# *LSM-Trees &*
# *its* **Read Optimizations**

*Subhadeep Sarkar*   *Niv Dayan*   *Manos Athanassoulis*

**BOSTON UNIVERSITY**

UNIVERSITY OF **TORONTO**

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

# LSM-tree

# LSM-tree

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

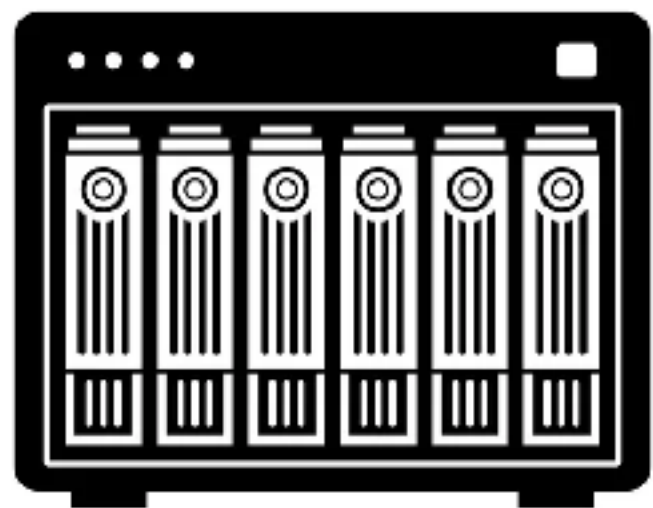**LSM-tree**

O'Neil *et al.*

1996

**LSM-tree**

O'Neil *et al.*

1996

**LSM-tree**
O'Neil *et al.*

Bigtable

APACHE
HBASE

1996

2006 2007

**LSM-tree**
O'Neil *et al.*

1996

2006 2007 2010

**LSM-tree**
O'Neil *et al.*

1996    2006  2007    2010  2011

**LSM-tree**
*O'Neil et al.*

1996               2006  2007    2010  2011  2013

**LSM-tree**
O'Neil *et al.*

1996    2006  2007    2010    2011  2013                    2023

