

COSI 167A

Advanced Data Systems

Class 12

Efficient Deletes in LSM-Engines

Prof. Subhadeep Sarkar

Class **logistics**

and administrivia

The **mid-semester project report** is due in **3 weeks** (**Nov 8**).

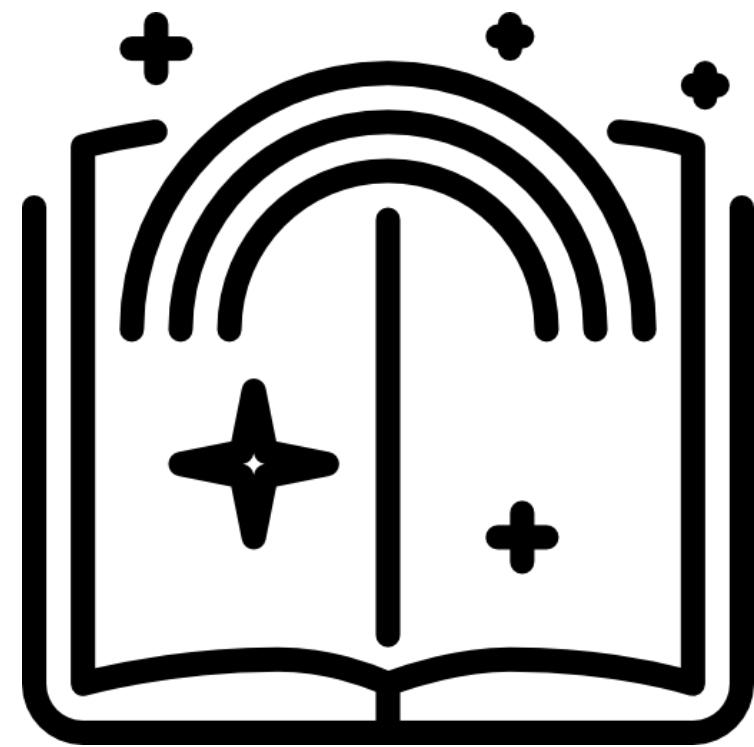
Start working on the project and use the **office hours**.

First student presentation on **Friday, Oct 18**.

Discuss the slides with me **a week before the presentation**.

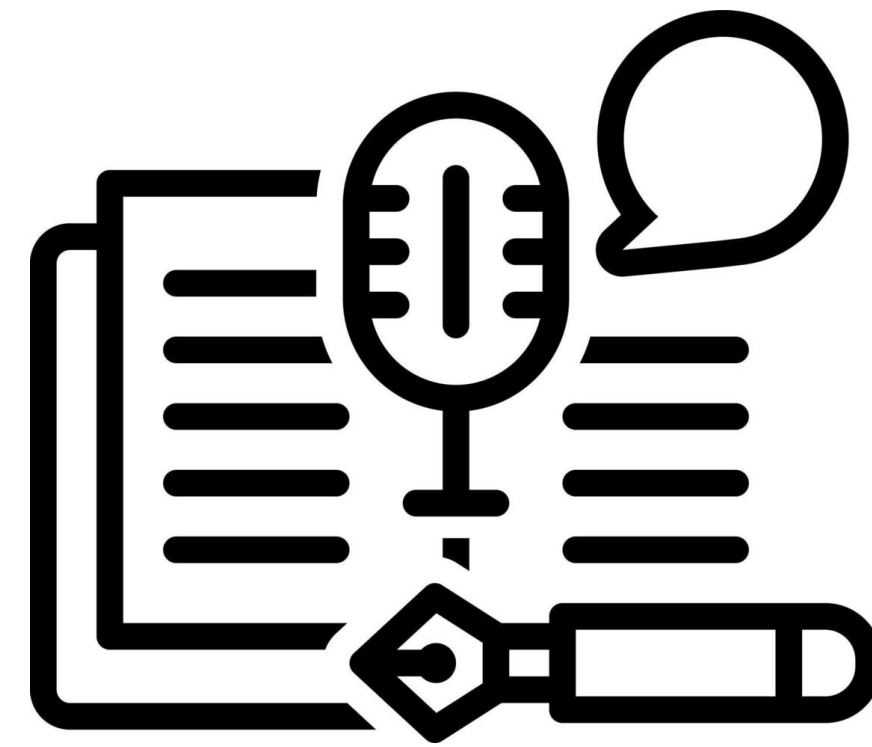
A good presentation

The three key ingredients



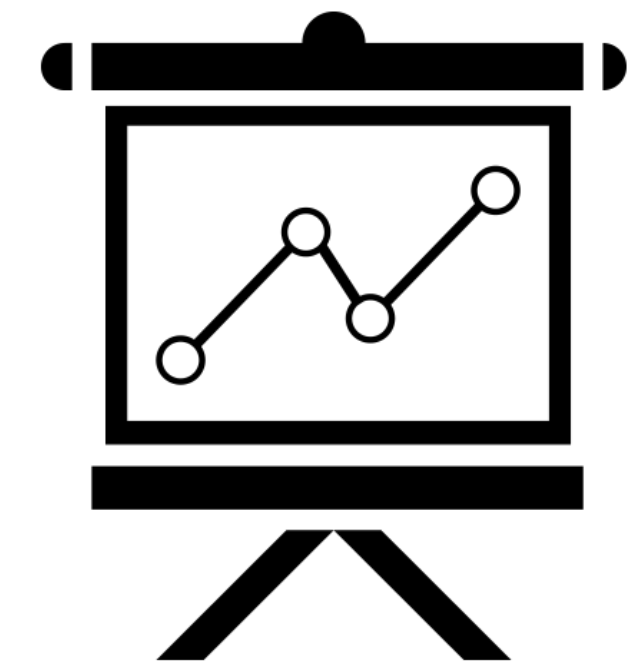
tell a **story**

highlight the problem
why it is worth solving
solution intuition & methodology
constructive criticism



prepare a **speech**

prepare a narrative
make sure it flows
rehearse!



presentation

make it engaging!

An engaging presentation

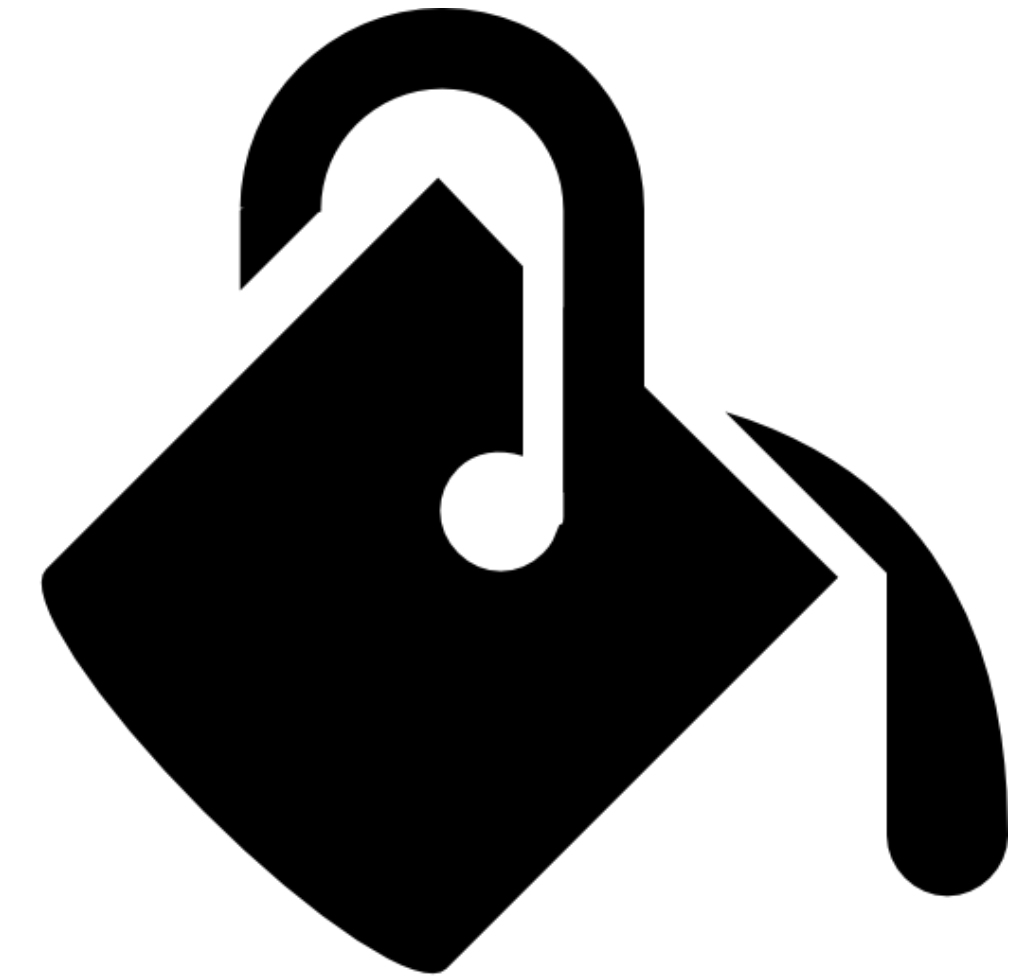
Think of what keeps you engaged!



don't use
bullets



1 message
per slide



1/2 colors

Today in COSI 167A

What's on the cards?

deletes in LSM

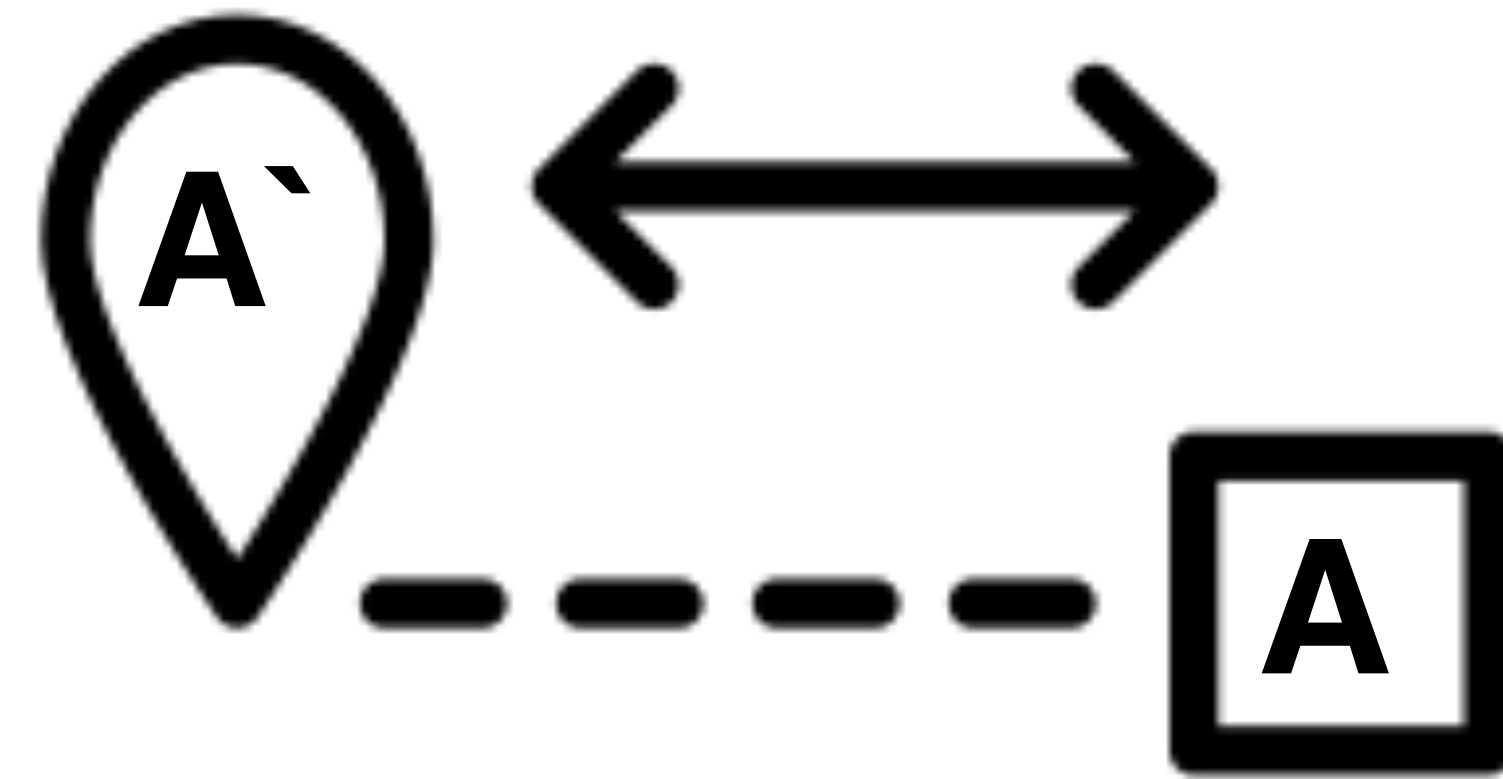
the problems with **out-of-place deletes**

building **deletion-compliant data systems**

Ingestion-Optimized Systems

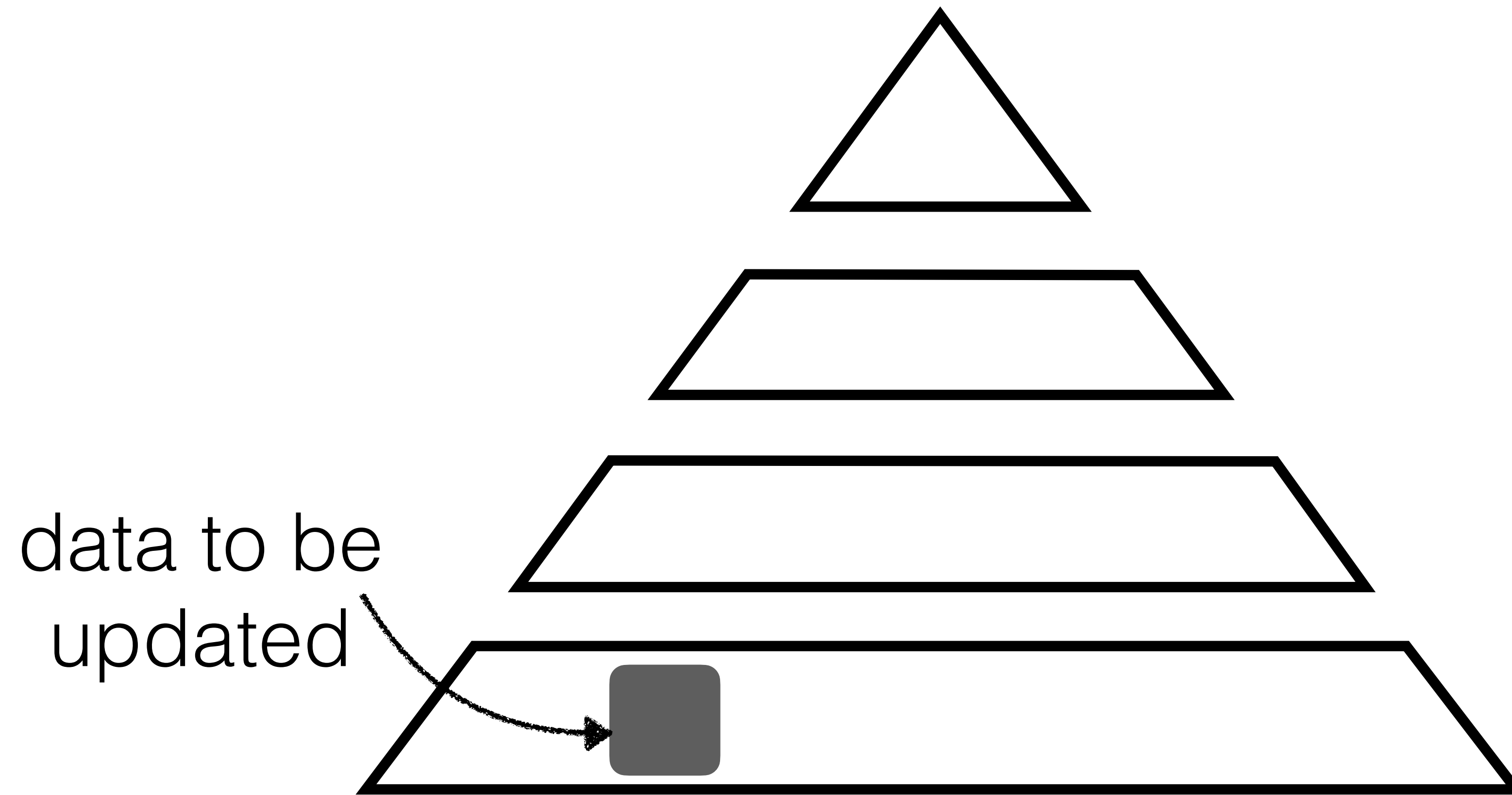


batched inserts

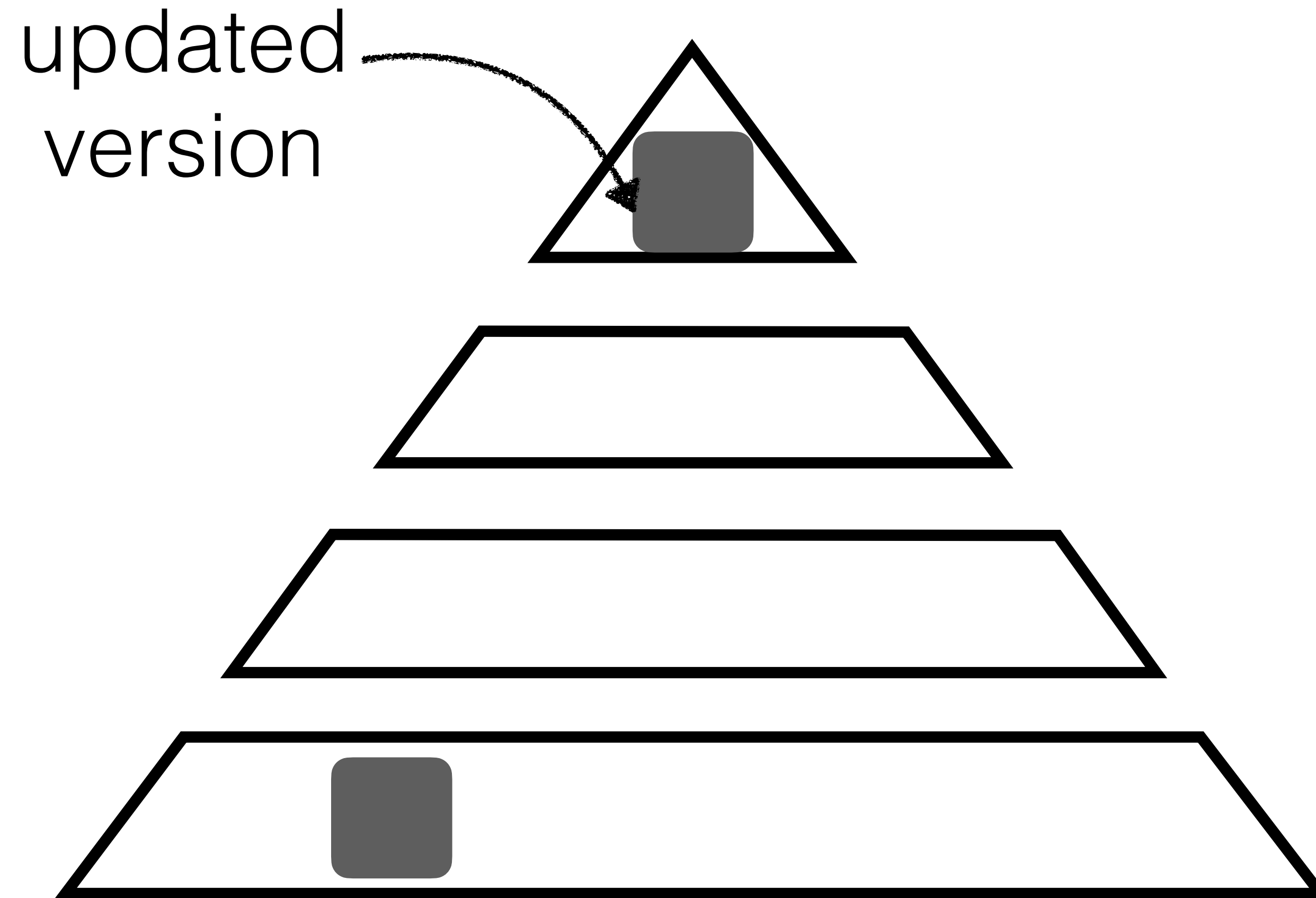


out-of-place
deletes/updates

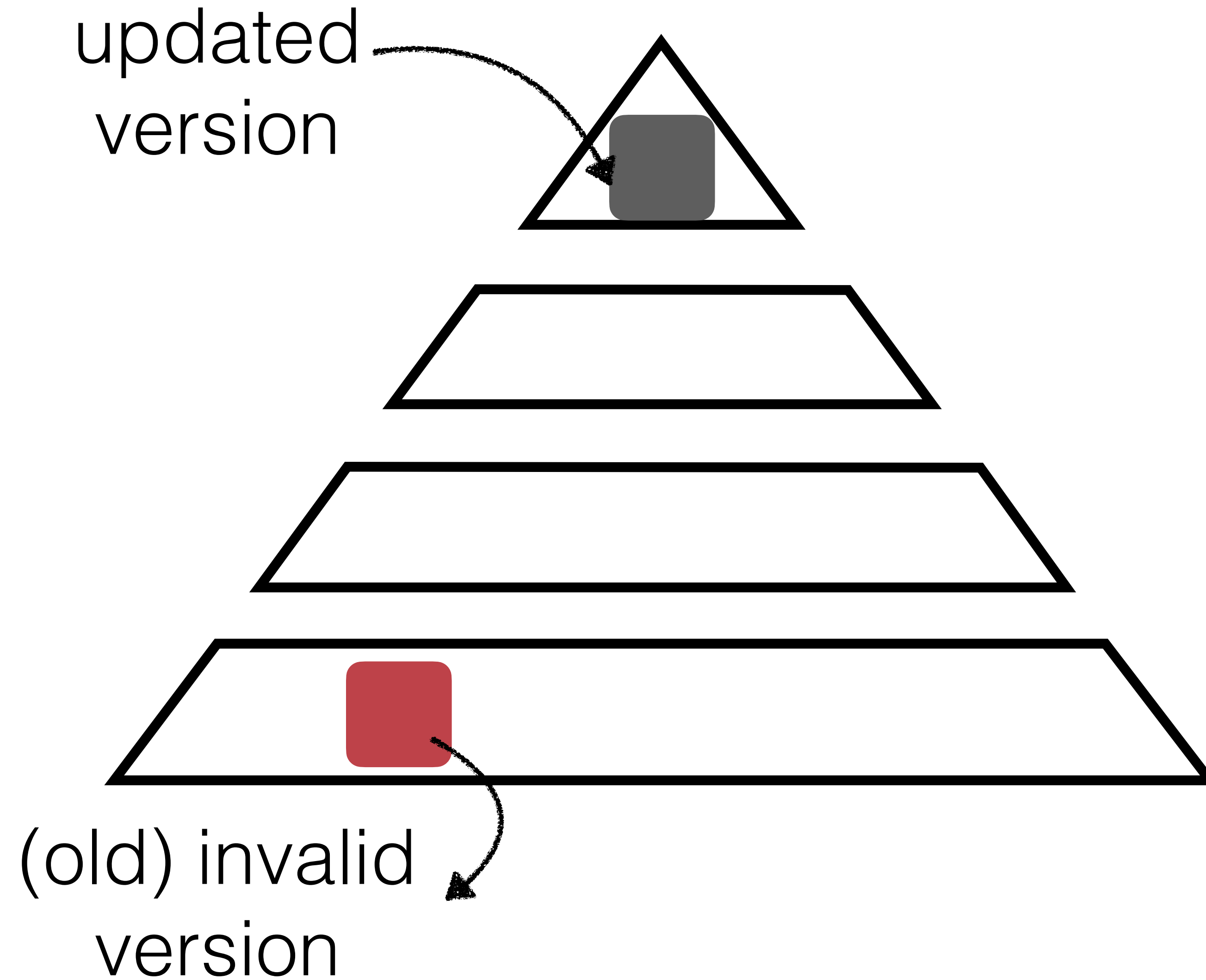
Out-of-place Deletes/Updates



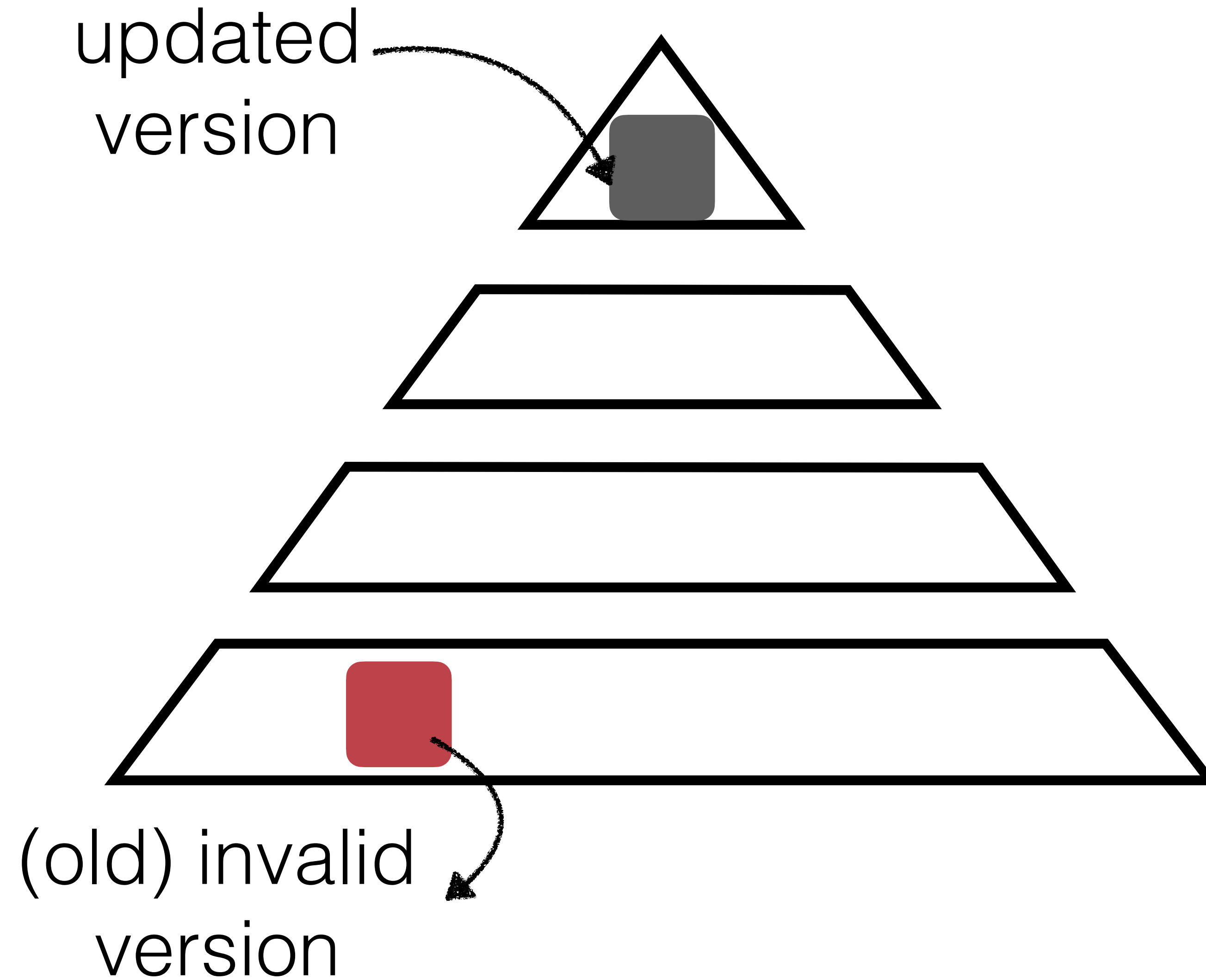
Out-of-place Deletes/Updates



Out-of-place Deletes/Updates

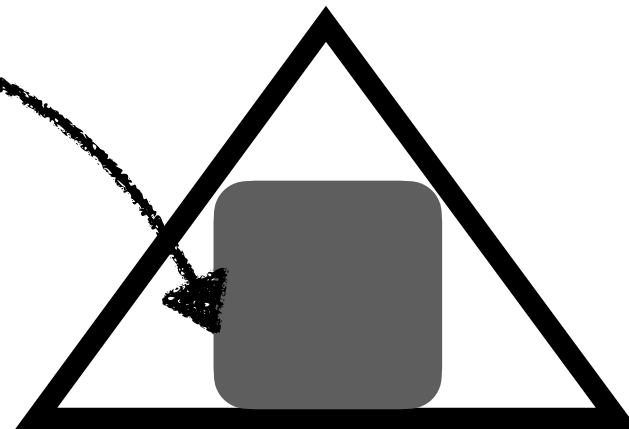


Out-of-place Deletes/Updates



Out-of-place Data Systems

updated
version



Google

∞ Meta

amazon



THE
APACHE[®]
SOFTWARE FOUNDATION



mongoDB

\$22B. growing at the

Hidden Cost: **does not scale** with deletes!

