



Dissecting, Designing, and Optimizing **LSM-Based Data Stores**

Subhadeep Sarkar

Manos Athanassoulis

Log-Structured Merge-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.*

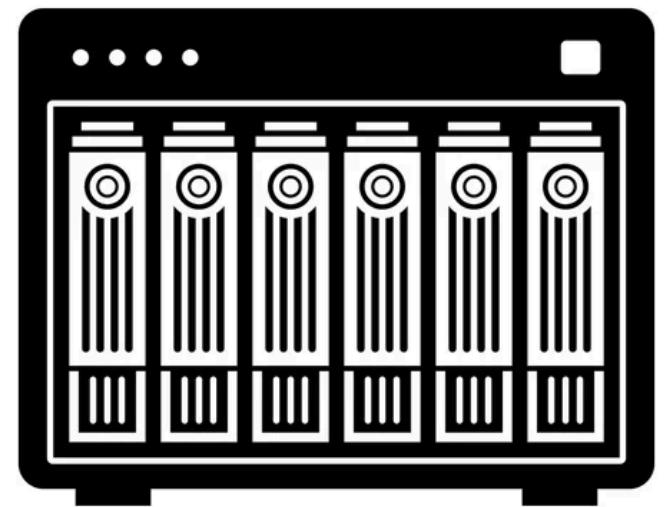


1996

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

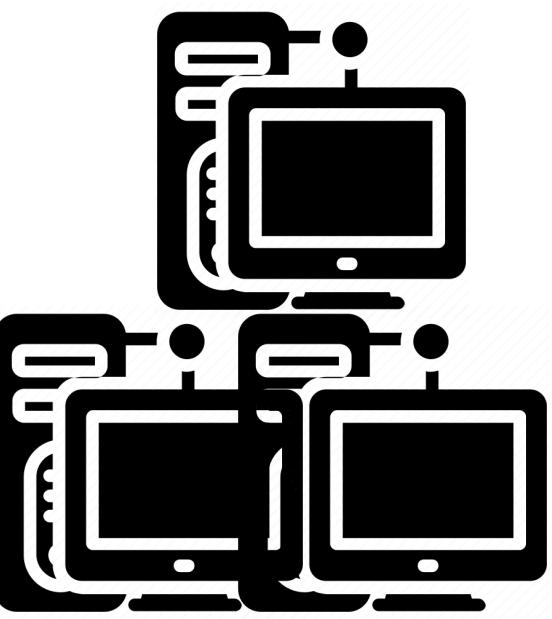
- good reads



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



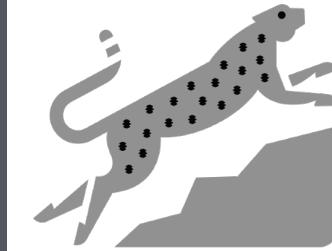
cassandra

2010



levelDB

2011



RocksDB

2013

LSM-tree

NoSQL



relational



time-series

2022

Why **LSM** ?

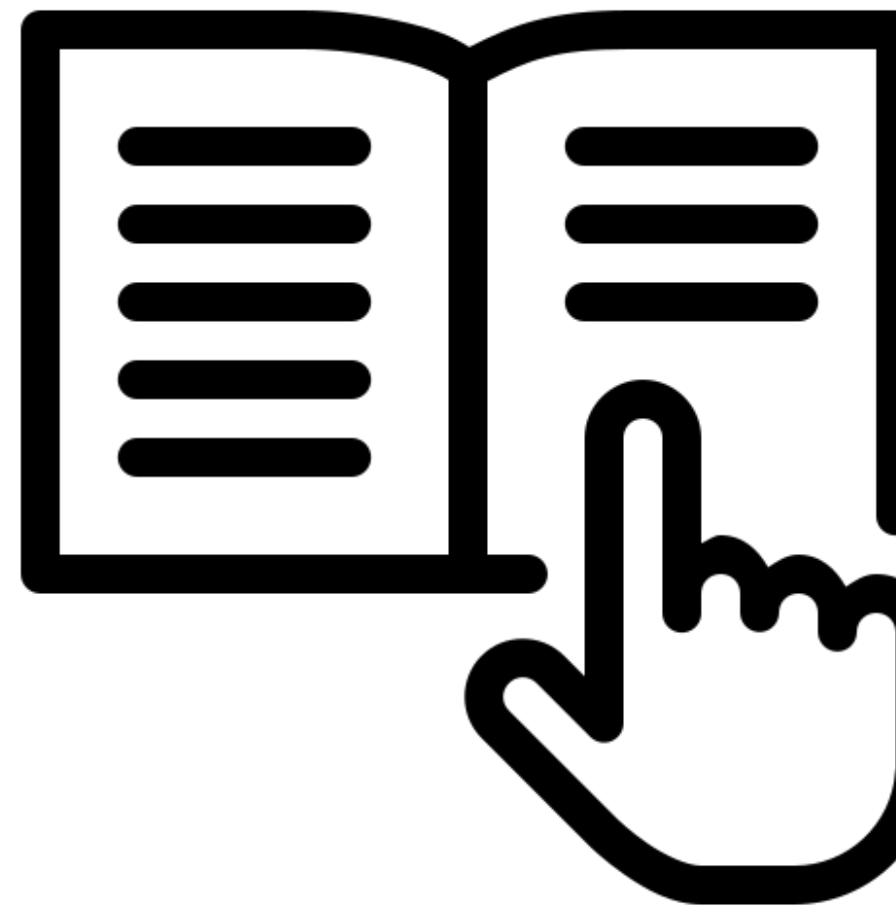


fast ingestion

Why **LSM** ?

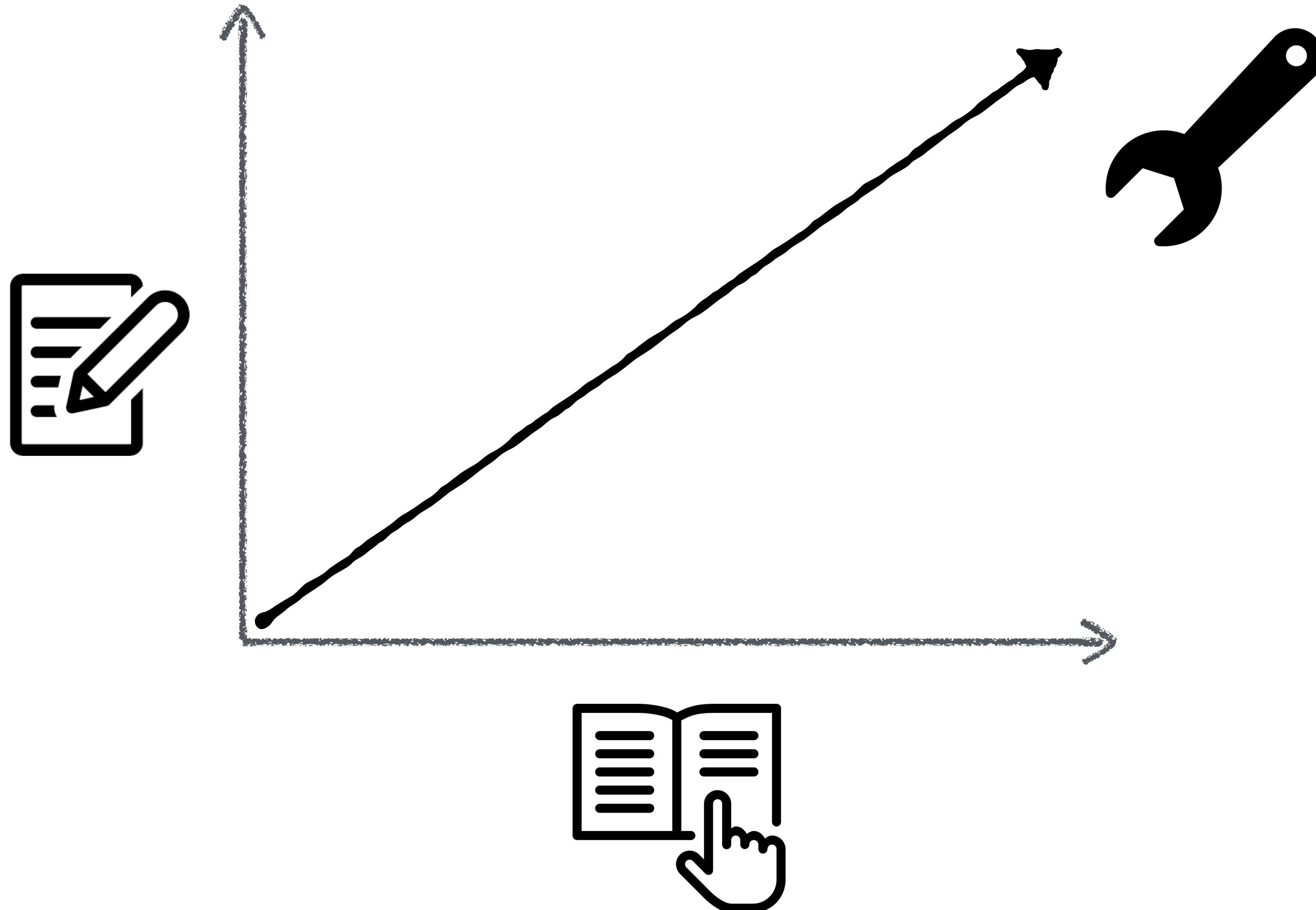


fast ingestion

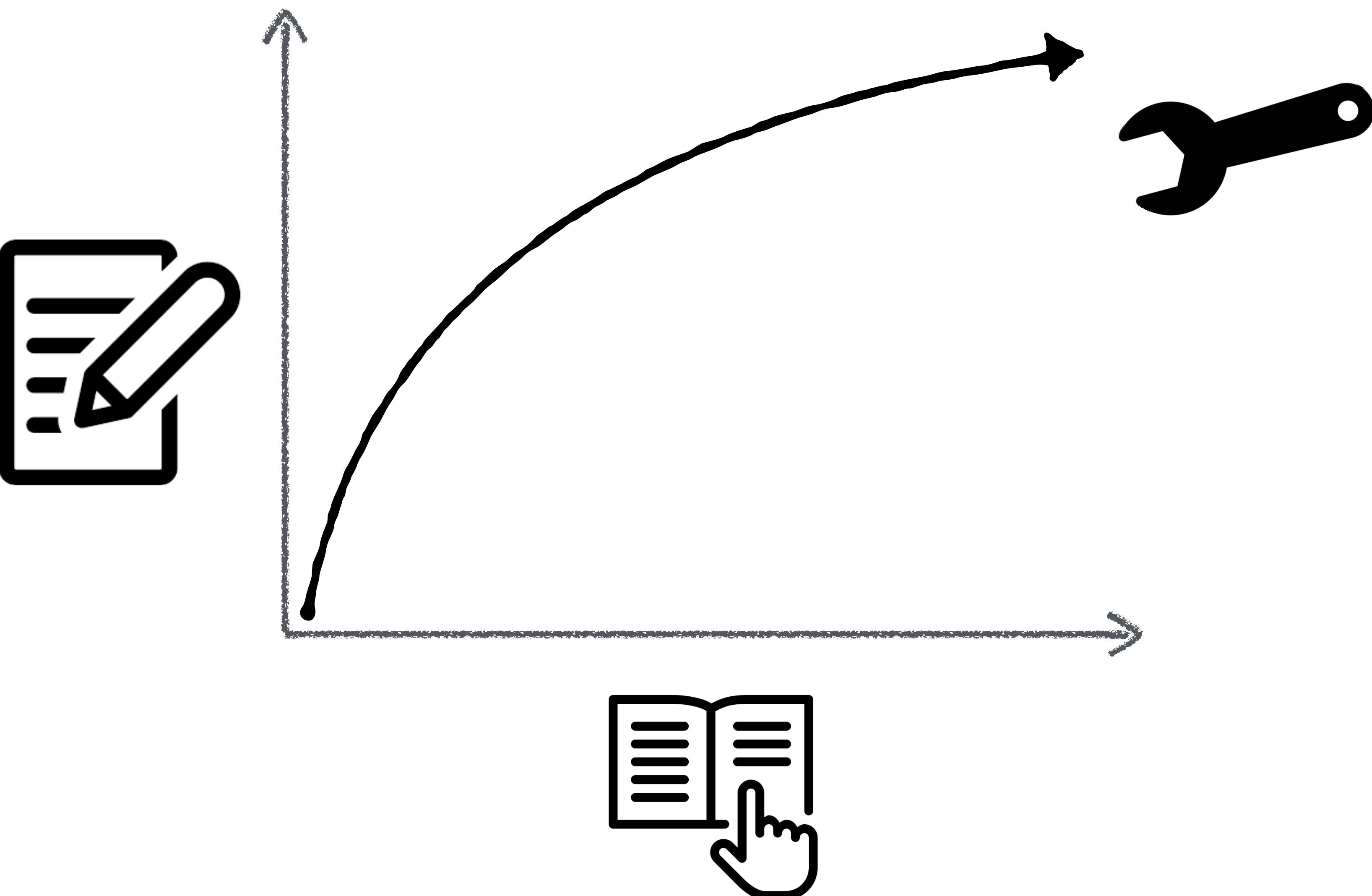


competitive reads

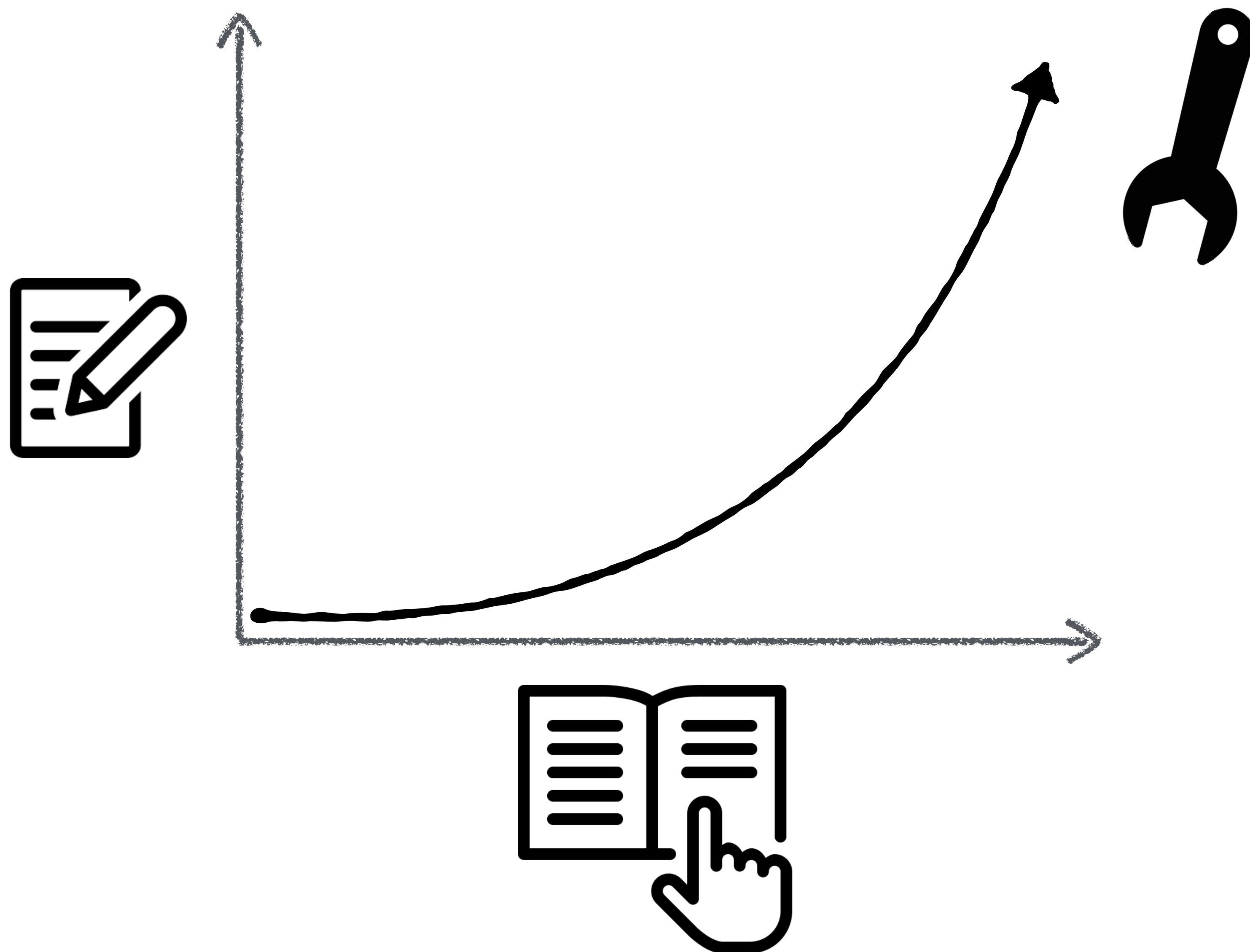
Why **LSM** ?



Why **LSM** ?



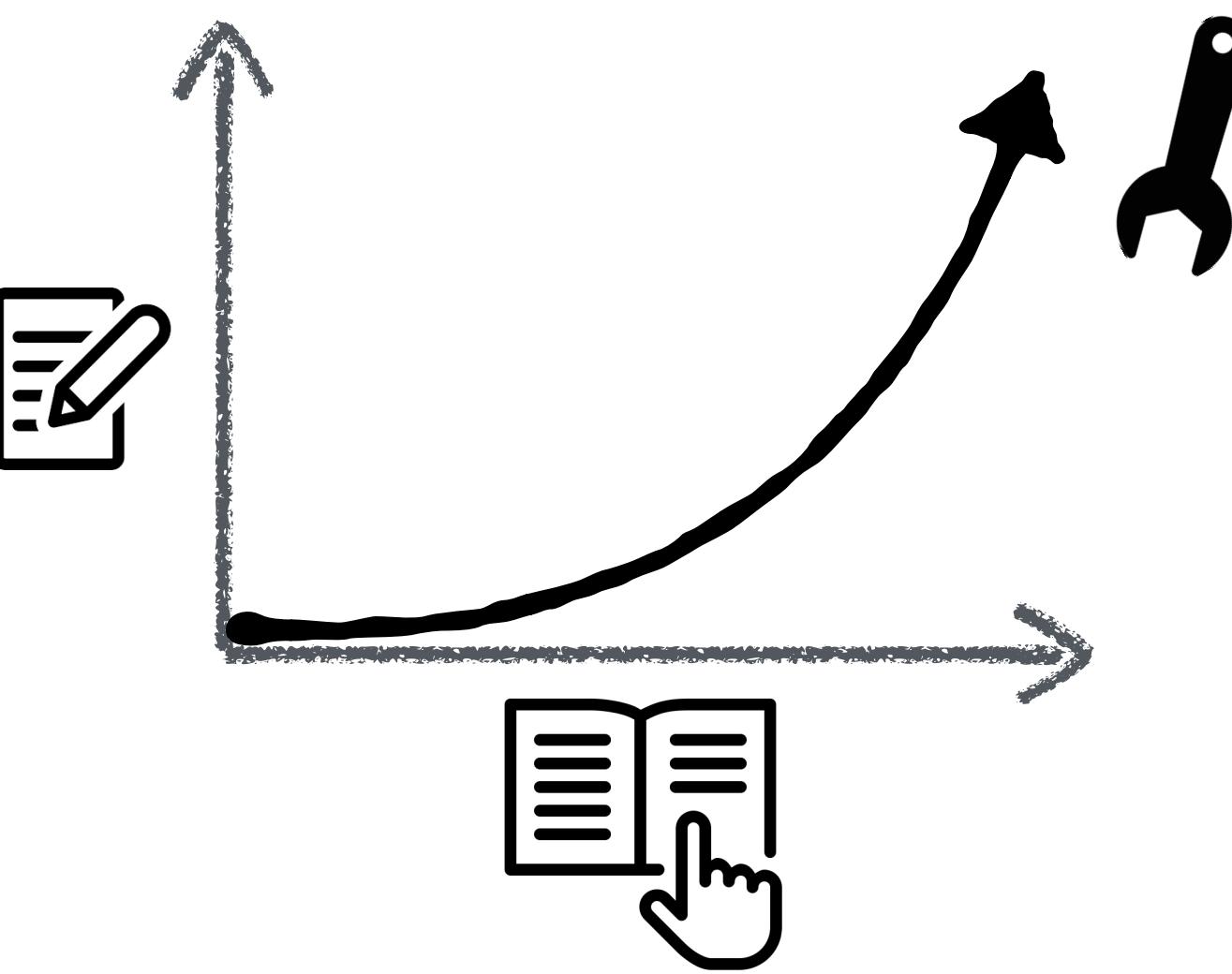
Why **LSM** ?



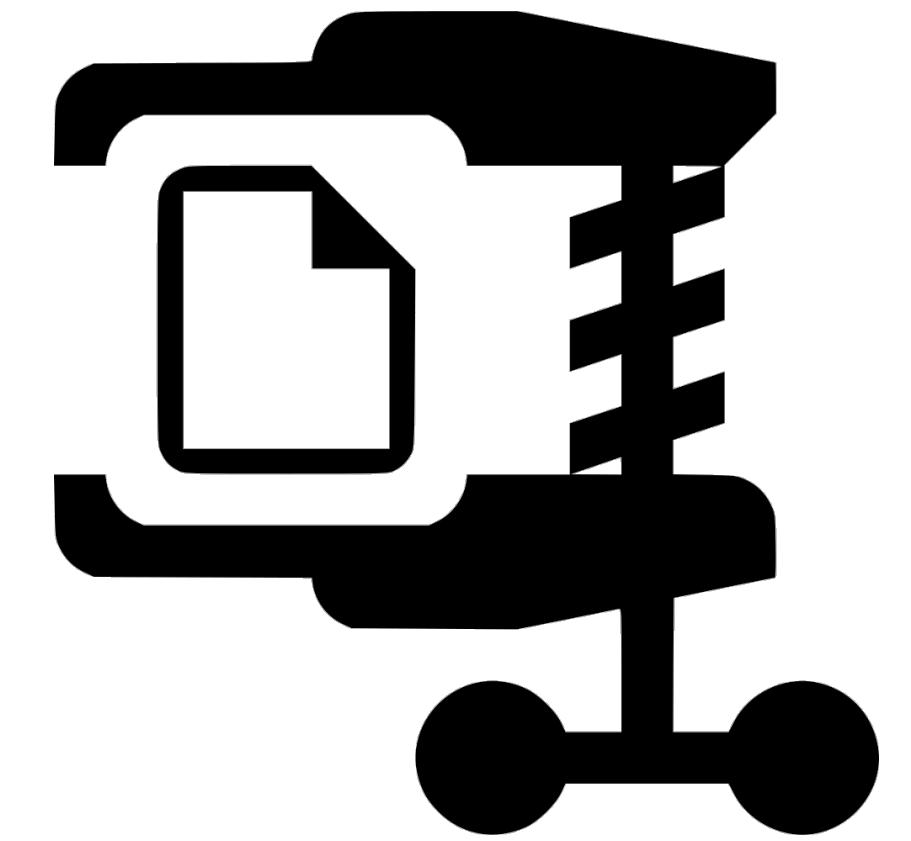
Why **LSM** ?



fast writes

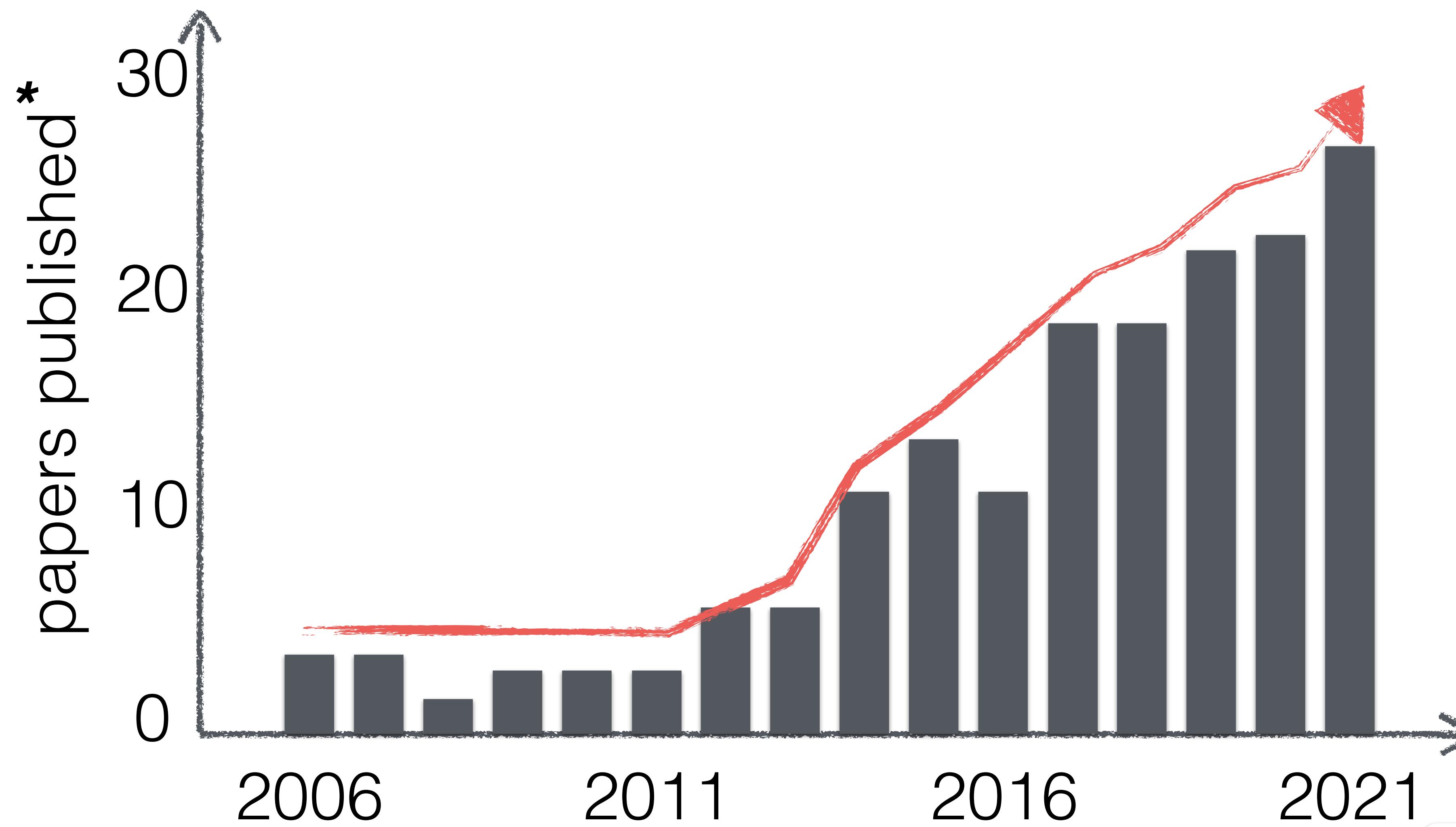


tunable read-write
performance



good space
utilization

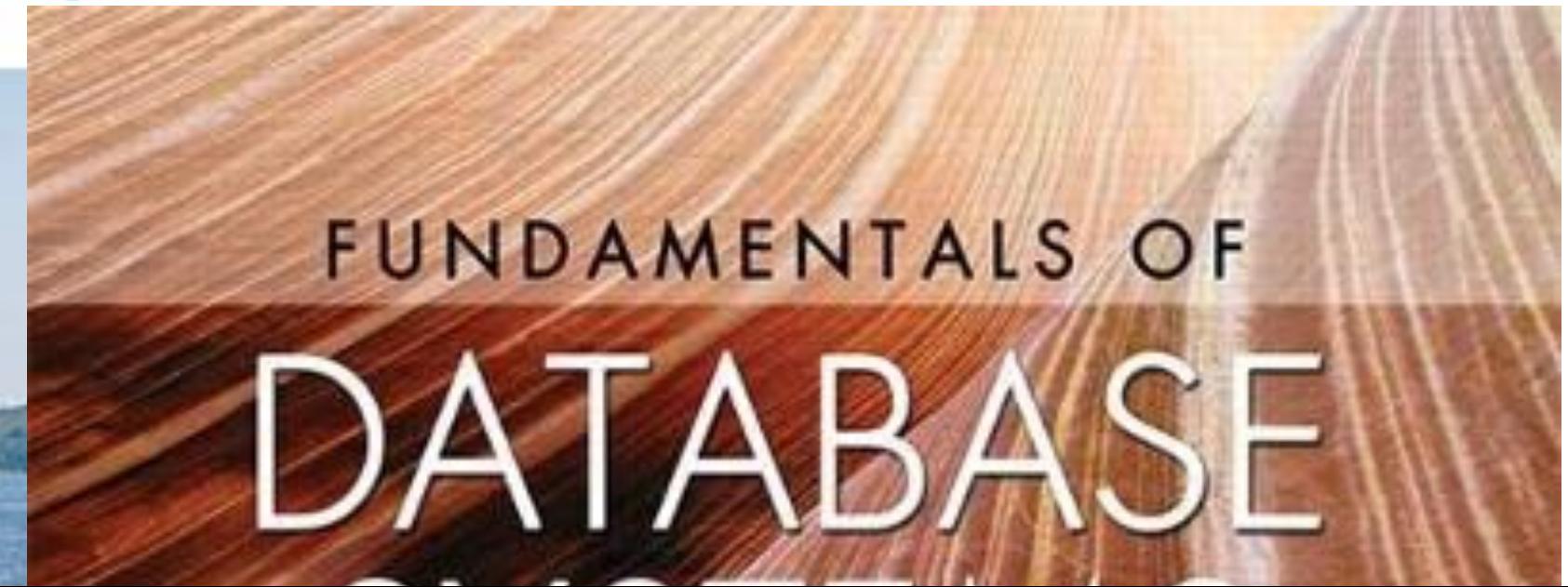
Research Trend



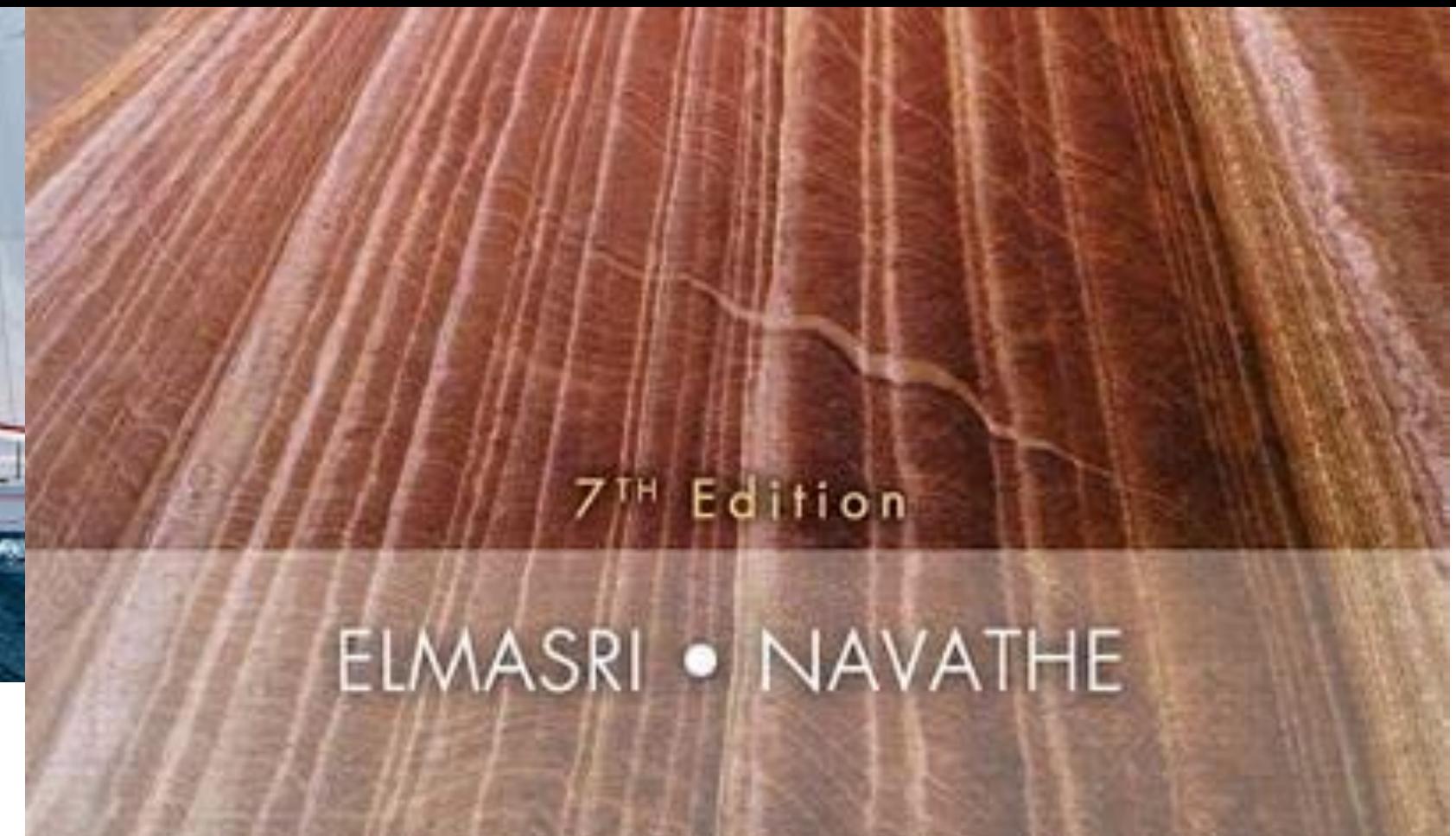
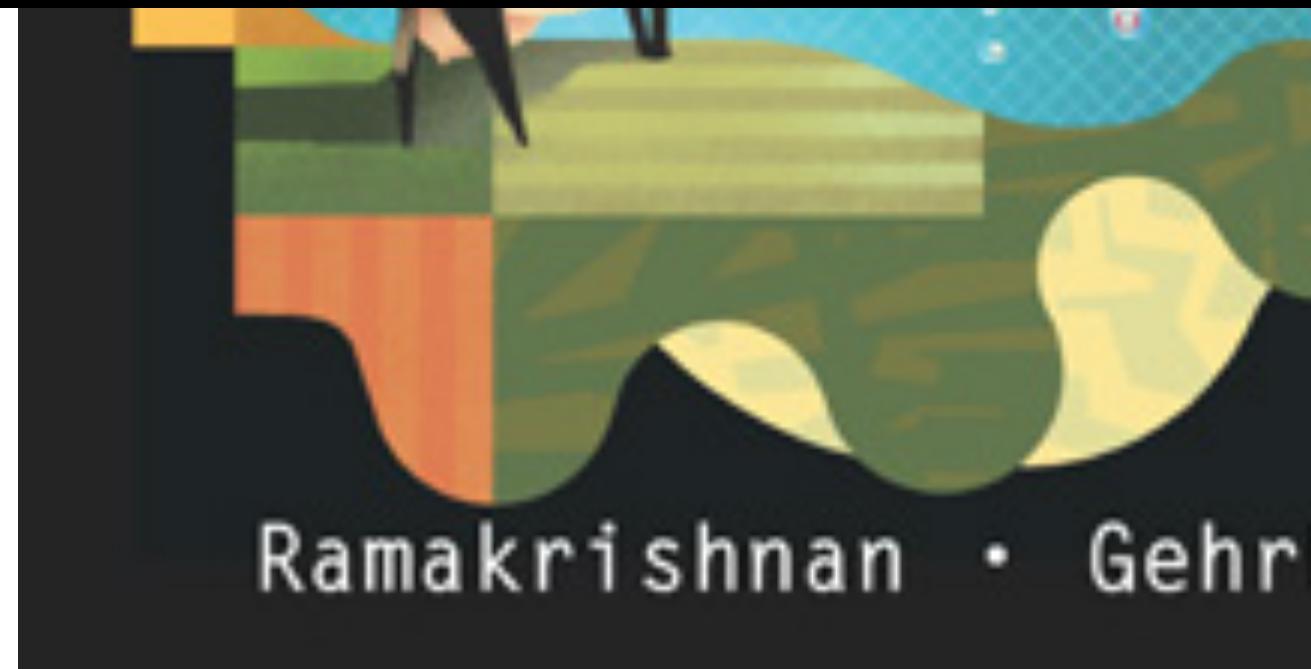
* data from DBLP



SEVENTH EDITION
Database System Concepts



No Textbook on LSMS !!