**LSM-tree**
O'Neil *et al.*

Bigtable

APACHE HBASE

cassandra

levelDB

RocksDB

1996    2006 2007    2010 2011 2013    2023

# LSM-tree

NoSQL

RocksDB  WT  levelDB  SCYLLA  riak

cassandra  tarantool  Bigtable  APACHE HBASE  DynamoDB  speedb  accumulo

SQLite

relational

influxdb  QuasarDB

time-series

2023

# LSM-tree

NoSQL

relational

time-series

2023

# Why **LSM** ?
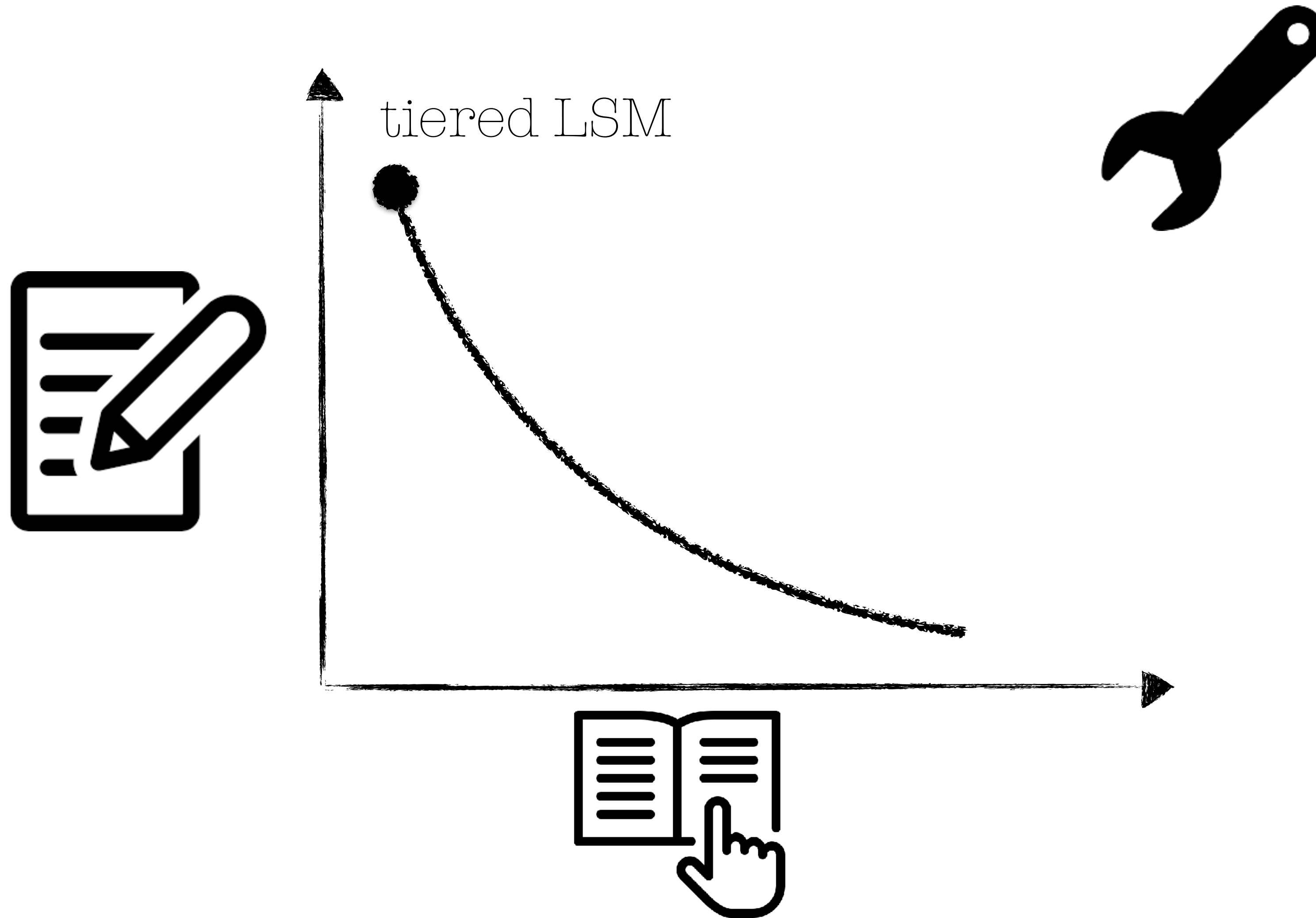


fast ingestion

# Why **LSM** ?
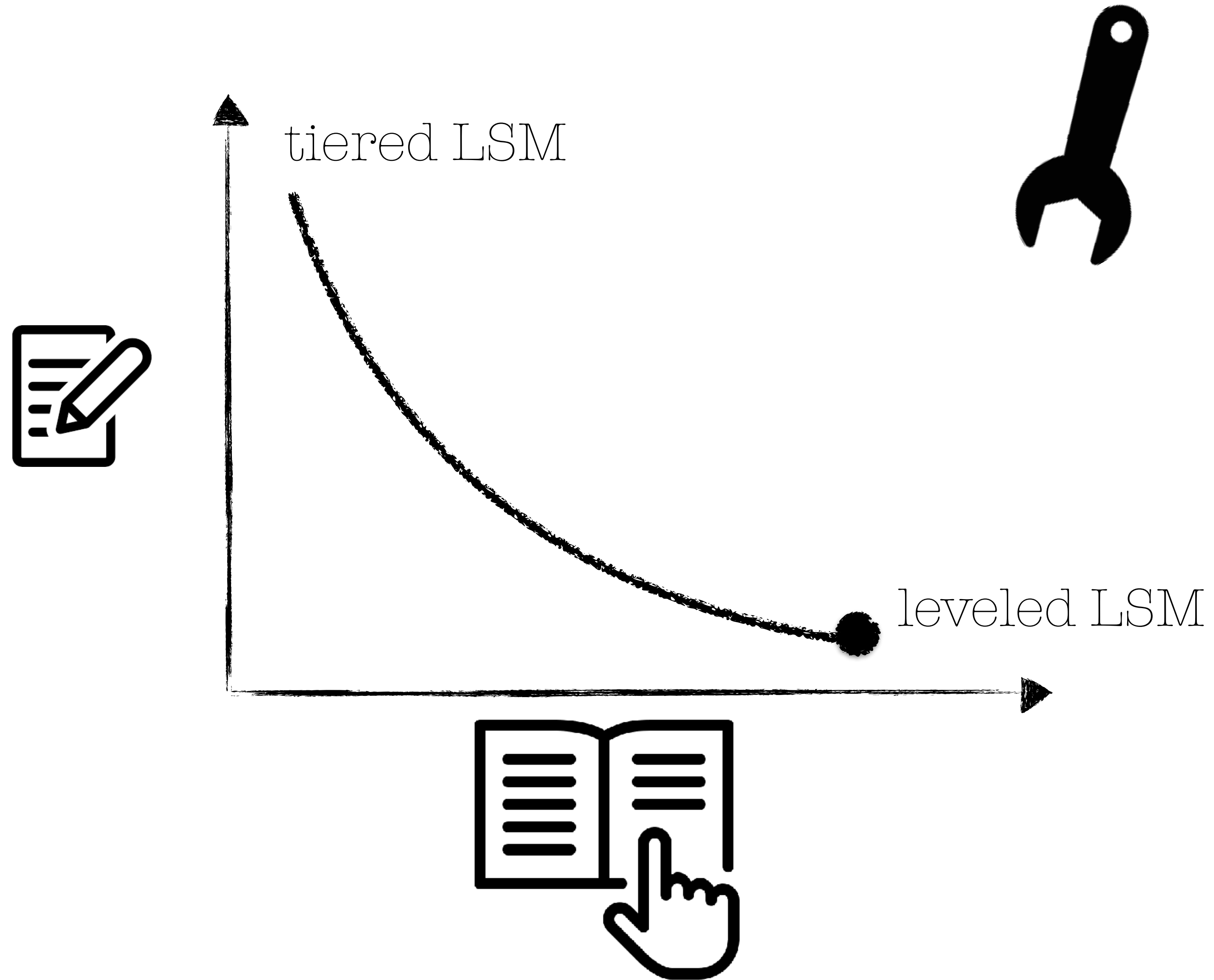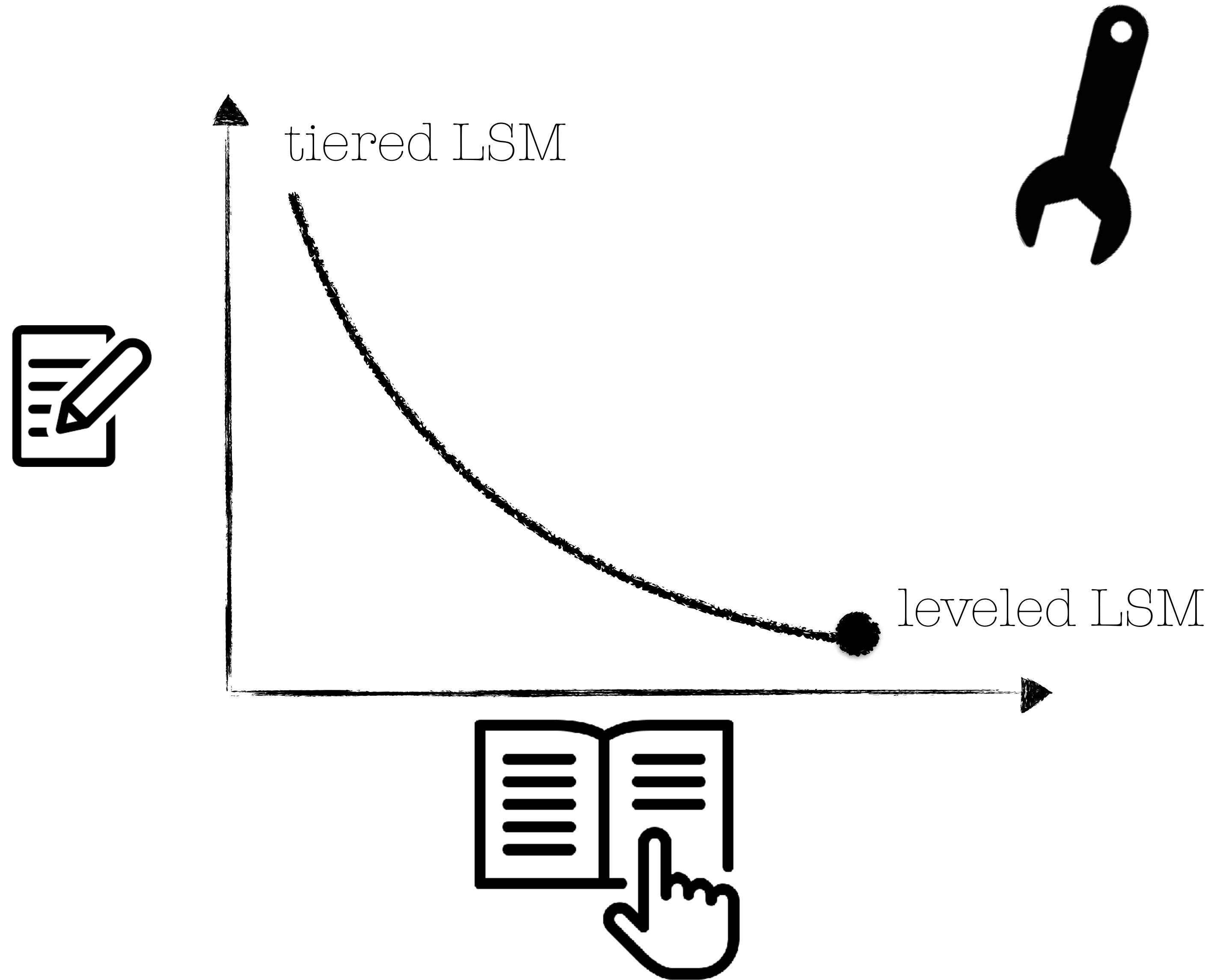
fast ingestion

competitive reads

# Why **LSM** ?

# Why **LSM** ?

# Why **LSM** ?



tiered LSM

# Why **LSM** ?

tiered LSM

leveled LSM

# Why **LSM** ?

tiered LSM

leveled LSM

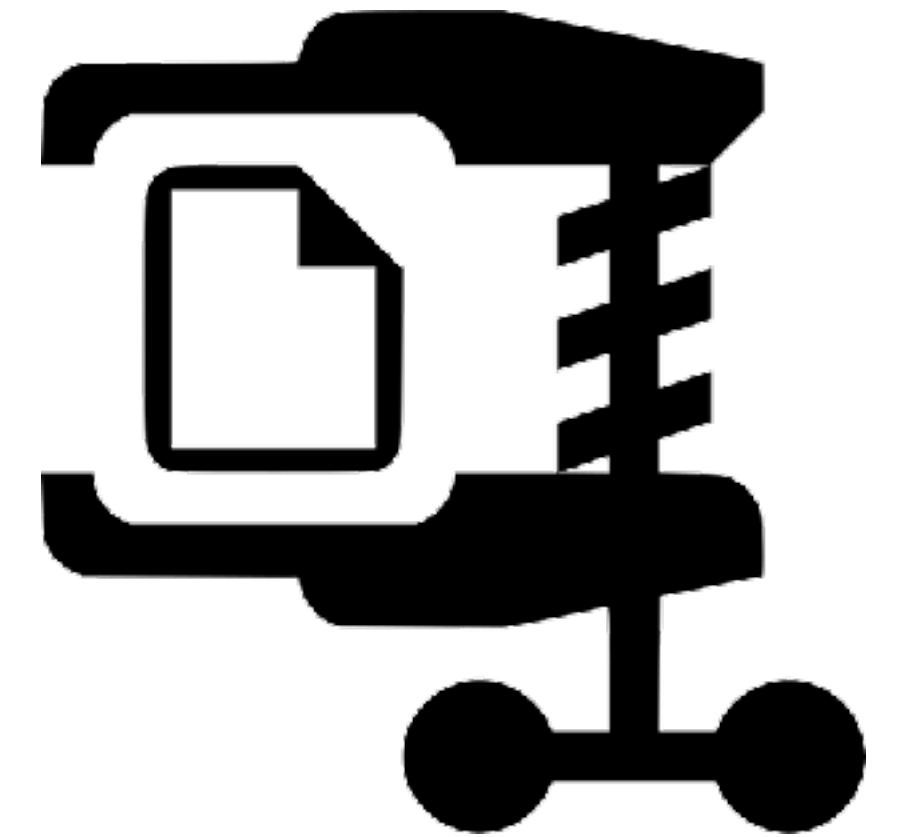# Why **LSM** ?



