diagnosing **data-centric applications**' storage I/O

   ‣ Provides the <u>flexibility</u> to narrow or broaden the <u>collection</u>, <u>analysis</u> and <u>visualization</u> of I/O behaviors

   ‣ Useful for <u>debugging</u>, <u>validating</u> and <u>exploring</u> both known and unknown storage patterns

◉ **CRIBA**, a practical tool for characterizing the I/O behavior of **cryptographic ransomware**

   ‣ <u>Automates</u> the analysis and visualization of <u>specific</u> ransomware behaviors

   ‣ Useful for better <u>understand</u> ransomware attacks and <u>enhance</u> detection tools

# Publications

## Core Publications

- **Tânia Esteves**, Francisco Neves, Rui Oliveira and João Paulo. **"CaT: Content-aware Tracing and Analysis of Distributed Systems"**. In *22nd International Middleware Conference*, 2021.

- **Tânia Esteves**, Ricardo Macedo, Rui Oliveira and João Paulo. **"Diagnosing Applications' I/O Behavior through System Call Observability"**. In *53rd Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops*, 2023.

- **Tânia Esteves**, Bruno Pereira, Rui Pedro Oliveira, João Marco and João Paulo. **"CRIBA: A Tool for Comprehensive Analysis of Cryptographic Ransomware's I/O Behavior"**. In *42nd Symposium on Reliable Distributed Systems*, 2023.

- **Tânia Esteves**, Ricardo Macedo, Rui Oliveira and João Paulo. **"Toward a Practical and Timely Diagnosis of Applications' I/O Behavior"**. In *IEEE Access*, 2023.

## Complementary Publications

- Mariana Miranda, **Tânia Esteves**, Bernardo Portela and João Paulo. **"S2Dedup: SGX-enabled Secure Deduplication"**. In *14th ACM International Conference on Systems and Storage*, 2021.

- **Tânia Esteves**, Ricardo Macedo, Alberto Faria, Bernardo Portela, João Paulo, José Pereira and Danny Harnik. **"TrustFS: An SGX-enable Stackable File System Framework"**. In *38th International Symposium on Reliable Distributed Systems Workshops*, 2019.