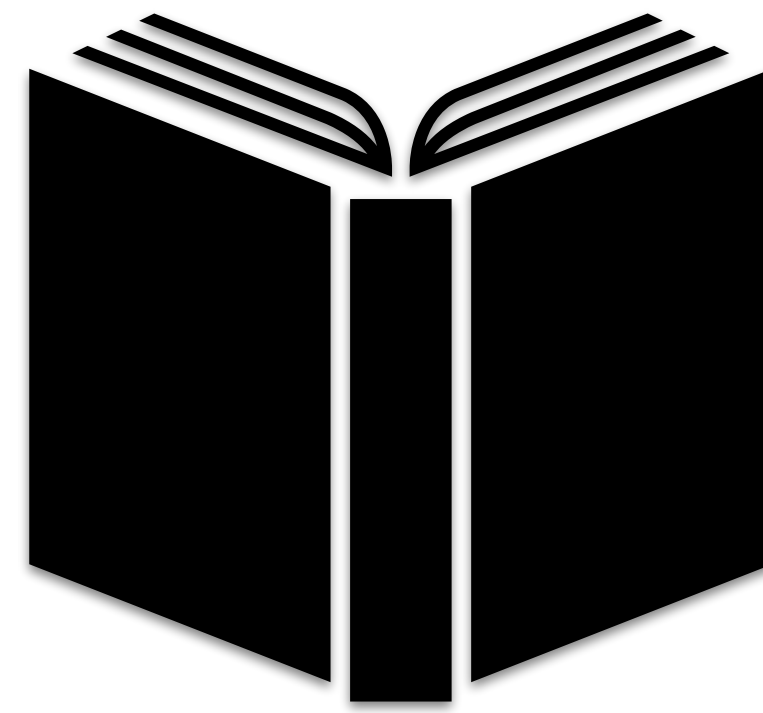
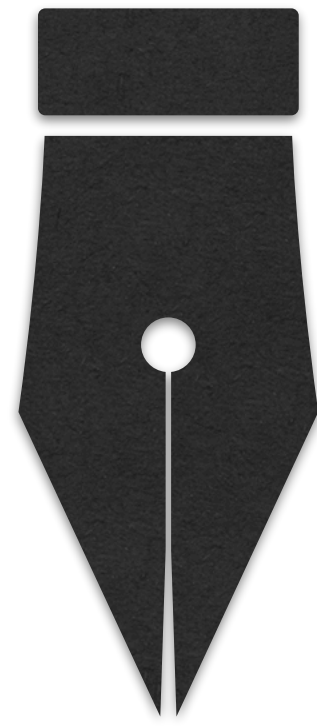


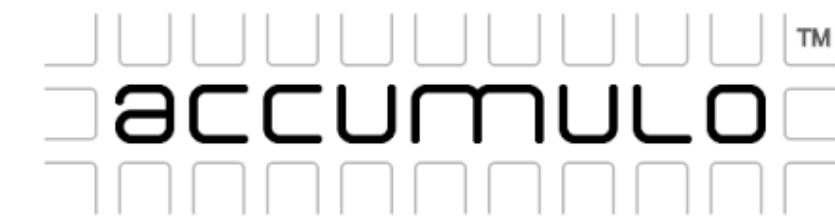
Lethe: A Tunable Delete-Aware LSM-Based Storage Engine



Subhadeep Sarkar
Tarikul Islam Papon
Dimitris Staratzis
Manos Athanassoulis

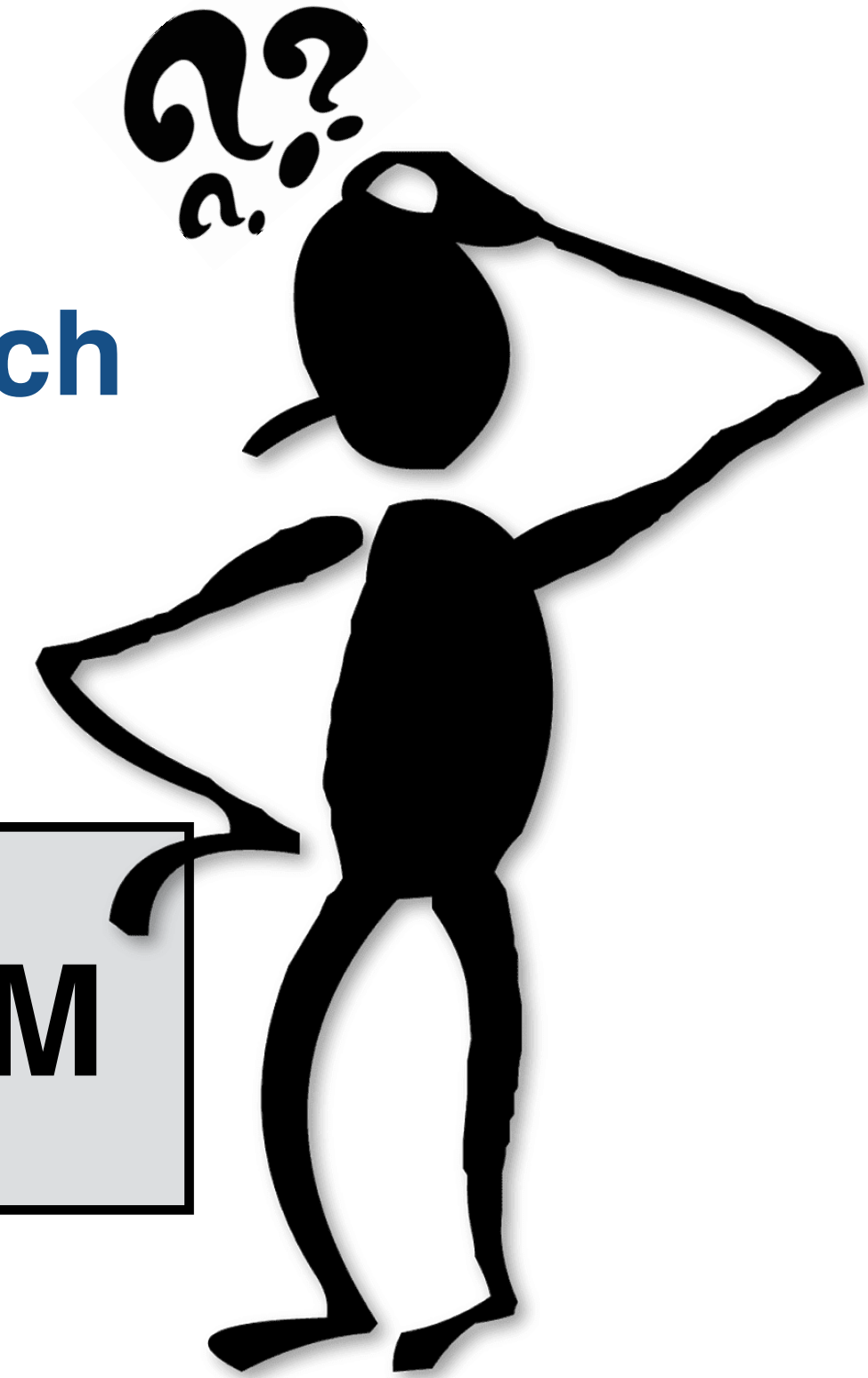


LSM-TREE



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

TechCrunch
Feb '19



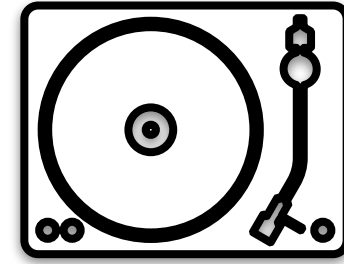
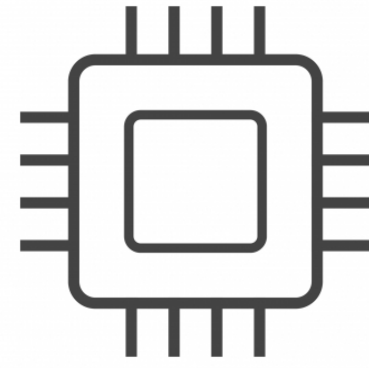
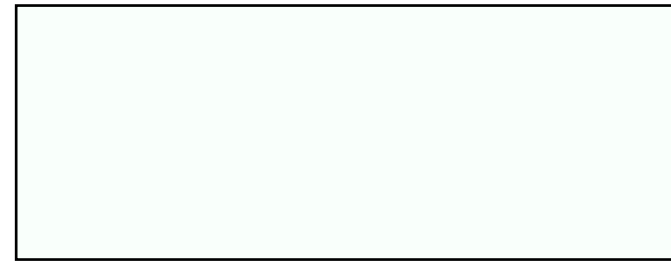
Small Datum
Jan '20

Deletes are fast and slow in an LSM

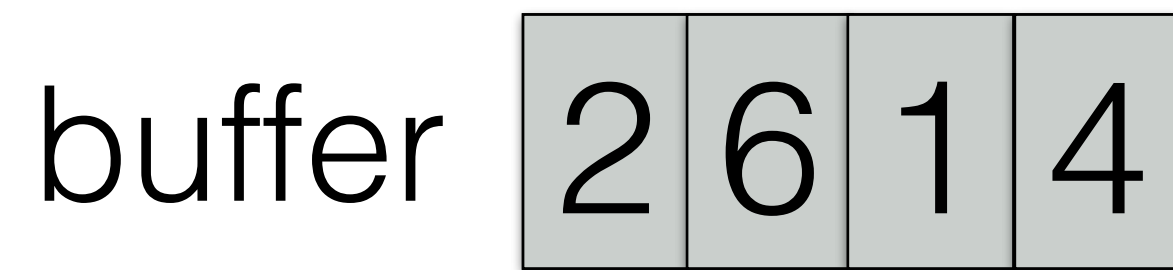
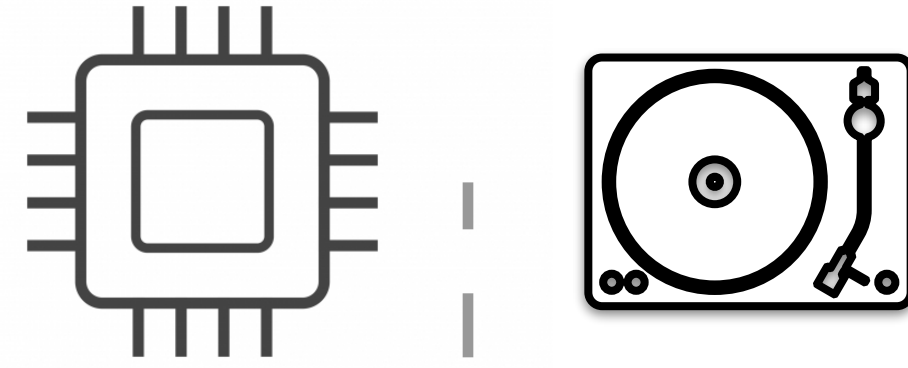
“LSM-based data stores perform suboptimally for workloads with deletes.”

log-structured merge-tree

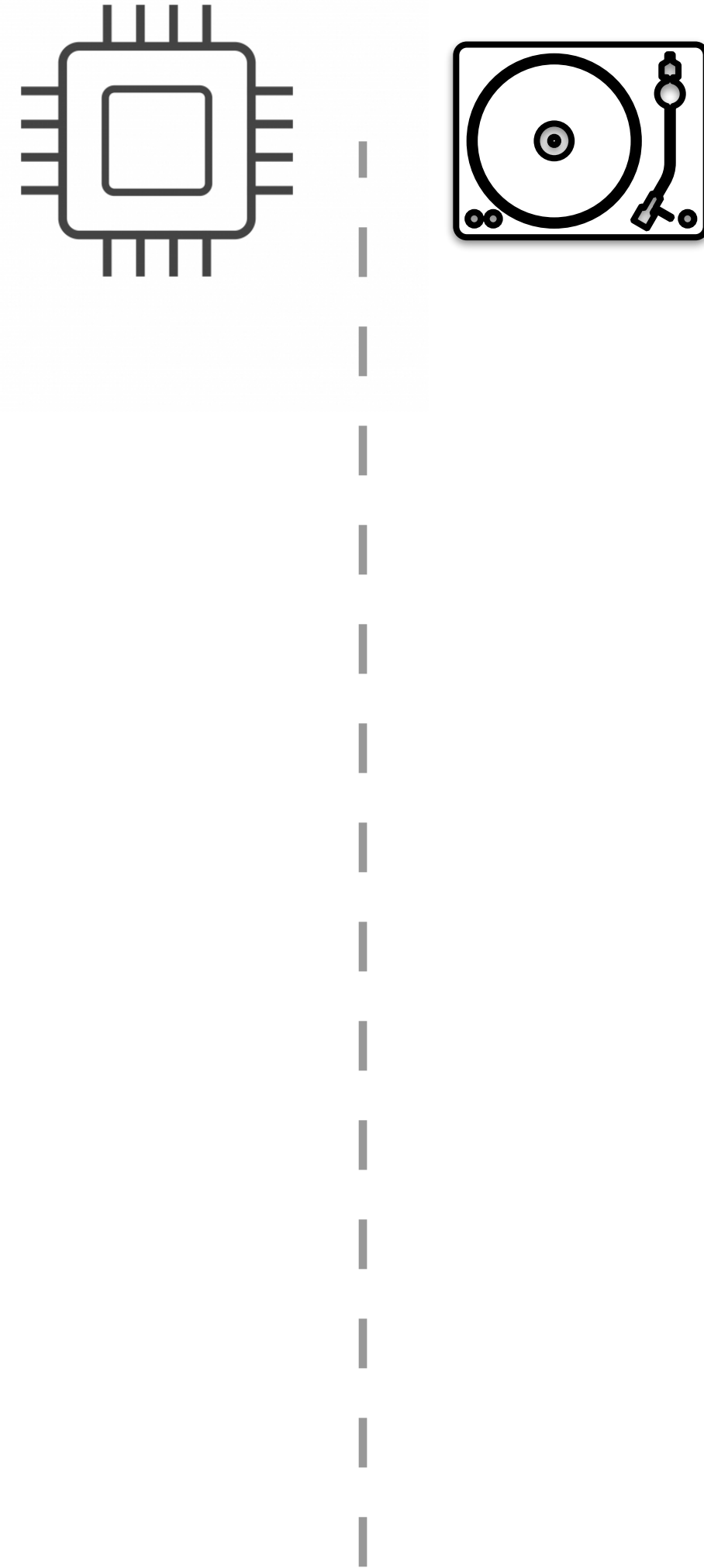
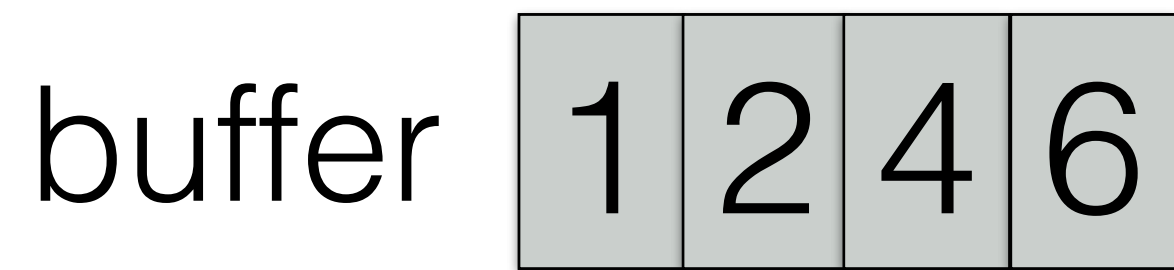
buffer



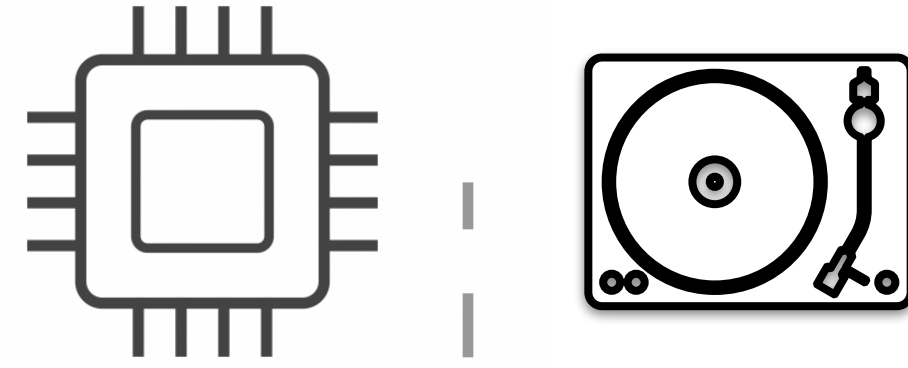
log-structured merge-tree



log-structured merge-tree



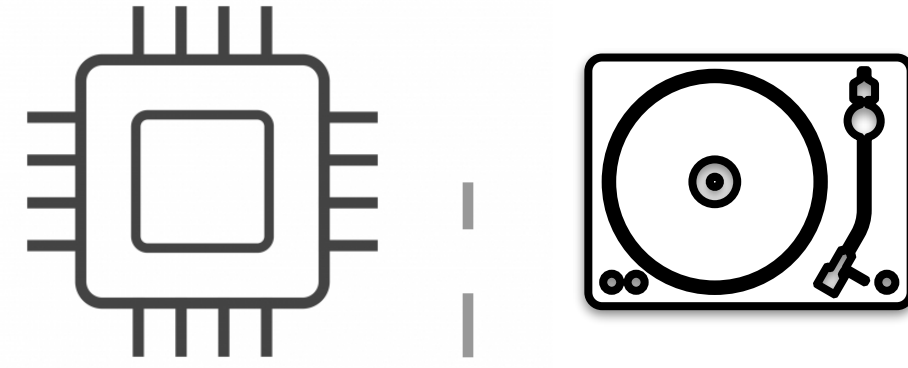
log-structured merge-tree



buffer



log-structured merge-tree



buffer 

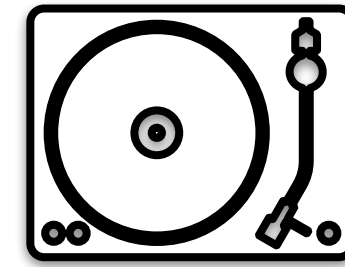
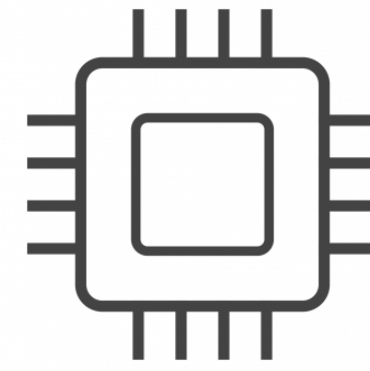
L1



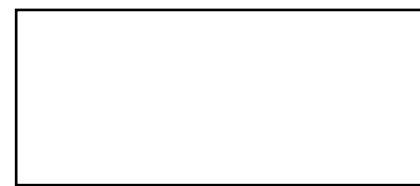
L2

L3

log-structured merge-tree



buffer



L1



size ratio = T

L2

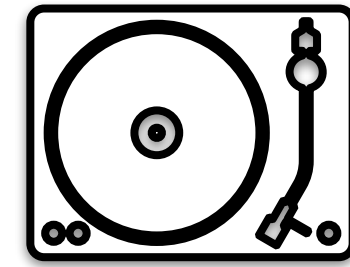
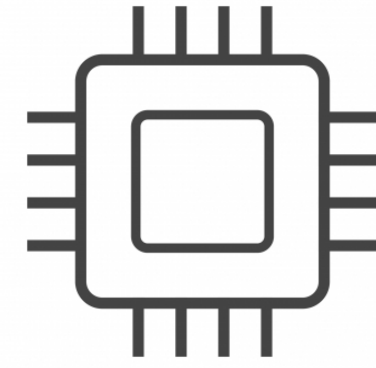


L3

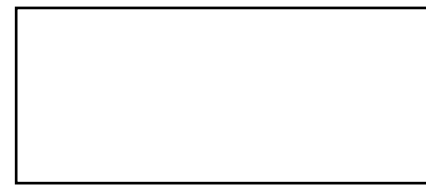


exponentially larger capacity

log-structured merge-tree



buffer



partial compaction

L1



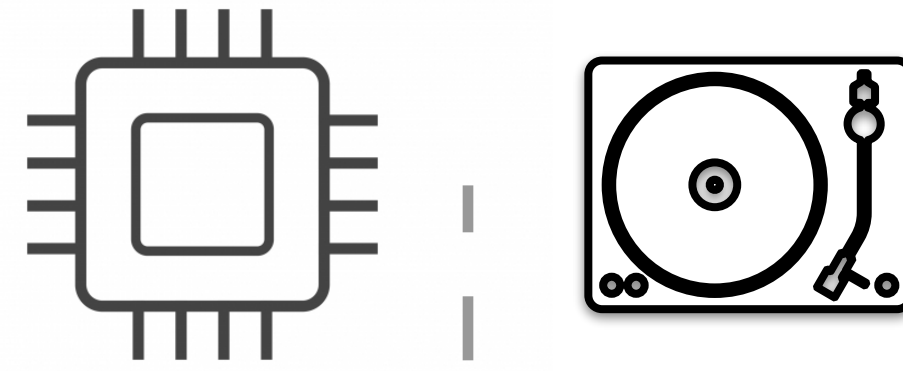
L2



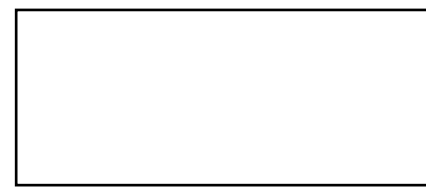
L3



log-structured merge-tree



buffer



partial compaction

L1



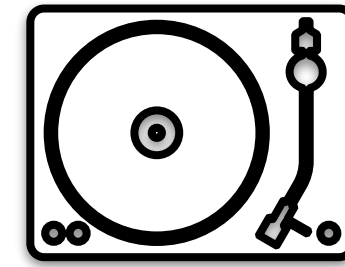
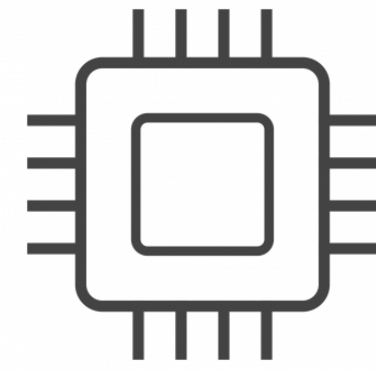
L2



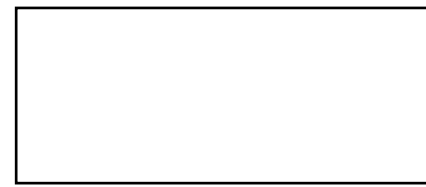
L3



log-structured merge-tree



buffer



partial compaction

L1



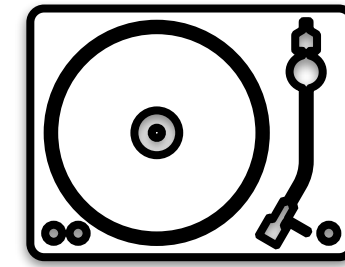
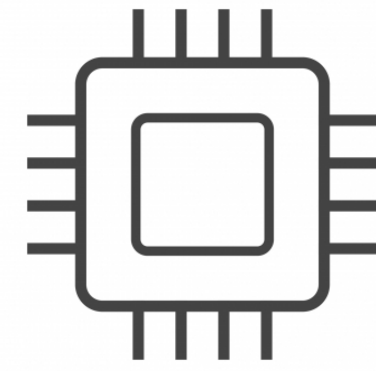
L2



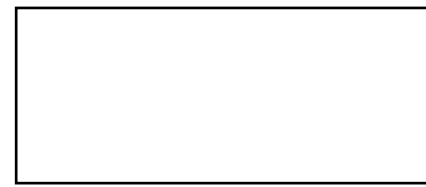
L3



log-structured merge-tree



buffer



partial compaction

L1



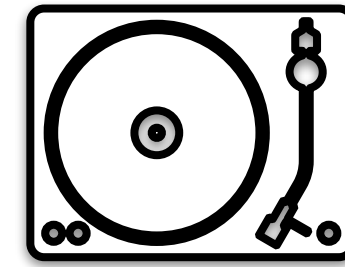
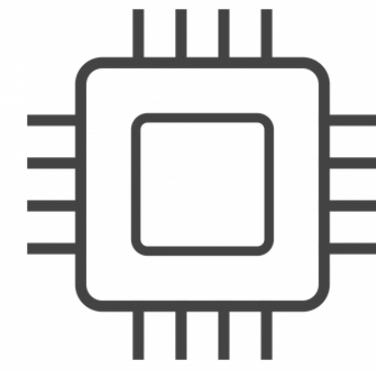
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



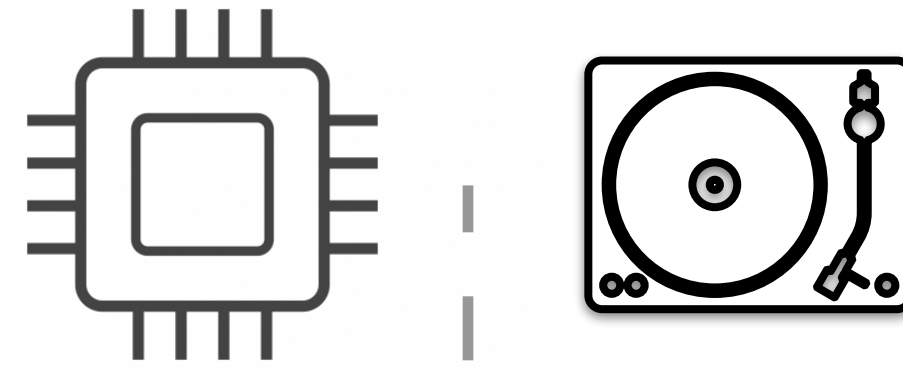
L2



L3

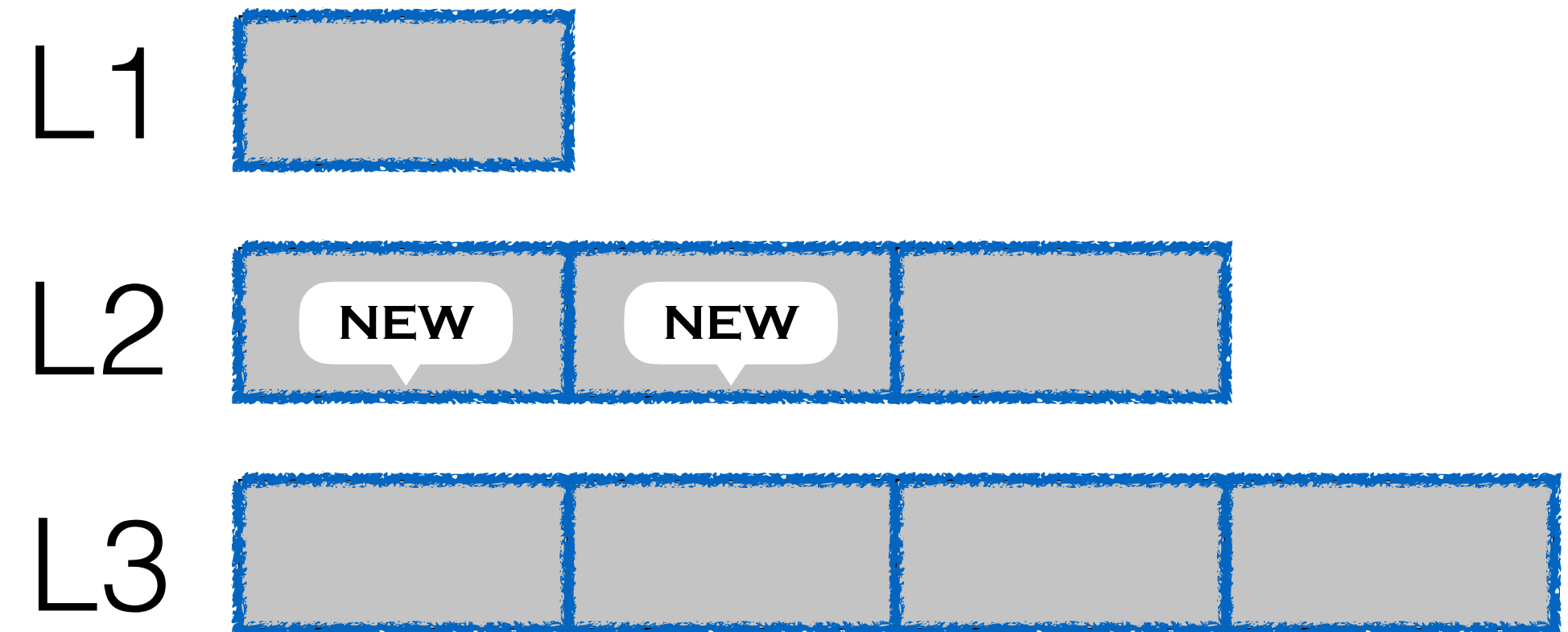


log-structured merge-tree

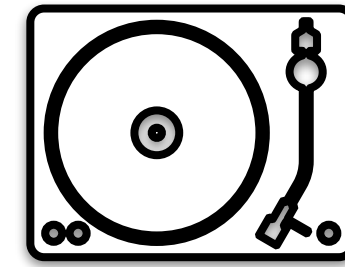
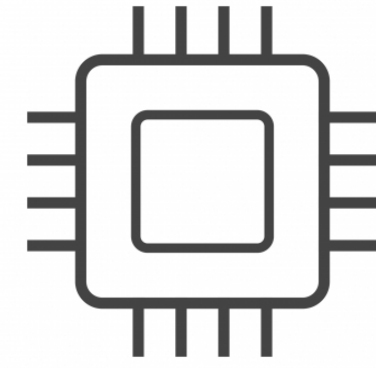


buffer 

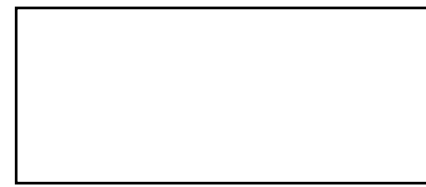
partial compaction



log-structured merge-tree



buffer



partial compaction

L1



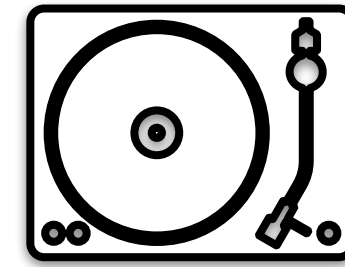
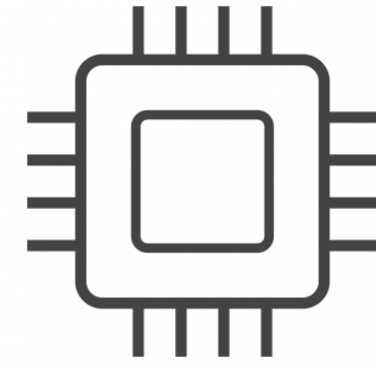
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