fast writes
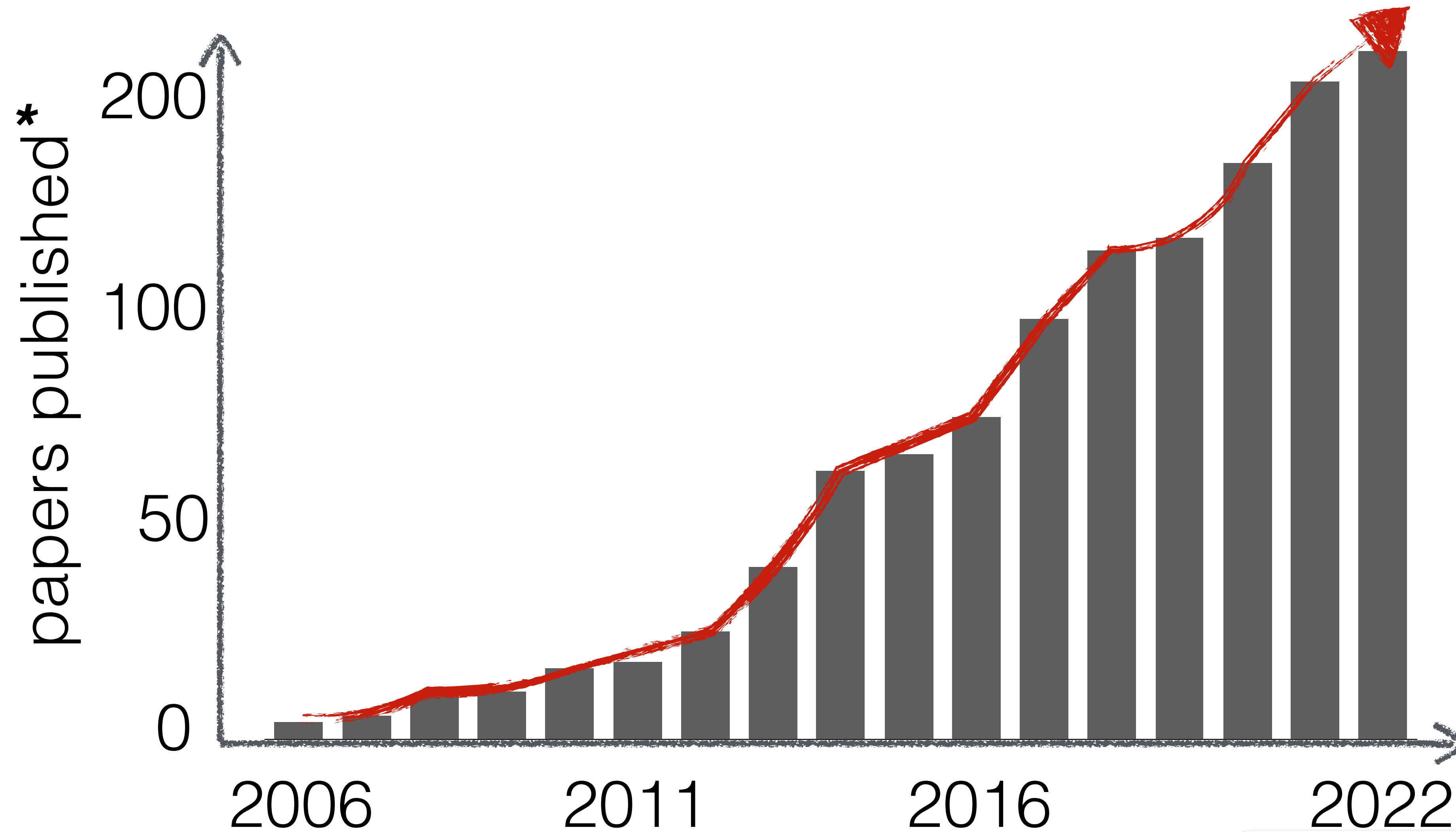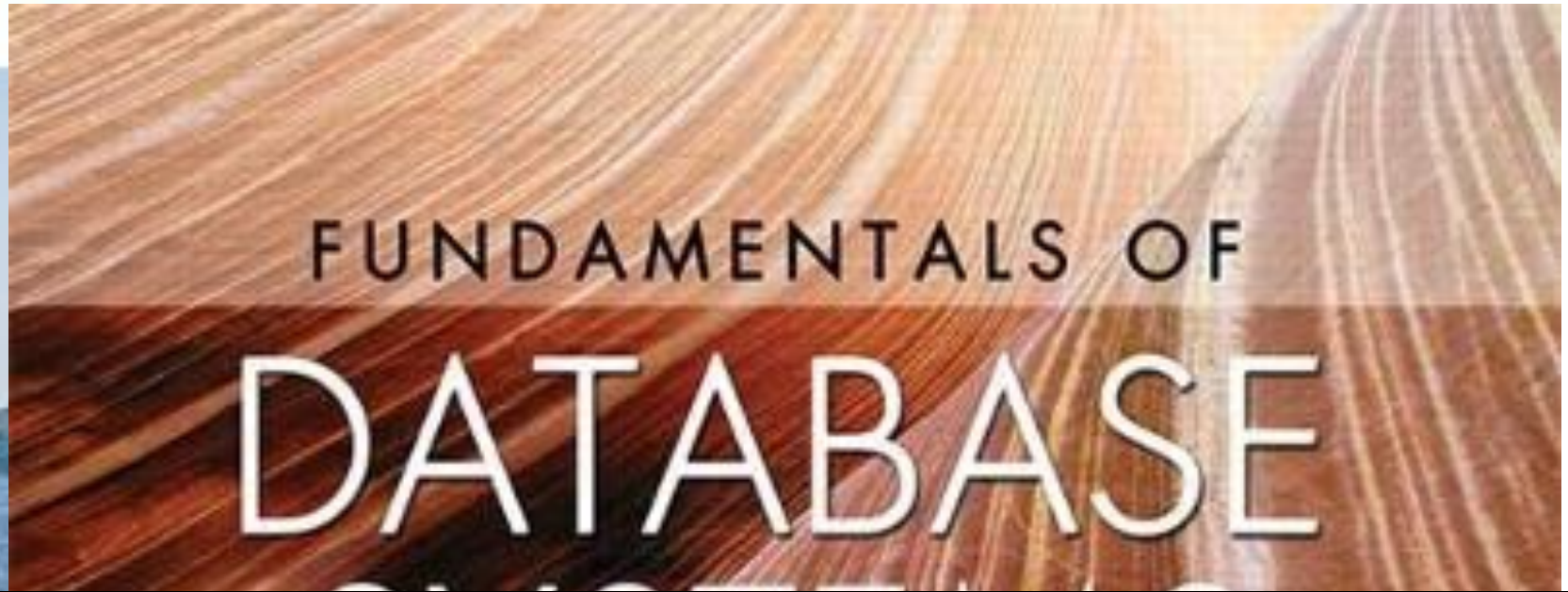
tunable read-write performance
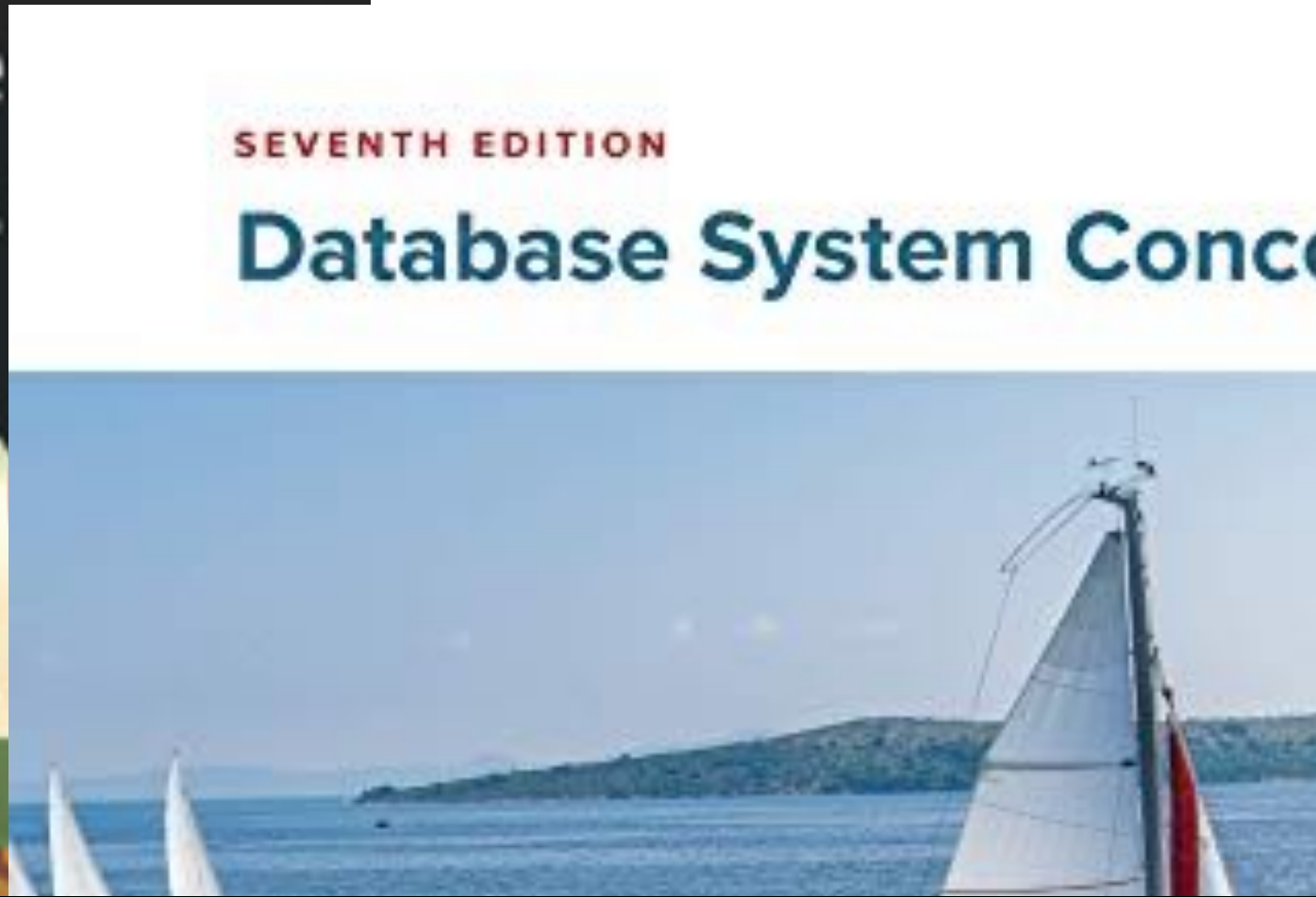
good space utilization