Today's **talk**

Deletes in LSMs



Lethe: A Tunable Delete-Aware LSM Engine

Subhadeep Sarkar, Tarikul Islam Papon, Dimitris Staratzis, Manos Athanassoulis

Boston University

{ssarkar1, papon, dstara, mathan}@bu.edu

presented at **SIGMOD 2020**

Today's talk

Deletes in LSMs



Lethe: A Tunable Delete-Aware LSM Engine

Subhadeep Sarkar, Tarikul Islam Papon, Dimitris Staratzis, Manos Athanassoulis

Boston University

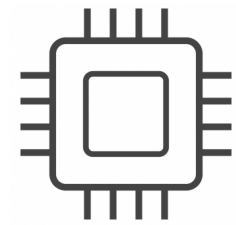
{ssarkar1, papon, dstara, mathan}@bu.edu

presented at **SIGMOD 2020**

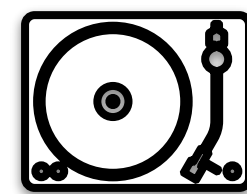


Lethe [\ 'lē-thē \]

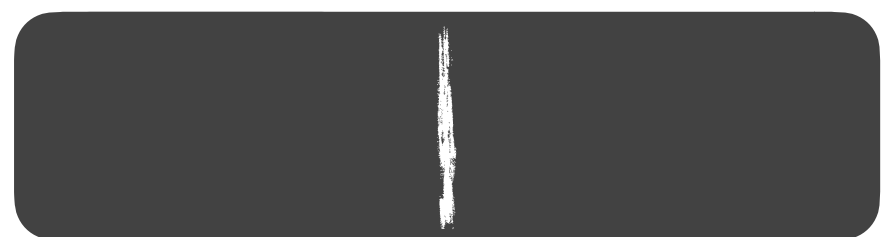
n: the goddess of forgetfulness



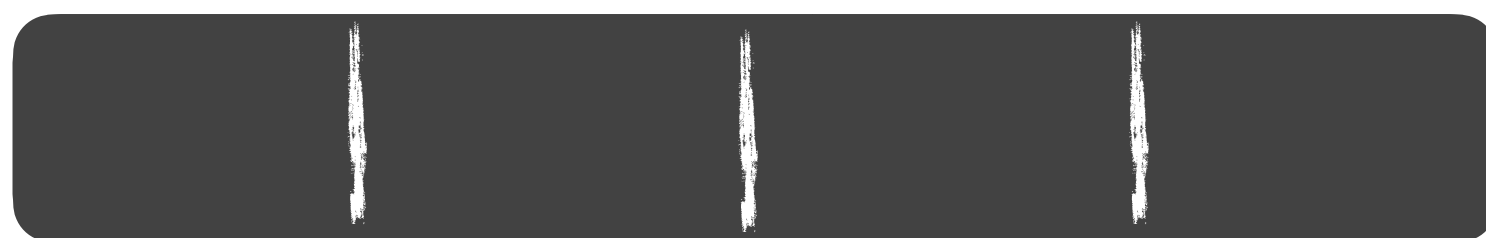
buffer



level 1



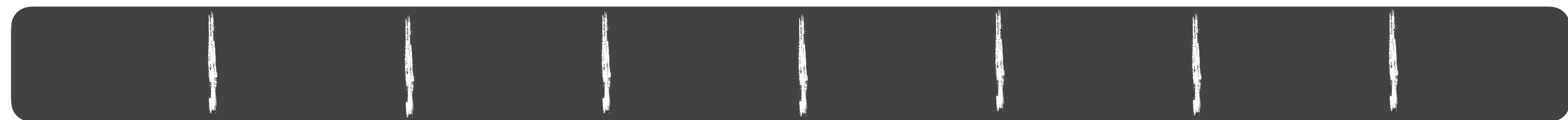
level 2



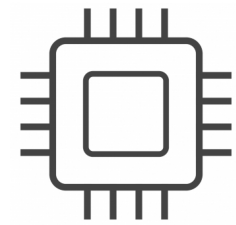
level 3



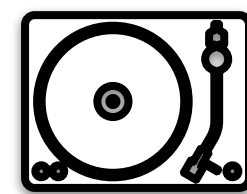
level 4



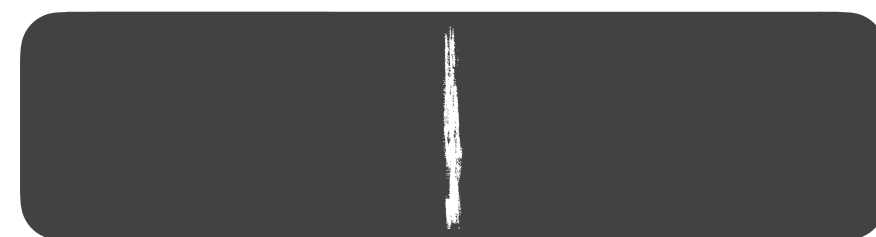
multiple files per
sorted run



buffer

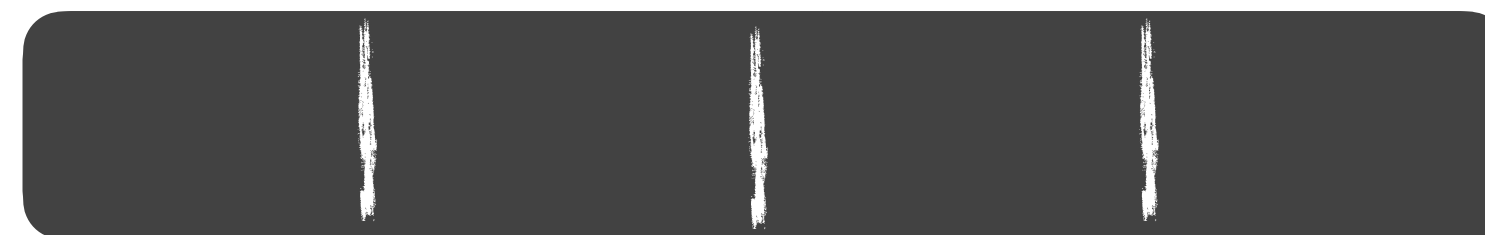


level 1



partial
compaction

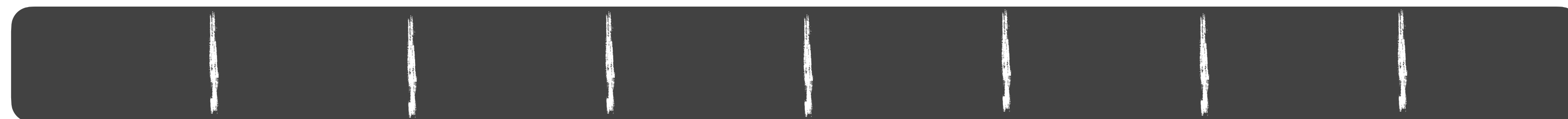
level 2

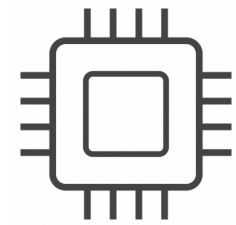


level 3

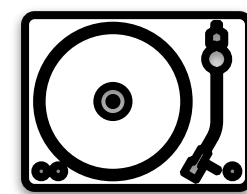


level 4

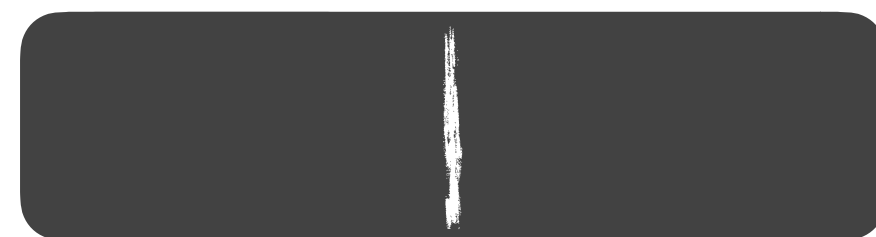




buffer

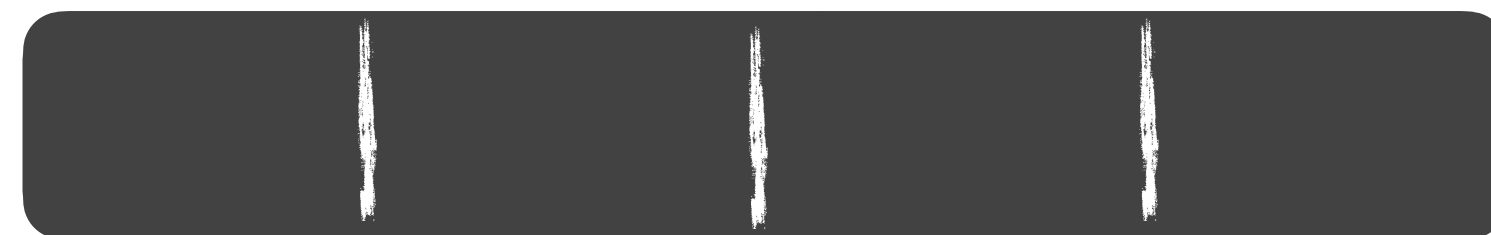


level 1



partial
compaction

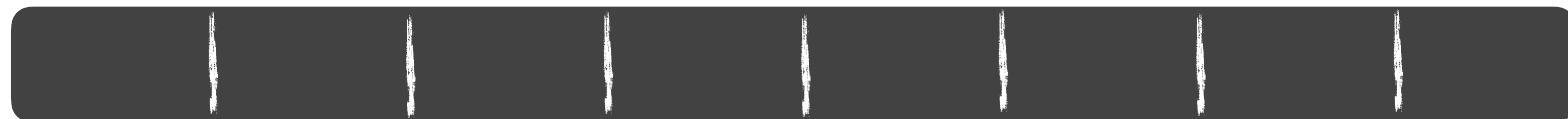
level 2

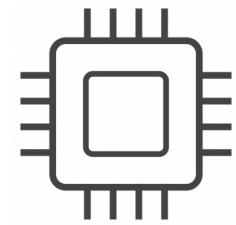


level 3

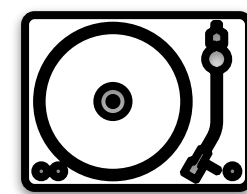


level 4





buffer



level 1



partial
compaction

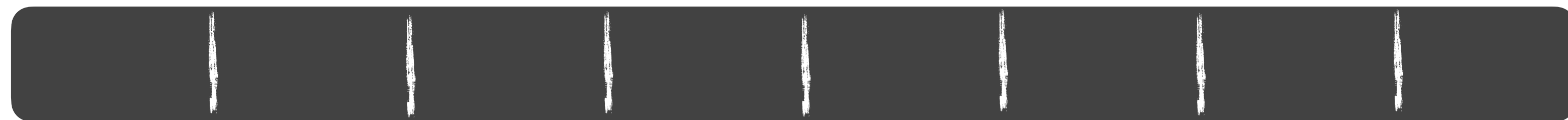
level 2

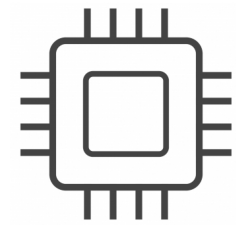


level 3

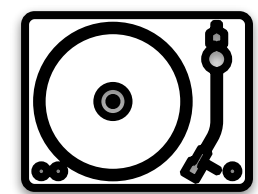


level 4

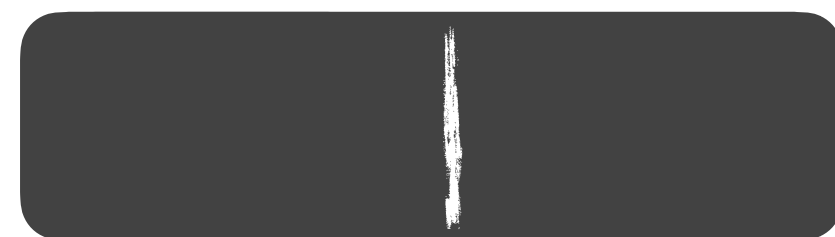




buffer



level 1



partial
compaction

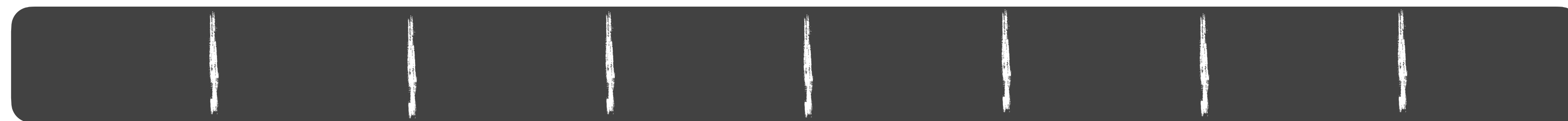
level 2

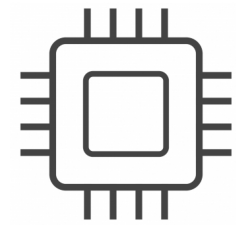


level 3

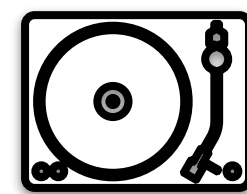


level 4

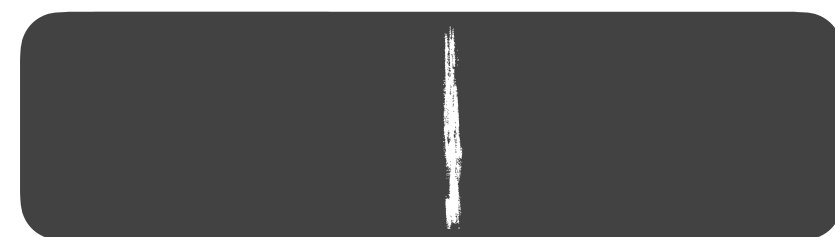




buffer



level 1



partial
compaction

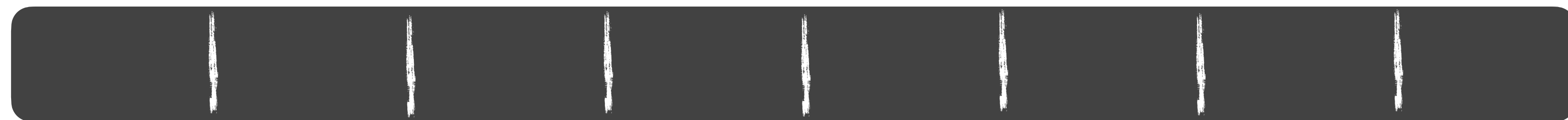
level 2

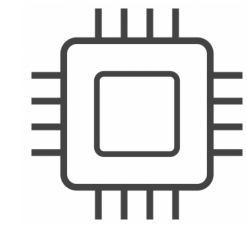


level 3

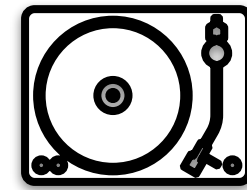


level 4

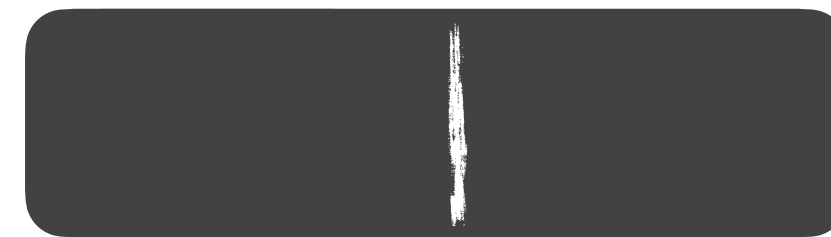




buffer

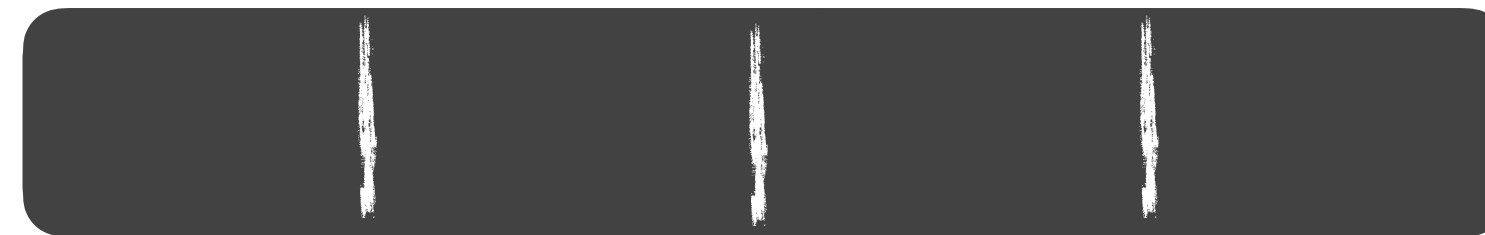


level 1

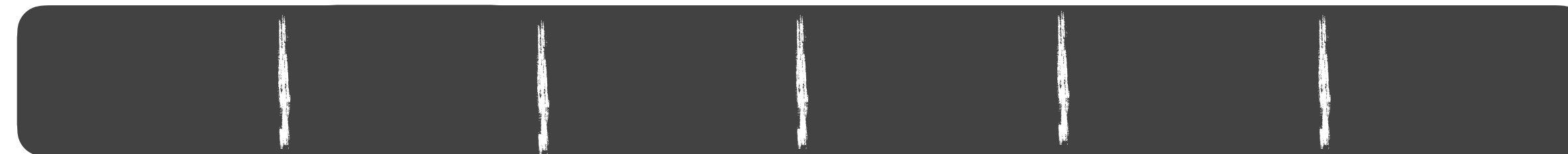


partial
compaction

level 2



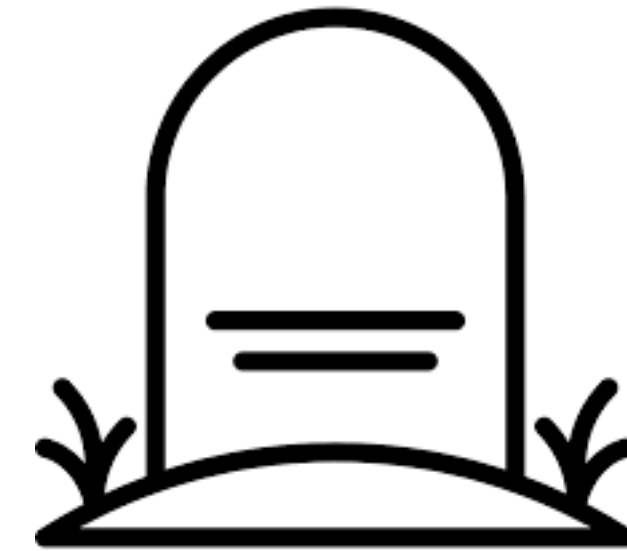
level 3



So, what about deletes?

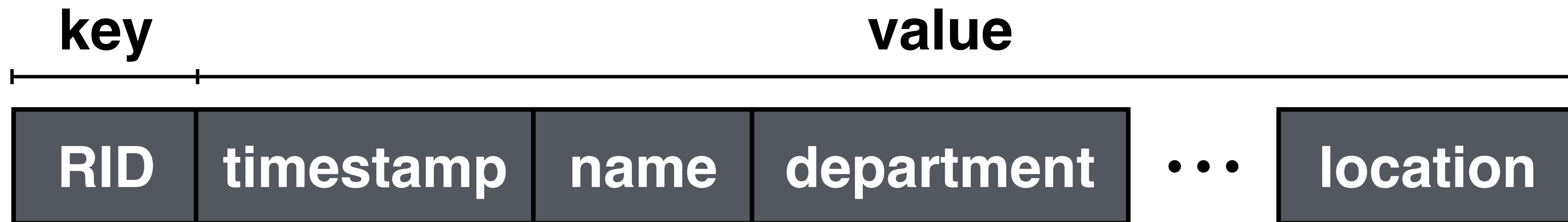
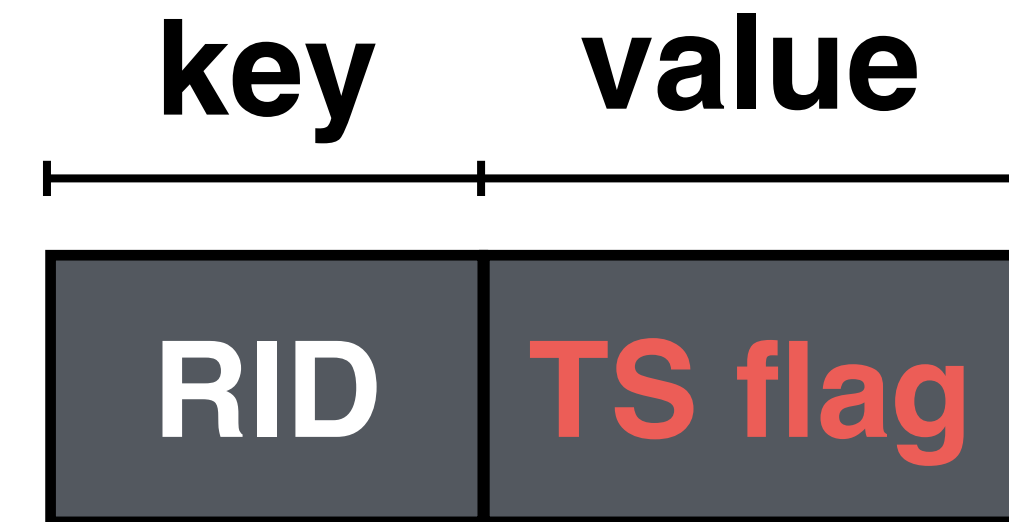
Deletes in LSMs

delete := insert tombstone



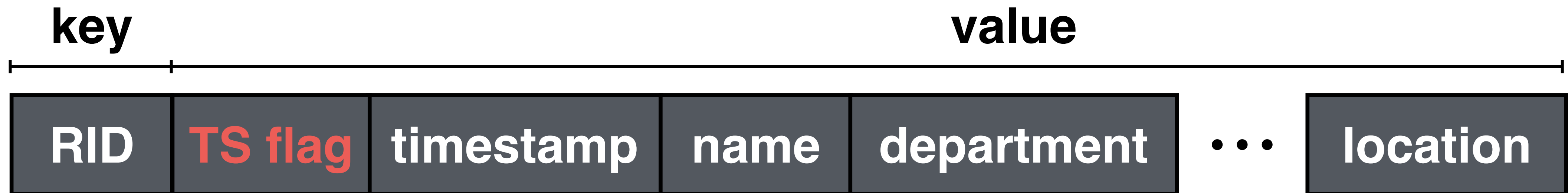
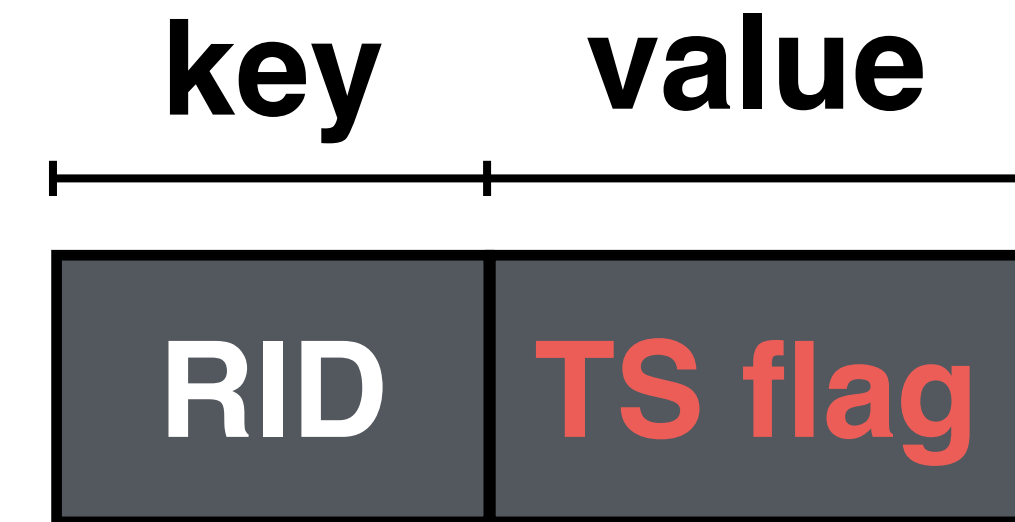
Deletes in LSMs

delete := insert tombstone



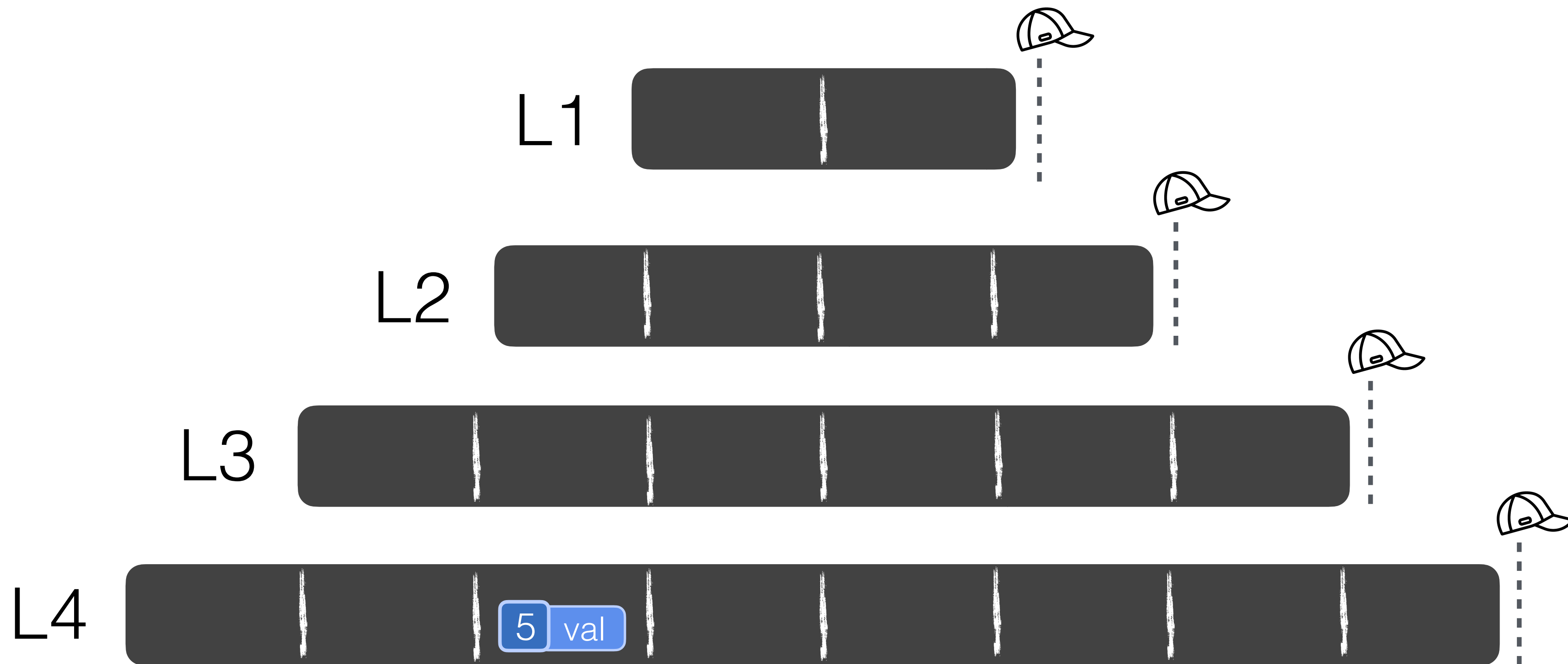
Deletes in LSMs

delete := insert tombstone



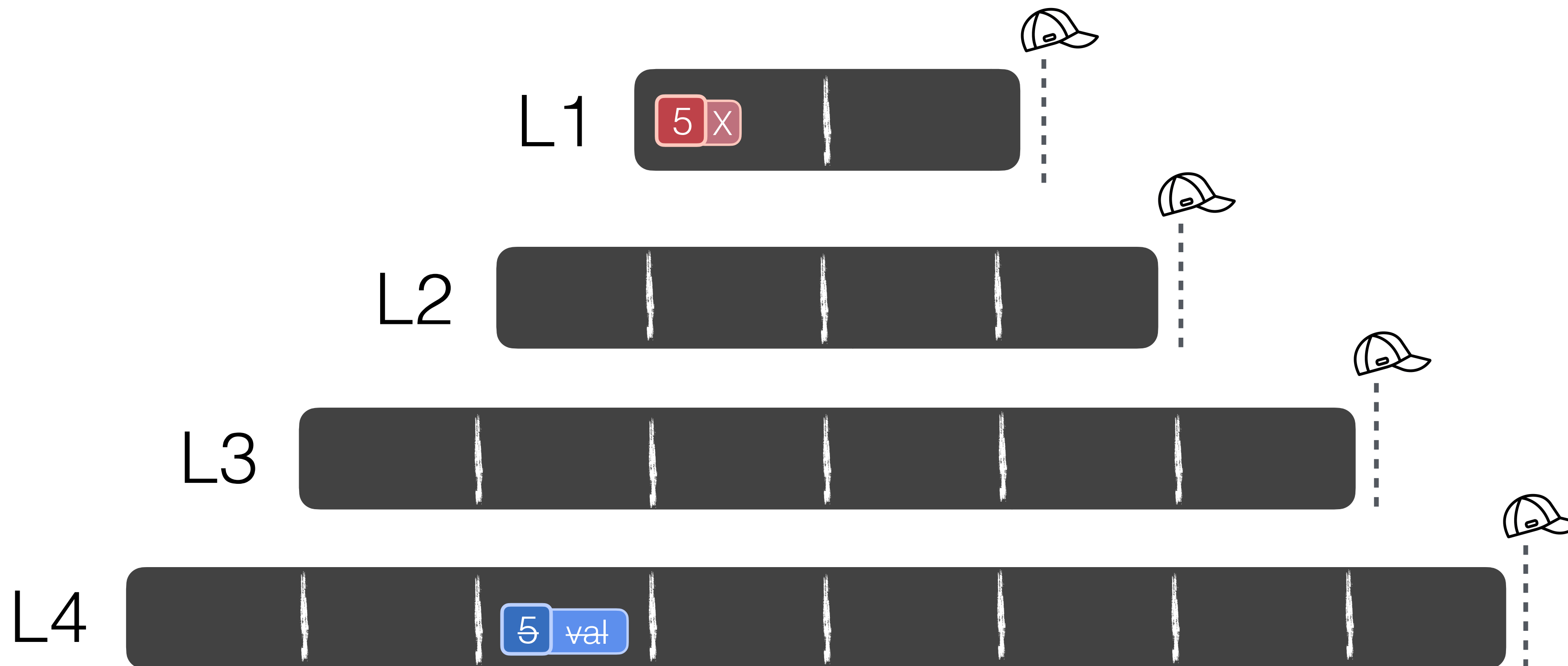
Deletes in LSMs

delete(5)



Deletes in LSMs

delete(5)

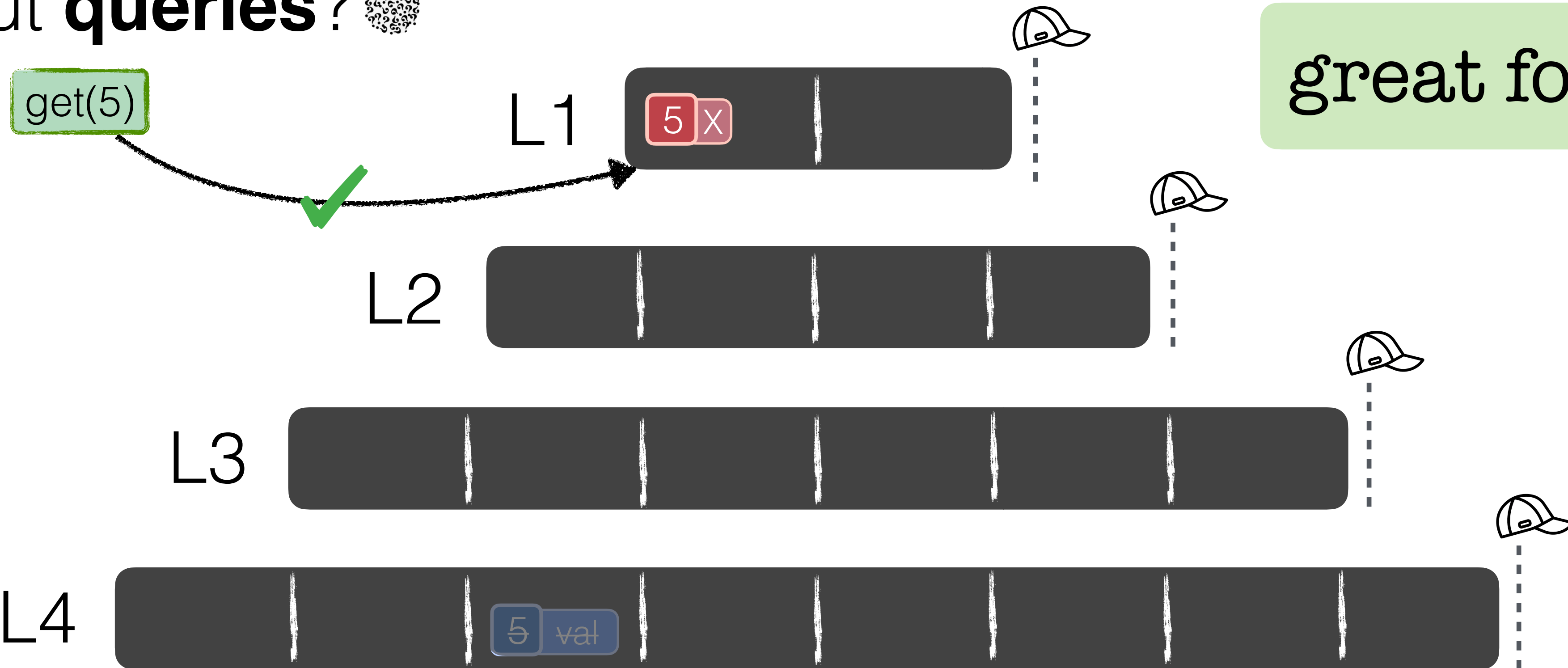


Deletes in LSMs



delete(5)

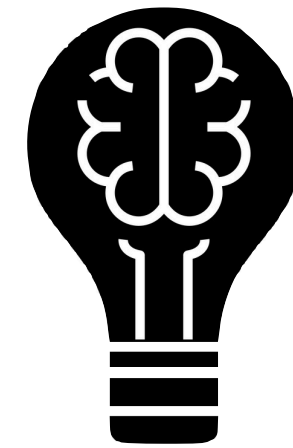
What about **queries**? 



great for inserts

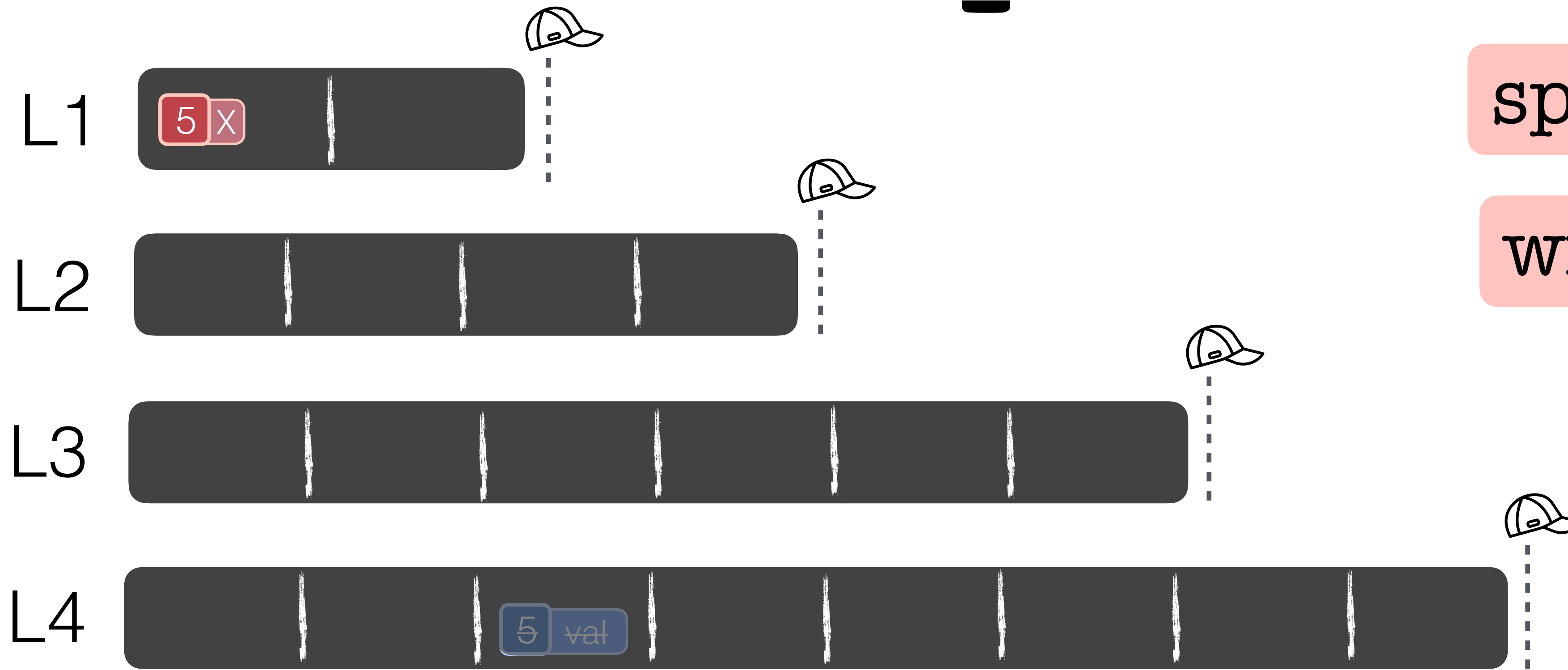
great for queries

Deletes in LSMs



Thought Experiment 1

What's the problem with **logical deletes**?



space amplification

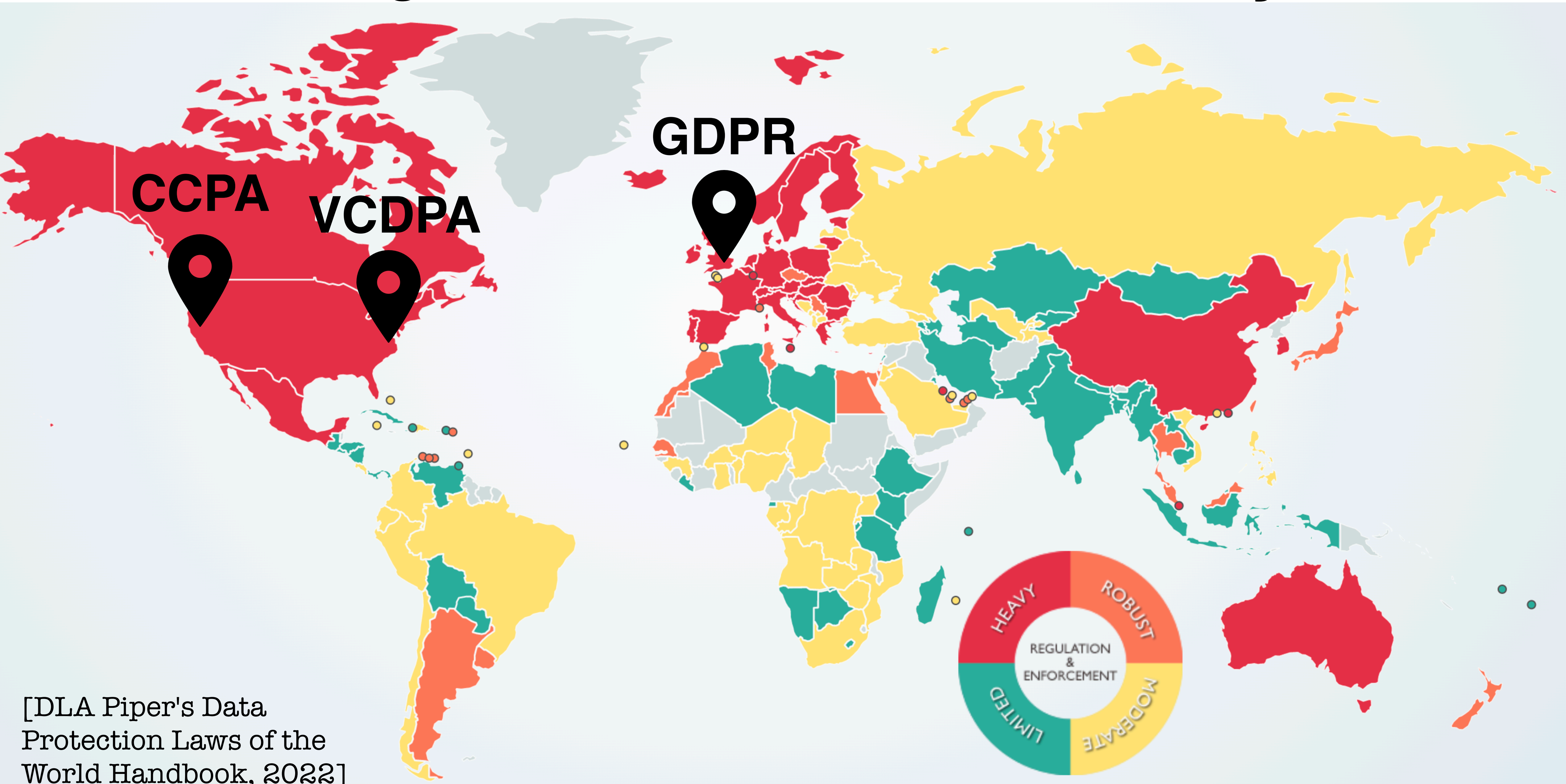
write amplification

poor reads

increased cost

loss of privacy!

