

LSM-Trees & its Read Optimizations

Subhadeep Sarkar

Niv Dayan

Manos Athanassoulis



UNIVERSITY OF
TORONTO

Log-Structured Merge-tree

LSM-tree

LSM-tree

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

1996

Patrick O'Neil¹, Edward Cheng²
Dieter Gawlick³, Elizabeth O'Neil¹
To be published: Acta Informatica

LSM-tree
O'Neil *et al.*

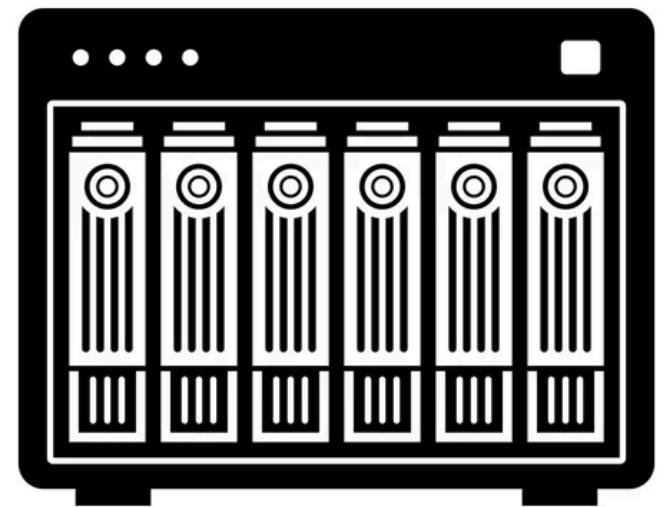


1996

A horizontal dashed grey line spans the width of the slide. A solid dark grey circle is positioned on the line at the right edge, with a short vertical line segment extending upwards from its top, indicating a specific year in the timeline.

- good random writes

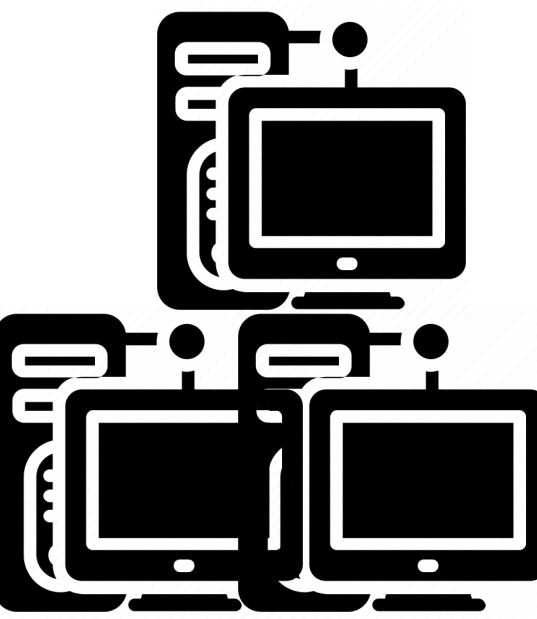
- good reads



array
of discs

- poor ingestion perf.

- poor query perf.



commodity
hardware

SSD wear-friendly

competitive rand. reads

fast ingestion



LSM-tree
O'Neil et al.

1980s

1996

2006

a decade



Bigtable

LSM-tree
O'Neil *et al.*

1996



Bigtable

2006



2007

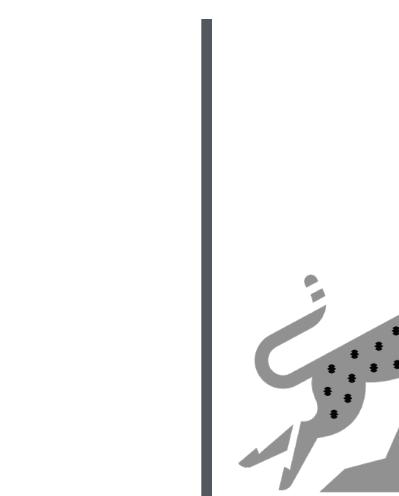


cassandra

2010



levelDB



RocksDB

2011

2013

LSM-tree
O'Neil *et al.*

1996



Bigtable

2006 2007

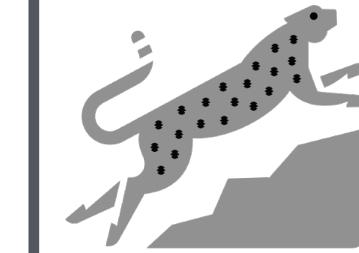


cassandra

2010 2011 2013



levelDB



RocksDB

2023

LSM-tree
O'Neil *et al.*

1996



Bigtable

2006 2007

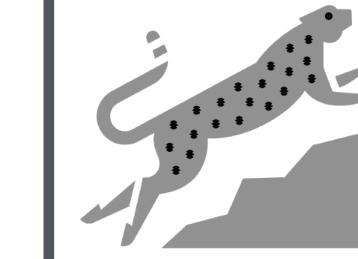


cassandra

2010 2011 2013



levelDB



RocksDB

2023

LSM-tree

NoSQL



relational



time-series

2023

Why **LSM** ?

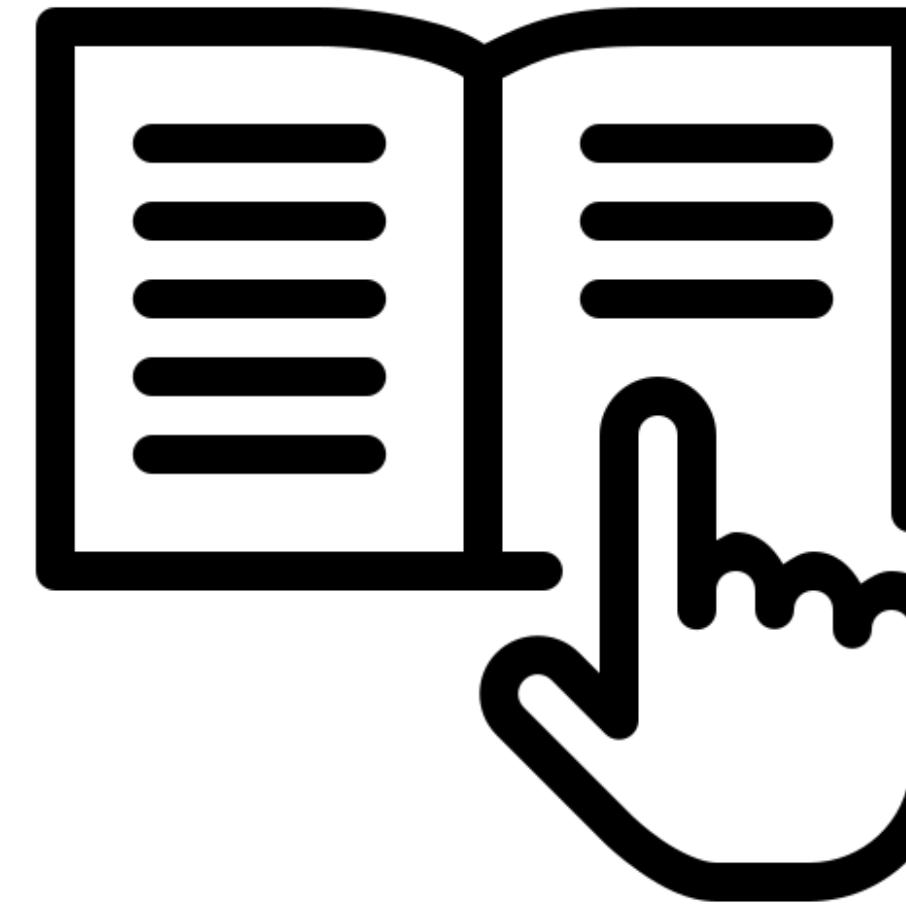


fast ingestion

Why **LSM** ?

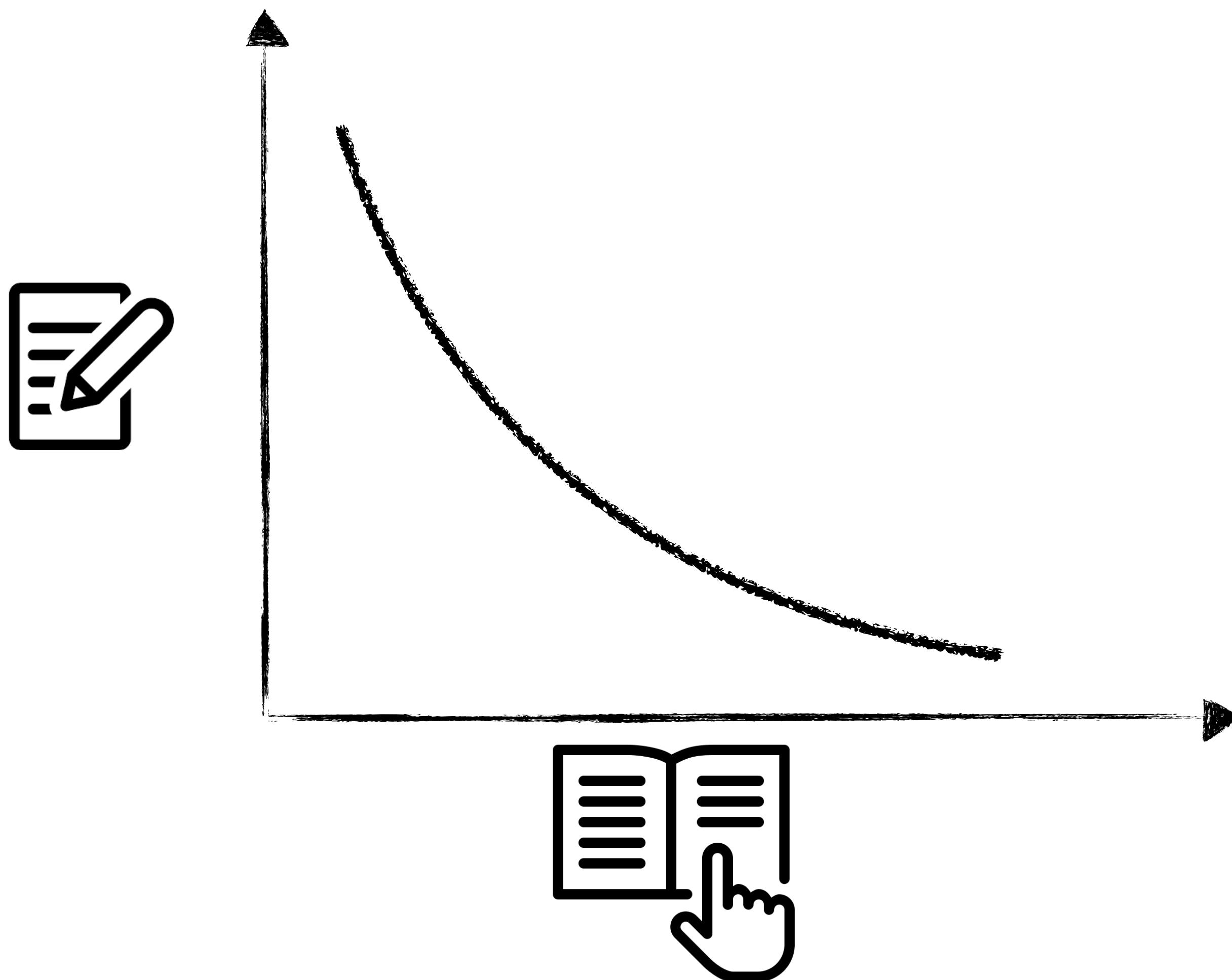


fast ingestion

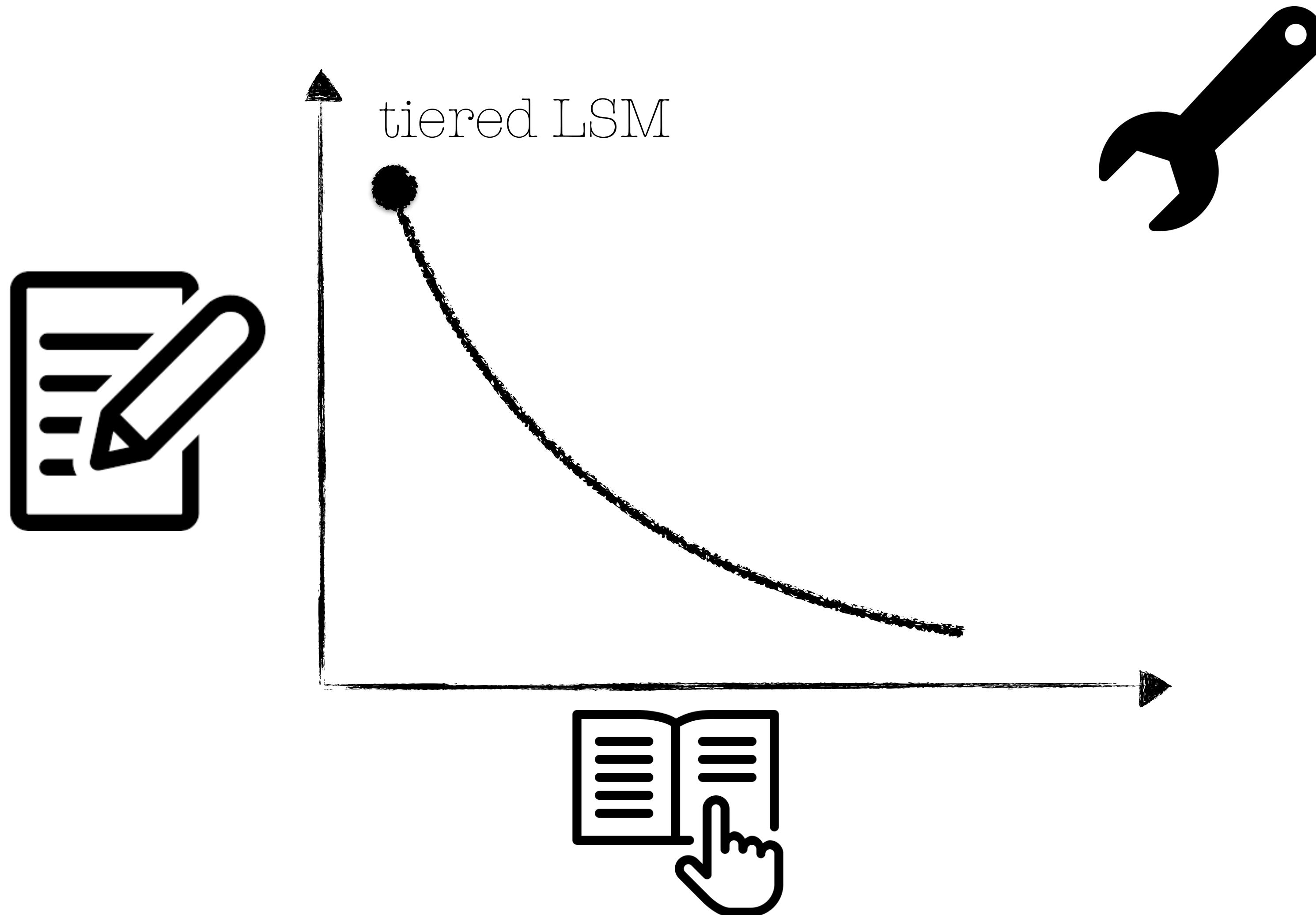


competitive reads

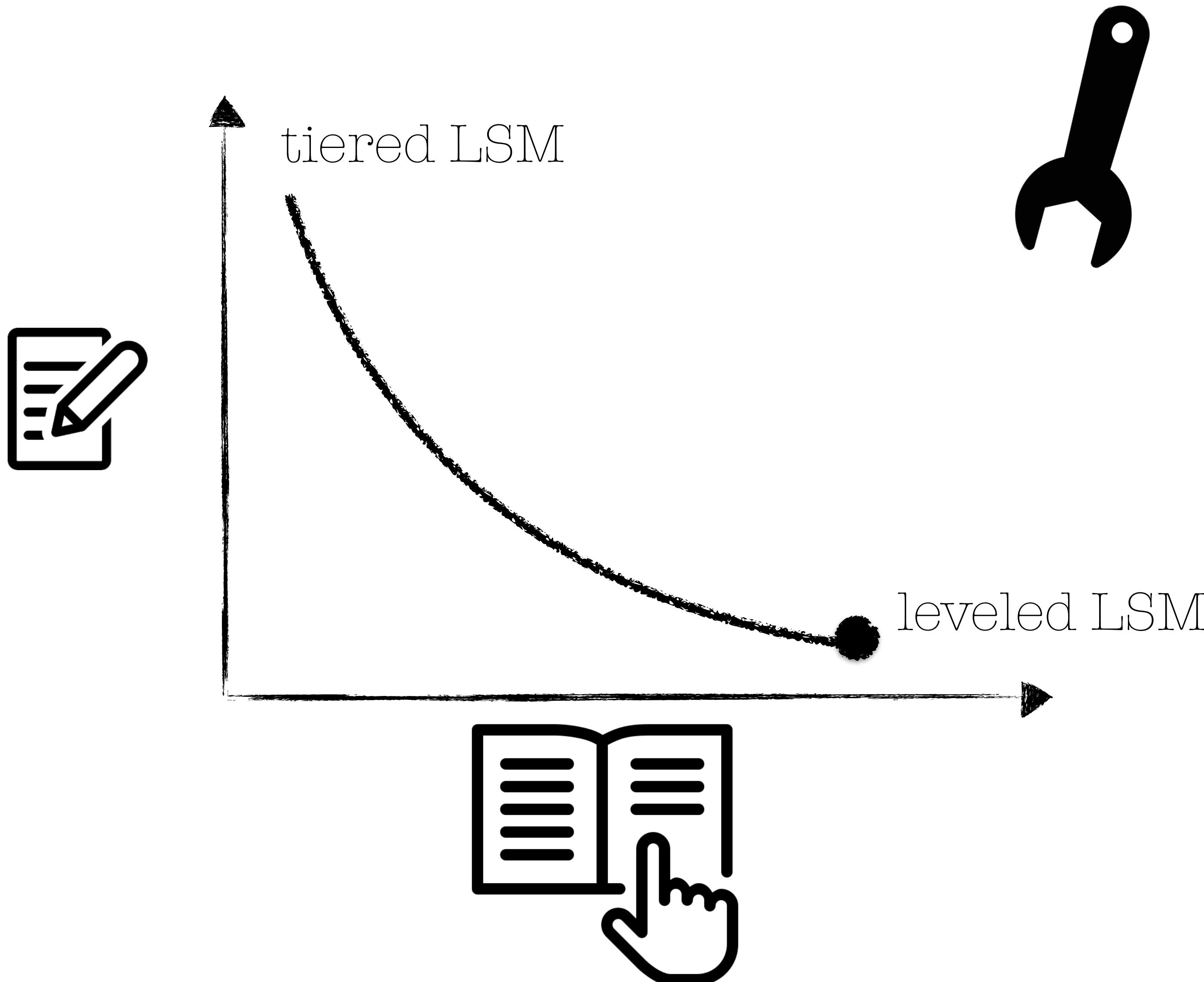
Why **LSM** ?



Why **LSM** ?



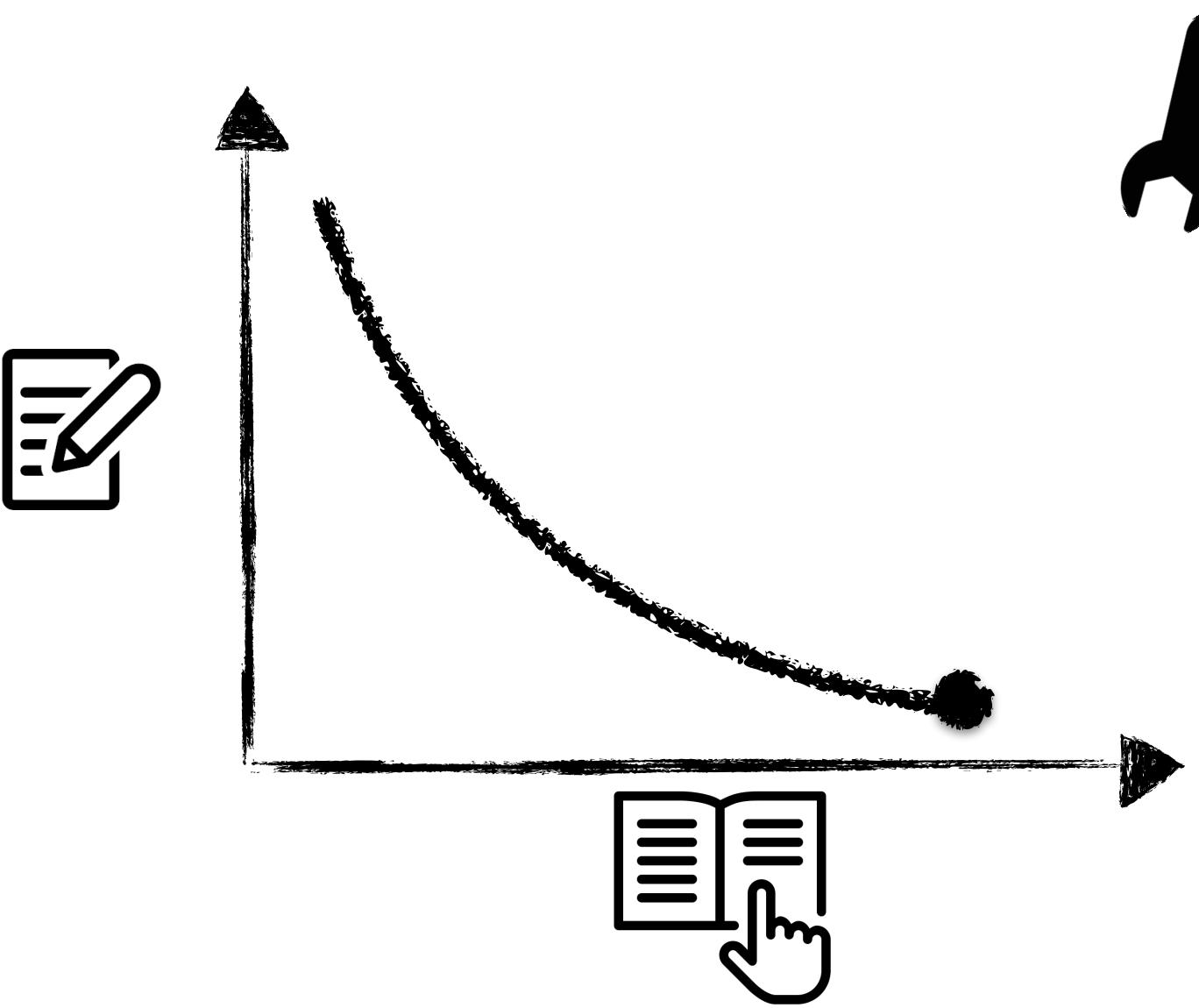
Why **LSM** ?



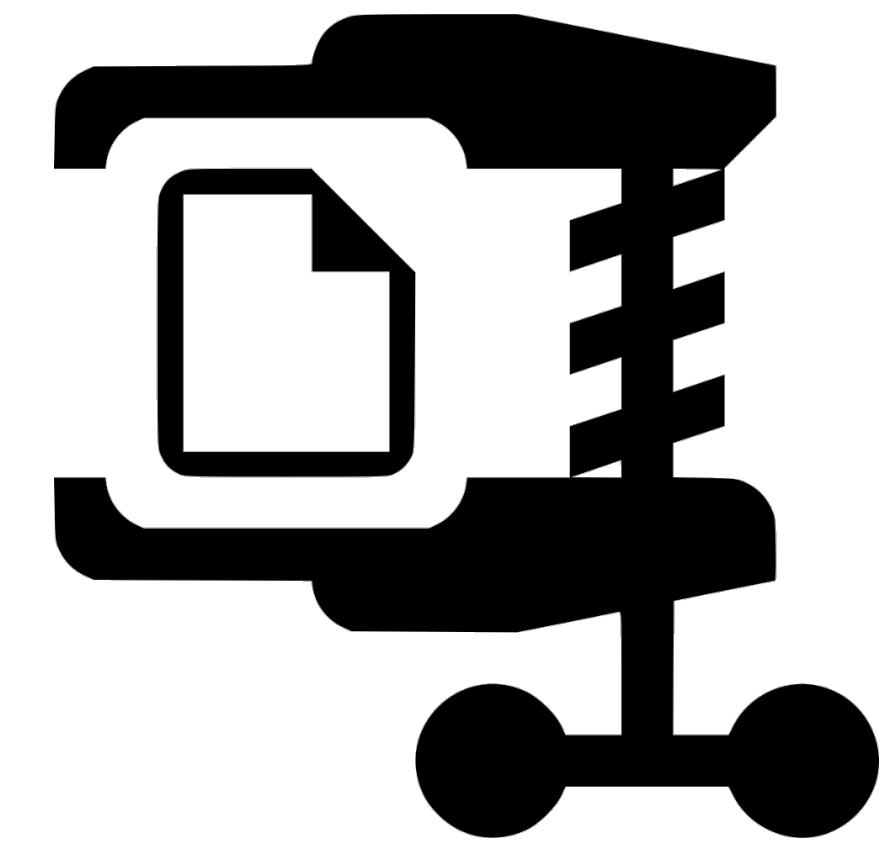
Why LSM ?



fast writes

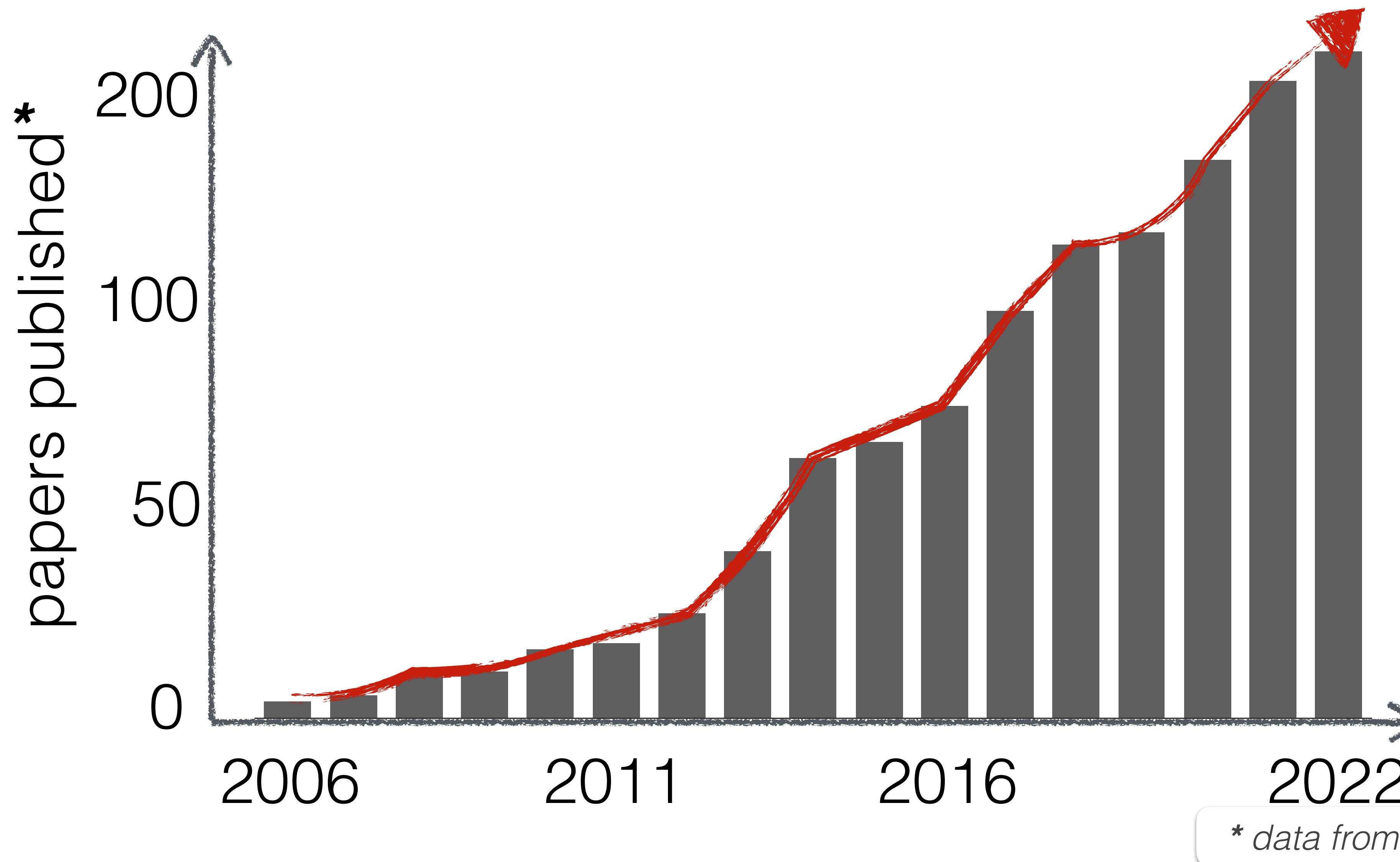


tunable read-write
performance



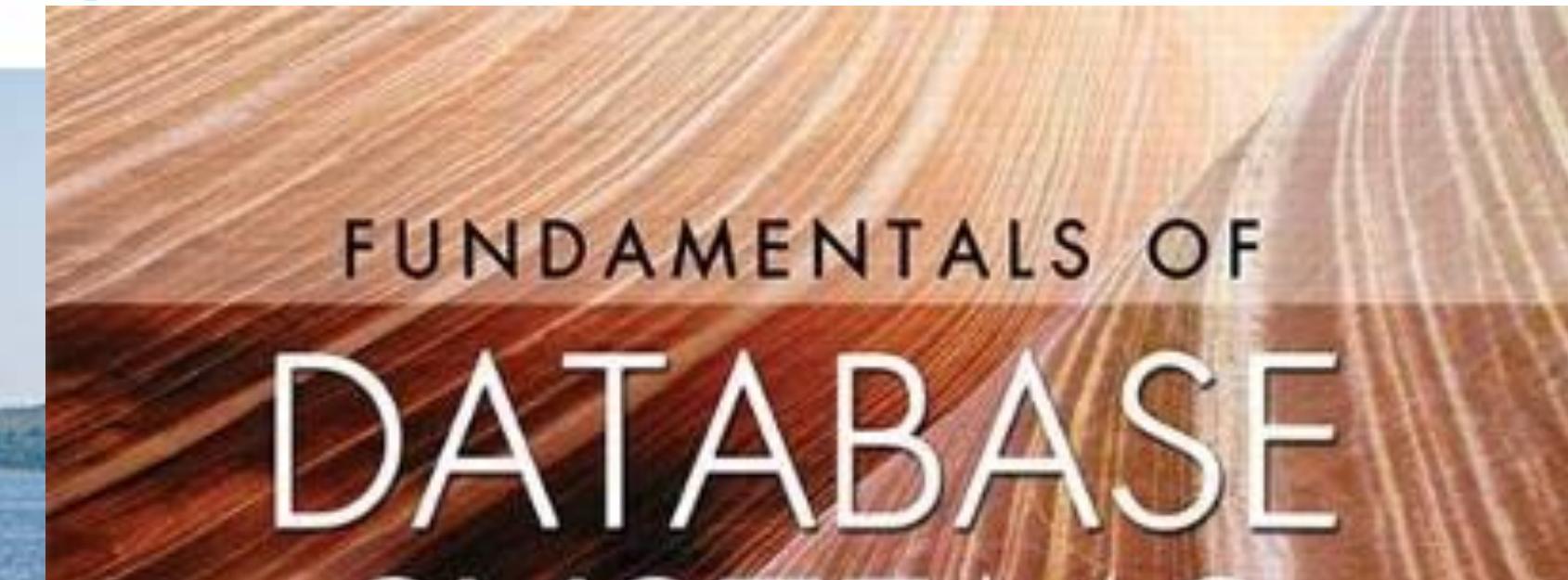
good space
utilization

Research Trend

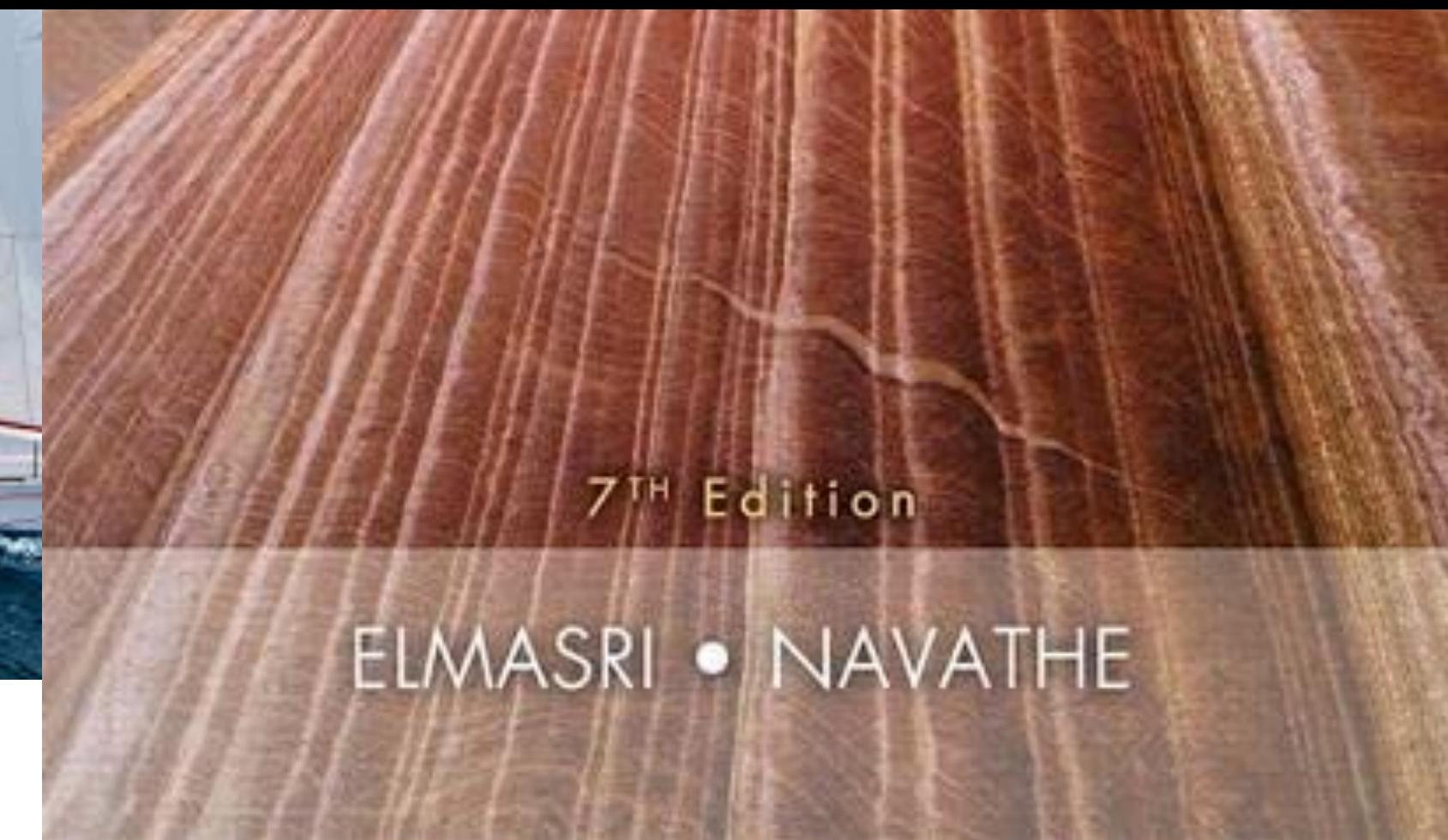
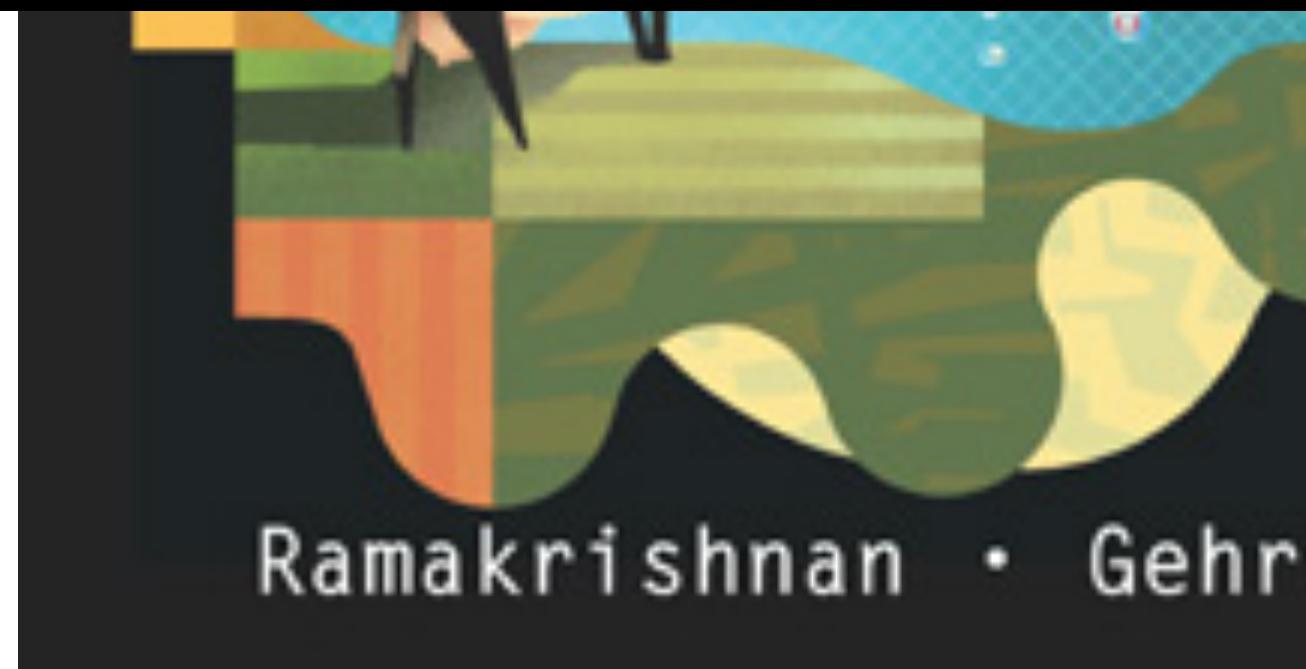




SEVENTH EDITION
Database System Concepts



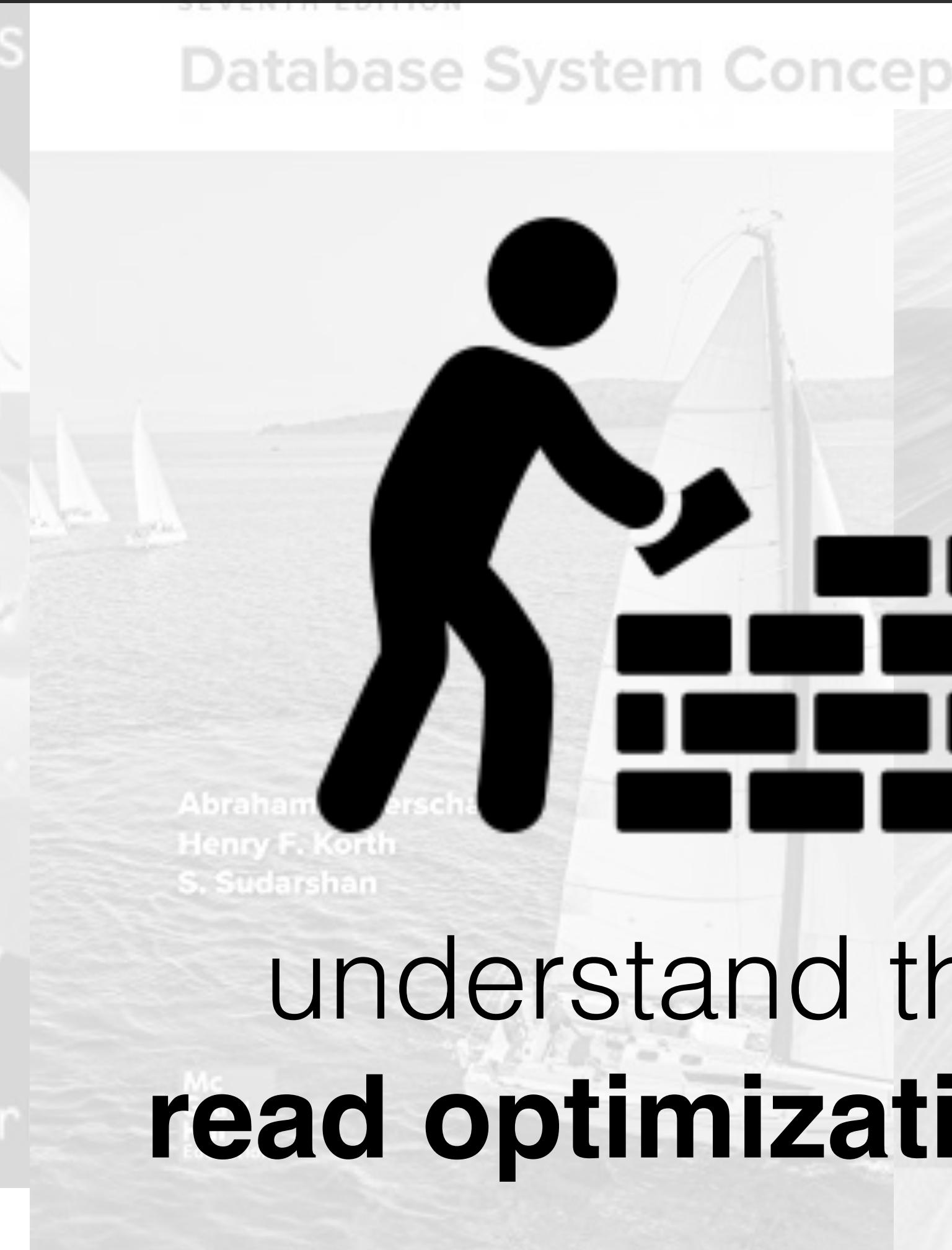
No Textbook on LSMs !!



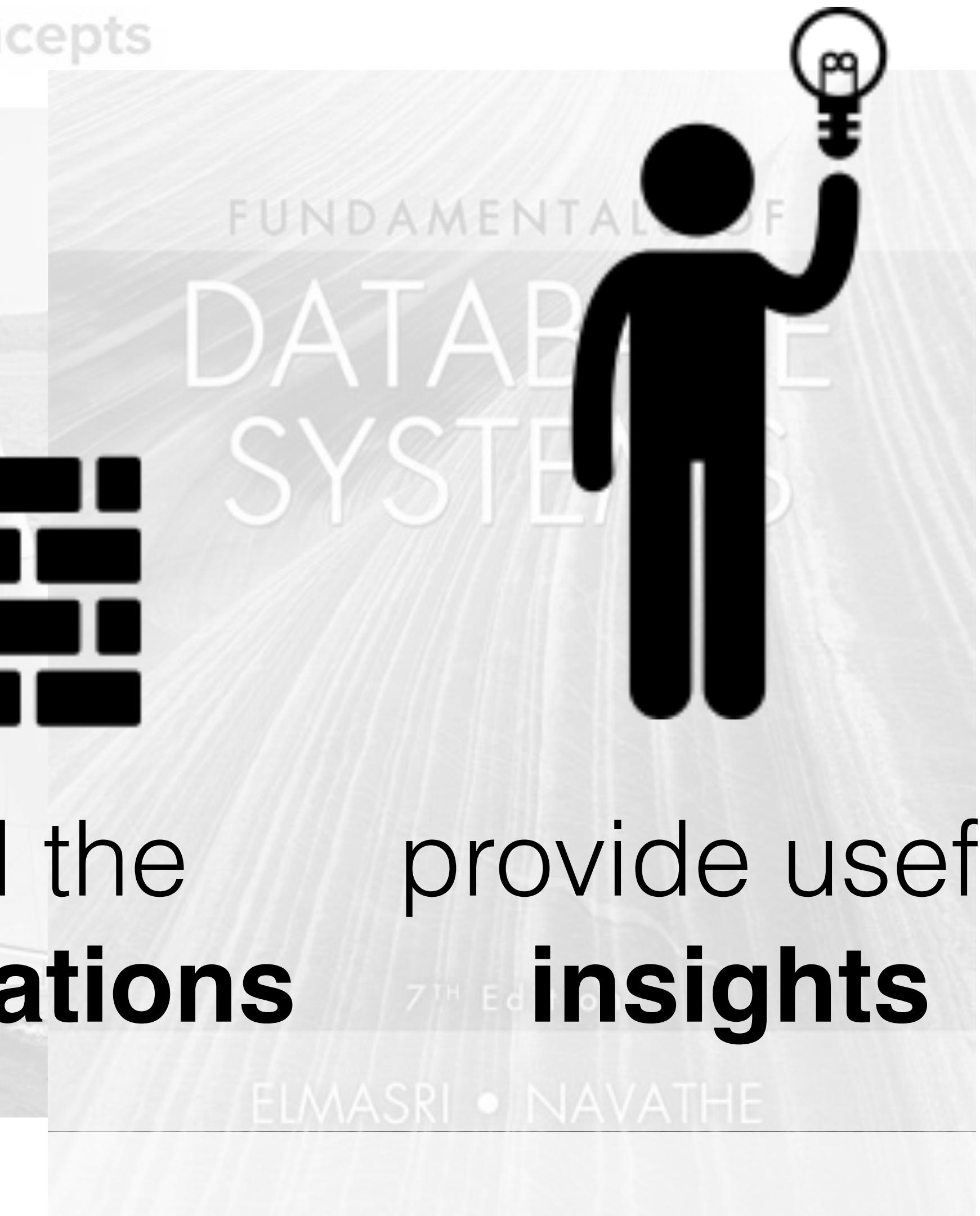
No Textbook on LSMs !!



explore the
LSM paradigm



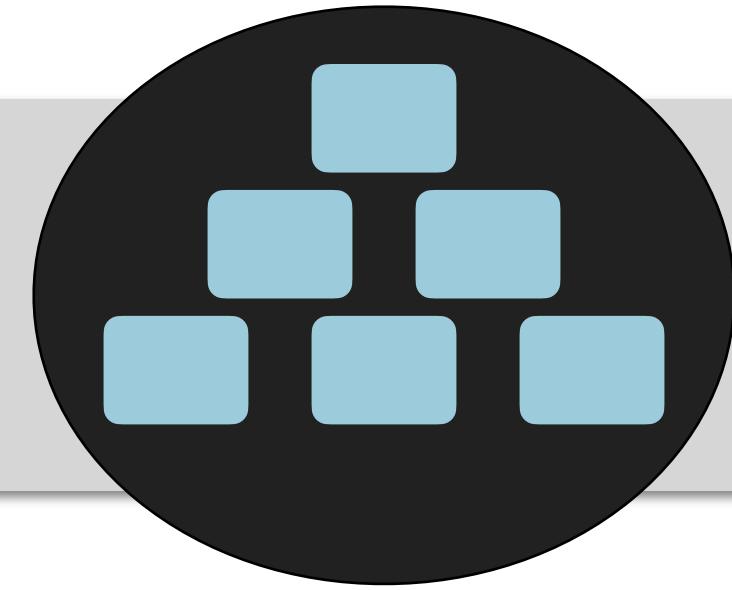
understand the
read optimizations



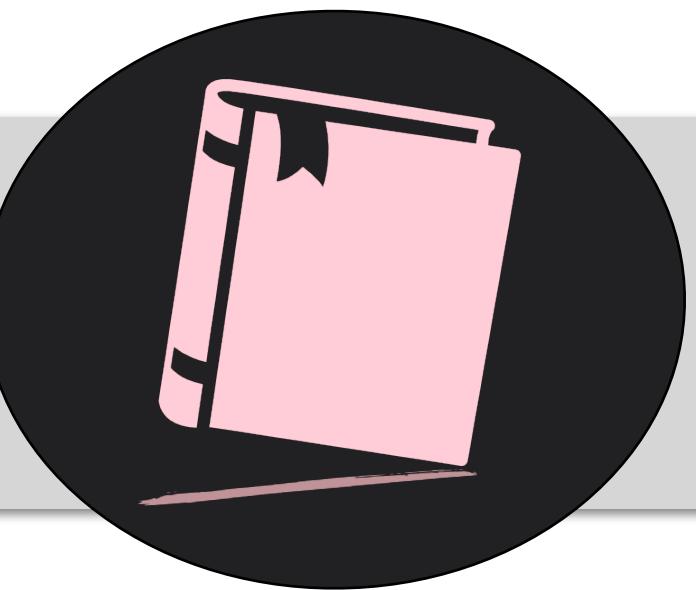
provide useful
insights

Outline

Part 1: **LSM Basics**



Part 2: **Read Optimizations in LSMS**

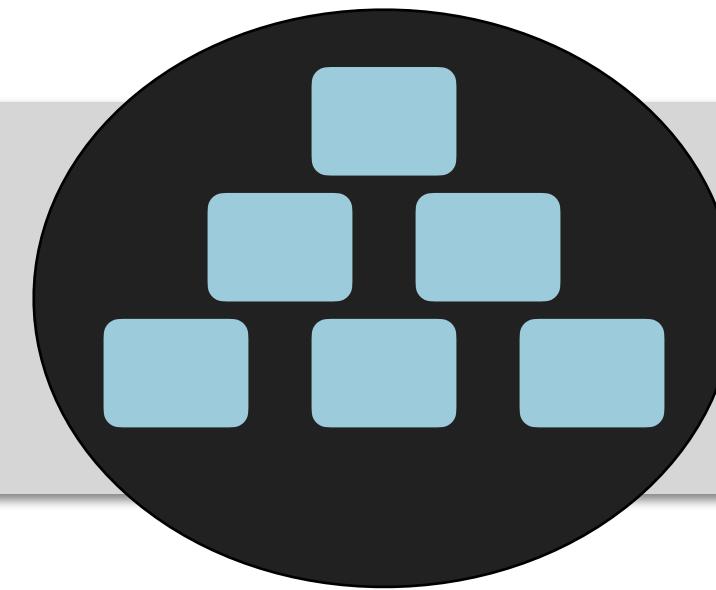


Part 3: **Navigating the LSM Design Space**



Outline

Part 1: **LSM Basics**



Part 2: Read Optimizations in LSMS

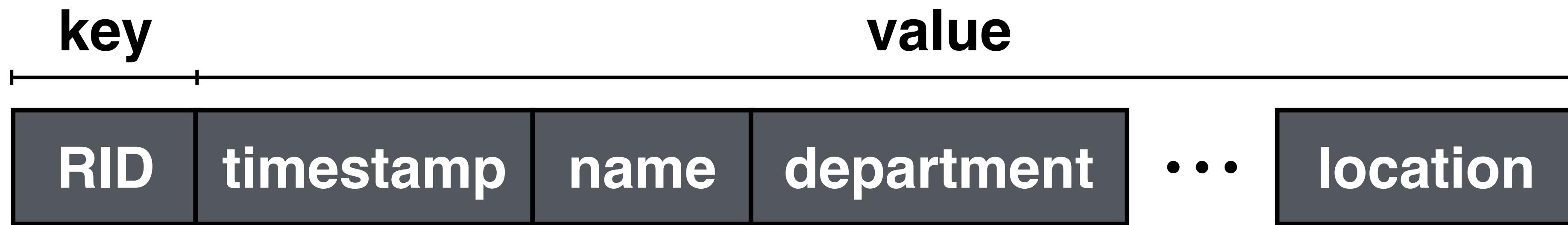


Part 3: Navigating the LSM Design Space



LSM Basics

key-value pairs

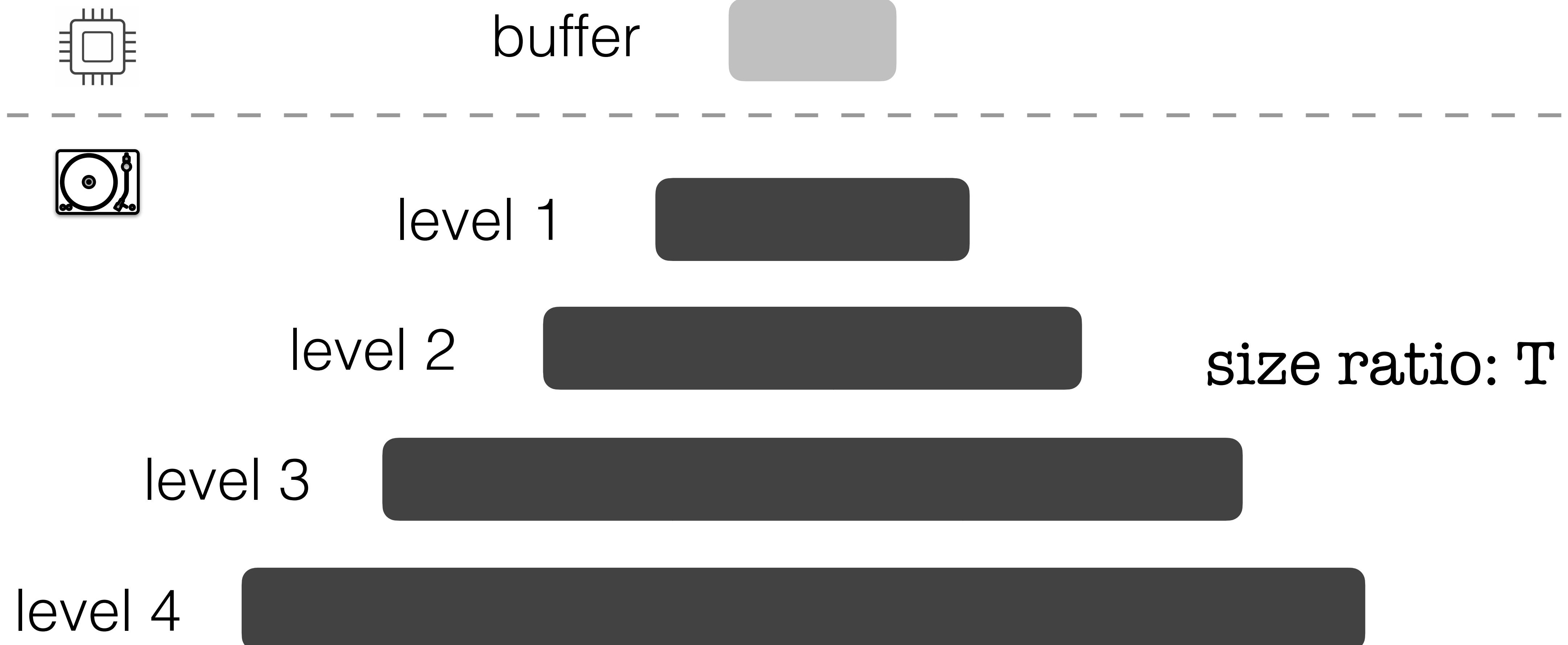


LSM Basics

key-value pairs

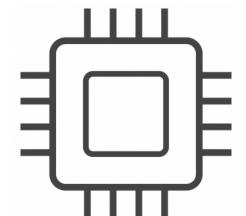
key	value

LSM Basics

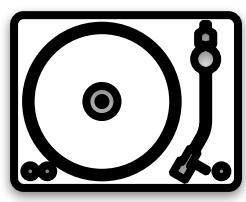
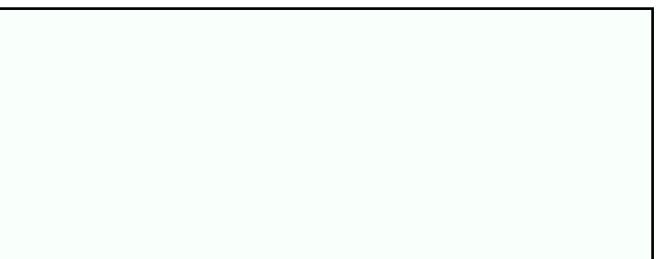


Buffering ingestion

put(6)
put(2)

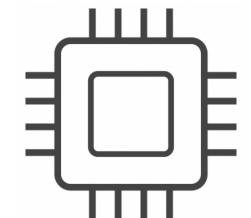


buffer

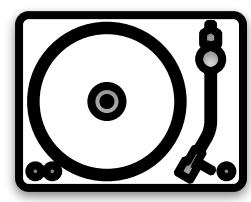
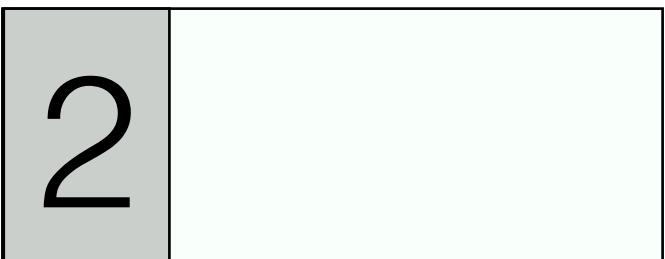


Buffering ingestion

put(1)
put(6)

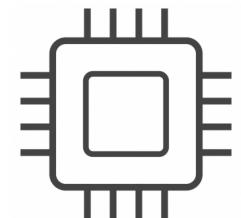


buffer

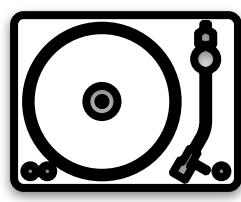
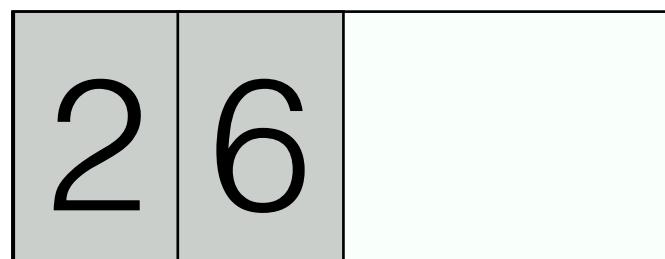


Buffering ingestion

put(4)
put(1)

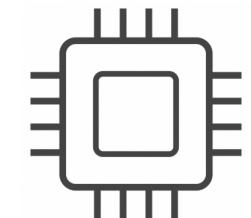


buffer



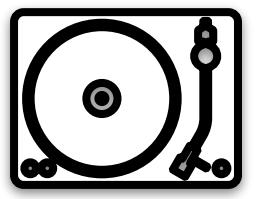
Buffering ingestion

put(4)

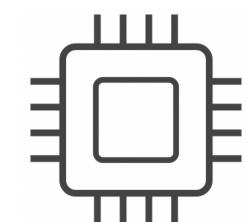


buffer

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

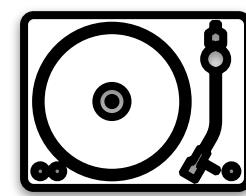


Buffering ingestion

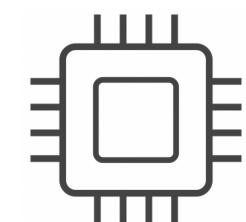


buffer

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

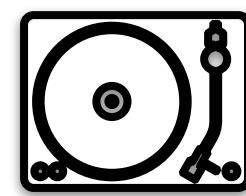


Buffering ingestion

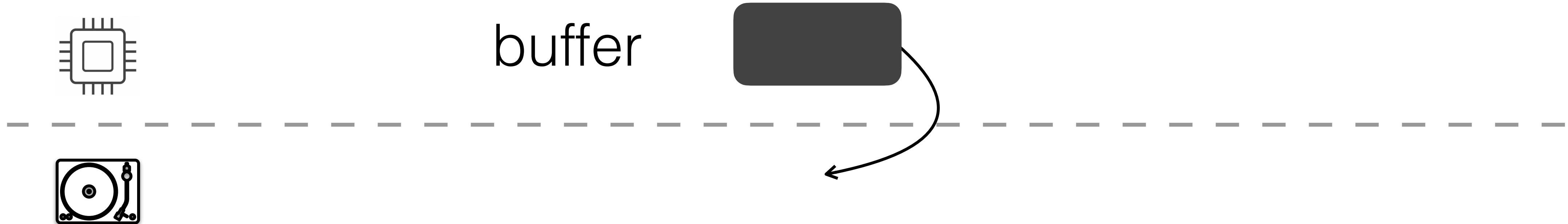


buffer

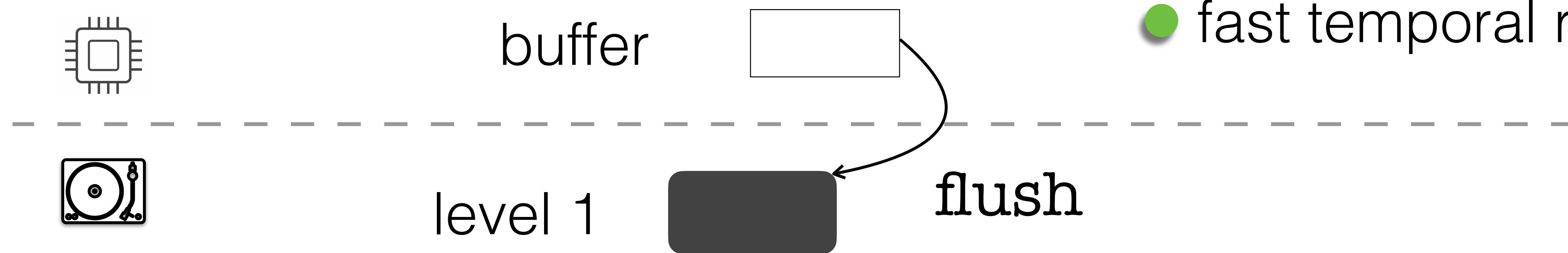
1	2	4	6
---	---	---	---



Buffering ingestion

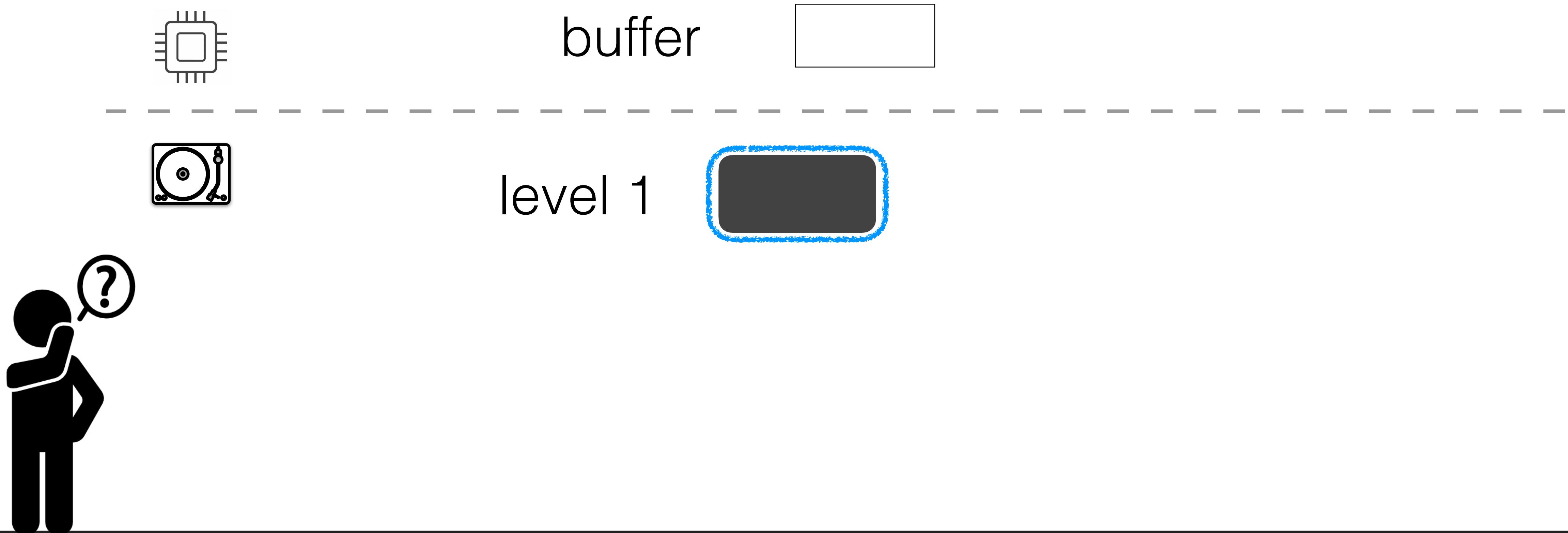


Buffering ingestion



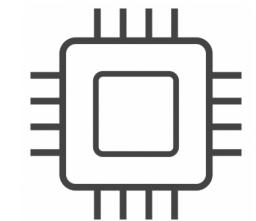
- low ingestion cost
- fast temporal reads

Immutable files on storage

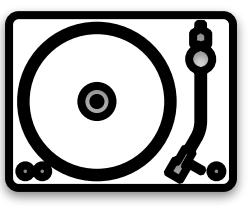
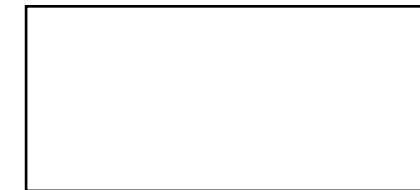


How do we update data?

put(6)

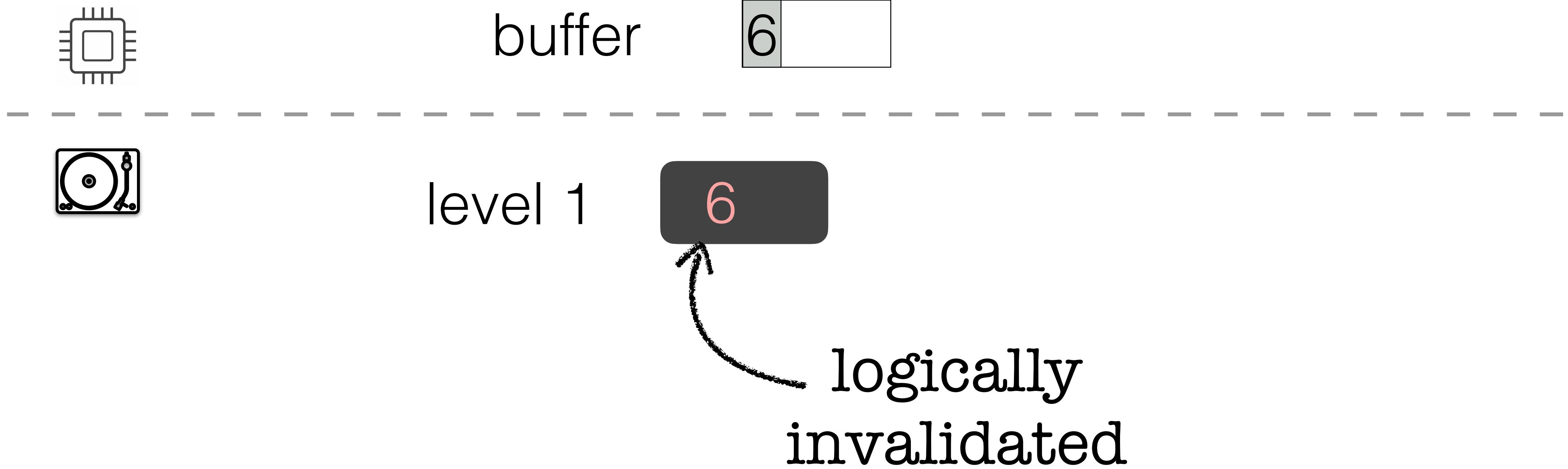


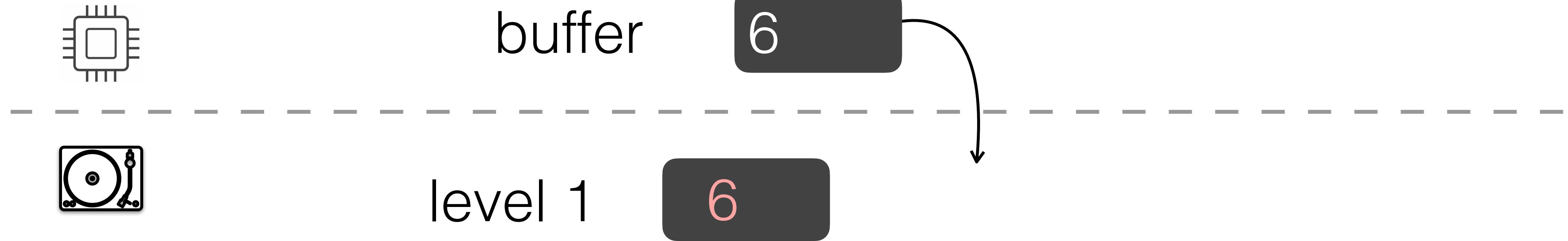
buffer



level 1

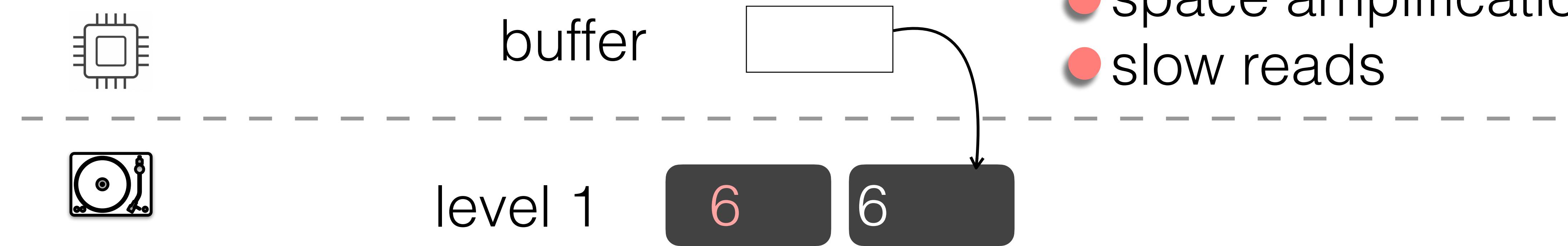
6



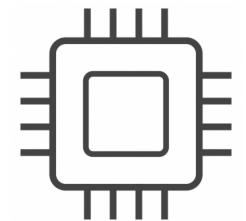


Out-of-place updates

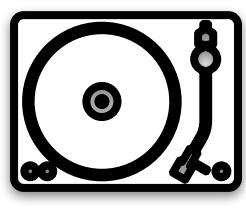
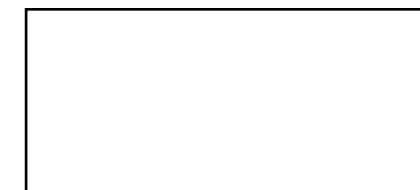
- fast ingestion
- space amplification
- slow reads



How do we reduce this space amplification?



