# CS561: Data Systems Architectures

class 8

# Efficient Deletes in Log-Structured Key-Value Storage

Prof. Manos Athanassoulis

# CS561: Data Systems Architectures

class 8

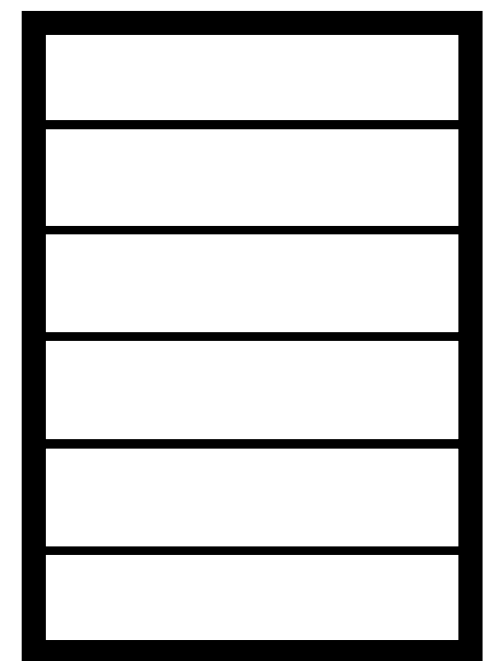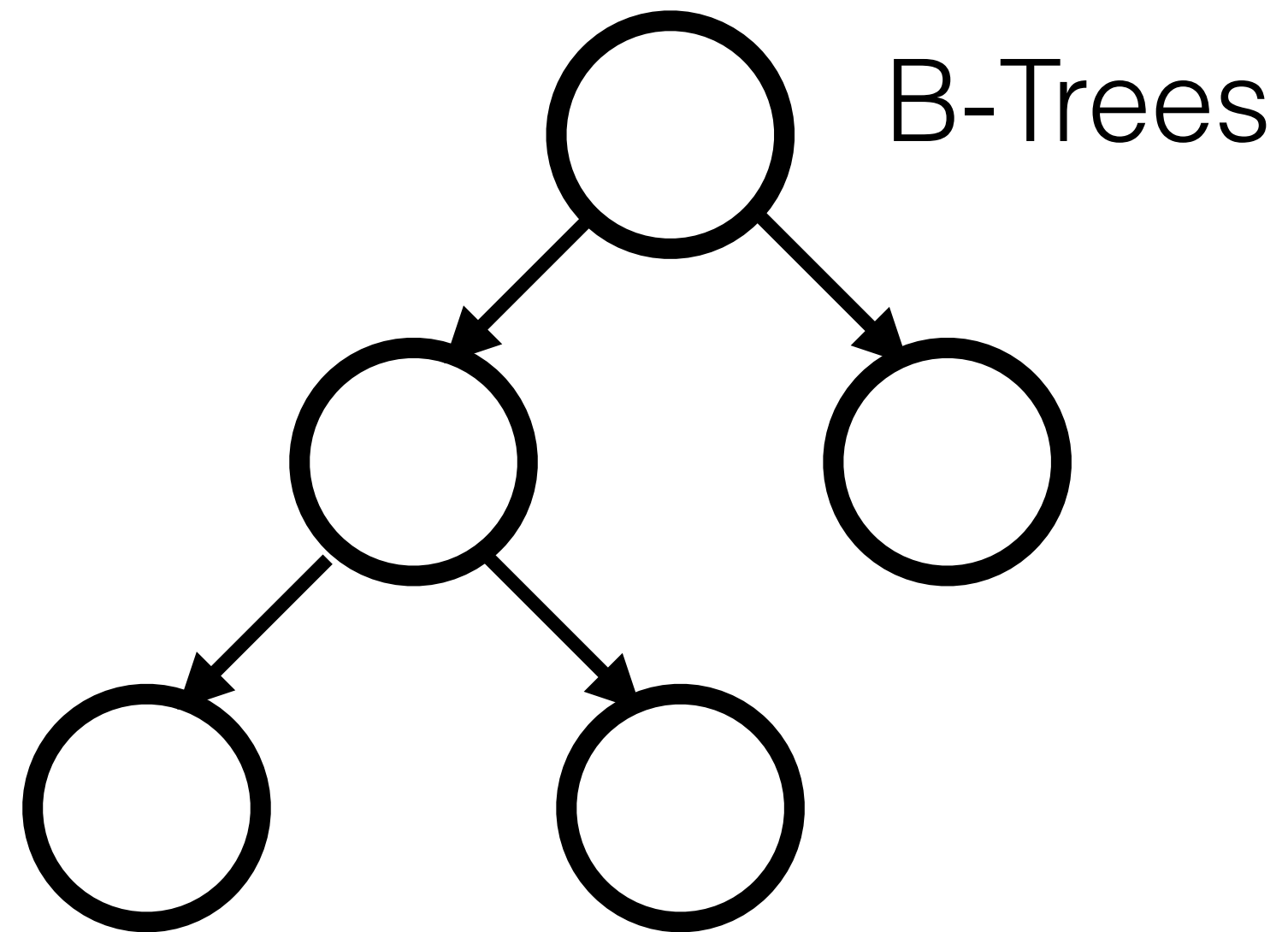# Delete: *the forgotten operator*

Prof. Manos Athanassoulis

```
<your_favorite_data_structure>::delete (key)
{
    //todo
}
```
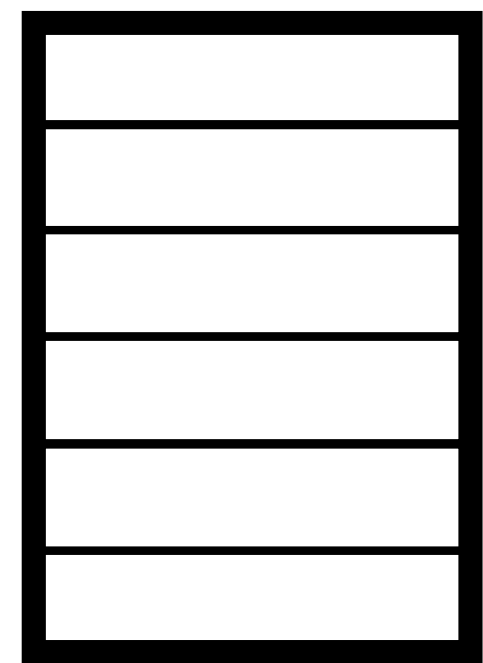
in-place                                        out-of-place

