# IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage

**Huaicheng Li**☆*, Martin L. Putra☆, Ronald Shi☆,
Xing Lin*, Gregory R. Ganger*, Haryadi S. Gunawi ☆

*The 28th ACM Symposium on Operating Systems Principles (SOSP'21)*

☆*University of Chicago,* **Carnegie Mellon University,** *NetApp*
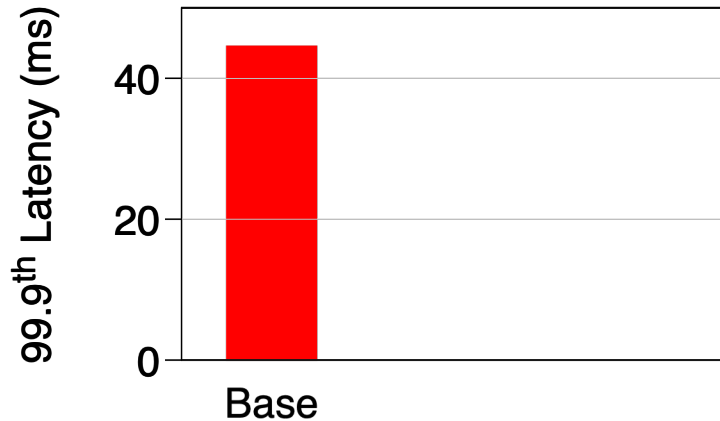
# IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage

# IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage
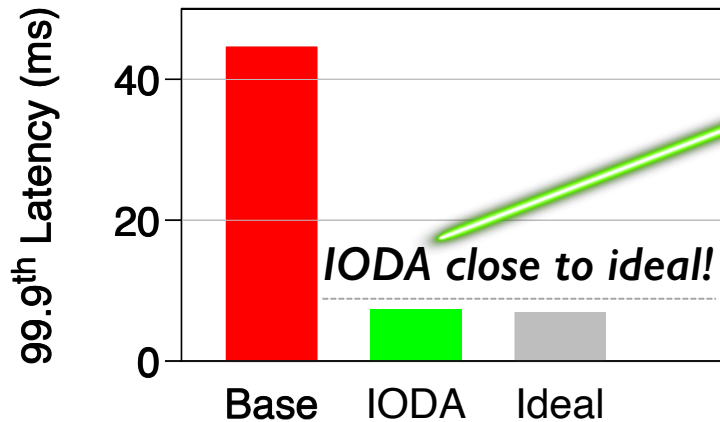


*"Small but powerful"*

# IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage
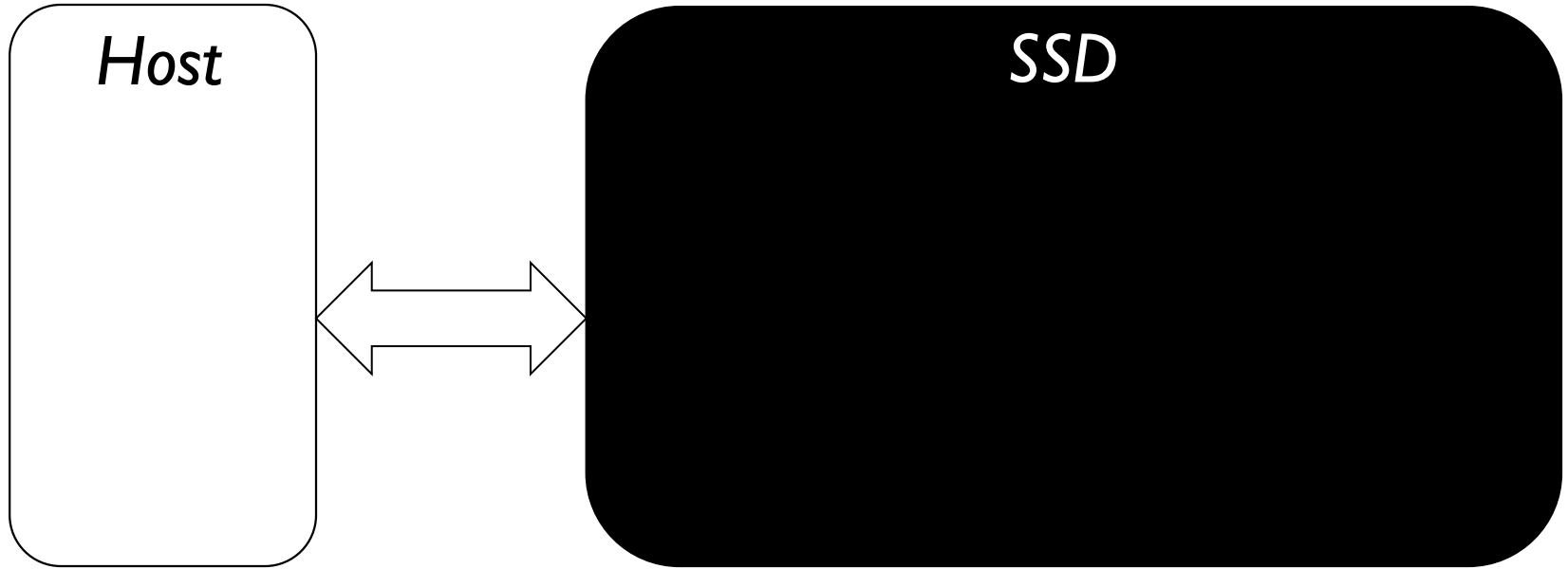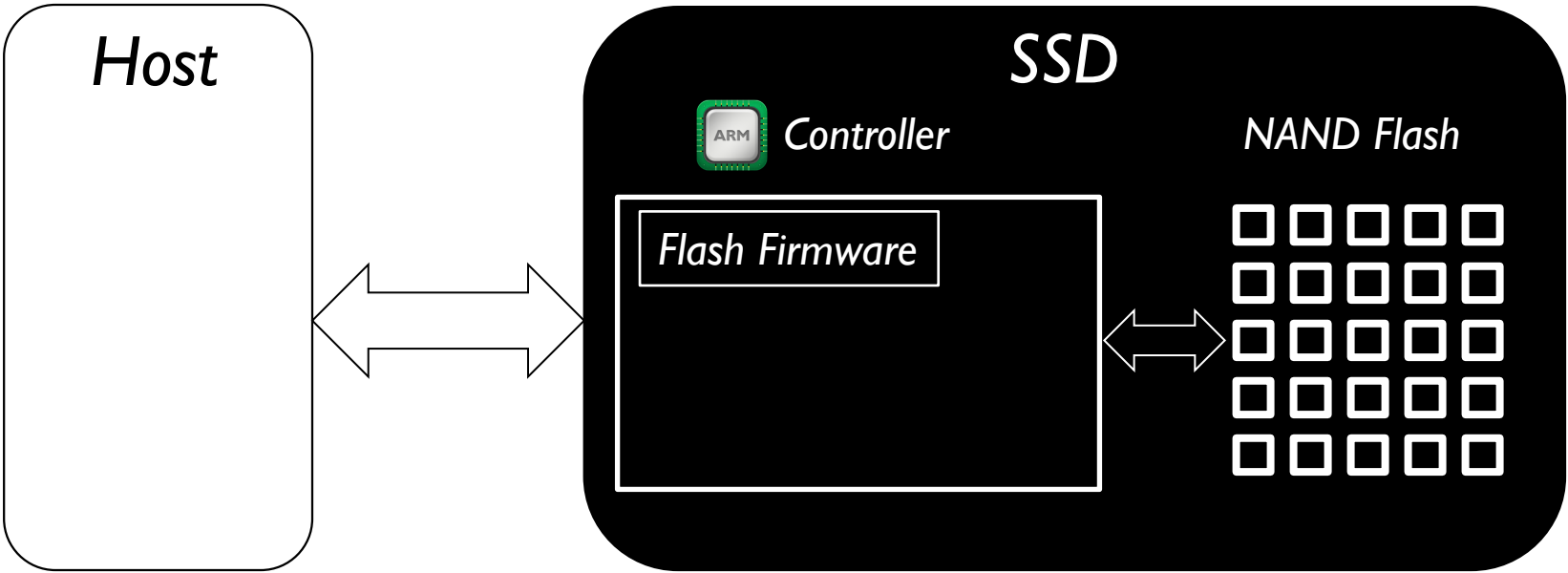


*"Small but powerful"*

# IODA: A Host/Device Co-Design for Strong Predictability Contract on Modern Flash Storage
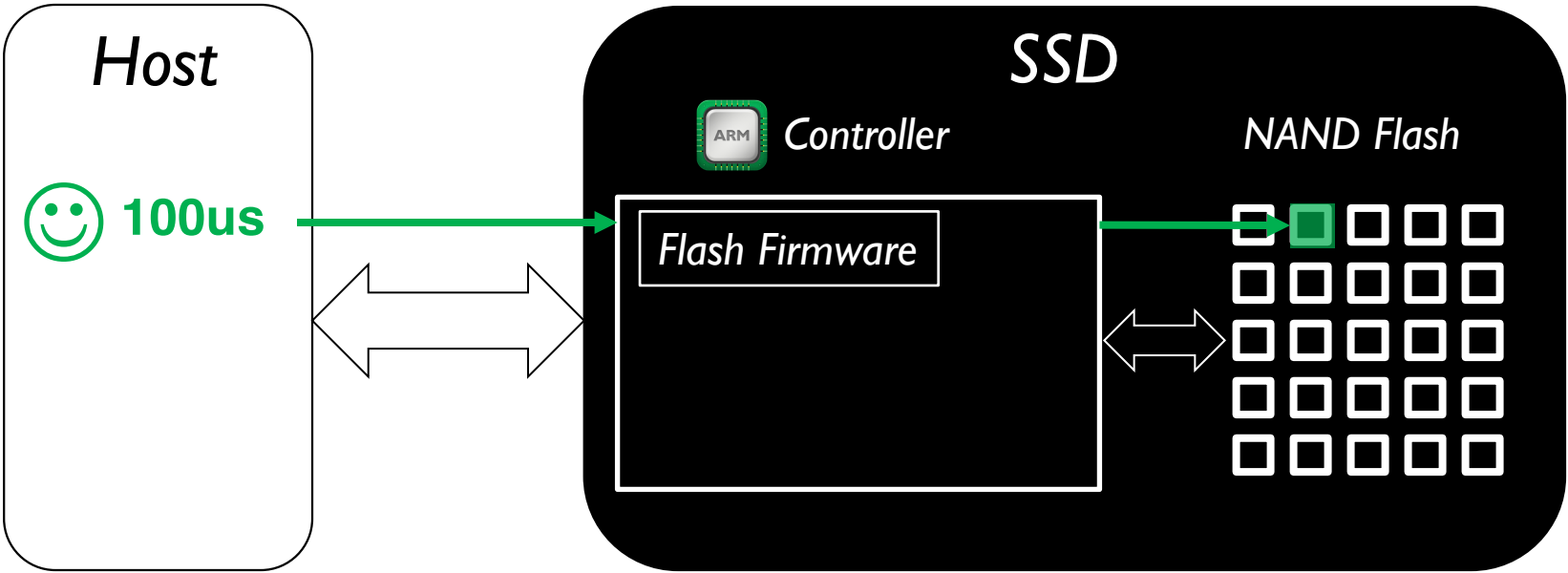


*"Small but powerful"*

**IODA close to ideal!**

99.9th Latency (ms)

