



Who moved my cheese?

How storage systems deal with changes in their media

Gala Yadgar



TECHNION

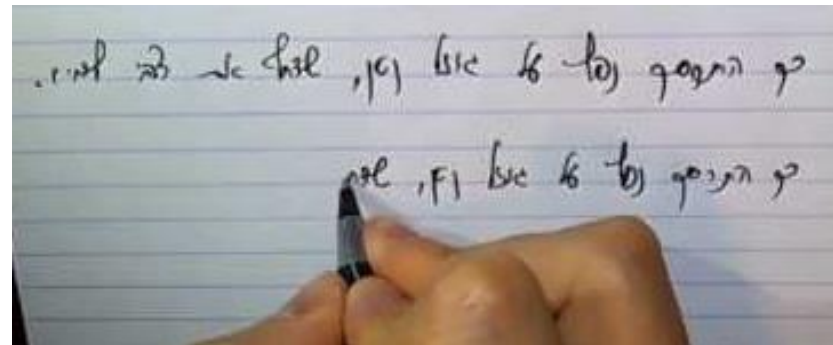


The Henry and Marilyn Taub
Faculty of Computer Science

“Why is Hebrew written backwards?”

→

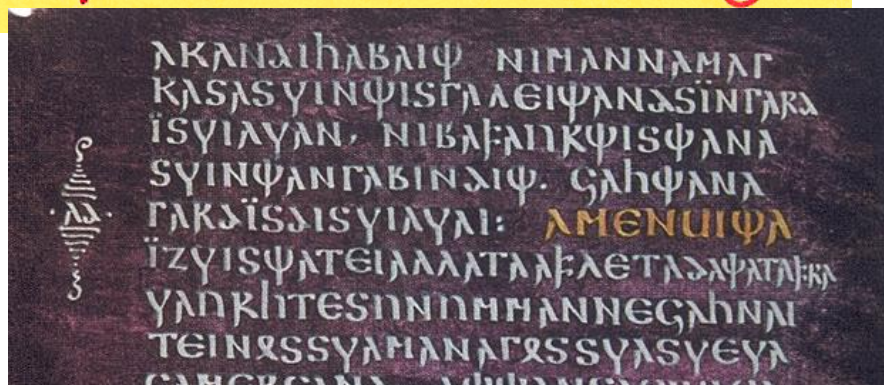
The quick brown fox
jumps over the lazy dog.



“Why is Hebrew written backwards?”

→

The quick brown fox
jumps over the lazy dog.



Codex Argenteus
(~500 AC)

←

Handwritten text in Hebrew script, showing the sentence "The quick brown fox jumps over the lazy dog." written from right to left.

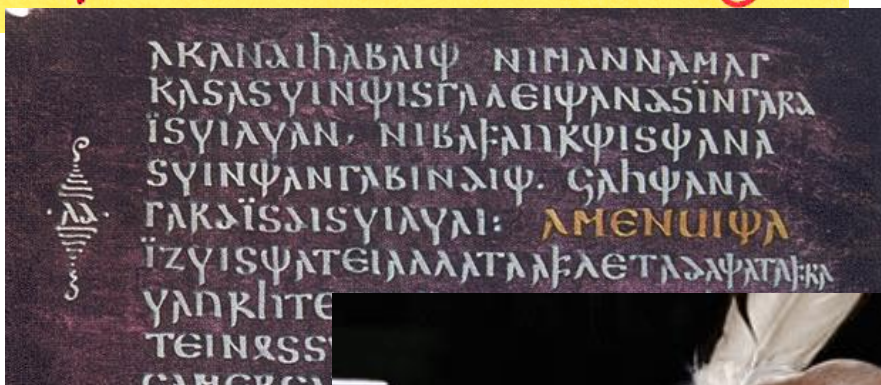


Cuneiform tablet
(~2000 BC)

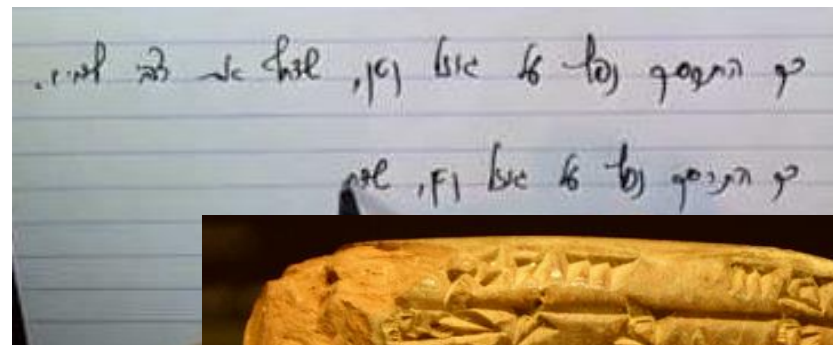
“Why is Hebrew written backwards?”

→

The quick brown fox
jumps over the lazy dog.



Codex Argenteus
(~500 AC)



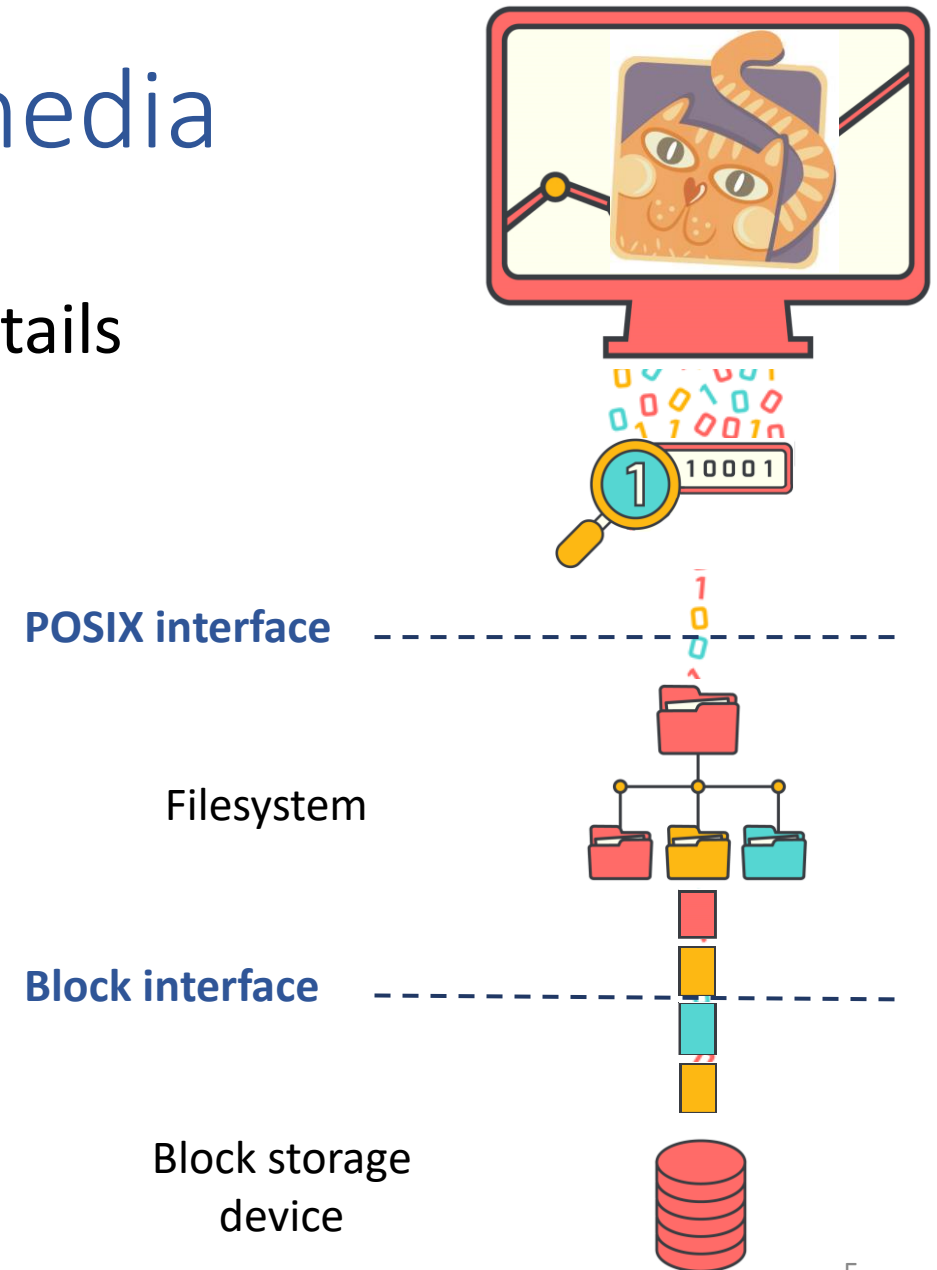
Cuneiform tablet
(~2000 BC)



It's the storage media!