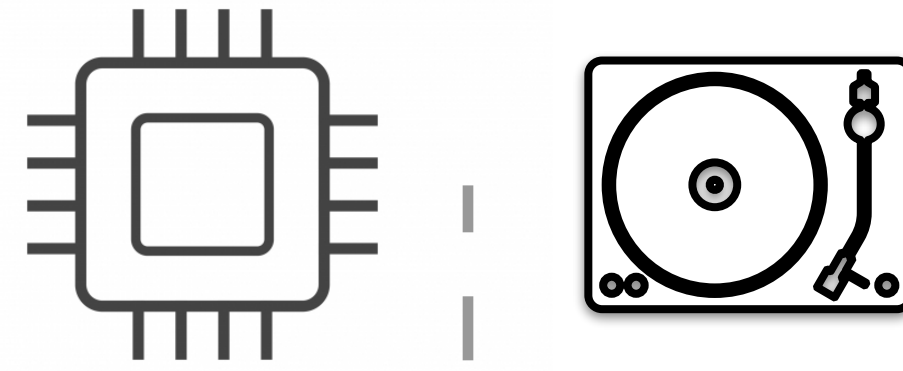
L2



L3

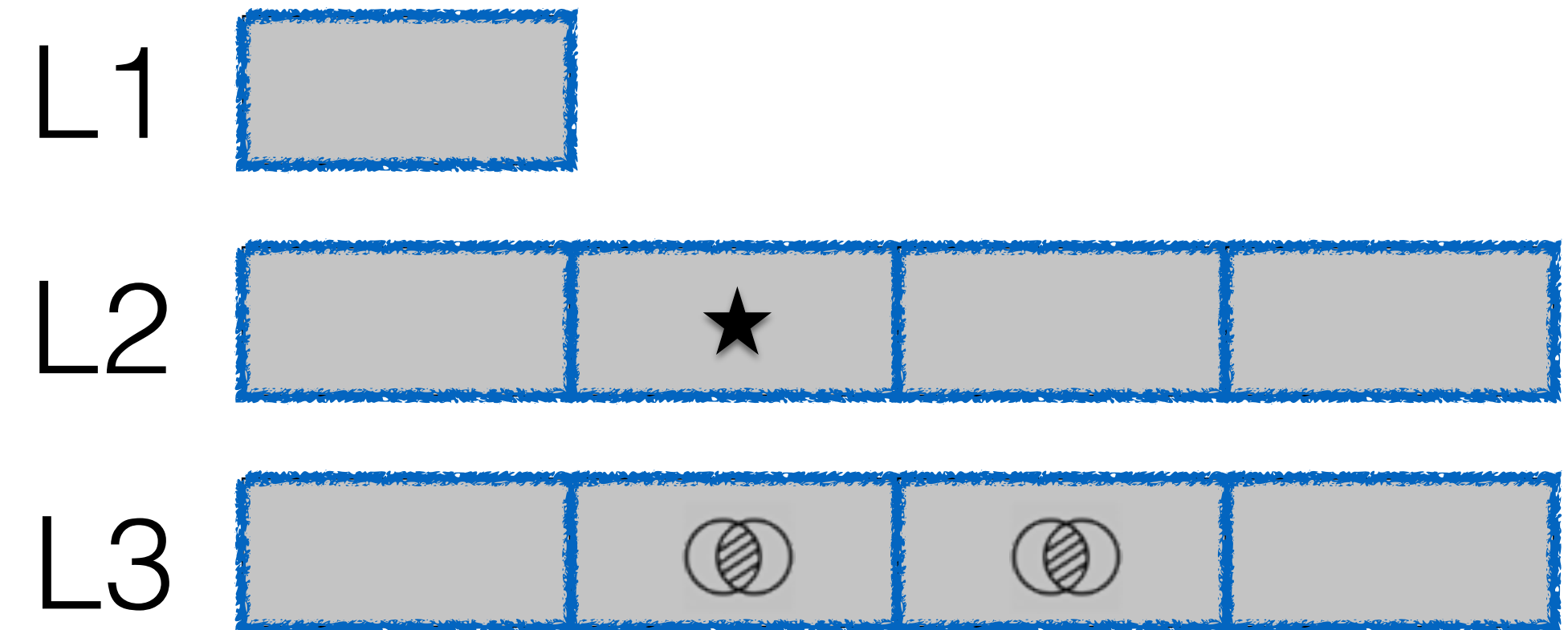


log-structured merge-tree

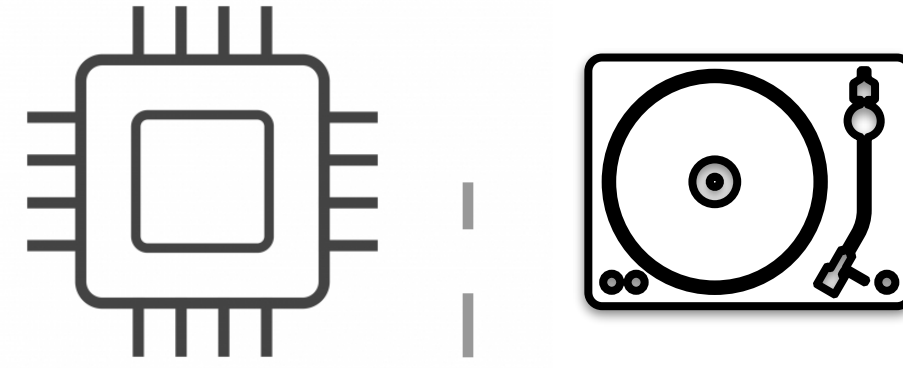


buffer 

partial compaction

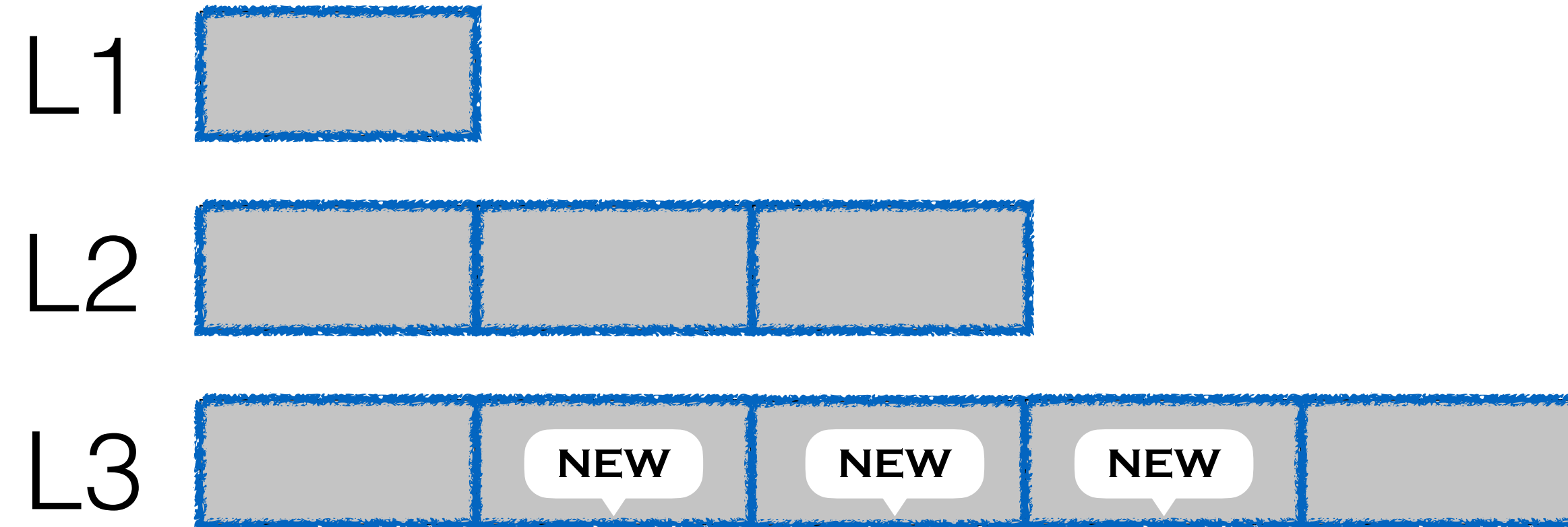


log-structured merge-tree



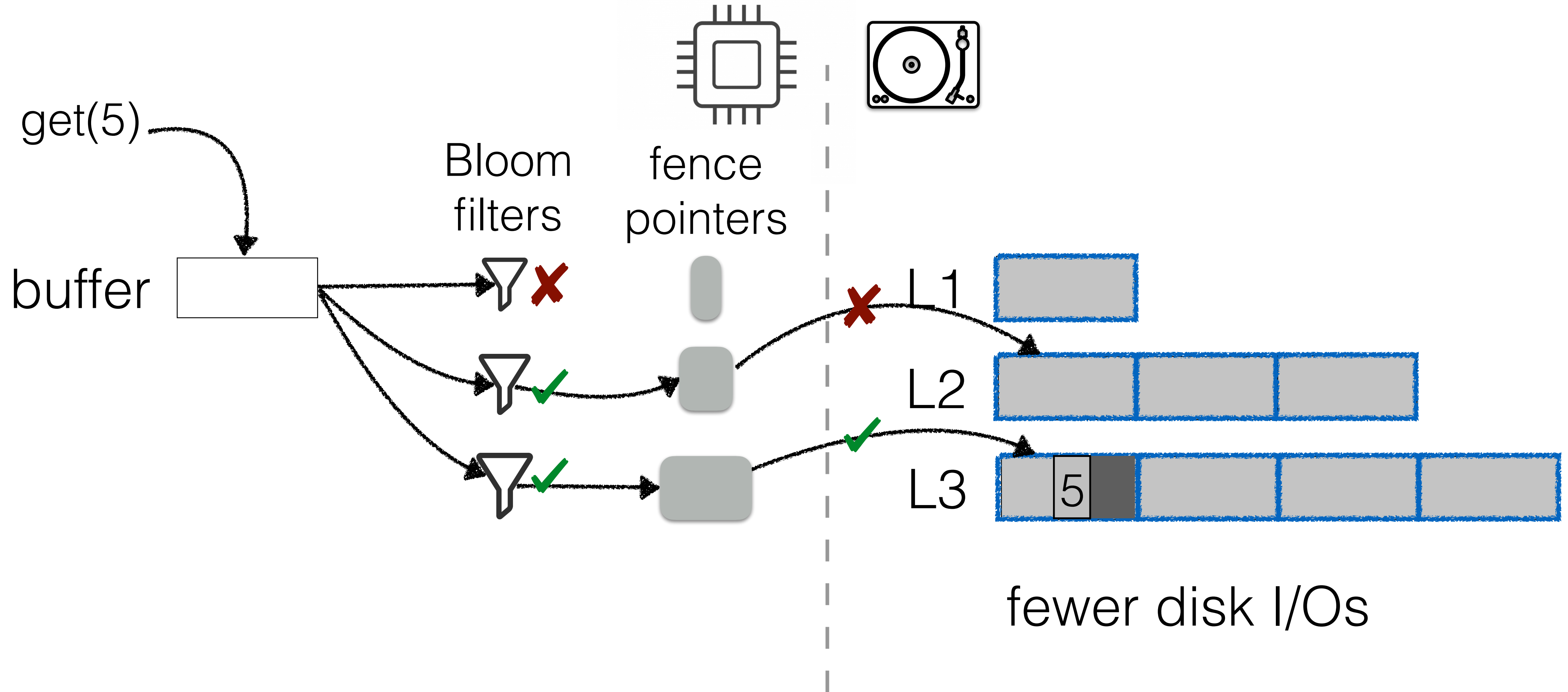
buffer 

partial compaction



amortized compaction cost

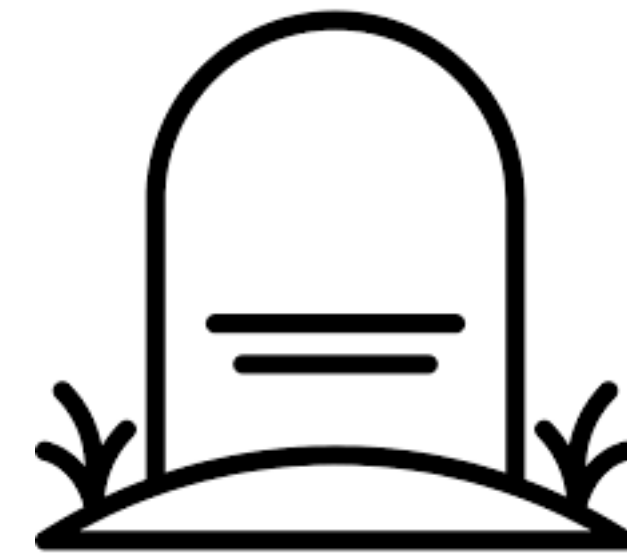
log-structured merge-tree



Now, let's talk about deletes!

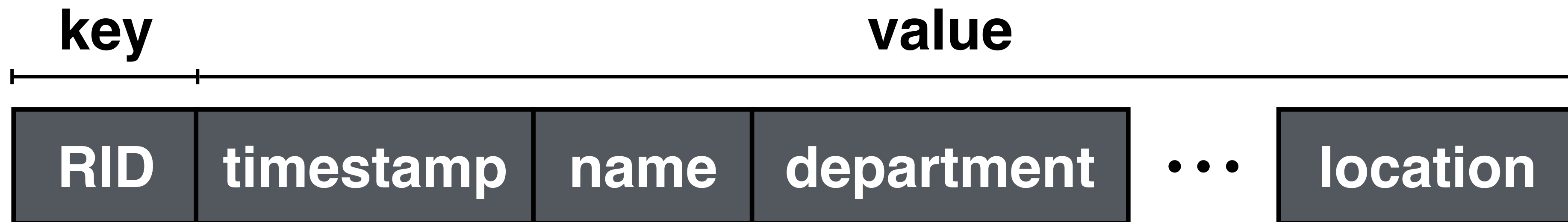
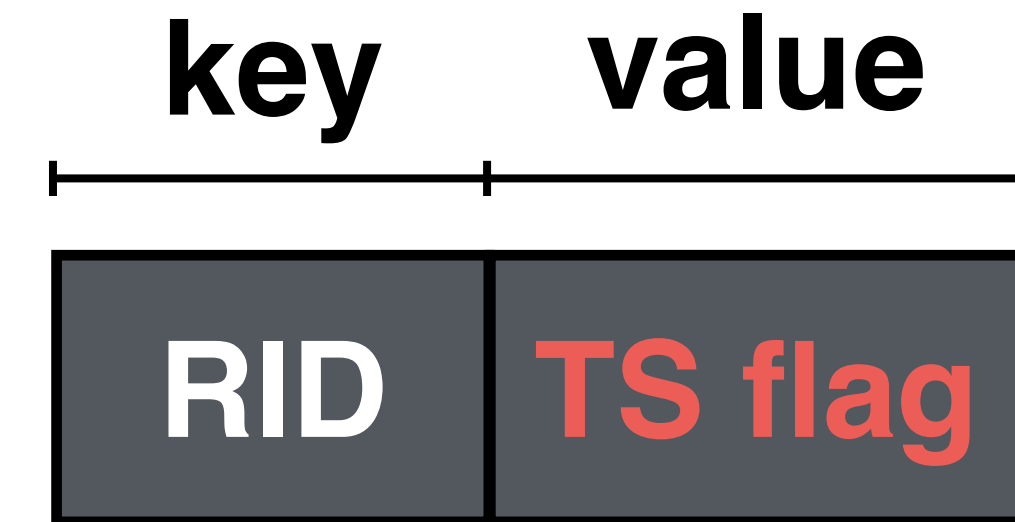
deletes in LSM-tree

delete := insert tombstone

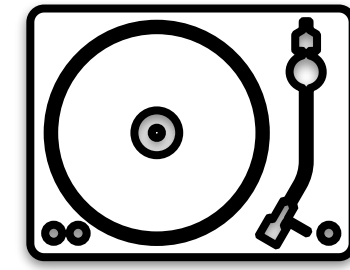
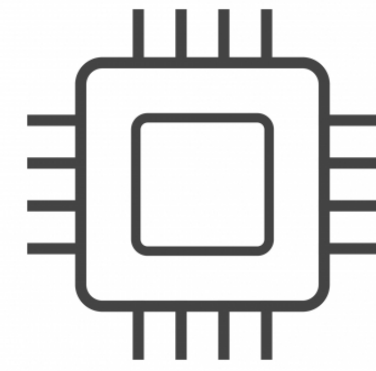
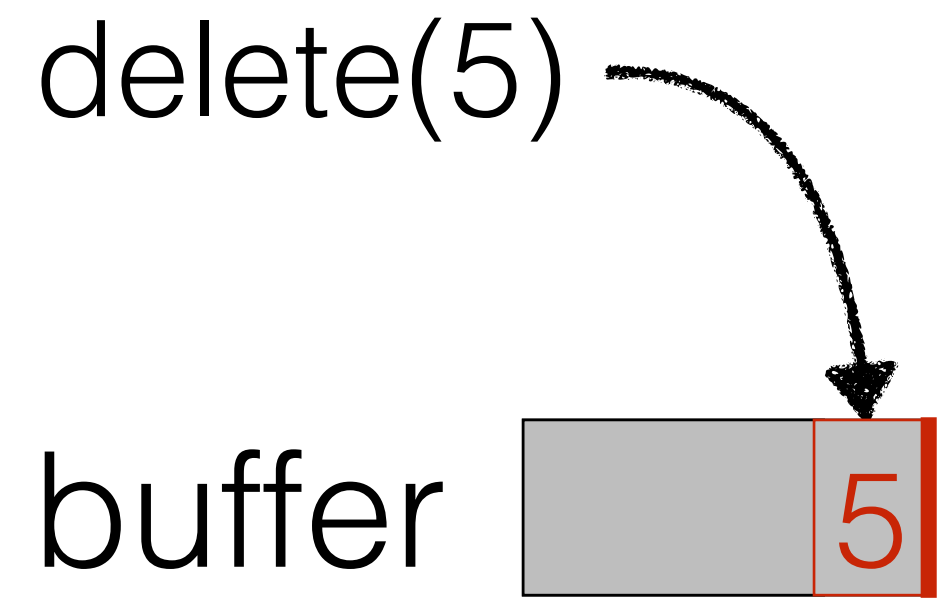


deletes in LSM-tree

delete := insert tombstone



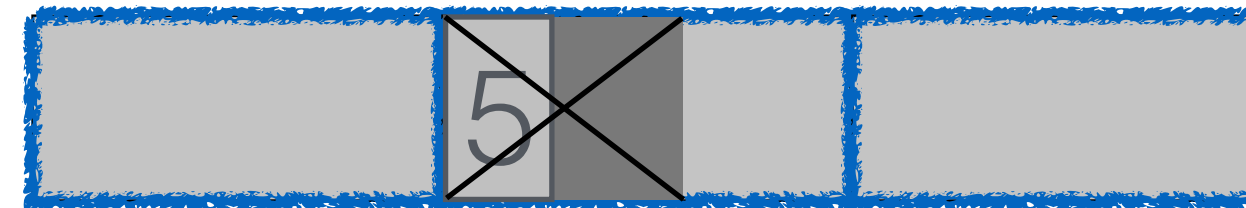
deletes in LSM-tree



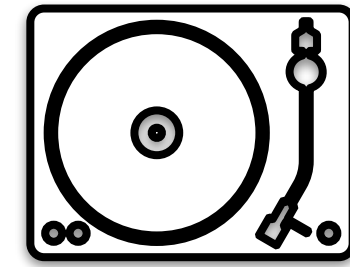
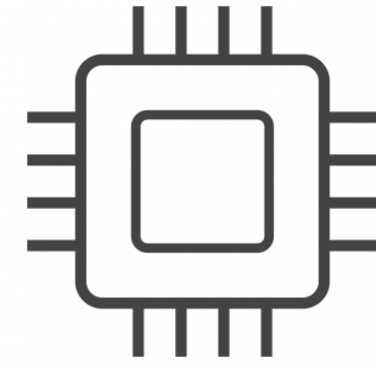
L1

L2

L3



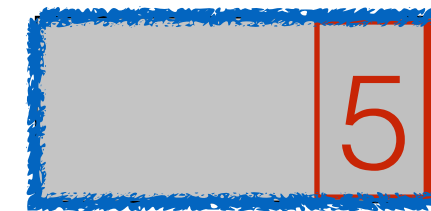
deletes in LSM-tree



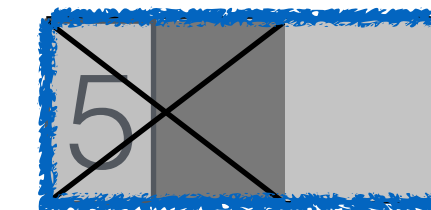
buffer



L1



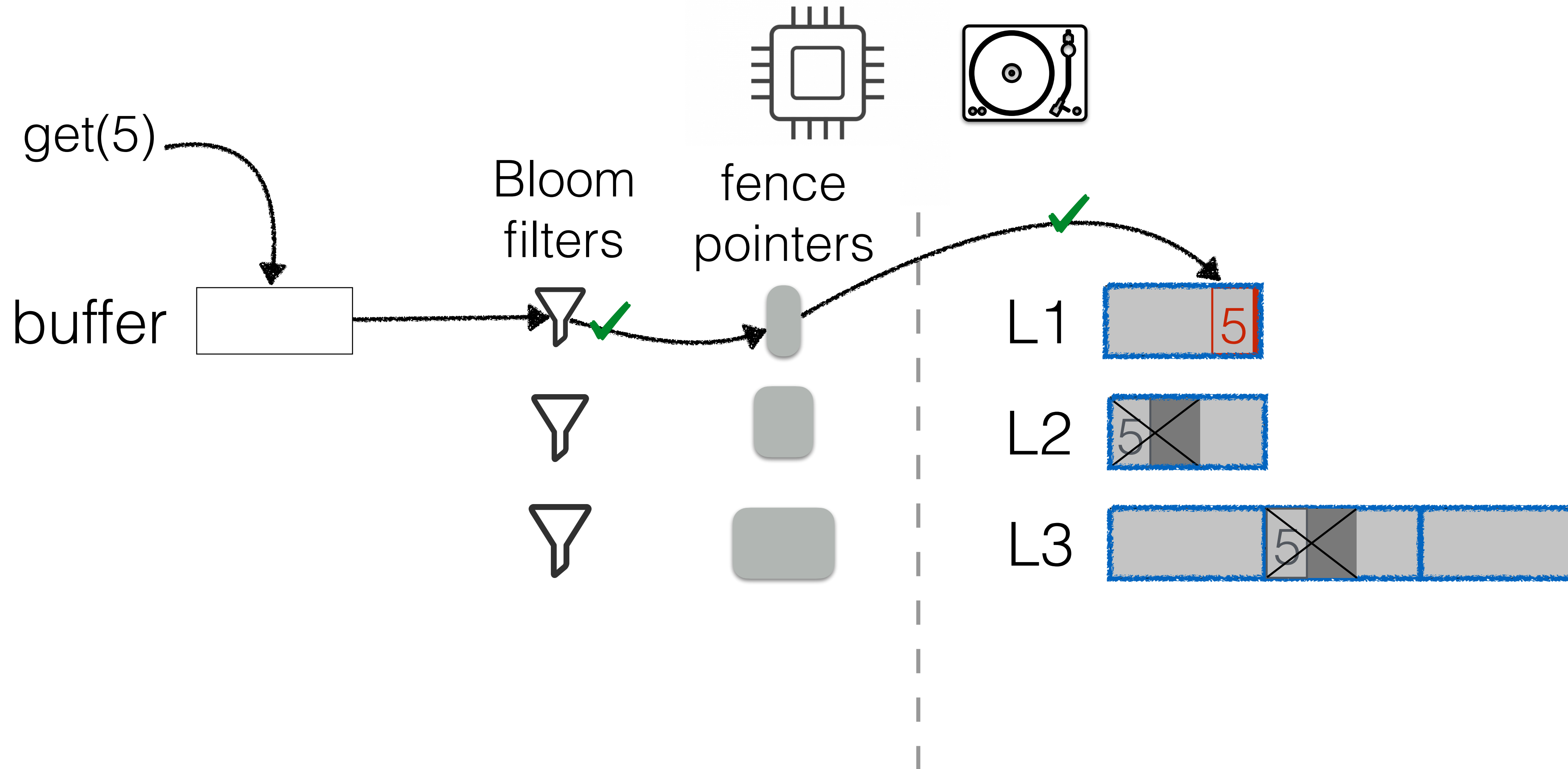
L2



L3



deletes in LSM-tree

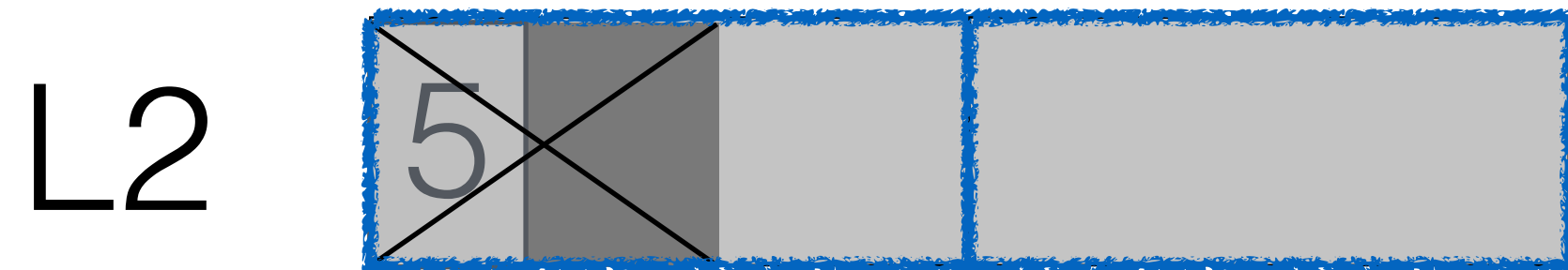


the problems



out-of-place deletes

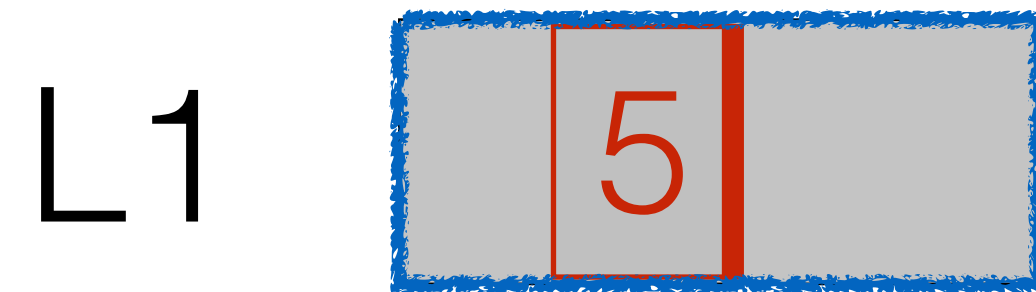
out-of-place deletes



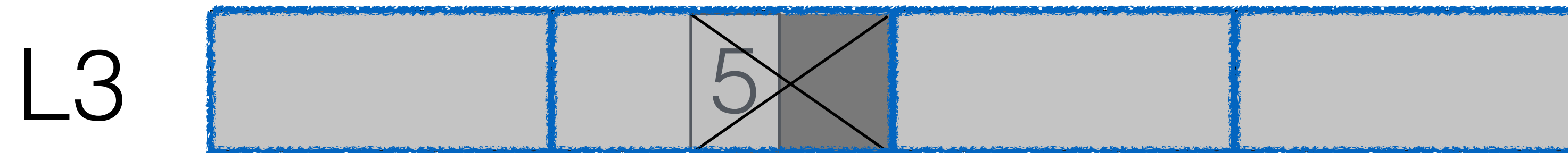
space amplification ✖



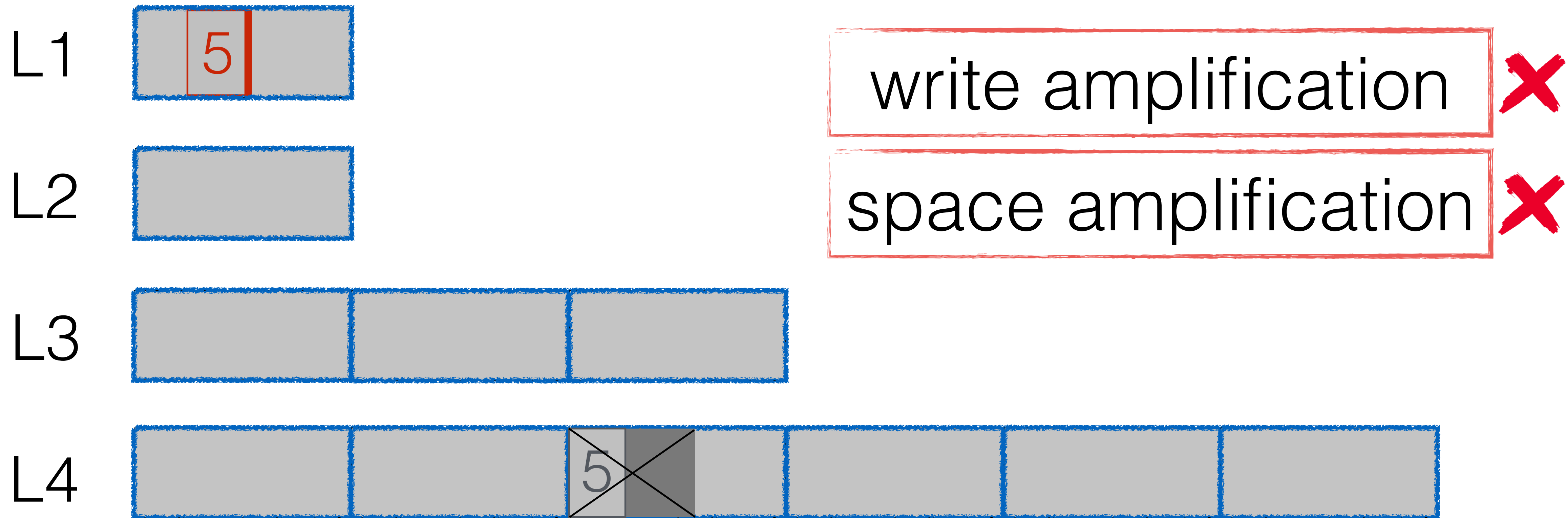
out-of-place deletes



space amplification ✖

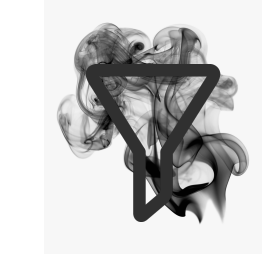


out-of-place deletes



out-of-place deletes

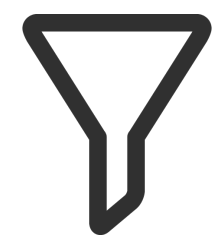
Bloom filters



L1



L2



L3



L4



poor read perf. ❌

write amplification ❌

space amplification ❌

the problems

poor read perf.