40

20

0

Base  IODA  Ideal

# *"Attack of GC"* – Unpredictability in SSDs

**Host**

**SSD**

# *"Attack of GC"* – Unpredictability in SSDs

# "Attack of GC" – Unpredictability in SSDs

# *"Attack of GC"* – Unpredictability in SSDs

# *"Attack of GC"* – Unpredictability in SSDs

# *"Attack of GC"* – Unpredictability in SSDs
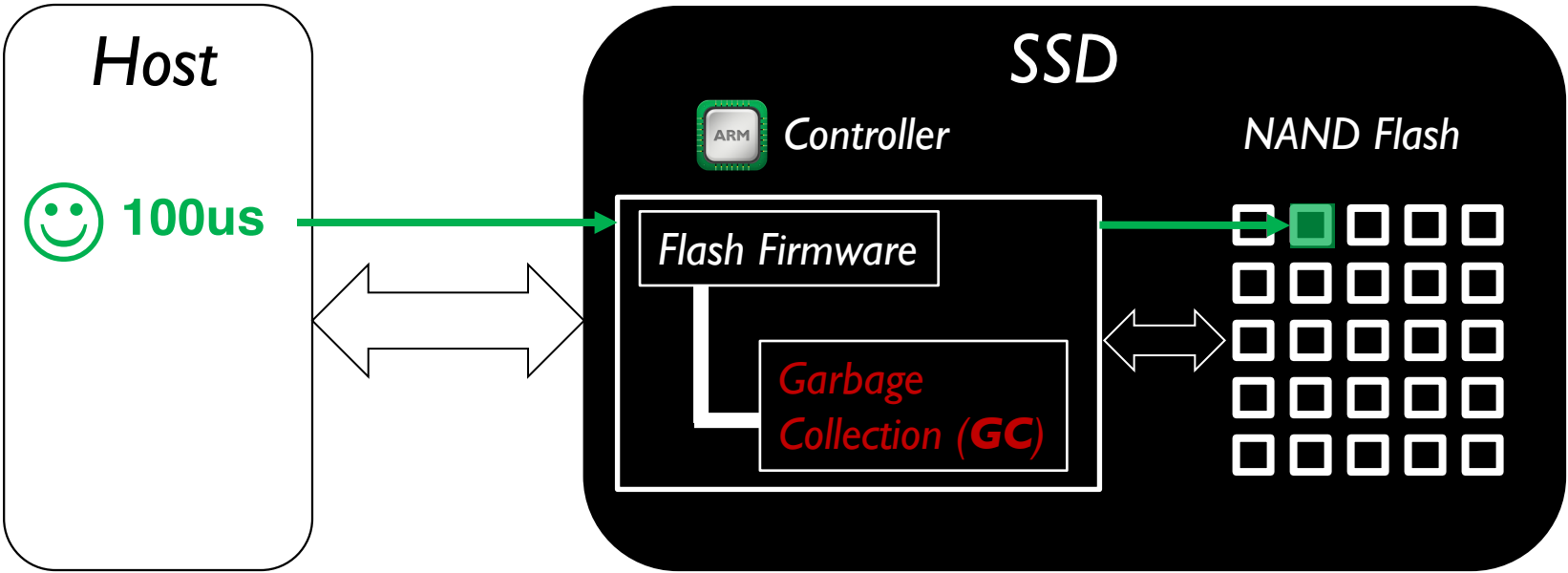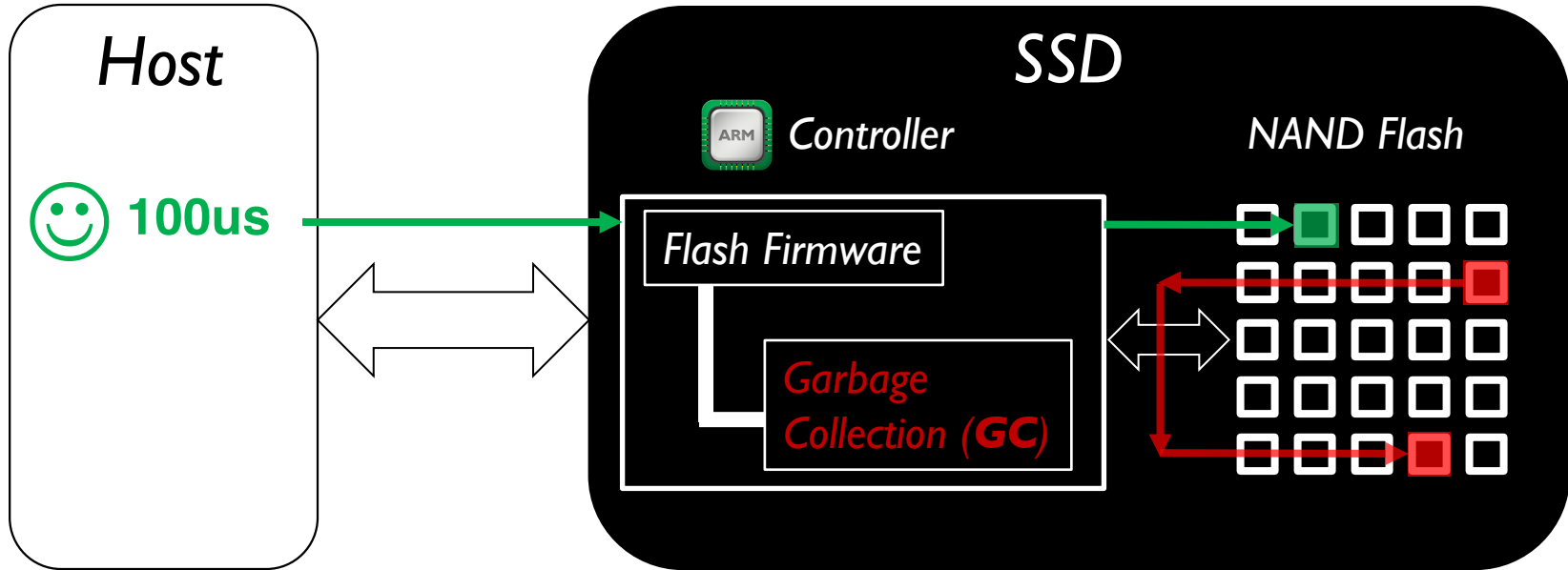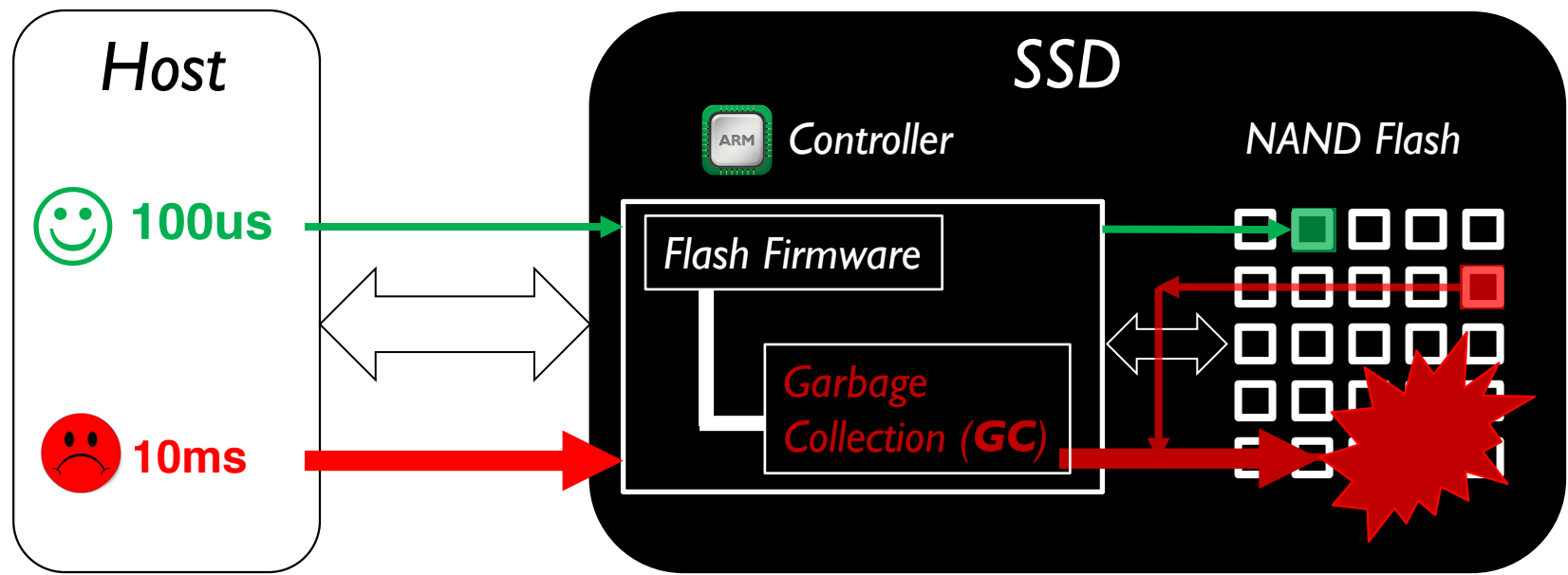
# *"Attack of GC"* – Unpredictability in SSDs
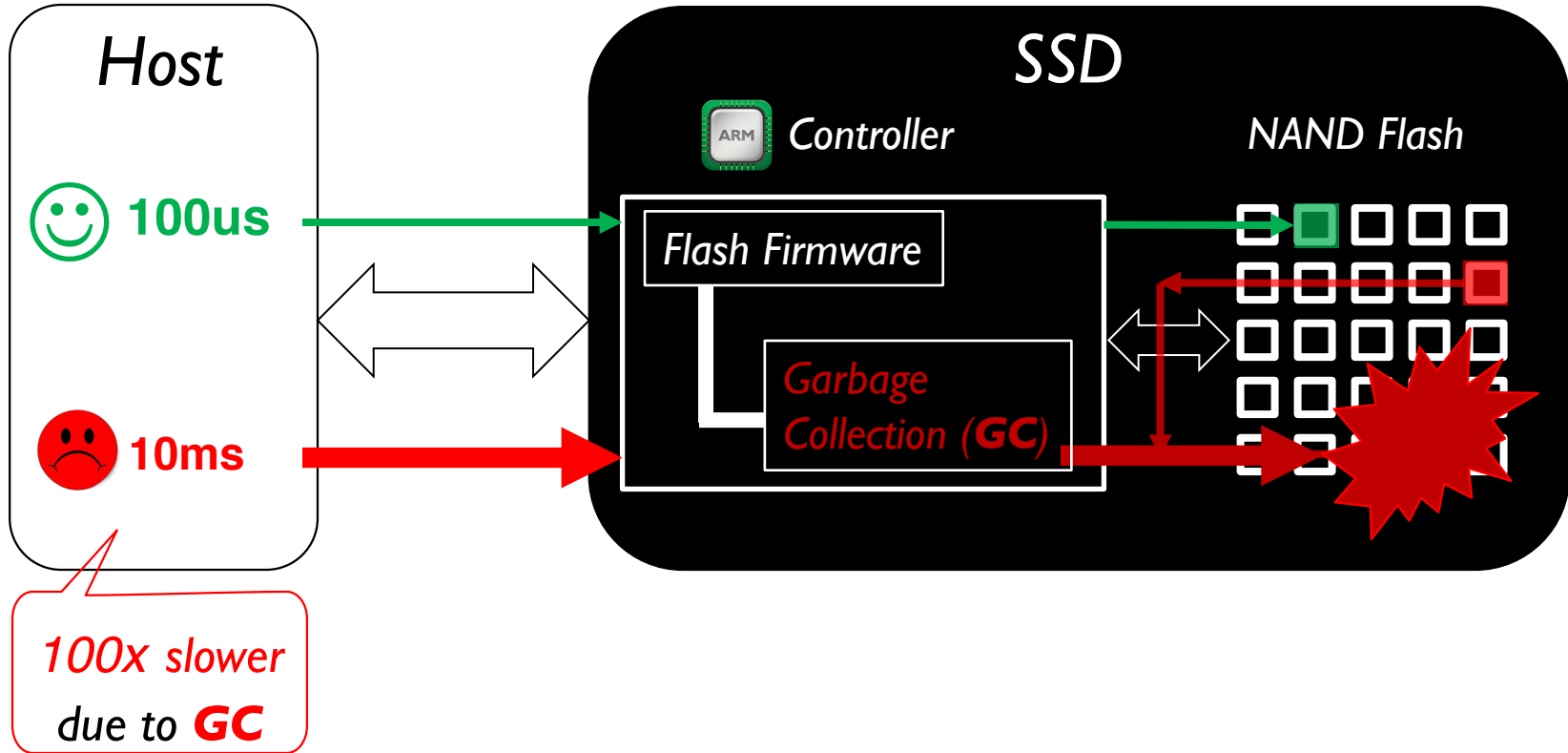
# *"Attack of GC"* – Unpredictability in SSDs



Host

100us

10ms

100x slower
due to **GC**

SSD

ARM  Controller

NAND Flash

Flash Firmware

Garbage
Collection (**GC**)

GC is ***Invisible*** to the Host

# *"The Tail Menace"* in Flash Arrays

Host

SSD0  SSD1  SSD2  SSD3

# *"The Tail Menace"* in Flash Arrays

Host

RAID

SSD0  SSD1  SSD2  SSD3

# *"The Tail Menace"* in Flash Arrays

Host

R

RAID

SSD0    SSD1    SSD2    SSD3

# *"The Tail Menace"* in Flash Arrays
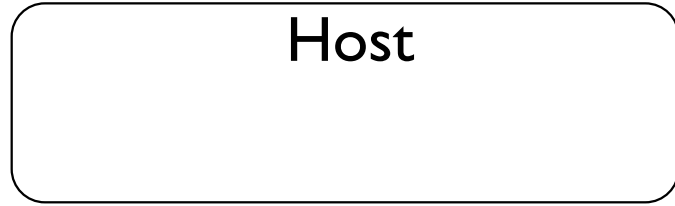
Host
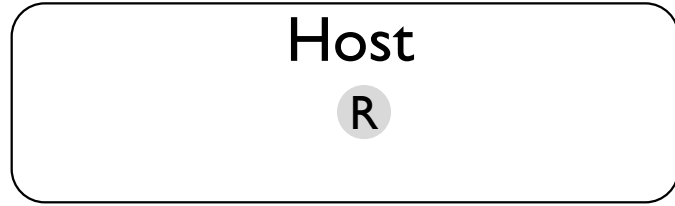
R

