MAP-i Doctoral Program in Computer Science

# User-level Software-Defined Storage Data Planes

**Ricardo Macedo**

Under the supervision of
**João Tiago Paulo**
**José Orlando Pereira**

INESCTEC    HASLab HIGH-ASSURANCE SOFTWARE LABORATORY    Universidade do Minho    MAPi DOCTORAL PROGRAMME IN COMPUTER SCIENCE

# Data-centric systems

- Data-centric systems have become an integral part of modern I/O stacks

- Good performance for these systems requires storage optimizations
  - Scheduling, caching, tiering, ...

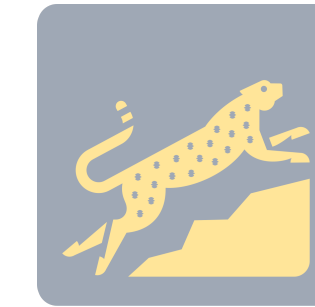- Optimizations are implemented sub-optimally

# Data-centric systems

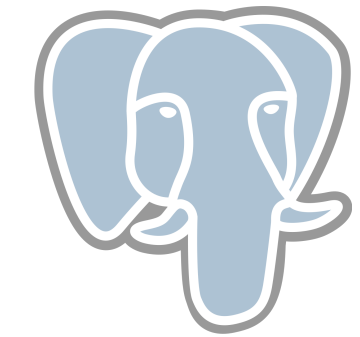- Data-centric systems have become an integral part of modern I/O stacks

- Good pe~~~~~~~~~~~~yTorch
  optimiza~~~~~

  - Sche~~~

- Optimizations are implemented sub-optimally

---

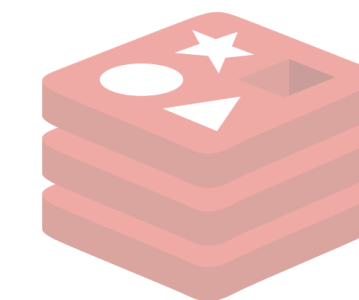**There is a better way to implement I/O optimizations**

# Challenge #1

❌ **Tightly coupled optimizations**

- I/O optimizations are single purposed

- Require deep understanding of the system's internal operation model

- Require profound system refactoring

- Have limited portability across systems



**Application**

**Key-Value Store**

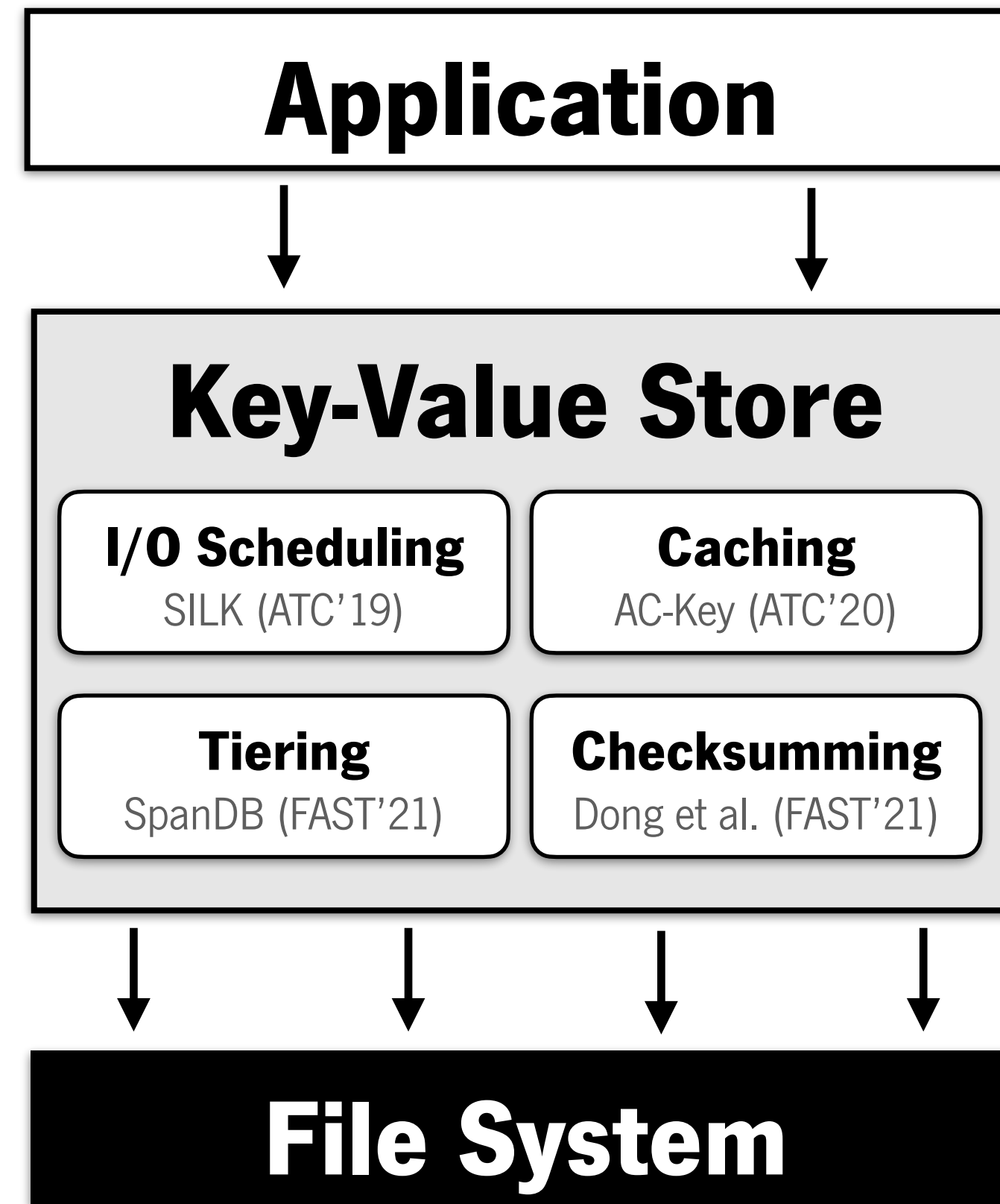| I/O Scheduling | Caching |
|---|---|
| SILK (ATC'19) | AC-Key (ATC'20) |
| **Tiering** | **Checksumming** |
| SpanDB (FAST'21) | Dong et al. (FAST'21) |

**File System**

# Challenge #1

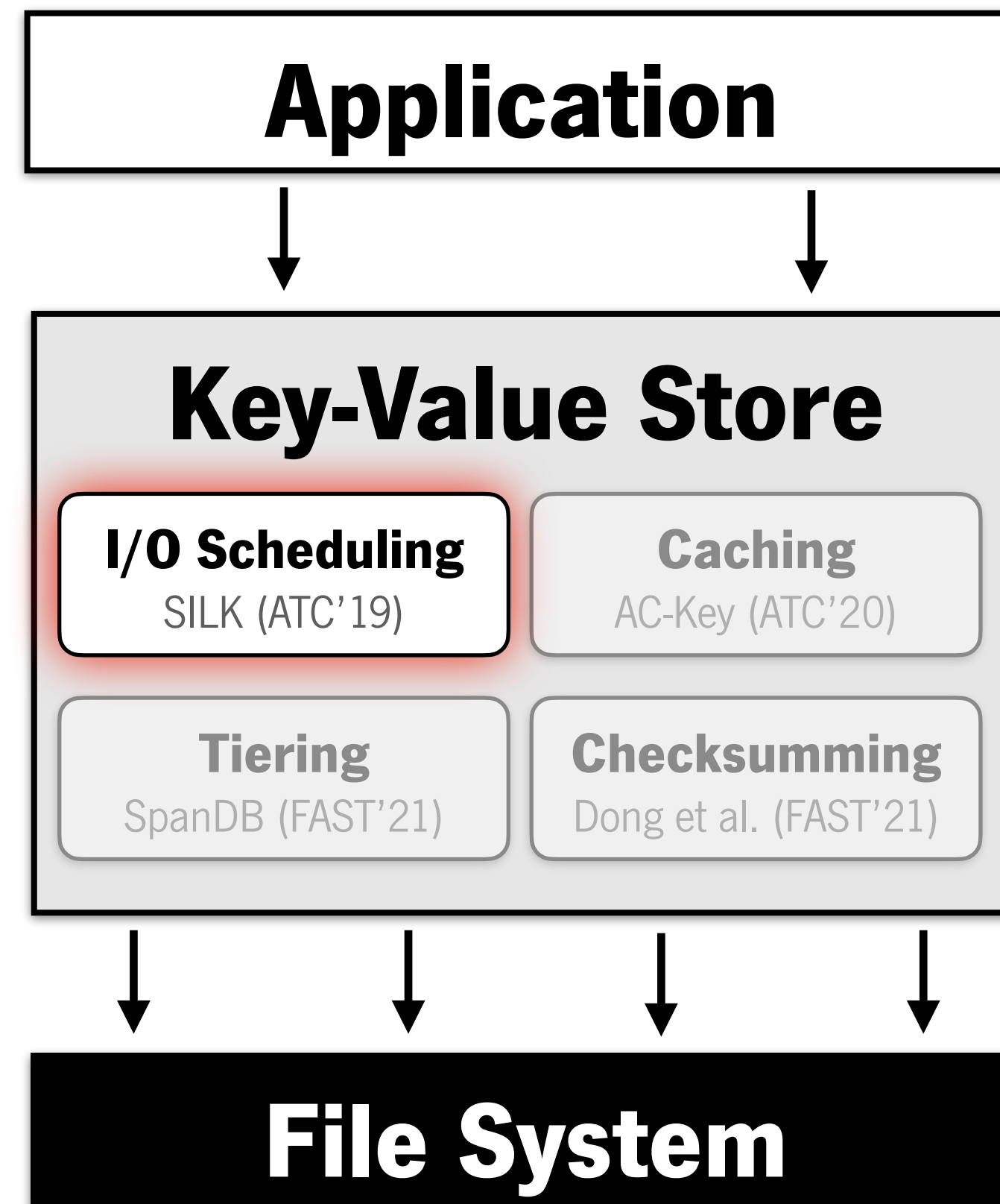❌ **Tightly coupled optimizations**

- I/O optimizations are single purposed

- Require deep understanding of the system's internal operation model

- Require profound system refactoring

- Have limited portability across systems

**Application**

↓      ↓

**Key-Value Store**

| I/O Scheduling<br>SILK (ATC'19) | Caching<br>AC-Key (ATC'20) |
| Tiering<br>SpanDB (FAST'21) | Checksumming<br>Dong et al. (FAST'21) |

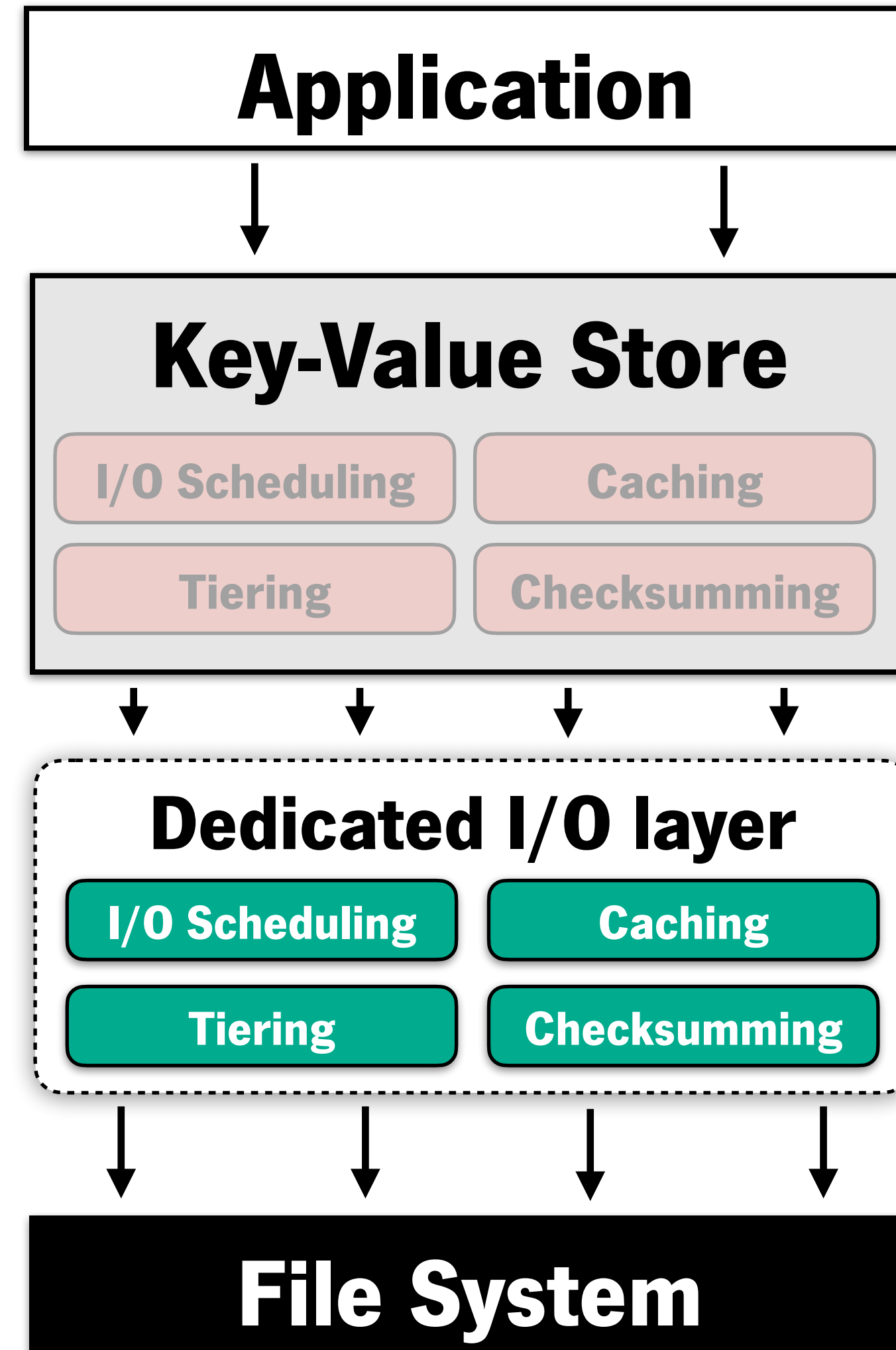↓    ↓    ↓    ↓

**File System**

**SILK's I/O Scheduler**

- Reduces tail latency spikes in RocksDB

- Controls the interference between foreground and background tasks

- Requires changing several modules, such as background operation handlers, internal queuing logic, and thread pools

