

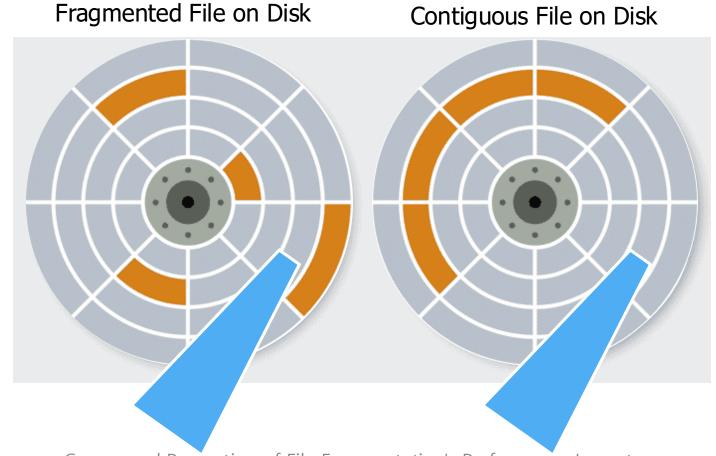
# We Ain't Afraid of No File Fragmentation: Causes and Prevention of Its Performance Impact on Modern Flash SSDs

**Yuhun Jun**<sup>1,2</sup>, Shinhyun Park<sup>1</sup>, Jeong-Uk Kang<sup>2</sup>, Sang-Hoon Kim<sup>3</sup> and Euiseong Seo<sup>1</sup>

<sup>1</sup>Sungkyunkwan University <sup>2</sup>Samsung Electronics <sup>3</sup>Ajou University

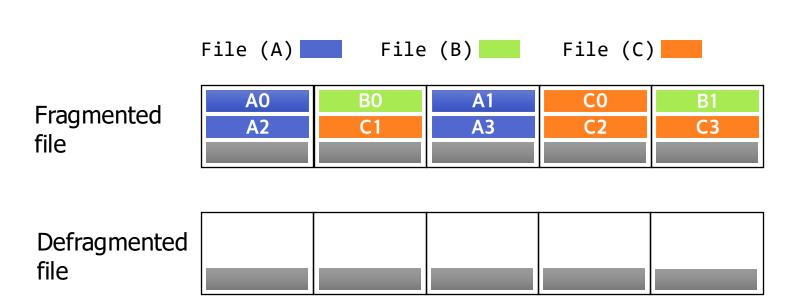
# **File Fragmentation**

Non-contiguously stored file → Degraded read performance

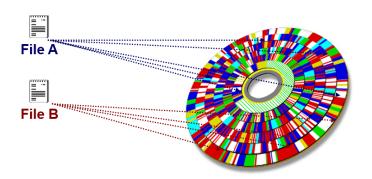


# **File Fragmentation**

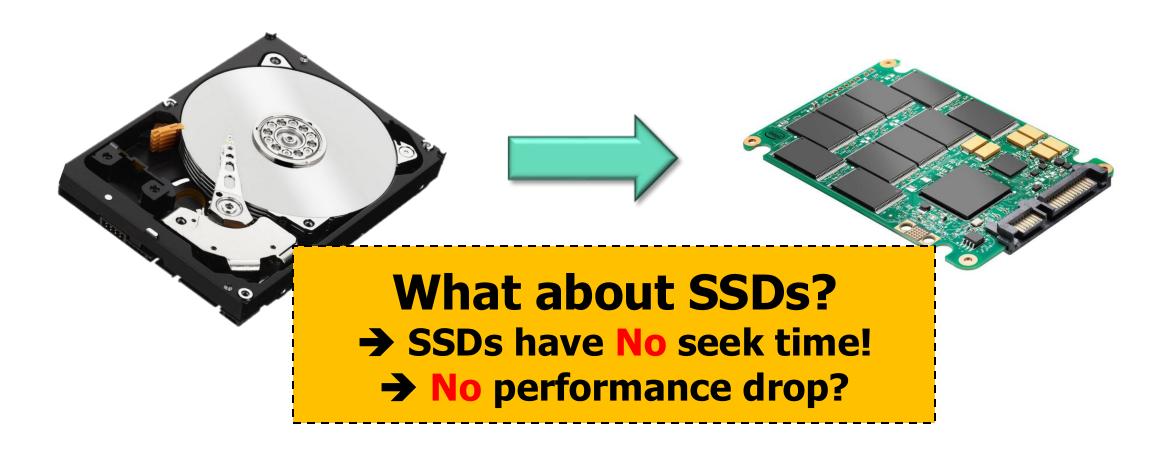
To recover, costly defragmentation should be performed





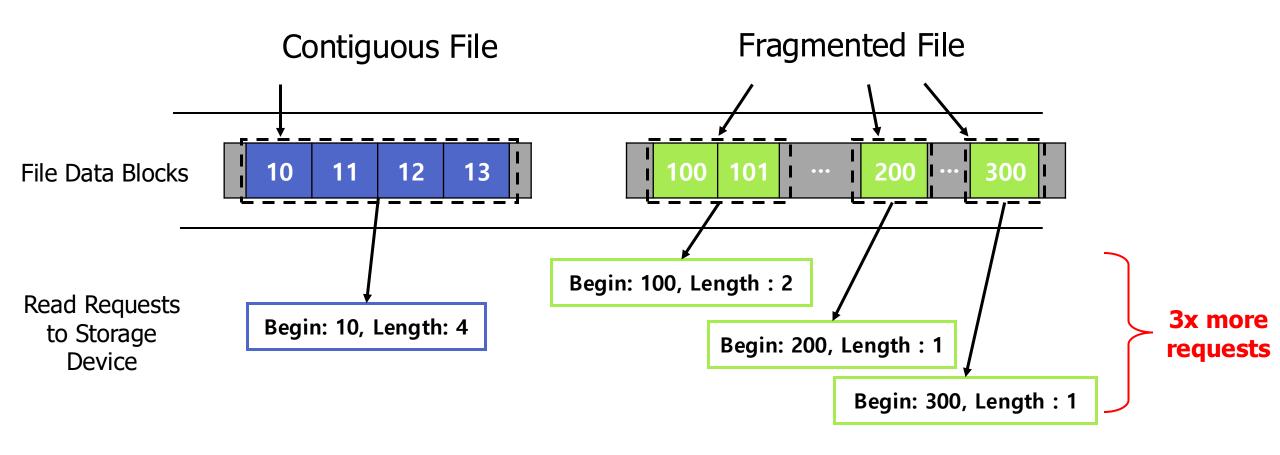


# **File Fragmentation**

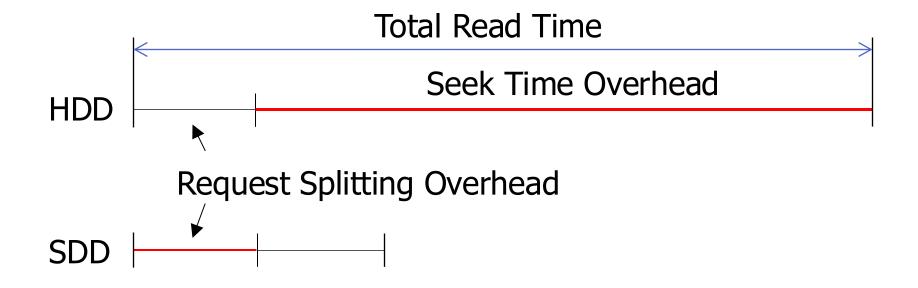


- Even in SSDs, still performance degradation occurs (observed reduction of 2x to 5x) \*
  - \* Conway et al. File systems fated for senescence? nonsense, says science!(FAST '17).
  - \* Kadekodi *et al*. Geriatrix: Aging what you see and what you don't see. A file system aging approach for modern storage systems (ATC '18).
  - \* Conway et al. Filesystem Aging: It's more usage than fullness (Hotstorage '19).
- request splitting caused by fragmentation increases kernel I/O stack overhead \*\*
  - \*\* Park and Eom. Fragpicker: A new defragmentation tool for modern storage devices (SOSP '21)