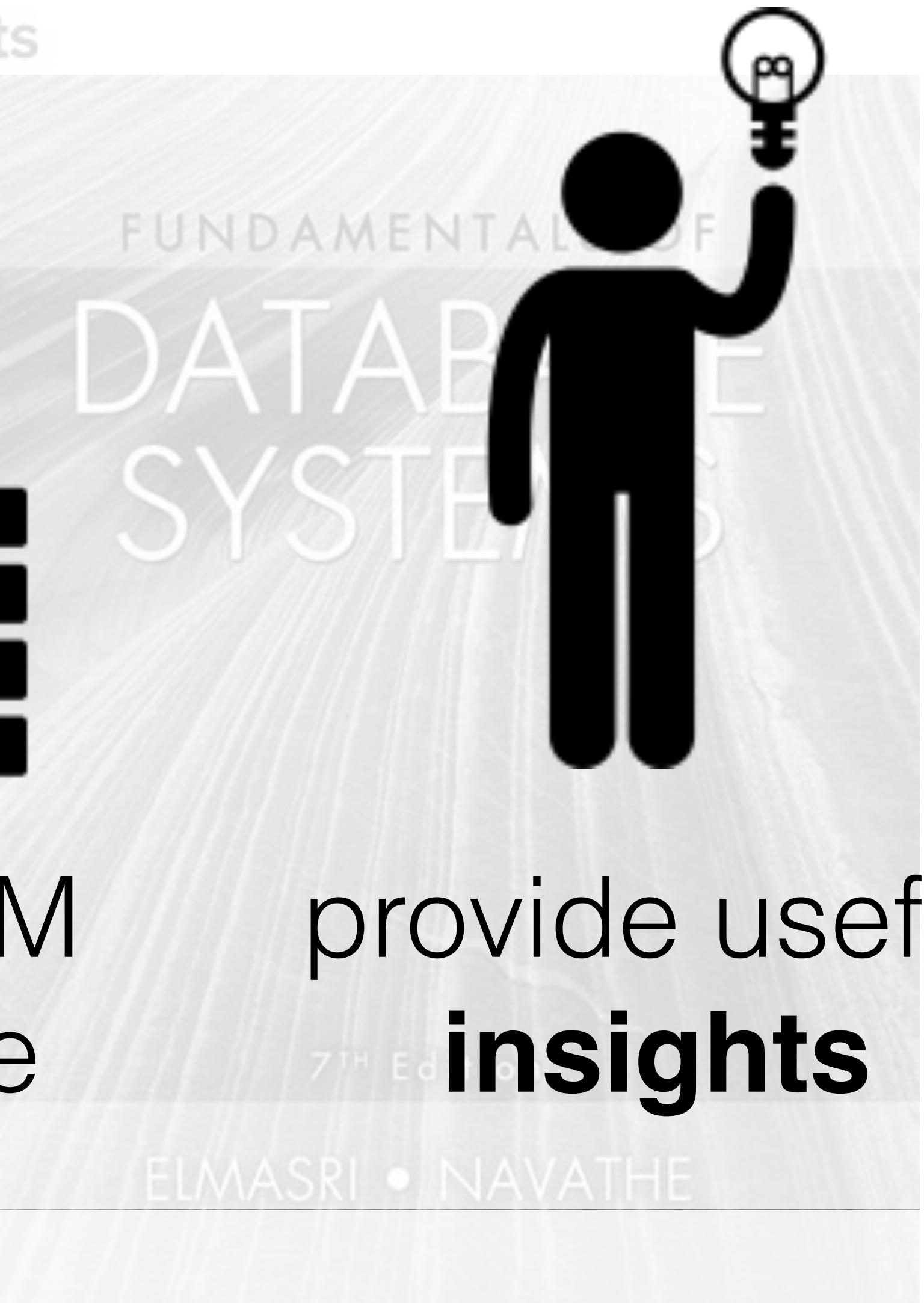
No Textbook on LMS !!



explore the
LSM paradigm



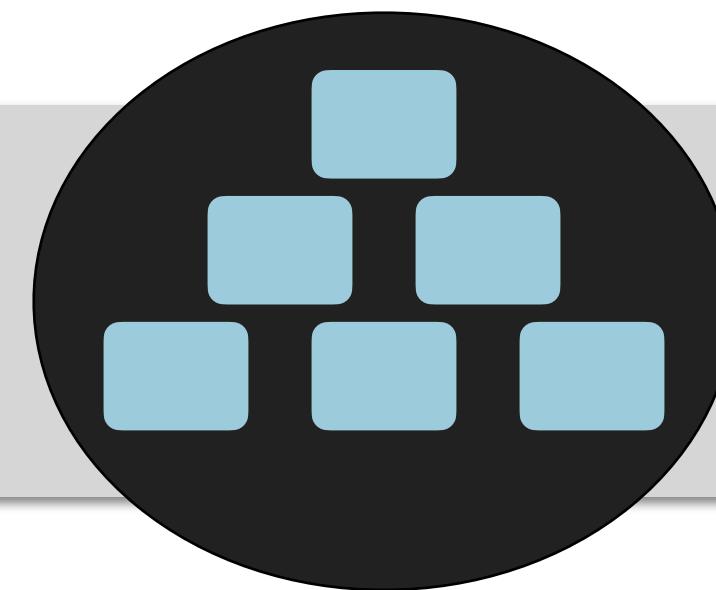
construct LSM
design space



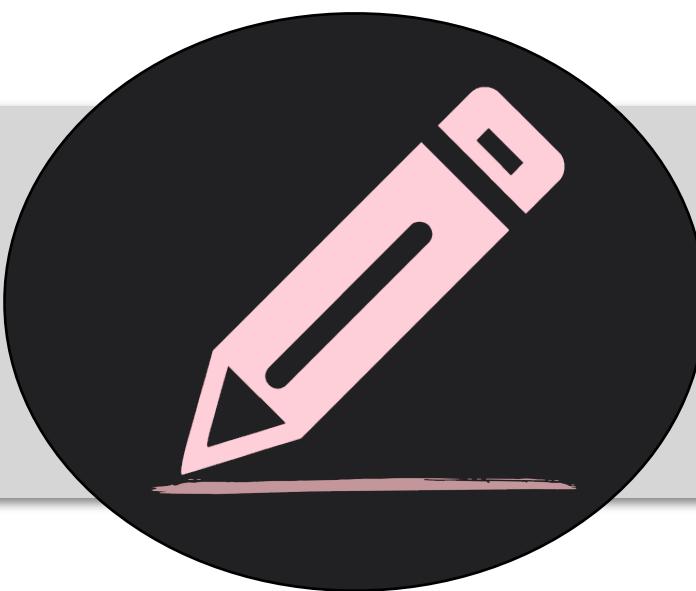
provide useful
insights

Outline

Part 1: **LSM Basics**



Part 2: **Optimizing Ingestion in LSMs**

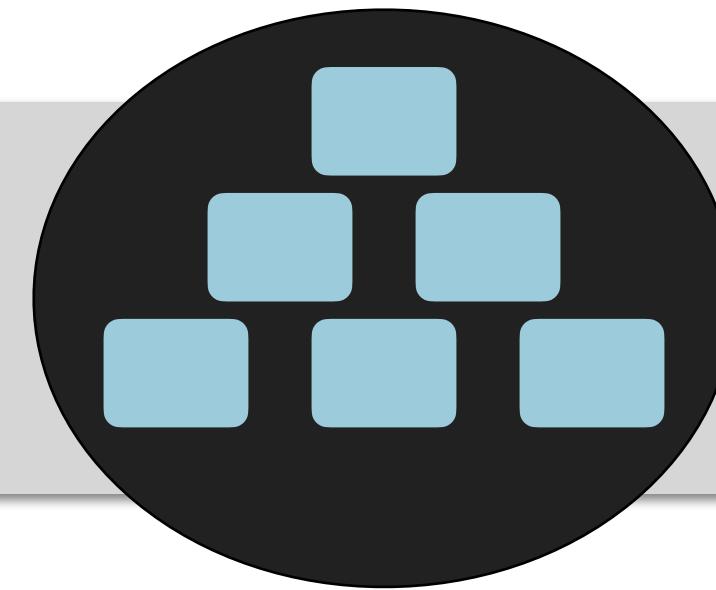


Part 3: **Navigating the LSM Design Space**

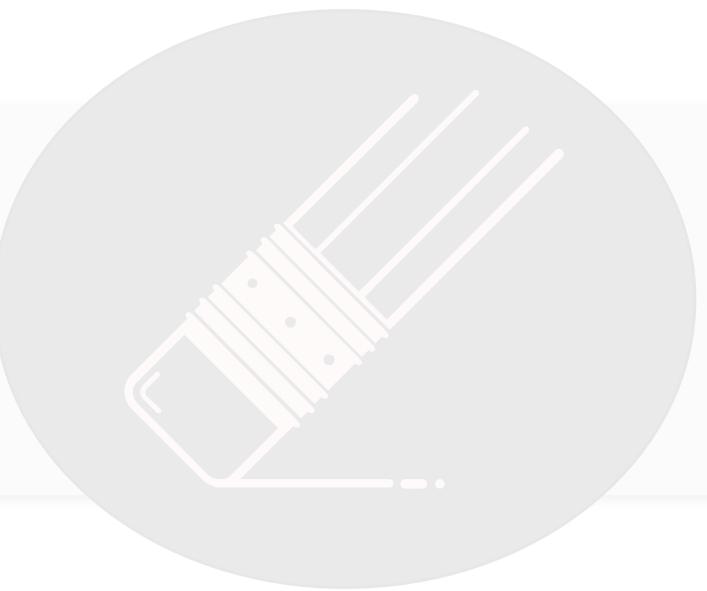


Outline

Part 1: **LSM Basics**



Part 2: Optimizing Ingestion in LSMS

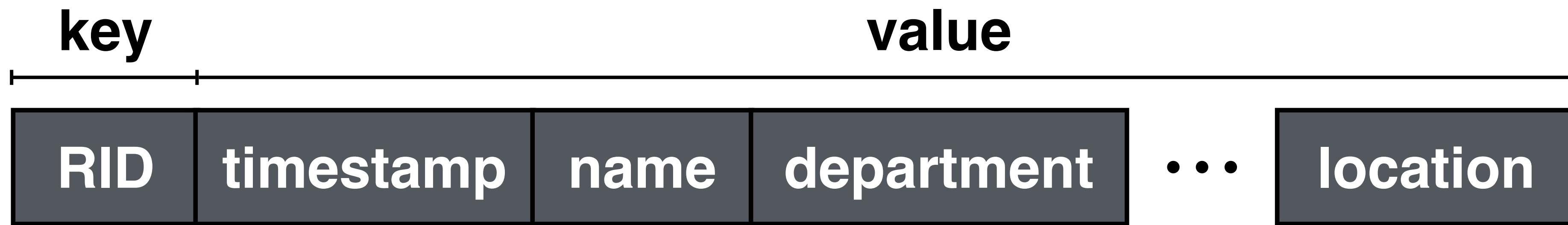


Part 3: Navigating the LSM Design Space



LSM Basics

key-value pairs



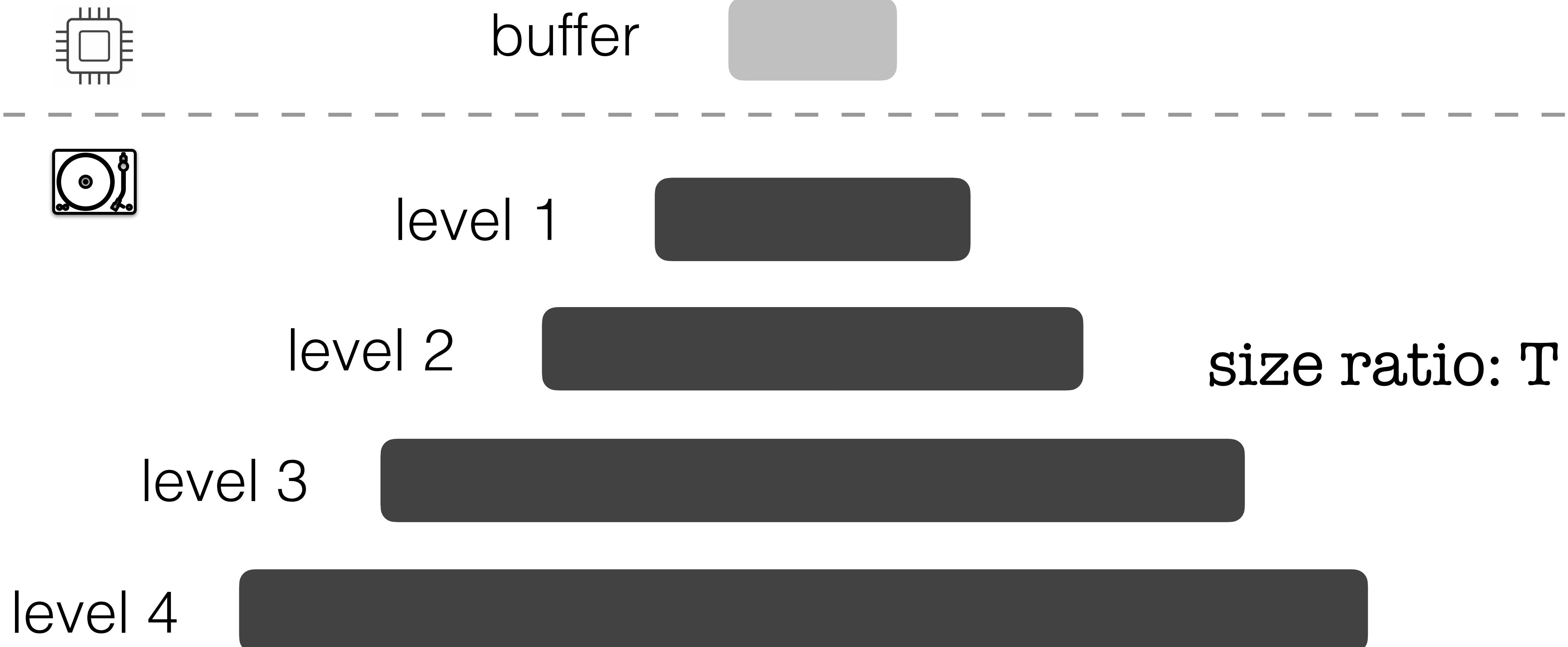
LSM Basics

key-value pairs



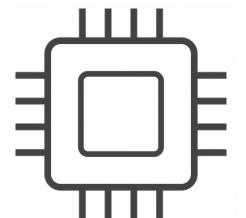
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio

LSM Basics

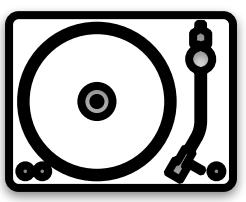
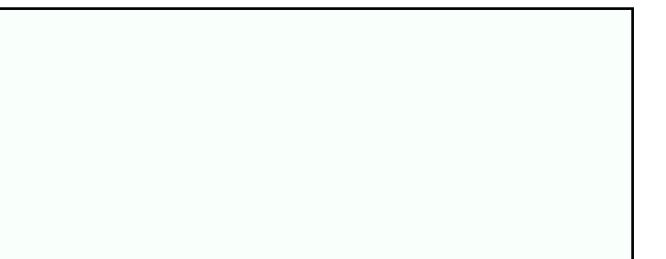


Buffering ingestion

put(6)
put(2)

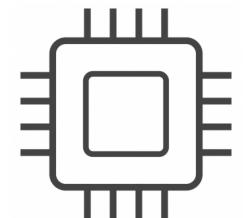


buffer

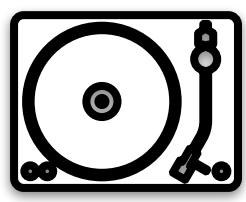
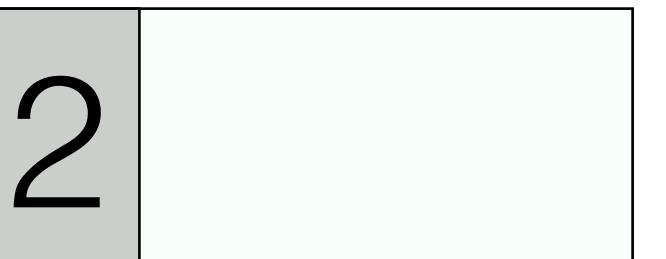


Buffering ingestion

put(1)
put(6)

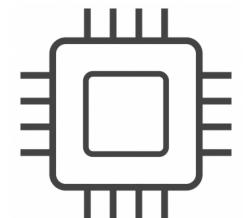


buffer

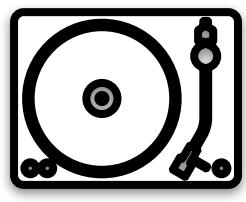
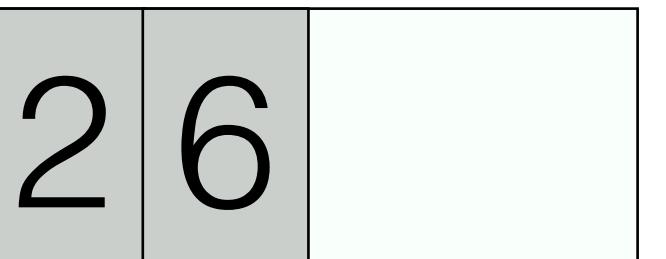


Buffering ingestion

put(4)
put(1)

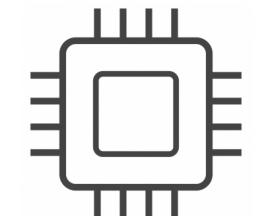


buffer

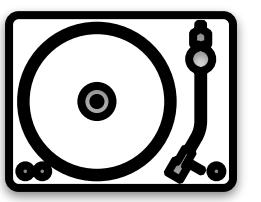
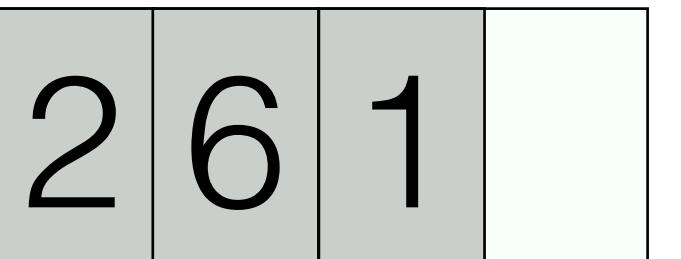


Buffering ingestion

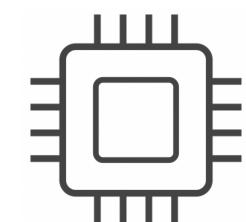
put(4)



buffer

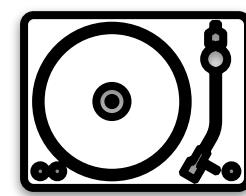


Buffering ingestion

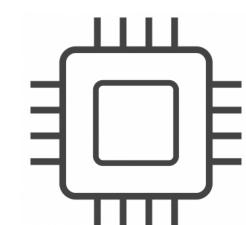


buffer

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

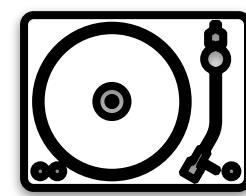


Buffering ingestion

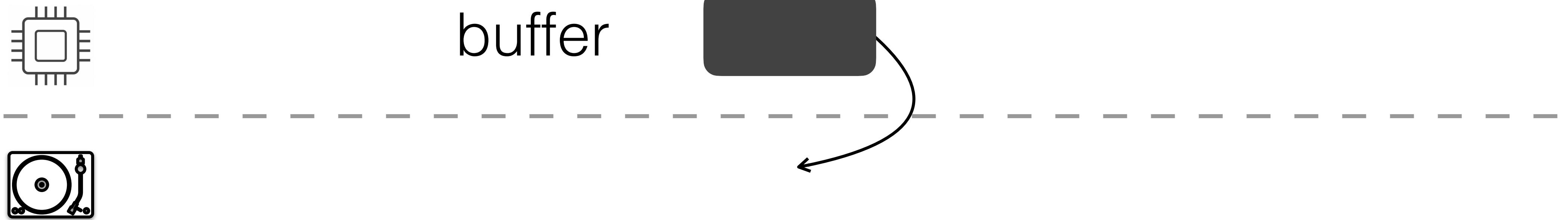


buffer

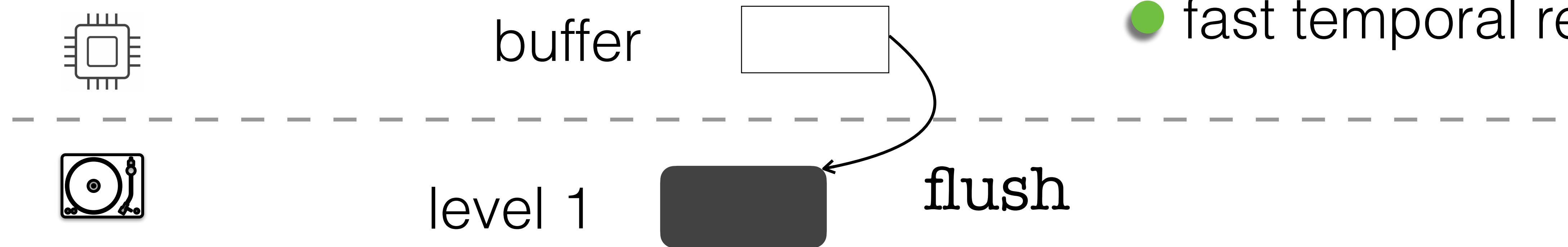
1	2	4	6
---	---	---	---



Buffering ingestion



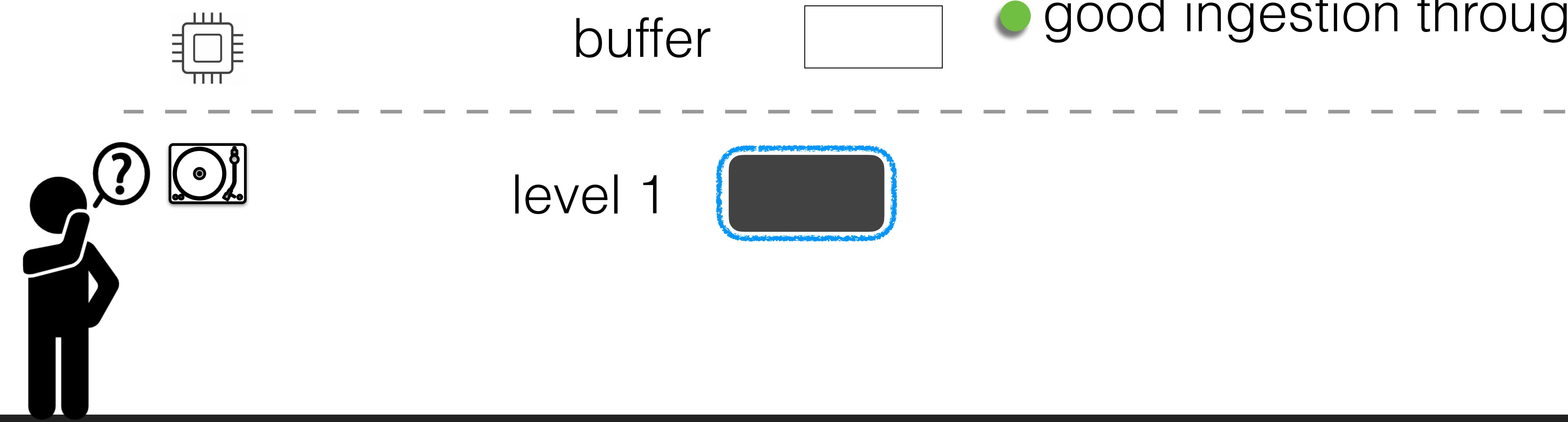
Buffering ingestion



- low ingestion cost
- fast temporal reads

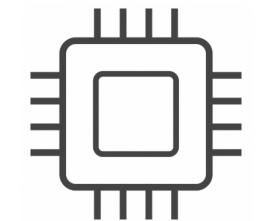
Immutable files on storage

- compact storage
- good ingestion throughput

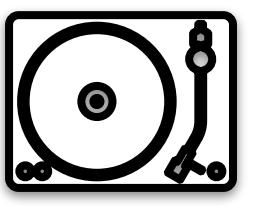
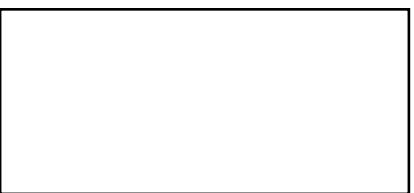


How do we update data?

put(6)

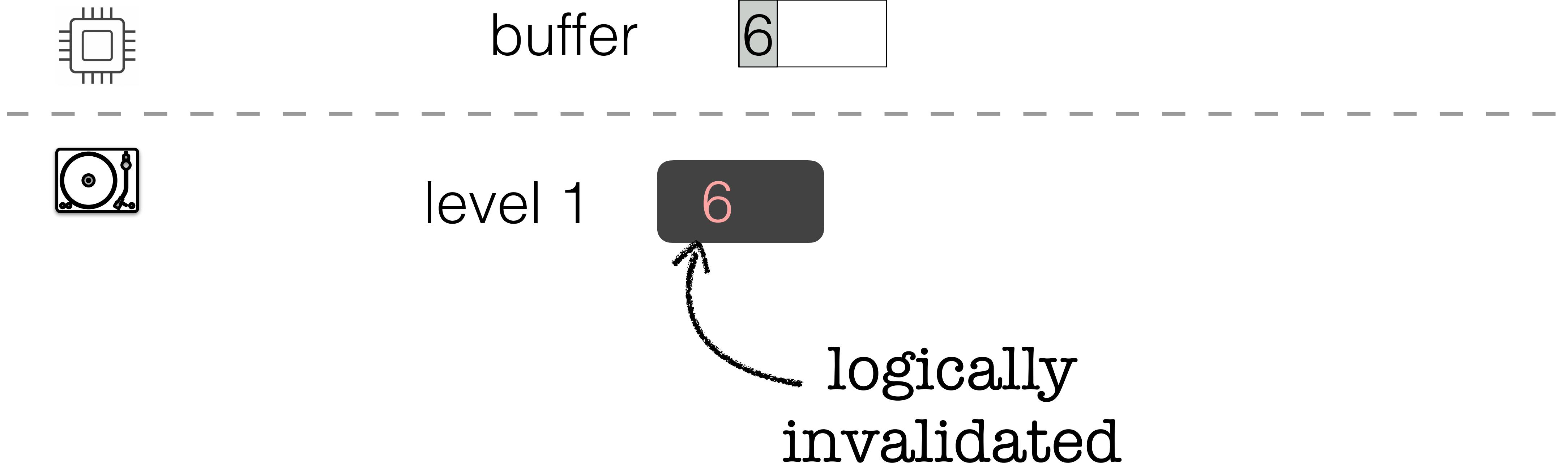


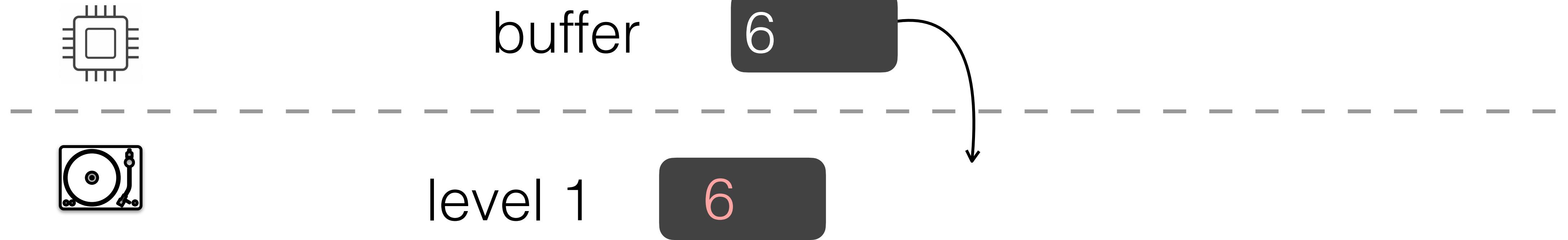
buffer

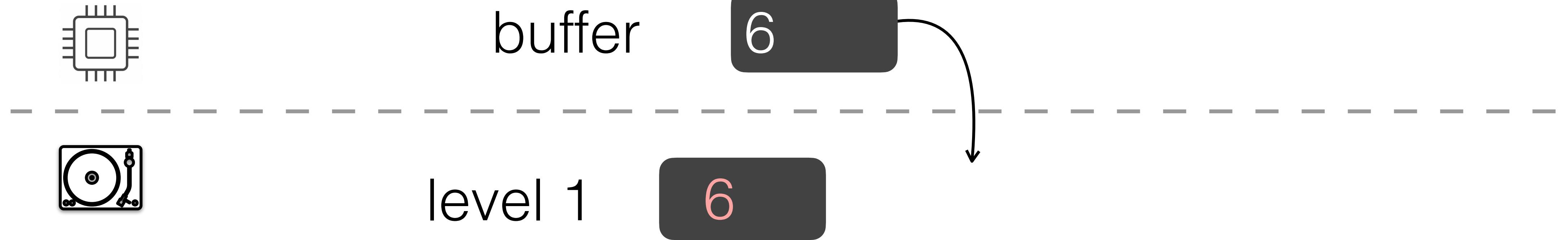


level 1

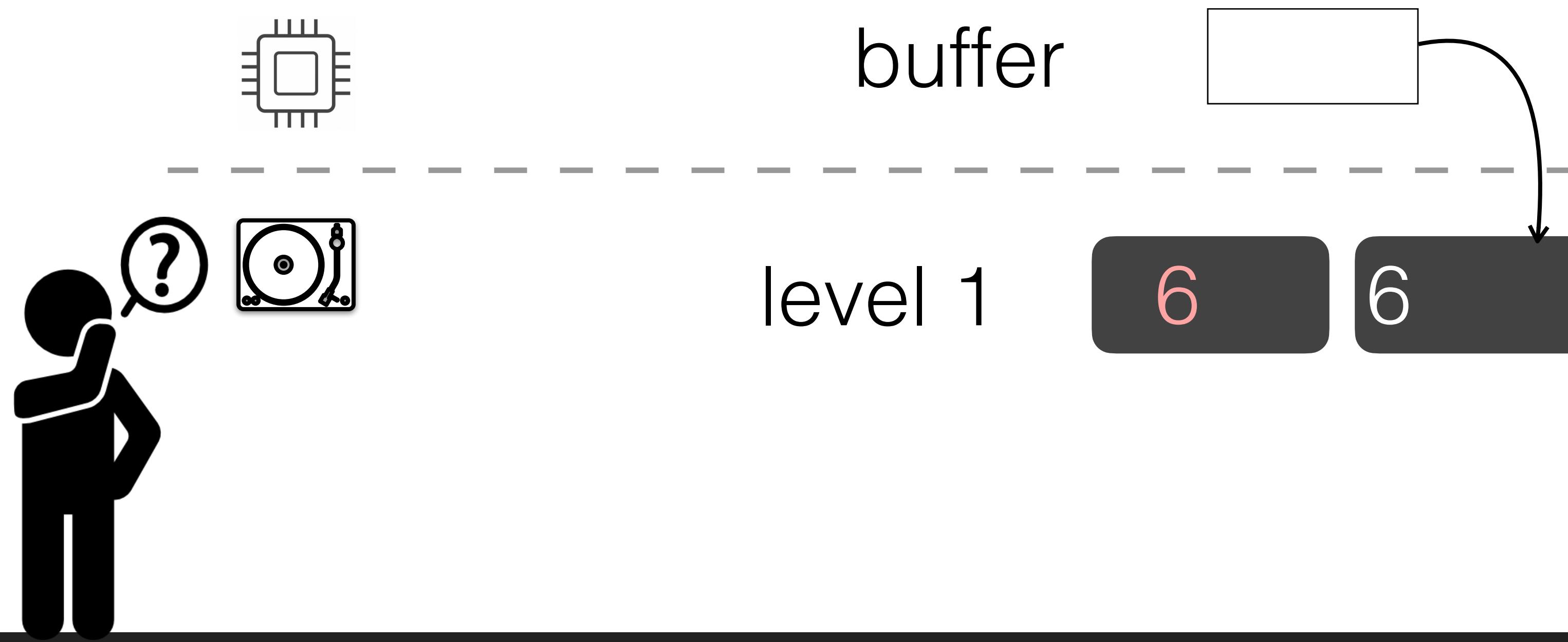
6





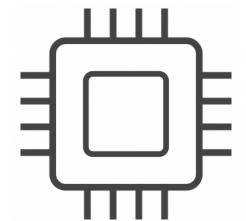


Out-of-place updates

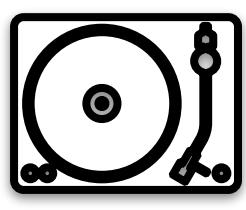


- fast ingestion
- space amplification
- slow reads

How do we reduce this space amplification?