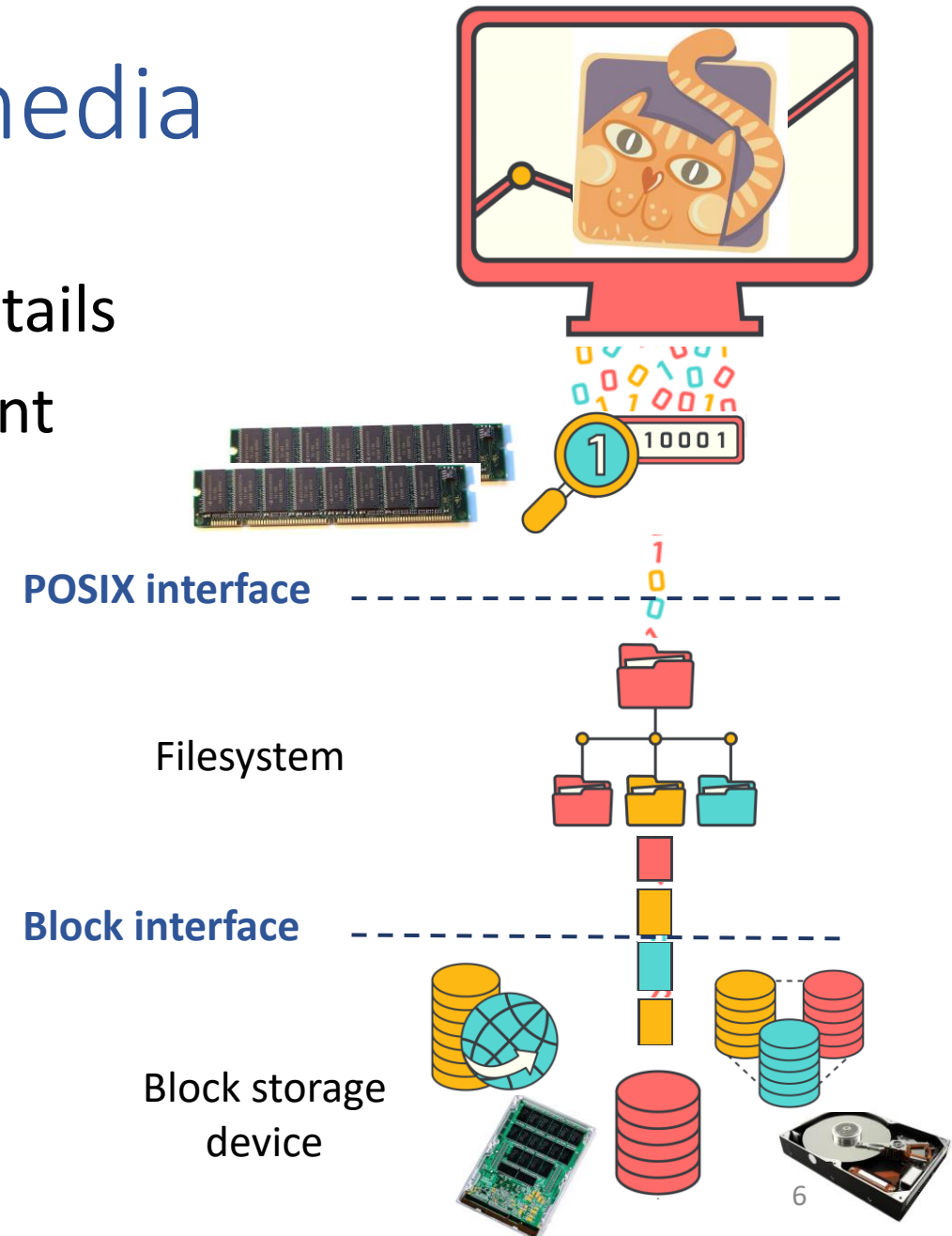
Storage systems and their media

- Abstraction layers hide complexity and details



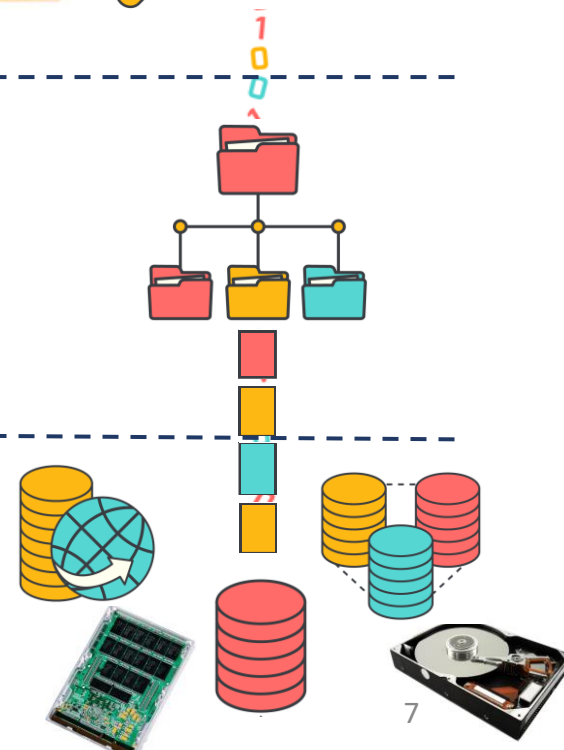
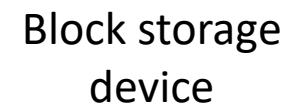
Storage systems and their media

- Abstraction layers hide complexity and details
- Add the new media as a direct replacement
 - Backwards compatibility
 - Little development effort
 - Quick adoption



A stylized illustration of a computer monitor with a red frame and base. The screen shows a square window with a white background. Inside the window is a cartoon orange cat face with large, expressive eyes and a small pink nose. A purple bandage is wrapped around the cat's head, covering its ears and part of its forehead. The cat's face is partially obscured by the bandage. The overall style is simple and colorful.

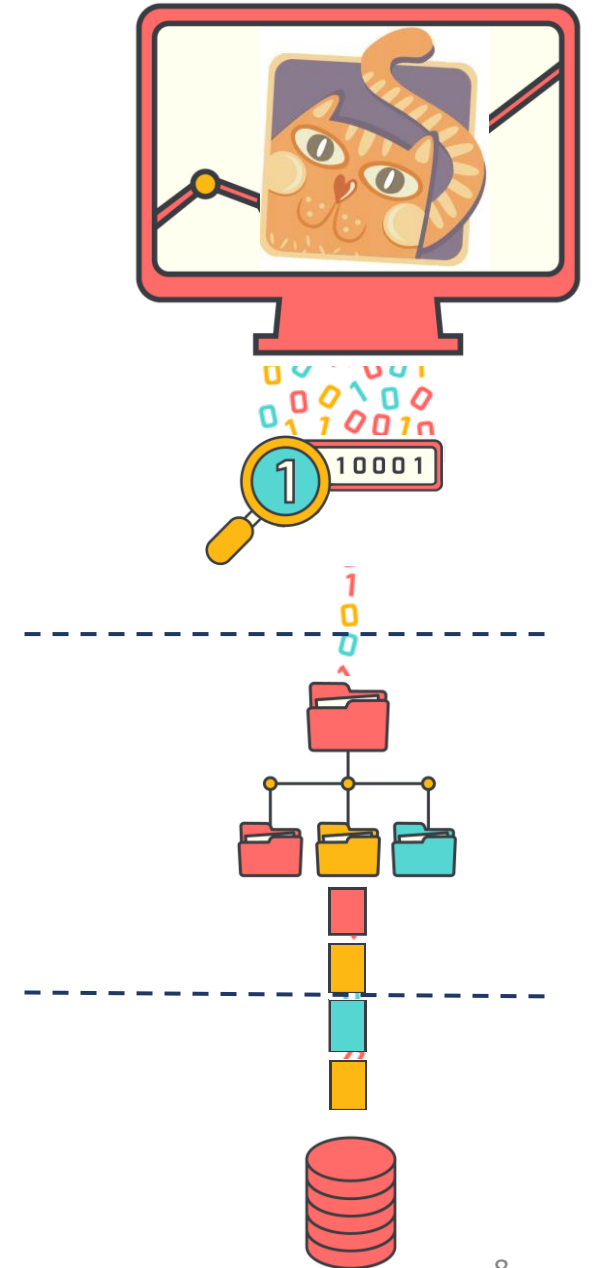
-
- An illustration showing two RAM modules on the left. On the right, a magnifying glass with a yellow handle and orange frame is focused on a large blue circle containing the number '1'. To the right of the magnifying glass is a white rectangular box with a red border containing the binary sequence '10001'. Above these elements, several colorful binary digits (0s and 1s) in red, yellow, and blue are scattered.