# Challenge #1

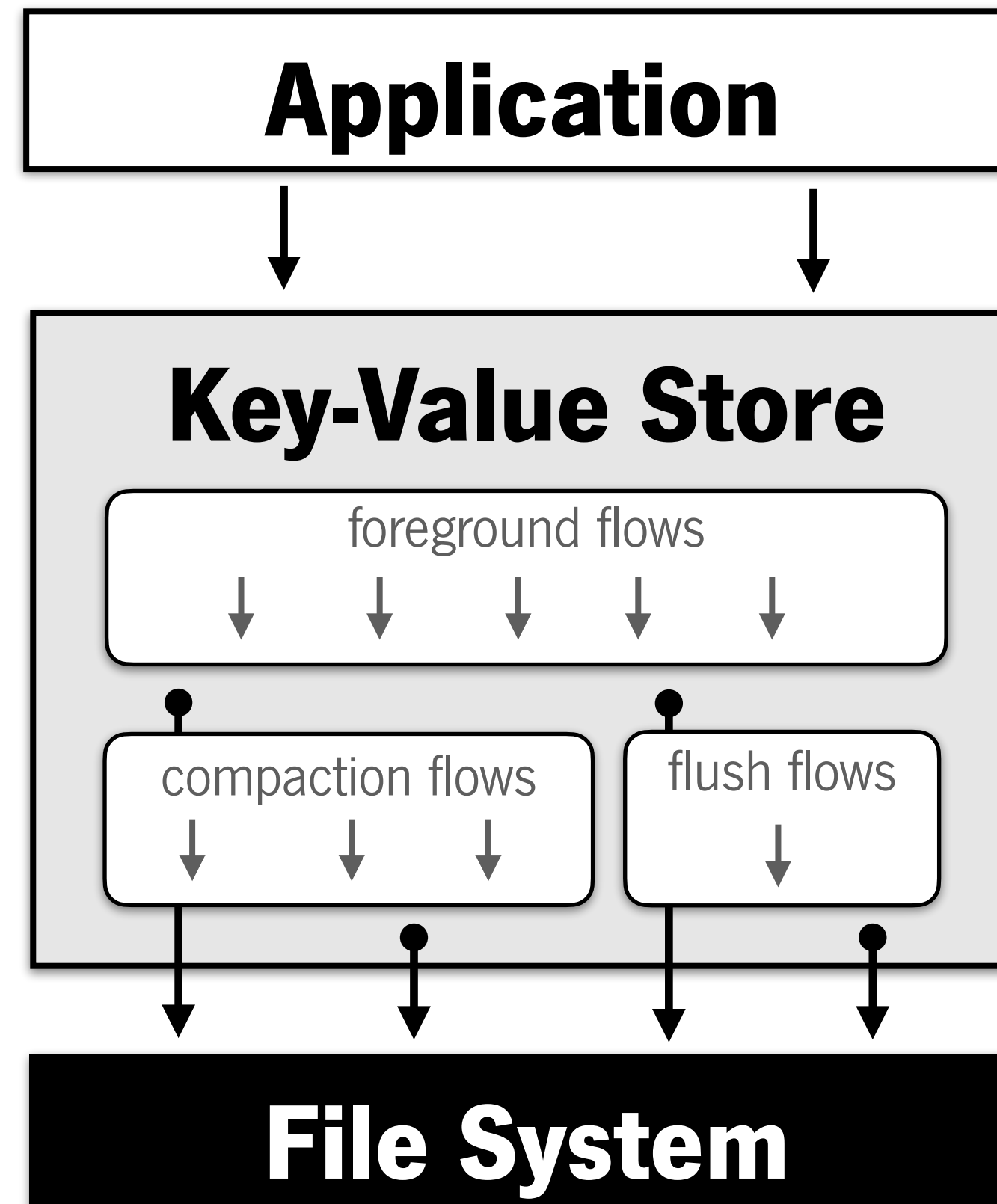✓ **Decoupled optimizations**

- I/O optimizations should be <u>disaggregated</u> from the internal logic of applications

- Moved to a <u>dedicated I/O layer</u>

- <u>Generally applicable</u>

- <u>Portable</u> across different scenarios

# Challenge #2

❌ **Rigid interfaces**

- Decoupled optimizations <u>lose granularity</u> and <u>internal</u> application <u>knowledge</u>

- I/O layers expose <u>rigid interfaces</u>

- <u>Discard information</u> that could be used to classify and differentiate requests

# Challenge #2

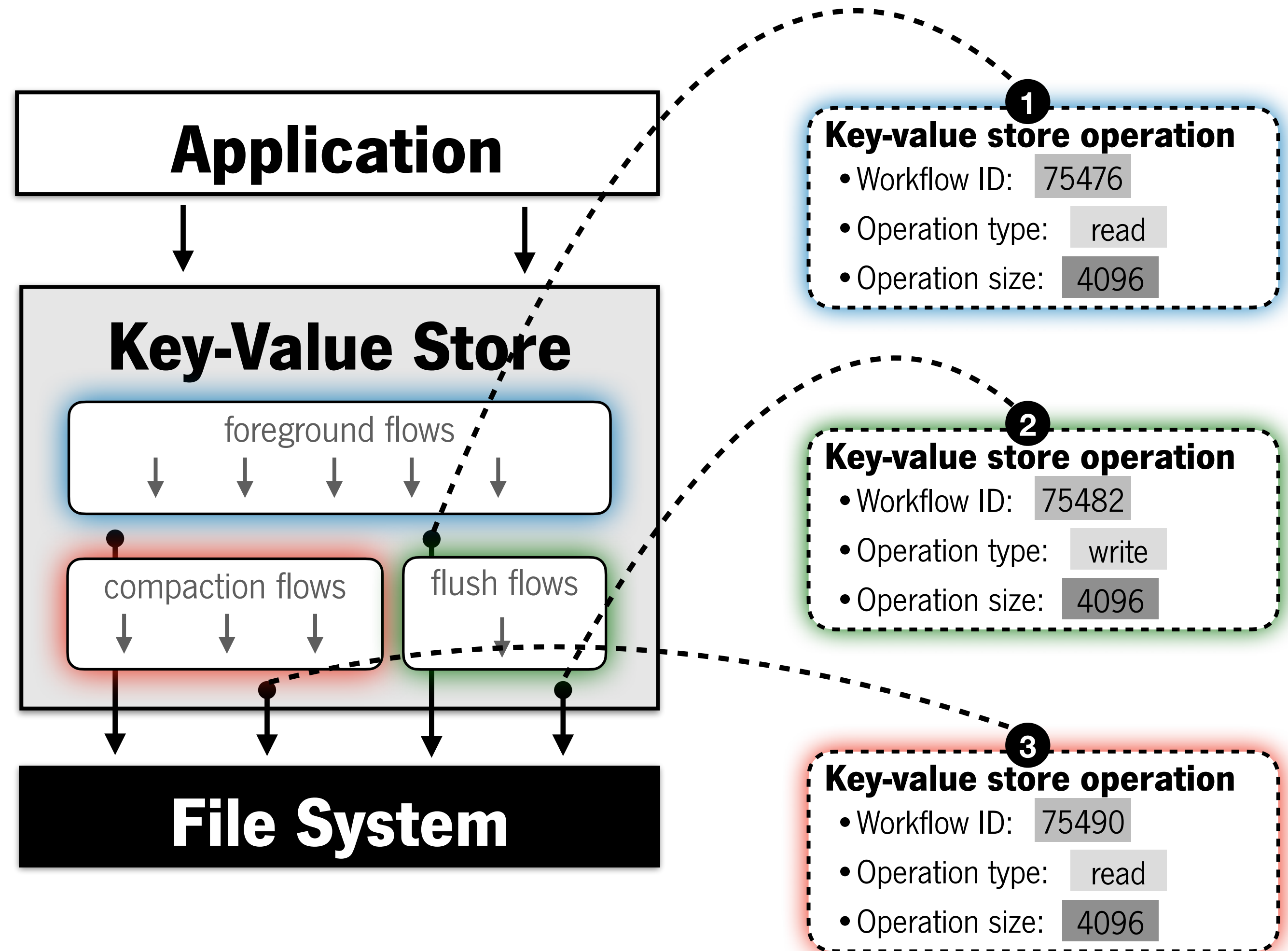❌ **Rigid interfaces**

- Decoupled optimizations <u>lose granularity</u> and <u>internal</u> application <u>knowledge</u>

- I/O layers expose <u>rigid interfaces</u>

- <u>Discard information</u> that could be used to classify and differentiate requests
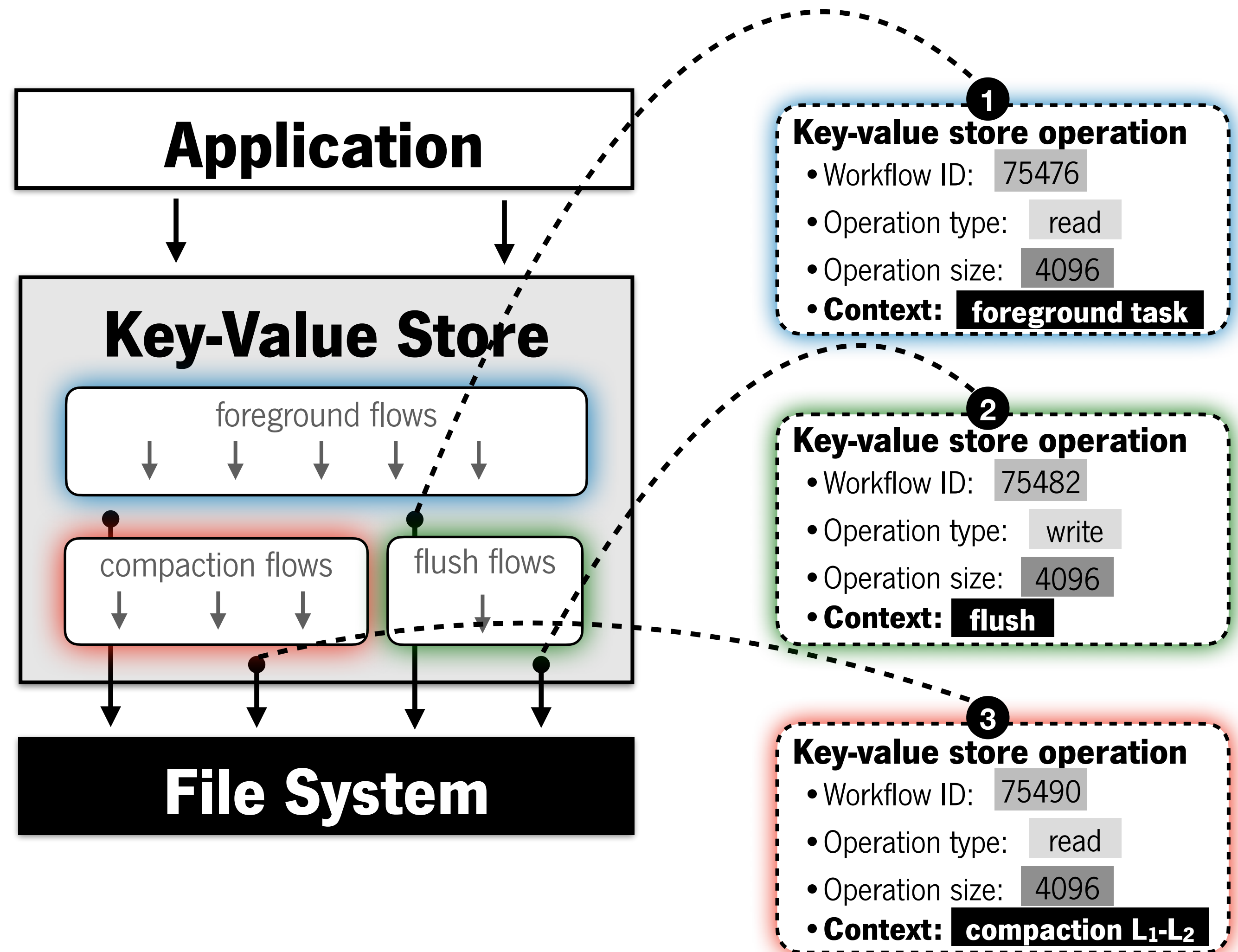


**Application**

**Key-Value Store**

foreground flows

compaction flows | flush flows

**File System**

**① Key-value store operation**
- Workflow ID: 75476
- Operation type: read
- Operation size: 4096

**② Key-value store operation**
- Workflow ID: 75482
- Operation type: write
- Operation size: 4096

**③ Key-value store operation**
- Workflow ID: 75490
- Operation type: read
- Operation size: 4096

# Challenge #2

✓ **Information propagation**
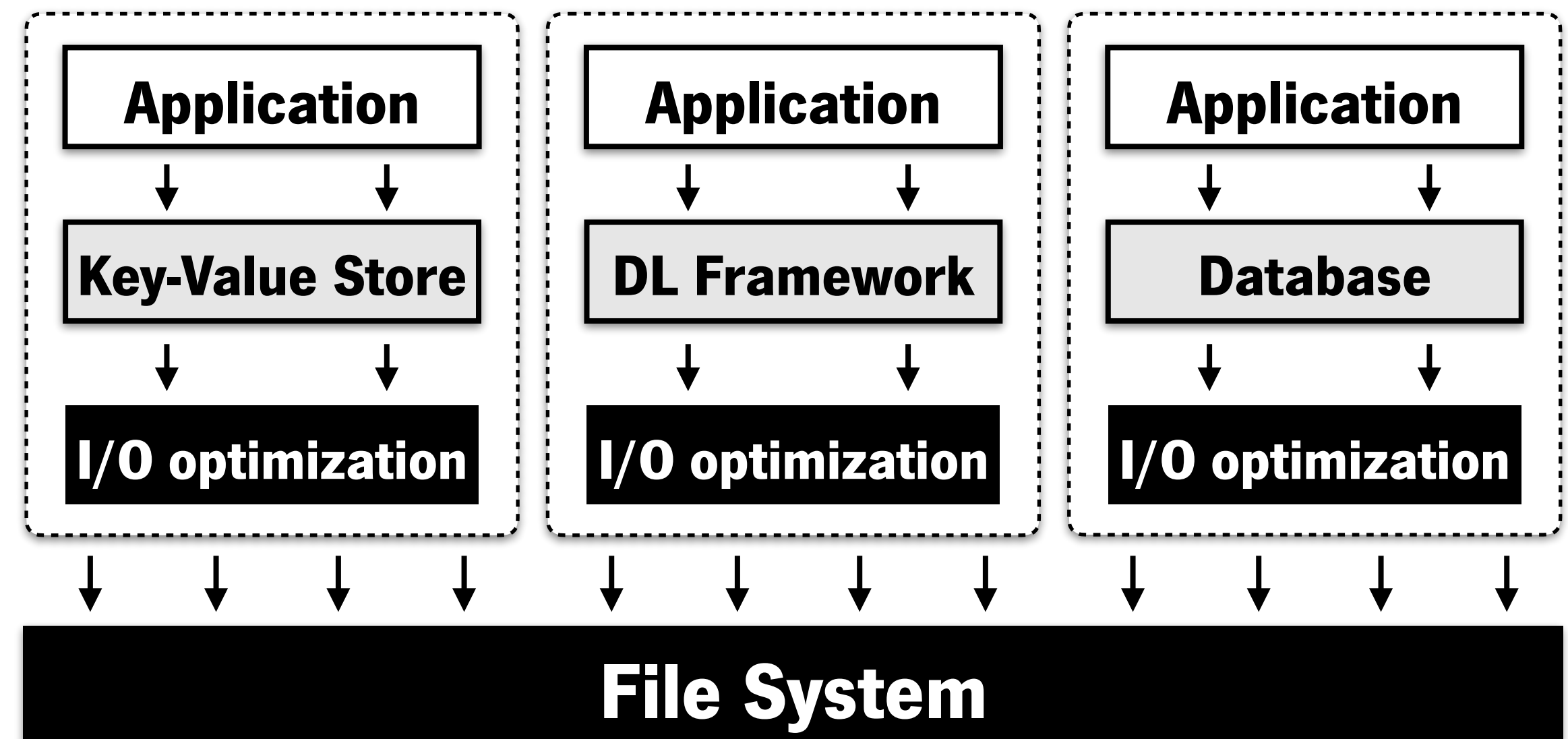
- Application-level information must be propagated throughout layers