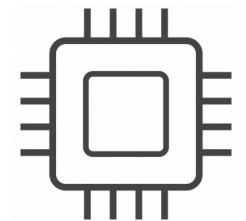
buffer



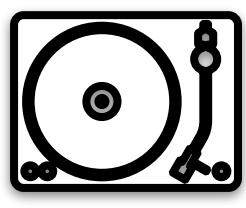
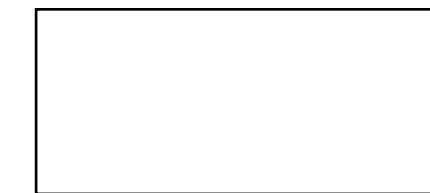
level 1

6

6



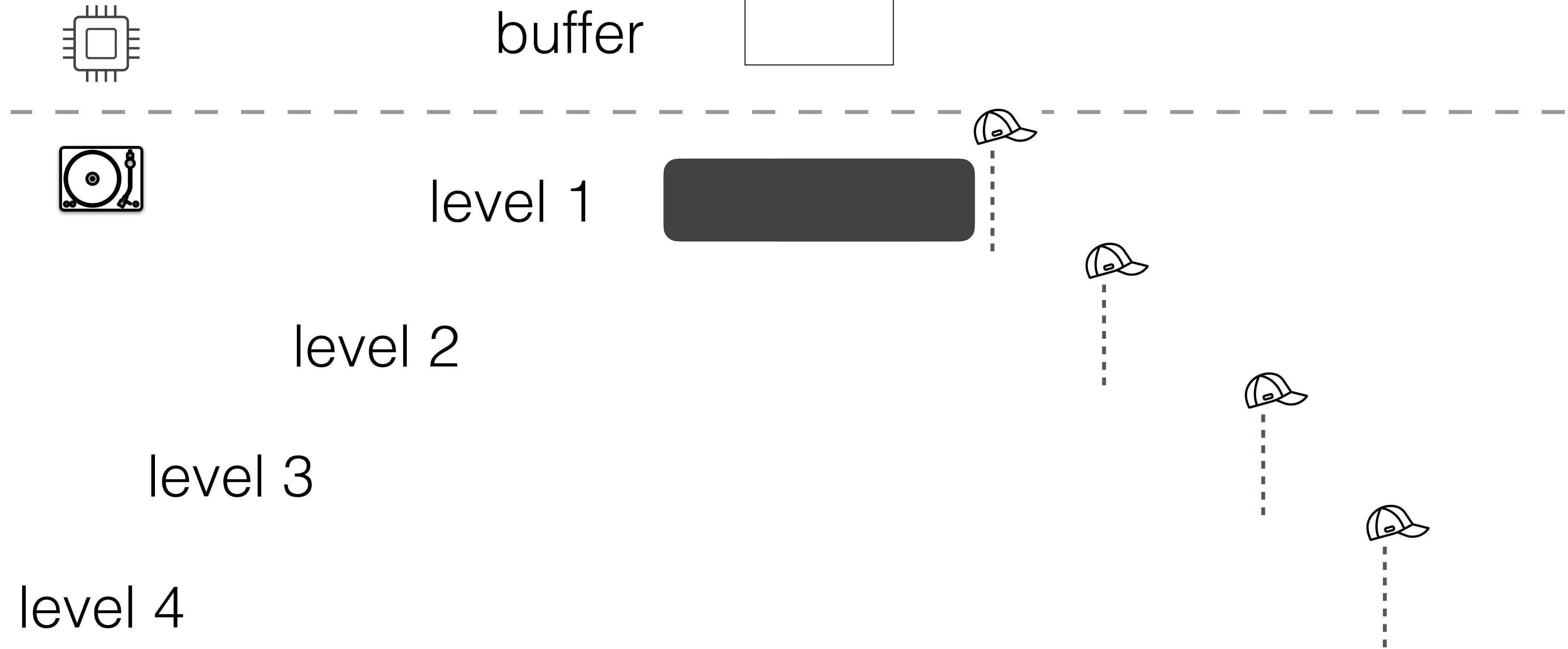
buffer



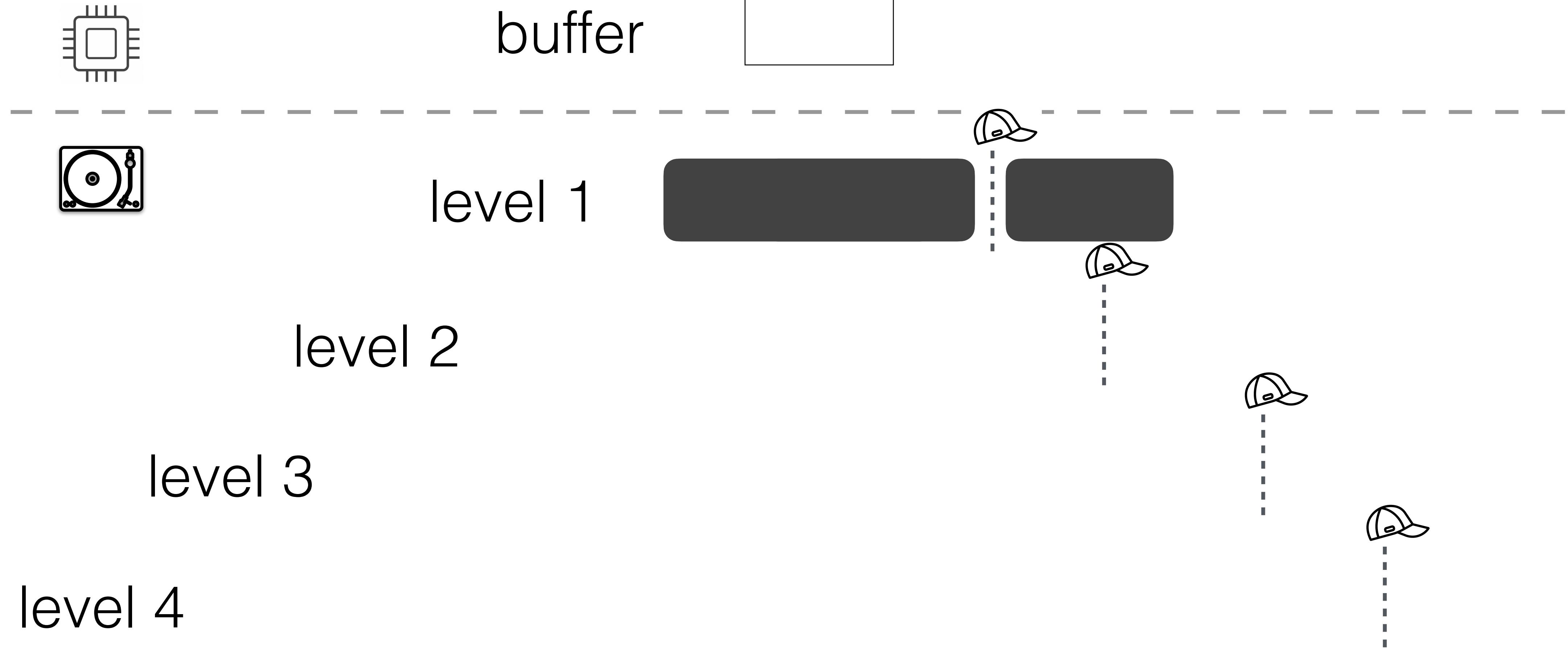
level 1

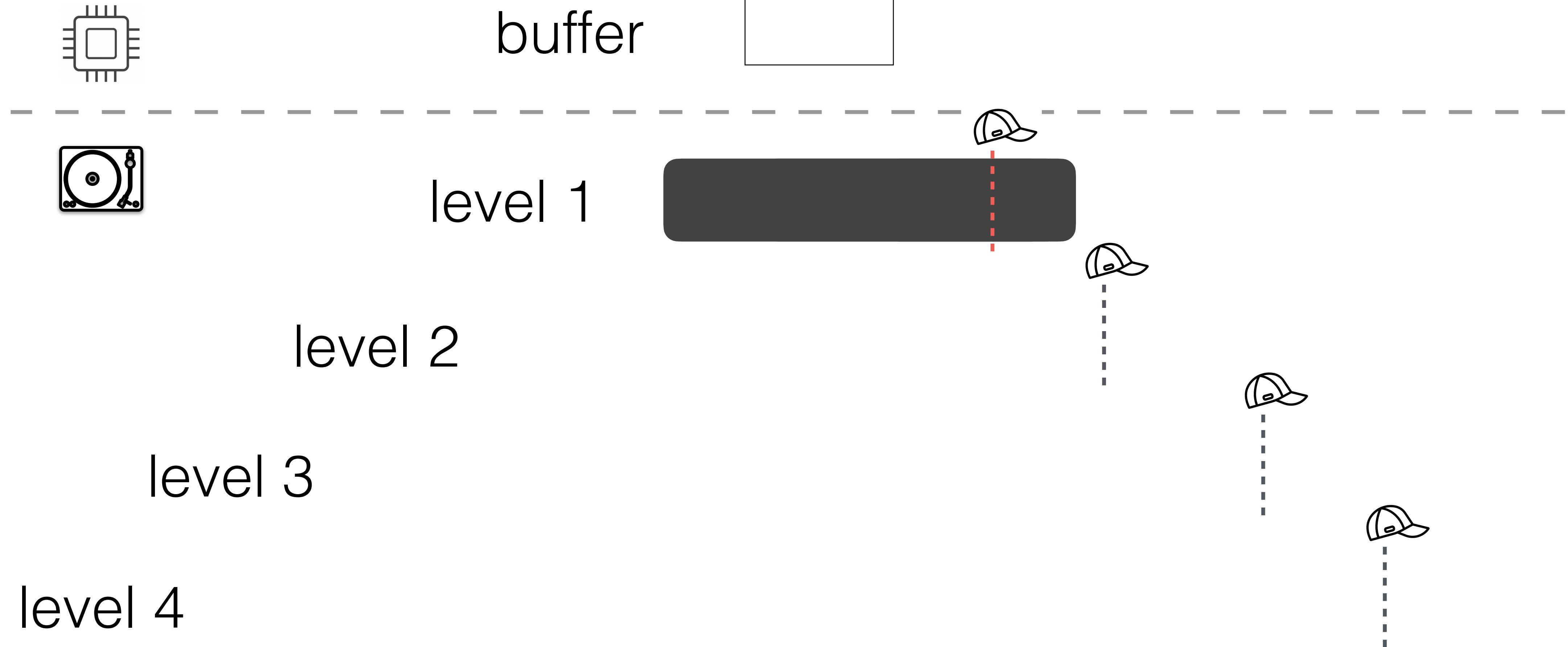
6

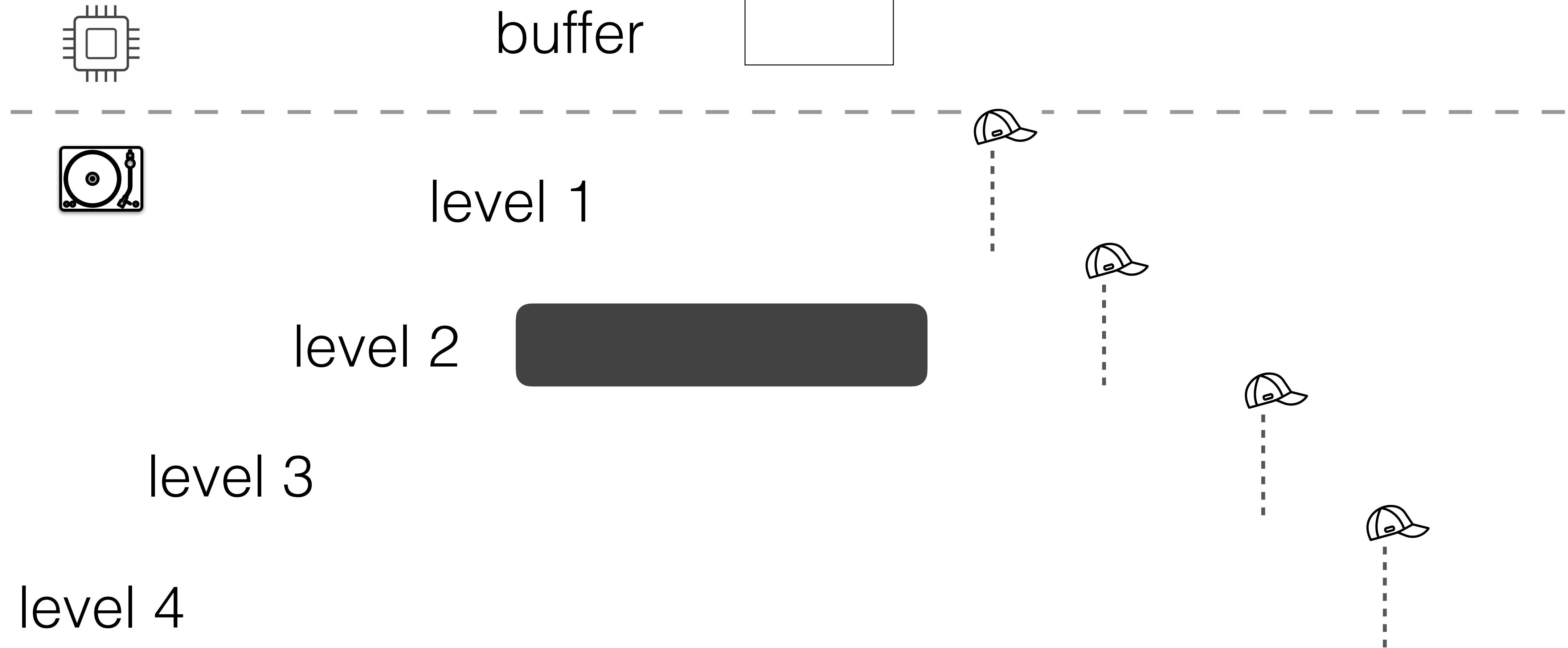
compaction





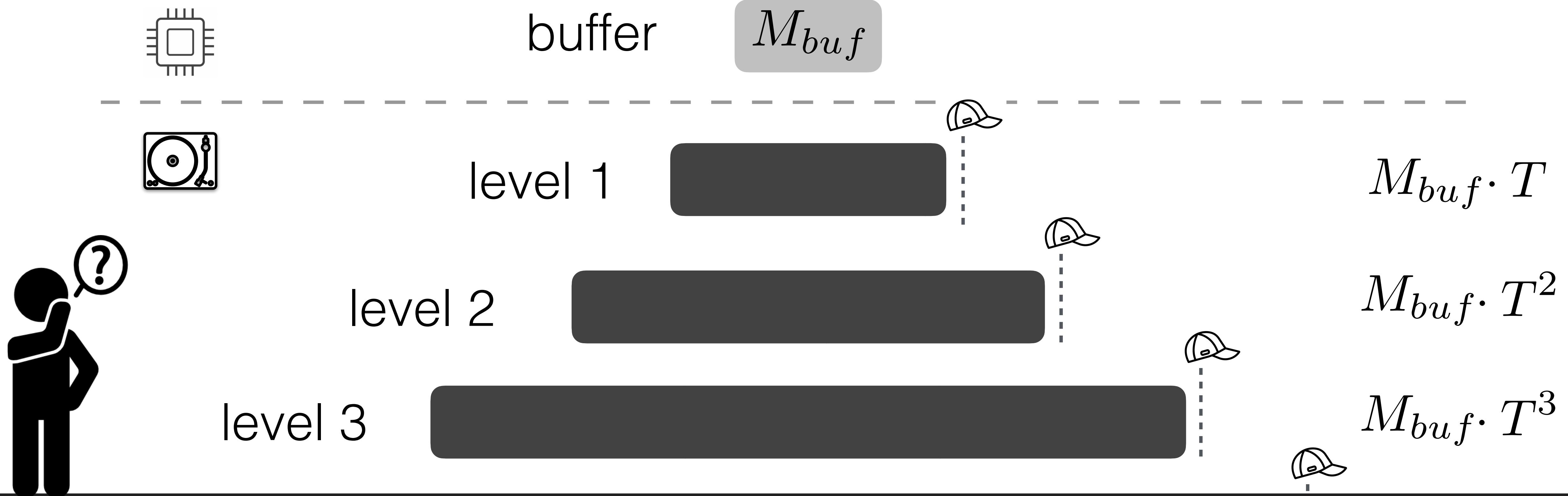






M_{buf} : buffer memory

T : size ratio

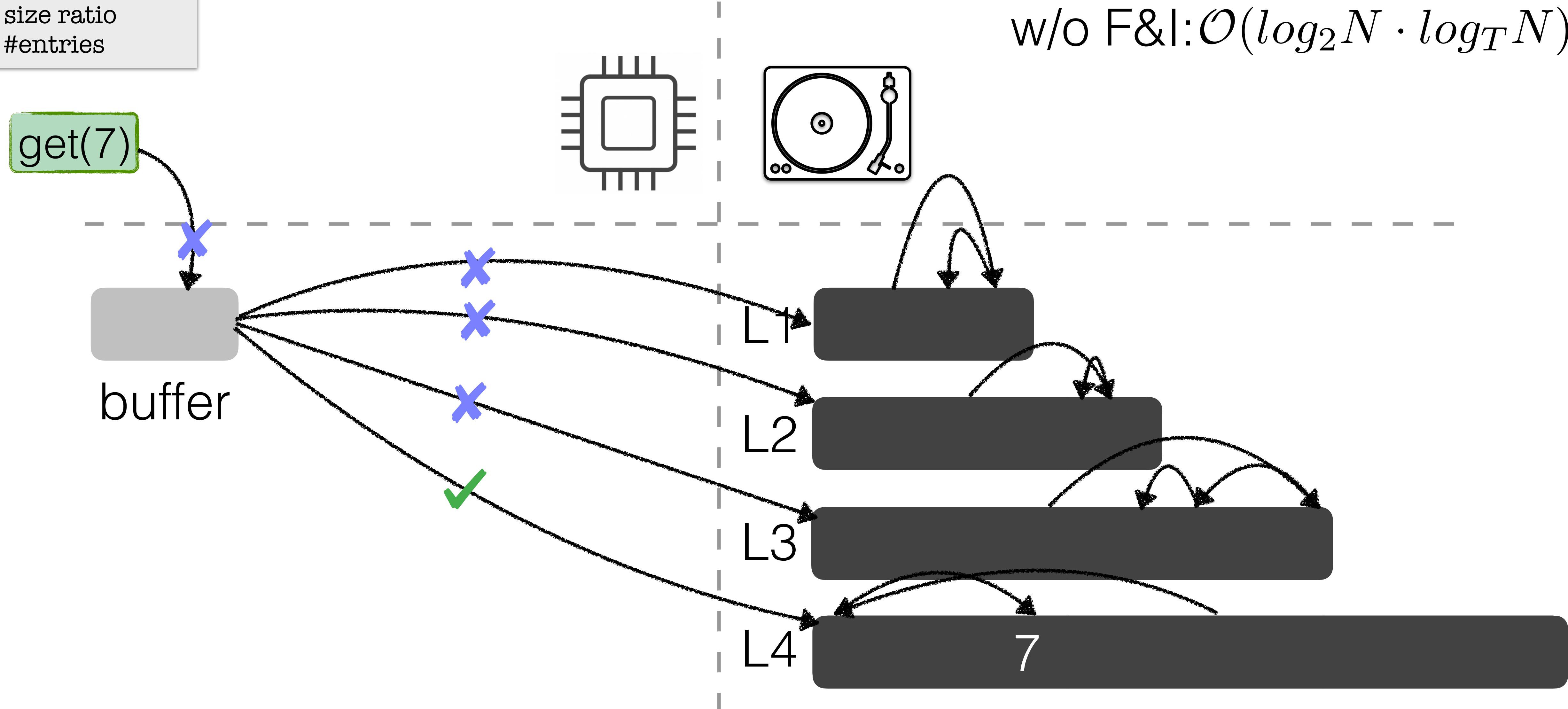


How about queries?

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

Cost analysis

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

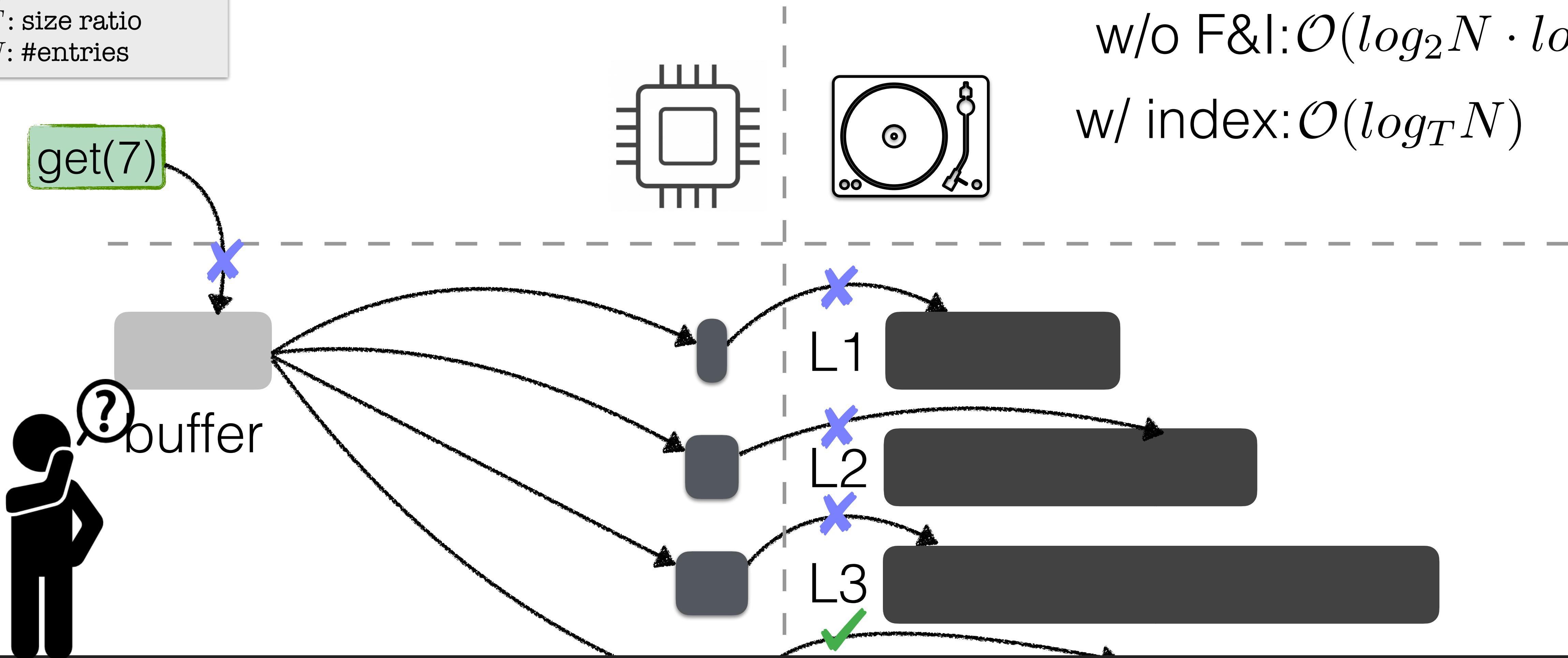


P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

Cost analysis

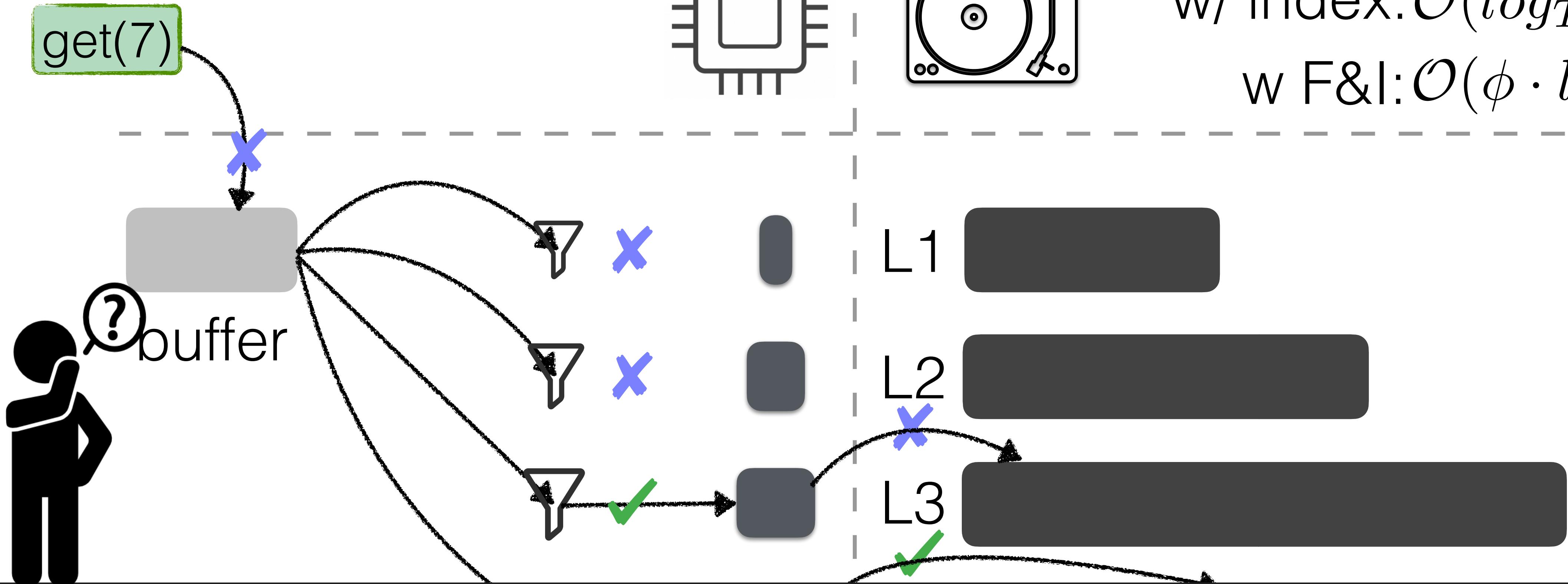
w/o F&I: $\mathcal{O}(\log_2 N \cdot \log_T N)$

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
 ϕ : FPR of BF



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: $\mathcal{O}(\phi \cdot \log_T N)$

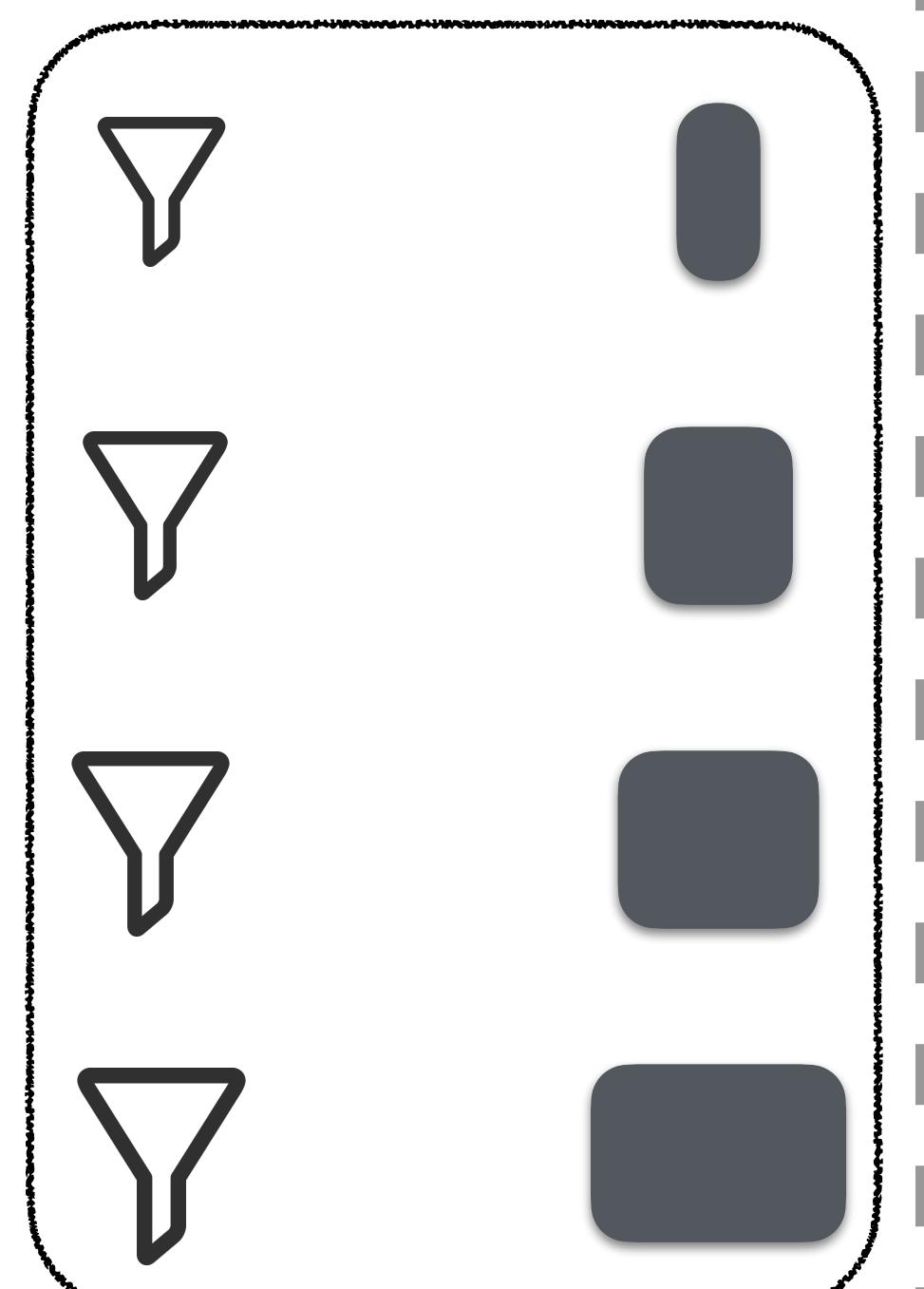
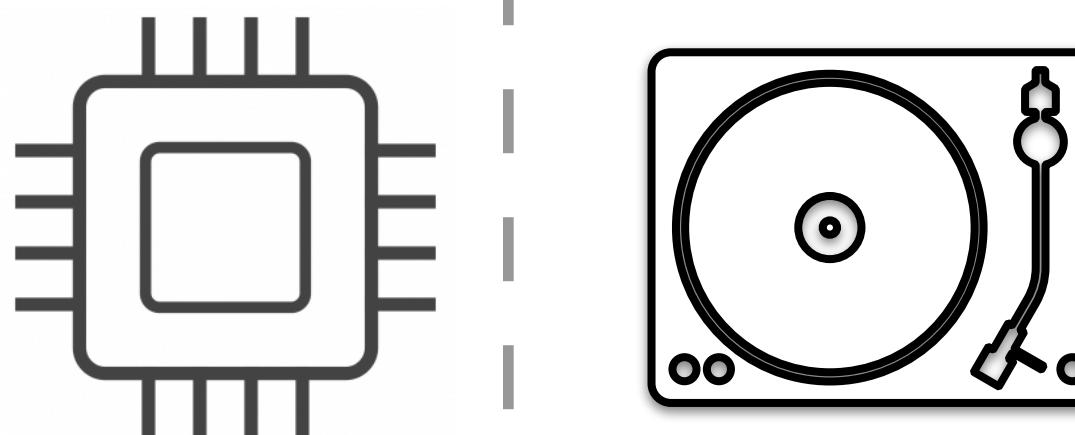
How to manage memory?



buffer



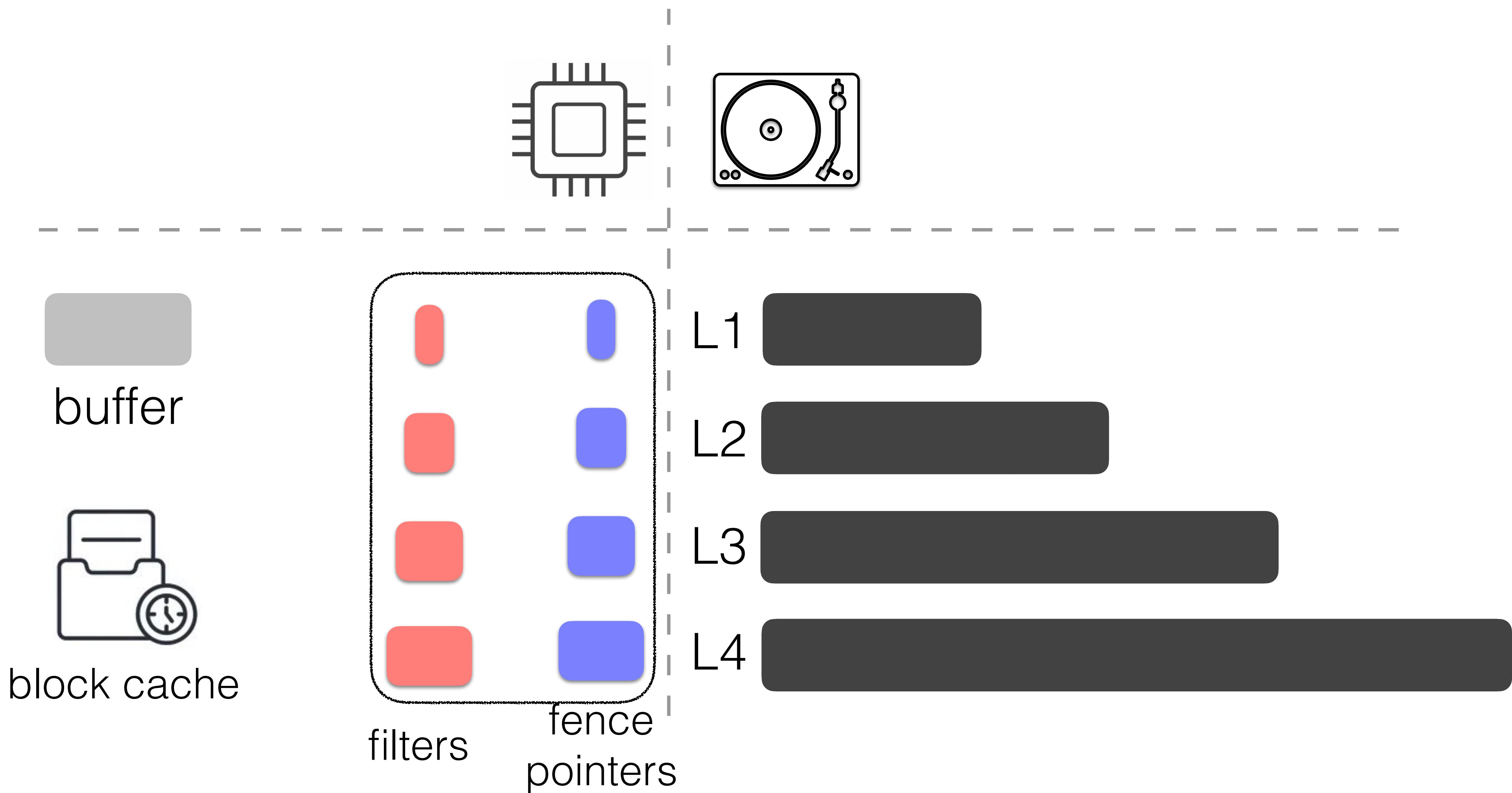
block cache



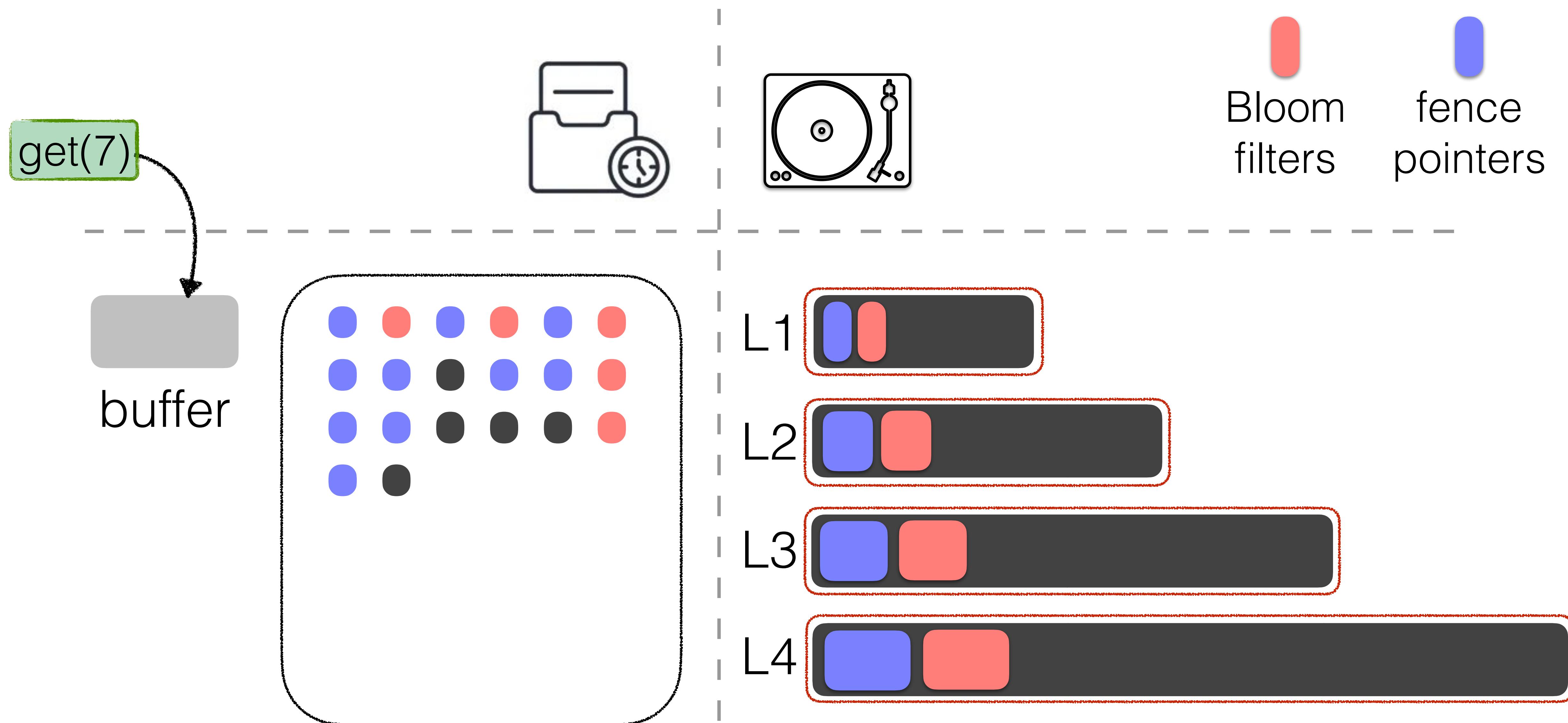
fence
pointers



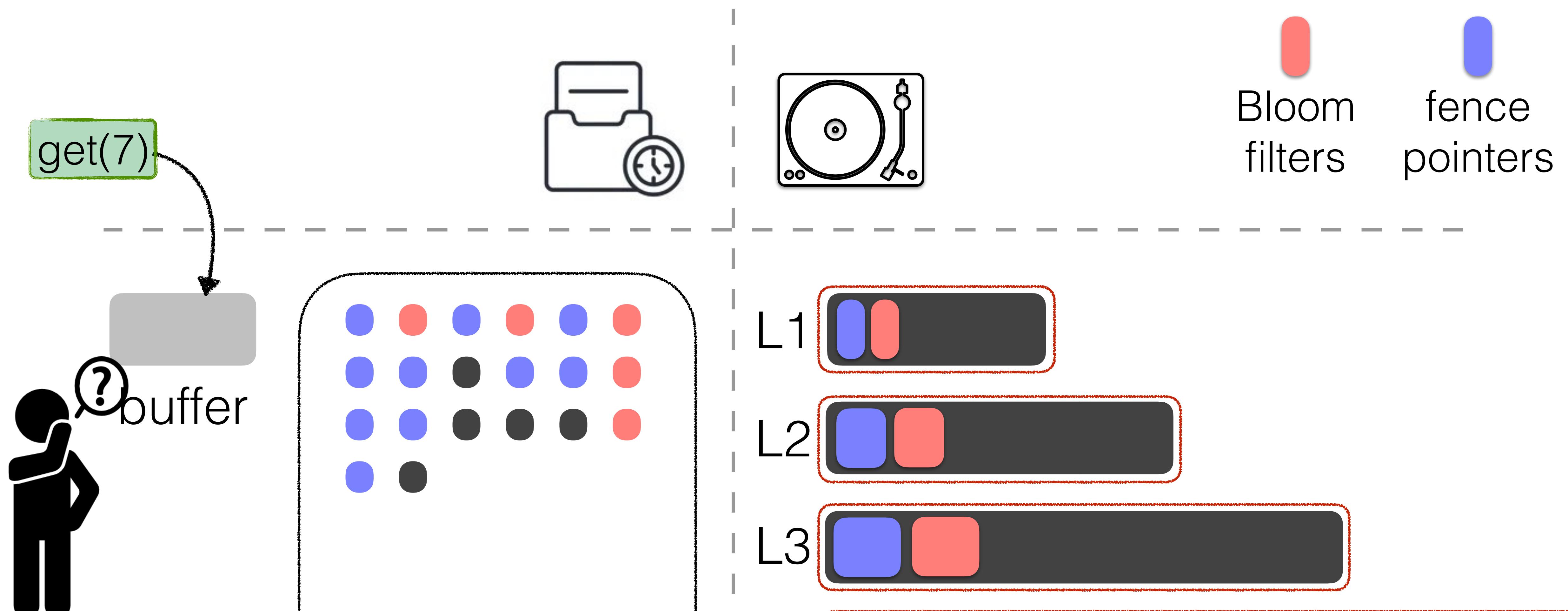
Block Cache



Block Cache



Block Cache

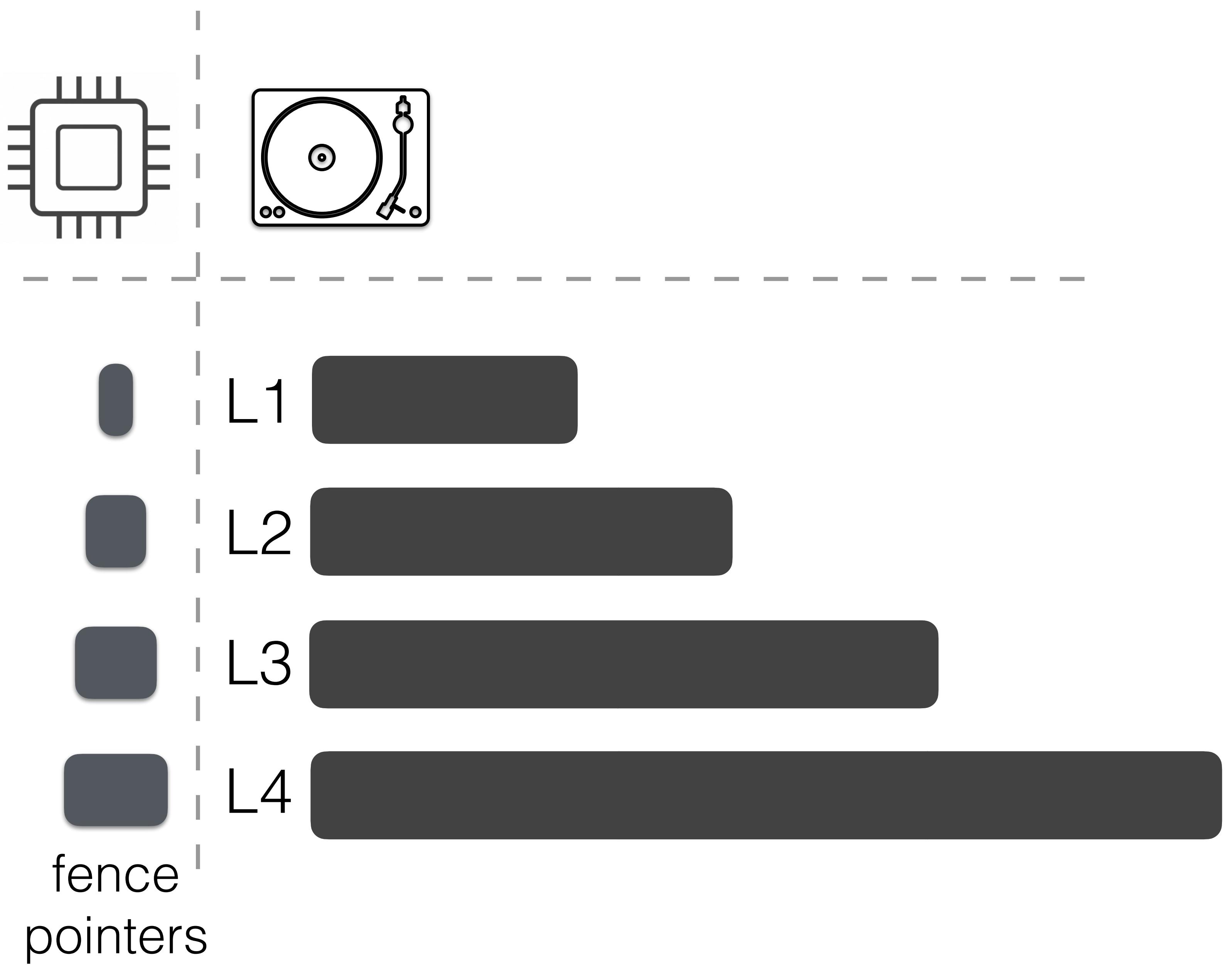


What about range queries?

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

s : selectivity LRQ

Range Queries



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

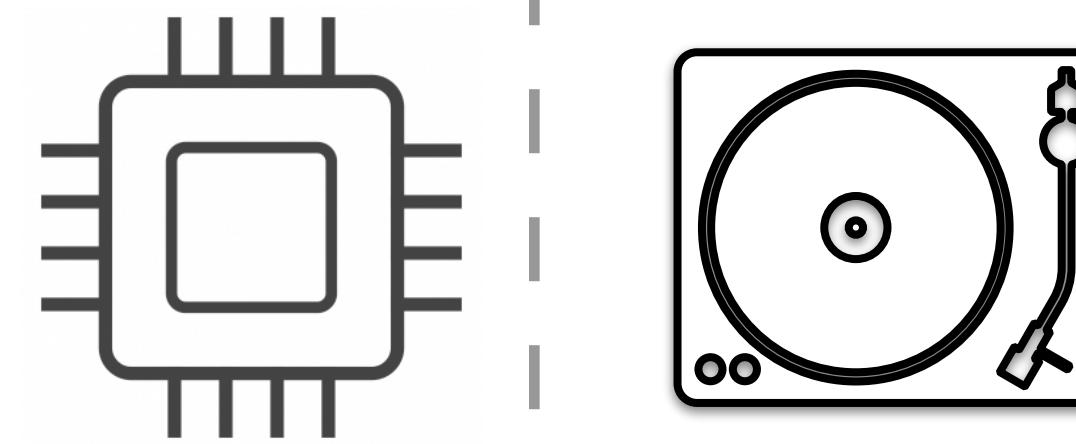
s : selectivity LRQ

Range Queries

Cost analysis

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

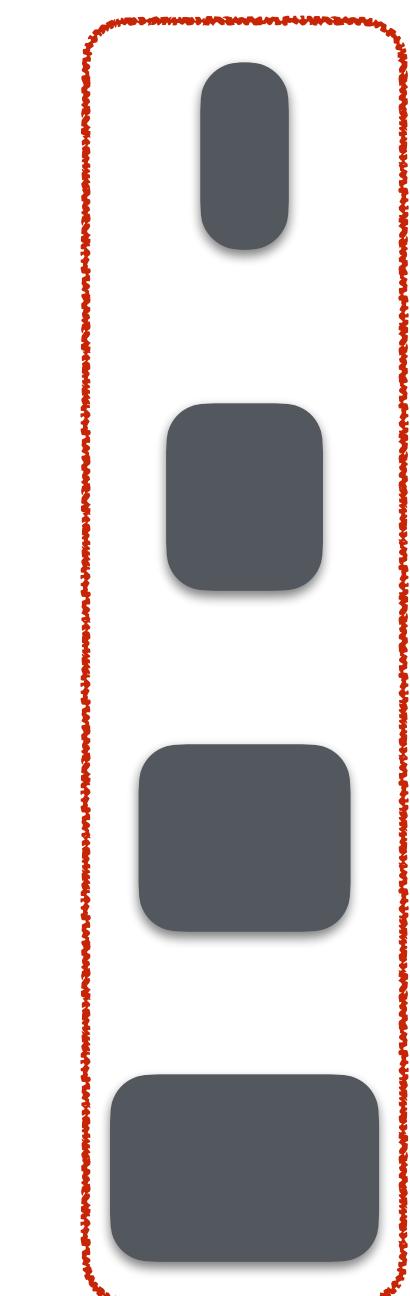
get(9,90)



buffer



filters



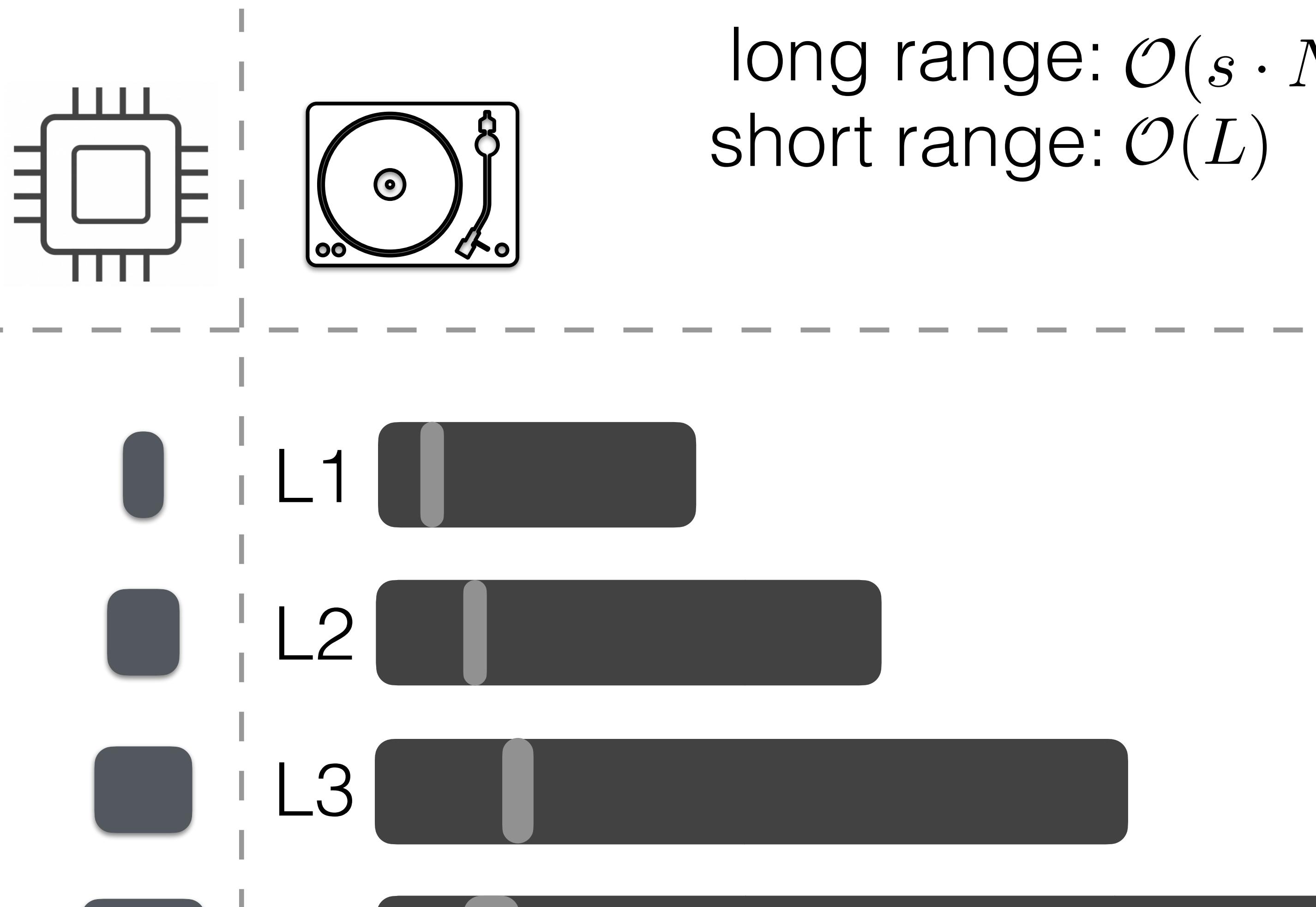
fence
pointers



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

s : selectivity SRQ

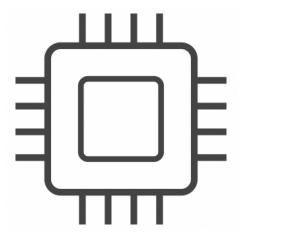
Range Queries



Cost analysis

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

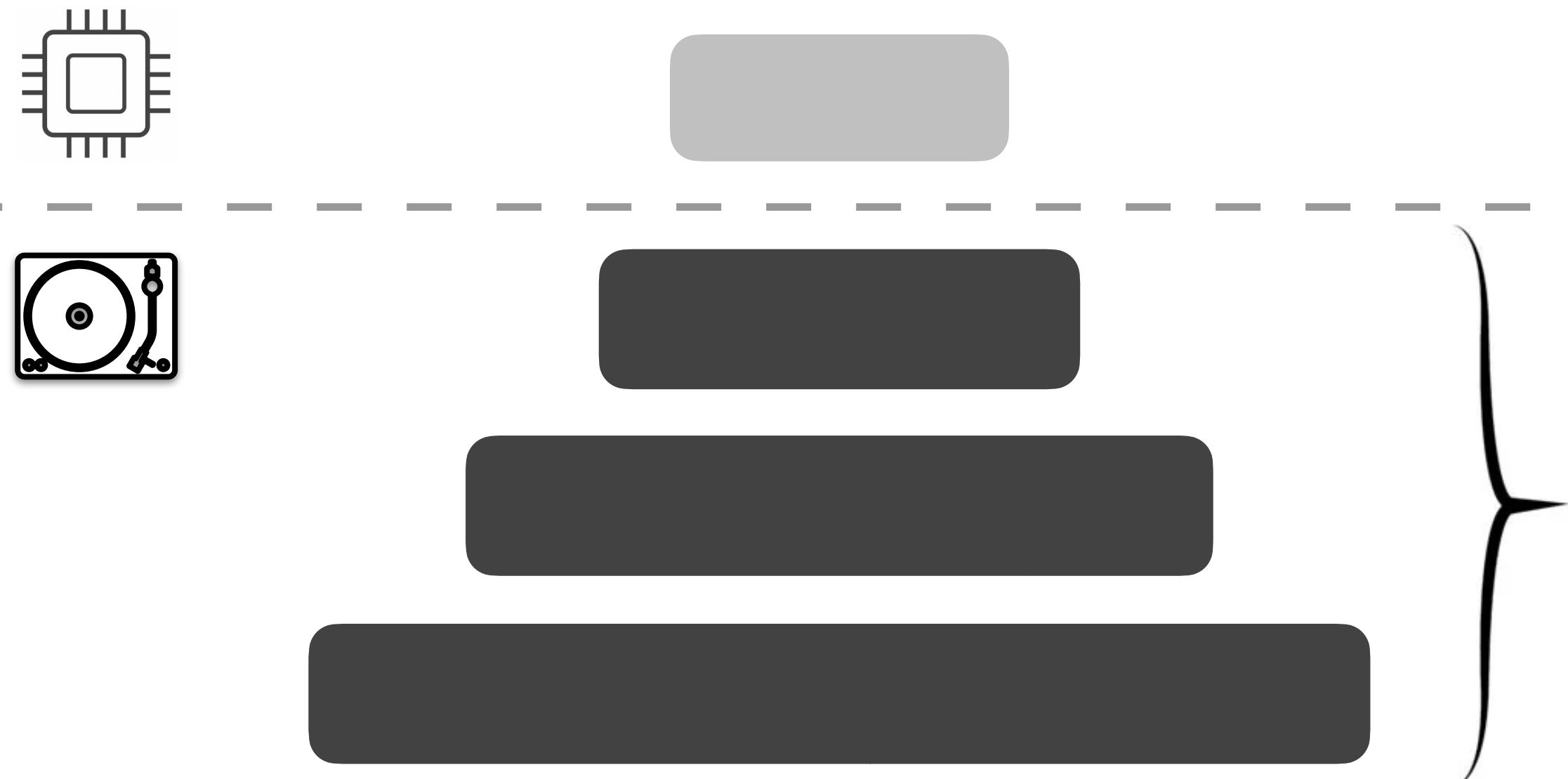
More on LSM Reads in Part 2.



most data
on storage

L : #levels

T : size ratio

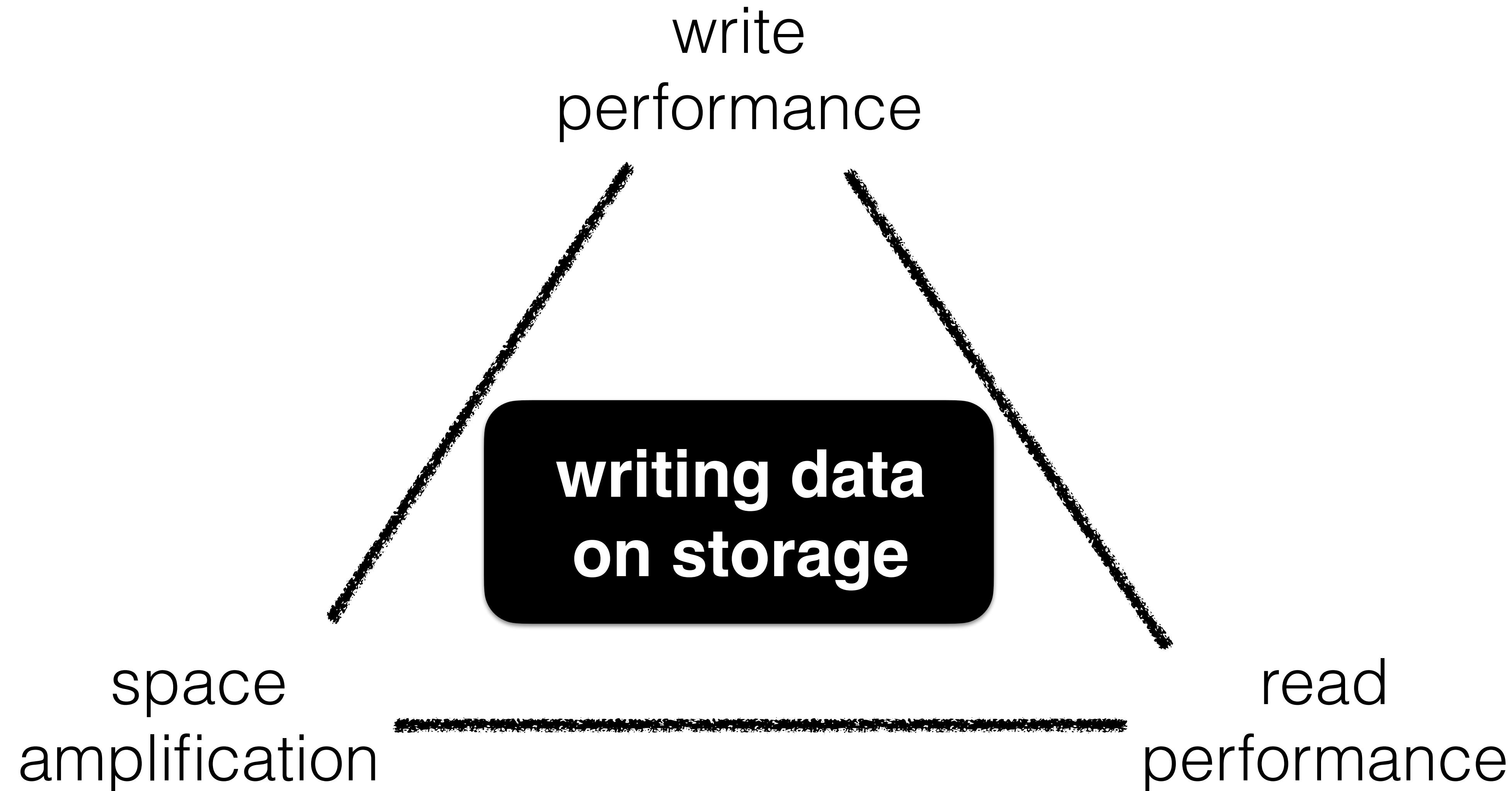


most data
on storage

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

99.9% on storage

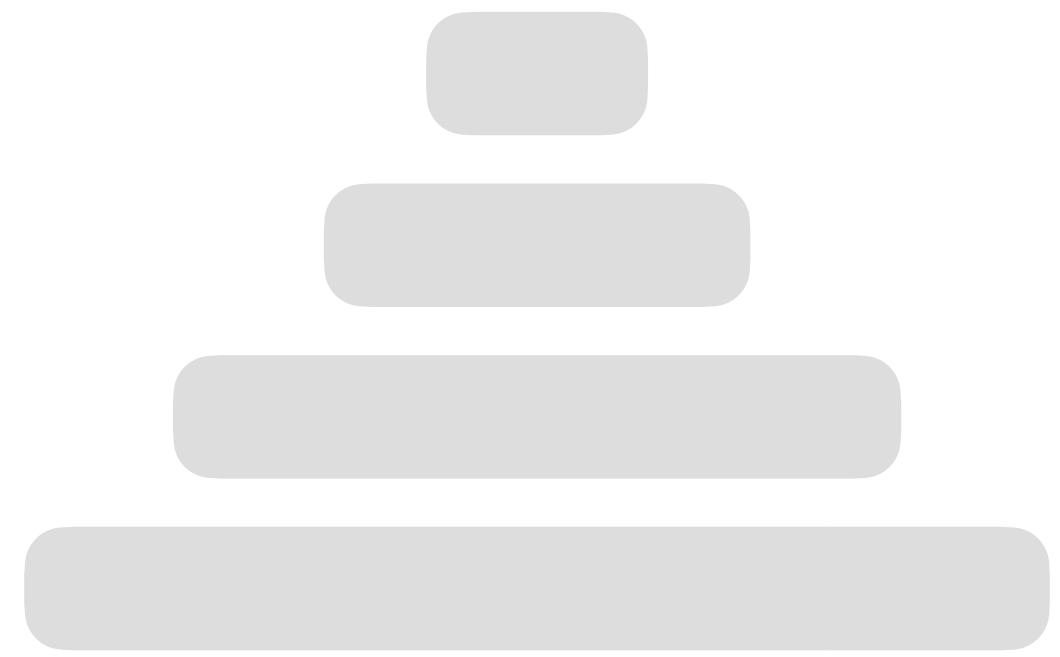
Performance Tradeoff



Data Layout

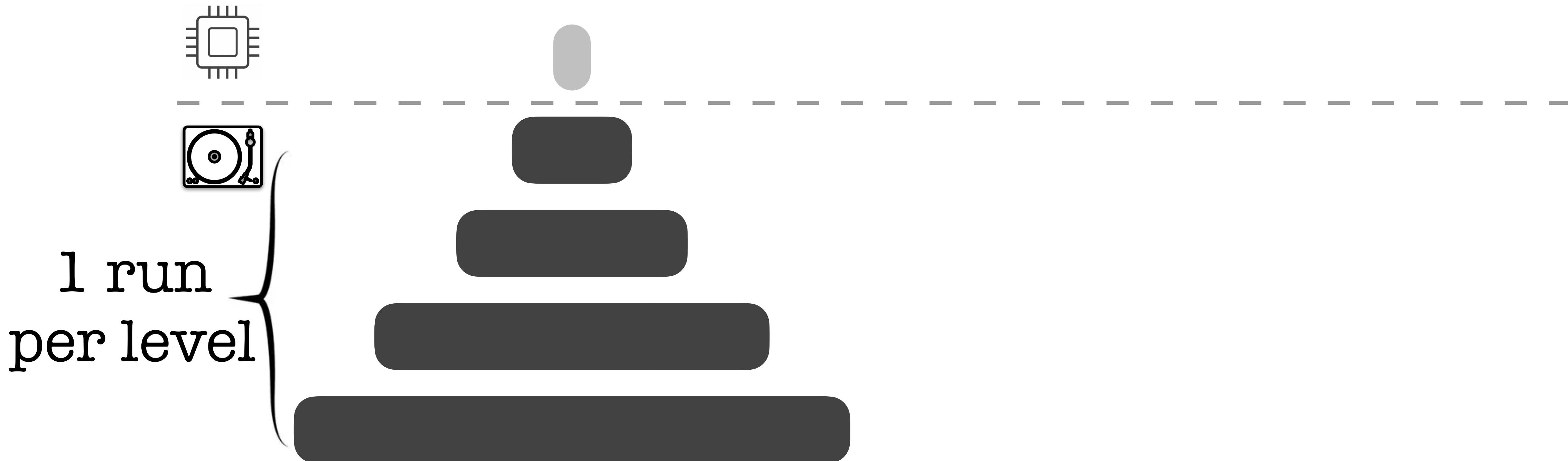
Classical LSM design: leveling

[eager merging]



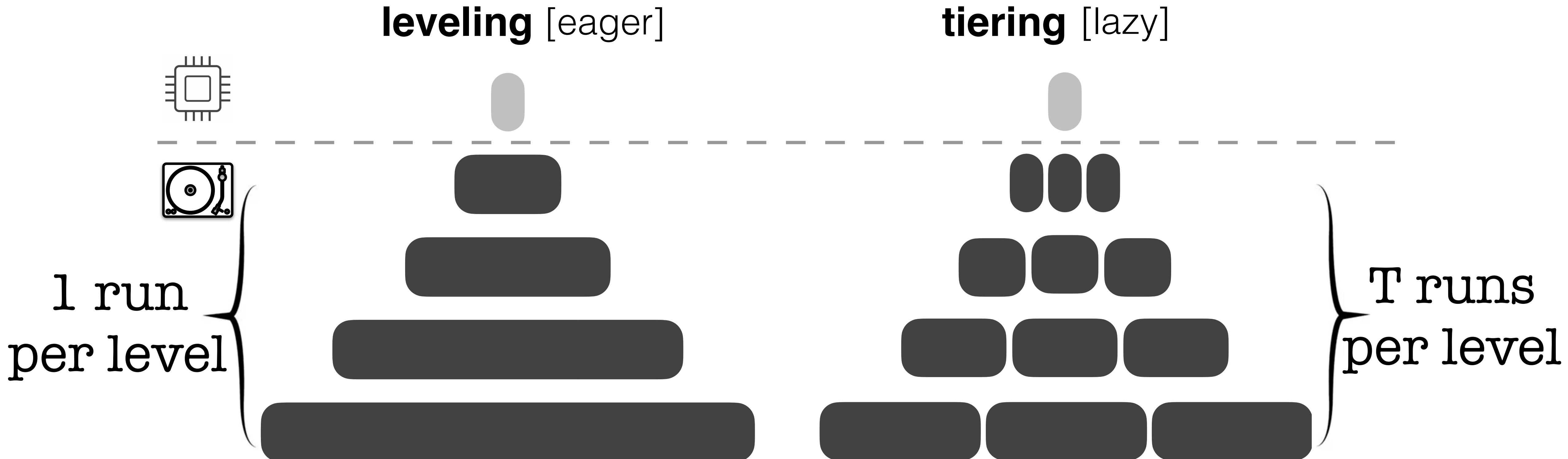
Data Layout

leveling [eager]



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

Data Layout



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

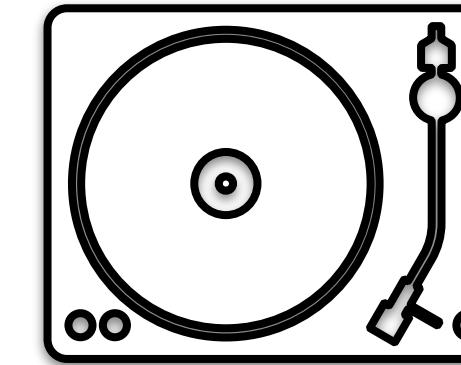
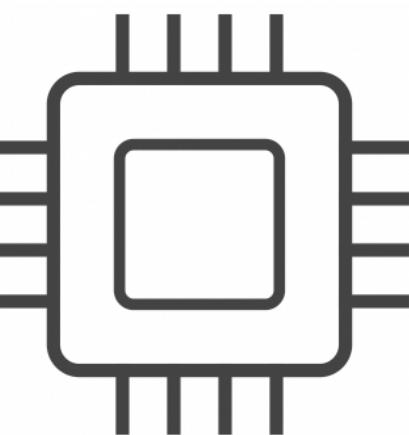
- poor read performance
- poor space amplification
- good ingestion performance

get(7)

buffer

Reduces
I/O cost
for
lookups

auxiliary data
structures



L1

L2

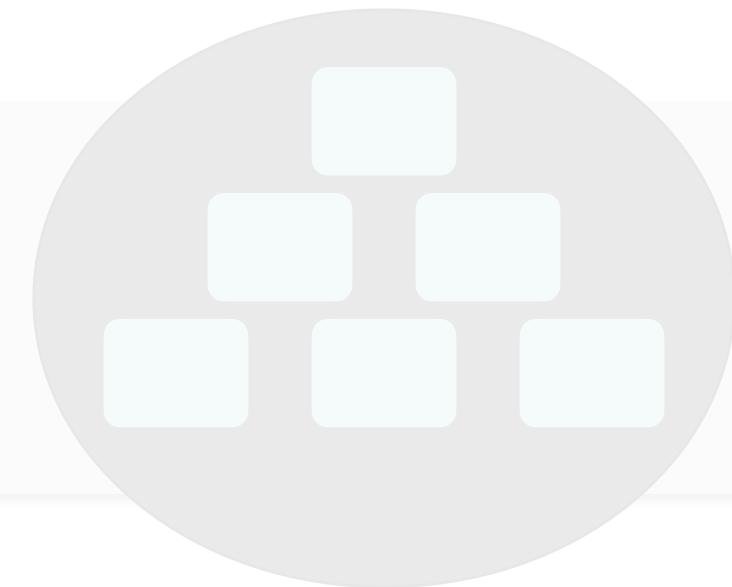
L3

L4

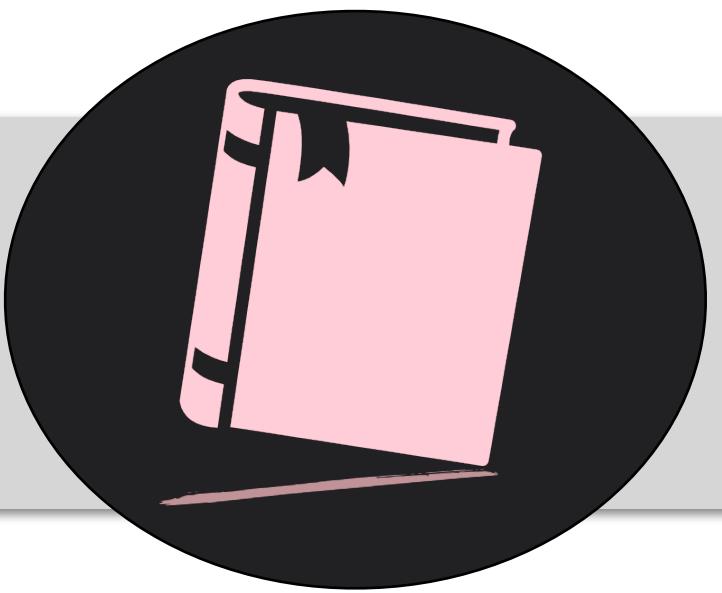
7

Outline

Part 1: **LSM Basics**



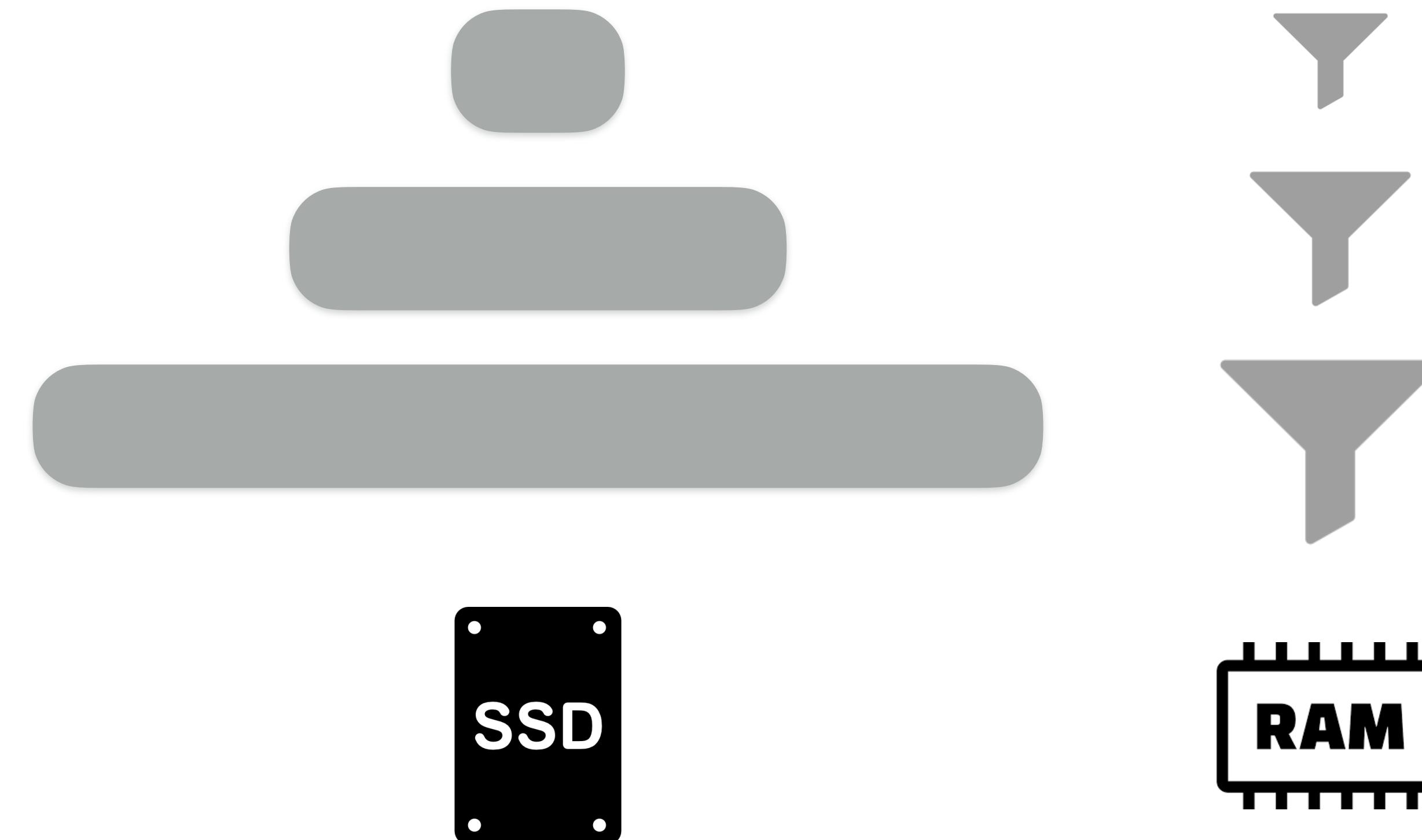
Part 2: Read Optimizations in LSMS



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



Filters to the Rescue



What is a filter

Does X exist?



