



A Parametric I/O Model for Modern Storage Devices

Tarikul Islam Papon papon@bu.edu Manos Athanassoulis mathan@bu.edu





Modeling Performance

"Algorithm/Data Structure **X** has O(f(N)) performance, where N is the number of data pages on disk"

... is probably one of the most commonly read phrases in SIGMOD papers.

















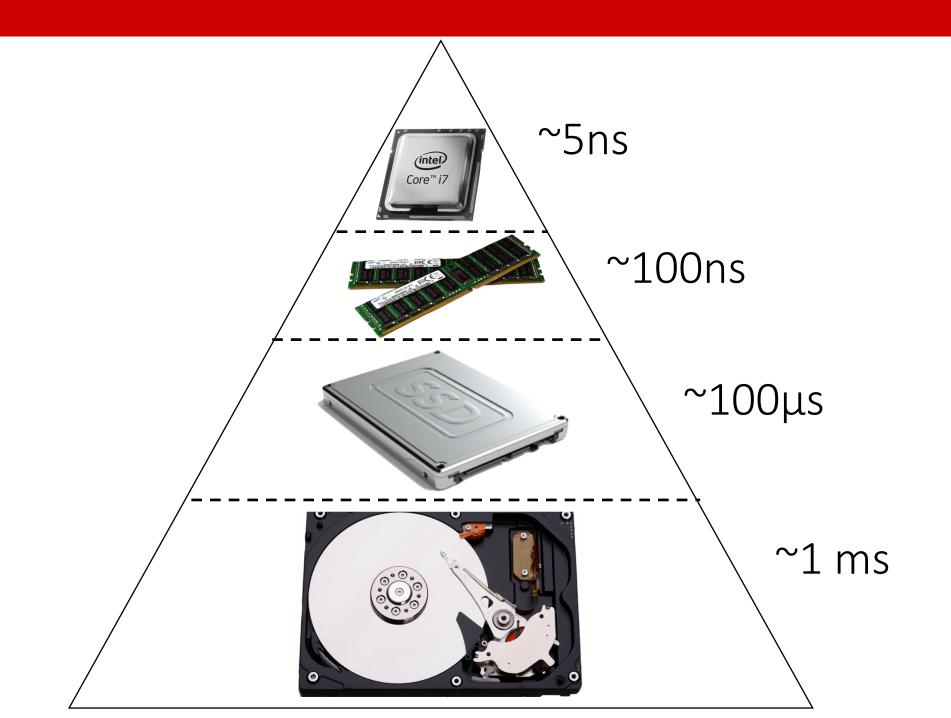






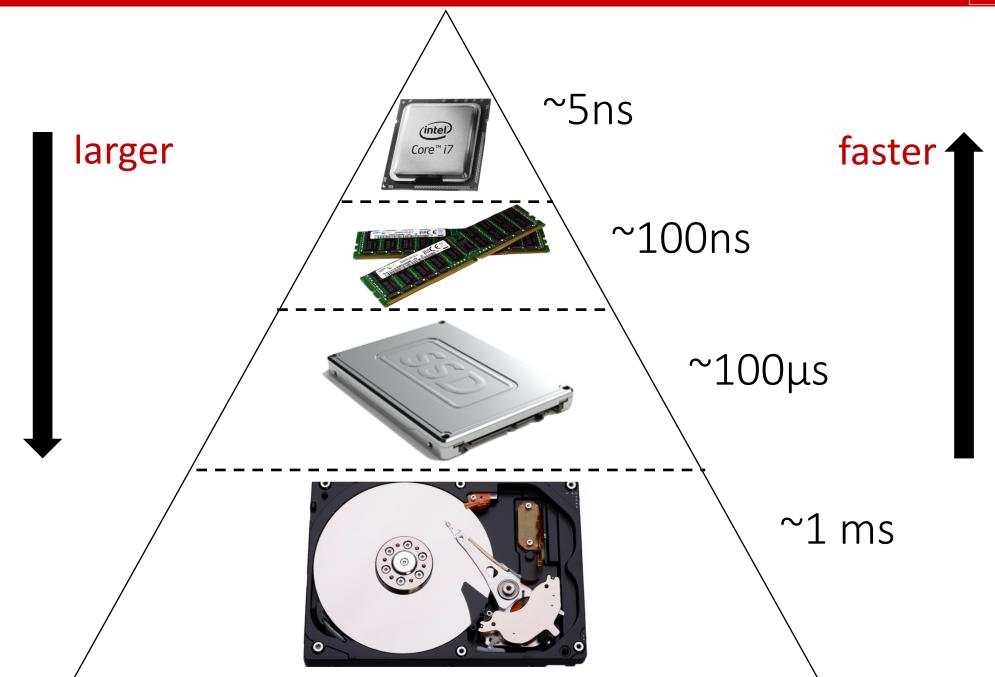






7











Small, fast main memory (size M)





망요 DiSC

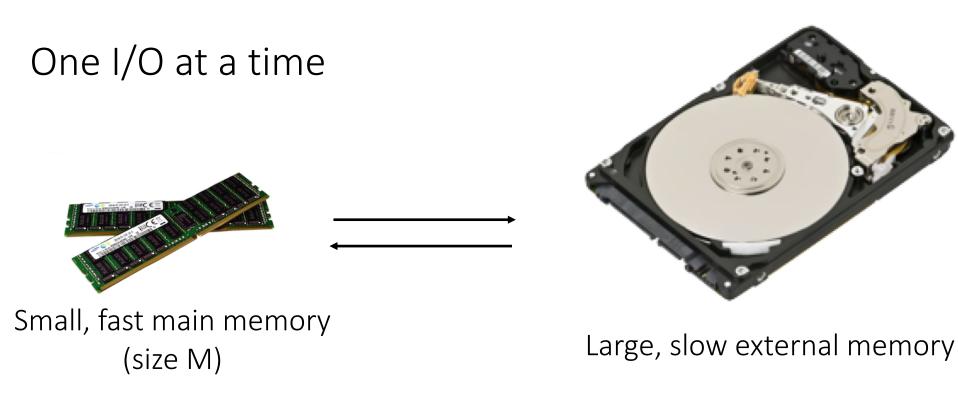
> Small, fast main memory (size M)



Large, slow external memory

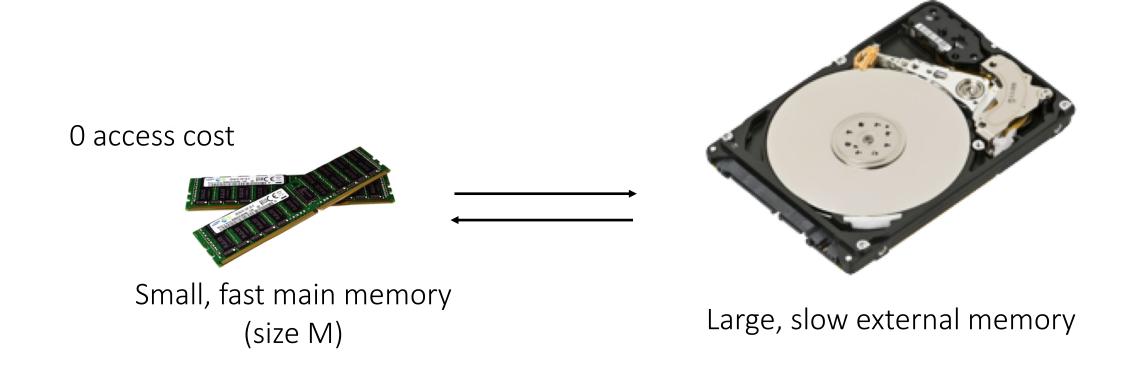


४१ हि DiSC



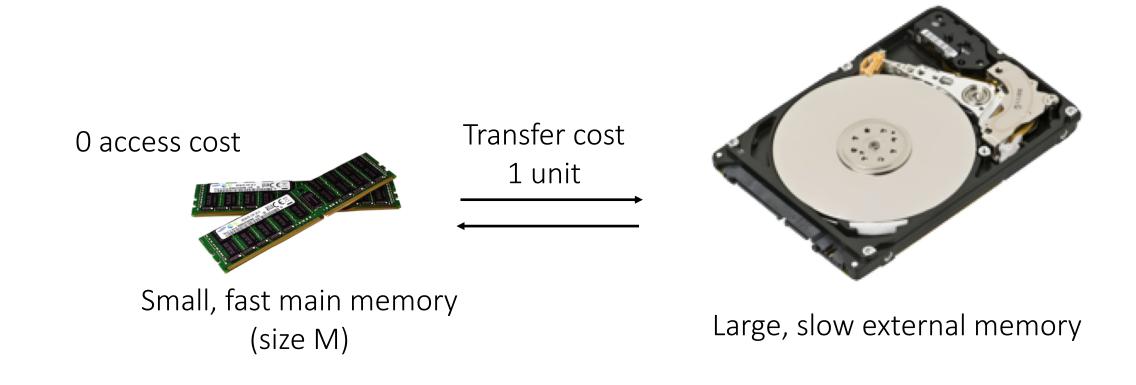


망요 DiSC





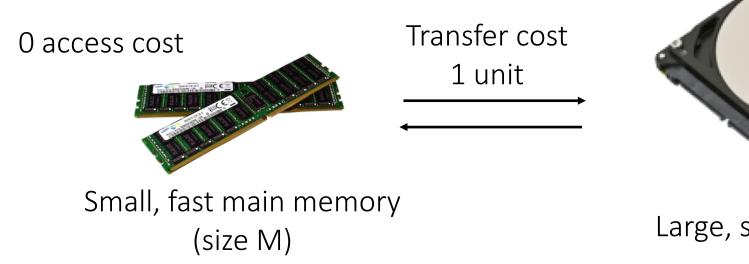
BS dg DiSC







total cost \cong total # reads/writes to disk





Large, slow external memory



Two (outdated) assumptions

- ✓ Symmetric cost for Read & Write to disk
- ✓ One I/O at a time

BS C



Small, fast main memory (size M)



Large, slow external memory





Hard Disk Drives

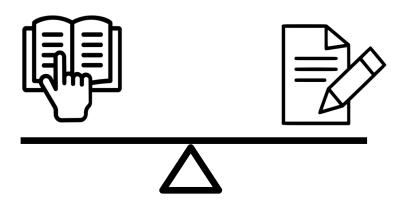


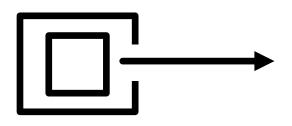




Hard Disk Drives

Two assumptions of the Traditional I/O Model





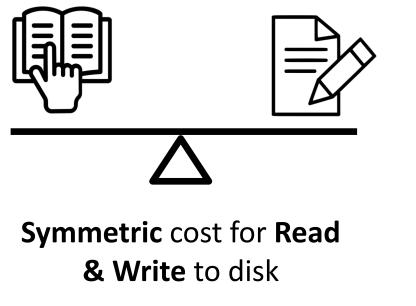
Symmetric cost for Read & Write to disk





Hard Disk Drives

Two assumptions of the Traditional I/O Model









ab SidD Sid

"Tape is Dead. Disk is Tape. Flash is Disk."

- Jim Gray



"Tape is Dead. Disk is Tape. Flash is Disk." - Jim Gray

Bb Iab OSiO

Device	Size	Seq B/W	Time to read
HDD 1980	100 MB	1.2 MB/s	~ 1 min
HDD 2020	4 TB	125 MB/s	~ 9 hours



"Tape is Dead. Disk is Tape. Flash is Disk." - Jim Gray

ि DiSC

Device	Size	Seq B/W	Time to read
HDD 1980	100 MB	1.2 MB/s	~ 1 min
HDD 2020	4 TB	125 MB/s	~ 9 hours

HDDs are moving deeper in the memory hierarchy, and new algorithms are designed for new faster storage devices



"Tape is Dead. Disk is Tape. Flash is Disk." - Jim Gray

<u>망</u>요 DiSC

Device	Size	Seq B/W	Time to read
HDD 1980	100 MB	1.2 MB/s	~ 1 min
HDD 2020	4 TB	125 MB/s	~ 9 hours

HDDs are moving deeper in the memory hierarchy, and new algorithms are designed for new faster storage devices

How do these modern storage devices perform?





ि <u>वि</u> Sid

No mechanical movement





ि S DisC



No mechanical movement

Fast access, High chip density, Low energy consumption











✓ SATA SSD

✓ PCIe SSD



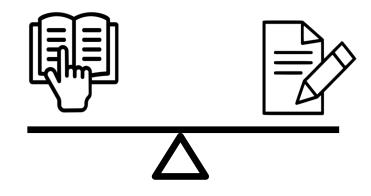


✓ SATA SSD

✓ PCIe SSD

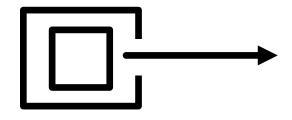
✓ Optane SSD



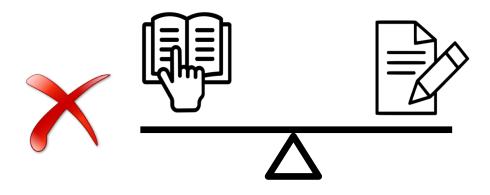


lab Sada DSiO

Symmetric cost for Read & Write

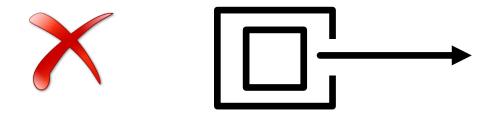






<u>ित्</u> <u>वि</u> DSiD

Symmetric cost for Read & Write

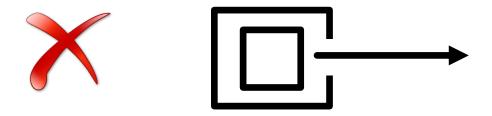






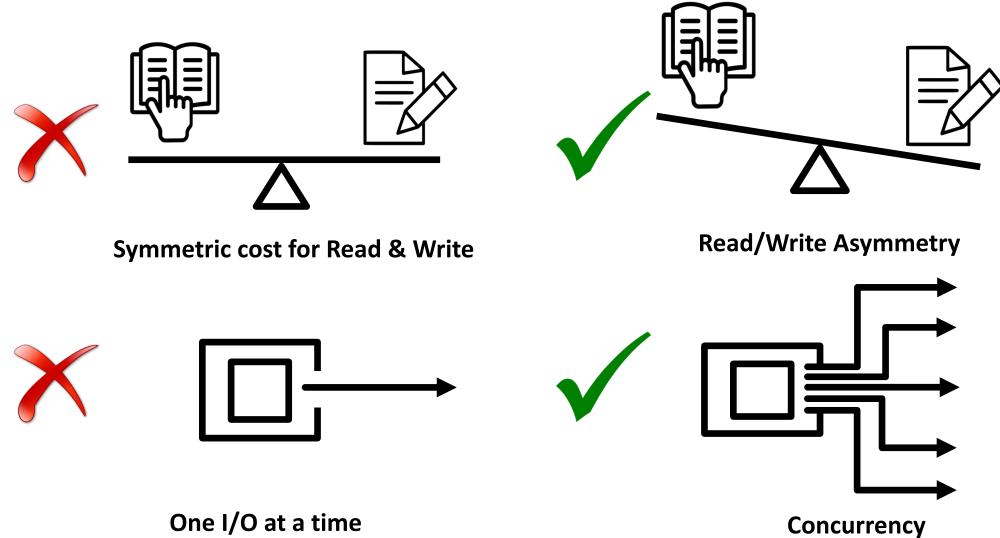
Symmetric cost for Read & Write

Read/Write Asymmetry



ि S S S S S





One I/O at a time

ि S S S S S

31



How should the I/O model be adapted in light of

read/write asymmetry and concurrency?



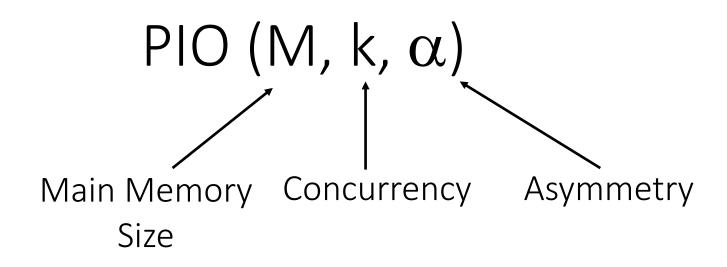


Parametric I/O Model



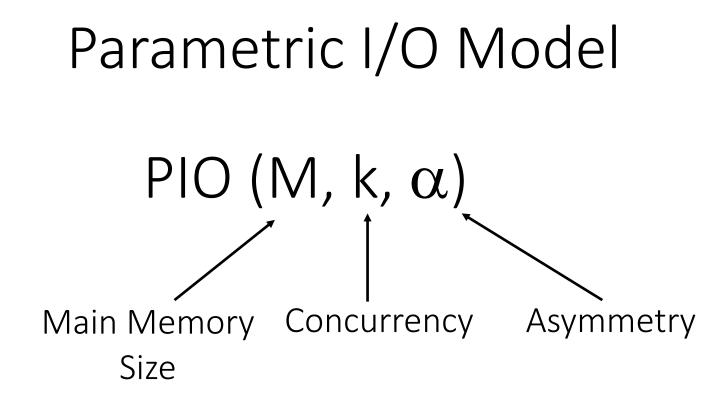
Parametric I/O Model

B8 <u>a</u>8 DiSC







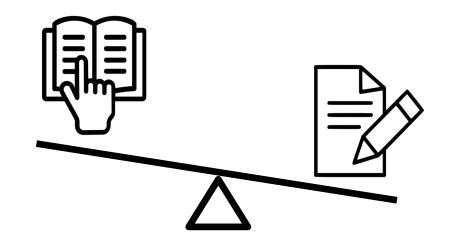


PIO(M, k_r, k_w, \alpha) assumes a fast main memory with capacity *M*, and storage of unbounded capacity that has **read/write asymmetry** α , and **read (write) concurrency** k_r (k_w).