- Decoupled optimizations can provide the same level of control and performance



**Application**

**Key-Value Store**

foreground flows

compaction flows    flush flows

**File System**

**1**
**Key-value store operation**
- Workflow ID: 75476
- Operation type: read
- Operation size: 4096
- **Context: foreground task**

**2**
**Key-value store operation**
- Workflow ID: 75482
- Operation type: write
- Operation size: 4096
- **Context: flush**

**3**
**Key-value store operation**
- Workflow ID: 75490
- Operation type: read
- Operation size: 4096
- **Context: compaction $L_1$-$L_2$**

# Challenge #3

## ⊗ Partial visibility

- <u>Optimizations are oblivious</u> of other systems

- <u>Lack of coordination</u>

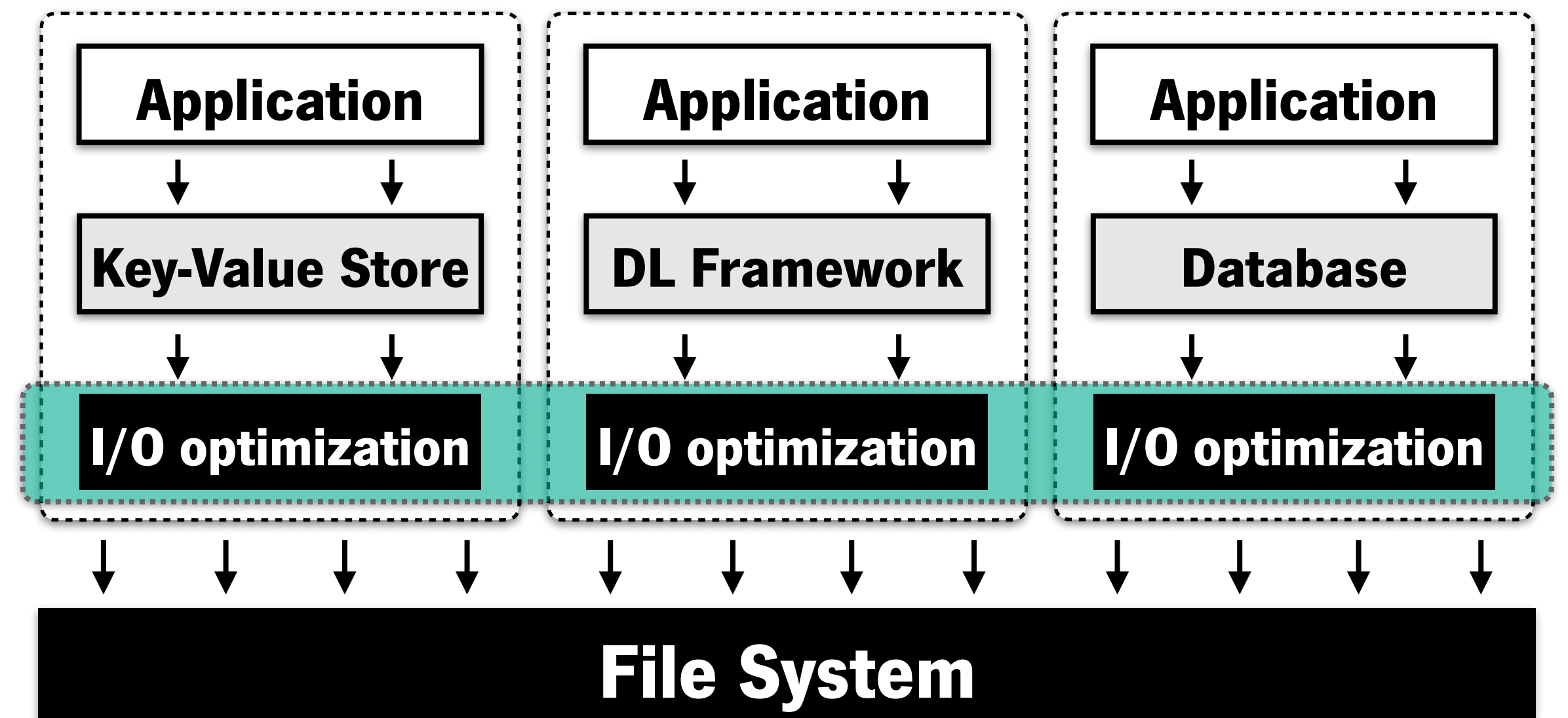- Conflicting optimizations, I/O contention, and performance variation



> ⚠ **Note:** the storage backend can either be local (e.g., ext4, xfs) or distributed (e.g., Lustre, GPFS), as well as the I/O layers on top

# Challenge #3

## ✅ Global I/O control

- Optimizations should be <u>aware</u> of the surrounding system stack

- Operate in <u>coordination</u>

- <u>Holistic control</u> of I/O workflows and shared resources



**Note:** the storage backend can either be local (e.g., ext4, xfs) or distributed (e.g., Lustre, GPFS), as well as the I/O layers on top
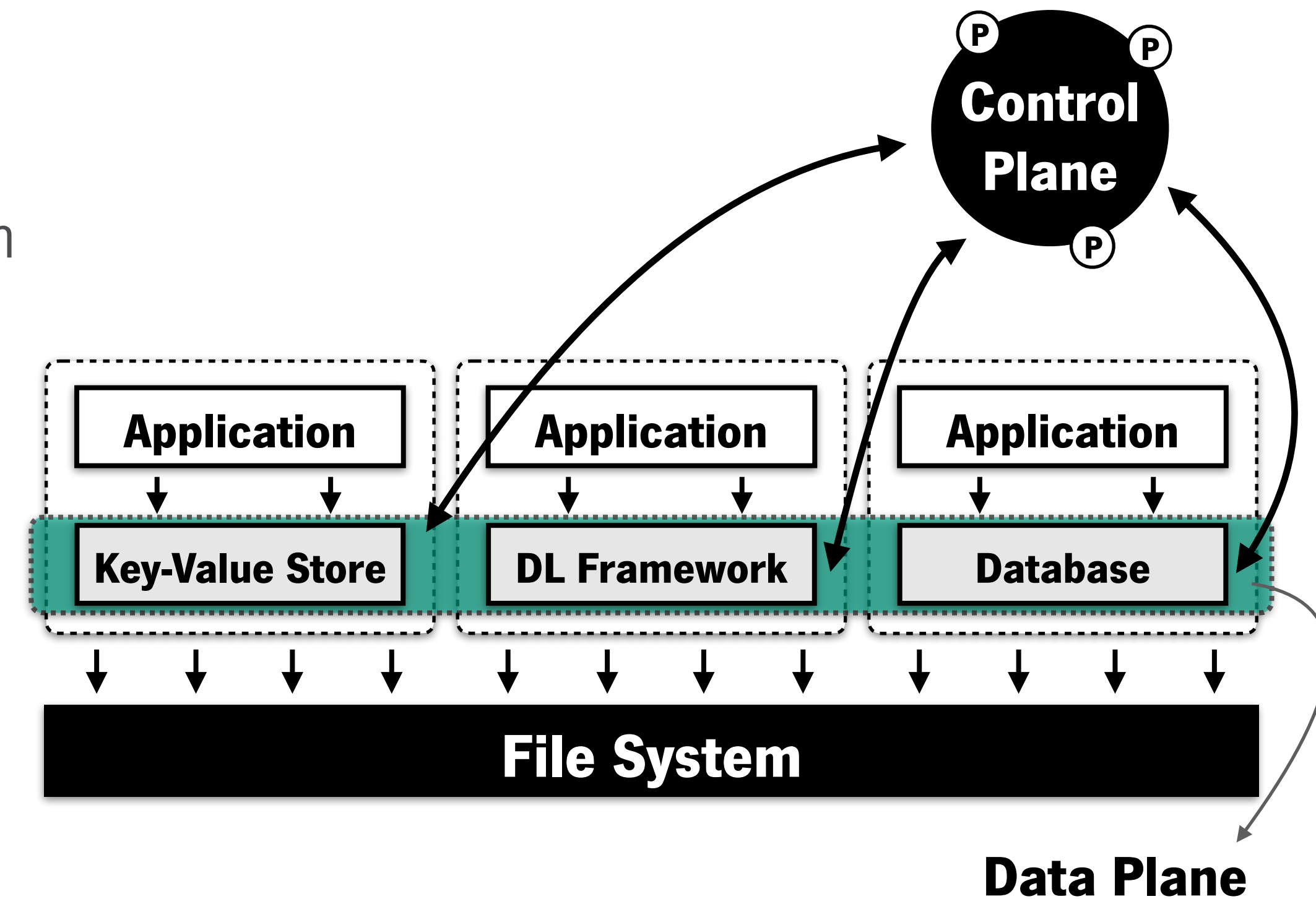
# Objectives

## Redefine how I/O optimizations are implemented

- <u>Decoupled</u> from the targeted system, minimizing intrusiveness

- Perform <u>coordinated</u> decisions over shared resources

- Impose <u>minimal performance overhead</u>

- <u>Programmable</u> and <u>adaptable</u> to different requirements and storage objectives
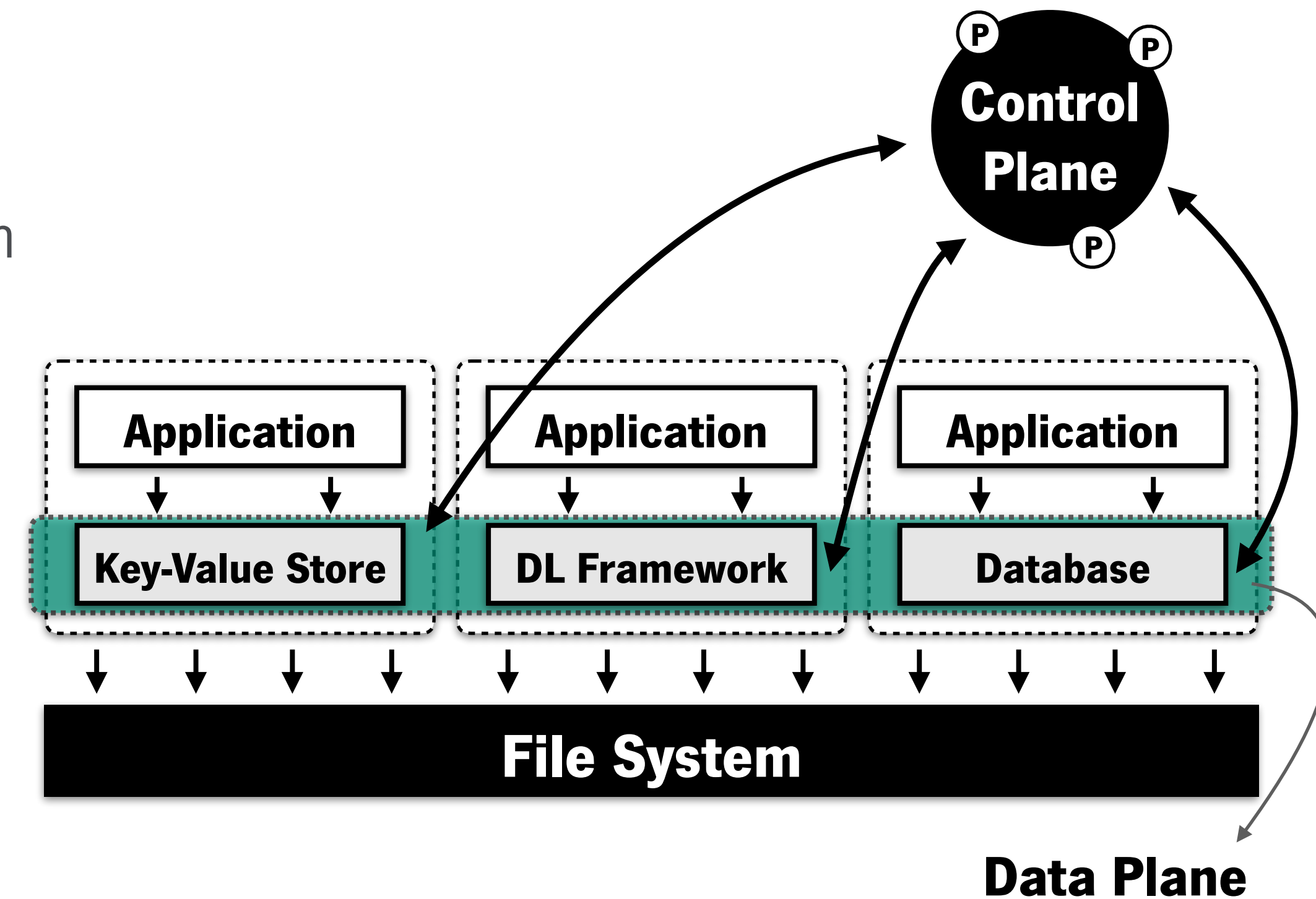
# Software-Defined Storage

