

Enabling Efficient Deletes in Log-Structured KV-Stores

BOSTON
UNIVERSITY

Subhadeep Sarkar

Enabling Efficient Deletes in Log-Structured KV-Stores

BOSTON
UNIVERSITY

Subhadeep Sarkar

Enabling Efficient Deletes in Log-Structured KV-Stores

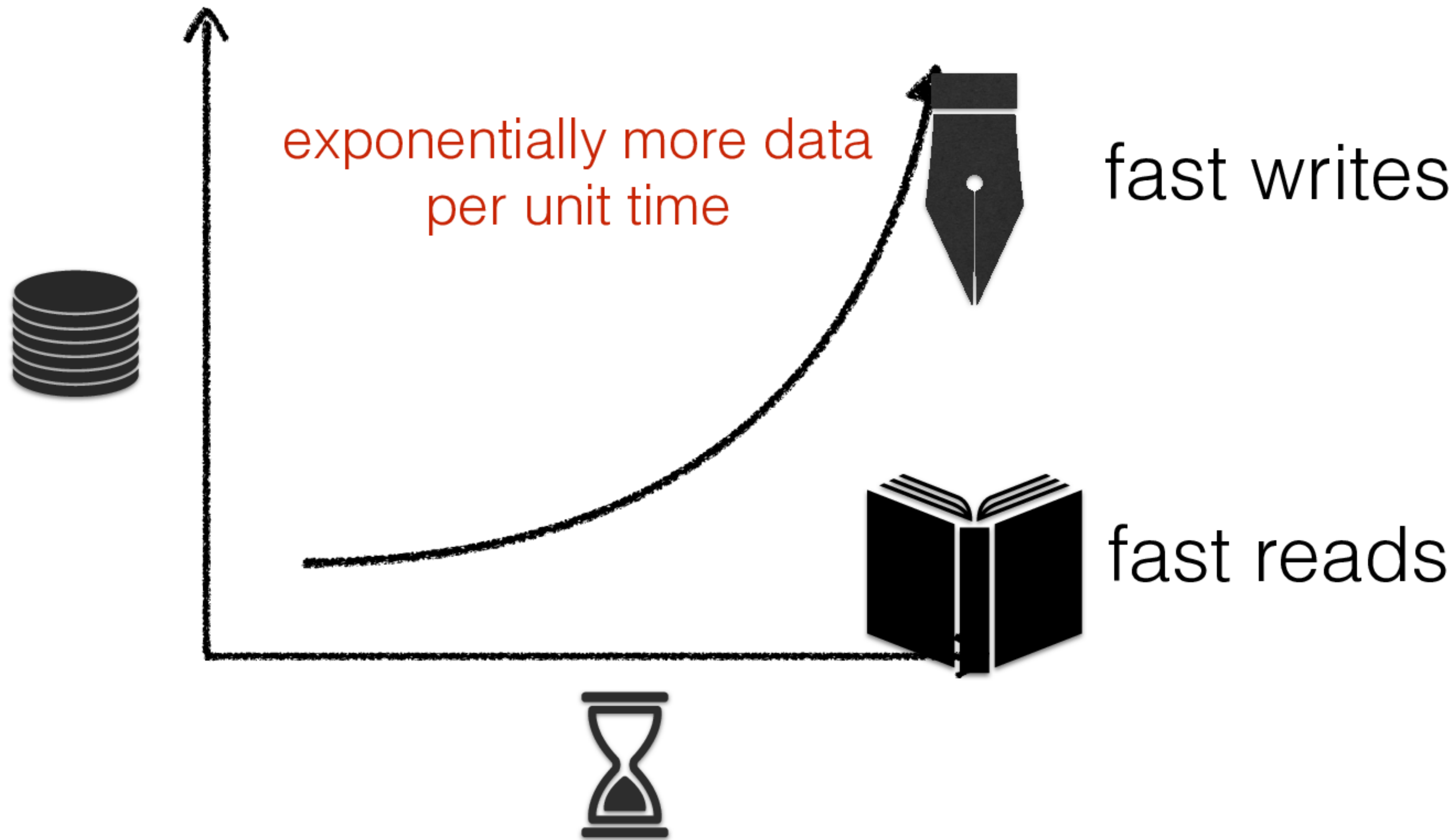
Enabling Efficient Deletes in Log-Structured KV-Stores

Enabling Efficient **Deletes** in **Log-Structured KV-Stores**

Enabling Efficient Deletes in Log-Structured KV-Stores

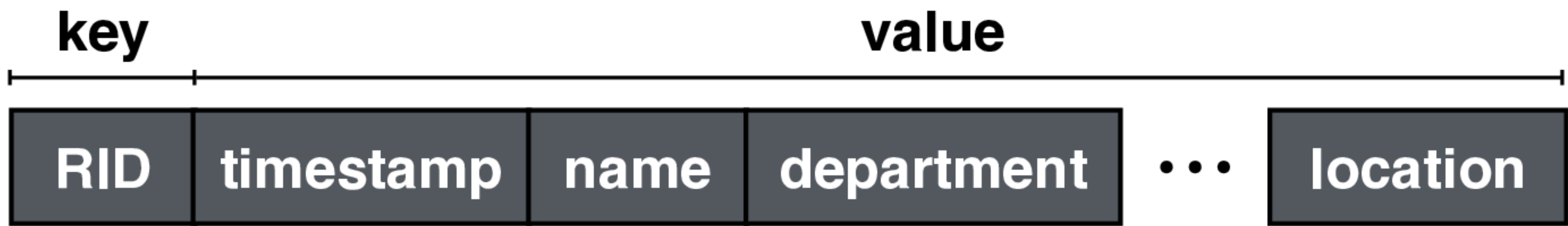


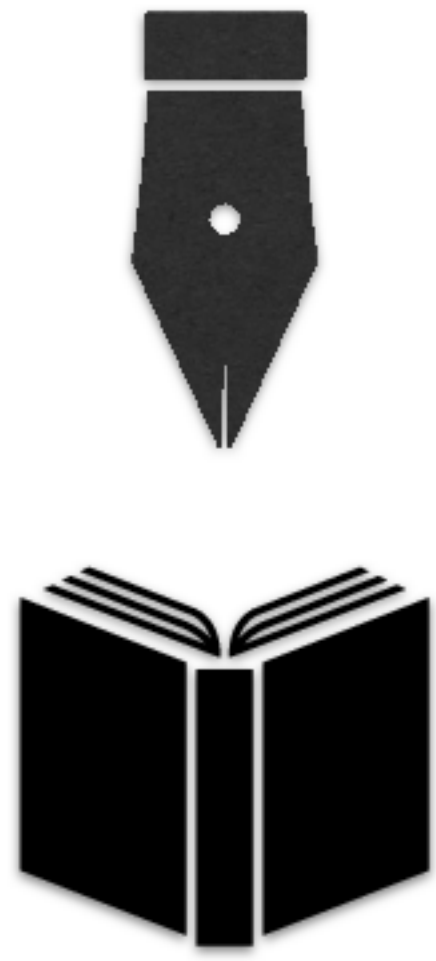
Lethe [*\ˈlē-thē*]
n: the goddess of
forgetfulness



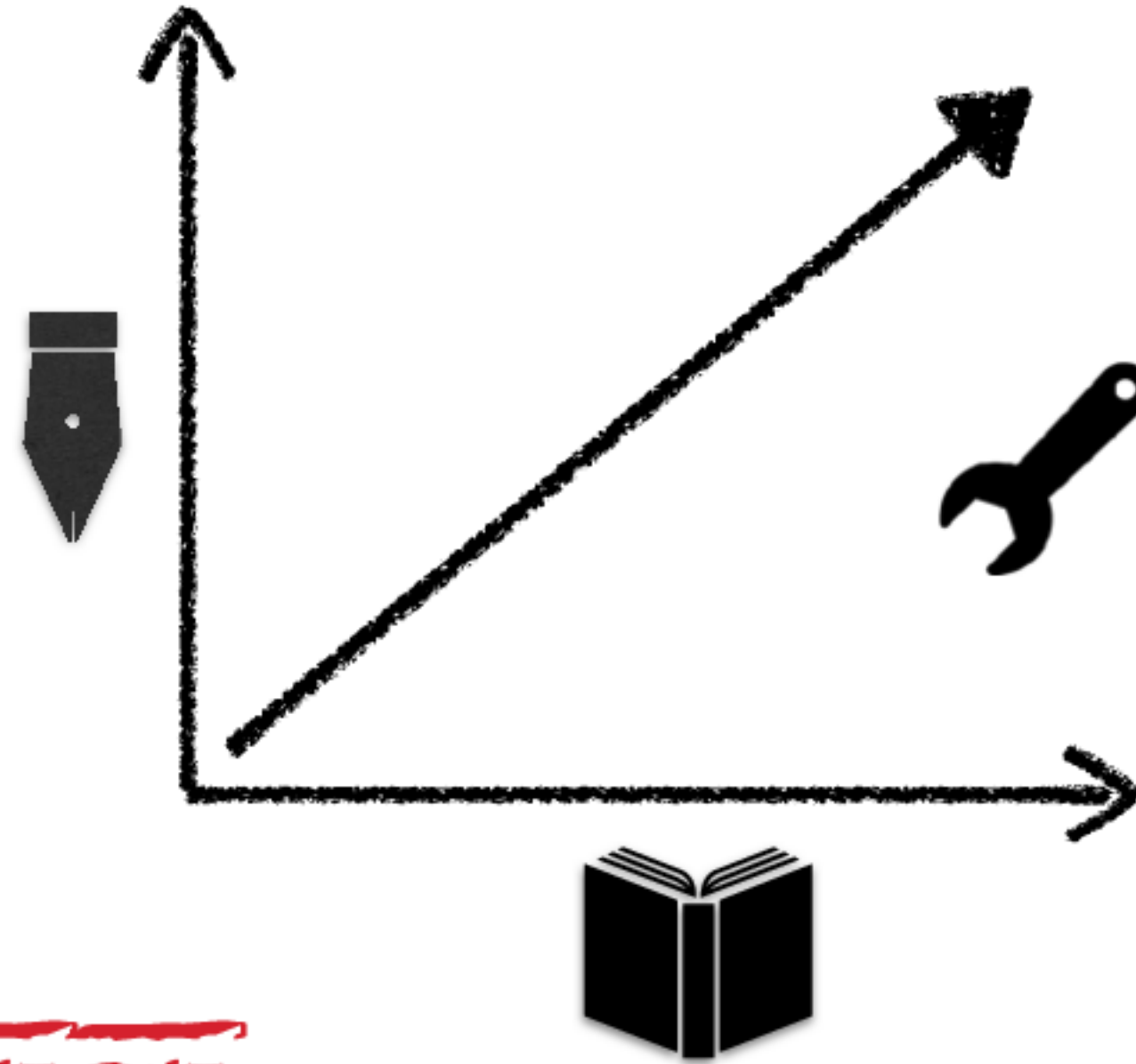
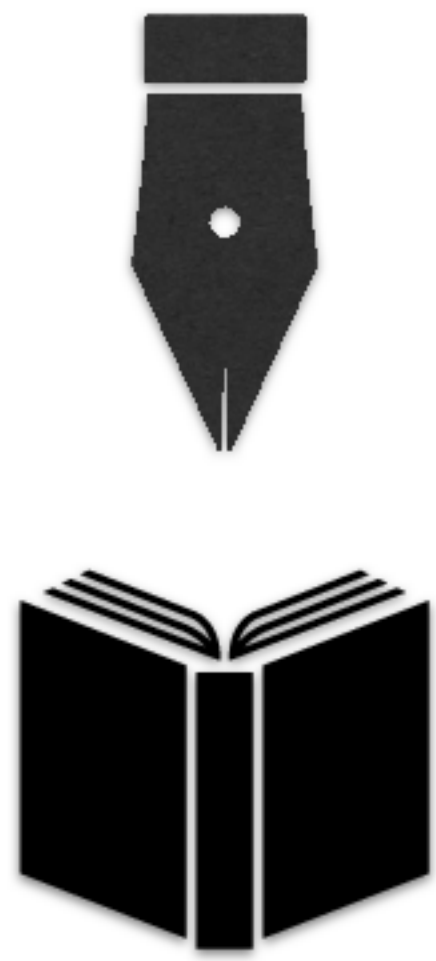
key-value pairs



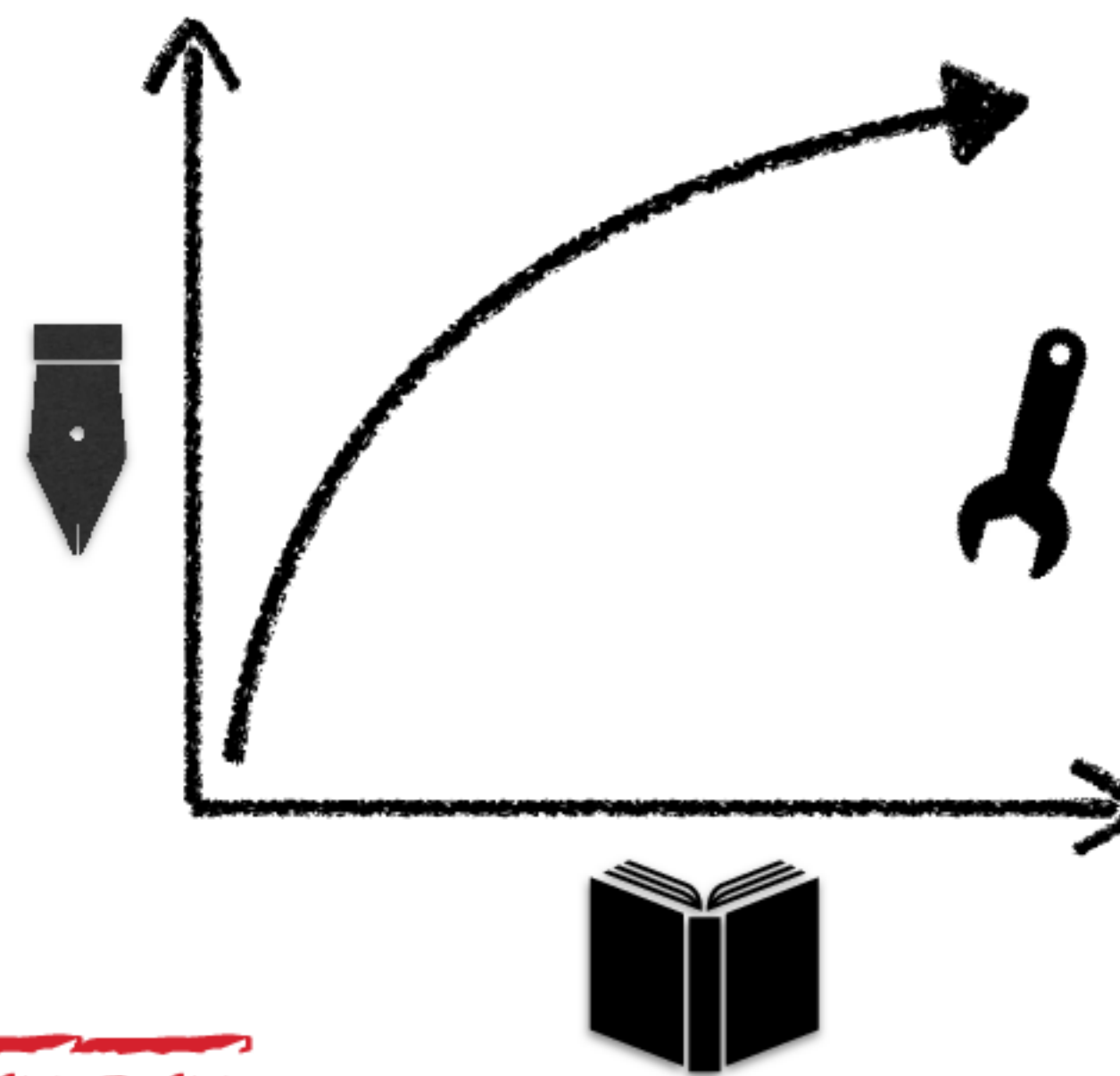
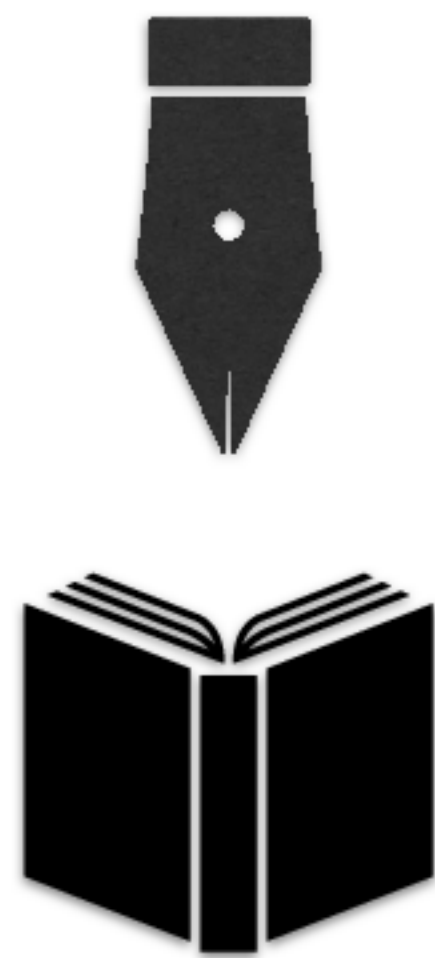




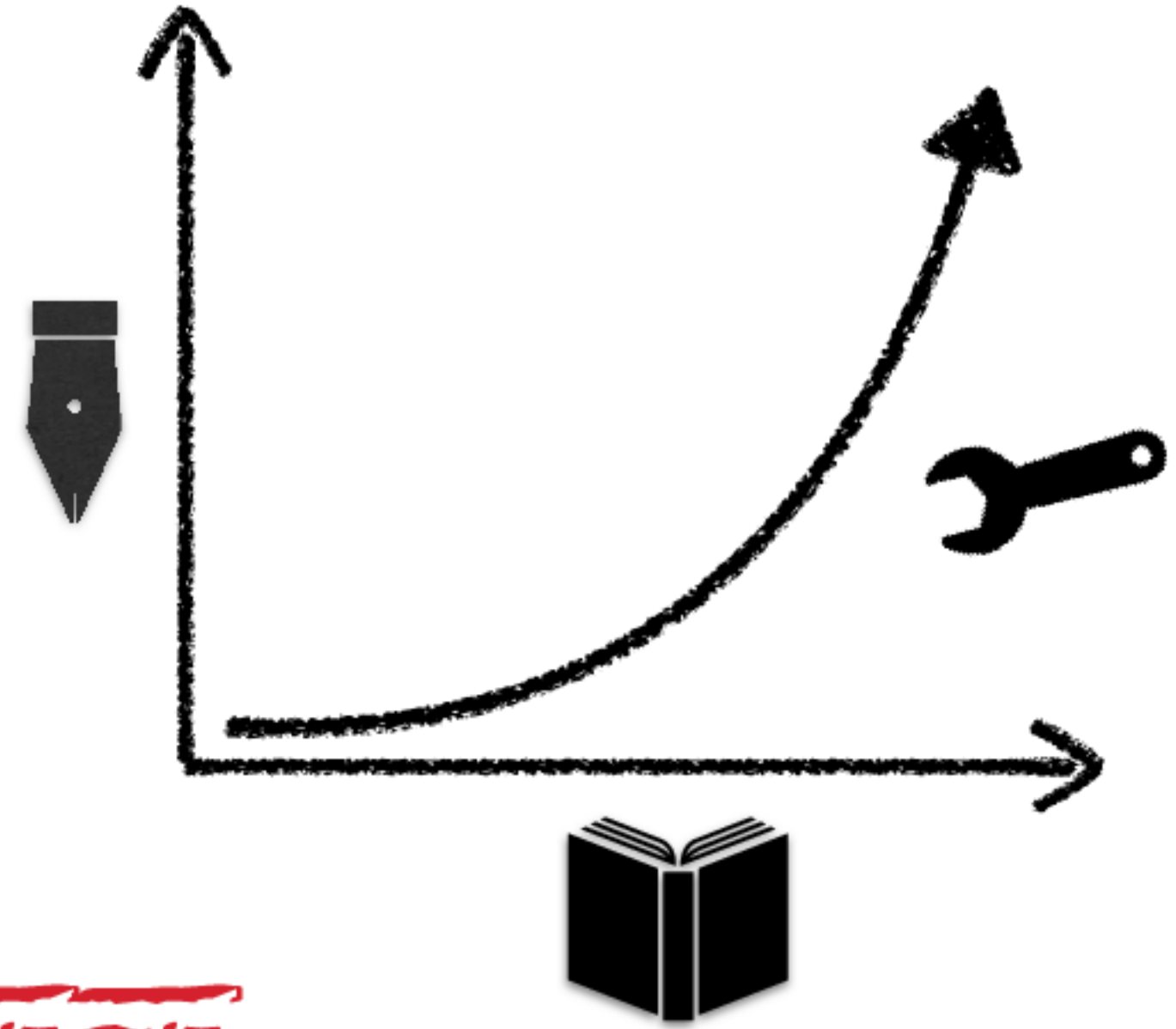
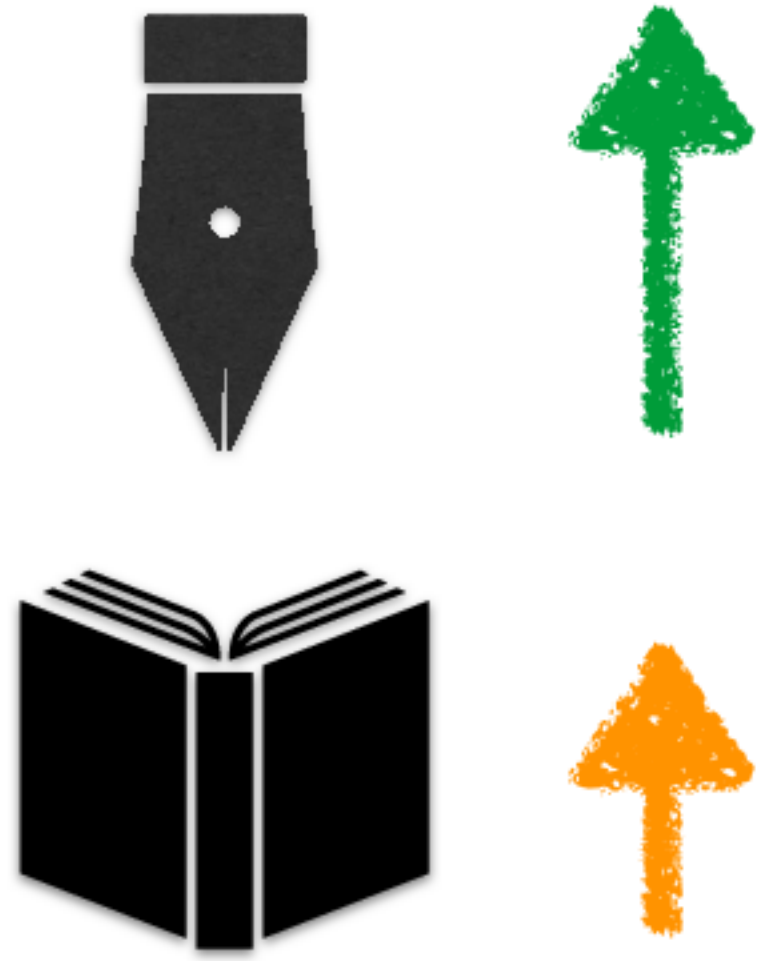
LSM-TREE



LSM-TIME



LSM-TIME



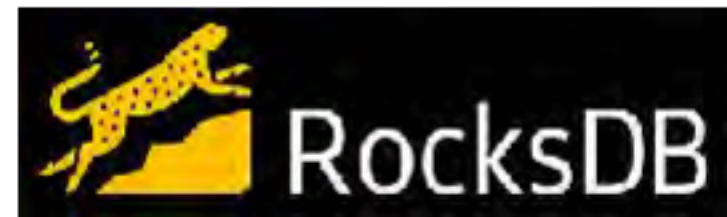
LSM-TREE



cassandra



levelDB



RocksDB

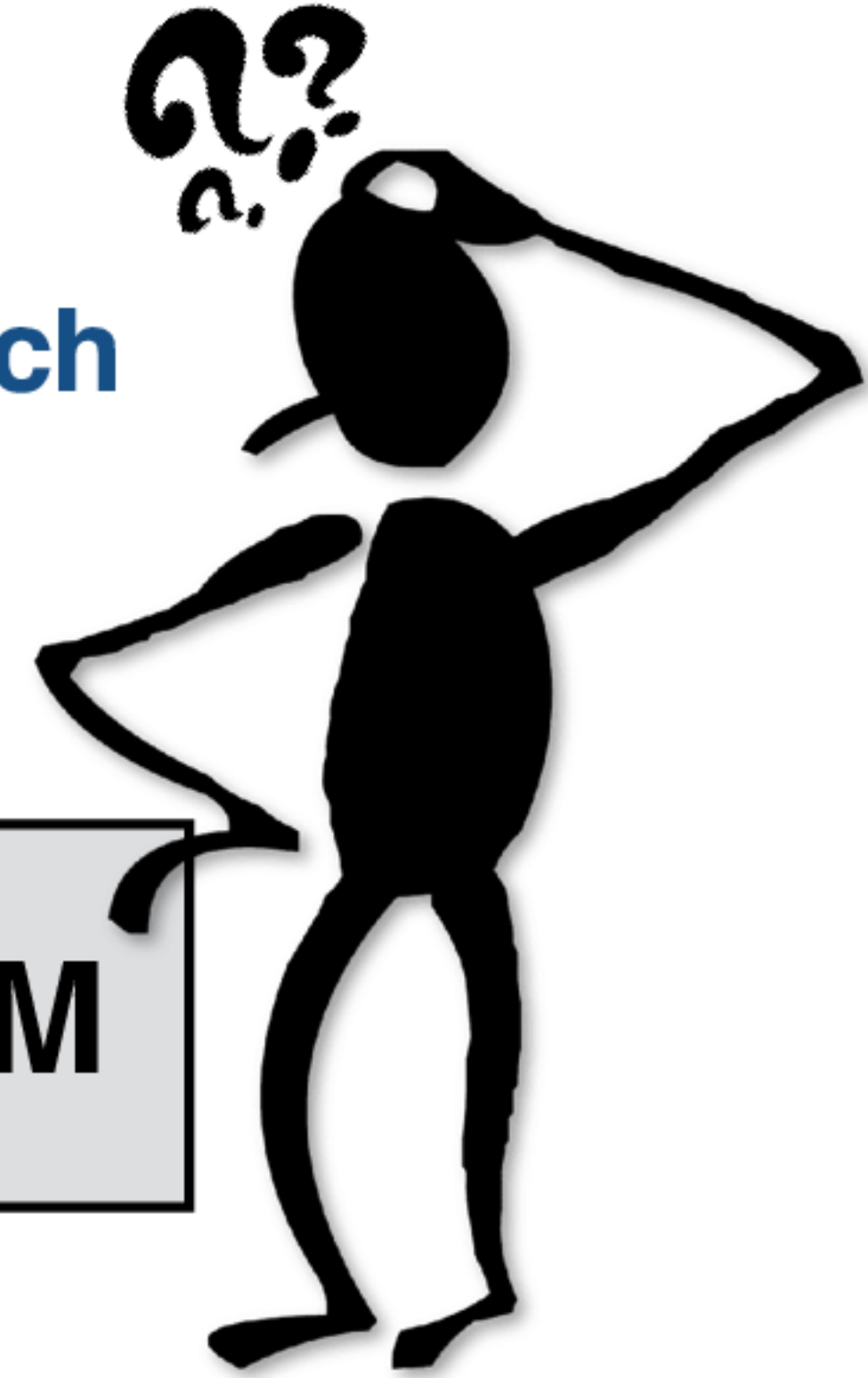


DynamoDB



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

TechCrunch
Feb '19



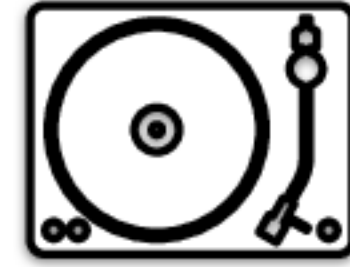
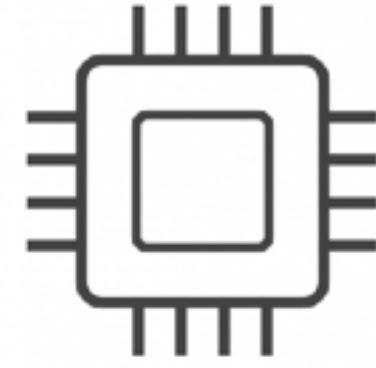
Small Datum
Jan '20

Deletes are fast and slow in an LSM

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

log-structured merge-tree

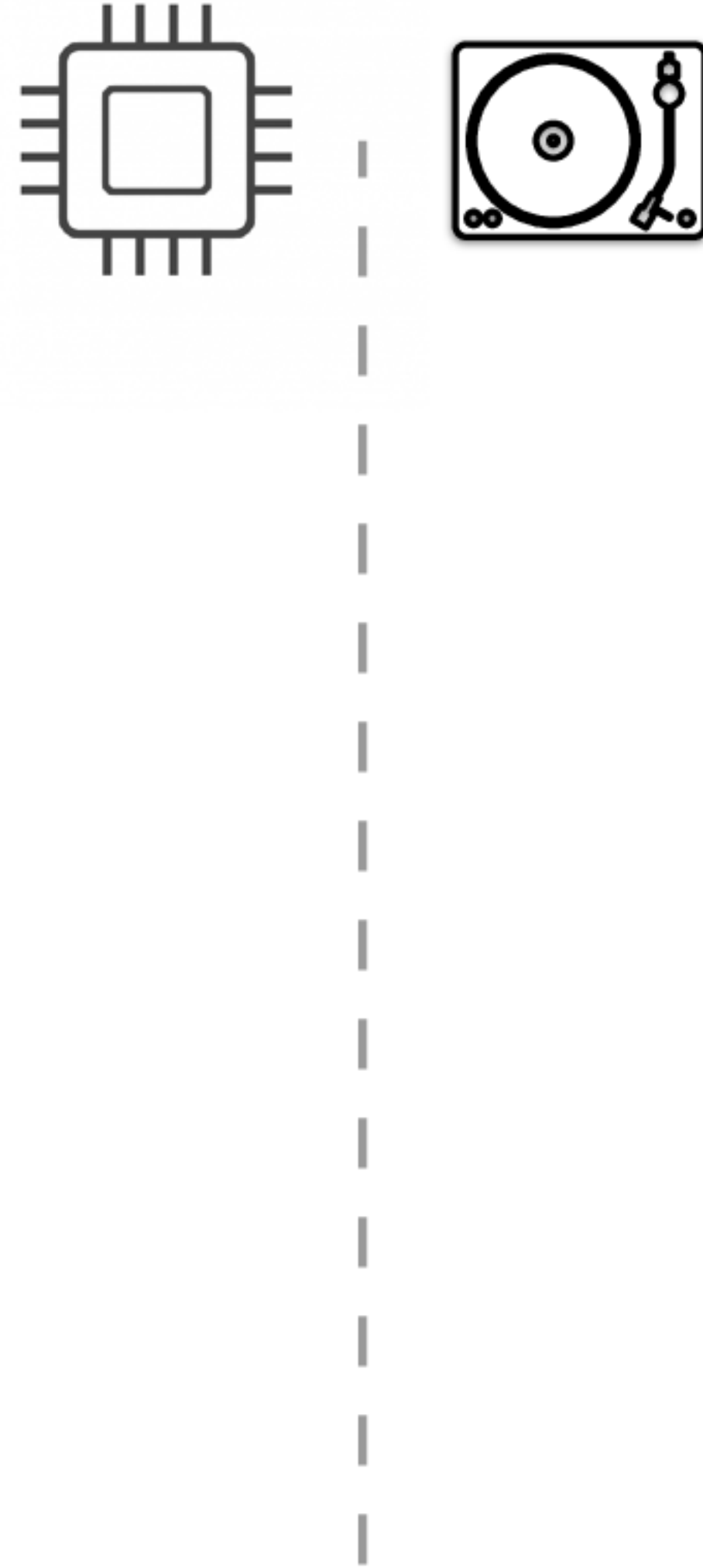
buffer



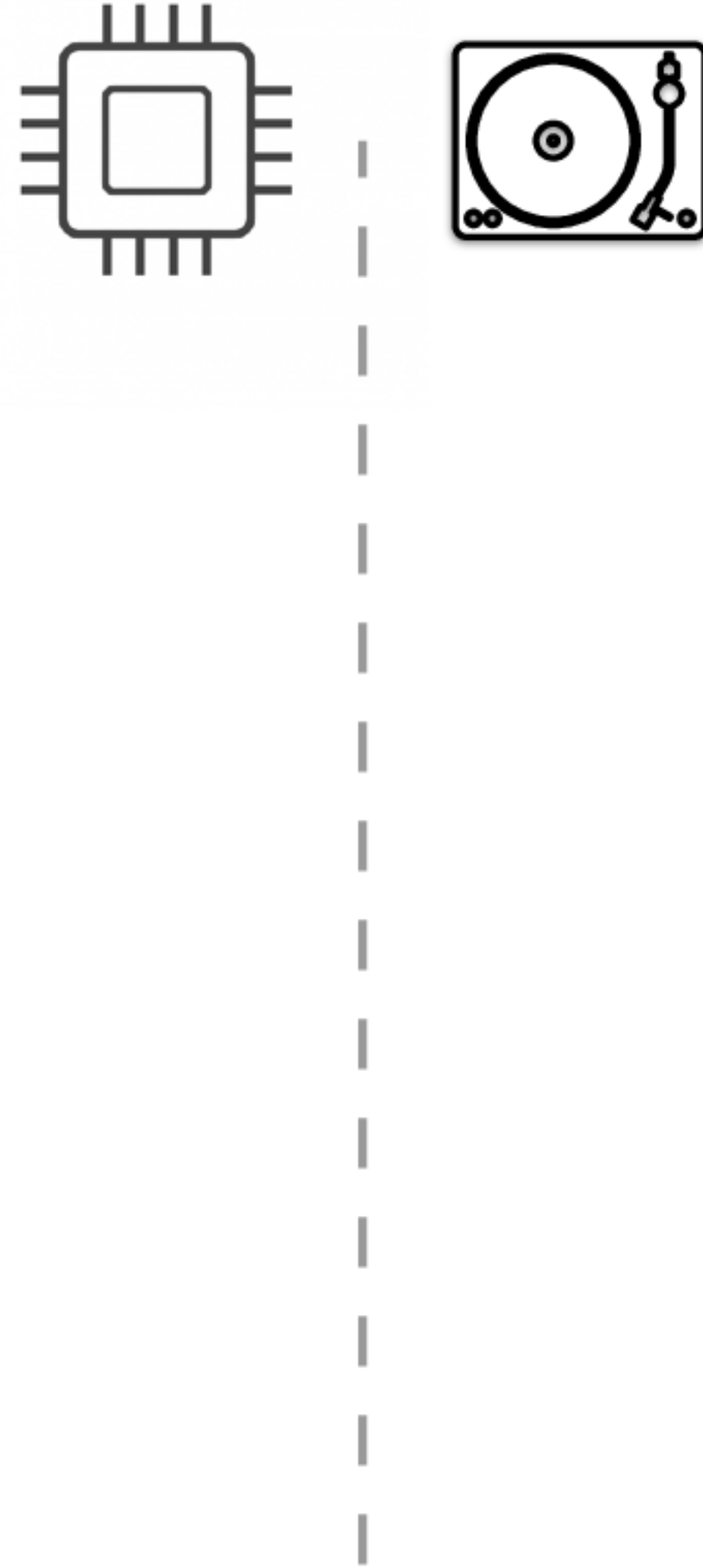
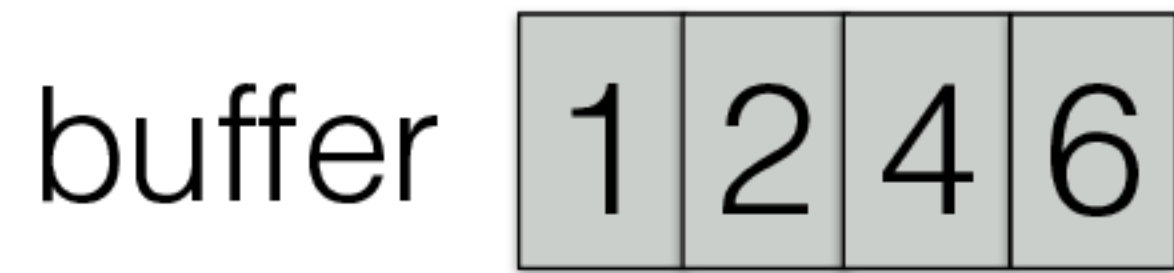
log-structured merge-tree

buffer

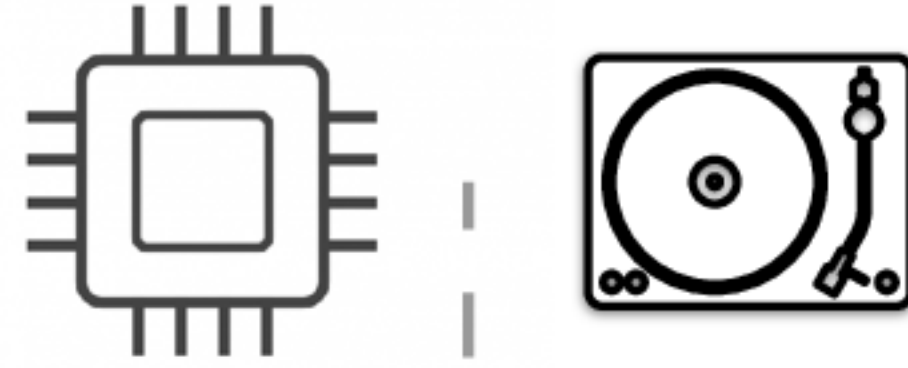
2	6	1	4
---	---	---	---



log-structured merge-tree



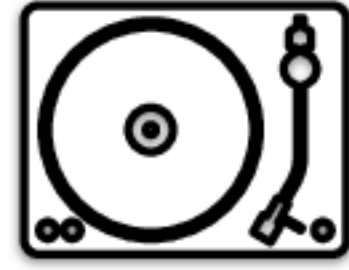
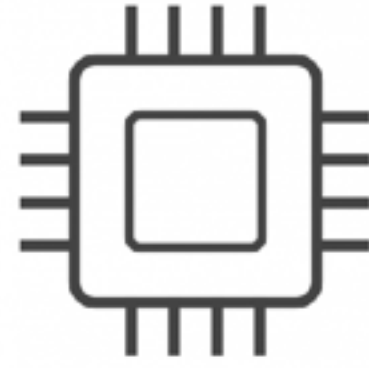
log-structured merge-tree



buffer



log-structured merge-tree



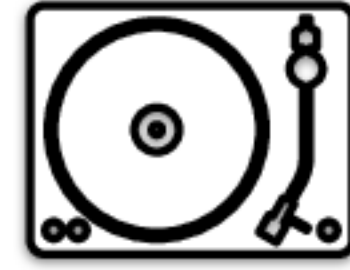
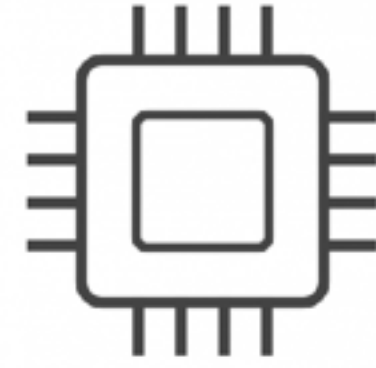
buffer



L1



log-structured merge-tree



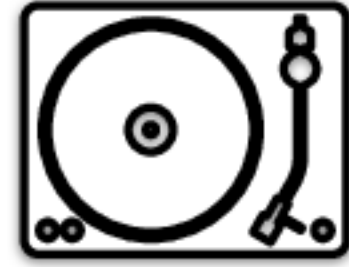
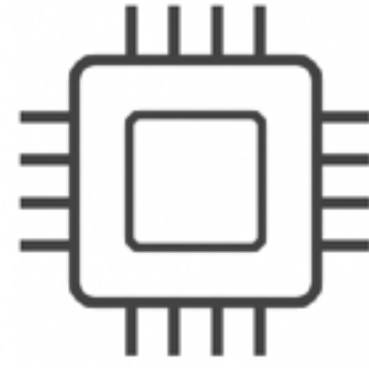
buffer



L1



log-structured merge-tree



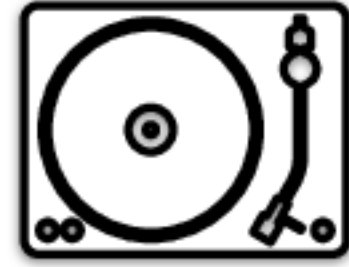
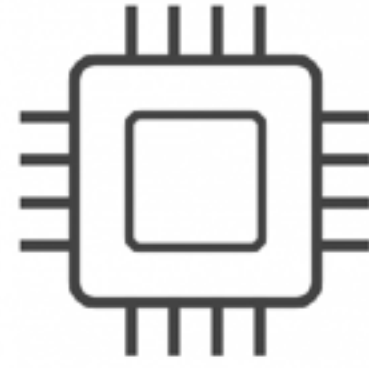
buffer



L1



log-structured merge-tree



buffer



L1

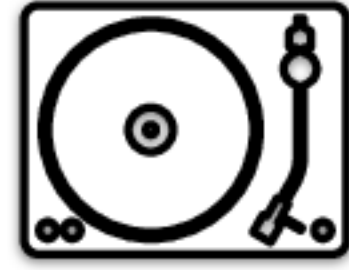
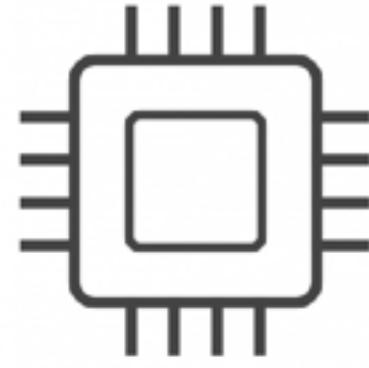
L2



compaction



log-structured merge-tree



buffer



L1



size ratio = T

L2

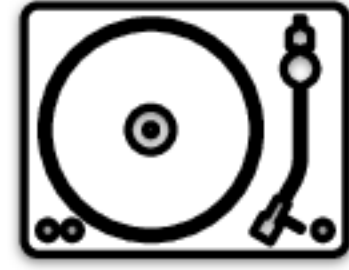
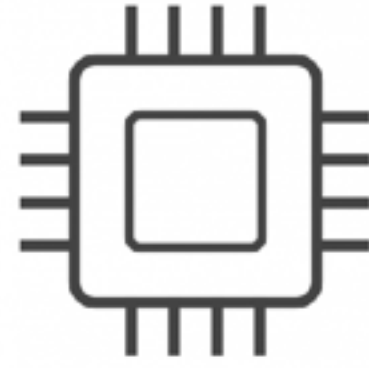


L3



exponentially larger capacity

log-structured merge-tree



buffer



L1



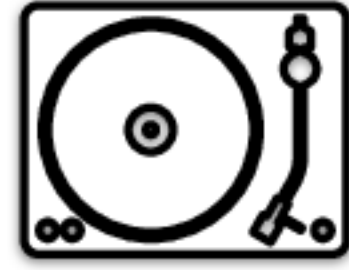
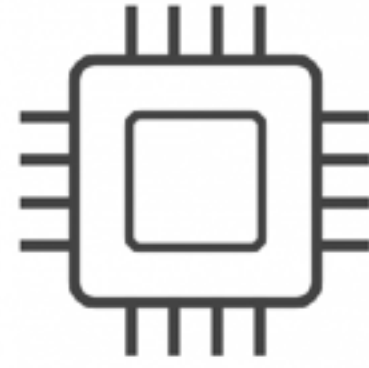
L2



L3



log-structured merge-tree



buffer



L1



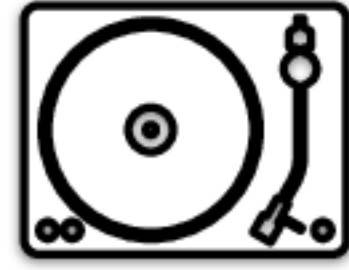
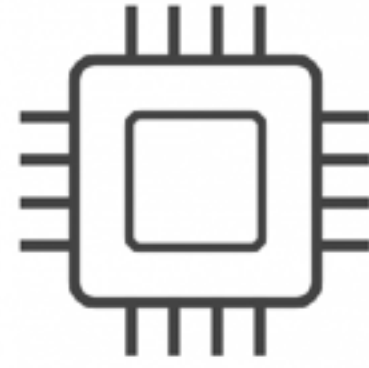
L2



L3



log-structured merge-tree



buffer



L1



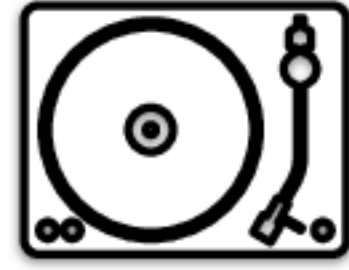
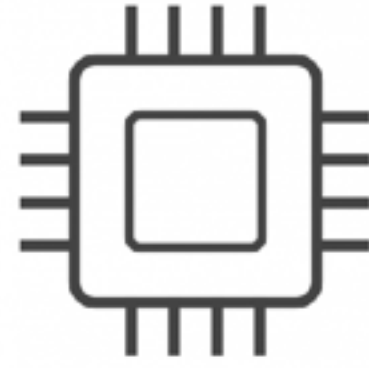
L2



