



**PURESTORAGE<sup>®</sup>**



# How to Build Exascale Storage with QLC

**Riley Thomasson**

Member of Technical Staff  
Pure Storage

# Agenda

---

**Context** Exascale Storage

---

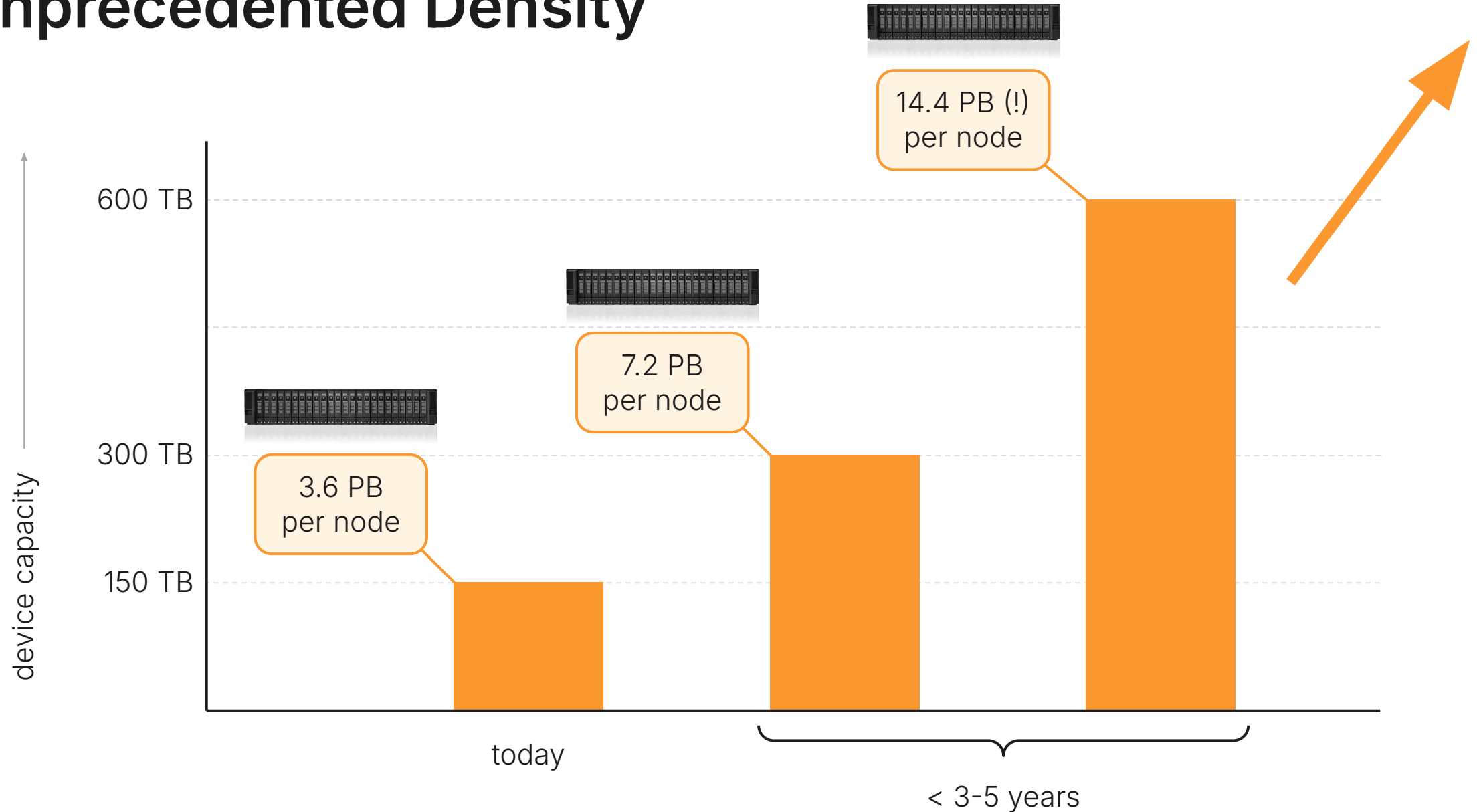
**Problem** QLC Challenges

---

**Solution** Architectural Perspective

---

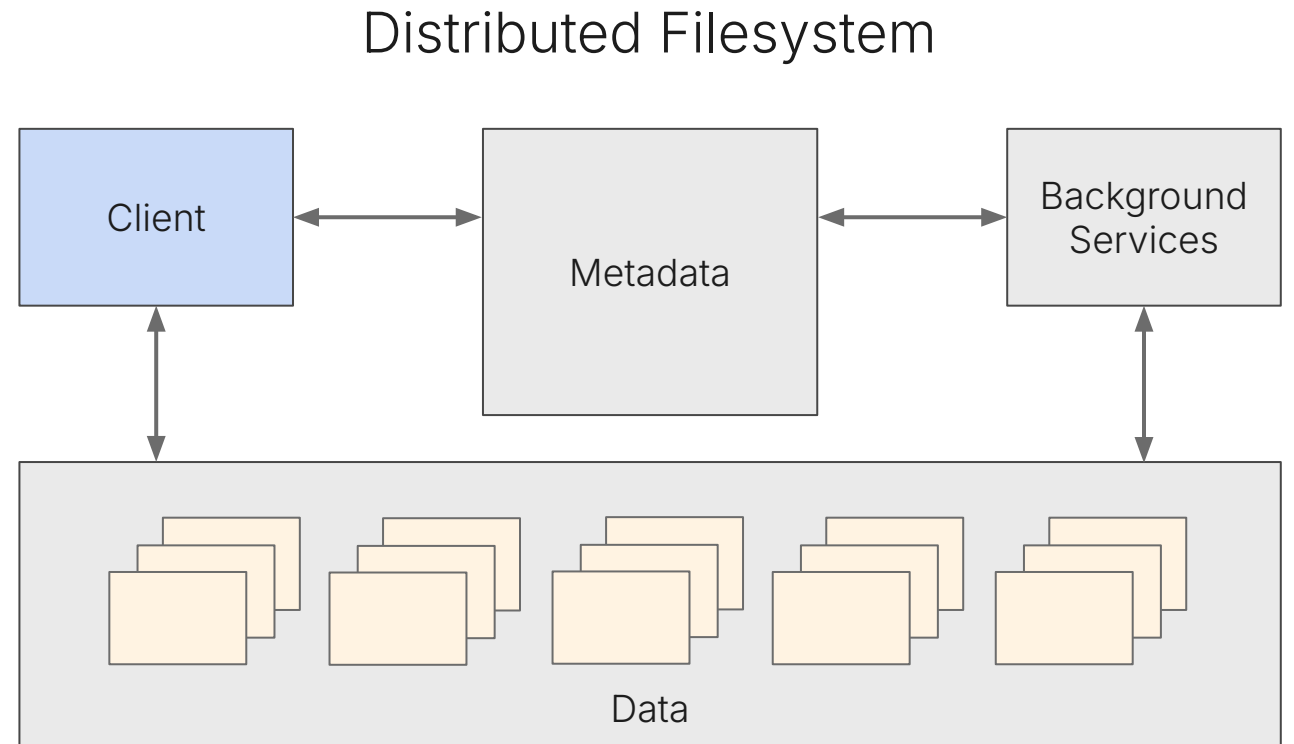
# Unprecedented Density



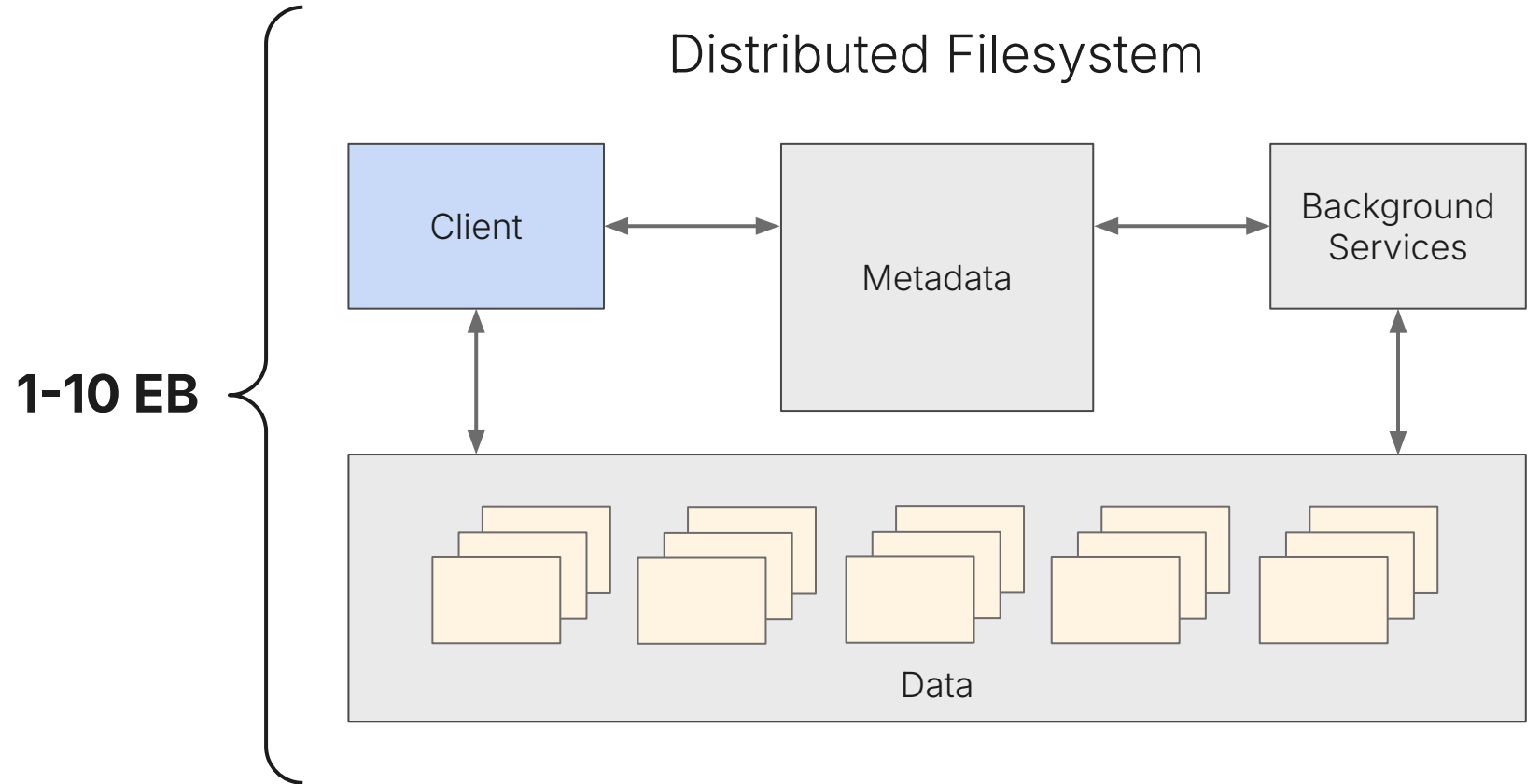
# Exascale Storage

## Use Cases

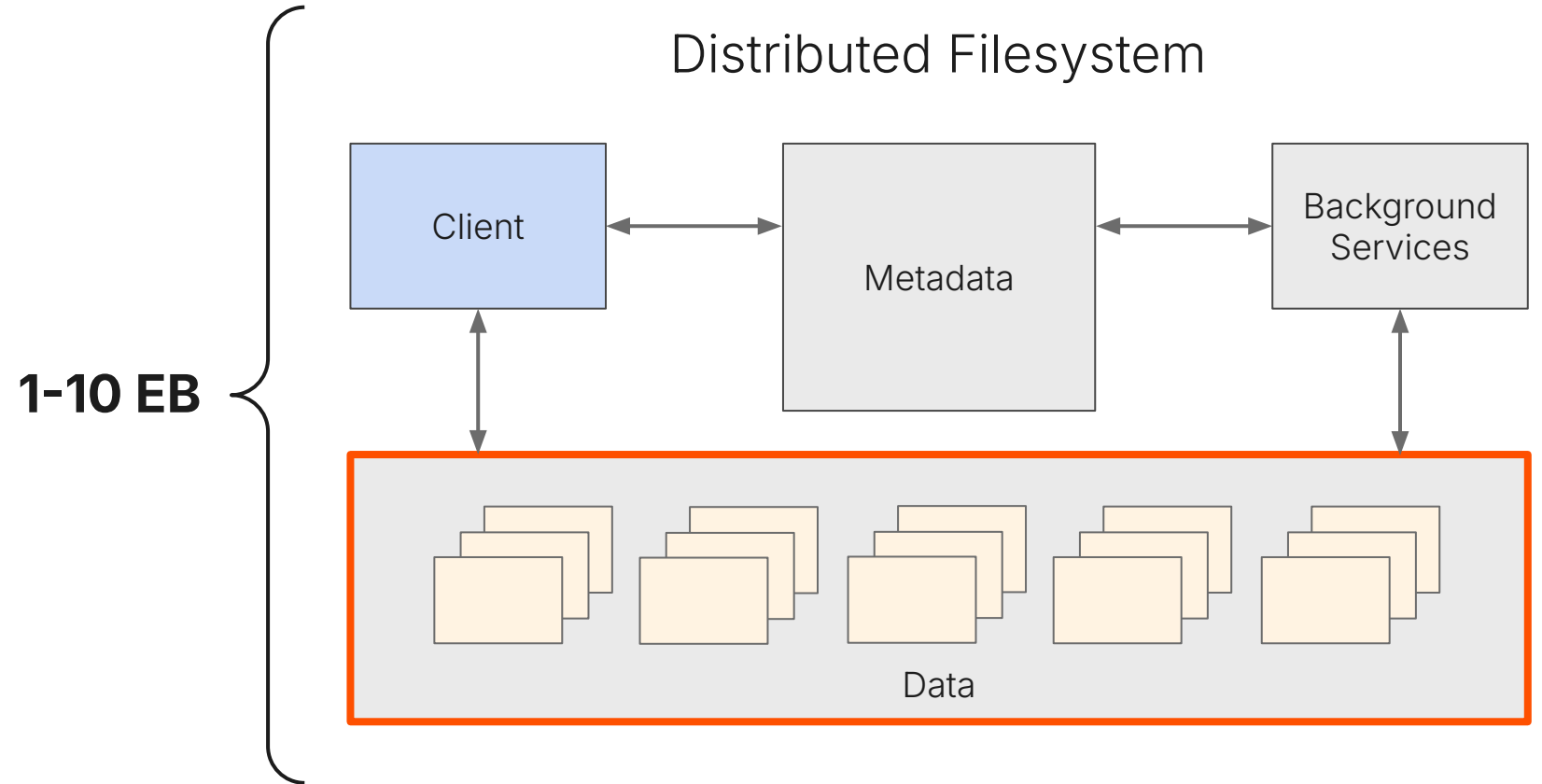
- GenAI
- Warm/Cold Object
- Archive
- Database



