# Tiny-Tail Flash

## Near-Perfect Elimination of Garbage Collection Tail Latencies in NAND SSDs

**Shiqin Yan**, Huaicheng Li, Mingzhe Hao,
Michael Hao Tong, Swaminathan Sundararaman[*],
Andrew Chien, and Haryadi S. Gunawi

CERES
Center for Unstoppable Computing

THE UNIVERSITY OF CHICAGO

P/A [*]

# Why SSDs don't perform

From their earliest days, people have reported that SSDs were not providing the performance they expected. As SSDs age, for instance, they get slower. Here's why.

## Google: Taming The Long Latency Tail - When More Machines Equals Worse Results

# Why it's hard to meet SLAs with SSDs

# Why SSDs don't perform

From their earliest days, people have reported that SSDs were not providing the performance they expected. As SSDs age, for instance, they get slower. Here's why.

## Google: Taming The Long Latency Tail - When More Machines Equals Worse Results

*"if your read is stuck behind an erase you may have wait **10s of milliseconds**. That's a **100x** increase in latency variance"*

# Why it's hard to meet SLAs with SSDs
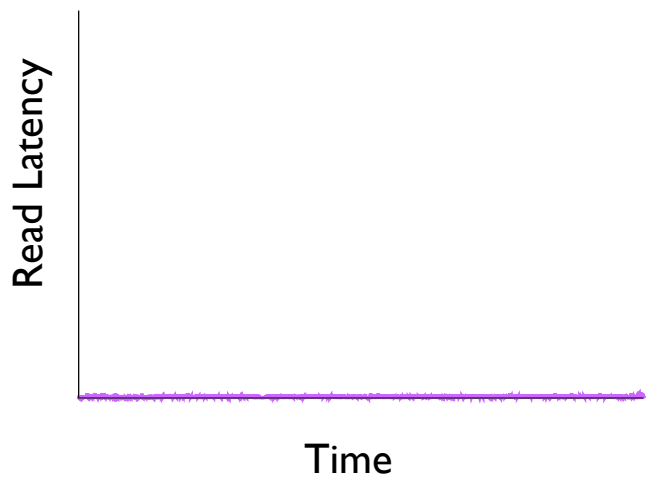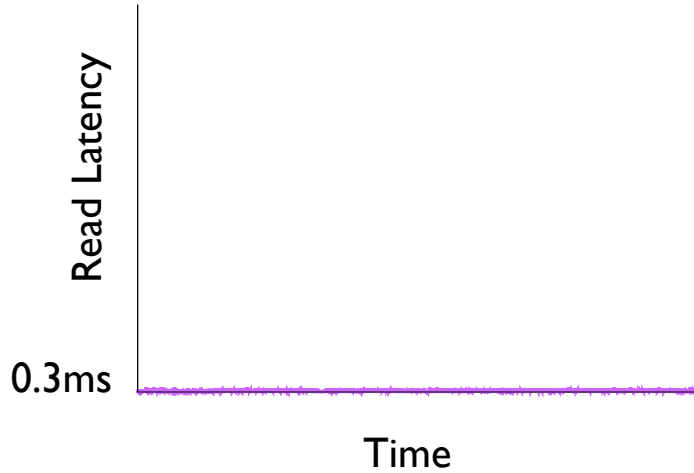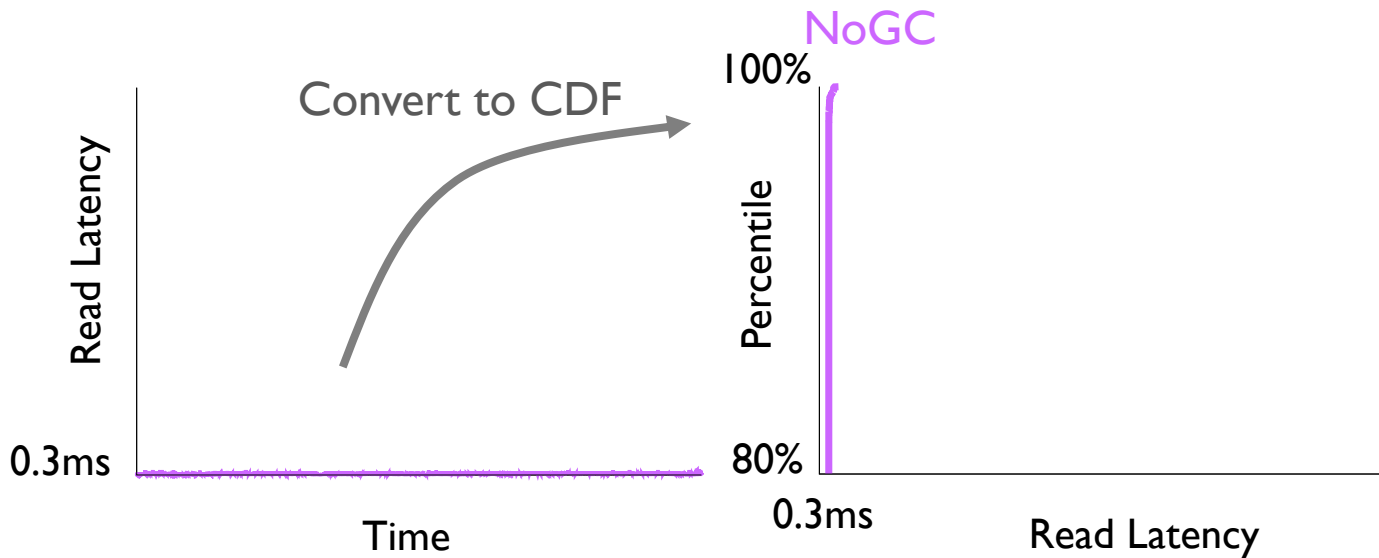
Reads + Writes