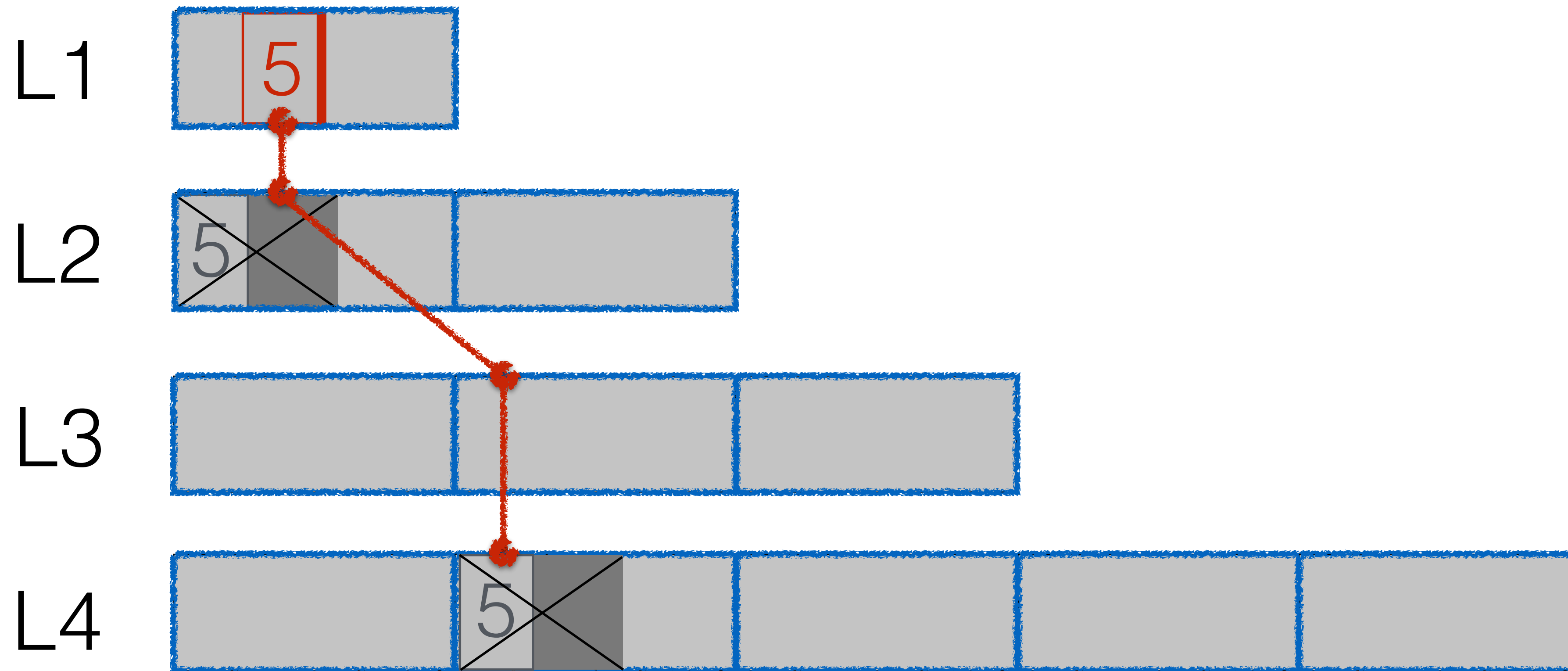
write amplification

space amplification



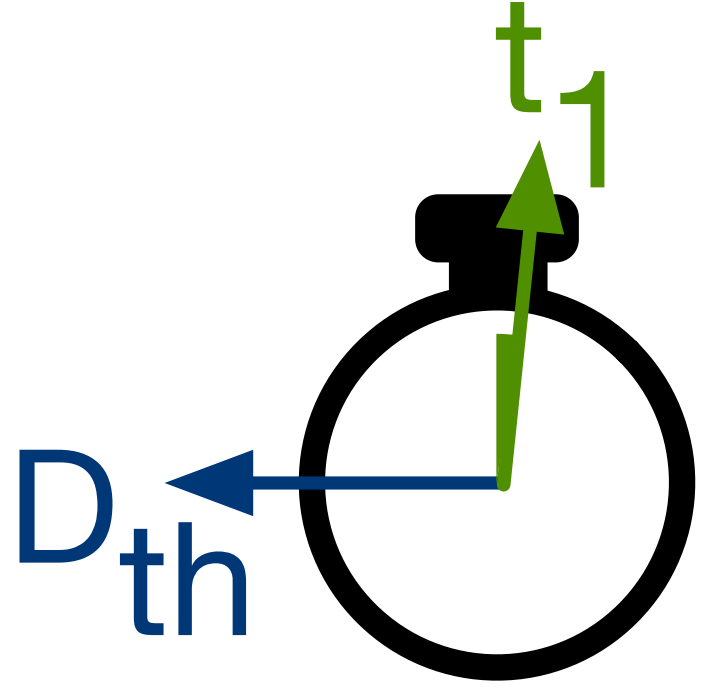
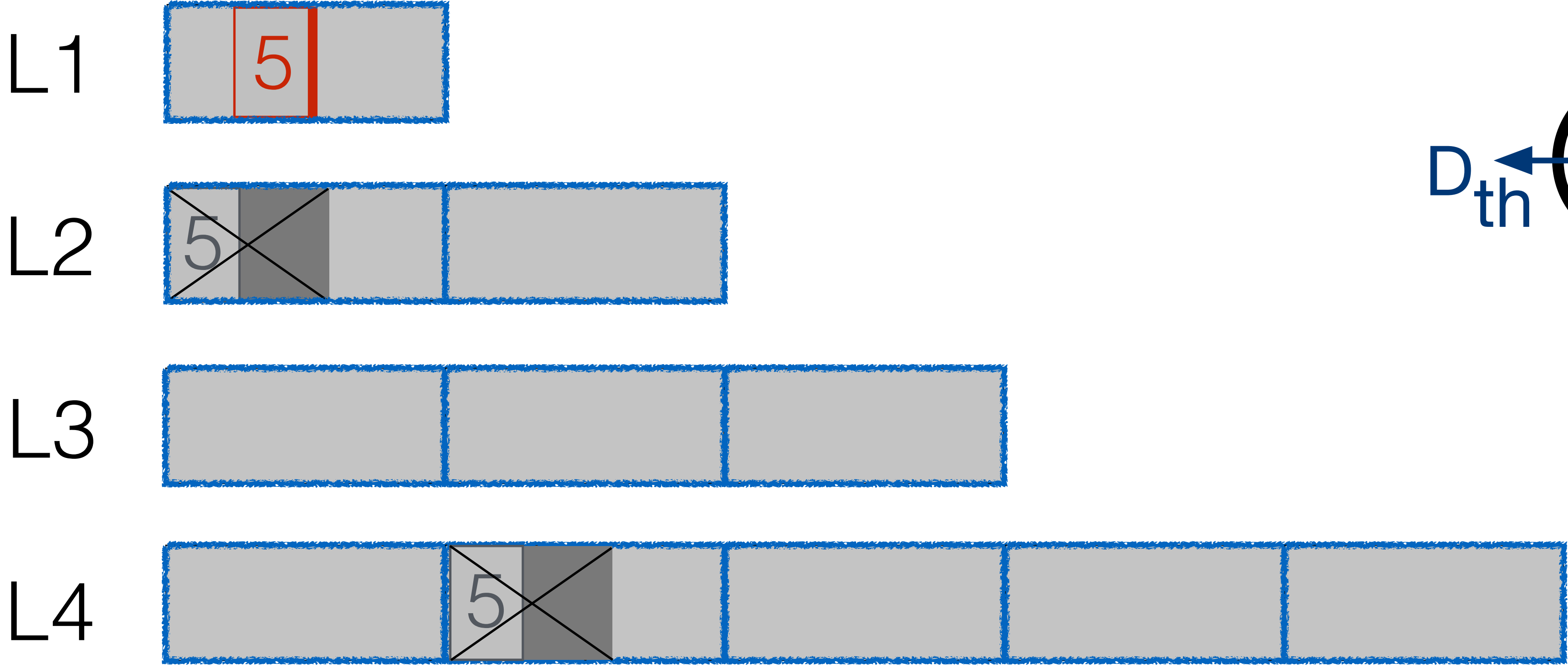
delete persistence latency

delete persistence latency



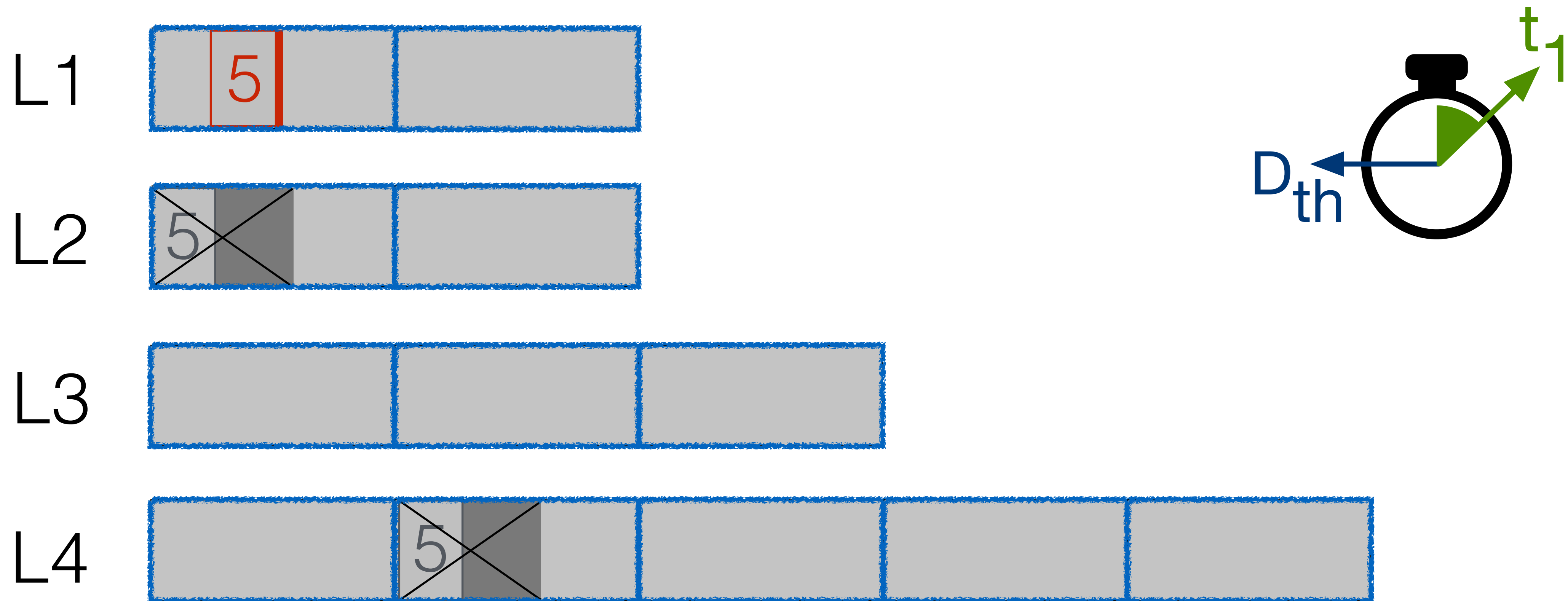
delete persistence latency

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



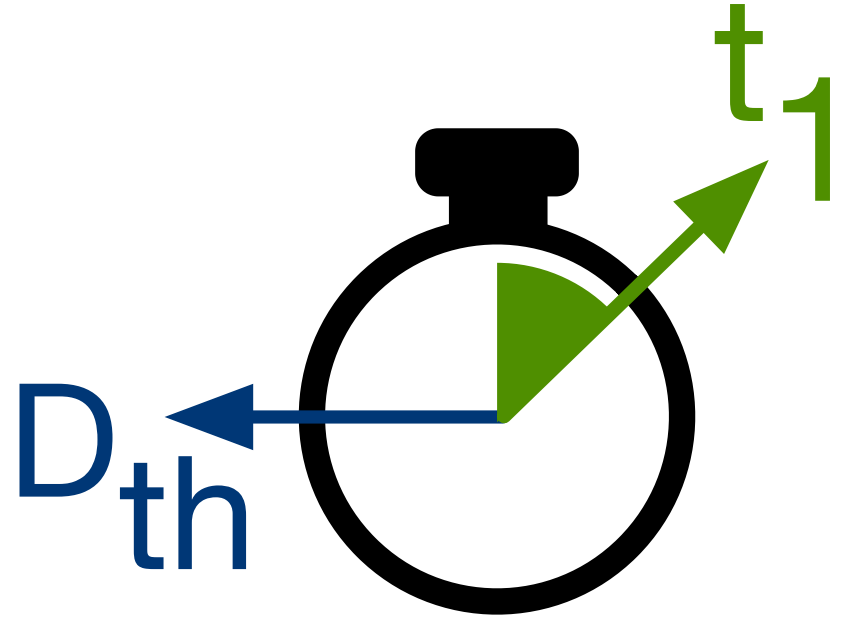
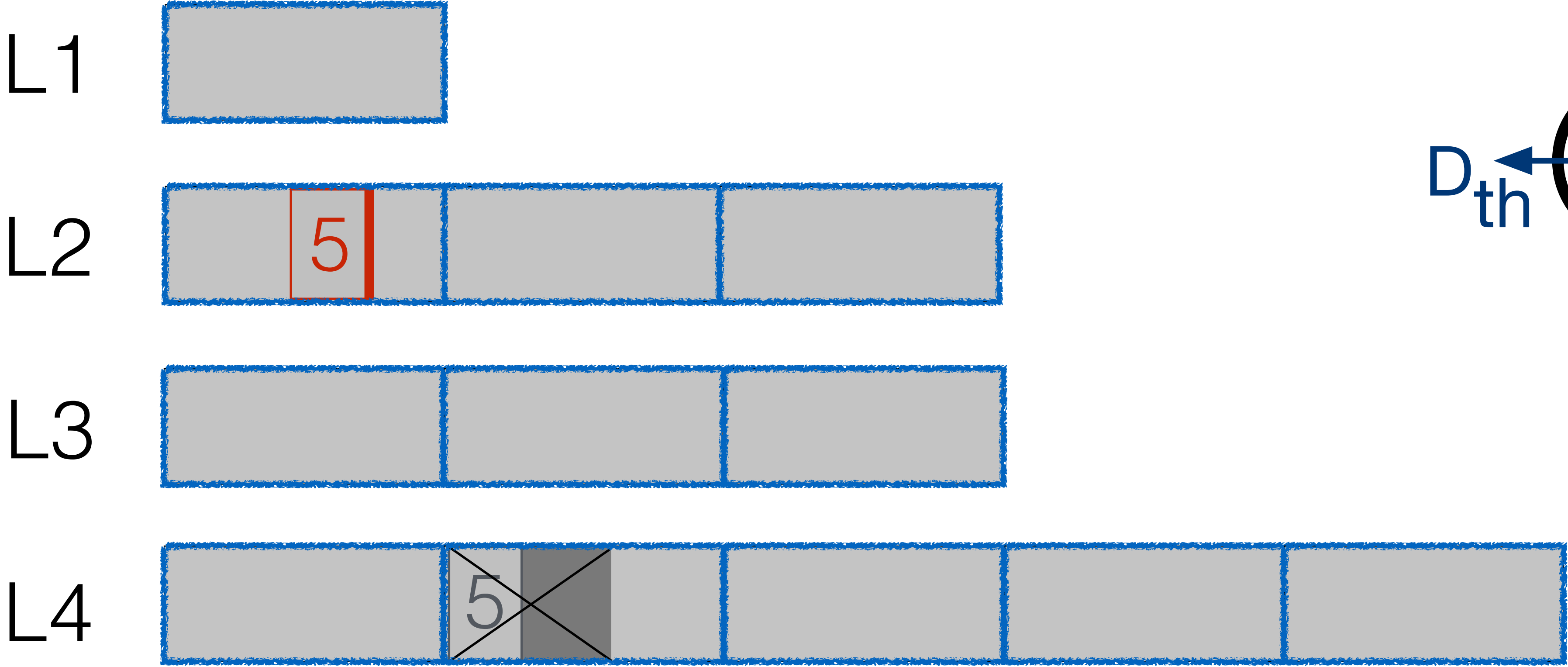
delete persistence latency

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



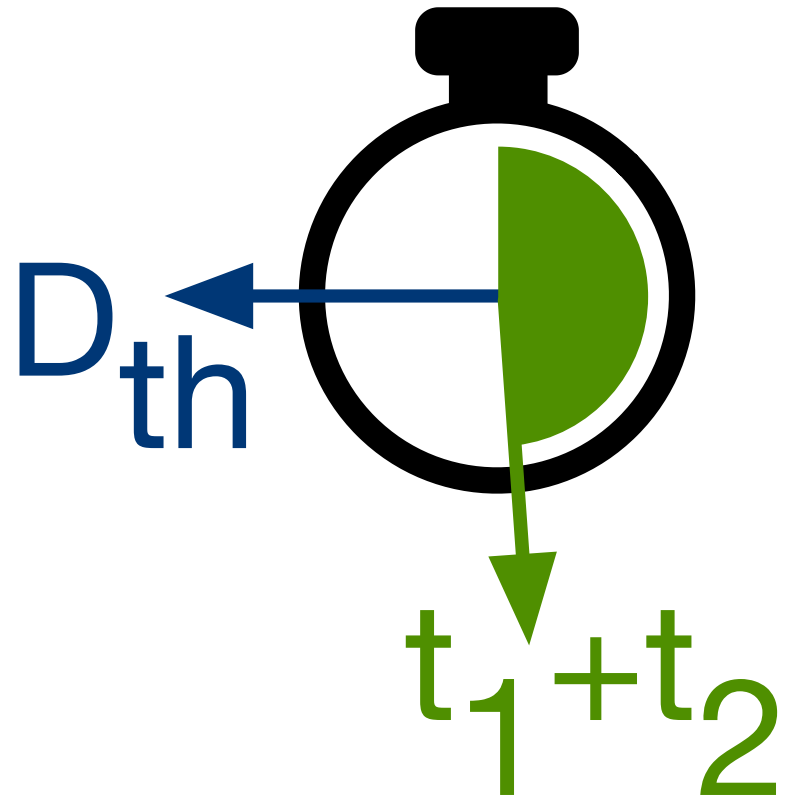
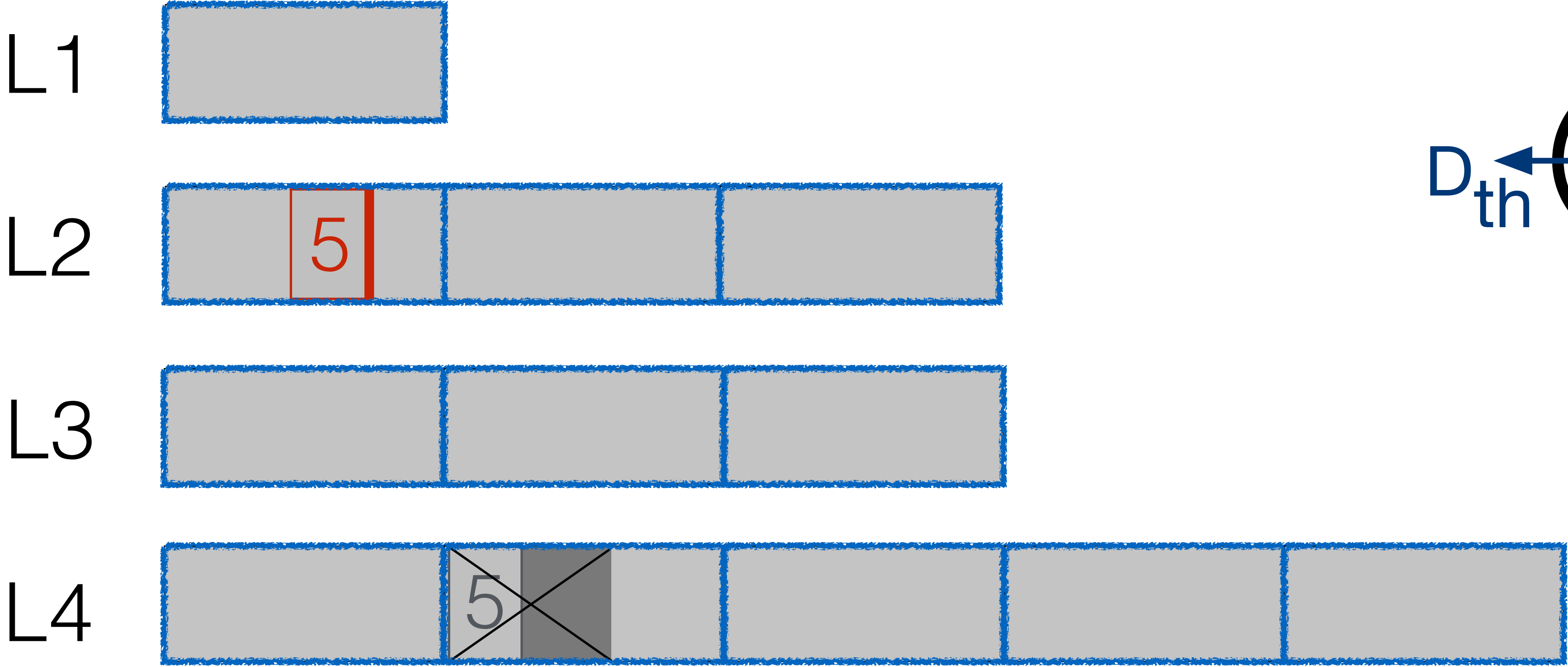
delete persistence latency

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



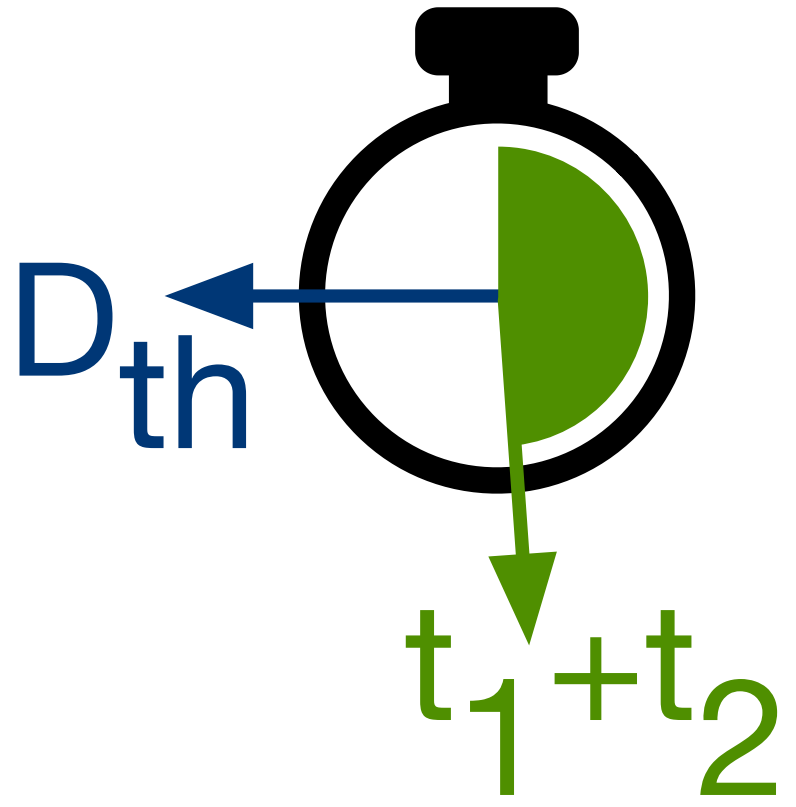
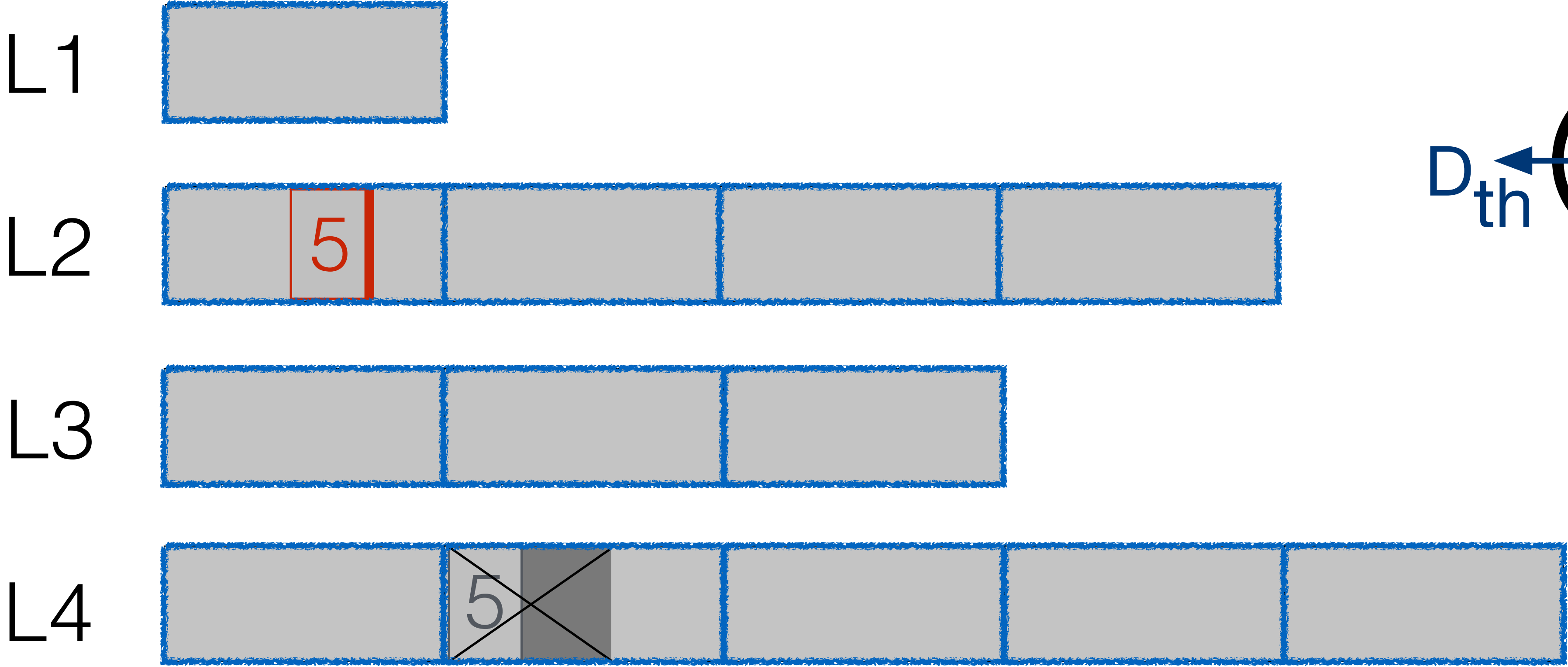
delete persistence latency

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



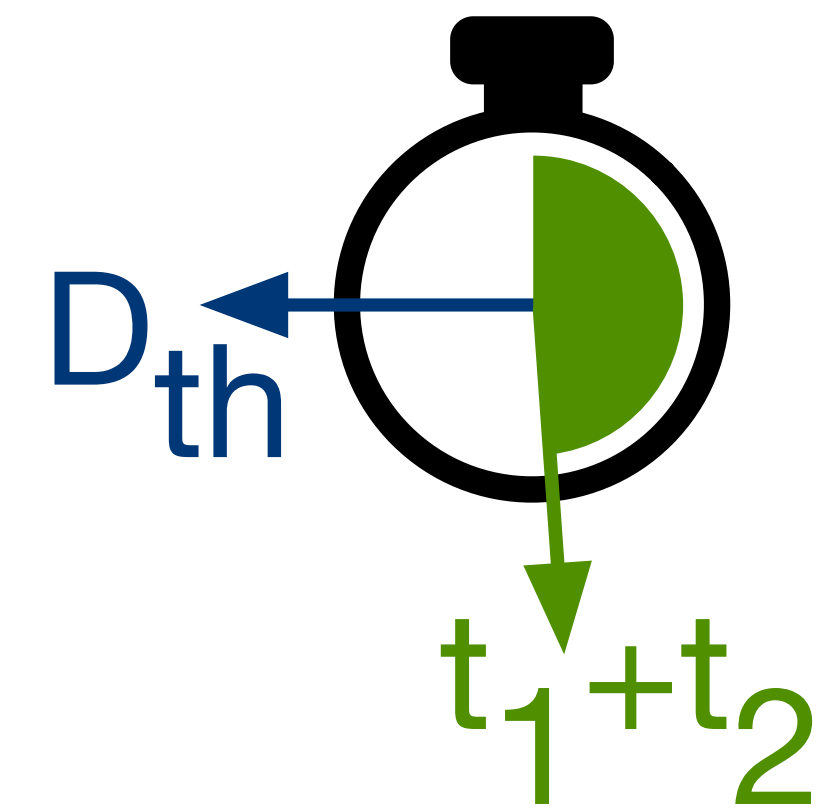
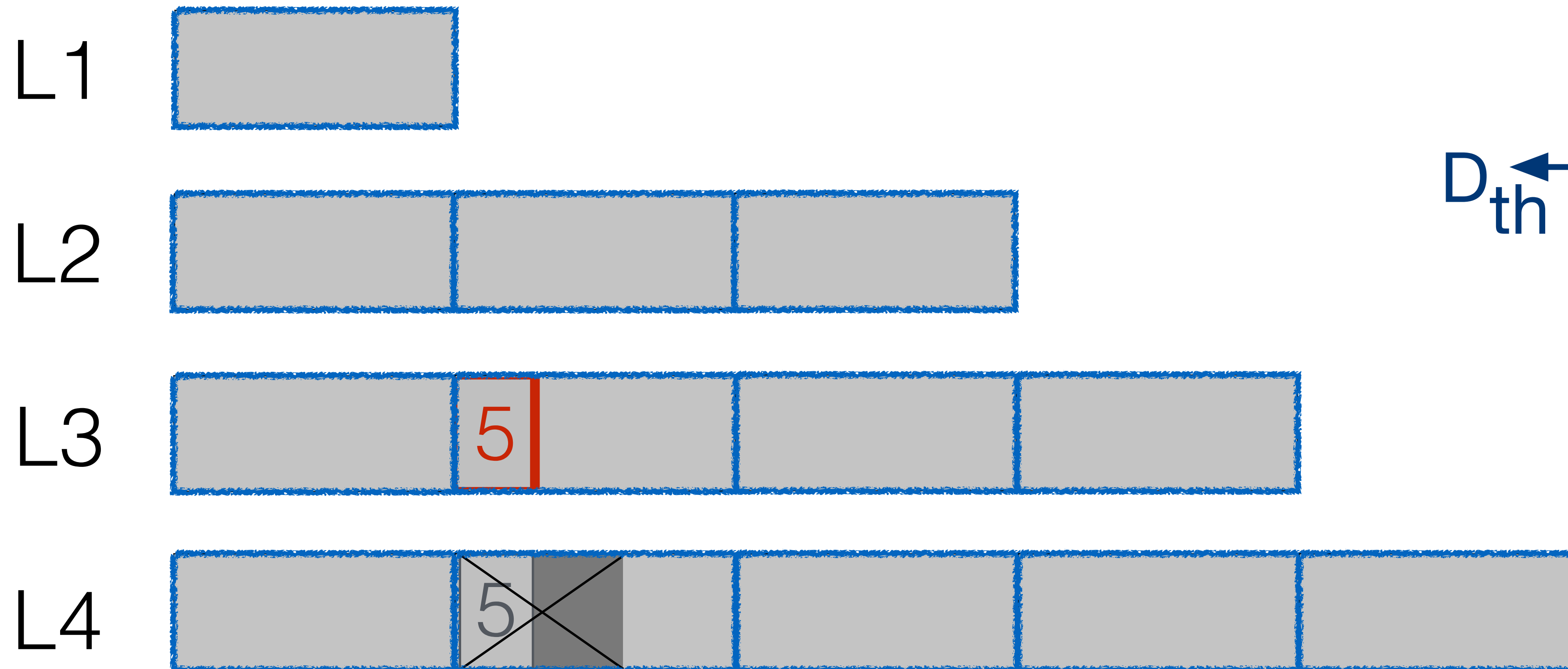
delete persistence latency