Set

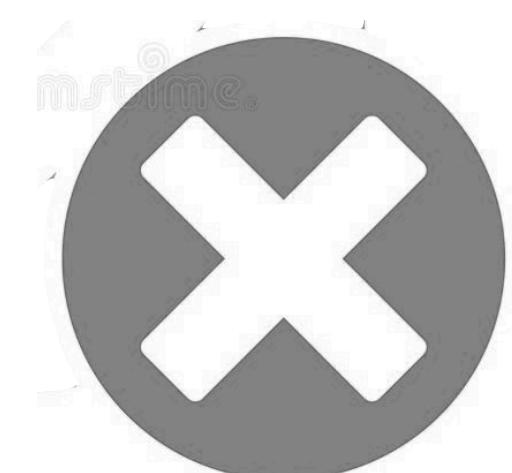
X Y Z

**Answers set
membership queries**

What is a filter



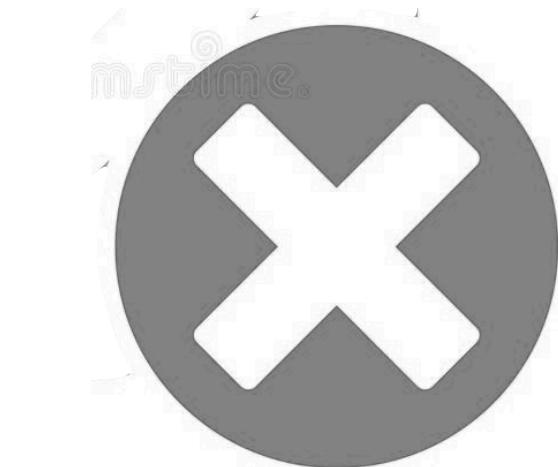
**No false
negatives**



What is a filter



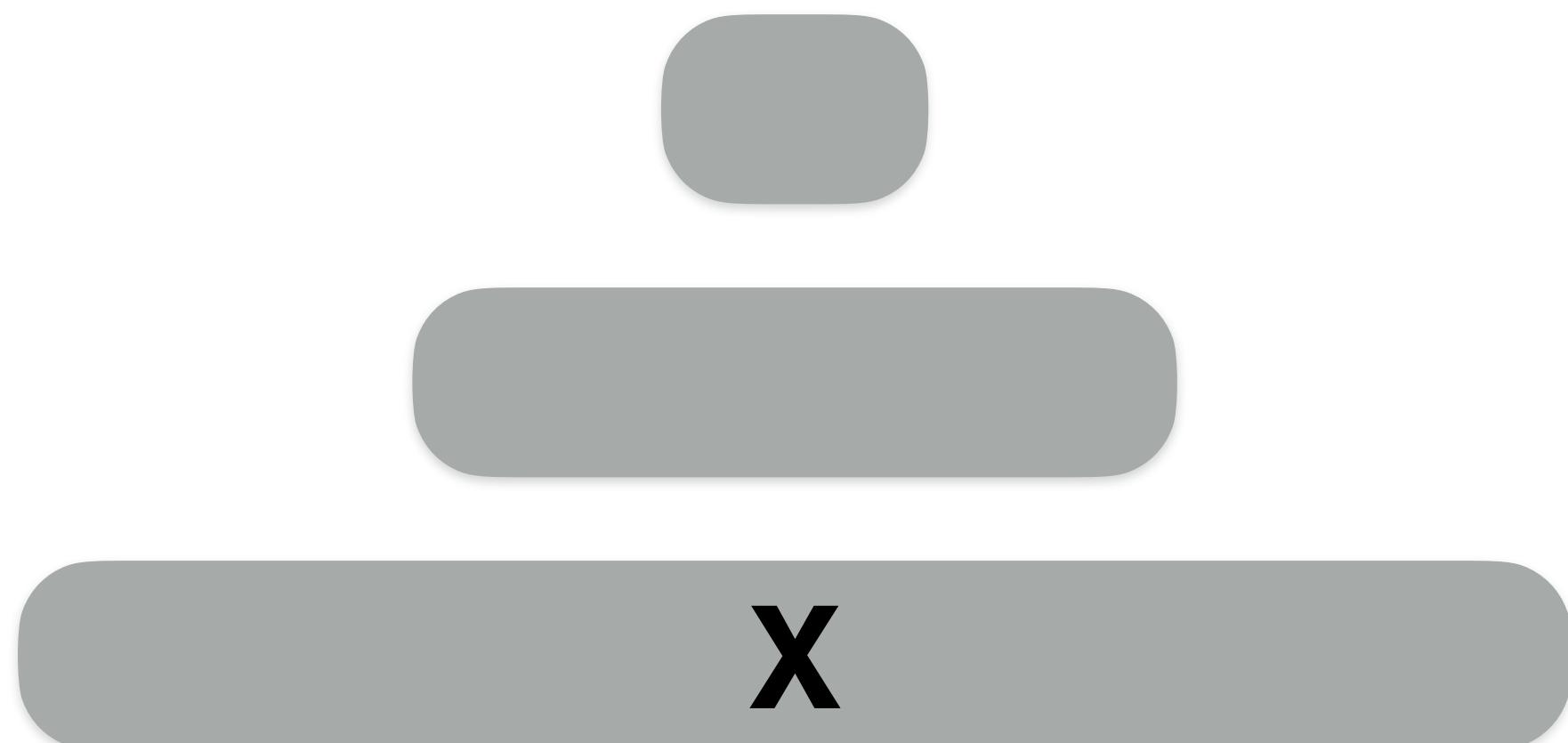
**No false
negatives**



**false positives with
tunable probability**



data



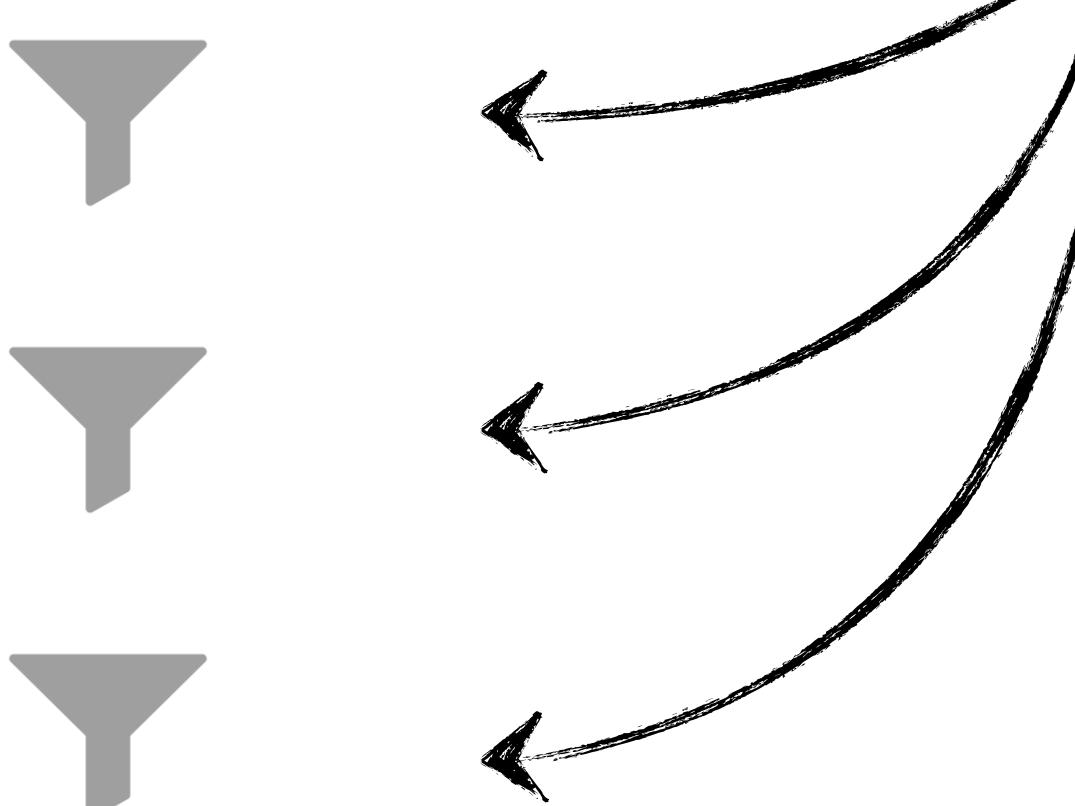
negative

false positive

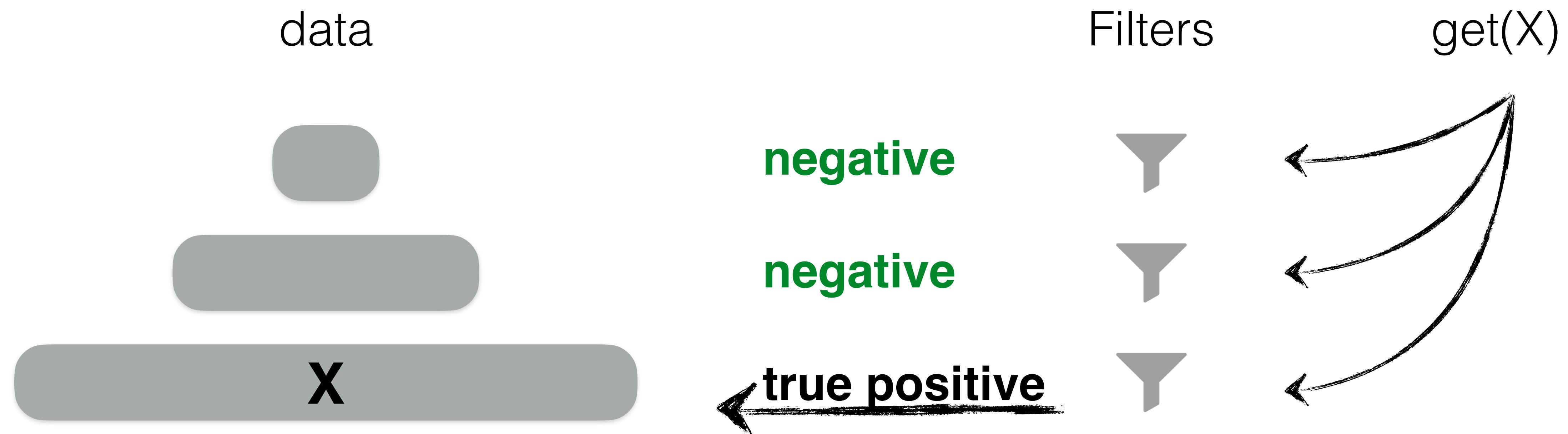
true positive

Filters
One per run

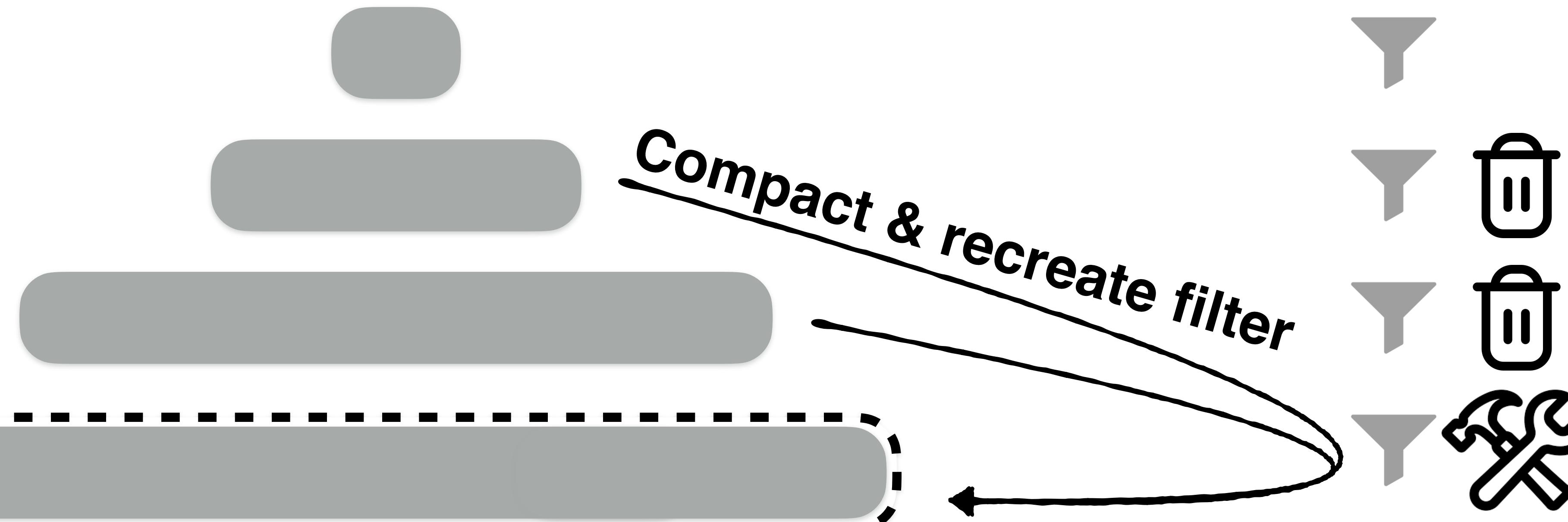
get(X)



more memory → fewer false positives



data

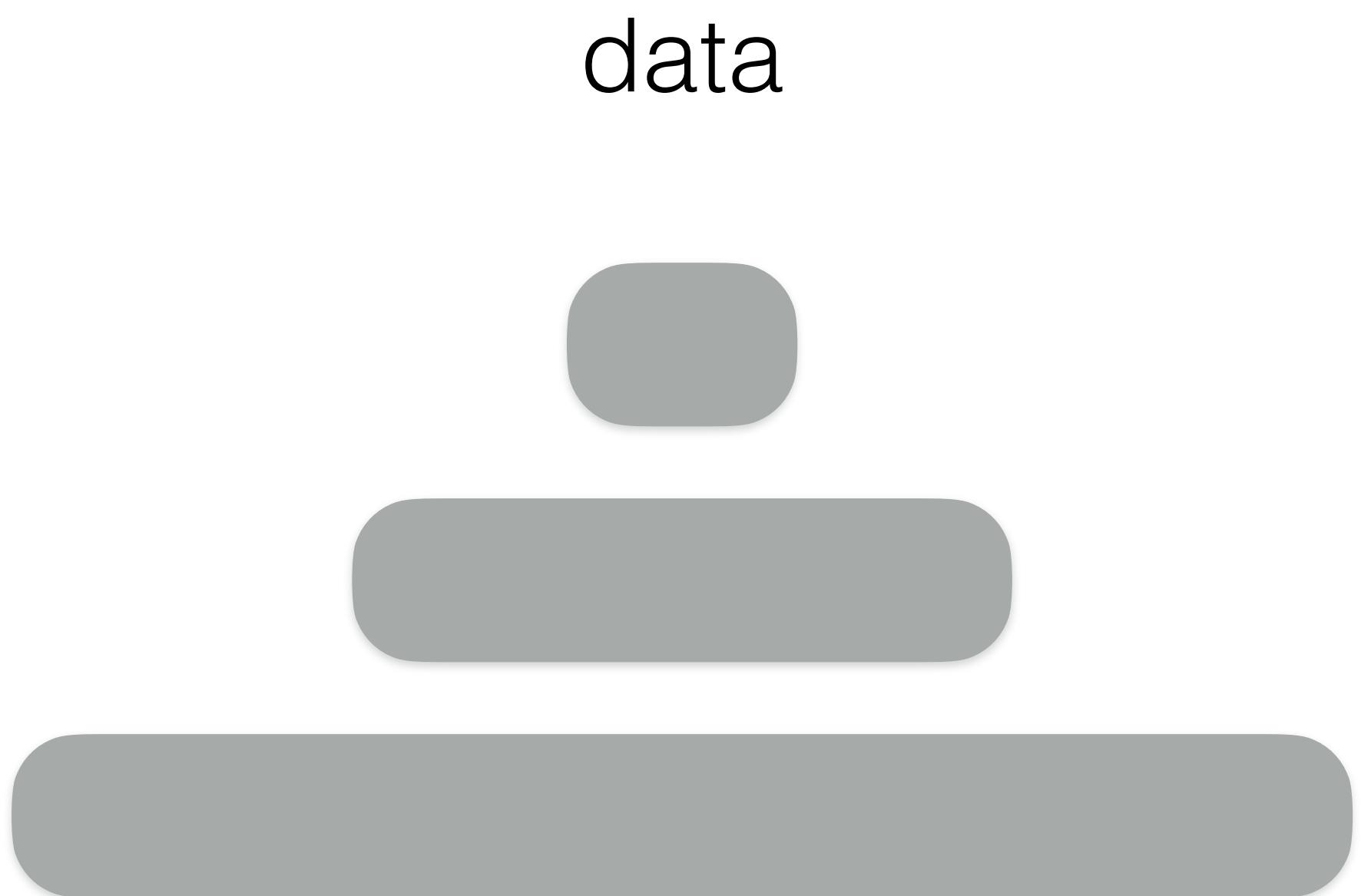




Bloom Filters

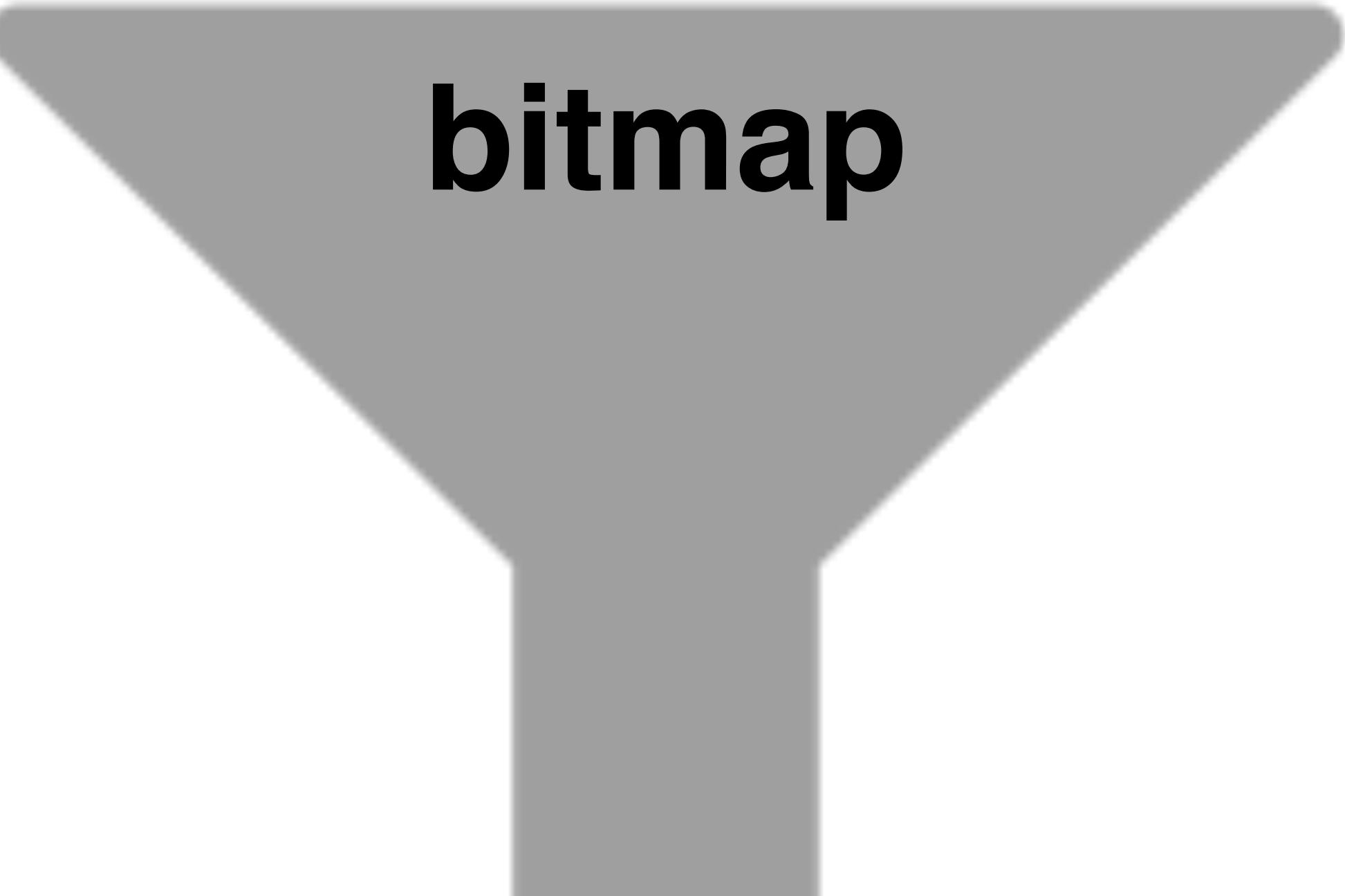


BloomCommunACM1970



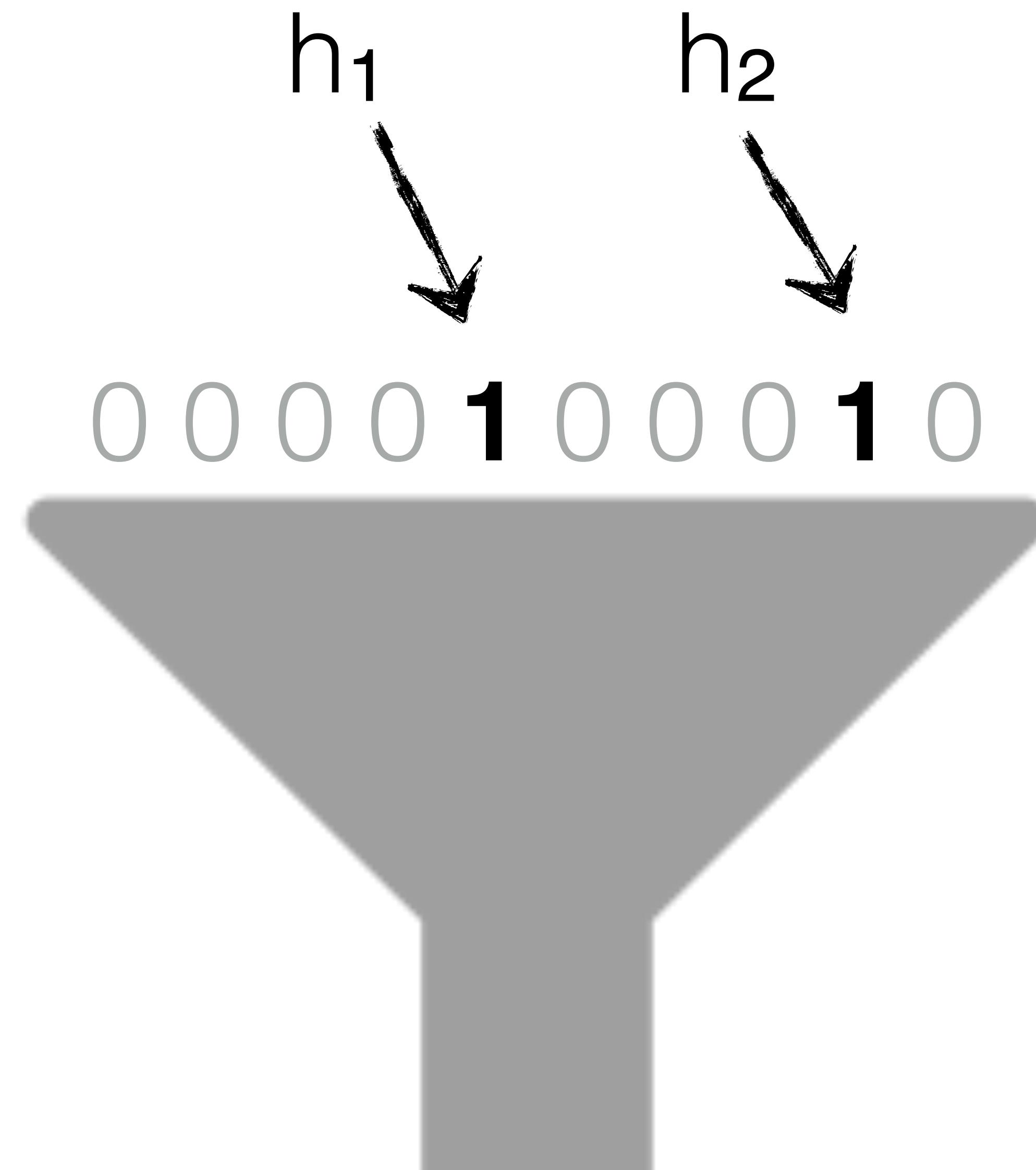
k hash functions

0 0 0 0 0 0 0 0 0

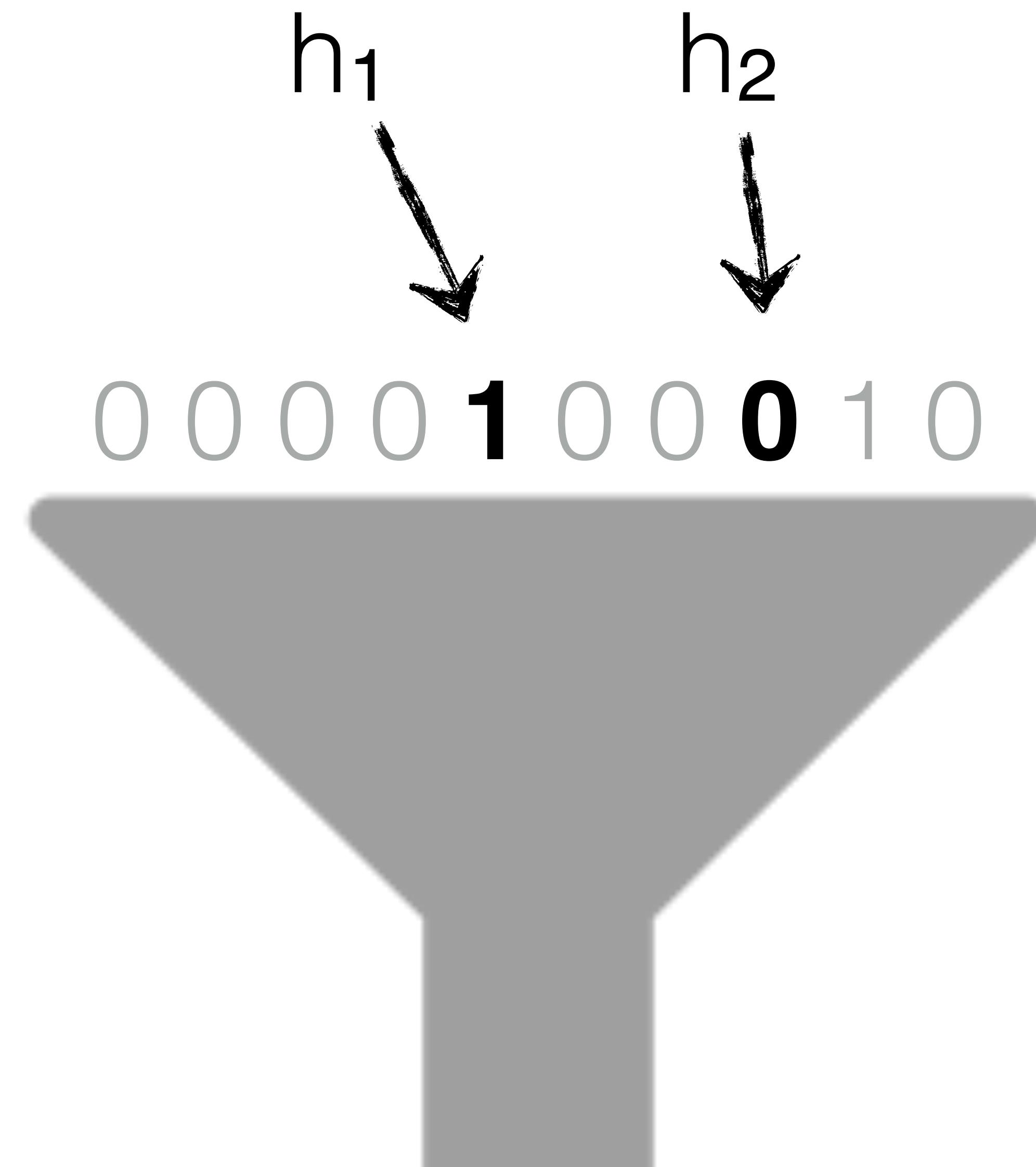


bitmap

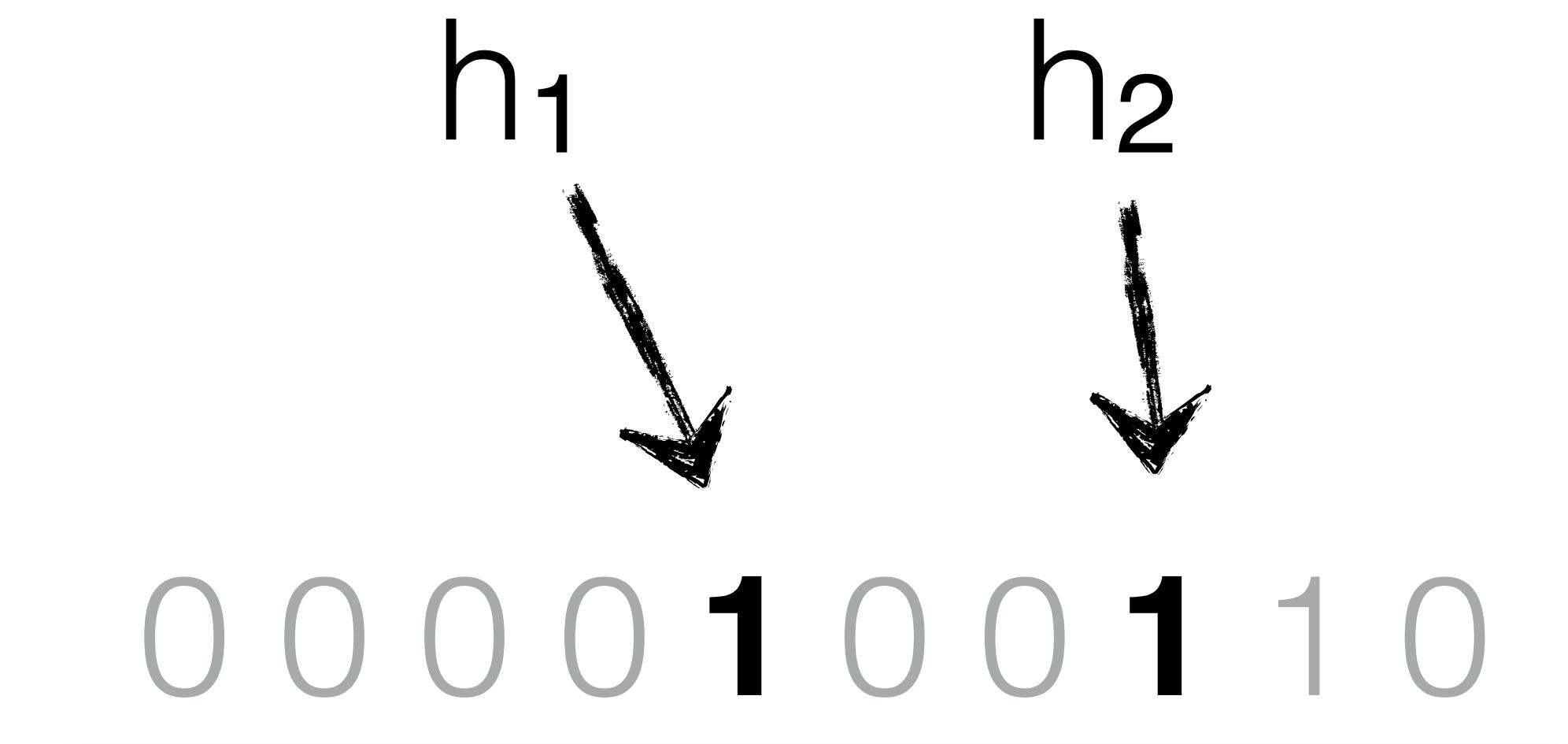
insert: Set from 0 to 1 or keep 1



negative lookup: at least one bit is zero



true or false positive lookup



Optimal number of hash functions

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

A horizontal double-headed arrow spans the width of the hash functions h_1, \dots, h_k . Above the arrow is the formula $= \ln(2) \cdot M \quad \leftarrow \text{bits / entry}$. Below the arrow, the hash functions h_1, \dots, h_k are listed. Three arrows point downwards from the hash functions towards a large grey downward-pointing arrow at the bottom.

With M bits / entry

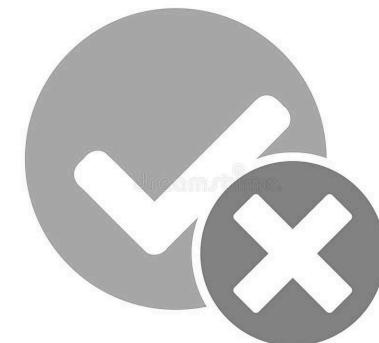
Optimal number of hash functions = $\ln(2) \cdot M$

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

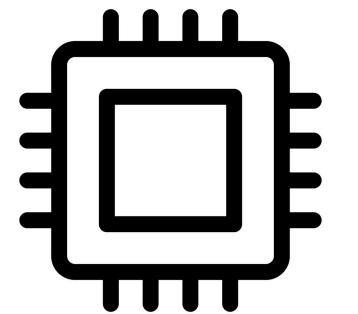


5 fronts

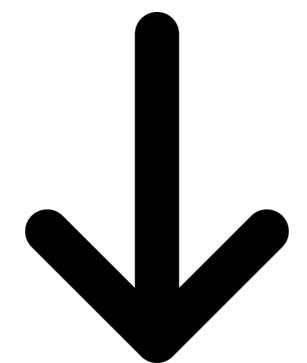
Holistic
Tuning



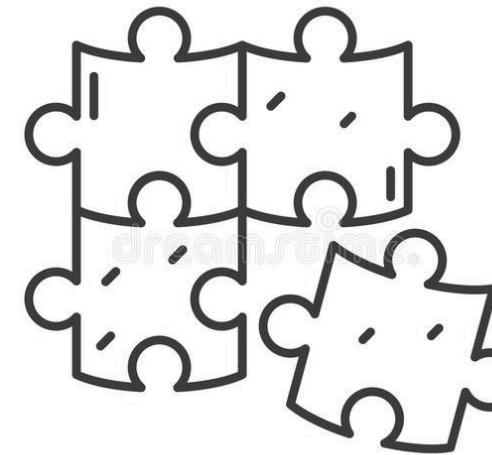
CPU



Lowering
Constants



Unification



Range



Holistic Tuning



Monkey



Dostoevsky



LSM-Bush



DayanSIGMOD17



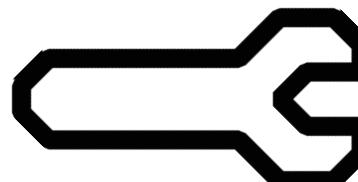
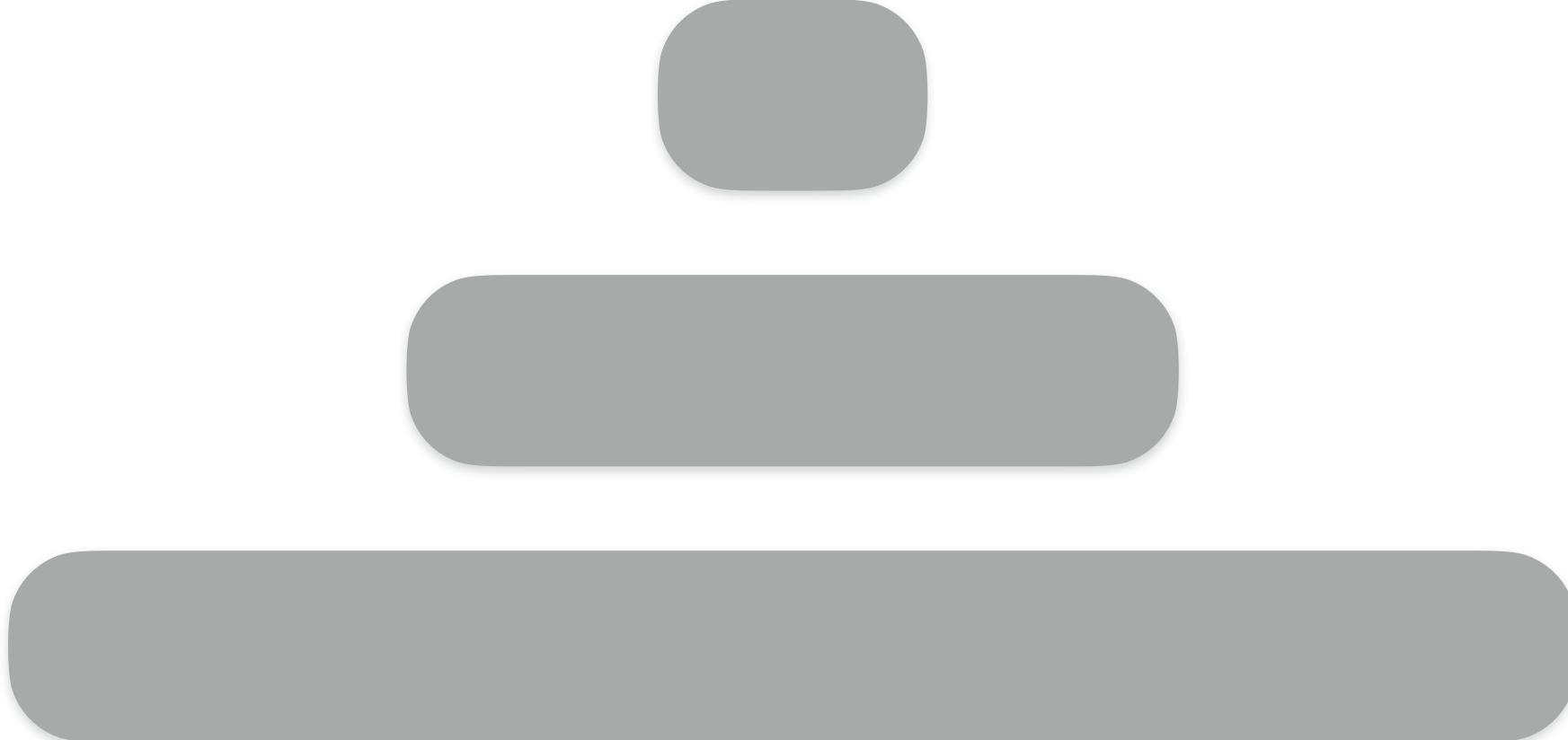
DayanSIGMOD18



DayanSIGMOD19

Monkey: Optimal Navigable **K**ey-Value Store

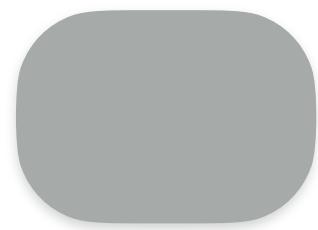
data



**Bloom
filters**



data



Bloom
filters



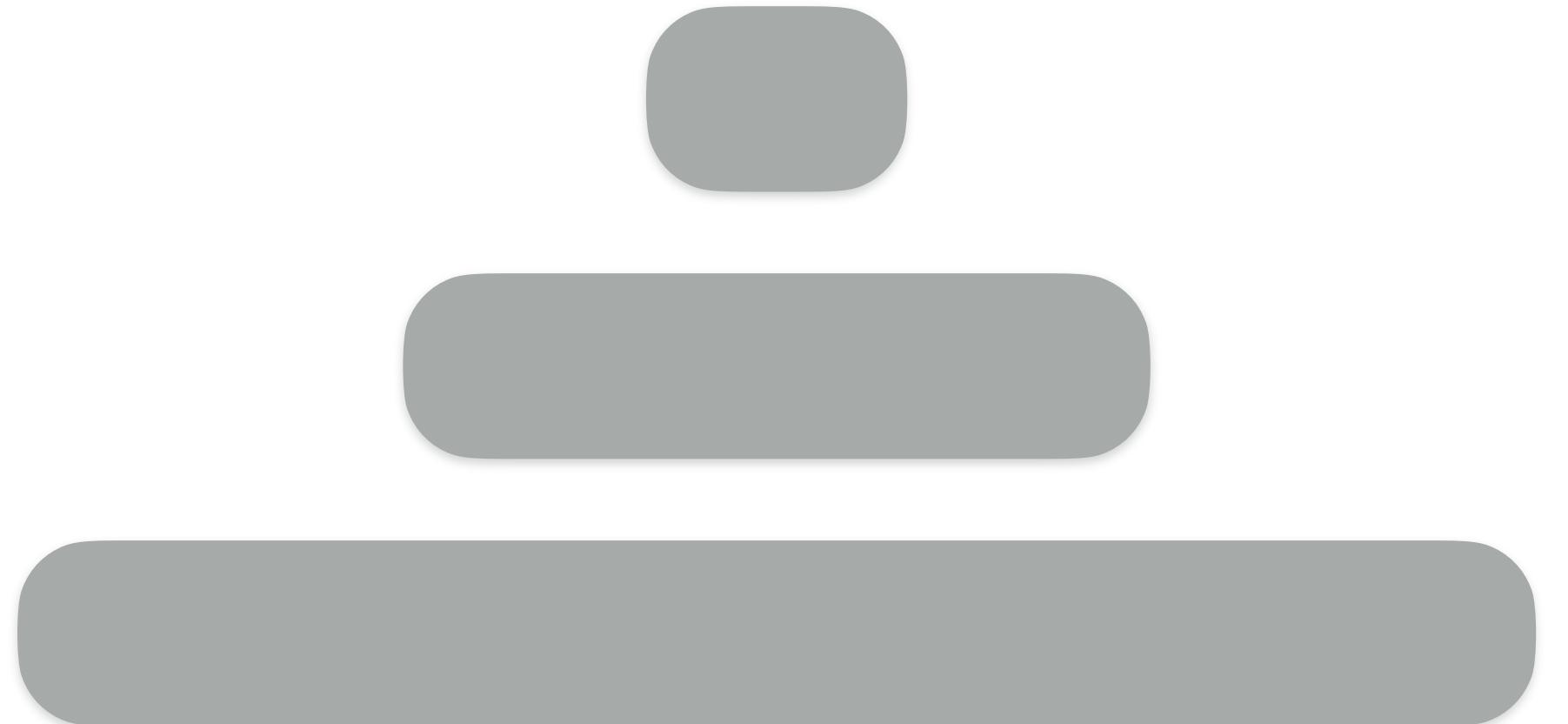
bits/entry

M

M

M

data



Bloom
filters



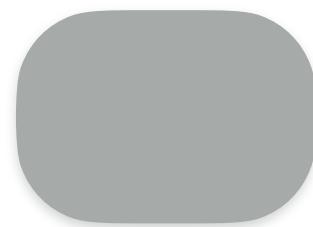
bits/entry

M

M

M

data



Bloom
filters



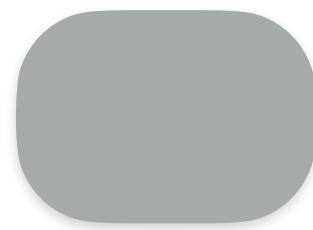
**false
positive rate**

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

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

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

data



Bloom
filters



**false
positive rate**

2^{-M}

2^{-M}

2^{-M}

Bloom
filters

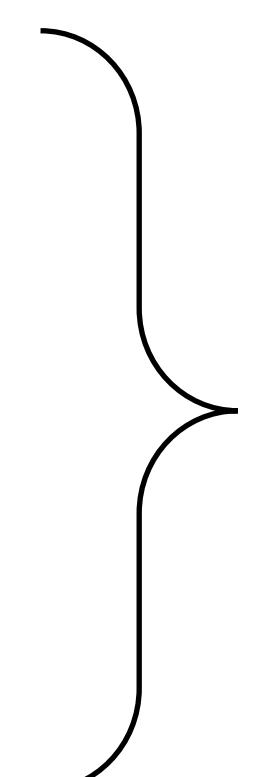


false
positive rate

2^{-M}

2^{-M}

2^{-M}



=

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

Bloom
filters

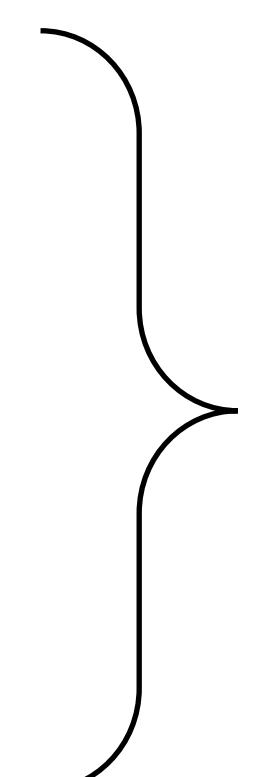


false
positive rate

2^{-M}

2^{-M}

2^{-M}



=

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

Bloom
filters

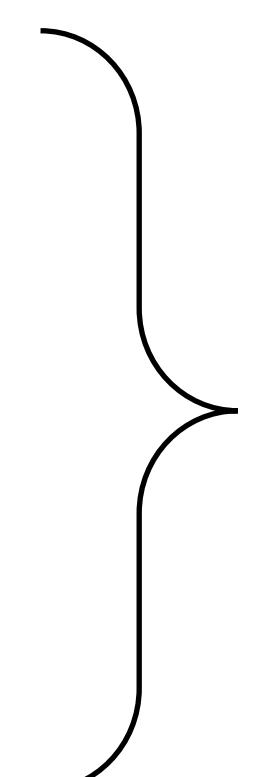


false
positive rate

2^{-M}

2^{-M}

2^{-M}



=

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



Bloom
filters



false
positive rate

2^{-M}



2^{-M}

**most
memory**



2^{-M}

Bloom
filters



**most
memory**



false
positive rate

2^{-M}

2^{-M}

2^{-M}

saves at most 1 access!

bits / entry



$M + 2$

$M + 1$

$M - 1$

reallocates



false
positive rates



$2^{-(M + 2)}$



$2^{-(M + 1)}$



$2^{-(M - 1)}$





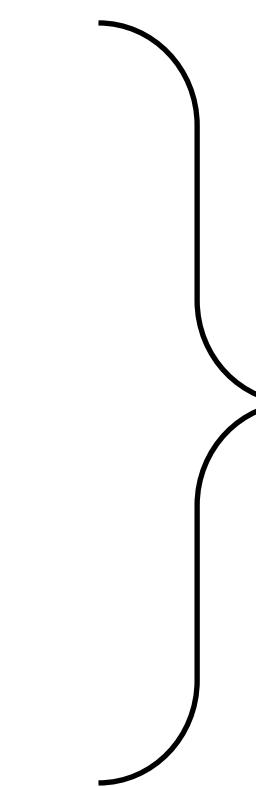
\propto $2^{-M} / T^2$



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



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

 $2^{-M} / T^2$ $2^{-M} / T^1$ $2^{-M} / T^0$ 

**geometric
progression**

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



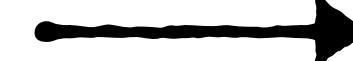
Faster worst case

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

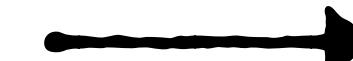
Monkey opens up new ways of optimizing write performance without sacrificing get performance



Monkey



Dostoevsky



LSM-Bush



DayanSIGMOD17



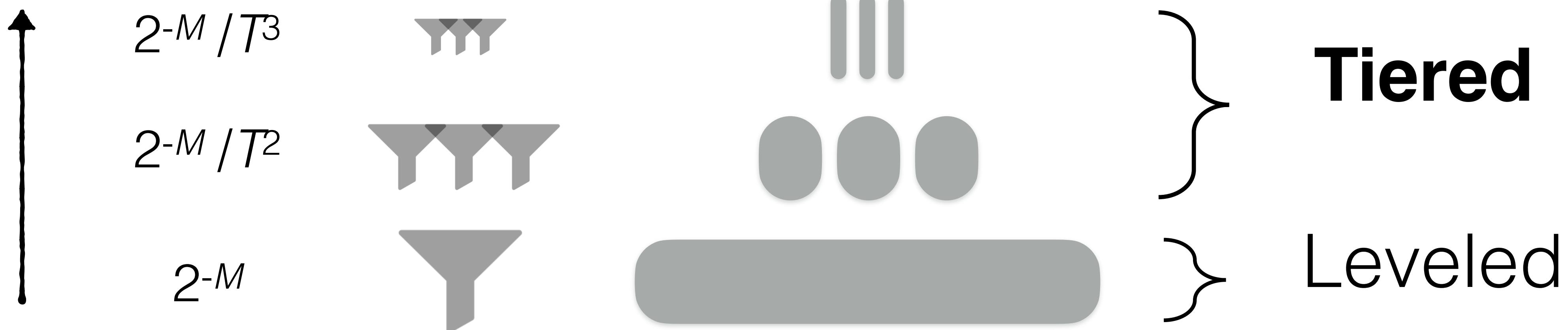
DayanSIGMOD18



DayanSIGMOD19

Dostoevsky

Smaller false positive rates



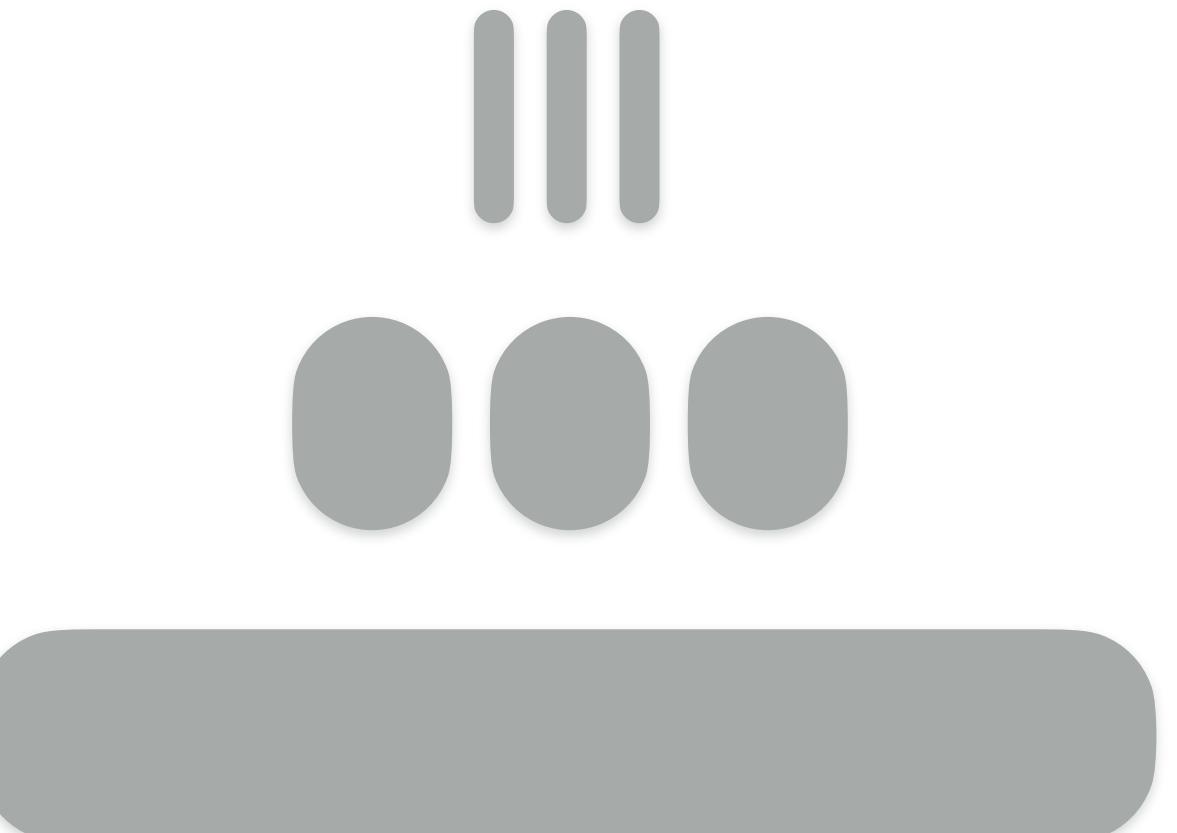
gets

$O(2^M)$

writes

$O(T + \log_T N)$

=



$O(1)$

+

$O(1)$

+

$O(T)$

gets

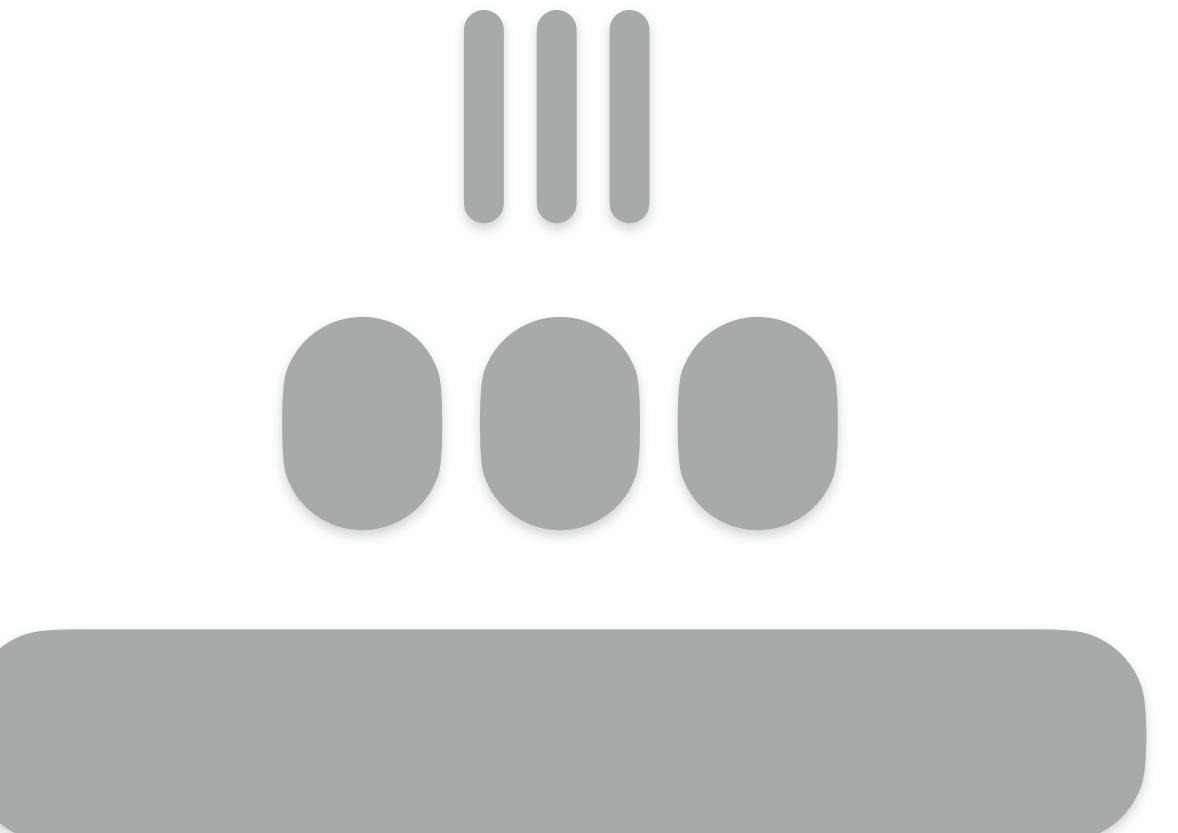
$$O(2^{-M})$$

writes

$$O(\textcolor{green}{T} + \log_{\tau} N) < O(\textcolor{red}{T} \cdot \log_{\tau} N)$$

leveling

=



$$O(1)$$

+

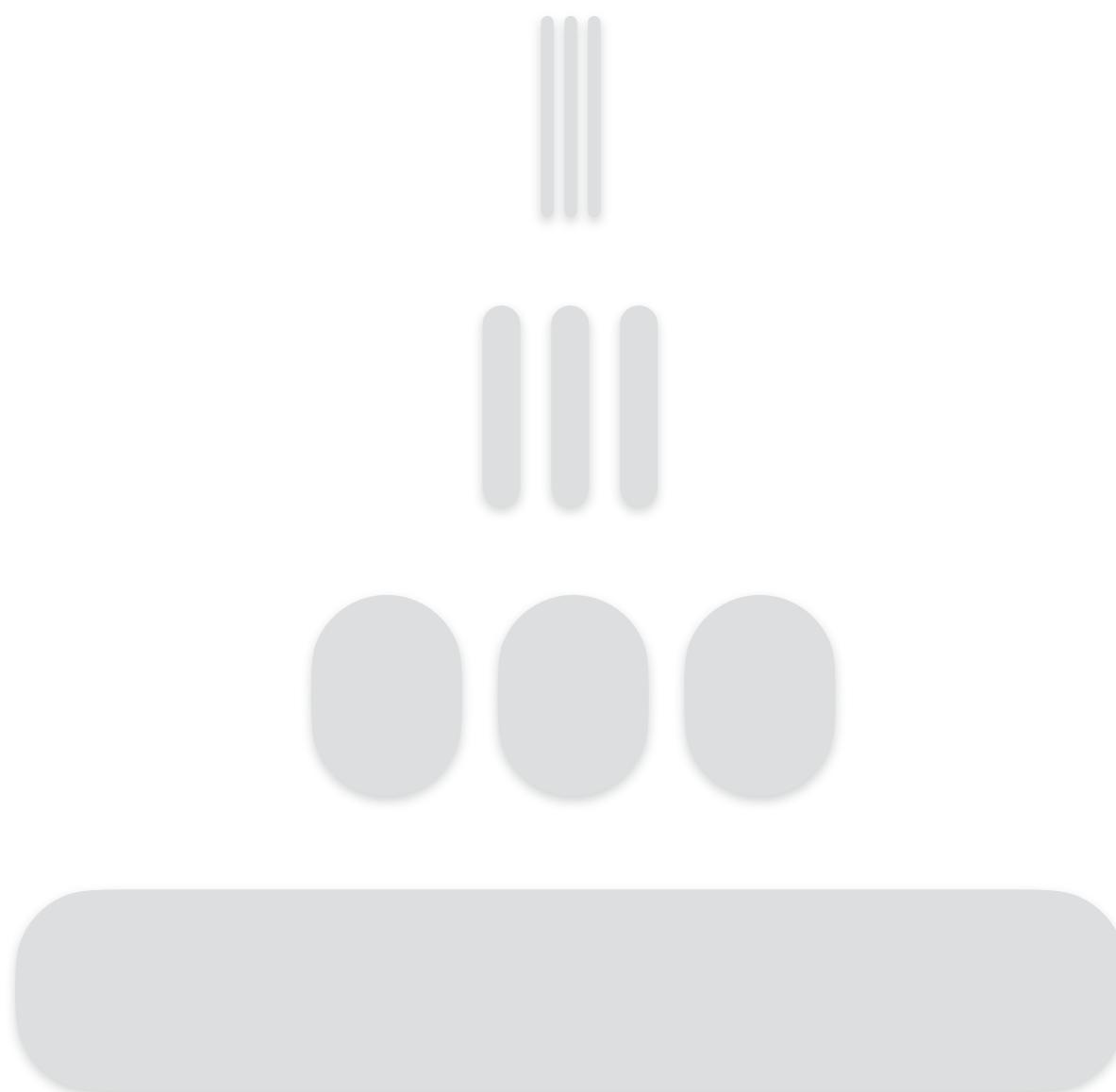
$$O(1)$$

+

$$O(T)$$

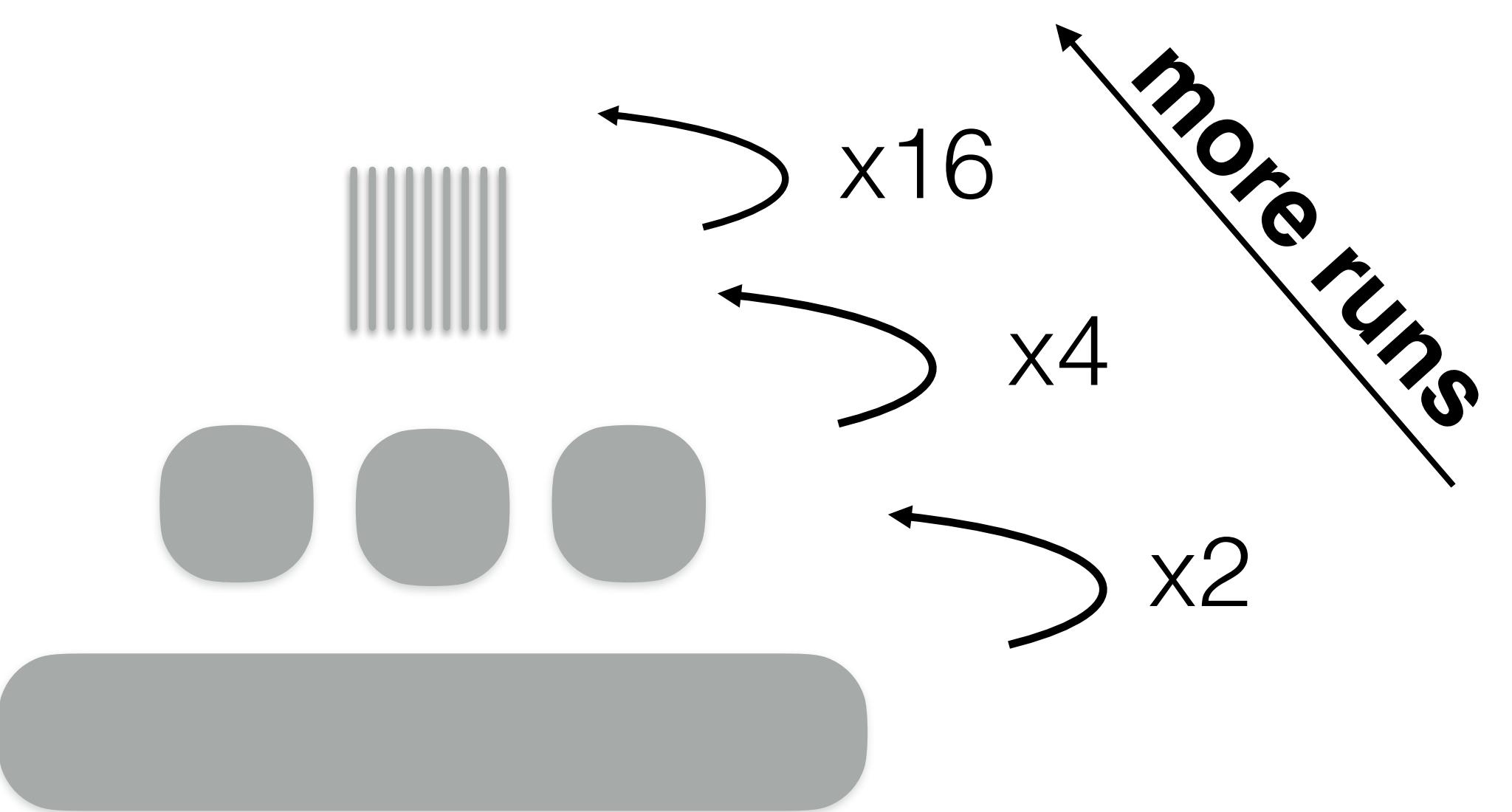
Dostoevsky

$O(T + \log_T N)$



LSM-Bush

$O(\log_2 \log_T N)$





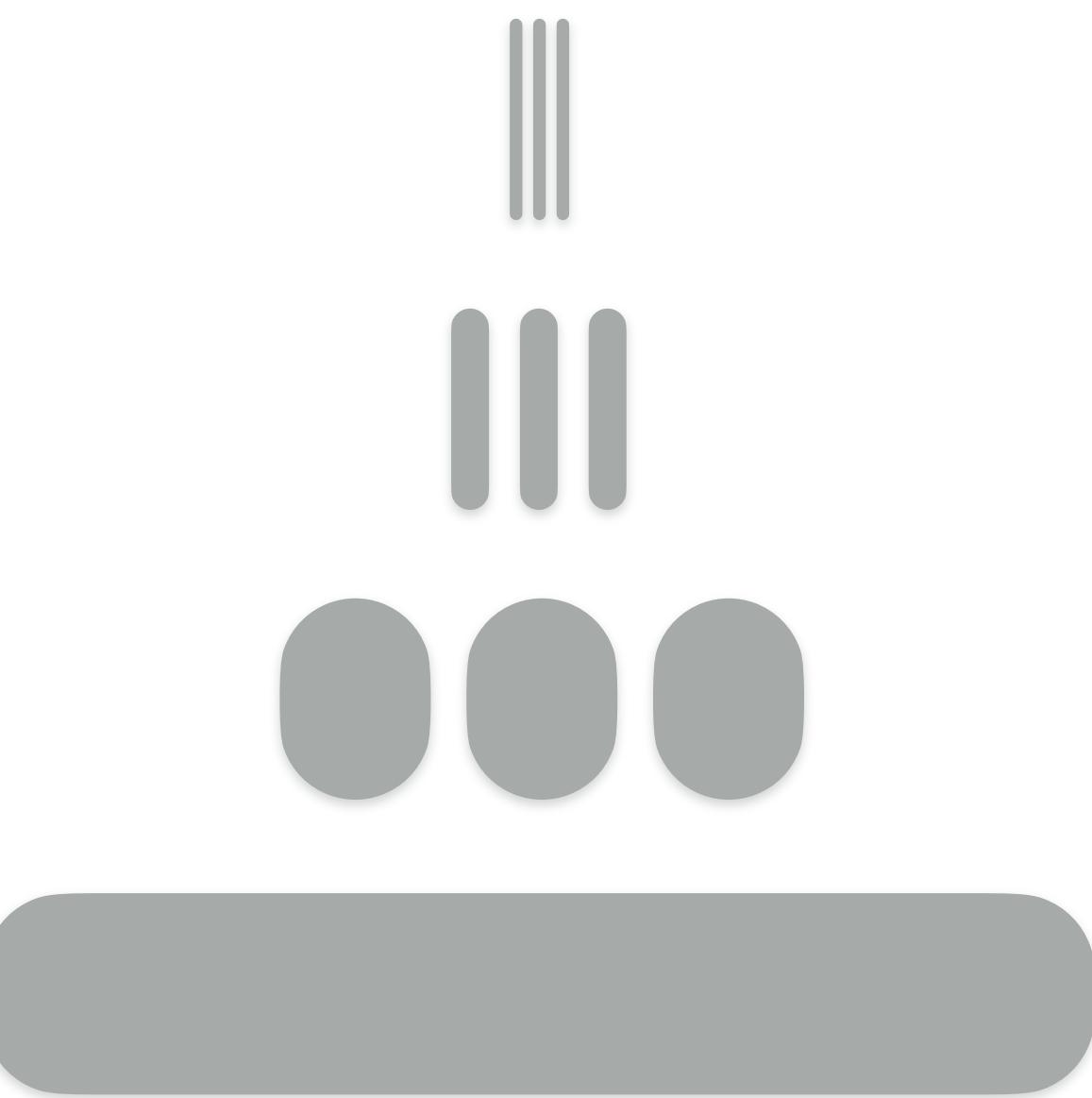
Monkey w. leveling



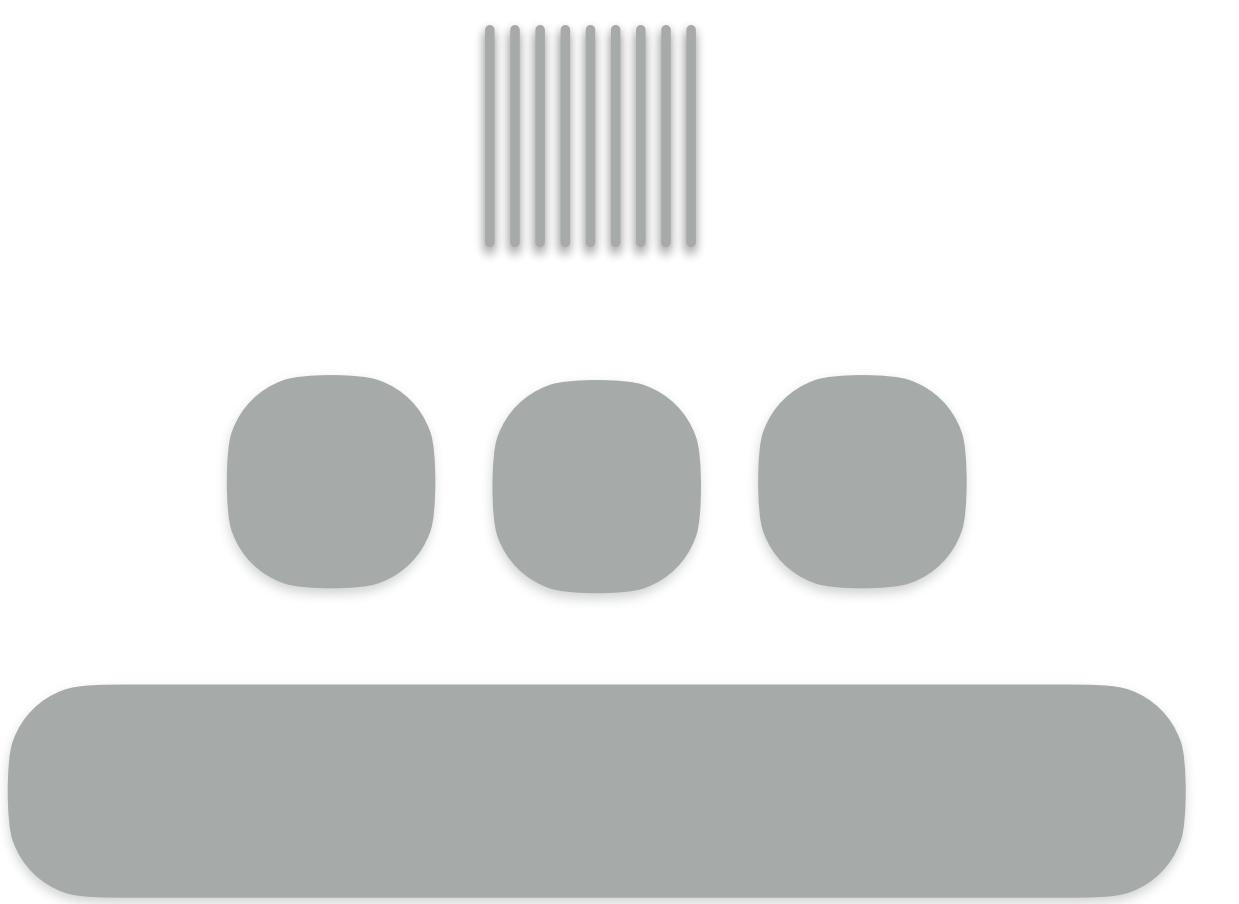
←→
Cheaper range



Dostoevsky



LSM-Bush



←→
Cheaper writes



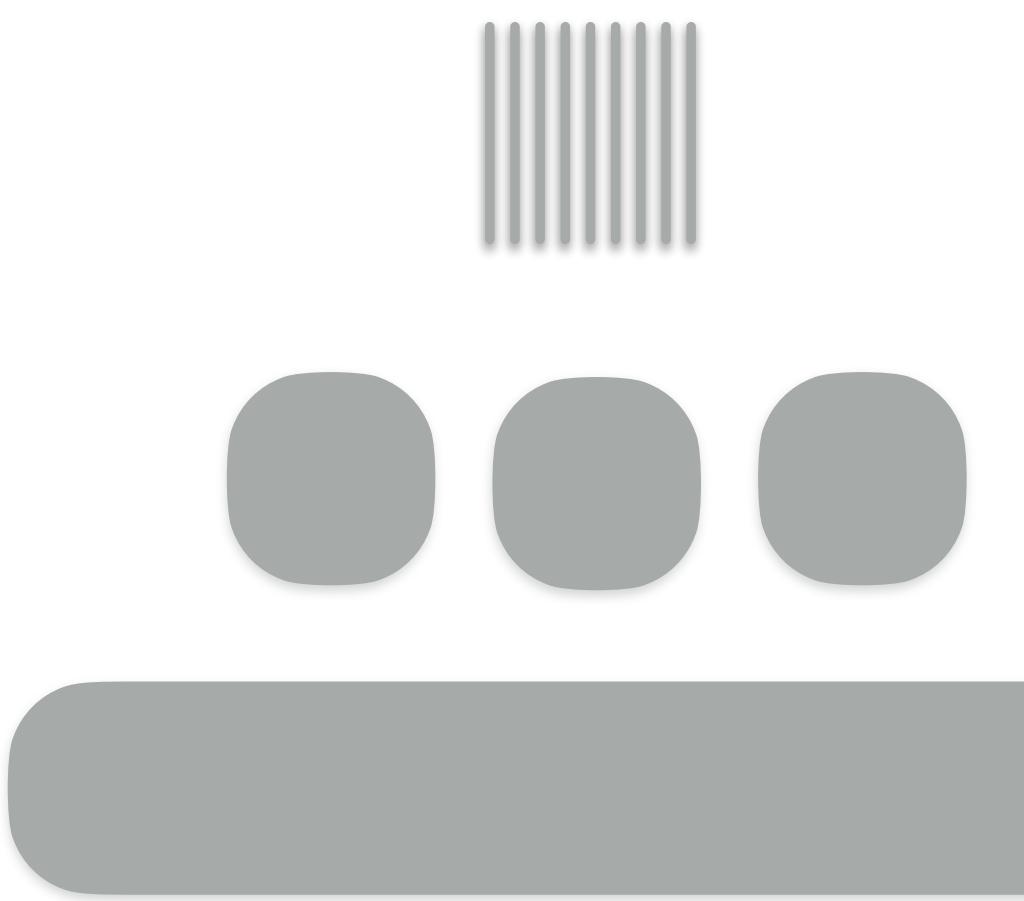
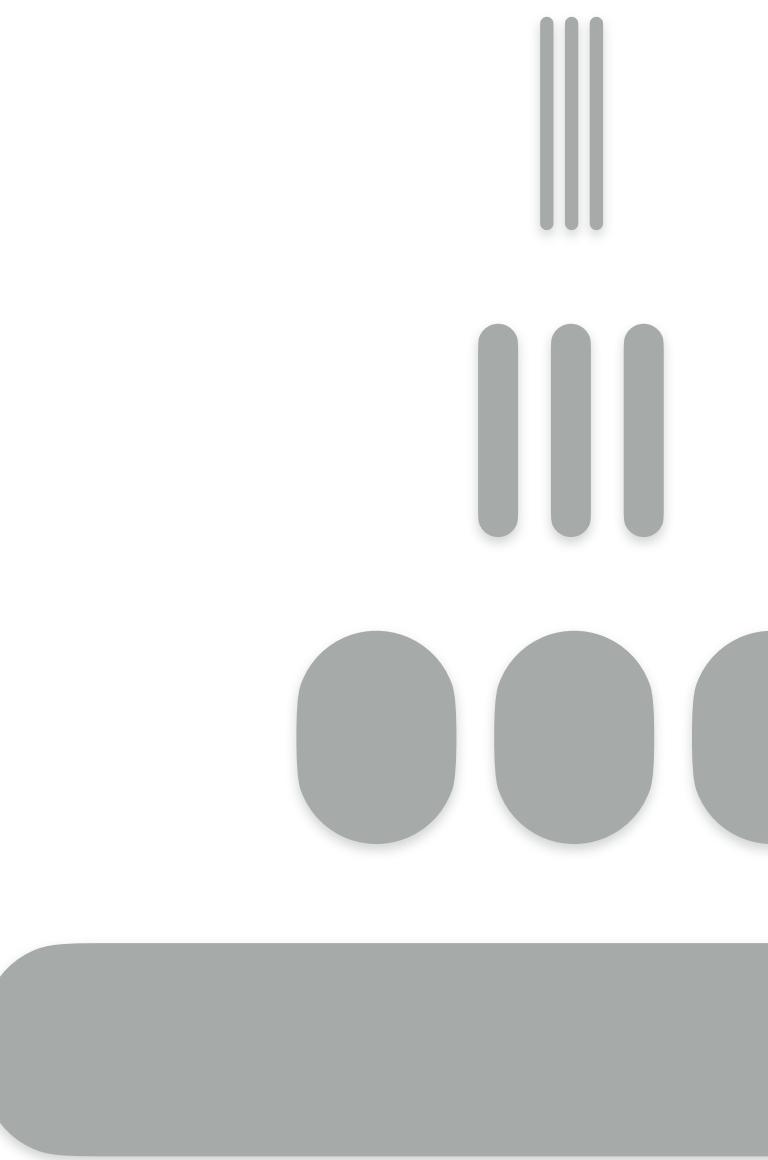
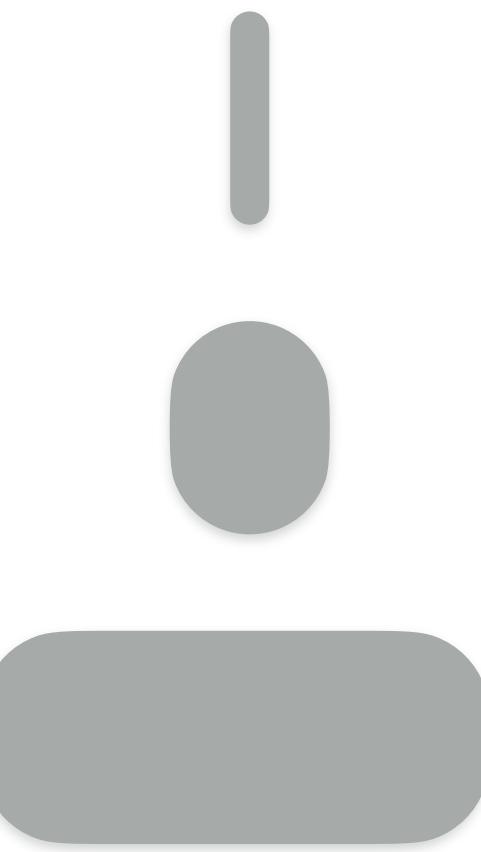
Monkey w. leveling