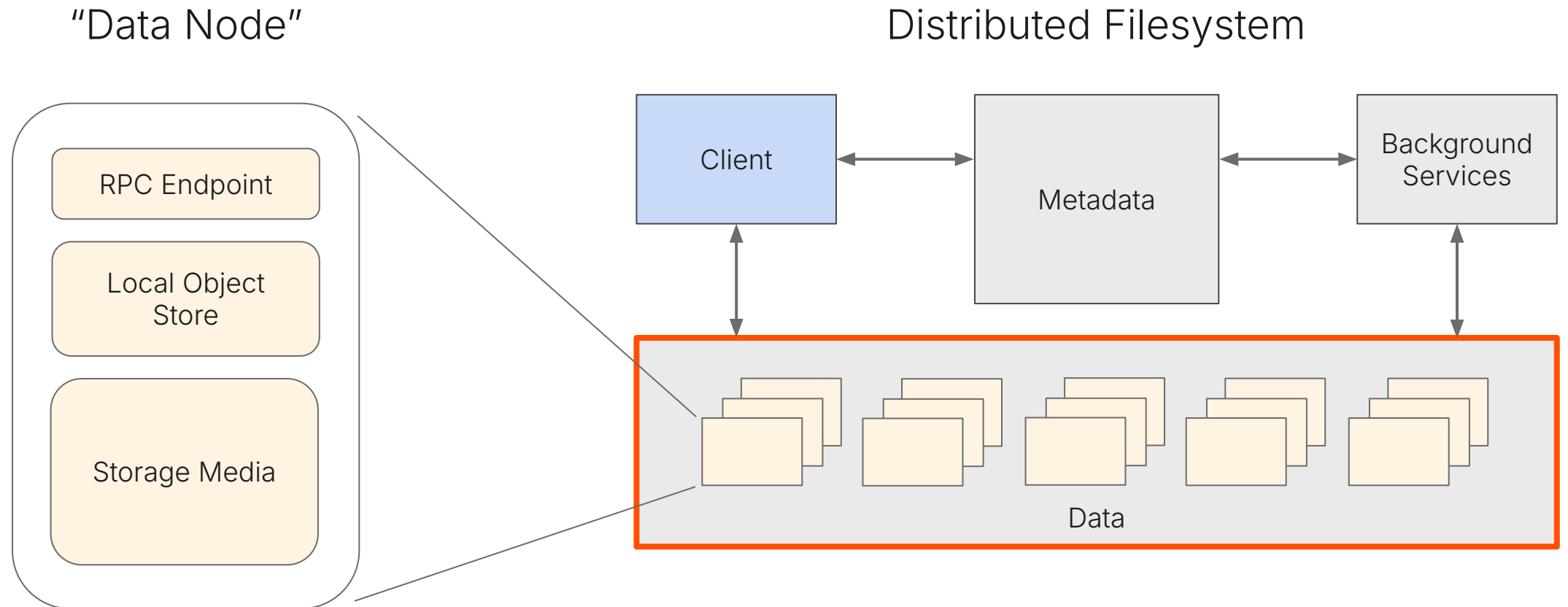
# Exascale Storage



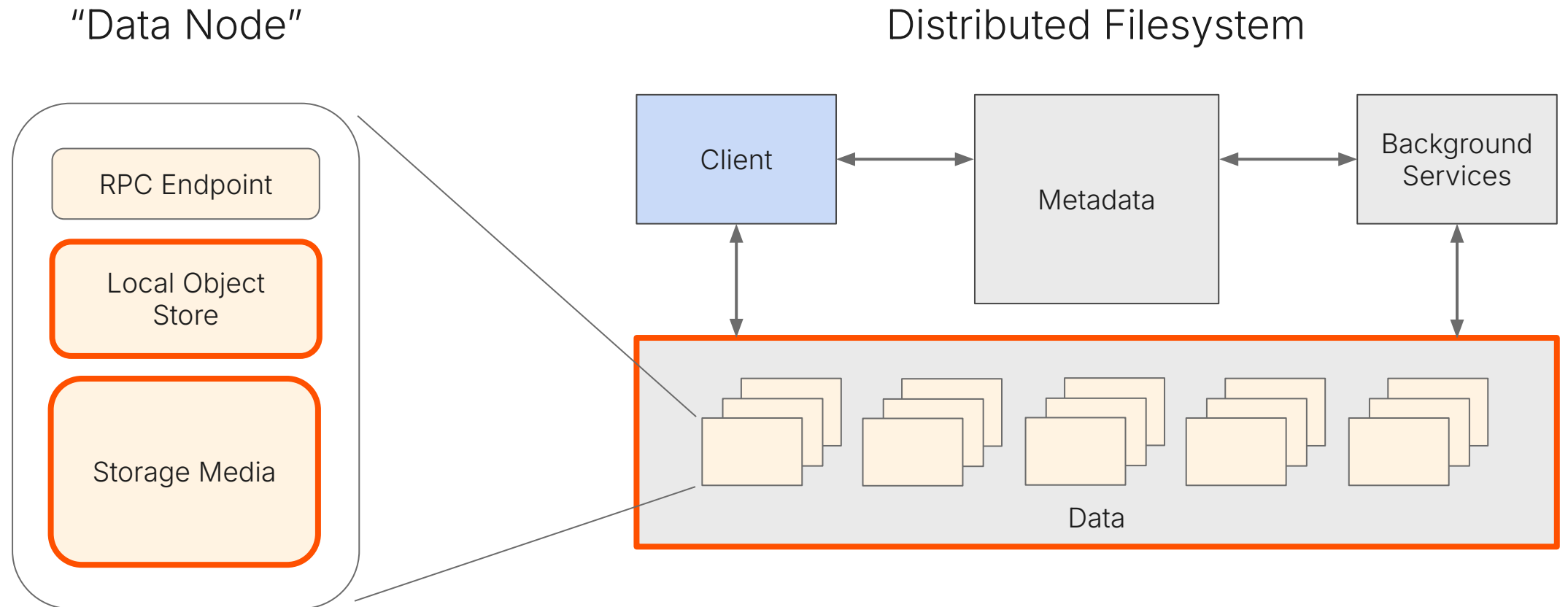
# Exascale Storage



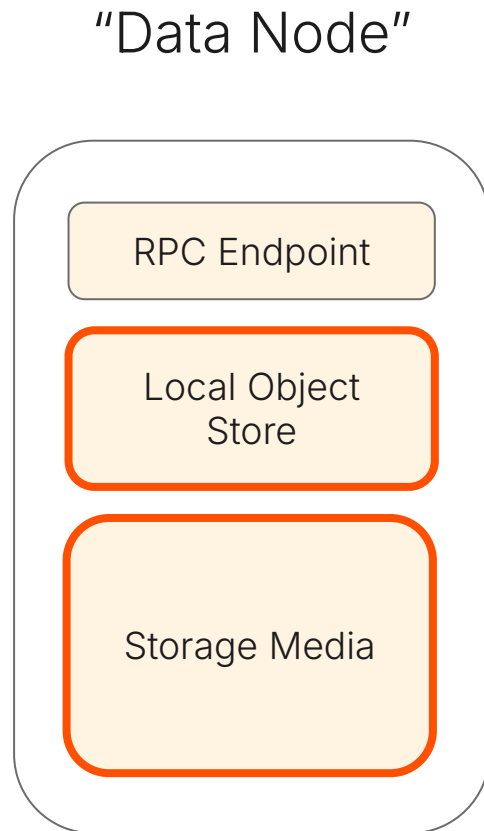
# Exascale Storage



# Exascale Storage



# Exascale Storage



## "Local Object Store"

- flat namespace
- append-only objects
- protected by cluster-level ECC
- no standard, many proprietary implementations

## Not: **general purpose storage**

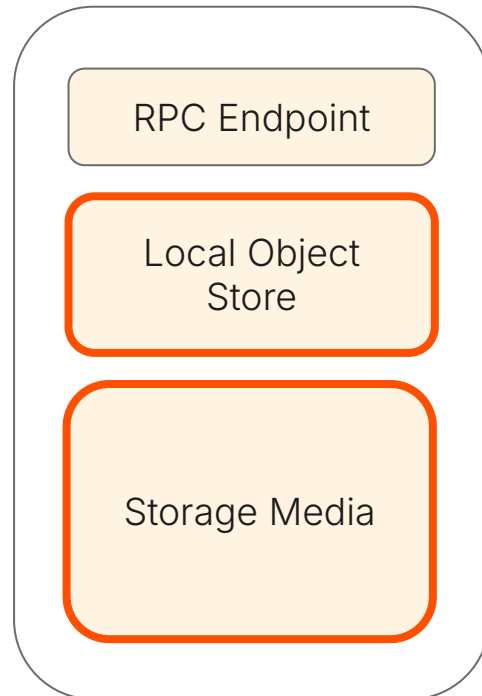
- filesystems (xfs, ext4, ntfs)
- fully-featured, hierarchical, overwrite-in-place
- standardized (POSIX)

## Not: **cloud "object storage"**

- HA
- external, standardized API
- served by a large distributed cluster

# Exascale Storage

"Data Node"



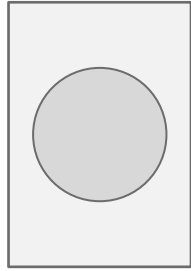
At exascale, efficiency is paramount!

- power/TB
- performance/TB
- cost (media, CPU, memory, racks...)
- end-to-end perspective

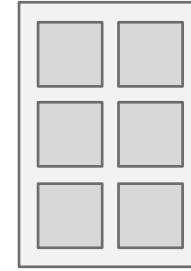
Opportunities for efficiency:

1. leverage QLC media
2. customize for the "data node" workload

# QLC



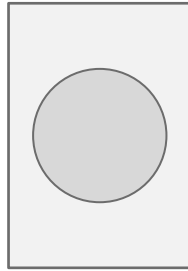
HAMR HDD  
30 TB



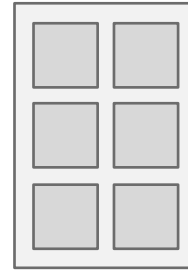
TLC  
30 TB

per TB	\$	1x	6x
	power	1x	2.8x
	read throughput	1x	280x
	write throughput	1x	90x
	IOPS	1x	30,000x

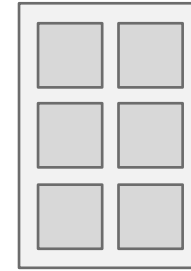
# QLC



HAMR HDD  
30 TB



QLC  
150 TB



TLC  
30 TB

per TB	\$	1x	3x	6x
	power	1x	0.5x	2.8x
	read throughput	1x	60x	280x
	write throughput	1x	7x	90x
	IOPS	1x	2,000x	30,000x

# Agenda

---

**Context** Exascale Storage

---

**Problem** QLC Challenges

---

**Solution** Architectural Perspective

---

# QLC Challenges



Scale



Endurance



Reliability

# QLC Challenges



Scale

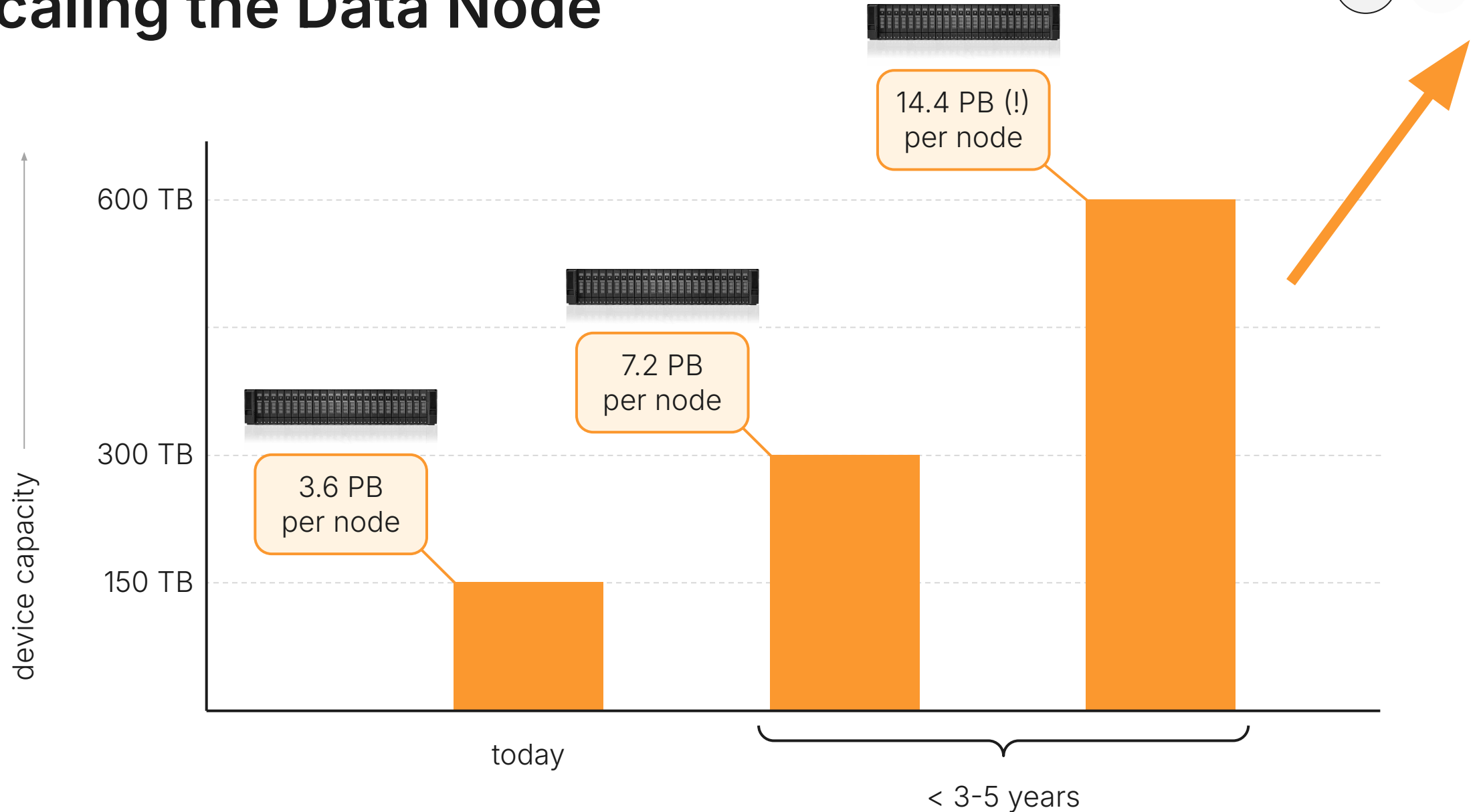


Endurance

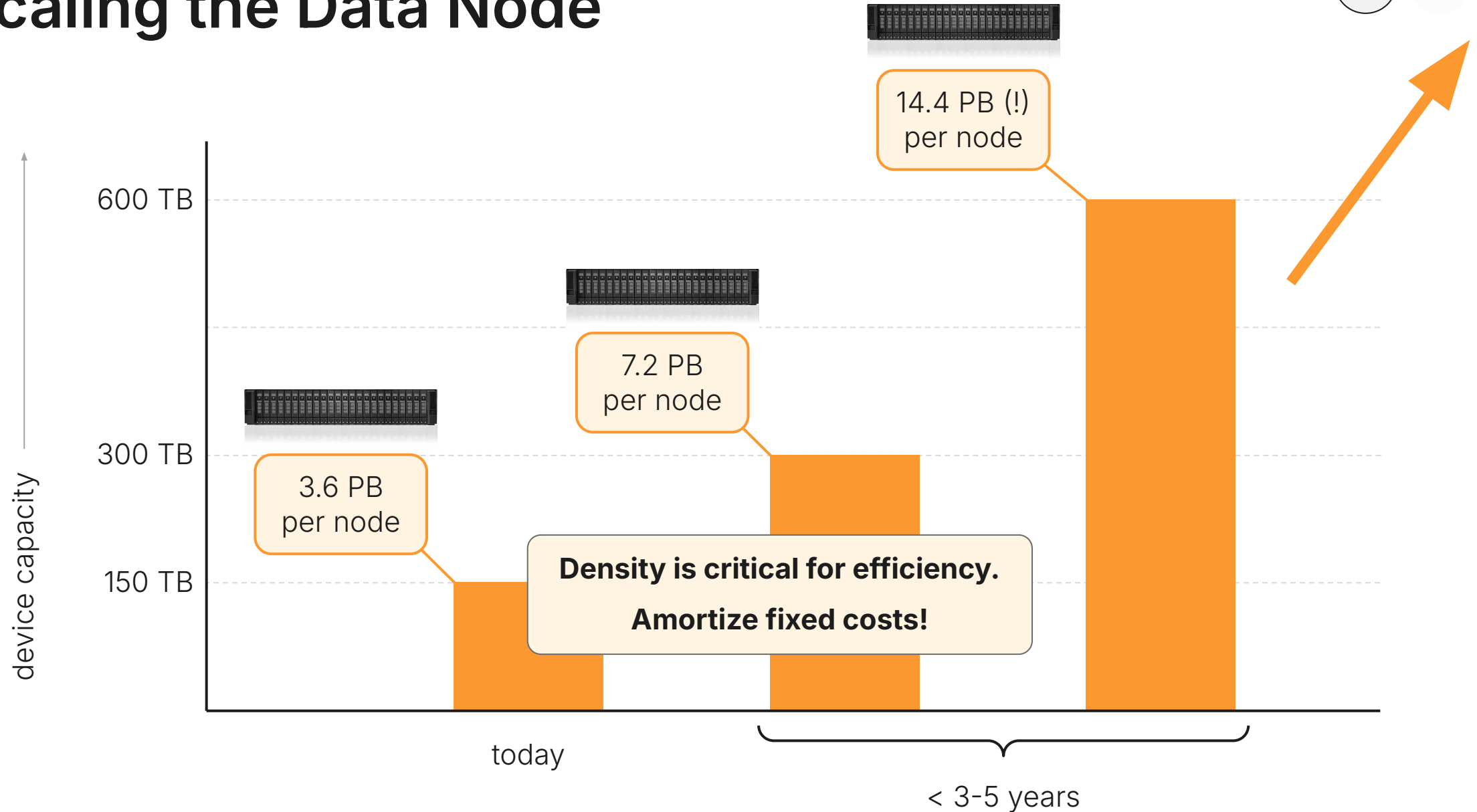


Reliability

# Scaling the Data Node

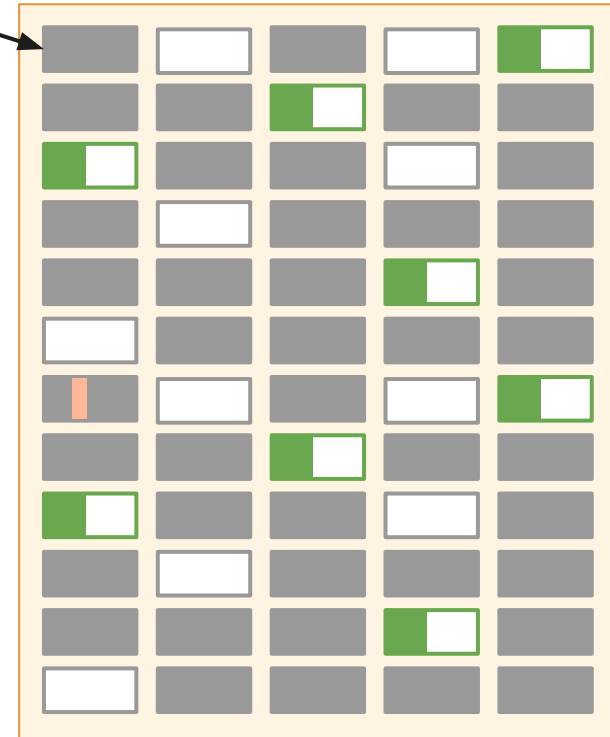


# Scaling the Data Node



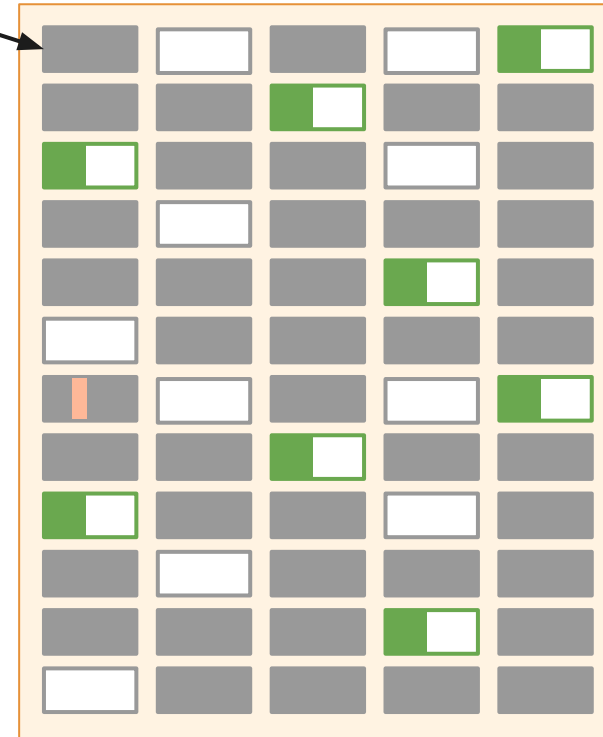
# Review: FTL




flash block



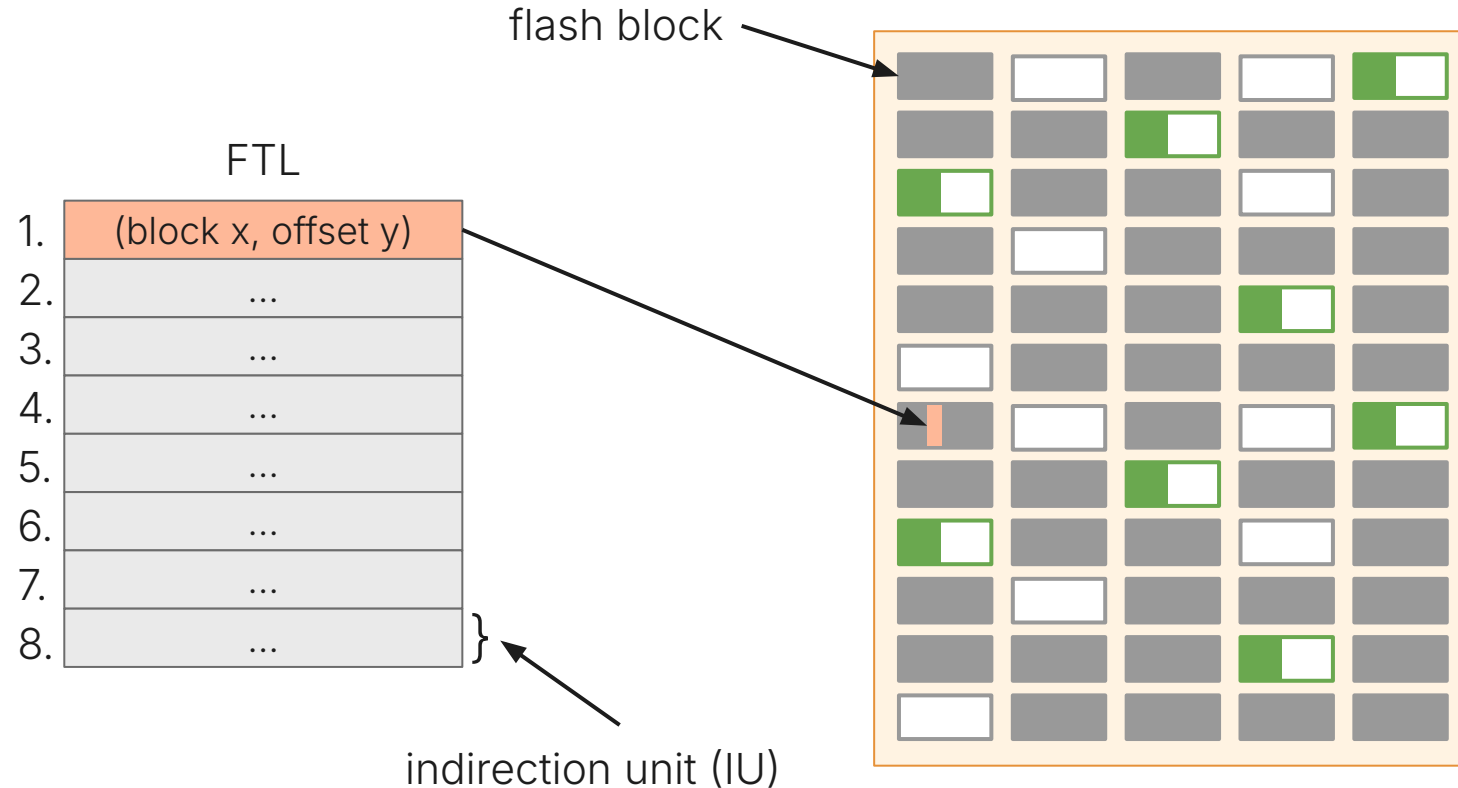
# Review: FTL




flash block



-  **Open:** accepting new writes via append
-  **Full:** no new writes, reads are ok
-  **Erased:** empty, ready to write

