

Manylogs

Improving CMR/SMR Disk
Bandwidth & Latency

Tirat Patana-anake, Vincentius Martin[†],
Nora Sandler, Cheng Wu, and Haryadi S. Gunawi

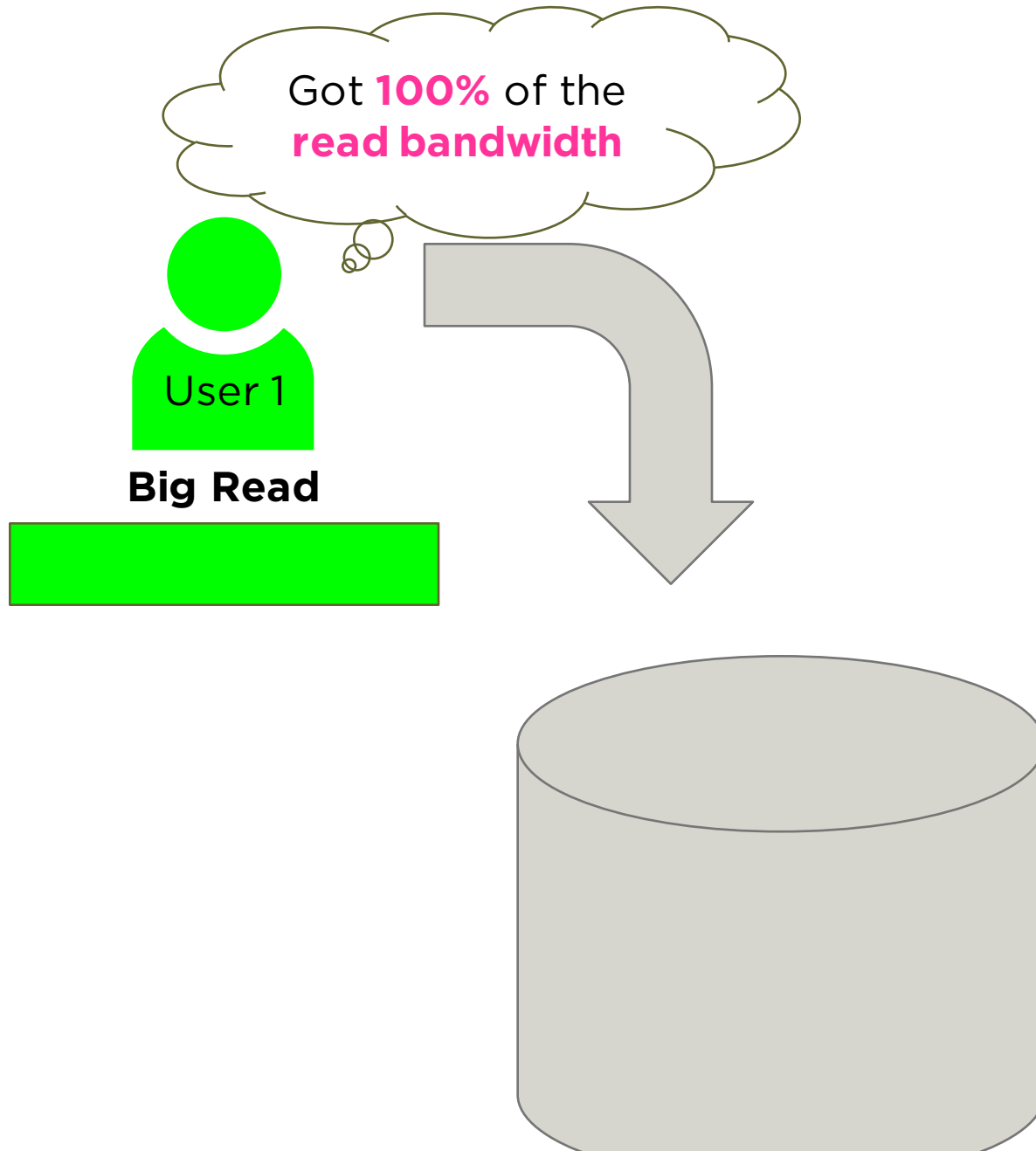


THE UNIVERSITY OF
CHICAGO

†

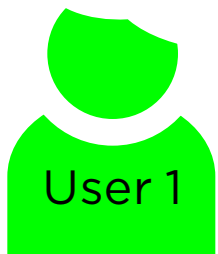


SURYA
UNIVERSITY



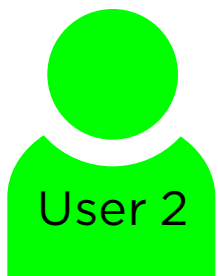


Got **50%** of the read bandwidth



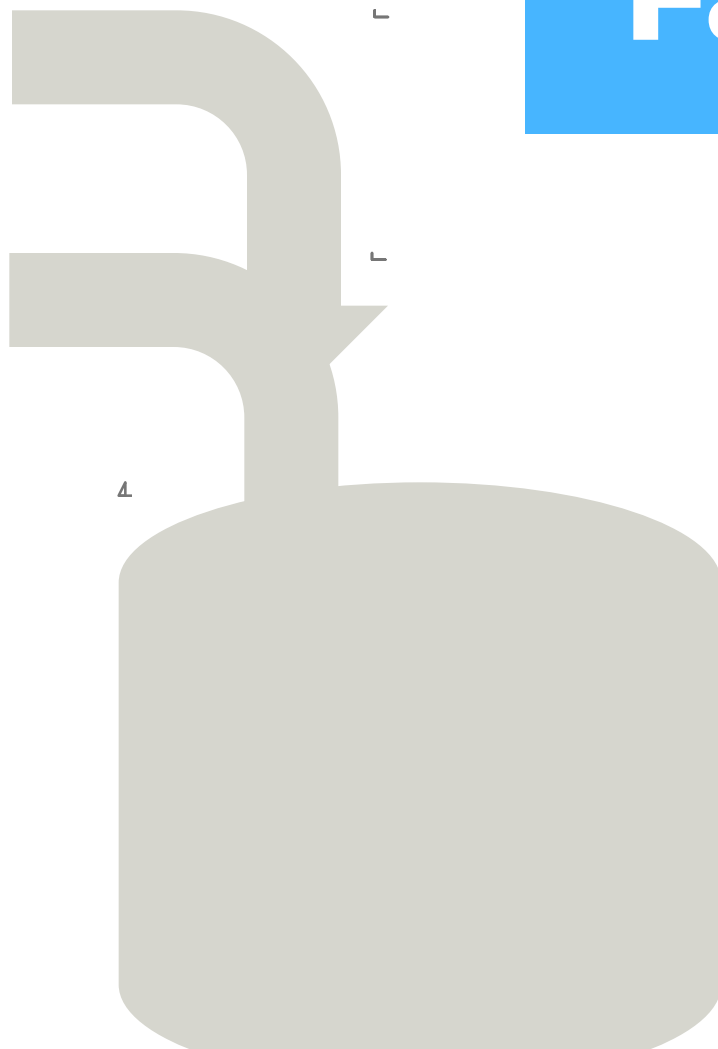
User 1

Big Read

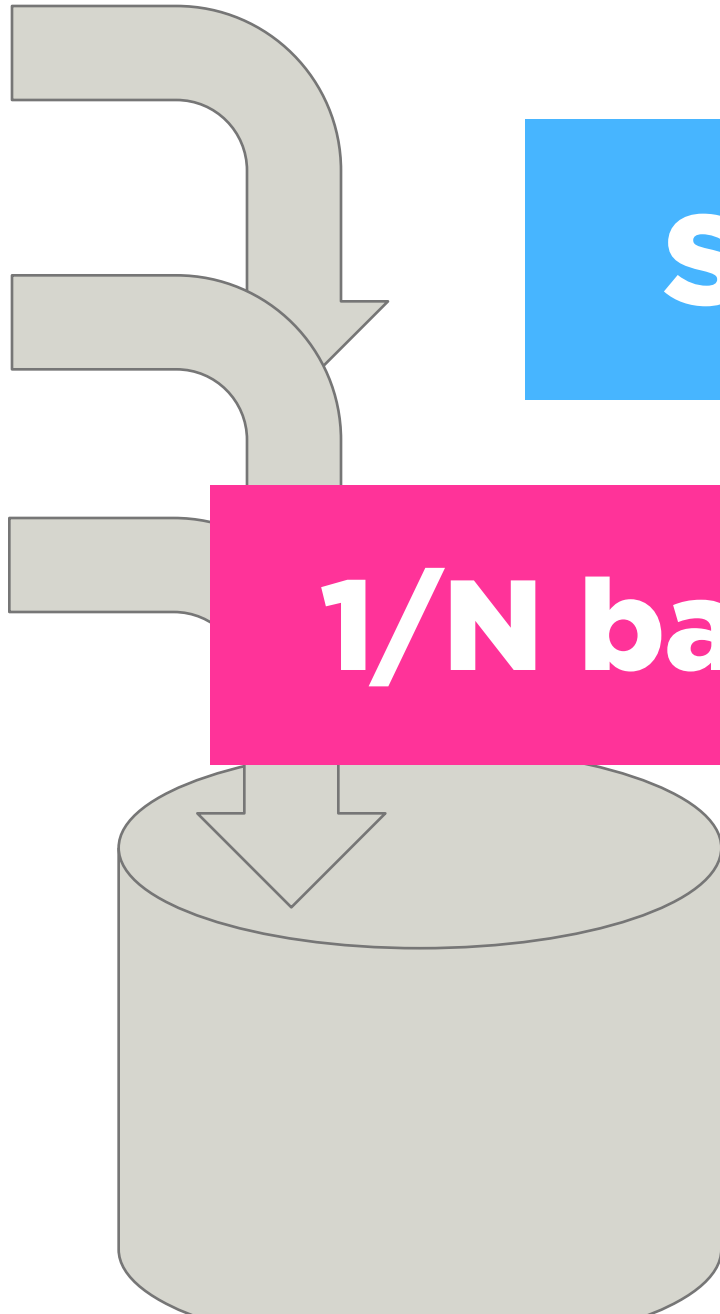
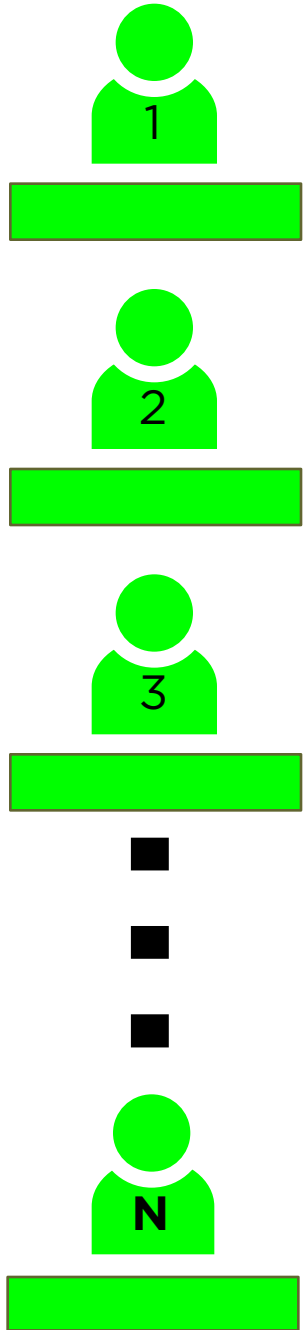


User 2

Big Read



Fair Share



Still Fair!

$1/N$ bandwidth

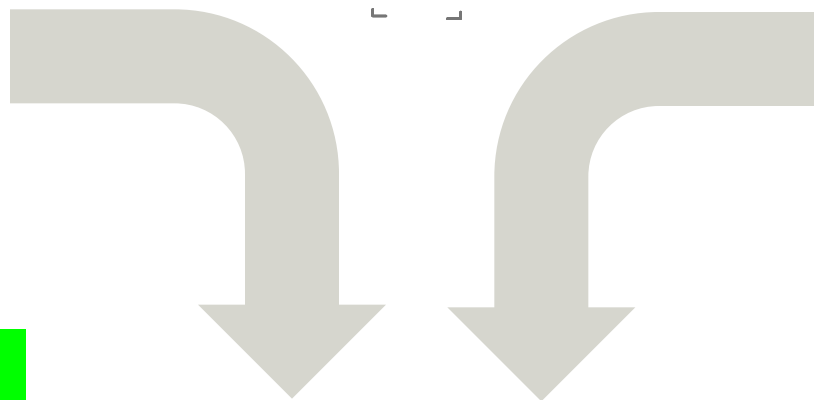
Our

Oh no!
5% bandwidth?!



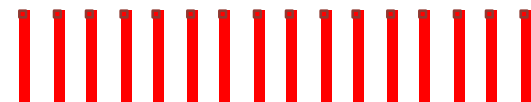
User 1

Big Read



User 2

Small Durable Writes

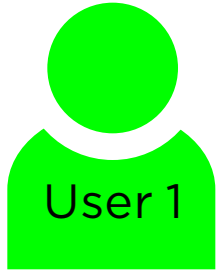


Our

More
Bandwidth
Please!

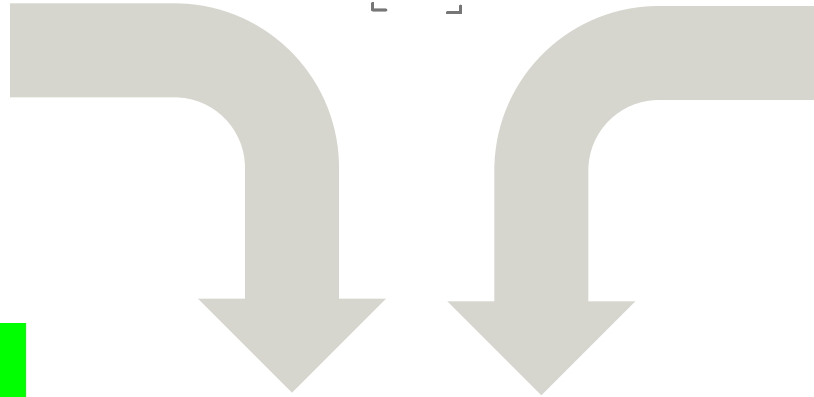
n

Faster
Latency
Please!



