# Toward Privacy-Aware Data Systems

## navigating the **privacy-performance tradeoff**

## Subhadeep Sarkar

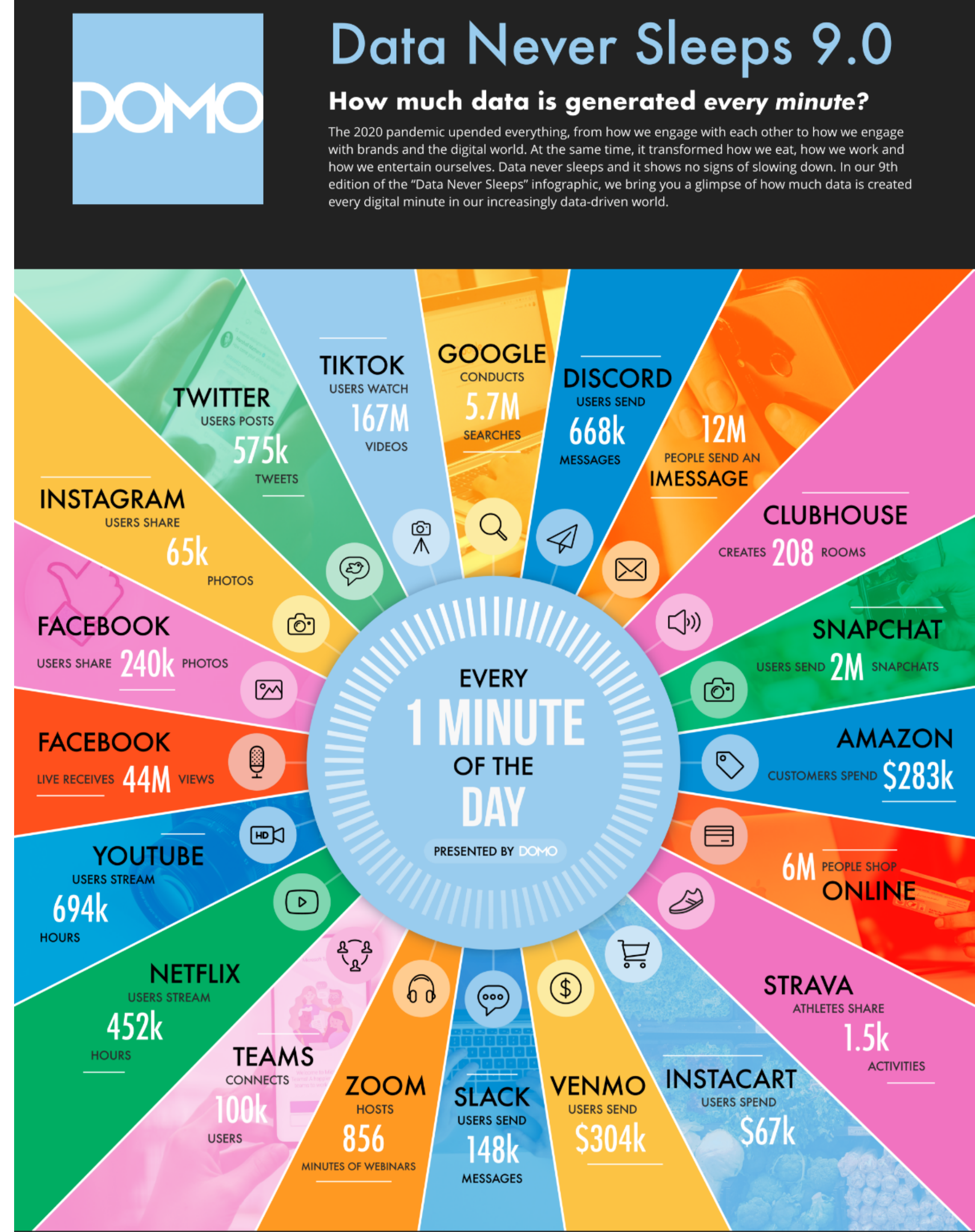https://subhadeep.net/

Brandeis
UNIVERSITY

> "*Every two days we generate as much data as