This talk

- How the **storage-systems community** addresses NVM
 - ⚠ Non-exhaustive, non-prioritized list of examples*
 - ⚠ Grossly over-simplified
 - Focus on insights and the adaptation process
- System model and characterization
- Unexpected bottlenecks
- Performance isolation / QoS

* Figures taken from respective cited papers and author presentations



What is NVM, exactly?

Non-volatile memory typically refers to storage in [semiconductor memory chips](#), which store data in [floating-gate memory cells](#) consisting of [floating-gate MOSFETs](#) ([metal–oxide–semiconductor field-effect transistors](#)), including [flash memory storage](#) such as [NAND flash](#) and [solid-state drives](#) (SSD).



WIKIPEDIA
The Free Encyclopedia

Other examples of non-volatile memory include [read-only memory](#) (ROM), [EPROM](#) (erasable [programmable ROM](#)) and [EEPROM](#) (electrically erasable programmable ROM), [ferroelectric RAM](#), most types of [computer data storage](#) devices (e.g. [disk storage](#), [hard disk drives](#), [optical discs](#), [floppy disks](#), and [magnetic tape](#)), and early computer storage methods such as [punched tape](#) and [cards](#).^[1]



What is NVM, exactly?

Non-volatile memory typically refers to storage in [semiconductor memory chips](#), which store data in [floating-gate memory cells](#) consisting of [floating-gate MOSFETs](#) ([metal–oxide–semiconductor field-effect transistors](#)), including [flash memory storage](#) such as [NAND flash](#) and [solid-state drives](#) (SSD).



WIKIPEDIA
The Free Encyclopedia

Other examples of non-volatile memory include [read-only memory](#) (ROM), [EPROM](#) (erasable [programmable ROM](#)) and [EEPROM](#) (electrically erasable programmable ROM), [ferroelectric RAM](#), most types of [computer data storage](#) devices (e.g. [disk storage](#), [hard disk drives](#), [optical discs](#), [floppy disks](#), and [magnetic tape](#)), and early computer storage methods such as [punched tape](#) and [cards](#).^[1]

- For many purposes, *fast* and *non-volatile* is NVM
 - Write-ahead logging
 - Metadata or index



What is NVM, exactly?

Non-volatile memory typically refers to storage in [semiconductor memory chips](#), which store data in [floating-gate memory cells](#) consisting of [floating-gate MOSFETs](#) ([metal–oxide–semiconductor field-effect transistors](#)), including [flash memory storage](#) such as [NAND flash](#) and [solid-state drives](#) (SSD).



WIKIPEDIA
The Free Encyclopedia

Other examples of non-volatile memory include [read-only memory](#) (ROM), [EPROM](#) (erasable [programmable ROM](#)) and [EEPROM](#) (electrically erasable programmable ROM), [ferroelectric RAM](#), most types of [computer data storage](#) devices (e.g. [disk storage](#), [hard disk drives](#), [optical discs](#), [floppy disks](#), and [magnetic tape](#)), and early computer storage methods such as [punched tape](#) and [cards](#).^[1]

- For many purposes, *fast* and *non-volatile* is NVM
 - Write-ahead logging
 - Metadata or index
- For our purpose, *byte-addressable* and *non-DRAM*



What is NVM, exactly?

- Intel Optane persistent memory

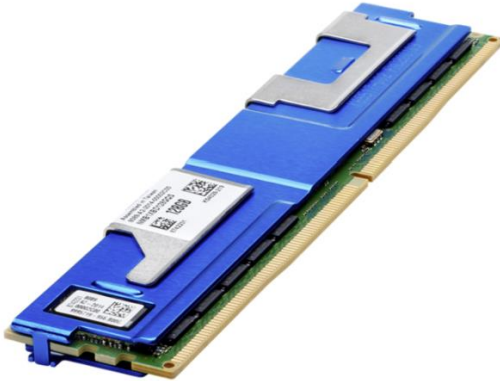


- For our purpose, *byte-addressable* and *non-DRAM*



What is NVM, exactly?

- Intel Optane persistent memory



- [Samsung's CMM-H \(CXL Memory Module – Hybrid\)](#)