Dostoevsky



LSM-Bush

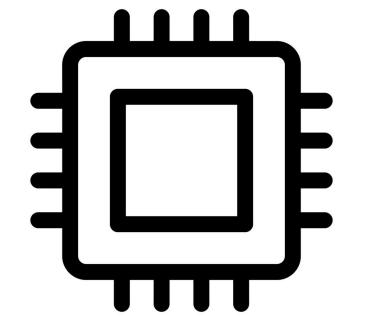


Great point reads all across

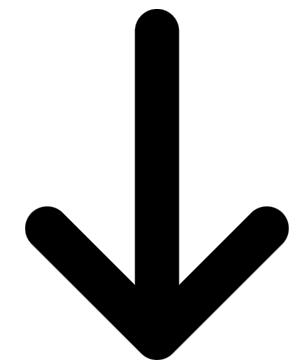
Holistic
Tuning



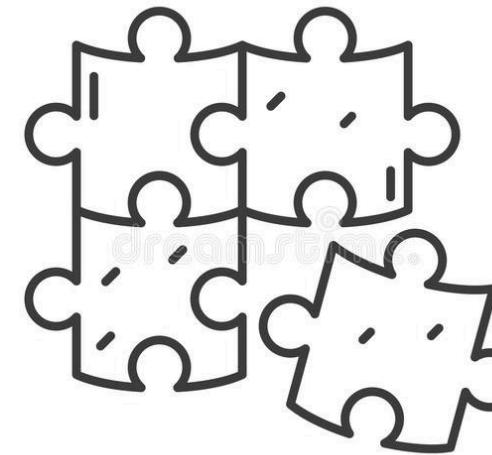
CPU



Lowering
Constants



Unification



Range



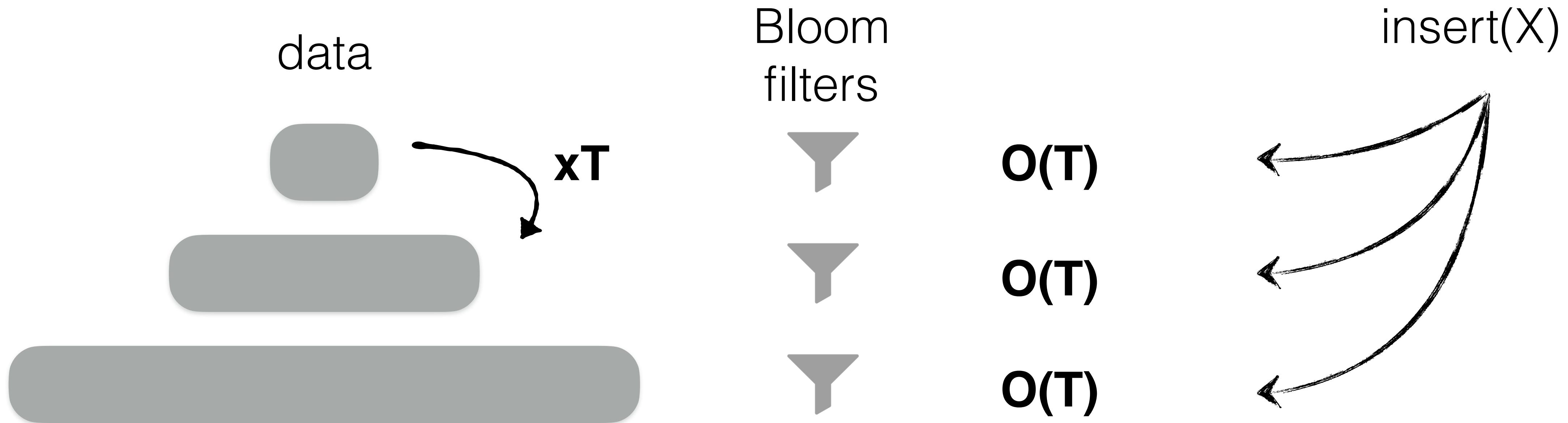
5 fronts

Bloom
filters



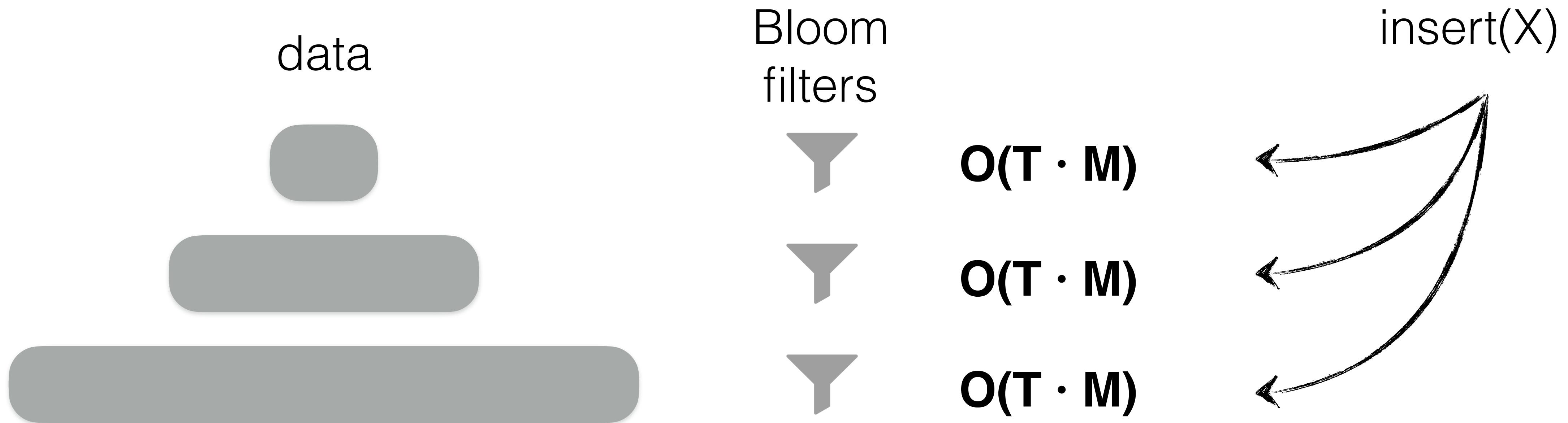
CPU overhead?

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



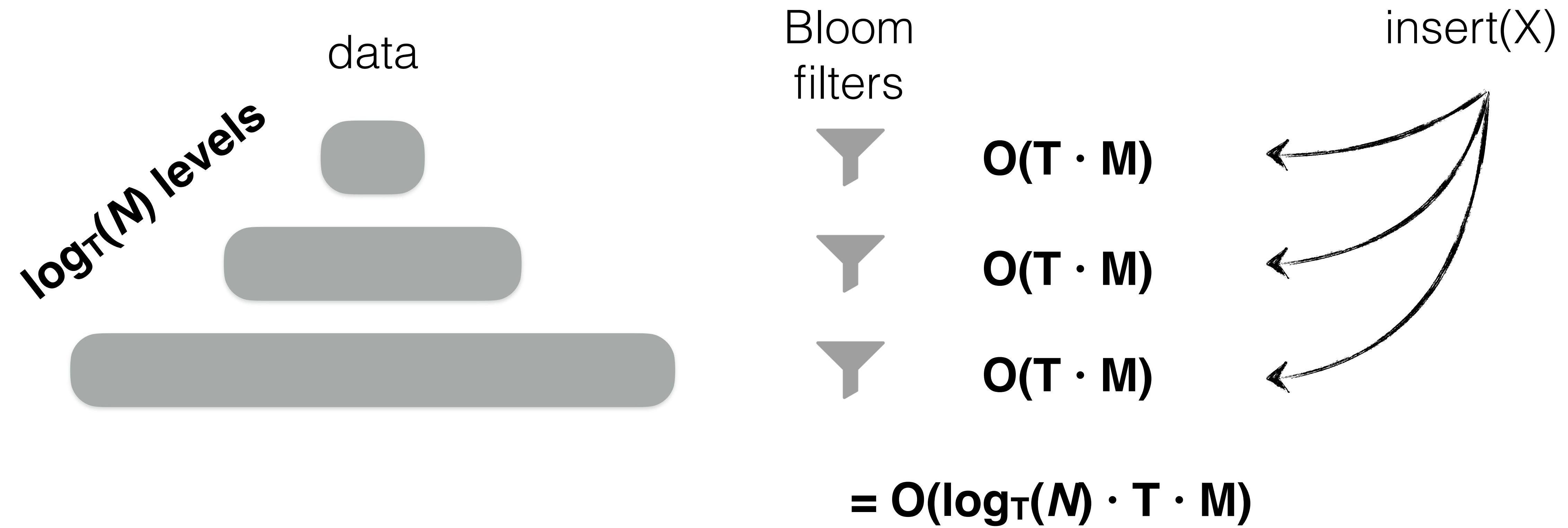
Each key is inserted $O(T)$ times per level into a filter

Each filter insertion uses $M \cdot \ln(2)$ hash functions

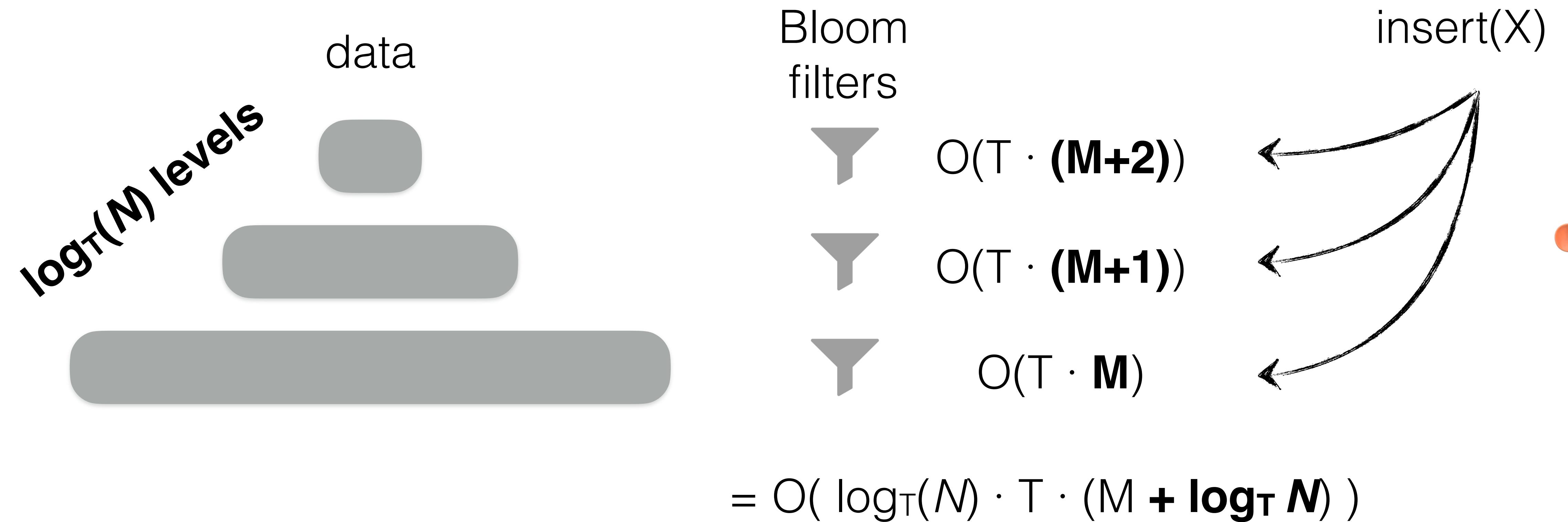


Each key is inserted $O(T)$ times per level into a filter

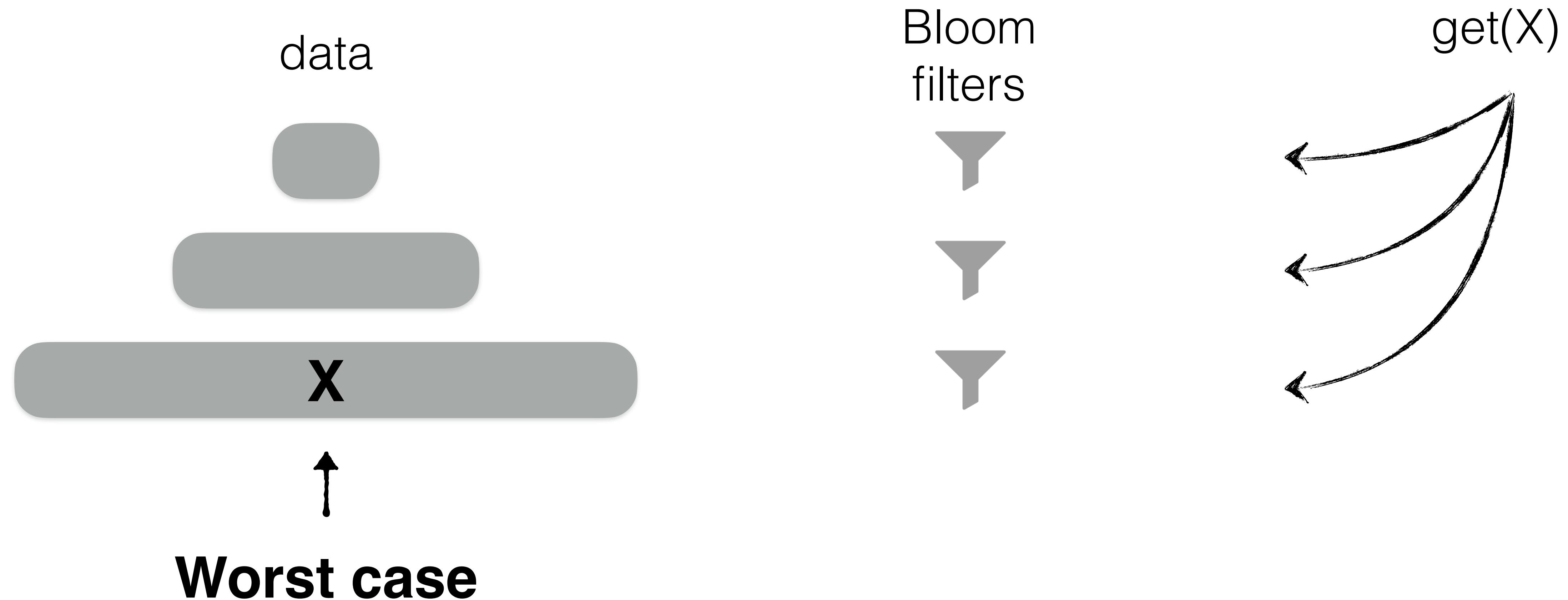
Each filter insertion uses $M \cdot \ln(2)$ hash functions



With Monkey more hash functions are used



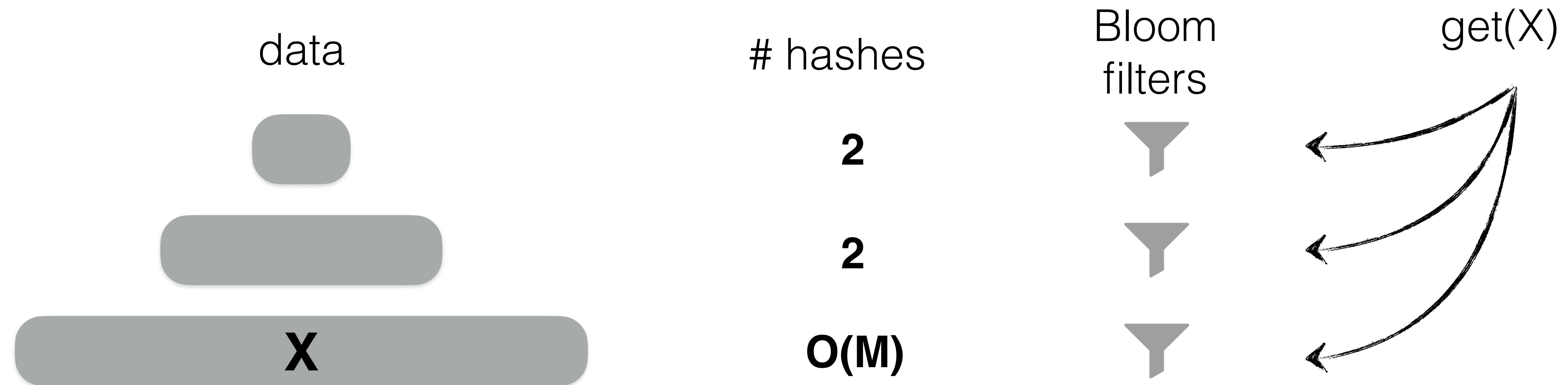
How about get cost?



Expected Negative Query Cost ≈ 2



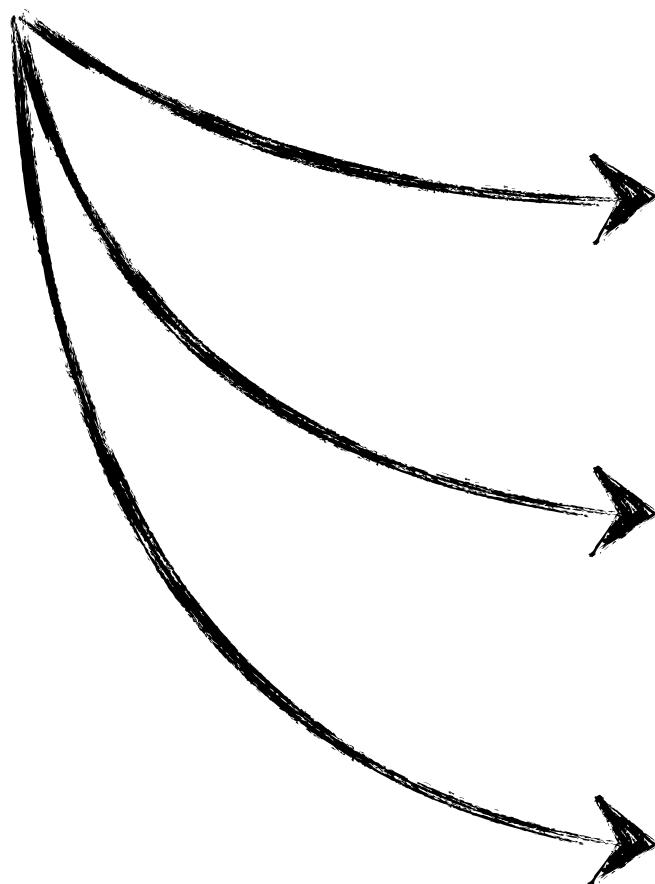
Positive Query Cost $\approx M \cdot \ln(2)$





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

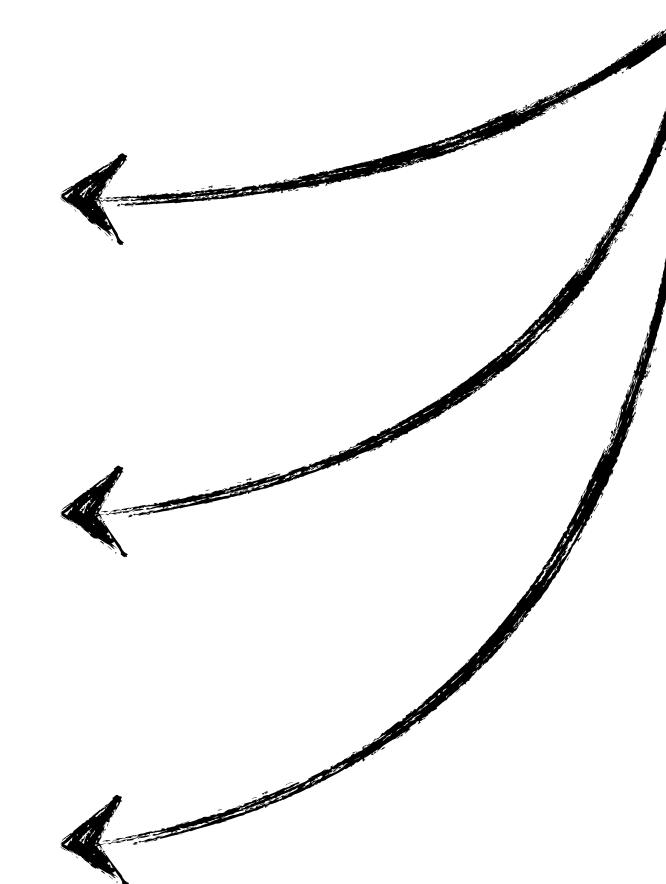
get
 $O(M + \log_T N)$



Bloom
filters

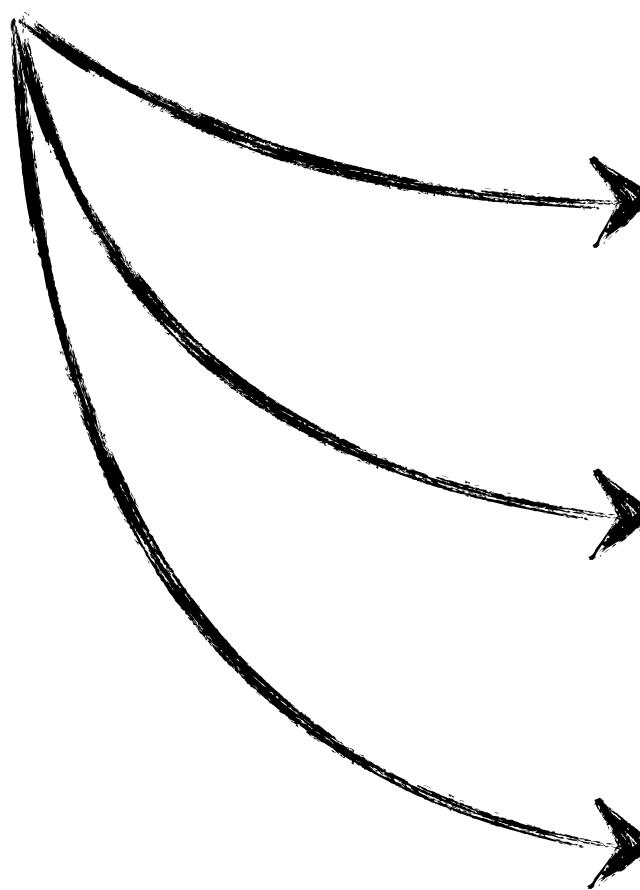


Insert
 $O(T \cdot M \cdot \log_T N)$



Address Using Blocking and SIMD

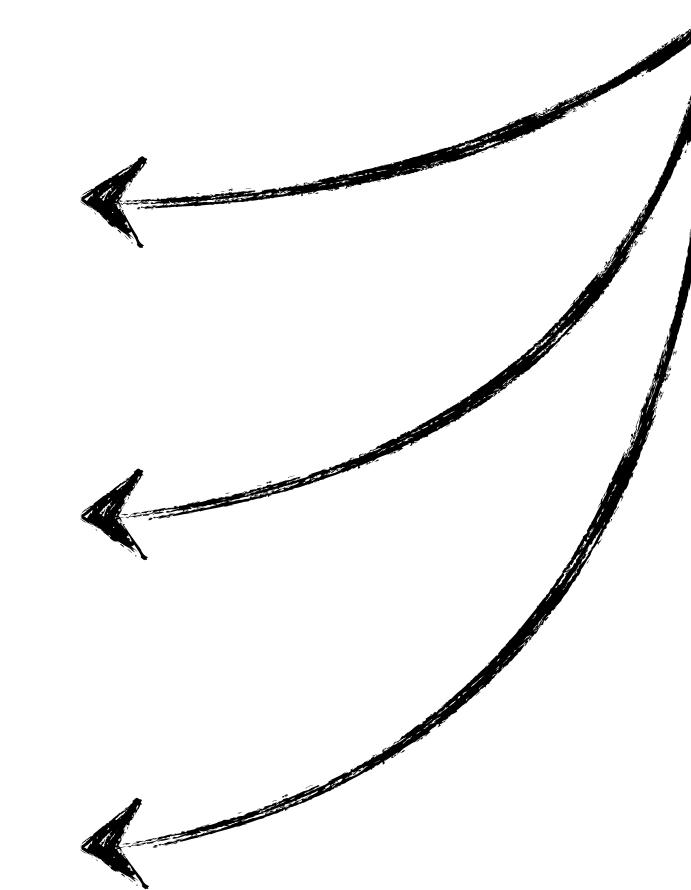
get
 $O(M + \log_T N)$



Bloom
filters

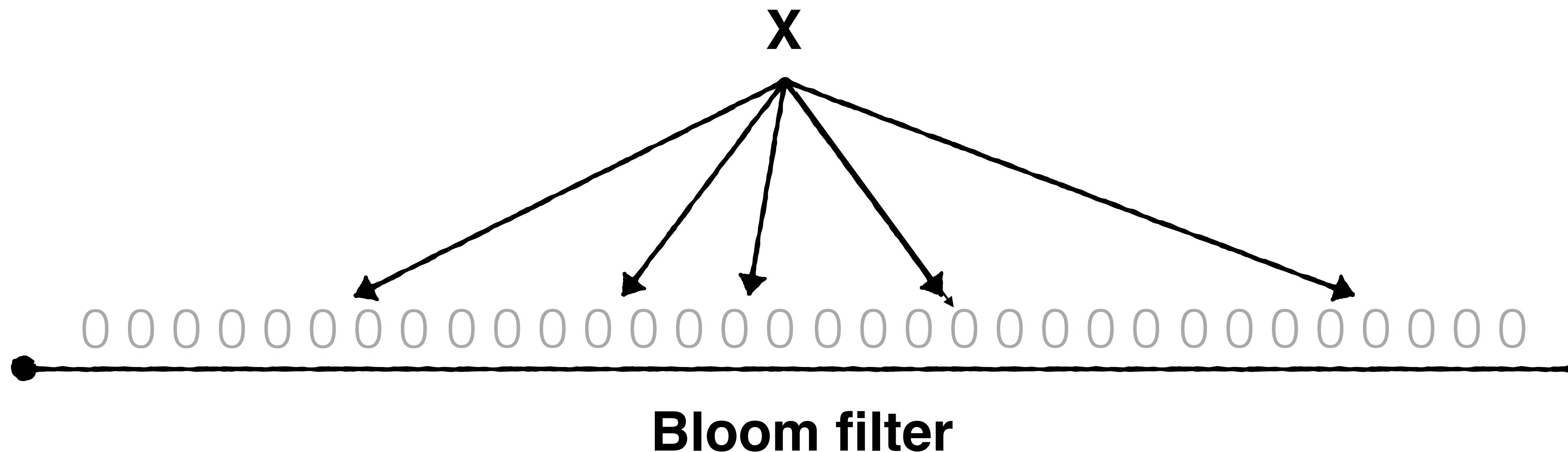


Insert
 $O(T \cdot M \cdot \log_T N)$



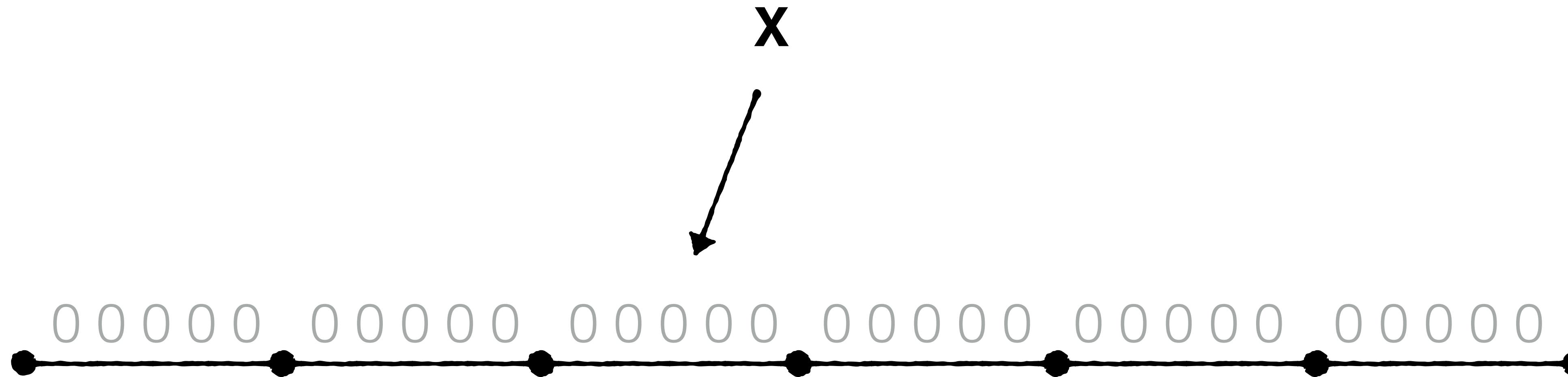
Blocking

PutzeJEA10

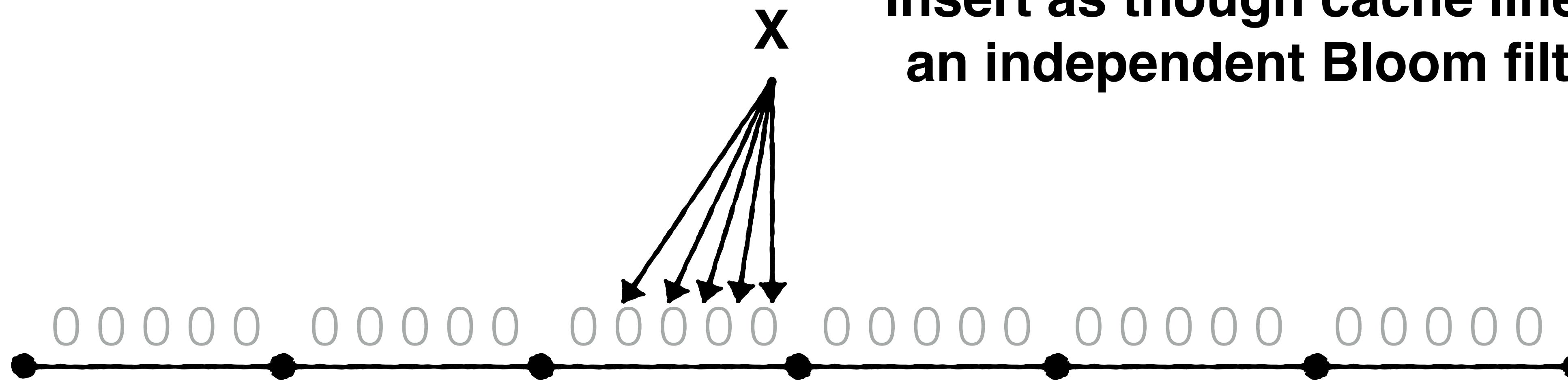


Blocking

Hash to one cache line

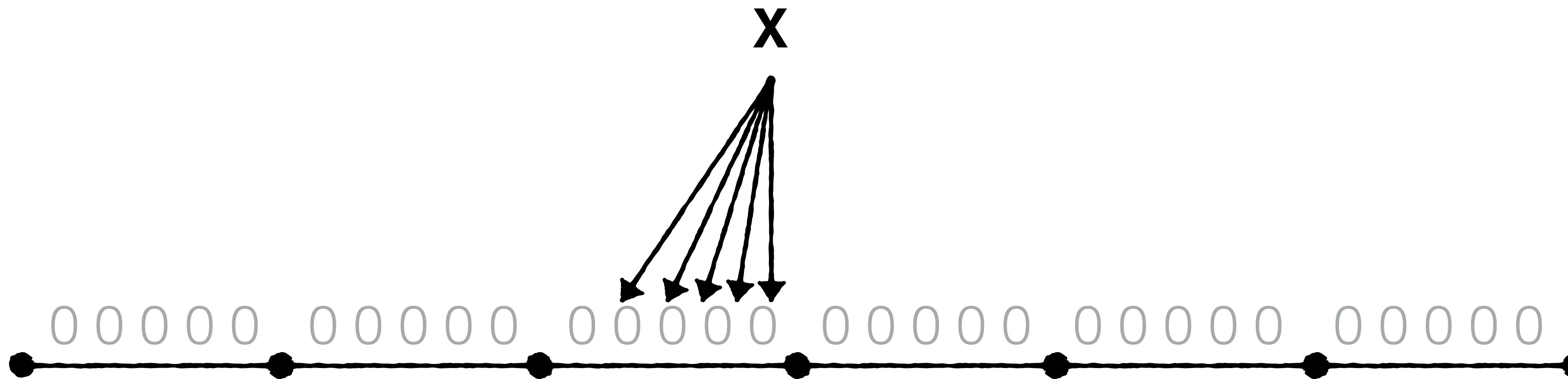


Blocking



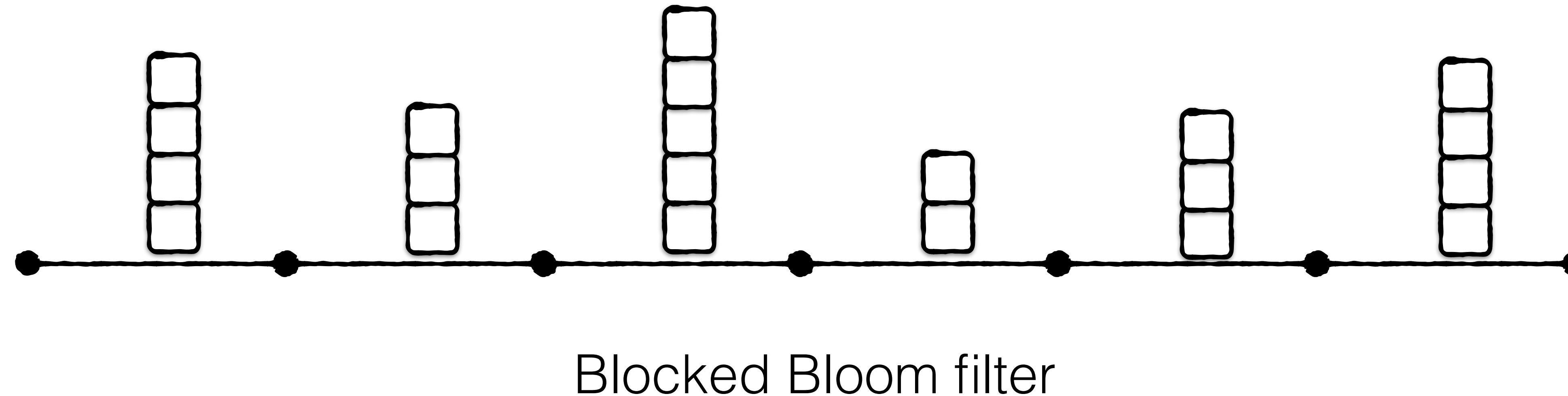
Blocking

Pro: one cache miss per insert/get



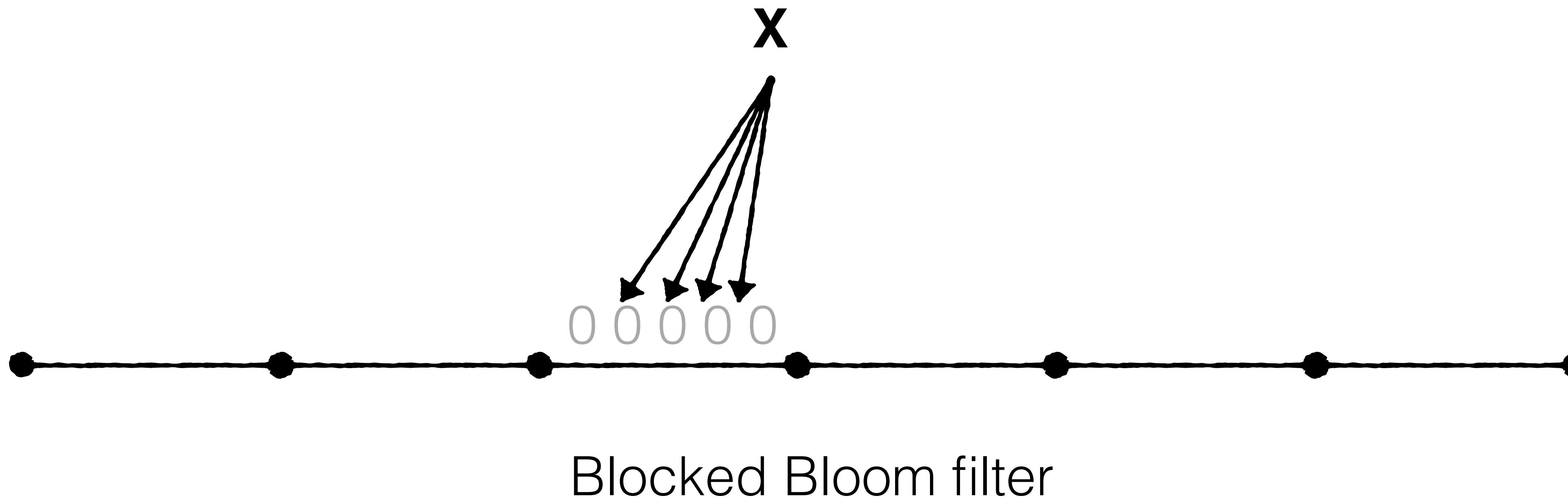
Blocking

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



Blocking

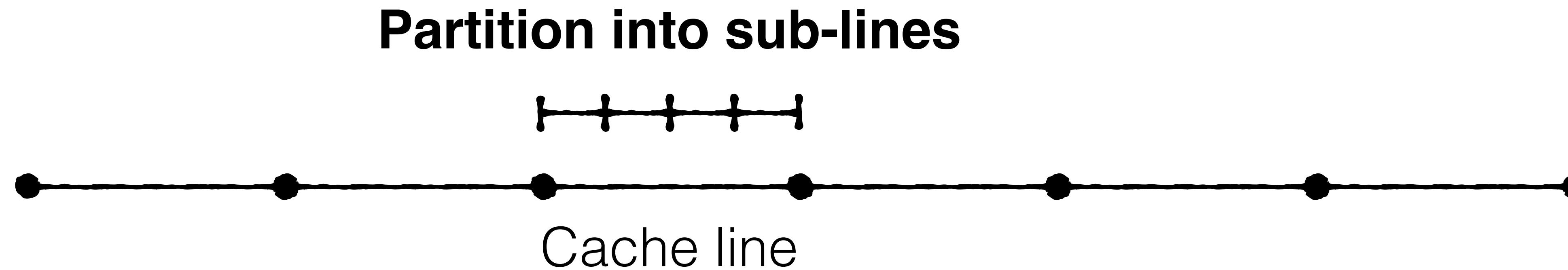
Con 2: still need to compute many hash functions per entry



SIMD

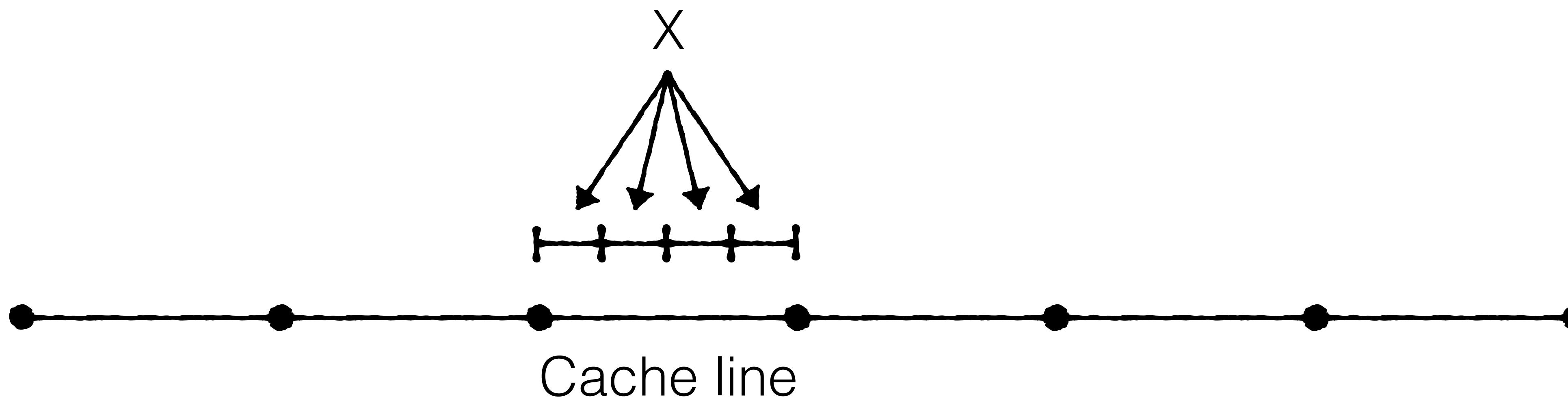
PolychroniouDAMON14

JianyuanTPDS18



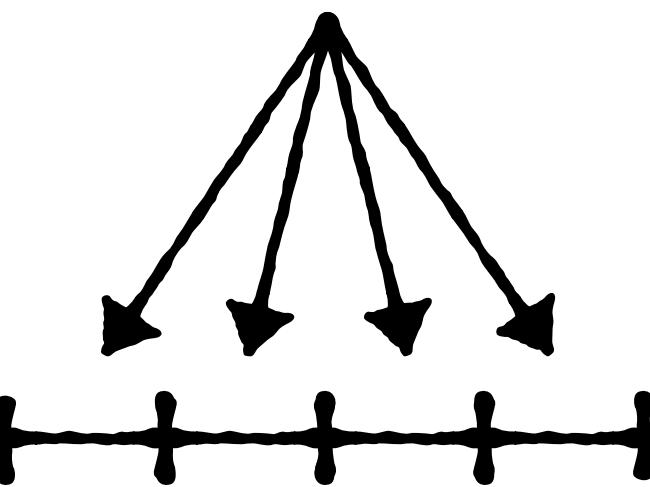
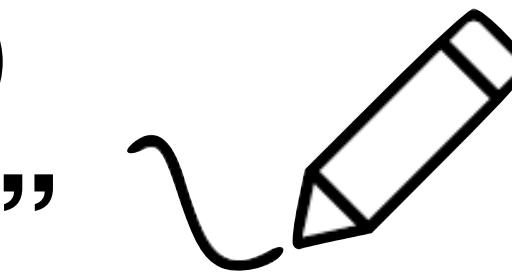
SIMD

Map one hash per sub-line



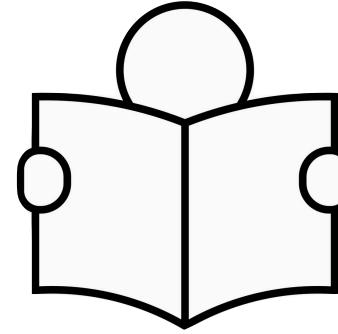
SIMD

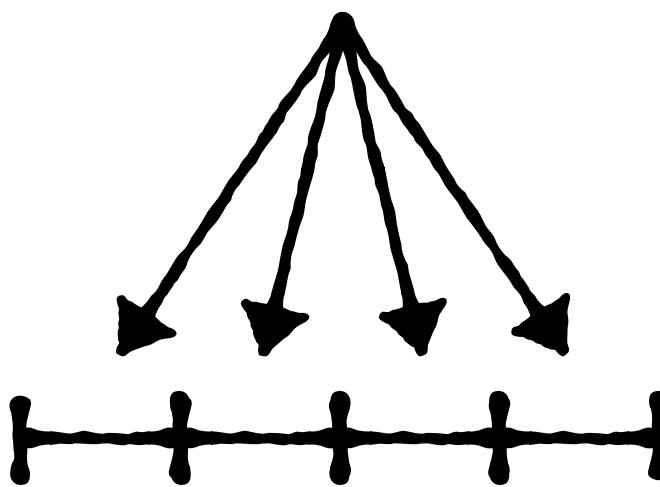
Insert(X)
bitwise “or”



Cache line

SIMD

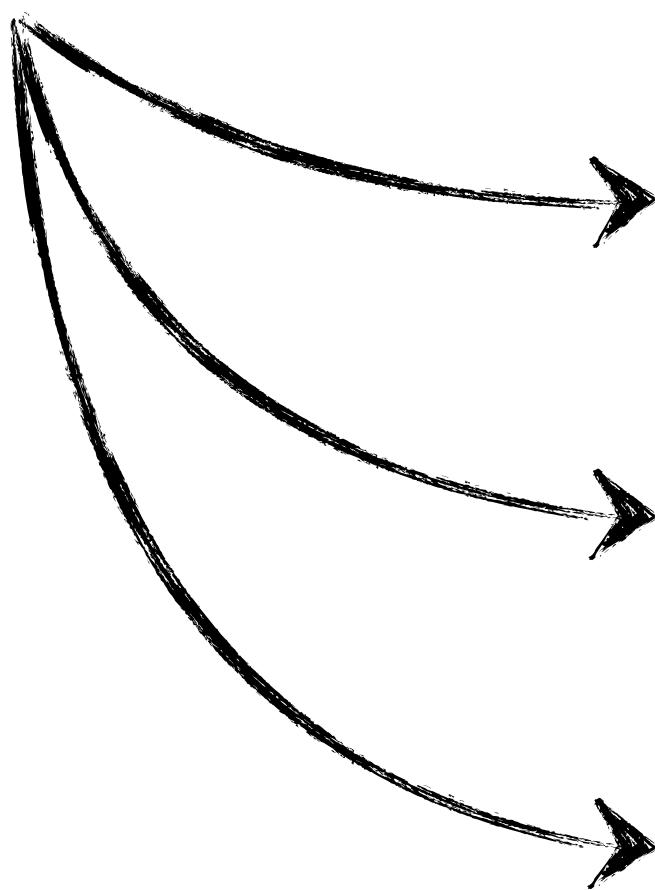
get(X)
bitwise “**and**” 



Cache line

Blocking and SIMD

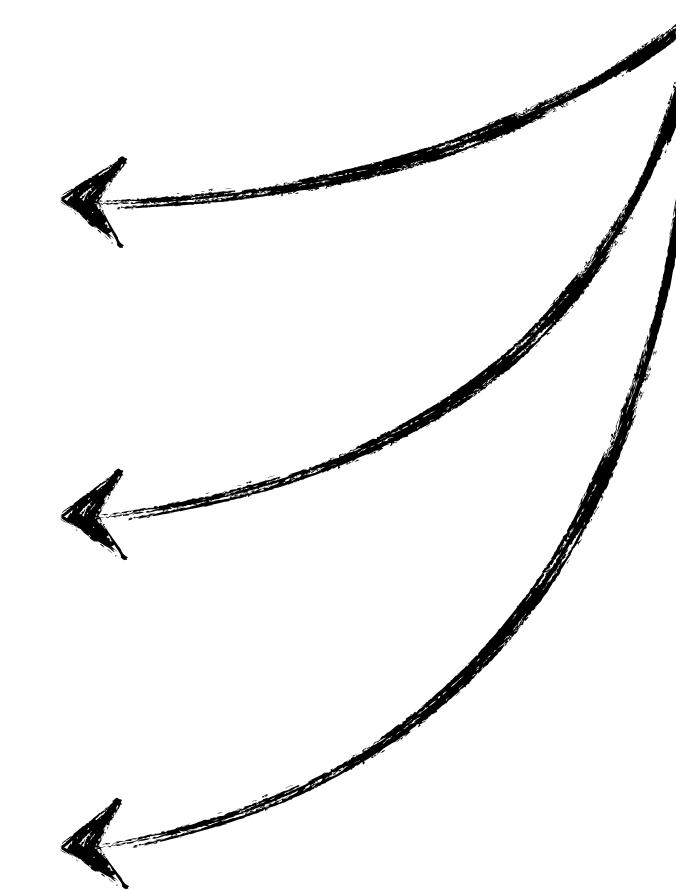
get
 $O(M + \log_T N)$



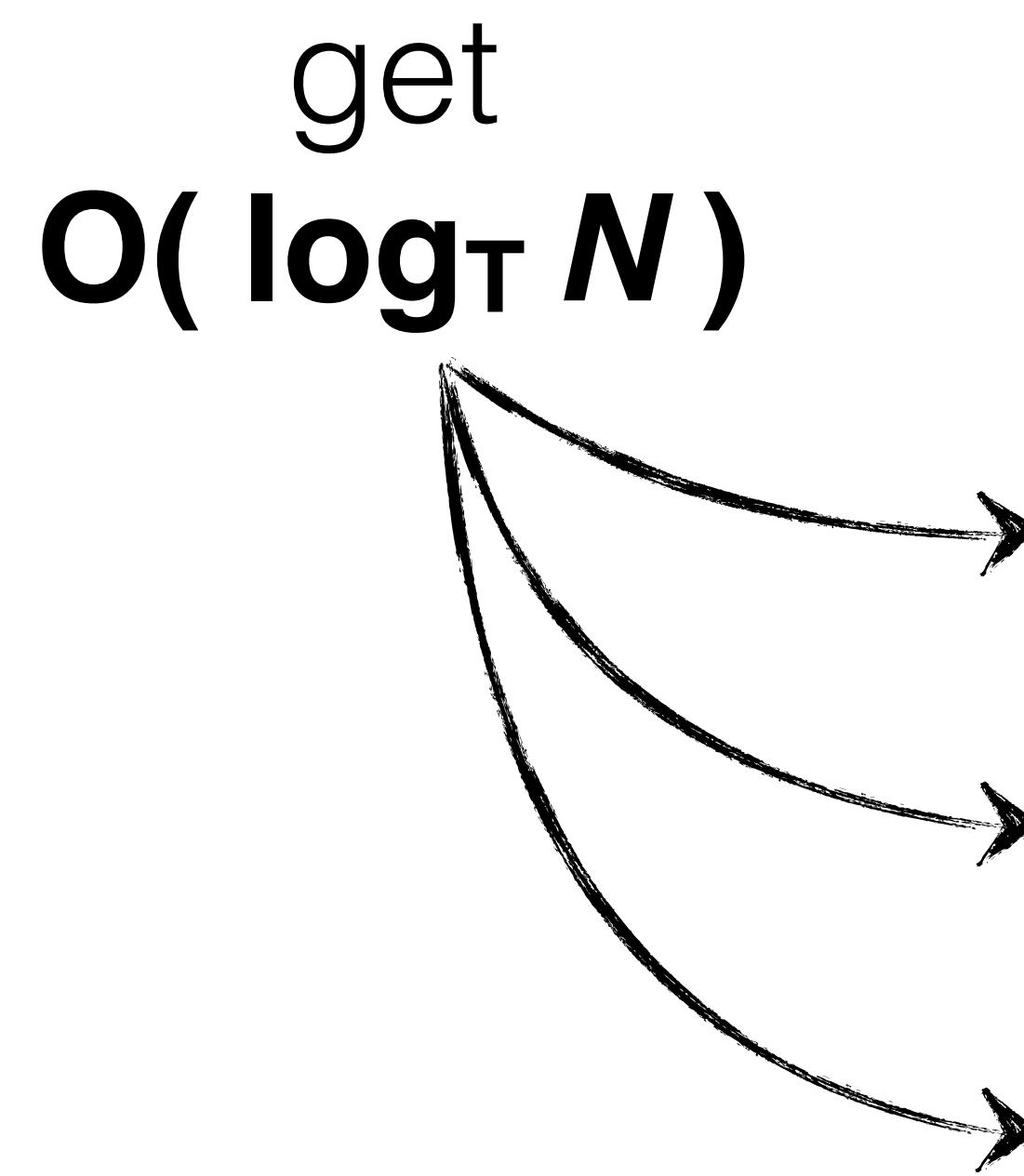
Bloom
filters



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



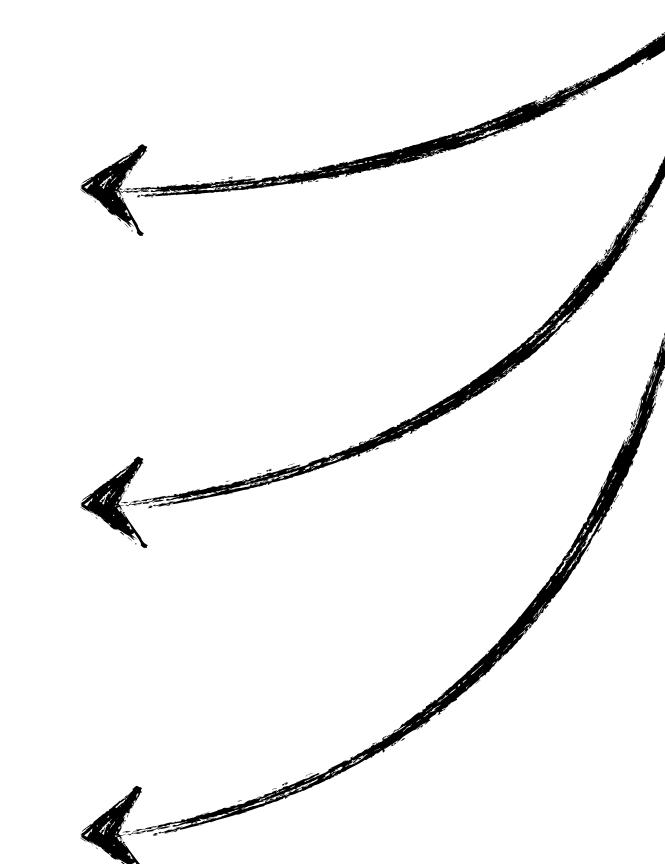
Blocking and SIMD



Bloom
filters



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



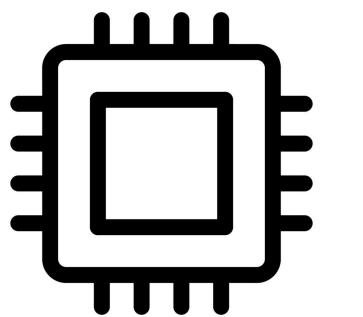
The diagram consists of three curved arrows drawn with a black marker. They start at different points on the right and curve downwards and to the left, ending at a single point. This visual metaphor represents the process of inserting multiple items into a data structure.

5 fronts

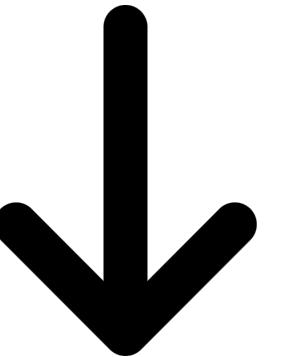
Holistic
Tuning



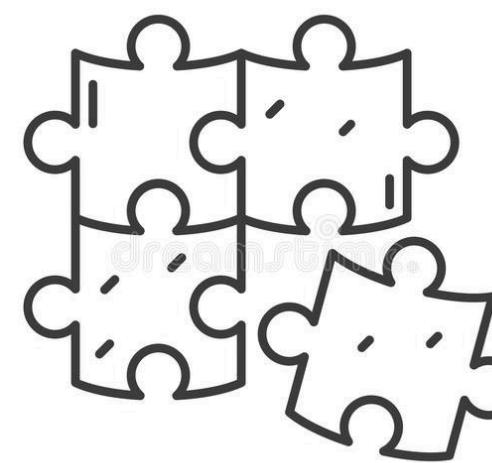
CPU



**Improving
Constants**



Unification



Range



Improving Constants

Bloom



False
positive rate

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

Ideal



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

False
positive rate

Bloom



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



Ideal



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

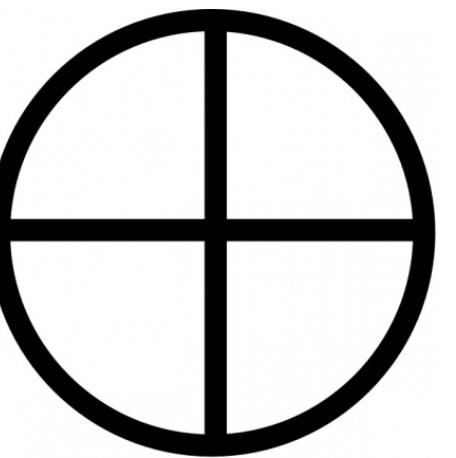
Can we improve this?