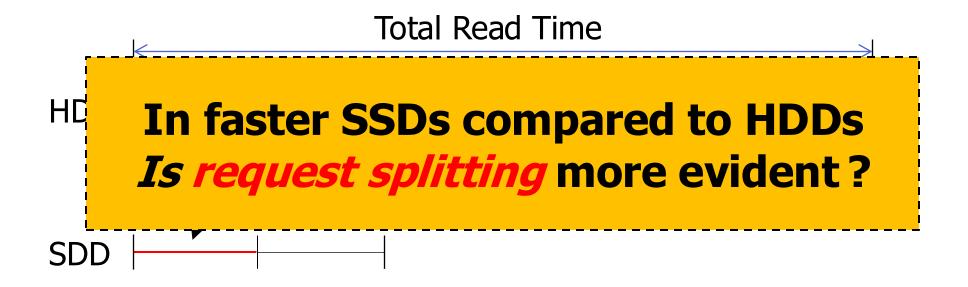
Request splitting



Request splitting

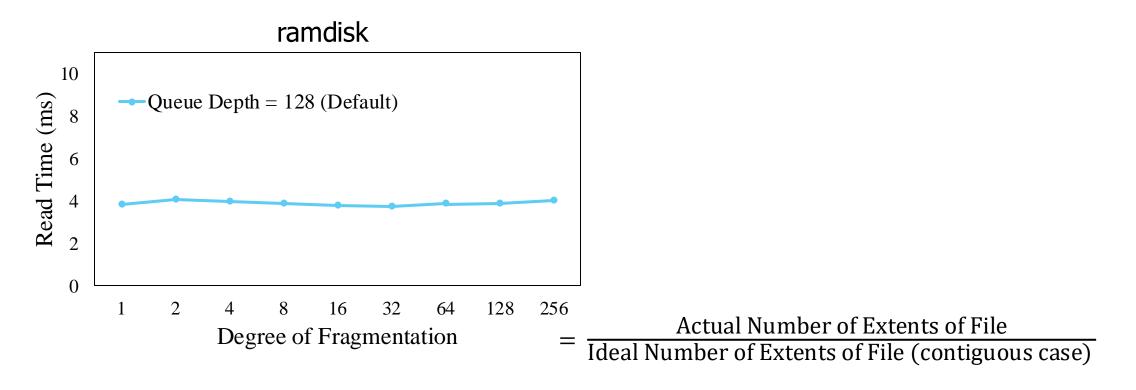


Request splitting

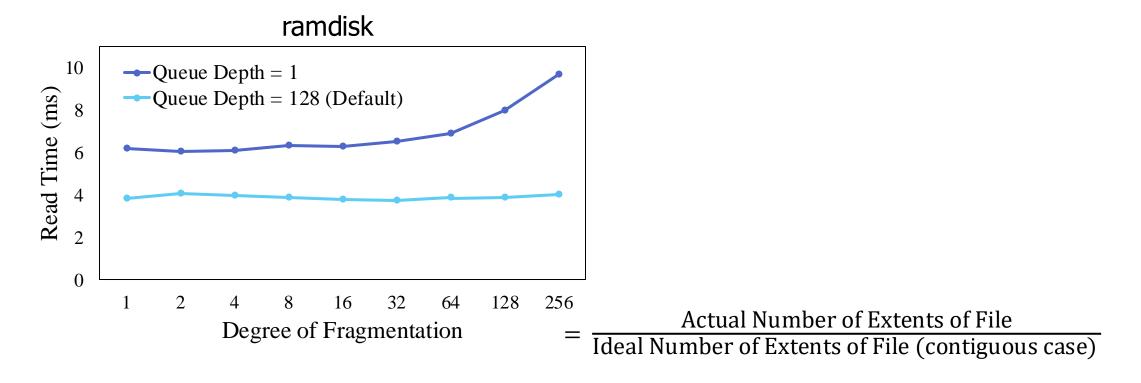


Does request splitting impact ramdisks more than SSDs?

Does request splitting impact ramdisks more than SSDs?

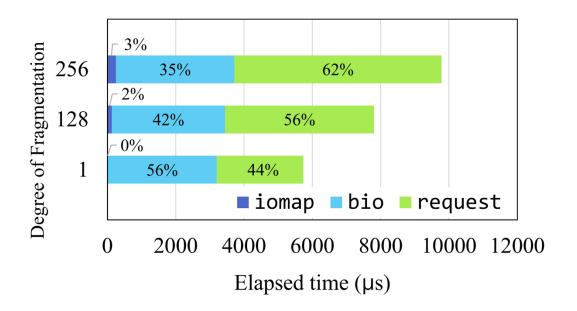


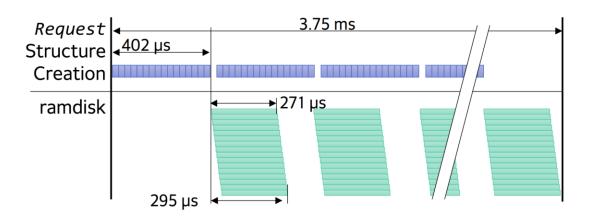
- Does request splitting impact ramdisks more than SSDs?
- Impact seen with forced queue depth of 1