we did since the dawn of humanity until 2003.*"
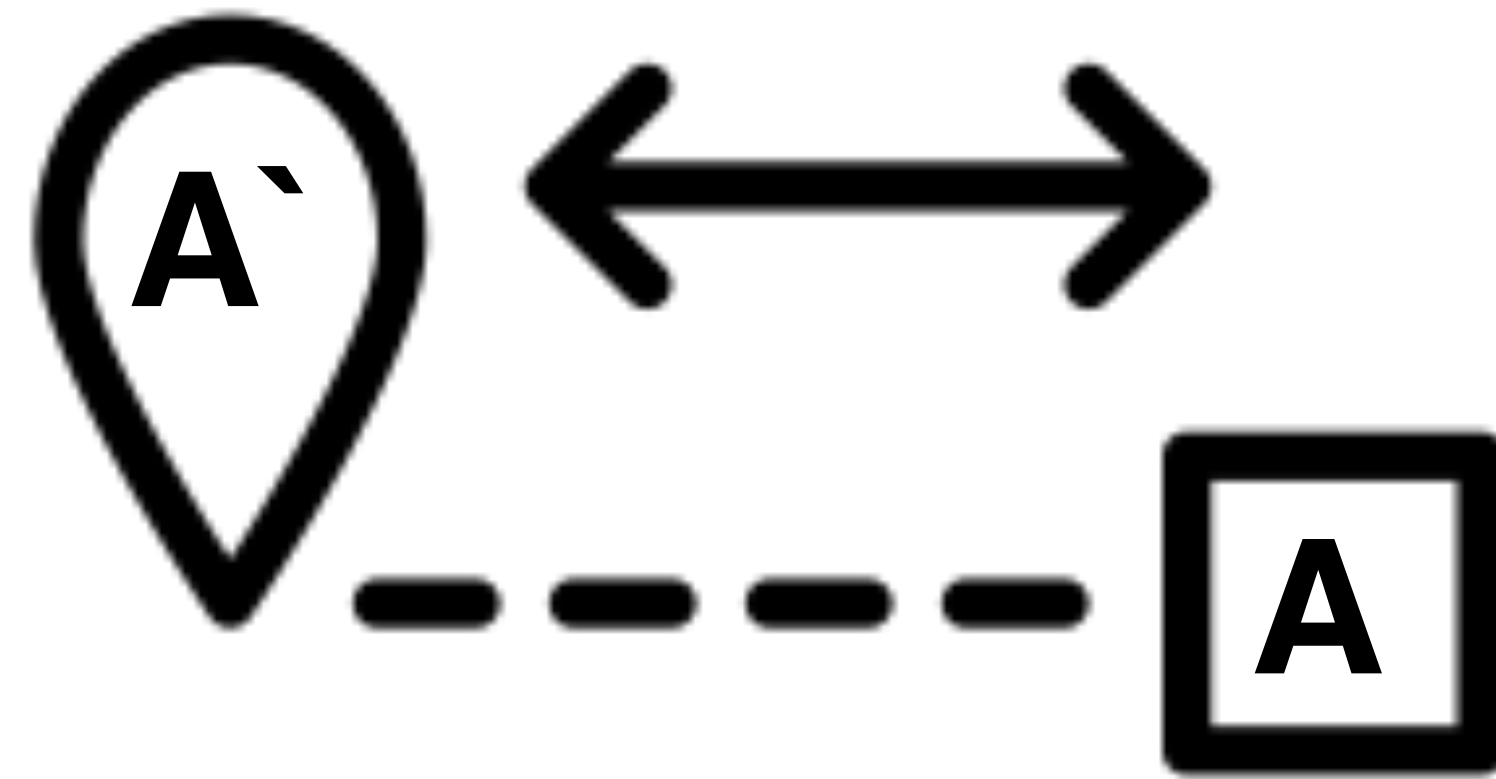>
> — **Eric Schmidt** (CEO, Google), 2010
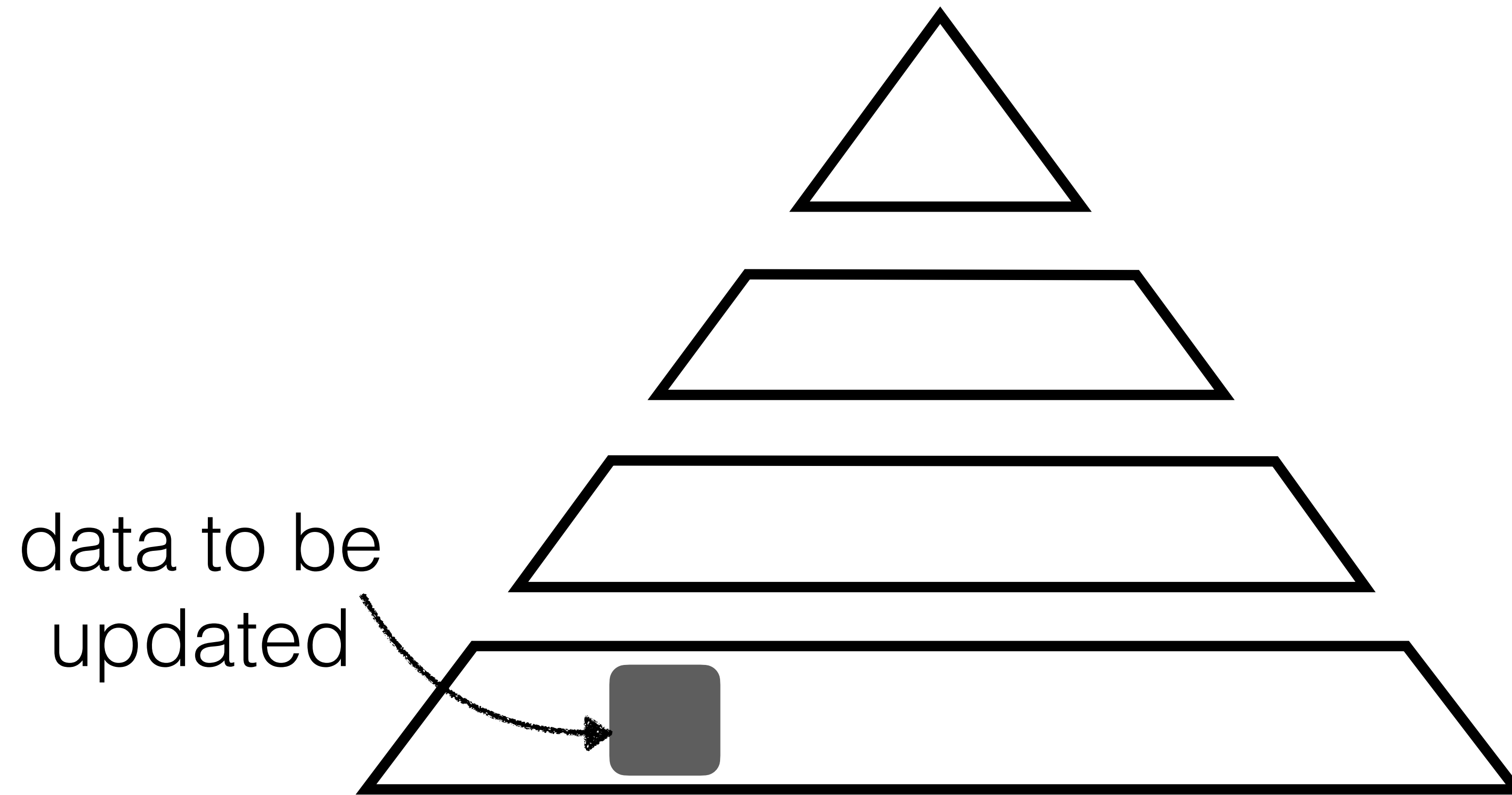
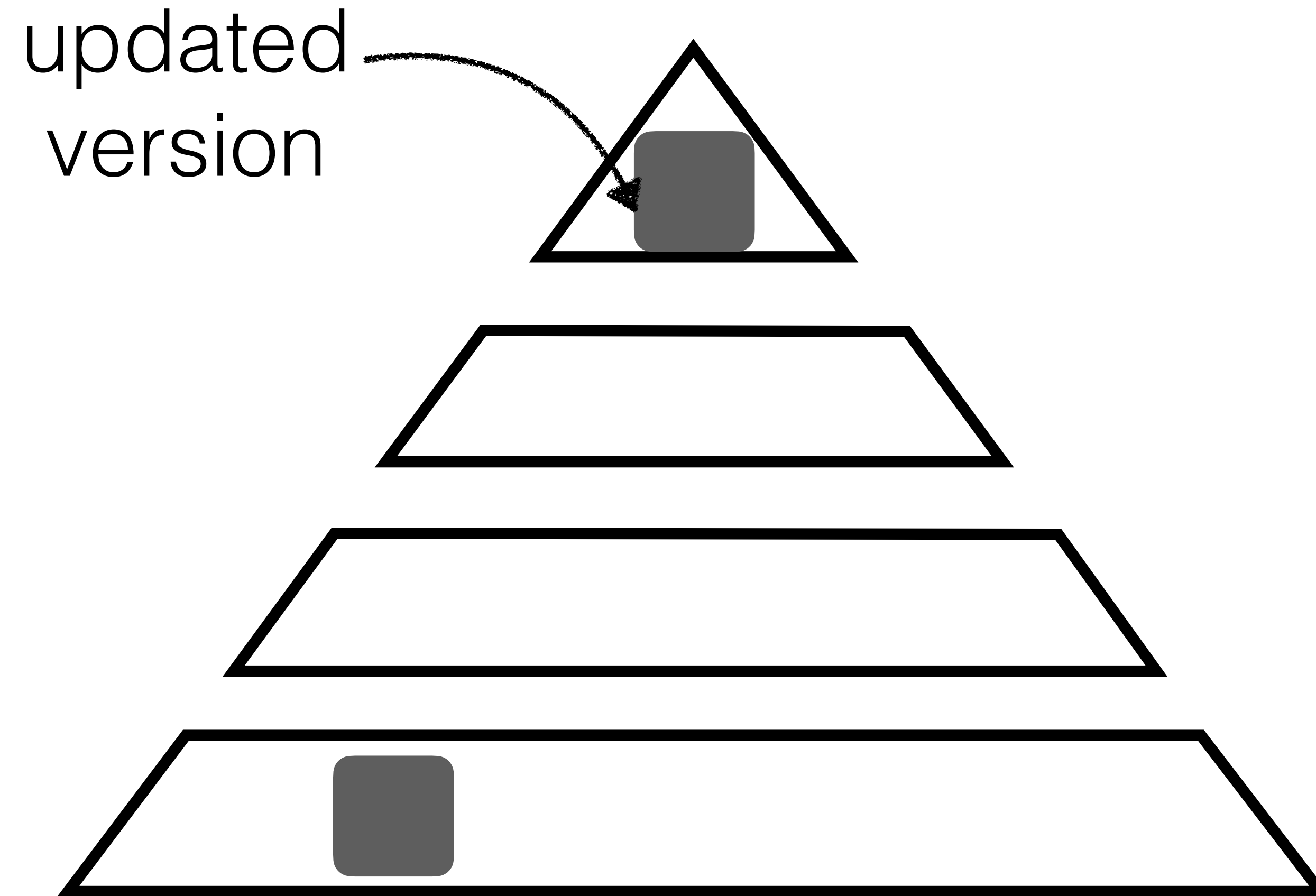# Ingestion-Optimized Systems