- **Software-Defined Storage** (SDS) decouples I/O mechanisms from the policies that govern them

- **Control plane** acts as a global coordinator that enforces policies holistically

  - QoS provisioning, performance control, resource fairness

- **Data plane** is a multi-stage component that implements custom I/O logic over requests

  - I/O schedulers, encryption, compression, and caching

# Software-Defined Storage

- **Software-Defined Storage** (SDS) decouples I/O mechanisms from the policies that govern them

- **Control plane** acts as a global coordinator that enforces policies holistically
  - QoS provisioning, performance control, resource fairness

- **Data plane** is a multi-stage component that implements custom I/O logic over requests
  - I/O schedulers, encryption, compression, and caching



**Control Plane**

| Application | Application | Application |
|---|---|---|
| Key-Value Store | DL Framework | Database |

**File System**

**Data Plane**

## Survey and classification of SDS systems

- Targeted for **specific I/O layers** or **storage objectives** (e.g., virtualization, file system, resource management)

- Tightly coupled design, driven by the architecture and specificities of the context they are applied

- Existing SDS systems follow a **similar path** as traditionally implemented I/O optimizations

**Macedo** *et al. "A Survey and Classification of Software-Defined Storage Systems".* ACM Computing Surveys, 2020.
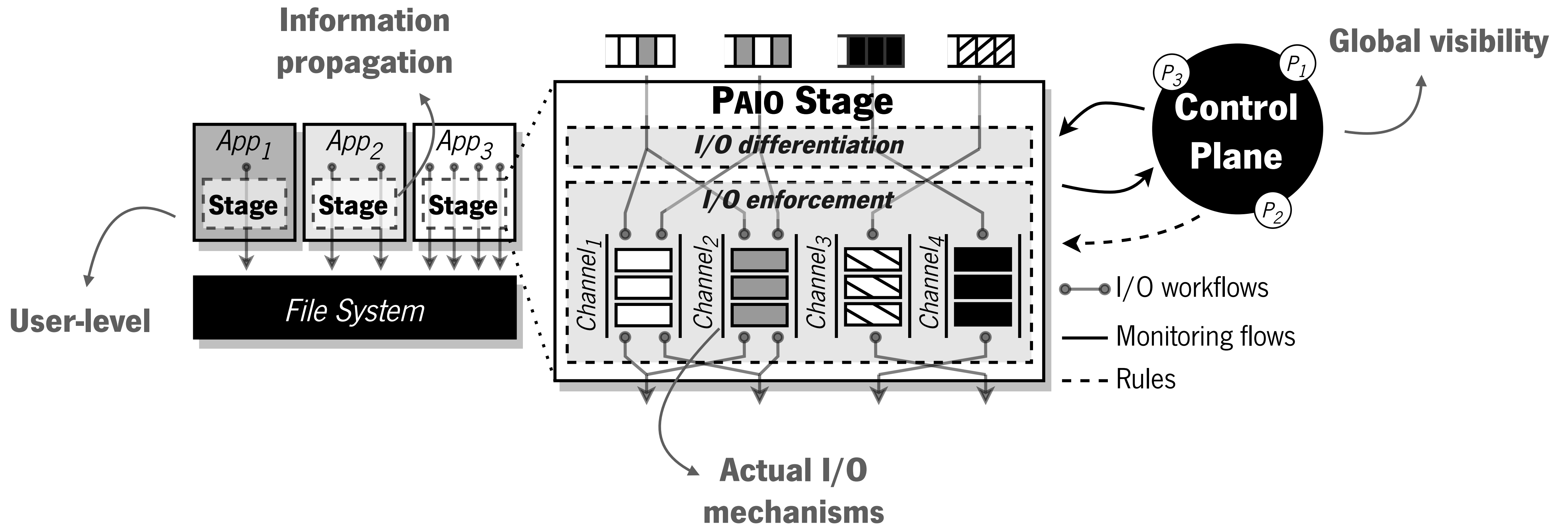
# Contributions

- **Software-Defined Storage survey**

  - Systematization of knowledge, taxonomy, and classification of existing SDS work

- **PAIO** data plane framework

  - Enables building user-level, portable, and generally applicable I/O optimizations

- **Data plane stages** built with PAIO

  - **Tail latency control** in LSM-based key-value stores

  - **Per-application bandwidth control** under shared storage environments
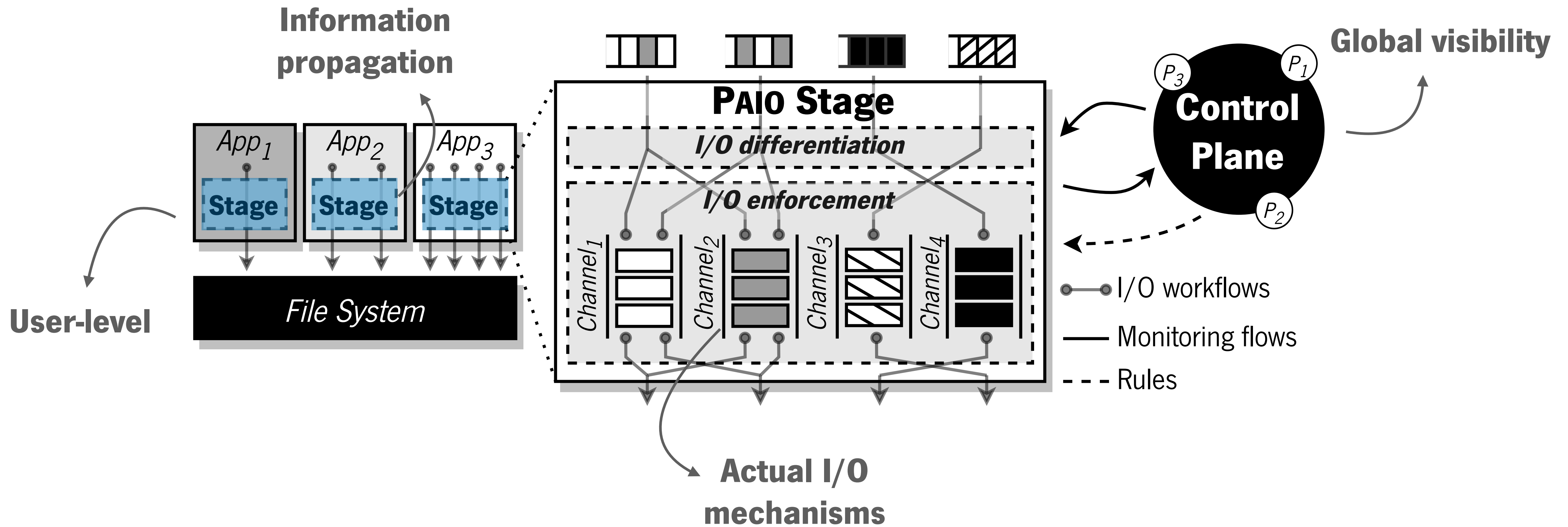
  - **Metadata control** in parallel file systems

# PAIO: Programmable and Adaptable I/O Workflows

- **User-level** framework for building **portable** and **generally applicable** I/O optimizations

- Follows a Software-Defined Storage design

  - I/O optimizations are implemented **outside** applications as **data plane stages**

  - **Stages** are controlled through a **control plane** for coordinated access to resources

- Enables the propagation of application-level information through **context propagation**

- Porting I/O layers to use PAIO requires **none to minor** code changes

**Macedo** *et al. "PAIO: General, Portable I/O Optimizations With Minor Application Modifications". USENIX FAST, 2022.*
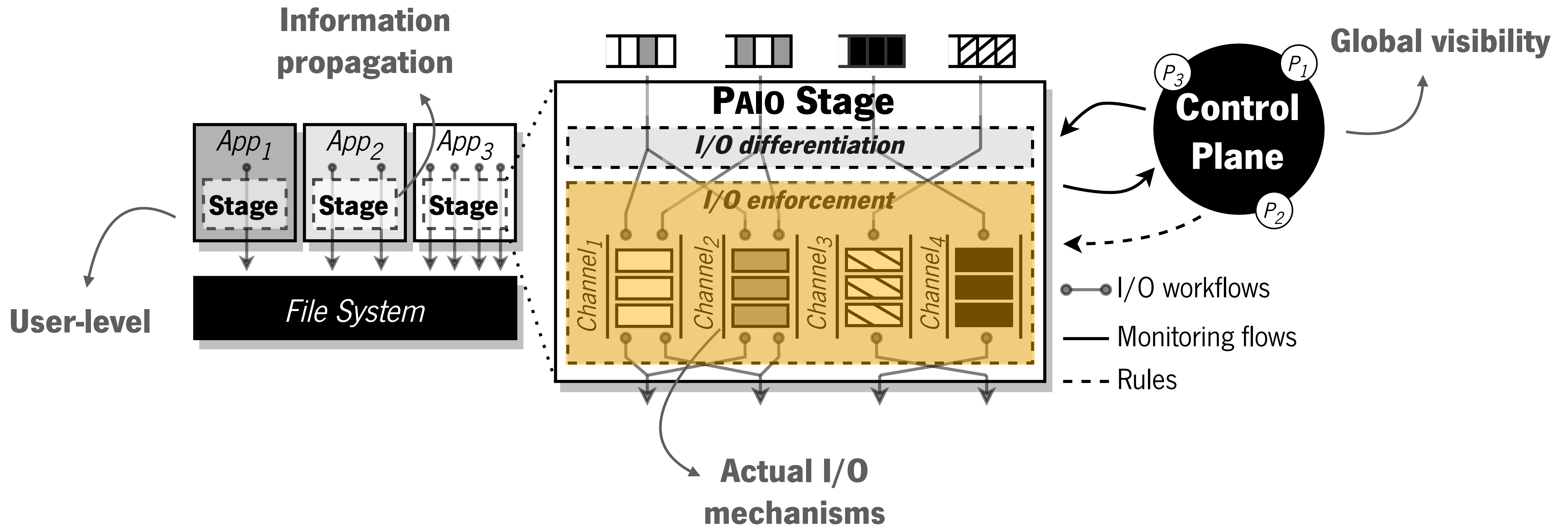
# PAIO design
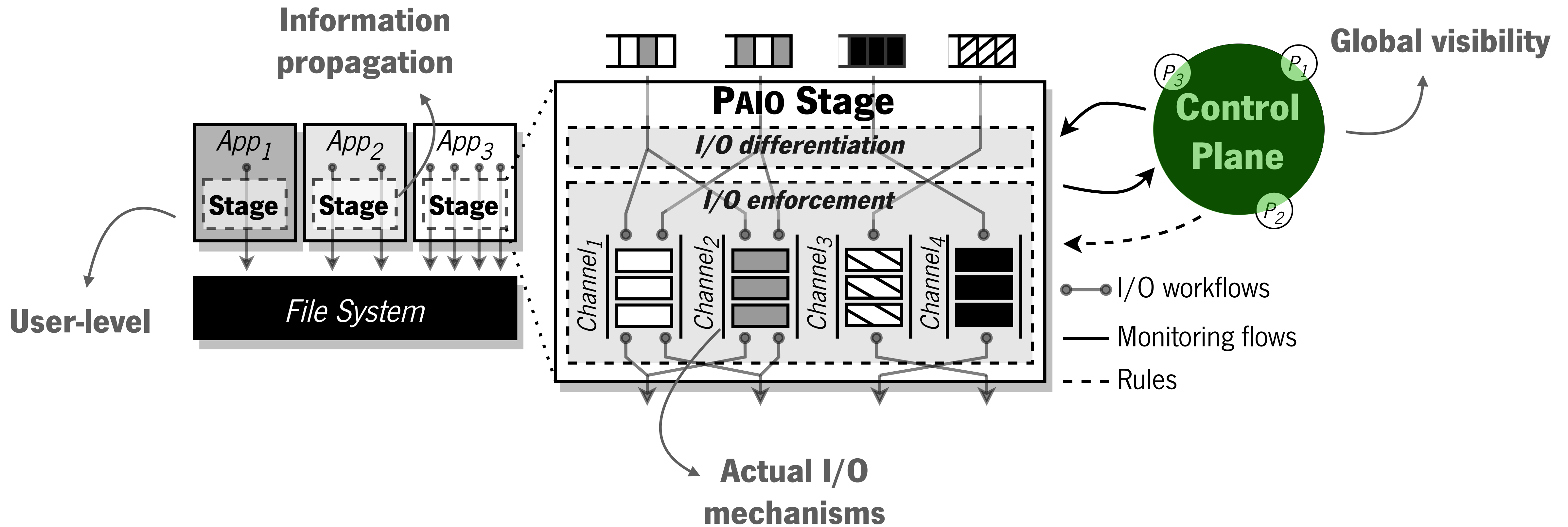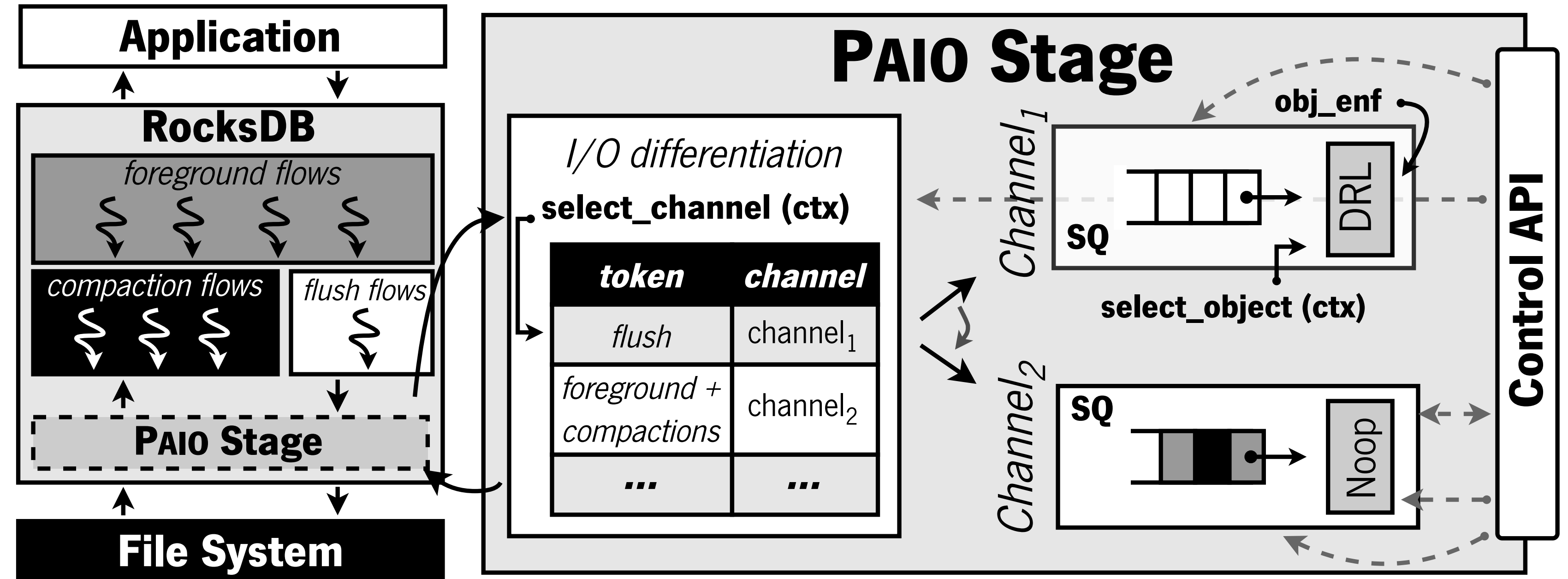
# PAIO design

