Advanced Storage

(Flash translation layers, SSD garbage-collection)

Niv Dayan - CSC2525: Research Topics in Database Management



SSDs & FTLs



Zoned
Namespaces /
KV SSDs



Spooky

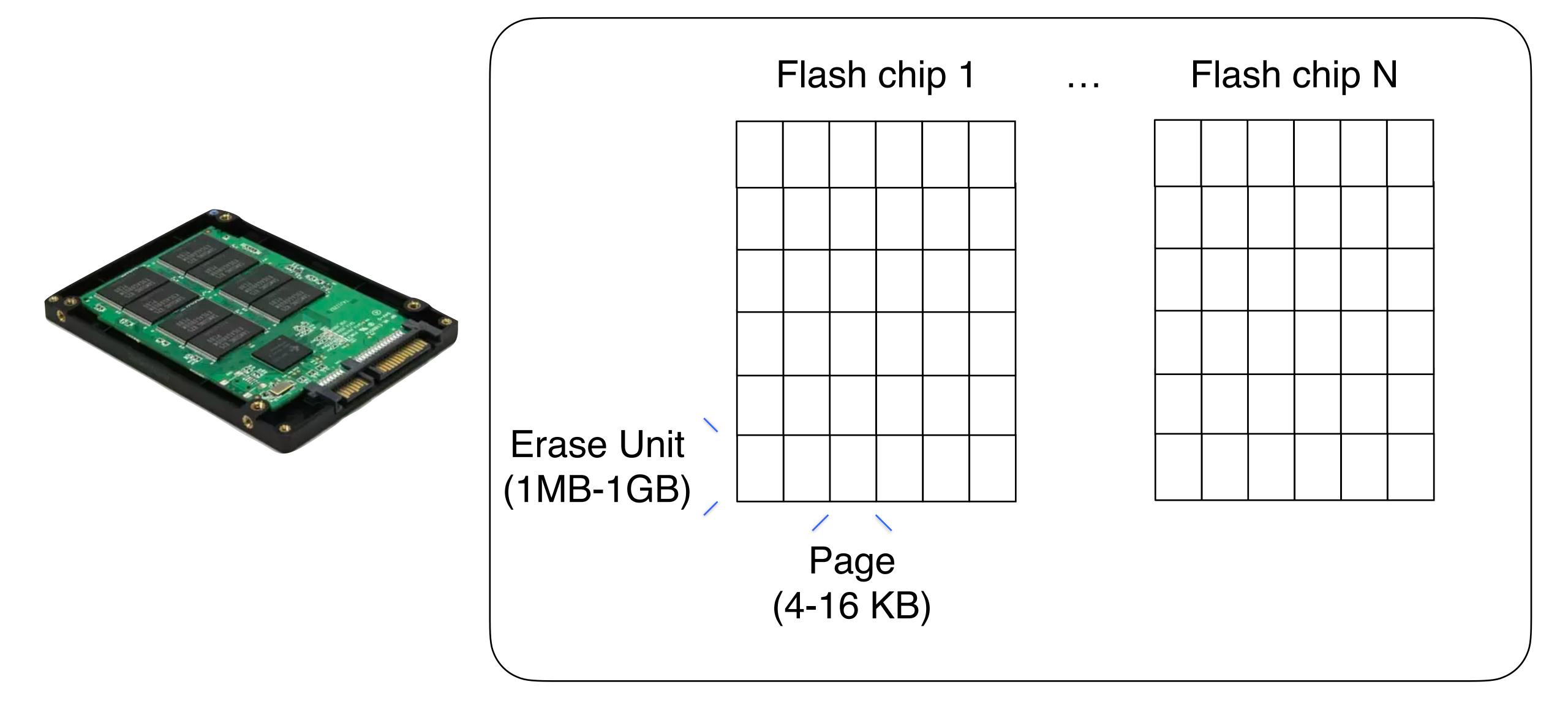


Shingled Magnetic Disks

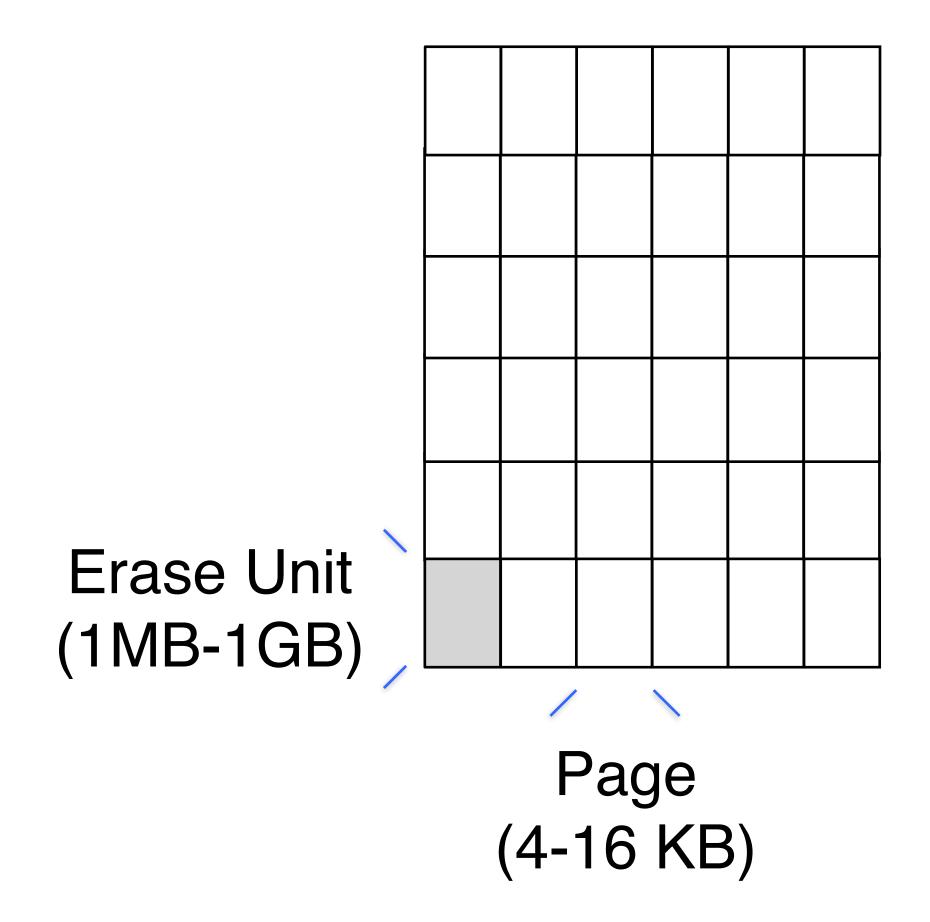
SSD Review



SSD Review

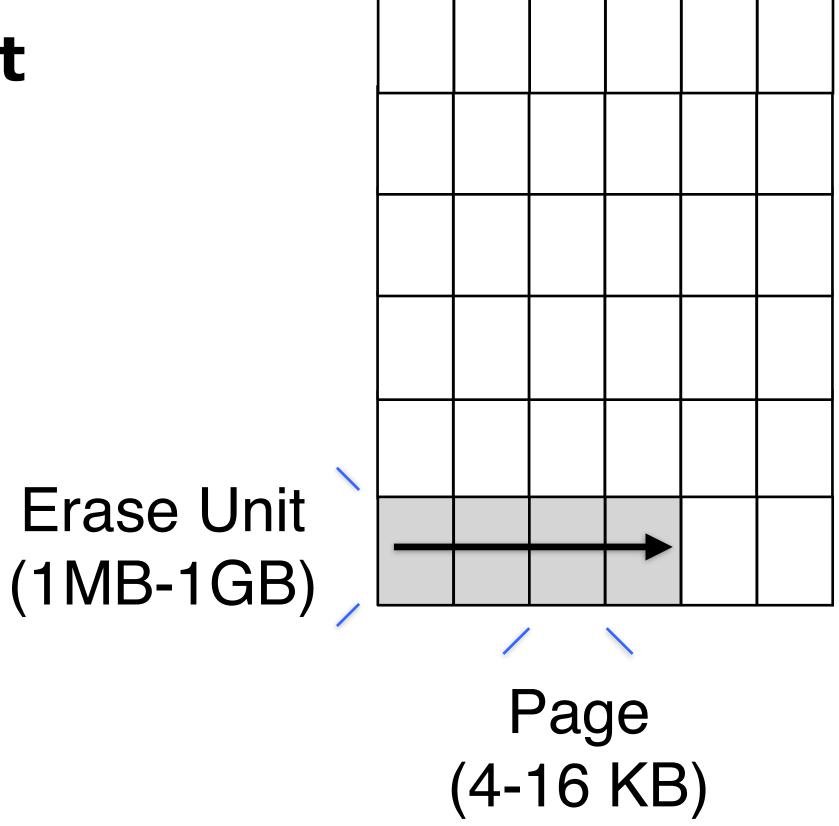


Flash chip



Flash chip

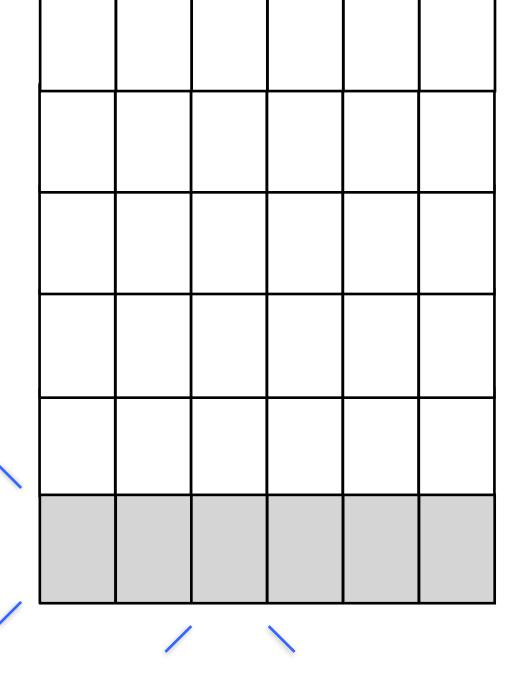
Pages must be written sequentially in an erase unit



Pages must be written sequentially in an erase unit

All data in an erase unit is erased at the same time

Flash chip



Page (4-16 KB)

Erase Unit

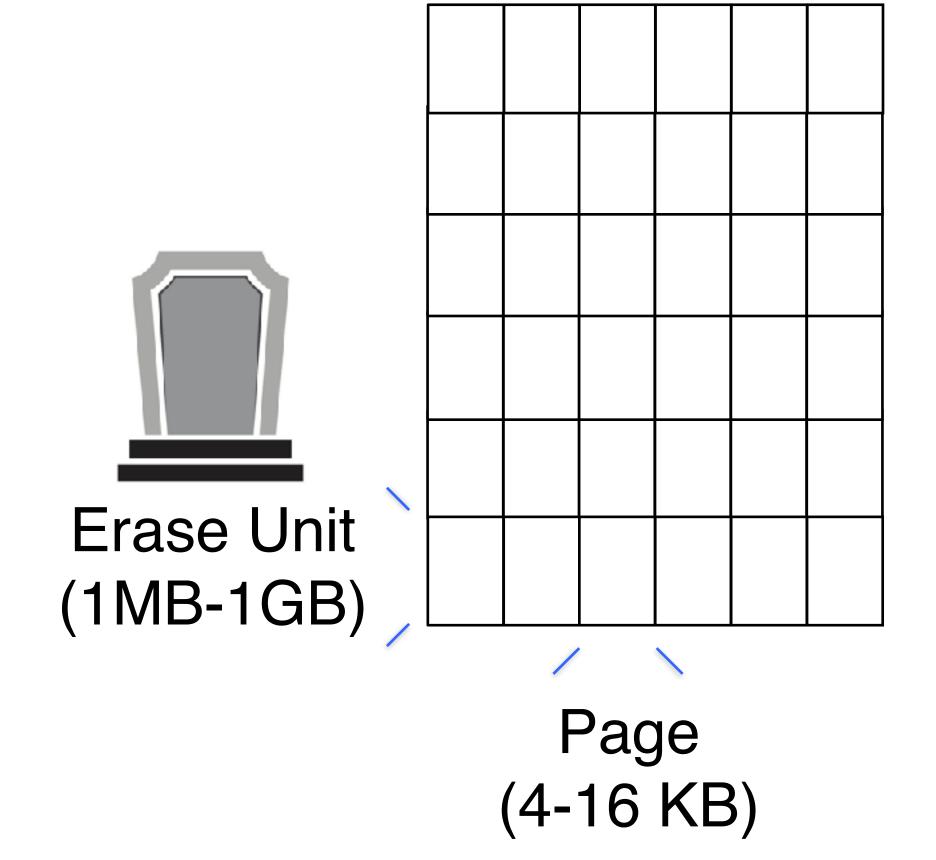
(1MB-1GB)

Pages must be written sequentially in an erase unit

All data in an erase unit is erased at the same time

Each erase unit has a lifetime (1-10K erases)

Flash chip



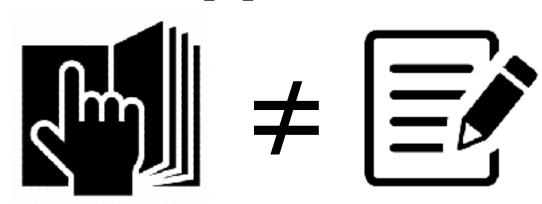
Pages must be written sequentially in an erase unit

All data in an erase unit is erased at the same time

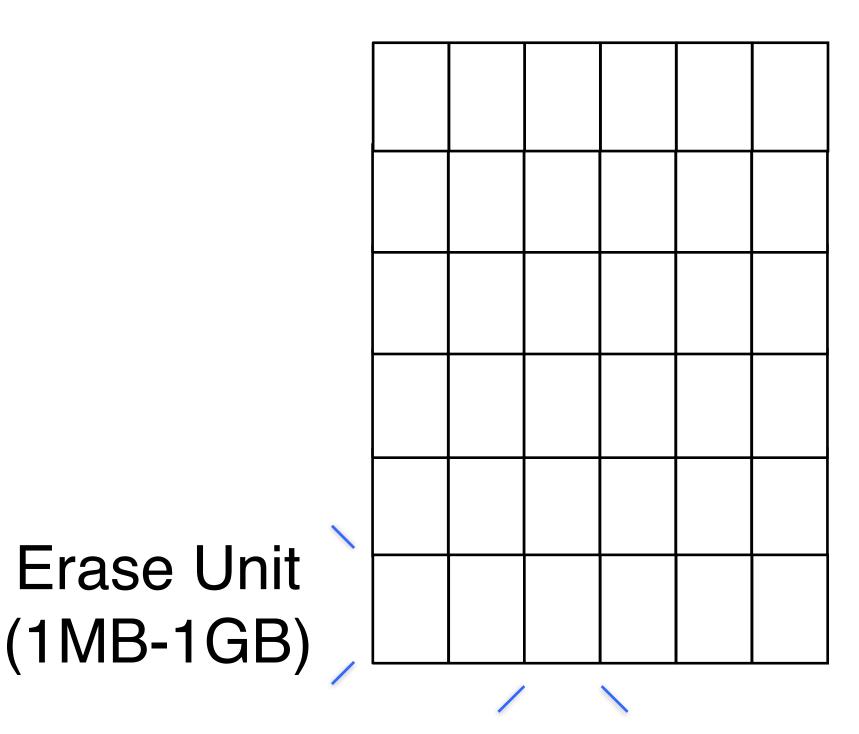
Each erase unit has a lifetime (1-10K erases)

Reading a page takes approx 50 us

Writing a page takes approx 100-200 us



Flash chip

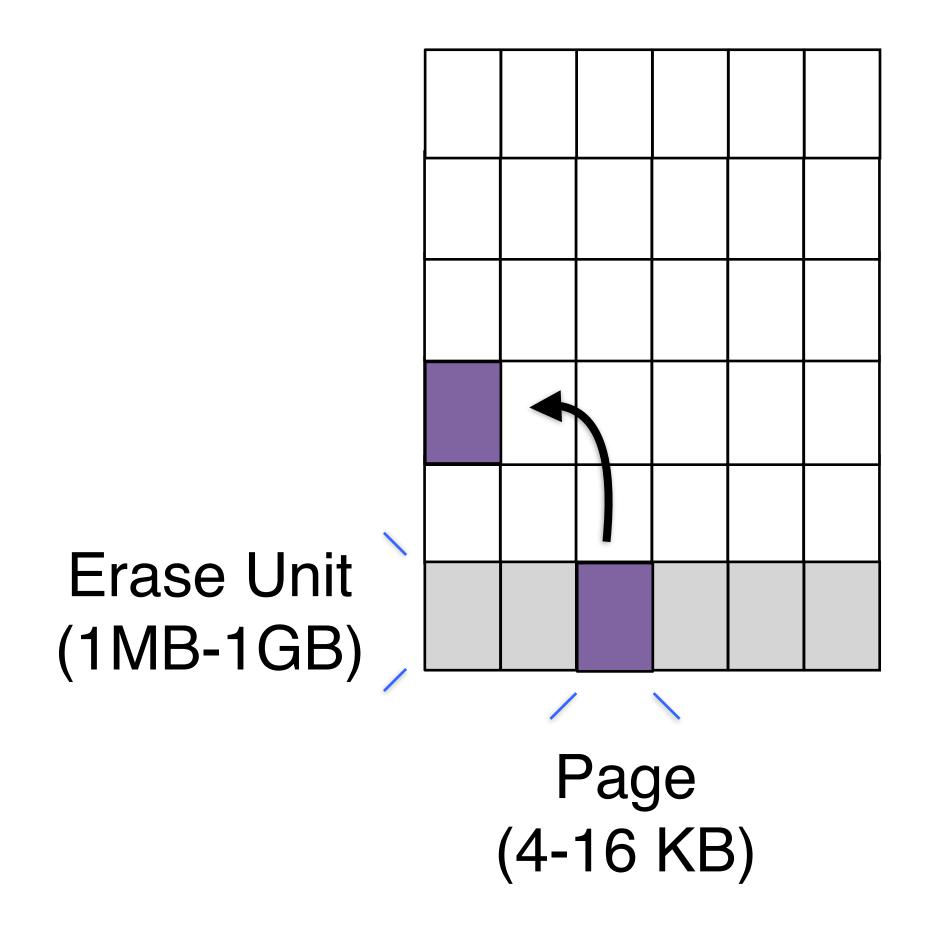


Page (4-16 KB)

Updating a Page Out-of-Place

Copy updated page to erase unit with free space

Flash chip

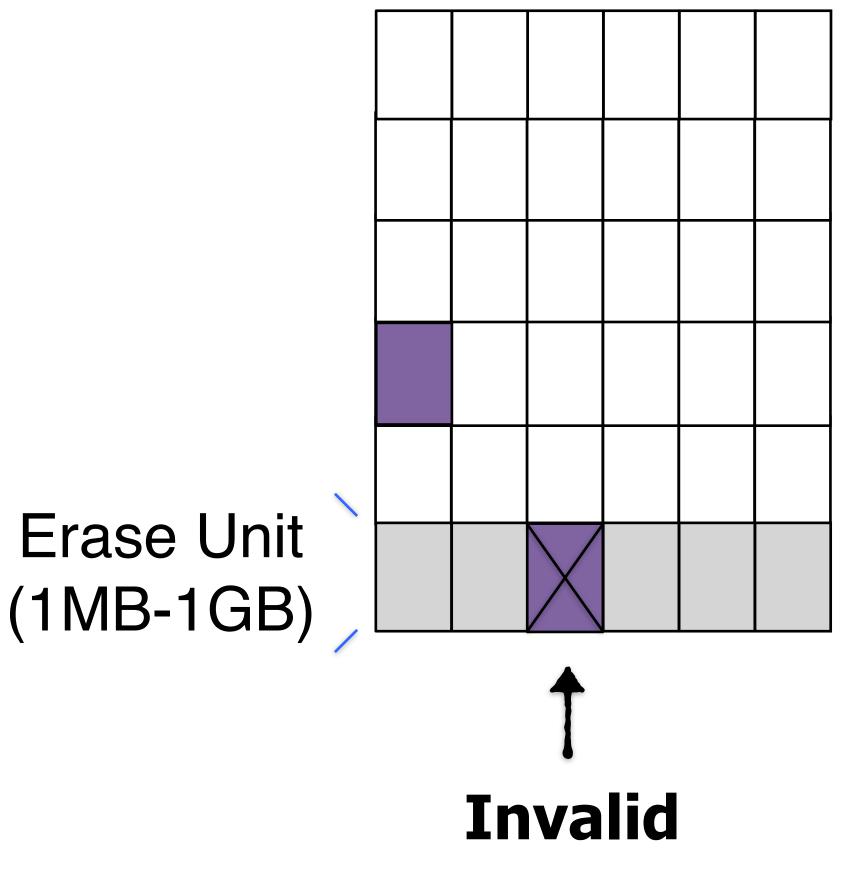


Updating a Page Out-of-Place

Copy updated page to erase unit with free space

Flash chip

Mark the original page as invalid (using a bitmap)



Garbage-Collection

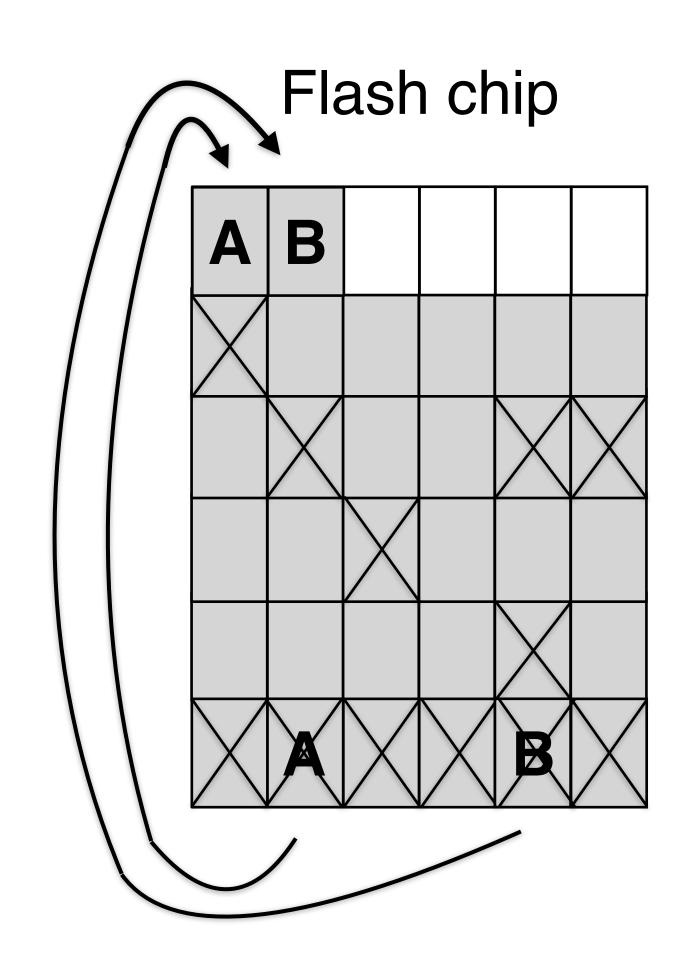
Flash chip

Find the erase unit with the least live data left.

Garbage-Collection

Find the erase unit with the least live data left.

Copy live pages to an erase unit with free space.



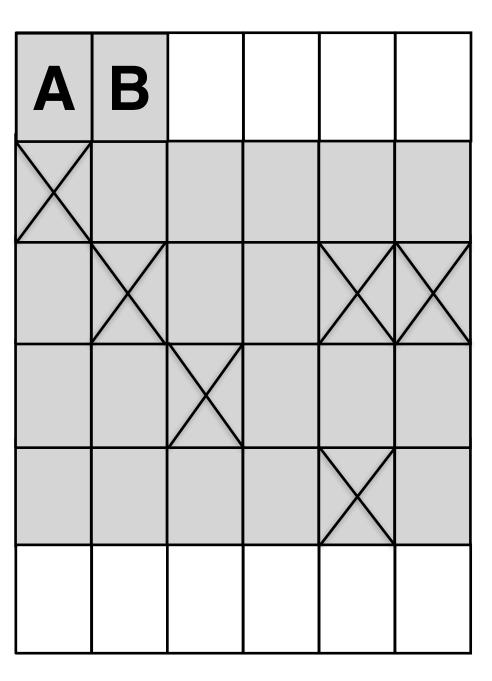
Garbage-Collection

Find the erase unit with the least live data left.

Copy live pages to an erase unit with free space.

Erase the block

Flash chip

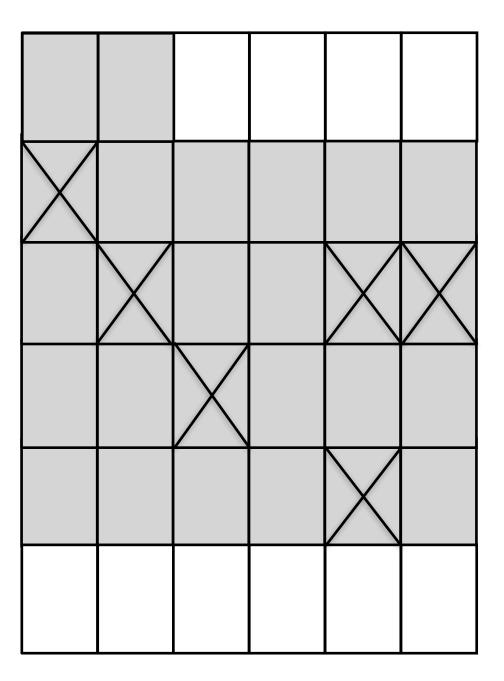


Over-Provisioning

SSD physical capacity must be greater than logical capacity

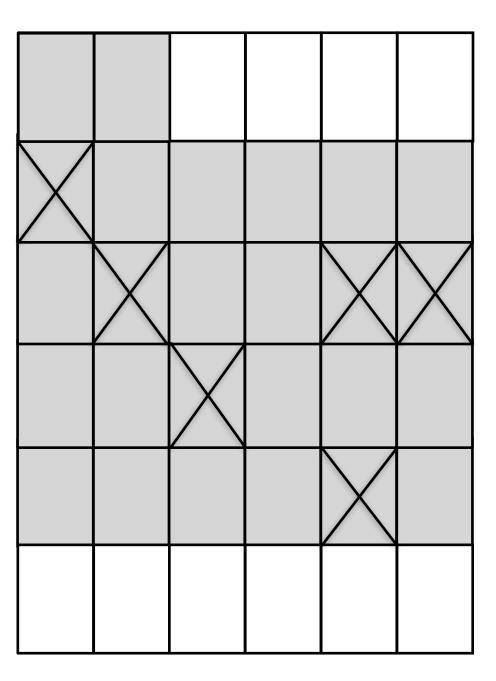
So invalid pages can accumulate to cheapen garbage-collection

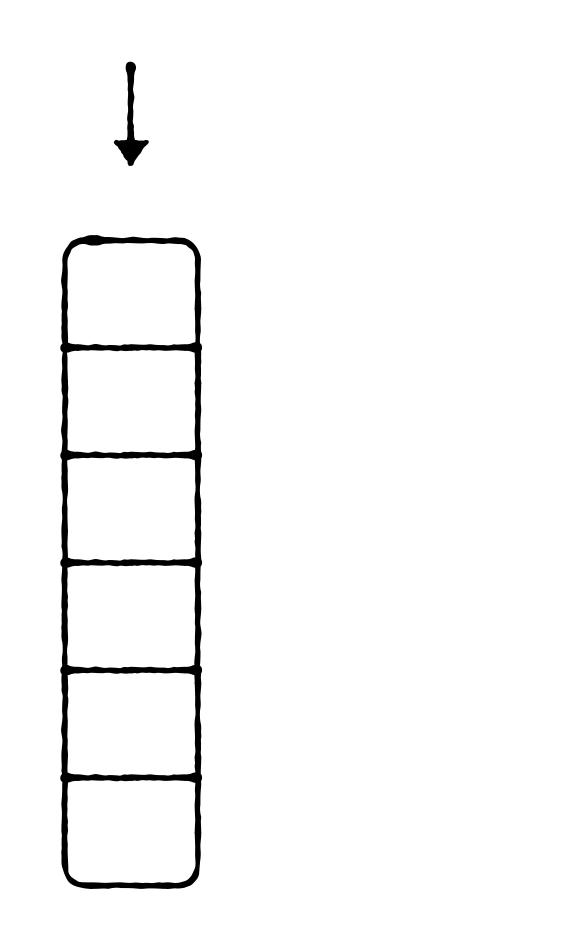
Flash chip



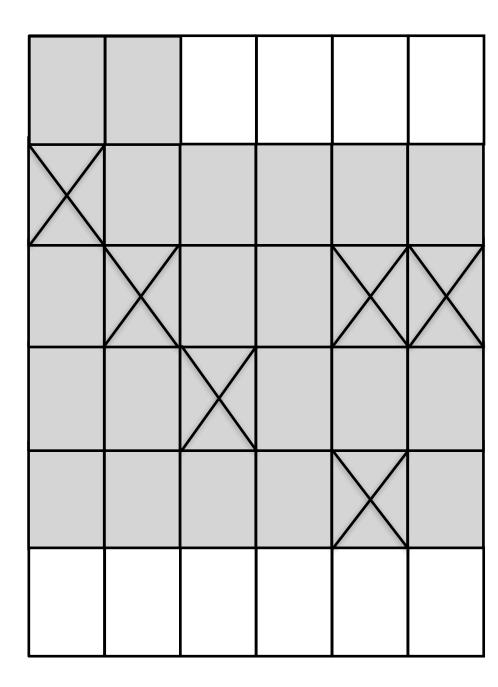
And now to new things:)

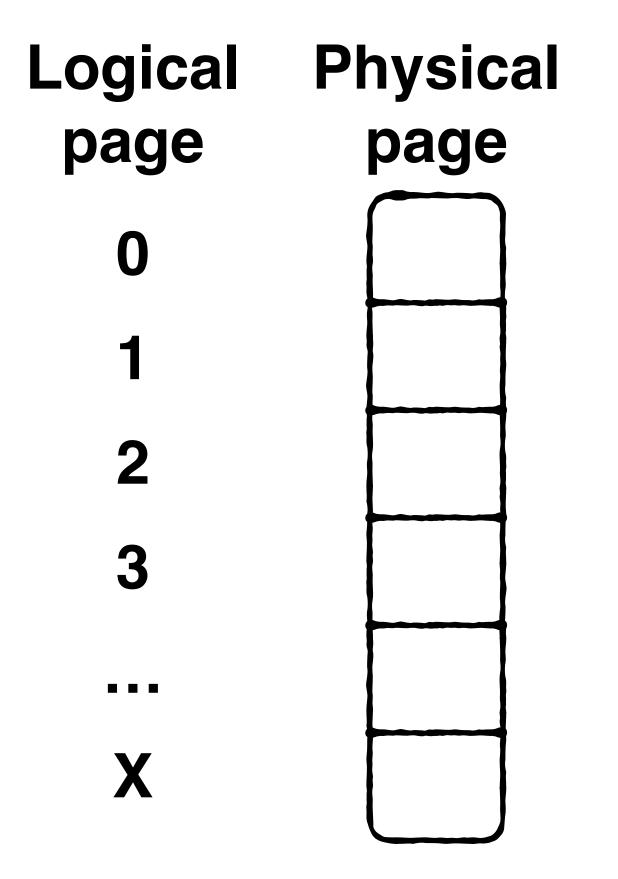
Flash chip





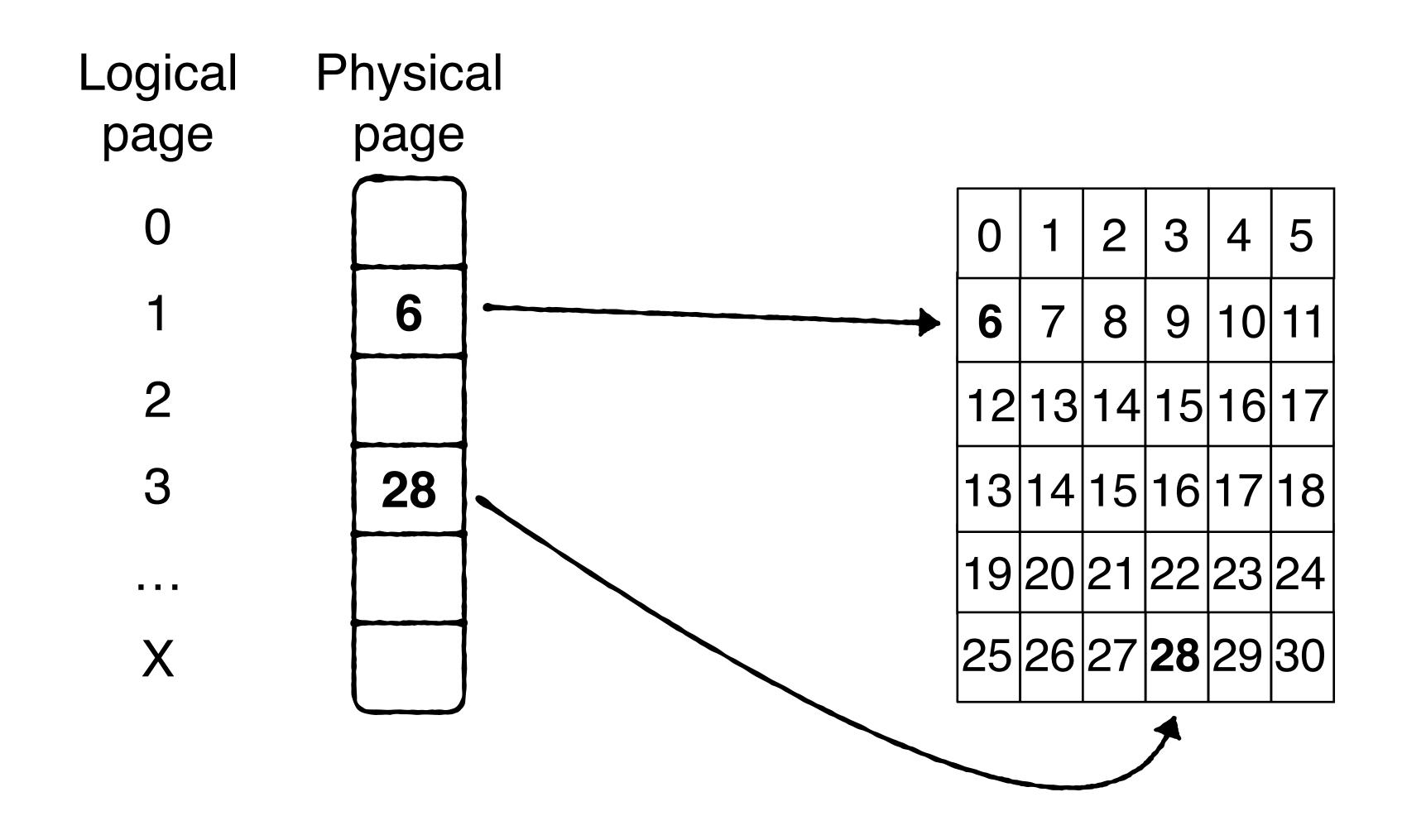
Flash chip

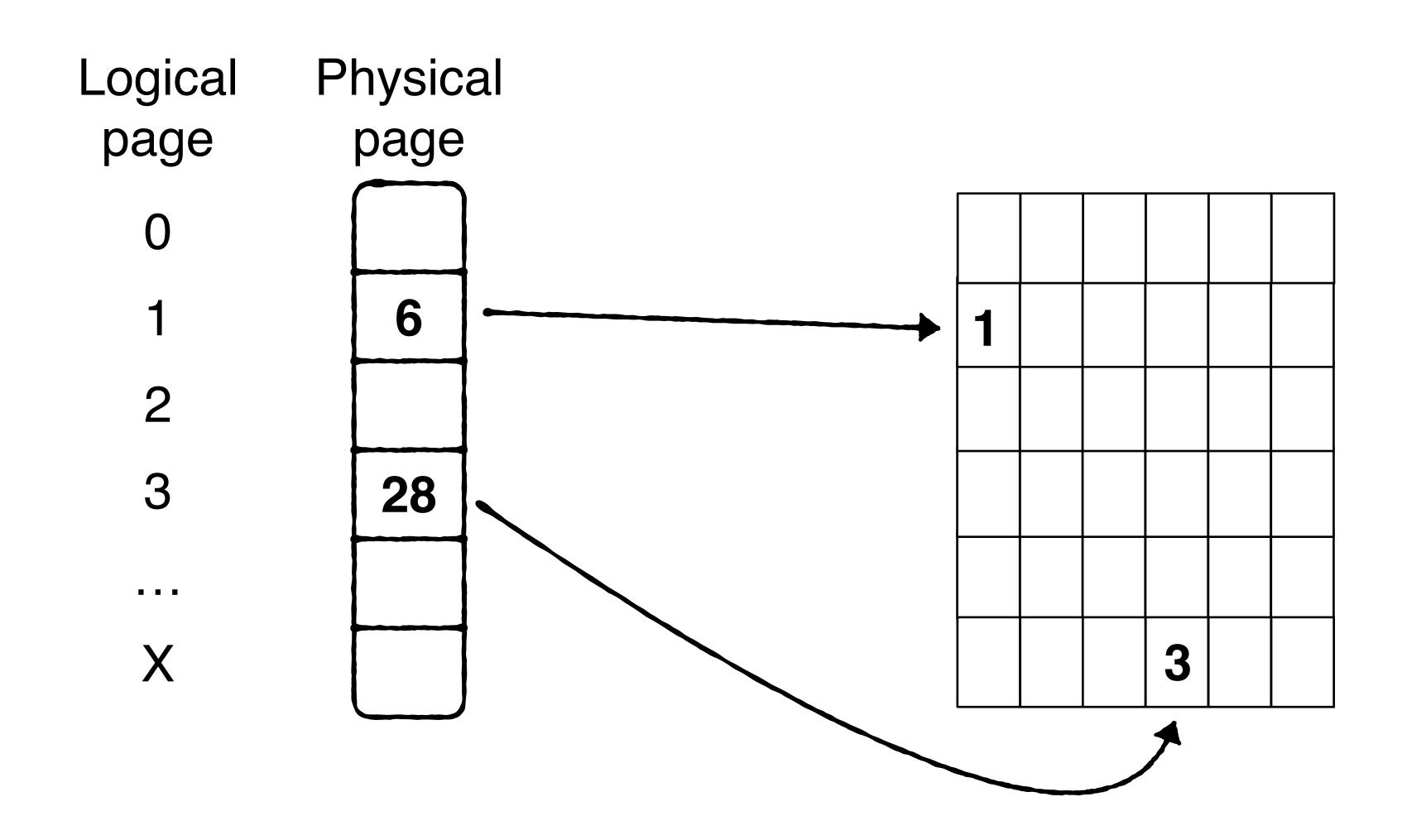




Flash chip

0	1	2	3	4	5
6	7	8	တ	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30





Flash translation layer (FTL)

mapping table

Garbagecollection Wear-Leveling Error-Correction

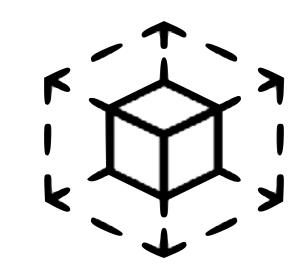
Over-Provisioning







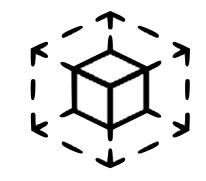








Flash translation layer (FTL)