Bb Iab OSiO

Read/Write Asymmetry

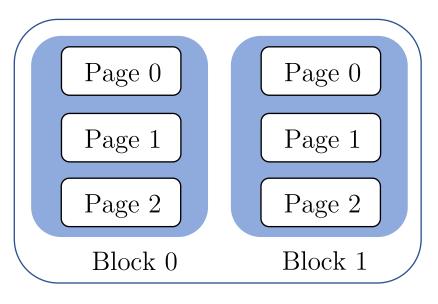




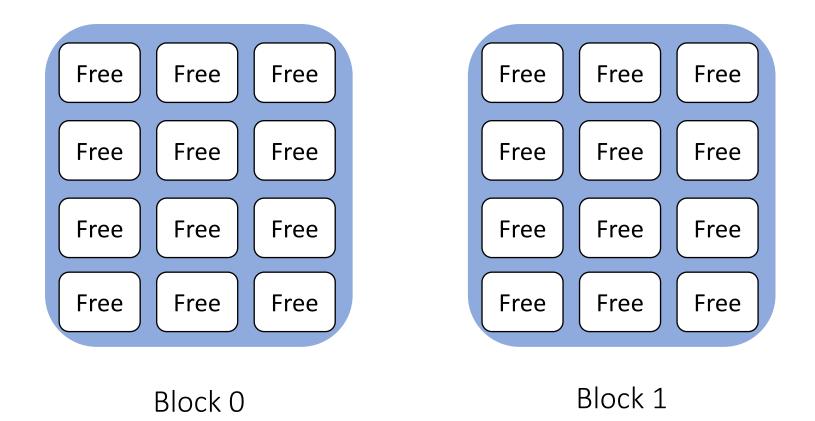


Out-of-place updates cause invalidation

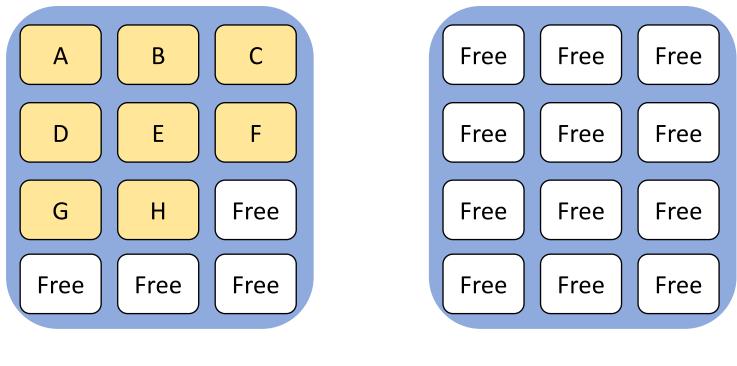
Invalidation causes garbage collection.











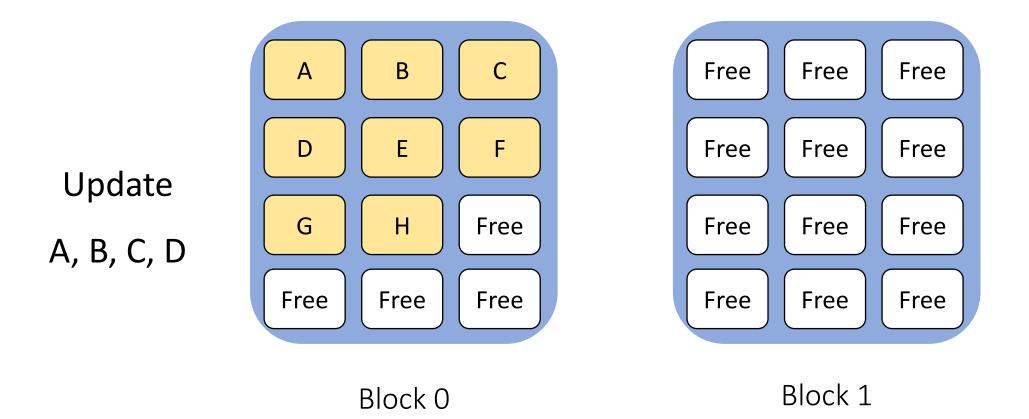
Block 0

lab **S**ad **S**ad

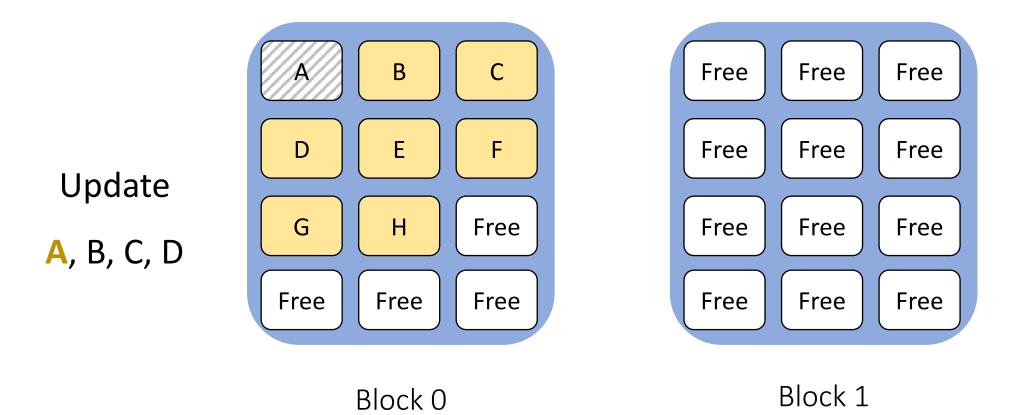
Block 1

Writing in a free page isn't costly!

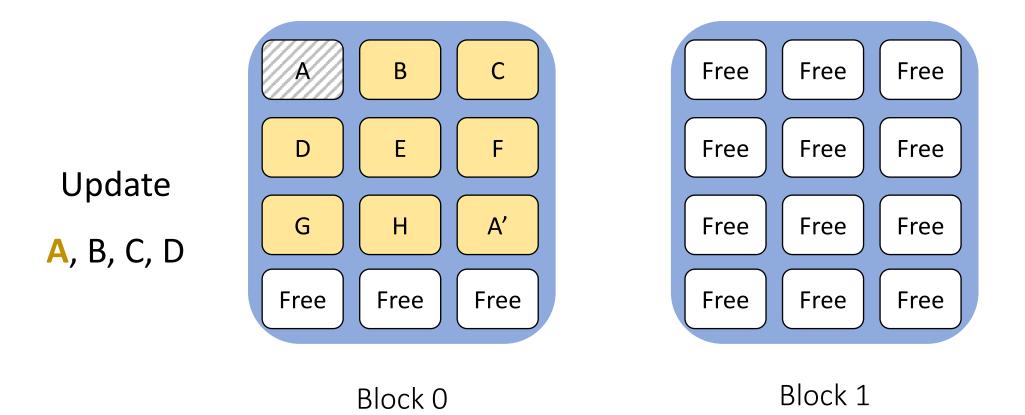




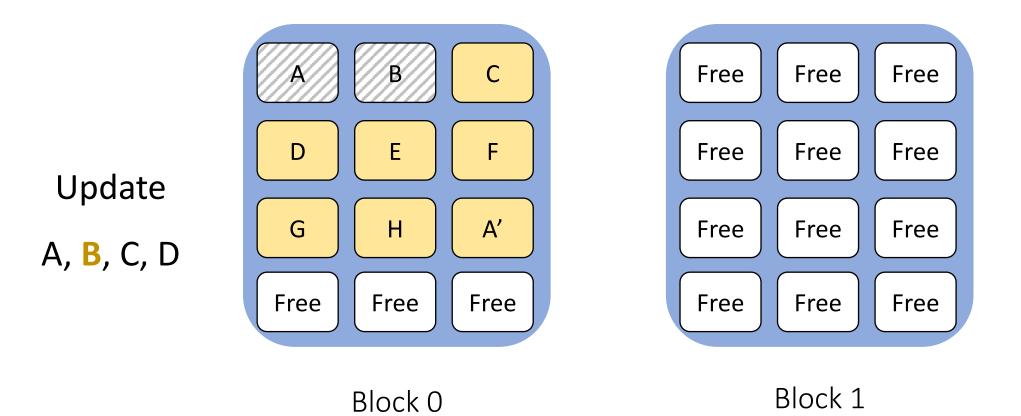




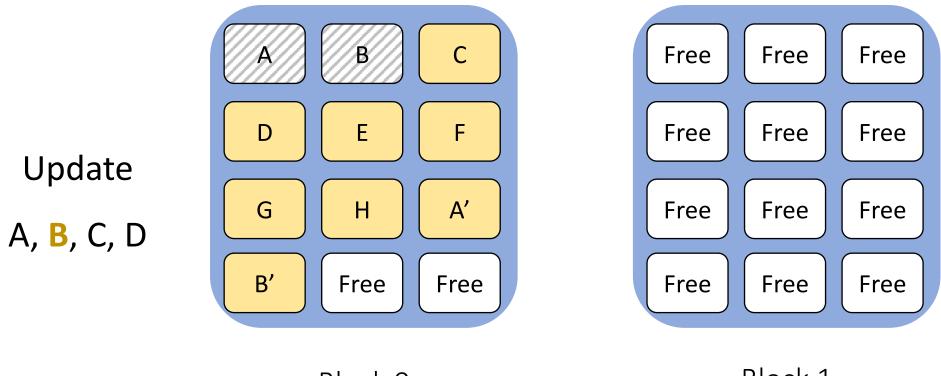










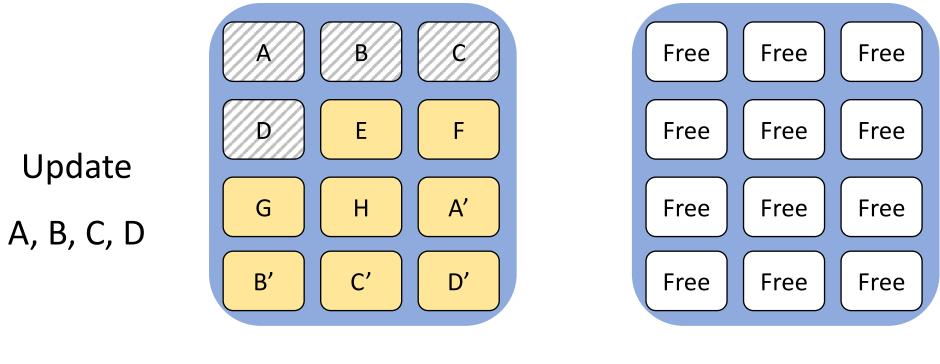


Block 0

lab Sada DSiO

Block 1





Block 0

lab **S**ad **OSiO**

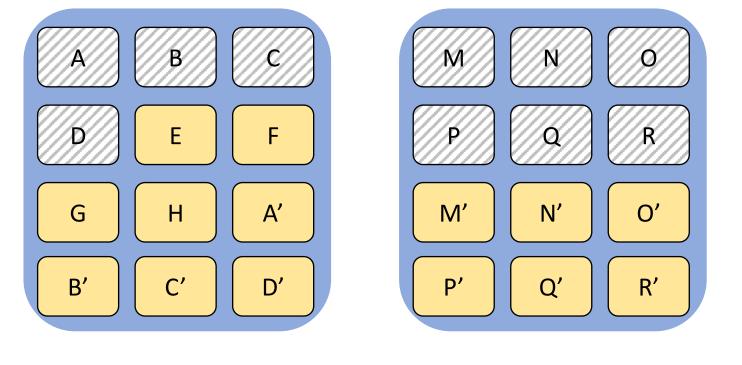
Block 1

Not all updates are costly!





What if there is no space?



. . .

Block 0

Block N

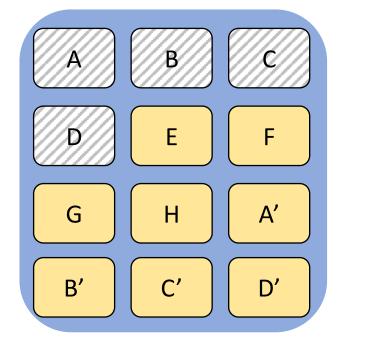


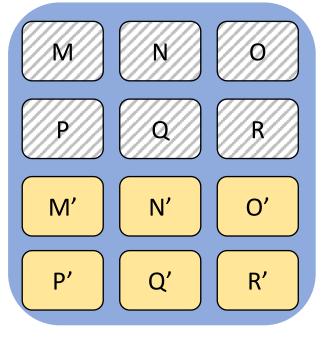
What if there is no space?

lab Sada DSiO



Garbage Collection!





Block 0

. . .

Block N

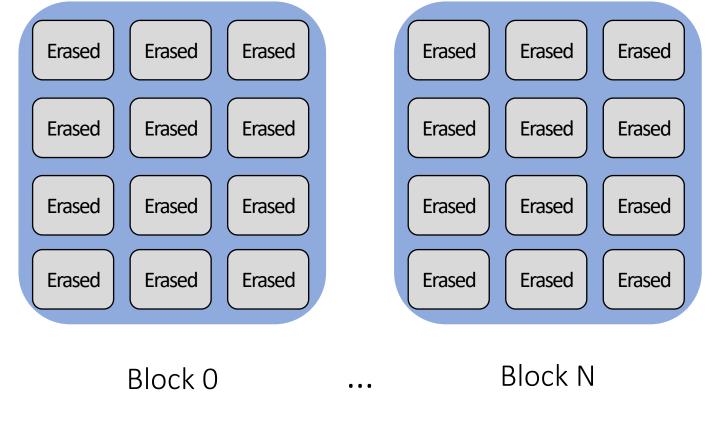


What if there is no space?

lab Sada DSiO



Garbage Collection!



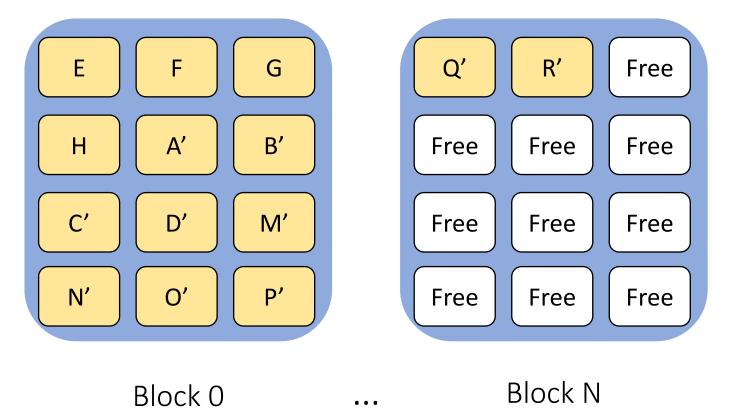


What if there is no space?

lab Sada DSiO

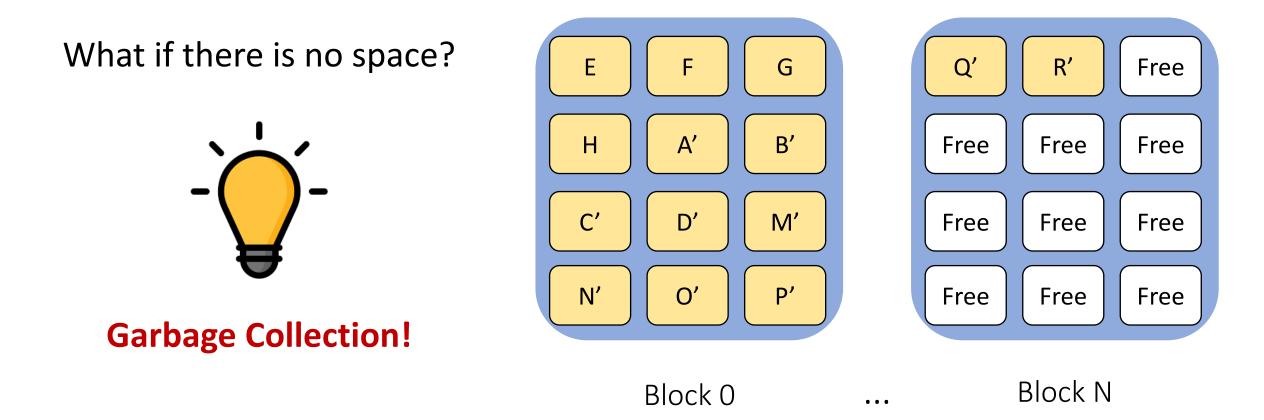


Garbage Collection!