L3



log-structured merge-tree



buffer



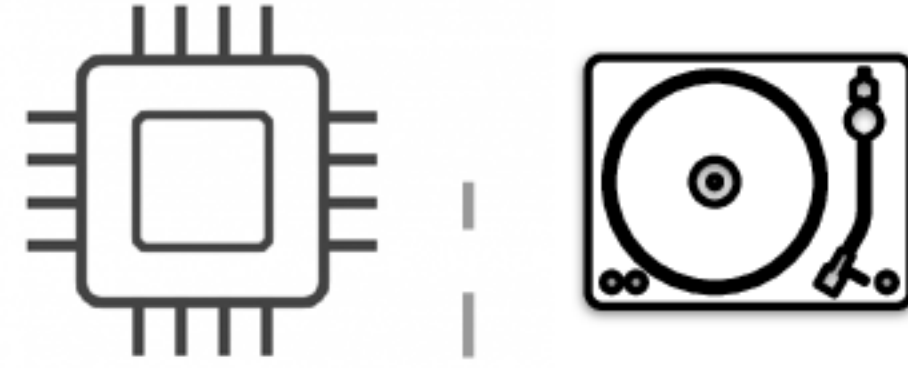
L1

L2

L3



log-structured merge-tree



buffer



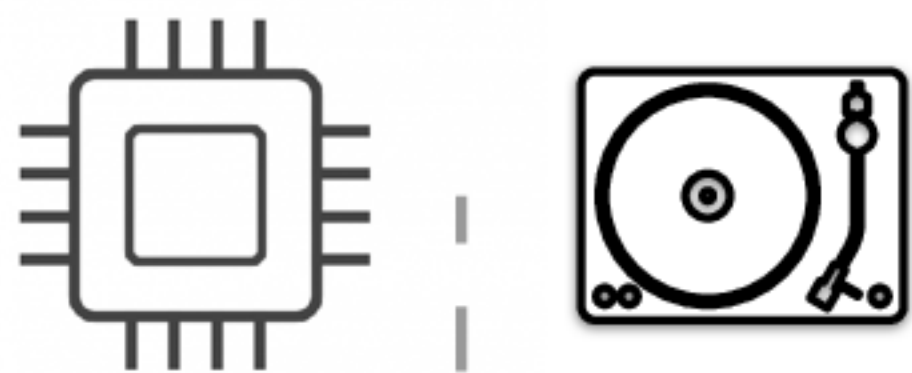
L1

L2

L3



log-structured merge-tree



buffer



L1

L2

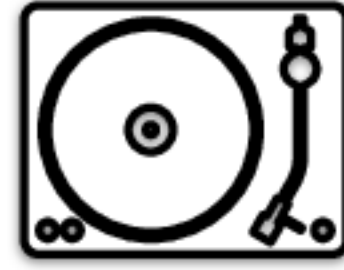
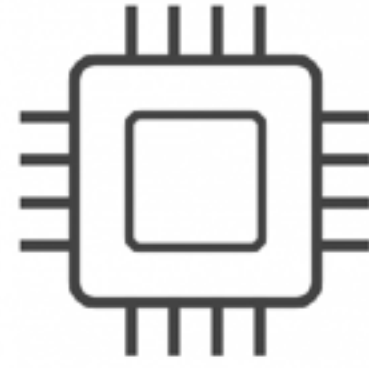
L3

L4

burst of I/Os
prolonged write stalls



log-structured merge-tree



partial compaction

buffer



L1



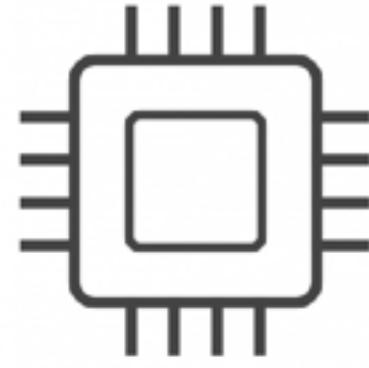
L2



L3



log-structured merge-tree



partial compaction

buffer



L1



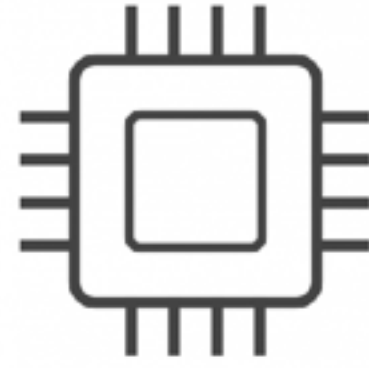
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



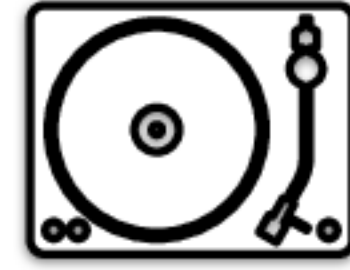
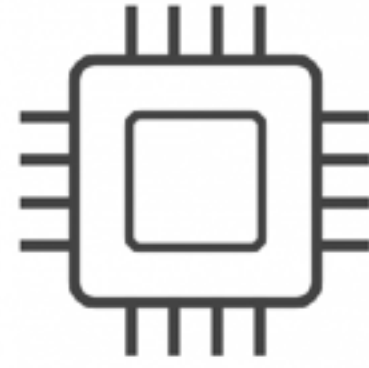
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



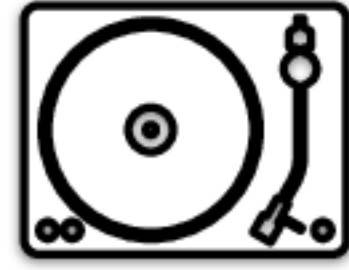
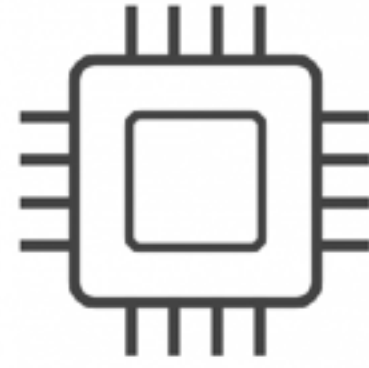
L2



L3



log-structured merge-tree



partial compaction

buffer



L1



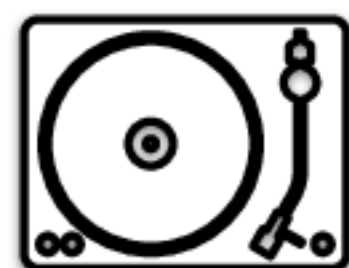
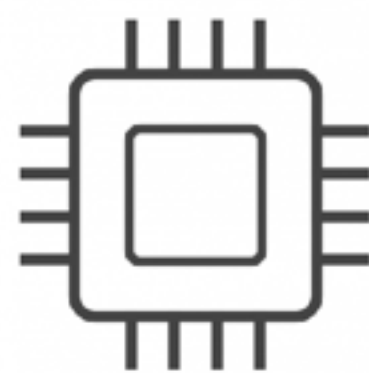
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



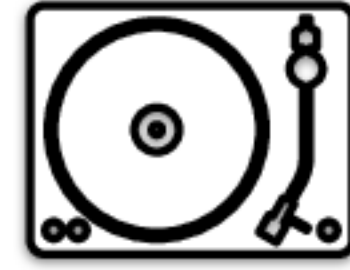
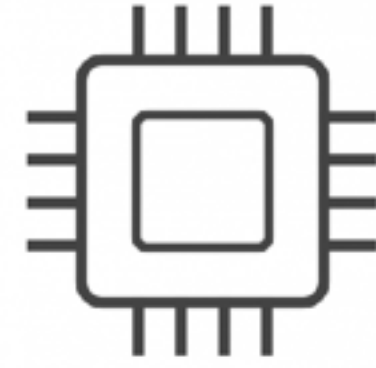
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



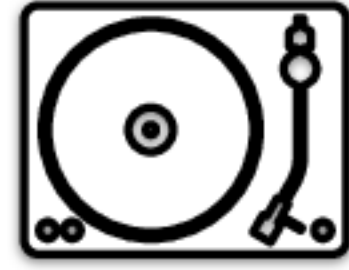
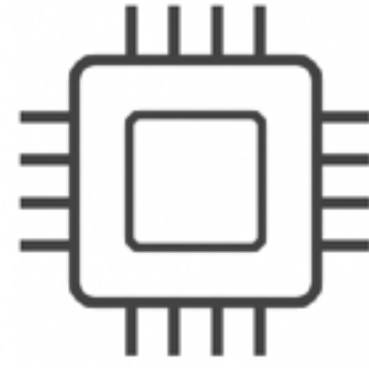
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



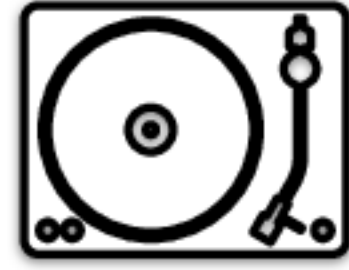
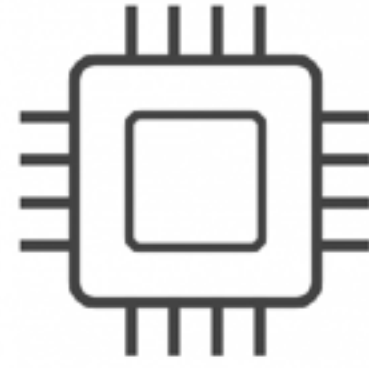
L2



L3



log-structured merge-tree



partial compaction

buffer



L1



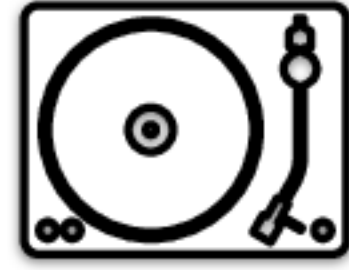
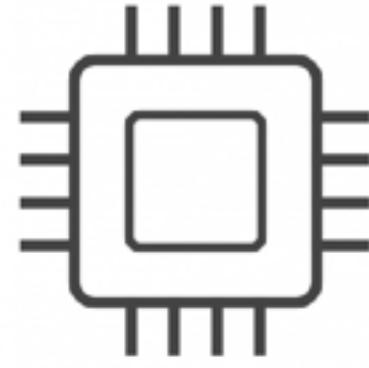
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



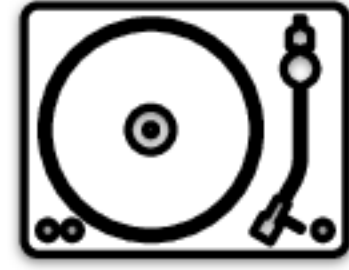
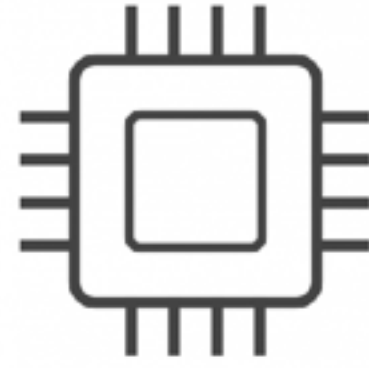
L2



L3



log-structured merge-tree



buffer



partial compaction

L1



L2

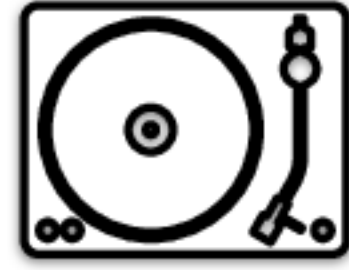
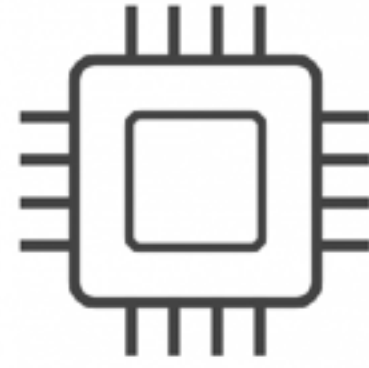


L3



amortized compaction cost

log-structured merge-tree



buffer



L1



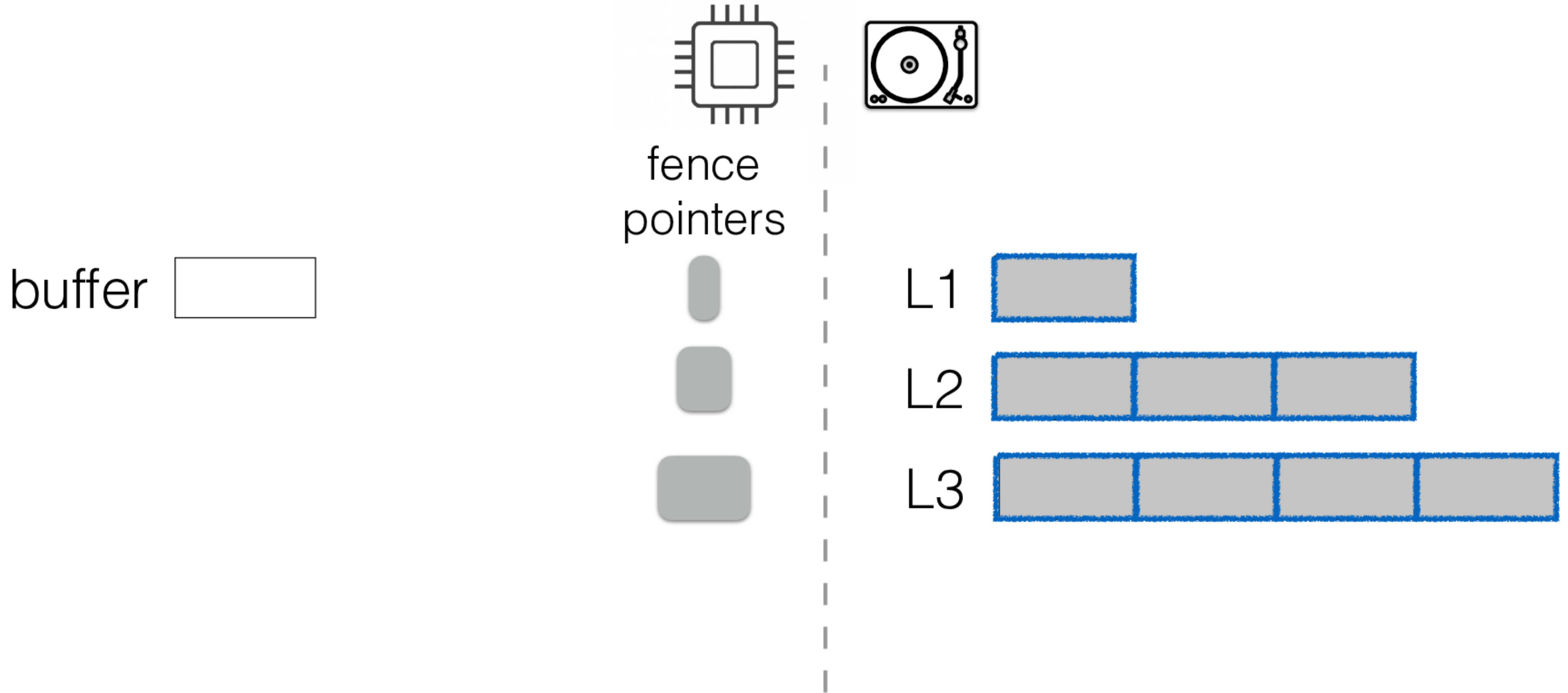
L2



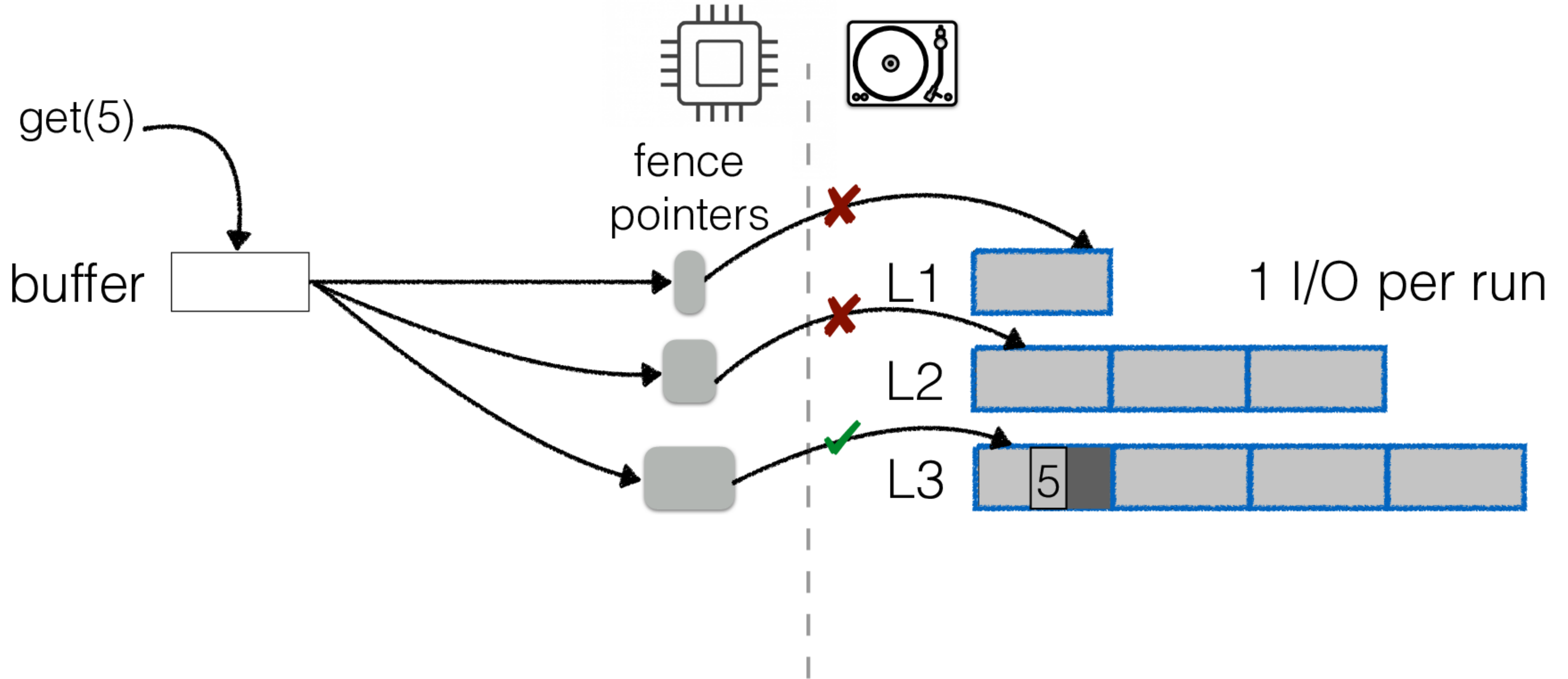
L3



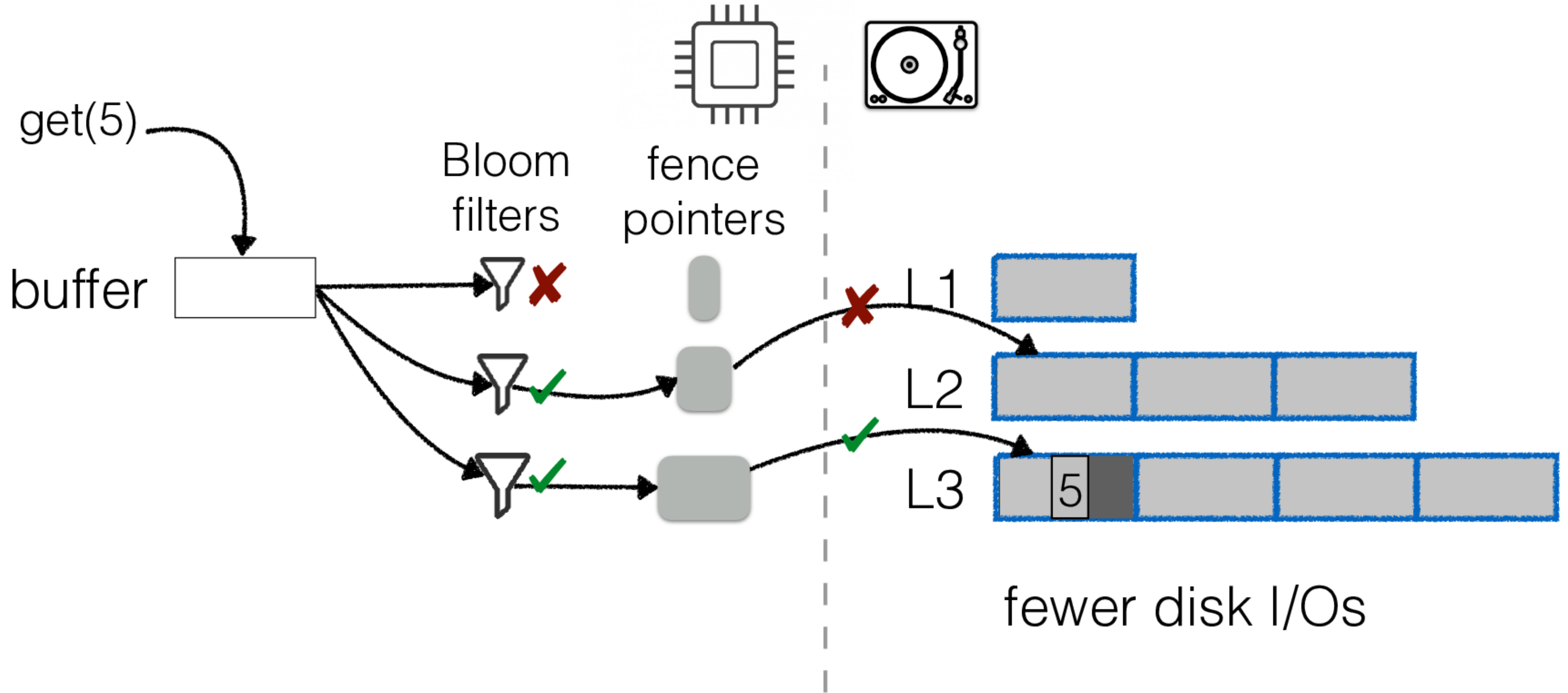
log-structured merge-tree



log-structured merge-tree



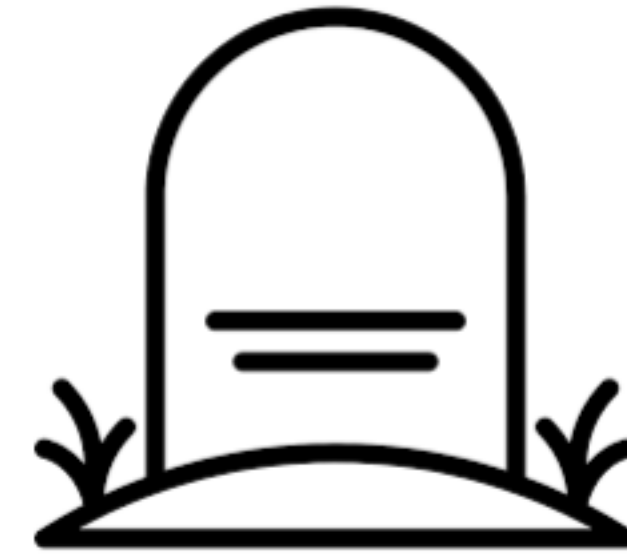
log-structured merge-tree



Now, let's talk about deletes!

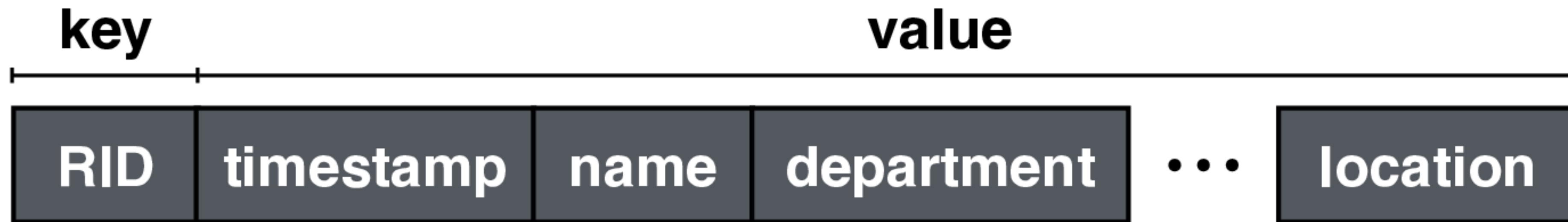
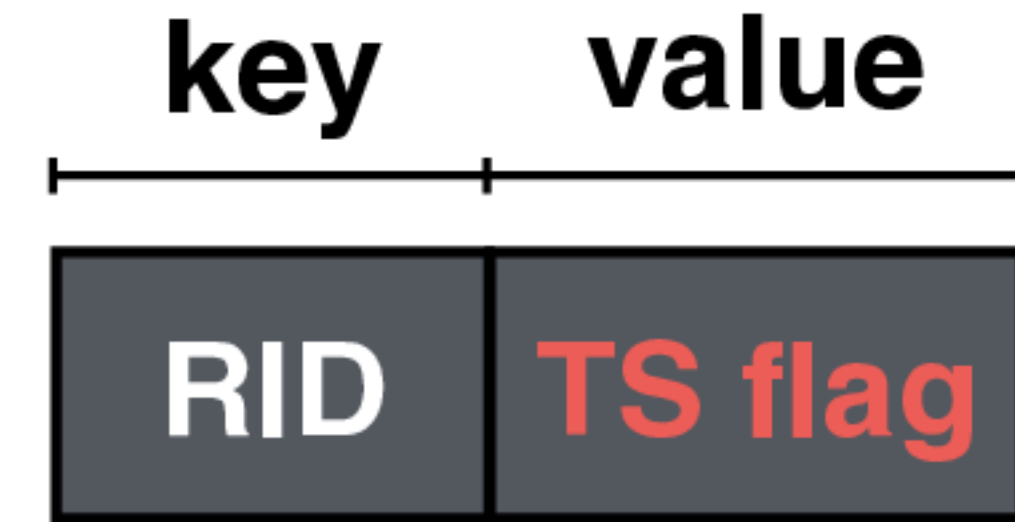
deletes in LSM-tree

delete



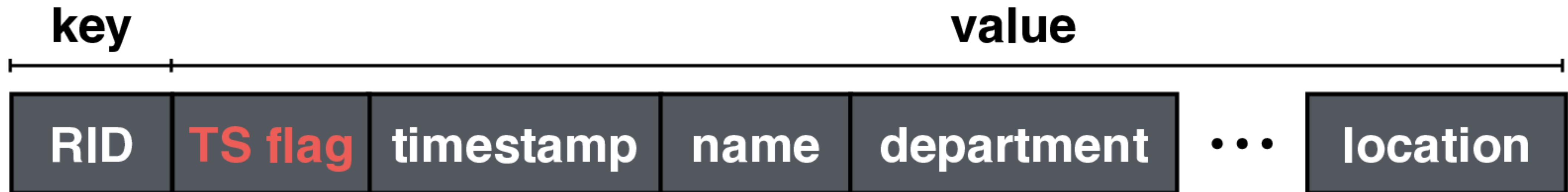
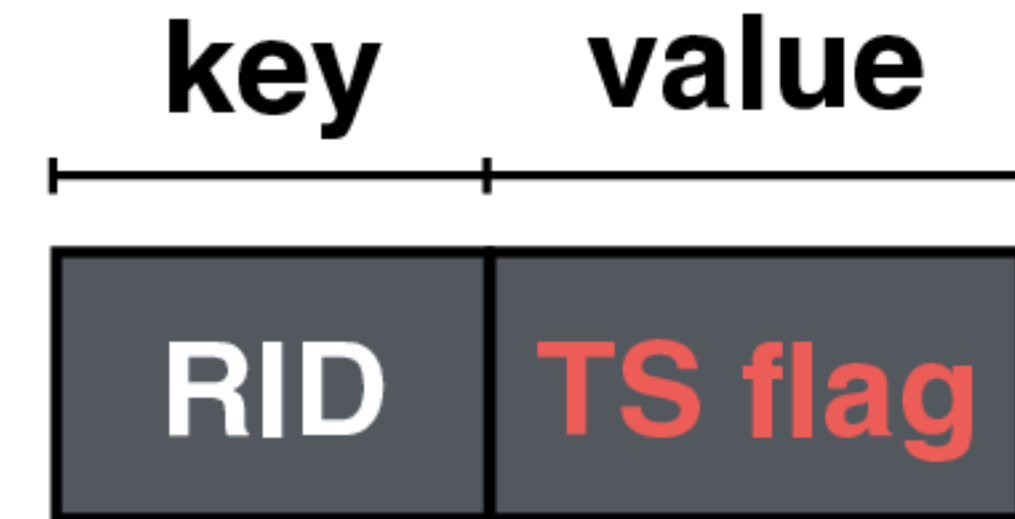
deletes in LSM-tree

delete := insert tombstone

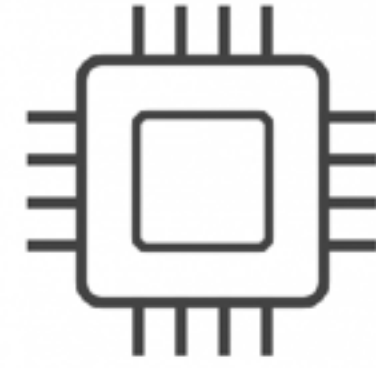
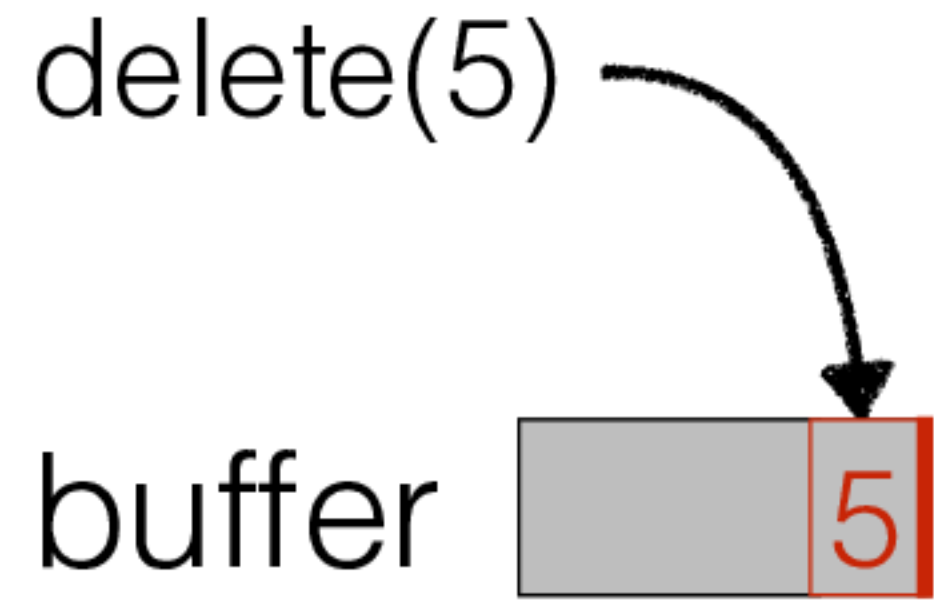


deletes in LSM-tree

delete := insert tombstone



deletes in LSM-tree



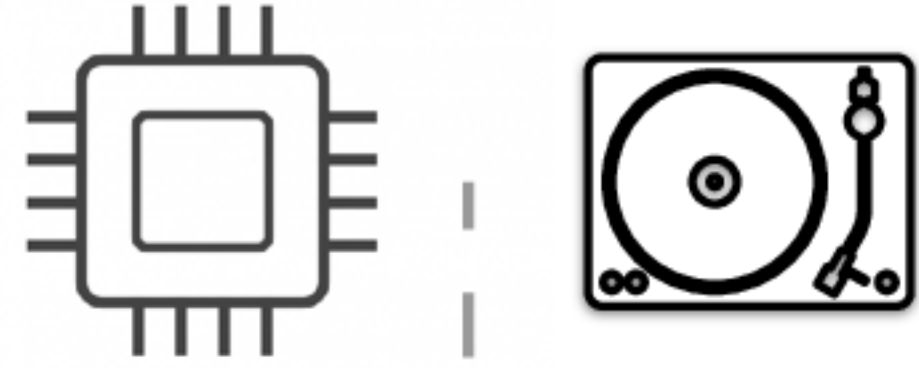
L1

L2

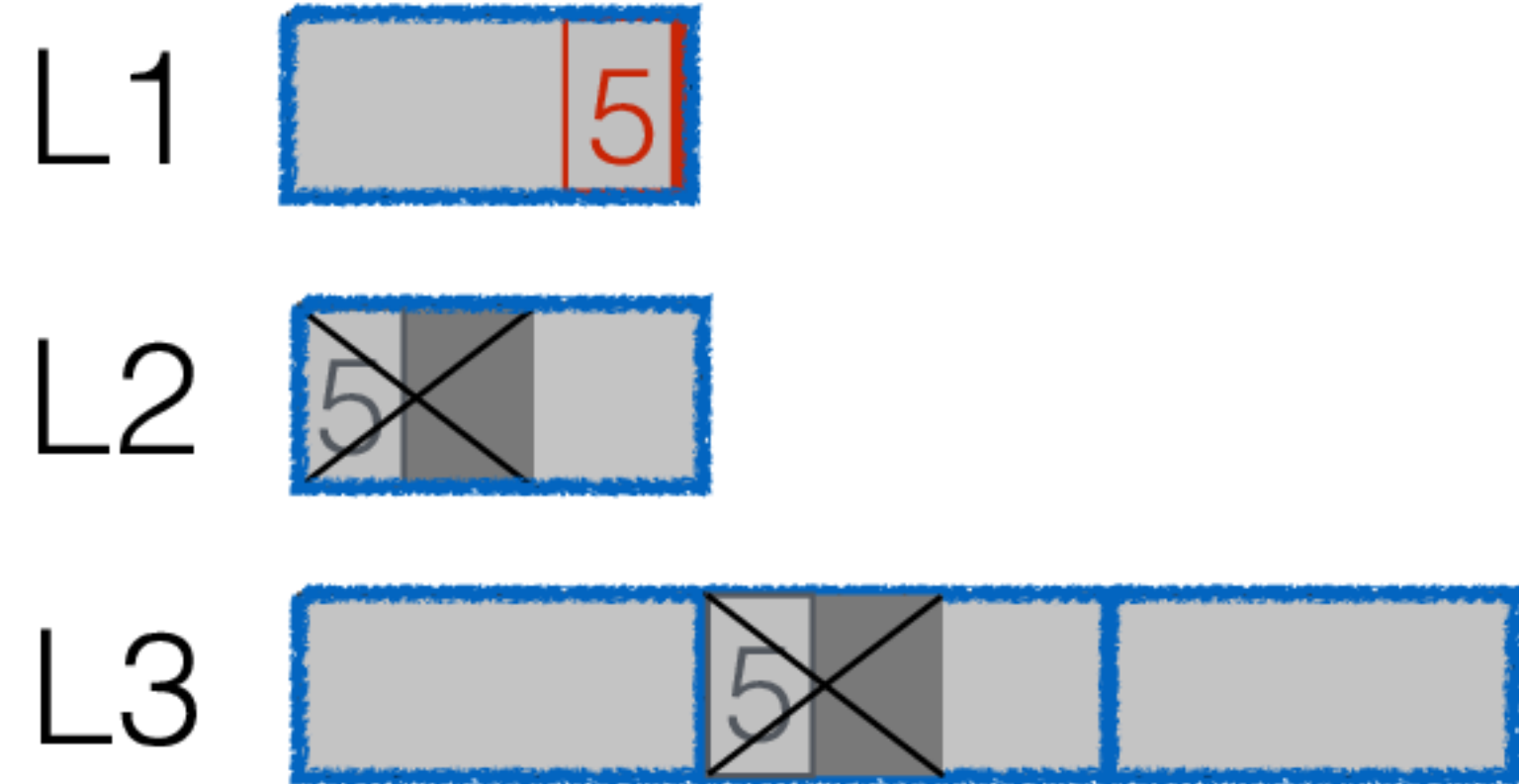
L3



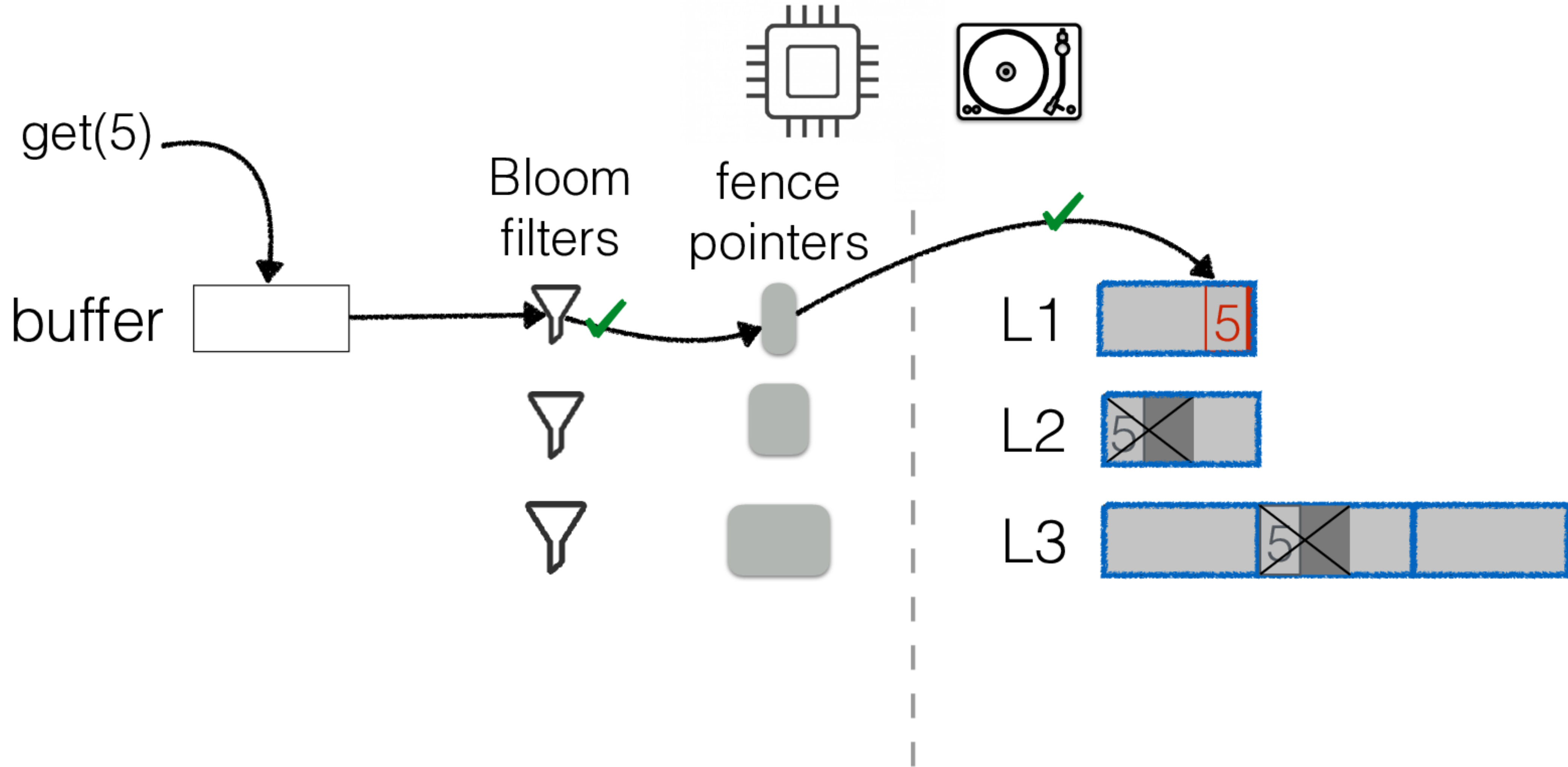
deletes in LSM-tree



buffer 



deletes in LSM-tree



the problems

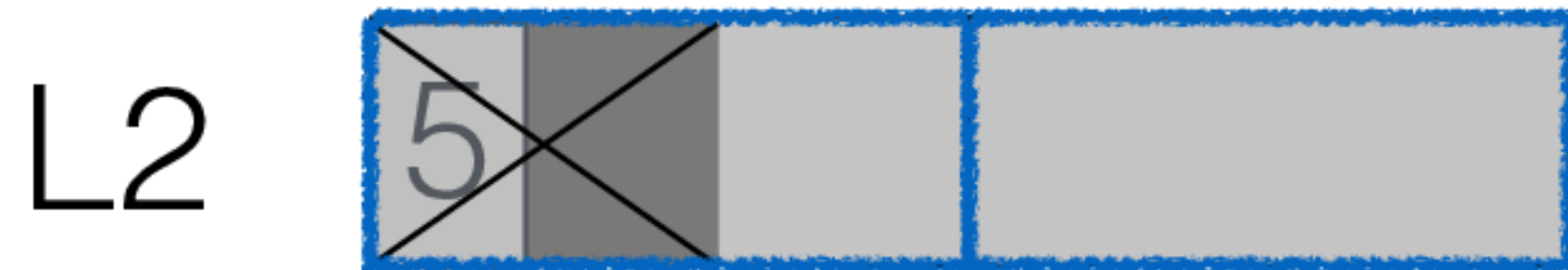


the problems



out-of-place deletes

out-of-place deletes



space amplification ✖



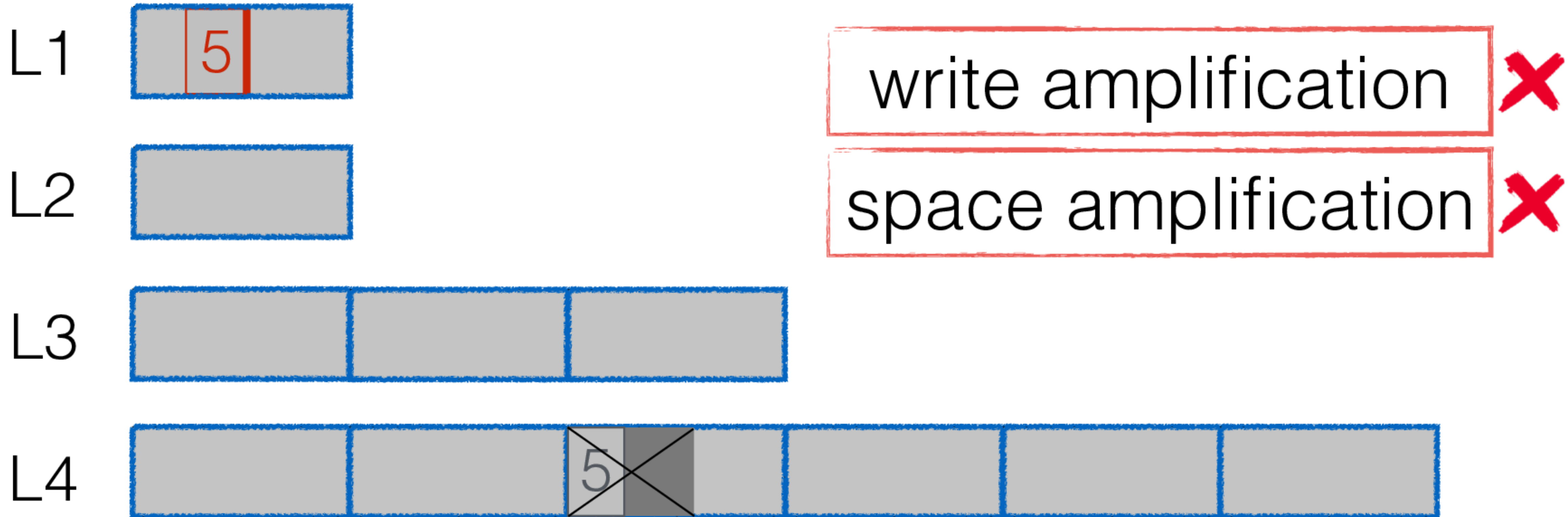
out-of-place deletes



space amplification ✖



out-of-place deletes

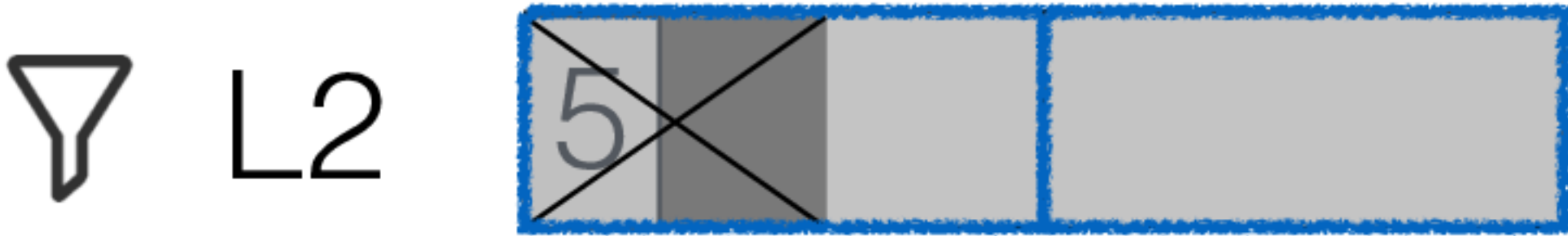


out-of-place deletes

Bloom filters



write amplification ❌



space amplification ❌

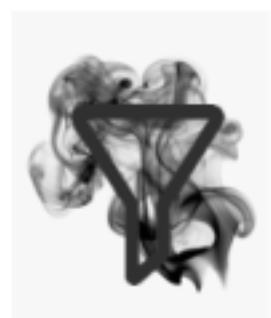


out-of-place deletes

Bloom filters



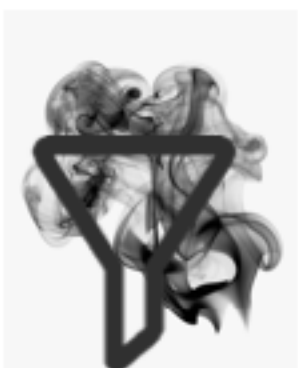
L1



L2



L3



L4



poor read perf. ❌

write amplification ❌

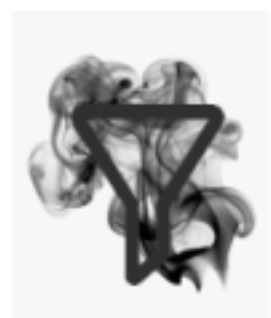
space amplification ❌

out-of-place deletes

Bloom filters



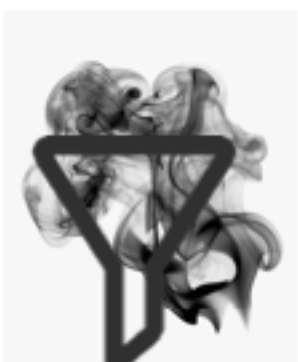
L1



L2



L3



L4



poor read perf. ❌

write amplification ❌

space amplification ❌

the problems

poor read perf.