Bloom



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

XOR



Ideal

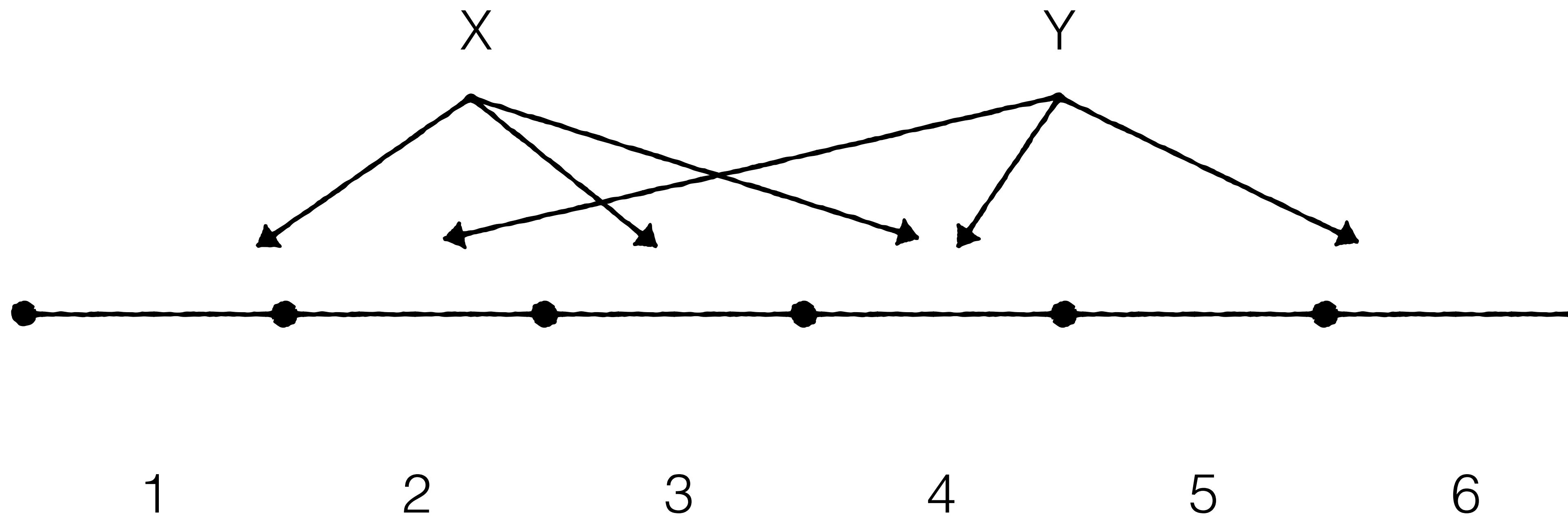


$$\approx 2^{-M}$$



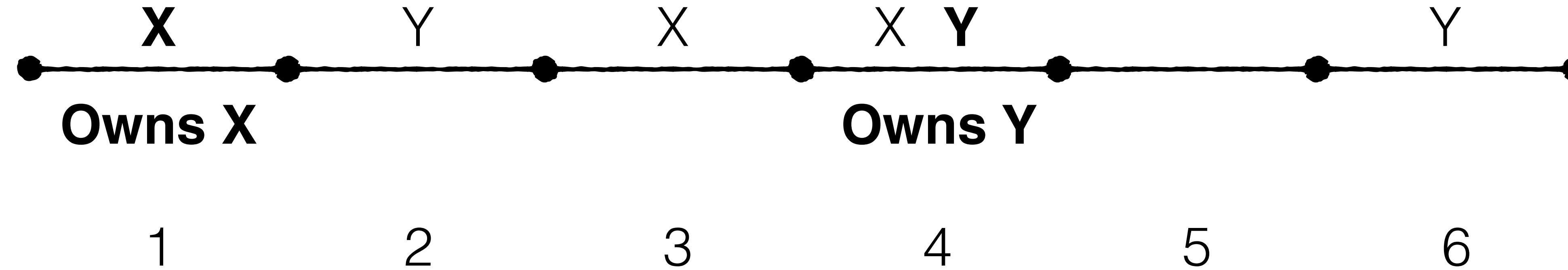
XOR Filter

Hash each entry to three buckets



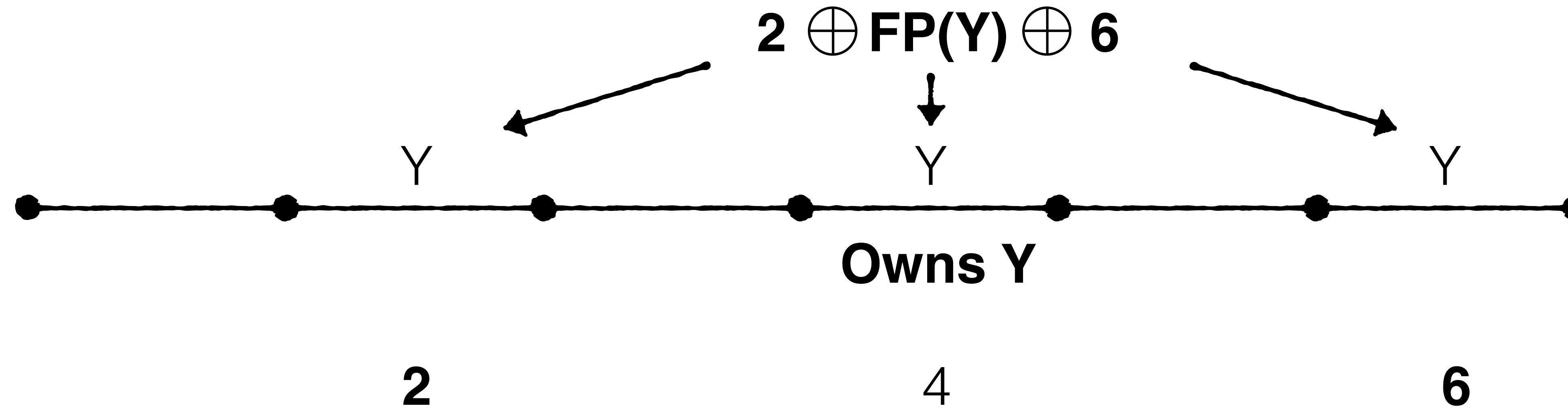
XOR Filter

Assign one bucket to own each entry



XOR Filter

Each bucket stores XOR of fingerprint and other two buckets

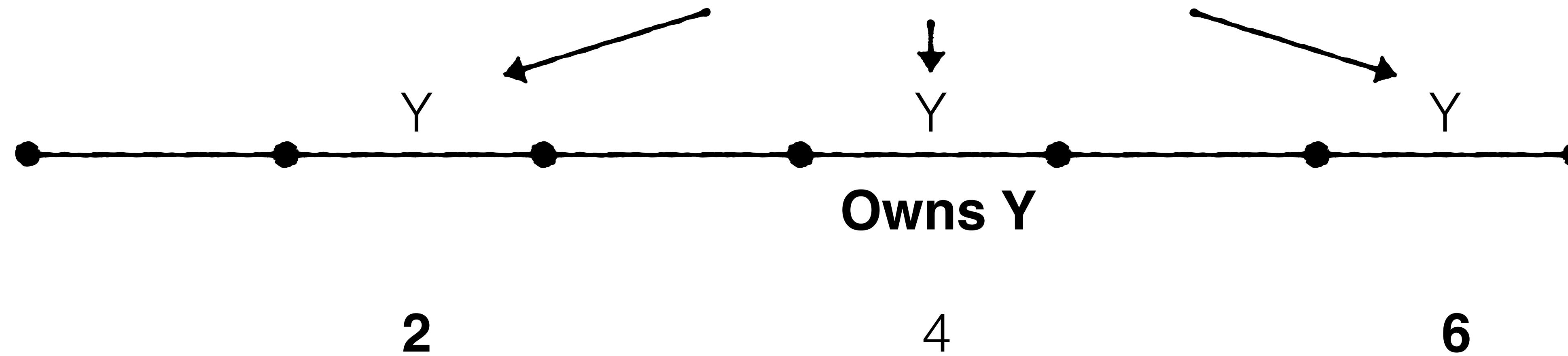


XOR Filter

During queries, recover fingerprints by xorring three buckets

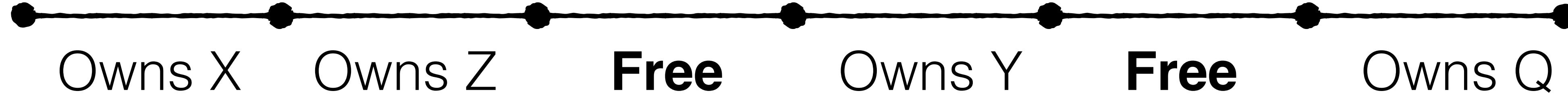
get(Y) returns true if

$$\text{FP}(Y) = 2 \oplus 4 \oplus 6$$



XOR Filter

free space ensures each bucket can own one entry

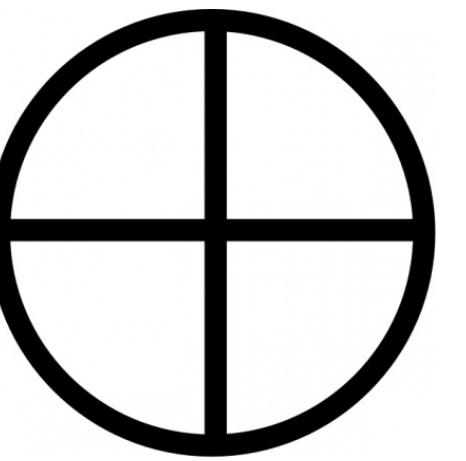


Bloom



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

XOR



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

Idealized



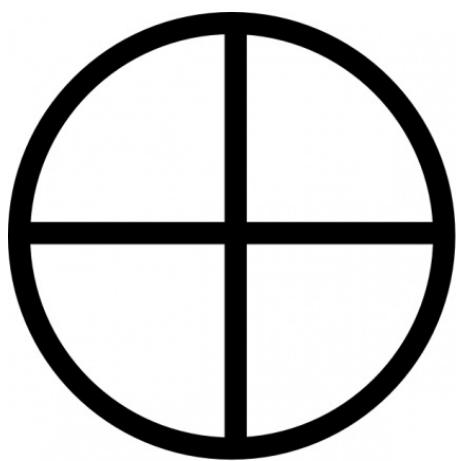
$$\approx 2^{-M}$$

Bloom



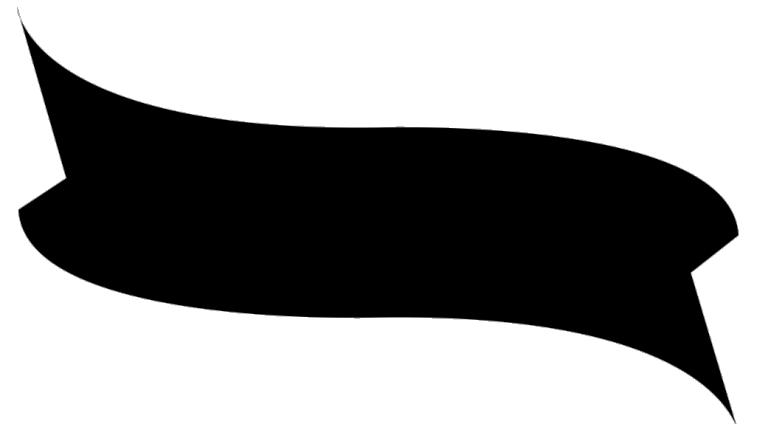
$$\approx 2^{-M} \cdot 0.69$$

XOR



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

Ribbon



$$\approx 2^{-M} \cdot 0.92$$

Idealized



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

Denser XOR filter



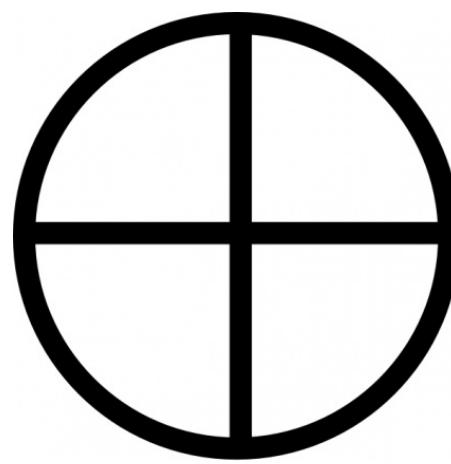
DillingerSEA22

Bloom



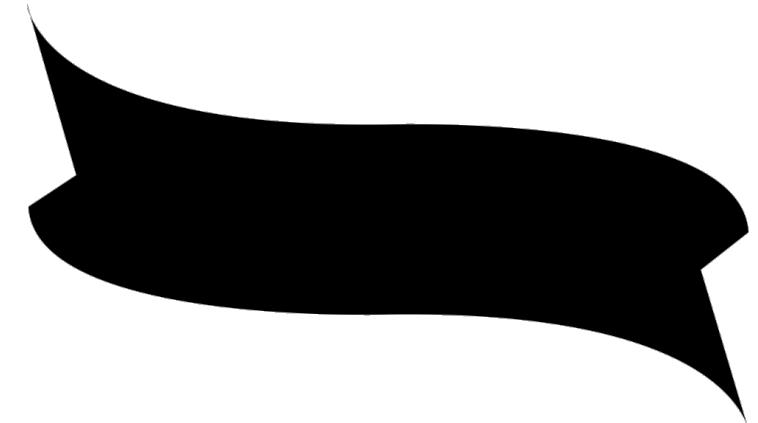
$\approx 2^{-M} \cdot 0.69$

XOR



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

Ribbon



$\approx 2^{-M} \cdot 0.92$

Idealized



$\approx 2^{-M}$

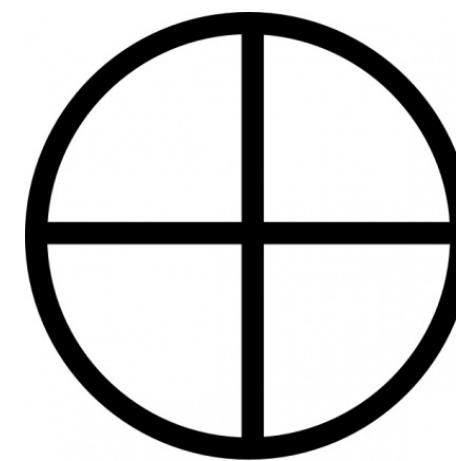
Denser XOR filter

In RocksDB since 2020

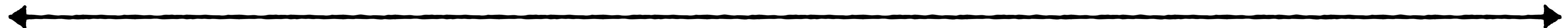
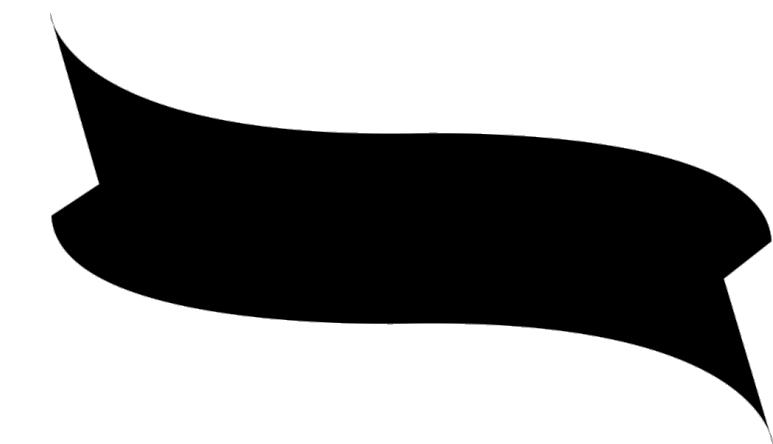
Bloom



XOR



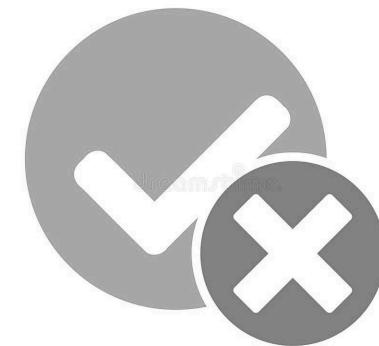
Ribbon



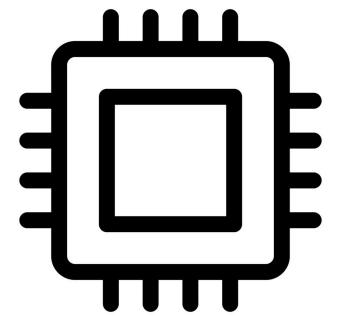
Lower CPU

**Lower false
positive rate**

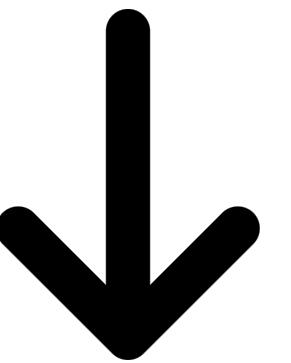
Holistic
Tuning



CPU

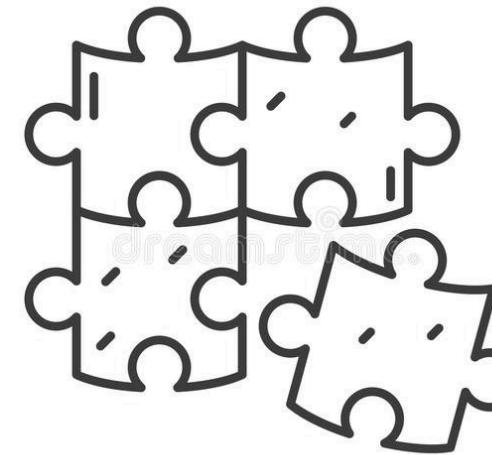


Improving
Constants



5 fronts

Unification



Range

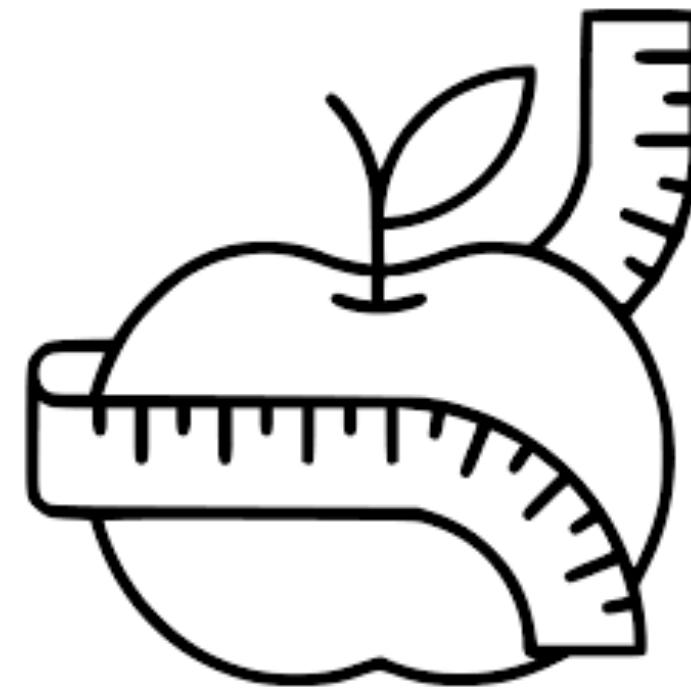


Unification



Chucky

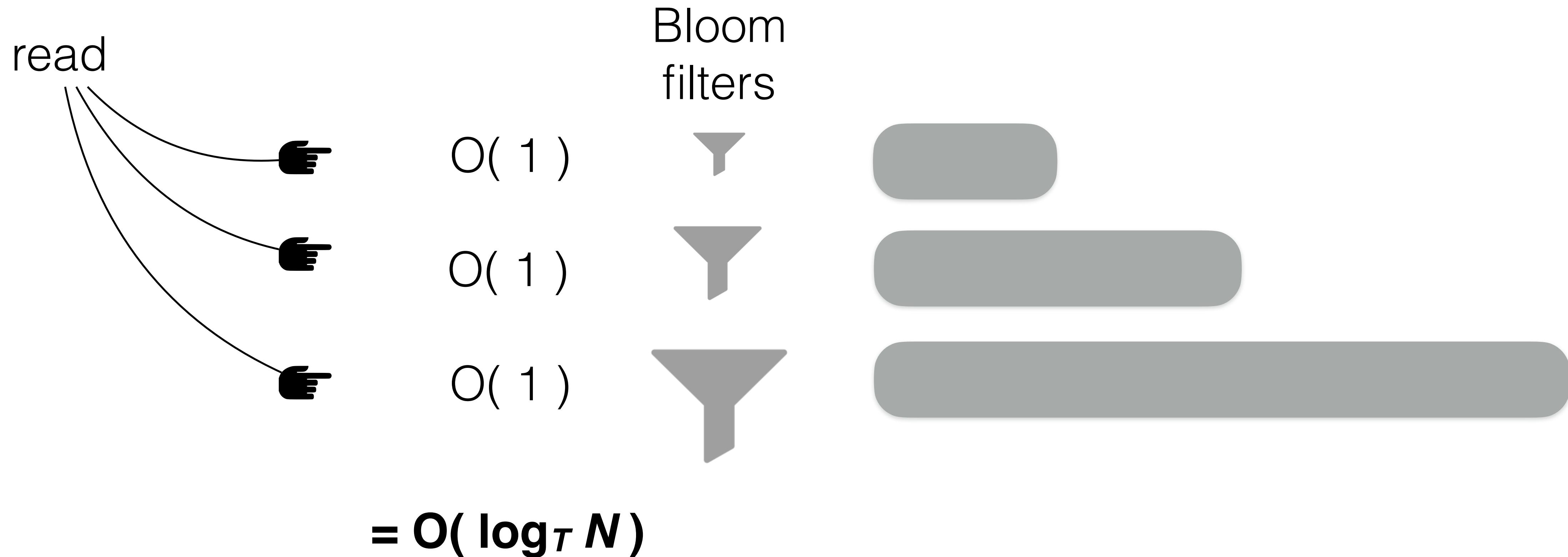
DayanSIGMOD21



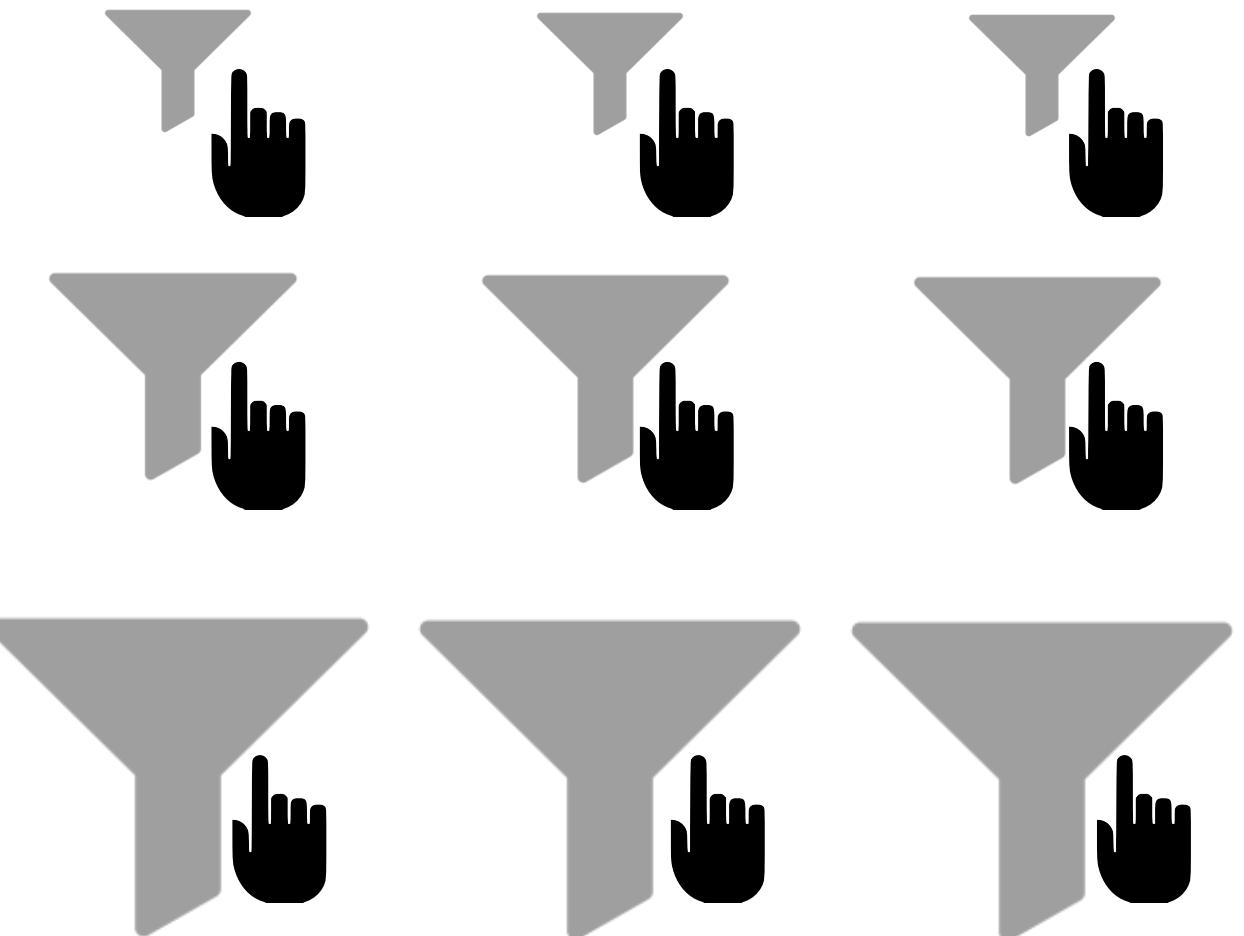
SlimDB

RenVLDB17

Unification



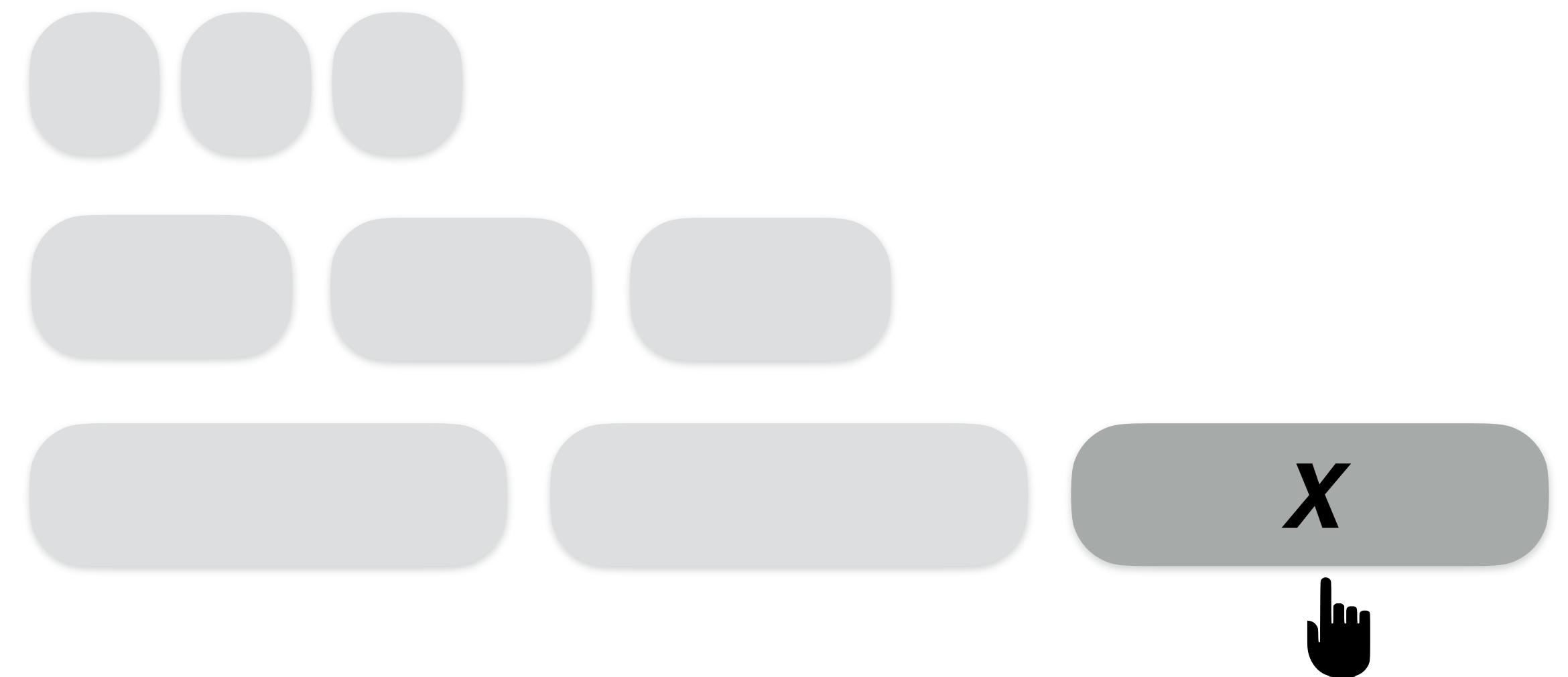
$O(T)$

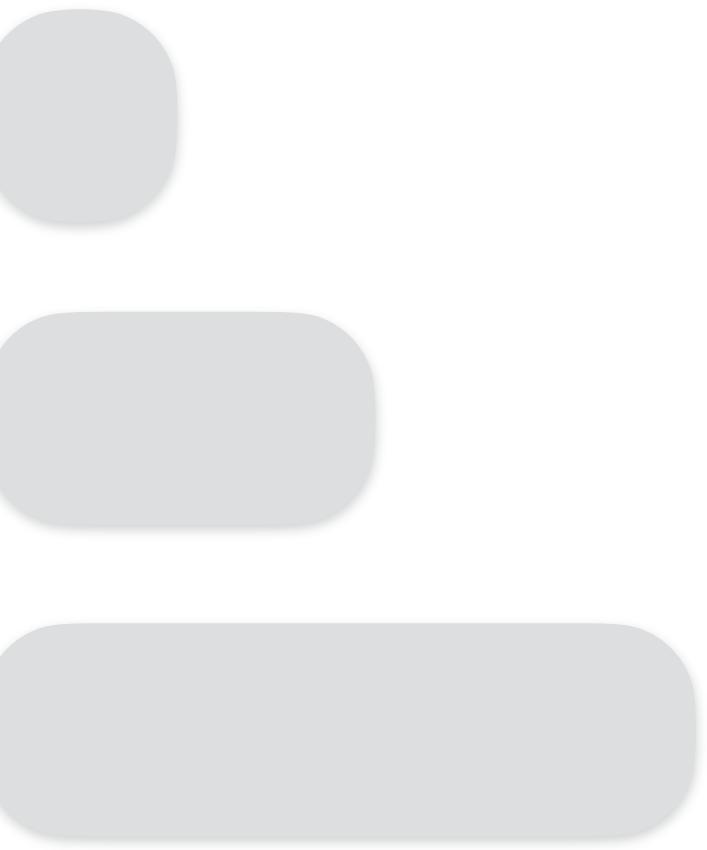
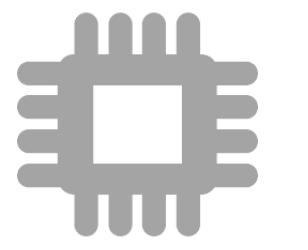
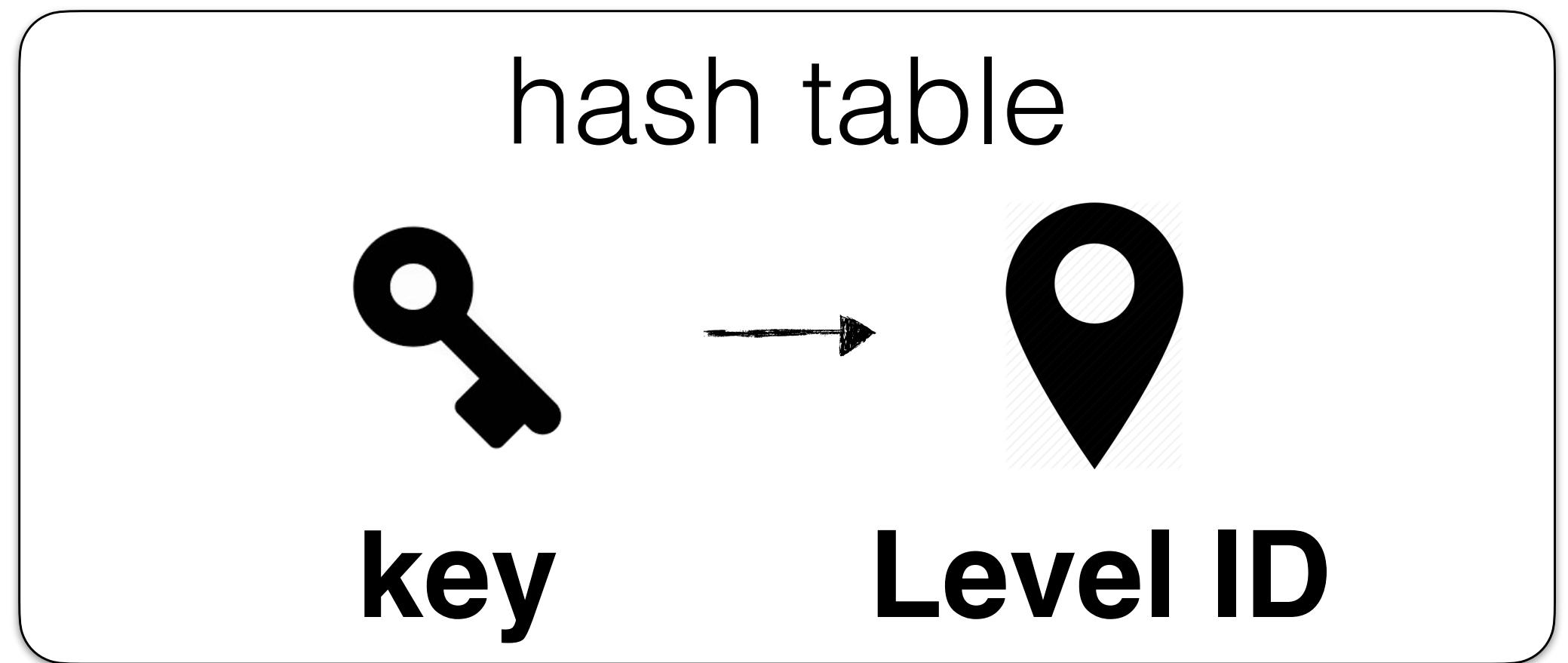


$O(T)$

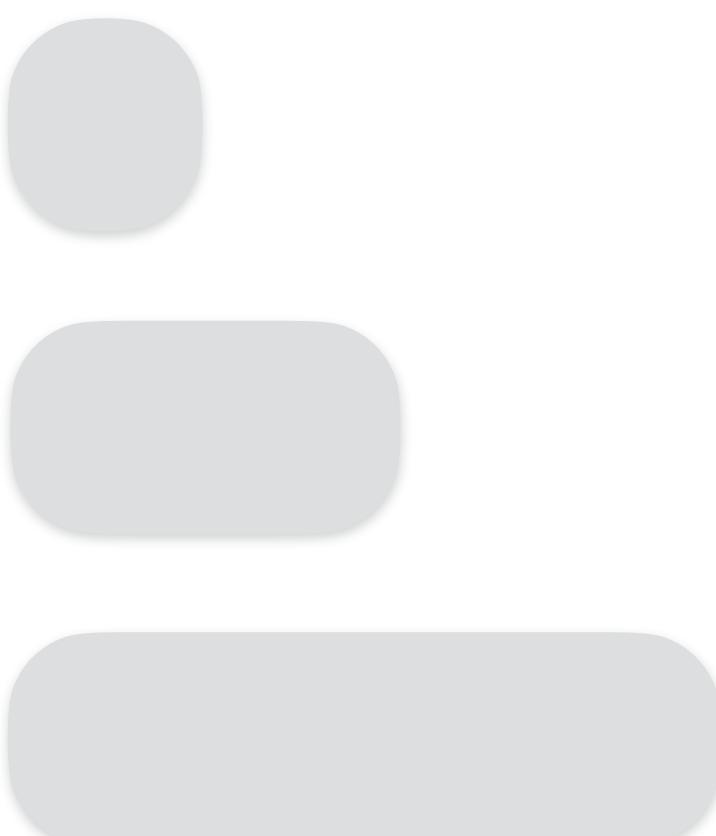
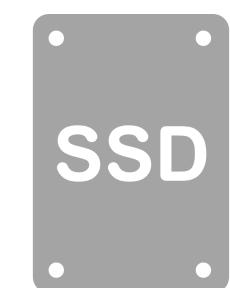
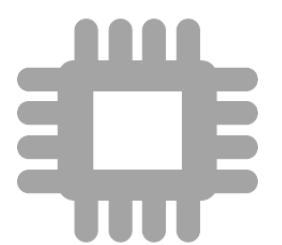
$O(T)$

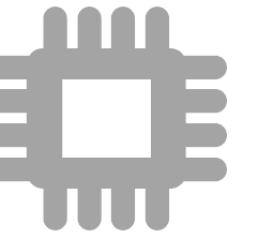
$= O(T \cdot \log_T N)$



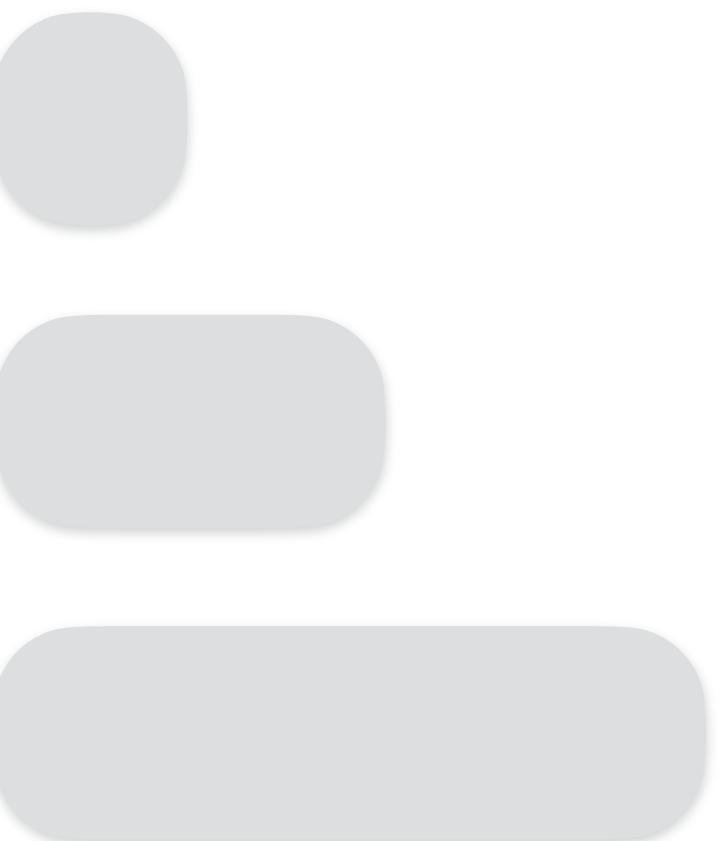


hash() =

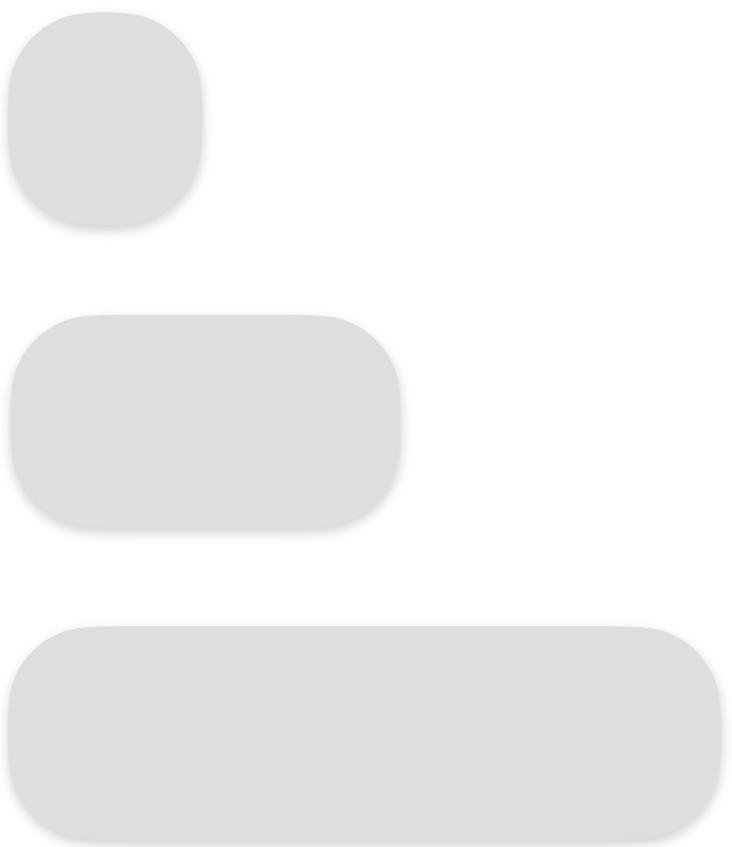
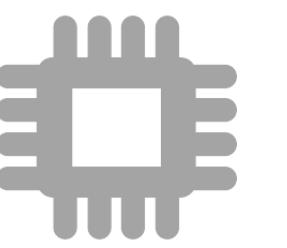
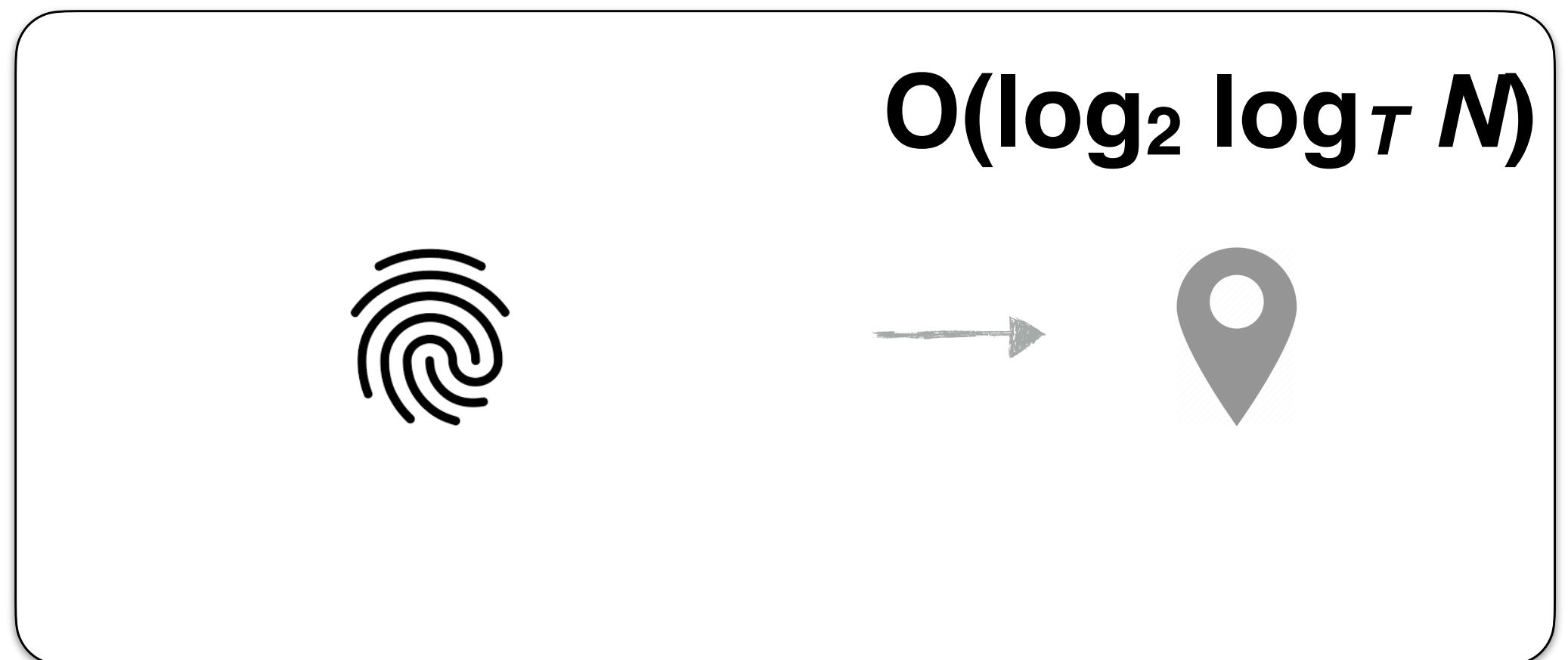




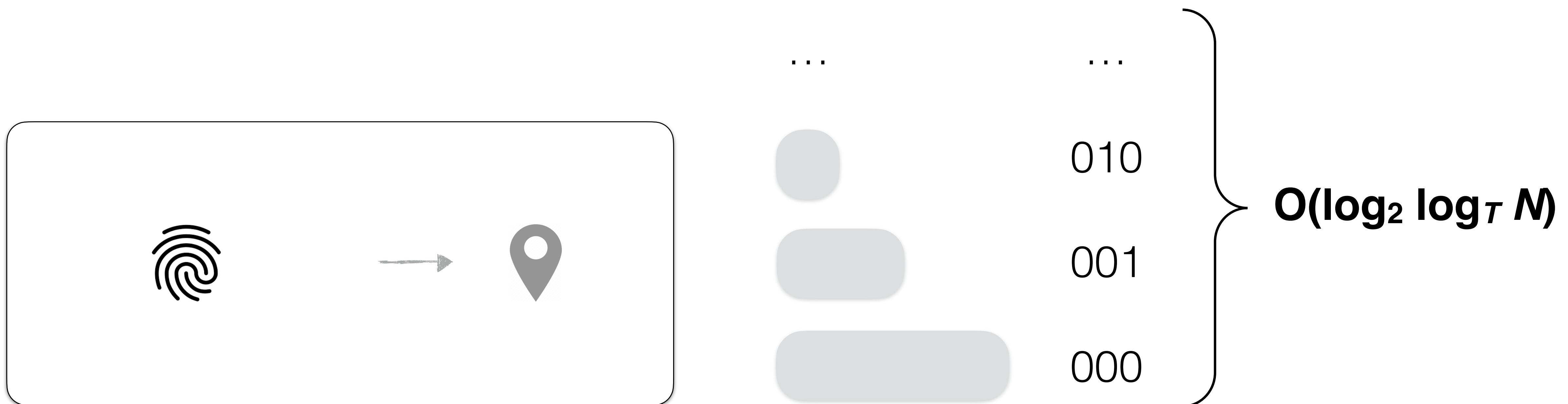
hash() =



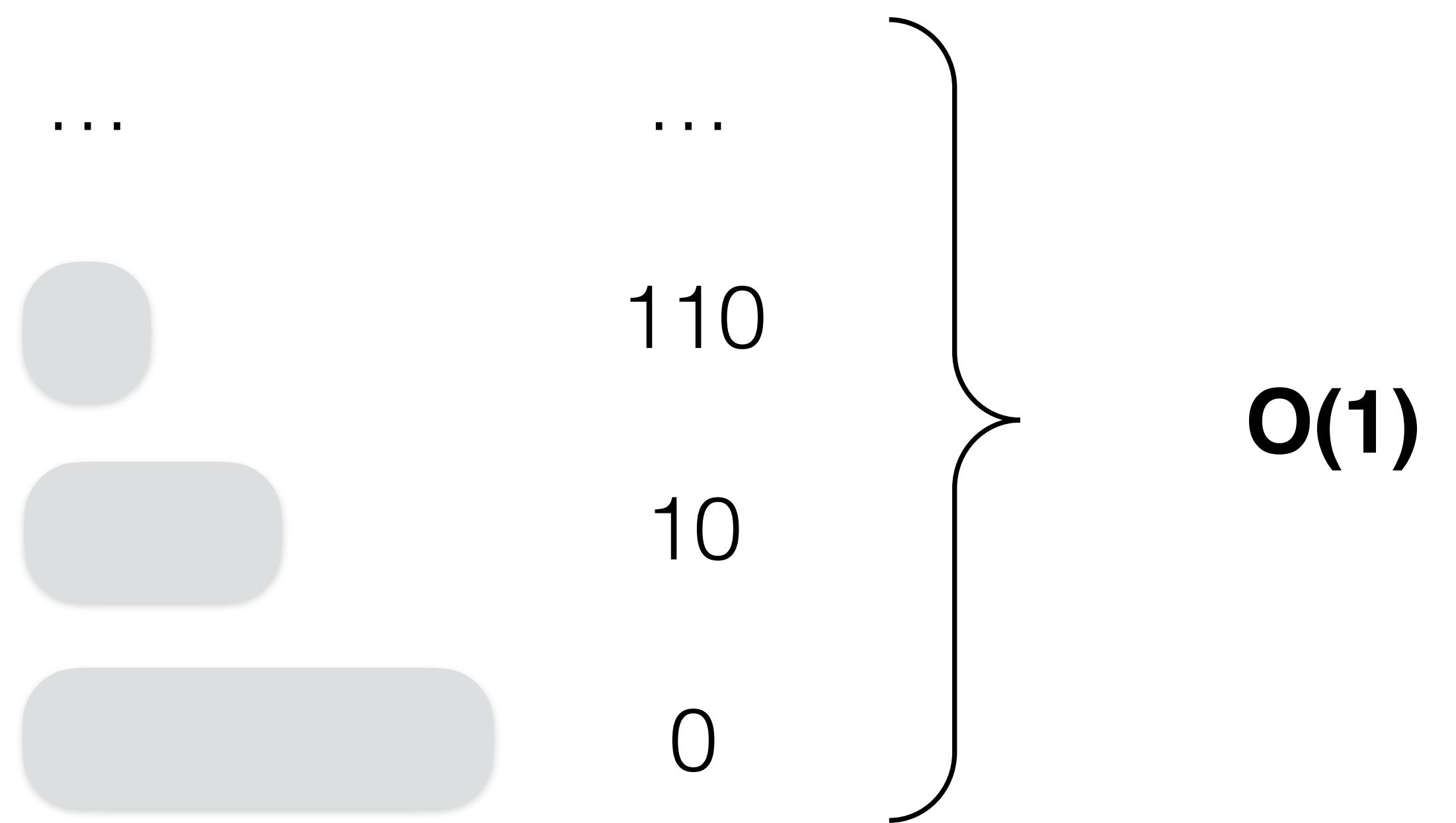
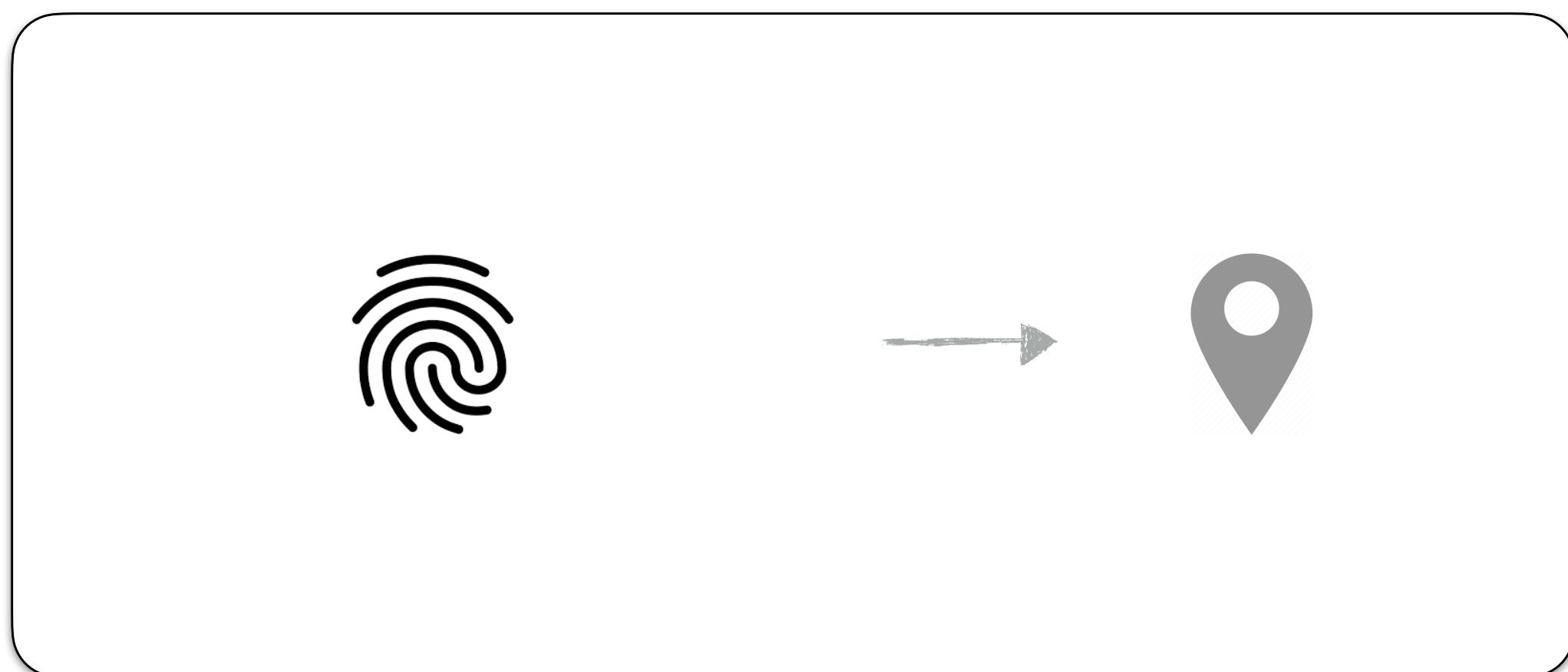
cuckoo filter



Binary encoding



e.g., unary encoding





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

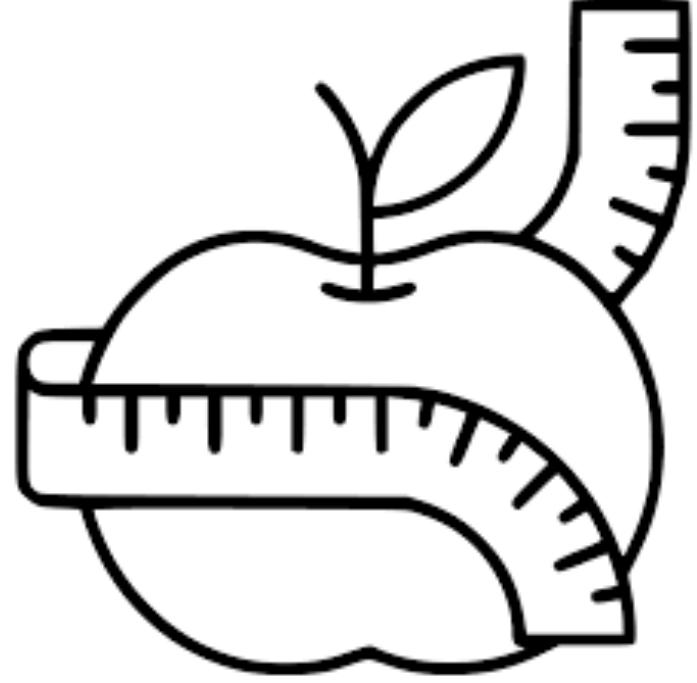
$O(\log_T N)$

$O(1)$

Insert CPU

$O(\textcolor{red}{T} \cdot \log_T N)$

$O(\textcolor{green}{\log}_T N)$



Chucky

Get I/O

$O(1 + 2^{-M+3})$

SlimDB

$O(1)$

Memory

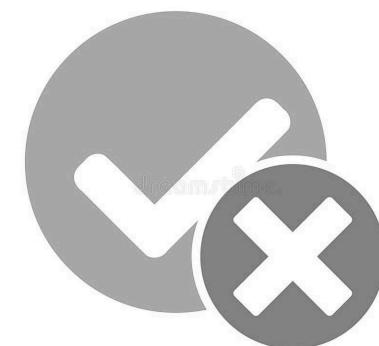
M

$M + 2^M \cdot \log_2(N)$

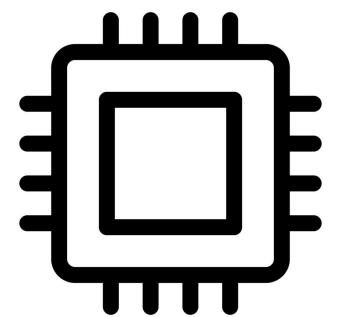
**Keeps full key in memory
whenever fingerprints collide**

5 fronts

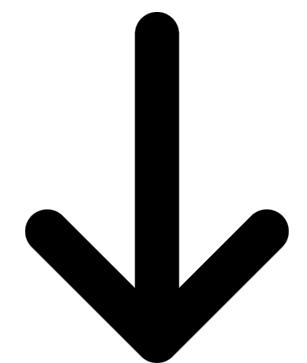
Holistic
Tuning



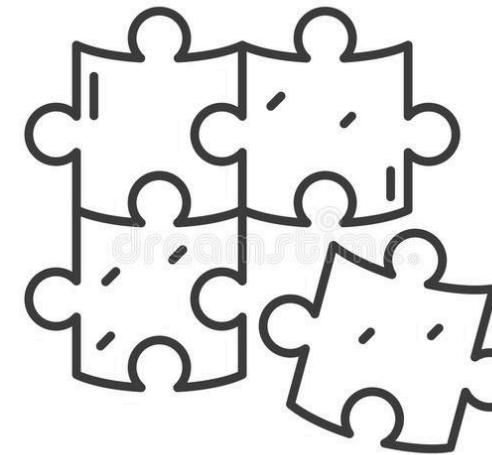
CPU



Improving
Constants



Unification



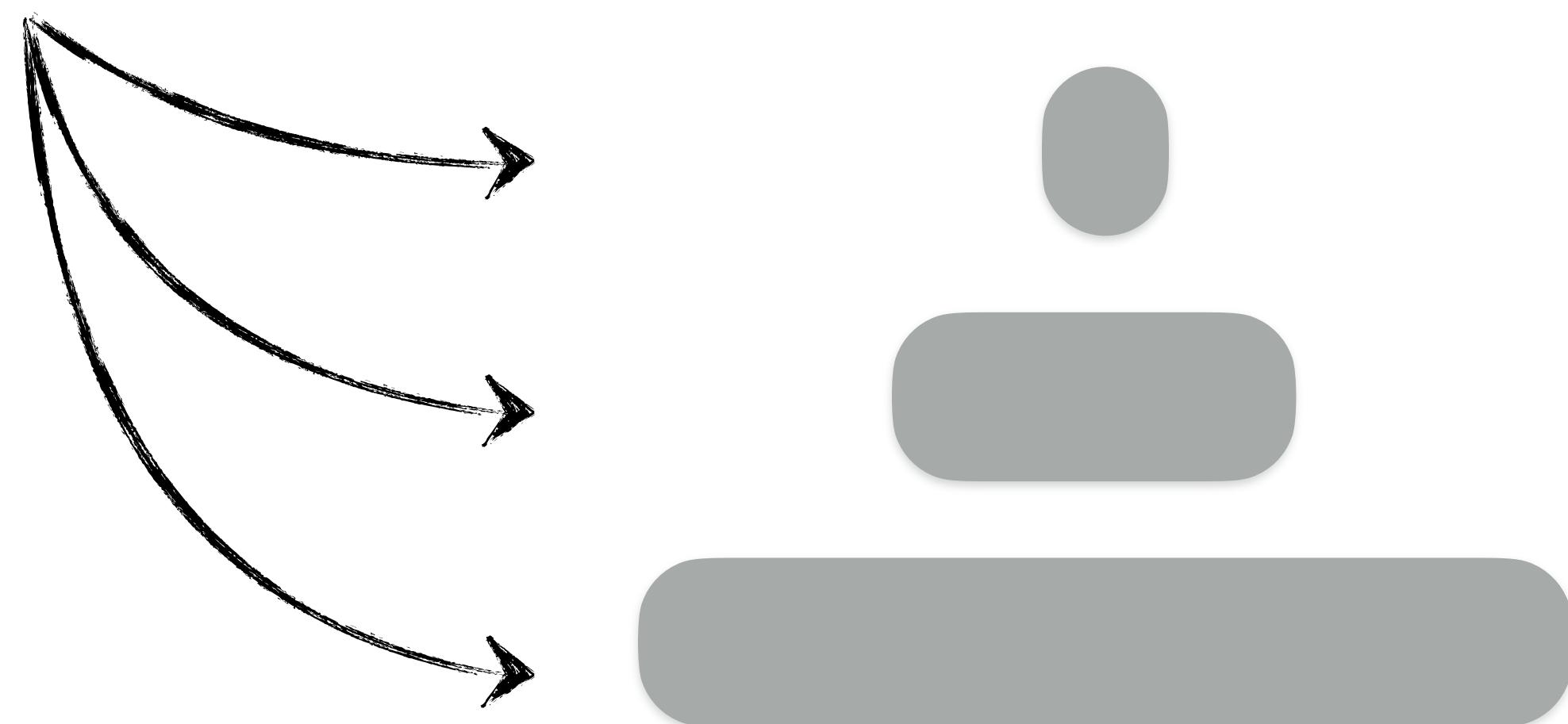
Range



Range Filtering

Traditional filters do not support ranges

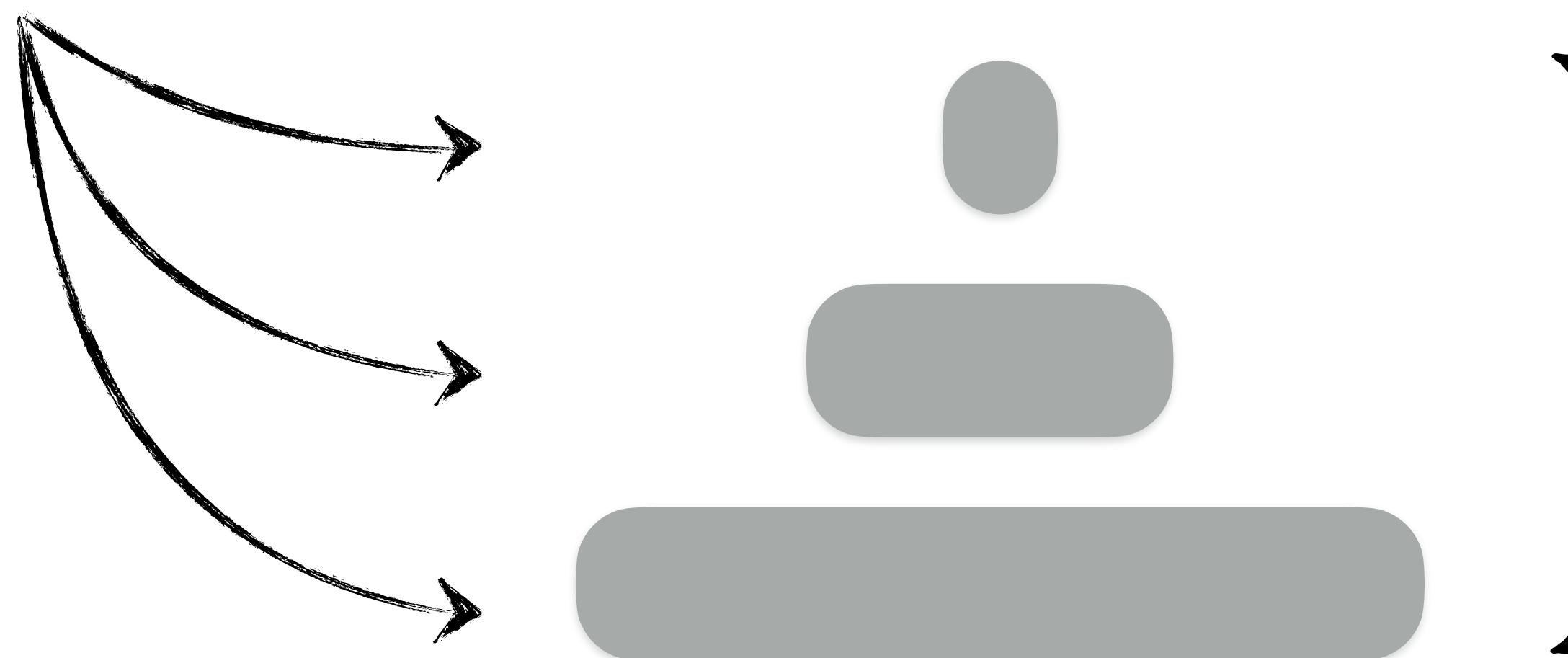
get(x, y)



Range Filtering

Traditional filters do not support ranges

`get(x, y)`



cost: $O(\log_T N)$

Range Filters



Prefix Filter

RocksDB20



Surf

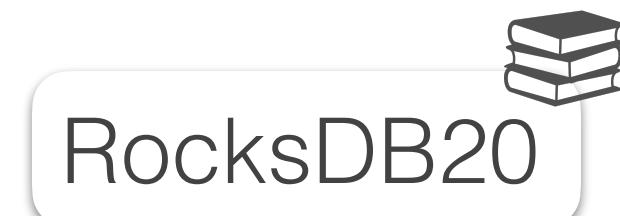
ZhangSIGMOD18



Rosetta

LouSIGMOD21

Prefix Filter



Users define prefix extraction method

Prefix Filter

Users define prefix extraction method

Country code



USA1234

CAN9876

Prefix Filter

Users define prefix extraction method

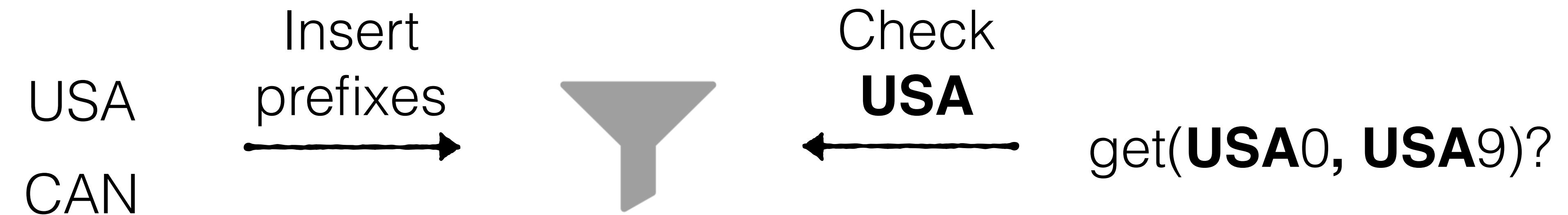
USA
CAN

**Insert
prefixes**



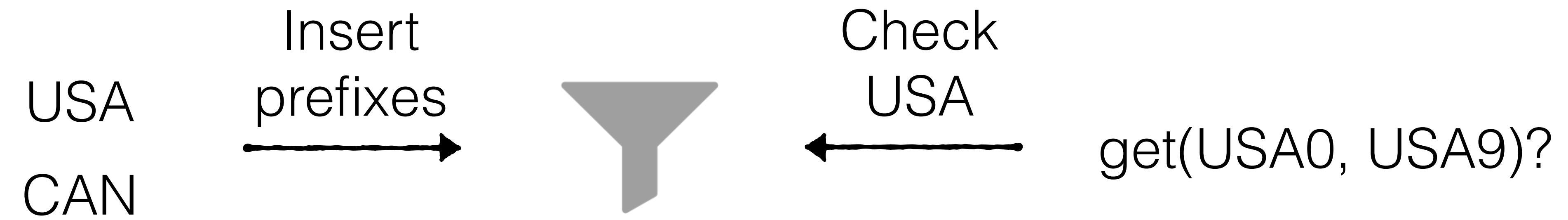
Prefix Filter

Users define prefix extraction method



Prefix Filter

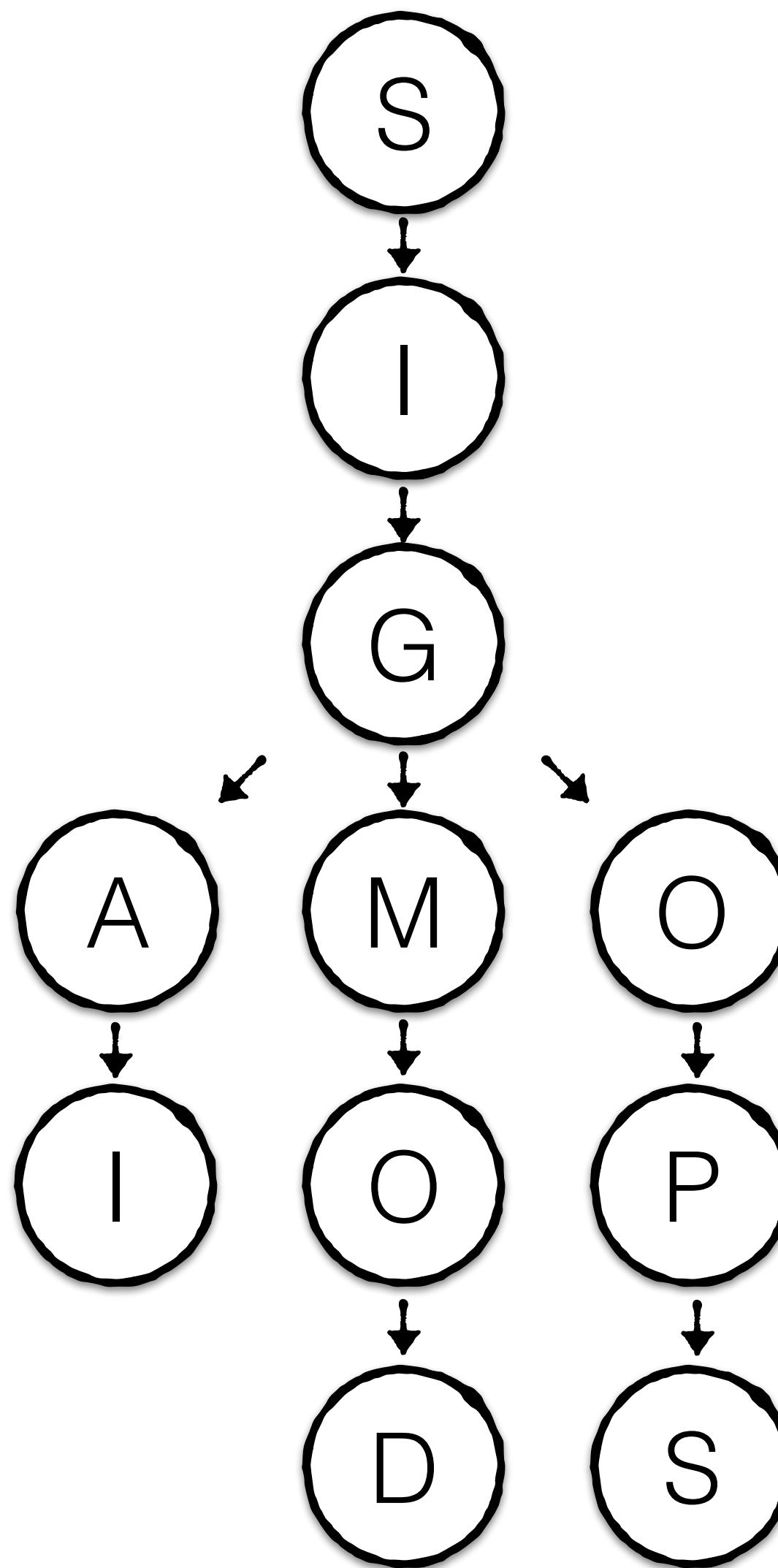
Users define prefix extraction method



Non-generic and requires API extension

Surf

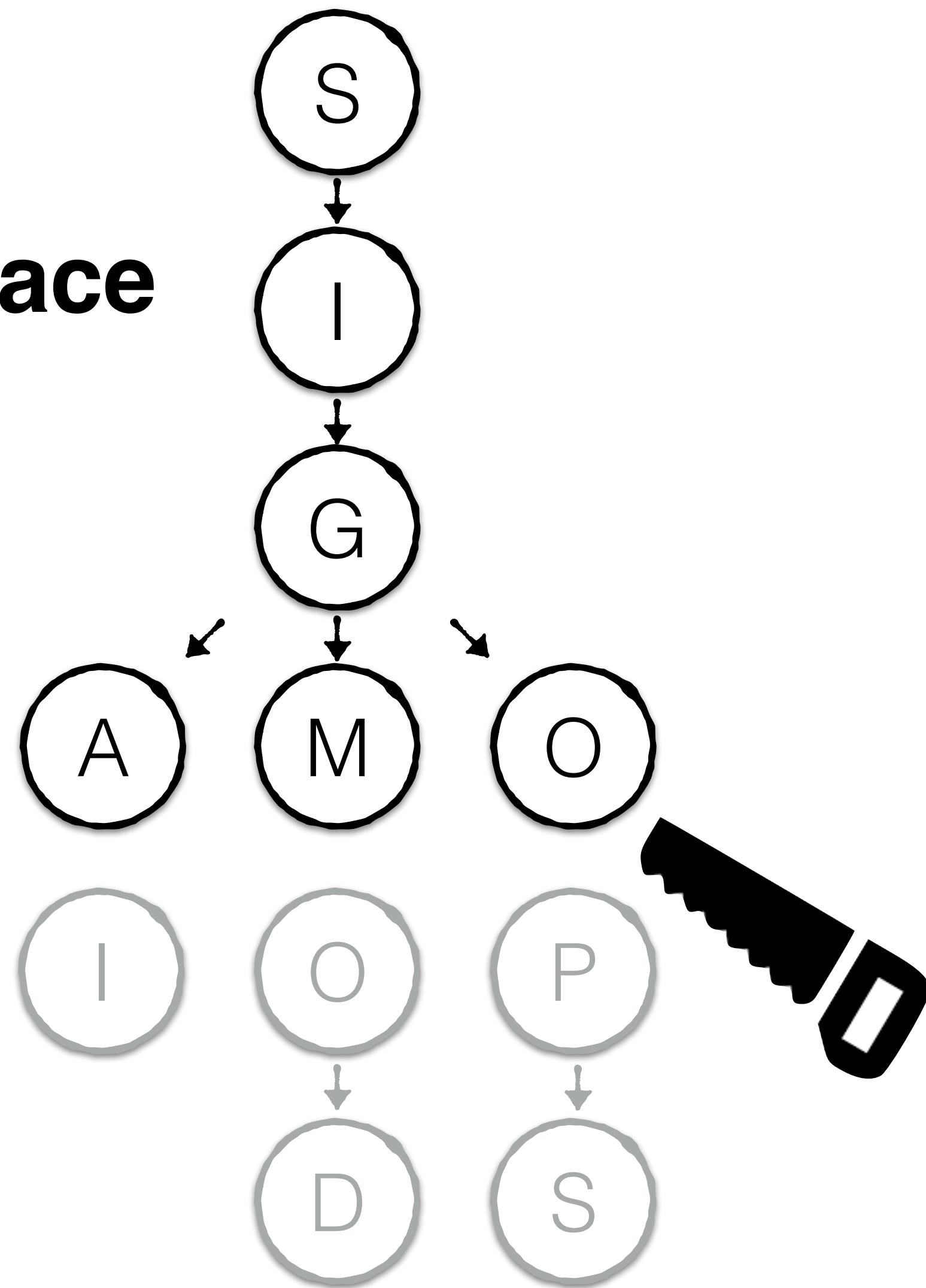
A trie of all keys



Surf

A trie of all keys

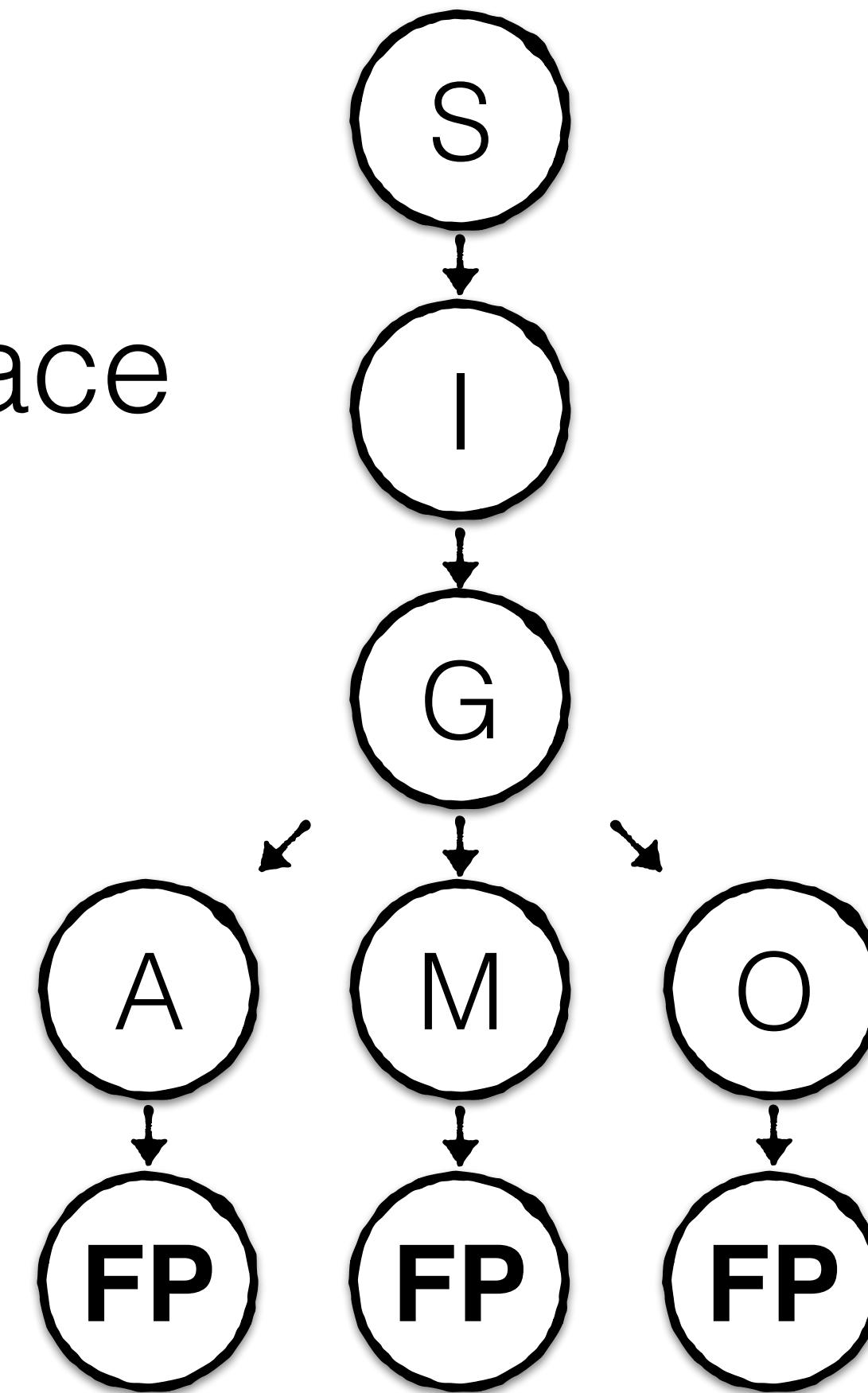
Truncated to reduce space



Surf

A trie of all keys

Truncated to reduce space



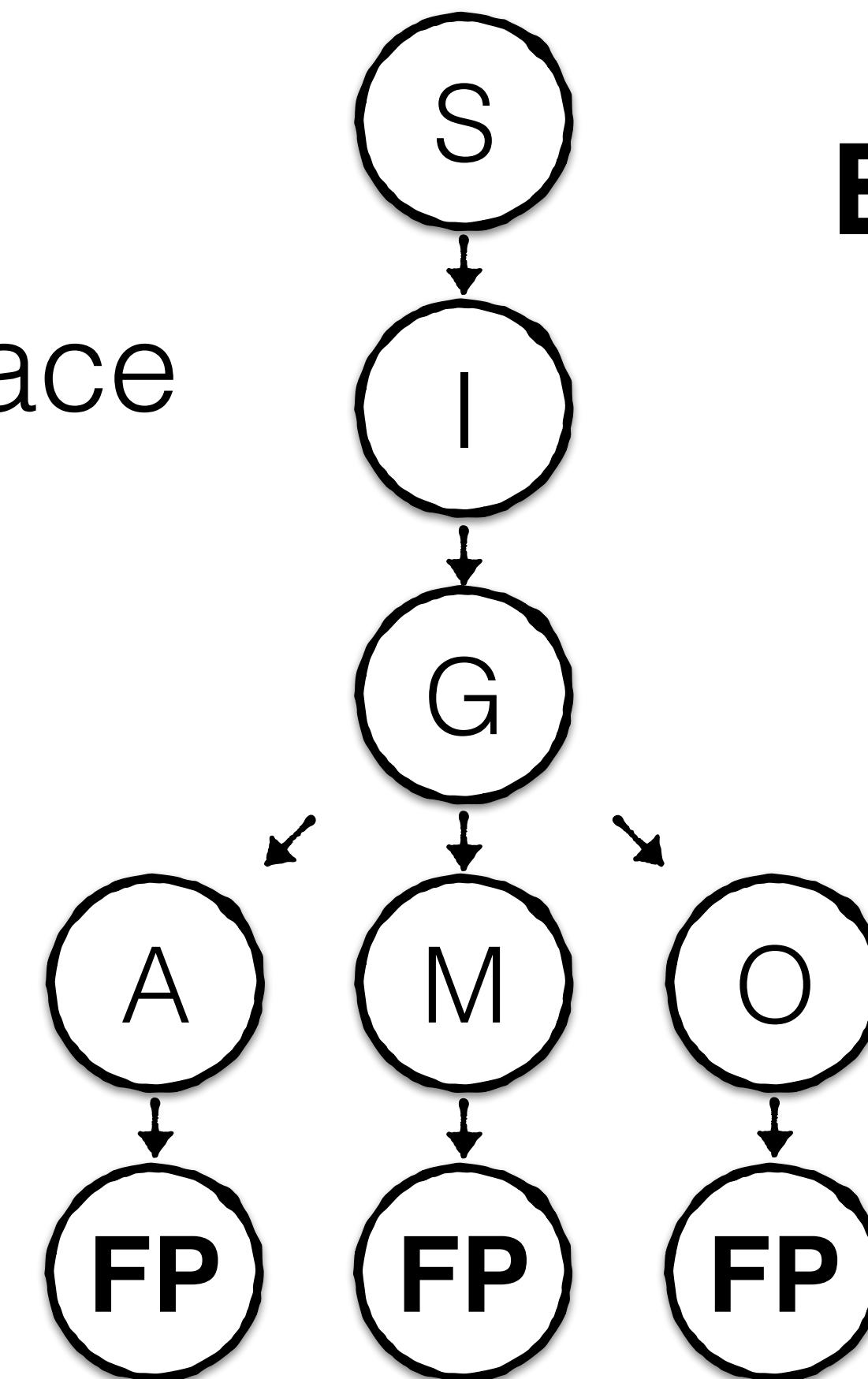
Add fingerprint for point reads

Surf

A trie of all keys

Truncated to reduce space

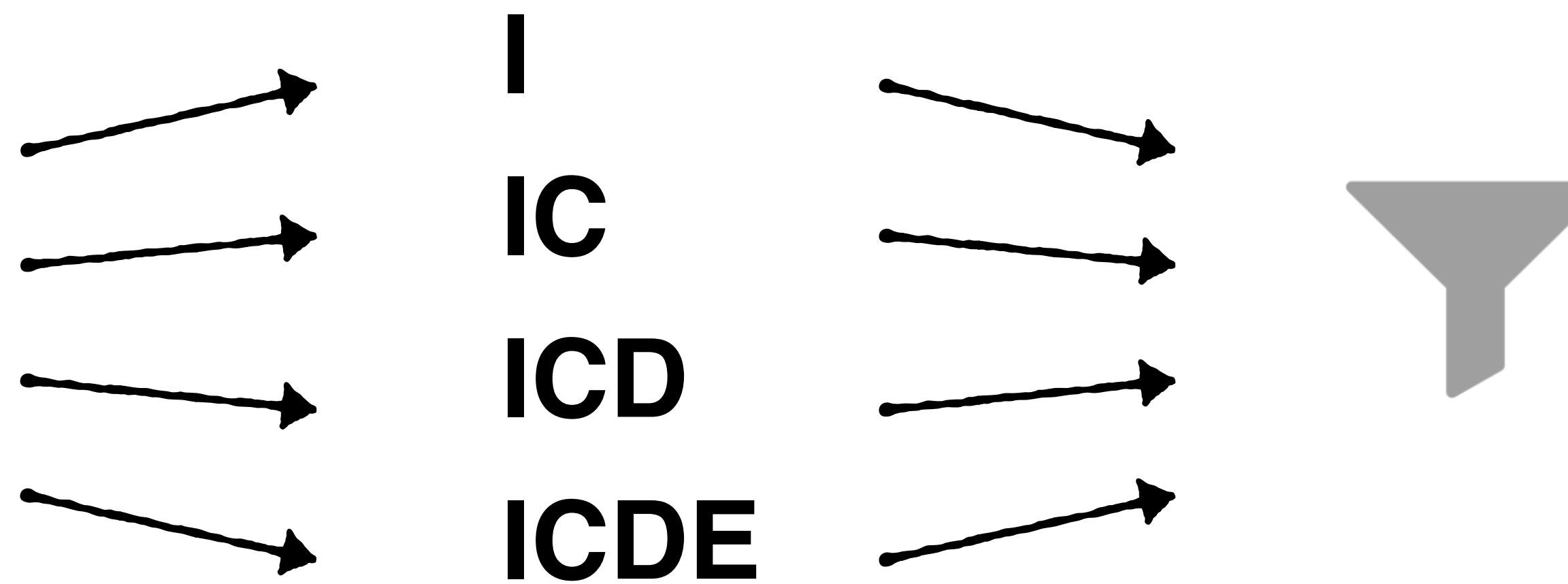
**Encoded as succinct trie
with rank & select**