Logical Deletes & **Data Privacy**



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



GDPR
(EU, UK)



CCPA
(California)



VCDPA
(Virginia)



Right to be forgotten



Right to delete

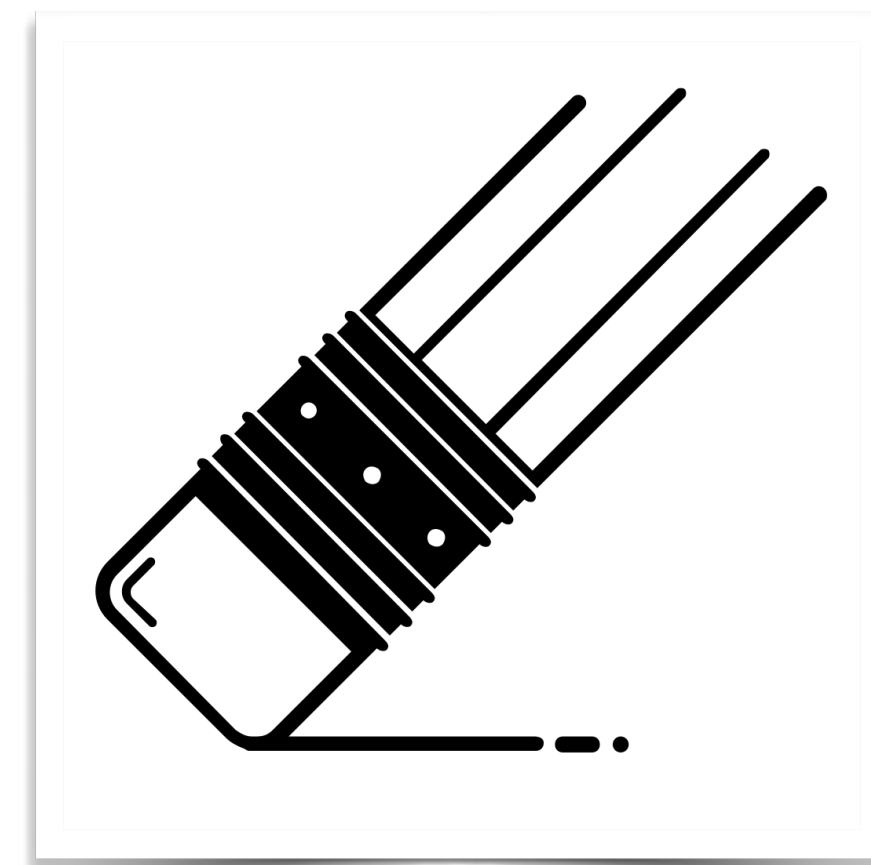


Deletion right




timely
deletes

+



persistent
deletes


Right to be forgotten


Right to delete


Deletion right

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

Zack Whittaker, Natasha Lomas / 1:57 PM EST • February 15, 2019


timely
deletes

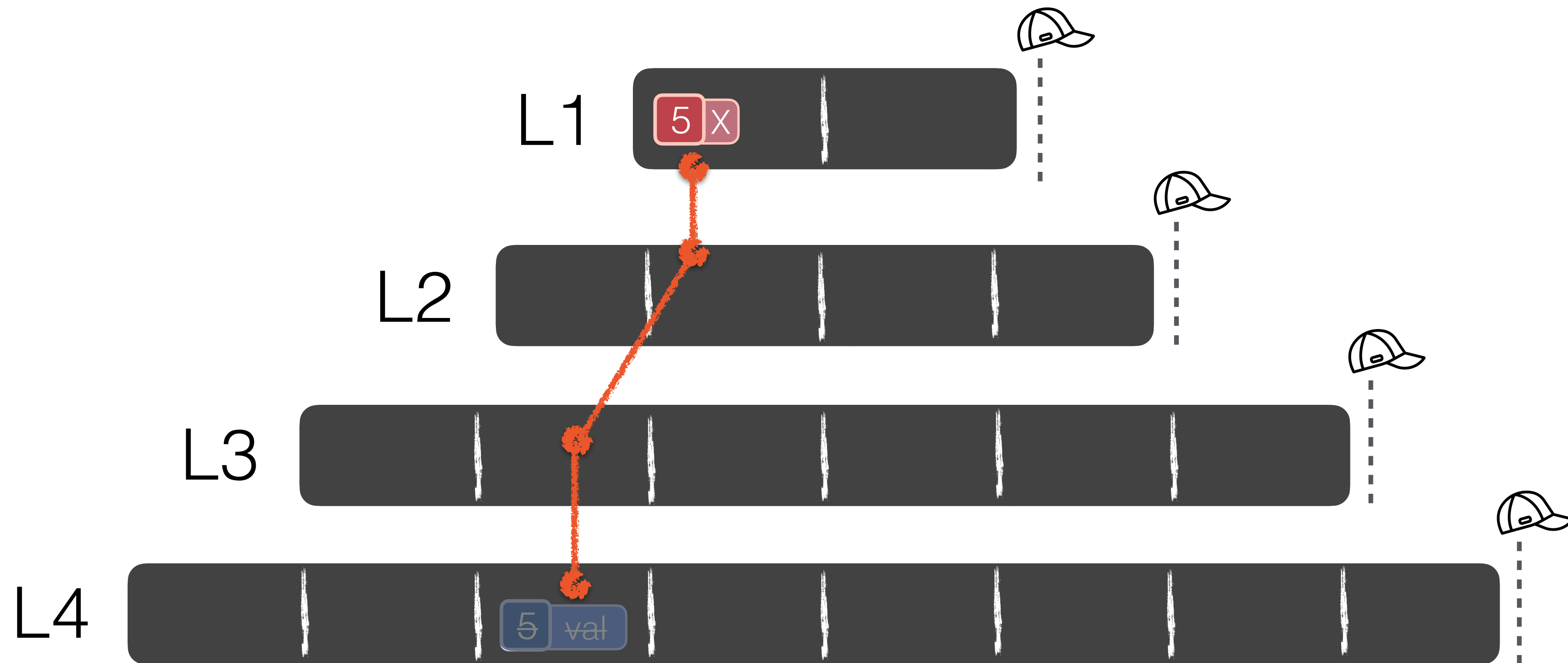
+


persistent
deletes



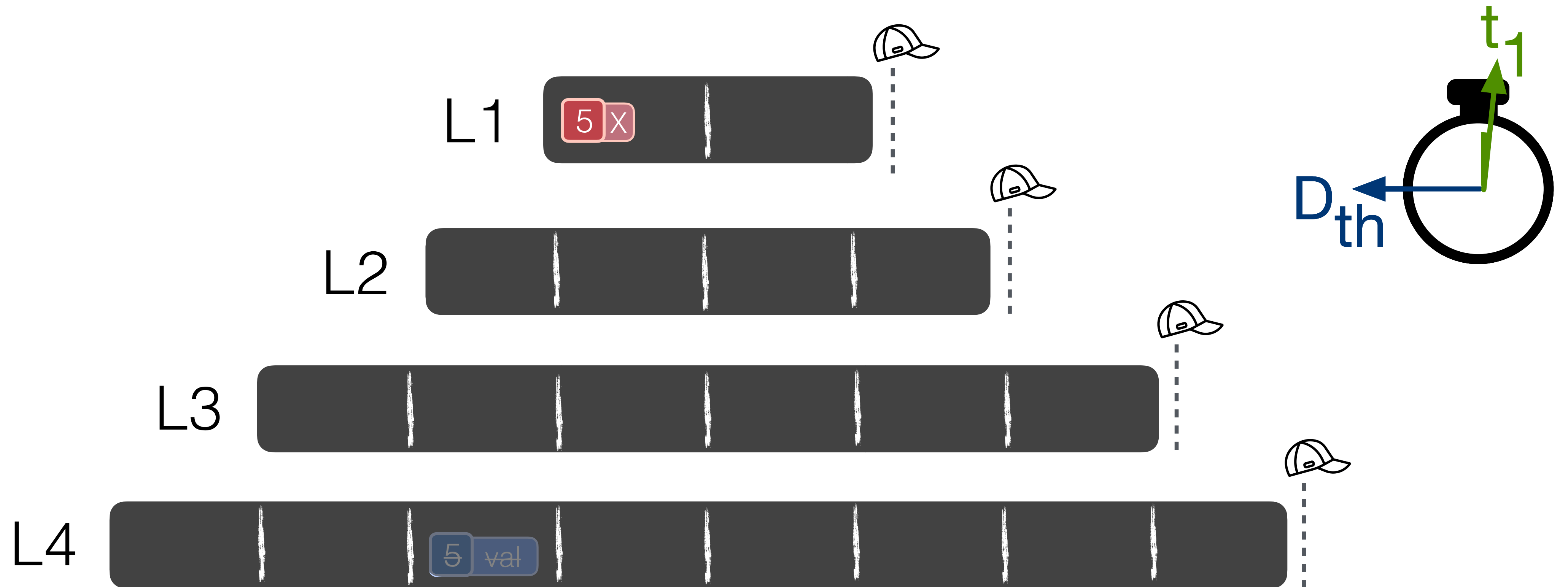
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



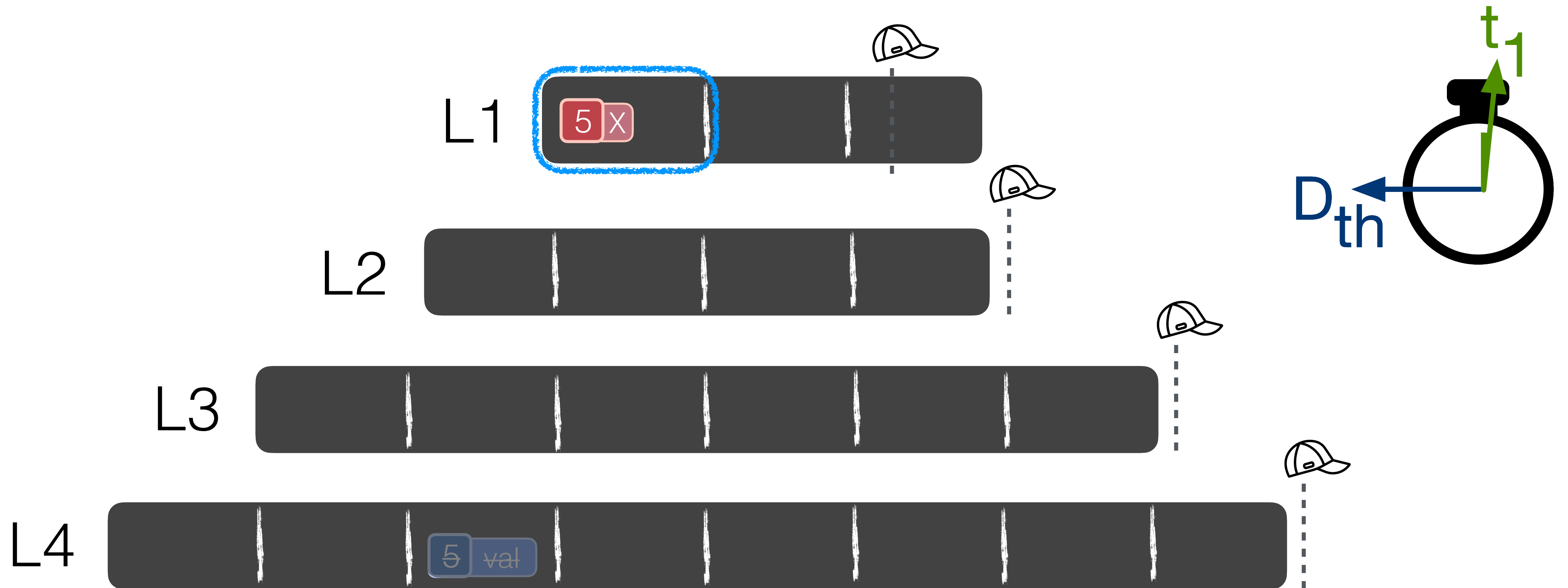
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



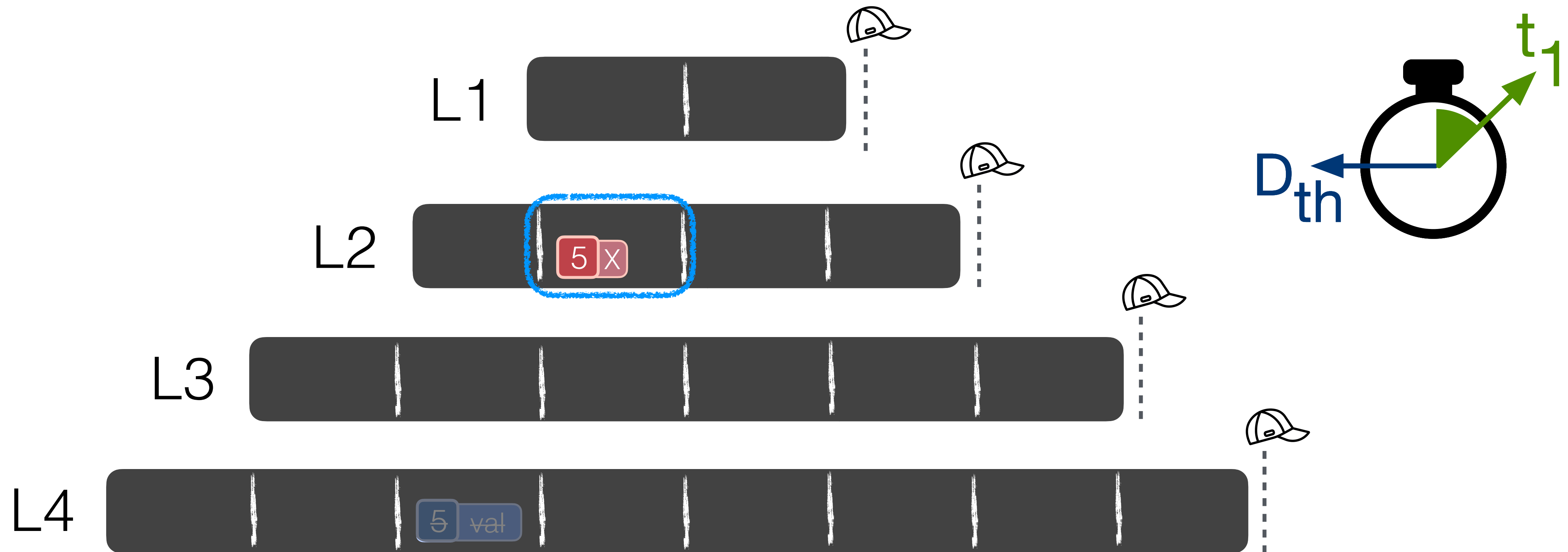
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



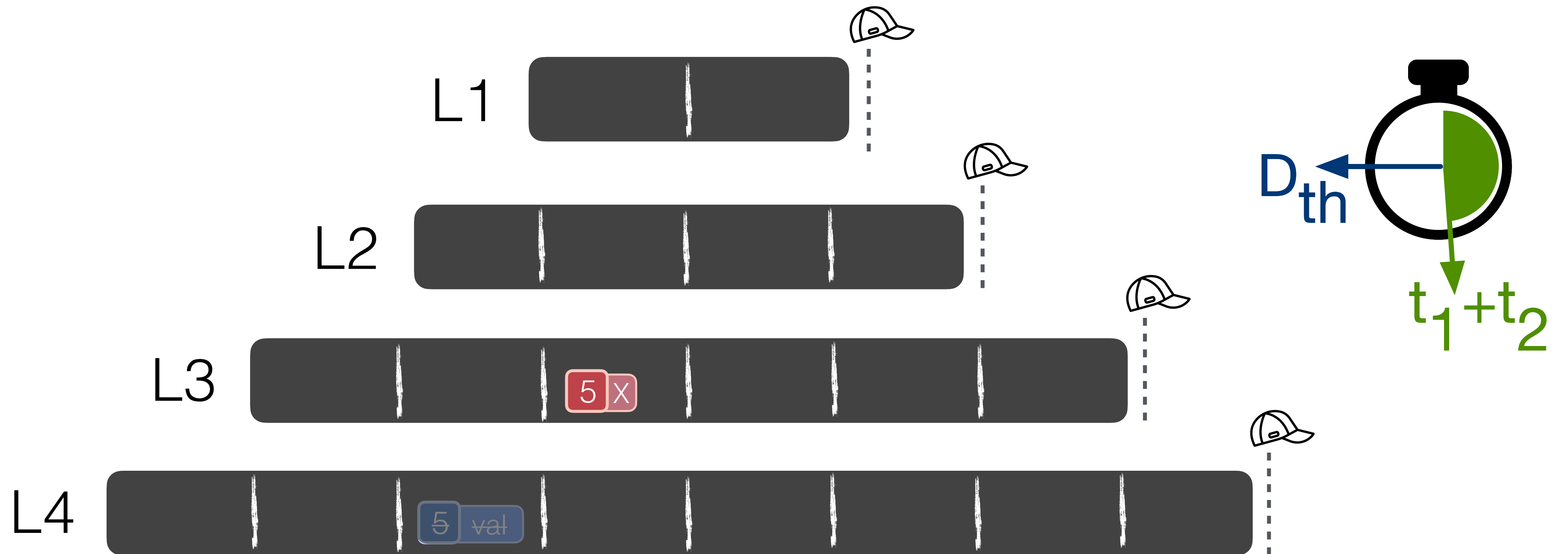
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



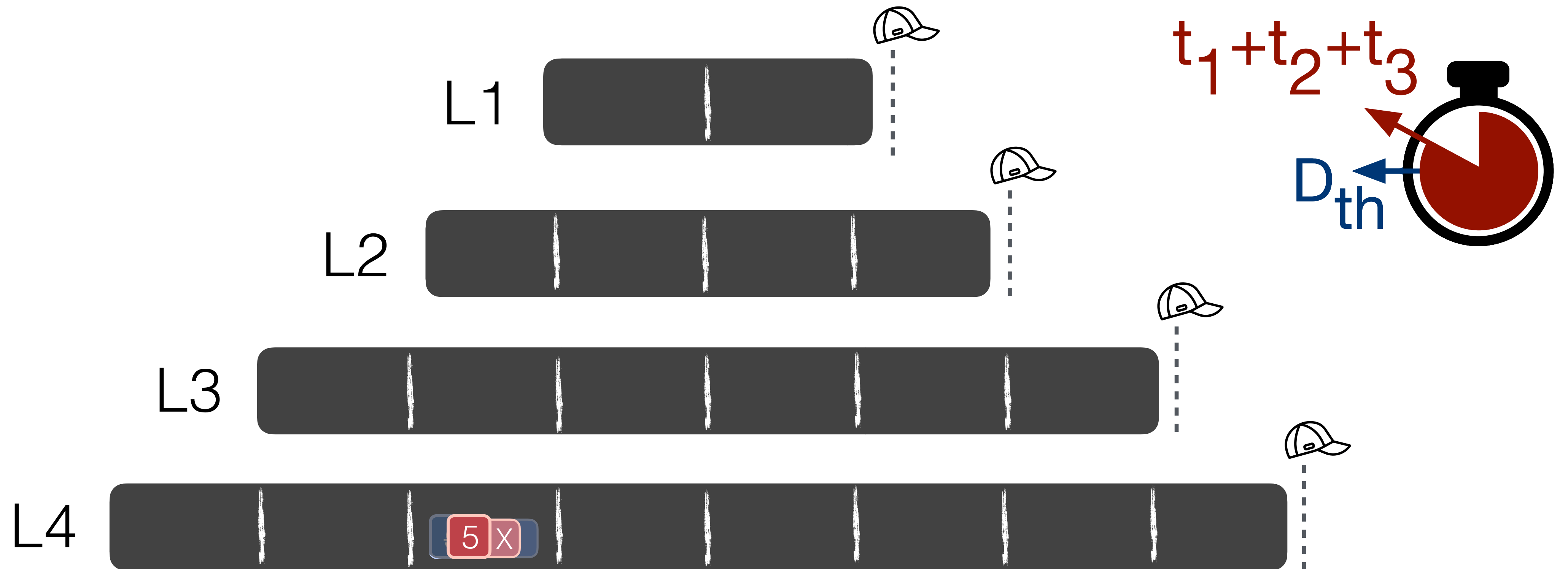
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



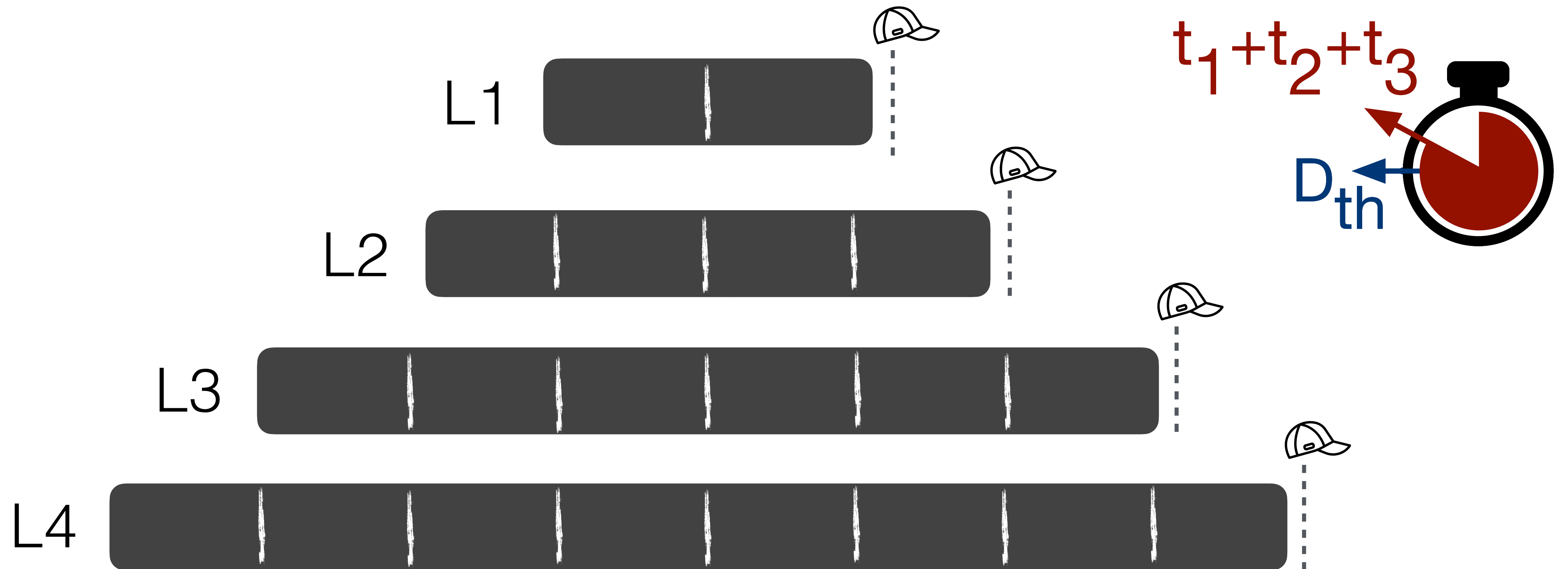
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



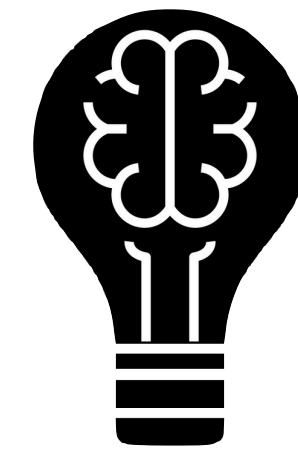
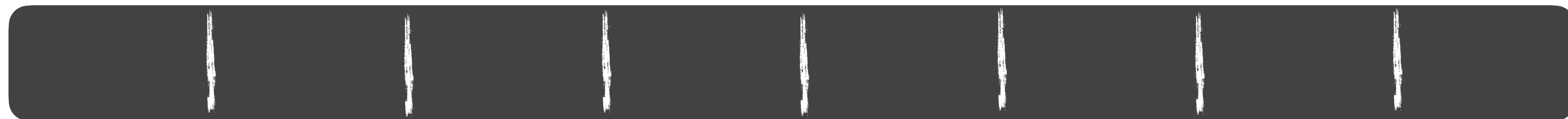
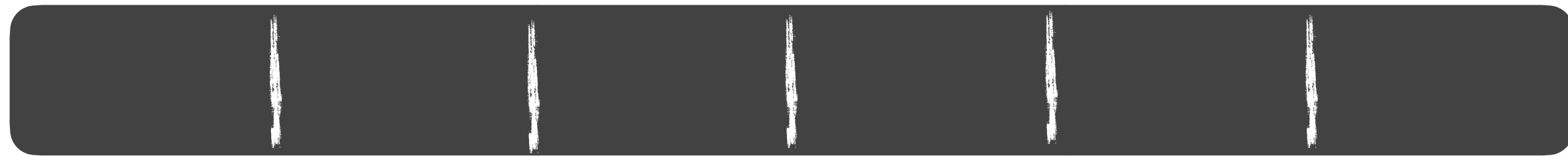
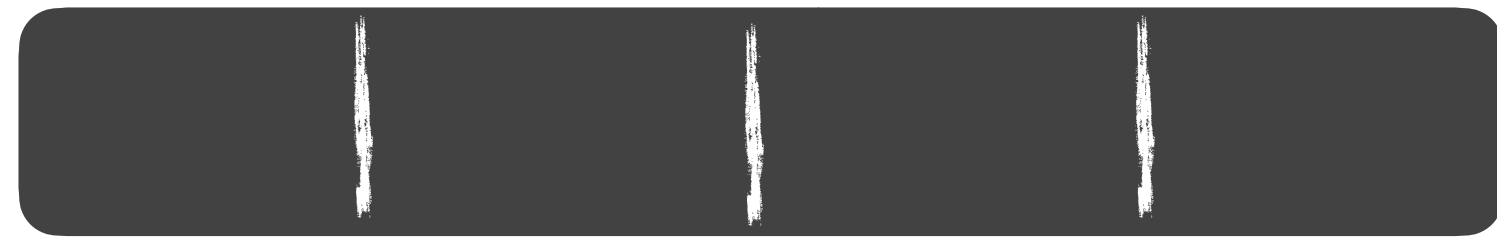
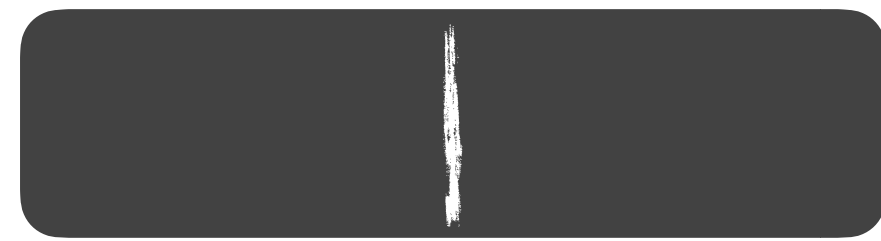
Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



Thought Experiment 2

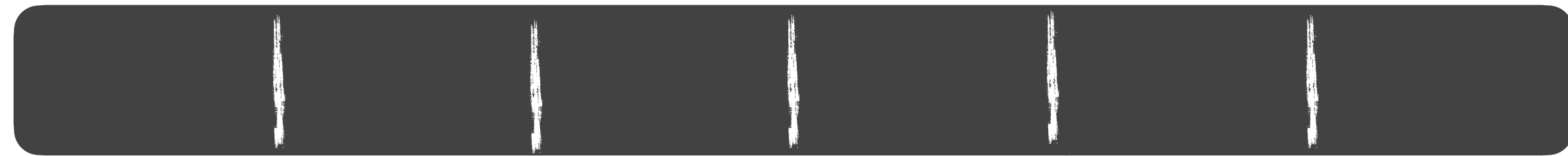
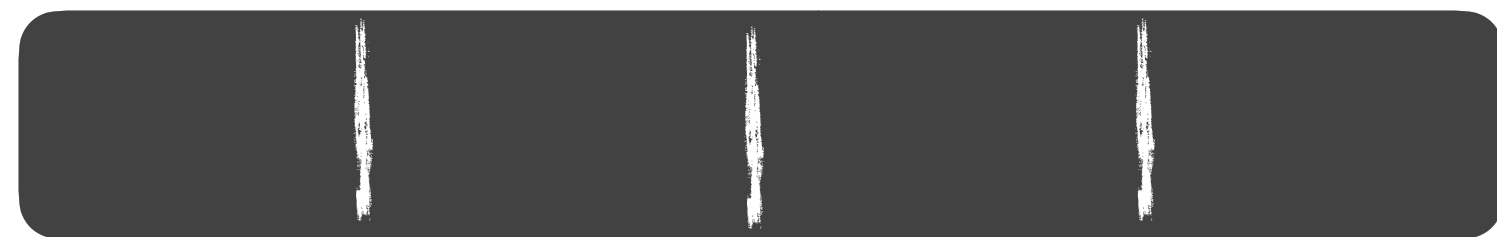
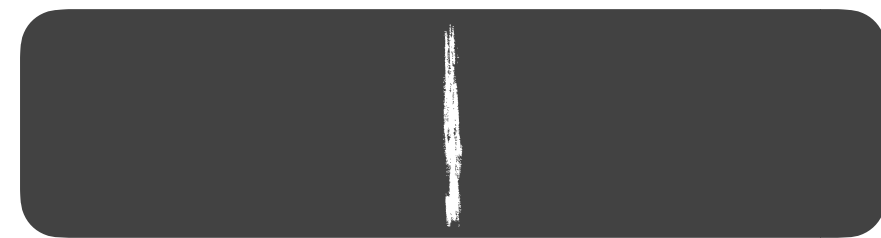
What drives the **delete persistence latency**?



unbounded delete persistence latency

Problem: Persisting Deletes Timely

delete(5) within threshold: D_{th}



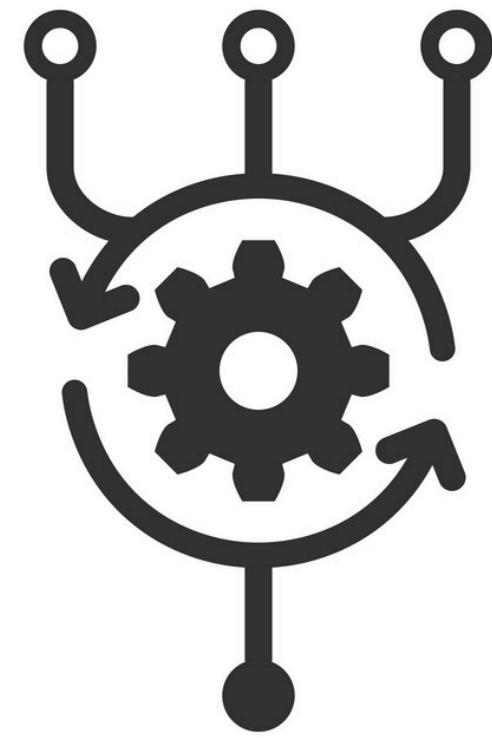
unbounded delete
persistence latency

- File picking policy
- Tree shape
- Ingestion rate

Intuition: Compaction holds the key!