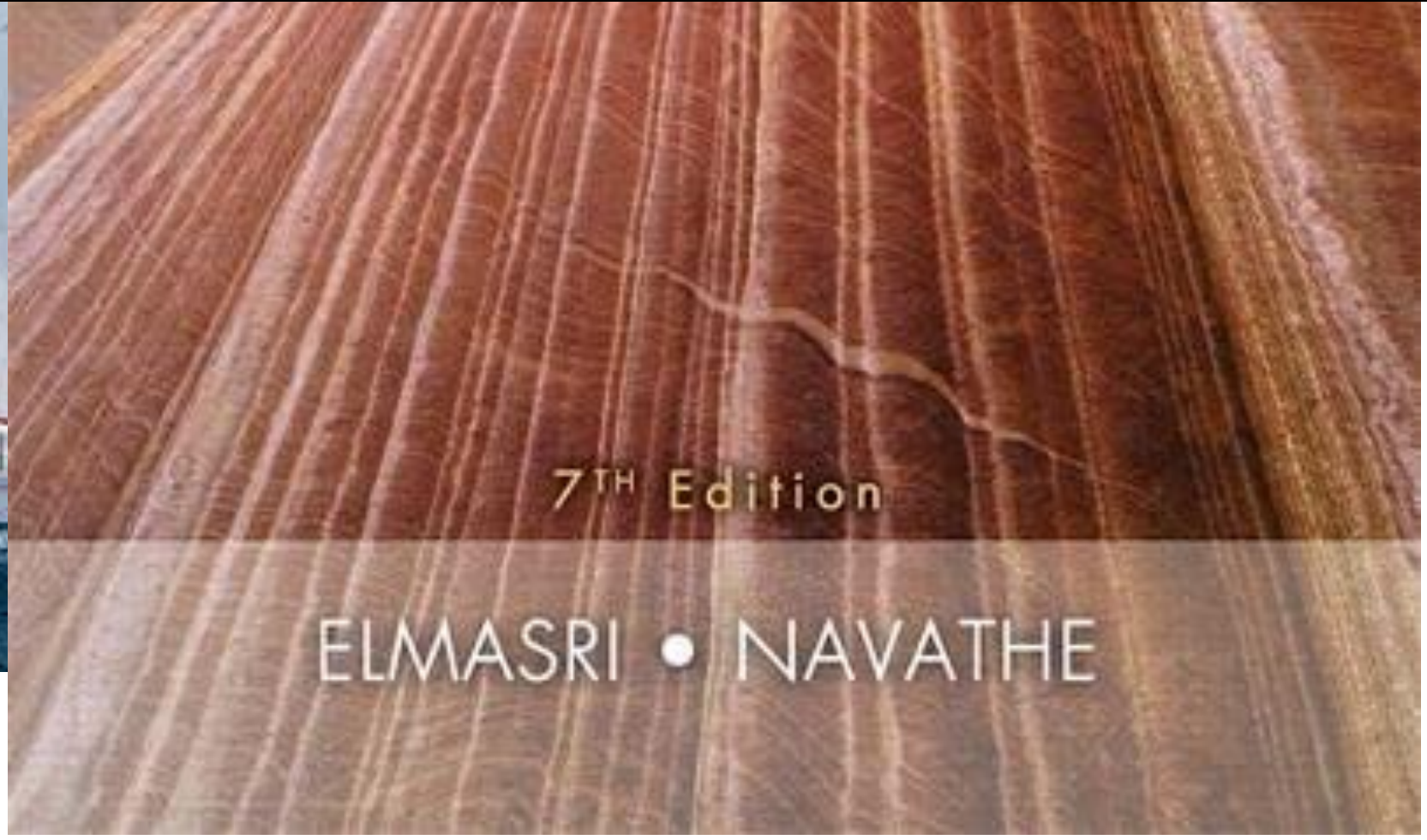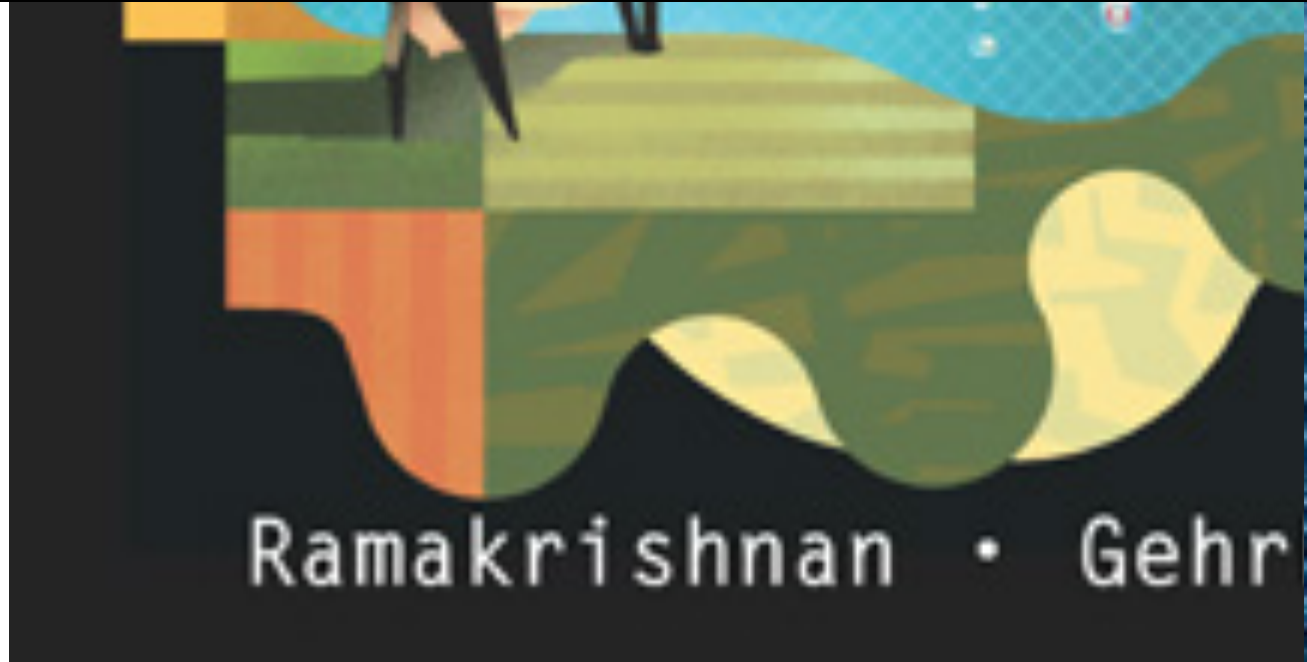# Research Trend



papers published*

200

100

50

0

2006    2011    2016    2022
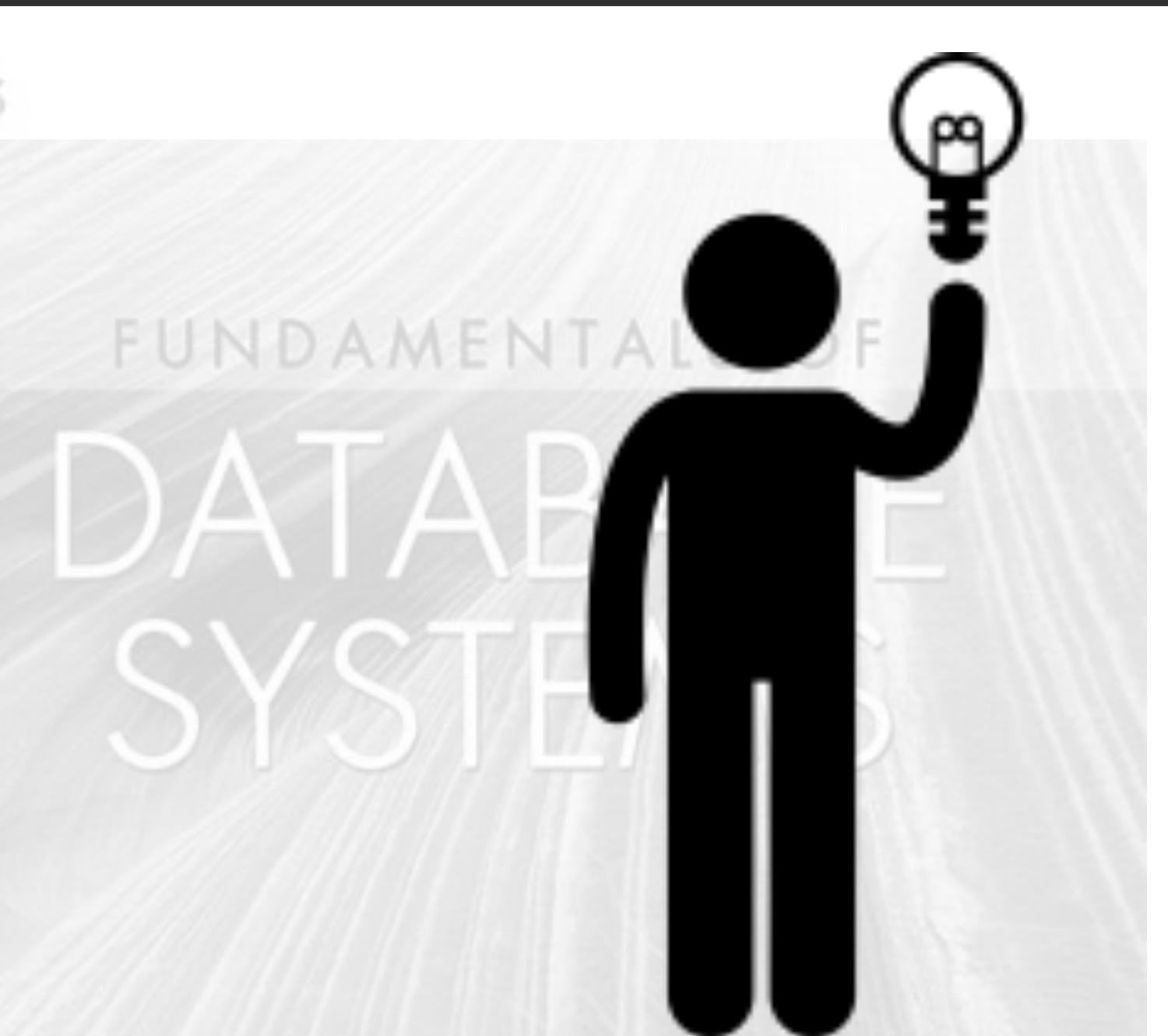
No Textbook on LSMs !!

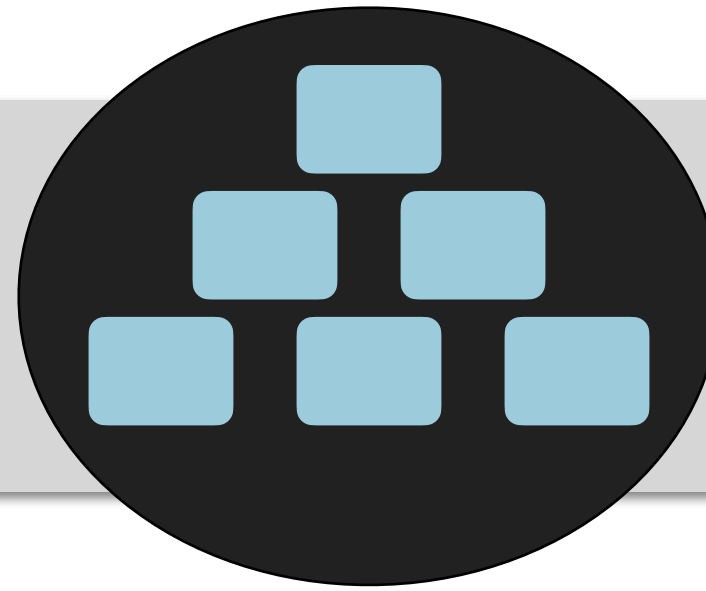# No Textbook on LSMs !!