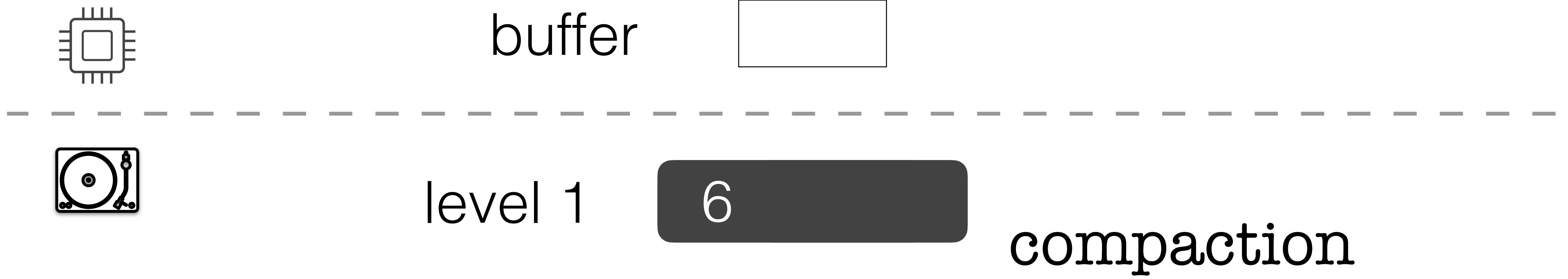


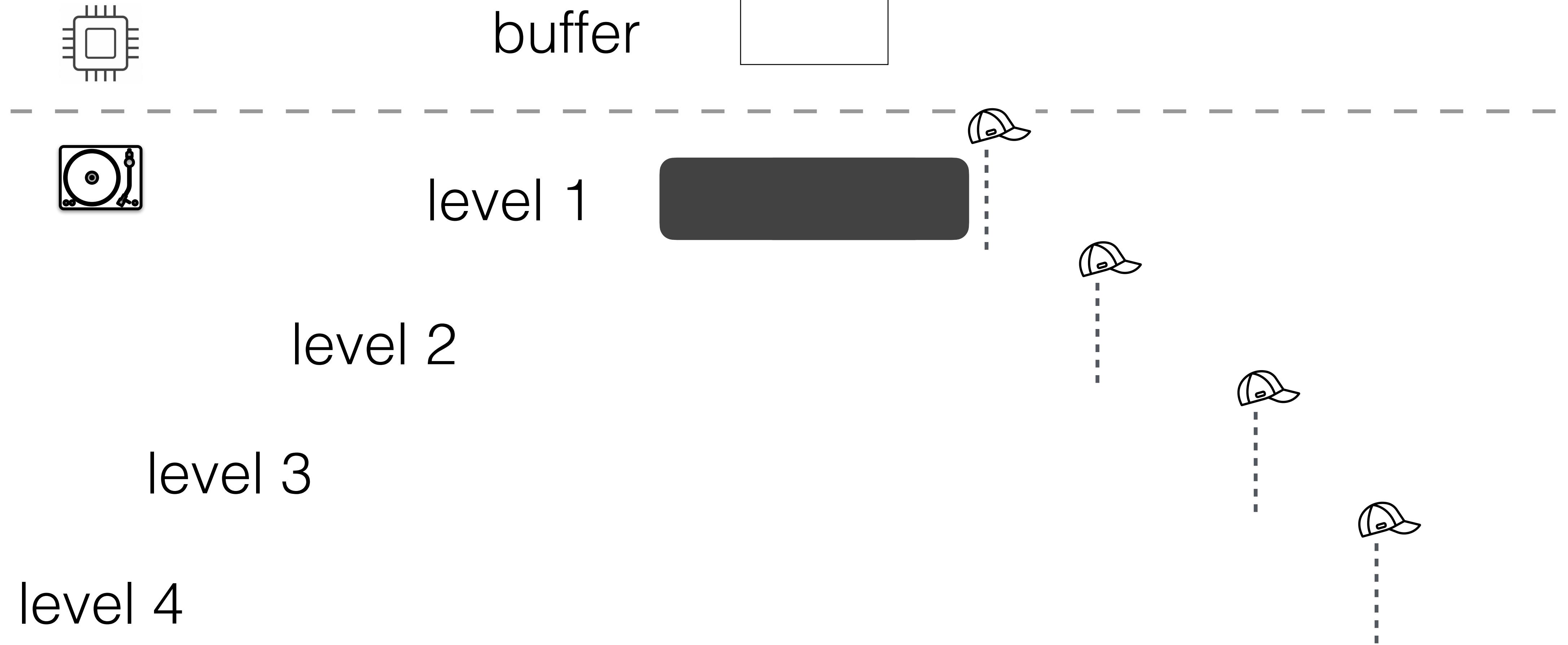
buffer

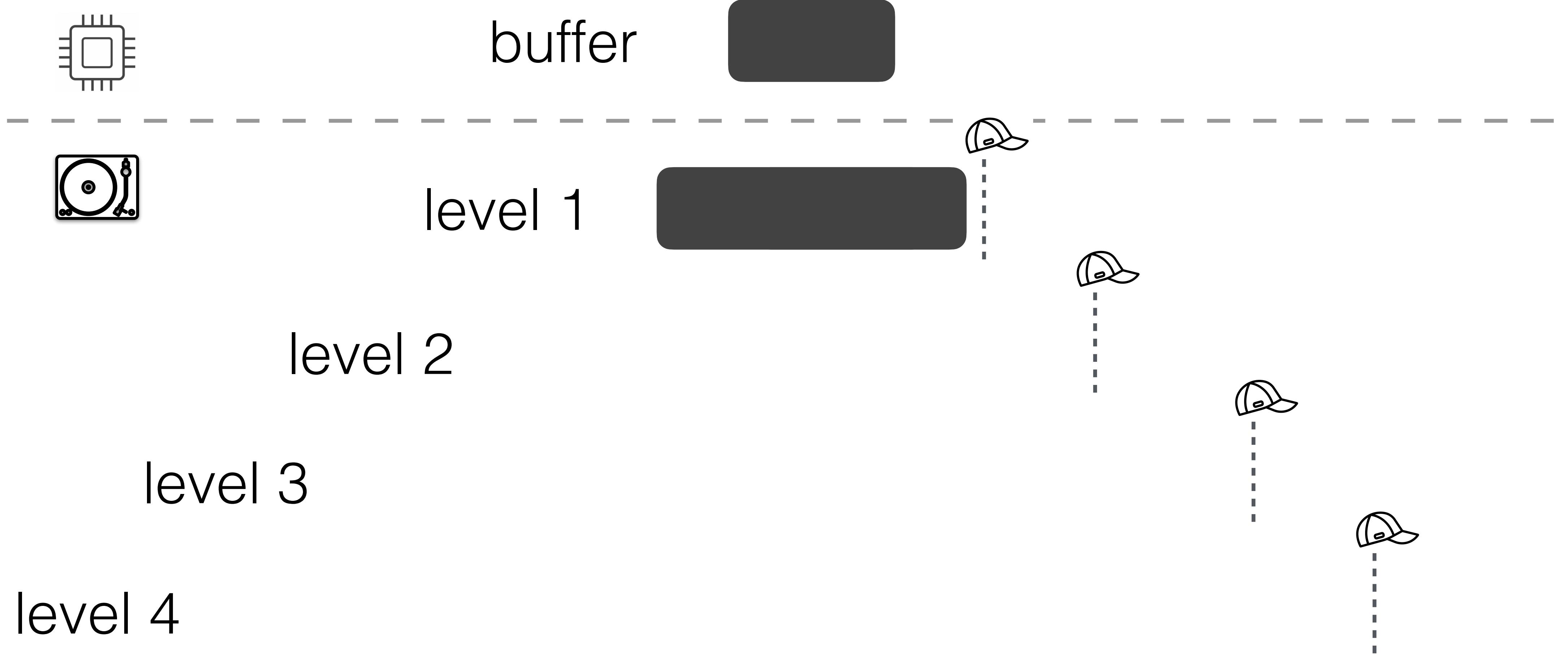


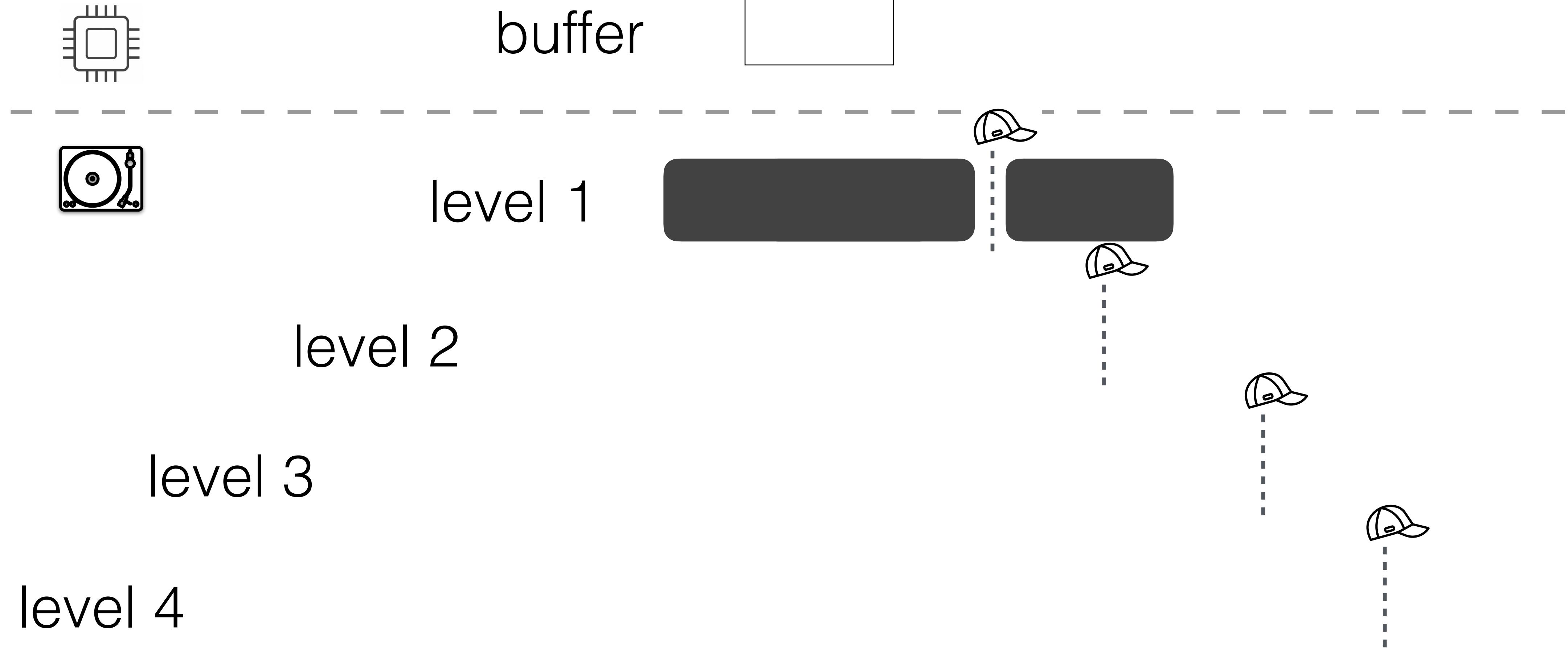
level 1

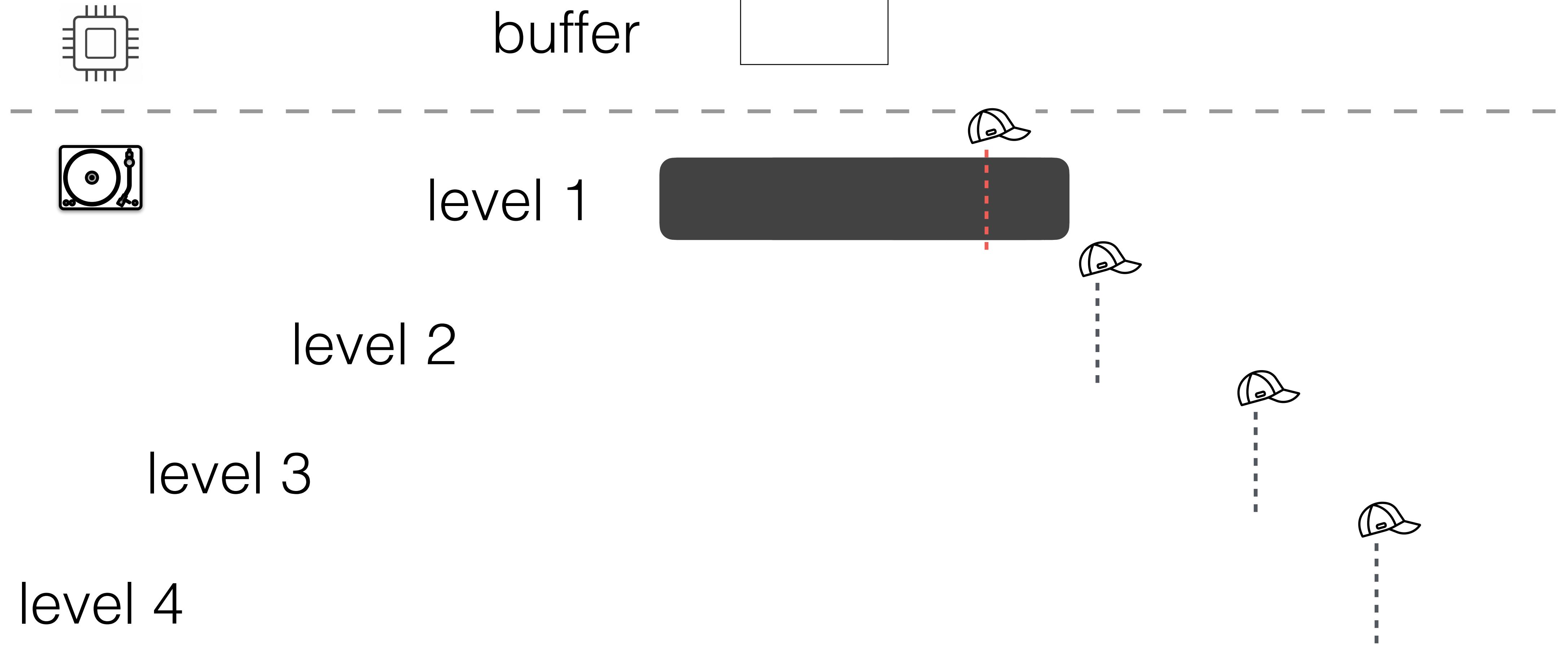


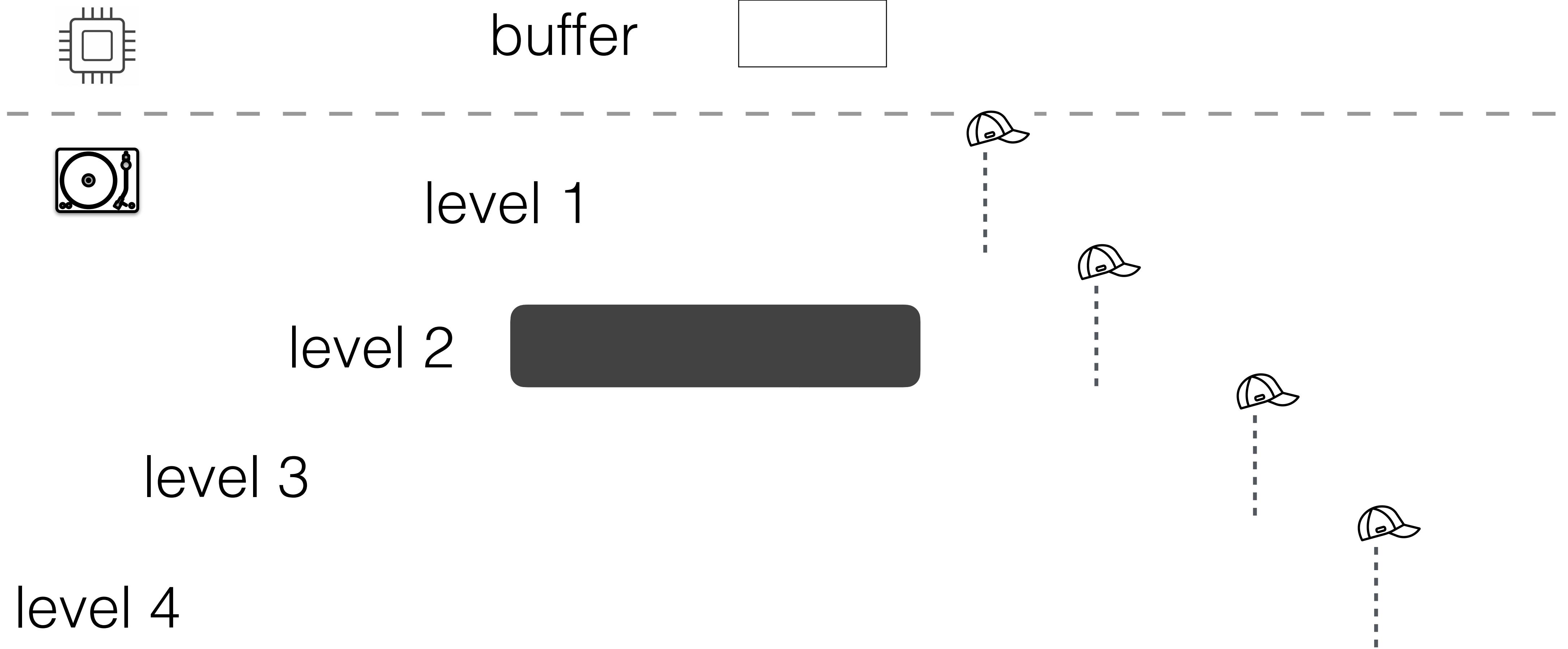












P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio

Periodic compactions

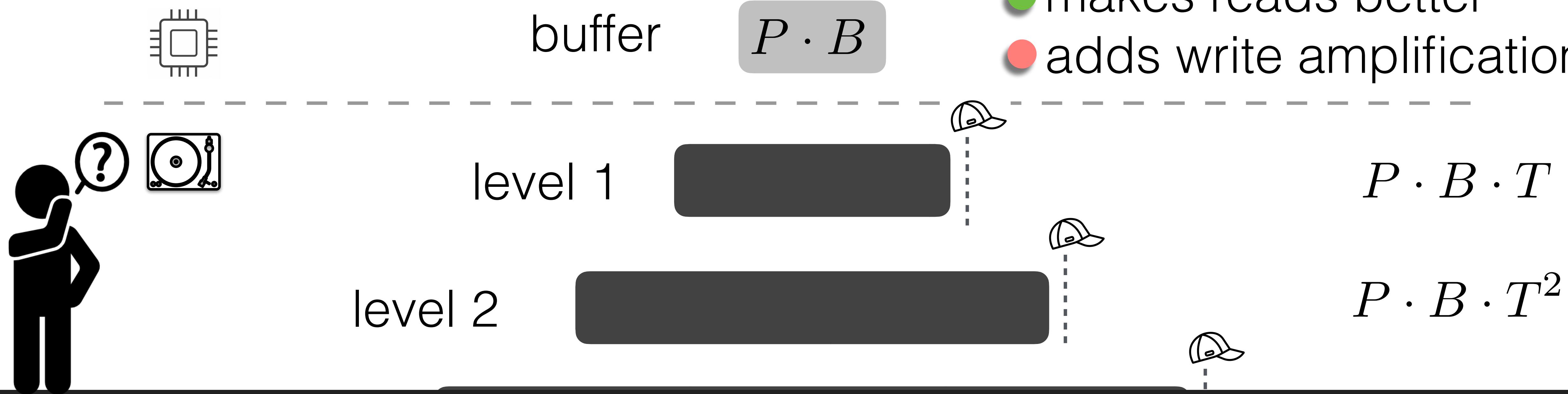
- low space amplification
- makes reads better
- adds write amplification



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

Periodic compactions

- low space amplification
- makes reads better
- adds write amplification



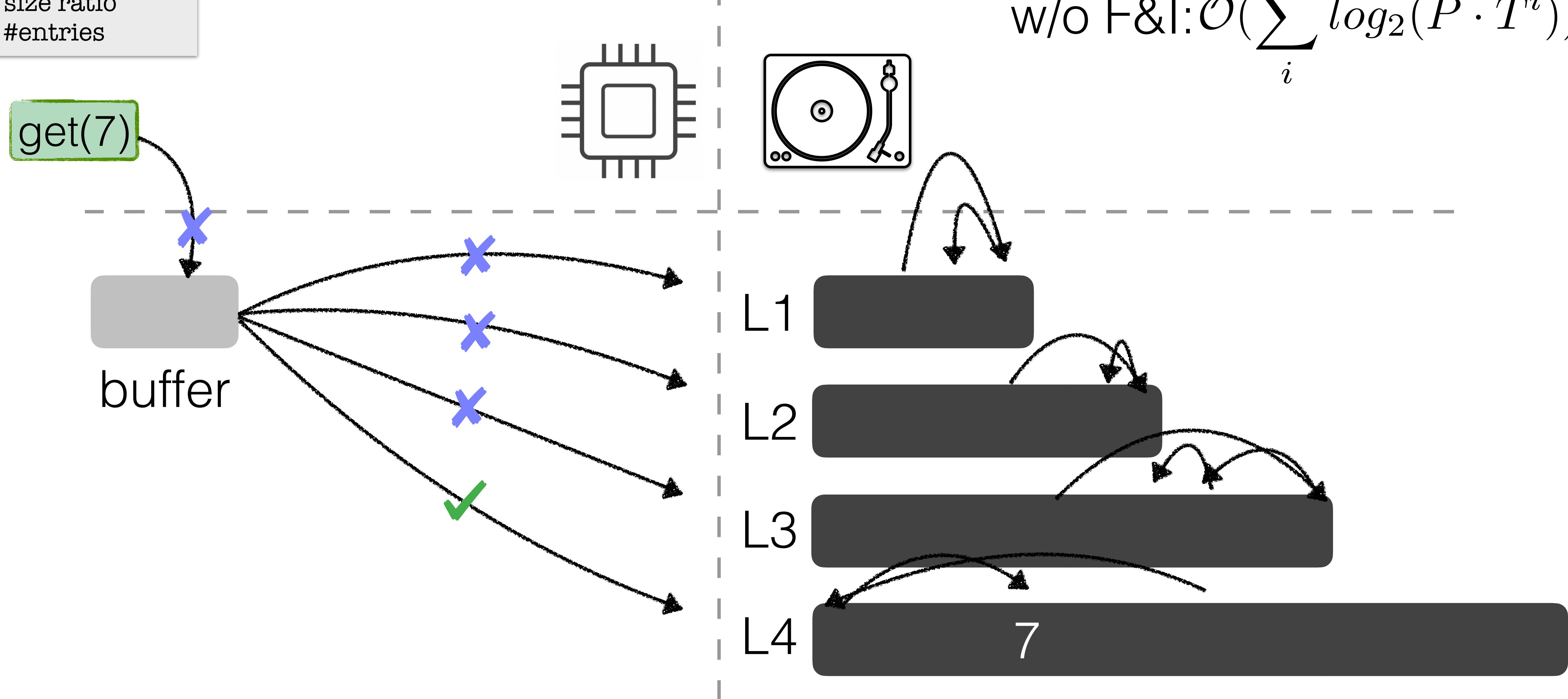
How about queries?

capacity
(entries)

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

Cost analysis

w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

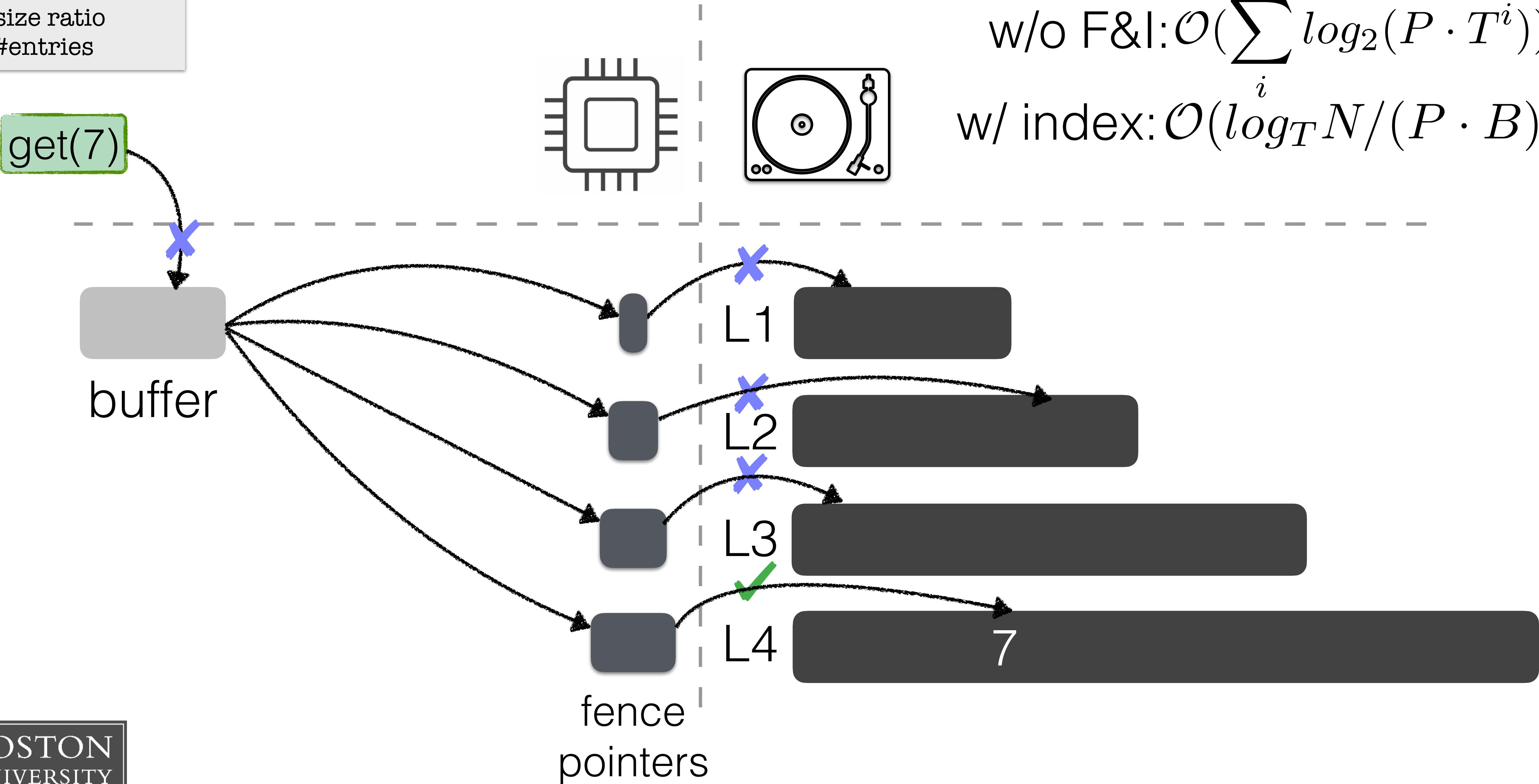


P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

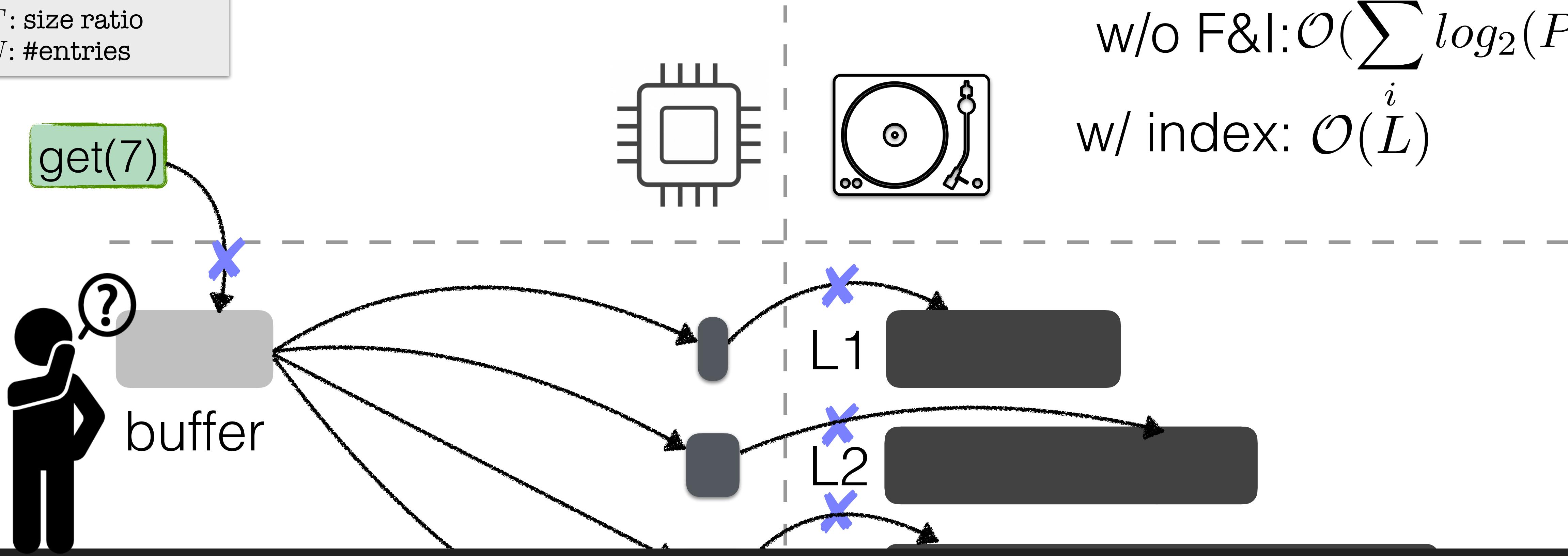
Cost analysis

w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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



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



Cost analysis

w/o F&I: $\mathcal{O}(\sum \log_2(P \cdot T^i))$

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

How to avoid unnecessary I/Os?

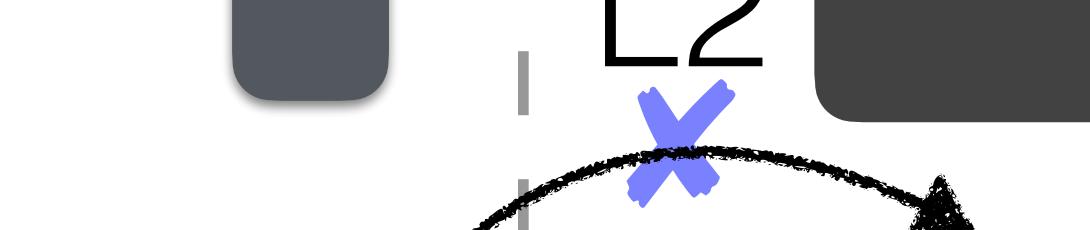
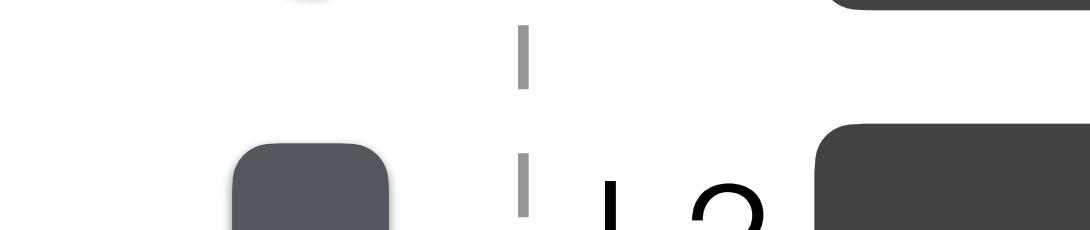
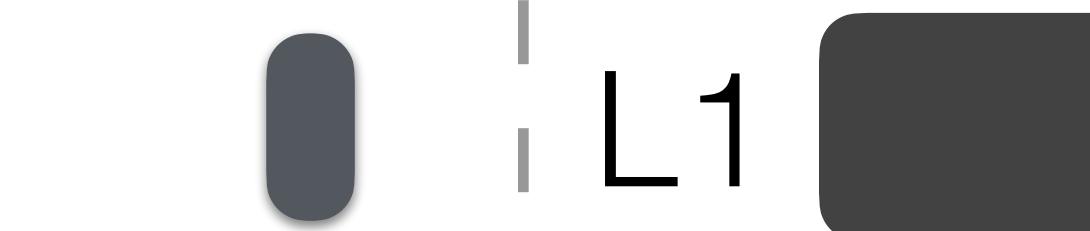
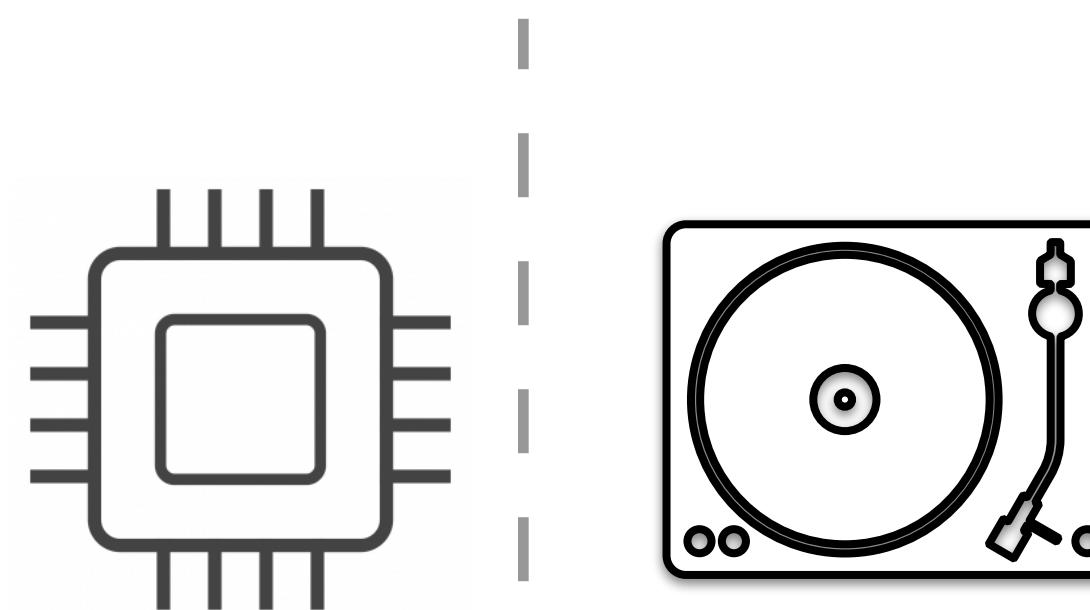
10.100
pointers

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

get(7)



buffer



Bloom
filters

fence
pointers

Cost analysis

w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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

w F&I: $\mathcal{O}(\phi \cdot L)$

$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2}$$

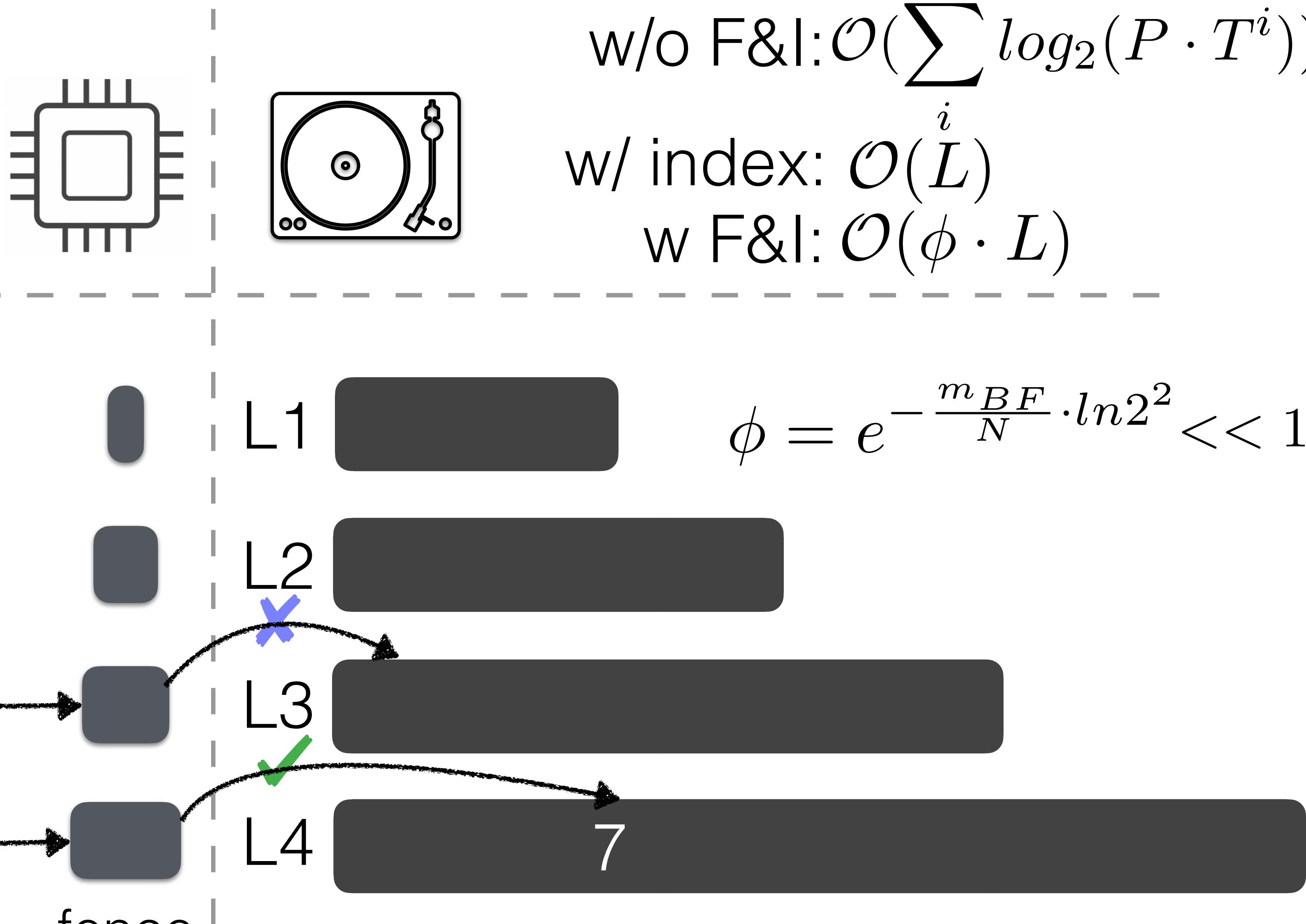
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

get(7)

buffer

Bloom filters

fence pointers



Cost analysis

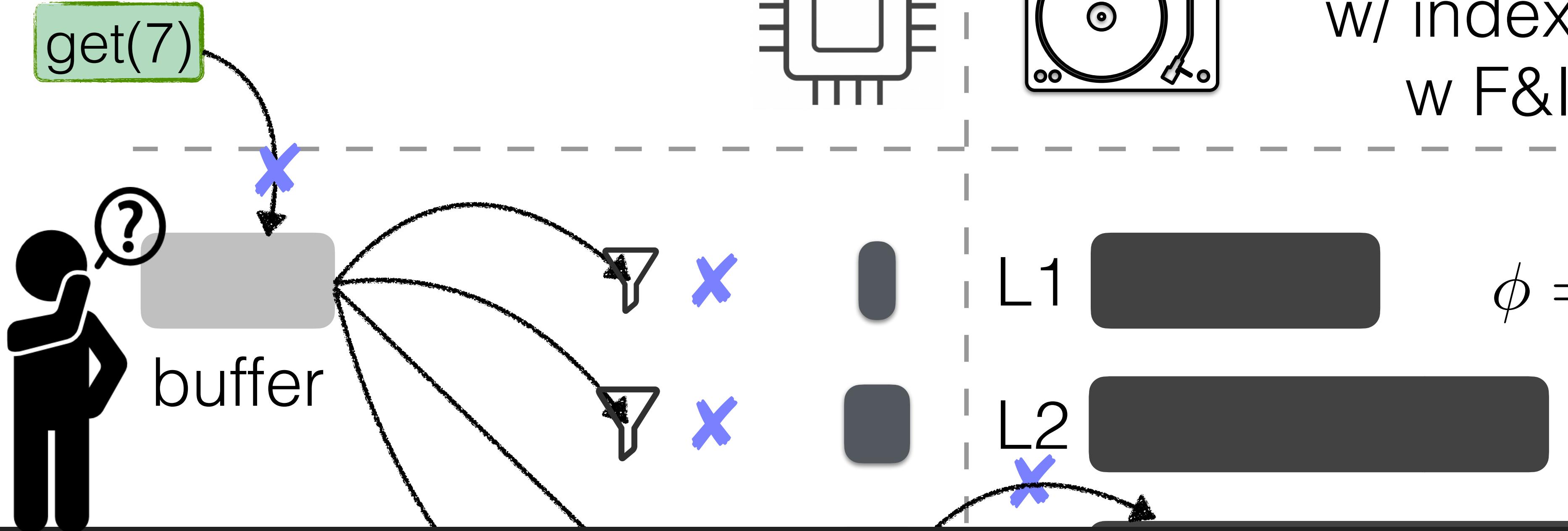
w/o F&I: $\mathcal{O}(\sum_i^i \log_2(P \cdot T^i))$

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

w F&I: $\mathcal{O}(\phi \cdot L)$

$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2} \ll 1$$

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



Cost analysis

w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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

w F&I: $\mathcal{O}(\phi \cdot L)$

$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2} \ll 1$$

How to manage memory?

Bloom
filters
pointers

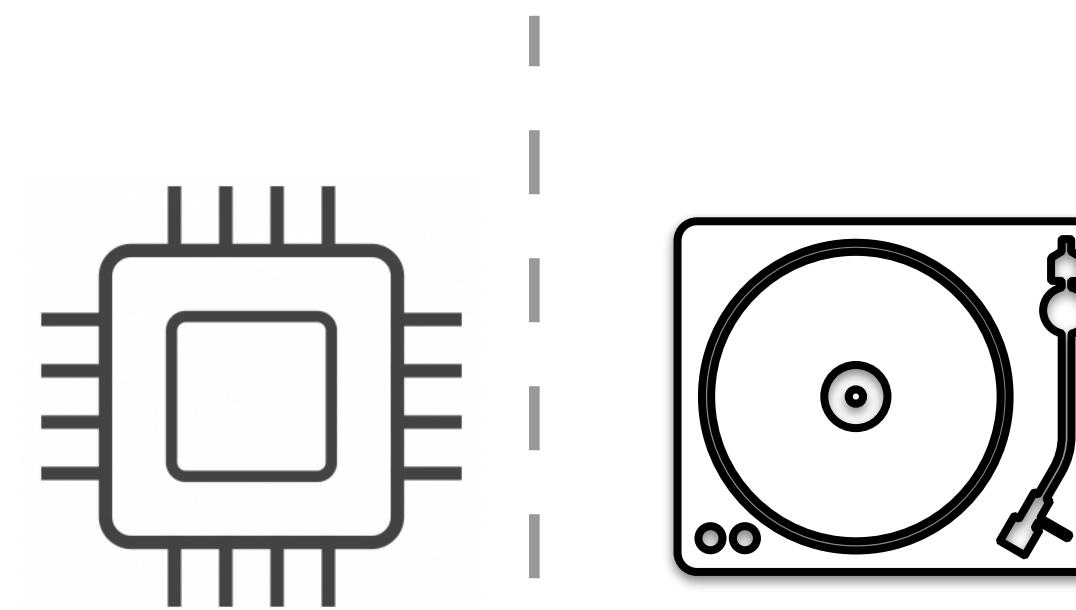
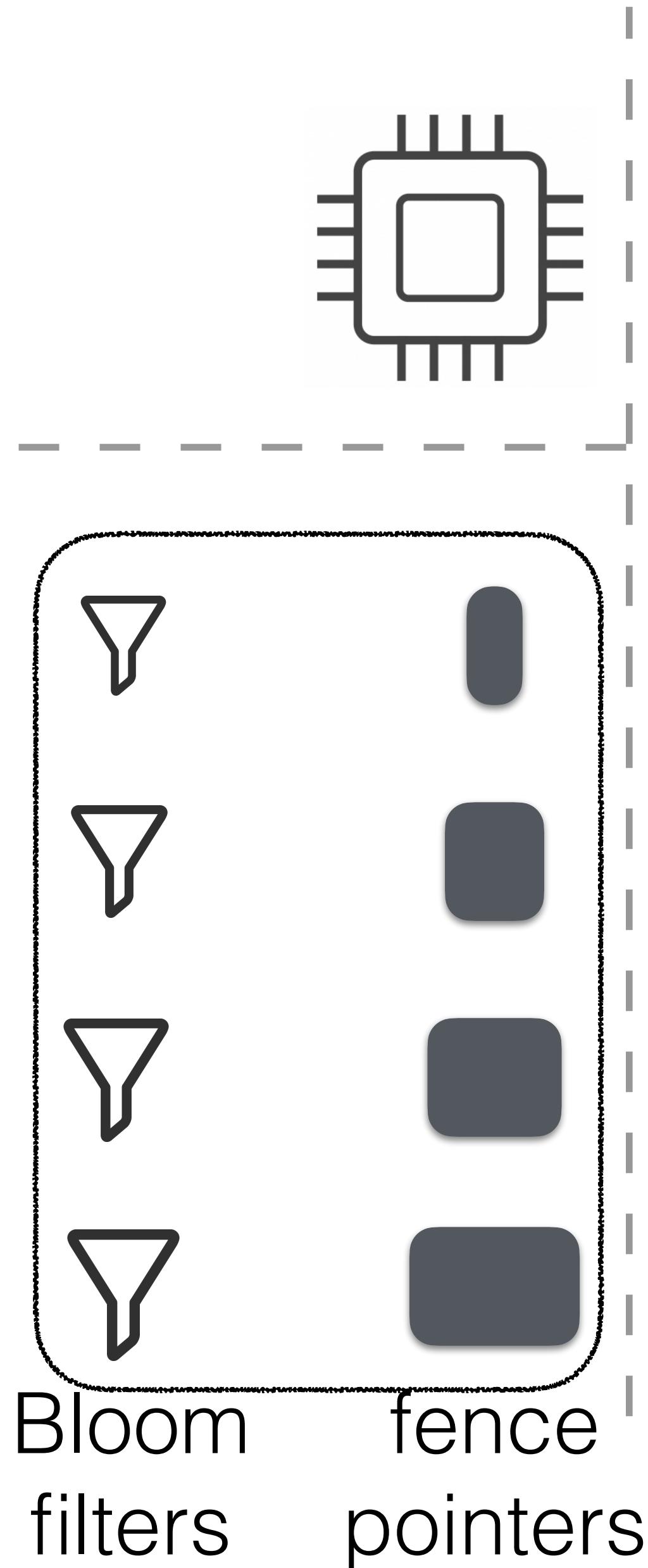
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF



buffer



block cache



Cost analysis

w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

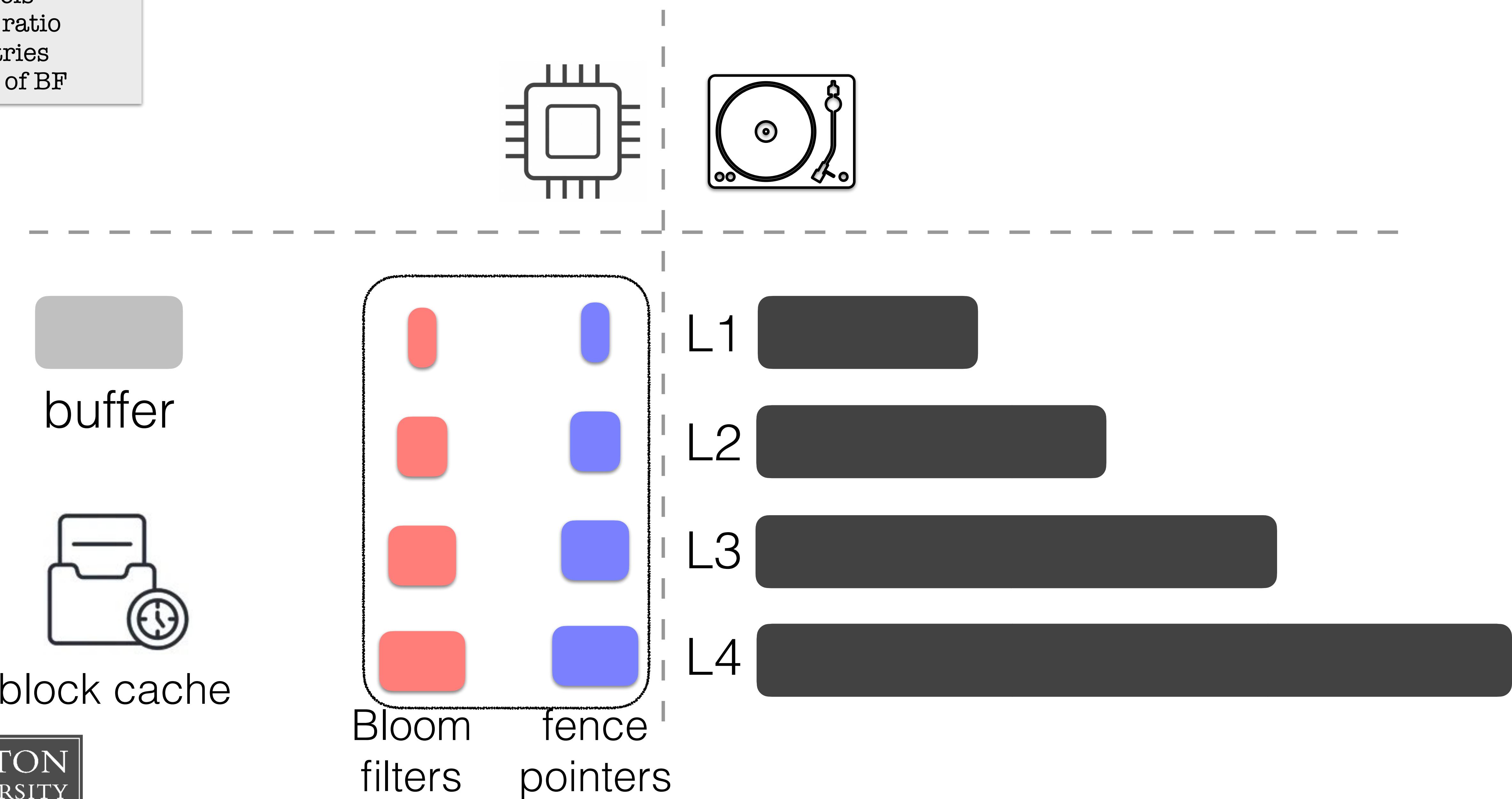
w/ index: $\mathcal{O}(L)$

w F&I: $\mathcal{O}(\phi \cdot L)$

$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2} \ll 1$$

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

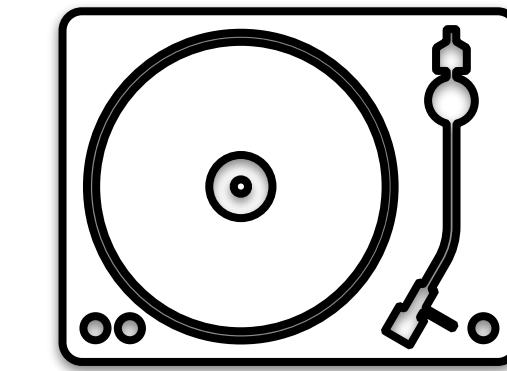
Block Cache



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

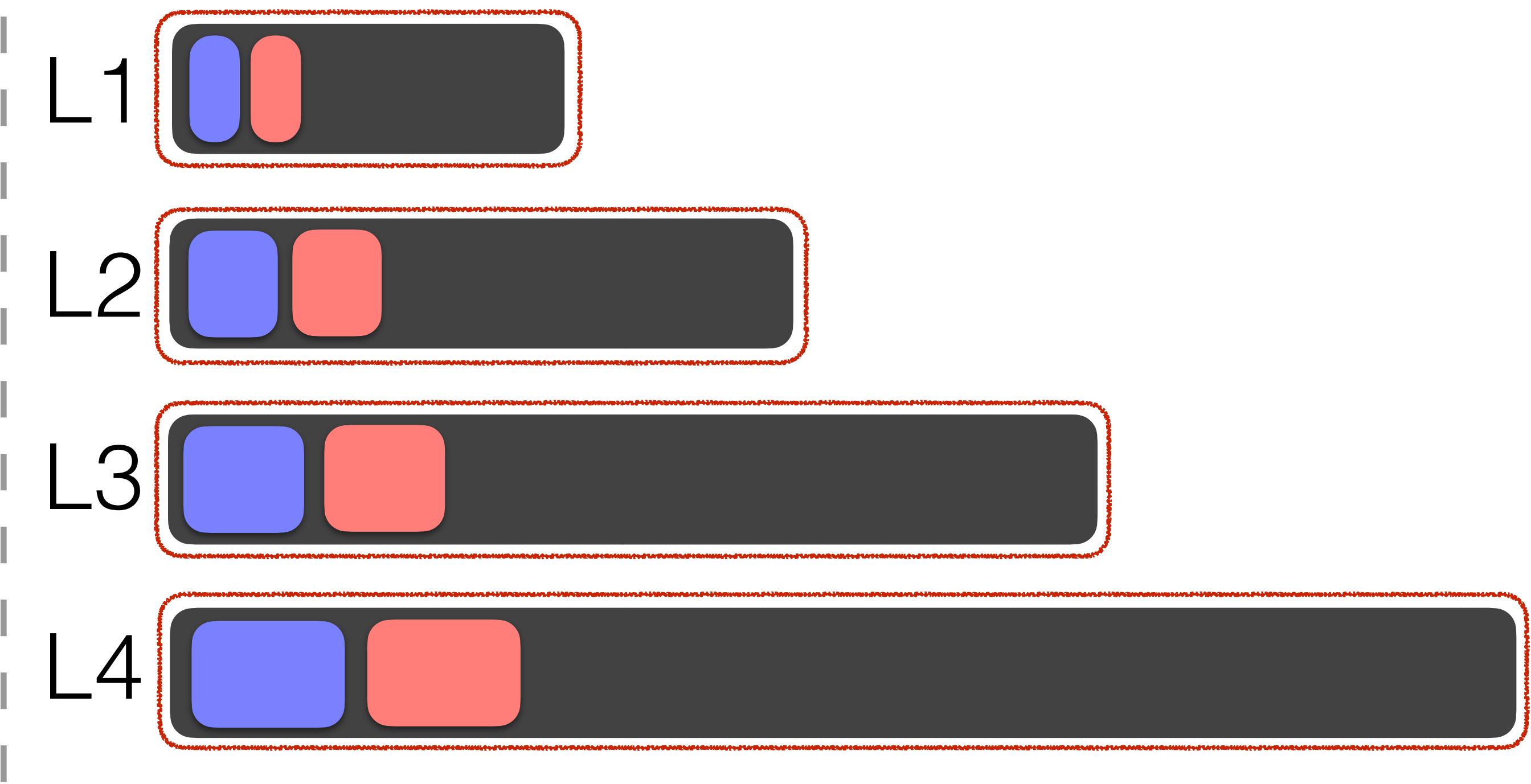
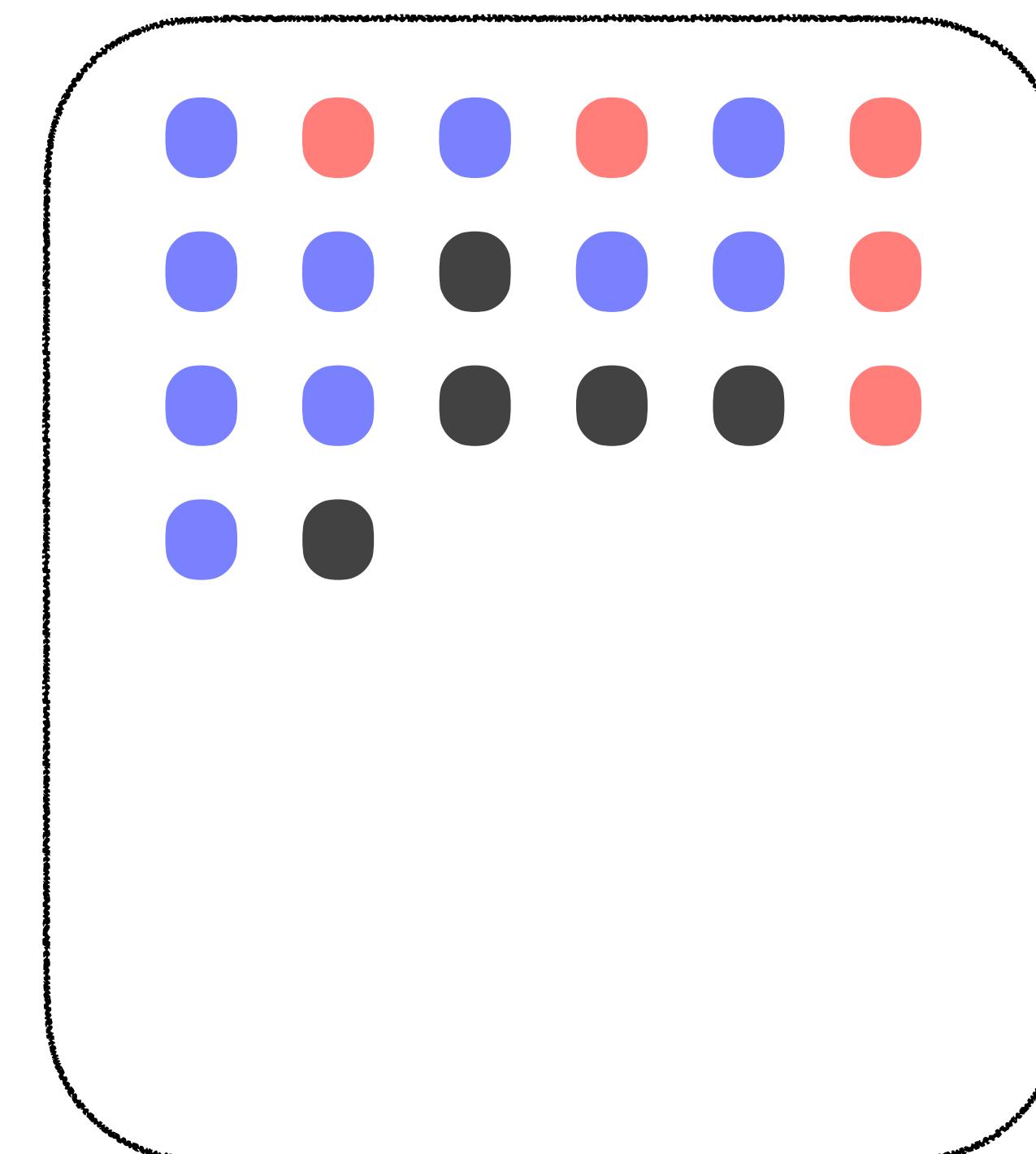
get(7)

buffer



Bloom filters fence pointers

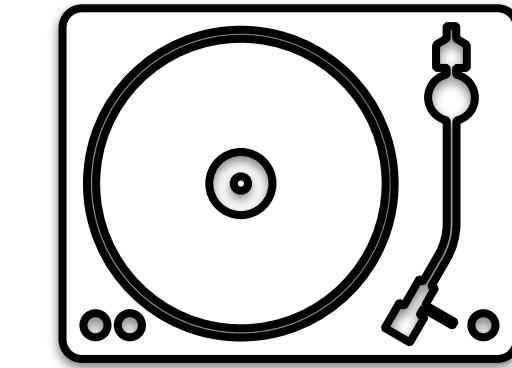
Block Cache



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

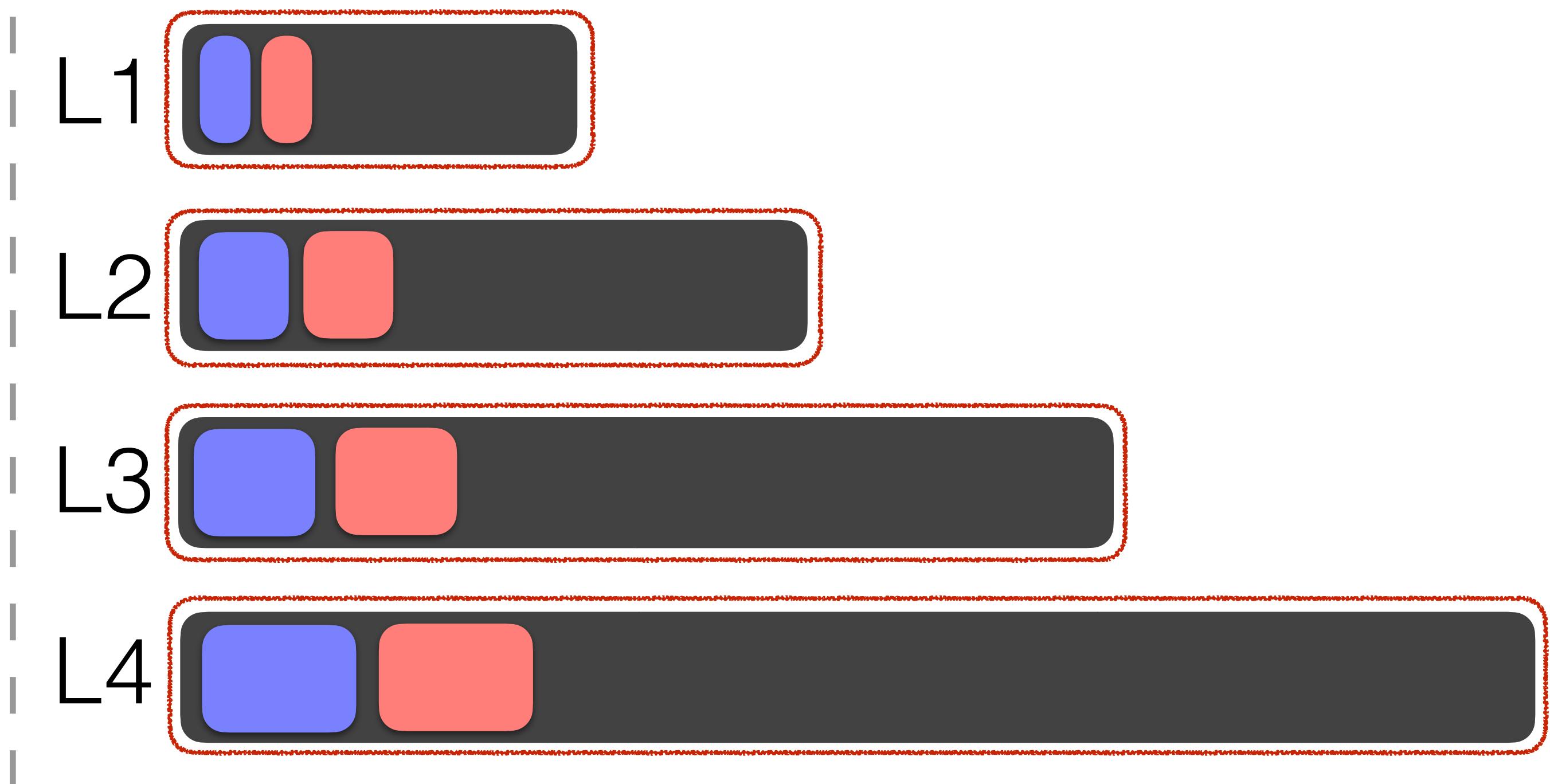
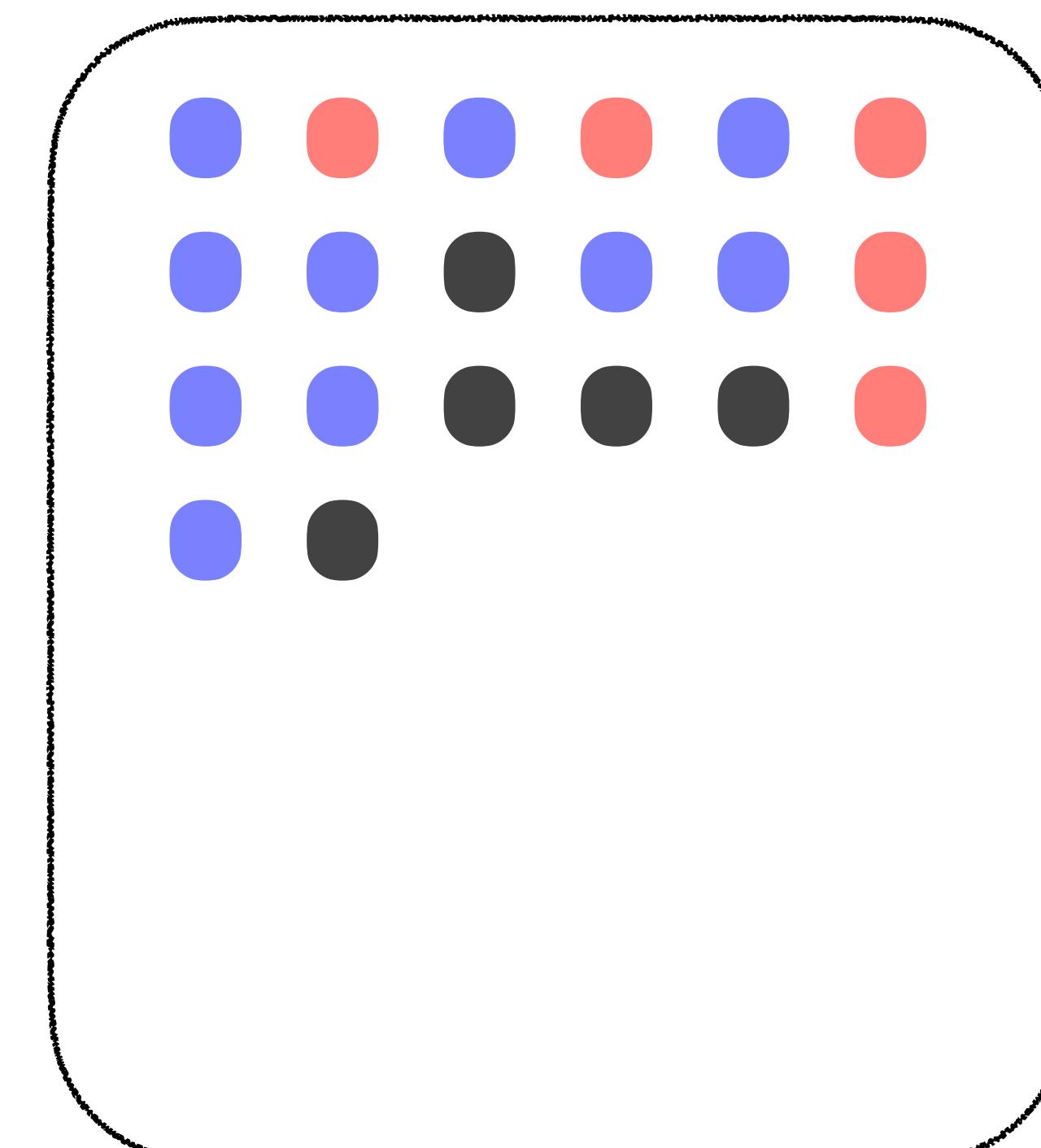
get(7)

buffer

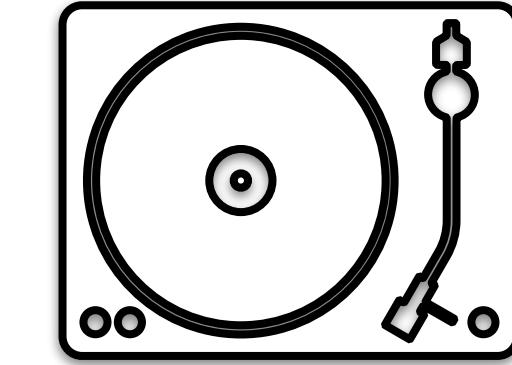
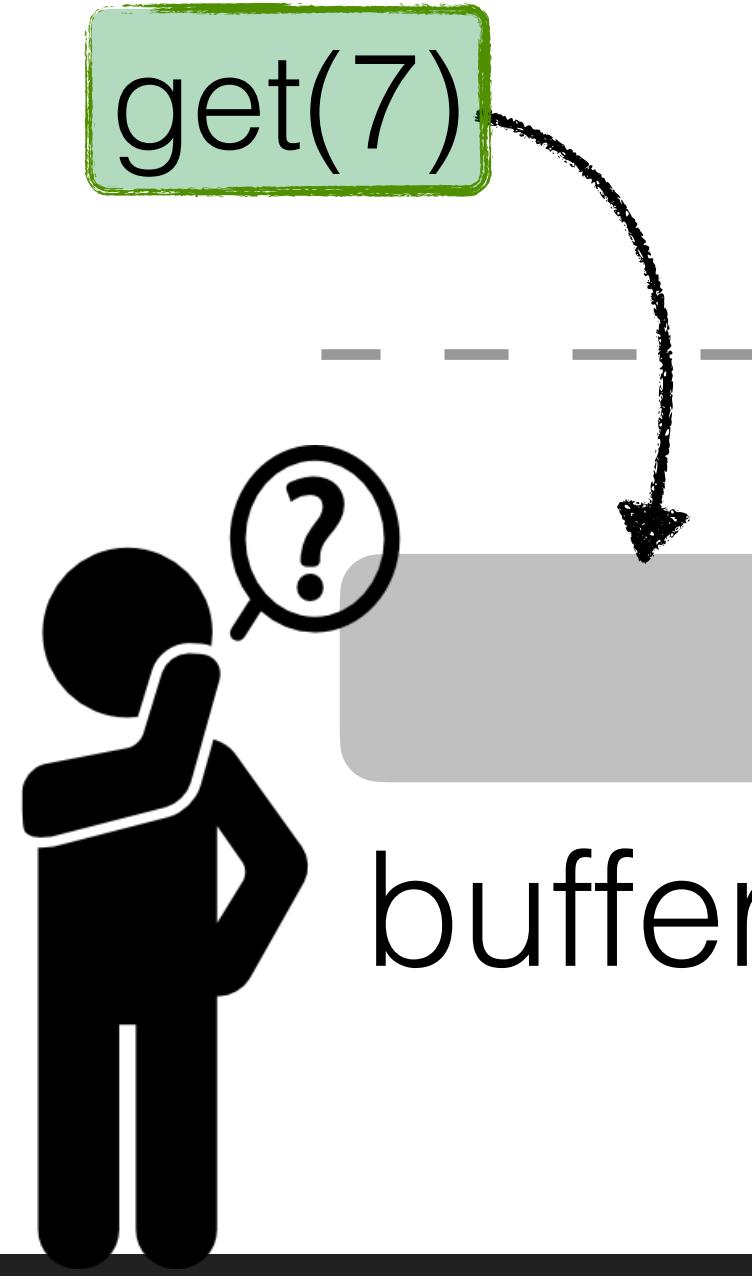


Bloom filters fence pointers

Block Cache



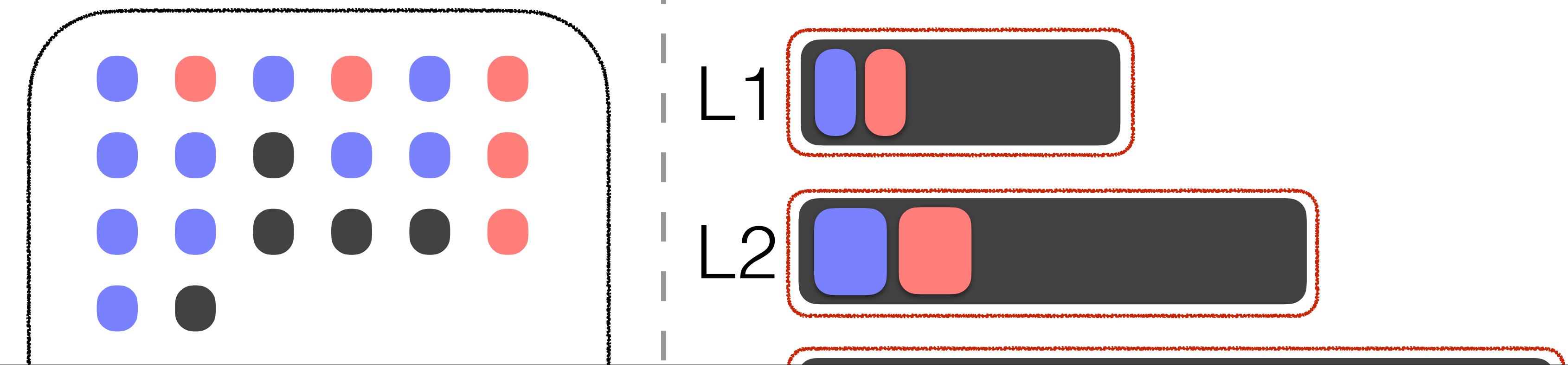
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF



Bloom filters fence pointers

Two small colored ovals are shown: a red one labeled "Bloom filters" and a blue one labeled "fence pointers".

Block Cache

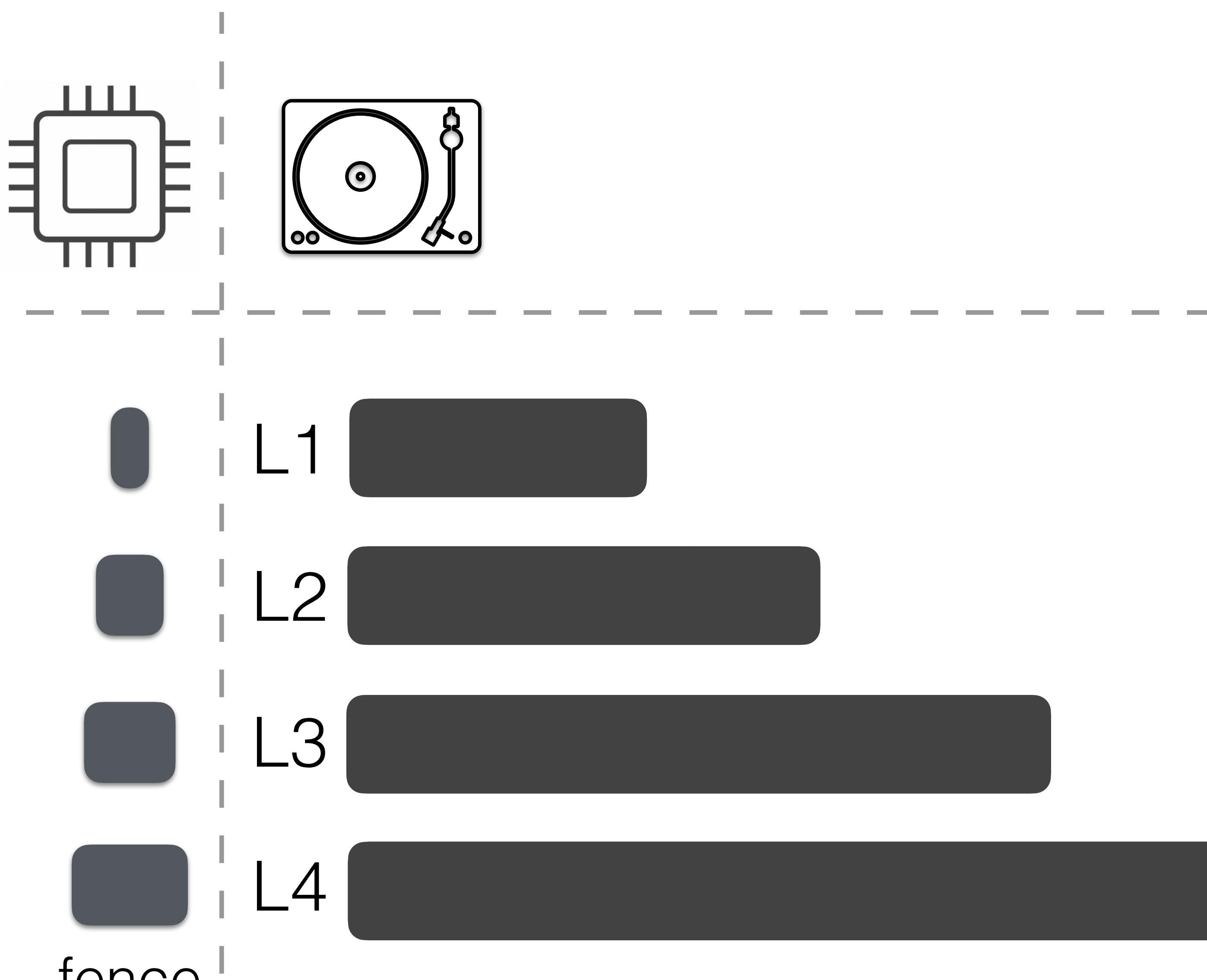


What about range queries?

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

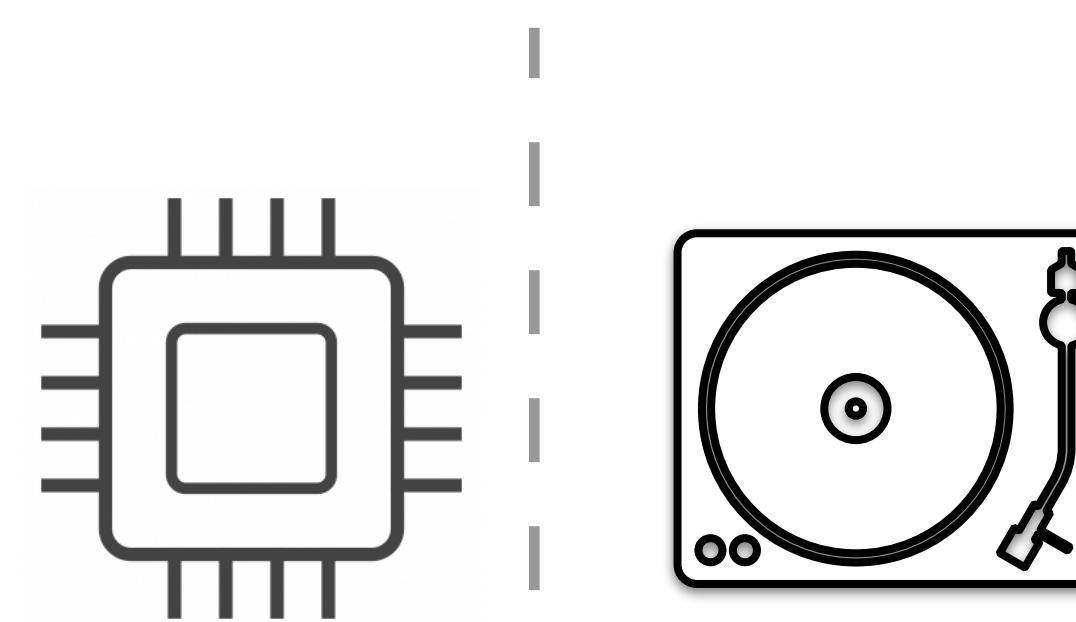
s : selectivity LRQ

Range Queries

Cost analysis

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

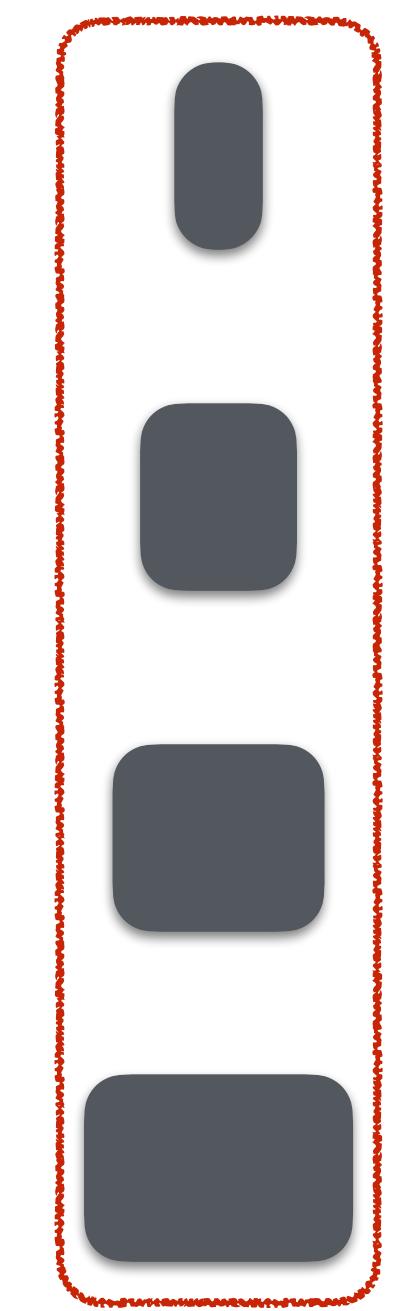
get(9,90)



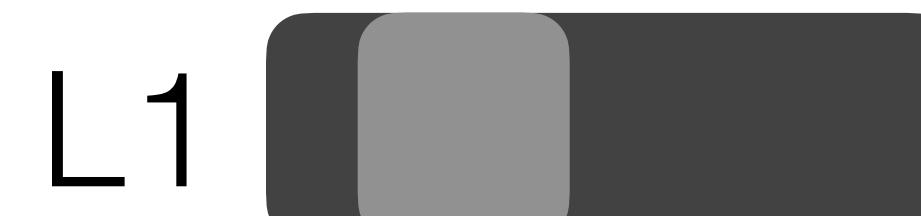
buffer



Bloom
filters



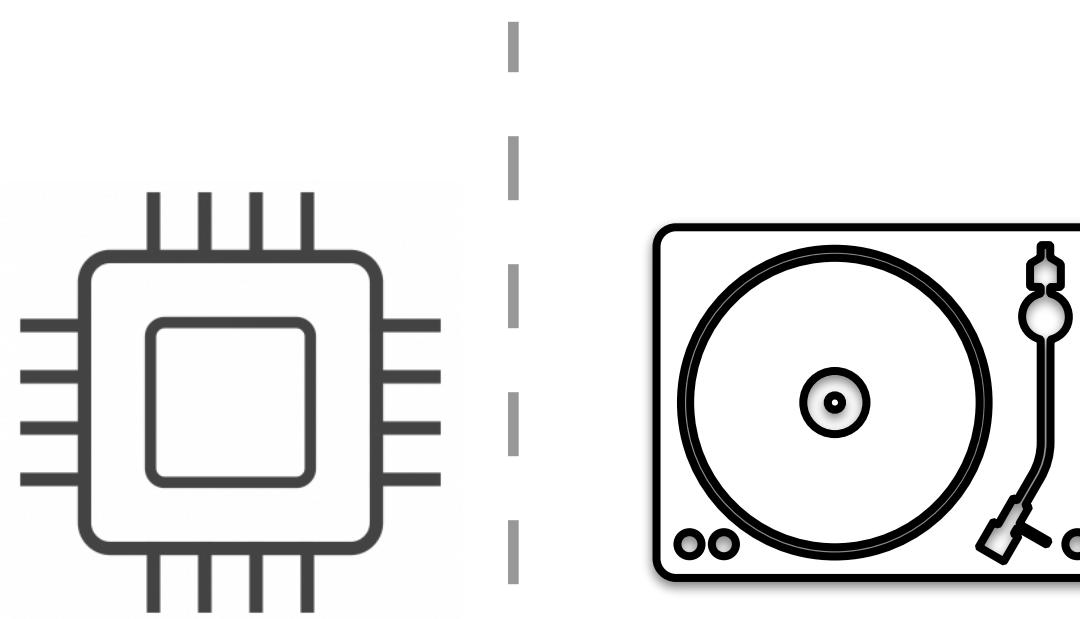
fence
pointers



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

s : selectivity SRQ

Range Queries

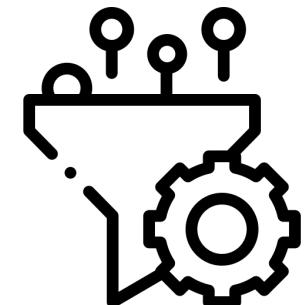


get(9,15)

Cost analysis

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

More Read Optimizations

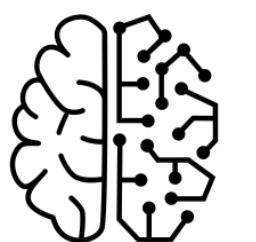


filter structures

Better performance

Chucky

DayanSIGMOD21



indexes

Better performance

Learned

DaiOSDI20



block cache

Less memory usage

Ribbon filter

DillingerArxiv21

Less memory usage

FerraginaVLDB20

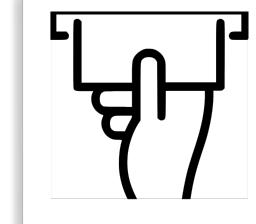
Better performance

Less memory usage

Less CPU usage

Shared Hashing

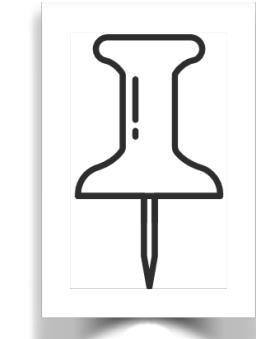
ZhuDAMON21



Prefetching



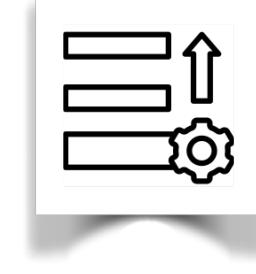
HuangVLDB20



Pinning



CallaghanSD18



Priority



CallaghanSD16

Support range queries

SuRF

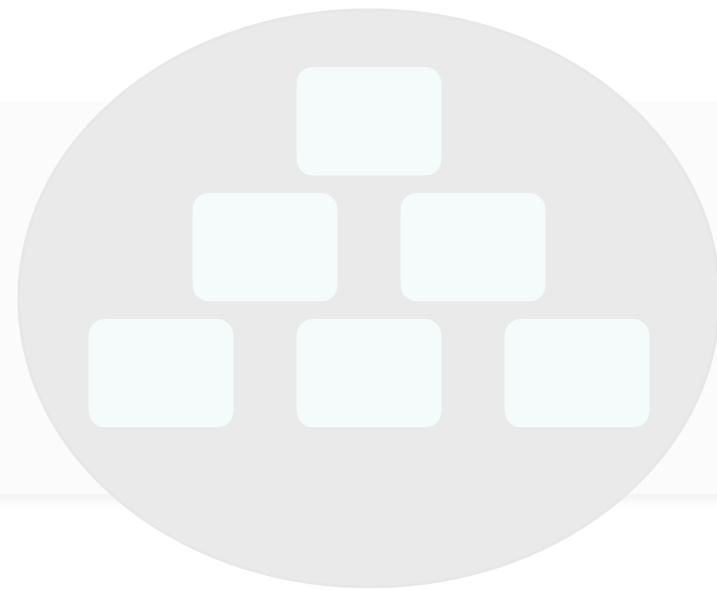
ZhangSIGMOD18

Rosetta

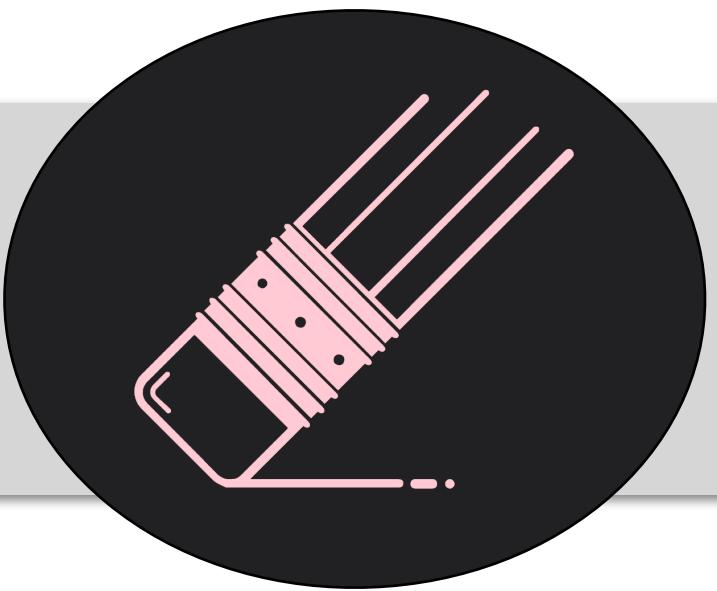
LuoSIGMOD20

Outline

Part 1: **LSM Basics**



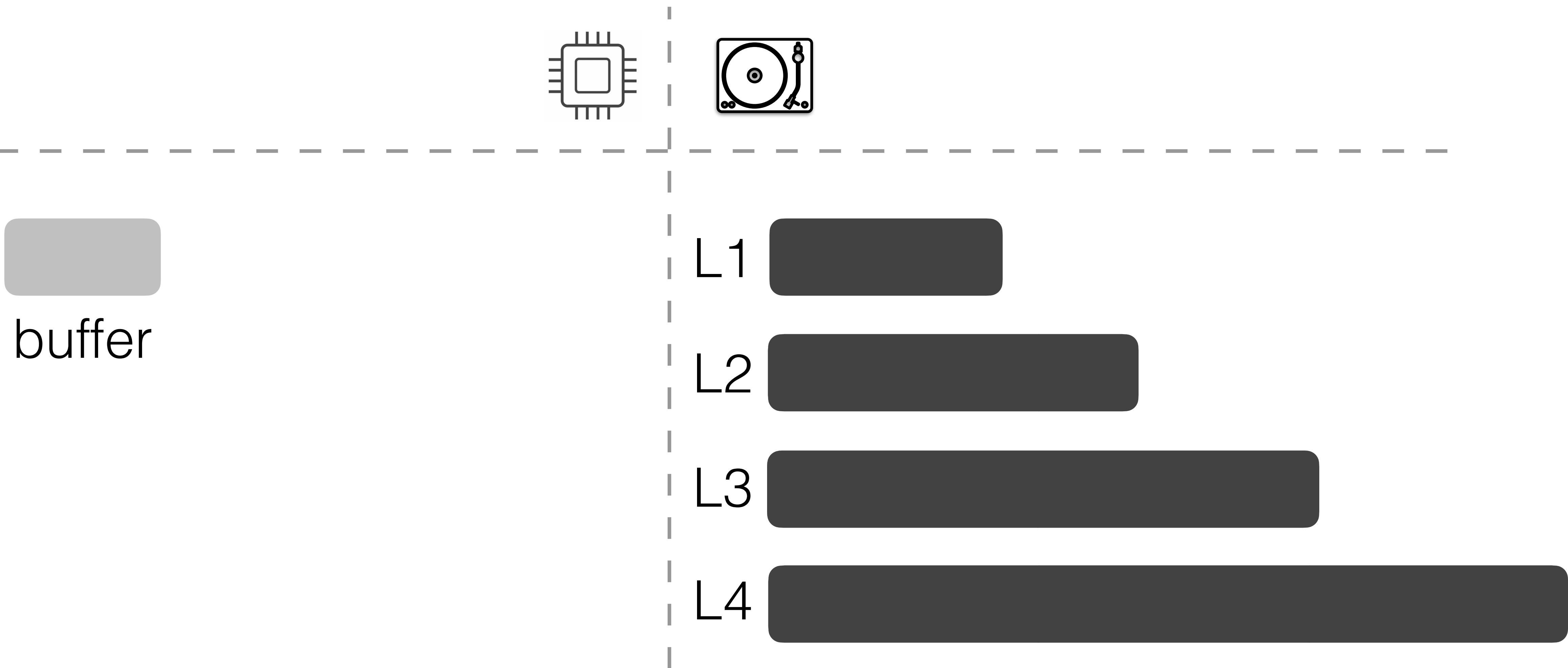
Part 2: **Optimizing Ingestion in LSMS**