# PAIO design

# PAIO design

# PAIO design

# PAIO design

- Context propagation
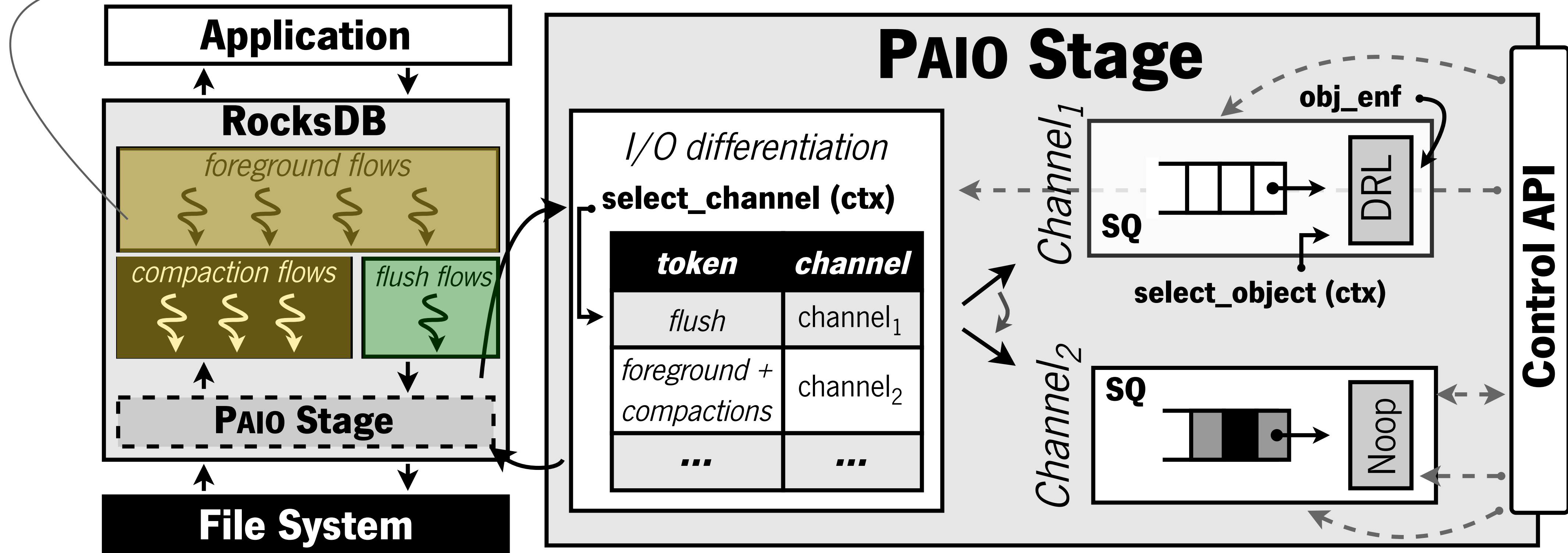- I/O differentiation
- I/O enforcement



**Policy:** limit the rate of RocksDB's flush operations to X MiB/s

# I/O differentiation



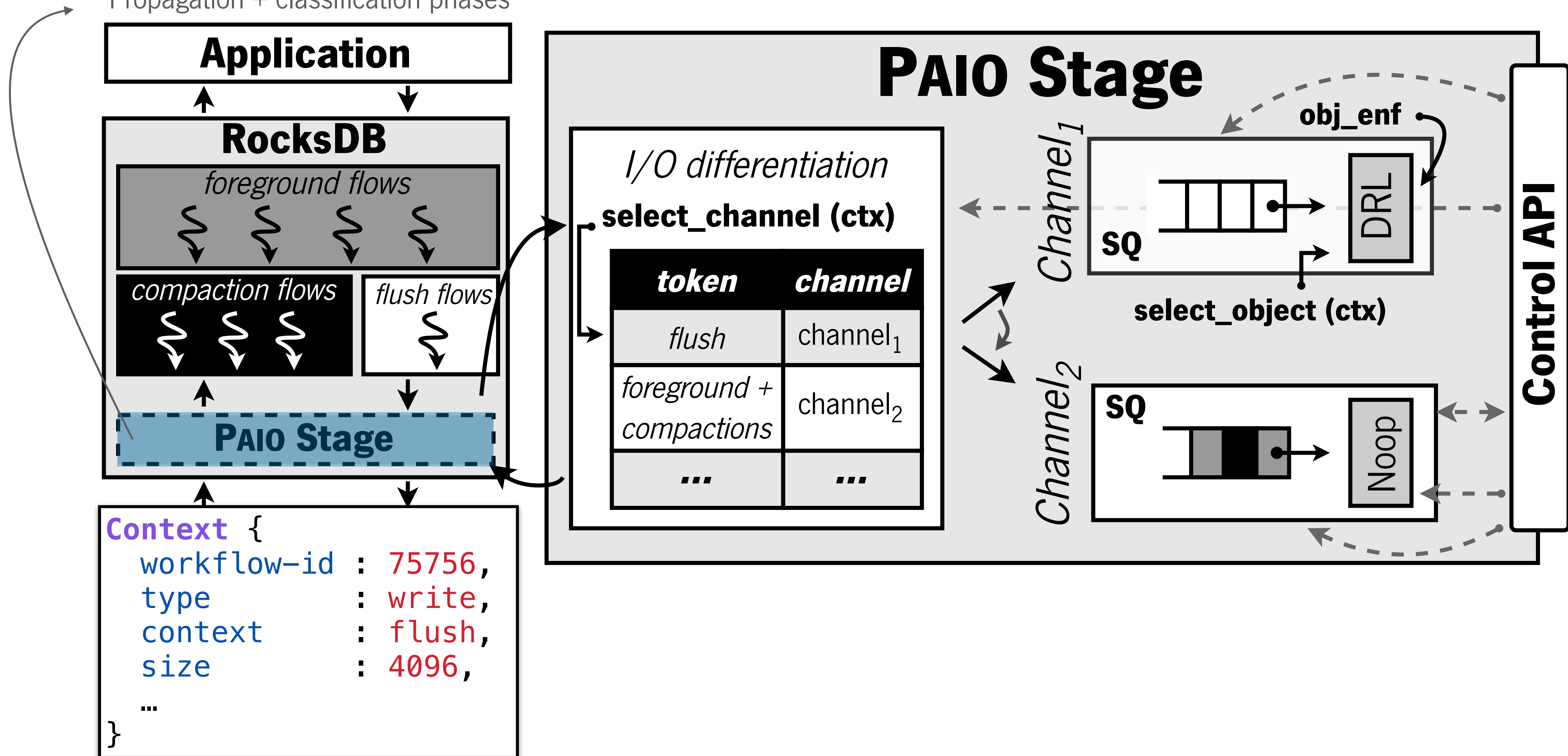**Context propagation:**
Instrumentation + propagation phases

Identify the origin of POSIX operations (i.e., **foreground**, **compaction**, or **flush** operations)
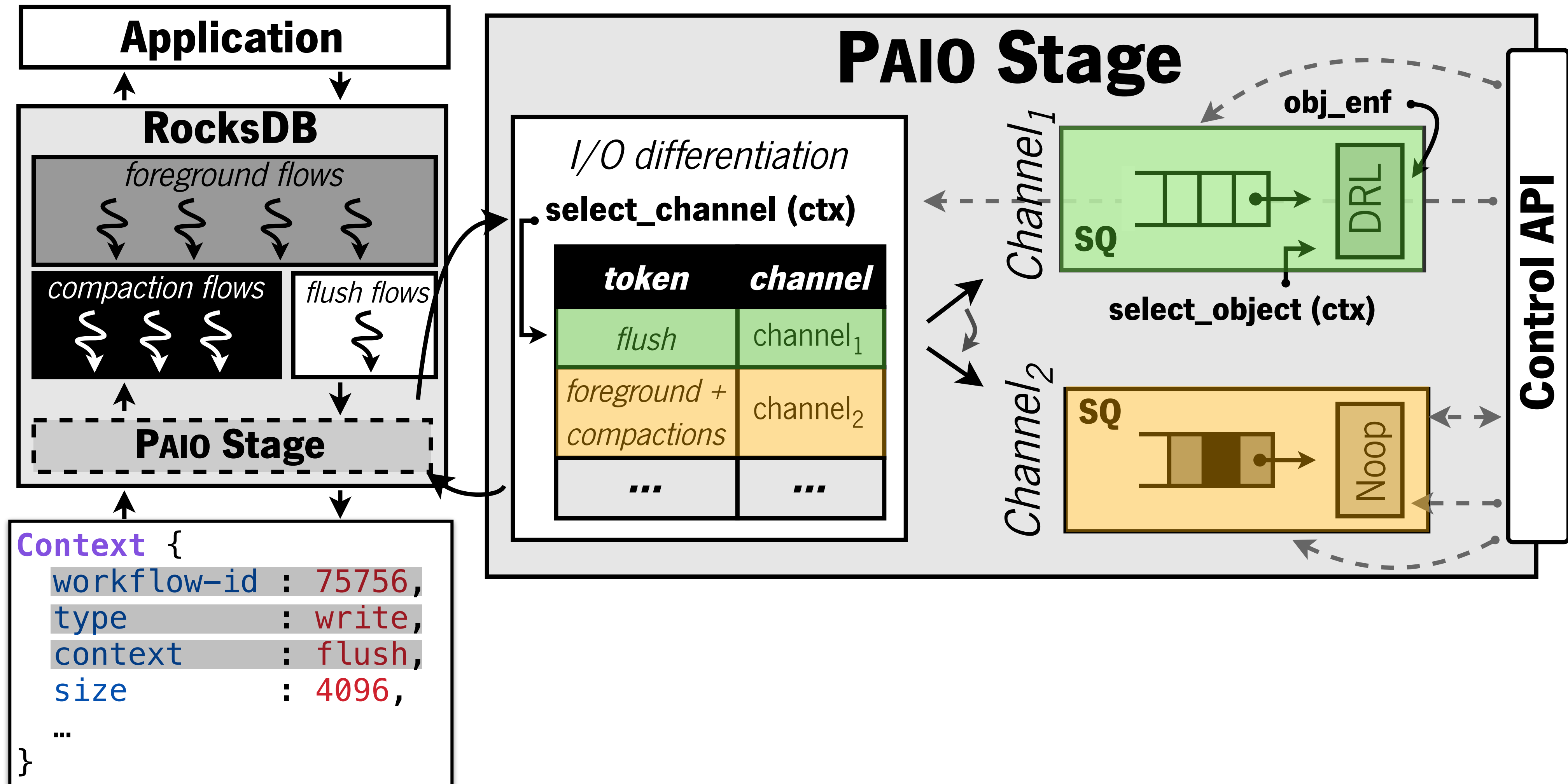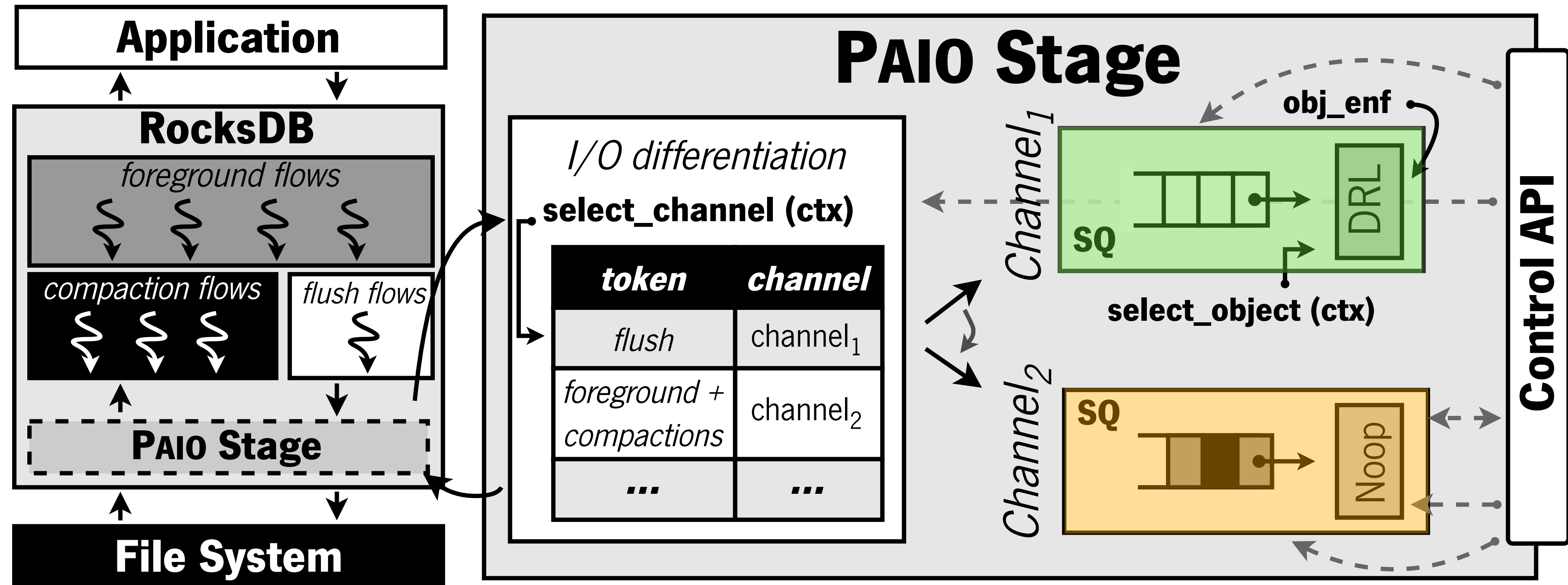
# I/O differentiation