B-Trees

Heap Files
(slotted pages)

in-place

out-of-place

B-Trees

invalidate the entry

Heap Files
(slotted pages)

in-place

out-of-place

B-Trees

invalidate the entry

Heap Files
(slotted pages)

in-place                                           out-of-place

B-Trees

invalidate the entry

Heap Files
(slotted pages)
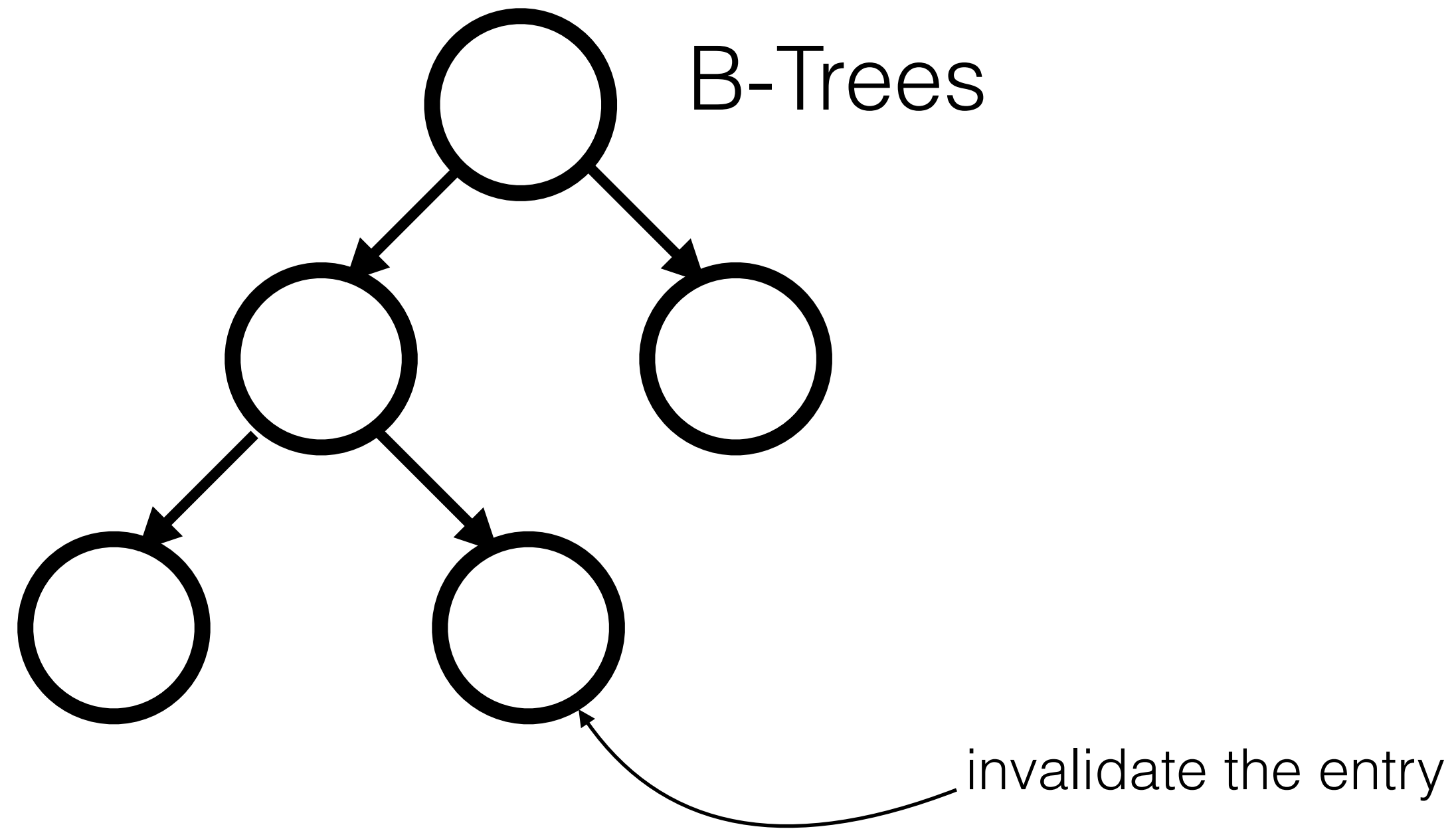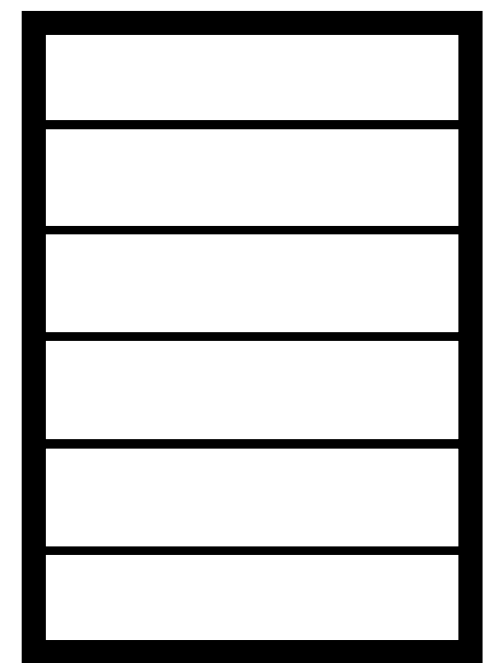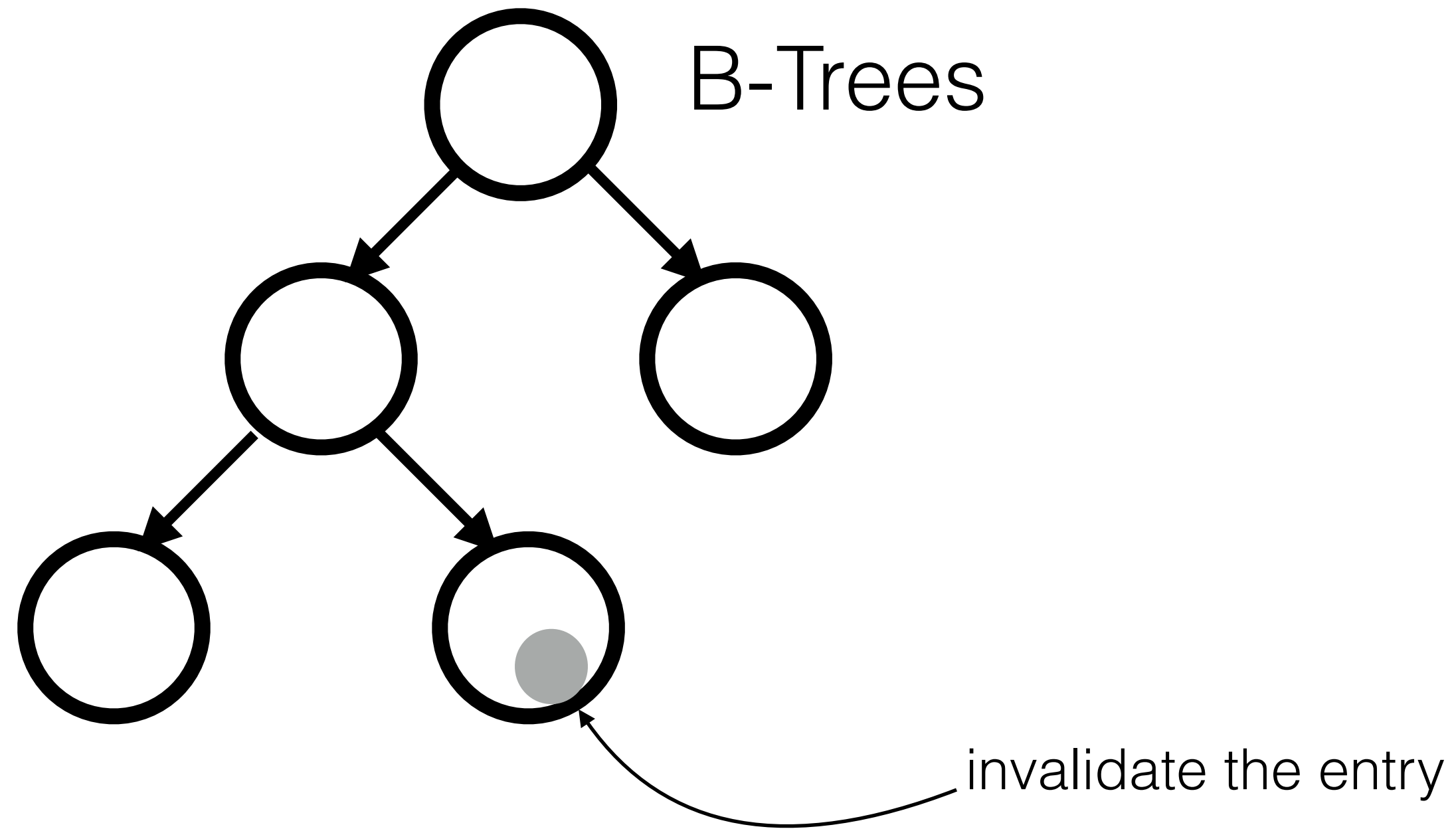
in-place