**batched** inserts

**out-of-place**
deletes/updates

# **Out-of-place** Deletes/Updates
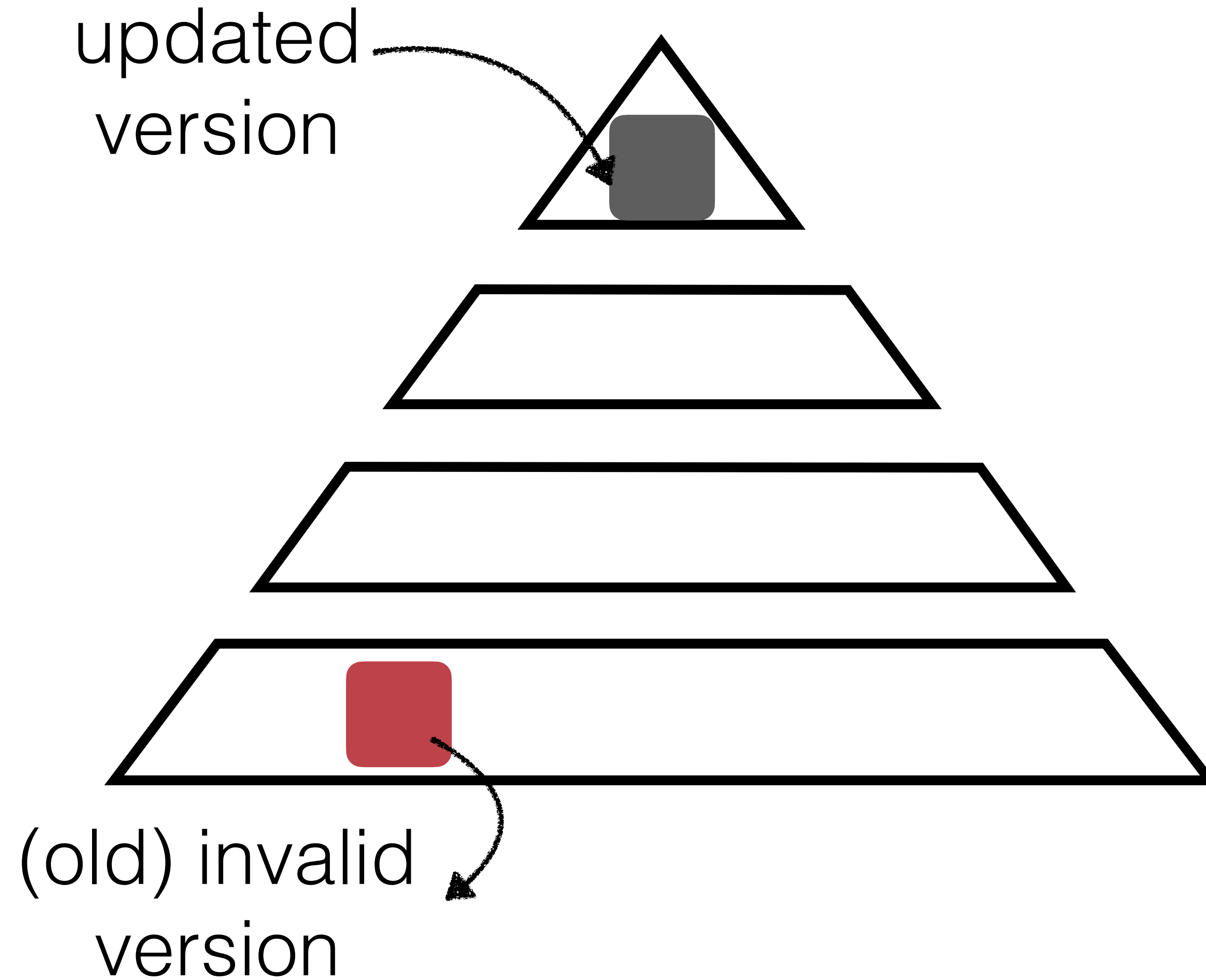


data to be updated

# **Out-of-place** Deletes/Updates

# **Out-of-place** Deletes/Updates



updated version

(old) invalid version

# **Out-of-place** Deletes/Updates



updated version

(old) invalid version

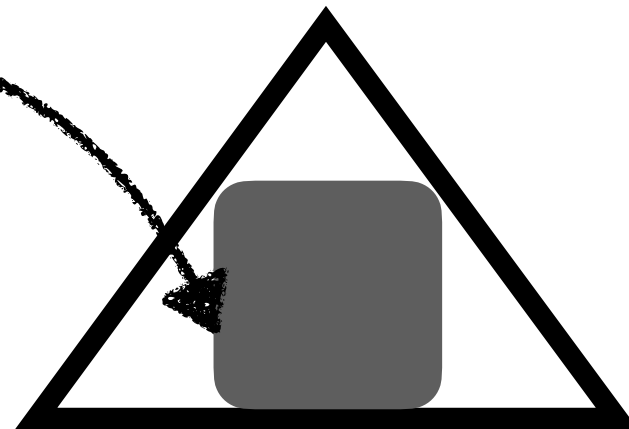# **Out-of-place** Data Systems

updated
version

Google ∞ Meta amazon

THE APACHE®
SOFTWARE FOUNDATION

SAP®

mongoDB

Hidden Cost: does not scale with deletes!

# Logical Deletes & **Data Privacy**



CCPA

VCDPA

