

# Indirection Systems for Shingled-Recording Disk Drives

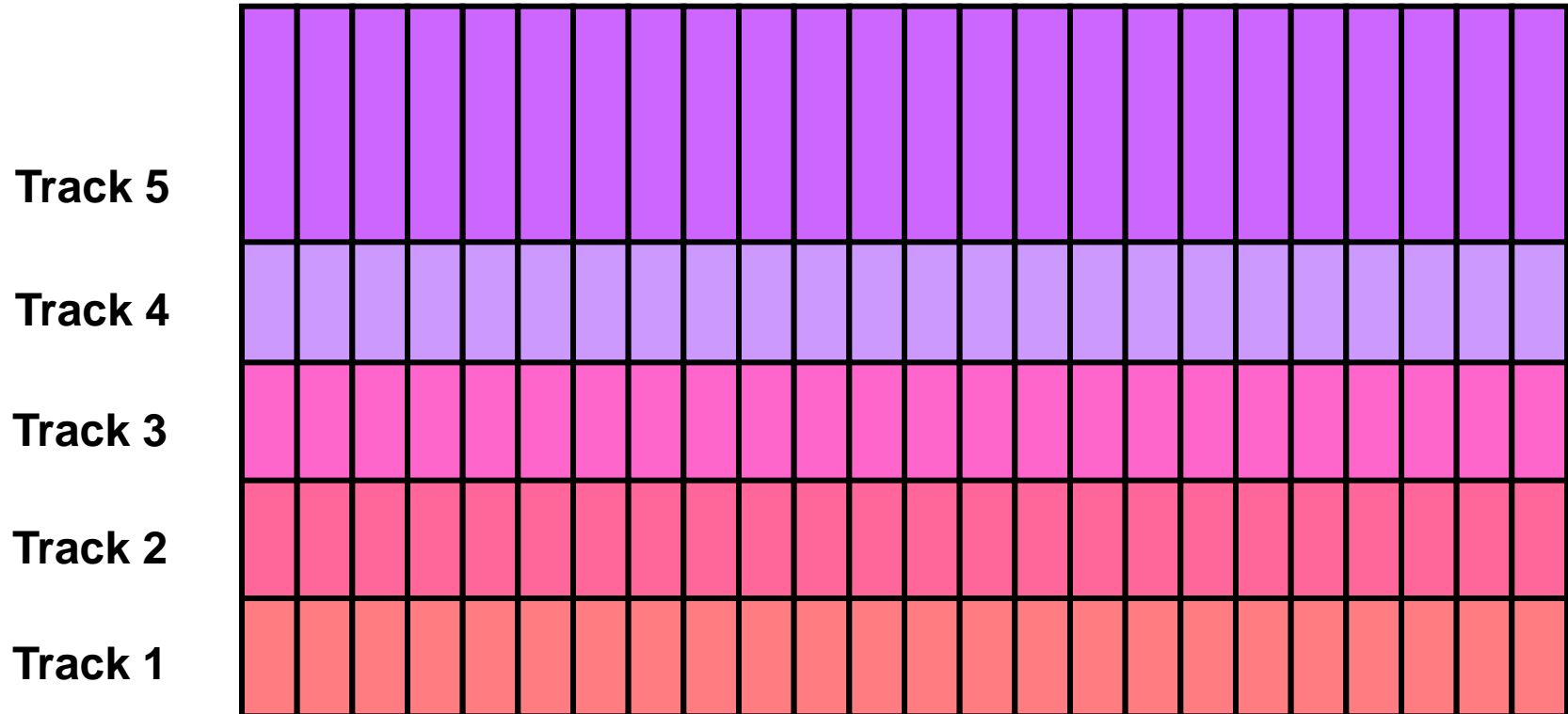
**Yuval Cassuto**

with co-authors:

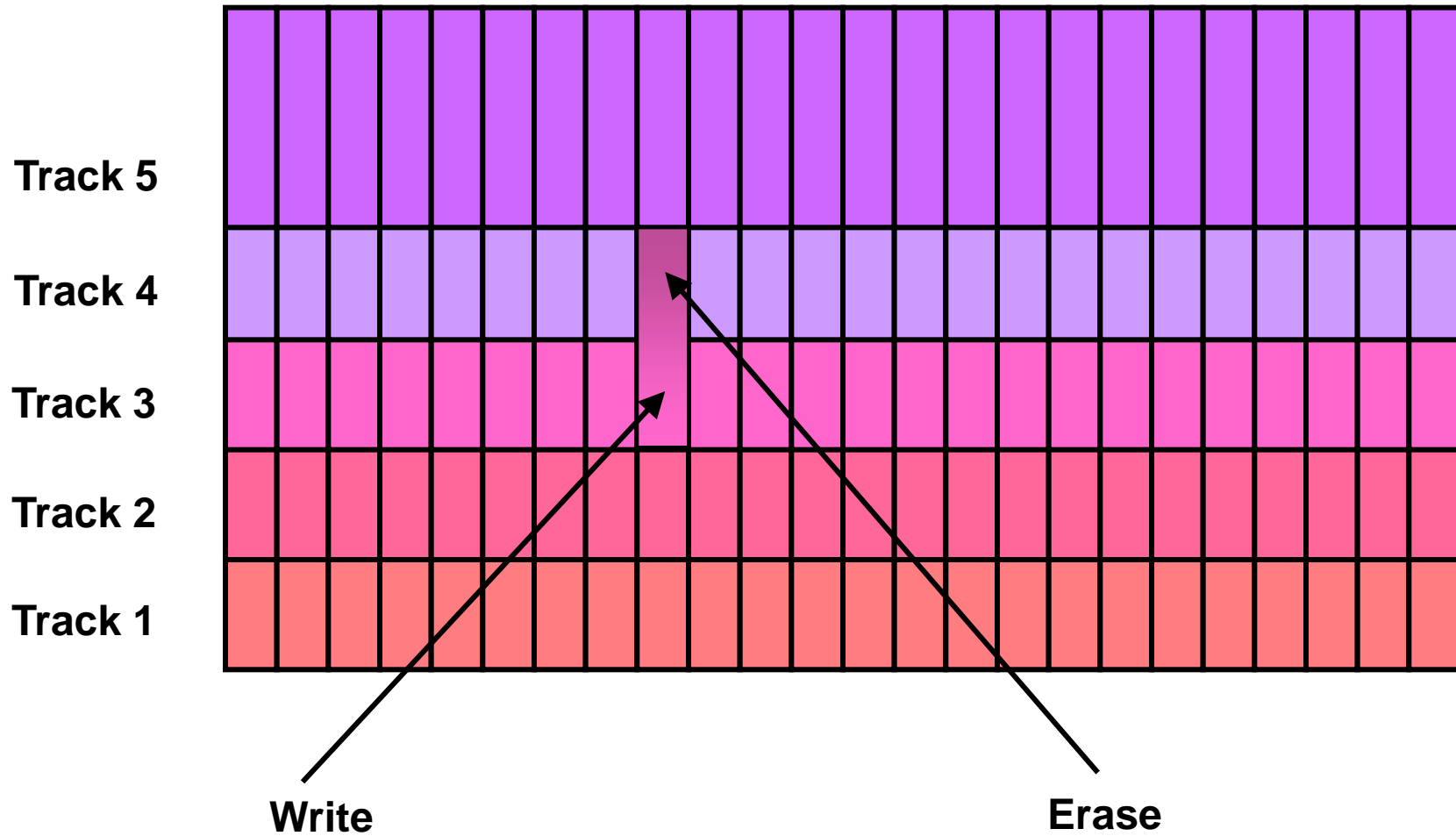
**M. Sanvido, C. Guyot, D. Hall and Z. Bandic**

 **Hitachi Global Storage Technologies**

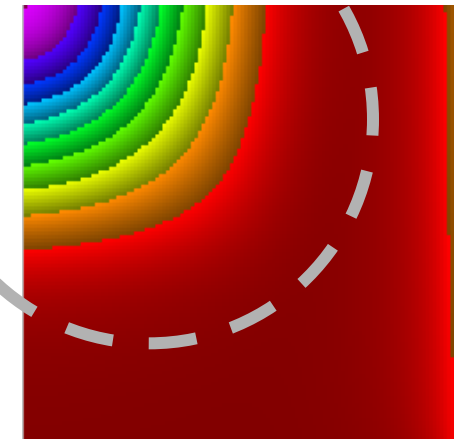
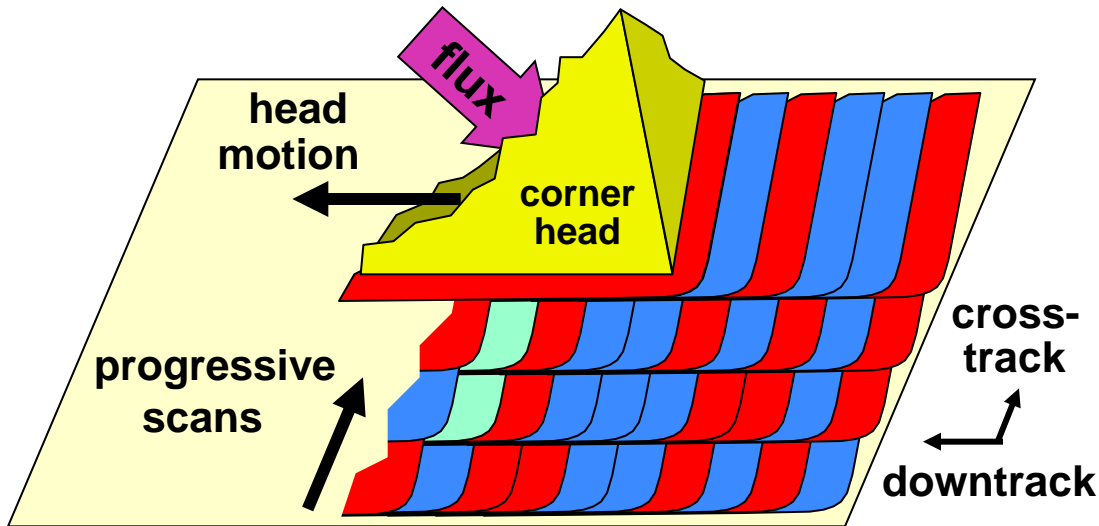
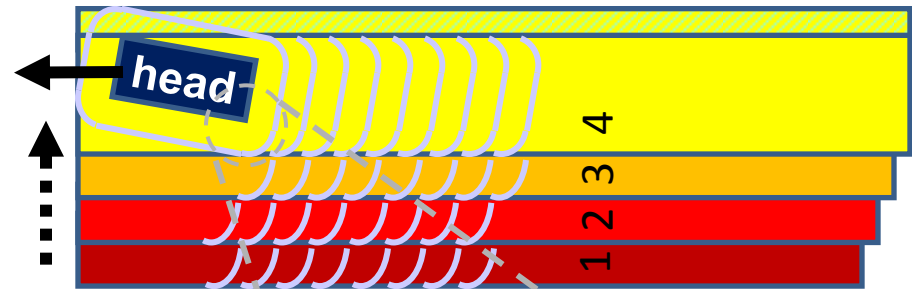
MSST 2010, May 7