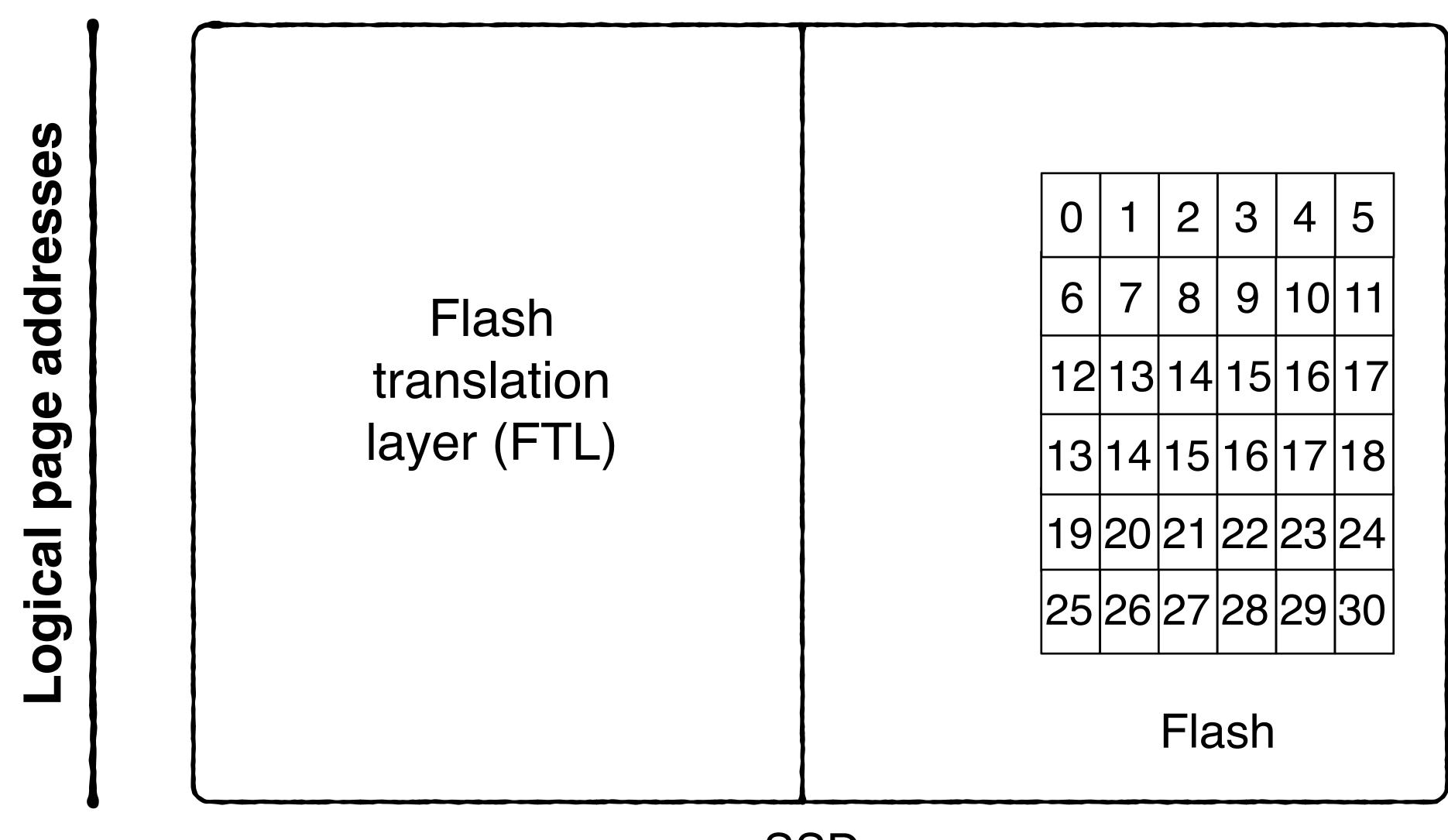




0	1	2	3	4	5
6	7	8	တ	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30

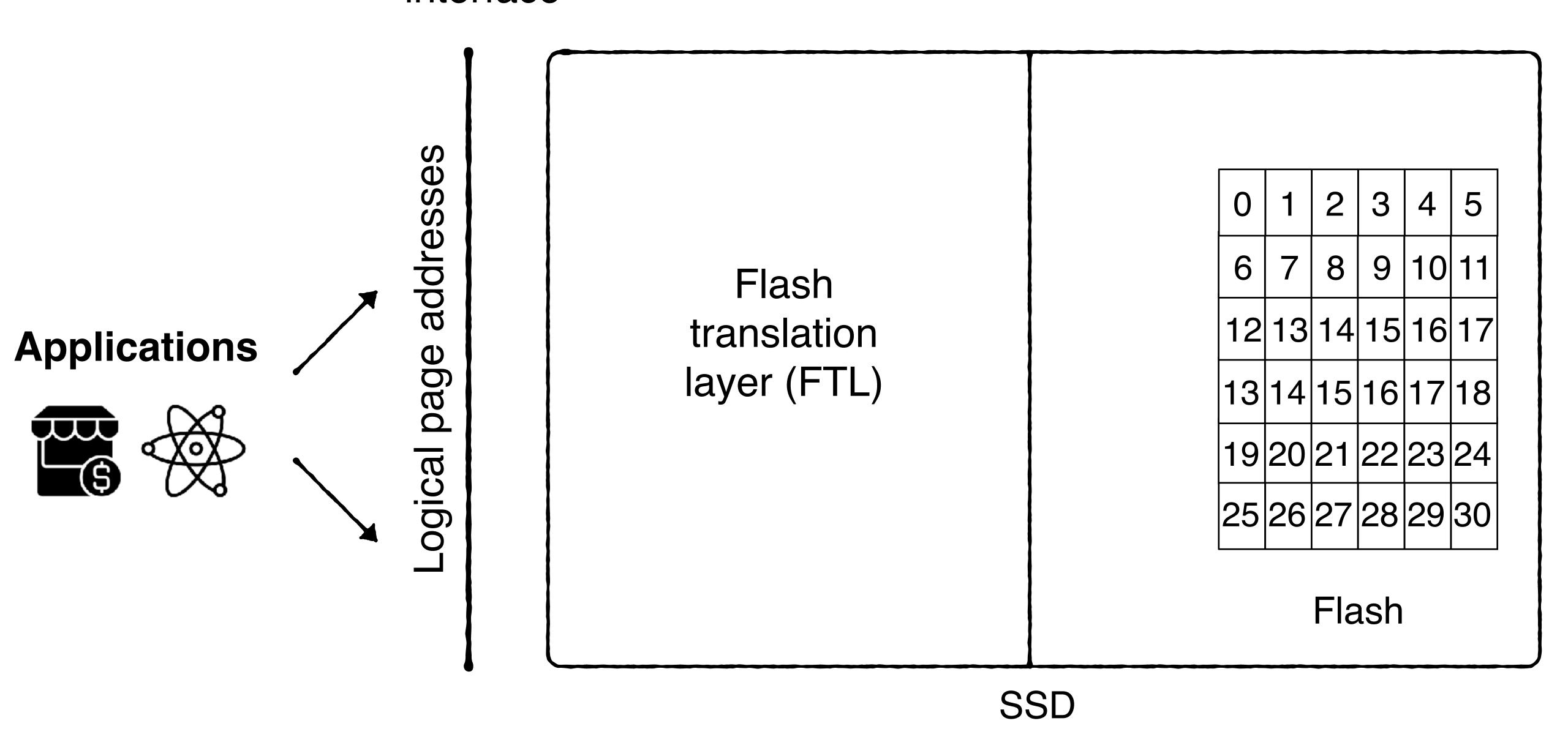
Flash

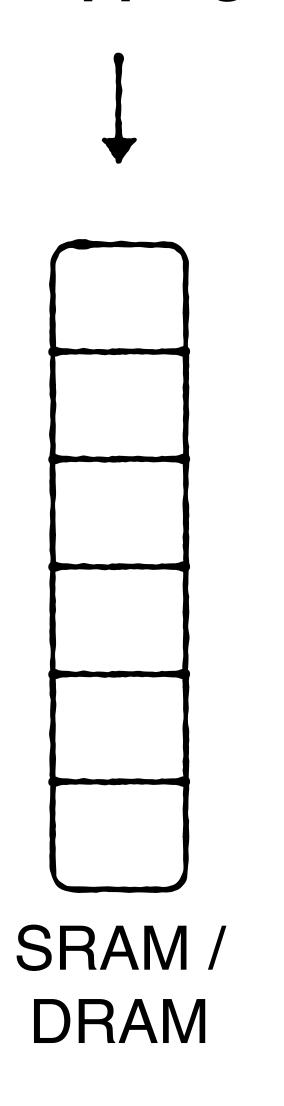
Block device interface



SSD

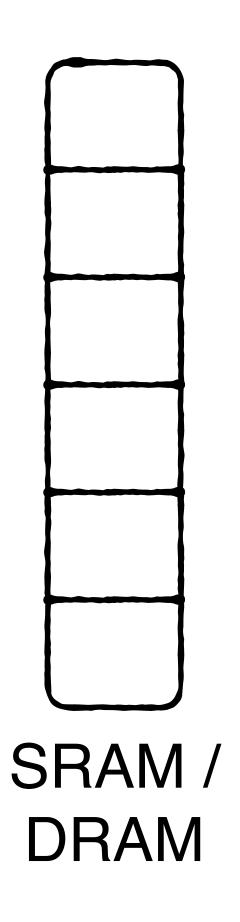
Block device interface





0	1	2	ര	4	5
6	7	8	9	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30

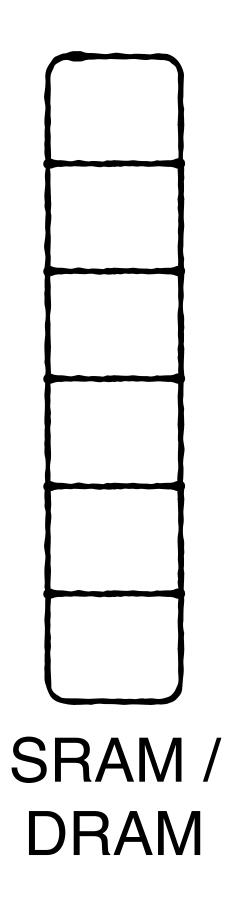
e.g., 1TB SSD with 4KB pages requires 1GB mapping table



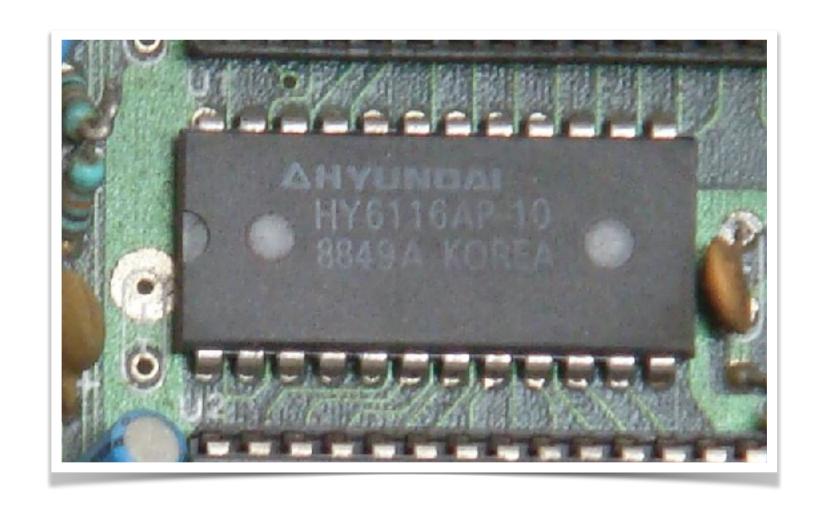
0	1	2	3	4	5
6	7	8	9	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30

e.g., 1TB SSD with 4KB pages requires 1GB mapping table

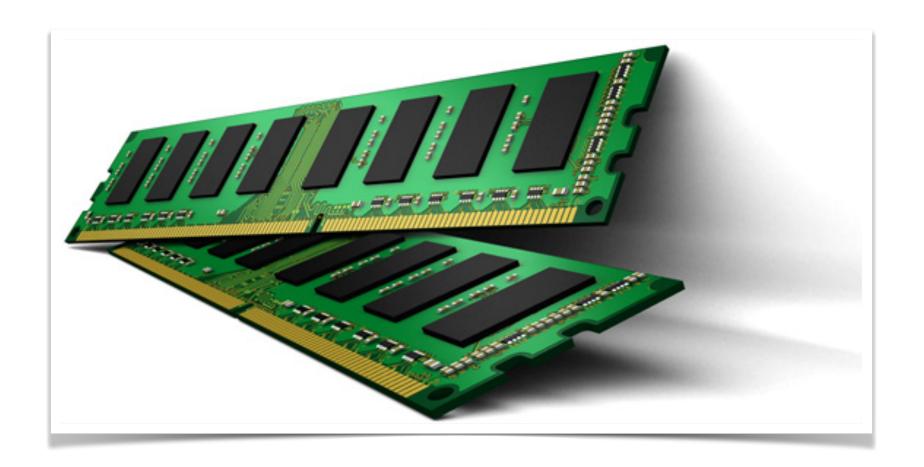
Doesn't sound like a lot but...



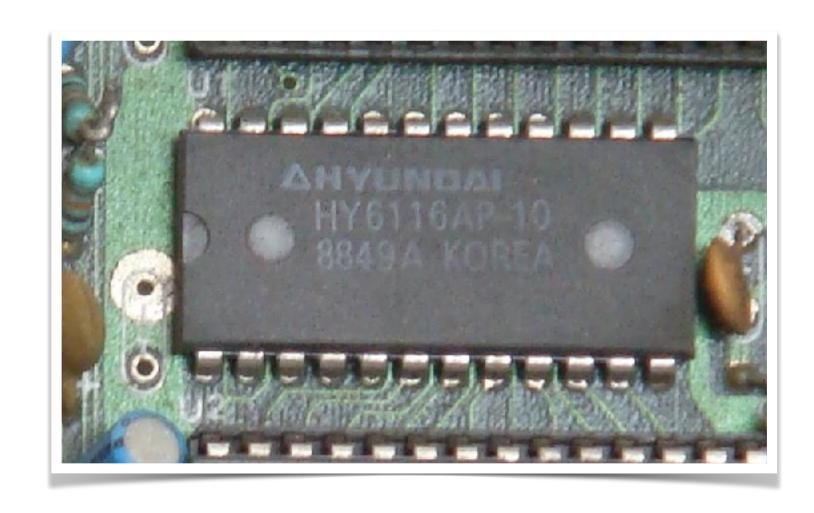
0	1	2	3	4	5
6	7	8	တ	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30



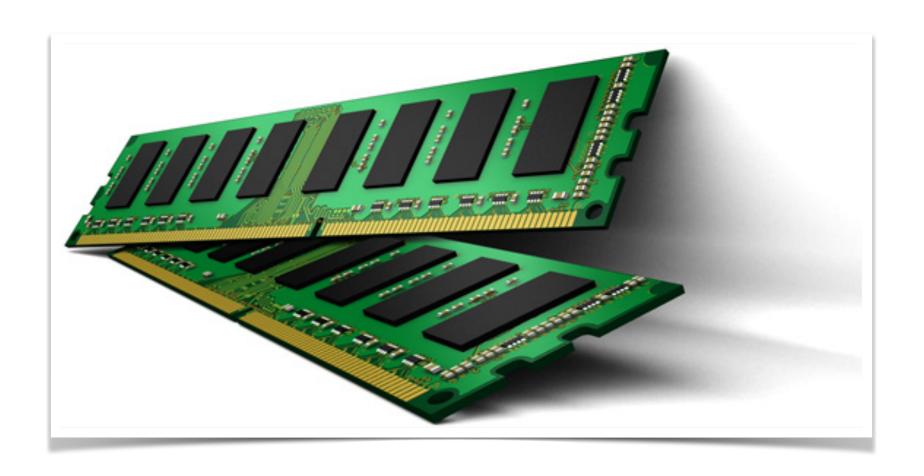
Static RAM (SRAM)



Dynamic RAM (DRAM)

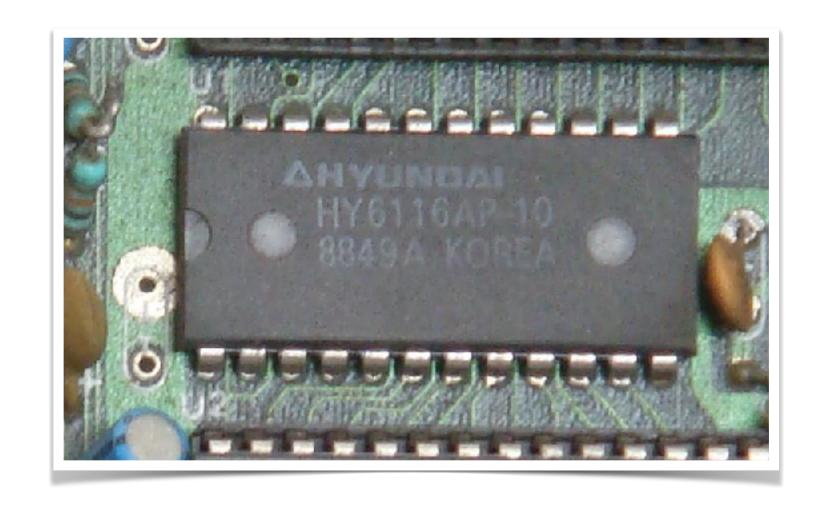


Static RAM (SRAM)
No refreshing



Dynamic RAM (DRAM)

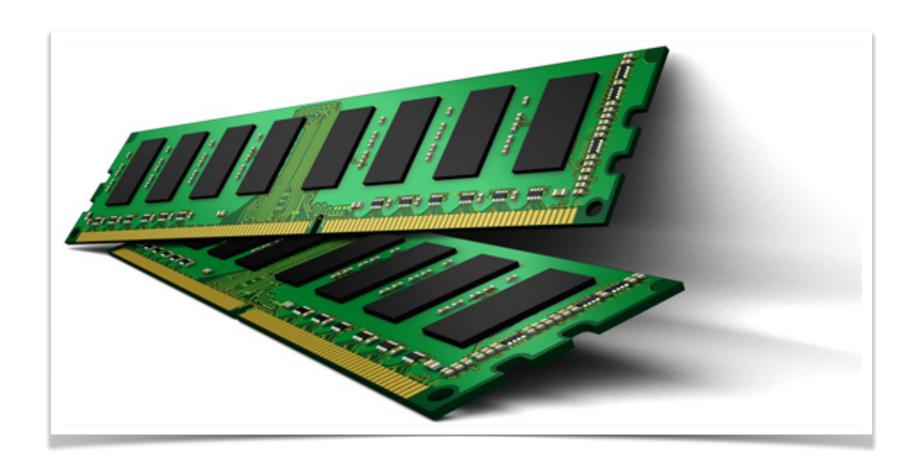
Requires refreshing even when idle



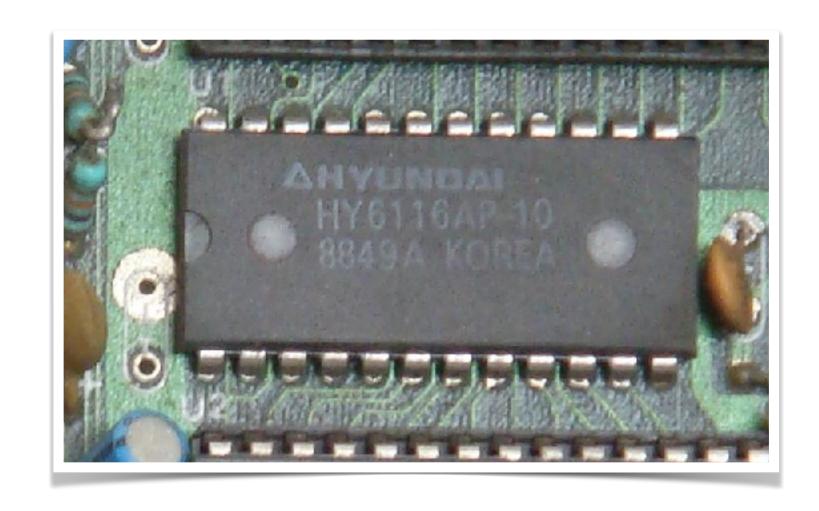
Static RAM (SRAM)

No refreshing

More energy efficient



Dynamic RAM (DRAM)
Requires refreshing even when idle
More power hungry

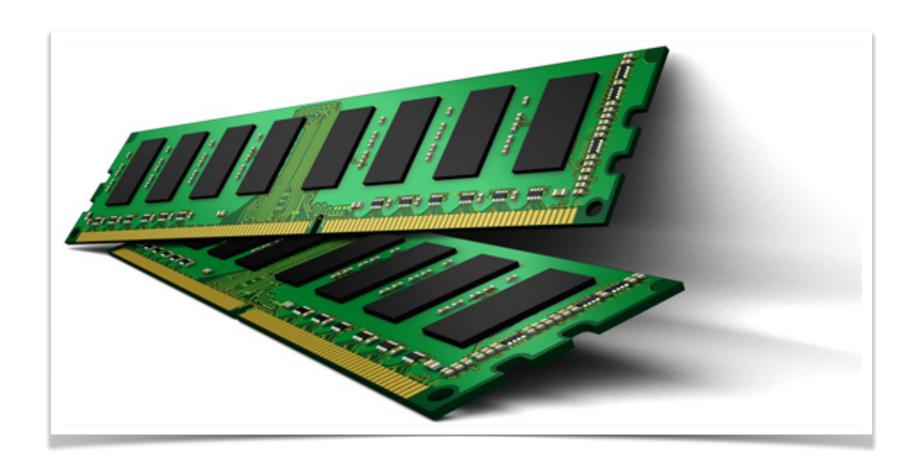


Static RAM (SRAM)

No refreshing

More energy efficient

Less error prone



Dynamic RAM (DRAM)

Requires refreshing even when idle

More power hungry

bit flips due to charge leakage (requires error correction codes)



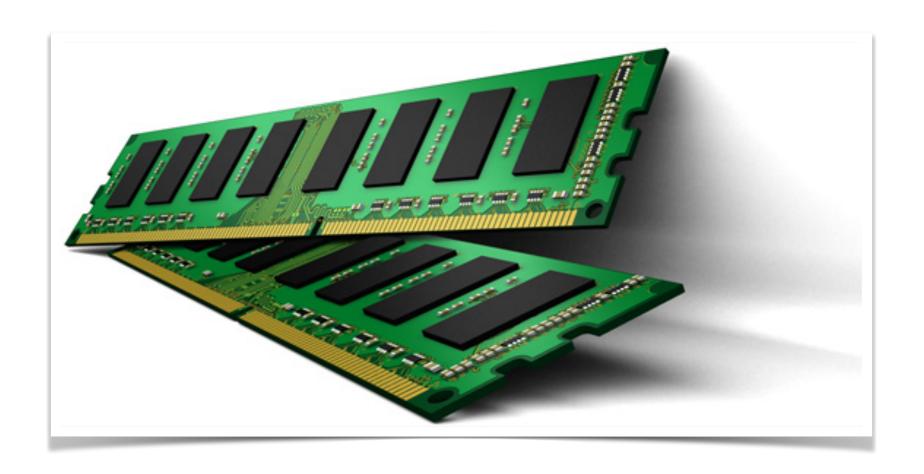
Static RAM (SRAM)

No refreshing

More energy efficient

Less error prone

Longer lifespan



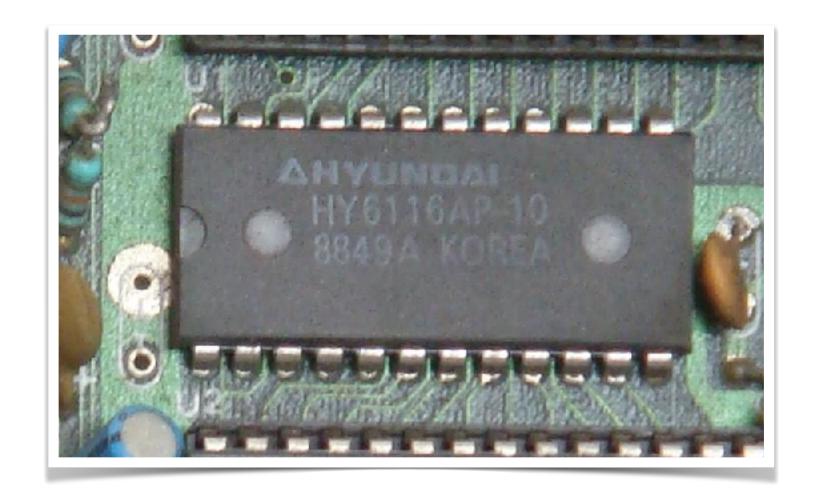
Dynamic RAM (DRAM)

Requires refreshing even when idle

More power hungry

bit flips due to charge leakage (requires error correction codes)

Shorter lifespan



Static RAM (SRAM)

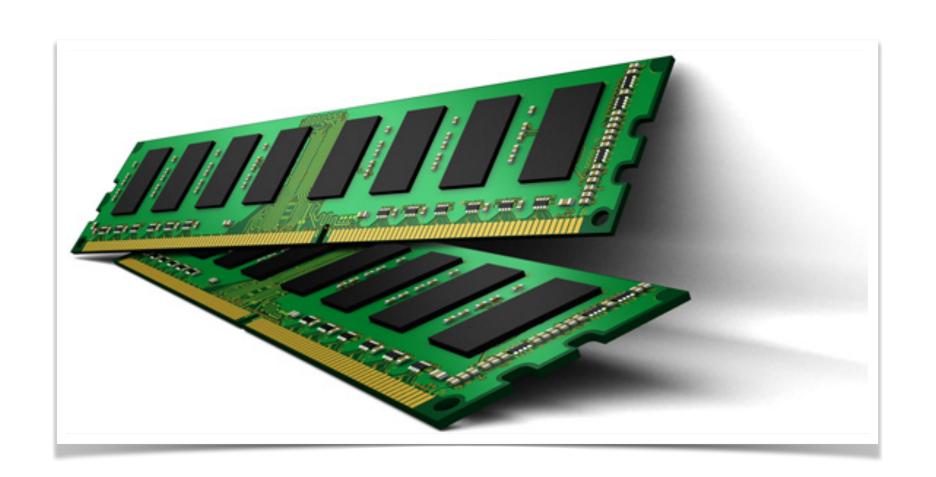
No refreshing

More energy efficient

Less error prone

Longer lifespan

\$1000 per GB



Dynamic RAM (DRAM)

Requires refreshing even when idle

More power hungry

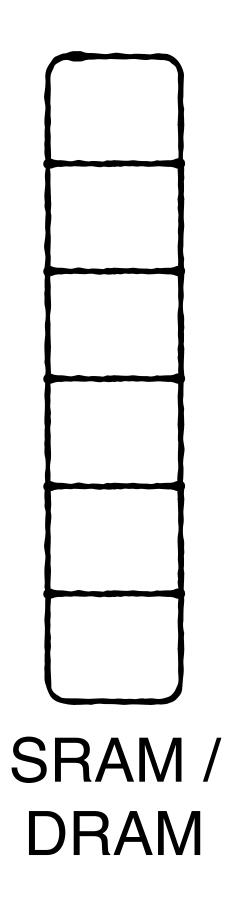
bit flips due to charge leakage (requires error correction codes)

Shorter lifespan

\$5 per GB

e.g., 1TB SSD with 4KB pages requires 1GB mapping table

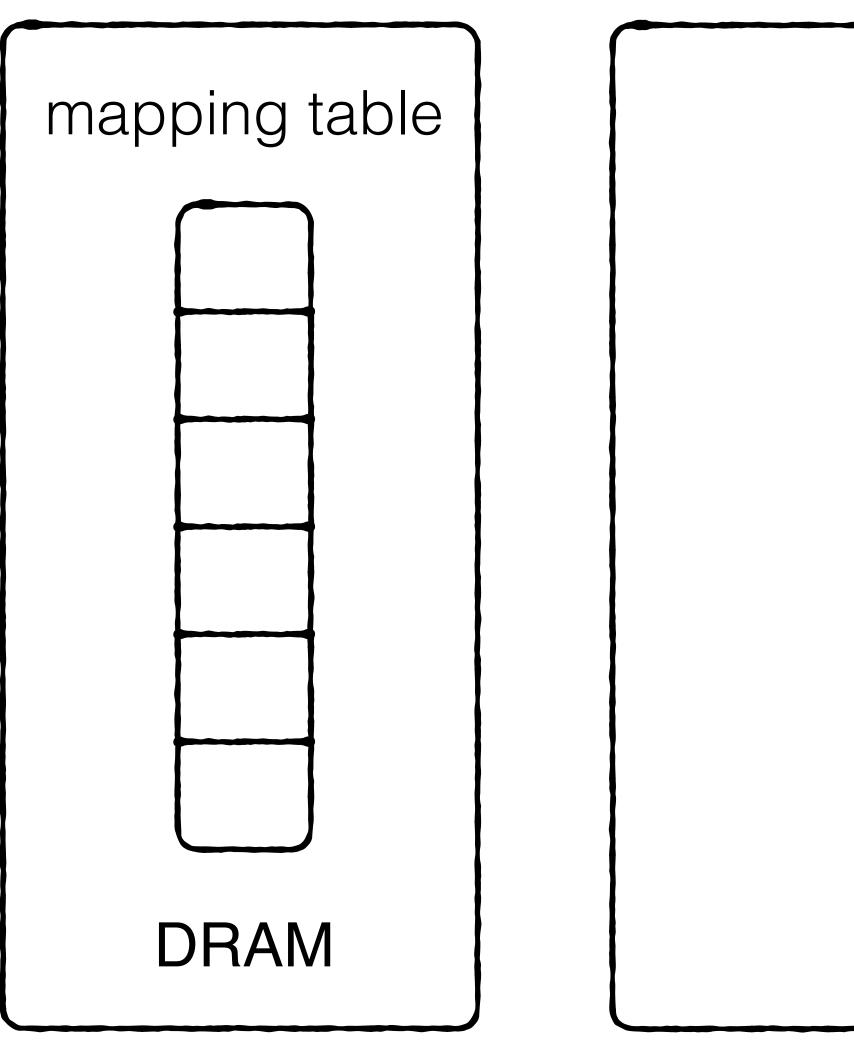
Doesn't sound like a lot but...



