

# *CS 561: Data Systems Architecture*

## Class 6

# **Efficient Deletes in LSM-Engines**

BOSTON  
UNIVERSITY

*Dr. Subhadeep Sarkar*

<https://bu-disc.github.io/CS561/>



# Updates: Logistics

First **technical question** is due on **02/07**.

First **review** is due on **02/14**.

**Project 1 is now online!** Deadline: **02/20**.

Project 1 is a **group project** (2-3 students per group).

The first **student presentation** is next week (on **02/14**)!

**A week before the presentation**, discuss the slides with me in OH.

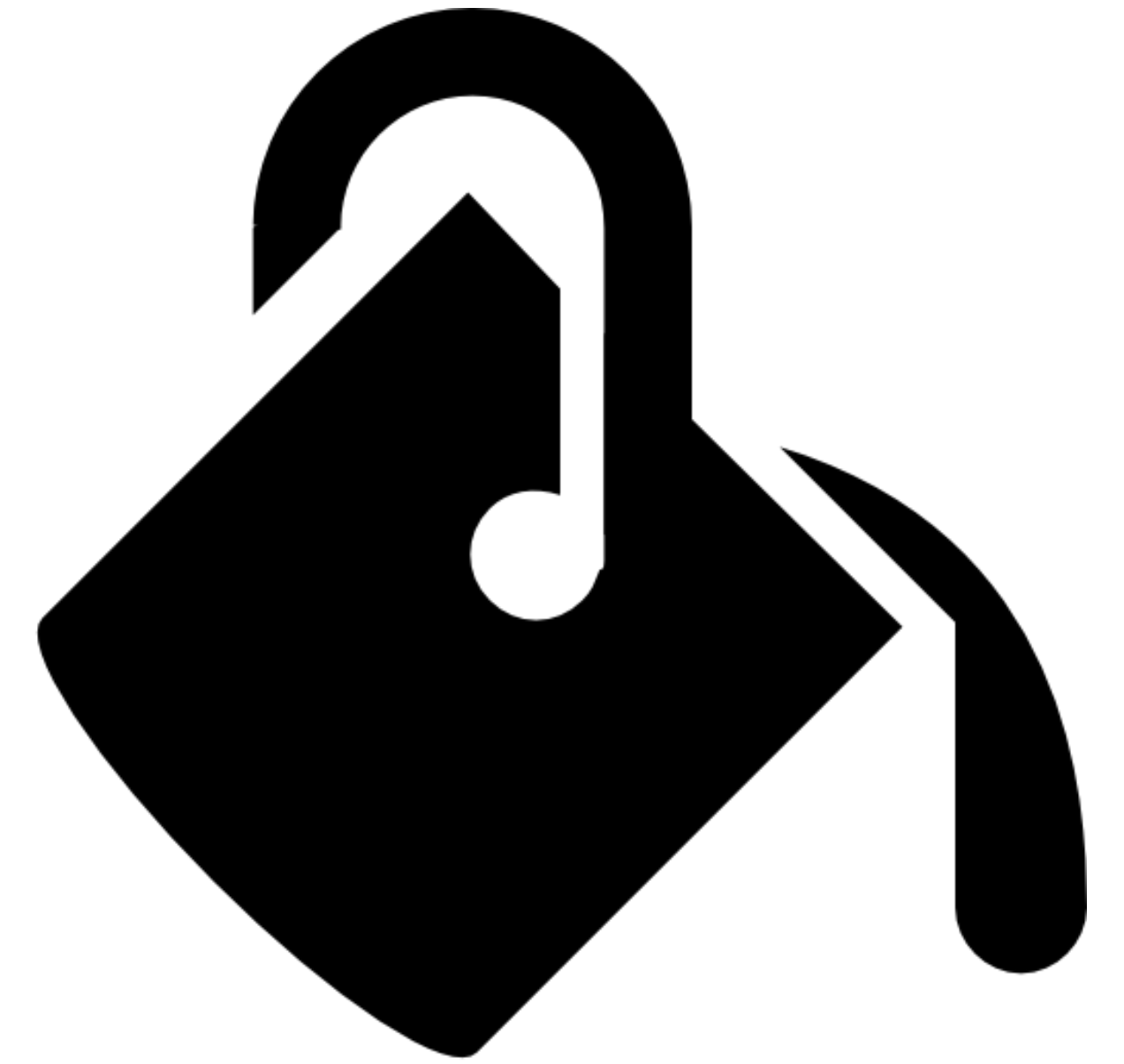
# How to prepare good slides



don't use  
bullets



1 message  
per slide



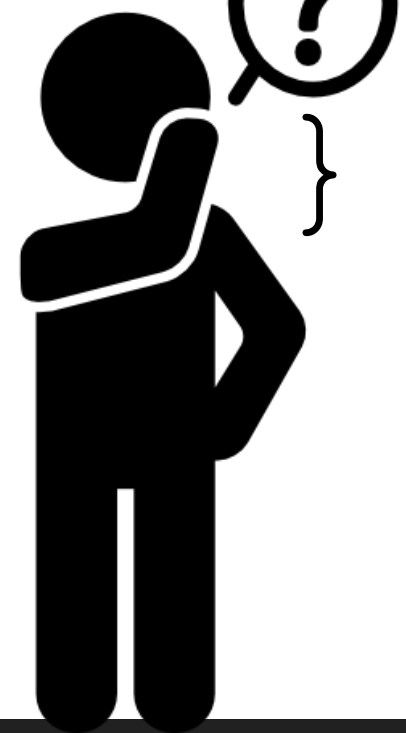
1/2 colors

```
<your_favorite_data_structure>::delete (key)
```

```
{
```

```
? //todo
```

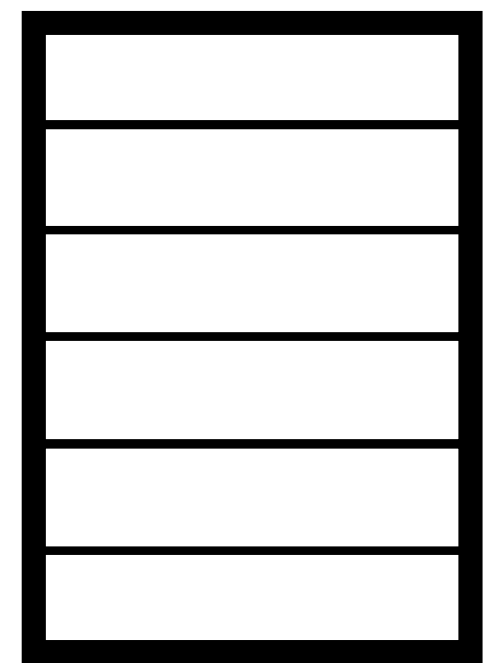
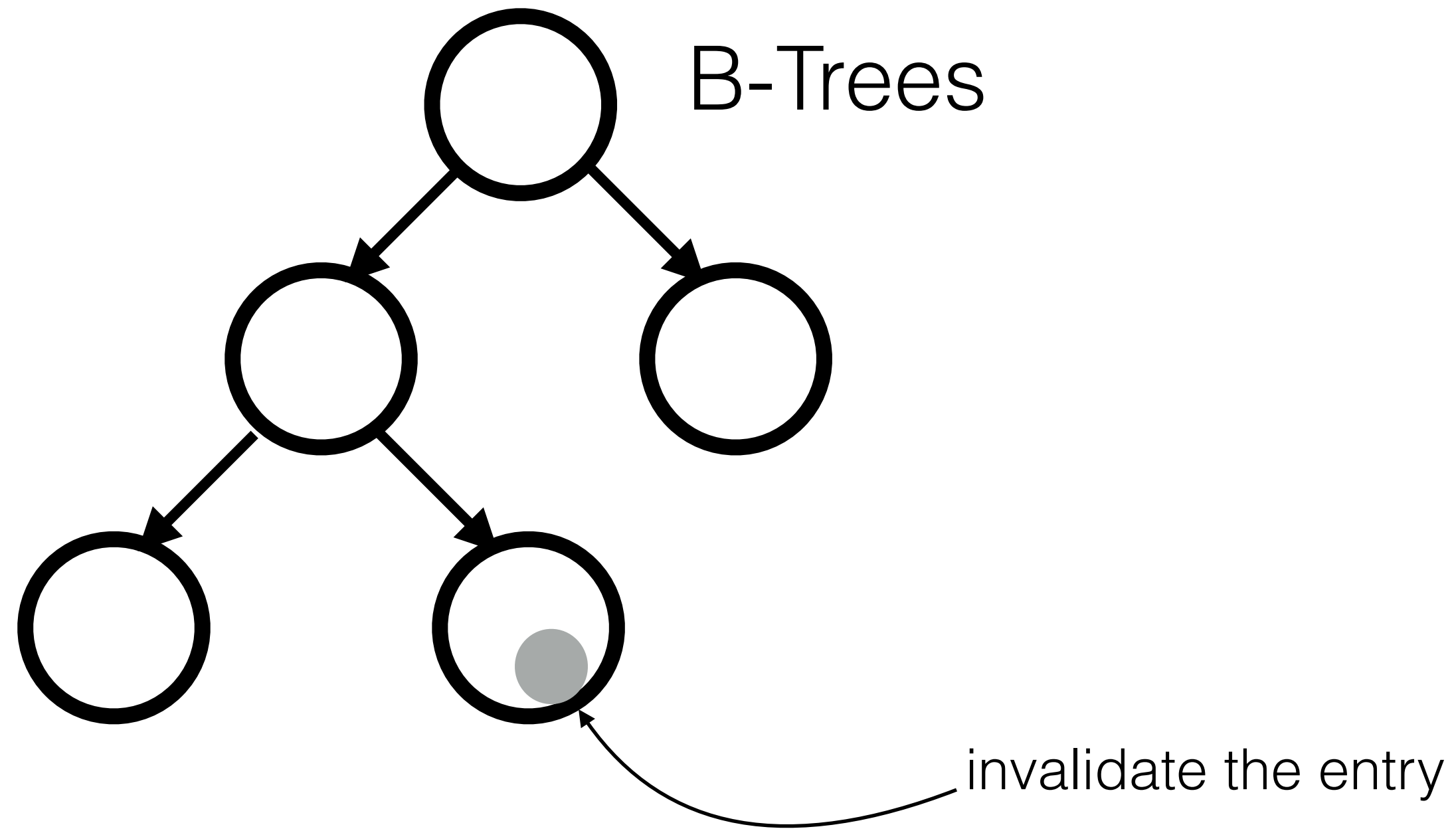
```
}
```



How do you delete data?

# IN-PLACE

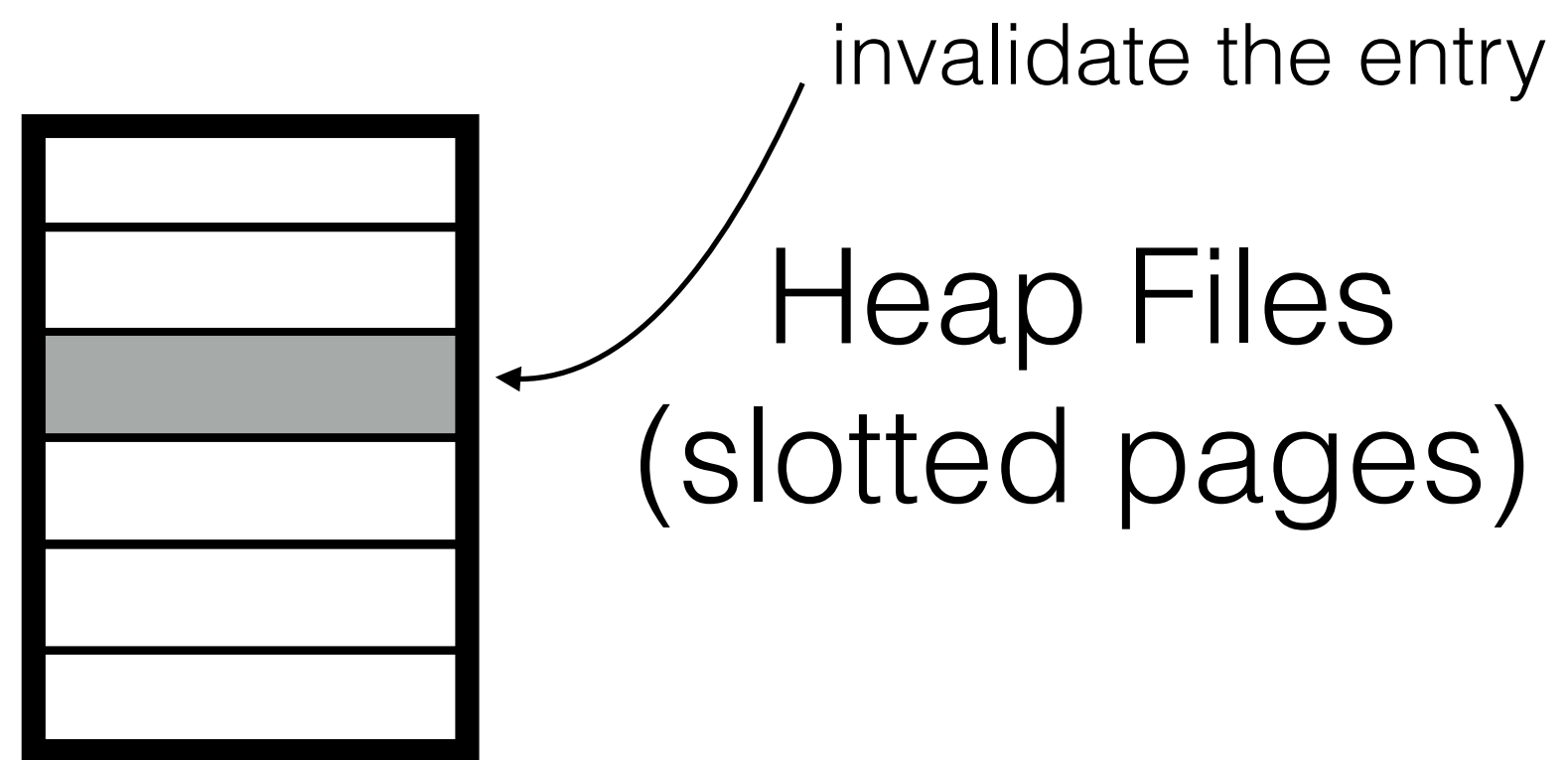
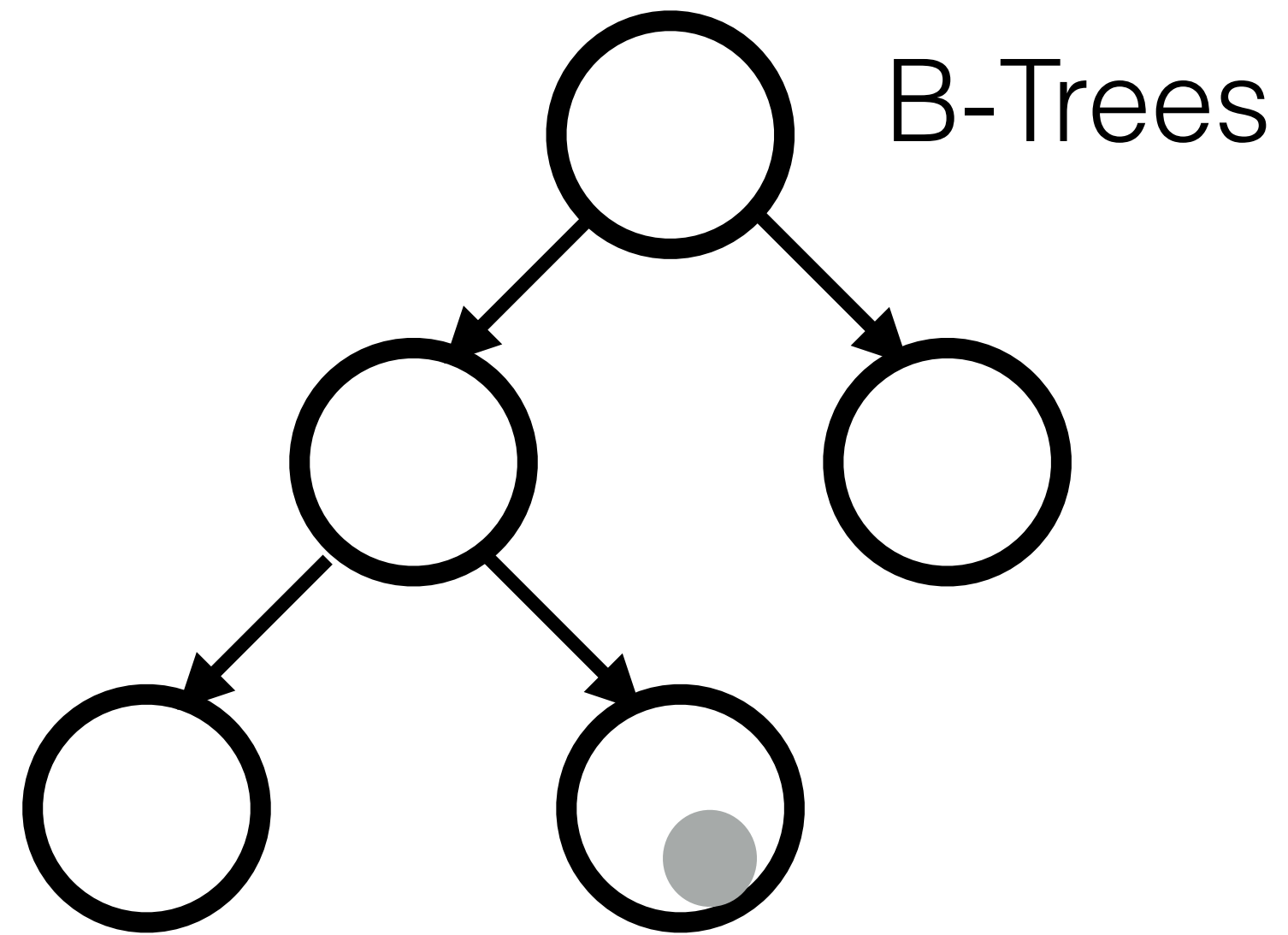
# OUT-OF-PLACE



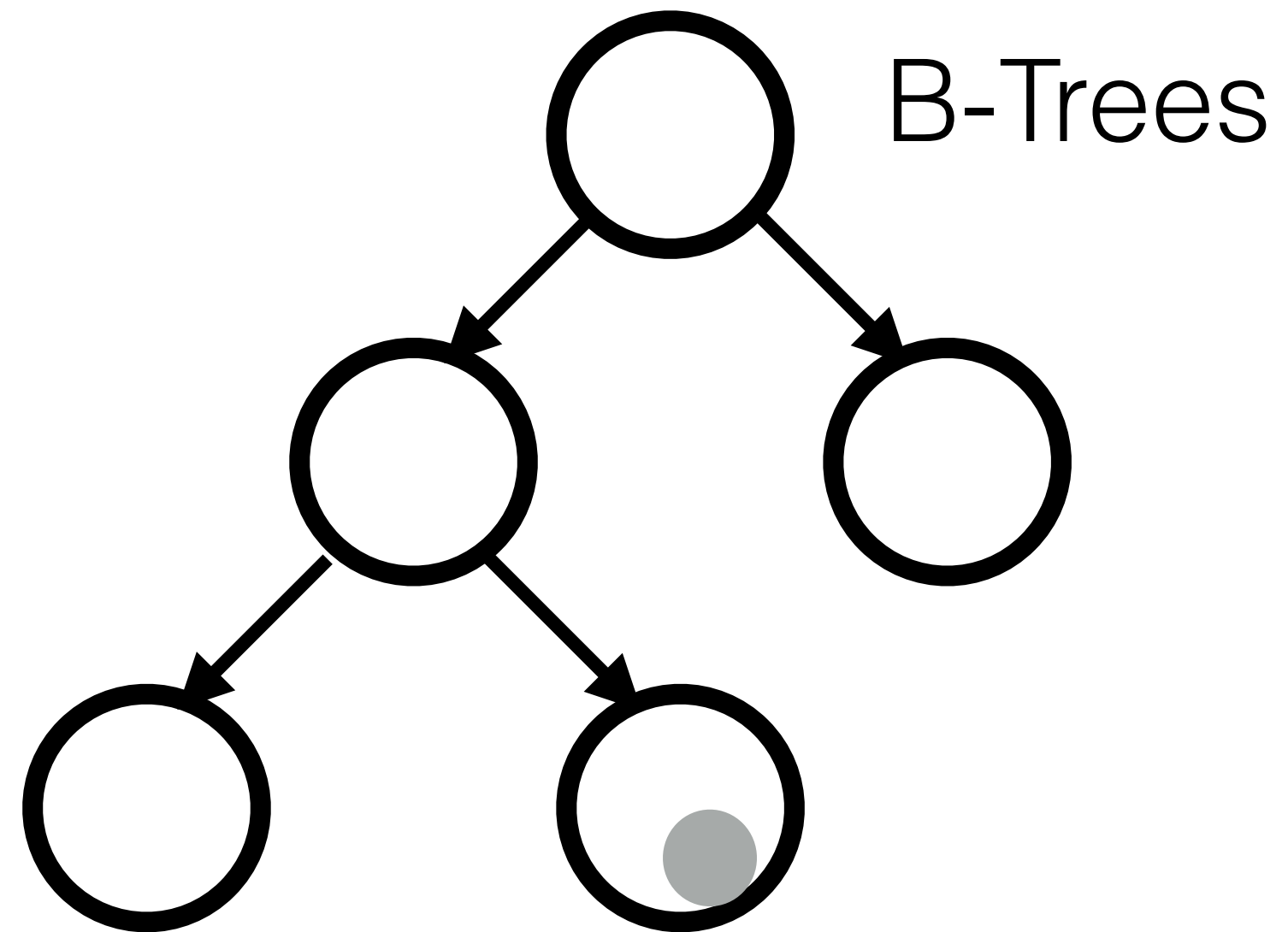
Heap Files  
(slotted pages)

# IN-PLACE

# OUT-OF-PLACE

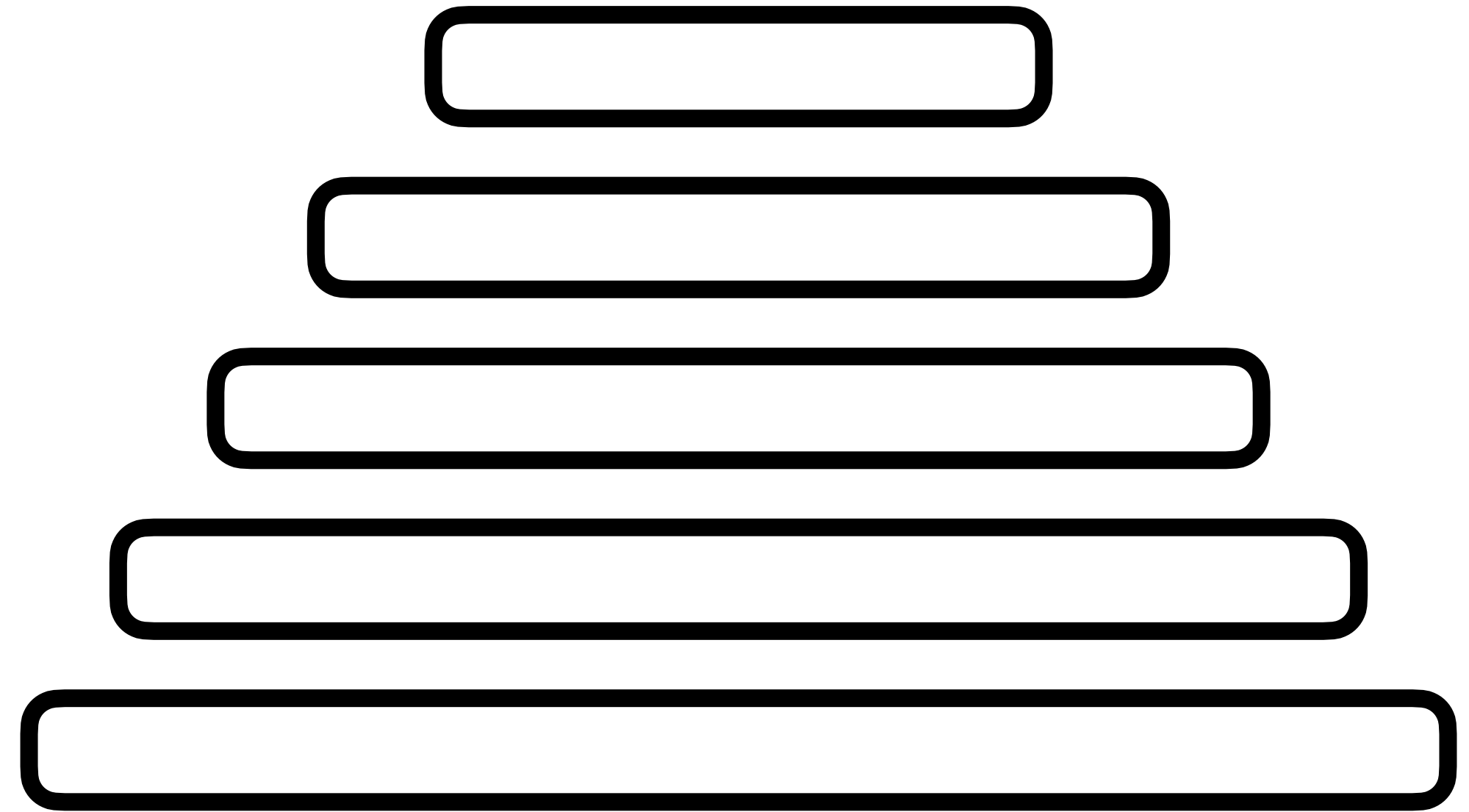


# IN-PLACE

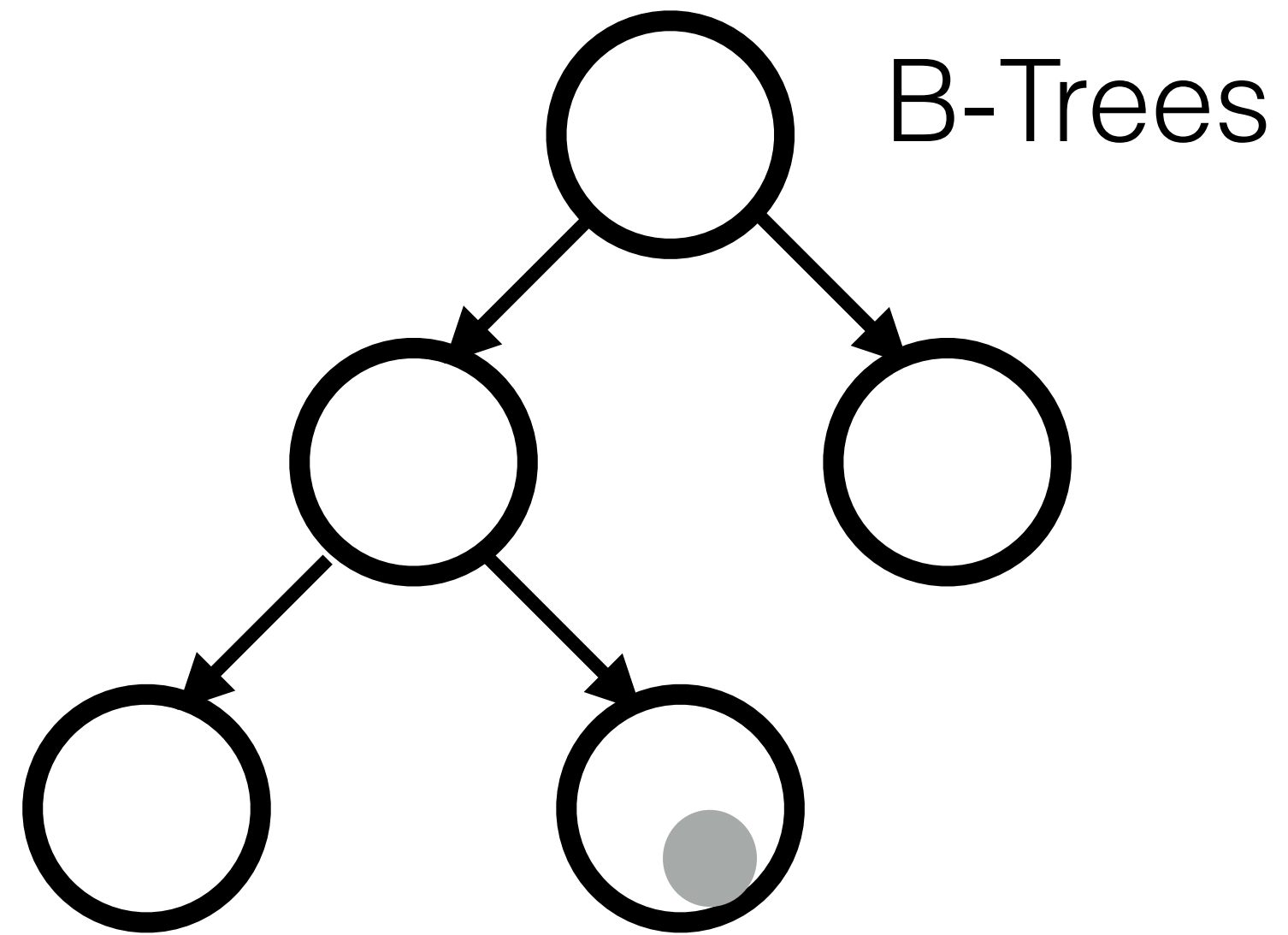


Heap Files  
(slotted pages)

# OUT-OF-PLACE

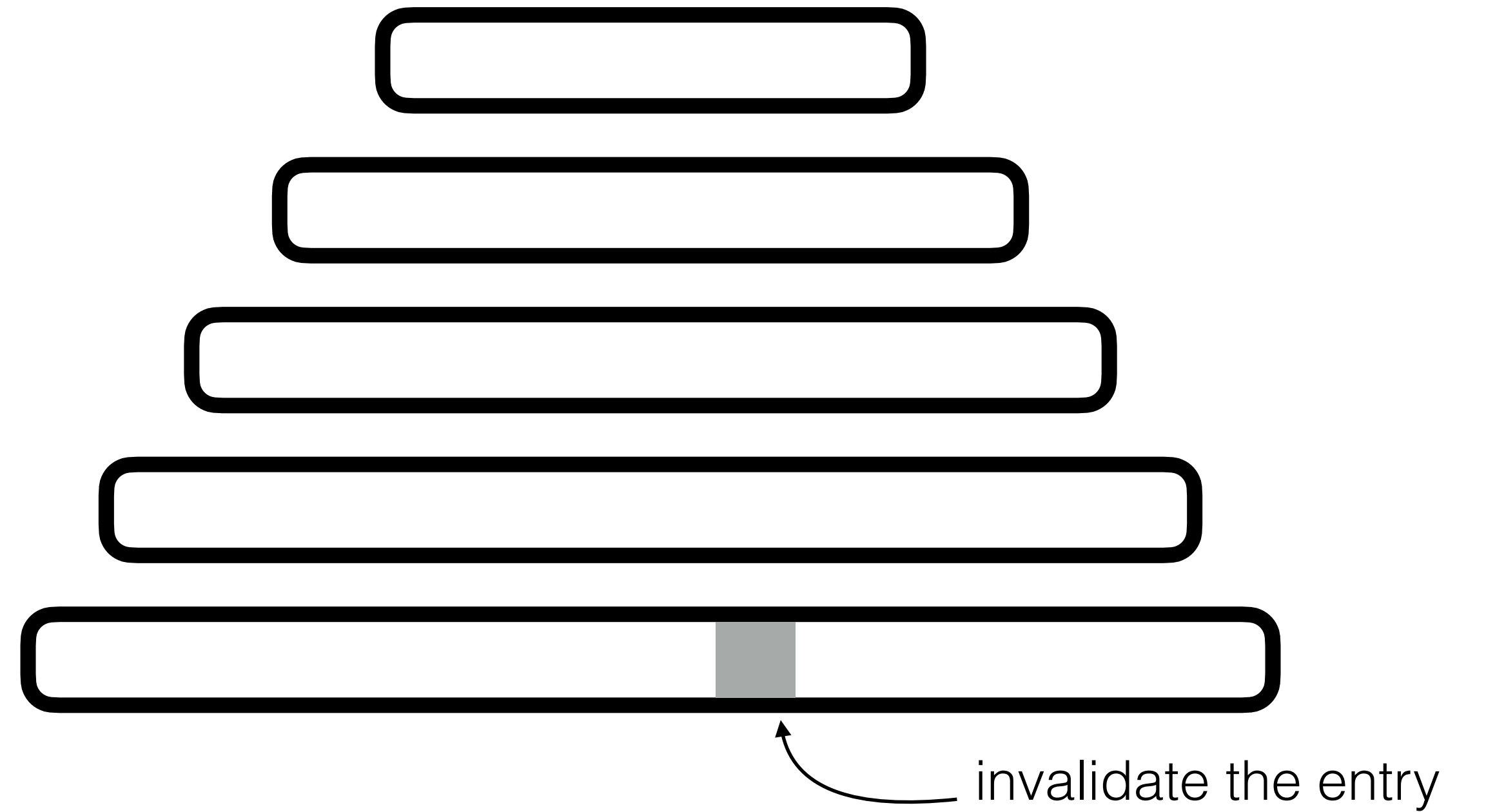


# IN-PLACE



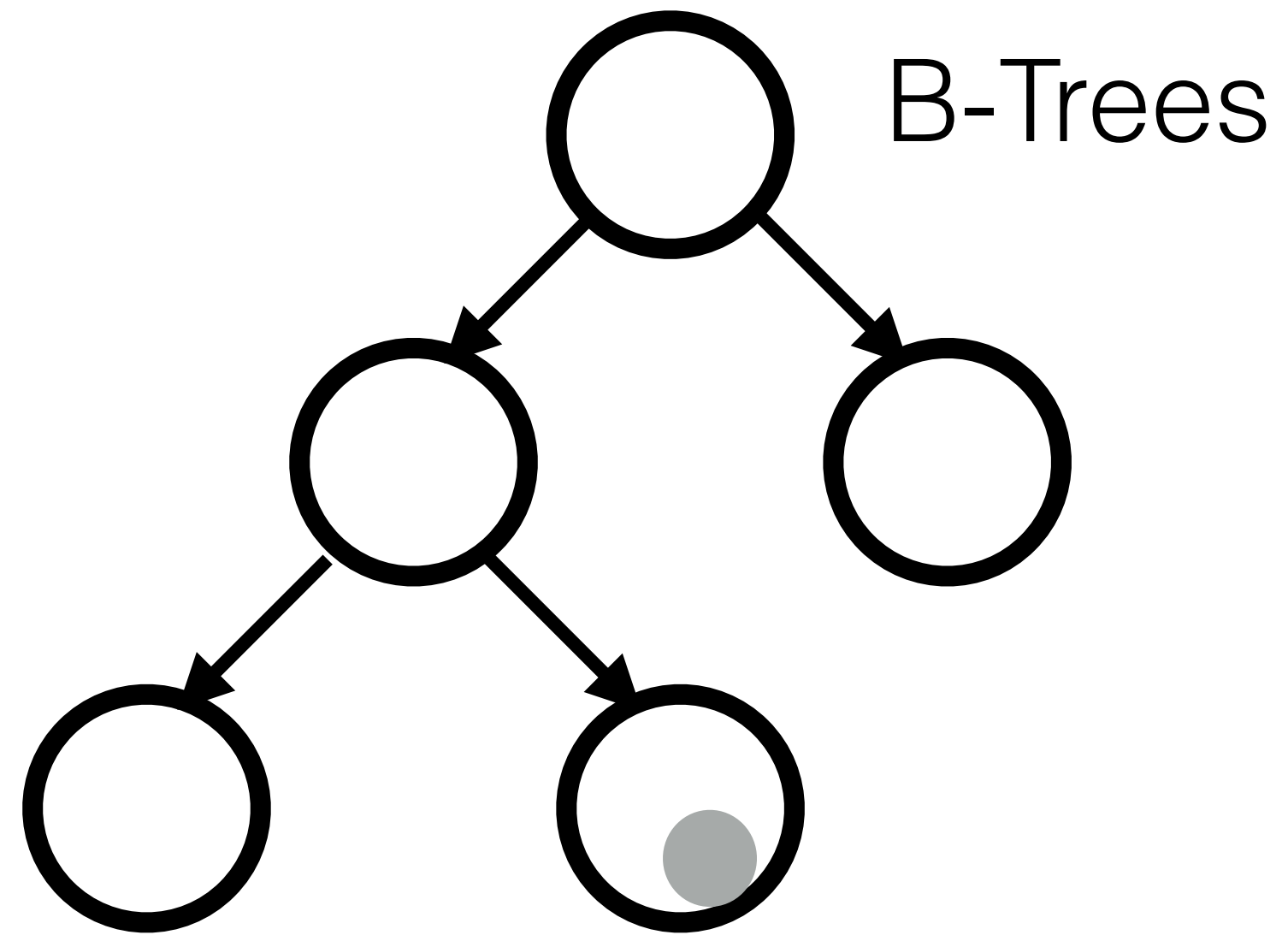
Heap Files  
(slotted pages)

# OUT-OF-PLACE



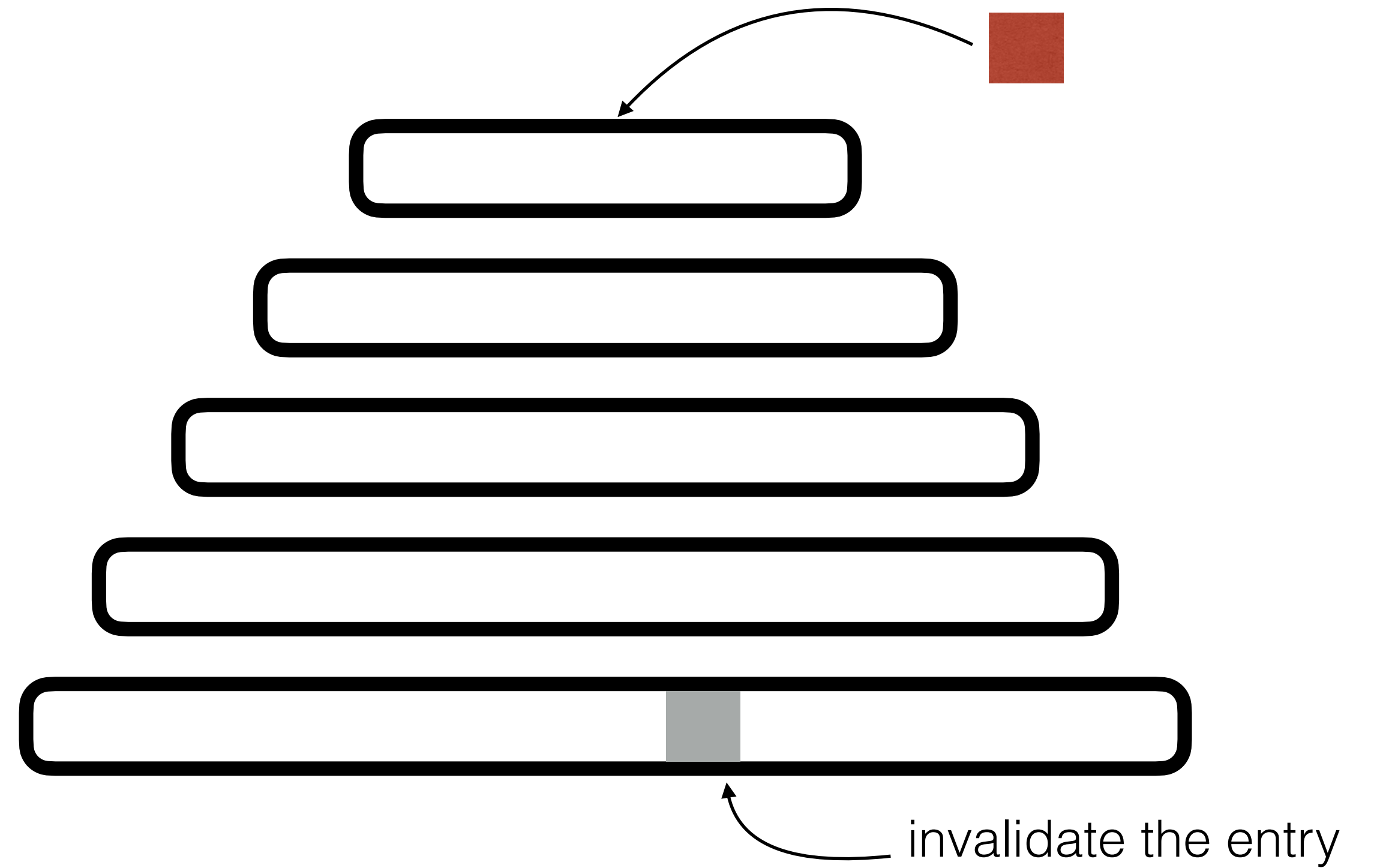


# IN-PLACE

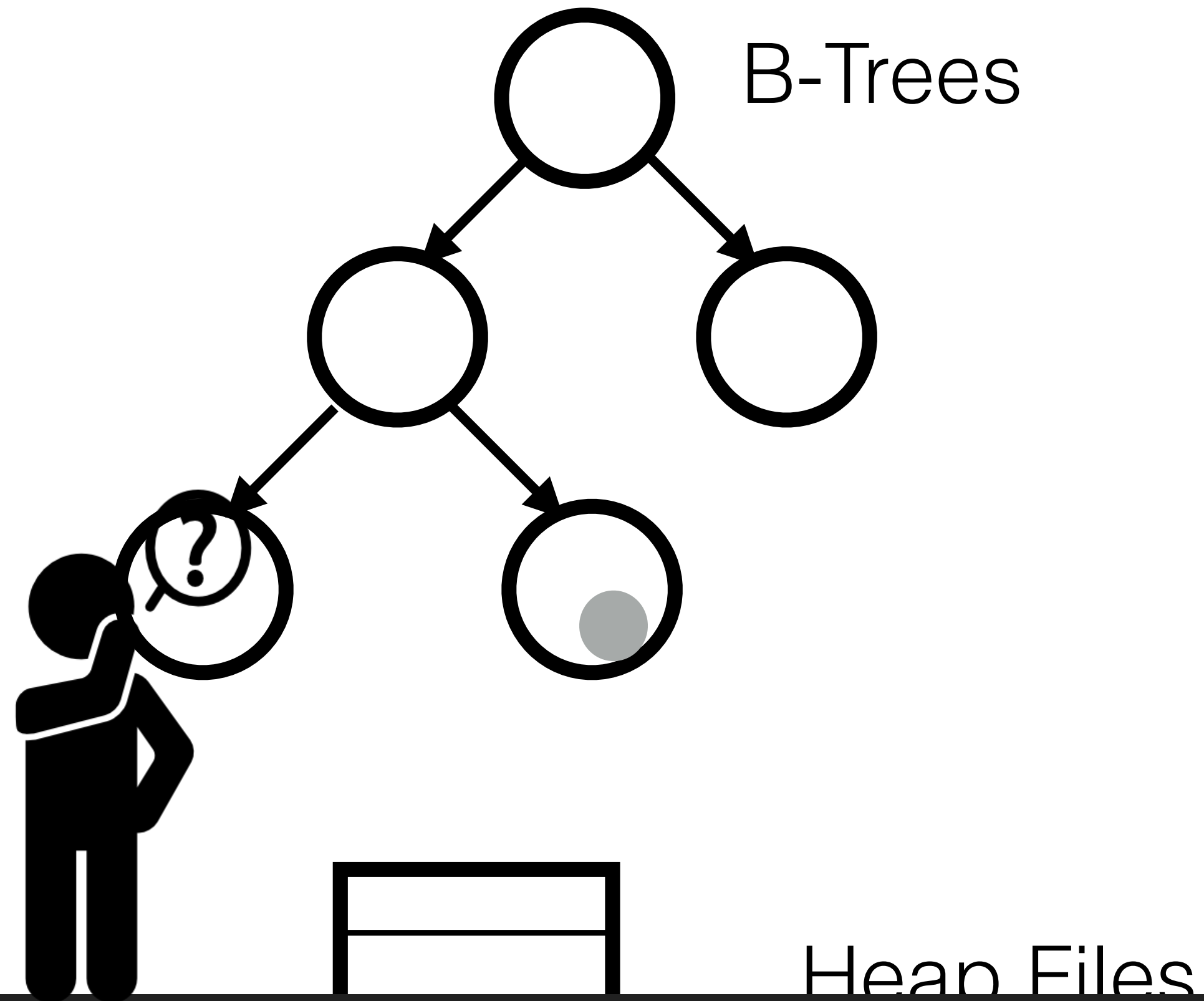


Heap Files  
(slotted pages)

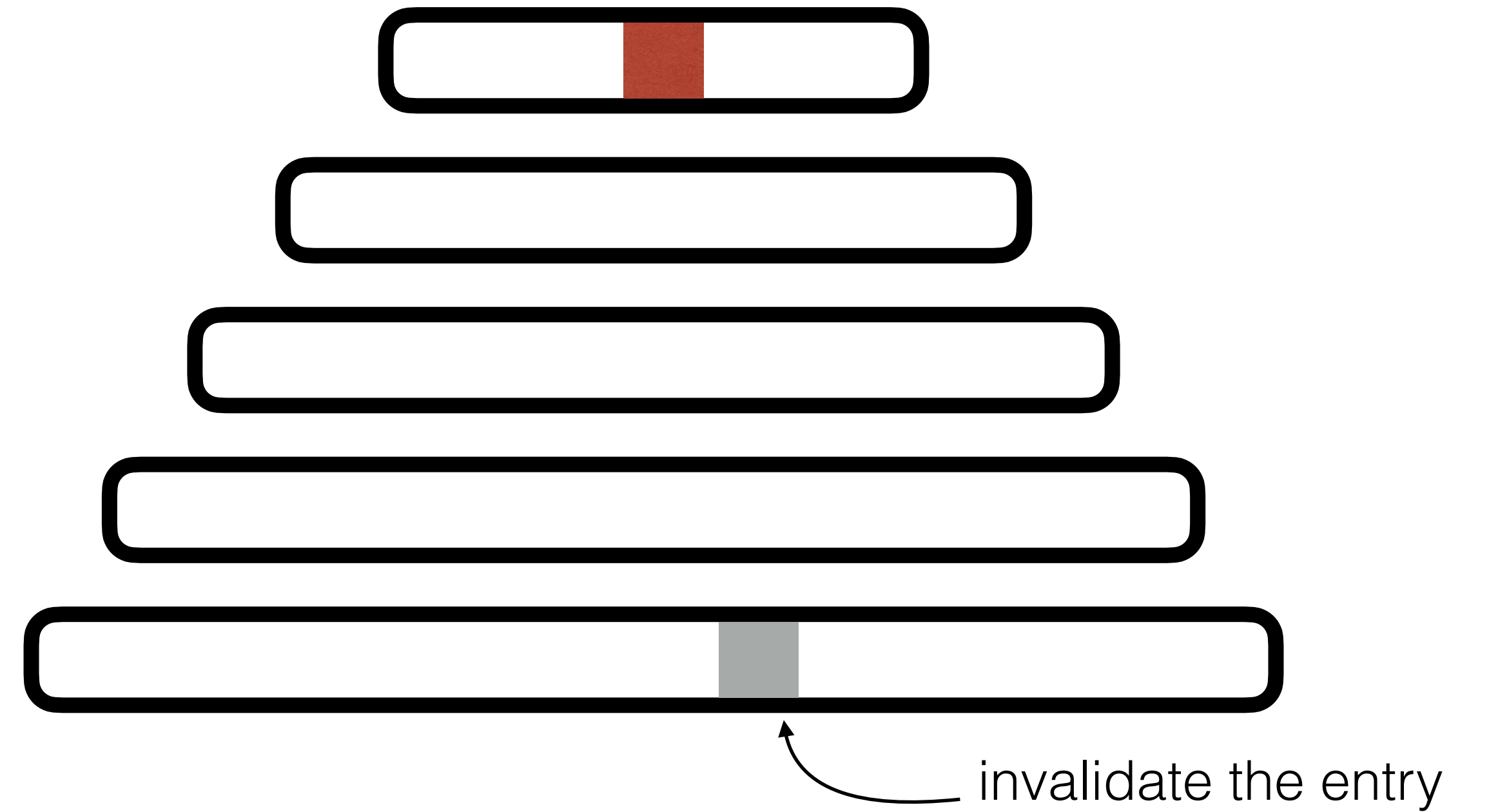
# OUT-OF-PLACE



# IN-PLACE

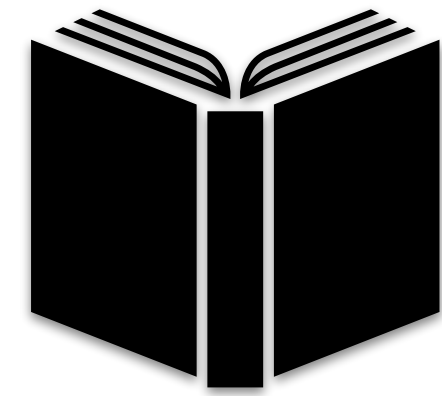


# OUT-OF-PLACE



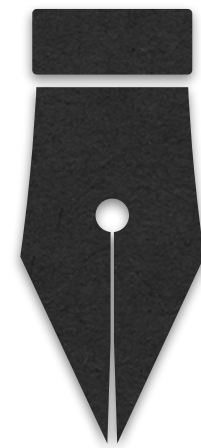
What is the tradeoff for deletes?