# Shingled Recording – No Random Write



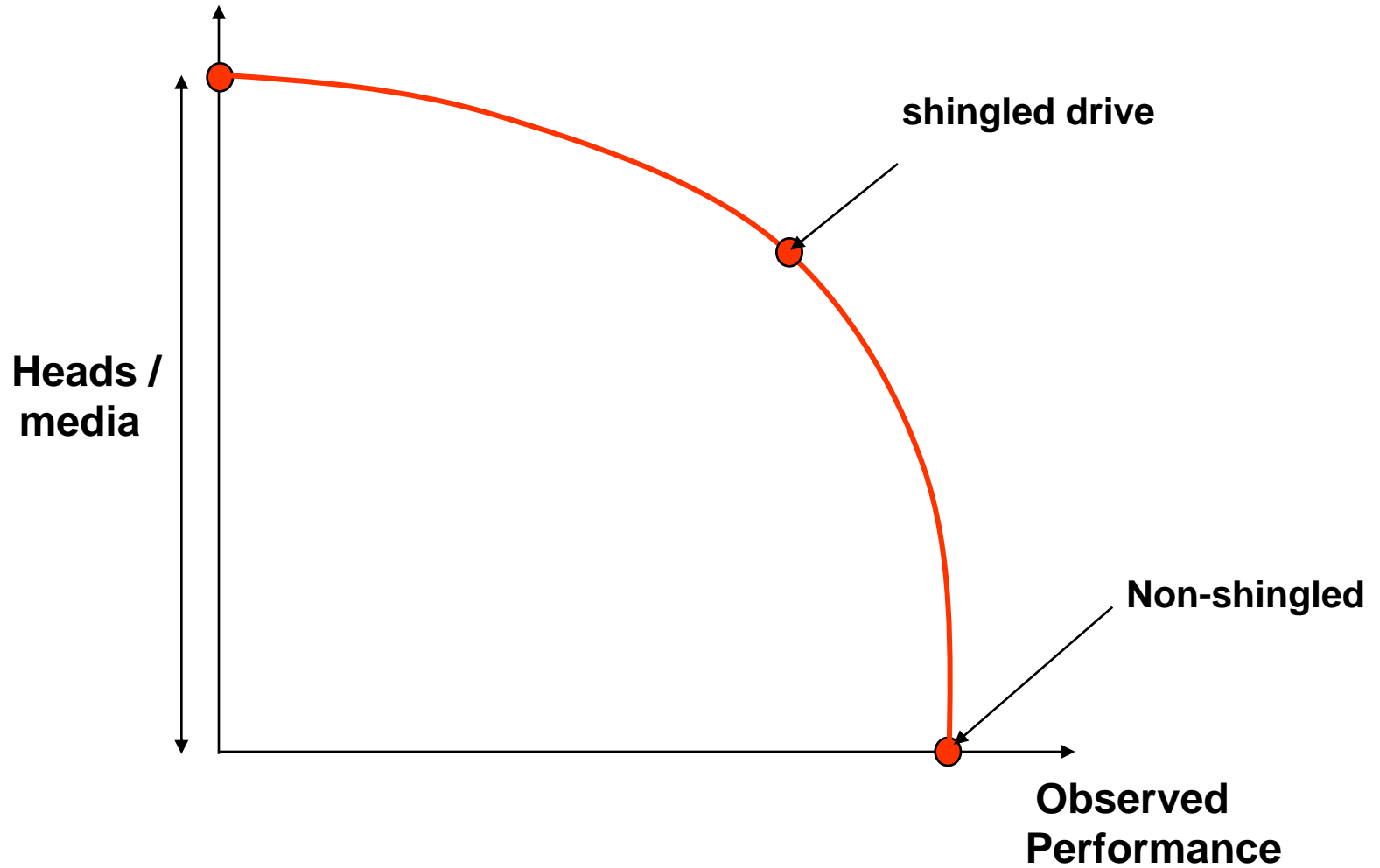
track layout for shingled-recording



head field contours

- **More:** Capacity
- **Less:** Functionality → Performance

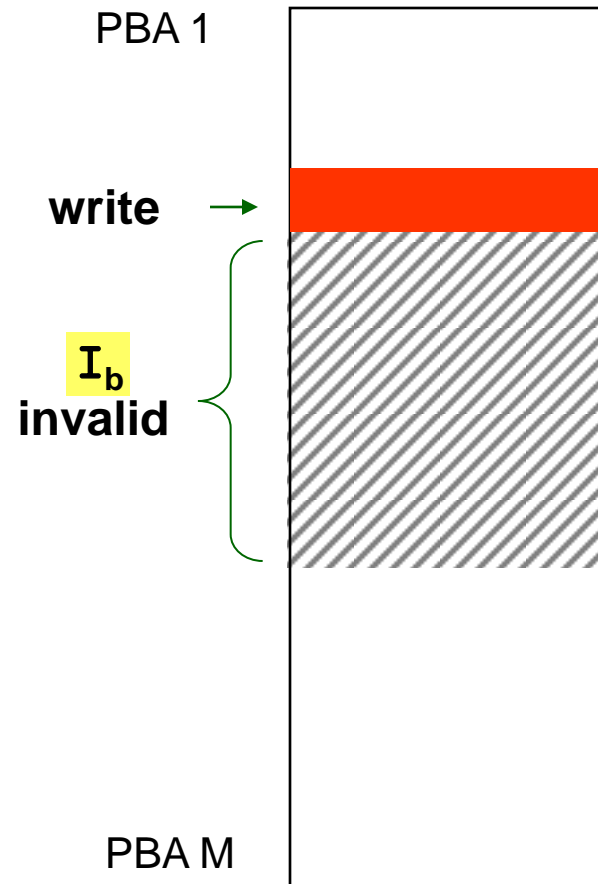
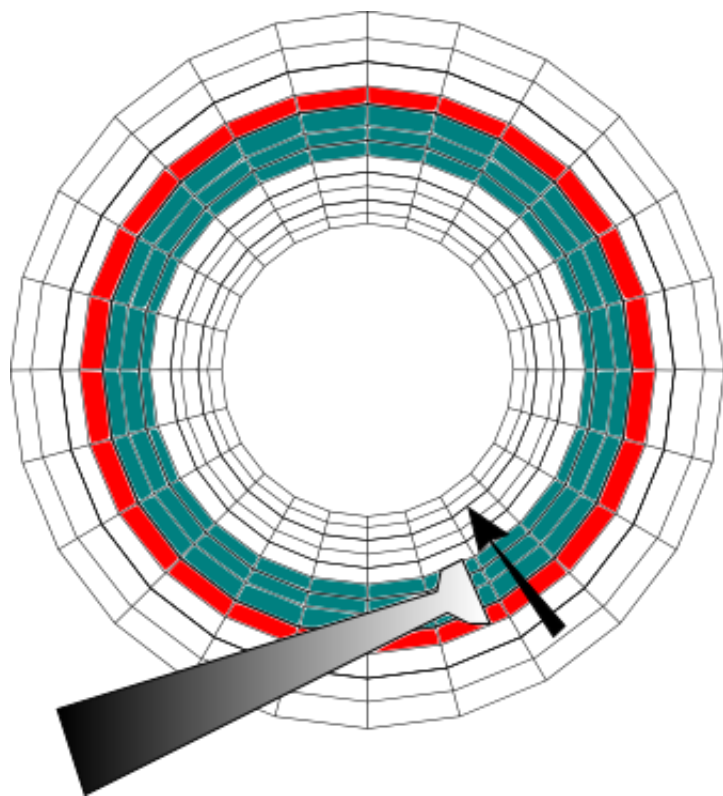
**Excess Capacity**

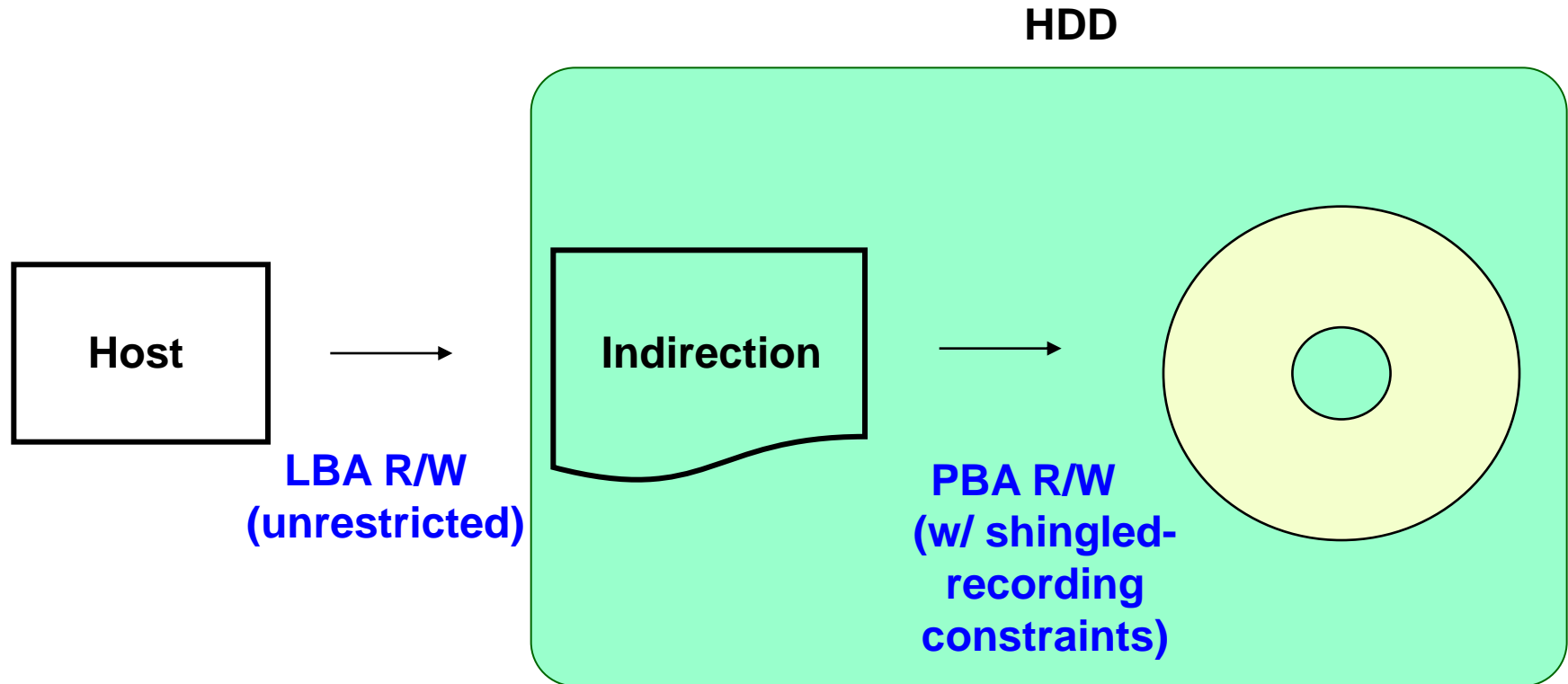


Track Interference



Block Interference







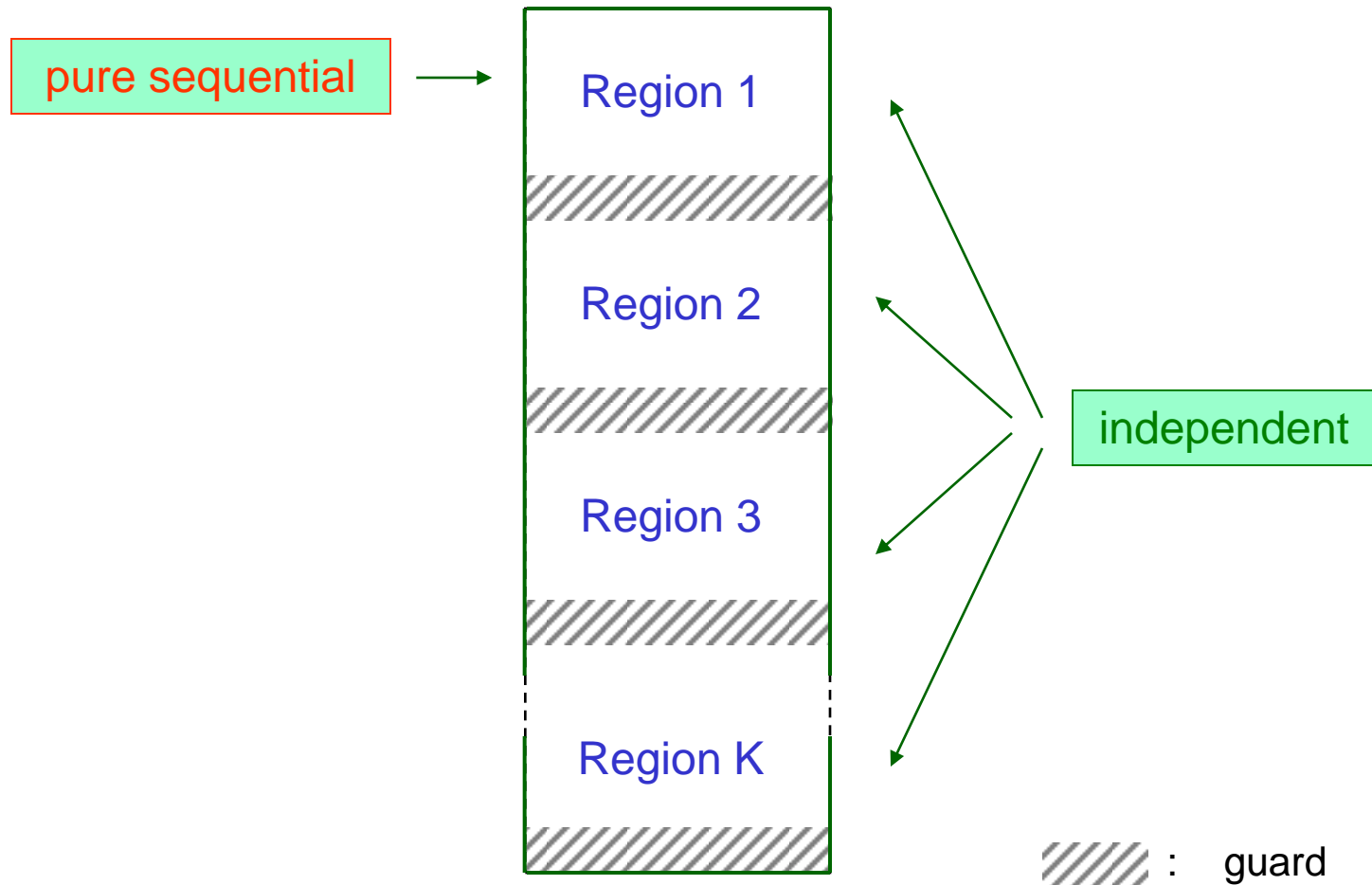
## Why on the drive?