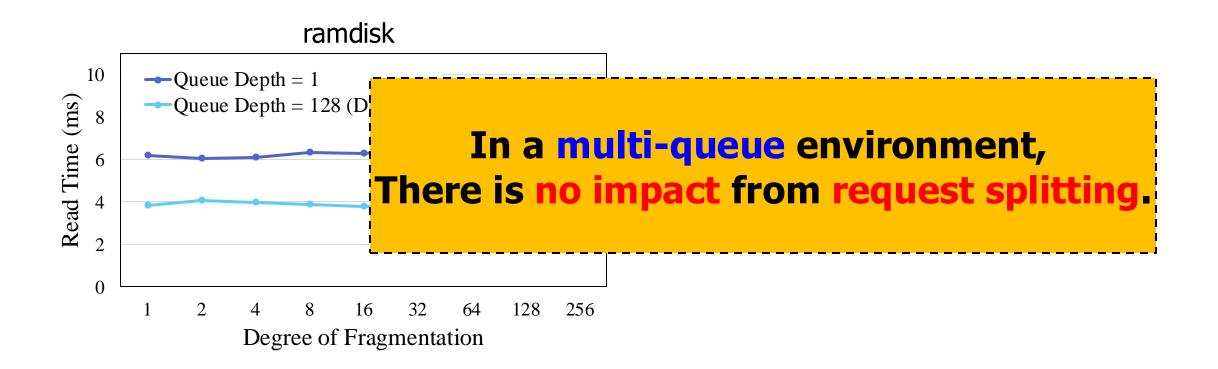


Request splitting overhead

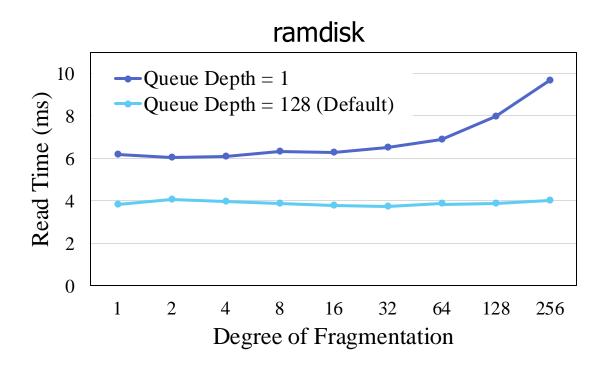


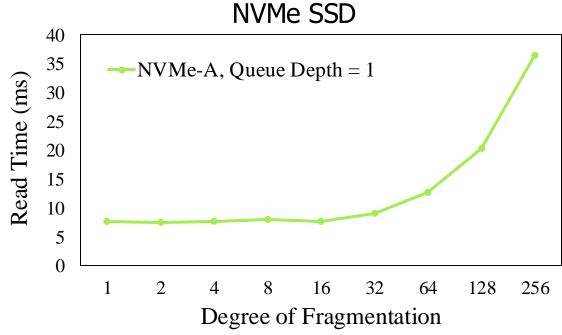


- Does request splitting impact ramdisks more than SSDs?
- → Impact seen with forced queue depth of 1

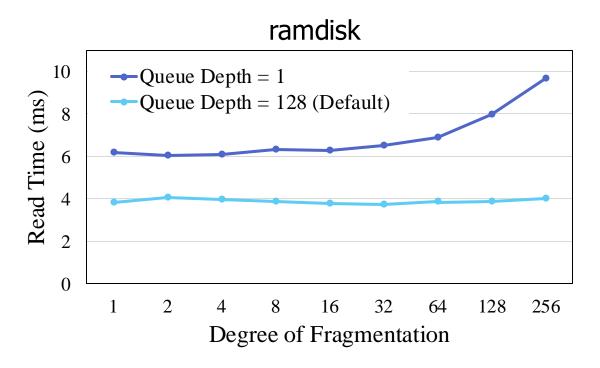


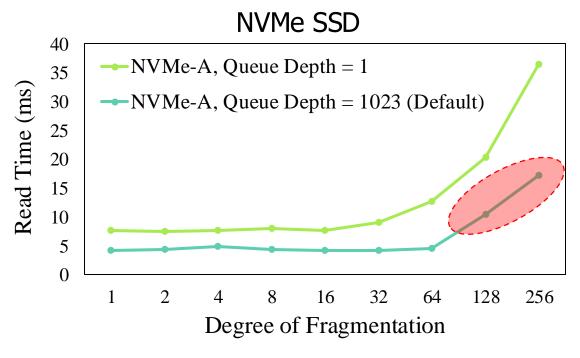
- Does request splitting impact ramdisks more than SSDs?
- → Impact seen with forced queue depth of 1; No request splitting impact in multi-queue (default)



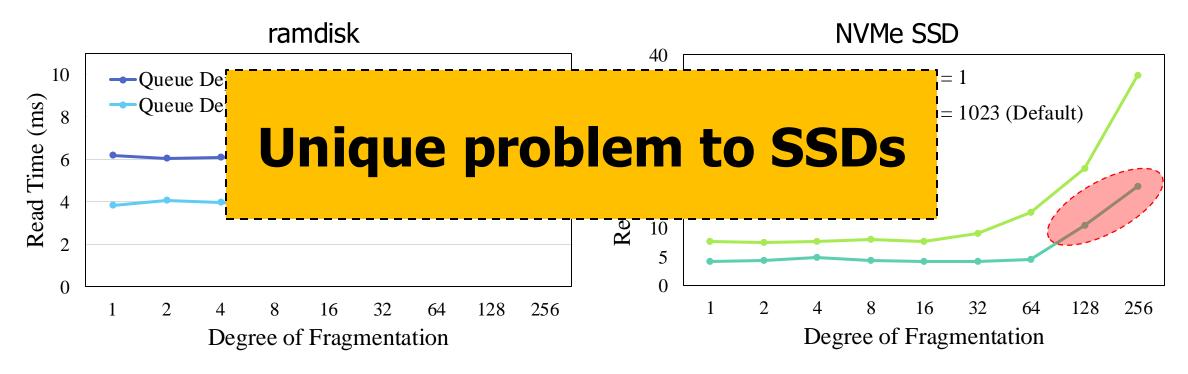


- Does request splitting impact ramdisks more than SSDs?
- → Impact seen with forced queue depth of 1; No request splitting impact in multi-queue (default) Commercial SSDs showed performance drop in fragmentation



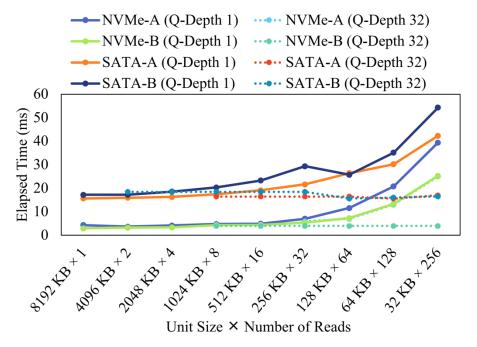


- Does request splitting impact ramdisks more than SSDs?
- → Impact seen with forced queue depth of 1; No request splitting impact in multi-queue (default) Commercial SSDs showed performance drop in fragmentation



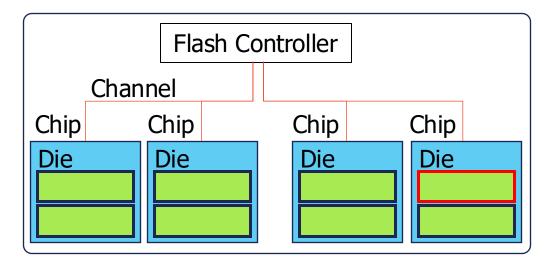
# **Analysis of Interface Overhead**

- Does PCIe Interface extend the read latency?
- → When we read the same number of sectors with different unit request sizes, the performance barely varied