# What is the delete tradeoff?



read

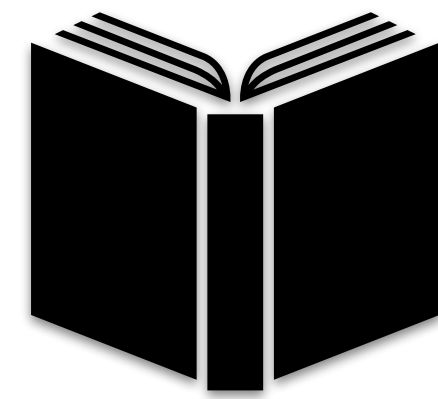
vs.



write

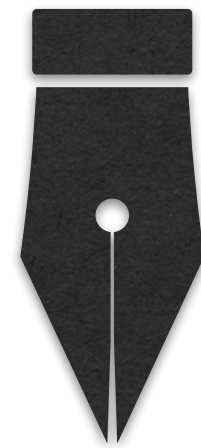


# What is the delete tradeoff?



read

vs.



write



Deletes are almost **exclusively** *logical*

**Today's talk:**

# Lethe: A Tunable Delete-Aware LSM-Based Storage Engine

Presented at SIGMOD 2020

# LSM-tree

NoSQL

This block contains logos for various NoSQL databases that utilize the LSM-tree architecture. The logos are arranged in two rows. The top row includes RocksDB (a yellow cheetah), WT (black and orange letters), levelDB (a green cylinder), SCYLLA (a blue alien head), and DynamoDB (a blue cylinder). The bottom row includes cassandra (an eye), tarantool (two red circles), Bigtable (a blue hexagon), APACHE HBASE (a black orca), riak (a grey starburst), and accumulo (a grid of squares).

This block contains two logos. On the left is SQLite, featuring a blue square with a white feather and the text 'SQLite'. On the right is a dark blue square containing a white silhouette of a dolphin.

relational

This block contains two logos. On the left is influxdb, featuring a blue cube and the text 'influxdb'. On the right is QuasarDB, featuring a blue grid pattern.

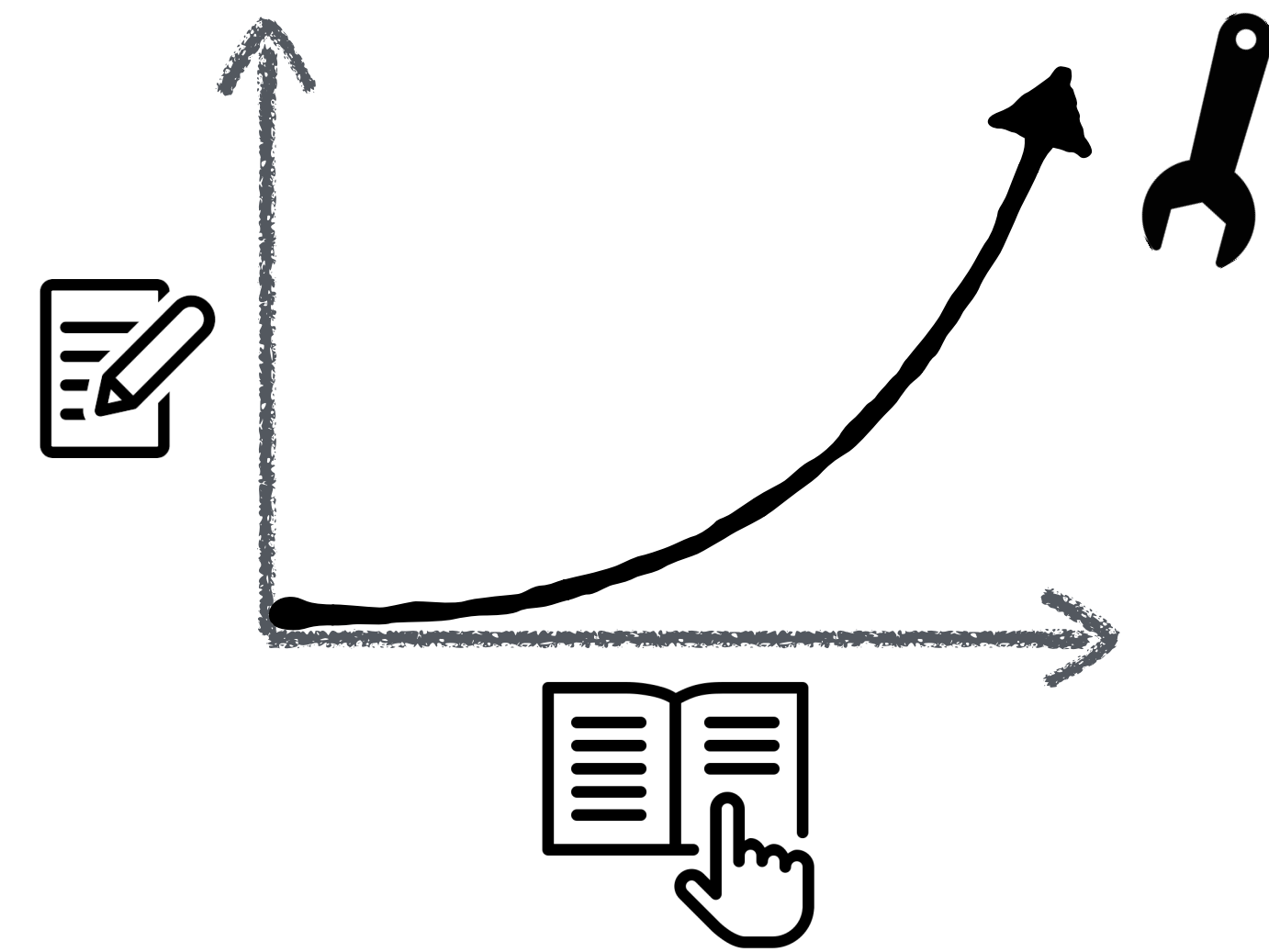
time-series

2023

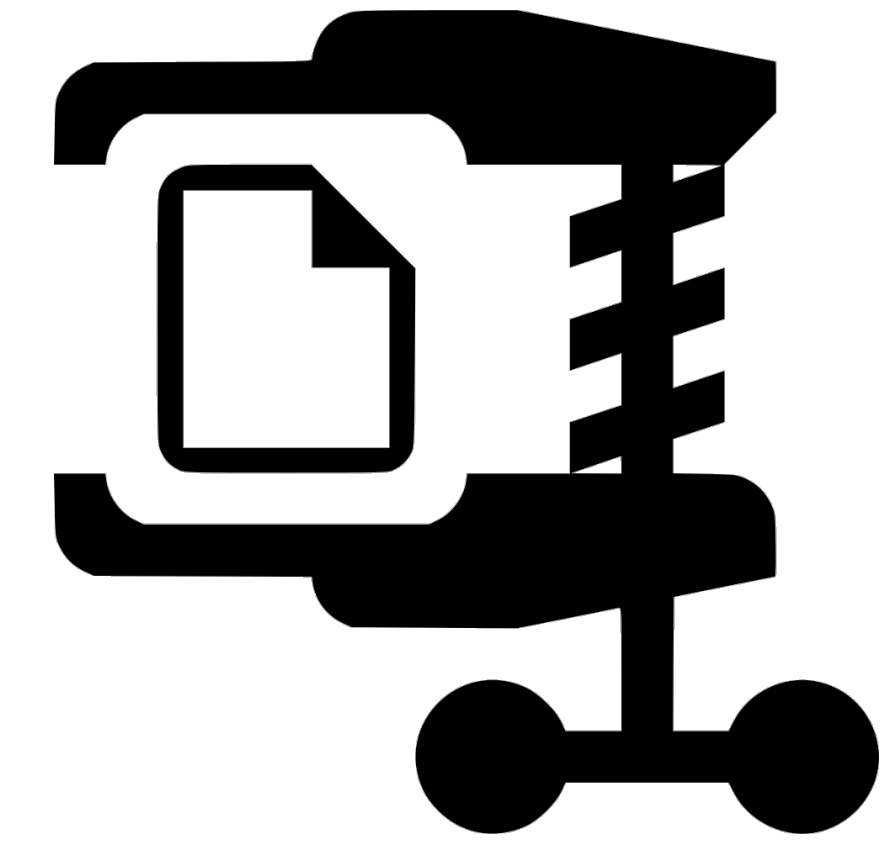
# Why **LSM** ?



fast writes



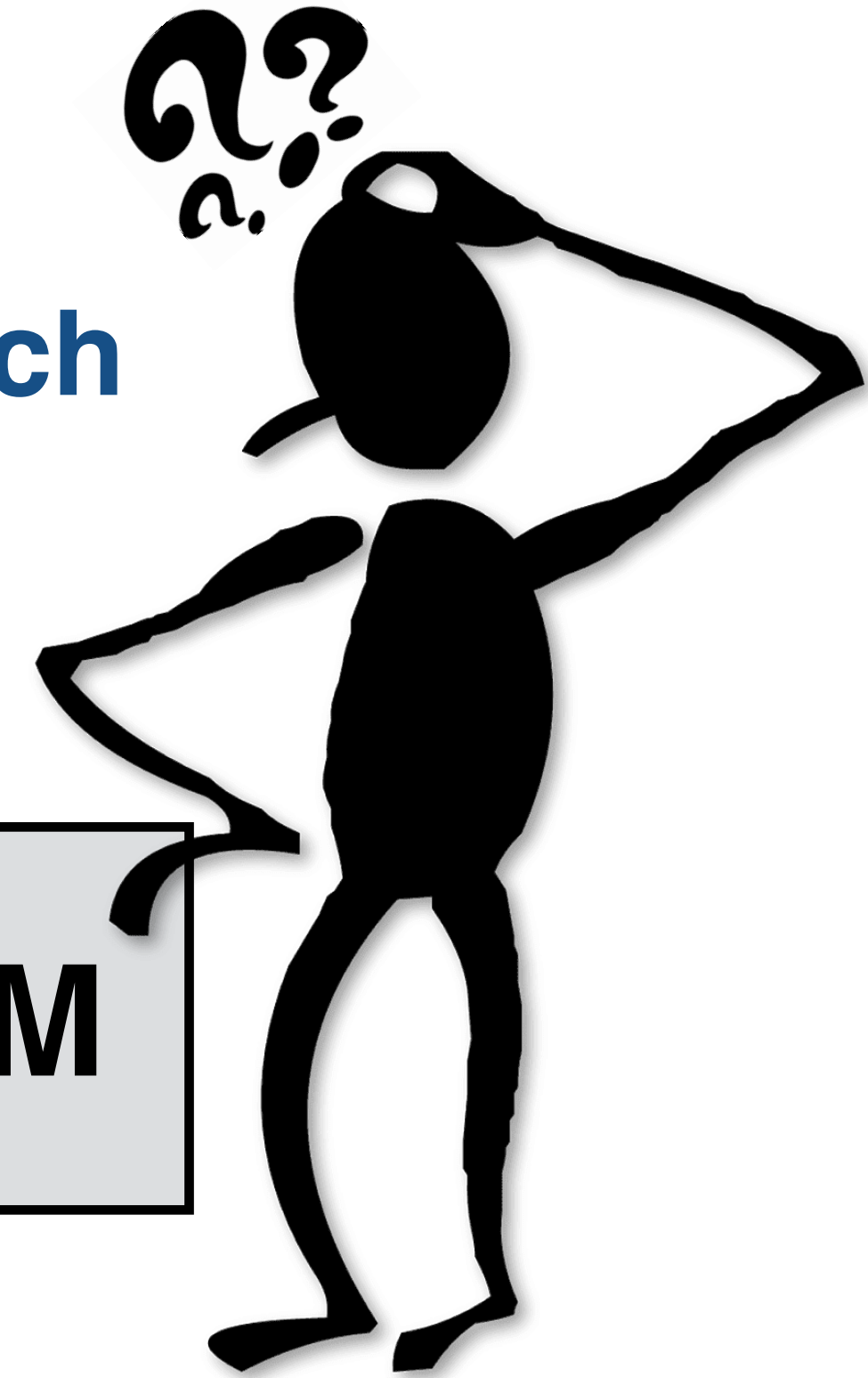
tunable read-write performance



good space utilization

**Even years later, Twitter doesn't delete your direct messages**

**TechCrunch**  
Feb '19



**Small Datum**  
Jan '20

**Deletes are fast and slow in an LSM**

“LSM-based data stores perform suboptimally for workloads with deletes.”

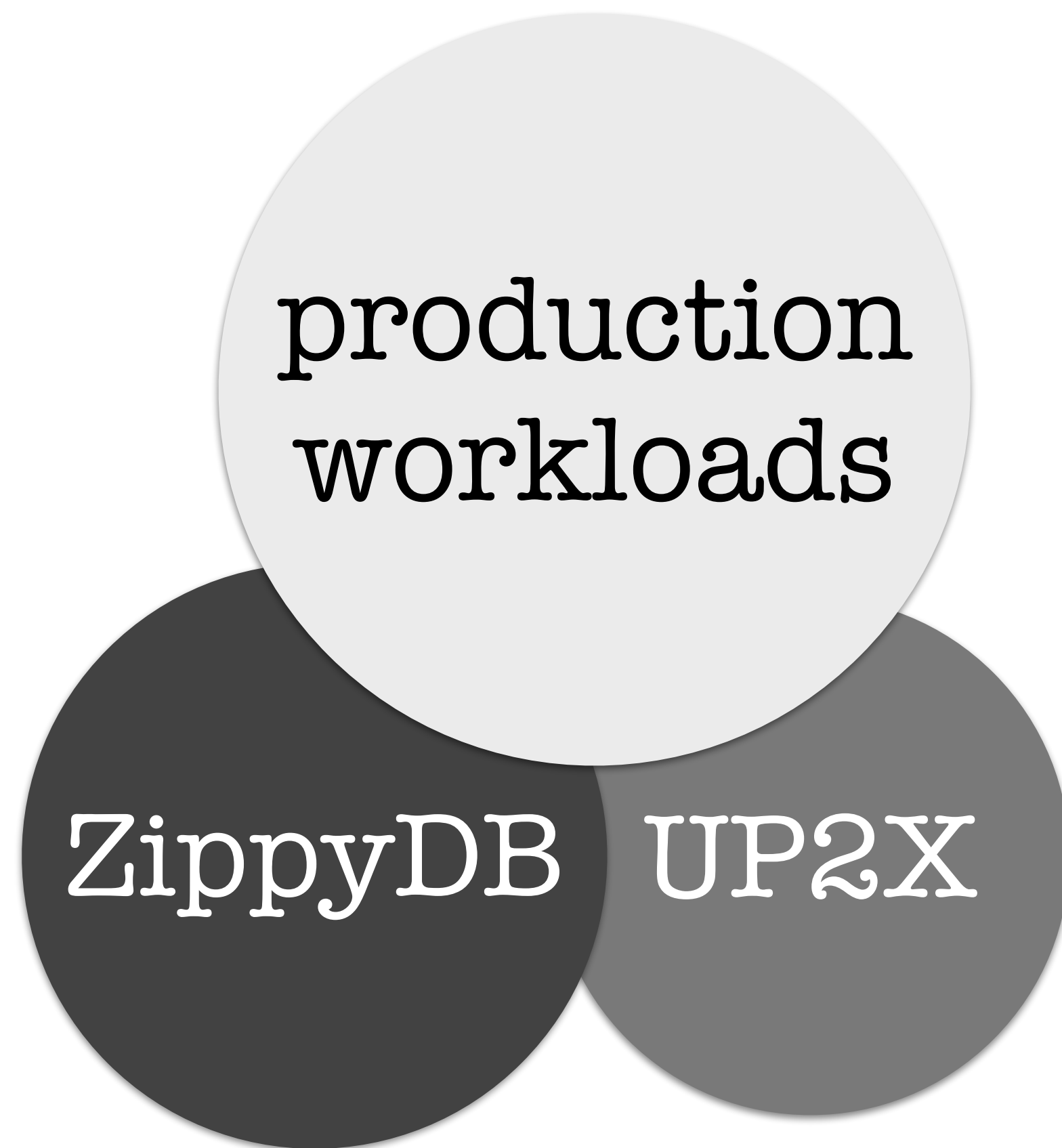


Some items will also be deleted from 11 albums.

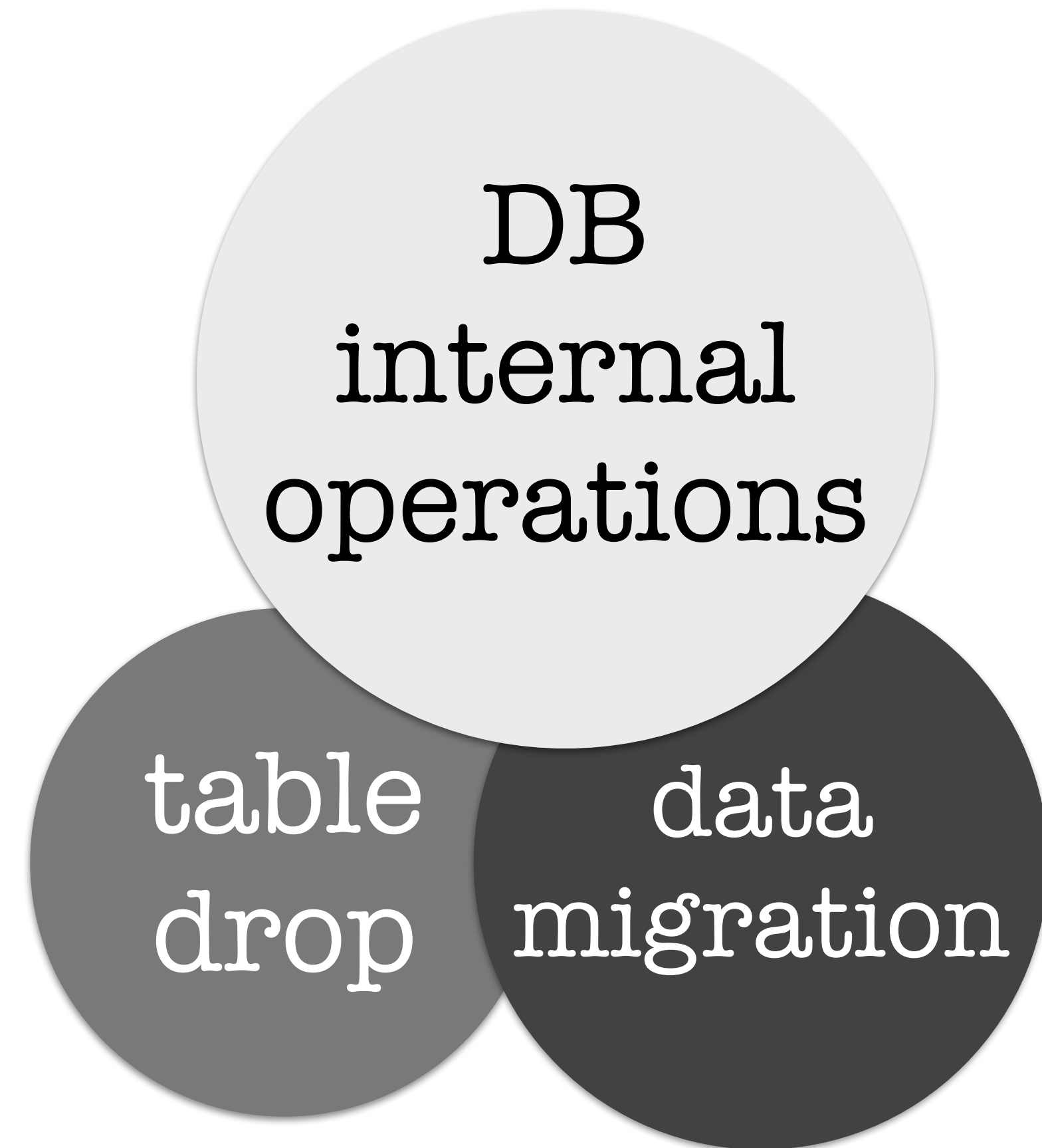
Delete 6,447 Items

Cancel

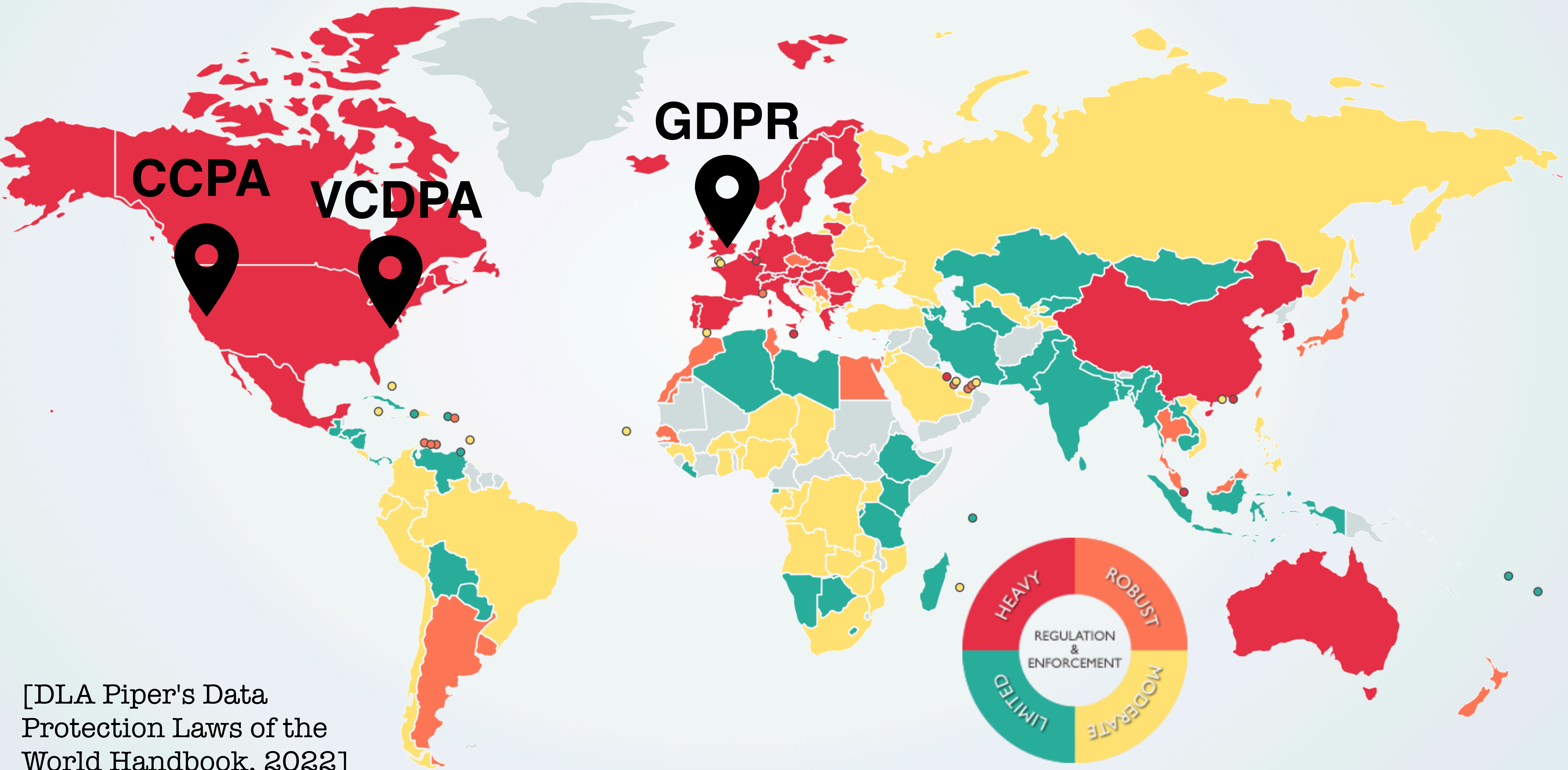
**100M+** deletes/day



deletes in  
**batches**



# Logical Deletes & Data Privacy



[DLA Piper's Data Protection Laws of the World Handbook, 2022]



GDPR  
(EU, UK)



CCPA  
(California)



VCPDA  
(Virginia)



on-demand



rolling

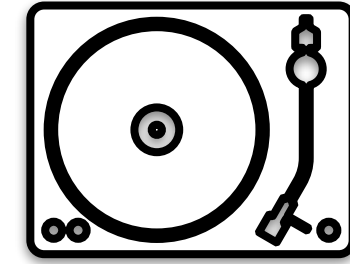
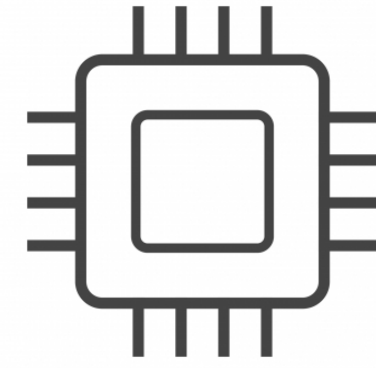
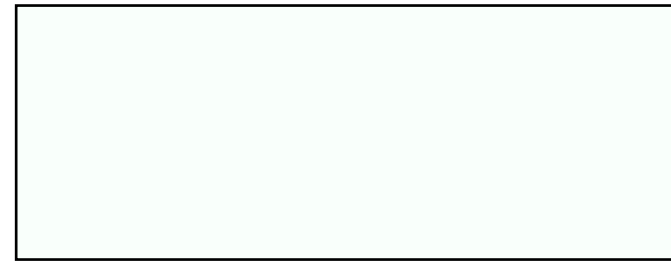
*delete all data for  
user X within D days*

*keep deleting all data  
older than D days*

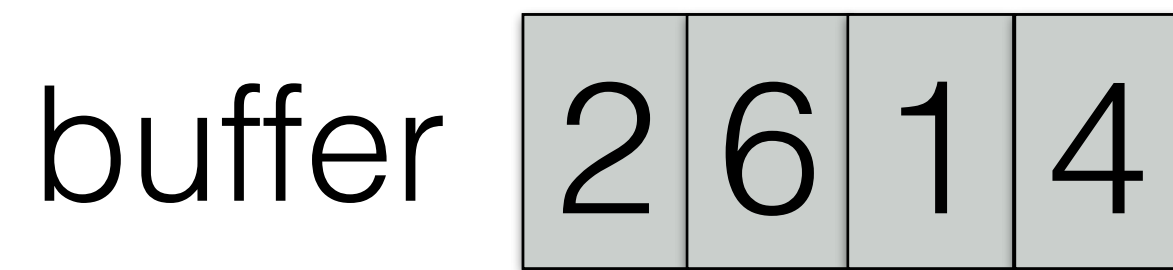
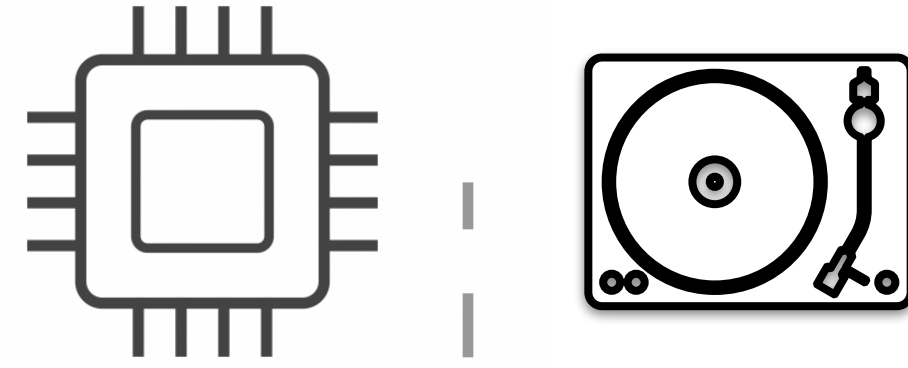
*A reminder on how LSM-trees work!*

# log-structured merge-tree

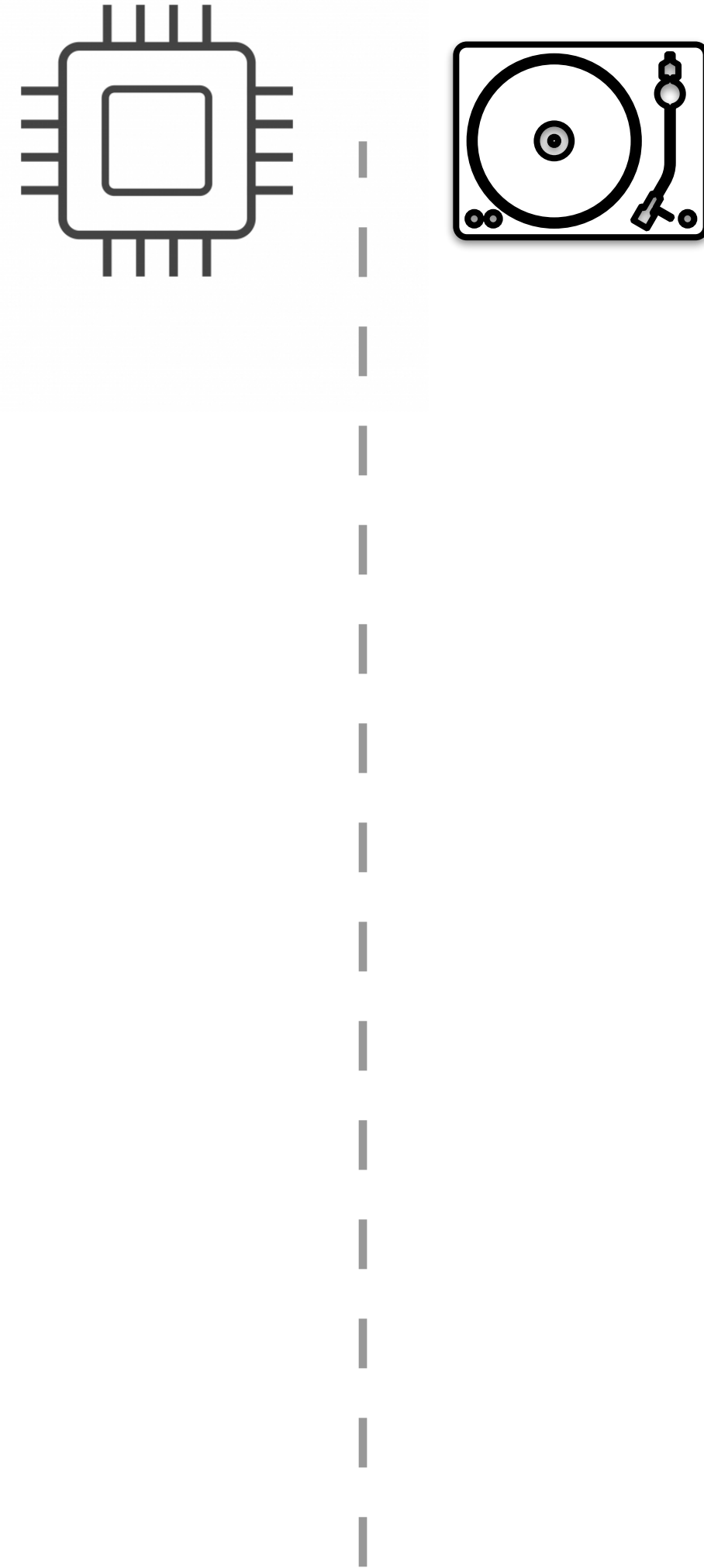
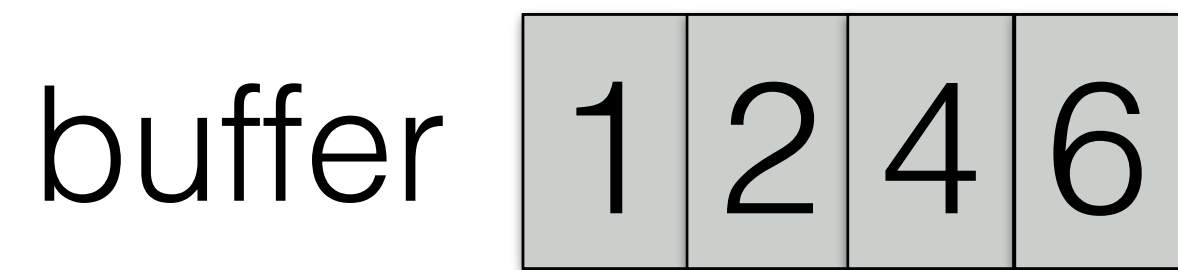
buffer



# log-structured merge-tree

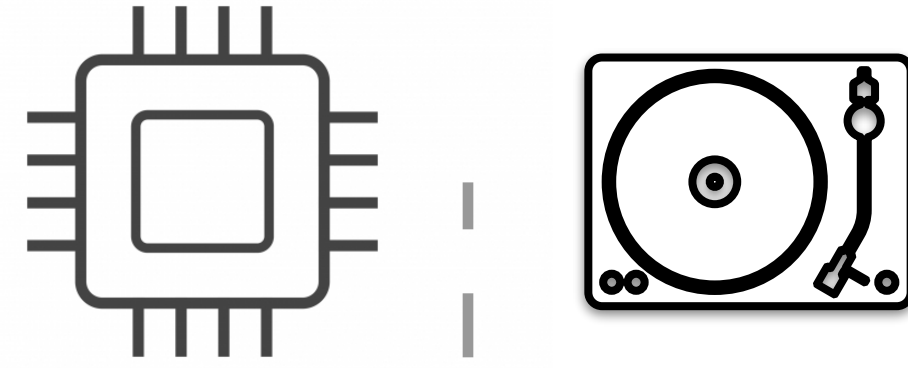


# log-structured merge-tree





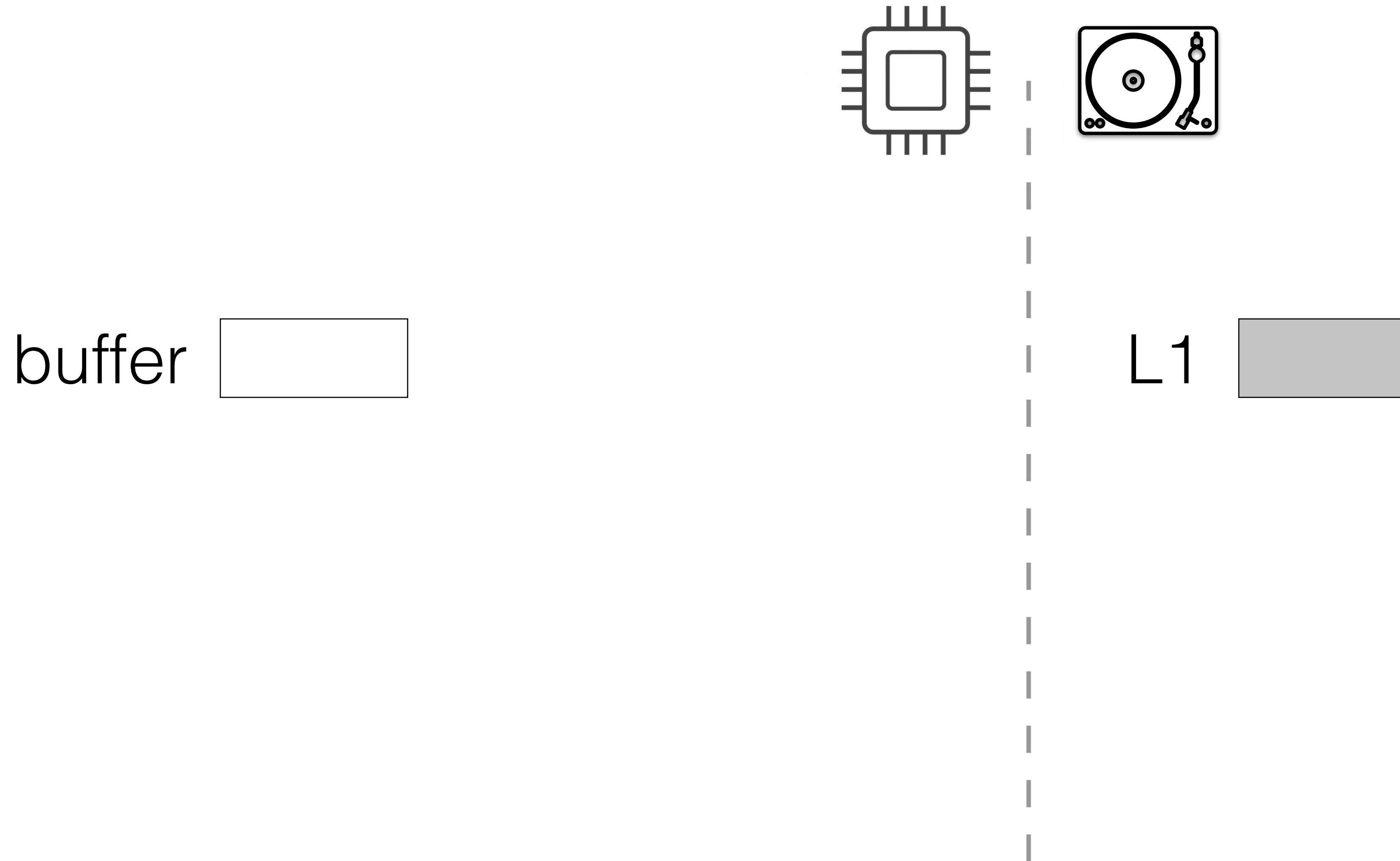
# log-structured merge-tree



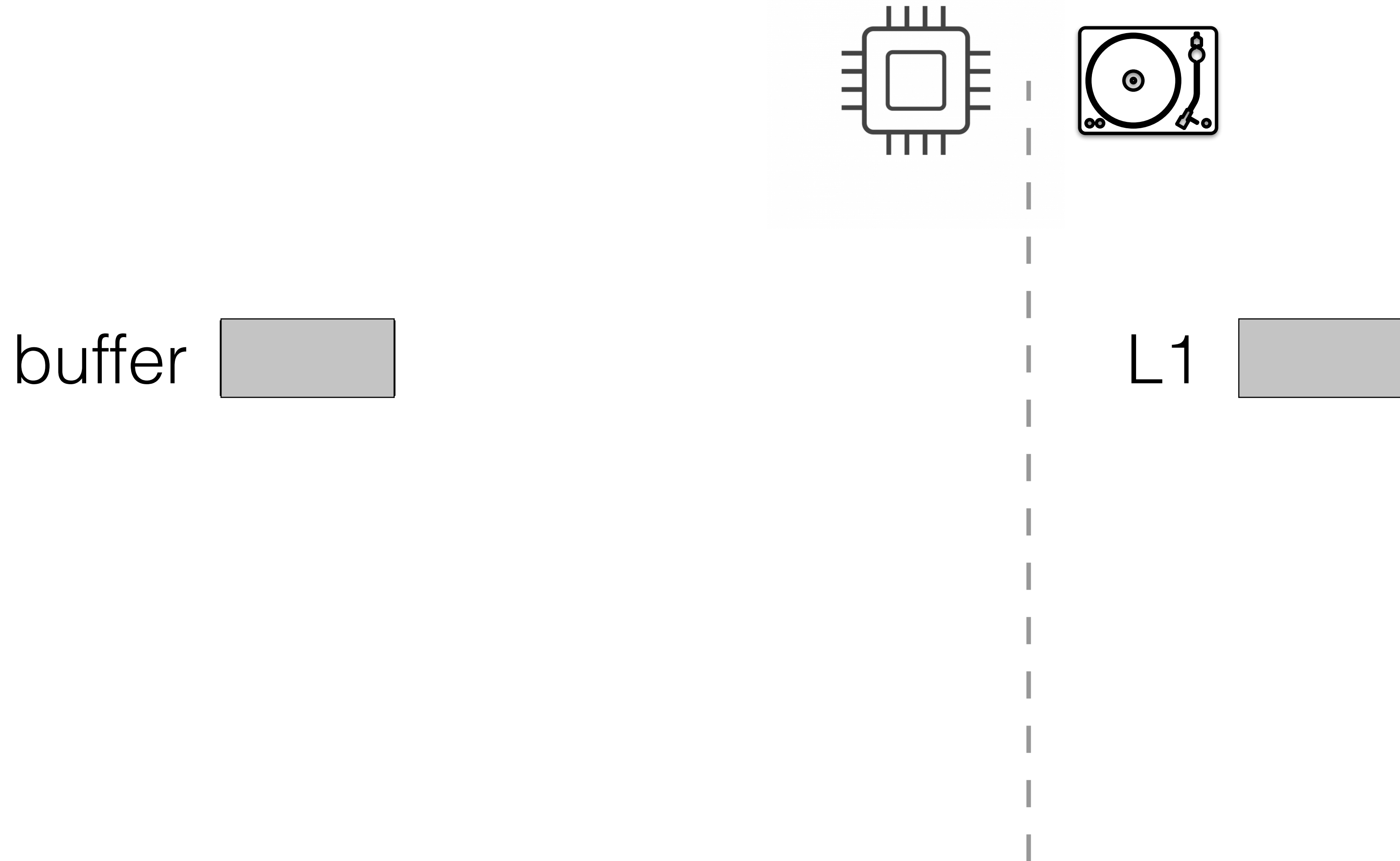
buffer



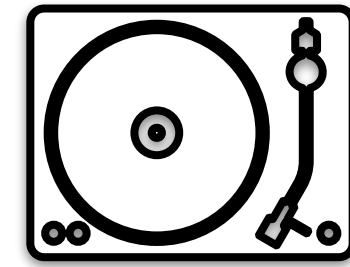
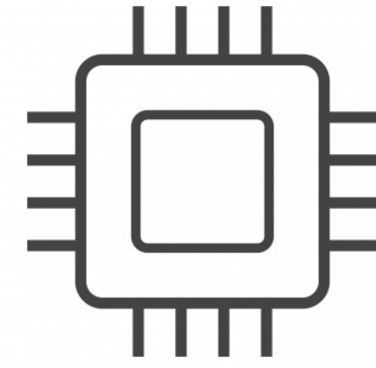
# log-structured merge-tree