GDPR

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

REGULATION & ENFORCEMENT

HEAVY

ROBUST

MODERATE

LIMITED

GDPR (EU, UK) — *Right to be forgotten*

CCPA (California) — *Right to delete*

VCDPA (Virginia) — *Deletion right*

timely deletes + persistent deletes

Right to be forgotten

Right to delete

Deletion right

timely deletes + persistent deletes

**Even years later, Twitter doesn't delete your direct messages**
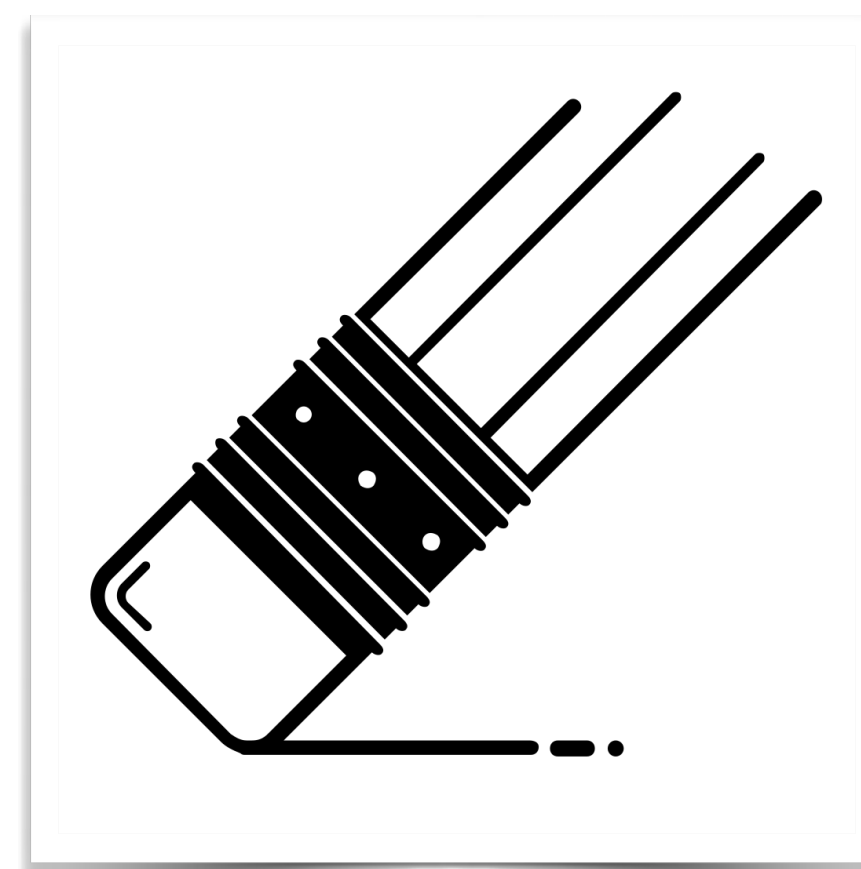
Zack Whittaker, Natasha Lomas / 1:57 PM EST • February 15, 2019
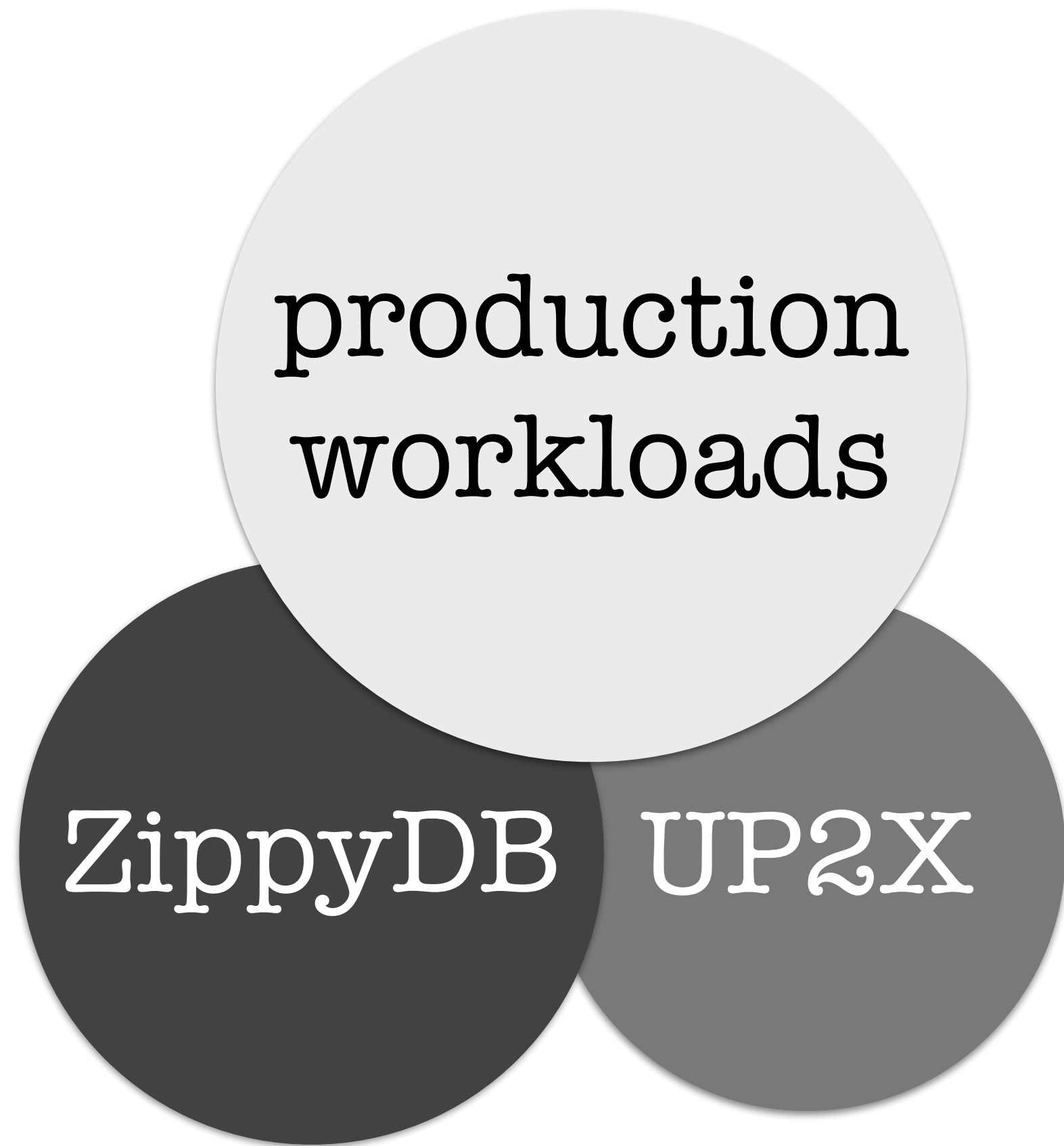
My Photo Stream