write amplification

space amplification





delete persistence latency

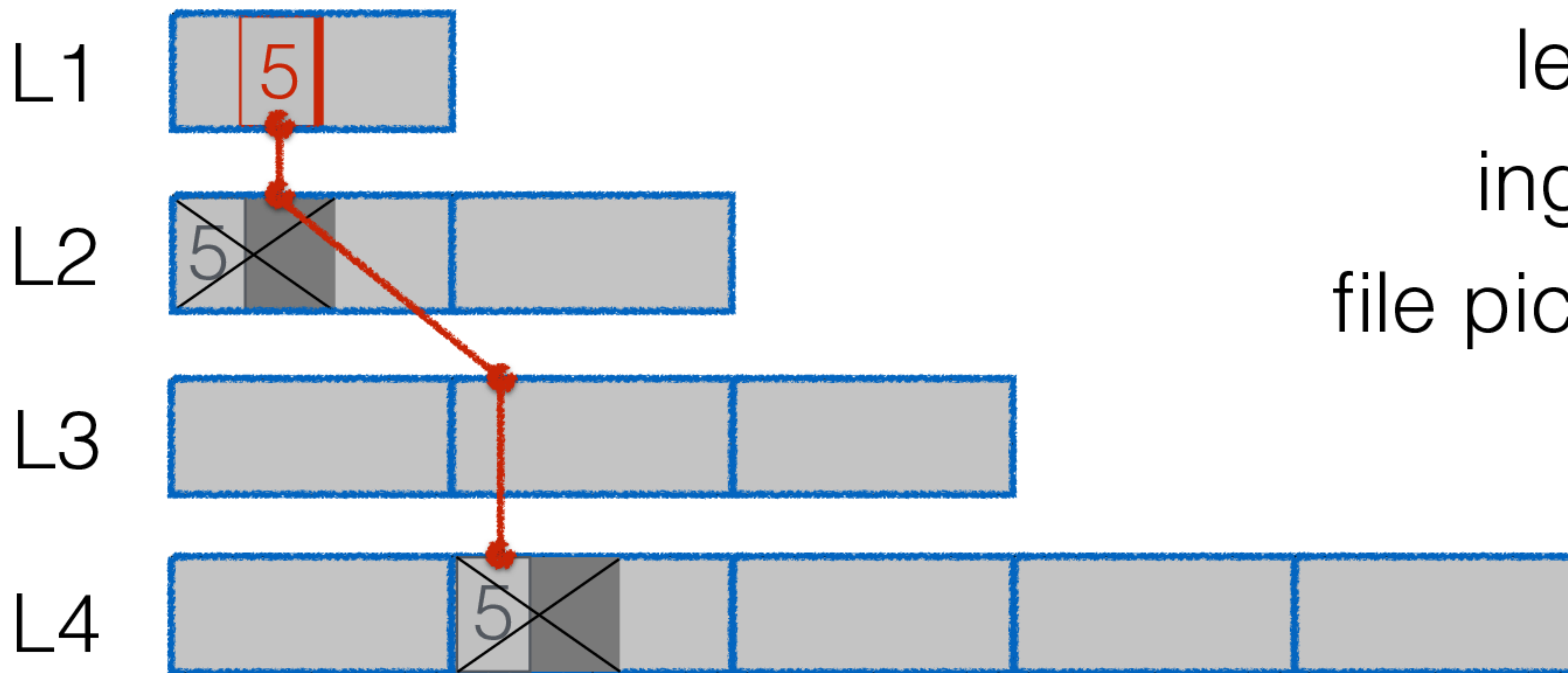
delete persistence latency

size ratio

levels in tree

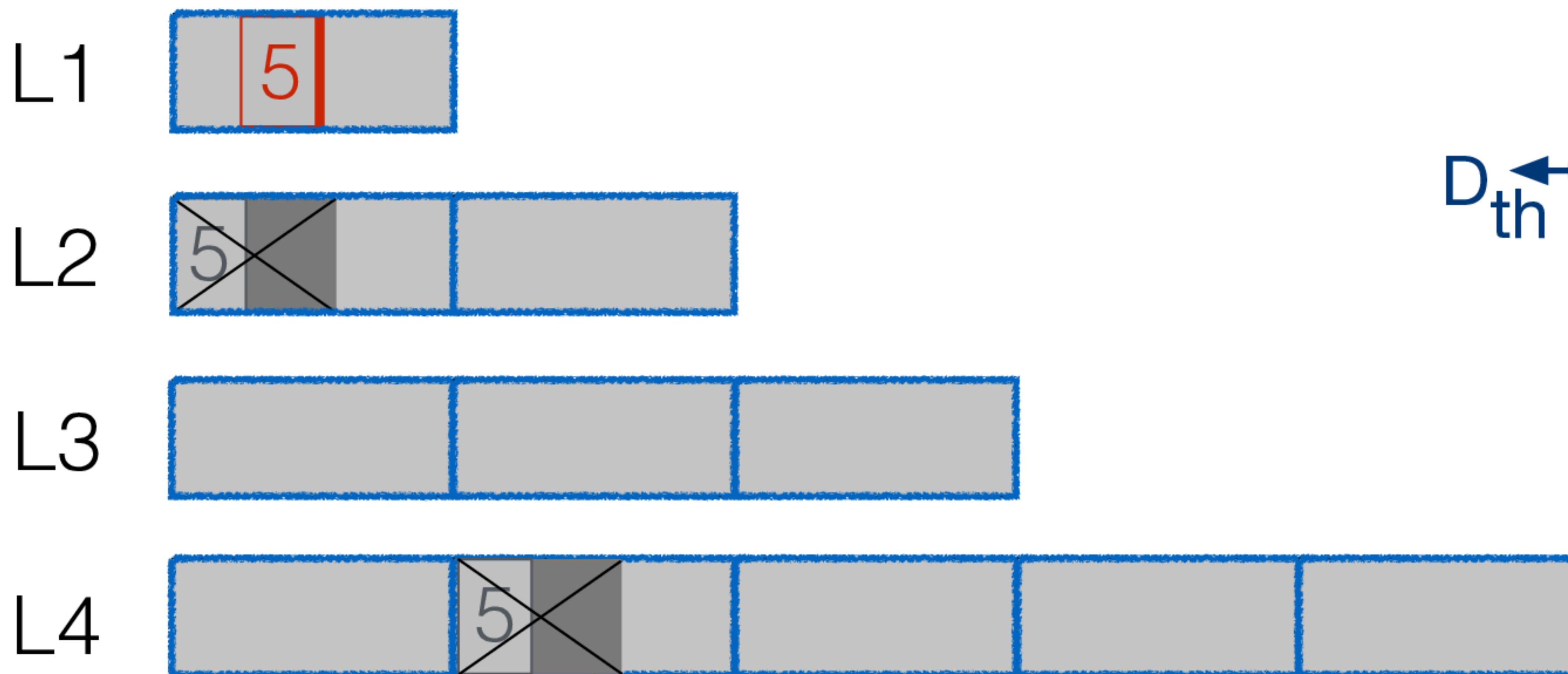
ingestion rate

file picking policy



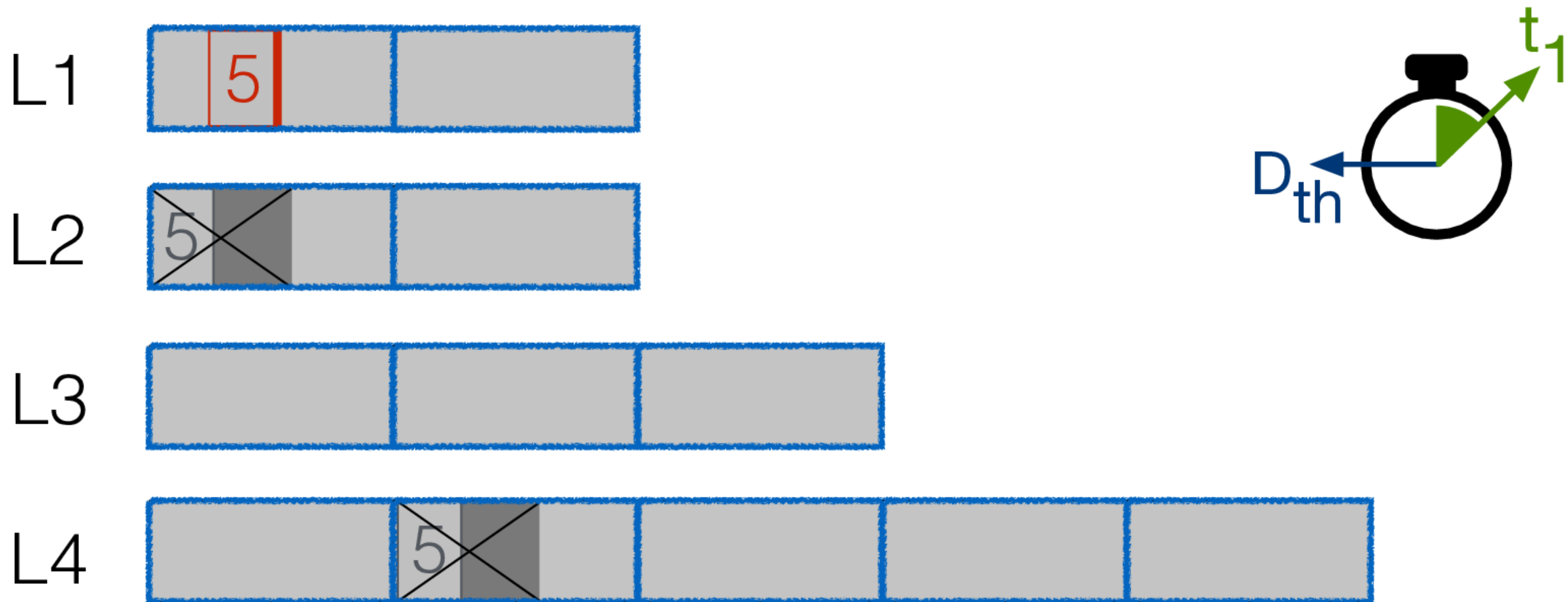
delete persistence latency

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



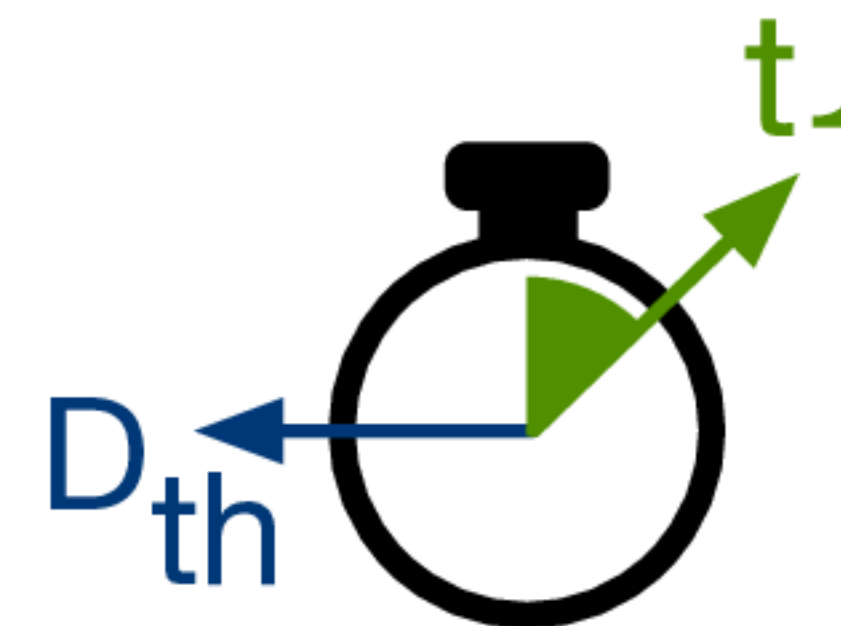
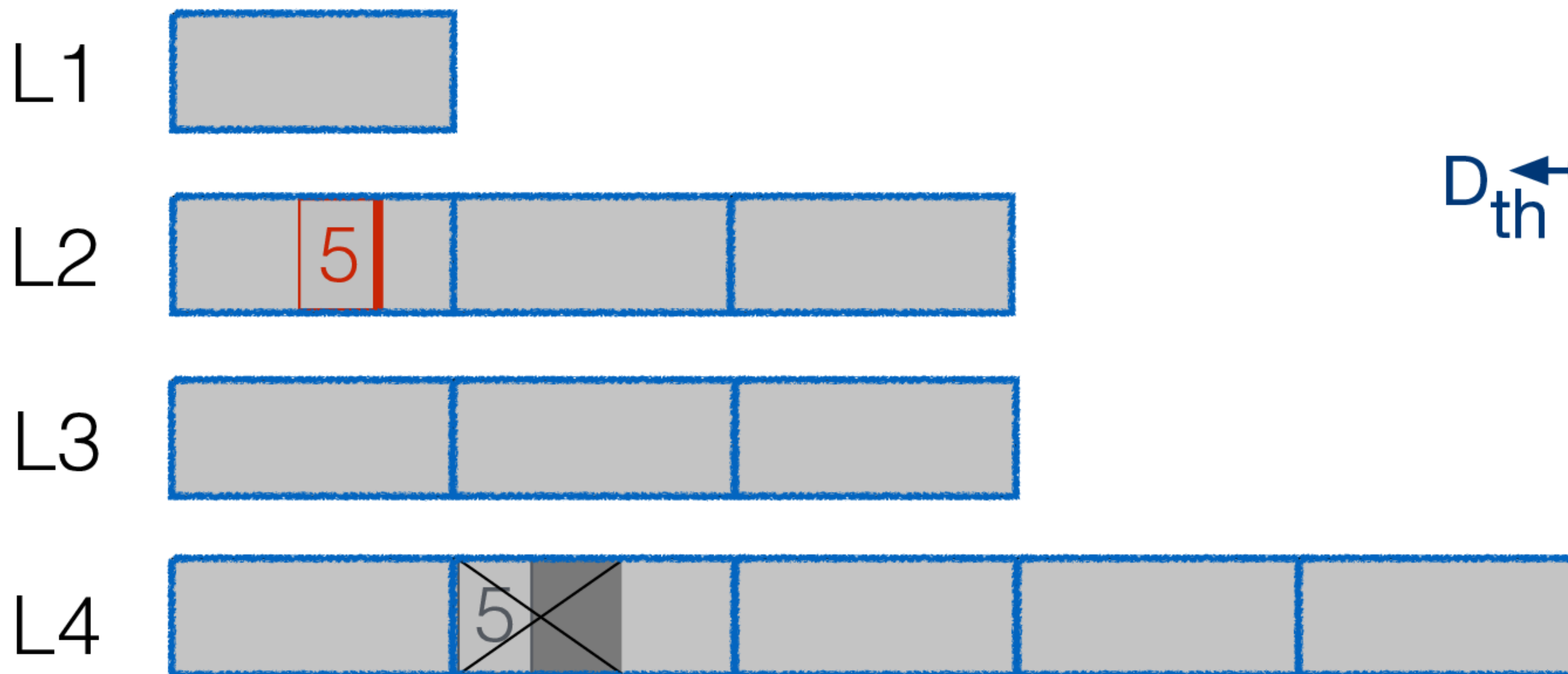
delete persistence latency

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



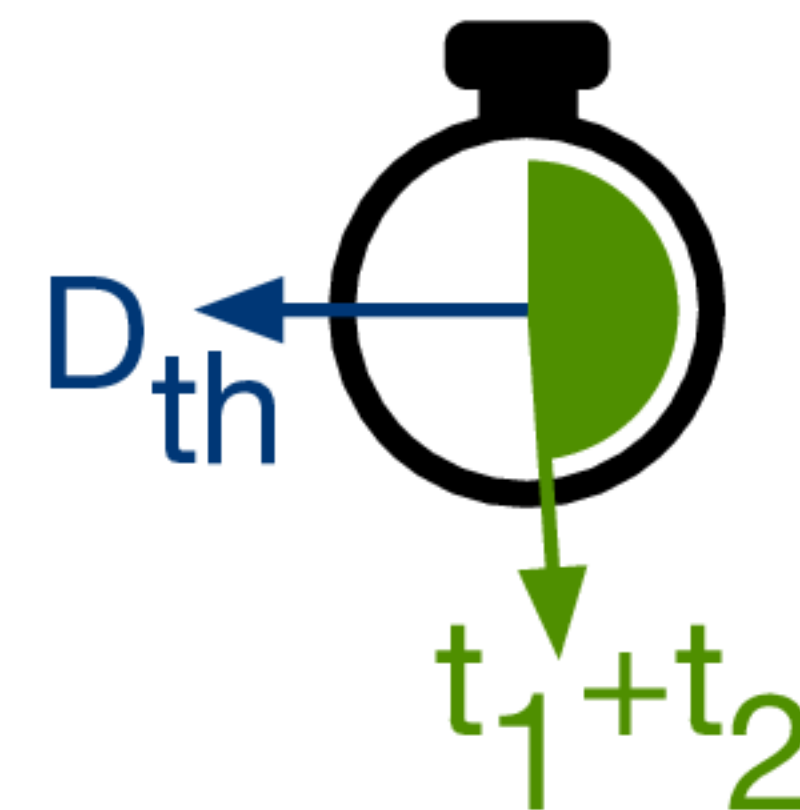
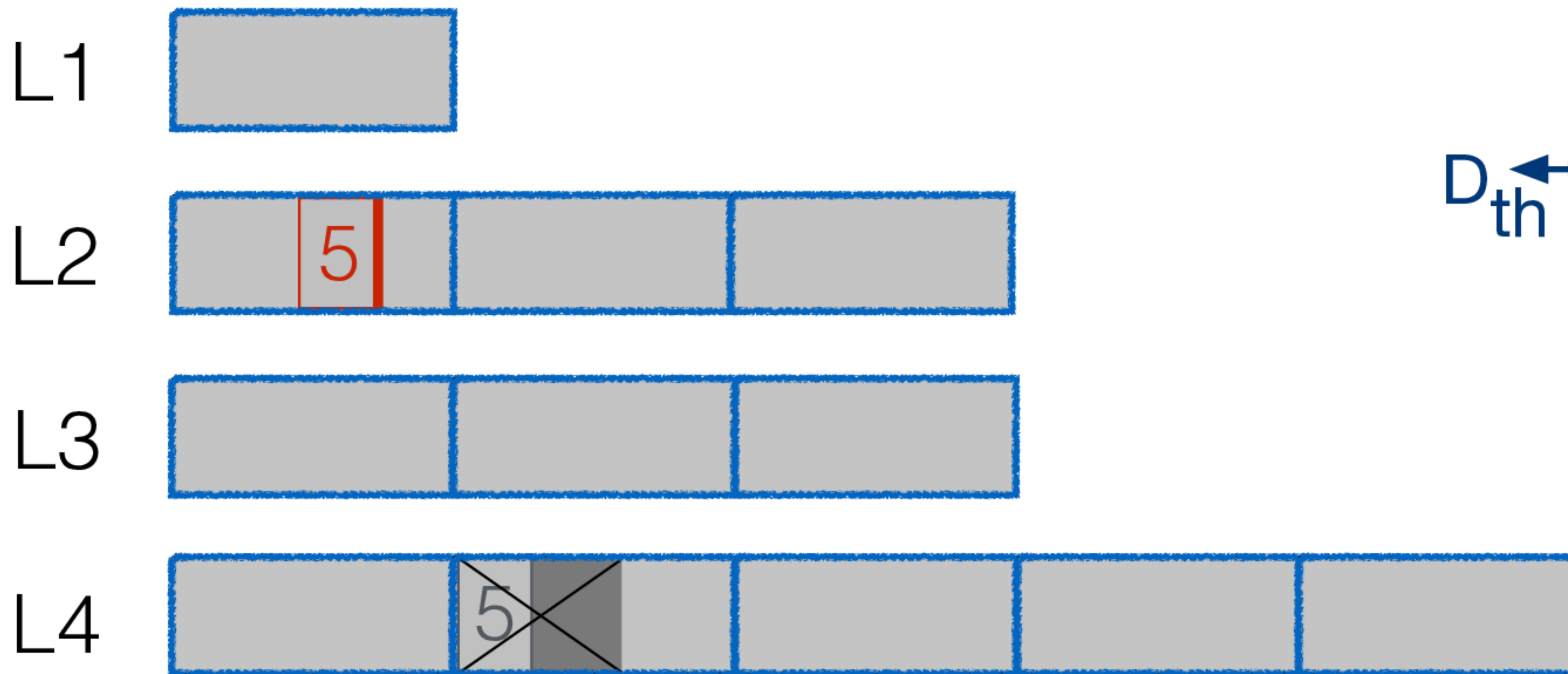
delete persistence latency

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



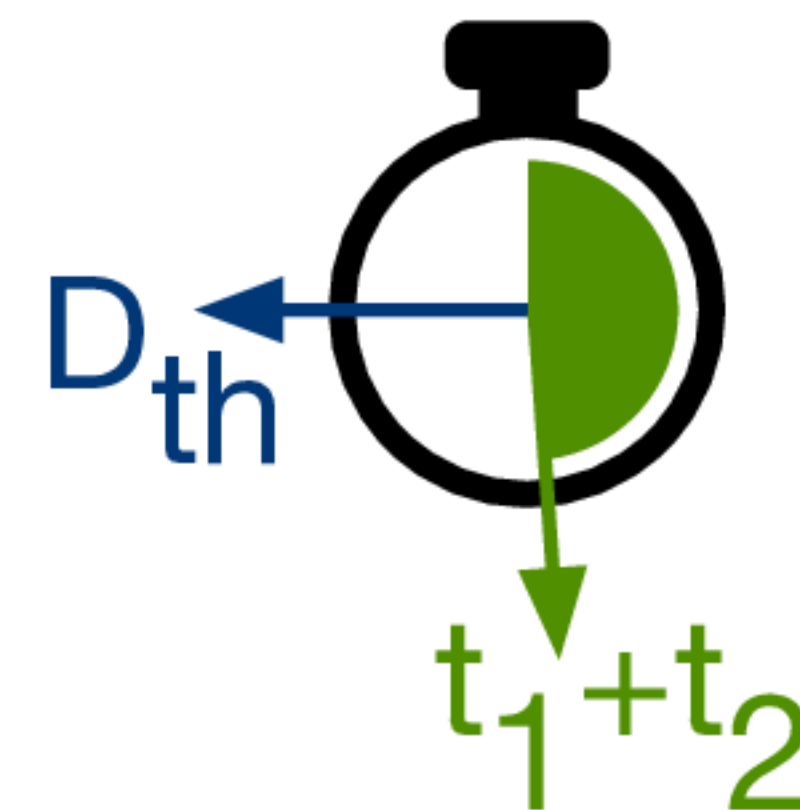
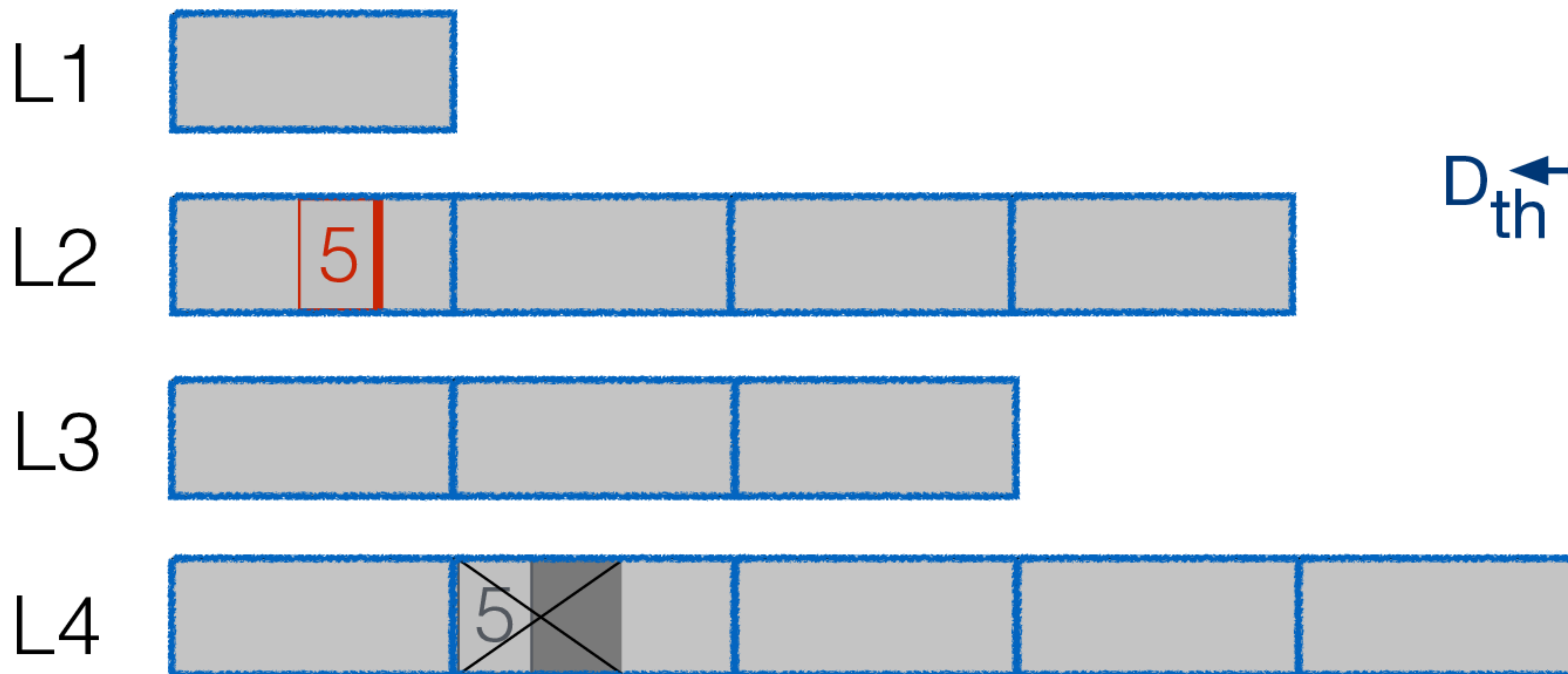
delete persistence latency

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



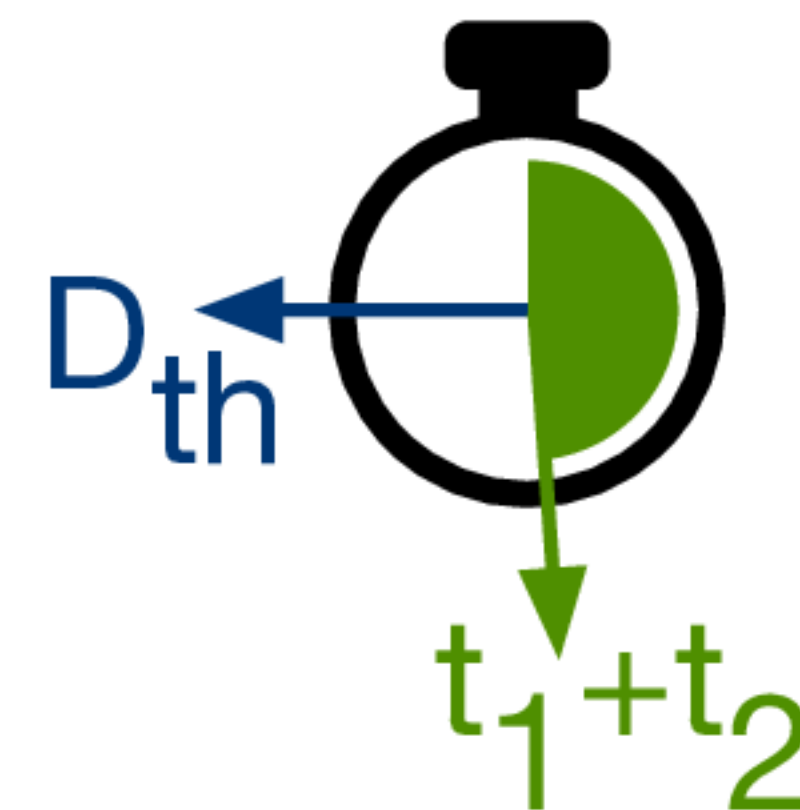
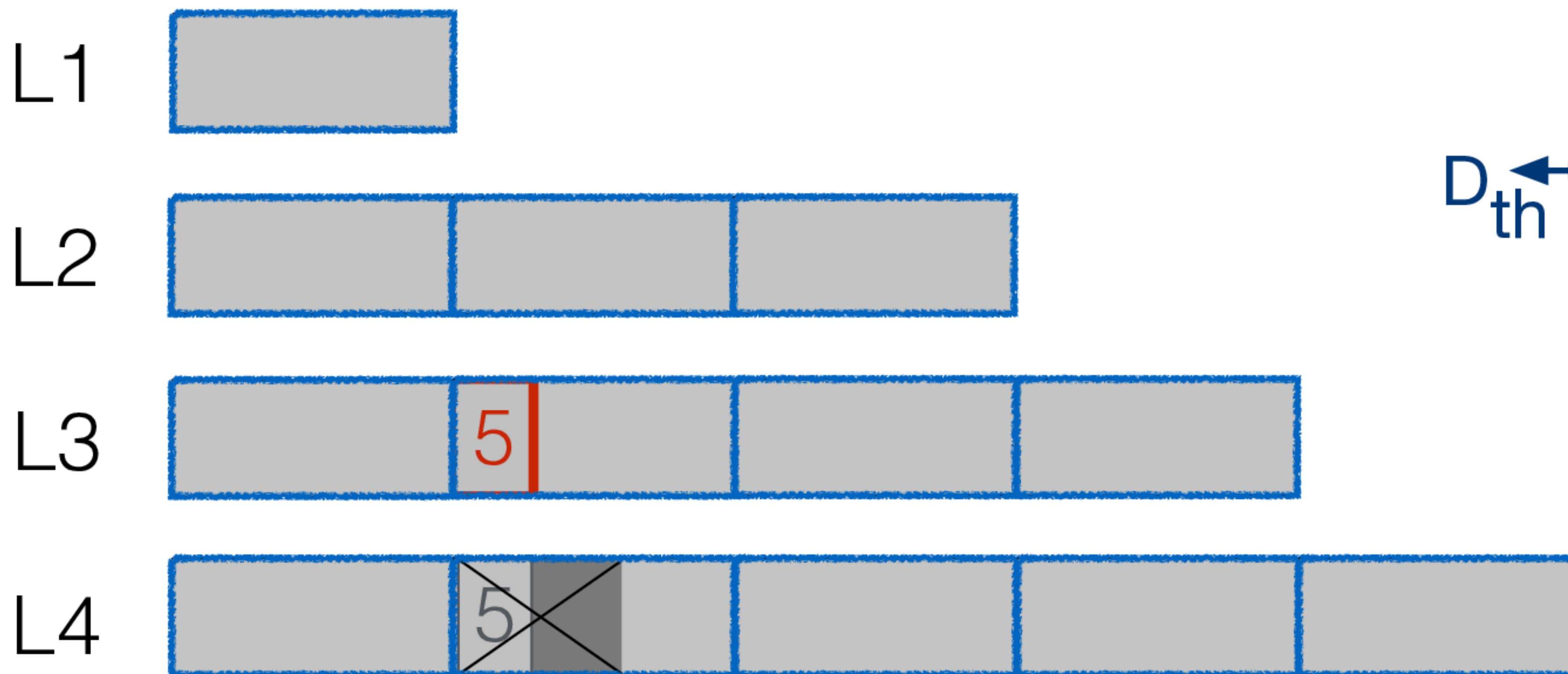
delete persistence latency

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



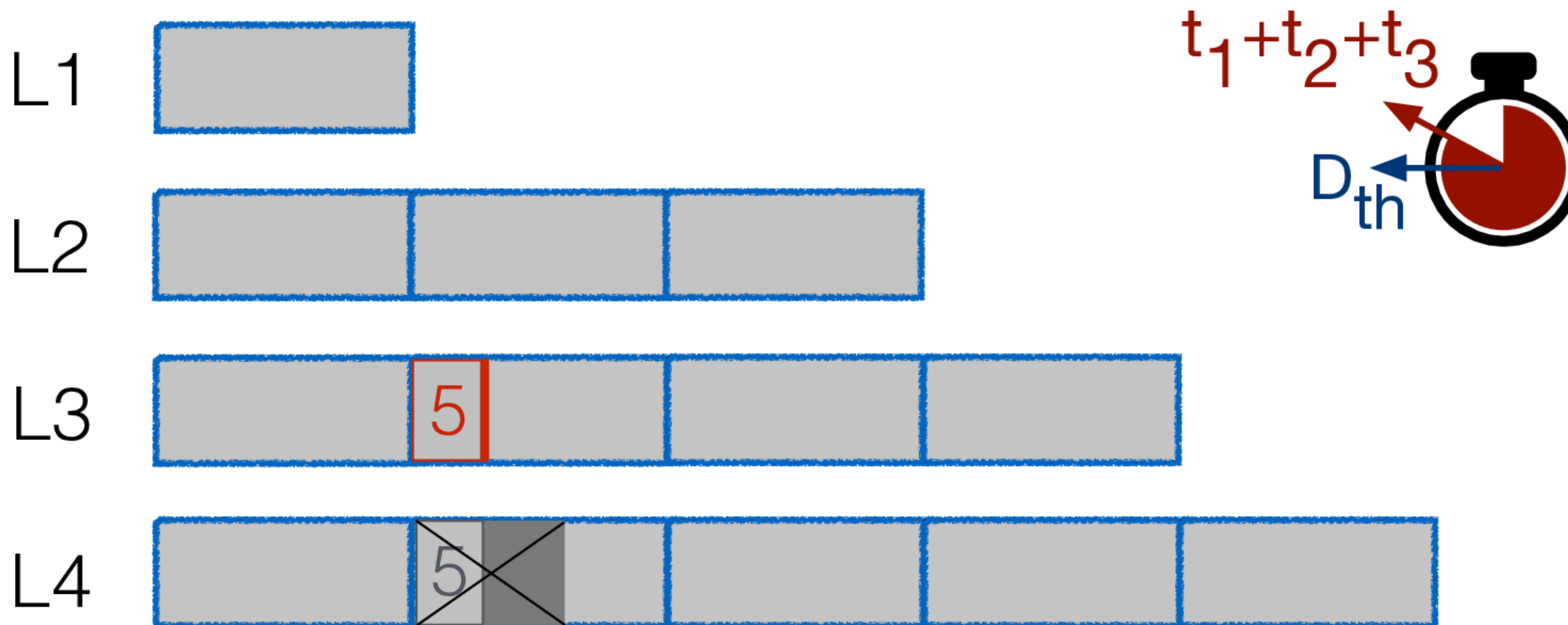
delete persistence latency

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



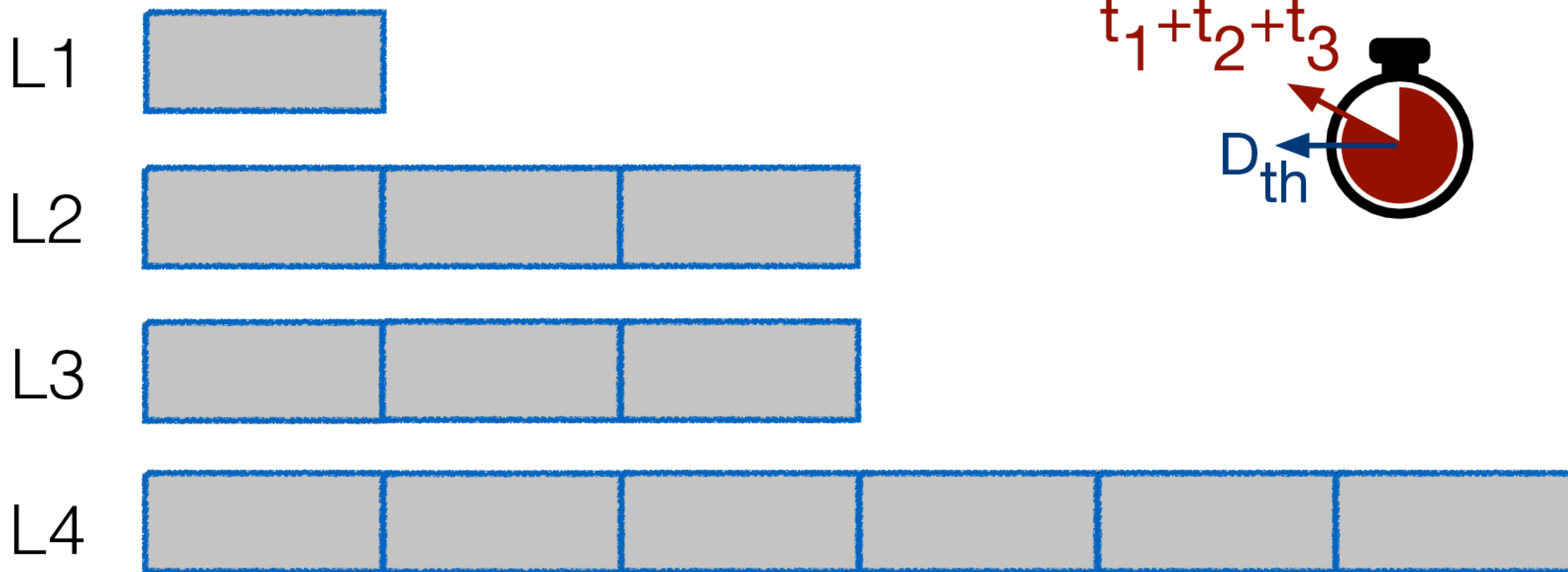
delete persistence latency

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



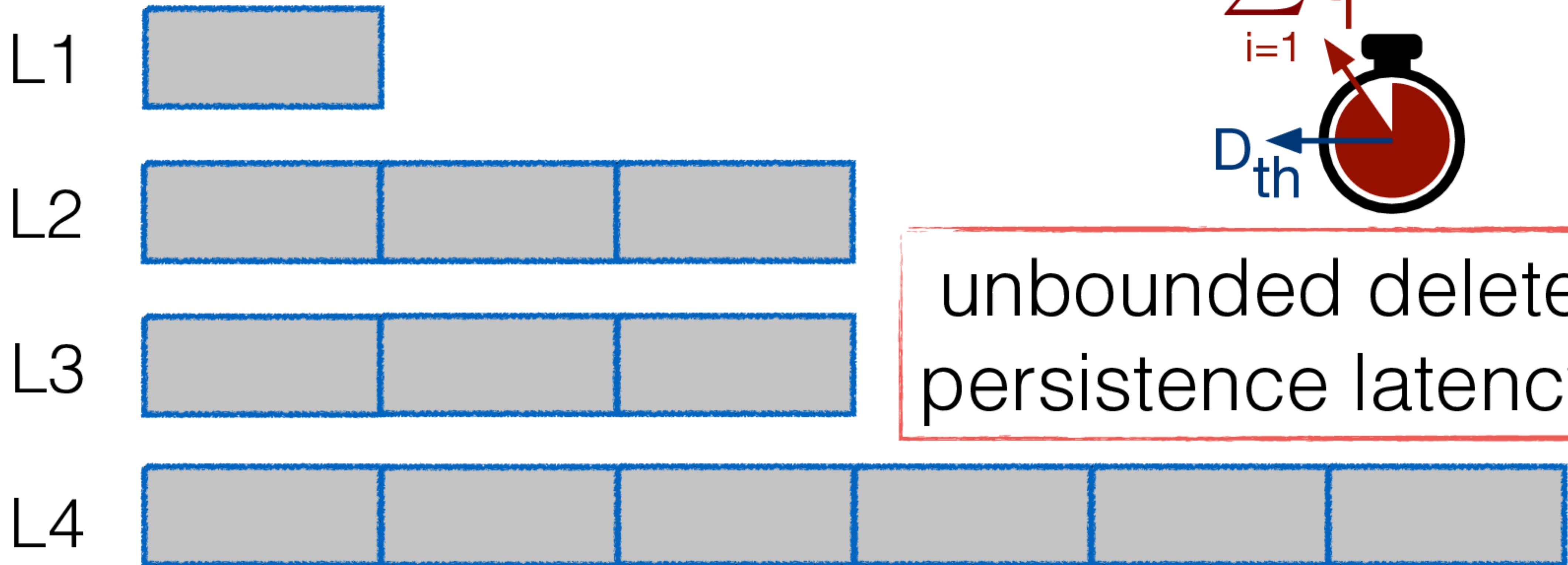
delete persistence latency

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

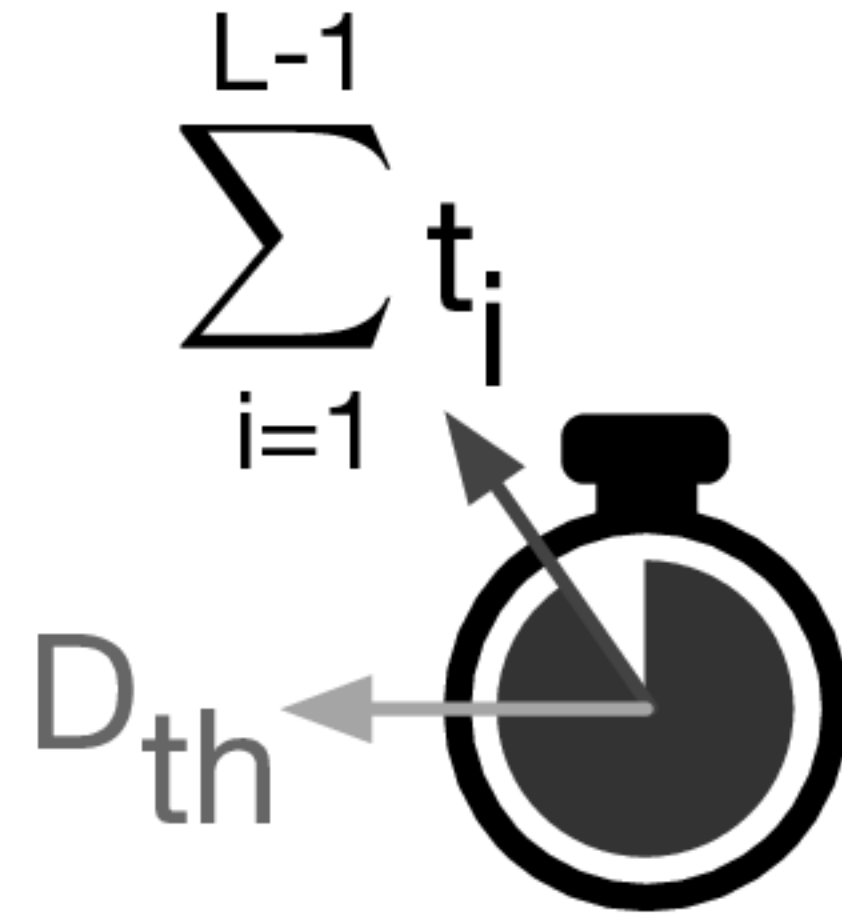


the problems

poor read perf.