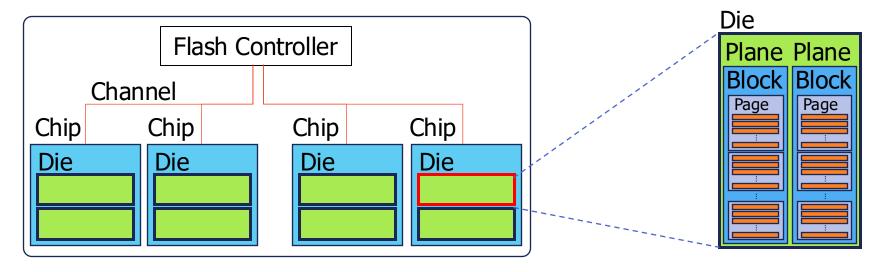
Time for reading 8 MB of data through raw device I/O operations

High performance from operating multiple NAND Flashes simultaneously



Structure of NAND Flash inside SSD (2-Channel 2-chip 2-Die SSD Total 8-Dies in an SSD)

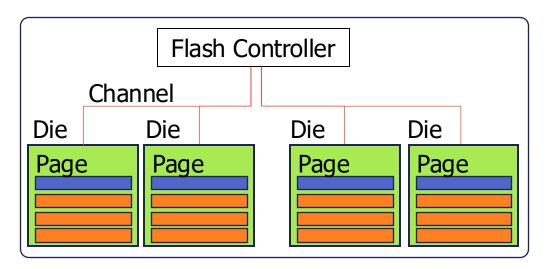
High performance from operating multiple NAND Flashes simultaneously



Structure of NAND Flash inside SSD (2-Channel 2-chip 2-Die SSD Total 8-Dies in an SSD)

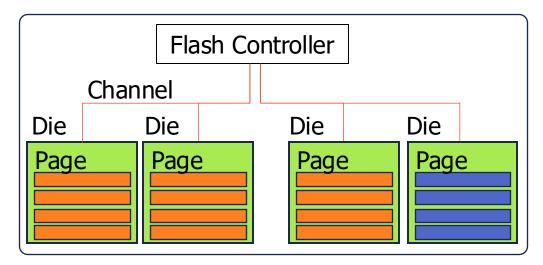
Page writing/reading suspends other operations issued to the same die

- High performance from operating multiple dies simultaneously
- For write and read operations, as many dies as possible should be utilized



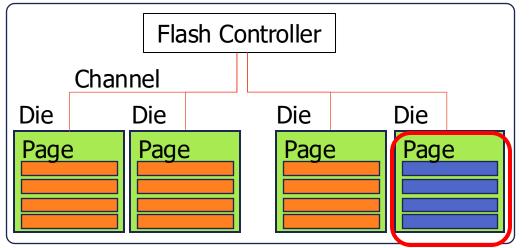
Structure of NAND **dies** inside SSD (2-Channel 2-Die SSD, Total 4-Dies in an SSD)

- High performance from operating multiple dies simultaneously.
- For write and read operations, as many dies as possible should be activated
- Focusing on a single die during reading → Reduced parallelism



Structure of NAND **dies** inside SSD (2-Channel 2-Die SSD, Total 4-Dies in an SSD)

- High performance from operating multiple dies simultaneously.
- For write and read operations, as many dies as possible should be activated
- Focusing on a single die during reading → Reduced parallelism

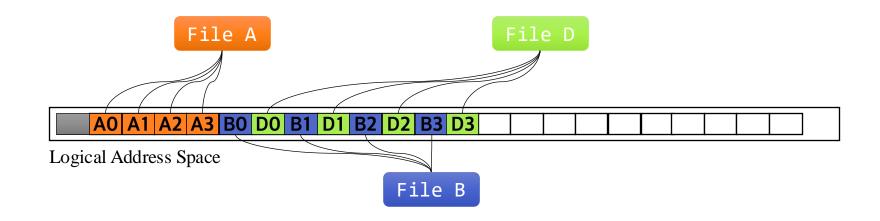


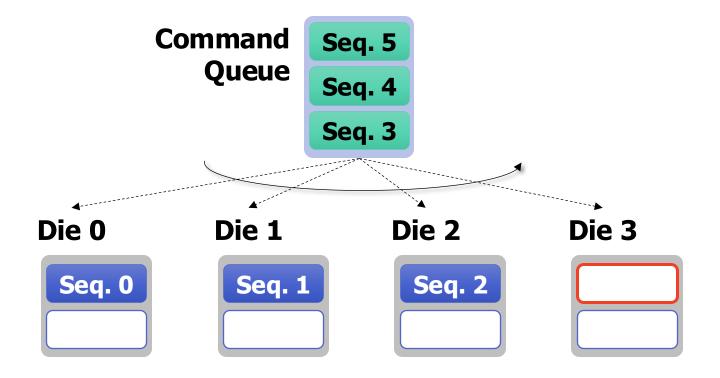
Read Collisions

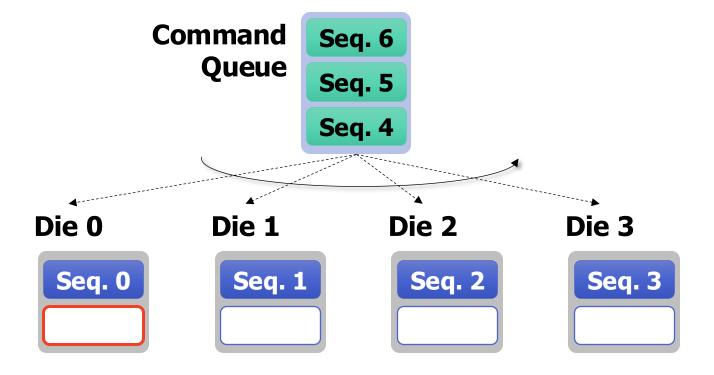
Structure of NAND **dies** inside SSD (2-Channel 2-Die SSD, Total 4-Dies in an SSD)

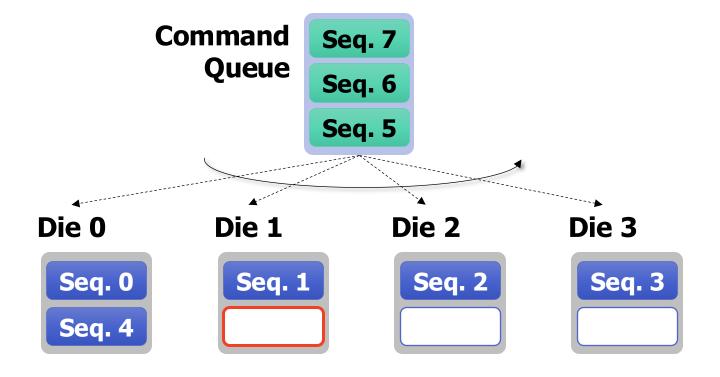
### **File Fragmentation Scenarios**

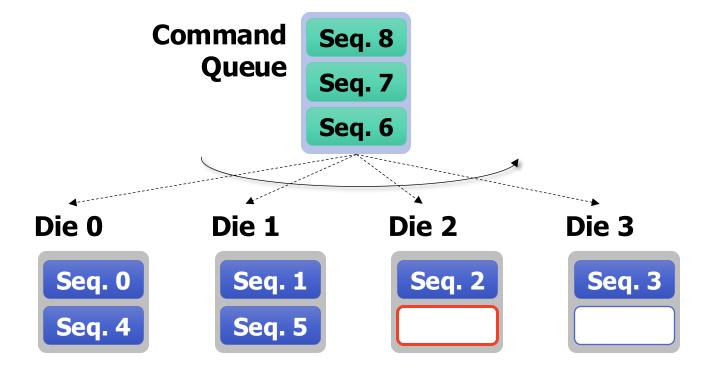
 File fragmentation occurs when multiple files are appended in an alternating manner