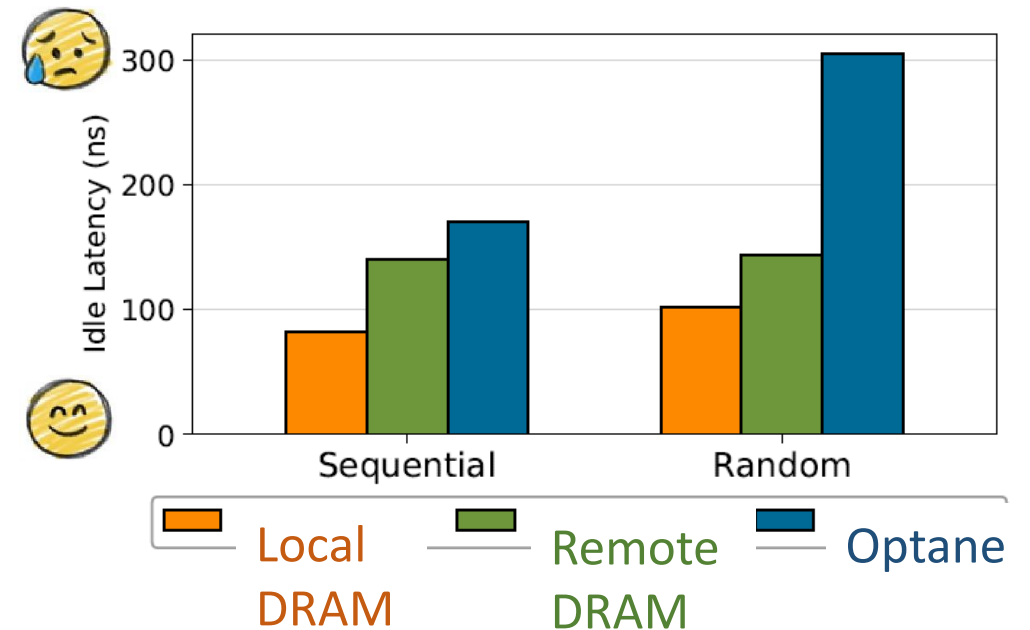


- For our purpose, *byte-addressable* and *non-DRAM*



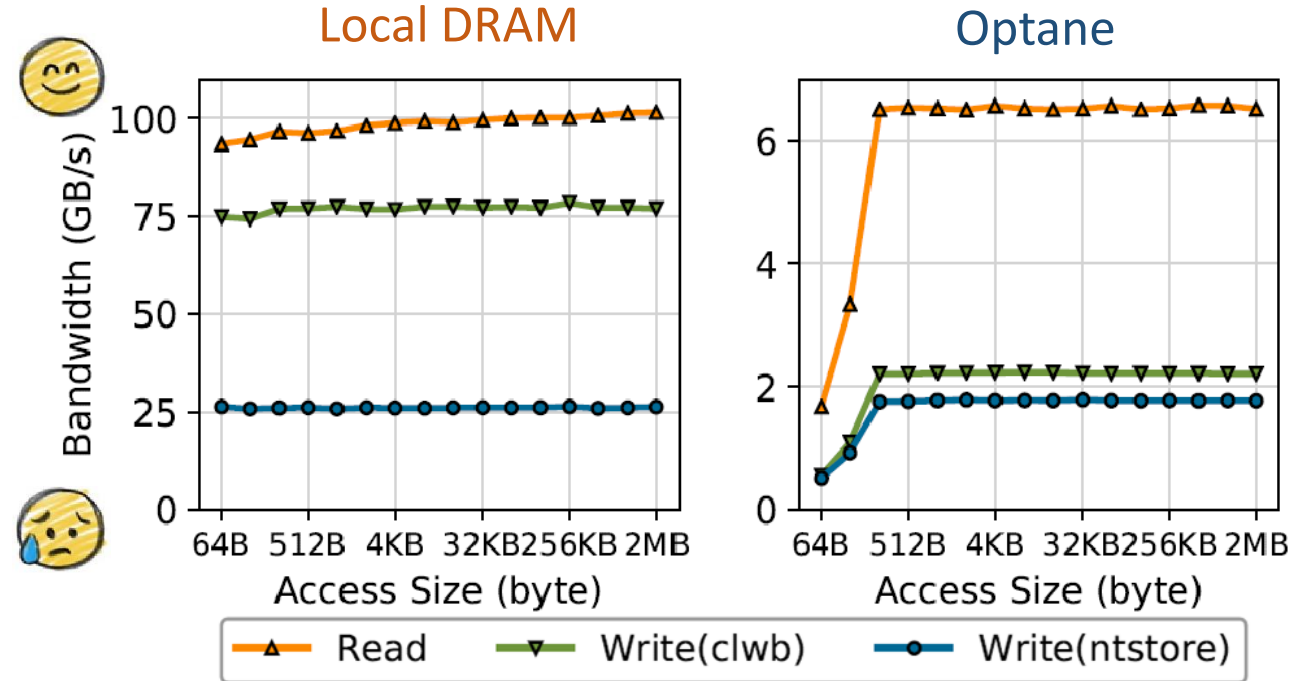
Characterization (Optane)

- Faster than HDD and SSD but slower than DRAM
- Does not behave like DRAM
 - Random \neq Sequential



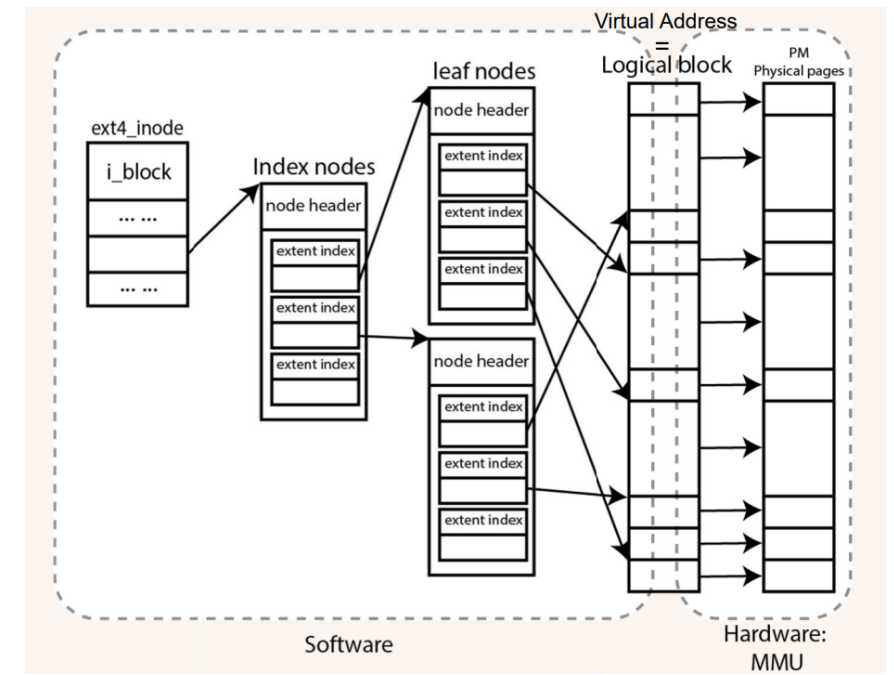
Characterization (Optane)

- Faster than HDD and SSD but slower than DRAM
- Does not behave like DRAM
 - Random \neq Sequential
 - Reads \neq Writes
 - Small \neq Large
 - Much lower bandwidth
 - (Much more interference)



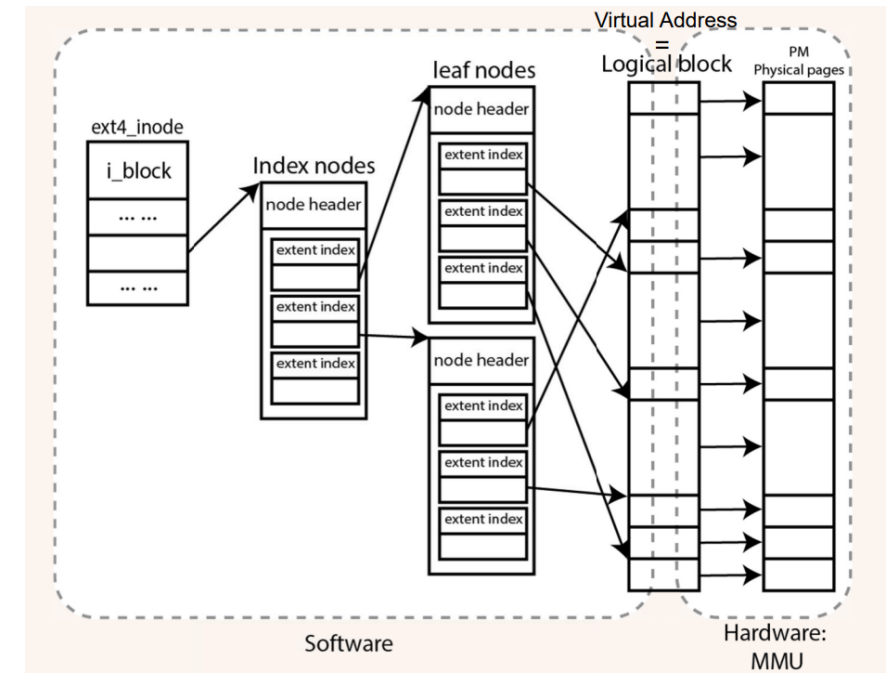
Rethinking file-system indexing

- **Bottleneck:** *tree-based inode index* incurs high overhead with NVM
 - Up to 63% of the time spend on FS operations (e.g., file append)



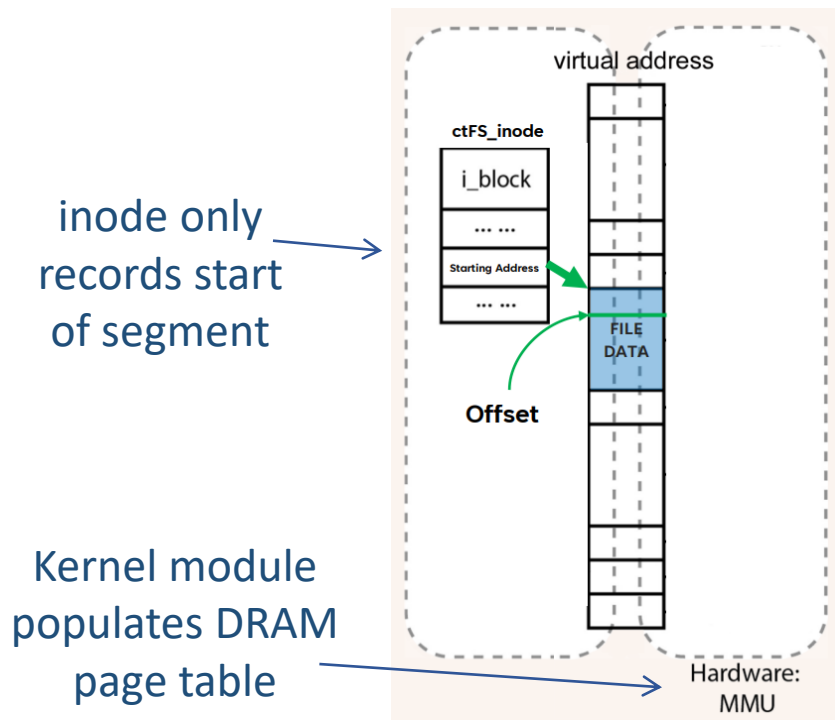
Rethinking file-system indexing

- **Bottleneck:** *tree-based inode index* incurs high overhead with NVM
 - Up to 63% of the time spend on FS operations (e.g., file append)
- **Insights:**
 - Use hardware-based translation
 - Don't care about physical contiguity
 - Persist small updates quickly