0	1	2	3	4	5
6	7	8	တ	10	11
12	13	14	15	16	17
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30

e.g., 1TB SSD with 4KB pages requires 1GB mapping table

Can store in host to save space (High-end SSDs)



Host

9 10 11

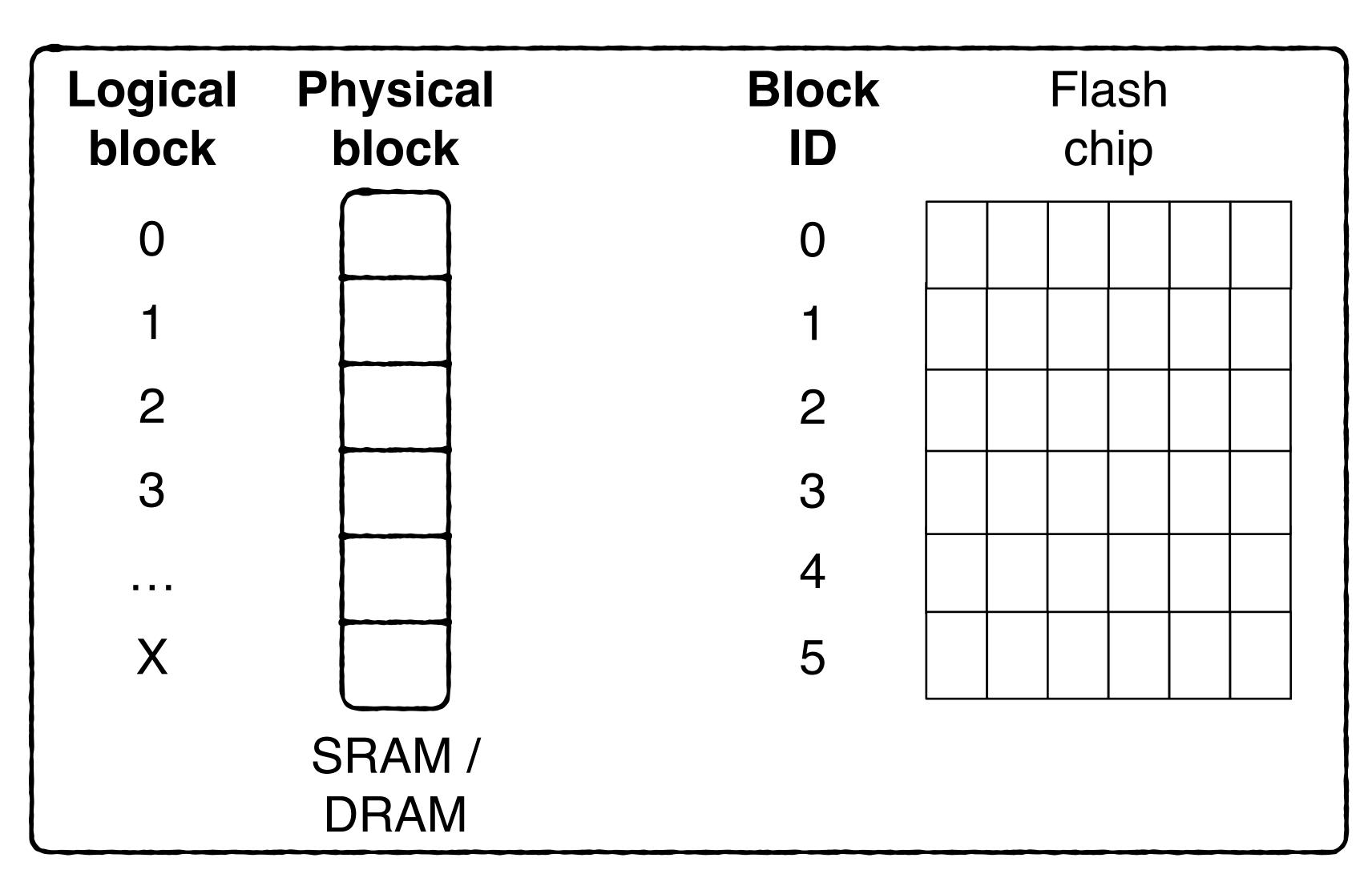
12 13 14 15 16 17

13 14 15 16 17 18

19 20 21 22 23 24

25 26 27 28 29 30

Alternative 1: use block mapping

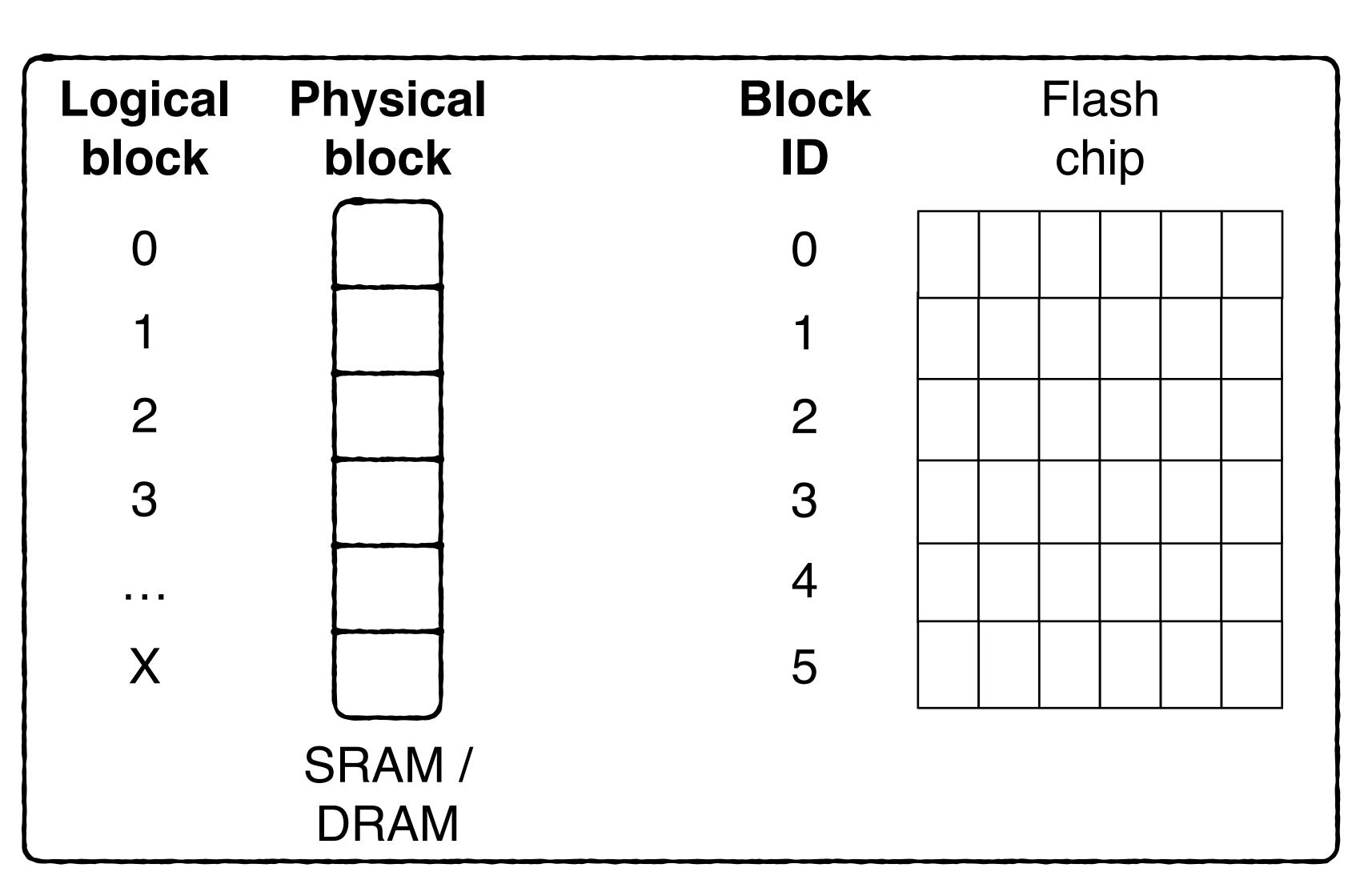


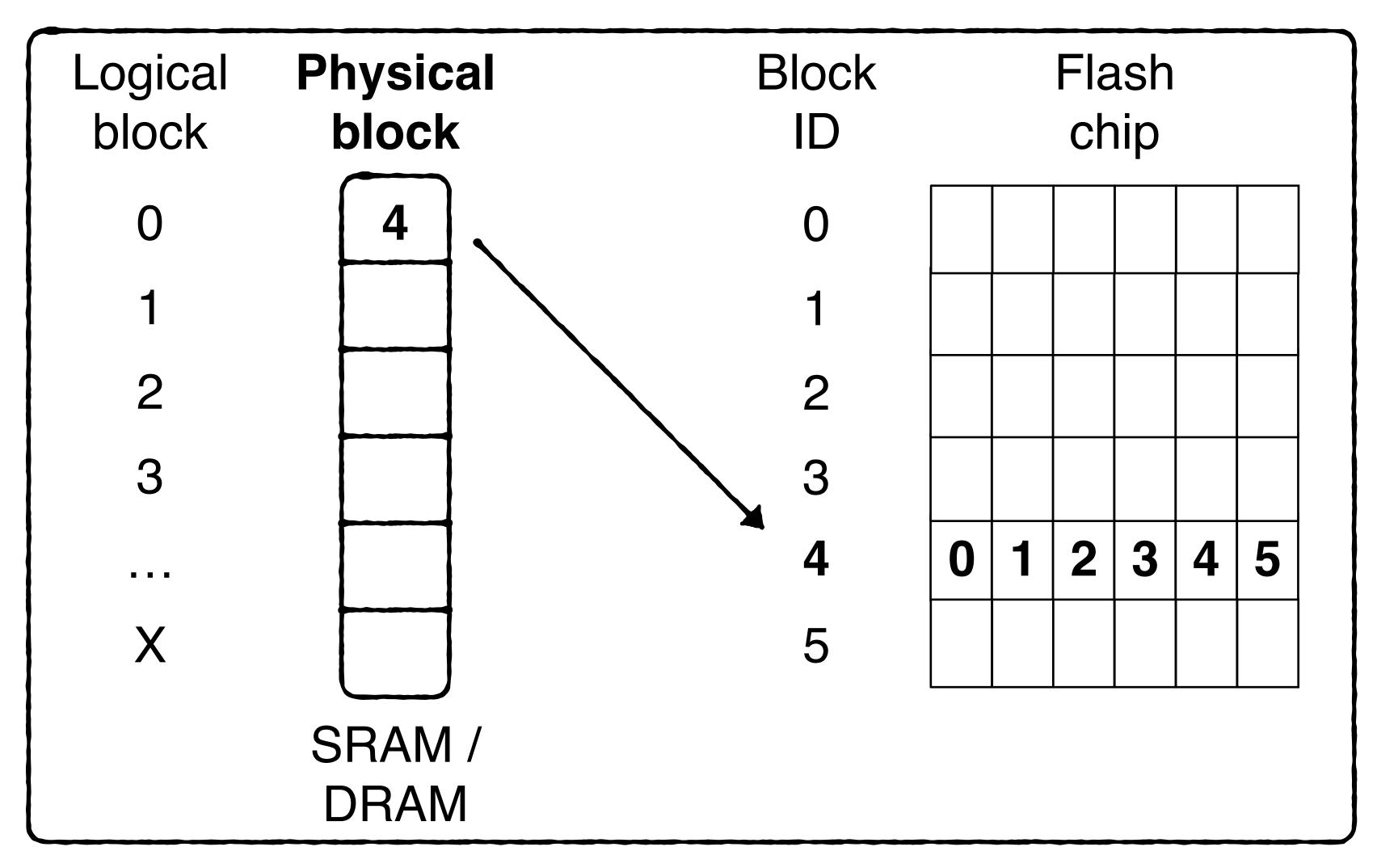
SSD

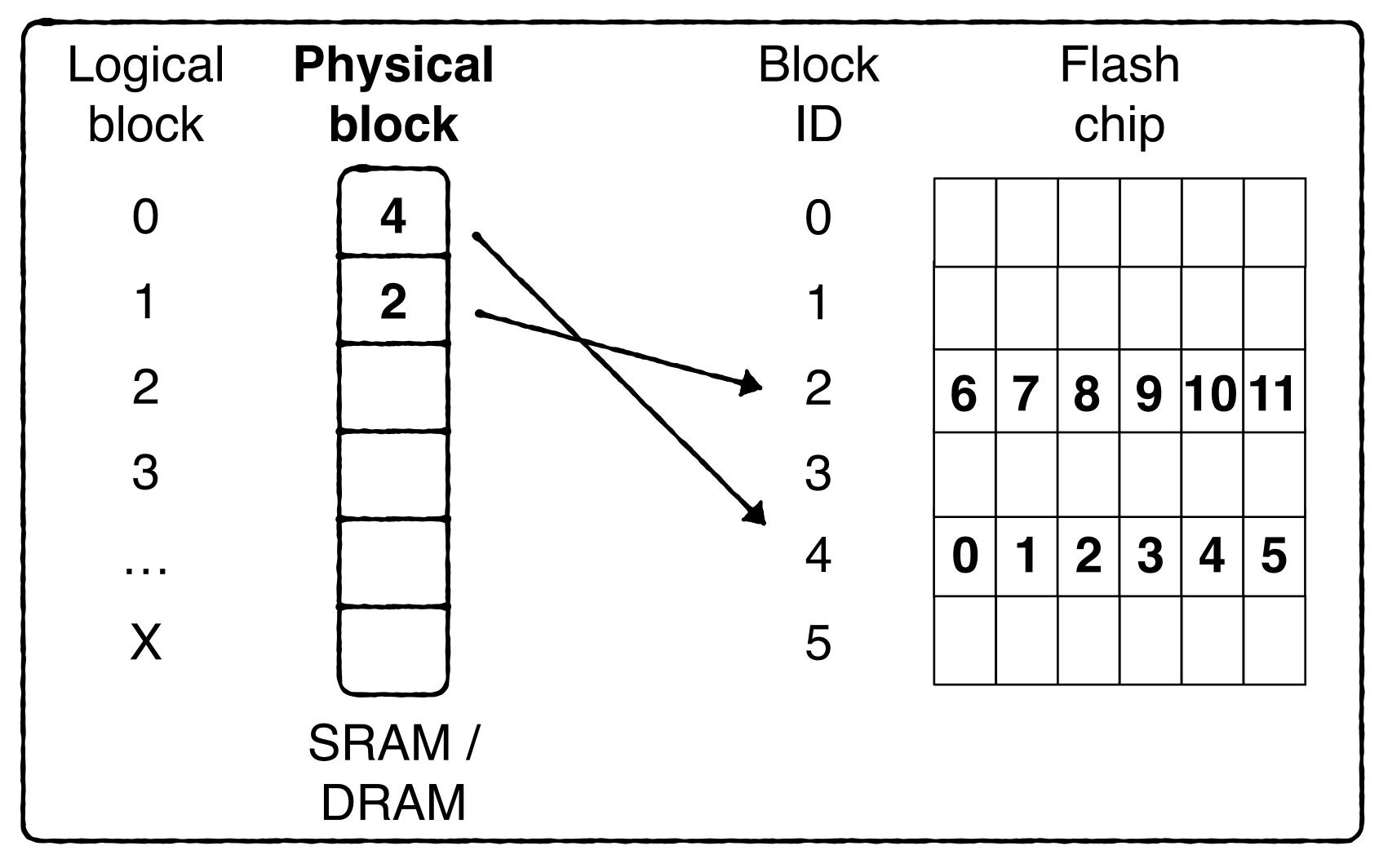
Alternative 1: use block mapping

FAST: A log buffer-based flash translation layer using fully-associative sector translation

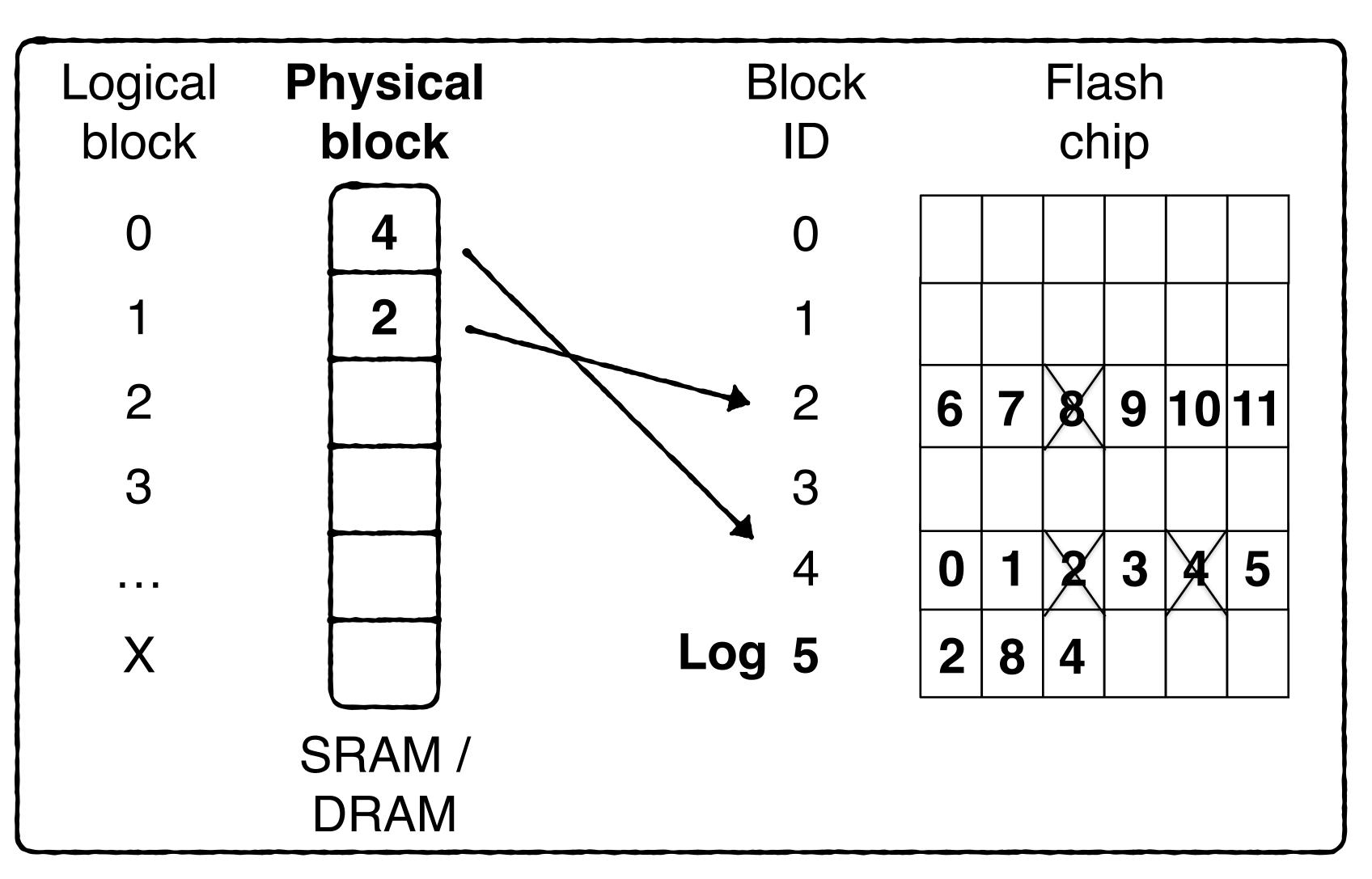
TECS 2007.



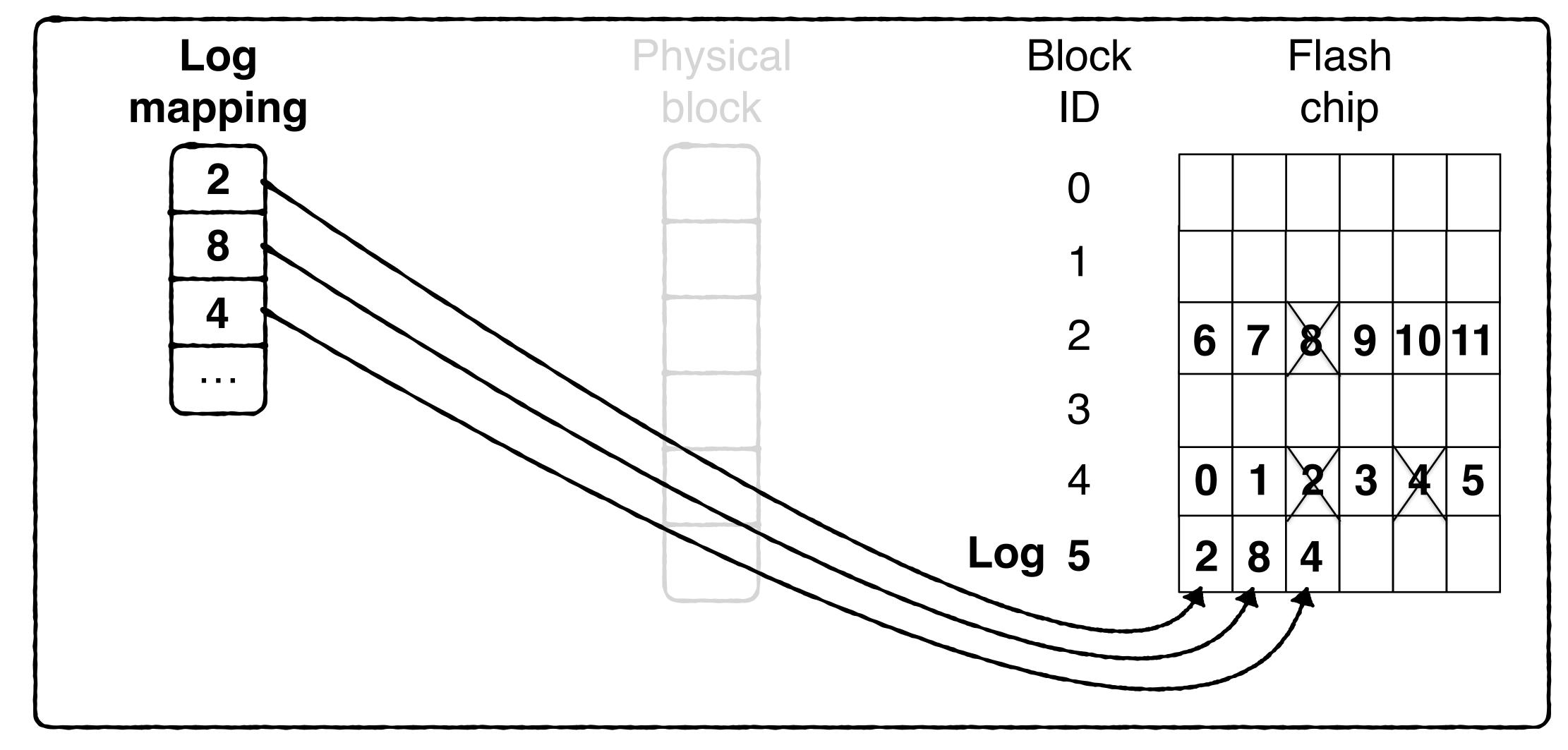


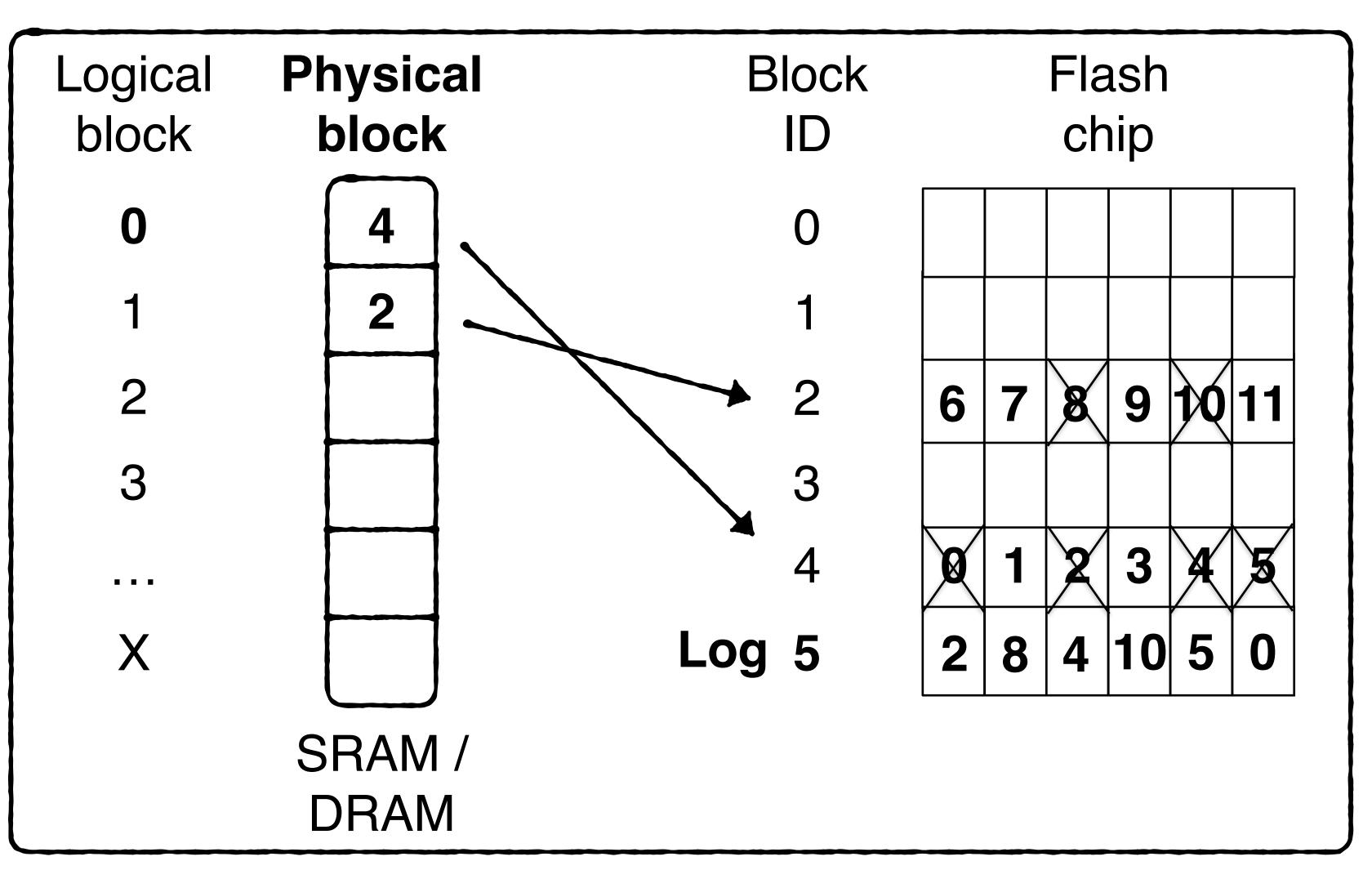


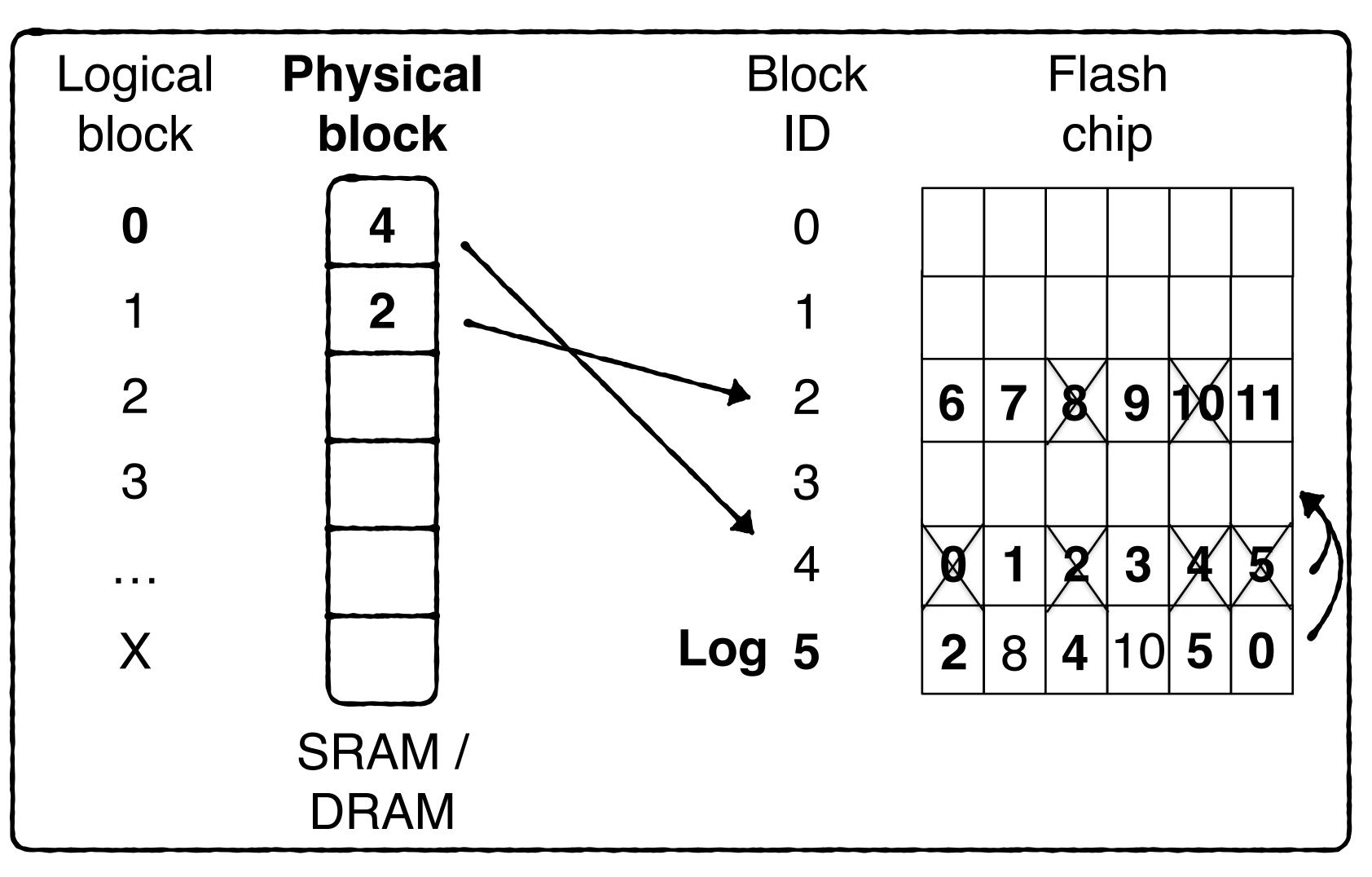
Store updates in log blocks

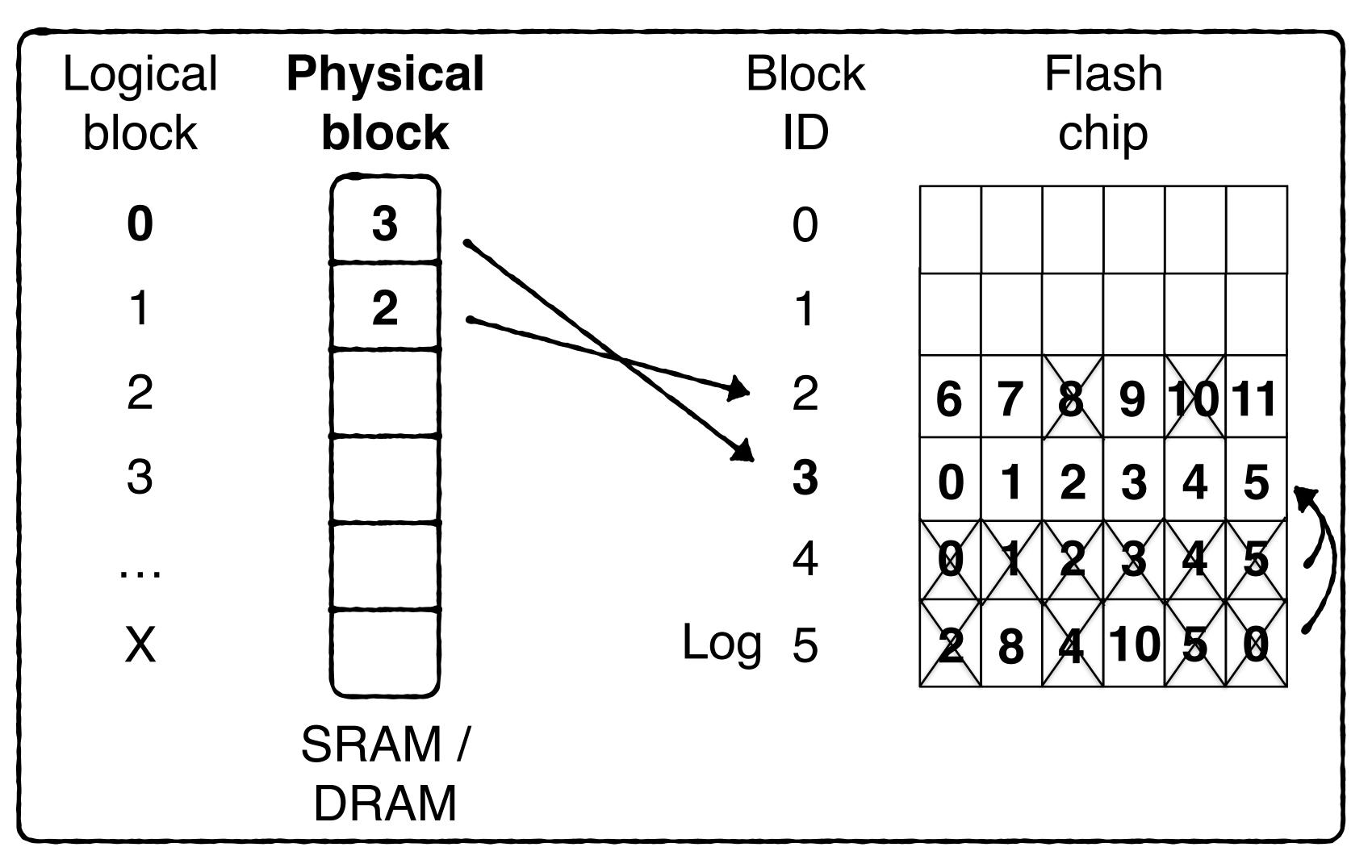


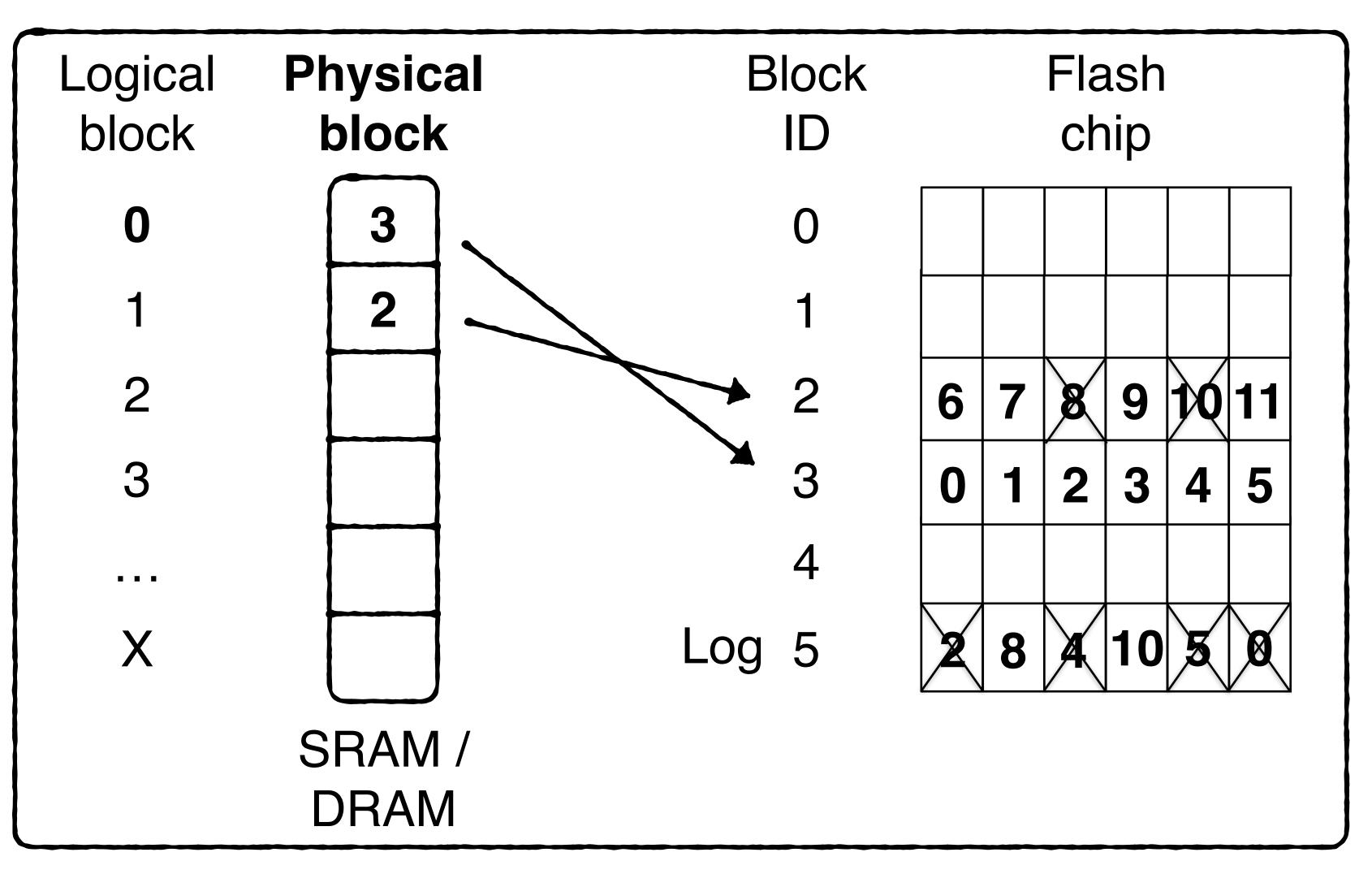
log blocks use page mapping. There are few of them