# Review: FTL

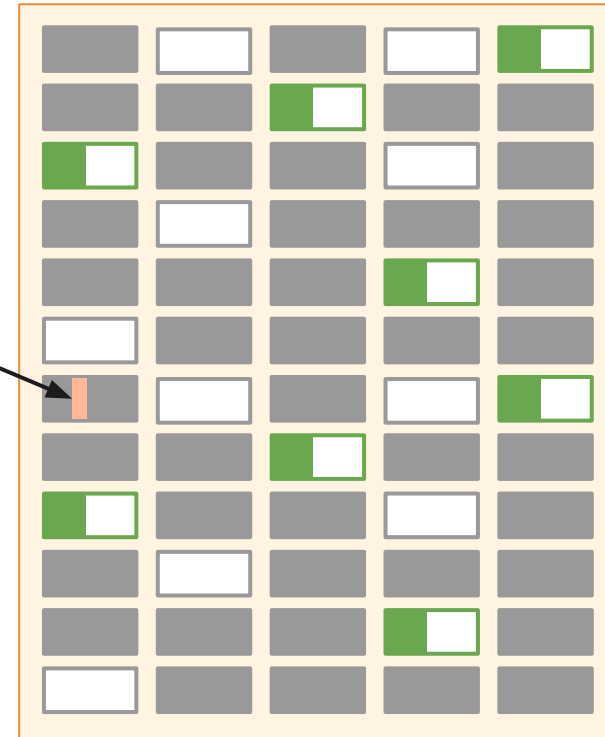
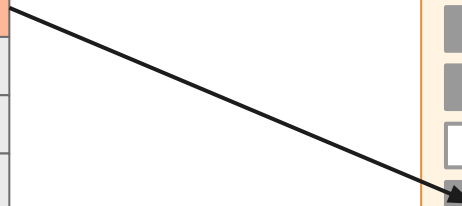





-  **Open:** accepting new writes via append
-  **Full:** no new writes, reads are ok
-  **Erased:** empty, ready to write

# Review: FTL

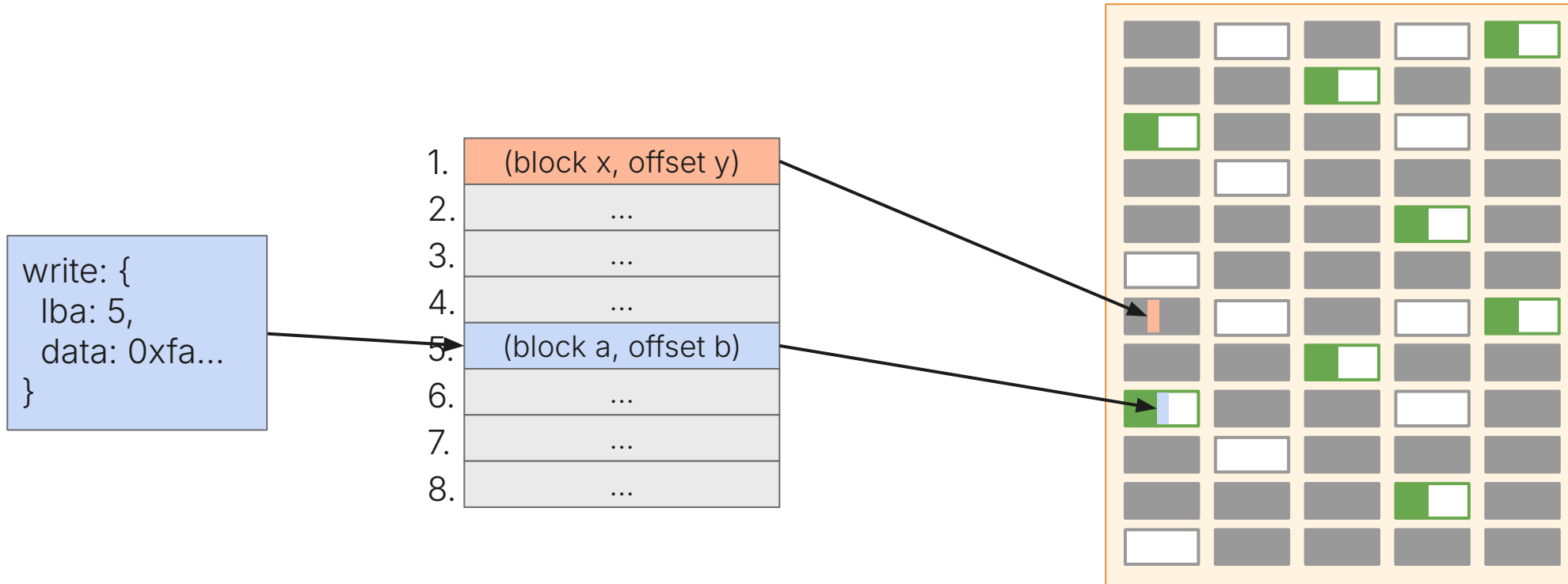
```
write: {  
  lba: 5,  
  data: 0xfa...  
}
```




1.	(block x, offset y)
2.	...
3.	...
4.	...
5.	...
6.	...
7.	...
8.	...



-  **Open:** accepting new writes via append
-  **Full:** no new writes, reads are ok
-  **Erased:** empty, ready to write

# Review: FTL



-  **Open:** accepting new writes via append
-  **Full:** no new writes, reads are ok
-  **Erased:** empty, ready to write

# Scaling the FTL

4B per IU mapping: addresses  $2^{32}$  IUs, aligns nicely in memory

<b>IU size</b>	<b>flash:memory ratio</b>	<b>total memory</b> <i>per 14 PB node</i>	<b>addressable capacity</b> <i>per FTL partition</i>	<b># FTL partitions</b>
4 KiB	1024:1	~14 TB	16 TiB	~800

# Scaling the FTL

4B per IU mapping: addresses  $2^{32}$  IUs, aligns nicely in memory

<b>IU size</b>	<b>flash:memory ratio</b>	<b>total memory</b> <i>per 14 PB node</i>	<b>addressable capacity</b> <i>per FTL partition</i>	<b># FTL partitions</b>
4 KiB	1024:1	~14 TB	16 TiB	~800

Expensive!

- Reduce mapping bits? More partitions!
- Page FTL into flash? Tail latency!

Inflexible!

- Increase mapping bits? More memory!

# Scaling the FTL

4B per IU mapping: addresses  $2^{32}$  IUs, aligns nicely in memory

IU size	flash:memory ratio	total memory <i>per 14 PB node</i>	addressable capacity <i>per FTL partition</i>	# FTL partitions
4 KiB	1024:1	~14 TB	16 TiB	~800

Expensive!

- Reduce mapping bits? More partitions!
- Page FTL into flash? Tail latency!

Inflexible!

- Increase mapping bits? More memory!

**Increase IU size!**

# Scaling the FTL

	<b>IU size</b>	<b>flash:memory ratio</b>	<b>total memory</b> <i>per 14 PB node</i>	<b>addressable capacity</b> <i>per FTL partition</i>	<b># FTL partitions</b>
status quo {	4 KiB	1024:1	~14 TB	16 TiB	~800

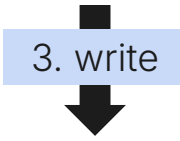
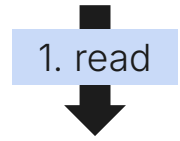
# Scaling the FTL

	<b>IU size</b>	<b>flash:memory ratio</b>	<b>total memory</b> <i>per 14 PB node</i>	<b>addressable capacity</b> <i>per FTL partition</i>	<b># FTL partitions</b>
status quo	4 KiB	1024:1	~14 TB	16 TiB	~800
new high-cap SSDs	16 KiB	4096:1	~3.4 TB	64 TiB	~200
	64 KiB	16384:1	~850 GB	256 TiB	~50

# Scaling the FTL

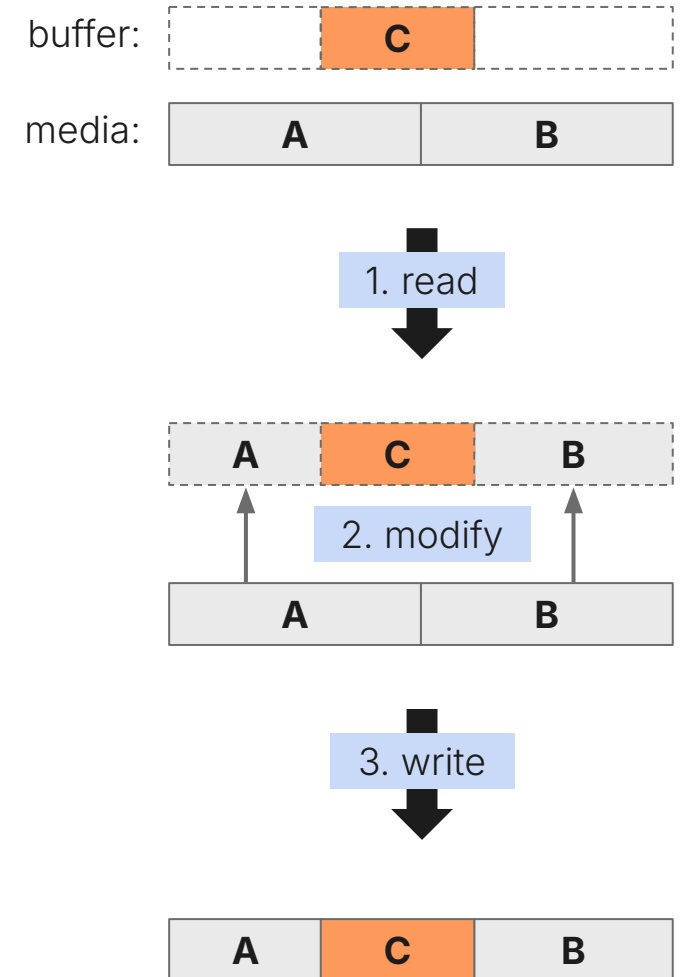
	<b>IU size</b>	<b>flash:memory ratio</b>	<b>total memory</b> <i>per 14 PB node</i>	<b>addressable capacity</b> <i>per FTL partition</i>	<b># FTL partitions</b>
status quo	4 KiB	1024:1	~14 TB	16 TiB	~800
new high-cap SSDs	16 KiB	4096:1	~3.4 TB	64 TiB	~200
	64 KiB	16384:1	~850 GB	256 TiB	~50
the future of ultra-dense QLC	512 KiB	131072:1	~100 GB	2 PiB	~6
	...				

# Pitfall: Read-Modify-Write



# Pitfall: Read-Modify-Write

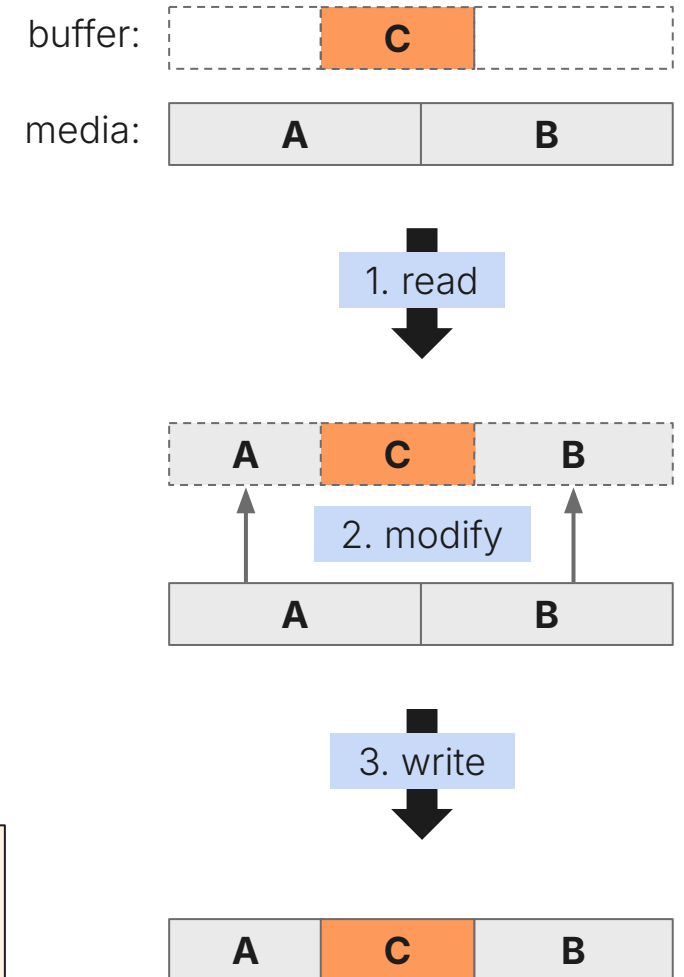
- extra read and write load → **reduced perf, endurance**
- extra media round-trip → **increased latency**
- write size **and alignment** are important
- overhead **scales with IU size**
- this is a real problem for **general-purpose** systems



# Pitfall: Read-Modify-Write

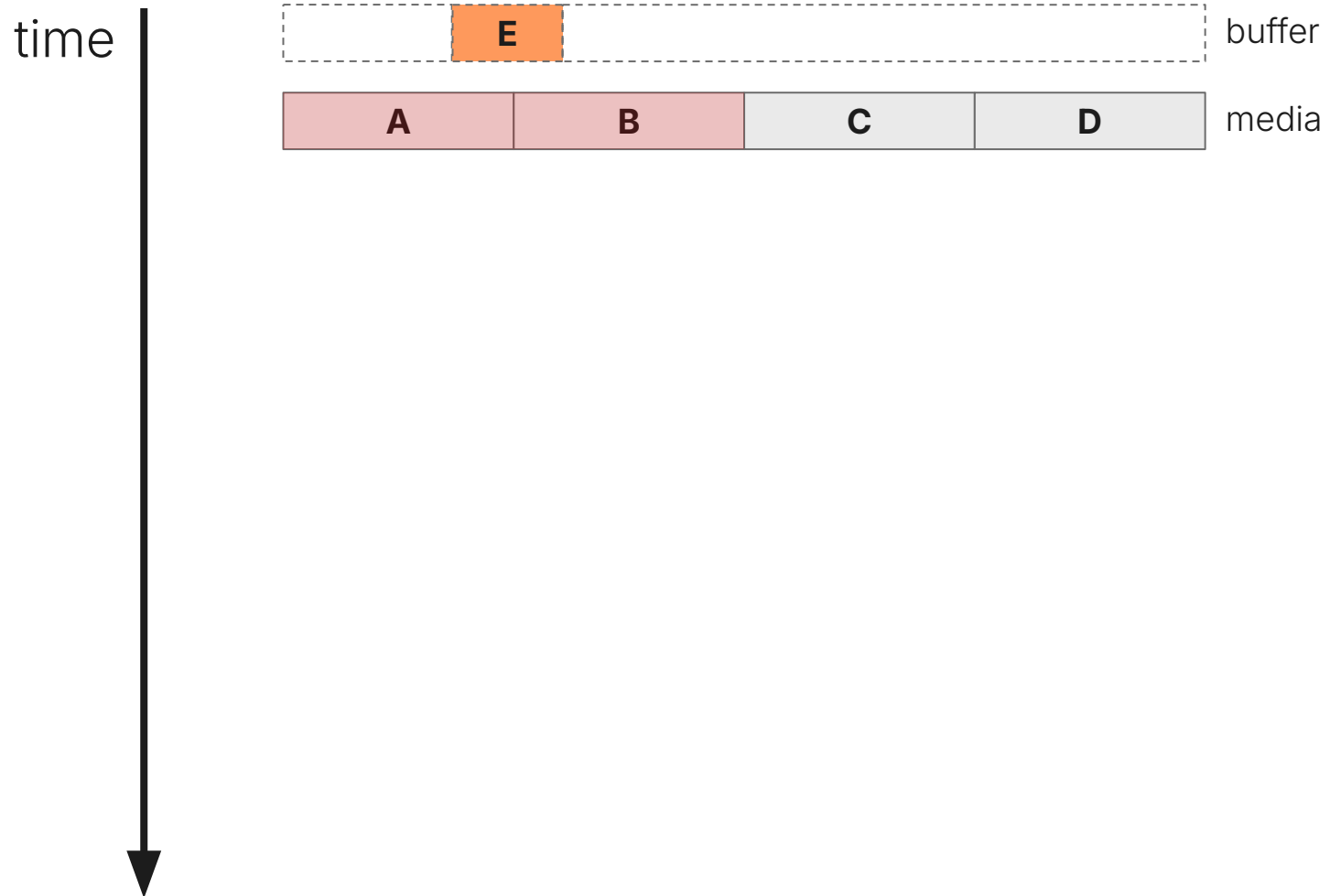


- extra read and write load → **reduced perf, endurance**
- extra media round-trip → **increased latency**
- write size **and alignment** are important
- overhead **scales with IU size**
- this is a real problem for **general-purpose** systems