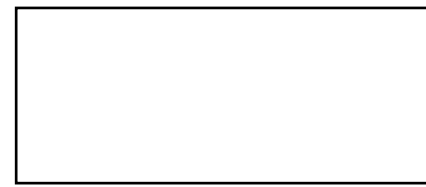
# log-structured merge-tree



# log-structured merge-tree



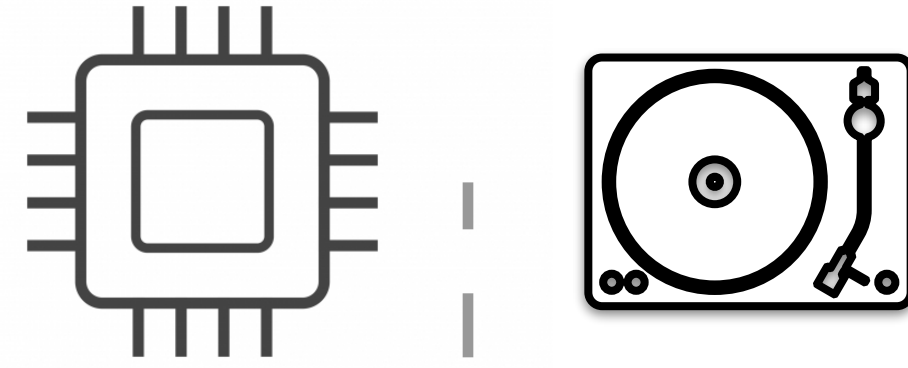
buffer



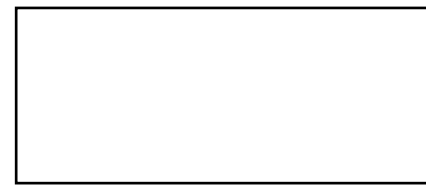
L1



# log-structured merge-tree



buffer



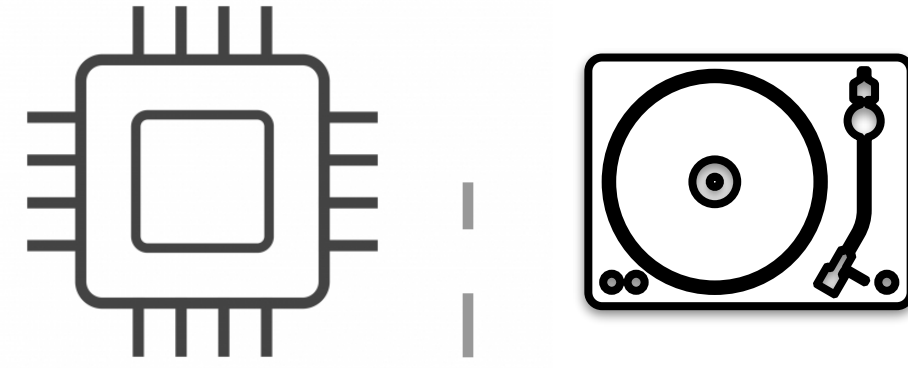
L1

L2

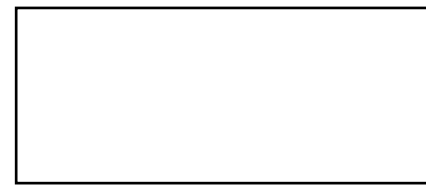


**compaction**

# log-structured merge-tree



buffer



L1



**size ratio = T**

L2

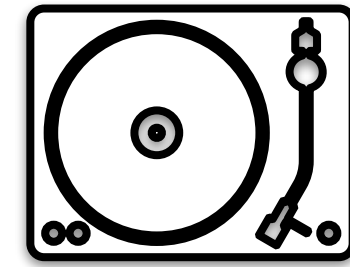
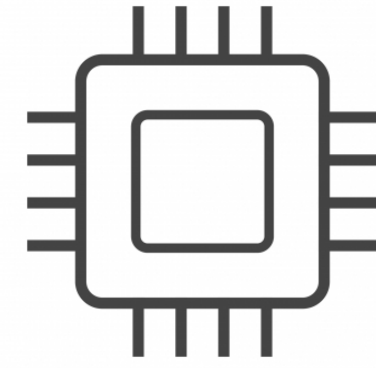


L3

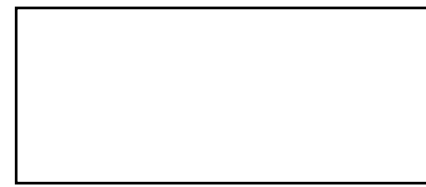


exponentially larger capacity

# log-structured merge-tree



buffer



L1



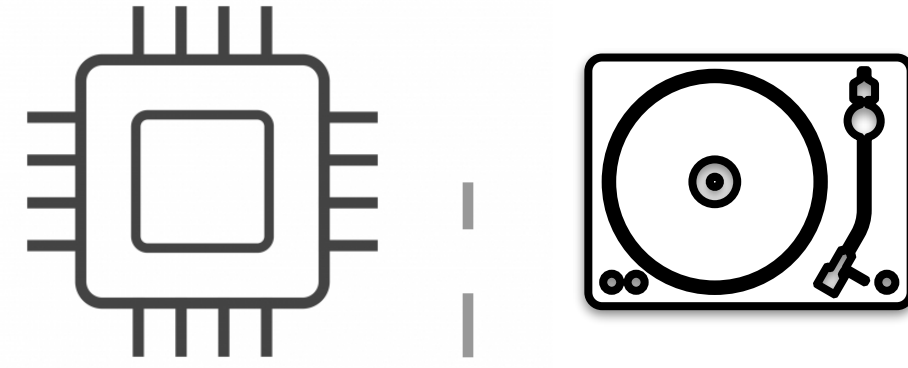
L2



L3



# log-structured merge-tree



buffer



L1



L2

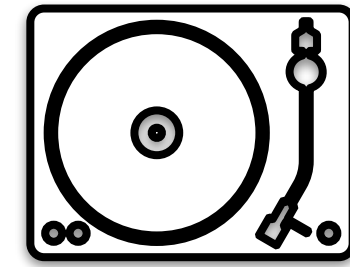
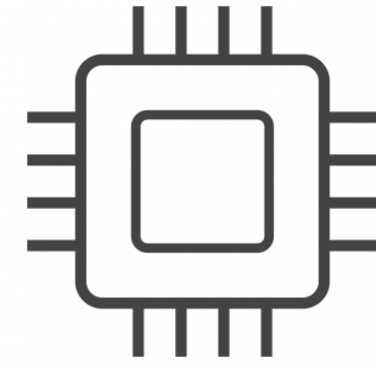


L3





# log-structured merge-tree



buffer



L1



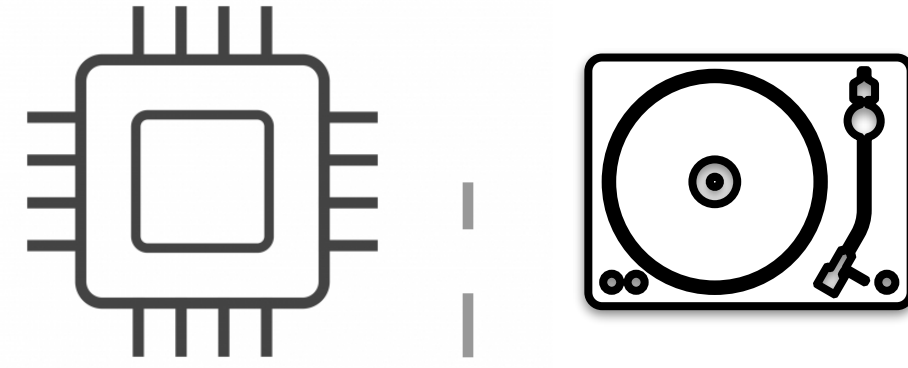
L2



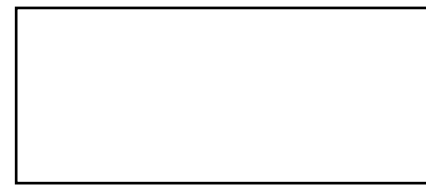
L3



# log-structured merge-tree



buffer



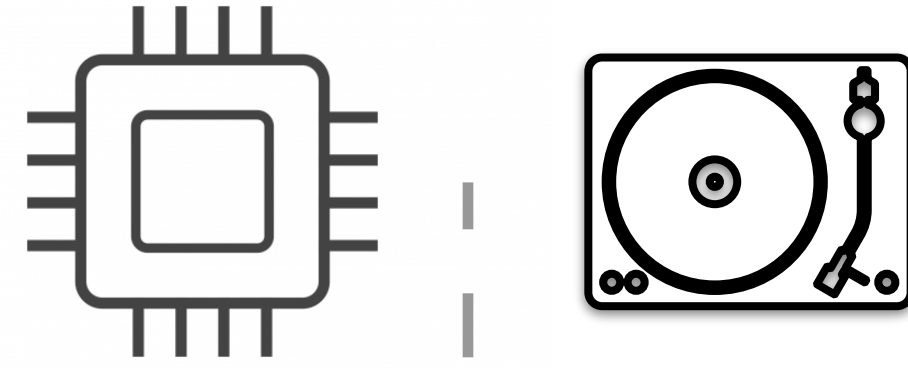
L1

L2

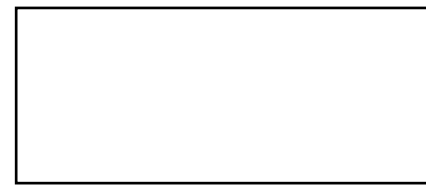
L3



# log-structured merge-tree



buffer



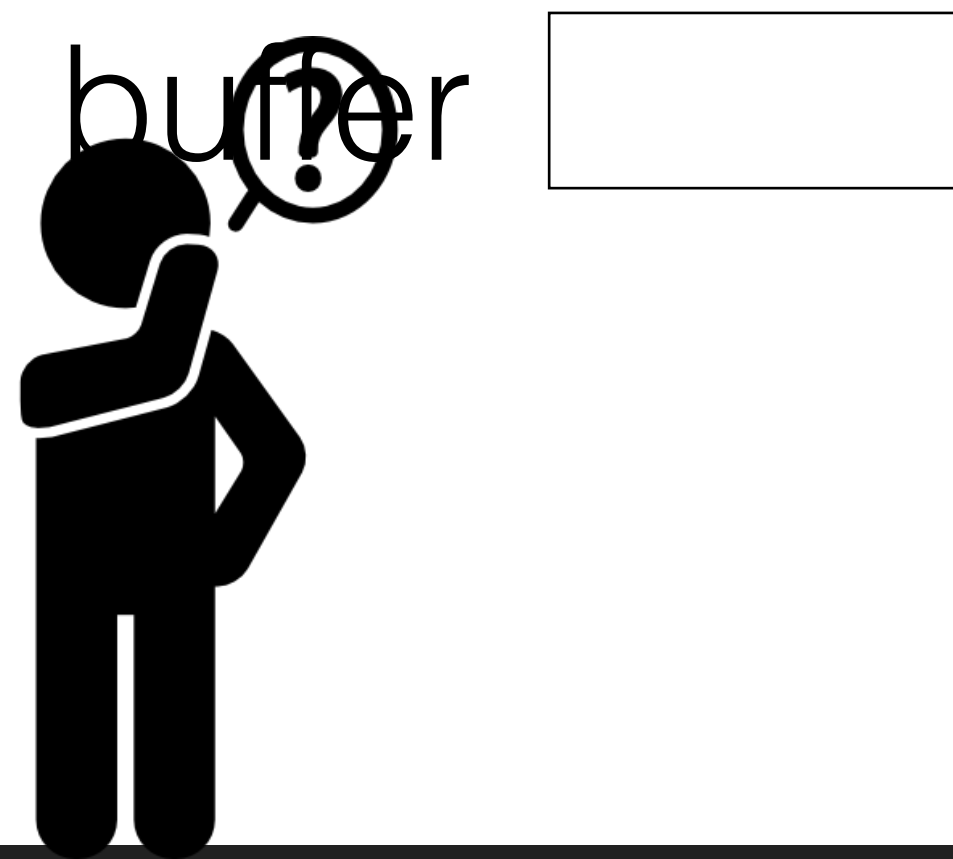
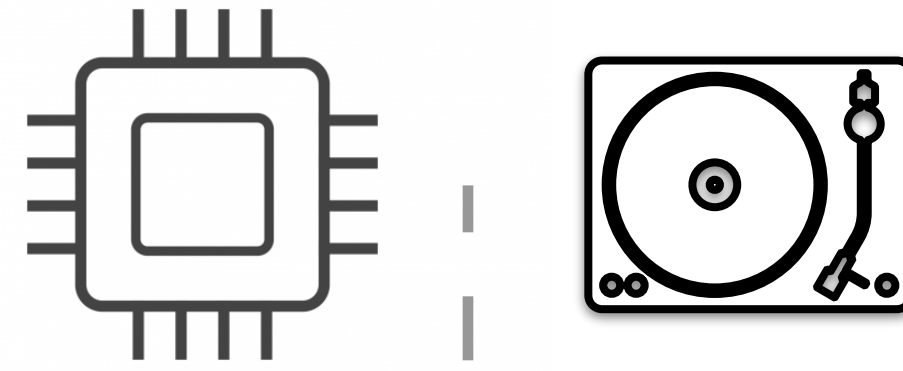
L1

L2

L3



# log-structured merge-tree



L1

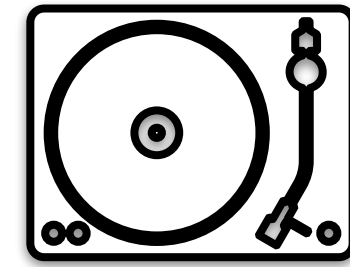
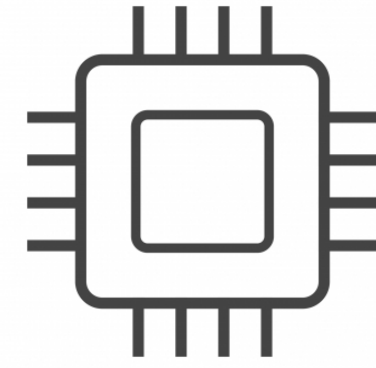
L2

L3

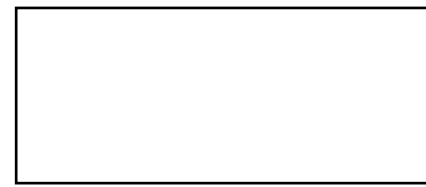
burst of I/Os  
prolonged write stalls

How do we avoid this?

# log-structured merge-tree



buffer



**partial compaction**

L1



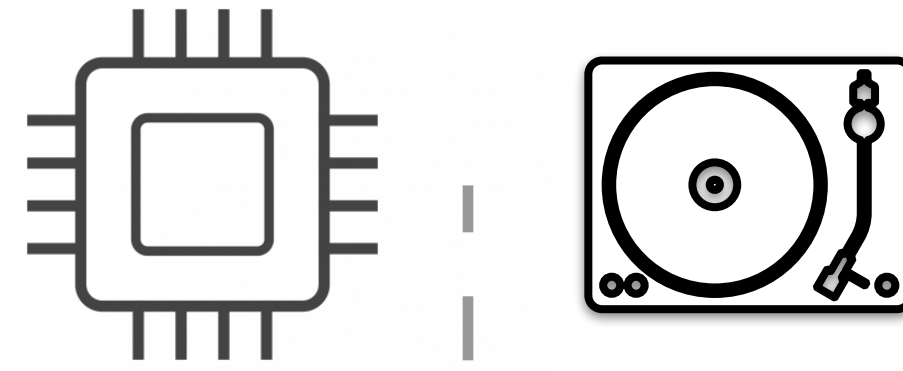
L2



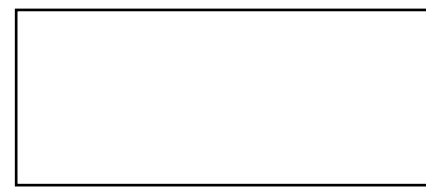
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



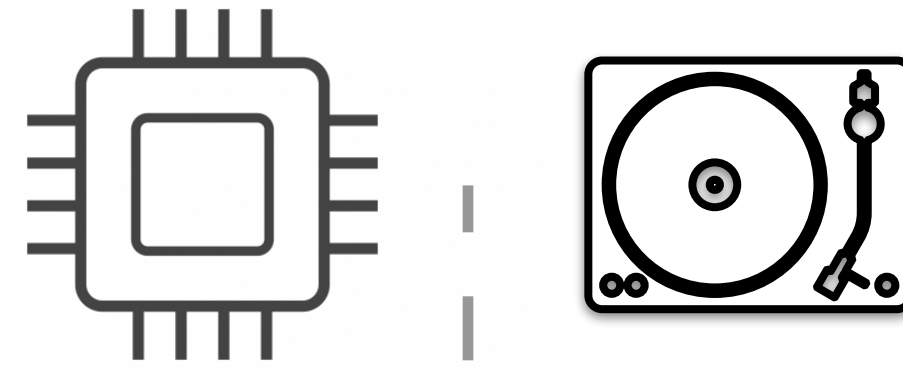
L2



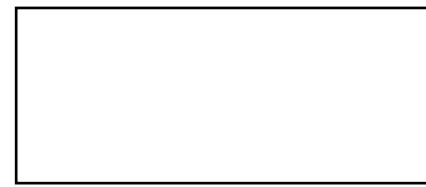
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



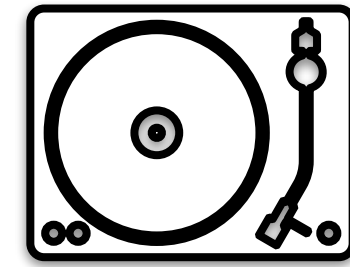
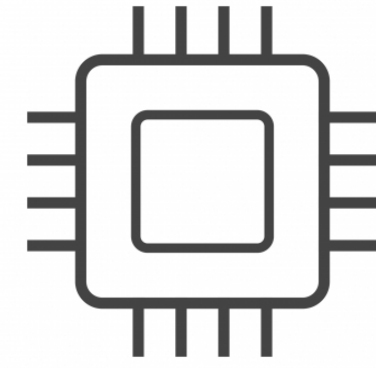
L2



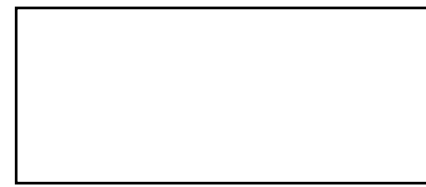
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



L2

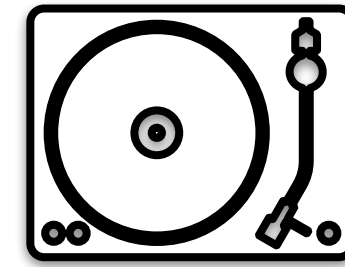
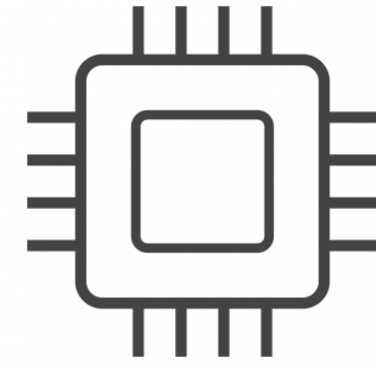


L3





# log-structured merge-tree



buffer



**partial compaction**

L1



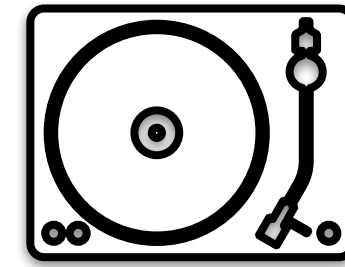
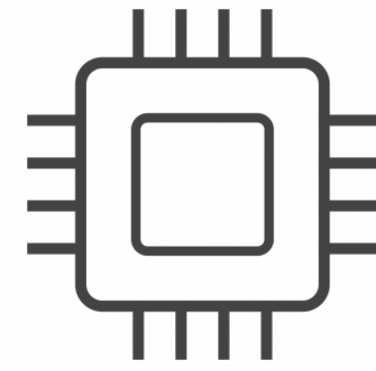
L2



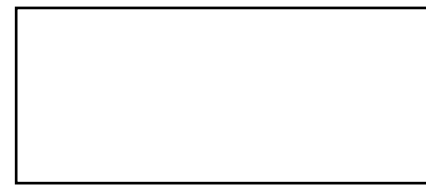
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



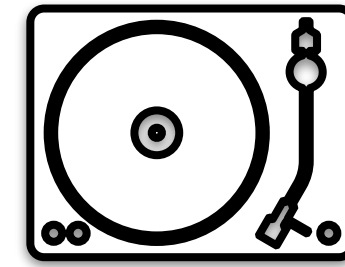
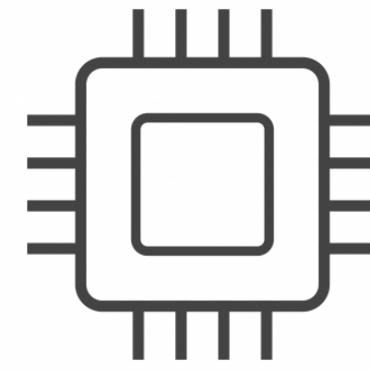
L2



L3



# log-structured merge-tree



buffer



**partial compaction**

L1



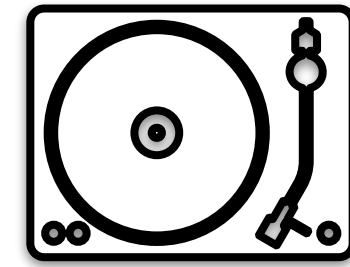
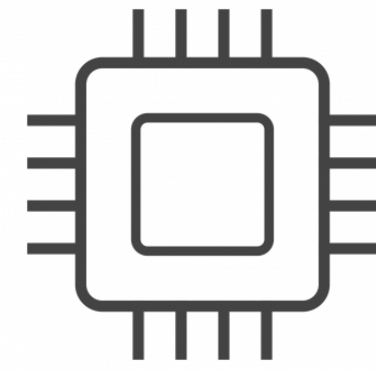
L2



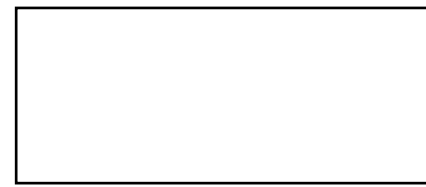
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



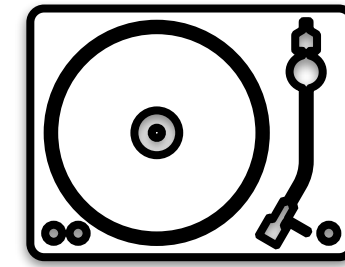
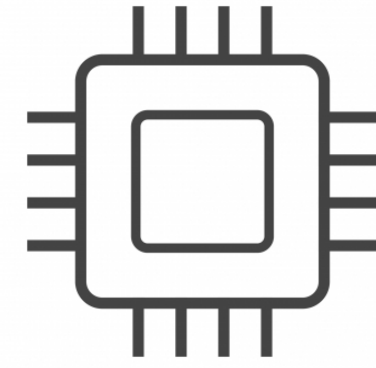
L2



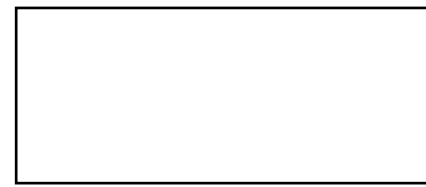
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



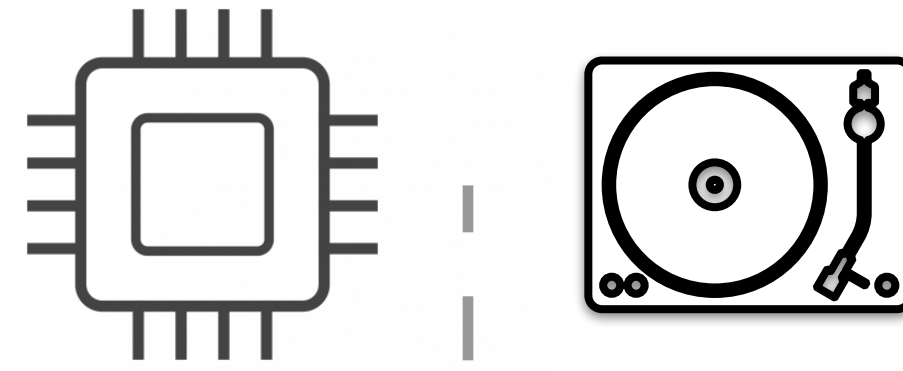
L2



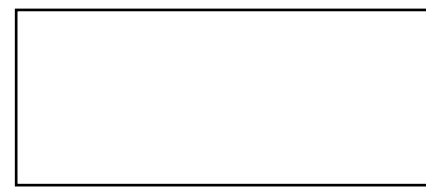
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



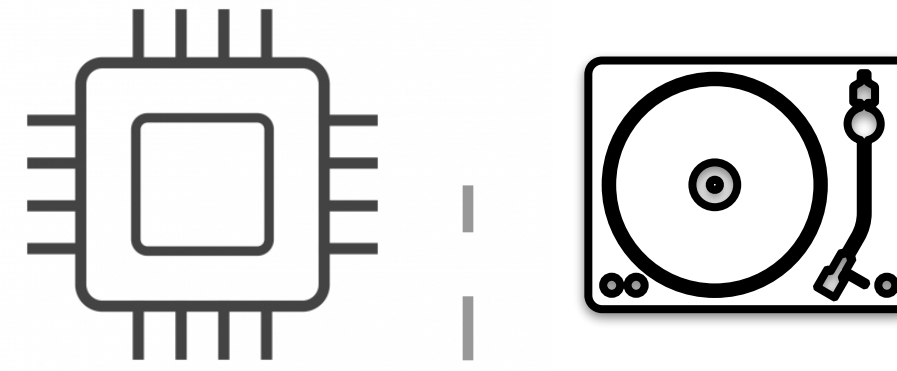
L2



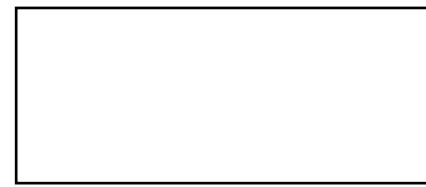
L3



# log-structured merge-tree



buffer



**partial compaction**

L1



L2



L3



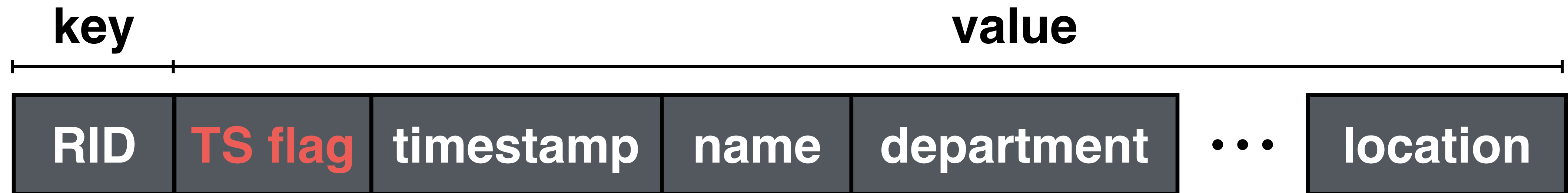
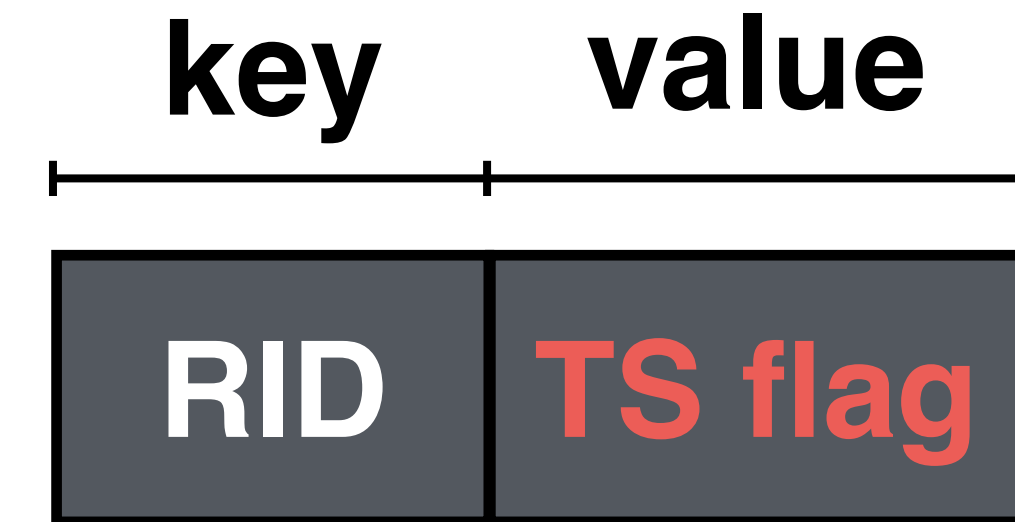
amortized compaction cost

*OK! But how to delete in LSM?*

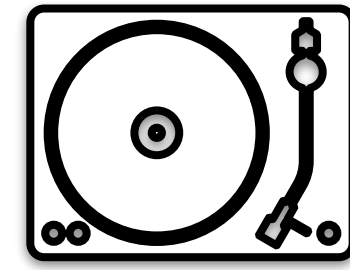
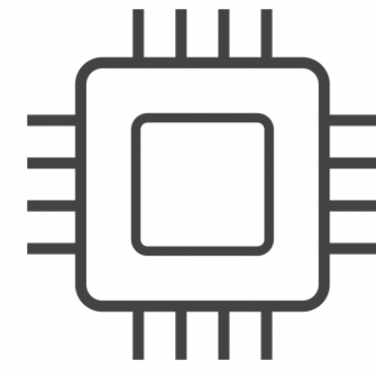
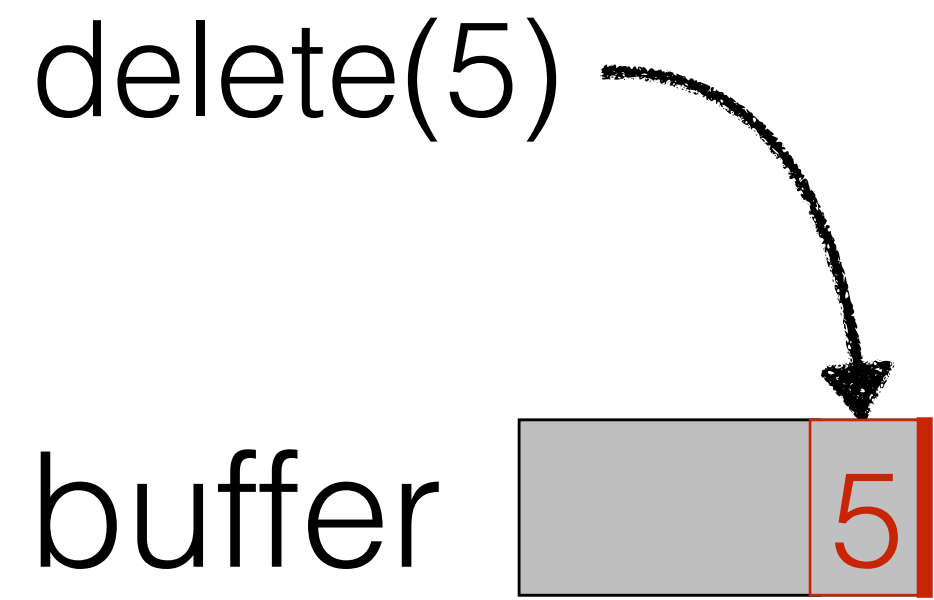


# deletes in LSM-tree

delete := insert tombstone



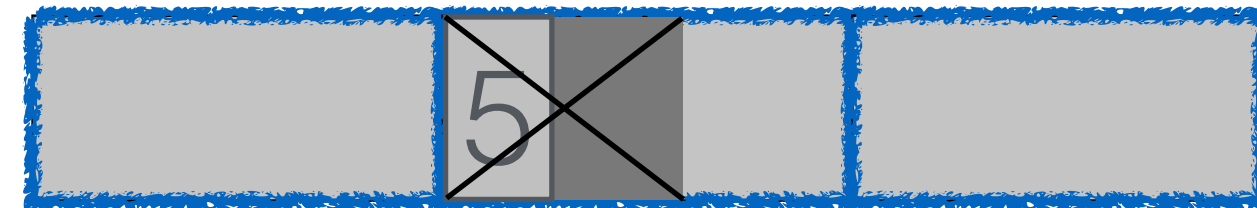
# deletes in LSM-tree



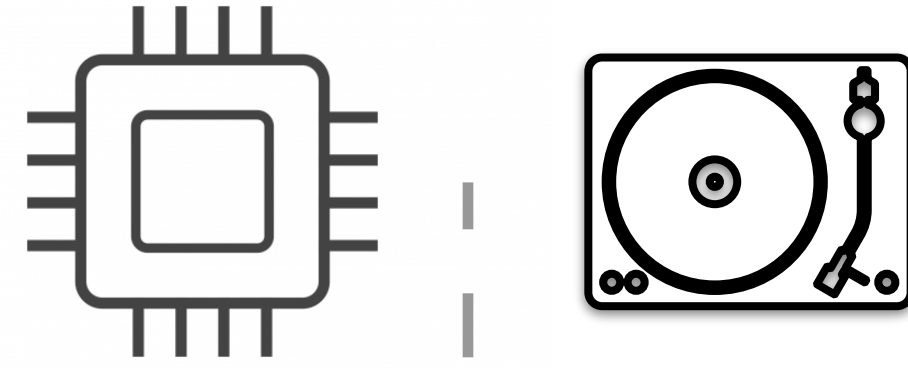
L1

L2

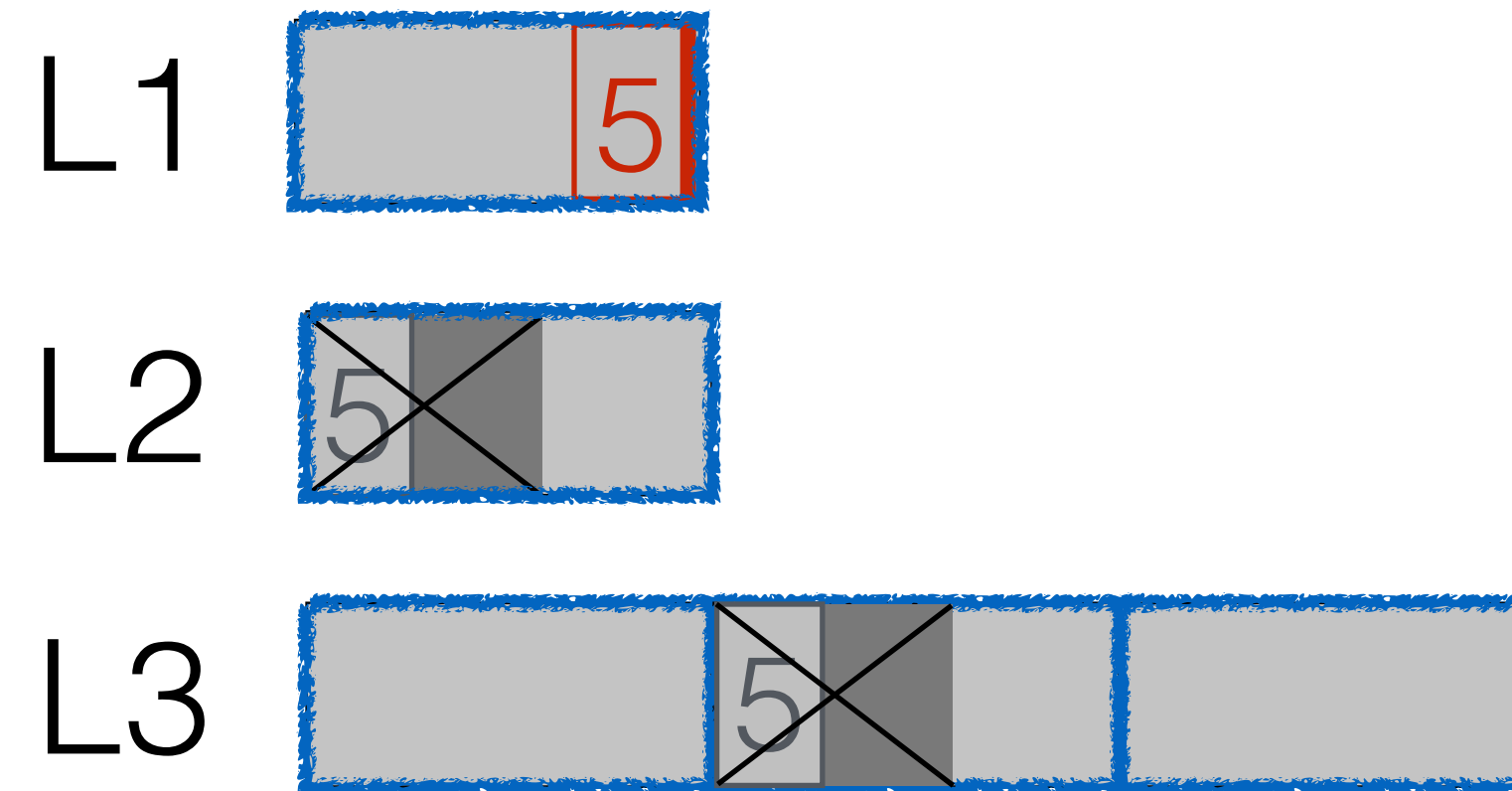
L3



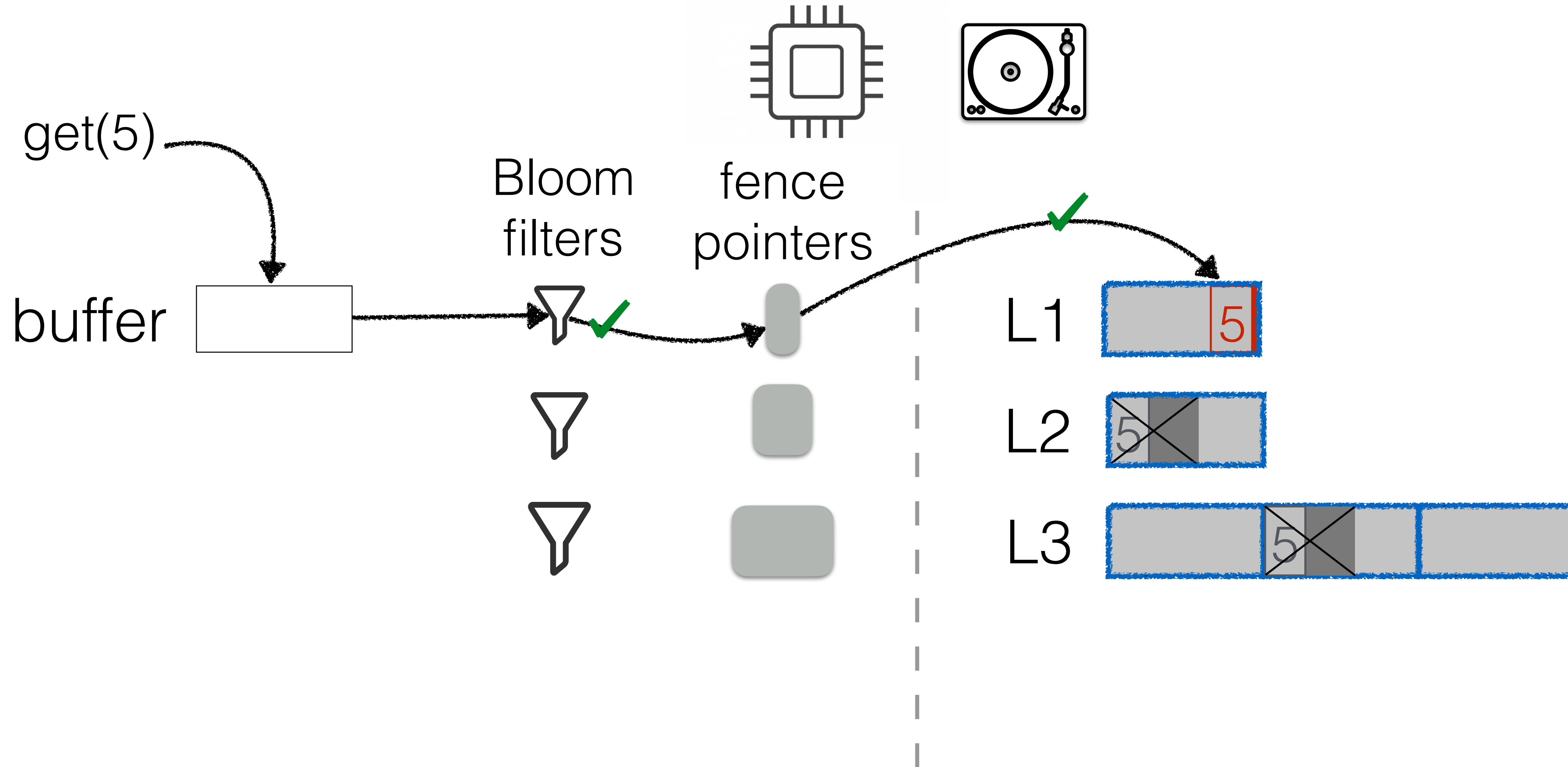
# deletes in LSM-tree



buffer



# deletes in LSM-tree

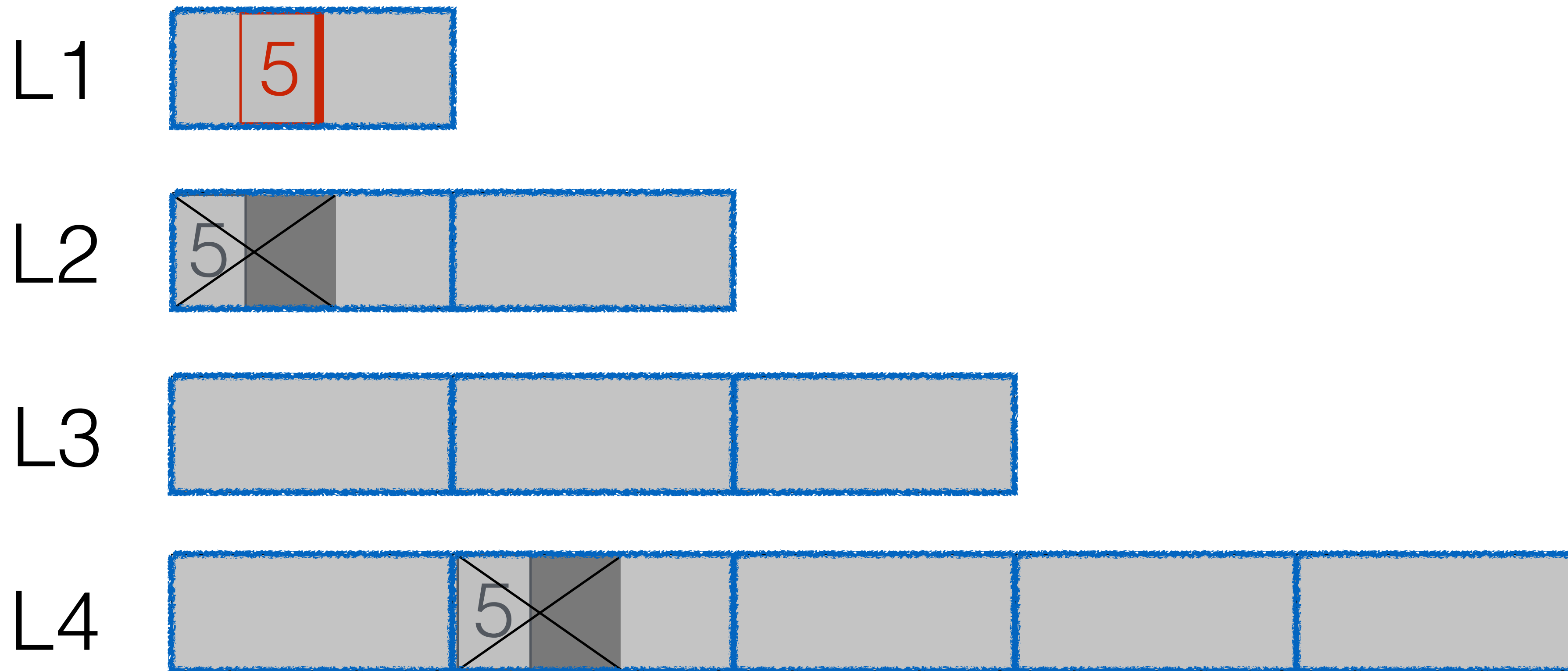


the problems



**out-of-place deletes**

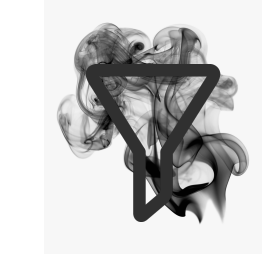
# out-of-place deletes



Problem?