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



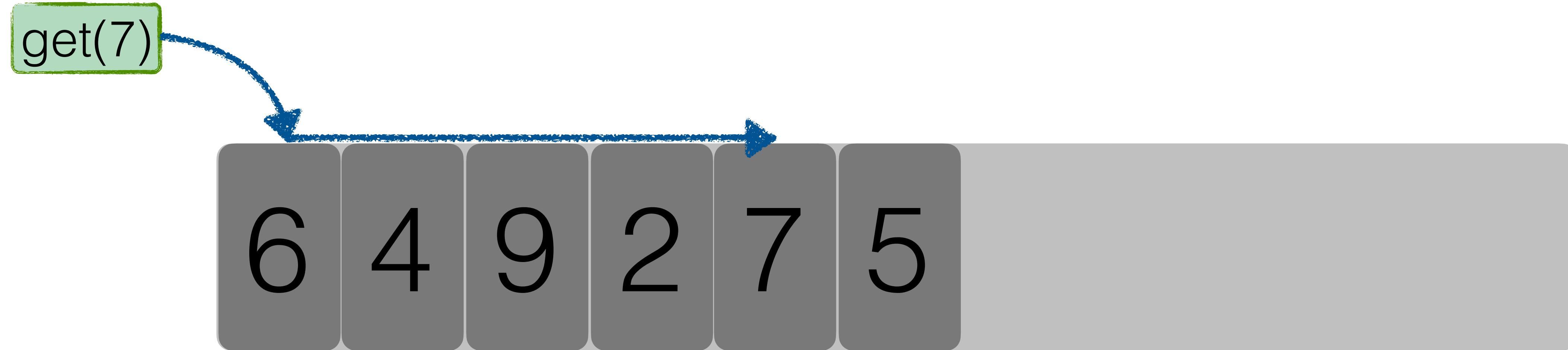
Buffer Optimizations



P : pages in buffer

B : entries/page

Buffer Implementation: **vector**



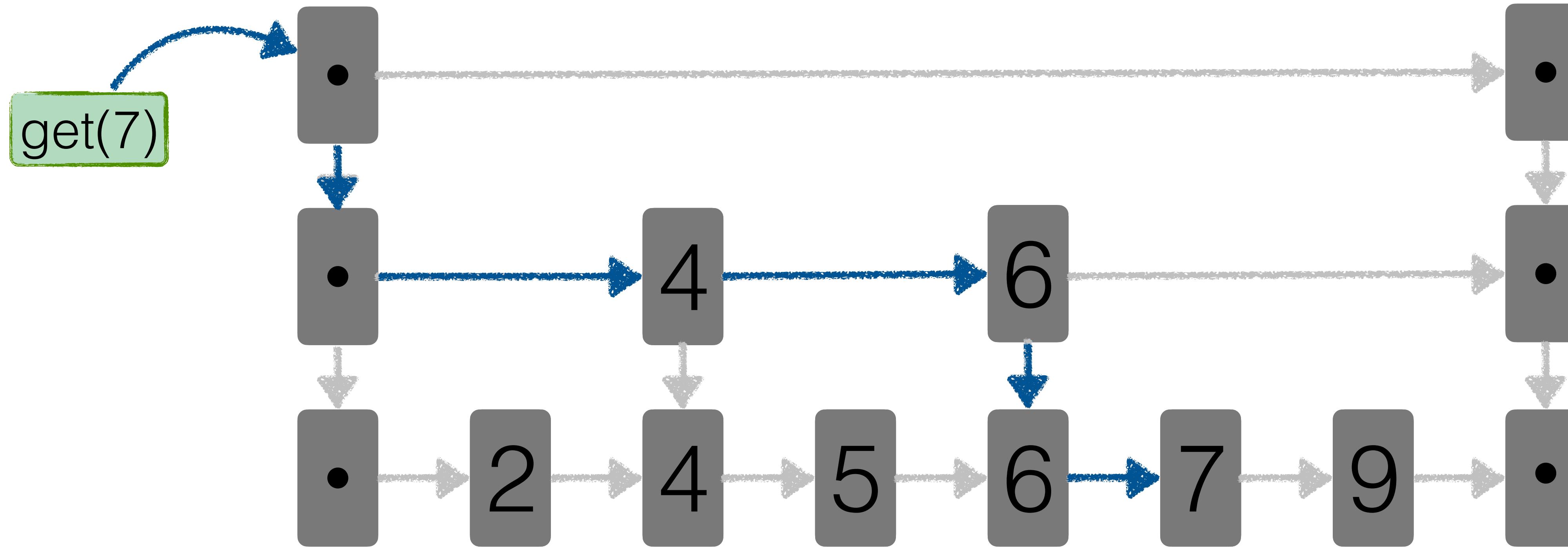
- great for ingestion-heavy w/l
- no extra space needed
- expensive points queries

ingestion cost: $\mathcal{O}(1)$

space complexity: $\mathcal{O}(P \cdot B)$

point query cost: $\mathcal{O}(P \cdot B)$

Buffer Implementation: **skiplist**



- great for mixed w/l
- some extra space needed
- good for points queries



P : pages in buffer

B : entries/page

Buffer Implementation

ingestion
cost

space
complexity

point query
cost

vector

$\mathcal{O}(1)$

$\mathcal{O}(P \cdot B)$

$\mathcal{O}(P \cdot B)$

skiplist

$\mathcal{O}(\log(P \cdot B))$

$\mathcal{O}(P \cdot B)$

$\mathcal{O}(\log(P \cdot B))$

hashmap

$\mathcal{O}(1)$

$\mathcal{O}(P \cdot B)$

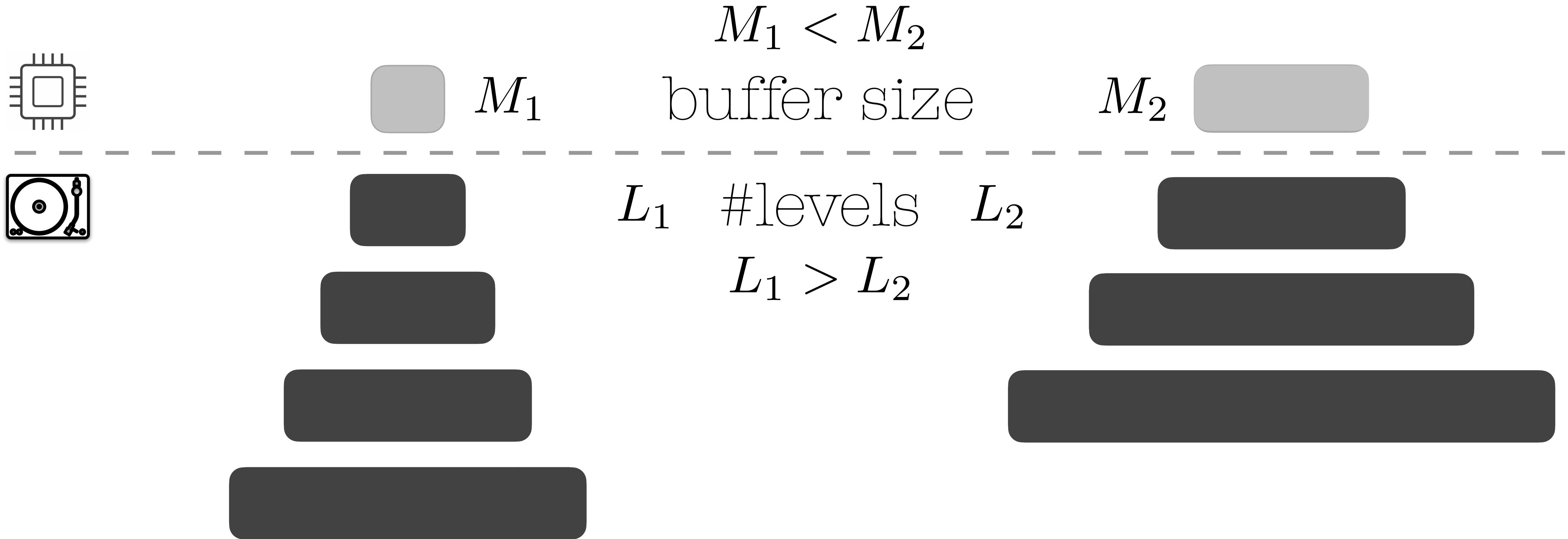
$\mathcal{O}(1)$

Ingestion-only
workloads

Mixed
workloads

I/O-bound
workloads

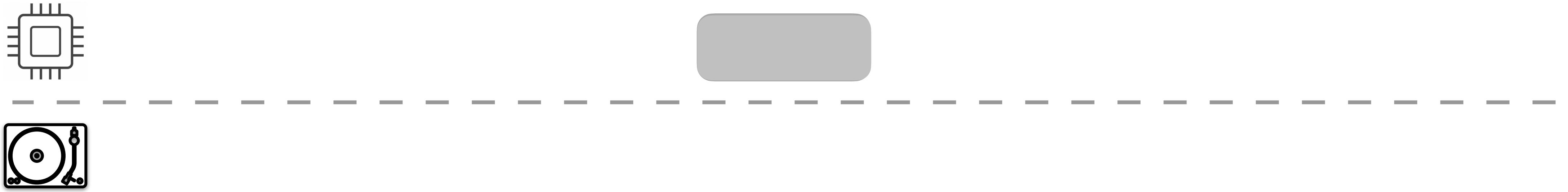
Size of the Buffer



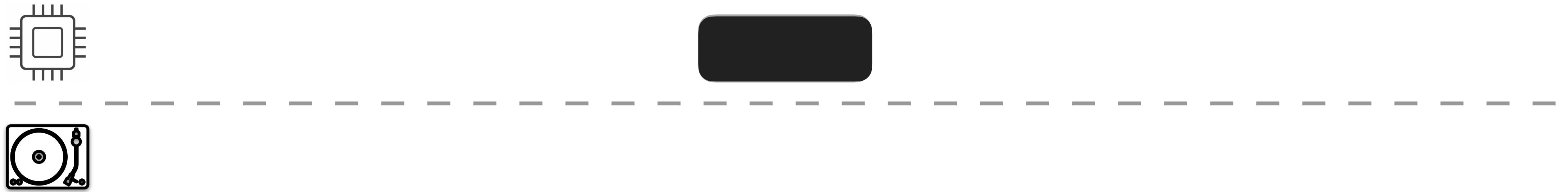
- frequent flushes
- smaller but more levels
- poor read performance

- fewer larger levels
- good for reads
- high tail latency

#Buffer Components

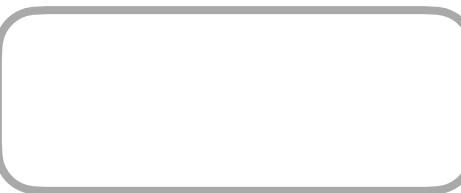
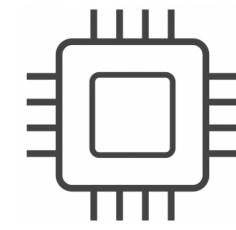


#Buffer Components



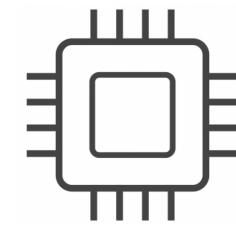
#Buffer Components

immutable
buffers



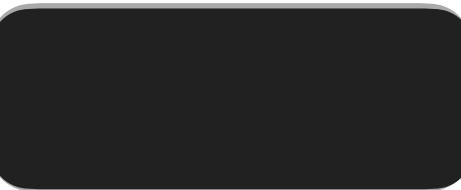
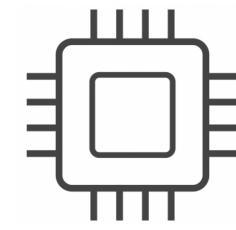
#Buffer Components

immutable
buffers



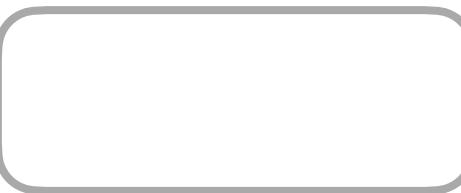
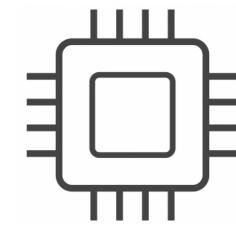
#Buffer Components

immutable
buffers



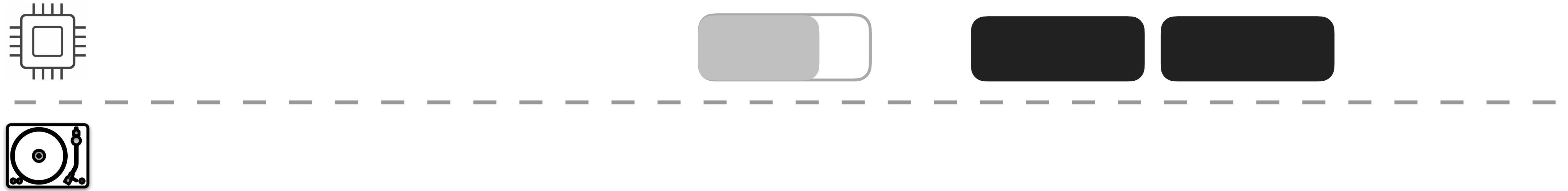
#Buffer Components

immutable
buffers



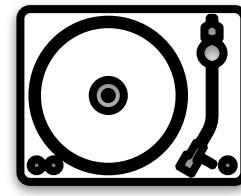
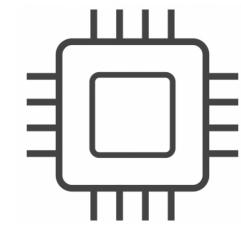
#Buffer Components

immutable
buffers



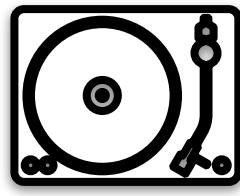
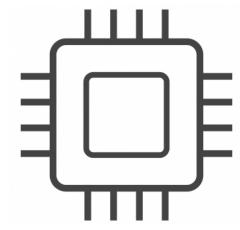
#Buffer Components

immutable
buffers



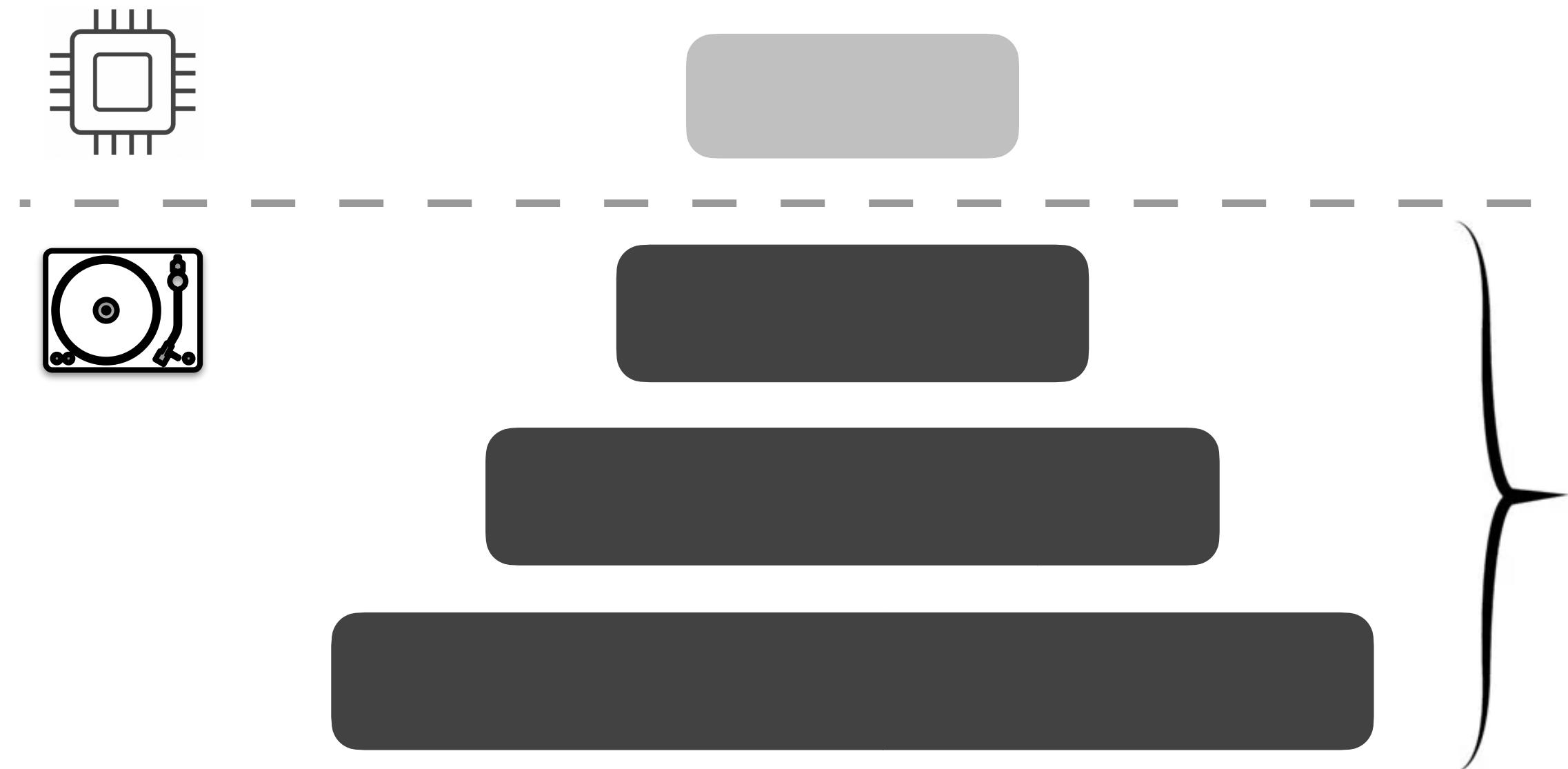
#Buffer Components

immutable
buffers



lazy flushing

- avoid write stalls
- improved ingestion throughput
- better bandwidth utilization
- requires more memory

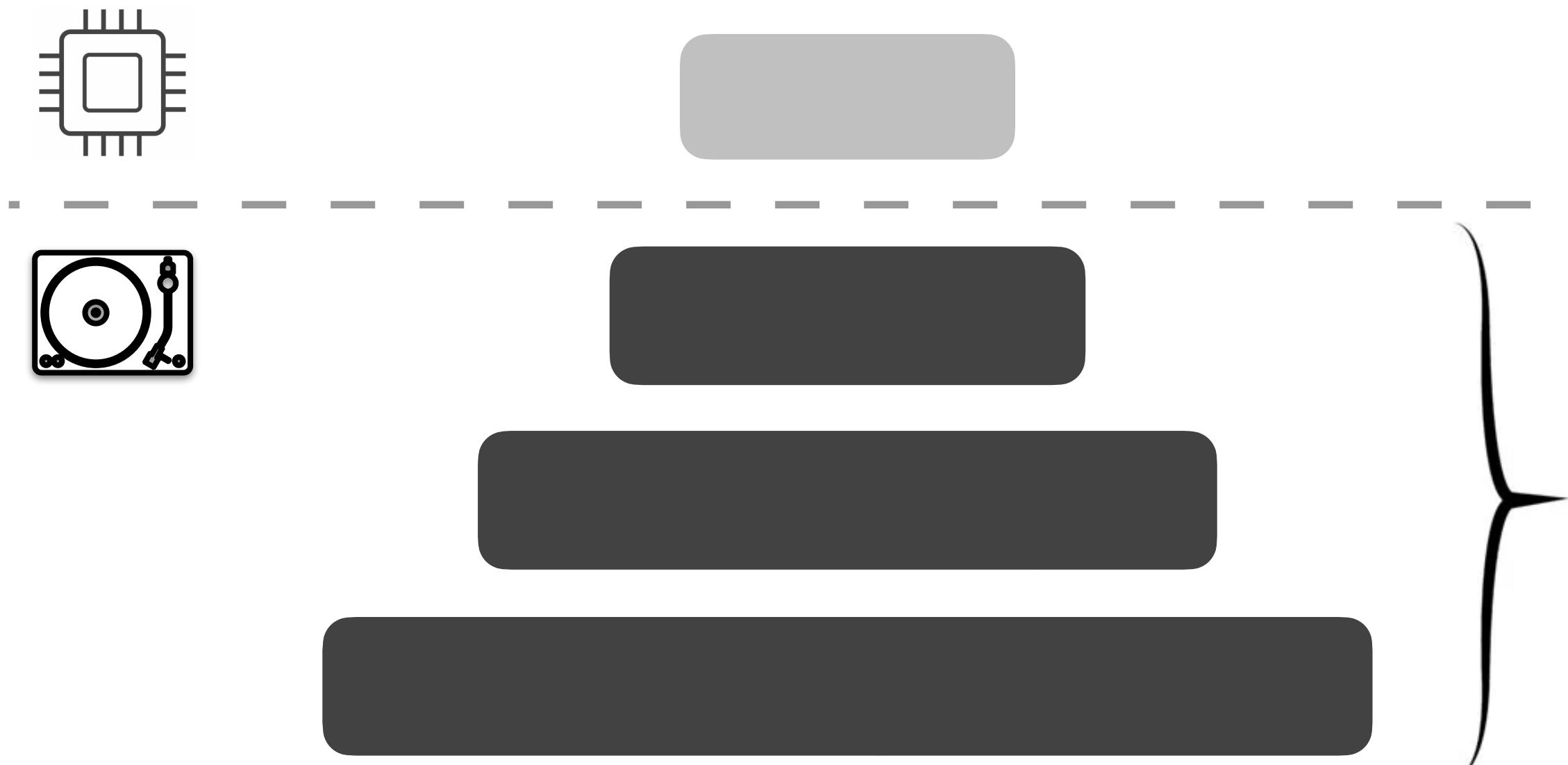


most data
on storage

How does the storage layer affect ingestion?

L : #levels

T : size ratio



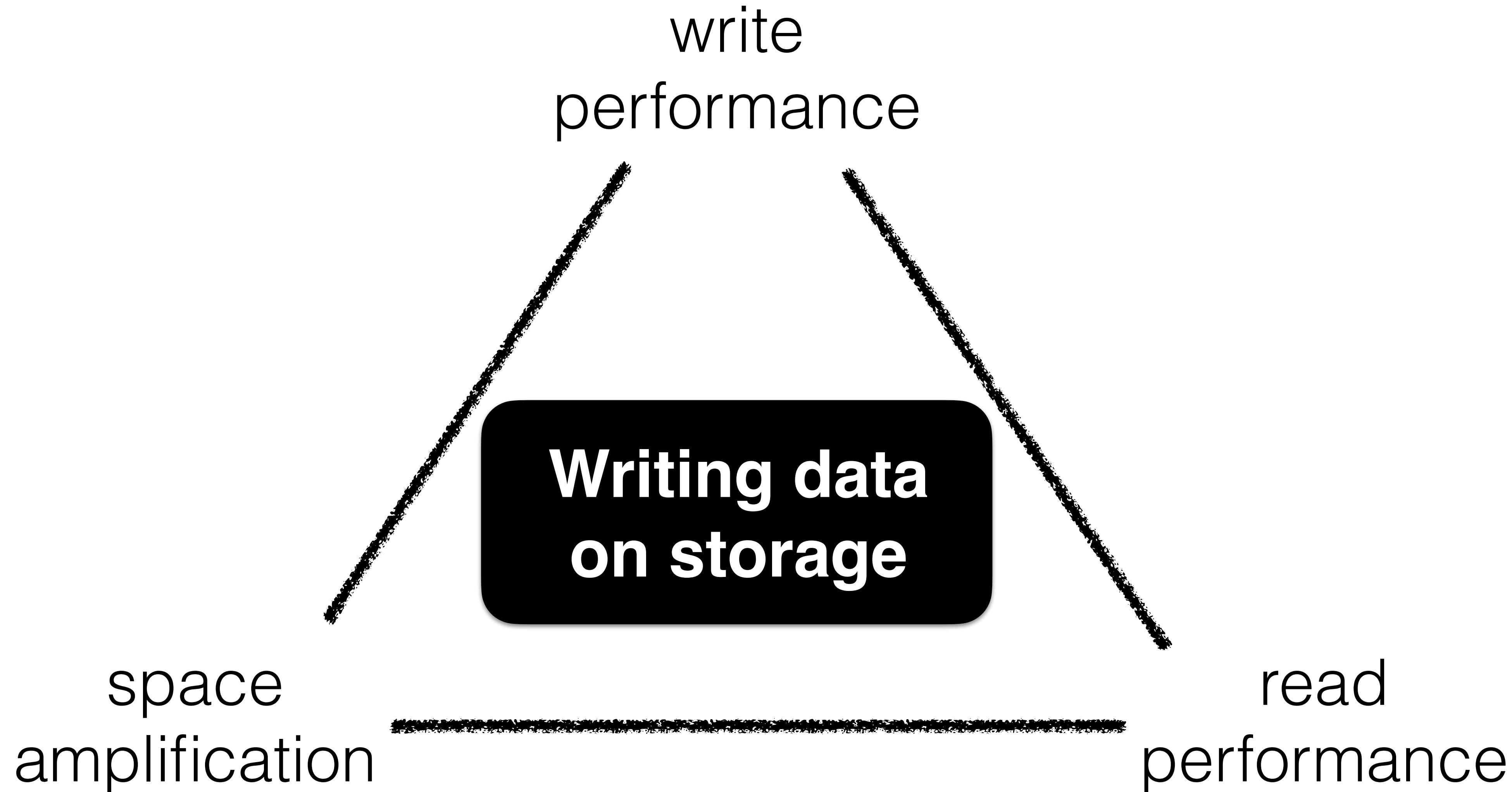
most data
on storage

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

99.9% on storage

How does the storage layer affect ingestion?

Storage Optimizations



Data Layout

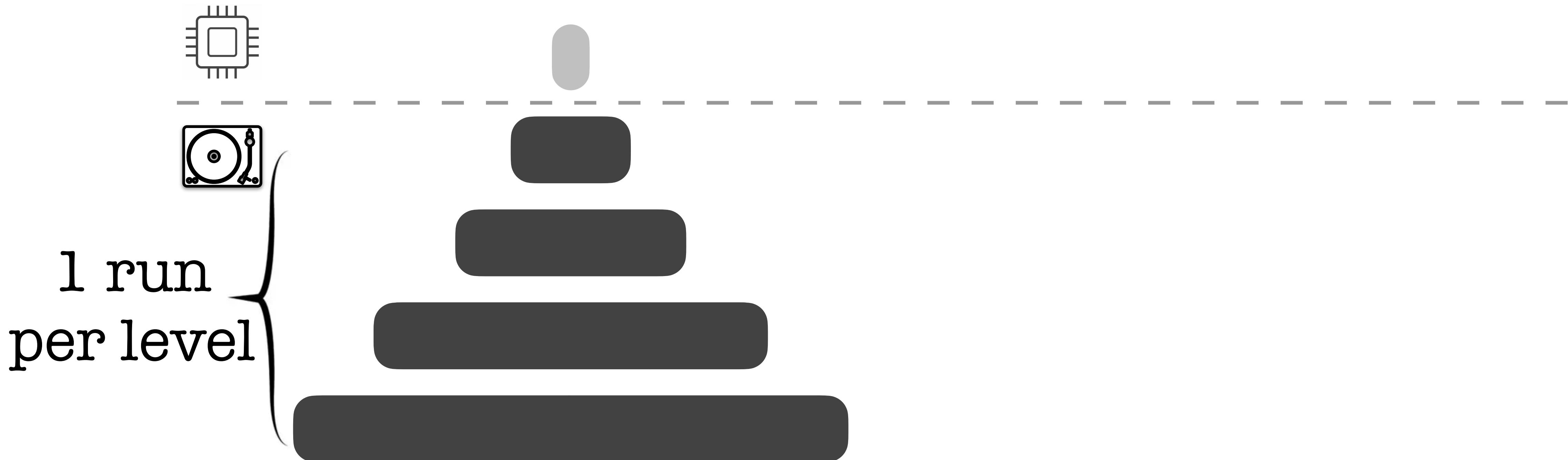
Classical LSM design: leveling

[eager merging]



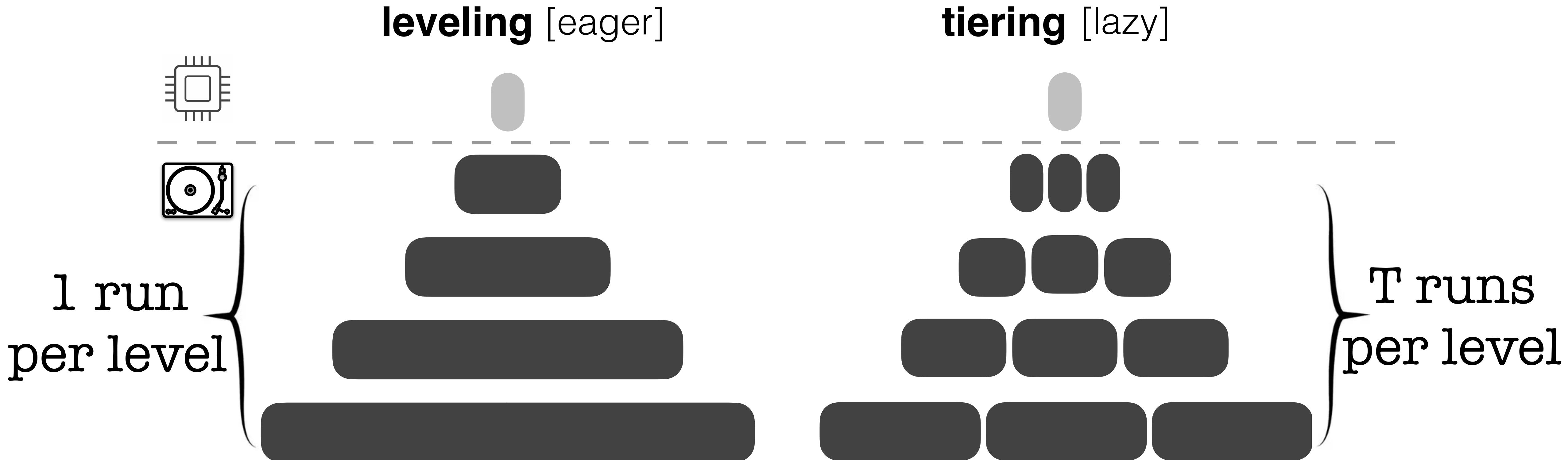
Data Layout

leveling [eager]



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

Data Layout

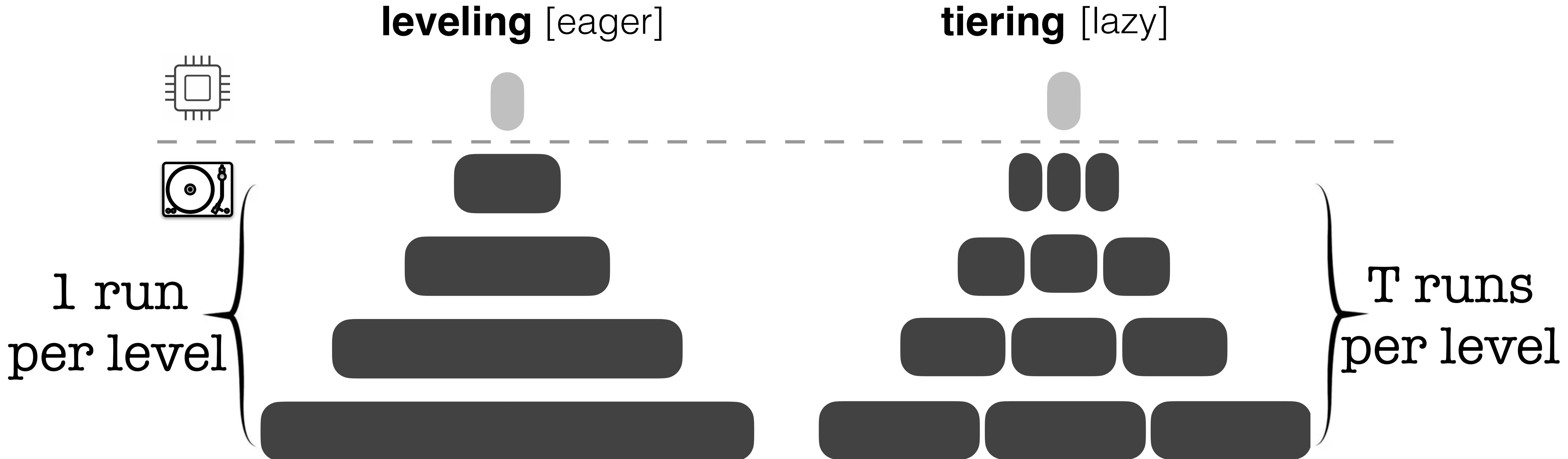


- good read performance
- good space amplification
- high write amplification

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

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

Data Layout



Read cost:

$$\mathcal{O}(L \cdot \phi)$$

Write cost:

$$\mathcal{O}(T \cdot L/B)$$

SA:

$$\mathcal{O}(1/T)$$

$$\mathcal{O}(T \cdot L \cdot \phi)$$

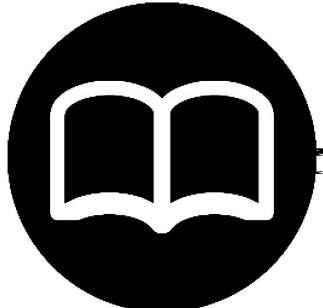
$$\mathcal{O}(L/B)$$

$$\mathcal{O}(T)$$

Data Layout

hybrid designs

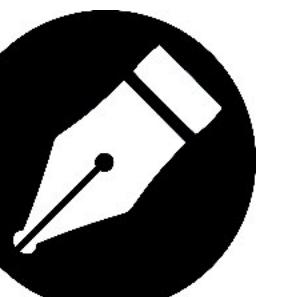
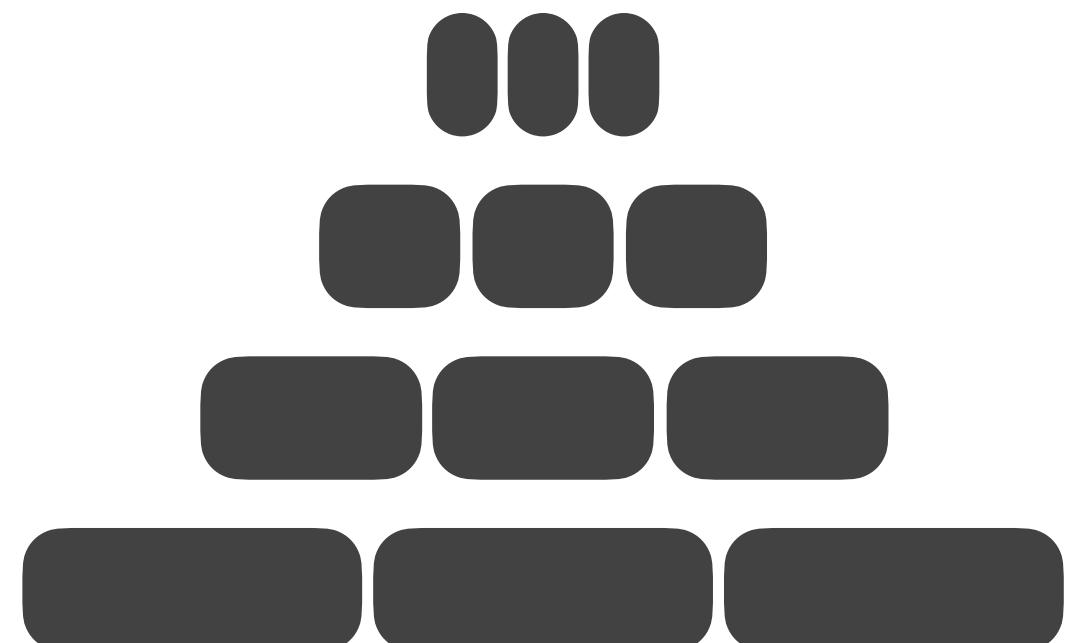
leveling



read

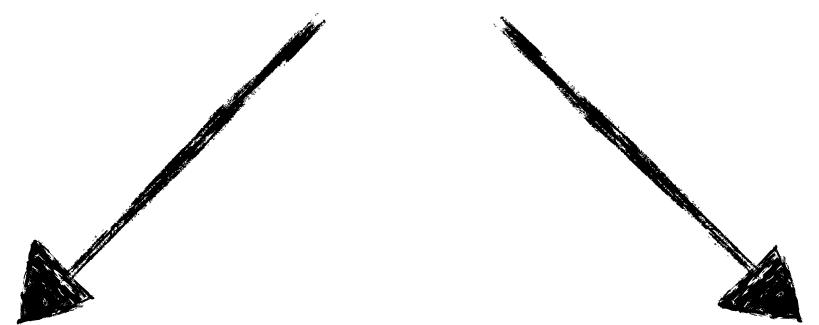
optimized

tiering

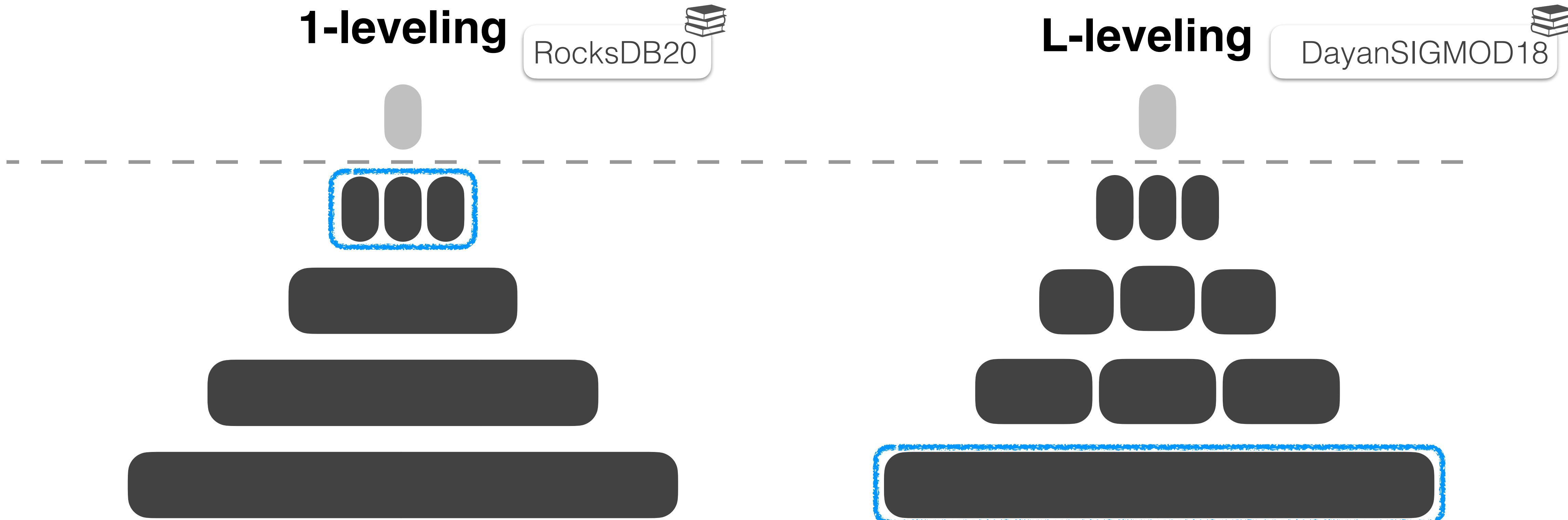


write

optimized



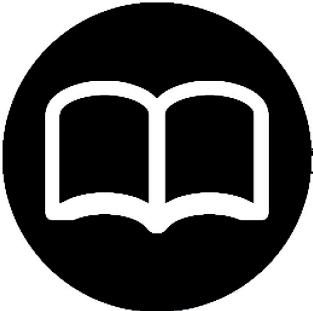
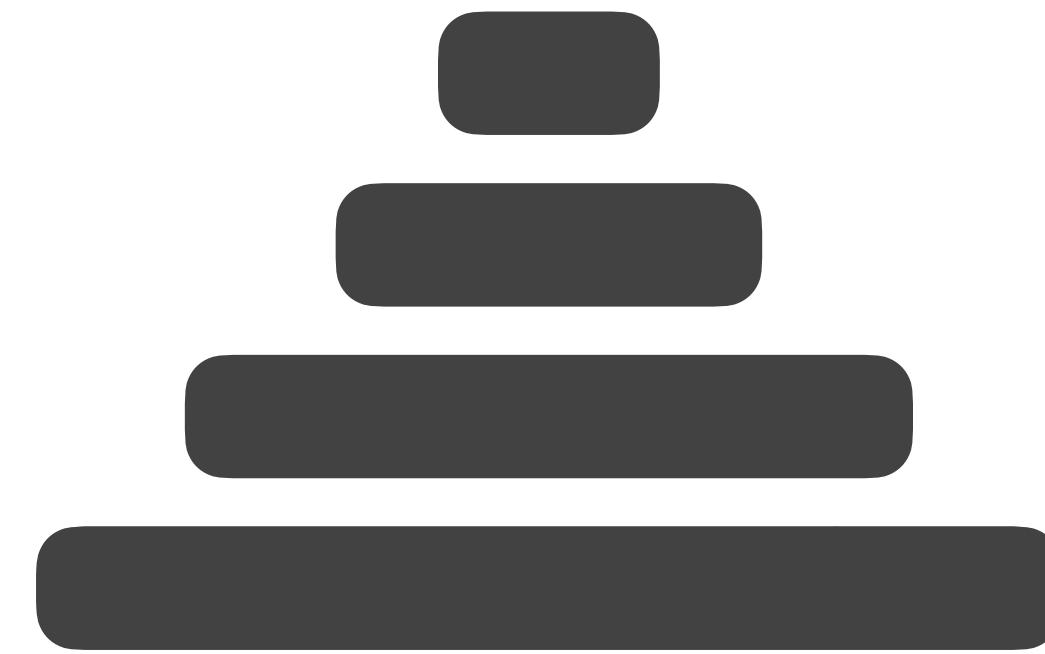
Data Layout