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



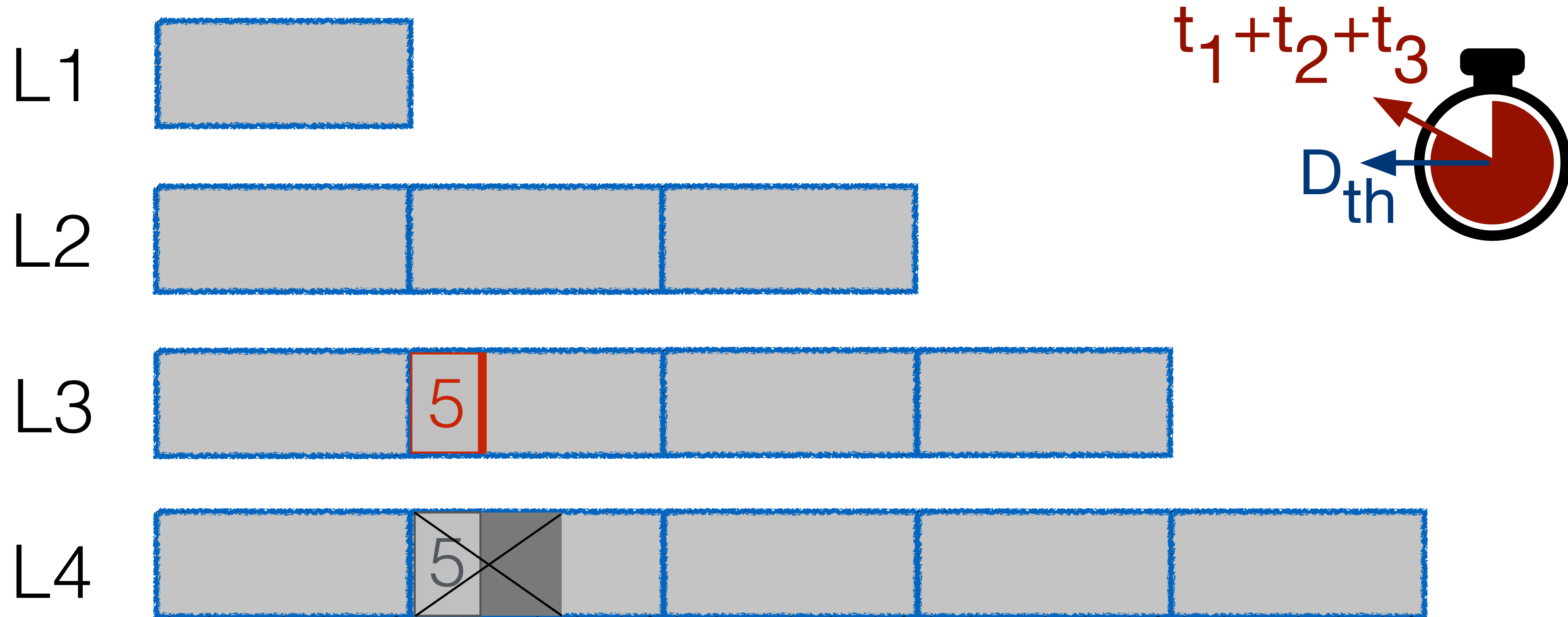
delete persistence latency

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



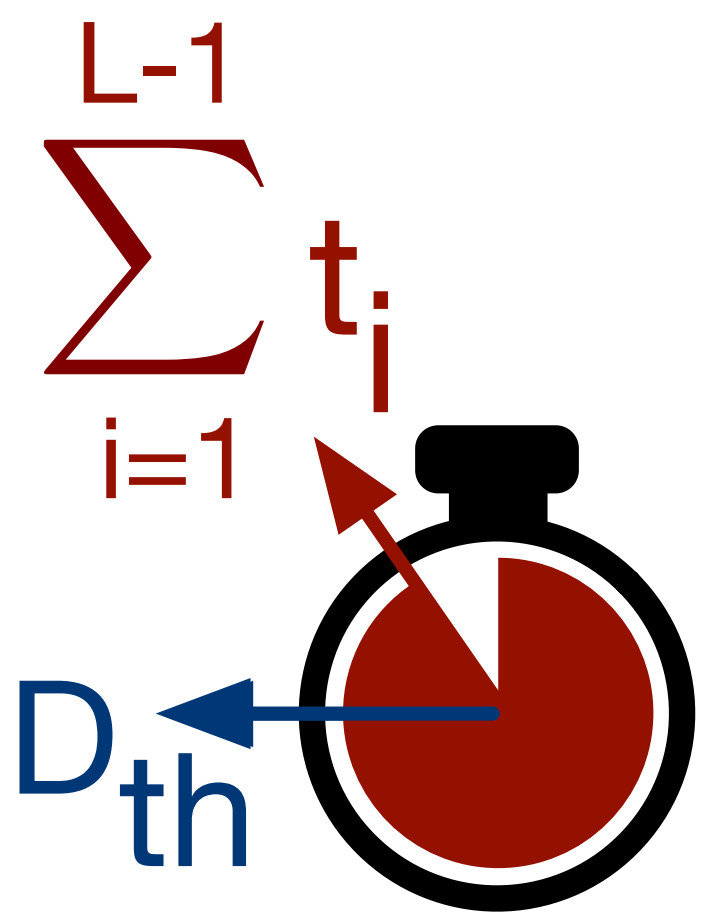
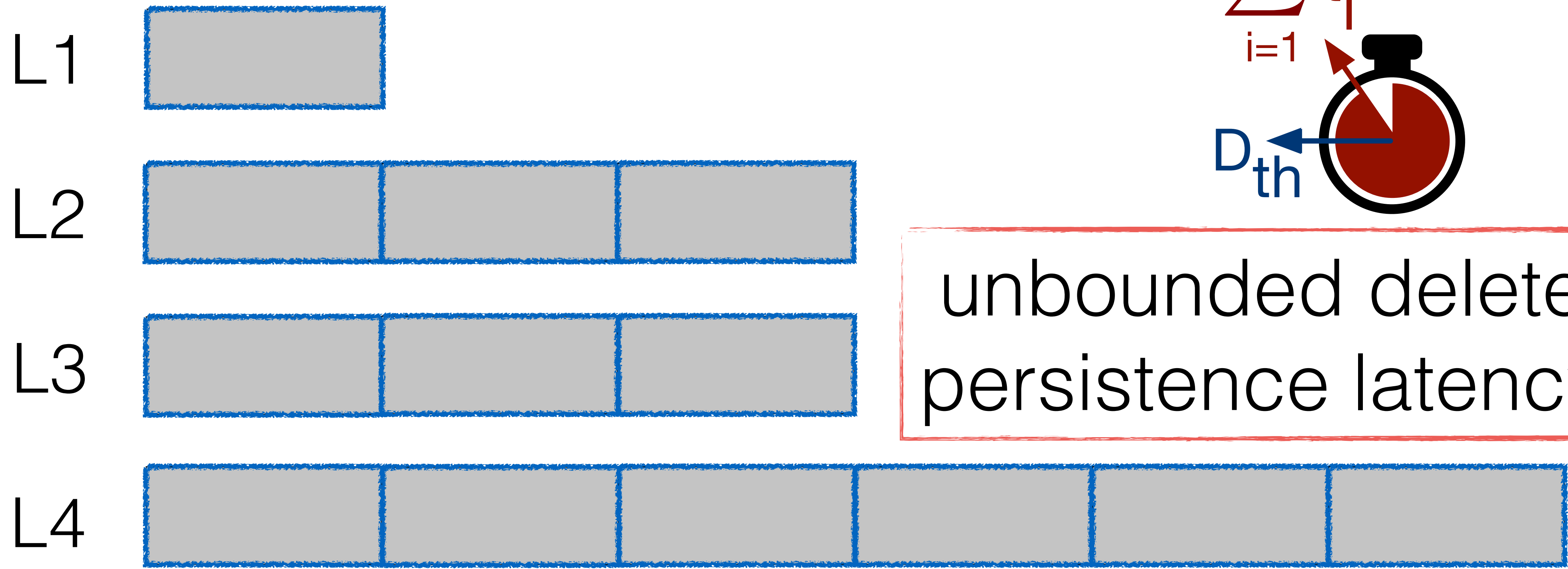
delete persistence latency

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



delete persistence latency

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

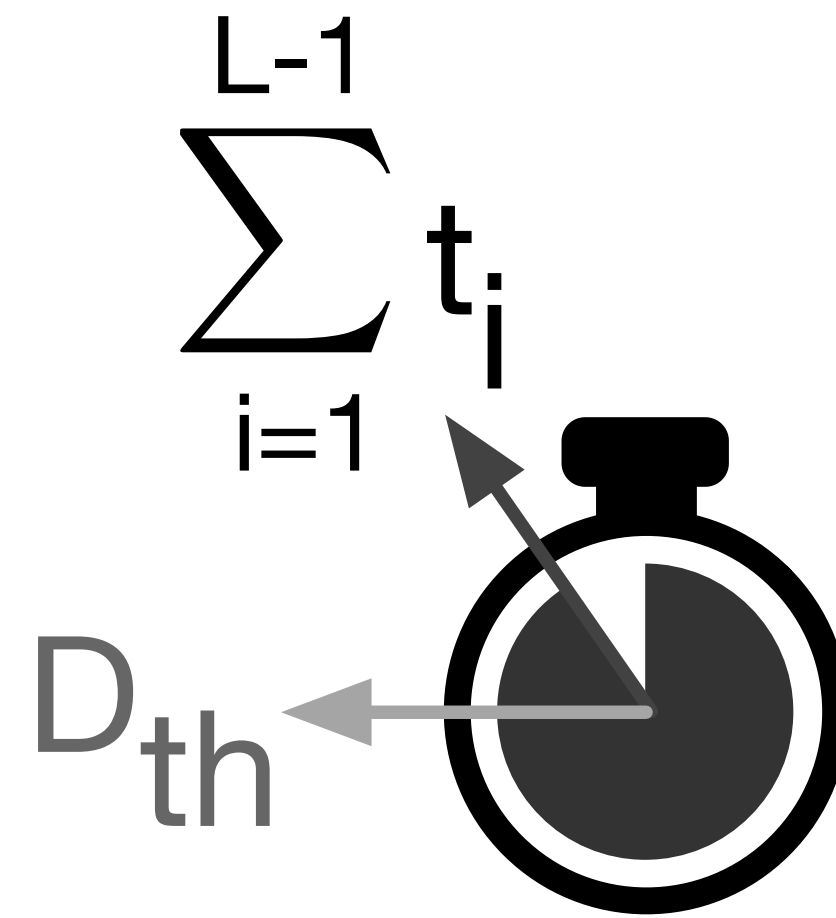


the problems

poor read perf.

write amplification

space amplification



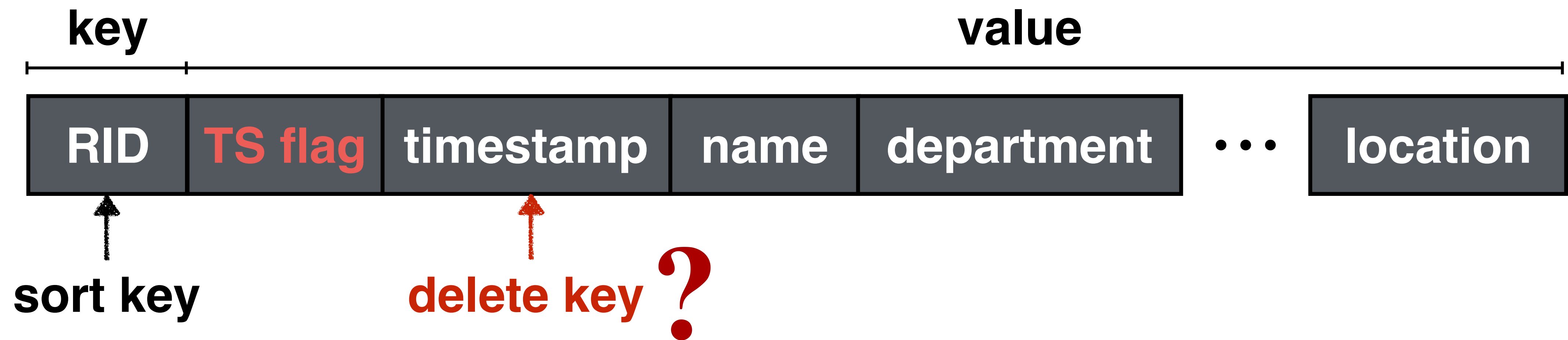
unbounded delete
persistence latency



deletes on a secondary attribute

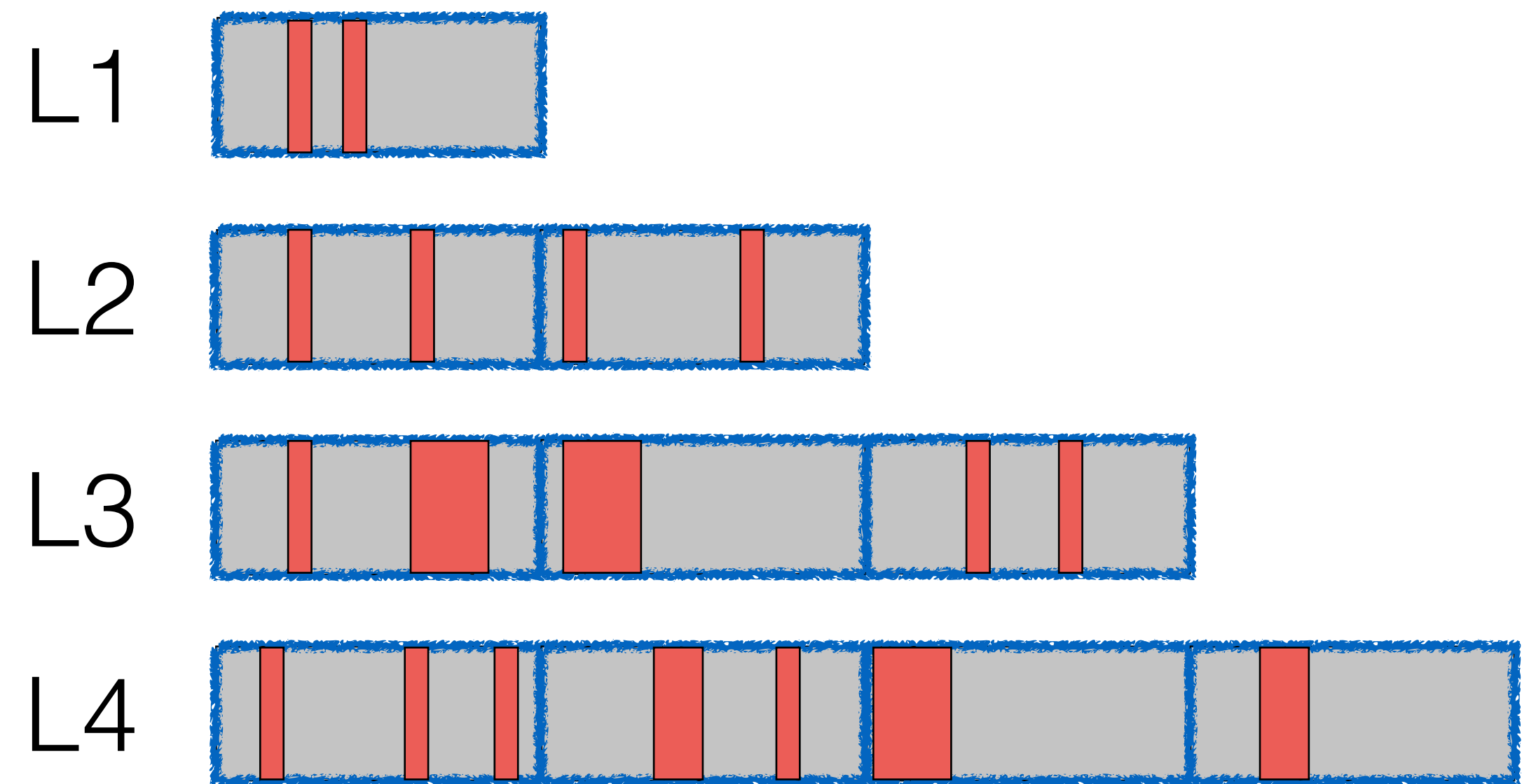
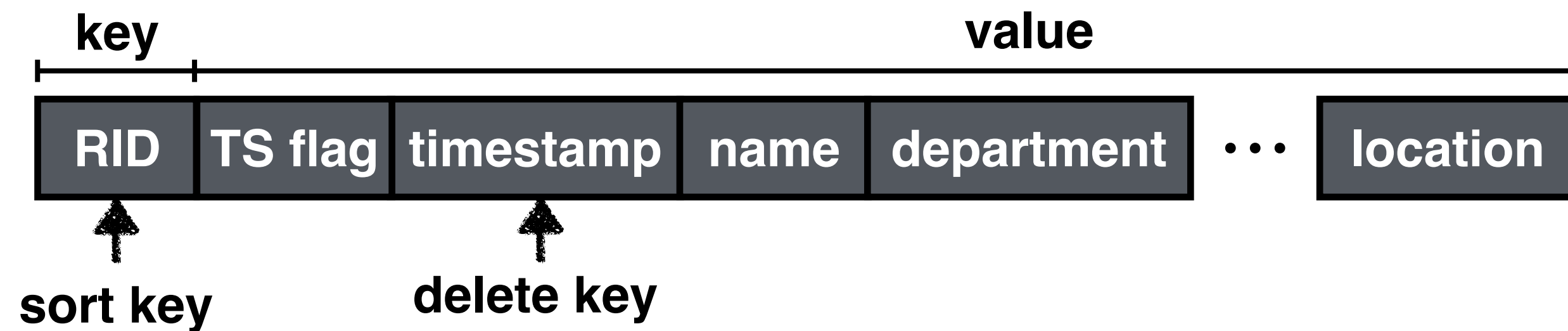
deletes on a secondary attribute

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



deletes on a secondary attribute

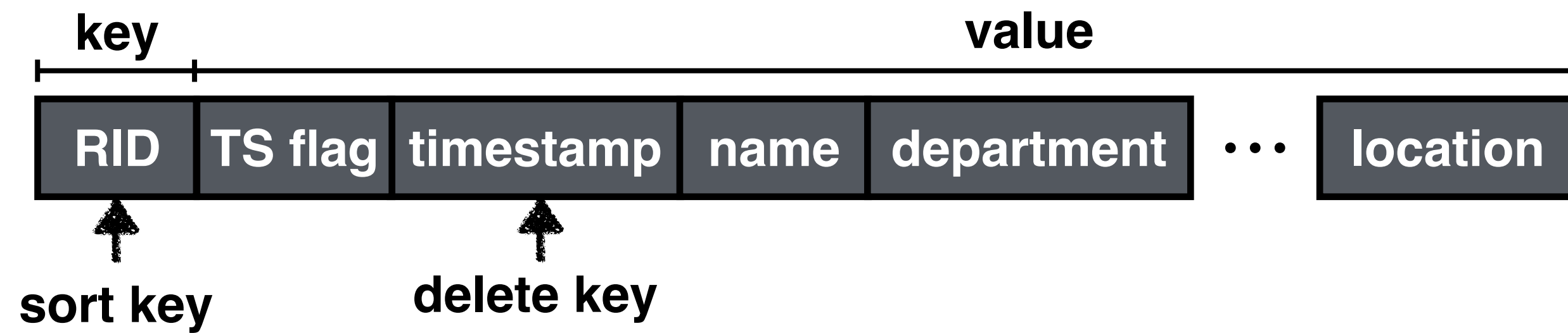
delete all entries older than: **D days**



scattered occurrences

deletes on a secondary attribute

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



L1

L2

L3

L4

latency spikes



superfluous I/Os

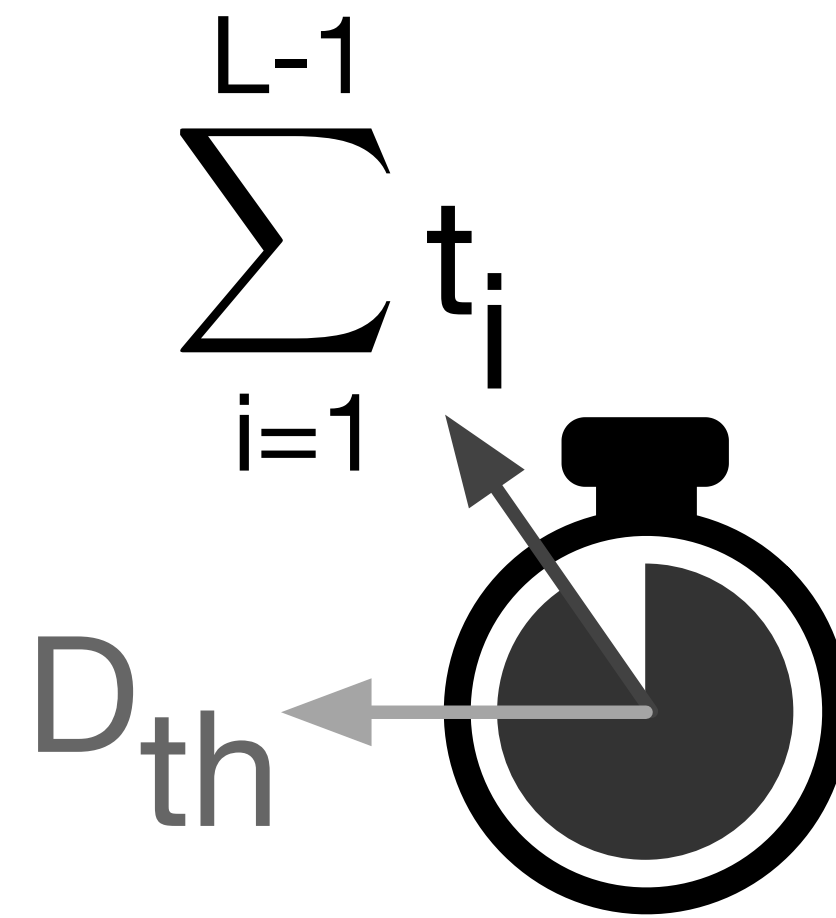


the problems

poor read perf.

write amplification

space amplification



unbounded delete
persistence latency

latency spikes

superfluous I/Os

the solution

poor read perf.