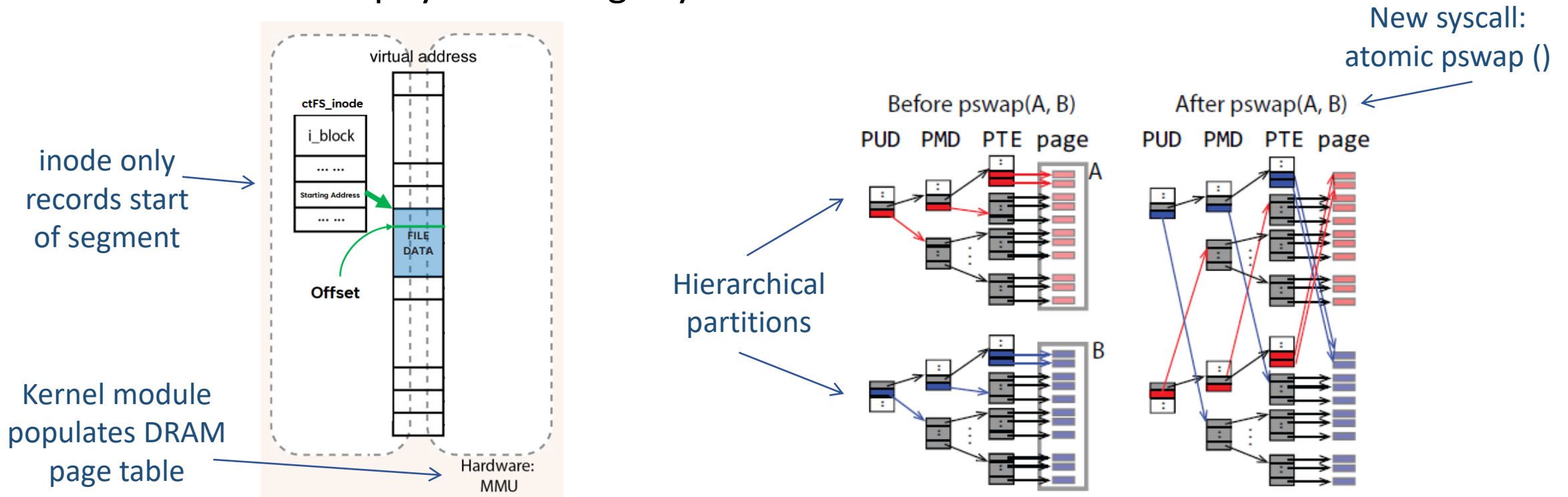
Rethinking file-system indexing: ctFS

- Use hardware-based translation
- Don't care about physical contiguity



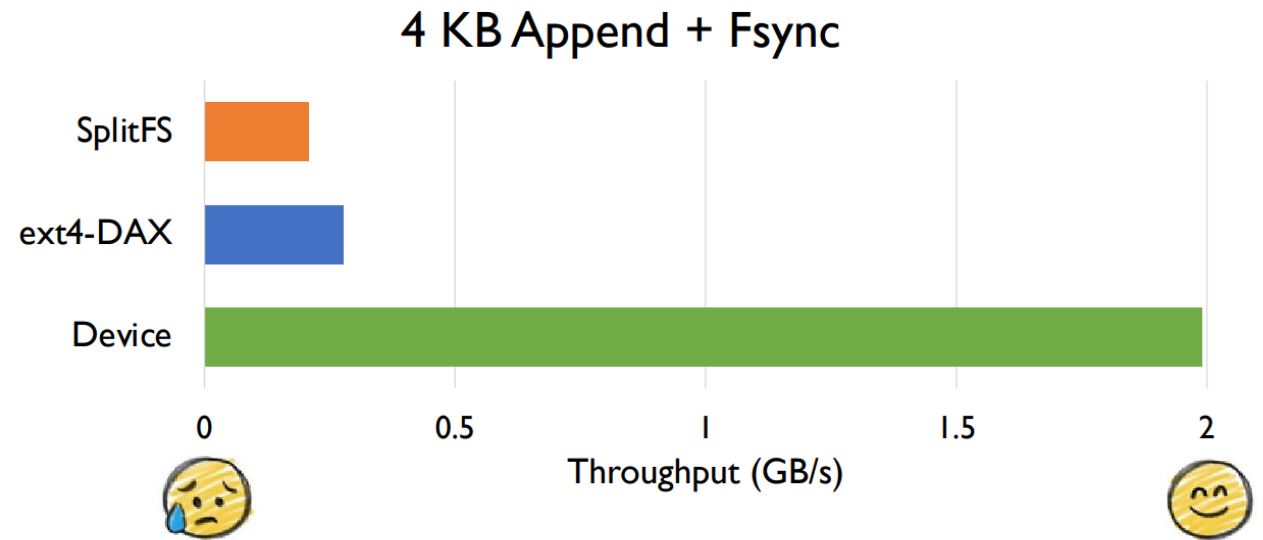
Rethinking file-system indexing: ctFS

- Use hardware-based translation
- Don't care about physical contiguity
- Persist small updates quickly