- fewer write stalls
- increased block cache hits
- low write amplification
- better read performance

Data Layout

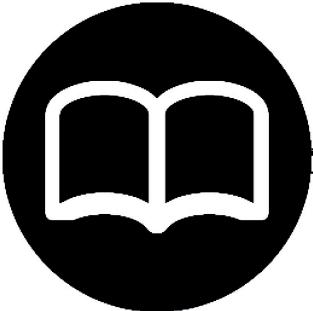
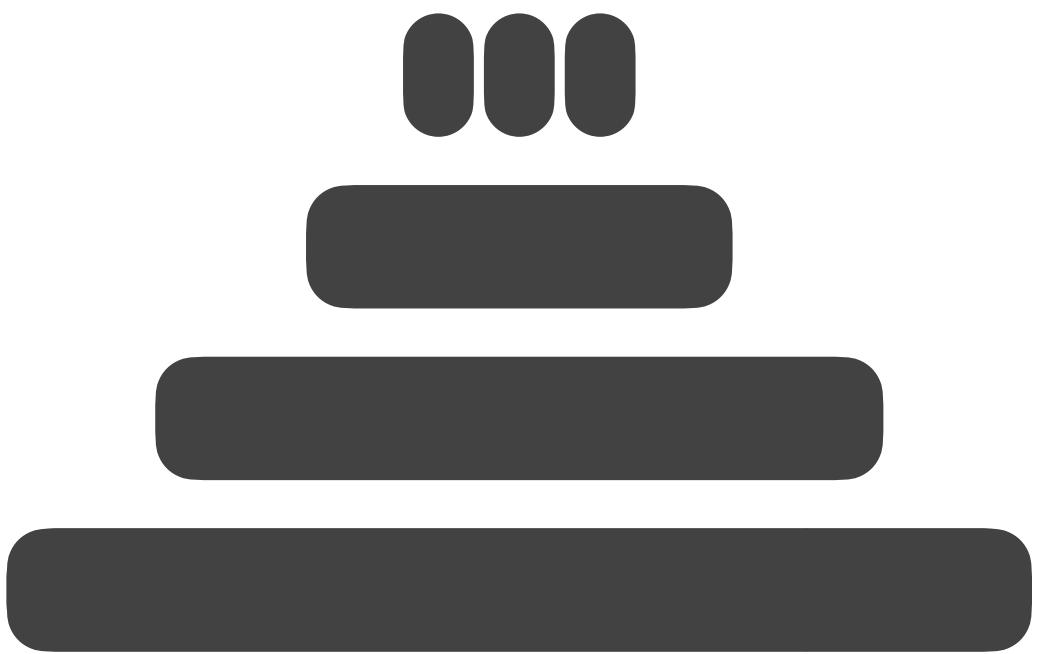
leveling



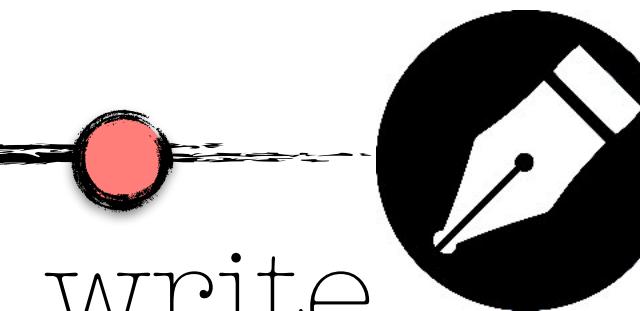
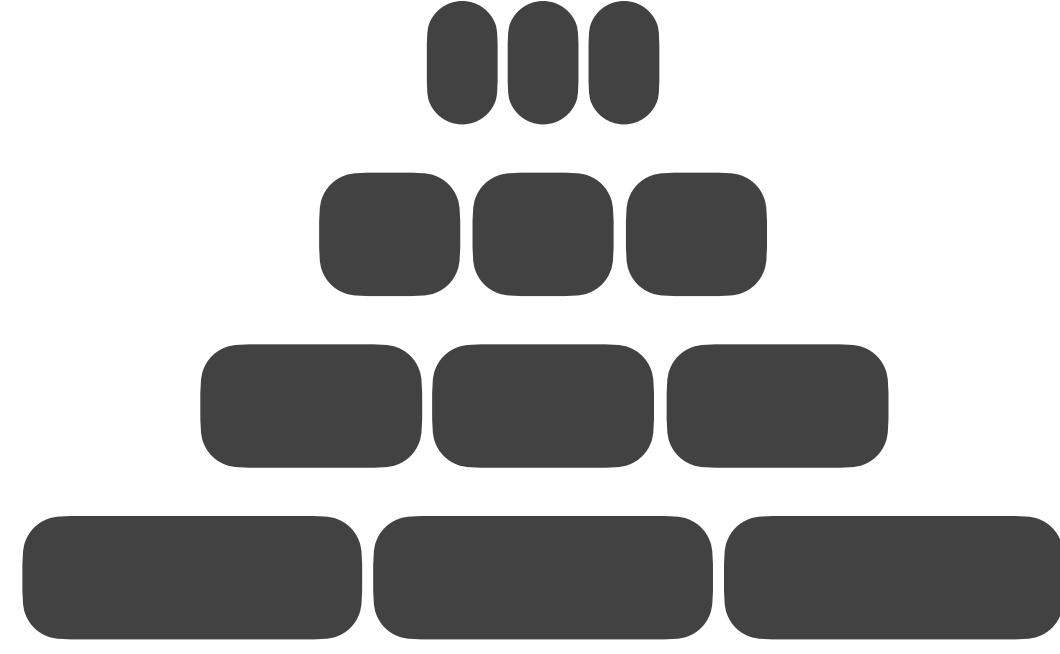
read

optimized

1-leveling

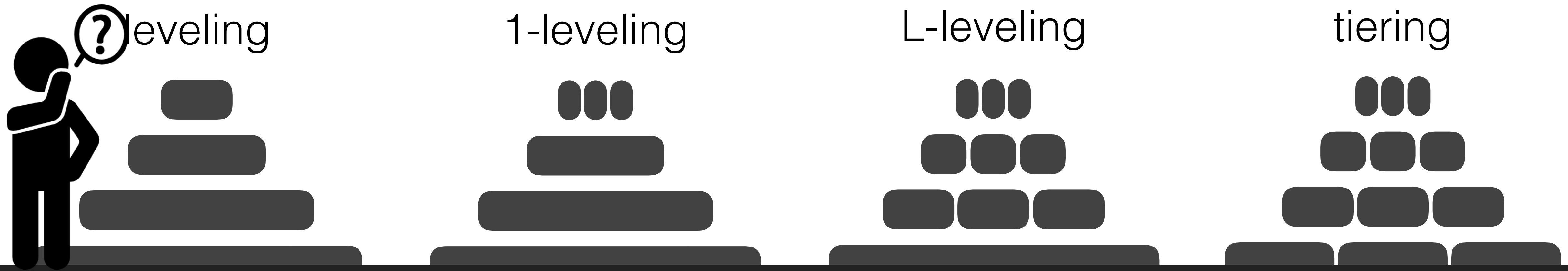


tiering

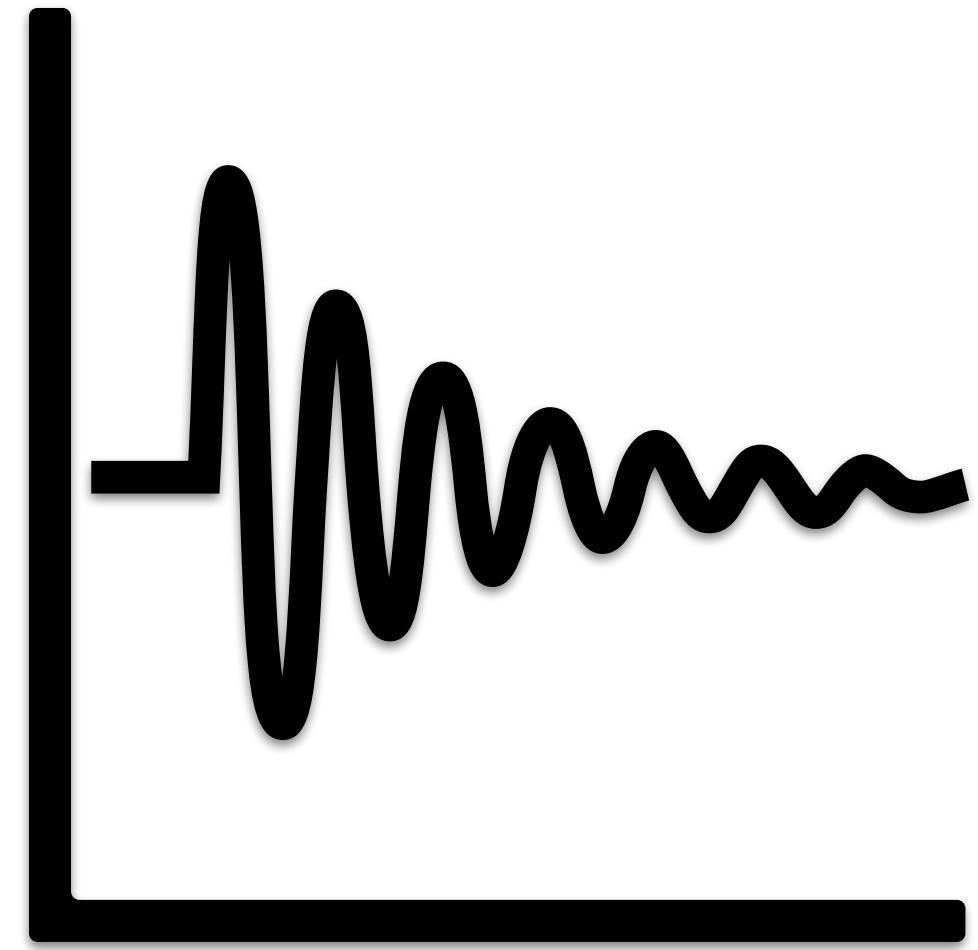


write
optimized

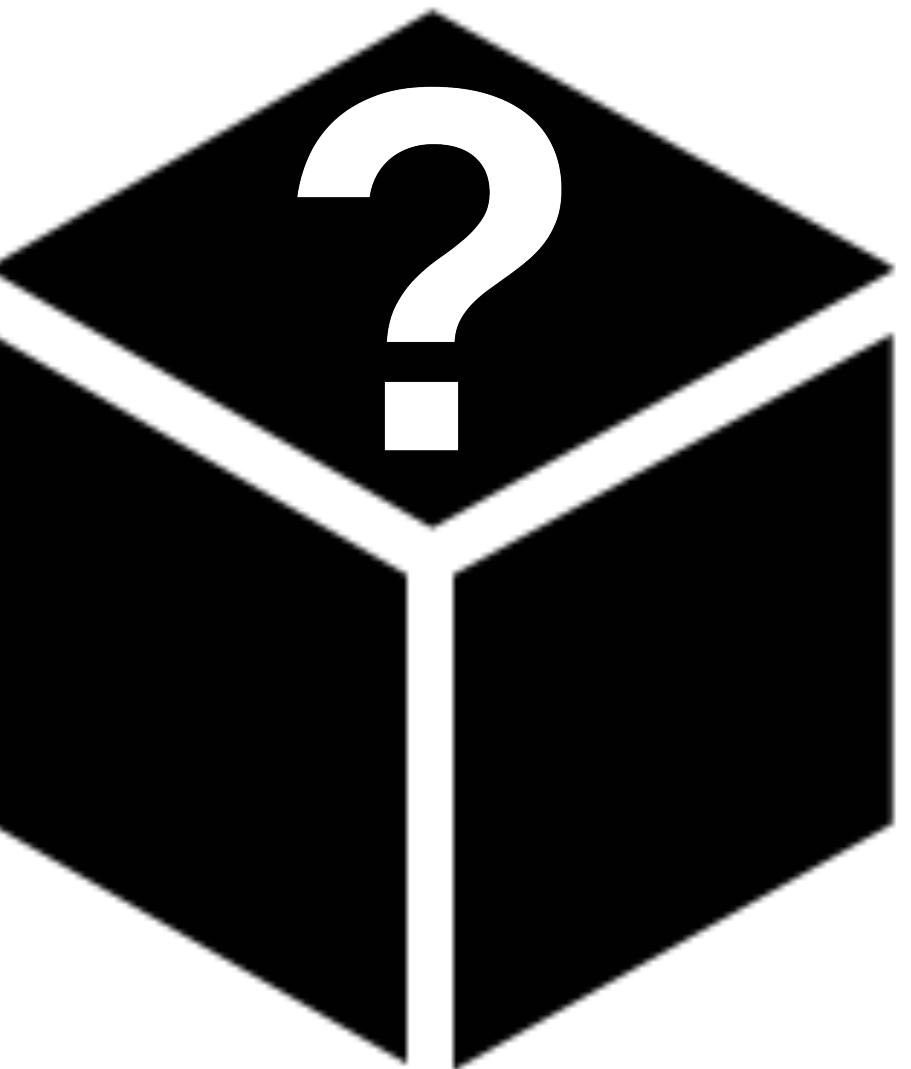
Data Layout



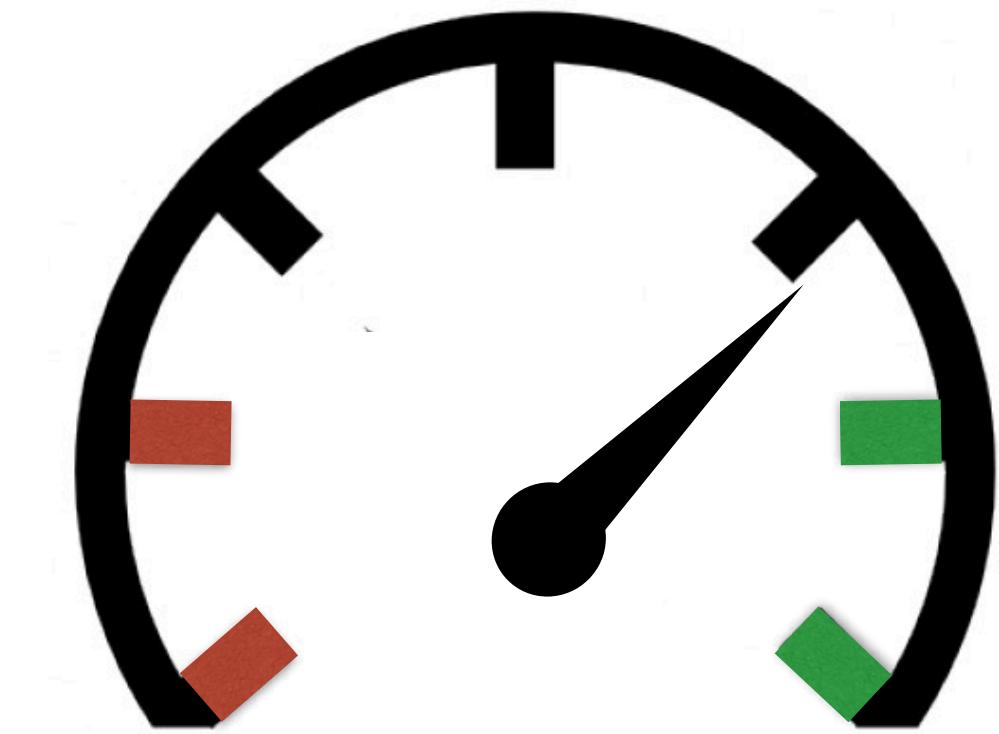
So, how do we reason about the data layout?



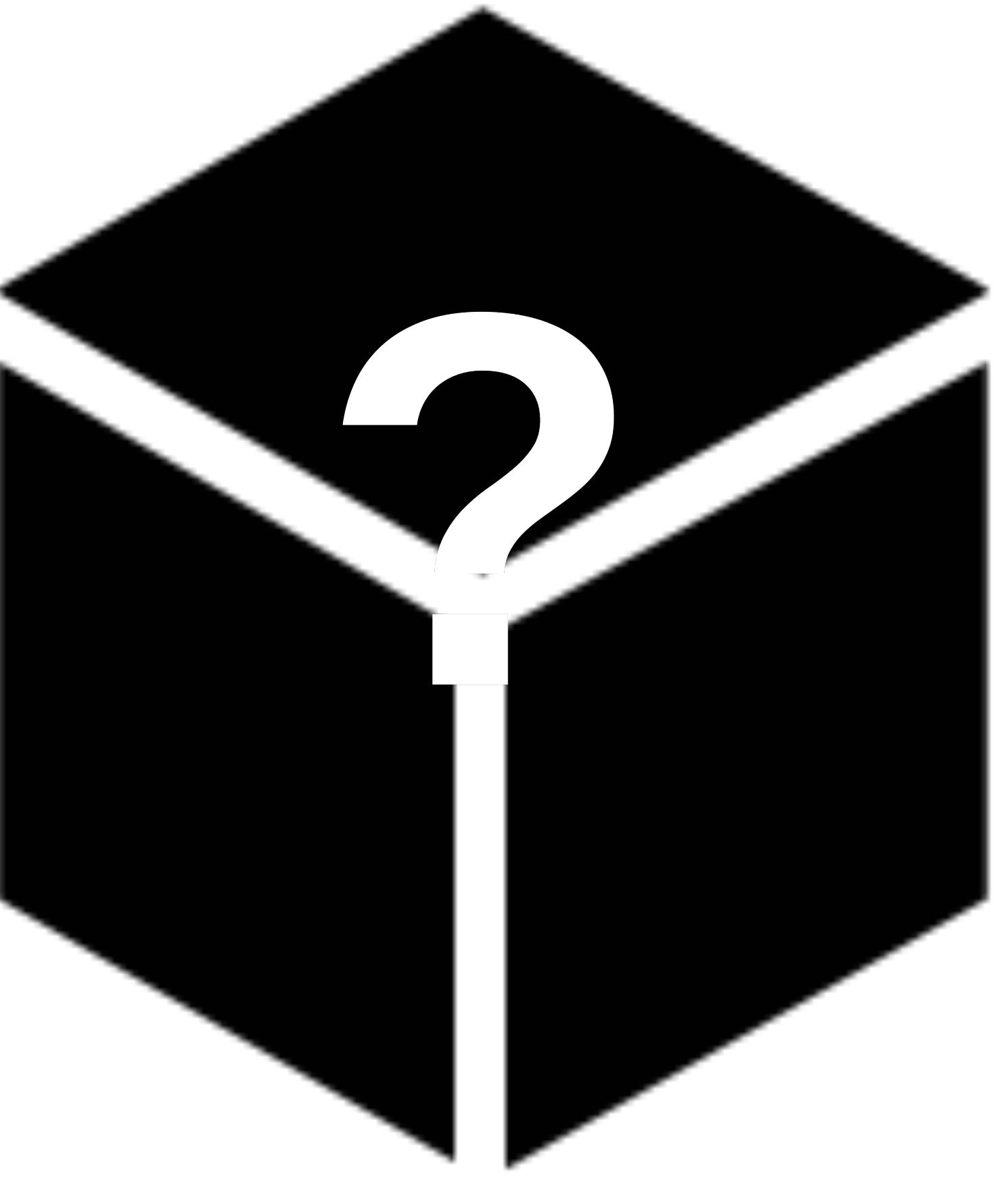
workload



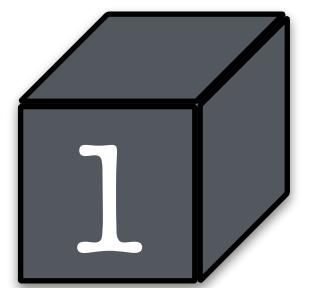
data layout



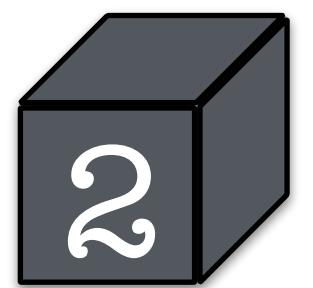
performance



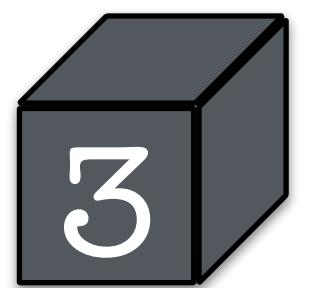
Compaction
black box



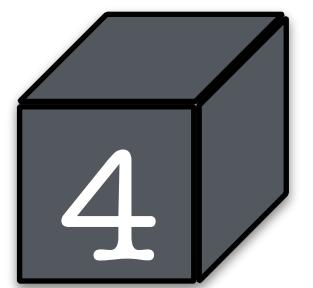
How to organize the data on device?



How much data to move at-a-time?



Which block of data to be moved?



When to re-organize the data layout?

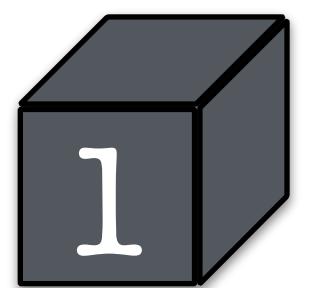


Data Layout

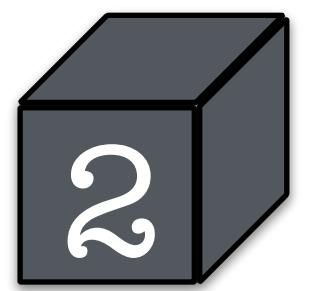
Compaction
Granularity

Data Movement
Policy

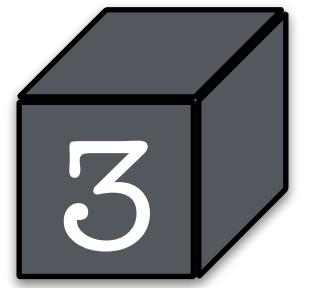
Compaction
Trigger



How to organize the data on device? 



How much data to move at-a-time?

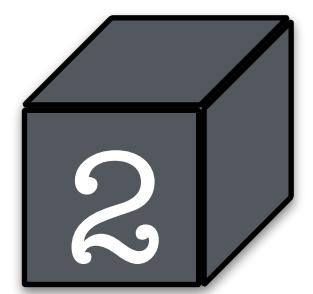


Which block of data to be moved?



When to re-organize the data layout?