Clean/Empty
SSD

Reads + Writes

Clean/Empty
SSD

Read Latency

Time

THE UNIVERSITY OF CHICAGO

Reads + Writes

Clean/Empty
SSD

Read Latency

0.3ms

Time

THE UNIVERSITY OF CHICAGO

Reads + Writes

Clean/Empty SSD

Read Latency

0.3ms

Time

Convert to CDF

NoGC

100%

Percentile

80%

0.3ms

Read Latency

Reads + Writes

Clean/Empty
SSD

Read Latency

Convert to CDF

0.3ms

Time

NoGC

100%

Percentile

80%

0.3ms

Read Latency

Reads + Writes

Clean/Empty
SSD

Read Latency

0.3ms

Time

Convert to CDF

NoGC

100%

Percentile

80%

0.3ms

Read Latency

Reads + Writes

Clean/Empty
SSD

Read Latency

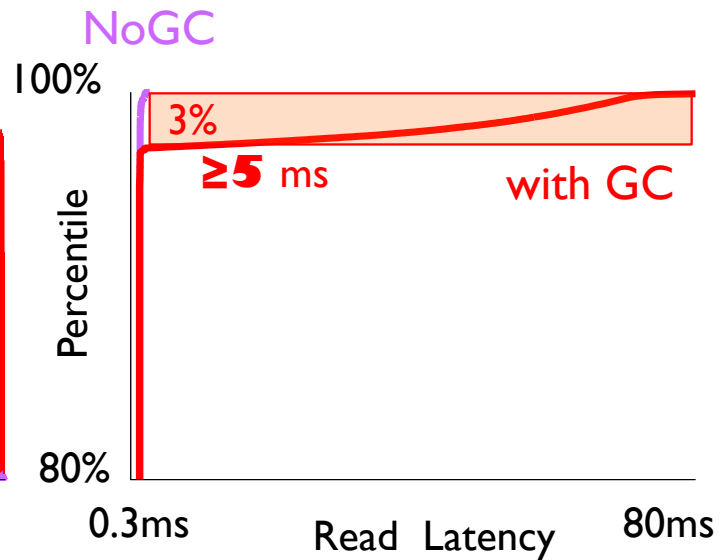Convert to CDF

0.3ms

Time

NoGC

100%

Percentile

80%

0.3ms

Read Latency
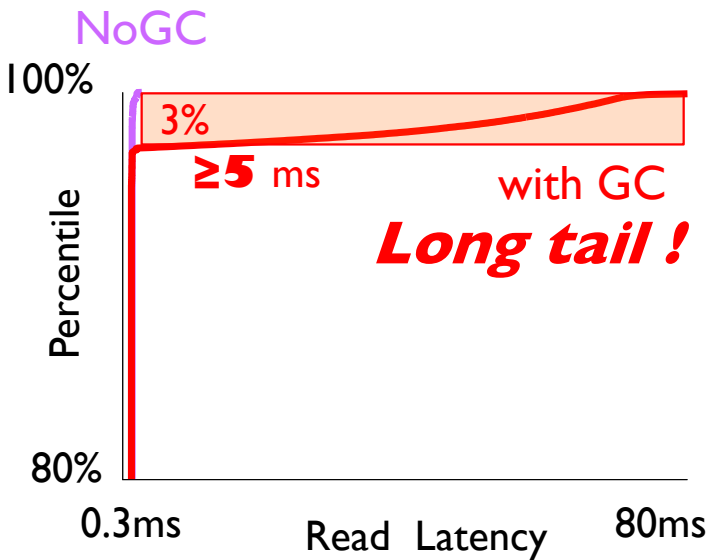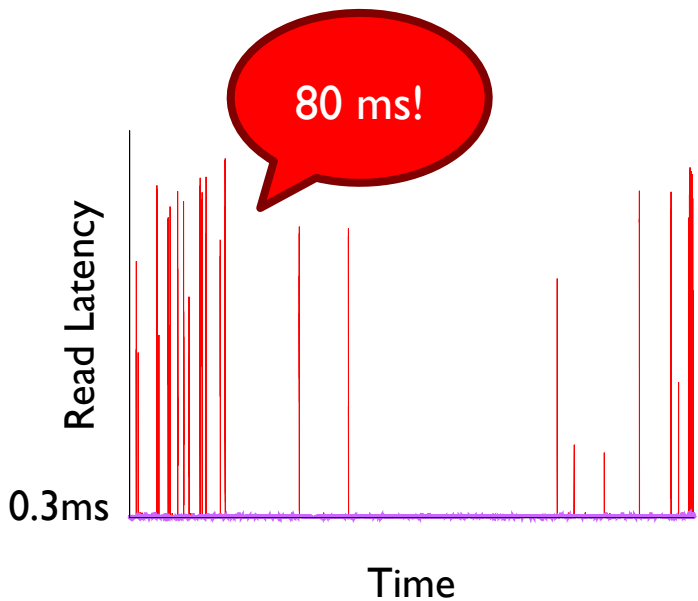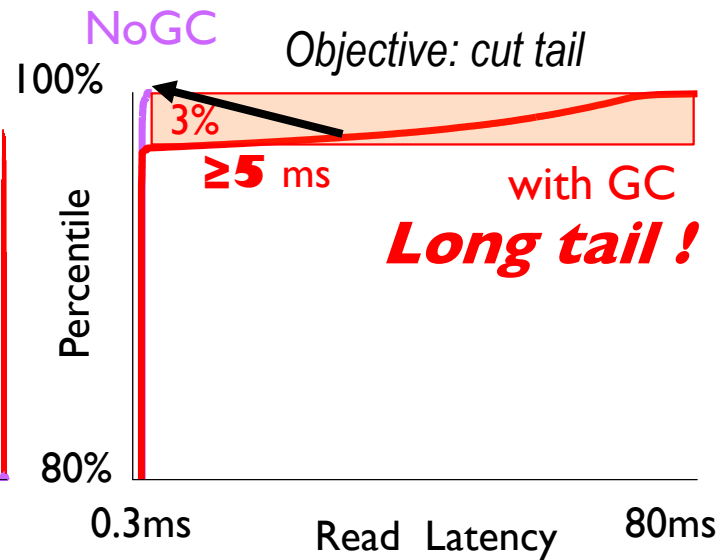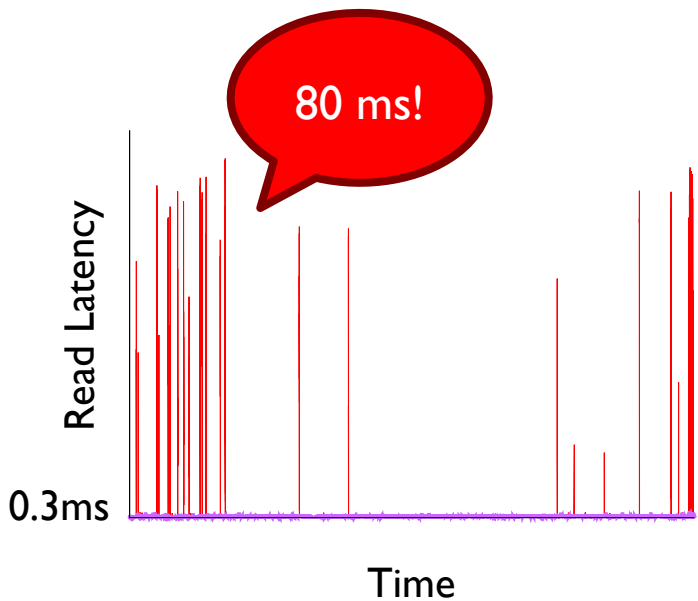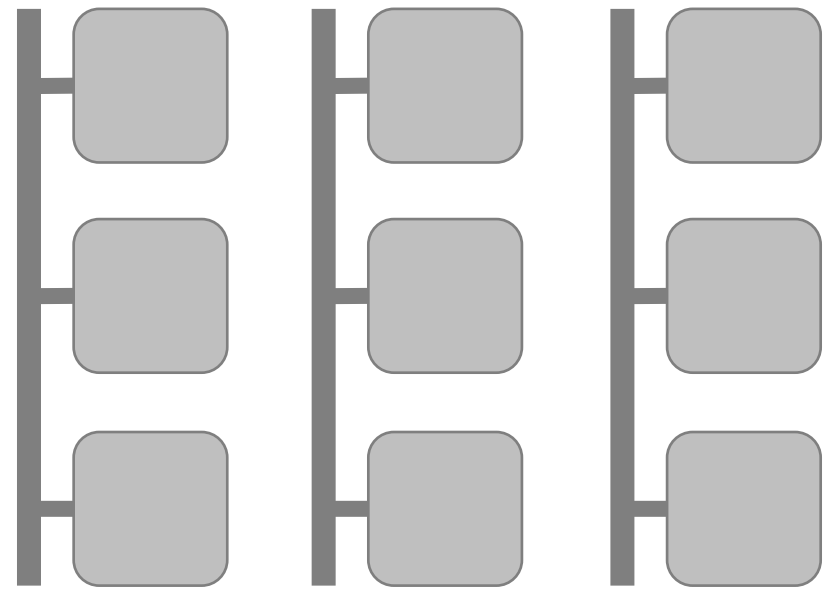
**Aged/Full
SSD**

Reads + Writes

**Aged/Full SSD**

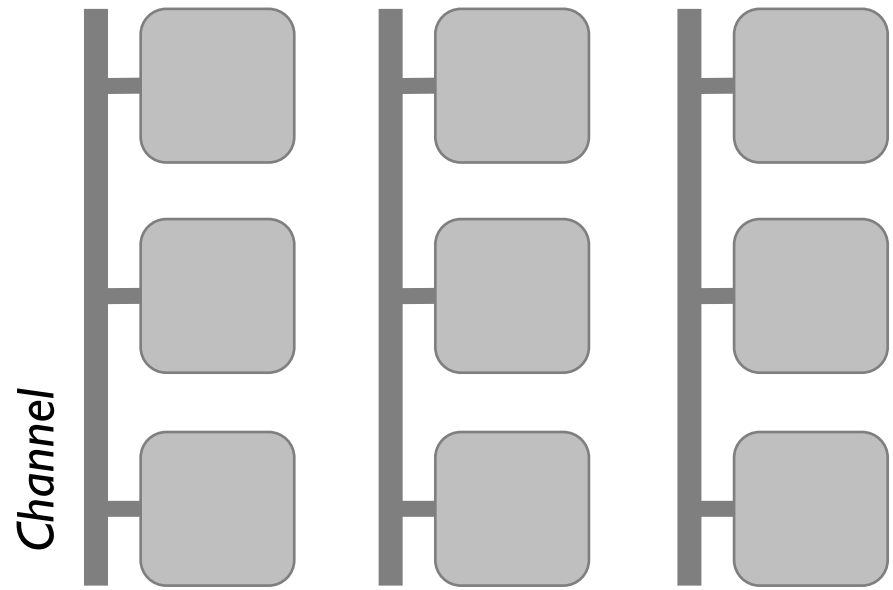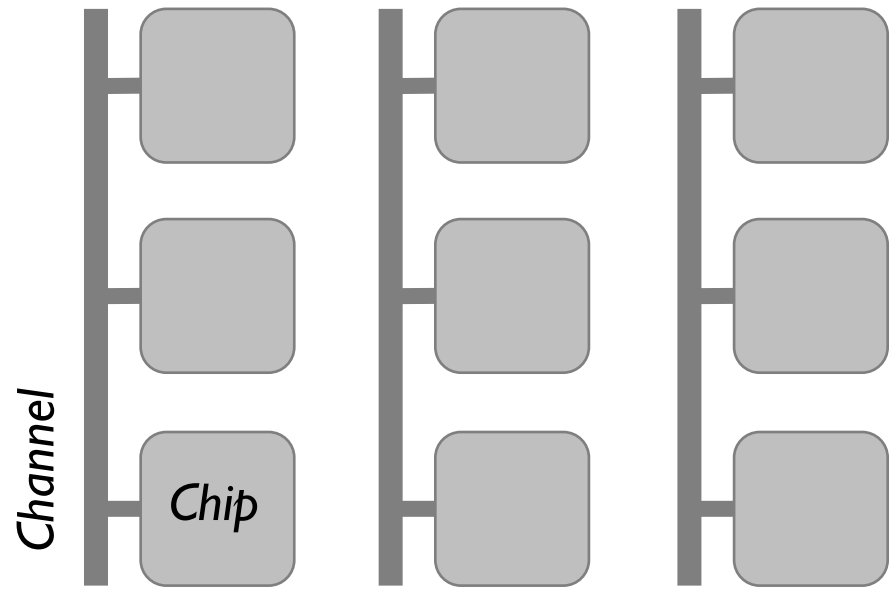# How GC delays read I/Os?

# How GC delays read I/Os?



*Channel*

# How GC delays read I/Os?



Channel

Chip

# How GC delays read I/Os?

# How GC delays read I/Os?



*A GC moves tens of valid pages!*

*Channel*

*Chip*

# How GC delays read I/Os?



*A GC moves tens of valid pages!*

*which makes* channel/chips busy *for* **tens of ms !**

# How GC delays read I/Os?

Read [ A ]

fast

Channel

Chip

A GC moves tens of valid pages!

which makes channel/chips busy for **tens of ms !**

# How GC delays read I/Os?



Read　A　B

*fast*

*delayed!*

*Channel*

*Chip*

A GC moves tens of valid pages!

which makes channel/chips busy for **tens of ms !**

# How to cut tail latencies?

# How to cut tail latencies?

Tail-tolerant techniques in distributed/storage systems:
*Leverage redundancy to cut tail!*

# How to cut tail latencies?

Tail-tolerant techniques in distributed/storage systems:
*Leverage redundancy to cut tail!*

*Full Stripe Read*

| A | B | C |
|---|---|---|

**RAID:**

# How to cut tail latencies?

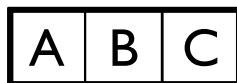Tail-tolerant techniques in distributed/storage systems:
*Leverage redundancy to cut tail!*

*Full Stripe Read*

**RAID:**
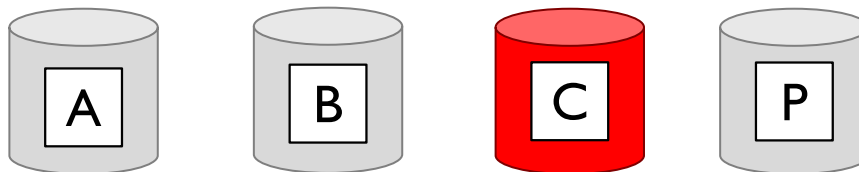


*fast*

**tail!**

Slow /
busy

# How to cut tail latencies?

Tail-tolerant techniques in distributed/storage systems:
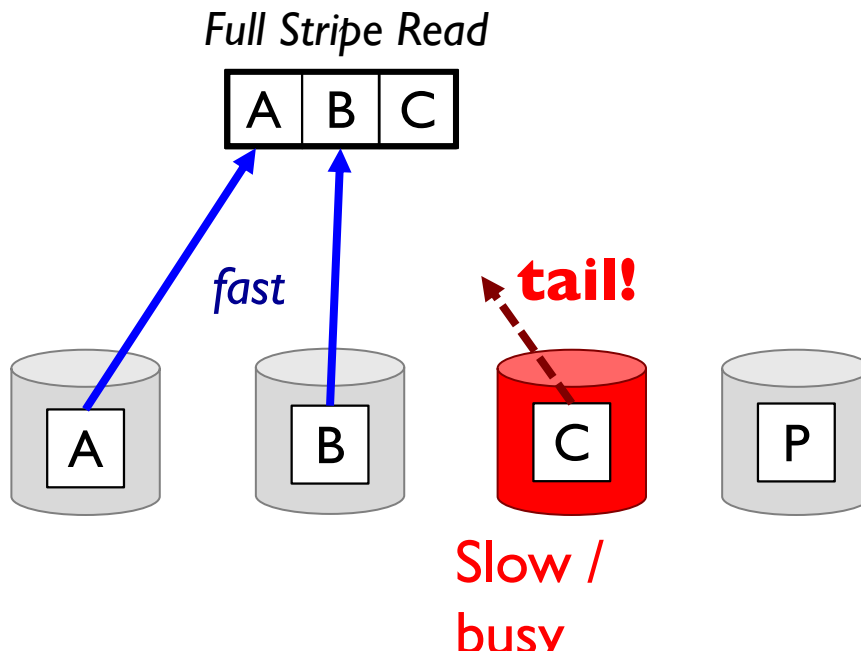*Leverage redundancy to cut tail!*