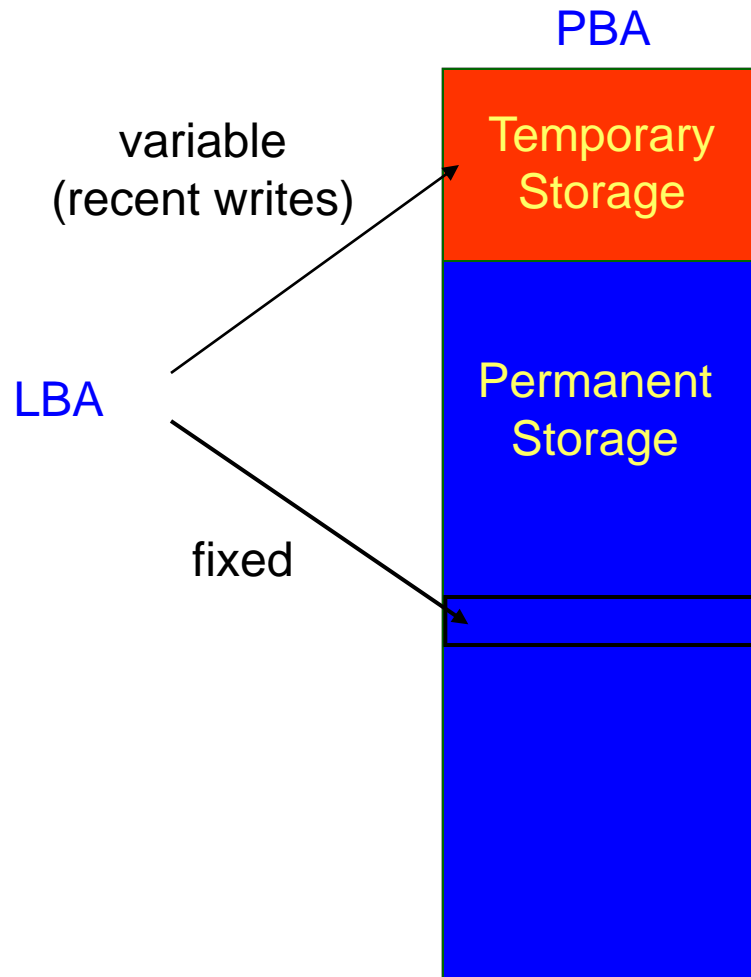
- Transparent to Host
- Complete knowledge of physical layout

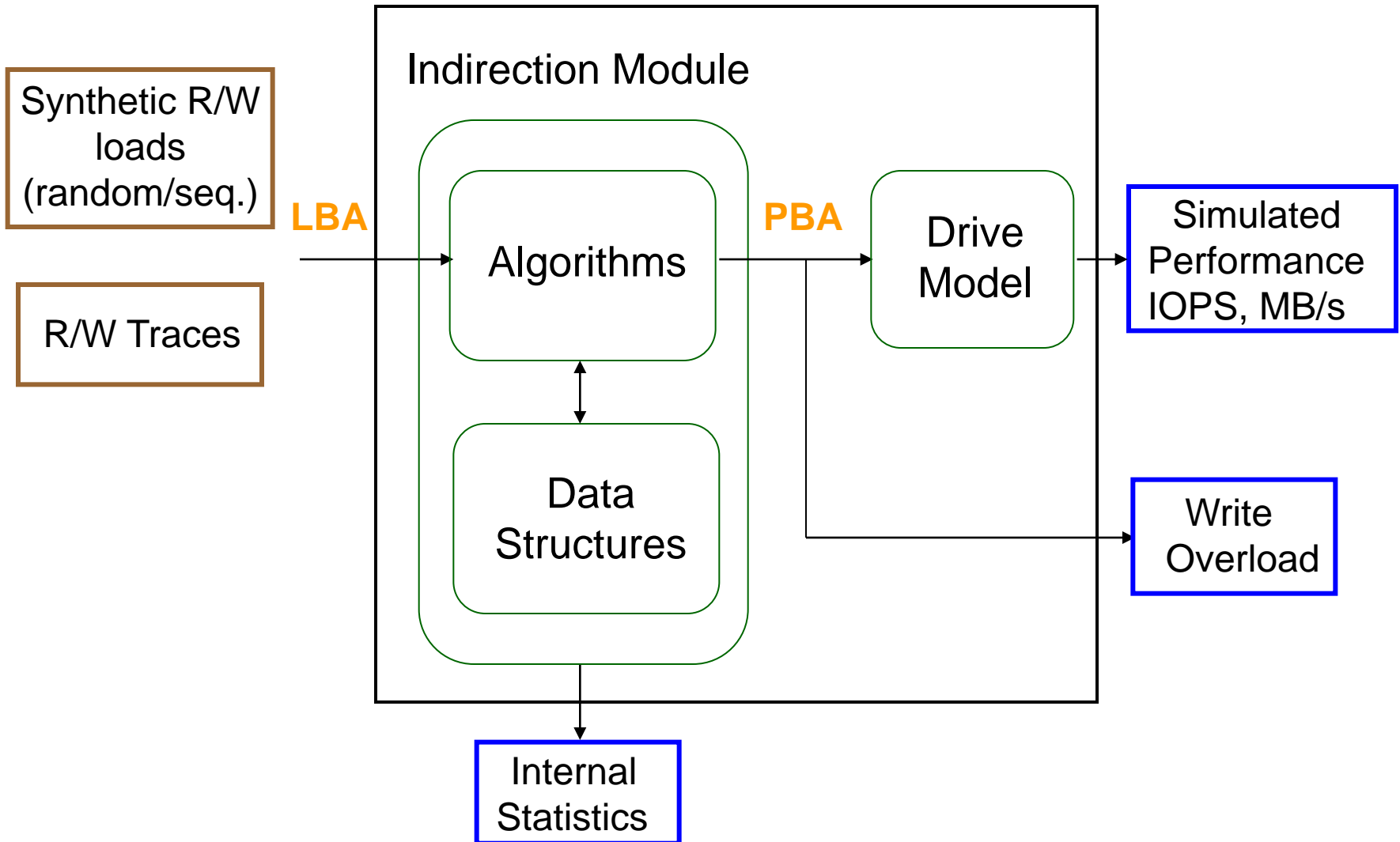
## Why on the host?

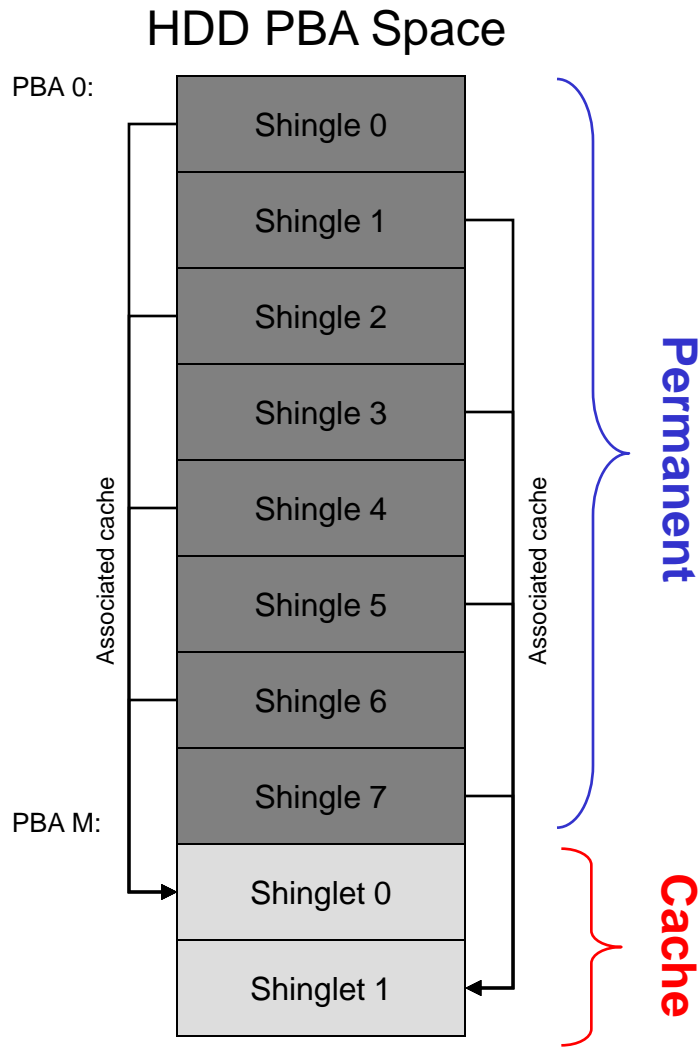
- “Shingle aware” access and allocation
- System specific performance optimization

- **Append/Read-Modify-Write** with **Shingled Regions**

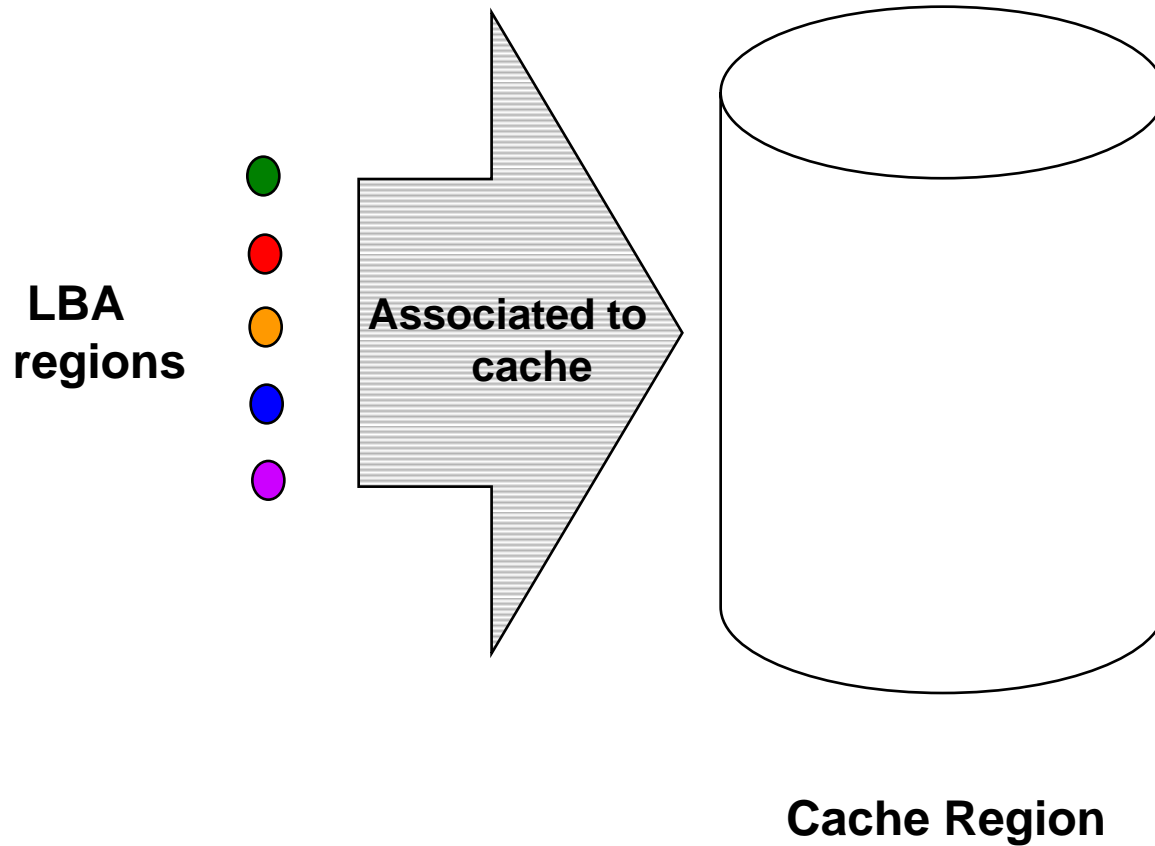




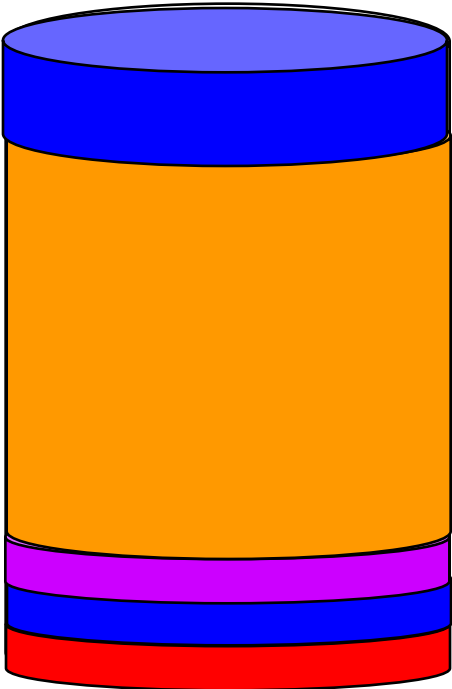




- Disk is partitioned to LBA shingled regions and cache shingled regions
- Each LBA region (shingle) is associated with one cache region (shinglet)
- Multiple LBA regions are associated with each cache region (set-associative cache)

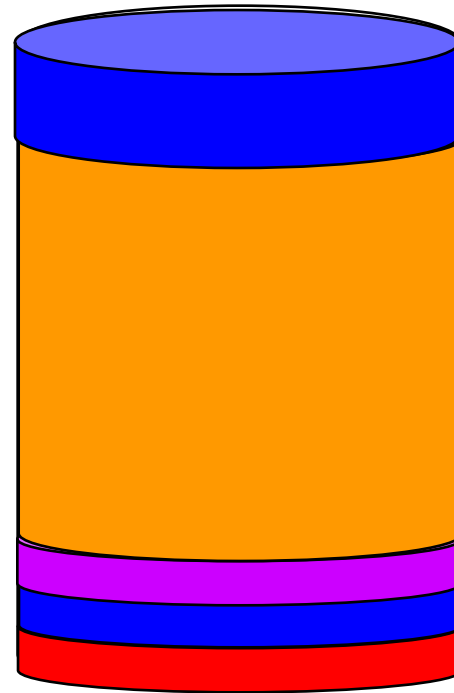
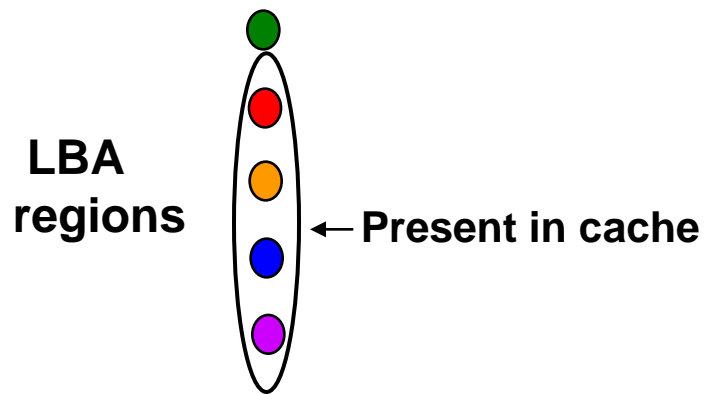