Some items will also be deleted from 11 albums.

Delete 6,447 Items

Cancel

**100M+** deletes/day

deletes in
**batches**

**user-generated**
deletes

production
workloads

ZippyDB UP2X

DB
internal
operations

table
drop

data
migration

privacy
regulations

GDPR CCPA

VCDPA

# Goal: Enabling **Privacy through Deletion**



dissecting privacy policies and regulations

fundamental design changes in data layouts and access methods

navigating the privacy-performance tradeoff

IEEE DEBull 2022

EDBT 2022

IEEE iThings 2018

ICDE 2023-a

TPCTC 2022

SIGMOD 2022-b

VLDB 2021

ICDE 2023-c

ACM TODS 2023-b

SIGMOD 2022-a

SIGMOD 2020

# Goal: Enabling **Privacy through Deletion**



<span style="color:#a52a2a">dissecting</span> privacy policies and regulations

<span style="color:#a52a2a">fundamental design changes</span> in data layouts and access methods

<span style="color:#a52a2a">navigating</span> the privacy-performance <span style="color:#a52a2a">tradeoff</span>

IEEE DEBull 2022

**EDBT** 2022

IEEE iThings 2018

ICDE 2023-a

TPCTC 2022

SIGMOD 2022-b

**VLDB** 2021

ICDE 2023-b

SIGMOD 2022-a

**SIGMOD** 2020

ACM TODS (under review)

# **L**og-**S**tructured **M**erge-tree

# LSM-tree

# LSM-tree

The Log-Structured Merge-Tree (LSM-Tree)

1996

Patrick O'Neil[1], Edward Cheng[2]
Dieter Gawlick[3], Elizabeth O'Neil[1]
To be published: Acta Informatica

# LSM-tree

NoSQL

RocksDB · WT · levelDB · SCYLLA · DynamoDB · cassandra · tarantool · Bigtable · APACHE HBASE · accumulo™ · riak

SQLite
relational

influxdb · QuasarDB
time-series

2023

# LSM-tree

NoSQL



relational

time-series

2023

# LSM **Basics**

key-value pairs

# LSM **Basics**

key-value pairs

| key | value |
|---|---|

buffer

put(6)

put(2)

buffer

put(1)

put(6)

buffer 2

put(4)

put(1)

buffer

| 2 | 6 |  |

put(4)

buffer | 2 | 6 | 1 |  |

buffer 2 6 1 4

buffer 1 2 4 6

buffer

buffer

level 1

immutable sorted run

buffer

level 1

buffer

level 1

buffer

level 1

buffer

level 1

buffer

level 1

buffer

level 1

buffer

level 1

level 2

level 3

level 4

sorted components

buffer

level 1

level 2

level 3

level 4

size ratio: T

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

compaction

level 2

NEW NEW NEW

level 3

level 4

buffer

level 1

compaction

level 2

level 3

level 4

buffer

level 1

compaction

level 2

level 3

So, what about deletes?

# **Deletes** in LSMs

delete := insert tombstone

# **Deletes** in LSMs

delete := insert tombstone

key | value

| RID | TS flag |

# **Deletes** in LSMs

## delete(5)

# **Deletes** in LSMs

delete(5)

# **Deletes** in LSMs

delete(5)

no support for timely delete persistence

great for inserts

L1

L2

L3

L4

5 X

5 val

# **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: $D_{th}$

L1

L2

L3

L4

$t_1$

$D_{th}$

5 X

5 val

# **Problem:** Persisting Deletes Timely

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

# **Problem:** Persisting Deletes Timely

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



L1

L2

L3

L4

5 X

$t_1+t_2+t_3$

$D_{th}$

# **Problem:** Persisting Deletes Timely

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



L1

L2

L3

L4

$t_1+t_2+t_3$

$D_{th}$

workload → **Compaction** → performance

What are the **design choices**?