<u>क</u> <u>वि</u> DisC



Higher average update cost (due to GC) → *Read/Write asymmetry*



Read/Write Asymmetry - Example

Device	Advertised Rand Read IOPS	Advertised Rand Write IOPS	Advertised Asymmetry
PCIe D5-P4320	427k	36k	11.9
PCIe DC-P4500	626k	51k	12.3
PCIe P4510	465k	145k	3.2
SATA D3-S4610	92k	28k	3.3
Optane P4800X	550k	500k	1.1



Read/Write Asymmetry - Example

Device	Advertised Rand Read IOPS	Advertised Rand Write IOPS	Advertised Asymmetry
PCIe D5-P4320	427k	36k	11.9
PCIe DC-P4500	626k	51k	12.3
PCIe P4510	465k	145k	3.2
SATA D3-S4610	92k	28k	3.3
Optane P4800X	550k	500k	1.1



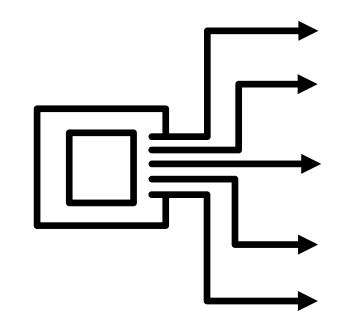
Read/Write Asymmetry - Example

Device	Advertised Rand Read IOPS	Advertised Rand Write IOPS	Advertised Asymmetry
PCIe D5-P4320	427k	36k	11.9
PCIe DC-P4500	626k	51k	12.3
PCIe P4510	465k	145k	3.2
SATA D3-S4610	92k	28k	3.3
Optane P4800X	550k	500k	1.1



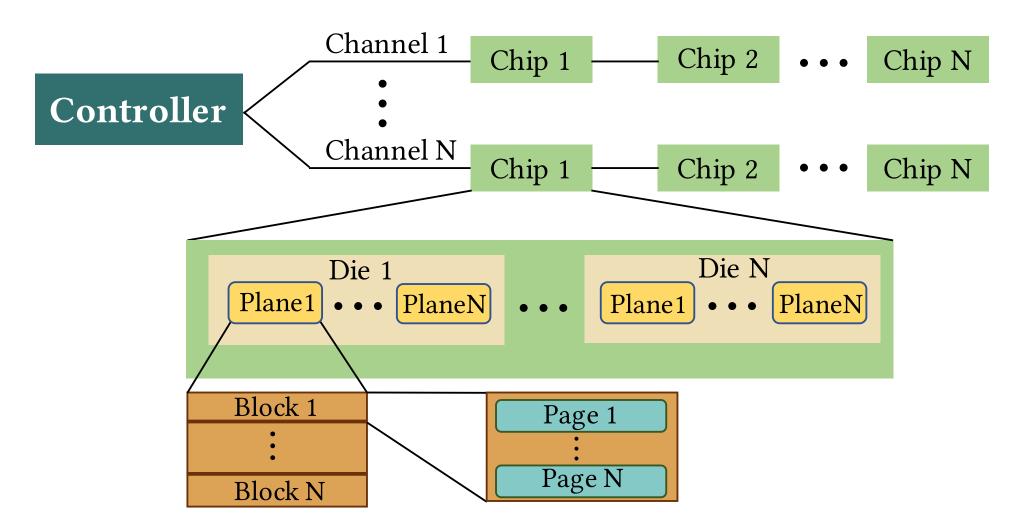
Concurrency

Dlab Siad



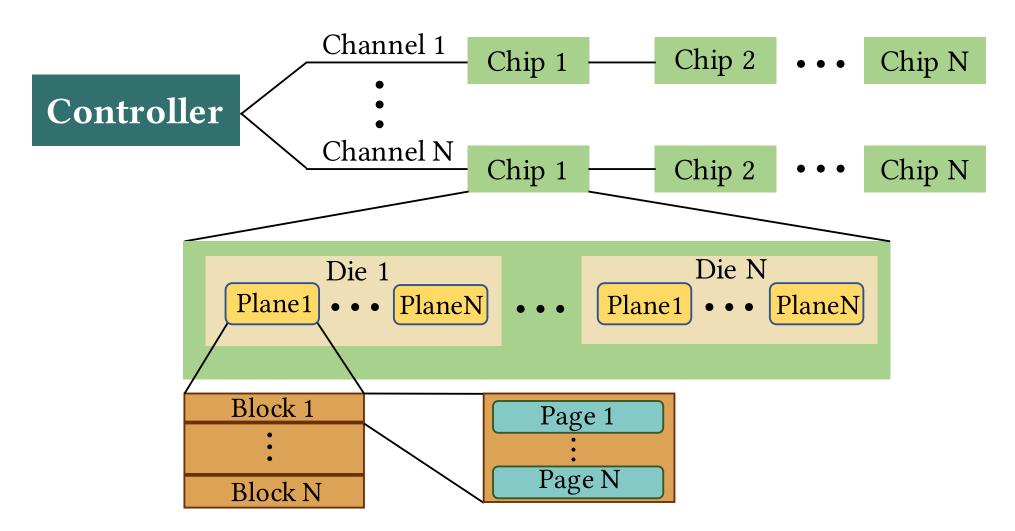
Internals of an SSD

lap de Calanda



Internals of an SSD

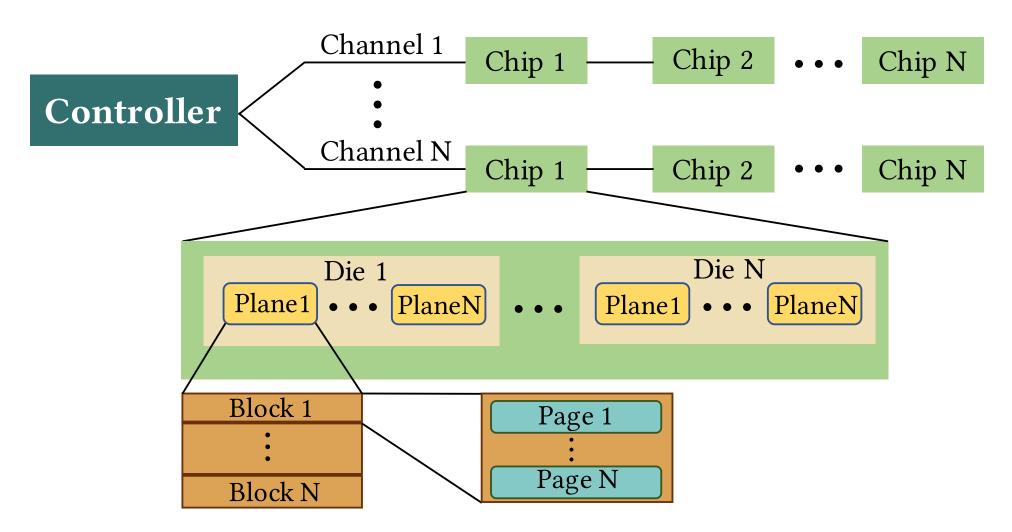
lap de Calanda



Parallelism at different levels (channel, chip, die, plane block, page)

Internals of an SSD

lab Sada DSiO



Essential to know the concurrency of a device!



lab Sada DSiO

How can we quantify

read/write asymmetry and concurrency?



lab **S**ad **OSiO**





Tools

- ✓ Custom micro-benchmarking
- ✓ Fio
- ✓ Intel's SPDK





Tools

Setup

✓ Custom micro-benchmarking

with and without FS

✓ Fio

✓ Intel's SPDK





Devices

- ✓ Optane SSD
- ✓ PCIe SSD with & w/o FS
- ✓ SATA SSD



Devices

B8 <u>9</u> DiSC

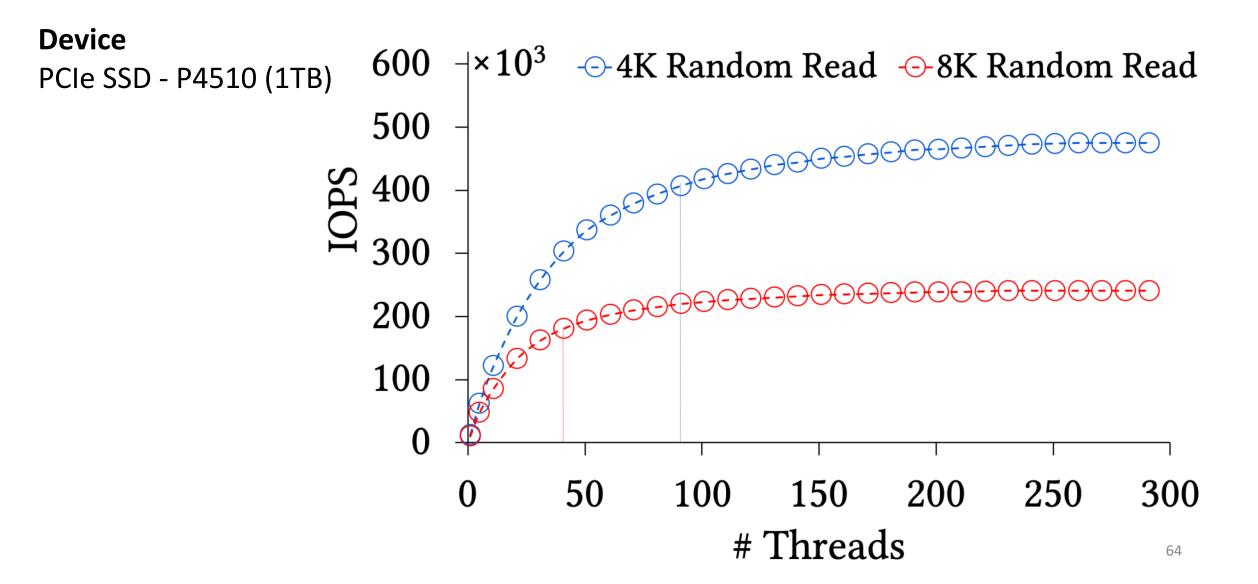
Setup

- ✓ Optane SSD
- ✓ PCIe SSD with & w/o FS ✓
- ✓ Increase #threads issuing I/Os
 - Report: latency and throughput