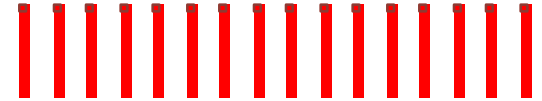
User 1

Big Read

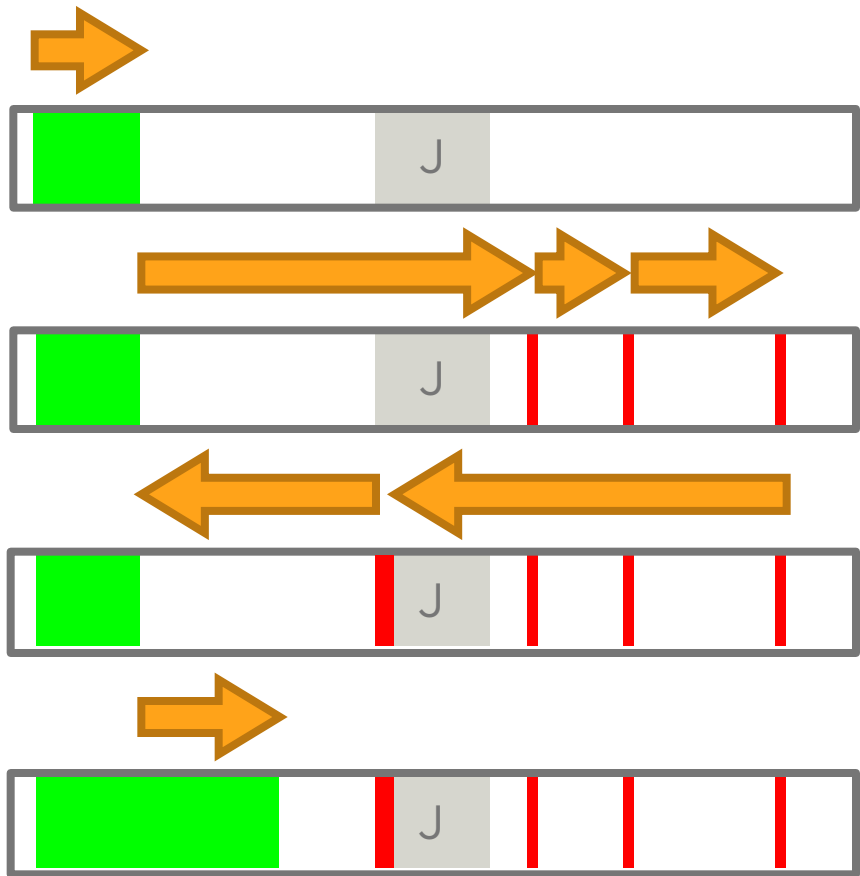


User 2

Small Durable Writes



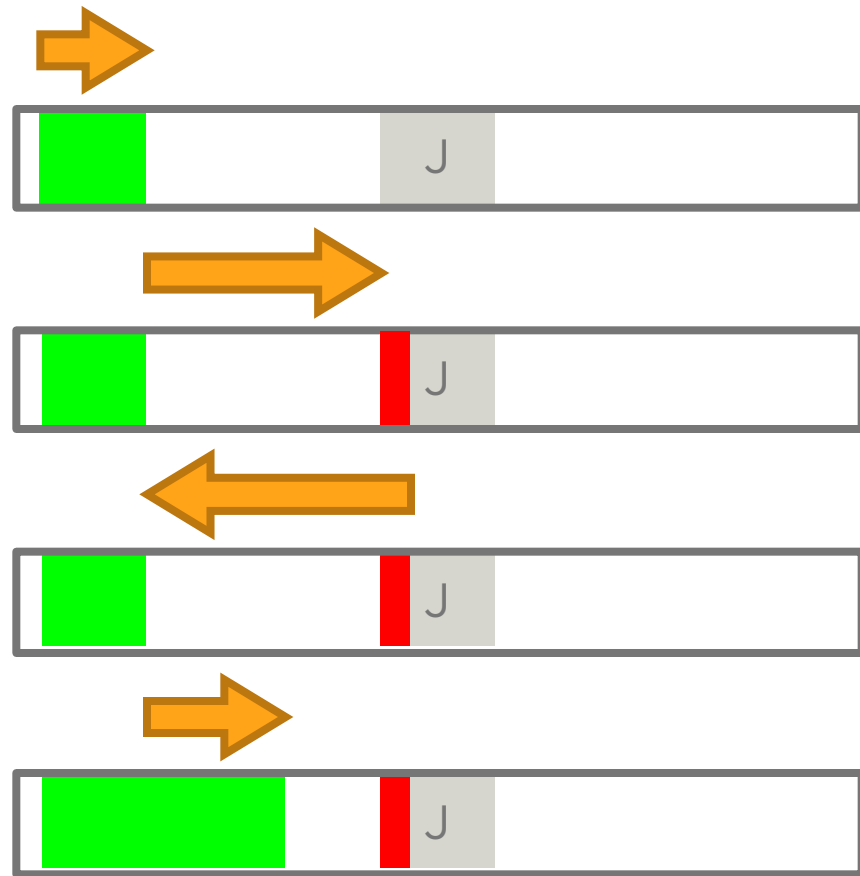
Ordered Journaling



Journal

Big Read

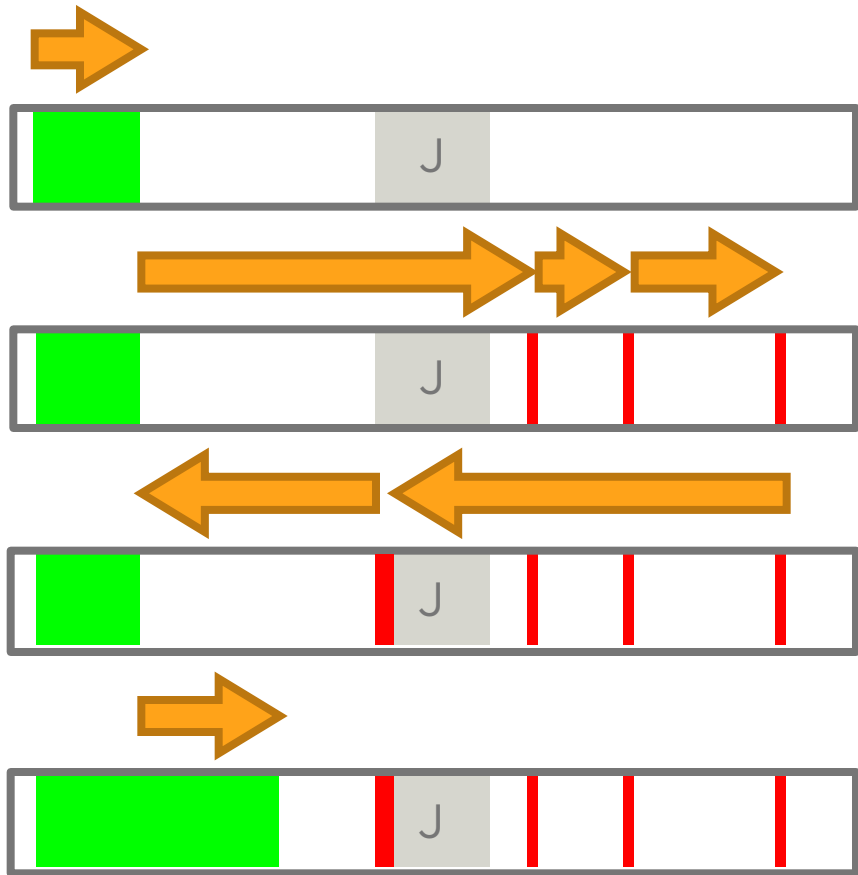
Data Journaling



Small Writes

Seek

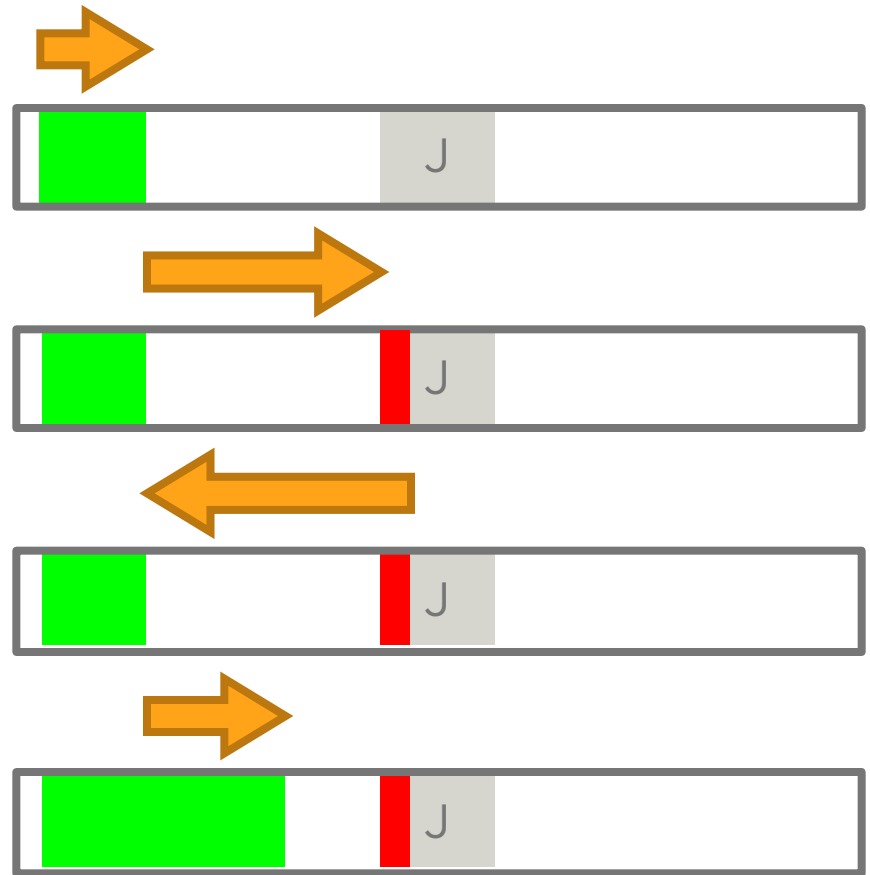
Ordered Journaling



Journal

Big Read

Data Journaling



Small Writes

Seek

Problems with Current Journaling

Ordered Journaling



Data Journaling



Both cannot handle random writes efficiently!

Problems with Current Journaling

Ordered Journaling



Data Jo

First Write



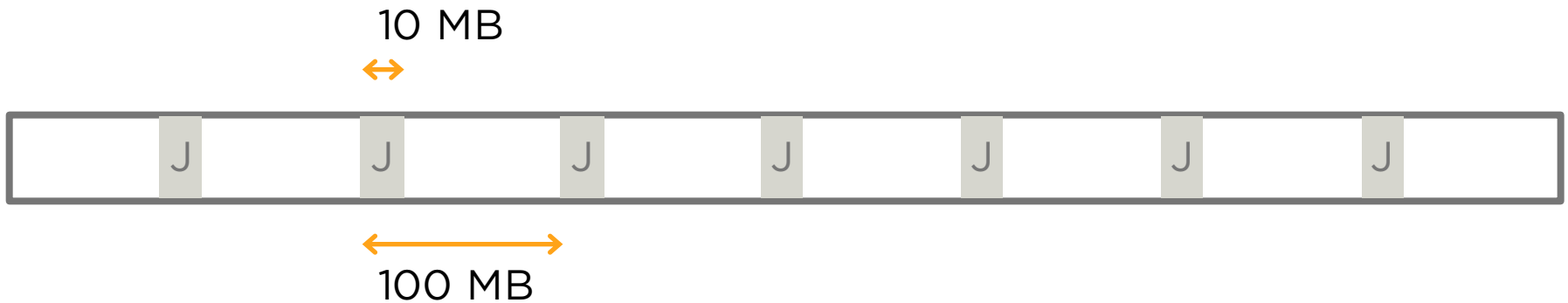
Both cannot handle random writes efficiently!

Introducing Manylogs

Single Log

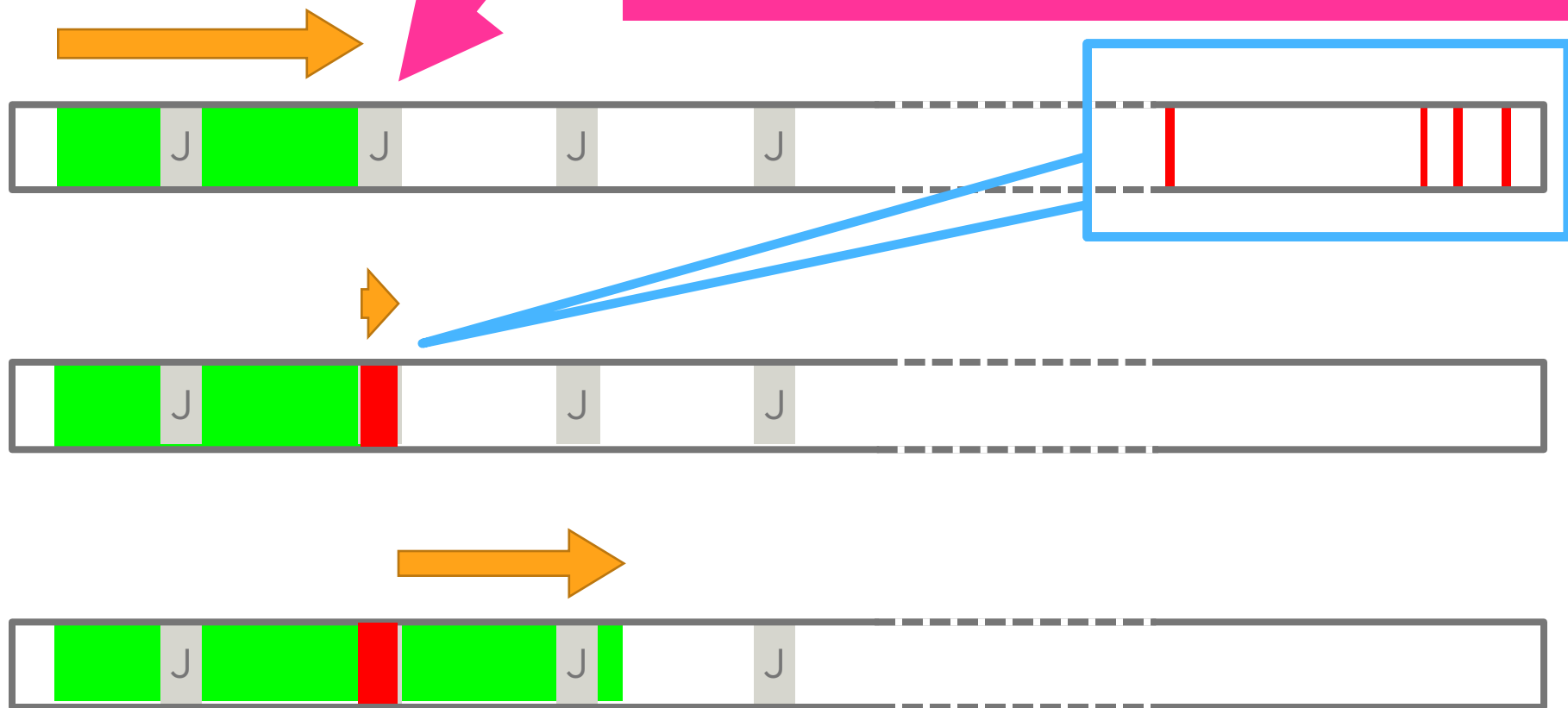


Manylogs



Manylogs

Small writes made durable to the nearest log without seeking



Journal

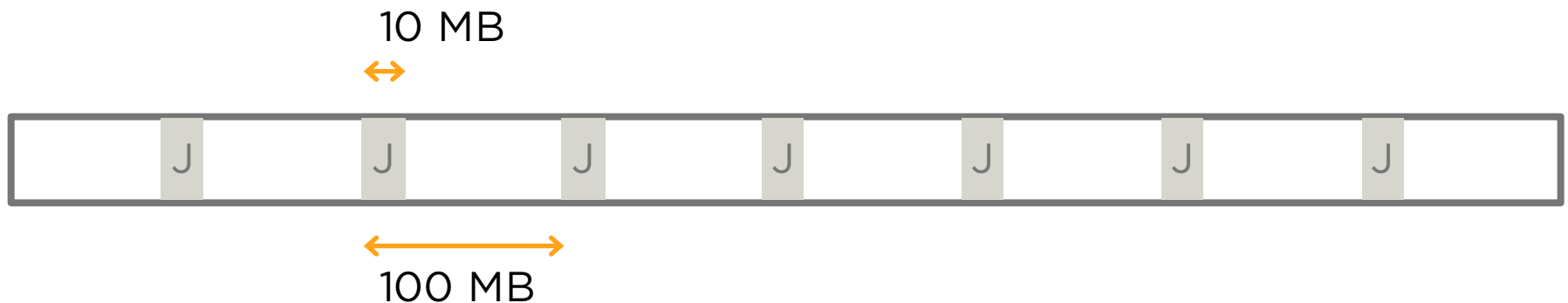
Big Read

Small Writes

Seek

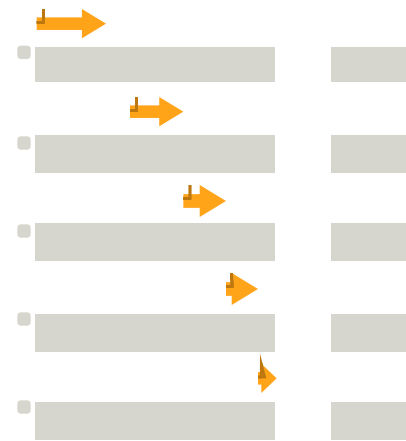
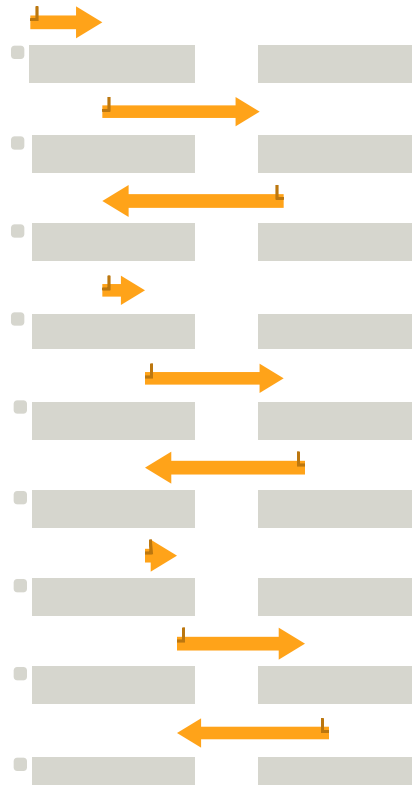
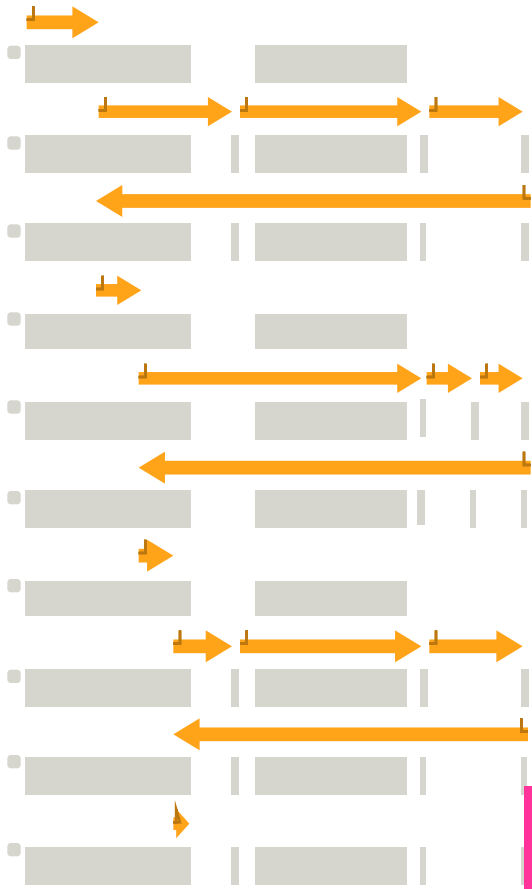
Manylogs

- ❑ Reserved log spaces uniformly across the disk
 - 10 MB every 100 MB
- ❑ Follow the disk head (last big I/O)
- ❑ Redirect Small Writes (e.g. ≤ 256 KB)
 - Nearest log: log closest to last big I/O
- ❑ Sequential Writes are left untouched