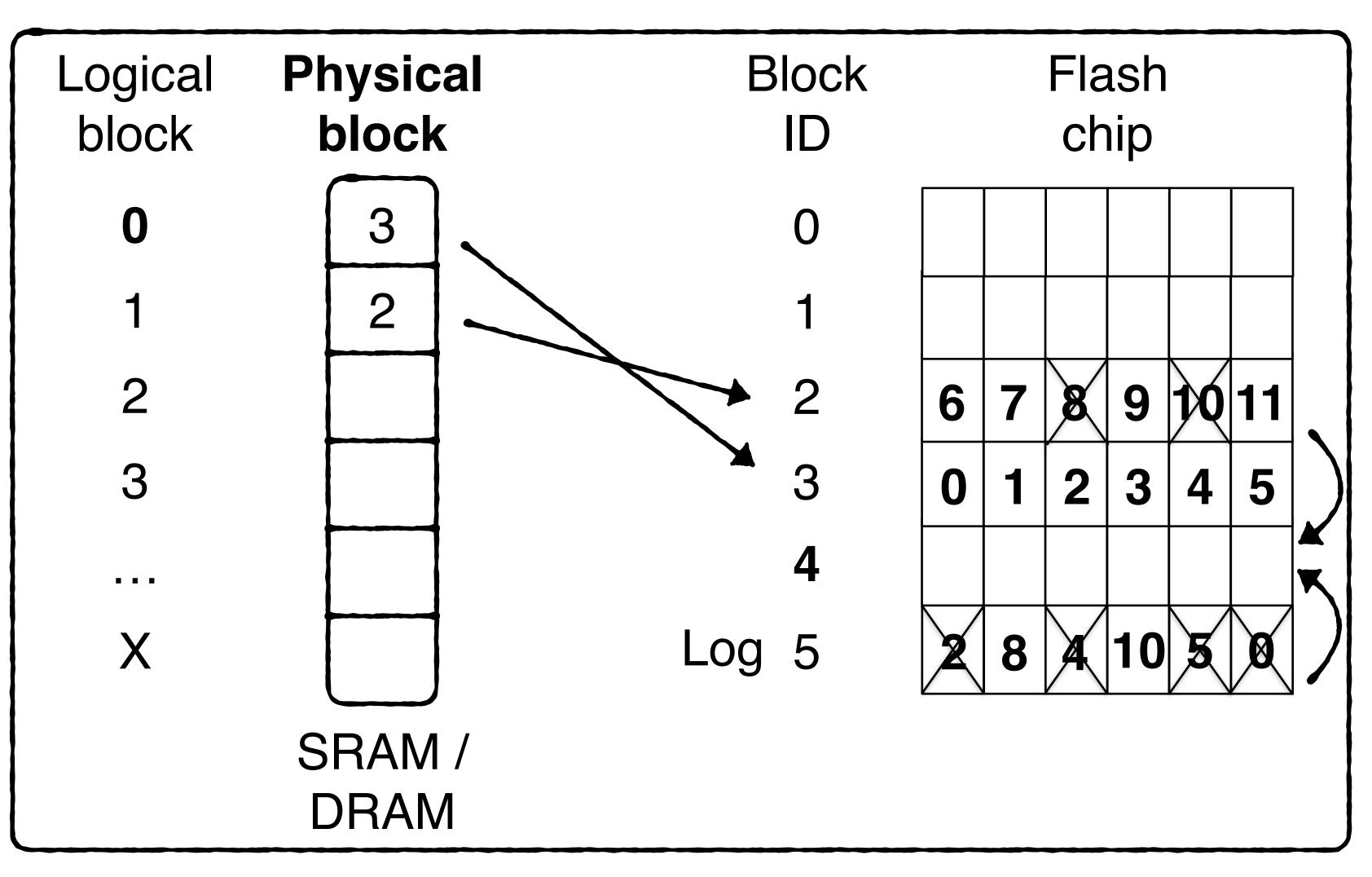


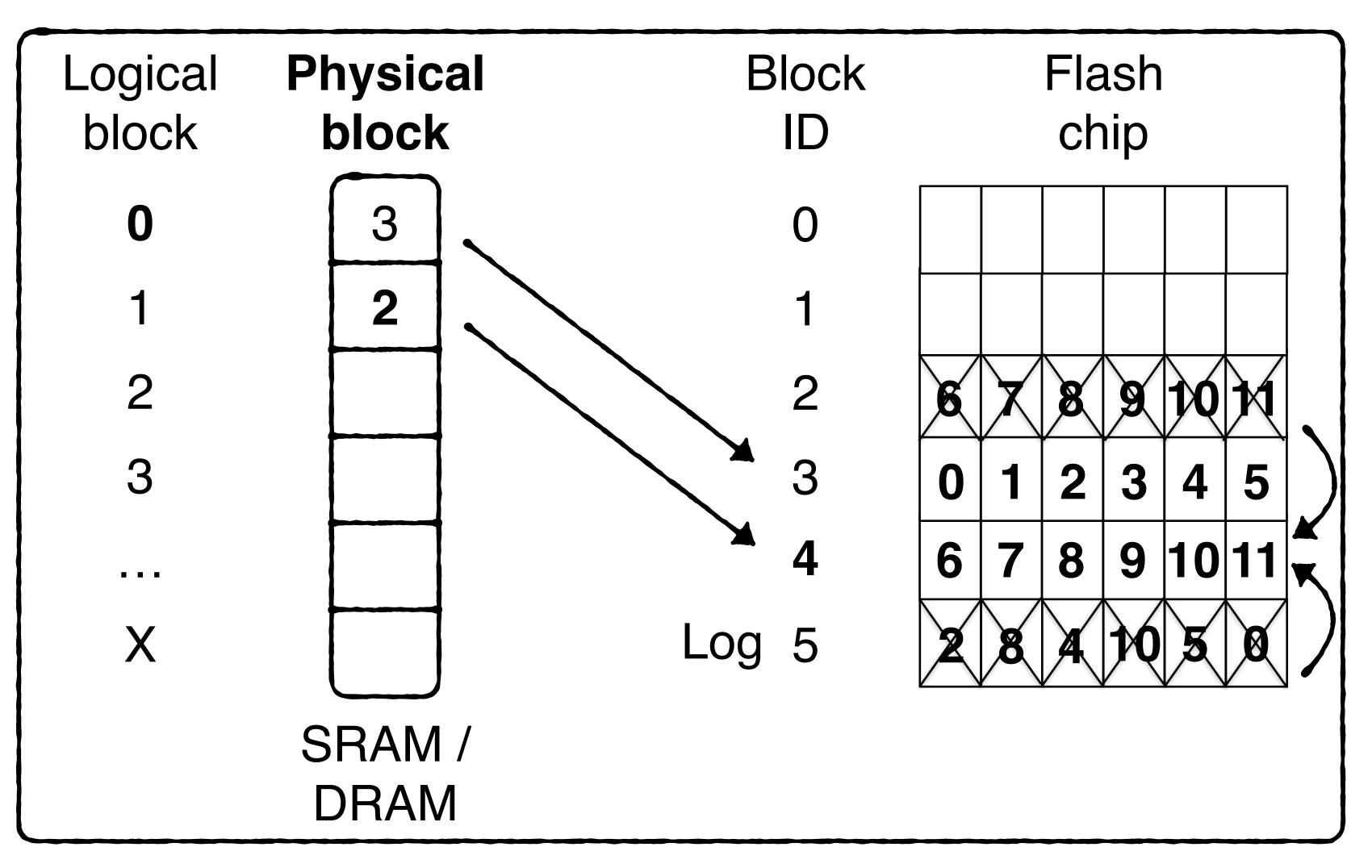


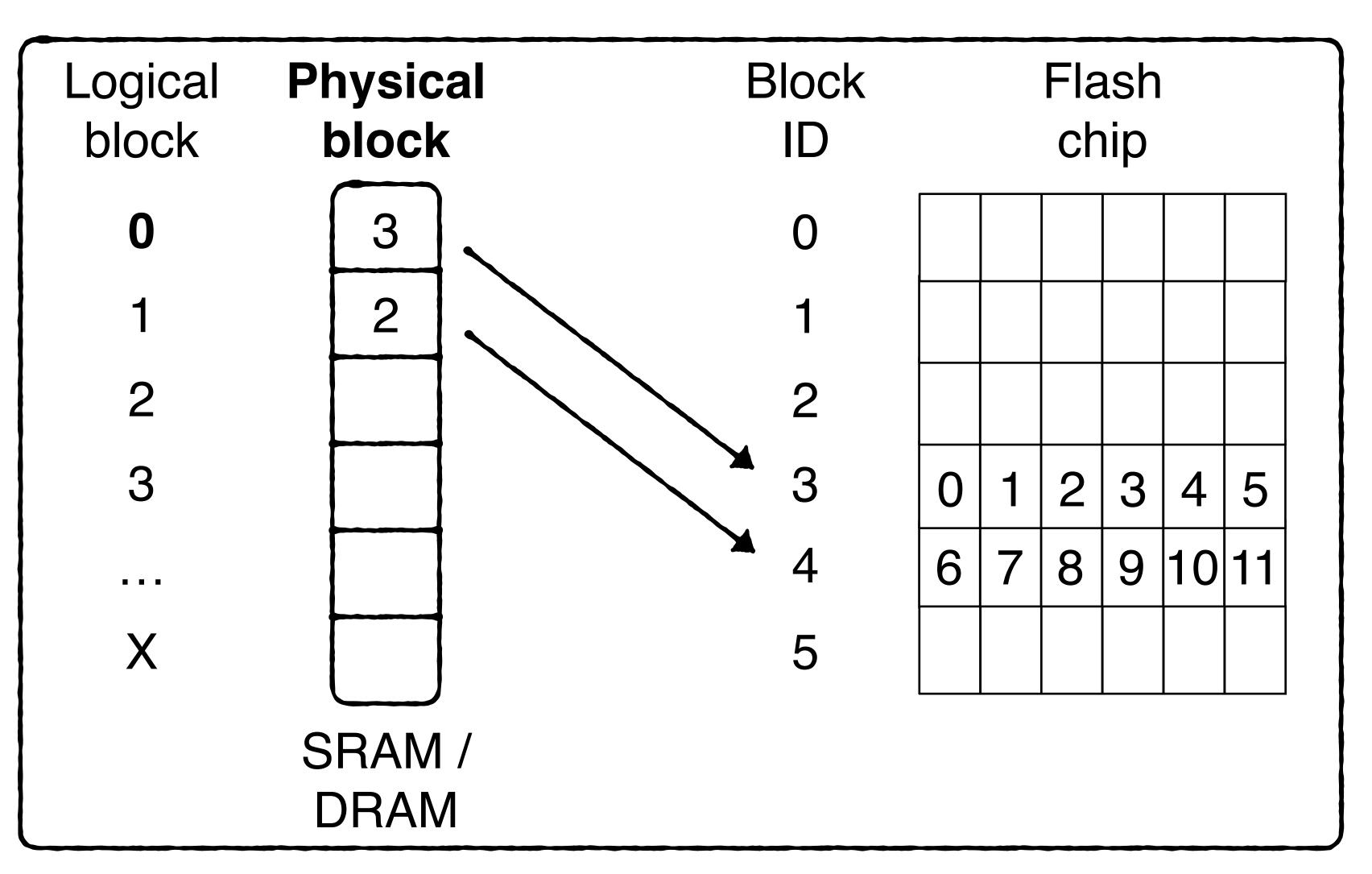




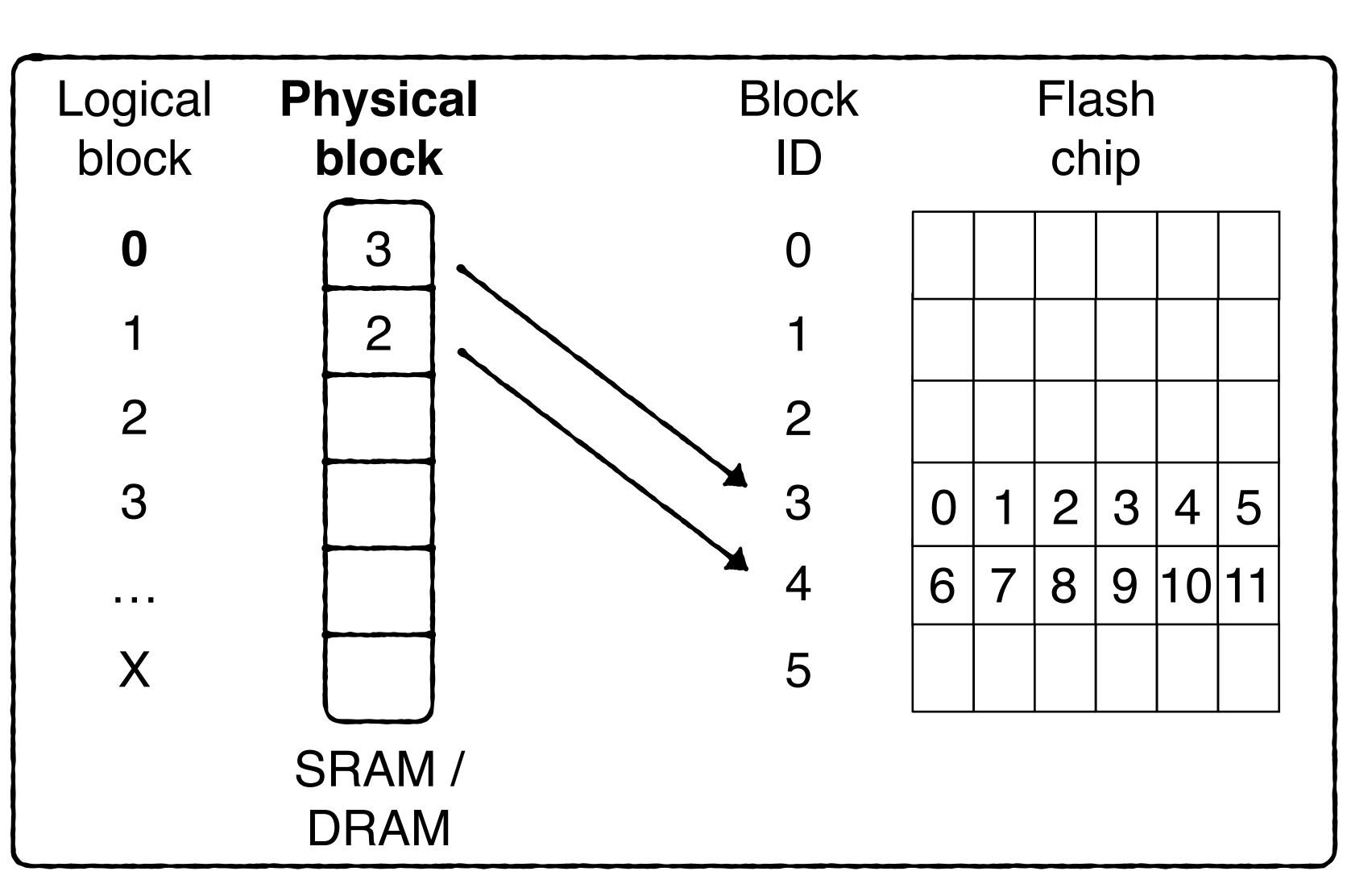






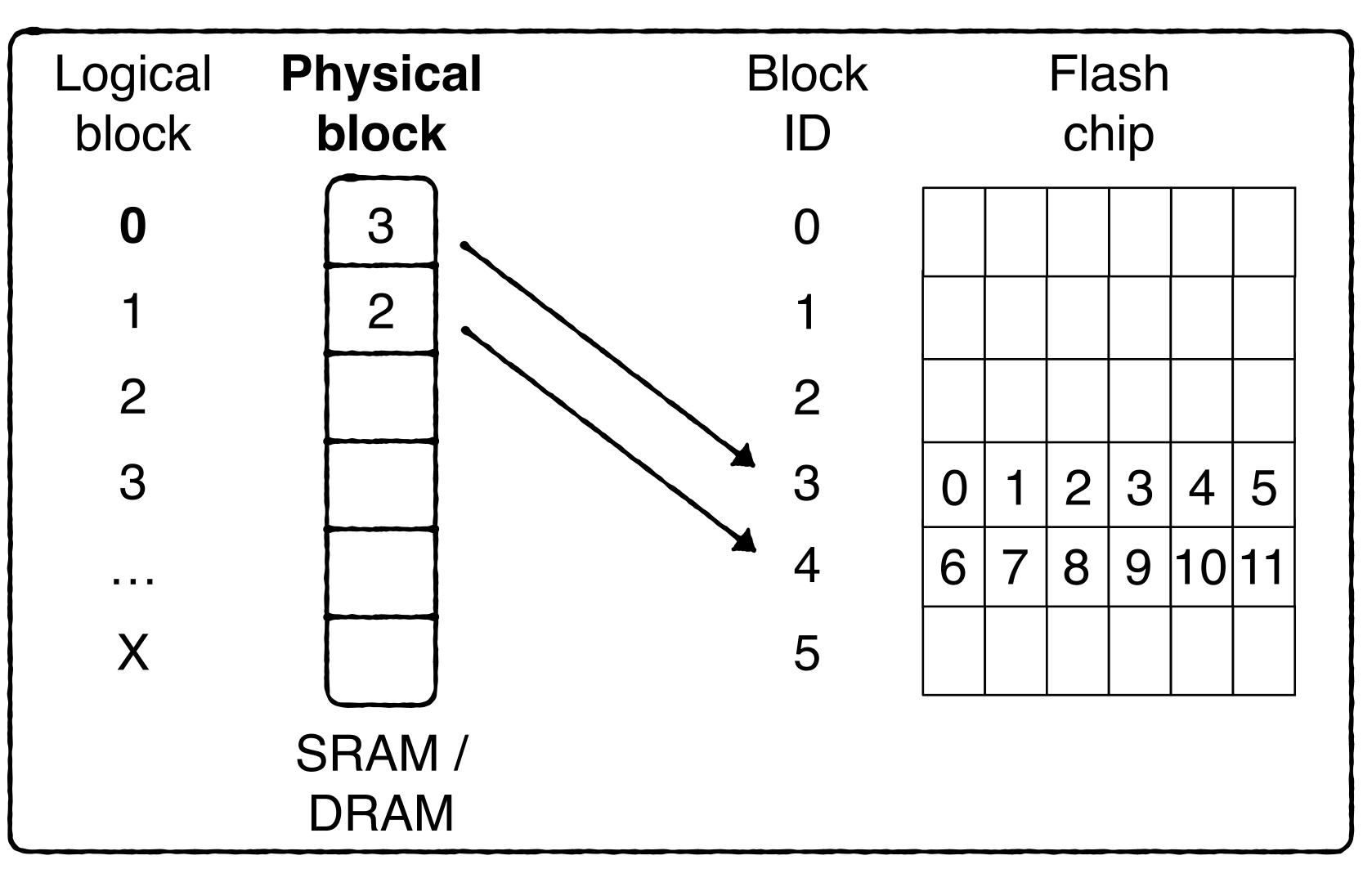


Coupling between mapping table and GC.



Coupling between mapping table and GC.

One log block can have updates from different data blocks, leading to long GC process that rewrites many blocks



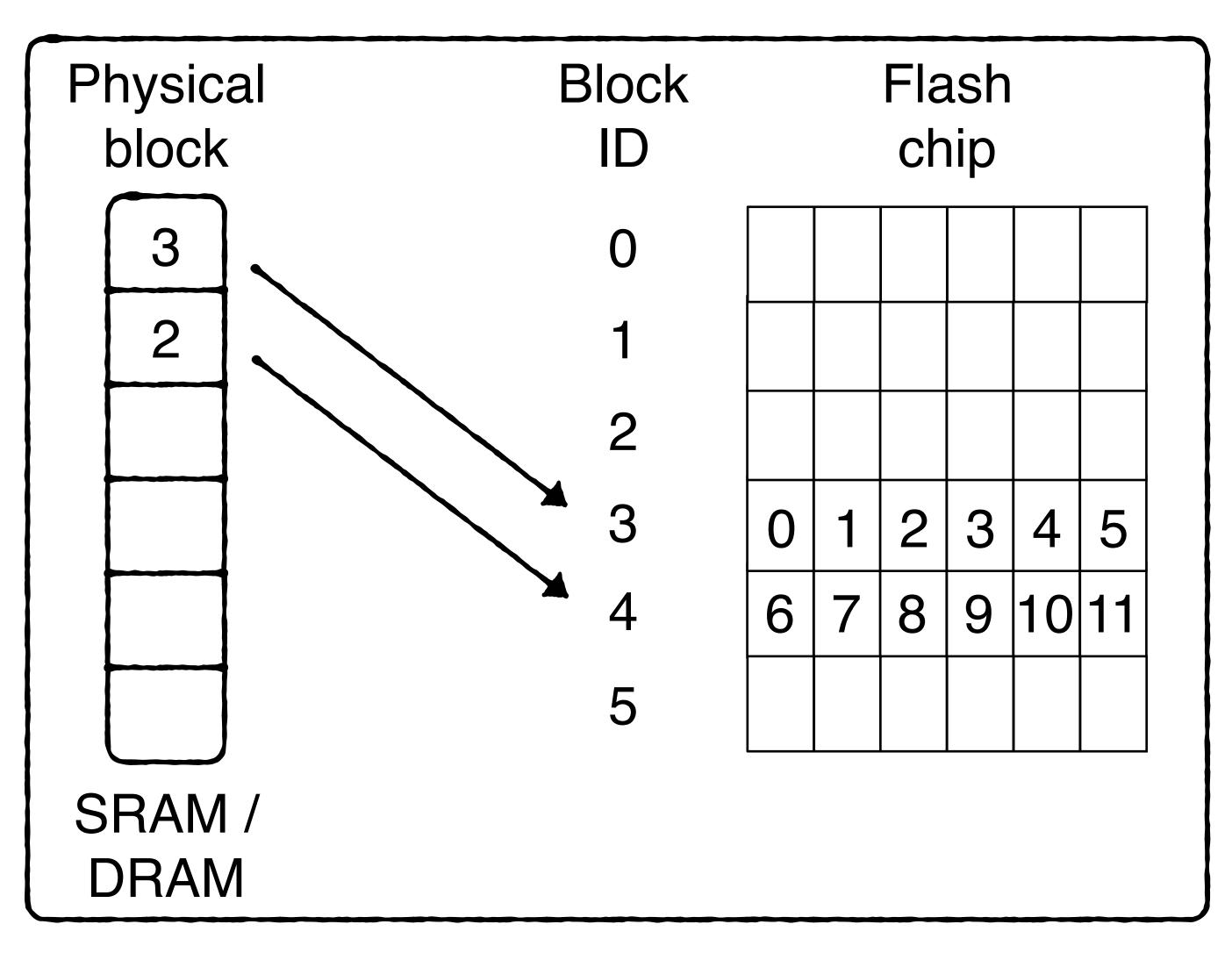
Coupling between mapping table and GC.

Sequential writes



Random writes





Better ideas?:)

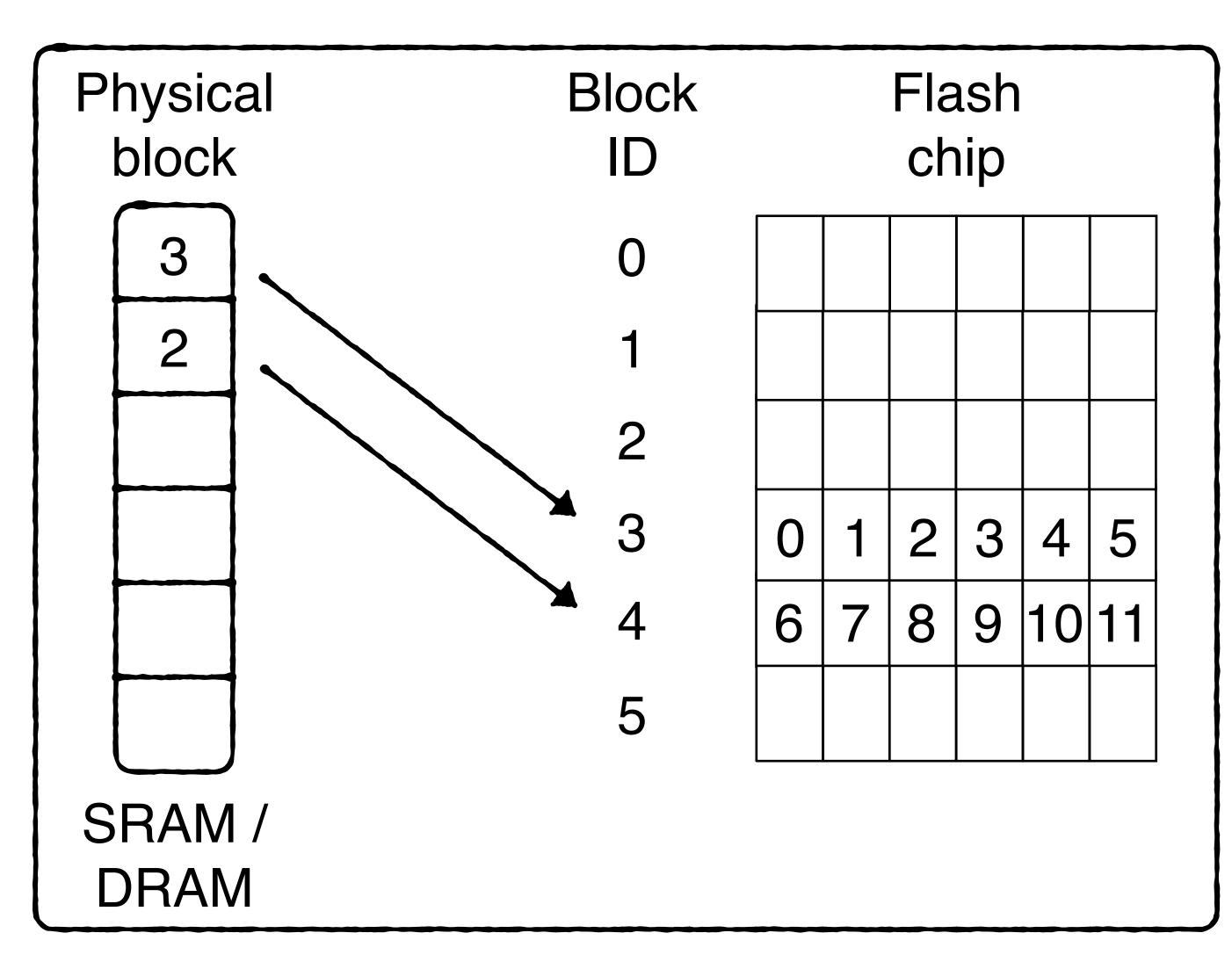
Downsides?

Coupling between mapping table and GC.

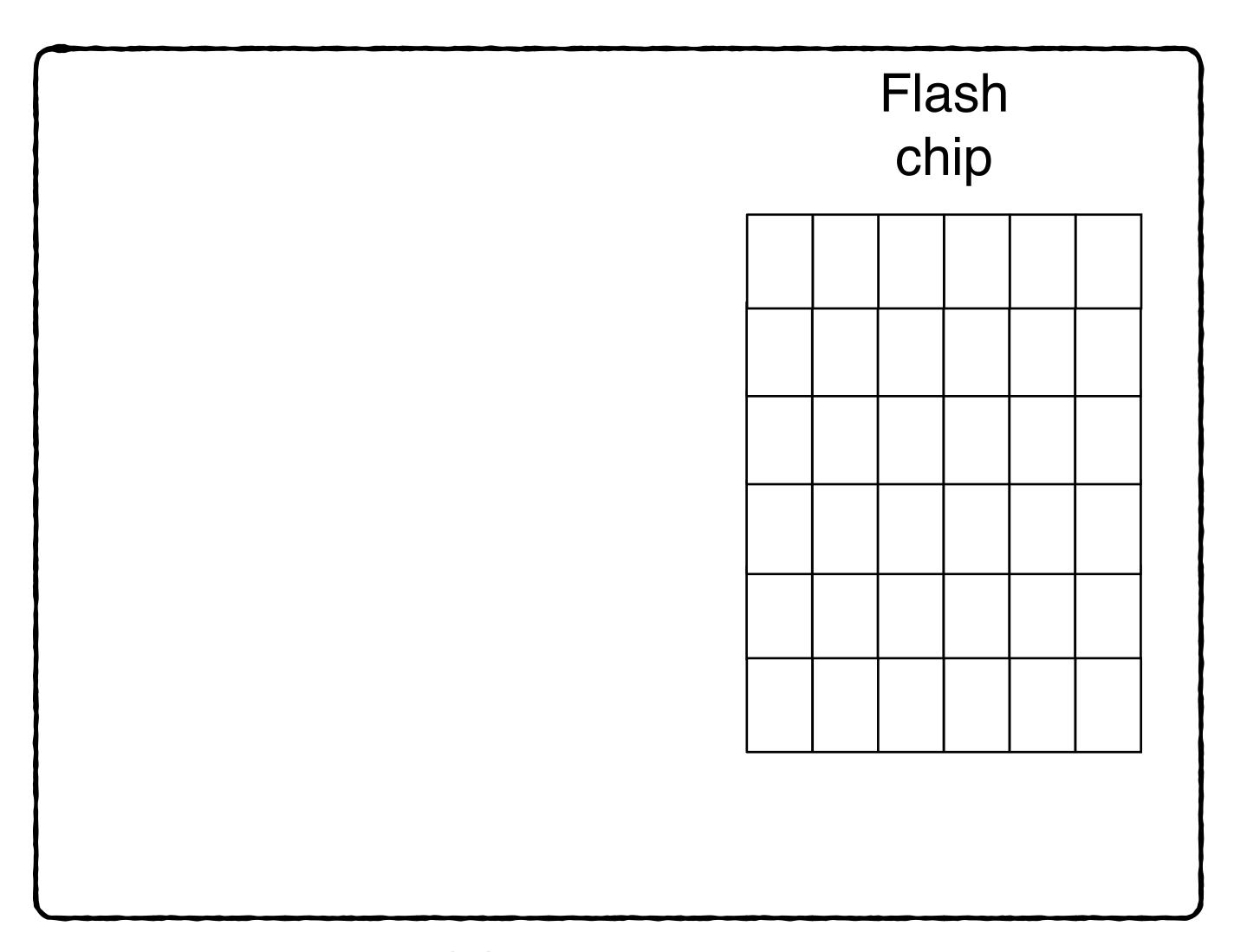
Sequential writes

Random writes





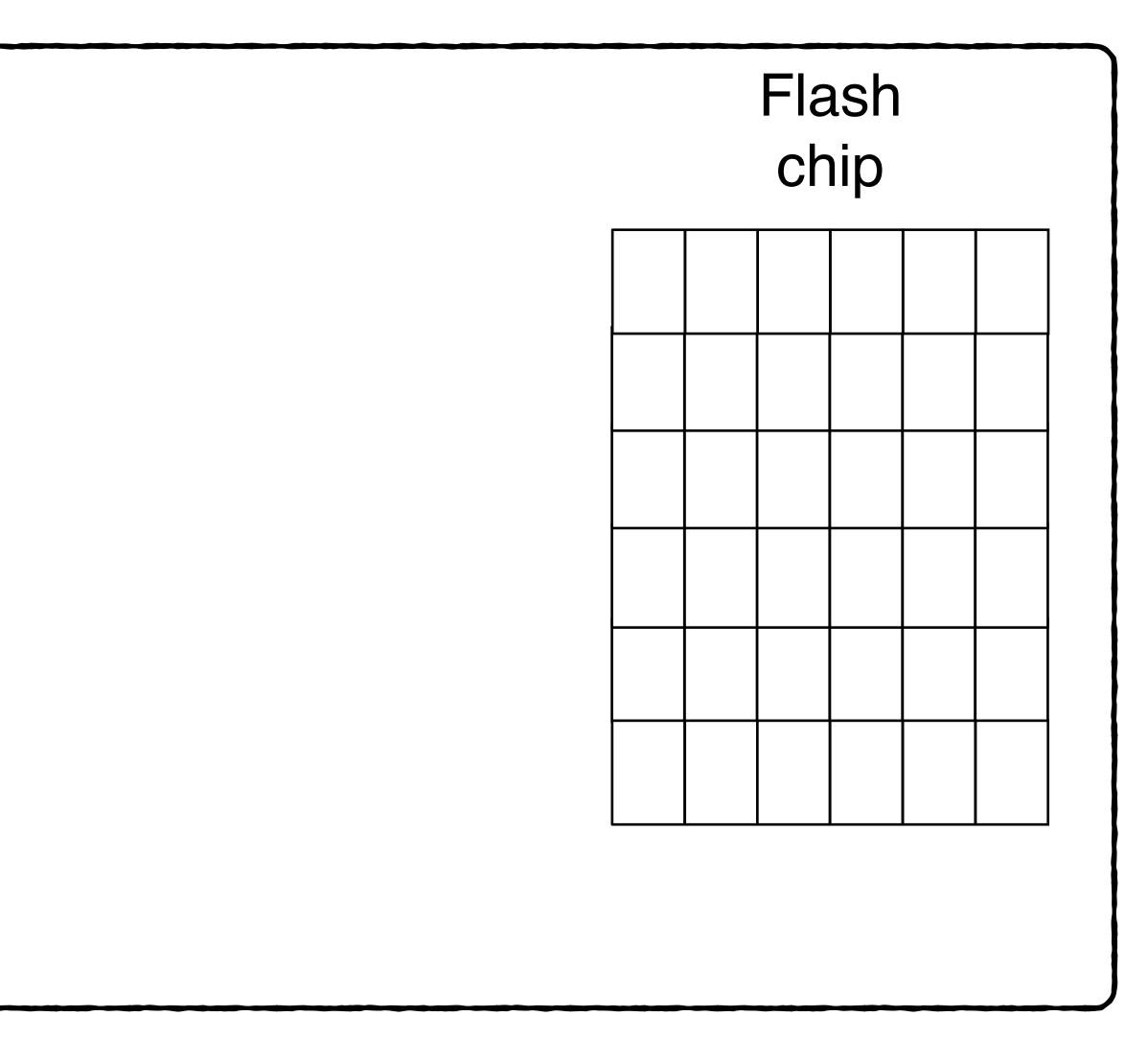
Alternative 2: Store Page Mapping in Flash & Cache Parts

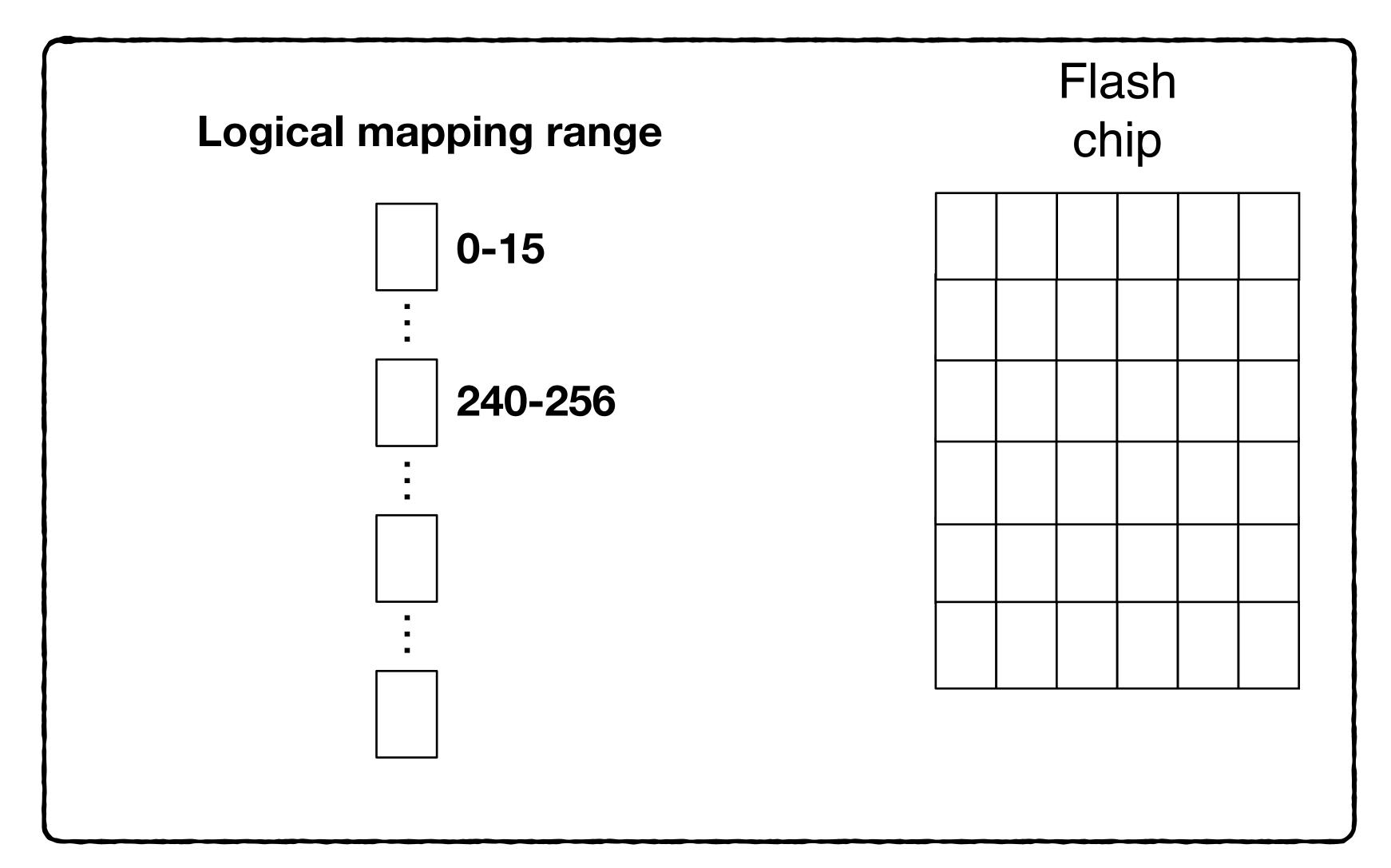


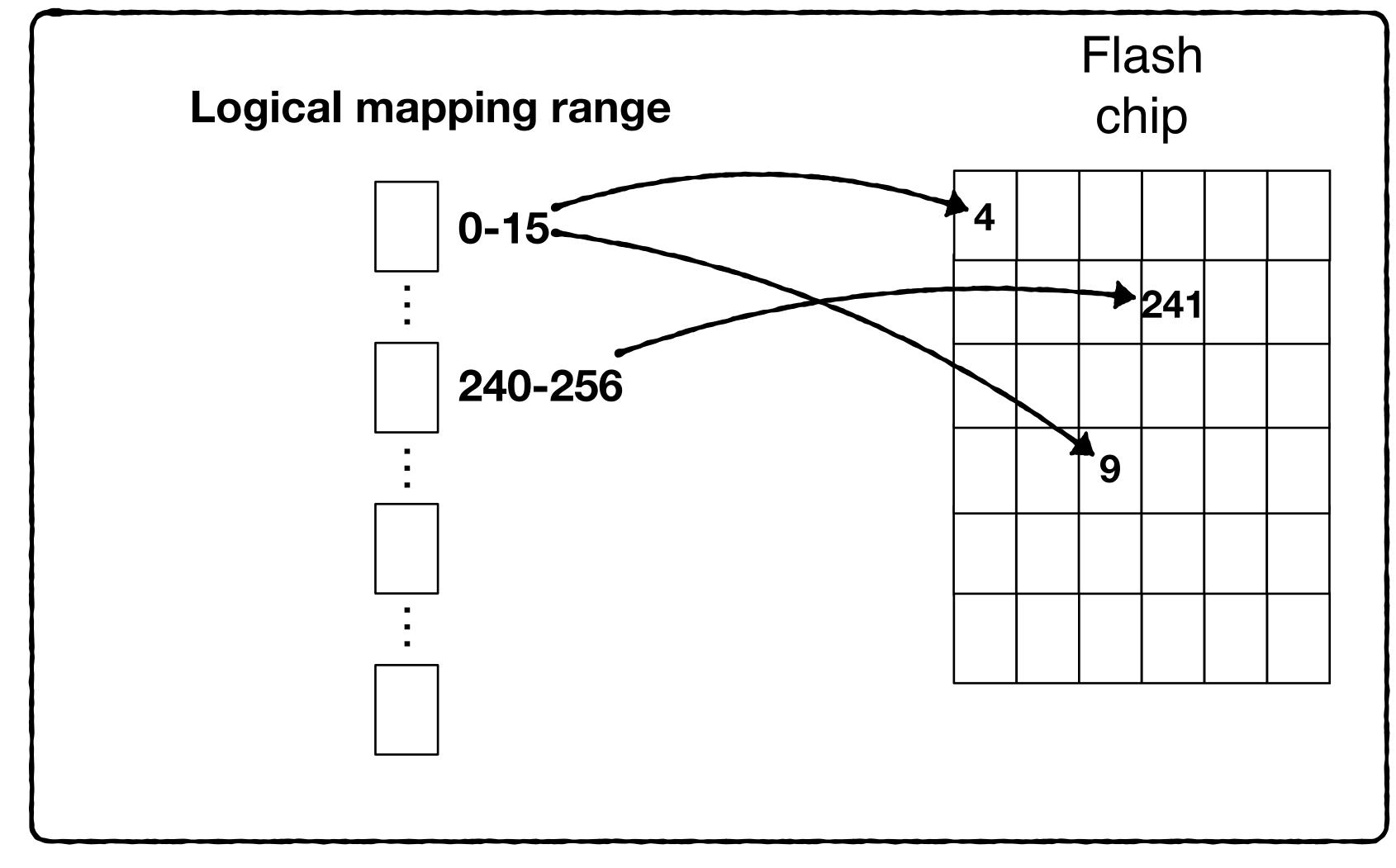
Alternative 2: Store Page Mapping in Flash & Cache Parts