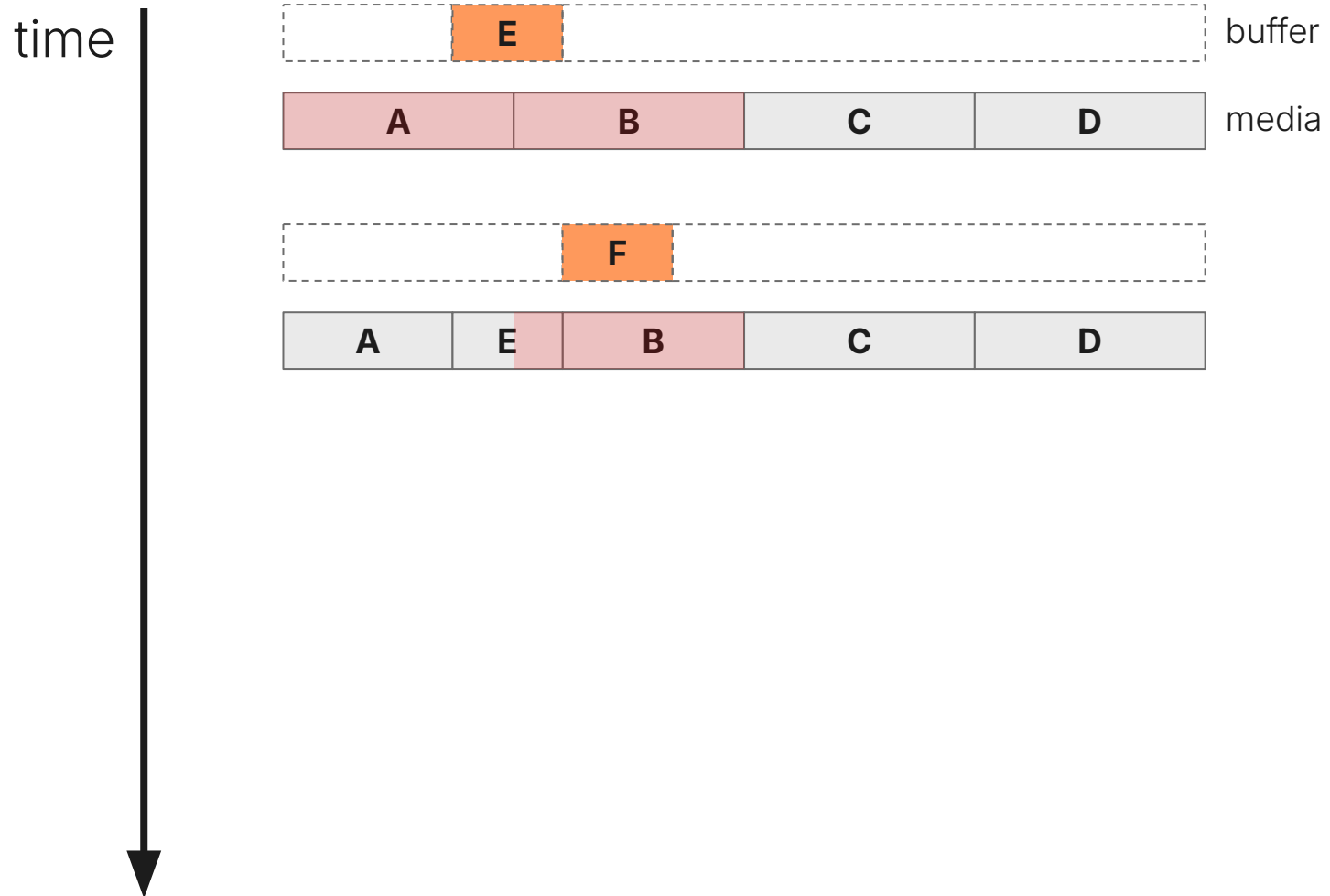


**Key Insight:** Take advantage of exascale workload patterns to deliver a large IU with minimal RMW.

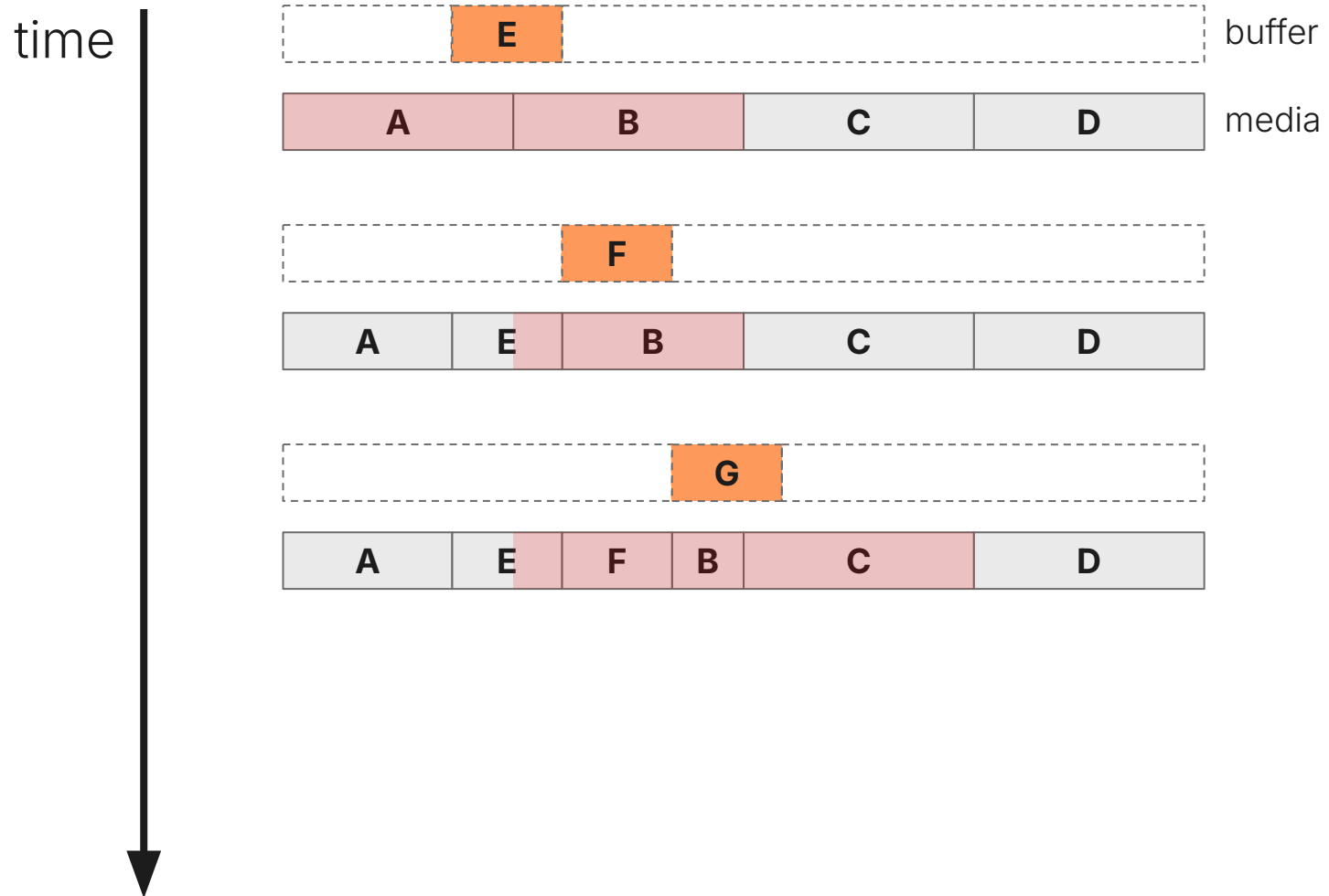
# Write Coalescing



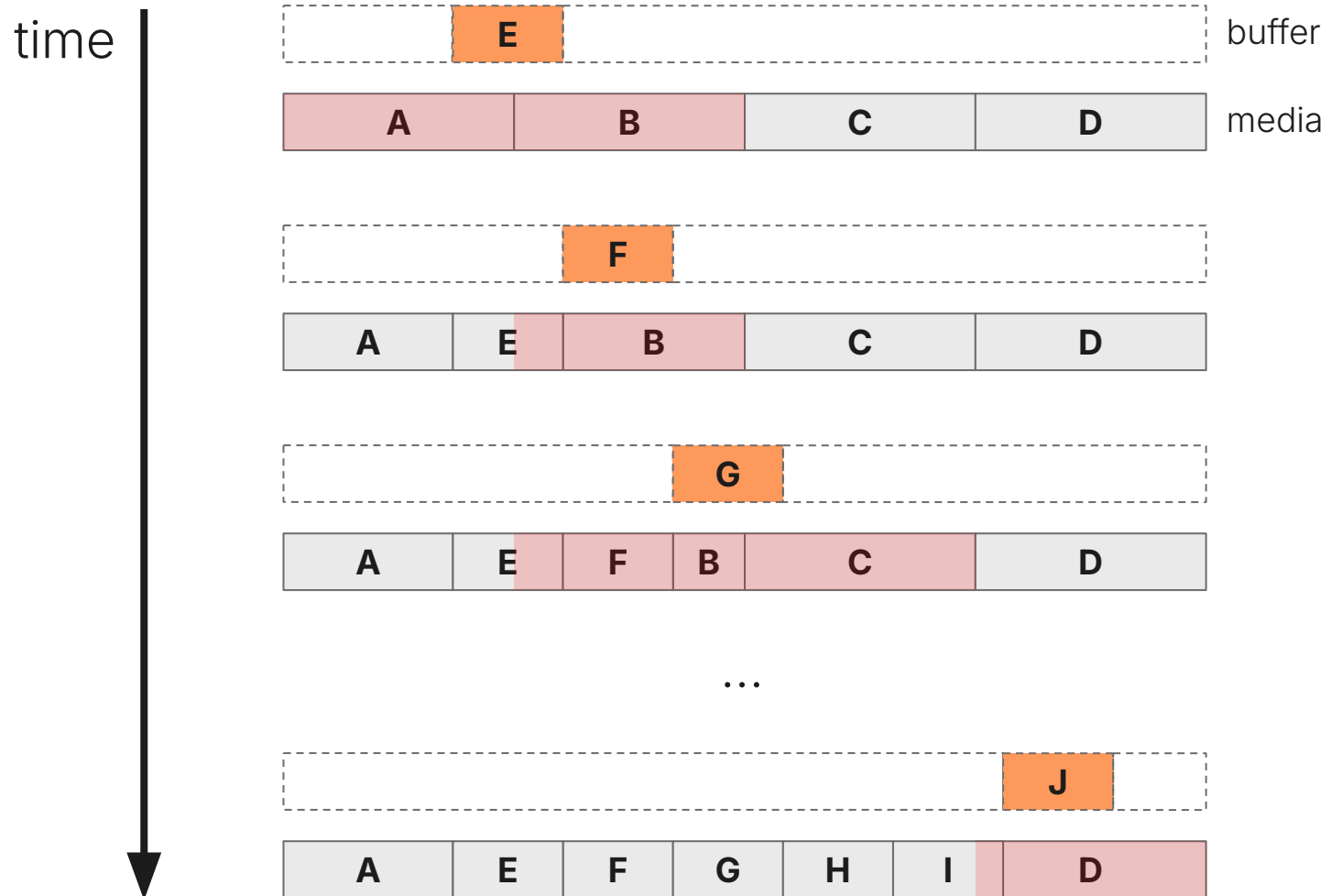
# Write Coalescing



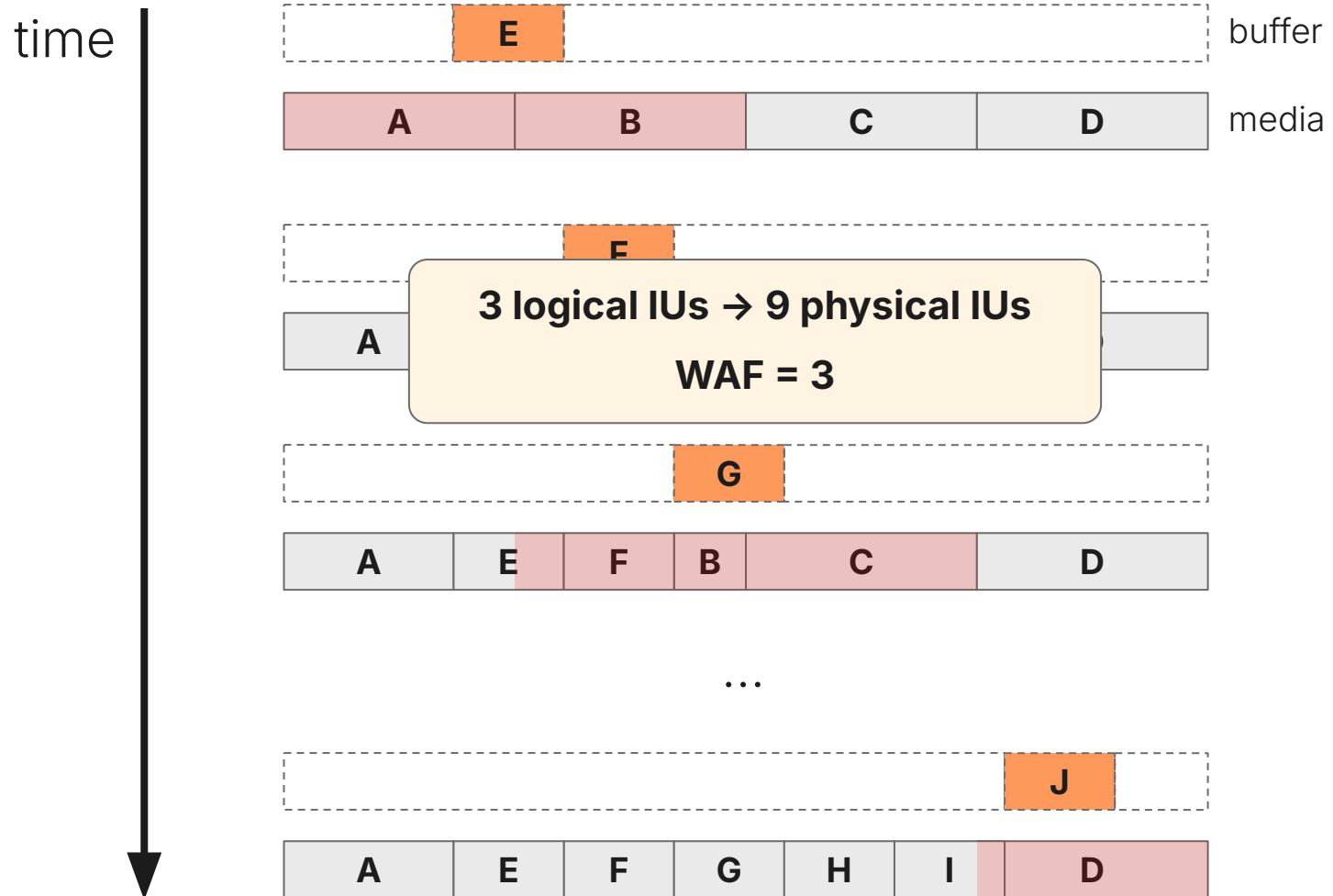
# Write Coalescing



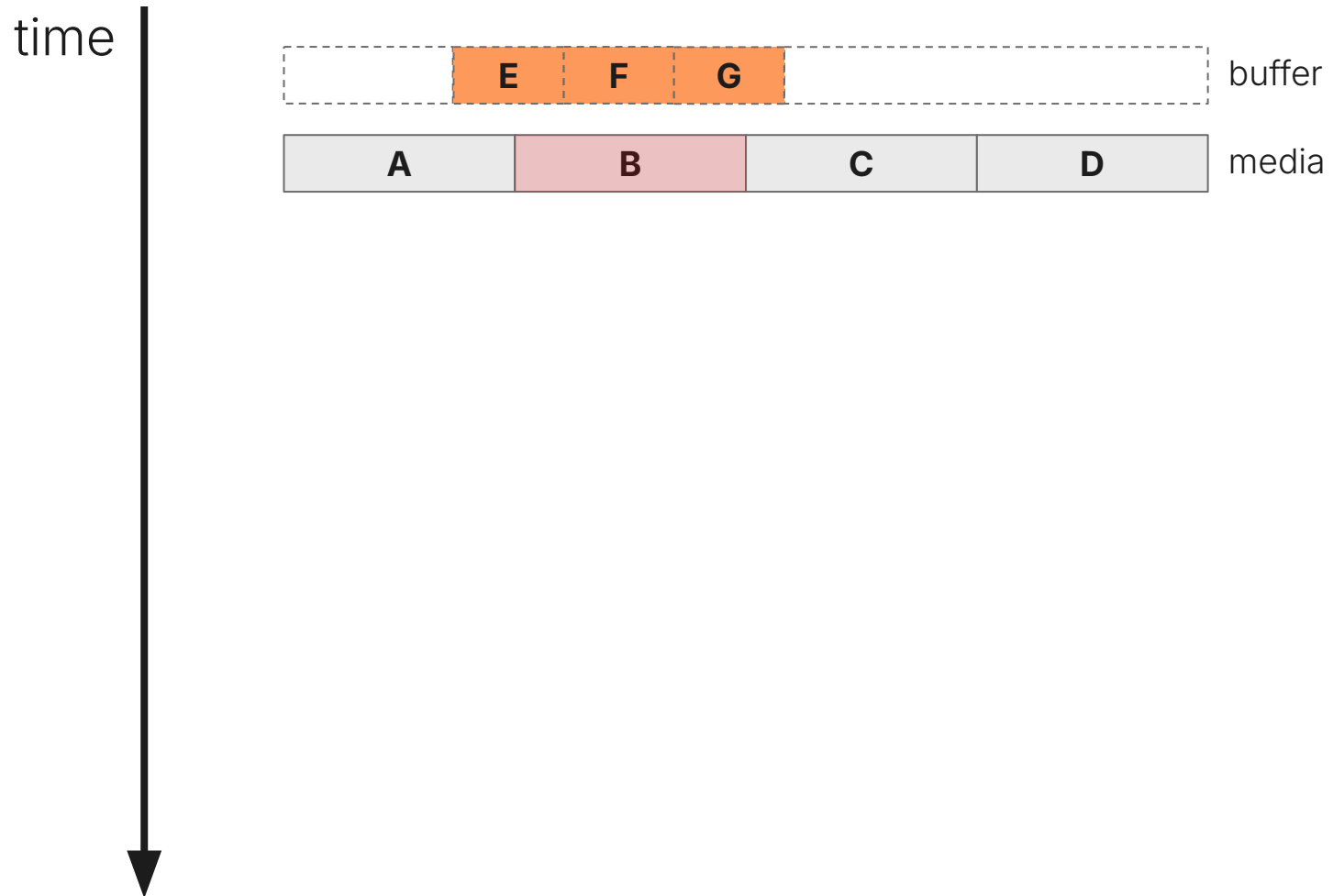
# Write Coalescing



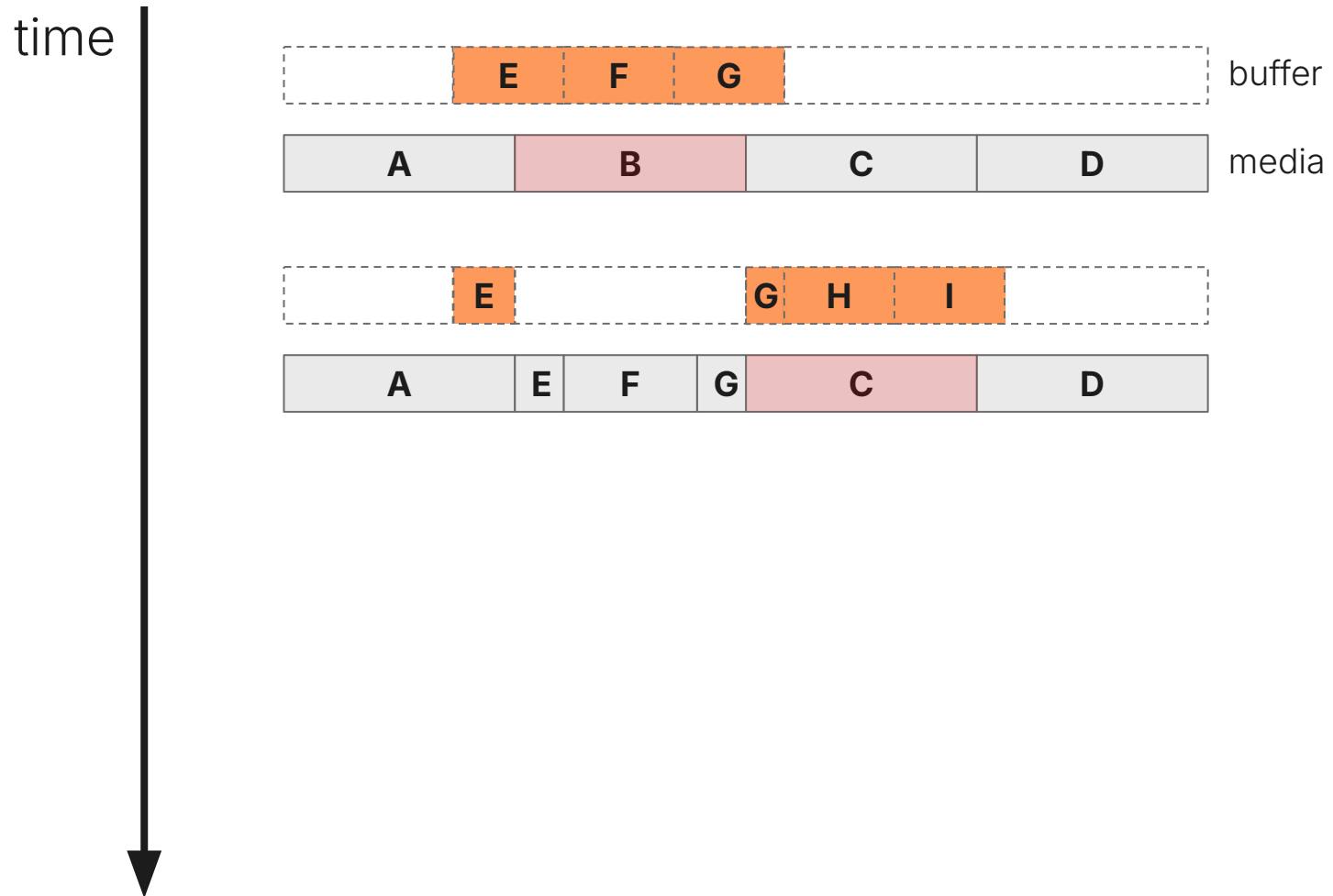
# Write Coalescing



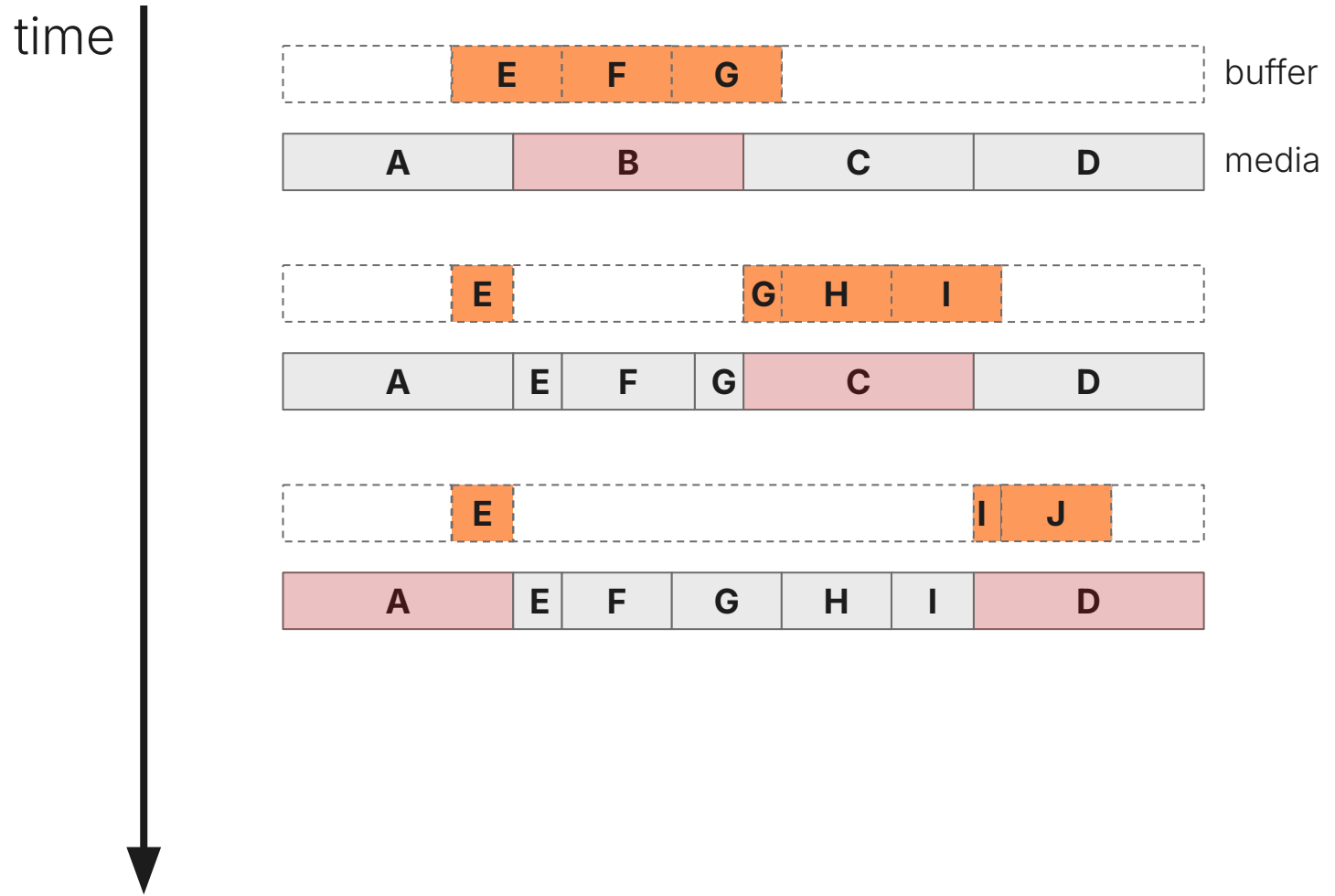
# Write Coalescing



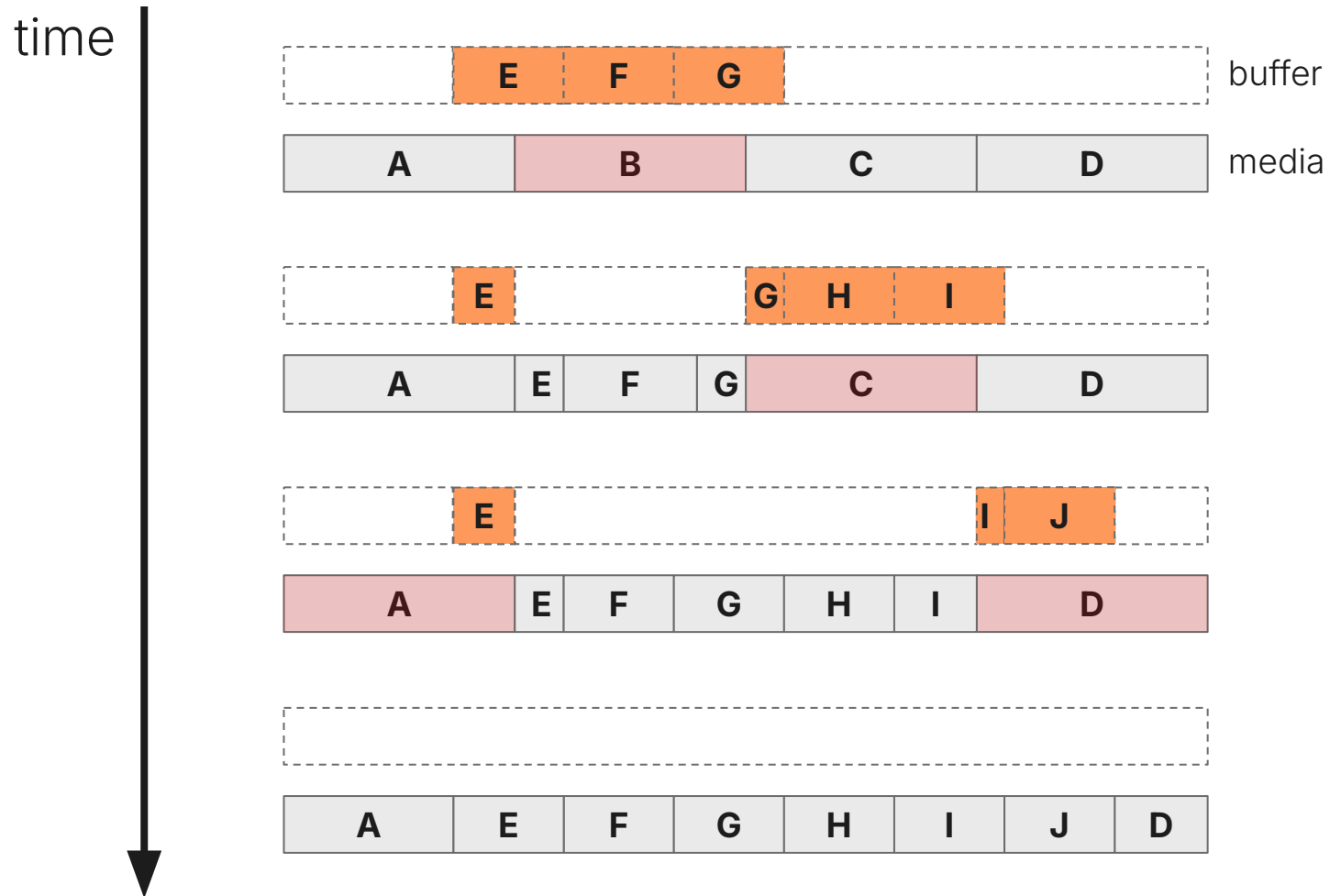
# Write Coalescing



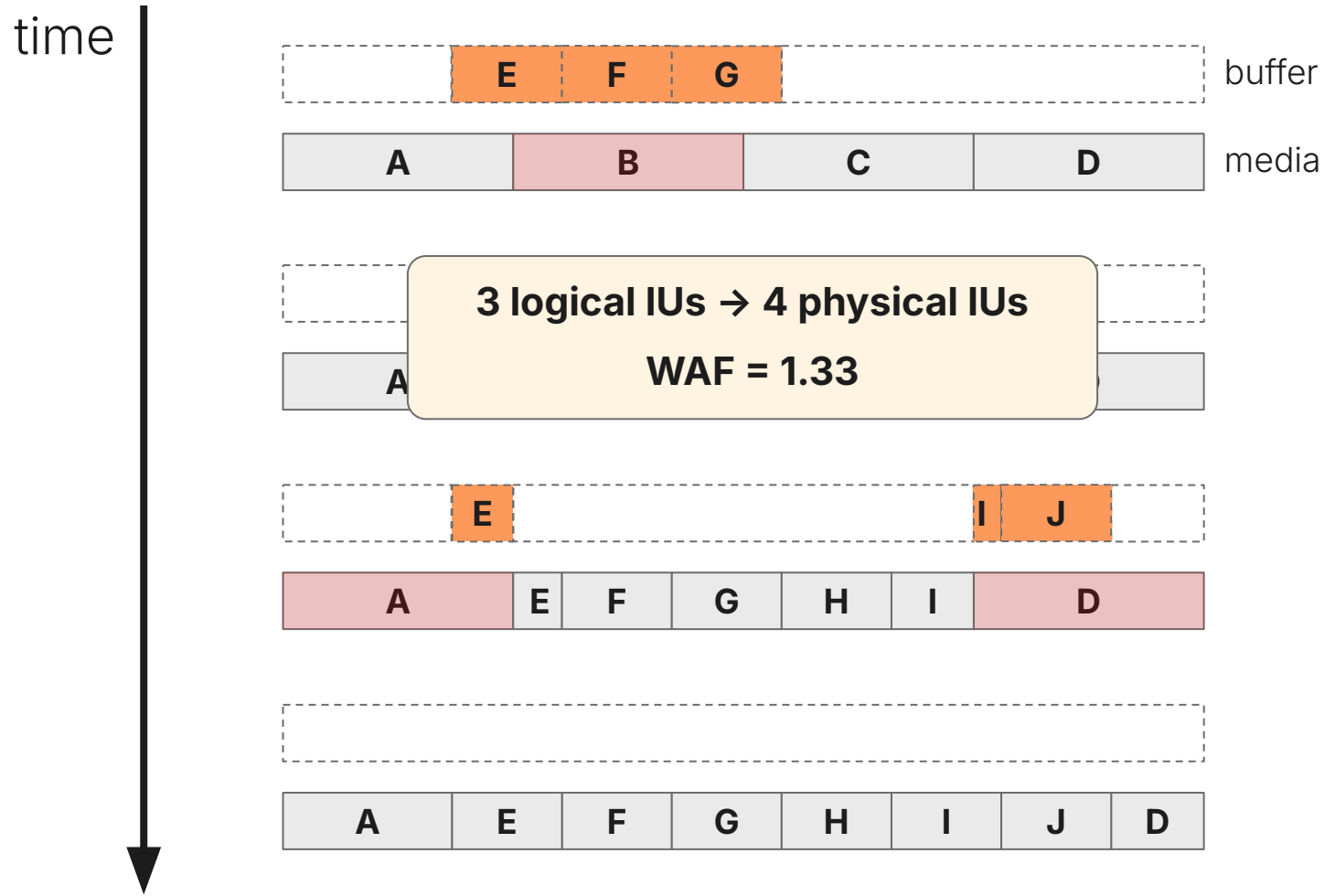
# Write Coalescing



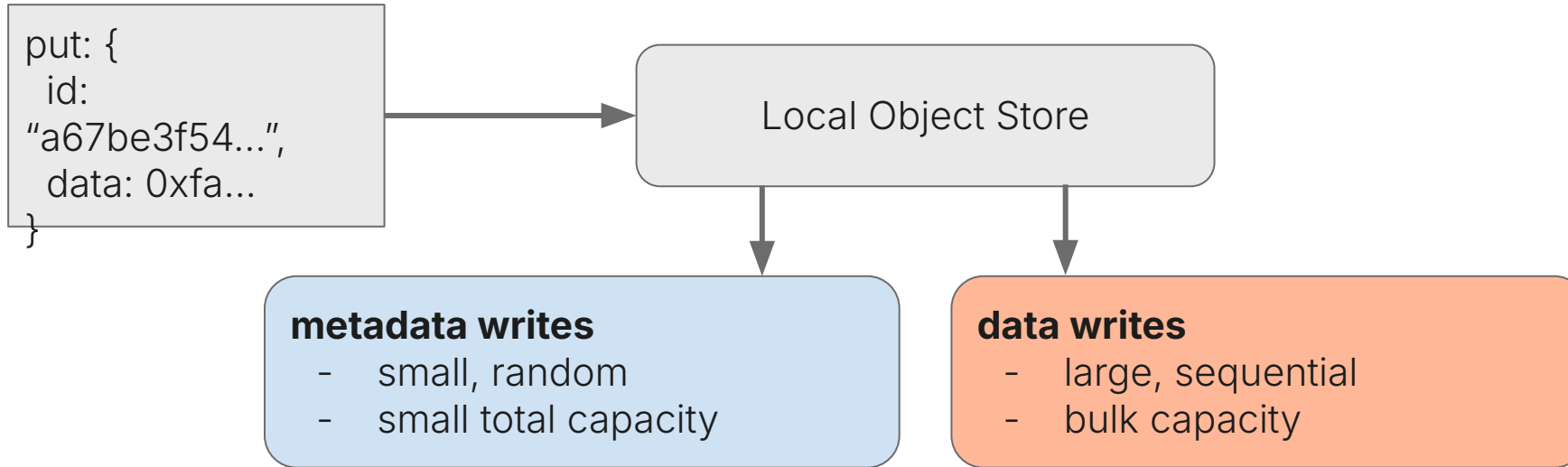
# Write Coalescing



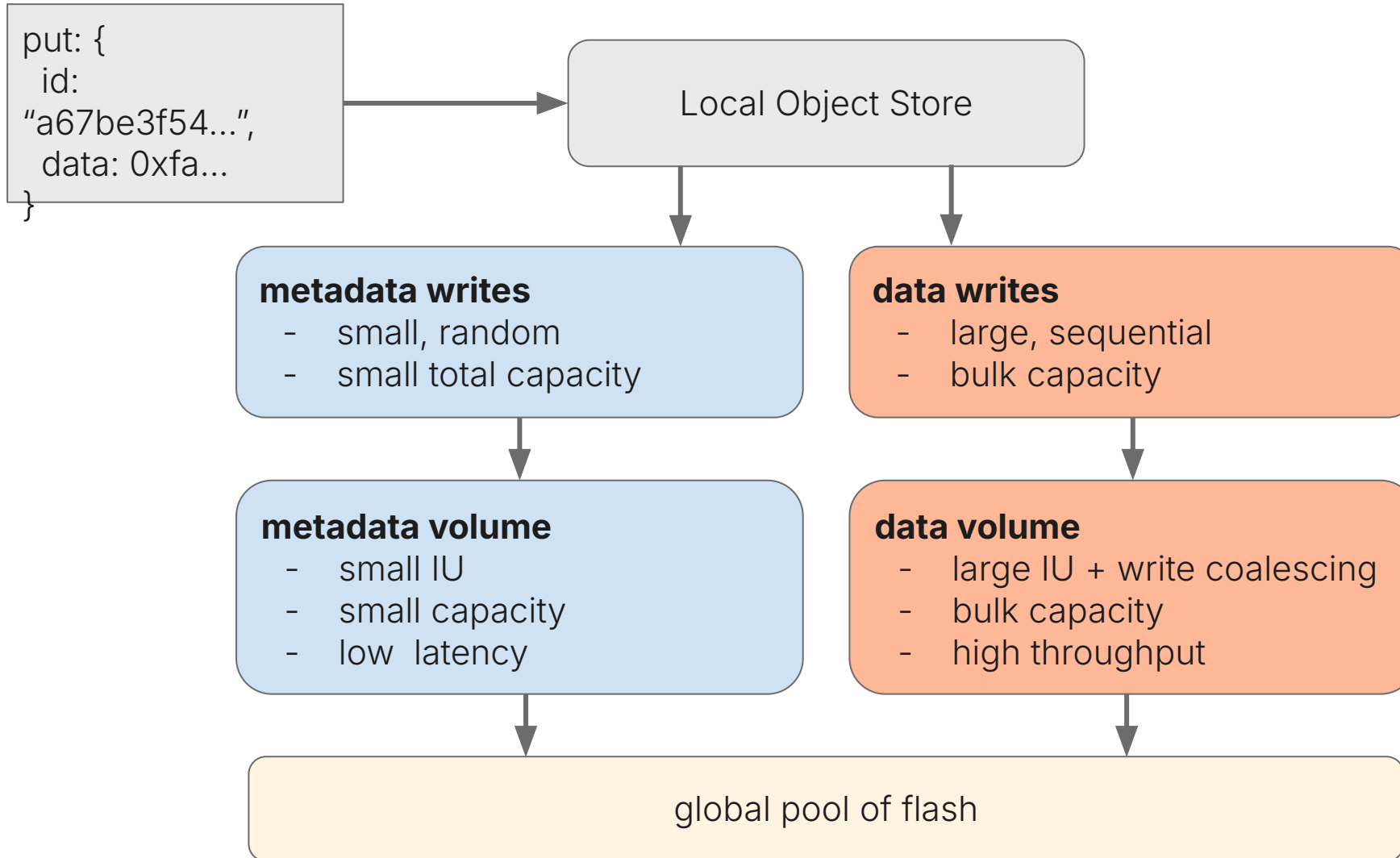
# Write Coalescing



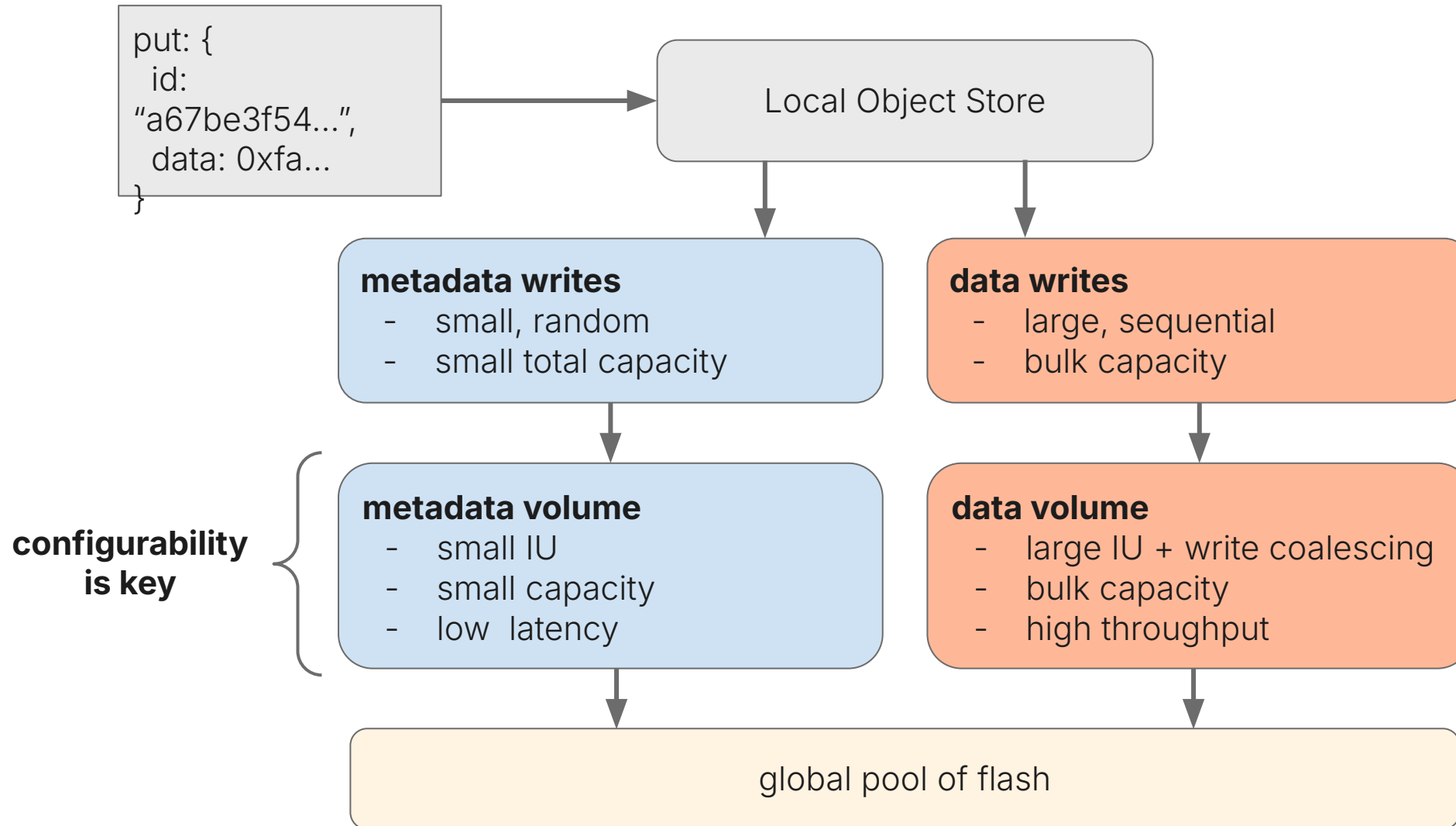