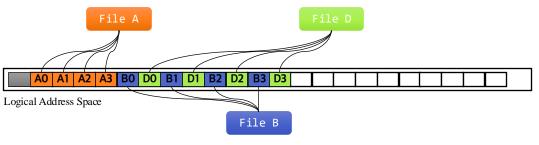


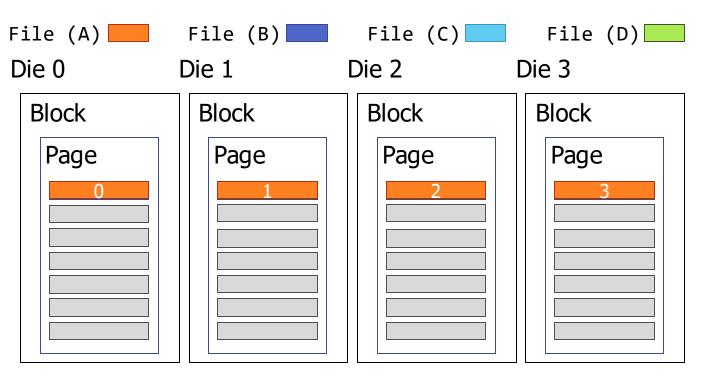


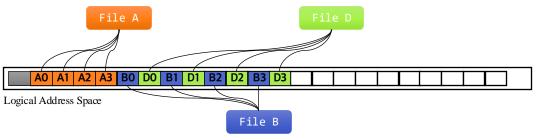


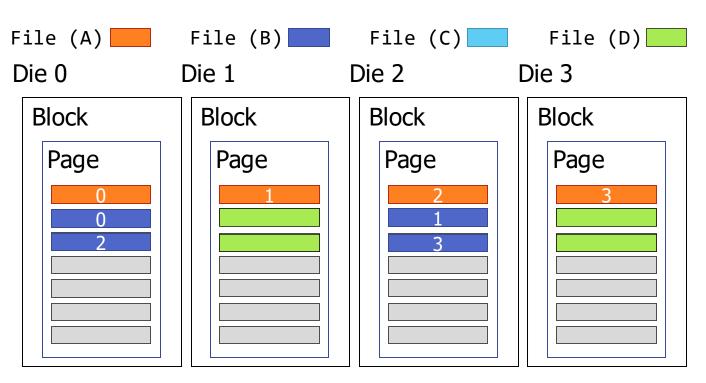


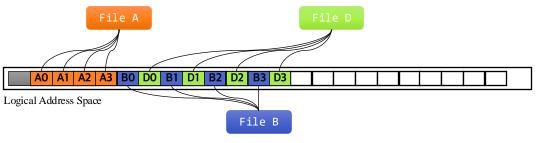


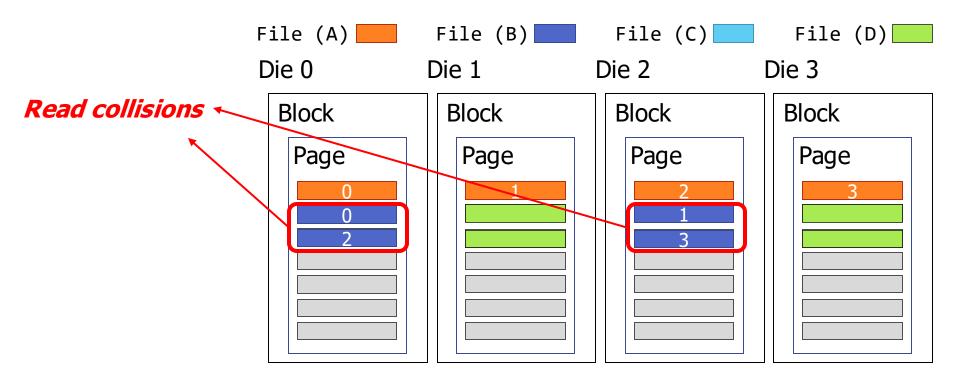


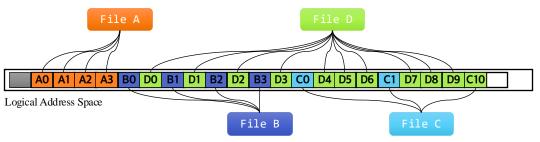


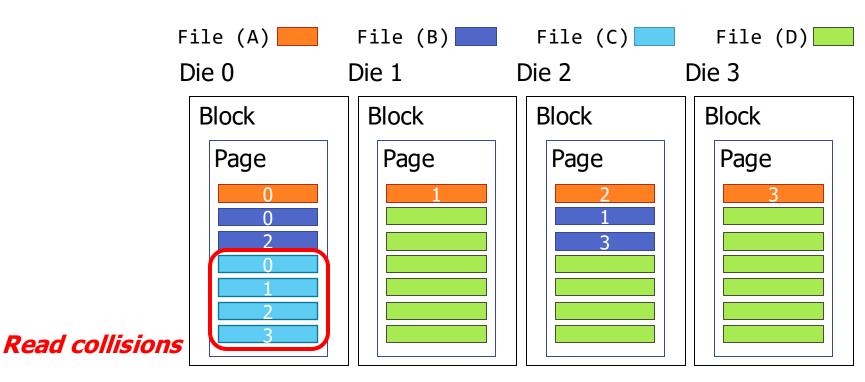




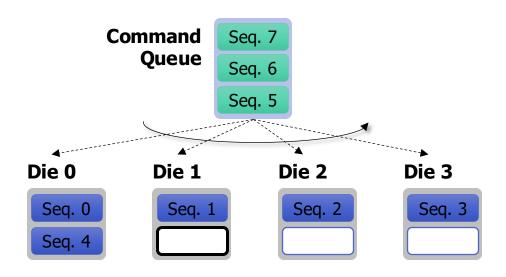




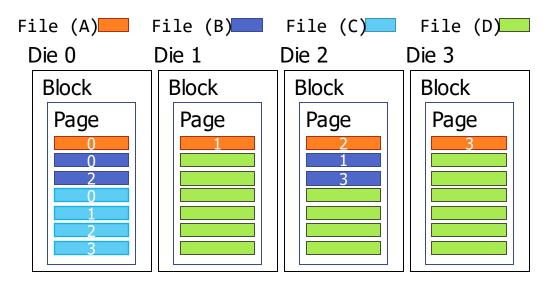




- Alternating appends causes misaligned die allocation
  - → Read collisions

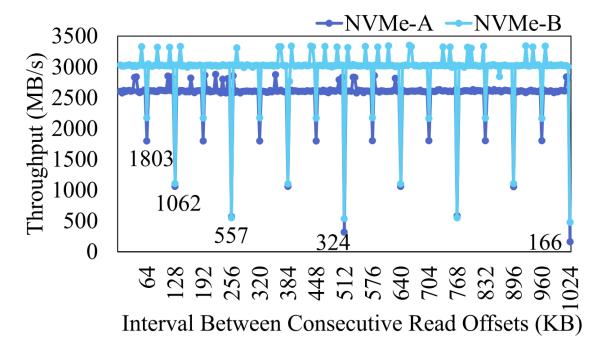


Die allocation in a round-robin manner



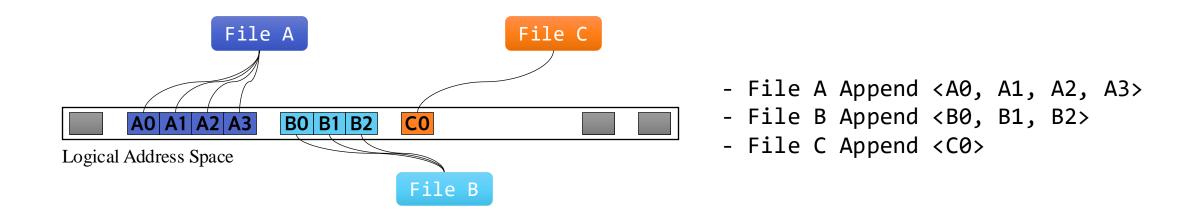
Concurrent writes to multiple files cause misaligned die allocation

Experiments with commercial SSDs clearly verified our conjecture



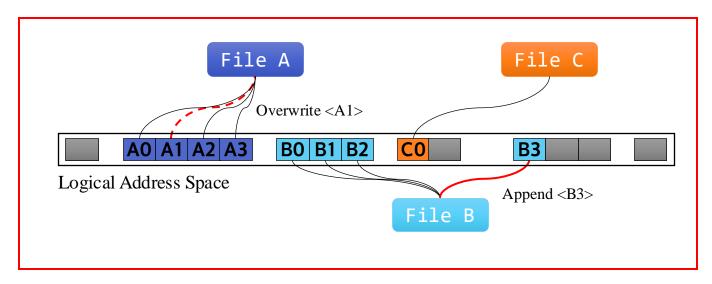
Throughput while varying the interval between starting points of consecutive read operations

### **Misaligned Die Allocation from Overwrites**

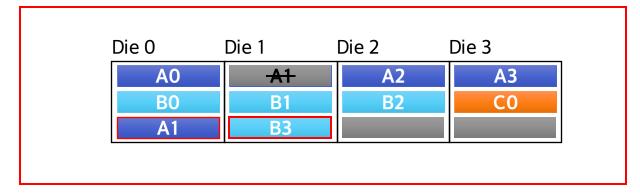




### **Misaligned Die Allocation from Overwrites**



- File A Append <A0, A1, A2, A3>
   File B Append <B0, B1, B2>
- File C Append <C0>
- File A Overwrite <A1>
- File B Append <B3>



### **Conventional Approach**

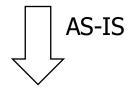