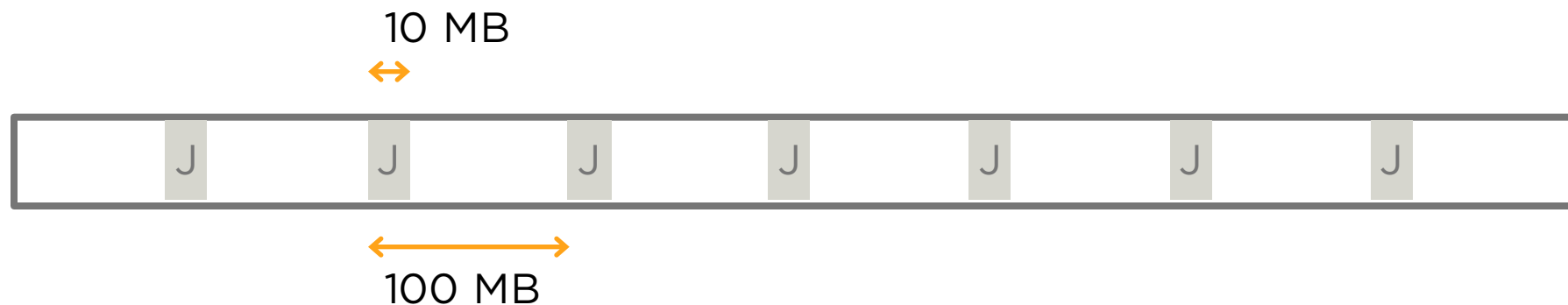
Increased Read Throughput

Manylogs

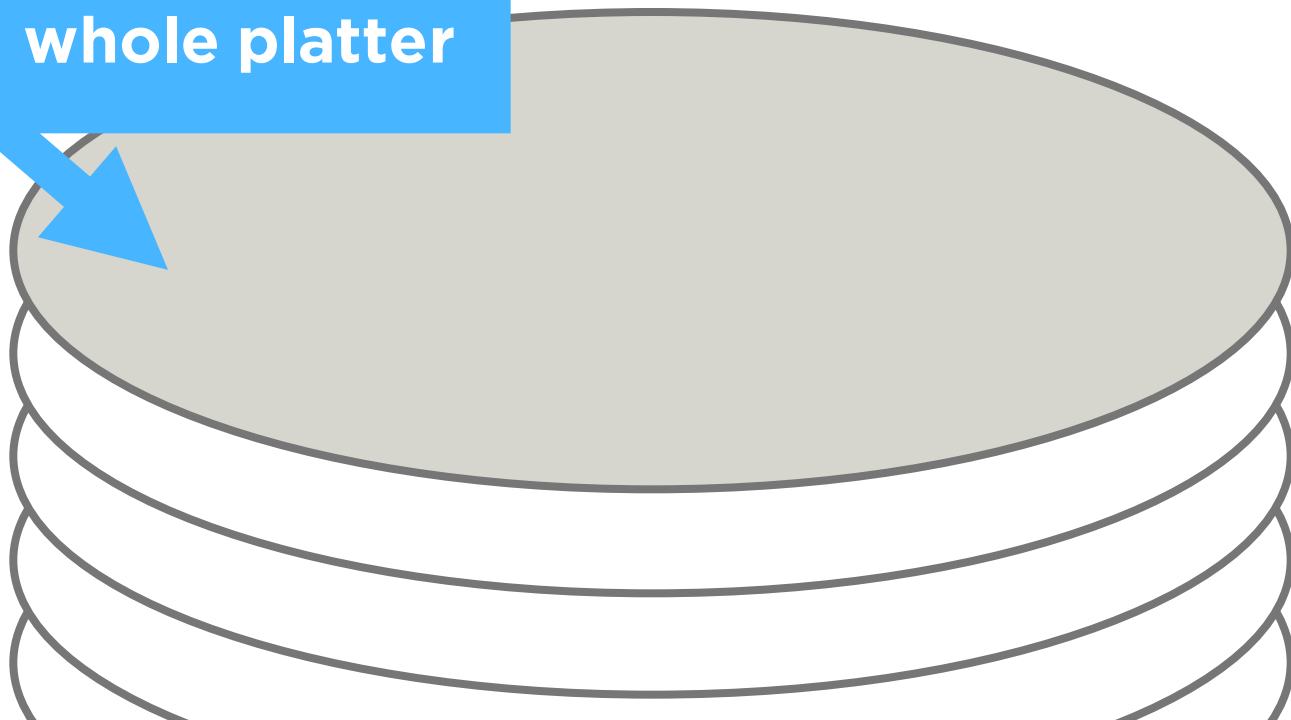


Reduced Write Latency

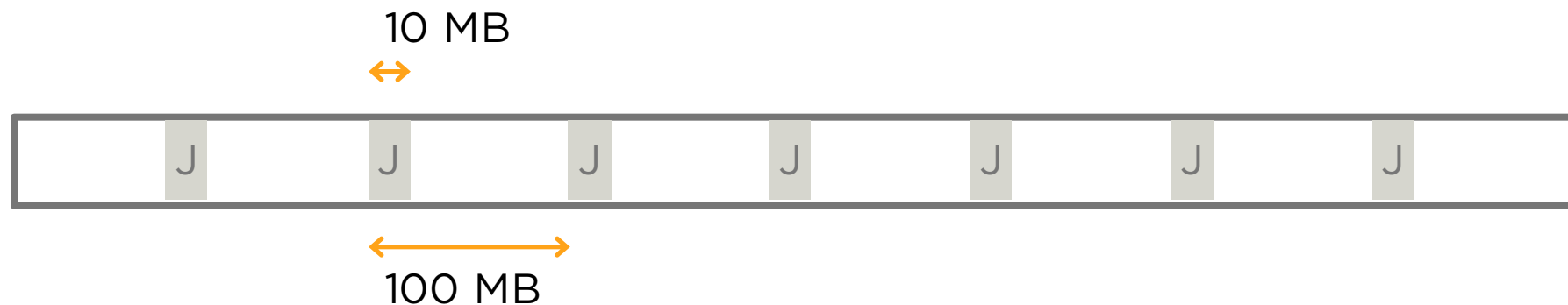
Where are logs on the disk?



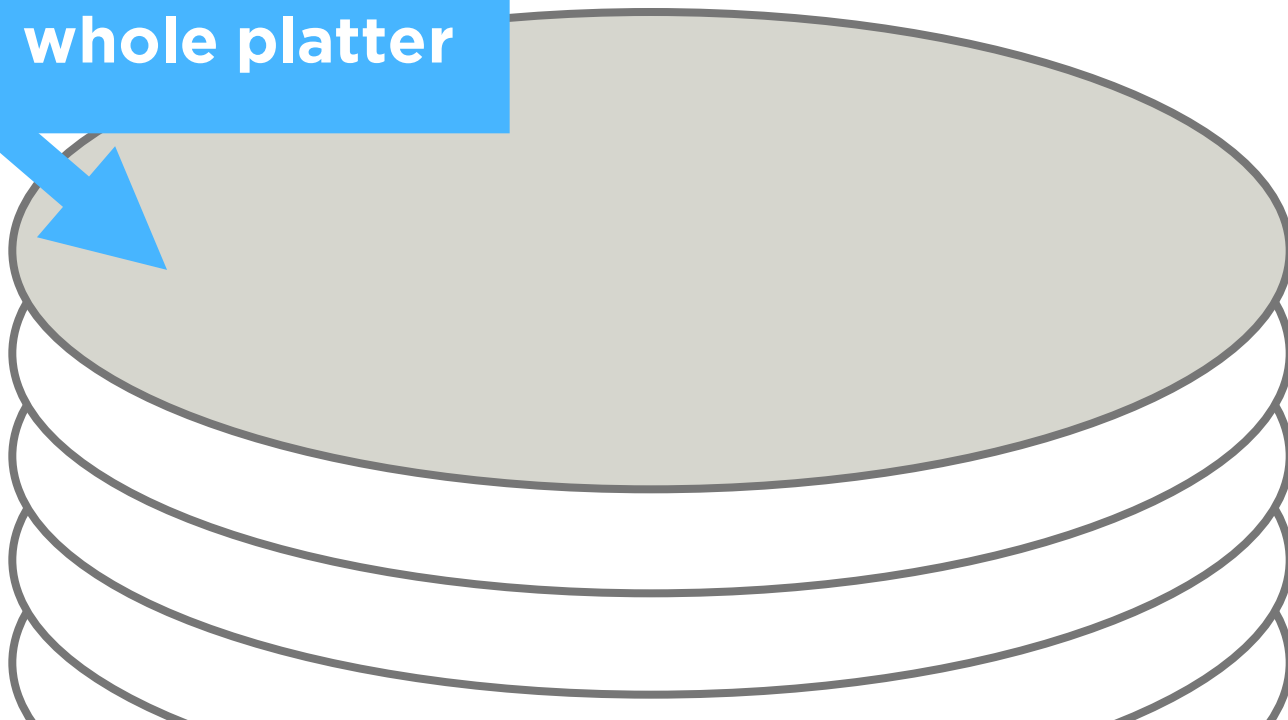
The log space = whole platter



Where are logs on the disk?



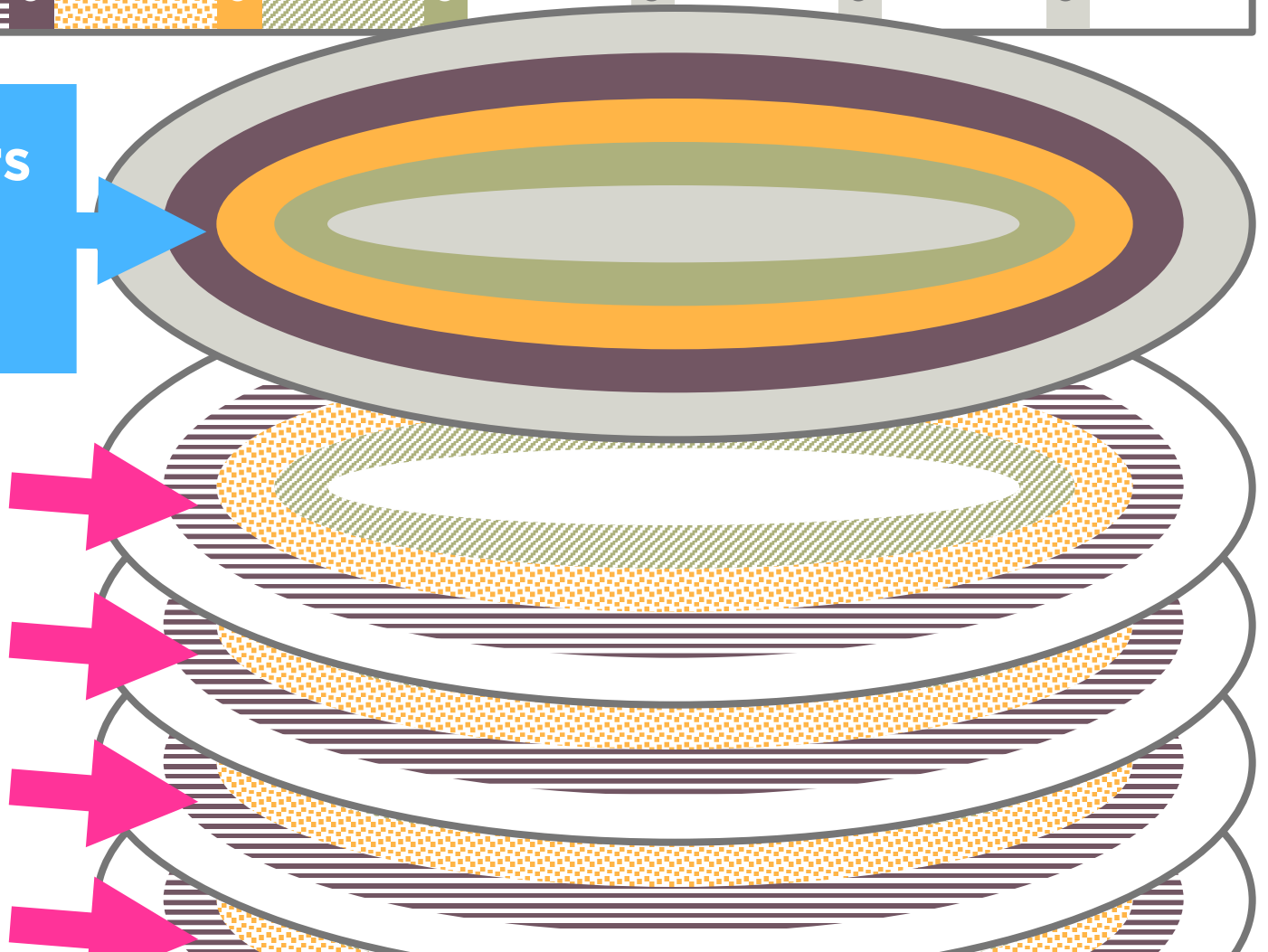
The log space = whole platter



Same cylinder = No seek!

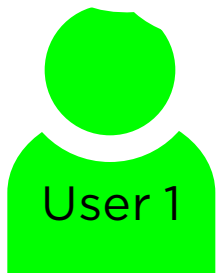


Log for others
in the same
cylinder



Ratio of Max Read Bandwidth

Latency (ms)



User 1

User 2

128MB Sequential Reads



4KB Random Writes



At different intensities

Ordered vs. Data vs. Adaptive vs. Manylogs

- 320 writes/s

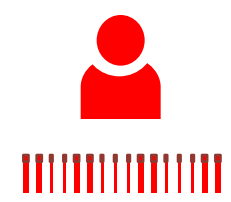
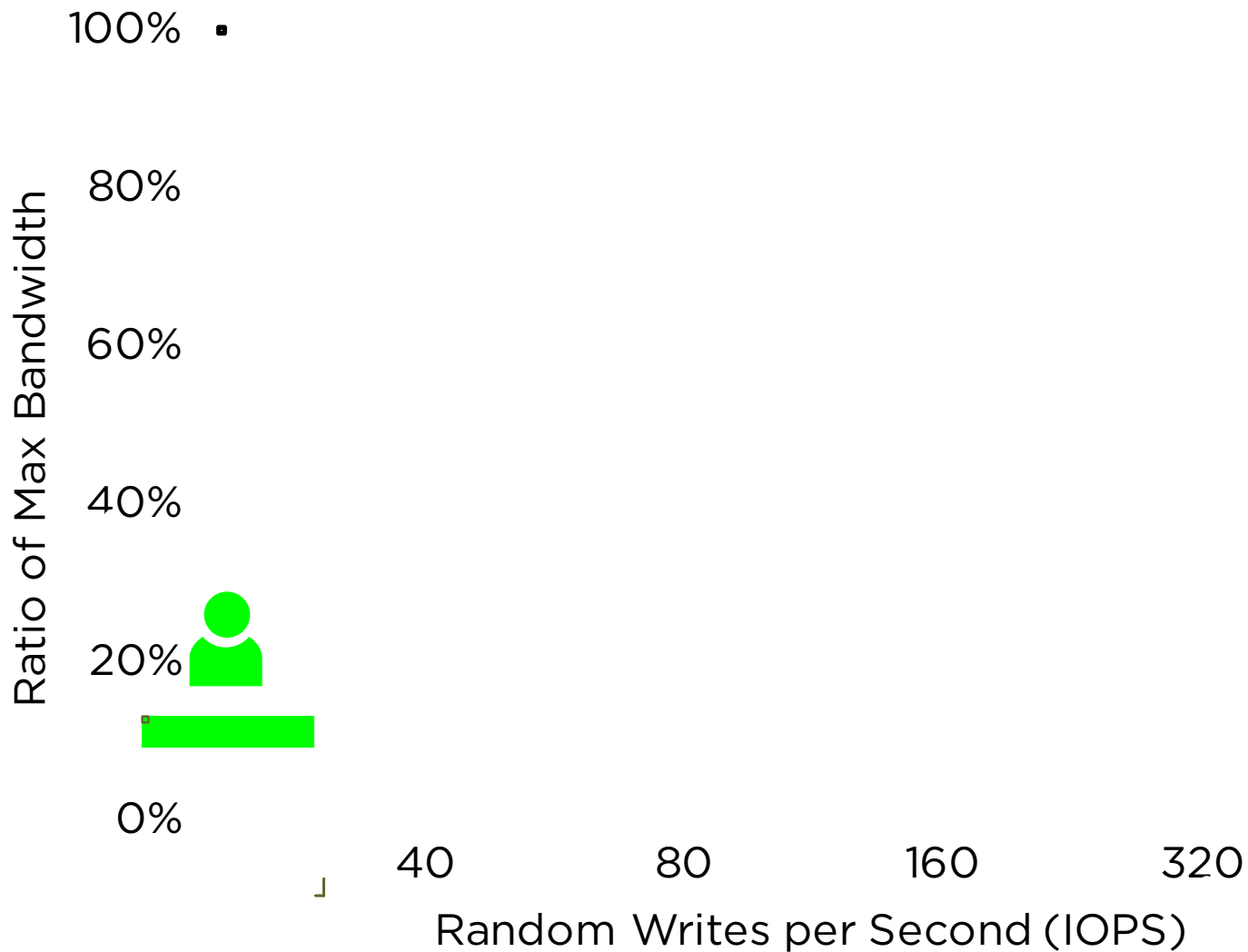