**explore** the LSM paradigm

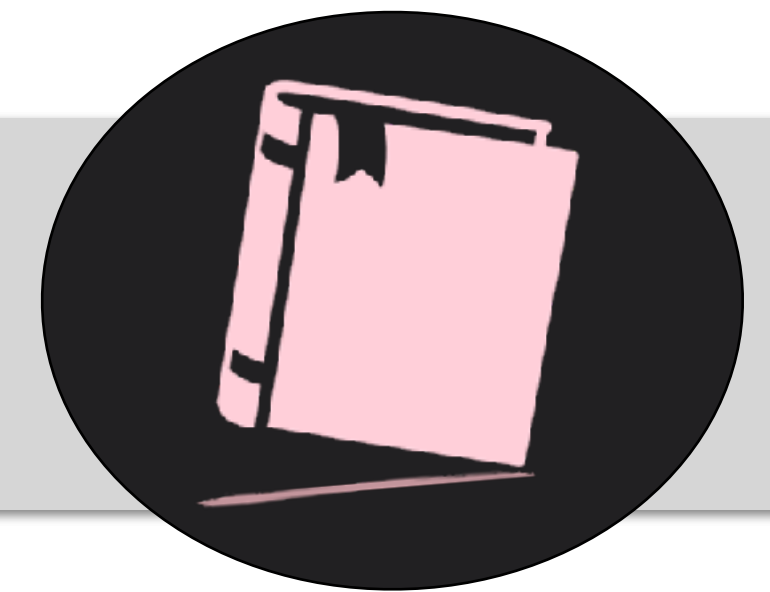understand the **read optimizations**

provide useful **insights**

# Outline

Part 1: **LSM Basics**
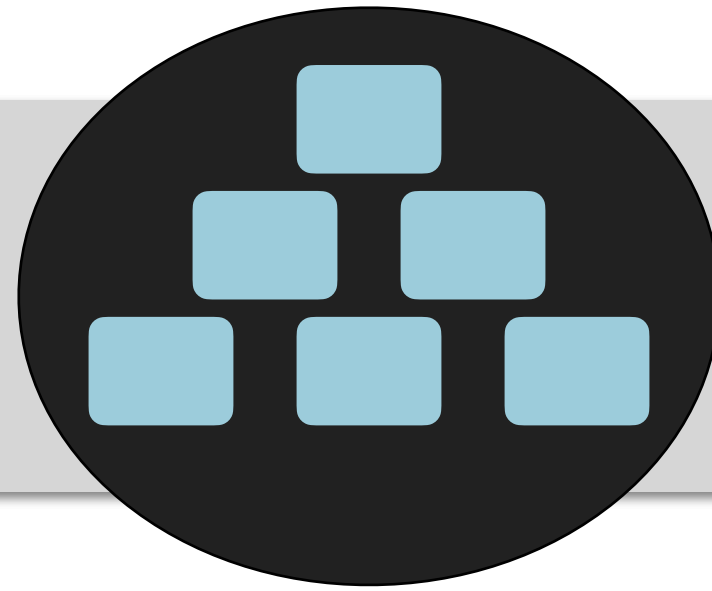
Part 2: **Read Optimizations in LSMs**

Part 3: **Navigating the LSM Design Space**

# Outline

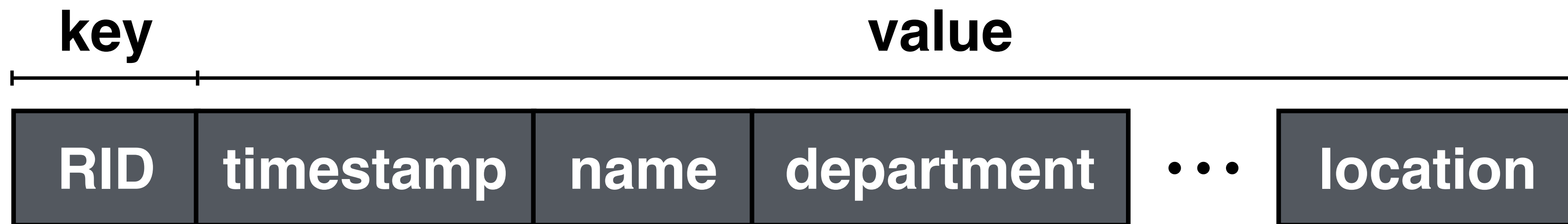Part 1: **LSM Basics**

Part 2: **Read Optimizations in LSMs**

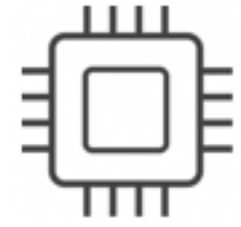Part 3: **Navigating the LSM Design Space**
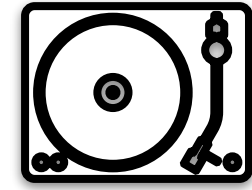
# LSM **Basics**

key-value pairs

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

# LSM **Basics**

key-value pairs

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

# LSM **Basics**
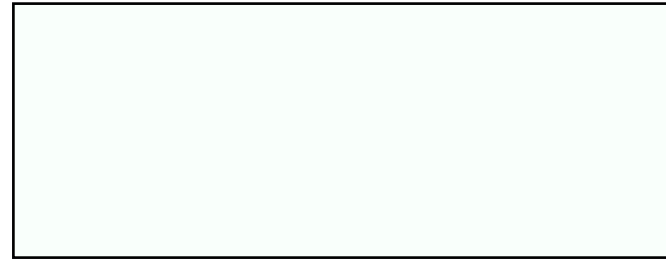
buffer

level 1

level 2

size ratio: T

level 3

level 4

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

buffer

level 1

How do we update data?

put(6)

buffer

level 1    6

buffer

6

level 1

6

logically
invalidated

buffer 6

level 1 6

buffer 6

level 1 6

Out-of-place updates

- fast ingestion
- space amplification
- slow reads

buffer

level 1   6   6

How do we reduce this space amplification?

buffer

level 1    6    6

buffer

level 1    6    compaction

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

buffer

level 1

level 2

level 3

level 4

$M_{buf}$: buffer memory

$T$: size ratio

buffer $M_{buf}$

level 1 $M_{buf} \cdot T$

level 2 $M_{buf} \cdot T^2$

level 3 $M_{buf} \cdot T^3$

How about queries?

$P$: pages in buffer
$B$: entries/page
$L$: #levels
$T$: size ratio
$N$: #entries

get(7)

buffer

Cost analysis
w/o F&I: $\mathcal{O}(log_2 N \cdot log_T N)$

L1

L2

L3

L4

7

$P$: pages in buffer
$B$: entries/page
$L$: #levels
$T$: size ratio
$N$: #entries

get(7)

buffer

L1
L2
L3

Cost analysis
w/o F&I: $\mathcal{O}(log_2 N \cdot log_T N)$

w/ index: $\mathcal{O}(log_T N)$

Can we do better?

$P$: pages in buffer
$B$: entries/page
$L$: #levels
$T$: size ratio
$N$: #entries
$\phi$: FPR of BF

get(7)

buffer

**Cost analysis**
w/o F&I: $\mathcal{O}(log_2 N \cdot log_T N)$
w/ index: $\mathcal{O}(log_T N)$
w F&I: $\mathcal{O}(\phi \cdot log_T N)$

L1
L2
L3

How to manage memory?

buffer

block cache

filters

fence pointers

L1

L2

L3

L4

# Block Cache



buffer

block cache

filters

fence
pointers

L1

L2

L3

L4

# Block Cache



get(7)

buffer

Bloom filters

fence pointers

L1

L2

L3

L4

# Block Cache

get(7)

buffer

Bloom filters

fence pointers

L1

L2

L3

**What about range queries?**

Range Queries

$P$: pages in buffer
$B$: entries/page
$L$: #levels
$T$: size ratio
$N$: #entries
$\phi$: FPR of BF

$s$: selectivity LRQ

buffer

filters

fence pointers

L1
L2
L3
L4

$P$: pages in buffer
$B$: entries/page
$L$: #levels
$T$: size ratio
$N$: #entries
$\phi$: FPR of BF

$s$: selectivity SRQ

## Range Queries

## Cost analysis

long range: $\mathcal{O}(s \cdot N)$
short range: $\mathcal{O}(L)$

get(9,15)

buffer

L1

L2

L3

More on LSM Reads in Part 2.

most data
on storage

$L$: #levels
$T$: size ratio

most data
on storage

if $T = 10$ & $L = 4$

99.9% on storage

# Performance **Tradeoff**

write
performance

**writing data
on storage**

space
amplification

read
performance

# Data **Layout**

Classical LSM design:  leveling

[eager merging]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

# Data **Layout**

**leveling** [eager]

1 run
per level

- good read performance
- good space amplification
- high write amplification

# Data **Layout**

**leveling** [eager]          **tiering** [lazy]

1 run
per level

- good read performance
- good space amplification
- high write amplification

# Data **Layout**

**leveling** [eager]

**tiering** [lazy]

1 run
per level

- 🟢 good read performance
- 🟢 good space amplification
- 🔴 high write amplification

# Data **Layout**

**leveling** [eager]　　　　　　　　**tiering** [lazy]

1 run
per level

- good read performance
- good space amplification
- high write amplification

# Data **Layout**

**leveling** [eager]

**tiering** [lazy]

1 run
per level

- good read performance
- good space amplification
- high write amplification

# Data **Layout**

**leveling** [eager]                    **tiering** [lazy]

1 run per level

T runs per level

- ● good read performance
- ● good space amplification
- ● high write amplification