- Appends and overwrites lead to misaligned die allocation
- While defragmentation addresses this issue, it incurs significant costs

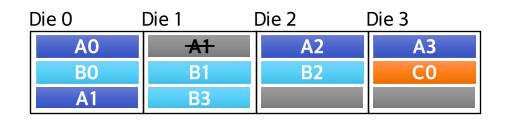


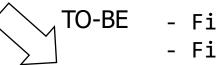
Reads and rewrites the entire file

Prevents misaligned allocation on the fly









- File A Overwrite <A1>

- File B Append <B3>

Die 0	Die 1	Die 2	Die 3
AO	<del>-A1</del>	A2	A3
ВО	B1	B2	CO
	A1		B3
	•	1	

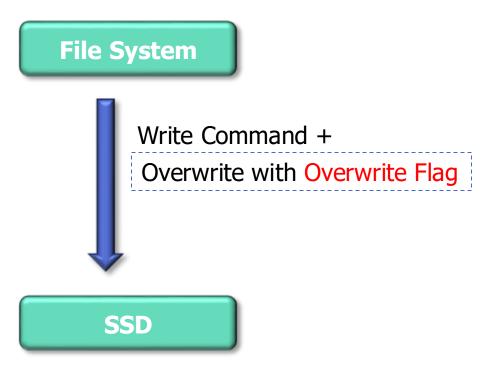
- Prevents misaligned allocation on the fly
  - File A Overwrite <A1>
  - File B Append <B3>

Die 0	Die 1	D	e 2	Die 3
A0	<del>A1</del>		A2	A3
В0	B1		B2	CO
	A1			B3
	•			

- File system provides file information to the SSD
- Overwrite to same die

- File A Overwrite <A1> (with OW flag)

Die 0	Die 1	Die 2	Die 3
A0	<del>A1</del>	A2	A3
В0	B1	B2	CO
	A1		B3



- File system provides file information to the SSD
- Overwrite to same die
- Append to the die next to last written one
  - File A Overwrite <A1> (with OW flag)
  - File B Append <B3> (with AP flag, B2)

Die 0	Die 1	Die 2	Die 3
AO	<del>A1</del>	A2	A3
B0	B1	B2	CO
	A1		B3

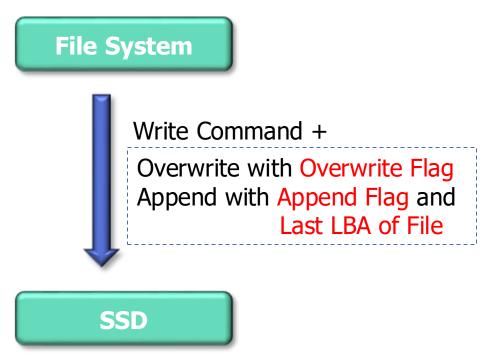


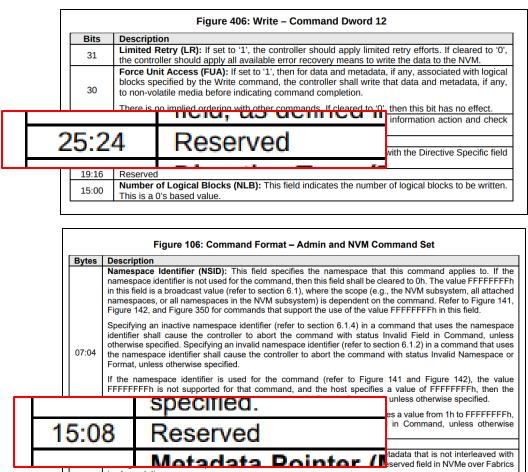


Overwrite with Overwrite Flag
Append with Append Flag and
Last Logical Block Address
of the file