out-of-place

B-Trees

invalidate the entry
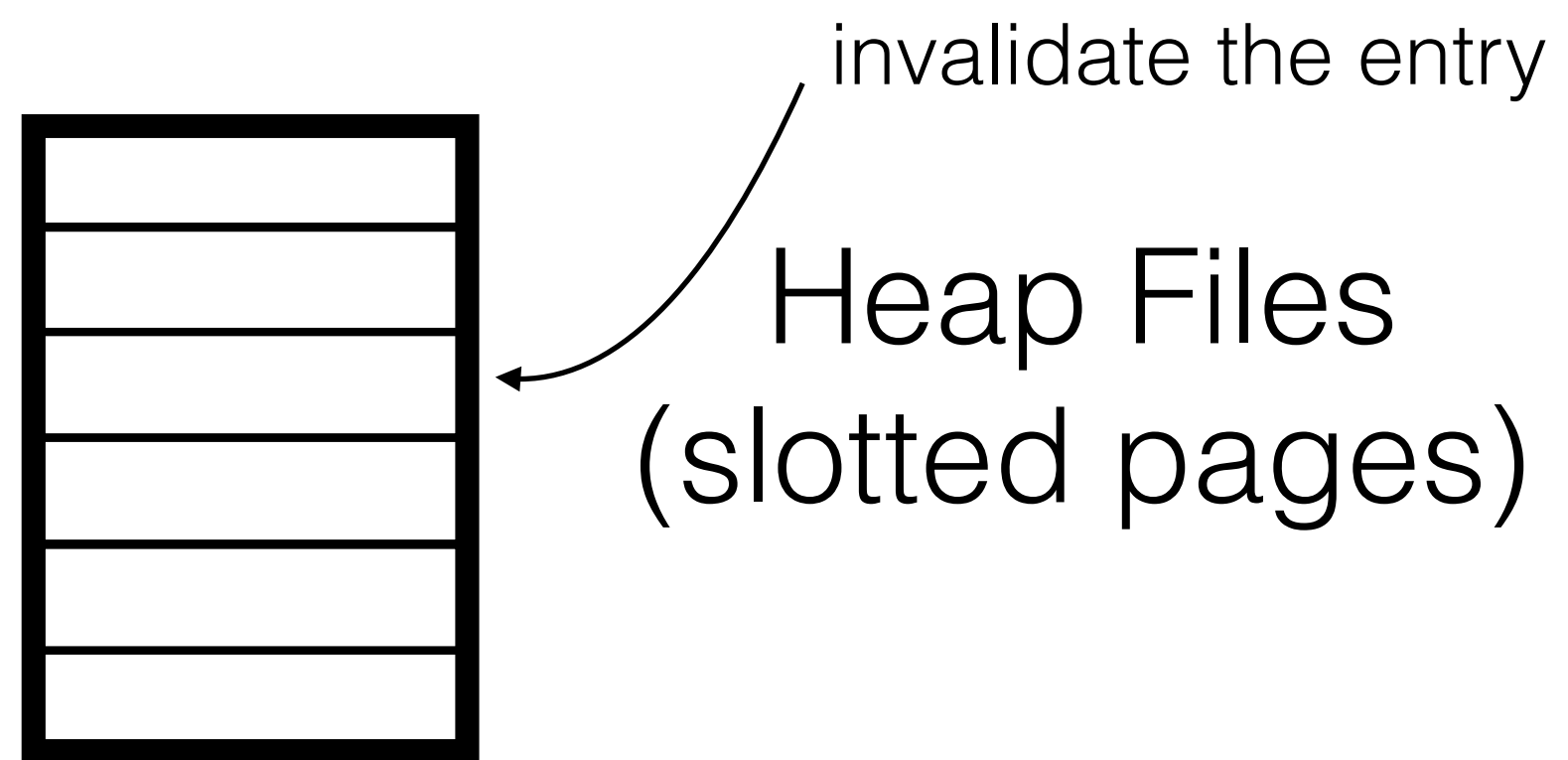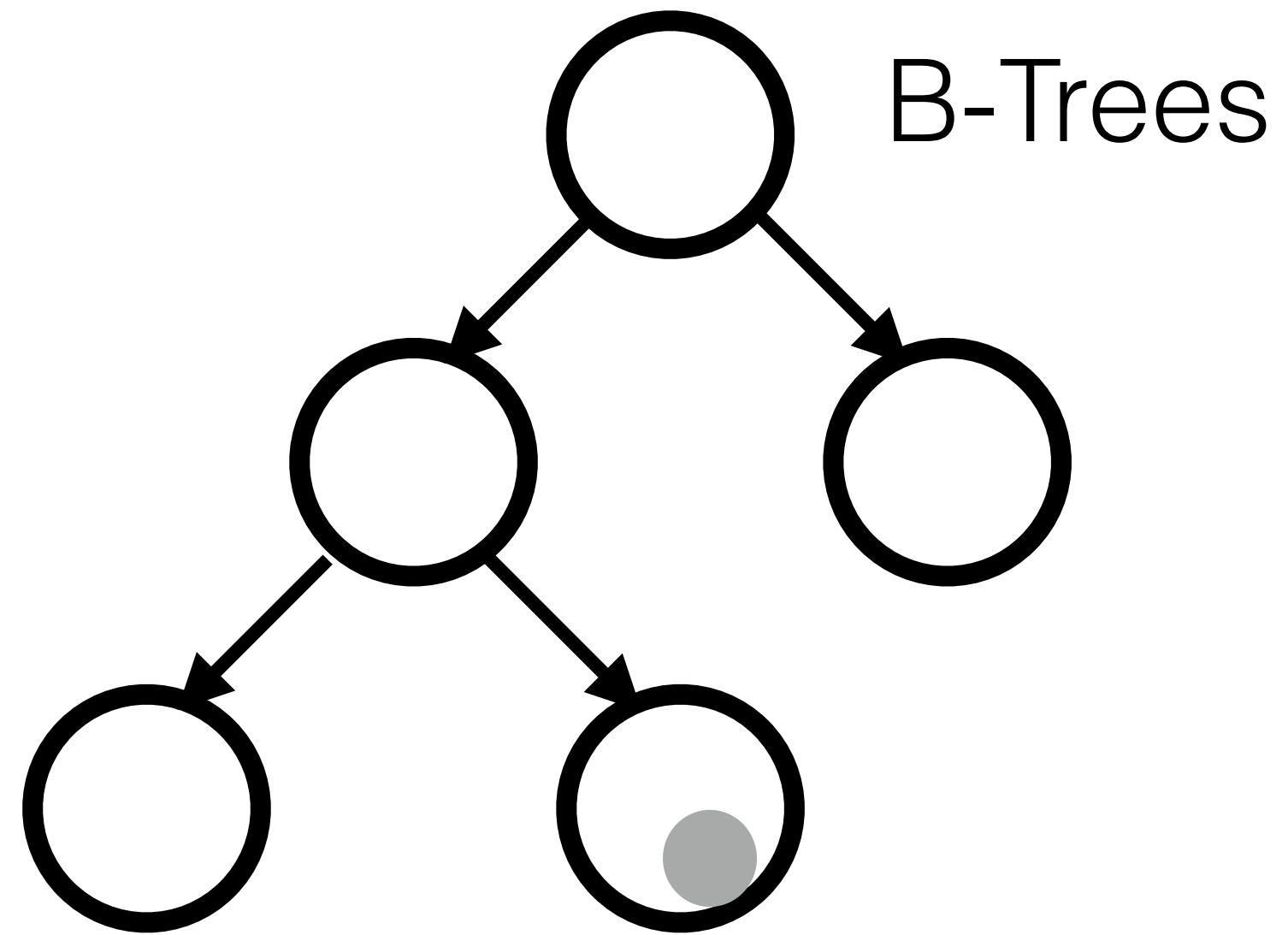
Heap Files
(slotted pages)