*Full Stripe Read*

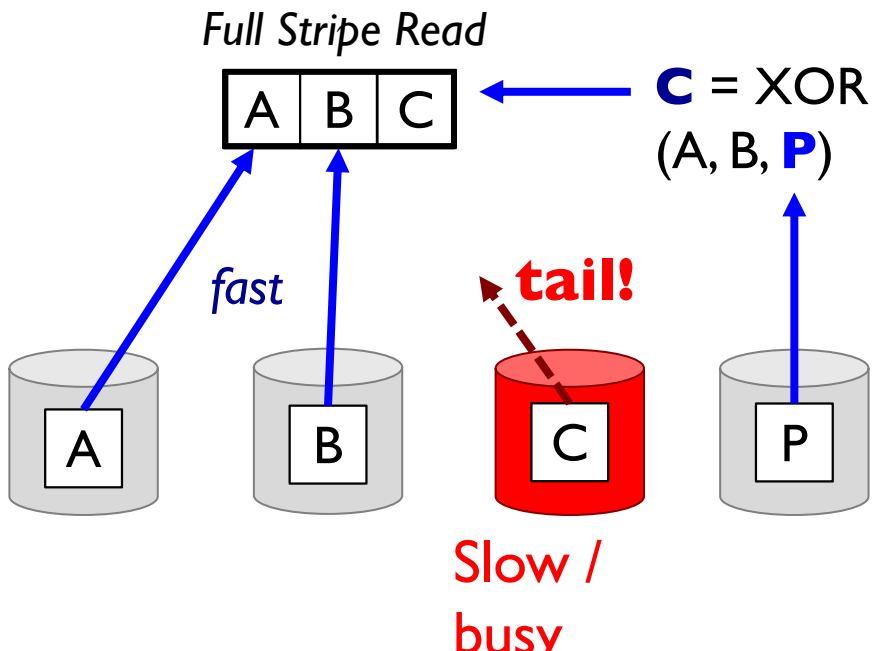A | B | C

**C** = XOR (A, B, **P**)

*fast*

**tail!**

**RAID:**

A    B    C    P

Slow / busy

# How to cut tails in SSD?

# How to cut tails in SSD?

*Error rate increases* ➔ **RAIN** (**R**edundant **A**rray of **I**ndependent **N**AND)
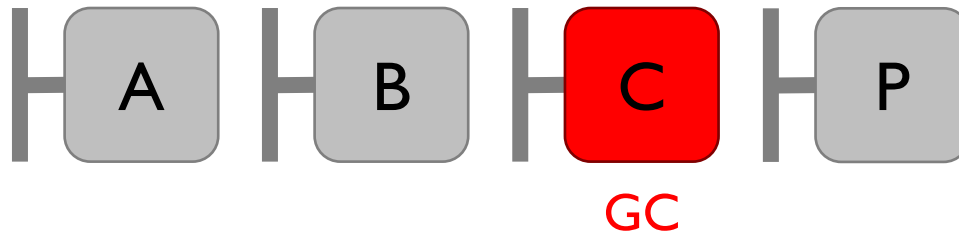
# How to cut tails in SSD?

*Error rate increases* → **RAIN** (**R**edundant **A**rray of **I**ndependent **N**AND)

*Similarly, we leverage RAIN to cut "tails"!*

Full Stripe Read

| A | B | C |

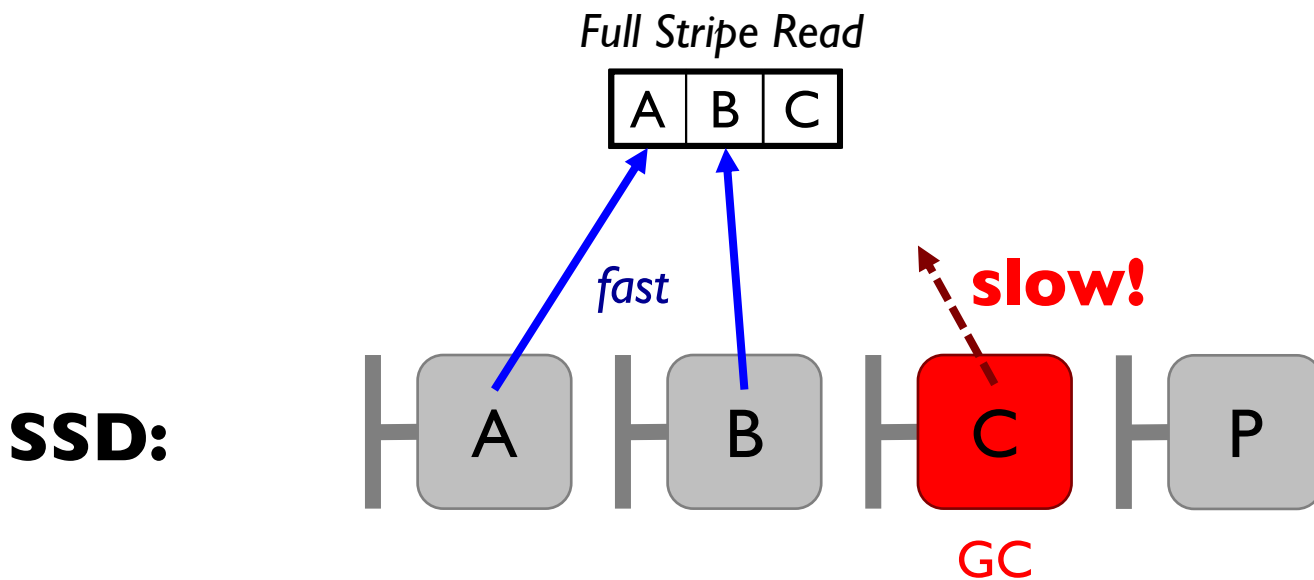**SSD:**    A    B    C    P

GC

# How to cut tails in SSD?

*Error rate increases* → **RAIN** (**R**edundant **A**rray of **I**ndependent **N**AND)

*Similarly, we leverage RAIN to cut "tails"!*

# How to cut tails in SSD?

*Error rate increases* → **RAIN** (**R**edundant **A**rray of **I**ndependent **N**AND)

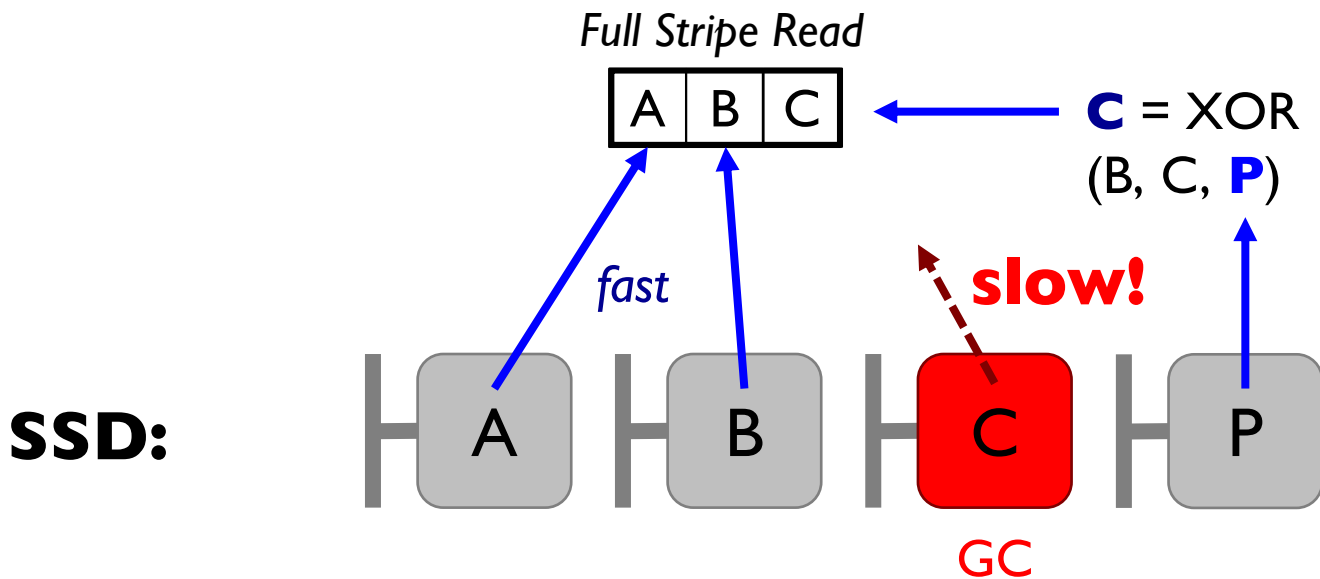*Similarly, we leverage RAIN to cut "tails"!*



*Full Stripe Read*

A B C ← **C** = XOR (B, C, **P**)

*fast*     **slow!**

**SSD:**     A     B     C     P

GC

# Contribution

*New techniques:*

*Current SSD technology:*

**RAIN**

(Parity-based Redundancy)