**Context propagation:**
Propagation + classification phases

Application

RocksDB

*foreground flows*

*compaction flows*

*flush flows*

PAIO Stage

PAIO Stage

*I/O differentiation*

**select_channel (ctx)**

| token | channel |
|---|---|
| *flush* | channel$_1$ |
| *foreground + compactions* | channel$_2$ |
| **...** | **...** |

Channel$_1$

SQ

DRL

**obj_enf**

**select_object (ctx)**

Channel$_2$

SQ

Noop

Control API

```
Context {
  workflow-id : 75756,
  type        : write,
  context     : flush,
  size        : 4096,
  ...
}
```
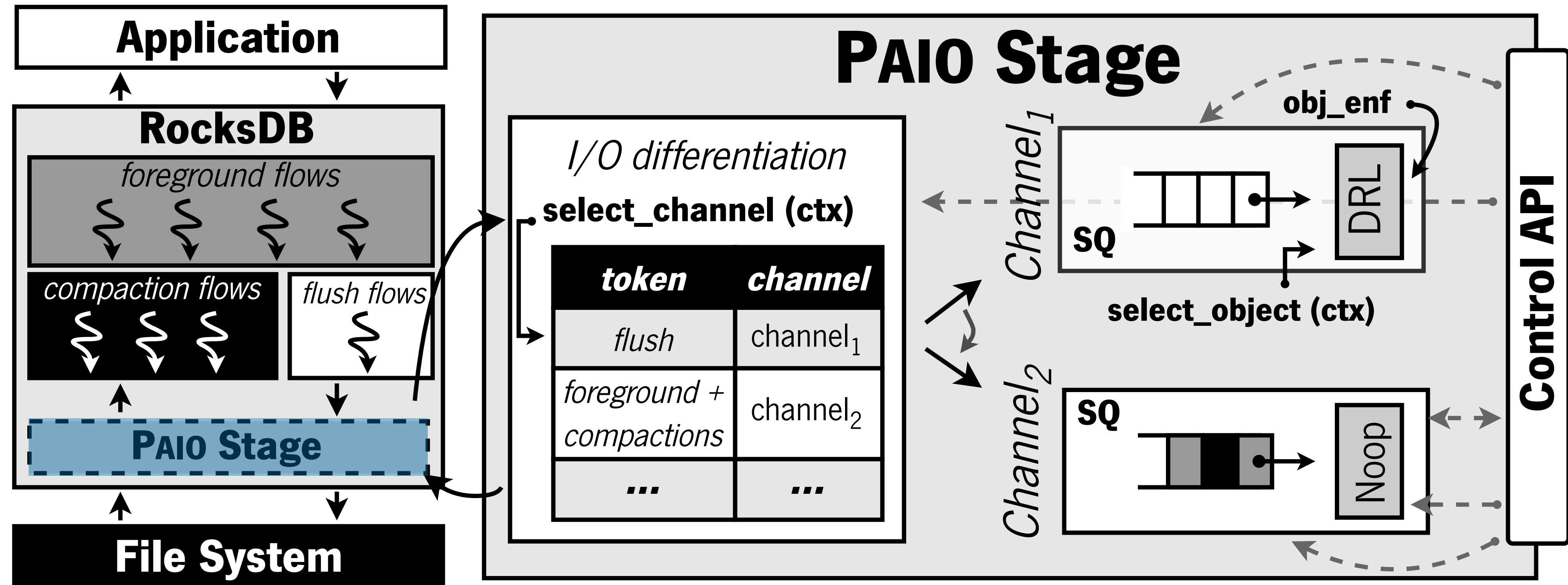
# I/O differentiation

# I/O enforcement



PAIO currently supports **Noop** (passthrough) and **DRL** (token-bucket) enforcement objects

# I/O enforcement



Requests return to their original I/O path

# Data plane stages built with PAIO

- **Per-application bandwidth control** under shared storage environments

  - Applied over multiple **TensorFlow** instances in the ABCI (AIST) supercomputer

- **Tail latency control** in Log-Structured Merge-tree key-value stores*

  - Applied over **RocksDB**, a production-ready key-value store from Meta

- **Metadata control** in Parallel File Systems*

  - Applied over **metadata-aggressive jobs** in Frontera (TACC) and ABCI supercomputers

**\* Discussed in this presentation.**

# Tail latency control in LSM-based key-value stores

**RocksDB**

- Interference between foreground and background tasks generates <u>high latency spikes</u>

- Latency spikes occur due to <u>$L_0$-$L_1$ compactions</u> and <u>flushes</u> being <u>slow or on hold</u>

**SILK**

- <u>I/O scheduler</u>

    - <u>Allocates bandwidth</u> for internal operations when client load is low

    - <u>Prioritizes</u> flushes and low level compactions

    - <u>Preempts</u> high level compactions with low level ones

- <u>Requires changing</u> several <u>core modules</u> made of thousands of LoC (≈335K LoC)

# Tail latency control in LSM-based key-value stores

**RocksDB**

- Interference between foreground and background tasks generates high latency spikes

- Latency spikes occur due to $L_0$-$L_1$ compactions and flushes being slow or on hold
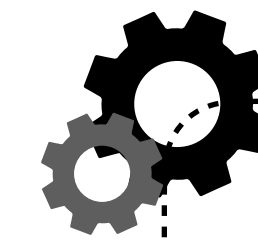
**SILK**

- I/O scheduler

  - **Allocates bandwidth for internal operations when client load is low**

  - **Prioritizes flushes and low level compactions**

  - Preempts high level compactions with low level ones

- Requires changing several core modules made of thousands of LoC (≈335K LoC)

**PAIO**

- Stage provides the I/O mechanisms for prioritizing and rate limiting background flows

  - Integrating PAIO in RocksDB only required adding 85 LoC

- Control plane provides a SILK-based I/O scheduling algorithm
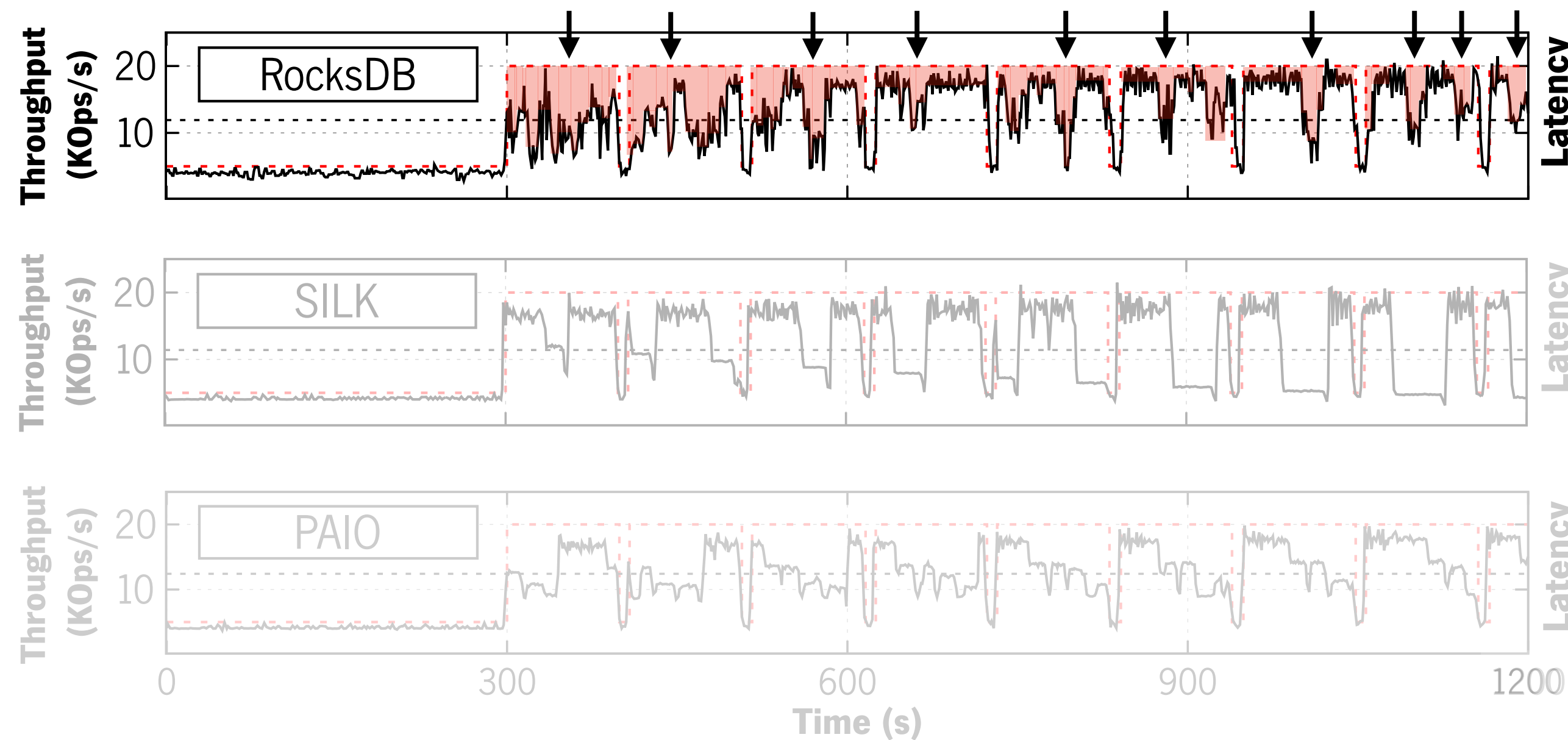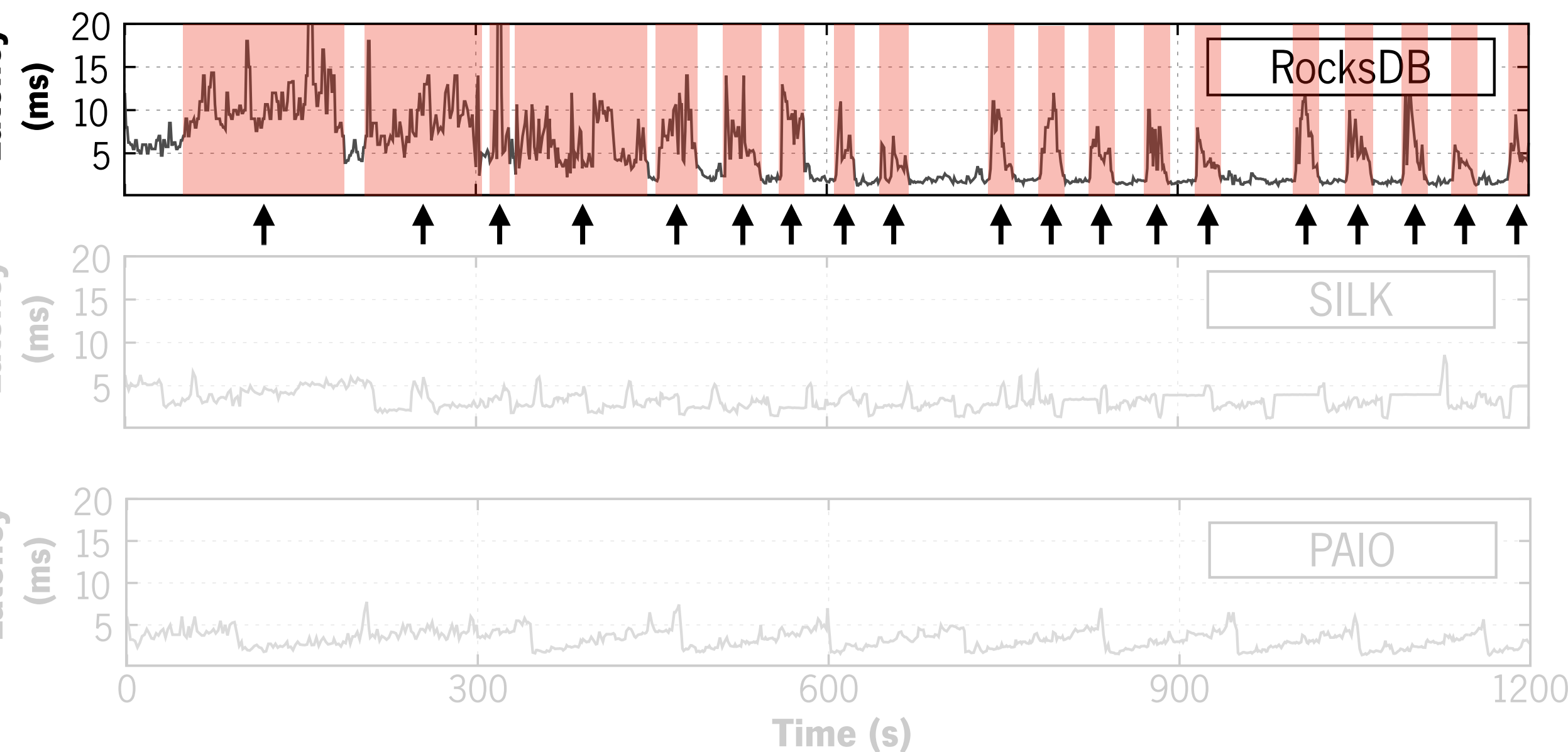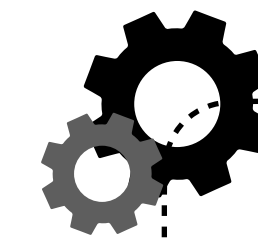
# Evaluation
## Mixture workload