✓ SATA SSD

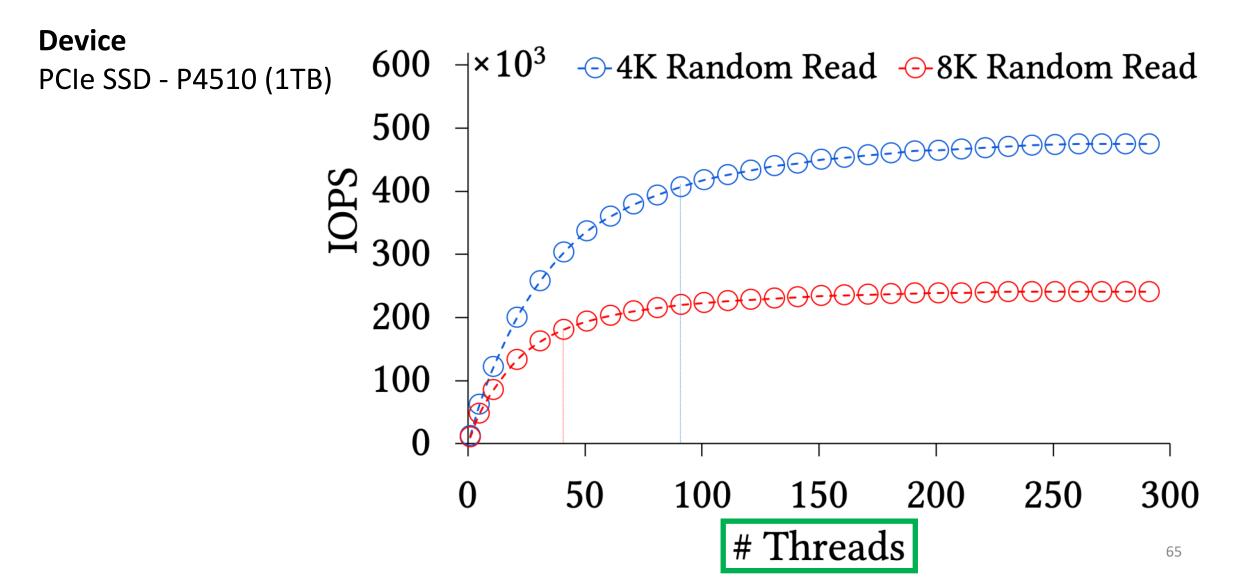






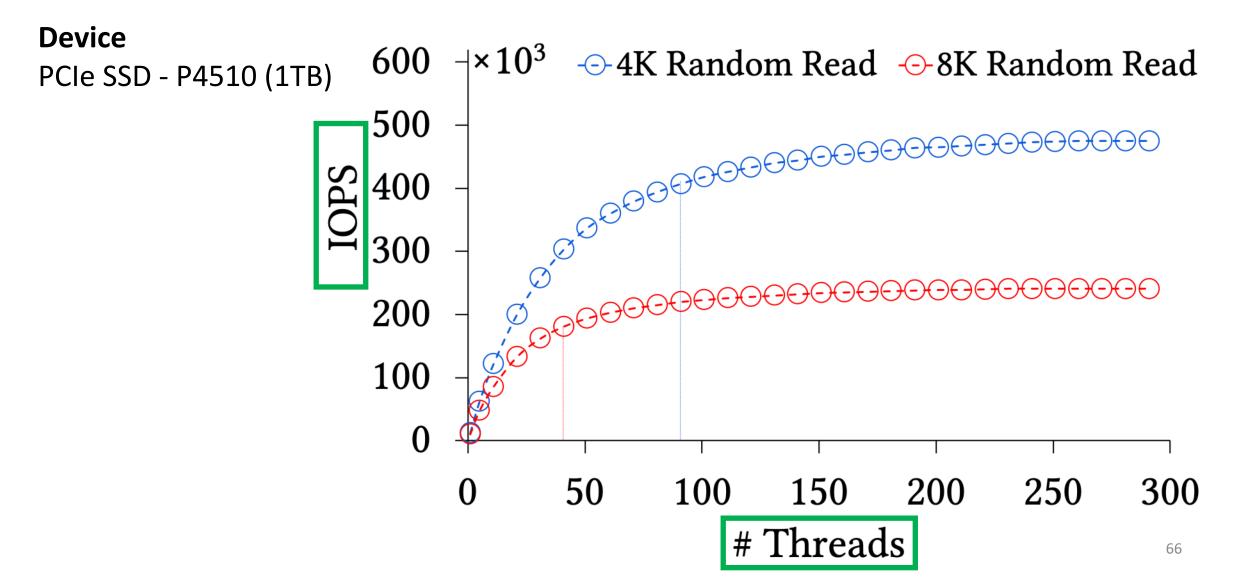






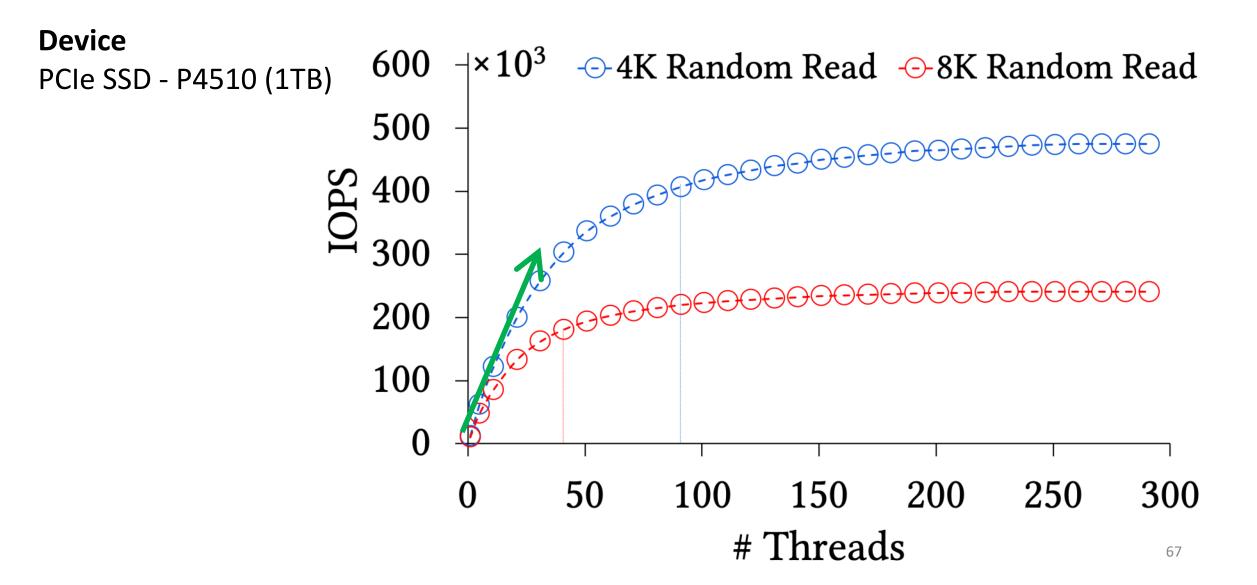






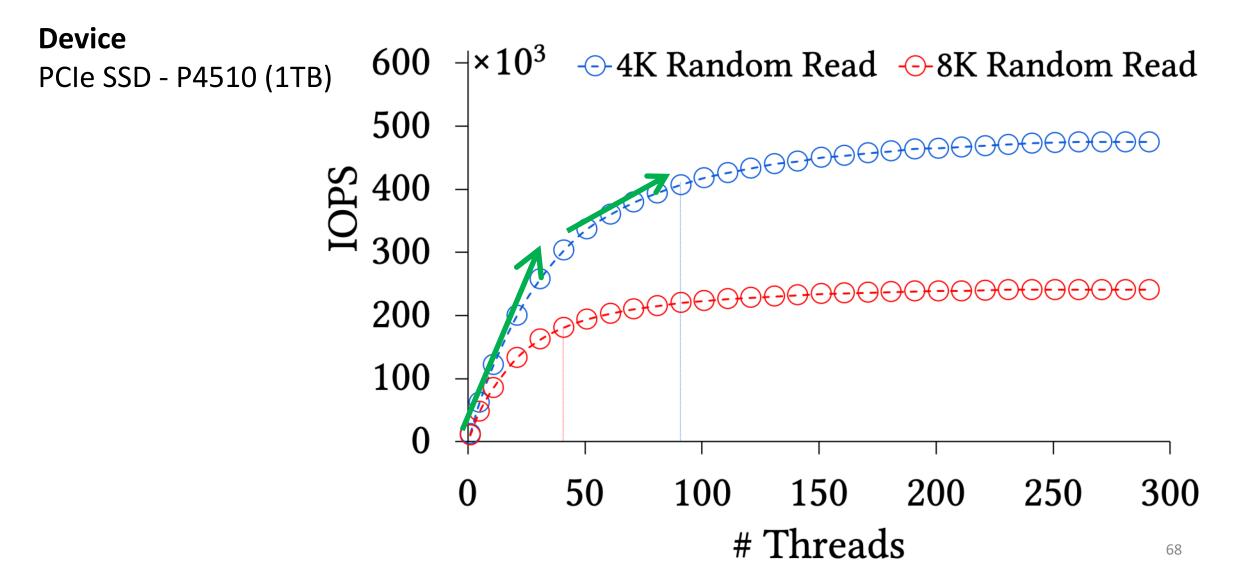






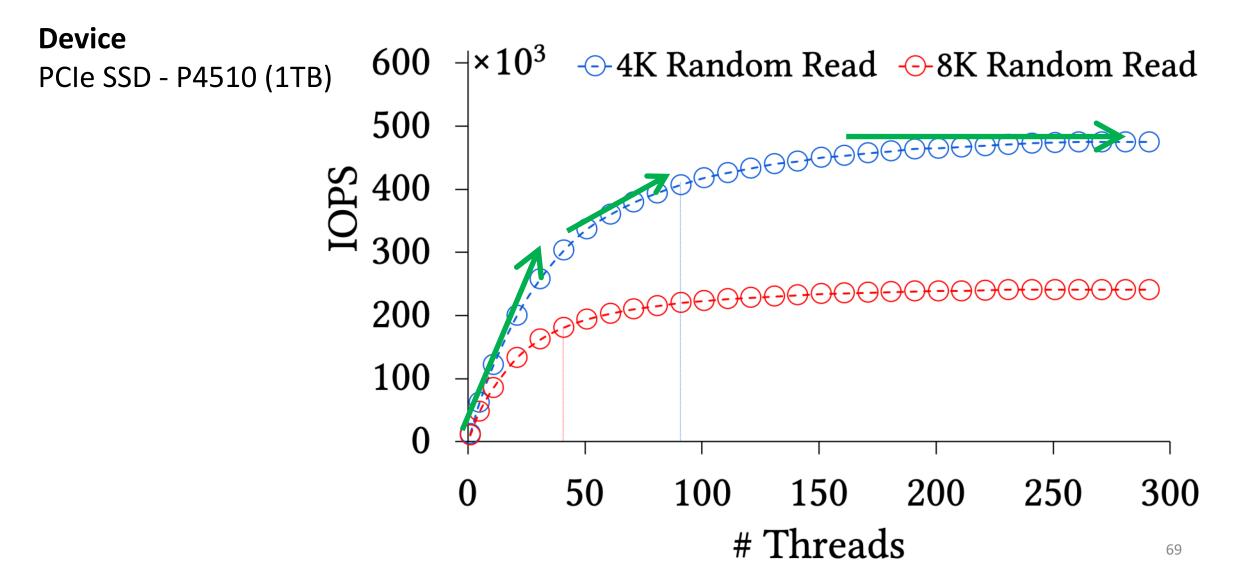






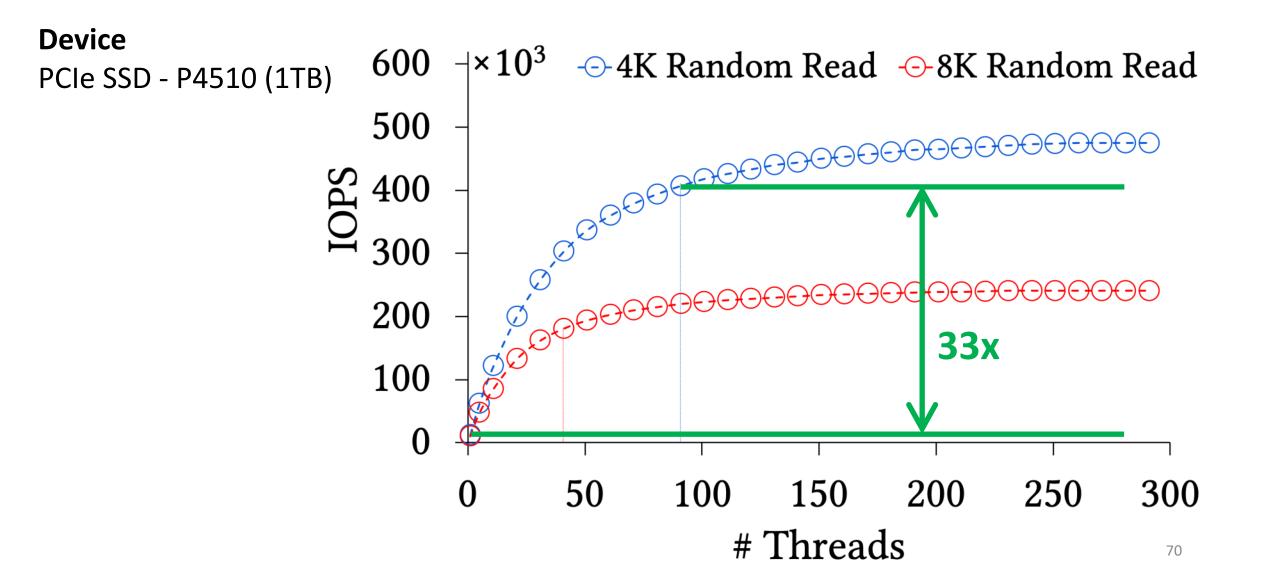








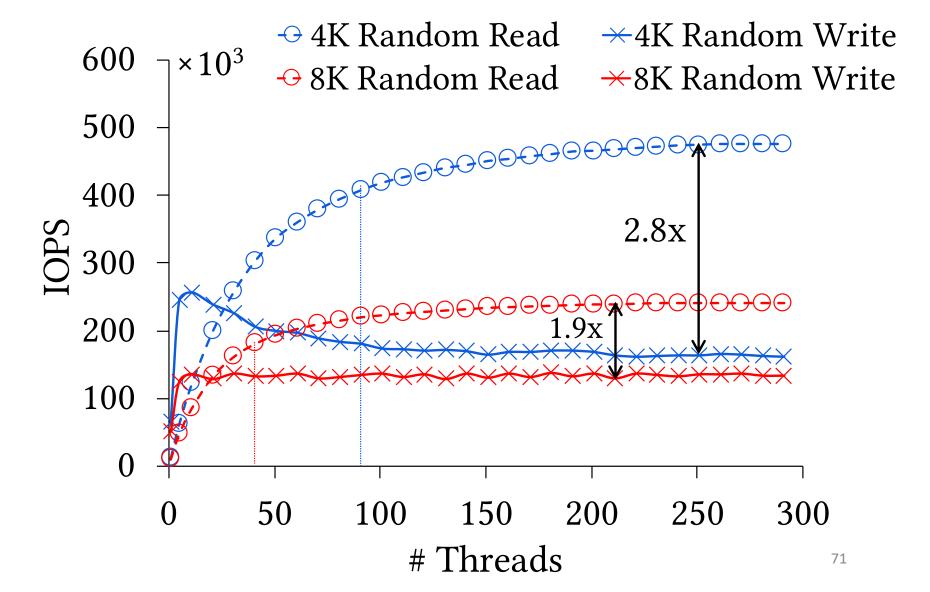






Device PCIe SSD - P4510 (1TB)

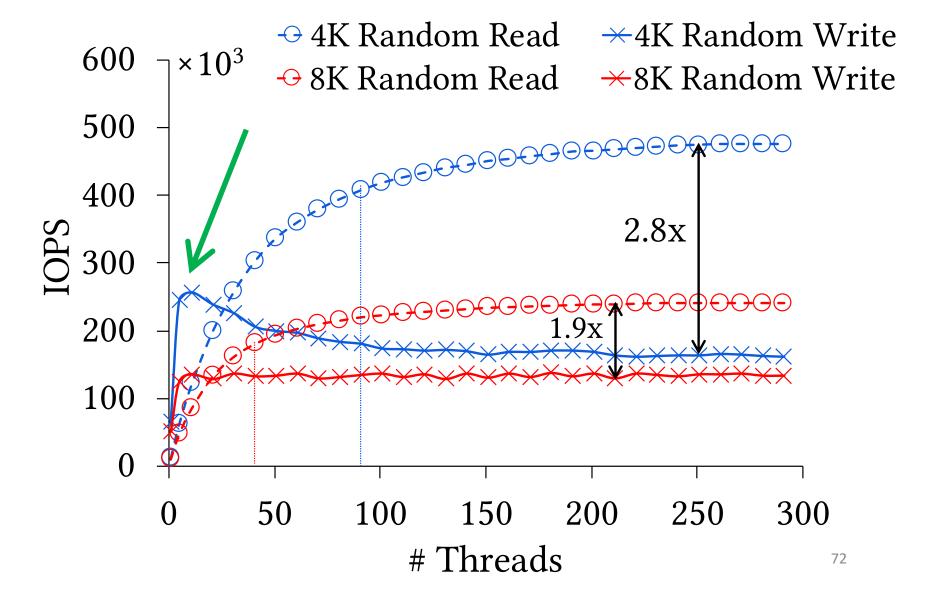
ab SidD Sid





Device PCIe SSD - P4510 (1TB)

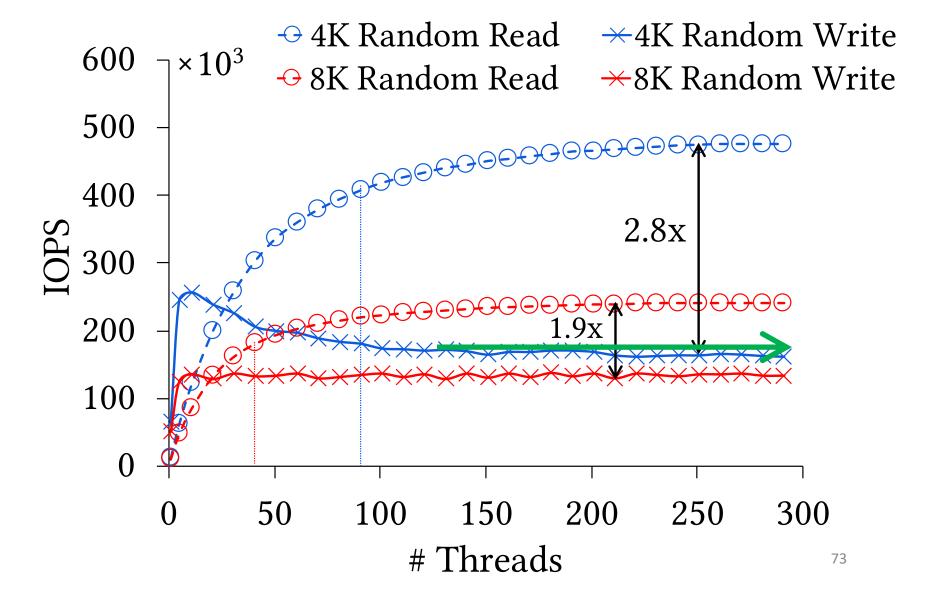
ab Iab **OSiO**





Device PCIe SSD - P4510 (1TB)

ab Iab **OSiO**

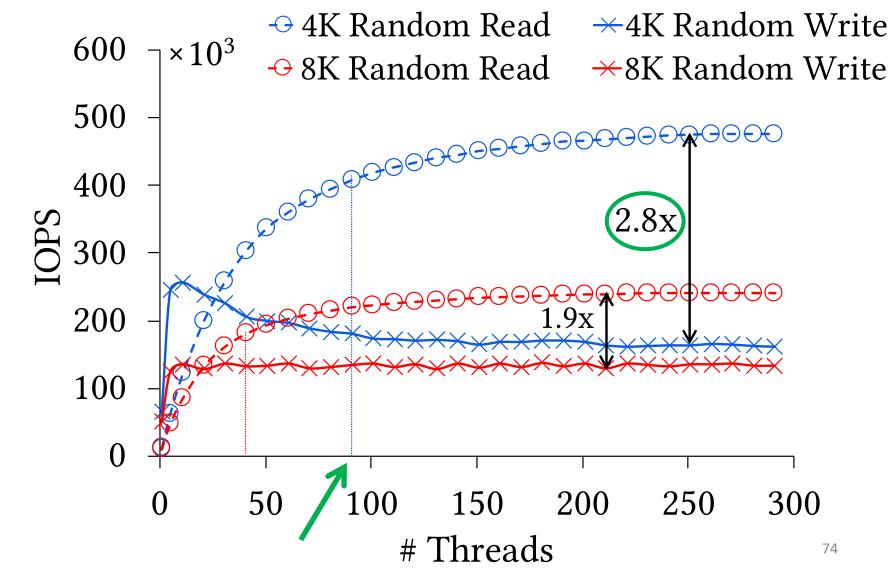




Device PCIe SSD - P4510 (1TB) For 4K random read,

Asymmetry: 2.8

Concurrency: 80





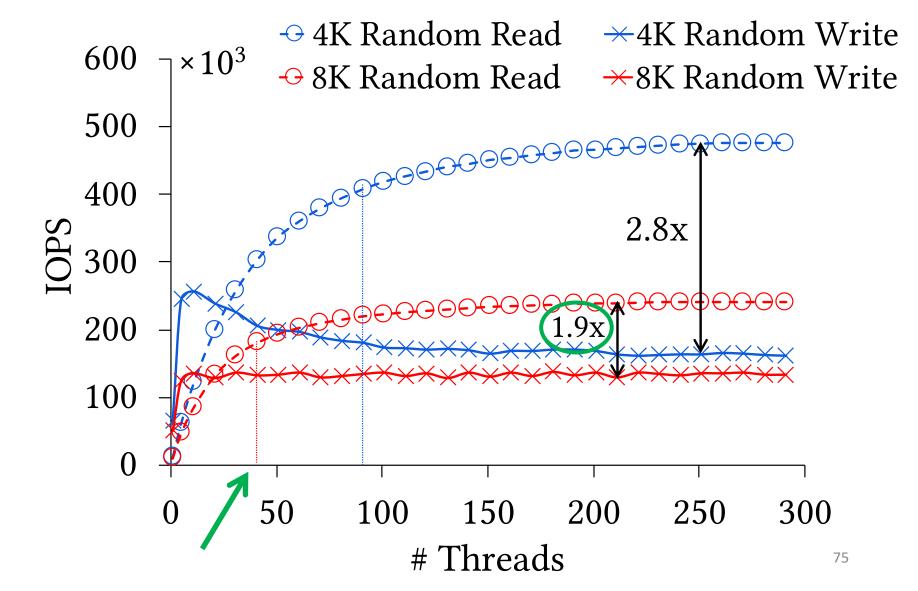


Device PCIe SSD - P4510 (1TB)