How does a choice **affect performance**?



**Compaction**

**1**    **How** to organize the data on device?

**2**    **How much** data to move at-a-time?

**3**    **Which** block of data to be moved?

**4**    **When** to re-organize the data layout?

**Data Layout**

**1** **How** to organize the data on device?

**Compaction granularity**

**2** **How much** data to move at-a-time?

**Data movement policy**

**3** **Which** block of data to be moved?

**Trigger**

**4** **When** to re-organize the data layout?

**any existing + completely new**
compaction strategies

| Database | Data layout | Compaction Trigger | | | | | Compaction Granularity | | | | Data Movement Policy | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Level saturation | #Sorted runs | File staleness | Space amp. | Tombstone-TTL | Level | Sorted run | File (single) | File (multiple) | Round-robin | Least overlap (+1) | Least overlap (+2) | Coldest file | Oldest file | Tombstone density | Expired TS-TTL | N/A (entire level) |
| RocksDB [30], Monkey [22] | Leveling / 1-Leveling | ✓ | | ✓ | | | | | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | | |
| | Tiering | | ✓ | | ✓ | ✓ | | ✓ | | | | | | | | | | ✓ |
| LevelDB [32], Monkey (J.) [21] | Leveling | ✓ | | | | | | ✓ | | | ✓ | ✓ | ✓ | | | | | |
| SlimDB [47] | Tiering | ✓ | | | | | | | ✓ | ✓ | | | | | | | | ✓ |
| Dostoevsky [23] | $L$-leveling | ✓$^L$ | ✓$^T$ | | | | ✓$^L$ | ✓$^T$ | | | | ✓$^L$ | | | | | | ✓$^T$ |
| LSM-Bush [24] | Hybrid leveling | ✓$^L$ | ✓$^T$ | | | | ✓$^L$ | ✓$^T$ | | | | ✓$^L$ | | | | | | ✓$^T$ |
| Silk [11], Silk+ [12] | Leveling | ✓ | | | | | | | ✓ | ✓ | ✓ | | | | | | | |
| HyperLevelDB [35] | Leveling | ✓ | | | | | | | | ✓ | ✓ | ✓ | ✓ | | | | | |
| PebblesDB [46] | Hybrid leveling | ✓ | | | | | | | ✓ | ✓ | | | | | | | | ✓ |
| Cassandra [8] | Tiering | | ✓ | ✓ | | ✓ | | ✓ | | | | | | | | | | ✓ |
| | Leveling | ✓ | | | | ✓ | | | ✓ | ✓ | | | | ✓ | | ✓ | ✓ | |
| WiredTiger [62] | Leveling | ✓ | | | | | ✓ | | | | | | | | | | | ✓ |
| X-Engine [34], Leaper [63] | Hybrid leveling | ✓ | | | | | | | ✓ | ✓ | | | | ✓ | | ✓ | | |
| HBase [7] | Tiering | | ✓ | | | | | ✓ | | | | | | | | | | ✓ |
| AsterixDB [3] | Leveling | ✓ | | | | | ✓ | | | | | | | | | | | ✓ |
| | Tiering | | ✓ | | | | | ✓ | | | | | | | | | | ✓ |
| Tarantool [57] | $L$-leveling | ✓$^L$ | ✓$^T$ | | | | ✓$^L$ | ✓$^T$ | | | | | | | | | | ✓ |

# FADE

# **FA**st **DE**lete

family of
**compaction
strategies**

# **FA**st **DE**lete



compaction
trigger

compaction file
picking policy

💡 decompose tasks       💡 piggyback

# **FA**st **DE**lete

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



L1

L2

L3

L4

5 val

$D_{th}$

# **FA**st **DE**lete

delete(5) within threshold: $\mathrm{D_{th}}$

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

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

L1

L2

L3

L4

5 val

$d_1$

$d_2$

$d_3$

# **FA**st **DE**lete



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

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

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

L1

5 X

L2

L3

L4

5 val

$a_1$  $d_1$

$d_2$

$d_3$

# **FA**st **DE**lete

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

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

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

# **FA**st **DE**lete

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

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

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



L1

L2

5 X

L3

L4

5 val

$a_1$ $d_1$

$d_2$

$d_3$

# **FA**st **DE**lete

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

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

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

# **FA**st **DE**lete



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

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

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

# **FA**st **DE**lete

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

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

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

# **FA**st **DE**lete

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

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

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



L1

L2

L3

5 X

L4

5 val

$d_1$

$d_2$

$d_3 = a_3$

# **FA**st **DE**lete

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

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

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

L1

L2

L3

L4

$d_1$

$d_2$

$a_3$

$d_3$

# **FA**st **DE**lete

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

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

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

persist all deletes within:

**300s** **150s** **600s**

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

$x10^3$

RocksDB

persists deletes timely

within threshold

#deletes

50
40
30
20
10
0

90    450    1800

age of data (s)

old ← → new

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

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

improved read performance

1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

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

read throughput (ops/s)

% deletes in workload

RocksDB
FADE/16%
FADE/25%
FADE/50%

higher write amplification

4% - 25%

improved read performance

1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

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

total data written (GB)

% deletes in workload

RocksDB
FADE/16%
FADE/25%
FADE/50%

higher write amplification

4% - 25%

improved read performance

1.2x - 1.4x

reduced space amplification

2.1x - 9.8x

persists deletes timely

within threshold

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

normalized bytes written

snapshot #

RocksDB
FADE/25%

higher write amplification

0.7%

improved read performance

1.2x - 1.4x

reduced space  amplification

2.1x - 9.8x

persists deletes timely

within threshold



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

FADE

RocksDB

on-demand

retention-based

persist all logical deletes within D days

delete all data older than T days

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$

**key**                                        **value**

| **RID** | attr-1 | attr-2 | attr-3 | $\cdots$ | attr-n |

**sort key**

sort key = delete key

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$



key | value

| RID | attr-1 | attr-2 | attr-3 | ··· | attr-n |

sort key

sort key = delete key

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$

key                           value

| RID | timestamp | attr-2 | attr-3 | ⋯ | attr-n |

sort key

**delete key**

sort key ≠ delete key

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$



key | value

RID | timestamp | attr-2 | attr-3 | ··· | attr-n

sort key

**delete key**

sort key ≠ delete key

scattered occurrences

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$

key | value

| RID | timestamp | attr-2 | attr-3 | $\cdots$ | attr-n |

sort key

**delete key**

sort key ≠ delete key

scattered occurrences

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$



key                    value

| RID | timestamp | attr-2 | attr-3 | ... | attr-n |

sort key

**delete key**

sort key ≠ delete key

scattered occurrences

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$



**key**          **value**

| **RID** | timestamp | attr-2 | attr-3 | ⋯ | attr-n |

**sort key**

**delete key**

sort key ≠ delete key

latency spikes

superfluous I/Os

# Realizing **Retention-Based Deletes**

delete all entries older than: $TS_x$



sort key ≠ delete ke

"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 $<= 65_D$



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

**file**

| 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 |
|---|---|---|----|----|----|----|----|
| $34_D$ | $69_D$ | $3_D$ | $79_D$ | $8_D$ | $80_D$ | $23_D$ | $24_D$ |