# Configurable IUs



# Configurable IUs



# Configurable IUs



# QLC Challenges



Scale



Endurance



Reliability

# Review: Garbage Collection

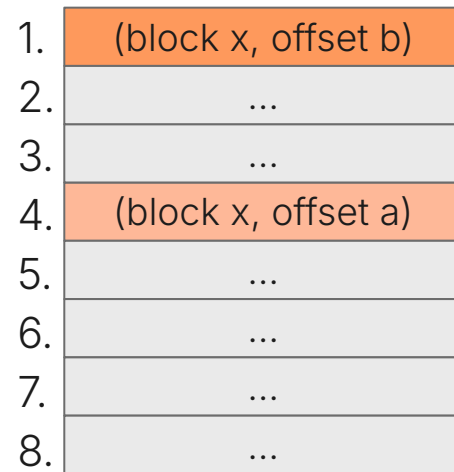


# Review: Garbage Collection

Endurance

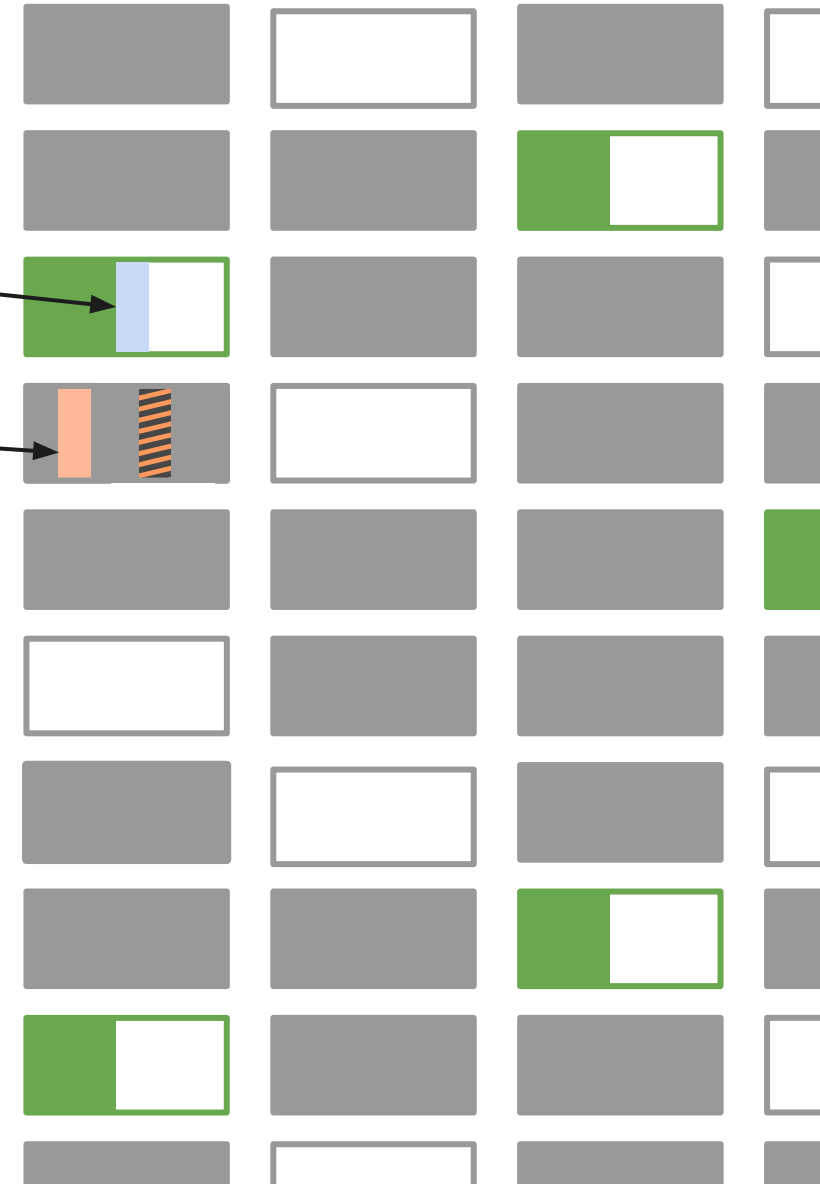
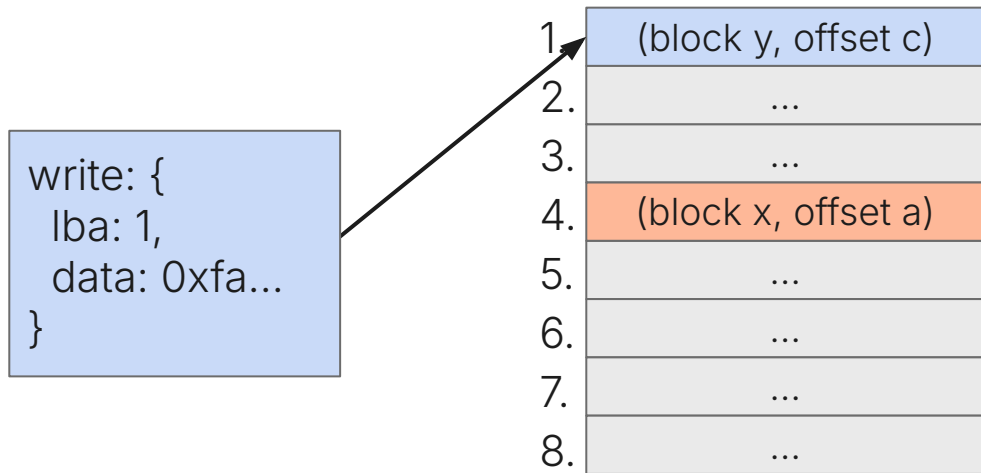


```
write: {  
  lba: 1,  
  data: 0xfa...  
}
```

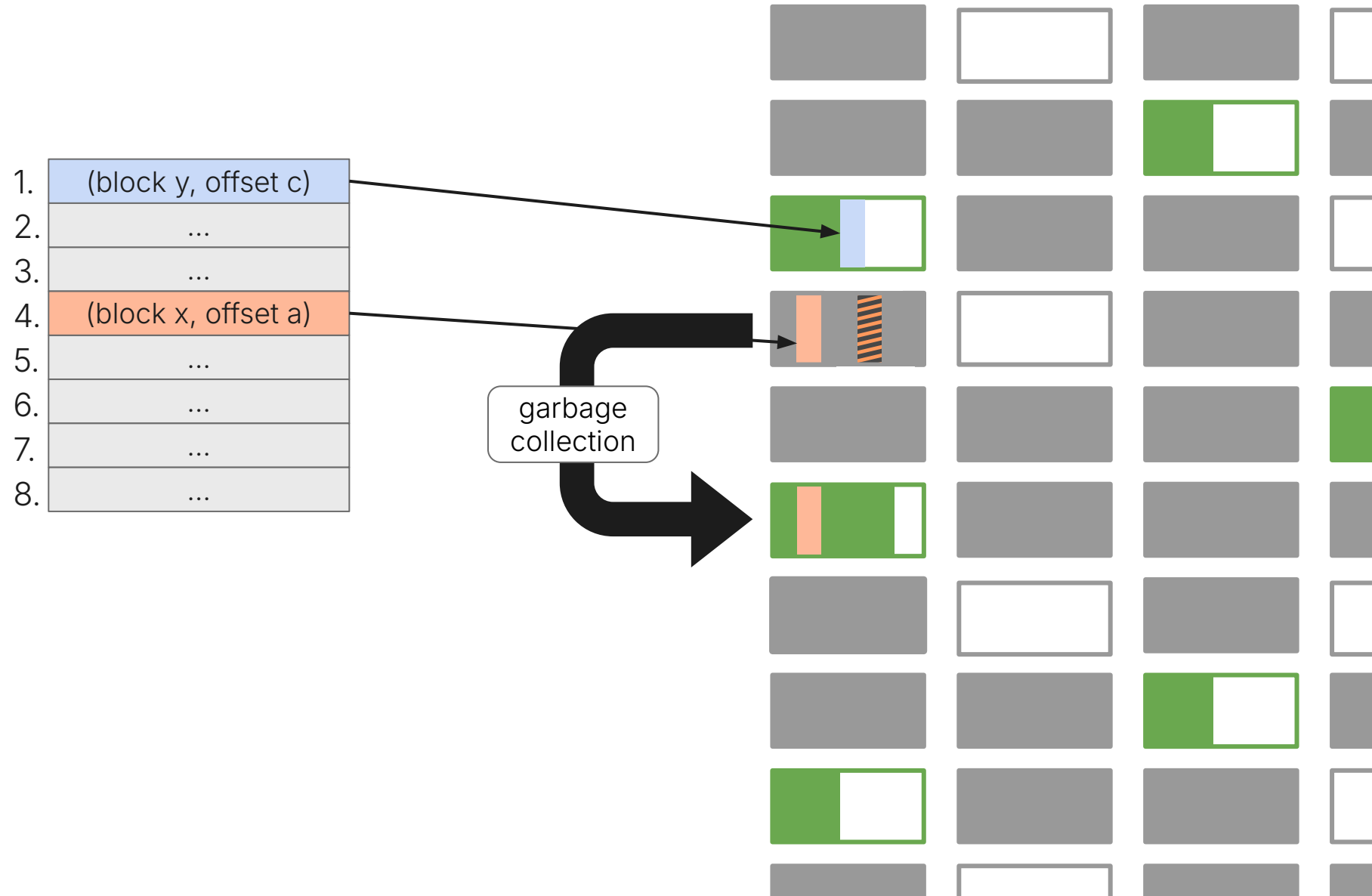


# Review: Garbage Collection

Endurance



# Review: Garbage Collection



# Review: Garbage Collection

- 1. (block y, offset c)
- 2. ...
- 3. ...
- 4. (block z, offset d)
- 5. ...
- 6. ...
- 7. ...
- 8. ...



# Write Amplification

Endurance



- QLC has relatively **low endurance**: O(thousands)
- WAF reduces media lifetime (and performance!) linearly  
*e.g. WAF=2 → ½ total lifetime bytes written by user*
- GC is the **dominant contributor** to WAF (others: RMW, retention)