Disk

Adaptive Journaling

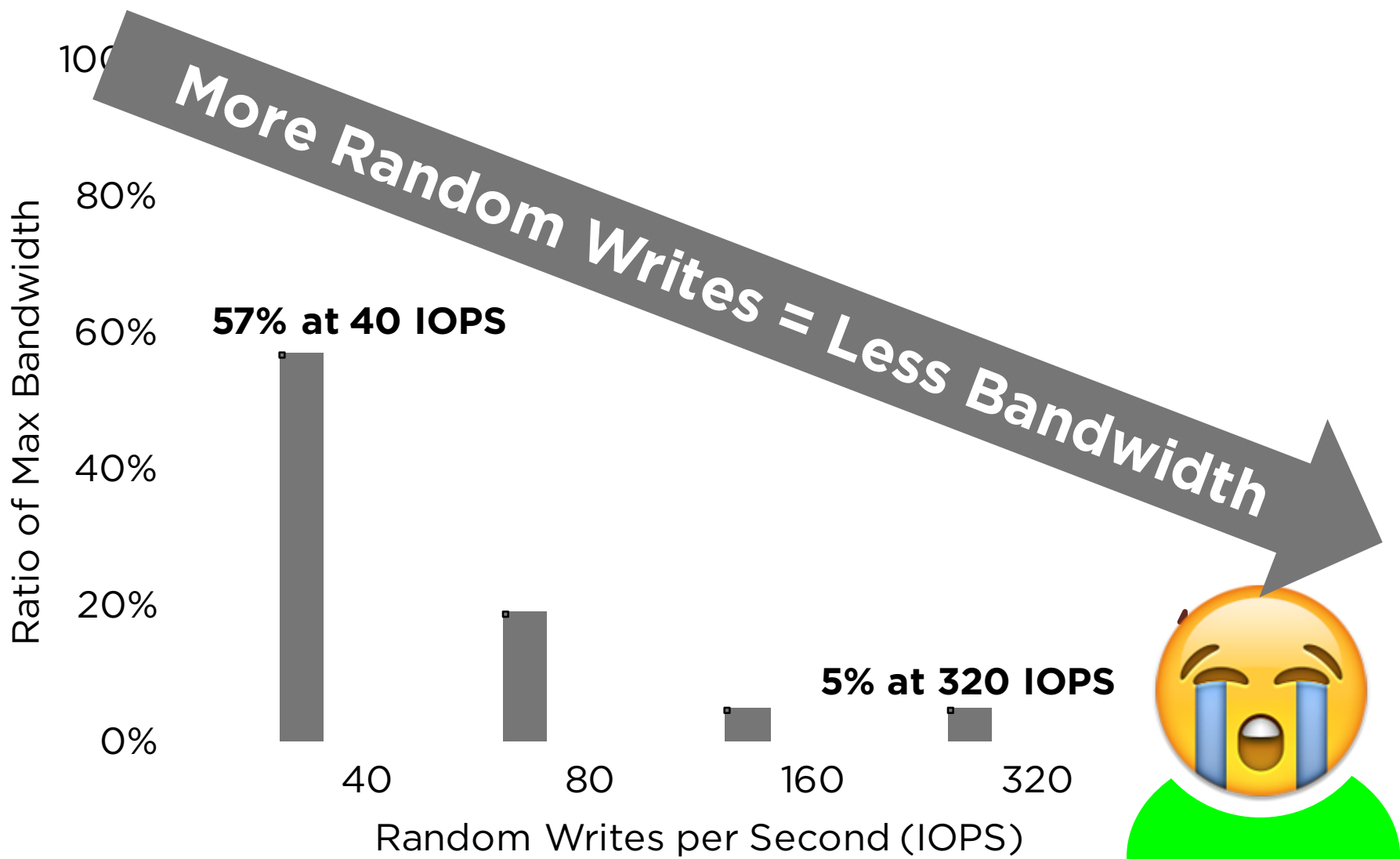
- ❑ Middle ground between ordered journaling and data journaling
- ❑ **Single-log** design
- ❑ Prabhakaran et al., ATC '05

Results

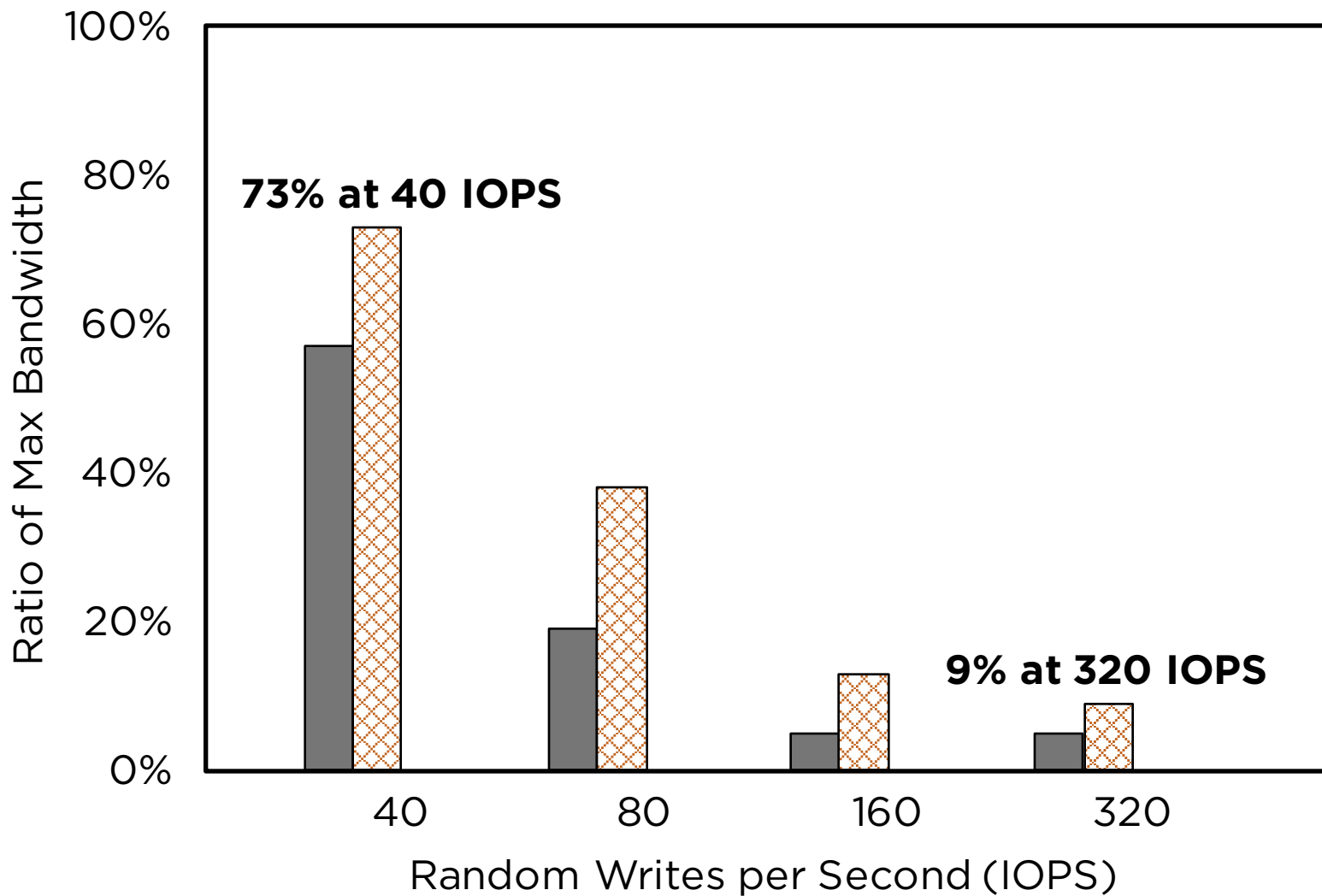


Results

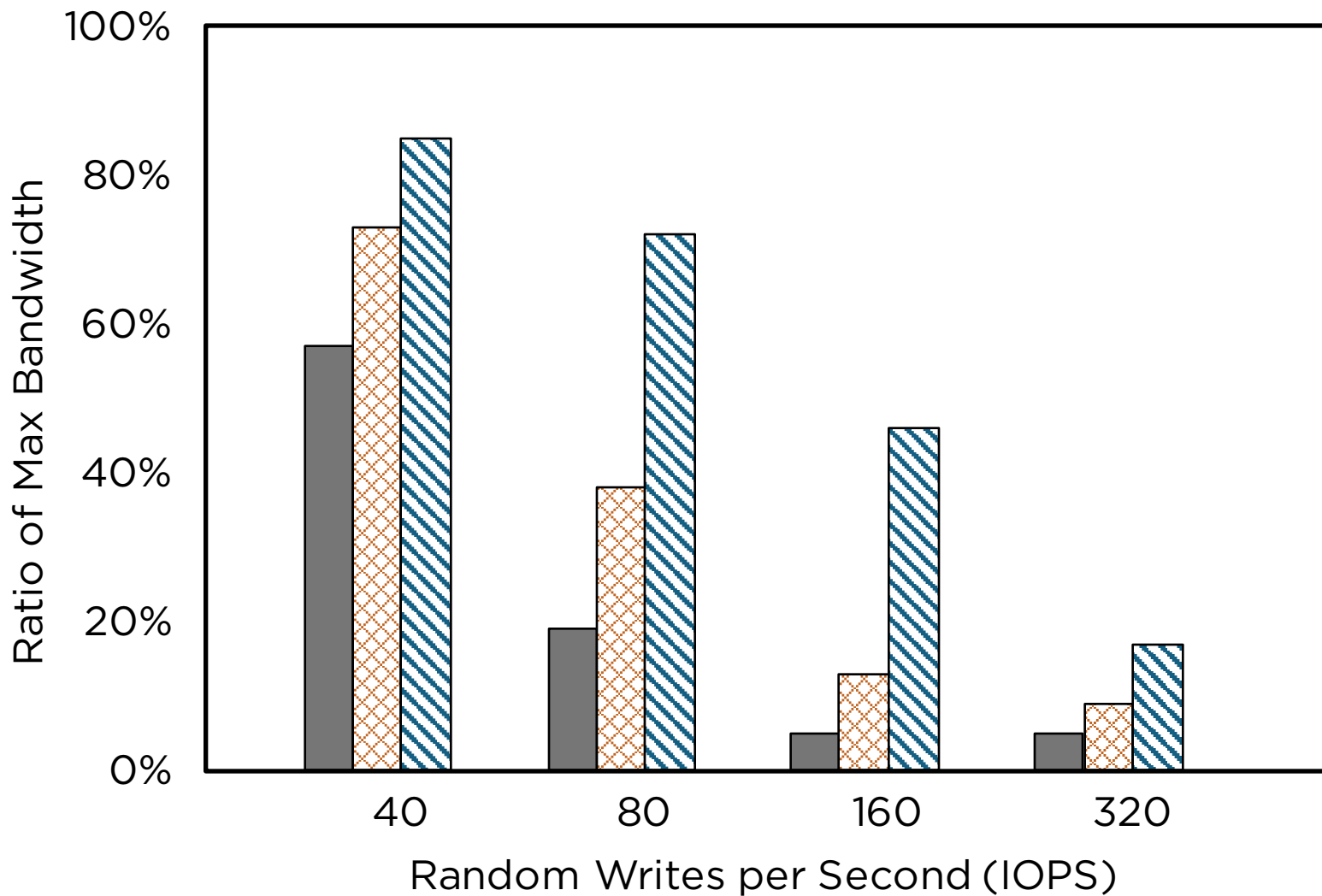
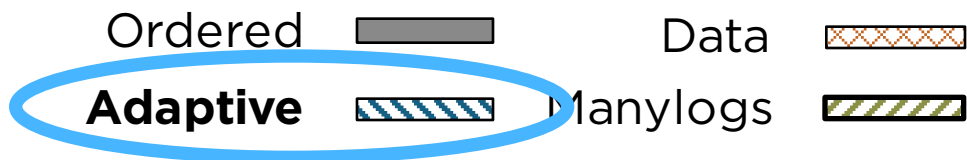
- Ordered 
- Adaptive 
- Data 
- Manylogs 



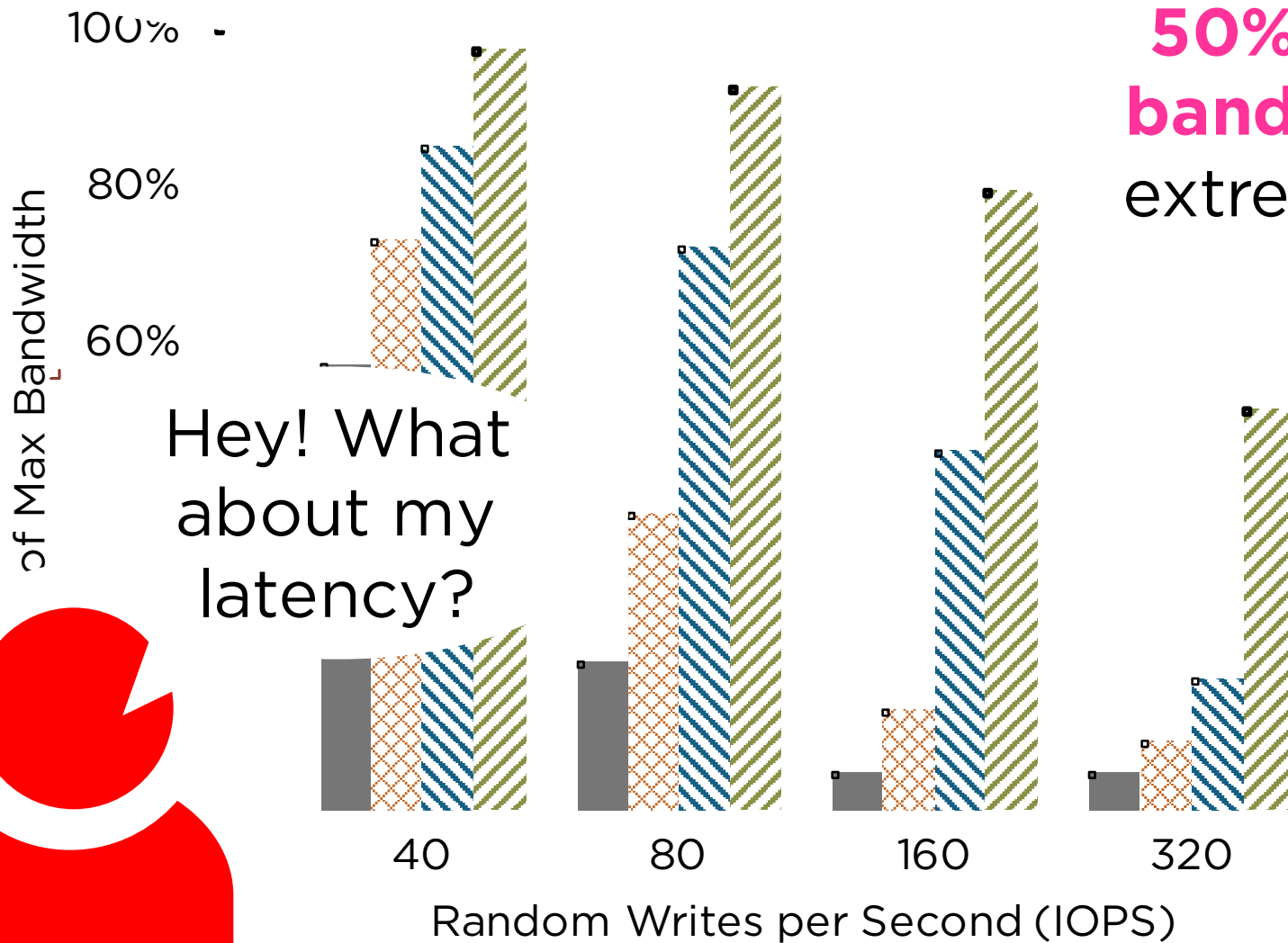
Results



Results

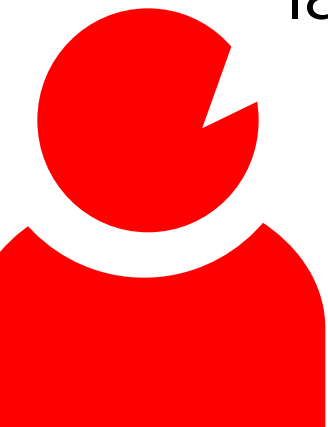


Manylogs gives the **most bandwidth**

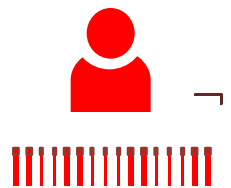
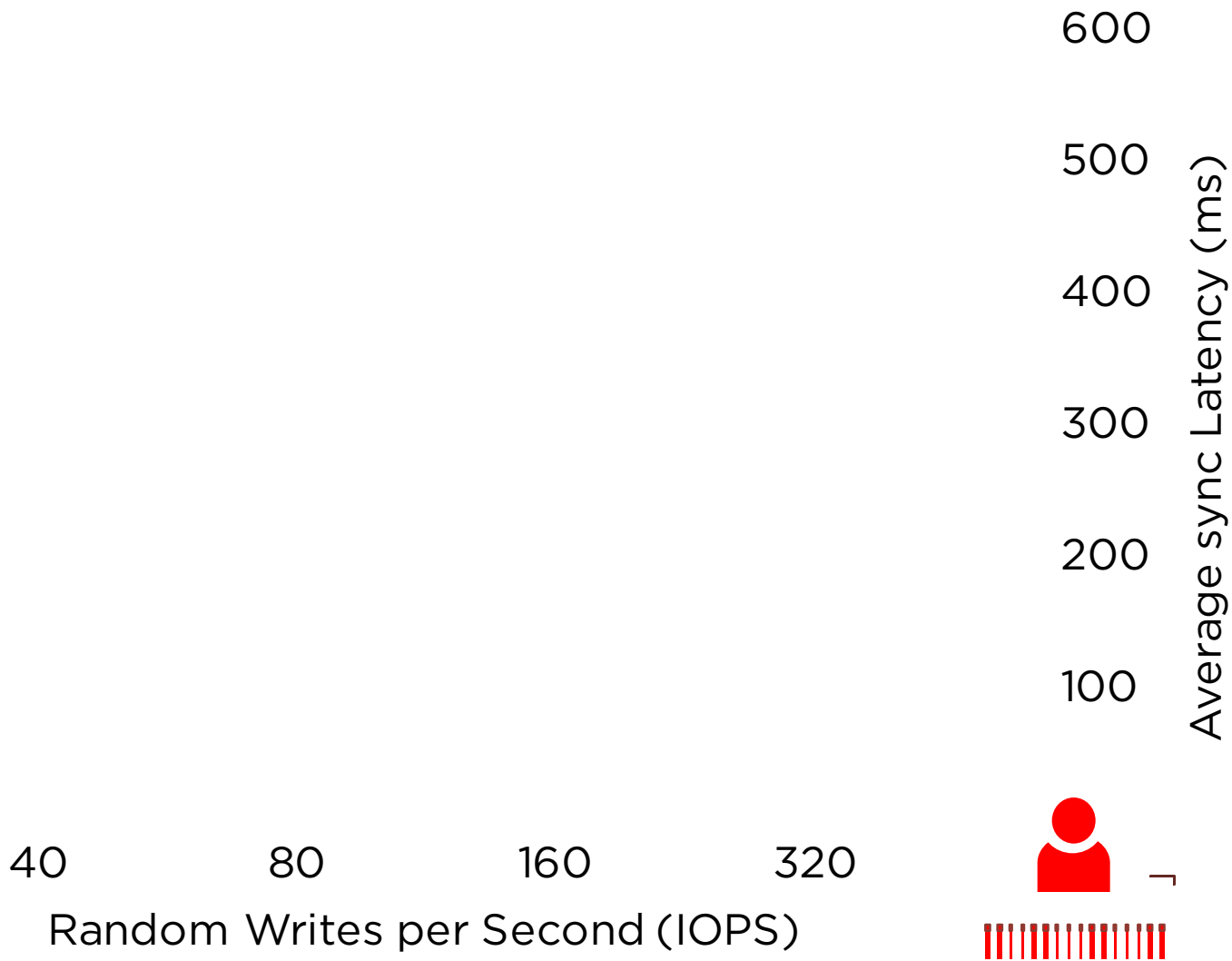


50% of max bandwidth at extreme IOPS

Hey! What about my latency?