**Throughput:** high variability due to constant flushes and compactions

**99th latency:** high tail latency with peaks with an average range between 3 and 15 ms

# Evaluation
## Mixture workload

**Throughput (higher is better)**
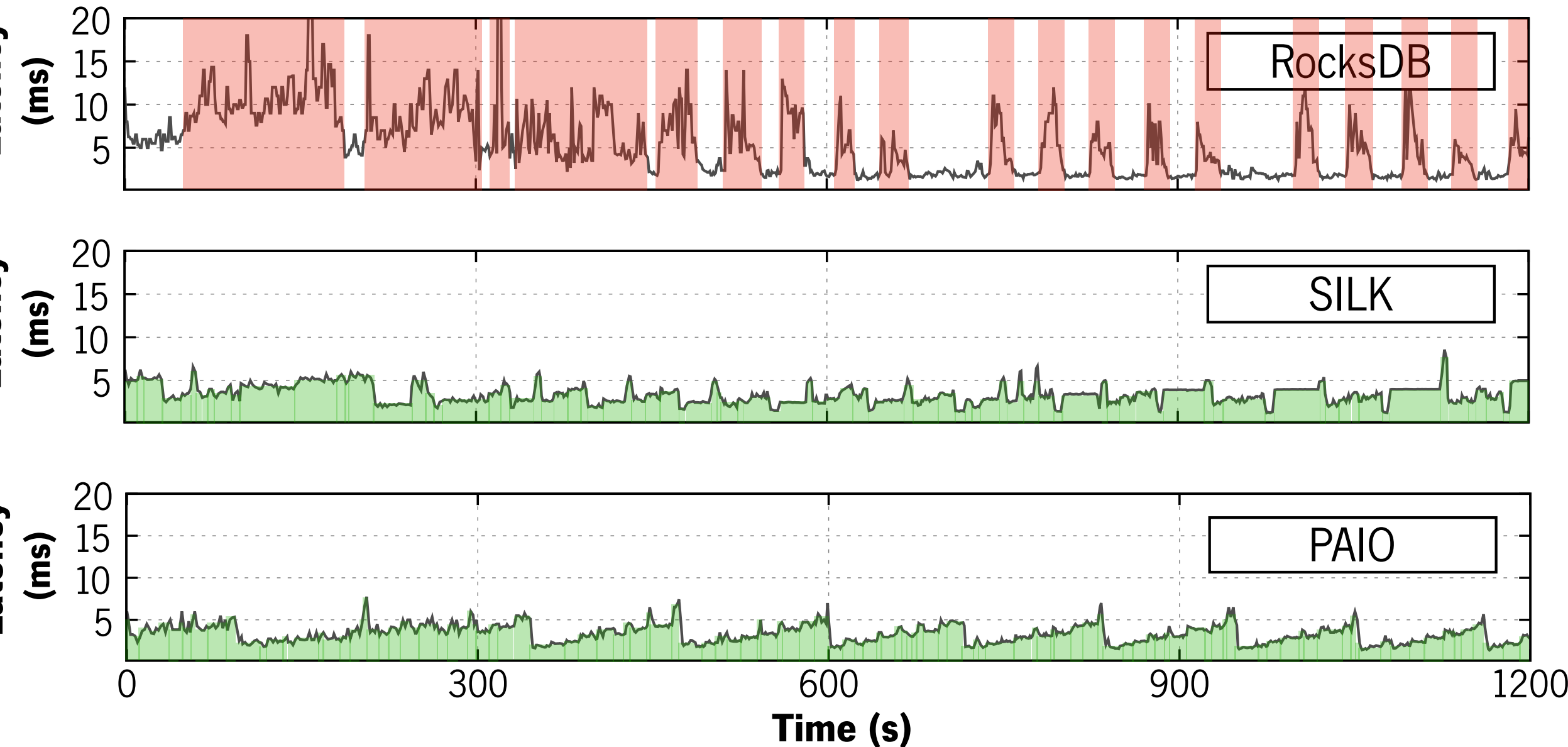
**Tail latency (lower is better)**



**Throughput:** suffers periodic throughput drops due to accumulated backlog

**99th latency:** low and sustained tail latency

## PAIO and SILK observe a 4x decrease in absolute tail latency

# Evaluation
## Mixture workload

Throughput (higher is better)

Tail latency (lower is better)



**By propagating application-level information to the stage, PAIO can enable similar control and performance as system-specific optimizations**

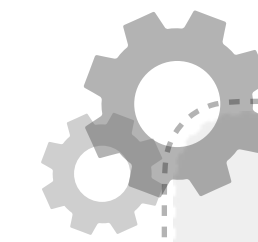**Throughput:** suffers periodic throughput drops due to accumulated backlog

**99th latency:** low and sustained tail latency

## PAIO and SILK observe a 4x decrease in absolute tail latency

# Metadata control in parallel file systems

- HPC workloads are no longer compute-bound and write-dominated

  - Modern workloads are **read-dominated** and with massive **bursts of metadata** operations

- Lustre-like parallel file systems (PFS) provide a **centralized metadata management** service

- Multiple jobs competing over **shared metadata resources**

  - Severe I/O contention

  - Overall performance degradation

# Metadata control in parallel file systems

- **PADLL**, an application and file system agnostic storage middleware that enables QoS of metadata workflows in HPC storage

- **Proactively** and **holistically controls** the rate at which POSIX requests are submitted to the PFS

- Data plane actuates at the **compute node level**

- Control plane follows a **hierarchical** organization

- New **max-min fair share algorithm** that prevents resource over-provisioning under volatile workloads

- PADLL does not require any code changes

**Macedo** *et al. "Protecting Metadata Servers From Harm Through Application-level I/O Control".* IEEE Cluster @ REX-IO, 2022.
**Macedo** *et al. Taming Metadata-intensive HPC Jobs Through Dynamic, Application-agnostic QoS Control".* 23rd IEEE/ACM CCGrid, 2023. In ***submission***.

# Metadata control in parallel file systems

- **PADLL**, an application and file system agnostic storage middleware that enables QoS of metadata workflows in HPC storage

- **Proactively** and **holistically controls** the rate at which POSIX requests are submitted to the PFS

- Data plane actuates at the **compute node level**

- Control plane follows a **hierarchical** organization

- New **max-min fair share algorithm** that prevents resource over-provisioning under volatile workloads

- PADLL does not require any code changes

**Macedo** *et al. "Protecting Metadata Servers From Harm Through Application-level I/O Control"*. IEEE Cluster @ REX-IO, 2022.
**Macedo** *et al. Taming Metadata-intensive HPC Jobs Through Dynamic, Application-agnostic QoS Control"*. 23rd IEEE/ACM CCGrid, 2023. In ***submission***.
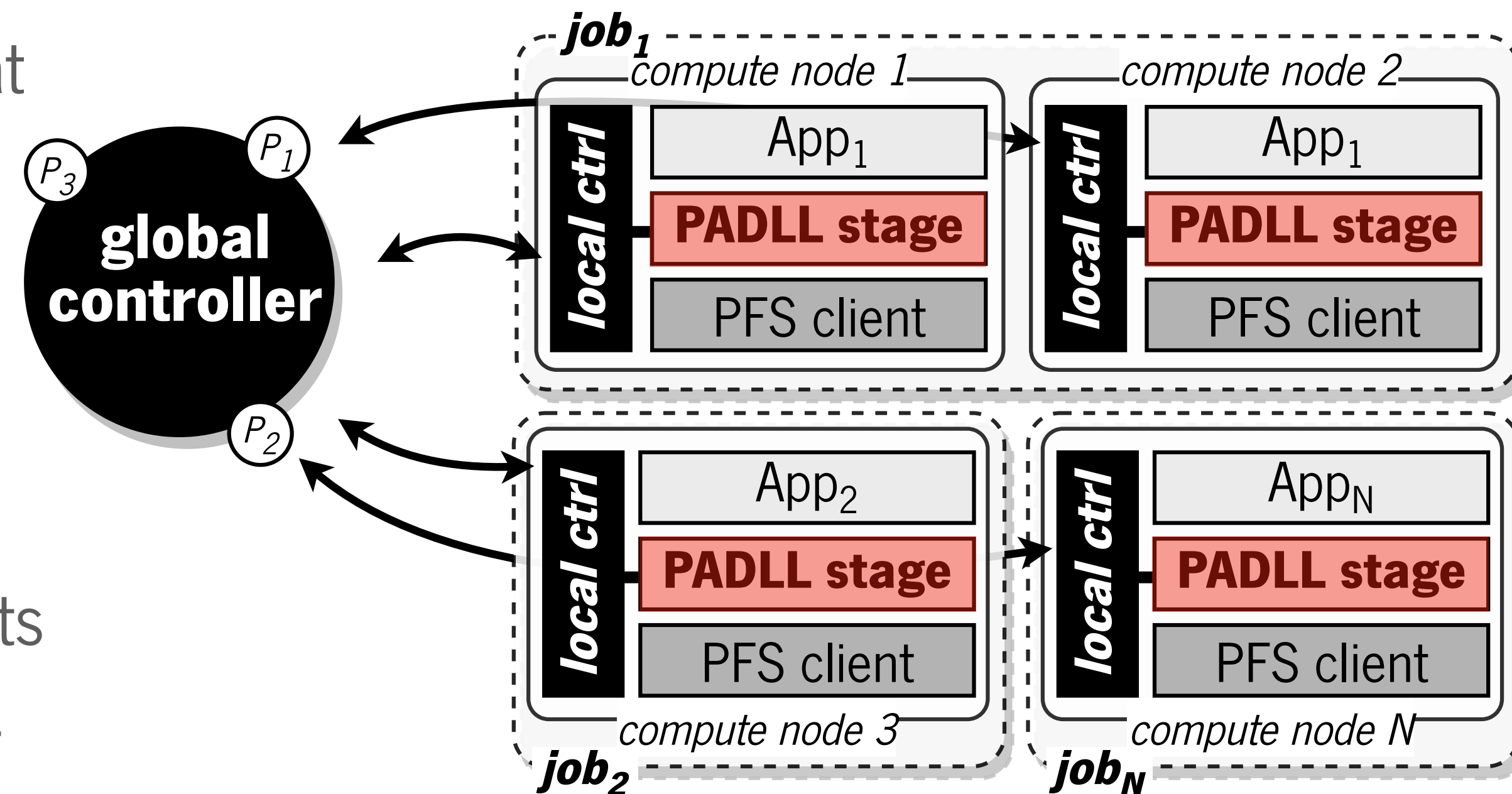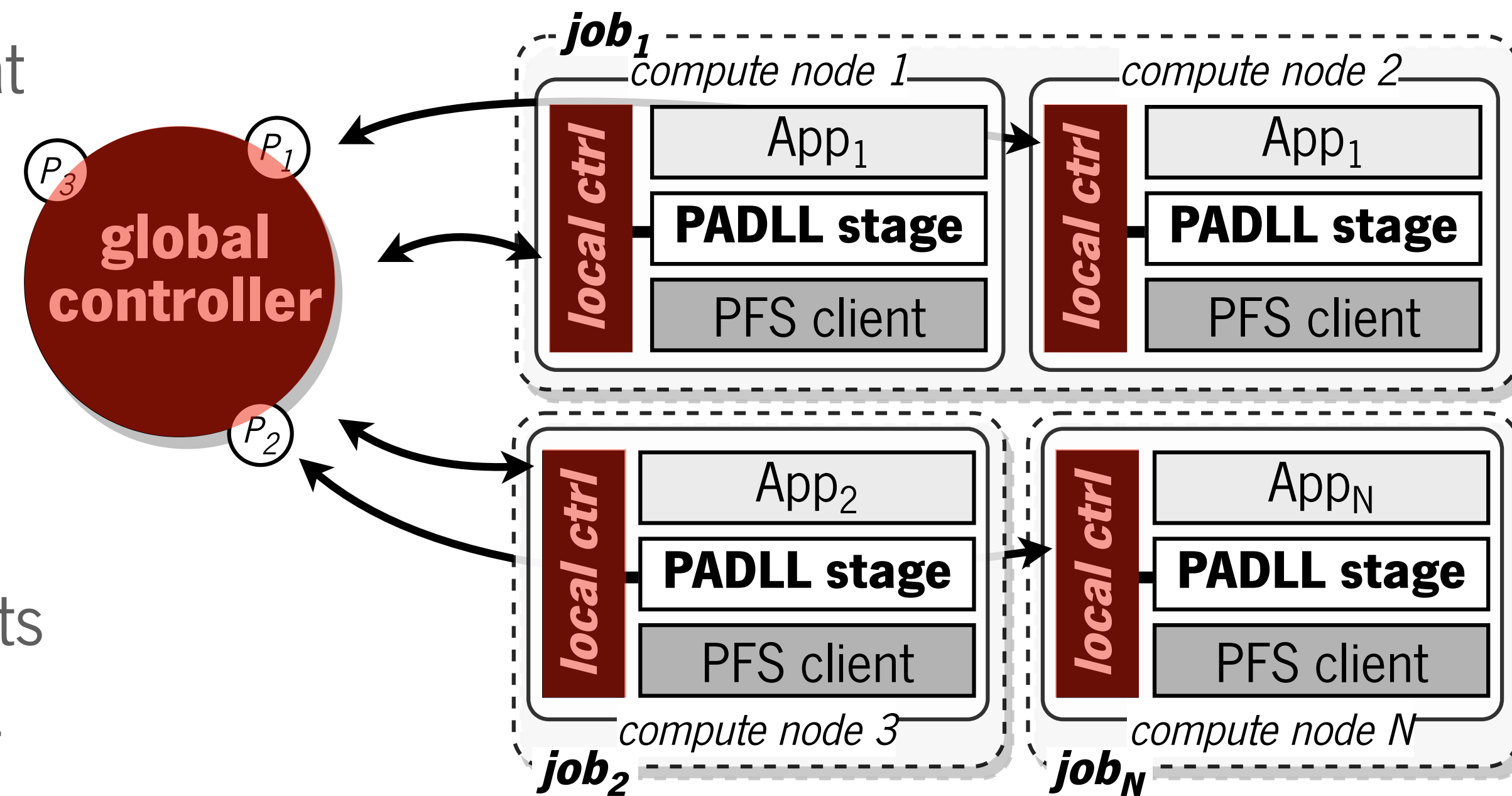
# Metadata control in parallel file systems

- **PADLL**, an application and file system agnostic storage middleware that enables QoS of metadata workflows in HPC storage

- **Proactively** and **holistically controls** the rate at which POSIX requests are submitted to the PFS

- Data plane actuates at the **compute node level**

- Control plane follows a **hierarchical** organization

- New **max-min fair share algorithm** that prevents resource over-provisioning under volatile workloads

- PADLL does not require any code changes

**Macedo** *et al. "Protecting Metadata Servers From Harm Through Application-level I/O Control"*. IEEE Cluster @ REX-IO, 2022.
**Macedo** *et al. Taming Metadata-intensive HPC Jobs Through Dynamic, Application-agnostic QoS Control"*. 23rd IEEE/ACM CCGrid, 2023. In ***submission***.