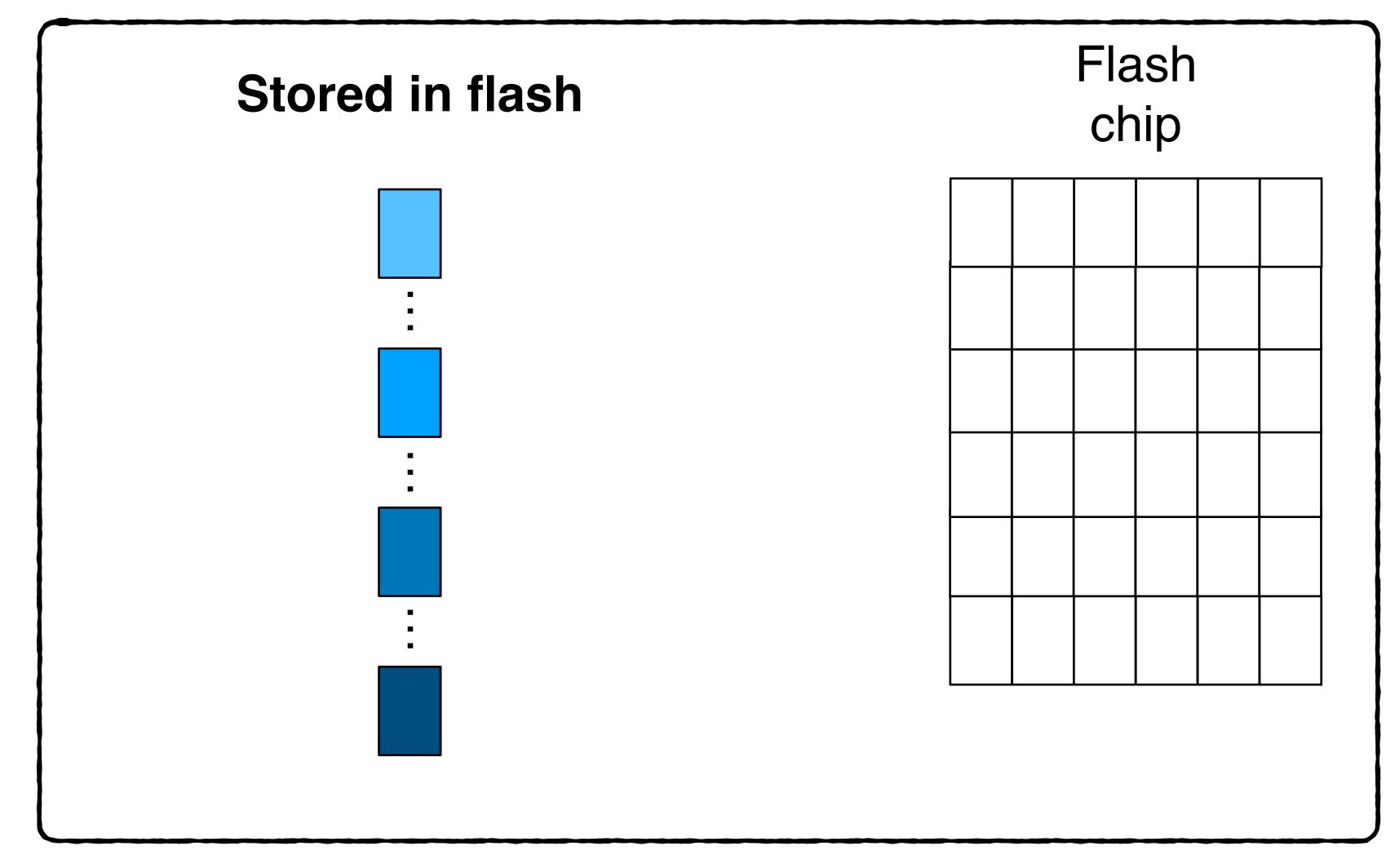
DFTL: a flash translation layer employing demand-based selective caching of page-level address mappings

ACM SIGPLAN Notices 2009

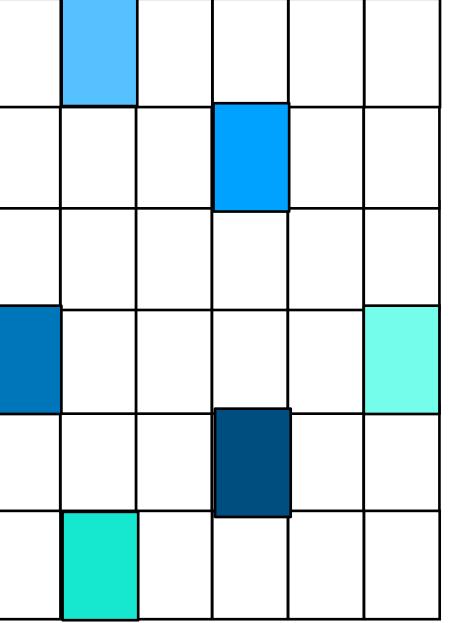




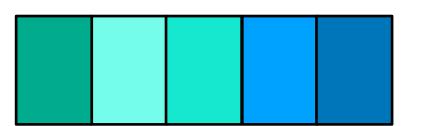




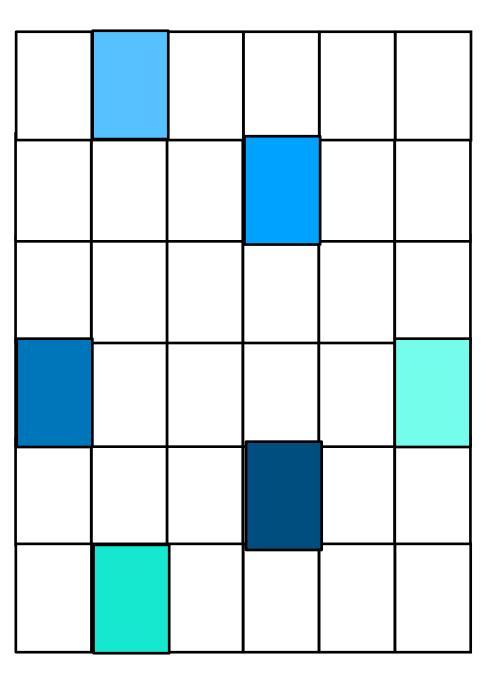
Flash chip



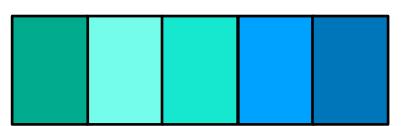
Buffer pool frequently accessed mapping pages (LRU/Clock)



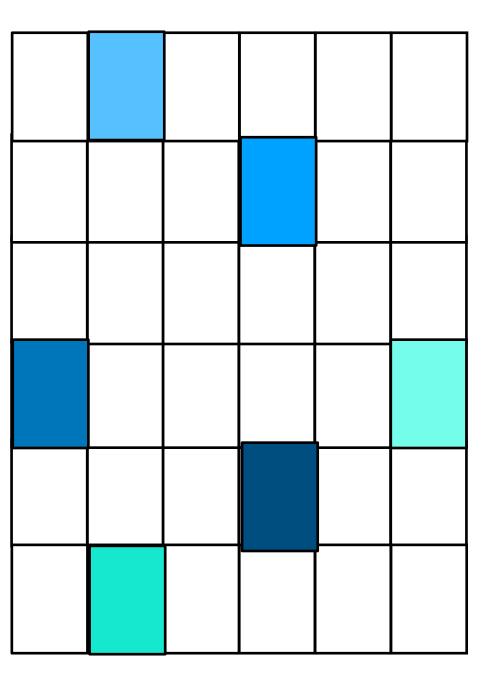
Flash chip



Buffer pool (LRU/Clock)

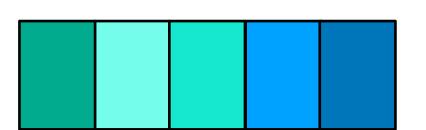


Flash chip

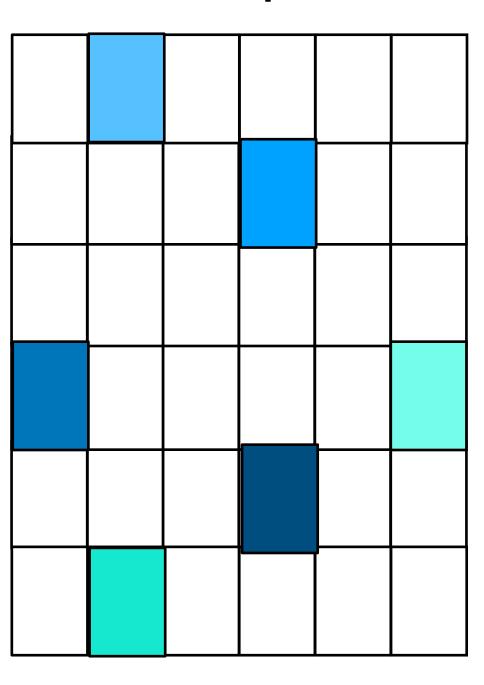


(1) read costs 2 I/Os if mapping page isn't cached

Buffer pool (LRU/Clock)



Flash chip

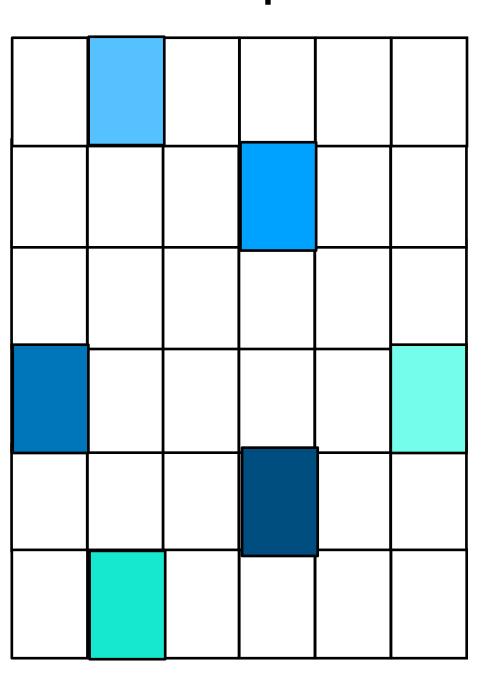


- (1) read costs 2 I/Os if mapping page isn't cached
- (2) writes may cost more due to buffer pool evictions

Buffer pool (LRU/Clock)



Flash chip



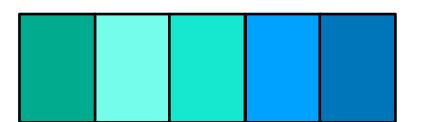
- (1) read costs 2 I/Os if mapping page isn't cached
- (2) writes may cost more due to buffer pool evictions

Sequential writes

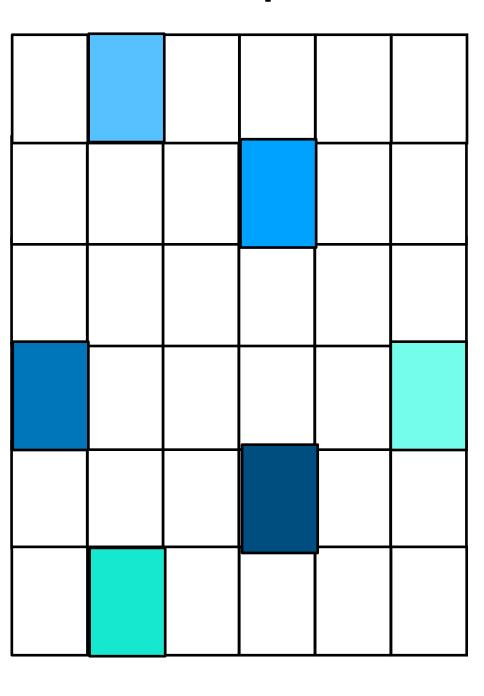


Random I/Os





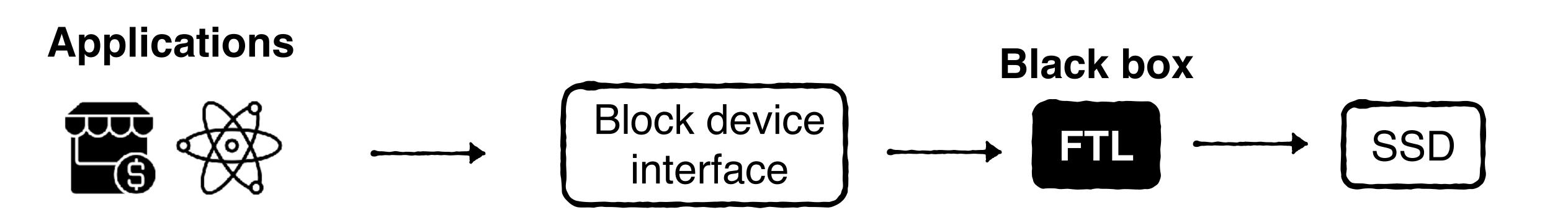
Flash chip



FTL implementation is opaque

Applications Black box Block device interface SSD

FTL implementation is opaque



Hard to reason about performance outcomes given random SSD

FTL Downsides

Memory overhead for mapping table

Storage space for over-provisioning

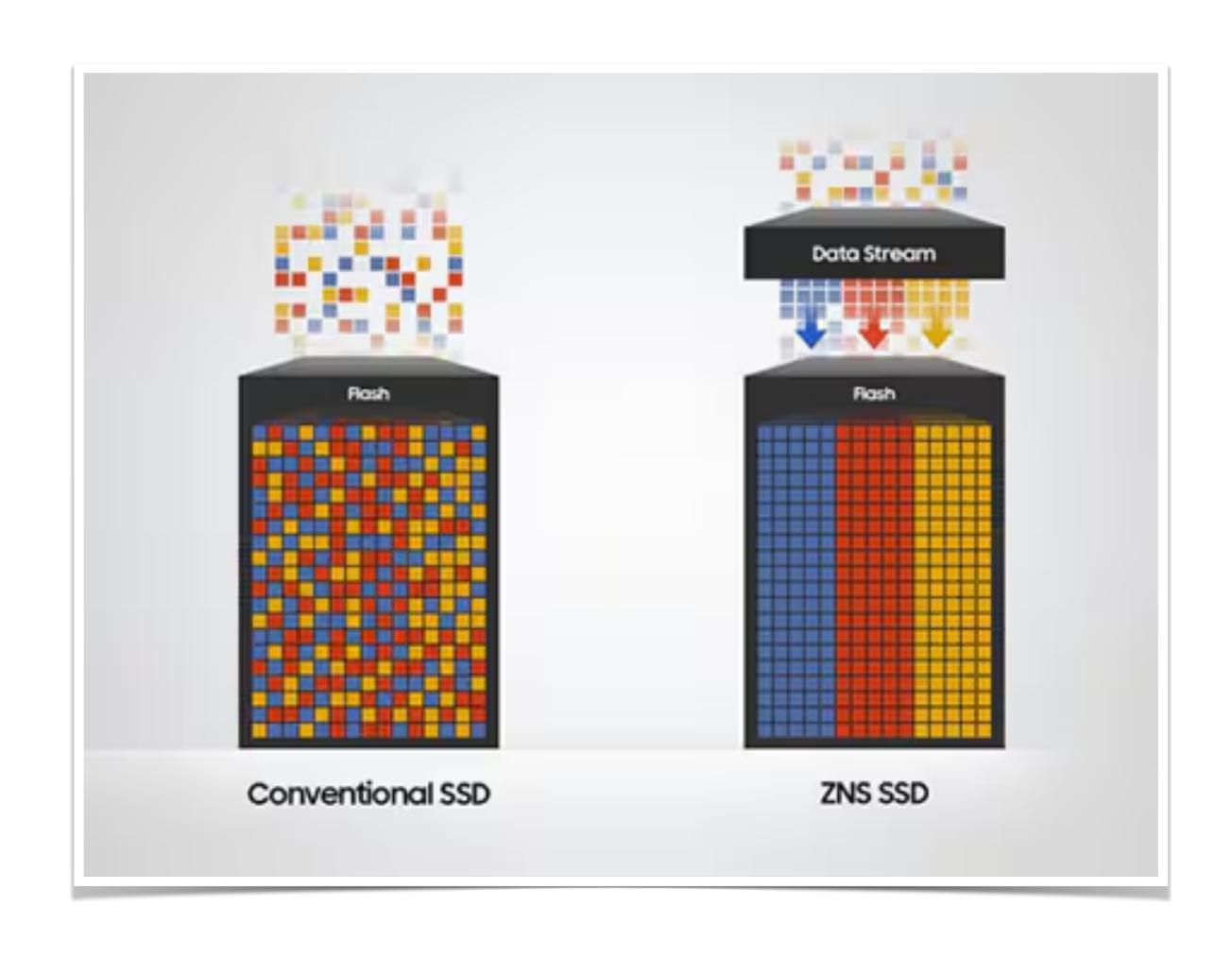
I/O & lifetime for garbage-collection

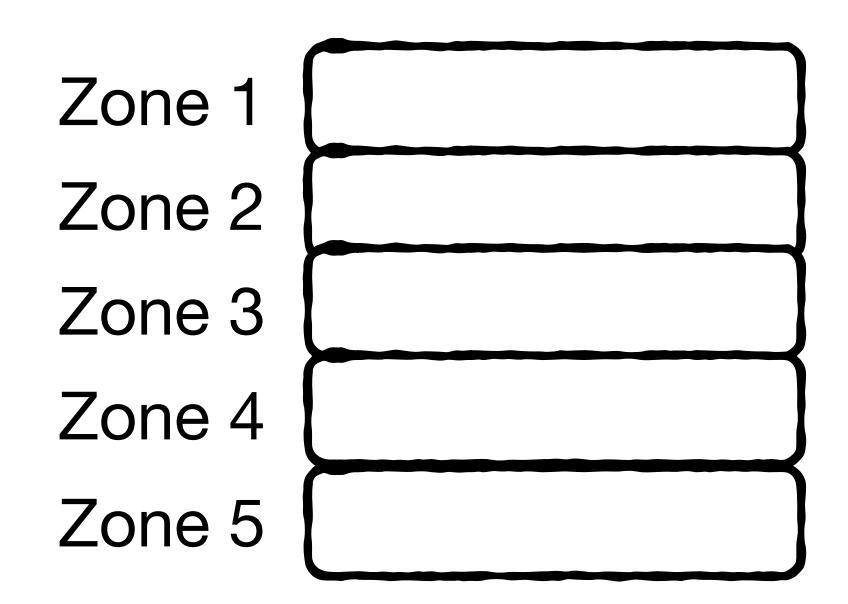






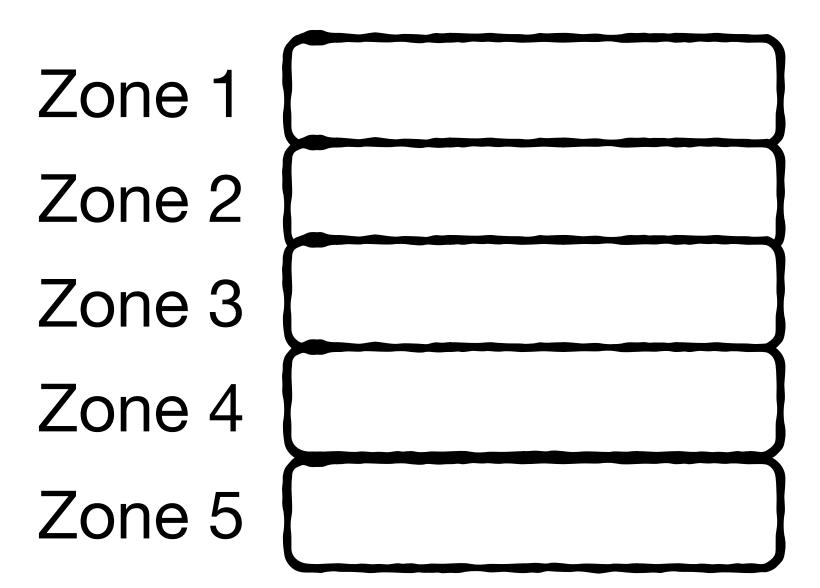






Applications write a zone sequentially

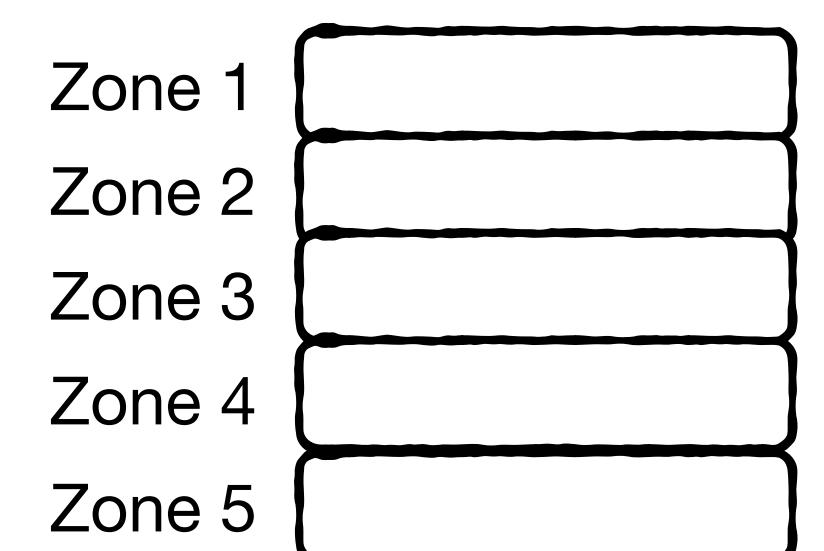
& later deletes it



Applications write a zone sequentially

& later deletes it

SSD performs no garbage-collection, mapping is small, and no over-provisioning



Applications write a zone sequentially

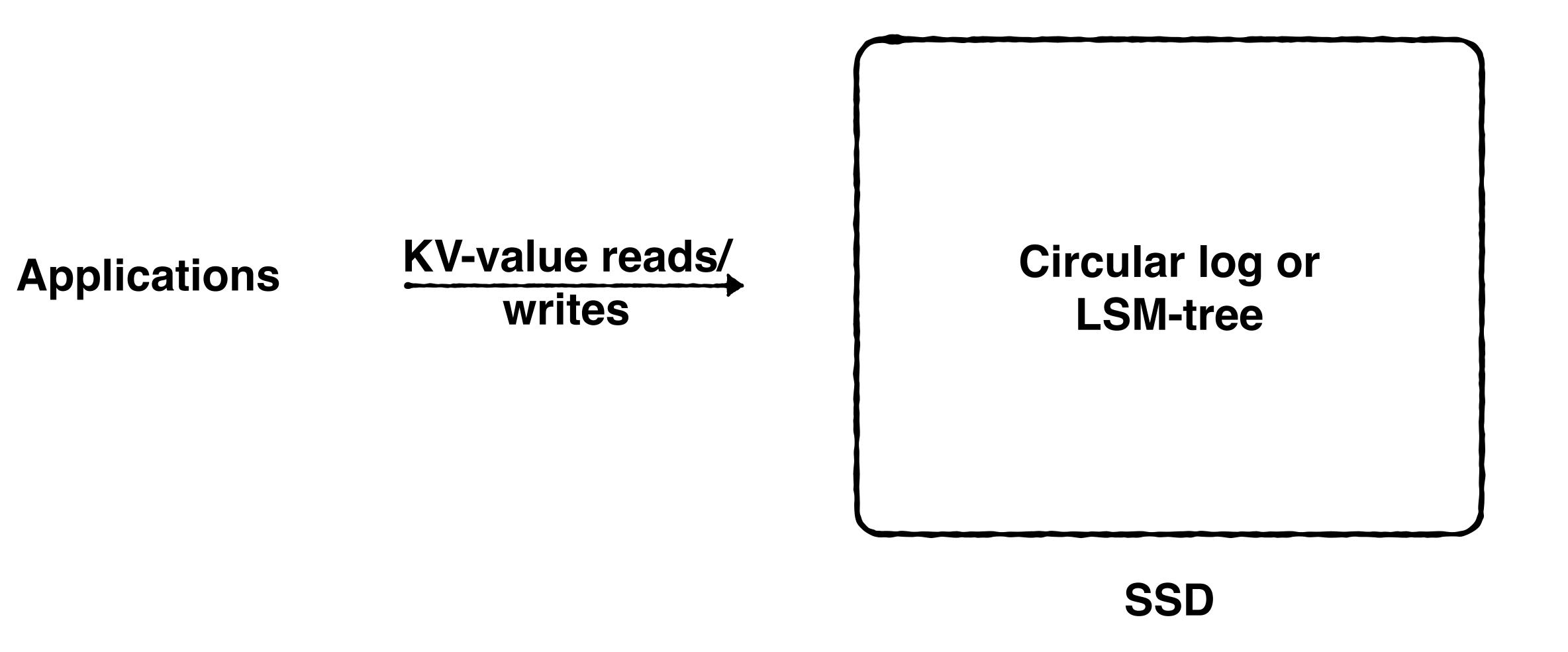
& later deletes it

SSD performs no garbage-collection, mapping is small, and no over-provisioning

Zone 1
Zone 2
Zone 3
Zone 4
Zone 5

Restrictive & application has to essentially implement FTL

Key-Value SSDs



Key-Value SSDs

No data movement or CPU for compaction or circular log garbage-collection

Reduces SSD over-provisioning

Circular log or LSM-tree

SSD

Spooky

Granulating LSM-Tree Compactions Correctly



Niv Dayan*
University of Toronto

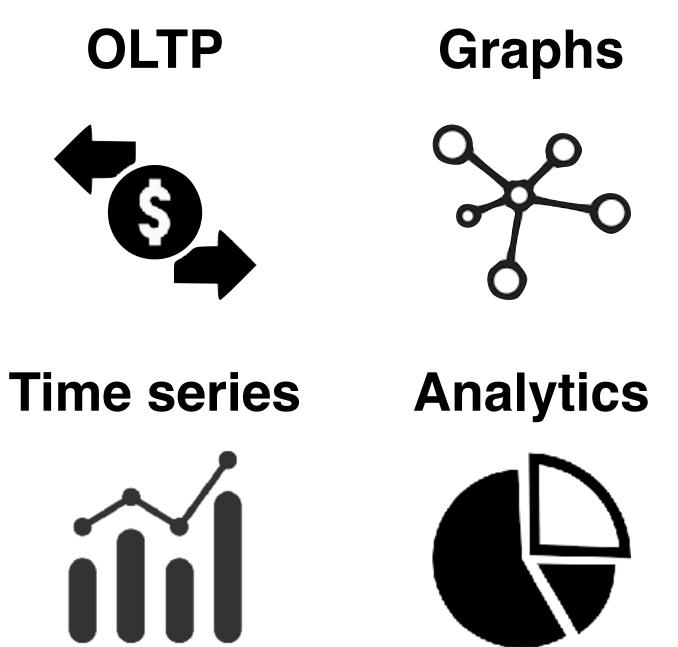
Tamar Weiss, Shmuel Dashevsky, Michael Pan, Edward Bortnikov, Moshe Twitto Pliops