Results



Results

Ordered

Adaptive



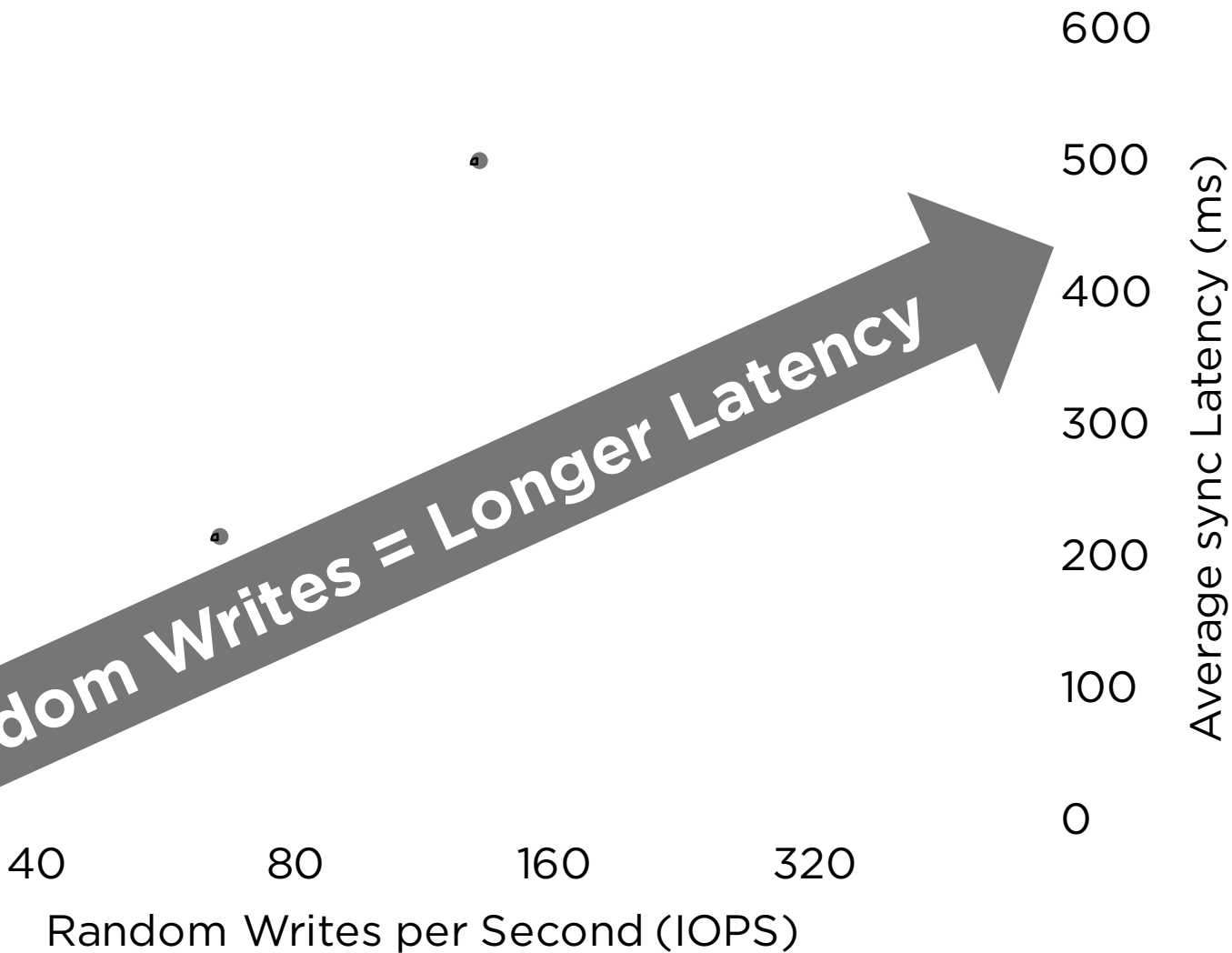
Data



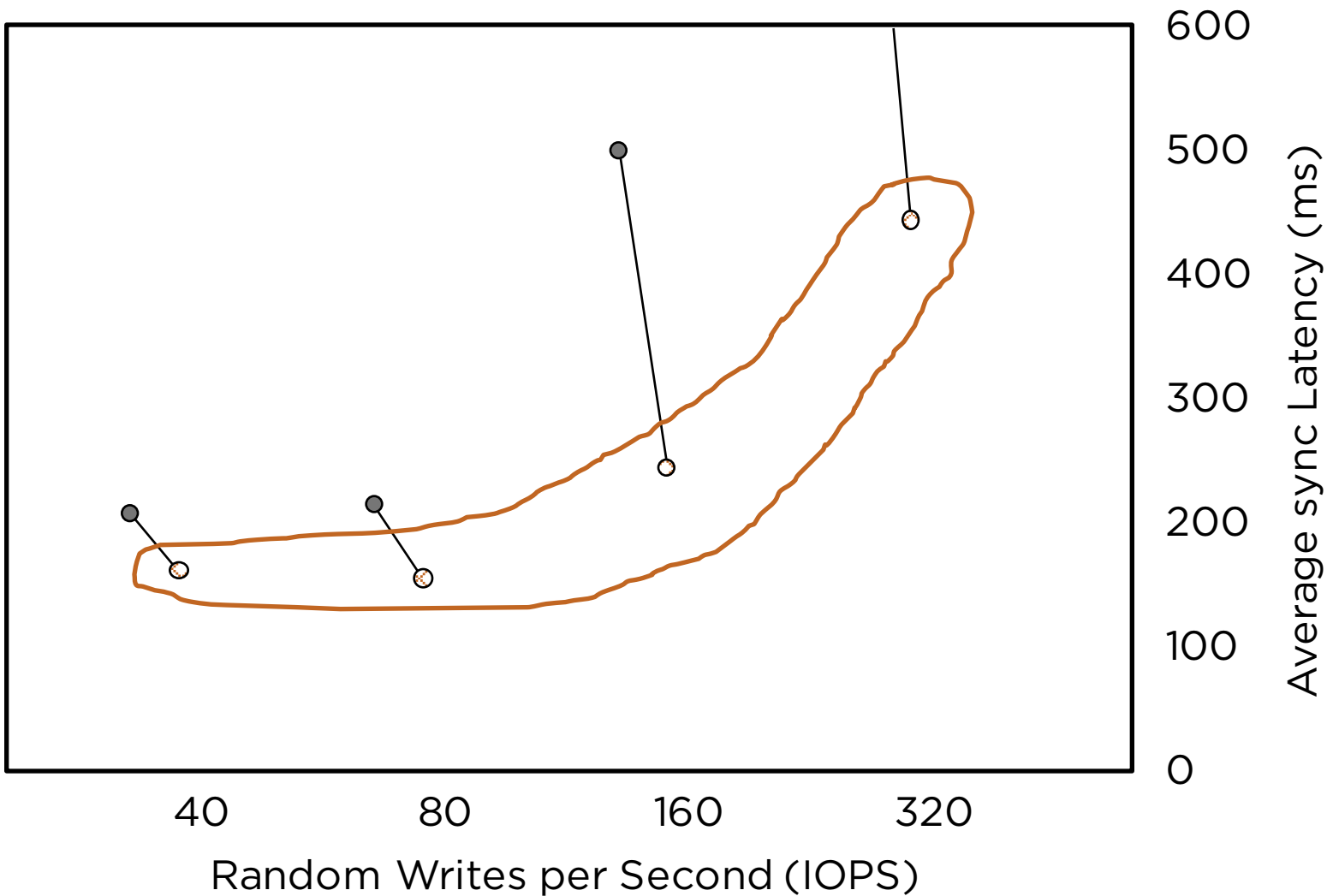
Manylogs



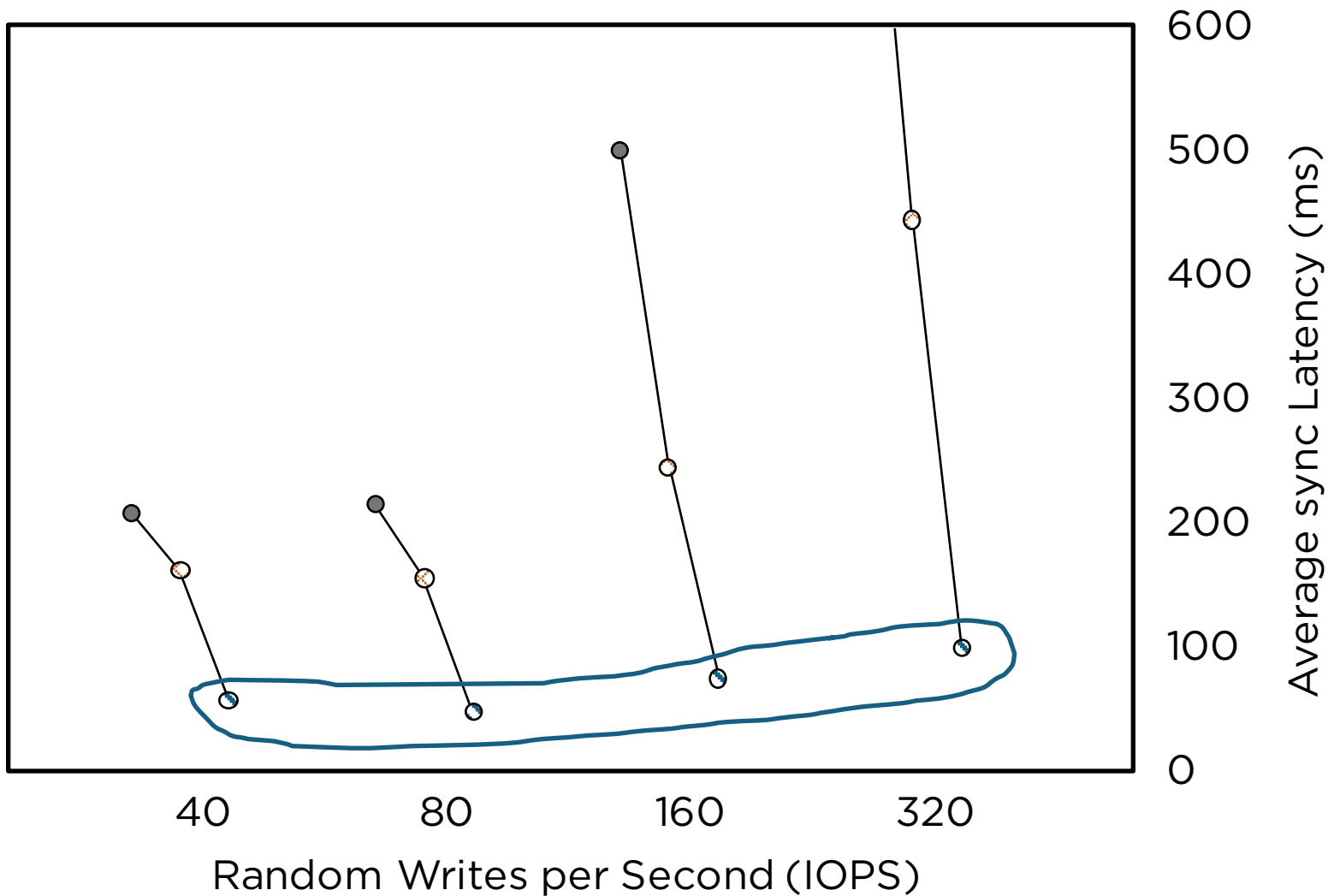
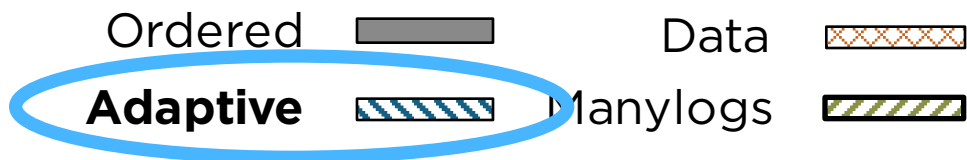
More Random Writes = Longer Latency