Add fingerprint for point reads

Rosetta

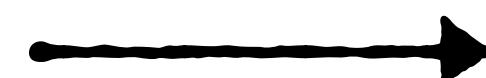
Insert(ICDE)



Add all prefixes of all keys to a Bloom filter

Rosetta

get(ICDE, ICDF)

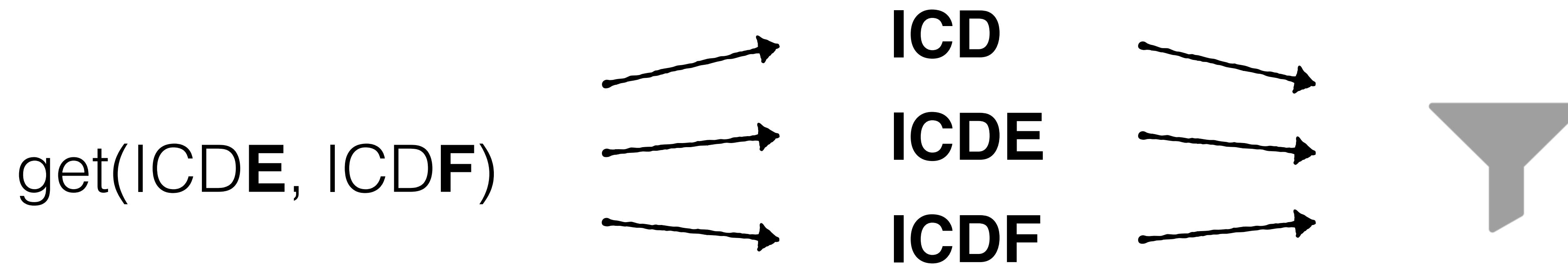


ICD



Check largest common prefixes

Rosetta



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



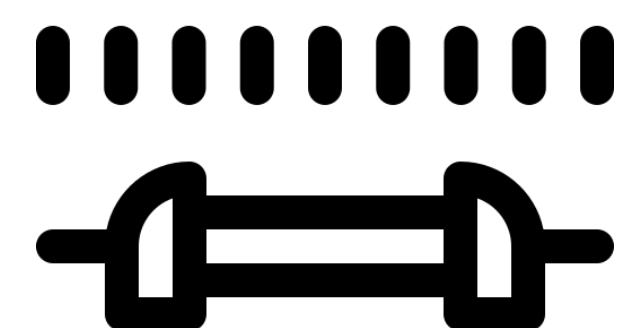
Proteus

KnorrSIGMOD22



BloomRF

MößnerEDBT23

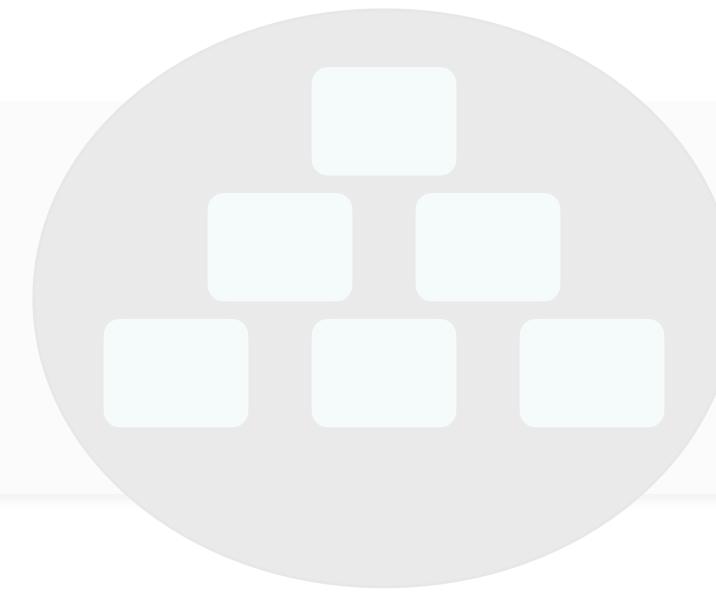


REncoder

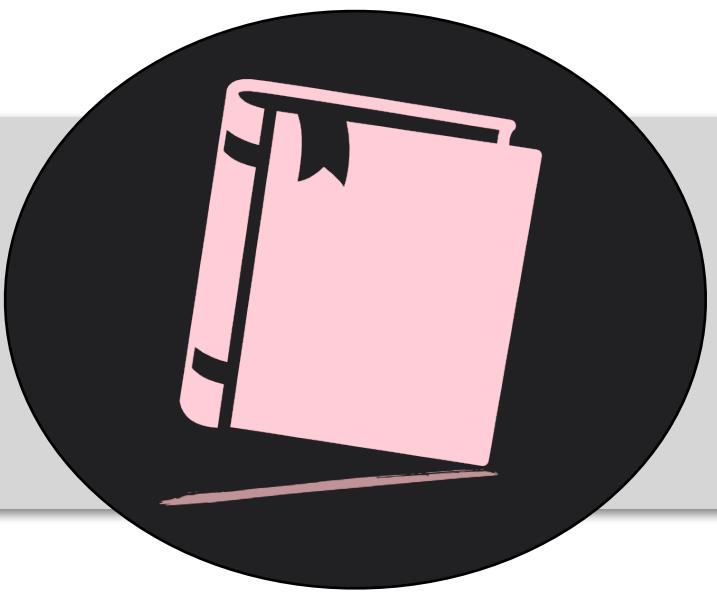
WangICDE23

Outline

Part 1: **LSM Basics**



Part 2B: Read Optimizations in LSMS

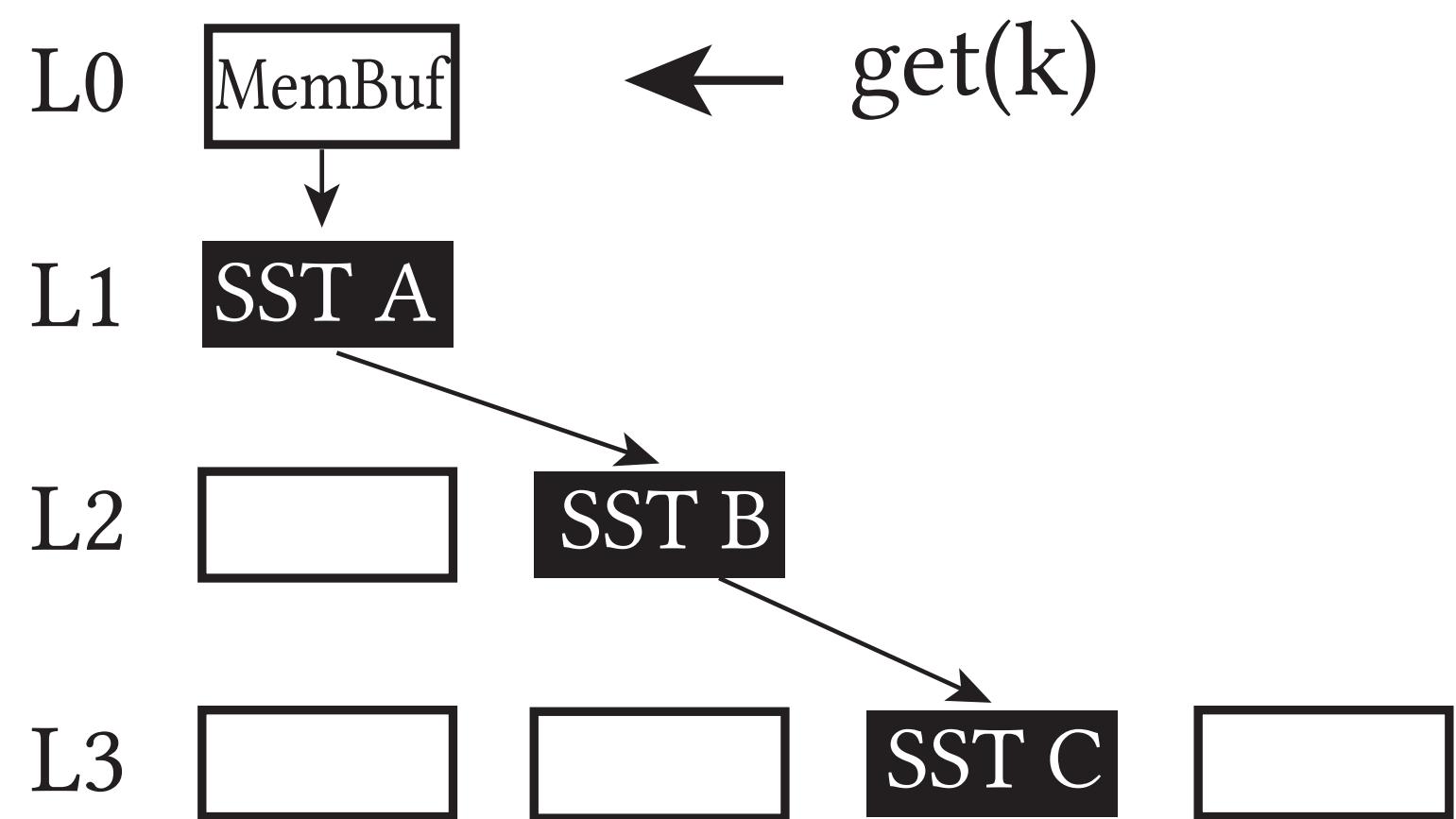


Part 3: **Navigating the LSM Design Space**



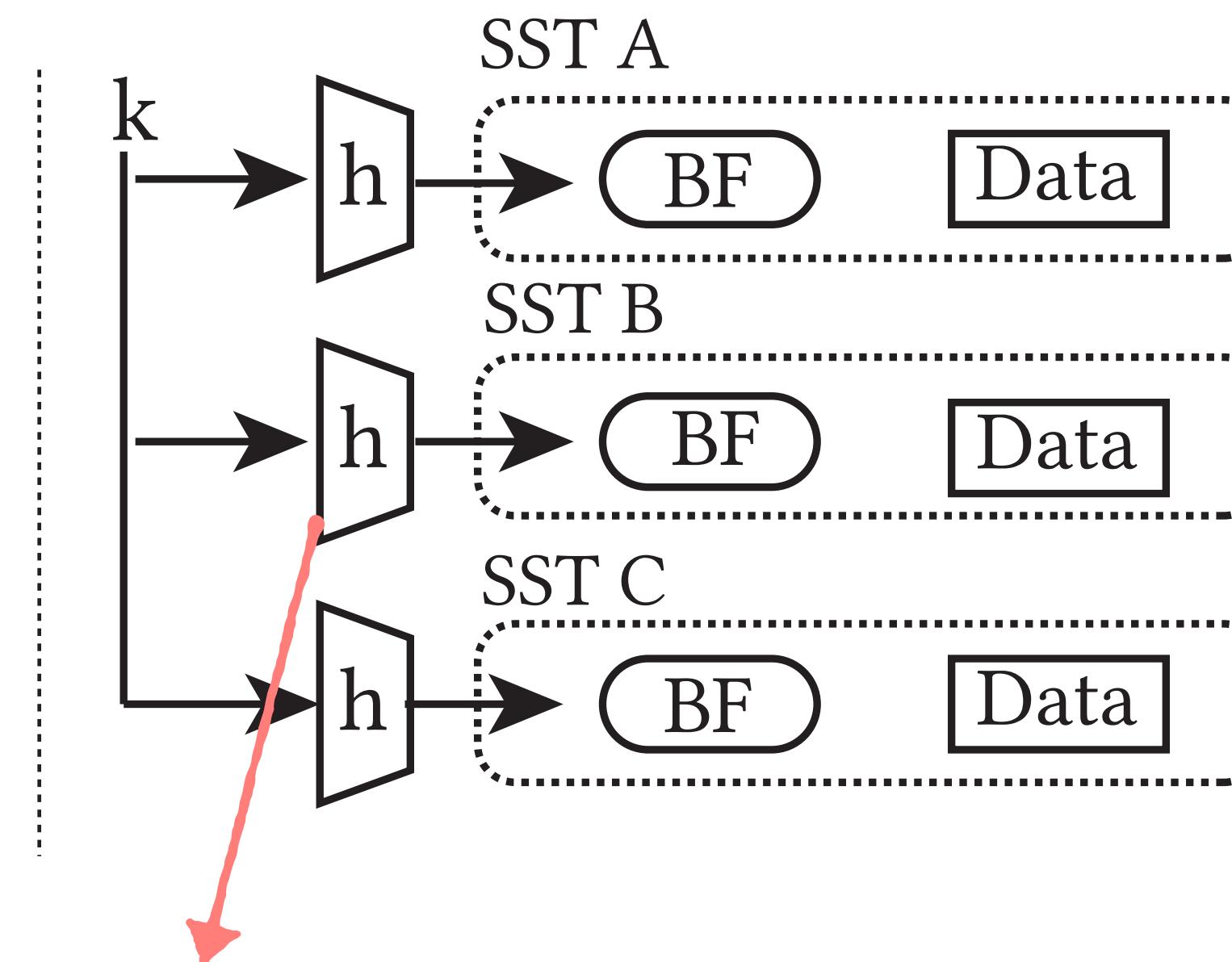
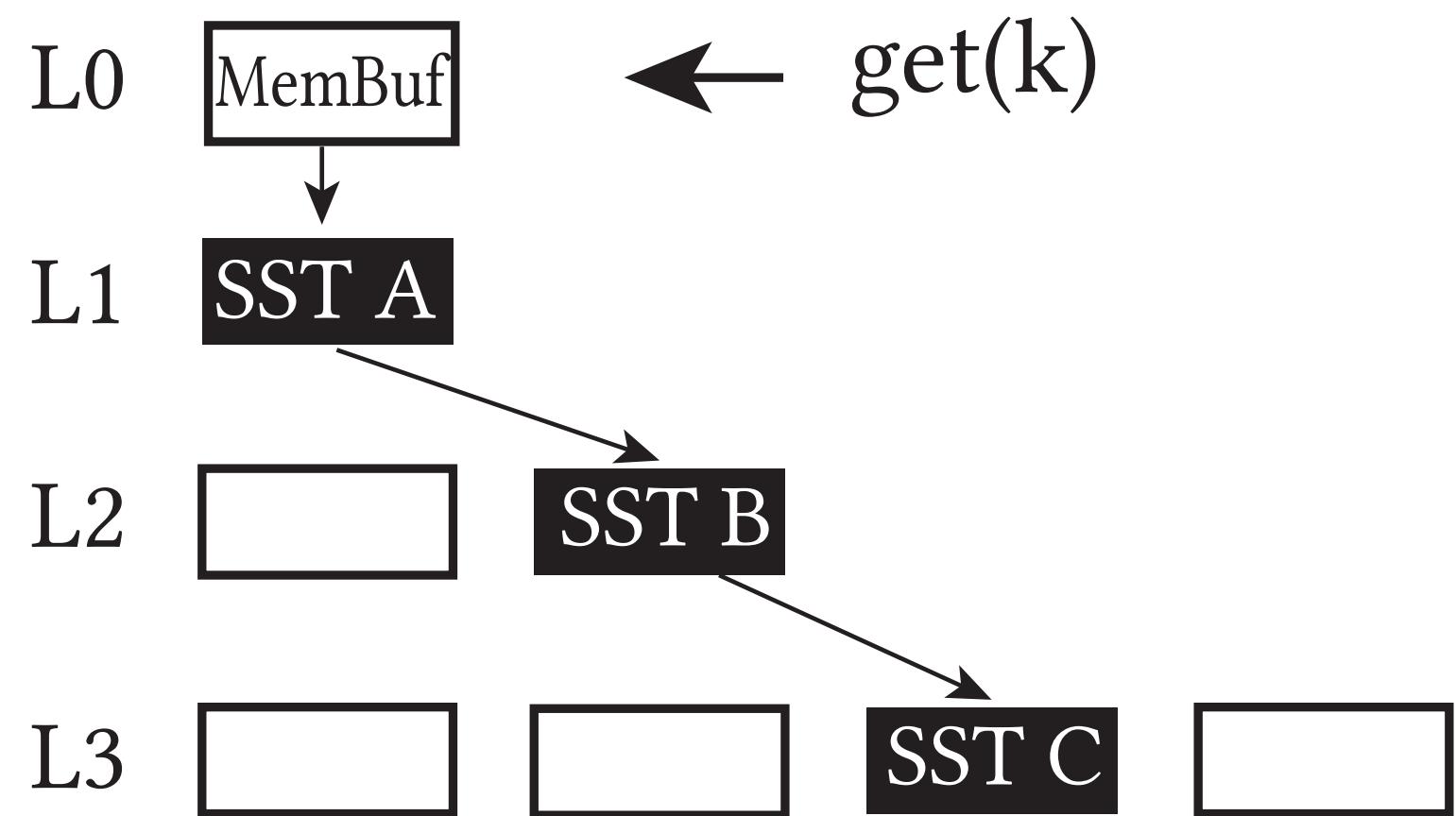
Reducing CPU Overheads in LSMs

For every query ...



Reducing CPU Overheads in LSMs

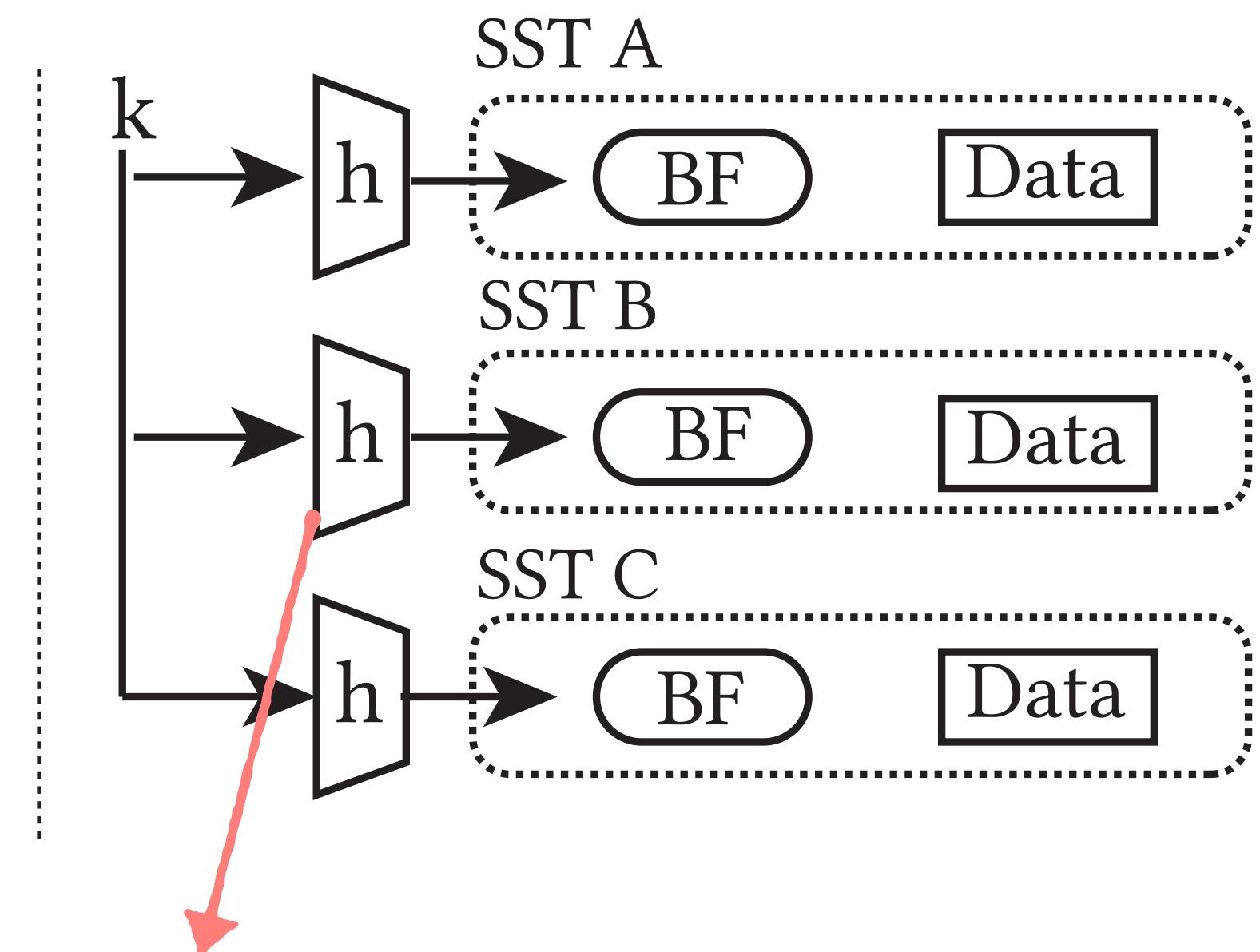
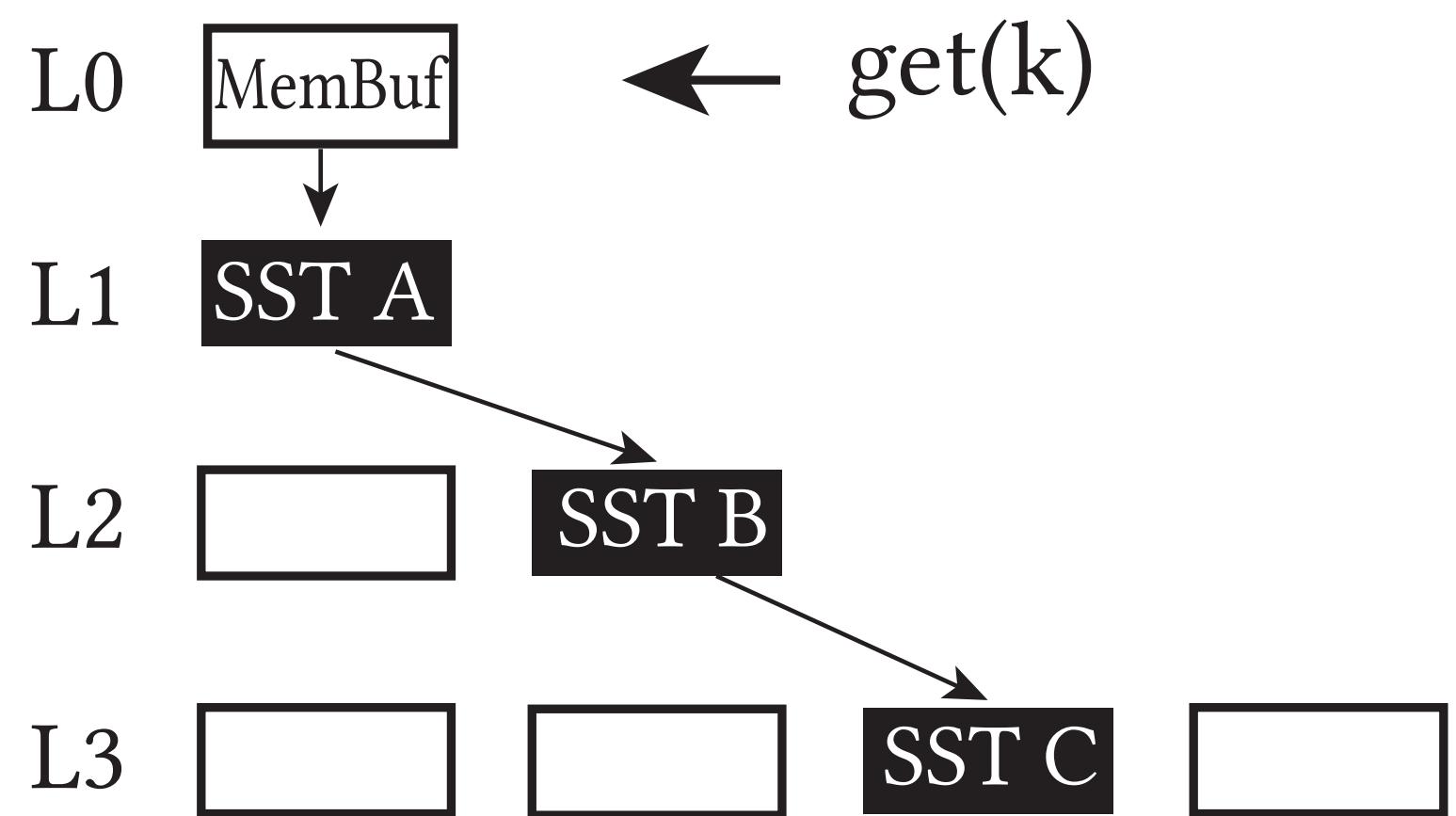
For every query ...



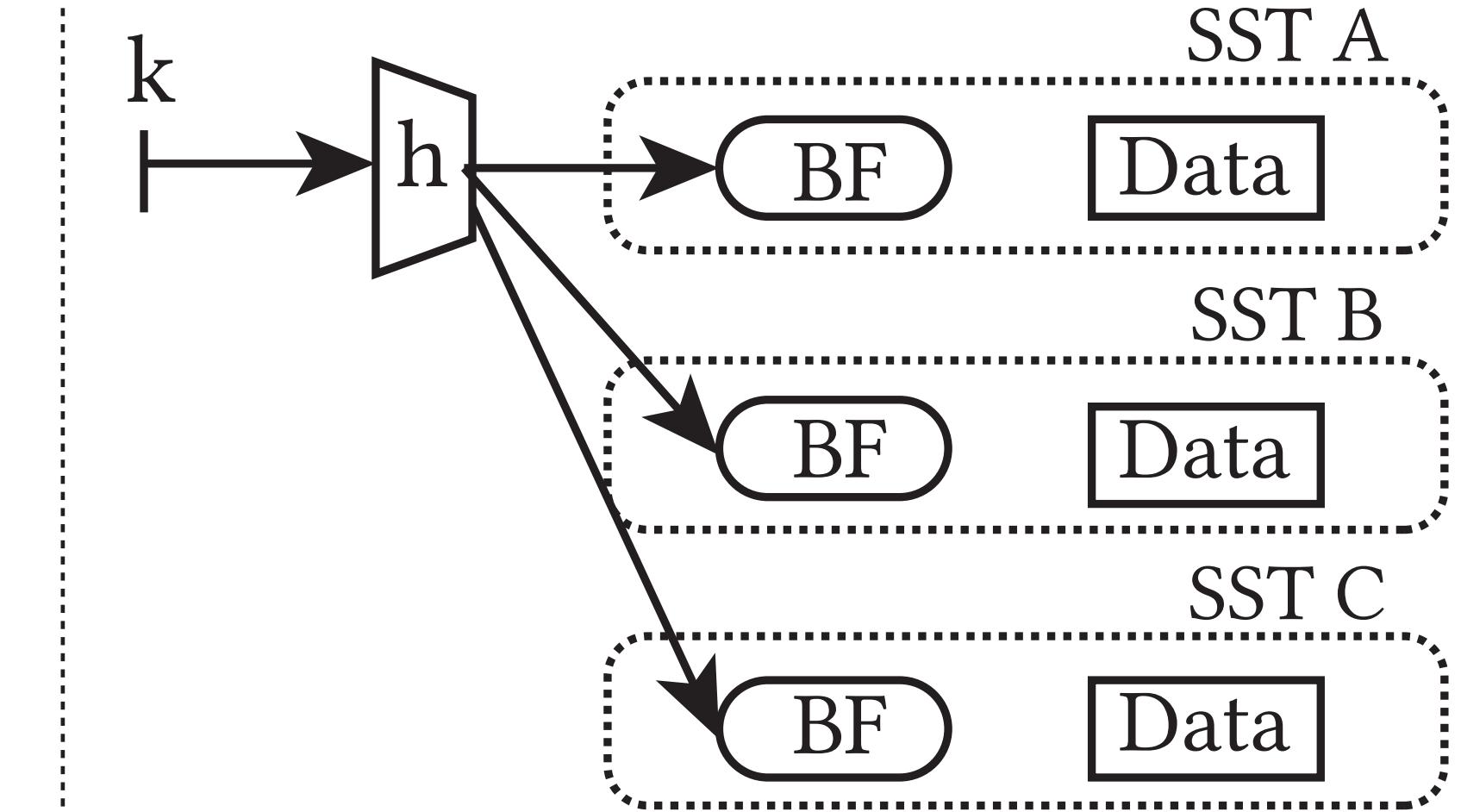
The same hash function is
calculated **O(L)** times

Reducing CPU Overheads in LSMs

For every query ...



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



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



Filters under Memory Pressure



data size ↑

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

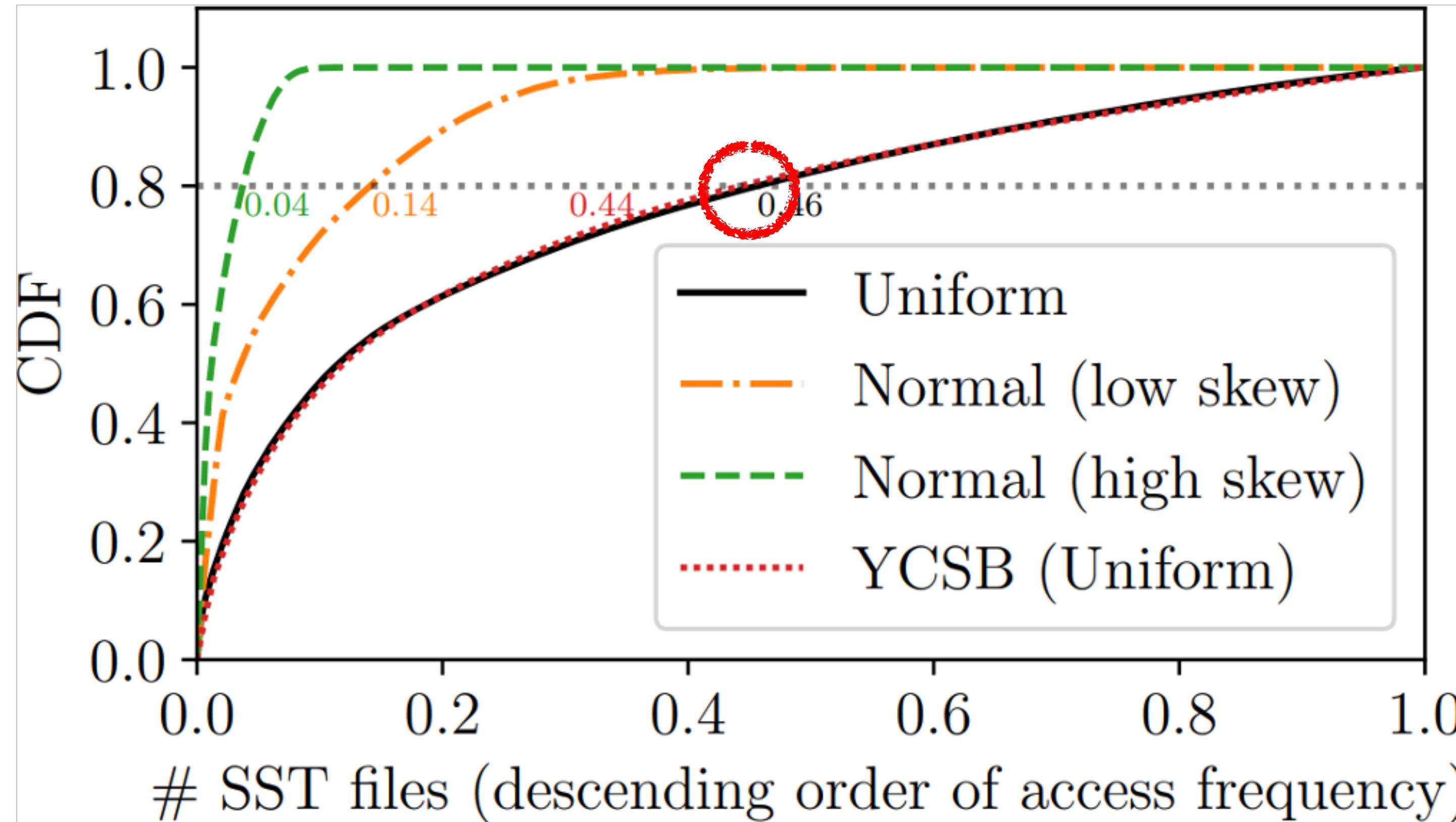
1KB entry, 64B key, BPK=10

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

Filters under Memory Pressure



MunADMS22

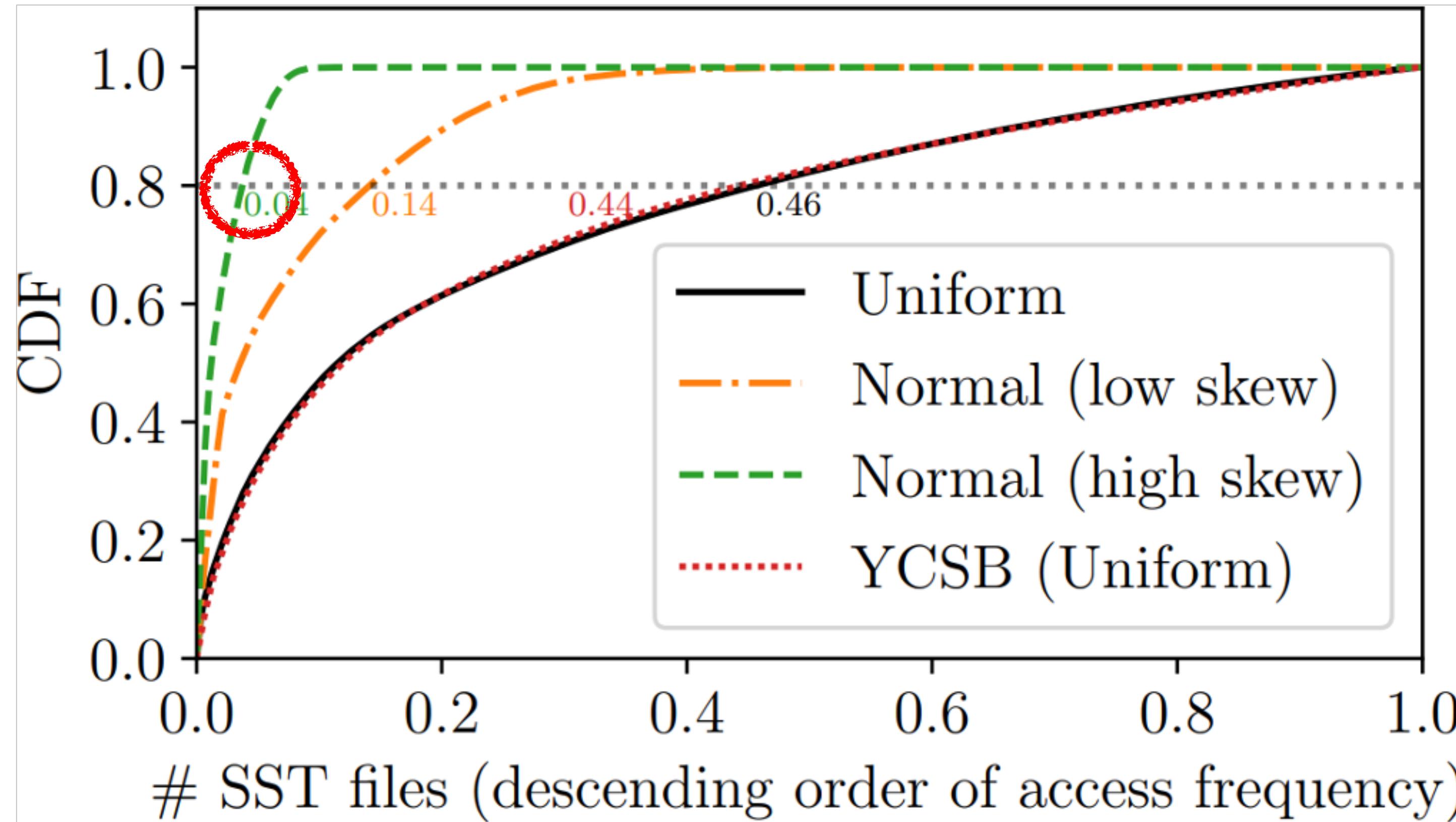


Even in a **perfectly uniform** workload,
80% of the queries access 45% of the files

Filters under Memory Pressure



MunADMS22

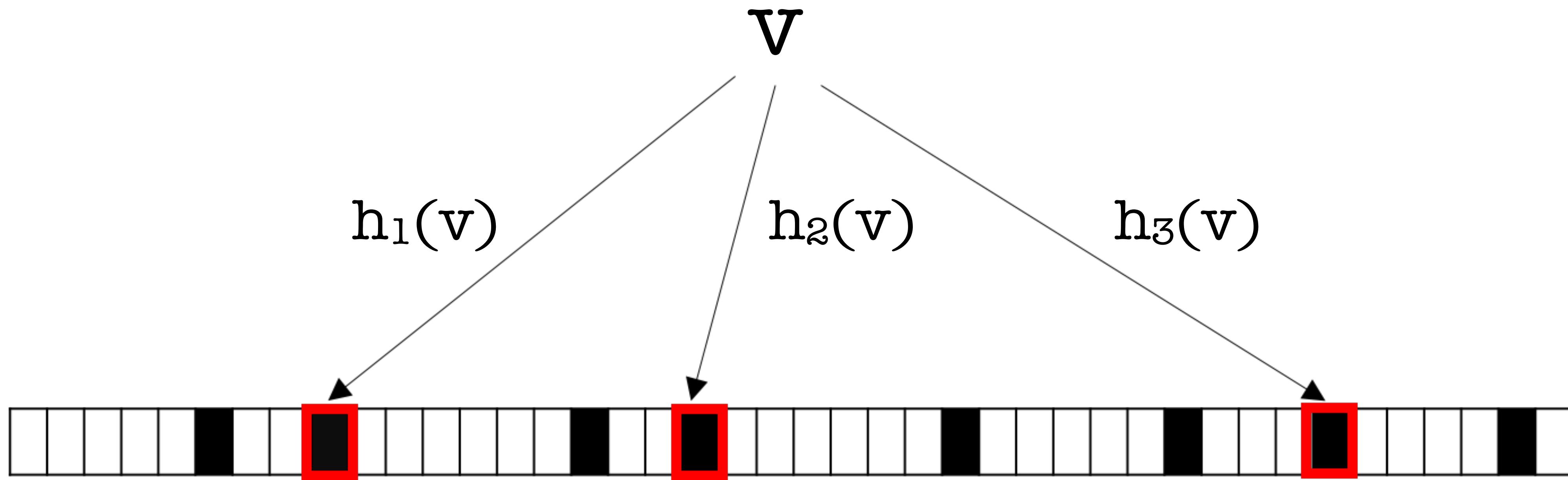


For a **skewed** workload,
80% of the queries access less than **5% of the files**

Filters under Memory Pressure



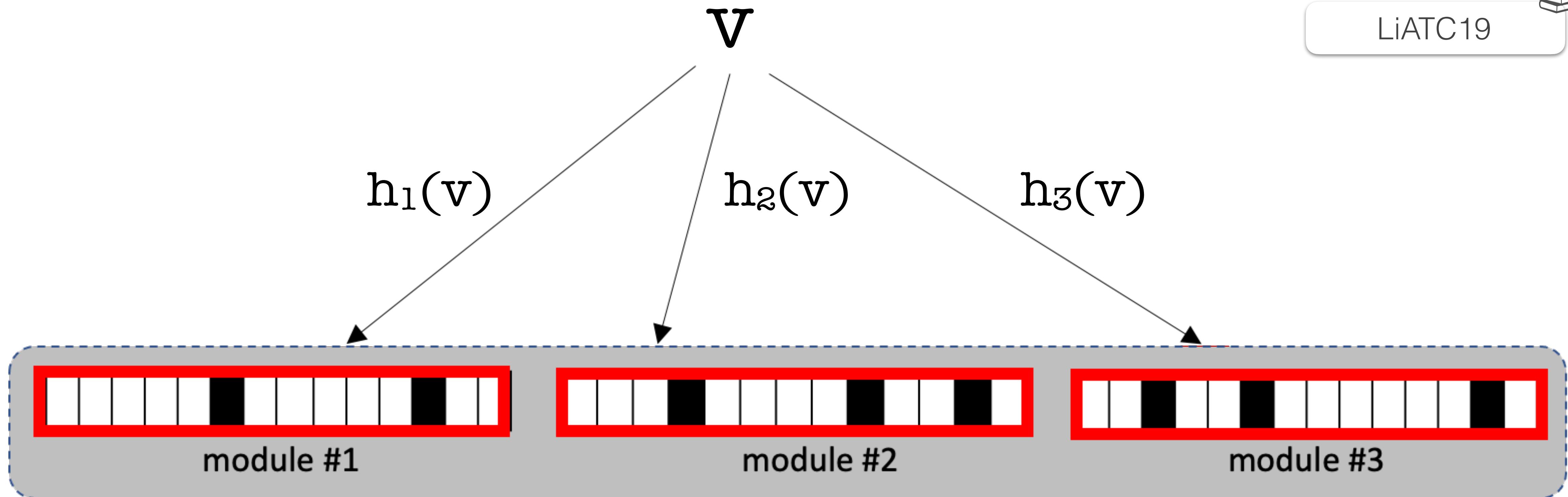
MunADMS22



Filters under Memory Pressure

MunADMS22

LiATC19



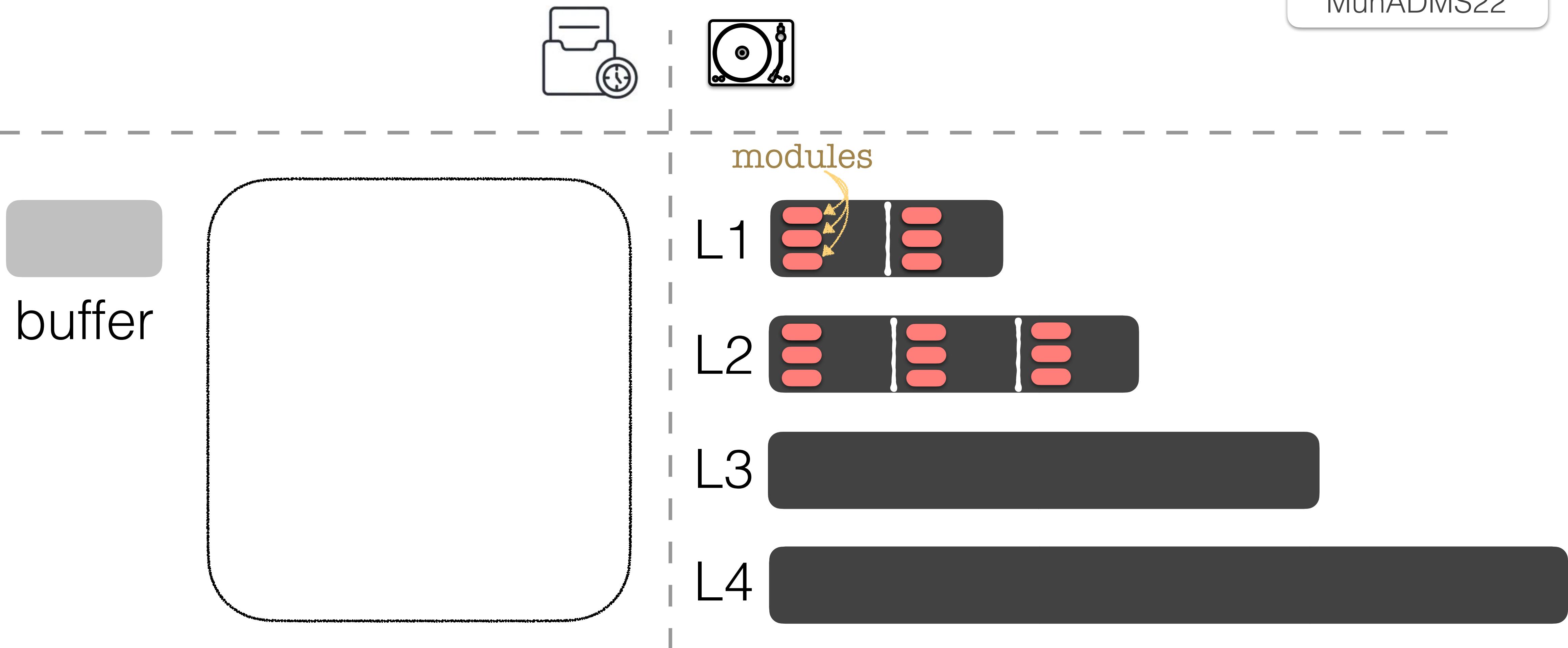
Modular Bloom filter is a collection of smaller Bloom filters