# out-of-place deletes

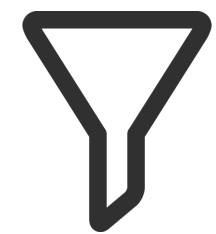
Bloom filters



L1



L2



L3



L4



poor read perf. ❌

write amplification ❌

space amplification ❌



# the problems

poor read perf.

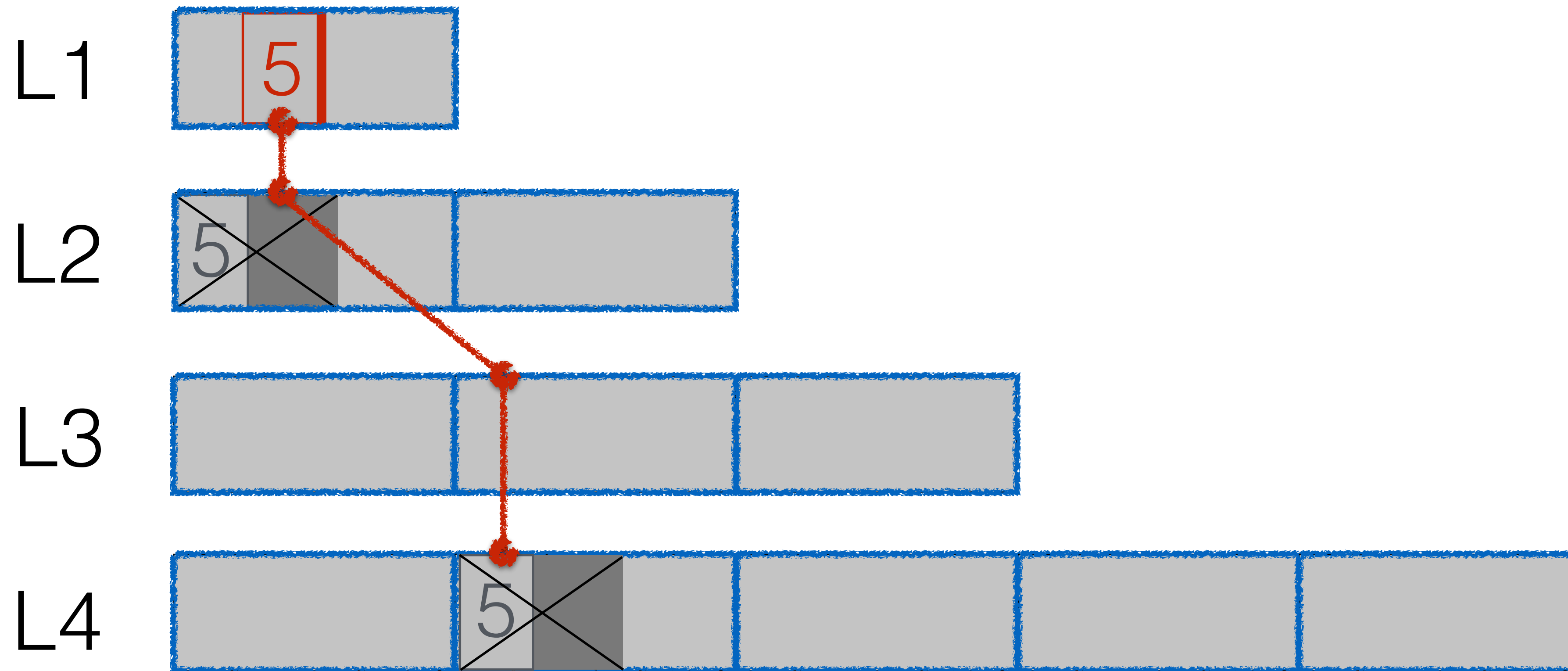
write amplification

space amplification



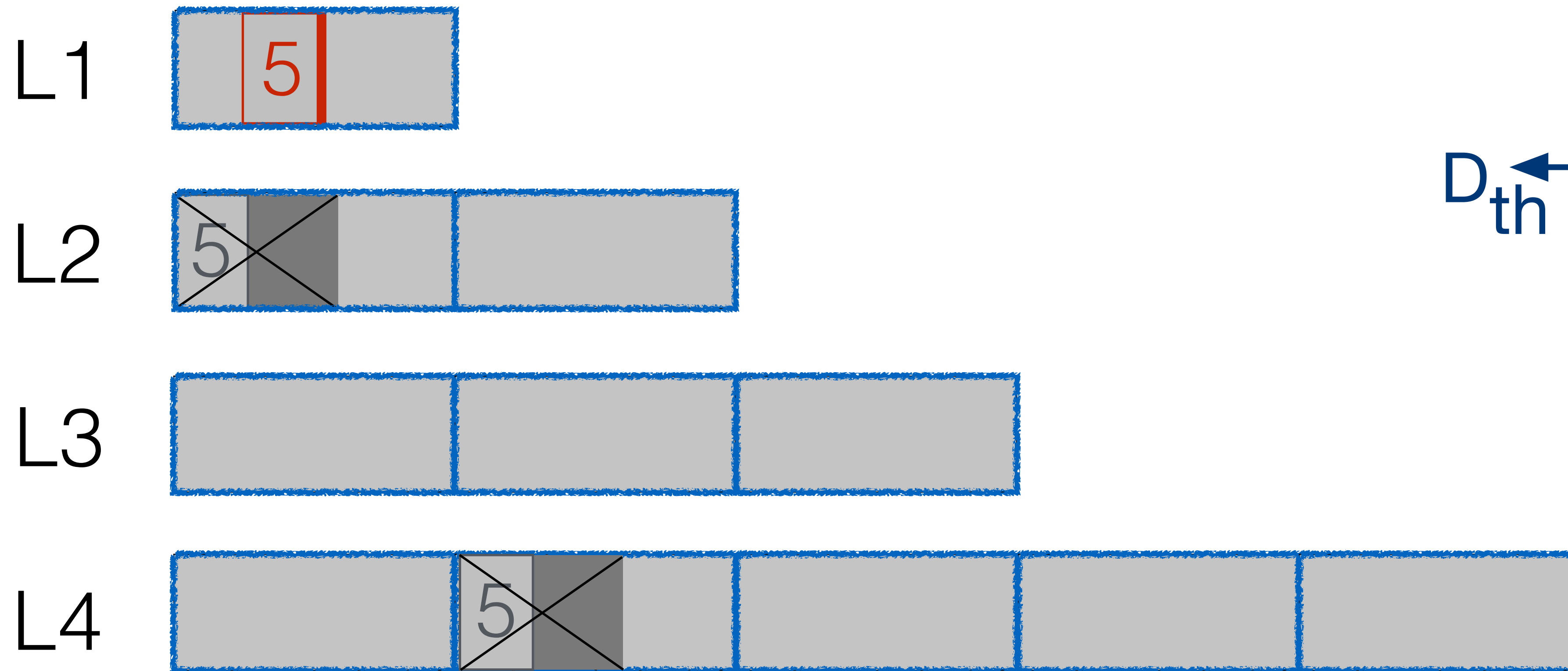
**delete persistence latency**

# delete persistence latency



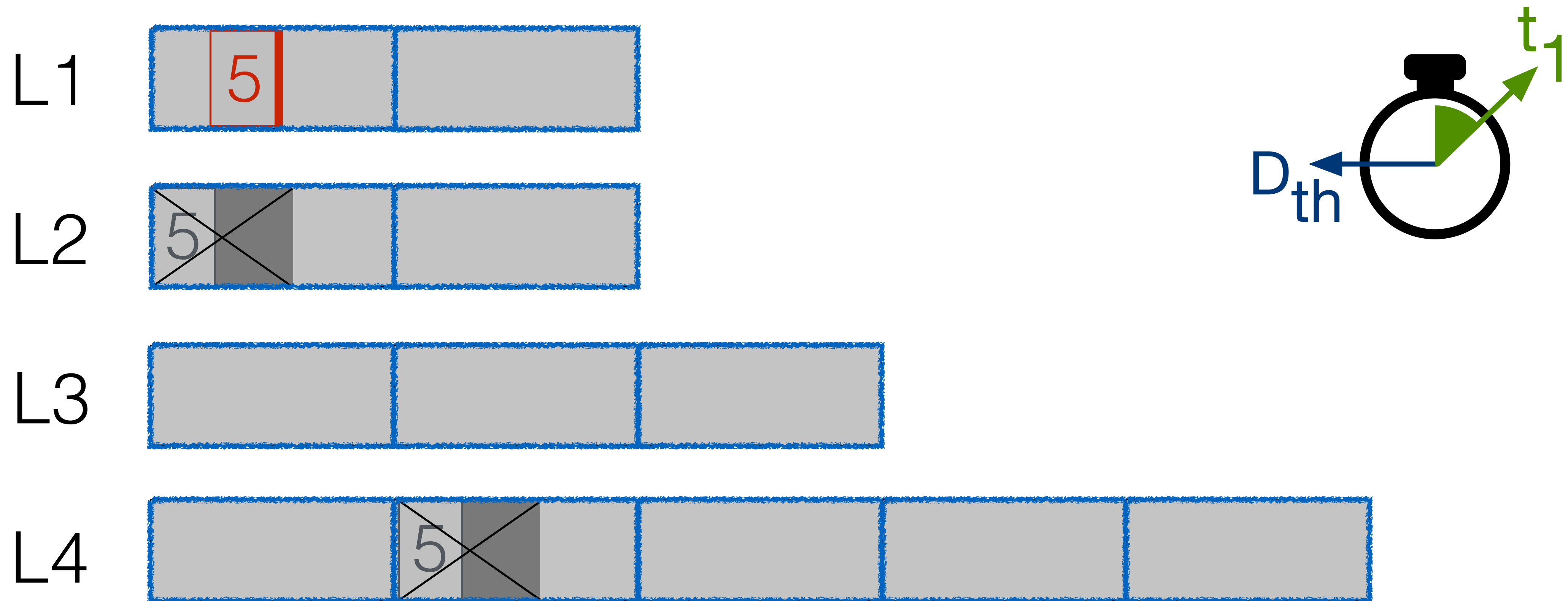
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



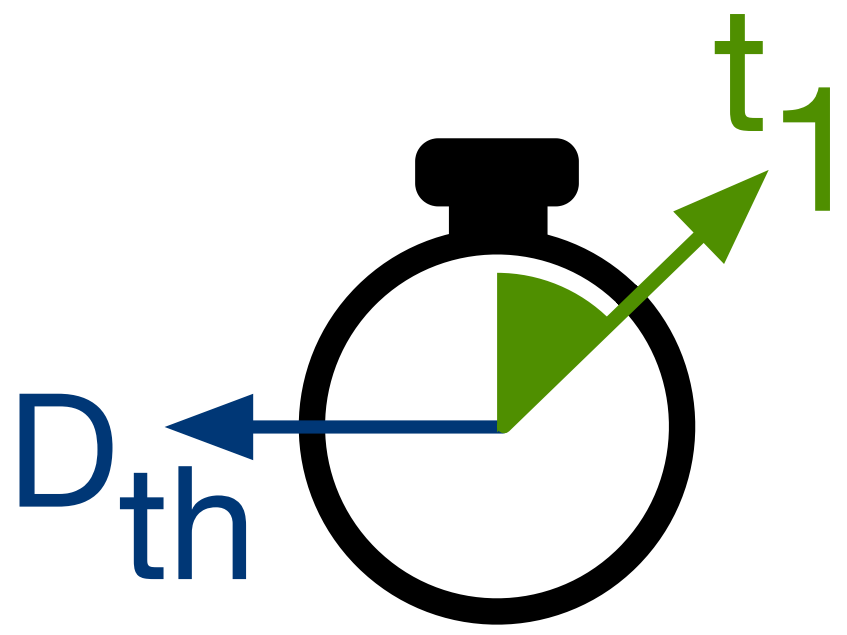
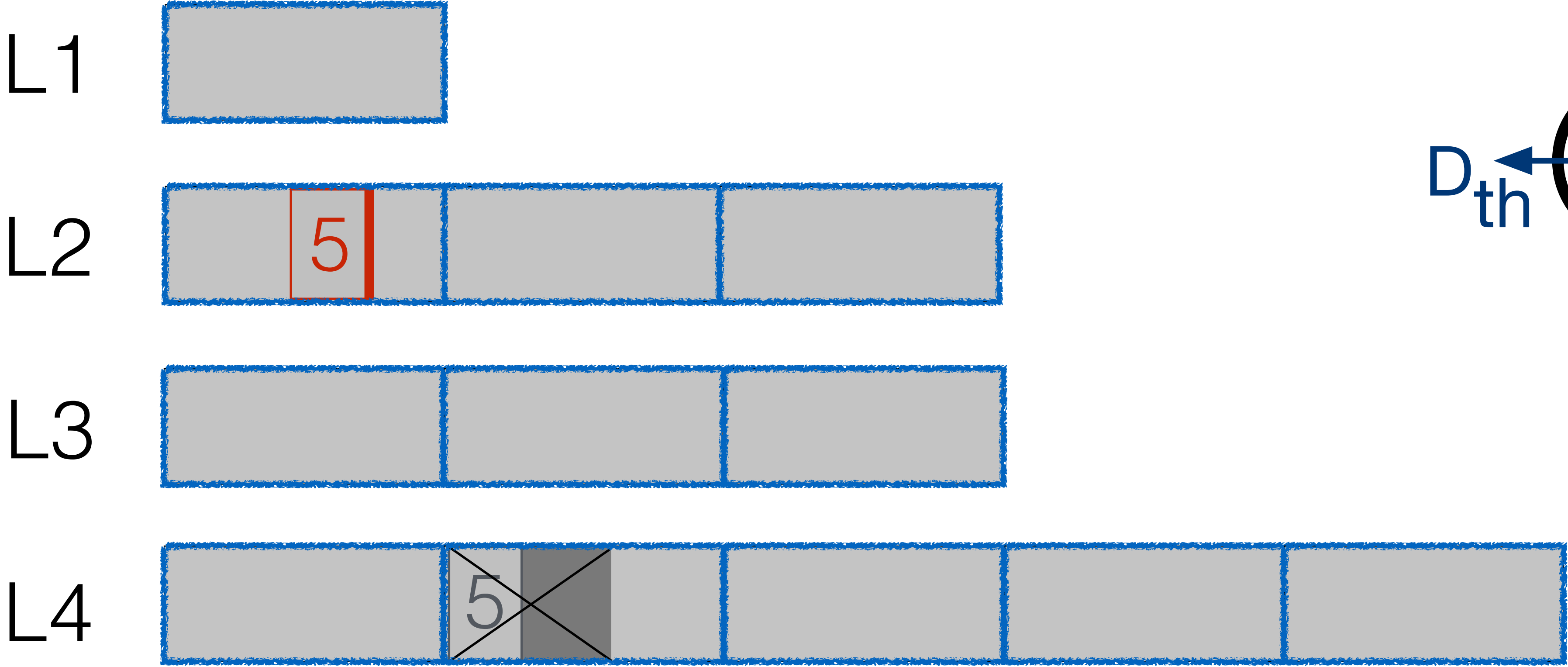
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



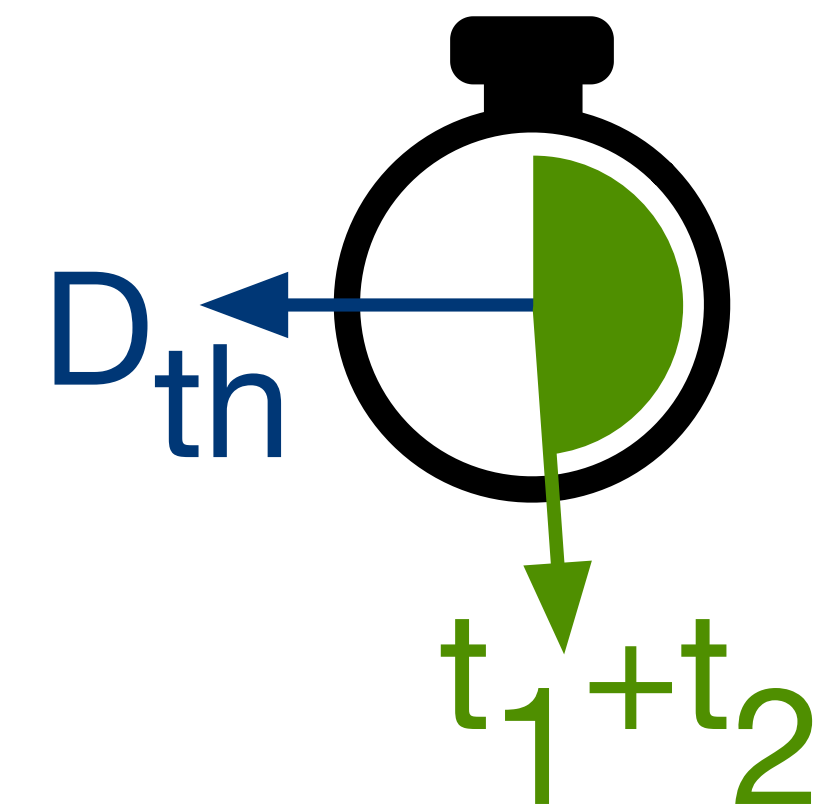
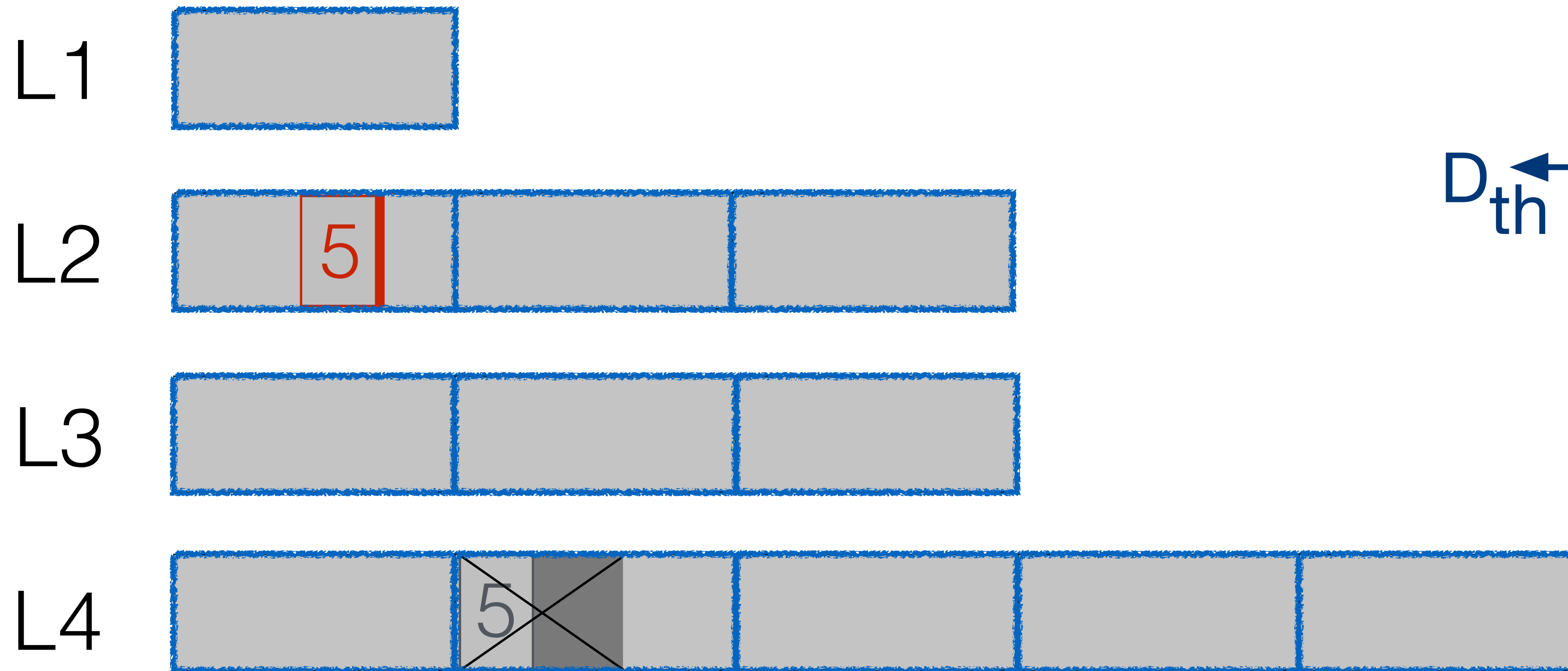
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



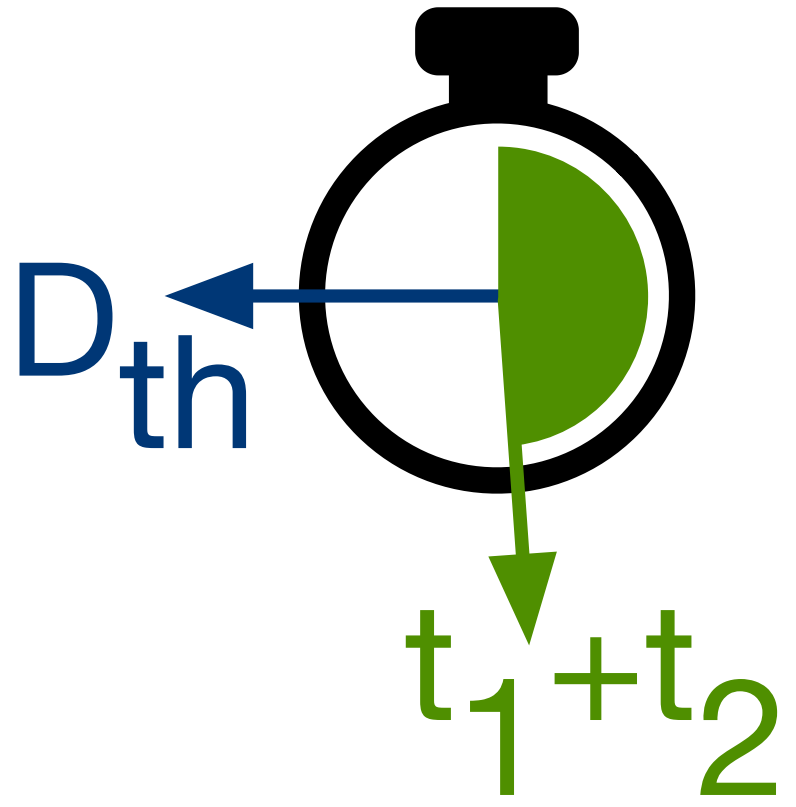
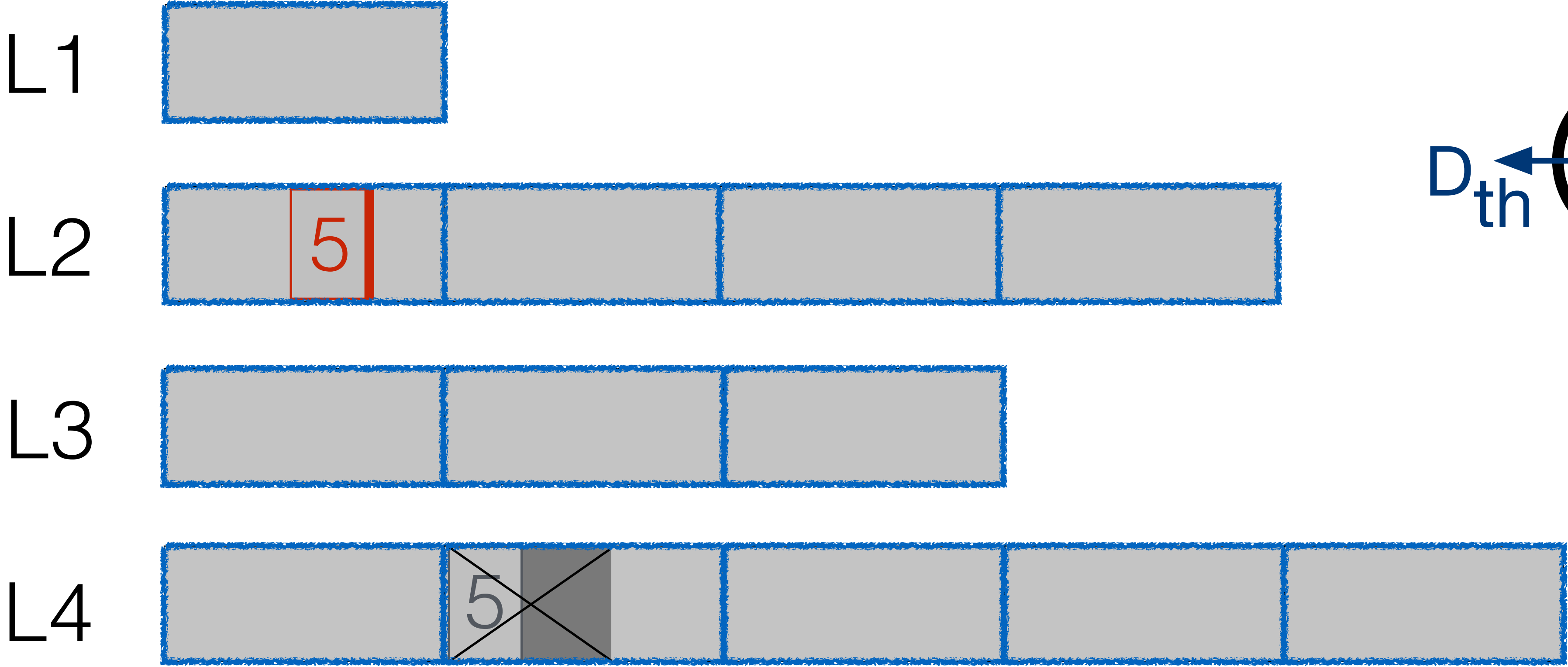
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



# delete persistence latency

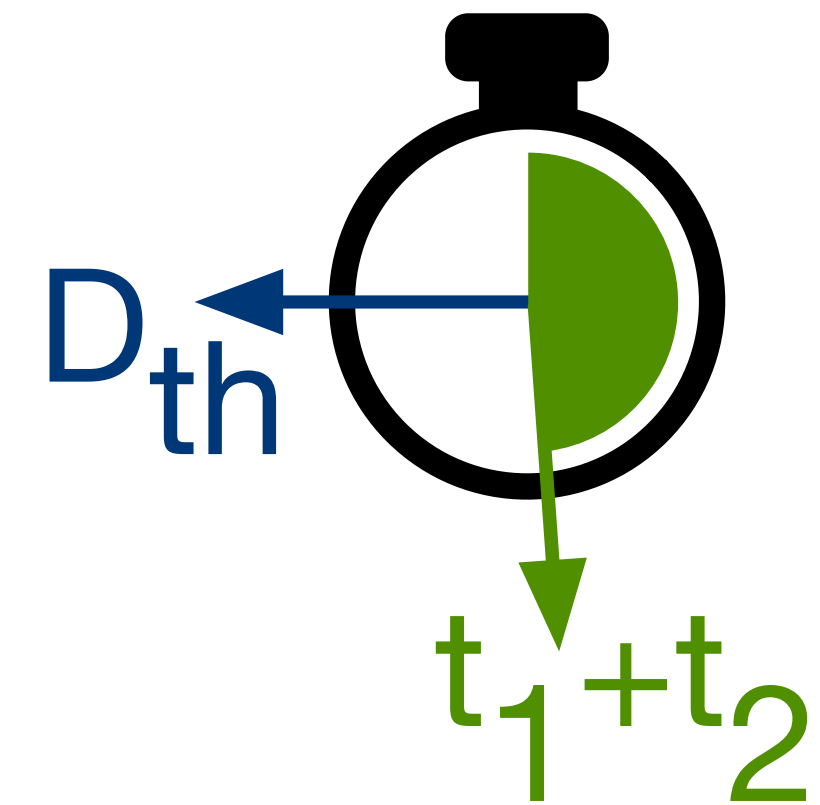
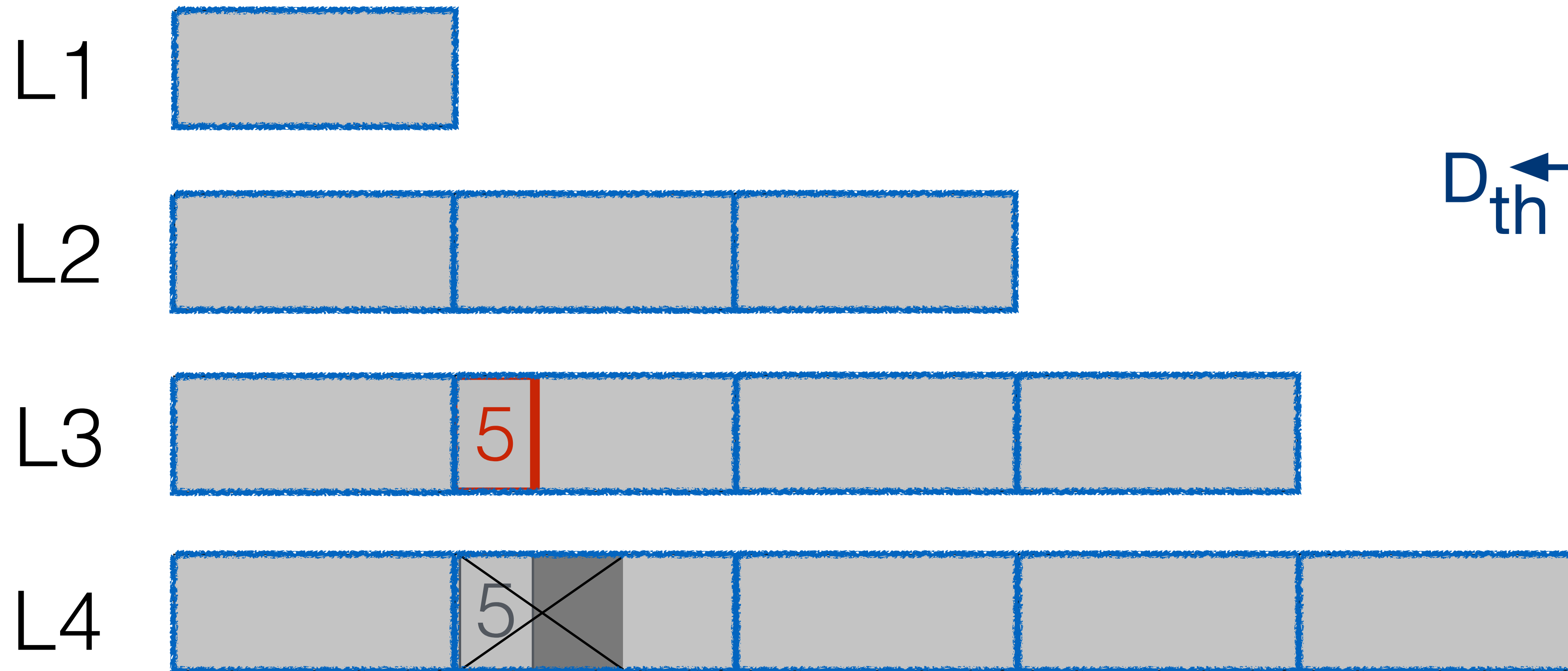
delete(5) within a threshold time:  $D_{th}$





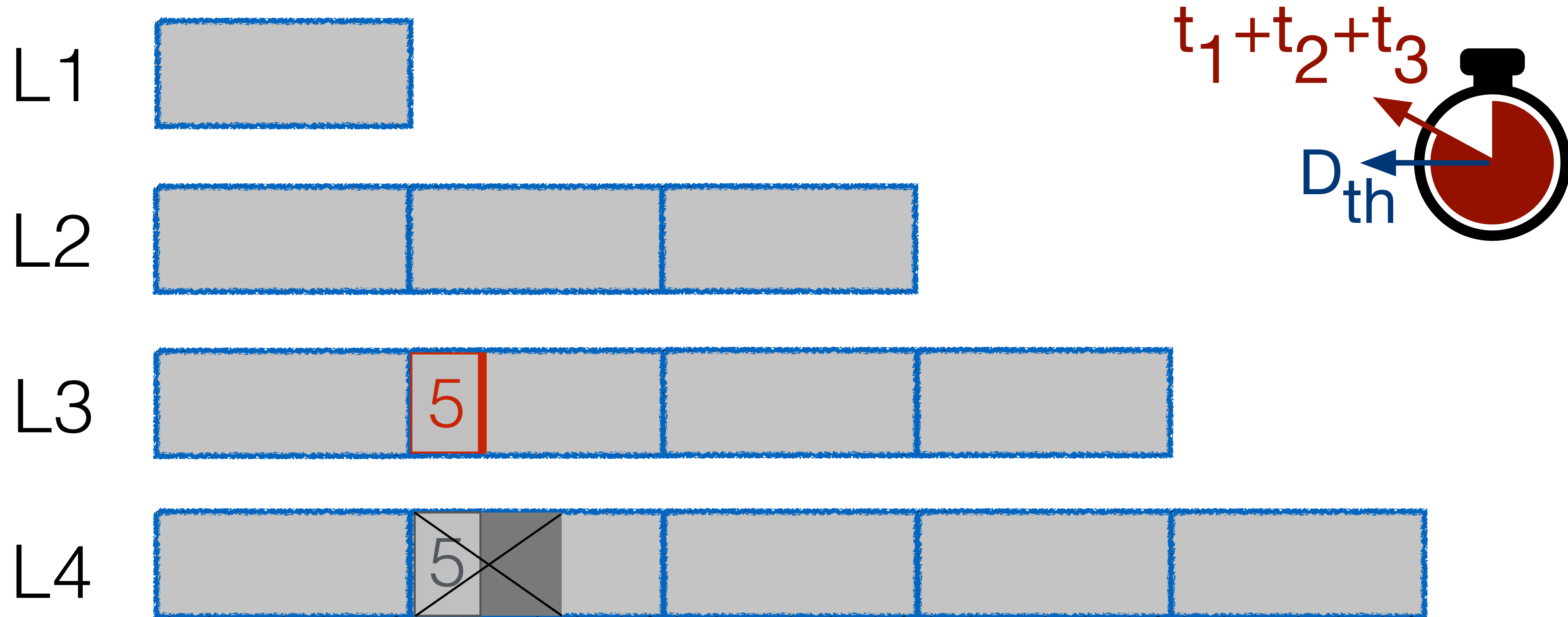
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



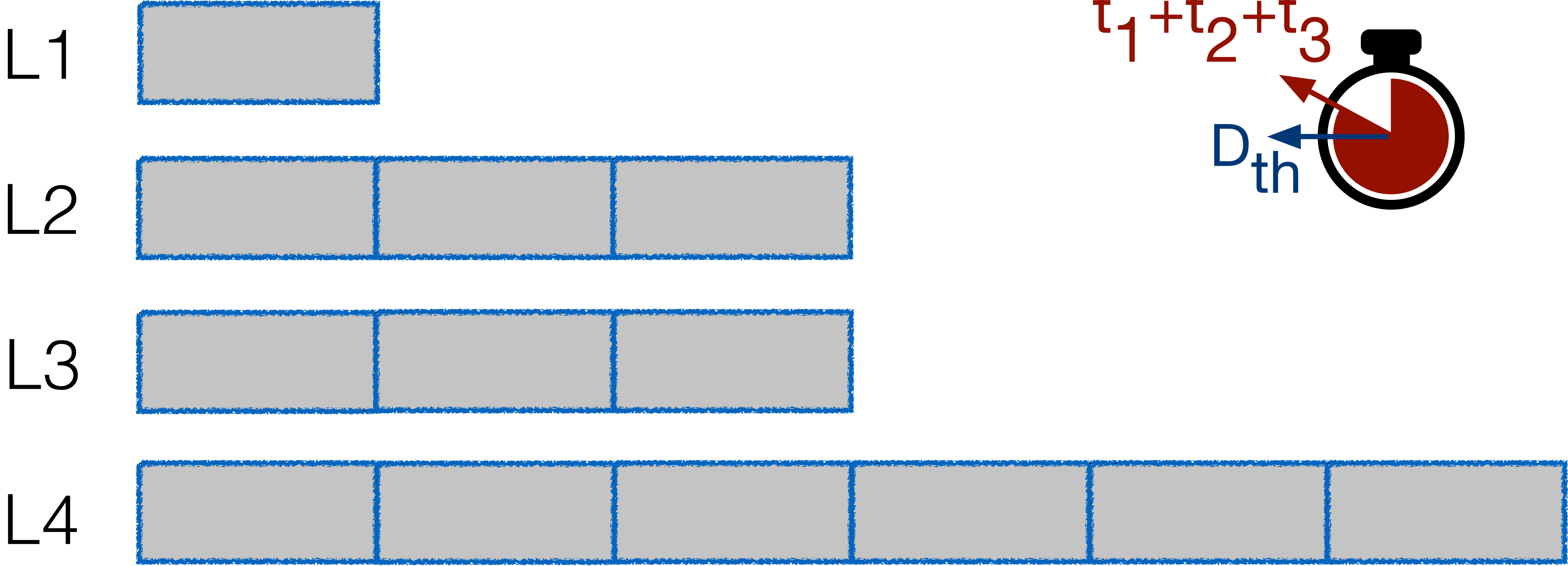
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



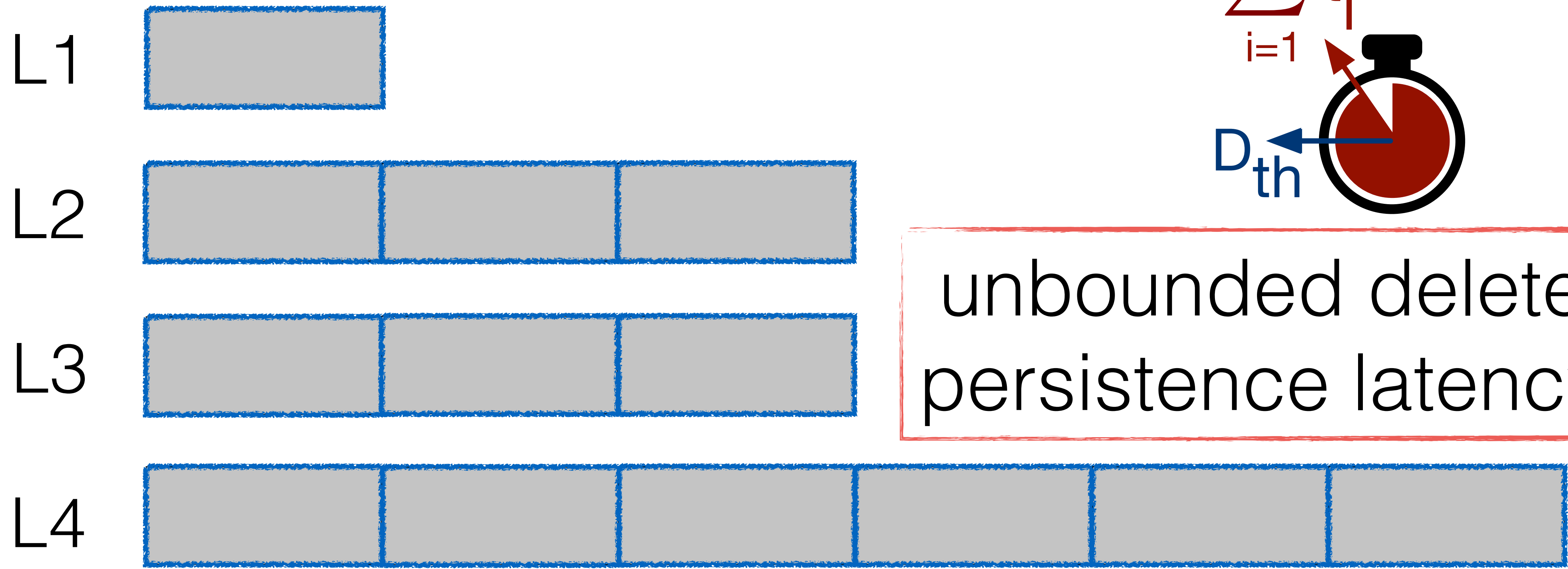
# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



# delete persistence latency

delete(5) within a threshold time:  $D_{th}$



$$\sum_{i=1}^{L-1} t_i$$

unbounded delete persistence latency



# the solution

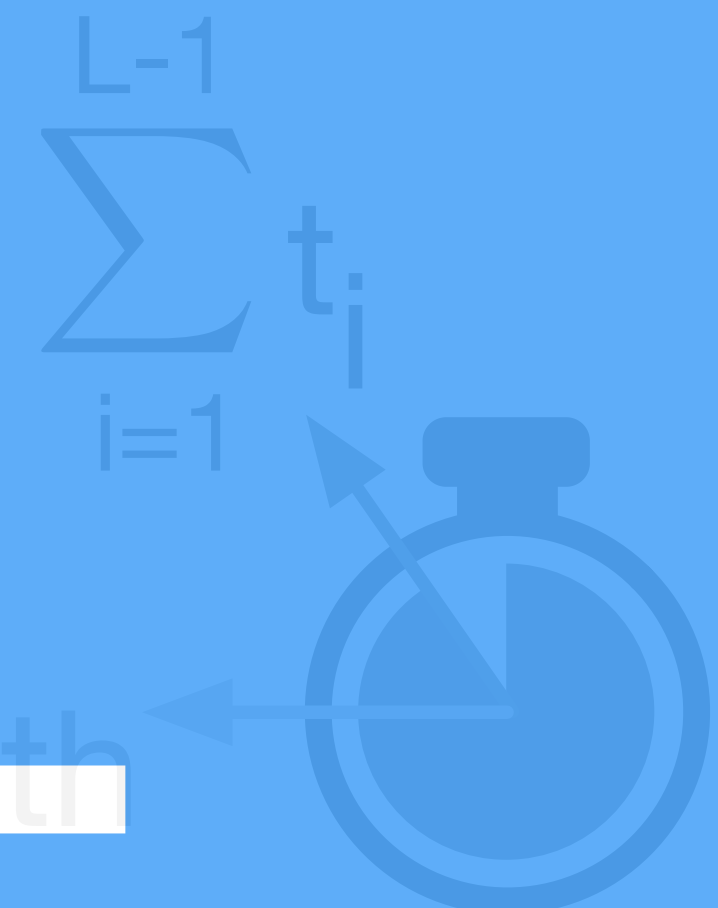
poor read perf.