Results



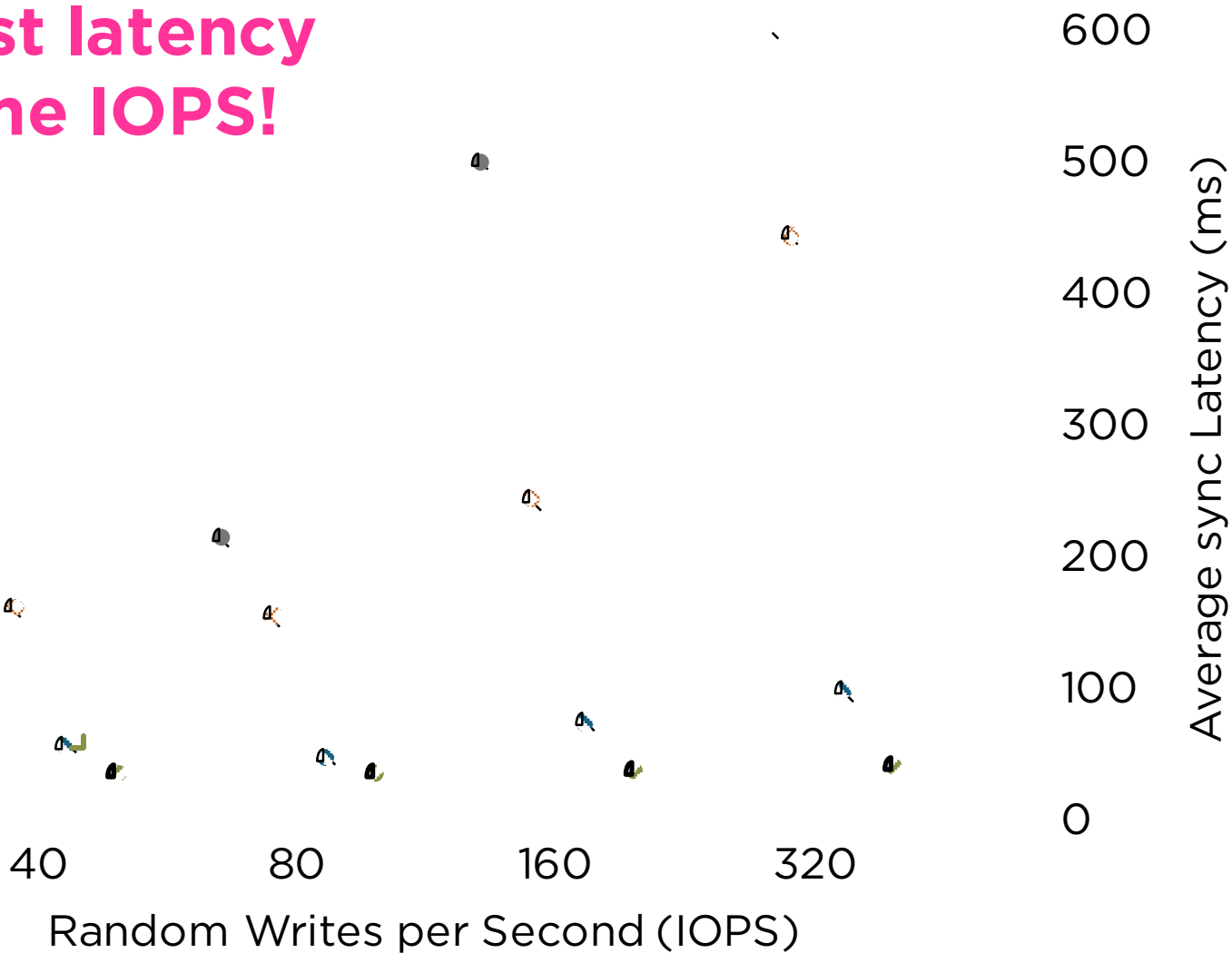
Results



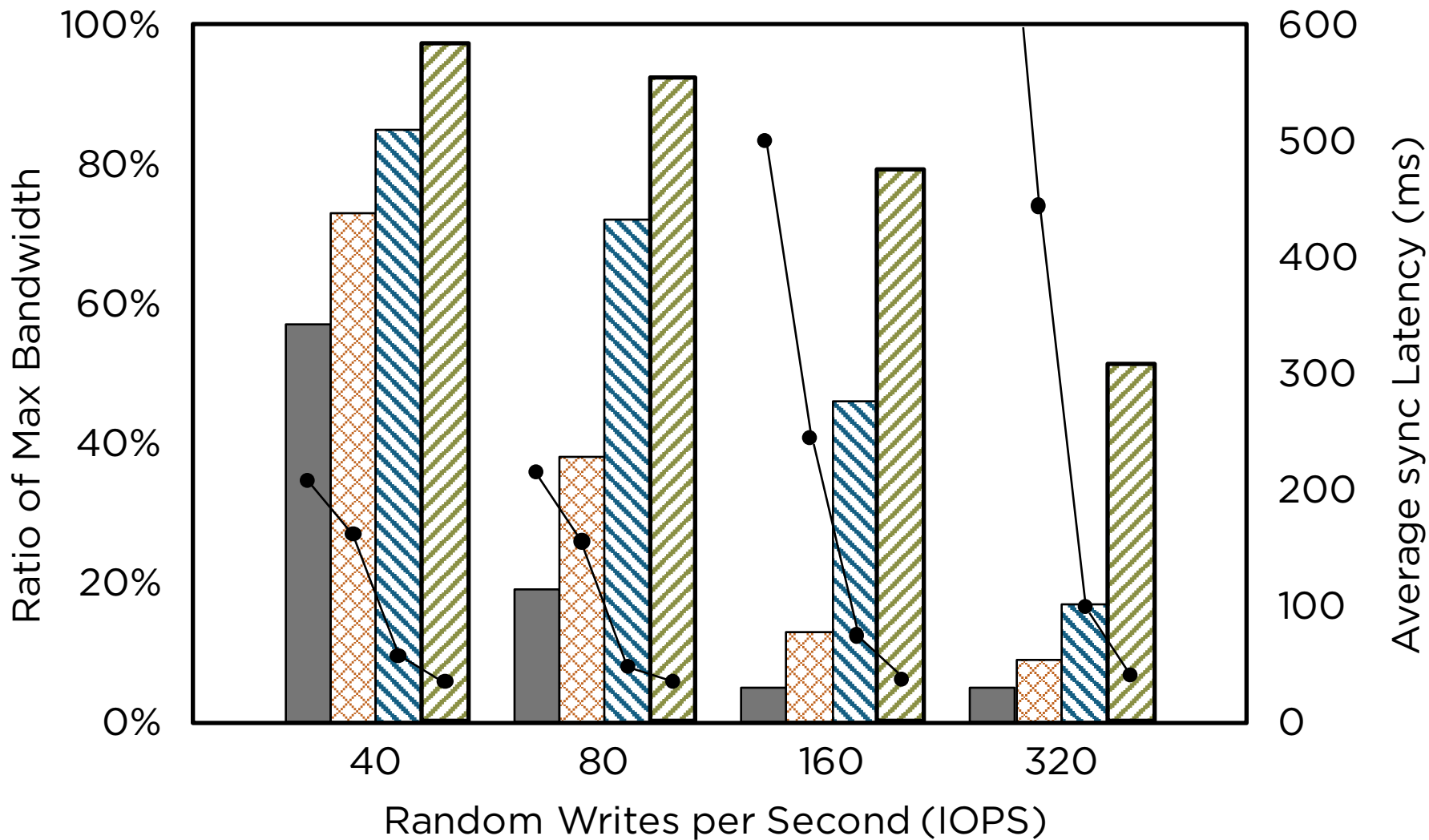
Results



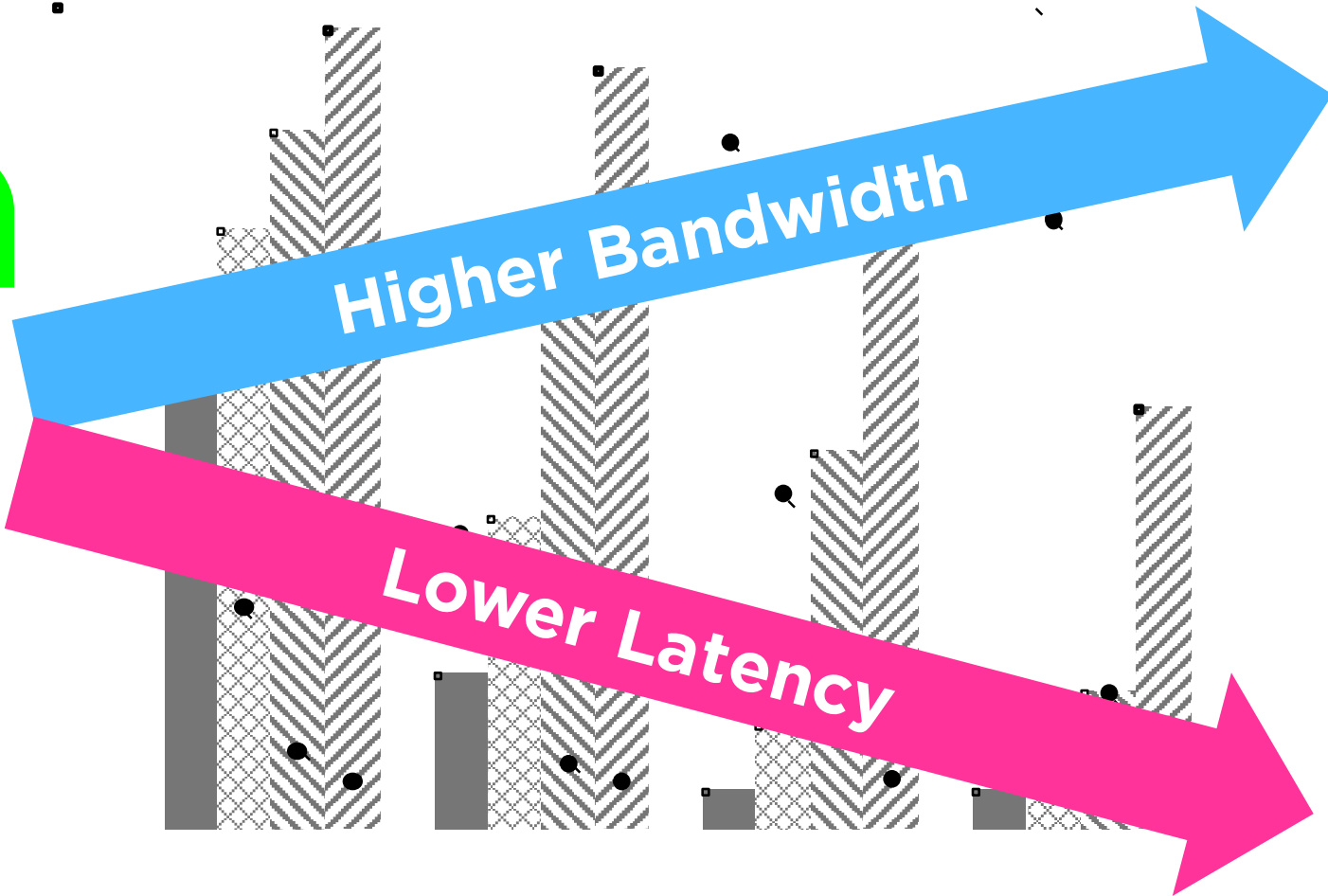
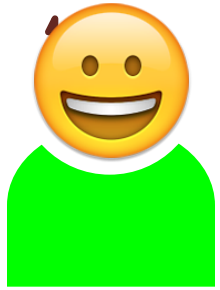
WOW! **Fast latency**
at **extreme IOPS!**

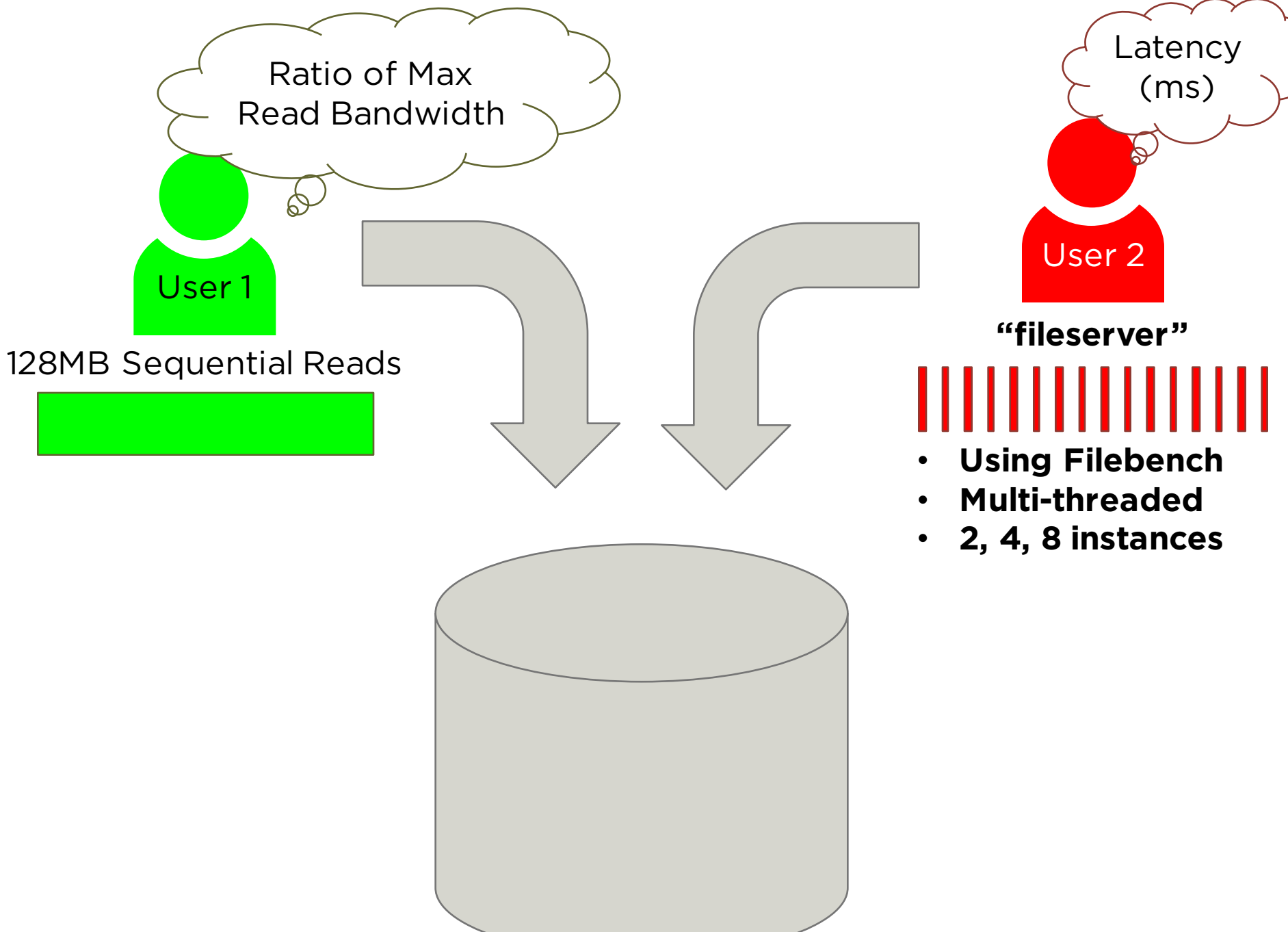


Results

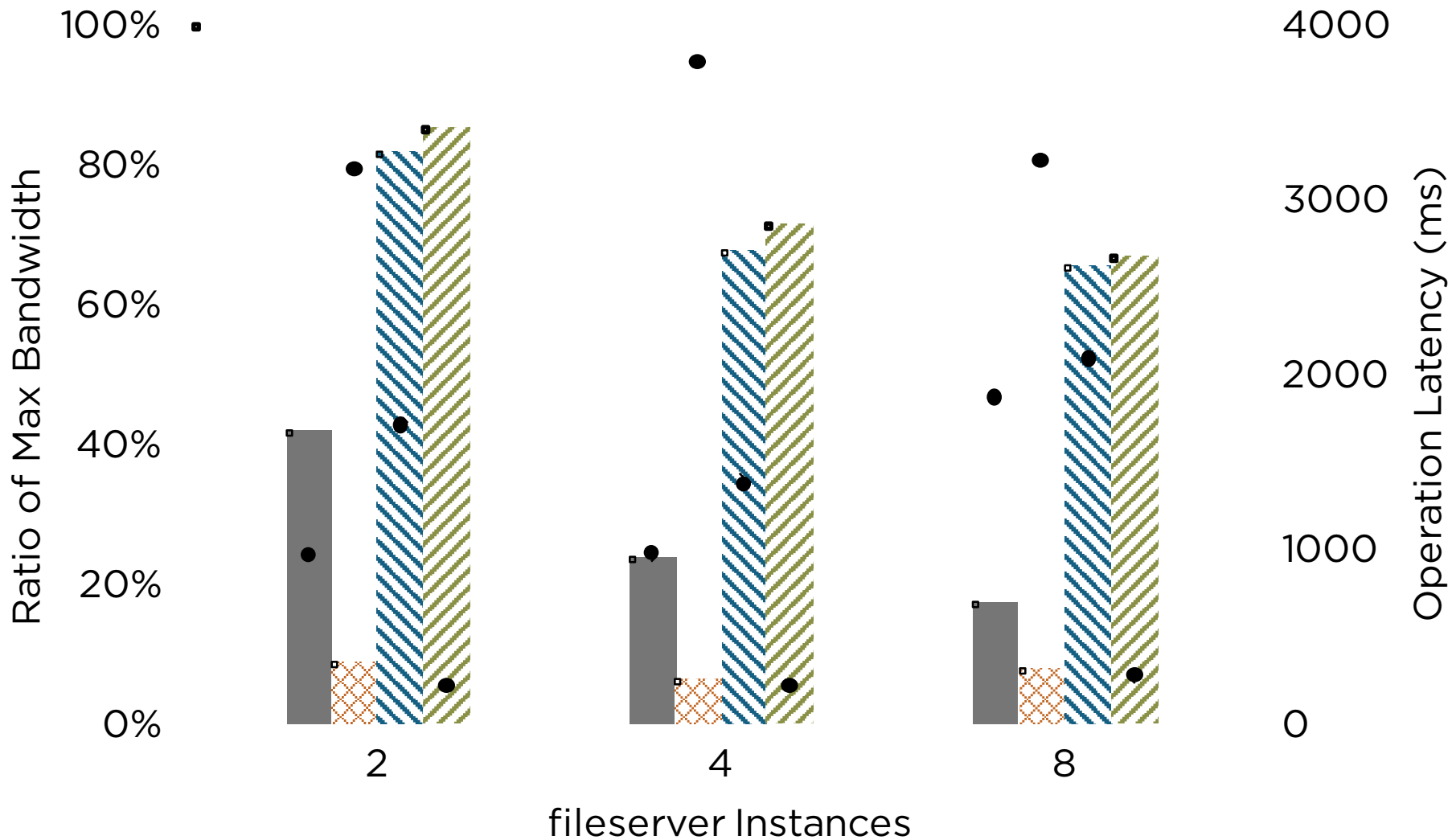


Results





Manylogs provides the best outcomes!



Checkpointing

Data Journaling

- ❑ Periodically
 - Usually every 5 secs
- ❑ Journal can get filled fast because **all writes** are in the journal!

Manylogs

- ❑ “Lazy” or “Off-hours”
- ❑ Rarely full because just **small writes** are redirected
- ❑ **Log Swapping**

Log Swapping

Log is almost full



Hot Area!

Cold Area!



Journal

Big Read

Small Writes



Integrations

- ❑ File System (MLFS)
 - Durability-Only Mode (O_DUR)
- ❑ SMR Disk (MLSMR)
- ❑ RAID