FADE

FAst DElete



family of
**compaction
strategies**

FAst DElete

compaction
trigger



compaction file
picking policy



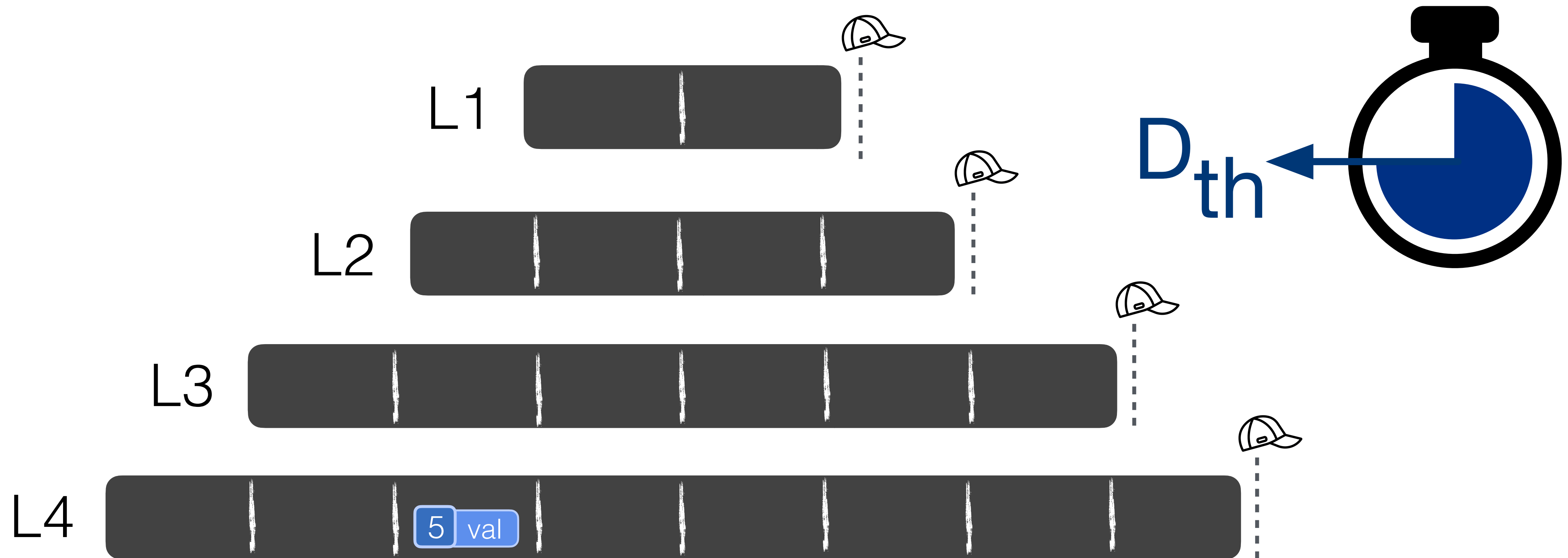
decompose tasks



piggyback

FAst DElete

delete(5) within threshold: D_{th}

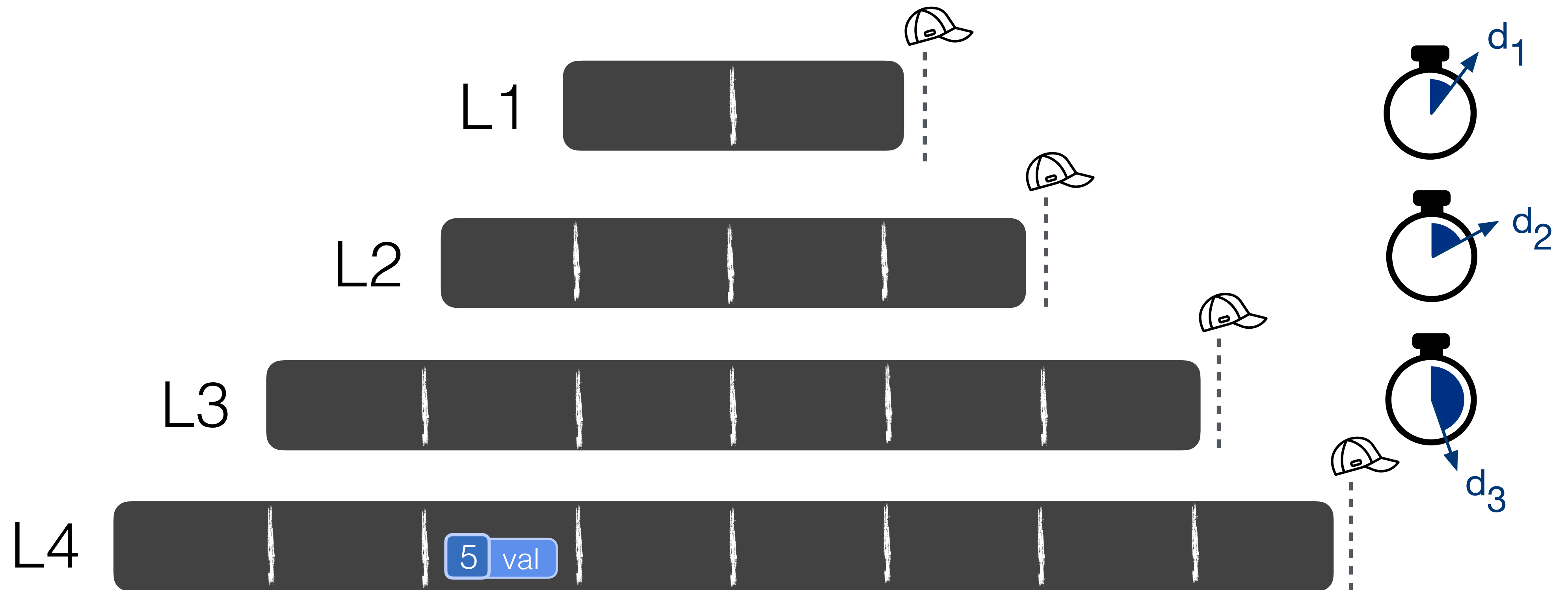


FAst DElete

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

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

delete(5) within threshold: D_{th}

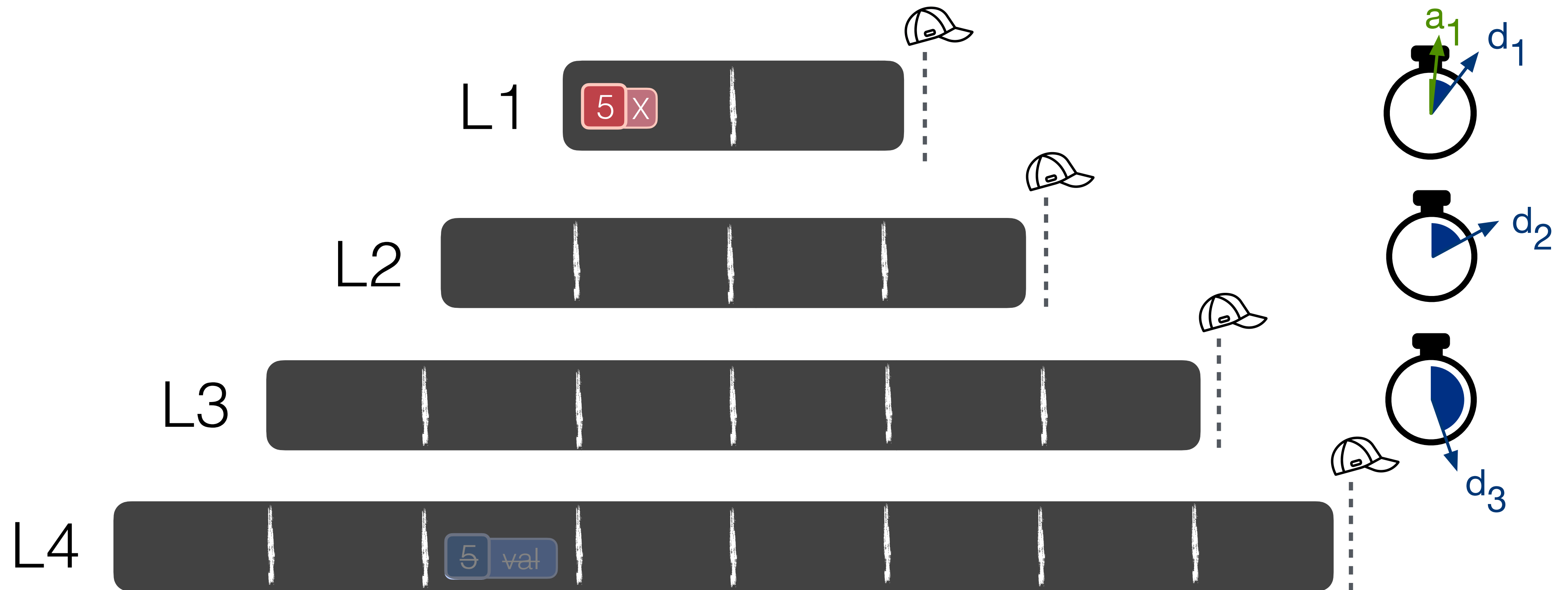


FASt DElete

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

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

delete(5) within threshold: D_{th}

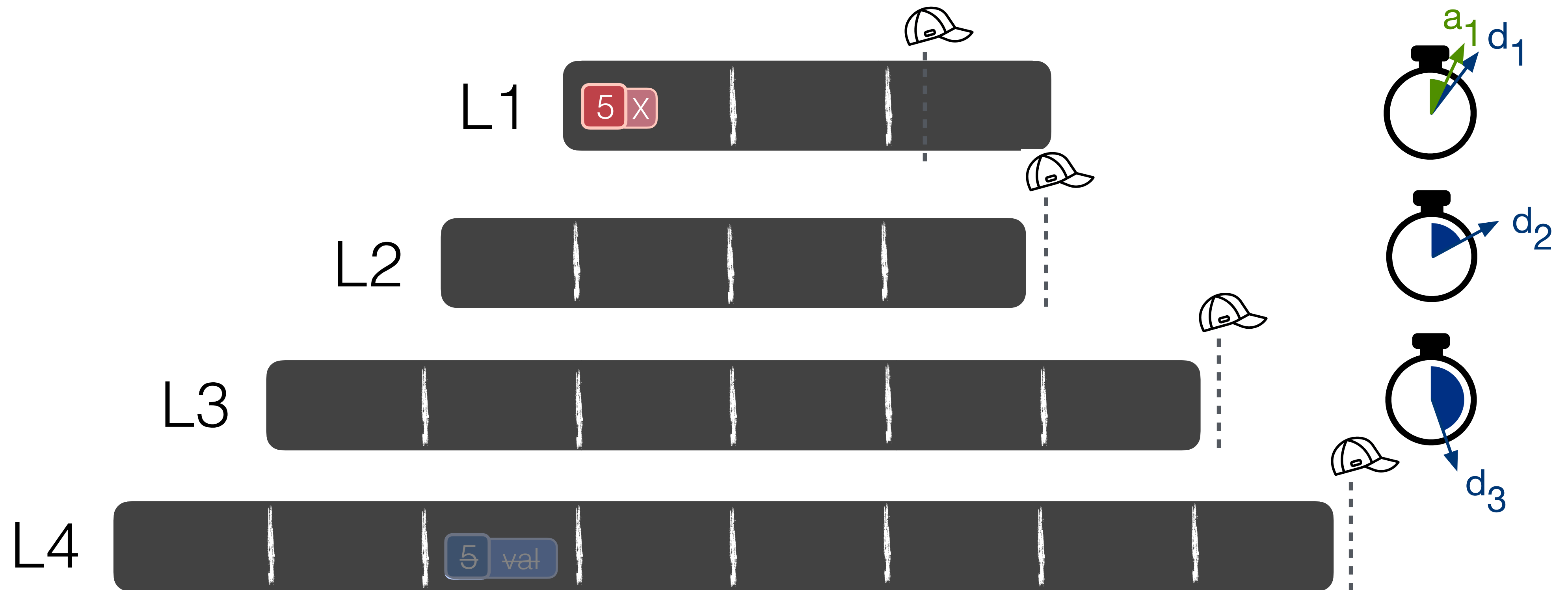


FASt DElete

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

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

delete(5) within threshold: D_{th}

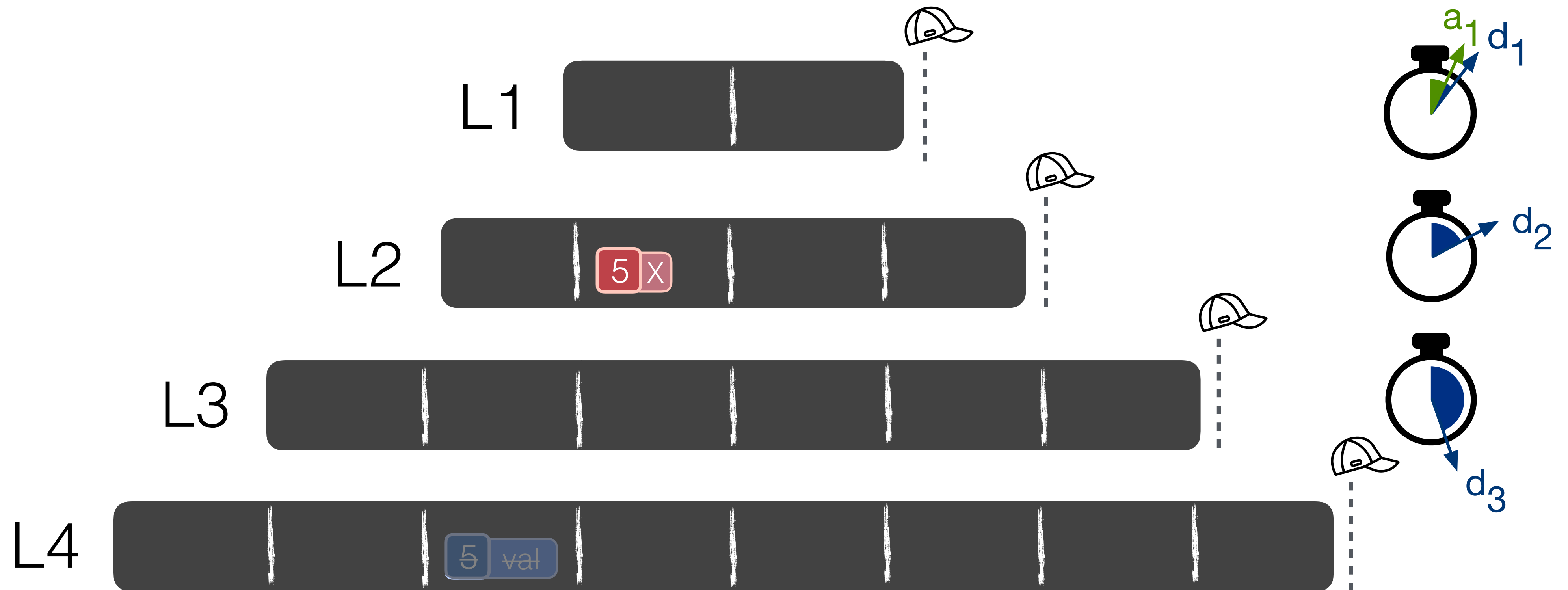


FASt DElete

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

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

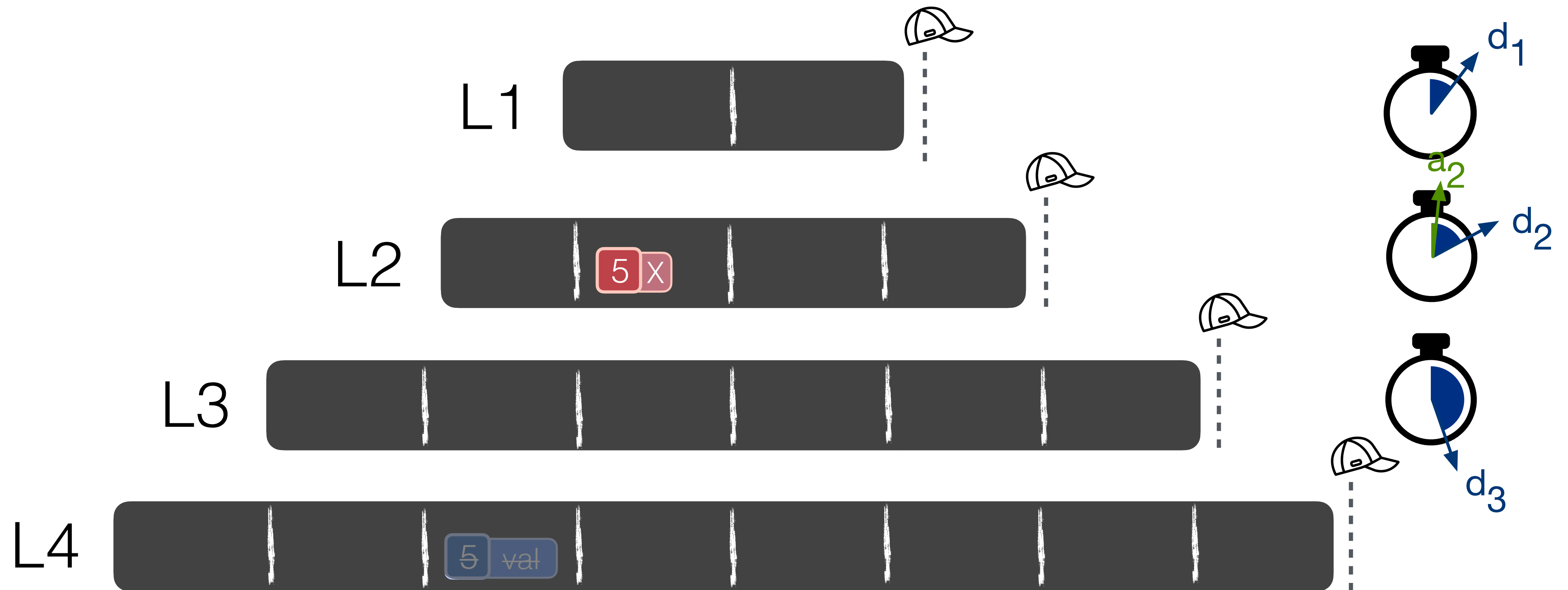
delete(5) within threshold: D_{th}



FAst DElete

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

delete(5) within threshold: D_{th}

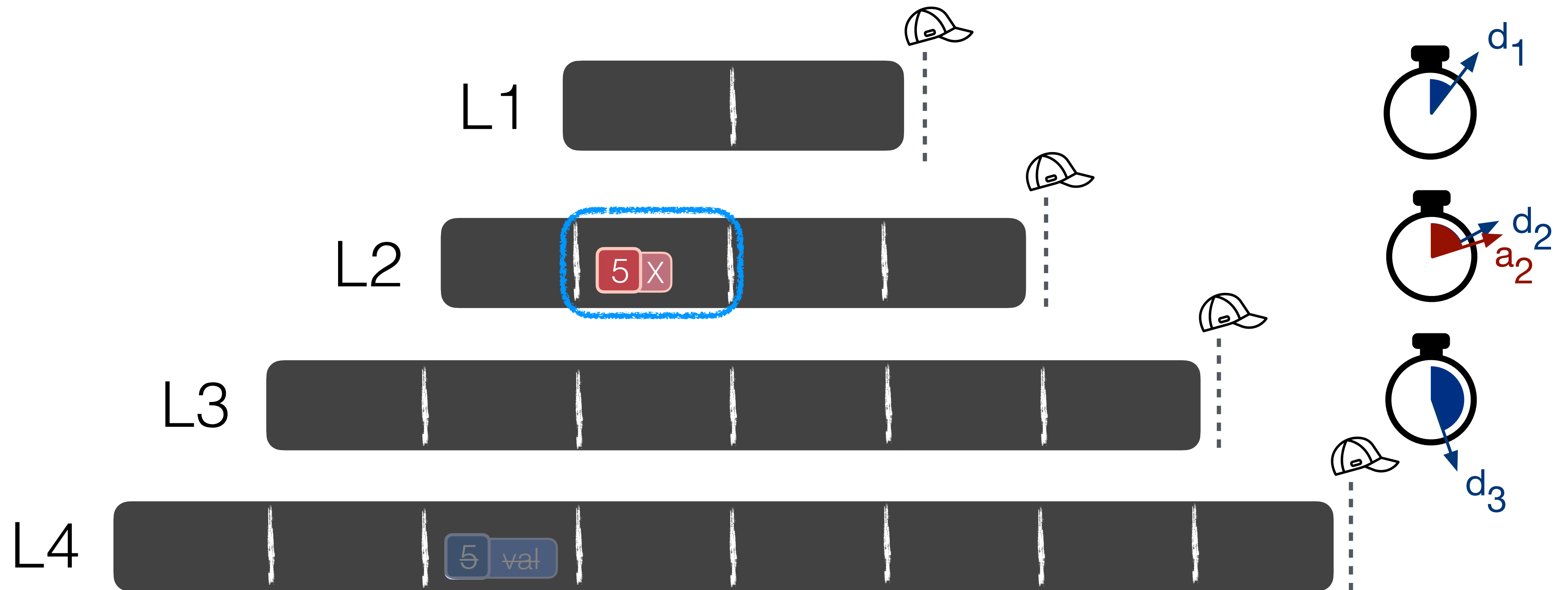


FASt DElete

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

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

delete(5) within threshold: D_{th}

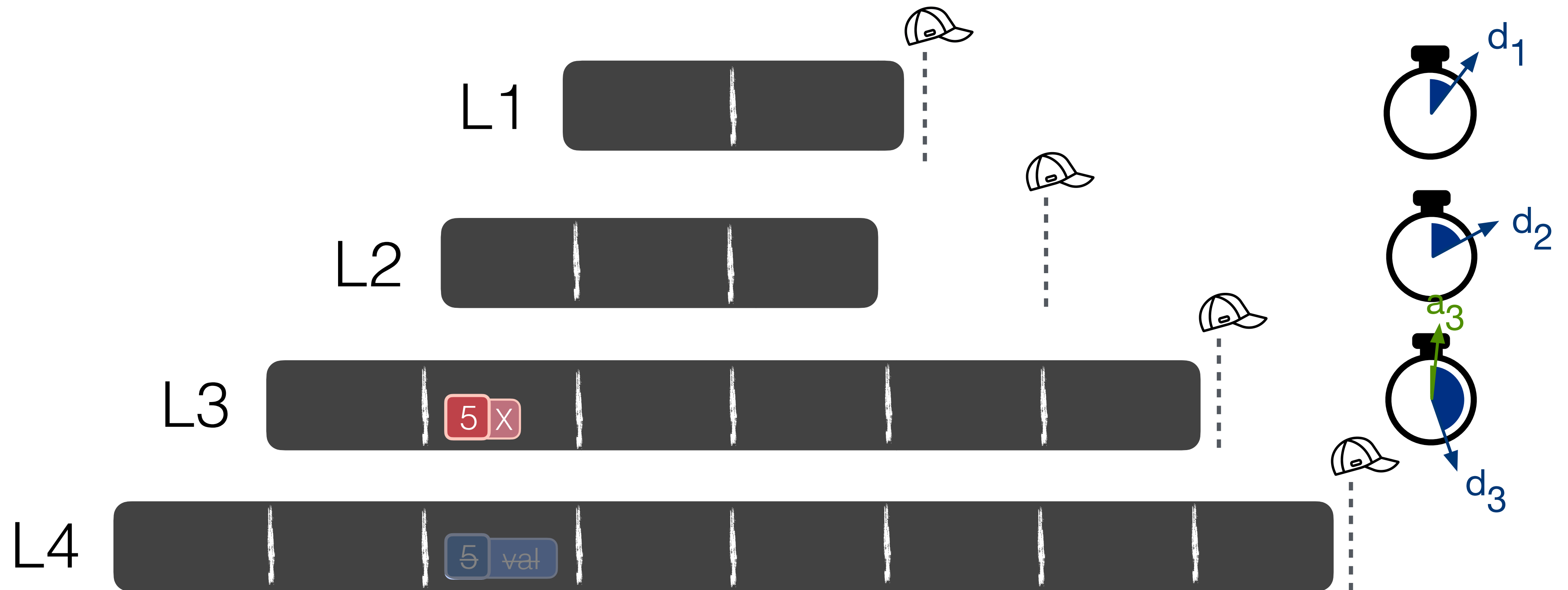


FASt DElete

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

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