# Contribution

*New techniques:*

I.  Plane-Blocking GC

*Current SSD technology:*

**RAIN**
(Parity-based Redundancy)

# Contribution

*New techniques:*

I. Plane-Blocking GC

II. GC-Tolerant Read

---

*Current SSD technology:*

**RAIN**
(Parity-based Redundancy)

# Contribution

*New techniques:*

I. Plane-Blocking GC

II. GC-Tolerant Read

III. Rotating GC

---

*Current SSD technology:*
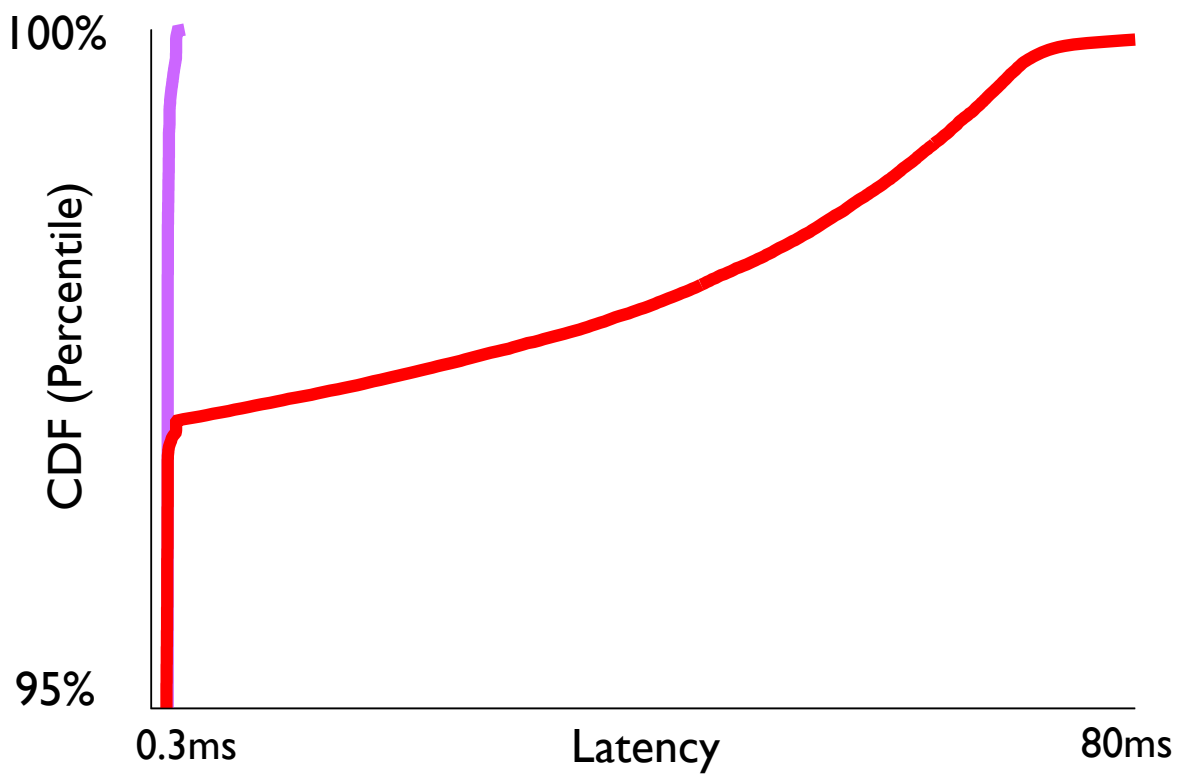
**RAIN**
(Parity-based Redundancy)