write amplification

space amplification



unbounded delete
persistence latency

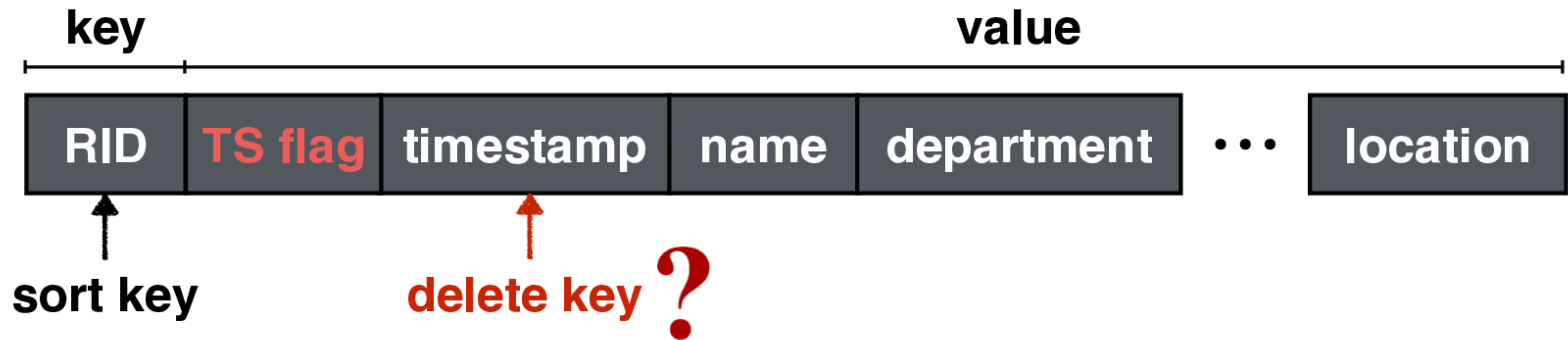




deletes on a secondary attribute

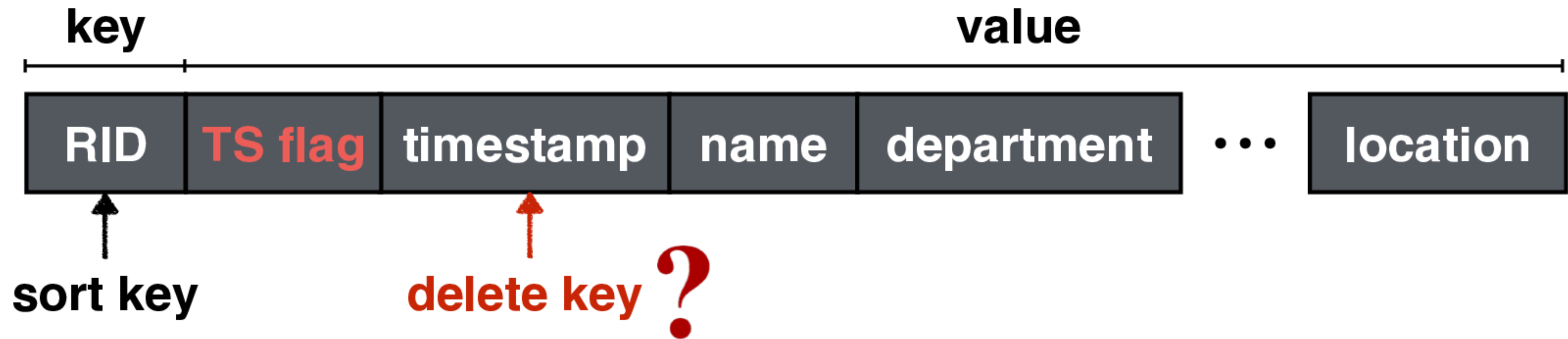
deletes on a secondary attribute

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



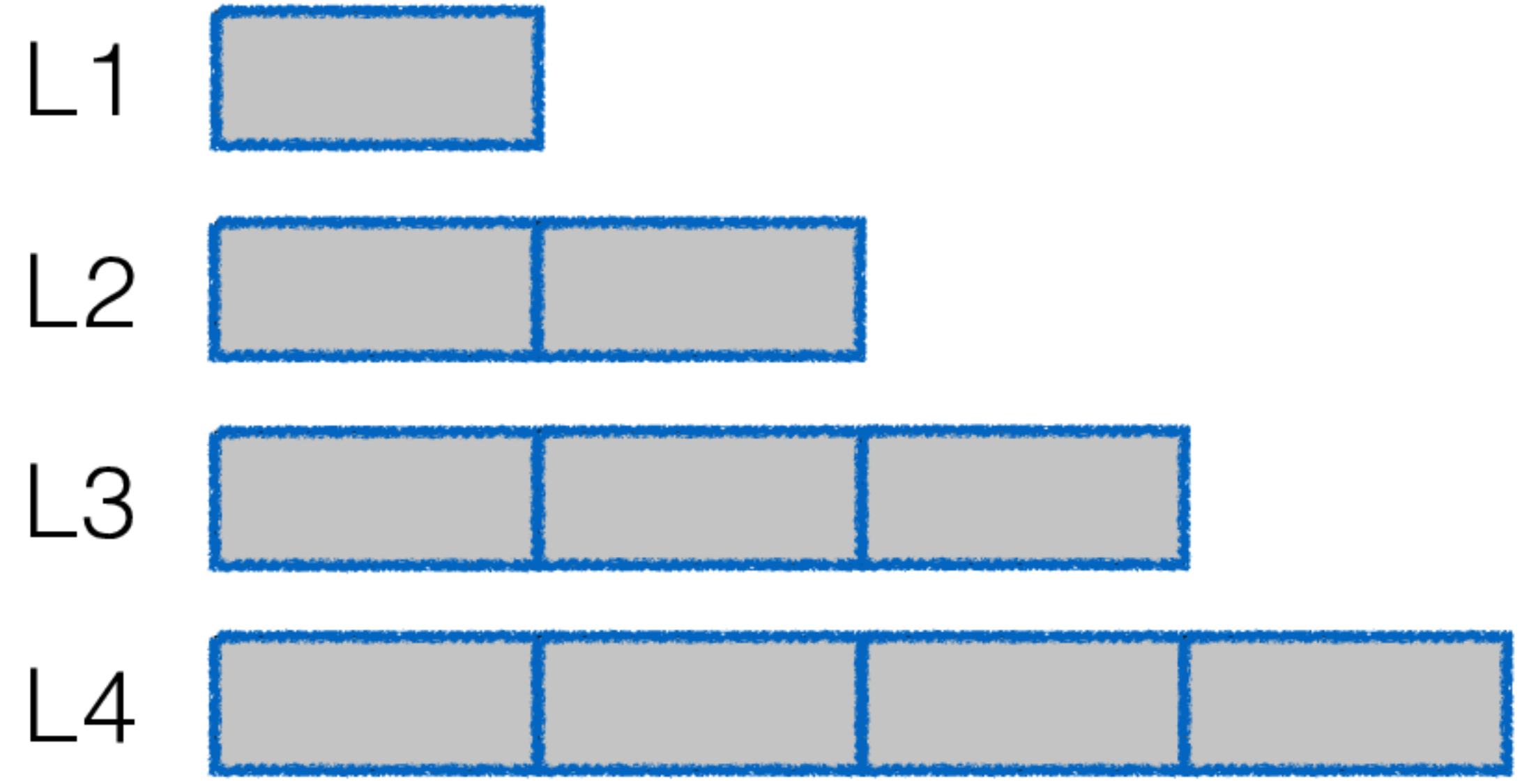
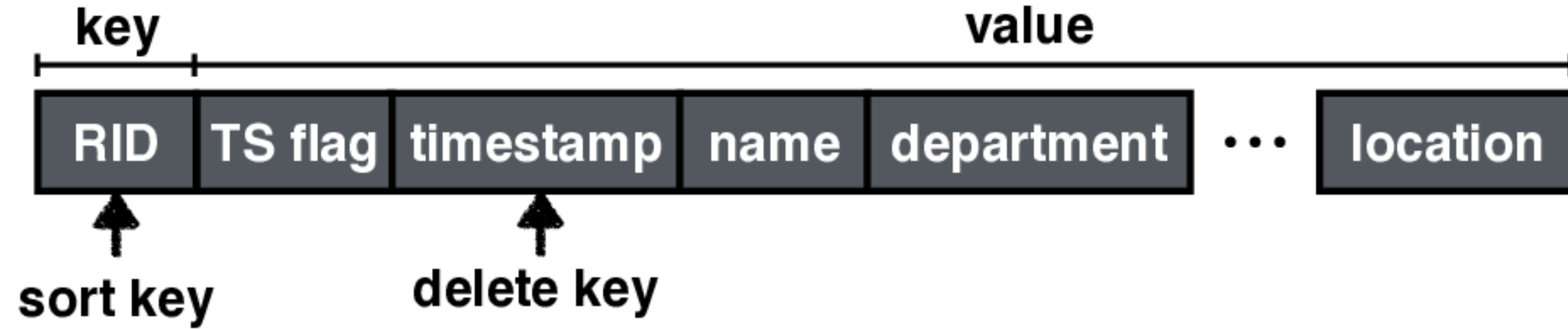
deletes on a secondary attribute

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



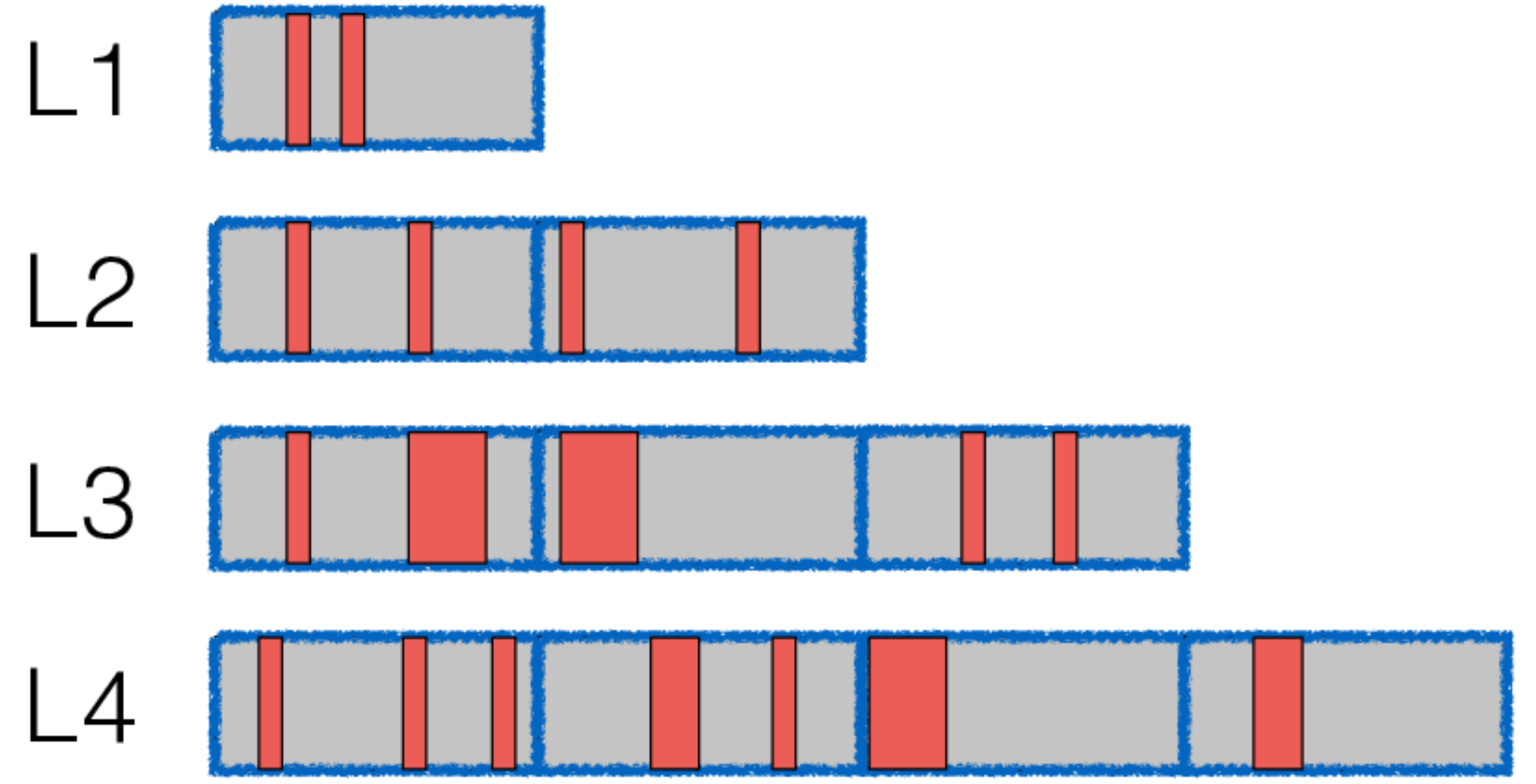
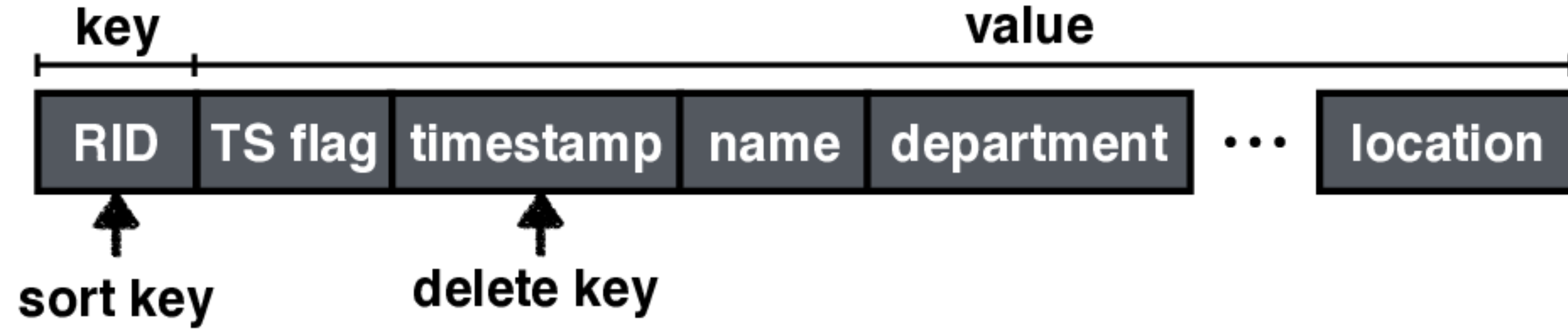
deletes on a secondary attribute

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



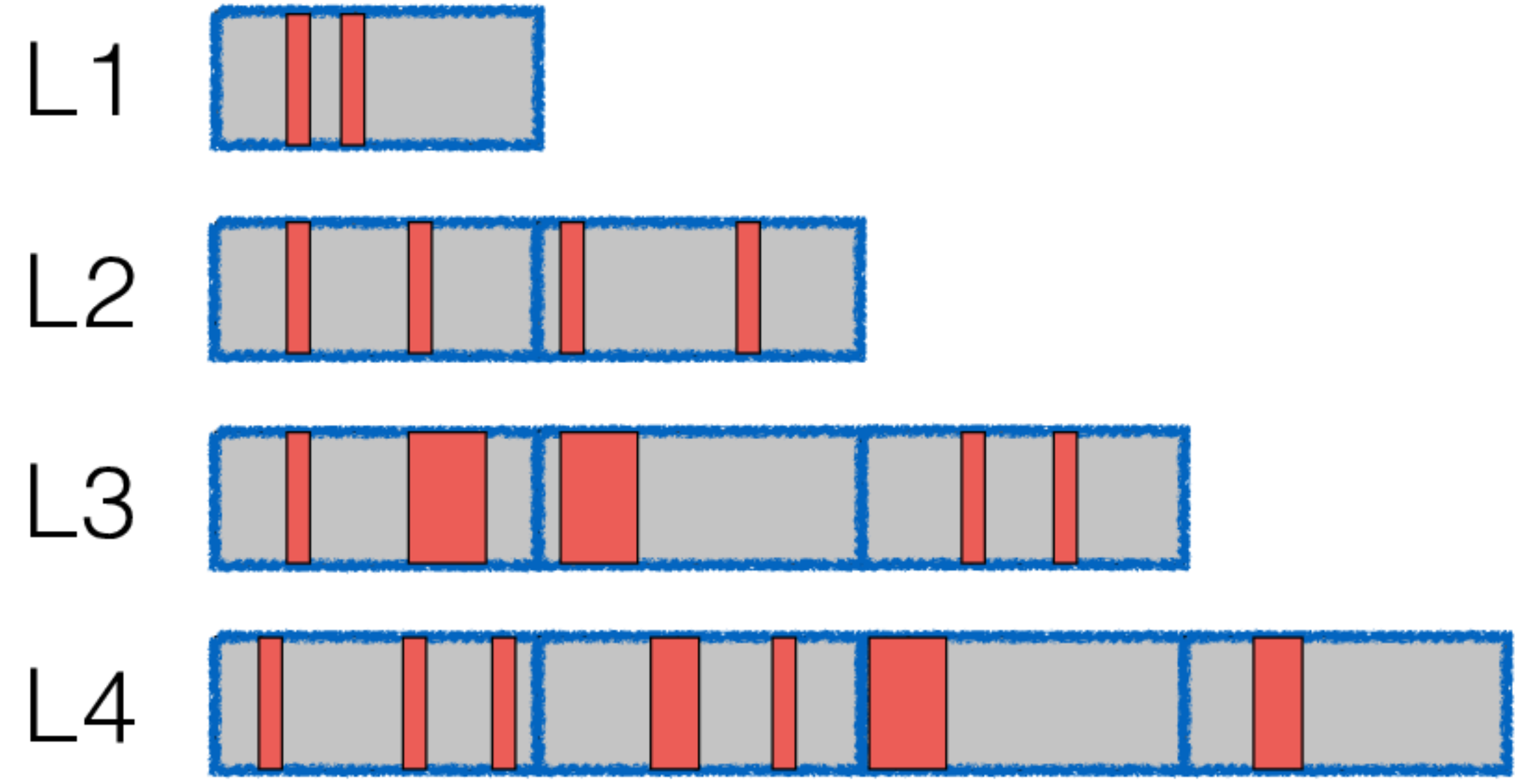
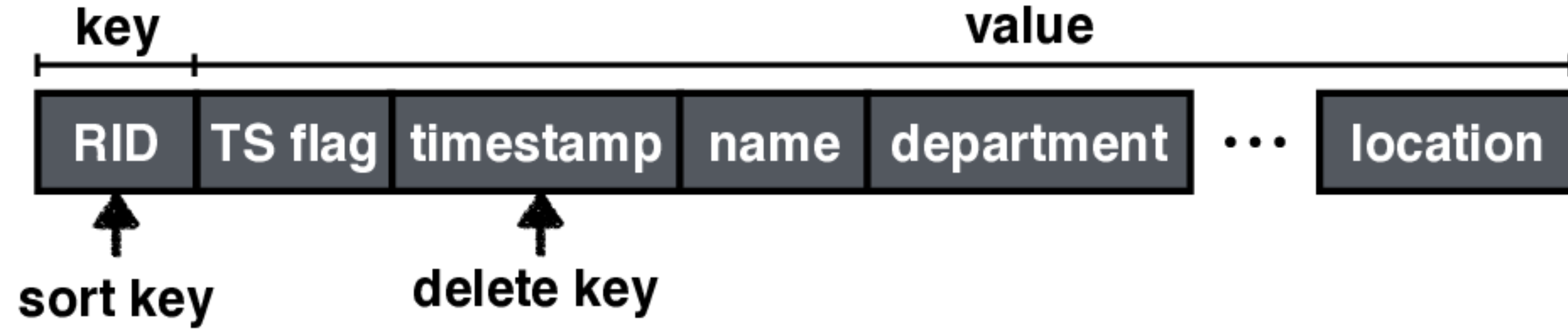
deletes on a secondary attribute

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



deletes on a secondary attribute

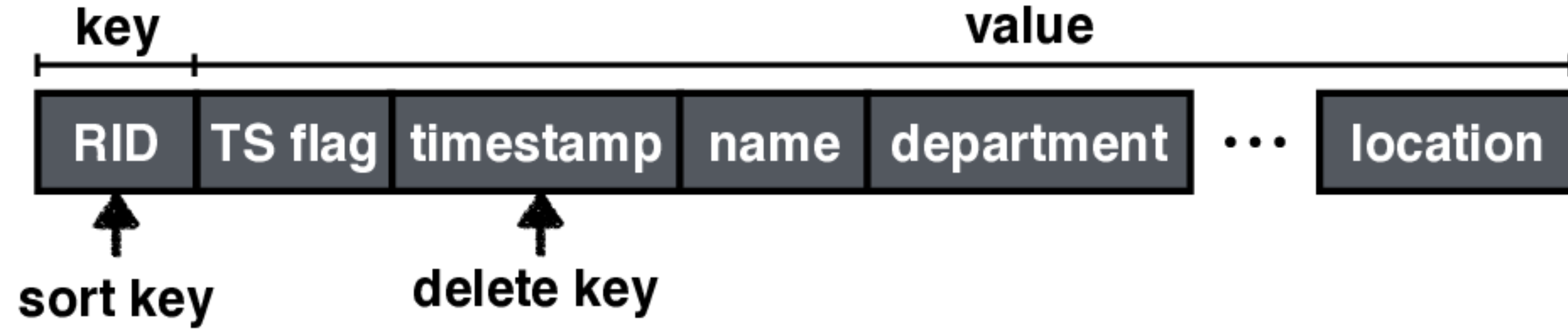
delete all entries older than: **D days**



scattered occurrences

deletes on a secondary attribute

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



L1

L2

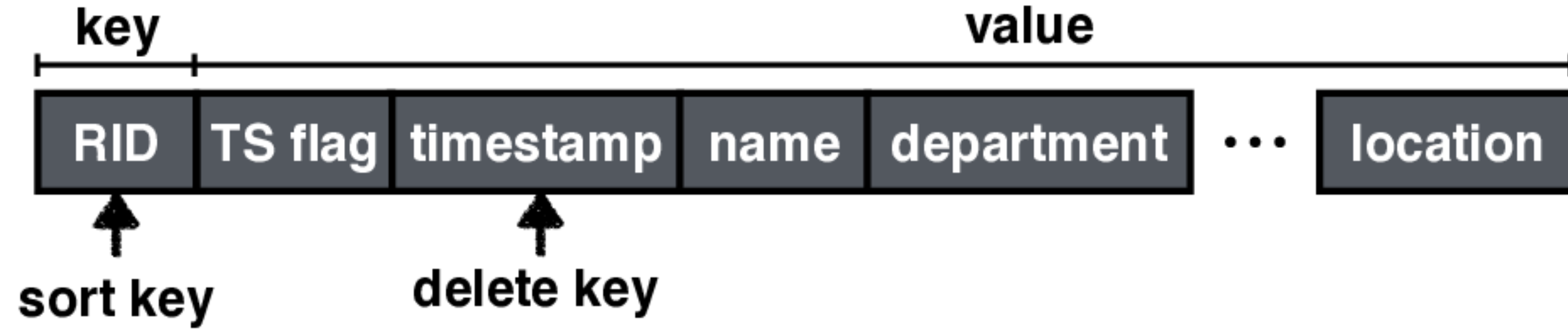
L3

L4



deletes on a secondary attribute

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



L1

L2

L3

L4

latency spikes

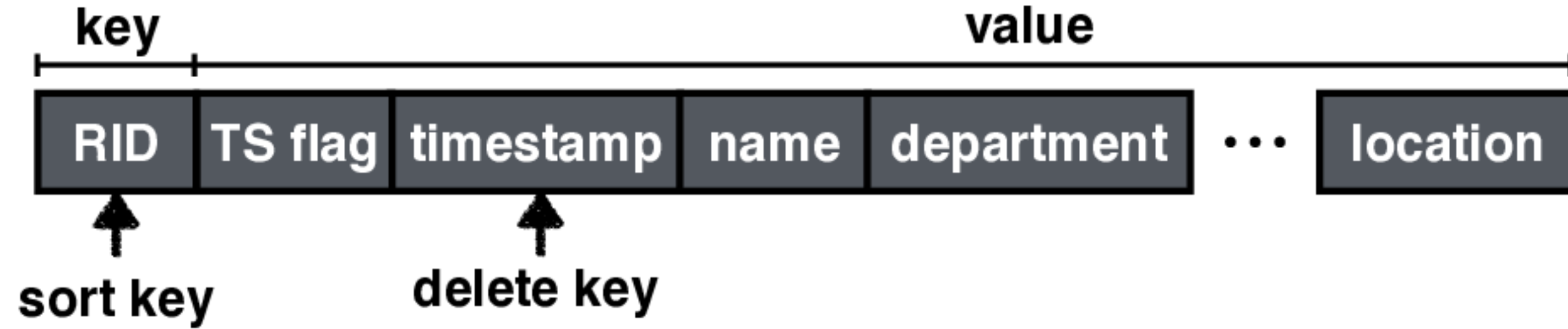


superfluous I/Os



deletes on a secondary attribute

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



L1

L2

L3

L4

latency spikes



superfluous I/Os

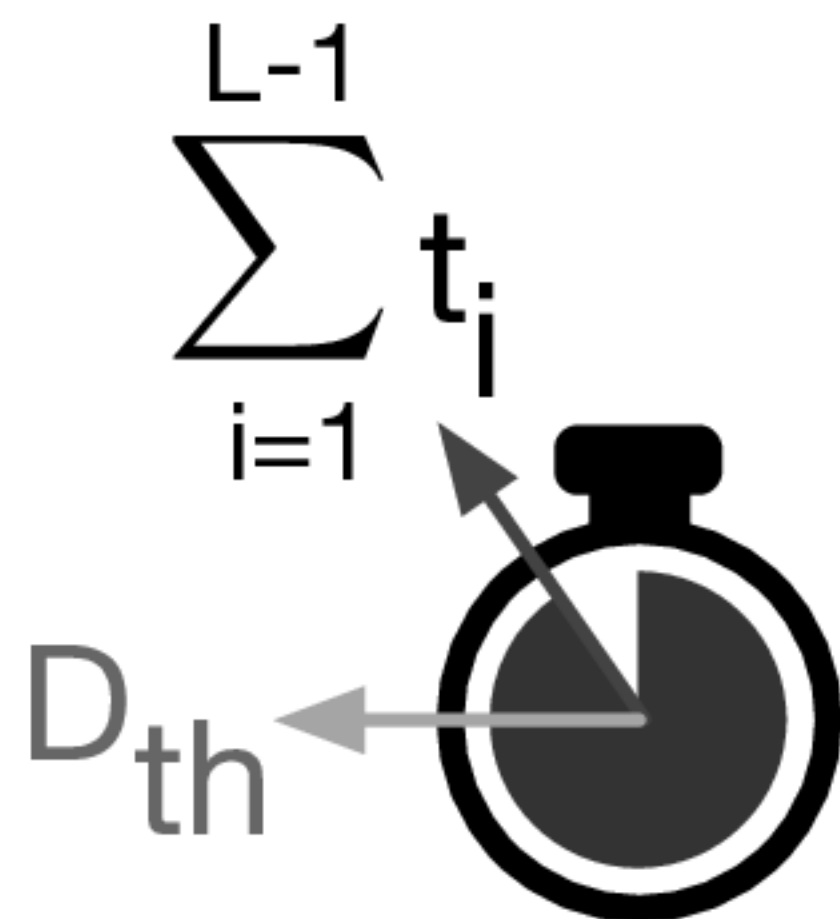


the problems

poor read perf.

write amplification

space amplification



unbounded delete
persistence latency

latency spikes

superfluous I/Os

the solution

poor read perf.