LBA  
regions



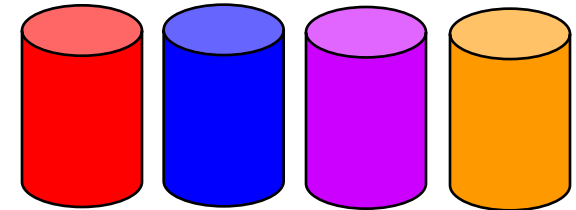
**FULL!**

**Cache Region**



Cache Region

- Read cache
- For each LBA region present in cache: RMW\* full region

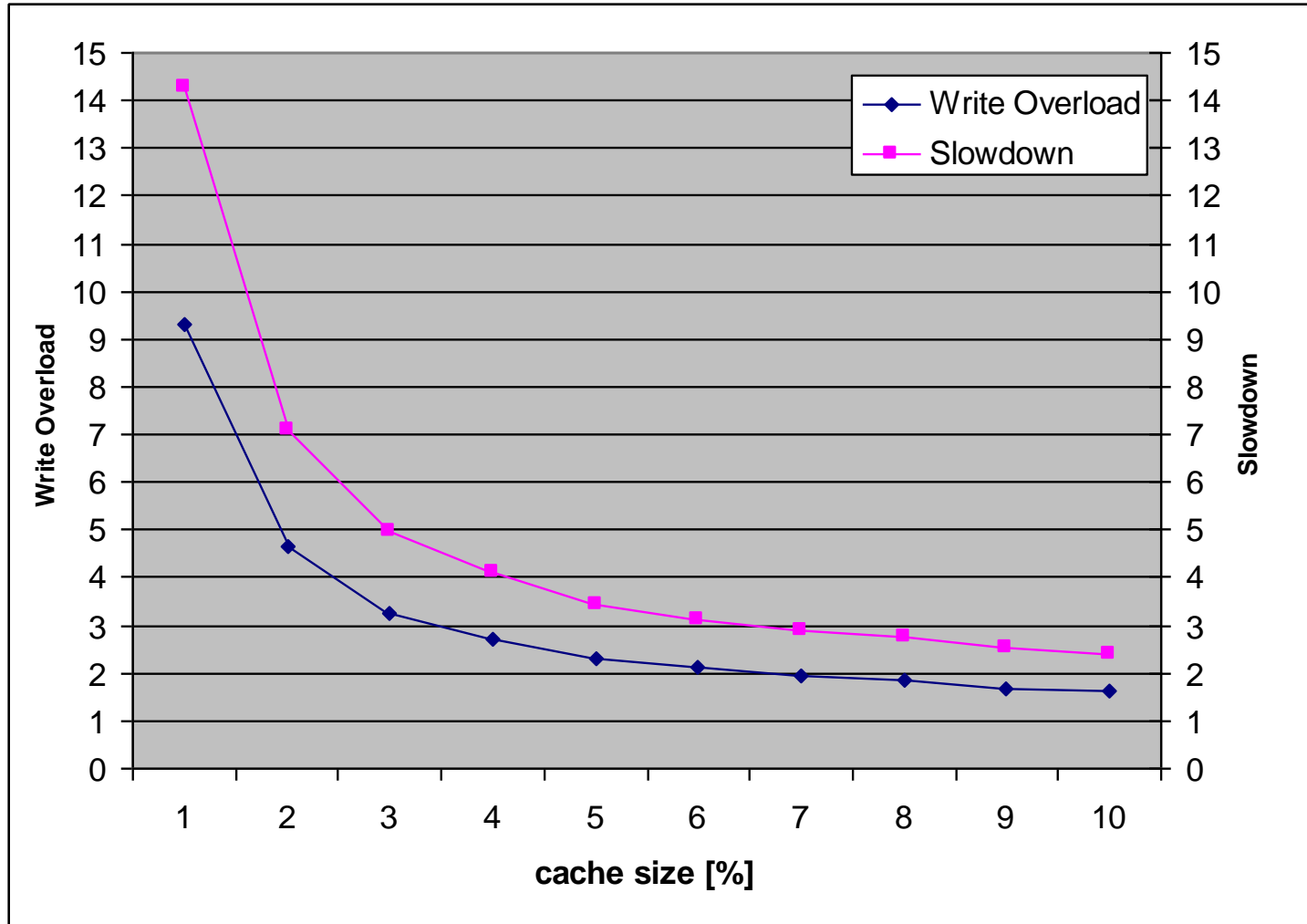


(Read+Write) x #Present

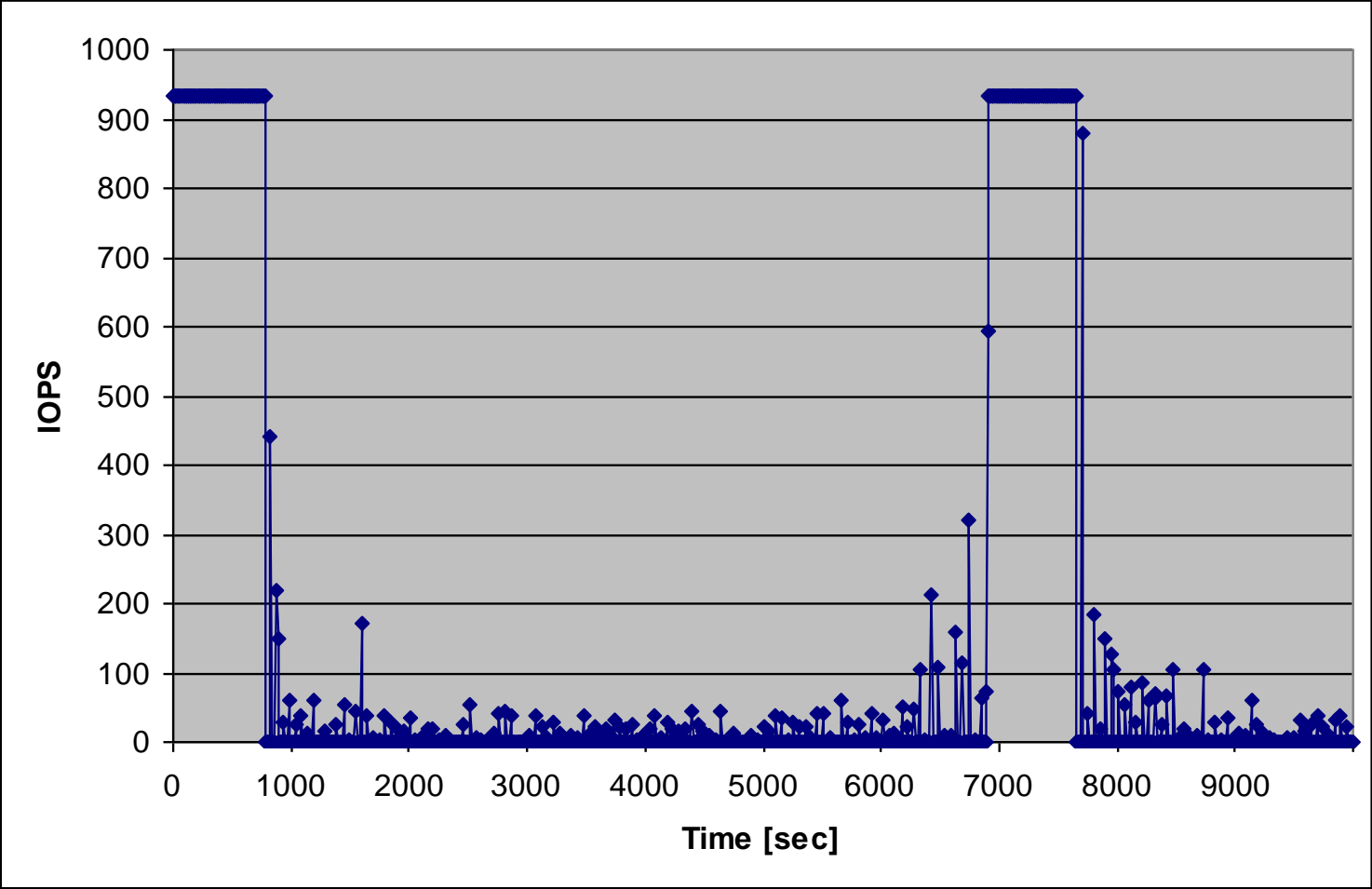
- Full cache available again

\* RMW = Read-Modify-Write

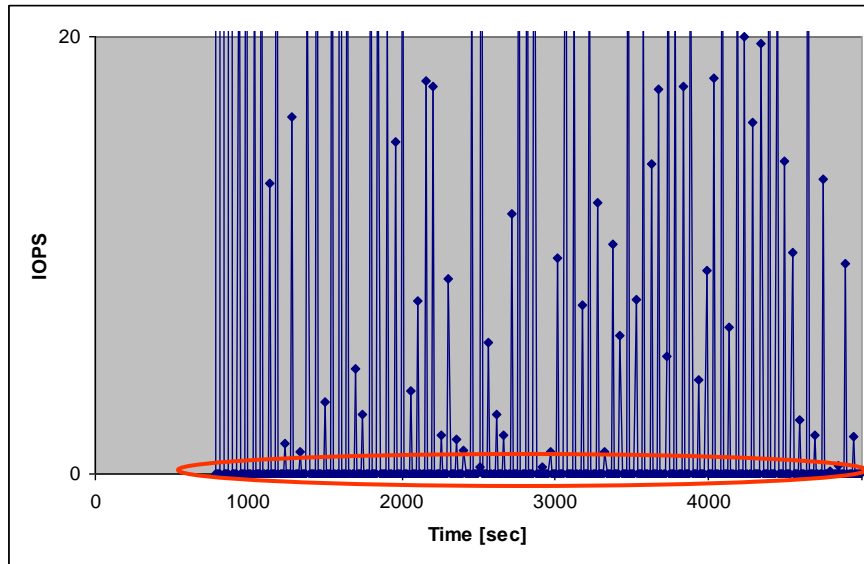




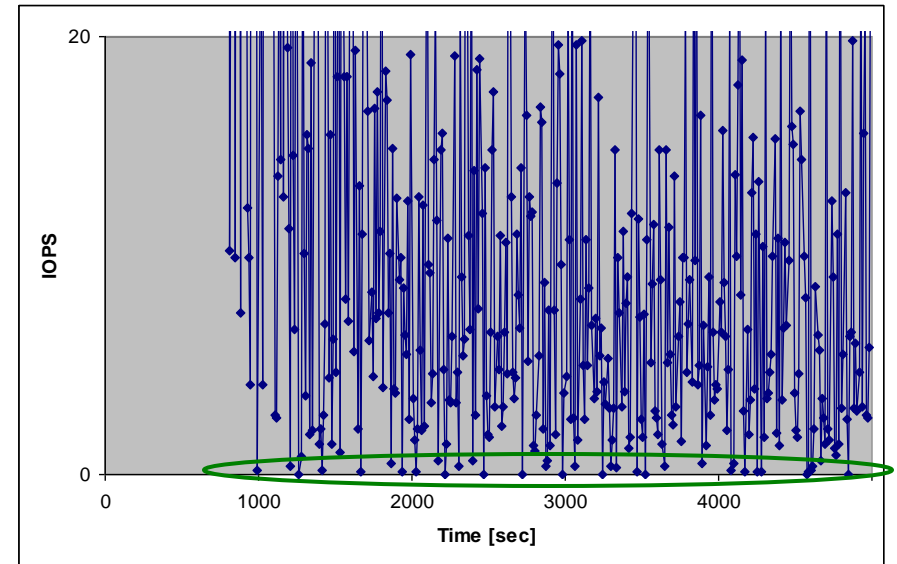
# 4K Random IOPS



Size 50000

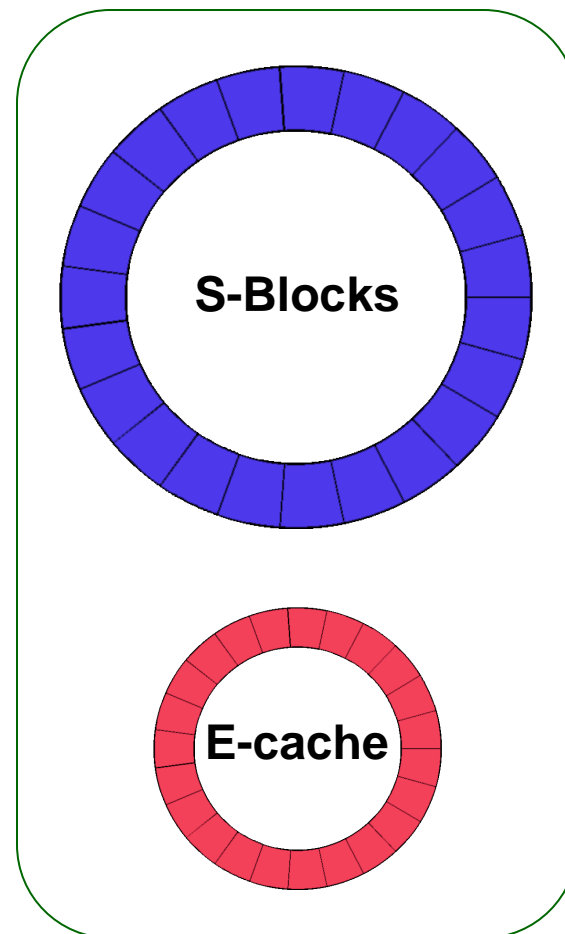


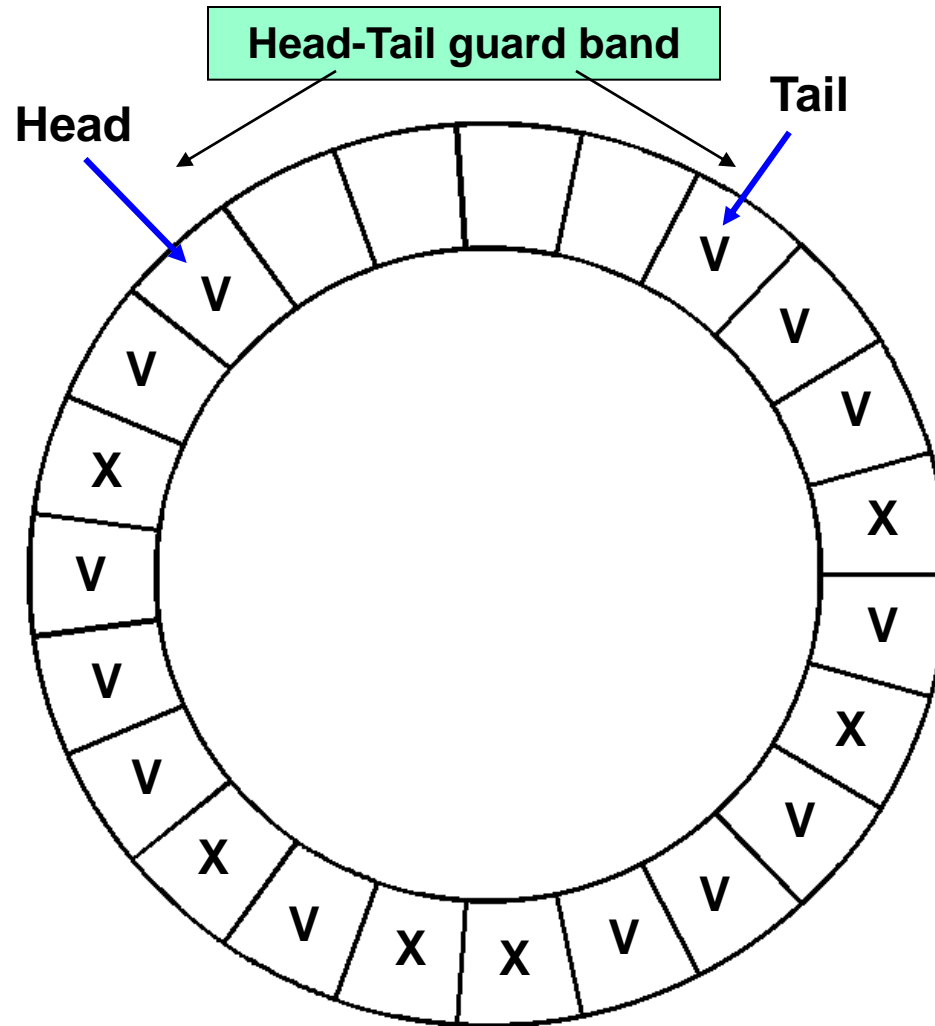
Size 10000



- Significant improvement over plain Read-Modify-Write
  - Simple to implement
- 
- Large dips in performance due to long garbage collections
  - Region updated in place → consistency issues