- ● poor read performance
- ● poor space amplification
- ● good ingestion performance

get(7)

buffer

Reduces
I/O cost
for
lookups

auxiliary data
structures

L1

L2

L3

L4

7

# Outline

# Filters to the Rescue

# What is a filter

# What is a filter

Does X exist? ⟶ 

**Answers set
membership queries**

Set

X Y Z

# What is a filter

Does X exist? → ▼ → Set

X Y Z

**No false negatives**

# What is a filter

Does Q exist? →

Set

X Y Z

**No false
negatives**

**false positives with
tunable probability**

data

X

Filters
**One per run**

get(X)

# more memory → fewer false positives

data

Filters

get(X)

negative

negative

**X**

**true positive**

data

Filters

**Compact & recreate filter**

data

**Bloom Filters**

BloomCommunACM1970

# *k* hash functions

0 0 0 0 0 0 0 0 0

**bitmap**

# insert: Set from 0 to 1 or keep 1

# **negative** lookup: at least one bit is zero

$h_1$  $h_2$

0 0 0 0 **1** 0 0 **0** 1 0

# true or false positive lookup

# Optimal number of hash functions

$$= \ln(2) \cdot M \quad \leftarrow \quad \text{bits / entry}$$

$$h_1 \;\; \ldots \;\; h_k$$

# Optimal number of hash functions

$$= \frac{\ln(2) \cdot M}{h_1 \ \ldots \ h_k}$$

← **bits / entry**

With M bits / entry

Optimal number of hash functions $= \ln(2) \cdot M$

**False positive rate $= 2^{-M \cdot \ln(2)}$**

# 5 fronts

| Holistic Tuning | CPU | Lowering Constants | Unification | Range |

# Holistic Tuning

**Monkey**

**Dostoevsky**

**LSM-Bush**

DayanSIGMOD17

DayanSIGMOD18

DayanSIGMOD19

# **M**onkey: **O**ptimal **N**avigable **Key**-Value Store

data

**Bloom filters**

data

Bloom
filters

**bits/entry**

*M*

*M*

*M*

data

Bloom
filters

**bits/entry**

*M*

*M*

*M*

data

Bloom
filters

**false
positive rate**

$2^{-M \cdot ln(2)}$

$2^{-M \cdot ln(2)}$

$2^{-M \cdot ln(2)}$

data

Bloom filters

**false positive rate**

$2^{-M}$

$2^{-M}$

$2^{-M}$

Bloom
filters

false
positive rate

$2^{-M}$

$2^{-M}$

$2^{-M}$

$= \quad O(2^{-M} \cdot \log_T N)$

Bloom
filters

false
positive rate



$2^{-M}$

$2^{-M}$

$2^{-M}$

$=$     $O(\mathbf{1+2^{-M} \cdot \log_T N})$

Bloom
filters

false
positive rate

$2^{-M}$

$2^{-M}$

$2^{-M}$

$\Bigg\}$ = $O(2^{-M} \cdot \log_T N)$

Bloom
filters

false
positive rate

$2^{-M}$

$2^{-M}$

**most
memory**

$2^{-M}$

Bloom
filters

false
positive rate

$2^{-M}$

$2^{-M}$

**most
memory**

**saves at most 1 access!**

$2^{-M}$

bits / entry

$M$ *+ 2*

$M$ *+ 1*

$M$ *- 1*

**reallocate**

false
positive rates

$2^{-(M + 2)}$ ↓

$2^{-(M + 1)}$ ↓

$2^{-(M - 1)}$ ↑

$\blacktriangledown$   $\propto$   $2^{-M} / T^2$

$\blacktriangledown$   $\propto$   $2^{-M} / T^1$

$\blacktriangledown$   $\propto$   $2^{-M} / T^0$

$$2^{-M} / T^2$$

$$2^{-M} / T^1$$

$$2^{-M} / T^0$$

**geometric progression**

$$= \quad O(\mathbf{2^{-M}})$$

**Faster worst case**

$$\bigcirc(2^{-M}) \quad < \quad \bigcirc(2^{-M} \cdot \log_T N)$$

# Monkey opens up new ways of optimizing write performance without sacrificing get performance



**Monkey** → **Dostoevsky** → **LSM-Bush**

DayanSIGMOD17  DayanSIGMOD18  DayanSIGMOD19

# Dostoevsky

**Smaller false positive rates**

$2^{-M}/T^3$

$2^{-M}/T^2$

$2^{-M}$

**Tiered**

Leveled

**gets**

$O(\ 2^{-M}\ )$

**writes**

$O(\ T + \log_T N\ )$

=

$O(1)$
+
$O(1)$
+
$O(T)$

gets

$O(\ 2^{-M}\ )$

**writes**

$O(\ \textbf{\textit{T} + log}_{\text{T}}\ \textbf{\textit{N}}\ )\ <\ O(\ \textbf{\textit{T} · log}_{\text{T}}\ \textbf{\textit{N}}\ )$

**leveling**

$=$

$O(1)$
$+$
$O(1)$
$+$
$O(T)$

# Dostoevsky

$O(T + \log_T N)$

# **LSM-Bush**

$O(\log_2 \log_T N)$

×16

×4

×2

**more runs**

Monkey w. leveling

Dostoevsky

LSM-Bush

**Cheaper range** ⟷ **Cheaper writes**

Monkey w. leveling　　　Dostoevsky　　　LSM-Bush

**Great point reads all across**

# 5 fronts

| Holistic Tuning | CPU | Lowering Constants | Unification | Range |
|:---:|:---:|:---:|:---:|:---:|

Bloom
filters

CPU overhead?

# Each key is inserted O(T) times per level into a filter

data

Bloom
filters

insert(X)

xT

O(T)

O(T)

O(T)

# Each key is inserted O(T) times per level into a filter

## Each filter insertion uses M · ln(2) hash functions

Each key is inserted O(T) times per level into a filter

Each filter insertion uses M · ln(2) hash functions

**log$_T$(*M*) levels**

data

Bloom filters

insert(X)

**O(T · M)**

**O(T · M)**

**O(T · M)**

**= O(log$_T$(*N*) · T · M)**

# With Monkey more hash functions are used

data

Bloom filters

insert(X)

**log$_T$(*M*) levels**

$O(T \cdot \textbf{(M+2)})$

$O(T \cdot \textbf{(M+1)})$

$O(T \cdot \textbf{M})$

$= O( \log_T(N) \cdot T \cdot (M \textbf{+ log}_T \textbf{\textit{N}}) )$

# How about get cost?

data

Bloom
filters

get(X)

**X**

↑

**Worst case**

# Expected Negative Query Cost ≈ 2

| data | # hashes | Bloom filters | get(X) |
|------|----------|---------------|--------|
| | 2 | | |
| | 2 | | |
| X | | | |

# Positive Query Cost ≈ M · ln(2)

data     # hashes     Bloom filters     get(X)

2

2

**X**     O(M)

**Avg. worst case = O( M + $\log_T N$ )**

**get**
$O(\ M + \log_T N\ )$

Bloom
filters

**Insert**
$O(\ T \cdot M \cdot \log_T N\ )$

# Address Using Blocking and SIMD

get
$O(\ M + \log_T N\ )$

Bloom
filters

Insert
$O(\ T \cdot M \cdot \log_T N\ )$

# Blocking

PutzeJEA10

**X**

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**Bloom filter**

# Blocking

**Hash to one cache line**

X

0 0 0 0 0  0 0 0 0 0  0 0 0 0 0  0 0 0 0 0  0 0 0 0 0  0 0 0 0 0

# Blocking

X

**Insert as though cache line is an independent Bloom filter**

0 0 0 0 0   0 0 0 0 0   0 0 0 0 0   0 0 0 0 0   0 0 0 0 0   0 0 0 0 0

# Blocking

**Pro: one cache miss per insert/get**

X

# Blocking

**Con 1: uneven distribution of entries across cache lines slightly harms the false positive rate**

Blocked Bloom filter

# Blocking

**Con 2: still need to compute many hash functions per entry**

X

O O O O O

Blocked Bloom filter

# SIMD

PolychroniouDAMON14

JianyuanTPDS18

## Partition into sub-lines

Cache line

# SIMD

**Map one hash per sub-line**

X

Cache line

# SIMD

**Insert( X )**
bitwise **"or"**

Cache line

# SIMD

**get( X )**
bitwise **"and"**

Cache line

# Blocking and SIMD

get
$O(\, M + \log_T N \,)$

Bloom
filters

Insert
$O(\, T \cdot M \cdot \log_T N \,)$

# Blocking and SIMD

get
**O( $\log_T N$ )**

Bloom
filters

Insert
**O( $T \cdot \log_T N$ )**

# 5 fronts

Holistic
Tuning

CPU

**Improving
Constants**

Unification

Range

# Improving Constants

Bloom

Ideal

False positive rate

$\approx 2^{-M \cdot \ln(2)}$

$\approx 2^{-M}$

Bloom

Ideal

False
positive rate

$\approx 2^{-M \cdot \mathbf{0.69}}$

$\mathbf{\approx 2^{-M}}$

**Can we improve this?**

# Bloom



$\approx 2^{-M} \cdot 0.69$

# XOR



GrafJEA20

# Ideal



$\approx 2^{-M}$

# XOR Filter

## Assign one bucket to own each entry

**X**  Y  X  X **Y**  Y

**Owns X**  **Owns Y**

1  2  3  4  5  6

# XOR Filter

**Each bucket stores XOR of fingerprint and other two buckets**

# XOR Filter

**During queries, recover fingerprints by xoring three buckets**

**get(Y) returns true if**    **FP(Y) =**   **2** $\oplus$ **4**   $\oplus$   **6**

Y                     Y                     Y

**Owns Y**

**2**                    4                    **6**

# XOR Filter

**free space ensures each bucket can own one entry**

Owns X   Owns Z   **Free**   Owns Y   **Free**   Owns Q

# Bloom

$$\approx 2^{-M \cdot 0.69}$$

# **XOR**

$$\approx 2^{-M \cdot 0.81}$$

# Idealized

$$\approx 2^{-M}$$

Bloom  XOR  **Ribbon**  Idealized

$\approx 2^{-M \cdot 0.69}$  $\approx 2^{-M \cdot 0.81}$  $\boldsymbol{\approx 2^{-M \cdot 0.92}}$  $\approx 2^{-M}$