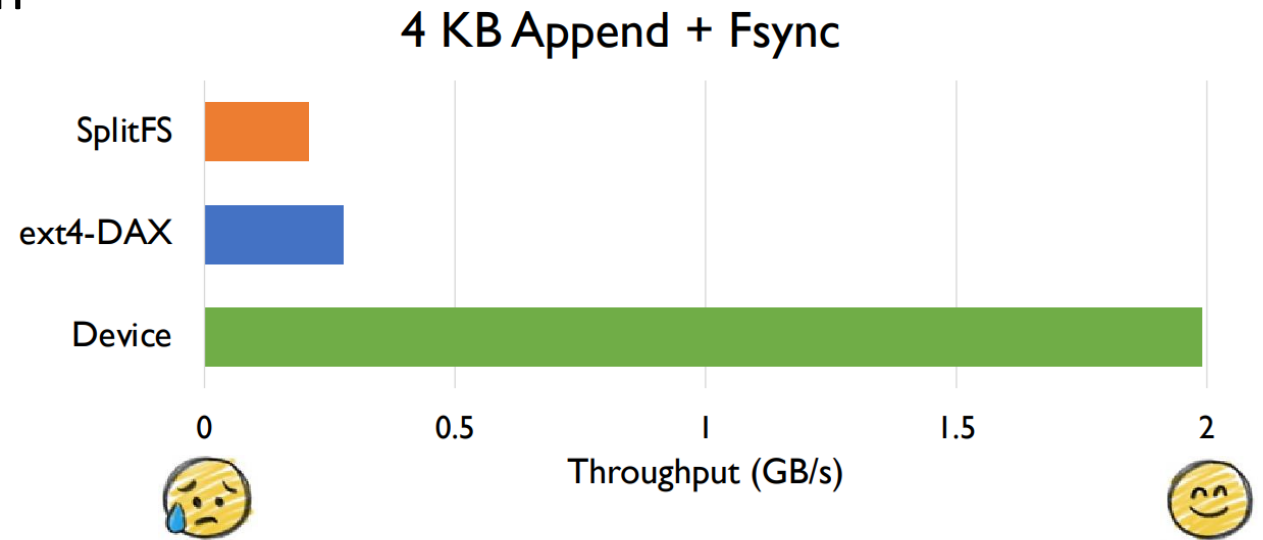
Rethinking kernel-space vs. user-space

- **Bottleneck:** metadata handling in *kernel space* is inefficient



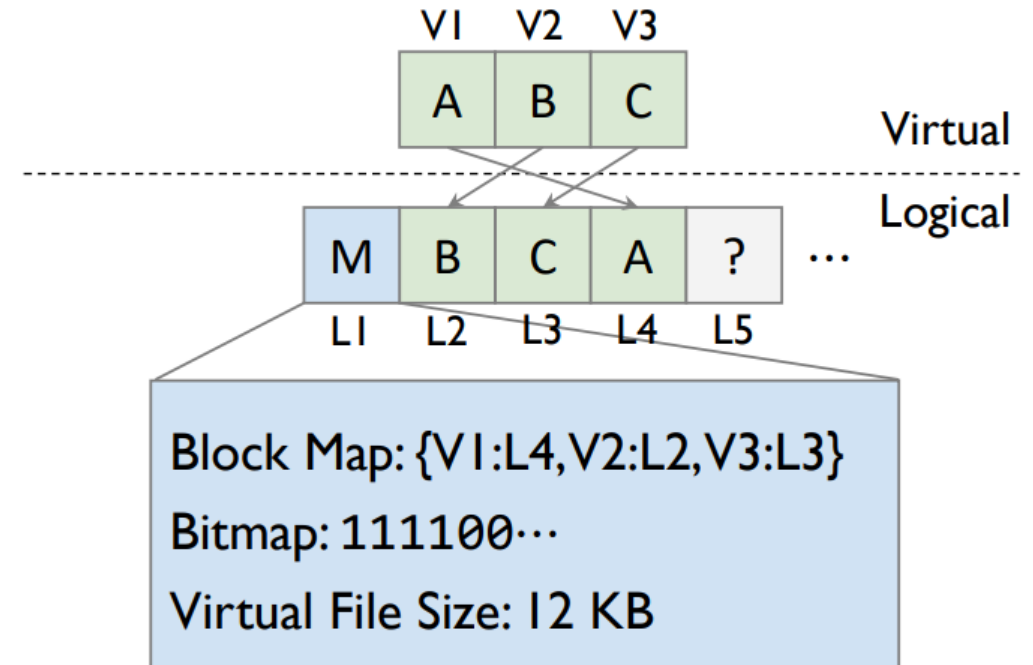
Rethinking kernel-space vs. user-space

- **Bottleneck:** metadata handling in *kernel space* is inefficient
- **Insights:**
 - Some metadata can be safely handled in user space
 - Rethink data/metadata separation



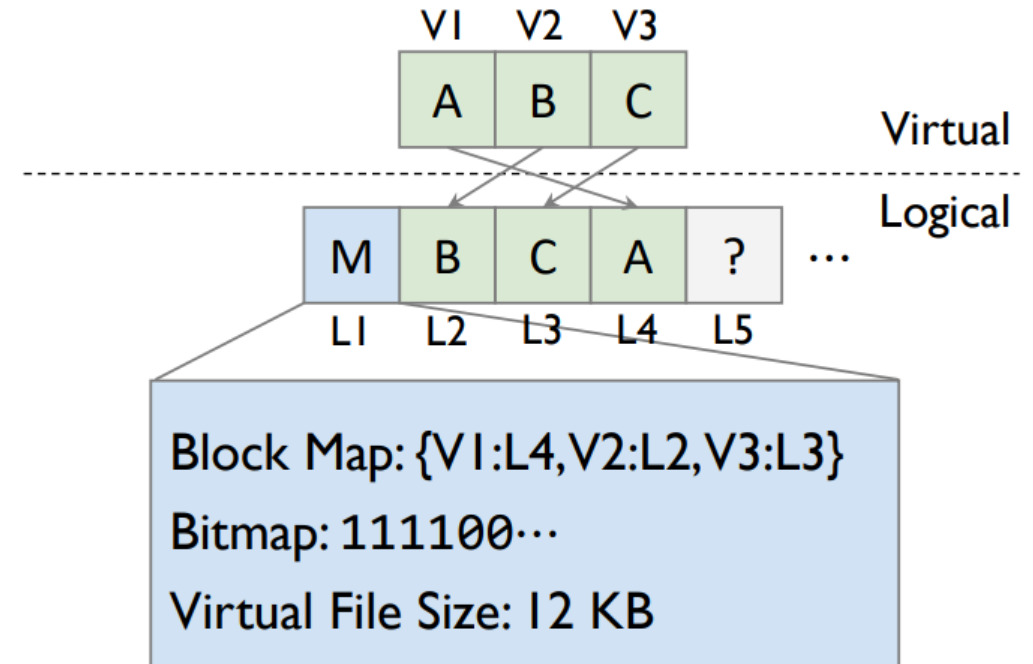
Rethinking kernel-space vs. user-space: MadFS

- Some metadata can be safely handled in user space
- Rethink data/metadata separation
- **Embedded metadata**
 - “Virtual” file block map and size



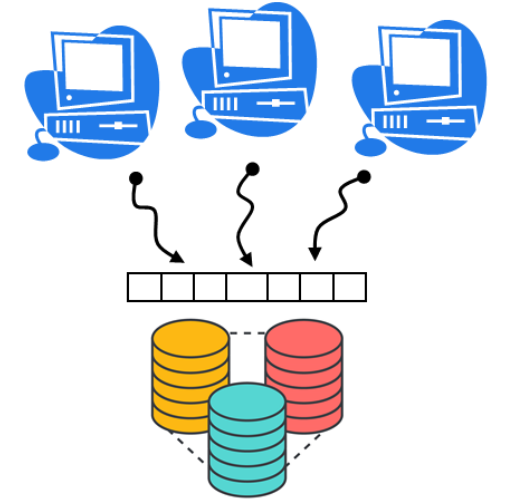
Rethinking kernel-space vs. user-space: MadFS

- Some metadata can be safely handled in user space
- Rethink data/metadata separation
- **Embedded metadata**
 - “Virtual” file block map and size
- **Kernel-managed metadata**
 - Logical-to-physical mapping
 - File permissions
 - Directory structures

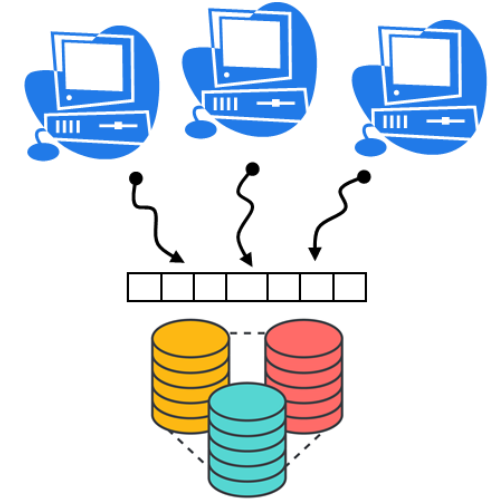


Rethinking shared caches

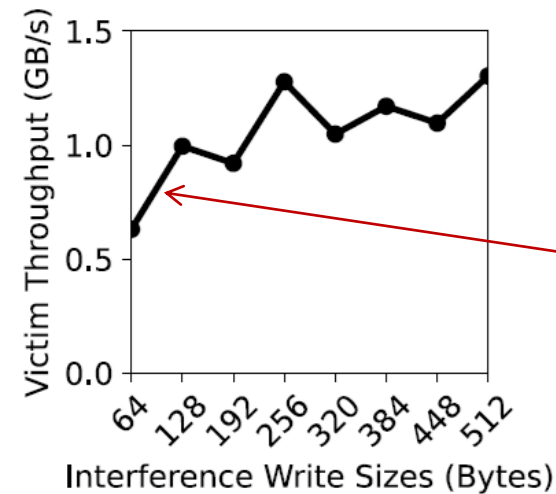
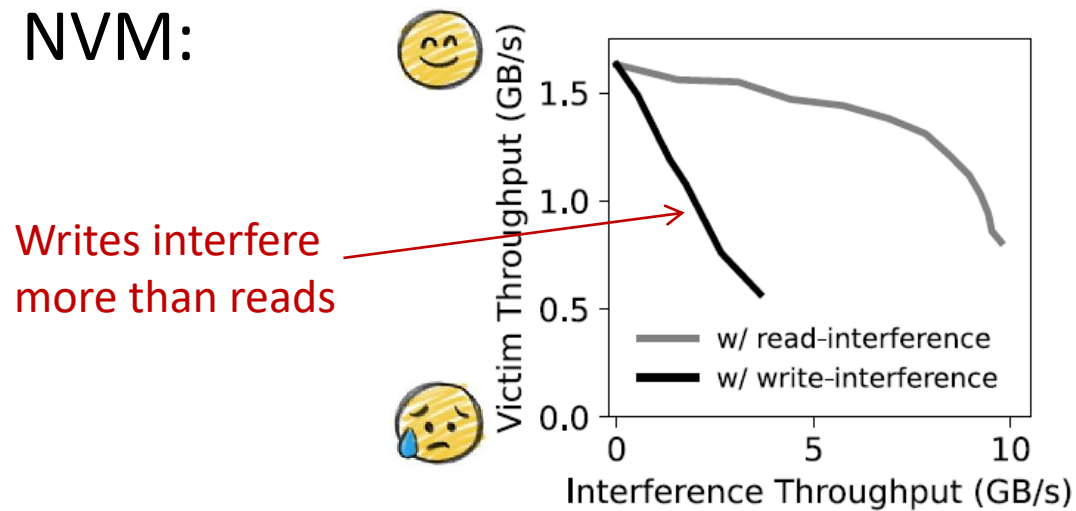
- DRAM:
 - More space allocation = higher hit rate = better performance
 - Higher bandwidth = more usage = more interference
 - Tenant A is too slow → throttle noisy neighbor B with max bandwidth