- Temporary (red) and permanent (blue) storage managed as ring buffers
- S-Block: intermediate unit between sector and region

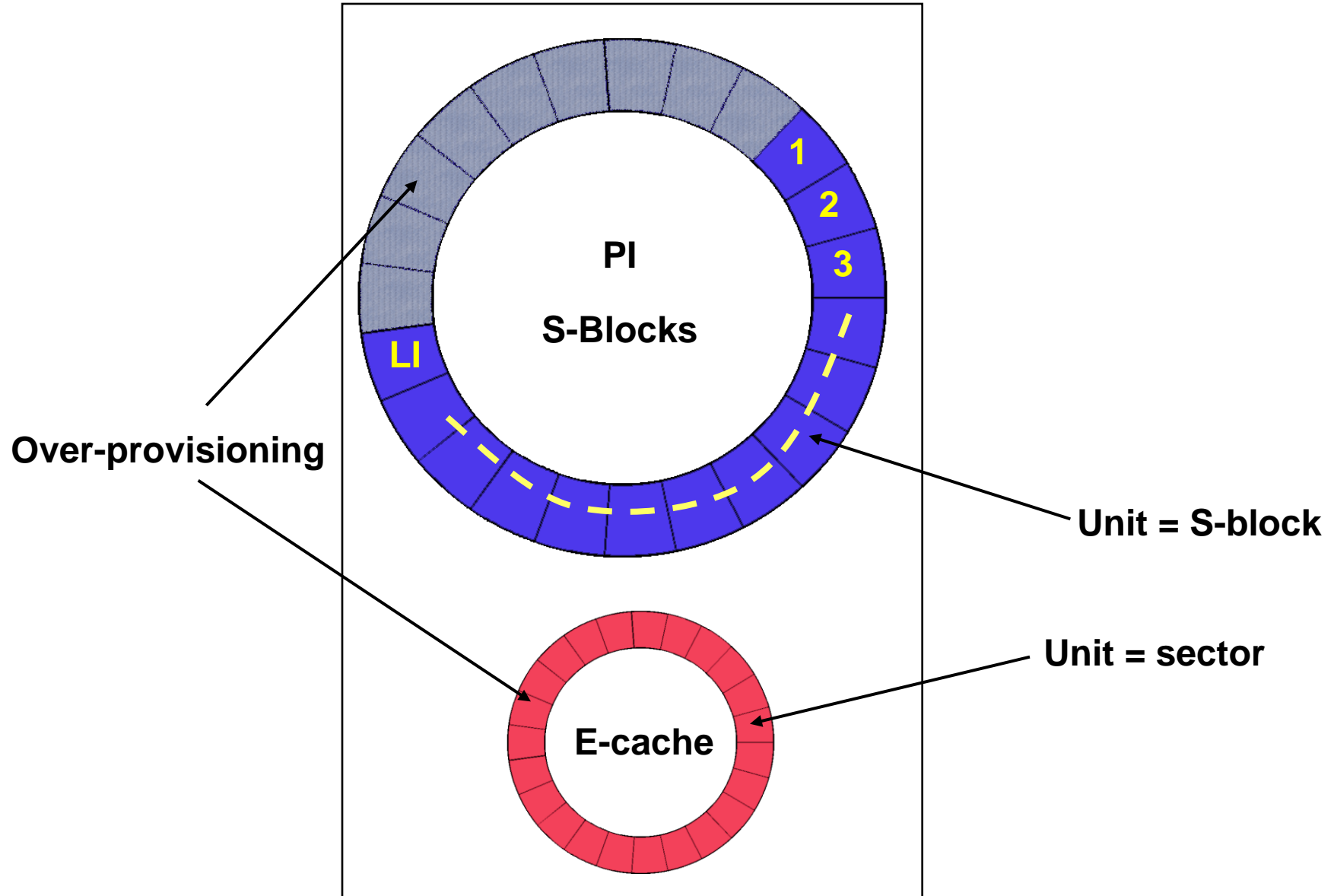


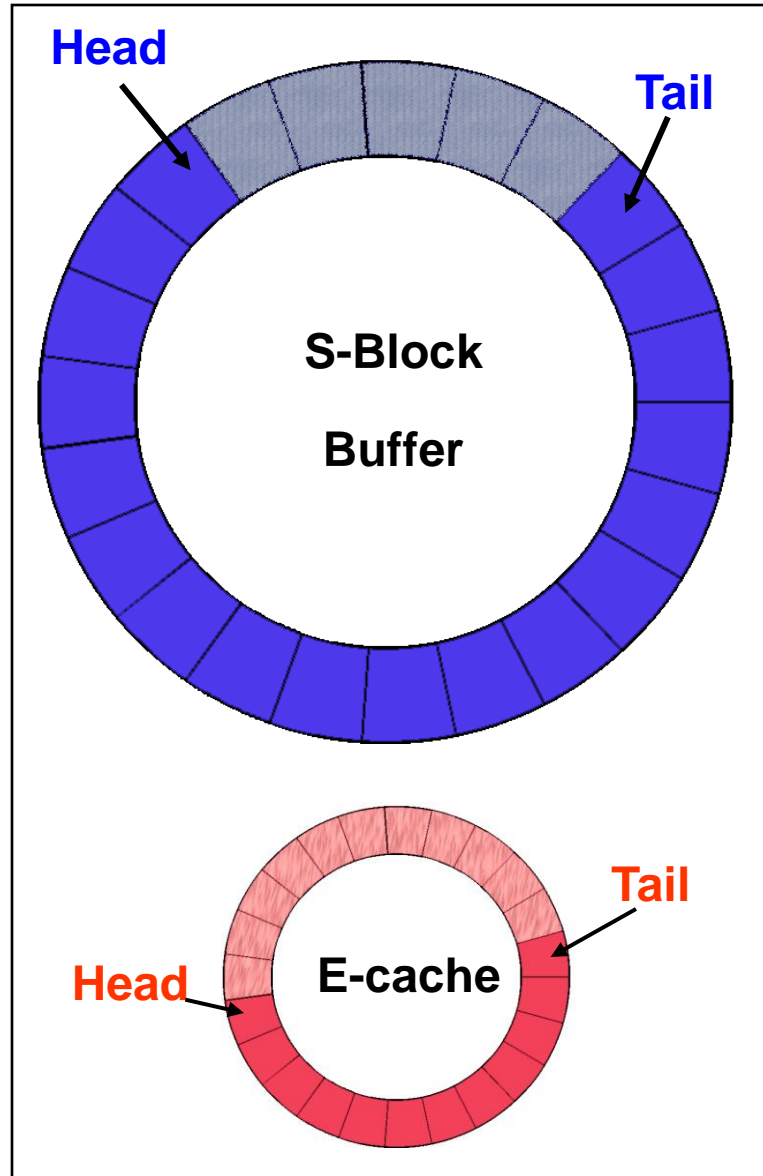


**V** Used, valid

**X** Used, invalid

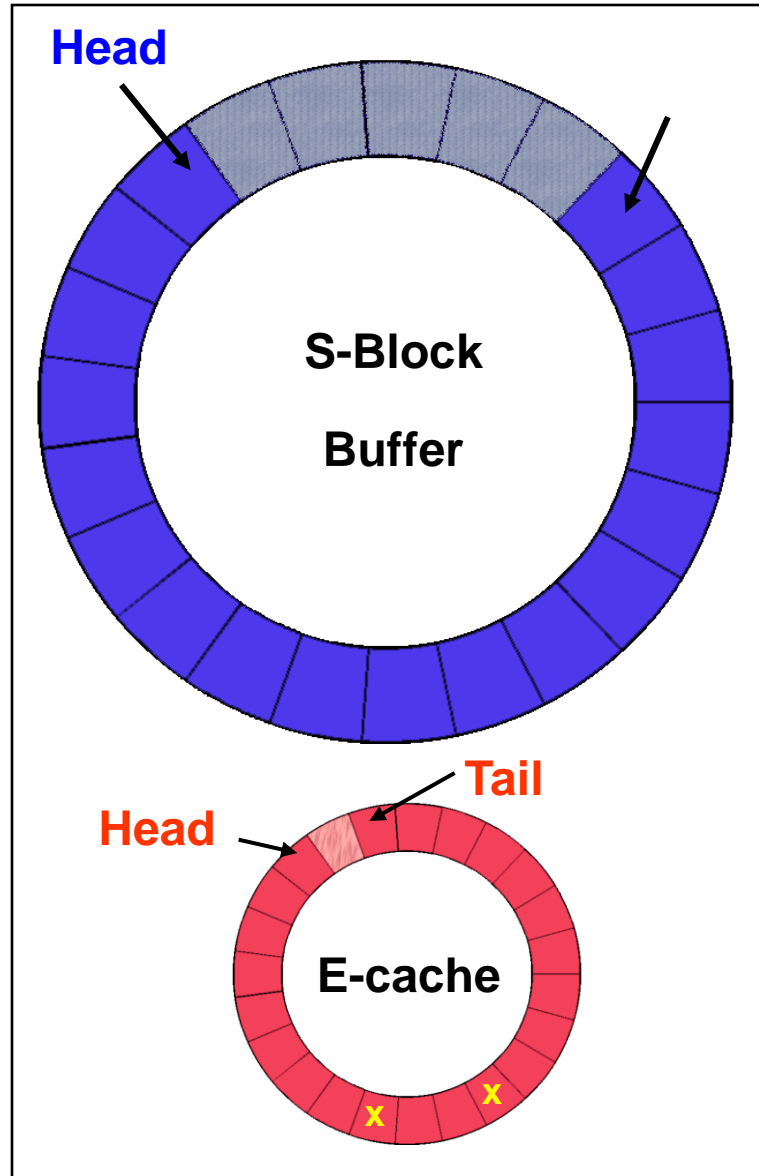
**□** Unused



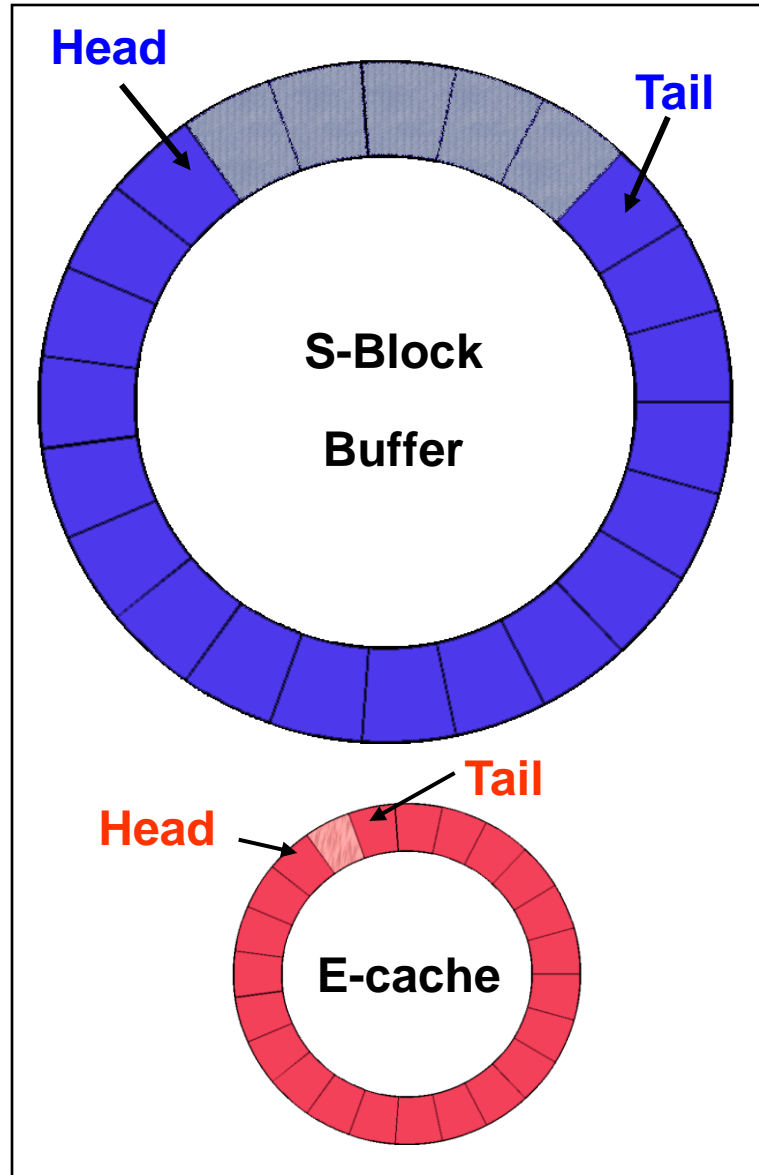


- If full S-Block: write directly in S-Block buffer
- If smaller than full S-Block: write in E-cache