write amplification

space amplification

FADE

unbounded delete persistence latency

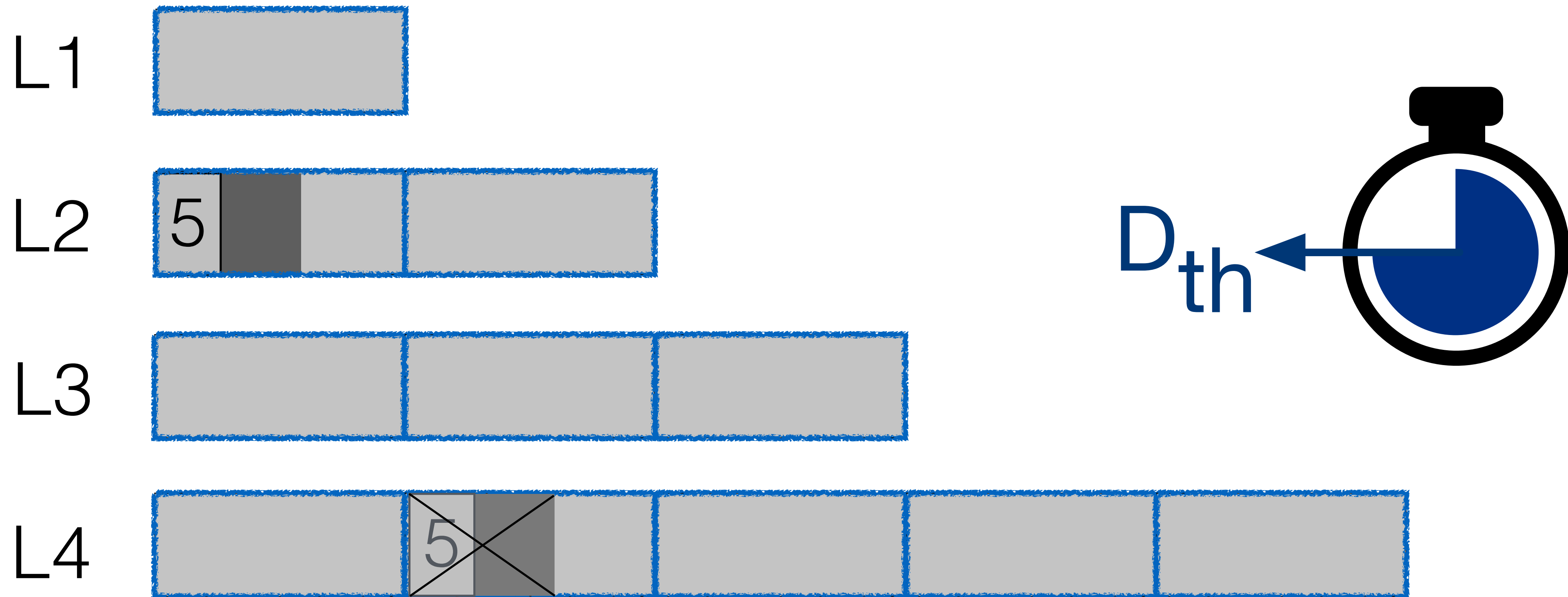


latency spikes

superfluous I/Os

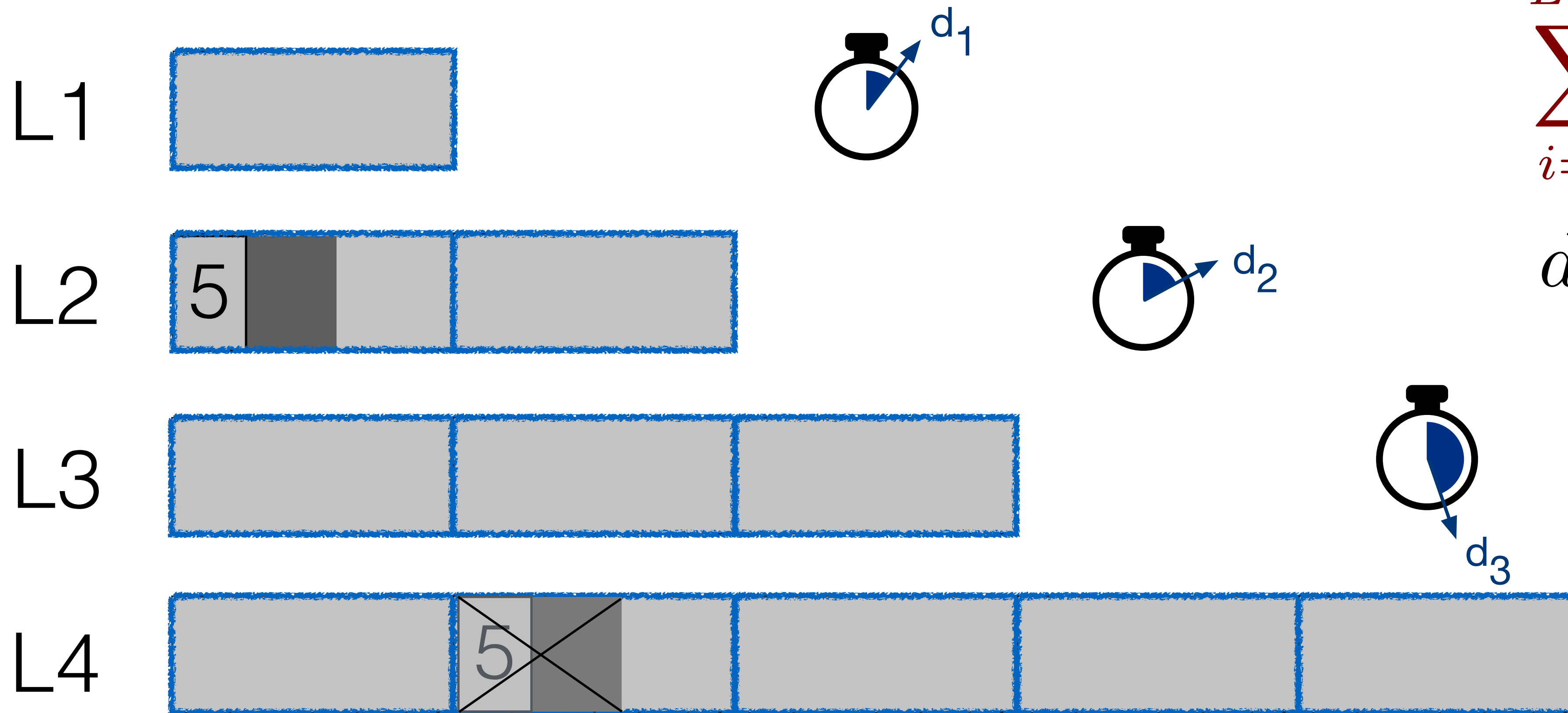
FAst DElete

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



FAst DElete

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

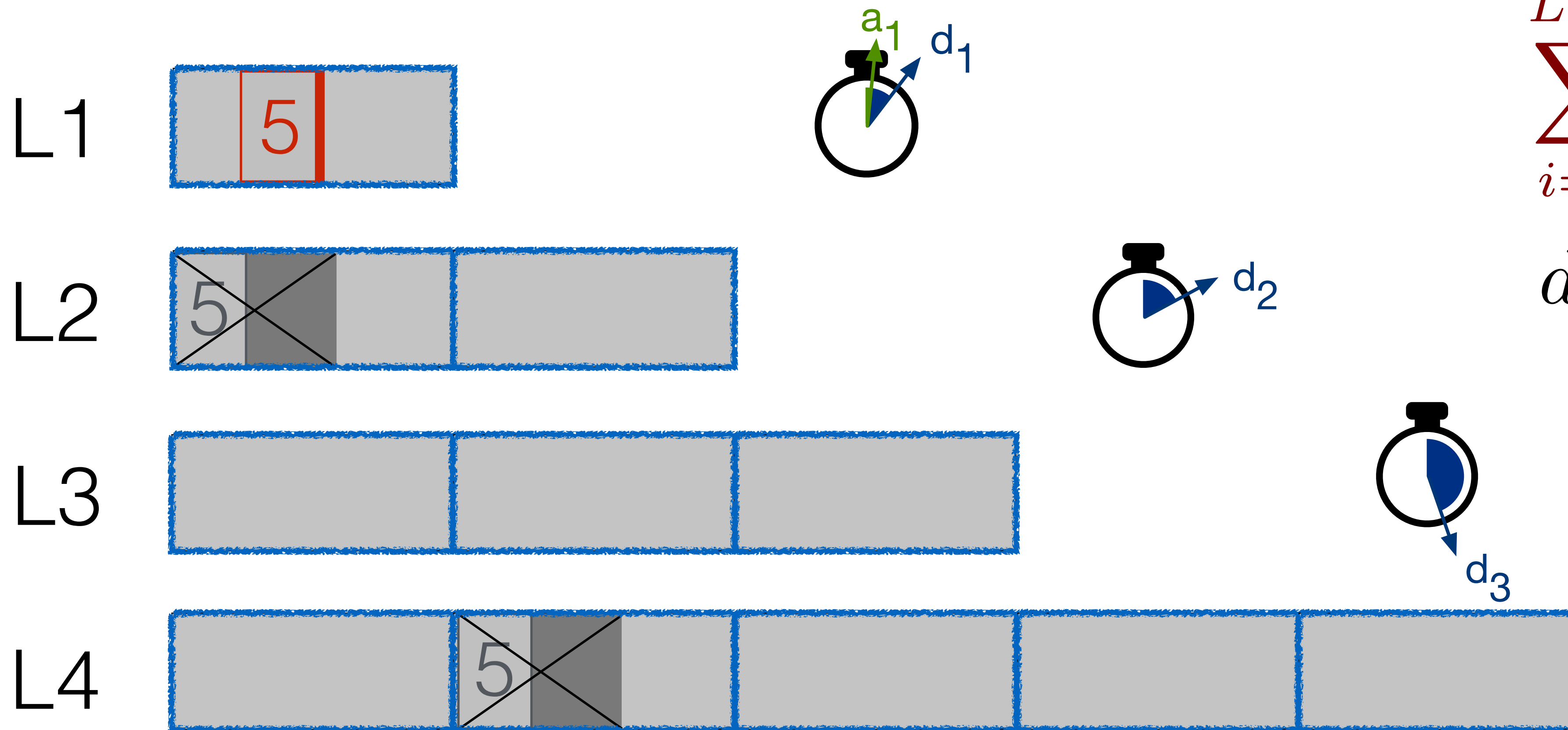


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

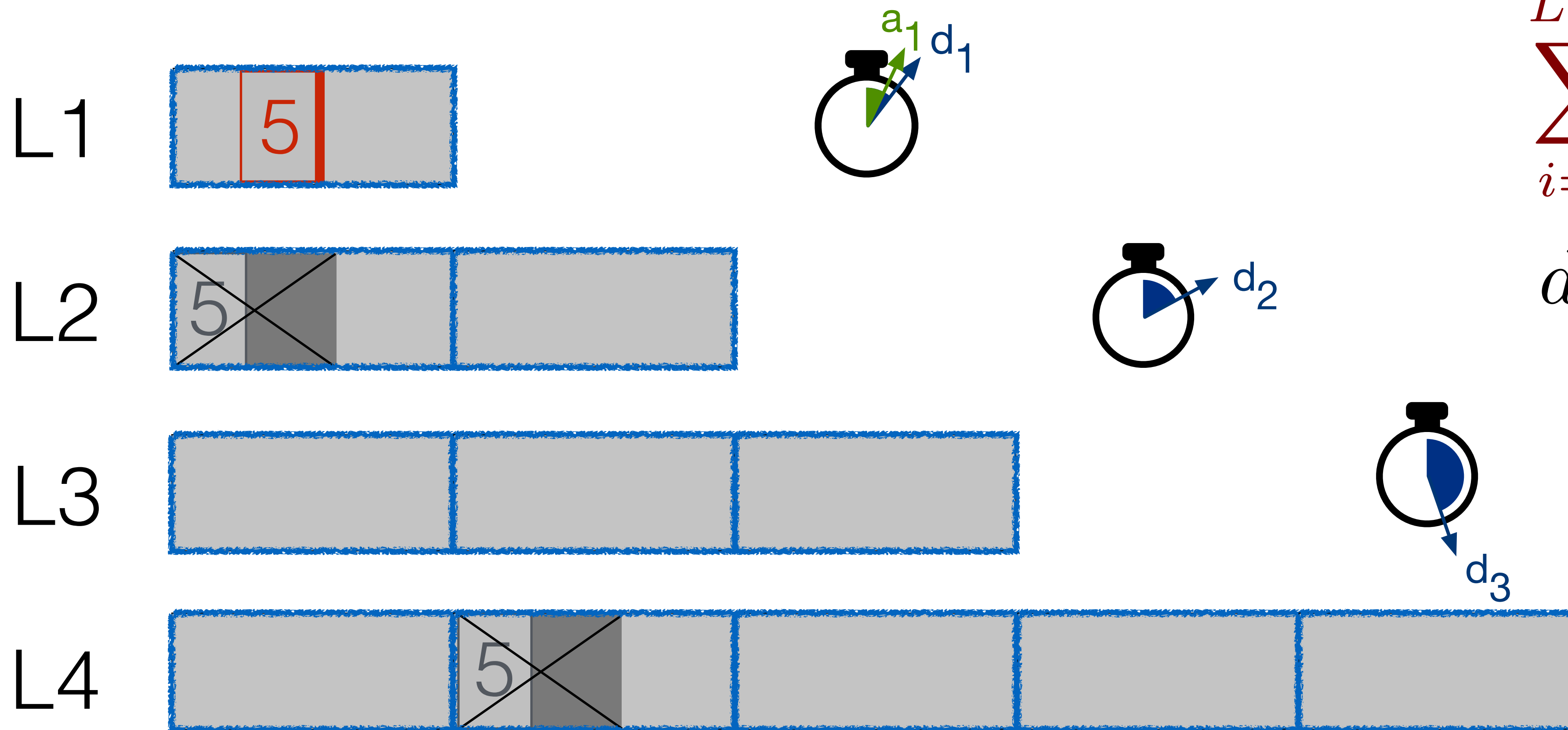


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

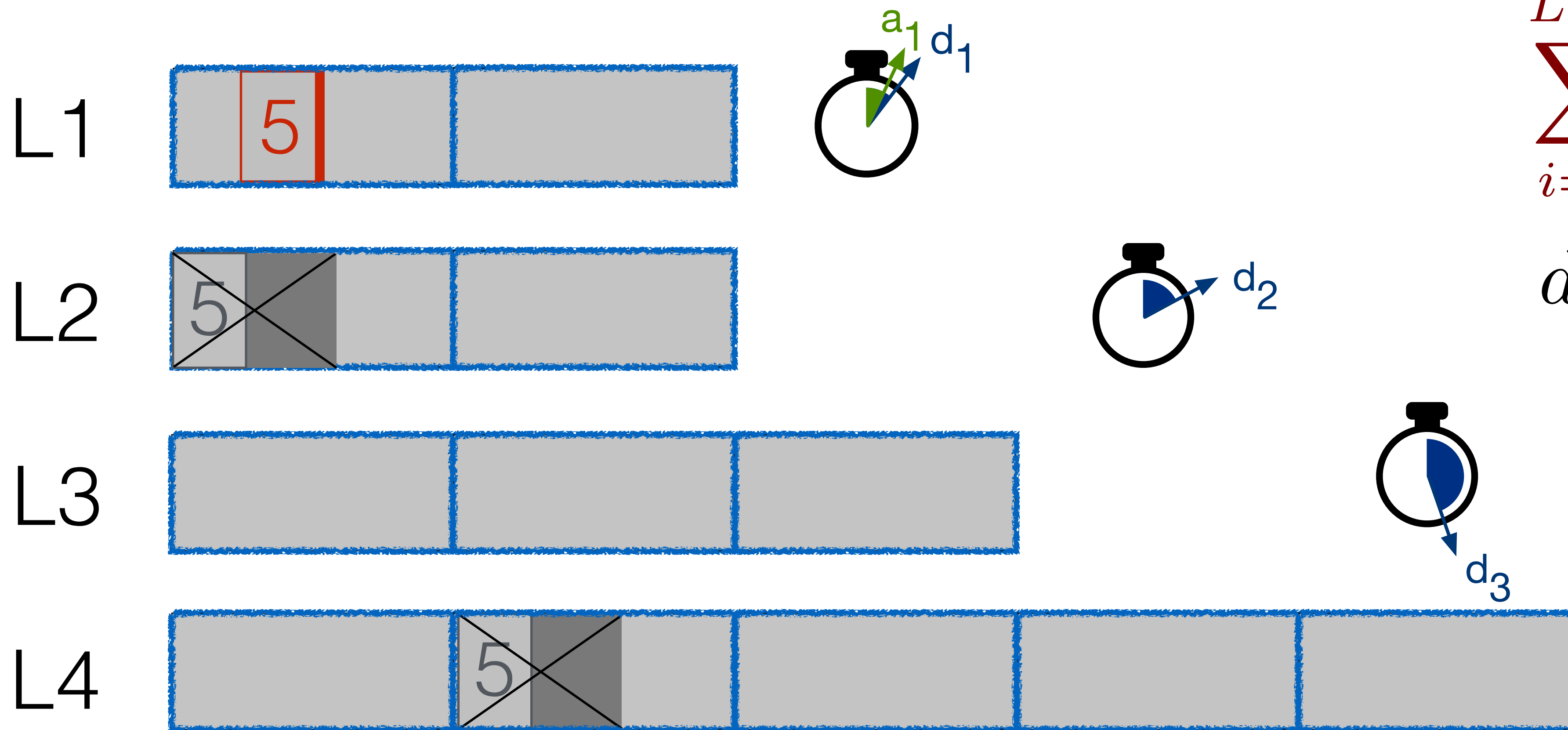


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

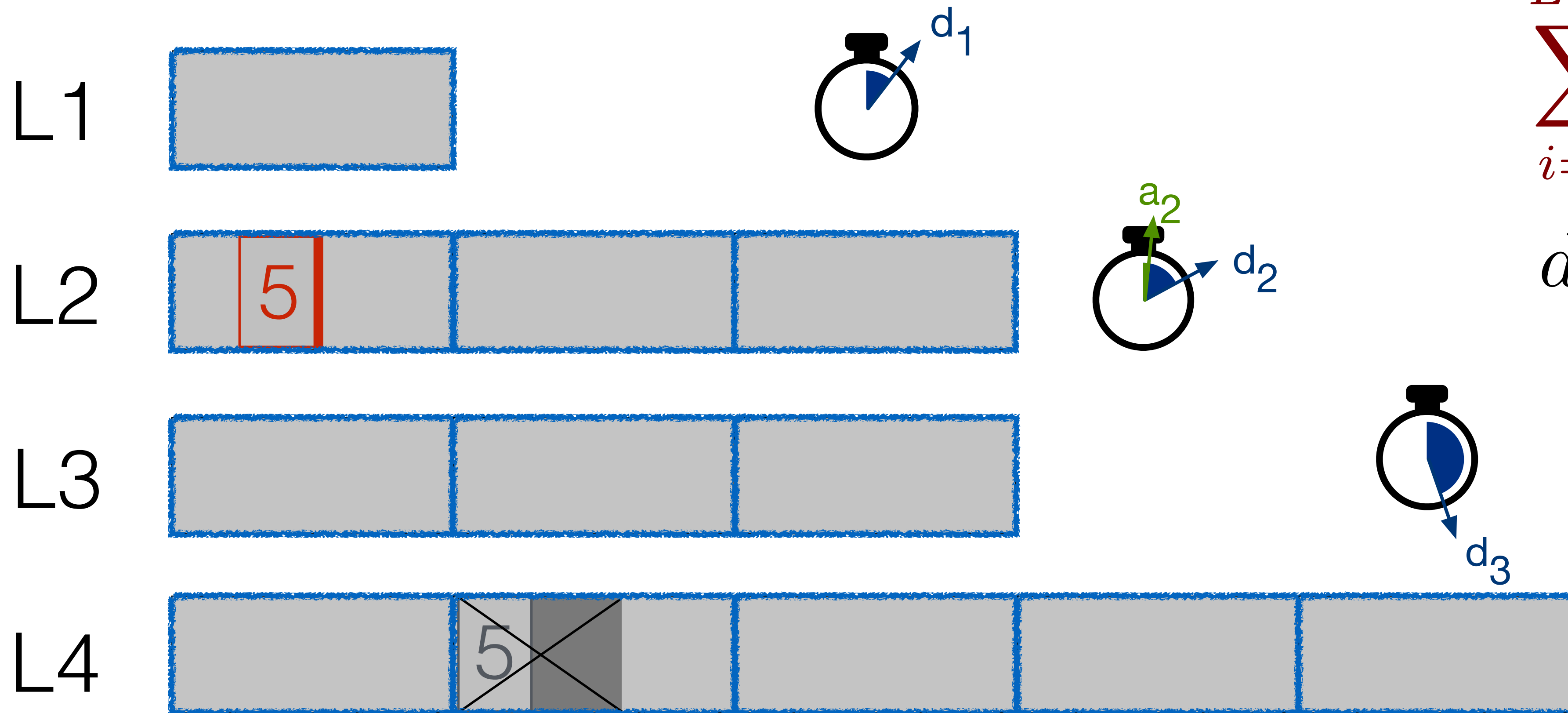


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

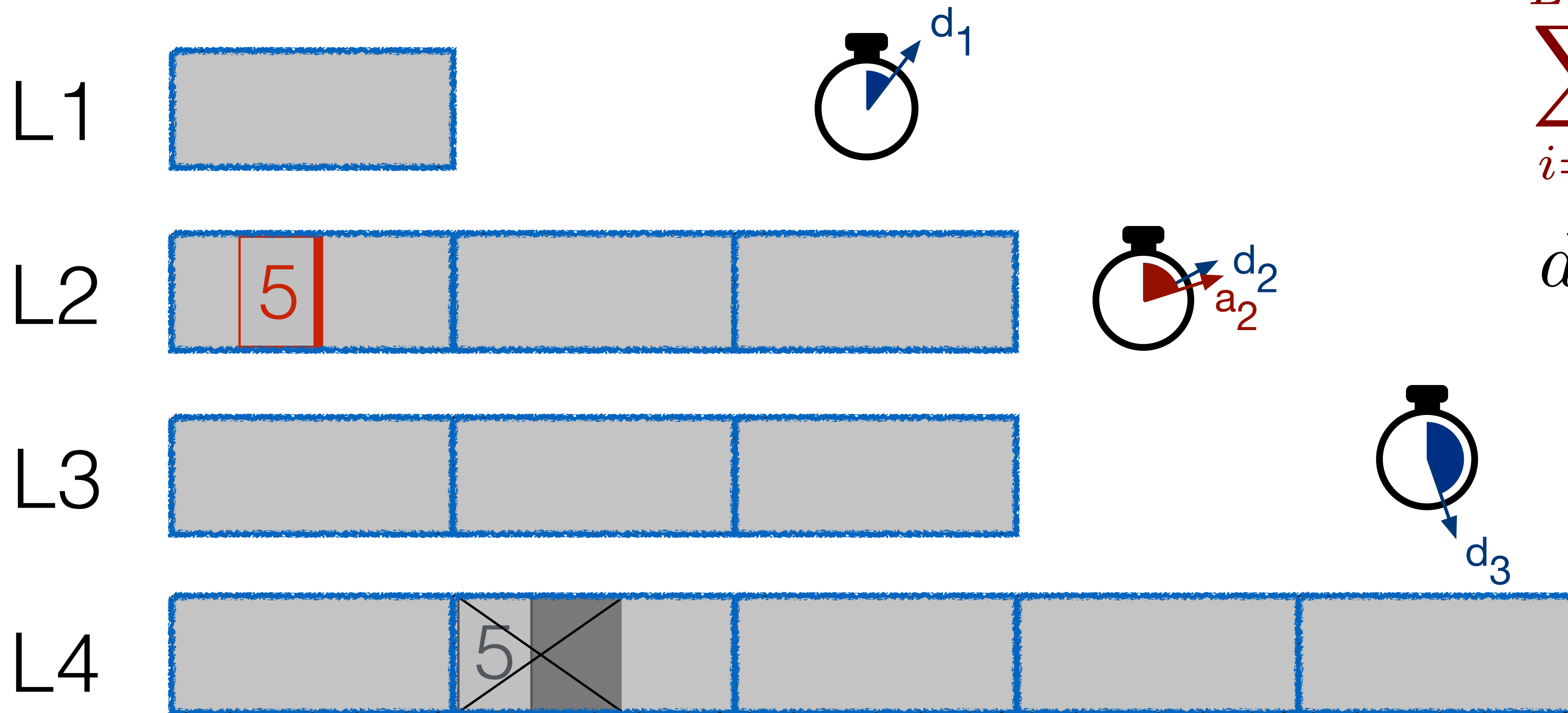


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

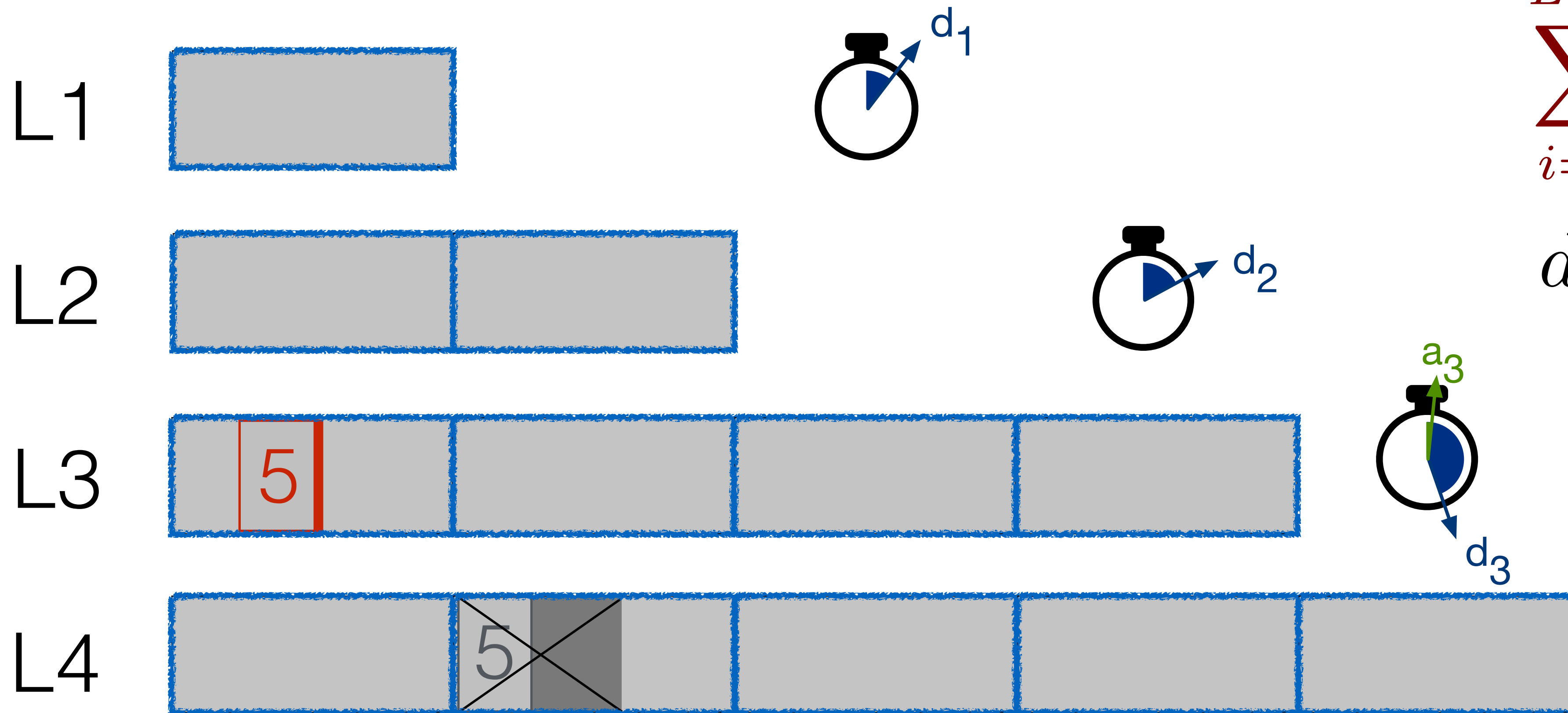


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

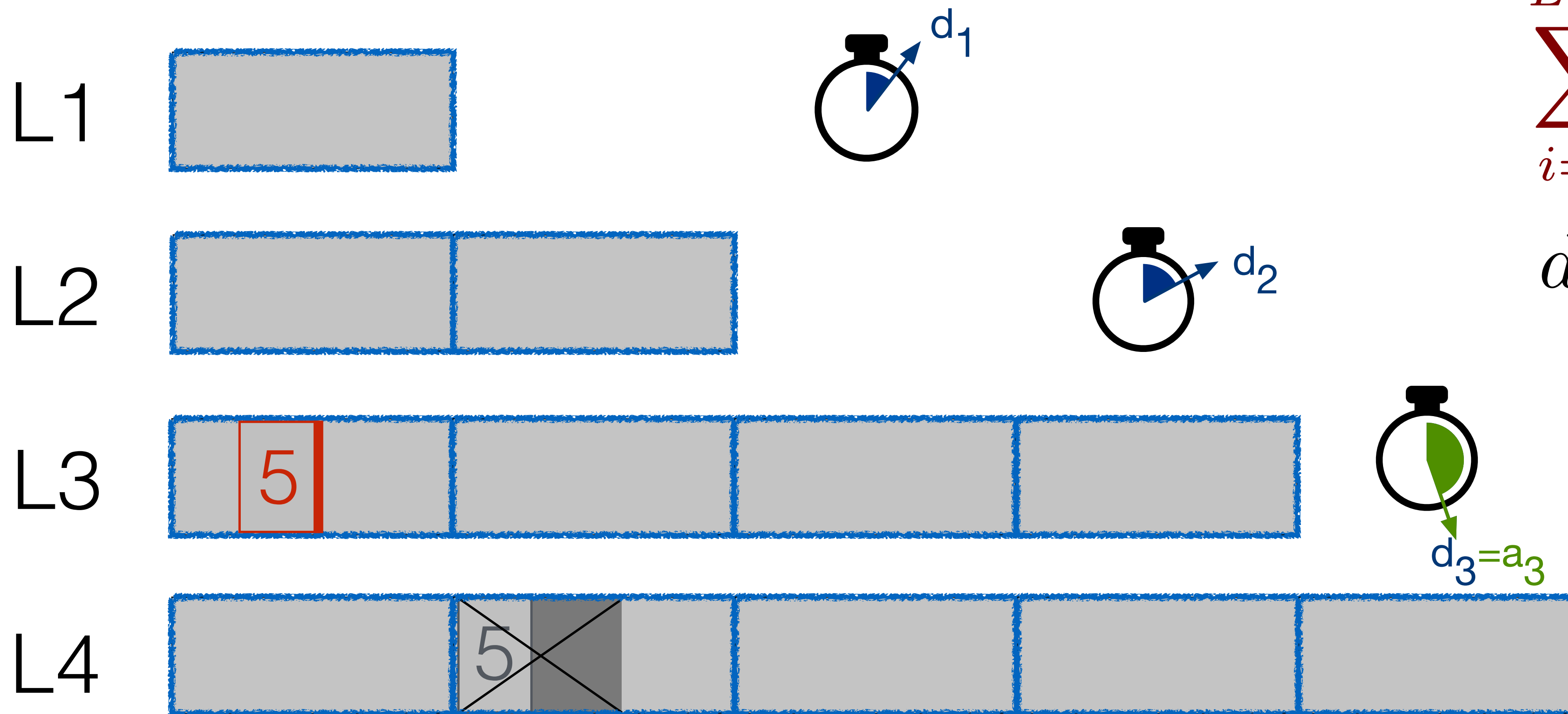


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

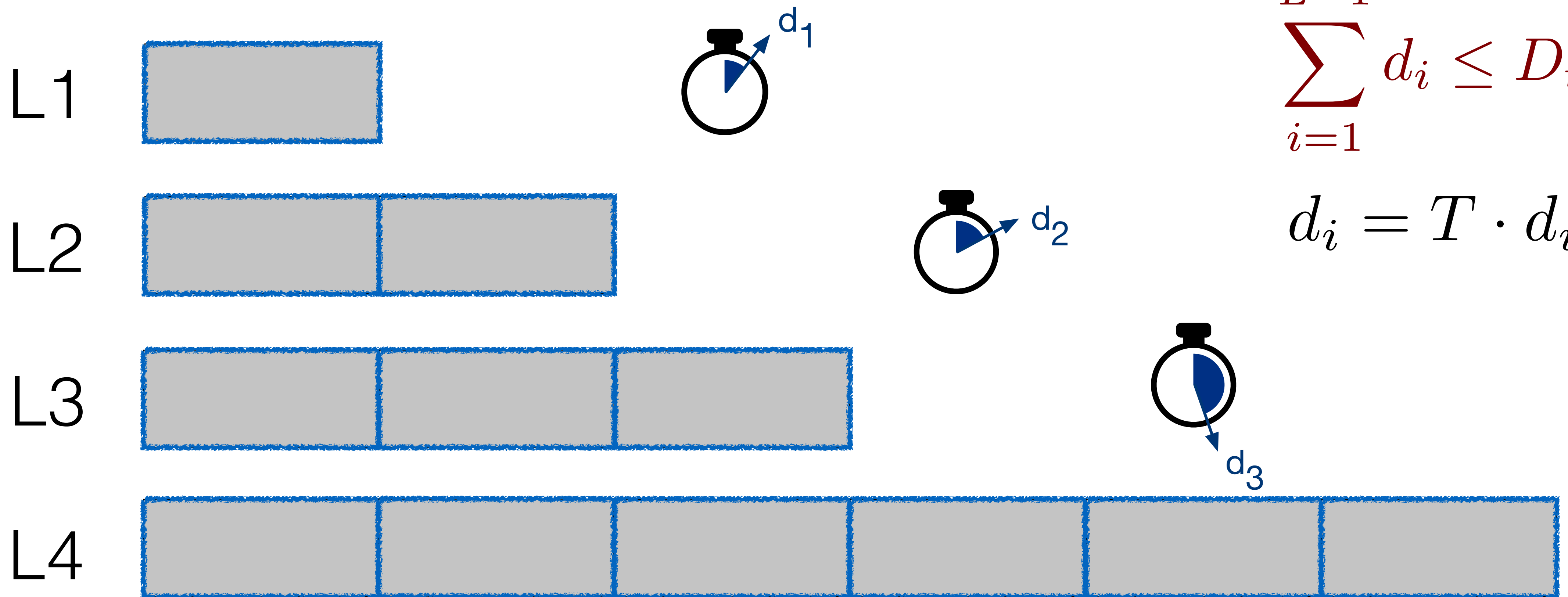


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

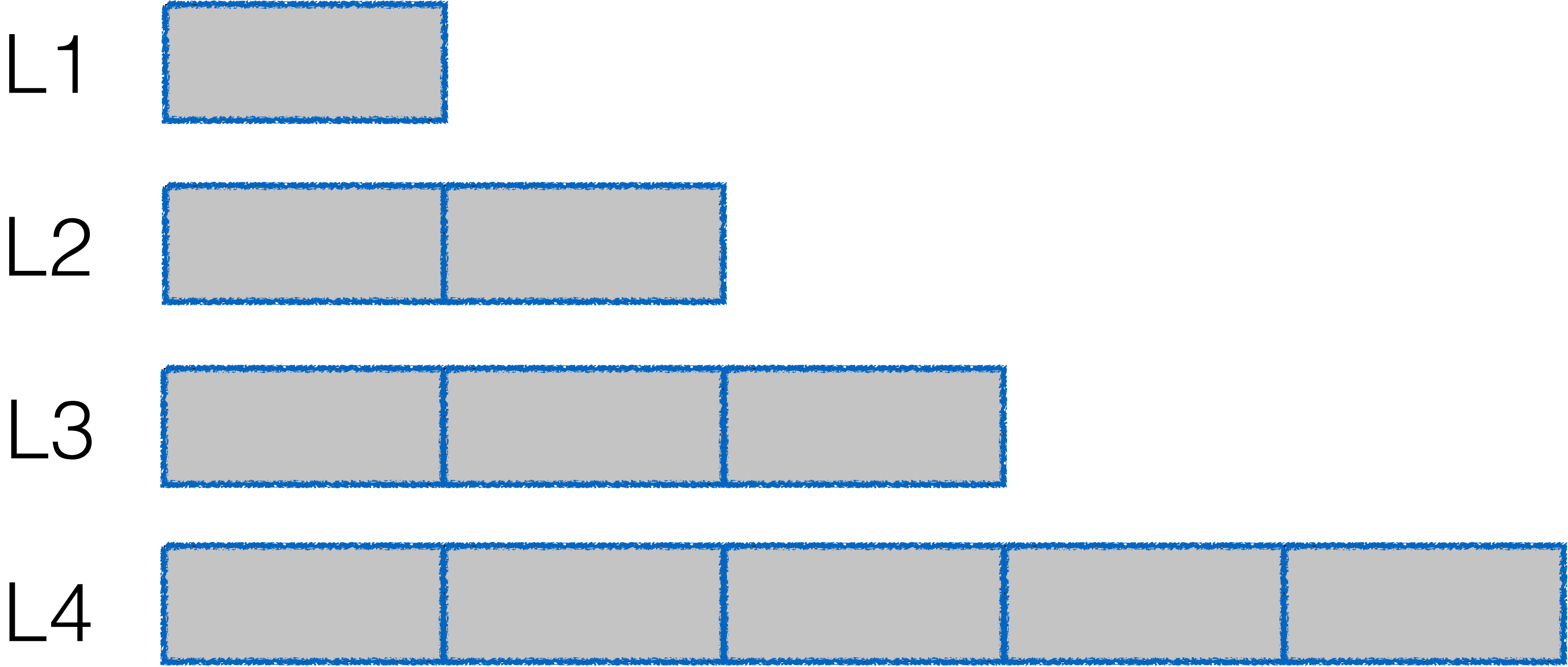


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

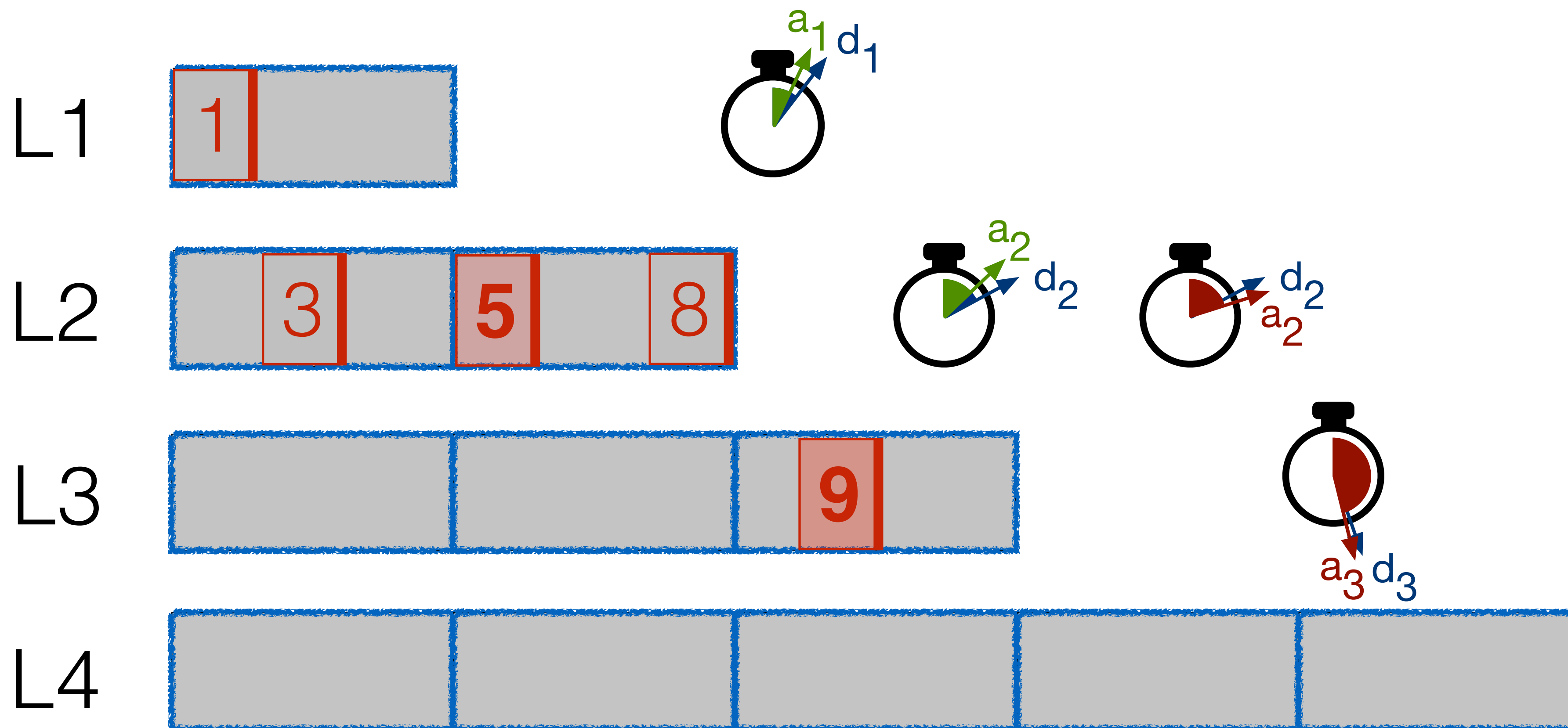
FAst DElete

breaking ties in practical workloads



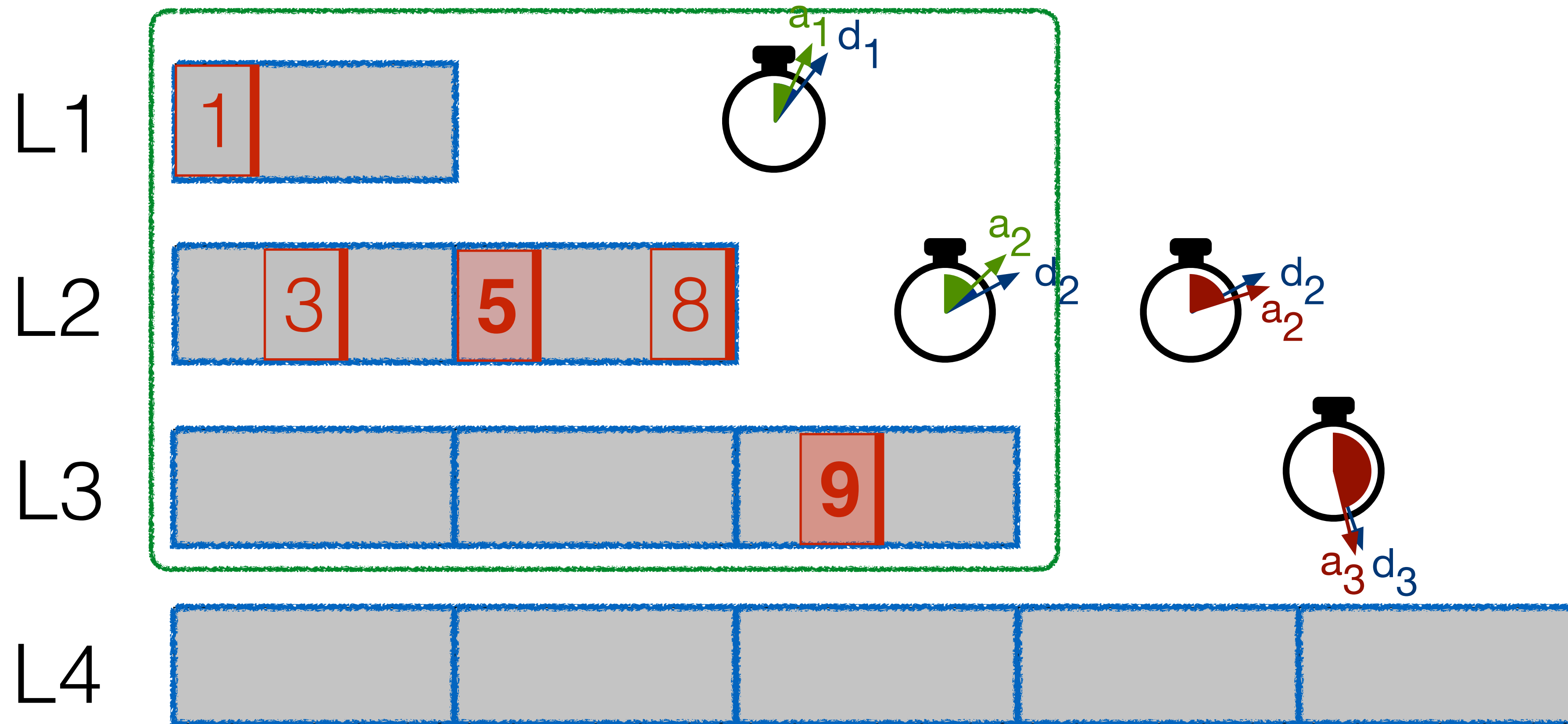