**Denser XOR filter**

| Bloom | XOR | **Ribbon** | Idealized |
|:---:|:---:|:---:|:---:|

$$\approx 2^{-M \cdot 0.69} \qquad \approx 2^{-M \cdot 0.81} \qquad \mathbf{\approx 2^{-M \cdot 0.92}} \qquad \approx 2^{-M}$$

Denser XOR filter

**In RocksDB since 2020**

**Bloom**

**XOR**

**Ribbon**

Lower CPU

Lower false
positive rate

5 fronts

| Holistic Tuning | CPU | Improving Constants | **Unification** | Range |

# Unification

**Chucky**

DayanSIGMOD21

**SlimDB**

RenVLDB17

# Unification

read

Bloom
filters

O( 1 )

O( 1 )

O( 1 )

$= O( \log_T N )$

O( $T$ )

O( $T$ )

O( $T$ )

**X**

**= O( T · log$_T$ N )**

hash table

key → Level ID

SSD

**hash(🔑) =**

hash( 🔑 ) =

**cuckoo filter**

$O(\log_2 \log_T M)$

**Binary encoding**

...                    ...

010

001                    $O(\log_2 \log_T M)$

000

**e.g., unary encoding**

110

10

0

O(1)

|  | **Monkey w. Bloom** | **Chucky w. Cuckoo** |
|---|---|---|
| Get I/O | $O(1+2^{-M \cdot \ln(2)})$ | $O(1+2^{-M+3})$ |
| **Get CPU** | $O(\log_T N)$ | $O(1)$ |

|  | **Monkey w. Bloom** | **Chucky w. Cuckoo** |
|---|---|---|
| Get I/O | $O(1 + 2^{-M \cdot ln(2)})$ | $O(1 + 2^{-M+3})$ |
| Get CPU | $O(\log_T N)$ | $O(1)$ |
| Insert CPU | $O(\mathbf{T \cdot \log_T N})$ | $O(\mathbf{\log_T N})$ |

|  | **Chucky** | **SlimDB** |
|---|---|---|
| **Get I/O** | $O(1 + 2^{-M+3})$ | **O(1)** |
| **Memory** | M | $M + 2^{-M} \cdot \log_2(N)$ |
|  |  | **Keeps full key in memory whenever fingerprints collide** |

# 5 fronts

Holistic Tuning

CPU

Improving Constants

Unification

**Range**

# Range Filtering

Traditional filters do not support ranges

get(x, y)

cost:  $O(\log_T M)$

# Range Filters



**Prefix Filter**

RocksDB20

**Surf**

ZhangSIGMOD18

**Rosetta**

LouSIGMOD21

# Prefix Filter

**Users define prefix extraction method**

# Prefix Filter

Users define prefix extraction method

**Country code**

↓

**USA**1234

**CAN**9876

# Prefix Filter

Users define prefix extraction method

USA

CAN

**Insert
prefixes** →

# Prefix Filter

Users define prefix extraction method

USA
CAN

Insert
prefixes →

Check
**USA**
←

get(**USA**0, **USA**9)?

# Prefix Filter

Users define prefix extraction method

USA

CAN

Insert prefixes →

Check USA ←

get(USA0, USA9)?

**Non-generic and requires API extension**

**Surf**

**A trie of all keys**

# Surf

A trie of all keys

**Truncated to reduce space**

# Surf

A trie of all keys

Truncated to reduce space



**Add fingerprint for point reads**

**Surf**

A trie of all keys

Truncated to reduce space

**Encoded as succinct trie
with rank & select**

Add fingerprint for point reads

# Rosetta

**Insert(ICDE)**

I

IC

ICD

ICDE

**Add all prefixes of all keys to a Bloom filter**

# Rosetta

**get(ICD**E**, ICD**F**)** $\longrightarrow$ **ICD** $\longrightarrow$

**Check largest common prefixes**

# Rosetta

get(ICD**E**, ICD**F**)

**ICD**

**ICDE**

**ICDF**

Check largest common prefixes

**Add more fine-grained checks to reduce false positive rate**

Surf

Rosetta

**Better long range** ⟷ **Better short range**

**Range Filters**

Prefix Filters
RocksDB20

Surf
ZhangSIGMOD18

Rosetta
LouSIGMOD21

**Remix**
ZhongFAST21

**Snarf**
VaidyaVLDB22

**Proteus**
KnorrSIGMOD22

**BloomRF**
MößnerEDBT23

**REncoder**
WangICDE23

# Outline

Part 1: **LSM Basics**

Part 2B: **Read Optimizations in LSMs**

Part 3: **Navigating the LSM Design Space**

# **Reducing CPU Overheads** in LSMs

For every query ...

L0  MemBuf   ←   get(k)

L1  SST A

L2  [ ]  SST B

L3  [ ] [ ] SST C [ ]

# **Reducing CPU Overheads** in LSMs

For every query …



The same hash function is calculated **O(L)** times

# **Reducing CPU Overheads** in LSMs

For every query ...



The same hash function is calculated **O(L)** times

Each key is hashed **O(1)** times

ZhuDaMoN21

# Filters under **Memory Pressure**



data size ↑

*for 1TB data,*
*1.3GB filter &17.2GB index*

*1KB entry, 64B key, BPK=10*

*price drop from **2010** to **today***
*SSD: 60x        DRAM: 10x*

# Filters under **Memory Pressure**

Even in a **perfectly uniform** workload,
**80% of the queries** access **45% of the files**

# Filters under **Memory Pressure**

For a **skewed** workload,
**80% of the queries** access less than **5% of the files**

# Filters under **Memory Pressure**

$$v$$

$$h_1(v) \qquad h_2(v) \qquad h_3(v)$$

# Filters under **Memory Pressure**

MunADMS22

LiATC19

$$v$$

$$h_1(v) \qquad h_2(v) \qquad h_3(v)$$

module #1      module #2      module #3

Modular Bloom filter is a collection of smaller Bloom filters

Elastic Bloom filter also works based on the same principle

# Filters under **Memory Pressure**

buffer

modules

L1

L2

L3

L4

# Filters under **Memory Pressure**

buffer

L1

L2

L3

L4

# Filters under **Memory Pressure**

buffer

L1

L2

L3

L4

# Filters under **Memory Pressure**

buffer

L1

L2

L3

L4

# Filters under **Memory Pressure**

buffer

L1

L2

L3

L4

# Filters under **Memory Pressure**



MunADMS22

buffer

L1

L2

L3

L4

**Overall, better performance with smaller memory budget**

# **Compactions** and **Caching**

# **Compactions** and **Caching**

filters    fence pointers

buffer

L1

L2

L3

L4

Leaper : A learned pre-fetcher that improves reads

YangVLDB20

# Outline

Part 1: **LSM Basics**

Part 2: **Read Optimizations in LSMs**

Part 3: **Navigating the LSM Design Space**

# LSM **Design Space**

# LSM **Design Space**

Update cost

Read cost

Memory/space footprint

AthanassoulisEDBT16

# LSM **Design Space**

Update cost

Read cost

Memory/space footprint

fixed Memory

Read cost

Update cost

LSM designs

# LSM **Design Space**

# LSM **Design Space**

Update cost

Read cost

Memory/space footprint

fixed Memory

Read cost

Update cost

**Navigating the RUM tradeoff**

# LSM **Design Space**



Update cost

Read cost

Memory/space

fixed Memory

Read cost

**Navigating the RUM tradeoff**

How to optimally allocate the available memory?

$M$: total memory
$M_{idx}$: index memory
$M_{ftr}$: filter memory
$M_{buf}$: buffer memory

# The **Optimal** Memory Allocation

available
memory

$$M = M_{idx} + M_{ftr} + M_{buf}$$

index    $M_{idx}$

**workload**

filter    $M_{ftr}$

$$read\_cost(M_{idx}, M_{ftr}, M_{buf})$$

$reads(R)$

vs.