Rethinking shared caches

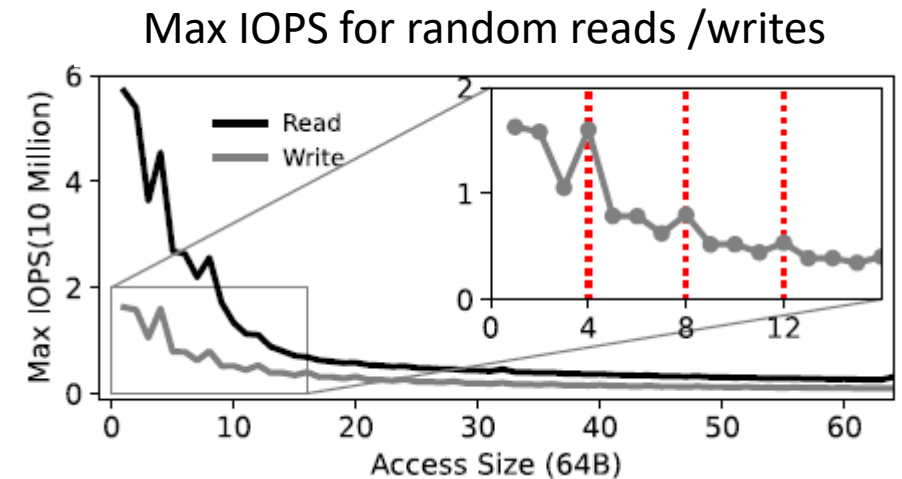


- DRAM:
 - More space allocation = higher hit rate = better performance
 - Higher bandwidth = more usage = more interference
 - Tenant A is too slow → throttle noisy neighbor B with max bandwidth
- NVM:



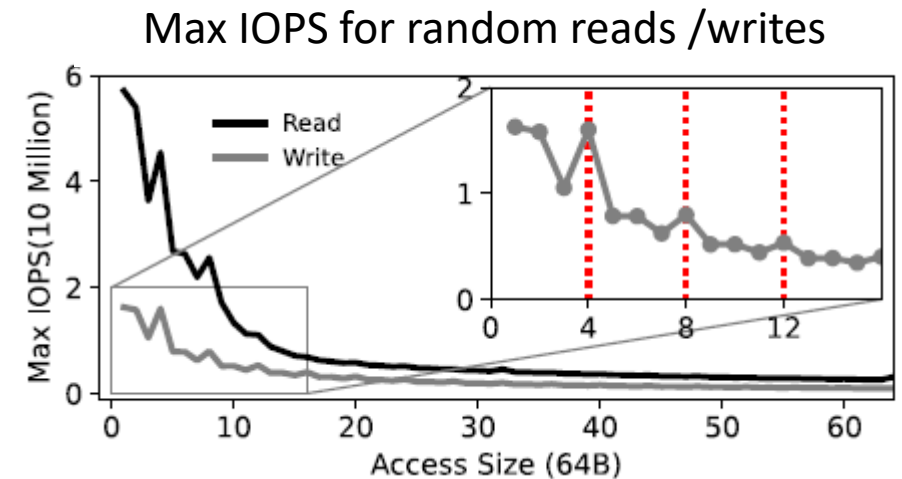
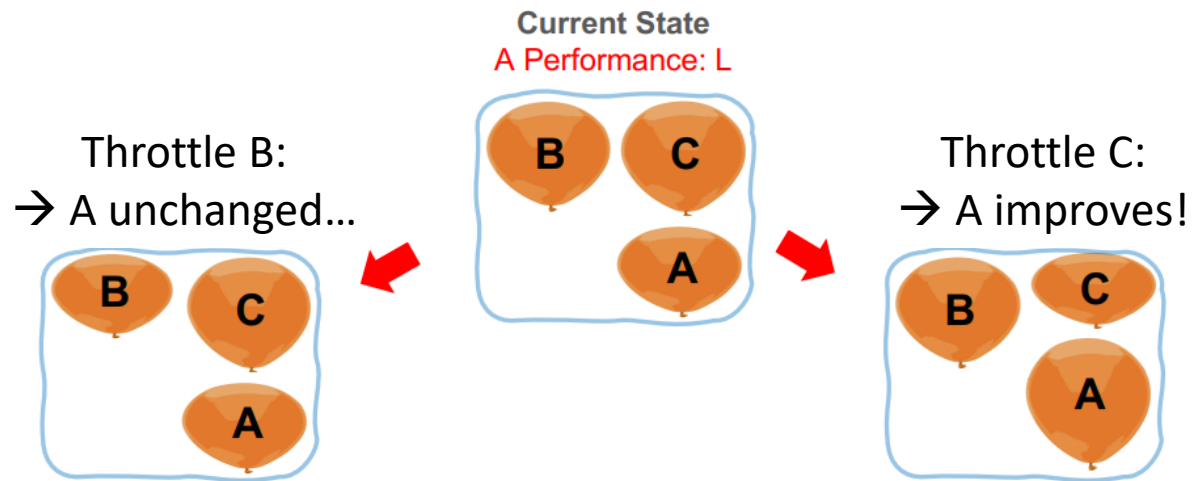
Rethinking shared caches: NyxCache

- Profile max performance of different access types



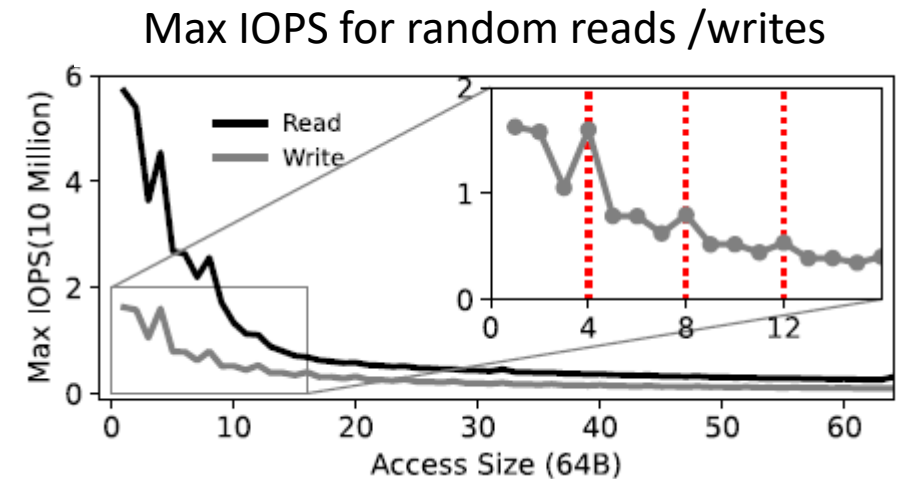
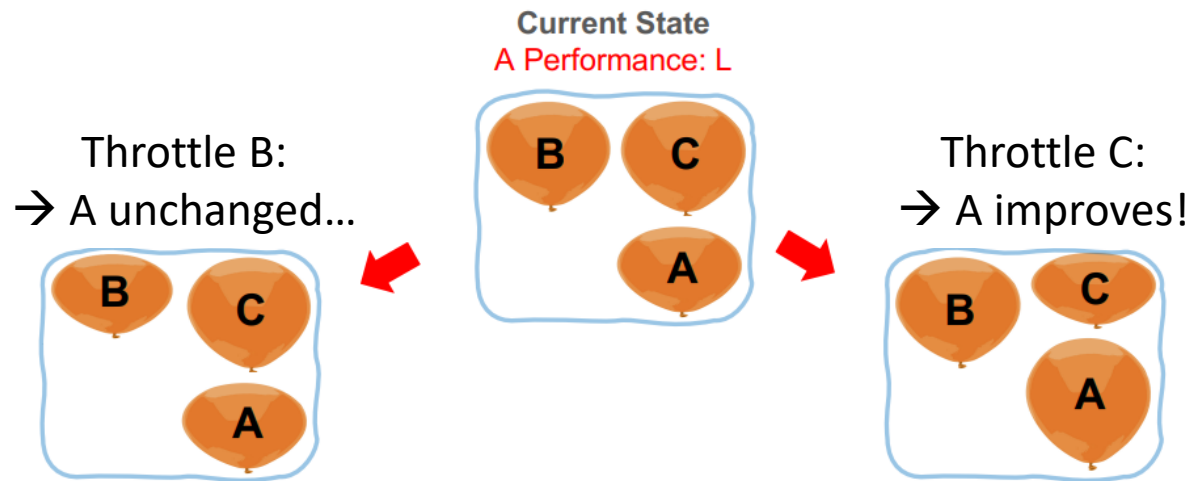
Rethinking shared caches: NyxCache