write amplification

space amplification

FADE

unbounded delete
persistence latency

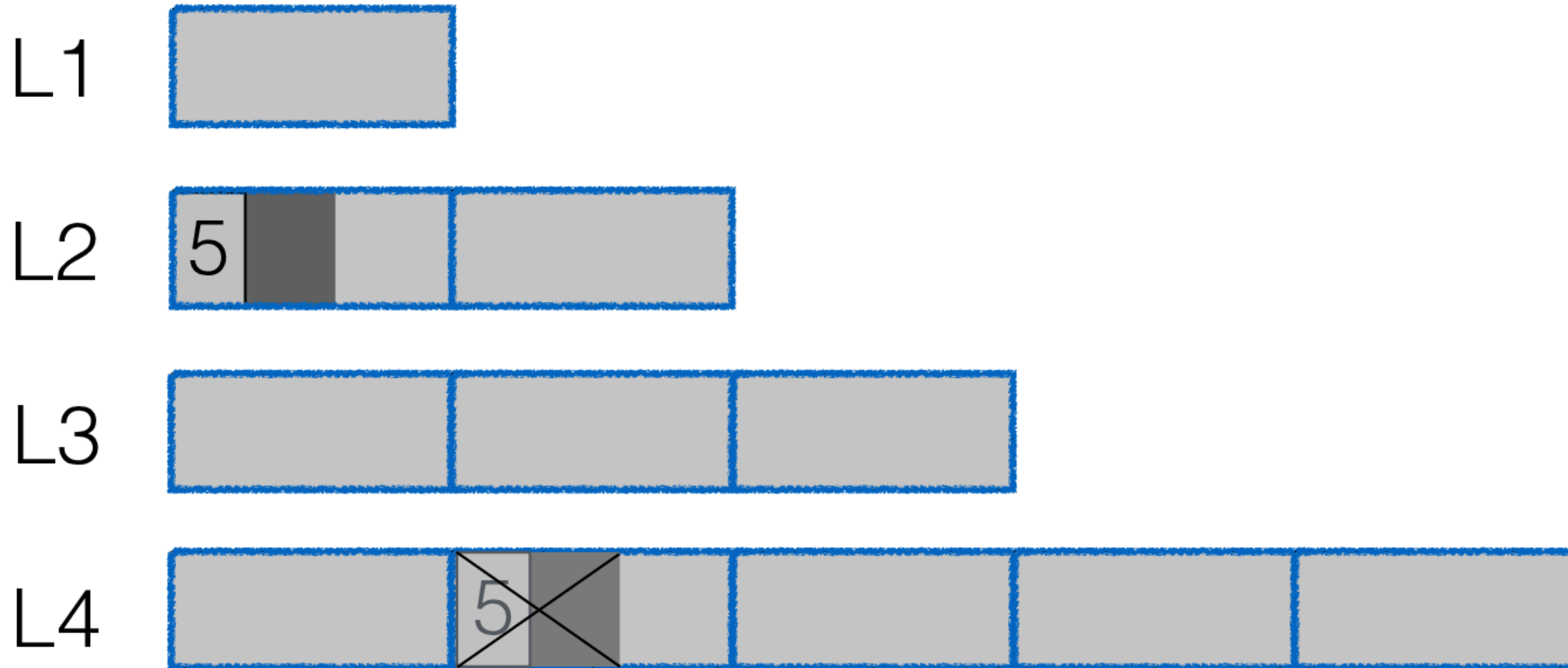


latency spikes

superfluous I/Os

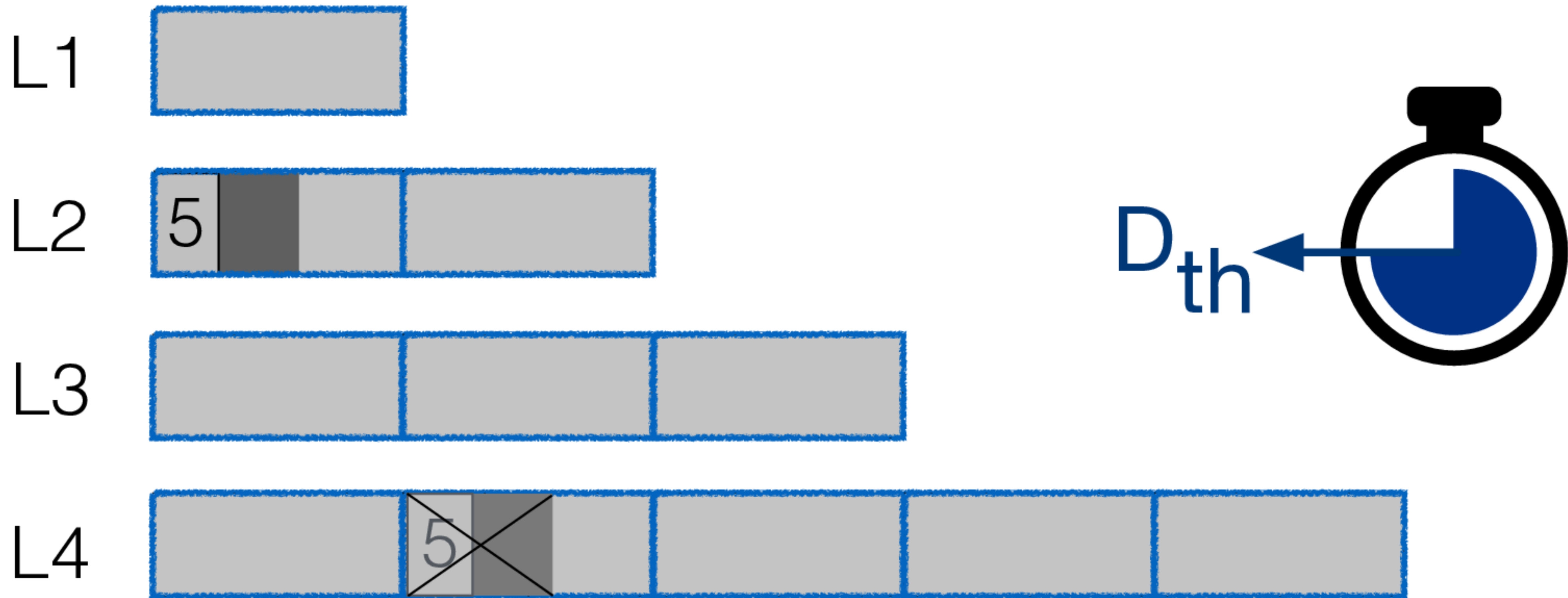
FAst DElete

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



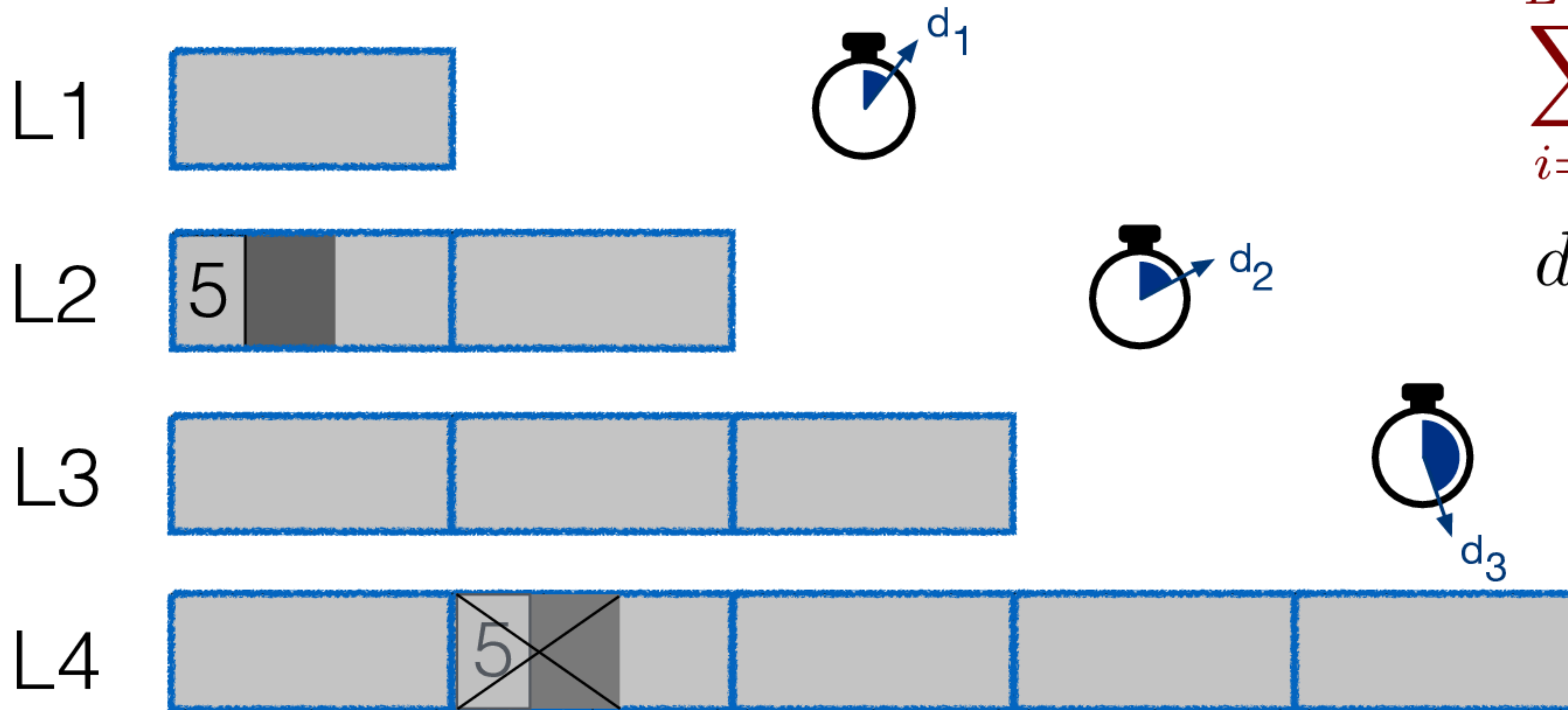
FAst DElete

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



FAst DElete

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

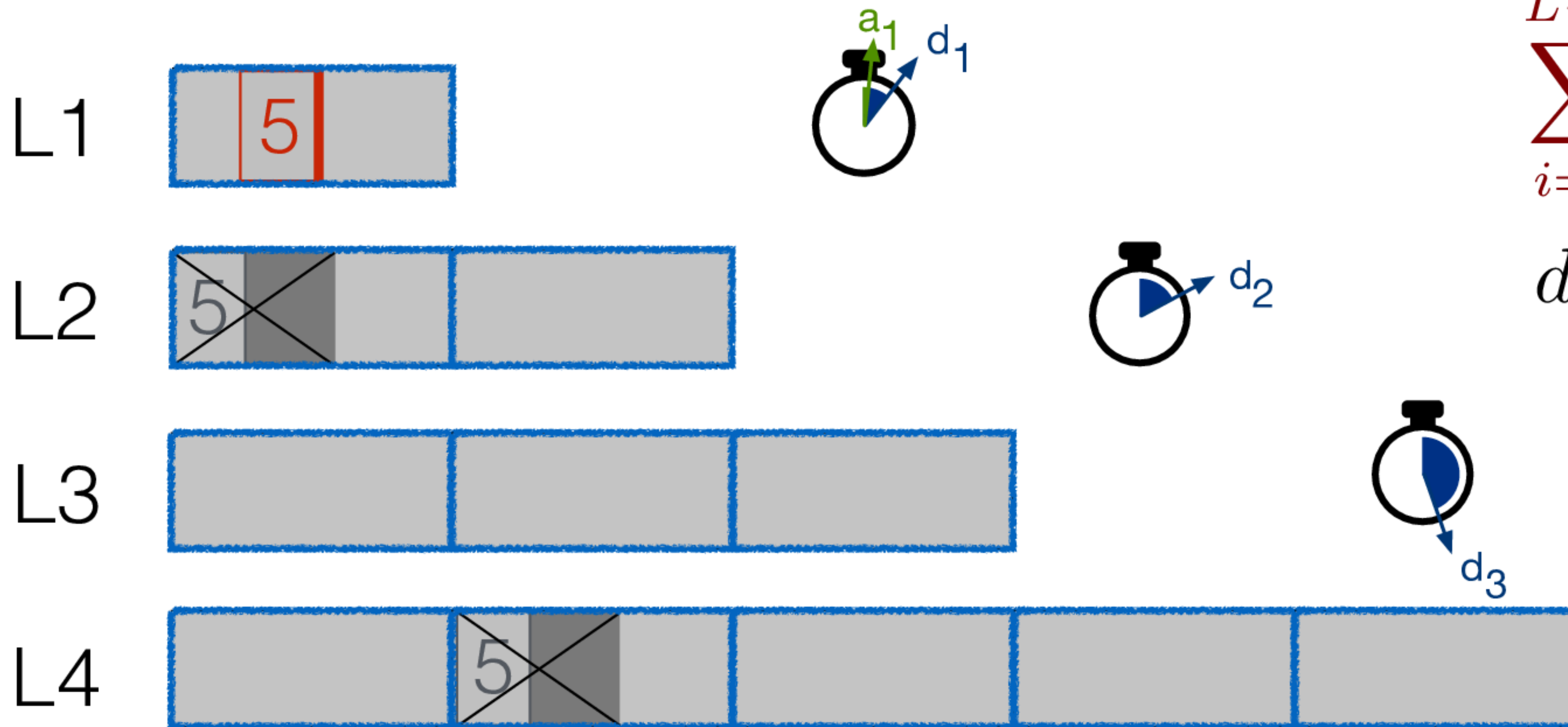


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

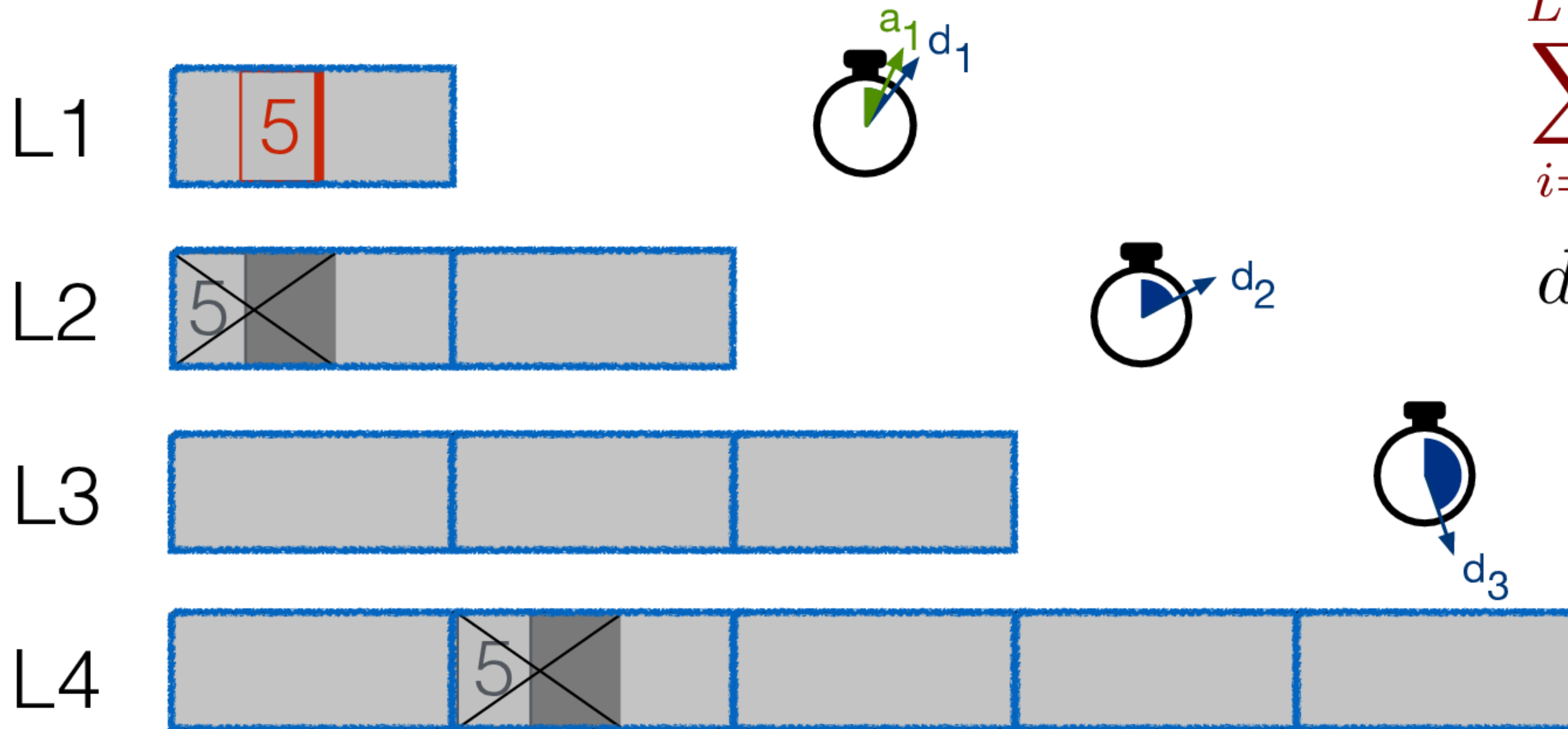


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

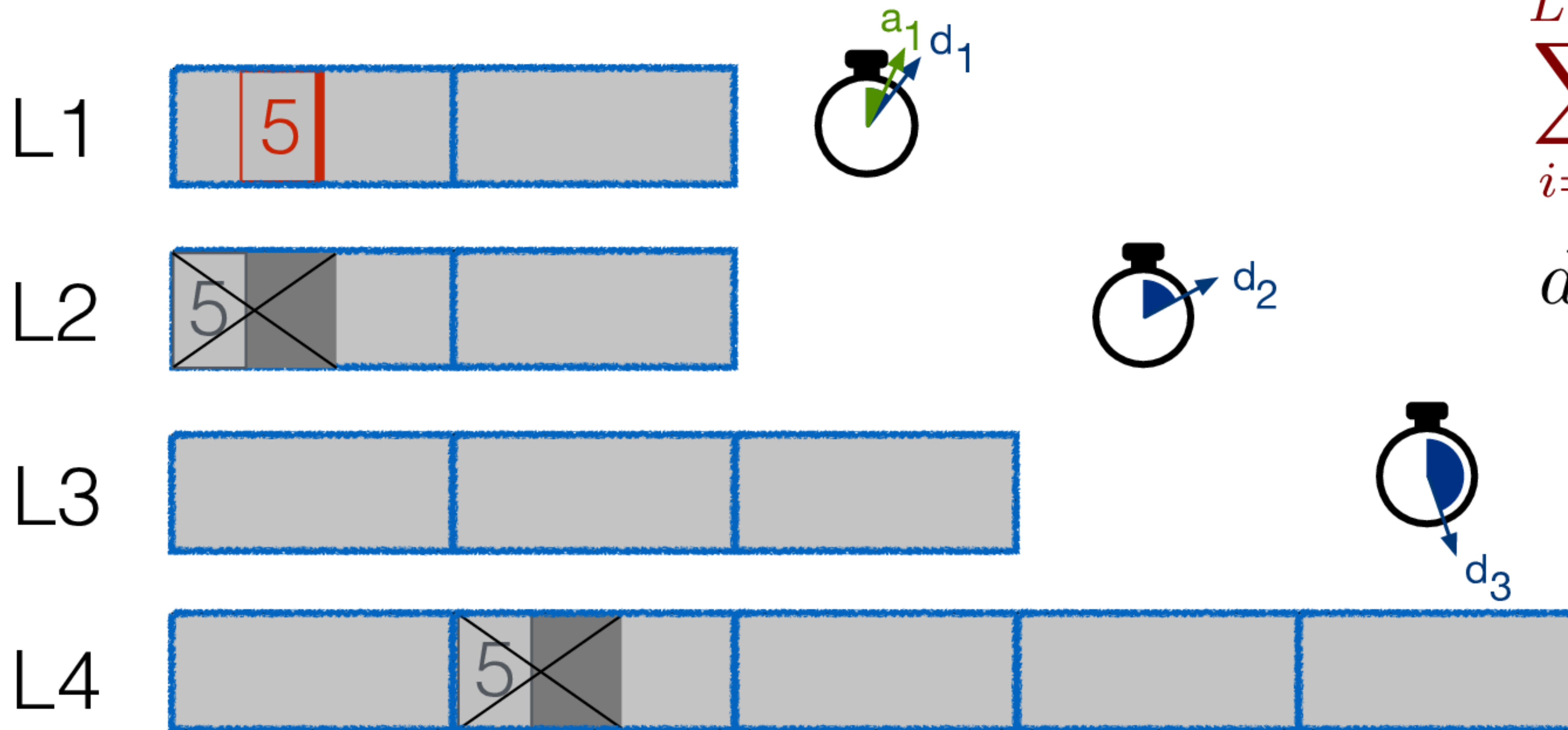


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

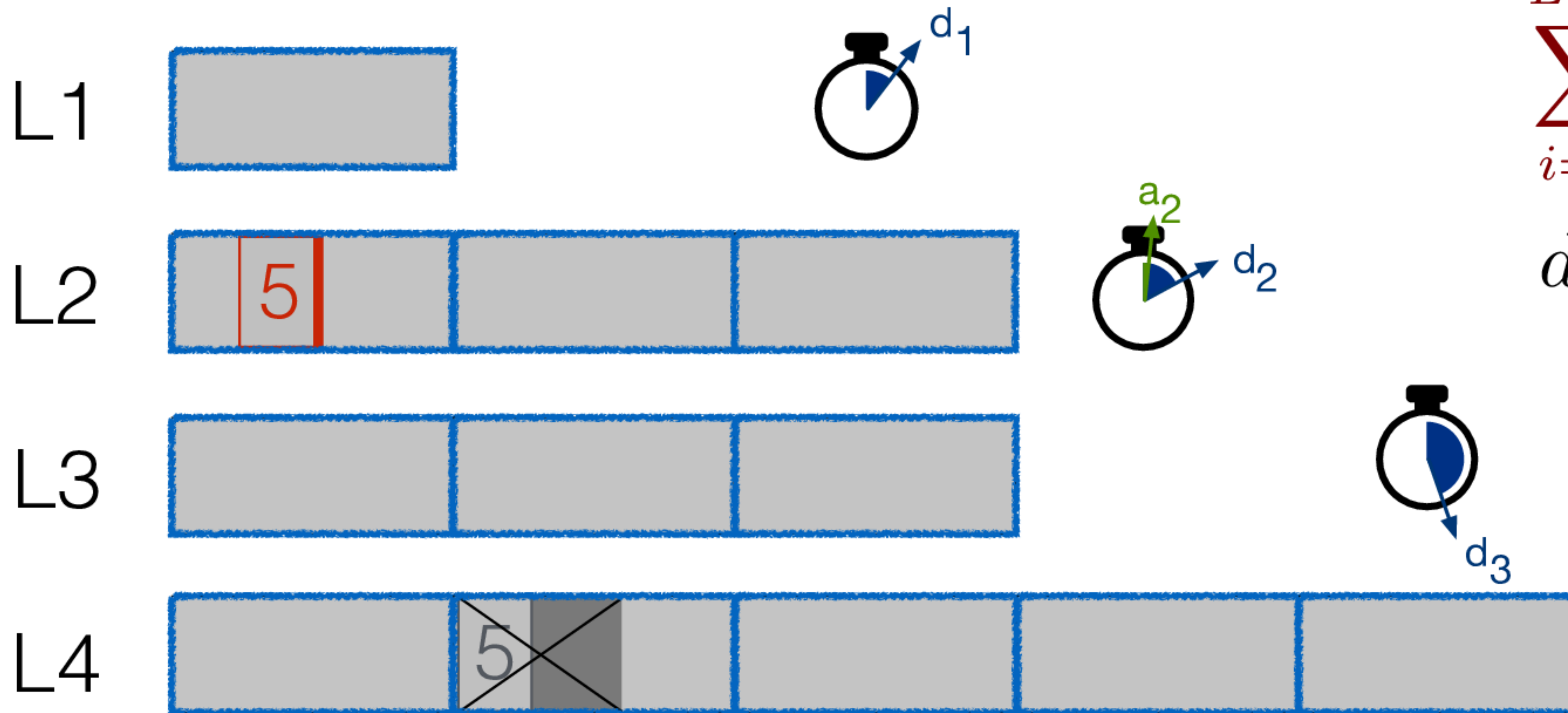


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

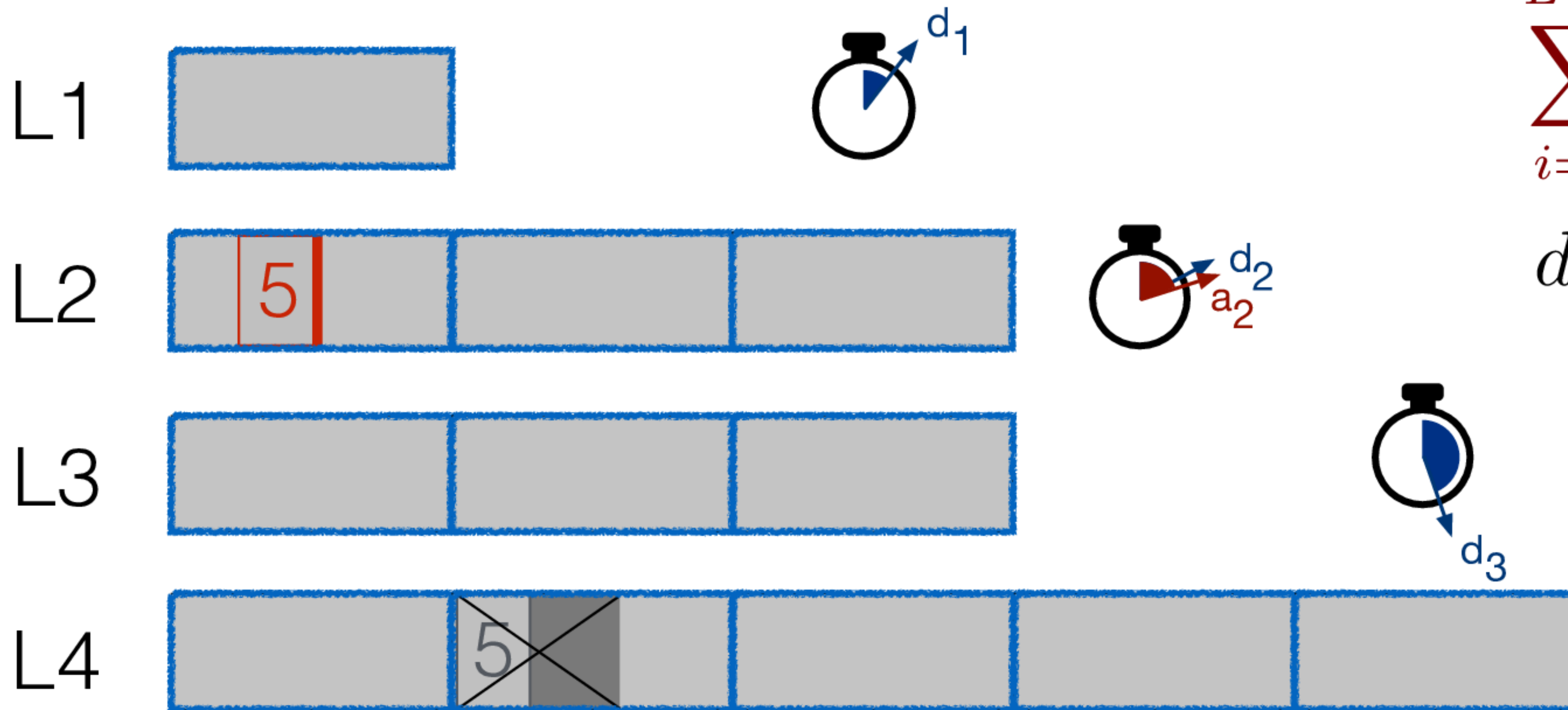


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

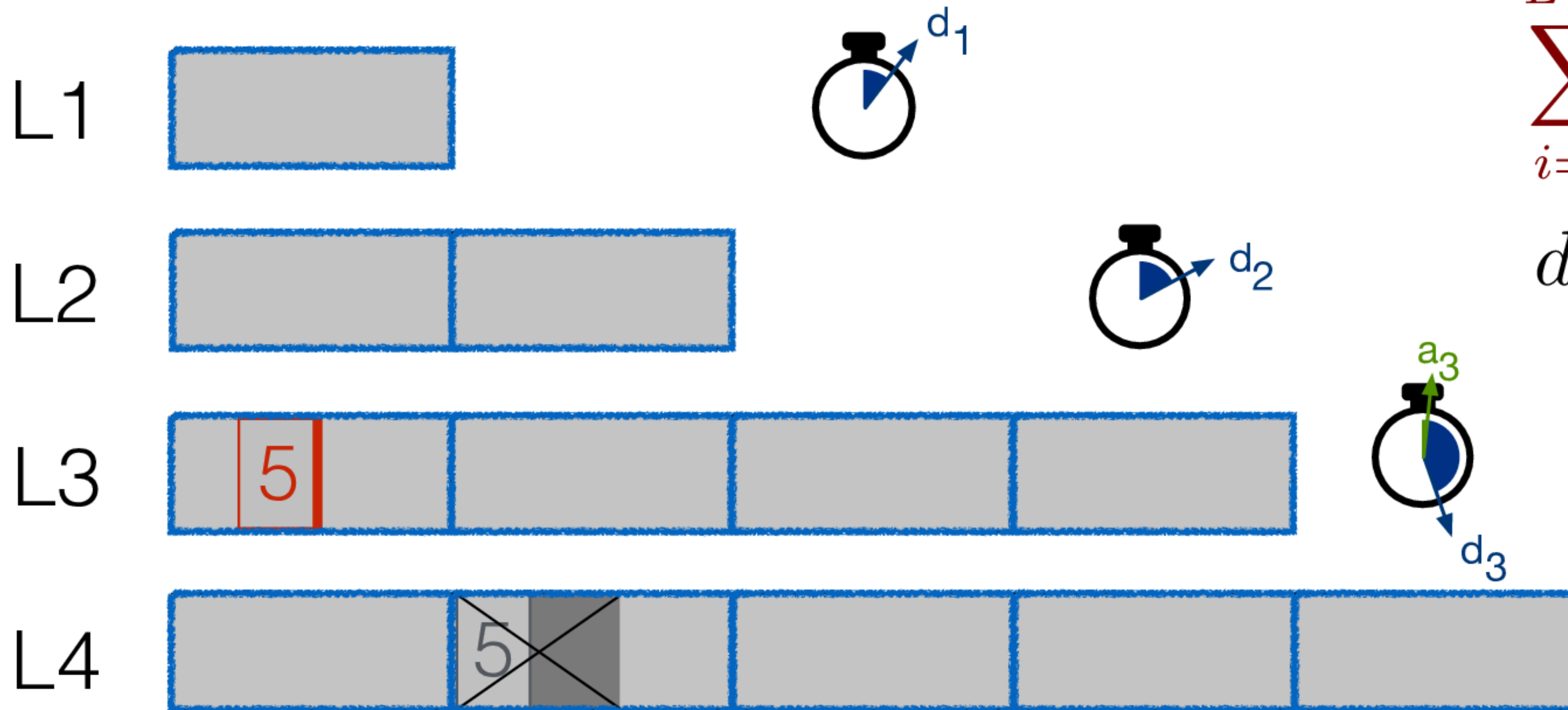


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

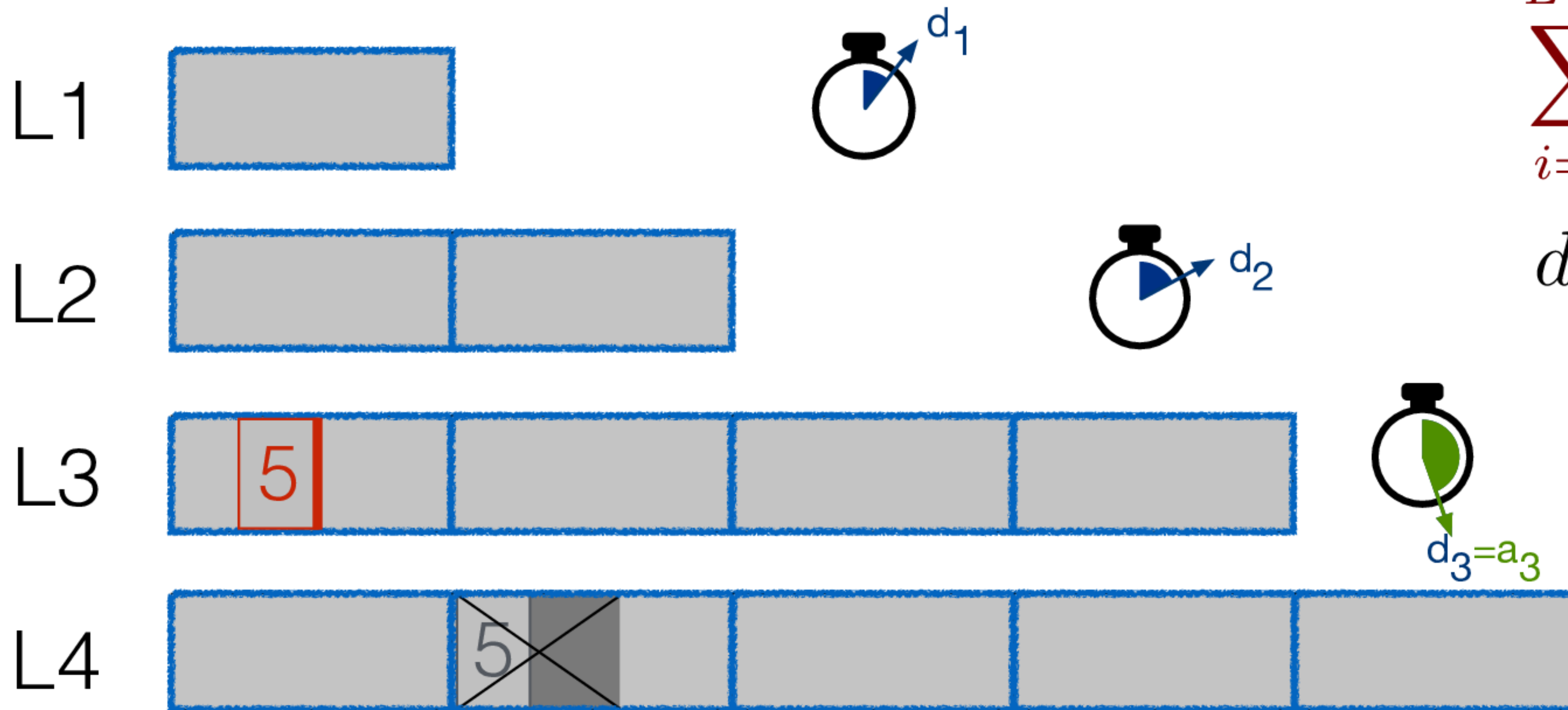


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

FAst DElete

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

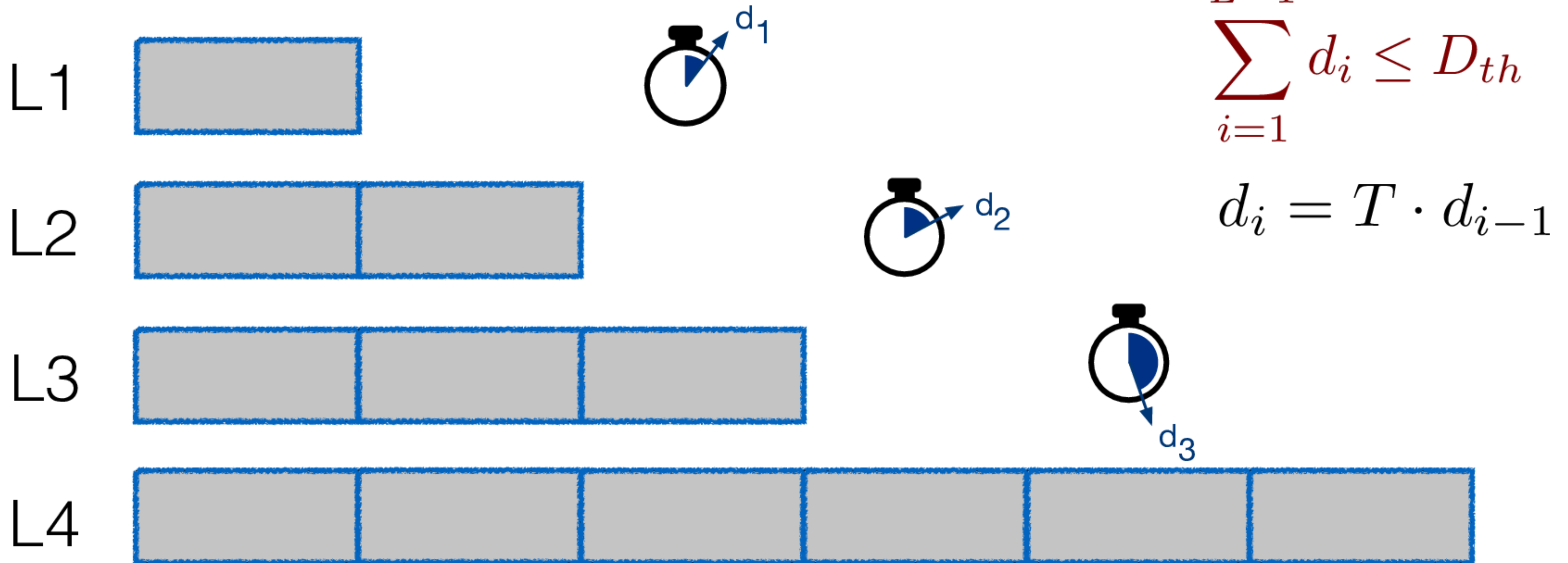


$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

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

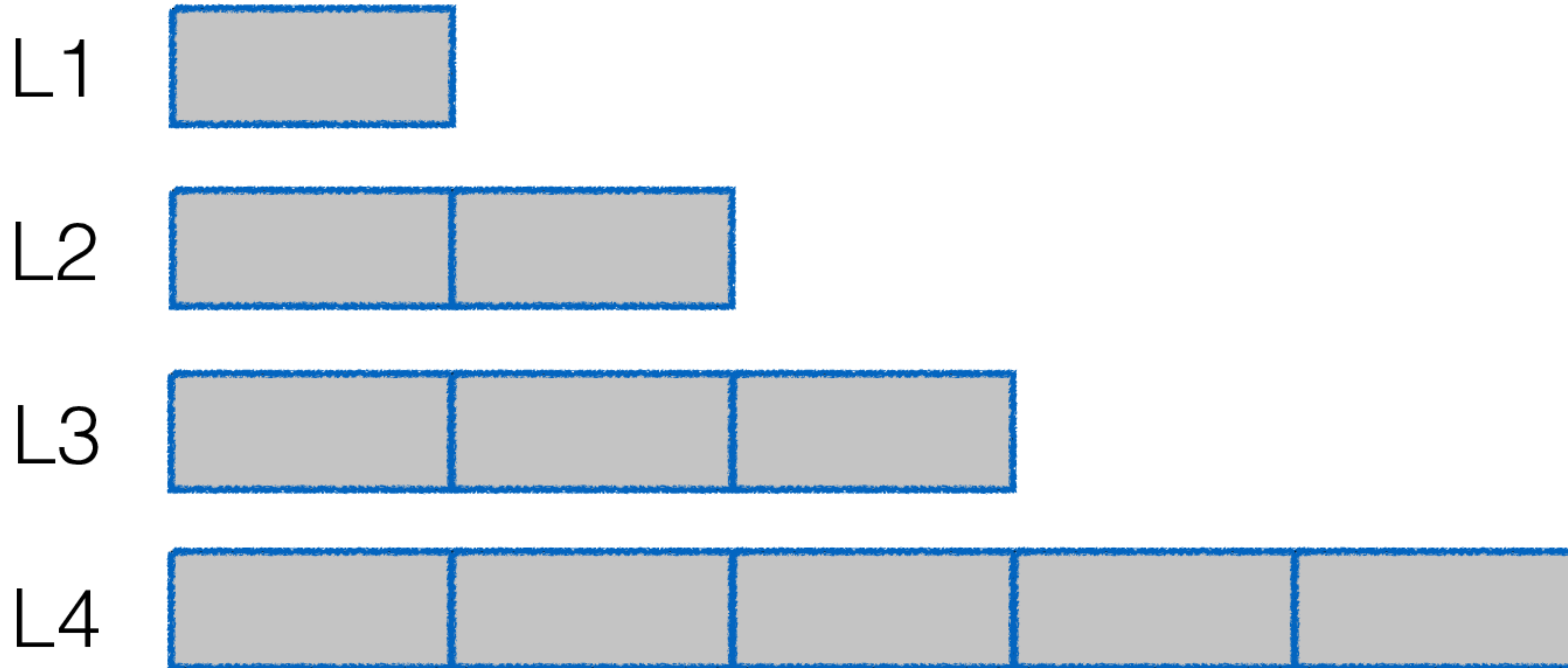
FAst DElete

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



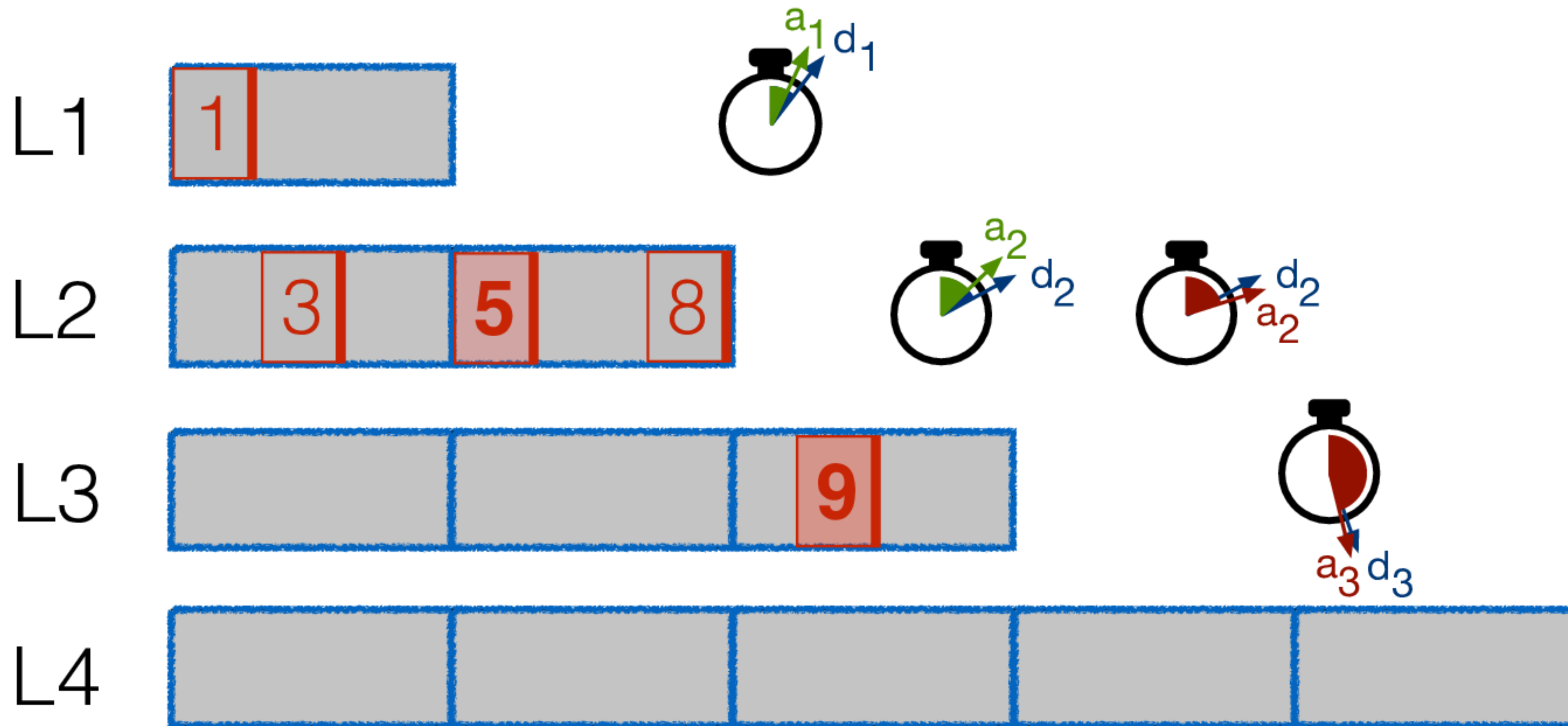
FAst DElete

breaking ties in practical workloads



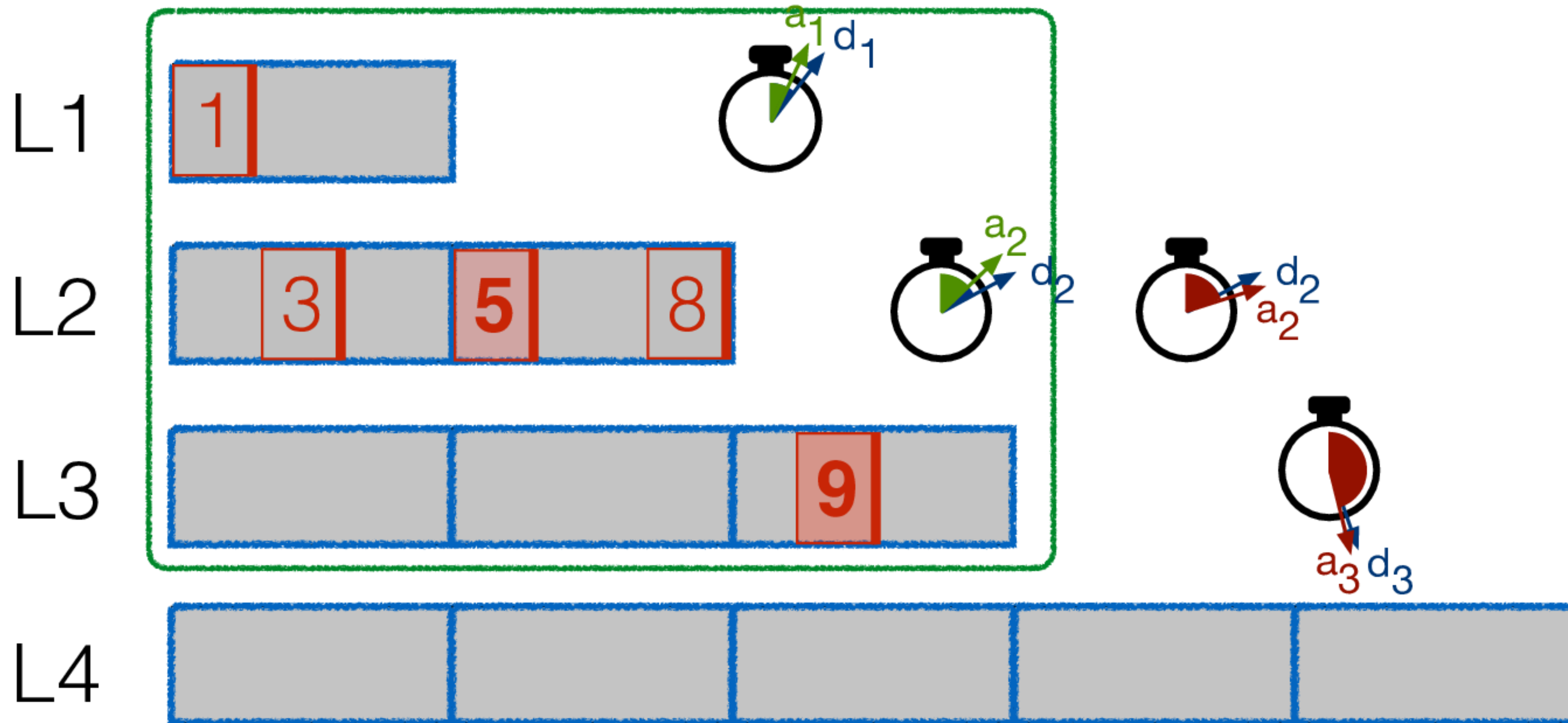
FASt DElete

breaking ties in practical workloads



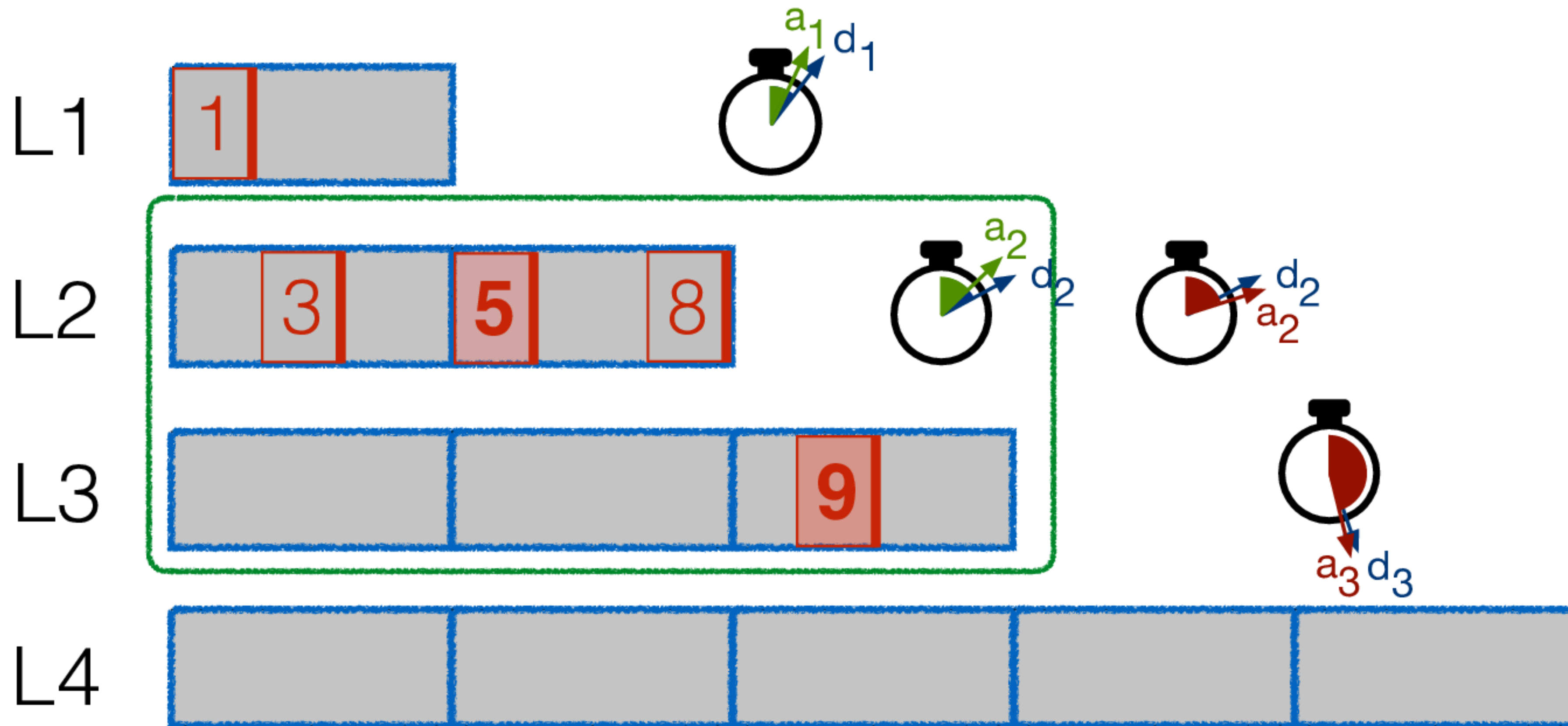
FAst DElete

breaking ties in practical workloads



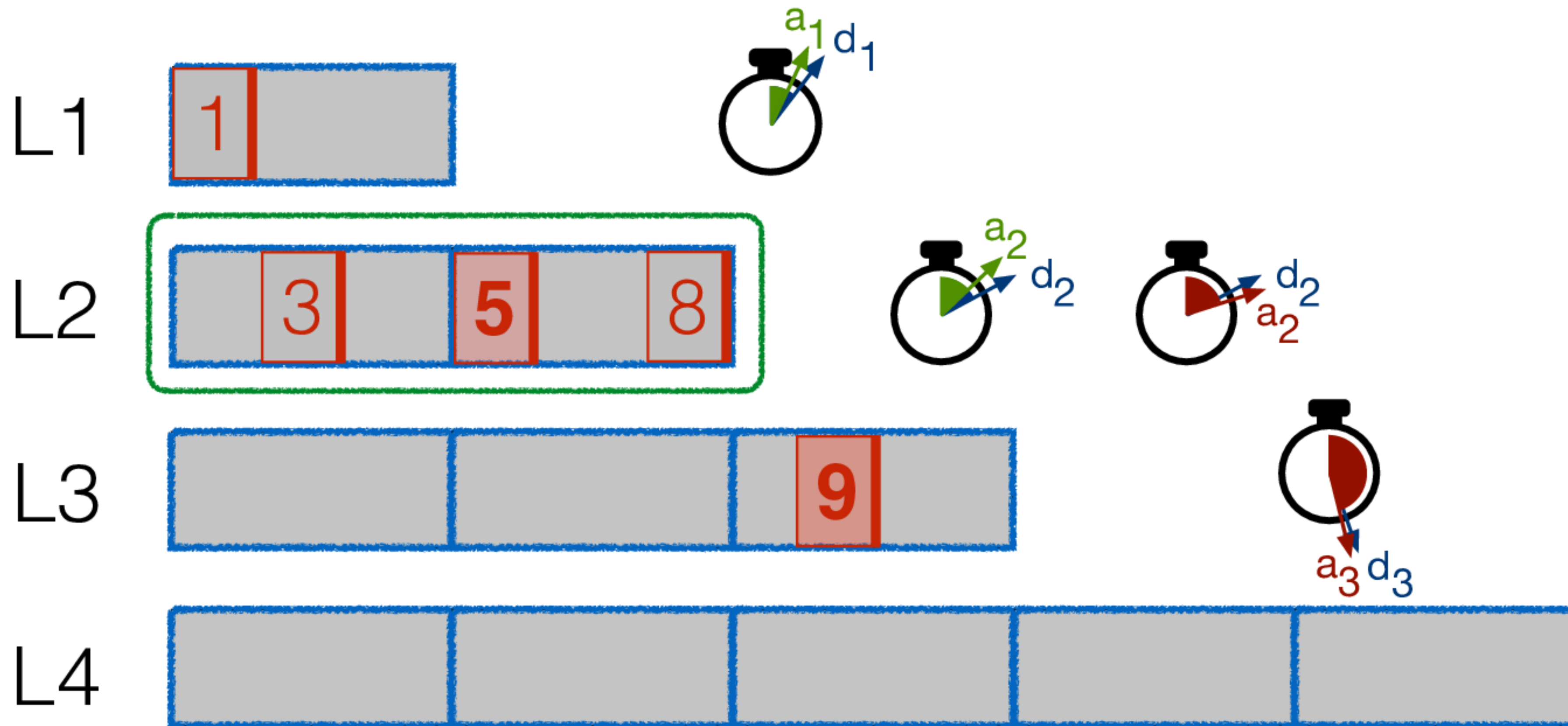
FAst DElete

breaking ties in practical workloads



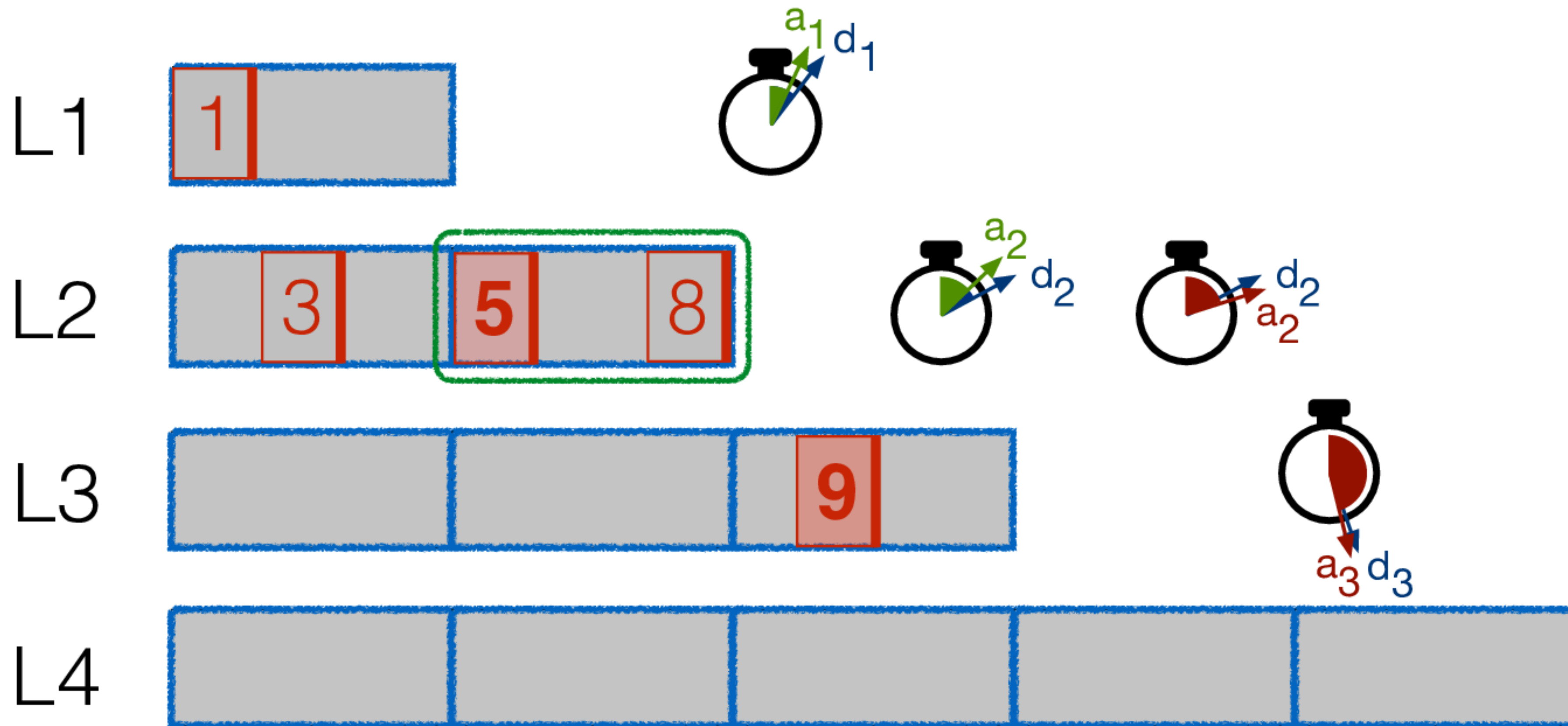
FASt DElete

breaking ties in practical workloads



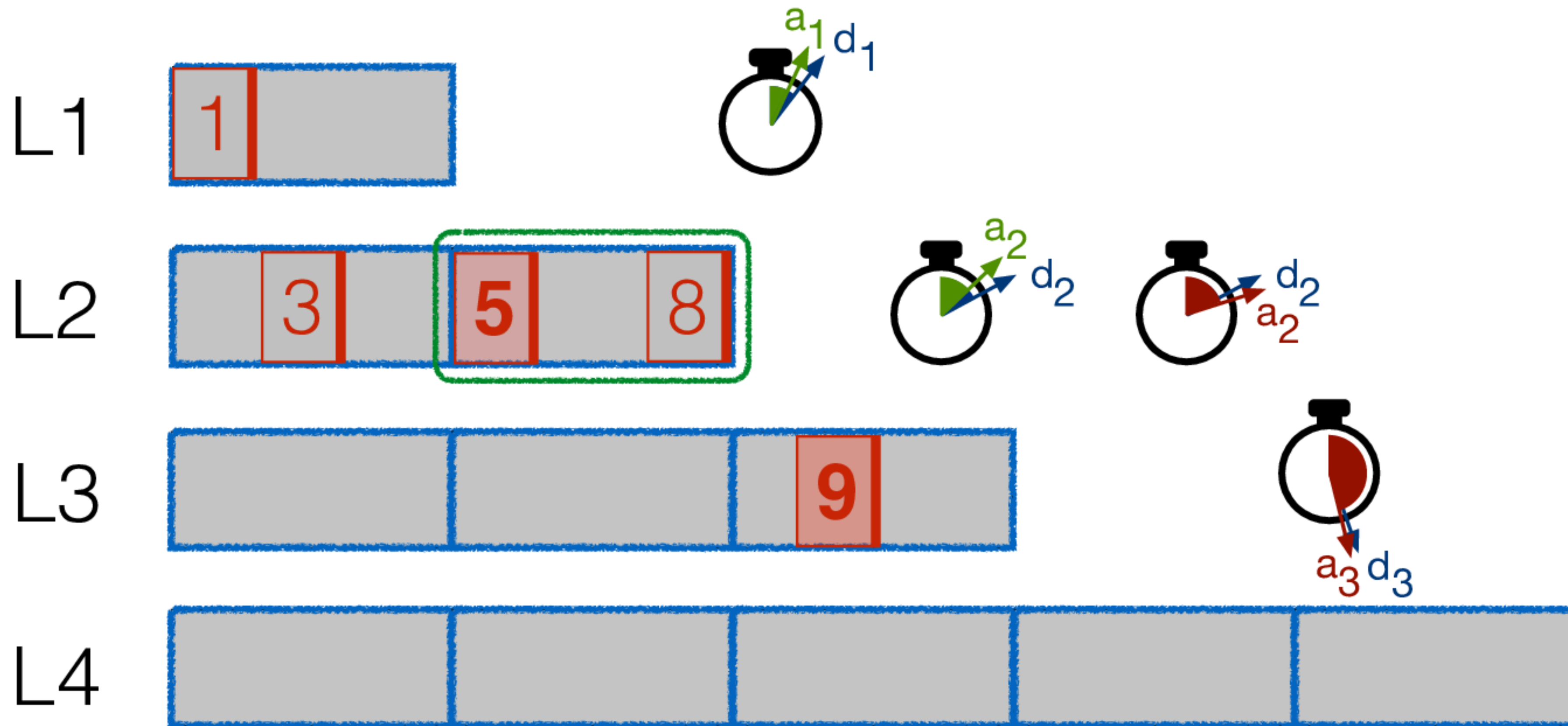
FASt DElete

breaking ties in practical workloads

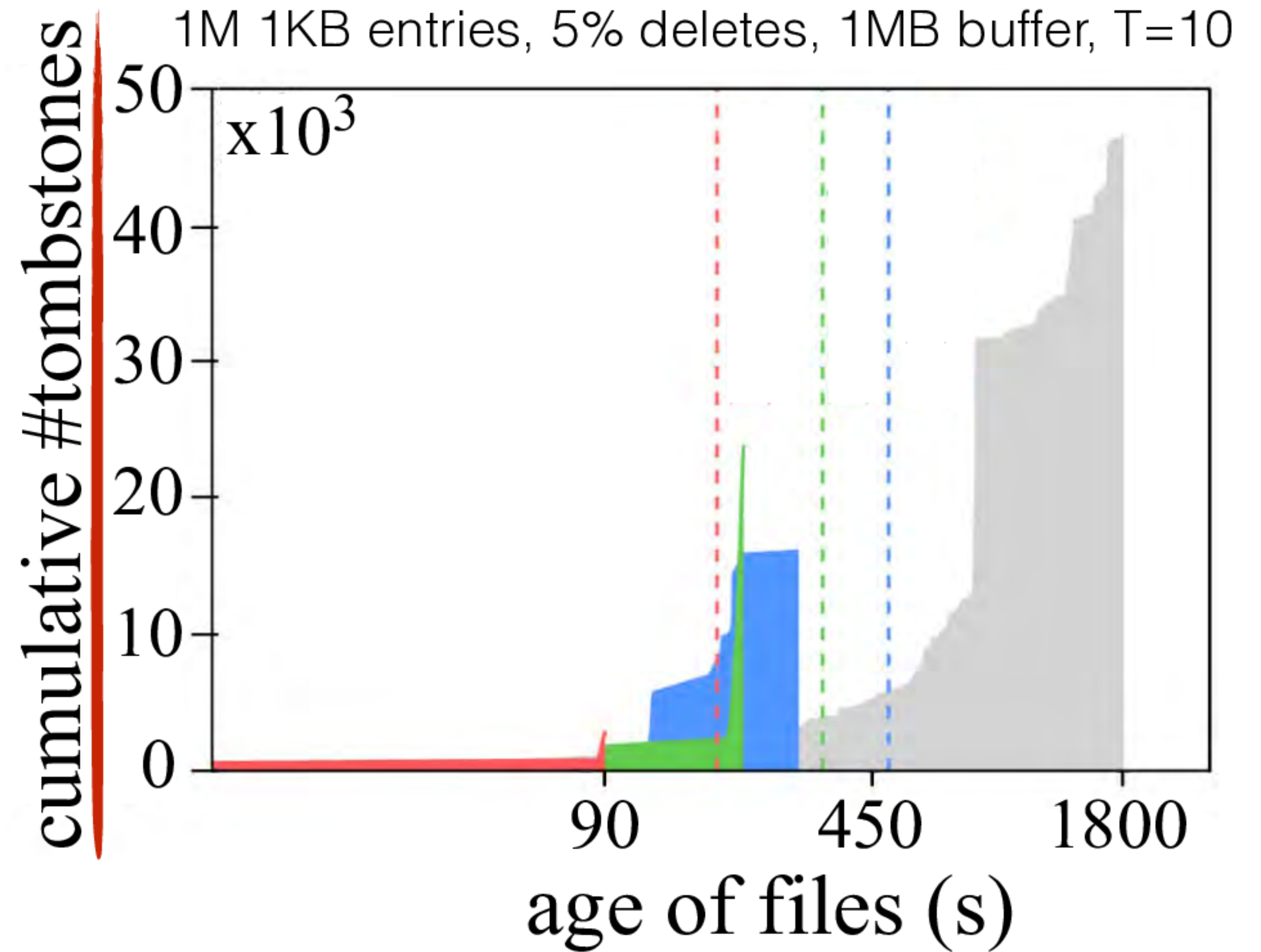


FASt DElete

breaking ties in practical workloads



FAst DElete



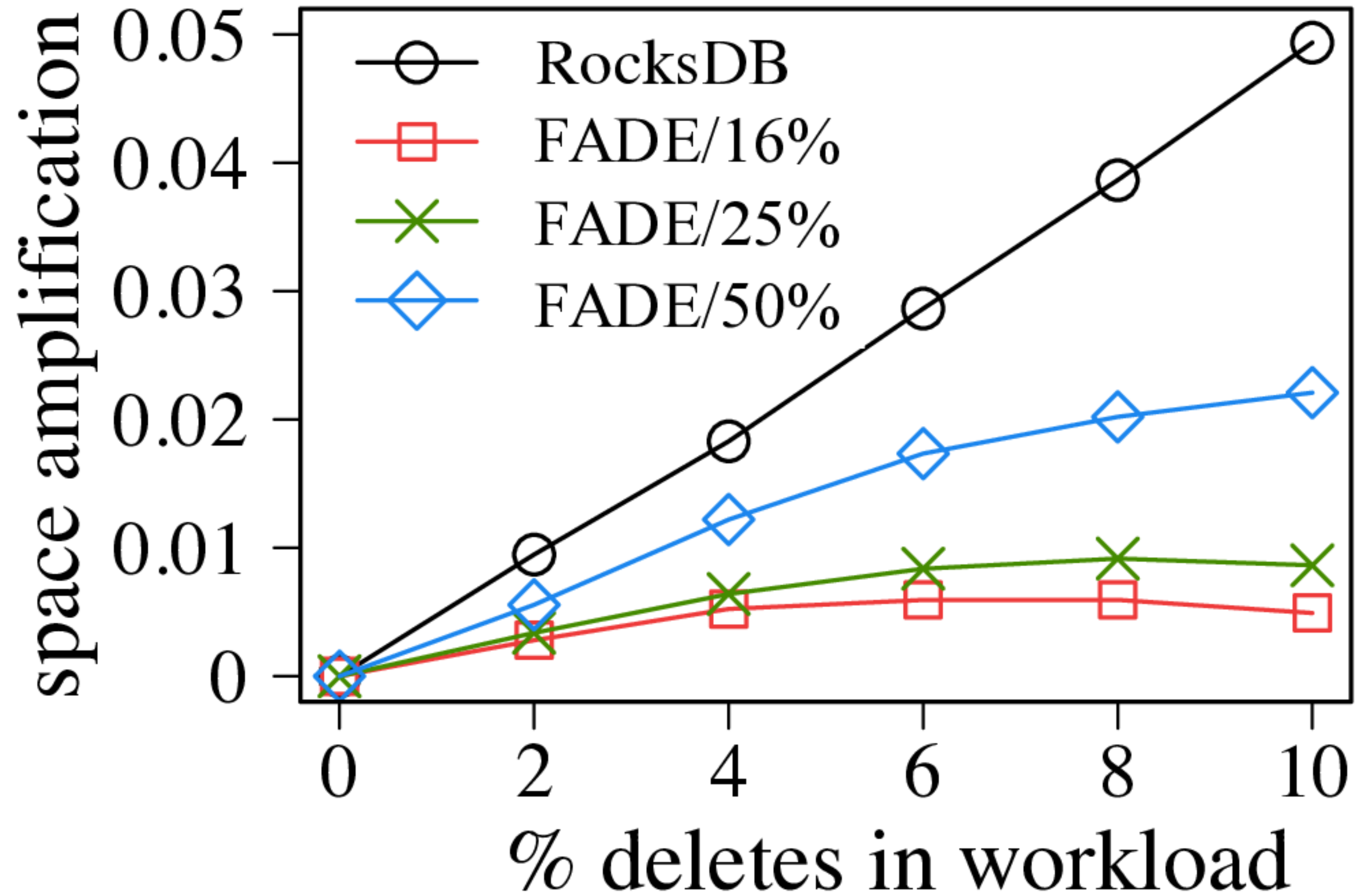
timely delete persistence

within D_{th}



FAst DElete

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