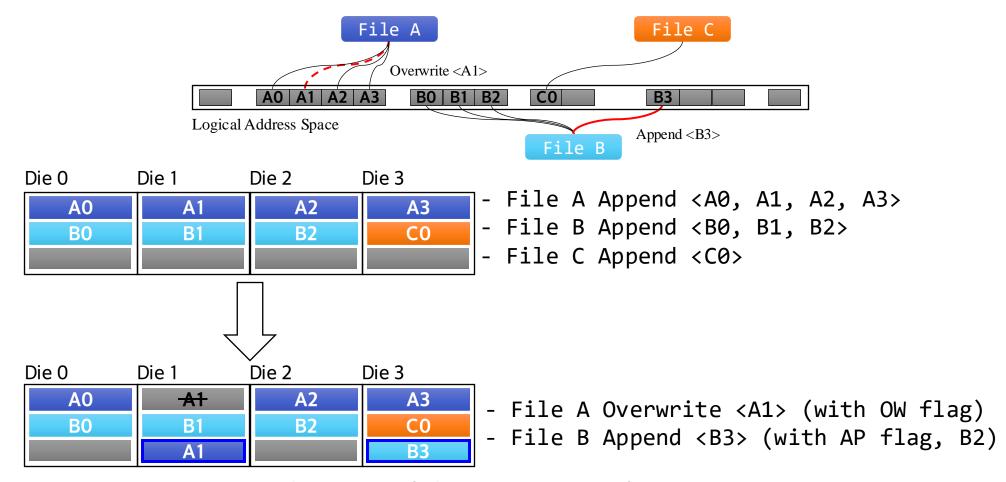


The hints are sent through an unused field of the NVMe write command





Prevents misaligned allocation on the fly





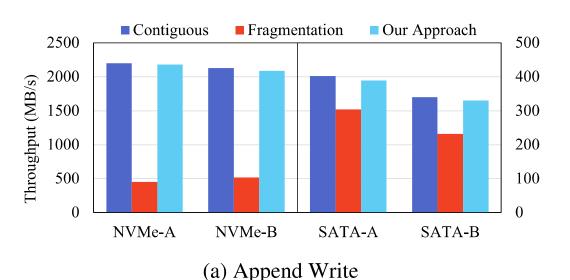
Environment

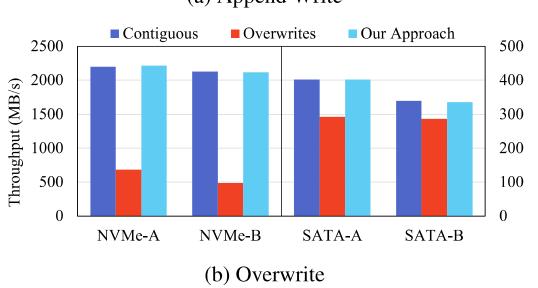
System Configuration		
Processor	Intel Xeon Gold 6138 2.0 GHz, 160-Core	
Chipset	Intel C621	
Memory	DDR4 2666 MHz, 32 GB x16	
OS	Ubuntu 20.04 Server (kernel v5.15.0)	
File system	Ext4	

NVMeVirt Emulator [22]				
Interface	PCIe Gen 3 x4			
Capacity	60 GB			
Channel Count	4			
Dies per Channel	2			
Read/Write Unit Size	32 KB			
Read Time	36 µs			
Write Time	185 µs			

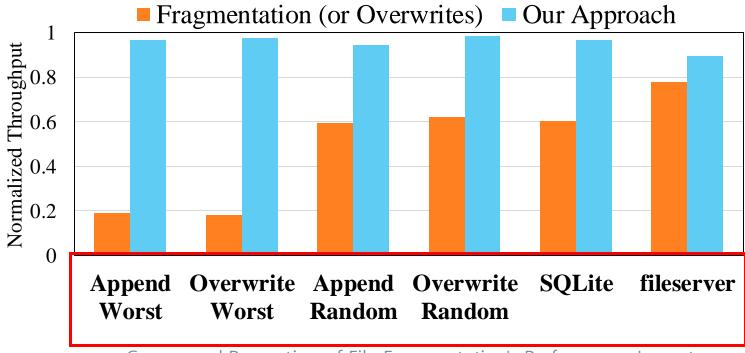
- Our approach was validated using commodity SSDs
- Evaluation our approach with SSD emulator
  - Modified Ext4 and NVMe driver to transmit info through NVMe Write Command
  - NVMeVirt adjusts die allocation policy using this information

- Reading 8MB file
  - Written by 32KB x 256 times
- Commodity SSDs
  - Samsung 980 Pro 1TB NVMe-A
  - WD Black SN850 1TB NVMe-B
  - Samsung 870 Evo 500GB SATA-A
  - WD Blue SA510 500GB— SATA-B

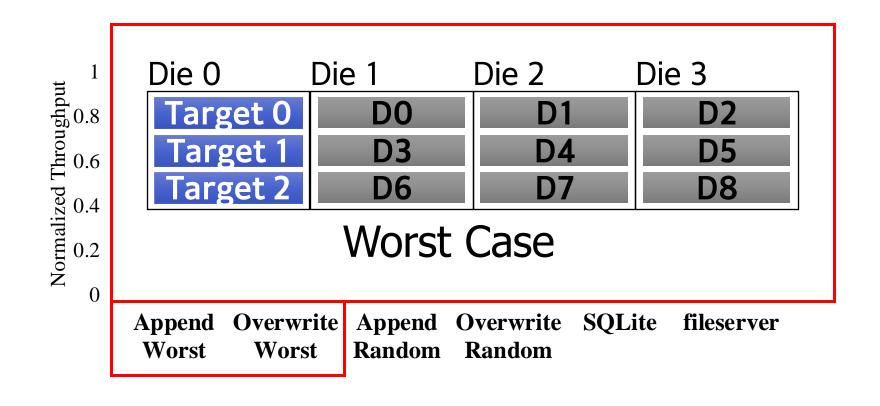




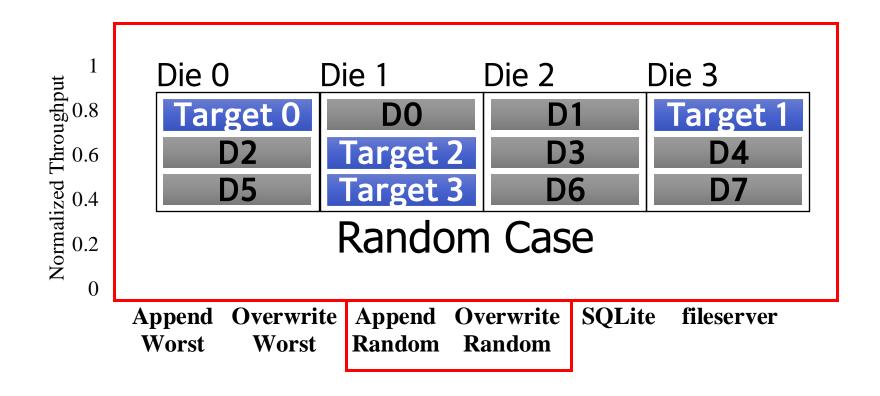
- Hypothetical workloads, SQLite and Filebench
  - SQLite: 16KB records x 10,000 while writing 100KB chunks to dummy files 5,005 DoF
  - Fileserver: 10 threads perform 32 KB rand. appends for 1 Min, and seq. reads over 128KB x 10,000 files