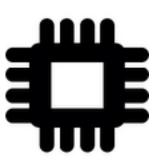
Spooky Getting LSM-trees right on SSDs





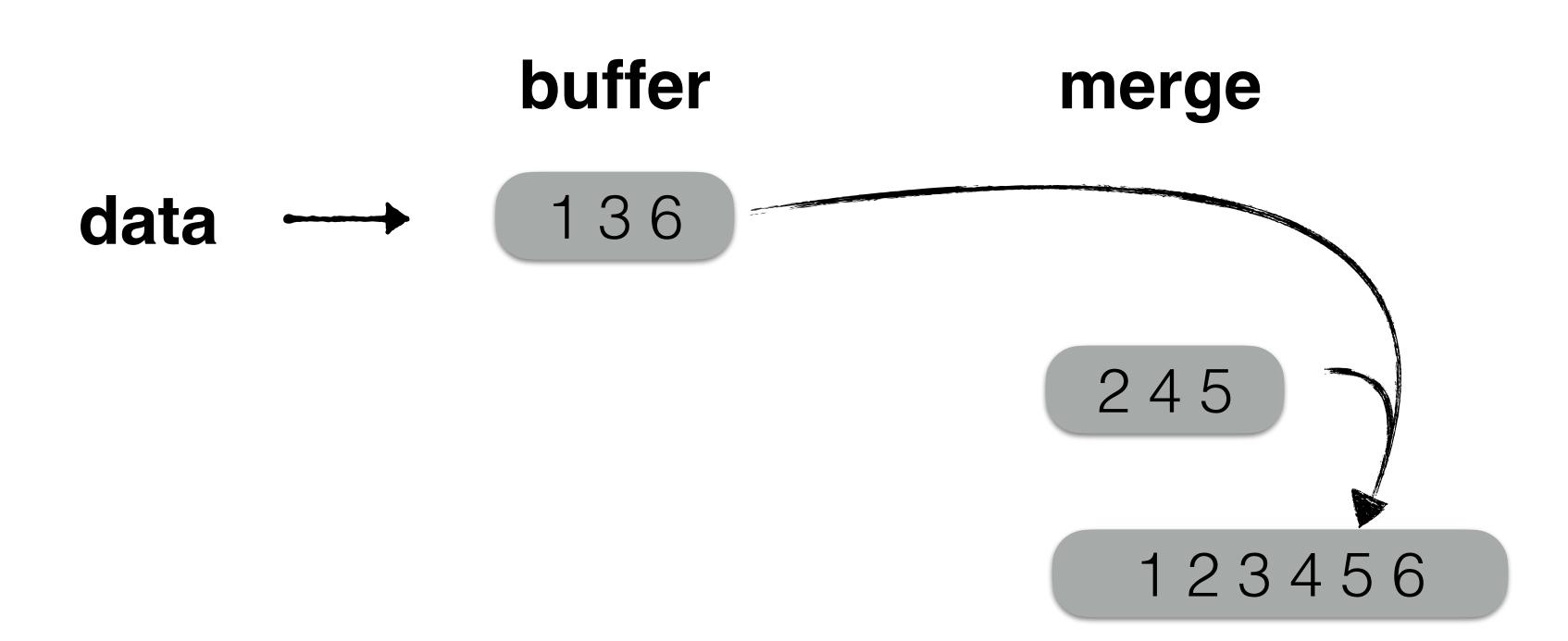








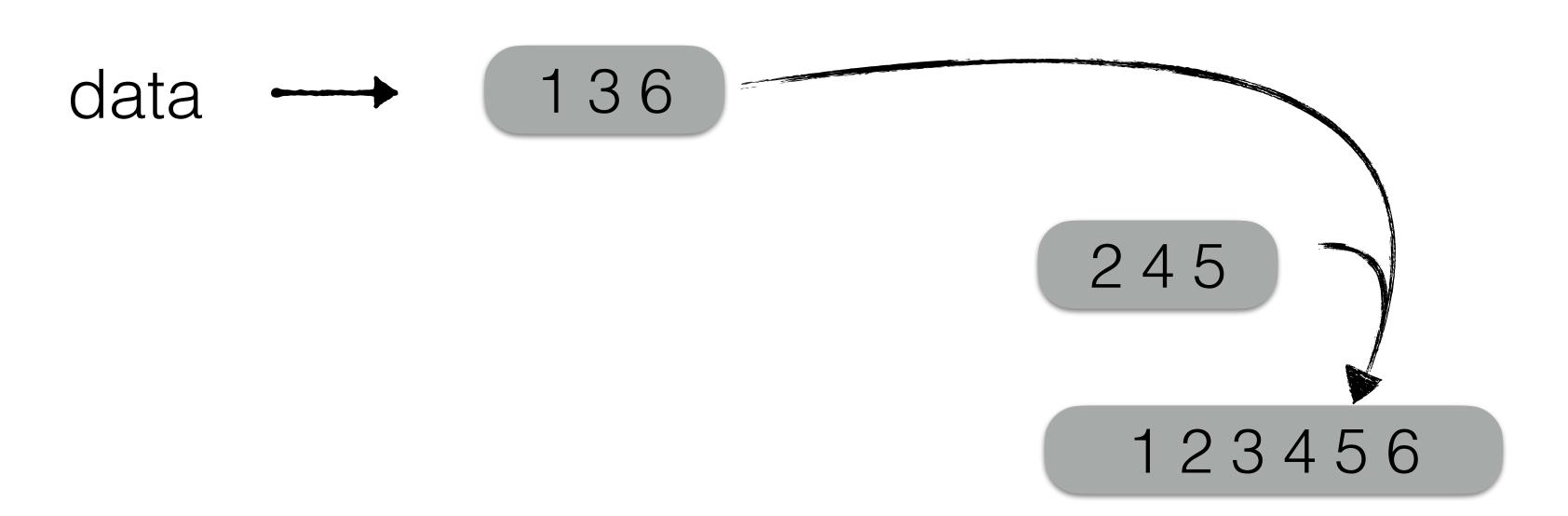




only sequential writes





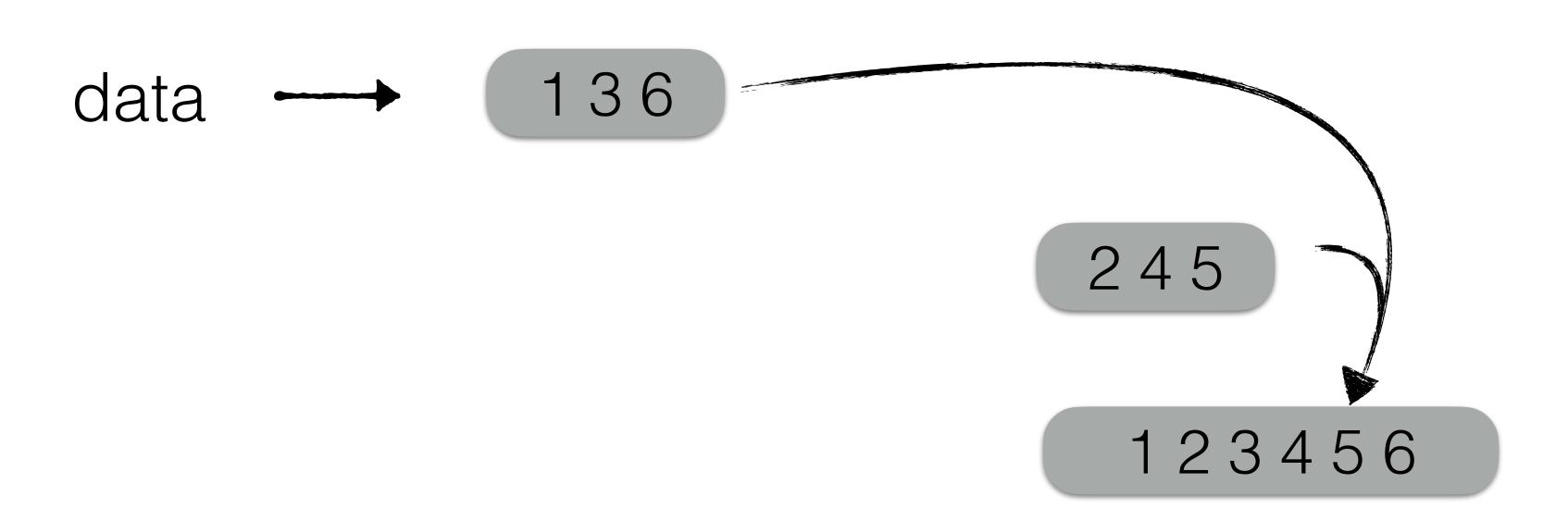


only sequential writes

write-amplification







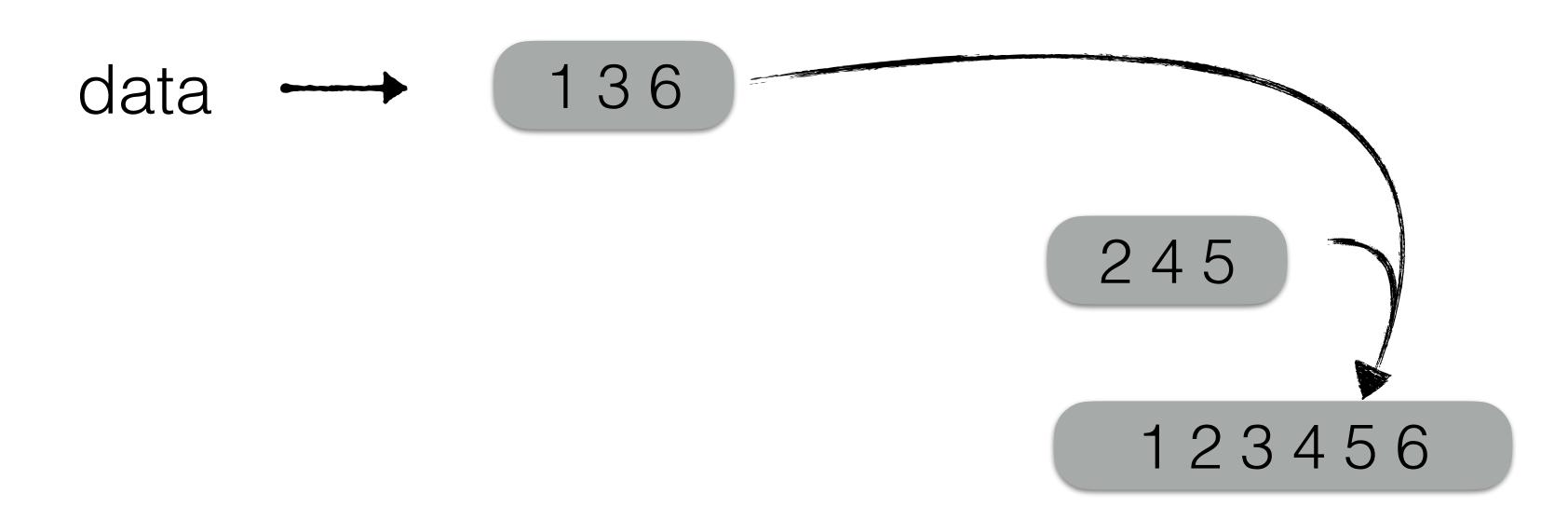
only sequential writes

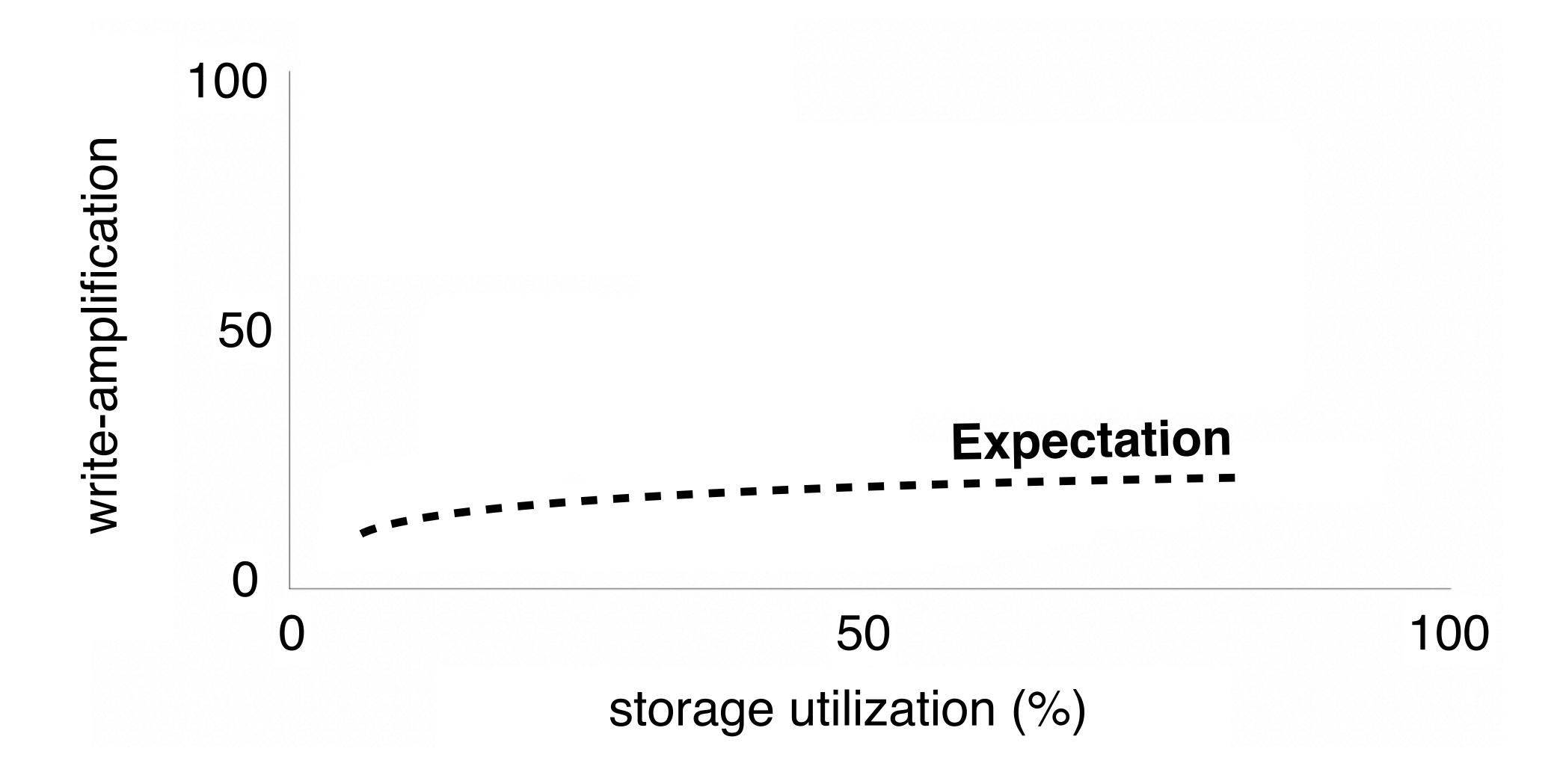
SSD

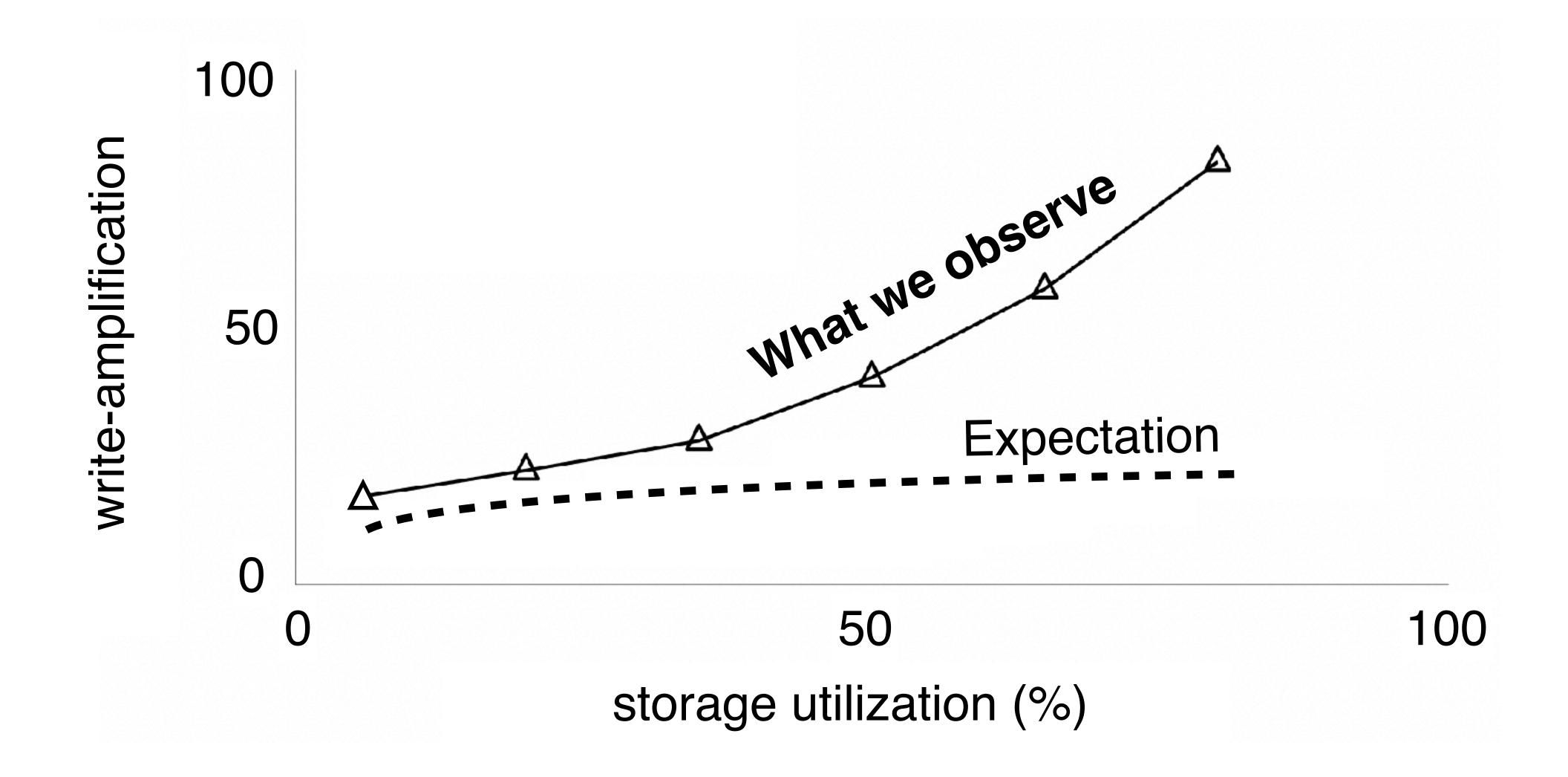
write-amplification

O(Log N)





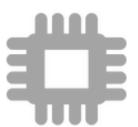




key-value pairs



buffer

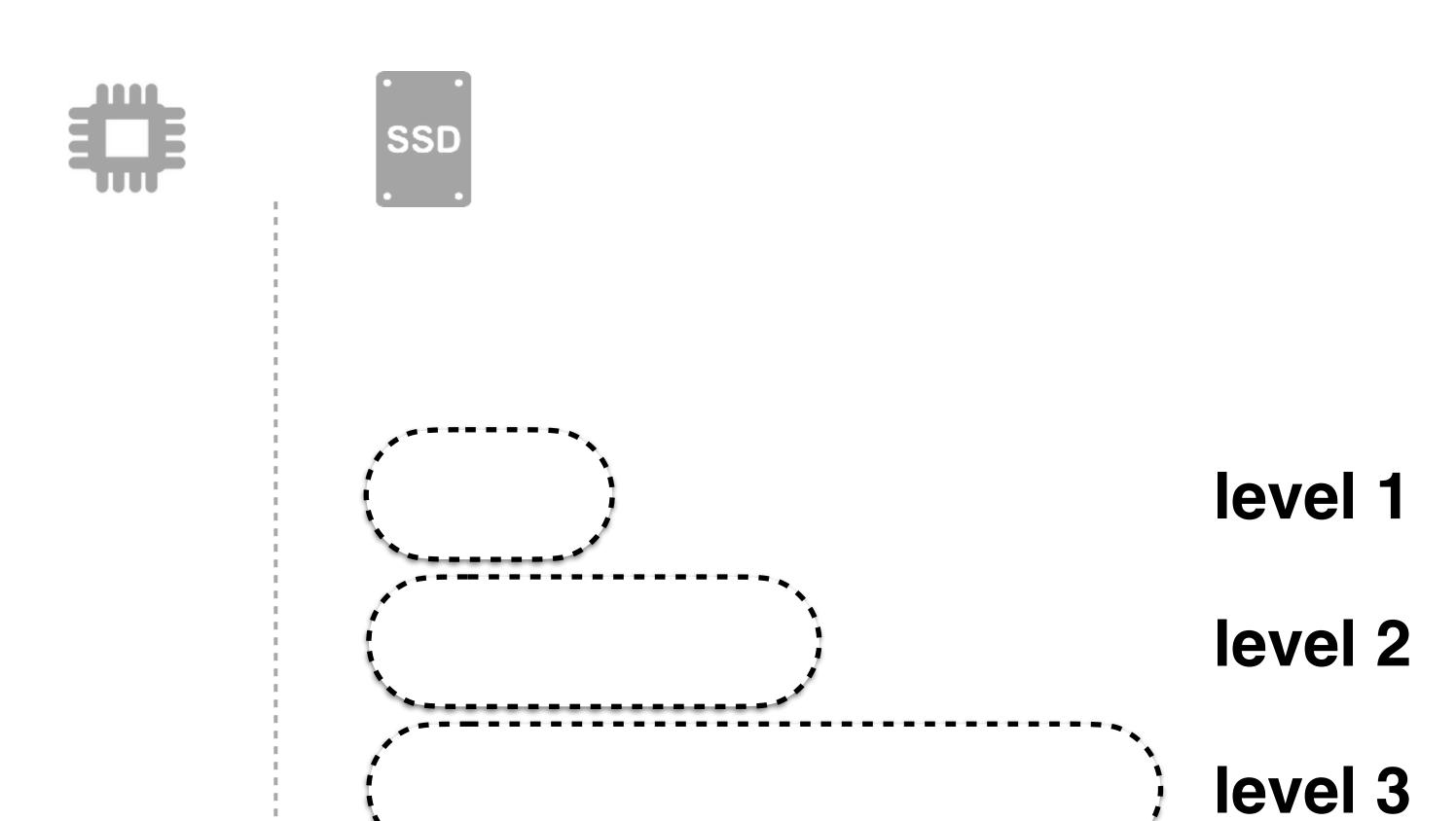


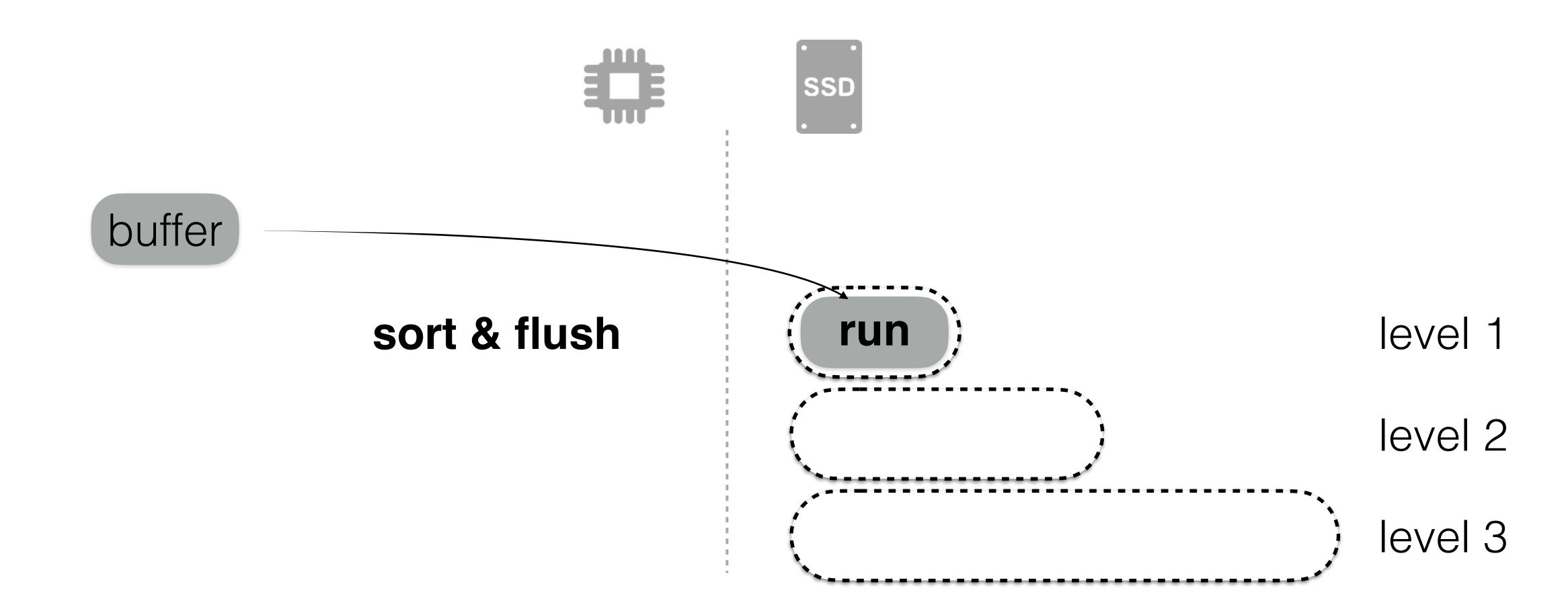


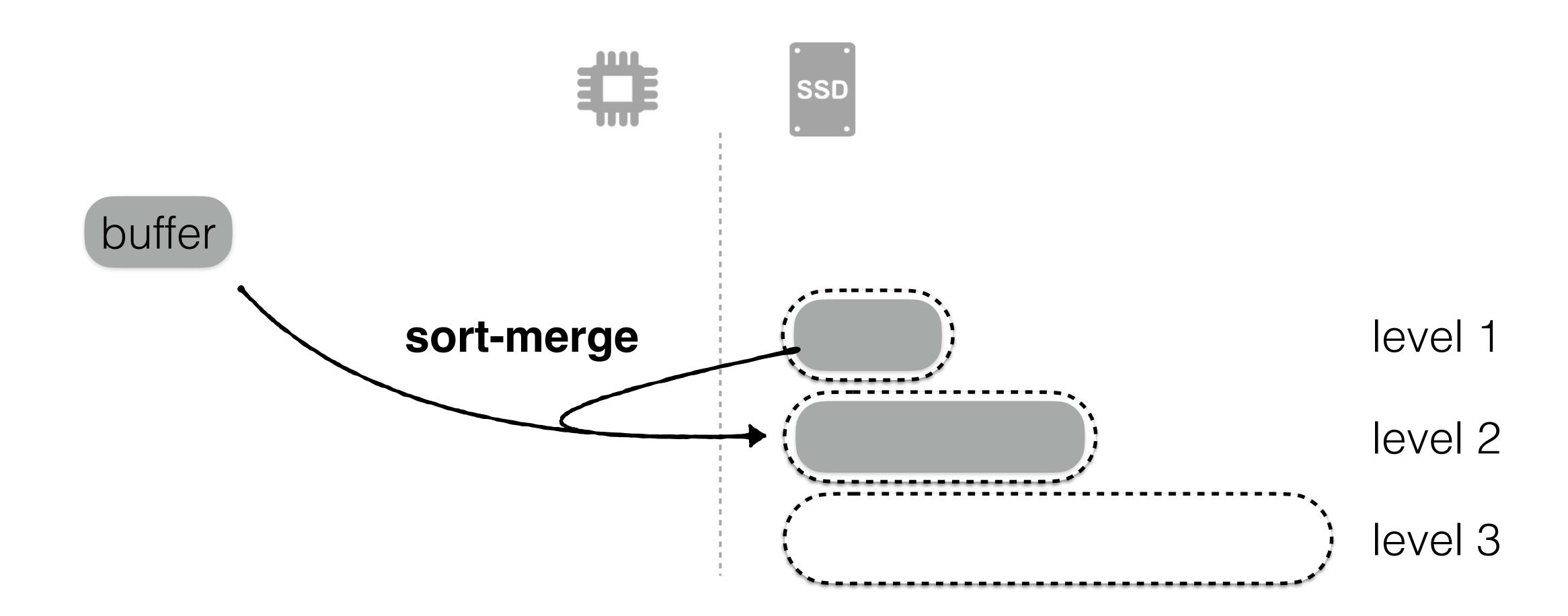
key-value pairs

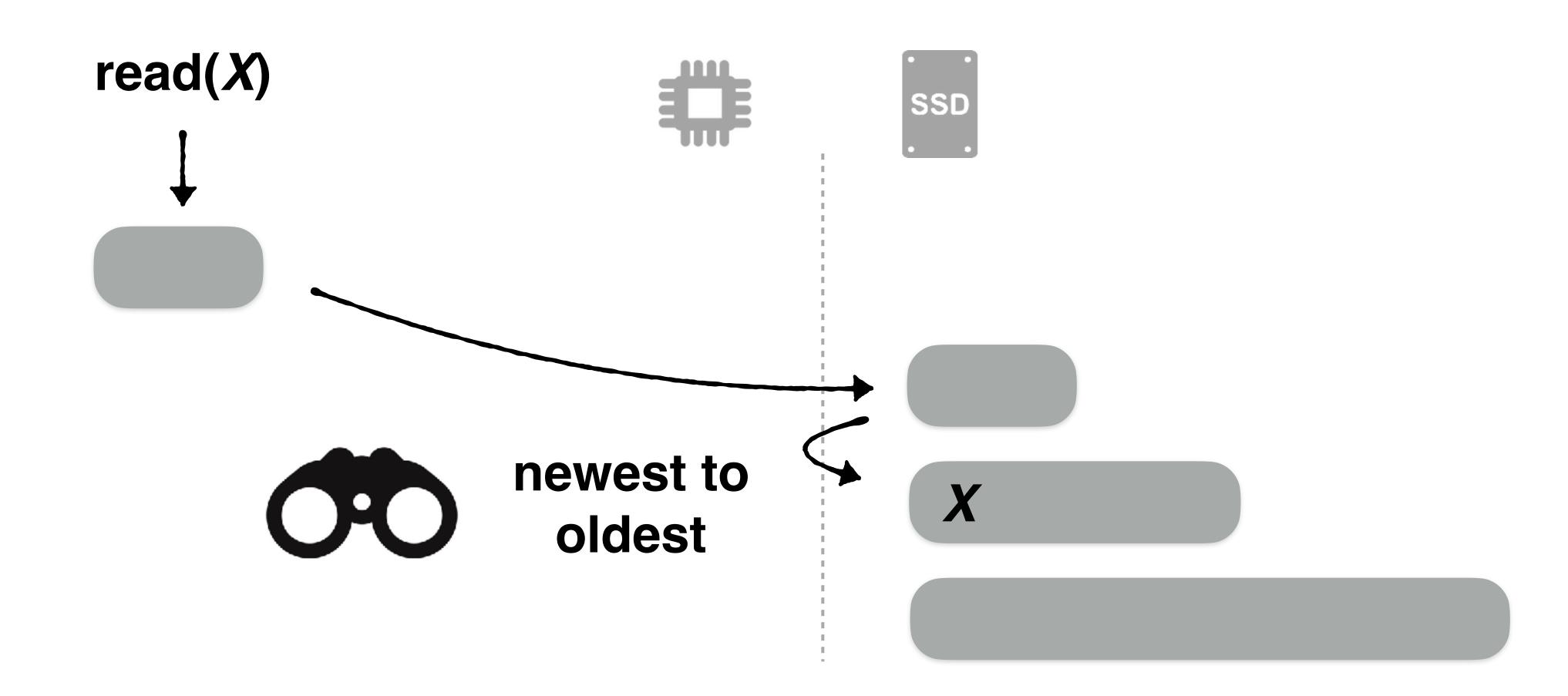
t

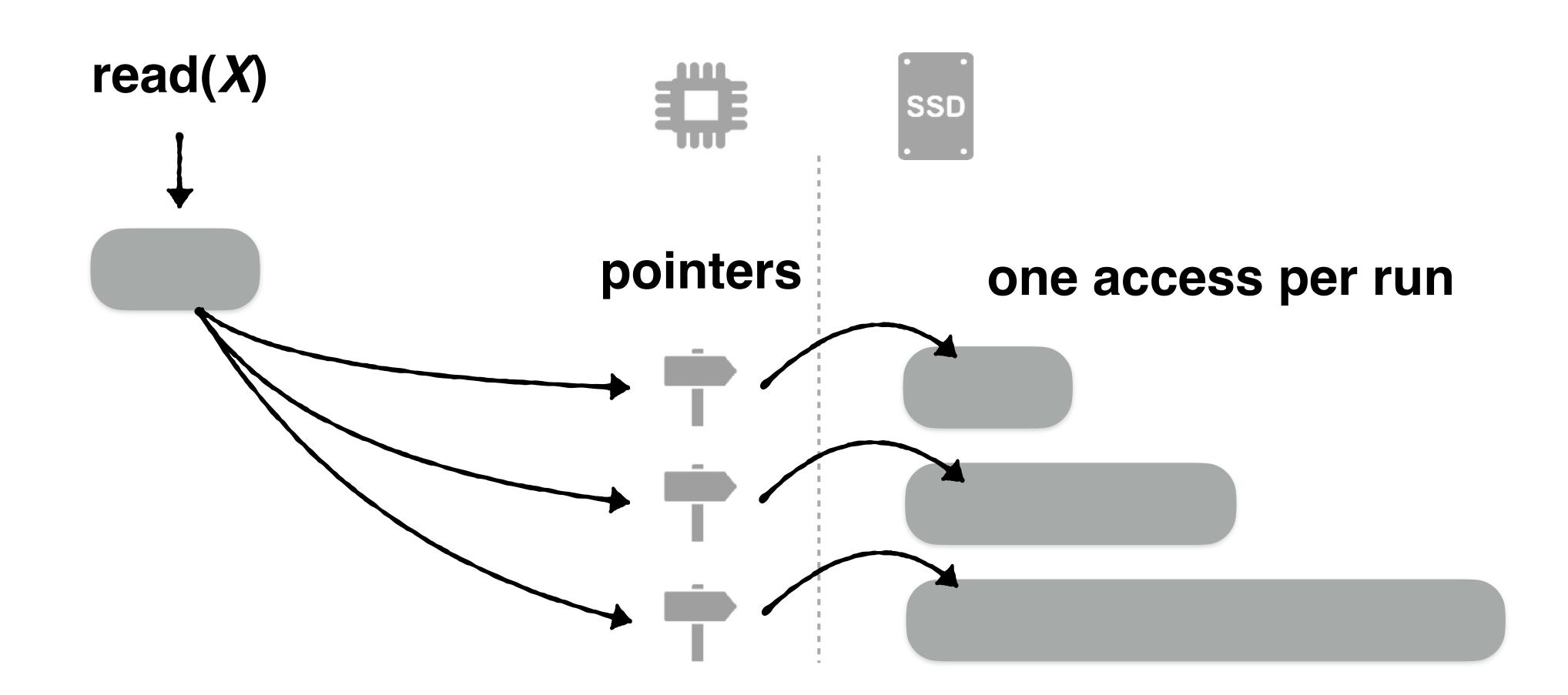
buffer

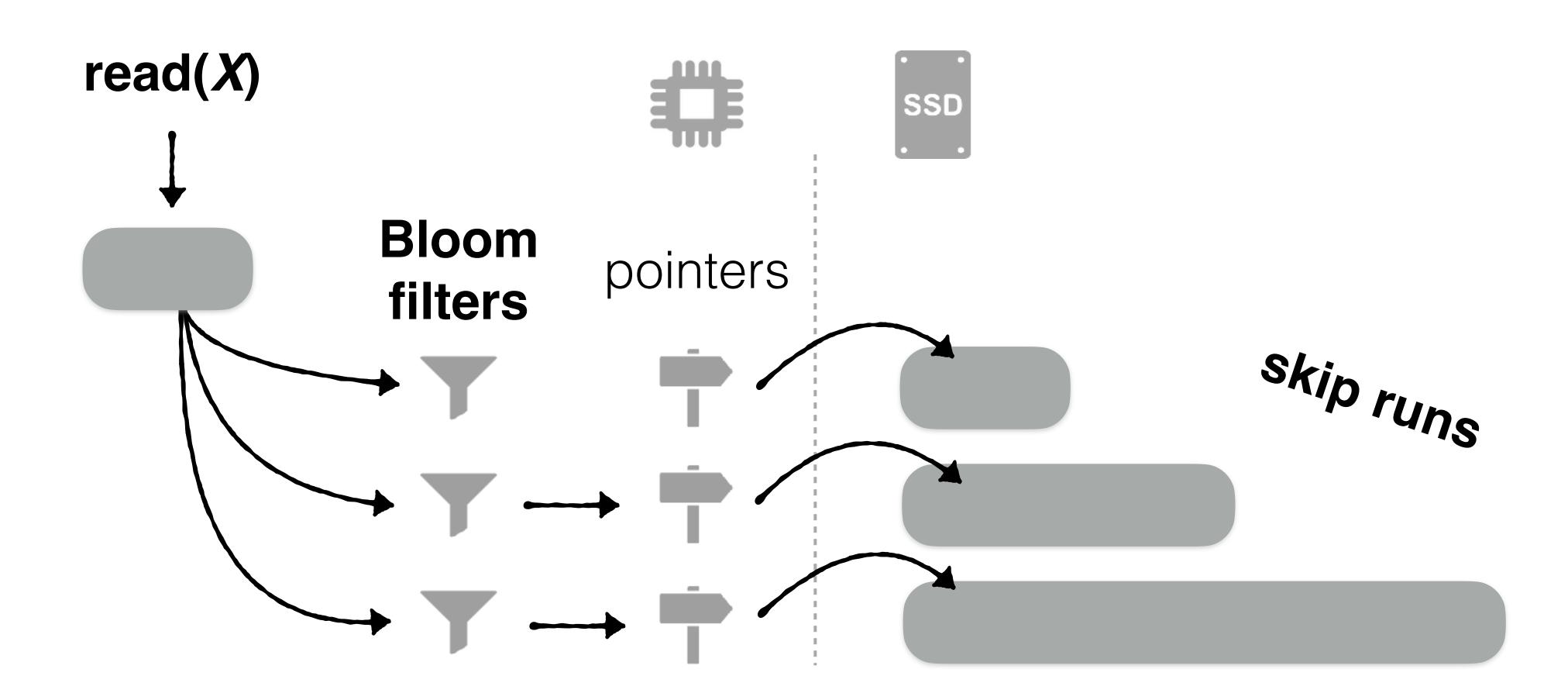


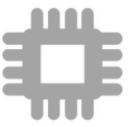




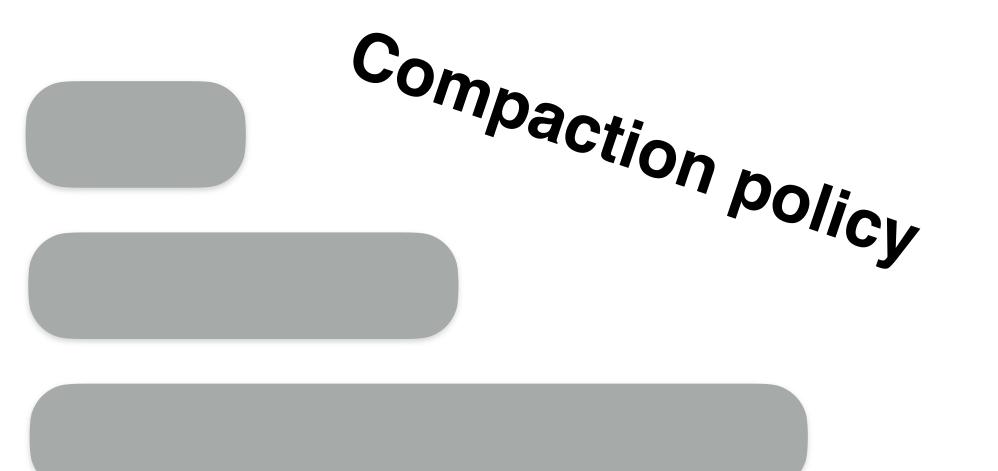












Compaction policy

Writeoptimized

Eagerness

Readoptimized





Compaction policy



Writeoptimized

Eagerness

Readoptimized



3 1 Tiering

2, 7 1, 5

1, 3, 6, 7

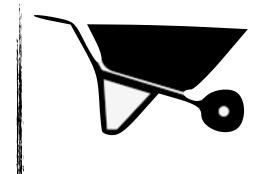
2, 3, 5, 9

2, 5 Leveling

1, 3, 6, 8

1, 2, 3, 5, 6, 8, 9

Compaction policy









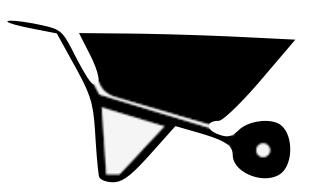


Compaction Granularity

Full Merge



Partial Merge



Compaction Granularity

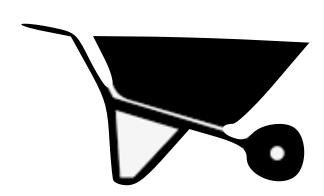
Full Merge







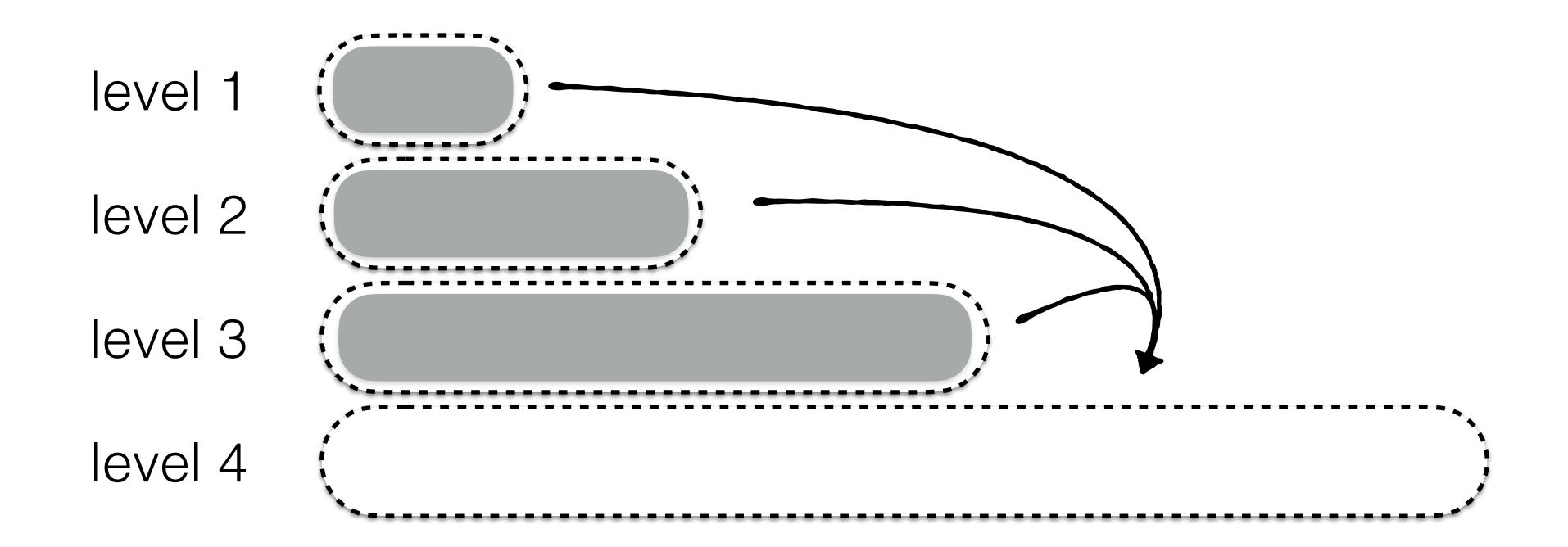
Partial Merge



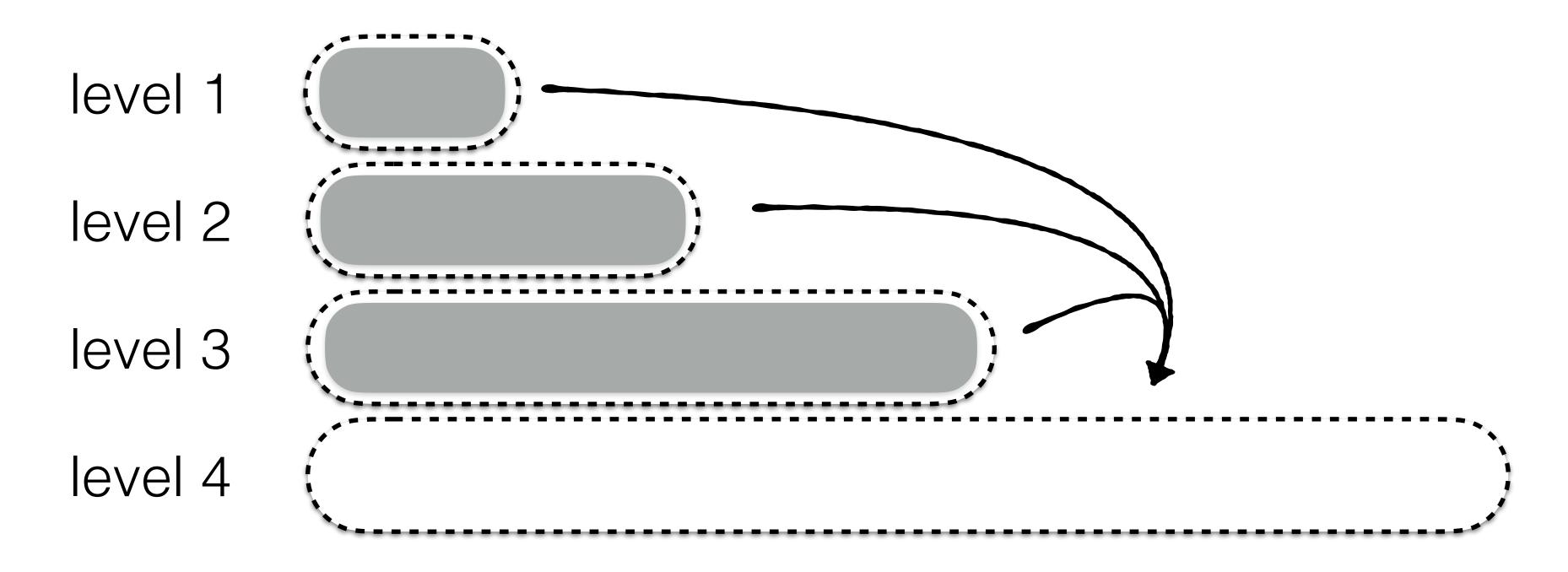




Merge consecutive full levels into first non-full level

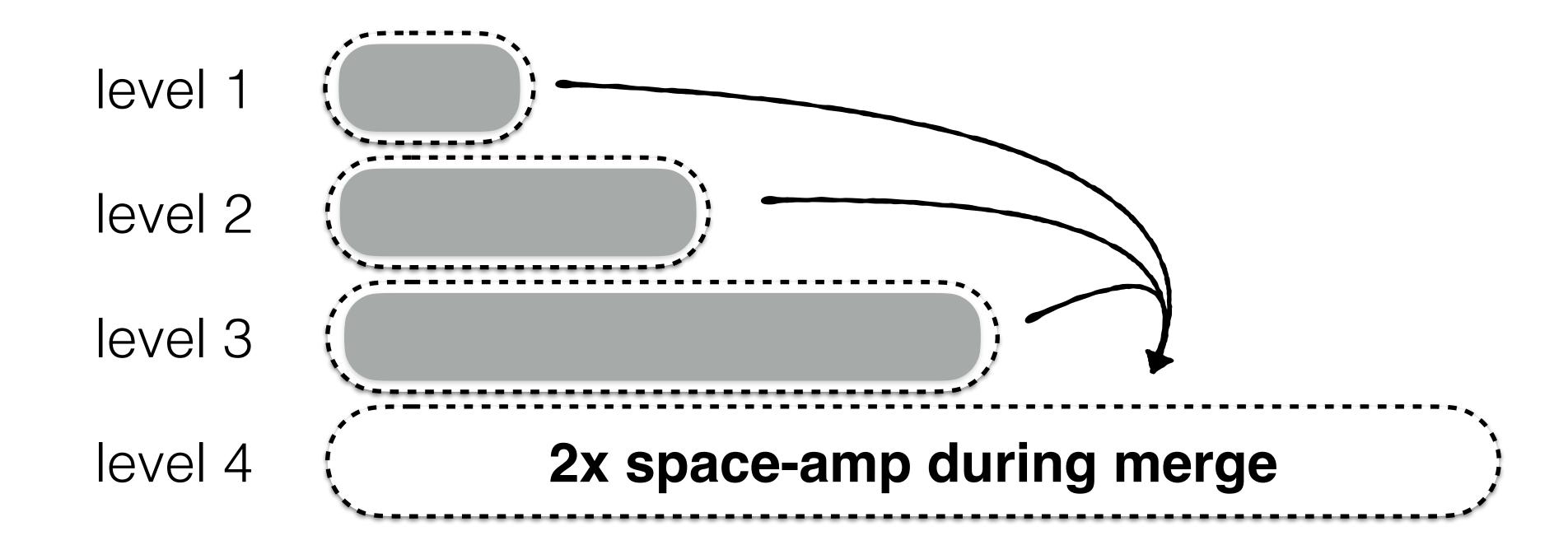


Merge consecutive full levels into first non-full level



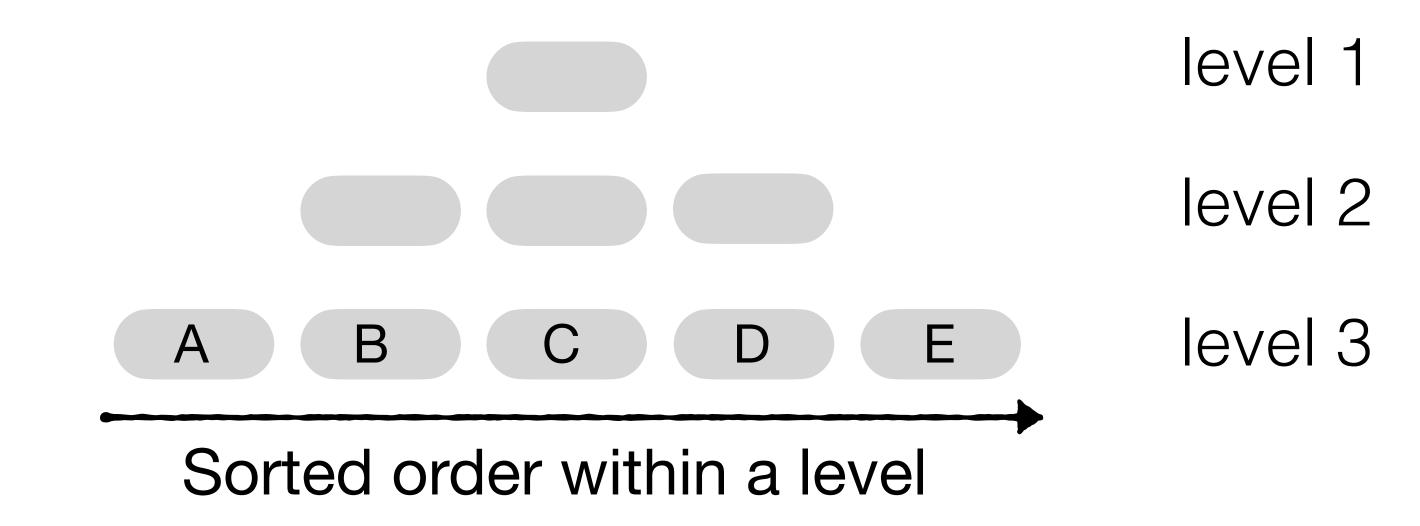
Problem?