write amplification

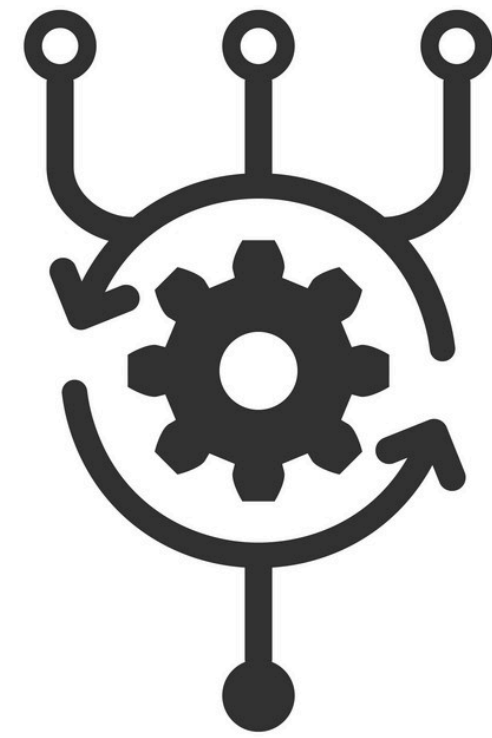
space amplification

# FADE

unbounded delete persistence latency



# FAst DElete



family of  
**compaction**  
**strategies**

# FAst DElete

compaction  
trigger

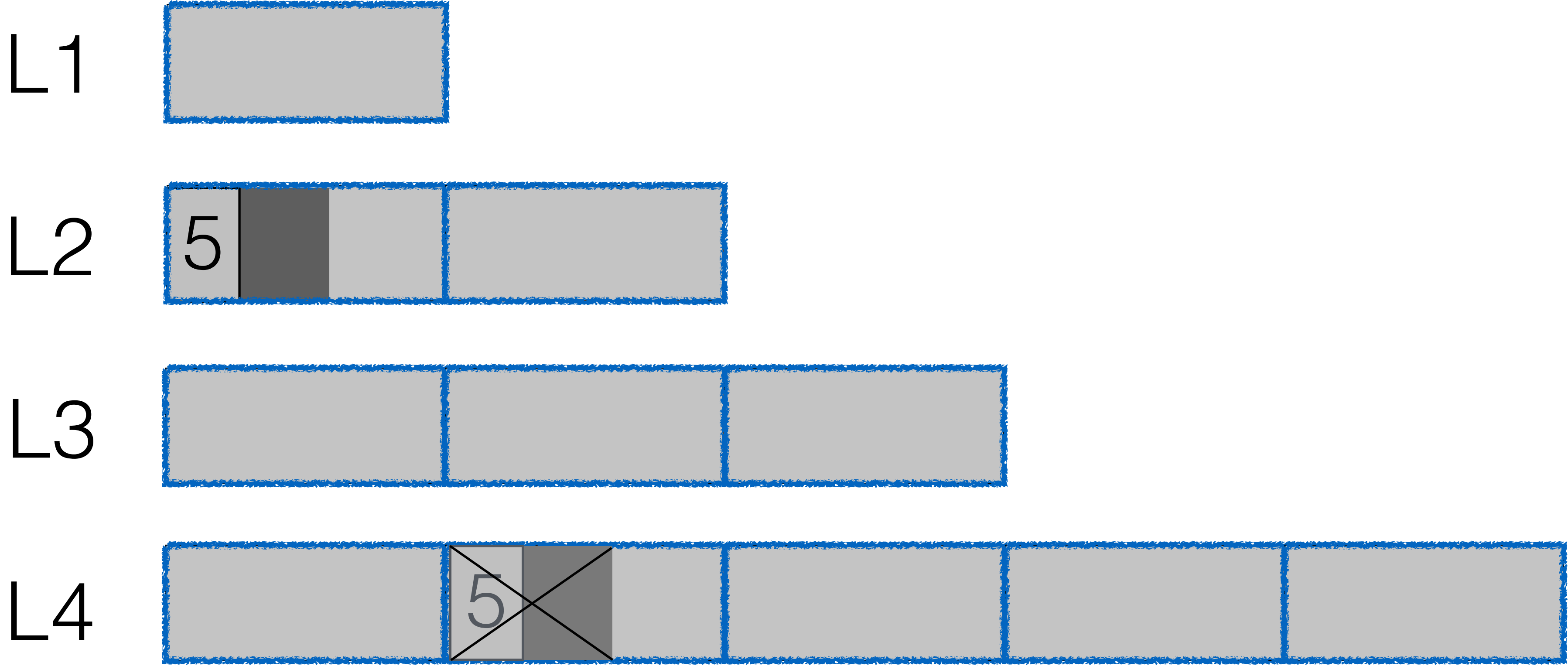


compaction file  
picking policy



# FAst DElete

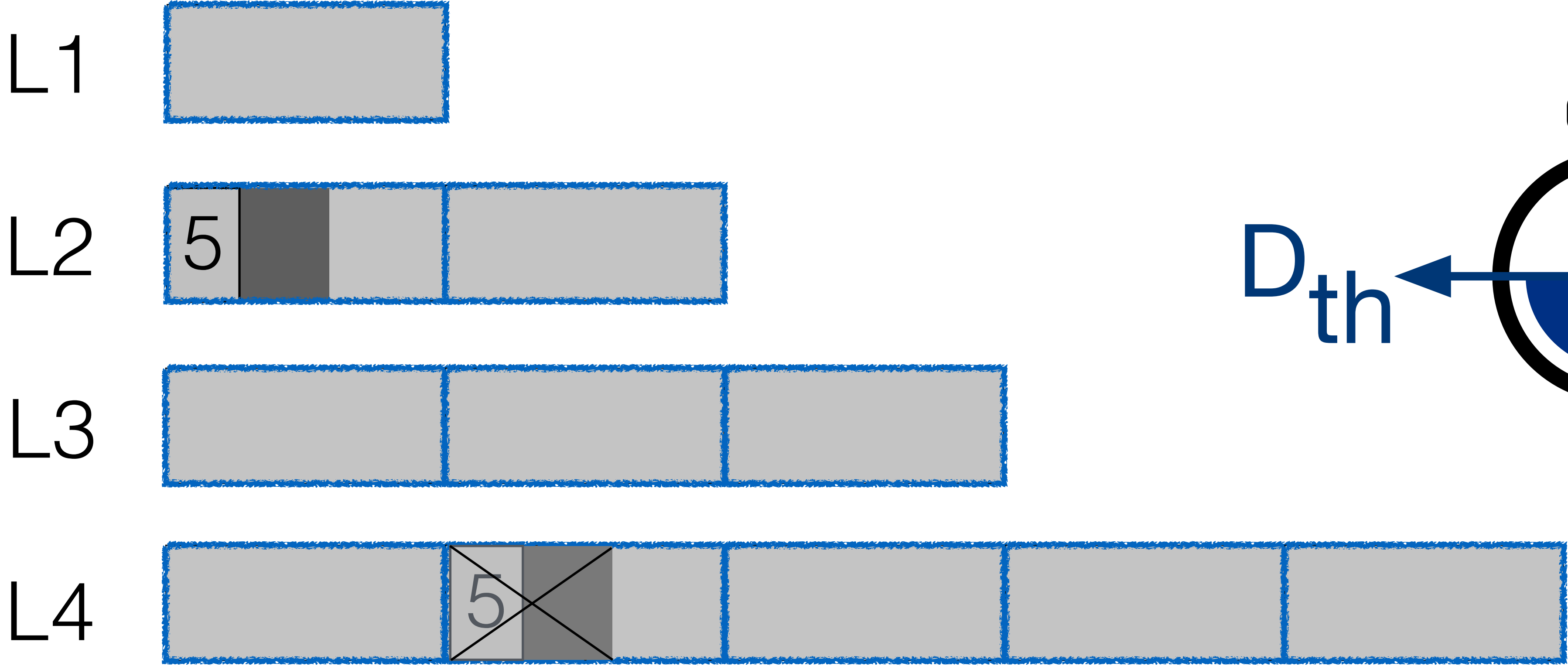
delete(5) within a threshold time:  $D_{th}$





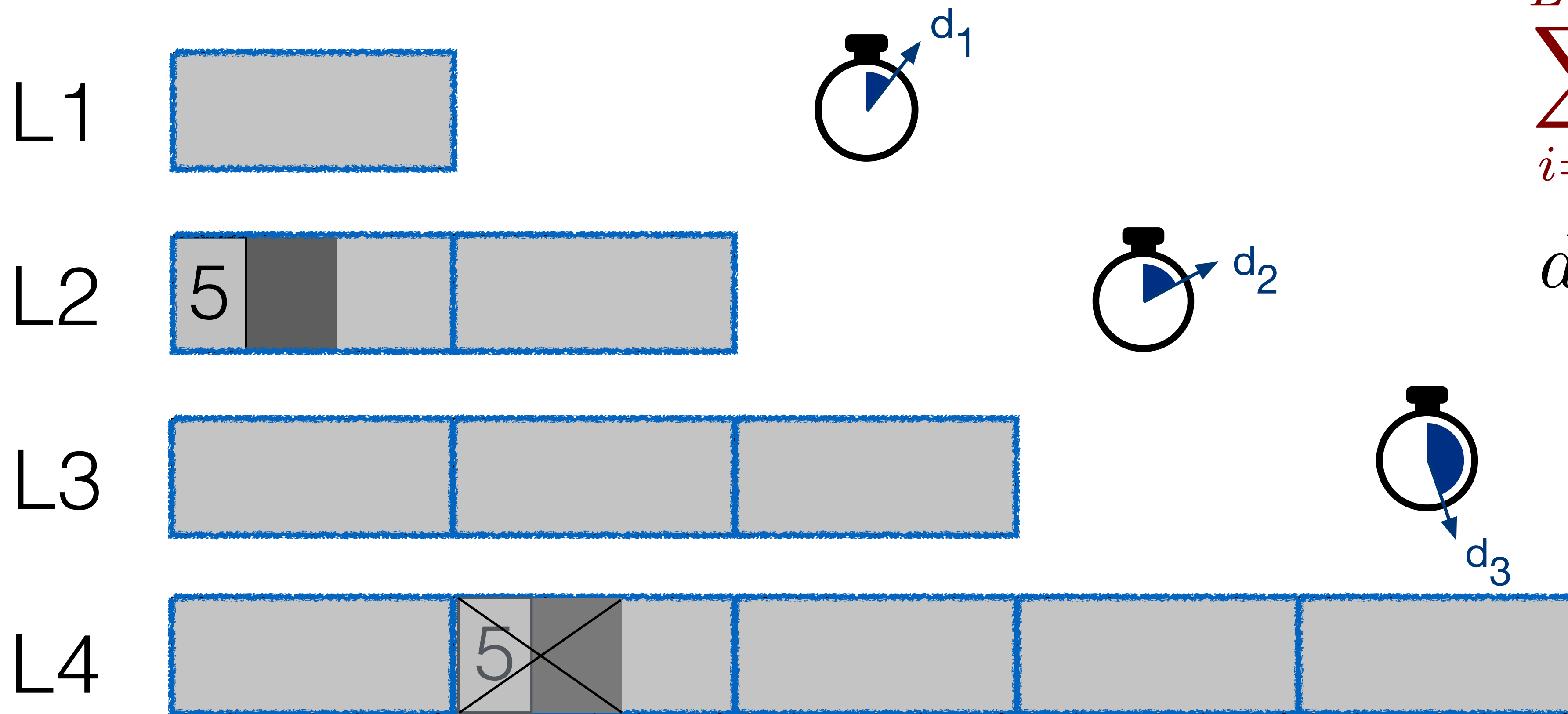
# FAst DElete

delete(5) within a threshold time:  $D_{th}$



# FAst DElete

delete(5) within a threshold time:  $D_{th}$

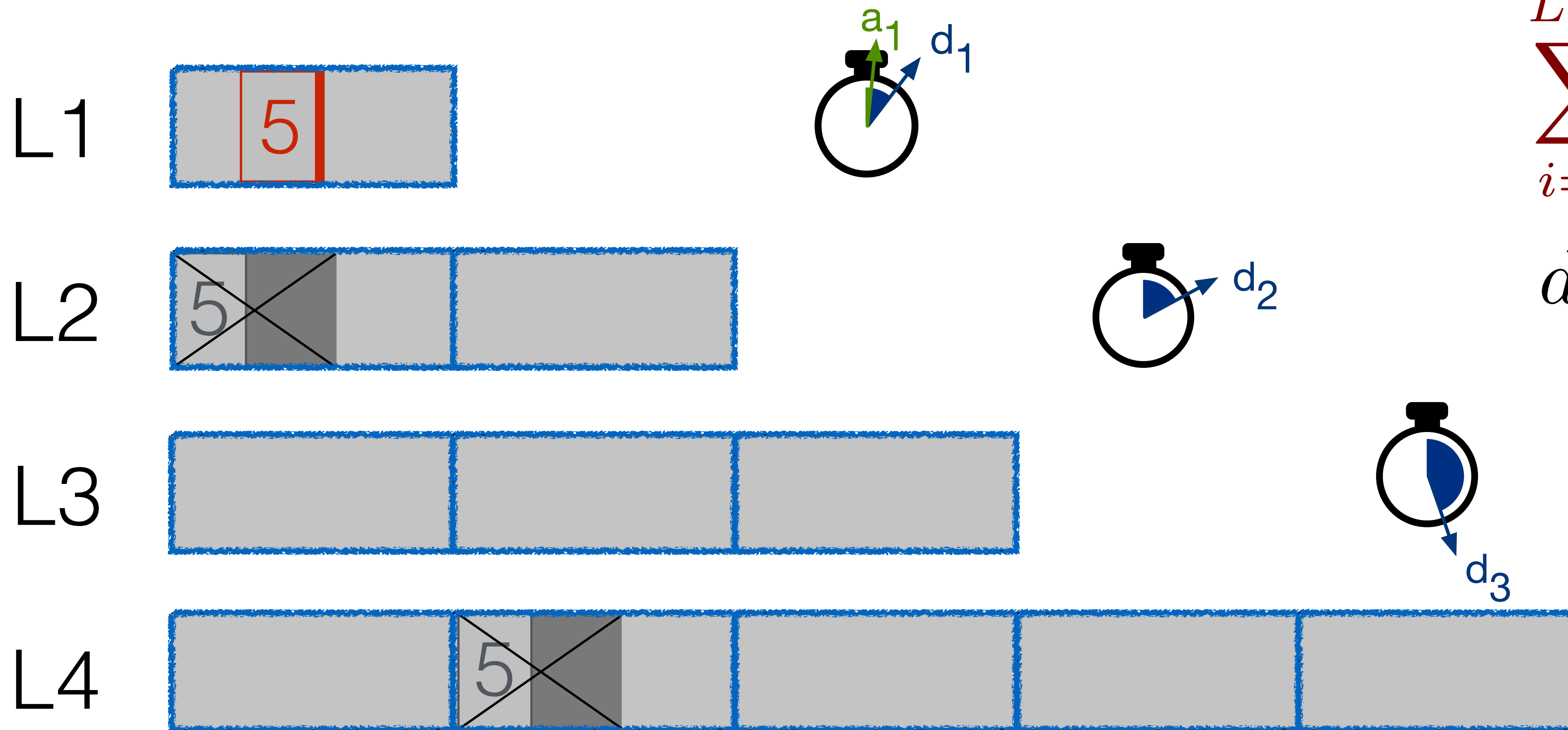


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

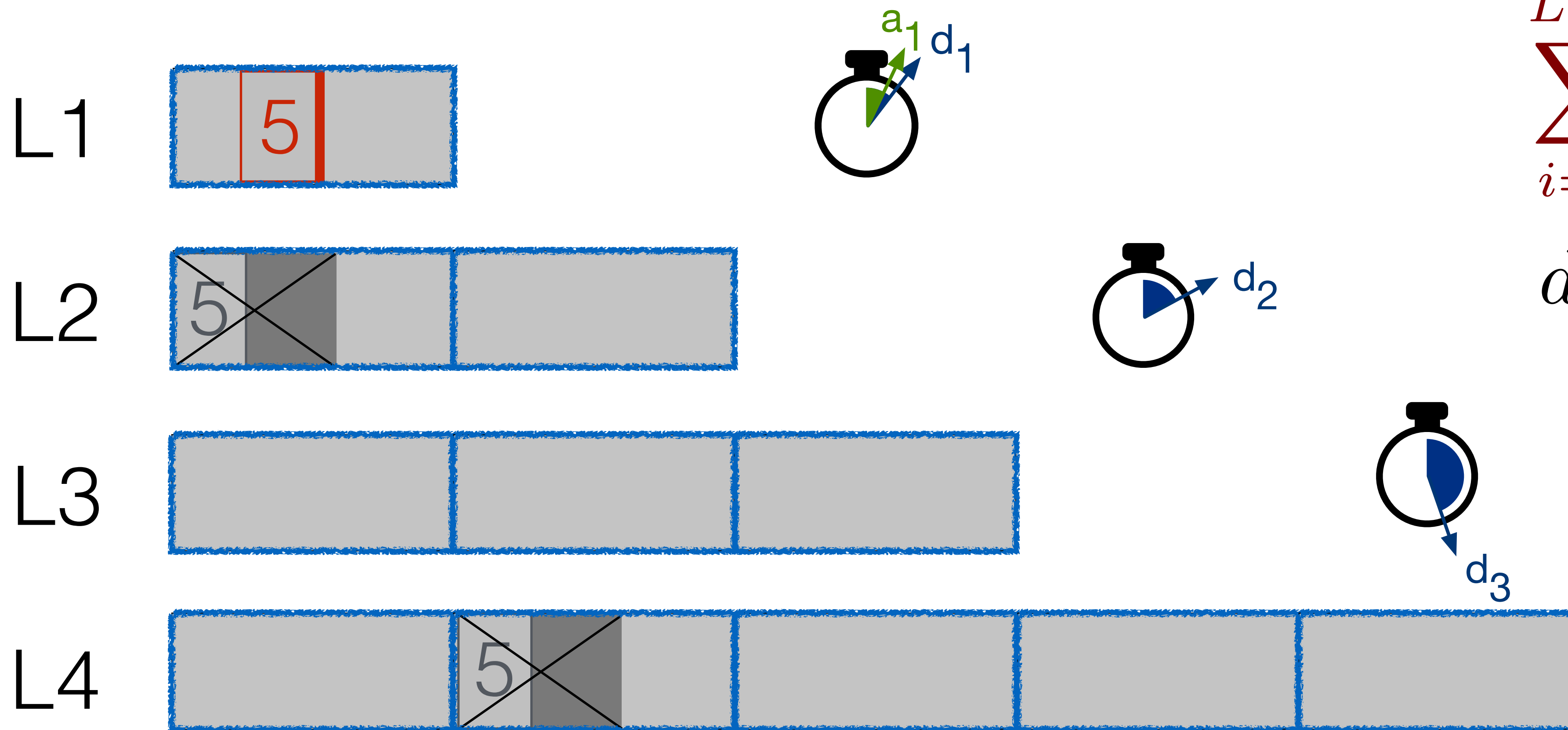


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

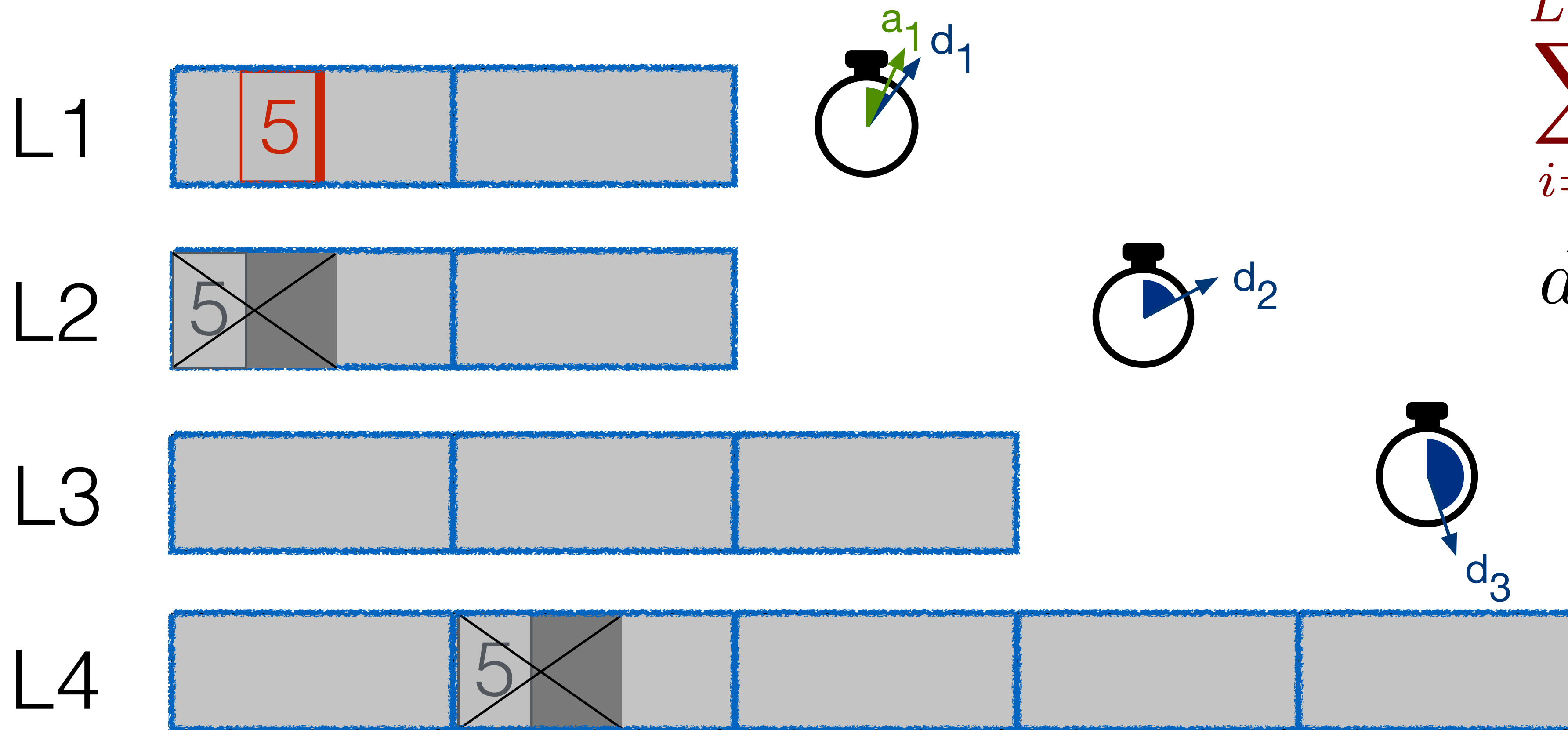


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

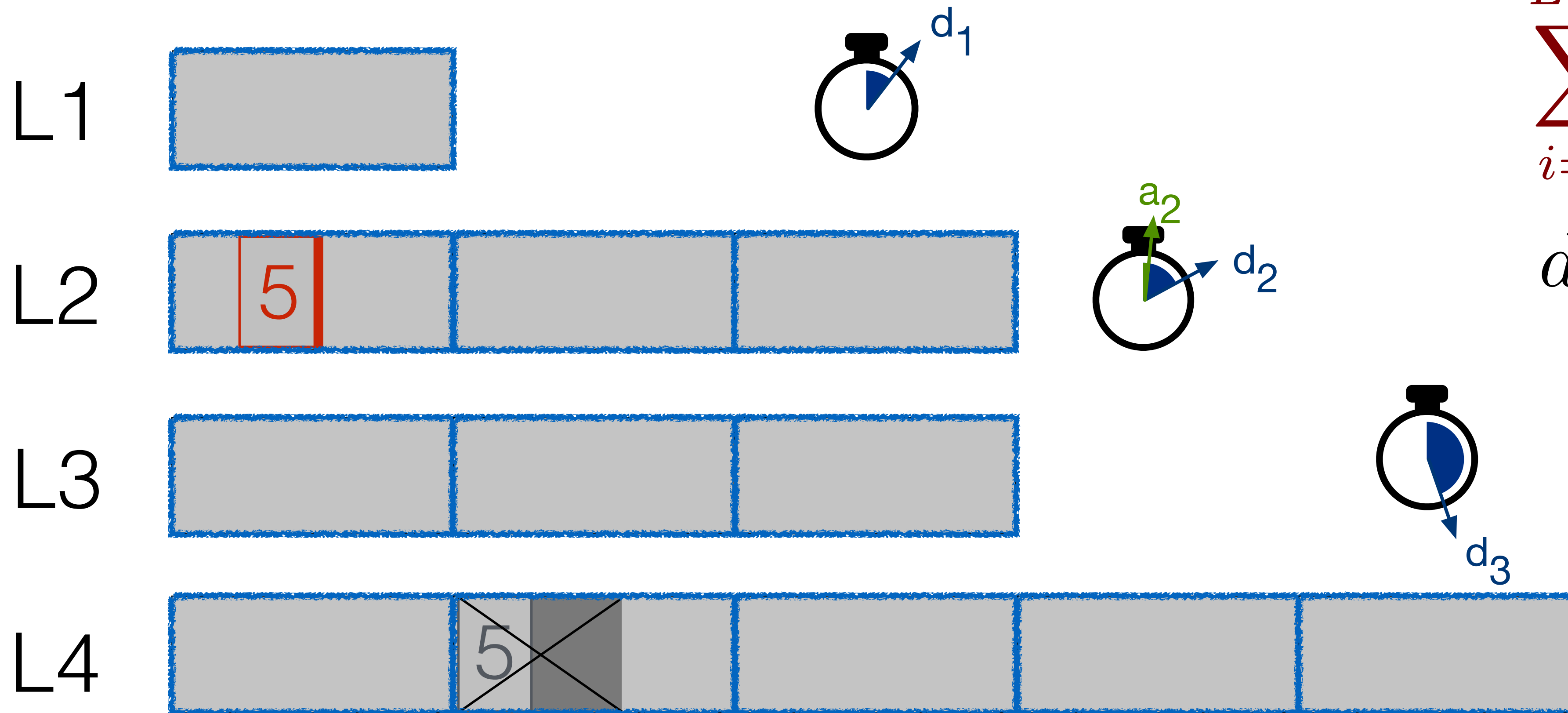


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

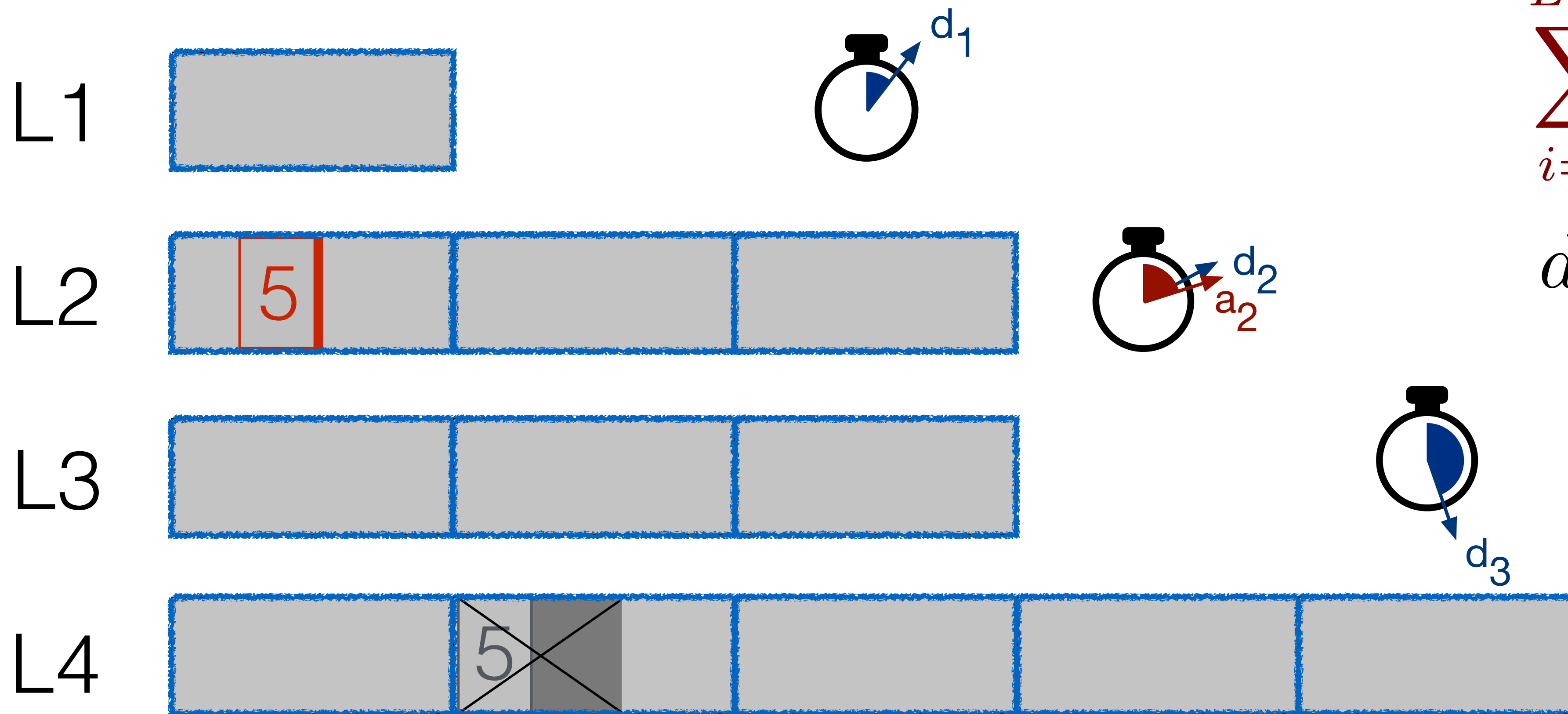


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

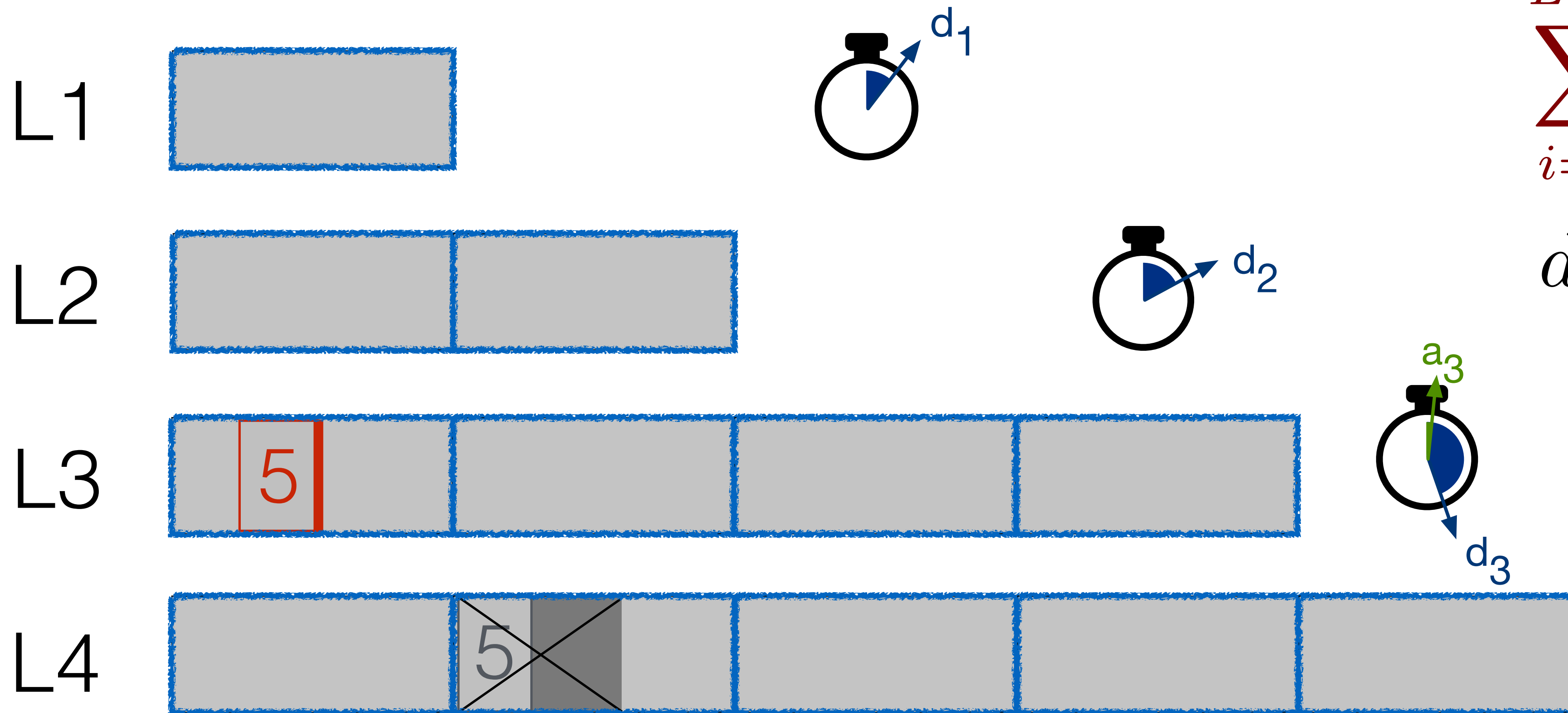


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$



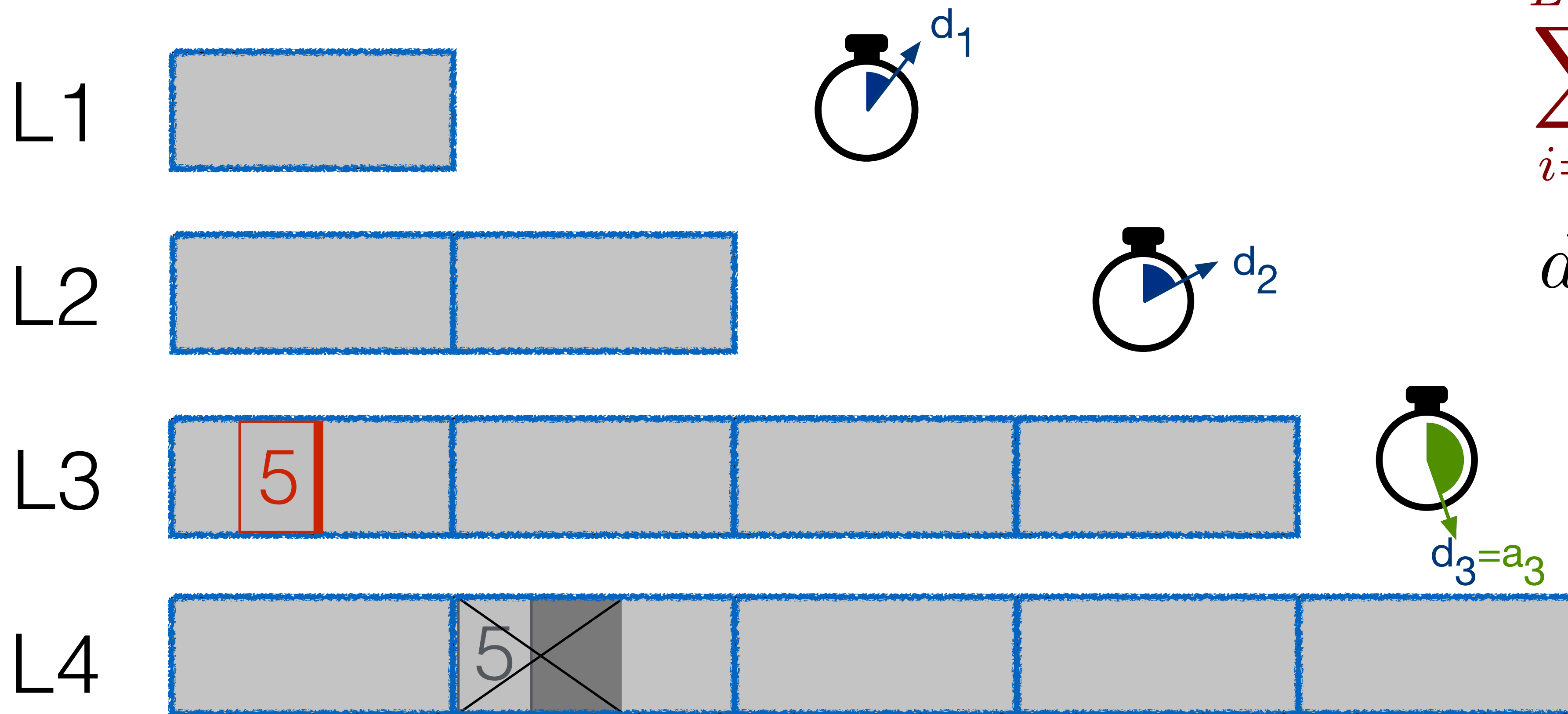
$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$