Cassandra Write Path

Writes



Memory

Disk



Cassandra Write Path

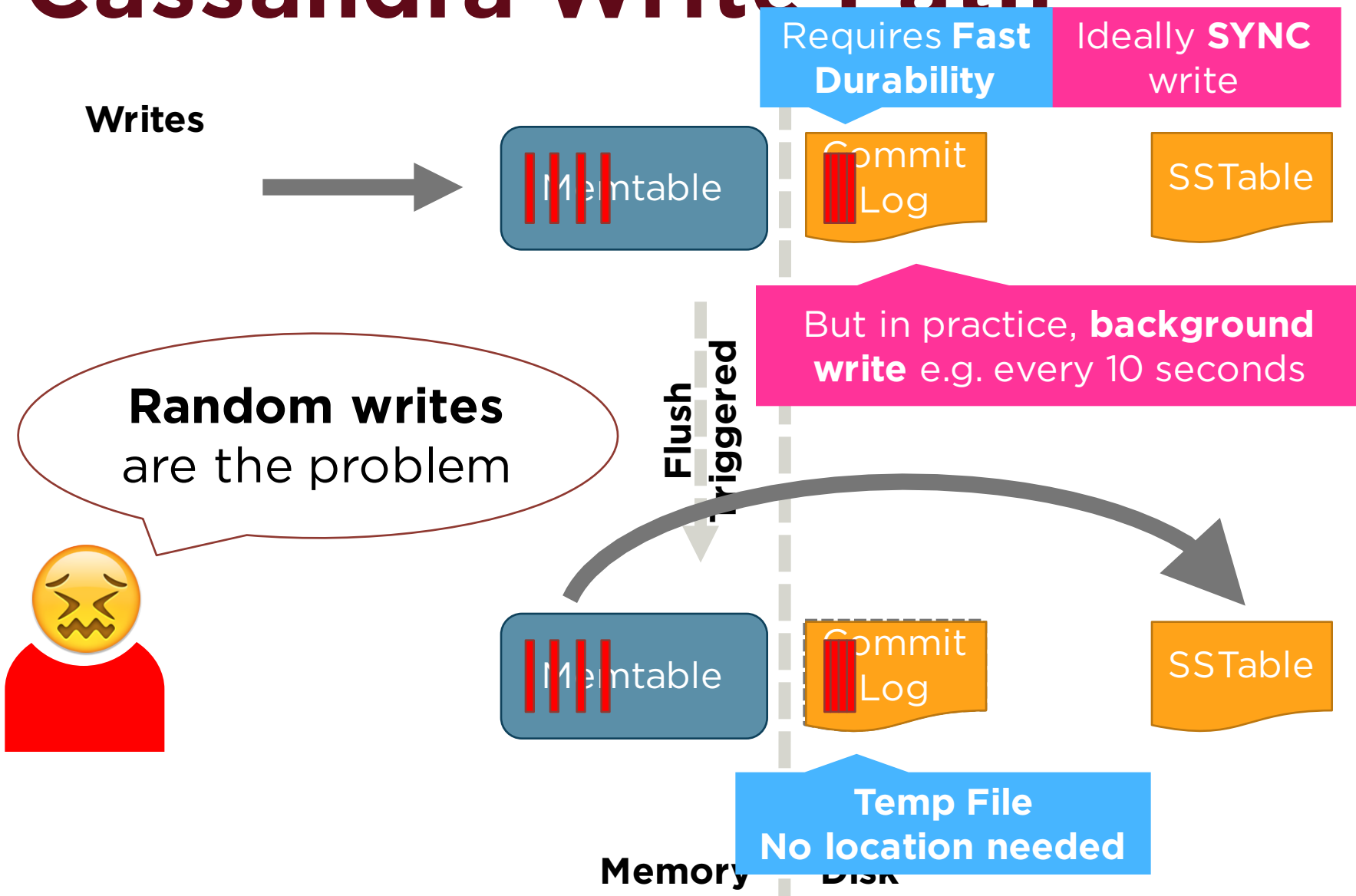
Writes



Memory

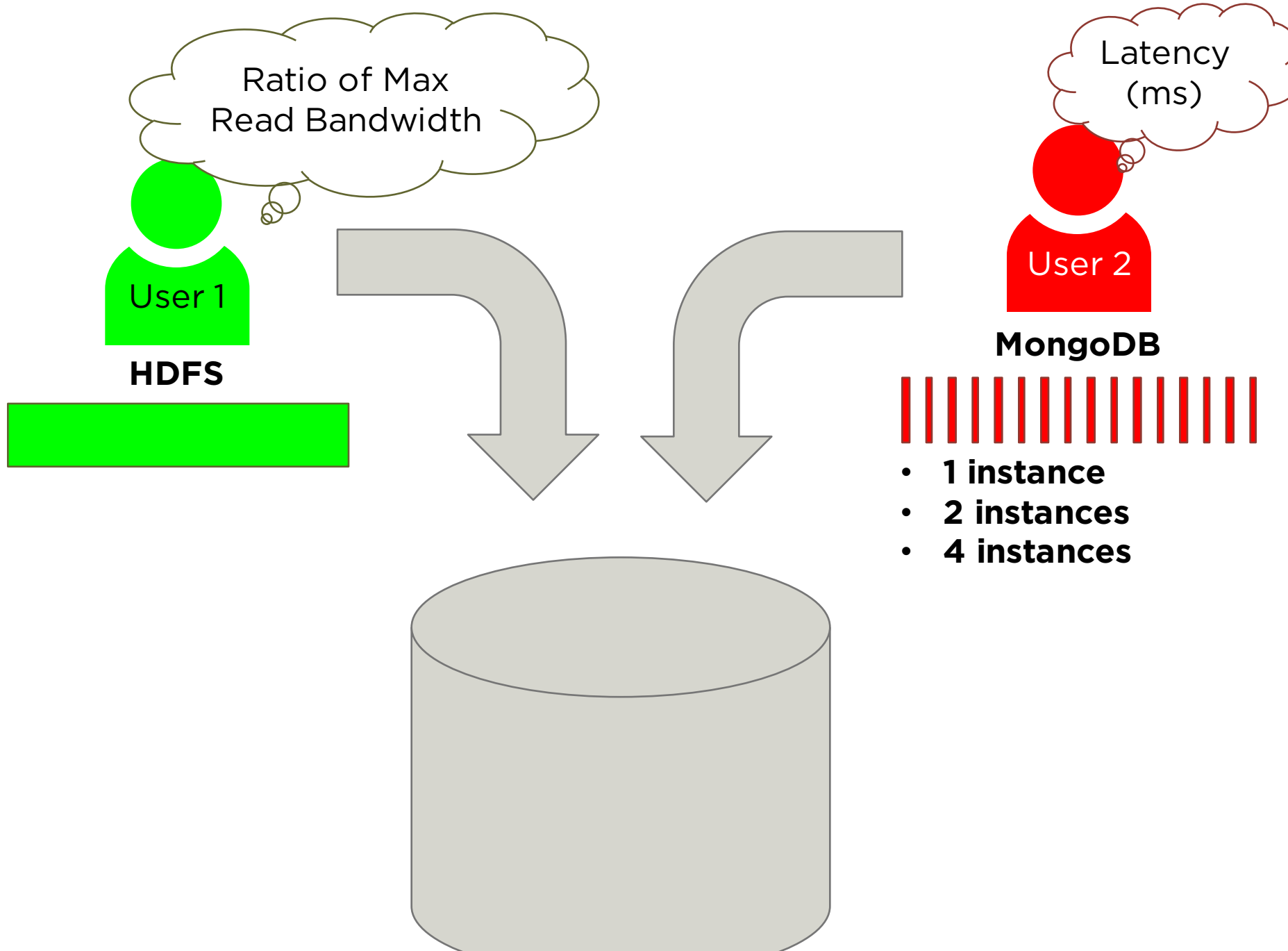
Disk

Cassandra Write Path

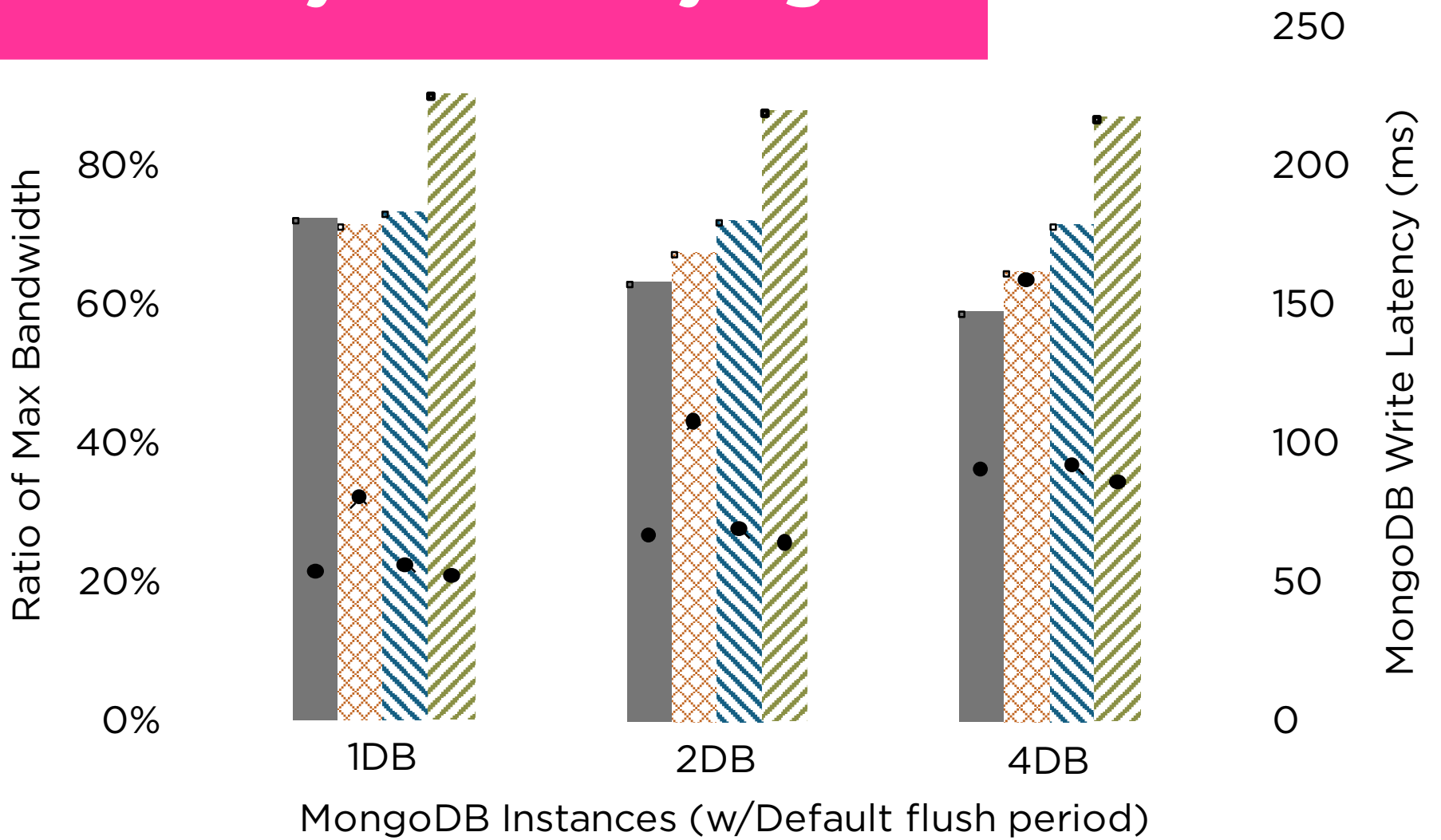


open(file, O_DUR);

- ❑ Need **fast durability** but not location constraints
- ❑ Content of files will be put in Manylogs regardless of the write size
- ❑ Never checkpoint their content
- ❑ **Random writes are not a problem anymore!**