Compaction Granularity

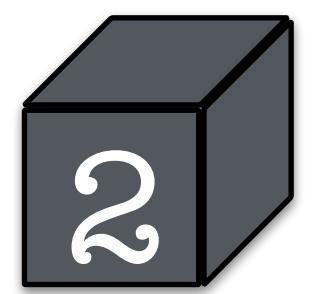
data moved per compaction



consecutive
levels

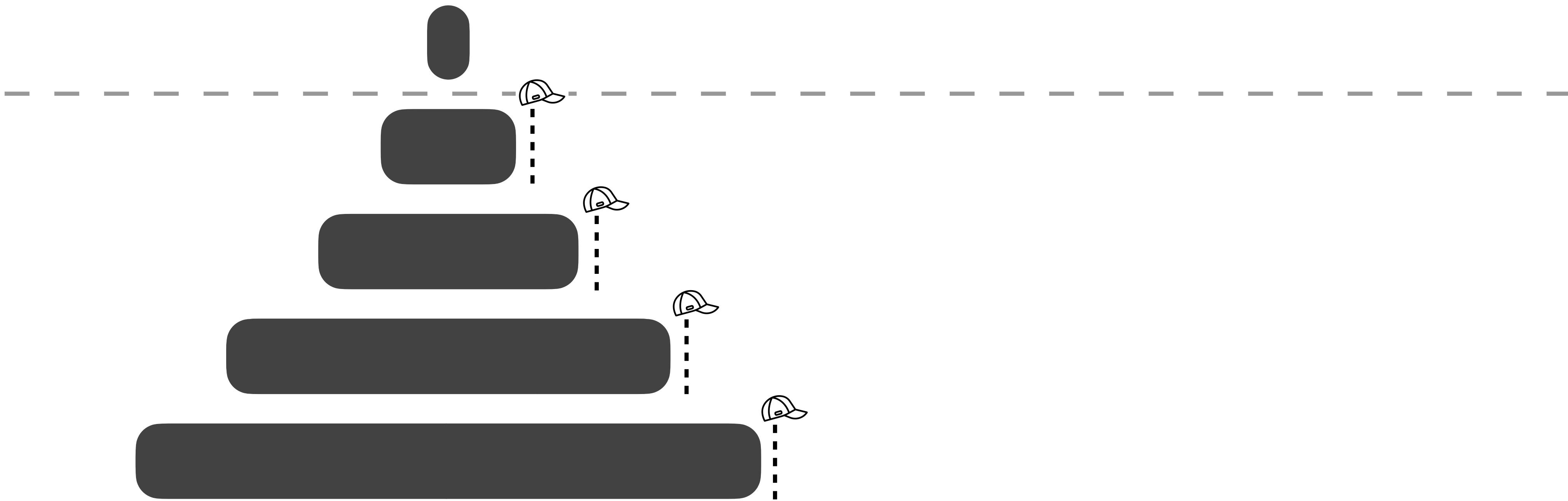
AsterixDB

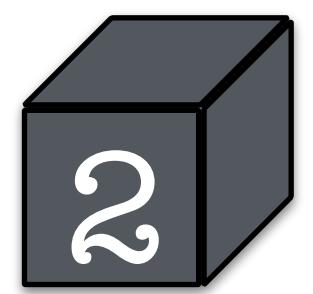




Compaction Granularity

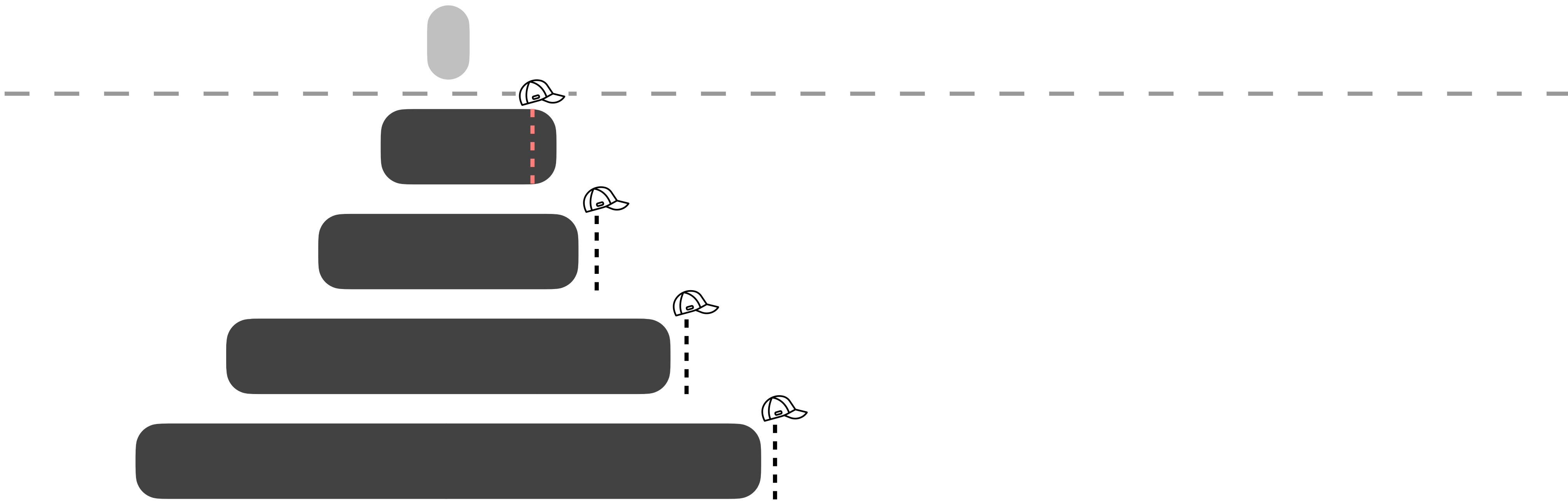
data moved per compaction

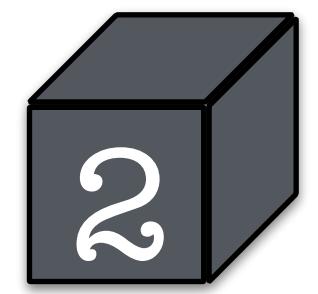




Compaction Granularity

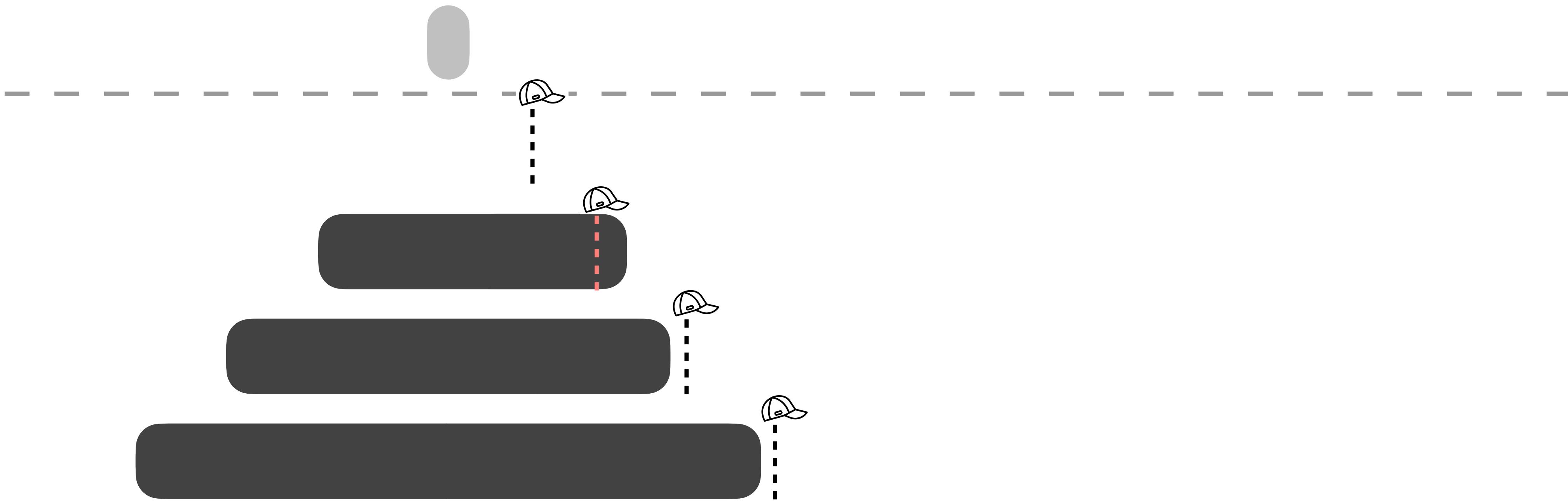
data moved per compaction

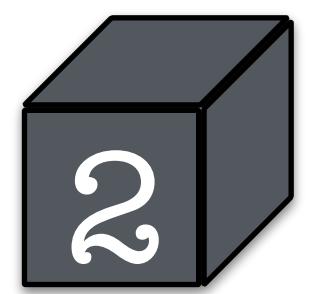




Compaction Granularity

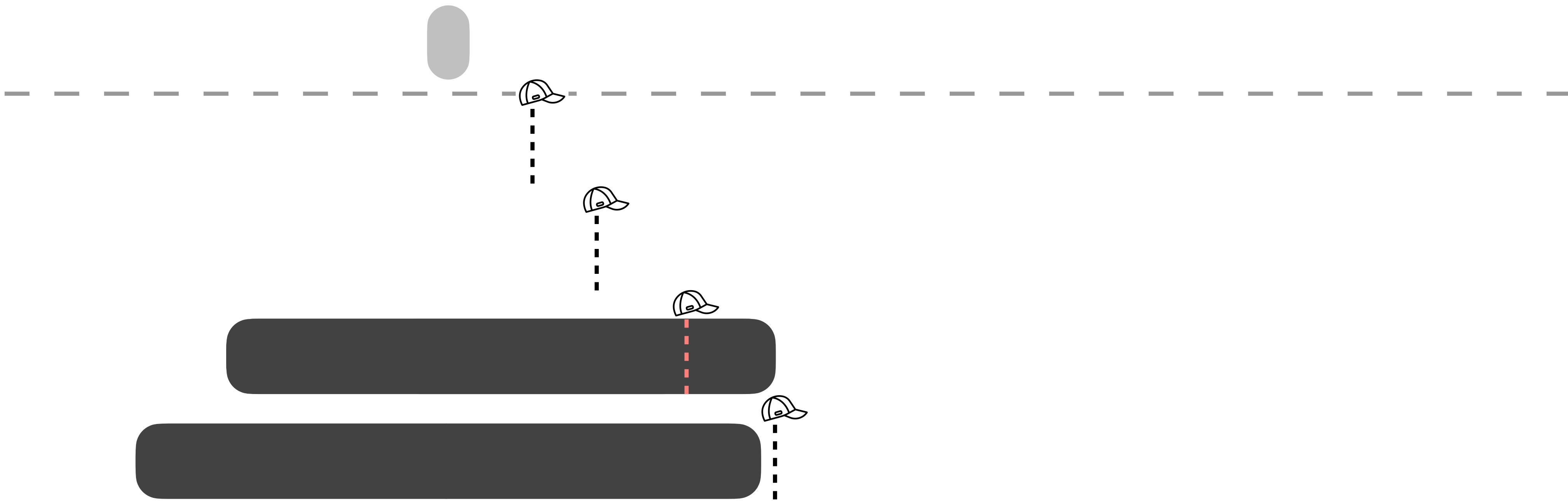
data moved per compaction

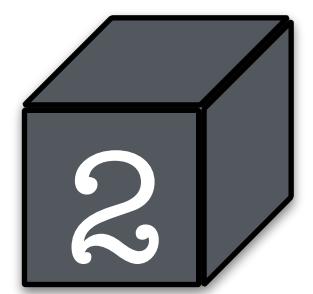




Compaction Granularity

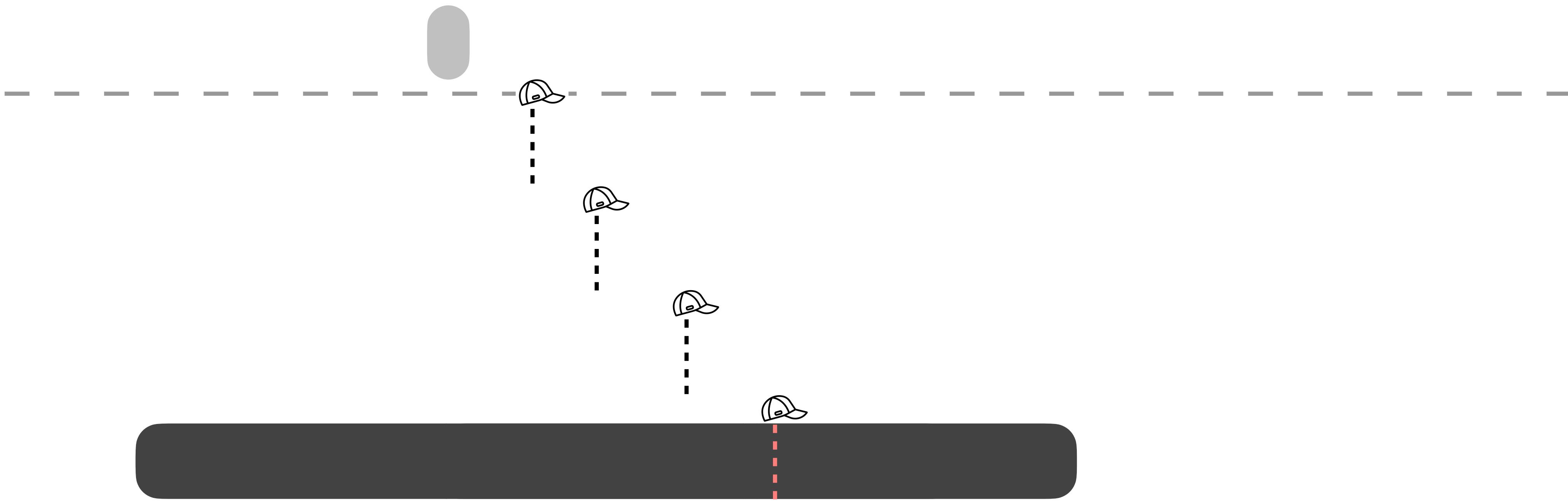
data moved per compaction

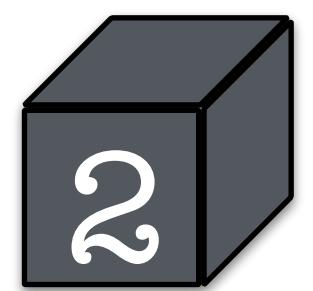




Compaction Granularity

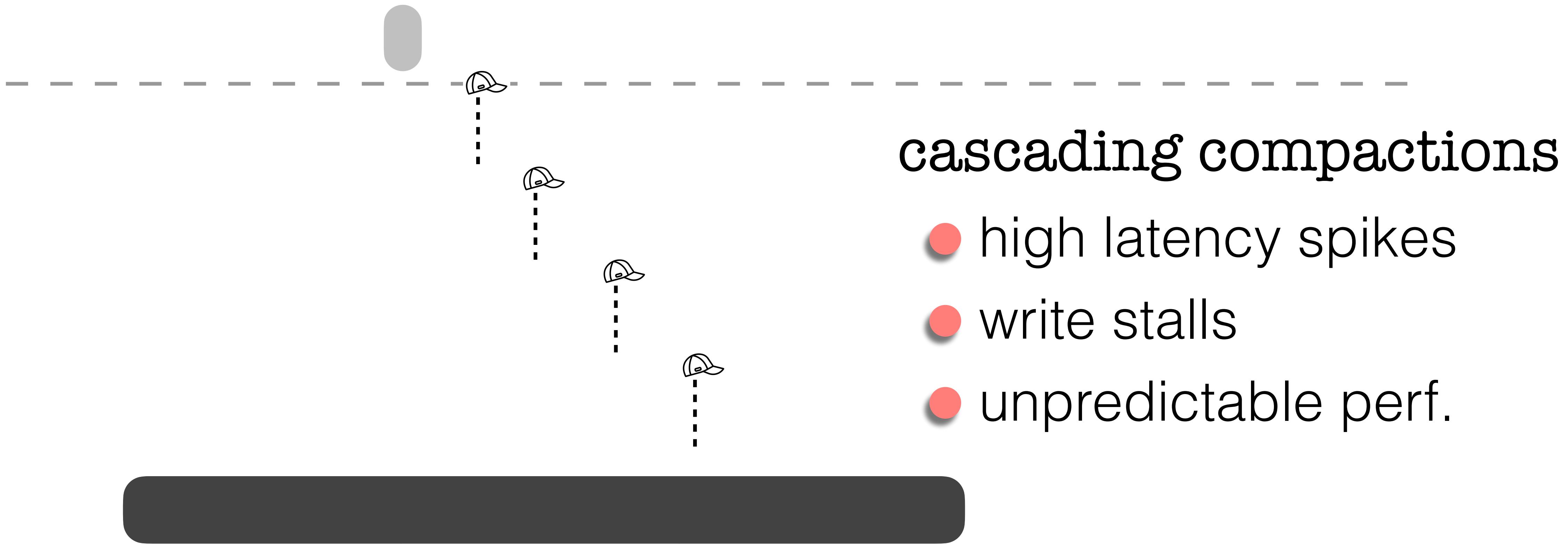
data moved per compaction

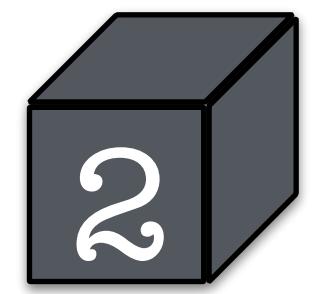




Compaction Granularity

data moved per compaction





Compaction **Granularity**

data moved per compaction

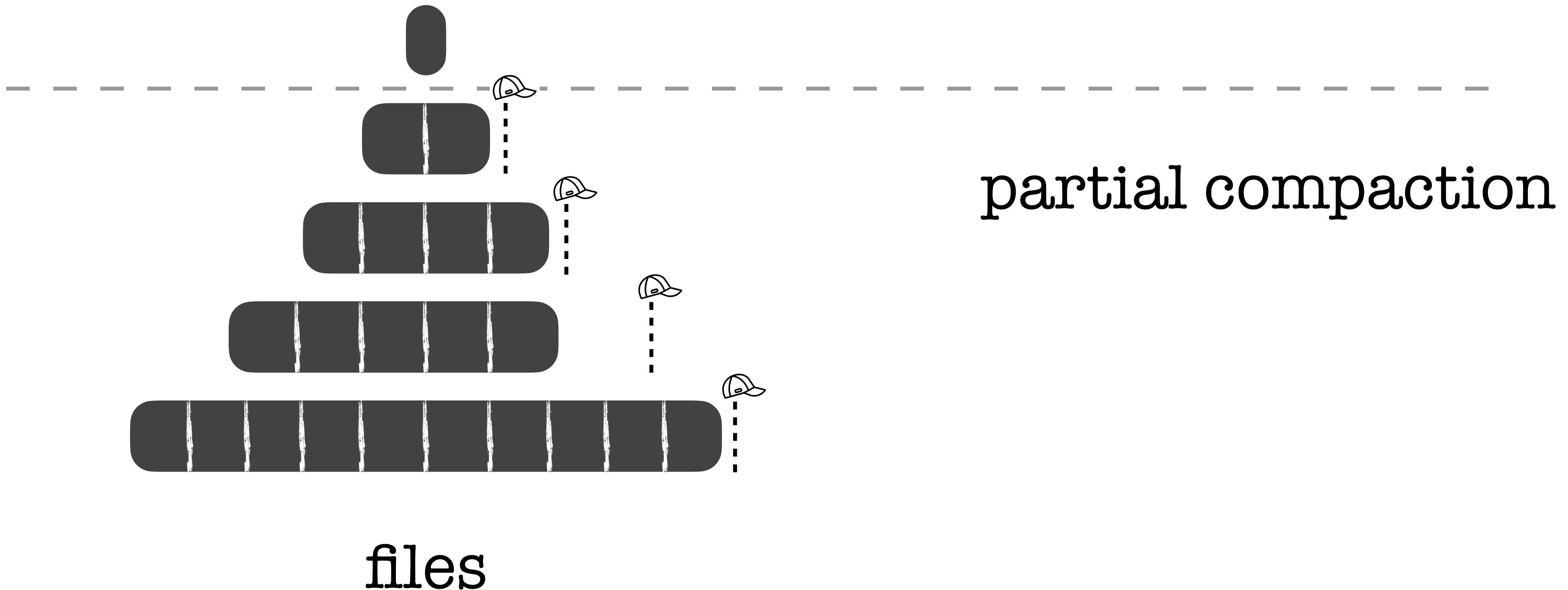
partial compaction

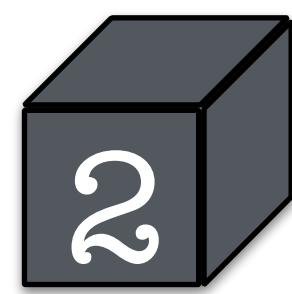
granularity: files

2

Compaction Granularity

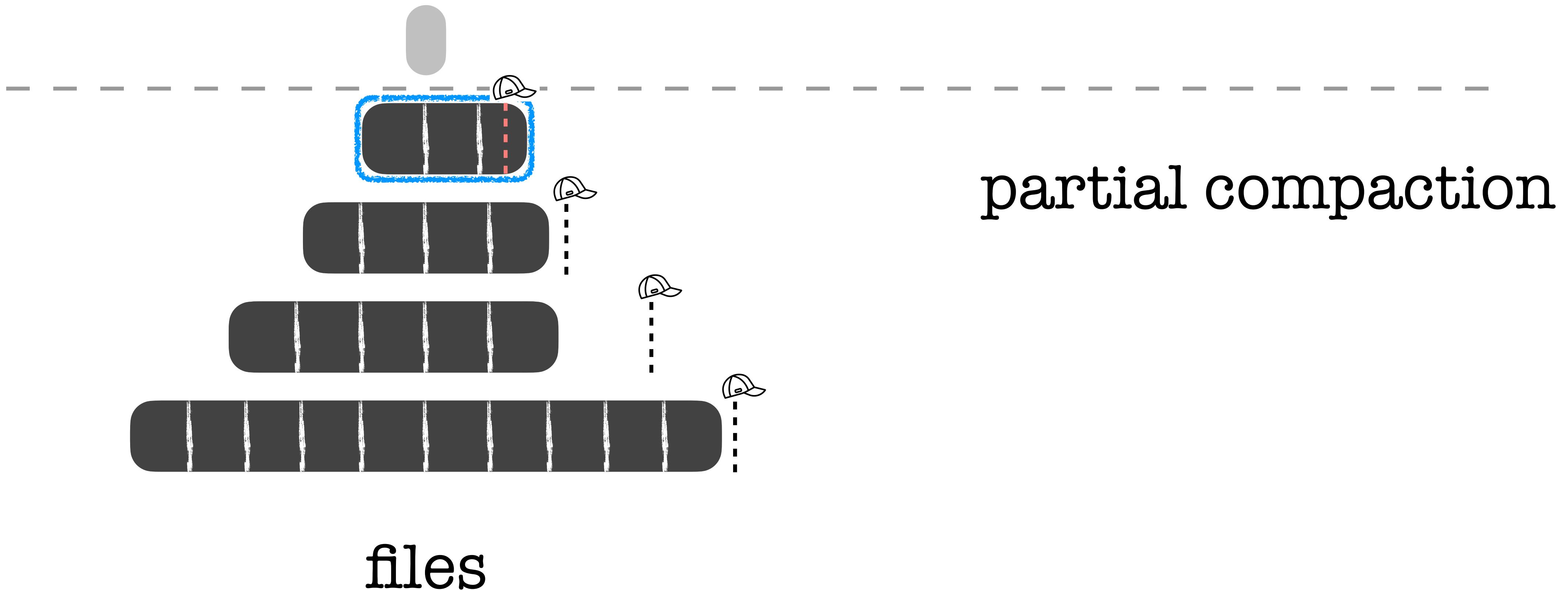
data moved per compaction





Compaction Granularity

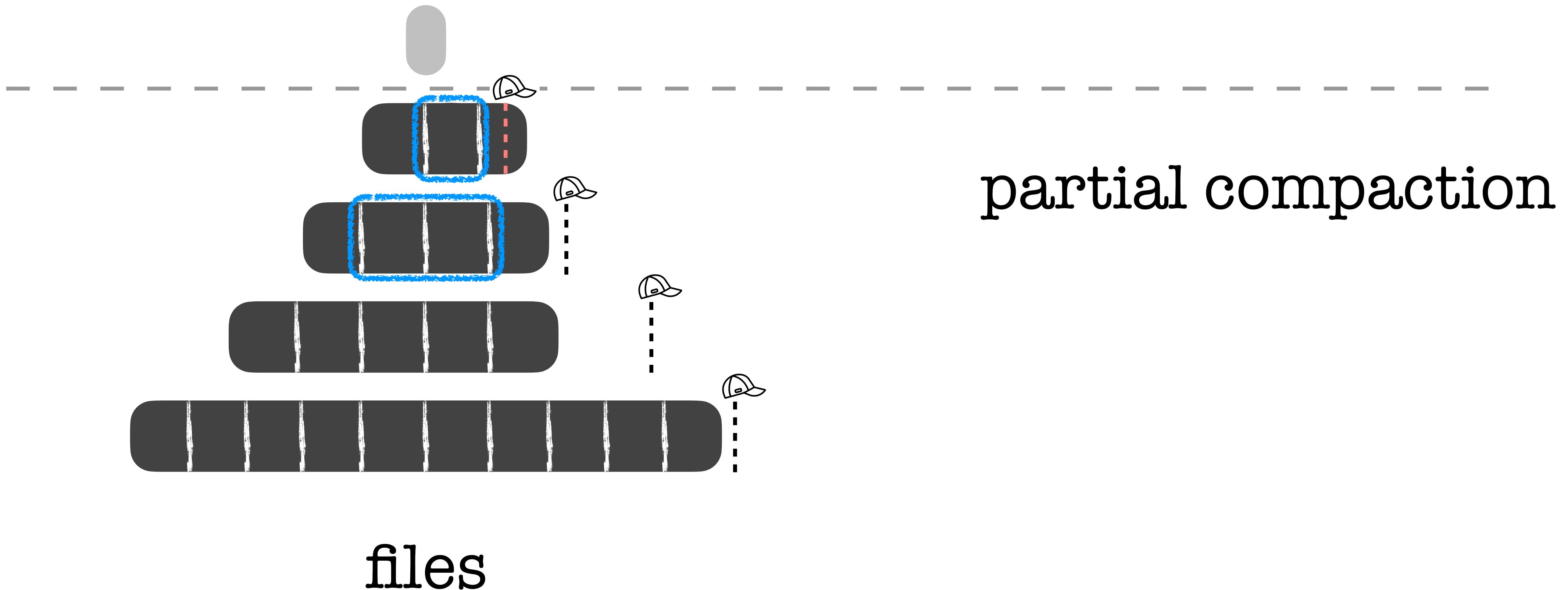
data moved per compaction



2

Compaction Granularity

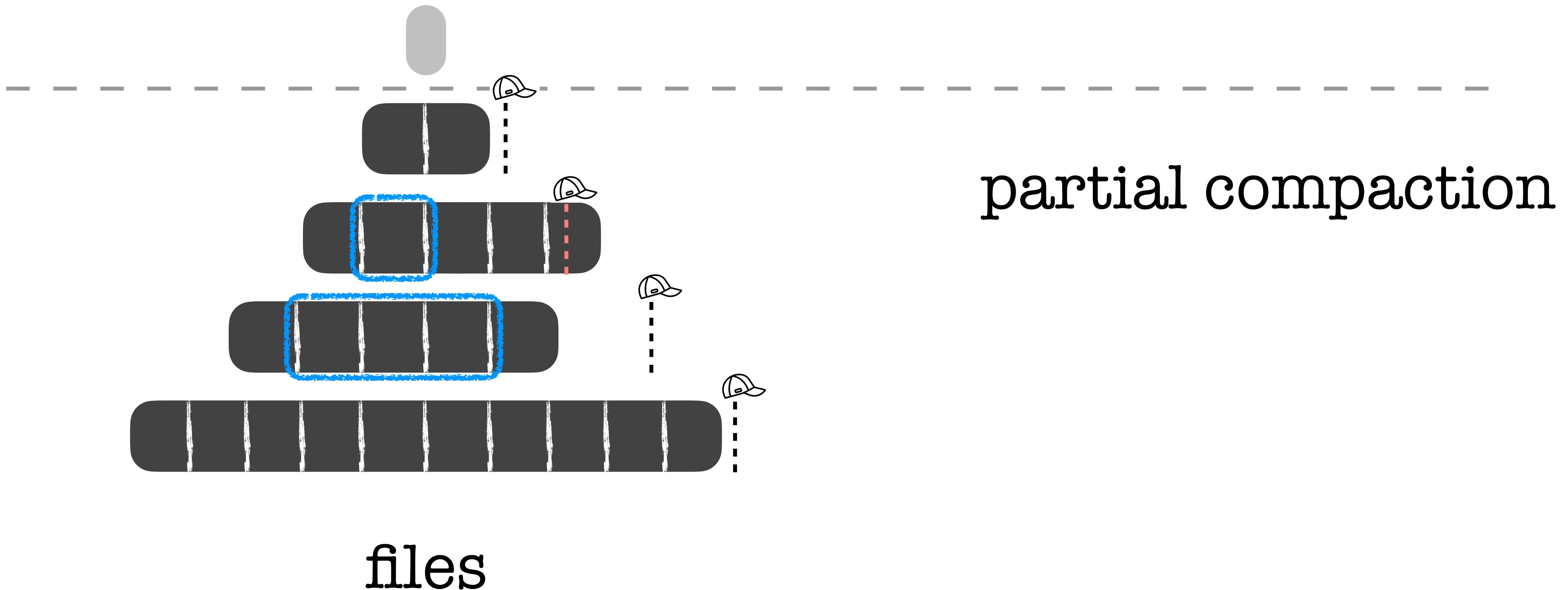
data moved per compaction

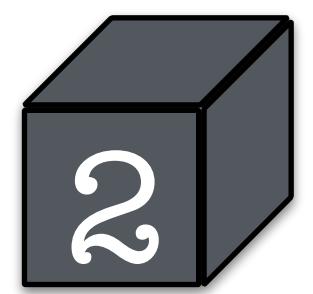


2

Compaction Granularity

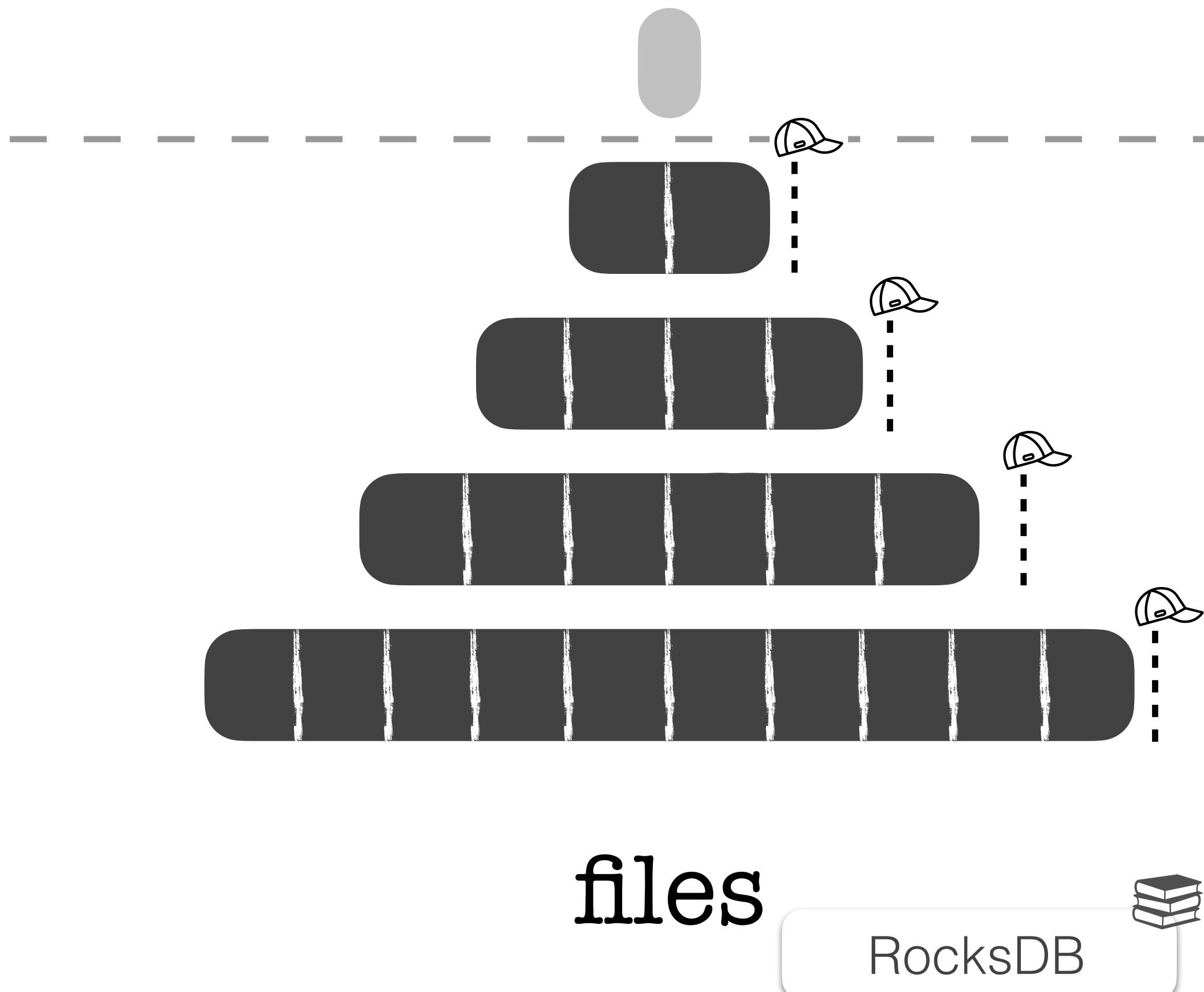
data moved per compaction





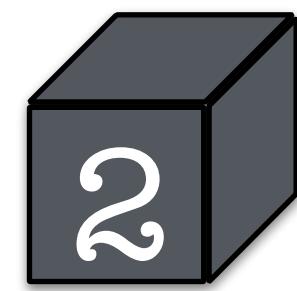
Compaction Granularity

data moved per compaction



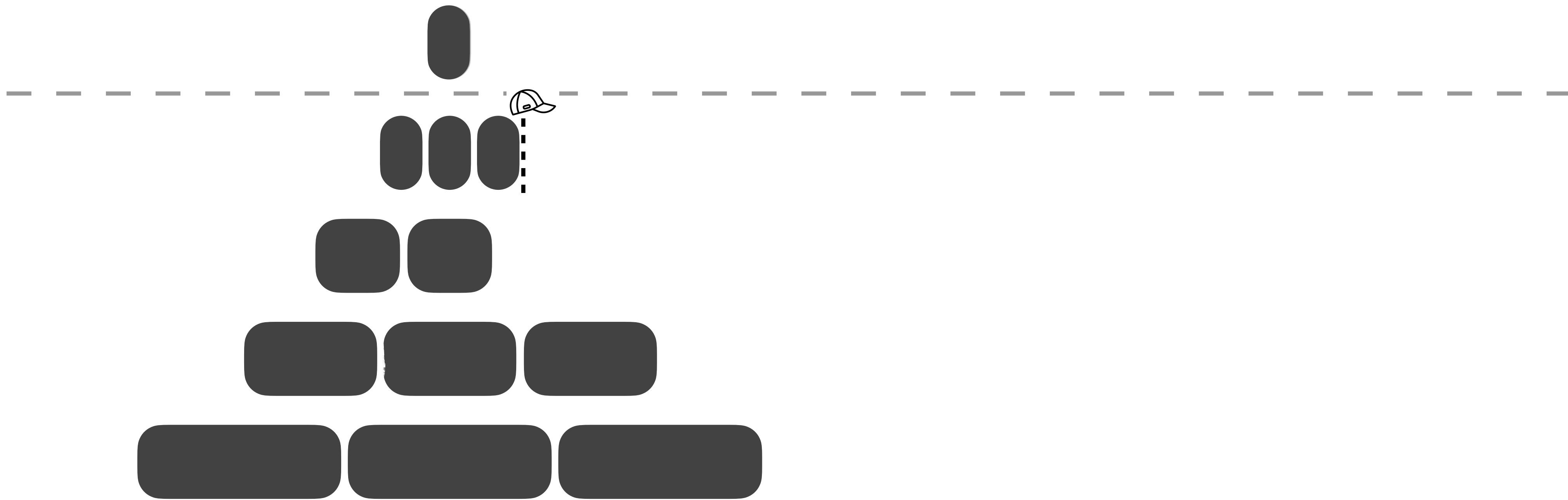
partial compaction

- ~same data movement
- amortized cost for compactions
- predictable perf.

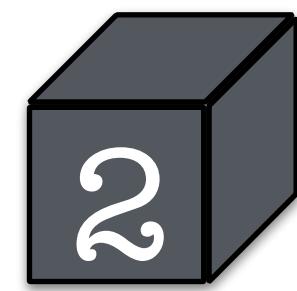


Compaction Granularity

data moved per compaction

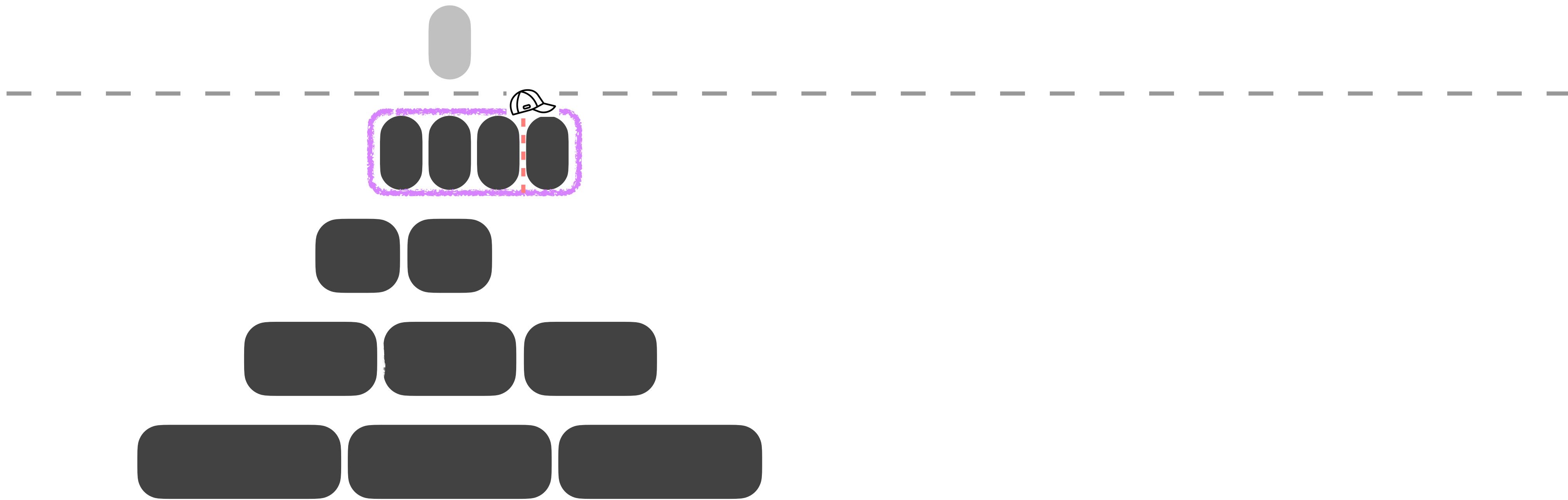


sorted runs in a level

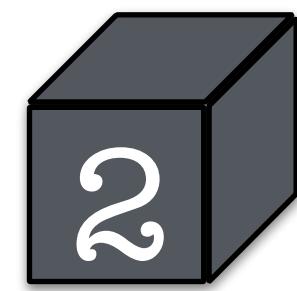


Compaction Granularity

data moved per compaction

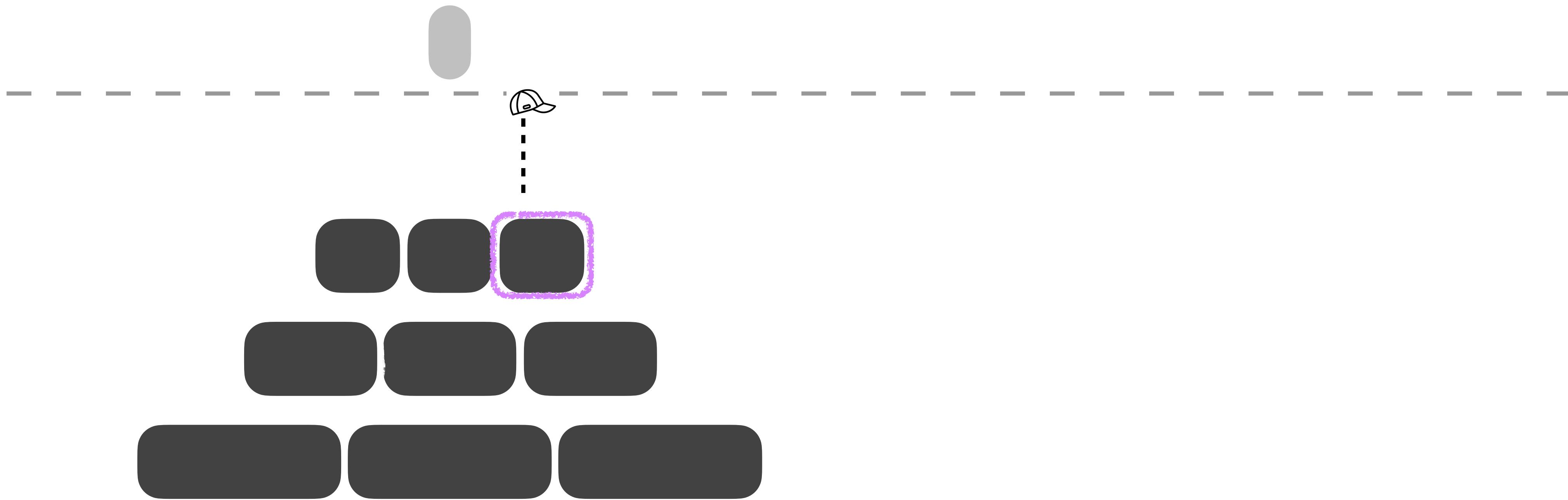


sorted runs in a level

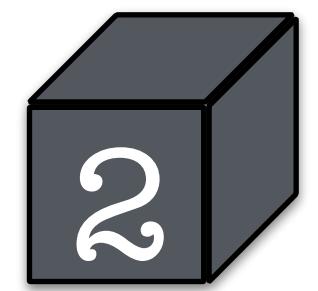


Compaction Granularity

data moved per compaction

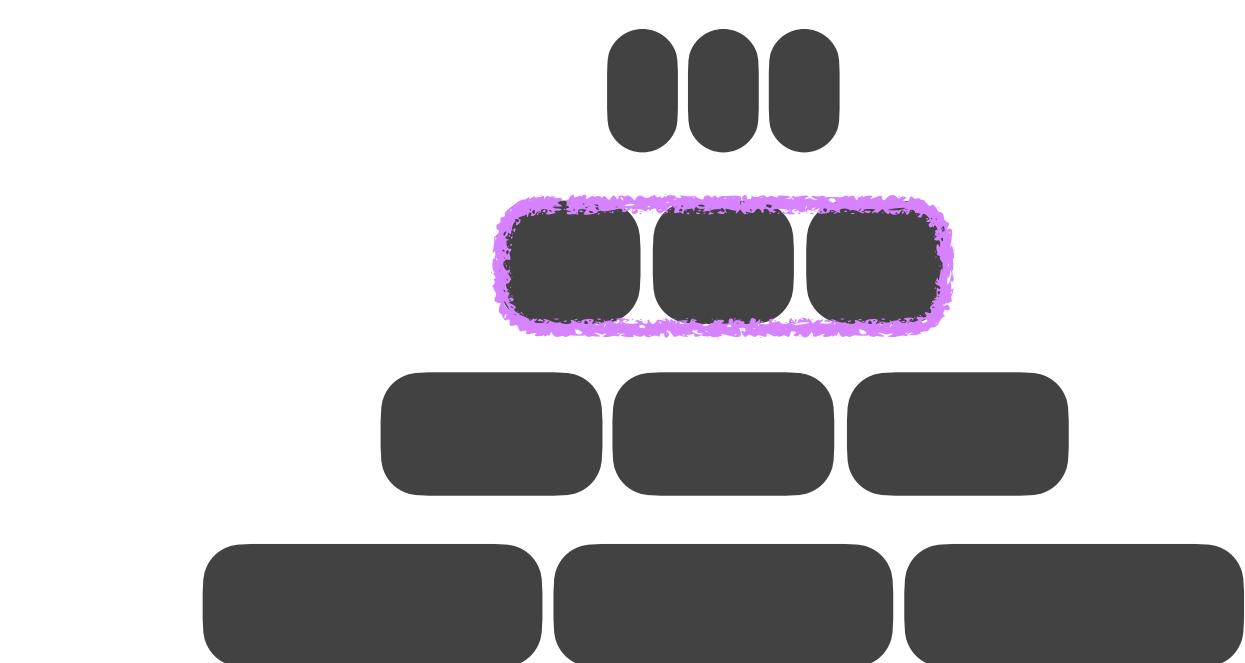
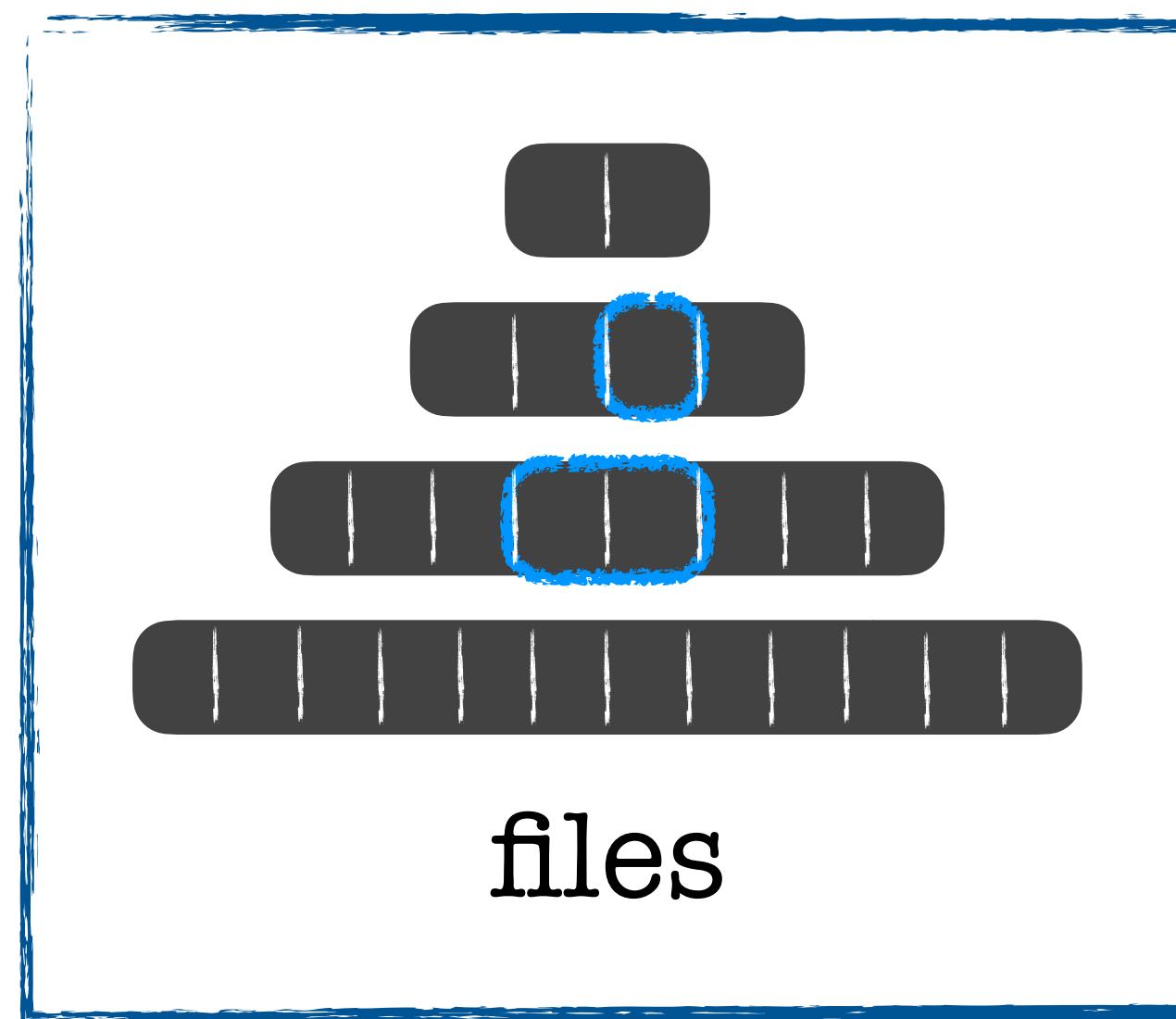
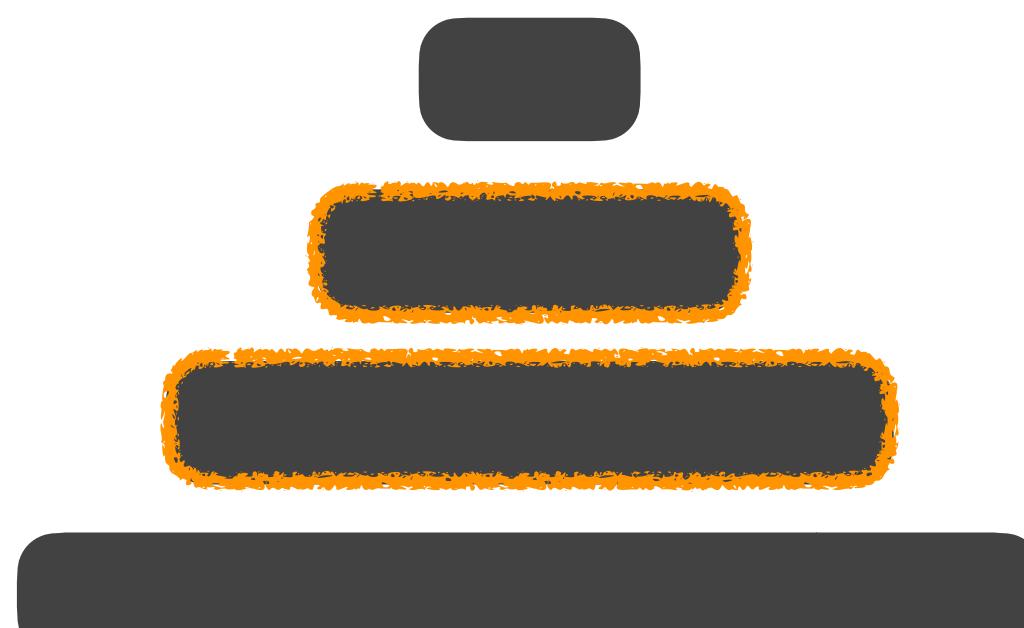