reduced space amplification ✓

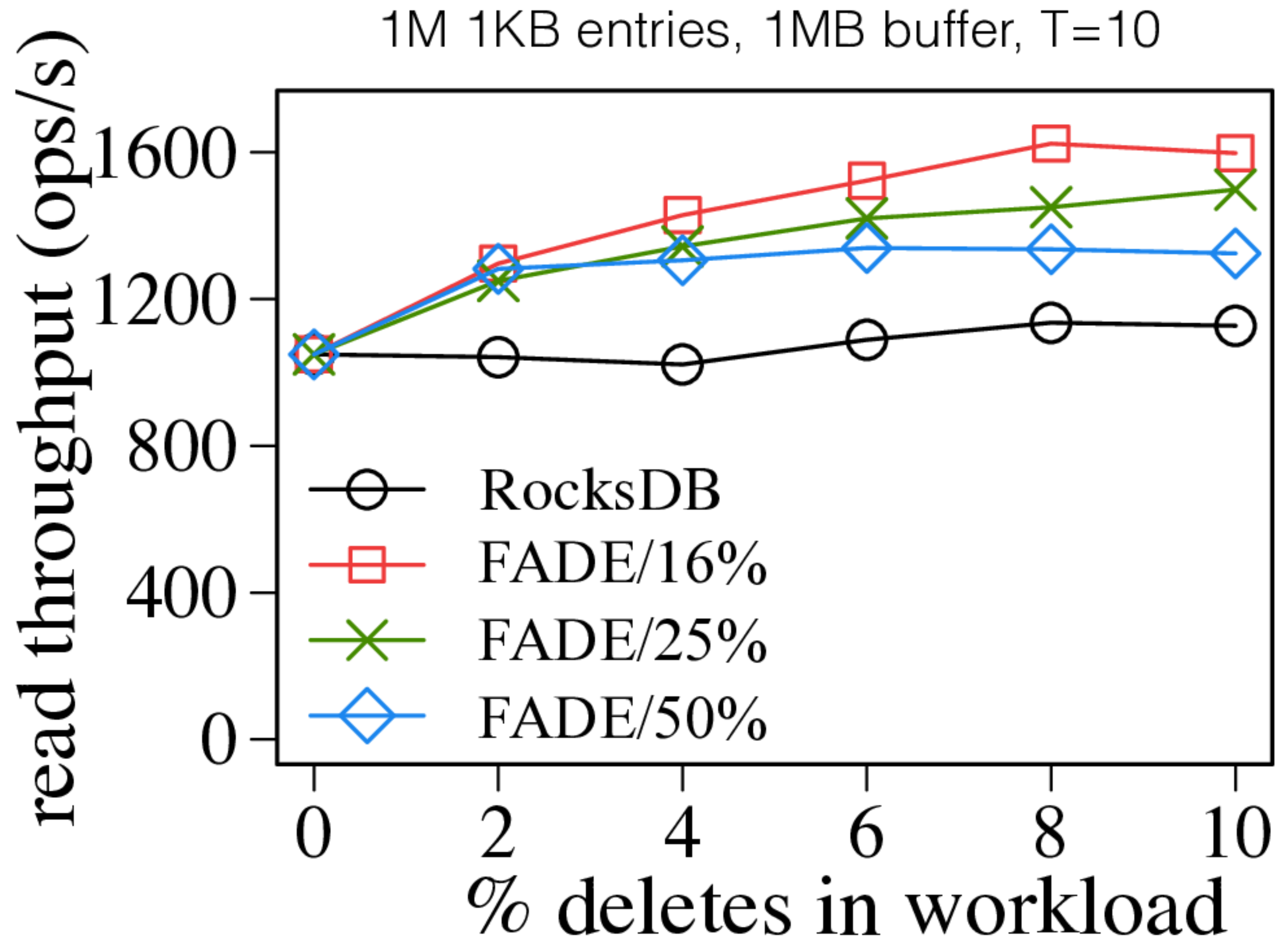
2.1 - 9.8x

timely delete persistence ✓

within D_{th}

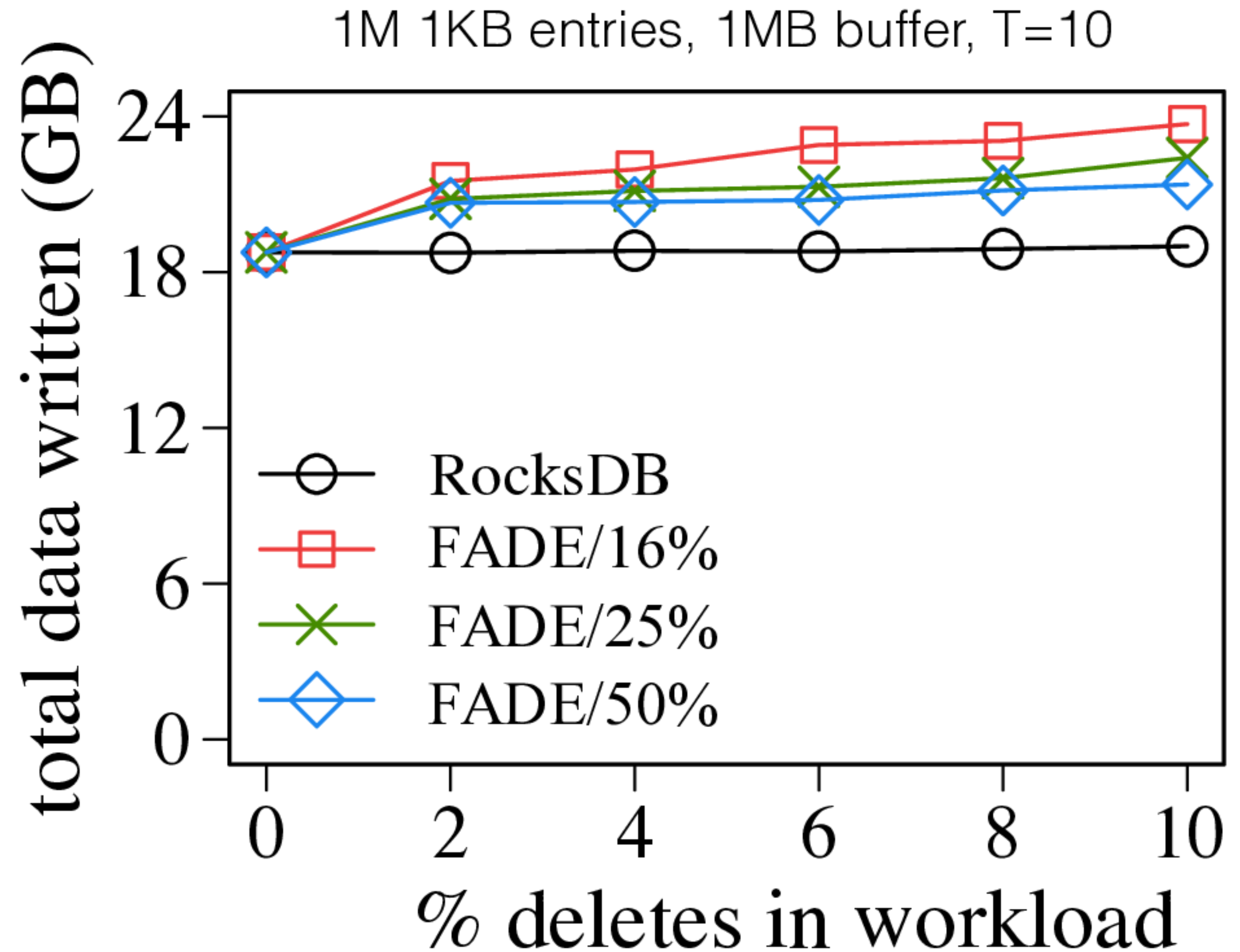
FAst DElete

- 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 ↑
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

4 - 25%



improved read performance

1.2 - 1.4x



reduced space amplification

2.1 - 9.8x

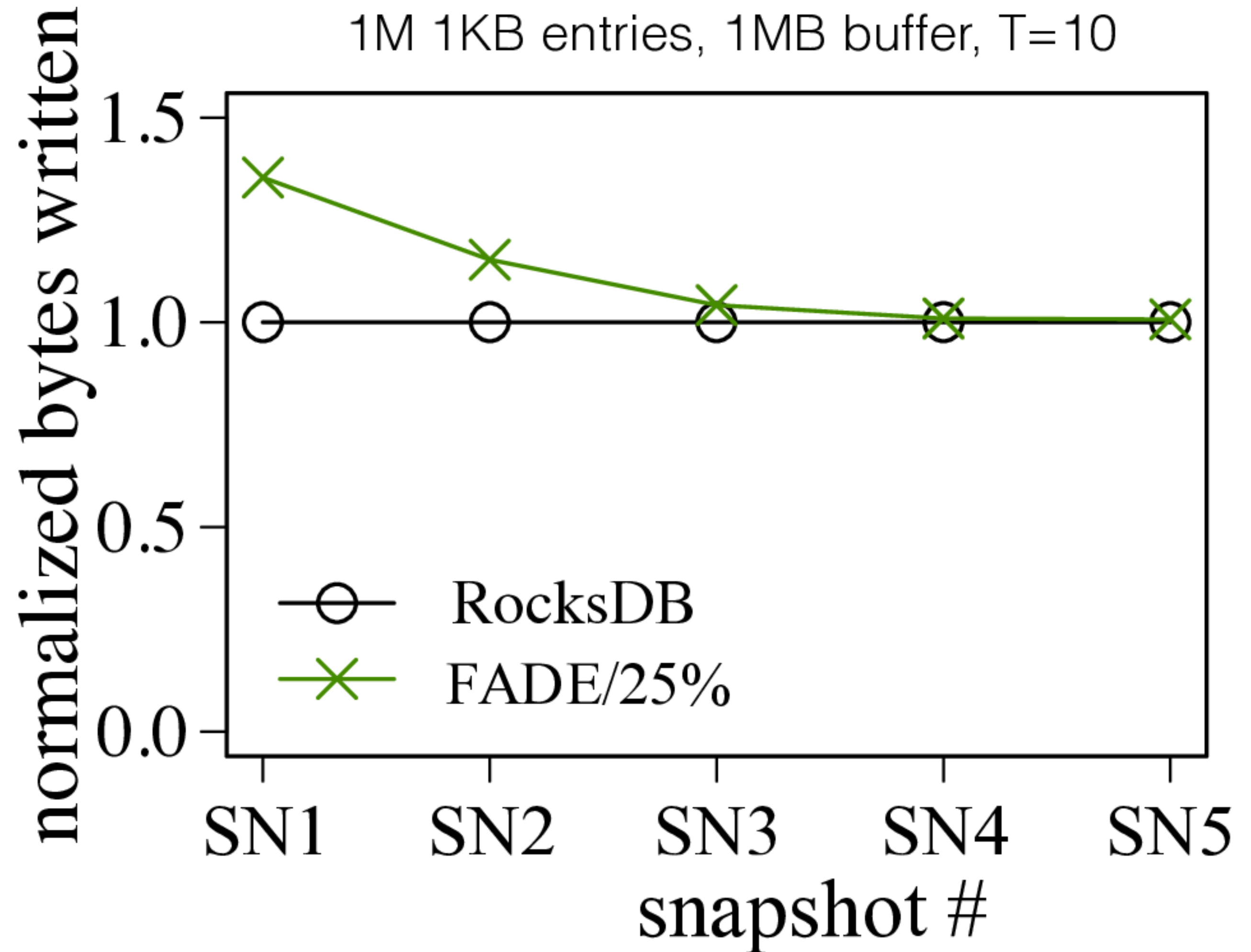


timely delete persistence

within D_{th}

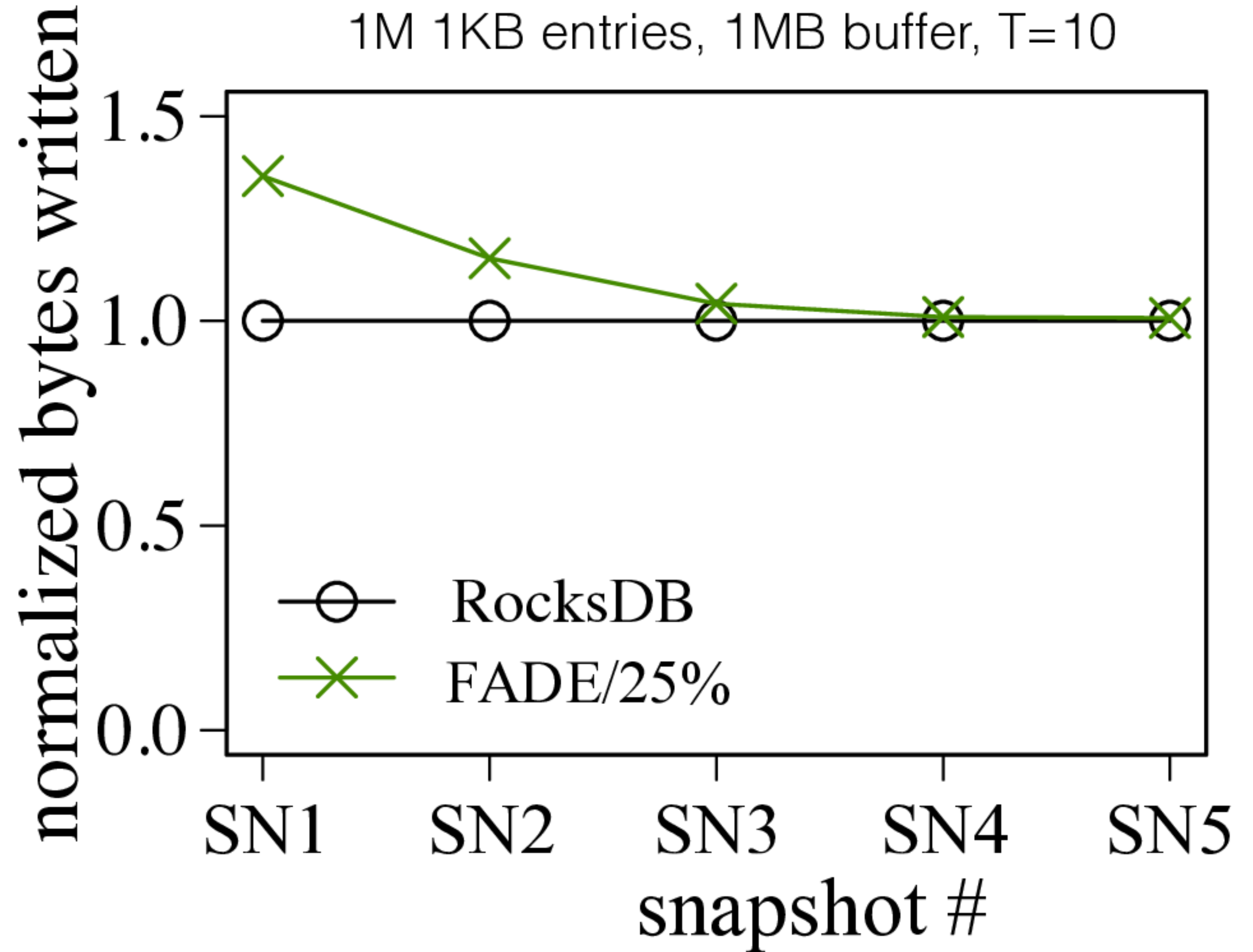


1M 1KB entries, 1MB buffer, T=10



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

improved read performance

reduced space amplification

timely delete persistence

FADE

latency spikes

Kiwi

superfluous I/Os

the solution

FASt DElete

higher write amplification



poor read performance



reduced space amplification



timely delete persistence



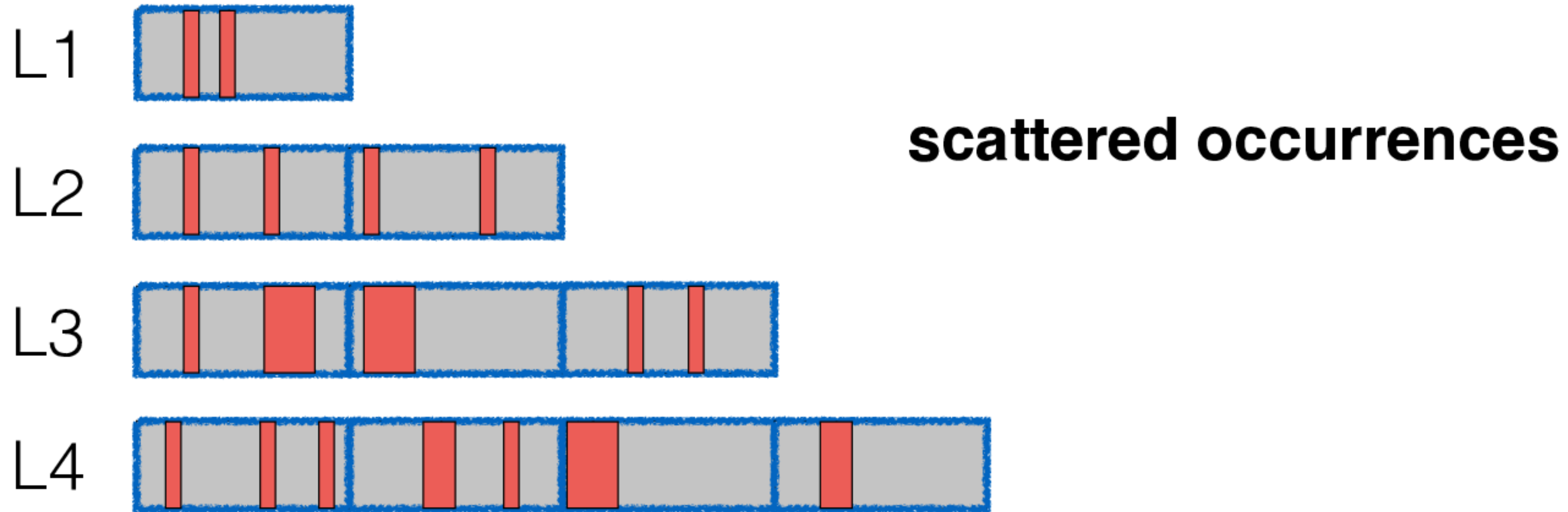
FADE

latency spikes

Kiwi

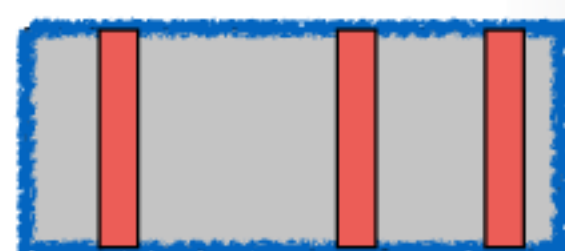
Key Weaving storage layout

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

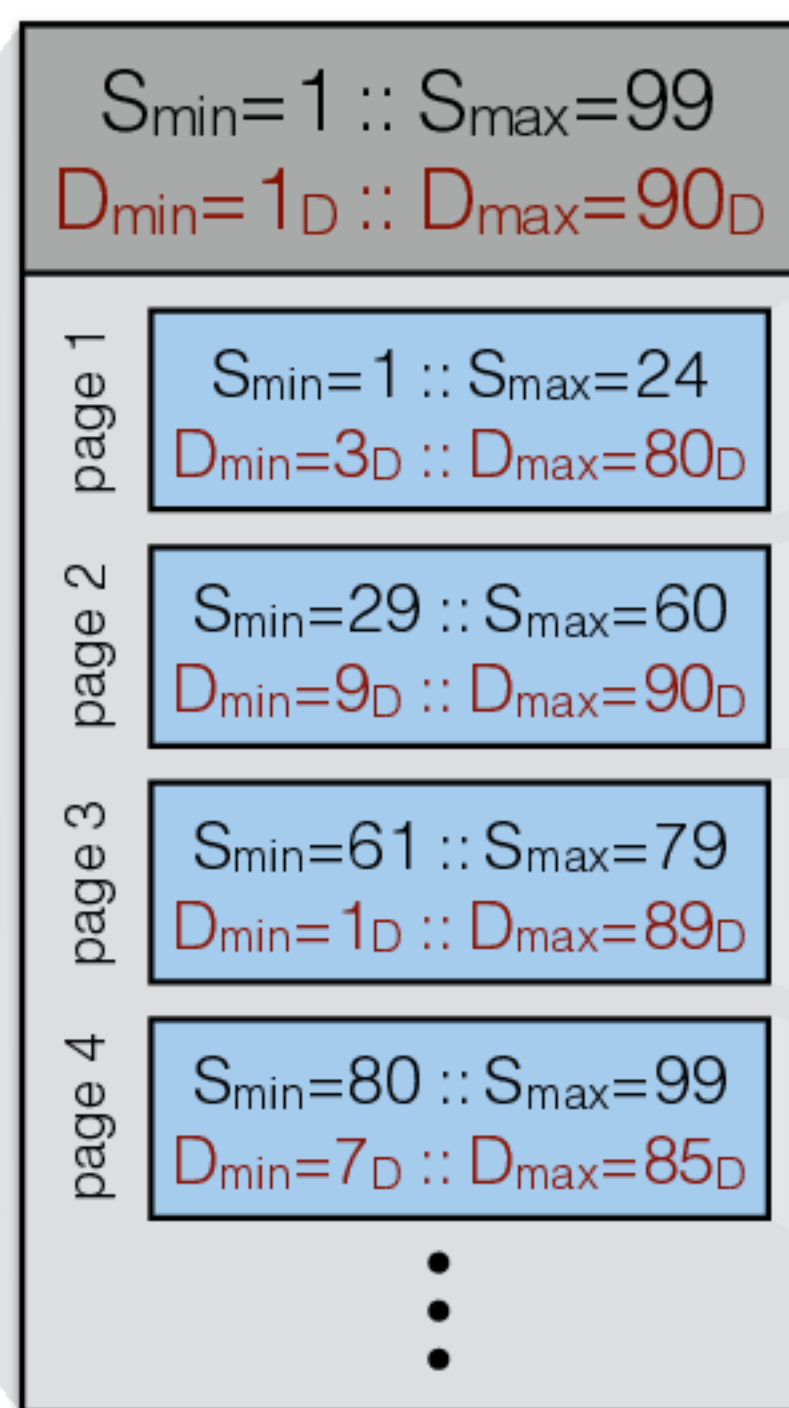


Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



SST file



page 1

1	4	9	14	15	19	20	24
34 _D	69 _D	3 _D	79 _D	8 _D	80 _D	23 _D	24 _D

page 2

29	32	33	40	44	52	56	60
88 _D	90 _D	28 _D	74 _D	9 _D	76 _D	81 _D	64 _D

page 3

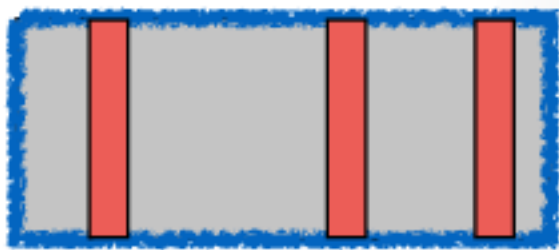
61	63	67	71	72	73	78	79
75 _D	82 _D	1 _D	67 _D	77 _D	89 _D	65 _D	12 _D

page 4

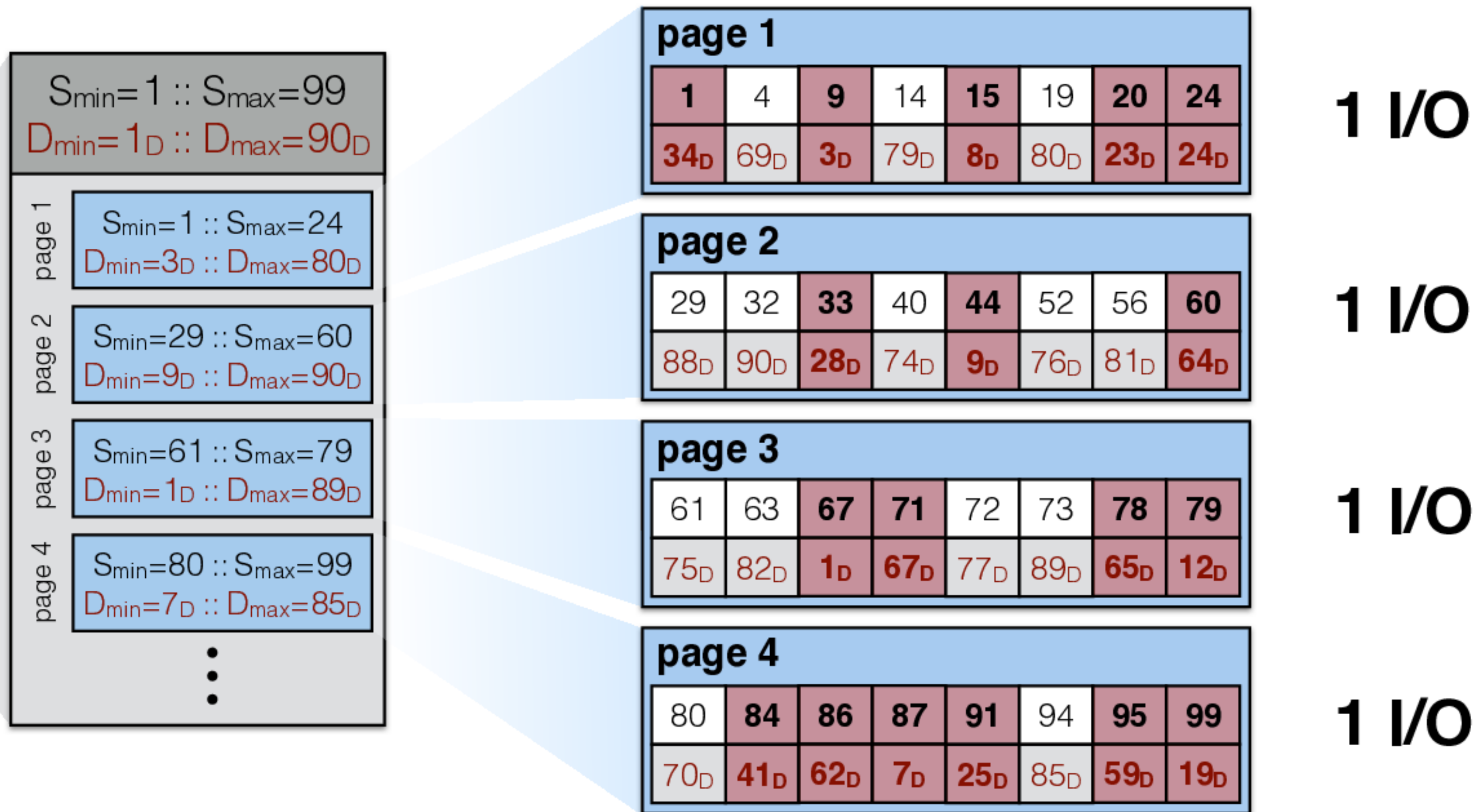
80	84	86	87	91	94	95	99
70 _D	41 _D	62 _D	7 _D	25 _D	85 _D	59 _D	19 _D

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

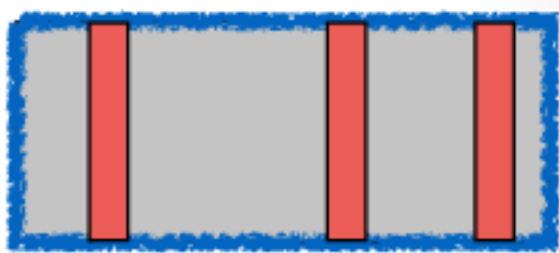


SST file

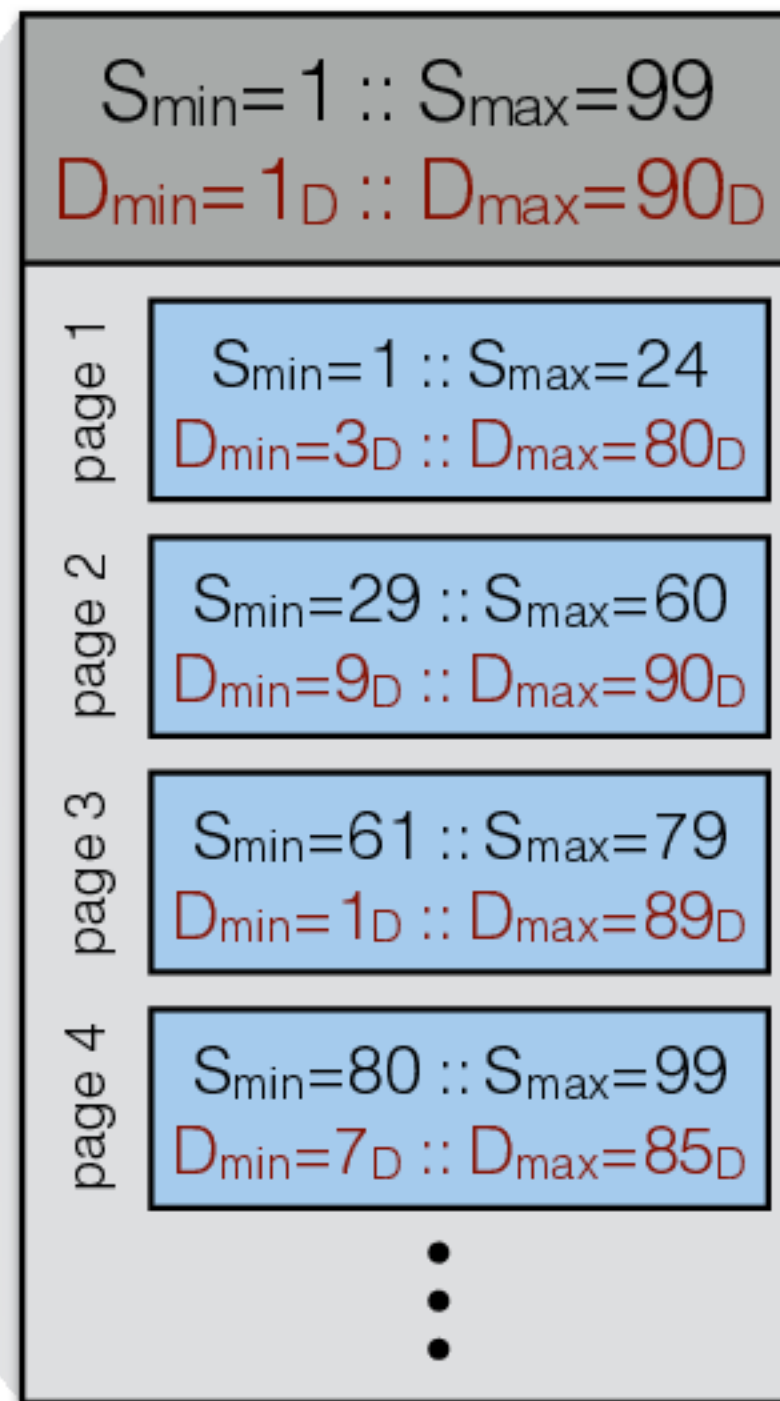


Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

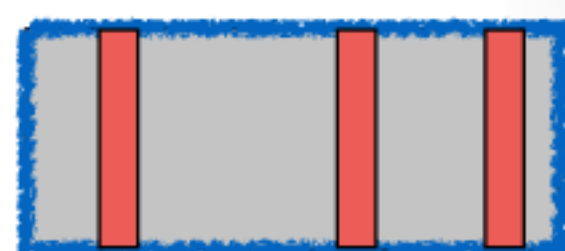


SST file

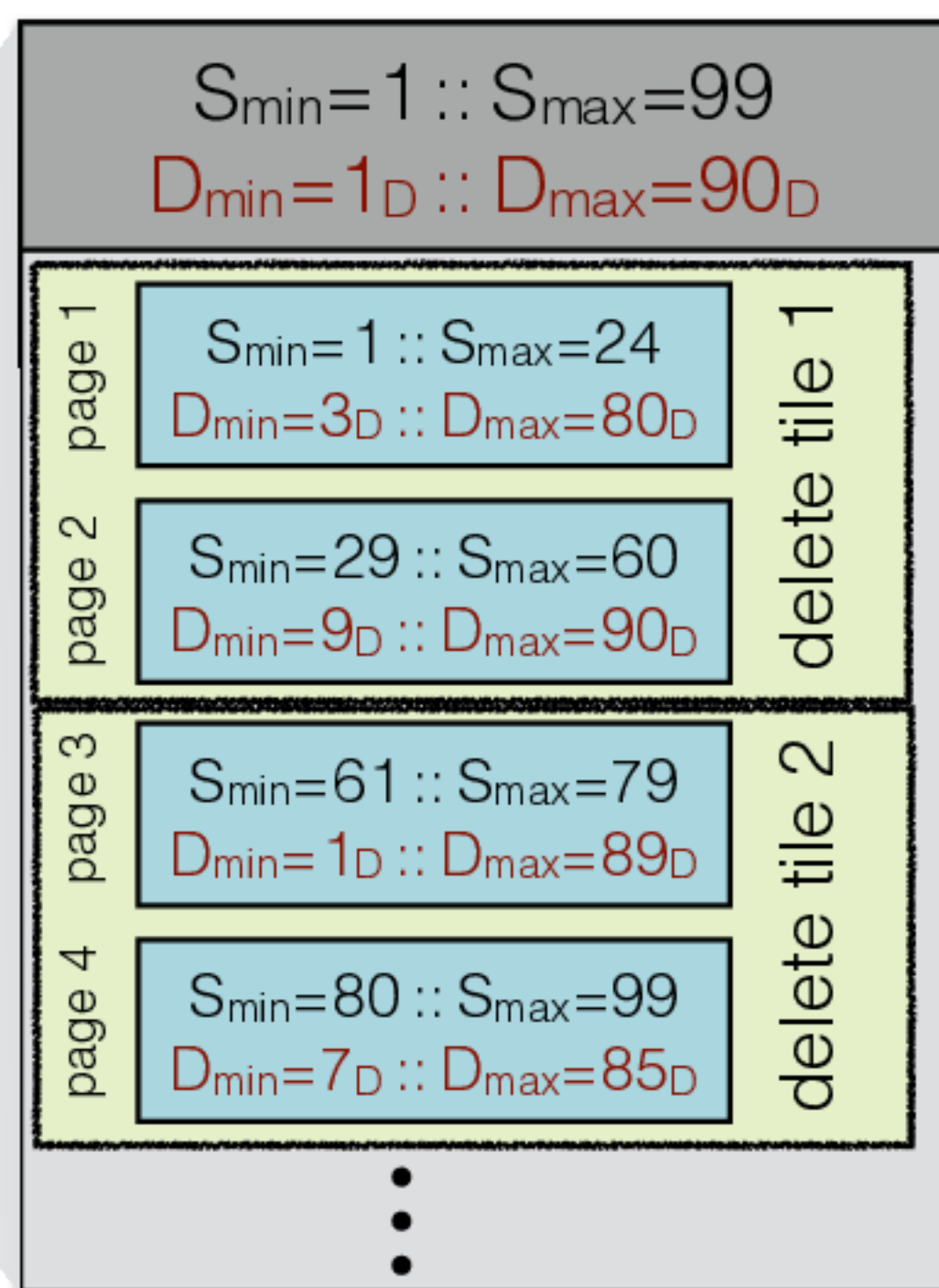


Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



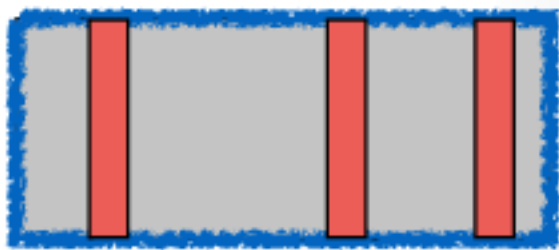
SST file



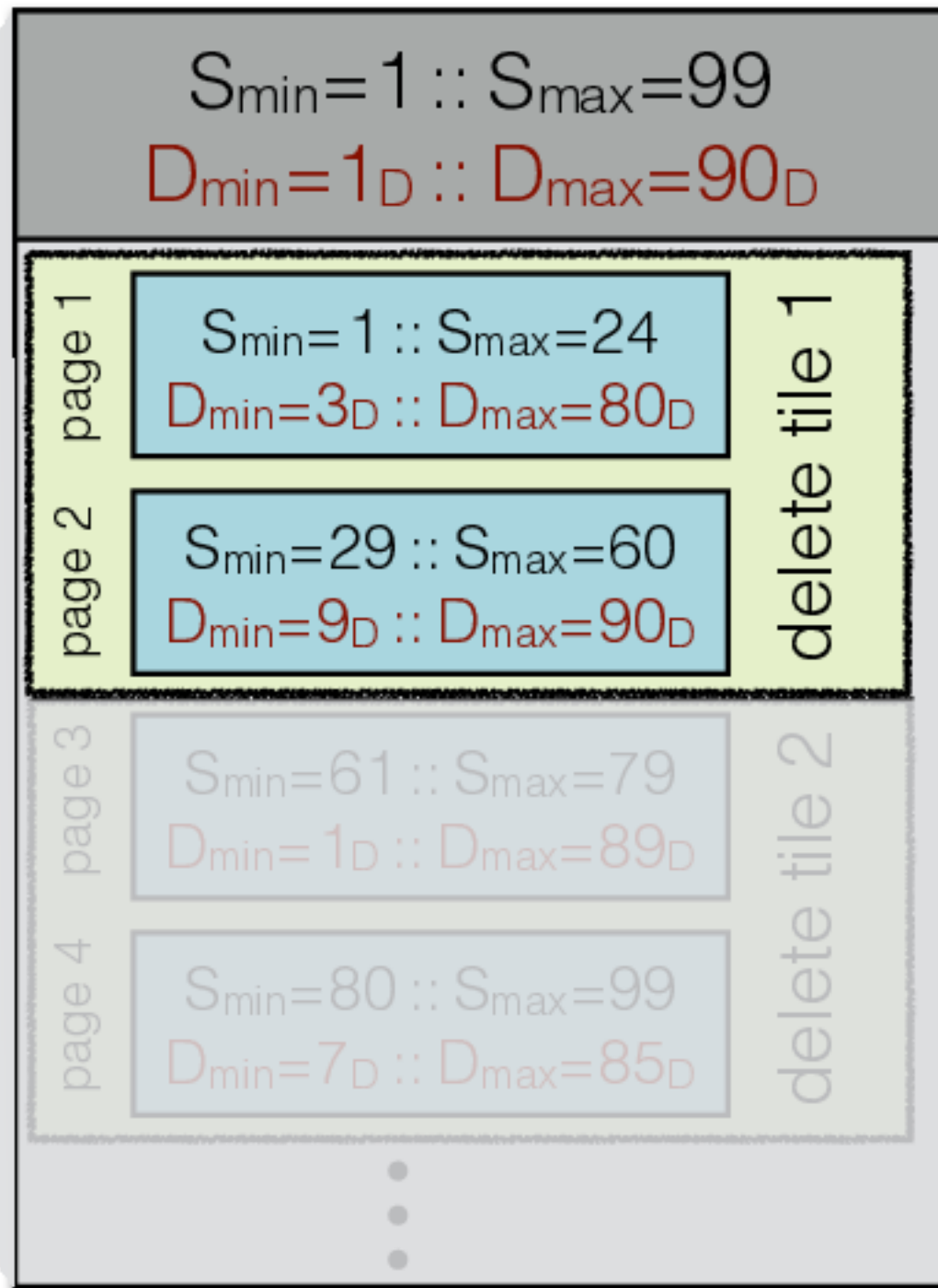
partitioned on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



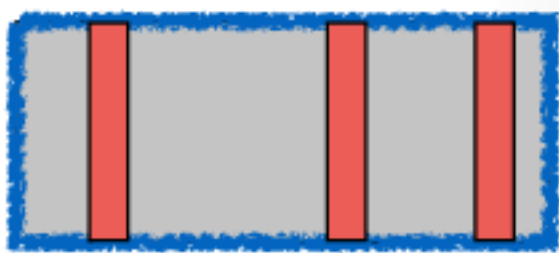
SST file



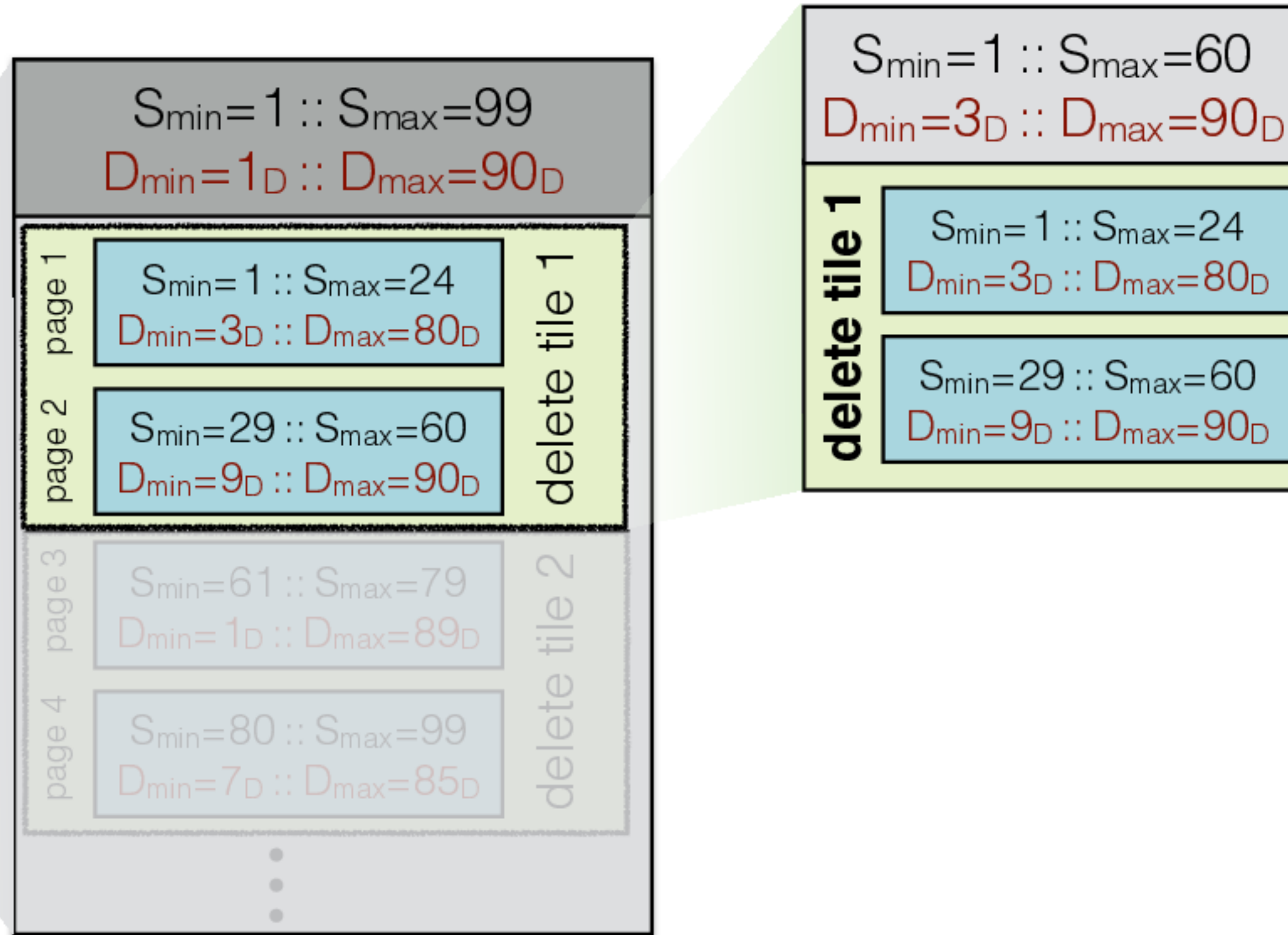
partitioned on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