# FAst DElete

delete(5) within a threshold time:  $D_{th}$

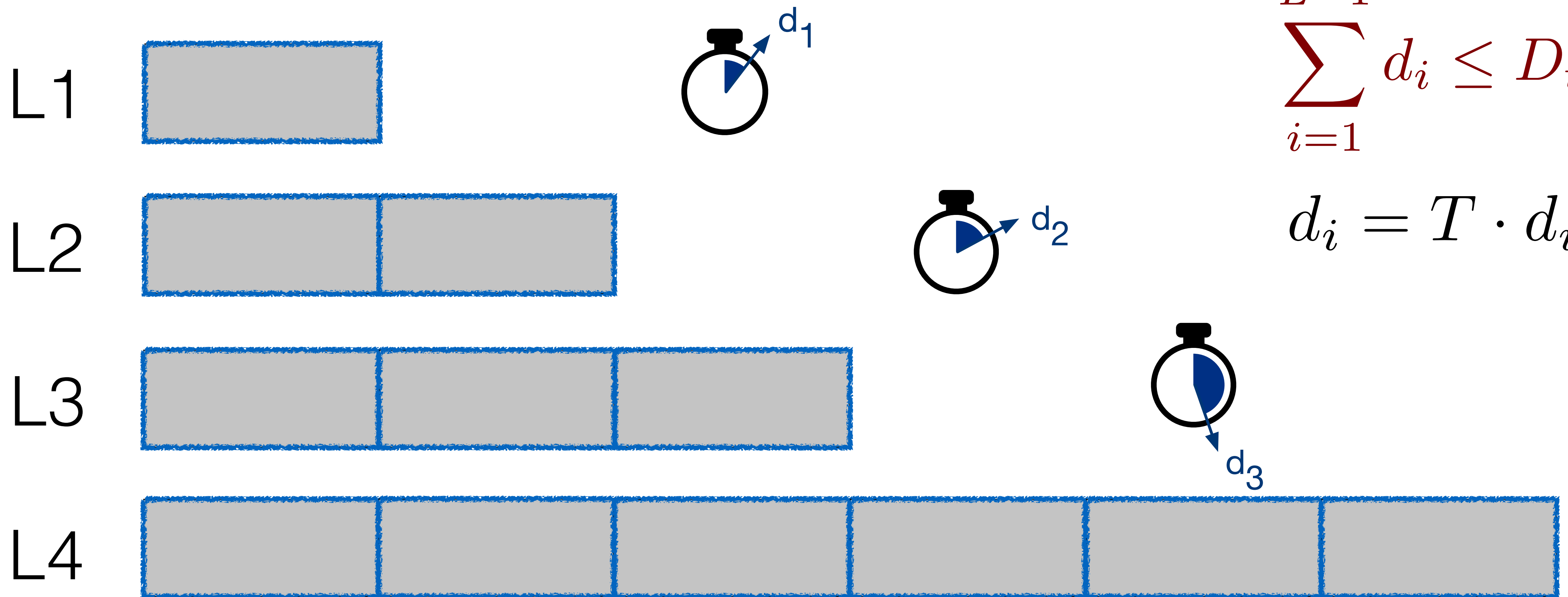


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

# FAst DElete

delete(5) within a threshold time:  $D_{th}$

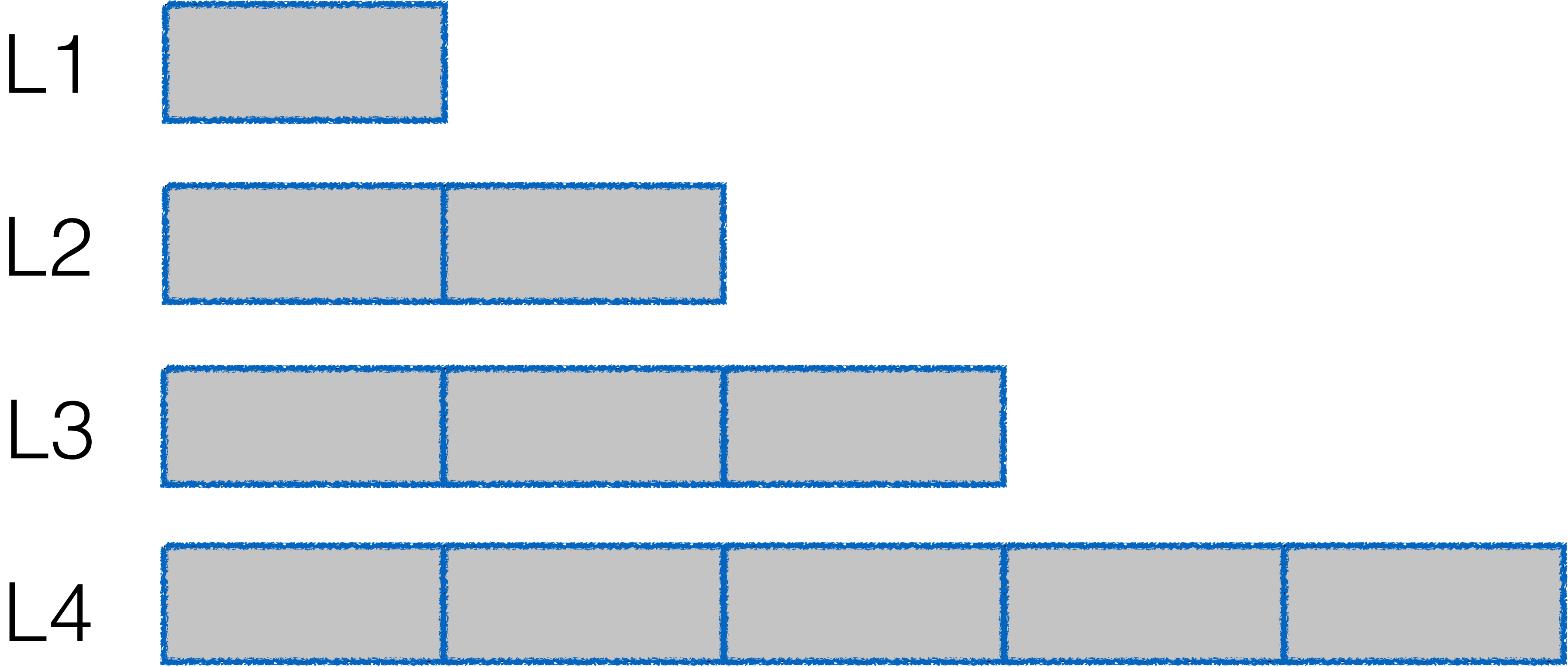


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

$$d_i = T \cdot d_{i-1}$$

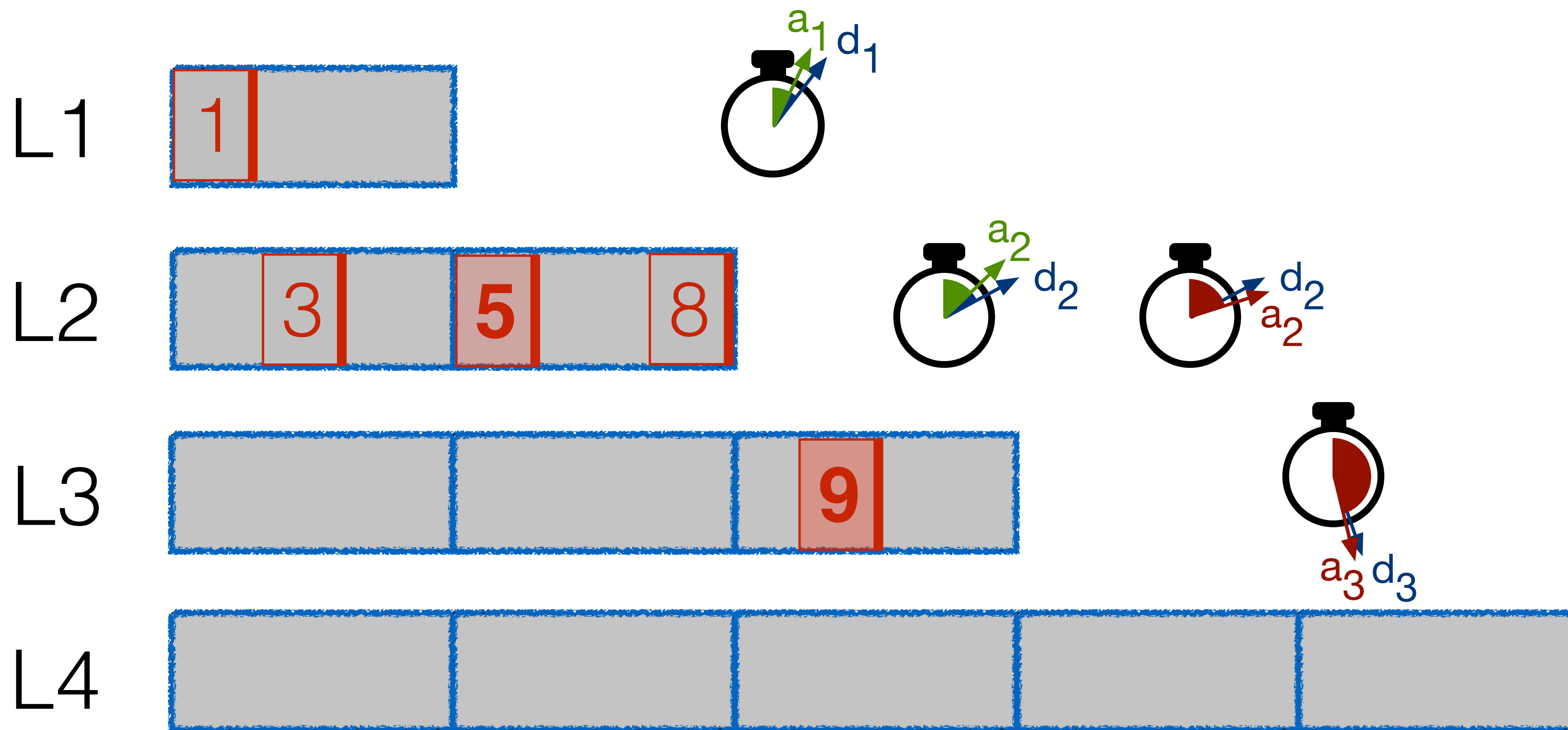
# FAst DElete

breaking ties in practical workloads



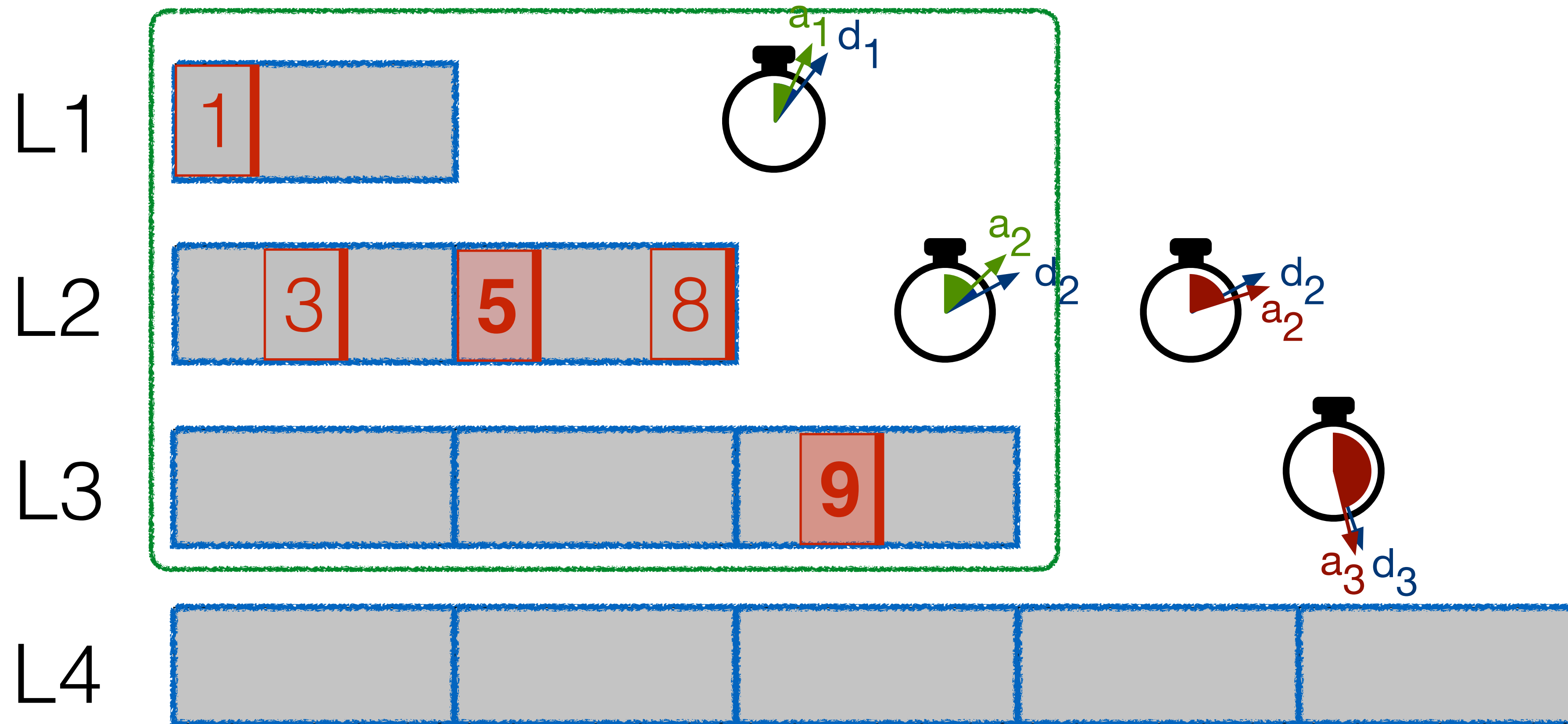
# FAst DElete

breaking ties in practical workloads



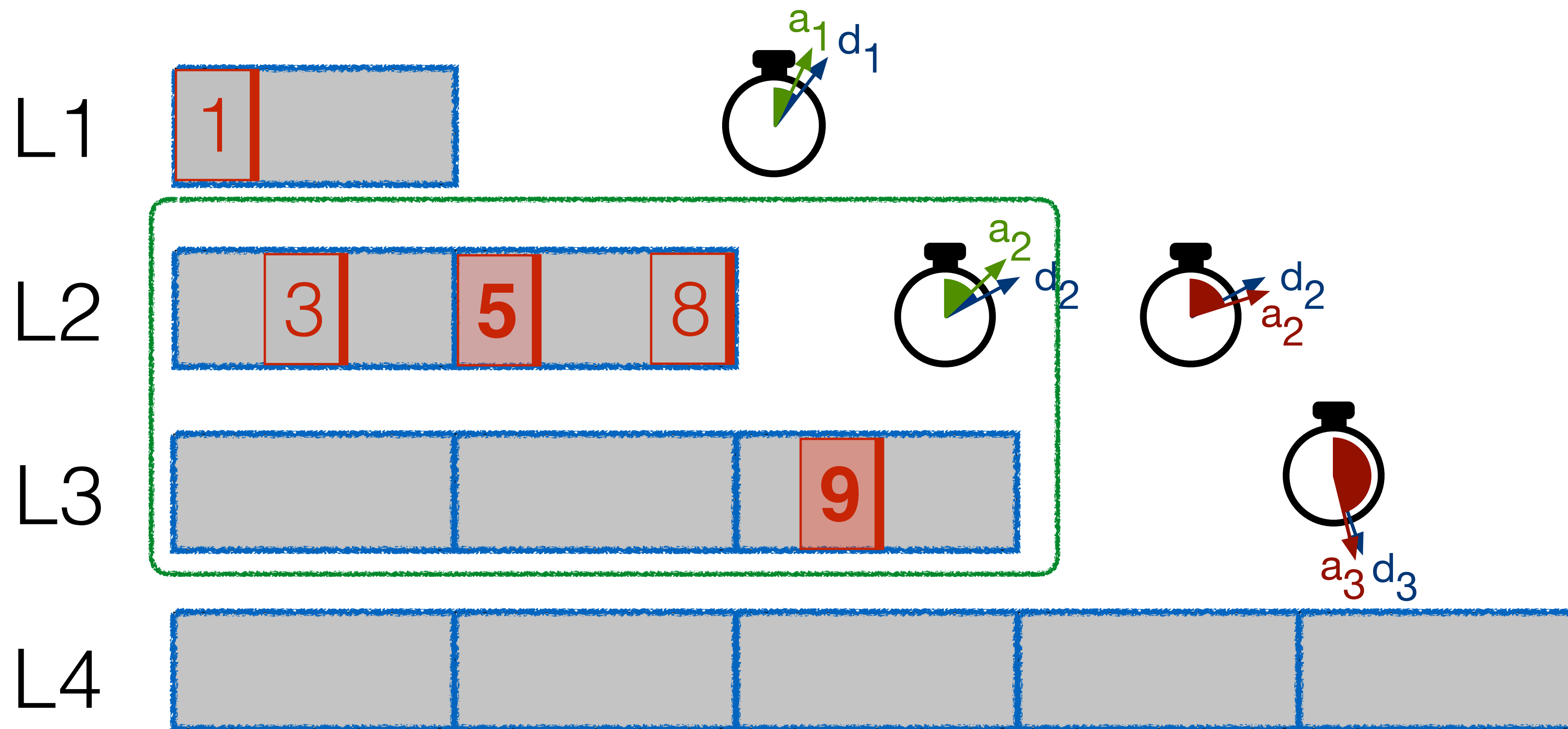
# FAst DElete

breaking ties in practical workloads



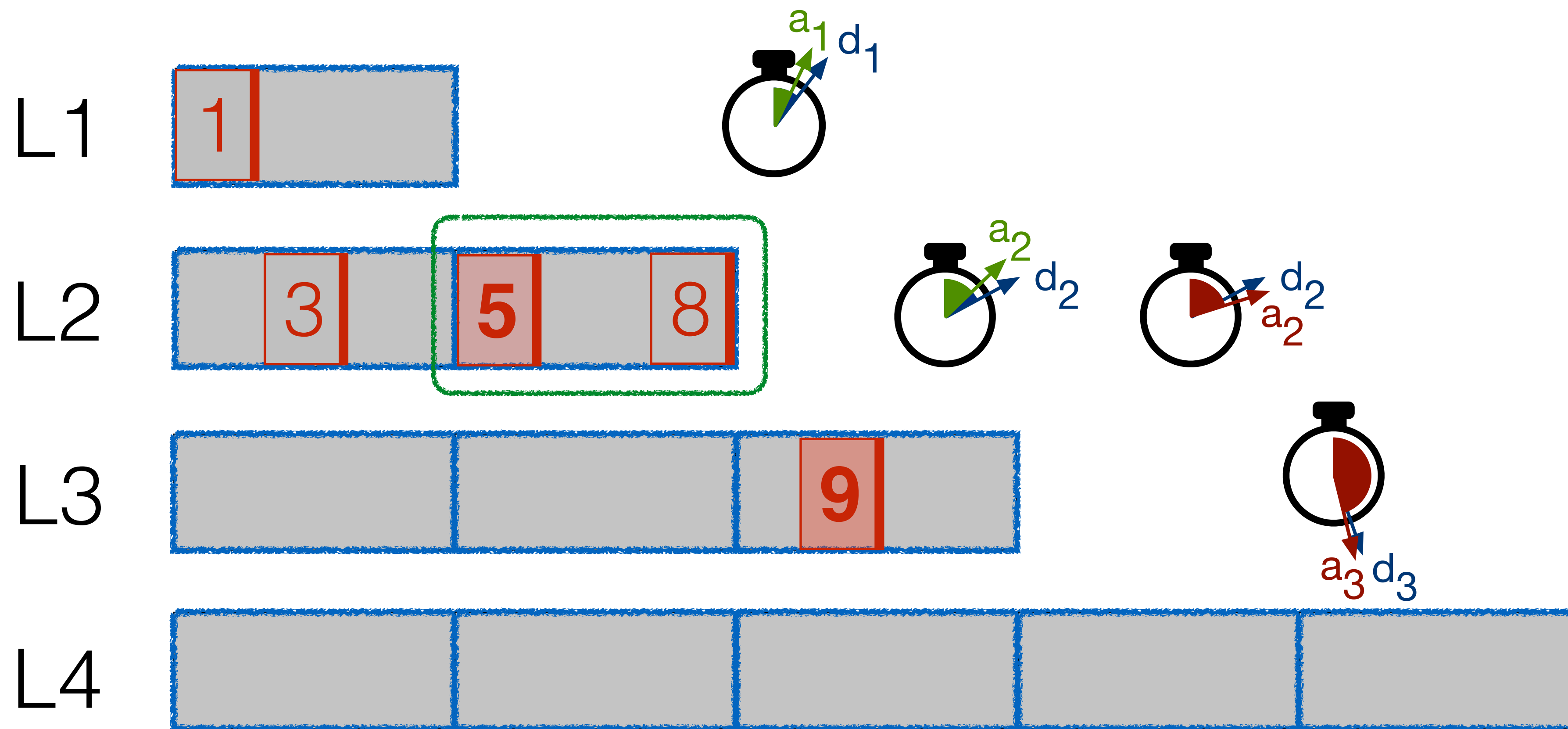
# FAst DElete

breaking ties in practical workloads



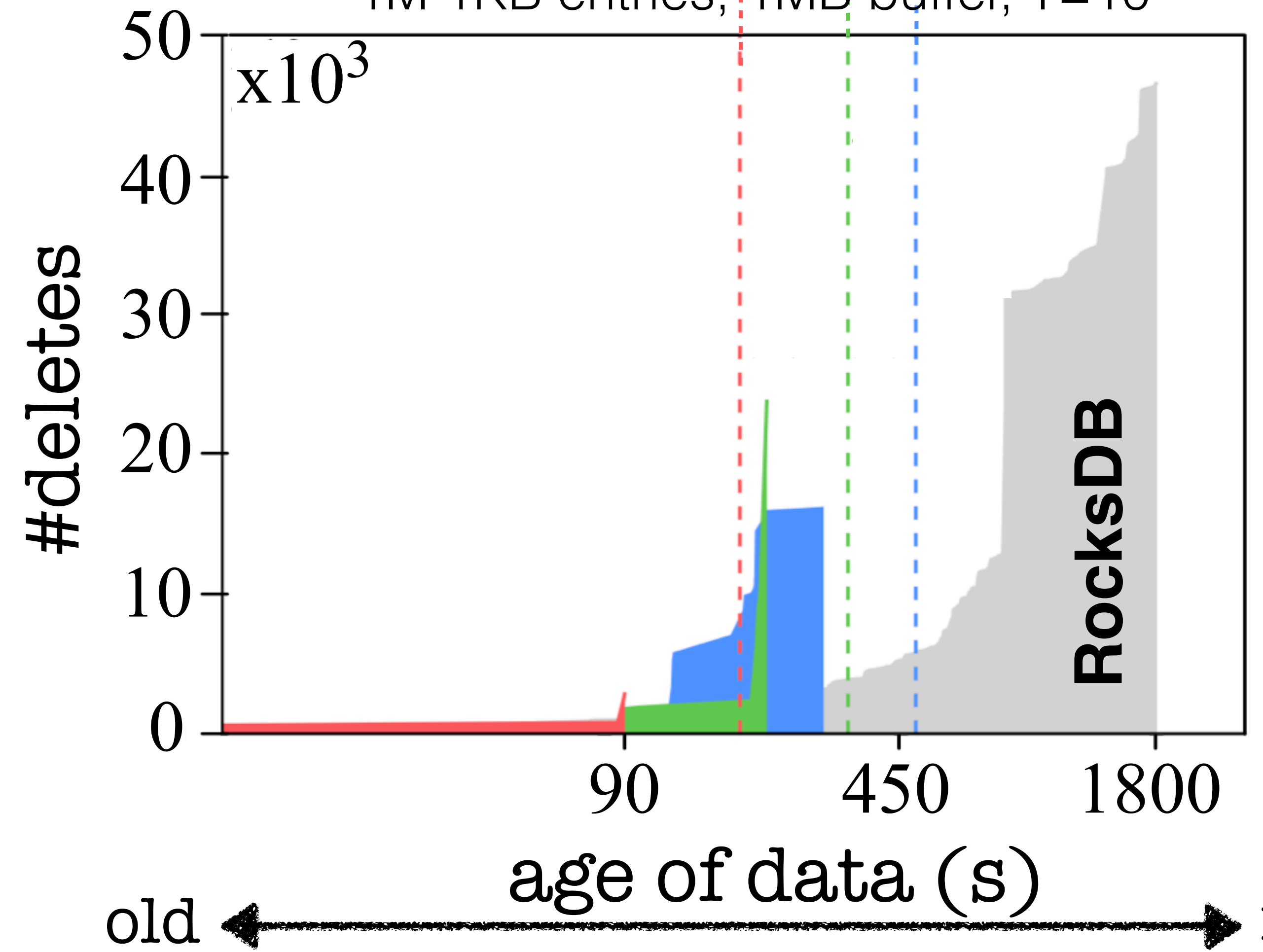
# FAst DElete

breaking ties in practical workloads



persist all deletes within: 300s  
150s 600s

1M 1KB entries, 1MB buffer, T=10



persists deletes timely

within threshold

old ← → new



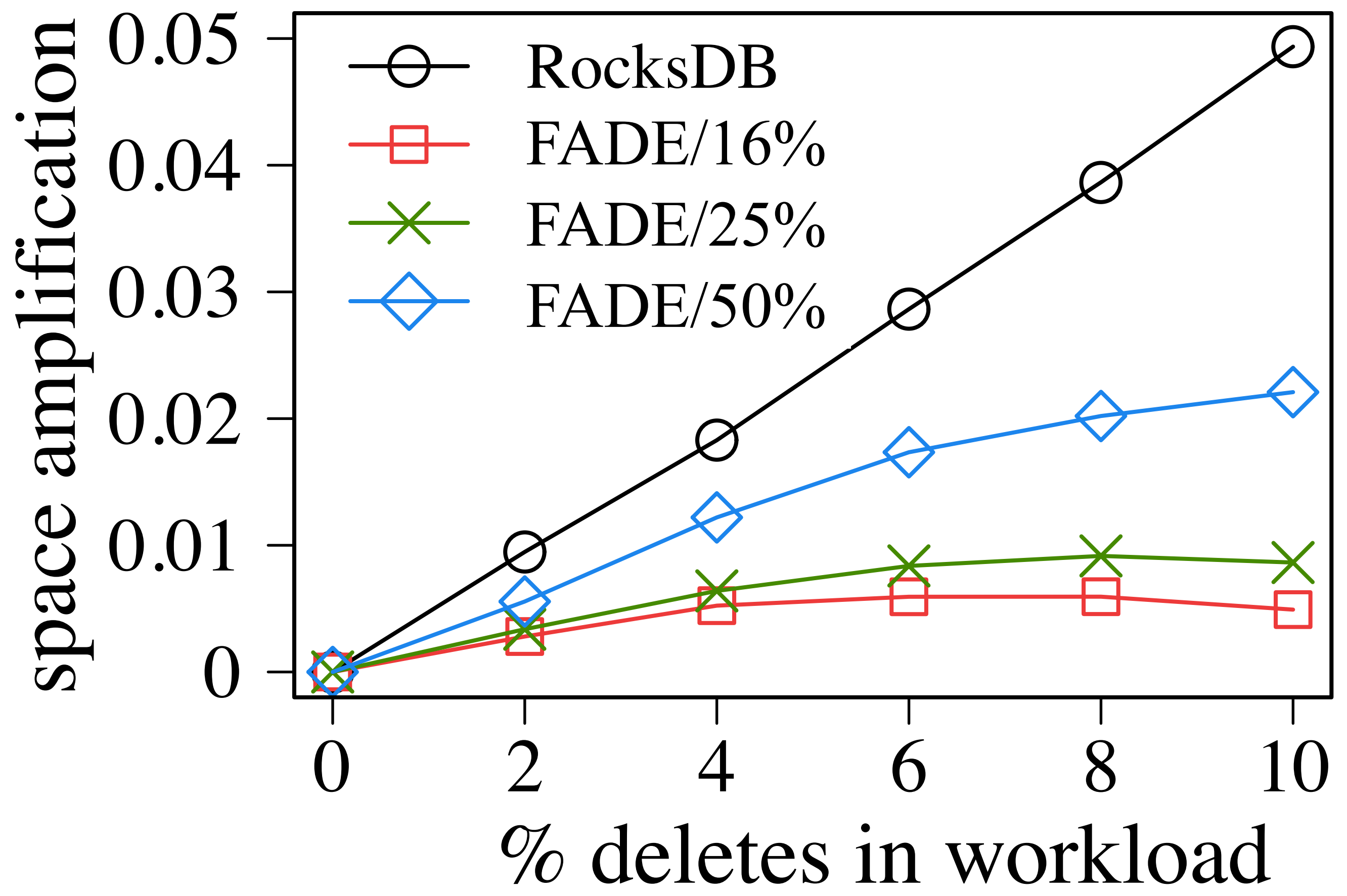
reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

1M 1KB entries, 1MB buffer, T=10



improved read performance

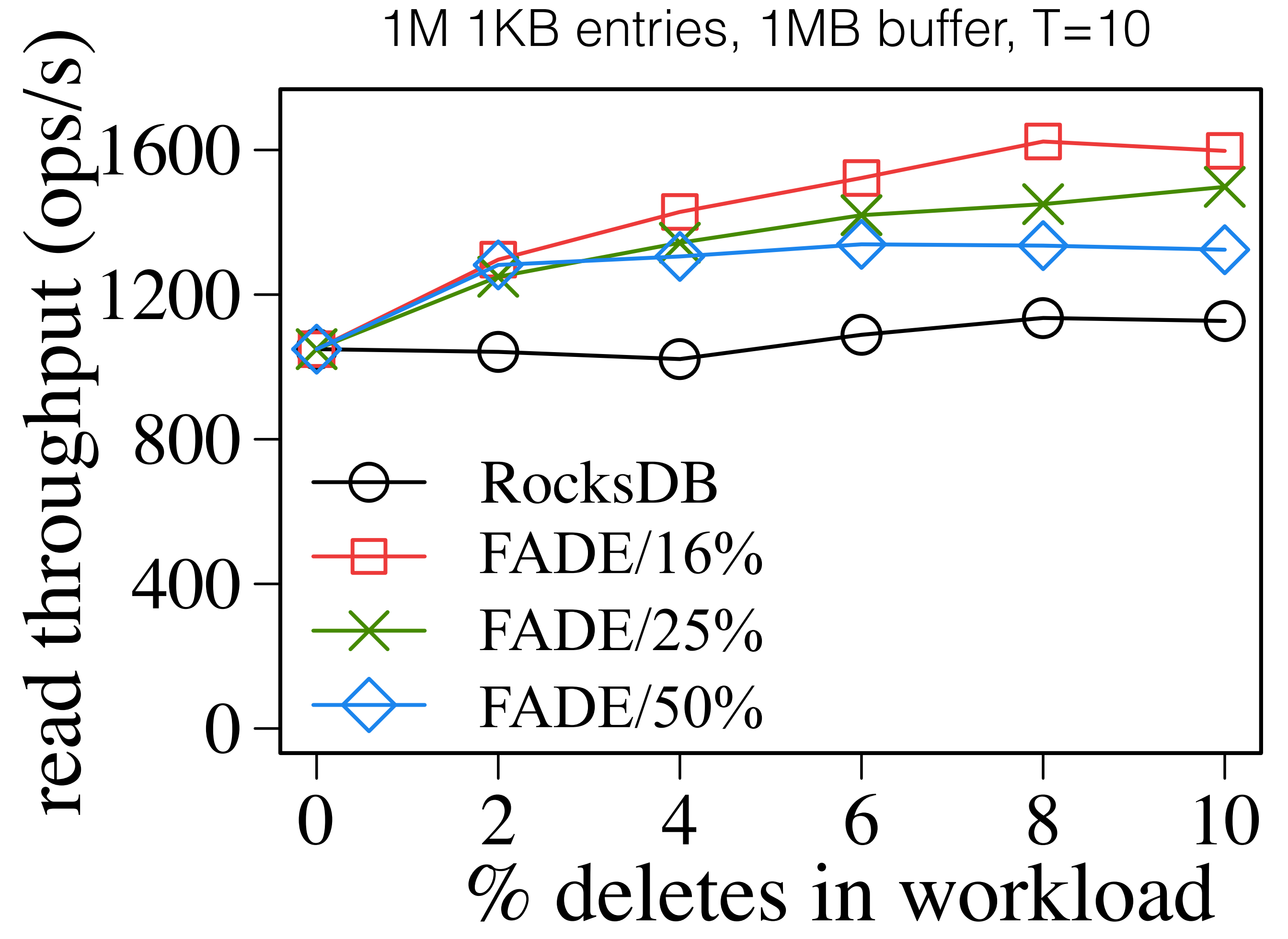
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

4% - 25%

improved read performance

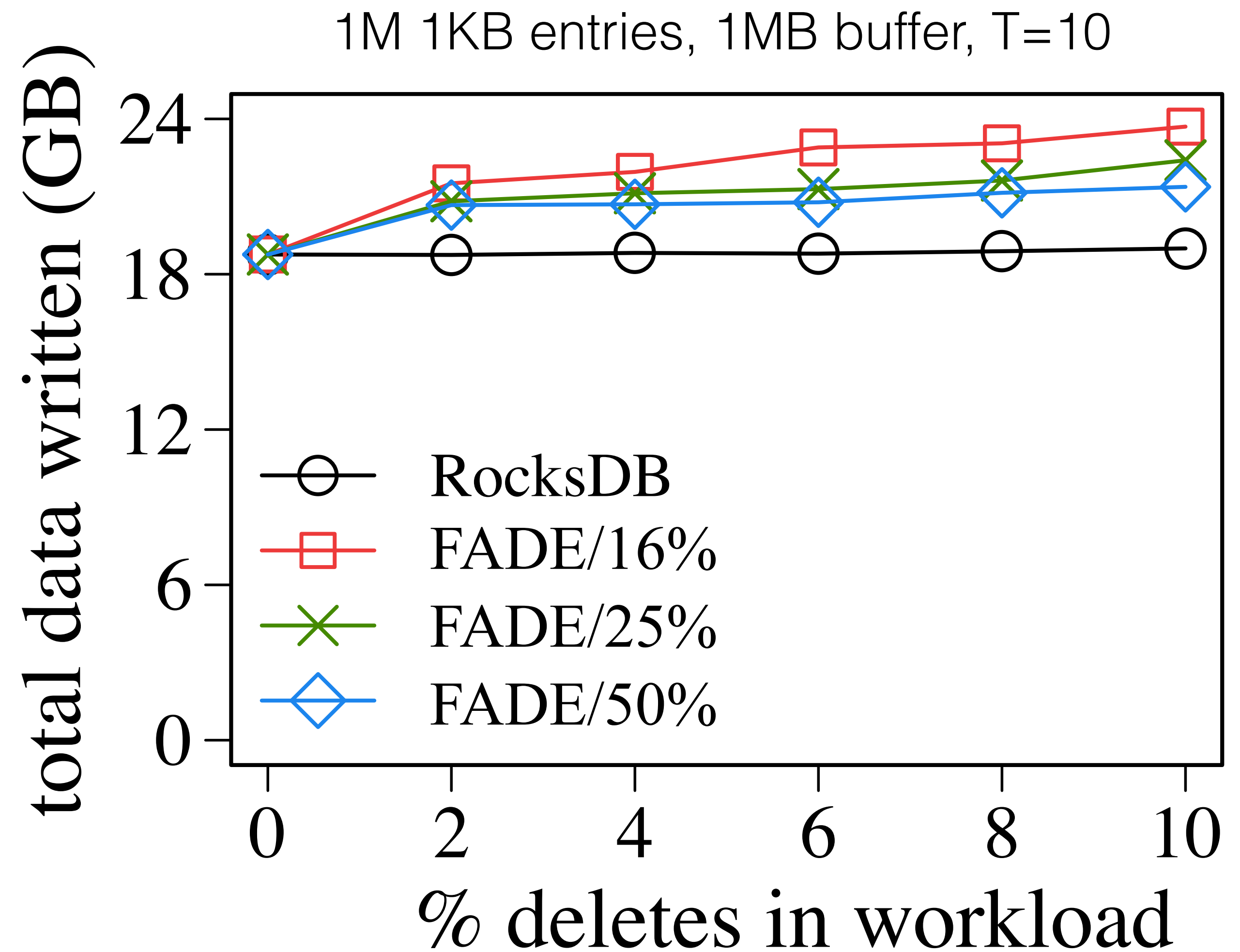
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

4% - 25%

improved read performance

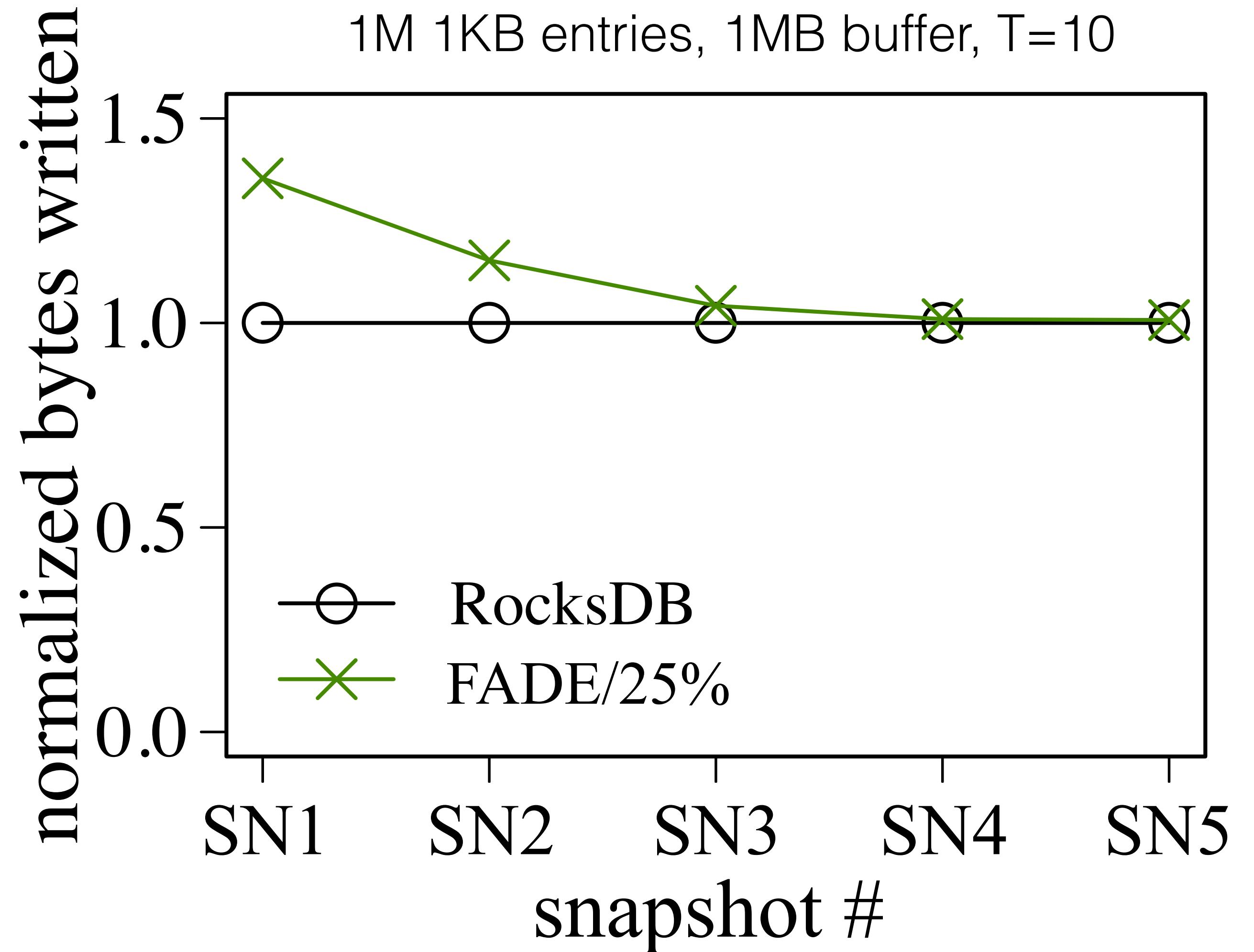
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

0.7%

improved read performance

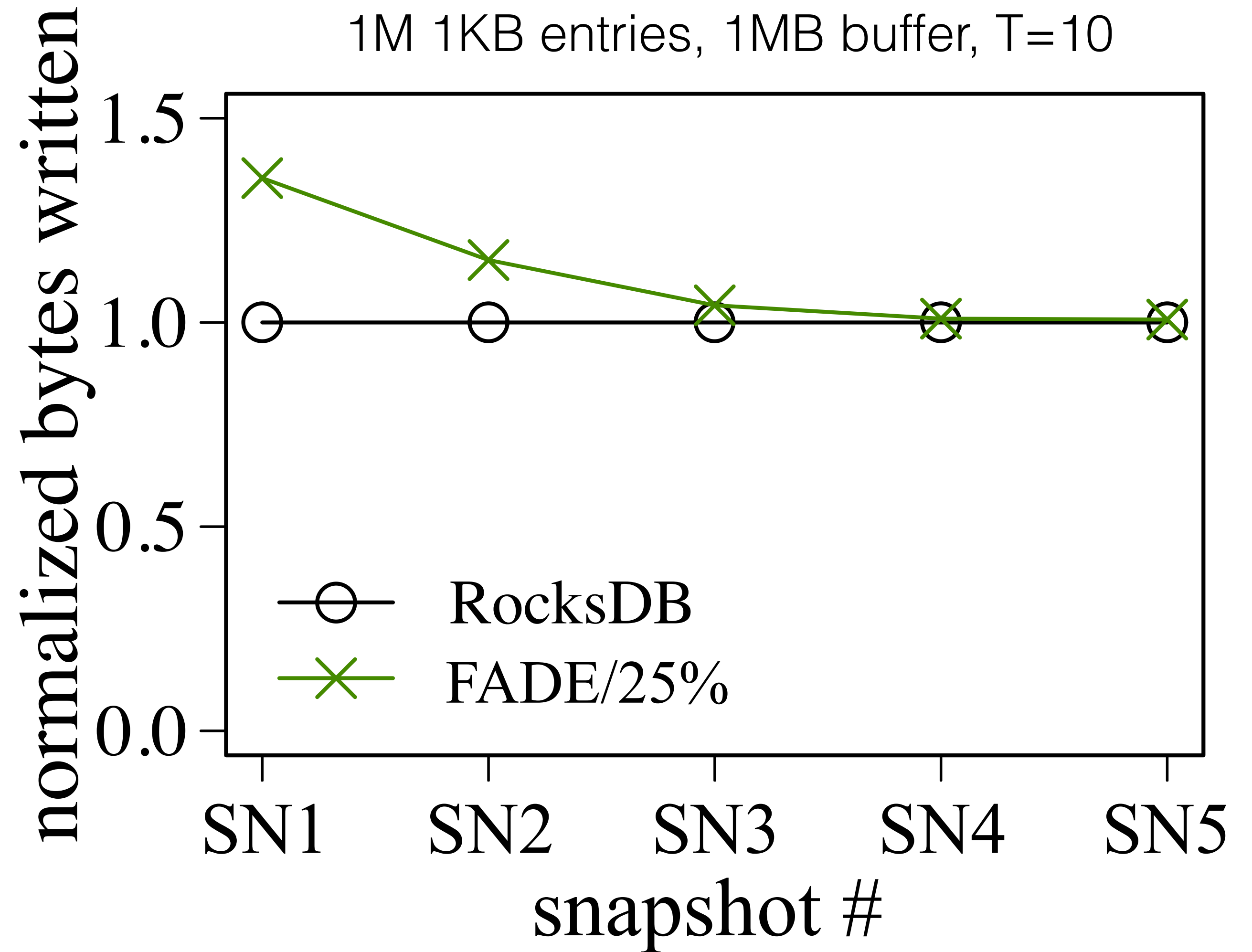
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



# the problems

poor read perf.

write amplification

space amplification

# FADE

unbounded delete  
persistence latency

$$\sum_{i=1}^{L-1} t_i$$



$D_{th}$



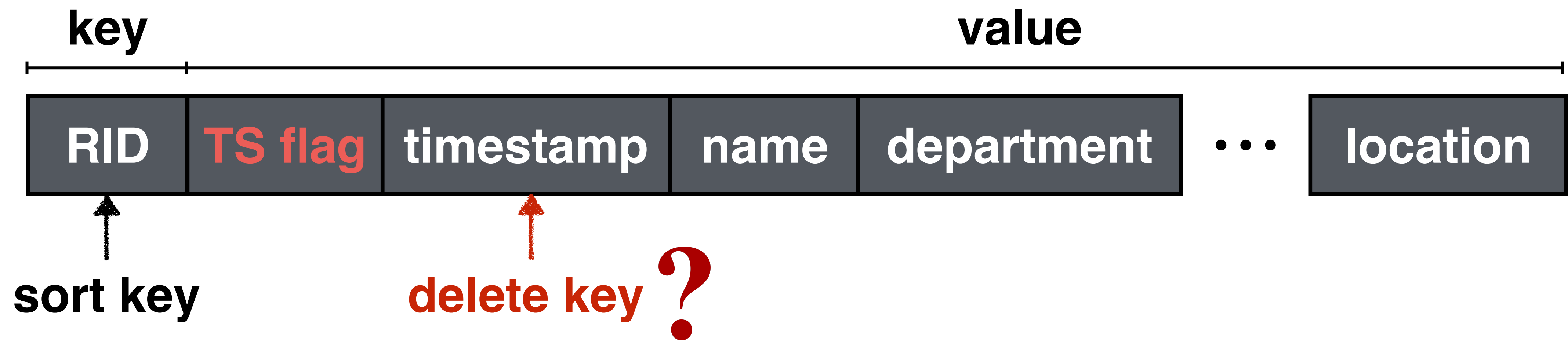


**deletes on a secondary attribute**



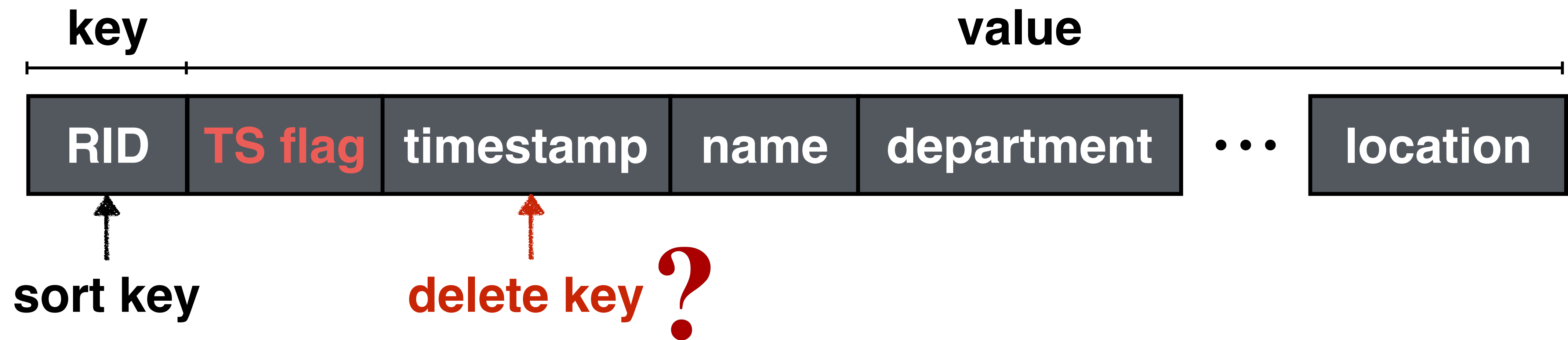
deletes on a secondary attribute

delete all entries older than: **D days**



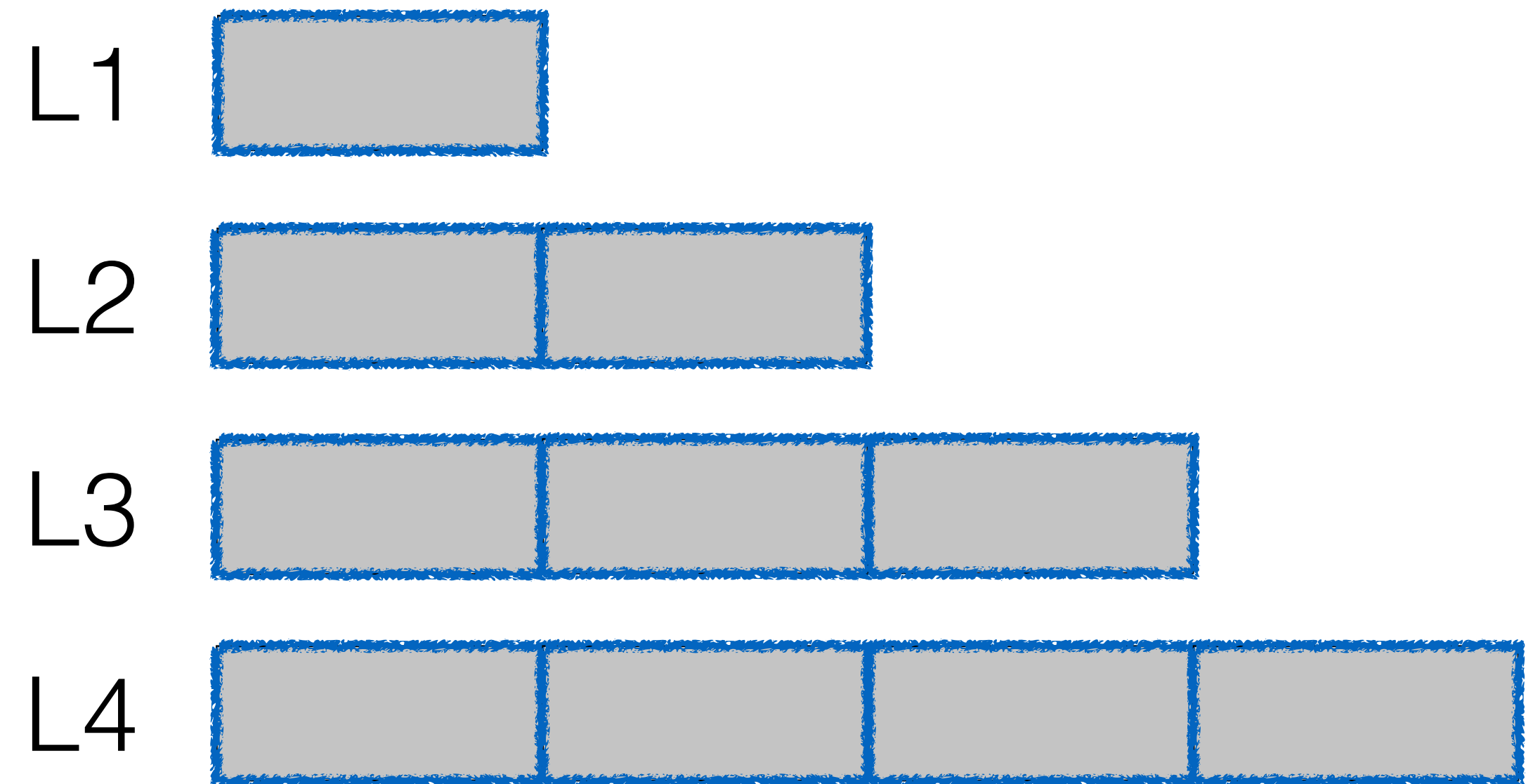
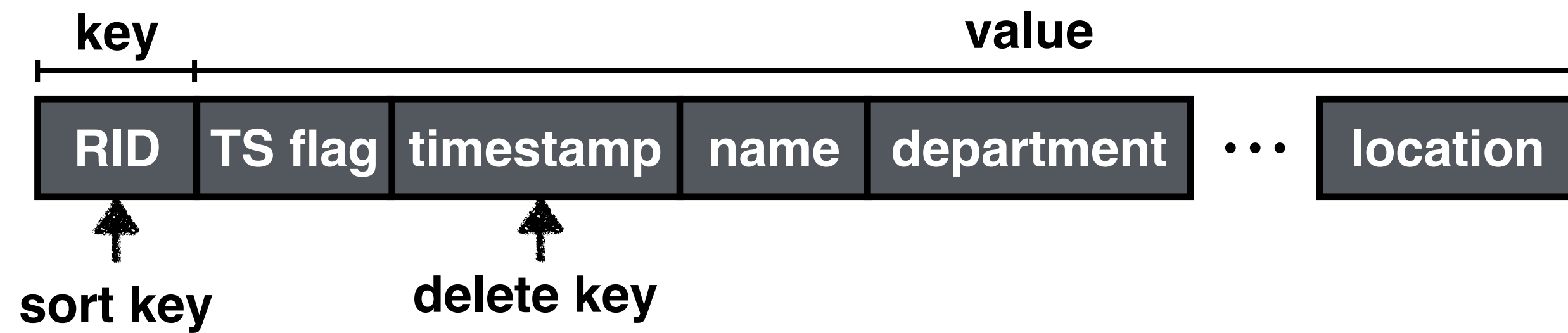
deletes on a secondary attribute