Elastic Bloom filter also works based on the same principle

Filters under Memory Pressure



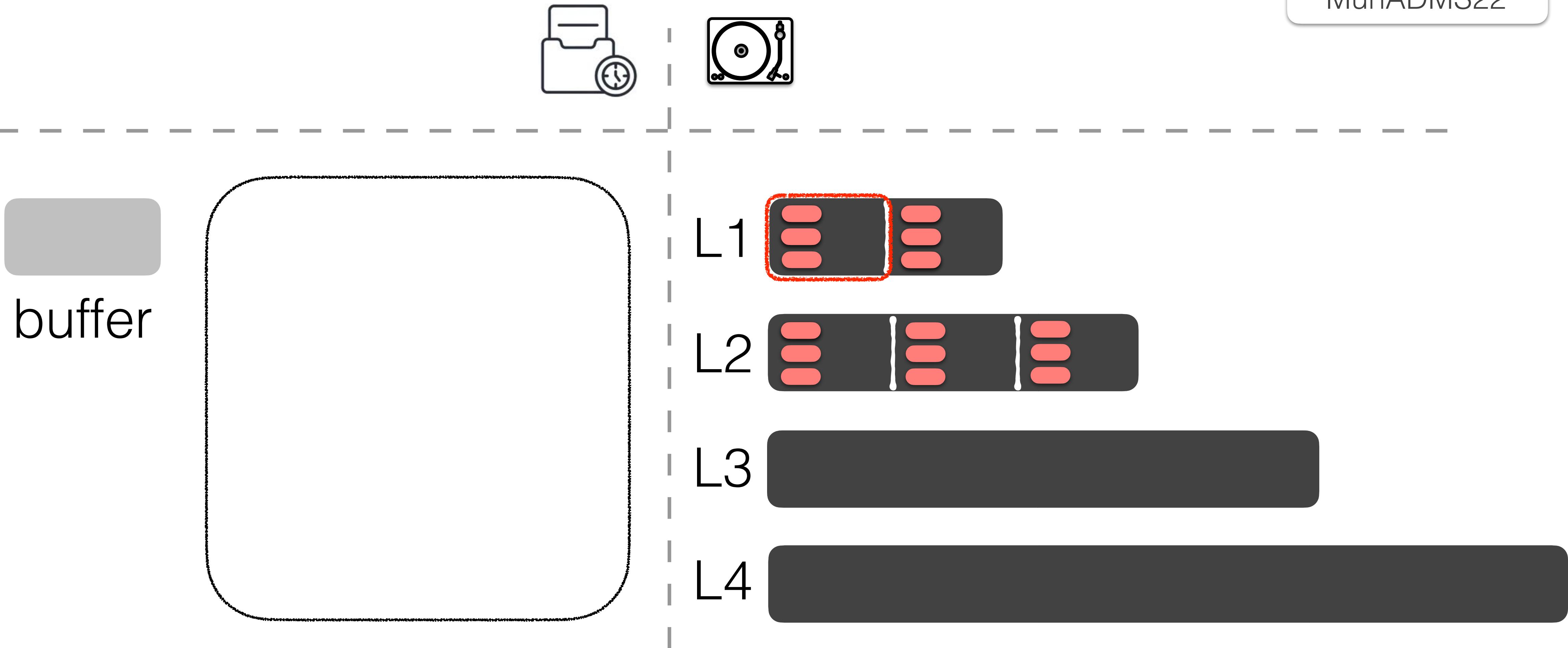
MunADMS22



Filters under Memory Pressure



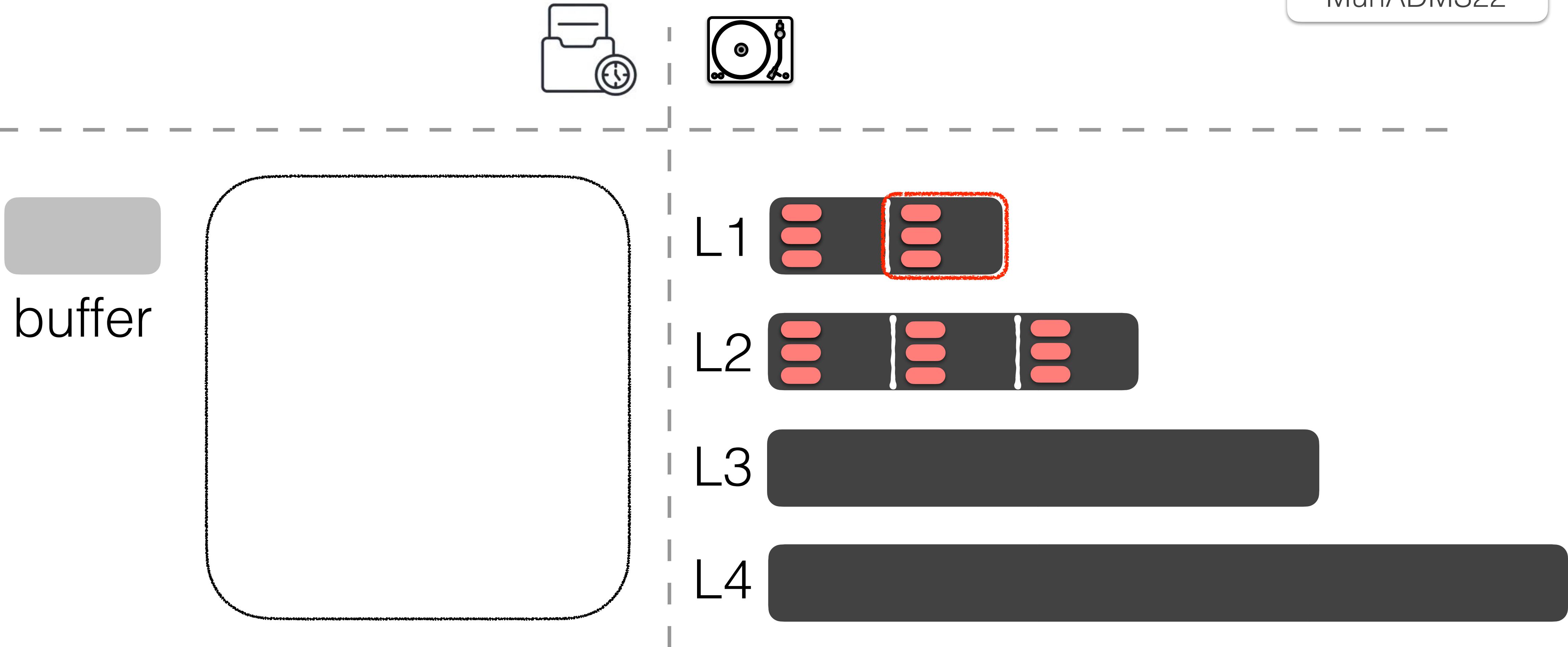
MunADMS22



Filters under Memory Pressure



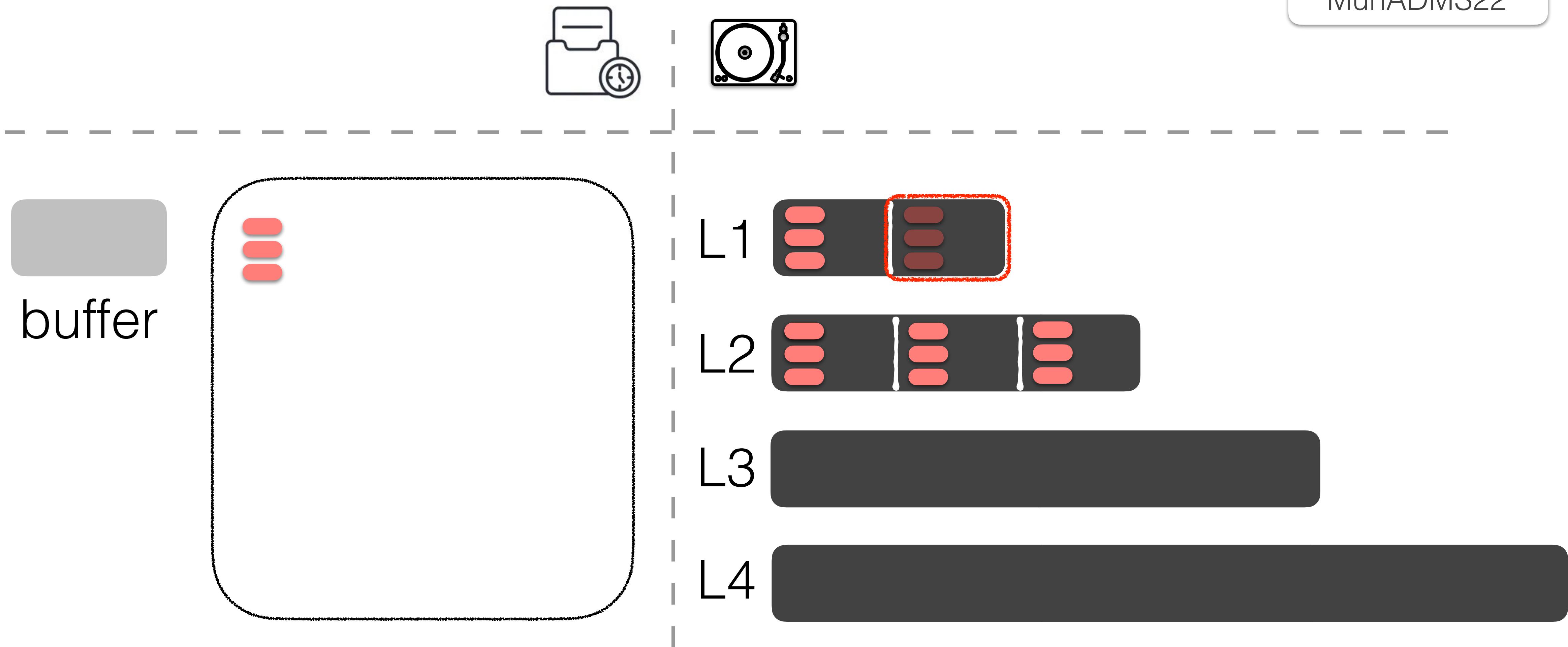
MunADMS22



Filters under Memory Pressure



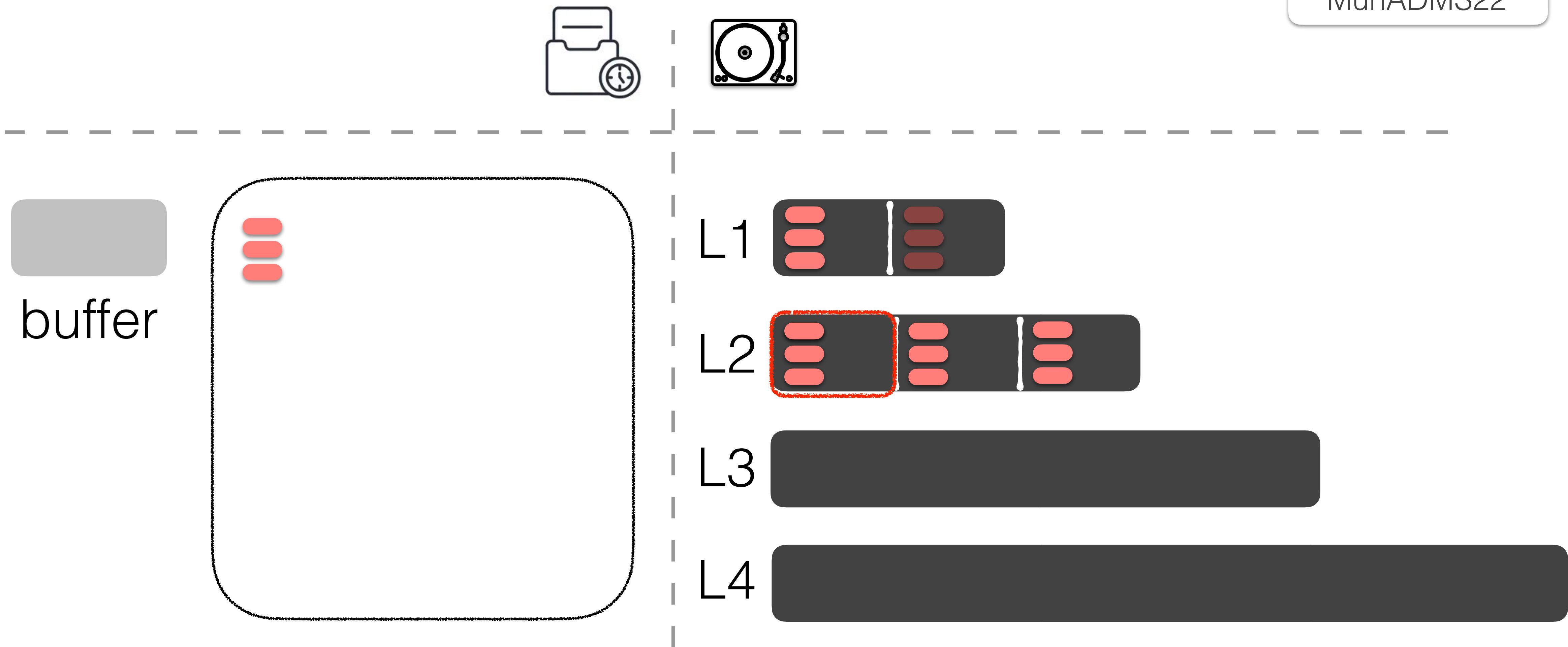
MunADMS22



Filters under Memory Pressure



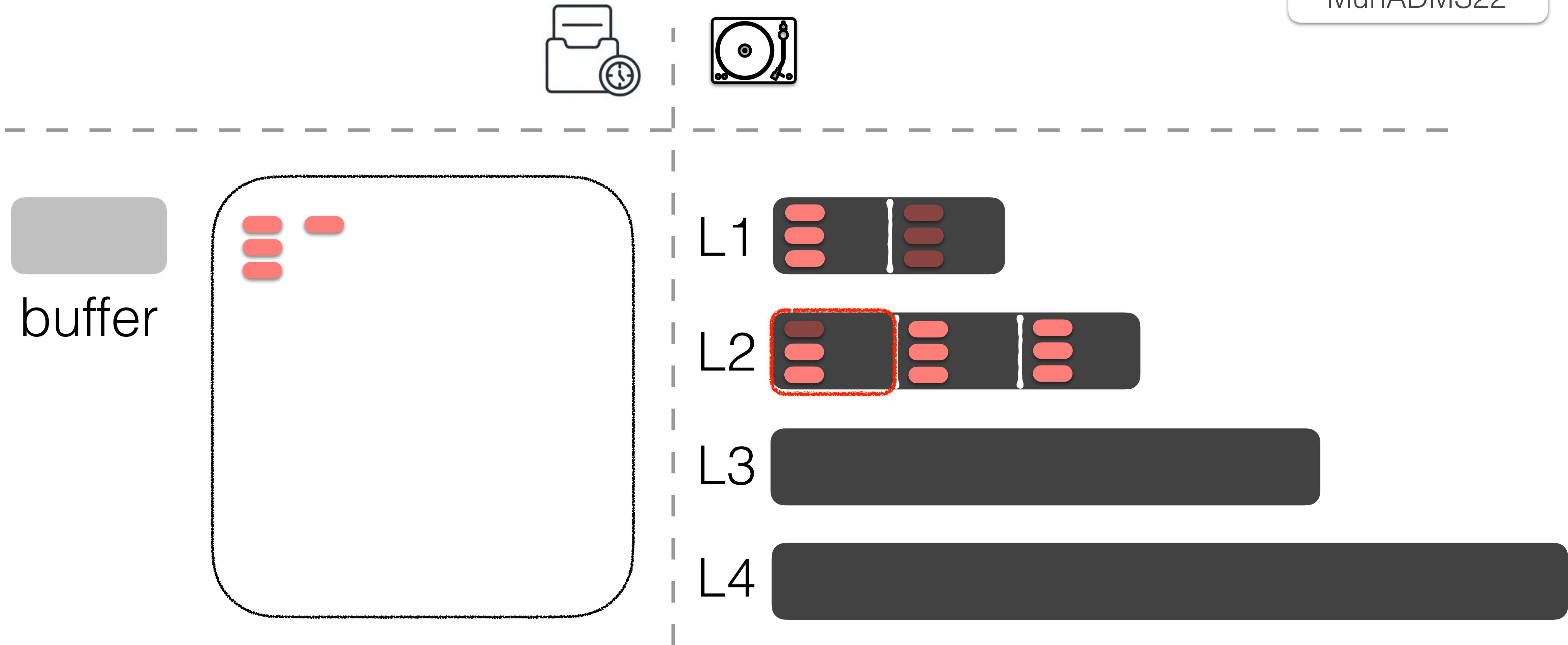
MunADMS22



Filters under Memory Pressure



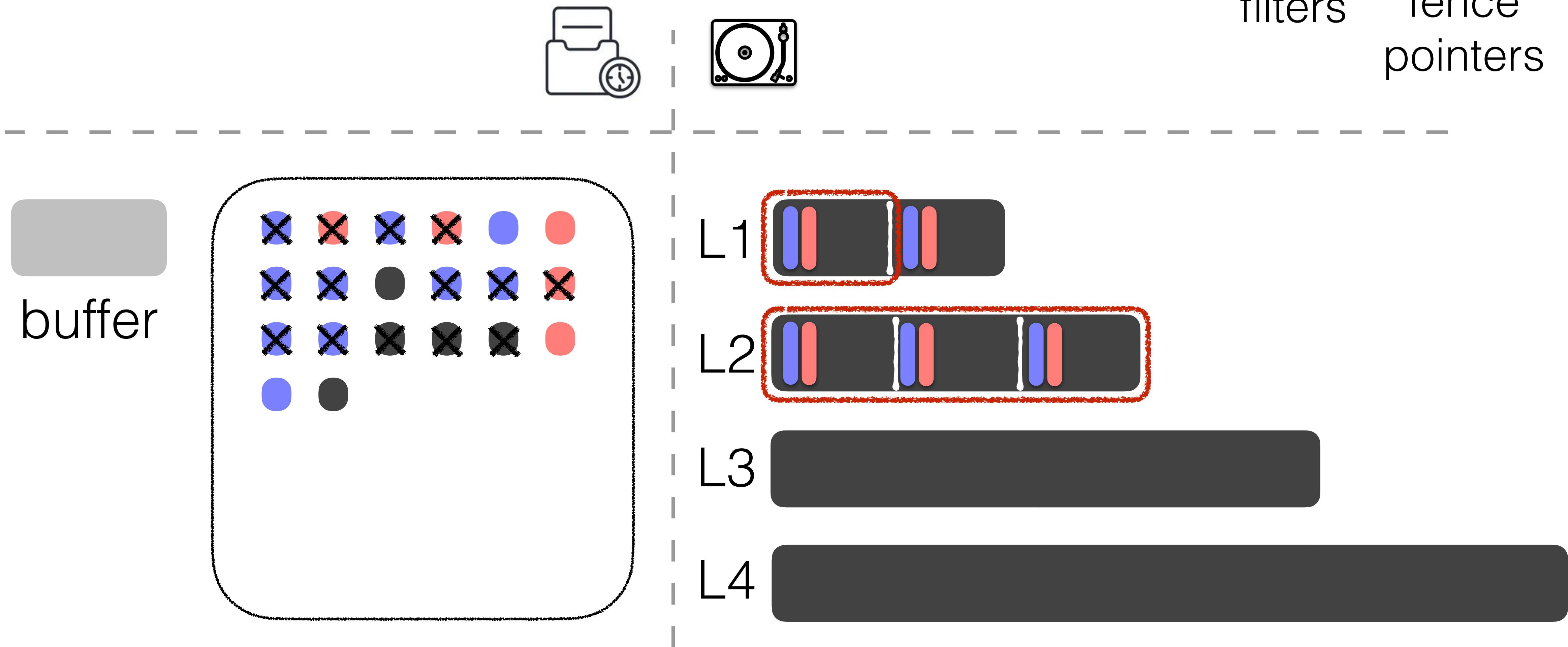
MunADMS22



Overall, better performance with smaller memory budget

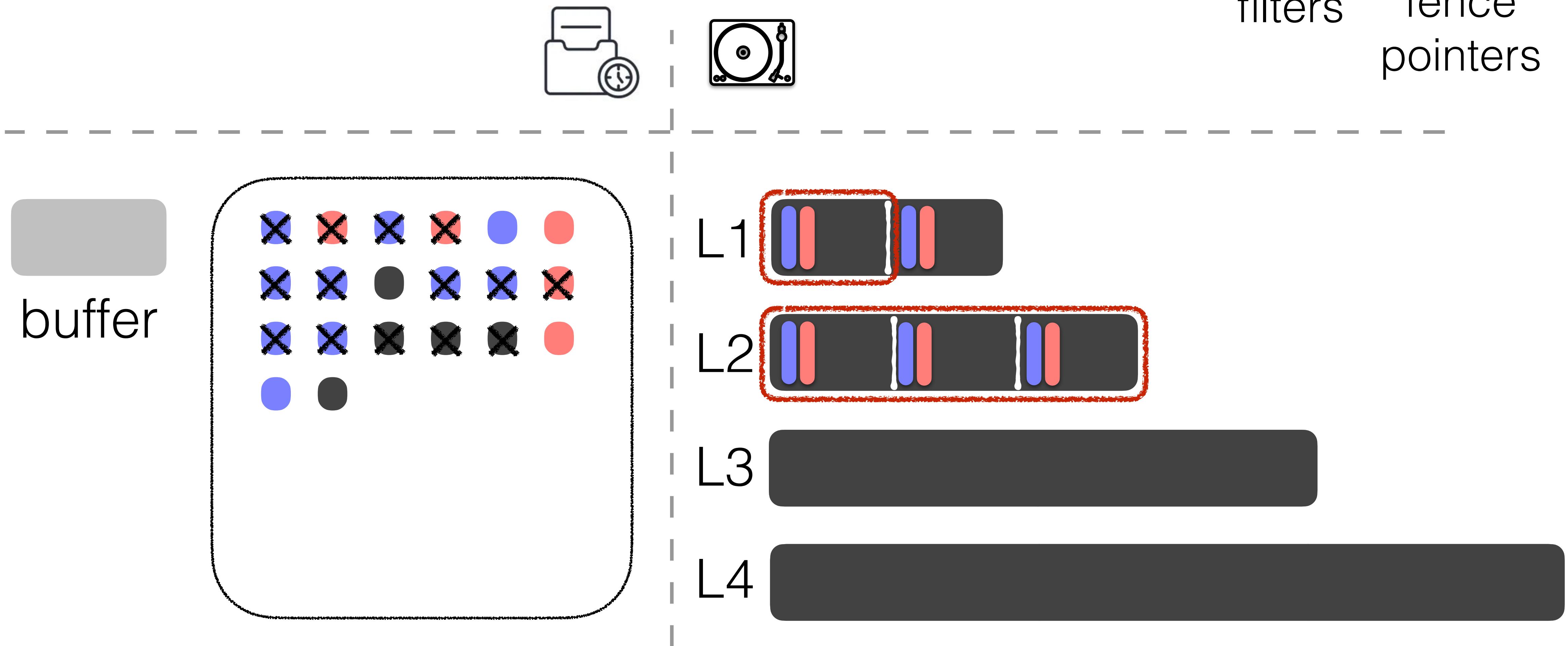
Compactions and Caching

 filters  fence pointers



Compactions and Caching

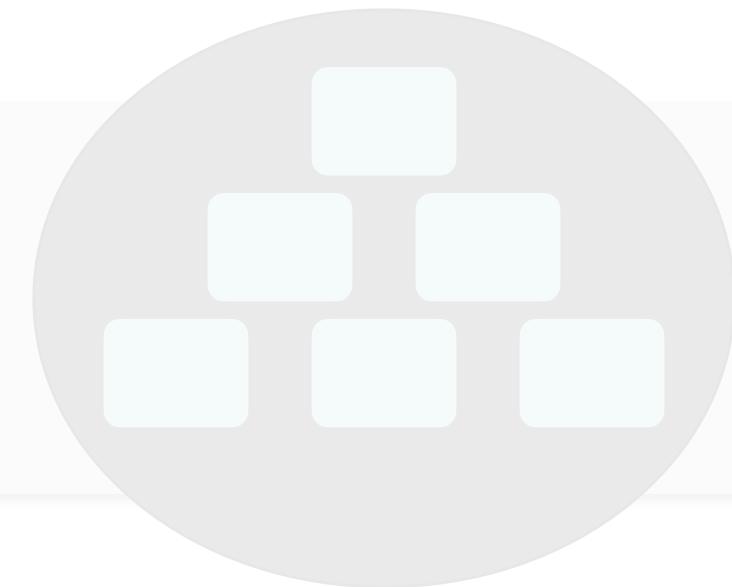
 filters  fence pointers



Leaper : A learned pre-fetcher that improves reads

Outline

Part 1: **LSM Basics**



Part 2: **Read Optimizations in LSMS**

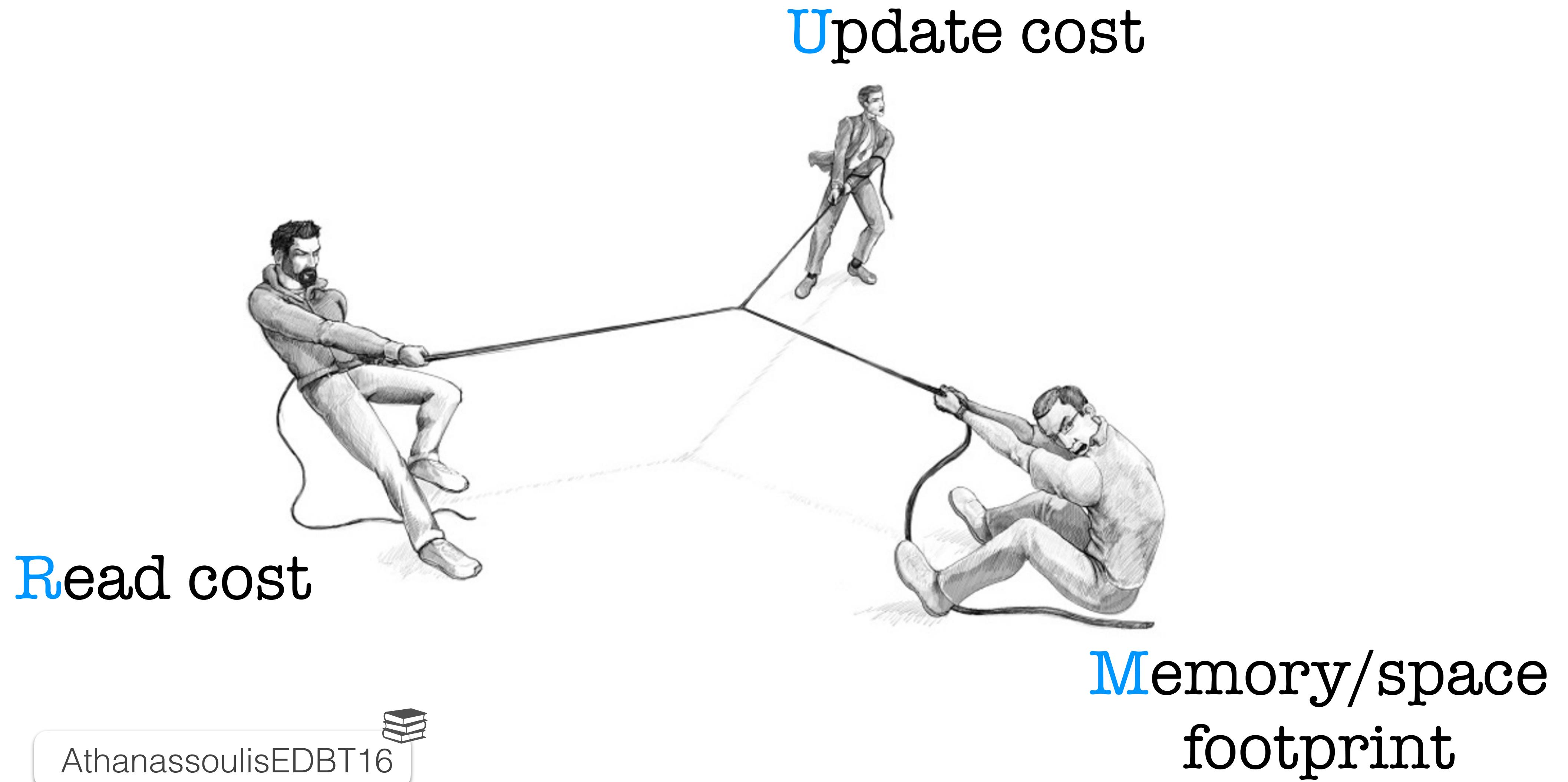


Part 3: **Navigating the LSM Design Space**

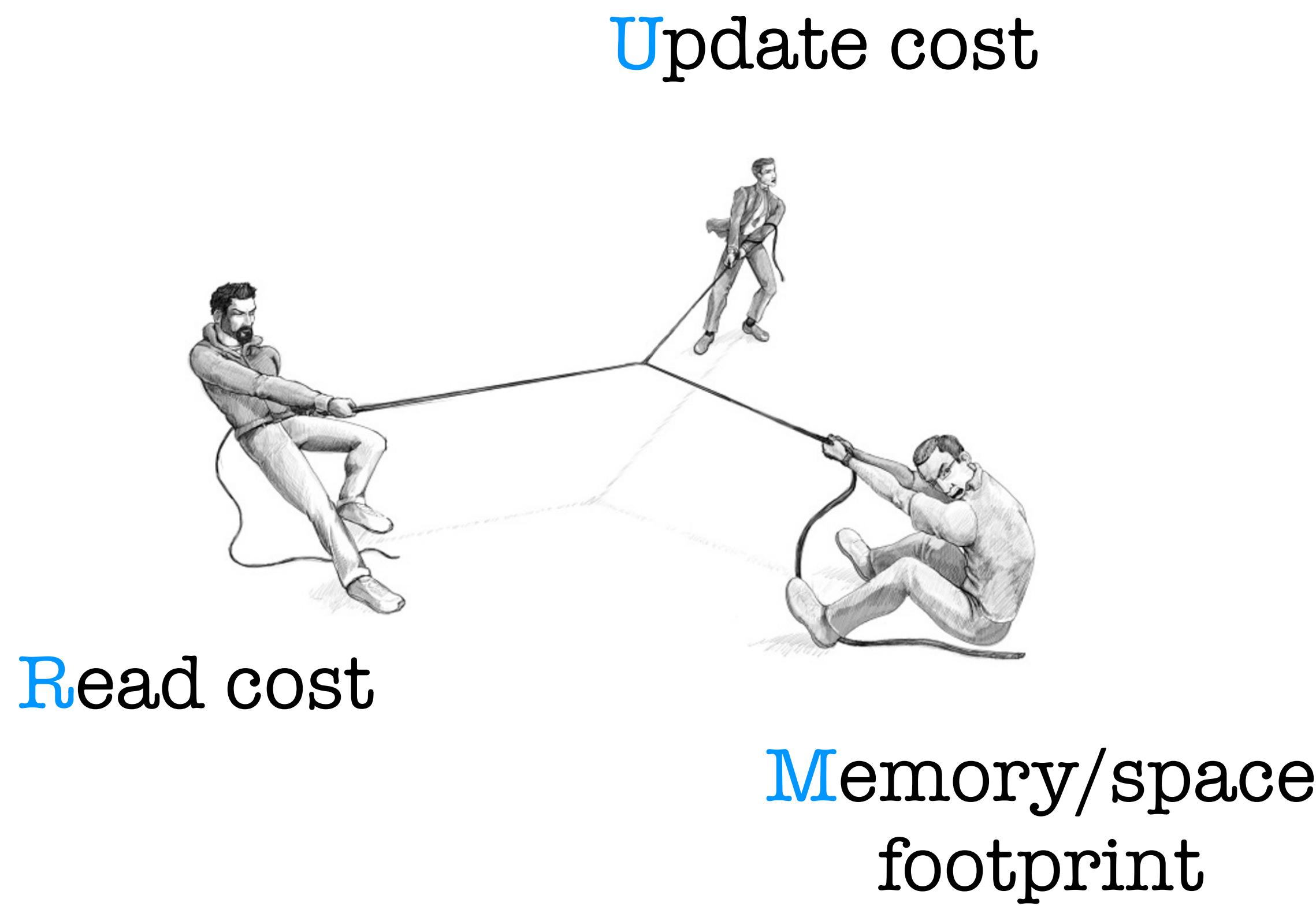


LSM Design Space

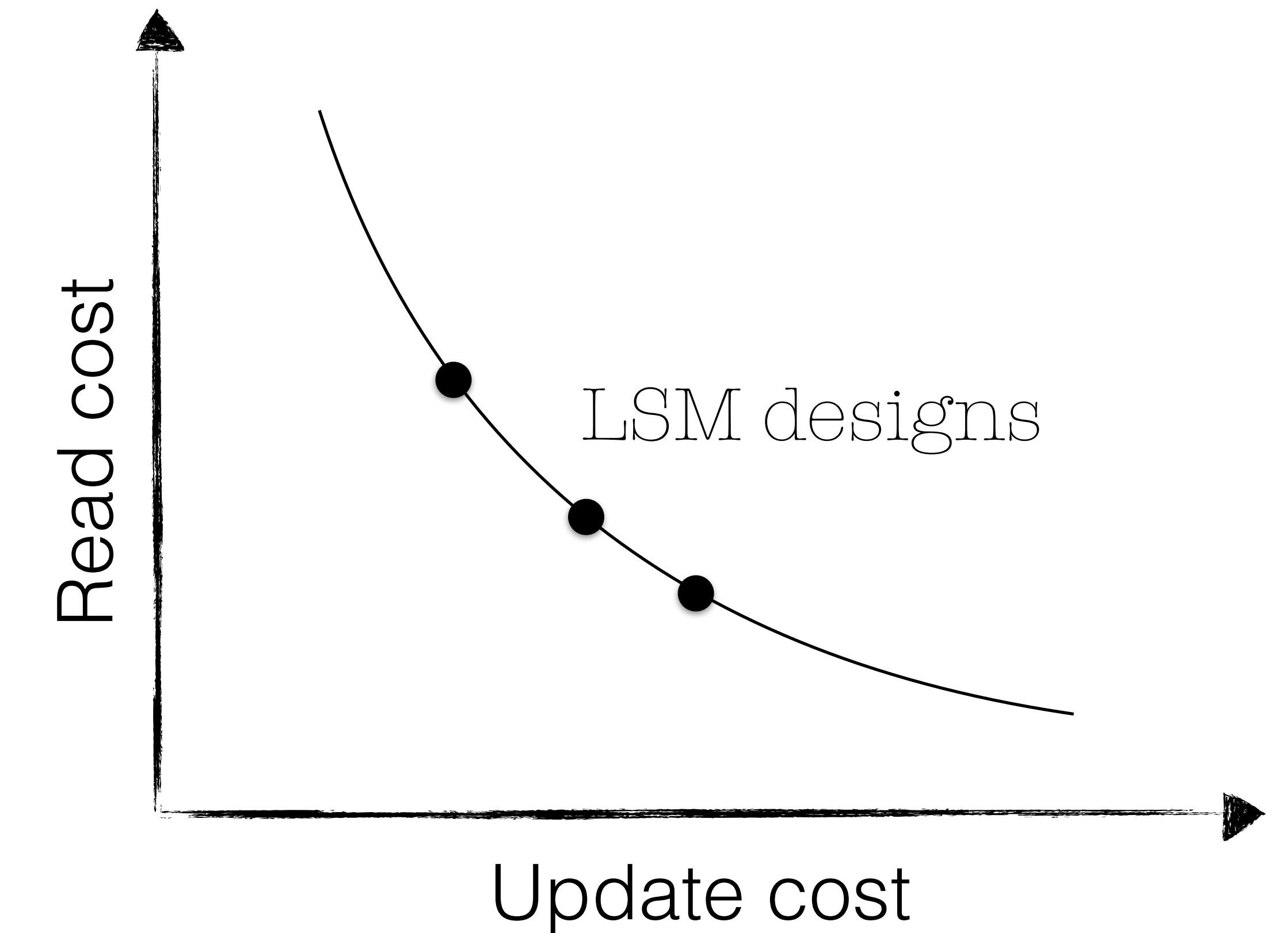
LSM Design Space



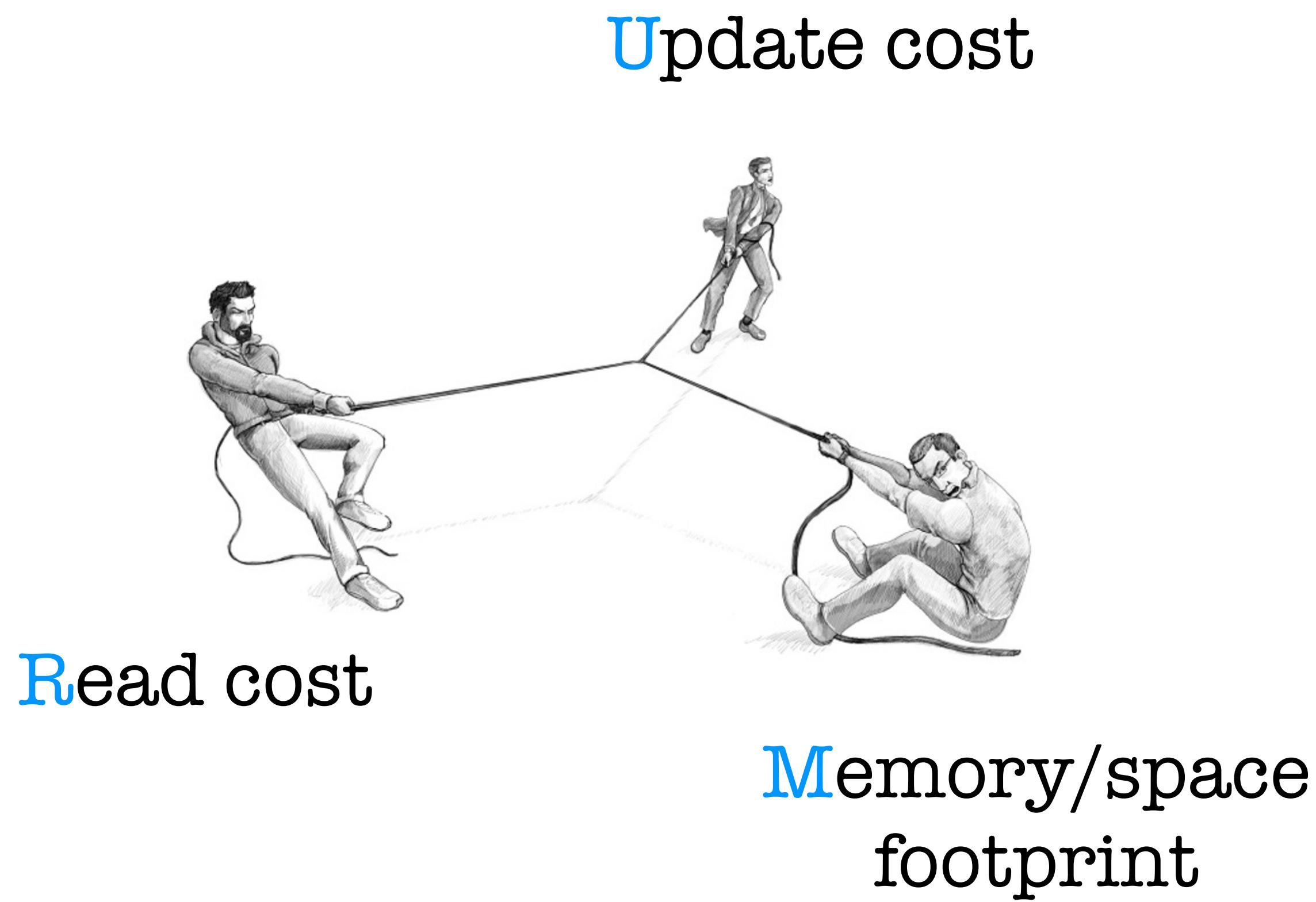
LSM Design Space



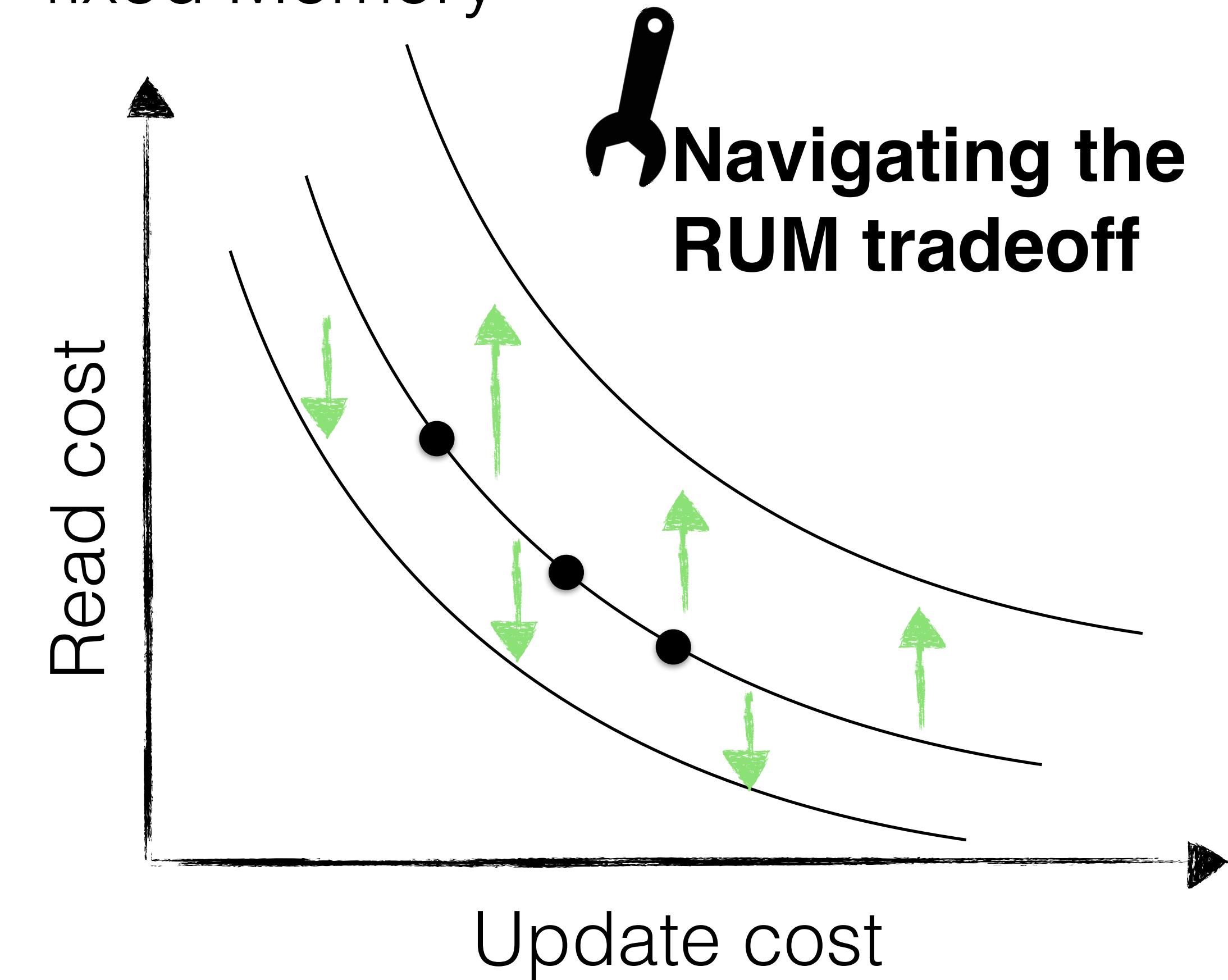
fixed Memory



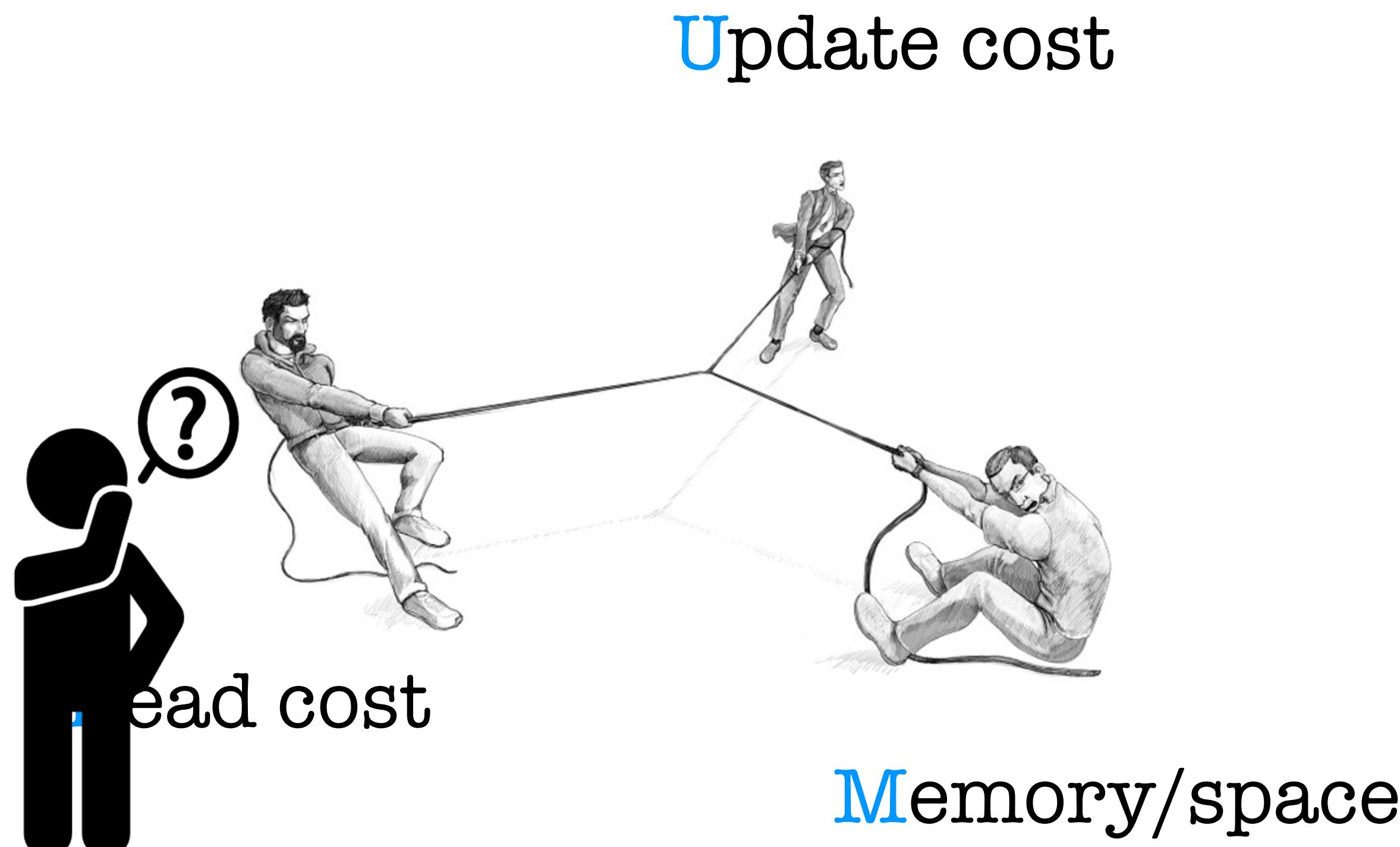
LSM Design Space



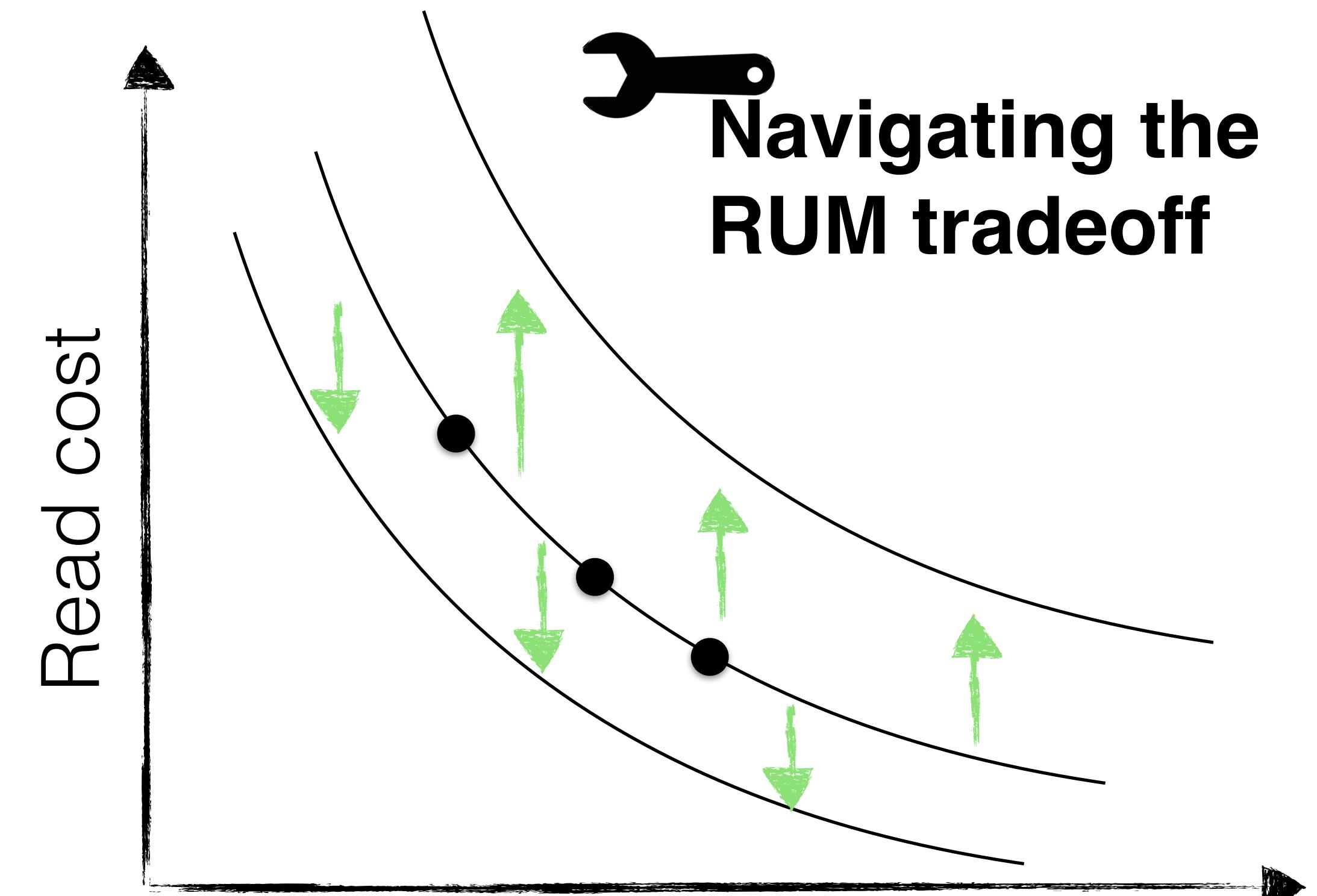
fixed Memory



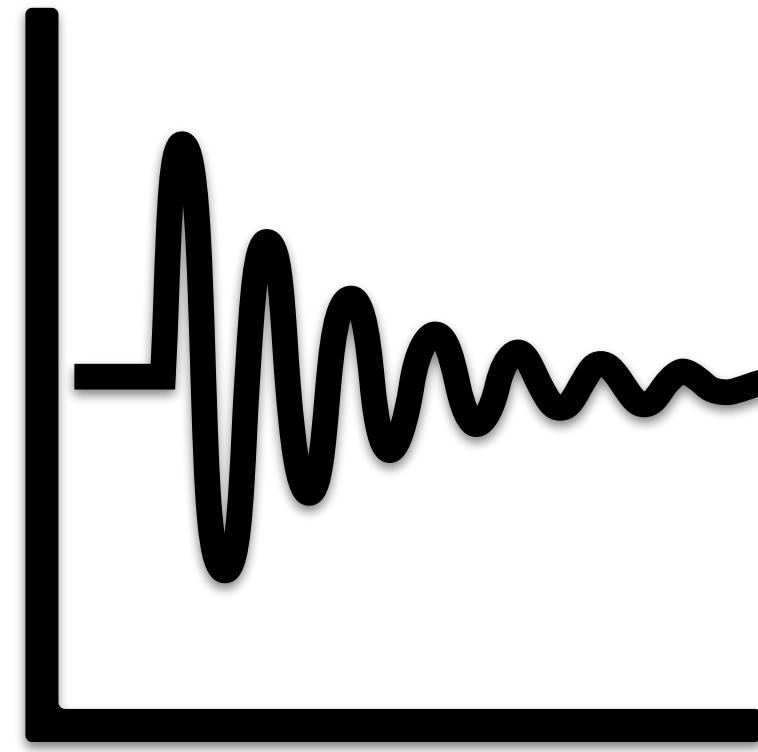
LSM Design Space



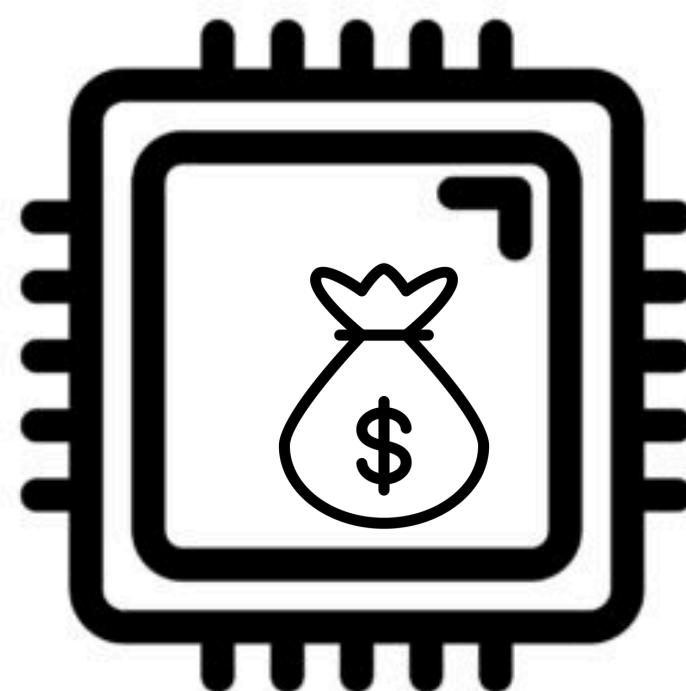
fixed Memory



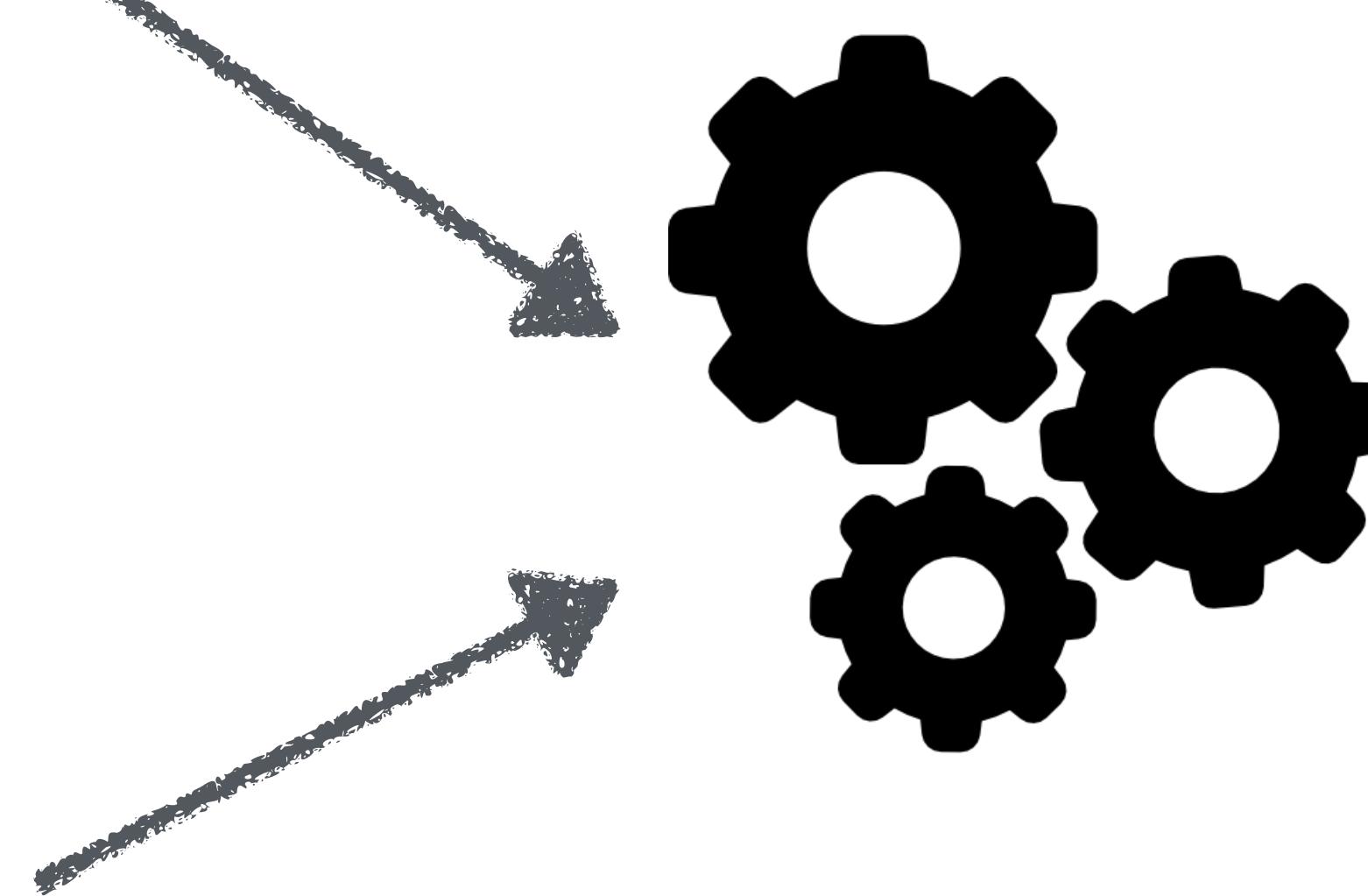
How to optimally allocate the available memory?



workload



memory
budget



How to allocate memory
between LSM components

How to allocate memory
among BFs in LSM



M : total memory

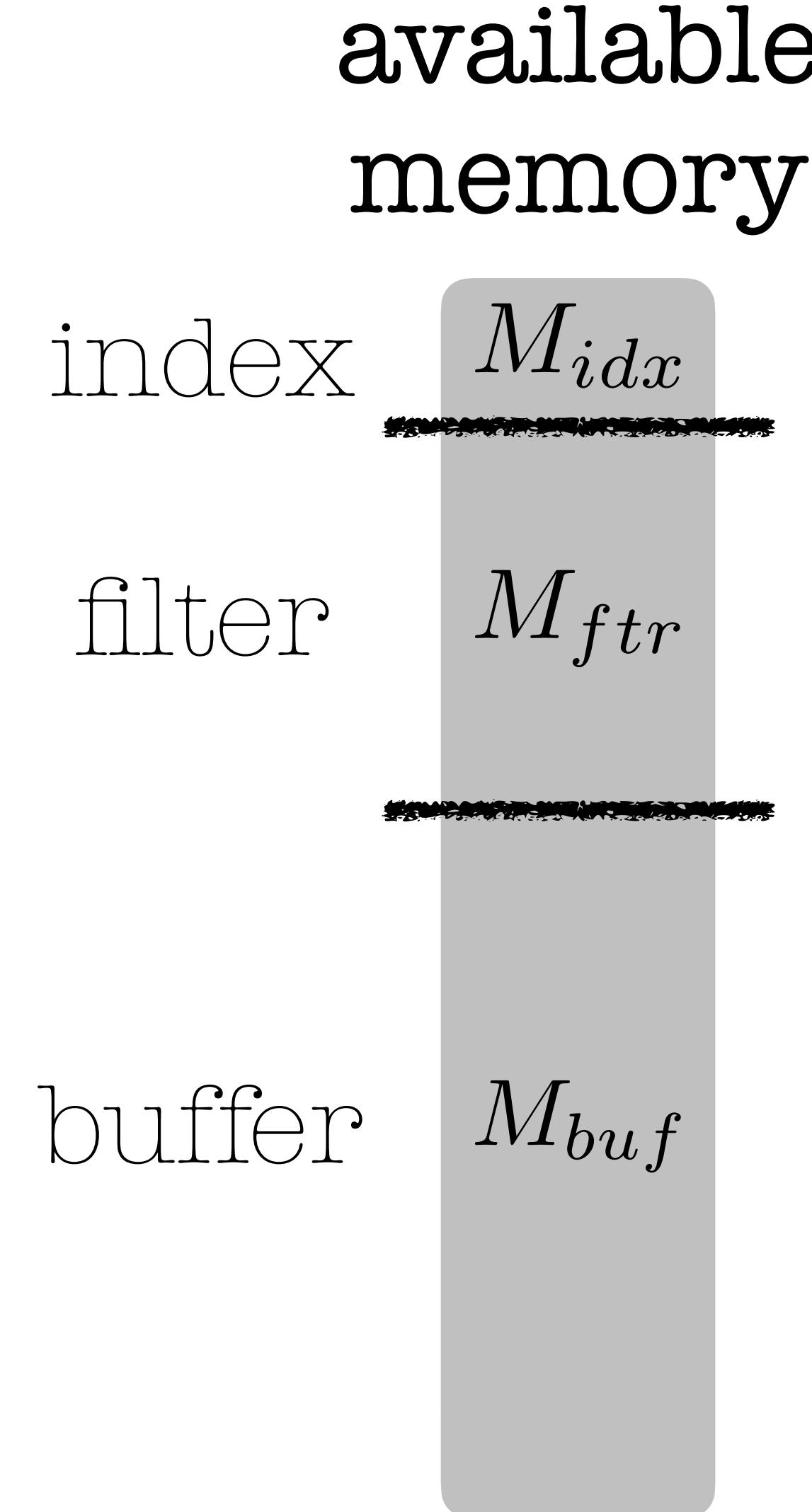
M_{idx} : index memory

M_{ftr} : filter memory

M_{buf} : buffer memory

The Optimal Memory Allocation

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



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

$read_cost(M_{idx}, M_{ftr}, M_{buf})$

$write_cost(M_{buf})$

$$cost = R \cdot read_cost + W \cdot write_cost$$

M : total memory

M_{idx} : index memory

M_{ftr} : filter memory

M_{buf} : buffer memory

M_{cache} : block cache
memory

workload

$reads(R)$
vs.

$writes(W)$

available
memory

block
cache

buffer



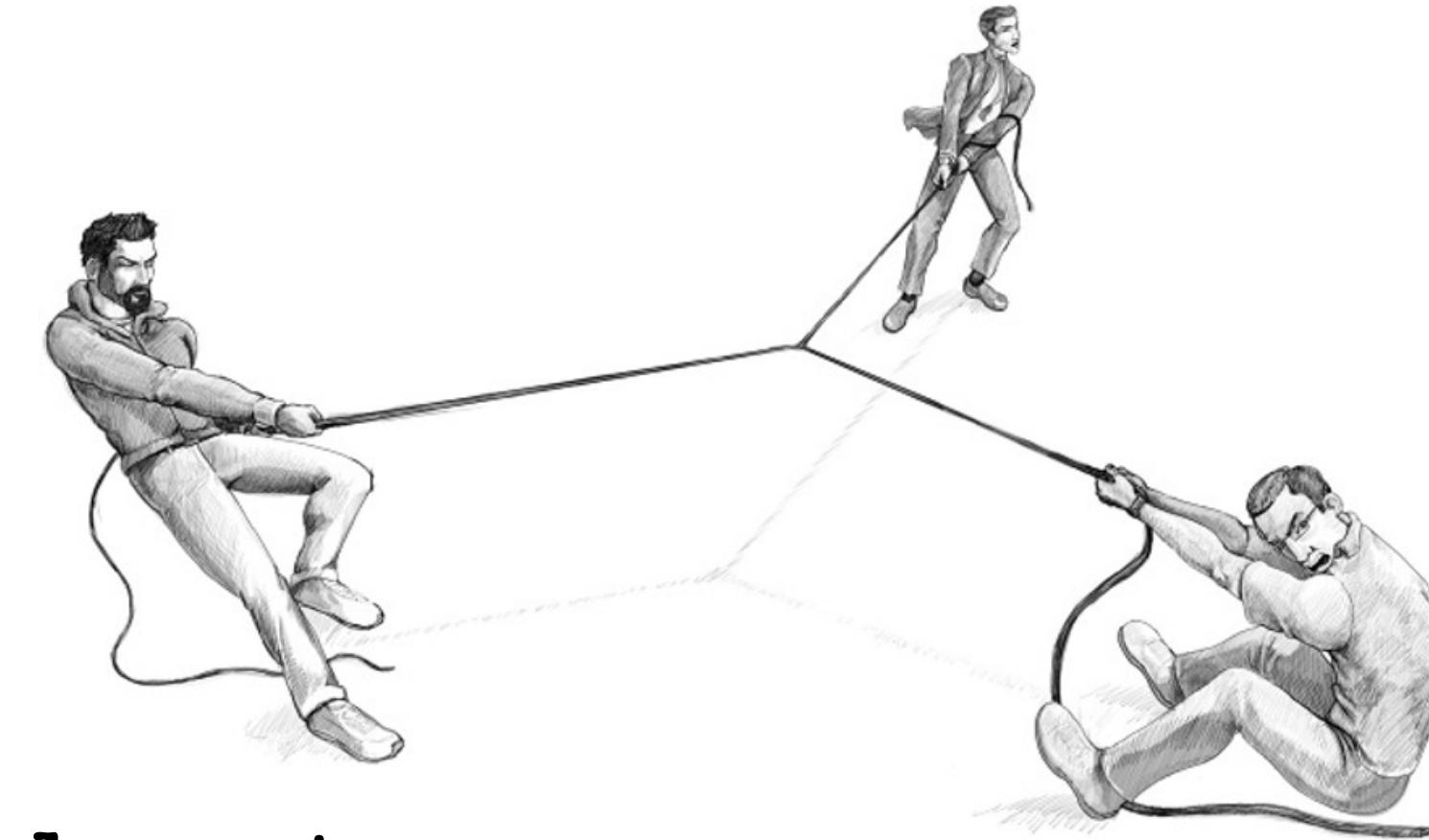
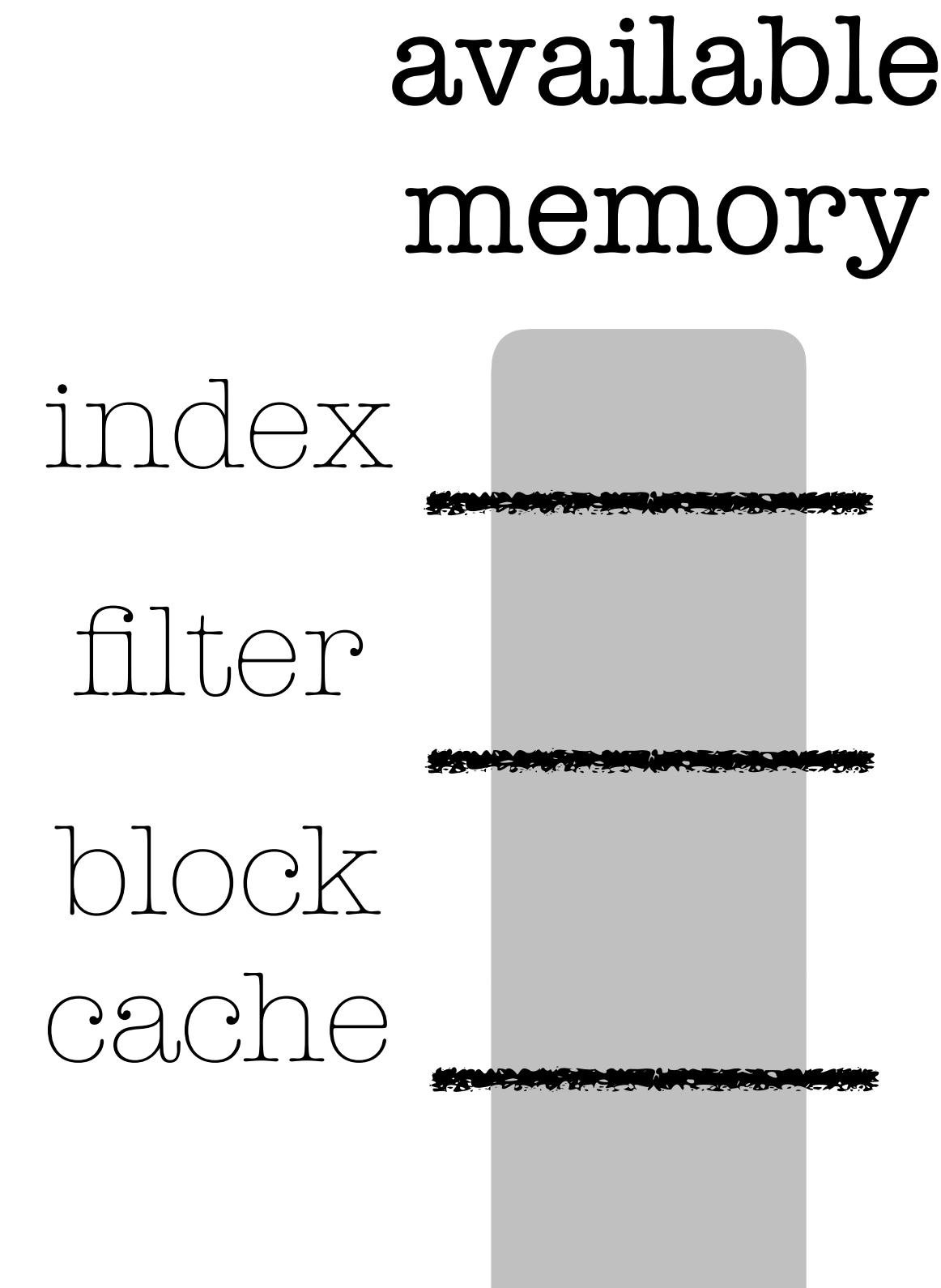
$$M = M_{cache} + M_{buf}$$

$read_cost(M_{cache})$

$write_cost(M_{buf})$

$$cost = R \cdot read_cost + W \cdot write_cost$$

The Optimal Memory Allocation



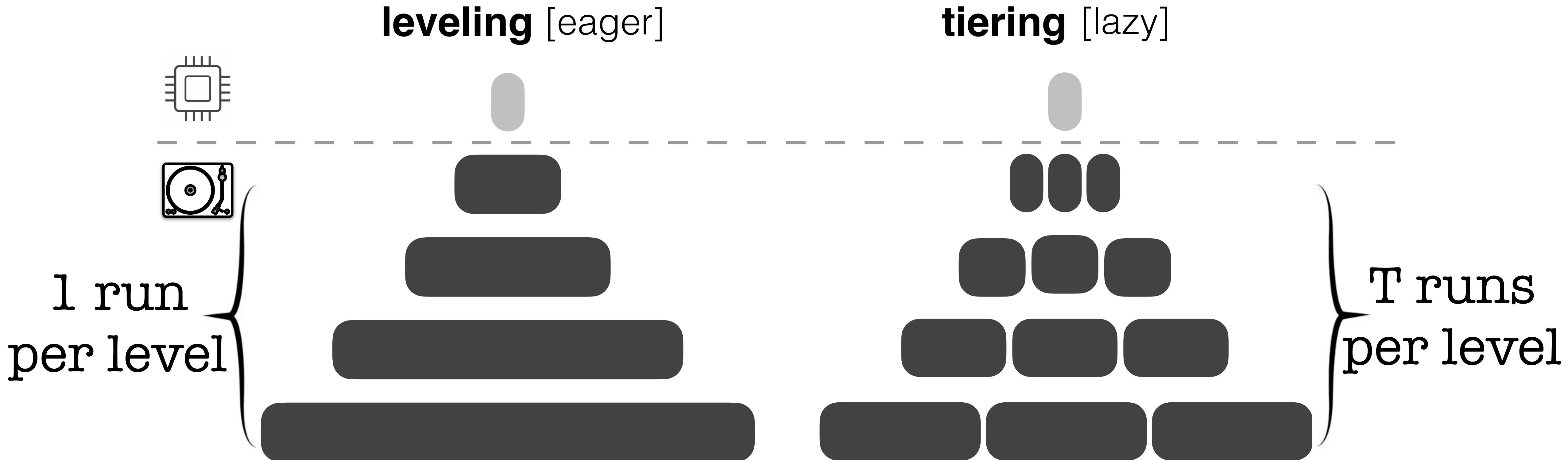
Update cost

Read cost

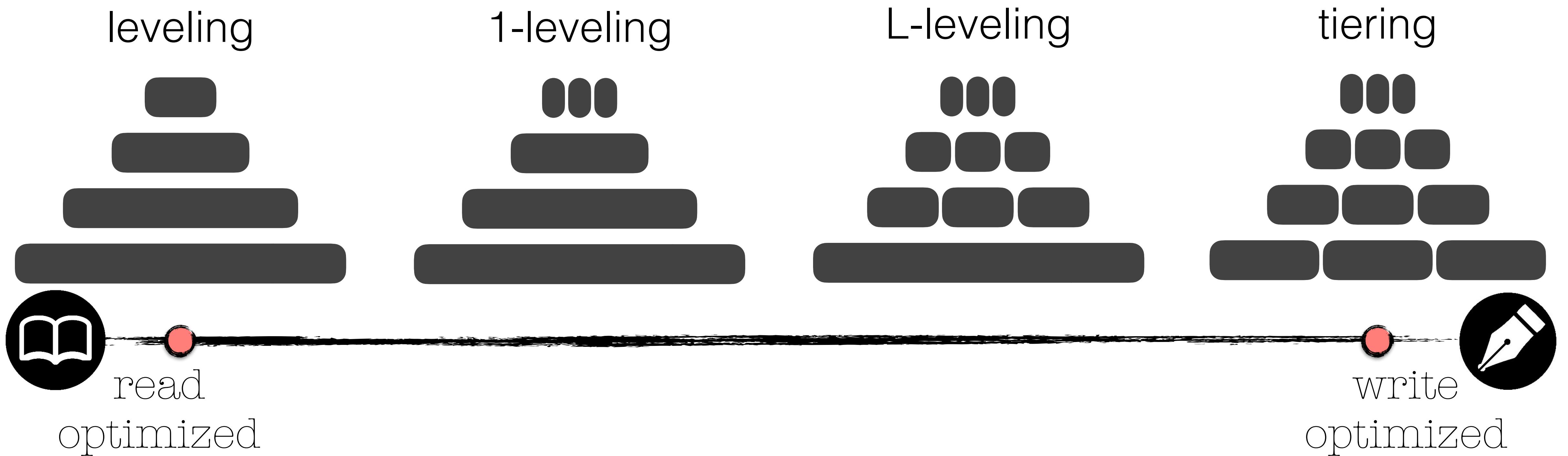
Memory/space
footprint

Navigating read vs. writes: data layouts

Data Layout

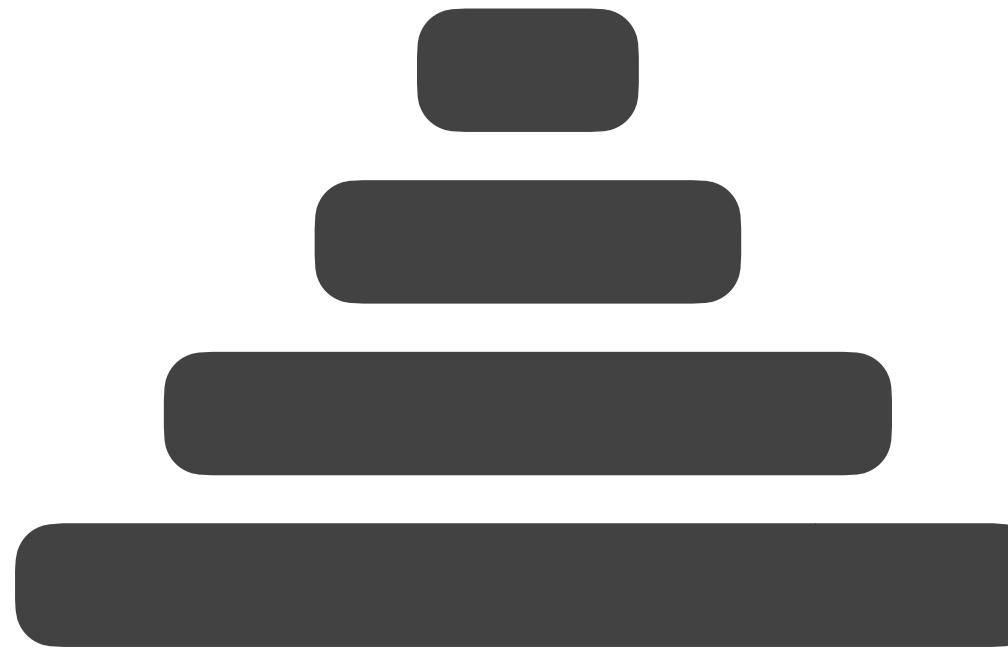


Data Layout

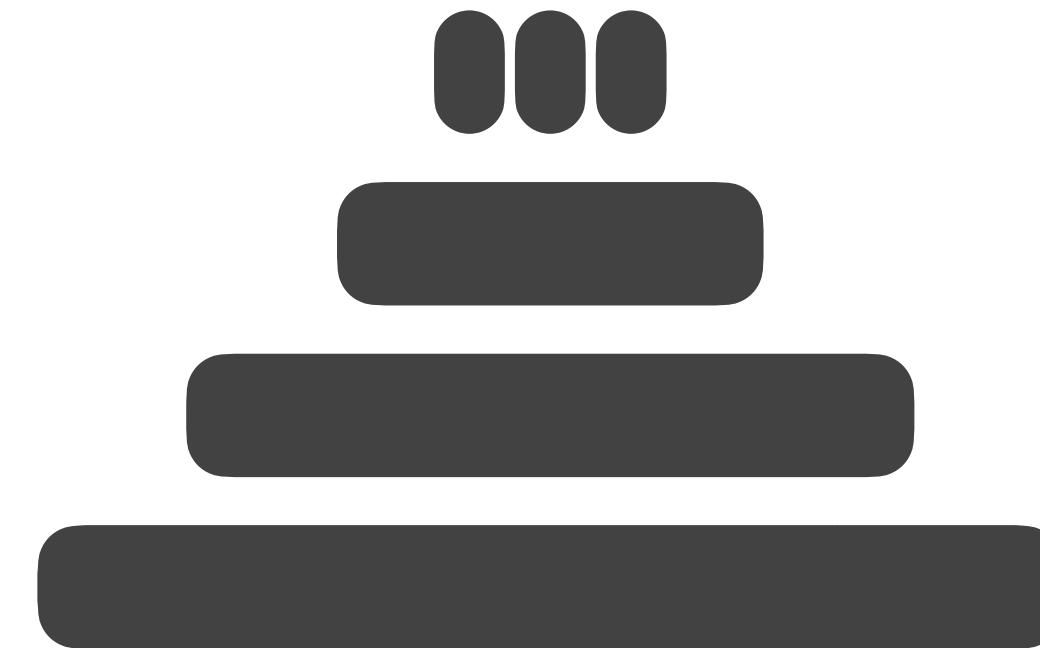


Storage Layer Design Continuum

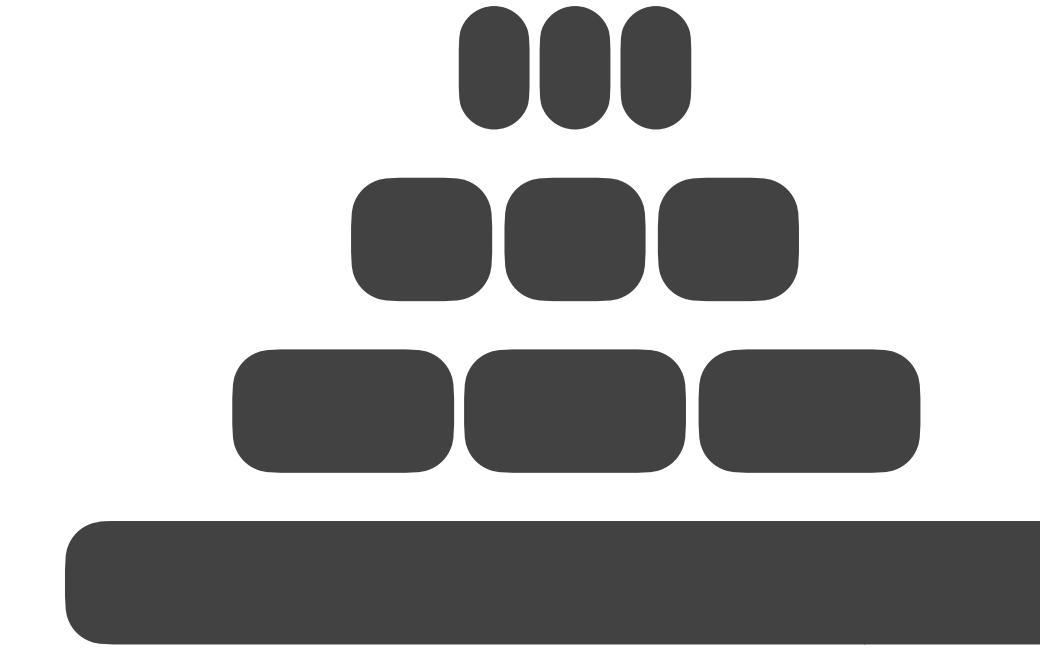
leveling



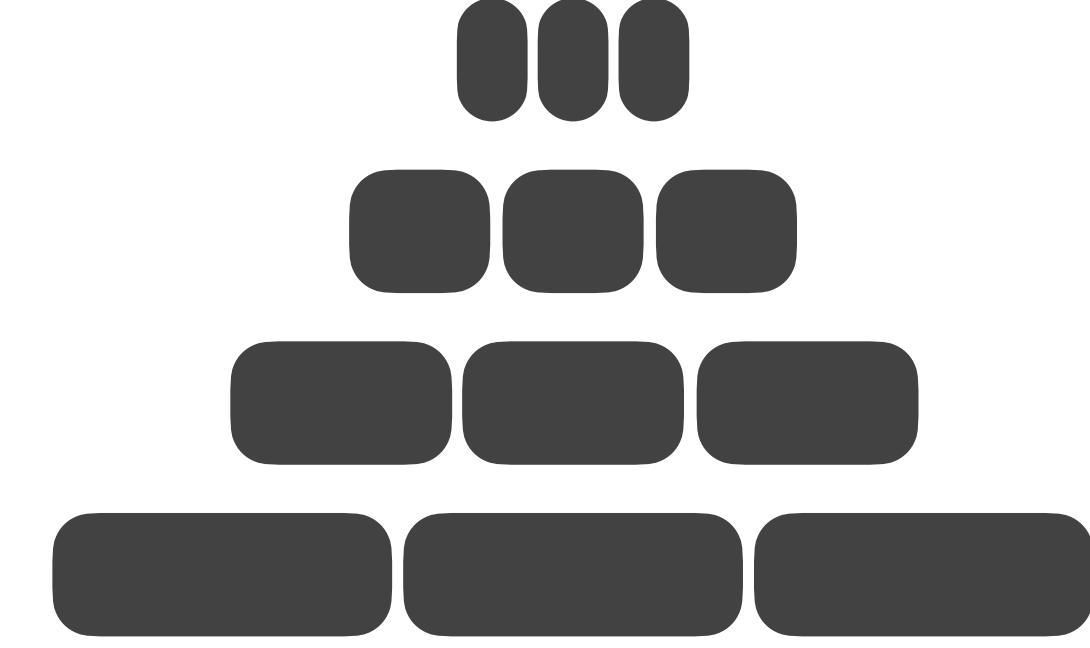
1-leveling



L-leveling



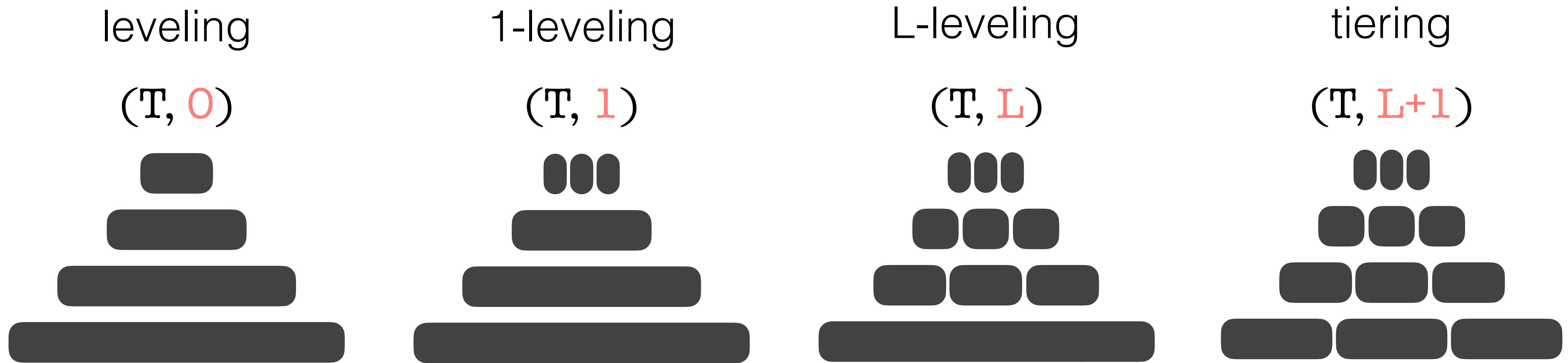
tiering



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



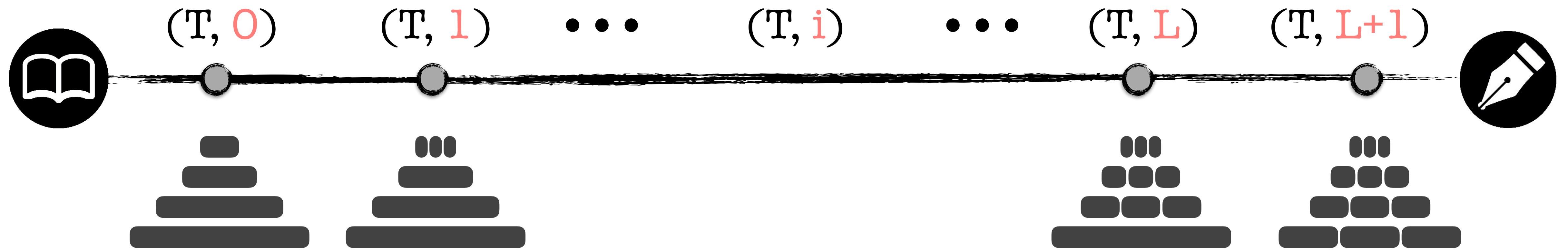
Storage Layer Design Continuum



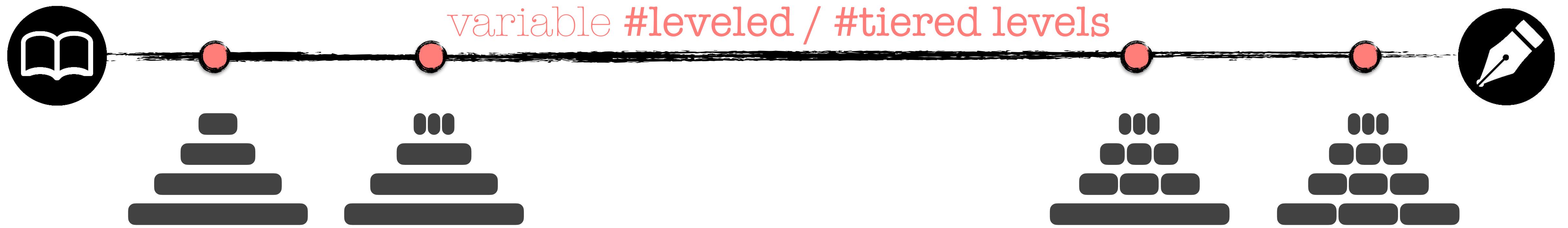
Any design can be defined by the tuple-set: (T, i)