in-place

out-of-place

B-Trees

Heap Files
(slotted pages)

in-place

B-Trees

Heap Files
(slotted pages)

out-of-place

in-place

out-of-place

B-Trees

Heap Files
(slotted pages)

in-place

B-Trees

Heap Files
(slotted pages)

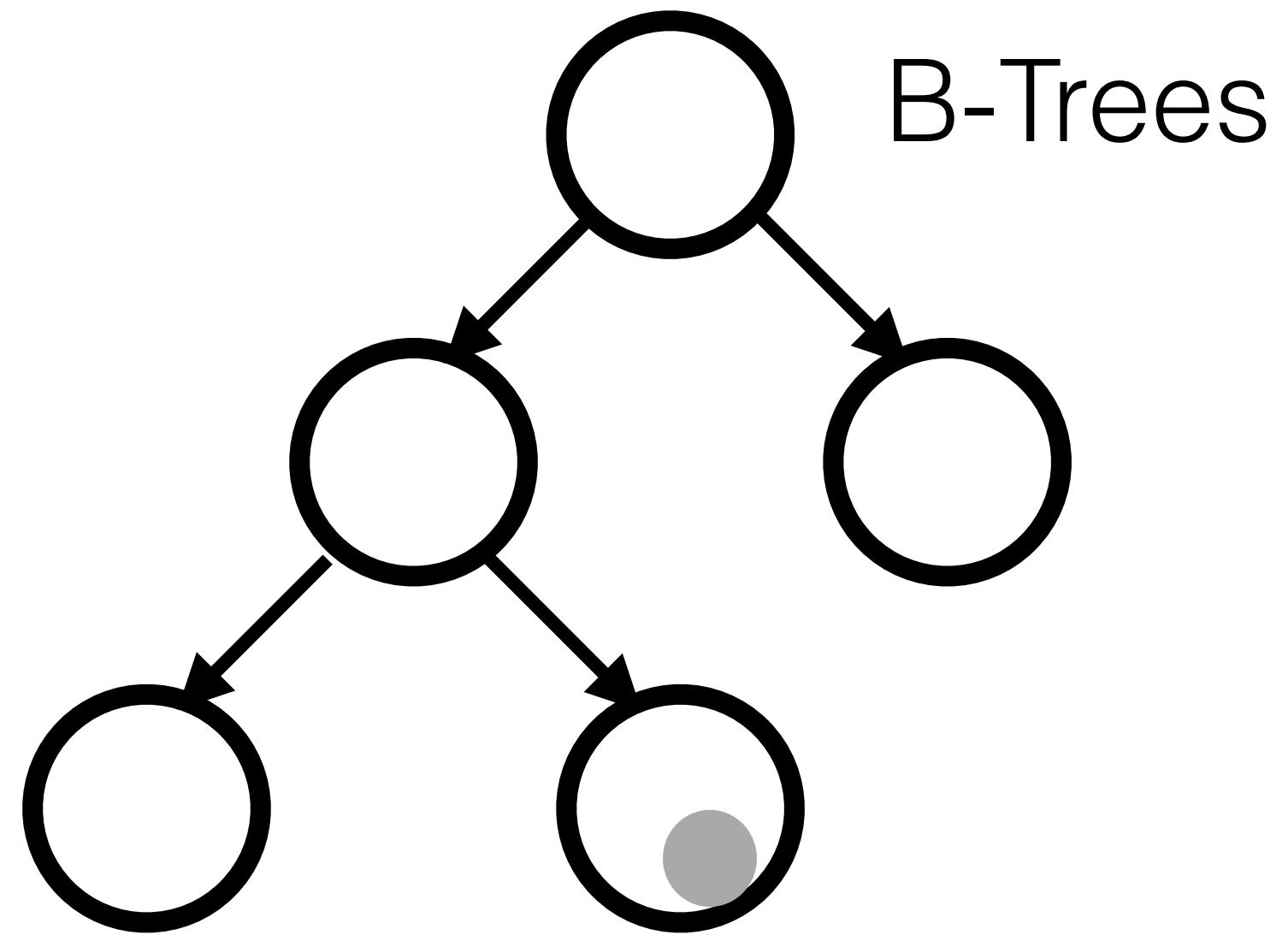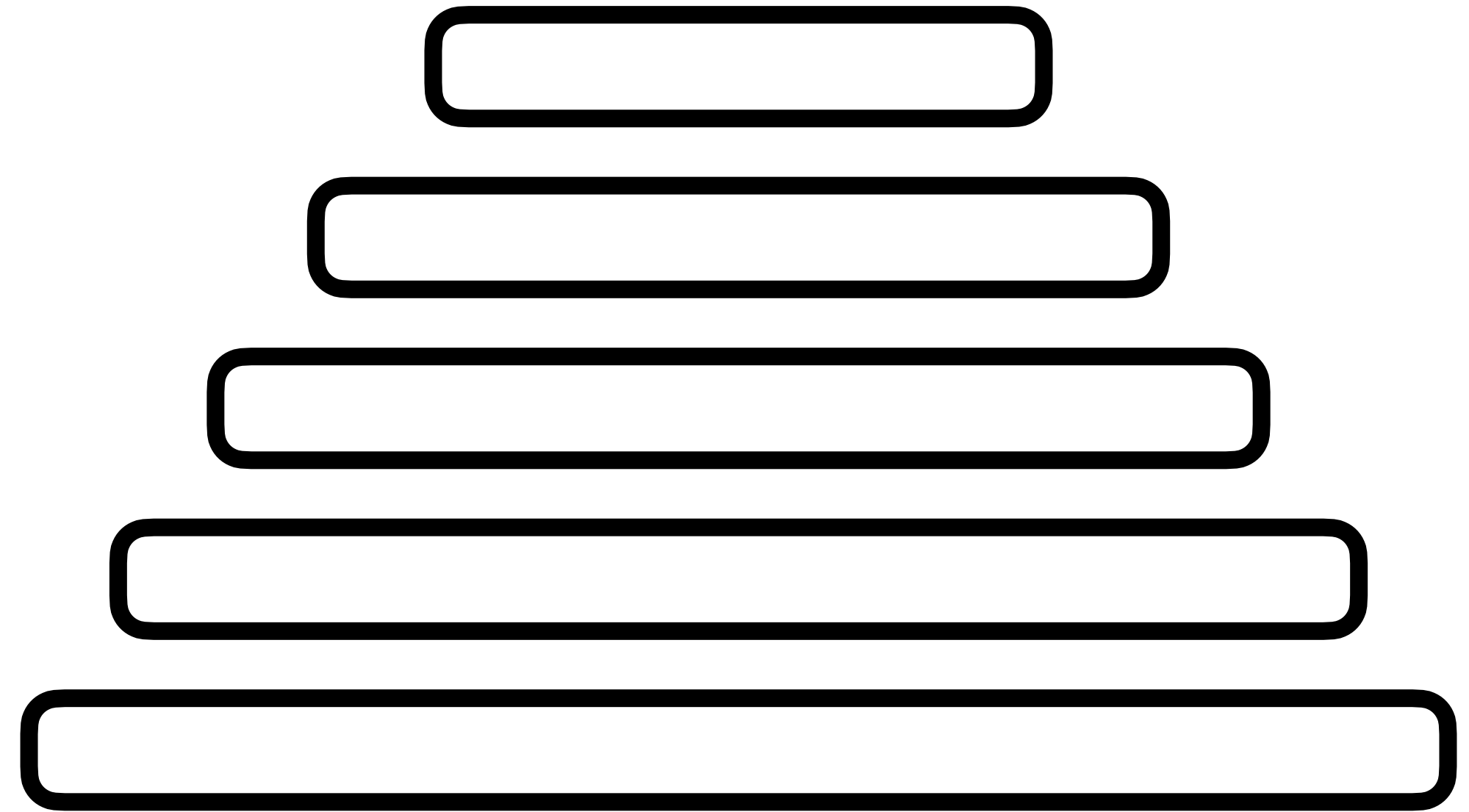out-of-place