delete all entries older than: **D days**



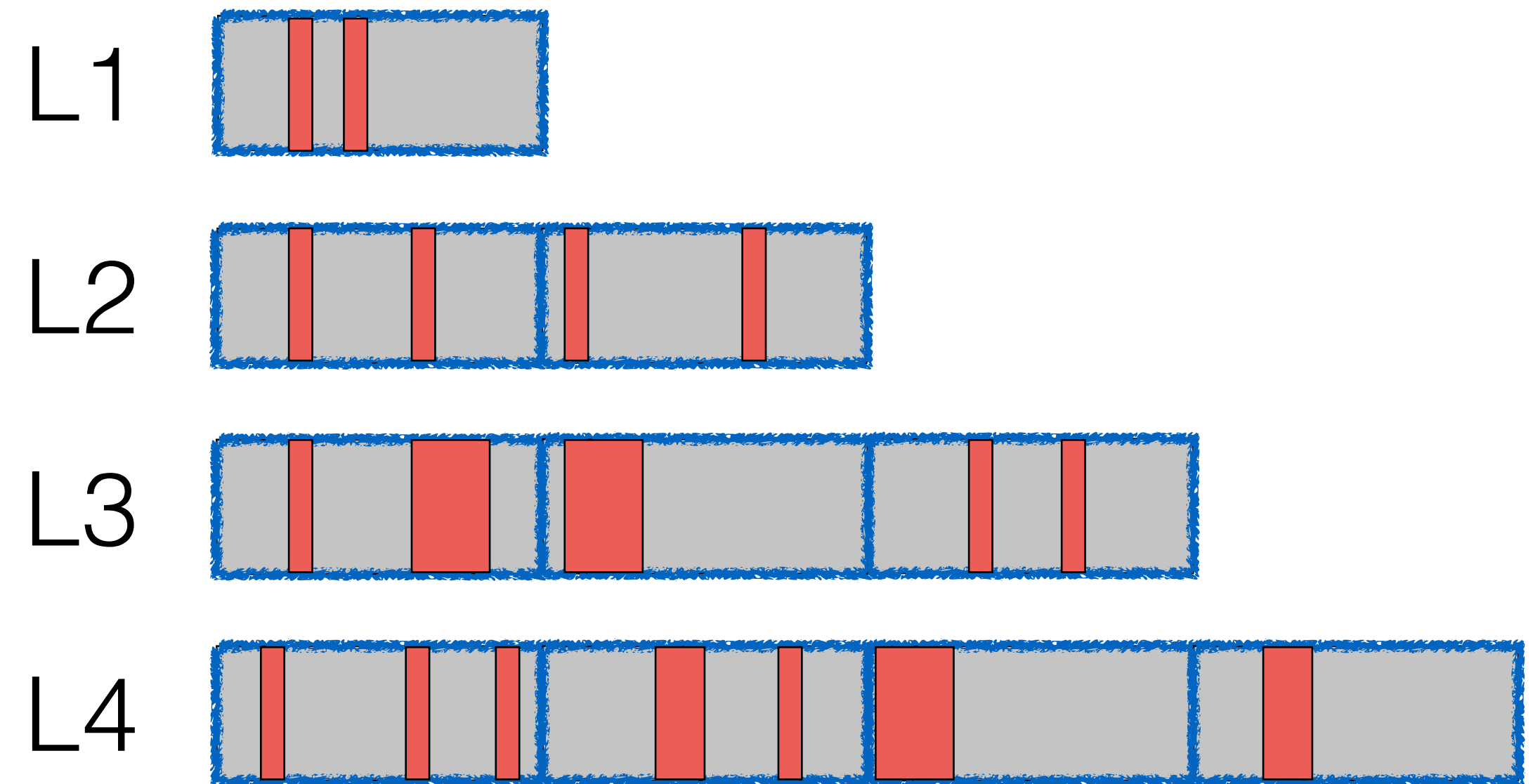
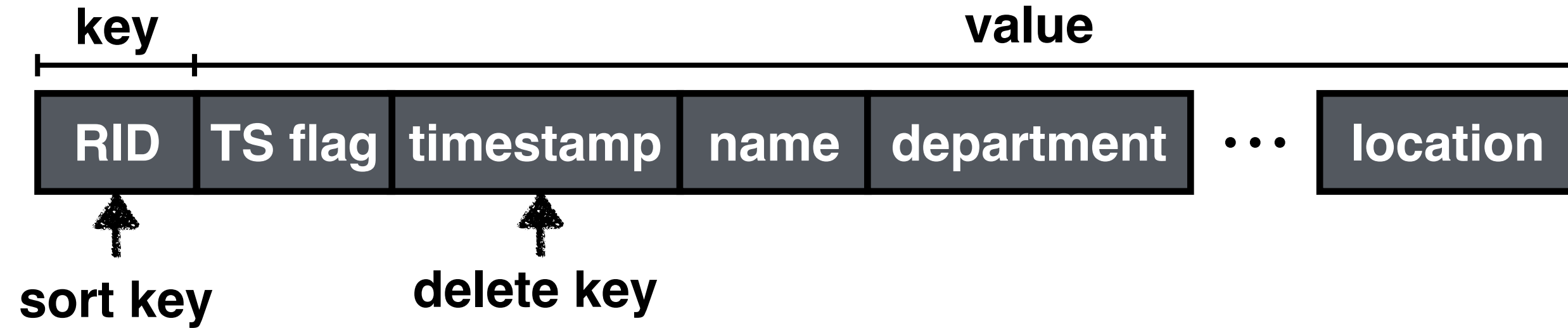
deletes on a secondary attribute

delete all entries older than: **D days**



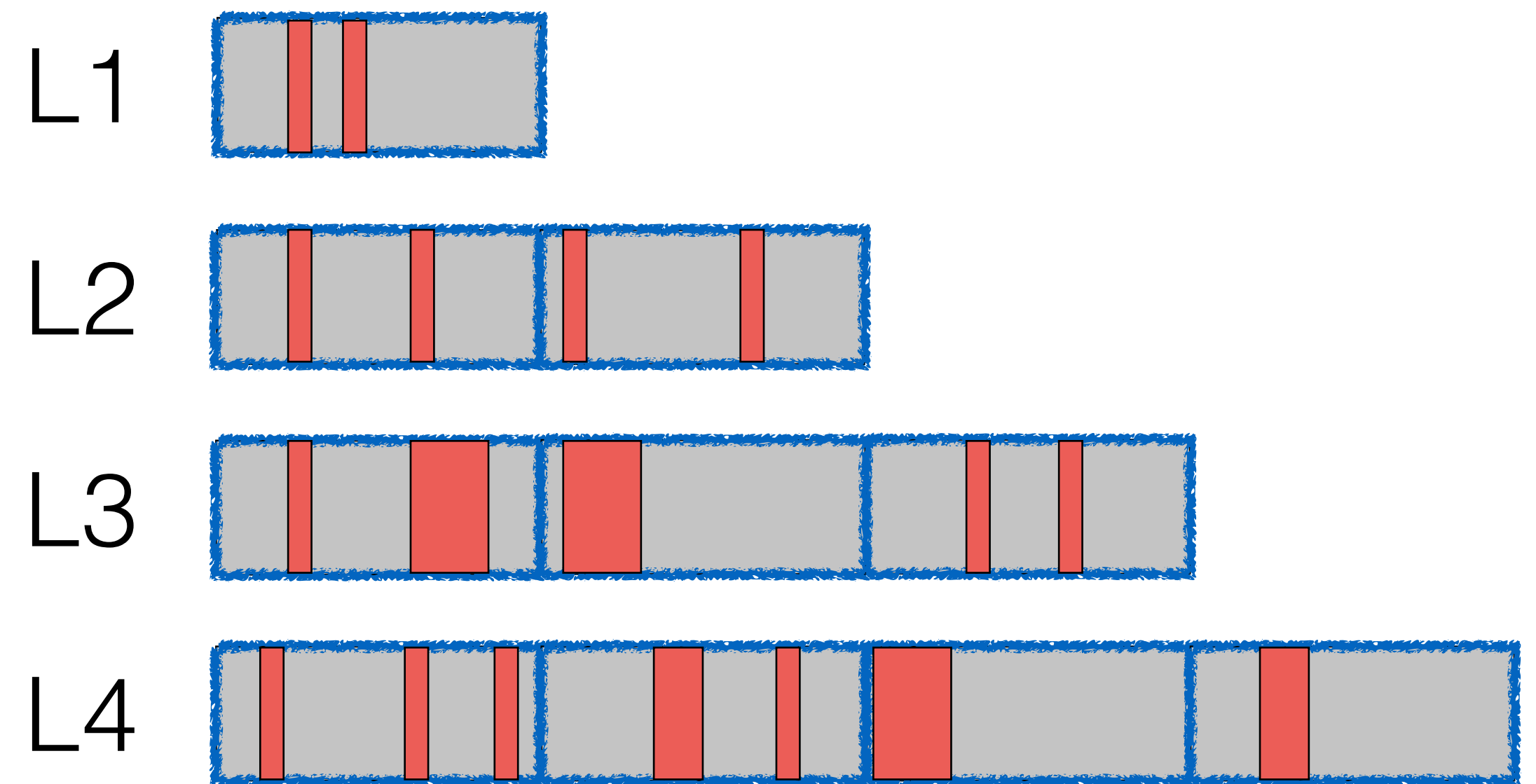
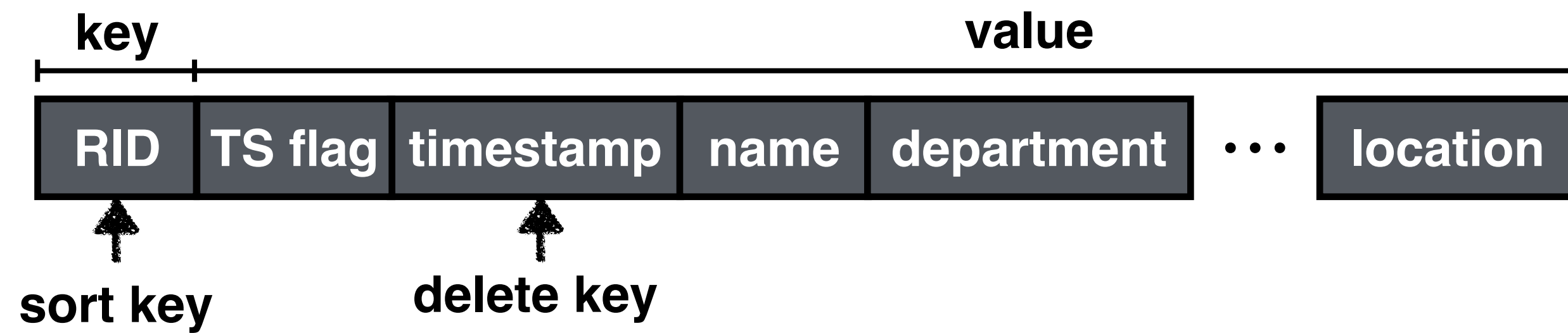
deletes on a secondary attribute

delete all entries older than: **D days**



deletes on a secondary attribute

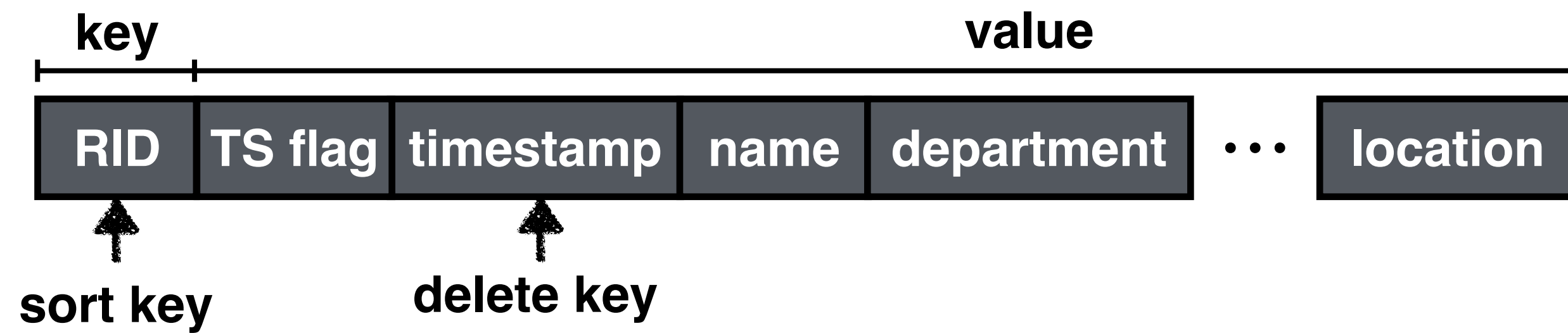
delete all entries older than: **D days**



**scattered occurrences**

deletes on a secondary attribute

delete all entries older than: **D days**



L1

L2

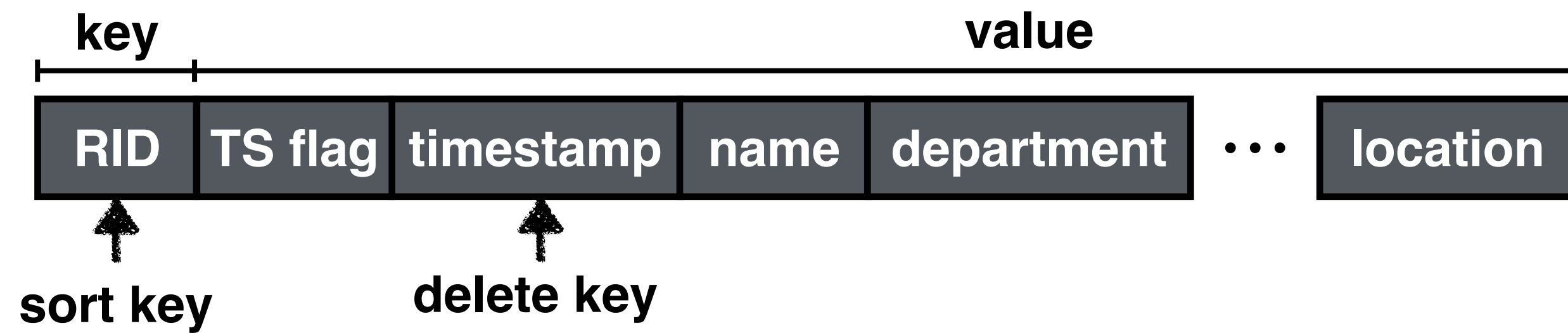
L3

L4



deletes on a secondary attribute

delete all entries older than: **D days**



L1

L2

L3

L4

latency spikes

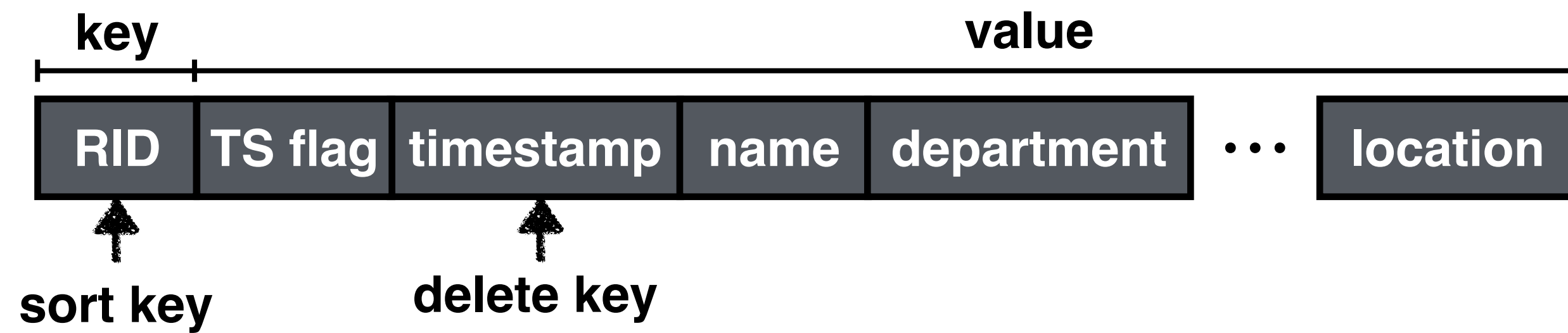


superfluous I/Os



deletes on a secondary attribute

delete all entries older than: **D days**



L1

L2

L3

L4

latency spikes



superfluous I/Os





# the problems

poor read perf.

write amplification

space amplification

# FADE

unbounded delete  
persistence latency



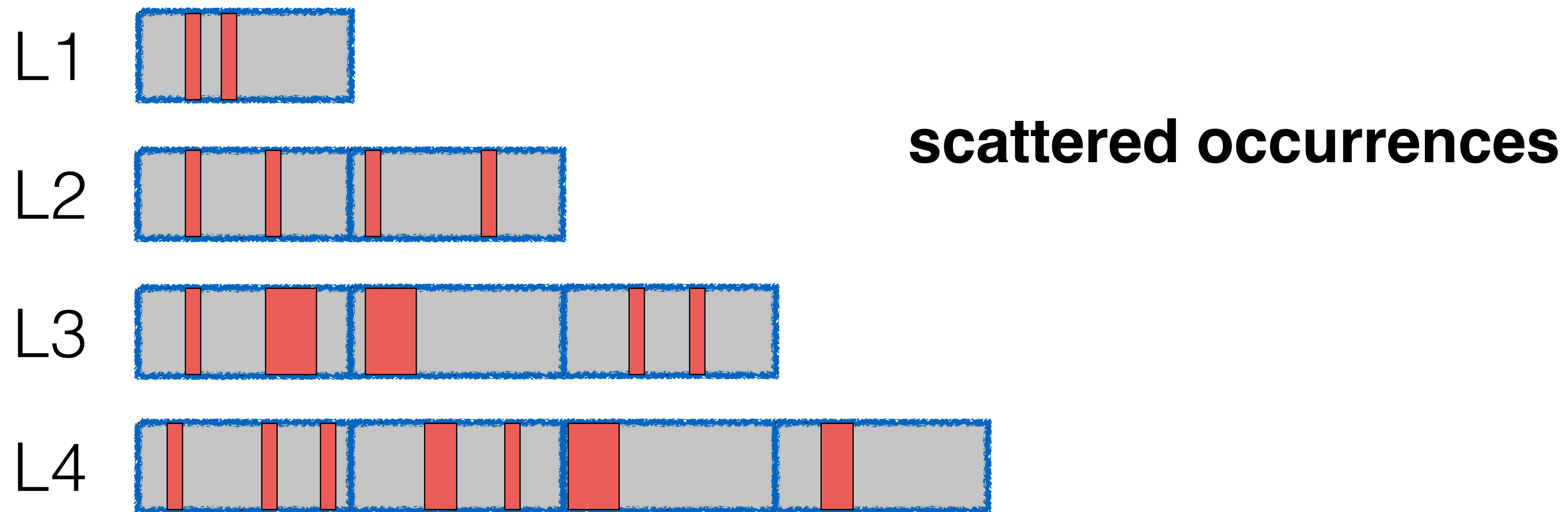
latency spikes

# Kiwi

superfluous I/Os

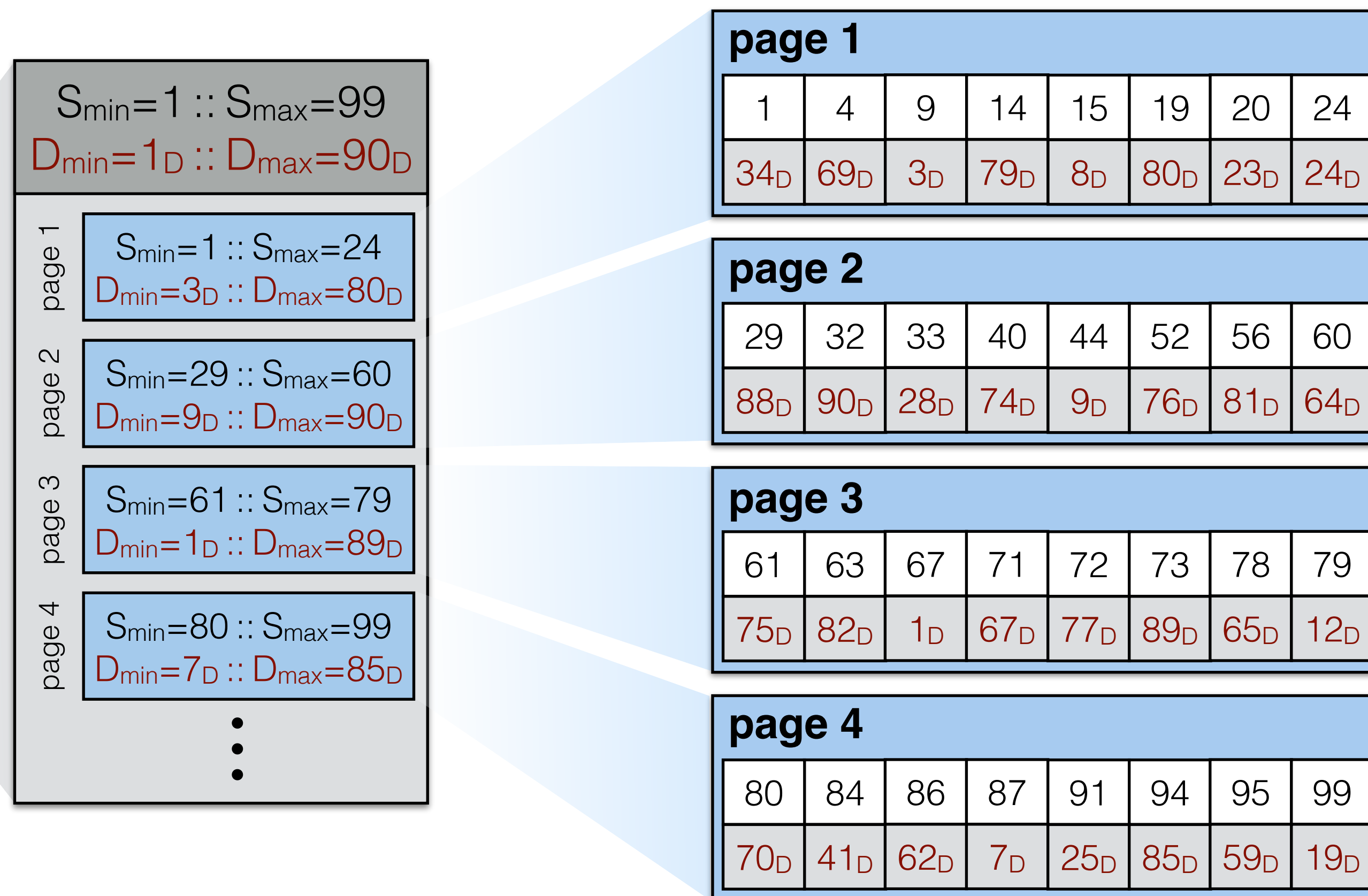
# Key Weaving storage layout

delete all entries older than: **D days**



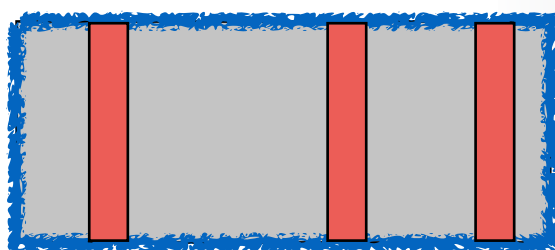
# Key Weaving storage layout

delete all entries with timestamp  $\leq 65_D$

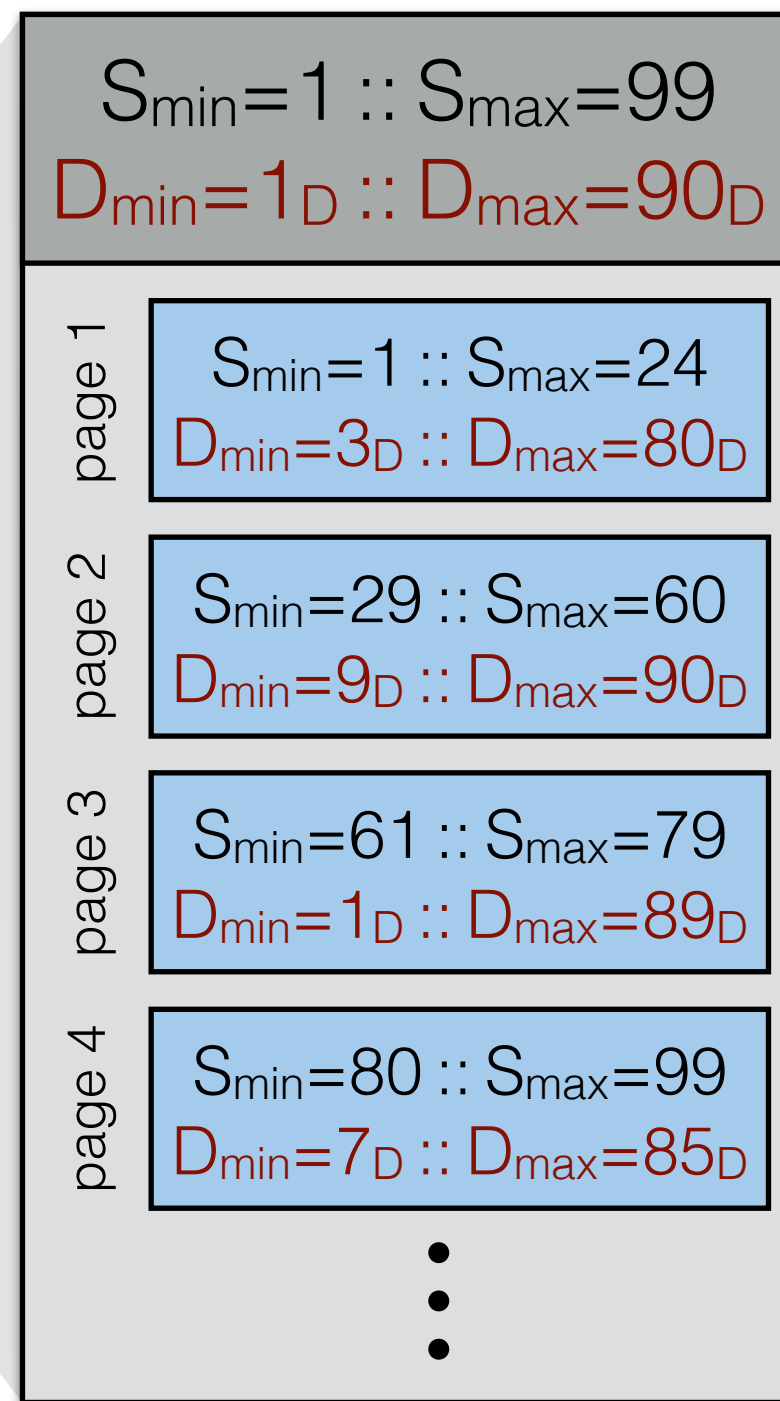


# Key Weaving storage layout

delete all entries with timestamp  $\leq 65_D$



file

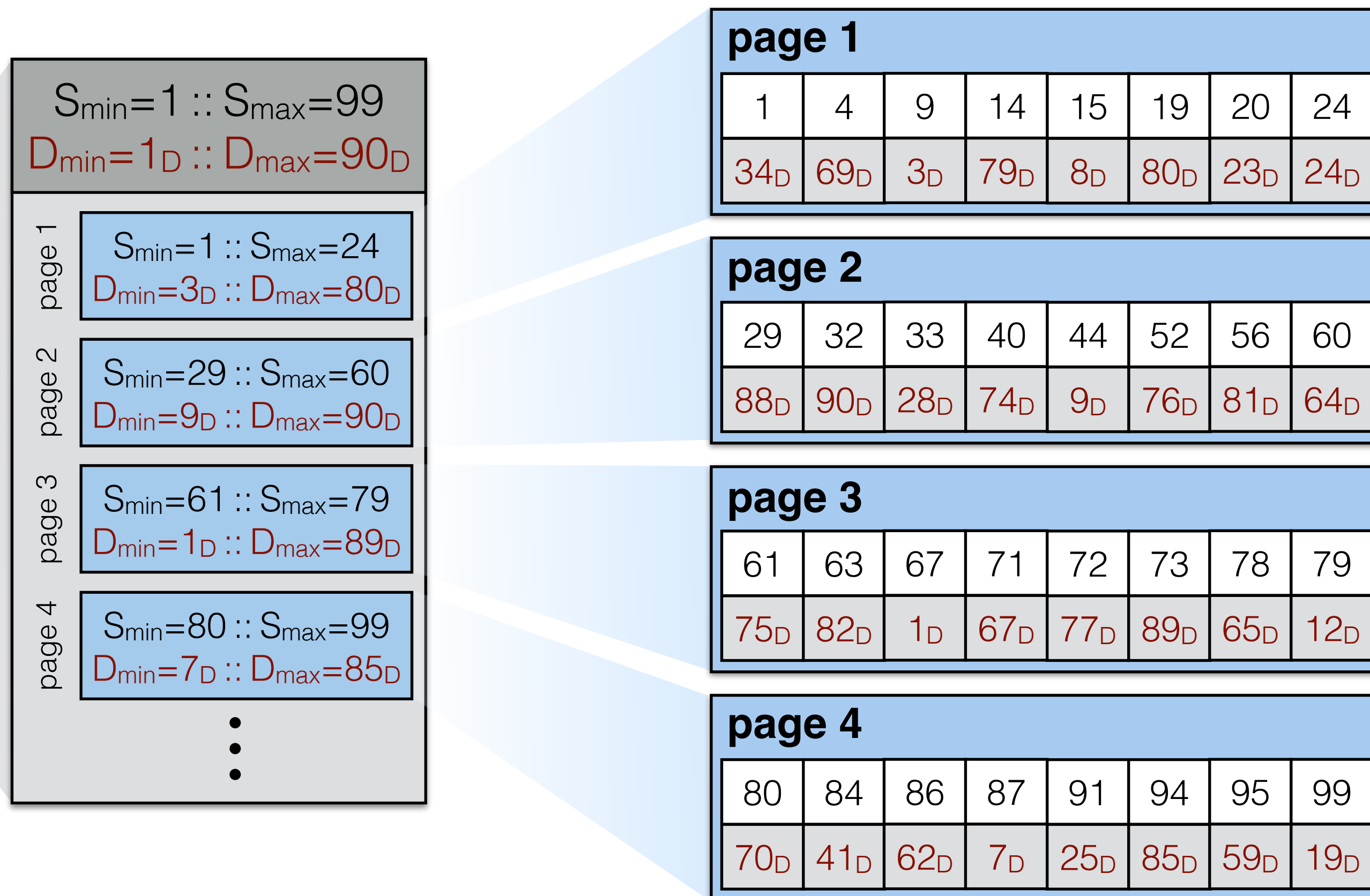


# Realizing Retention-Based Deletes

delete all entries older than  $\leq 65_D$



file

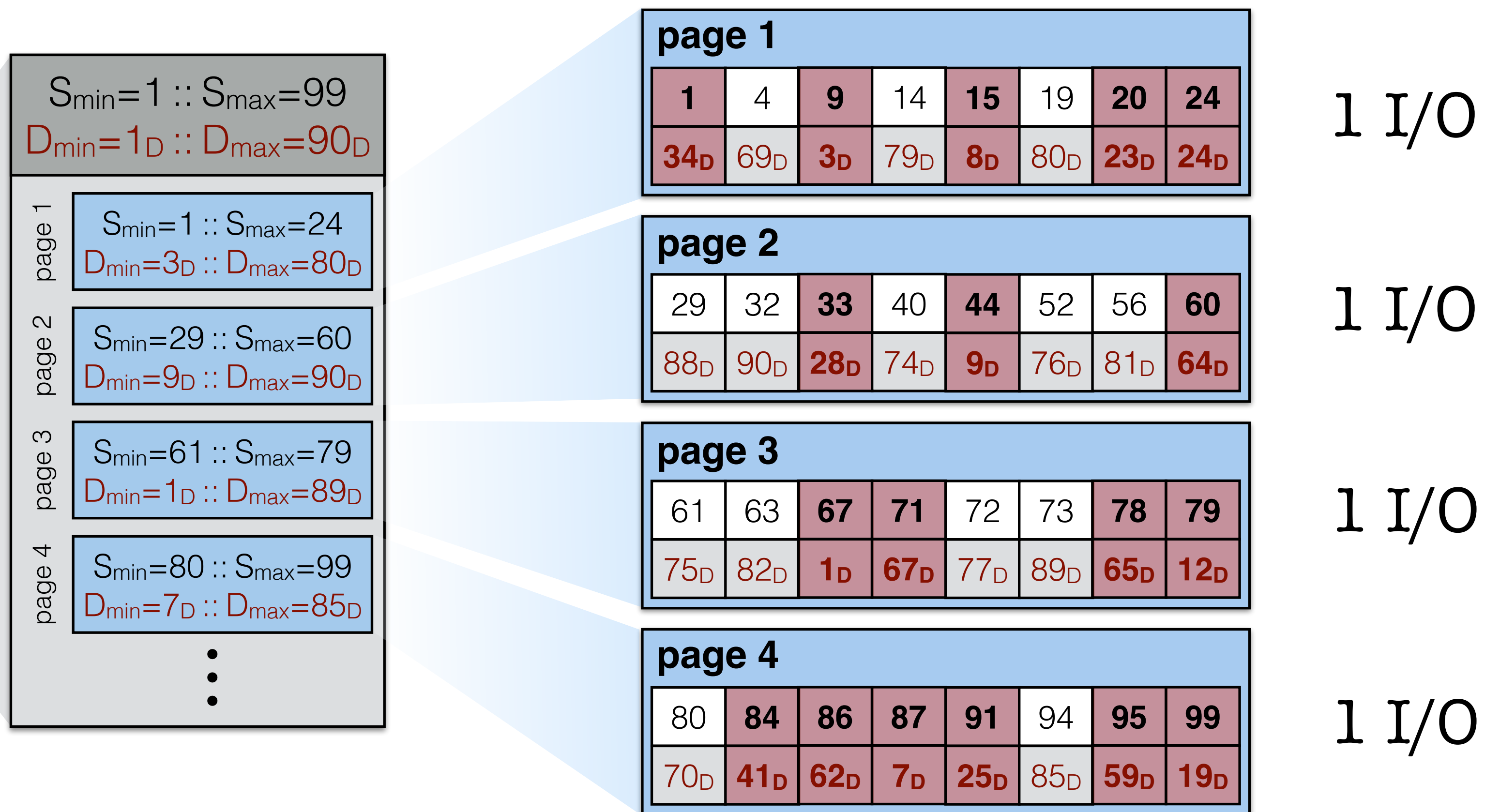


# Realizing Retention-Based Deletes

delete all entries older than  $\leq 65_D$



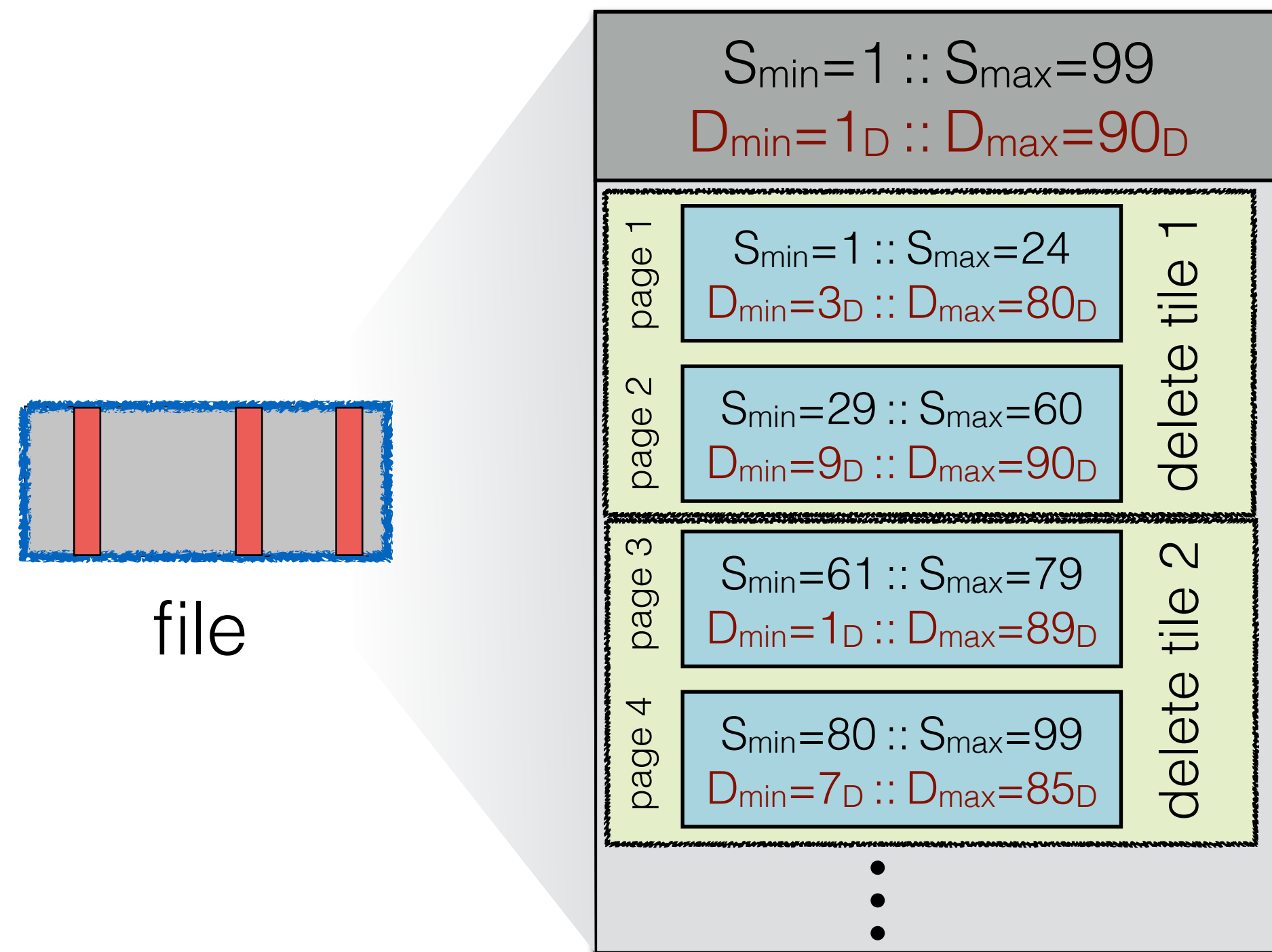
file



++

# Key Weaving storage layout

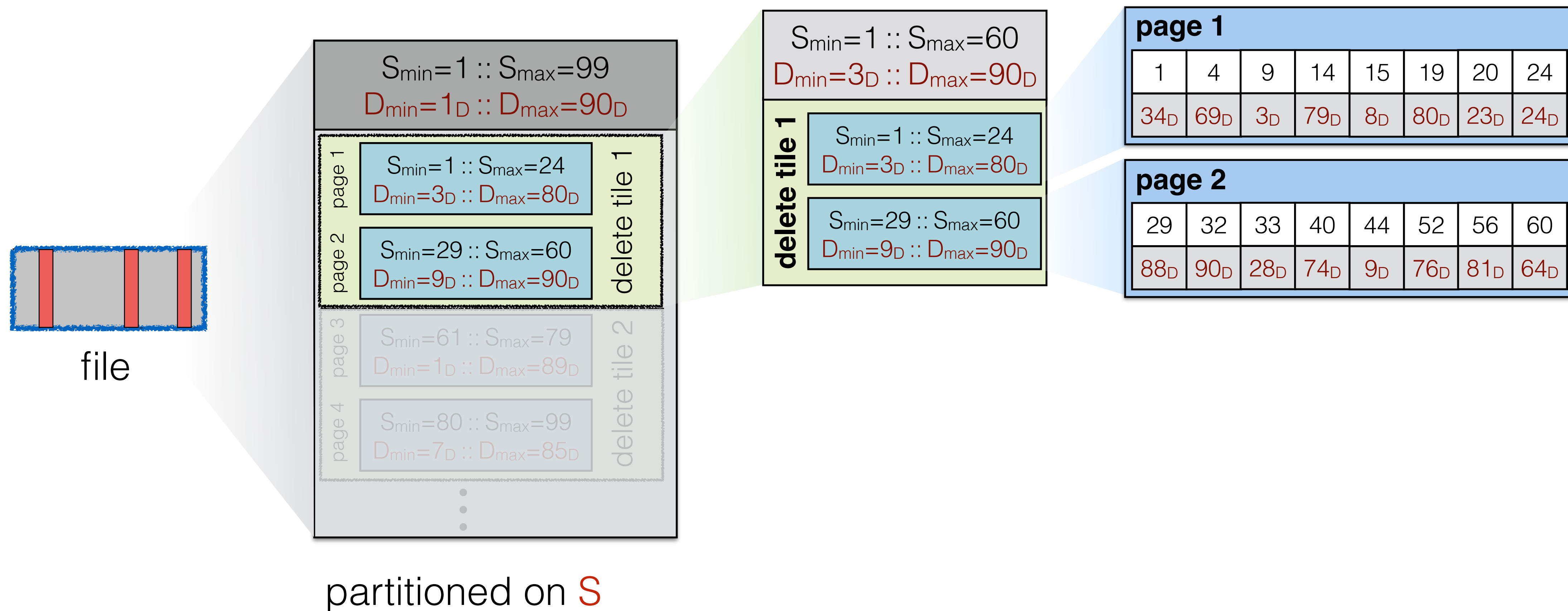
delete all entries with timestamp  $\leq 65_D$



partitioned on  $S$

# Key Weaving storage layout

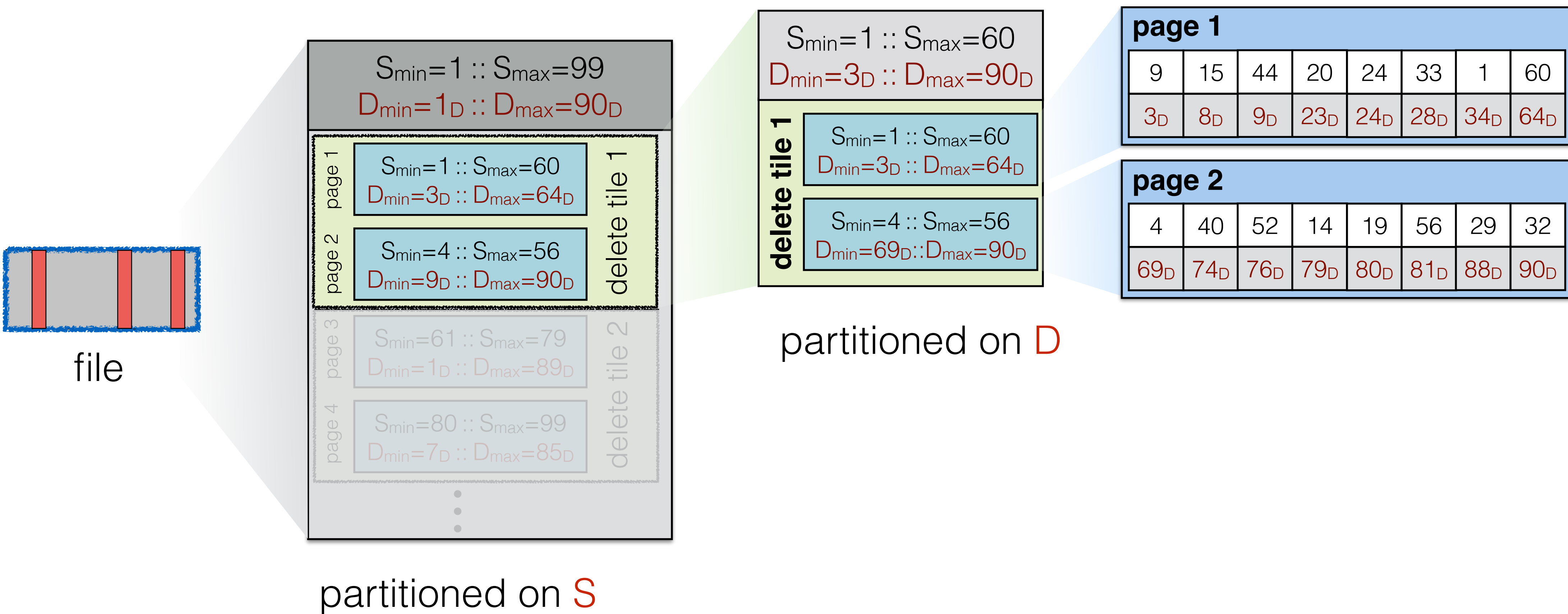
delete all entries with timestamp  $\leq 65_D$





# Key Weaving storage layout

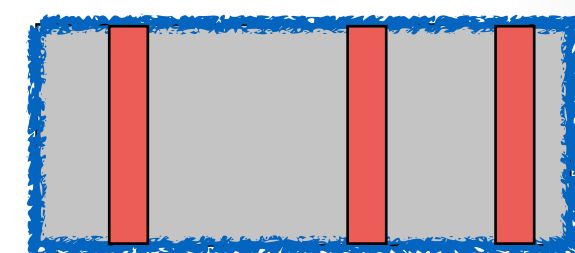
delete all entries with timestamp  $\leq 65D$