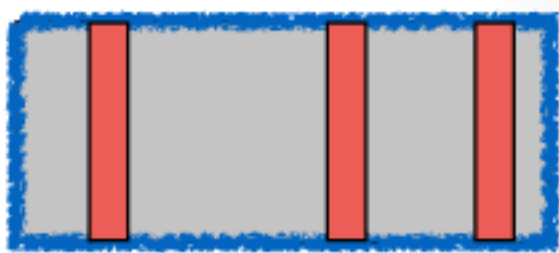
SST file



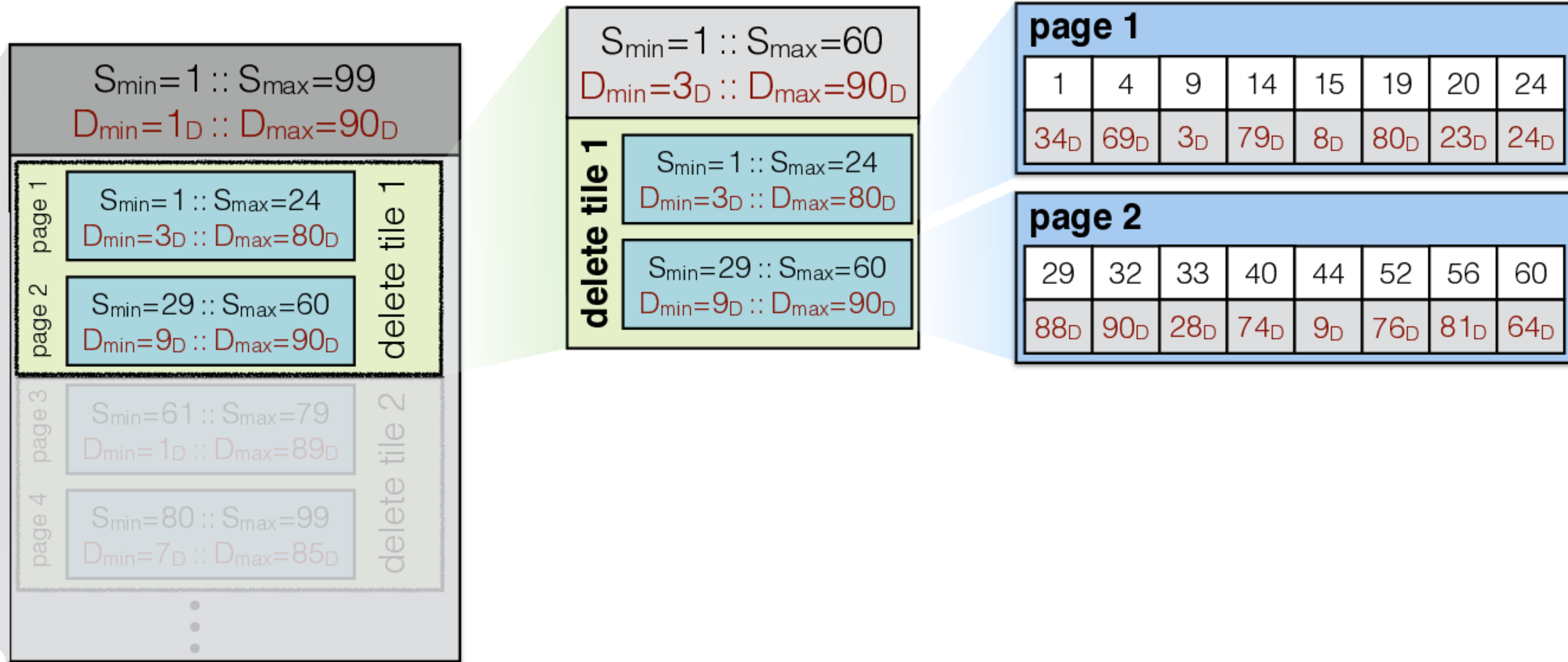
partitioned on S

Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$



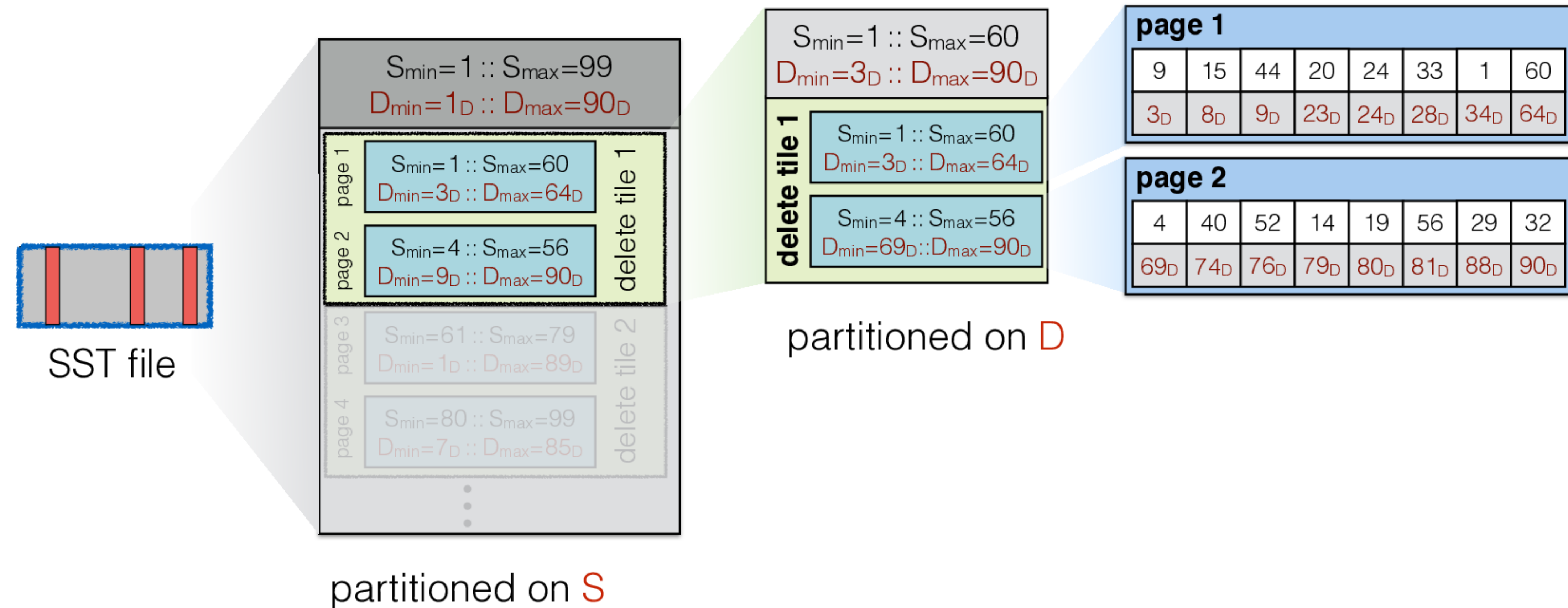
SST file



partitioned on S

Key Weaving storage layout

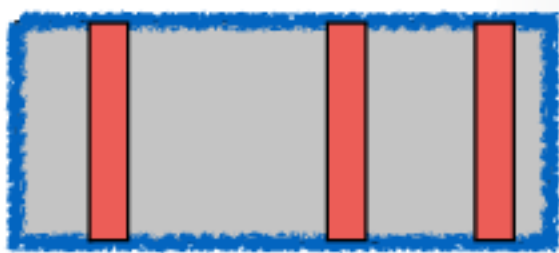
delete all entries with timestamp $\leq 65_D$



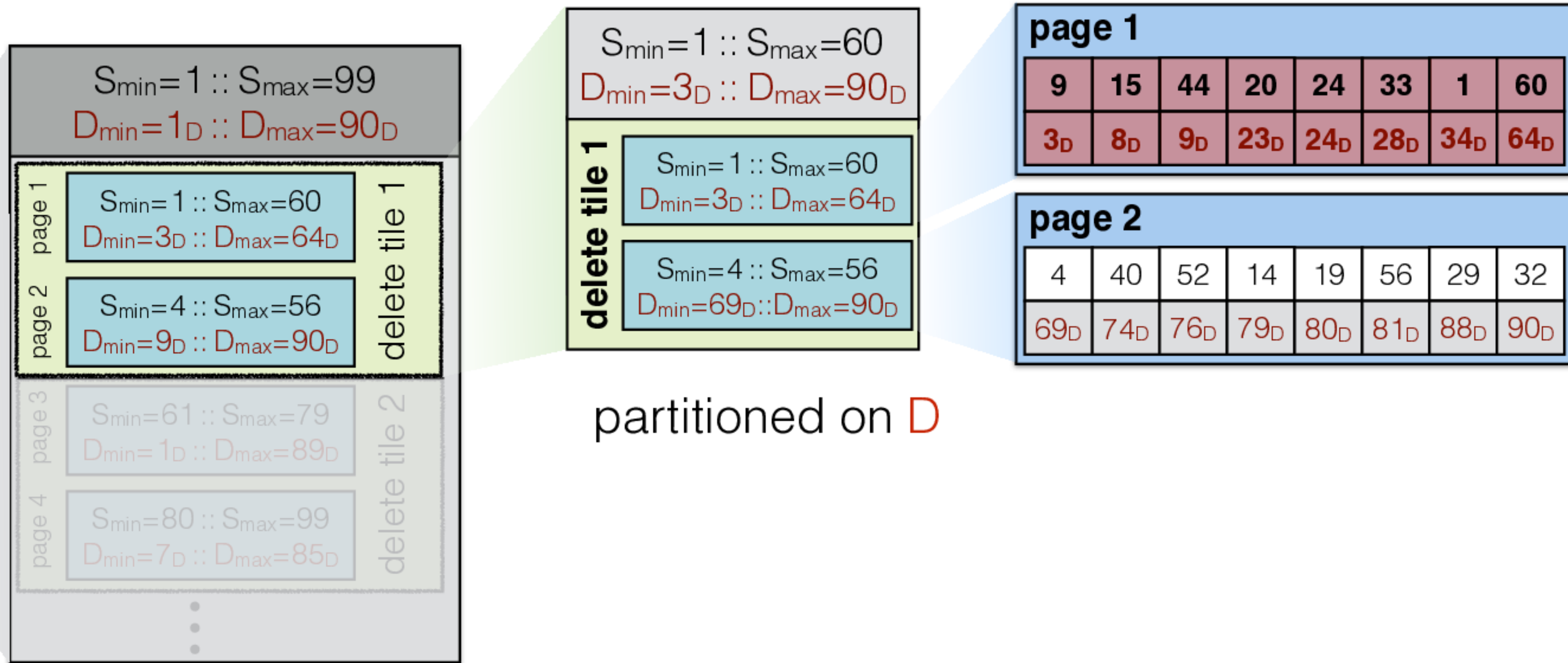
Key Weaving storage layout

delete all entries with timestamp $\leq 65_D$

drop
page



SST file

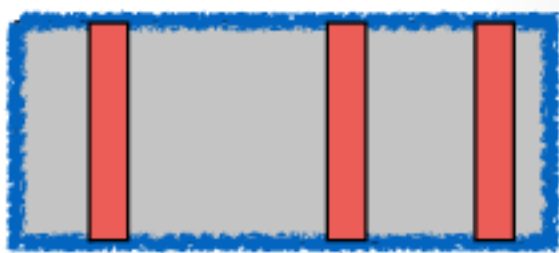


partitioned on S

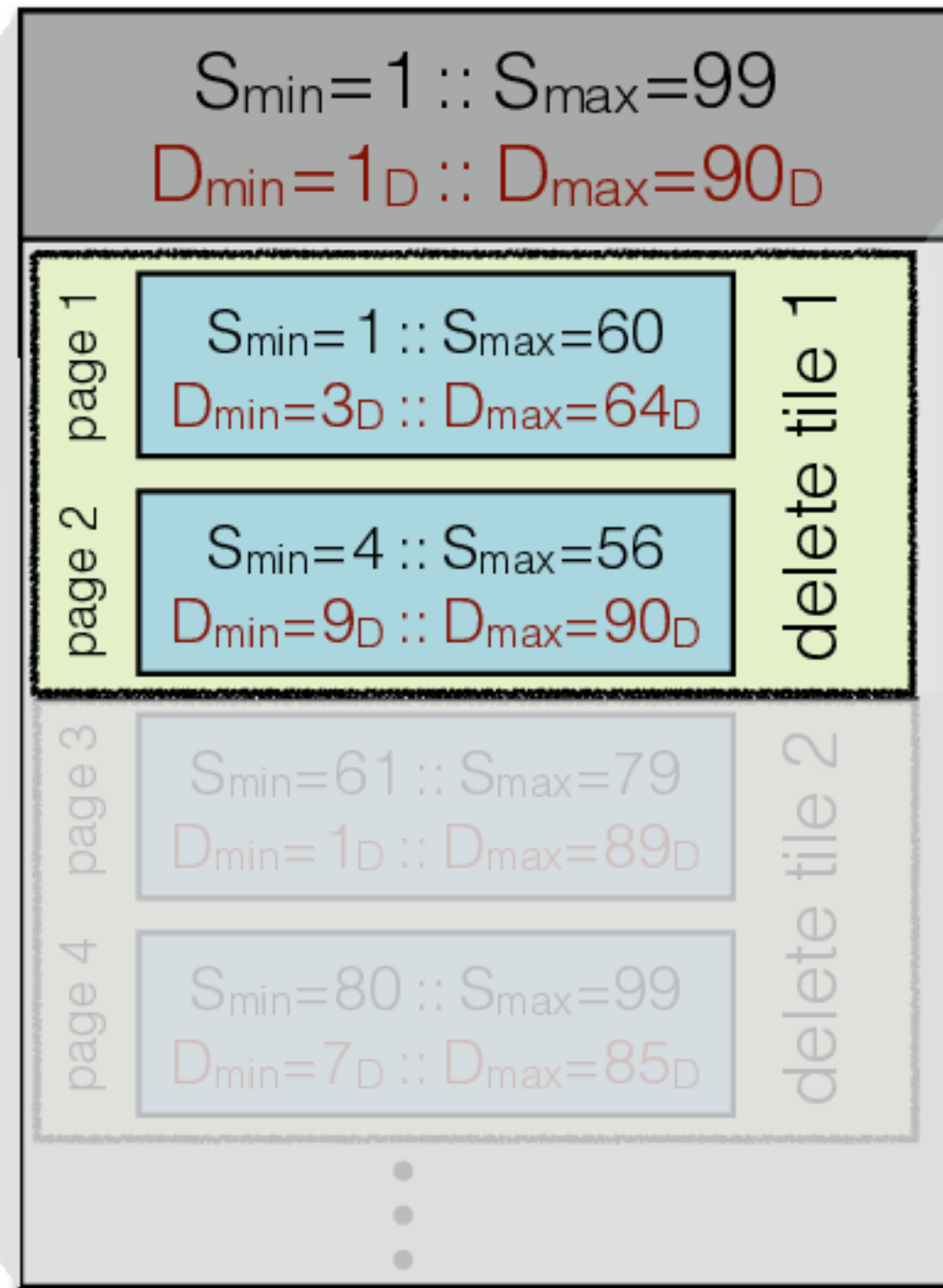
partitioned on D

Key Weaving storage layout

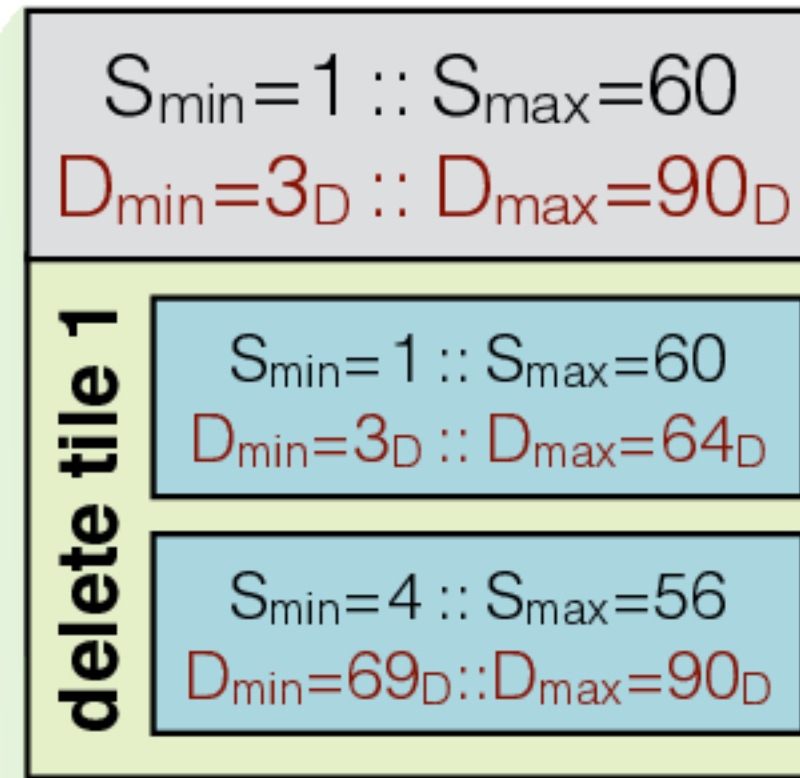
delete all entries with timestamp $\leq 65_D$



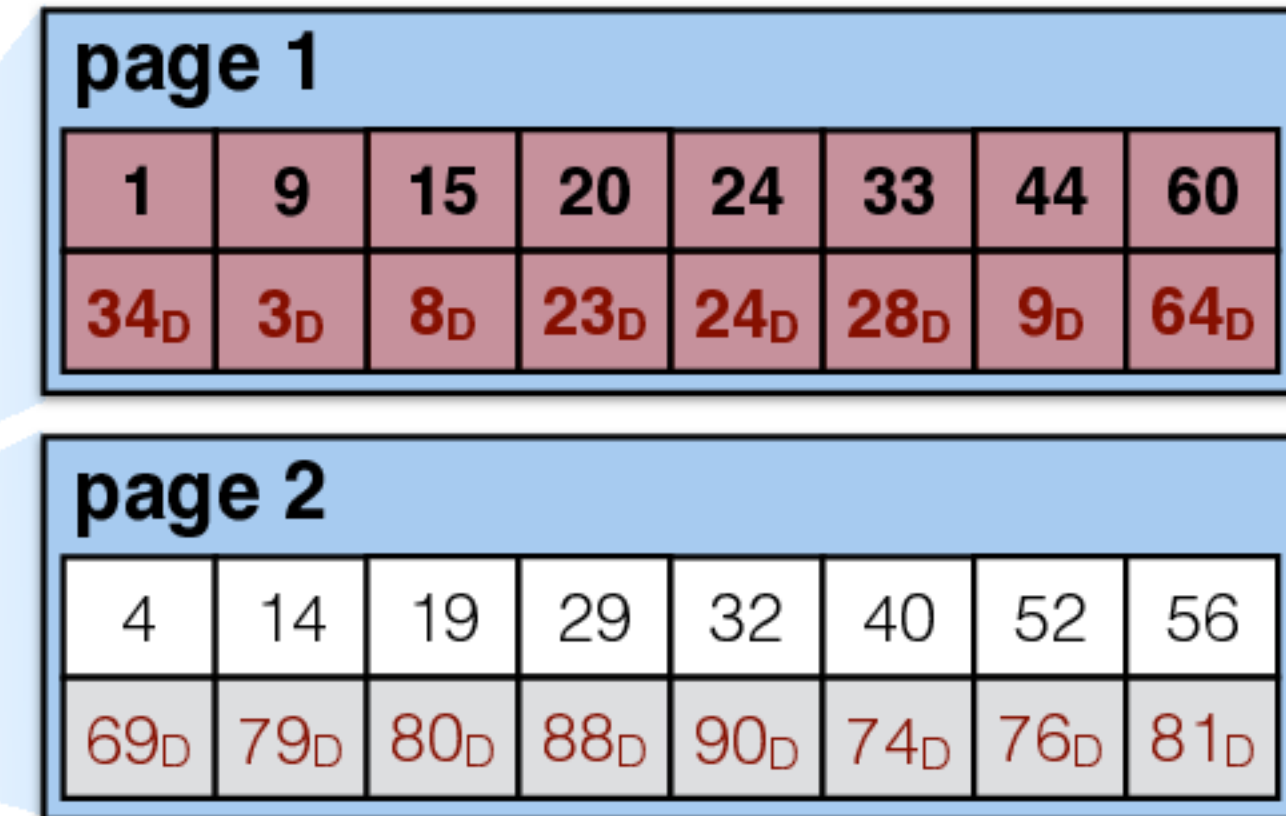
SST file



partitioned on S



partitioned on D

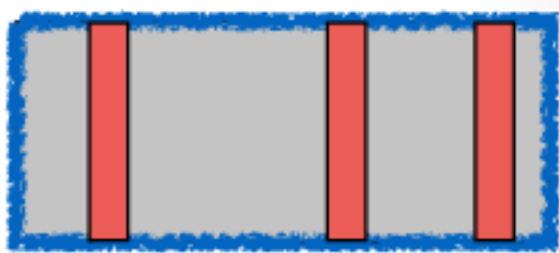


sorted on S

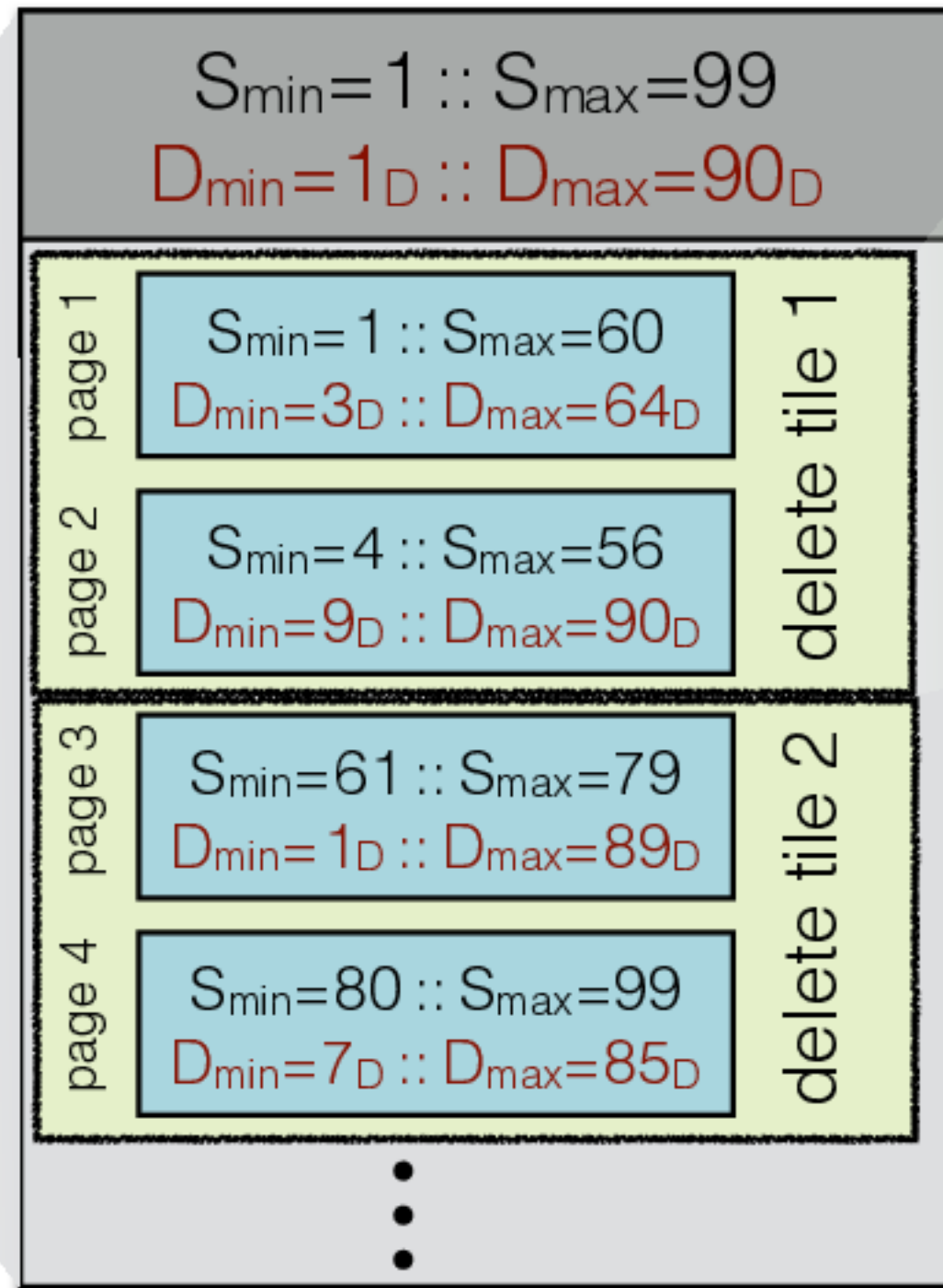
drop
page

Key Weaving storage layout

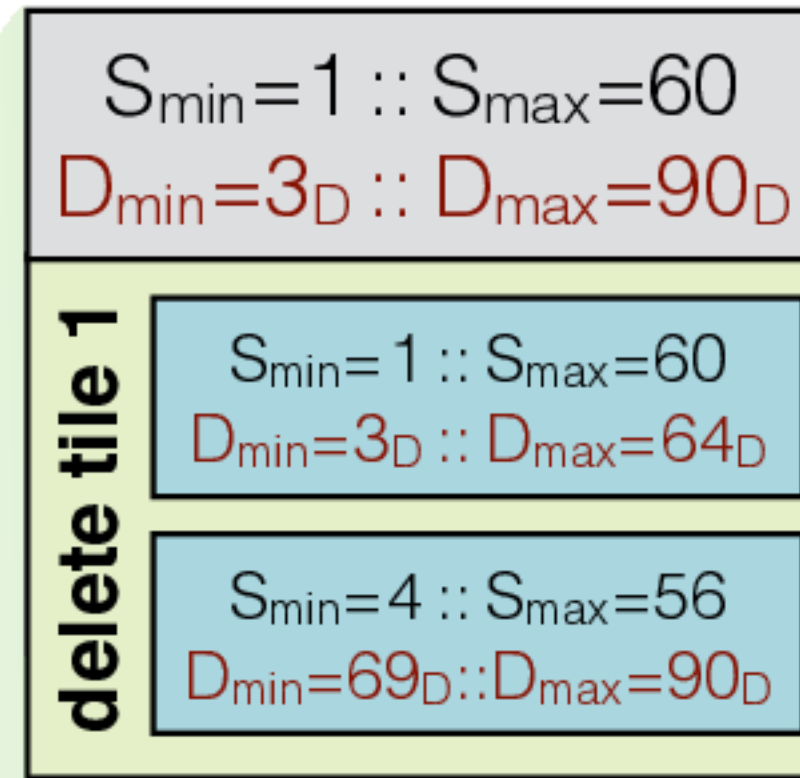
delete all entries with timestamp $\leq 65_D$



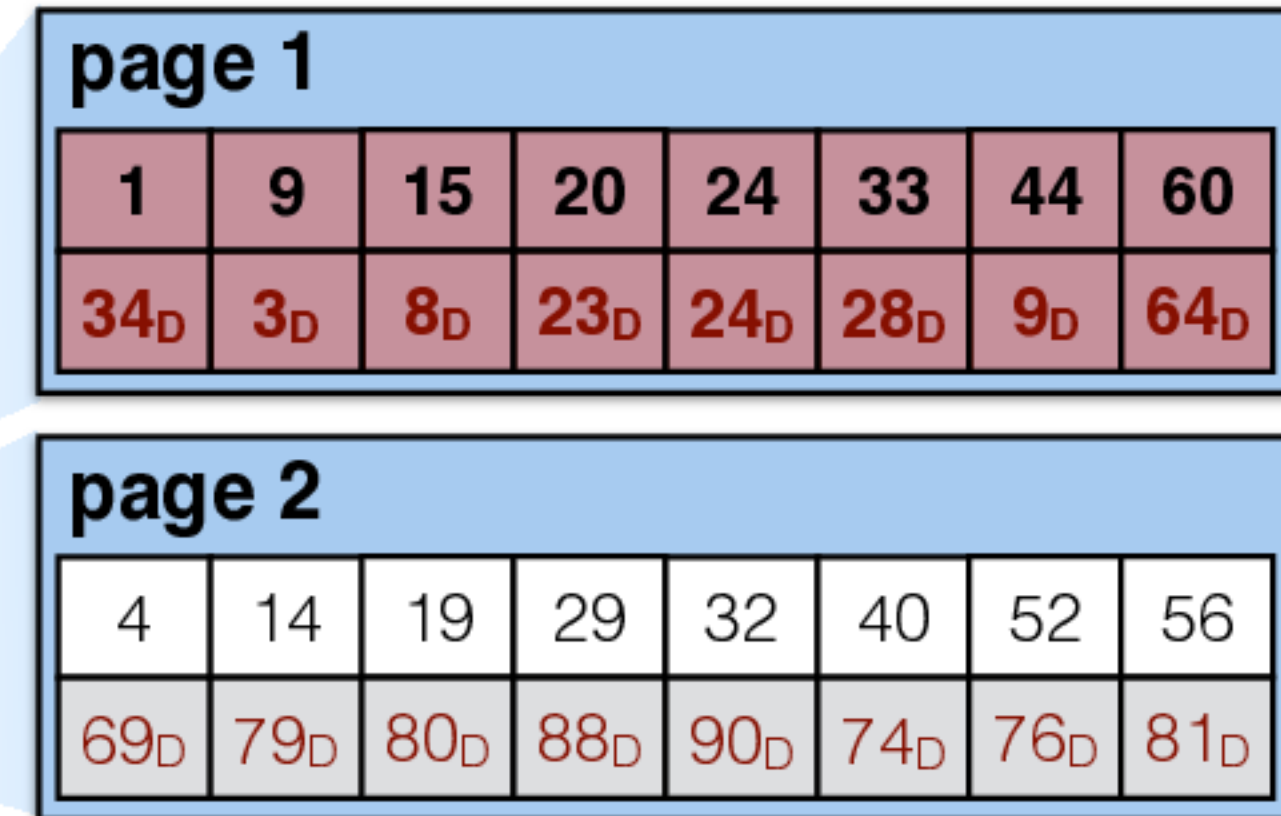
SST file



partitioned on S



partitioned on D

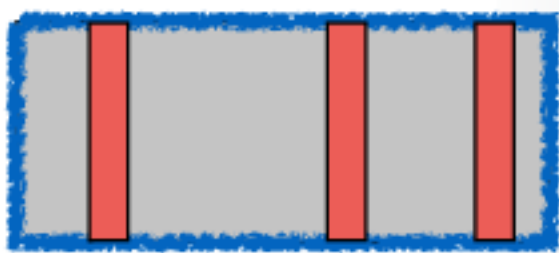


sorted on S

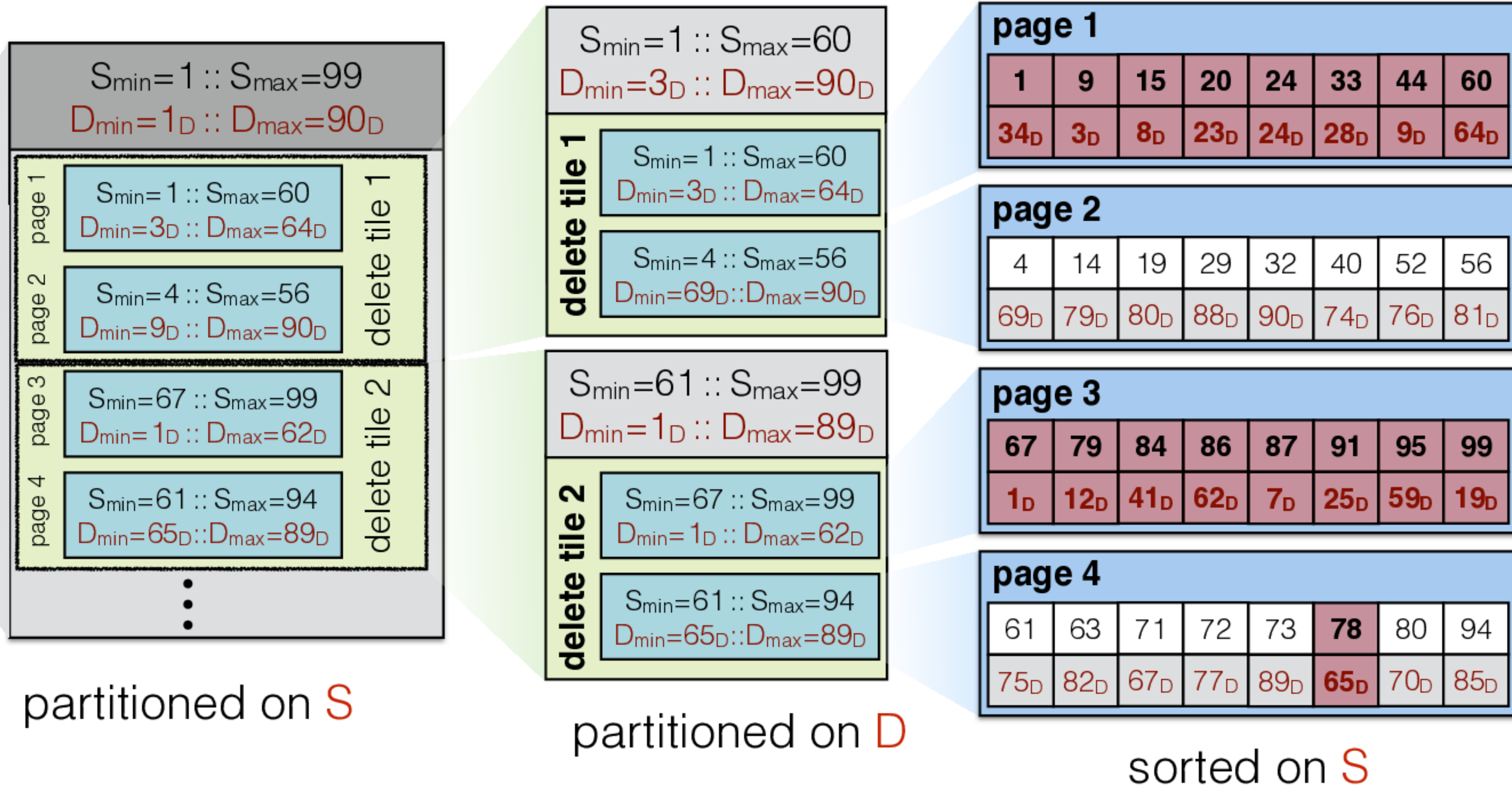
drop
page

Key Weaving storage layout

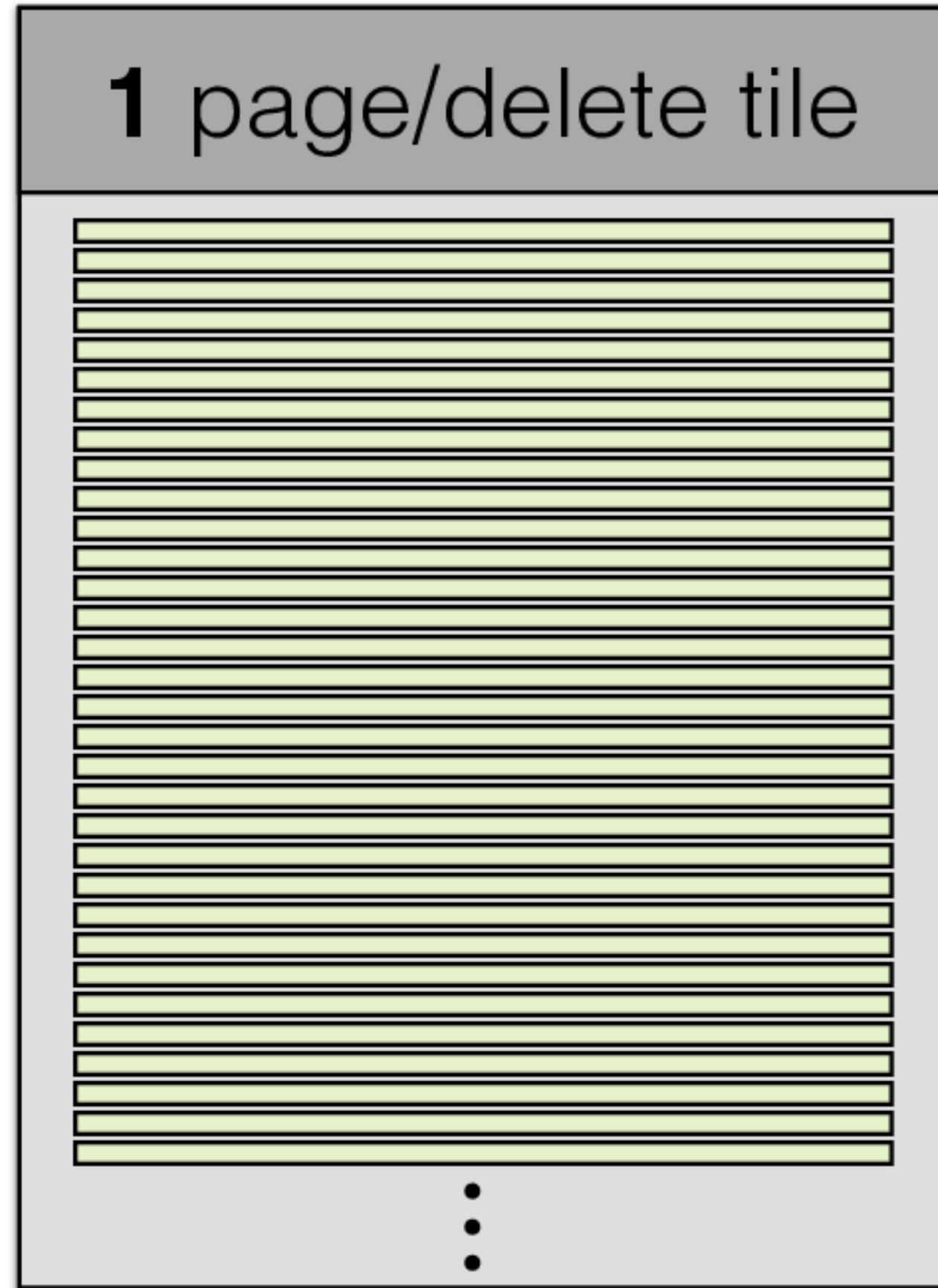
delete all entries with timestamp $\leq 65_D$