# What is the delete tradeoff?

# What is the delete tradeoff?

read

vs.

write

?

# What is the delete tradeoff?

read

vs.

write

?

Deletes are almost **exclusively** *logical*

# What is the delete tradeoff?

read

vs. ?

write

Deletes are almost **exclusively** *logical*

b-trees, slotted pages, LSM-trees **invalidate** the entry under deletion

# delete tradeoffs

**other tradeoffs?**

# delete tradeoffs

delete latency vs. **_future_** read performance

# delete tradeoffs

delete latency vs. *future* read performance

e.g., tree re-org, page re-org, more metadata in LSM

# delete tradeoffs

delete latency vs. *future* read performance

e.g., tree re-org, page re-org, more metadata in LSM

delete latency vs. data *privacy*

# delete tradeoffs

delete latency vs. *future* read performance

e.g., tree re-org, page re-org, more metadata in LSM

delete latency vs. data *privacy*

logical deletes keep around deleted entries, what if they leak?

# delete tradeoffs

delete latency vs. *future* read performance

e.g., tree re-org, page re-org, more metadata in LSM

delete latency vs. data *privacy*

logical deletes keep around deleted entries, what if they leak?

delete latency vs. *storage amplification*

# delete tradeoffs

delete latency vs. *future* read performance

e.g., tree re-org, page re-org, more metadata in LSM

delete latency vs. data *privacy*

logical deletes keep around deleted entries, what if they leak?

delete latency vs. *storage amplification*

logical deletes keep around deleted entries and metadata!!

# delete tradeoffs

delete latency vs. **_future_** read performance

e.g., tree re-org, page re-org, more metadata in LSM

delete latency vs. data **_privacy_**

logical deletes keep around deleted entries, what if they leak?

## what if we persisted the deletes immediately?

# delete tradeoffs

delete latency vs. **persistent** delete latency

# delete tradeoffs

**logical** delete latency vs. **persistent** delete latency

**Today's talk:**

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

## Presented at SIGMOD 2020

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

# **key-value** pairs

| RID | timestamp | name | department | ... | location |
|-----|-----------|------|------------|-----|----------|

LSM-TREE

LSM-TREE

4

LSM-TREE

# LSM-TREE

4

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

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

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

2

# Large-scale production

|  | Internal db ops | | Privacy |
|---|---|---|---|
| **ZippyDB** 📌 **25.2M/day** | table drop 📌 | | CCPA 📌 (California) |
| | data migration 📌 | | GDPR 📌 (EU, UK) |
| **UP2X** 📌 **92.5M merge through deletes** | cleanup /gc 📌 | | VCPDA 📌 (Virginia) |

| Large-scale production | Internal db ops | Privacy |
|---|---|---|
| ZippyDB 📌 **25.2M/day** | table drop 📌 | CCPA 📌 (California) |
| UP2X 📌 **92.5M merge through deletes** | data migration 📌 | GDPR 📌 (EU, UK) |
| | cleanup /gc 📌 | VCPDA 📌 (Virginia) |

| Large-scale production | Internal db ops | Privacy |
|---|---|---|
| ZippyDB 📌 **25.2M/day** | table drop 📌 | 📜 GDPR 📌 (EU, UK) |
| UP2X 📌 **92.5M merge through deletes** | data migration 📌 | 📜 CCPA 📌 (California) |
| | cleanup /gc 📌 | 📜 VCPDA 📌 (Virginia) |

GDPR (EU, UK)   CCPA (California)   VCPDA (Virginia)

on-demand

*delete all data for user X within D days*

rolling

*keep deleting all data older than D days*

*A reminder on how LSM-trees work!*

# log-structured merge-tree

buffer

# log-structured merge-tree

buffer | 2 | 6 | 1 | 4 |

# log-structured merge-tree

buffer | 1 | 2 | 4 | 6 |

# 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

L2

L3

**size ratio = T**

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

**How to reduce those?**

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



**partial compaction**

buffer

L1

L2

L3

# log-structured merge-tree



**partial compaction**

buffer

L1

L2

L3

# log-structured merge-tree

**partial compaction**

buffer

L1

L2

L3

# log-structured merge-tree



**partial compaction**

buffer

L1

L2   NEW   NEW

L3

# log-structured merge-tree

**partial compaction**

buffer

L1

L2 NEW NEW

L3

# log-structured merge-tree

**partial compaction**

buffer

L1

L2

L3

# log-structured merge-tree

**partial compaction**

buffer

L1

L2

L3

# log-structured merge-tree

**partial compaction**

buffer

L1

L2 ★

L3

# log-structured merge-tree
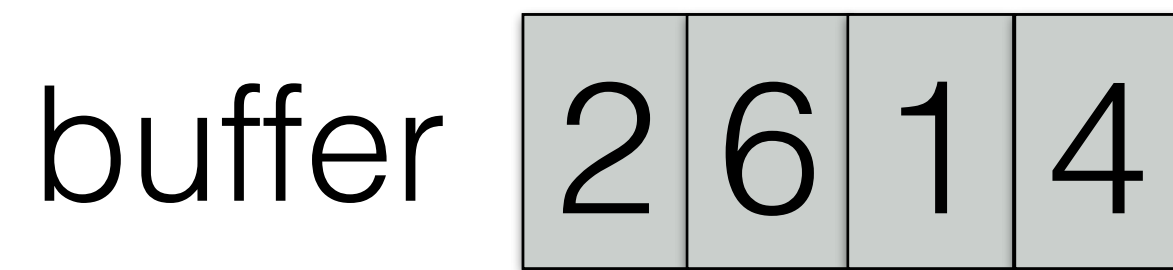
**partial compaction**

buffer

L1

L2

L3    NEW    NEW    NEW

amortized compaction cost

# log-structured merge-tree



buffer

fence
pointers

L1

L2

L3

# log-structured merge-tree

**What is the cost of a read?**

fence
pointers

buffer

L1

L2

L3

# log-structured merge-tree



get(5)

buffer

fence pointers

**one I/O per run**

L1

L2

L3

5

# log-structured merge-tree



get(5)

buffer

Bloom filters

fence pointers

L1

L2

L3 5

fewer disk I/Os

# *Now, let's talk about deletes!*

# deletes in LSM-tree

delete

# deletes in LSM-tree

delete := insert tombstone

**key** **value**

| RID | TS flag |
|-----|---------|

**key** **value**

| RID | timestamp | name | department | ⋯ | location |
|-----|-----------|------|------------|---|----------|