**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 |
|----|----|----|----|----|----|----|----|
| $70_D$ | $41_D$ | $62_D$ | $7_D$ | $25_D$ | $85_D$ | $59_D$ | $19_D$ |

# Realizing **Retention-Based Deletes**

delete all entries older than $<= 65_D$



**file**

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

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

# Intuition: Data Layout holds the key!

# Realizing **Retention-Based Deletes**



# KiWi
Key Weaving storage layout

# Key Weaving storage layout

delete all entries older than <= $65_D$



file

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

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 older than <= $65_D$



file

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

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 older than <= $65_D$



file

partitioned on S

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

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

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

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

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 older than <= $65_D$



file

partitioned on S

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

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 D

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

delete tile 1

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

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

**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 |
|---|---|---|---|---|---|---|---|
| $69_D$ | $74_D$ | $76_D$ | $79_D$ | $80_D$ | $81_D$ | $88_D$ | $90_D$ |

# Key Weaving storage layout

## delete all entries older than <= $65_D$



file

partitioned on S

partitioned on D

drop page

**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 |
|---|----|----|----|----|----|----|----|
| $69_D$ | $74_D$ | $76_D$ | $79_D$ | $80_D$ | $81_D$ | $88_D$ | $90_D$ |

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

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

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

# Key Weaving storage layout

## delete all entries older than $\leq 65_D$

file

partitioned on S

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

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 D

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

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

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

sorted on S

drop page

**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 |
|---|----|----|----|----|----|----|----|
| $69_D$ | $79_D$ | $80_D$ | $88_D$ | $90_D$ | $74_D$ | $76_D$ | $81_D$ |

# Key Weaving storage layout

delete all entries older than <= $65_D$

**file**

partitioned on S

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

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

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

sorted on S

**page 1** — drop page

| 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** — drop page

| 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** — 1 I/O

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

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

superior delete performance

up to 2.5x

# Key Weaving storage layout

get(14)

**file**

partitioned on S

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

delete tile 1
- 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 2
- page 3: $S_{min}=67 :: S_{max}=99$ / $D_{min}=1_D :: D_{max}=62_D$
- page 4: $S_{min}=61 :: S_{max}=94$ / $D_{min}=65_D :: D_{max}=89_D$

partitioned on 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**
- $S_{min}=67 :: S_{max}=99$ / $D_{min}=1_D :: D_{max}=62_D$
- $S_{min}=61 :: S_{max}=94$ / $D_{min}=65_D :: D_{max}=89_D$

sorted on S

**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 |
|----|----|----|----|----|----|----|----|
| $75_D$ | $82_D$ | $67_D$ | $77_D$ | $89_D$ | $65_D$ | $70_D$ | $85_D$ |

# Key Weaving storage layout

get(14)

**file**

partitioned on S

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

**delete tile 1**

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 2**

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

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

partitioned on 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**

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

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

sorted on S

**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 |
|---|---|---|---|---|---|---|---|
| $75_D$ | $82_D$ | $67_D$ | $77_D$ | $89_D$ | $65_D$ | $70_D$ | $85_D$ |

data structure

workload

retention-based delete performance

large delete tiles

KiWi navigates

small delete tiles

lookup performance

$$h^* = \sqrt{\dfrac{f_{\mathrm{SRD}} \cdot \frac{N}{B}}{\phi \cdot L \cdot (f_{\mathrm{EPQ}} + f_{\mathrm{PQ}}) \cdot + L \cdot f_{\mathrm{SRQ}}}}$$