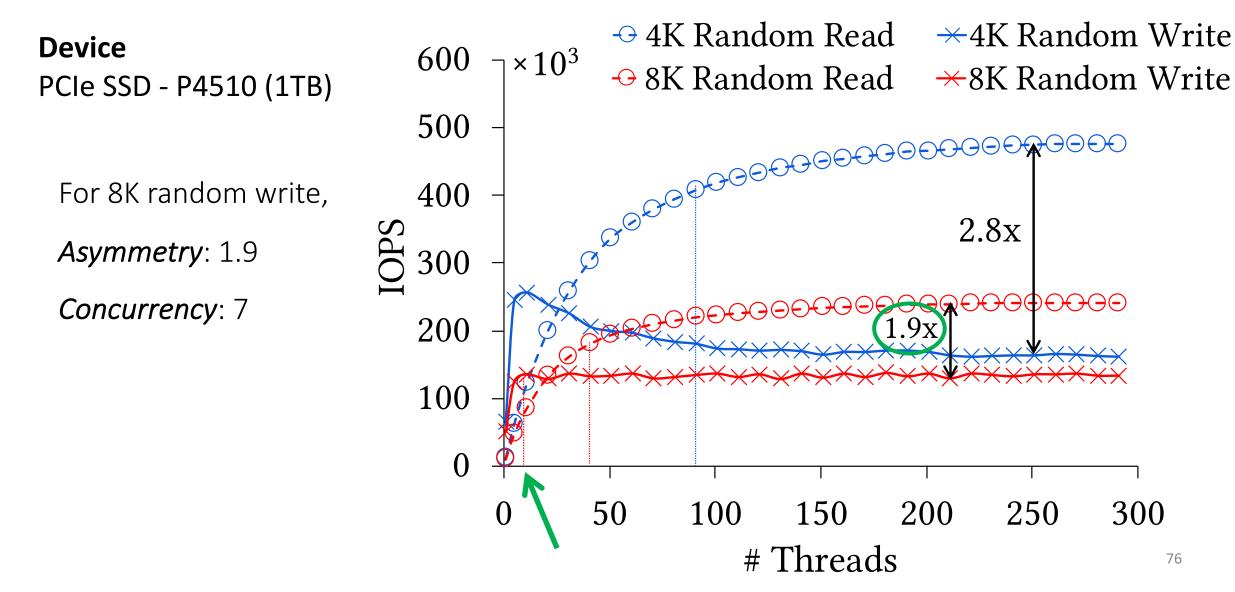
For 8K random read,

Asymmetry: 1.9

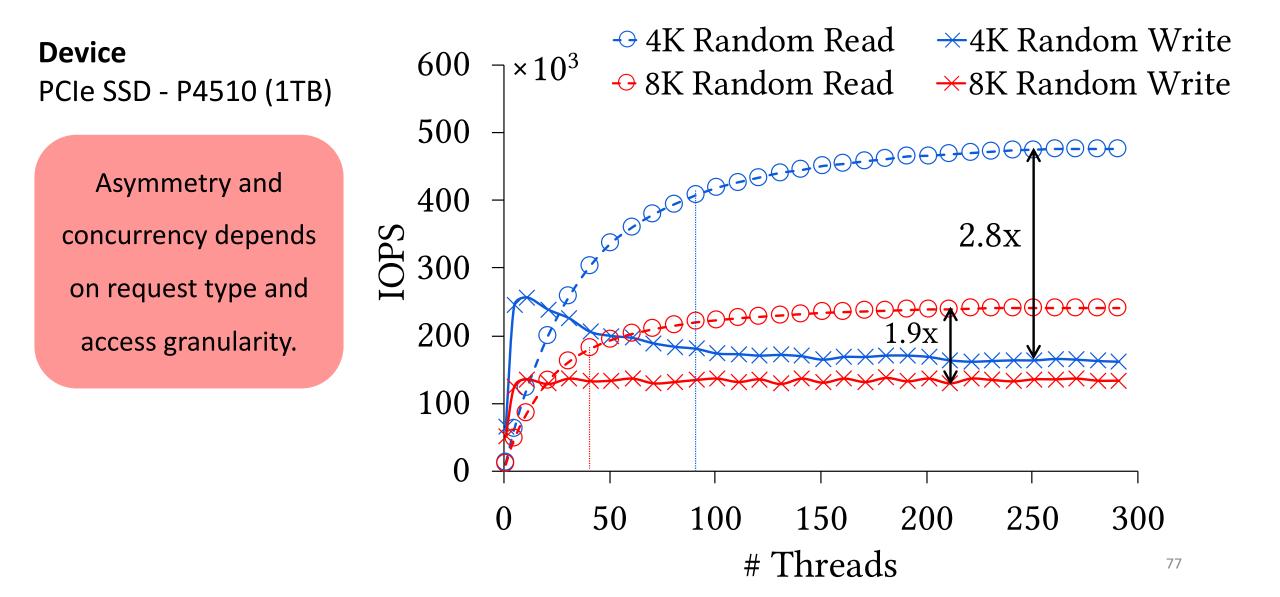
Concurrency: 40





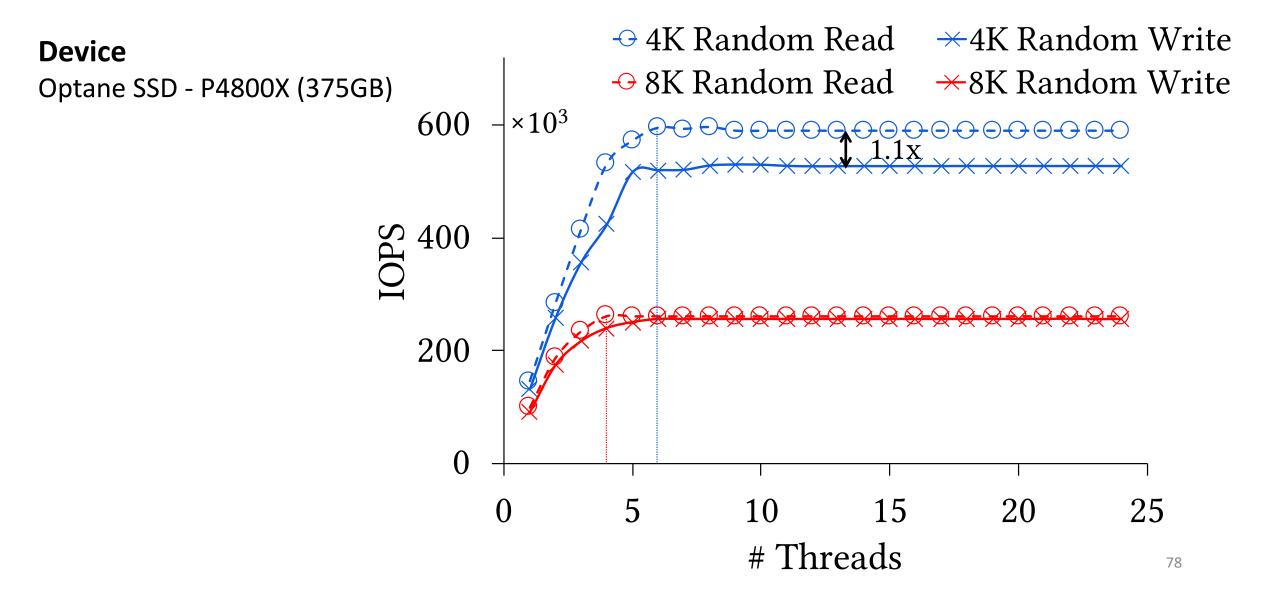






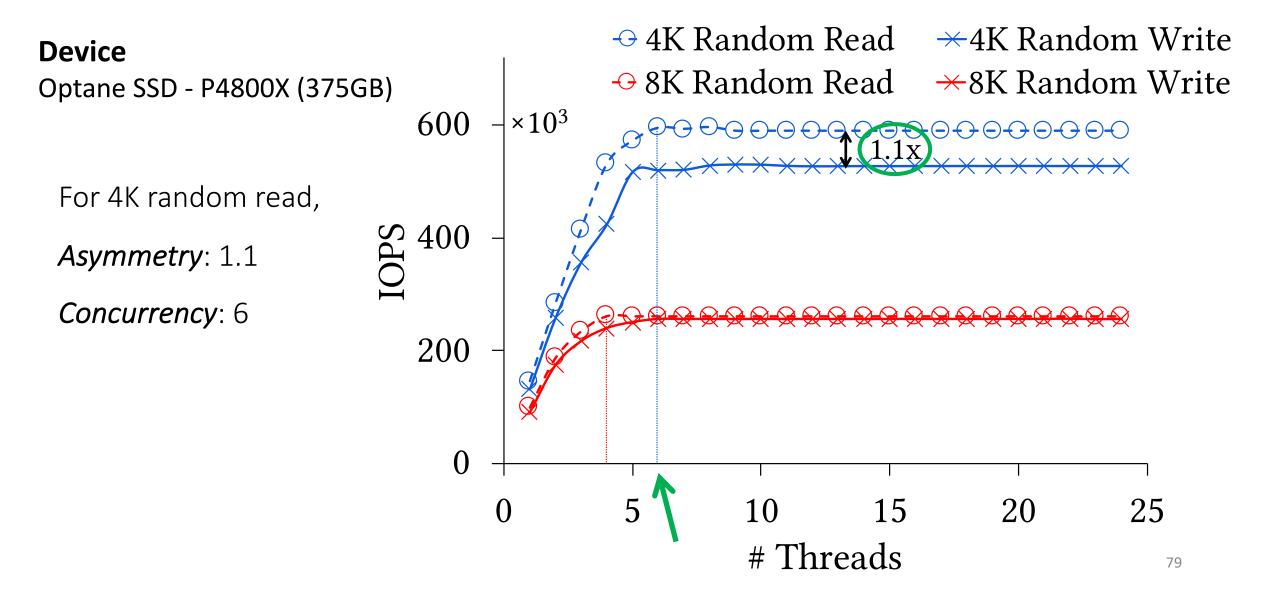


B8 <u>9</u> DiSC



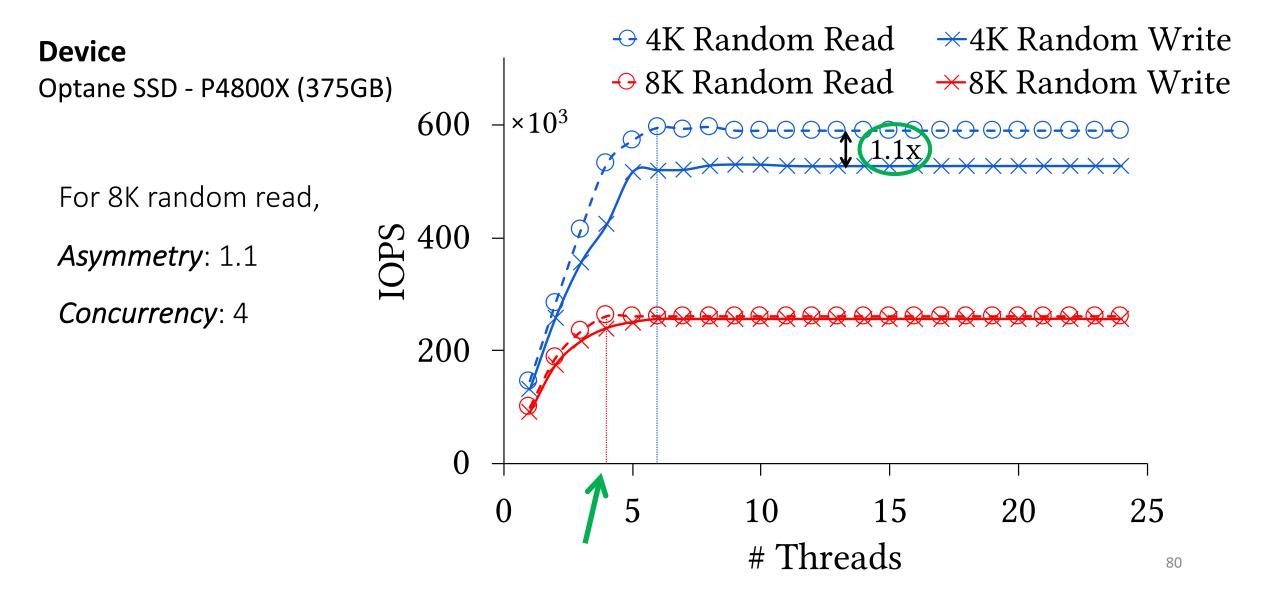


BS dg DiSC



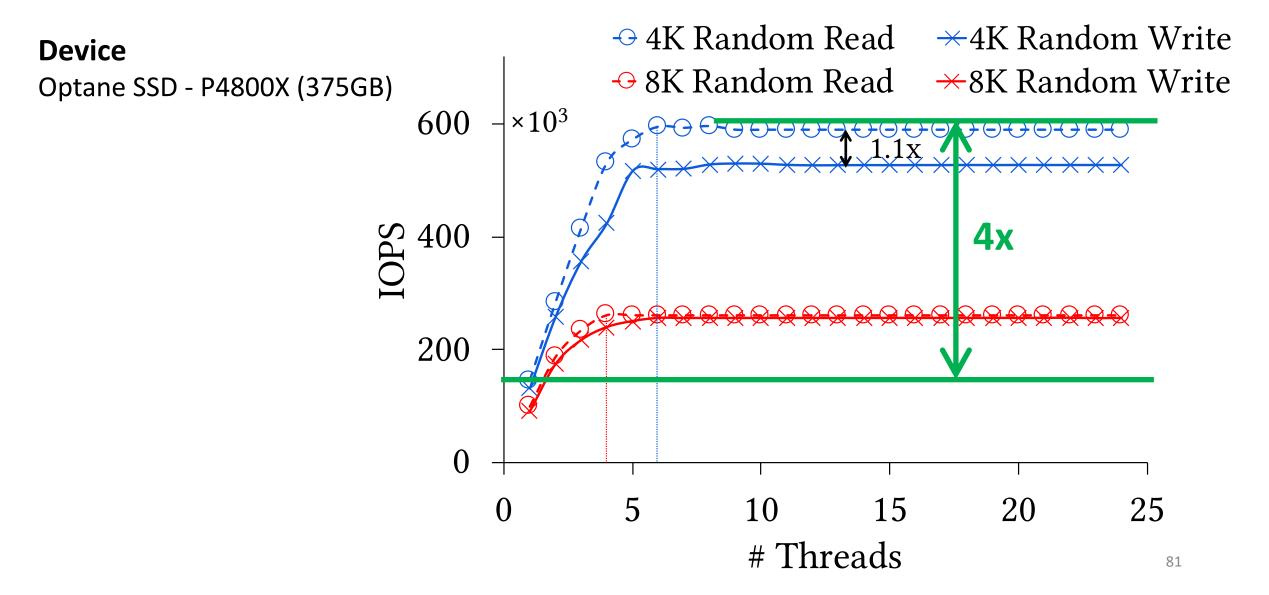


BS dg DiSC





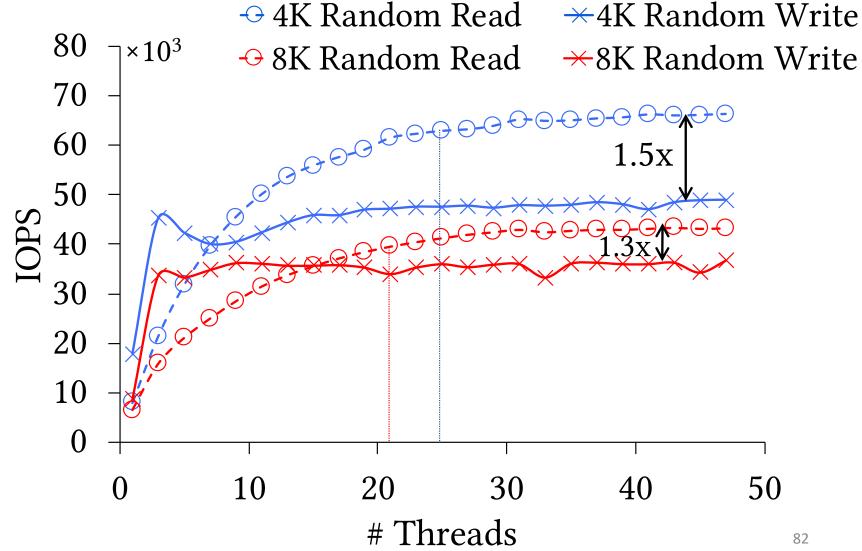
lab SSID





Device SATA SSD - S4610 (240GB)

ab Iab **OSiO**





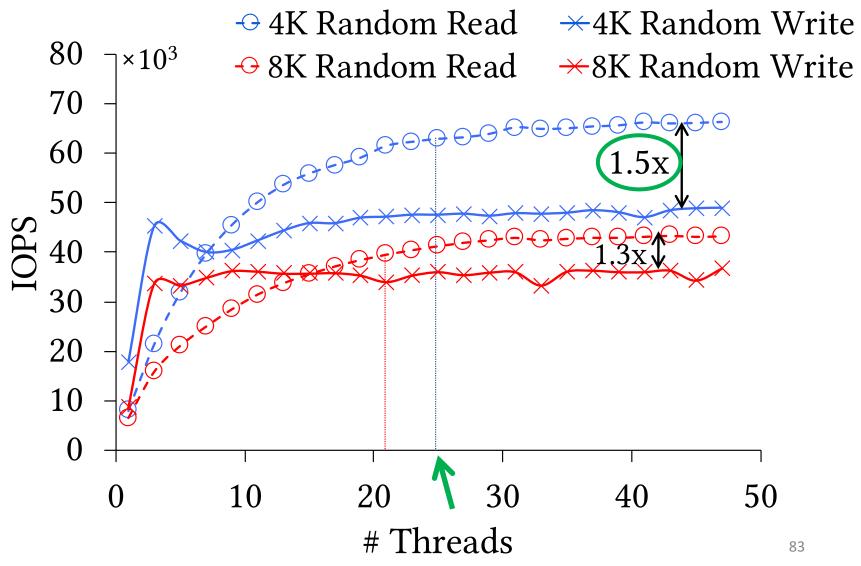
Device SATA SSD - S4610 (240GB)