FASt DElete

breaking ties in practical workloads



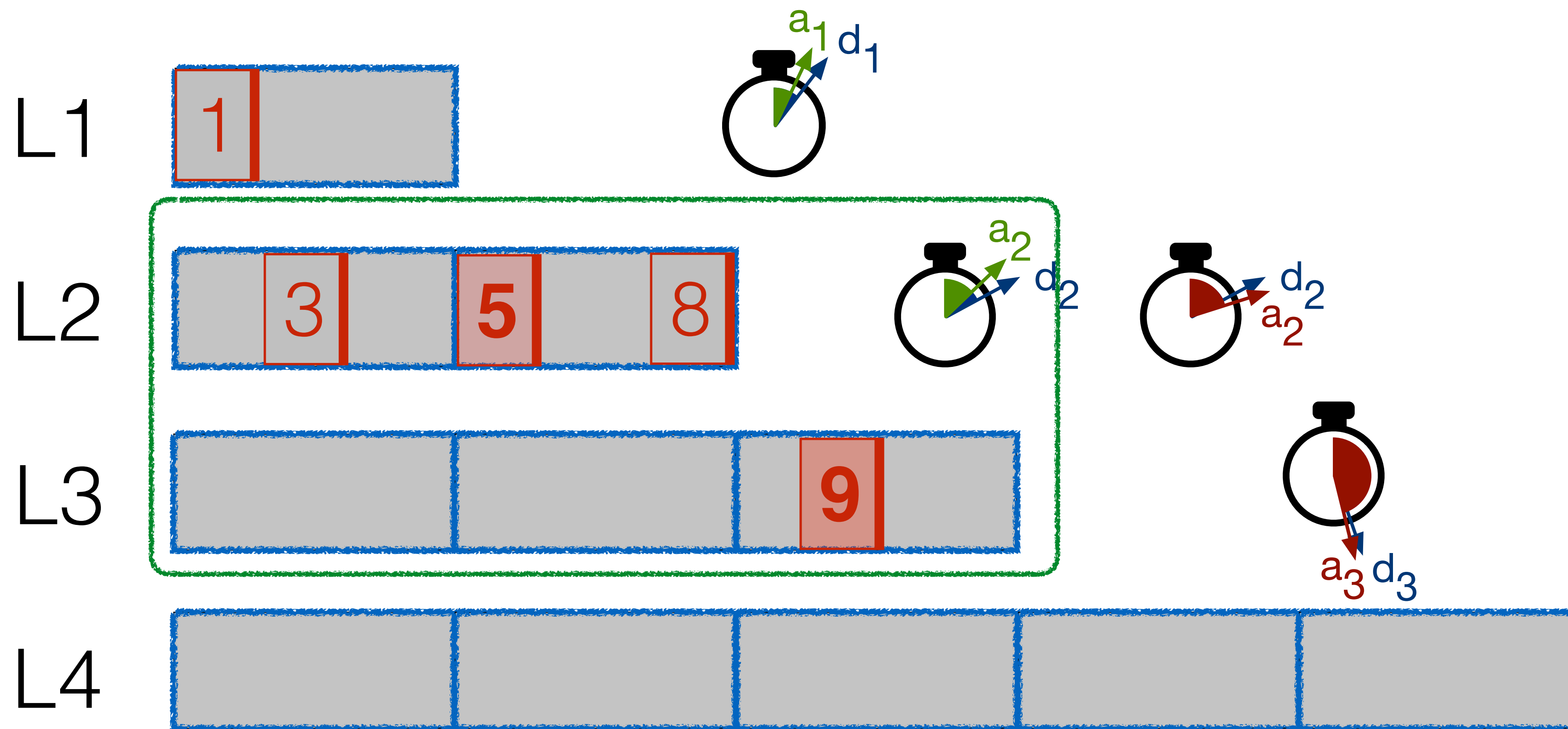
FAst DElete

breaking ties in practical workloads



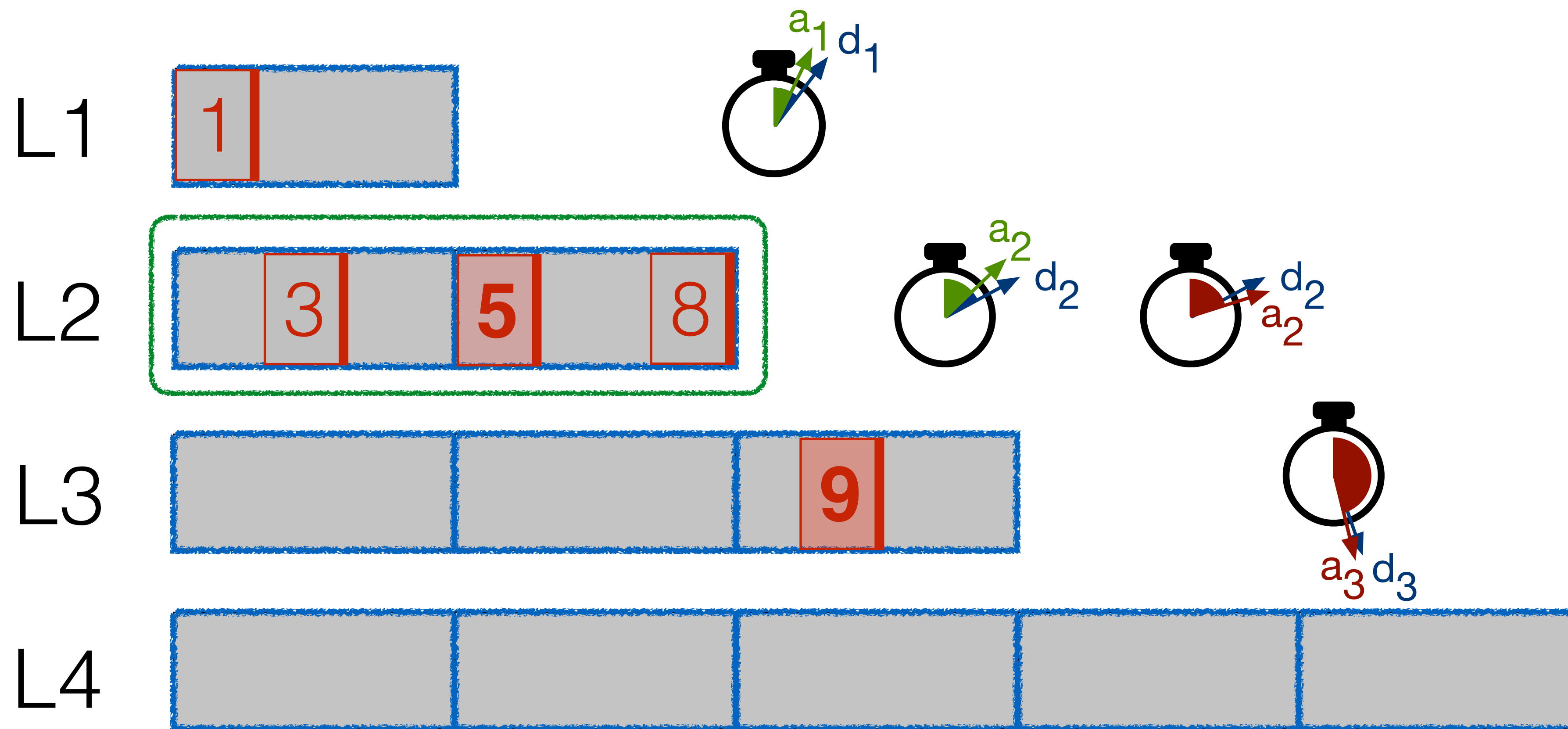
FAst DElete

breaking ties in practical workloads



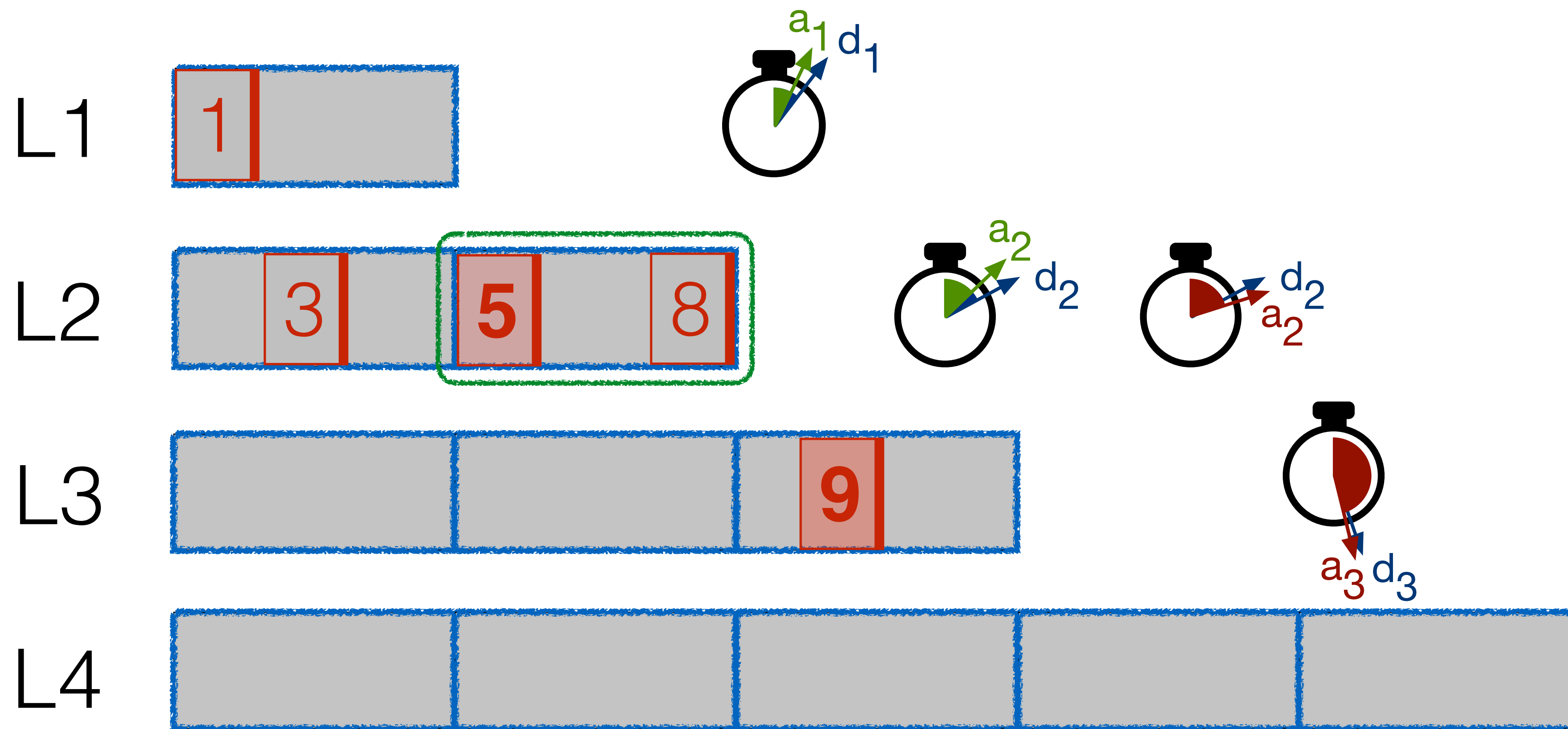
FAst DElete

breaking ties in practical workloads



FAst DElete

breaking ties in practical workloads



FAst DElete

higher write amplification

4 - 25%



improved read performance

1.2 - 1.4x



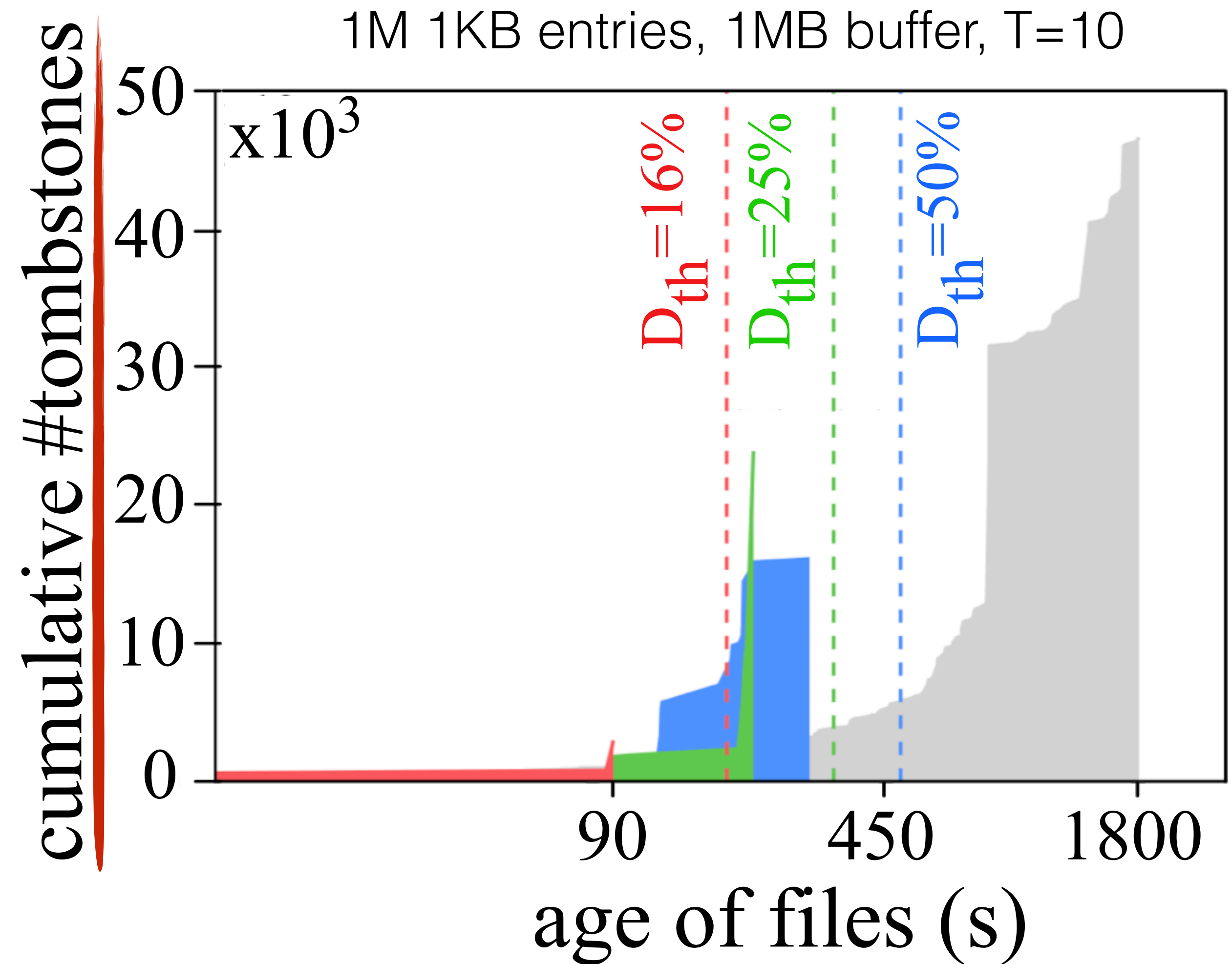
reduced space amplification

2.1 - 9.8x



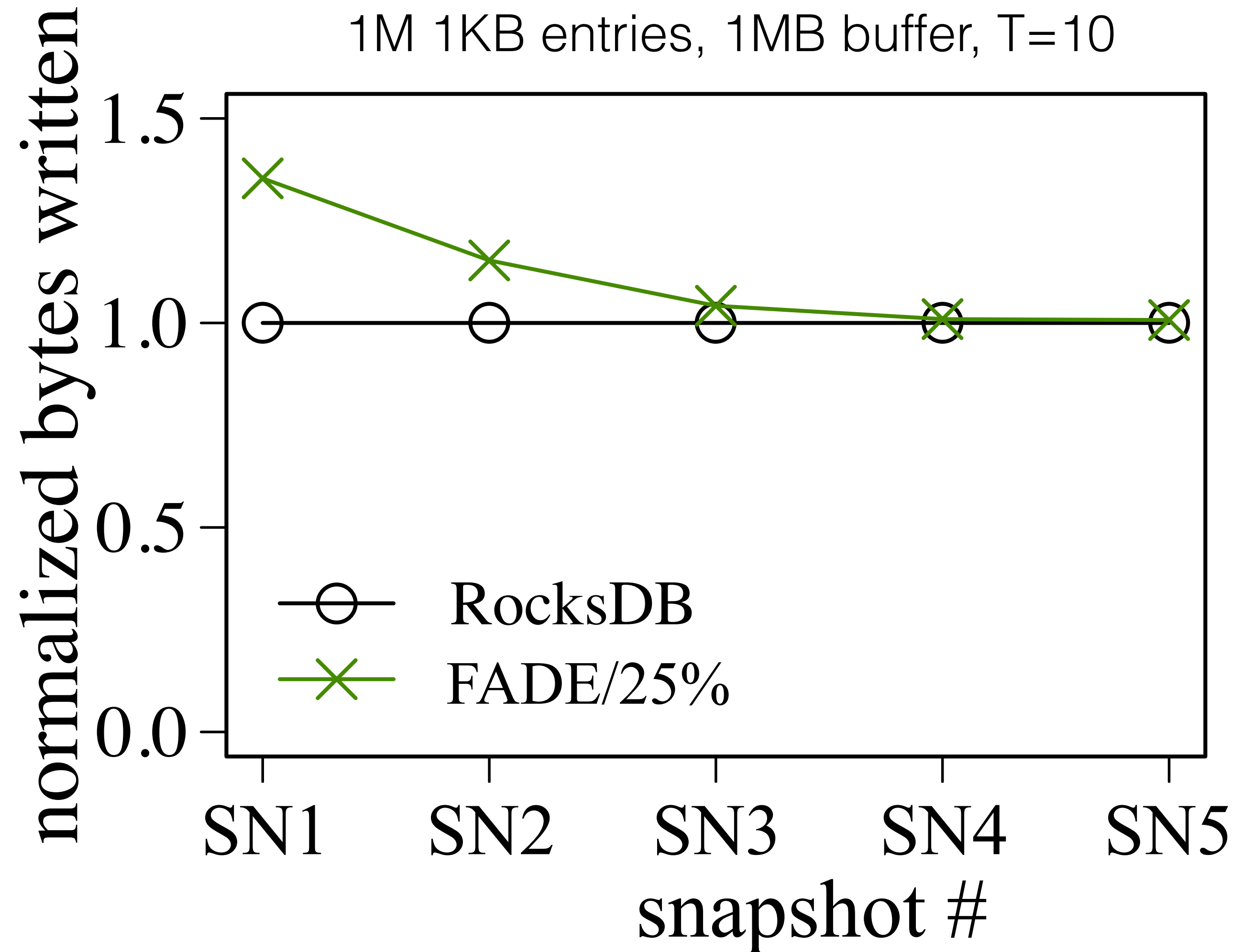
timely delete persistence

within D_{th}



FAst DElete

- higher write amplification ●
0.7%
- improved read performance ✓
1.2 - 1.4x
- reduced space amplification ✓
2.1 - 9.8x
- timely delete persistence ✓
within D_{th}



the solution

FASt DElete

higher write amplification

FADE

improved read performance

reduced space amplification

timely delete persistence

latency spikes

Kiwi

superfluous I/Os

the solution

FASt DElete

higher write amplification

FADE

parallel read performance

reduced space amplification

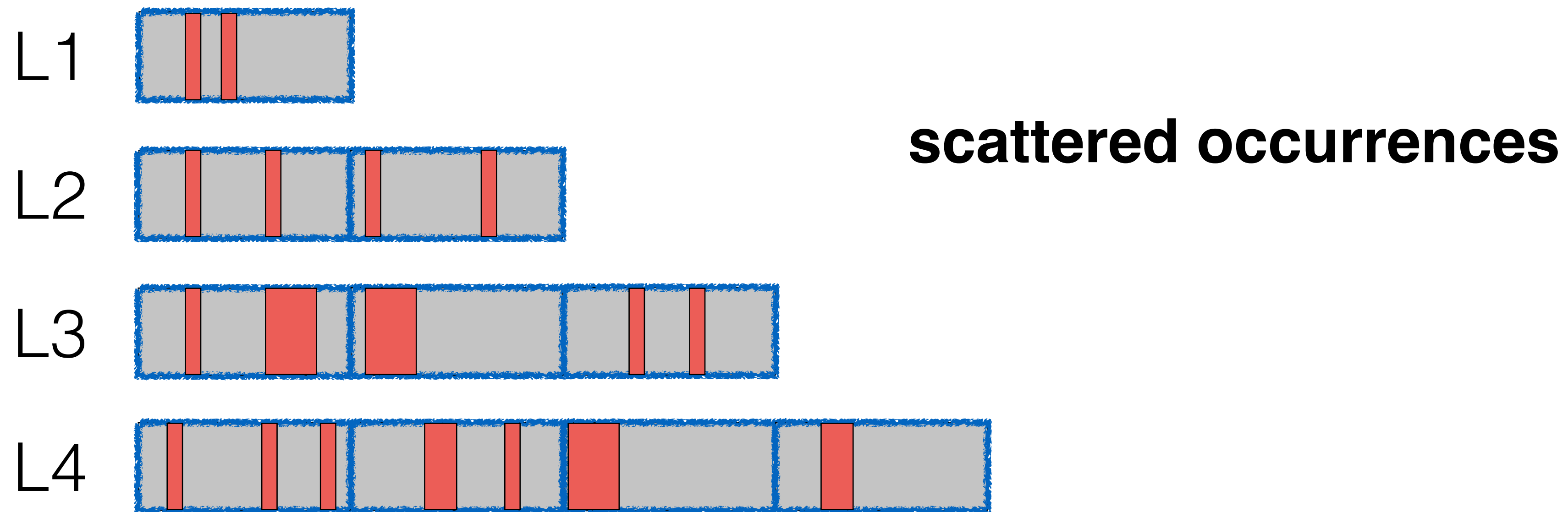
timely delete persistence

latency spikes

Kiwi

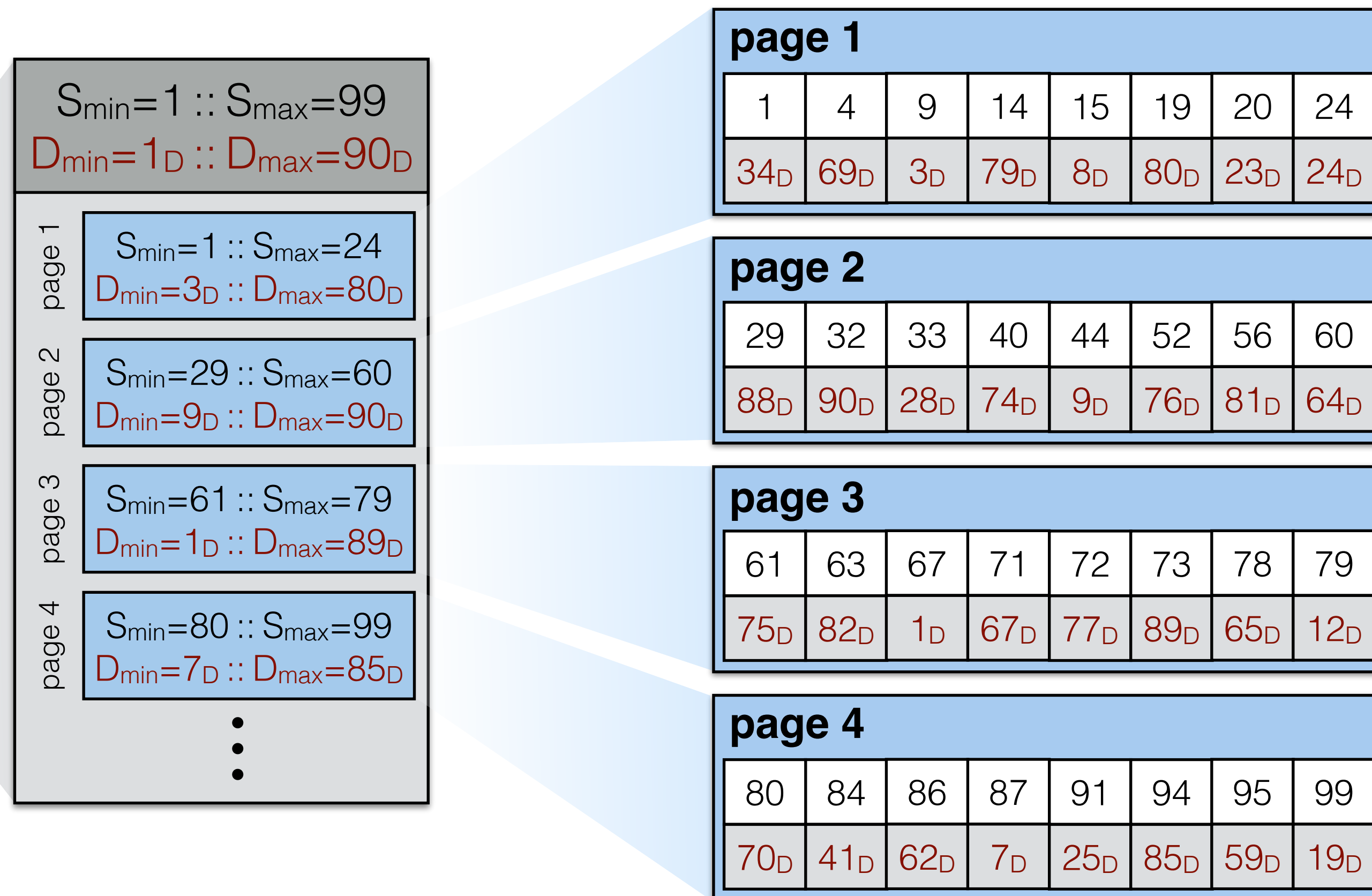
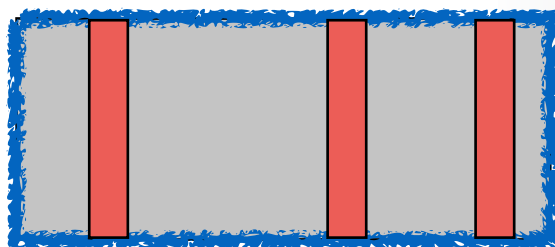
Key Weaving storage layout