- Option: increase over-provisioning? → direct impact to cost!
- Option: use smarter GC algorithms? → no help to random-write benchmarks!

$$\text{WAF} = \frac{\text{physical written}}{\text{logical written}}$$

# Write Amplification

Endurance



- QLC has relatively **low endurance**: O(thousands)
- WAF reduces media lifetime (and performance!) linearly  
*e.g. WAF=2 → ½ total lifetime bytes written by user*
- GC is the **dominant contributor** to WAF (others: RMW, retention)
- Option: increase over-provisioning? → direct impact to cost!
- Option: use smarter GC algorithms? → no help to random-write benchmarks!

$$\text{WAF} = \frac{\text{physical written}}{\text{logical written}}$$

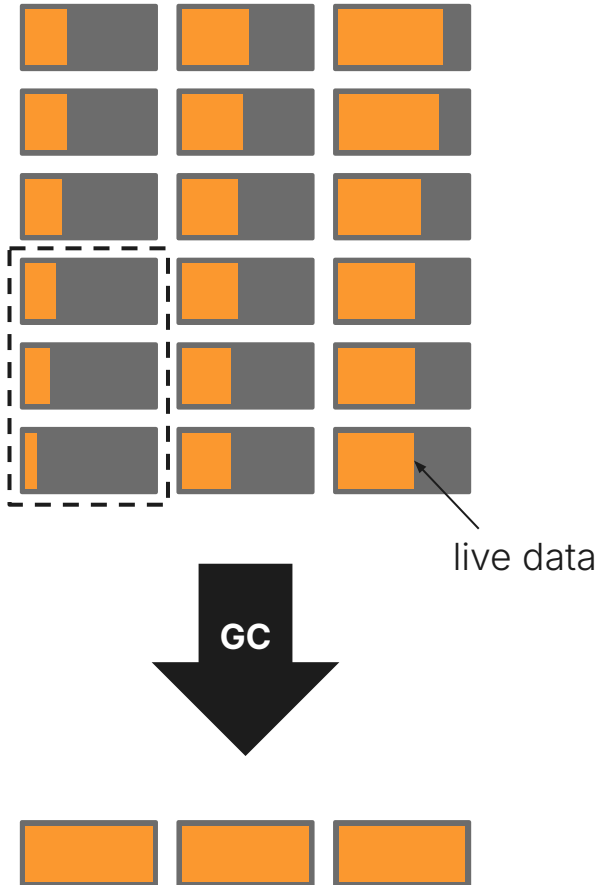
**Key Insight:** Target real-world workloads with smarter algorithms to minimize WAF.