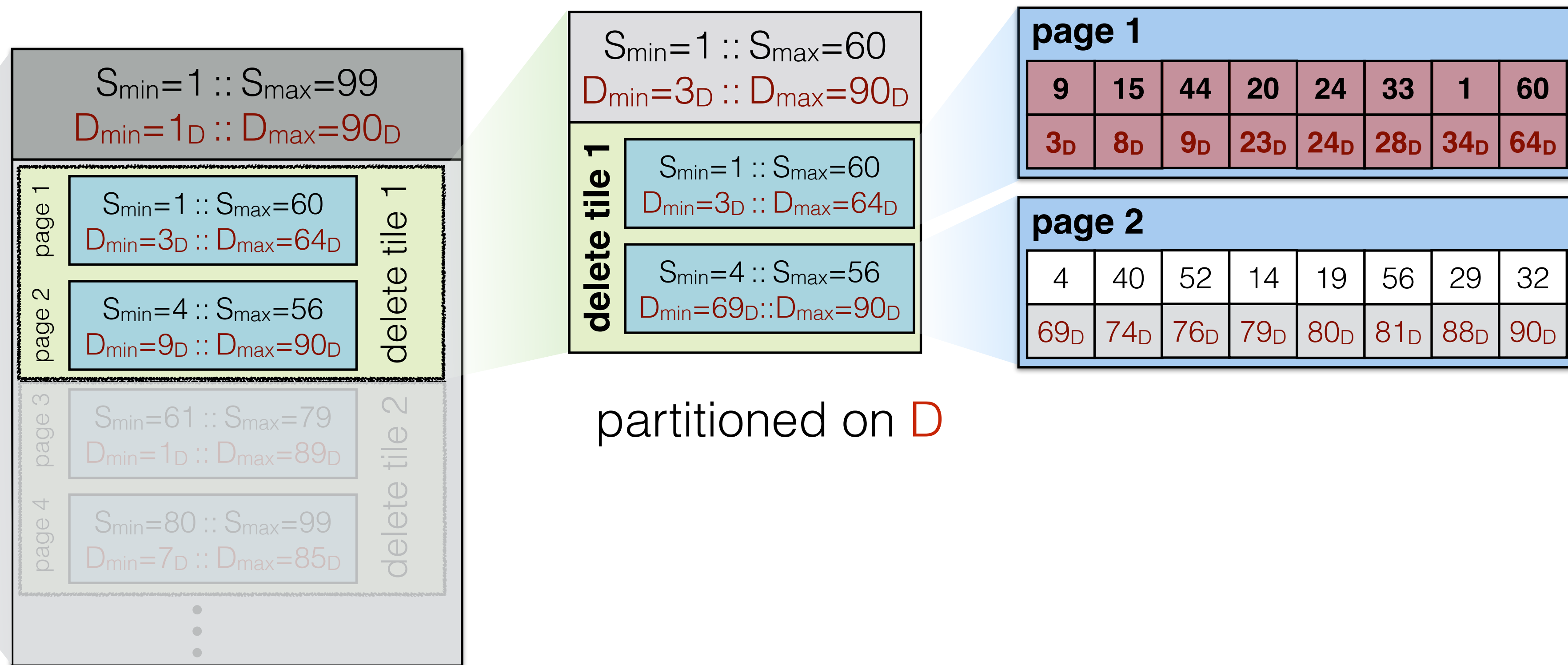
# Key Weaving storage layout

delete all entries with timestamp  $\leq 65_D$

drop  
page



file



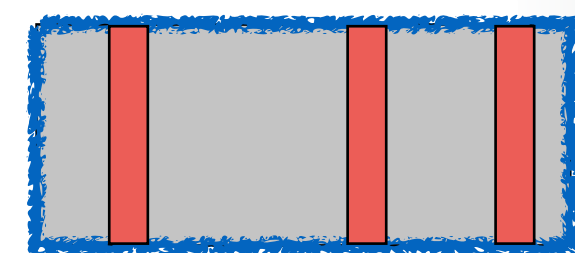
partitioned on D

partitioned on S

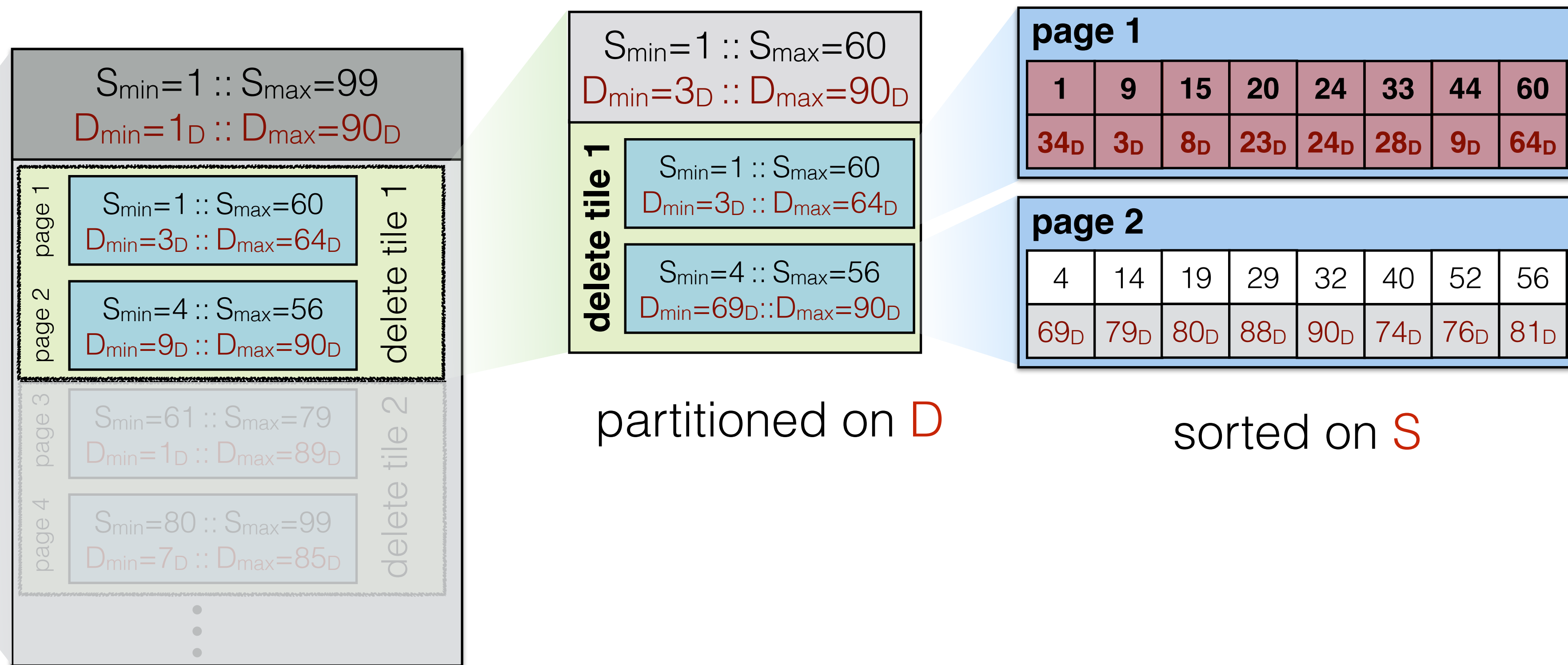
# Key Weaving storage layout

delete all entries with timestamp  $\leq 65_D$

drop  
page



file



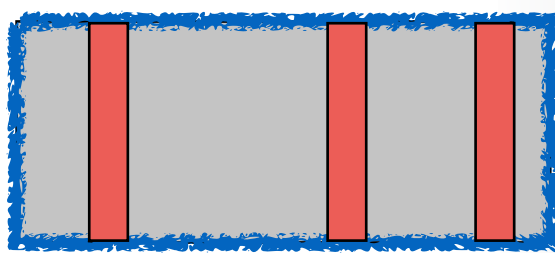
partitioned on  $D$

sorted on  $S$

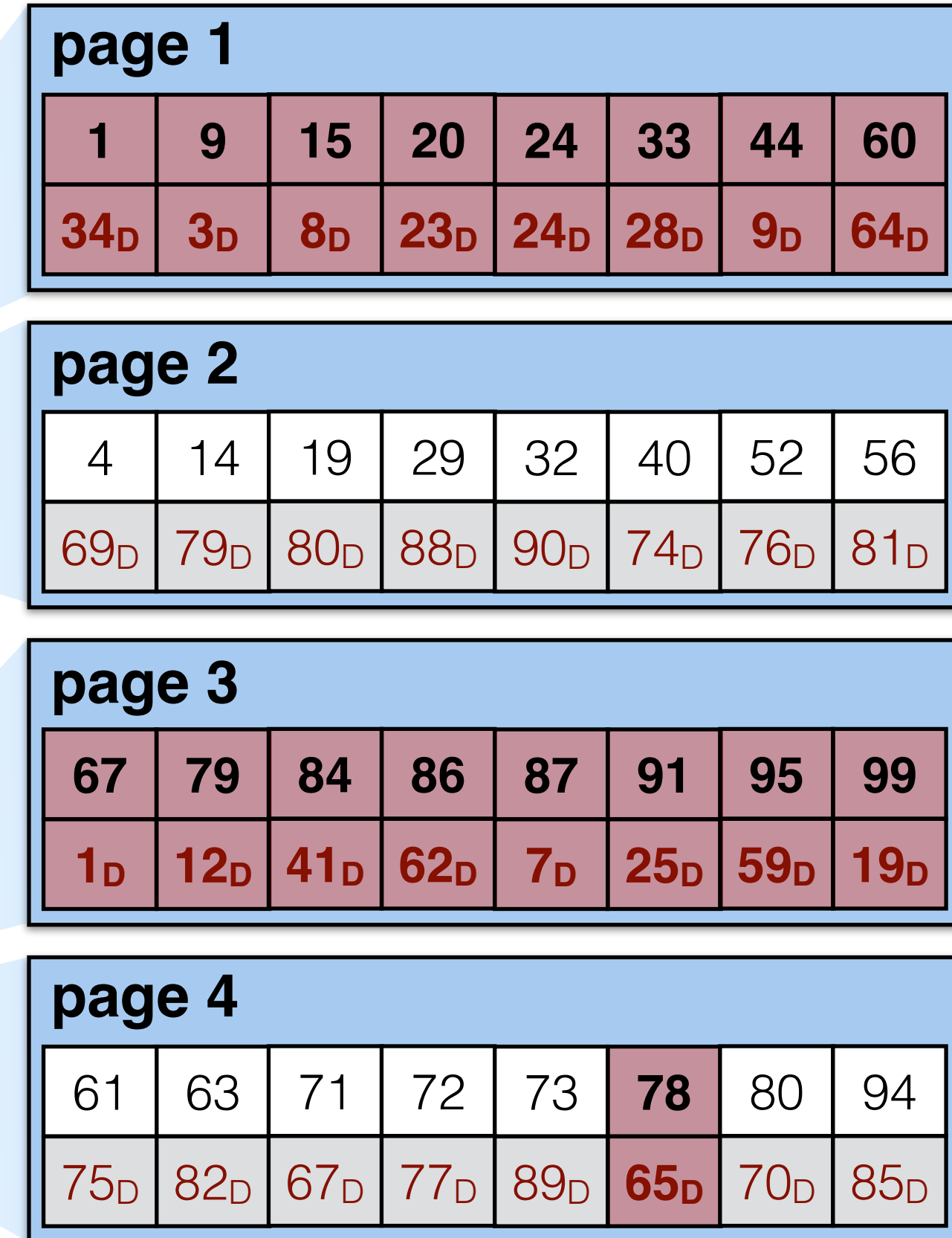
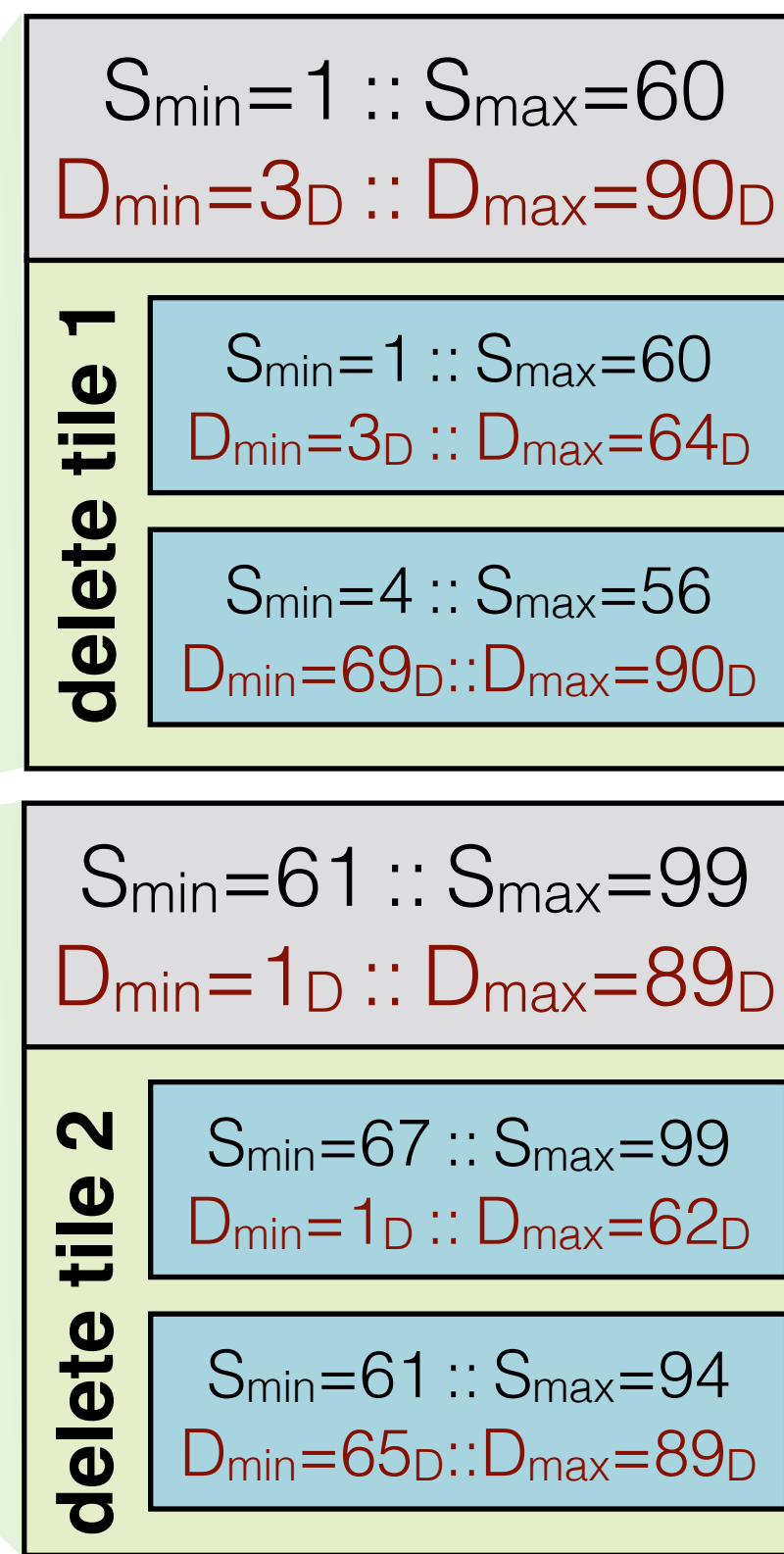
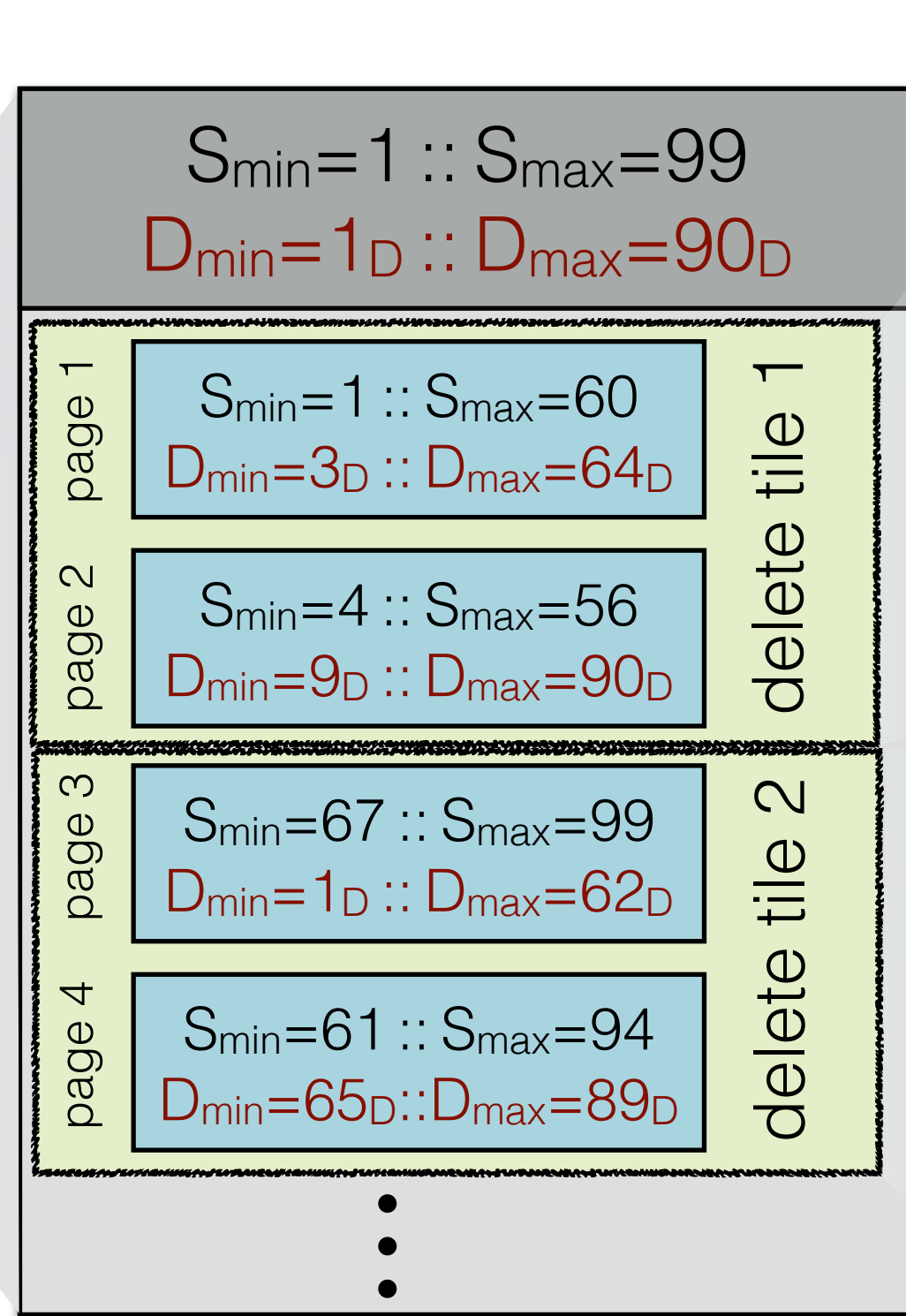
partitioned on  $S$

# Key Weaving storage layout

delete all entries with timestamp  $\leq 65_D$



file

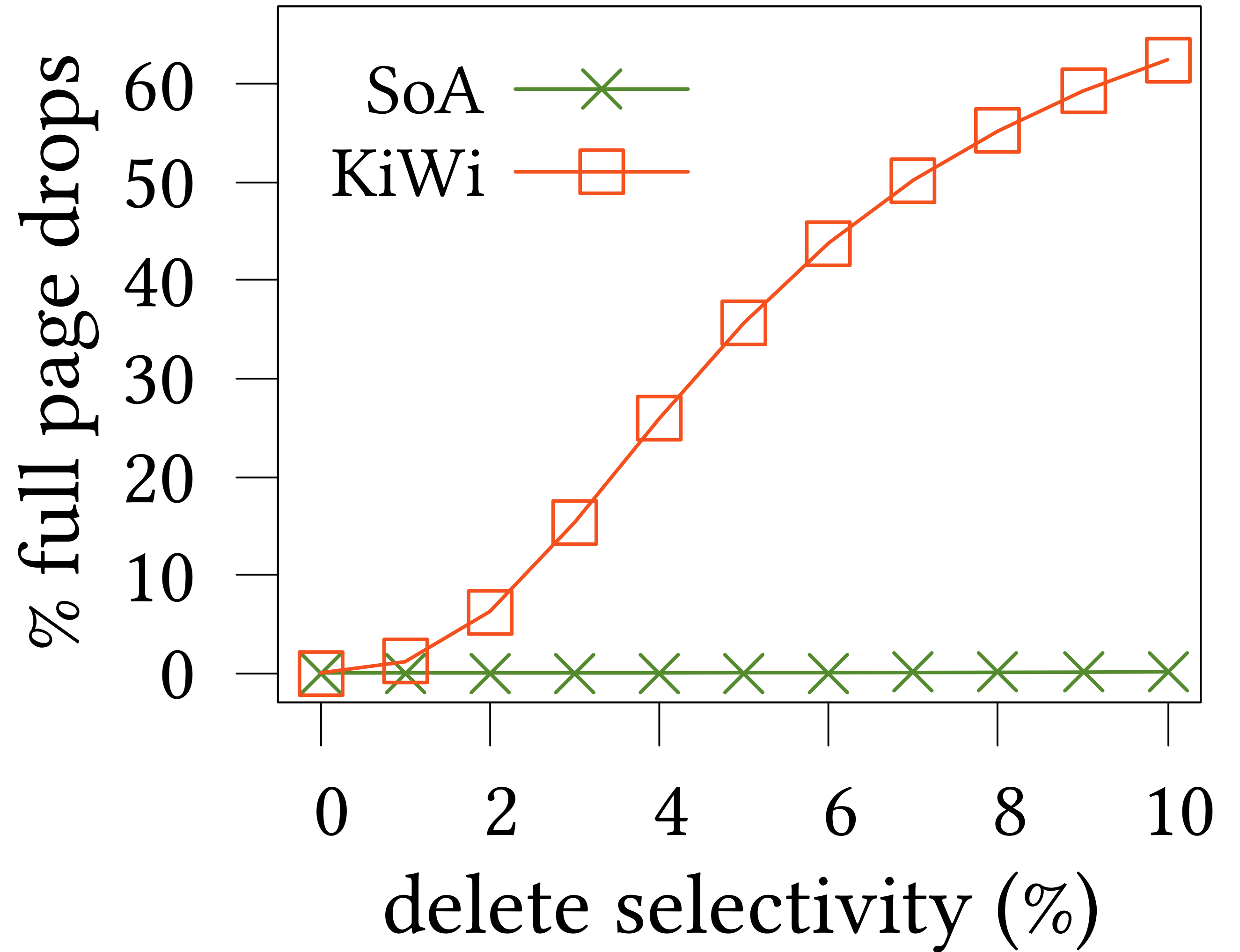


drop page

drop page

1 I/O

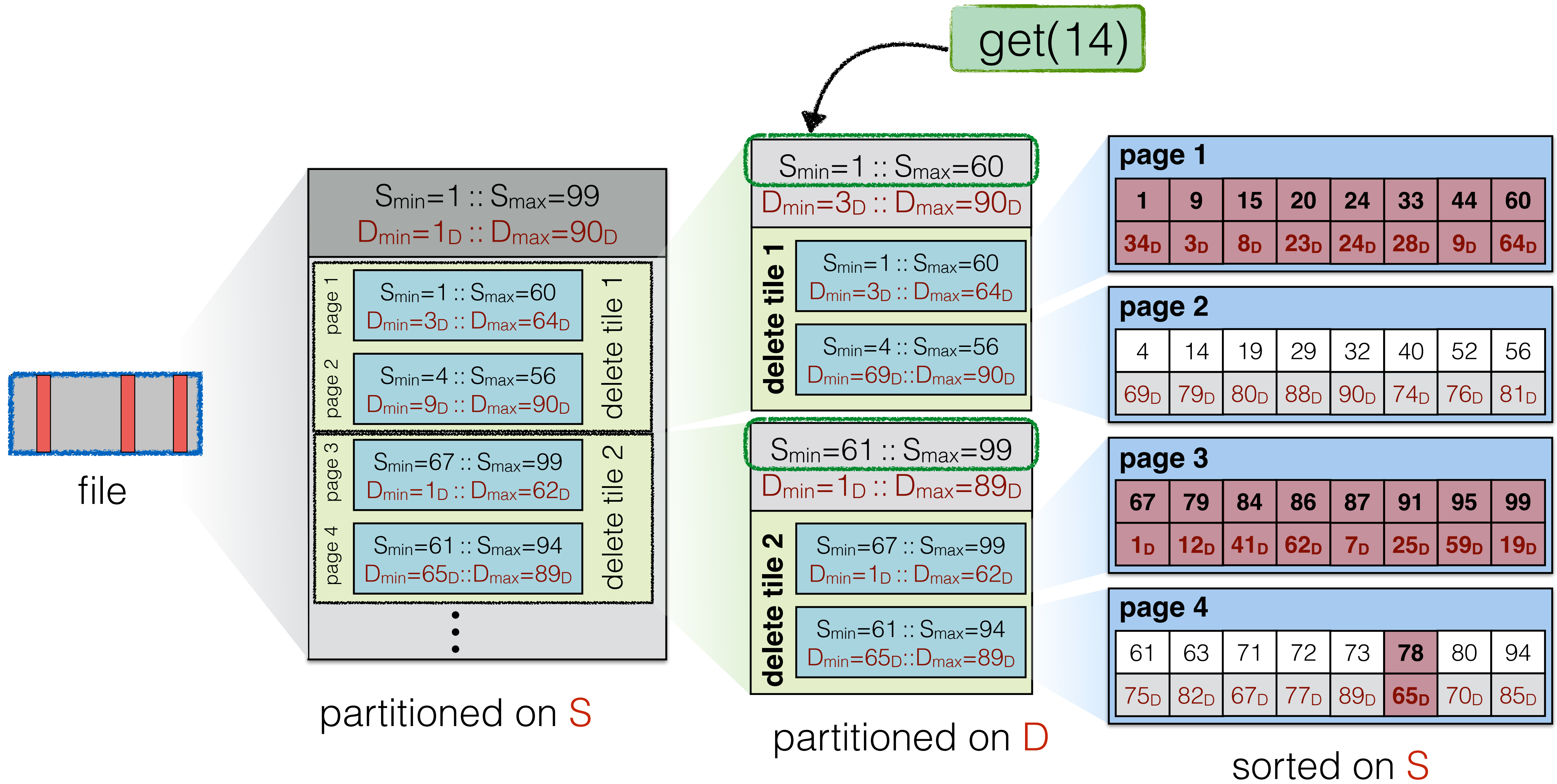
1M point lookups, buffer = file = 256 pages, T=10



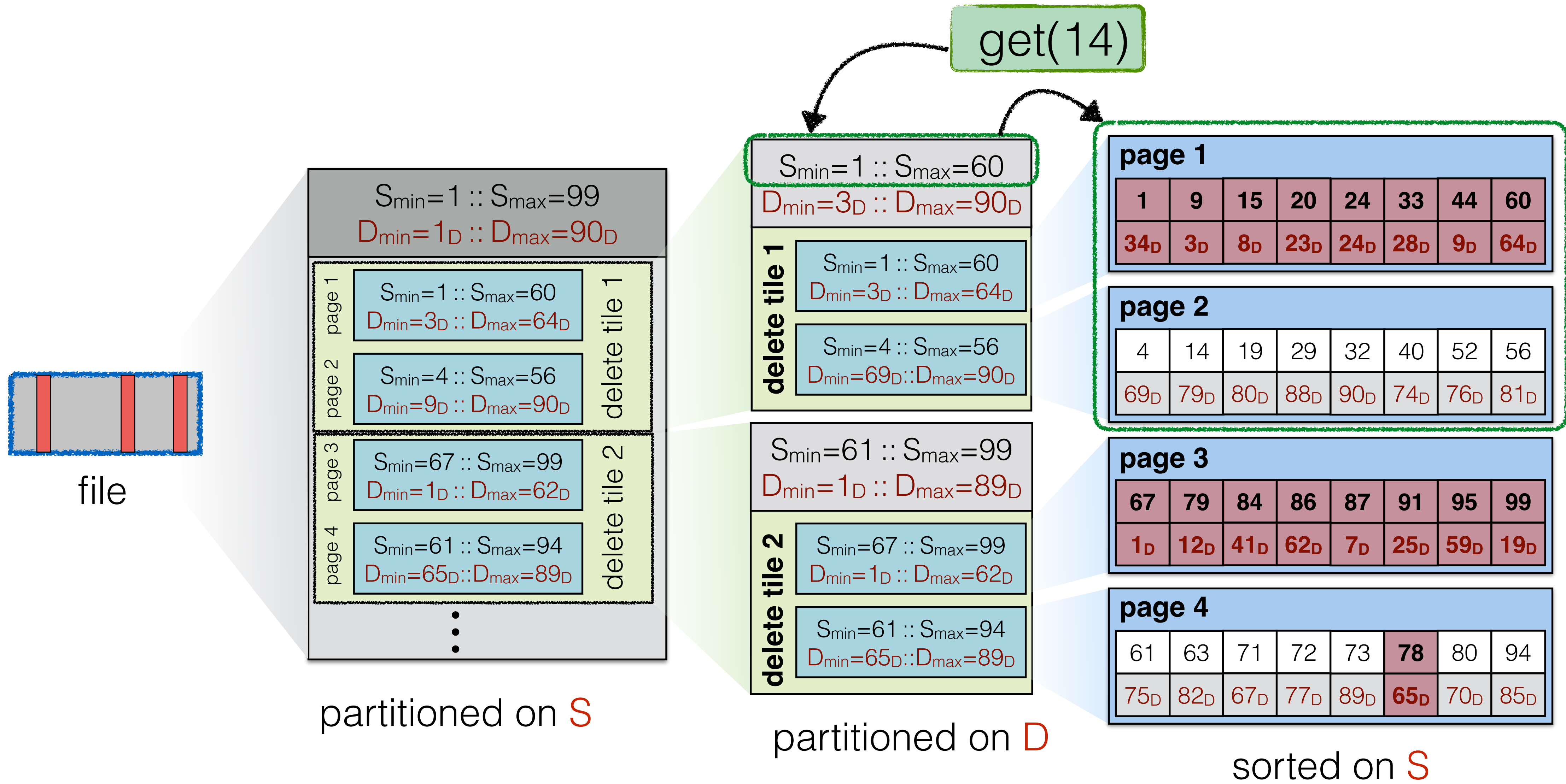
superior delete performance

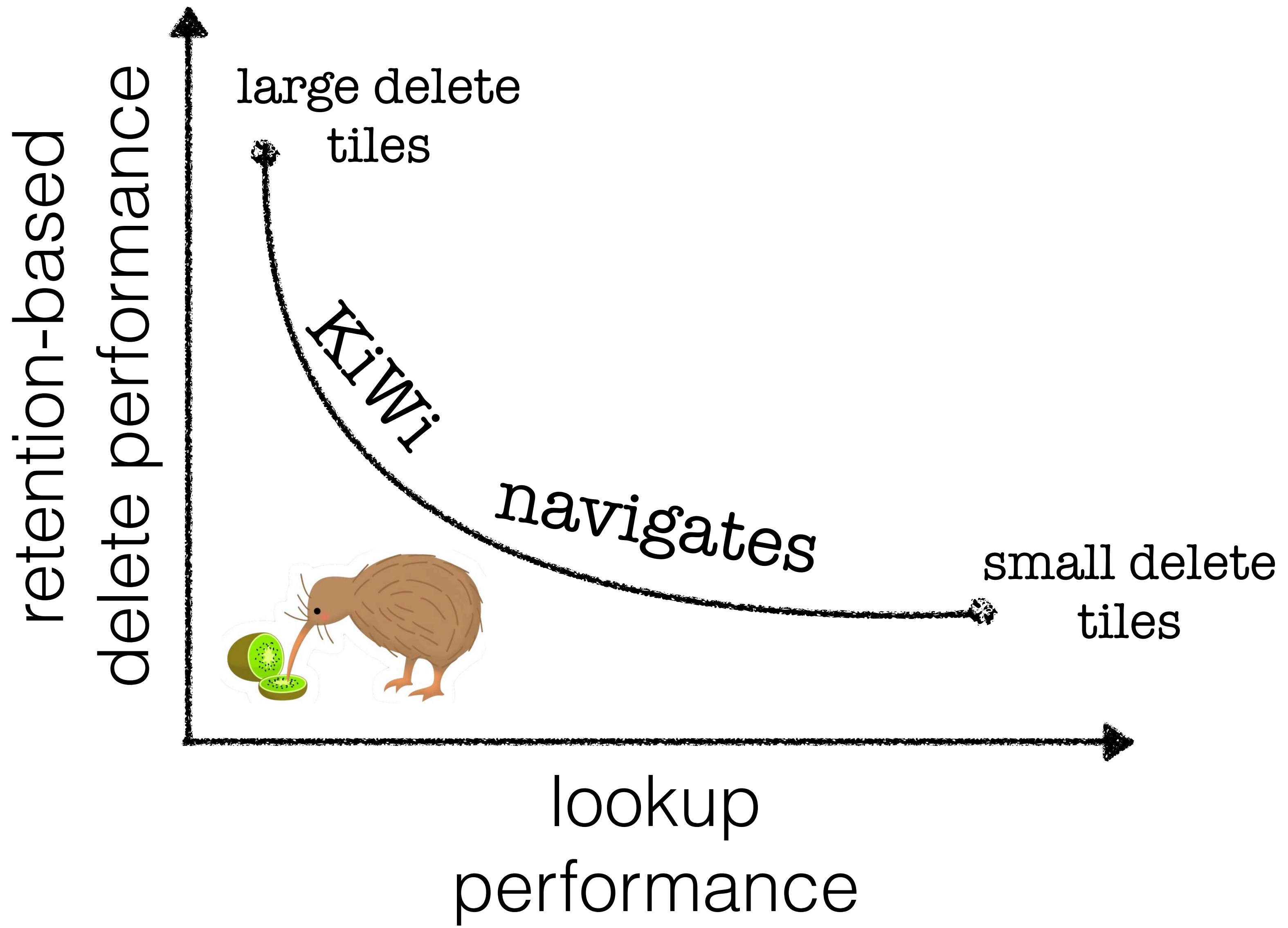
up to 2.5x

# Key Weaving storage layout

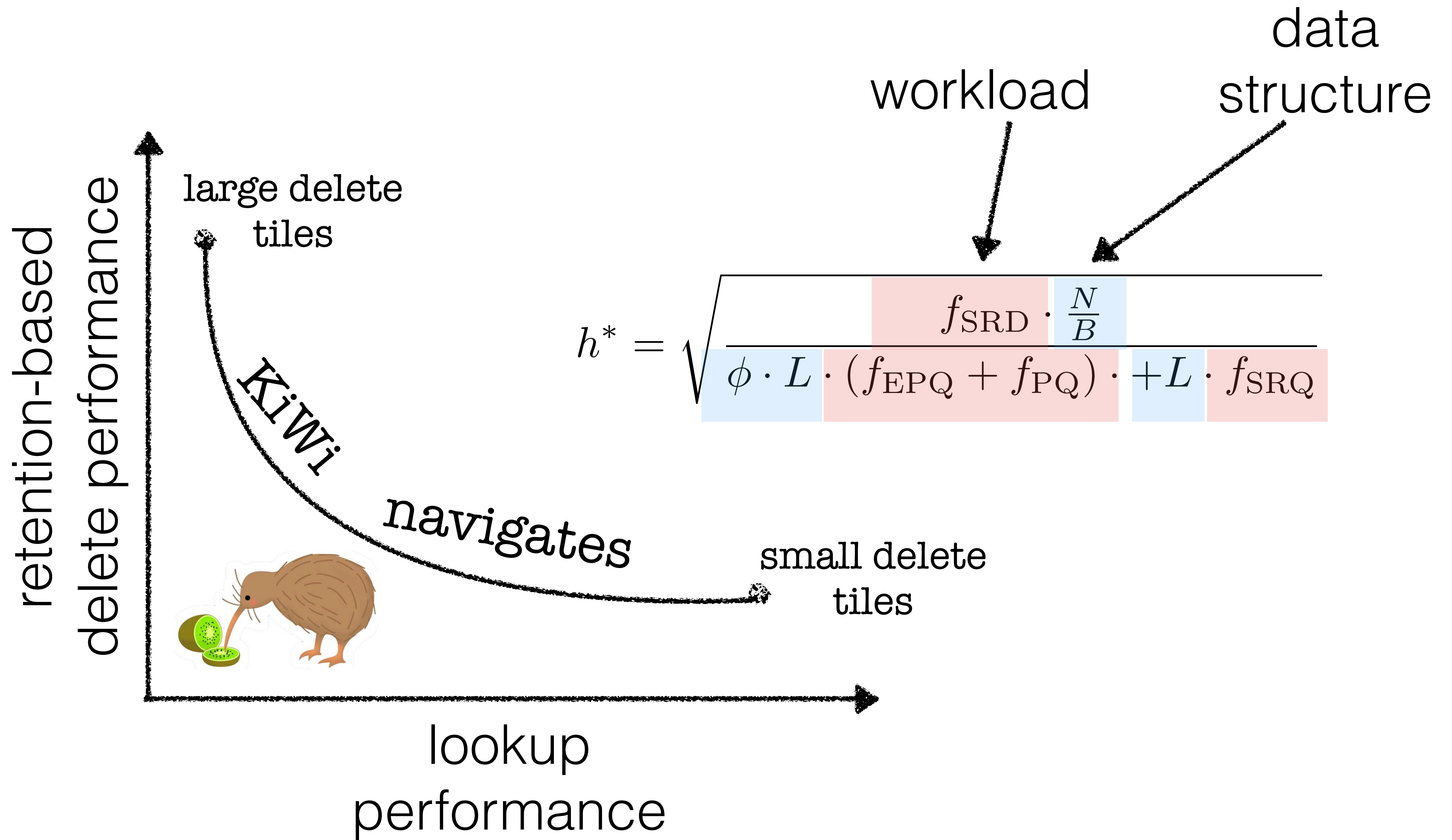


# Key Weaving storage layout









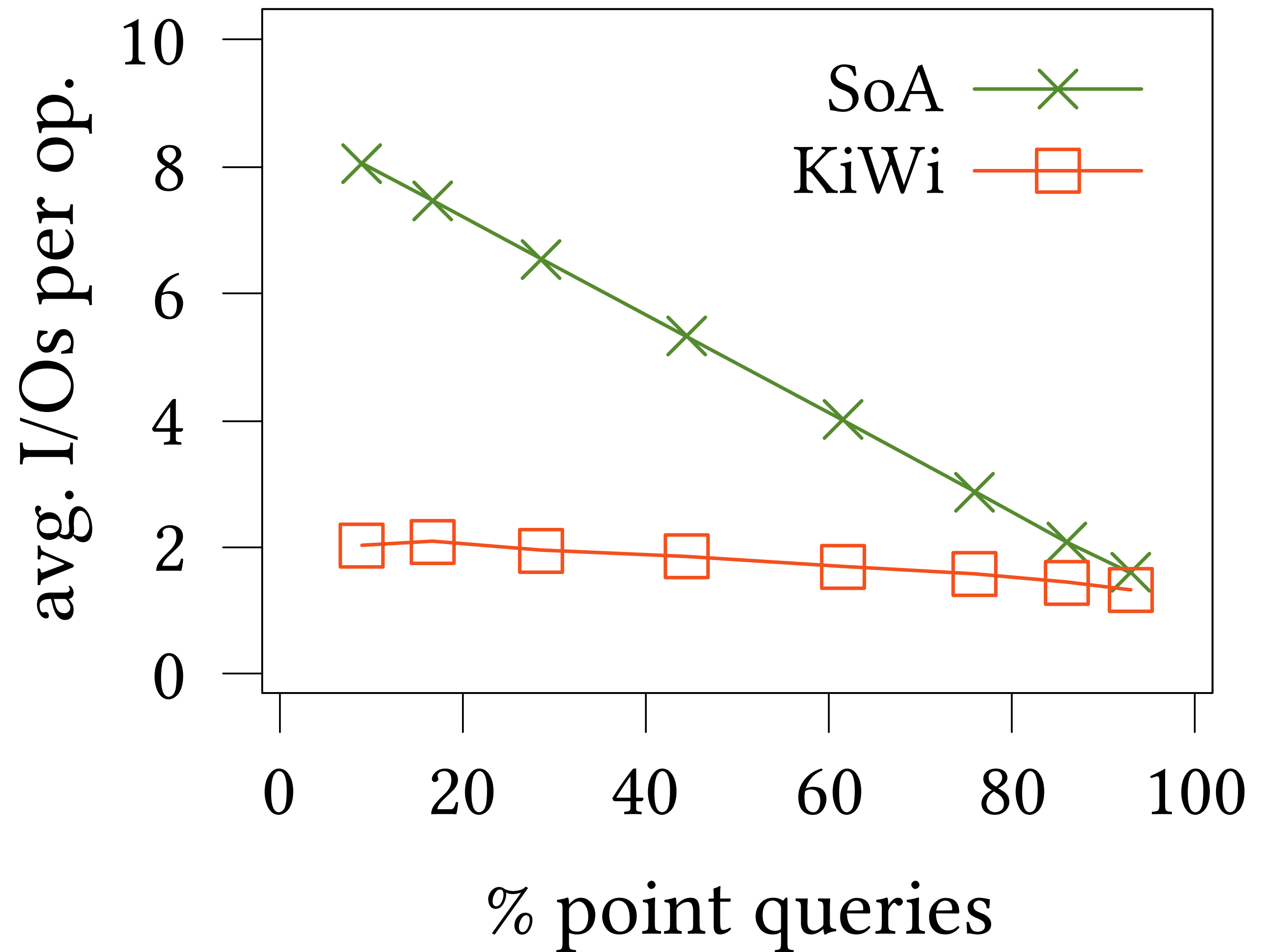
better overall performance

up to 4x

superior delete performance

up to 2.5x

1M point lookups, buffer = file = 256 pages, T=10



# the solution

FAst DElete

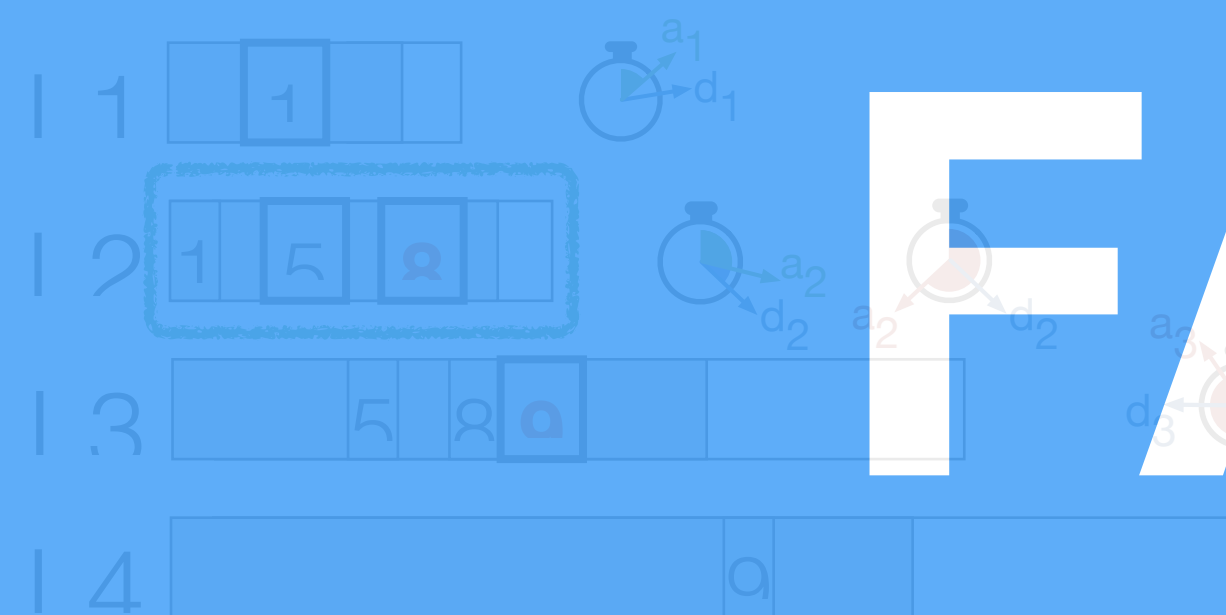
amortized write

reduced space

improved read

timely delete persistence

# FADE

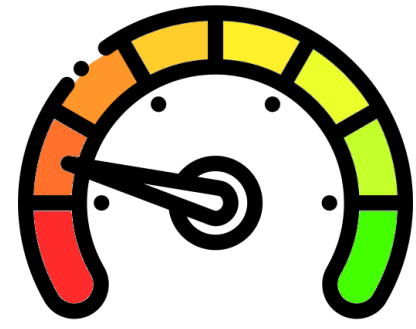


higher lookup cost

# Kiwi

reduced latency spikes

full page drops reduces  
superfluous I/Os



suboptimal state-of-the-art design  
for workloads with deletes



FADE persists deletes timely  
using latency-driven compactions



KiWi supports efficient  
secondary range deletes  
using key-interweaved data storage

# *CS 561: Data Systems Architecture*

## Class 6

# **Efficient Deletes in LSM-Engines**

BOSTON  
UNIVERSITY

*Dr. Subhadeep Sarkar*

<https://bu-disc.github.io/CS561/>