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



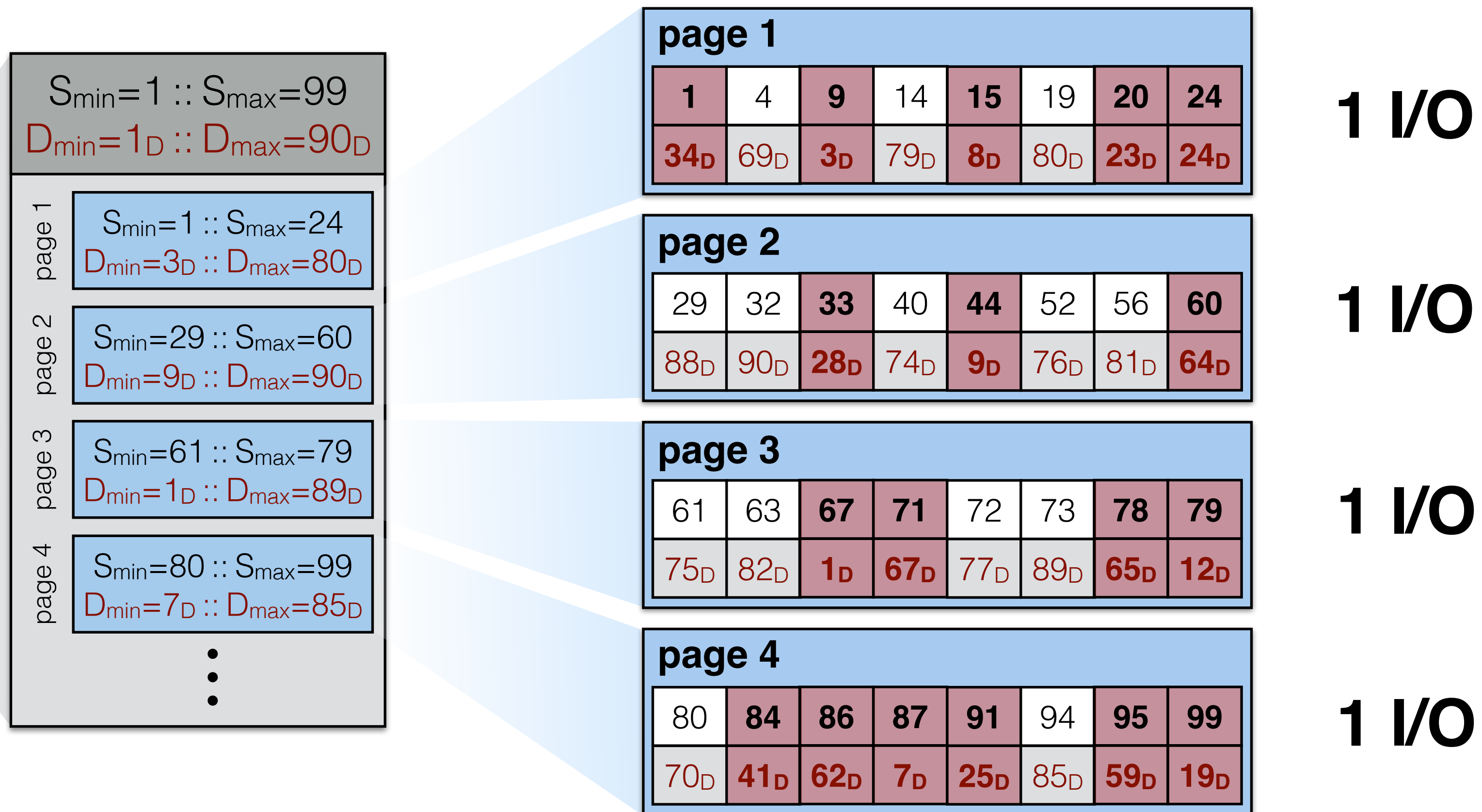
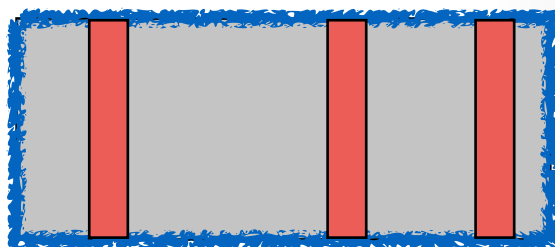
Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



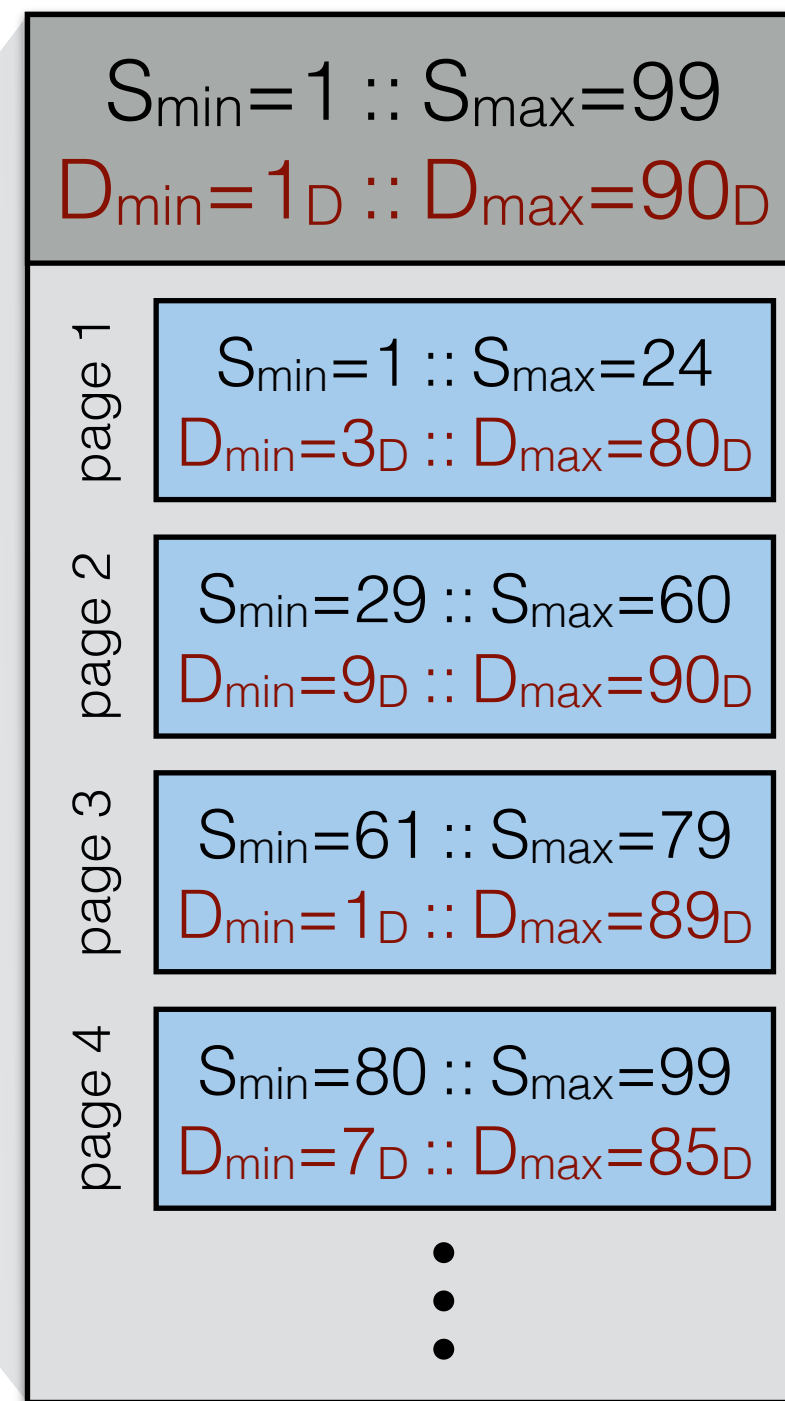
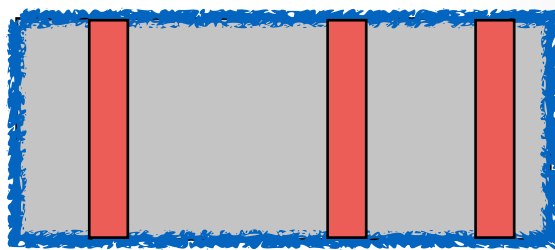
Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



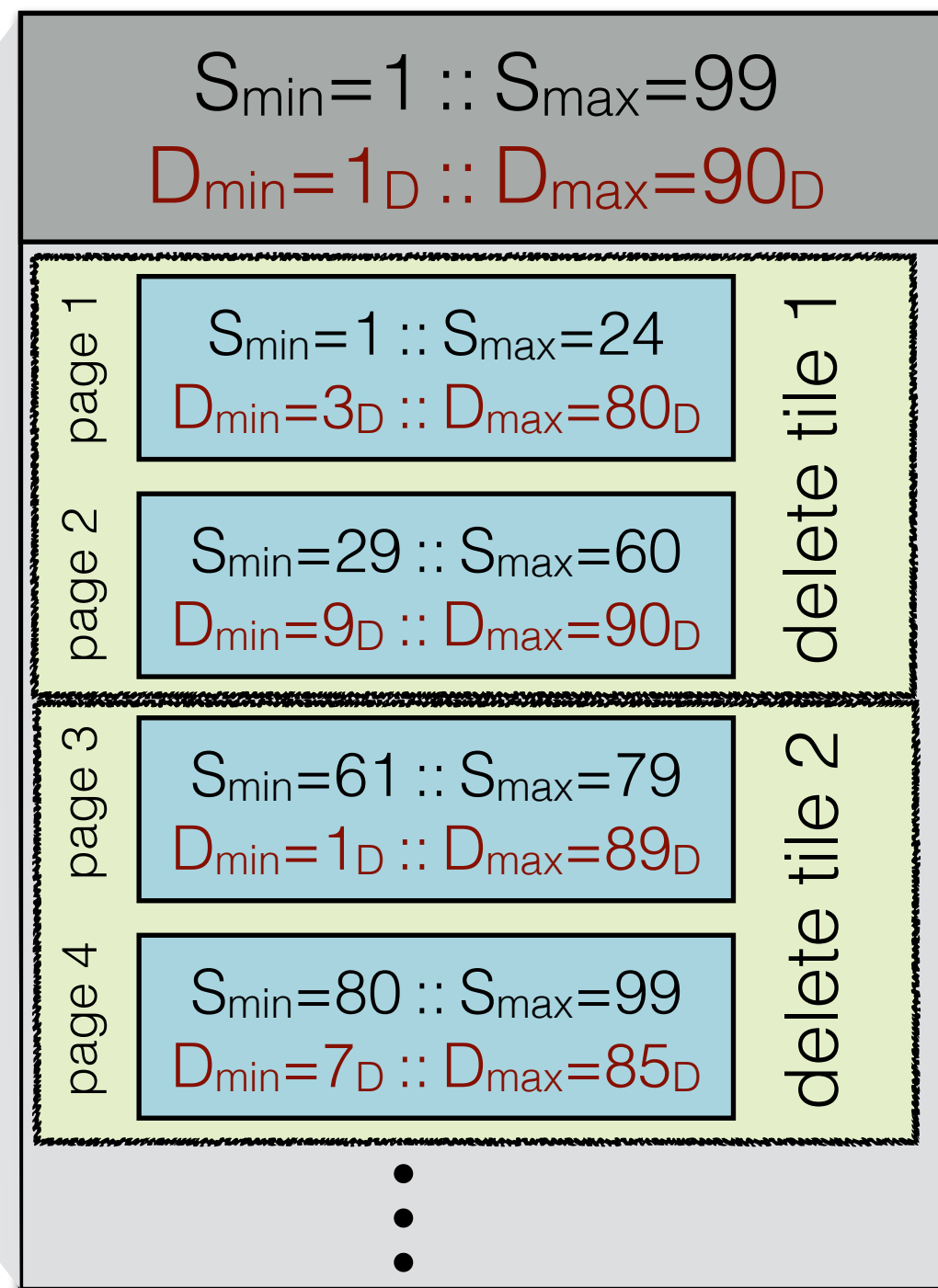
Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



Key Weaving storage layout

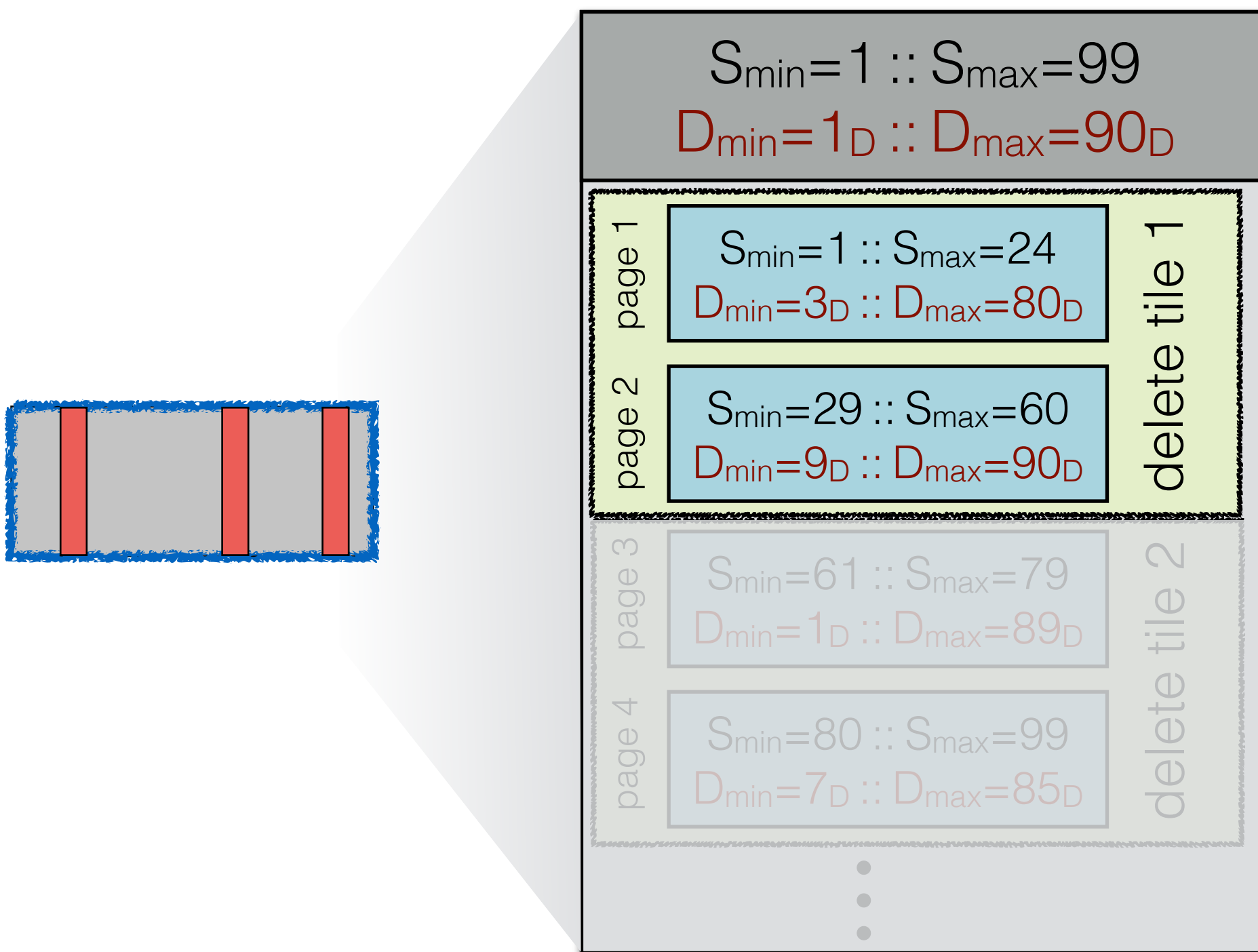
delete all entries with timestamp $\leq 65_D$



partitioned on S

Key Weaving storage layout

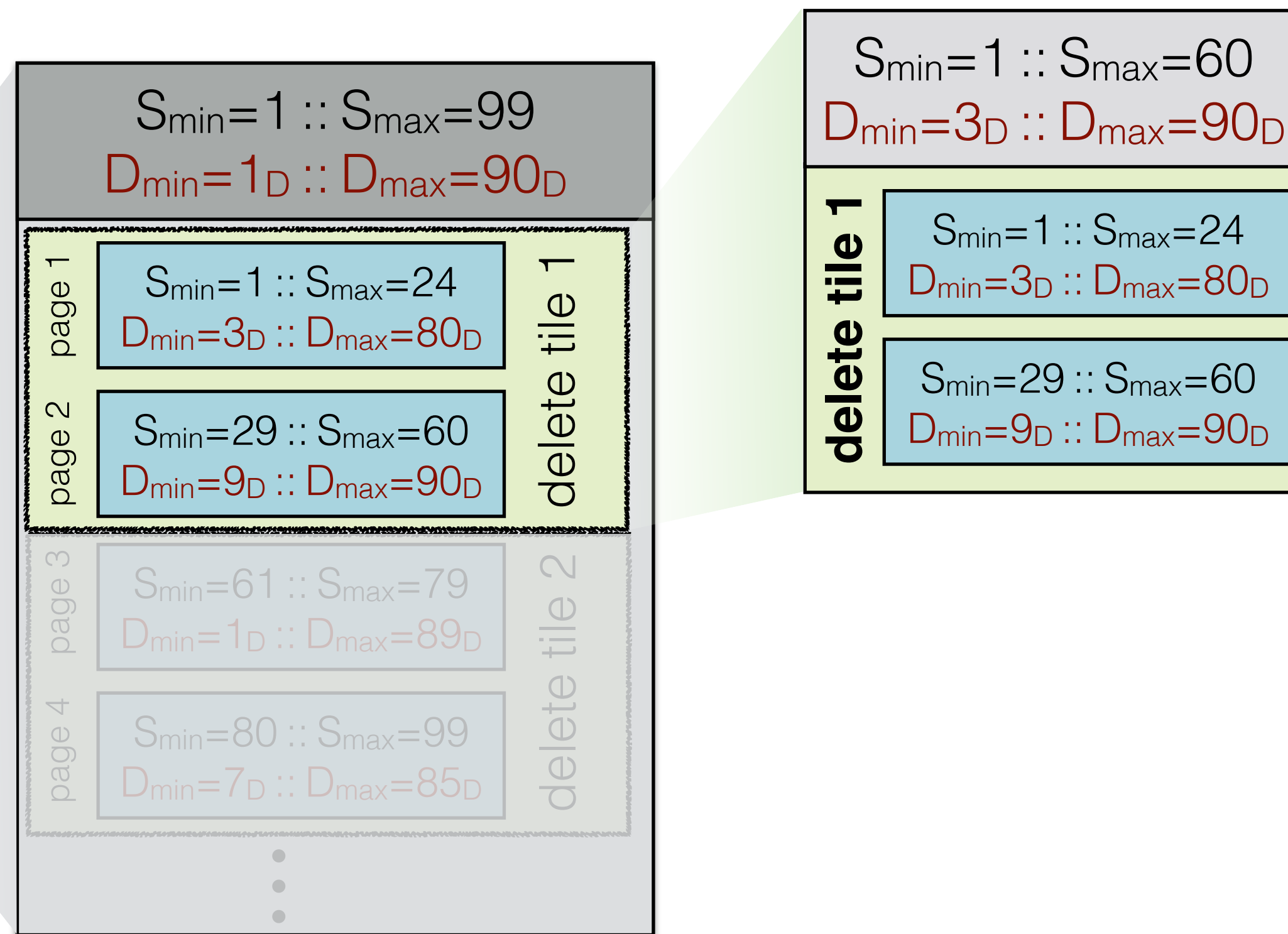
delete all entries with timestamp $\leq 65_D$



partitioned on S

Key Weaving storage layout

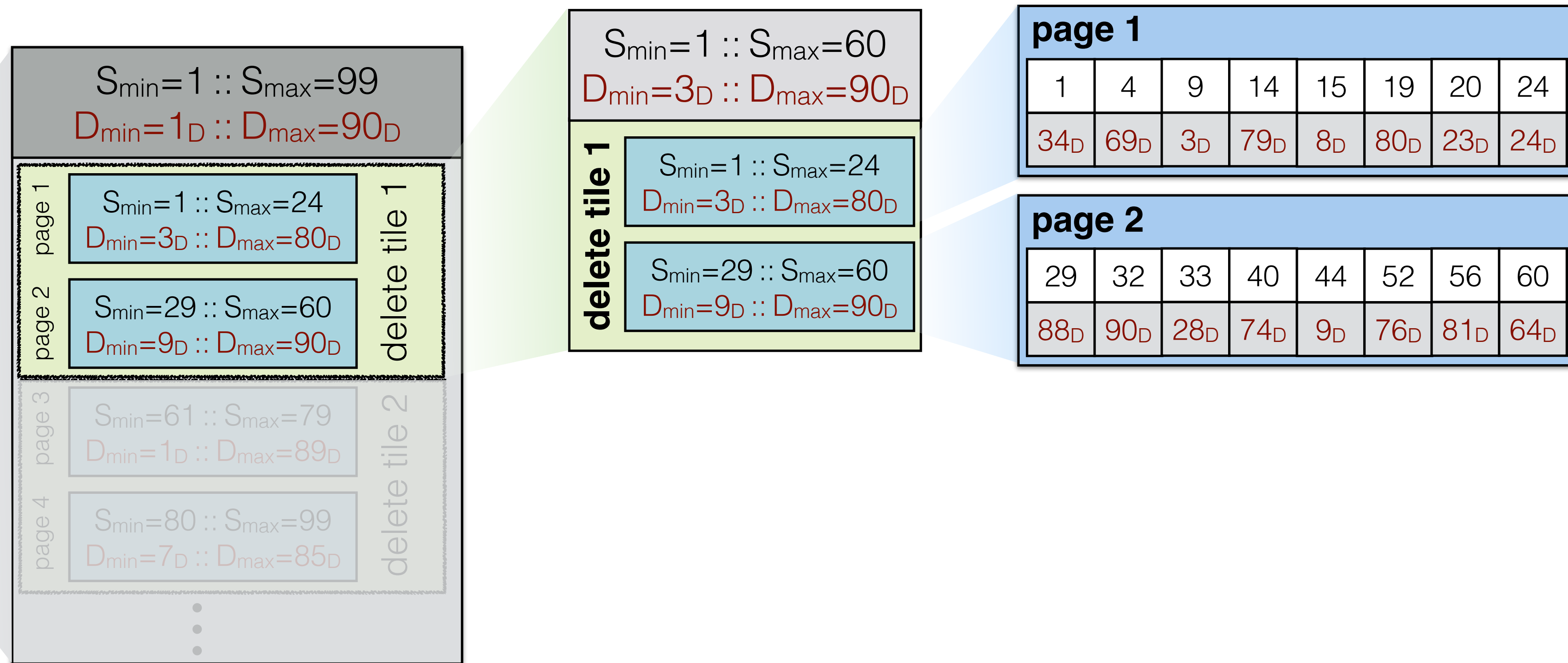
delete all entries with timestamp $\leq 65_D$



partitioned on S

Key Weaving storage layout

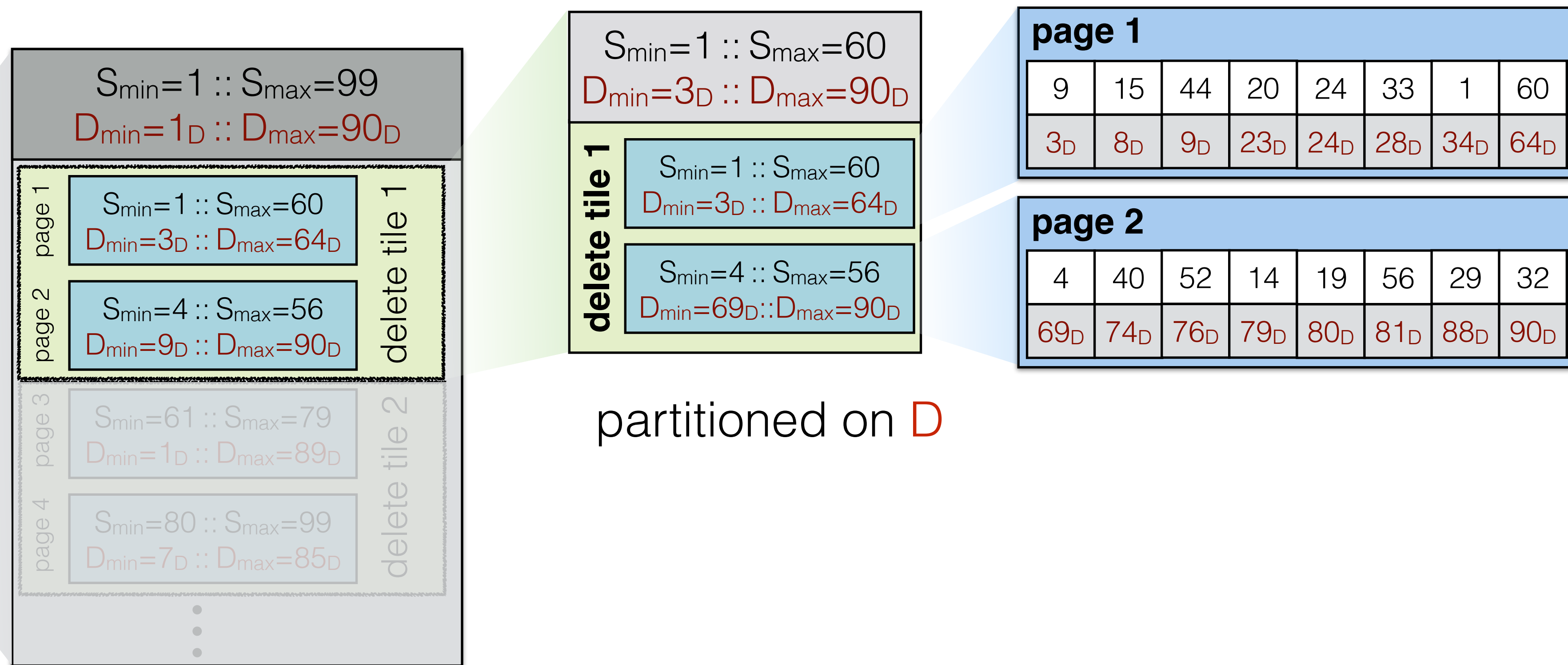
delete all entries with timestamp $\leq 65_D$