# Contribution

*New techniques:*

I.   Plane-Blocking GC

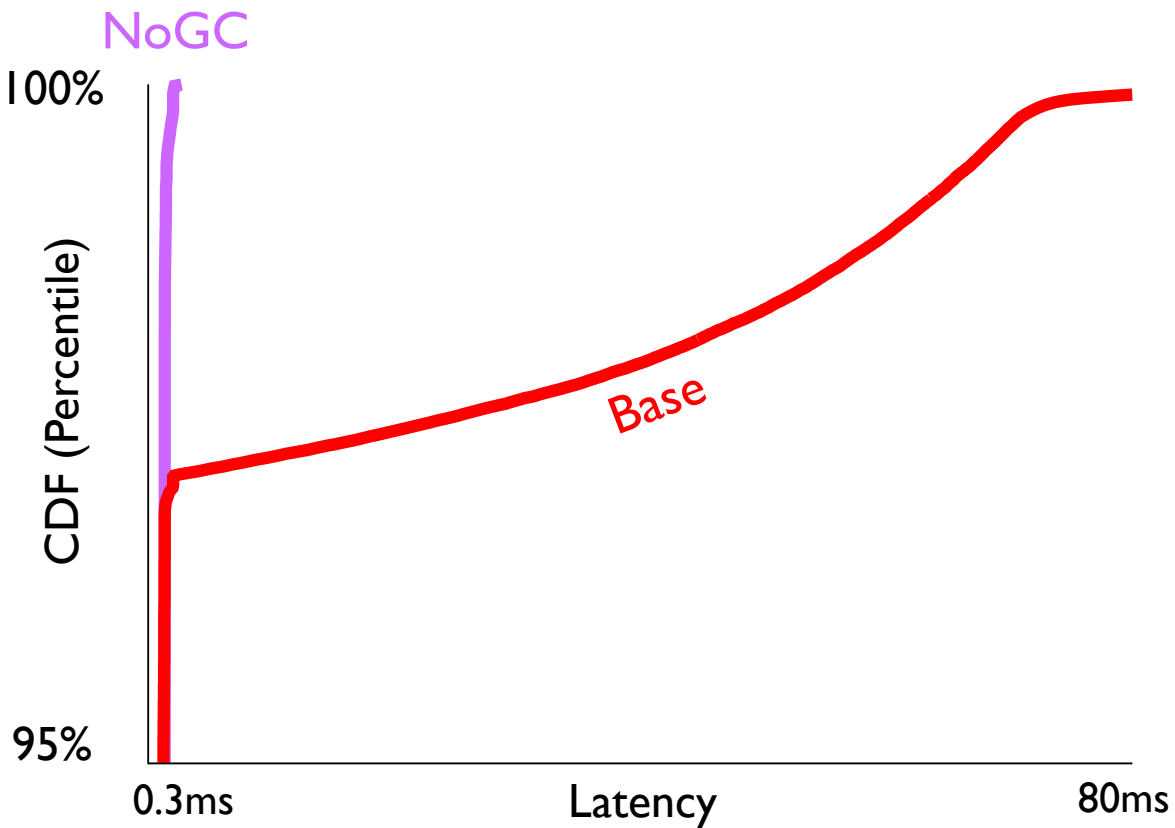II.  GC-Tolerant Read

III. Rotating GC

IV.  GC-Tolerant Flush

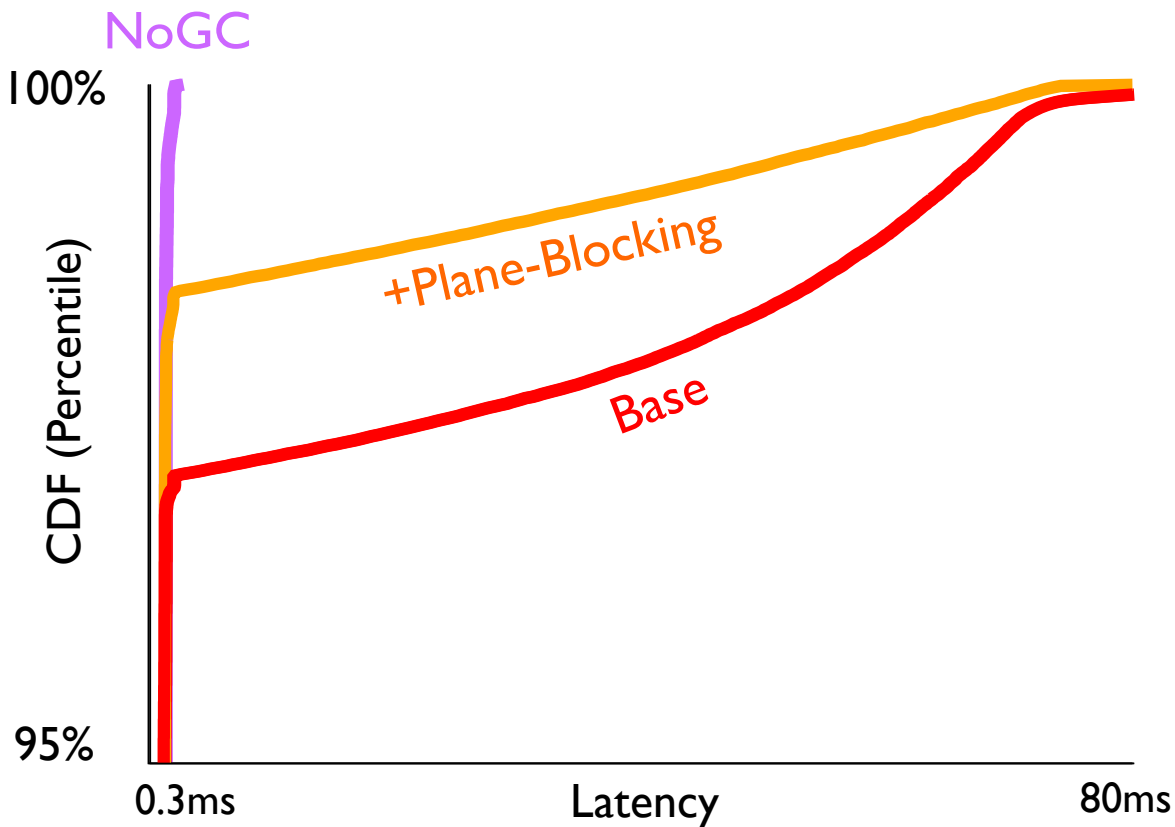*Current SSD technology:*
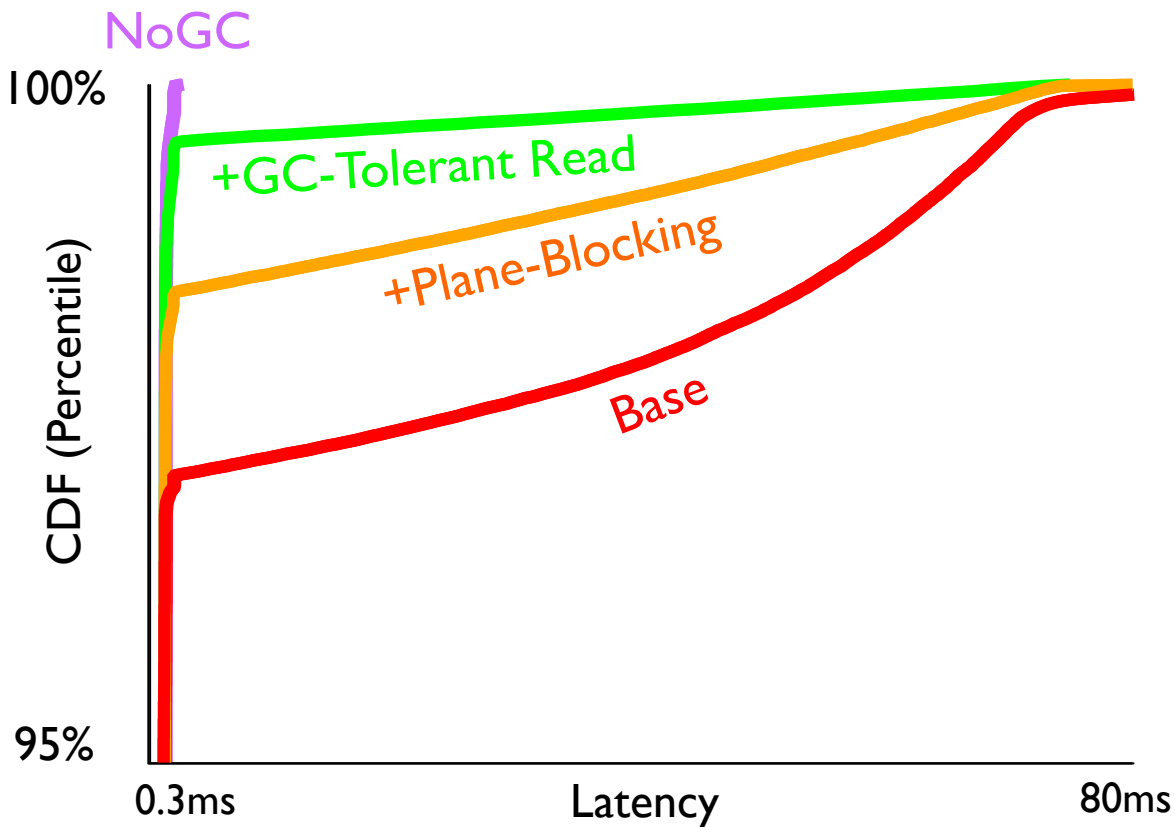
**RAIN**
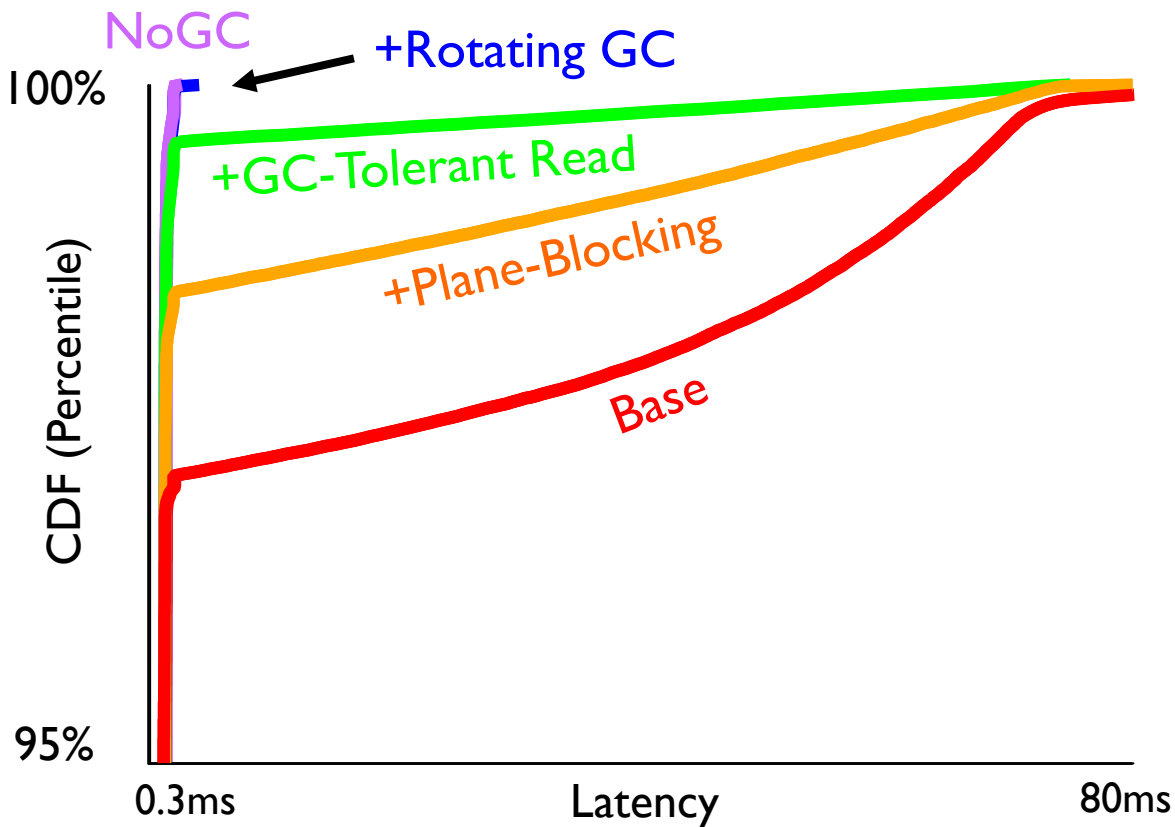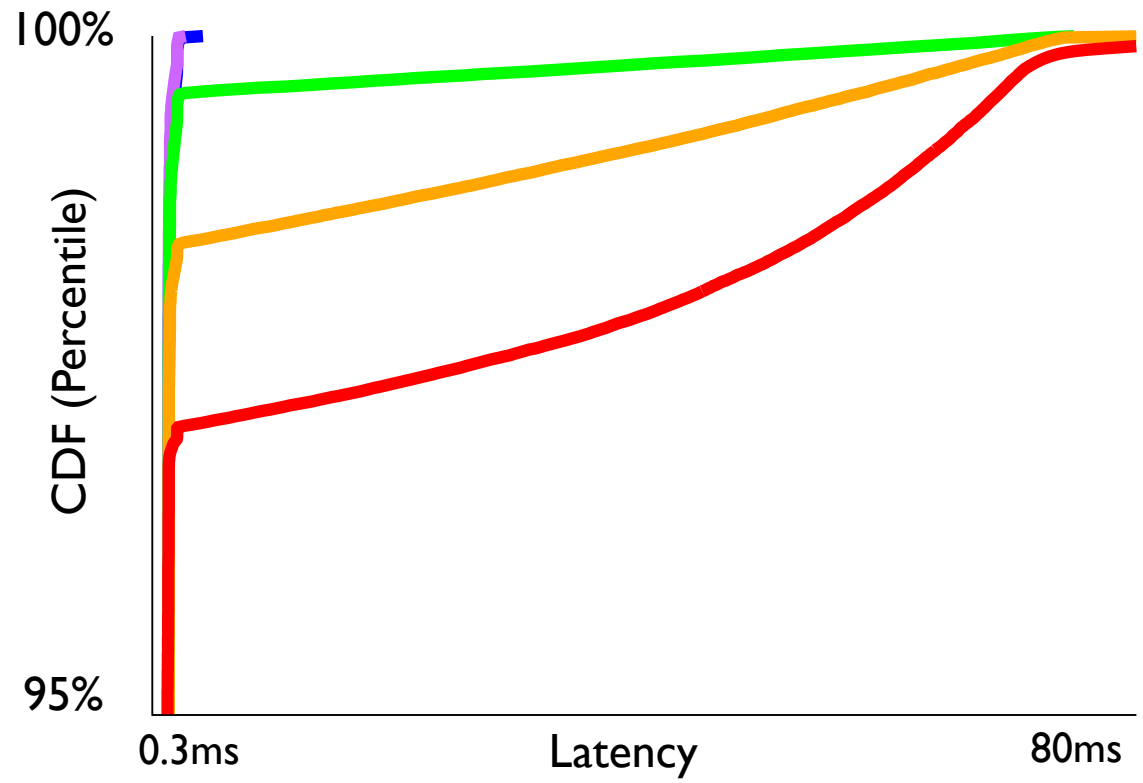(Parity-based Redundancy)

# Results

# Results

# Results



NoGC

+Plane-Blocking

Base

100%

95%

CDF (Percentile)
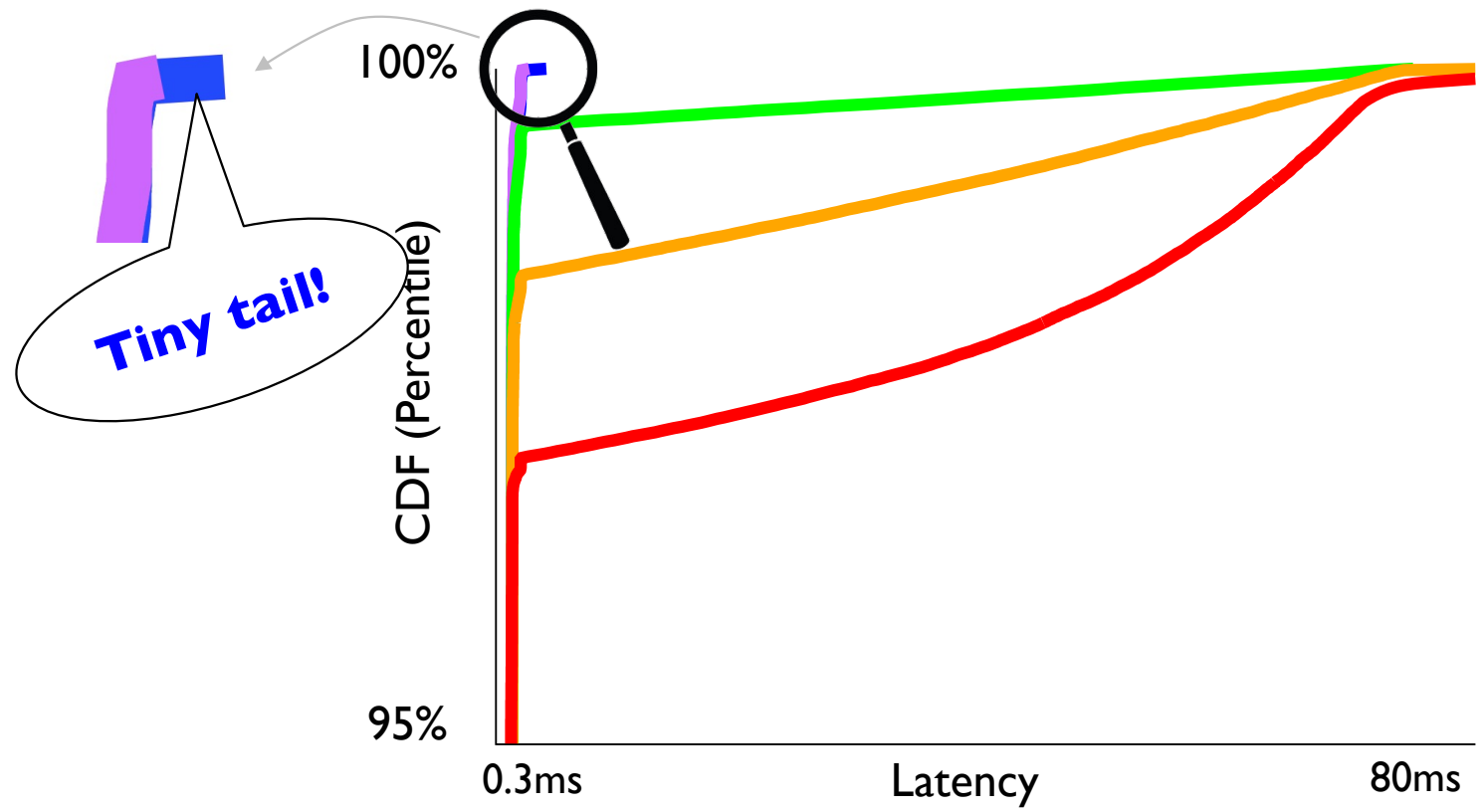
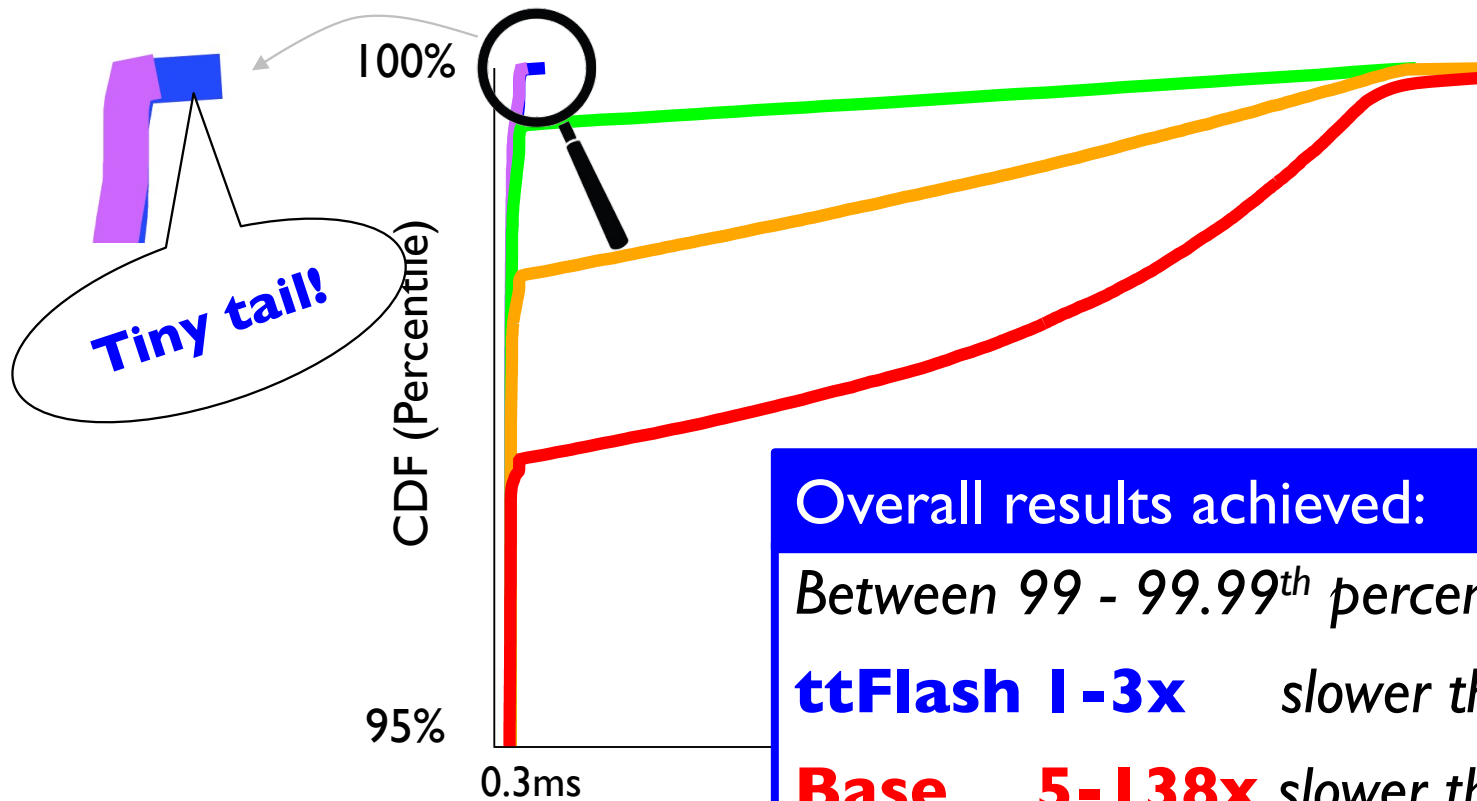0.3ms          Latency          80ms