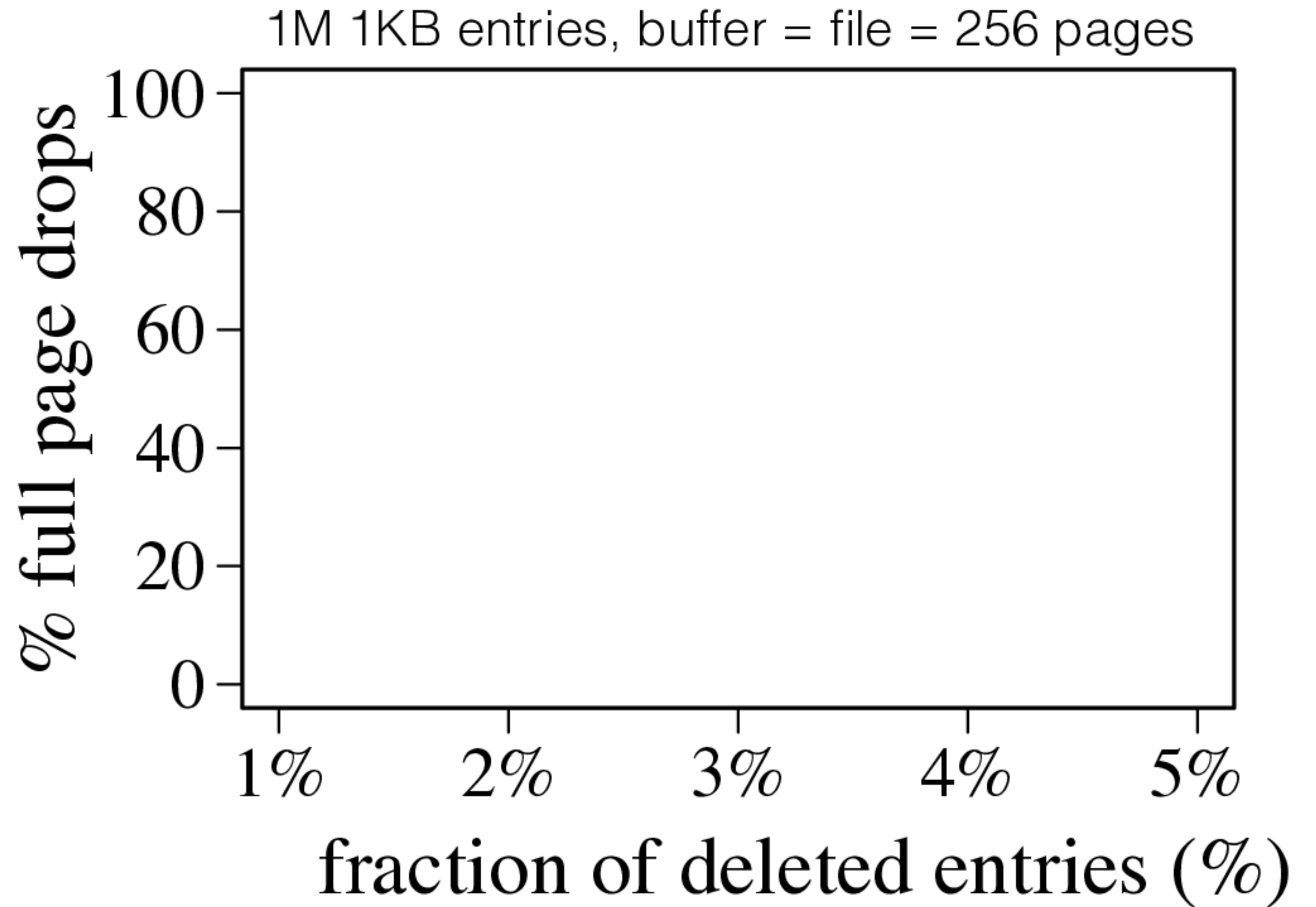
SST file



Key Weaving storage layout



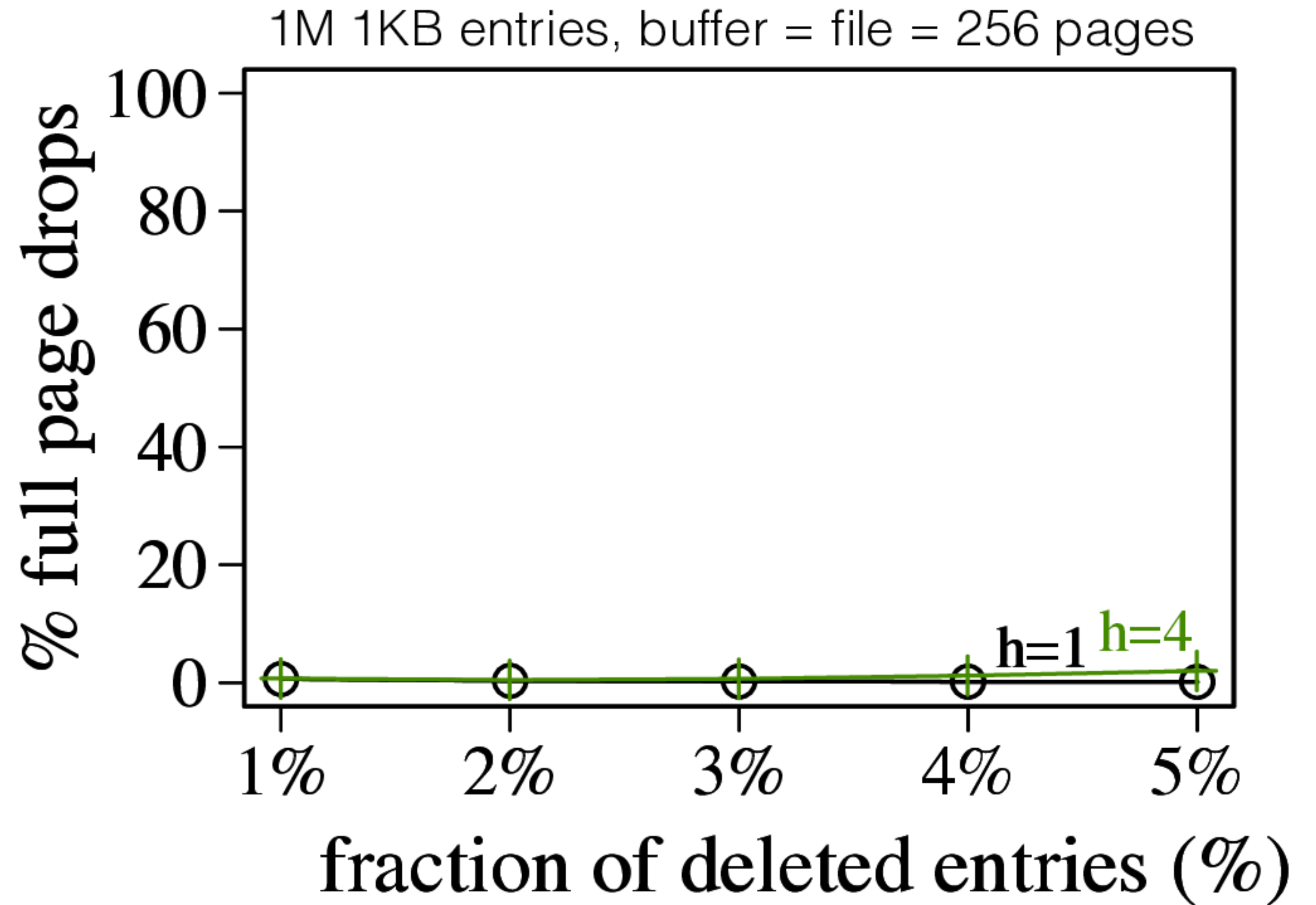
Internals of an SST file in **KiWi**



Key Weaving storage layout



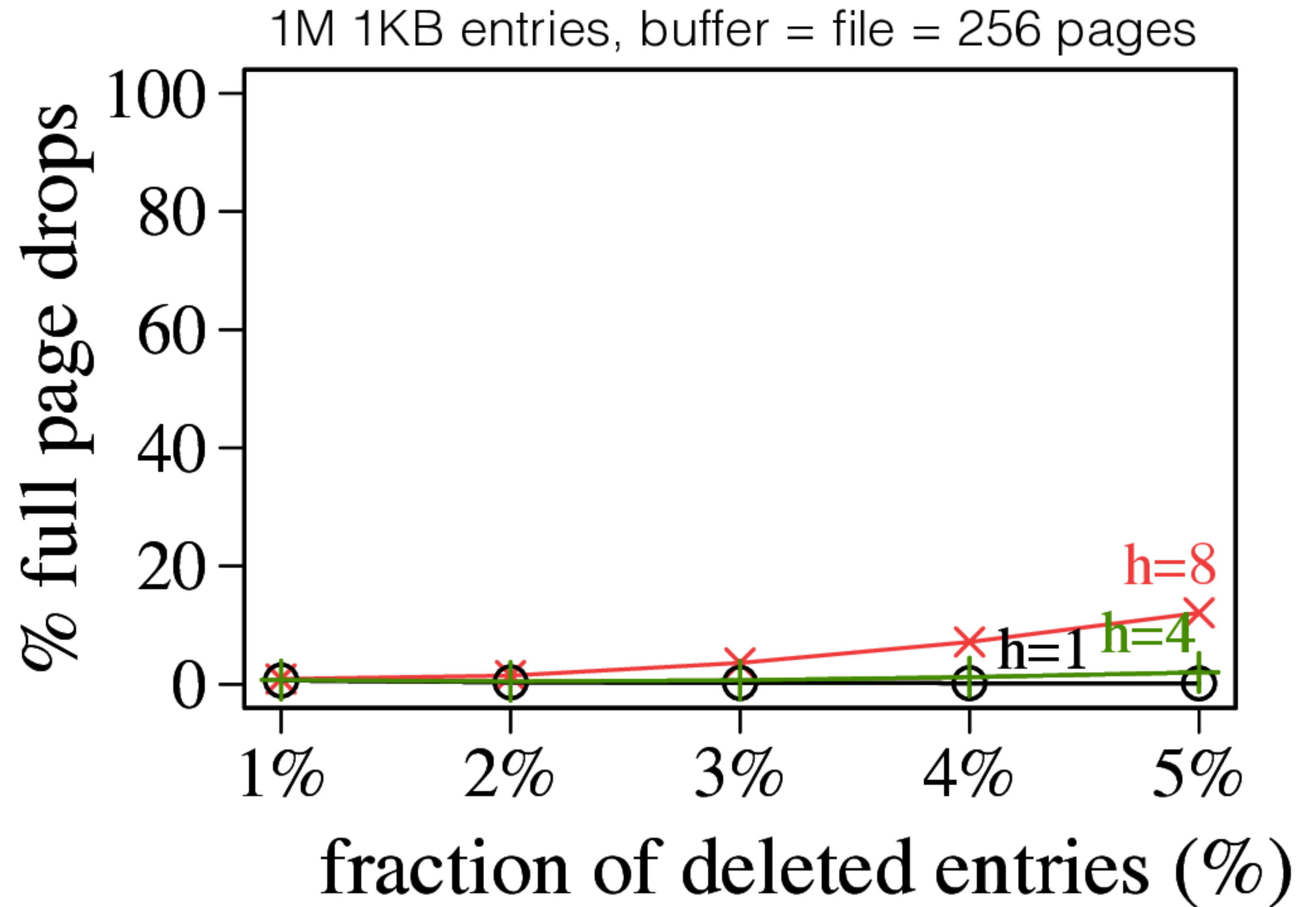
Internals of an SST file in **KiWi**



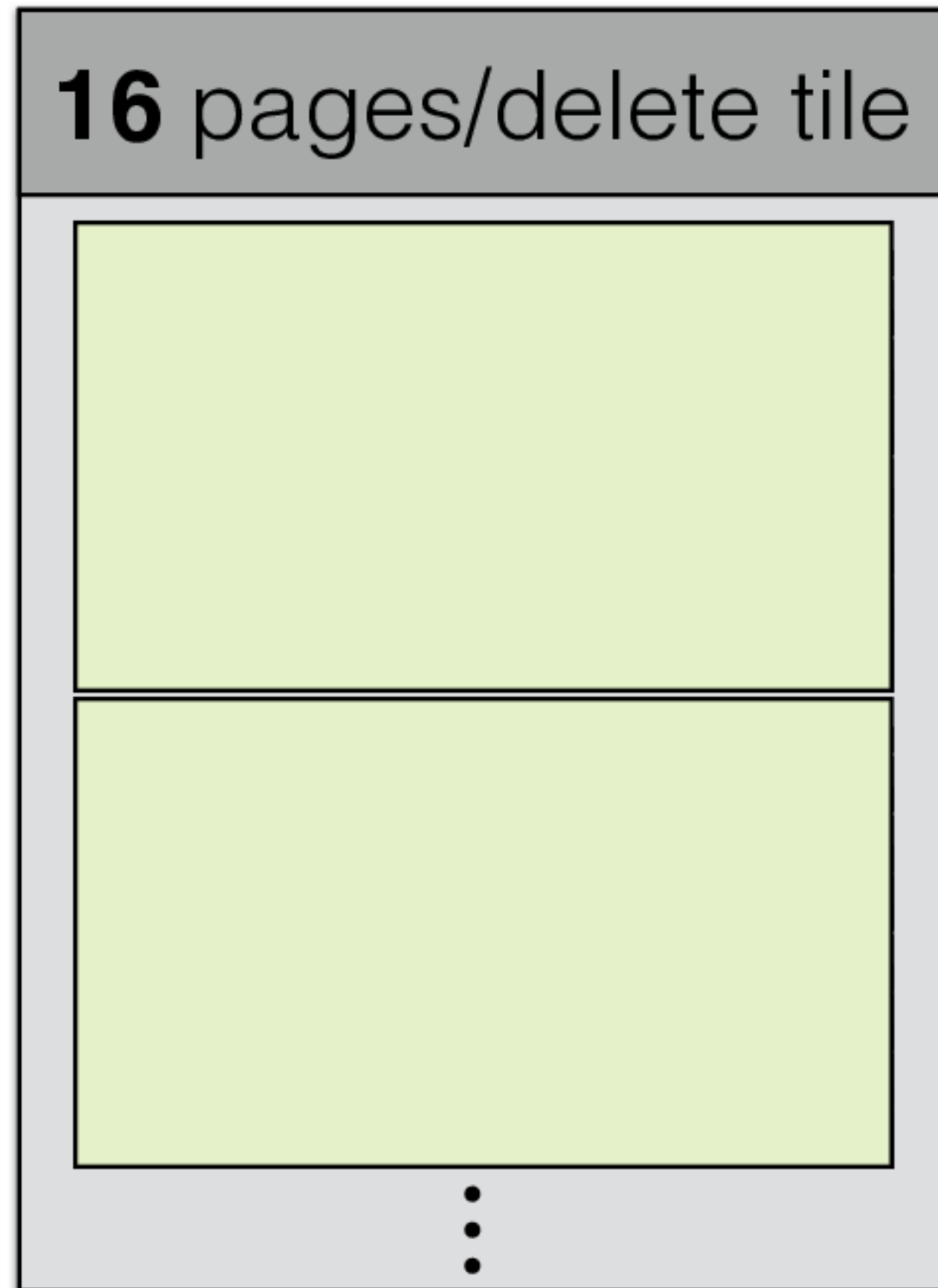
Key Weaving storage layout



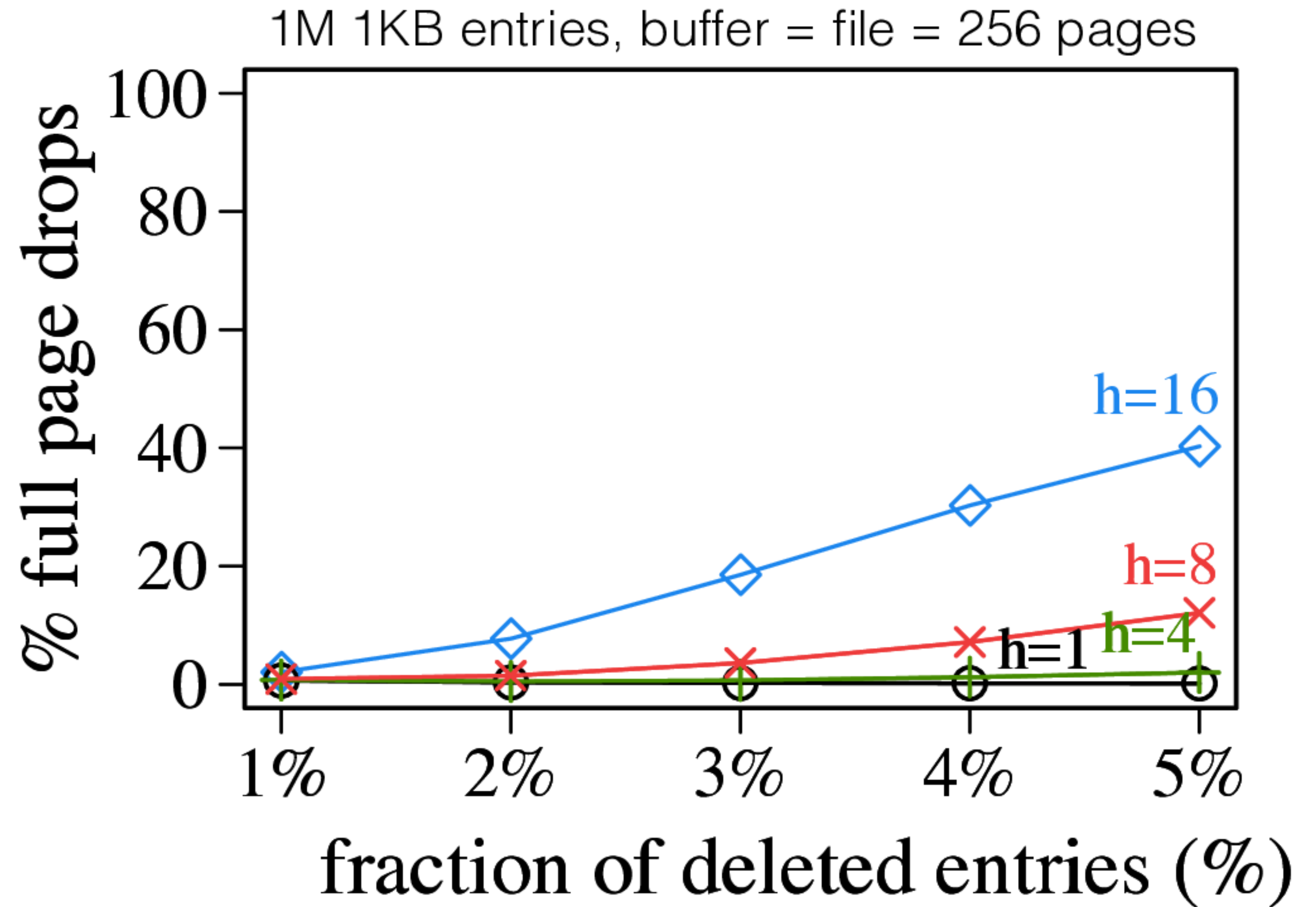
Internals of an SST file in **KiWi**



Key Weaving storage layout



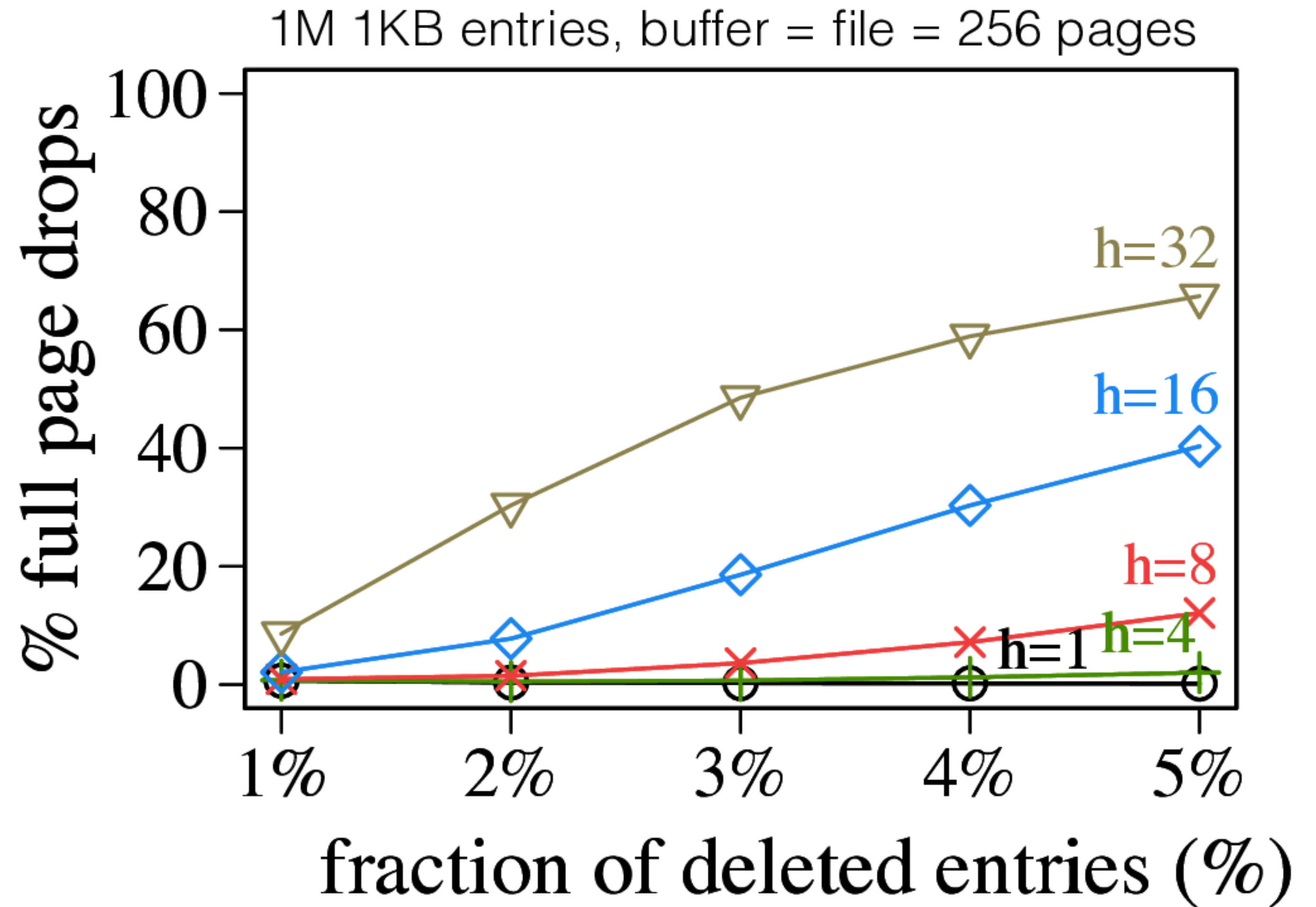
Internals of an SST file in **KiWi**



Key Weaving storage layout



Internals of an SST file in **KiWi**

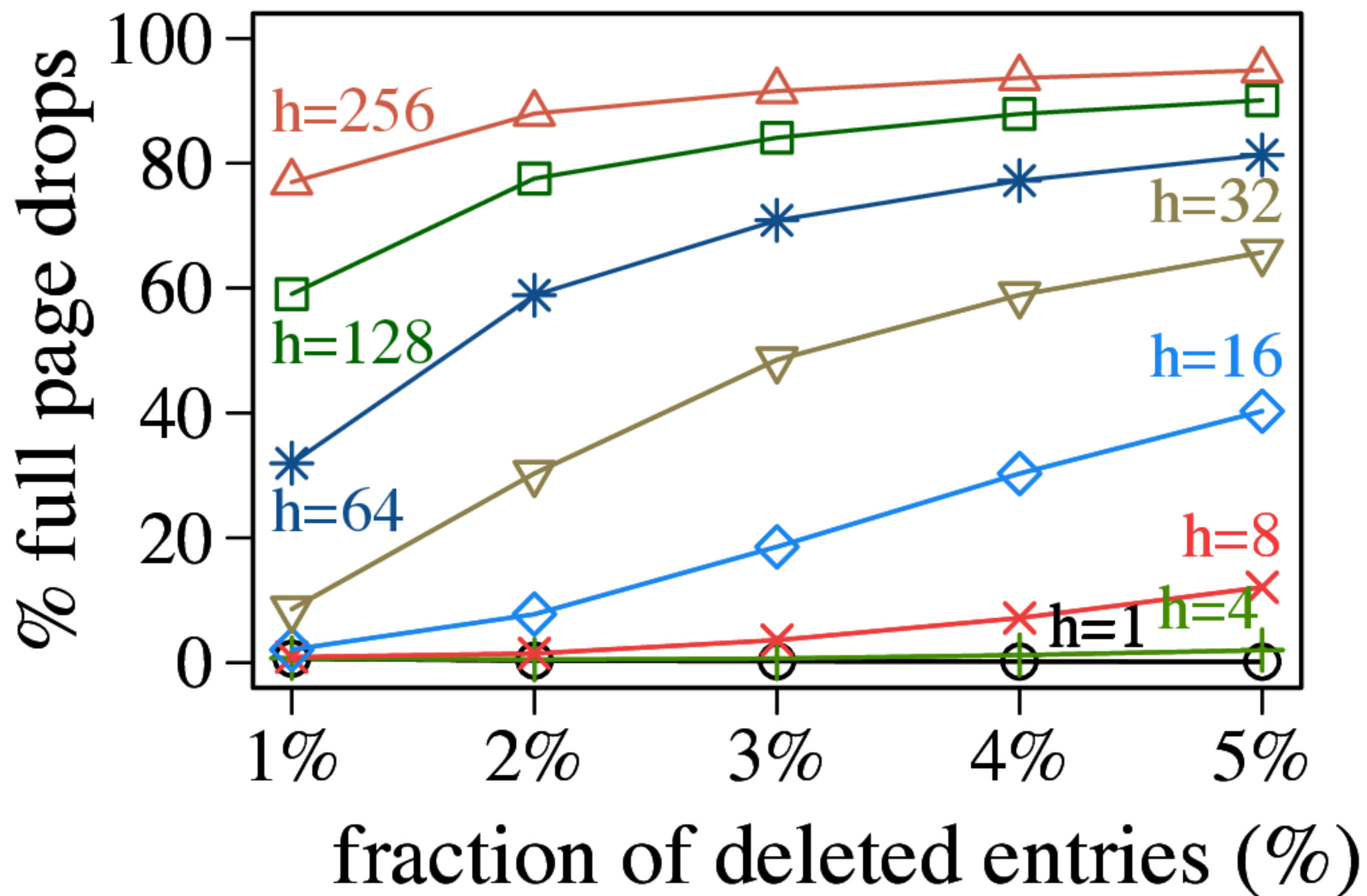


Key Weaving storage layout

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

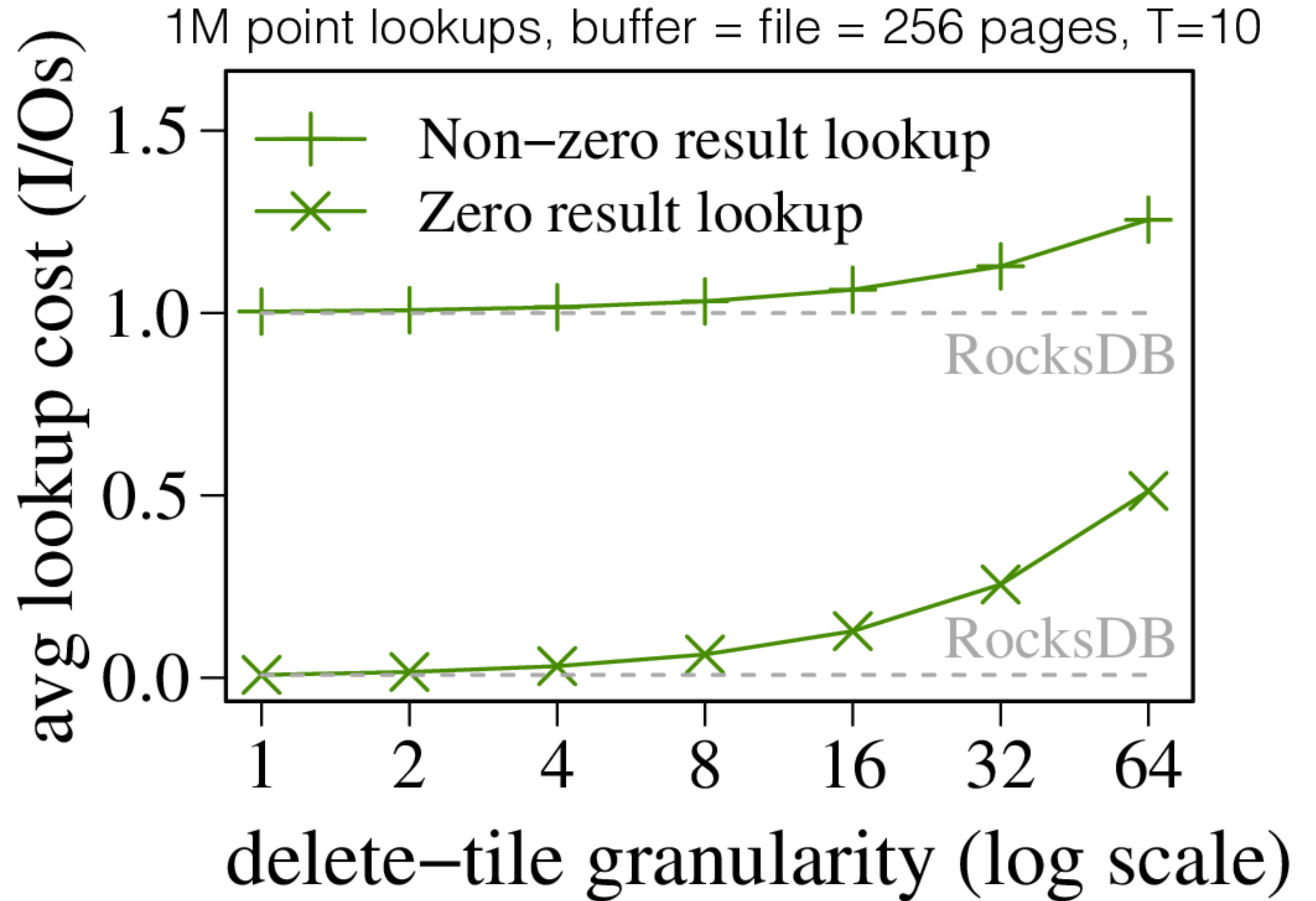
reduced latency spikes ✓

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



Lethe



delete tile size



lookup
cost

secondary range
delete cost





lookup
cost

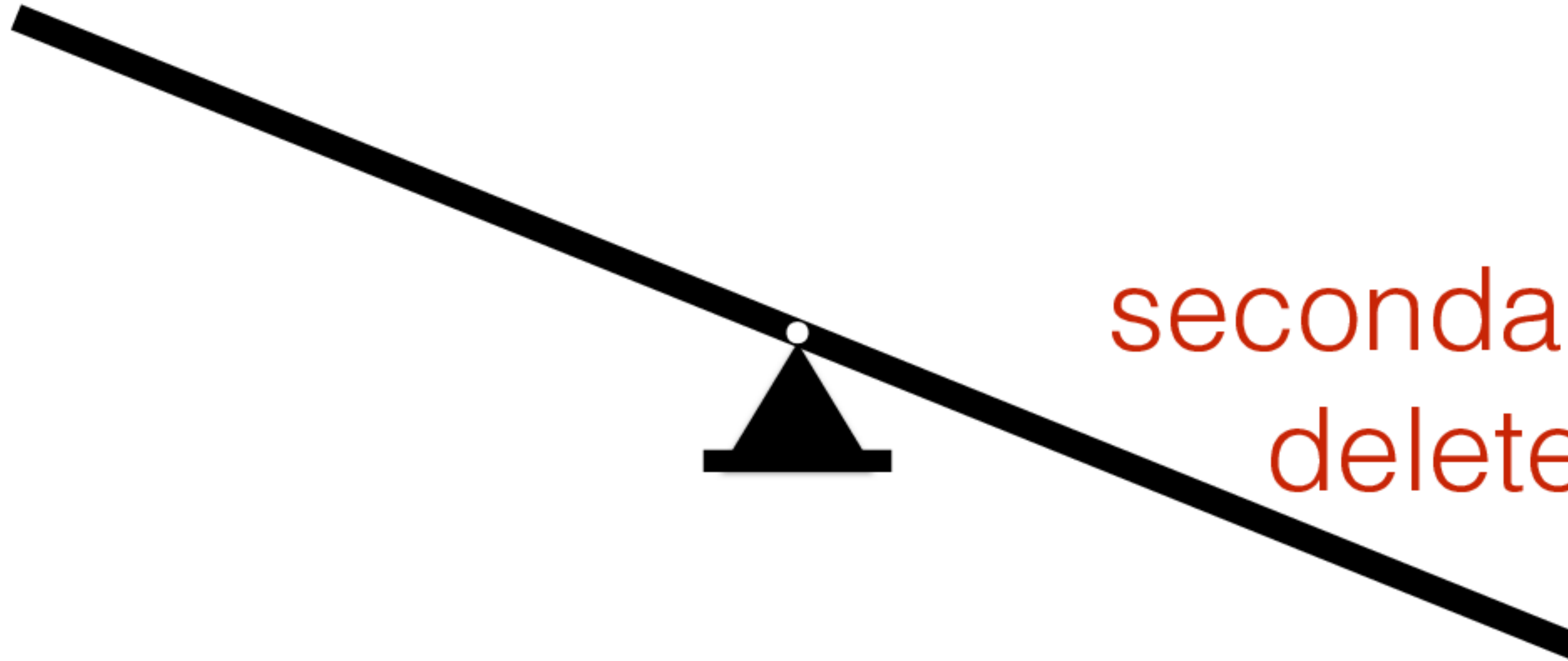
secondary range
delete cost



delete tile size



lookup
cost

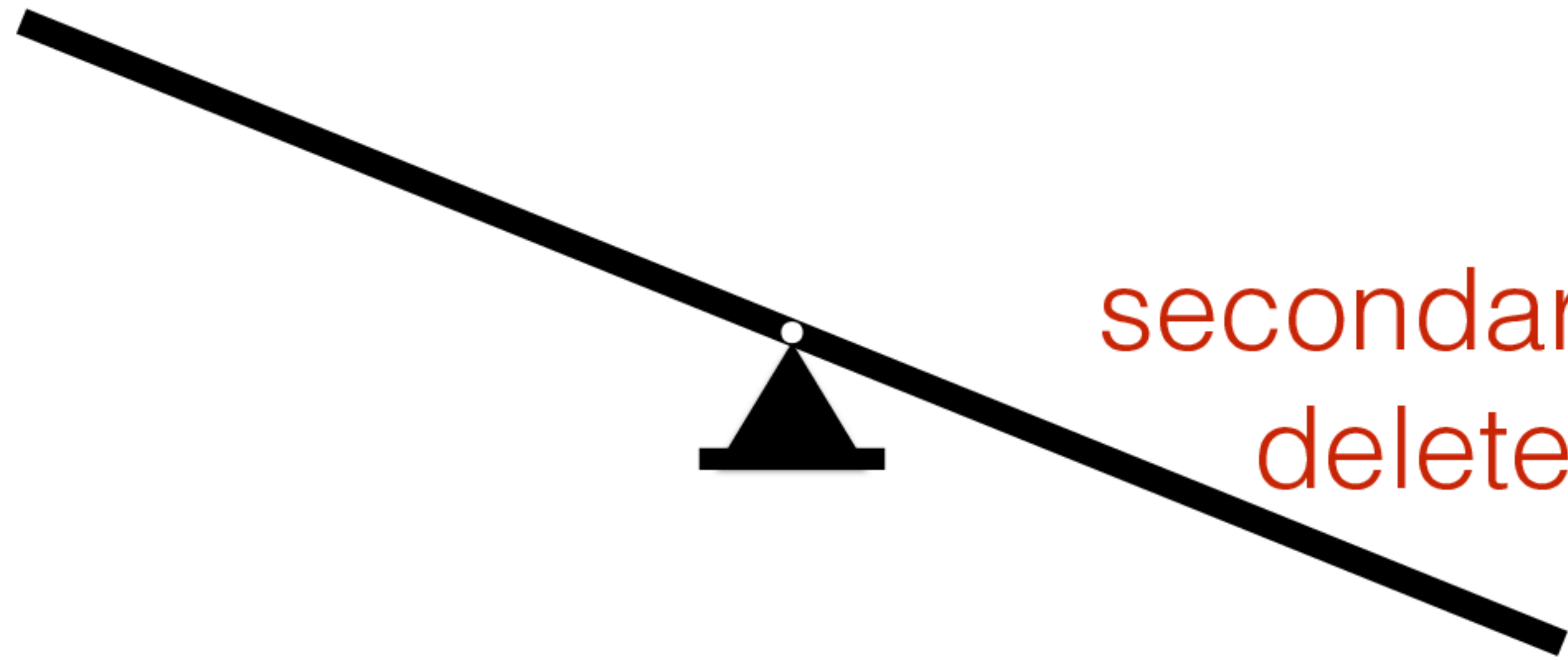


secondary range
delete cost

delete tile size



lookup
cost



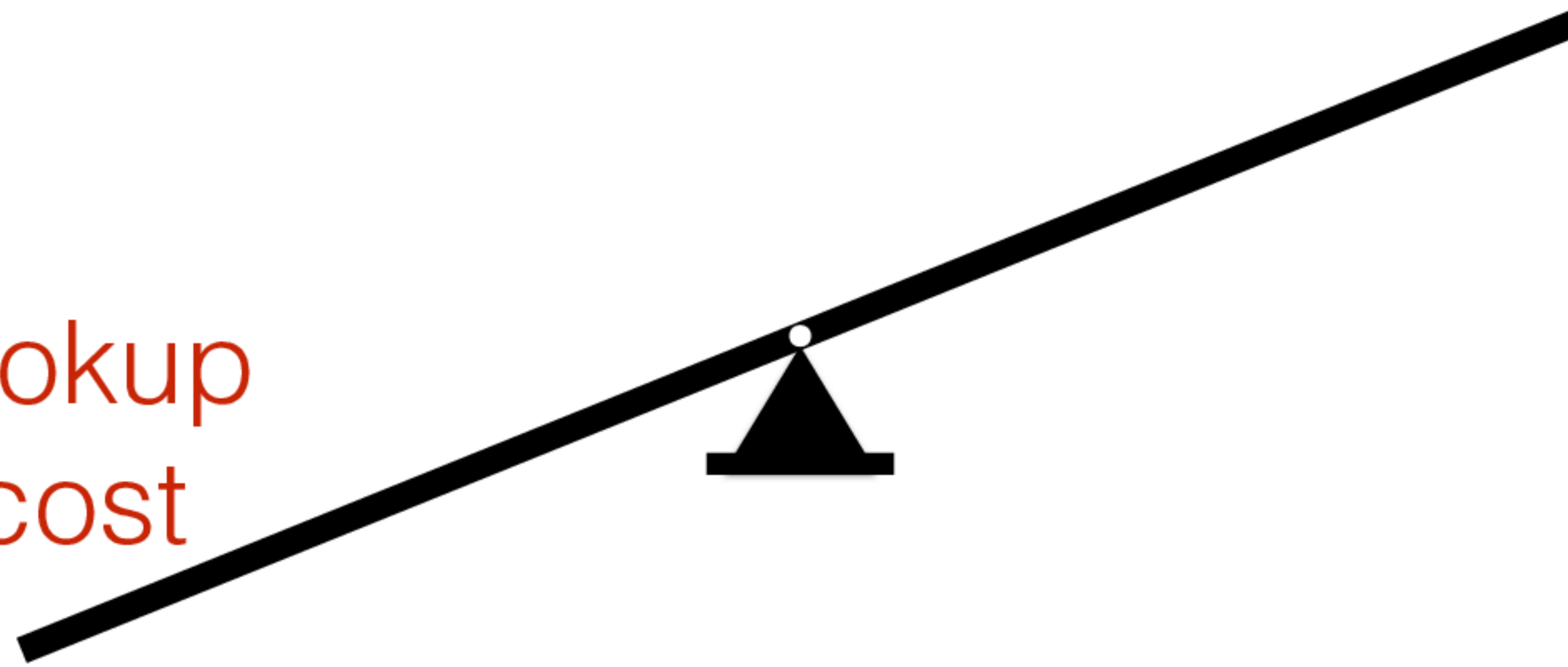
secondary range
delete cost

delete tile size



secondary range
delete cost

lookup
cost

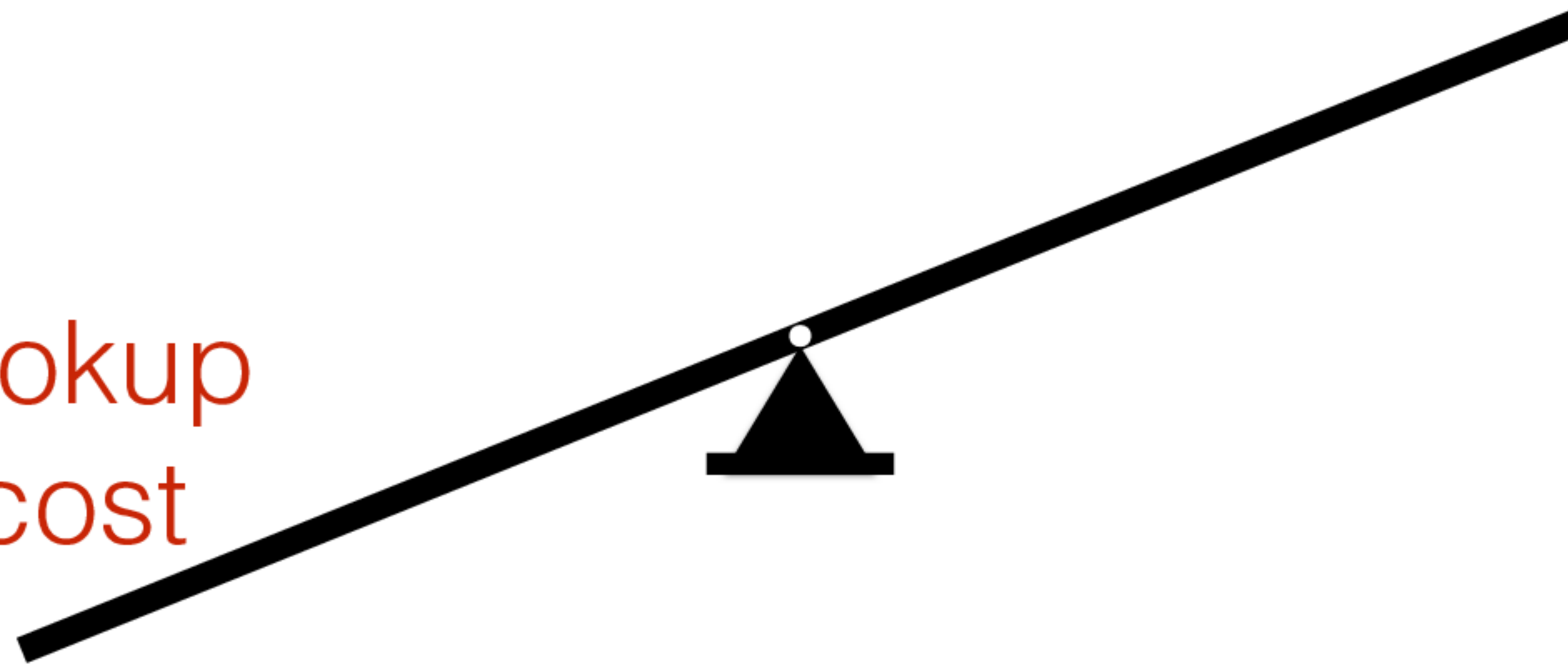


delete tile size



secondary range
delete cost

lookup
cost

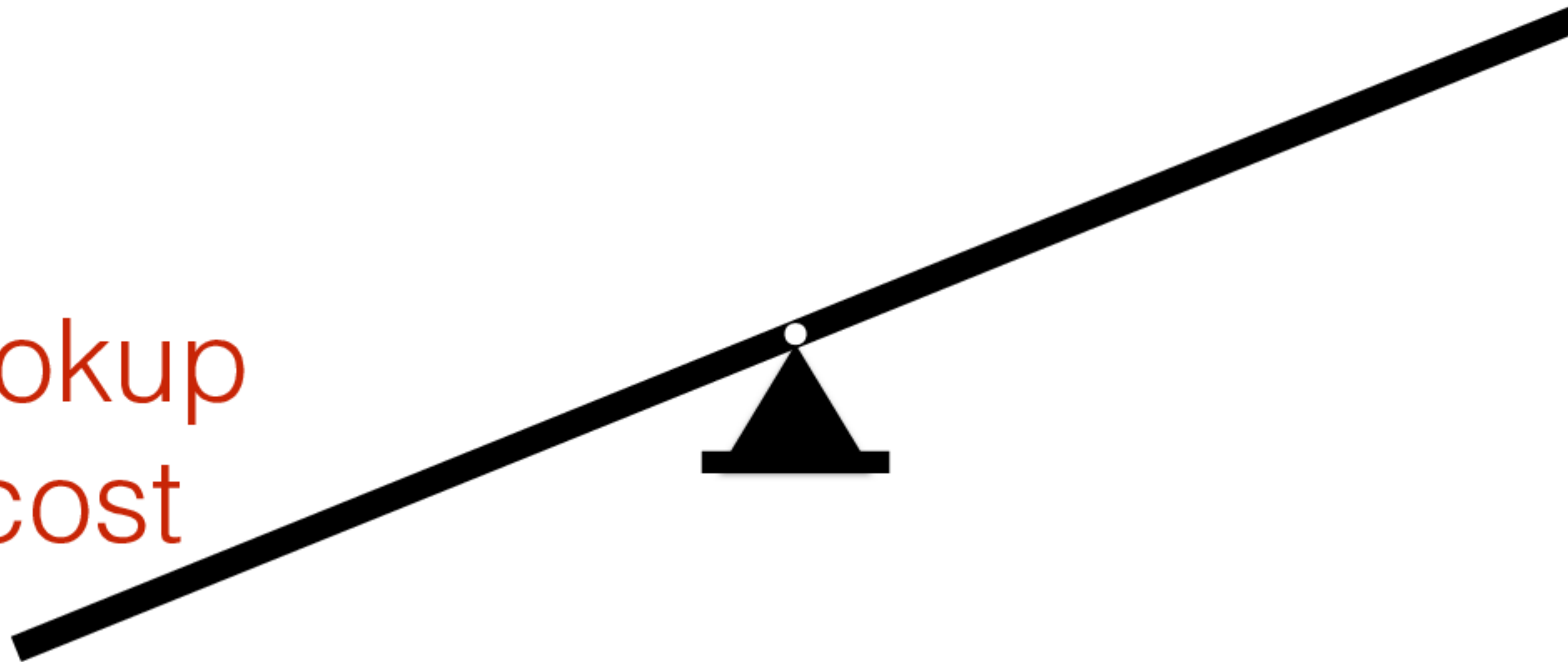


delete tile size



secondary range
delete cost

lookup
cost

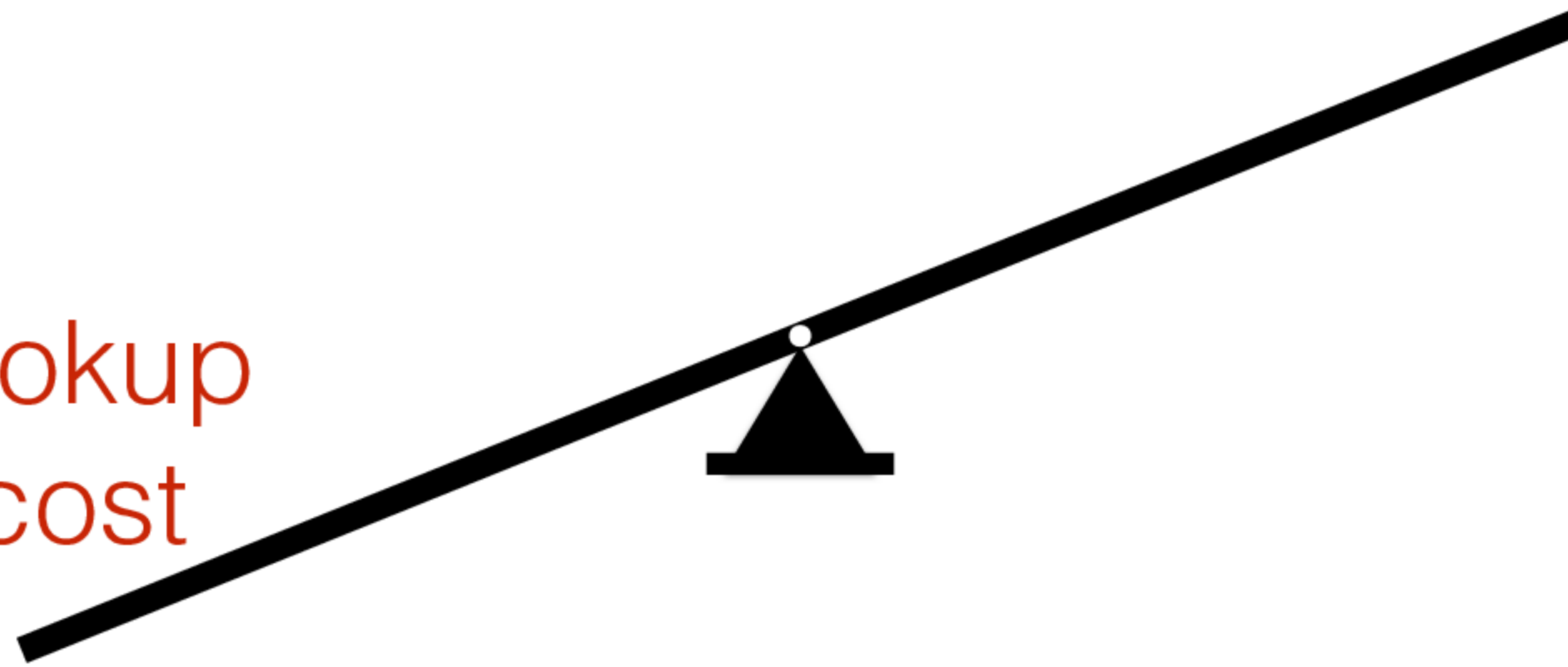


delete tile size



secondary range
delete cost

lookup
cost



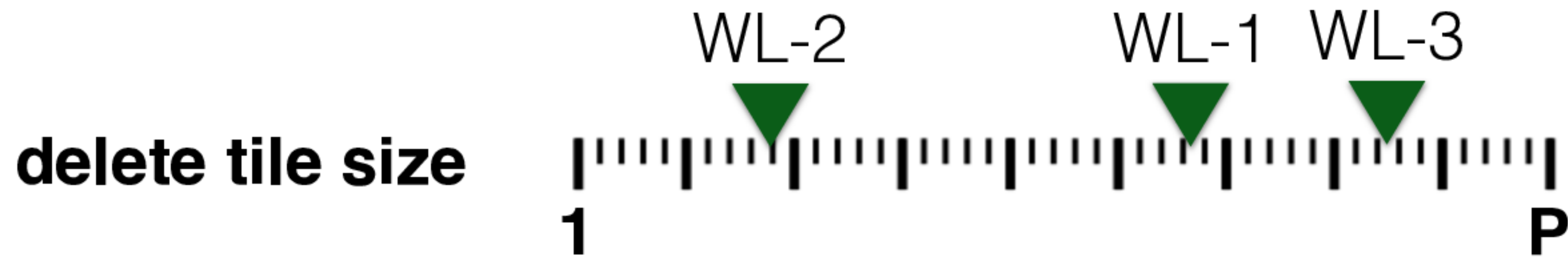
delete tile size



lookup
cost

secondary range
delete cost





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
using key-interweaved data storage



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

Thank You!

disc-projects.bu.edu/lethe/



Lethe strikes balance between
cost, performance, and latency