Worst case of hypothetical workloads

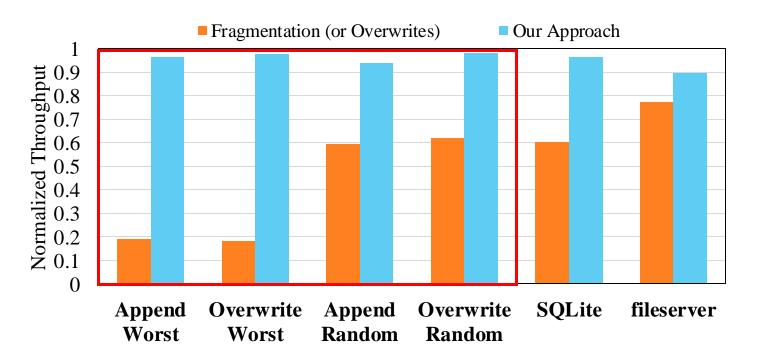


Random case of hypothetical workloads



- In the worst case, 20% of contiguous file's
- In the random case, 60% of contiguous file's
- In SQLite, 60% of contiguous file's
- In fileserver, 77% of contiguous file's

- → Improved to within 6%.
- → Improved to within 10%.



Die 0	Die 1	Die 2	Die 3
TO	D0	D1	D2
T1	D3	D4	D5
T2	D6	D7	D8

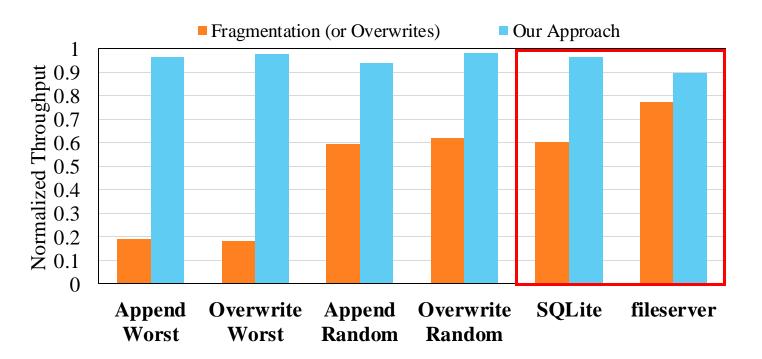
**Worst Case** 

Die 0	Die 1	Die 2	Die 3
TO	D0	D1	Ţ
D2	T2	D3	D4
D5	<b>T3</b>	D6	D7

Random Case

- In the worst case, 20% of contiguous file's
- In the random case, 60% of contiguous file's
- In SQLite, 60% of contiguous file's
- In fileserver, 77% of contiguous file's

- → Improved to within 6%.
- → Improved to within 10%.



Die 0	Die 1	Die 2	Die 3
TO	D0	D1	D2
T1	D3	D4	D5
T2	D6	D7	D8

**Worst Case** 

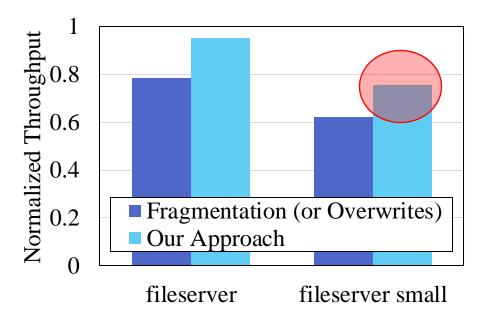
Die 0	Die 1	Die 2	Die 3
TO	D0	D1	T1
D2	T2	D3	D4
D5	<b>T3</b>	D6	D7

Random Case



## **Issues Occurring When Write Sizes are Small**

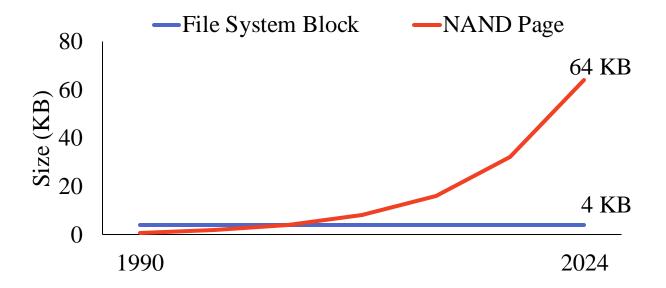
Discovered a performance issue when the write chunk size is small





## **Page Bloating Caused by Small Writes**

- Due to cost and performance considerations, the size of NAND pages is increasing
  - Originally functioning at 512 bytes<sup>1</sup>, they are now up to 64 KB<sup>2</sup> and continue to grow
- In contrast, the file system block size has not increased in over 20 years<sup>3</sup>

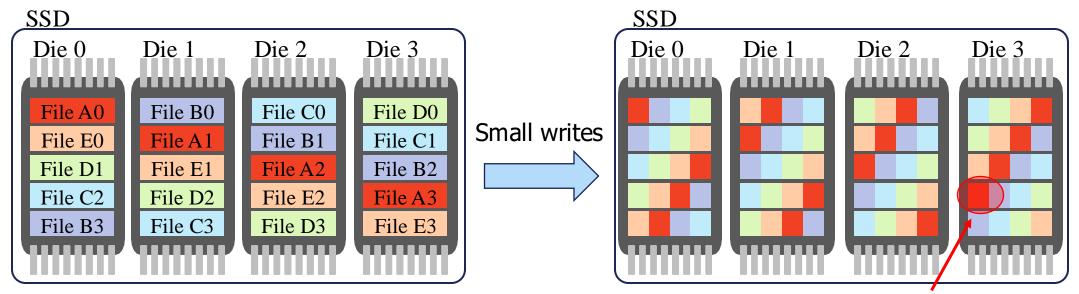


- 1. Suh et al. A 3.3 V 32 Mb NAND flash memory with incremental step pulse programming scheme. IEEE Journal of Solid-State Circuits, 1995
- 2. Kim *et al.* A 1Tb 3b/Cell 8th-generation 3D-NAND flash memory with 164MB/s write throughput and a 2.4Gb/s interface. IEEE International Solid-State Circuits Conference, 2022
- 3. Tweedie et al. Journaling the linux ext2fs filesystem. In Proceedings of The Fourth Annual Linux Expo. Durham, North Carolina, 1998.



### **Page Bloating Caused by Small Writes**

- Writing multiple files smaller than NAND page size simultaneously can cause bloating across several NAND pages
- Inefficiency results from size mismatch between SSD's page and file system's block



Die throughput decreased to 1/4

### **Conclusion**

- We identify the true cause of performance degradation due to file fragmentation
  - Request splitting overhead is concealed in a multi-queue environment
  - Primary cause is read collisions due to misaligned die allocation
- We proposed an approach to mitigate the misalignment
  - By providing filesystem information to the SSD,
     it maintains the proper die allocations even under adverse conditions
  - Addressing not only append write cases, but also overwrite cases