# Smarter Algorithms

Endurance



greedy

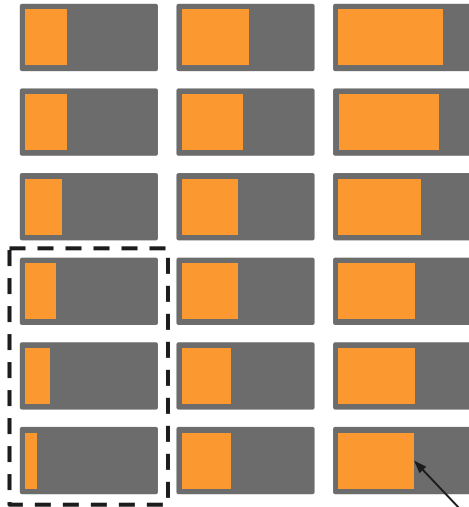


# Smarter Algorithms

Endurance



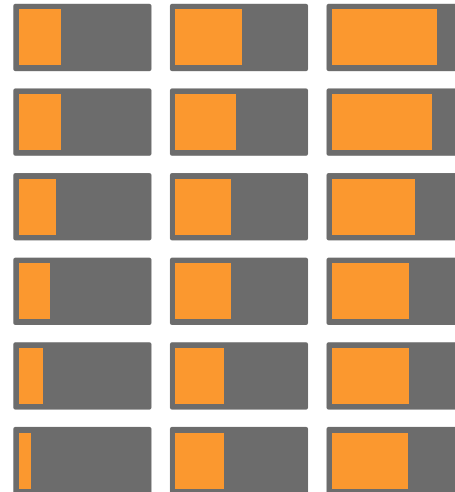
greedy



live data



consider: age

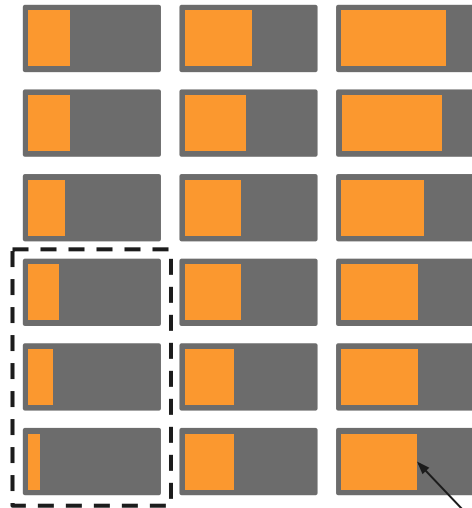


# Smarter Algorithms

Endurance



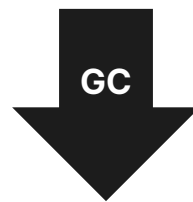
greedy



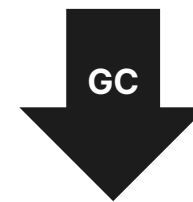
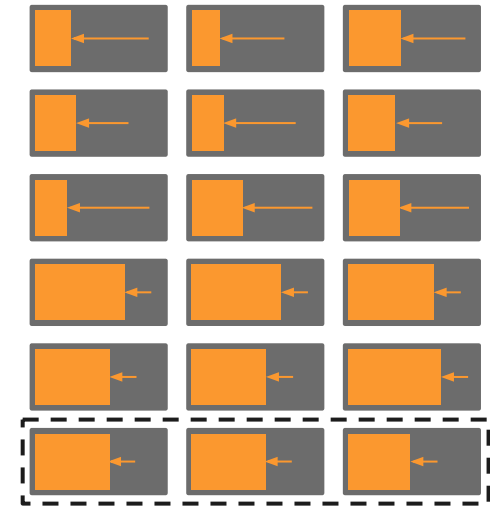
live data



consider: age



consider: stability



# Smarter Algorithms



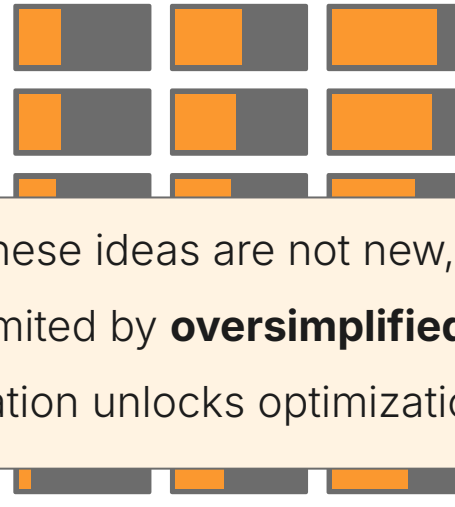
greedy



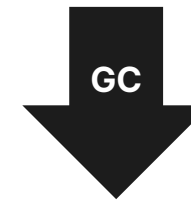
live data



consider: age



consider: stability



These ideas are not new, *but*  
adoption is limited by **oversimplified benchmarks**.  
Workload specialization unlocks optimization without overfitting!



# QLC Challenges



Scale



Endurance



Reliability

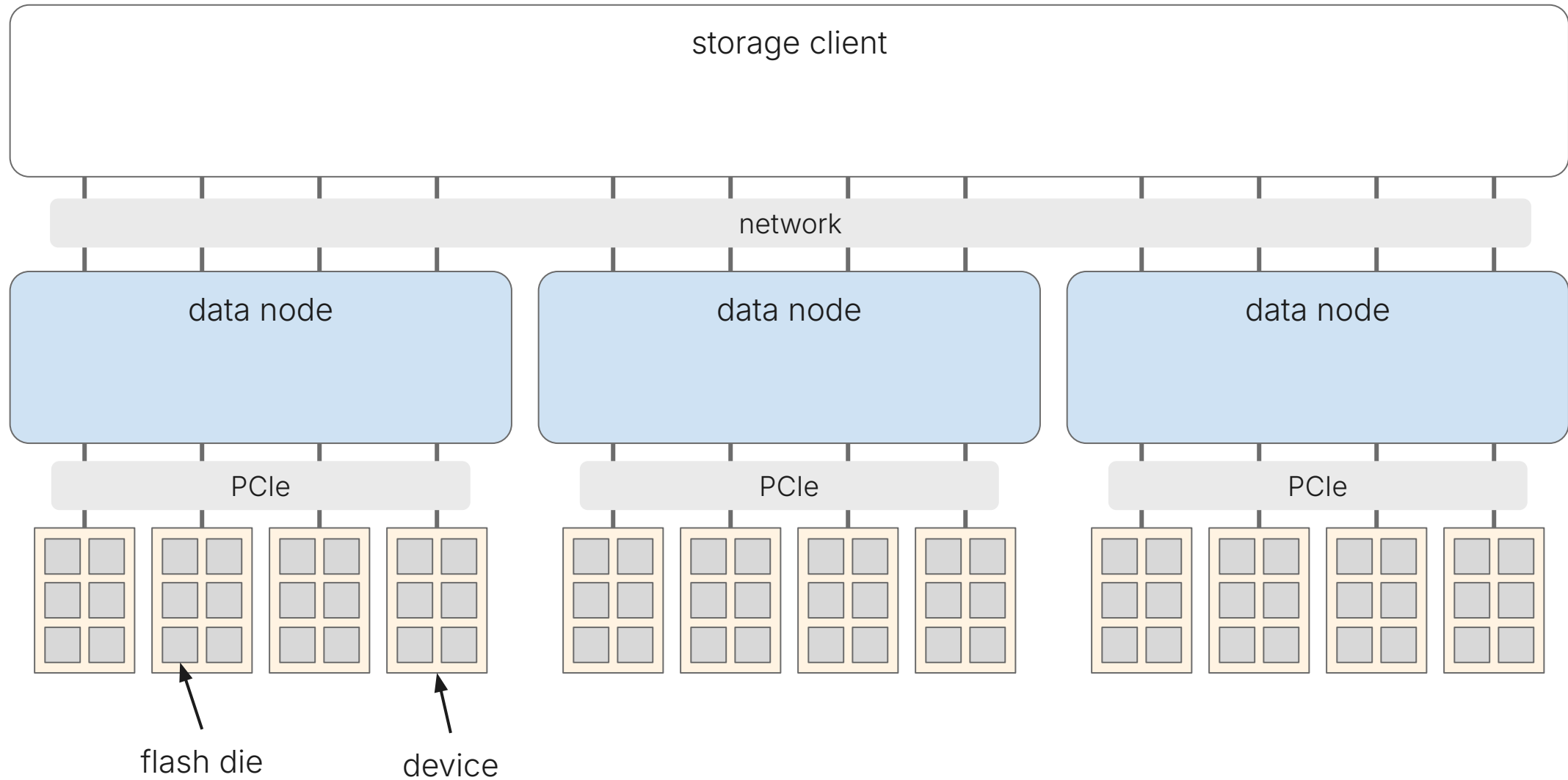
# Blast Radius

Reliability

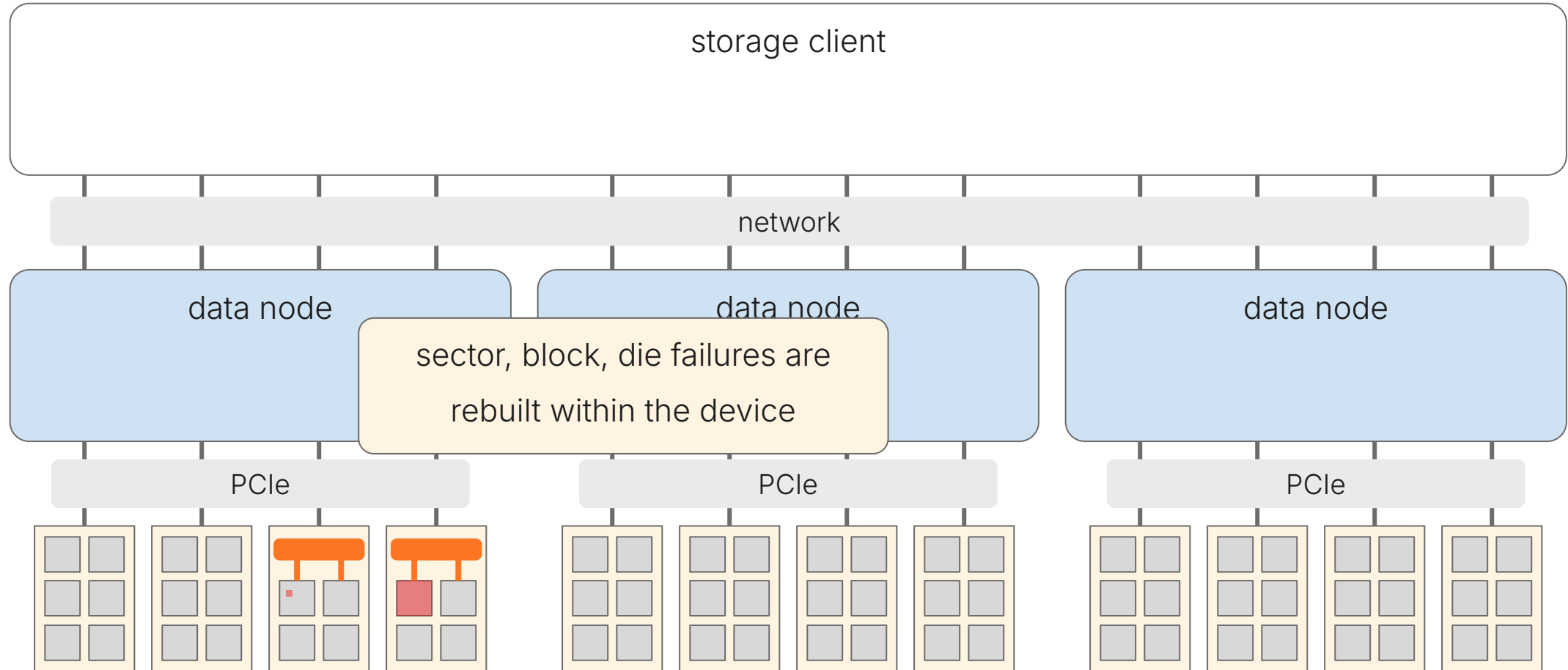


- failure modes:
  - sector ~KB
  - flash block ~100s MB
  - flash die ~100s GB
  - device ~100s TB
  - data node ~PB
  
- QLC density → large devices/nodes → increased impact of failures

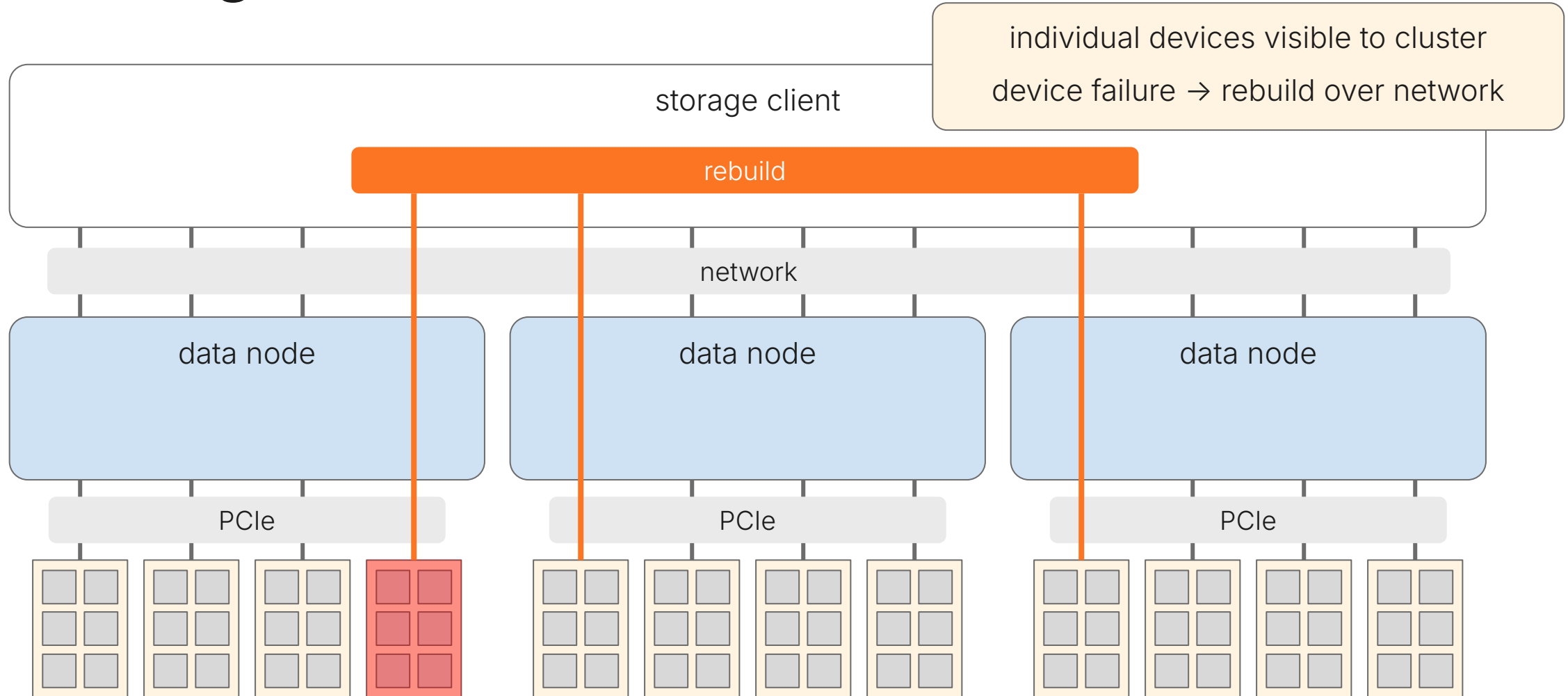
# Rebuilding Failures



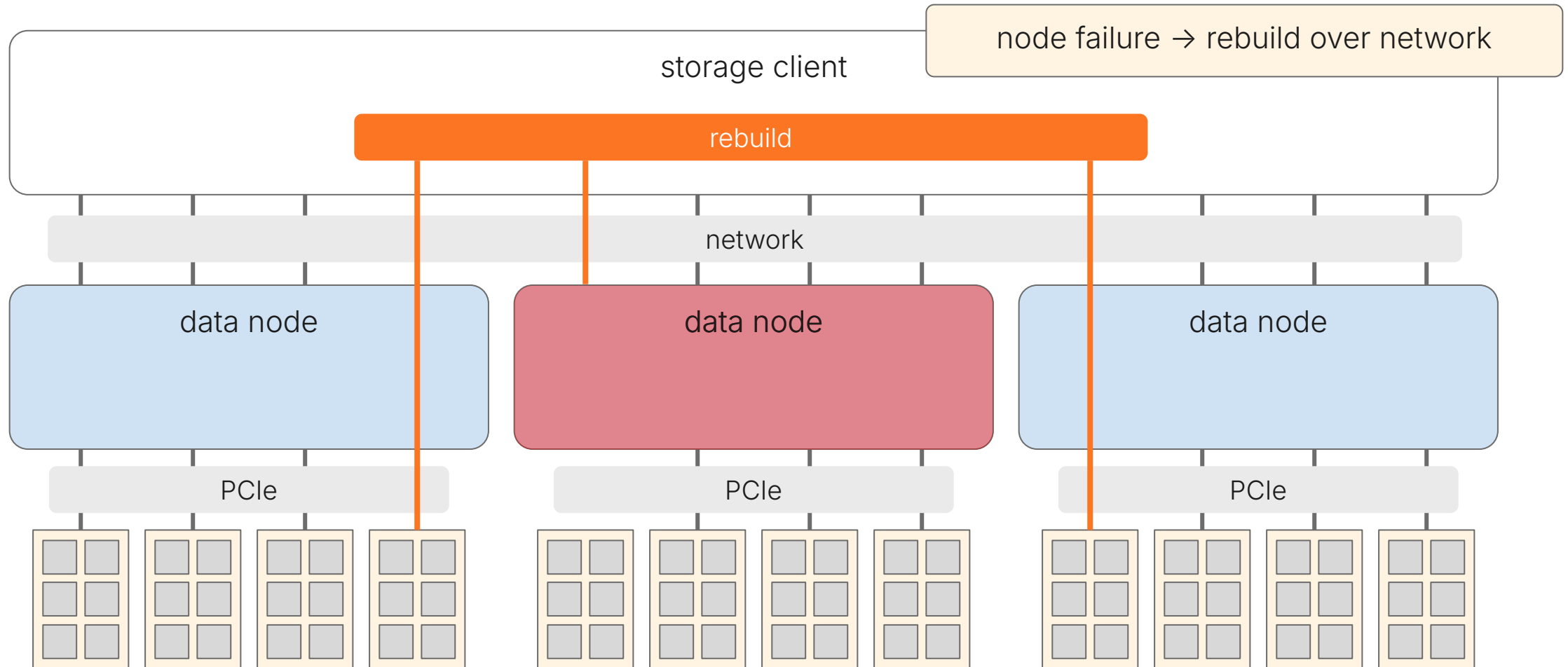
# Rebuilding Failures



# Rebuilding Failures



# Rebuilding Failures



# Blast Radius

Reliability



- failure modes:
  - sector ~KB
  - flash block ~100s MB
  - flash die ~100s GB
  - device ~100s TB
  - data node ~PB
  