$write\_cost(M_{buf})$

$writes(W)$

$$cost = R \cdot read\_cost + W \cdot write\_cost$$

buffer    $M_{buf}$

# The **Optimal** Memory Allocation

available
memory

$$M = M_{cache} + M_{buf}$$

block
cache $\quad M_{cache}$

## workload

$reads(R)$
vs.
$writes(W)$

$read\_cost(M_{cache})$

$write\_cost(M_{buf})$

$cost = R \cdot read\_cost + W \cdot write\_cost$

buffer $\quad M_{buf}$

# The **Optimal** Memory Allocation

available
memory

Update cost

index

filter

block

cache



Read cost

Memory/space
footprint

Navigating read vs. writes: **data layouts**

# Data **Layout**

**leveling** [eager]  **tiering** [lazy]

1 run per level

T runs per level

# Data **Layout**



leveling        1-leveling        L-leveling        tiering

read
optimized

write
optimized

# Storage Layer **Design Continuum**



leveling      1-leveling      L-leveling      tiering

Any design can be defined by the tuple-set: $(T, i)$

# Storage Layer **Design Continuum**

leveling $(T, 0)$

1-leveling $(T, 1)$

L-leveling $(T, L)$

tiering $(T, L+1)$



Any design can be defined by the tuple-set: $(T, i)$

# Storage Layer **Design Continuum**

$(T, 0)$     $(T, 1)$     $\bullet\bullet\bullet$     $(T, i)$     $\bullet\bullet\bullet$     $(T, L)$     $(T, L+1)$

# Storage Layer **Design Continuum**

variable #leveled / #tiered levels

SarkarVLDB21

# Storage Layer **Design Continuum**

variable #leveled / #tiered levels

size ratio

T

T

T

T

# Storage Layer **Design Continuum**

variable #leveled / #tiered levels

size ratio

T

T

T

T

# Storage Layer **Design Continuum**



variable #leveled / #tiered levels

| size ratio | #runs |
|:---:|:---:|
| T | 4 |
| T | 3 |
| T | 2 |
| T | 1 |

# Storage Layer **Design Continuum**

variable **#runs/level**

variable **#leveled / #tiered levels**

| size ratio | #runs |
|:---:|:---:|
| T | 4 |
| T | 3 |
| T | 2 |
| T | 1 |

Storage Layer **Design Continuum**

variable **#runs/level**

variable **#leveled / #tiered levels**

| size ratio | #runs |
|---|---|
| 2 | 4 |
| 2.5 | 3 |
| 3 | 2 |
| 4 | 1 |

DayanSIGMOD19

# Storage Layer **Design Continuum**

variable **#runs/level**

variable **#leveled / #tiered levels**

variable **size ratio**

The LSM storage layer
design continuum

workload

performance
target

performance
modeling

LSM designs

workload

performance target

**worst-case** performance modeling

LSM designs

**worst-case**
performance
modeling

**worst-case**
performance
modeling

**worst-case** read cost: $1 + \displaystyle\sum_{i=1}^{L-1} \phi_i$

7

**worst-case**
performance
modeling

**average-case**
performance
modeling

7

$$\sum_{i=1}^{L} (\mathbb{P}[\text{query in } L_i] \cdot (1 + \sum_{j=1}^{i-1} \phi_i))$$

**worst-case** read cost: $1 + \sum_{i=1}^{L-1} \phi_i$

HuynhVLDB22      ChatterjeeVLDB22

**average-case** performance modeling

workload

$ budget

Cosine

optimal configuration

SE design · h/w · cloud provider

ChatterjeeVLDB22

workload

average-case

What if the workload comes with **unpredictability**?

# Workload-based Tuning



$$\pi_{w_0} = argmin_\pi(cost(w_0, \pi)) \qquad cost(w_0, \pi_{w_0})$$

optimal configuration for $w_0$

# **Nominal** Tuning



$$\pi_{w_0} = argmin_\pi(cost(w_0, \pi))$$

optimal configuration for $w_0$

$$cost(w_0, \pi_{w_0})$$

**same configuration**

**due to unpredictability**

$$\pi_{w_0}$$

$$cost(w_1, \pi_{w_0})$$

**... but not optimal!**

# **Robust** Tuning



$$\pi_{w_0}^{robust}$$

$$= argmin_\pi max_{w\prime \in R(w_0)}(cost(w_0, \pi_{w_0}^{robust}))$$

$$cost(w_1, \pi_{w_0}^{robust})$$

**… close-to-optimal!**

# **Robust** Tuning



Nominal     Robust     Robust (for more noise)

Speedup

4x
3x
2x
1x

**10K workloads**

**observed drift**    **observed drift**    **observed drift**

expected workload drift

HuynhVLDB22

# The **Key Takeaways**

The LSM design space is **vast and complex**.

**Read optimizations are crucial** to make LSMs better.

A **tuned LSM** engine can offer superior performance.

# **Open** Research **Challenges**

Reduce write amplification

Workload-aware compactions & layout transformation

Performance Stability & Holistic Tuning

Automatic Tuning & Adaptive Behavior

Privacy-aware LSM designs

# Please see our manuscript for all references!

## REFERENCES

[1] H. Abu-Libdeh, D. Altınbüken, A. Beutel, E. H. Chi, L. Doshi, T. Kraska, Xiaozhou, Li, A. Ly, and C. Olston. Learned Indexes for a Google-scale Disk-based Database. In *Proceedings of the Workshop on ML for Systems at NeurIPS*, 2020.

[2] S. Alsubaiee, Y. Altowim, H. Altwaijry, A. Behm, V. R. Borkar, Y. Bu, M. J. Carey, I. Cetindil, M. Cheelangi, K. Faraaz, E. Gabrielova, R. Grover, Z. Heilbron, Y.-S. Kim, C. Li, G. Li, J. M. Ok, N. Onose, P. Pirzadeh, V... AsterixDB: A S... *Endowment*, 7(...

[3] Apache. Accum...

[4] Apache. HBase...

[5] Apache. Cassar...

[6] M. Athanassou... Methods. In *Pr... on Managemen...

[7] M. Athanassou... A. Ailamaki, a... RUM Conjectu... *Extending Data...

[8] O. Balmau, F. ... and D. Didona... Merge Key-Va... *Trans. Comput...

[9] O. Balmau, R.... Unlocking Men... the *ACM Europ... 80–94, 2017.

[10] M. A. Bender, ... maul, D. Medje... Don't Thrash: ... *VLDB Endowm...

[11] A. Breslow an...

[22] N. Dayan, Y. Rochman, I. Naiss, S. Dashevsky, N. Rabinovich, E. Bortnikov, I. Maly, O. Frishman, I. B. Zion, Avraham, M. Twitto, U. Beitler, E. Ginzburg, and M. Mokryn. The End of Moore's Law and the Rise of The Data Processor. *Proceedings of the VLDB Endowment*, 14(12):2932–2944, 2021.

[23] N. Dayan and M. Twitto. Chucky: A Succinct Cuckoo Filter for LSM-Tree. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 365–378, 2021.

[24] N. Dayan, T. Weiss, S. Dashevsky, M. Pan, E. Bortnikov, and M. Twitto.

[44] J. Kim, S. Lee, and J. S. Vetter. PapyrusKV: a high-performance parallel key-value store for distributed NVM architectures. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, pages 57:1—-57:14, 2017.

[45] Y.-S. Kim, T. Kim, M. J. Carey, and C. Li. A Comparative Study of Log-Structured Merge-Tree-Based Spatial Indexes for Big Data. In *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, pages 147–150, 2017.

[46] A. Kipf, R. Marcus, A. van Renen, M. Stoian, A. Kemper, T. Kraska, and T. Neumann. RadixSpline: a single-pass learned index. In *Proceedings of the International Workshop on Exploiting Artificial Intelligence Techniques for Data Management (aiDM@SIGMOD)*, pages 5:1—-5:5, 2020.

[47] H. Kondylakis, N. Dayan, K. Zoumpatianos, and T. Palpanas. Coconut Palm: Static and Streaming Data Series Exploration Now in your Palm. In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD)*, 2019.

[48] T. Kraska, A. Beutel, E. H. Chi, J. Dean, and N. Polyzotis. The Case for Learned Index Structures. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 489–504, 2018.
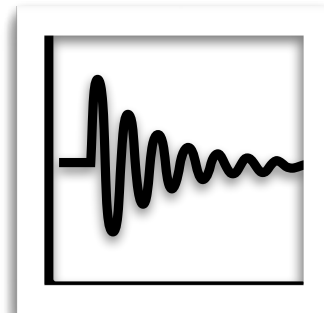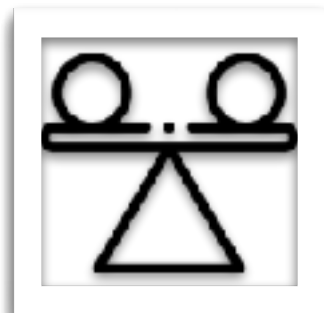
[49] Y. Li, Z. Liu, P. P. C. Lee, J. Wu, Y. Xu, Y. Wu, L. Tang, Q. Liu, and Q. Cui. Differentiated Key-Value Storage Management for Balanced I/O Performance. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 673–687, 2021.

[67] P. Raju, R. Kadekodi, V. Chidambaram, and I. Abraham. PebblesDB: Building Key-Value Stores using Fragmented Log-Structured Merge Trees. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP)*, pages 497–514, 2017.

[68] K. Ren, Q. Zheng, J. Arulraj, and G. Gibson. SlimDB: A Space-Efficient Key-Value Storage Engine For Semi-Sorted Data. *Proceedings of the VLDB Endowment*, 10(13):2037–2048, 2017.

[69] RocksDB. Leveled Compaction. *https://github.com/facebook/rocksdb/wiki/Leveled-Compaction*, 2020.

[70] RocksDB. Prefix Bloom Filter. *https://github.com/facebook/rocksdb/wiki/Prefix-Seek#configure-prefix-bloom-filter*, 2020.

[71] RocksDB. Block Cache. *https://github.com/facebook/rocksdb/wiki/Block-Cache*, 2021.

[72] S. Sarkar and M. Athanassoulis. Dissecting, Designing, and Optimizing LSM-based Data Stores. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2022.

[73] S. Sarkar, J.-P. Banâtre, L. Rilling, and C. Morin. Towards Enforcement of the EU GDPR: Enabling Data Erasure. In *Proceedings of the IEEE International Conference of Internet of Things (iThings)*, pages 1–8, 2018.

[74] S. Sarkar, K. Chen, Z. Zhu, and M. Athanassoulis. Compactionary: A Dictionary for LSM Compactions. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2022.

[75] S. Sarkar, T. I. Papon, D. Staratzis, and M. Athanassoulis. Lethe:

# LSM-Trees &
# *its* **Read Optimizations**

*Subhadeep Sarkar*    *Niv Dayan*    *Manos Athanassoulis*

Thank You!

Questions?

BOSTON UNIVERSITY

UNIVERSITY OF TORONTO