RAID

$R_0$    $R_1$    $R_3$
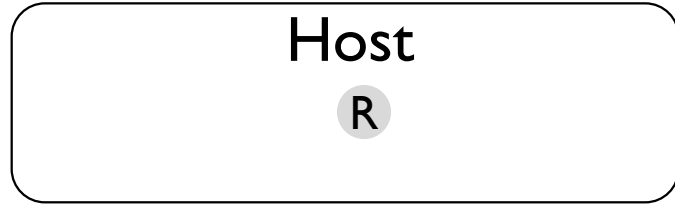
SSD0    SSD1    SSD2    SSD3

# *"The Tail Menace"* in Flash Arrays

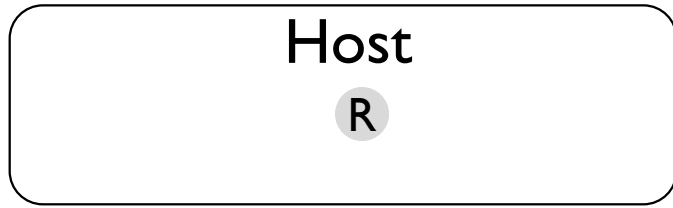# "The Tail Menace" in Flash Arrays

# *"The Tail Menace"* in Flash Arrays

# *"The Tail Menace"* in Flash Arrays

# *"The Tail Menace"* in Flash Arrays



RAID

Percentage

Better    Worse

Long Tails

100us          10ms

**A slow SSD makes the entire flash array slow!**

R₀    R₁    R₃

SSD0  SSD1  SSD2  SSD3

GC

# *"A New Hope"* – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

# "*A New Hope*" – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

*A major leap*

+ Predictable/Busy *Time Window (TW)*

+ Device status *query & toggling*

# "A New Hope" – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

A major leap

**+** Predictable/Busy *Time Window (TW)*

**+** Device status *query & toggling*

TW

Predictable | Busy | Predictable | Busy • • •

*Time*

# "A New Hope" – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

A major leap

- **+** Predictable/Busy *Time Window (TW)*
- **+** Device status *query & toggling*

*TW*  *Are you busy?*  *Go busy!*

Predictable | Busy | Predictable | Busy •••

*Time*

# *"A New Hope"* – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

**A major leap**

➕ Predictable/Busy *Time Window (TW)*

➕ Device status *query & toggling*

**But insufficient**

➖ *Coarse-grained* device-level predictability

➖ *"Soft-contract"* breaking predictability

➖ Requiring *complex* status tracking

➖ •••

*TW*   *Are you busy?*   *Go busy!*

Predictable | Busy | Predictable | Busy •••

*Time*

# *"A New Hope"* – NVMe Predictable Latency Mode

NVMe Predictable Latency Mode (**PLM**)

**A major leap**

**+** Predictable/Busy *Time Window (TW)*
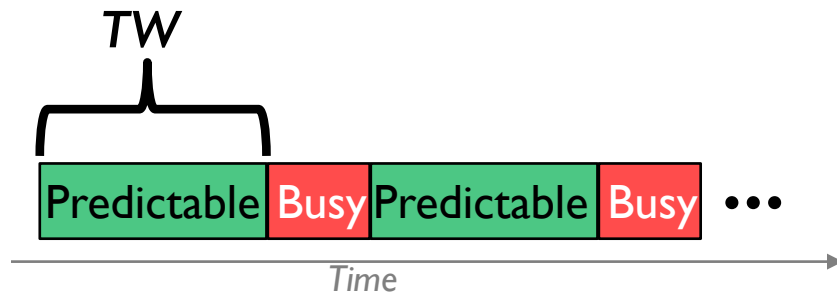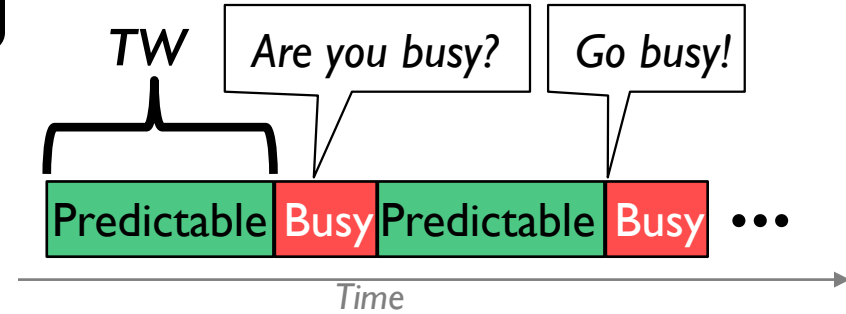
**+** Device status *query & toggling*

**But insufficient**

**−** *Coarse-grained* device-level predictability

**−** *"Soft-contract"* breaking predictability

**−** Requiring *complex* status tracking

**−** • • •

*TW*  *Are you busy?*  *Go busy!*

Predictable | Busy | Predictable | Busy  • • •

*Time*

*How to leverage NVMe PLM and enhance it for predictable latencies?*

# The IODA Story

❑ **Goal**: ***Tail-free*** flash array system on top of *slightly-extended* PLM interface

# The IODA Story

❑ **Goal**: ***Tail-free*** flash array system on top of *slightly-extended* PLM interface

❑ Design Principles:

- ▪ ***Simple*** policies for efficiency
- ▪ ***Minimal changes*** for easy deployment

# The IODA Story

❑ **Goal**: *Tail-free* flash array system on top of *slightly-extended* PLM interface

❑ Design Principles:
- ▪ ***Simple*** policies for efficiency
- ▪ ***Minimal changes*** for easy deployment

❑ IODA Approach/Techniques:

**NVMe** *PLM*

- ➖ *Coarse-grained*
- ➖ *"Soft-contract"*
- ➖ *Complex*

# The IODA Story

❑ **Goal**: *Tail-free* flash array system on top of *slightly-extended* PLM interface

❑ Design Principles:

▪ *Simple* policies for efficiency

▪ *Minimal changes* for easy deployment

❑ IODA Approach/Techniques:

✚ *Per-I/O* latency predictability

✚ Busy Remaining Time (*BRT*) Exposure

**NVMe** *PLM*

➖ *Coarse-grained*

➖ *"Soft-contract"*

➖ *Complex*

# The IODA Story

❑ **Goal**: *Tail-free* flash array system on top of *slightly-extended* PLM interface

❑ Design Principles:

  ▪ ***Simple*** policies for efficiency

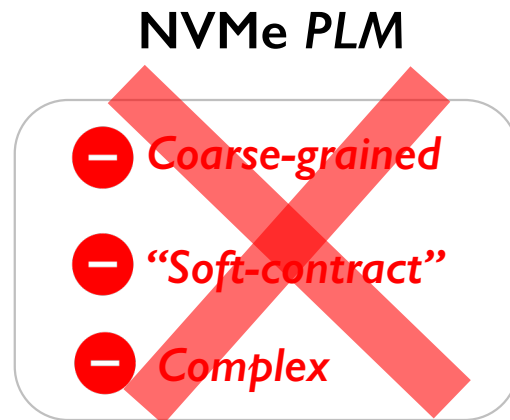  ▪ ***Minimal changes*** for easy deployment

❑ IODA Approach/Techniques:

  ➕ *Per-I/O* latency predictability

  ➕ Busy Remaining Time (*BRT*) Exposure

  ➕ **Time Window** (*TW*) Formulation

**NVMe** *PLM*

➖ *Coarse-grained*
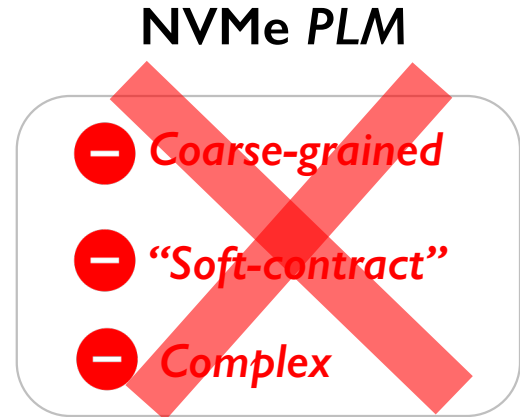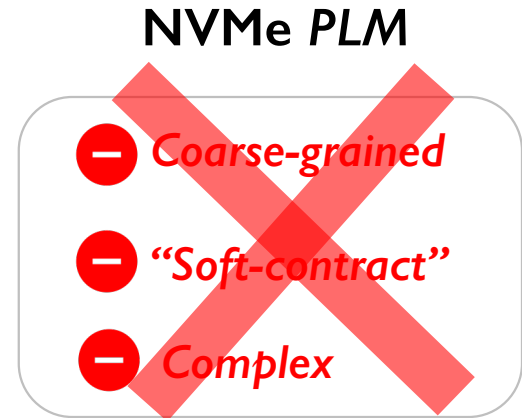
➖ *"Soft-contract"*

➖ *Complex*

# The IODA Story

❑ **Goal**: *Tail-free* flash array system on top of *slightly-extended* PLM interface

❑ Design Principles:
  ▪ *Simple* policies for efficiency
  ▪ *Minimal changes* for easy deployment

❑ IODA Approach/Techniques:
  ➕ *Per-I/O* latency predictability
  ➕ Busy Remaining Time (*BRT*) Exposure
  ➕ **Time Window** (*TW*) Formulation
  ➕ An end-to-end design exploiting above extensions

**NVMe** *PLM*

  ➖ *Coarse-grained*
  ➖ *"Soft-contract"*
  ➖ *Complex*

❑ Background & Motivation

❑ IODA Overview

❑ IODA Design
  ▪ Predictable latency flagged I/Os
  ▪ Busy remaining time
  ▪ Time window formulation
  ▪ Relaxed TW for better write amplification

❑ Evaluation

❑ Summary

# Leverage Redundancy for Performance

*An old, effective idea;*

# Leverage Redundancy for Performance

*An old, effective idea;*

**Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs**

Trimming the Tail for Deterministic Read Performance in SSDs

**Latency Reduction and Load Balancing in Coded Storage Systems**

**RAIL: Predictable, Low Tail Latency for NVMe Flash**

**MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface**

**EC-Cache: Load-balanced, Low-latency Cluster Caching with Online Erasure Coding**

K. V. Rashmi[1], Mosharaf Chowdhury[2], Jack Kosaian[2], Ion Stoica[1], Kannan Ramchandran[1]

[1]UC Berkeley   [2]University of Michigan

**Abstract**

Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any $k$ out of $(k + r)$ splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance

pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures.

Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-
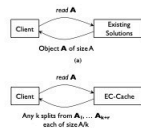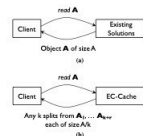
Figure 1: EC-Cache splits individual objects and encodes them using an erasure code to enable read parallelism and late binding during individual reads.

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

**When to issue the parity reads?**

**Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs**

**Trimming the Tail for Deterministic Read Performance in SSDs**

**Latency Reduction and Load Balancing in Coded Storage Systems**

**RAIL: Predictable, Low Tail Latency for NVMe Flash**

1

**MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface**

**EC-Cache: Load-balanced, Low-latency Cluster Caching with Online Erasure Coding**

K. V. Rashmi[1], Mosharaf Chowdhury[2], Jack Kosaian[2], Ion Stoica[1], Kannan Ramchandran[1]
[1]UC Berkeley   [2]University of Michigan

**Abstract**

Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any $k$ out of $(k + r)$ splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance
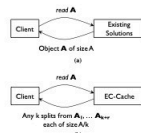
pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures.

Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-

**Figure 1:** EC-Cache splits individual objects and encodes them using an erasure code to enable read parallelism and late binding during individual reads.

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

**Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs**

Trimming the Tail for Deterministic Read Performance in SSDs

**Latency Reduction and Load Balancing in Coded Storage Systems**

**RAIL: Predictable, Low Tail Latency for NVMe Flash**

**MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface**

**EC-Cache: Load-balanced, Low-latency Cluster Caching with Online Erasure Coding**

K. V. Rashmi[1], Mosharaf Chowdhury[2], Jack Kosaian[2], Ion Stoica[1], Kannan Ramchandran[1]
[1]UC Berkeley   [2]University of Michigan

**Abstract**
Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any $k$ out of $(k + r)$ splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance
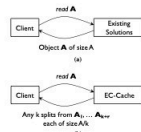
pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures.

Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-

**When to issue the parity reads?**

(1) *Wait* for timeout

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

**When to issue the parity reads?**

(1) *Wait* for timeout

➖ Best threshold? *Tricky*



**Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs**

**Trimming the Tail for Deterministic Read Performance in SSDs**

**Latency Reduction and Load Balancing in Coded Storage Systems**

**RAIL: Predictable, Low Tail Latency for NVMe Flash**

**MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface**

**EC-Cache: Load-balanced, Low-latency Cluster Caching with Online Erasure Coding**

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

**When to issue the parity reads?**

(1) *Wait* for timeout

⊖ Best threshold? *Tricky*

(2) *Always* Proactive (always send full-stripe)

**Tiny-Tail Flash: Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs**

**Trimming the Tail for Deterministic Read Performance in SSDs**

**Latency Reduction and Load Balancing in Coded Storage Systems**

**RAIL: Predictable, Low Tail Latency for NVMe Flash**

**MittOS: Supporting Millisecond Tail Tolerance with Fast Rejecting SLO-Aware OS Interface**

**EC-Cache: Load-balanced, Low-latency Cluster Caching with Online Erasure Coding**

K. V. Rashmi[1], Mosharaf Chowdhury[2], Jack Kosaian[2], Ion Stoica[1], Kannan Ramchandran[1]
[1]*UC Berkeley* [2]*University of Michigan*

**Abstract**

Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any $k$ out of $(k + r)$ splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance.
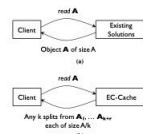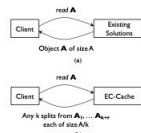
pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures.

Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

Tiny-Tail Flash: Near-Perfect Elimination of
Garbage Collection Tail Latencies in NAND SSDs

Trimming the Tail for Deterministic Read
Performance in SSDs

Latency Reduction and Load Balancing in Coded
Storage Systems

RAIL: Predictable, Low Tail Latency for NVMe
Flash

MittOS: Supporting Millisecond Tail Tolerance with
Fast Rejecting SLO-Aware OS Interface

EC-Cache: Load-balanced, Low-latency Cluster Caching
with Online Erasure Coding

K. V. Rashmi[1], Mosharaf Chowdhury[2], Jack Kosaian[2], Ion Stoica[1], Kannan Ramchandran[1]
[1]UC Berkeley   [2]University of Michigan

**Abstract**

Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any k out of (k + r) splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance
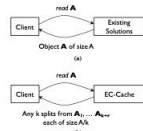
Figure 1: EC-Cache splits individual objects and encodes them using an erasure code to enable read parallelism and late binding during individual reads.

pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures. Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-

## When to issue the parity reads?

(1) *Wait* for timeout
  ⊖ Best threshold? *Tricky*

(2) *Always* Proactive (always send full-stripe)
  ⊖ Increased load ➔ *Inefficient*

# Leverage Redundancy for Performance

**An old, effective idea;** **Yet, challenging for PLM**

**When to issue the parity reads?**

(1) *Wait* for timeout
  − Best threshold? *Tricky*

(2) *Always* Proactive (always send full-stripe)
  − Increased load ⟶ *Inefficient*

*Semantic gap* between the Host and SSD to communicate the *"busyness"*

Tiny-Tail Flash: Near-Perfect Elimination of
Garbage Collection Tail Latencies in NAND SSDs

Trimming the Tail for Deterministic Read
Performance in SSDs

Latency Reduction and Load Balancing in Coded
Storage Systems

RAIL: Predictable, Low Tail Latency for NVMe
Flash

MittOS: Supporting Millisecond Tail Tolerance with
Fast Rejecting SLO-Aware OS Interface

**Abstract**

Data-intensive clusters and object stores are increasingly relying on in-memory object caching to meet the I/O performance demands. These systems routinely face the challenges of popularity skew, background load imbalance, and server failures, which result in severe load imbalance across servers and degraded I/O performance. Selective replication is a commonly used technique to tackle these challenges, where the number of cached replicas of an object is proportional to its popularity. In this paper, we explore an alternative approach using erasure coding.

EC-Cache is a load-balanced, low latency cluster cache that uses online erasure coding to overcome the limitations of selective replication. EC-Cache employs erasure coding by: (i) splitting and erasure coding individual objects during writes, and (ii) late binding, wherein obtaining any $k$ out of $(k + r)$ splits of an object are sufficient, during reads. As compared to selective replication, EC-Cache improves load balancing by more than 3× and reduces the median and tail read latencies by more than 2×, while using the same amount of memory. EC-Cache does so using 10% additional bandwidth and a small increase in the amount of stored metadata. The benefits offered by EC-Cache are further amplified in the presence of background network load imbalance

pling [12, 16, 52] and compression [15, 27, 53, 79] are some of the popular approaches employed to increase the effective memory capacity. (iii) Ensuring good I/O performance for the cached data in the presence of skewed popularity, background load imbalance, and failures.

Typically, the popularity of objects in cluster caches are heavily skewed [20, 47], and this creates signifi-

Figure 1: EC-Cache splits individual objects and encodes them using an erasure code to enable read parallelism and late binding during individual reads.