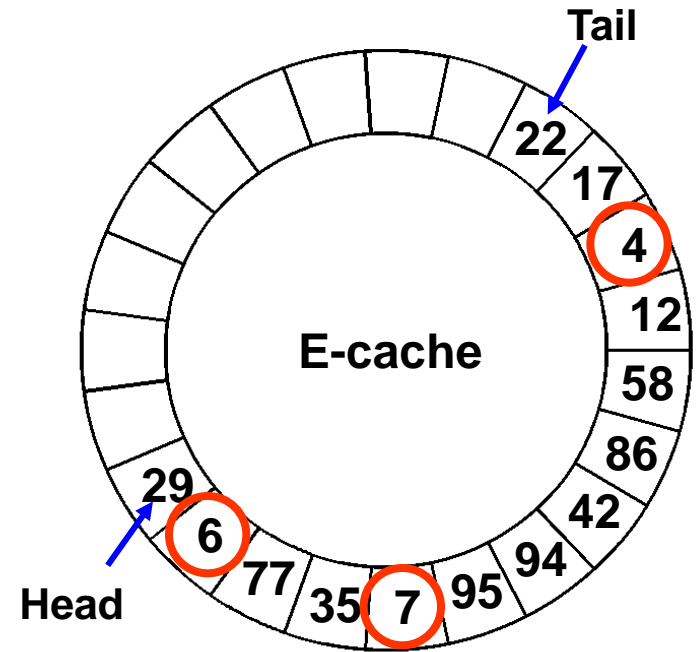
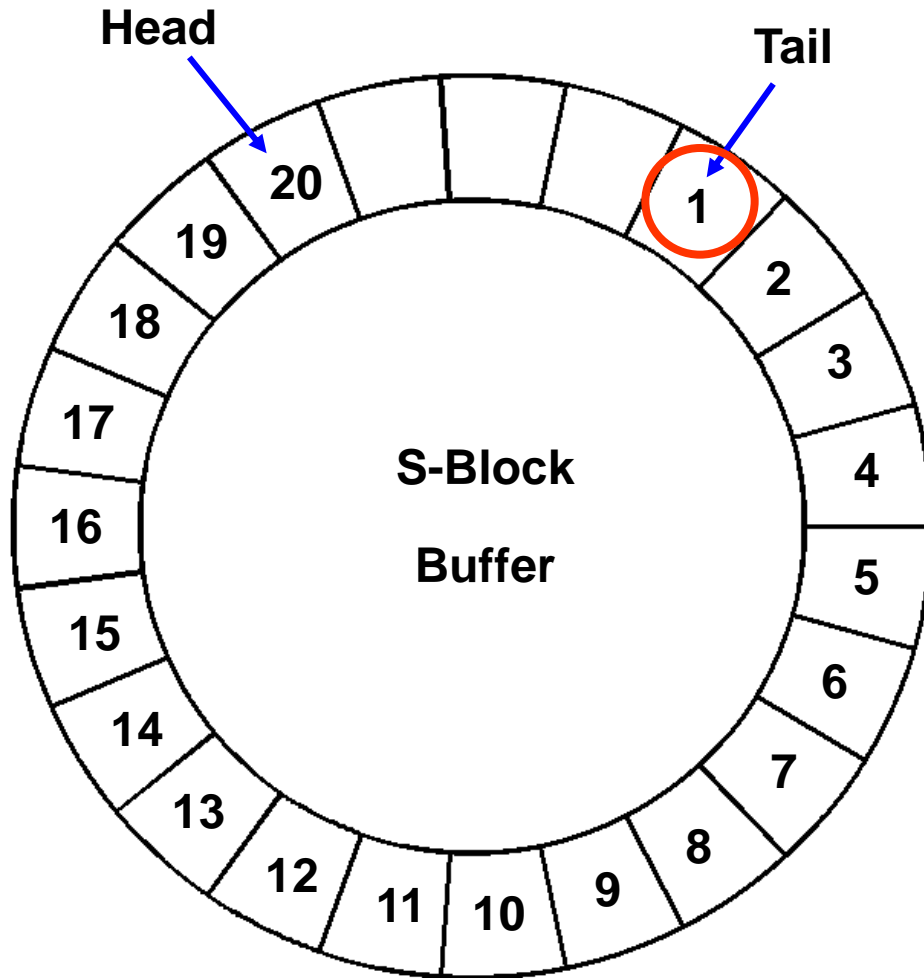




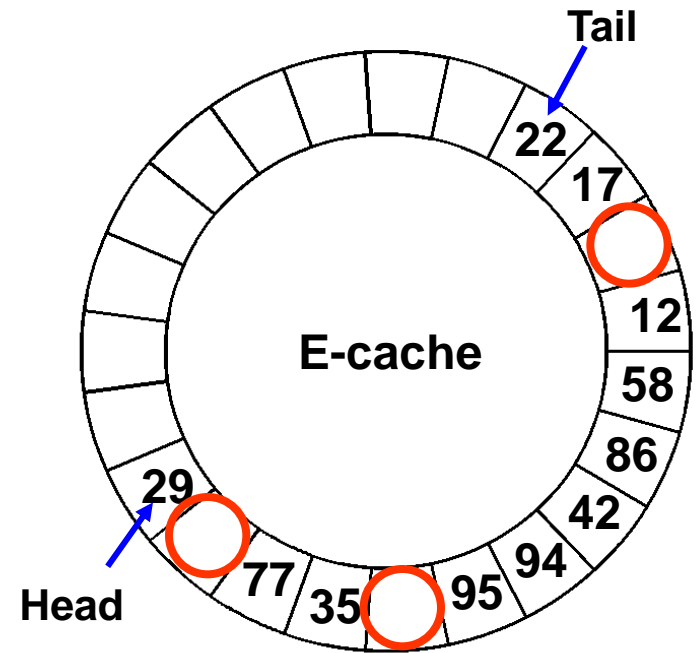
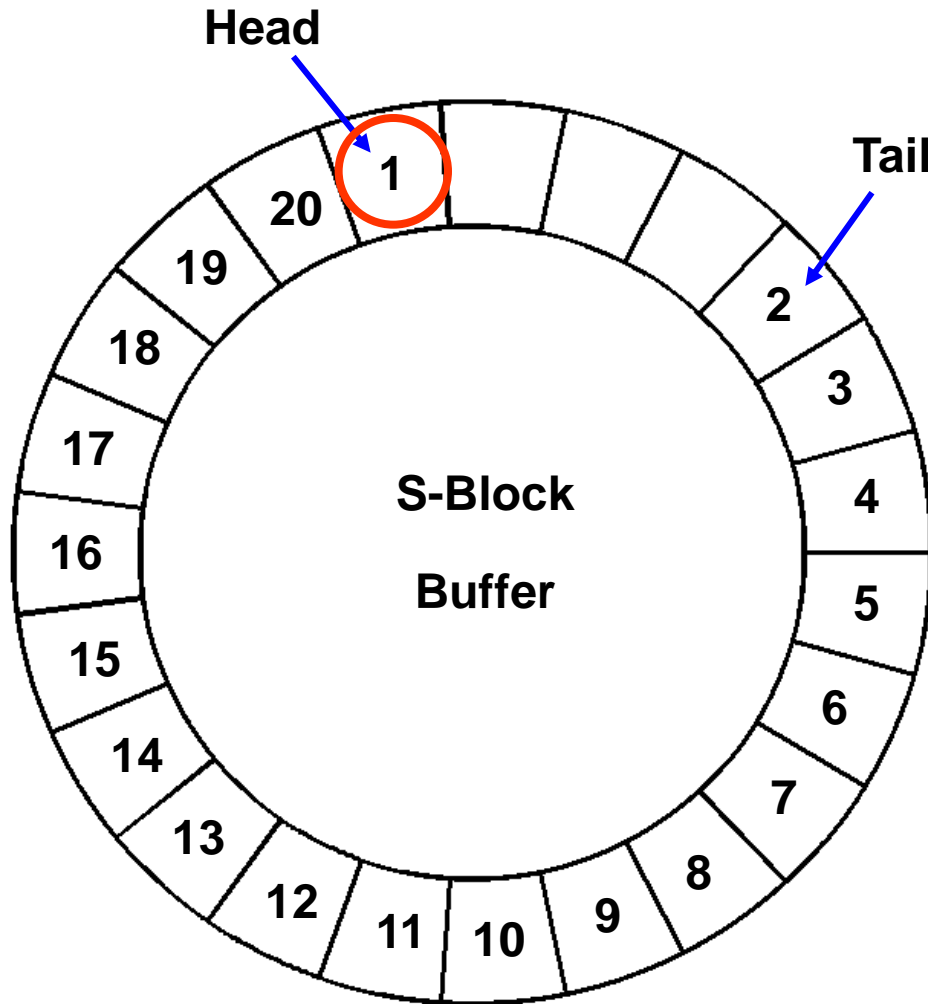
- If no room for incoming write: reclaim invalid exceptions in E-cache (defrag)



- If not enough invalid exceptions:
  - Choose S-Block
  - Destage(S-Block)

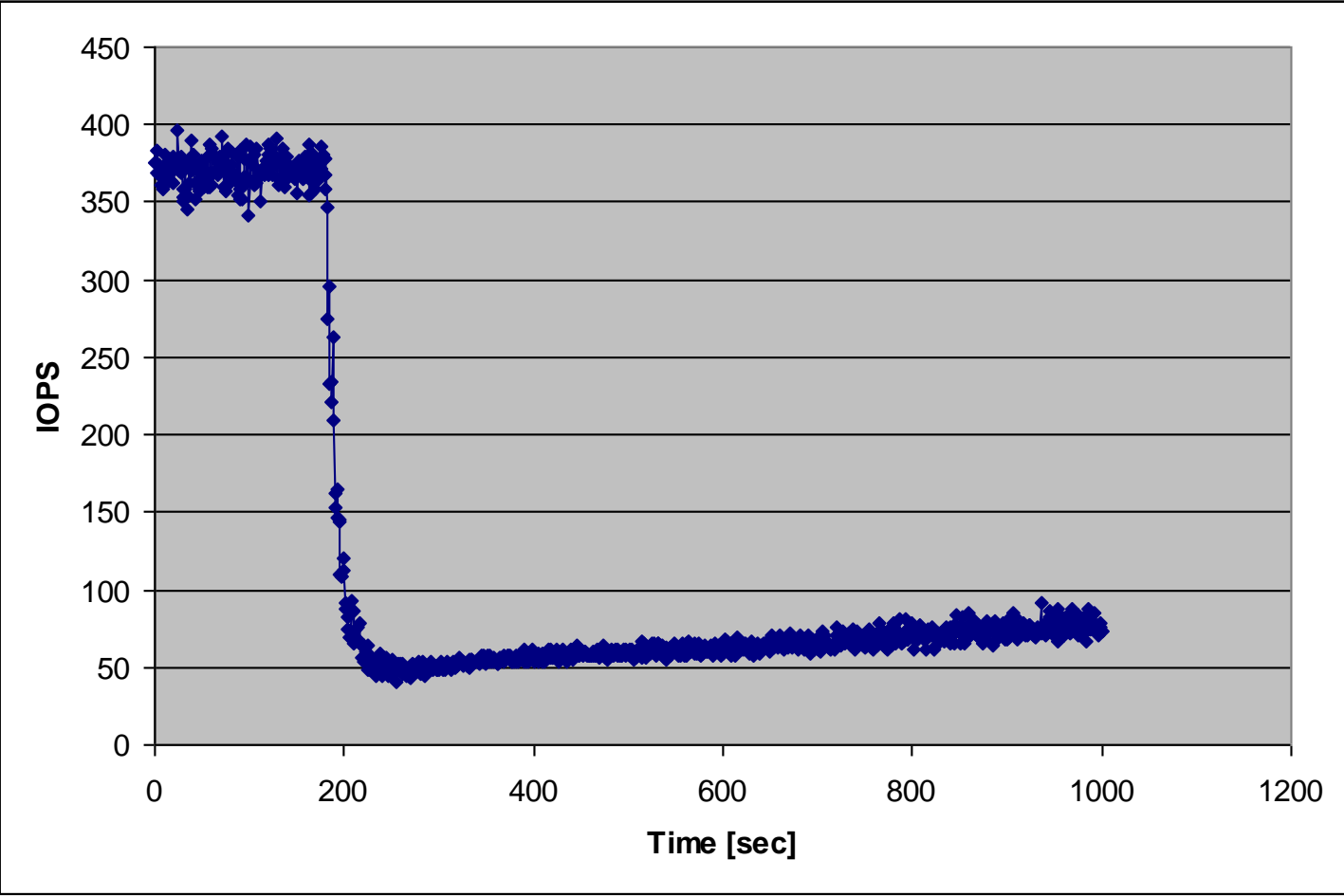


**S-Block 1 → sectors 1-10**



**S-block writes in E-cache are not needed after destage**

# 4K Random IOPS



## Optimal Destage

- Choose S-Block with highest exception count
- Best amortization of S-Block rewrite
- Good for biased workloads (e.g. hotspots)

## Optimal Defrag

- Choose S-Block closest to the tail
- Least amount of copying in S-Block defrag
- Good for random workloads