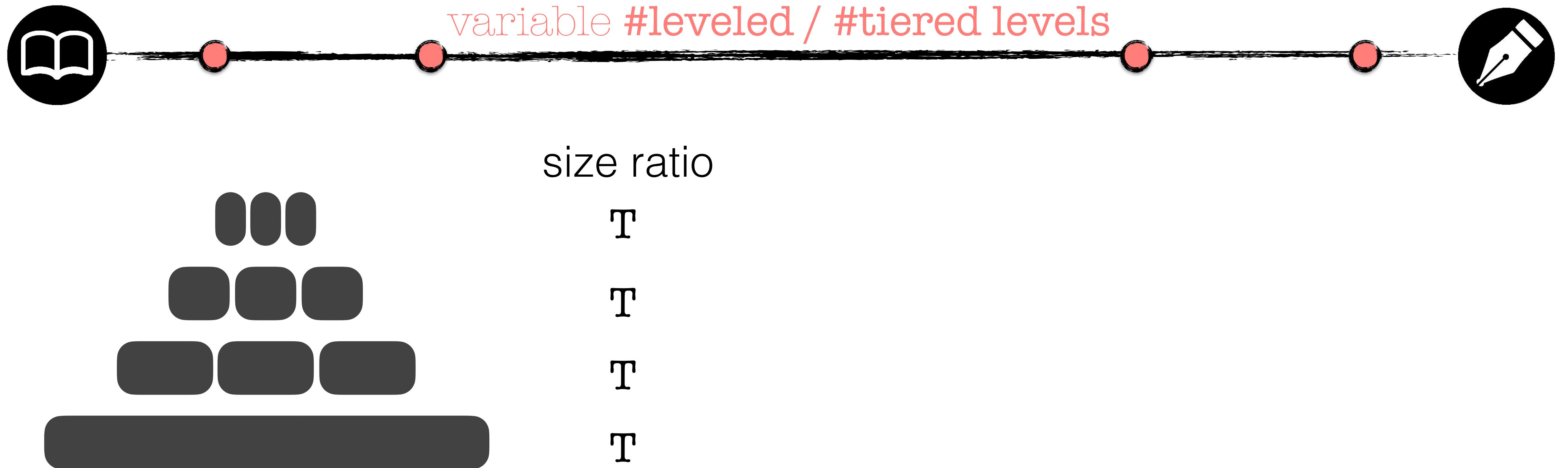
Storage Layer Design Continuum



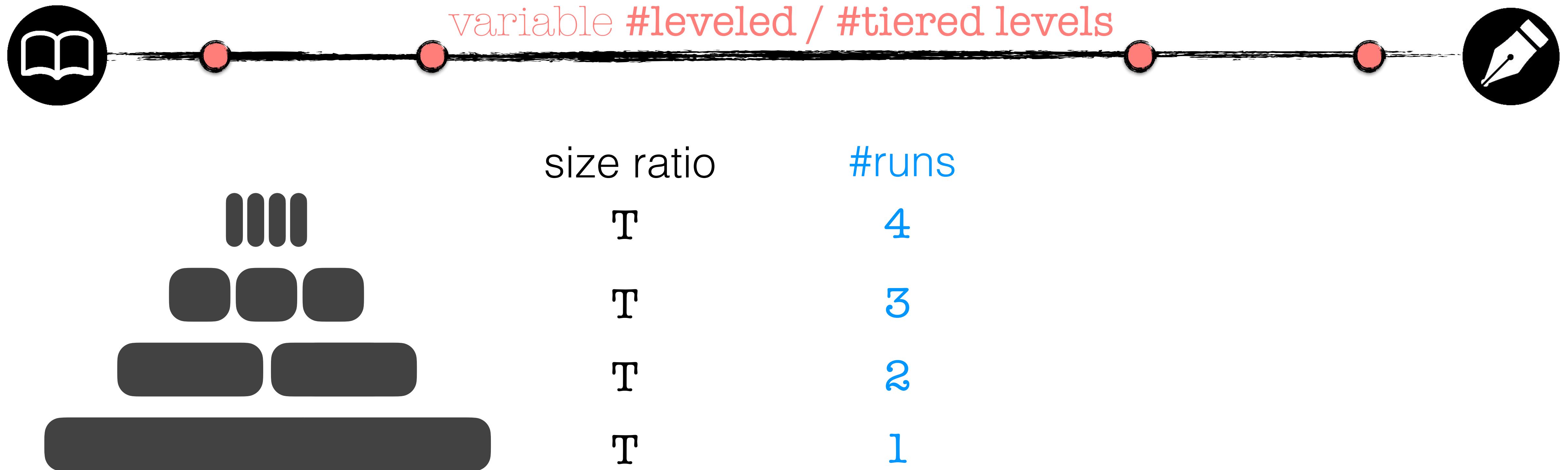
Storage Layer Design Continuum



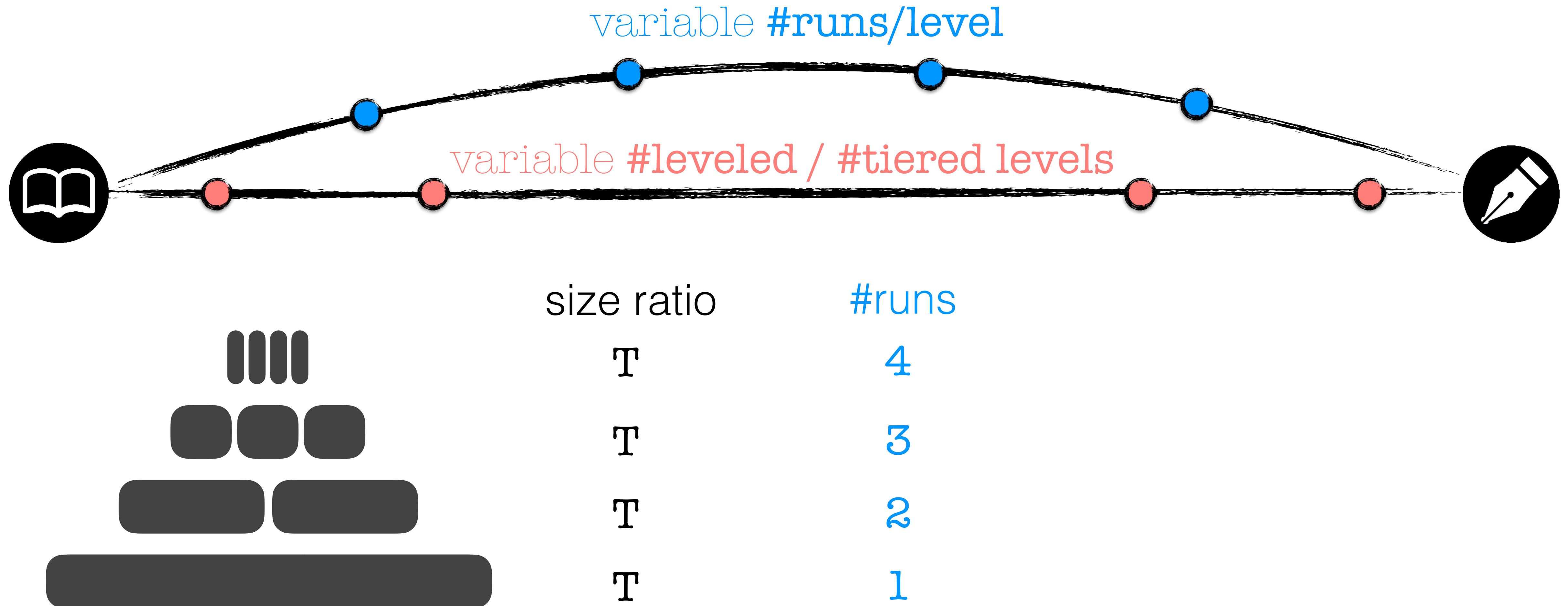
Storage Layer Design Continuum



Storage Layer Design Continuum



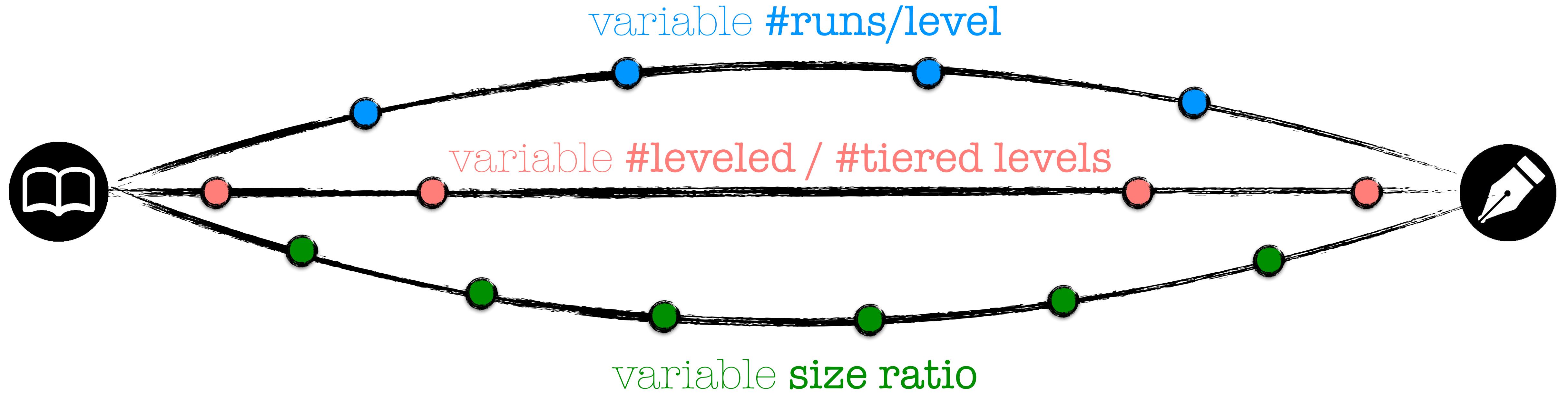
Storage Layer Design Continuum



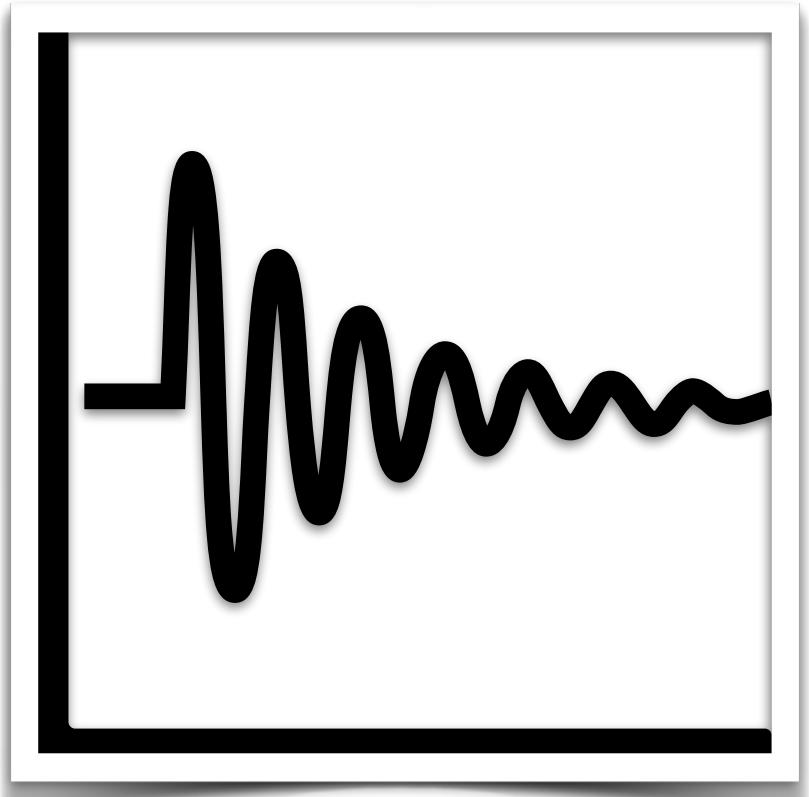
Storage Layer Design Continuum



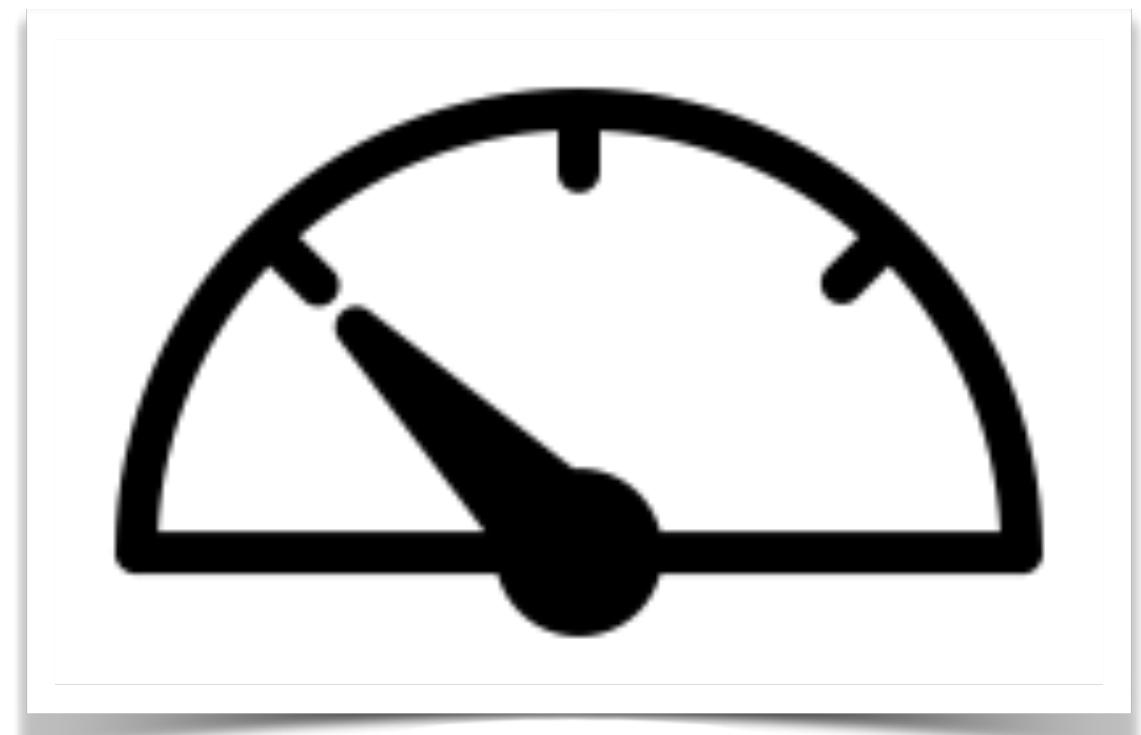
Storage Layer Design Continuum



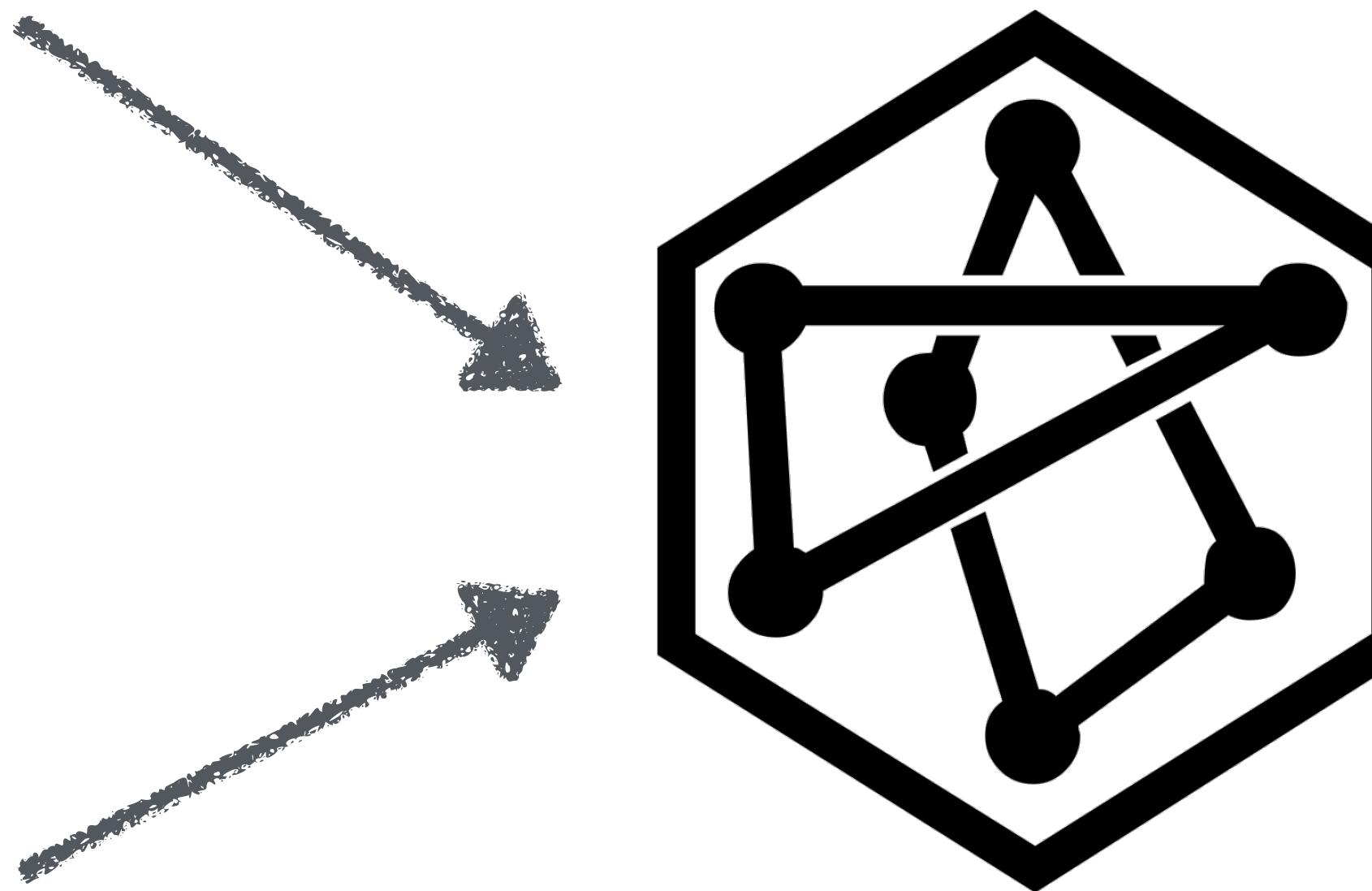
The LSM storage layer
design continuum



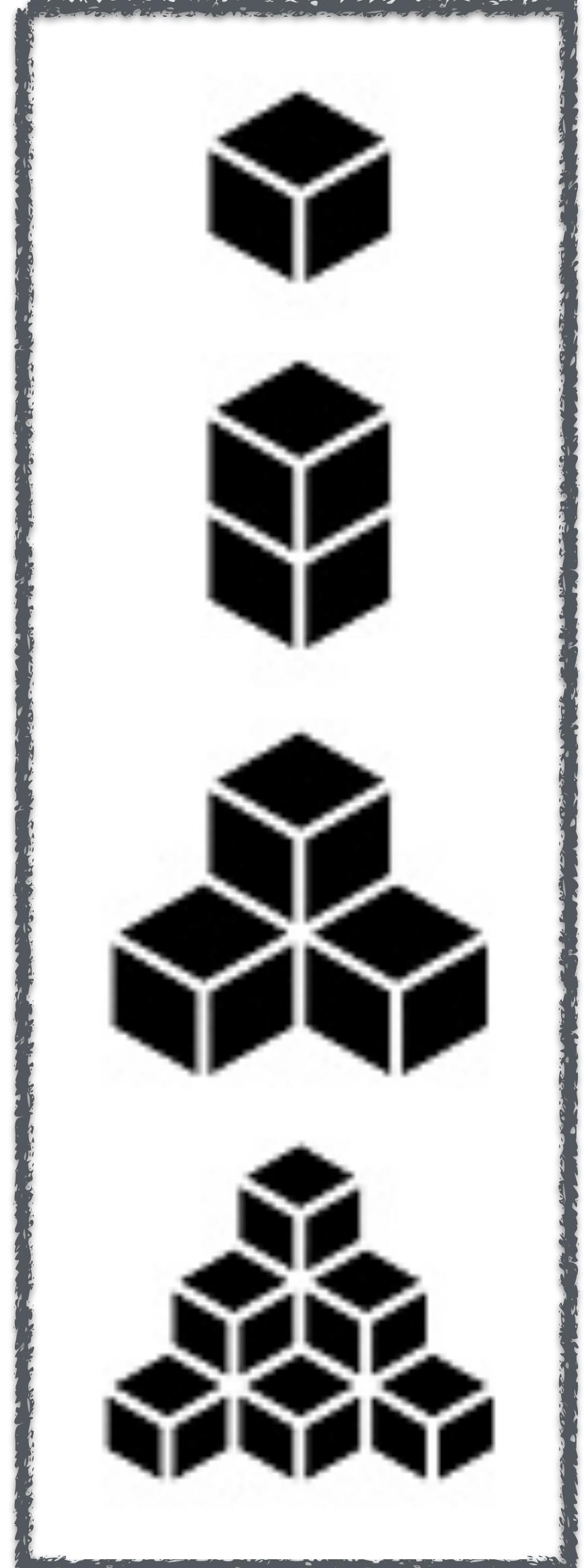
workload



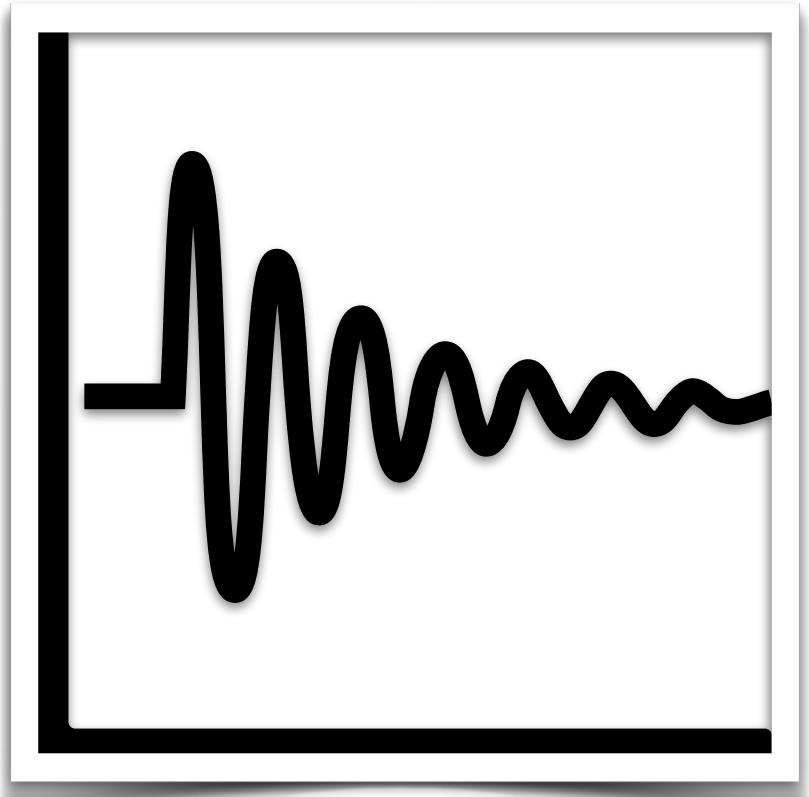
performance
target



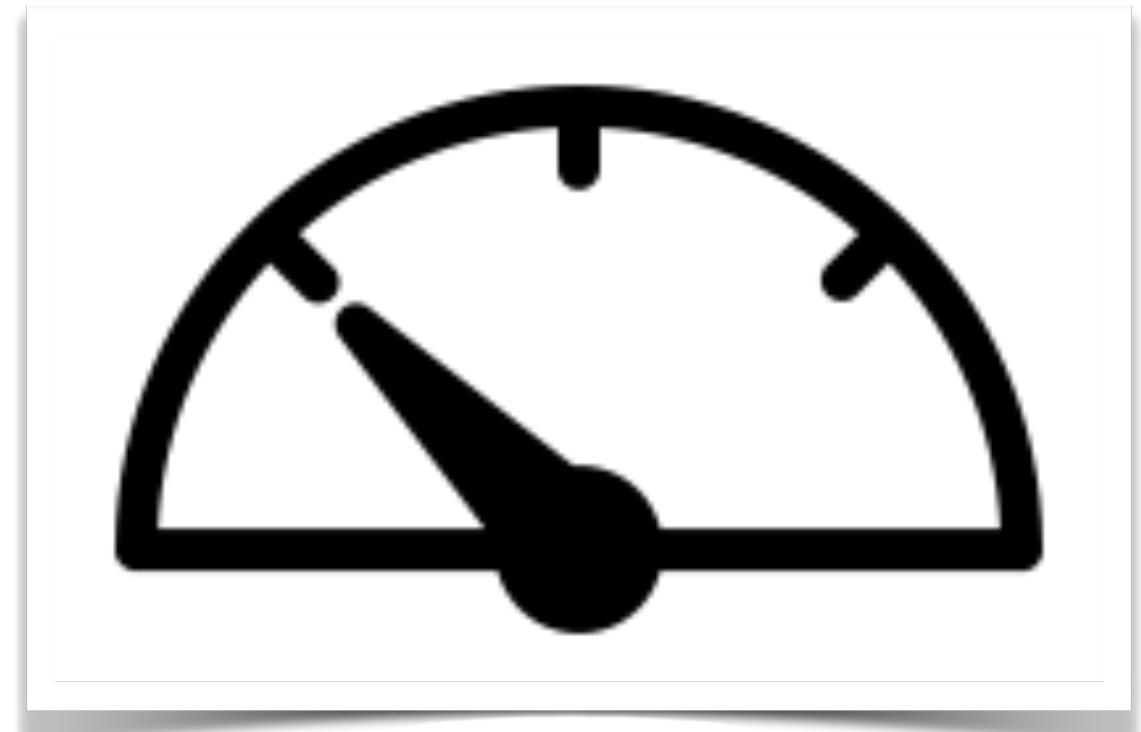
performance
modeling



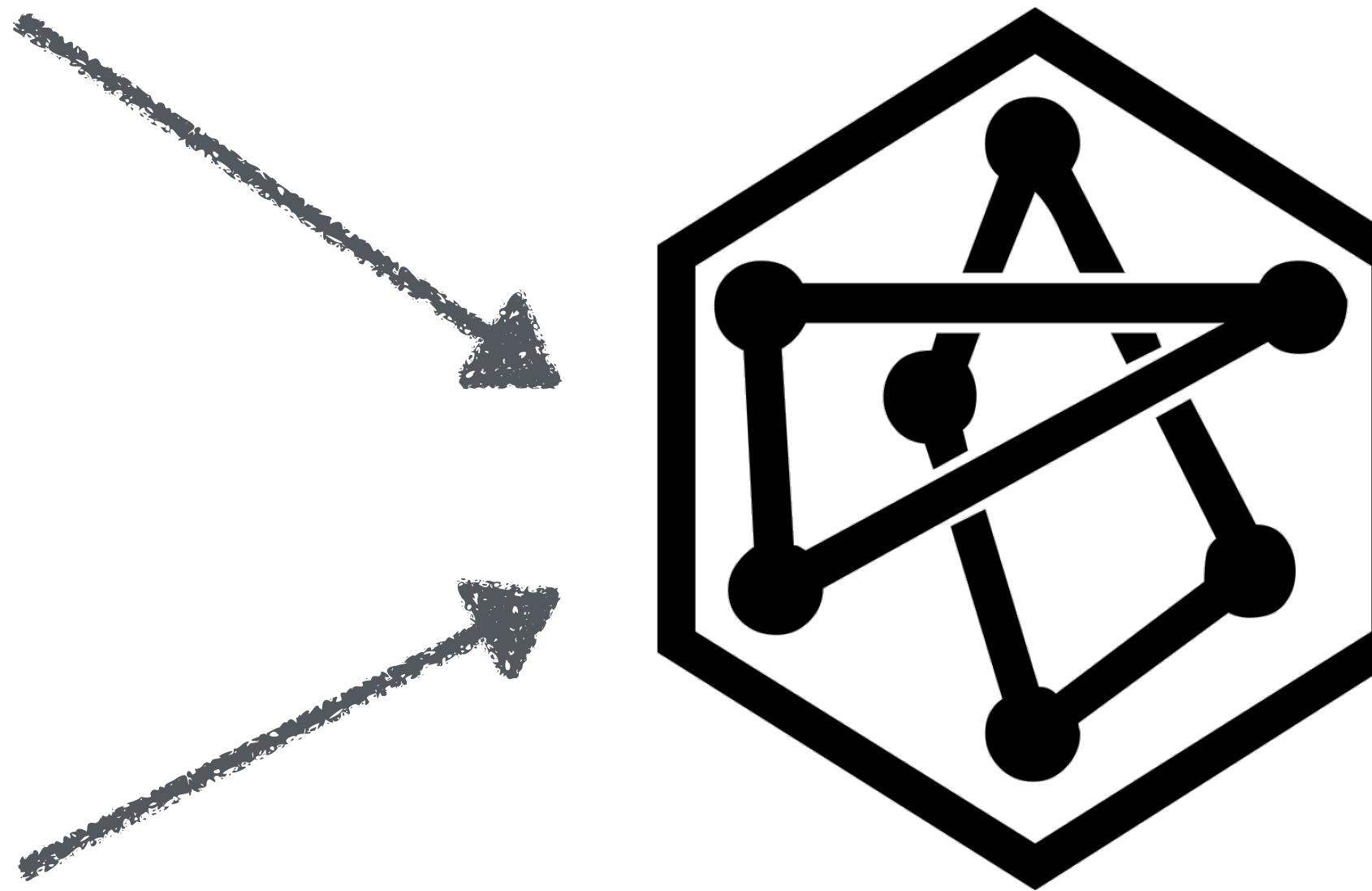
LSM designs



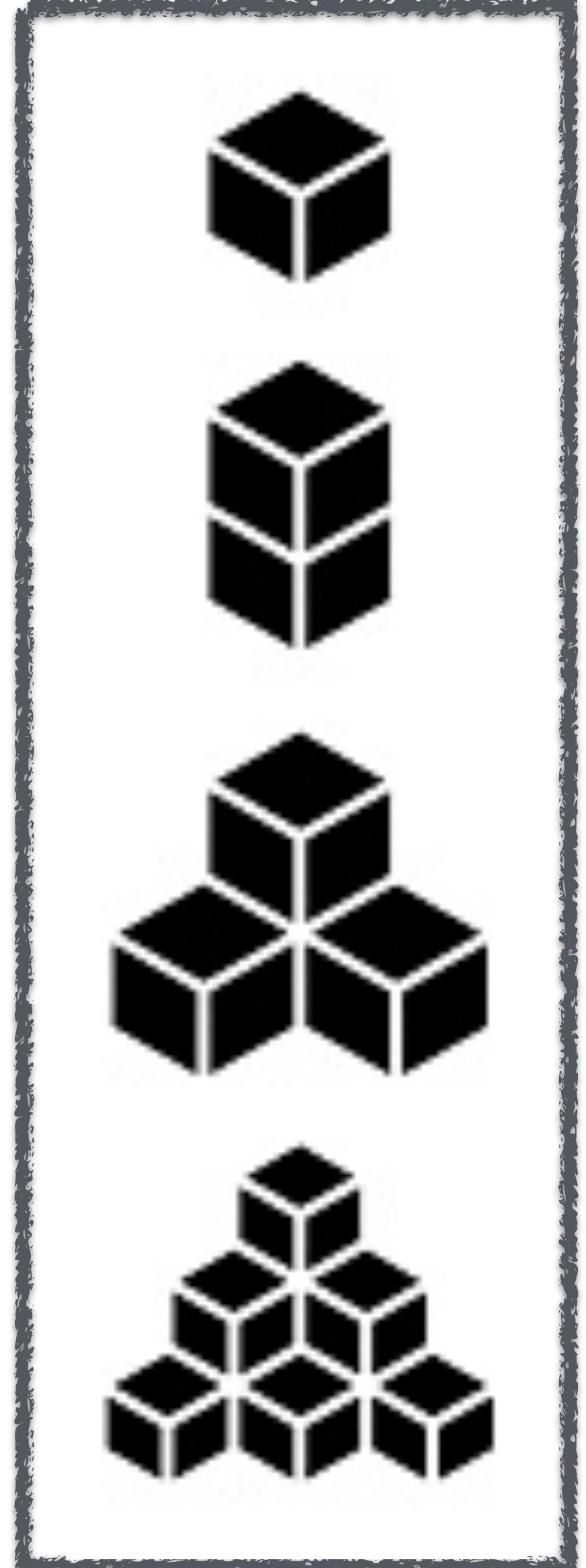
workload



performance
target



worst-case
performance
modeling

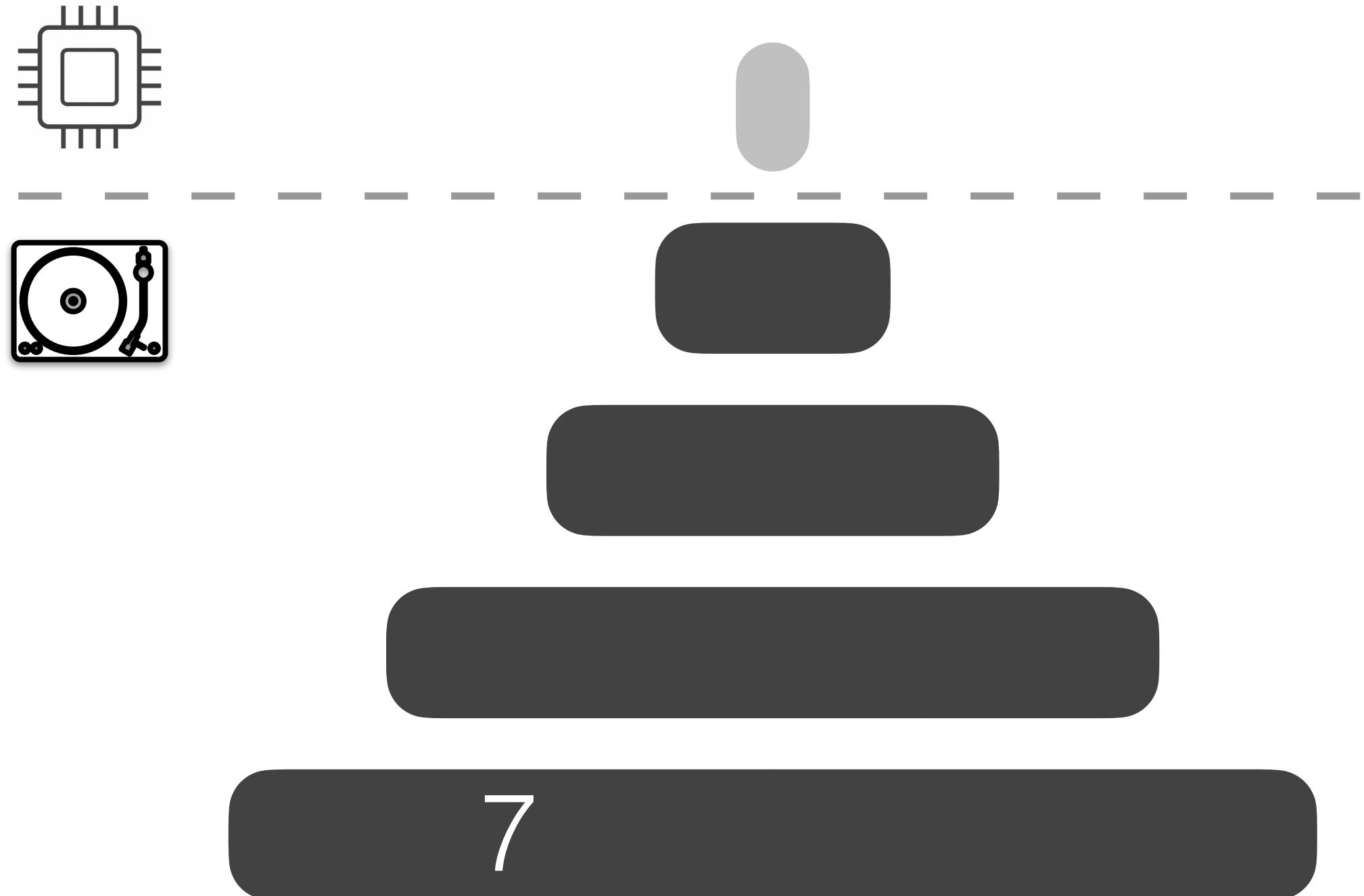


LSM designs

L : #levels

ϕ : FPR of BF

worst-case performance modeling

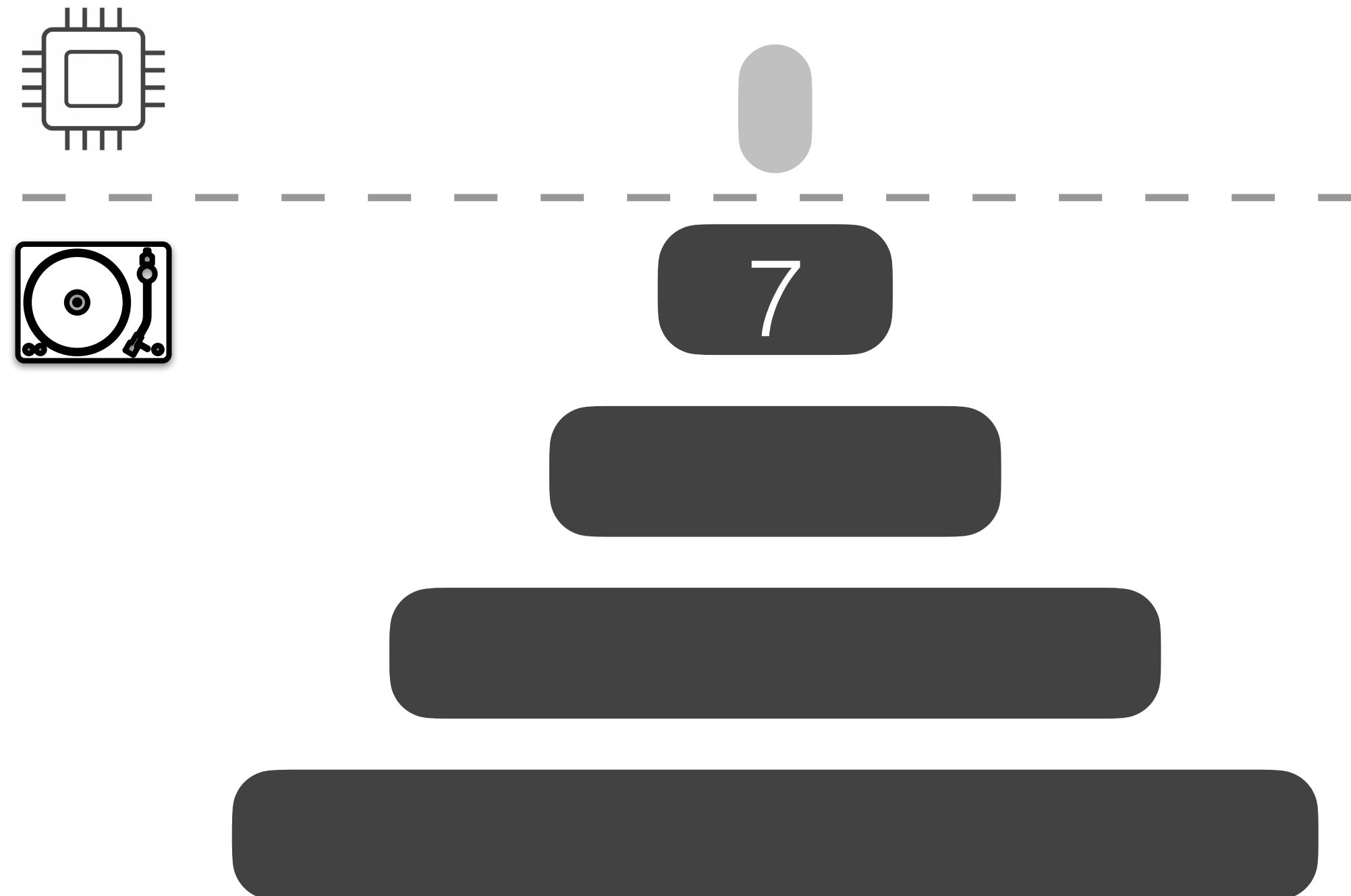


worst-case read cost: $1 + \sum_{i=1}^{L-1} \phi_i$

L : #levels

ϕ : FPR of BF

worst-case performance modeling

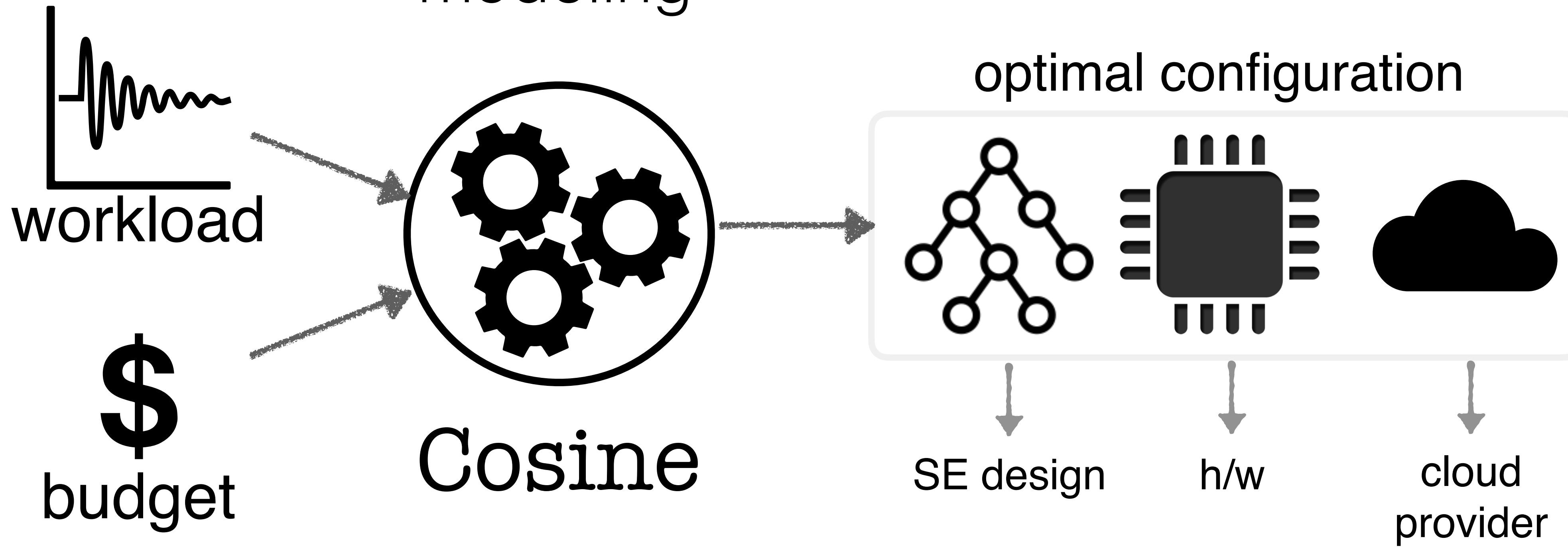


worst-case read cost: $1 + \sum_{i=1}^{L-1} \phi_i$

average-case performance modeling

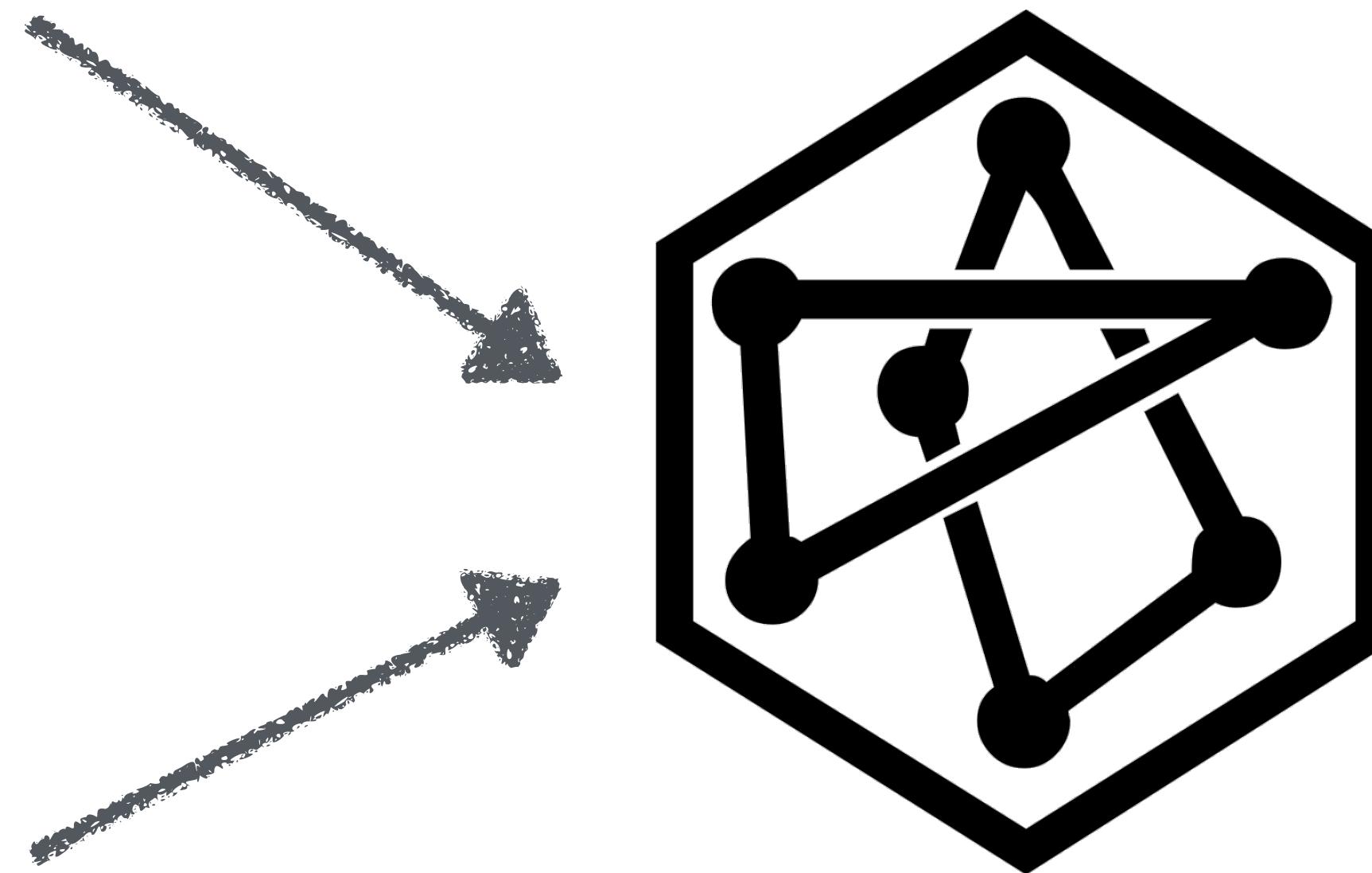
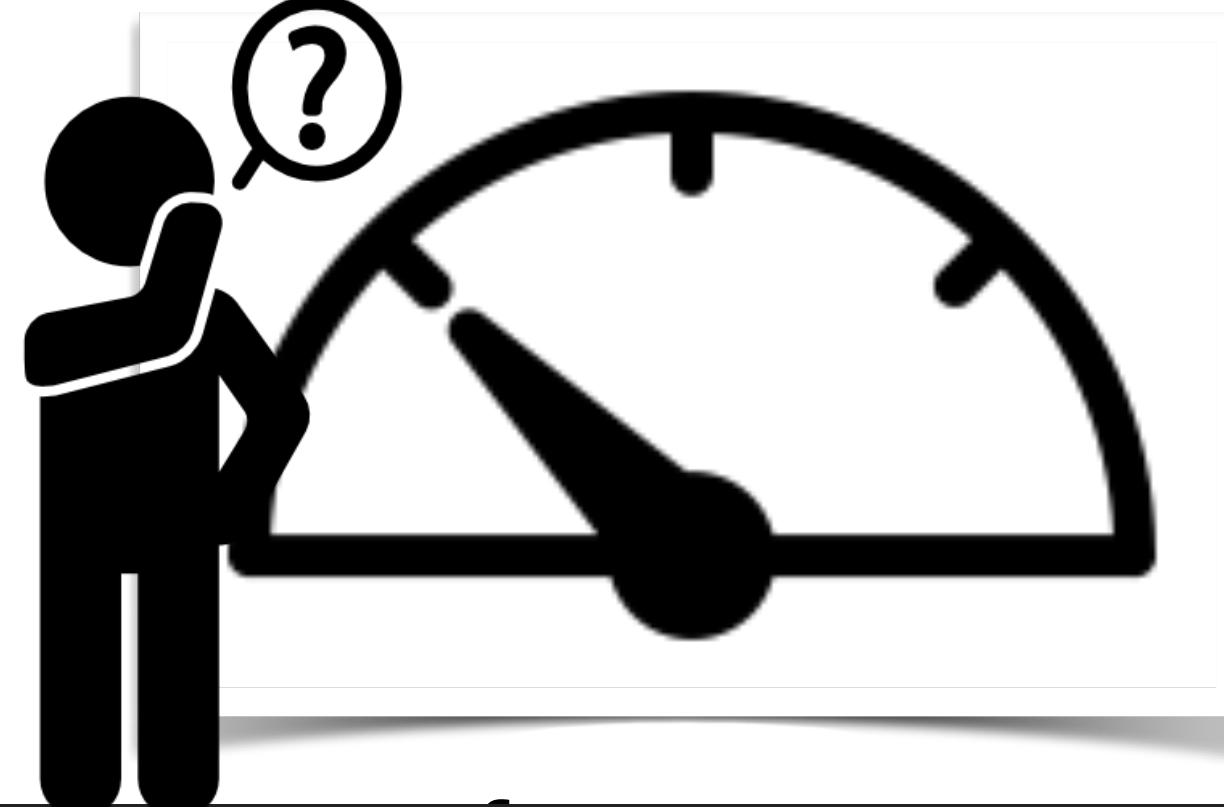
$$\sum_{i=1}^L (\mathbb{P}[\text{query in } L_i] \cdot (1 + \sum_{j=1}^{i-1} \phi_i))$$

average-case performance modeling

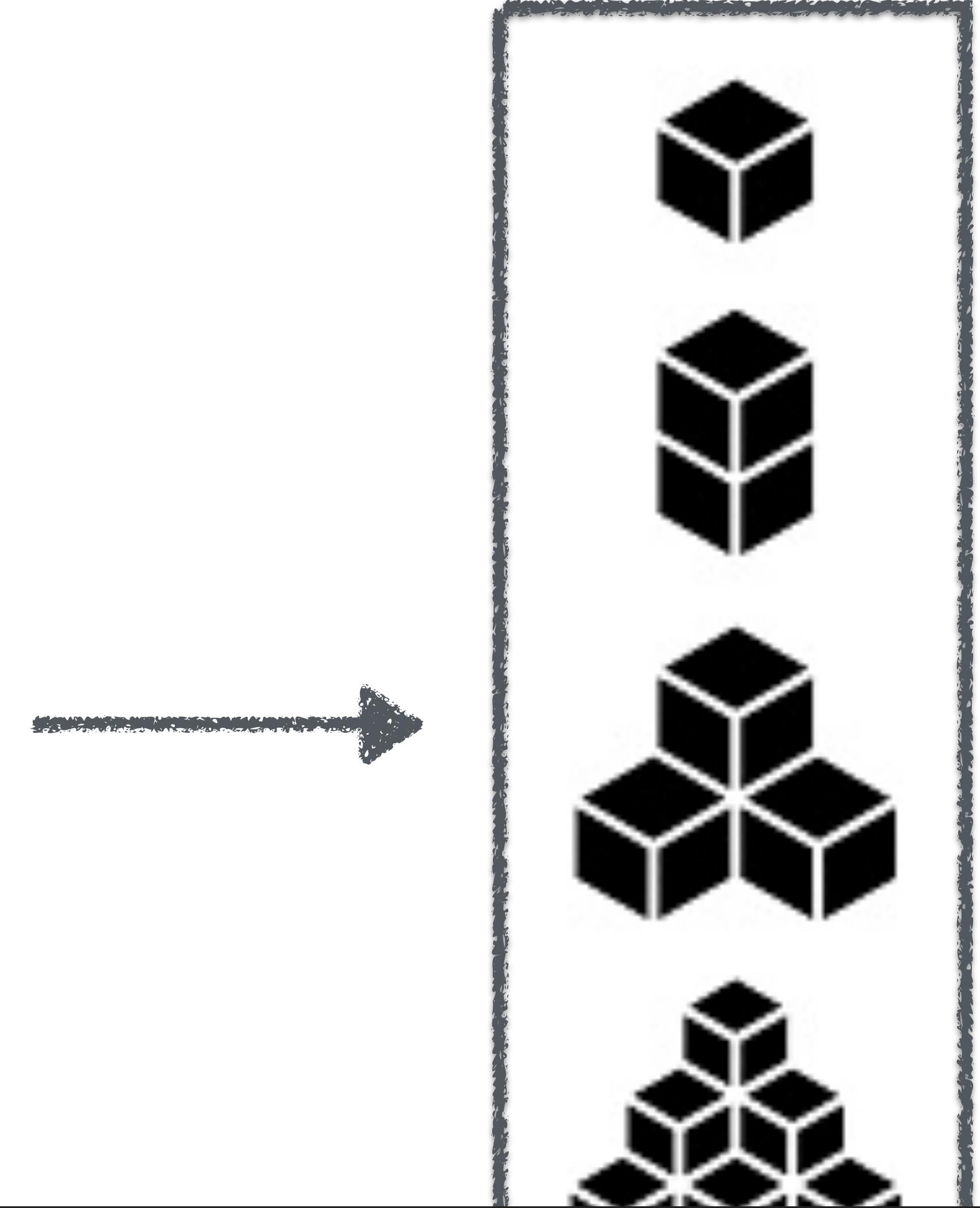




workload

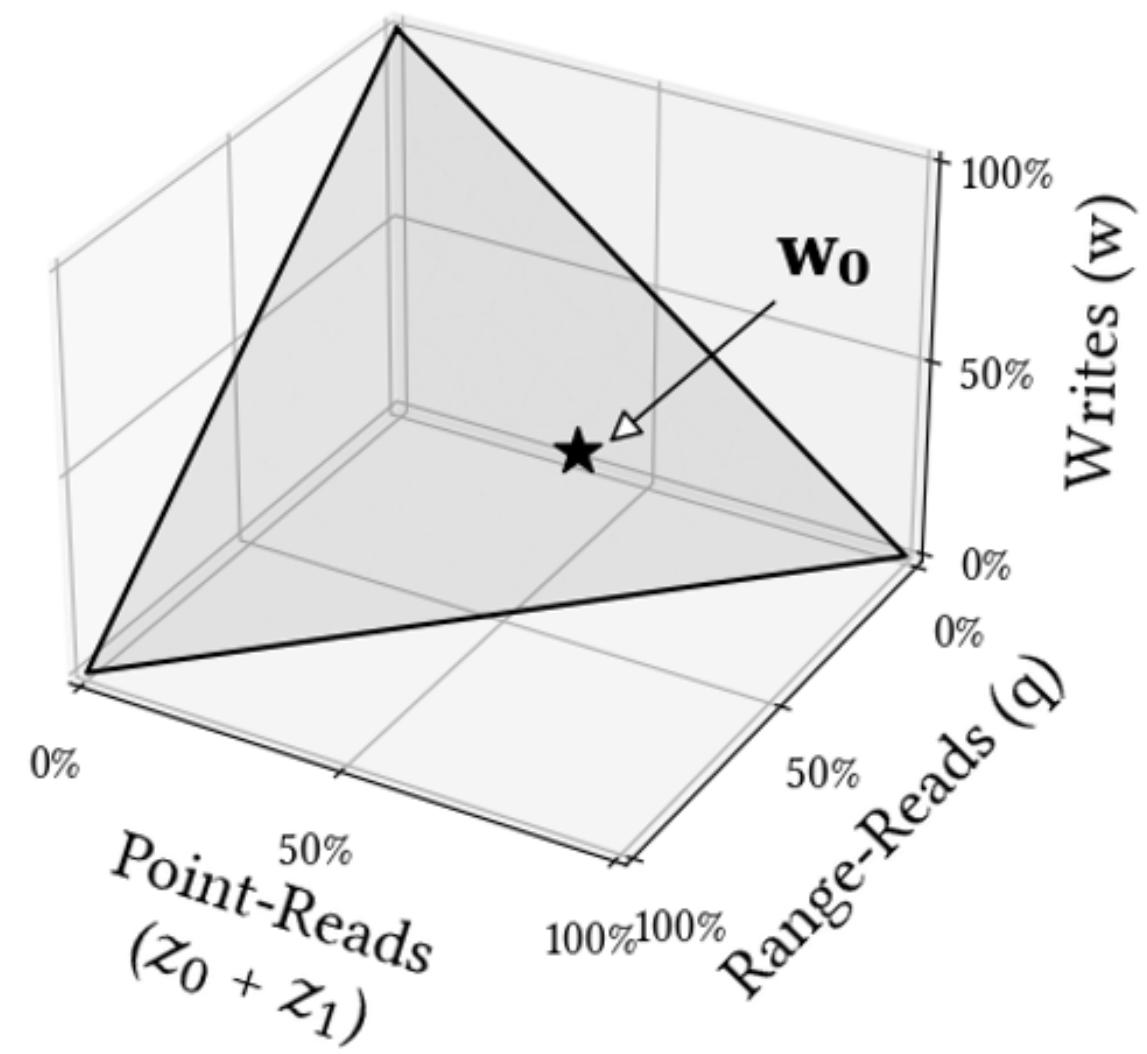


average-case



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

Workload-based Tuning

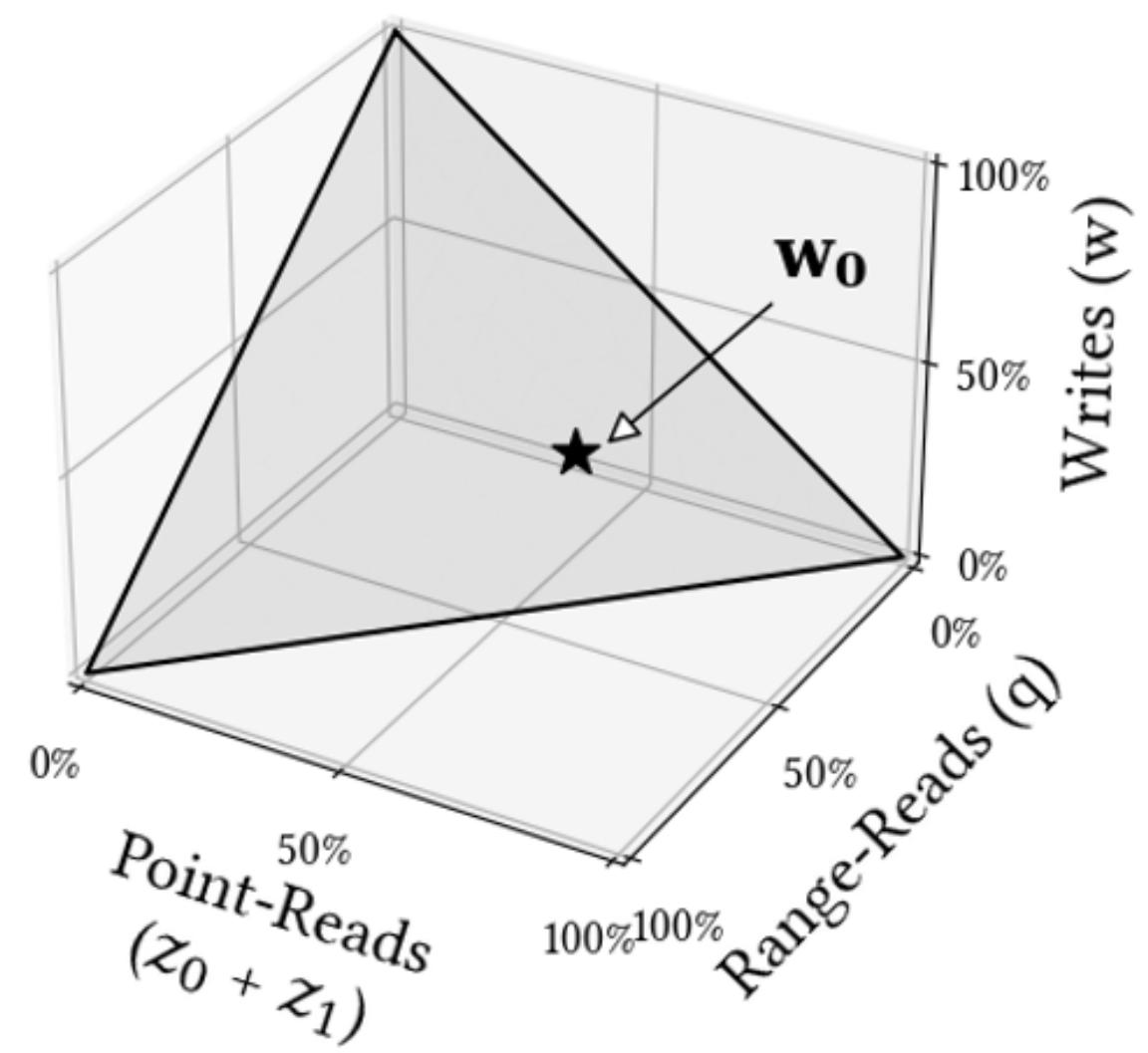


$$\pi_{w_0} = \operatorname{argmin}_{\pi} (\operatorname{cost}(w_0, \pi))$$

optimal configuration for w_0

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

Nominal Tuning



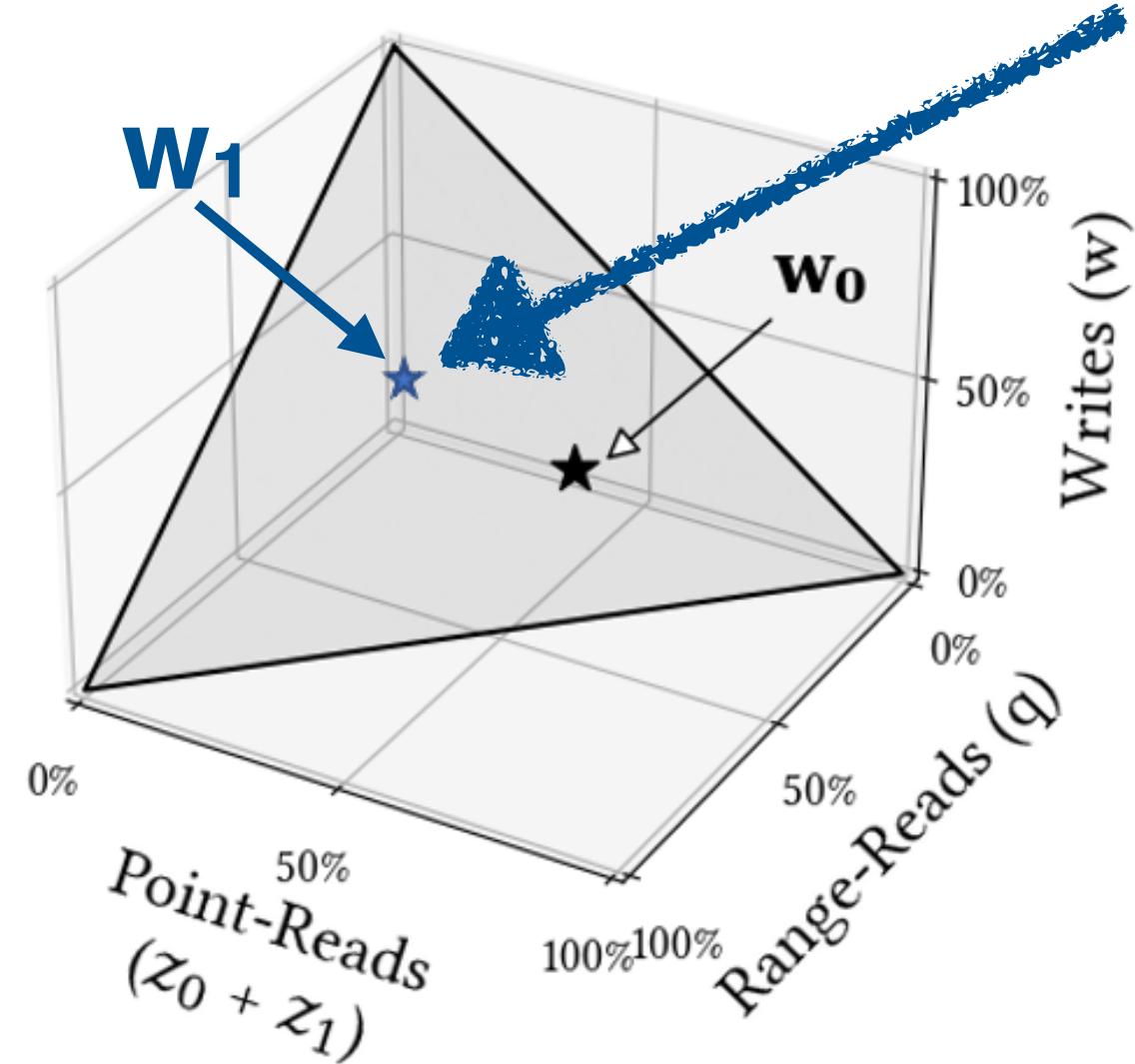
$$\pi_{w_0} = \operatorname{argmin}_{\pi} (\operatorname{cost}(w_0, \pi))$$

optimal configuration for w_0

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

same configuration

due to unpredictability

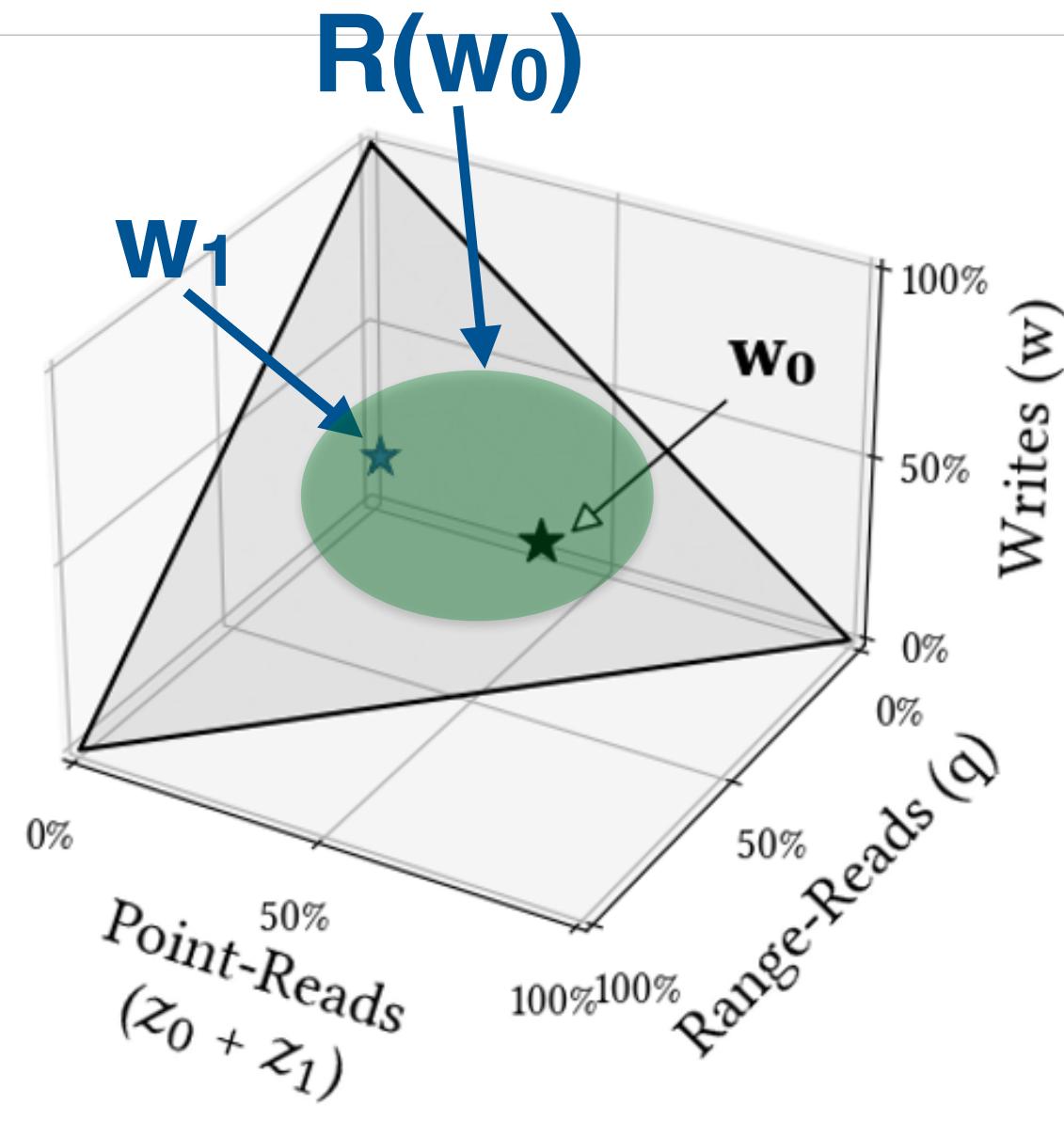


$$\pi_{w_0}$$

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

... but not optimal!

Robust Tuning

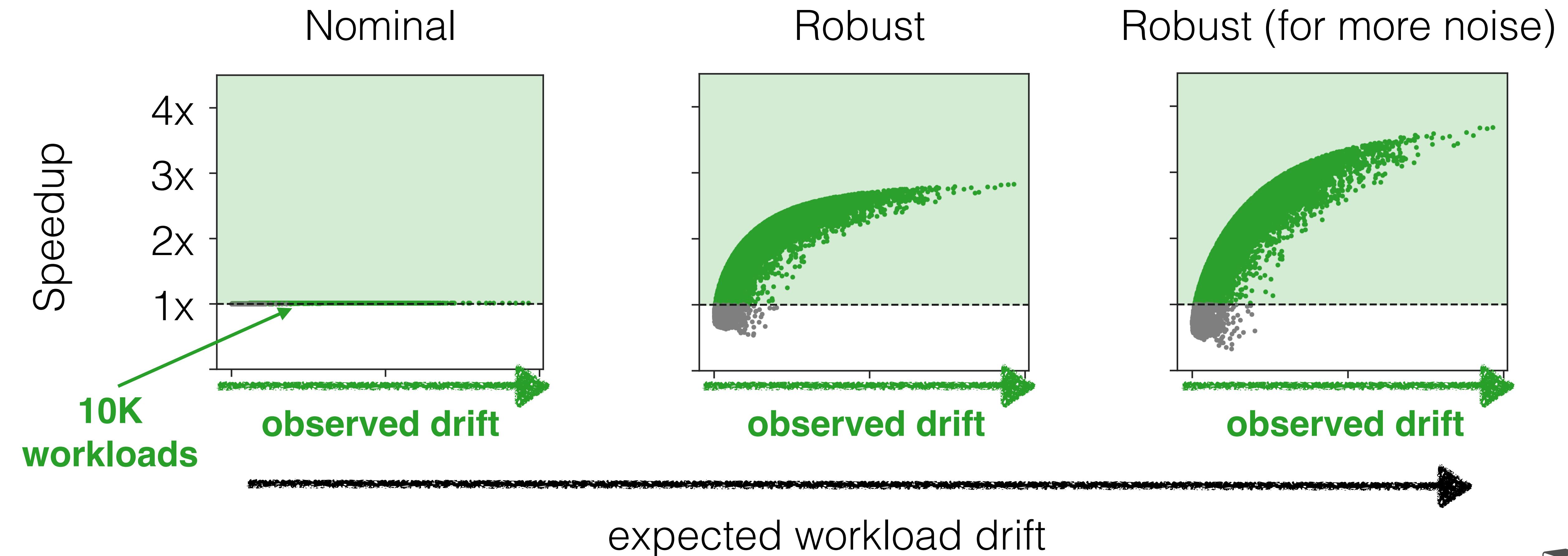


$$\pi_{w_0}^{robust} = \operatorname{argmin}_{\pi} \max_{w' \in R(w_0)} (\text{cost}(w_0, \pi_{w_0}^{robust}))$$

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

... close-to-optimal!

Robust Tuning



The Key Takeaways

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

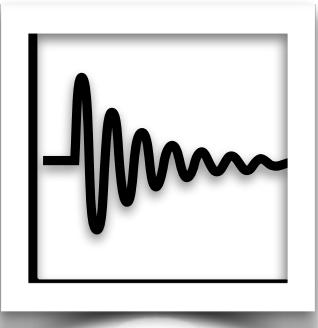
Read optimizations are crucial to make LSMs better.

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

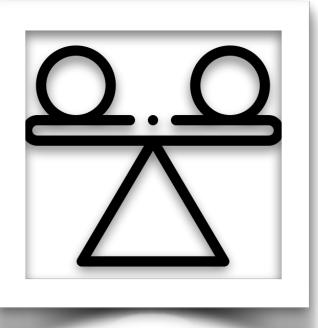
Open Research Challenges



Reduce write amplification



Workload-aware compactations & 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. Altinbüken, 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 System. In *VLDB Endowment*, 7(1).
- [3] Apache. Accumulo.
- [4] Apache. HBase.
- [5] Apache. Cassandra.
- [6] M. Athanassoulis. Merge Methods. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2017.
- [7] M. Athanassoulis, A. Ailamaki, and A. RUM Conjecture: Extending Data Structures to Handle Big Data. In *Proceedings of the VLDB Endowment*, 11(1).
- [8] O. Balmau, F. Medjedai, and D. Didona. Merge Key-Value Stores. In *Journal of Computer Systems, Storage and Trans. Comput.*, 2020.
- [9] O. Balmau, R. Donati, and M. Athanassoulis. Unlocking Memory Efficiency in Data Stores. In *Proceedings of the ACM European Conference on Data Management Systems*, pages 80–94, 2017.
- [10] M. A. Bender, O. Balmau, D. Medjedai, and D. Staratzis. Don't Thrash: Learning Index Structures. In *VLDB Endowment*, 11(1).
- [11] A. Broslow, and M. Athanassoulis. The Case for Learned Index Structures. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 489–504, 2018.
- [12] 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.
- [13] 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.
- [14] N. Dayan, T. Weiss, S. Dashevsky, M. Pan, E. Bortnikov, and M. Twitto. 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.
- [15] 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.
- [16] RocksDB. Leveled Compaction. <https://github.com/facebook/rocksdb/wiki/Leveled-Compaction>, 2020.
- [17] RocksDB. Prefix Bloom Filter. <https://github.com/facebook/rocksdb/wiki/Prefix-Seek#configure-prefix-bloom-filter>, 2020.
- [18] RocksDB. Block Cache. <https://github.com/facebook/rocksdb/wiki/Block-Cache>, 2021.
- [19] 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.
- [20] 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.
- [21] 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.
- [22] 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?



UNIVERSITY OF
TORONTO