# IOD$_1$: Predictable Latency Flagged I/Os

💡 *"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

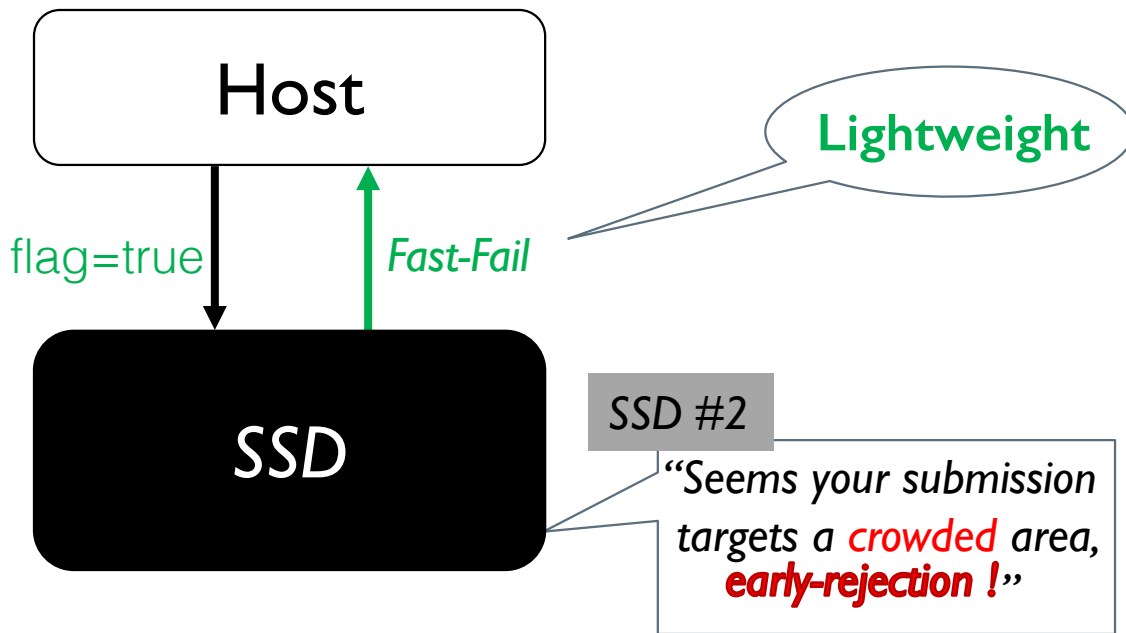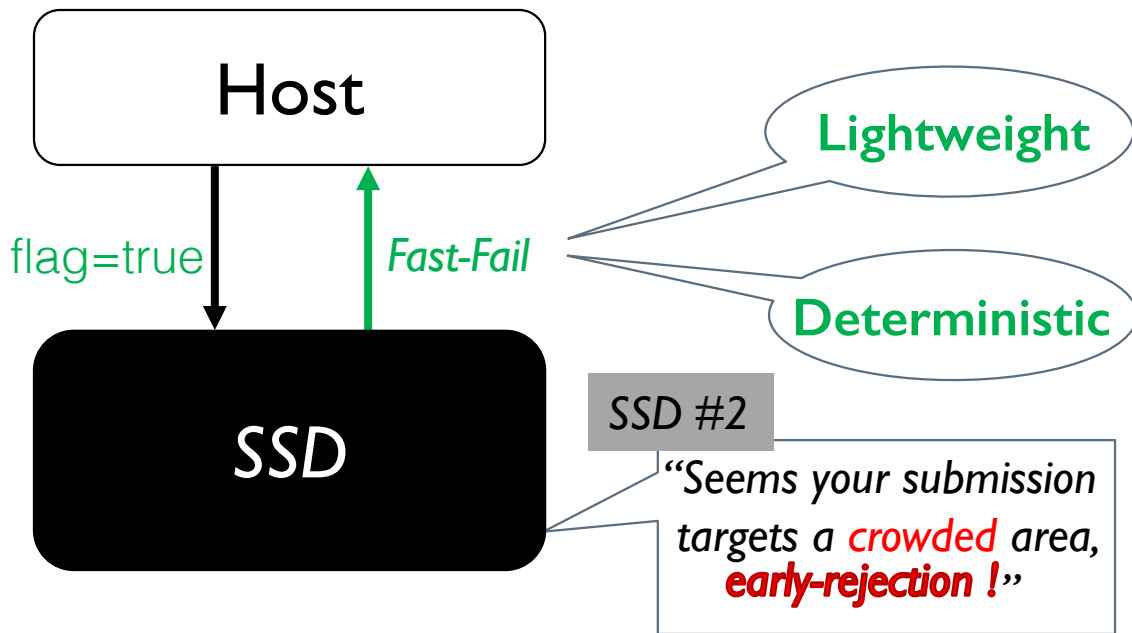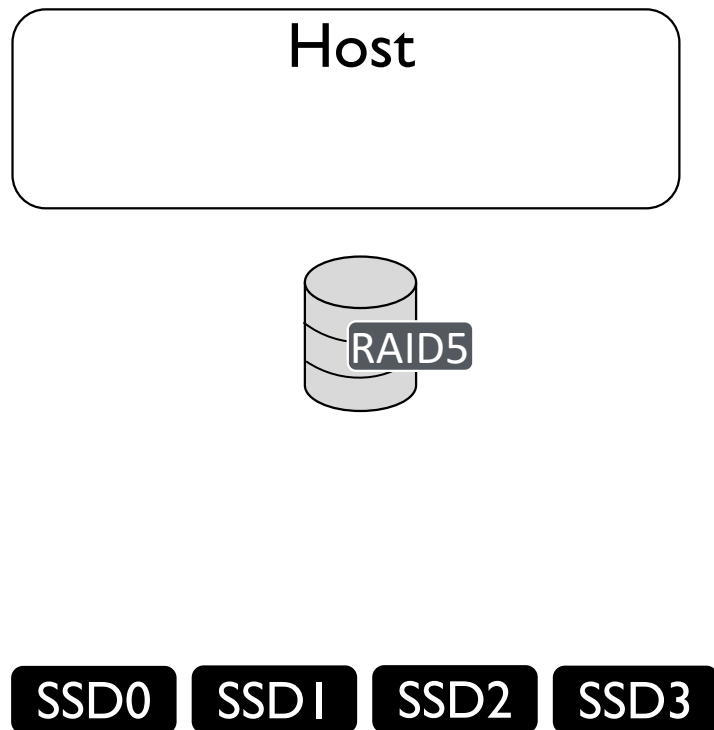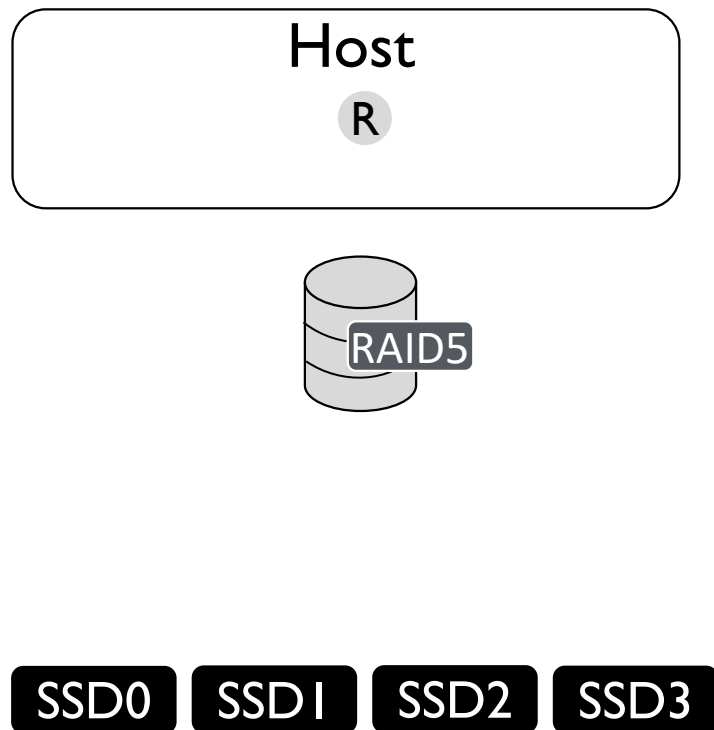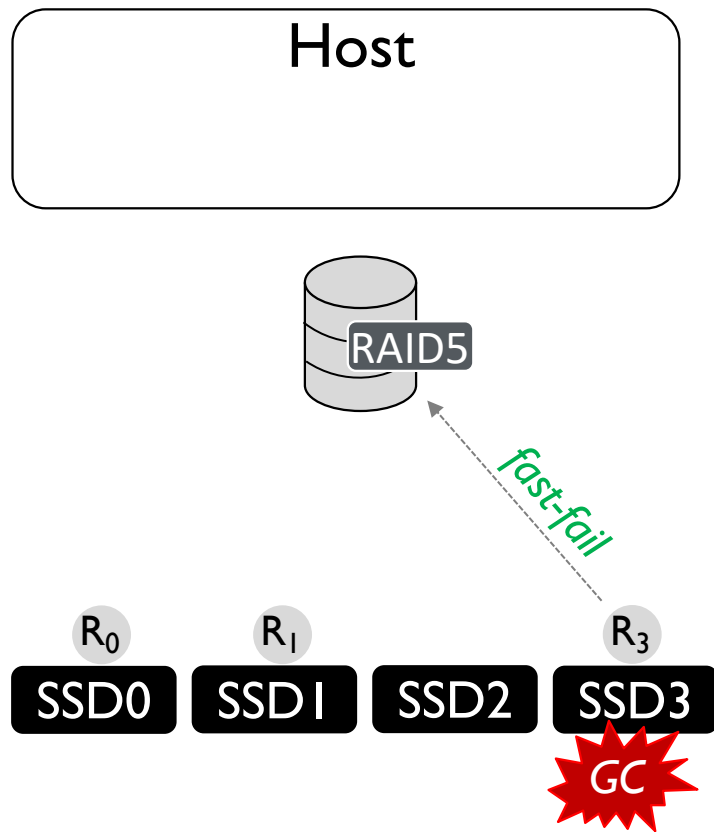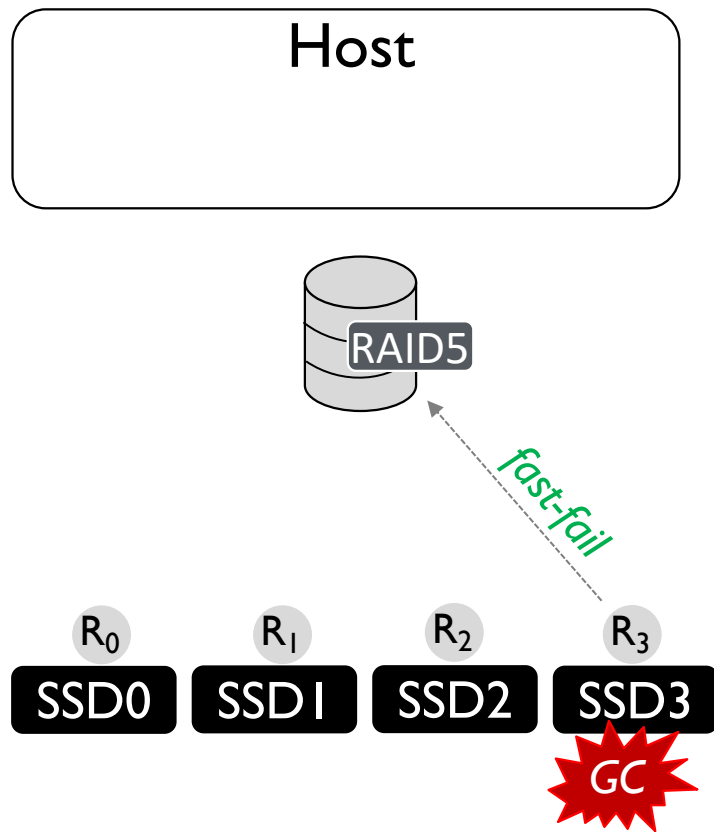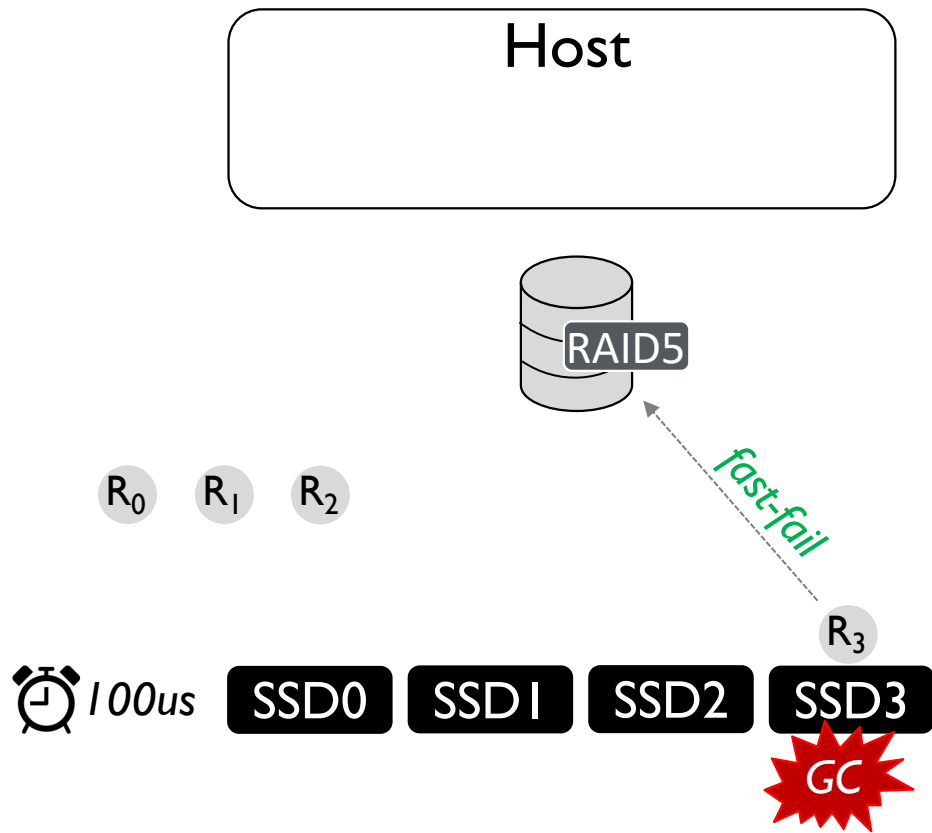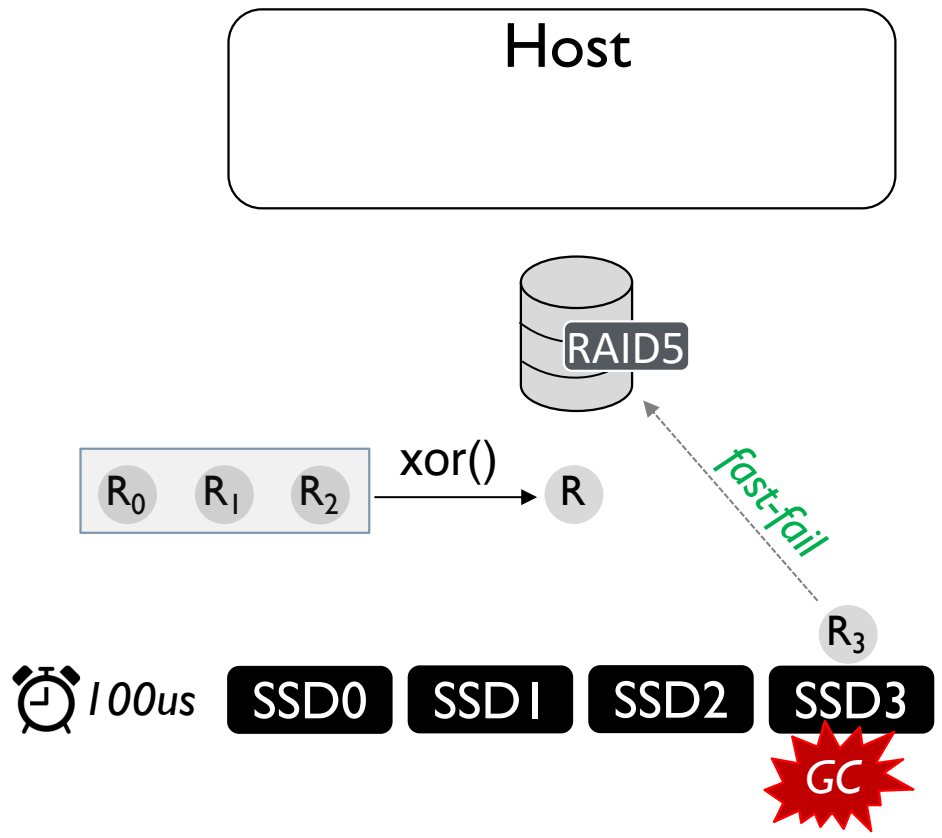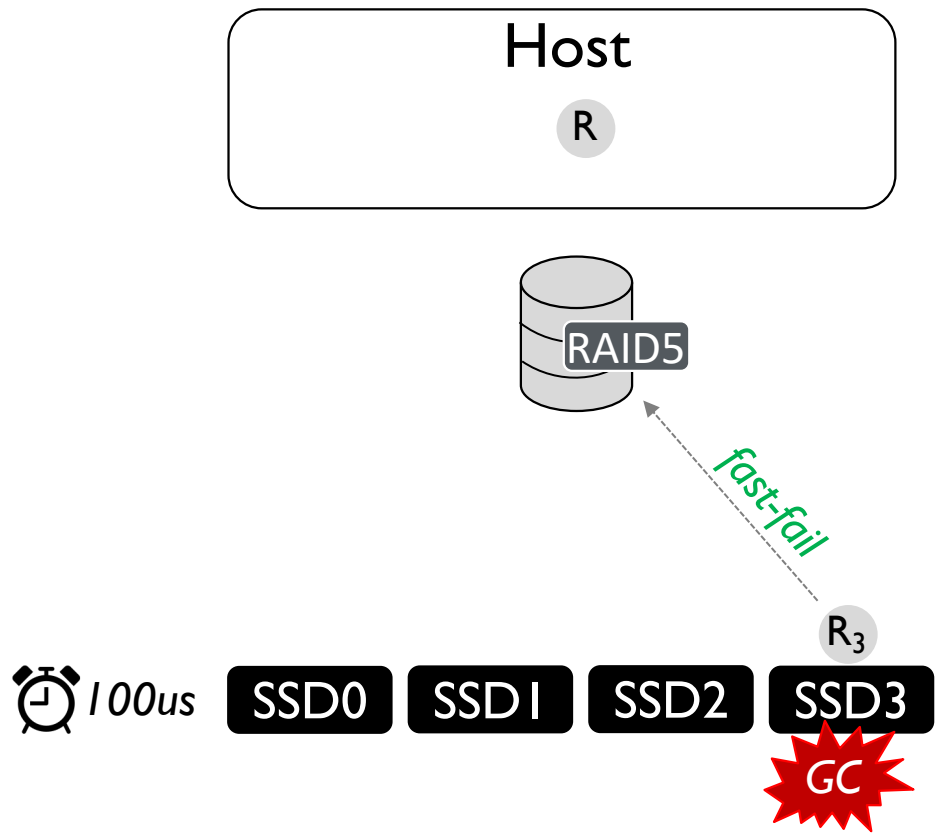💡 *"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