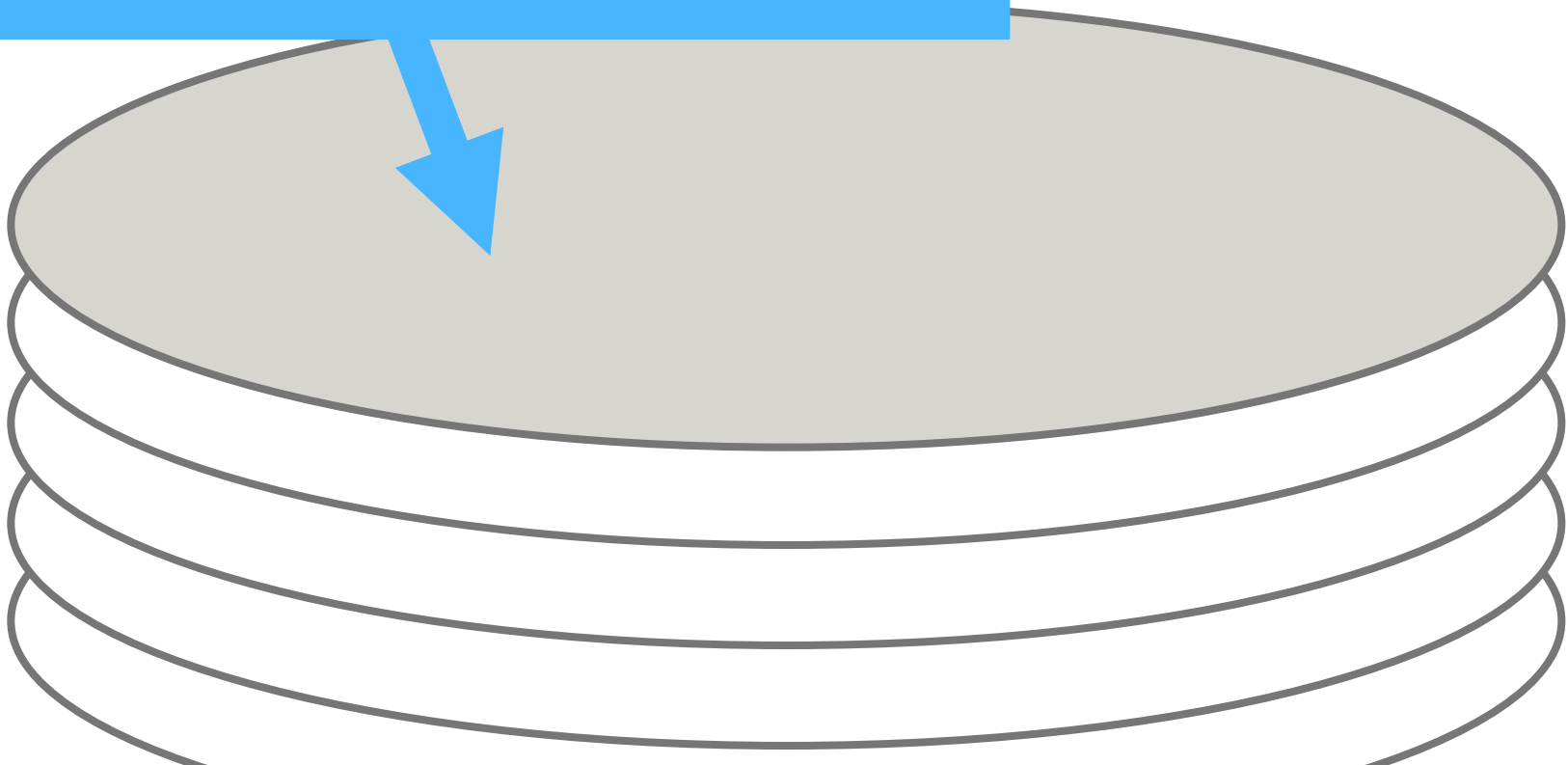


Most Bandwidth & Lowest Latency with Manylogs



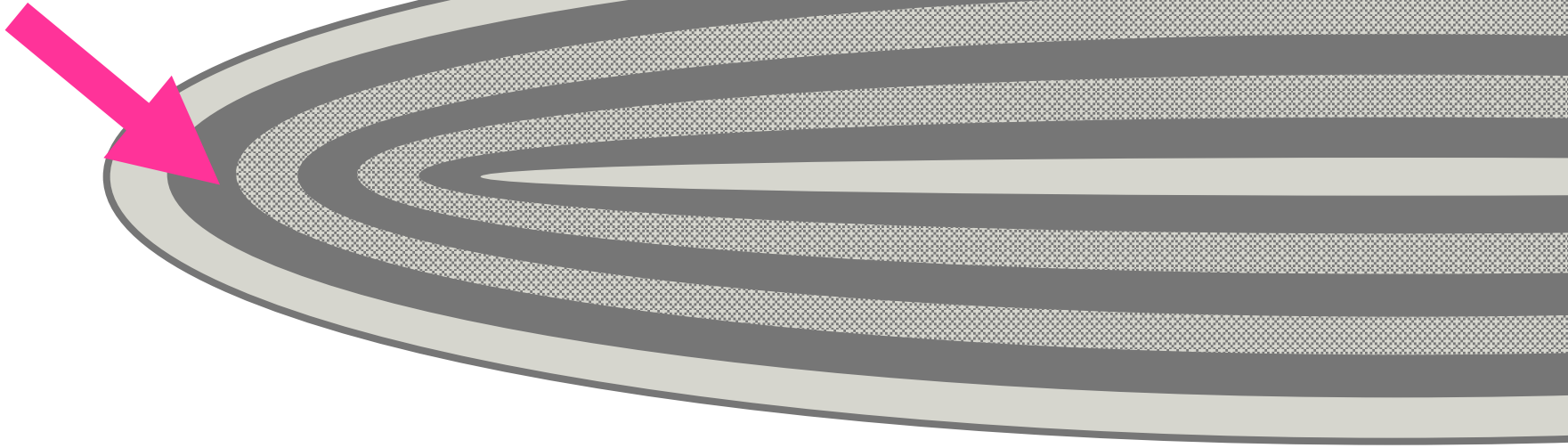
Manylogs & SMR

One non-shingled surface
= log space



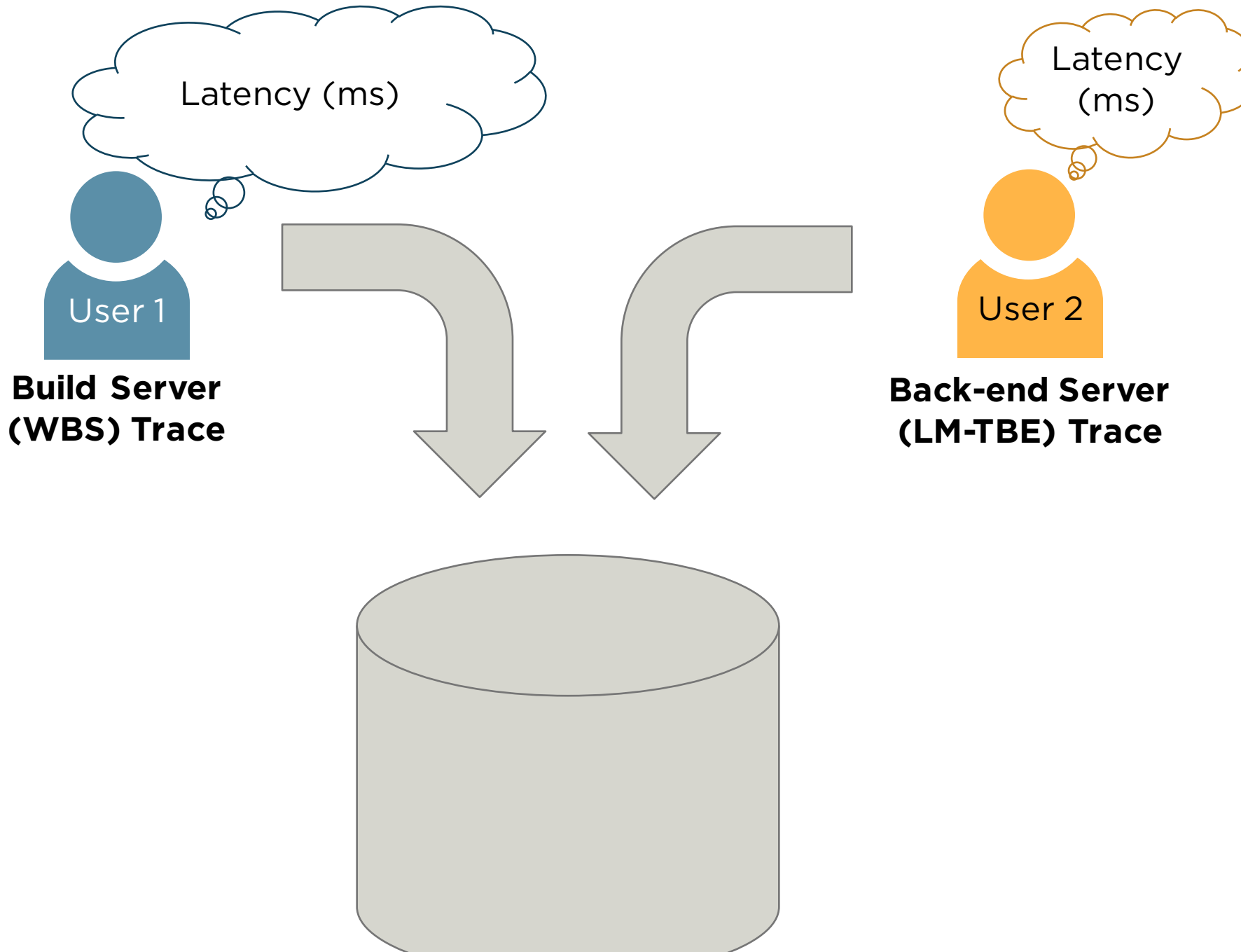
Manylogs & SMR

Non-Shingled Tracks





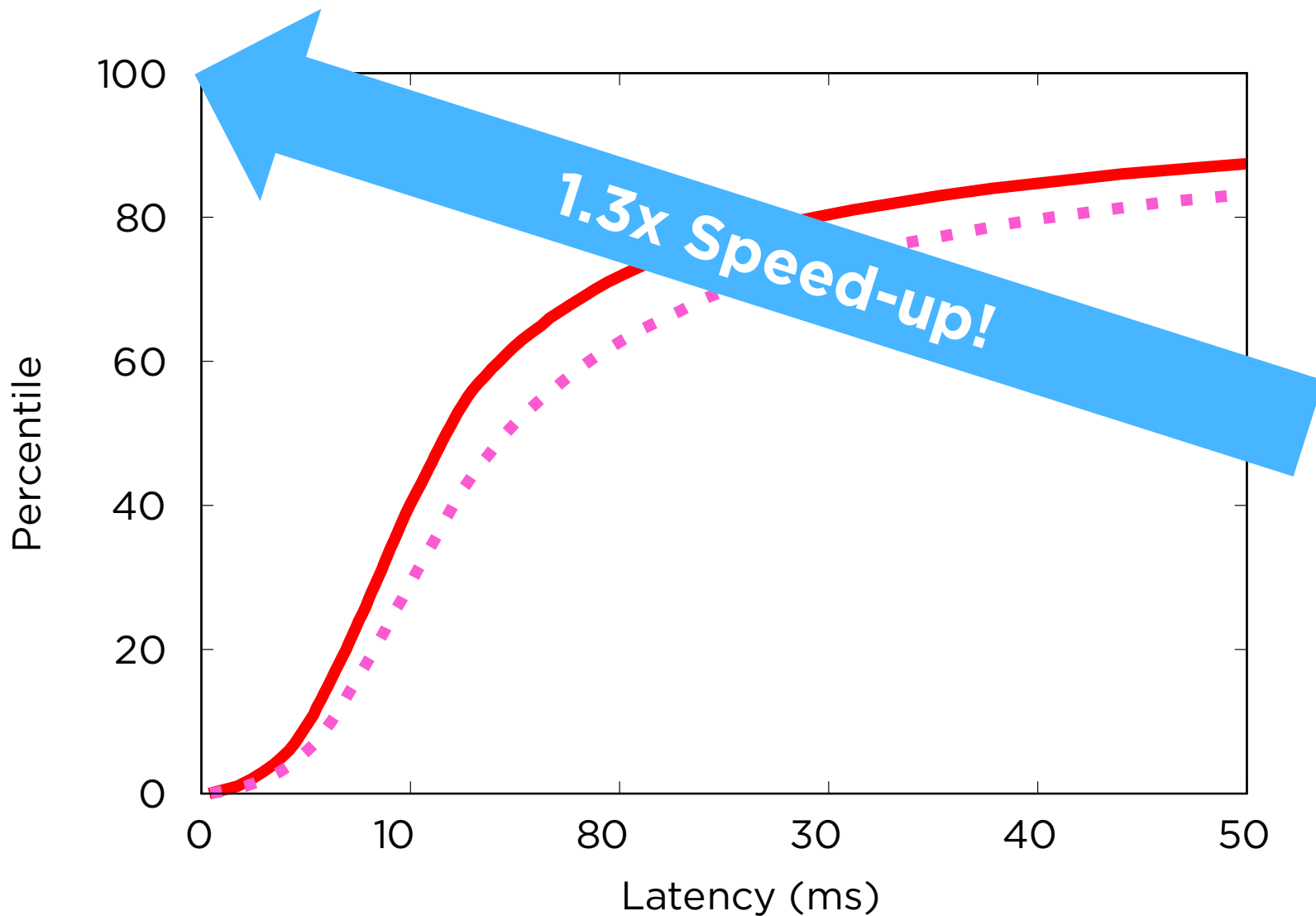
Shingled Band



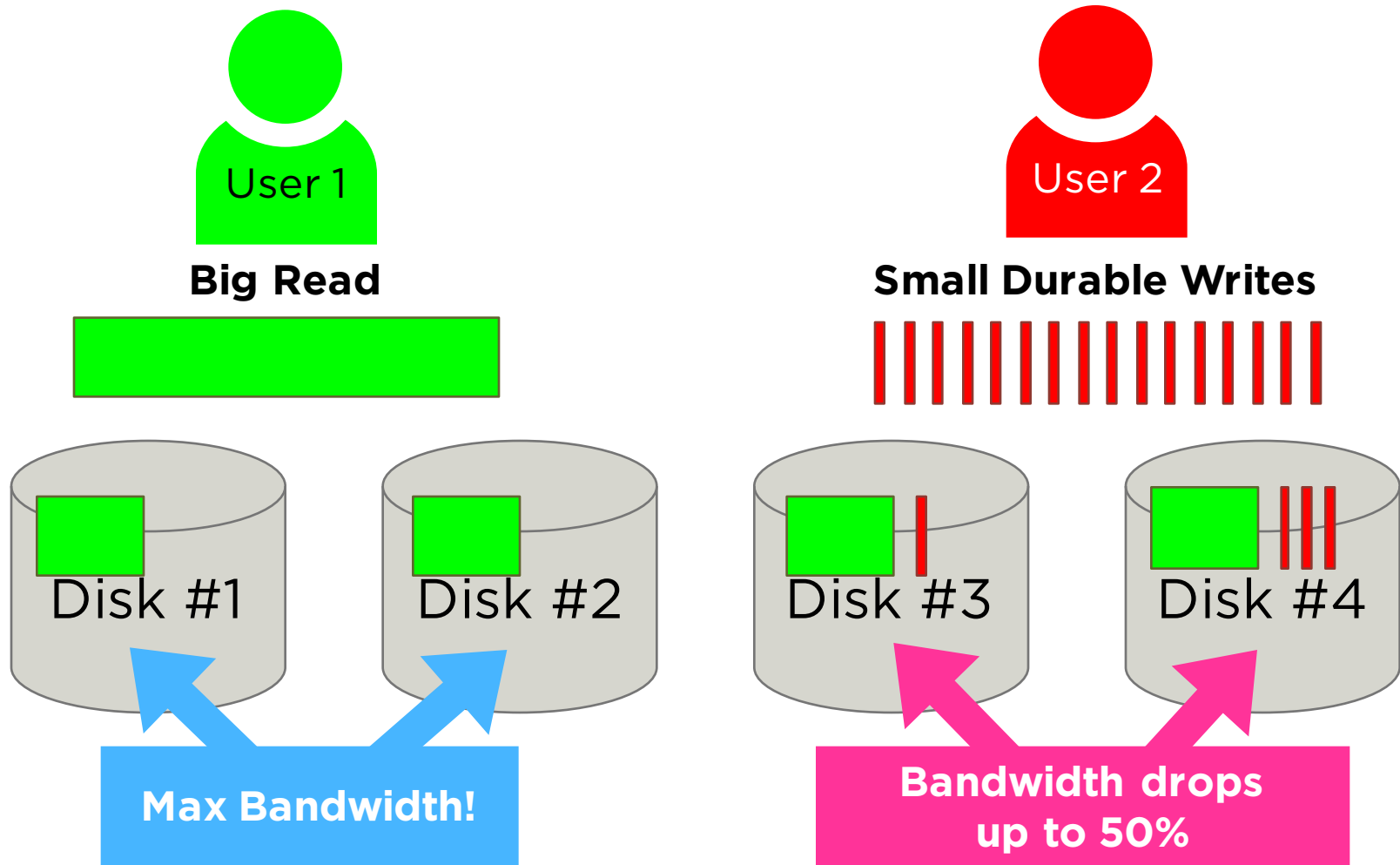


SMR

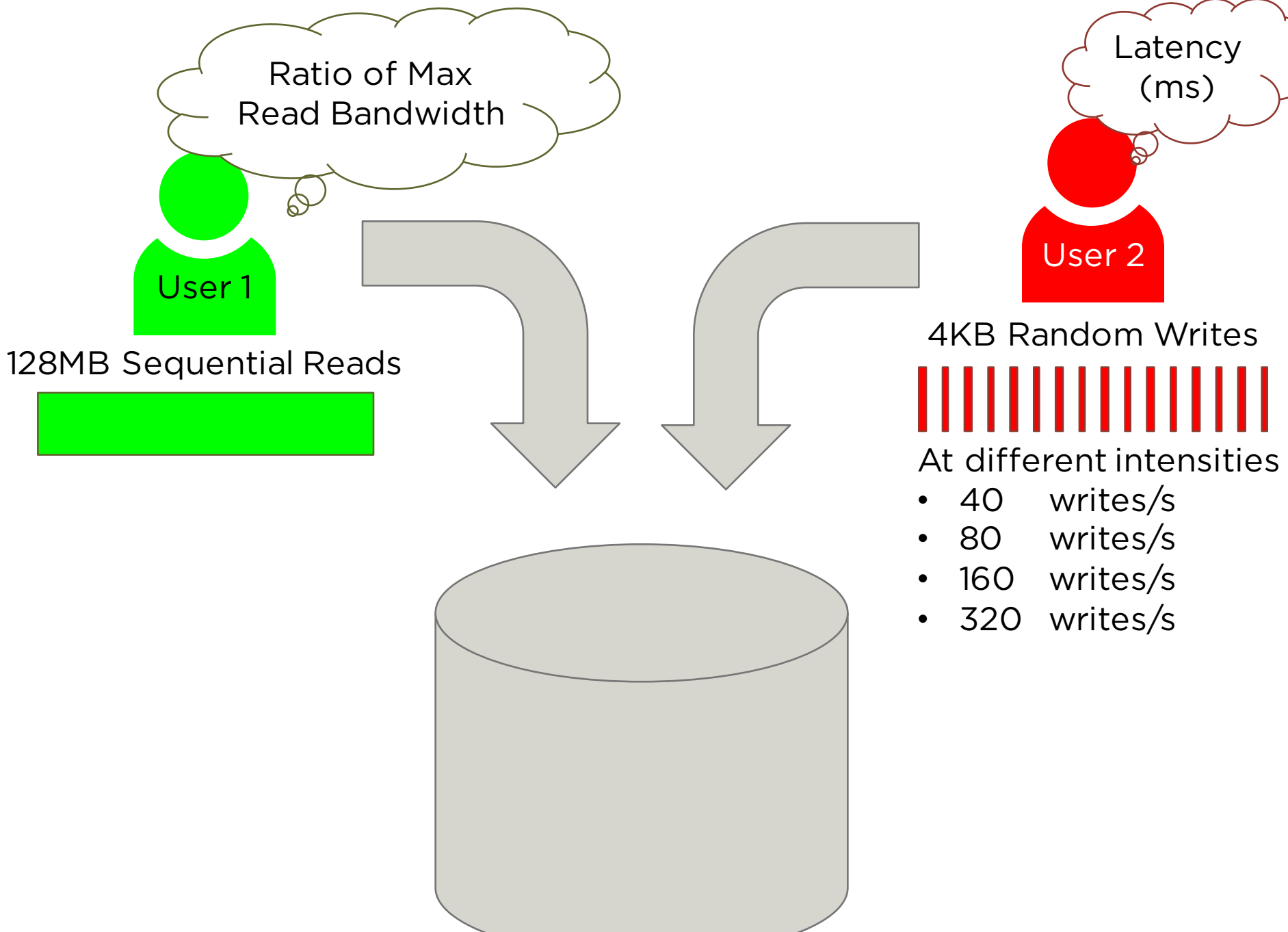
Manylogs SMR (MLSMR) 
Single-log SMR (SLSMR) 



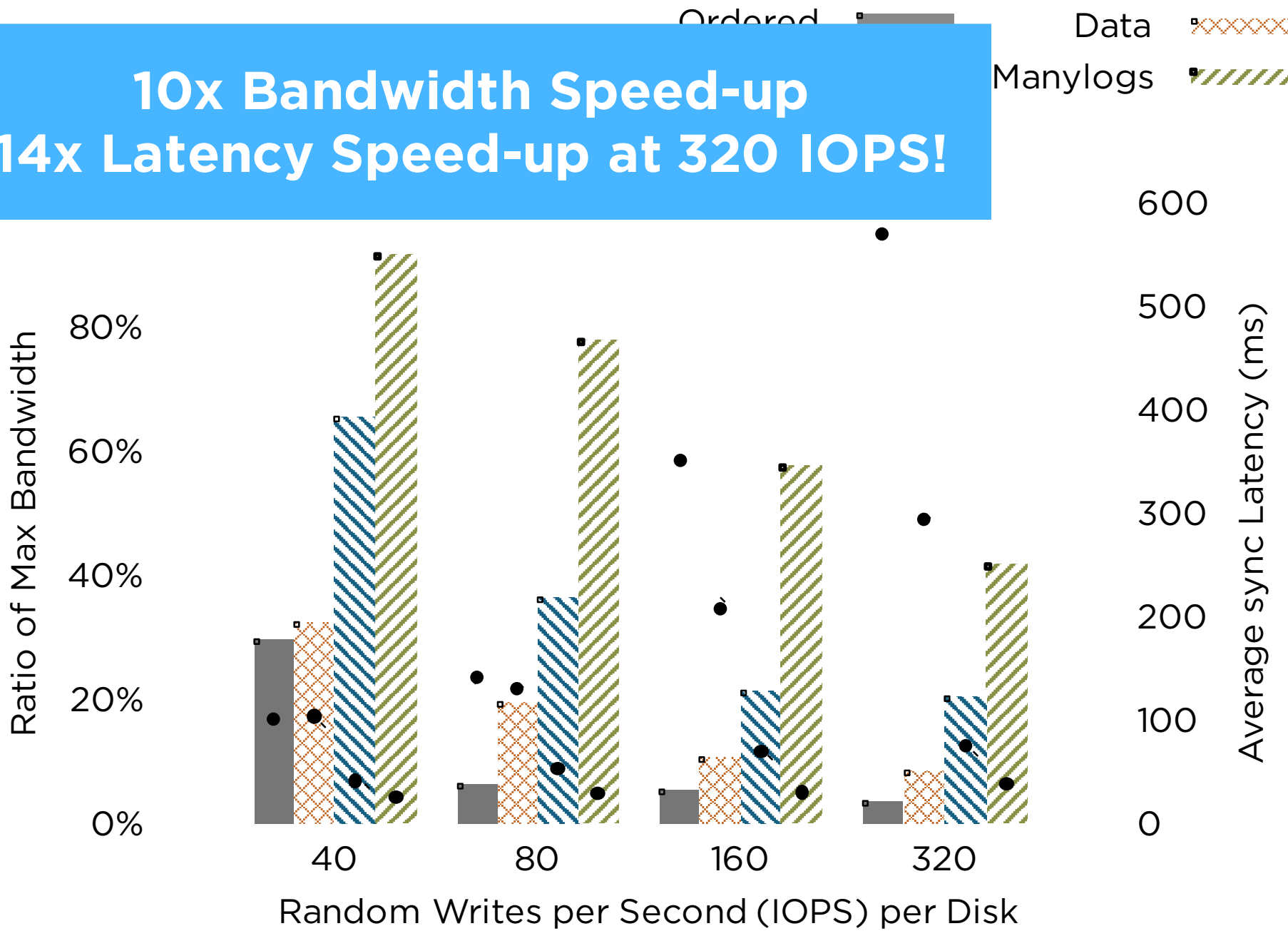
Manylogs & RAID



Mingzhe Hao, Gokul Soundararajan, Deepak Kenchamma-Hosekote, Andrew A. Chien, and Haryadi S. Gunawi. "The Tail at Store: A Revelation from Millions of Hours of Disk and SSD Deployments." FAST'16.



**10x Bandwidth Speed-up
14x Latency Speed-up at 320 IOPS!**

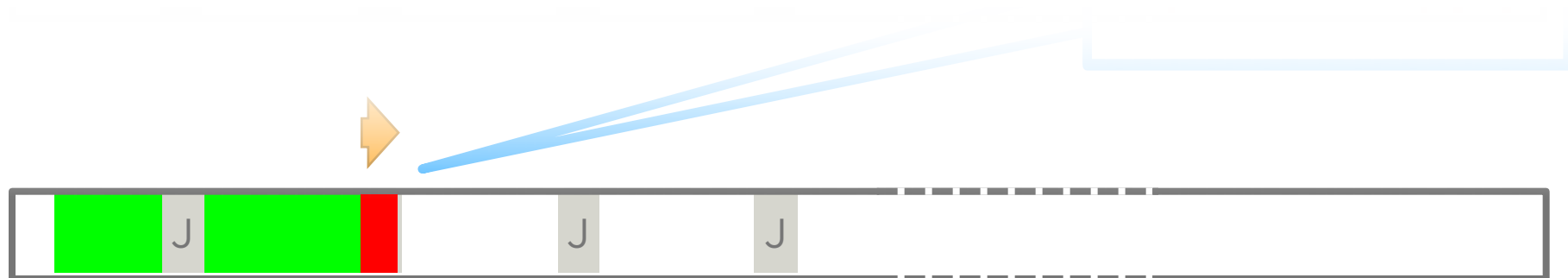


More in the paper

- ❑ Block-Level Manylogs
- ❑ Other workloads
 - Sequential Writes
 - “varmail”
 - More Traces
- ❑ Log Size
- ❑ Logged Write Size
- ❑ Mapping Table

Manylogs

- ❑ Reserved log spaces uniformly across the disk
- ❑ Redirect small writes to the nearest log
- ❑ Can help with **NoSQL**, **SMR**, **RAID**, and more!
- ❑ Provide up to **5x speed-up** on average



Manylogs

	Bandwidth Speed-up	Latency Speed-up
vs. Ordered	3.7x	5.7x
vs. Adaptive	2.7x	2.0x
vs. Single-log SMR		1.3x

Thank you!

Questions?



THE UNIVERSITY OF
CHICAGO

<http://ucare.cs.uchicago.edu>