delete(5) within threshold: D_{th}

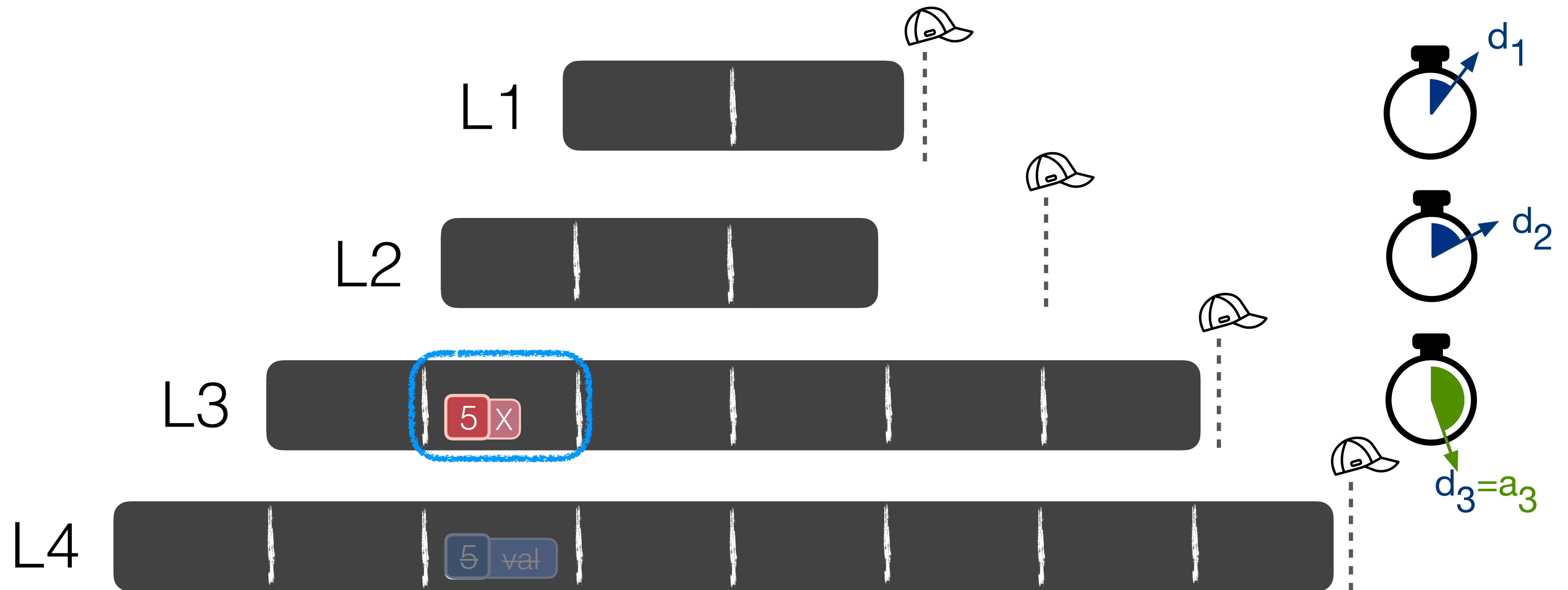


FAst DElete

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

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

delete(5) within threshold: D_{th}

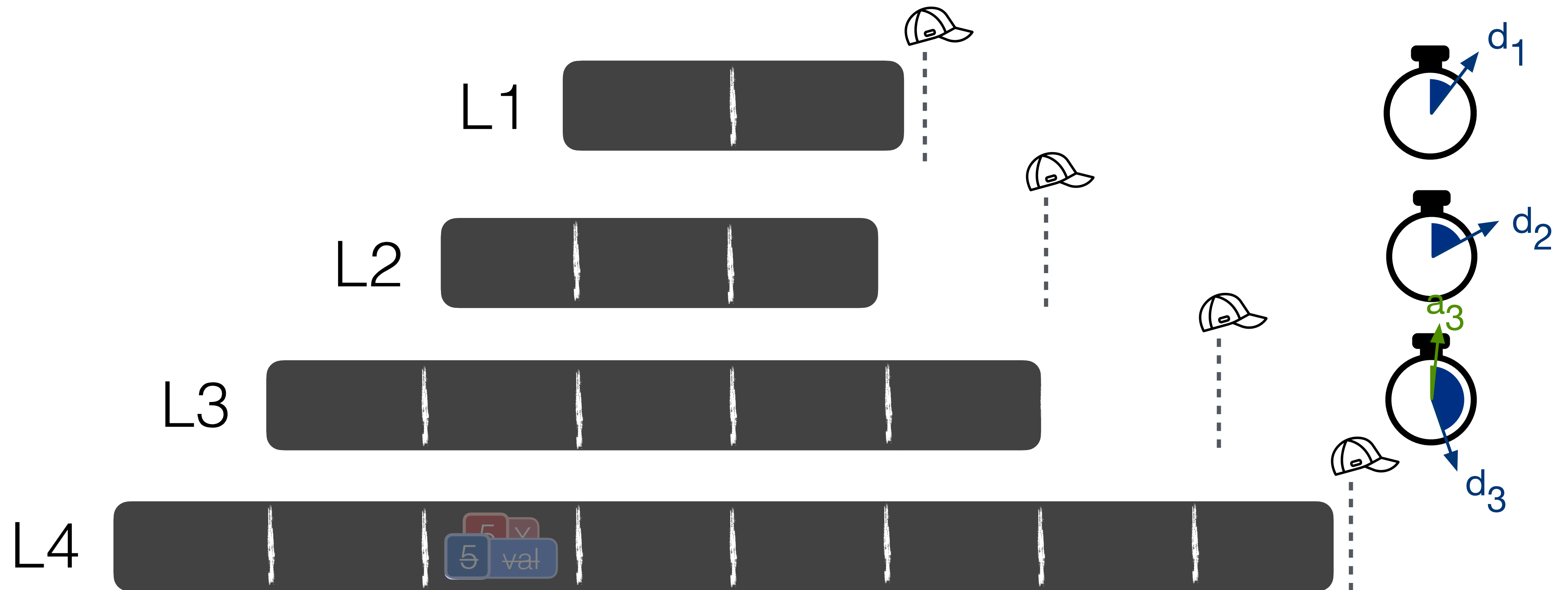


FASt DElete

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

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

delete(5) within threshold: D_{th}

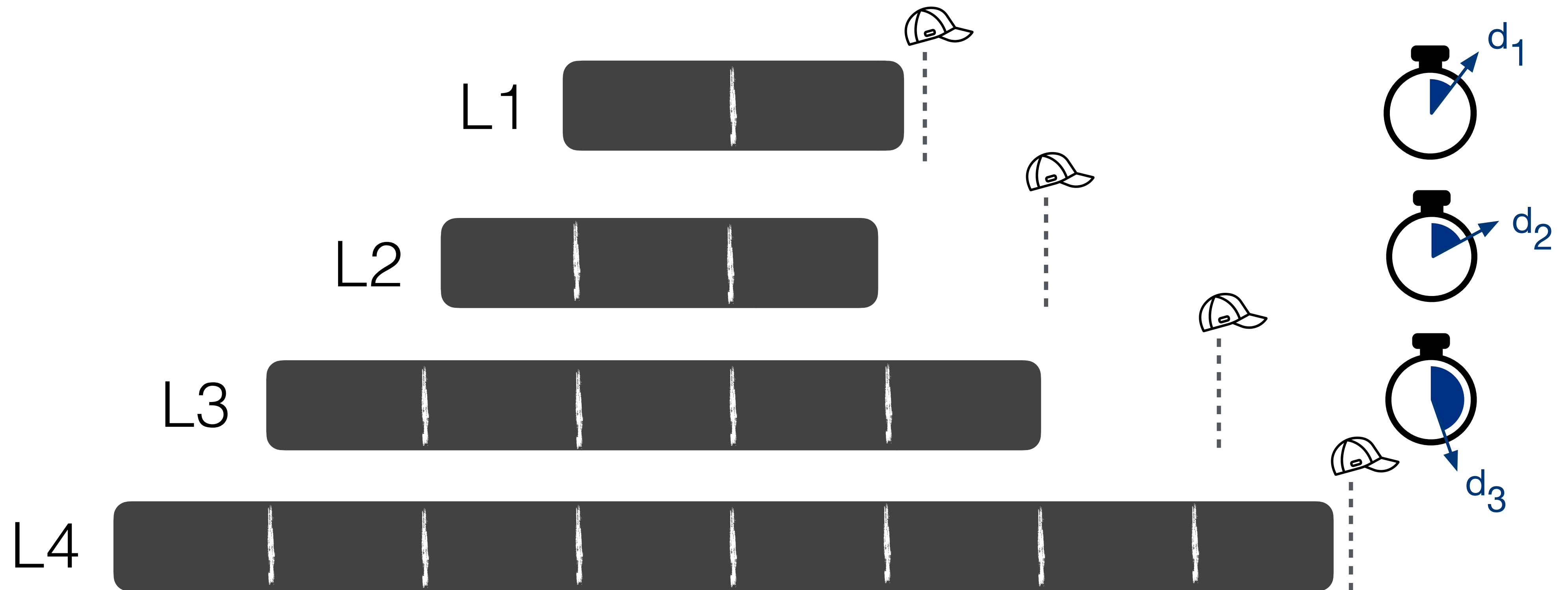


FAst DElete

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

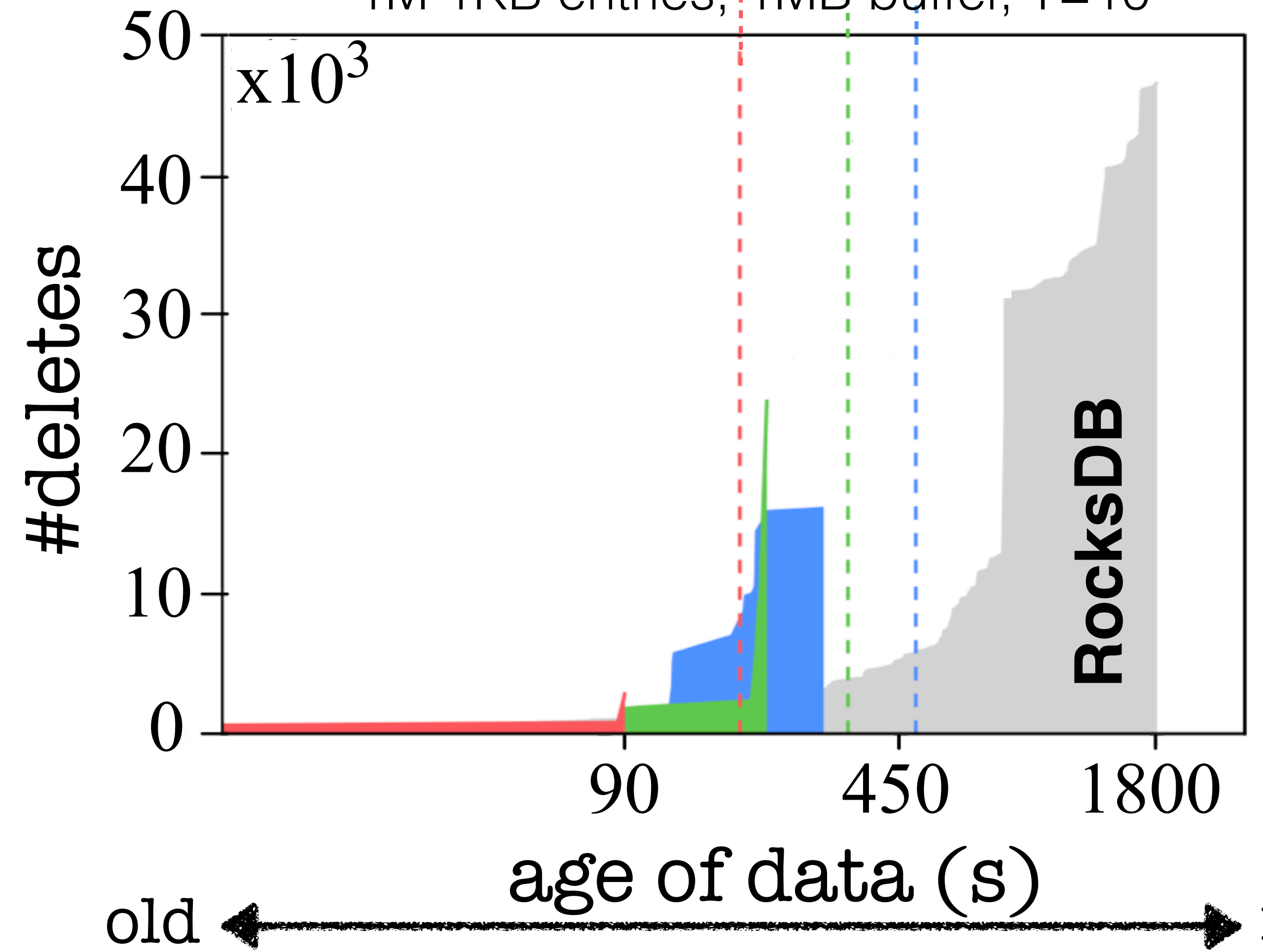
$$d_i = T \cdot d_{i-1}$$

delete(5) within threshold: D_{th}



persist all deletes within: 300s
150s 600s

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



persists deletes timely

within threshold

old ← → new

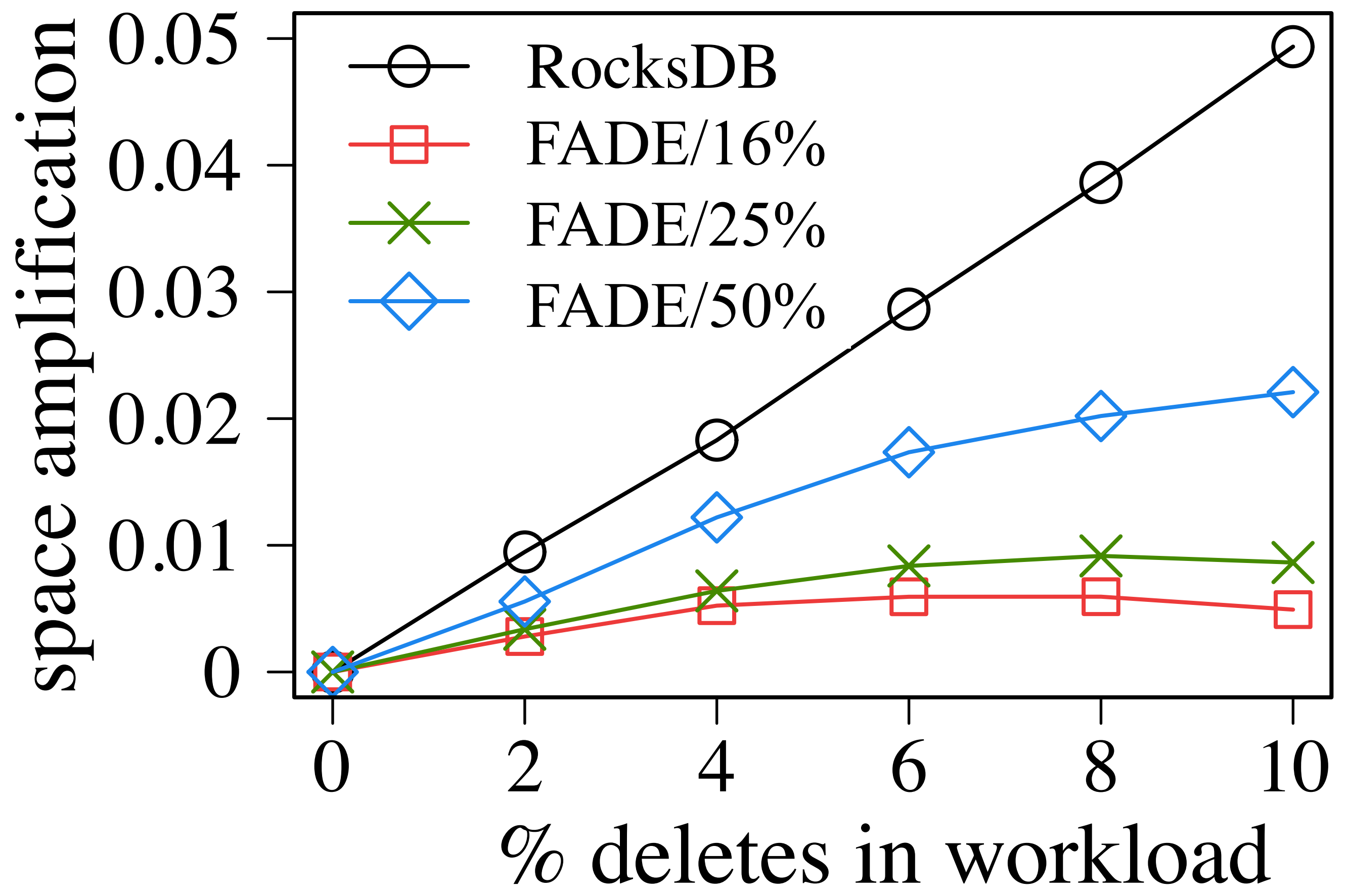
reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

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



improved read performance

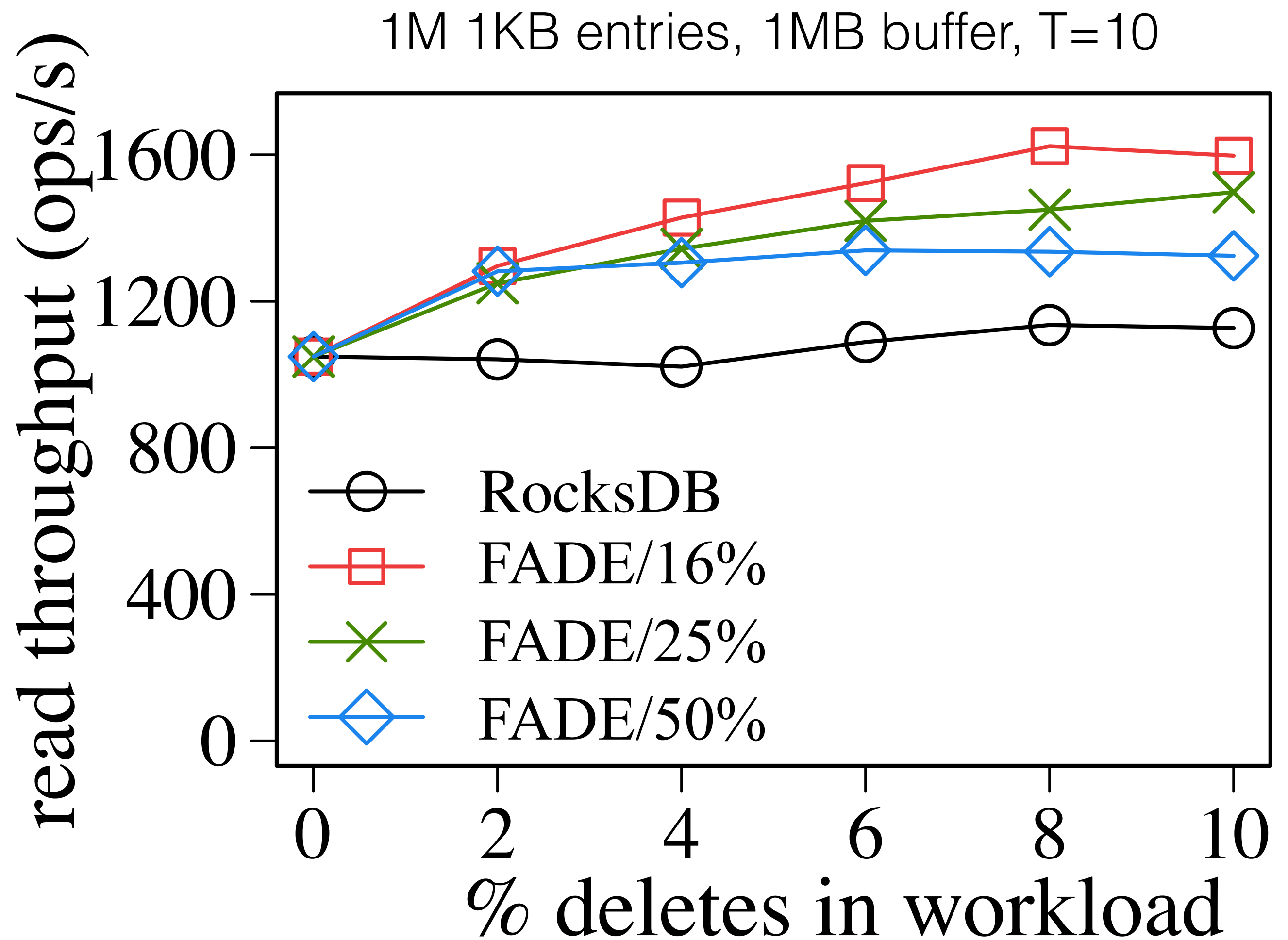
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

4% - 25%

improved read performance

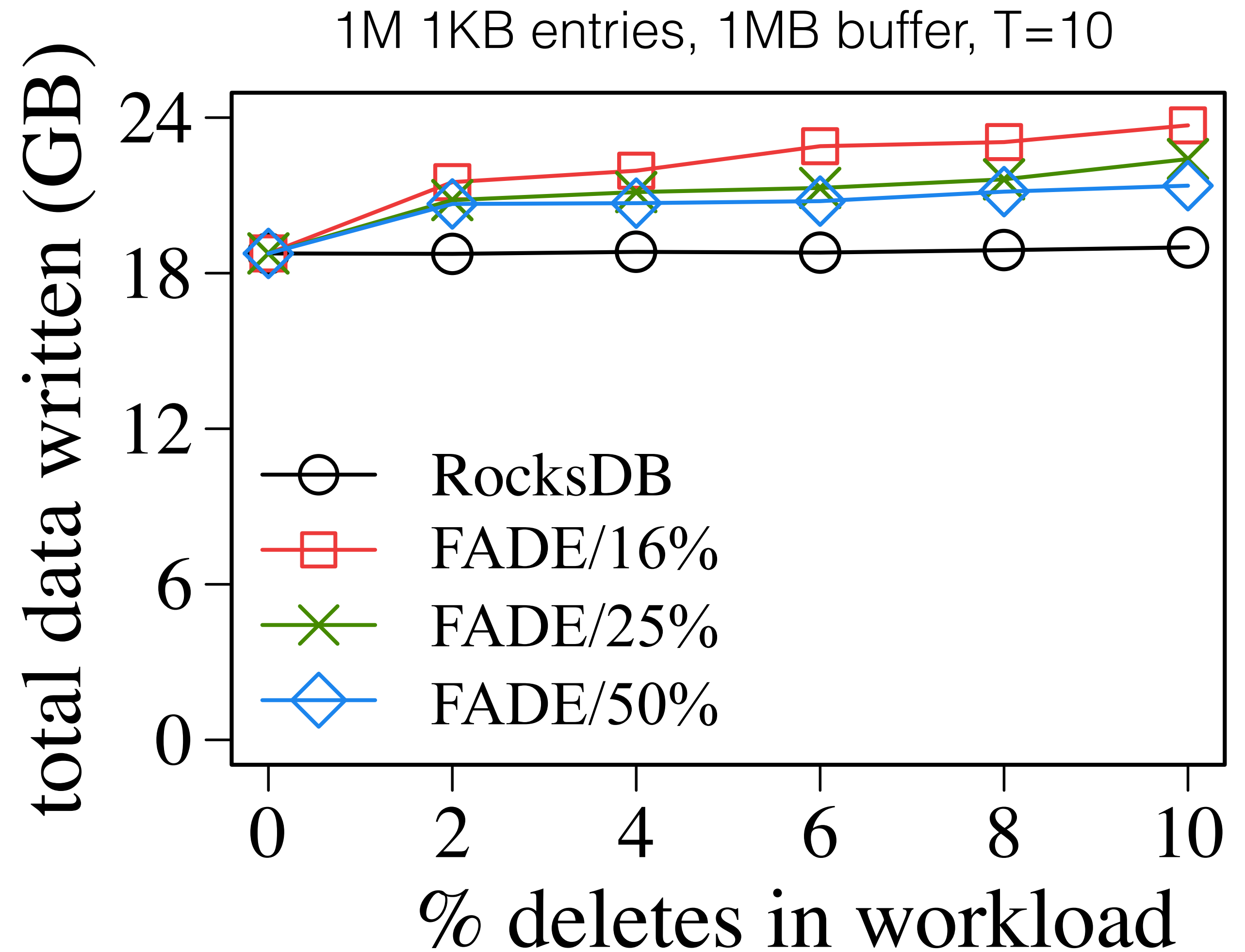
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

4% - 25%

improved read performance

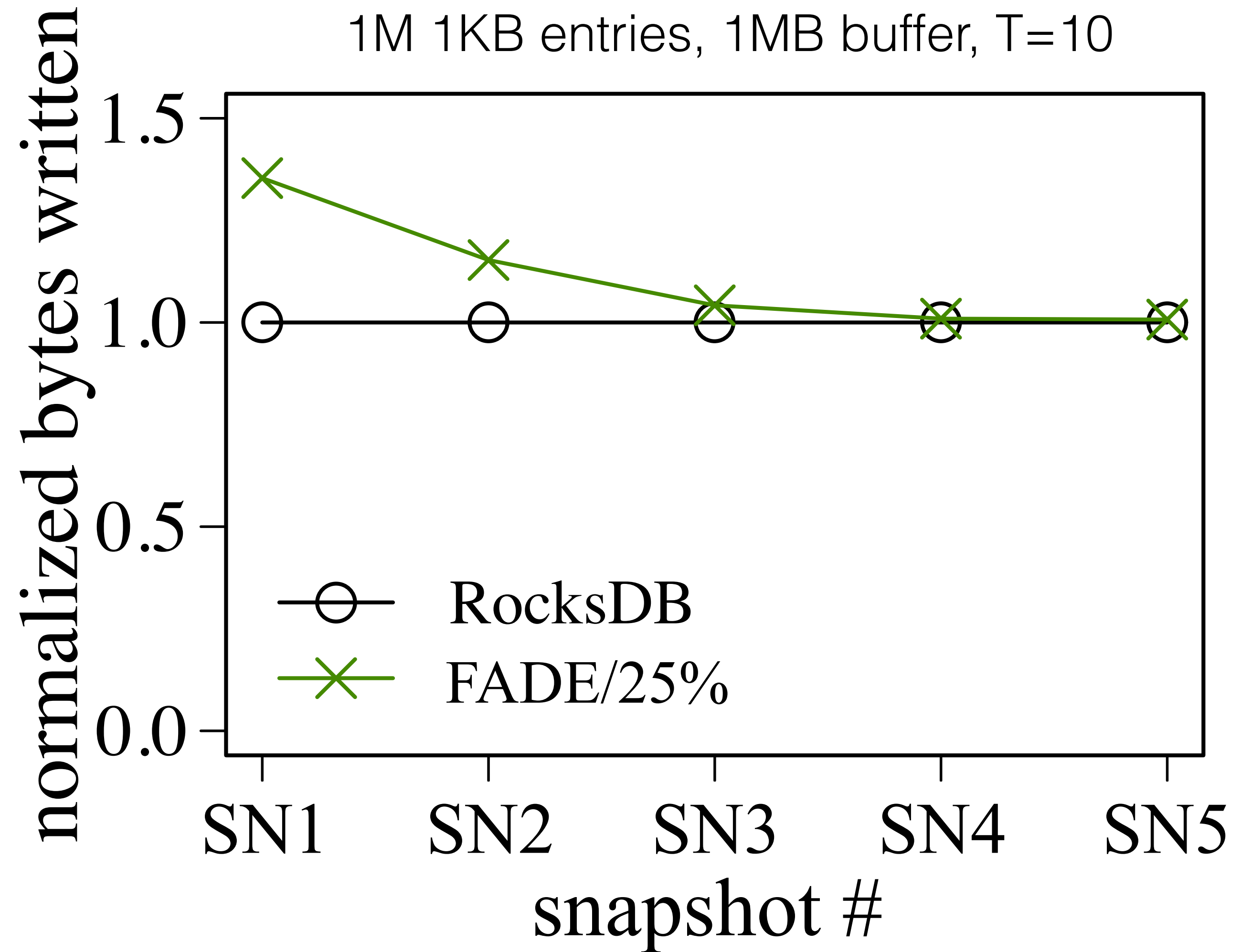
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



higher write amplification

0.7%

improved read performance

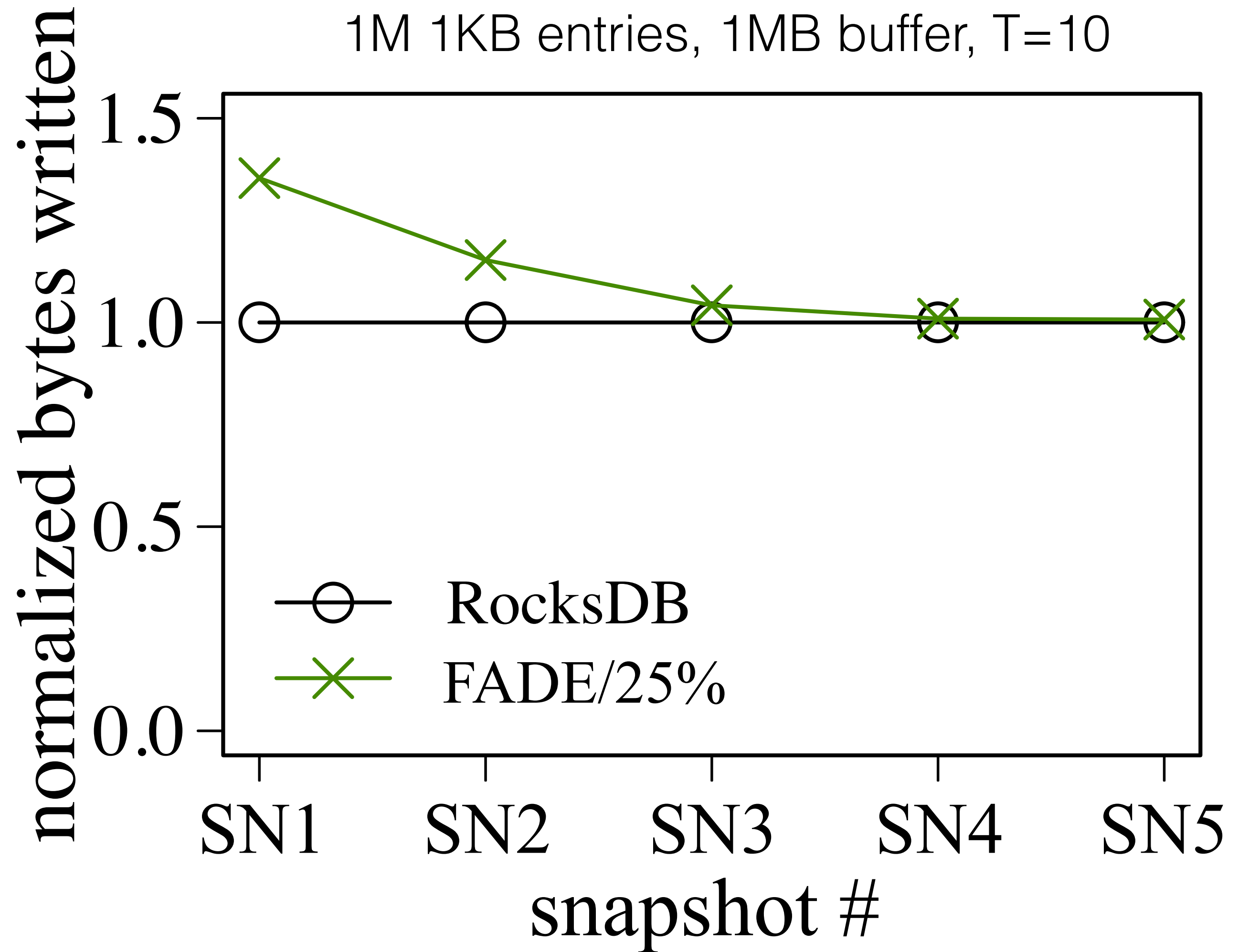
1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold



FADE

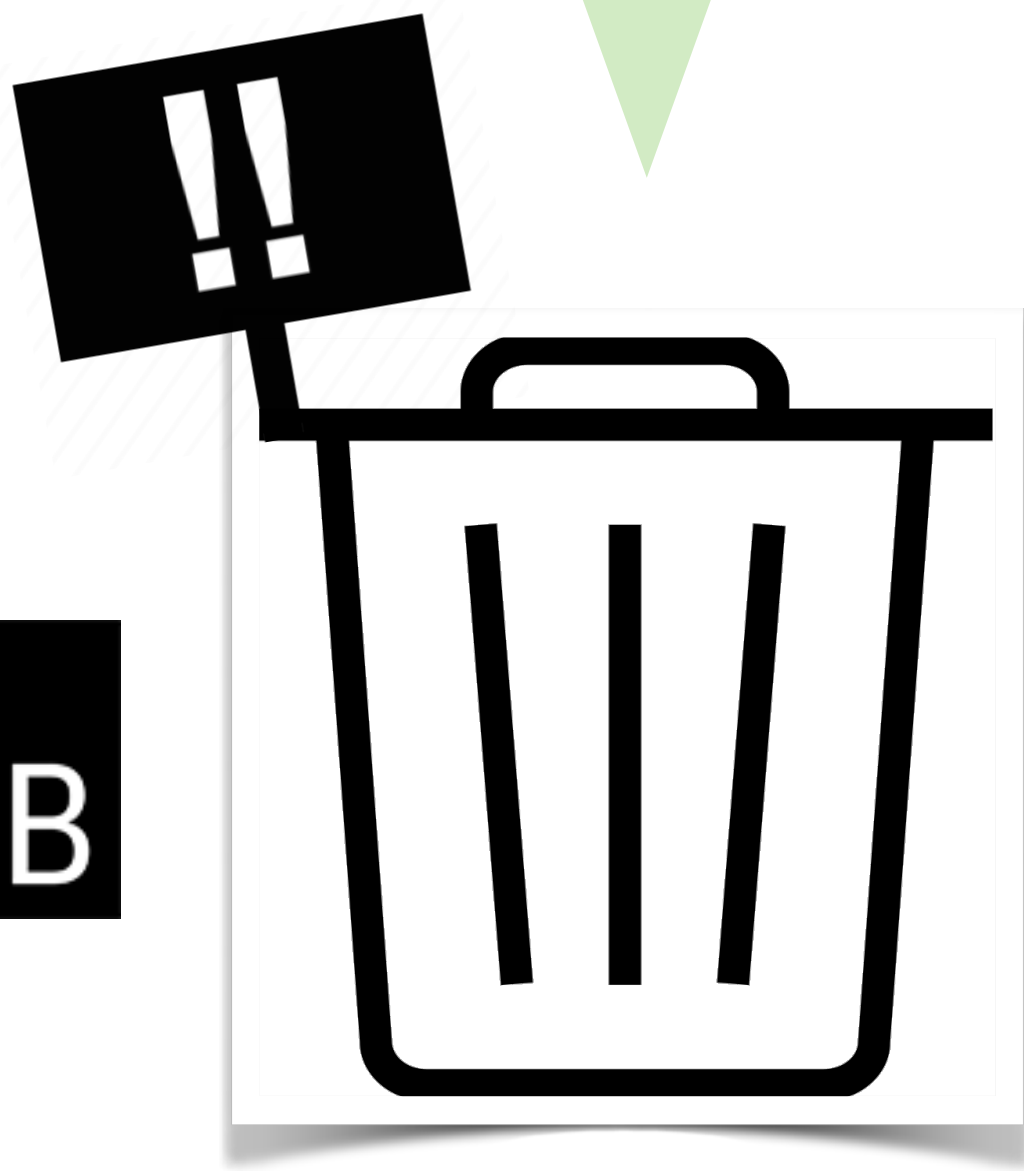


on-demand

persist all logical
deletes within **D** days



FADE



on-demand

persist all logical
deletes within D days

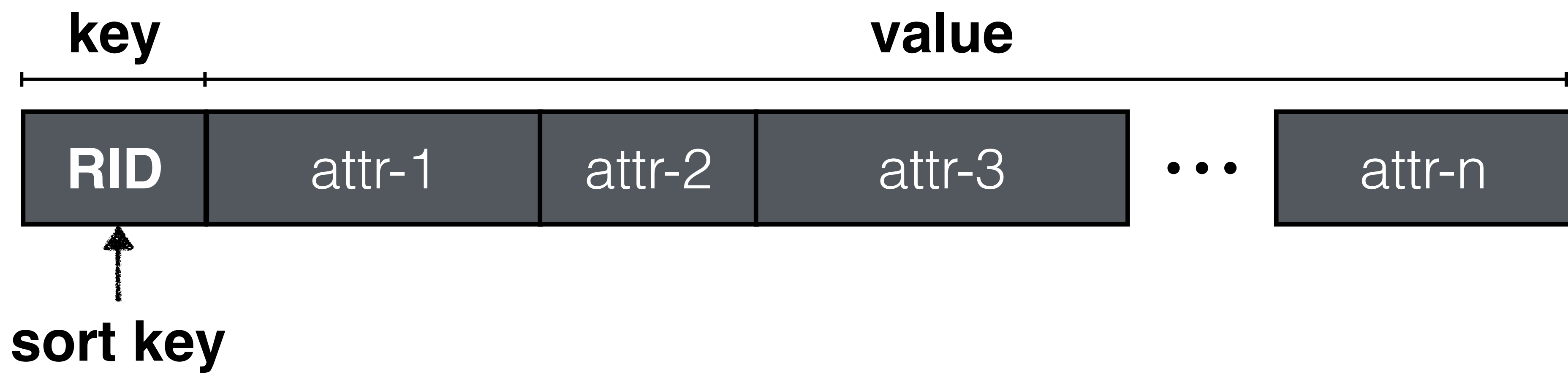


retention-based

delete all data older
than T days

Realizing **Retention-Based Deletes**

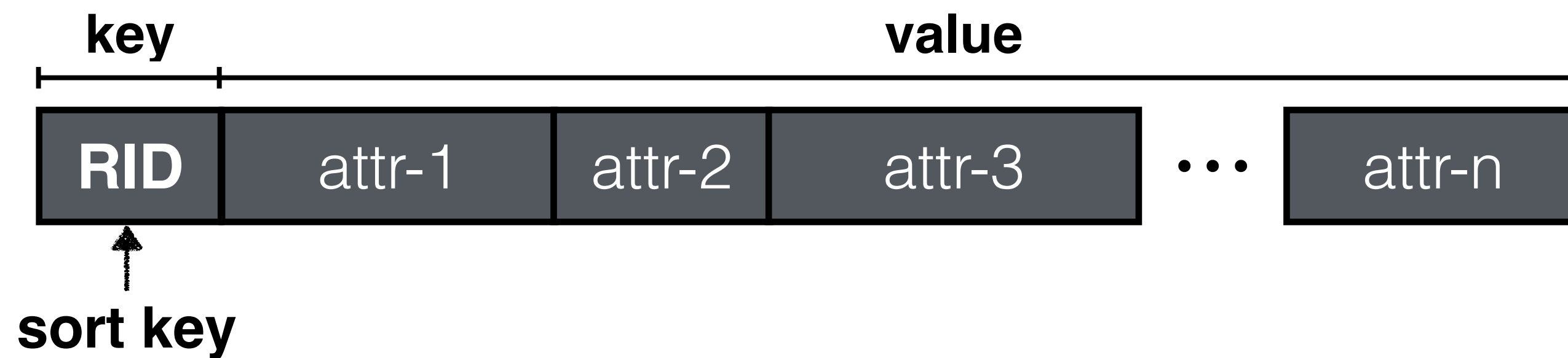
delete all entries older than: TS_x



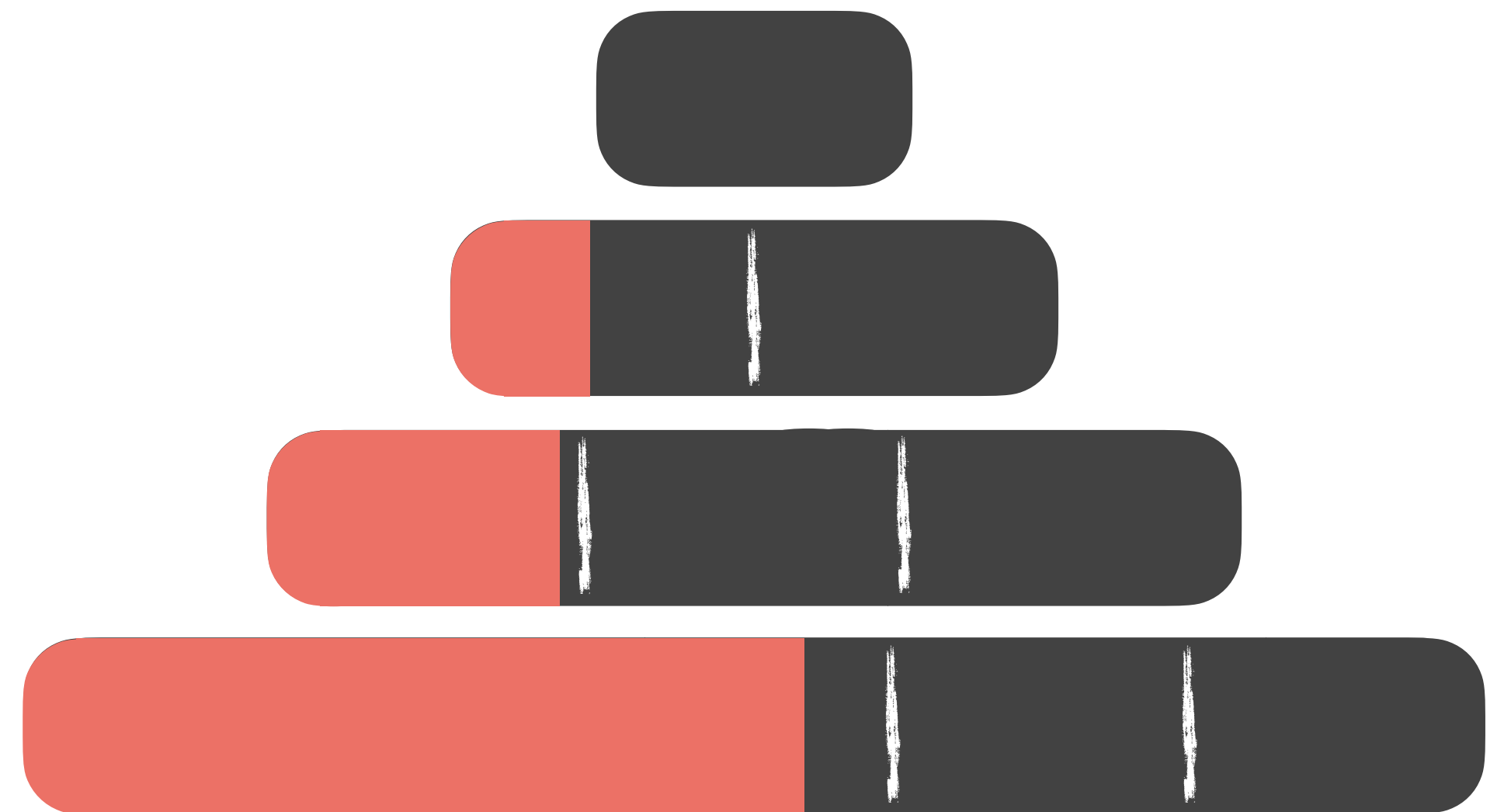
sort key = delete key

Realizing Retention-Based Deletes

delete all entries older than: TS_x

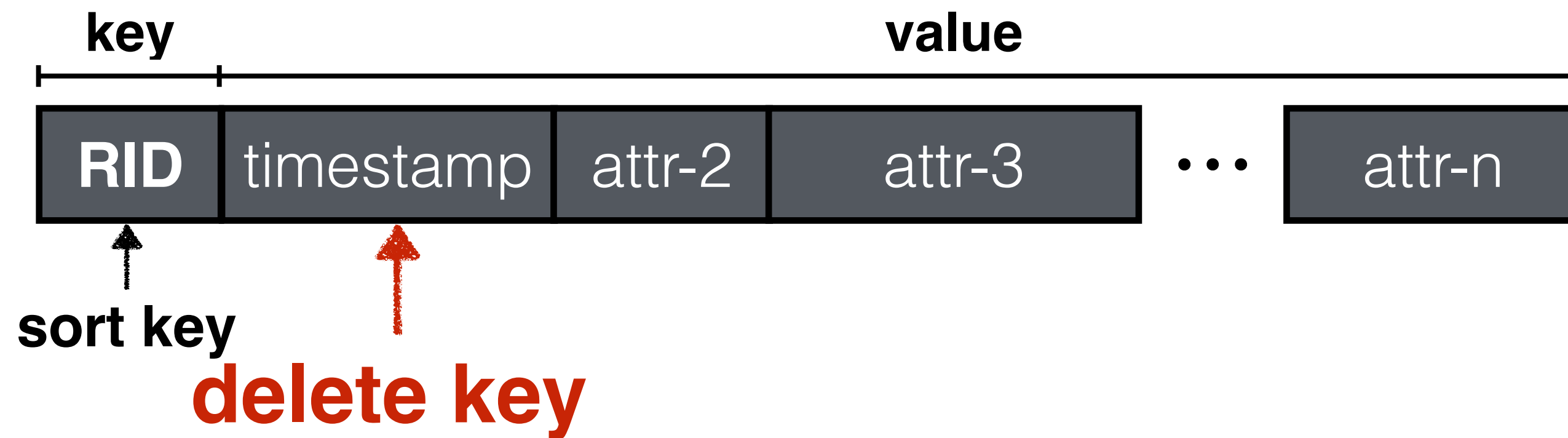


sort key = delete key

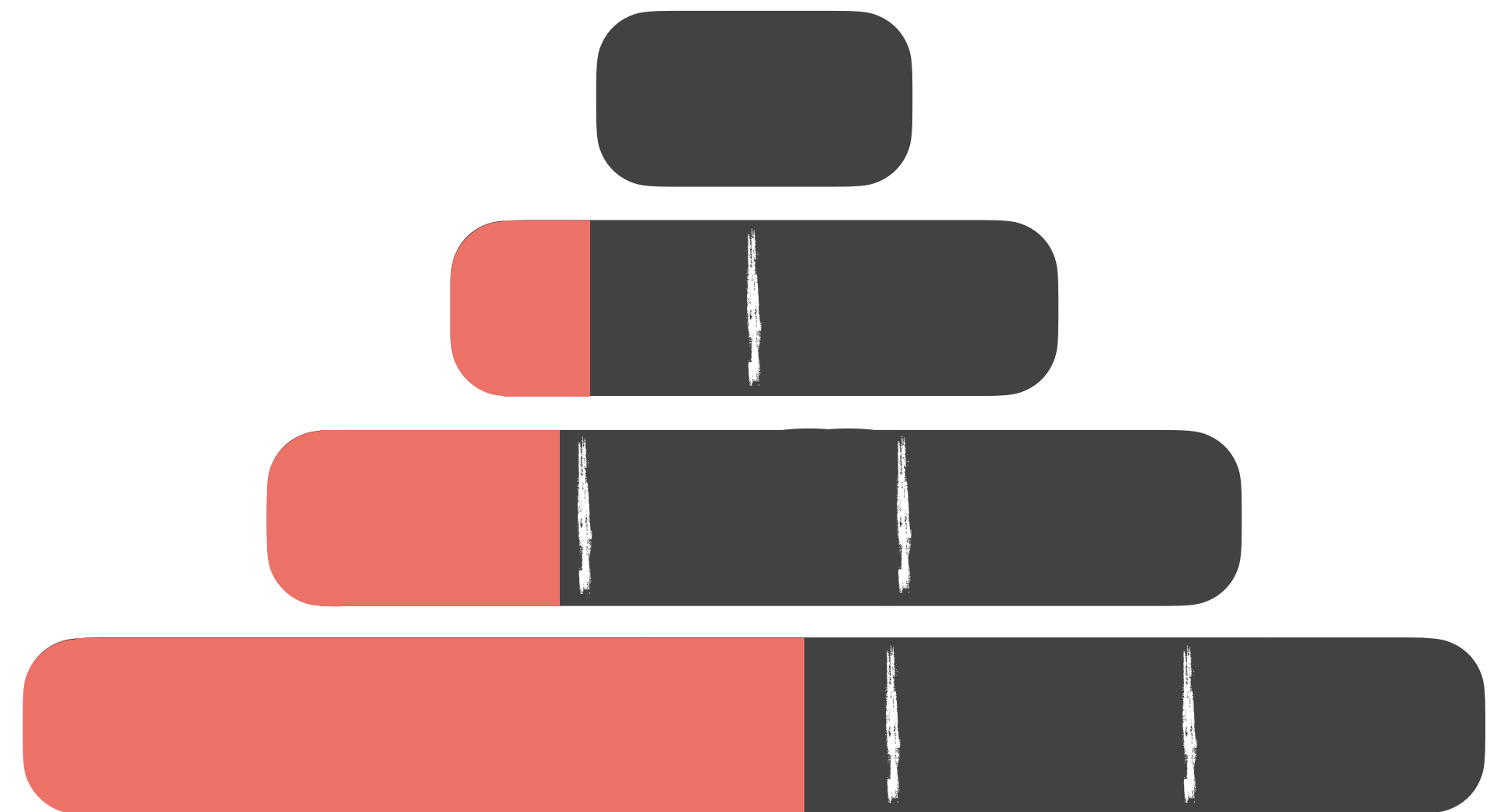


Realizing Retention-Based Deletes

delete all entries older than: TS_x

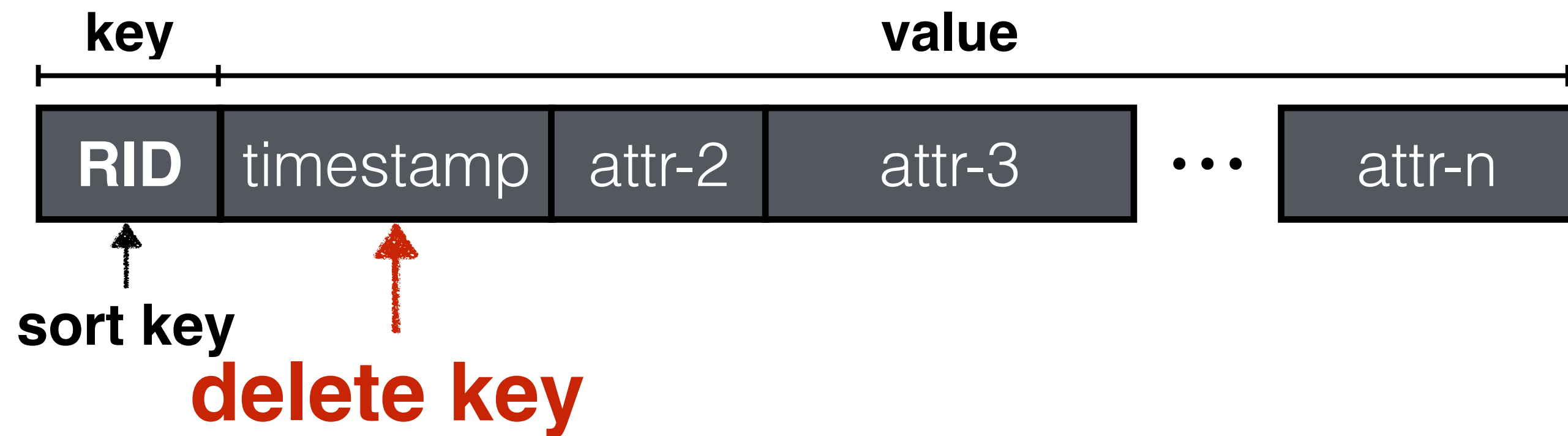


sort key \neq delete key

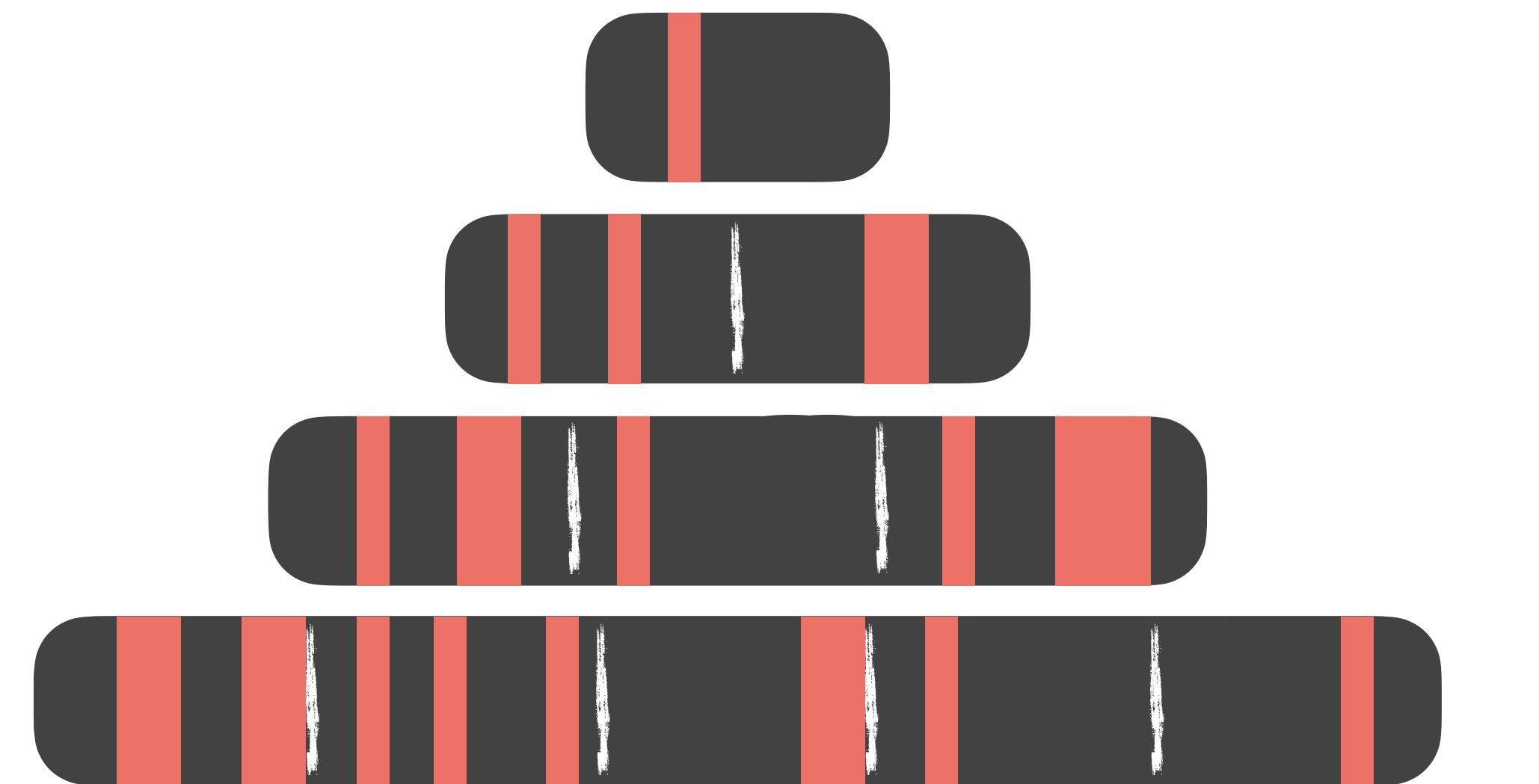


Realizing Retention-Based Deletes

delete all entries older than: TS_x



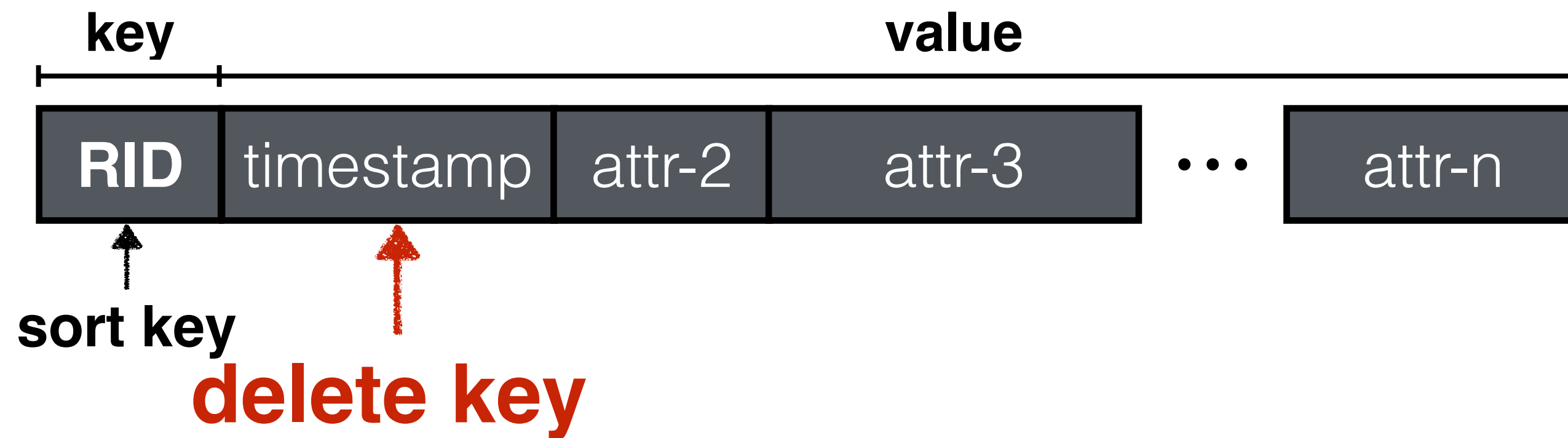
sort key \neq delete key



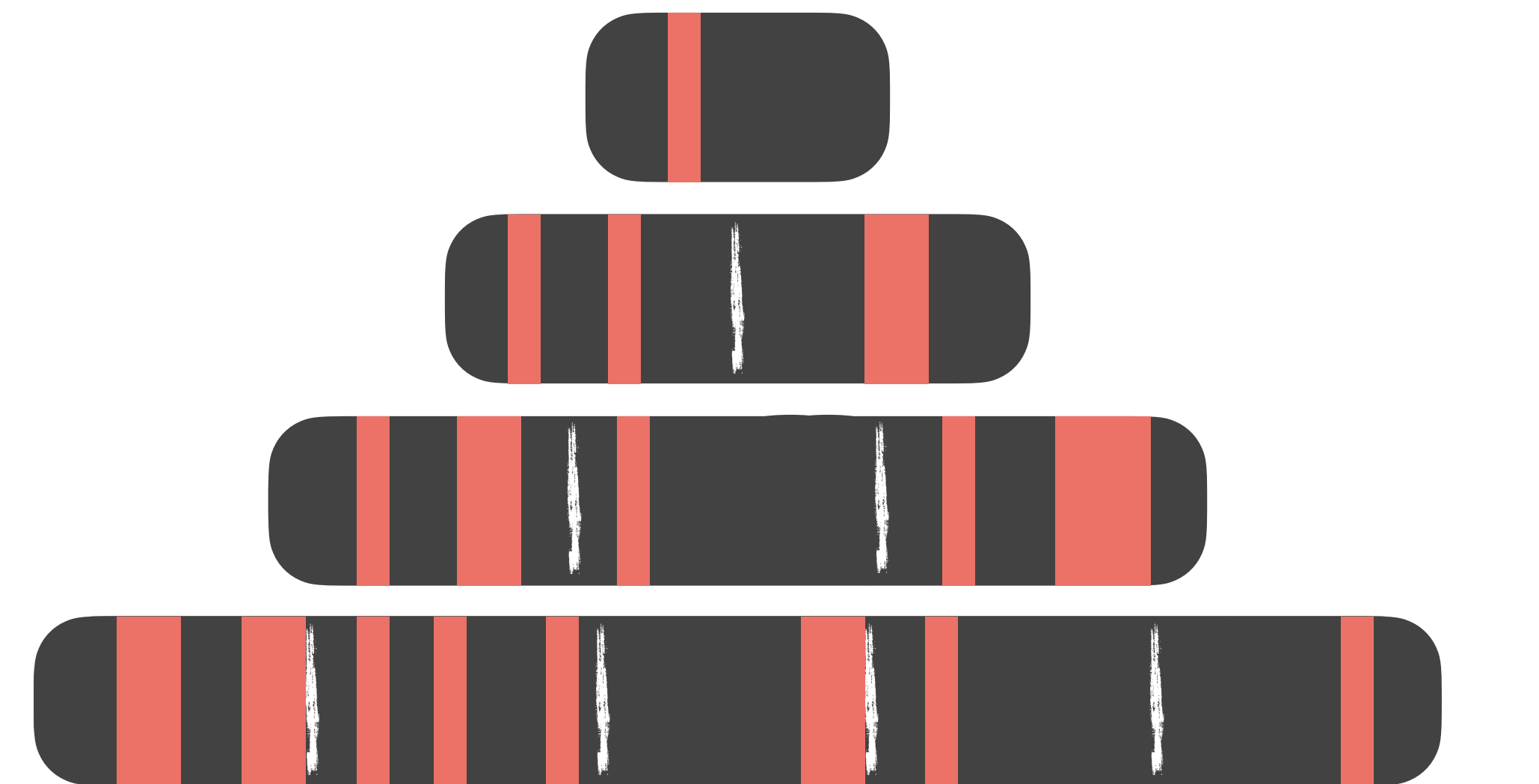
scattered occurrences

Realizing Retention-Based Deletes

delete all entries older than: TS_x



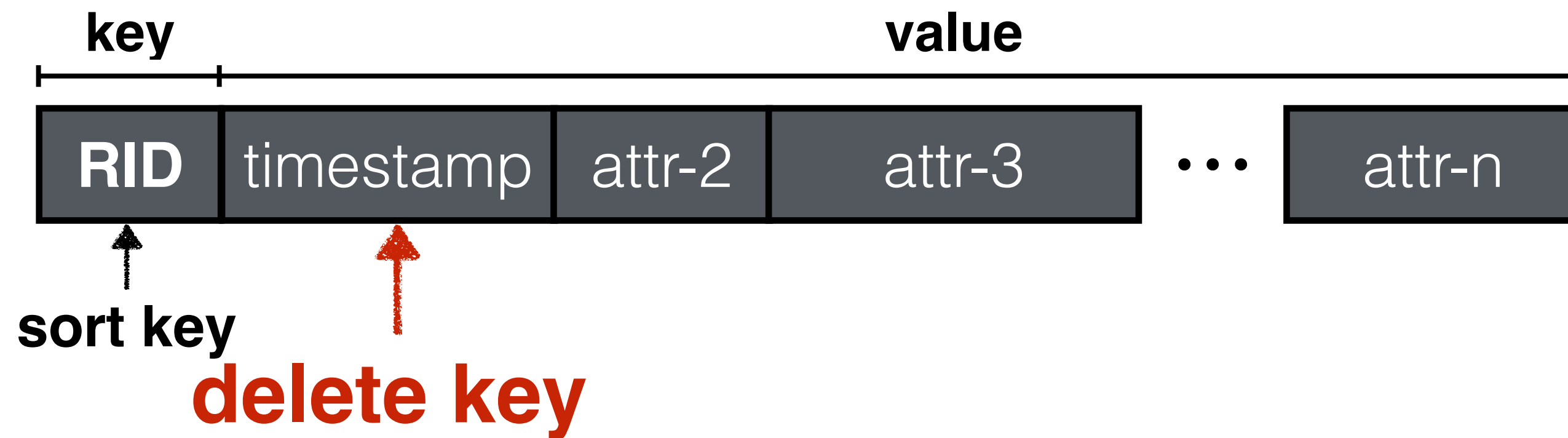
sort key \neq delete key



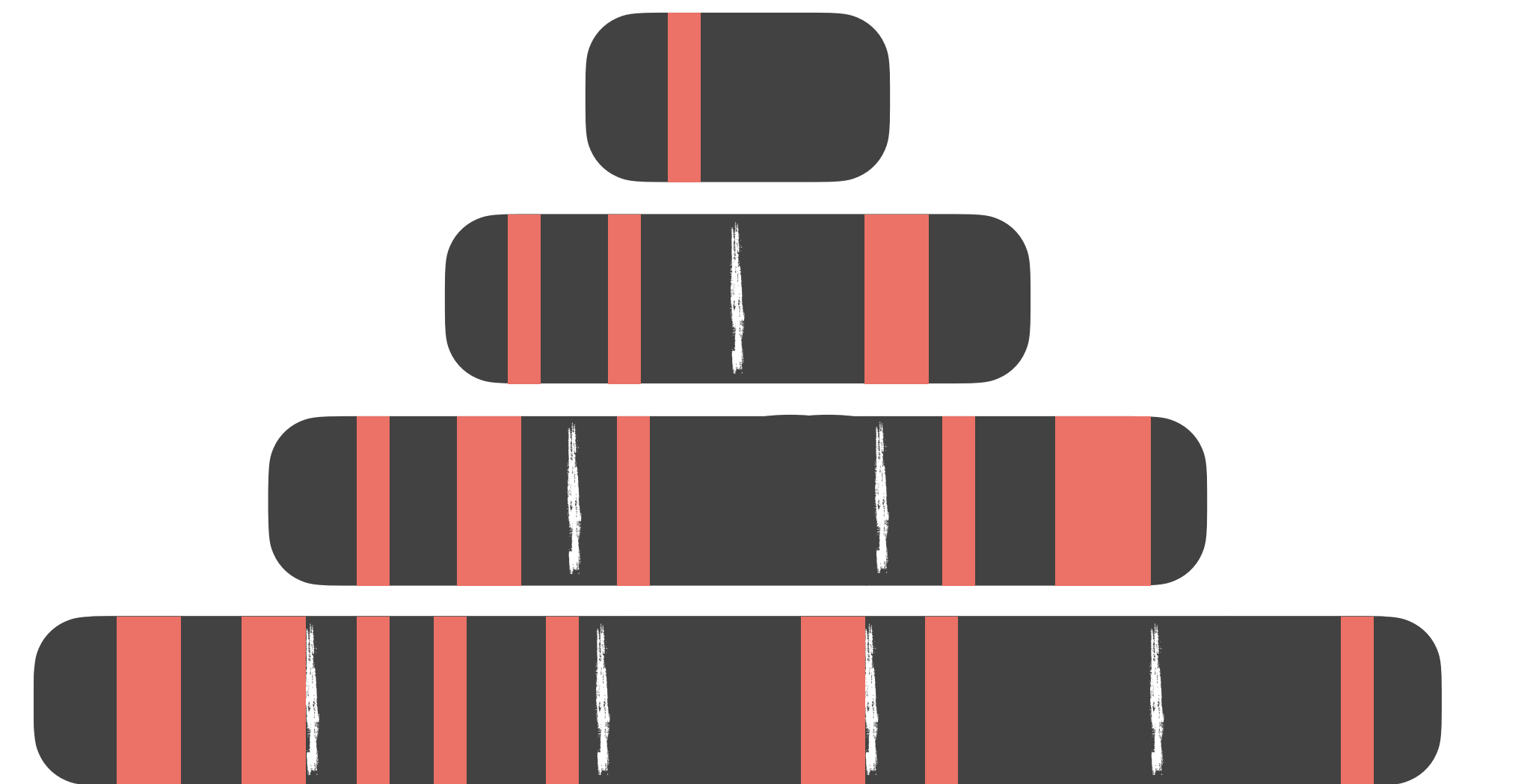
scattered occurrences

Realizing Retention-Based Deletes

delete all entries older than: TS_x



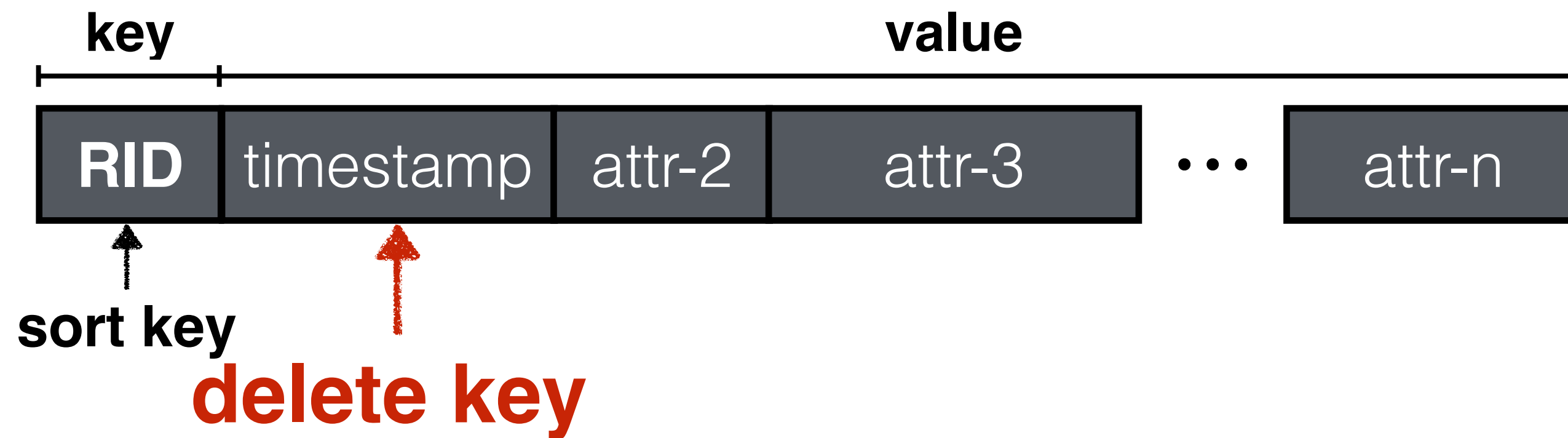
sort key \neq delete key



scattered occurrences

Realizing Retention-Based Deletes

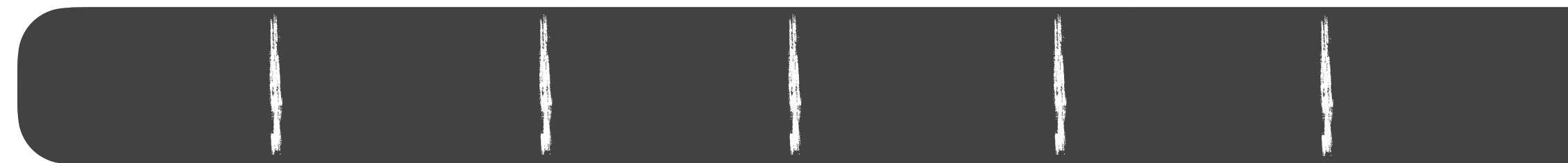
delete all entries older than: TS_x



sort key \neq delete key

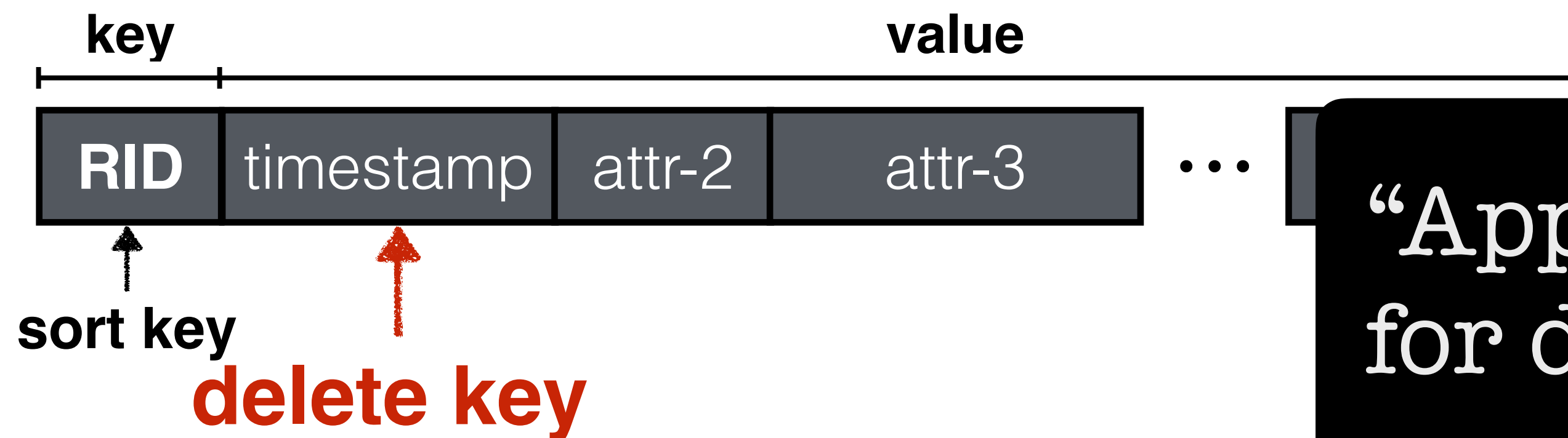
latency spikes

superfluous I/Os



Realizing Retention-Based Deletes

delete all entries older than: TS_x



sort key \neq delete key

“Applications have requirements for deletes every day. E.g., they may keep data for 30 days, ... effectively purging 1/30 of the database every day.