partitioned on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



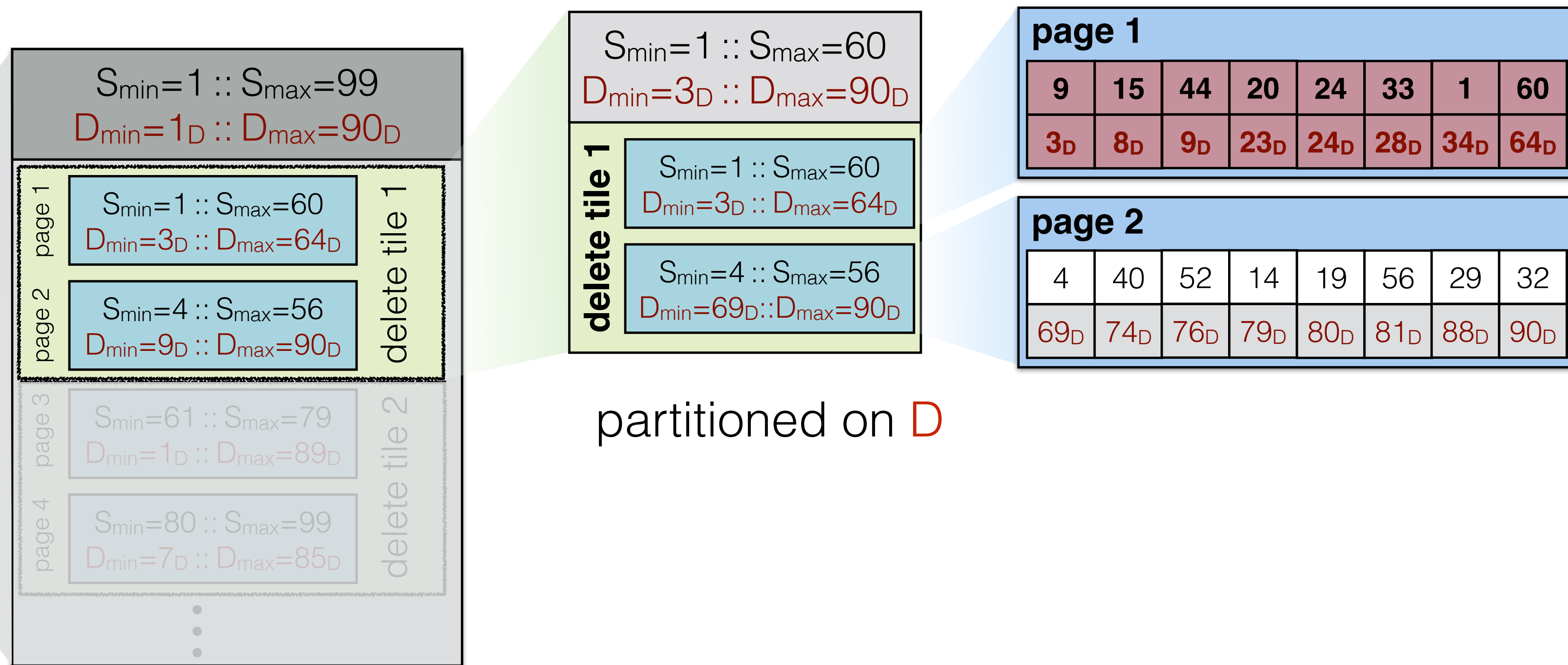
partitioned on S

partitioned on D

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

drop
page

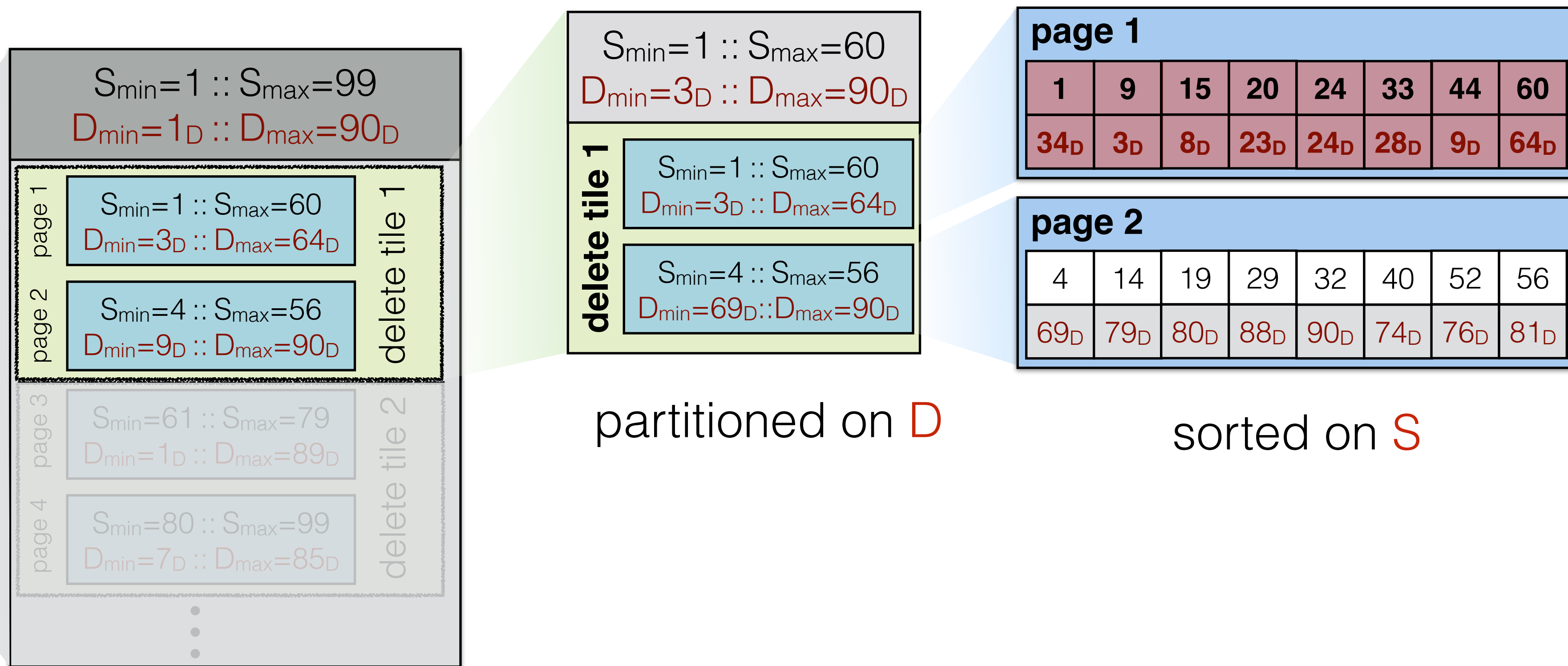


partitioned on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

drop
page



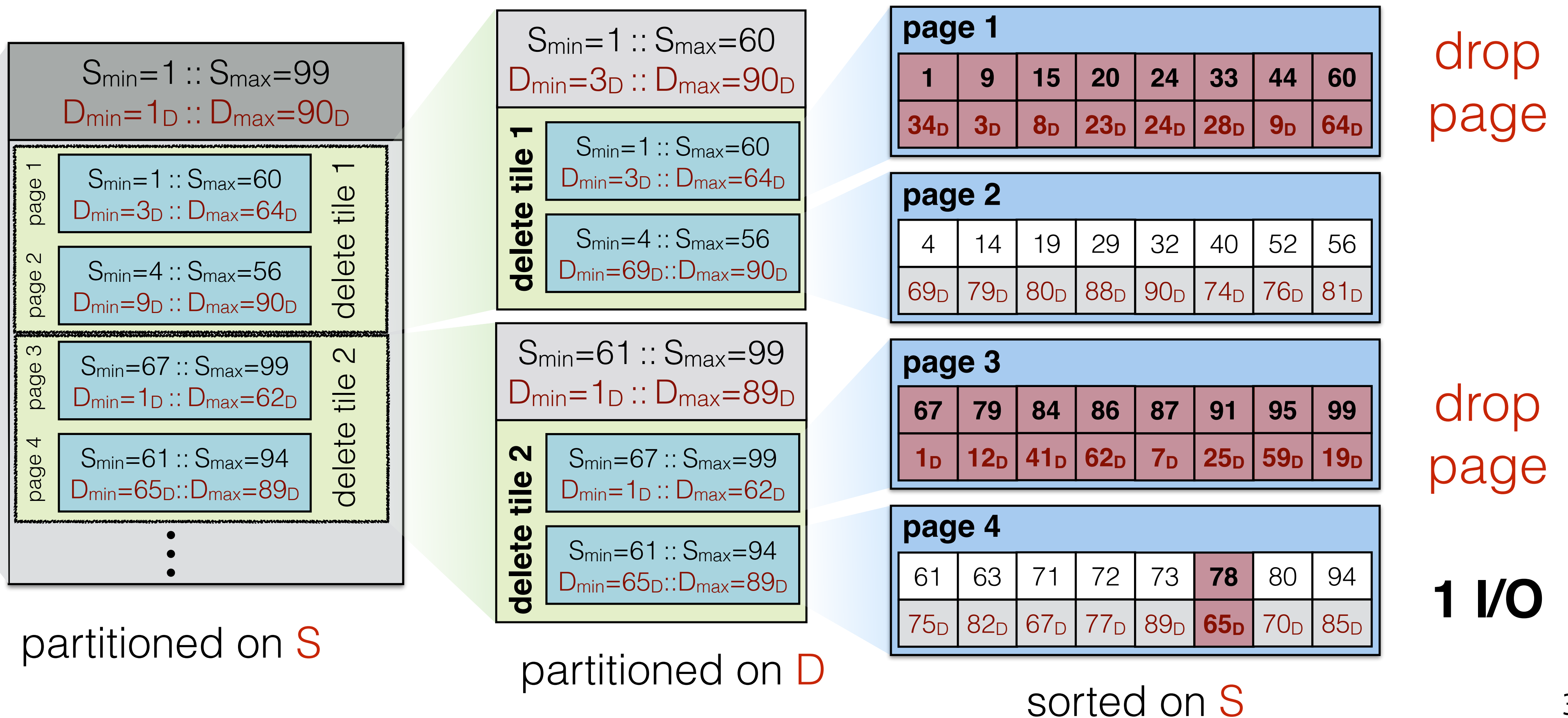
partitioned on S

partitioned on D

sorted on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

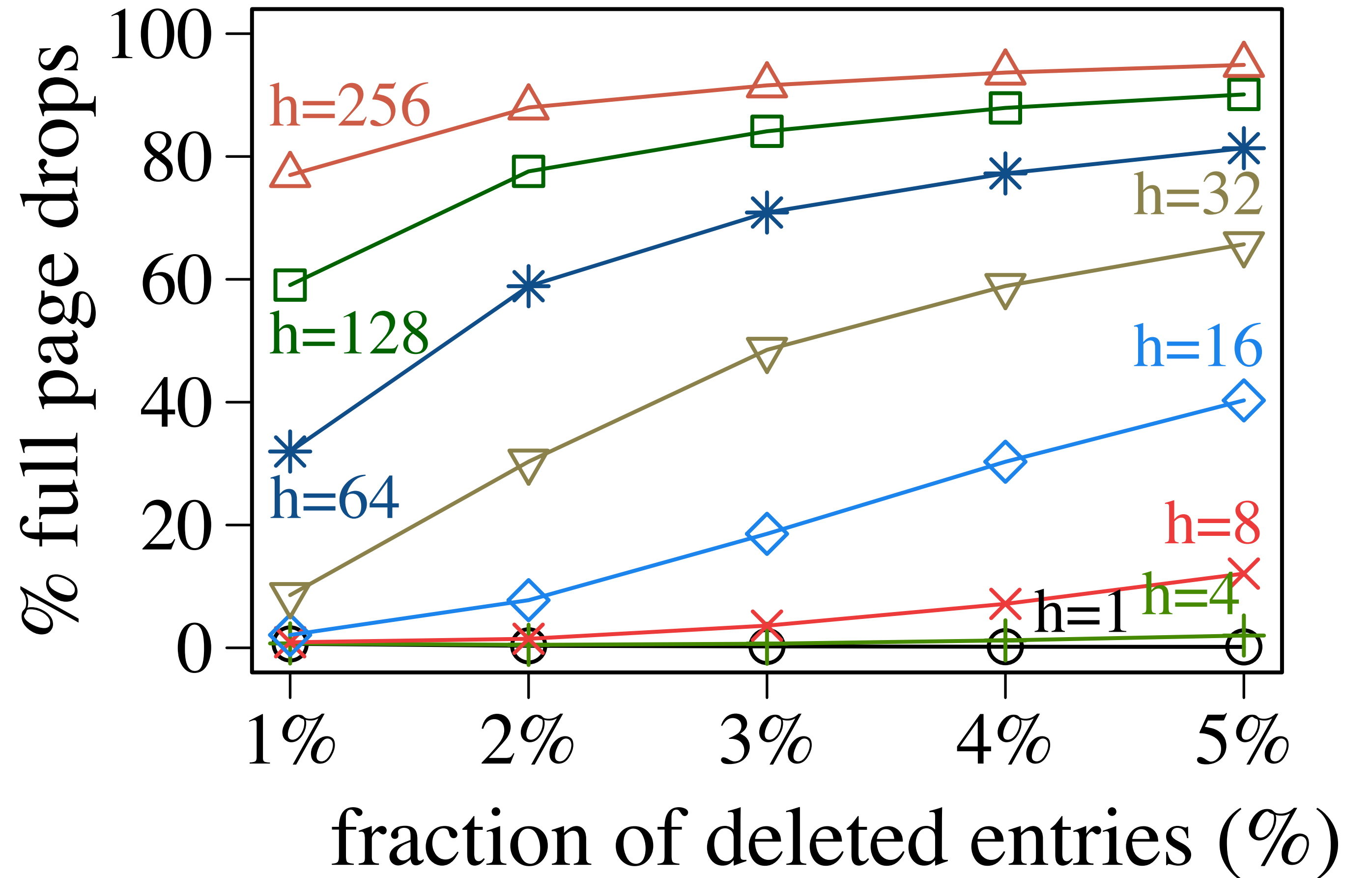


Key Weaving storage layout




1M 1KB entries, buffer = file = 256 pages

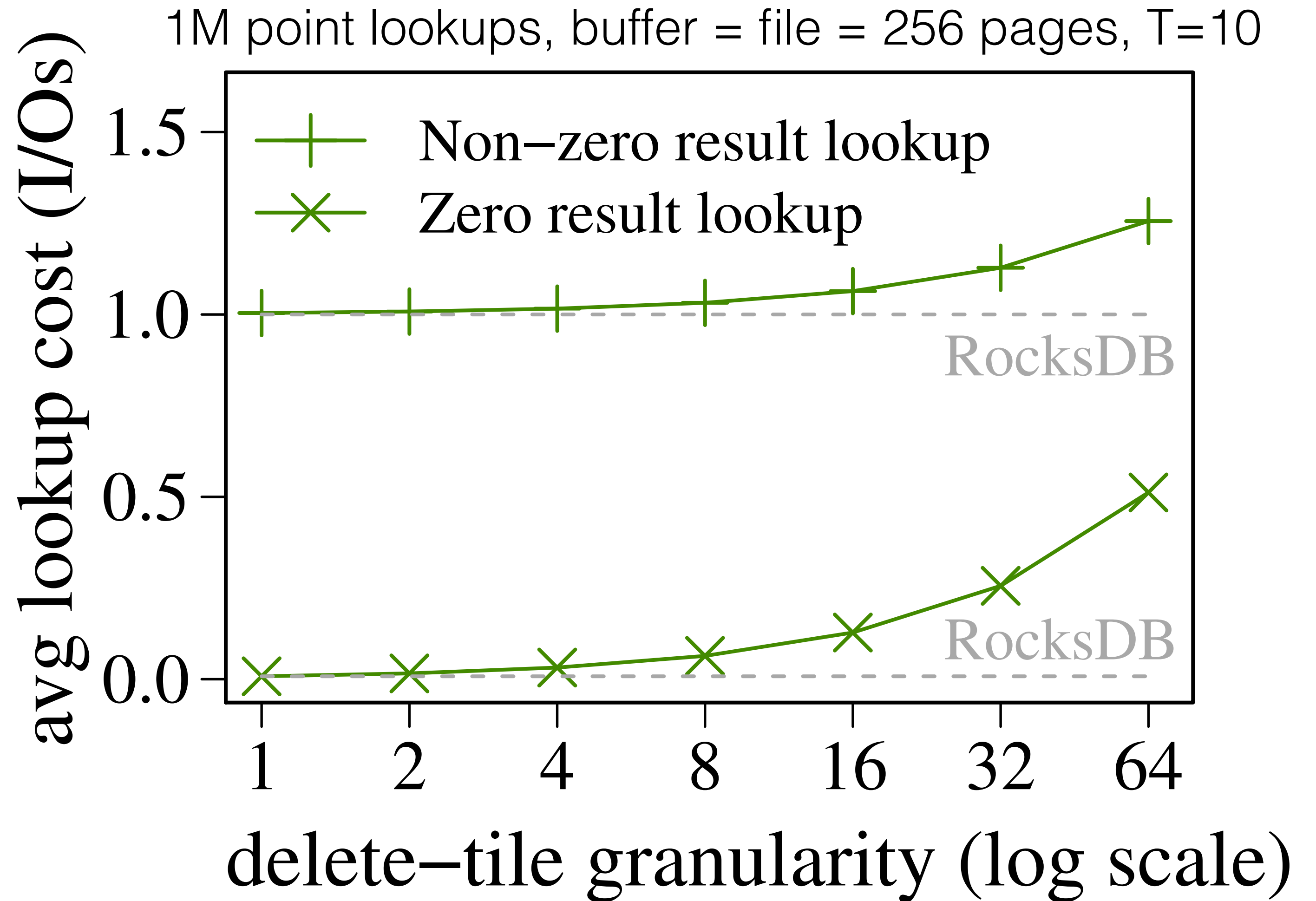
reduced latency spikes ✓

full page drops reduces superfluous I/Os ✓



Key Weaving storage layout

- higher lookup cost 
- reduced latency spikes 
- full page drops reduces superfluous I/Os 



the solution

FAst DElete

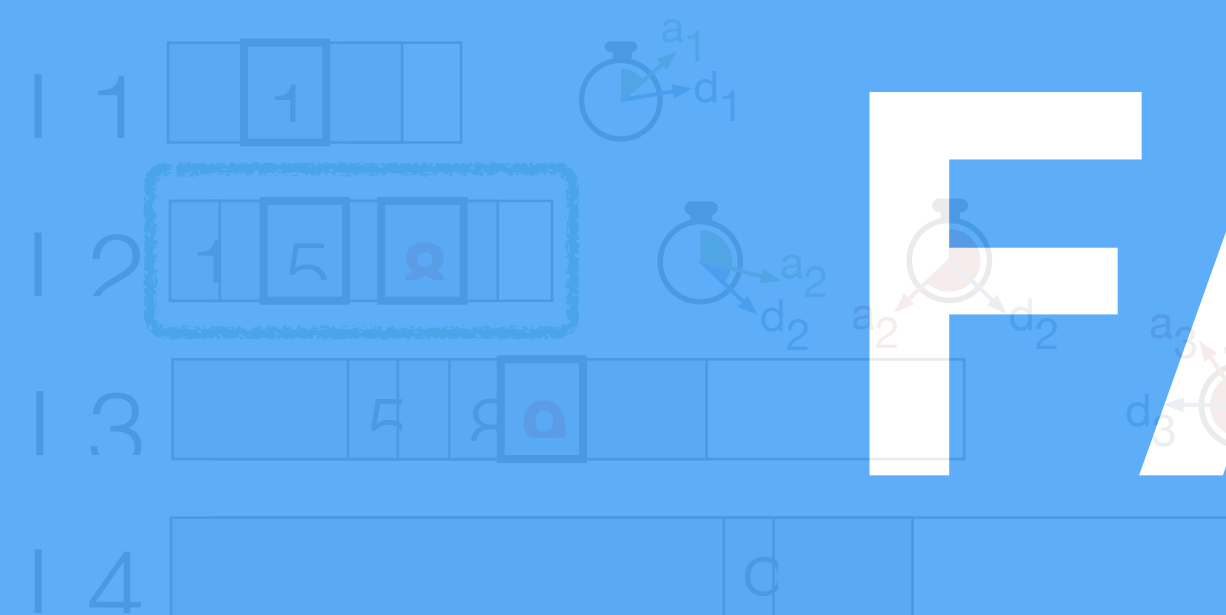
amortized write

reduced space

improved read

timely delete persistence

FADE



higher lookup cost

Kiwi

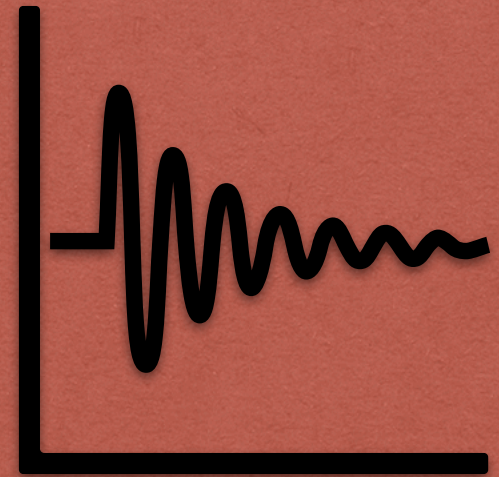
reduced latency spikes

full page drops reduces
superfluous I/Os

the solution

FADE

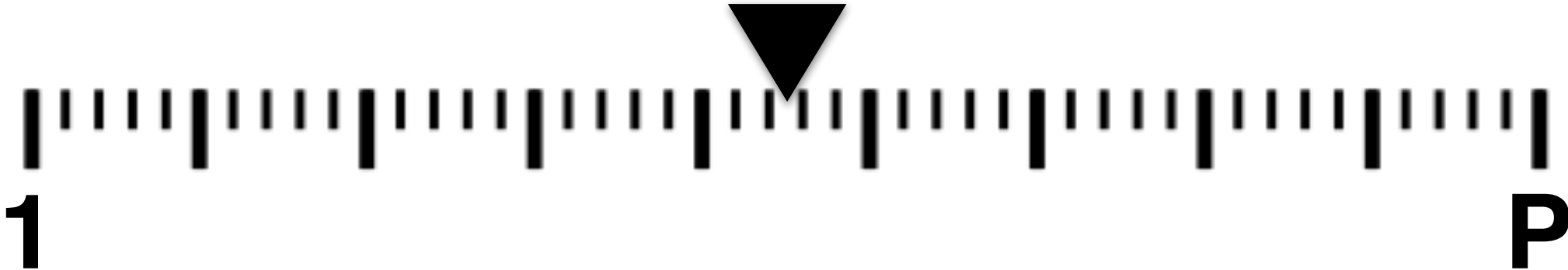
Kiwi



Letha

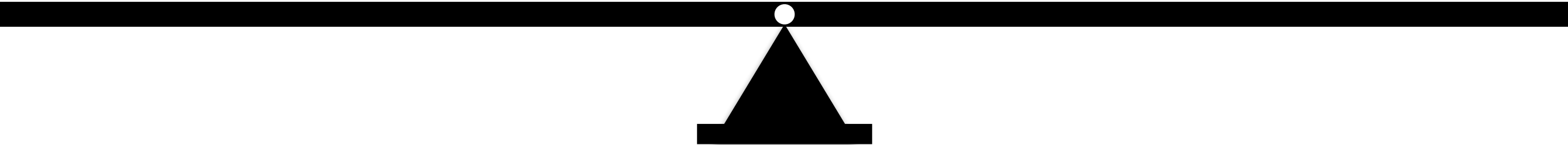


delete tile size



lookup
cost

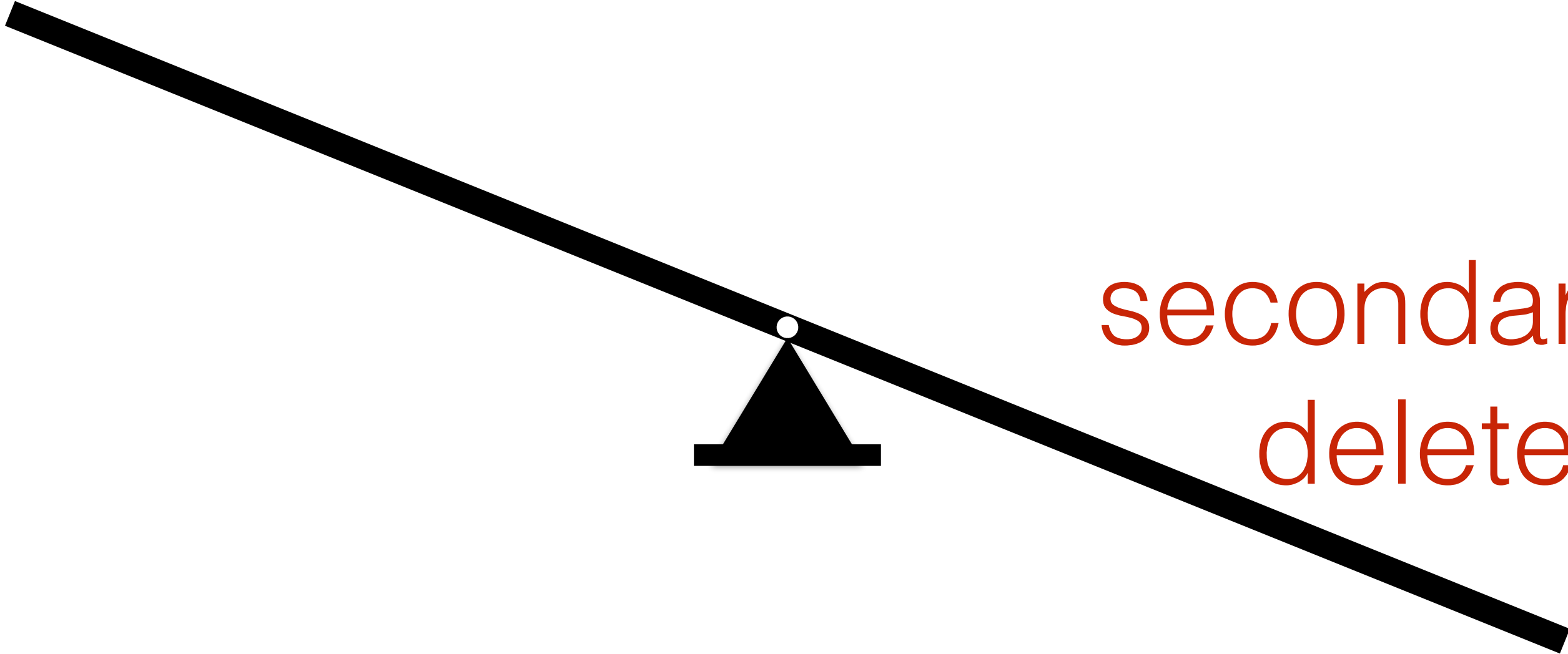
secondary range
delete cost



delete tile size



lookup
cost



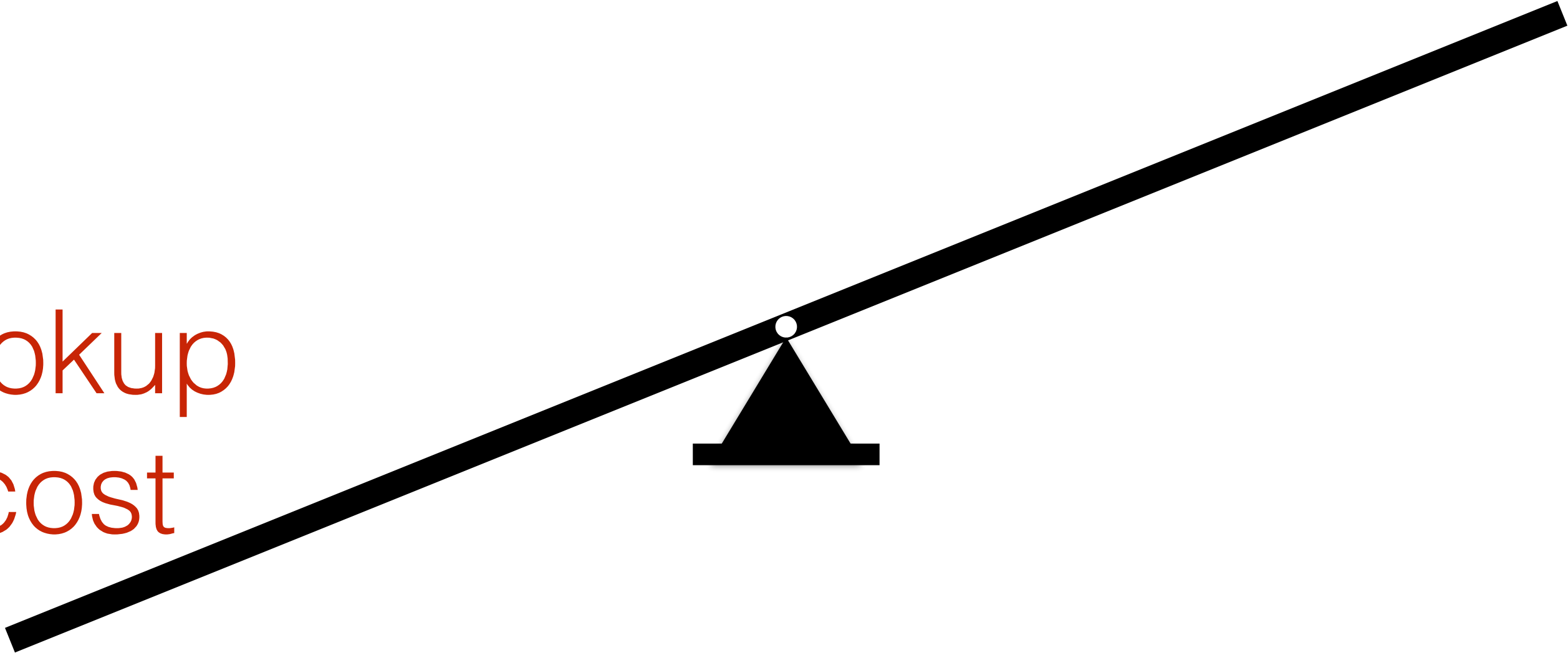
secondary range
delete cost

delete tile size



secondary range
delete cost

lookup
cost

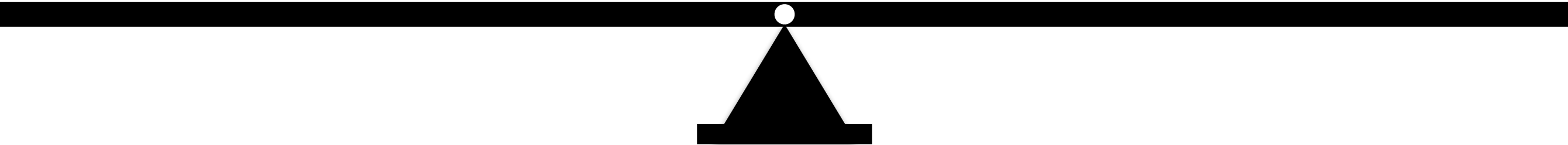


delete tile size

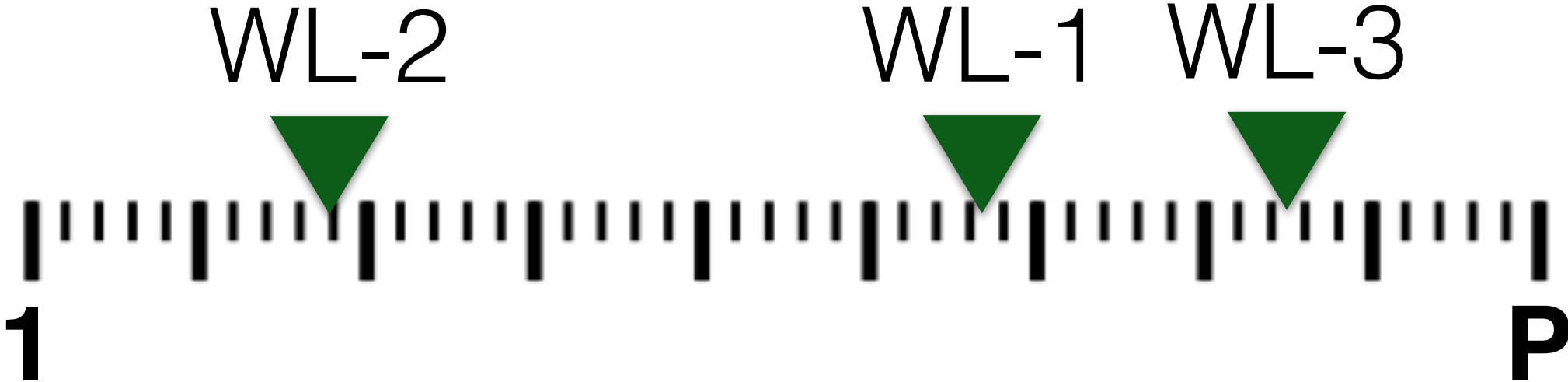


lookup
cost

secondary range
delete cost

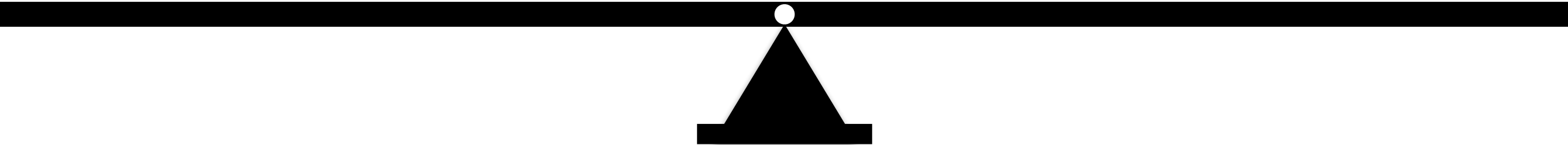


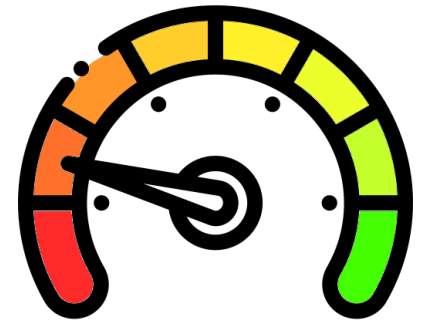
delete tile size



lookup
cost

secondary range
delete cost





suboptimal state of the art design
for workloads with deletes



FADE persists deletes timely
using latency-driven compactions



KiWi supports efficient
secondary range deletes
by key-interweaved data layout



Lethe strikes balance between
cost, performance, and latency

Thank You!

BOSTON
UNIVERSITY



midas.bu.edu/lethe