## Optimal Destage

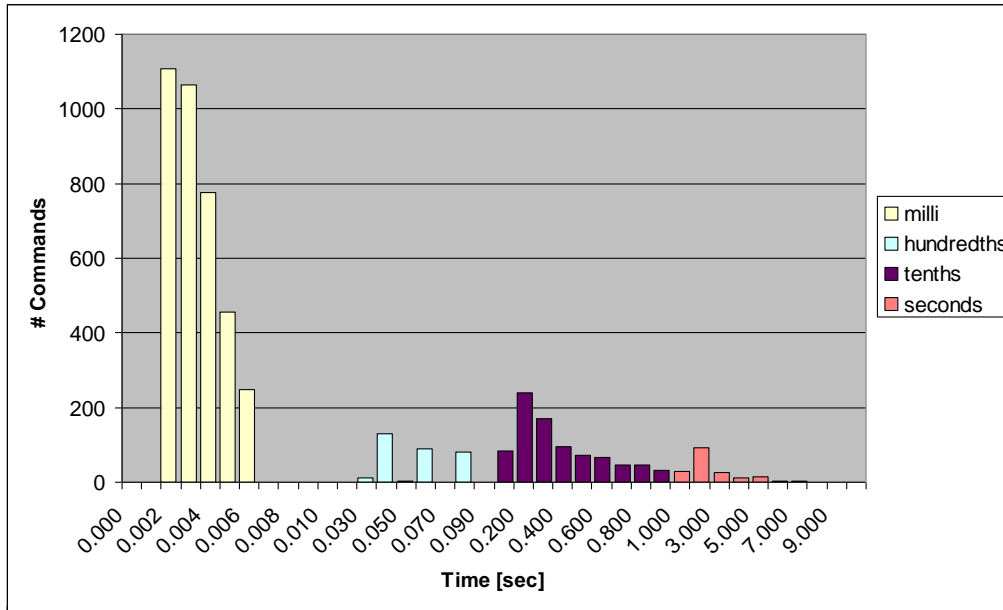
- S-Block destage average # invalidated exceptions **55.26**
- S-Block defrag average copy count **1744.64**
- average IOPS = **73.44**

## Optimal Defrag

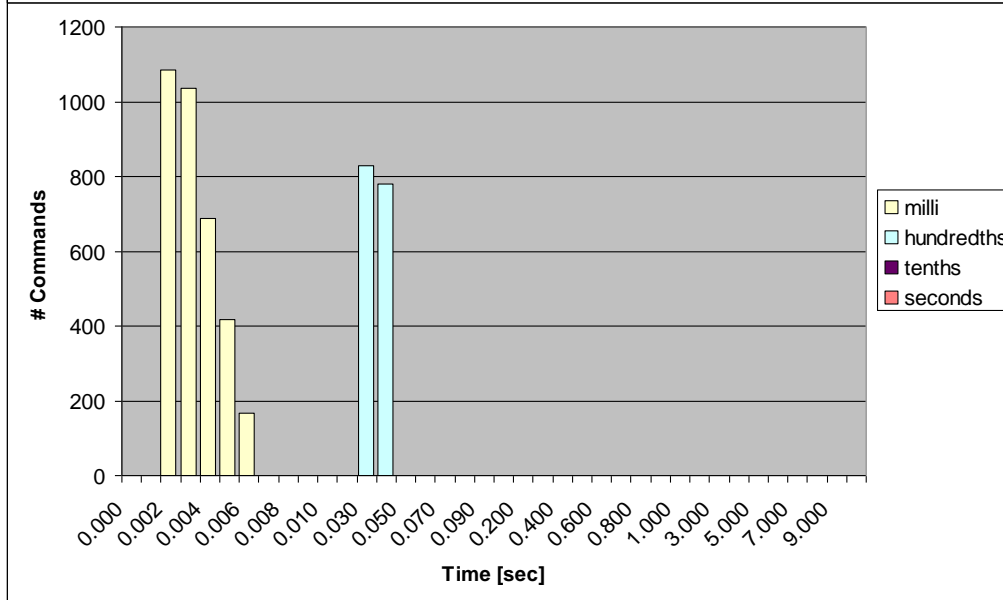
- S-Block destage average # invalidated exceptions **12.52**
- S-Block defrag average copy count **0.00**
- average IOPS = **121.54**

# Effect of S-Block Choice on Performance

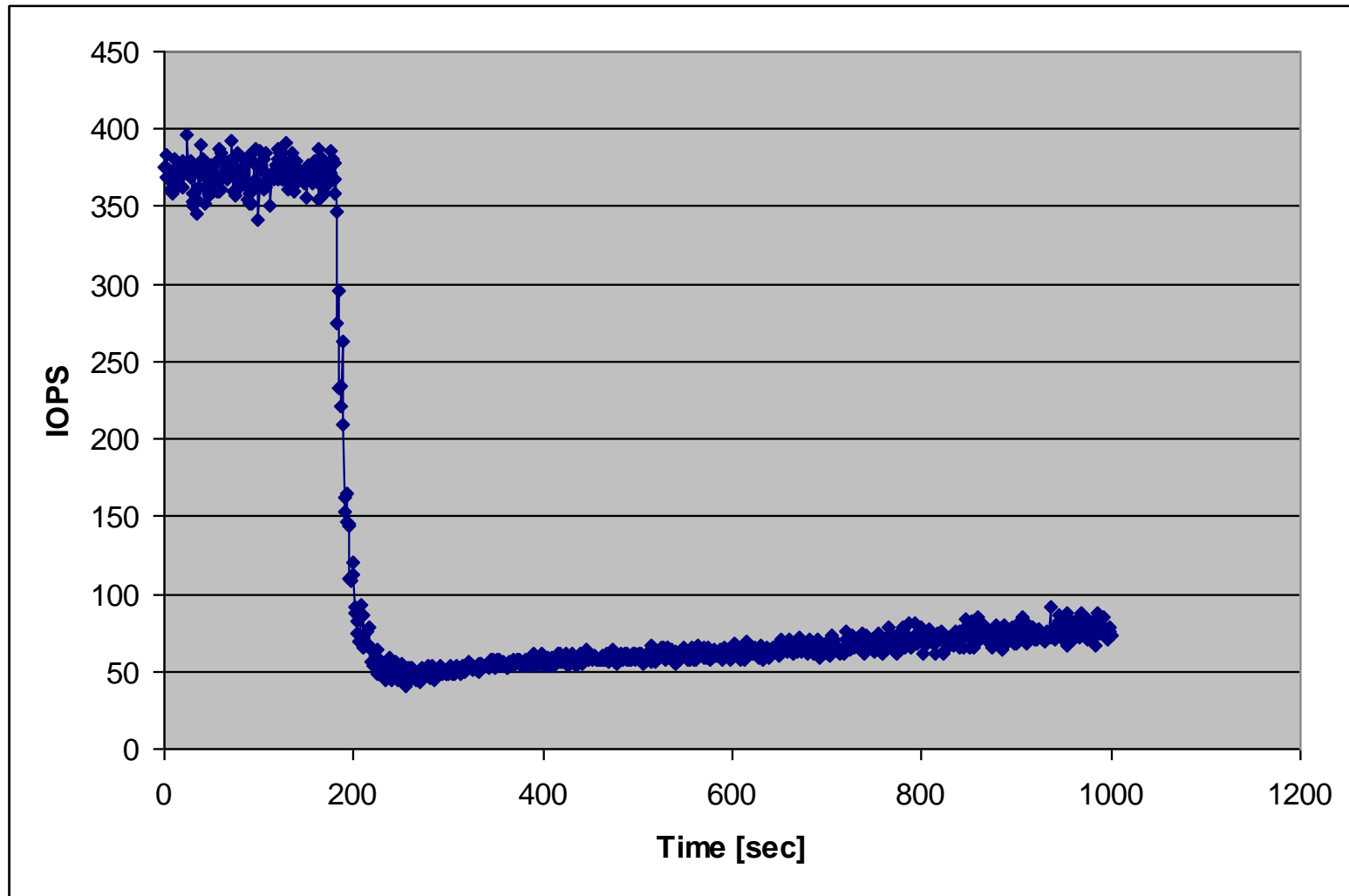
Optimal Destage



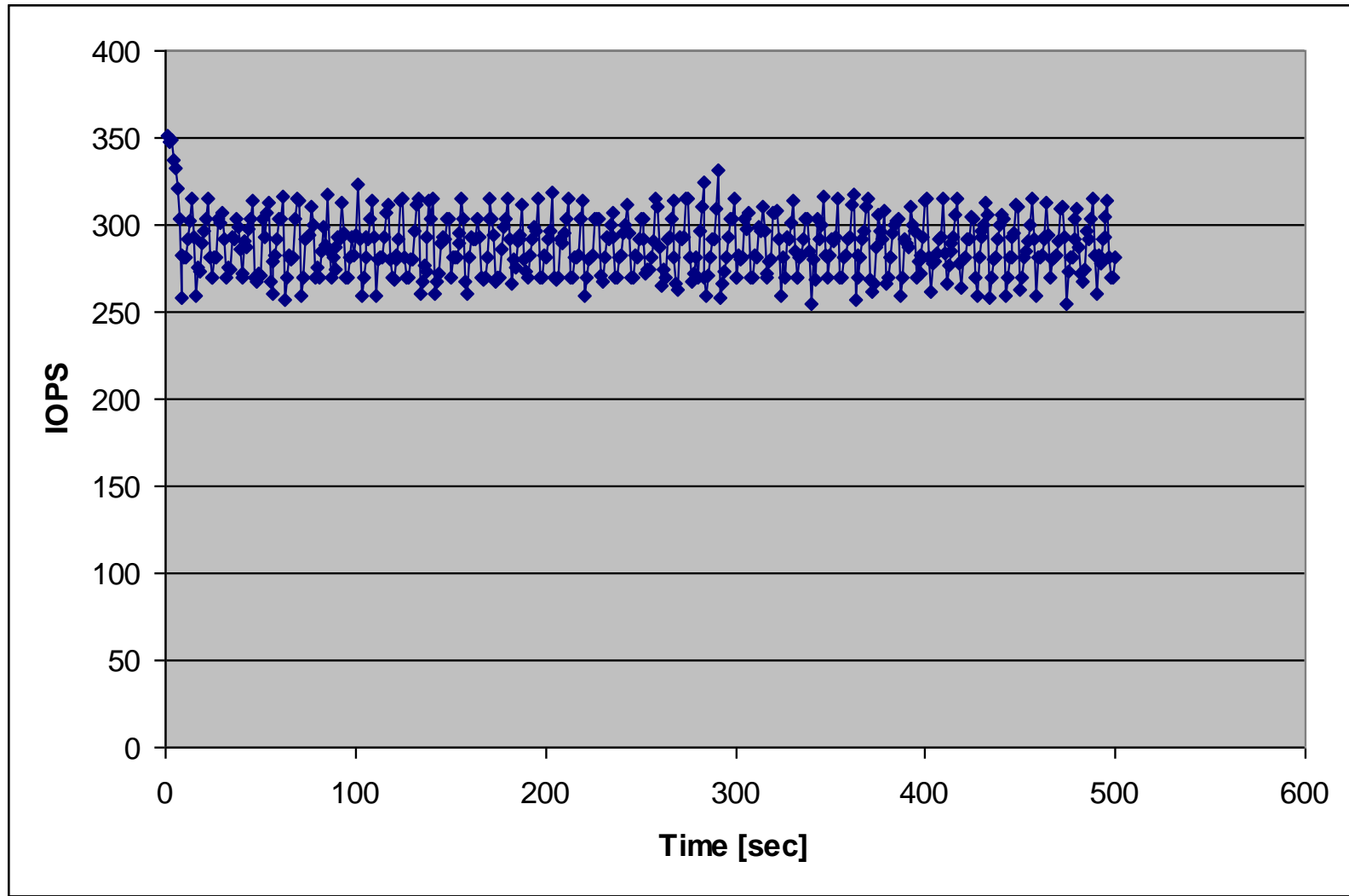
Optimal Defrag







# Constant Garbage Collection



- Good sustained random write performance
  - Append update → no consistency issues
  - Flexible to handle different workload types
  - Good sequential performance with direct S-block writes
- 
- Non-trivial implementation

- Shingled magnetic recording opens a rich area of systems research
- Good understanding of main issues and tradeoffs
- Proposed architectures likely basis for real implementations
- Significant research needed in performance optimization

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**Thank YOU !**

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

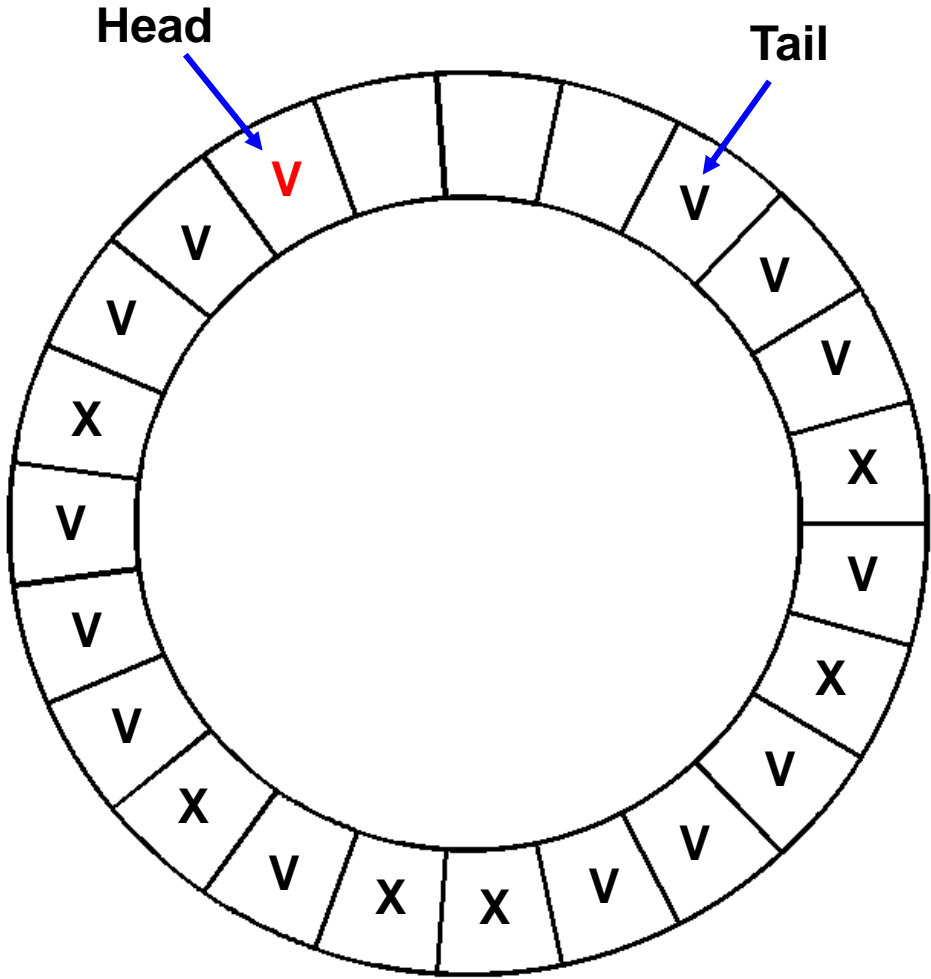
**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