# deletes in LSM-tree

delete := insert tombstone

**key** | **value**

| RID | TS flag |

**key** | **value**

| RID | TS flag | timestamp | name | department | ... | location |

# deletes in LSM-tree

delete(5)

buffer

L1

L2

L3

# deletes in LSM-tree



buffer

L1   5

L2   5

L3   5

# deletes in LSM-tree

get(5)

buffer

Bloom filters

fence pointers

L1   5

L2   5

L3   5

# the problems

# the problems

# out-of-place deletes

# out-of-place deletes



L1

L2

space amplification ✗

L3

L4

# out-of-place deletes



L1

L2

space amplification ✗

L3

L4

# out-of-place deletes

L1 [ 5 ]

L2 [ ]

L3 [ ]

L4 [ 5 ]

write amplification ✗

space amplification ✗

# out-of-place deletes

Bloom
filters

L1

5

L2

5

L3

L4

5

write amplification ✖

space amplification ✖

# out-of-place deletes

Bloom
filters

L1

L2

L3

L4

poor read perf. ✗

write amplification ✗

space amplification ✗

5

5

5

# out-of-place deletes

Bloom filters

L1

L2

L3

L4

poor read perf. ✗

write amplification ✗

space amplification ✗

# the problems

poor read perf.

write amplification

space amplification

?

?

# delete persistence latency

# delete persistence latency

# delete persistence latency

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

delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

# delete persistence latency

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

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

$D_{th}$

L1

L2

L3

L4

unbounded delete persistence latency ✗

# the problems

poor read perf.

write amplification

space amplification

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

$D_{th}$

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**

| key | | value | | | | |
|-----|---|-------|---|---|---|---|
| **RID** | **TS flag** | **timestamp** | **name** | **department** | ⋯ | **location** |

**sort key**

**delete key** ?

17

# deletes on a secondary attribute

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



key | value

| RID | TS flag | timestamp | name | department | ⋯ | location |

↑ sort key

↑ delete key

L1
L2
L3
L4

# 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**



| key | value |
|-----|-------|
| **RID** | **TS flag** **timestamp** **name** **department** ⋯ **location** |

↑ **sort key**

↑ **delete key**

L1

L2

L3

L4

# deletes on a secondary attribute

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

| key | value | | | | | | |
|-----|-------|--|--|--|--|--|--|
| **RID** | **TS flag** | **timestamp** | **name** | **department** | ⋯ | **location** | |

↑ sort key

↑ delete key

L1

L2

latency spikes ✖

L3

superfluous I/Os ✖

L4

deletes on a secondary attribute

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

| key | value |

RID | TS flag | timestamp | name | department | ⋯ | location

sort key
delete key

L1

L2

latency spikes ✗

L3

superfluous I/Os ✗

L4

# the problems

poor read perf.

write amplification

space amplification

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

$D_{th}$

unbounded delete
persistence latency

latency spikes

superfluous I/Os

the solution

poor read perf.

write amplification

**FADE**

space amplification

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

$D_{th}$

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**



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

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

L1

L2

L3

L4

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



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

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

L1

L2

L3

L4

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



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

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

L1

L2

L3

L4

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



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

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

L1

L2    | 5 |

L3

L4

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



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

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

L1

L2  5

L3

L4  5

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



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

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

L1

L2

L3  5

L4  5

d$_1$

d$_2$

a$_3$

d$_3$

# FAst DElete

delete(5) within a threshold time: **D<sub>th</sub>**



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

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

L1

L2

L3

5

$d_1$

$d_2$

$d_3 = a_3$

L4

5

# FAst DElete

delete(5) within a threshold time: **D$_{th}$**



L1

L2

L3

L4

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

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

$d_1$

$d_2$

$d_3$

# FAst DElete

## breaking ties in practical workloads

L1

L2

L3

L4

# 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

1M 1KB entries, 5% deletes, 1MB buffer, T=10



timely delete persistence ✔

**within D$_{th}$**

# FAst DElete



reduced space amplification ✓

**2.1 - 9.8 x**

timely delete persistence ✓

**within D$_{th}$**

# FAst DElete



improved read performance ✔
**1.2 - 1.4 x**

reduced space amplification ✔
**2.1 - 9.8 x**

timely delete persistence ✔
**within D_th**

# FAst DElete



higher write amplification ↑
**4 - 25 %**

improved read performance ✔
**1.2 - 1.4 x**

reduced space amplification ✔
**2.1 - 9.8 x**

timely delete persistence ✔
**within $D_{th}$**

# FAst DElete

higher write amplification ⬆
**4 - 25 %**

improved read performance ✔
**1.2 - 1.4 x**

reduced space amplification ✔
**2.1 - 9.8 x**

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 - 1.4 x**

reduced space amplification ✔

**2.1 - 9.8 x**

timely delete persistence ✔

**within $D_{th}$**



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

RocksDB
FADE/25%

normalized bytes written

snapshot #

# 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**

improved read performance ✓

reduced space amplification ✓

timely delete persistence ✓

latency spikes

**KiWi**

superfluous SMOs

# Key Weaving storage layout

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

L1

L2

**scattered occurrences**

L3

L4

# Key Weaving storage layout

delete all entries with timestamp **<= 65$_D$**

# Key Weaving storage layout

## delete all entries with timestamp **<= 65$_D$**



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1: $S_{min}=1 :: S_{max}=24$ | $D_{min}=3_D :: D_{max}=80_D$

page 2: $S_{min}=29 :: S_{max}=60$ | $D_{min}=9_D :: D_{max}=90_D$

page 3: $S_{min}=61 :: S_{max}=79$ | $D_{min}=1_D :: D_{max}=89_D$

page 4: $S_{min}=80 :: S_{max}=99$ | $D_{min}=7_D :: D_{max}=85_D$

**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$_D$ |

**1 I/O**

**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$_D$ |

**1 I/O**

**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$_D$ |

**1 I/O**

**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$_D$ |

**1 I/O**

29

# Key Weaving storage layout

delete all entries with timestamp **<= 65$_D$**



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1
$S_{min}=1 :: S_{max}=24$
$D_{min}=3_D :: D_{max}=80_D$

page 2
$S_{min}=29 :: S_{max}=60$
$D_{min}=9_D :: D_{max}=90_D$

page 3
$S_{min}=61 :: S_{max}=79$
$D_{min}=1_D :: D_{max}=89_D$

page 4
$S_{min}=80 :: S_{max}=99$
$D_{min}=7_D :: D_{max}=85_D$

# Key Weaving storage layout

delete all entries with timestamp **<= 65$_D$**



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1
$S_{min}=1 :: S_{max}=24$
$D_{min}=3_D :: D_{max}=80_D$