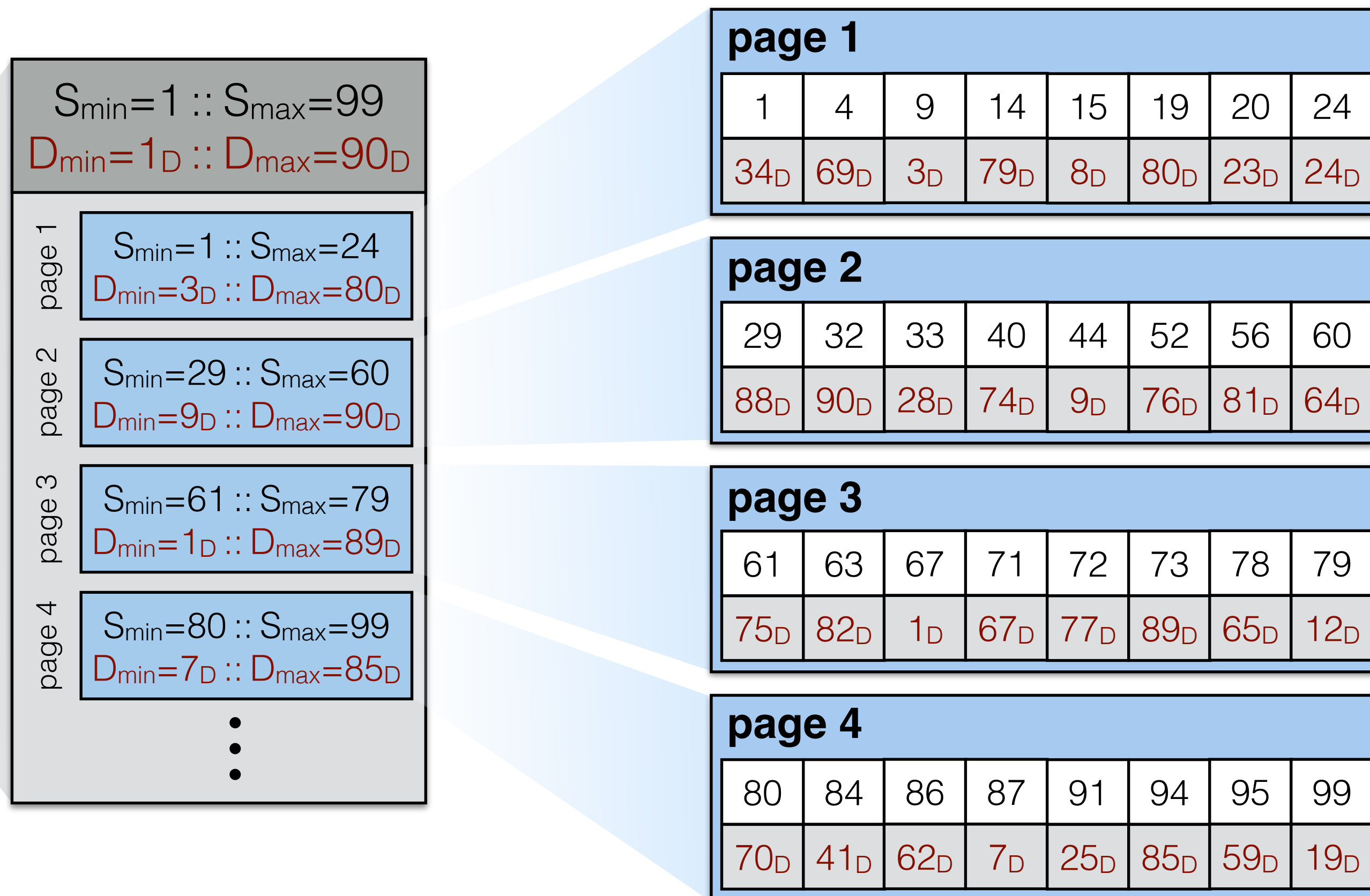
This induces performance pains!”

Realizing Retention-Based Deletes

delete all entries older than $\leq 65_D$



file

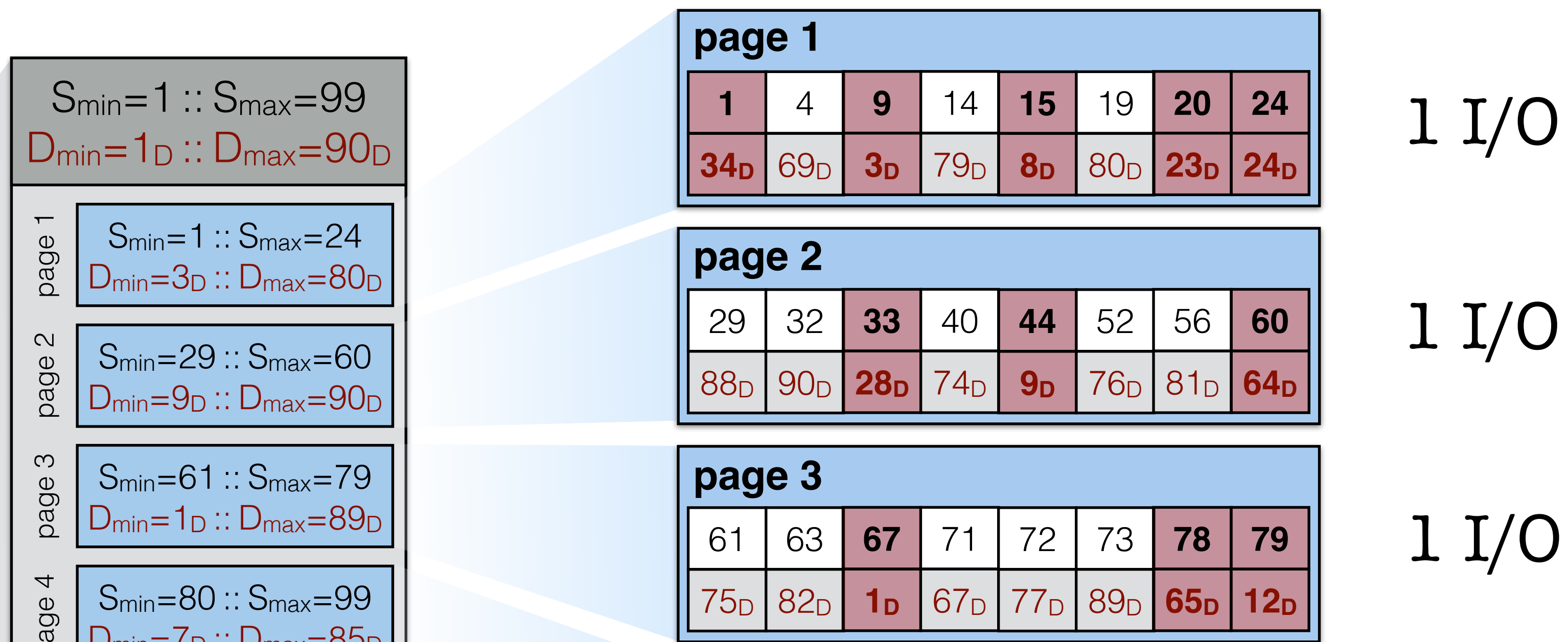


Realizing Retention-Based Deletes

delete all entries older than $\leq 65_D$



file



Intuition: Data Layout holds the key!

Realizing **Retention-Based Deletes**



KiWi

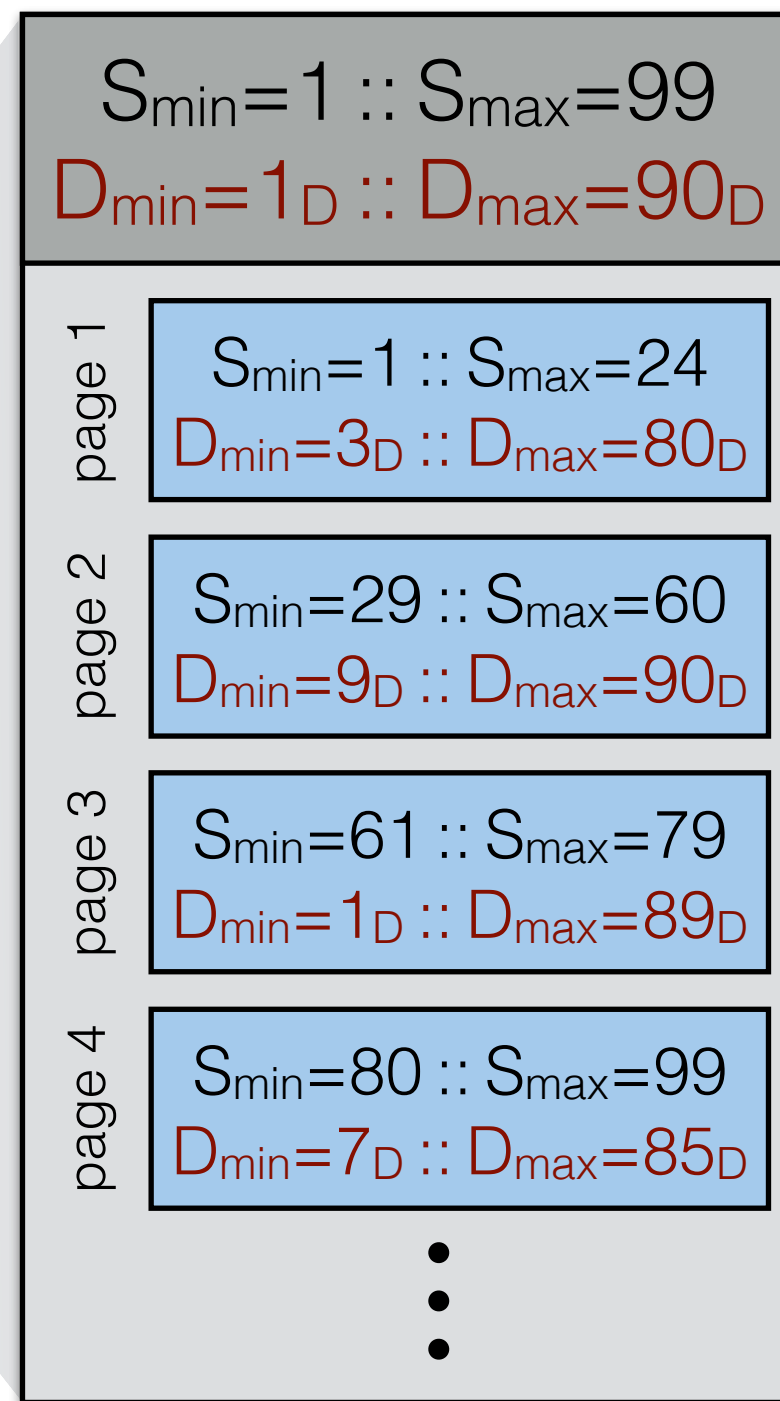
Key Weaving storage layout

Key Weaving storage layout

delete all entries older than $\leq 65_D$

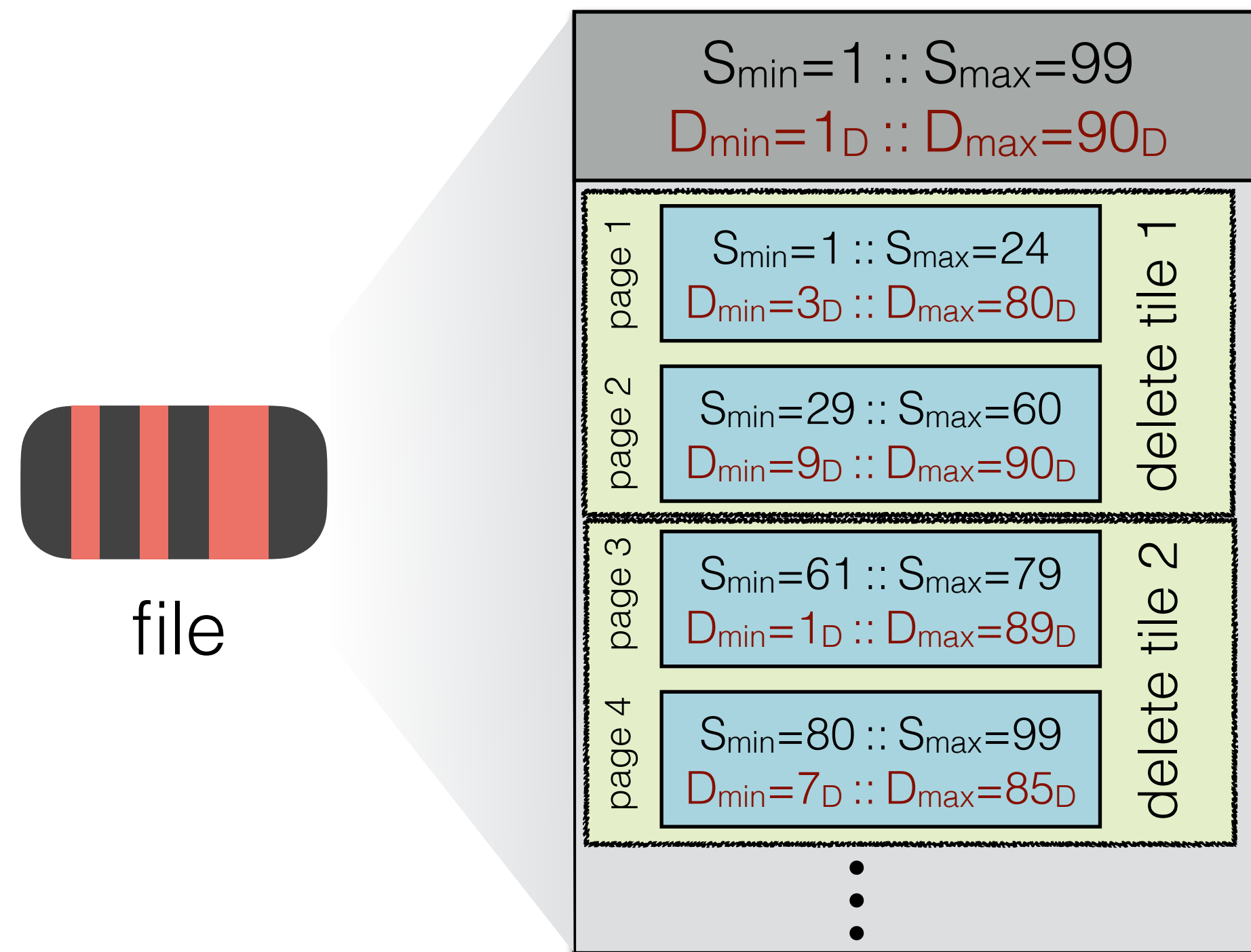


file



Key Weaving storage layout

delete all entries older than $\leq 65D$



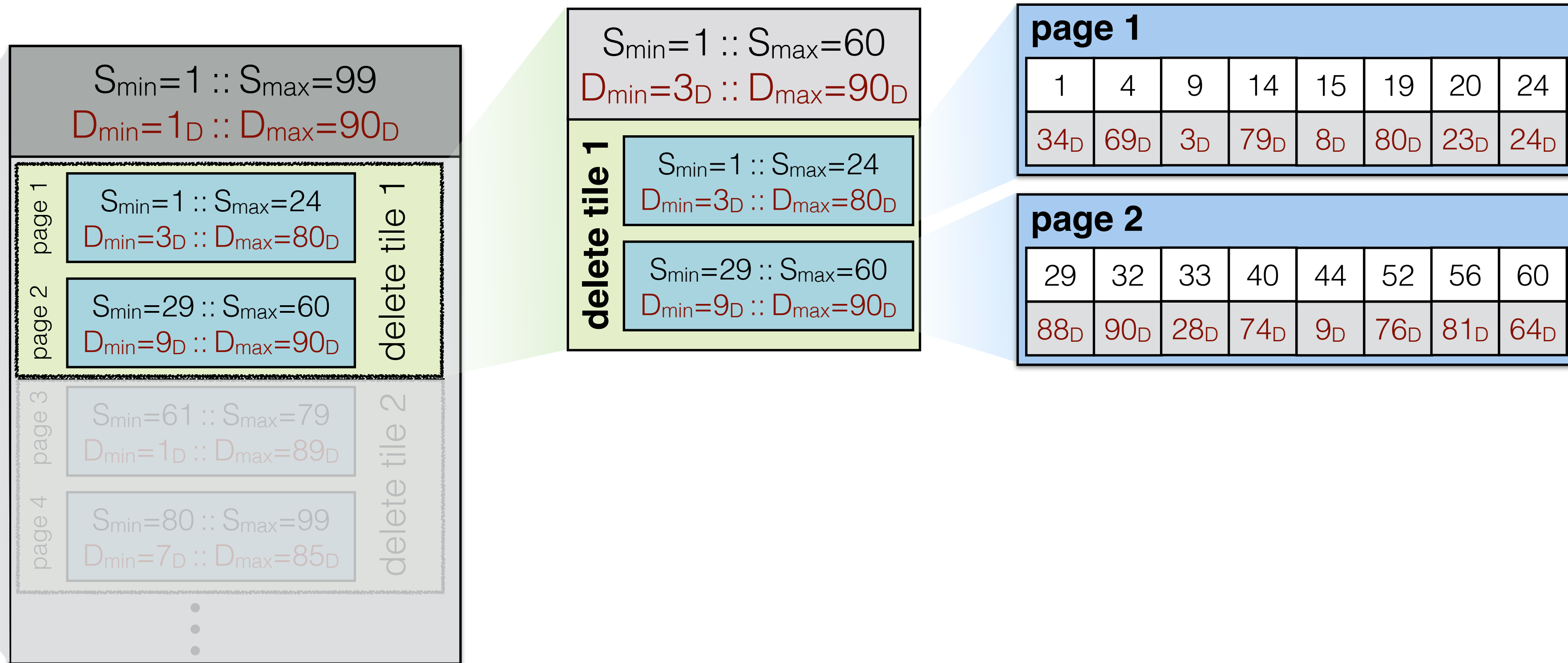
partitioned on S

Key Weaving storage layout

delete all entries older than $\leq 65_D$



file



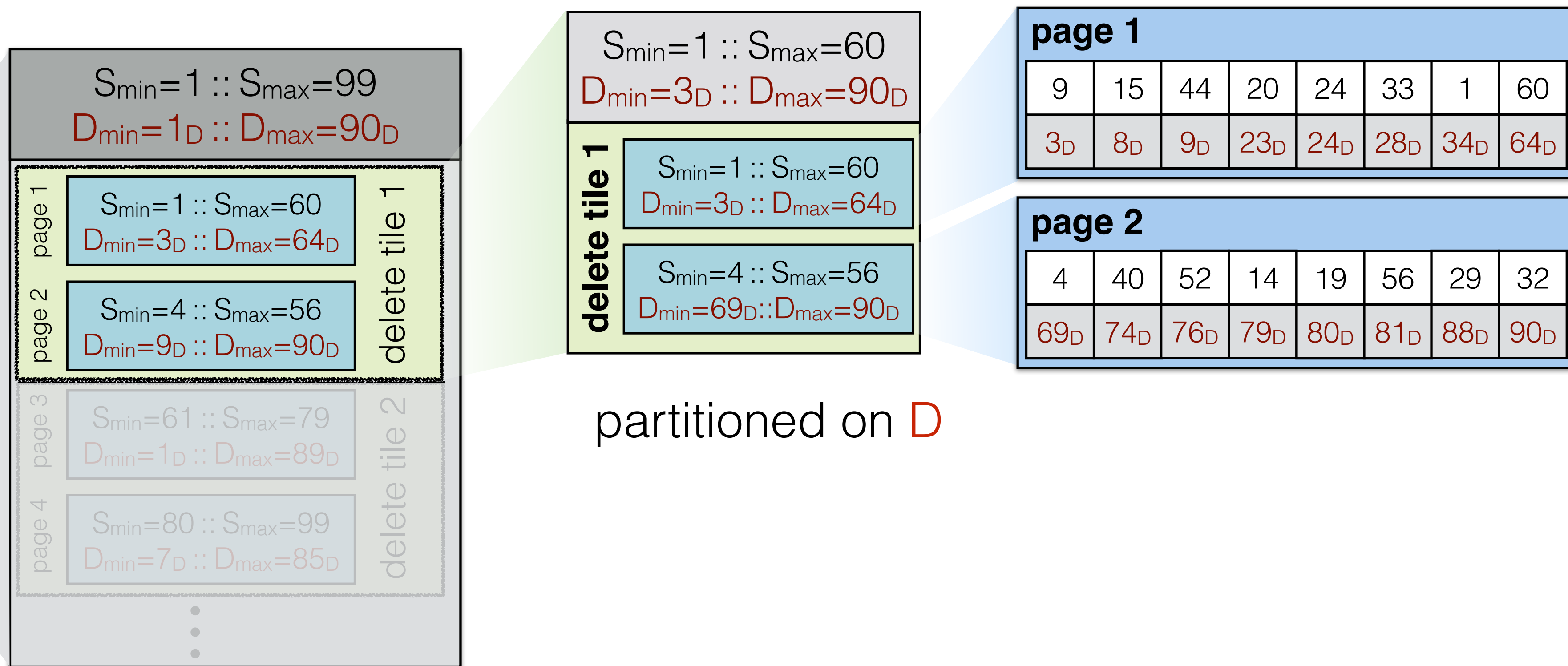
partitioned on S

Key Weaving storage layout

delete all entries older than $\leq 65D$



file



partitioned on S

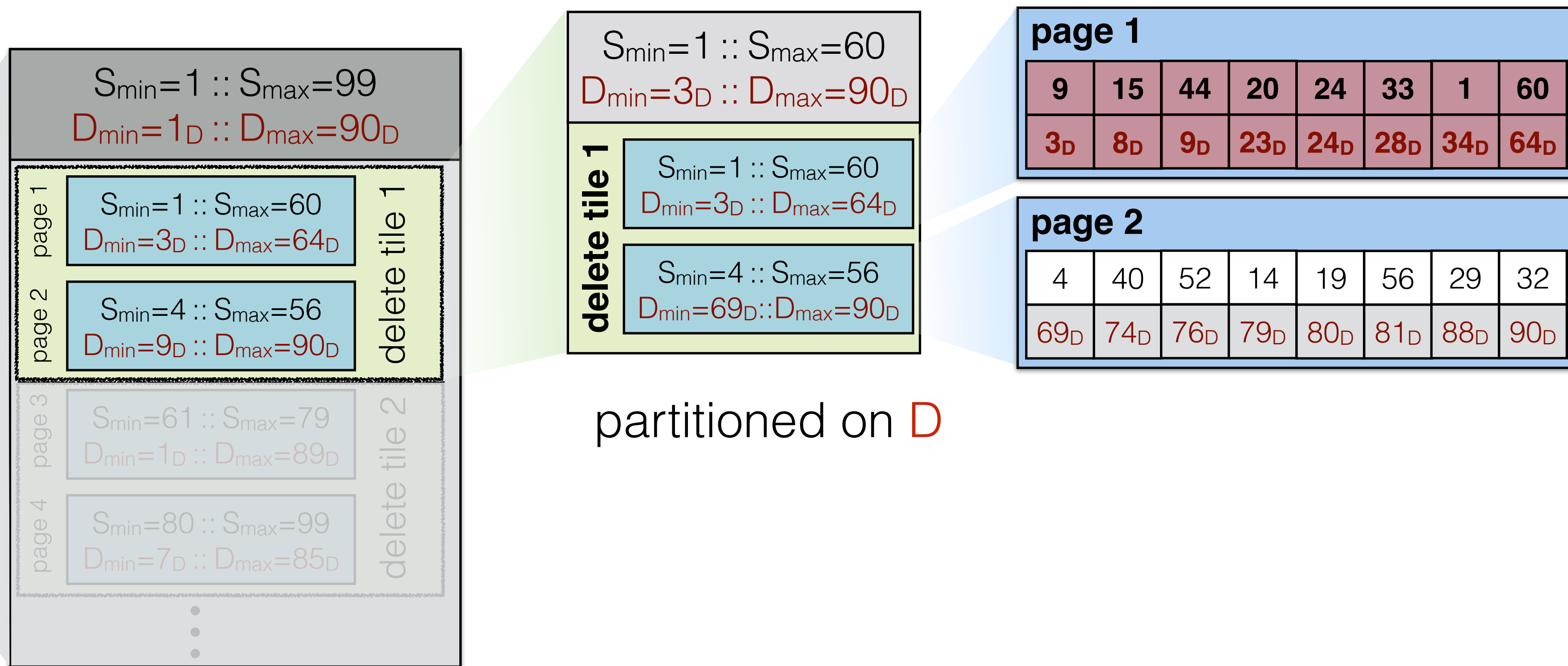
partitioned on D

Key Weaving storage layout

delete all entries older than $\leq 65D$



file



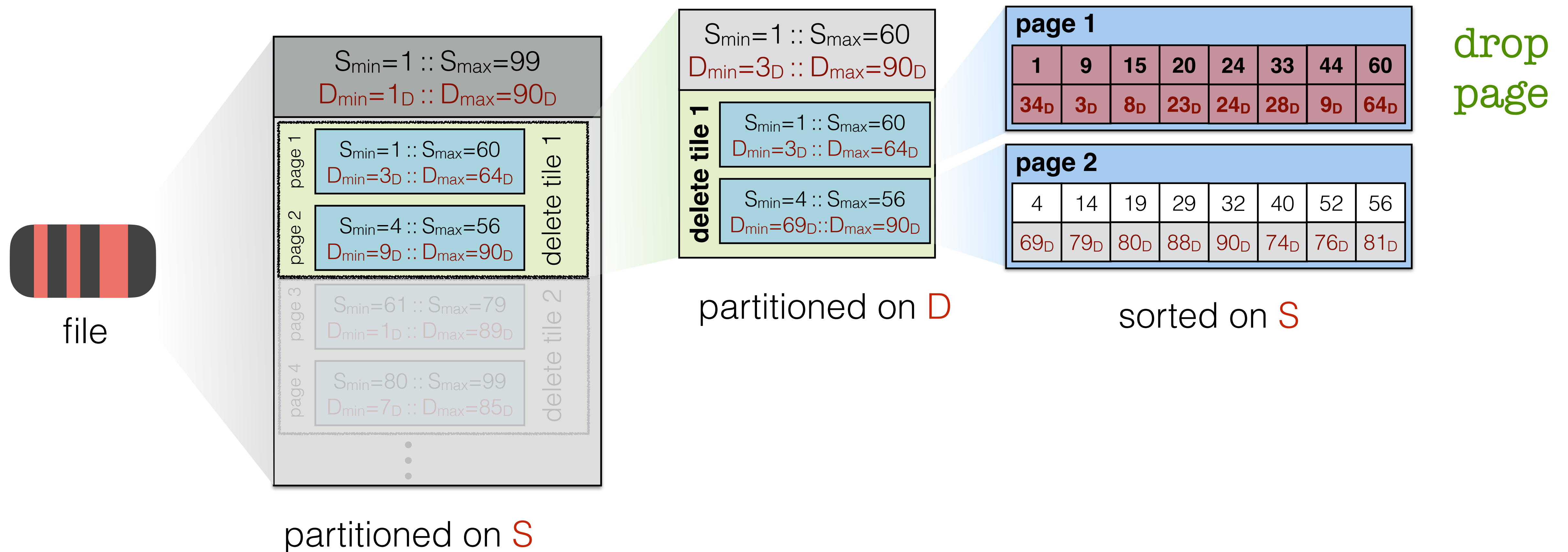
drop
page

partitioned on D

partitioned on S

Key Weaving storage layout

delete all entries older than $\leq 65_D$

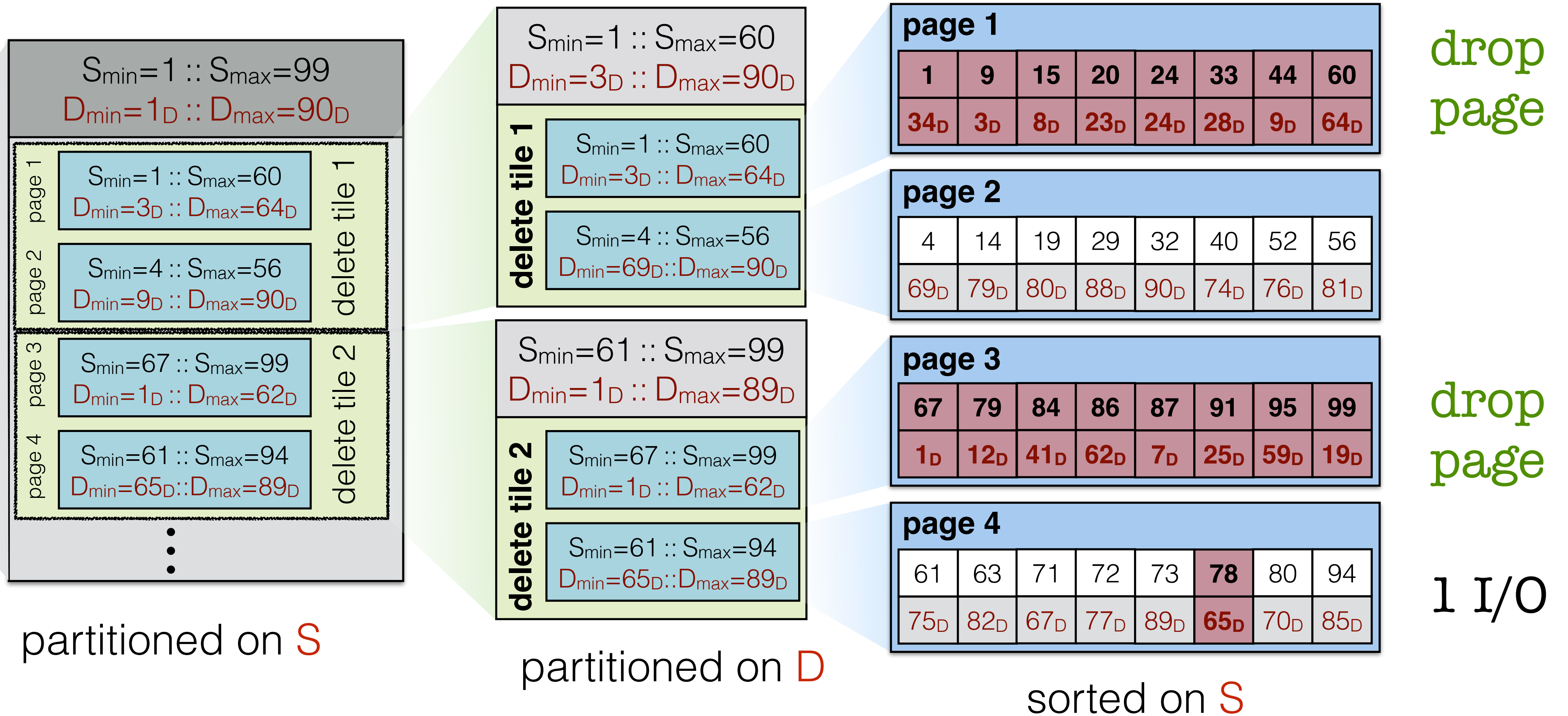


Key Weaving storage layout

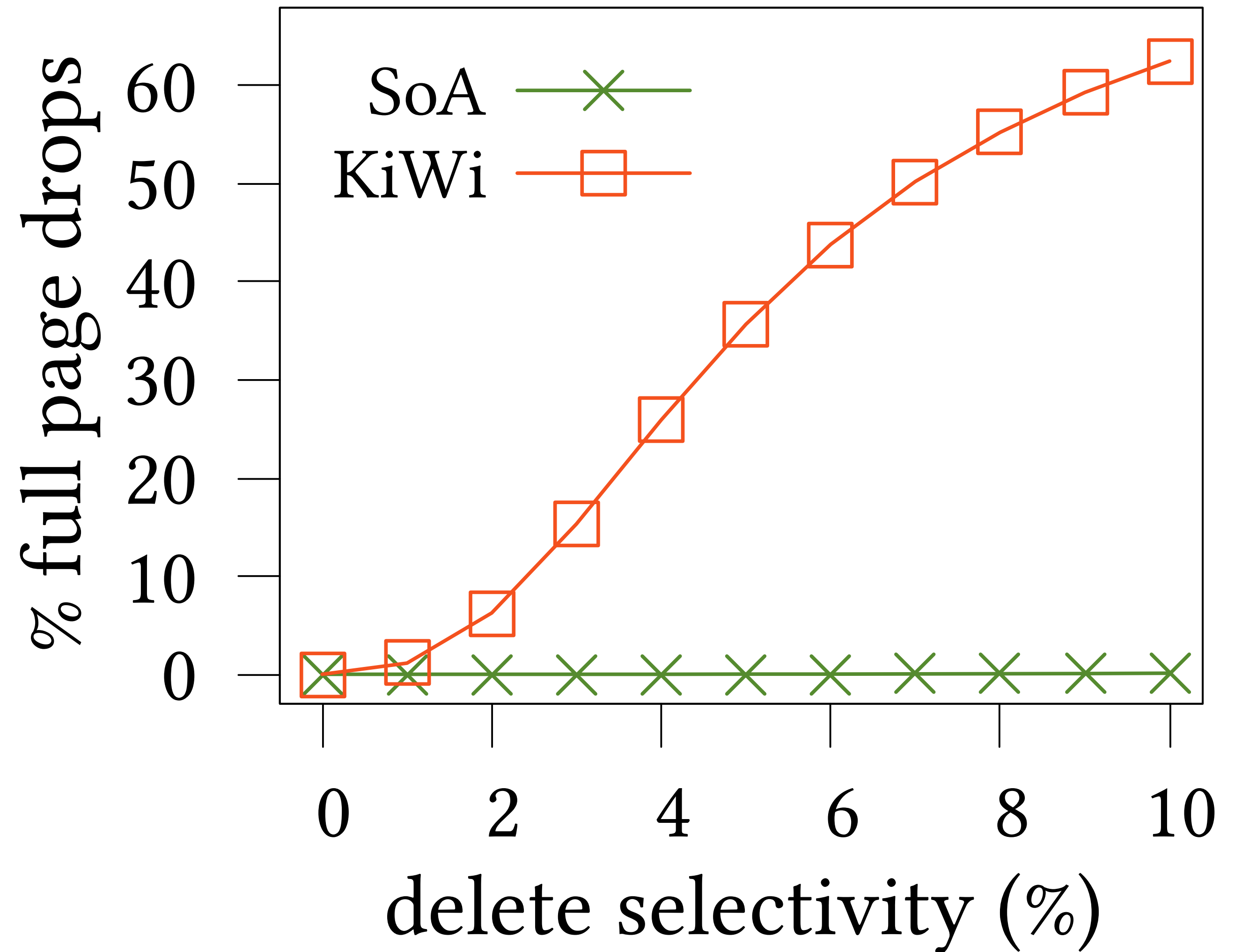
delete all entries older than $\leq 65_D$