For 4K random read,

Asymmetry: 1.5

B8 <u>9</u> DiSC

Concurrency: 25





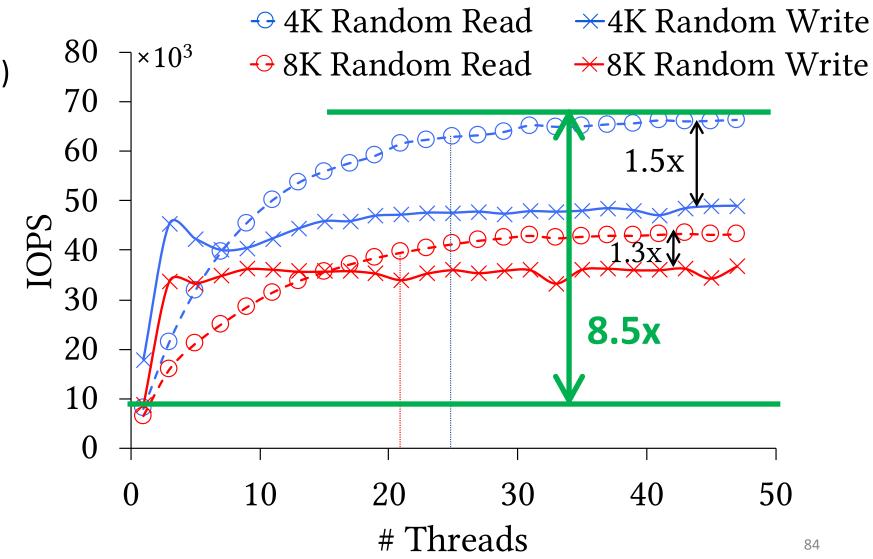
Device SATA SSD - S4610 (240GB)

For 4K random read,

Asymmetry: 1.5

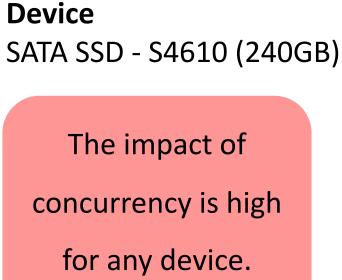
B8 <u>9</u> DiSC

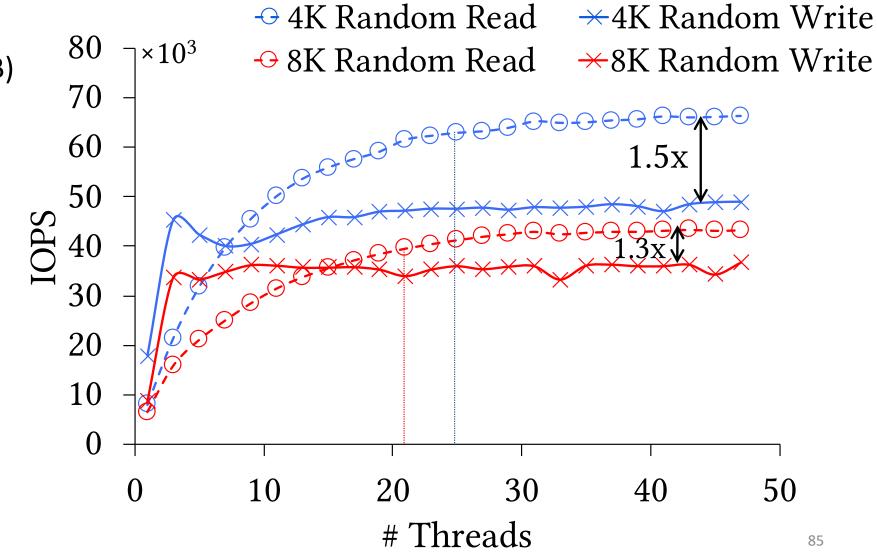
Concurrency: 25







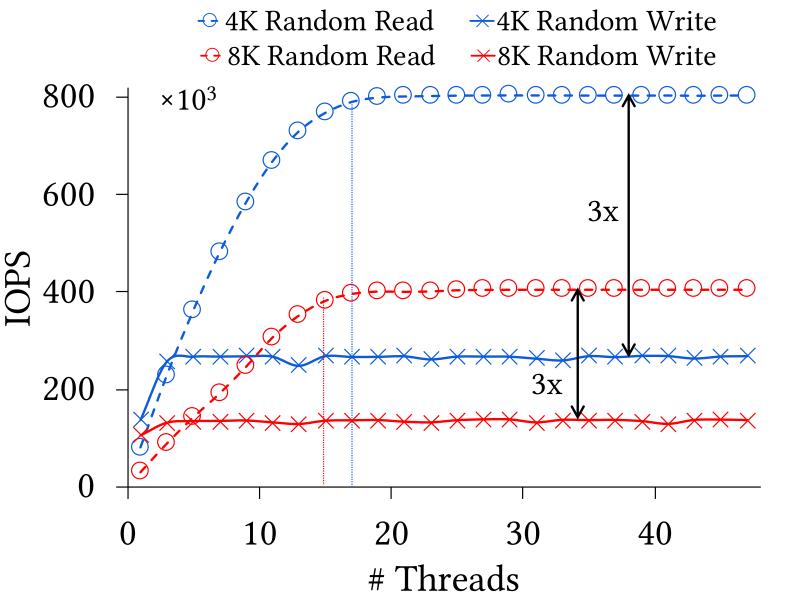




Device

lab **Sid**

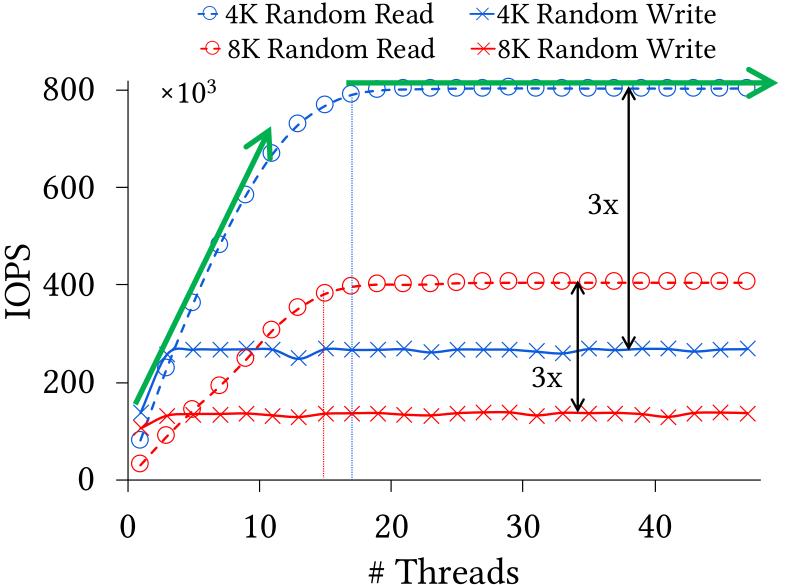
PCIe SSD - P4510 (1TB)



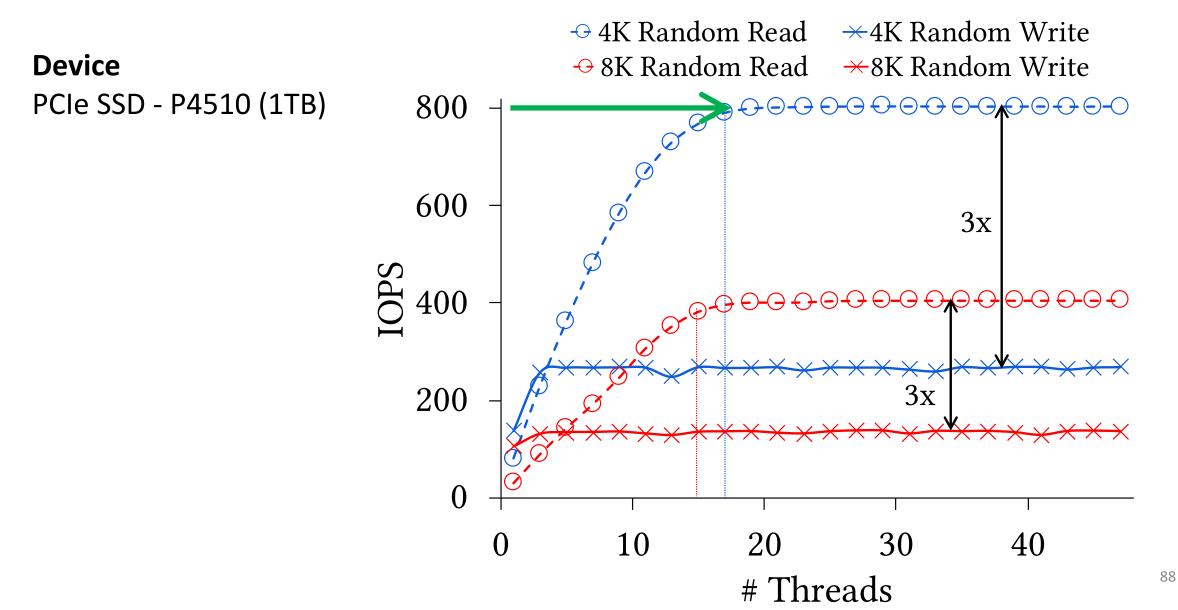
Device

lab OSIC

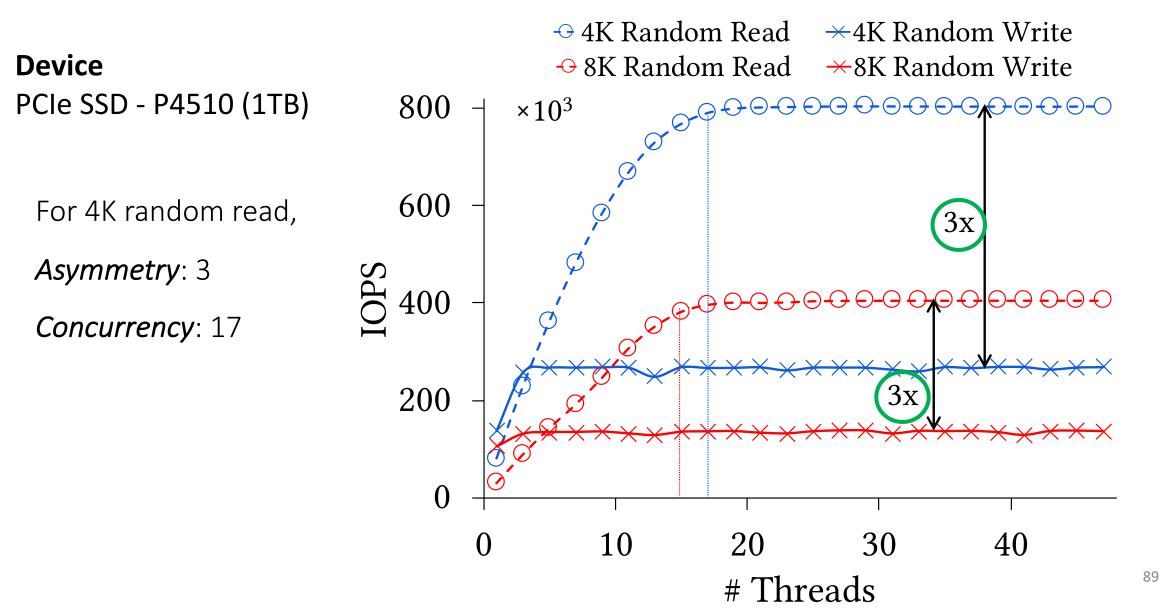
PCIe SSD - P4510 (1TB)