page 2
$S_{min}=29 :: S_{max}=60$
$D_{min}=9_D :: D_{max}=90_D$

delete tile 1

page 3
$S_{min}=61 :: S_{max}=79$
$D_{min}=1_D :: D_{max}=89_D$

page 4
$S_{min}=80 :: S_{max}=99$
$D_{min}=7_D :: D_{max}=85_D$

delete tile 2

partitioned on S

# Key Weaving storage layout

delete all entries with timestamp **<= 65$_D$**



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1

$S_{min}=1 :: S_{max}=24$
$D_{min}=3_D :: D_{max}=80_D$

page 2

$S_{min}=29 :: S_{max}=60$
$D_{min}=9_D :: D_{max}=90_D$

delete tile 1

page 3

$S_{min}=61 :: S_{max}=79$
$D_{min}=1_D :: D_{max}=89_D$

page 4

$S_{min}=80 :: S_{max}=99$
$D_{min}=7_D :: D_{max}=85_D$

delete tile 2

partitioned on S

# Key Weaving storage layout

delete all entries with timestamp **<= 65$_D$**



$S_{min}$=1 :: $S_{max}$=99
$D_{min}$=1$_D$ :: $D_{max}$=90$_D$

page 1

$S_{min}$=1 :: $S_{max}$=24
$D_{min}$=3$_D$ :: $D_{max}$=80$_D$

page 2

$S_{min}$=29 :: $S_{max}$=60
$D_{min}$=9$_D$ :: $D_{max}$=90$_D$

delete tile 1

page 3

$S_{min}$=61 :: $S_{max}$=79
$D_{min}$=1$_D$ :: $D_{max}$=89$_D$

page 4

$S_{min}$=80 :: $S_{max}$=99
$D_{min}$=7$_D$ :: $D_{max}$=85$_D$

delete tile 2

$S_{min}$=1 :: $S_{max}$=60
$D_{min}$=3$_D$ :: $D_{max}$=90$_D$

delete tile 1

$S_{min}$=1 :: $S_{max}$=24
$D_{min}$=3$_D$ :: $D_{max}$=80$_D$

$S_{min}$=29 :: $S_{max}$=60
$D_{min}$=9$_D$ :: $D_{max}$=90$_D$

partitioned on S

# Key Weaving storage layout

## delete all entries with timestamp **<= 65$_D$**



$S_{min}$=1 :: $S_{max}$=99
$D_{min}$=1$_D$ :: $D_{max}$=90$_D$

**page 1**

$S_{min}$=1 :: $S_{max}$=24
$D_{min}$=3$_D$ :: $D_{max}$=80$_D$

**page 2**

$S_{min}$=29 :: $S_{max}$=60
$D_{min}$=9$_D$ :: $D_{max}$=90$_D$

**delete tile 1**

**page 3**

$S_{min}$=61 :: $S_{max}$=79
$D_{min}$=1$_D$ :: $D_{max}$=89$_D$

**page 4**

$S_{min}$=80 :: $S_{max}$=99
$D_{min}$=7$_D$ :: $D_{max}$=85$_D$

**delete tile 2**

partitioned on S

$S_{min}$=1 :: $S_{max}$=60
$D_{min}$=3$_D$ :: $D_{max}$=90$_D$

**delete tile 1**

$S_{min}$=1 :: $S_{max}$=24
$D_{min}$=3$_D$ :: $D_{max}$=80$_D$

$S_{min}$=29 :: $S_{max}$=60
$D_{min}$=9$_D$ :: $D_{max}$=90$_D$

**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$_D$ |

# Key Weaving storage layout

## delete all entries with timestamp <= 65$_D$



$S_{min}$=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$_D$

delete tile 1

page 3
$S_{min}$=61 :: $S_{max}$=79
$D_{min}$=1$_D$ :: $D_{max}$=89$_D$

page 4
$S_{min}$=80 :: $S_{max}$=99
$D_{min}$=7$_D$ :: $D_{max}$=85$_D$

delete tile 2

partitioned on S

$S_{min}$=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$_D$

partitioned on D

**page 1**

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

**page 2**

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

31

# Key Weaving storage layout

## delete all entries with timestamp <= 65$_D$

$S_{min}$=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$_D$

**delete tile 1**

**page 3**

$S_{min}$=61 :: $S_{max}$=79
$D_{min}$=1$_D$ :: $D_{max}$=89$_D$

**page 4**

$S_{min}$=80 :: $S_{max}$=99
$D_{min}$=7$_D$ :: $D_{max}$=85$_D$

**delete tile 2**

$S_{min}$=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$_D$

partitioned on D

**page 1**

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

**page 2**

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

drop page

partitioned on S

31

# Key Weaving storage layout

## delete all entries with timestamp <= 65$_D$

$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1
$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=64_D$

page 2
$S_{min}=4 :: S_{max}=56$
$D_{min}=9_D :: D_{max}=90_D$

delete tile 1

page 3
$S_{min}=61 :: S_{max}=79$
$D_{min}=1_D :: D_{max}=89_D$

page 4
$S_{min}=80 :: S_{max}=99$
$D_{min}=7_D :: D_{max}=85_D$

delete tile 2

partitioned on S

$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=90_D$

delete tile 1
$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=64_D$

$S_{min}=4 :: S_{max}=56$
$D_{min}=69_D :: D_{max}=90_D$

partitioned on 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$_D$ |

drop page

sorted on S

31

# Key Weaving storage layout

## delete all entries with timestamp <= 65$_D$



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1
$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=64_D$