**HITACHI**  
Inspire the Next

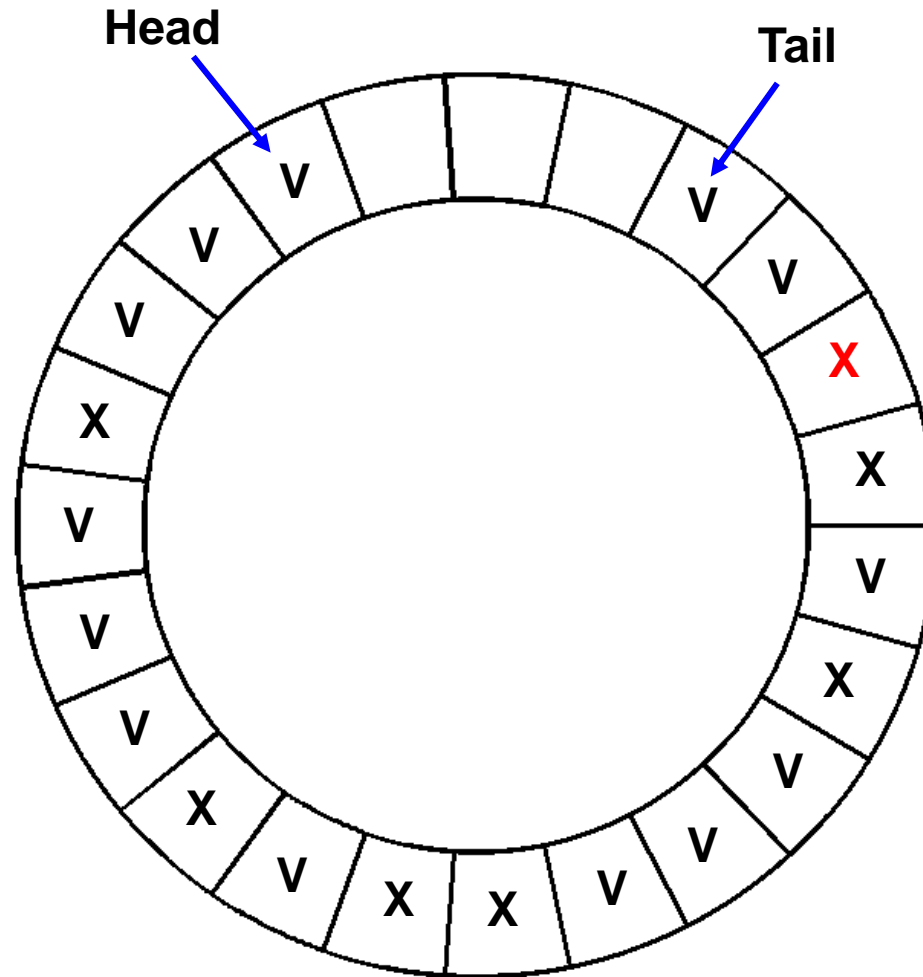




**V** Used, valid

**X** Used, invalid

**□** Unused

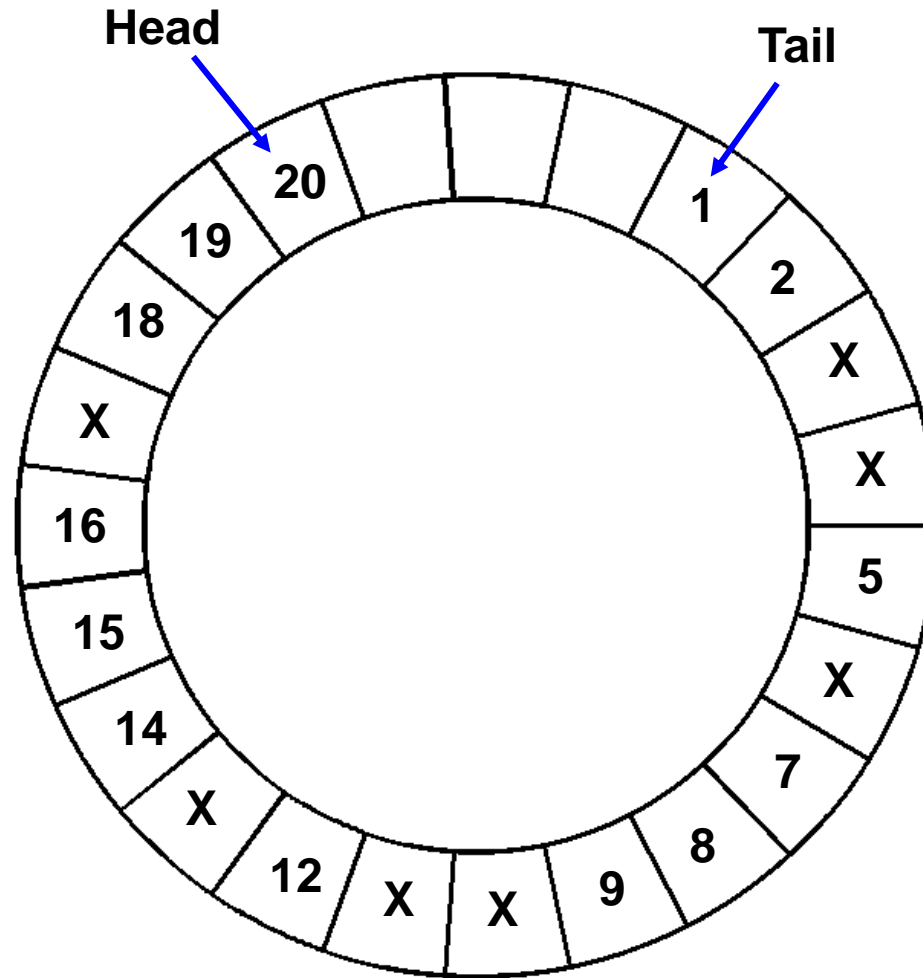


**V** Used, valid

**X** Used, invalid

**□** Unused



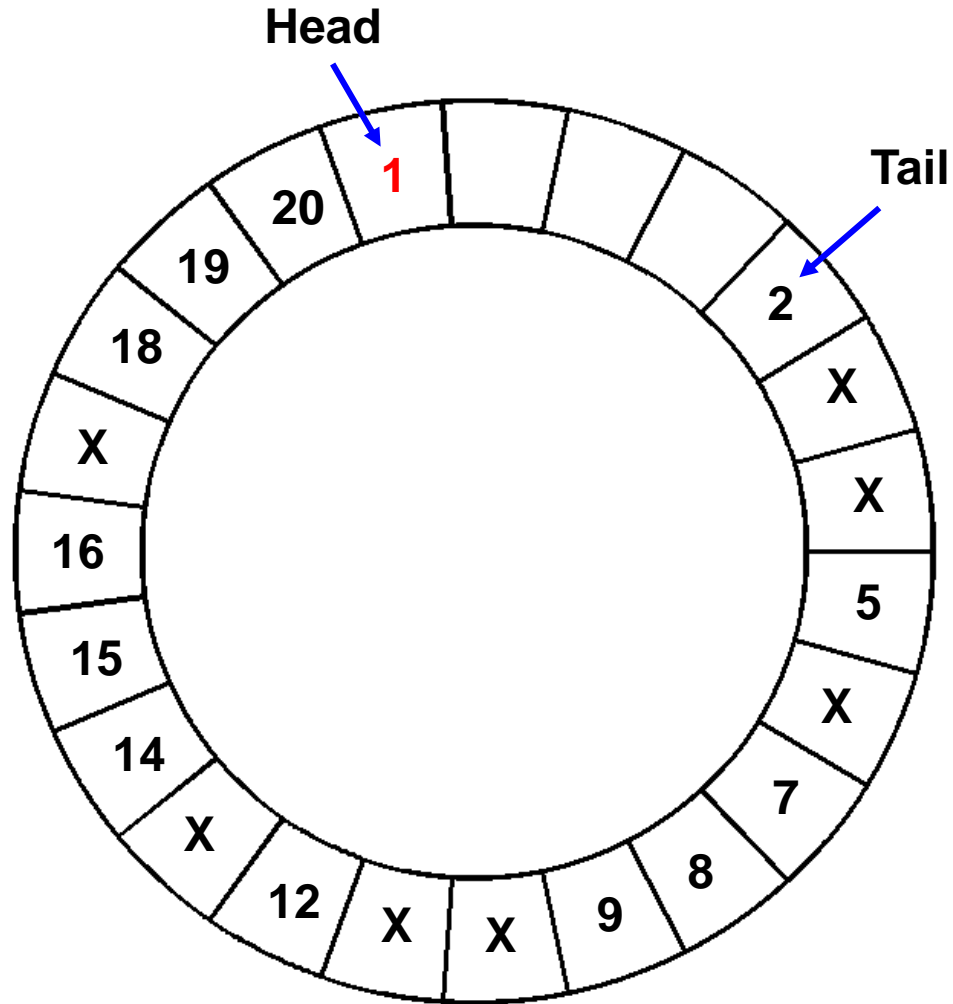


# Used, valid

X Used, invalid

Unused

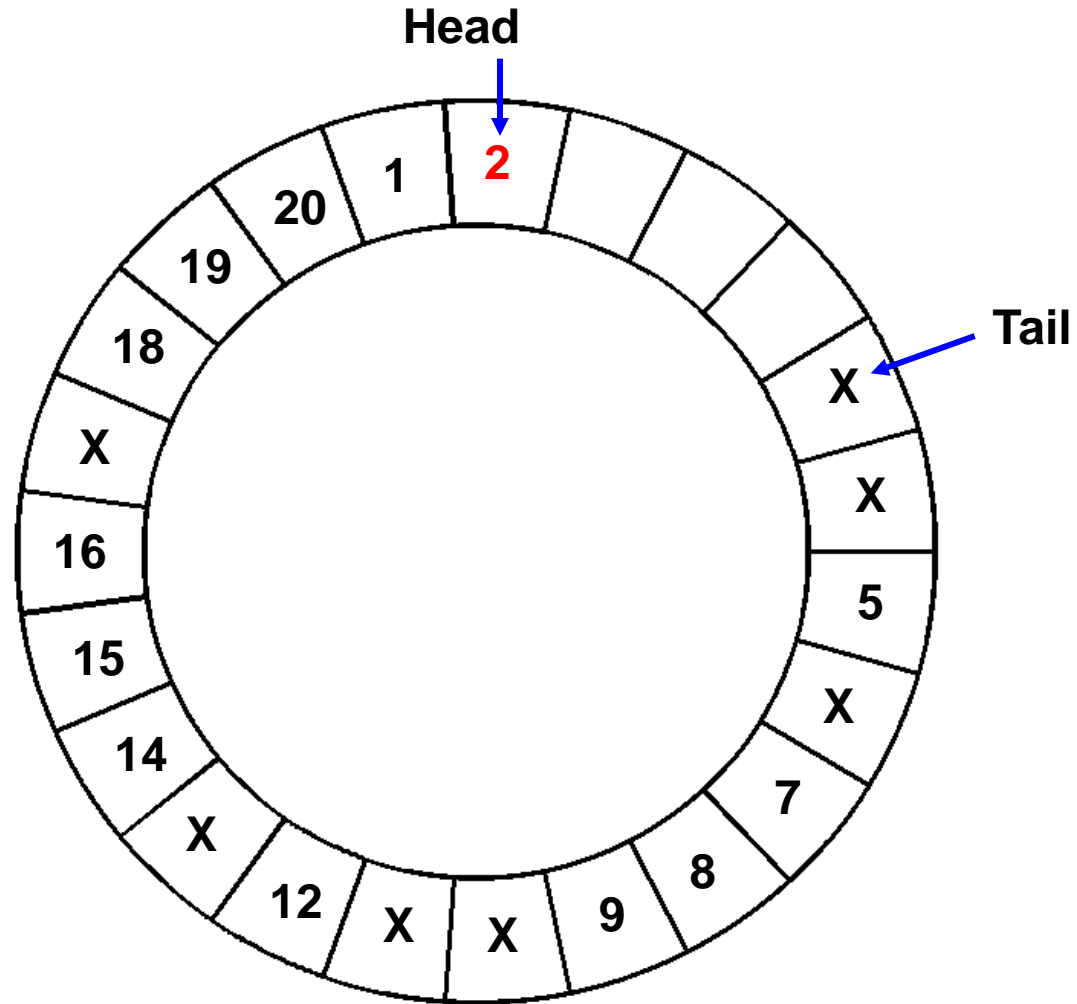
# Defrag Ring Buffer



# Used, valid

X Used, invalid

Unused



# Used, valid

X Used, invalid

Unused

# Defrag Ring Buffer