# Results

# Results

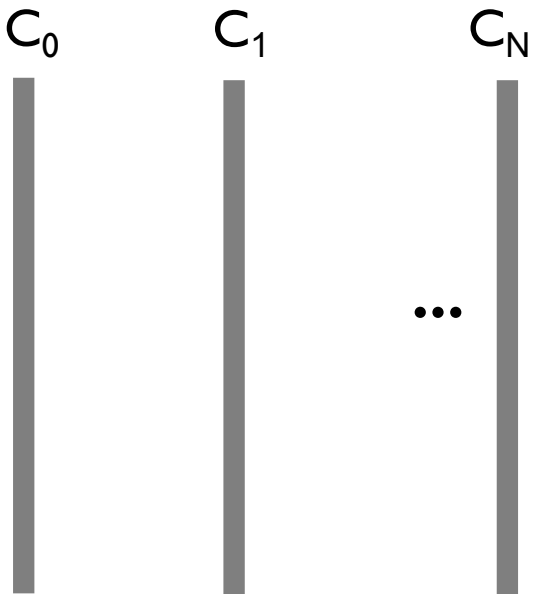# Outline

☐ Introduction

☐ **Background**

☐ Tiny-Tail Flash Design

☐ Evaluation, limitations, conclusion

# SSD Internals

# SSD Internals

$C_0$      $C_1$      $C_N$

...

# SSD Internals

# SSD Internals

# SSD Controller

Old block

Empty block

# SSD Controller

*1.* **read** *to controller (check with ECC)*

Old block    Empty block

# SSD Controller

1. **read** *to controller*
   *(check with ECC)*
2. **write** *to another block*

*Old block*    *Empty block*

# SSD Controller

for (1 … # of valid pages):
   1. **read** to controller
      (check with ECC)
   2. **write** to another block

1  2

Old block    Empty block

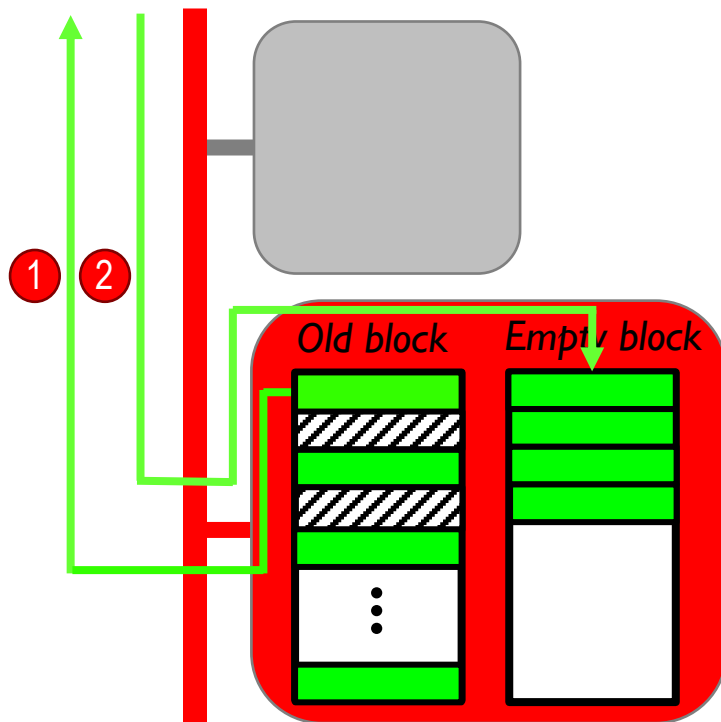# SSD Controller

for (1 … # of valid pages):
  1. **read** to controller
     (check with ECC)
  2. **write** to another block

Old block    Empty block

# SSD Controller

*3. Erase the old block*

Old block        Empty block

Erase!

SSD Controller

3. *Erase* the old block

*Old block*    *Empty block*

Erase!

*Erase* operation
**block** *the plane*

blocked!

GCing plane

blocked!

Channel blocking GC

"**Base**" approach

GCing plane

# Outline

❑Introduction

❑Background

❑**Tiny-Tail Flash Design**

- ▪ Plane-Blocking GC
- ▪ GC-Tolerant Read
- ▪ Rotating GC
- ▪ GC-Tolerant Flush

❑Evaluation, limitations, conclusion

Base: Channel Blocking

blocked!

**Plane** Blocking

GCing plane

GCing plane

# Plane Blocking

**Base GC Logic:**

*for (every valid page)*
   *1. flash read+write*
     *(over channel)*
   *2. wait*



SSD Controller

Old block    Empty block

# Plane Blocking

**Base GC Logic:**

*for (every valid page)*
    *1. flash r...*
       *(over ...*
    *2. wait*

**Plane Blocking GC Logic:**
*for (every valid page)*
    **(1)** *flash read+write*
        *(inside plane)*

**SSD Controller**

*"Intra-plane copyback"*

**(1)**

Old block    Empty block

# Plane Blocking

**Base GC Logic:**

*for (every valid page)*
   *1. flash r...*
      *(over...*
   *2. wait*

**Plane Blocking GC Logic:**
*for (every valid page)*
   **1** *flash read+write (inside plane)*

   **2** *serve other user I/Os*

## SSD Controller

*"Intra-plane copyback"*

**1**

Old block    Empty block

# Plane Blocking

**Base GC Logic:**

*for (every valid page)*
  *1. flash r...*
     *(over ...)*
  *2. wait*

**Plane Blocking
GC Logic:**
*for (every valid page)*
  **1** *flash read+write
     (inside plane)*

  **2** *serve other
     user I/Os*

## SSD Controller

**2**

Read
Page

A
B
C

*"Intra-plane
copyback"*  **1**

Old block   Empty block

# Plane Blocking

**Base GC Logic:**

for (every valid page)
 1. flash read+write
    (over ...
 2. wait

**Plane Blocking
GC Logic:**
for (every valid page)
 **1** flash read+write
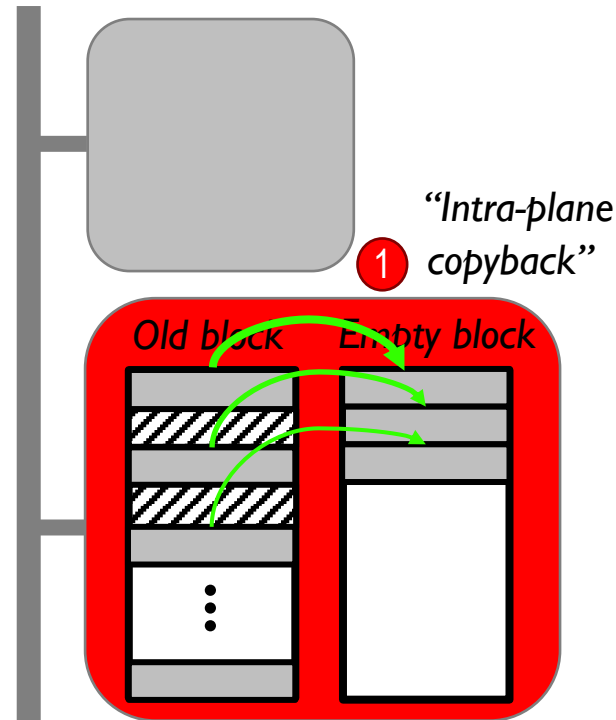    (inside plane)

 **2** serve other
    user I/Os

**SSD Controller**

**2** Read
Page

A
B
C

**1** *"Intra-plane
copyback"*

Old block    Empty block

**Overlap**
 **1** *intra-plane copyback* **with**
 **2** *channel usage for other non-GCing planes*

❑ **Issue 1:** No ECC check *for garbage-collected pages*
   ▪ *(will discuss later)*

☐ **Issue 1:** No ECC check *for garbage-collected pages*
  ▪ *(will discuss later)*

☐ **Issue 2:**



*GC-ing plane*
**still blocks**

☐ **Issue 1:** No ECC check *for garbage-collected pages*
  ▪ *(will discuss later)*

☐ **Issue 2:**



Read  X

Read  Y

Read  Z

*GC-ing plane*
**still blocks**

delayed!