- Profile max performance of different access types
- Monitor individual tenants and their interaction

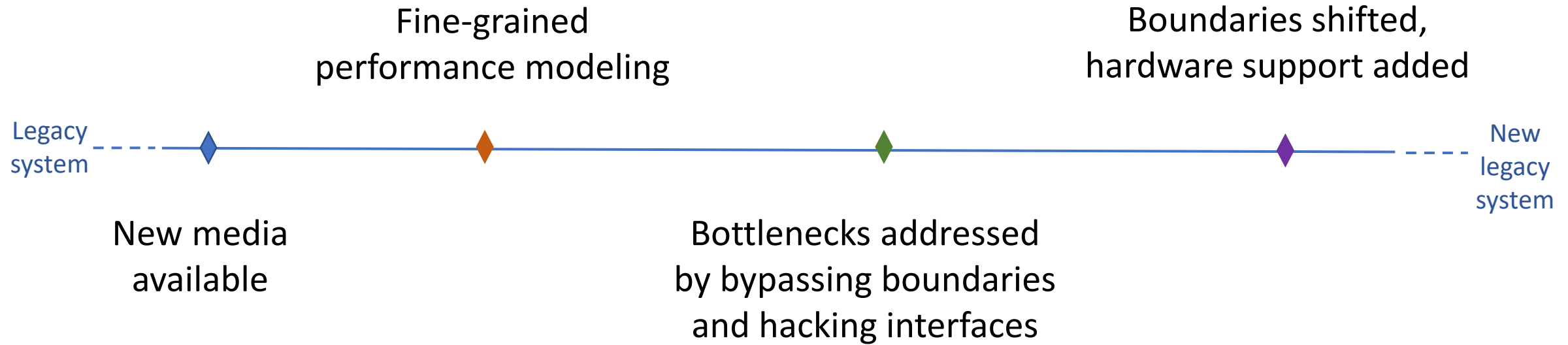


Rethinking shared caches: NyxCache

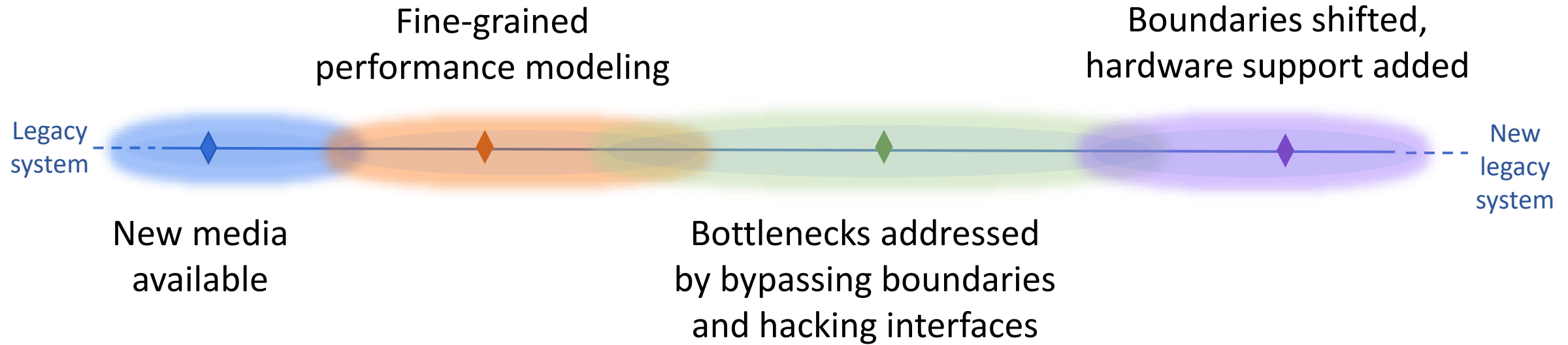
- Profile max performance of different access types
- Monitor individual tenants and their interaction
- Library throttles tenant using most resources / causing most interference



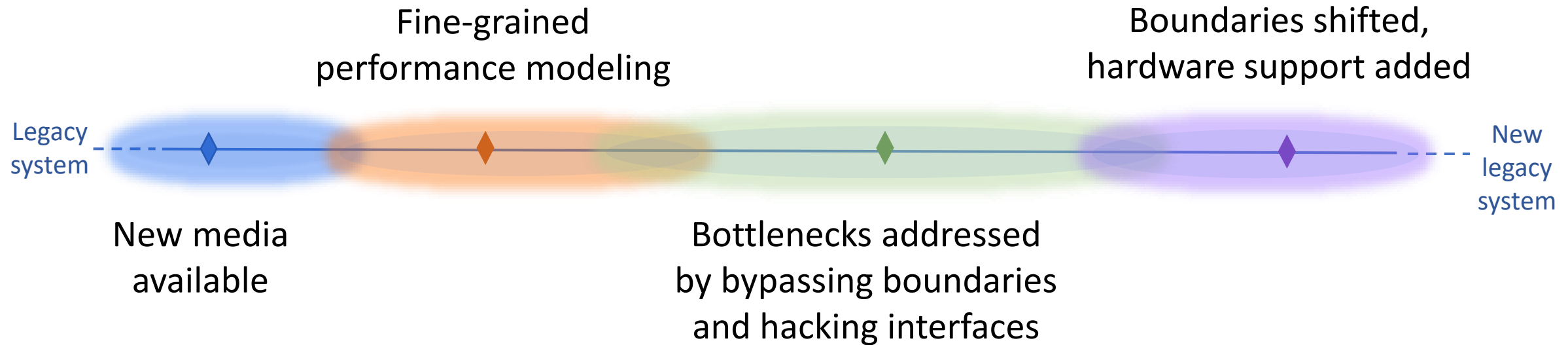
Conclusions



Conclusions

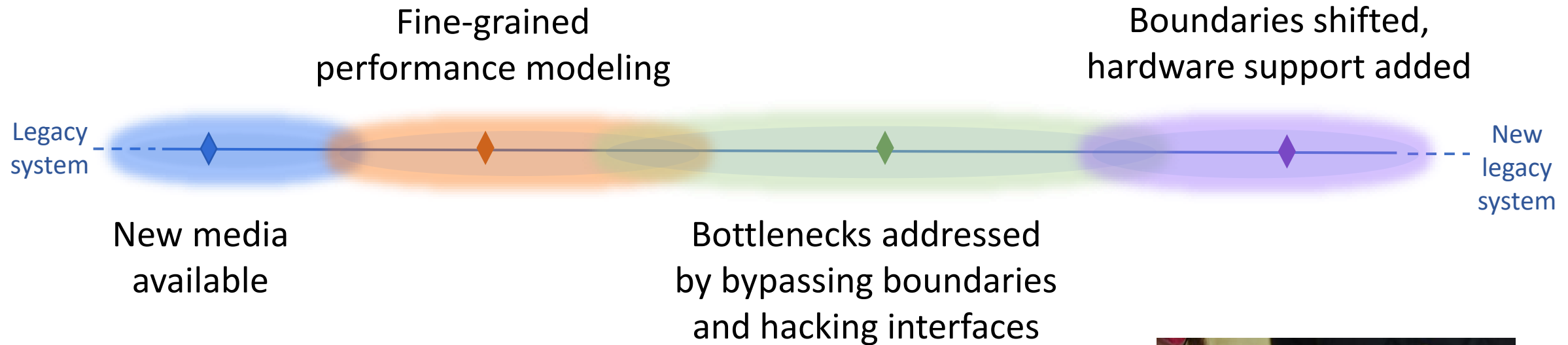


Conclusions



- And what about our written language...?

Conclusions



- And what about our written language...?



?