Host
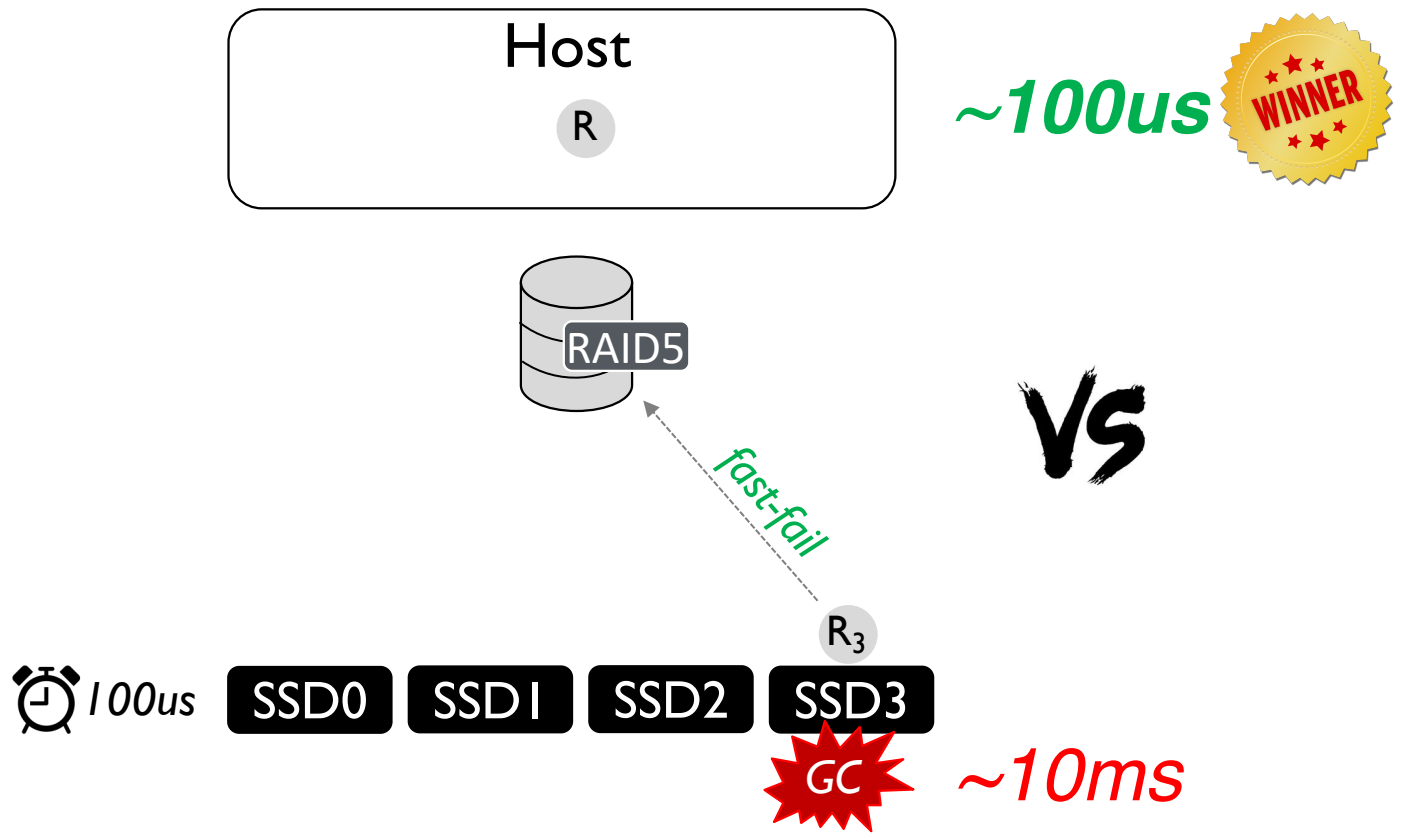
*SSD*

# IOD$_1$: Predictable Latency Flagged I/Os

💡 **"Fail-if-Slow":** the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

# IOD$_1$: Predictable Latency Flagged I/Os

💡 *"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

Host

R

RAID5

SSD0  SSD1  SSD2  SSD3

Host

RAID5

$R_0$ $R_1$ $R_2$

fast-fail

$R_3$

100us SSD0 SSD1 SSD2 SSD3

GC

Host

RAID5

$R_0$ $R_1$ $R_2$ xor() → R

fast-fail

$R_3$

⏰ 100us  SSD0  SSD1  SSD2  SSD3

GC

Host

R

RAID5

fast-fail

R₃

100us   SSD0   SSD1   SSD2   SSD3

GC

# The Effectiveness of "*Fail-if-Slow*" Interface

TPCC Read Latency CDF

# The Effectiveness of "*Fail-if-Slow*" Interface

TPCC Read Latency CDF

# The Effectiveness of "*Fail-if-Slow*" Interface

TPCC Read Latency CDF

# The Effectiveness of "*Fail-if-Slow*" Interface



TPCC Read Latency CDF

@p99 ~5x

@p95 ~7x

# The Effectiveness of "*Fail-if-Slow*" Interface

TPCC Read Latency CDF



@p99 ~5x

@p95 ~7x

Latency (ms)

*Cut tails up to ~99th percentile*

# The Effectiveness of "*Fail-if-Slow*" Interface



TPCC Read Latency CDF

@p99 ~5x

@p95 ~7x

Latency (ms)

Percentage (%)

# of busy sub-IOs

*Cut tails up to ~99th percentile*

# The Effectiveness of "*Fail-if-Slow*" Interface



TPCC Read Latency CDF

@p99 ~5x

@p95 ~7x

Latency (ms)

Percentage (%)

# of busy sub-IOs

*Cut tails up to ~99$^{th}$ percentile*

# The Effectiveness of "*Fail-if-Slow*" Interface

**TPCC Read Latency CDF**



@p99 ~5x

@p95 ~7x

Latency (ms)

Percentage (%)

# of busy sub-IOs

*Cut tails up to ~99th percentile*

# The Effectiveness of "*Fail-if-Slow*" Interface

NoGC

TPCC Read Latency CDF

@p99 ~5x

@p95 ~7x

Latency (ms)

Percentage (%)

# of busy sub-IOs

*Cut tails up to ~99$^{th}$ percentile*

# The Effectiveness of "*Fail-if-Slow*" Interface



Cut tails up to ~99<sup>th</sup> percentile

# A Case Against Proactive Reconstruction

Host

RAID5

SSD0    SSD1    SSD2    SSD3

# A Case Against Proactive Reconstruction

# A Case Against Proactive Reconstruction

# A Case Against Proactive Reconstruction



Host

RAID5

$R_0$ $R_1$ ✗ xor() → R

fast-fail

$R_2$ $R_3$

SSD0 SSD1 SSD2 SSD3

GC GC

*Semantic Gap*: the host doesn't know how long SSD "busyness" will last

End up waiting for the busiest SSD

# Busy Remaining Time (BRT) Exposure

*"Fail-if-Slow"*: the SSD should *fast-fail* an I/O if it contends with GC

**+**

Piggybacking **BRT** to reconstruct data from less busy SSDs

# Busy Remaining Time (BRT) Exposure

💡 *"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

**+**

💡 Piggybacking ***BRT*** to reconstruct data from less busy SSDs

Host

flag=true        *Fast-Fail*

*SSD*

# Busy Remaining Time (BRT) Exposure

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

**+**

Piggybacking **BRT** to reconstruct data from less busy SSDs

Host

flag=true    *Fast-Fail*    *"BRT: 60ms"*

*SSD*

# The Effectiveness of "*BRT*" Interface

# The Effectiveness of "*BRT*" Interface

TPCC Read Latency CDF

# The Effectiveness of "*BRT*" Interface



TPCC Read Latency CDF

# The Effectiveness of "*BRT*" Interface



TPCC Read Latency CDF

BRT helps a little

# The Effectiveness of "*BRT*" Interface

TPCC Read Latency CDF



BRT helps a little

Latency (ms)

# The Effectiveness of "*BRT*" Interface

TPCC Read Latency CDF



BRT helps a little

*Can we do better?*

# IODA Busy Latency Windows

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

**+**

*TW Coordination*: SSDs take turns to perform GCs

# IODA Busy Latency Windows

💡 *"Fail-if-Slow"*: the SSD should *fast-fail* an I/O if it contends with GC

**+**

💡 *TW Coordination*: SSDs take turns to perform GCs

| | | | | |
|---|---|---|---|---|
| **SSD#3** | Predictable | Predictable | Predictable | Busy |
| **SSD#2** | Predictable | Predictable | Busy | Predictable |
| **SSD#1** | Predictable | Busy | Predictable | Predictable |
| **SSD#0** | Busy | Predictable | Predictable | Predictable |

$t$      $t+TW$      $t+2{\times}TW$      $t+3{\times}TW$      $t+4{\times}TW$

# IODA Busy Latency Windows

*"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

**+**

*TW Coordination*: SSDs take turns to perform GCs

**IODA: Always Predictable Latencies!** ☺

| | $t$ | $t+TW$ | $t+2{\times}TW$ | $t+3{\times}TW$ | $t+4{\times}TW$ |
|---|---|---|---|---|---|
| SSD#3 | Predictable | Predictable | Predictable | Busy | |
| SSD#2 | Predictable | Predictable | Busy | Predictable | |
| SSD#1 | Predictable | Busy | Predictable | Predictable | |
| SSD#0 | Busy | Predictable | Predictable | Predictable | |

# IODA Busy Latency Windows

💡 *"Fail-if-Slow":* the SSD should *fast-fail* an I/O if it contends with GC

➕

💡 *TW Coordination*: SSDs take turns to perform GCs

**IODA: Always Predictable Latencies!** ☺

How long should *TW* be?

| | | | | |
|---|---|---|---|---|
| SSD#3 | Predictable | Predictable | Predictable | Busy |
| SSD#2 | Predictable | Predictable | Busy | Predictable |
| SSD#1 | Predictable | Busy | Predictable | Predictable |
| SSD#0 | Busy | Predictable | Predictable | Predictable |

$t$     $t+TW$     $t+2 \times TW$     $t+3 \times TW$     $t+4 \times TW$

# IODA Time Window (*TW*) Formulation

# IODA Time Window (*TW*) Formulation

*SSD free space* >= *User load*

# IODA Time Window (*TW*) Formulation

*SSD free space* >= *User load*

SSD

Busy

t0    t1    t2    t3

# IODA Time Window (*TW*) Formulation

SSD free space >= User load

$B_{burst}$ : User load



SSD

Busy

t0    t1    t2    t3

# IODA Time Window (*TW*) Formulation

SSD *free space* >= *User load*

$B_{burst}$ : *User load*

SSD

Busy

t0   t1   t2   t3

$B_{gc}$ : *GC reclamation speed*

$S_p$ : *Over-provisioning space*

# IODA Time Window (*TW*) Formulation

*SSD free space* **>=** *User load*

$B_{burst}$ : *User load*

SSD

Busy

t0     t1     t2     t3

$$TW \leq S_p \ / \ ((N_{ssd} \times B_{burst}) - B_{gc})$$

$B_{gc}$ : *GC reclamation speed*

$S_p$ : *Over-provisioning space*

# IODA Time Window (*TW*) Formulation

*SSD free space* >= *User load*

$B_{burst}$ : *User load*

SSD

Busy

t0    t1    t2    t3

$TW \leq S_p \ / \ ((N_{ssd} \times B_{burst}) - B_{gc})$

$B_{gc}$ : *GC reclamation speed*

$S_p$ : *Over-provisioning space*

$$TW \leq \frac{R_p \times S_t}{(N_{ssd} \times Min(B_{pcie}, Max(\frac{N_{dwpd} \times (1 - R_p) \times S_t}{8 hours/day}))) - (\frac{(1 - R_v) \times N_{ch} \times S_{pg} \times N_{pg}}{(t_r + t_w + 2 \times t_{cpt}) \times R_v \times N_{pg} + t_e})}$$

*TW Upper Bound*

TPCC Read Latency CDF

No Tails!

TPCC Read Latency CDF

No Tails!