# Outline

❑Introduction

❑Background

❑**Tiny-Tail Flash Design**

- Plane-Blocking GC

- **RAIN + GC-Tolerant Read**

- Rotating GC

- GC-Tolerant Flush

❑Evaluation, limitations, conclusion

# RAIN

$C_0$     $C_1$     $C_2$     $C_3$

# RAIN

$C_0$ $\qquad$ $C_1$ $\qquad$ $C_2$ $\qquad$ $C_3$

$PG_0$

$PG_1$

$PG_2$

# RAIN

LPN (Logical Page #)

# RAIN



LPN (Logical Page #)

Static mapping:
LPN0 → C[0]PG[0]
LPN1 → C[1]PG[0]
…

# RAIN

LPN (Logical Page #)

C_0          C_1          C_2          C_3

$PG_0$

$PG_1$

$PG_2$

In the grid: $PG_0$ row contains 0, 1, 2, $P_{0,1,2}$

Static mapping:
LPN0 → C[0]PG[0]
LPN1 → C[1]PG[0]
…

Add parity:

LPN 0, 1, 2 → $P_{0,1,2}$

# RAIN



LPN (Logical Page #)

Static mapping:
LPN0 → C[0]PG[0]
LPN1 → C[1]PG[0]
…

Add parity:

LPN 0, 1, 2 → $P_{0,1,2}$

Rotating parity as RAID 5

# RAIN enables GC-Tolerant Read

*Full Stripe Read*

| 0 | 1 | 2 |
|---|---|---|

# **RAIN** enables **GC-Tolerant Read**



*Full Stripe Read*

# **RAIN** enables **GC-Tolerant Read**



*Full Stripe Read*

$\mathbf{2} = XOR$
$(0, 1, \mathbf{P_{0,1,2}})$

*fast*

**tail**

0

1

2

$\mathbf{P_{0,1,2}}$

GC

# **RAIN** enables **GC-Tolerant Read**

# **RAIN** enables **GC-Tolerant Read**

# GC-Tolerant Read
## *Issue: **partial** stripe read*

# GC-Tolerant Read

## Issue: *partial* stripe read

Partial stripe read:  2

**slow!**

# GC-Tolerant Read

## Issue: ***partial*** *stripe read*

*Partial stripe read:*  2

**slow!**

*Must generate extra*
**N-1 reads!**

# GC-Tolerant Read

## Issue: ***partial*** *stripe read*



*Partial stripe read:* 2 ← 2 = XOR

**slow!**

$(0, 1, P_{0,1,2})$

*Must generate extra* **N-1 reads!**

0    1    2    $P_{0,1,2}$

# GC-Tolerant Read

## *Issue: **partial** stripe read*



*Partial stripe read:*   $2$   ← $2 = XOR$

**slow!**   $(0, 1, P_{0,1,2})$

$0$   $1$   $2$   $P_{0,1,2}$

*Must generate extra*
**N-1 reads!**

*Add **contention** to other*
N -1 *channels and planes*

# GC-Tolerant Read

## Issue: *partial* stripe read



Partial stripe read:  2  ← 2 = XOR (0, 1, $P_{0,1,2}$)

**slow!**

0    1    2    $P_{0,1,2}$

Must generate extra
**N-1 reads!**

Add **contention** to other
N-1 channels and planes

Convert to full stripe if:
$T_{extra-reads} < T_{GC}$

# Issue: **more than 1 GCs**
# in a plane group?

*Full-stripe read*  | 0 | 1 | 2 |

# Issue: **more than 1 GCs**
# in a plane group?

*One parity → cut one tail*
**Can't cut two tails!**

*Full-stripe read*  | 0 | 1 | 2 |

# Issue: **more than 1 GCs** in a plane group?

*One parity → cut one tail*
**Can't cut two tails!**

*Full-stripe read*  | 0 | 1 | 2 |

# Issue: **more than 1 GCs** in a plane group?

*One parity → cut one tail*
**Can't cut two tails!**

# Issue: **more than 1 GCs** in a plane group?

*One parity → cut one tail*
**Can't cut two tails!**

*Full-stripe read*

# Issue: **more than 1 GCs** in a plane group?

*One parity → cut one tail*
**Can't cut two tails!**

*Full-stripe read*  | 0 | 1 | 2 |

**2 tails!**

*DOES NOT HELP!*

$PG_0$

0    1    2    $P_{0,1,2}$

GC    **GC**

# **Outline**

❑Introduction

❑Background

❑**Tiny-Tail Flash Design**

- ▪ Plane-Blocking GC
- ▪ GC-Tolerant Read
- ▪ **Rotating GC**
- ▪ GC-Tolerant Flush

❑Evaluation, limitations, conclusion

# Rotating GC:
*Anytime, **at most 1** plane per plane group can perform GC*

**Rotating GC:**

*Anytime,* **at most 1** *plane per plane group can perform GC*

# Rotating!



PG$_0$

| 0 | 1 | 2 | P$_{0,1,2}$ |

# Rotating GC:

*Anytime,* **at most 1** *plane per plane group can perform GC*

**Rotating!**



# Rotating GC:

*Anytime, **at most 1** plane per plane group can perform GC*

# Rotating!



PG$_0$

0  1  2  P$_{0,1,2}$

# Rotating GC:

*Anytime,* **at most 1** *plane per plane group can perform GC*

# Rotating GC:

*Anytime, **at most 1** plane per plane group can perform GC*

*Concurrent GCs in **different** PGs are permitted.*

Tiny tail!

Why still tiny tails?

Small/partial-stripe read
→ *Sometimes may be better to* **wait for GC** *than adding extra reads/contentions!*

+Rotating GC

CDF (Percentile)

100%

95%

0.3ms          Latency          80ms

# Outline

❏ **Tiny-Tail Flash Design**

- Plane-Blocking GC
- GC-Tolerant Read
- Rotating GC
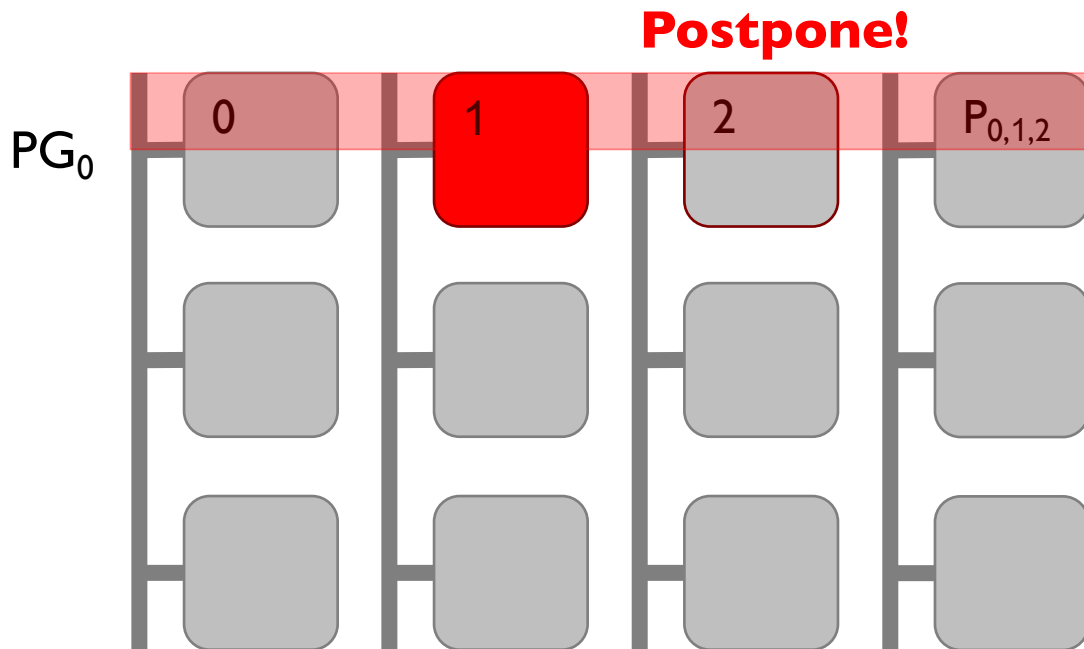- GC-Tolerant Flush (in paper)

❏ Evaluation

❏ Limitations

❏ conclusion

# Implementation

- ❏ **SSDsim** (*~2500 LOC*)
  - ▪ *Device simulator*

- ❏ **VSSIM** (*~900 LOC*)
  - ▪ *QEMU/KVM-based*
  - ▪ *Run Linux and applications*

- ❏ **OpenSSD**
  - ▪ *Many limitations of the simple programming model*

- ❏ Future: ttFlash on **OpenChannel SSD**

# Evaluation

❑ Simulator: **SSDsim** (verified against hardware)

❑ Workload: 6 real-world traces from Microsoft Windows

❑ Settings and SSD parameters:

▪ SSD size: 256GB, plane group width = 8 planes (1 parity, 7 data)

| Sizes | | Latencies | |
|---|---|---|---|
| SSD Capacity | 256 GB | Page Read | $40\mu s$ |
| #Channels | 8 | (flash-to-register) | |
| #Planes/channel | 8 | Page Write | $800\mu s$ |
| Plane size | 4 GB | (register-to-flash) | |
| #Planes/chip | ** 1 | Page data transfer | $100\mu s$ |
| #Blocks/plane | 4096 | (via channel) | |
| #Pages/block | 256 | Block erase | 2 ms |
| Page size | 4 KB | | |

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace

# Developer Tools Release Server Trace



NoGC

+Rotating GC

100%

+GC-Tolerant Read

+Plane-Blocking

Base

CDF (Percentile)

95%

0.3ms     Latency     80ms

ttFlash

99.99th

Result:

99.99th percentile:

**ttFlash 3x**  *slower than NoGC*

**Base    138x** *slower than NoGC*

# Evaluated on 6 windows workload traces with various characteristics



Display Ads Server (DAPPS)  Dev. Tools Release (DTRS)  Exchange Server (Exch)

Live Maps Server (LMBE)  MSN File Server (MSNFS)  TPC-C

*Reduced blocked I/Os (total) from 2 – 7% to 0.003 – 0.05%*
*99 – 99.99%: 1.0 – 2.6x slower for ttFlash and 5.6 – 138.2x for Base*

# Other Evaluations

# Other Evaluations

❑ Filebench on VSSIM+ttFlash

  ▪ *ttFlash achieves better average latency than base case*

# Other Evaluations

❑ Filebench on VSSIM+ttFlash

  ■ *ttFlash achieves better average latency than base case*

❑ Vs. Preemptive GC

  ■ *ttFlash is more stable than semi-preemptive GC*

    – *(If no idle time, preemptive GC will create GC backlogs, creating latency spikes)*



Base  +PB  +GTR  +RGC

Avg Latency (ms)