- QLC density → large devices/nodes → increased impact of failures

# Blast Radius



- failure modes:
  - sector ~KB
  - flash block ~100s MB ← **most common, small impact**
  - flash die ~100s GB ← **rare, moderate impact**
  - device ~100s TB ← **moderate frequency, huge impact!**
  - data node ~PB
  
- QLC density → large devices/nodes → increased impact of failures

# Blast Radius

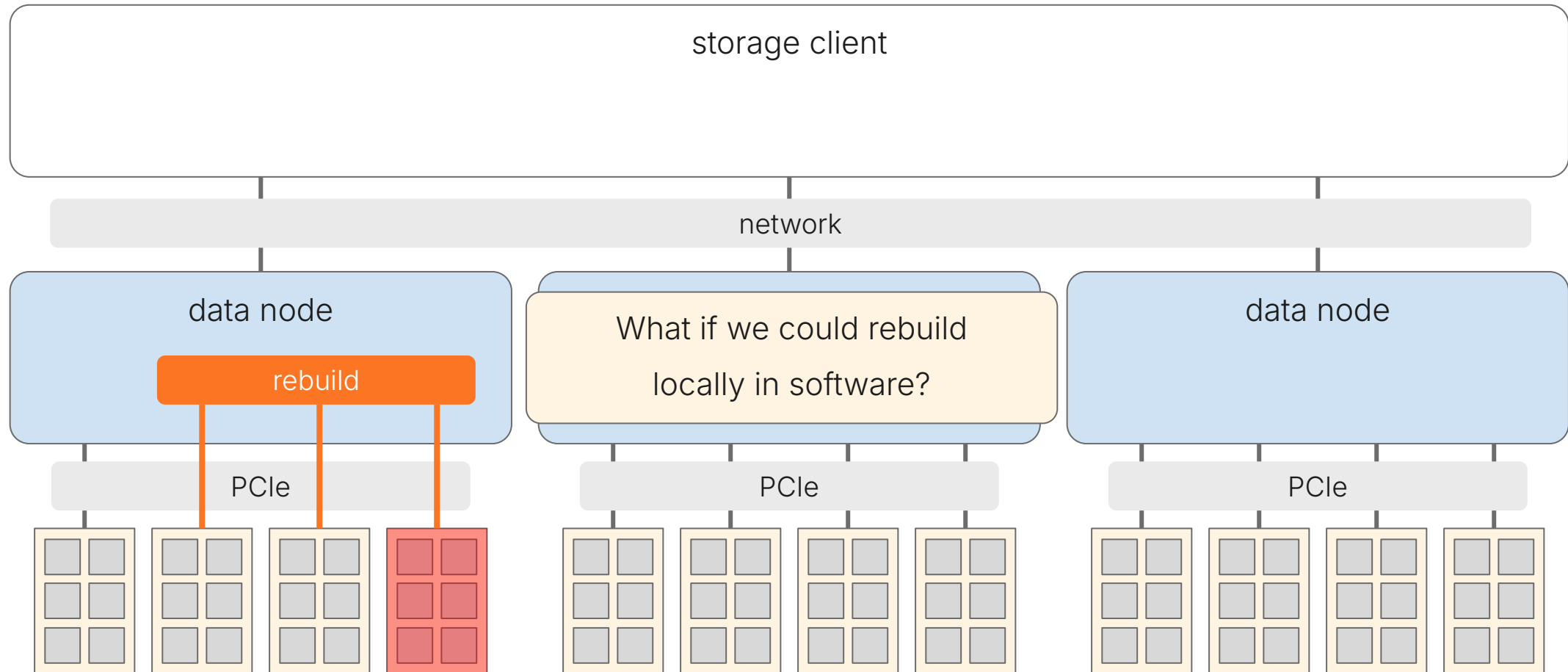


- failure modes:
  - sector ~KB
  - flash block ~100s MB ← **most common, small impact**
  - flash die ~100s GB ← **rare, moderate impact**
  - device ~100s TB ← **moderate frequency, huge impact!**
  - data node ~PB
  
- QLC density → large devices/nodes → increased impact of failures

**Key Insight:** Protect efficiently against dominant failure modes.

# Node-Local Parity?

Reliability



# Recap: QLC Challenges



## Scale

- large IUs
- write coalescing
- data/metadata



## Endurance

- real workloads
- smarter algos
- minimize WAF



## Reliability

- failure modes
- blast radius
- local parity

# Recap: QLC Challenges



## Scale

- large IUs
- write coalescing
- data/metadata



## Endurance

- real workloads
- smarter algos
- minimize WAF



## Reliability

- failure modes
- blast radius
- local parity

**configurability, customization, control**

# Agenda

---

**Context** Exascale Storage

---

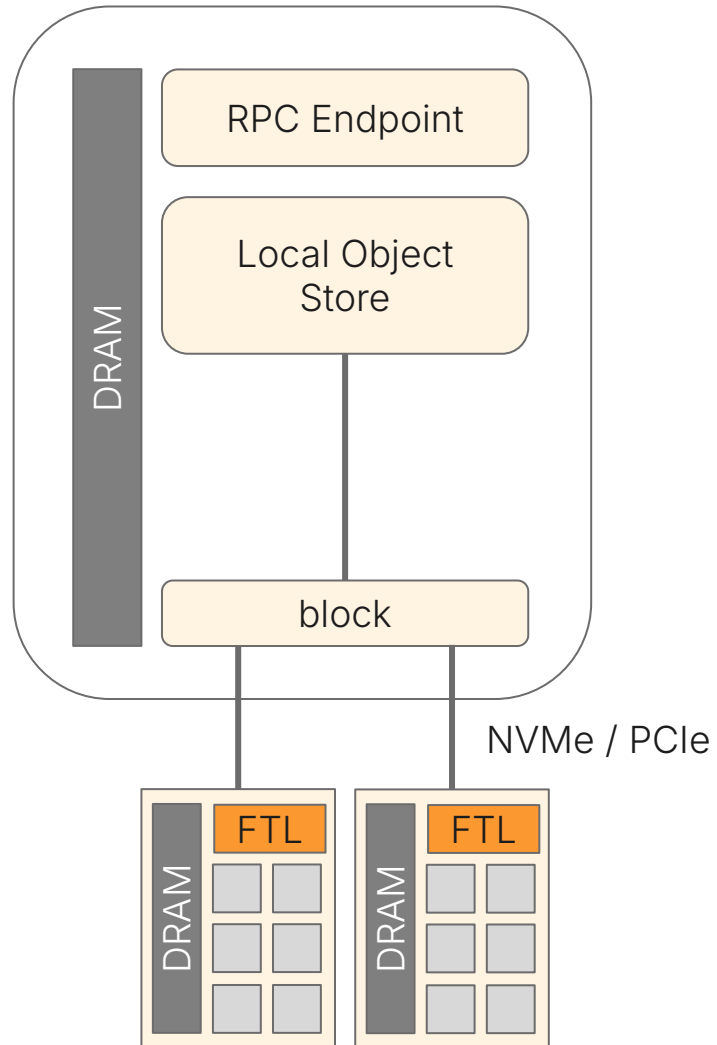
**Problem** QLC Challenges

---

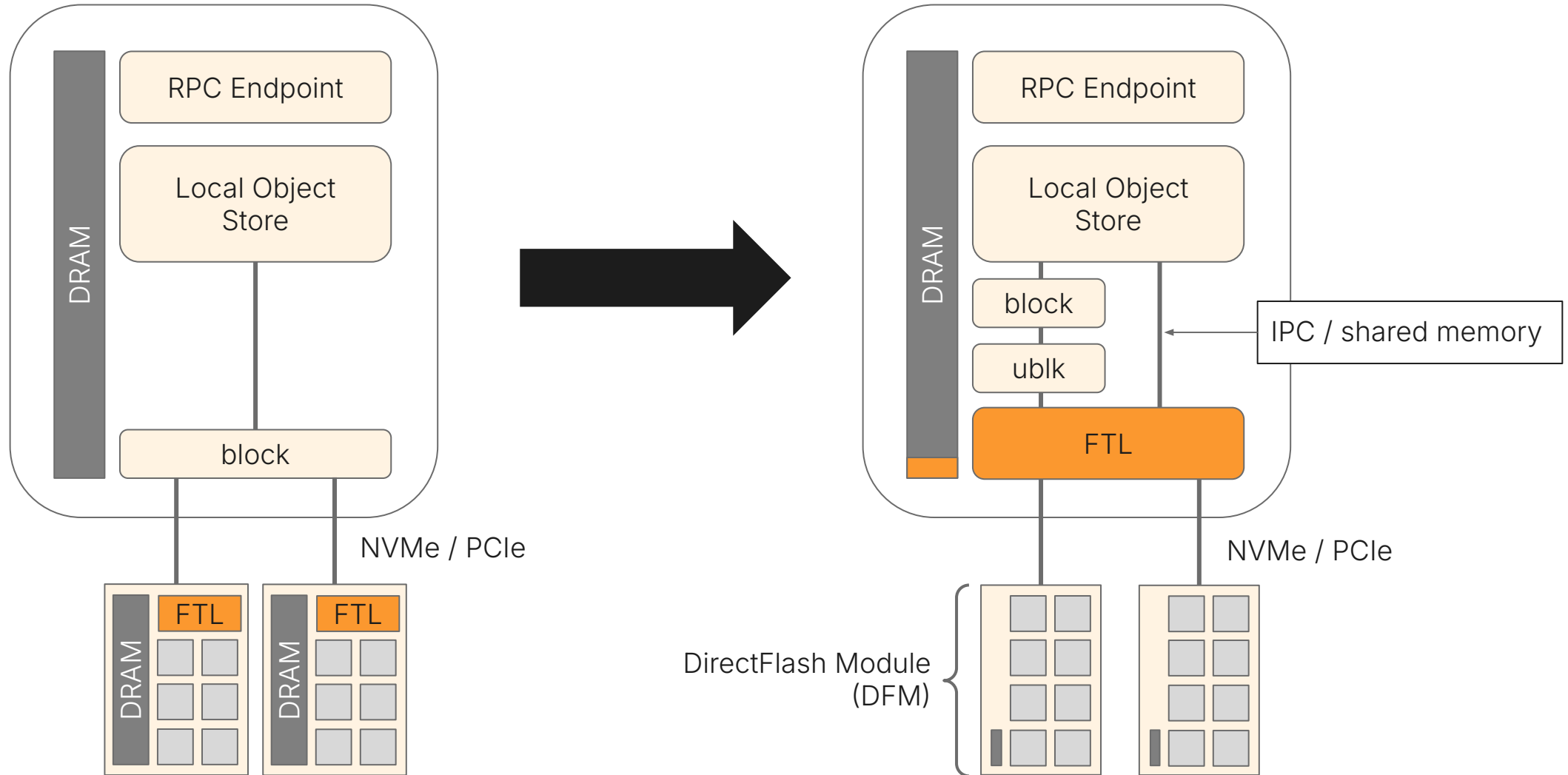
**Solution** Architectural Perspective

---

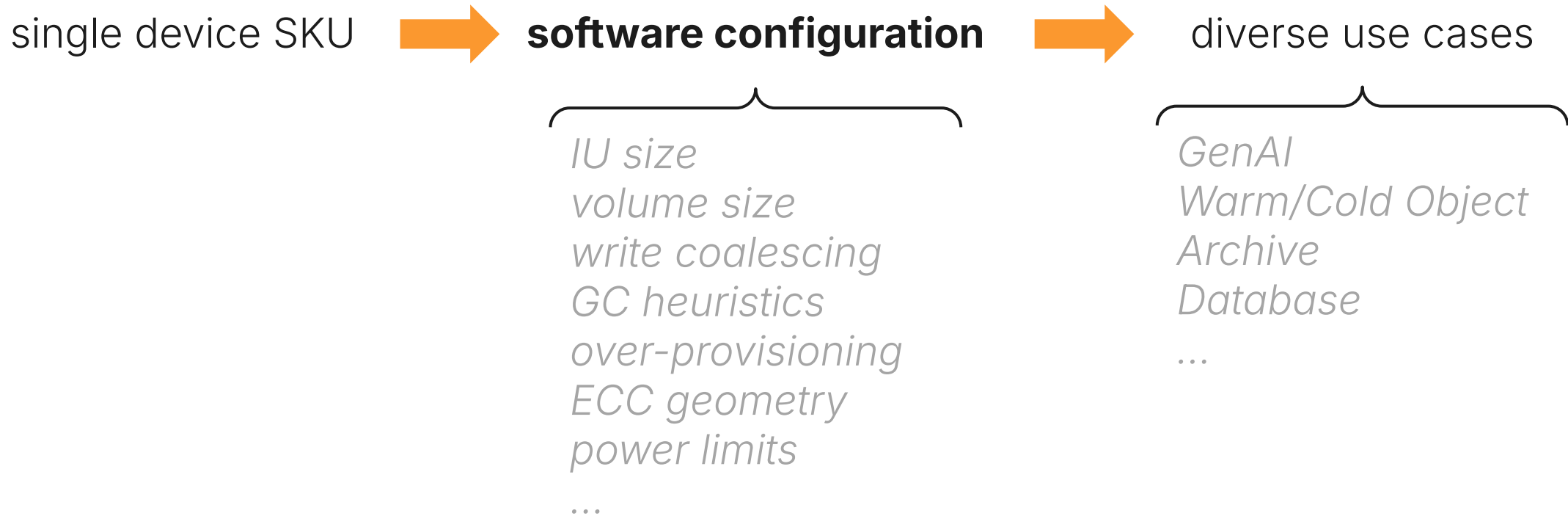
# Host-Based FTL



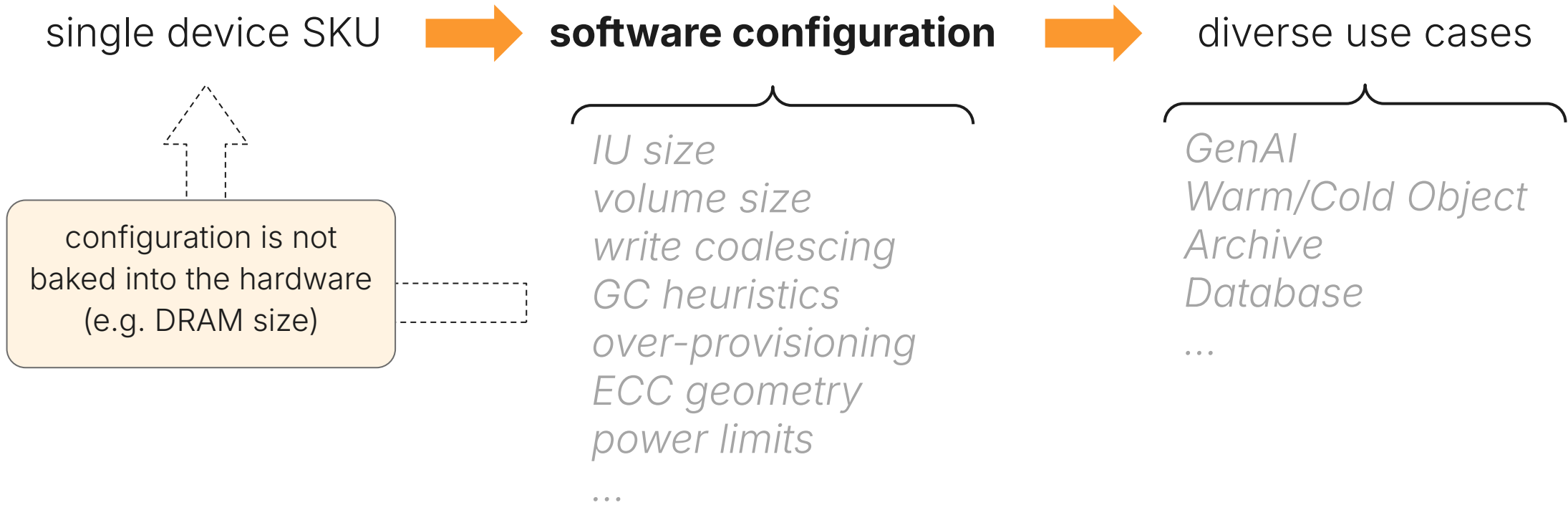
# Host-Based FTL



# Late-Binding Configuration

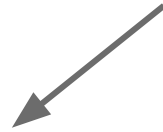


# Late-Binding Configuration



# What Next?

driving efficiency in  
exascale storage



## **configurability**

Scale drives the balance from abstraction to full-stack optimization.

Configure storage media to match the natural patterns of the local object store.

# What Next?

driving efficiency in  
exascale storage

## configurability

Scale drives the balance from abstraction to full-stack optimization.

Configure storage media to match the natural patterns of the local object store.

## media lifetime

low WAF → media lasts **10 years**

How can we fully leverage this, technically and operationally?



**PURESTORAGE®**