Merge consecutive full levels into first non-full level

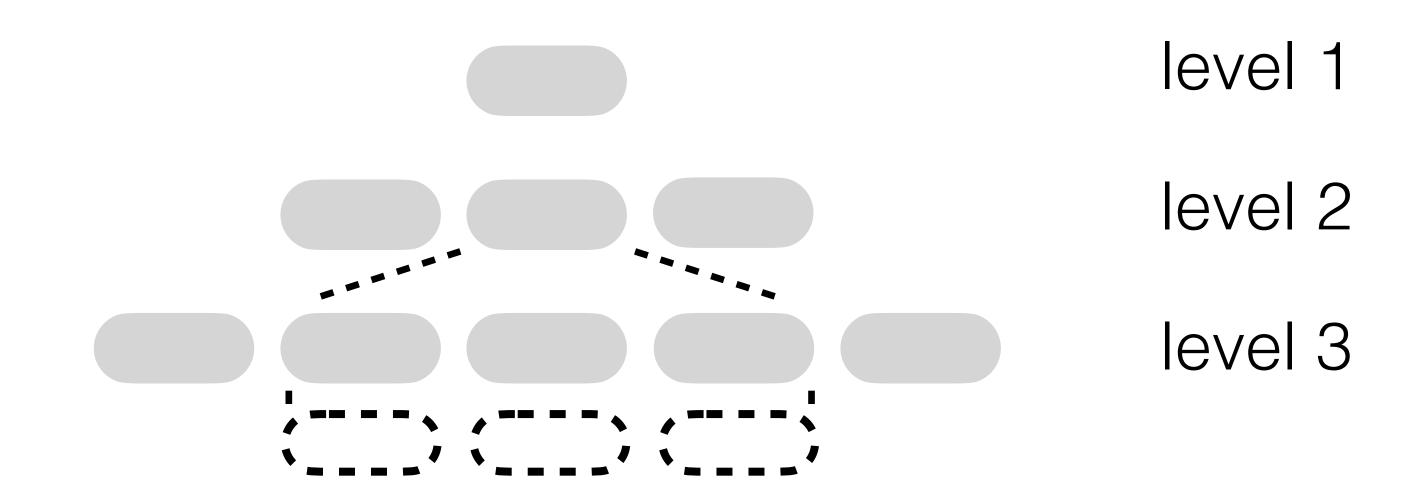


Partial Merge

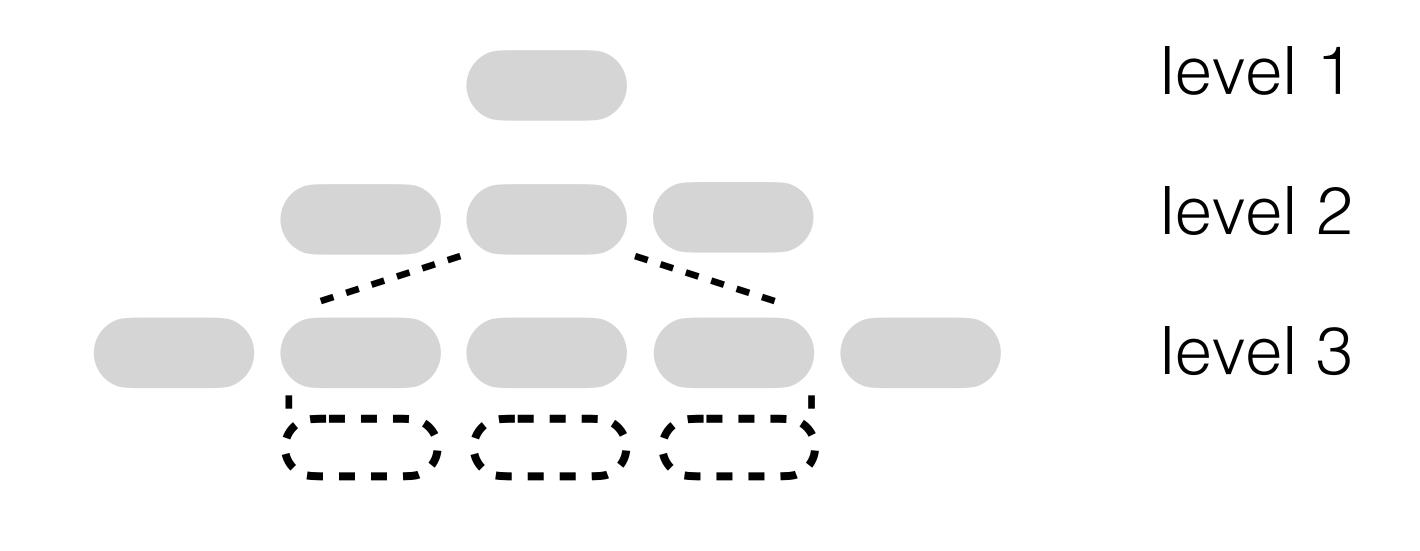
1. split runs into many files (SSTs) in each level



- 1. split runs into many files (SSTs) in each level
- 2. When a level is full, pick SST with smallest intersection into next level

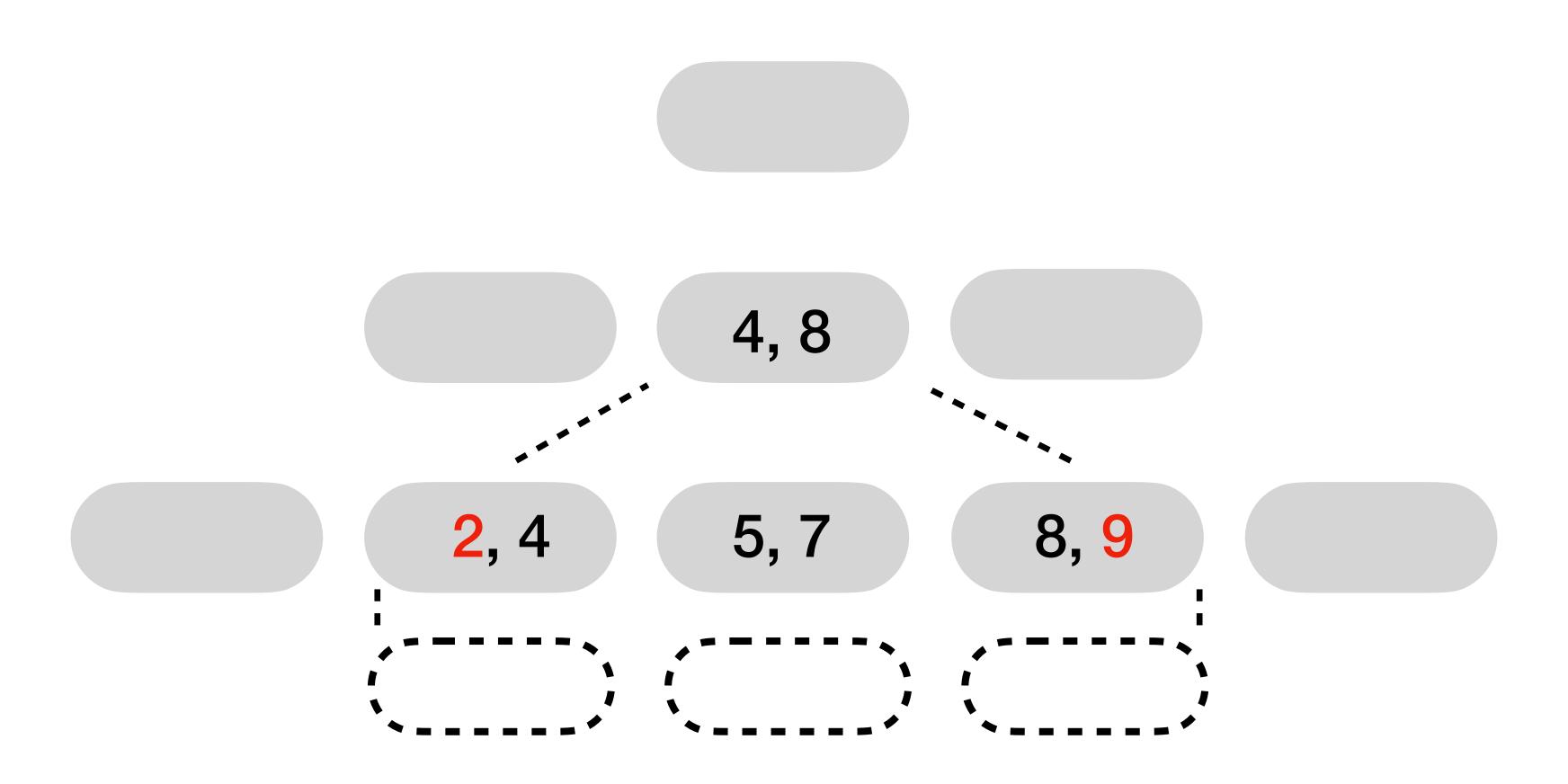


- 1. split runs into many files (SSTs) in each level
- 2. When a level is full, pick SST with smallest intersection into next level

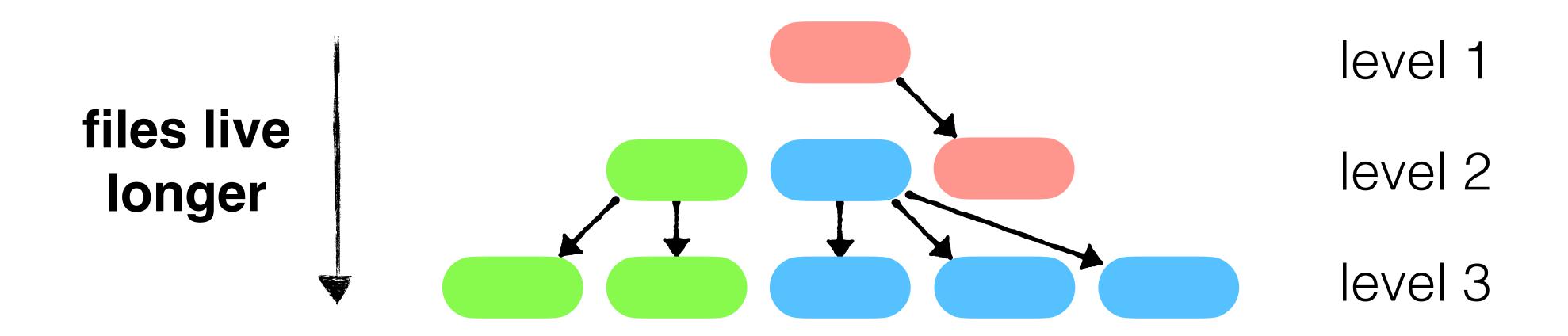


Problem?

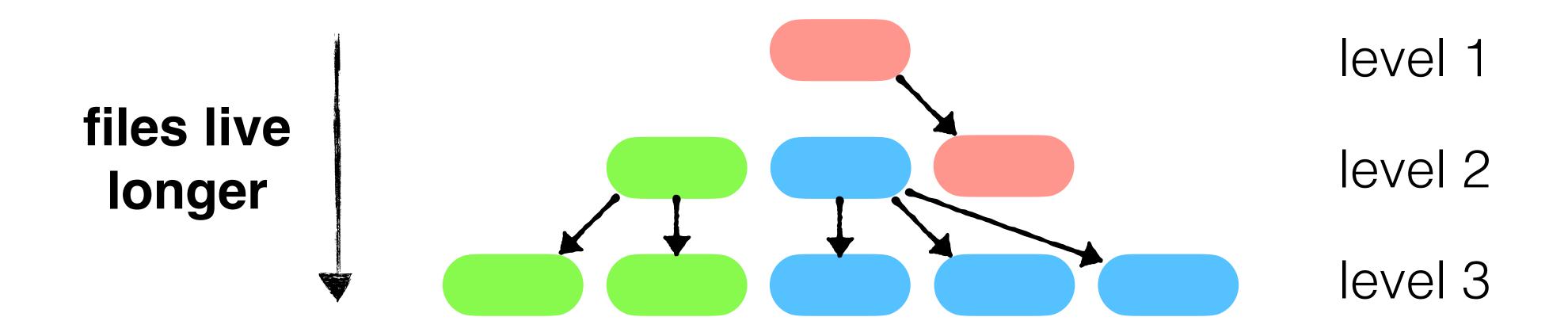
Problem 1: non-intersecting entries increase write-amplification

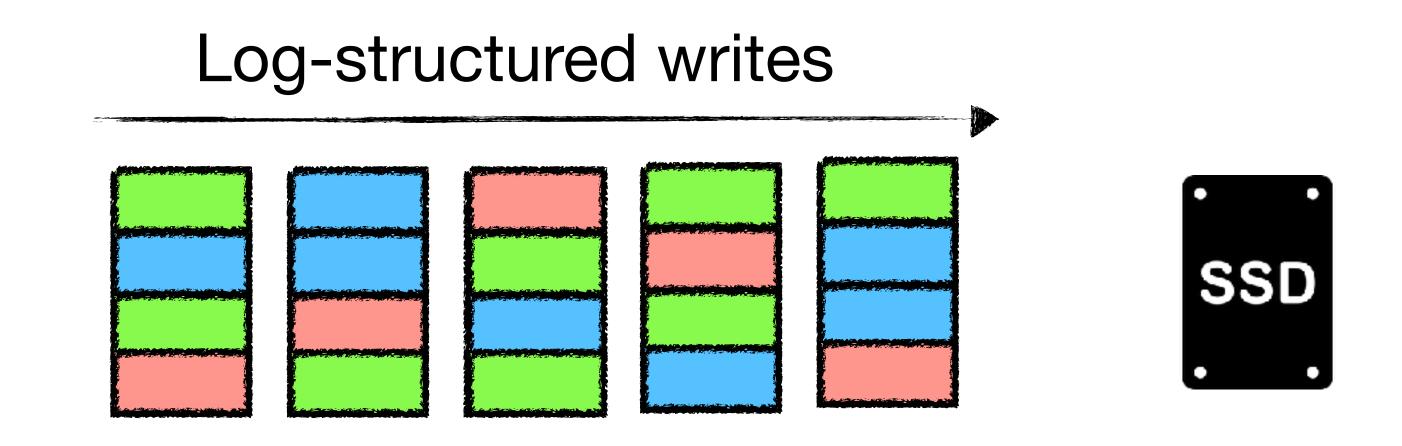


Problem 2: many small simultaneous compactions

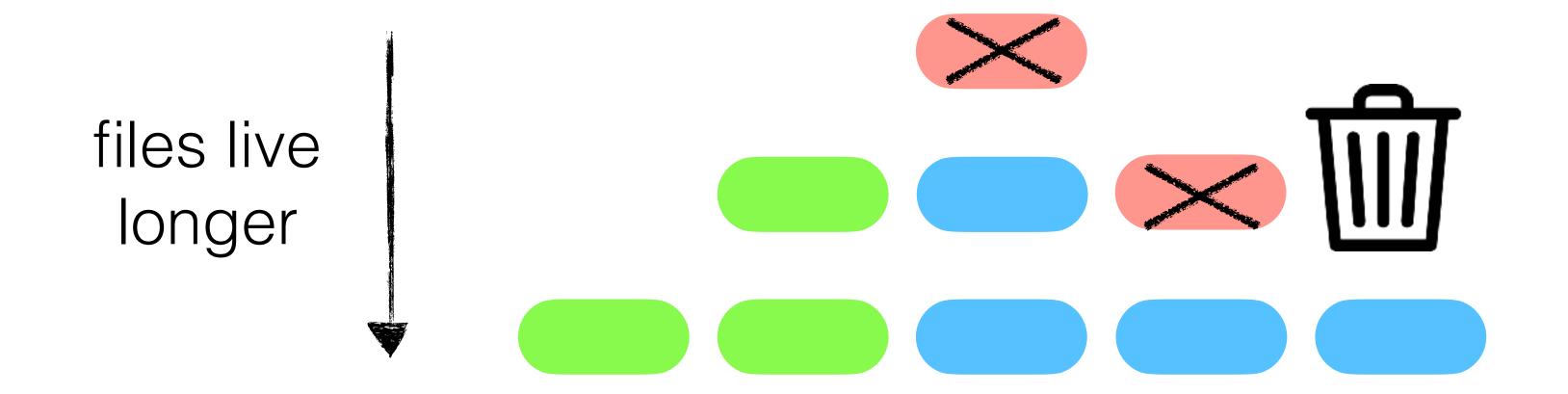


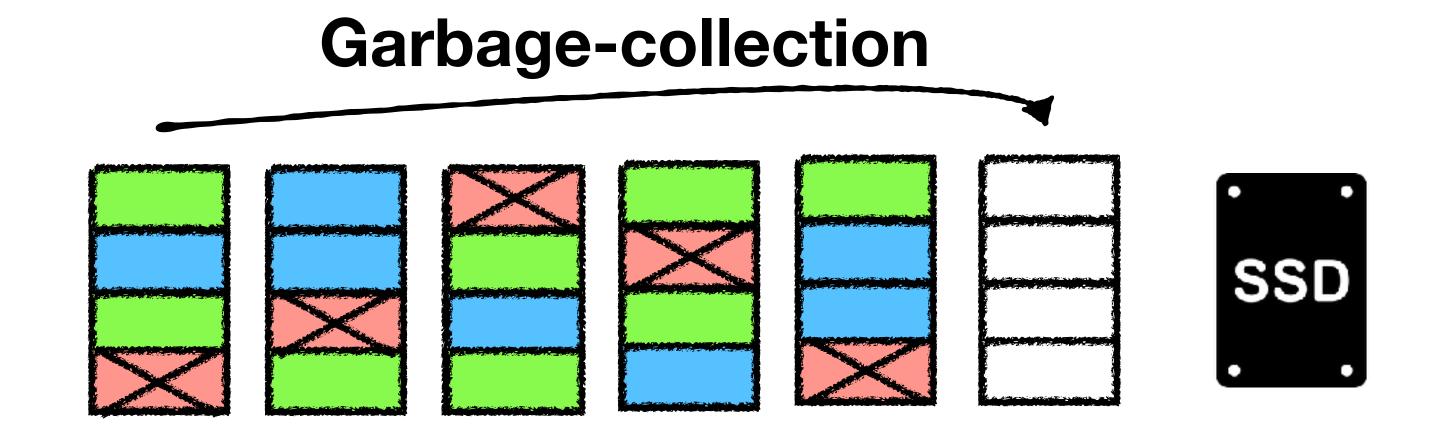
Problem 2: many small simultaneous compactions





Problem 2: many small simultaneous compactions



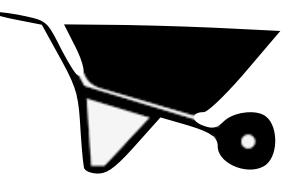




space amplification

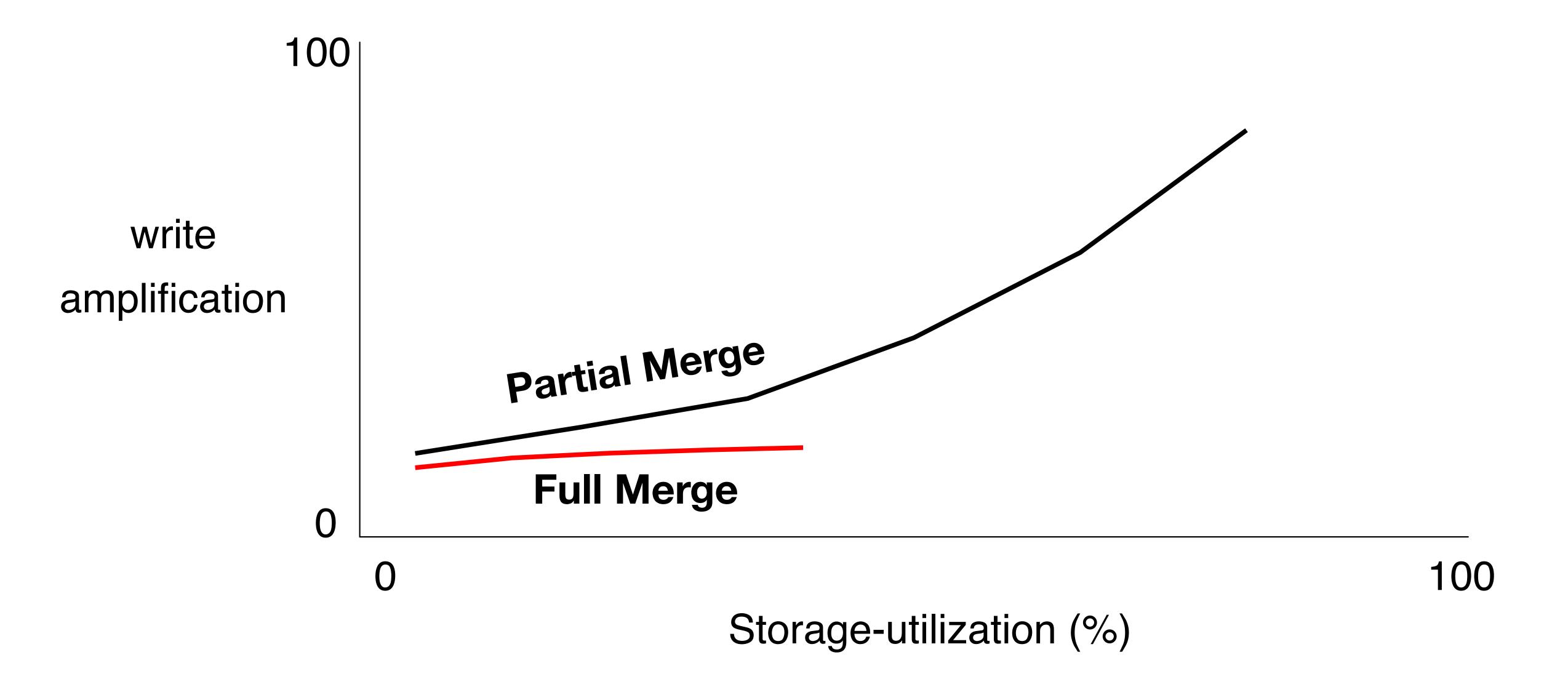


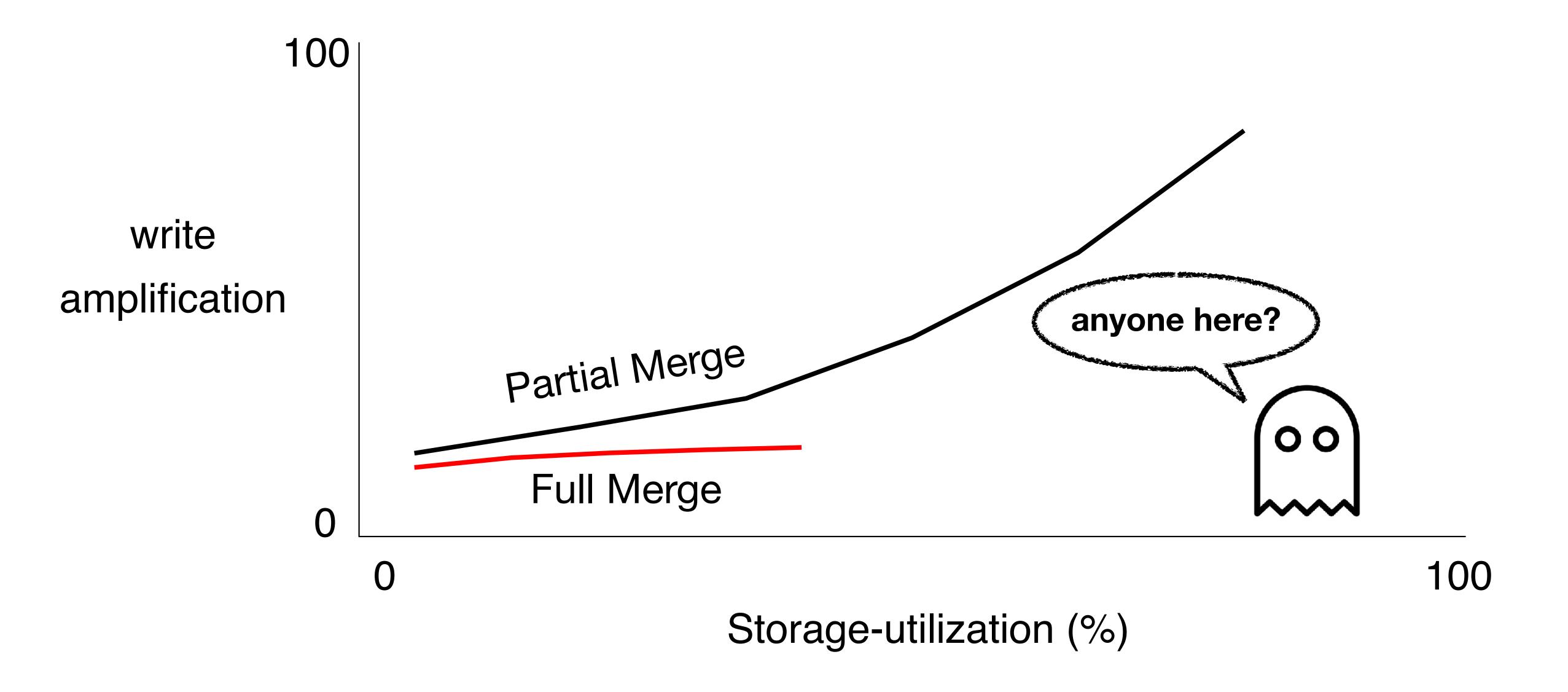
Partial Merge



write amplification



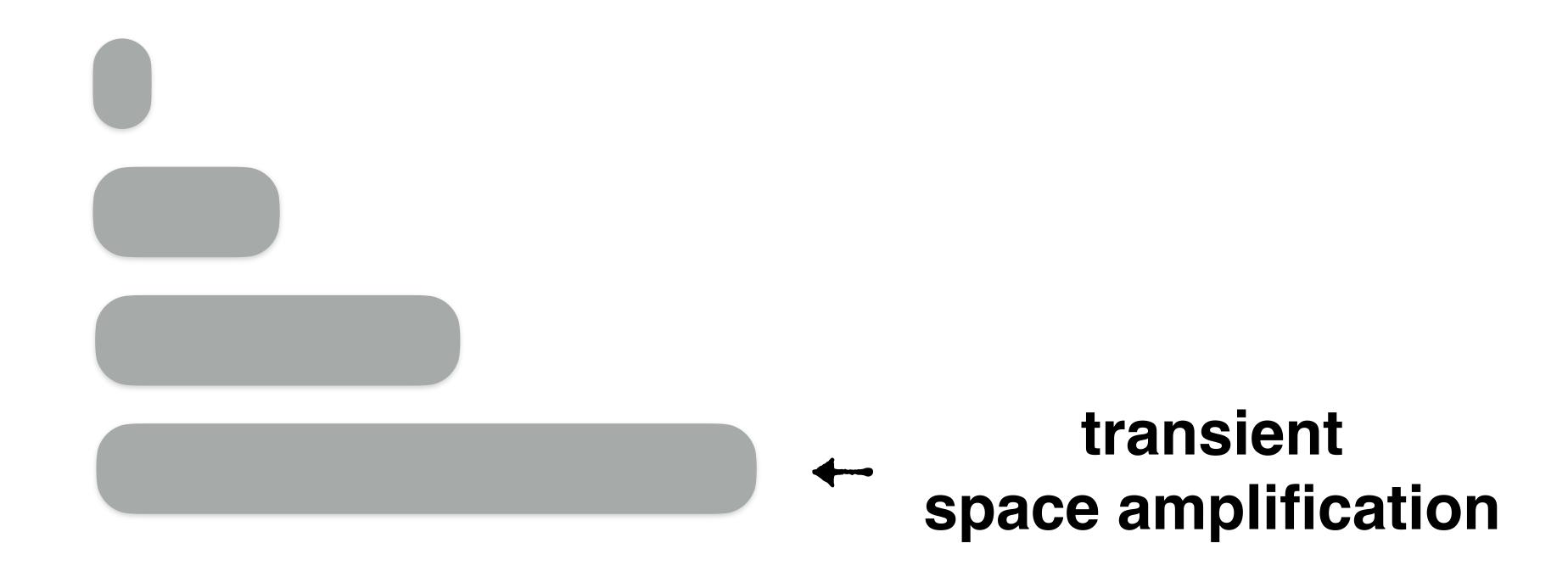


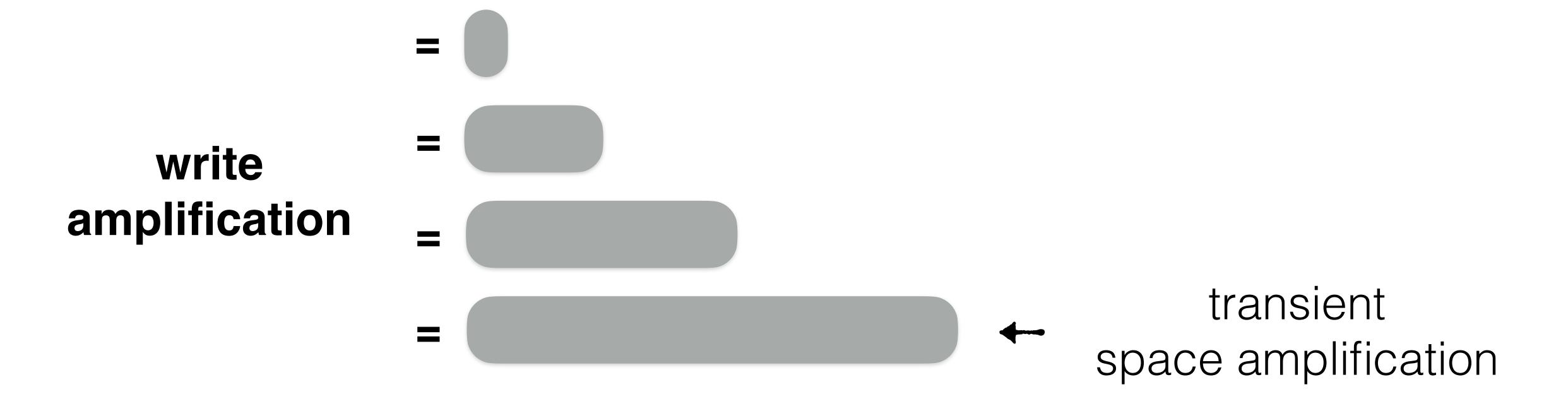


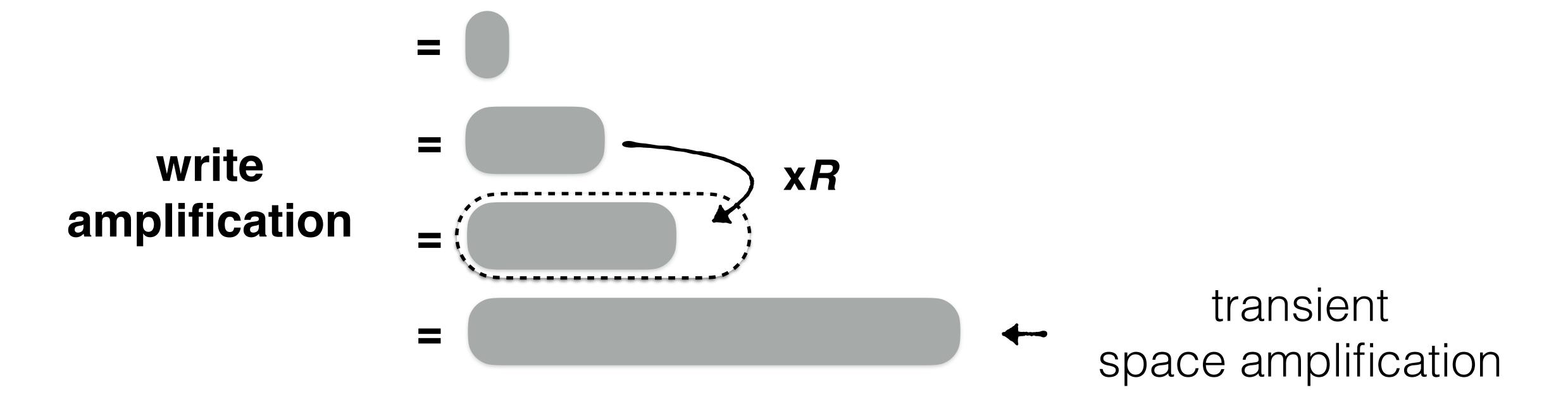


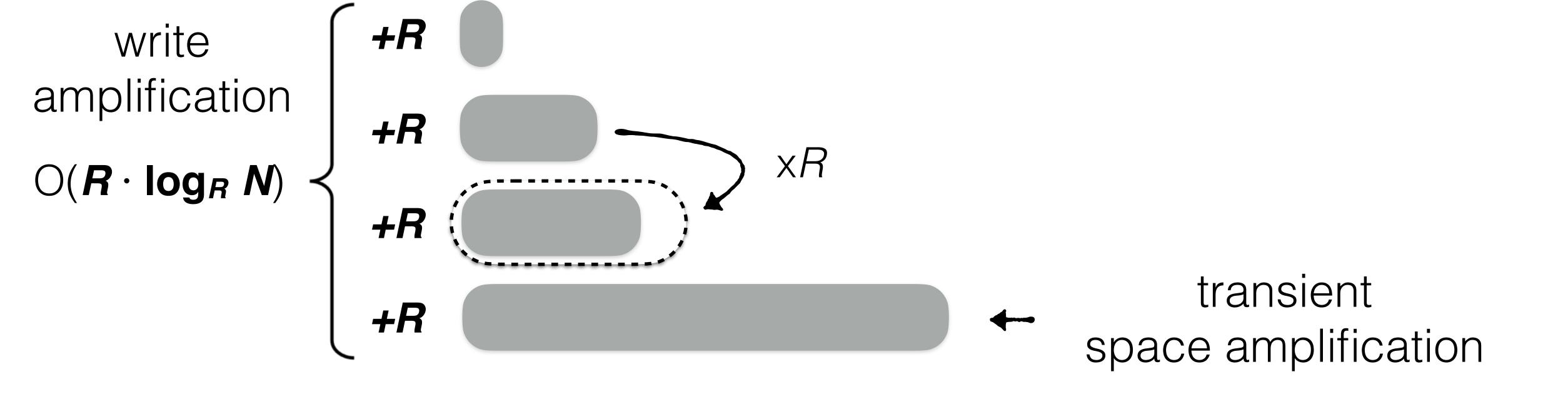
Spooky: partitioned compaction for key-value stores