TPCC Read Latency CDF

**IODA** closes the gap between **Base** and **NoGC**

# More in the paper!

❑ IODA *TW* analysis
- ▪ **6** SSD models
- ▪ Relaxed *TW*
- ▪ *TW vs. WAF* tradeoffs

❑ Implementation
- ▪ Platforms: FEMU + OpenChannel-SSD
- ▪ Kernel: Linux Software-RAID + NVMe

❑ More evaluation results
- ▪ **9** datacenter block traces + *21* real applications
- ▪ IODA vs. **7** State-of-the-art approaches
- ▪ IODA on OpenChannel-SSD
- ▪ IODA throughput and write latency
- ▪ …

# IODA Stack and Evaluation Setup

Kernel

- - - - - - - - - - - - - - - - - - - - - - - - - -

SSDs
OpenChannel-SSD    FEMU

# IODA Stack and Evaluation Setup

Kernel

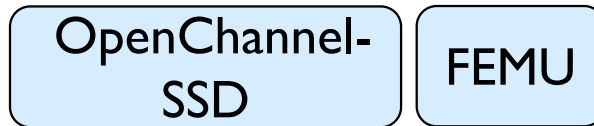| Software-RAID |
| NVMe Driver |

SQ      CQ

SSDs

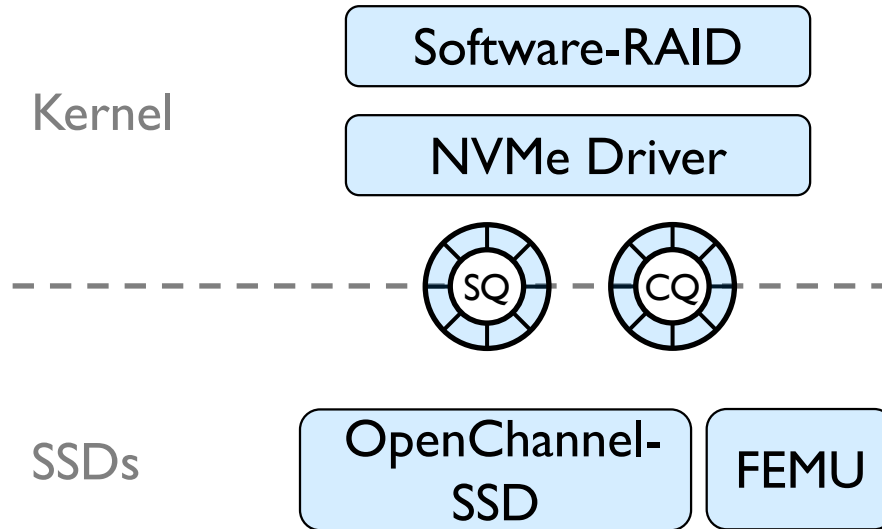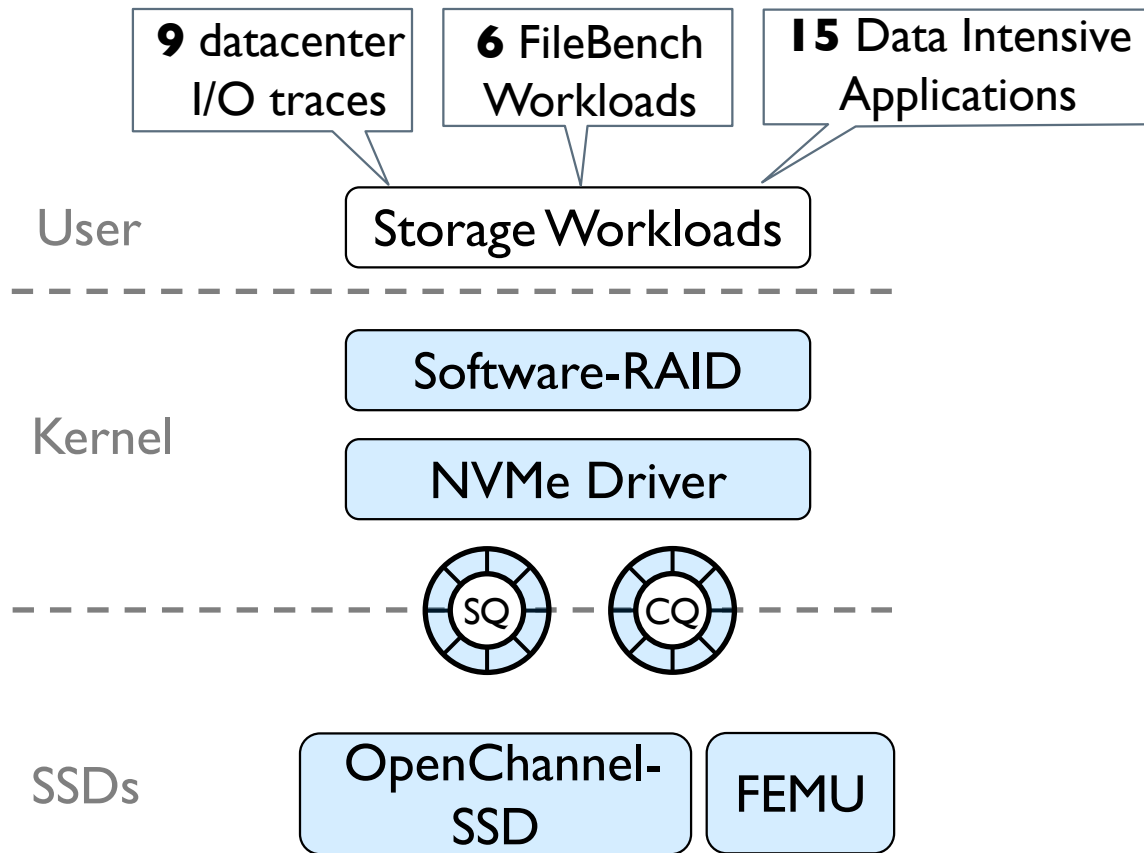| OpenChannel-SSD |  | FEMU |

# IODA Stack and Evaluation Setup

# IODA Stack and Evaluation Setup



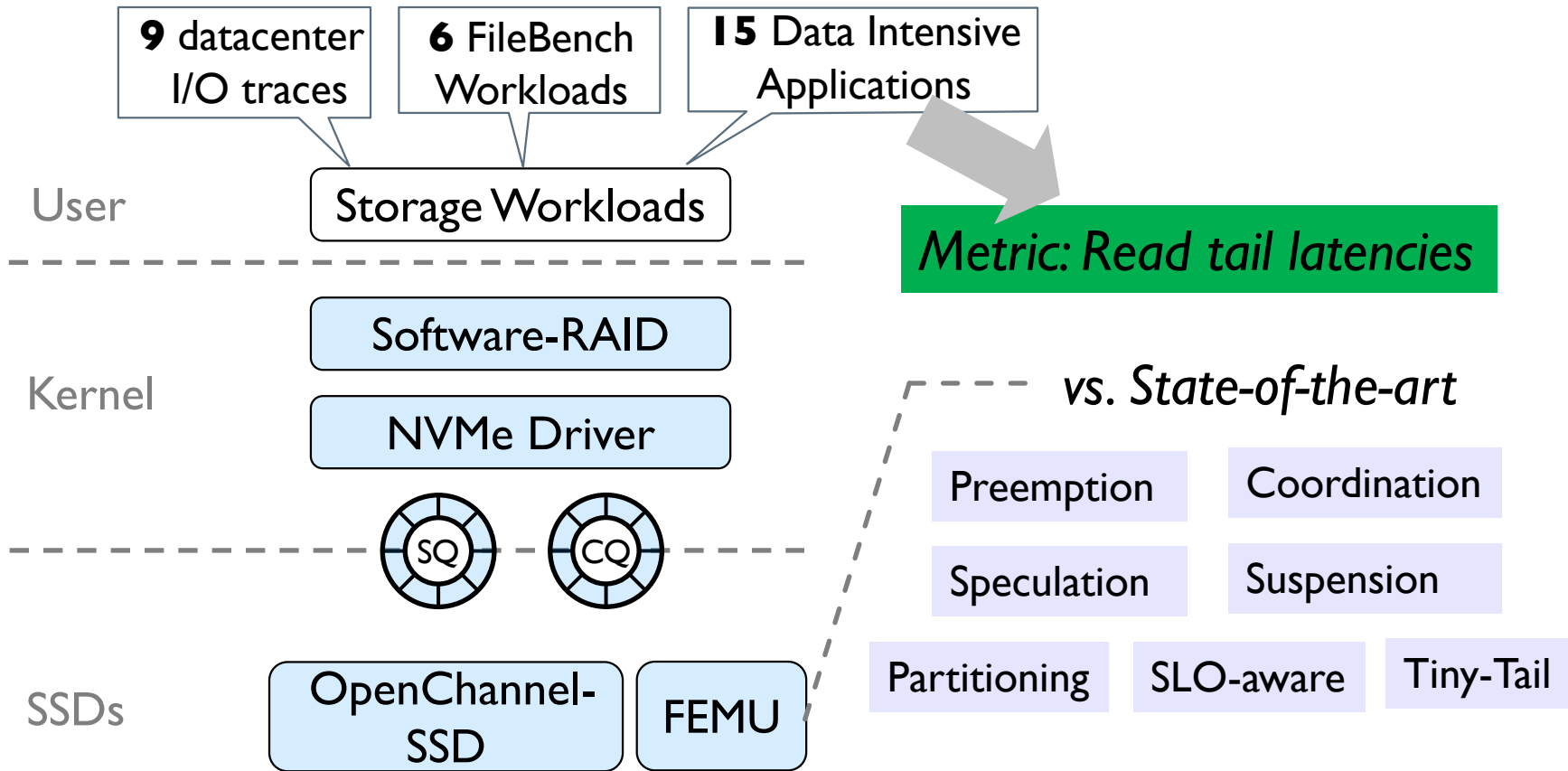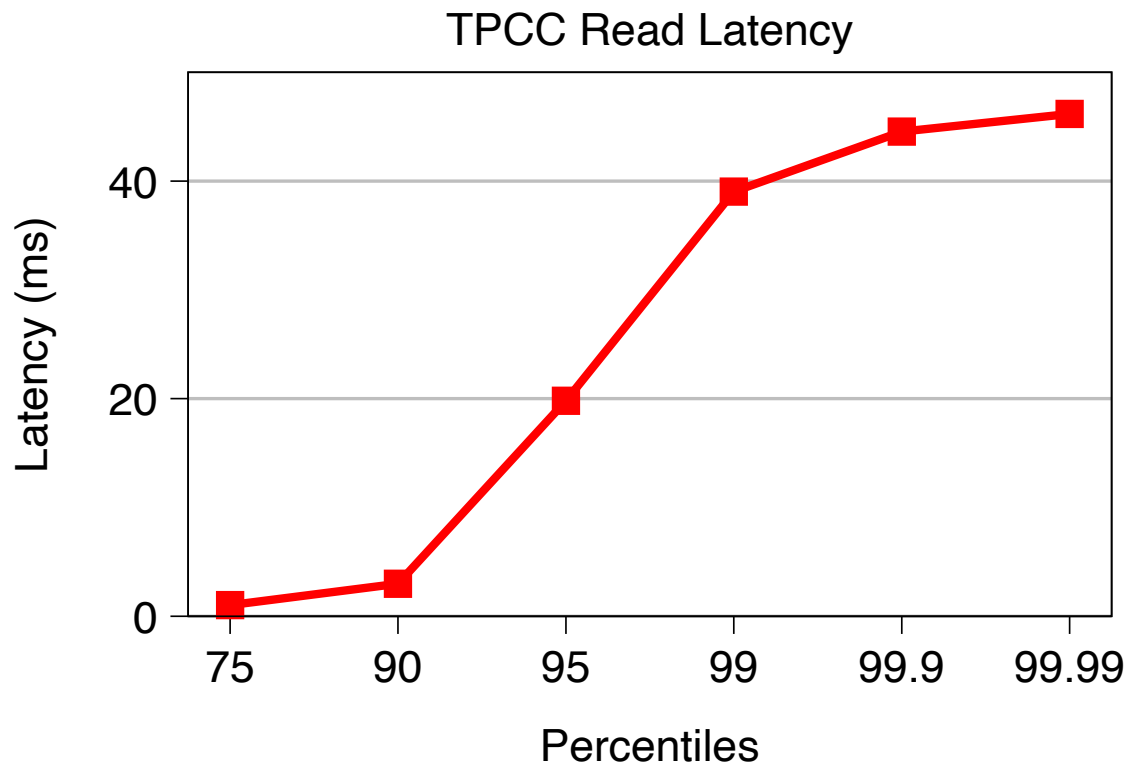**9** datacenter I/O traces