lab OSIC

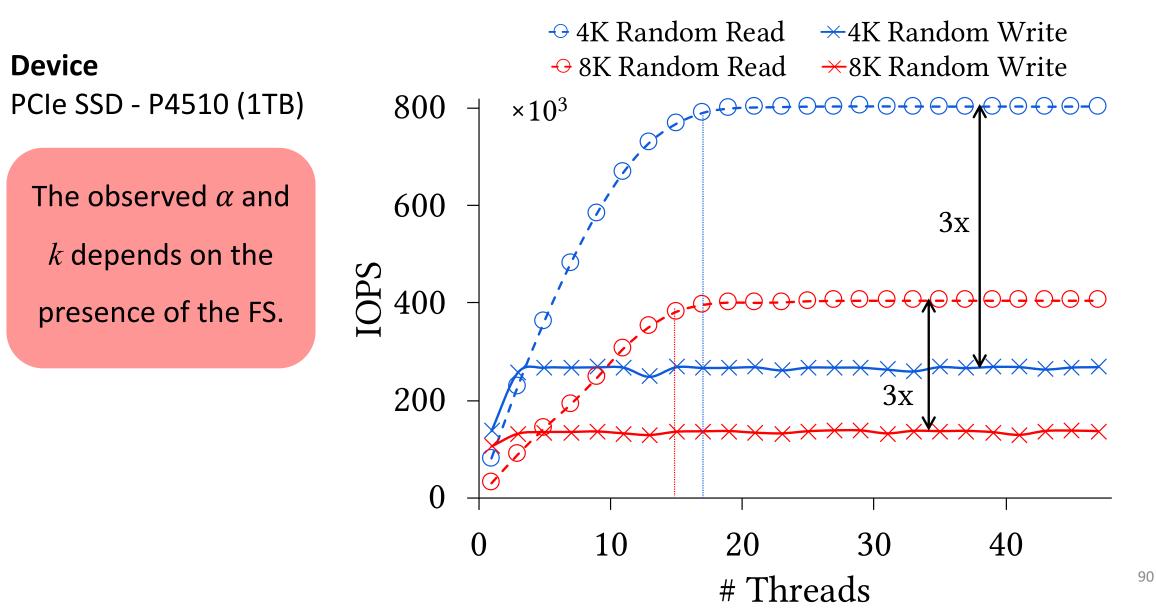


<u>क</u> <u>वि</u> DisC





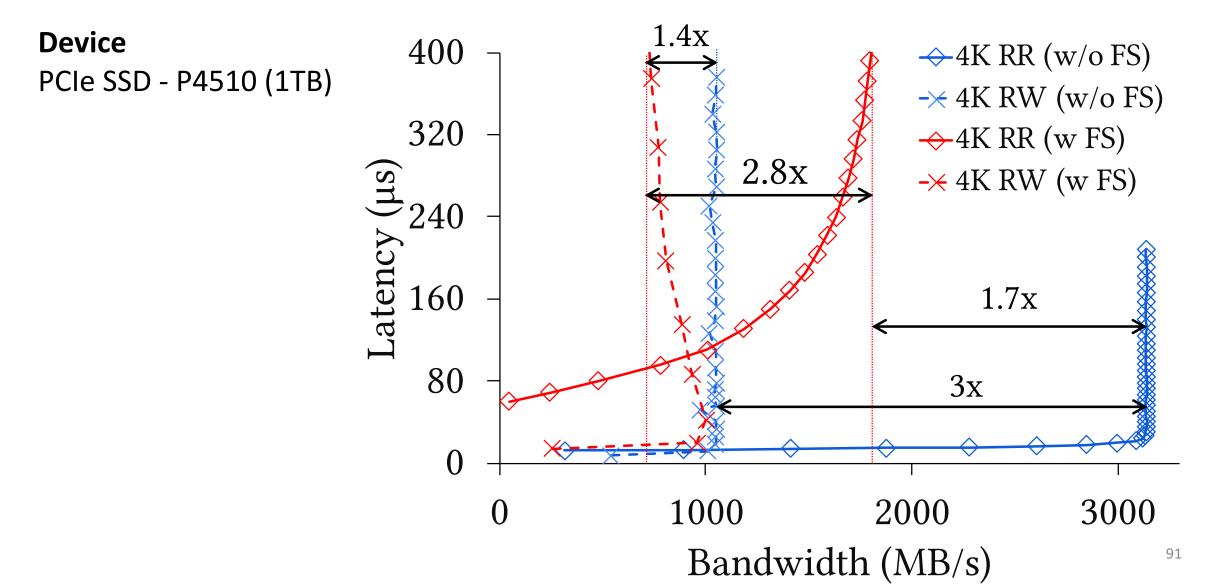
ि DiSC





Impact of the File System

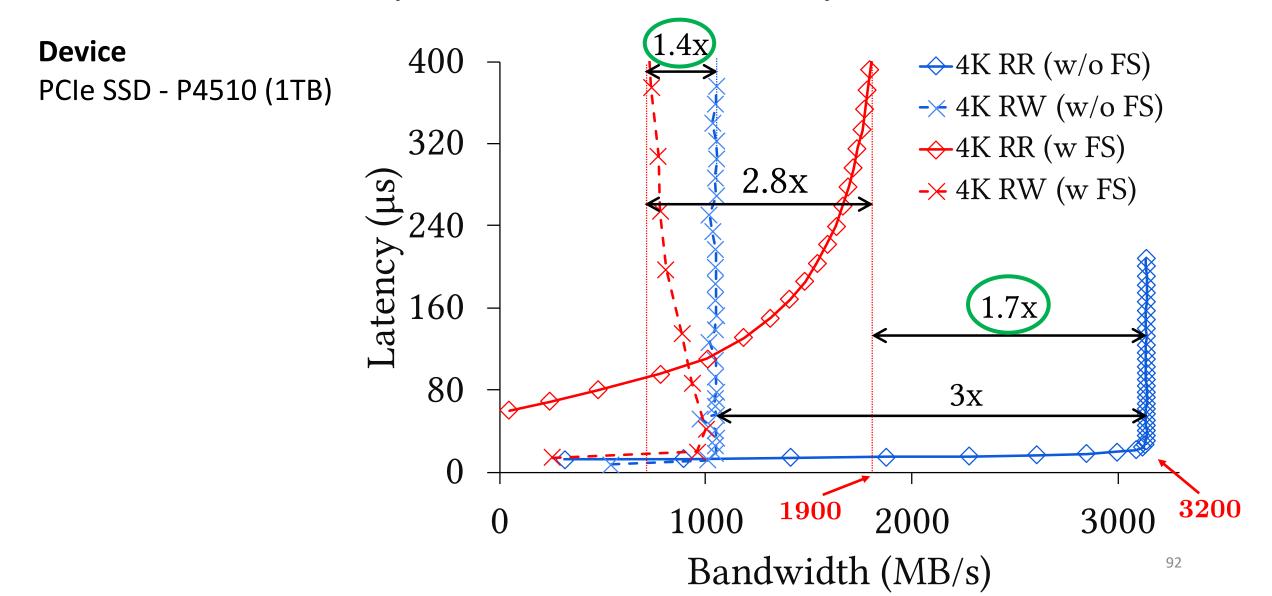
lab **S**ad **OSiO**





Impact of the File System

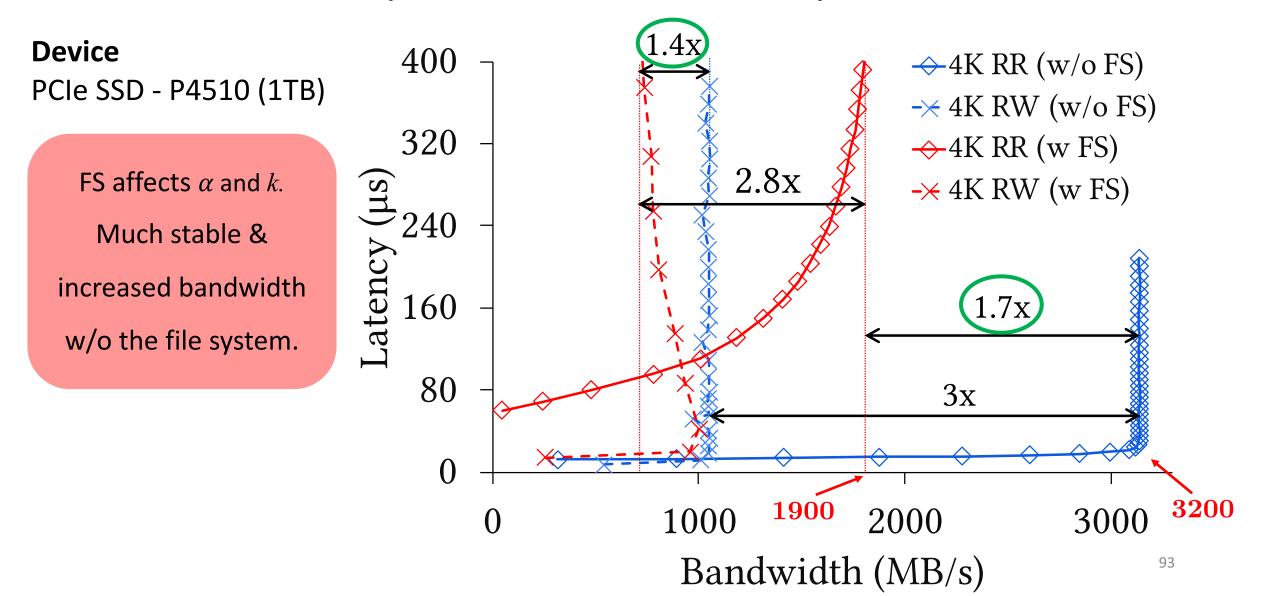
lab **S**ad **OSiO**





Impact of the File System

lab OSIC





Empirical Asymmetry and Concurrency

lab Sada DSiO

	4KB			8KB		
Devices	α	k_r	k_w	α	k_r	k_w
Optane SSD	1.1	6	5	1.0	4	4
PCle SSD (with FS)	2.8	80	8	1.9	40	7
PCIe SSD (w/o FS)	3.0	16	6	3.0	15	4
SATA SSD	1.5	25	9	1.3	21	5



Empirical Asymmetry and Concurrency

lab Sada DSiO

	4KB			8KB		
Devices	α	k_r	k_w	α	k _r	k_w
Optane SSD	1.1	6	5	1.0	4	4
PCIe SSD (with FS)	2.8	80	8	1.9	40	7
PCIe SSD (w/o FS)	3.0	16	6	3.0	15	4
SATA SSD	1.5	25	9	1.3	21	5



Empirical Asymmetry and Concurrency

lab **S**ad **OSiO**

	4KB			8KB		
Devices	α	k_r	k_w	α	k_r	k_w
Optane SSD	1.1	6	5	1.0	4	4
PCIe SSD (with FS)	2.8	80	8	1.9	40	7
PCIe SSD (w/o FS)	3.0	16	6	3.0	15	4
SATA SSD	1.5	25	9	1.3	21	5



ab SidD Sid

Performance Analysis following PIO (M, k, α)





Performance Analysis

We classify storage-intensive applications into three classes

- ✓ Batchable Writes
- ✓ Batchable Reads
- ✓ Batchable Reads and Writes





 \checkmark Can exploit write concurrency (k_w) by batching writes



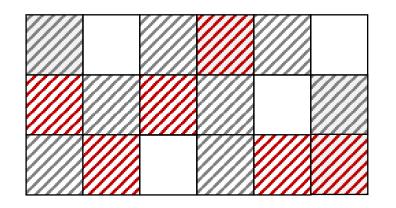


- \checkmark Can exploit write concurrency (k_w) by batching writes
- ✓ Amortized cost per write following PIO is $\frac{\alpha}{k_w}$





- \checkmark Can exploit write concurrency (k_w) by batching writes
- ✓ Amortized cost per write following PIO is $\frac{\alpha}{k_w}$
- ✓ Example: DBMS bufferpool



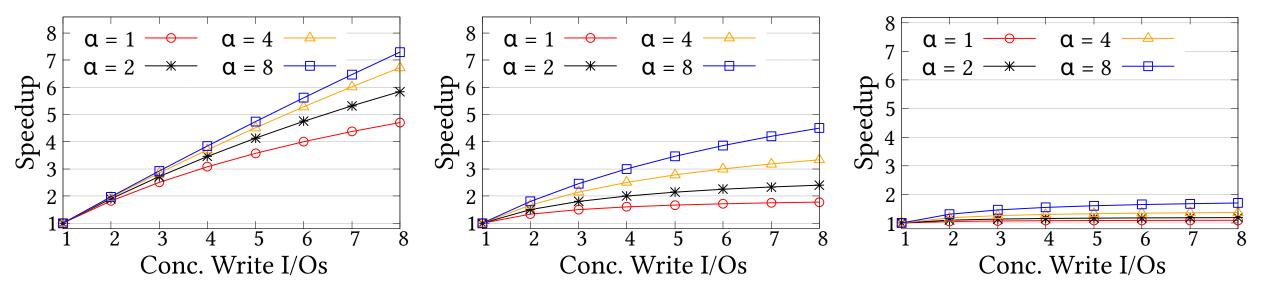




10% reads







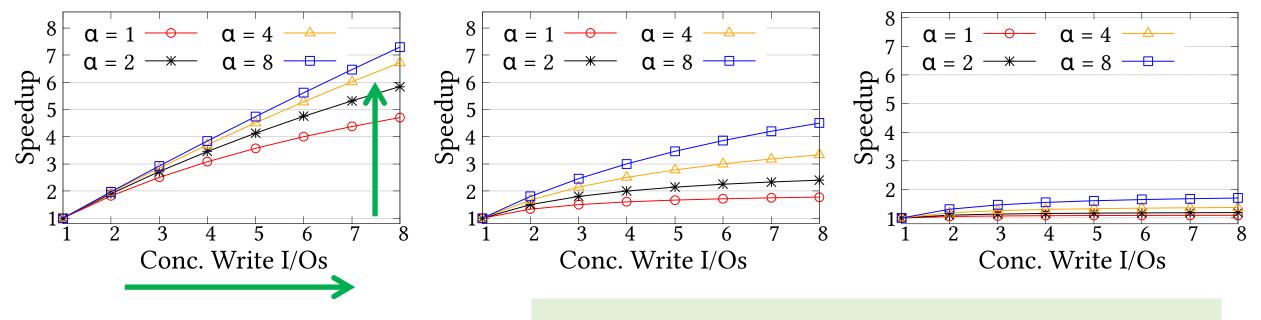




10% reads







Speedup increases with increasing concurrent I/Os

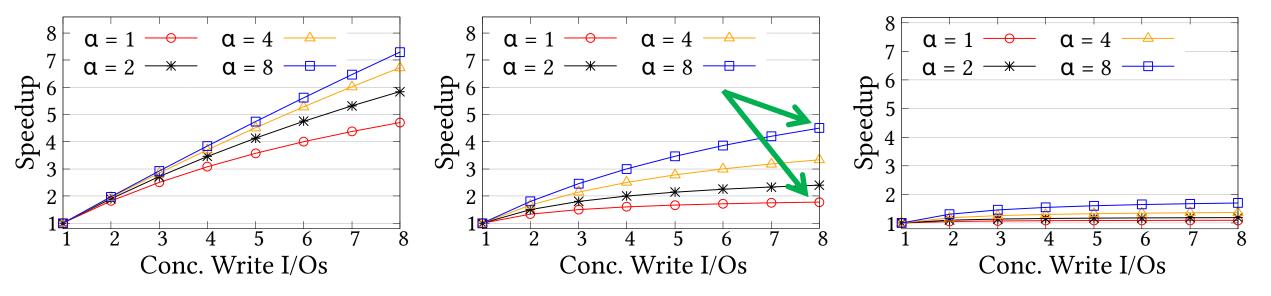




10% reads

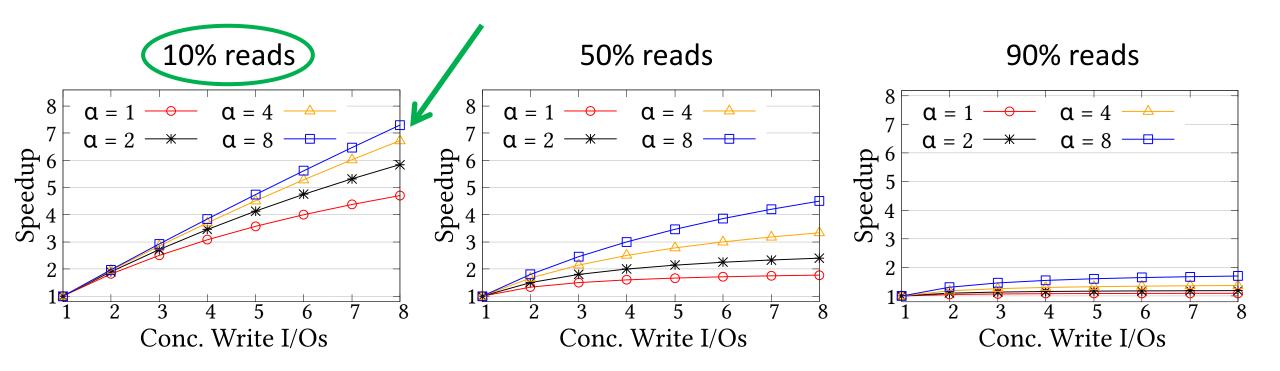






Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry ($gain \propto \alpha$)





Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry ($gain \propto \alpha$) Speedup is highest for write-intensive workloads





Batchable Reads

 \checkmark Can exploit read concurrency (k_r) by batching reads





Batchable Reads

- \checkmark Can exploit read concurrency (k_r) by batching reads
- ✓ Amortized cost per read following PIO is $\frac{1}{k_r}$





Batchable Reads

- \checkmark Can exploit read concurrency (k_r) by batching reads
- ✓ Amortized cost per read following PIO is $\frac{1}{k_r}$
- Example: Graph/tree traversal



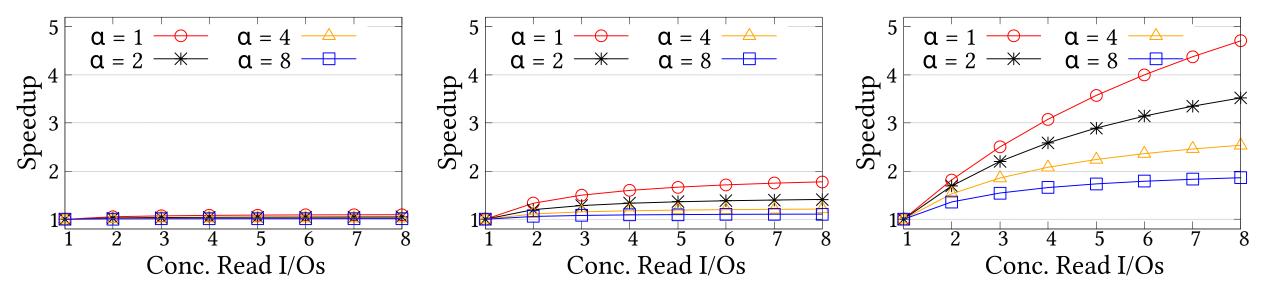


Batchable Reads











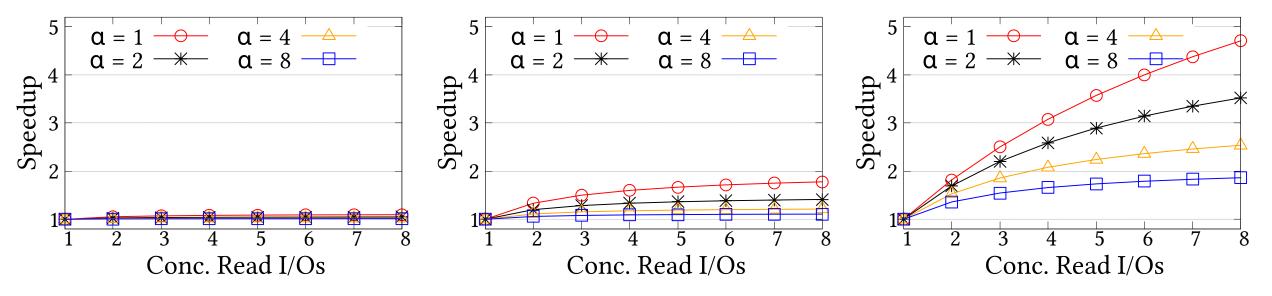


Batchable Reads









Speedup increases with increasing concurrent I/Os



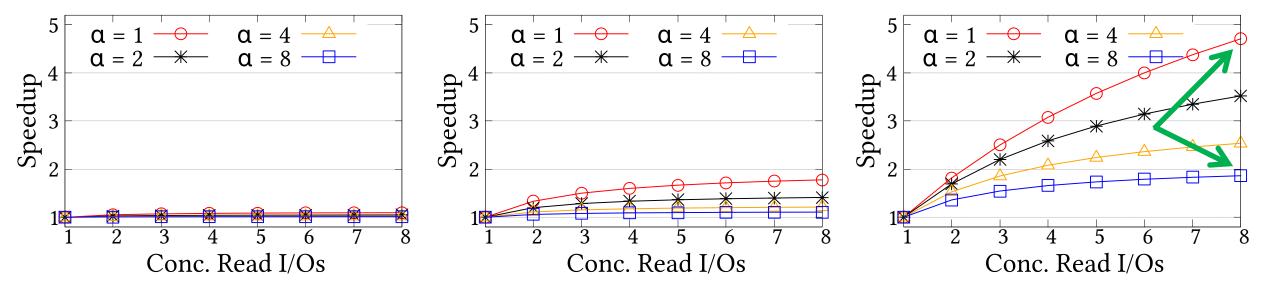


Batchable Reads









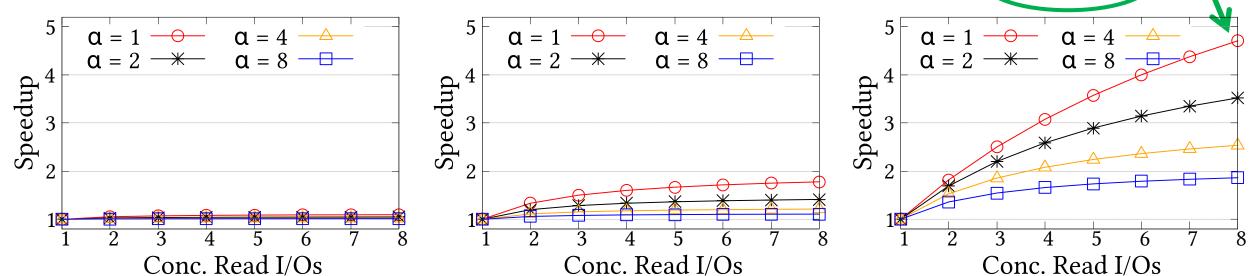
Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry ($gain \propto 1/\alpha$)