$h^*$ = 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
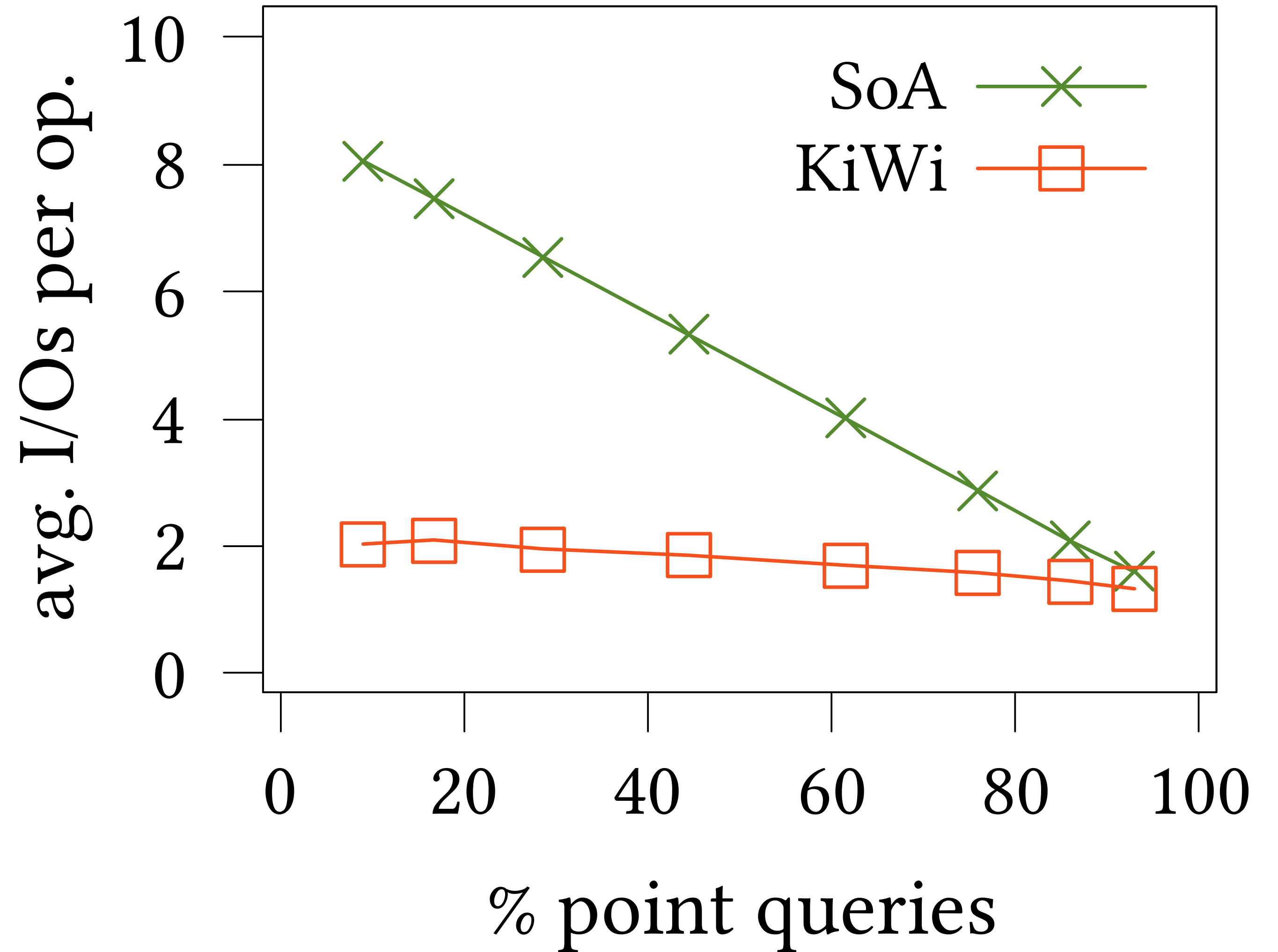
$\phi$ = false positive rate of query filter
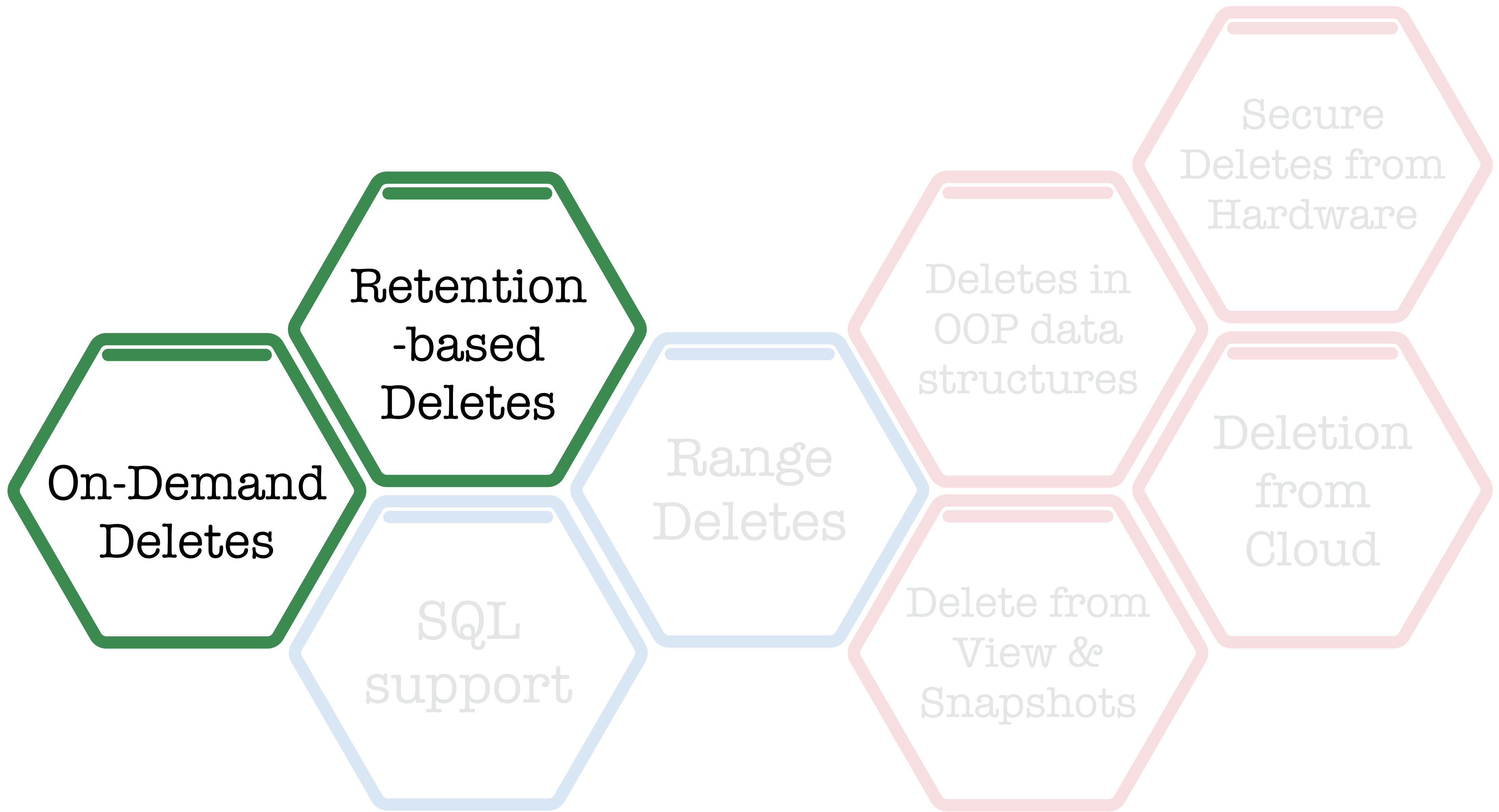
better overall performance

up to 4x

superior delete performance

up to 2.5x

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

avg. I/Os per op.

% point queries

SoA
KiWi

FADE

SQL support

Range Deletes

Deletes in OOP data structures

Secure Deletes from Hardware

Delete from View & Snapshots

Deletion from Cloud

FADE

Deletes in OOP data structures

Secure Deletes from Hardware

Deletion from Cloud

Delete from View & Snapshots