**6** FileBench Workloads

**15** Data Intensive Applications

User

Storage Workloads

Kernel

Software-RAID

NVMe Driver

SQ   CQ

SSDs

OpenChannel-SSD   FEMU

*vs. State-of-the-art*

Preemption   Coordination

Speculation   Suspension

Partitioning   SLO-aware   Tiny-Tail

# IODA Stack and Evaluation Setup

**9** datacenter I/O traces

**6** FileBench Workloads

**15** Data Intensive Applications

User

Storage Workloads

*Metric: Read tail latencies*

Kernel

Software-RAID

NVMe Driver

SQ   CQ

*vs. State-of-the-art*

Preemption

Coordination

Speculation

Suspension

SSDs

OpenChannel-SSD

FEMU

Partitioning

SLO-aware

Tiny-Tail

# IODA Evaluation



TPCC Read Latency

# IODA Evaluation



TPCC Read Latency
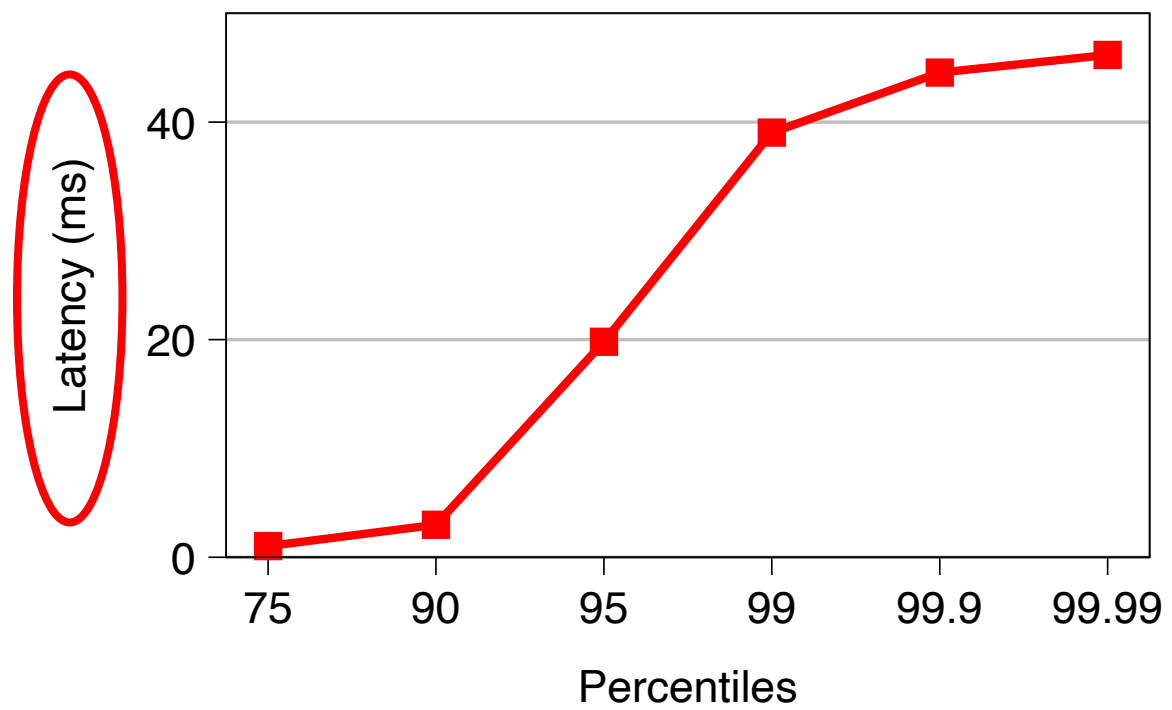
# IODA Evaluation

TPCC Read Latency

# IODA Evaluation

TPCC Read Latency

# IODA Evaluation



TPCC Read Latency

# IODA Evaluation



TPCC Read Latency

*Predictable Latency Flag*
*+ Reconstruction*

Base

# IODA Evaluation



TPCC Read Latency

*Predictable Latency Flag*
*+ Reconstruction*

*Predictable Latency Flag*
*+ Busy Remaining Time*

Base

# IODA Evaluation



TPCC Read Latency

Predictable Latency Flag
+ Reconstruction

Predictable Latency Flag
+ Busy Remaining Time

Predictable Latency Flag
+ Time Window

Base

IODA is close to Ideal!

Base ■ IOD₁ ■ IOD₂ ■ IOD₃ ■ IODA ■ Ideal ■

[a] Azure

[b] BingIdx — IODA close to Ideal

[c] BingSel

[d] Cosmos

[e] DTRS

[f] Exch

[g] LMBE

[h] MSNFS

[i] TPCC

Latency (ms)

Legend: Base, $IOD_1$, $IOD_2$, $IOD_3$, IODA, Ideal

[a] Azure

[b] BingIdx — IODA close to Ideal

[c] BingSel

[d] Cosmos

[e] DTRS

[f] Exch

[g] LMBE

[h] MSNFS

[i] TPCC

Latency (ms)

Base ⬛ IOD₁ ⬛ IOD₂ ⬛ IOD₃ ⬛ IODA ⬛ Ideal ⬛

[a] Azure
[b] BingIdx — IODA close to Ideal
[c] BingSel
[d] Cosmos
[e] DTRS
[f] Exch
[g] LMBE
[h] MSNFS
[i] TPCC

Latency (ms)

# IODA Results: (95th – 99.99th)

Up to 75x improvement *over Base*

Base ▬ IOD₁ ▬ IOD₂ ▬ IOD₃ ▬ IODA ▬ Ideal ▬

[a] Azure
[b] BingIdx — IODA close to Ideal
[c] BingSel
[d] Cosmos
[e] DTRS
[f] Exch
[g] LMBE
[h] MSNFS
[i] TPCC

Latency (ms)

# IODA Results: ($95^{th} - 99.99^{th}$)
Up to 75x improvement *over Base*
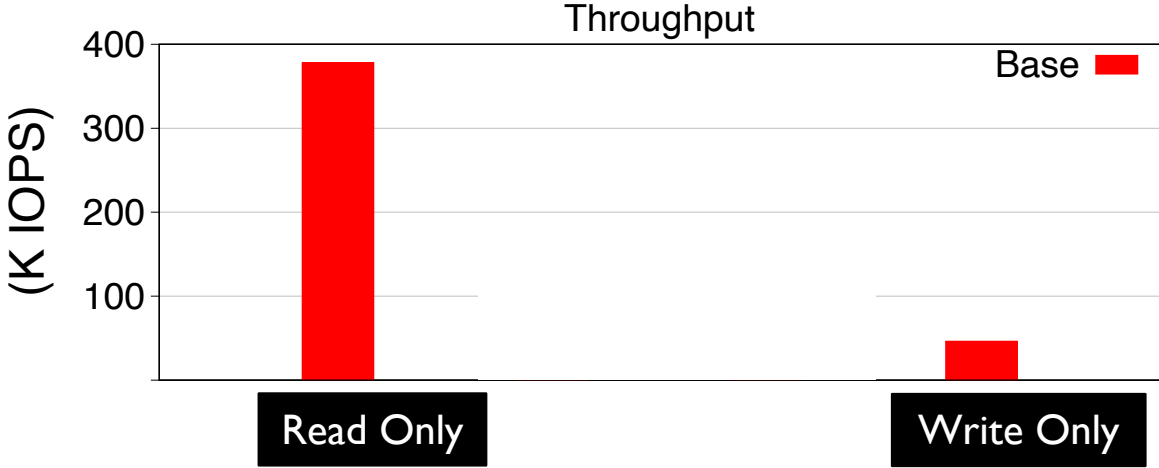
## vs.

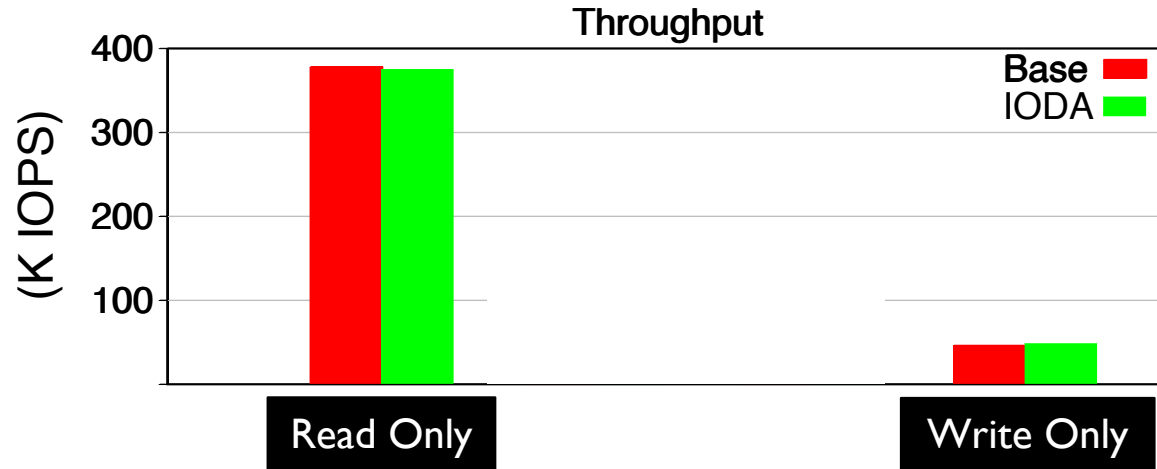| | |
|---|---|
| Preemption | Coordination |
| Speculation | Suspension |
| Partitioning | SLO-aware | Tiny-Tail |

*IODA is more deterministic and efficient in cutting tail latencies!*
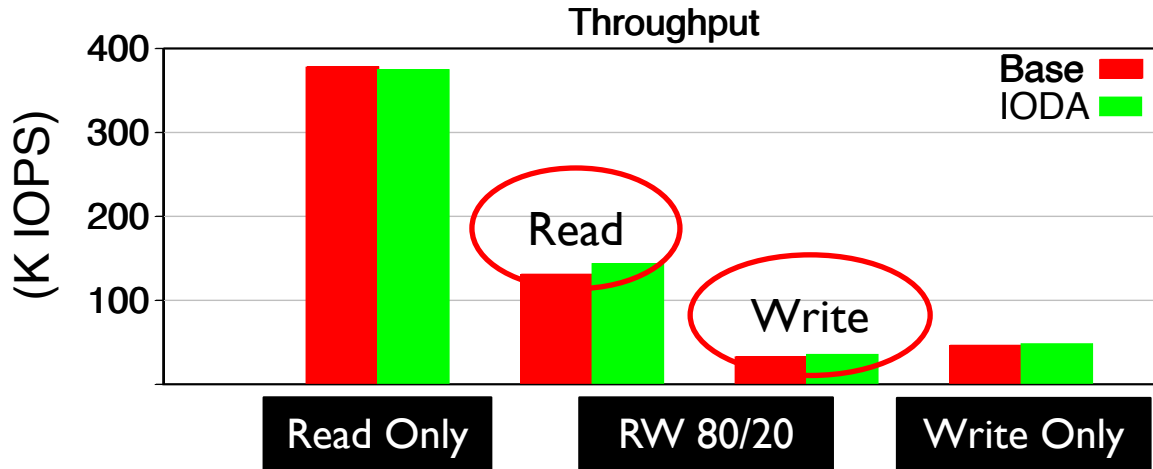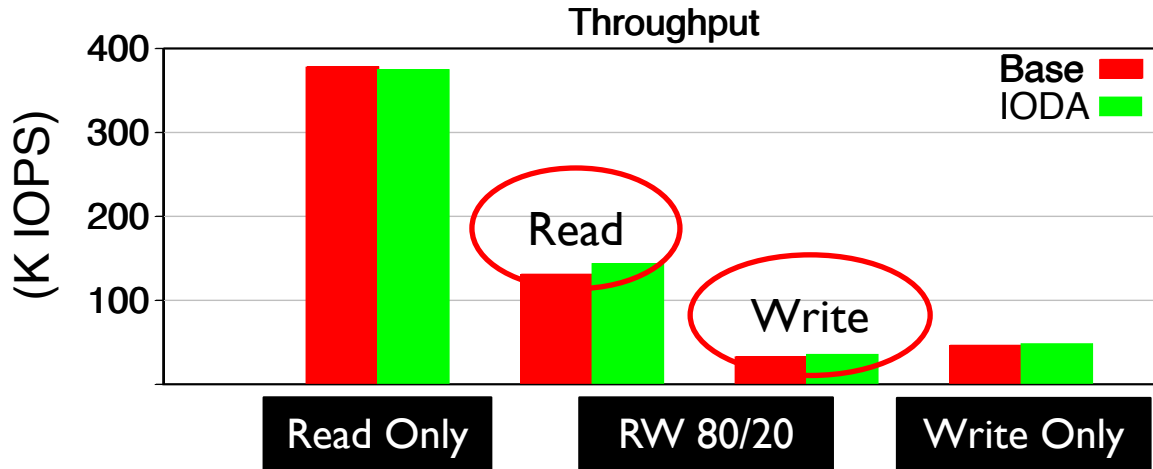
# IODA Throughput

# IODA Throughput

# IODA Throughput

# IODA Throughput



*IODA doesn't sacrifice the array's aggregate bandwidth*

# IODA Takeaways

❑ A *Co-Design* Approach for Performance Predictability
- Proactive *reconstruction* via *fast-fail* interface
- *BRT* for improved latencies
- *TW* formulation to program the window length
- Cross-device synchronization

*I'm on the job market.*

IODA: https://github.com/huaicheng/IODA

# IODA Takeaways

❑ A *Co-Design* Approach for Performance Predictability

- Proactive *reconstruction* via *fast-fail* interface
- *BRT* for improved latencies
- *TW* formulation to program the window length
- Cross-device synchronization

# Thank you!

*I'm on the job market.*

IODA: https://github.com/huaicheng/IODA