90% reads

Batchable Reads

10% reads





Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry ($gain \propto 1/_{\alpha}$) Speedup is highest for read-intensive workloads





✓ Can exploit both read and write concurrency (k_r , k_w)

ab SidD Sid



✓ Can exploit both read and write concurrency (k_r, k_w)

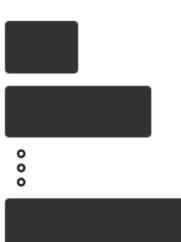
✓ Amortized cost per read following PIO is
$$\frac{1}{k_r}$$

ab Iab **OSiO**

✓ Amortized cost per write following PIO is
$$\frac{\alpha}{k_w}$$



- ✓ Can exploit both read and write concurrency (k_r, k_w)
- ✓ Amortized cost per read following PIO is $\frac{1}{k_r}$
- ✓ Amortized cost per write following PIO is $\frac{\alpha}{k_w}$
- ✓ Example: LSM compaction



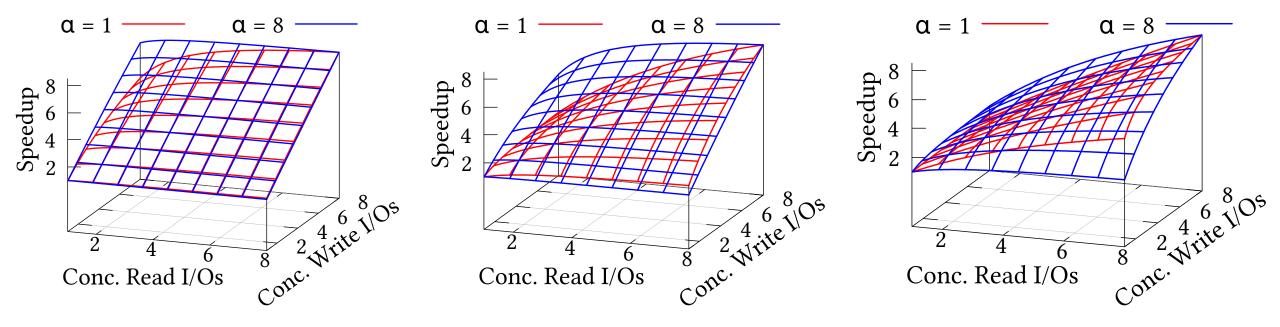


10% reads

ab SidD Sid

50% reads





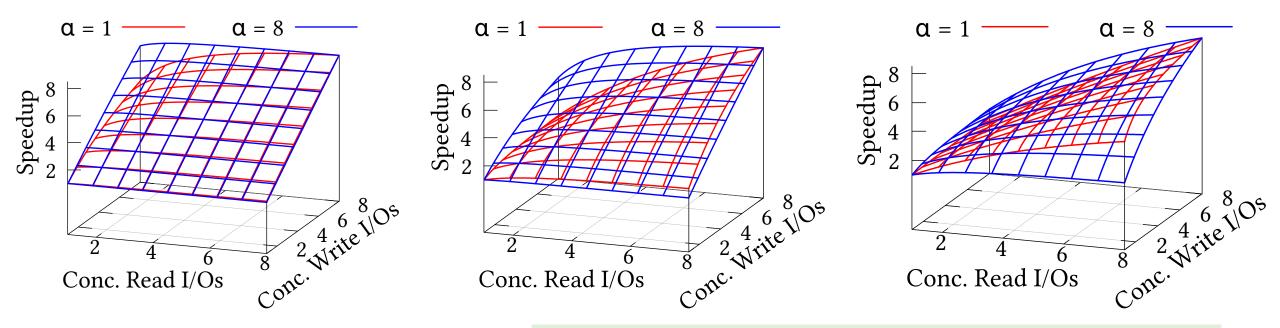


10% reads

B8 <u>9</u> DiSC

50% reads





Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry

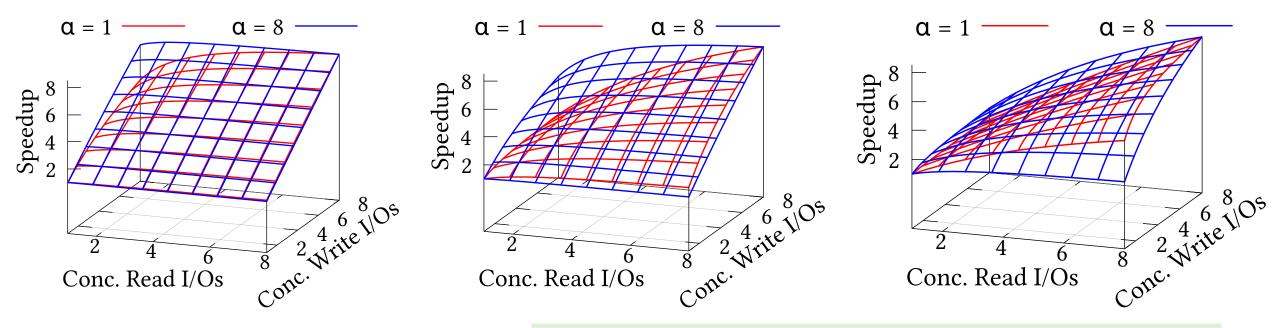


10% reads

B8 <u>9</u> DiSC

50% reads





Speedup increases with increasing concurrent I/Os Speedup depends on asymmetry

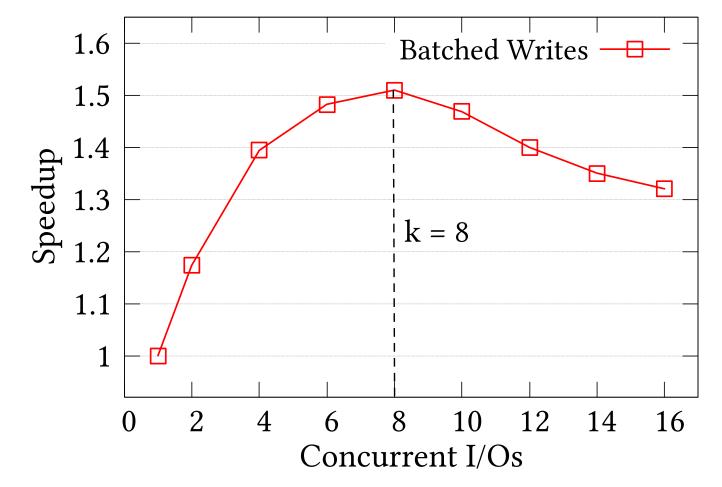
Impact of utilizing write concurrency is higher



Importance of using Proper k

Sample application: batchable writes

ab SidD Sid



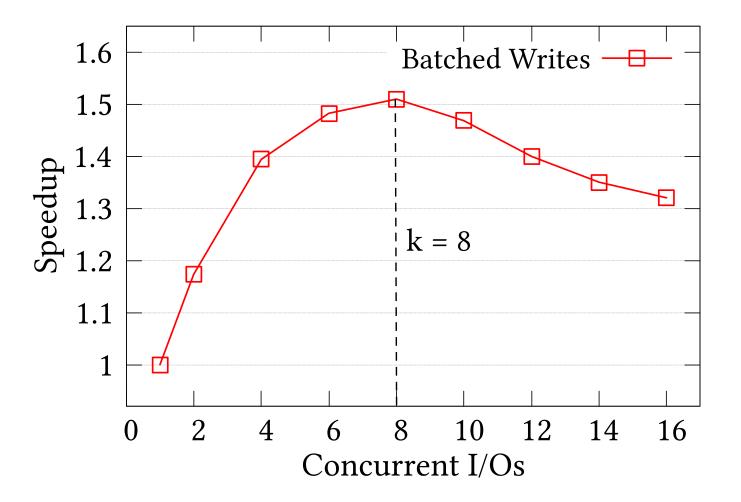


Importance of using Proper k

Sample application: batchable writes

lab OSIC

We use PCIe SSD ($k_w = 8$) to run this concurrency-aware application





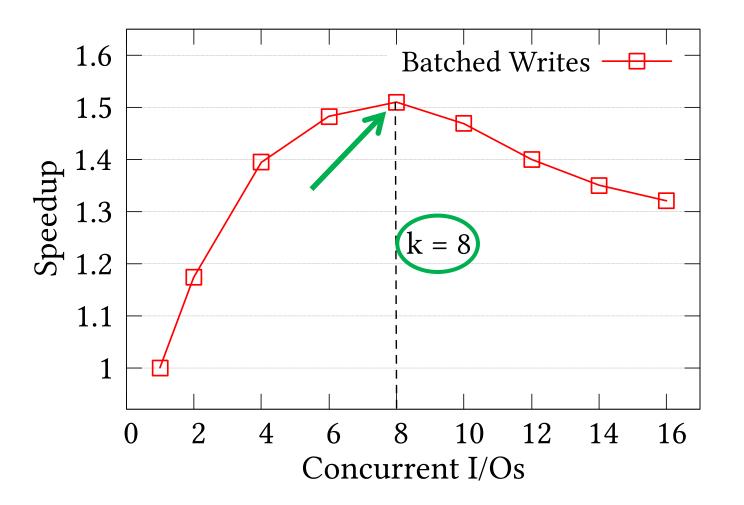
Importance of using Proper k

Sample application: batchable writes

<u>क</u> <u>वि</u> DisC

We use PCIe SSD ($k_w = 8$) to run this concurrency-aware application

Optimal speedup at the device concurrency.









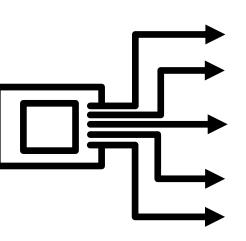
ि <u>वि</u> Sid

Know Thy Device





<u>क</u> <u>वि</u> DisC



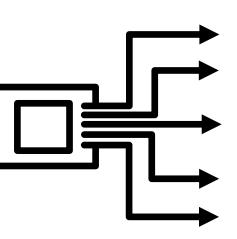
Know Thy Device

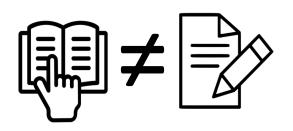
Exploit concurrency (with care)





ि S DisC





Know Thy Device

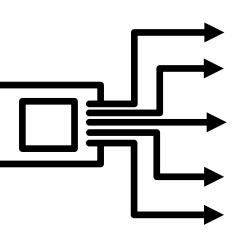
Exploit concurrency (with care)

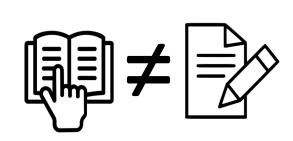
Treat read and write differently.

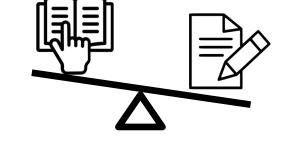




BS dg DiSC







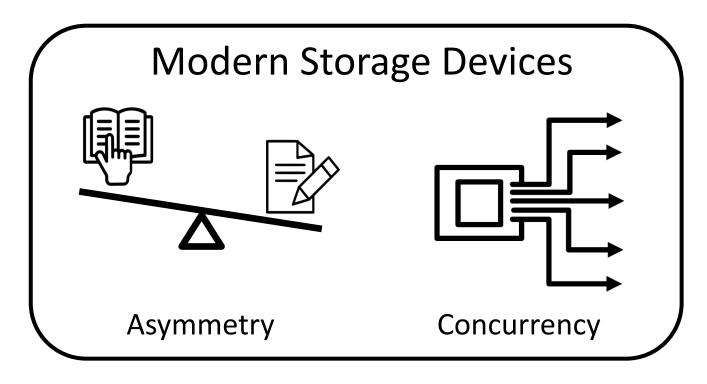
Know Thy Device

Exploit concurrency (with care)

Treat read and write differently.

 α controls performance

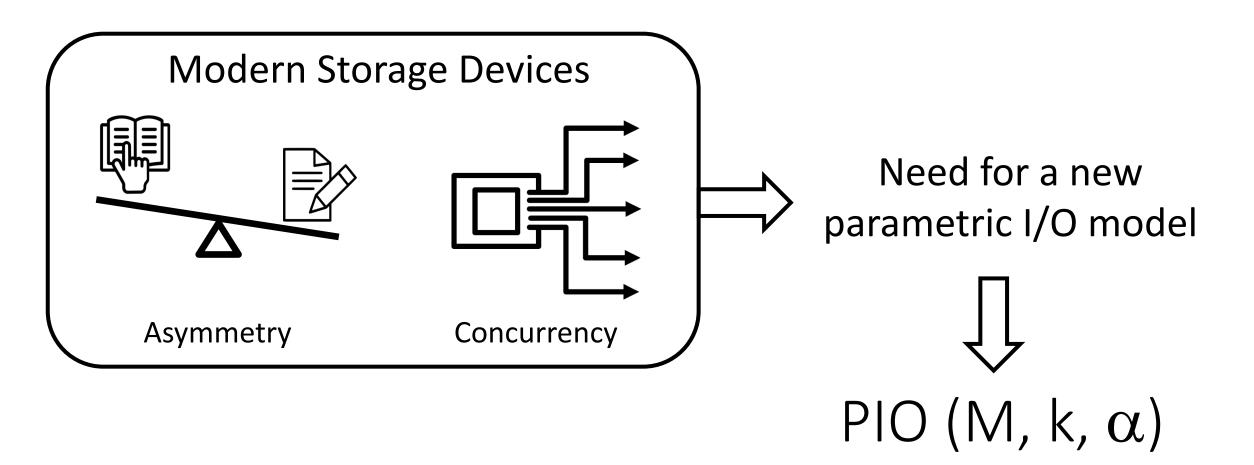




ि <u>वि</u> Sid



ि S DisC





DiaD Sid

Make *asymmetry* and *concurrency* part of *algorithm design*

... not simply an engineering optimization



B8 <u>9</u> DiSC

Make *asymmetry* and *concurrency* part of *algorithm design*

... not simply an engineering optimization

Build algorithms/data structures for storage devices with asymmetry α and concurrency k



Make *asymmetry* and *concurrency* part of *algorithm design*

... not simply an engineering optimization

Build algorithms/data structures for storage devices with asymmetry α and concurrency k

index structures

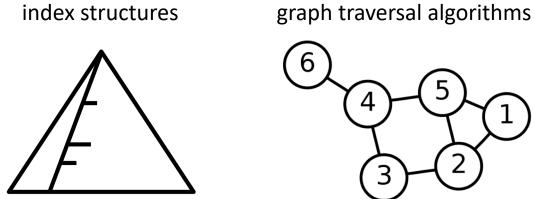
BS da DiSC



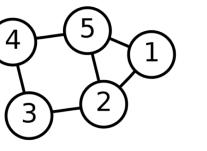
Make *asymmetry and concurrency* part of *algorithm design*

... not simply an engineering optimization

Build algorithms/data structures for storage devices with asymmetry α and concurrency k



BS da DisC

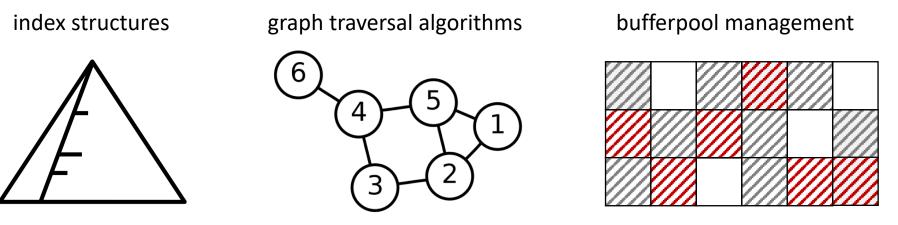




Make *asymmetry* and *concurrency* part of *algorithm design*

... not simply an engineering optimization

Build algorithms/data structures for storage devices with asymmetry α and concurrency k





Make *asymmetry and concurrency* part of *algorithm design*

... not simply an engineering optimization

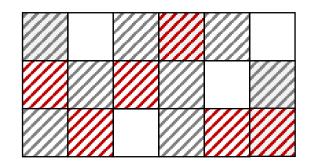
Build algorithms/data structures for storage devices stay runed! with asymmetry α and concurrency k

index structures

BS da DiSC

graph traversal algorithms

bufferpool management





Thank You!

lab **SS** da DSIC

dísc.bu.edu/pío