page 2
$S_{min}=4 :: S_{max}=56$
$D_{min}=9_D :: D_{max}=90_D$

delete tile 1

page 3
$S_{min}=61 :: S_{max}=79$
$D_{min}=1_D :: D_{max}=89_D$

page 4
$S_{min}=80 :: S_{max}=99$
$D_{min}=7_D :: D_{max}=85_D$

delete tile 2

partitioned on S

$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=90_D$

delete tile 1
$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=64_D$

$S_{min}=4 :: S_{max}=56$
$D_{min}=69_D :: D_{max}=90_D$

partitioned on 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**

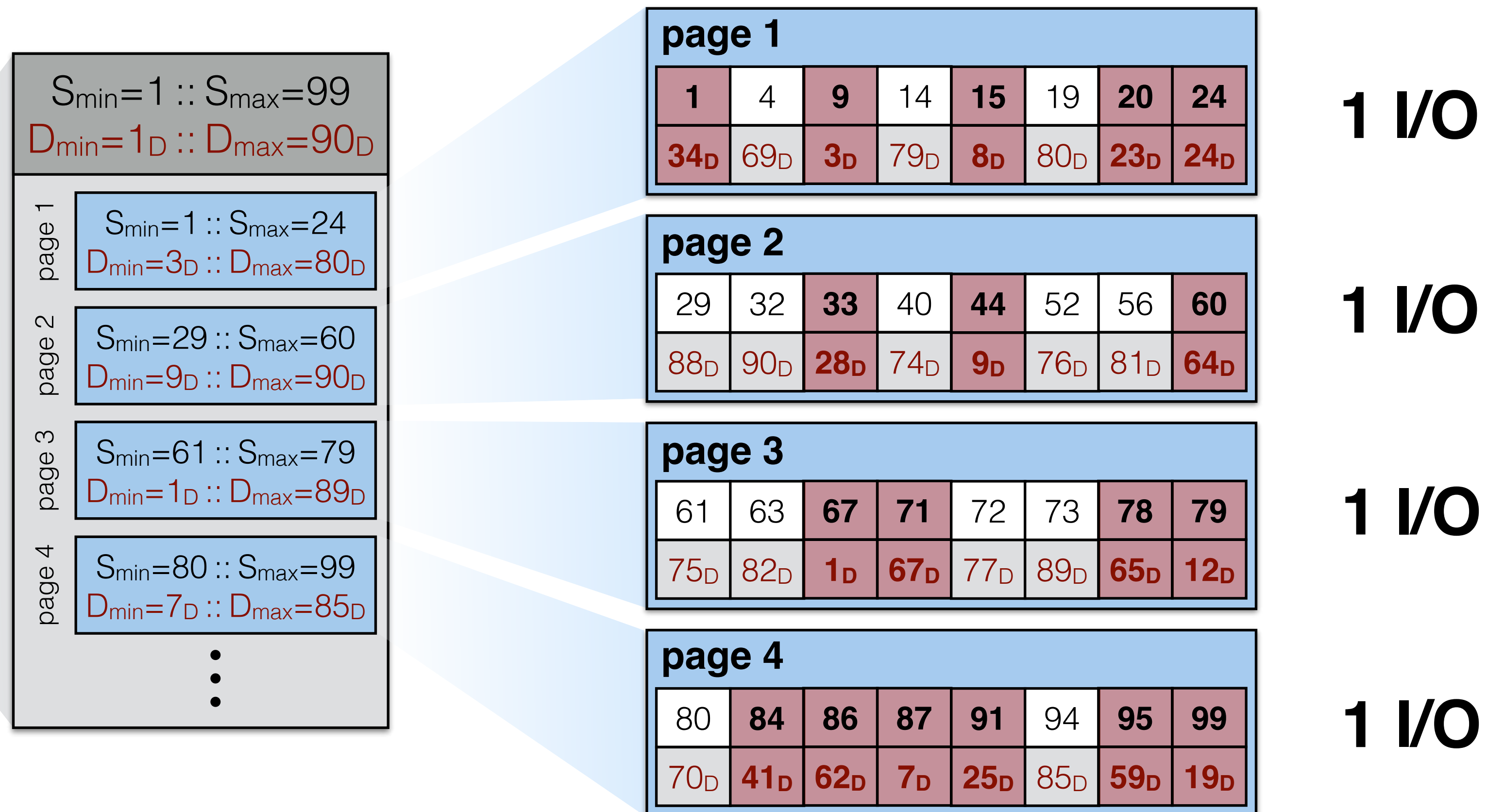| 4 | 14 | 19 | 29 | 32 | 40 | 52 | 56 |
|---|----|----|----|----|----|----|----|
| 69$_D$ | 79$_D$ | 80$_D$ | 88$_D$ | 90$_D$ | 74$_D$ | 76$_D$ | 81$_D$ |

sorted on S

drop page

31

# Key Weaving storage layout
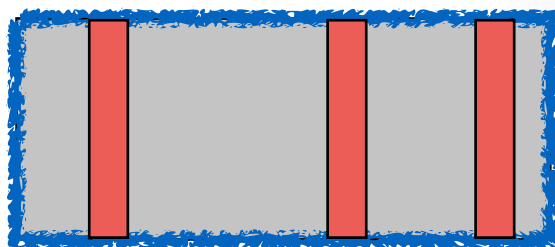
## delete all entries with timestamp <= 65$_D$



$S_{min}=1 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=90_D$

page 1
$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=64_D$

page 2
$S_{min}=4 :: S_{max}=56$
$D_{min}=9_D :: D_{max}=90_D$

delete tile 1

page 3
$S_{min}=67 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=62_D$

page 4
$S_{min}=61 :: S_{max}=94$
$D_{min}=65_D :: D_{max}=89_D$

delete tile 2

partitioned on S

$S_{min}=1 :: S_{max}=60$
$D_{min}=3_D :: D_{max}=90_D$

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
$S_{min}=67 :: S_{max}=99$
$D_{min}=1_D :: D_{max}=62_D$

$S_{min}=61 :: S_{max}=94$
$D_{min}=65_D :: D_{max}=89_D$

partitioned on 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$_D$ |

drop page

### 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$_D$ |

drop page

### 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$_D$ |

**1 I/O**

sorted on S

# Key Weaving storage layout

**1** page/delete tile

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

% full page drops

100

80

60

40

20

0

1%   2%   3%   4%   5%

fraction of deleted entries (%)

# Key Weaving storage layout

**4** pages/delete tile

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



h=1  h=4

% full page drops

fraction of deleted entries (%)

# Key Weaving storage layout



**8** pages/delete tile

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

% full page drops

100

80

60

40

20

0

1%   2%   3%   4%   5%

fraction of deleted entries (%)

h=8
h=1   h=4

# Key Weaving storage layout

**16** pages/delete tile

h=16

h=4

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



% full page drops

h=16

h=8

h=1    h=4

fraction of deleted entries (%)

# Key Weaving storage layout



**32** pages/delete tile

h=32
h=16
h=8
h=4
h=1

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

% full page drops

fraction of deleted entries (%)

# Key Weaving storage layout

reduced latency spikes ✓

full page drops reduces superfluous I/Os ✓

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

% full page drops

h=256
h=128
h=64
h=32
h=16
h=8
h=4
h=1

fraction of deleted entries (%)

# Key Weaving storage layout

higher lookup cost ↑

reduced latency spikes ✔

full page drops reduces superfluous I/Os ✔

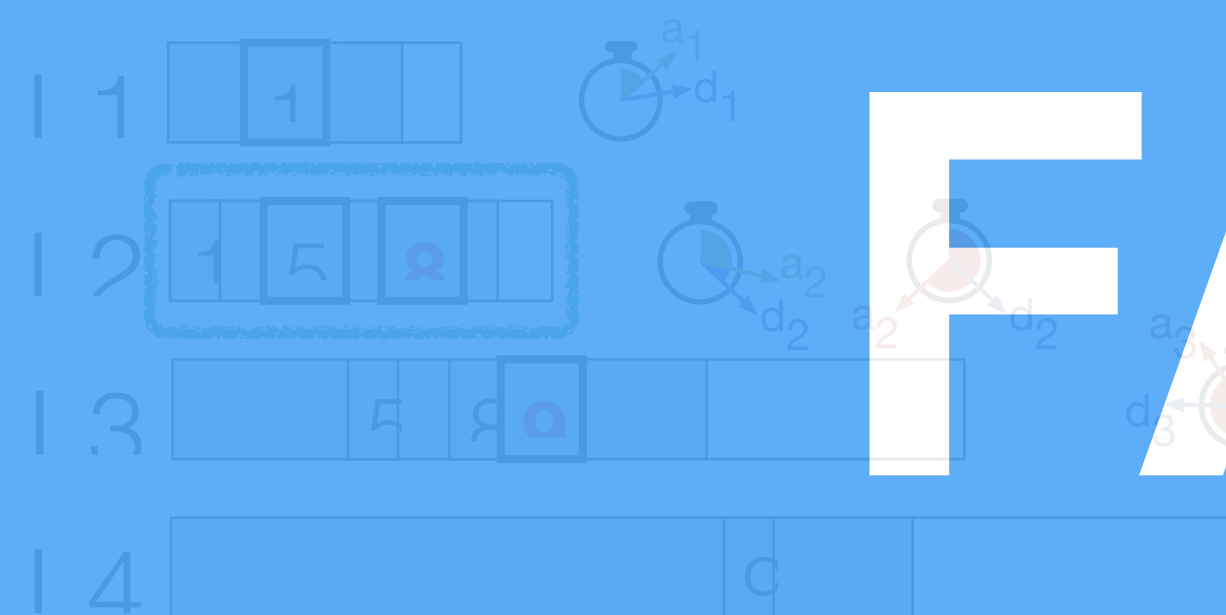1M point lookups, buffer = file = 256 pages, T=10



avg lookup cost (I/Os)

- ✚ Non-zero result lookup
- ✖ Zero result lookup

RocksDB

RocksDB

delete−tile granularity (log scale)

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

# CS561: Data Systems Architectures

class 8

# Efficient Deletes in Log-Structured Key-Value Storage

Prof. Manos Athanassoulis