file



5M entries, buffer = file = 256 pages, T=10



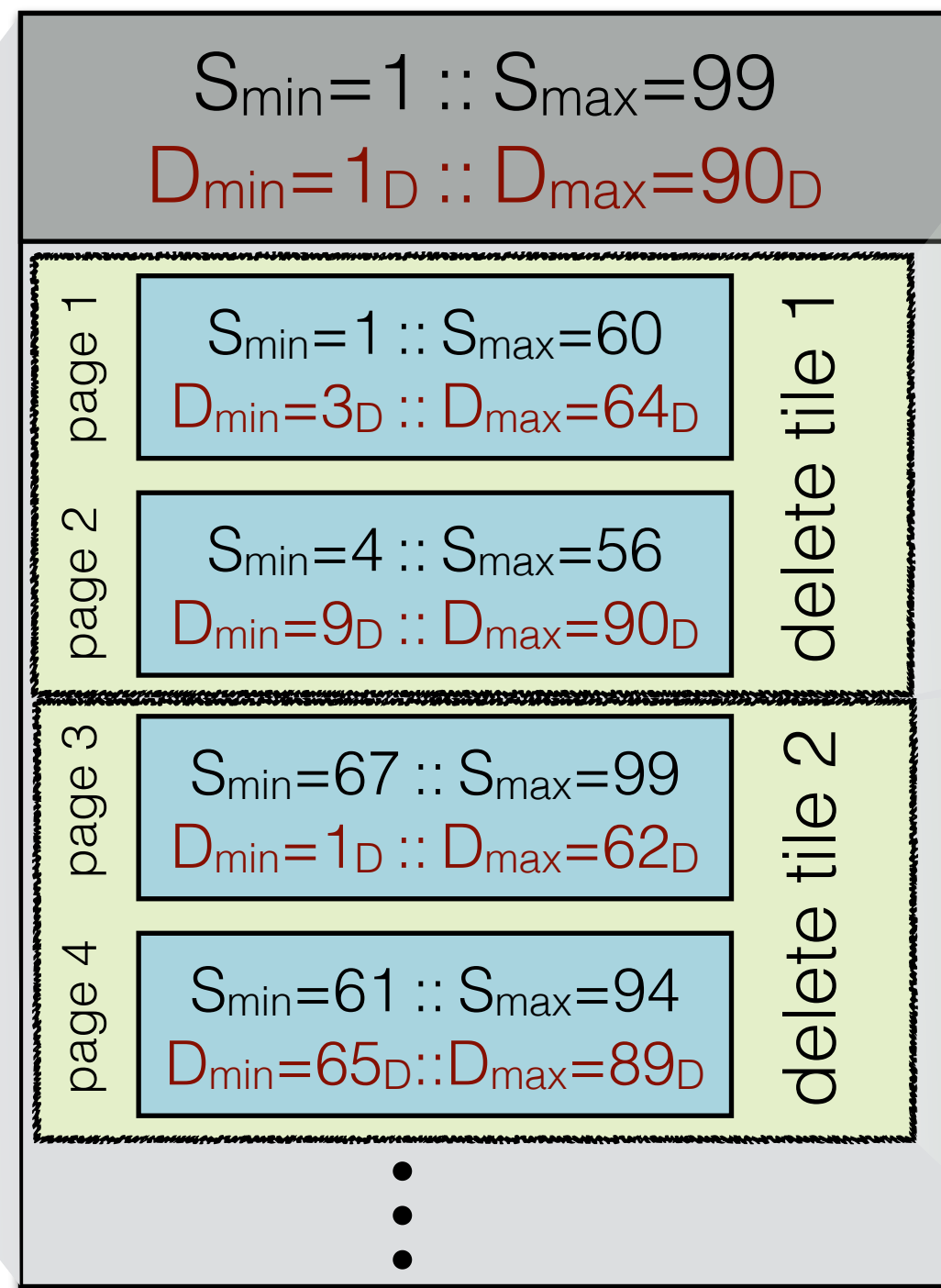
superior delete performance

up to 2.5x

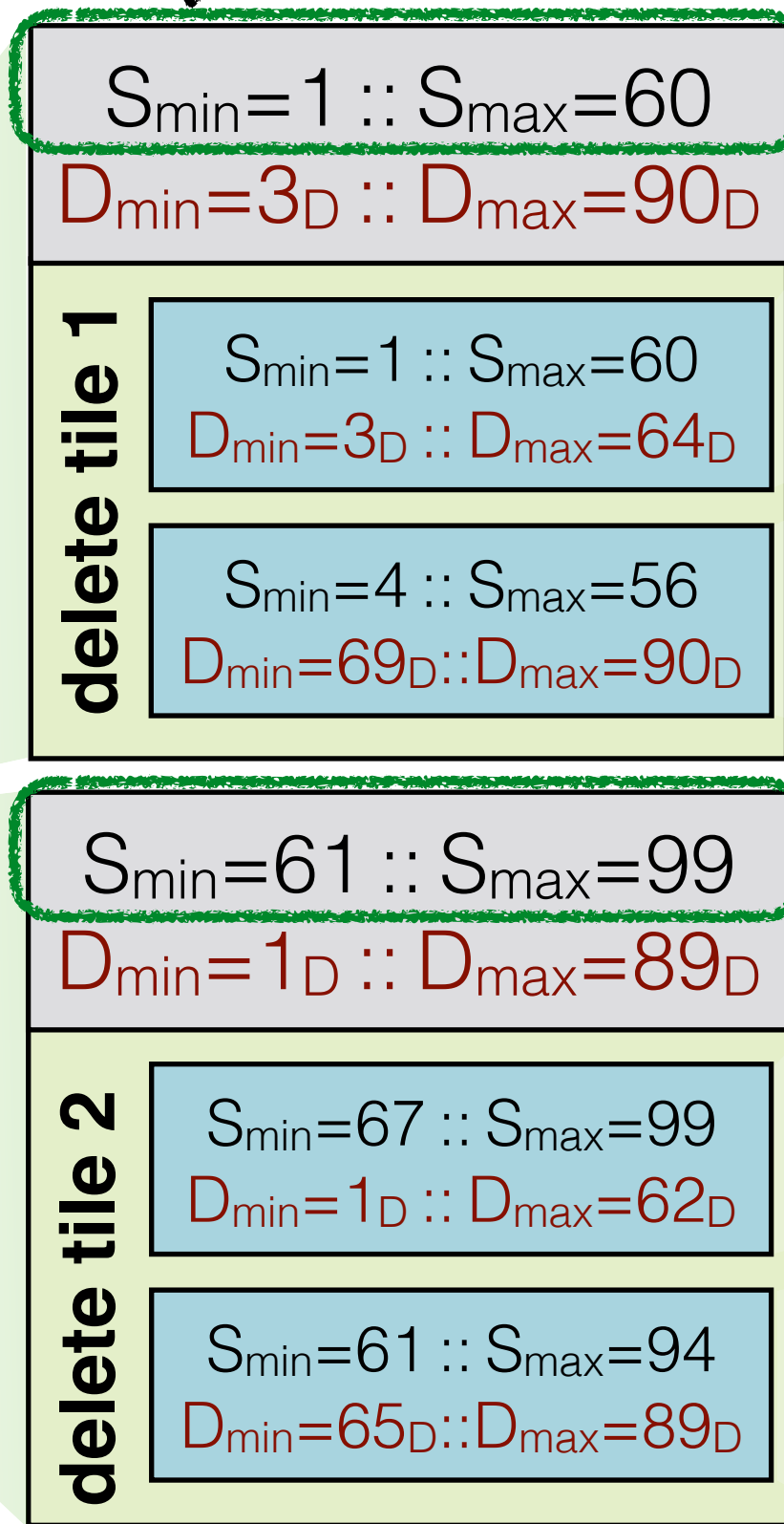
Key Weaving storage layout



file



get(14)



page 1

| | | | | | | | |
|-----------------|----------------|----------------|-----------------|-----------------|-----------------|----------------|-----------------|
| 1 | 9 | 15 | 20 | 24 | 33 | 44 | 60 |
| 34 _D | 3 _D | 8 _D | 23 _D | 24 _D | 28 _D | 9 _D | 64 _D |

page 2

| | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| 4 | 14 | 19 | 29 | 32 | 40 | 52 | 56 |
| 69 _D | 79 _D | 80 _D | 88 _D | 90 _D | 74 _D | 76 _D | 81 _D |

page 3

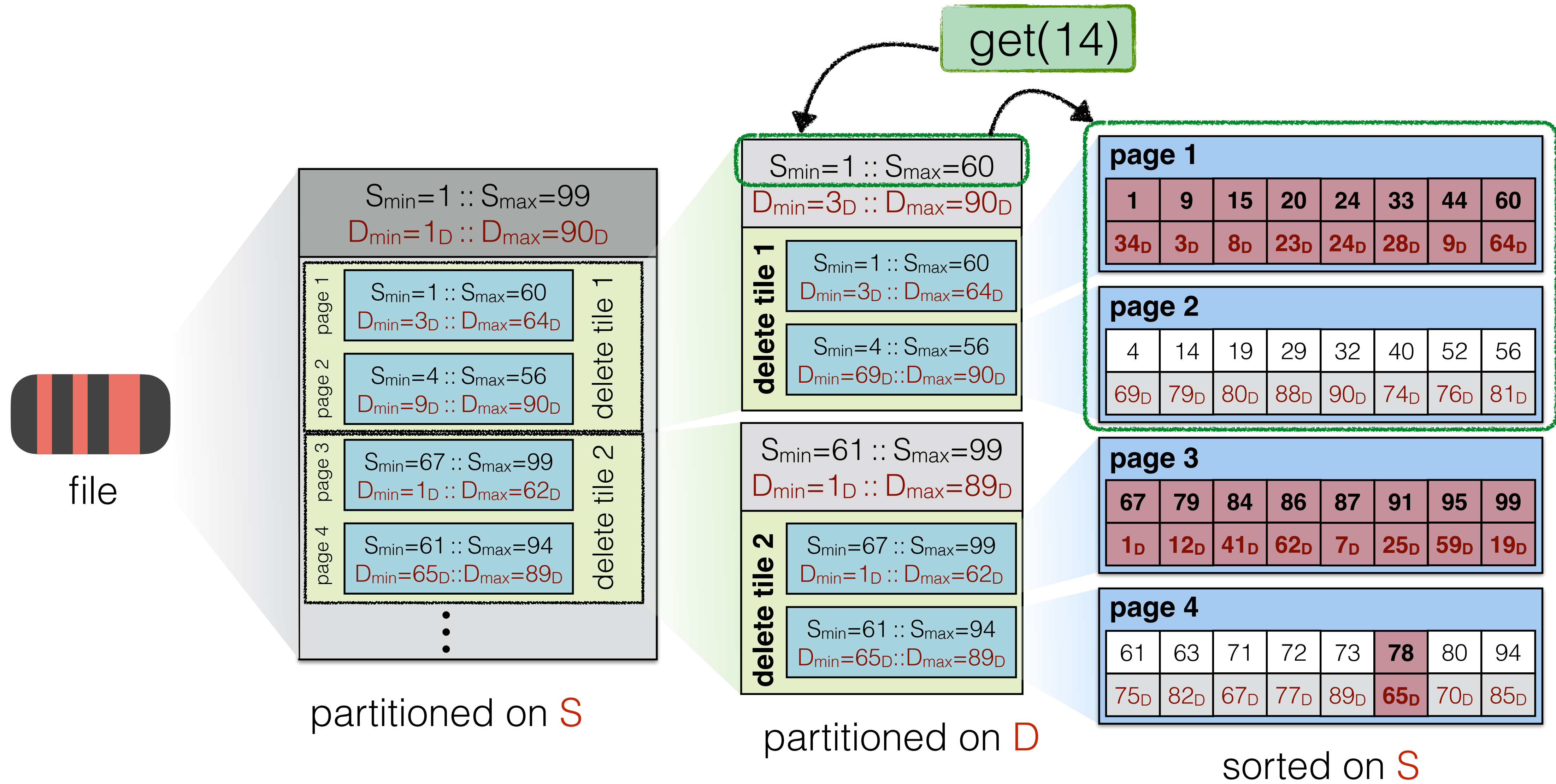
| | | | | | | | |
|----------------|-----------------|-----------------|-----------------|----------------|-----------------|-----------------|-----------------|
| 67 | 79 | 84 | 86 | 87 | 91 | 95 | 99 |
| 1 _D | 12 _D | 41 _D | 62 _D | 7 _D | 25 _D | 59 _D | 19 _D |

page 4

| | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| 61 | 63 | 71 | 72 | 73 | 78 | 80 | 94 |
| 75 _D | 82 _D | 67 _D | 77 _D | 89 _D | 65 _D | 70 _D | 85 _D |

sorted on S

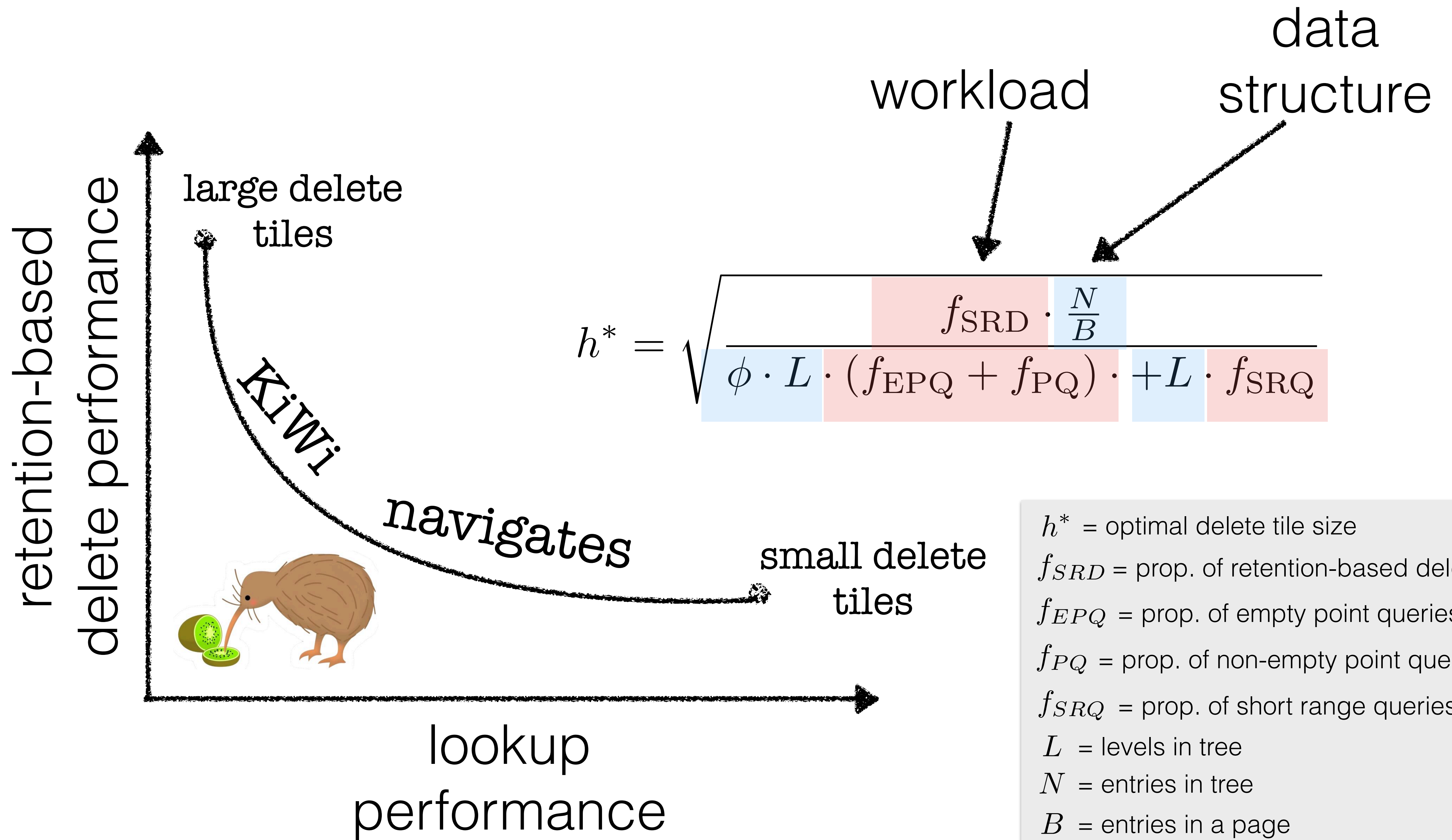
Key Weaving storage layout



| | |
|--|--|
| $S_{min}=1 :: S_{max}=99$ $D_{min}=1_D :: D_{max}=90_D$ | |
| page 1 | $S_{min}=1 :: S_{max}=60$ $D_{min}=3_D :: D_{max}=64_D$ |
| page 2 | $S_{min}=4 :: S_{max}=56$ $D_{min}=9_D :: D_{max}=90_D$ |
| page 3 | $S_{min}=67 :: S_{max}=99$ $D_{min}=1_D :: D_{max}=62_D$ |
| page 4 | $S_{min}=61 :: S_{max}=94$ $D_{min}=65_D :: D_{max}=89_D$ |
| ⋮ | |

| | |
|---|--|
| $S_{min}=1 :: S_{max}=60$ $D_{min}=3_D :: D_{max}=90_D$ | |
| delete tile 1 | $S_{min}=1 :: S_{max}=60$ $D_{min}=3_D :: D_{max}=64_D$ |
| | $S_{min}=4 :: S_{max}=56$ $D_{min}=69_D :: D_{max}=90_D$ |
| $S_{min}=61 :: S_{max}=99$ $D_{min}=1_D :: D_{max}=89_D$ | |
| delete tile 2 | $S_{min}=67 :: S_{max}=99$ $D_{min}=1_D :: D_{max}=62_D$ |
| | $S_{min}=61 :: S_{max}=94$ $D_{min}=65_D :: D_{max}=89_D$ |

| | | | | | | | |
|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| page 1 | | | | | | | |
| 1 | 9 | 15 | 20 | 24 | 33 | 44 | 60 |
| 34 _D | 3 _D | 8 _D | 23 _D | 24 _D | 28 _D | 9 _D | 64 _D |
| page 2 | | | | | | | |
| 4 | 14 | 19 | 29 | 32 | 40 | 52 | 56 |
| 69 _D | 79 _D | 80 _D | 88 _D | 90 _D | 74 _D | 76 _D | 81 _D |
| page 3 | | | | | | | |
| 67 | 79 | 84 | 86 | 87 | 91 | 95 | 99 |
| 1 _D | 12 _D | 41 _D | 62 _D | 7 _D | 25 _D | 59 _D | 19 _D |
| page 4 | | | | | | | |
| 61 | 63 | 71 | 72 | 73 | 78 | 80 | 94 |
| 75 _D | 82 _D | 67 _D | 77 _D | 89 _D | 65 _D | 70 _D | 85 _D |



h^* = optimal delete tile size
 f_{SRD} = prop. of retention-based deletes
 f_{EPQ} = prop. of empty point queries
 f_{PQ} = prop. of non-empty point queries
 f_{SRQ} = prop. of short range queries
 L = levels in tree
 N = entries in tree
 B = entries in a page
 ϕ = false positive rate of query filter

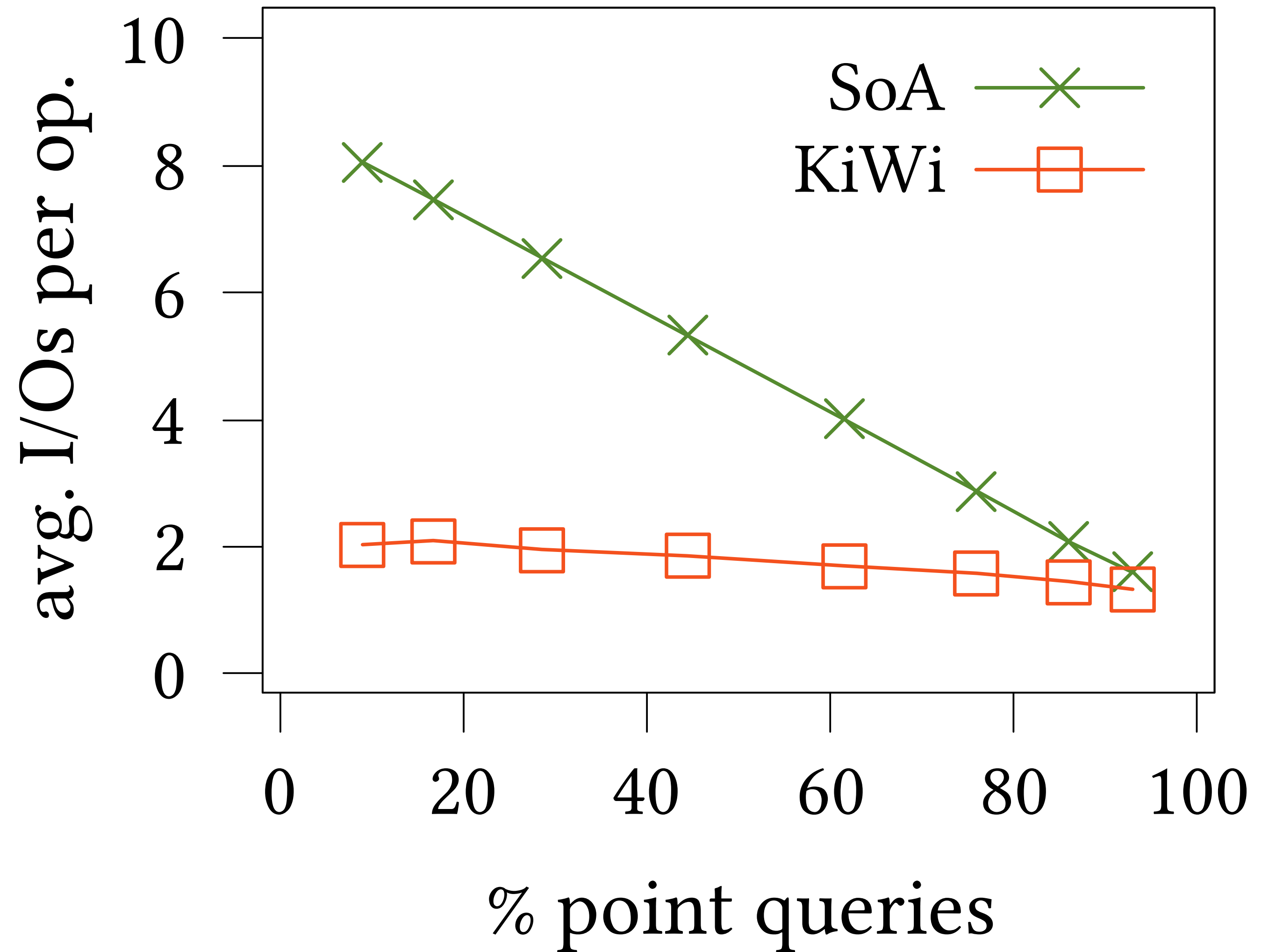
better overall performance

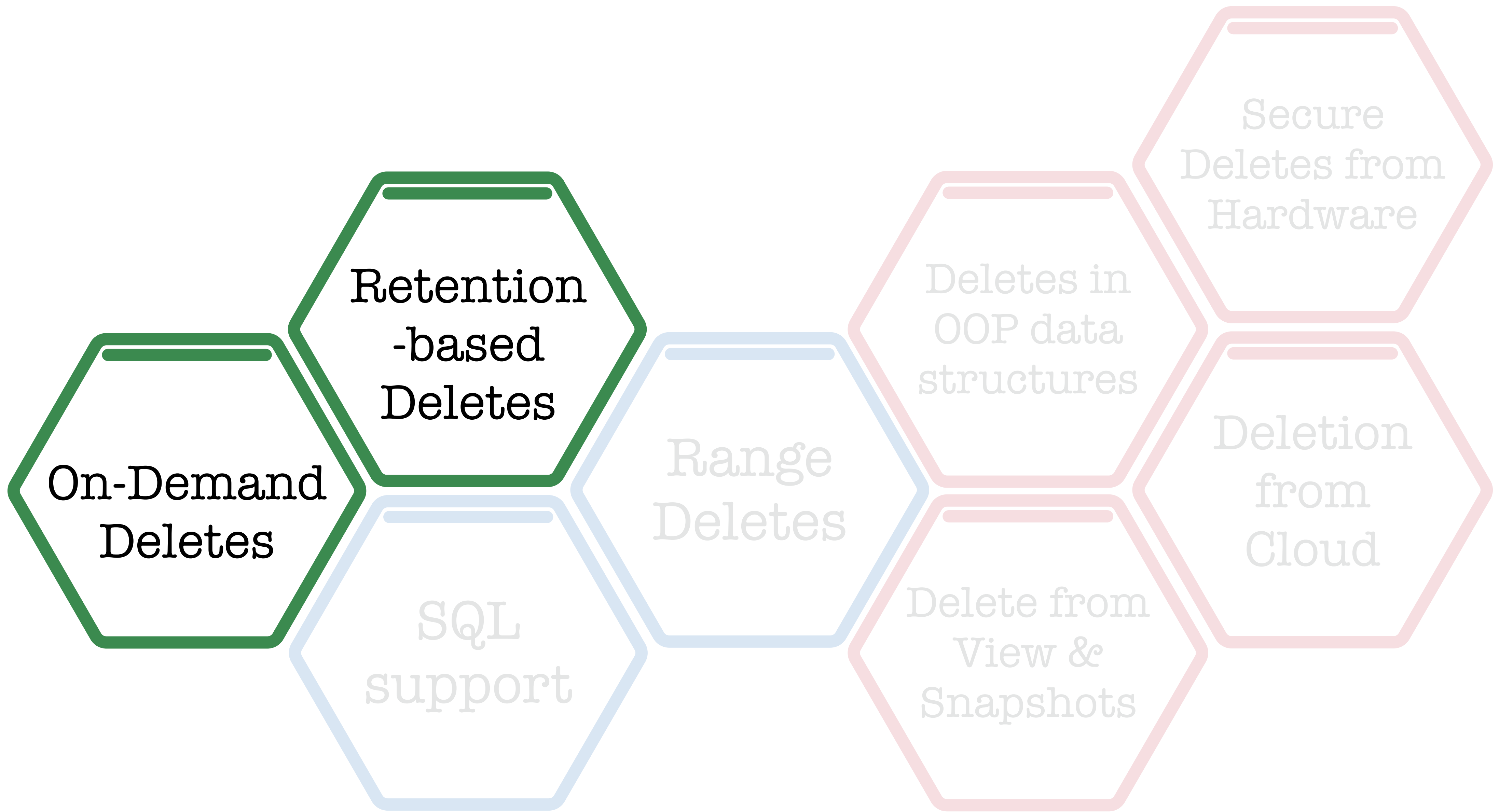
up to 4x

superior delete performance

up to 2.5x

5M entries, buffer = file = 256 pages, T=10





FADE



SQL
support

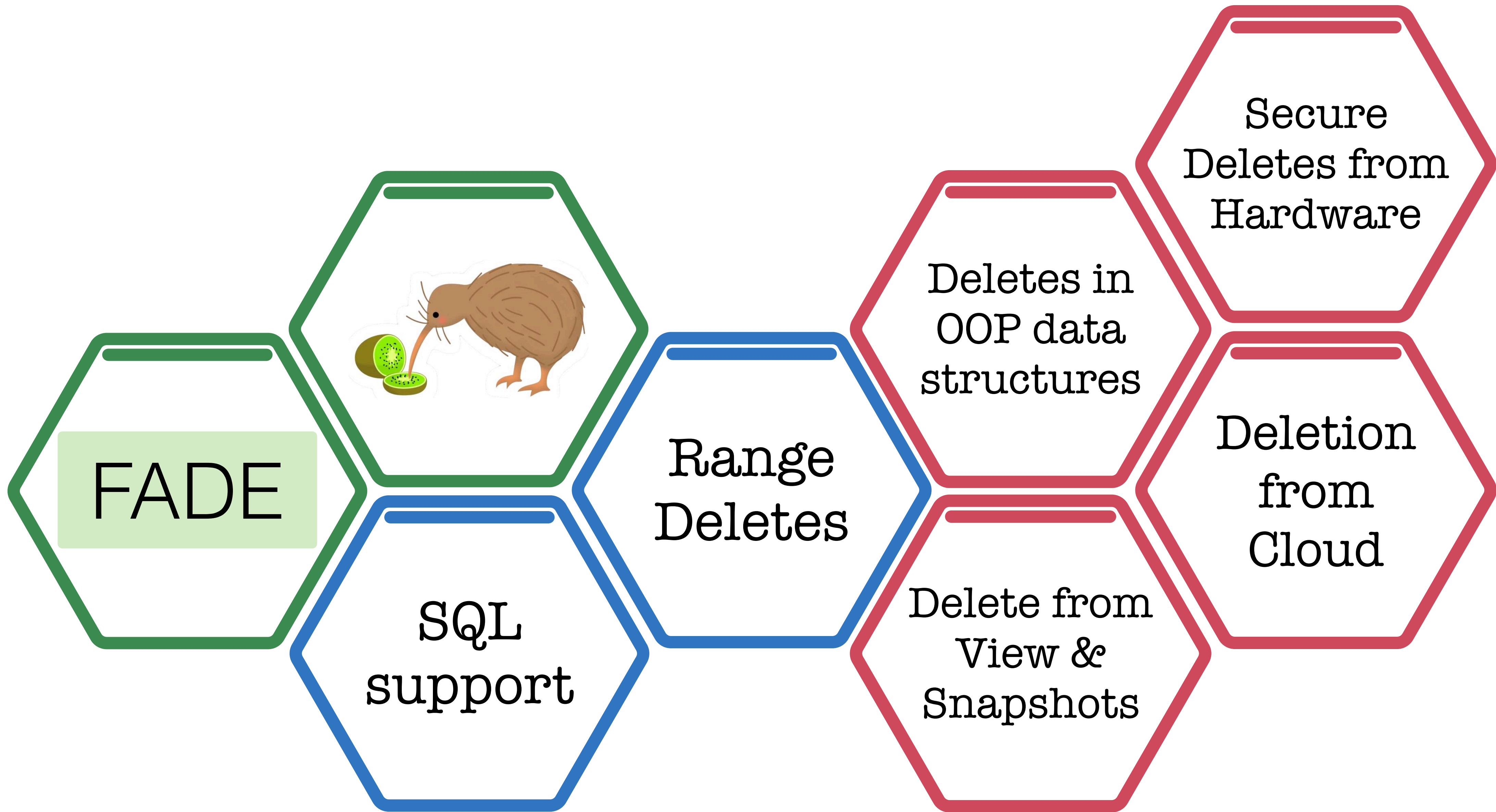
Range
Deletes

Deletes in
OOP data
structures

Delete from
View &
Snapshots

Secure
Deletes from
Hardware

Deletion
from
Cloud



COSI 167A

Advanced Data Systems

Class 12

Efficient Deletes in LSM-Engines

Prof. Subhadeep Sarkar