# Evaluation
## Metadata-aggressive jobs

- **Objective**
  - **Limit overall metadata load in the PFS, while assigning different I/O priorities to jobs**

- **Experimental environment**
  - Multi-job QoS control in the Frontera supercomputer
  - Trace replayer with metadata traces from the ABCI production cluster

- **Setups**
  - Baseline
  - Proportional Sharing (state-of-the-art QoS algorithm)
  - Proportional Sharing Without False Resource Allocation (new QoS algorithm)

# Evaluation
## Metadata-aggressive jobs

**System configuration and workload**
- Maximum metadata rate is set to **220 KOps/s**
- New job is added every 3 minutes
- Baseline execution time is 36 minutes (per job)
- Jobs execute with different loads {15%,20%,20%,45%}

**PFS QoS objective (maximum metadata rate)**

❌ Volatile and bursty workload
❌ Peaks reaching over 600 KOps/s

Throughput (KOps/s)

Time (minutes)

Job₁   Job₂   Job₃   Job₄

# Evaluation
## Metadata-aggressive jobs

**Proportional sharing:** enforce per-job metadata rate reservations, while assigning leftover rate when available

**Job₁ starts** **Job₂ starts** **Job₃ starts** **Job₄ starts** — All jobs are executing — **Job₁ ends** **Job₂ ends** **Job₃ ends** **Job₄ ends**

**Baseline**
300 225

**Prop. sharing**
240 180 120 60 0

Throughput (KOps/s)

✓ Maximum metadata limit is respected, <u>eliminating burstiness</u>
✓ Each reservation of metadata is respected
✓ Leftover rate is assigned to jobs whenever available
✗ Executes 5 minutes longer than Baseline

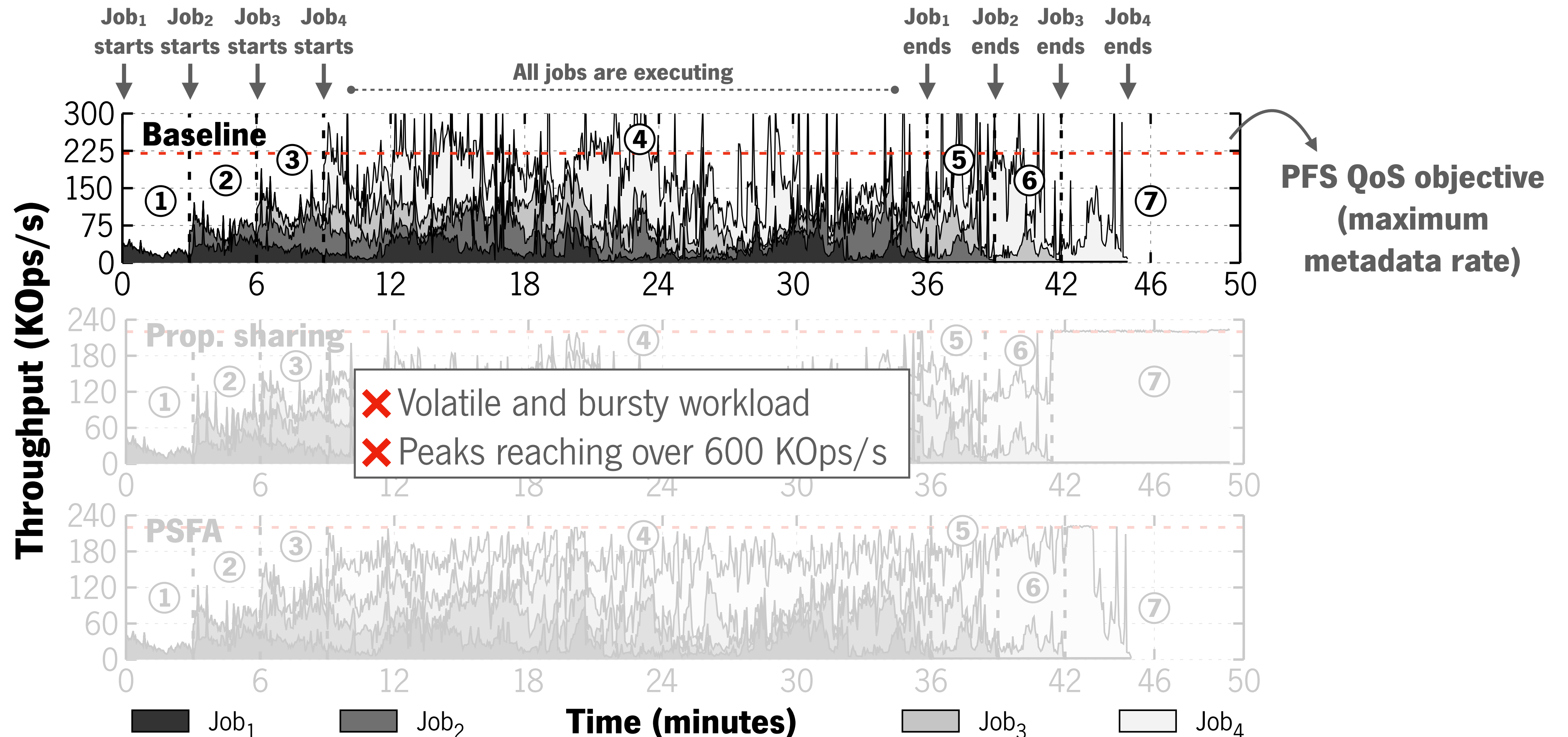Job₁　Job₂　**Time (minutes)**　Job₃　Job₄
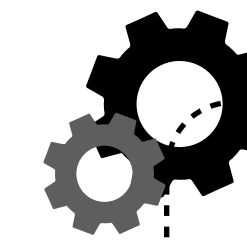
# Evaluation
## Metadata-aggressive jobs

**System configuration and workload**
- Maximum metadata rate is set to **220 KOps/s**
- New job is added every 3 minutes
- Baseline execution time is 36 minutes (per job)
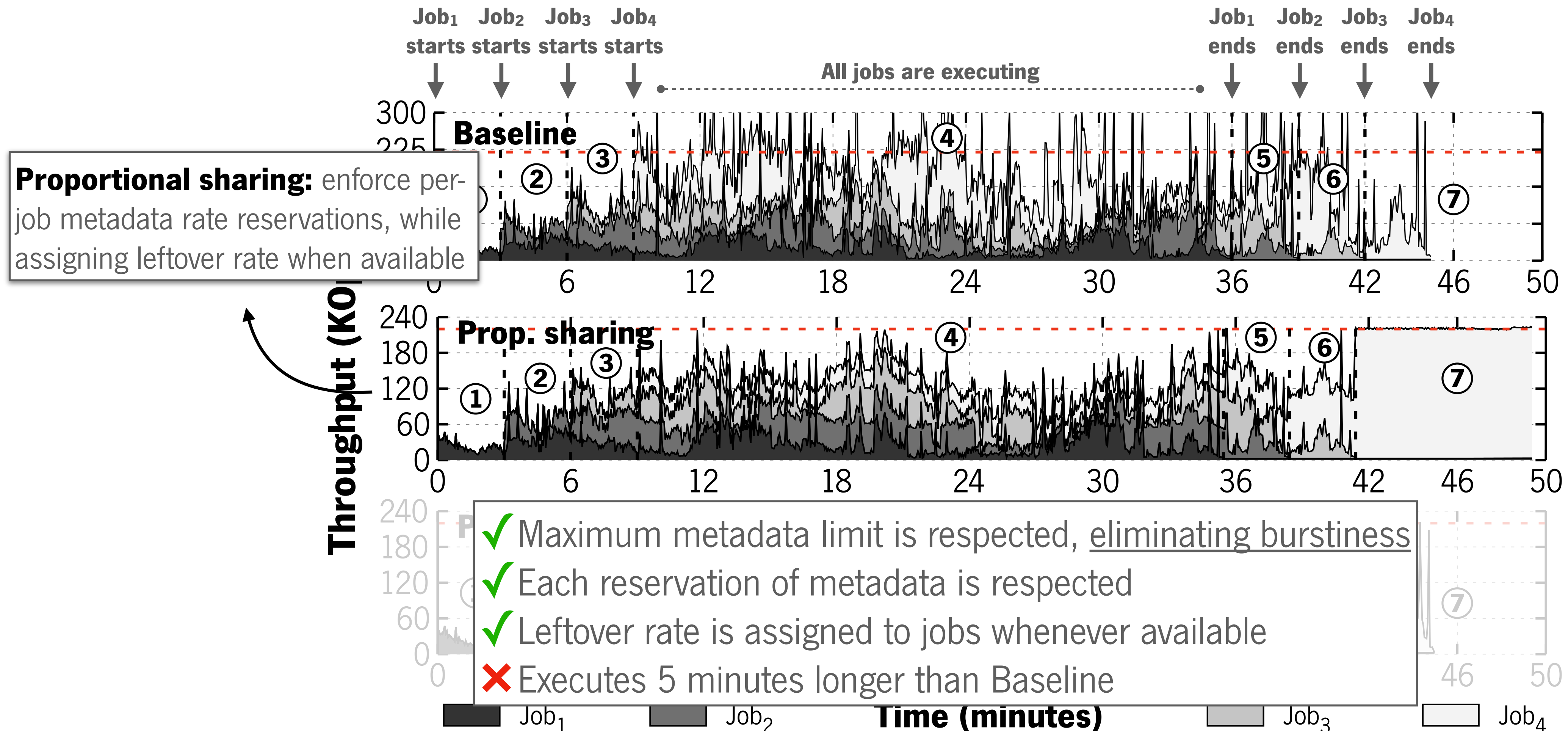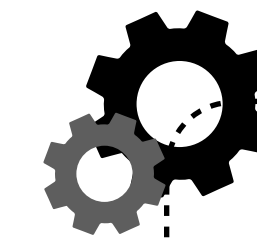- Jobs execute with different loads {15%,20%,20%,45%}

Job₁ starts  Job₂ starts  Job₃ starts  Job₄ starts

All jobs are executing

Job₁ ends  Job₂ ends  Job₃ ends  Job₄ ends

**Baseline**

**Proportional sharing:** enforce per-job metadata rate reservations, while assigning leftover rate when available

**Long periods of over-provisioning (resources assigned but not used)**

**Prop. sharing**

**Accumulated backlog!**

**Throughput (KOps/s)**

✓ Maximum metadata limit is respected, eliminating burstiness
✓ Each reservation of metadata is respected
✓ Leftover rate is assigned to jobs whenever available
✗ Executes 5 minutes longer than Baseline

Job₁    Job₂    **Time (minutes)**    Job₃    Job₄

# Evaluation
## Metadata-aggressive jobs

**Job₁ starts**  **Job₂ starts**  **Job₃ starts**  **Job₄ starts**

**Job₁ ends**  **Job₂ ends**  **Job₃ ends**  **Job₄ ends**

All jobs are executing

✓ Maximum metadata limit is respected, <u>eliminating burstiness</u>
✓ Each reservation of metadata is respected
✓ Unused I/O resources are reassigned, <u>preventing over-provisioning</u>
✓ All jobs finish under the same time as Baseline

**PSFA:** enforce per-job metadata rate reservations based on the actual I/O usage

Throughput (KOps/s)

Time (minutes)

Job₁    Job₂    Job₃    Job₄

# Summary

- **Survey** and **classification** of **SDS** systems

  - Systematization of knowledge and taxonomy of existing SDS solutions

  - Uncovers open research challenges in the field

- **PAIO**, a novel **SDS system** that enables building **complex I/O optimizations**

  - Decoupled from the targeted system

  - Perform coordinated control decisions over shared resources

  - Programmable and adaptable

- Data plane **stages** built with **PAIO**

  - Reimplement complex I/O optimizations that achieve similar performance as system-specific ones

  - New optimizations that address unsolved challenges present in modern I/O infrastructures

  - Currently working with leading HPC centers in the integration of PAIO and PADLL in production

# Publications

## Core publications

- **R. Macedo**, Y. Tanimura, J. Haga, V. Chidambaram, J. Pereira, J. Paulo. **"PAIO: General, Portable I/O Optimizations With Minor Application Modifications"**. *20th USENIX Conference on File and Storage Technologies*, 2022.

- **R. Macedo**, J. Paulo, J. Pereira, A. Bessani. **"A Survey and Classification of Software-Defined Storage Systems"**. *ACM Computing Surveys 53, 3 (48)*, 2020.

- **R. Macedo**, A. Faria, J. Paulo, J. Pereira. **"A Case for Dynamically Programmable Storage Background Tasks"**. *38th International Symposium on Reliable Distributed Systems Workshops*, 2019.

- **R. Macedo**, C. Correia, M. Dantas, C. Brito, W. Xu, Y. Tanimura, J. Haga, J. Paulo. **"The Case for Storage Optimization Decoupling in Deep Learning Frameworks"**. *IEEE Cluster @ REX-IO Workshop*, 2021.

- **R. Macedo**, M. Miranda, Y. Tanimura, J. Haga, A. Ruhela, S. Harrell, R. Evans, J. Paulo. **"Protecting Metadata Servers From Harm Through Application-level I/O Control"**. *IEEE Cluster @ REX-IO Workshop*, 2022.

- **R. Macedo**, M. Miranda, Y. Tanimura, J. Haga, A. Ruhela, S. Harrell, R. Evans, J. Pereira, J. Paulo. **"Taming Metadata-intensive HPC Jobs Through Dynamic, Application-agnostic QoS Control"**. *23rd IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing*, 2023. In ***submission***.

## Complementary publications

- M. Dantas, D. Leitão, P. Cui, **R. Macedo**, X. Liu, W. Xu, J. Paulo. **"Accelerating Deep Learning Training Through Transparent Storage Tiering"**. *22nd IEEE/ACM International Symposium on Cluster, Cloud and Internet Computing*, 2022.

- A. Faria, **R. Macedo**, J. Pereira, J. Paulo. **"BDUS: Implementing Block Devices in User Space"**. *14th ACM International System and Storage Conference*, 2021. ***Best paper runner-up.***

- M. Dantas, D. Leitão, C. Correia, **R. Macedo**, W. Xu, J. Paulo. **"Monarch: Hierarchical Storage Management for Deep Learning Frameworks"**. *IEEE Cluster @ REX-IO Workshop*, 2021.

- A. Faria, **R. Macedo**, J. Paulo. **"Pods-as-Volumes: Effortlessly Integrating Storage Systems and Middleware into Kubernetes"**. *ACM/IFIP Middleware @ WoC*, 2021.

- T. Esteves, **R. Macedo**, A. Faria, B. Portela, J. Paulo, J. Pereira, D. Harnik. **"TrustFS: An SGX-enabled Stackable File System Framework"**. *38th International Symposium on Reliable Distributed Systems Workshops*, 2019.