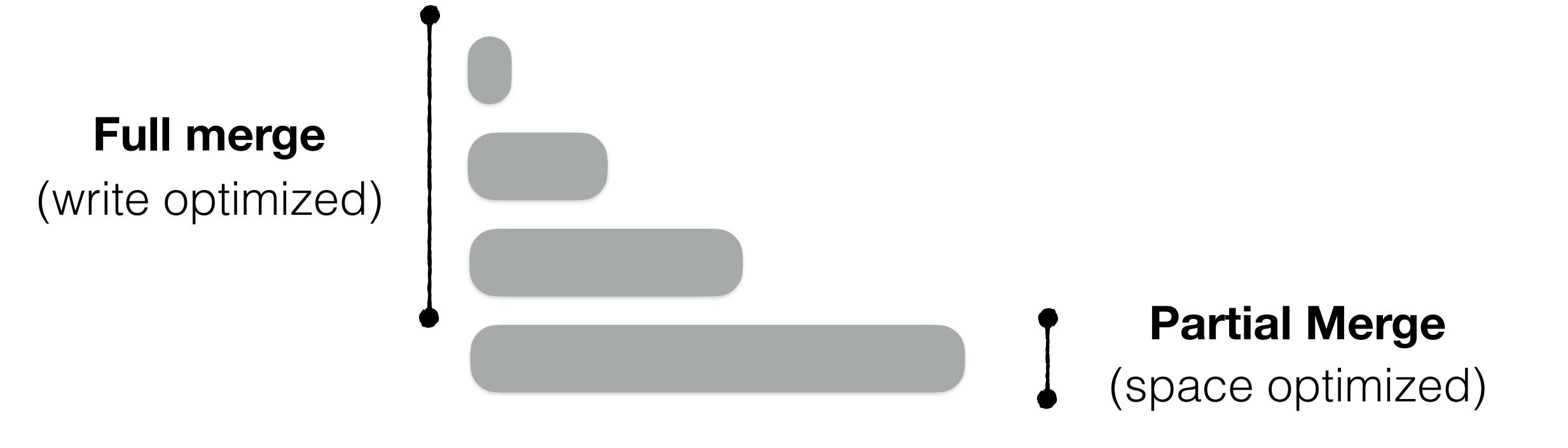




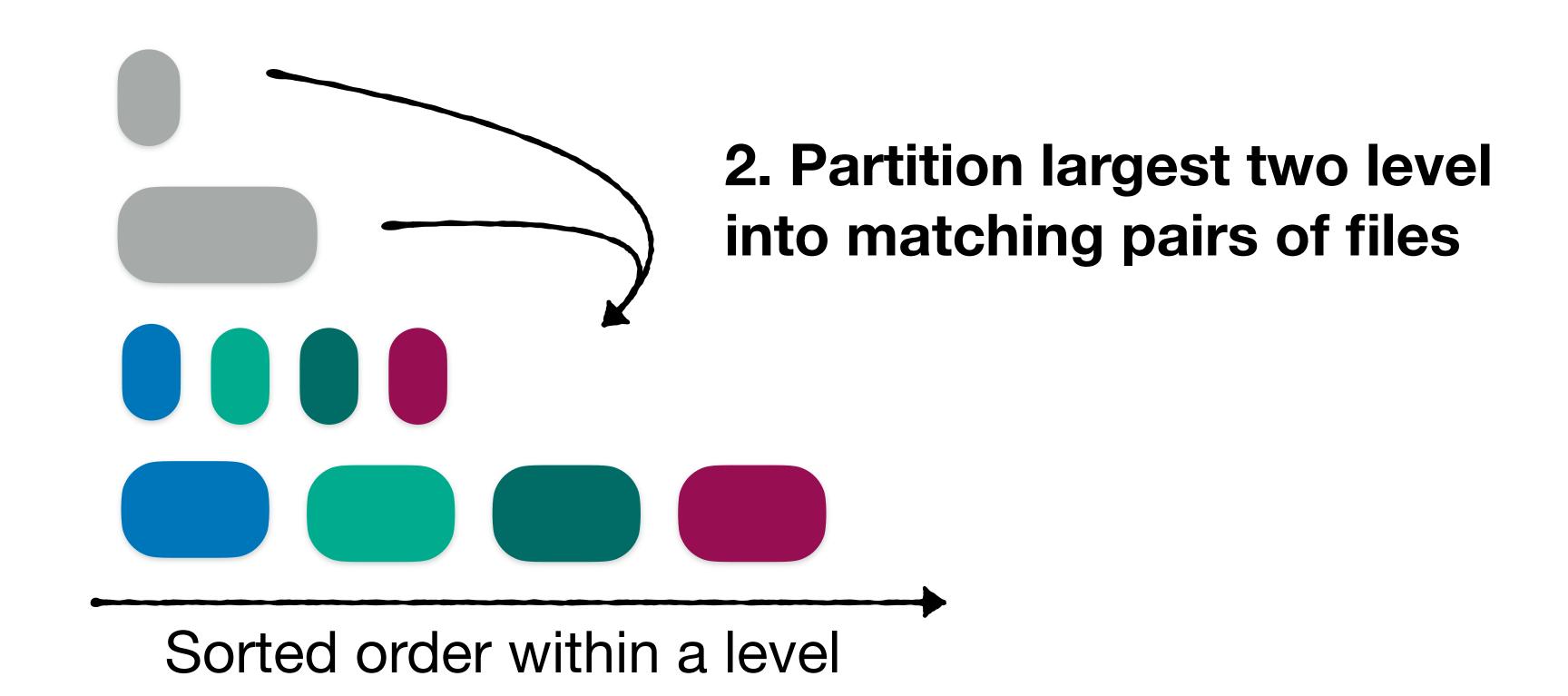




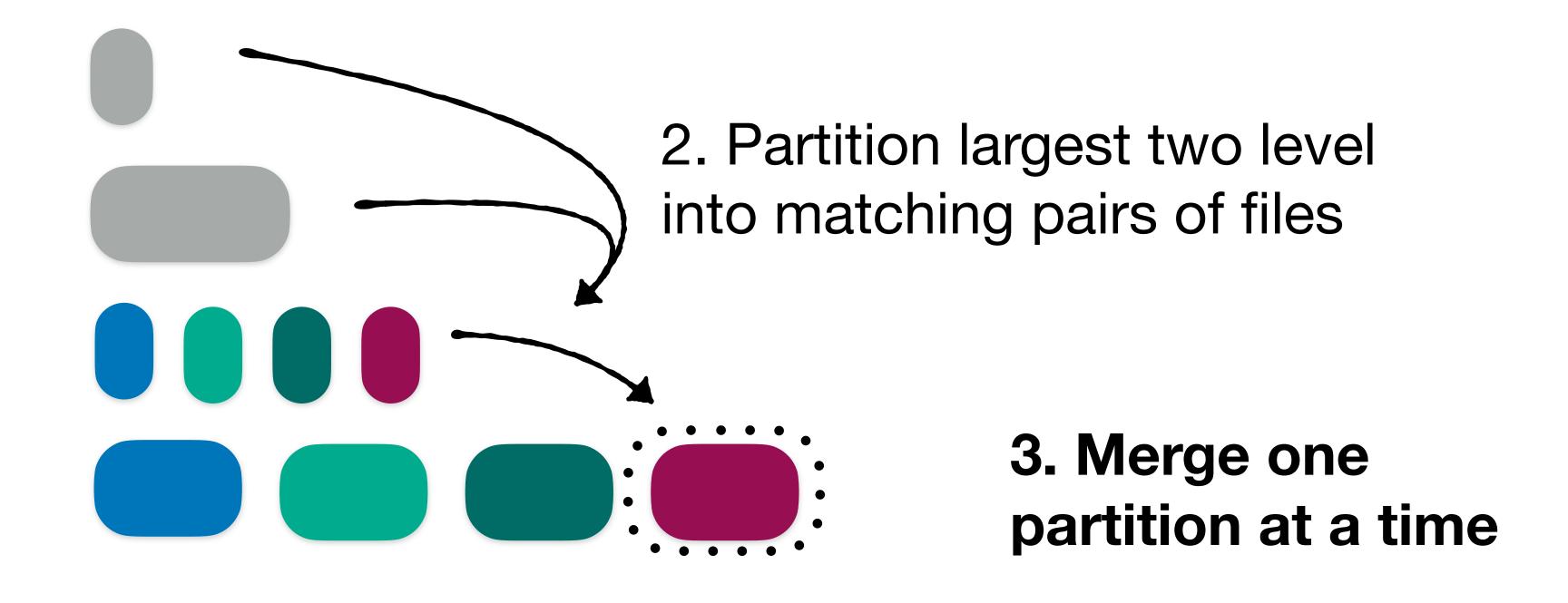




1. Full merge at up to second largest level

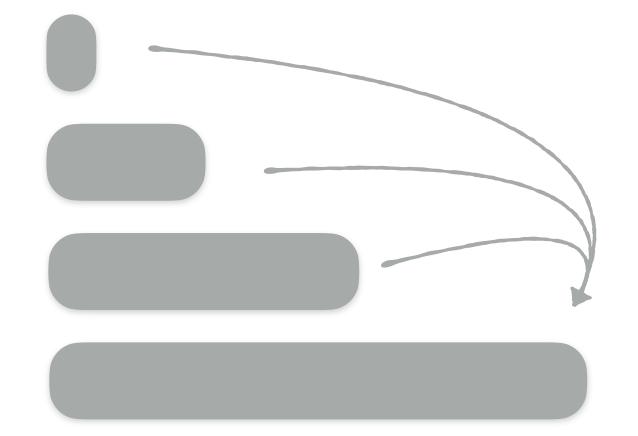


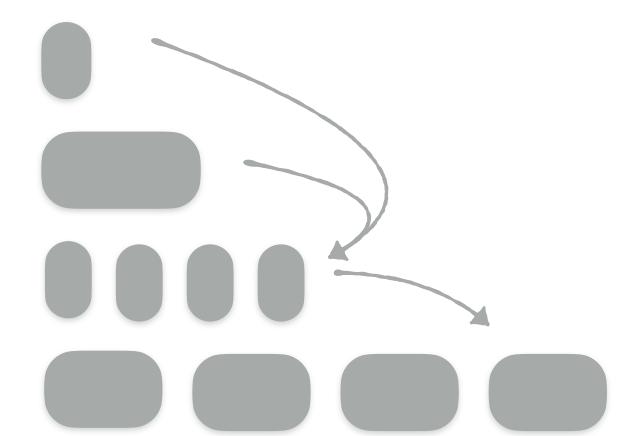
1. Full merge at up to second largest level

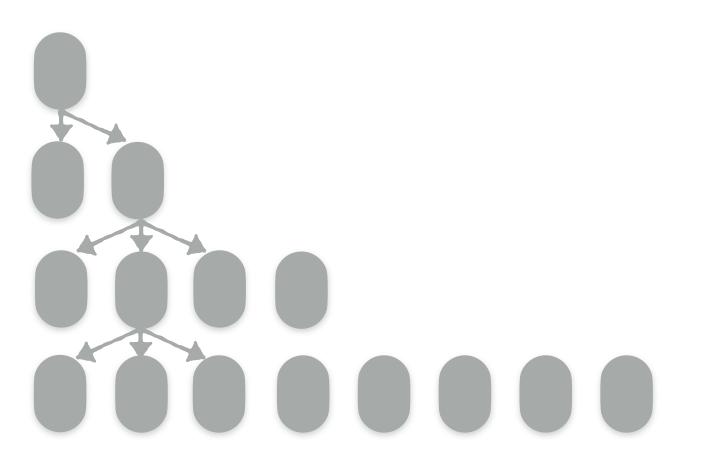


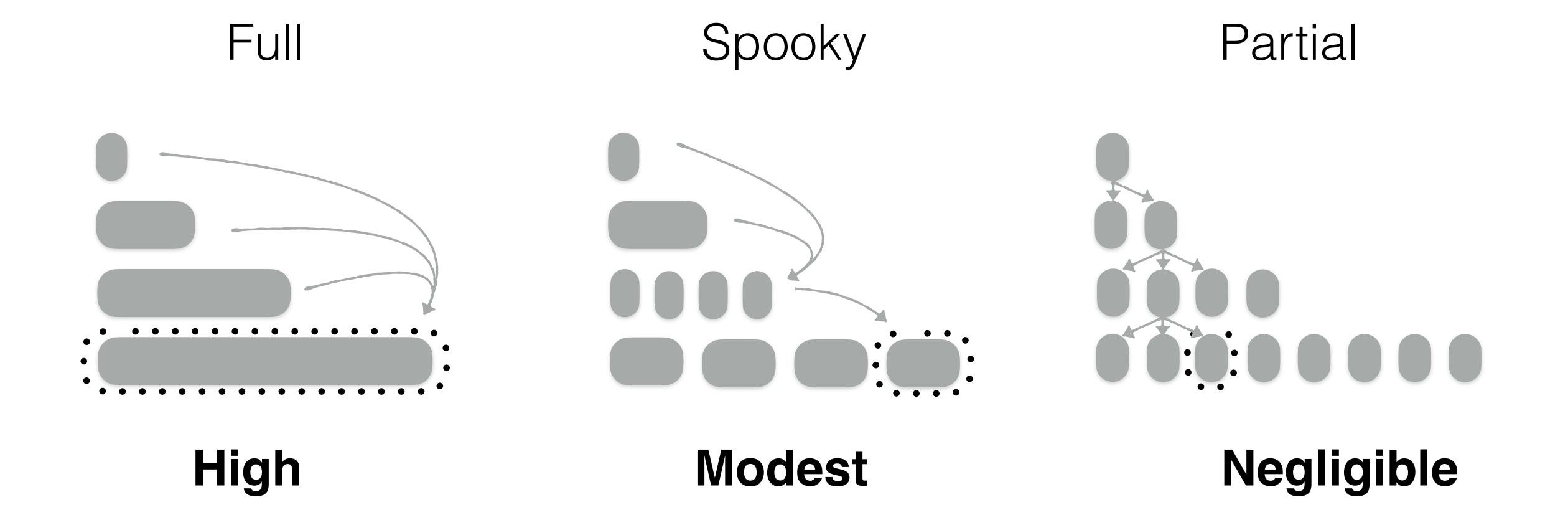
Full

Partial

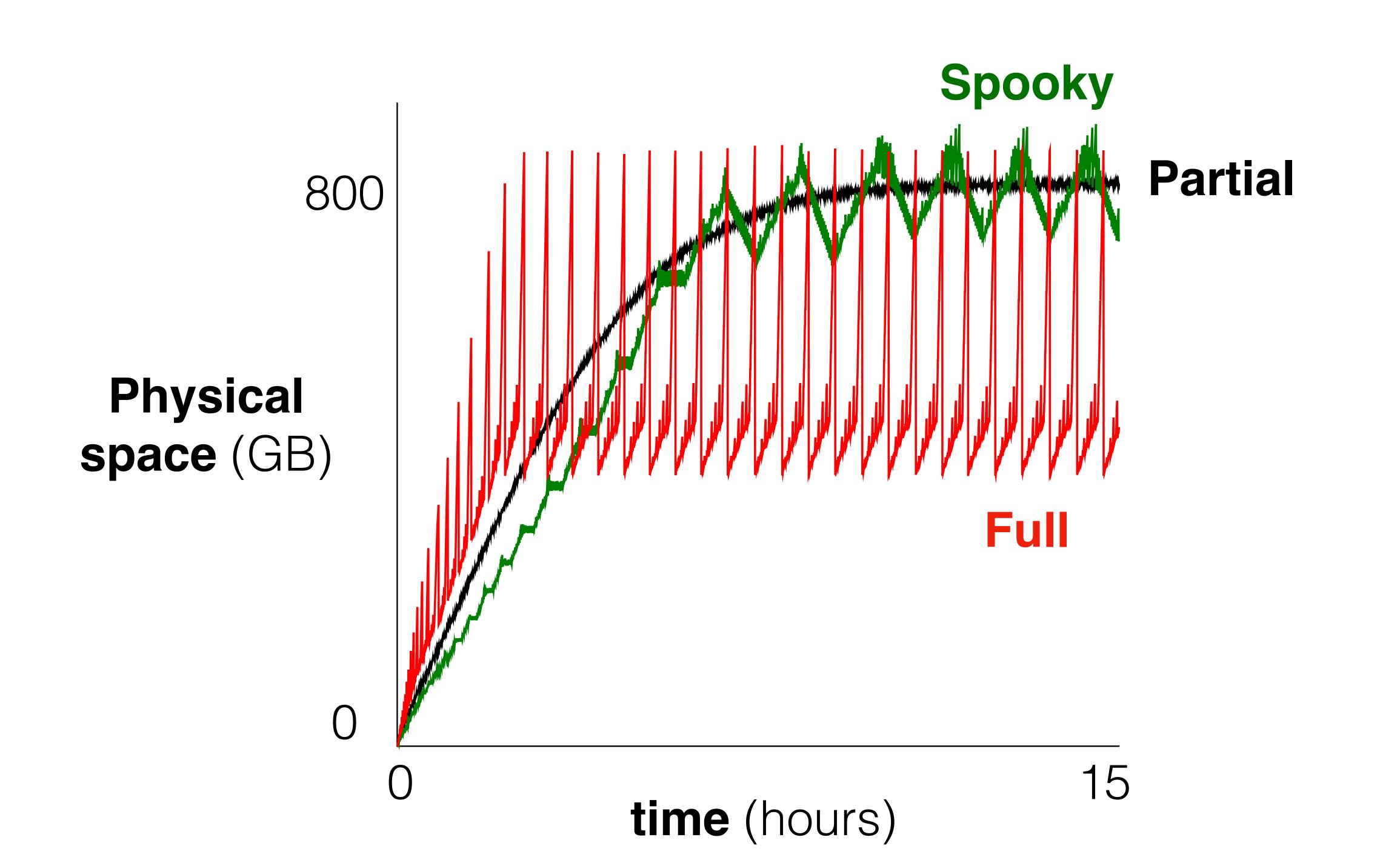






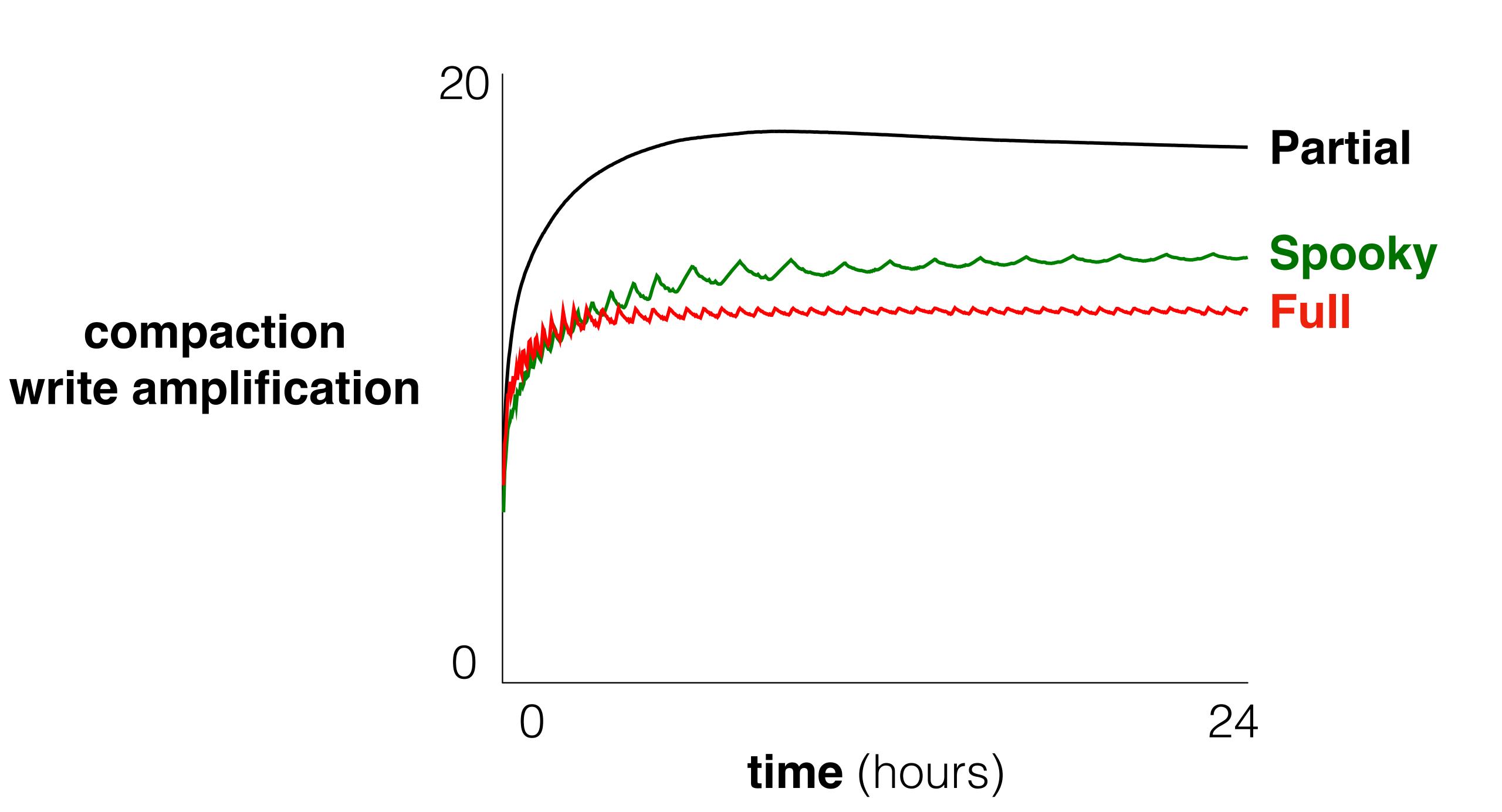


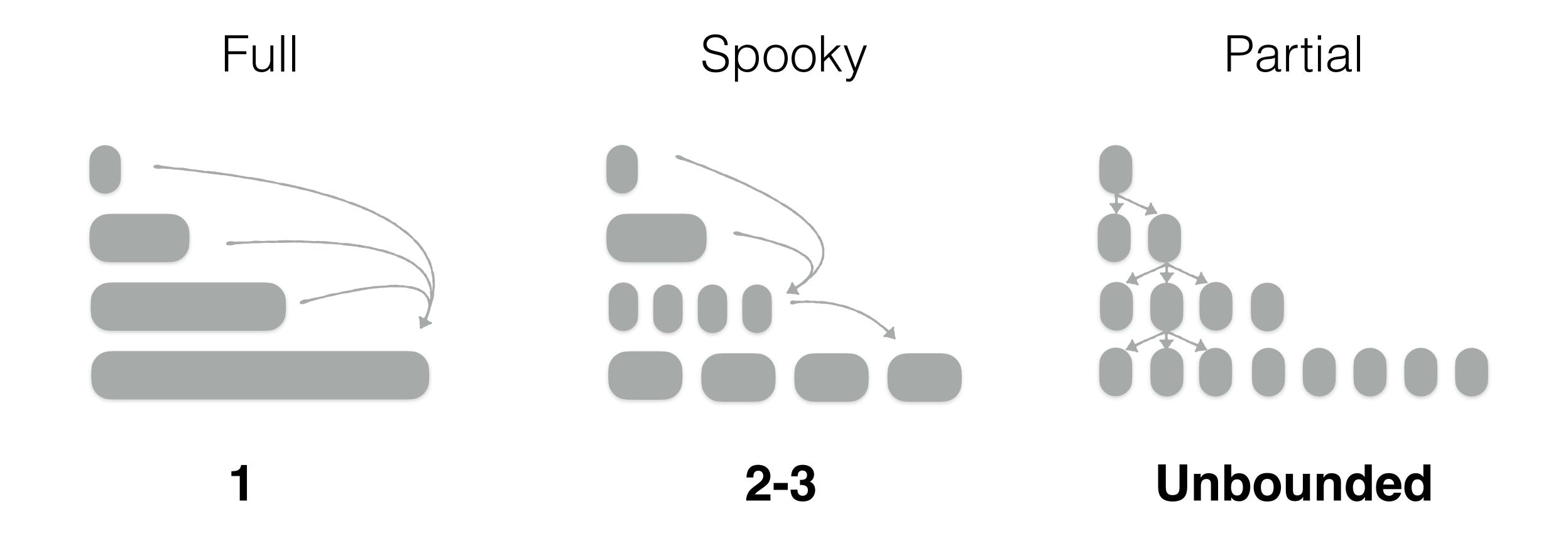
space-amplification



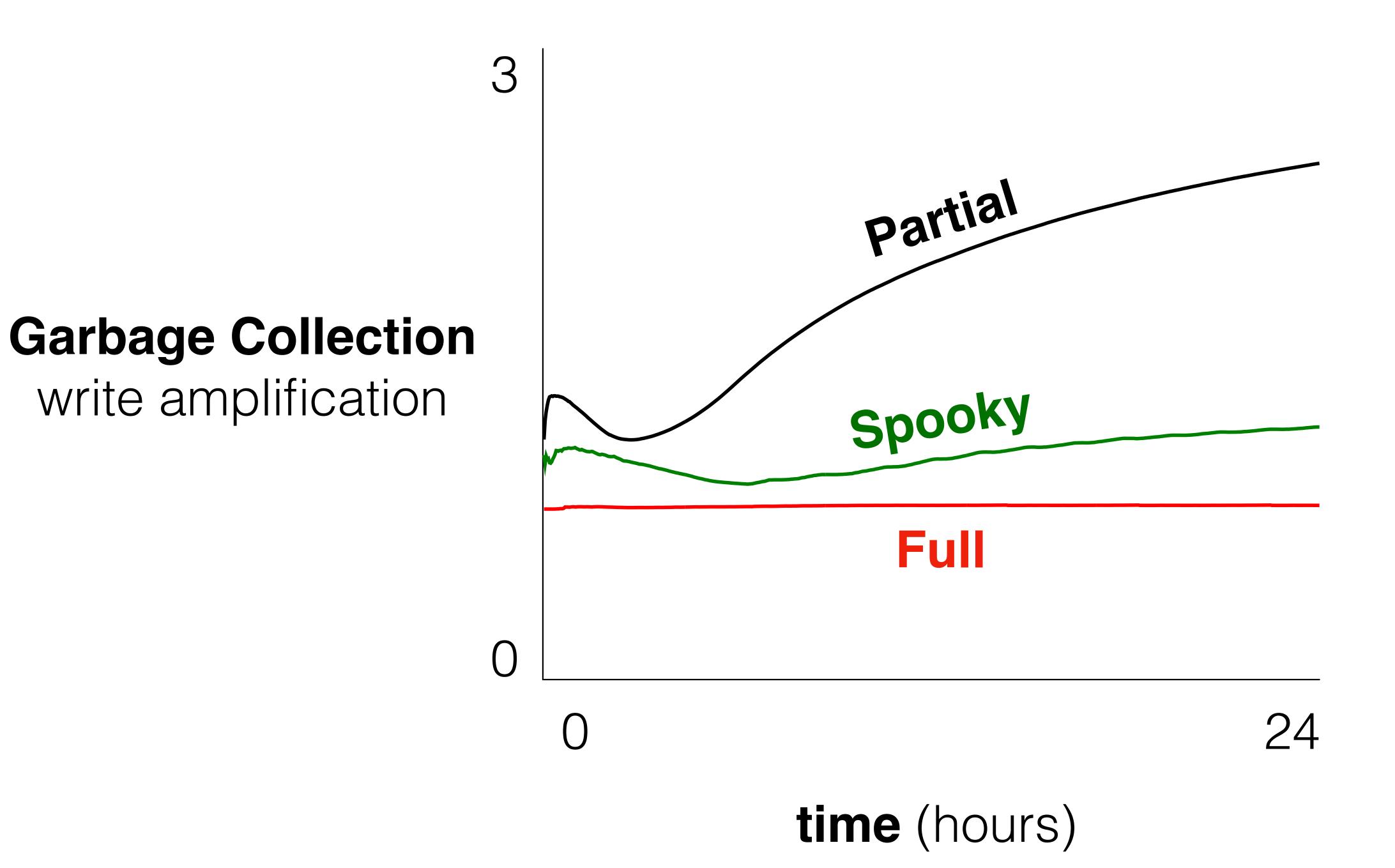
Full Spooky Partial Highest Modest Lowest

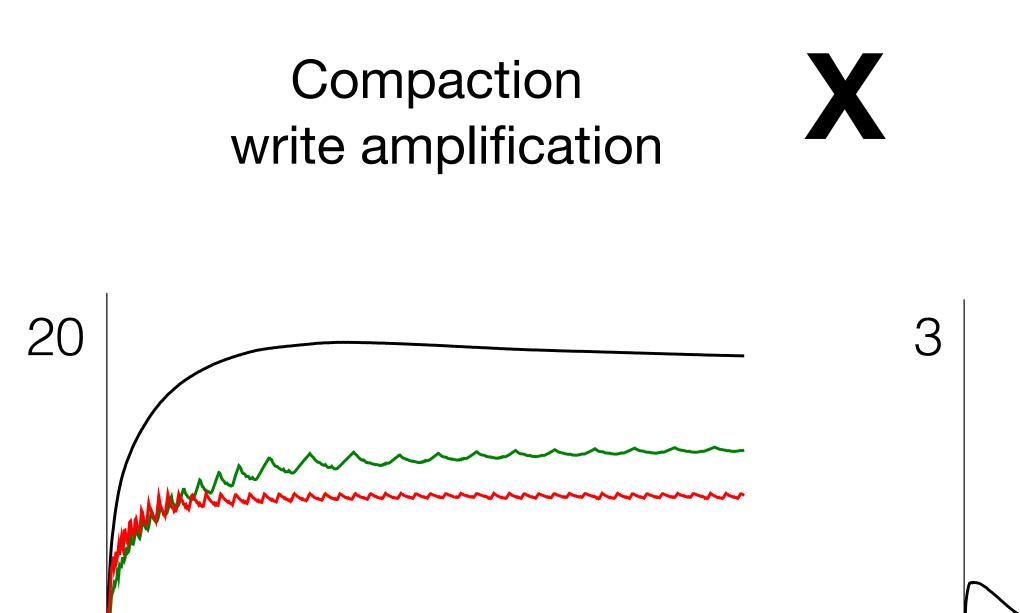
compaction write amplification





Simultaneous compactions

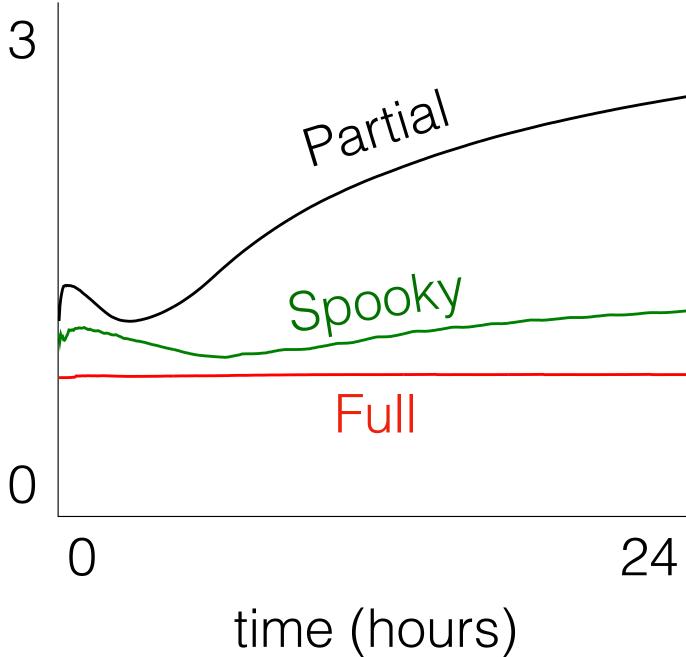


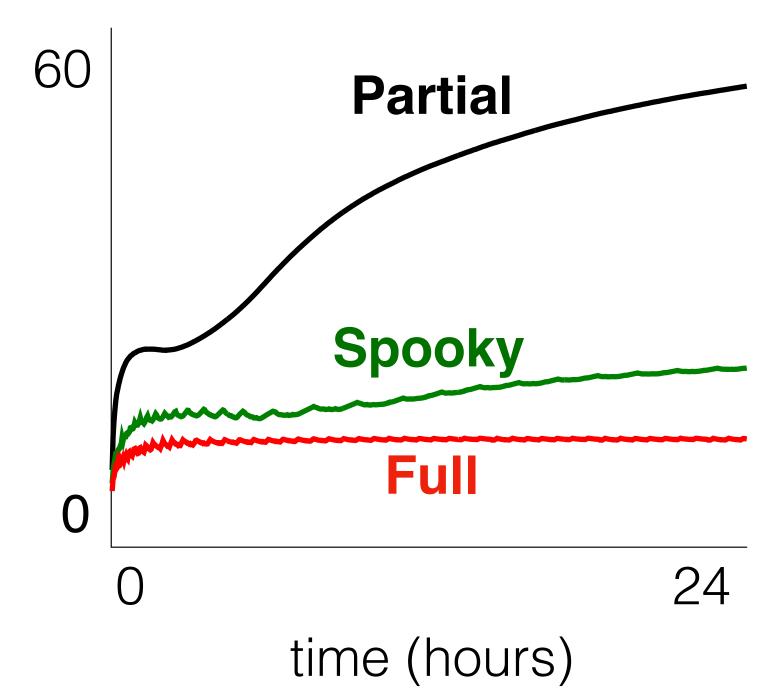


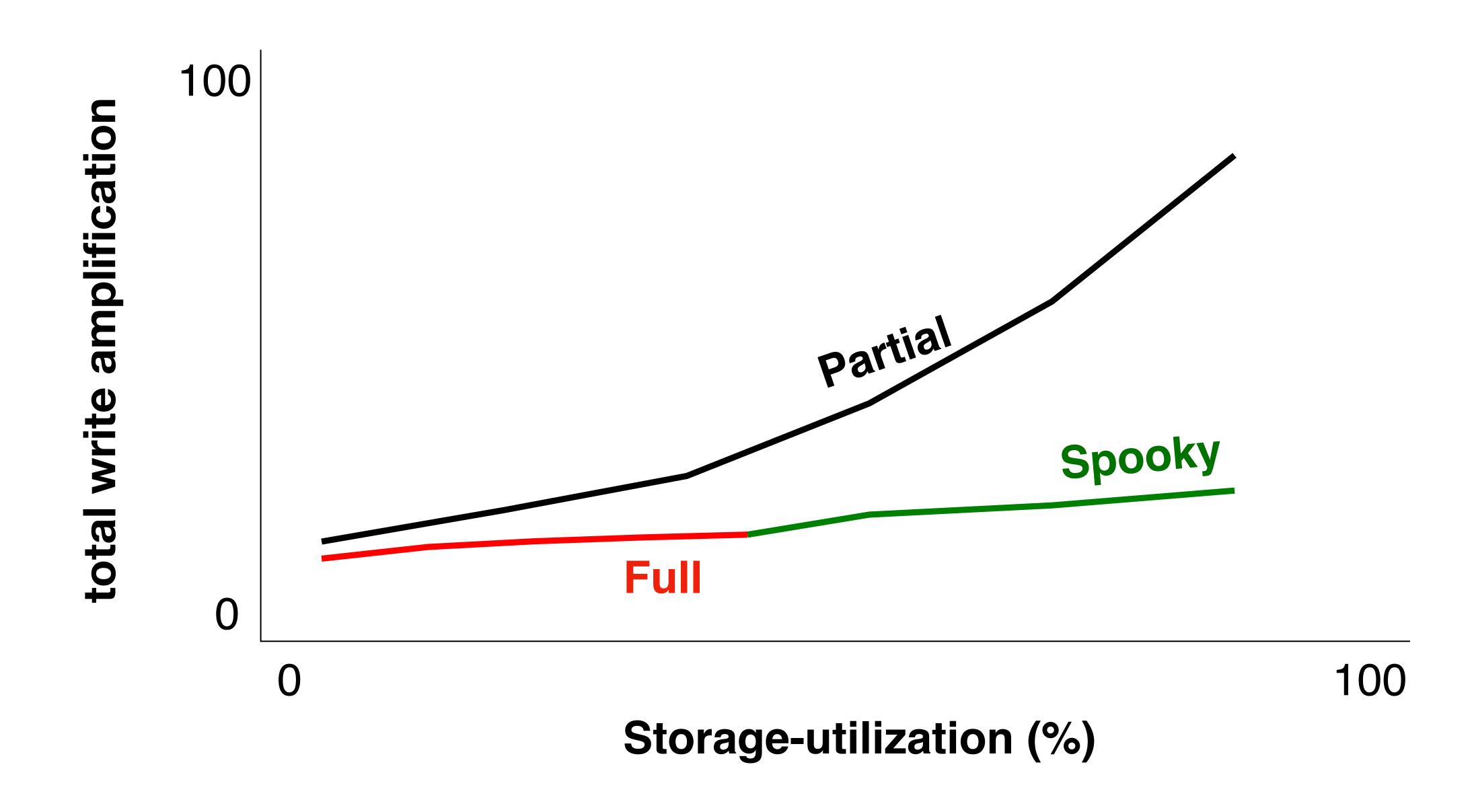
time (hours)

















space-amplification