File Server | Network FS | OLTP | Varmail | Video Server | Web Proxy

ttFlash
Preempt

Latency (s)

4386    Elapsed time (s)    4522

ttflash stable

# Tradeoffs/Limitations

# Tradeoffs/Limitations

❑ ttFlash depends on RAIN

- ▪ *1 parity for N parallel pages/channels*
- ▪ *We set N = 8, so we lose one channel out of 8 channels.*
- ▪ *Average latencies are* **1.09 – 1.33x** *slower than NoGC, No-RAIN case*

# Tradeoffs/Limitations

❑ ttFlash depends on RAIN

- *1 parity for N parallel pages/channels*
- *We set N = 8, so we lose one channel out of 8 channels.*
- *Average latencies are* **1.09 – 1.33x** *slower than NoGC, No-RAIN case*

❑ RAID → more writes (P/E cycles)

- *ttFlash increases P/E cycles by* **15 – 18%** *for most of workloads*
- *Incur > 53% P/E cycles for TPCC, MSN (random write)*

# Tradeoffs/Limitations

❑ ttFlash depends on RAIN

- *1 parity for N parallel pages/channels*
- *We set N = 8, so we lose one channel out of 8 channels.*
- *Average latencies are* **1.09 – 1.33x** *slower than NoGC, No-RAIN case*

❑ RAID → more writes (P/E cycles)

- *ttFlash increases P/E cycles by* **15 – 18%** *for most of workloads*
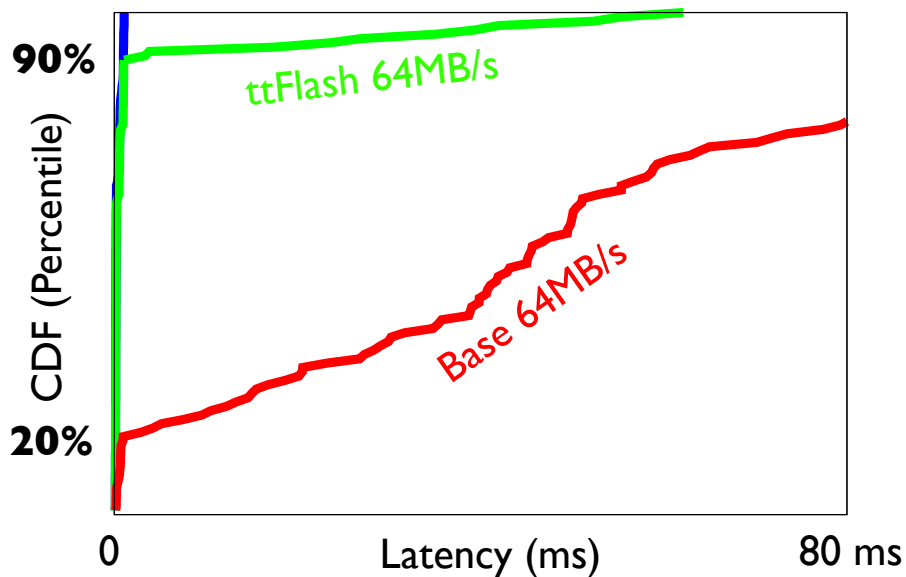- *Incur > 53% P/E cycles for TPCC, MSN (random write)*

❑ ECC is **not** checked during GC

- *Suggest background scrubbing (read is fast & not as urgent as GC)*
- *Important note: in ttFlash, foreground/user reads are still ECC checked*

# Tails under Write Bursts



Latency CDF w/ Write Bursts

ttFlash 55MB/s

ttFlash 64MB/s

Base 64MB/s

90%

20%

CDF (Percentile)

0

Latency (ms)

80 ms

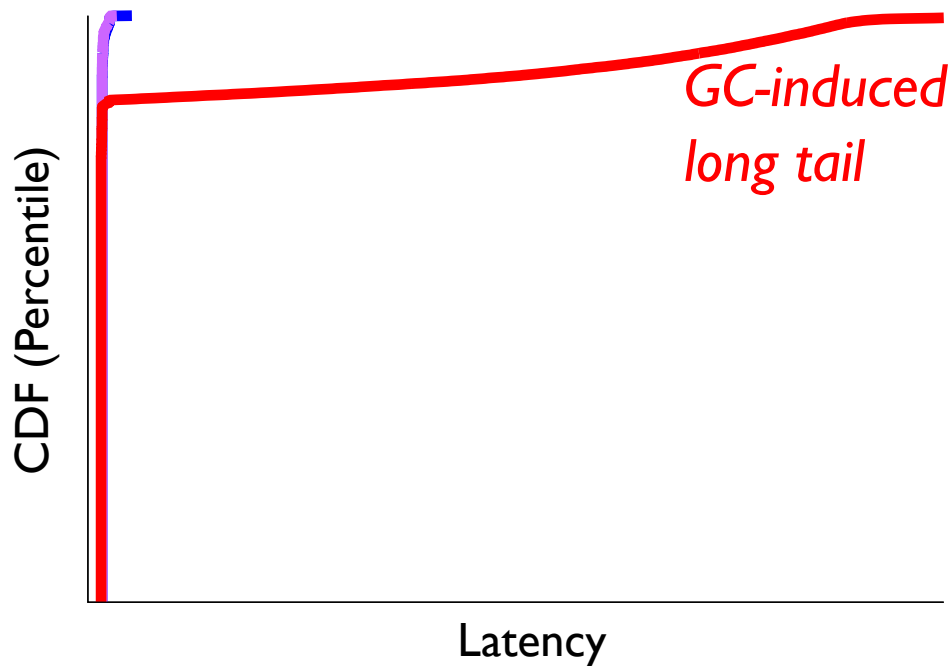*Under **write burst** and **at high watermark**, ttFlash must dynamitcally **disable** Rotating GC to ensure there is always enough number of free pages.*

# Conclusion

# Conclusion



*GC-induced long tail*

CDF (Percentile)

Latency

# Conclusion



*GC-induced long tail*

CDF (Percentile)

Latency
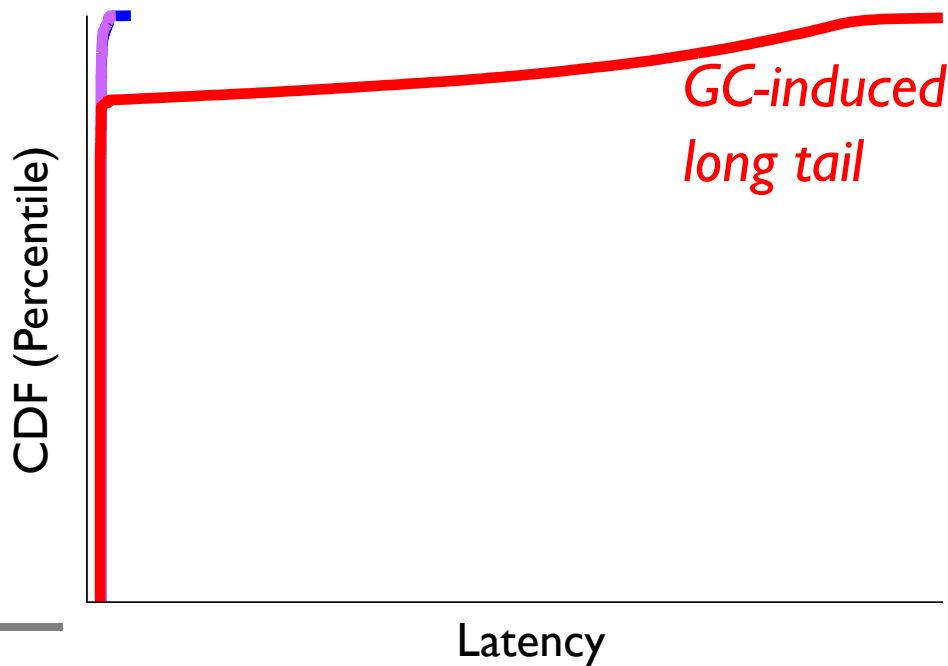
*technology:* *Powerful Controller*
**RAIN** (parity-based redundancy)
*Capacitor-backed RAM*

# Conclusion

*New techniques:*

I.  Plane-Blocking GC

II.  GC-Tolerant Read

III.  Rotating GC

IV.  GC-Tolerant Flush

*technology:* *Powerful Controller*
**RAIN** (parity-based redundancy)
*Capacitor-backed RAM*



*GC-induced long tail*

CDF (Percentile)

Latency
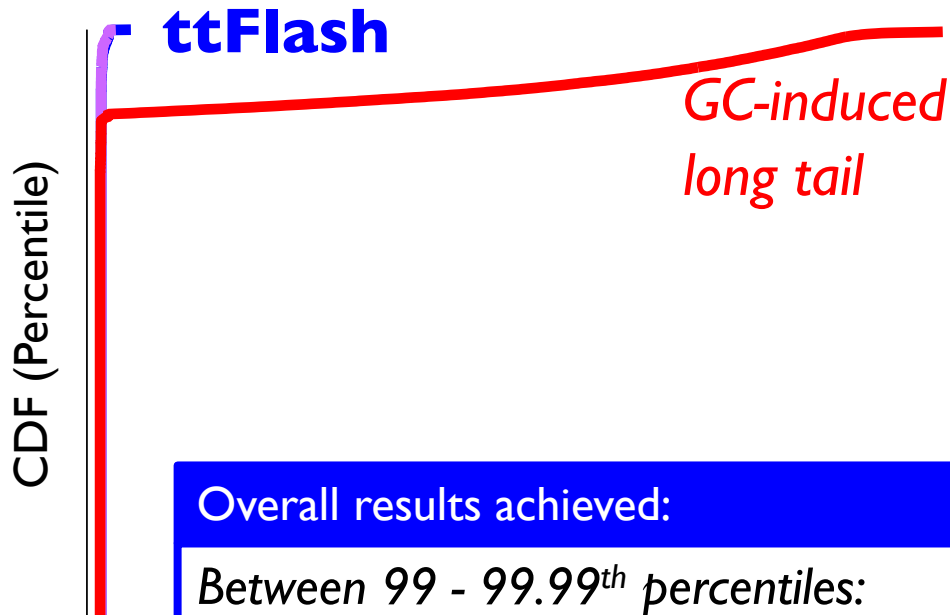
# Conclusion

*New techniques:*

   I.     Plane-Blocking GC

   II.    GC-Tolerant Read

   III.   Rotating GC

   IV.   GC-Tolerant Flush

*technology:* *Powerful Controller*
            **RAIN** (parity-based redundancy)
            *Capacitor-backed RAM*

**- ttFlash**

*GC-induced long tail*

CDF (Percentile)

Overall results achieved:

*Between 99 - 99.99$^{th}$ percentiles:*
**ttFlash 1-3x**   *slower than NoGC*
**Base**   **5-138x** *slower than NoGC*

# Thank you!
# Questions?

http://ucare.cs.uchicago.edu

CERES
Center for Unstoppable Computing

https://ceres.uchicago.edu