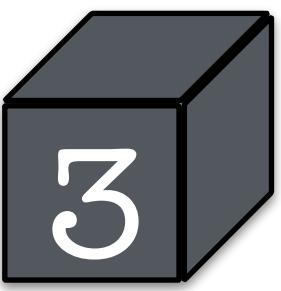
sorted runs in a level



Compaction Granularity

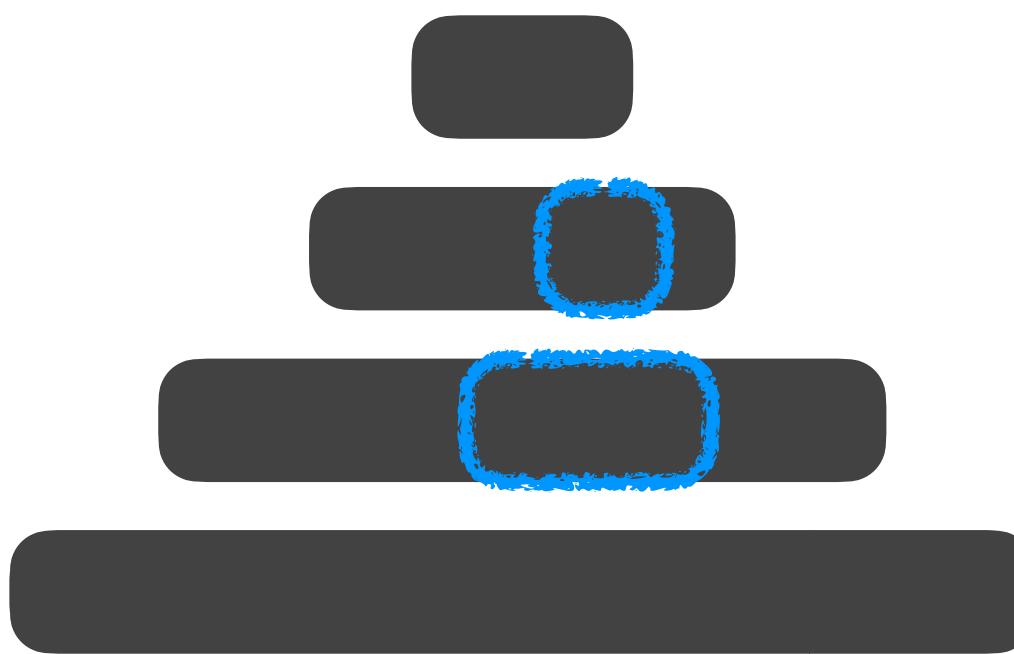
data moved per compaction



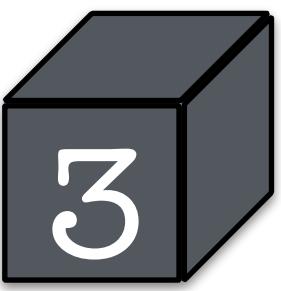


Data Movement Policy

which data to compact

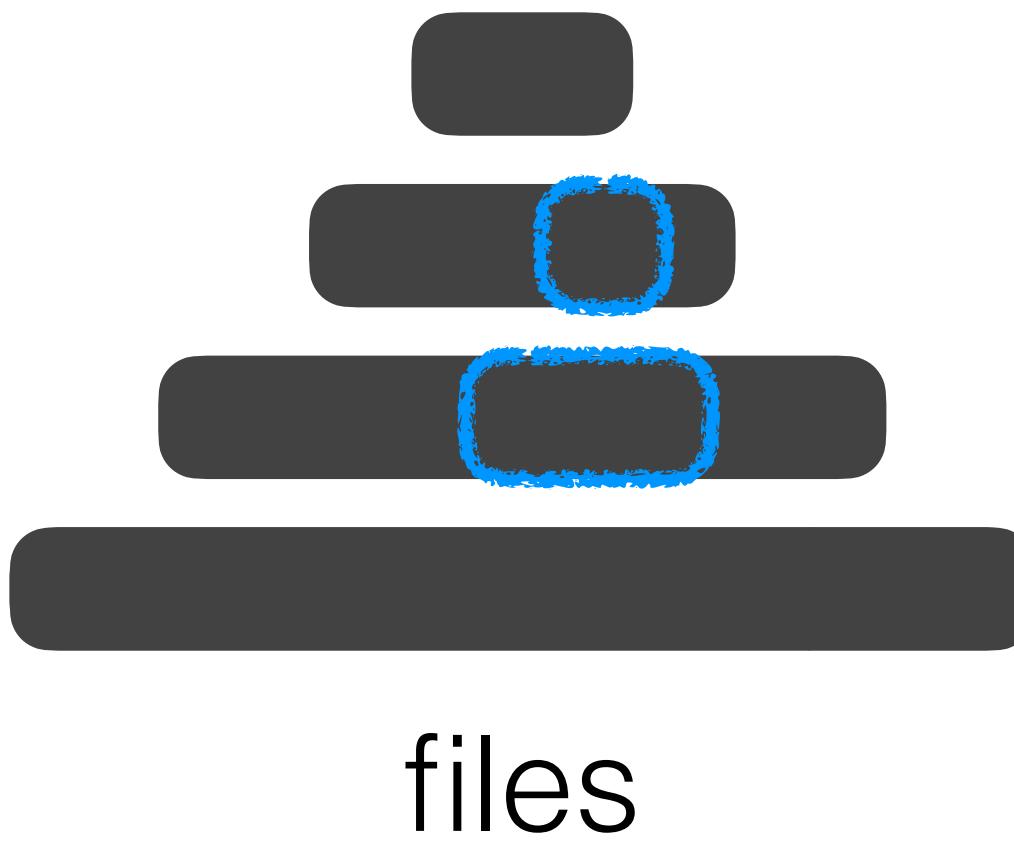


files



Data Movement Policy

which data to compact



round-robin



minimum **overlap with parent level**



file with most **tombstones**



coldest file

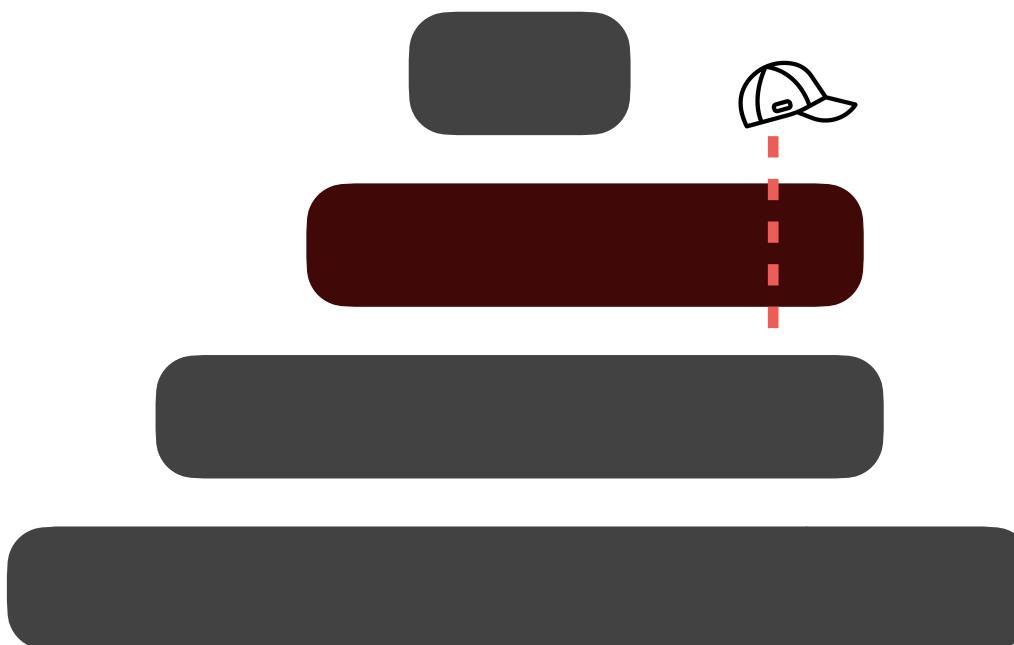




Compaction Trigger

invoking the compaction routine

level **saturation**

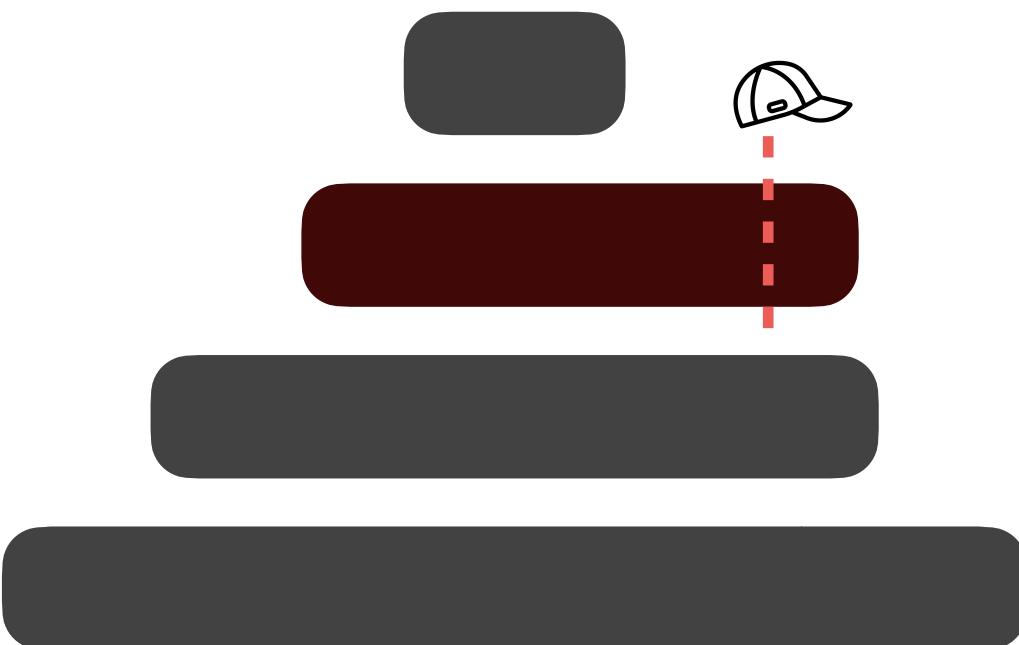




Compaction Trigger

invoking the compaction routine

level **saturation**



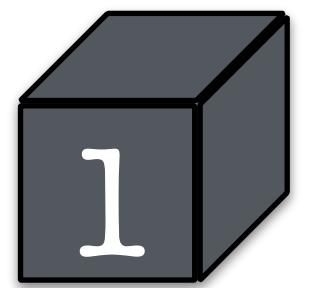
number of **sorted runs**

space amplification

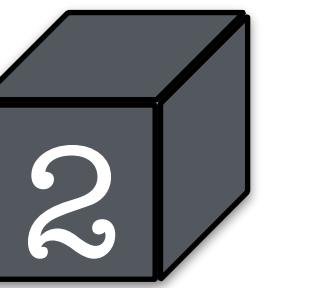
age of a file

SA

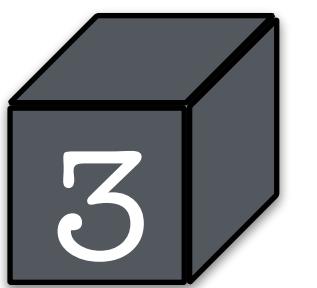
De



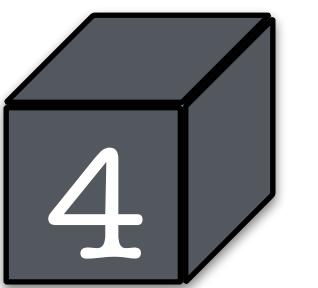
Data Layout



Compaction
Granularity



Data Movement
Policy



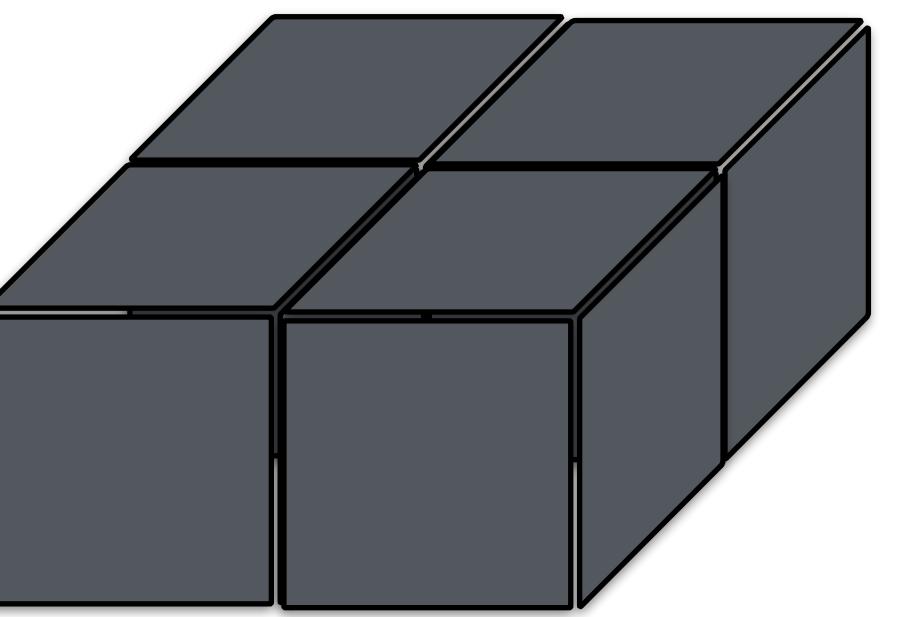
Compaction
Trigger

Data Layout

Compaction
Granularity

Data Movement
Policy

Compaction
Trigger



Any Compaction Algorithm



SarkarVLDB21

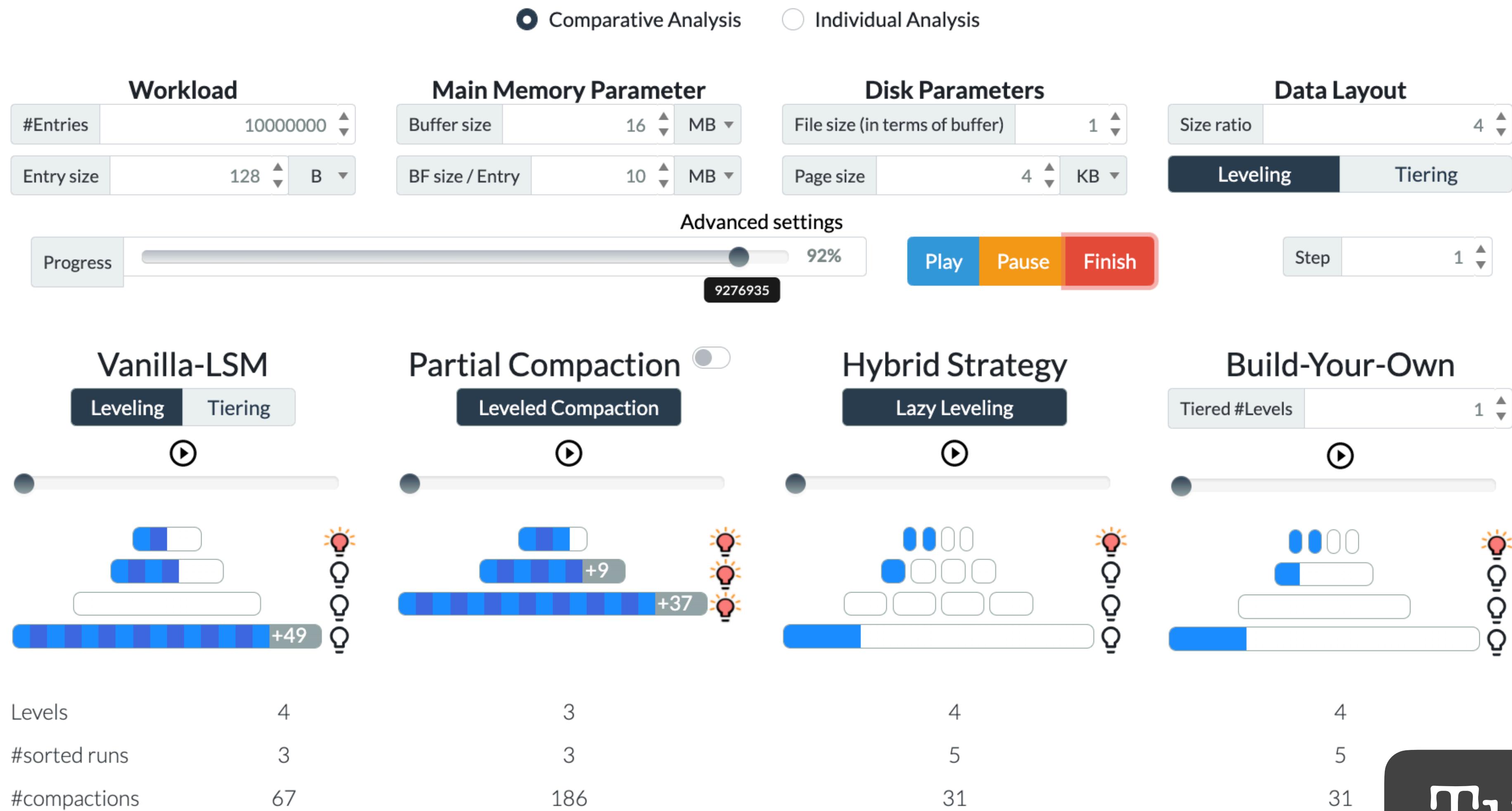
Database	Data layout	Compaction Trigger				Compaction Granularity		Data Movement Policy									
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone density	Expired TS-TTL
RocksDB [30], Monkey [22]	Leveling / 1-Leveling Tiering	✓	✓					✓	✓		✓	✓	✓	✓	✓	✓	✓
LevelDB [32], Monkey (J.) [21]	Leveling	✓						✓			✓	✓	✓				
SlimDB [47]	Tiering	✓						✓	✓							✓	
Dostoevsky [23]	L-leveling	✓ ^L	✓ ^T					✓ ^L	✓ ^T		✓ ^L					✓ ^T	
LSM-Bush [24]	Hybrid leveling	✓ ^L	✓ ^T					✓ ^L	✓ ^T		✓ ^L					✓ ^T	
Lethe [51]	Leveling	✓		✓				✓	✓		✓					✓	
Silk [11], Silk+ [12]	Leveling	✓						✓	✓		✓						
HyperLevelDB [35]	Leveling	✓						✓			✓	✓	✓				
PebblesDB [46]	Hybrid leveling	✓						✓	✓							✓	
Cassandra [8]	Tiering		✓	✓	✓	✓		✓								✓	✓
	Leveling	✓		✓				✓	✓		✓				✓	✓	✓
WiredTiger [62]	Leveling	✓					✓									✓	
X-Engine [34], Leaper [63]	Hybrid leveling	✓						✓	✓		✓				✓	✓	
HBase [7]	Tiering		✓					✓								✓	
AsterixDB [3]	Leveling	✓					✓									✓	
	Tiering	✓					✓									✓	

Database	Data layout	Compaction Trigger				Compaction Granularity		Data Movement Policy								
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone density
RocksDB [30], Monkey [22]	Leveling / 1-Leveling	✓	✓					✓	✓		✓	✓	✓	✓	✓	
	Tiering		✓	✓	✓		✓								✓	
LevelDB [32], Monkey (J.) [21]	Leveling	✓						✓			✓	✓	✓			
SlimDB [47]	Tiering	✓						✓	✓						✓	
Dostoevsky [23]	L-leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
LSM-Bush [24]	Hybrid leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
Lethe [51]	Leveling	✓		✓				✓	✓		✓				✓	
Silk [11], Silk+ [12]	Leveling	✓						✓	✓		✓					
HyperLevelDB [35]	Leveling	✓						✓			✓	✓	✓			
PebblesDB [46]	Hybrid leveling	✓						✓	✓						✓	
Cassandra [8]	Tiering		✓	✓	✓	✓		✓							✓	
	Leveling	✓		✓				✓	✓		✓			✓	✓	
WiredTiger [62]	Leveling	✓				✓									✓	
X-Engine [34], Leaper [63]	Hybrid leveling	✓						✓	✓		✓				✓	
HBase [7]	Tiering		✓			✓									✓	
AsterixDB [3]	Leveling	✓			✓										✓	
	Tiering		✓			✓									✓	

Database	Data layout	Compaction Trigger				Compaction Granularity		Data Movement Policy								
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone density
RocksDB [30], Monkey [22]	Leveling / 1-Leveling	✓	✓					✓	✓		✓	✓	✓	✓	✓	
	Tiering		✓	✓	✓		✓								✓	
LevelDB [32], Monkey (J.) [21]	Leveling	✓						✓			✓	✓	✓			
SlimDB [47]	Tiering	✓						✓	✓						✓	
Dostoevsky [23]	L-leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
LSM-Bush [24]	Hybrid leveling	✓ ^L	✓ ^T				✓ ^L	✓ ^T			✓ ^L				✓ ^T	
Lethe [51]	Leveling	✓		✓				✓	✓		✓				✓	
Silk [11], Silk+ [12]	Leveling	✓						✓	✓		✓					
HyperLevelDB [35]	Leveling	✓						✓			✓	✓	✓			
PebblesDB [46]	Hybrid leveling	✓						✓	✓						✓	
Cassandra [8]	Tiering		✓	✓	✓		✓								✓	
	Leveling	✓		✓				✓	✓		✓				✓	✓
WiredTiger [62]	Leveling	✓				✓										✓
X-Engine [34], Leaper [63]	Hybrid leveling	✓						✓	✓		✓				✓	
HBase [7]	Tiering		✓				✓								✓	
AsterixDB [3]	Leveling	✓				✓									✓	
	Tiering	✓				✓									✓	

Database	Data layout	Compaction Trigger				Compaction Granularity		Data Movement Policy								
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone density
RocksDB [30], Monkey [22]	Leveling / 1-Leveling	✓	✓					✓	✓		✓	✓	✓	✓	✓	
	Tiering		✓	✓	✓		✓								✓	
LevelDB [32], Monkey (J.) [21]	Leveling	✓						✓			✓	✓	✓			
SlimDB [47]	Tiering	✓						✓	✓						✓	
Dostoevsky [23]	<i>L</i> -leveling	✓ ^L	✓ ^T			✓ ^L	✓ ^T			✓ ^L					✓ ^T	
LSM-Bush [24]	Hybrid leveling	✓ ^L	✓ ^T			✓ ^L	✓ ^T			✓ ^L					✓ ^T	
Lethe [51]	Leveling	✓		✓				✓	✓		✓					✓
Silk [11], Silk+ [12]	Leveling	✓						✓	✓		✓					
HyperLevelDB [35]	Leveling	✓						✓			✓	✓	✓			
PebblesDB [46]	Hybrid leveling	✓						✓	✓						✓	
Cassandra [8]	Tiering		✓	✓	✓	✓		✓							✓	✓
	Leveling	✓		✓				✓	✓		✓				✓	✓
WiredTiger [62]	Leveling	✓				✓										✓
X-Engine [34], Leaper [63]	Hybrid leveling	✓						✓	✓		✓				✓	
HBase [7]	Tiering		✓				✓									✓
AsterixDB [3]	Leveling	✓				✓										✓
	Tiering	✓				✓										✓

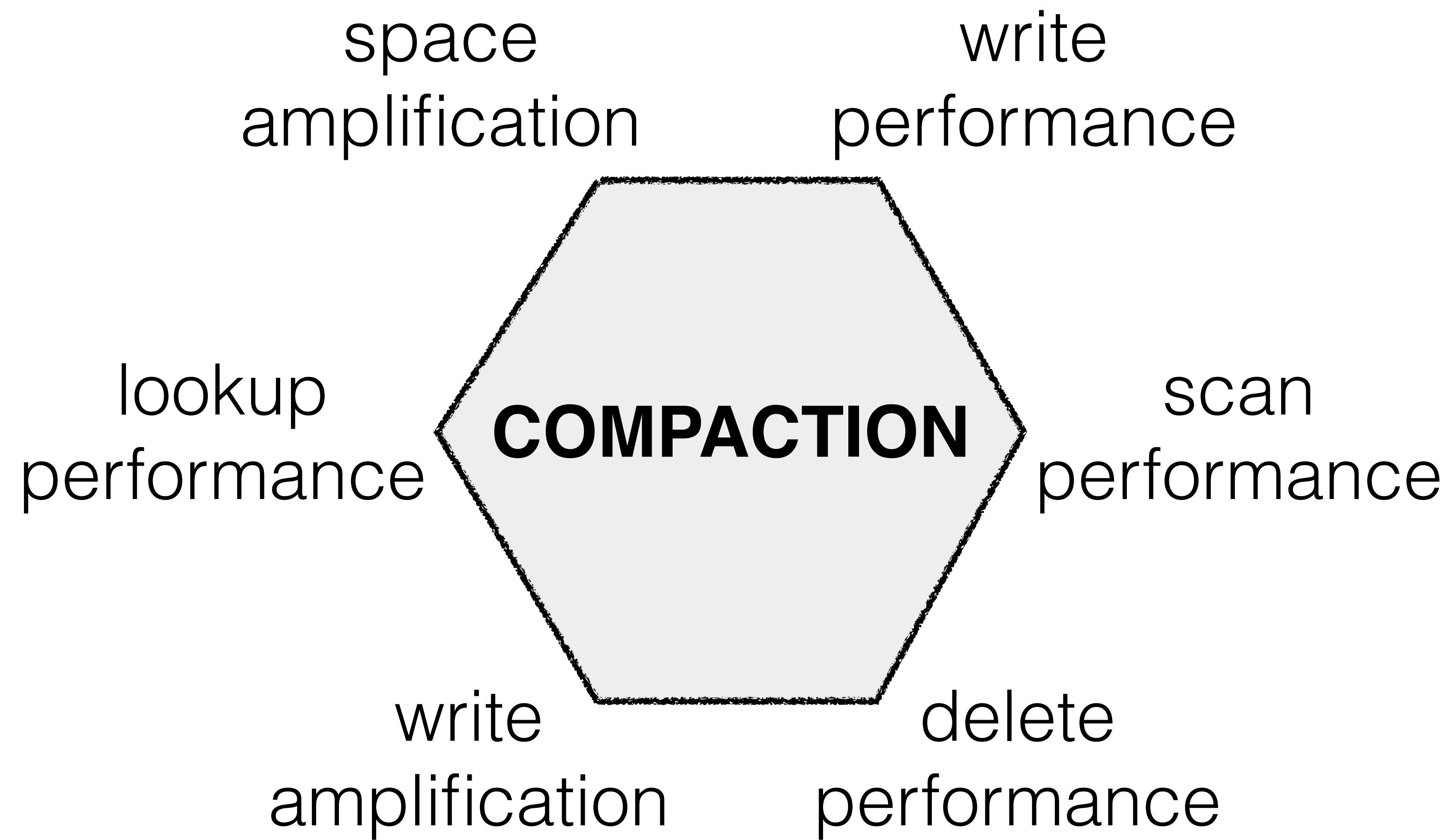
SIGMOD Demo: Compactionary



Tuesday
4-6PM

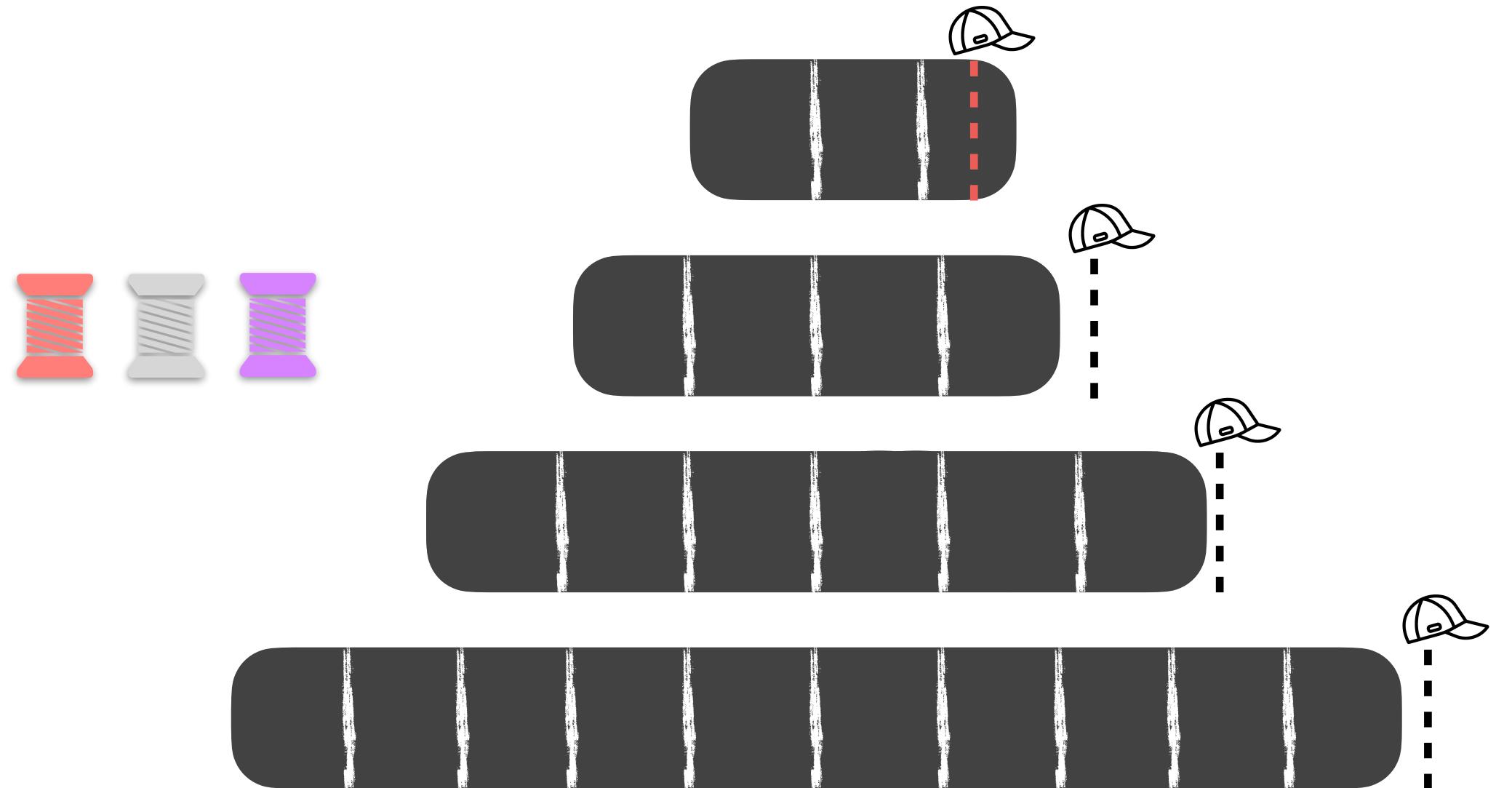
BOSTON UNIVERSITY

SarkarSIGMOD22



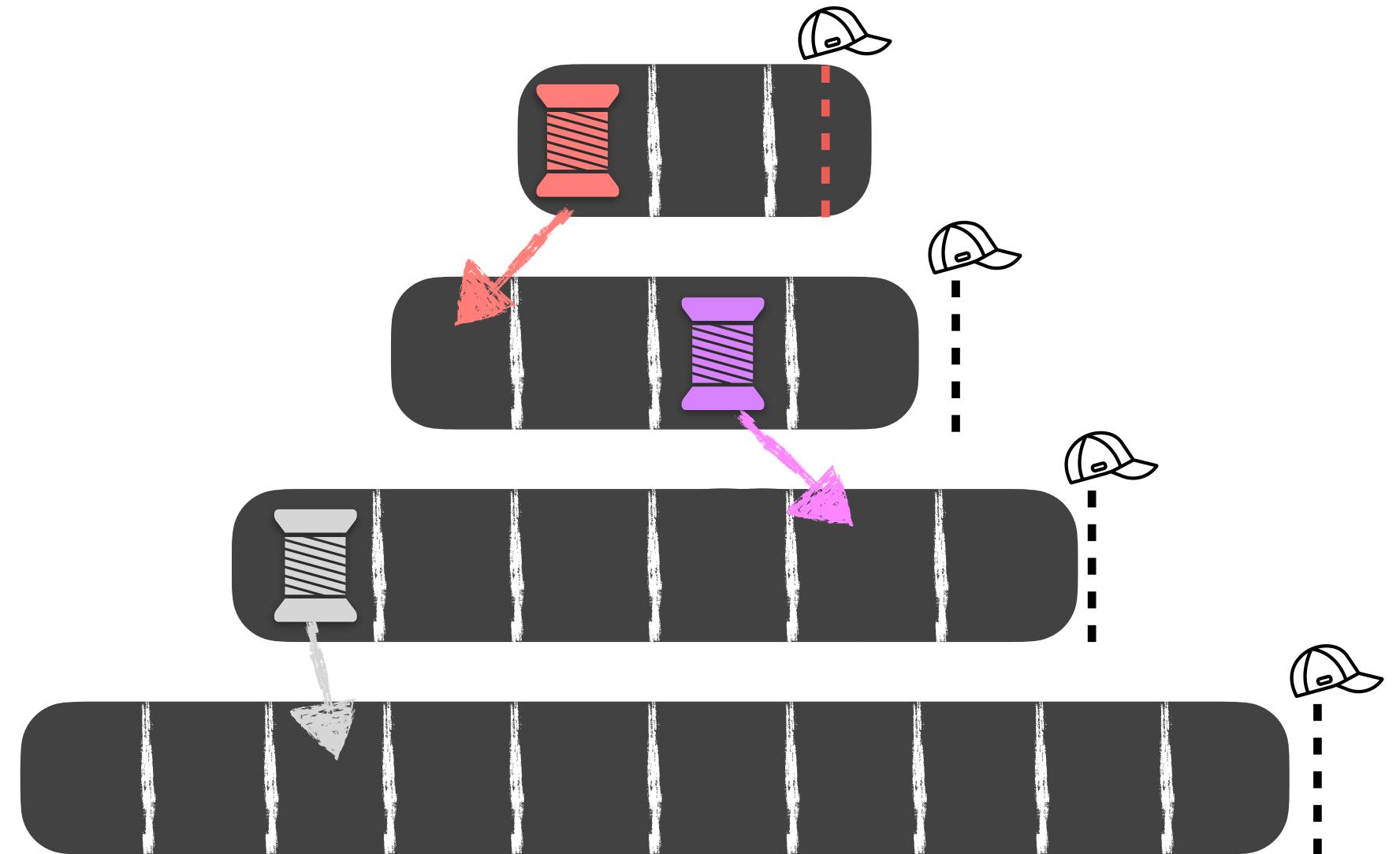
Optimizing Compactions

Background
Compactions



Optimizing Compactions

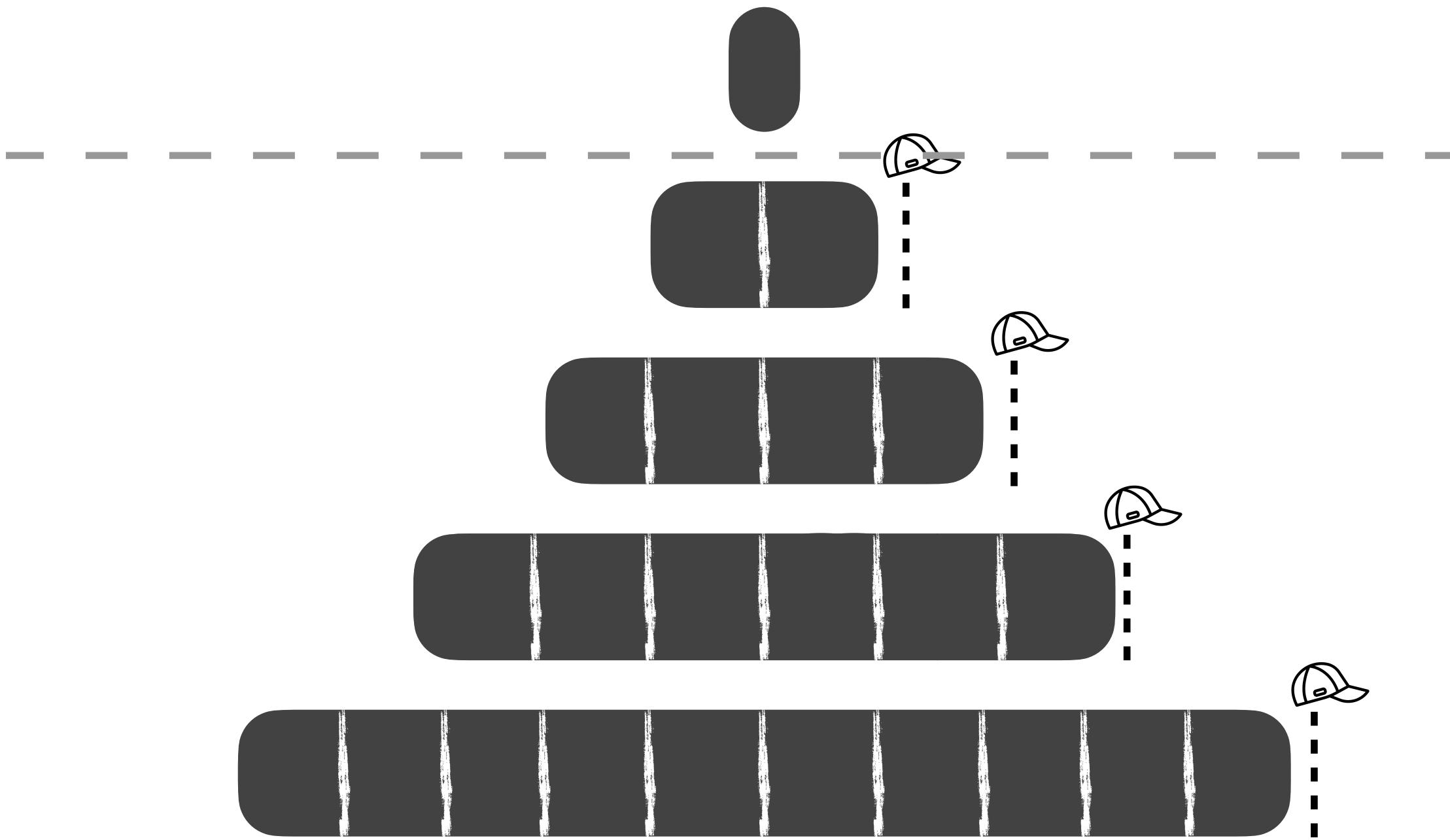
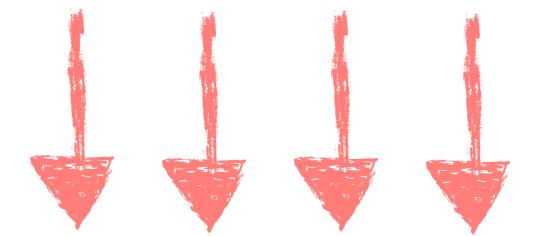
Background Compactions



- non-blocking reads/writes
- improves write throughput

Optimizing Compactions

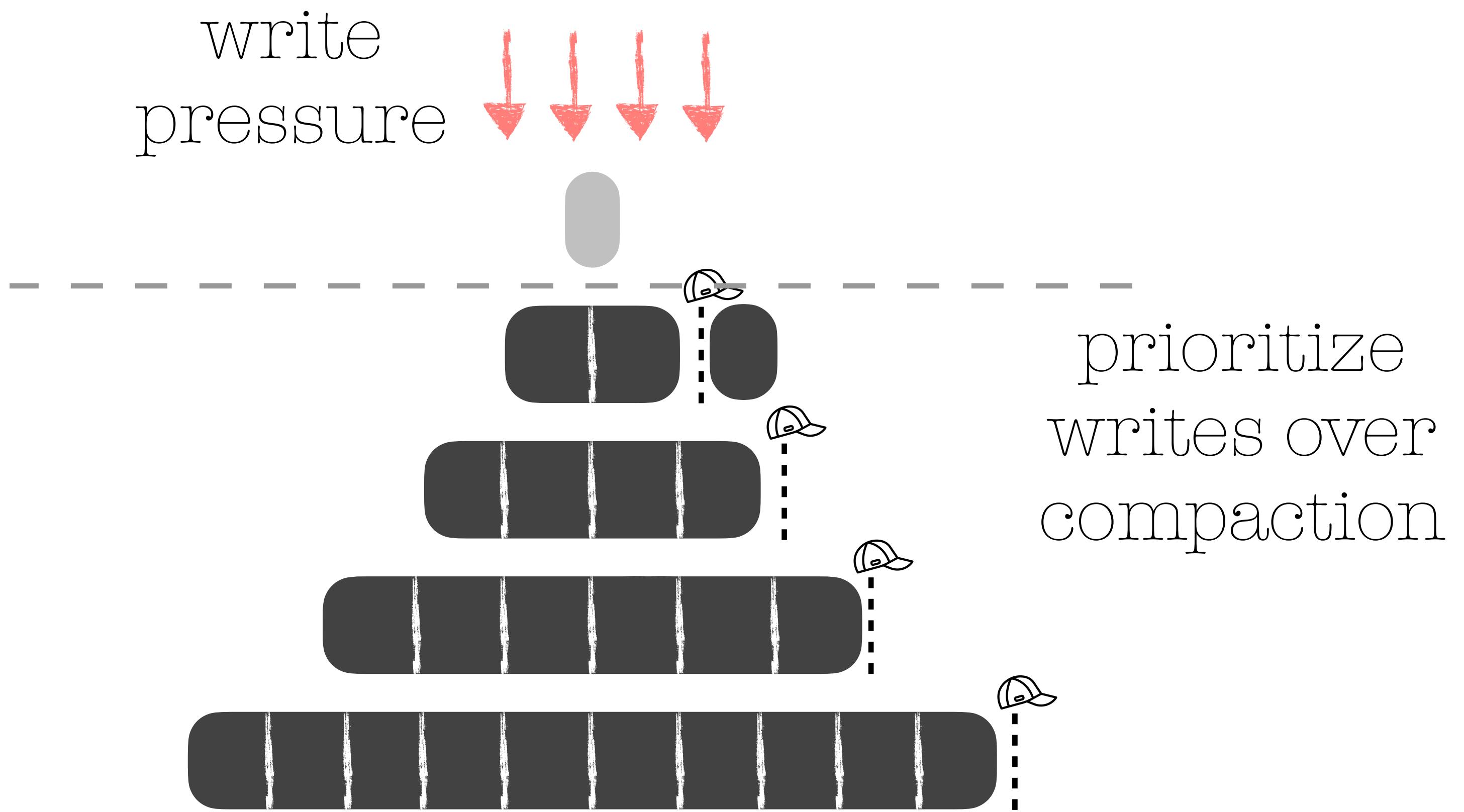
write
pressure



Background
Compactions

Compaction
Priority

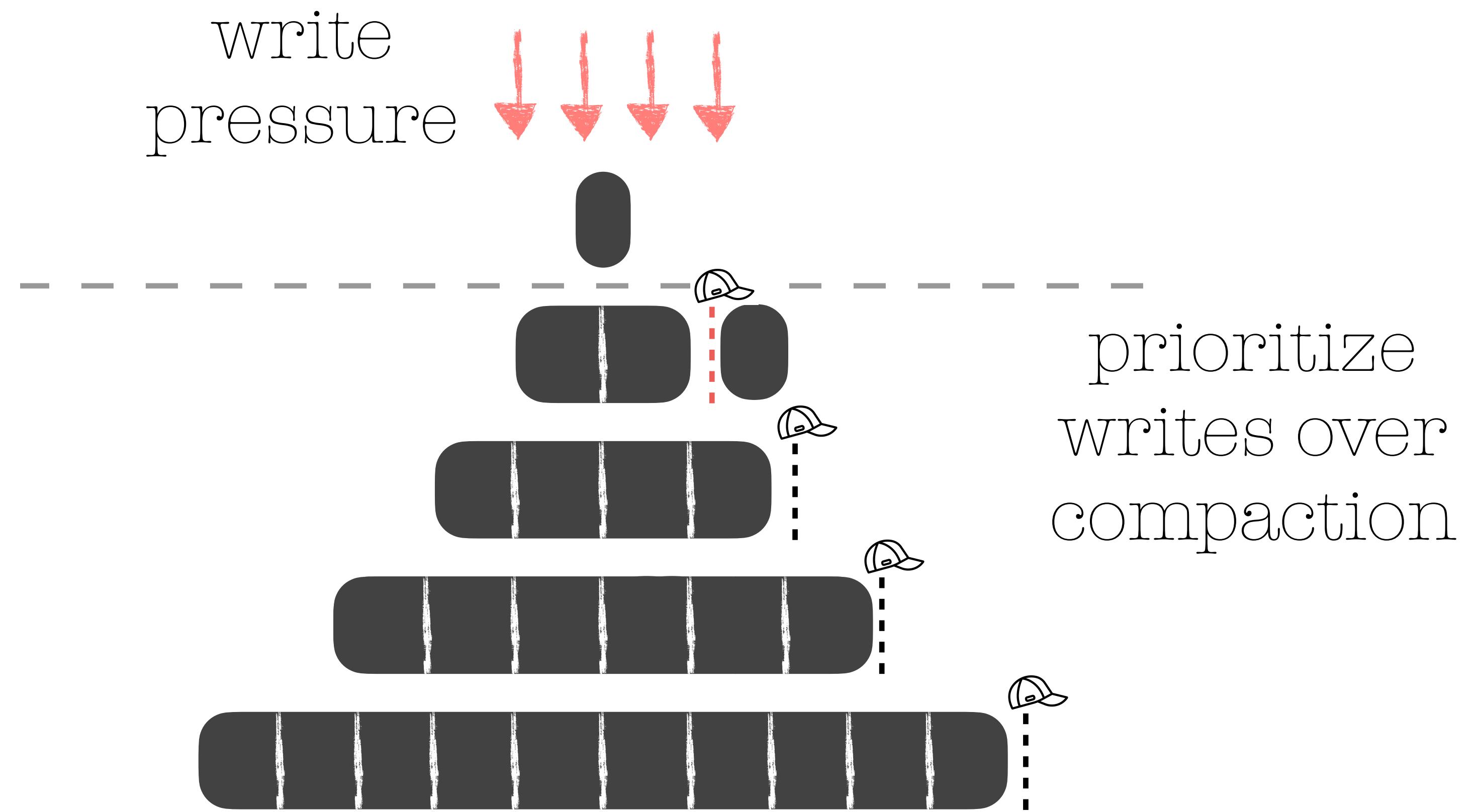
Optimizing Compactions



Background
Compactions

Compaction
Priority

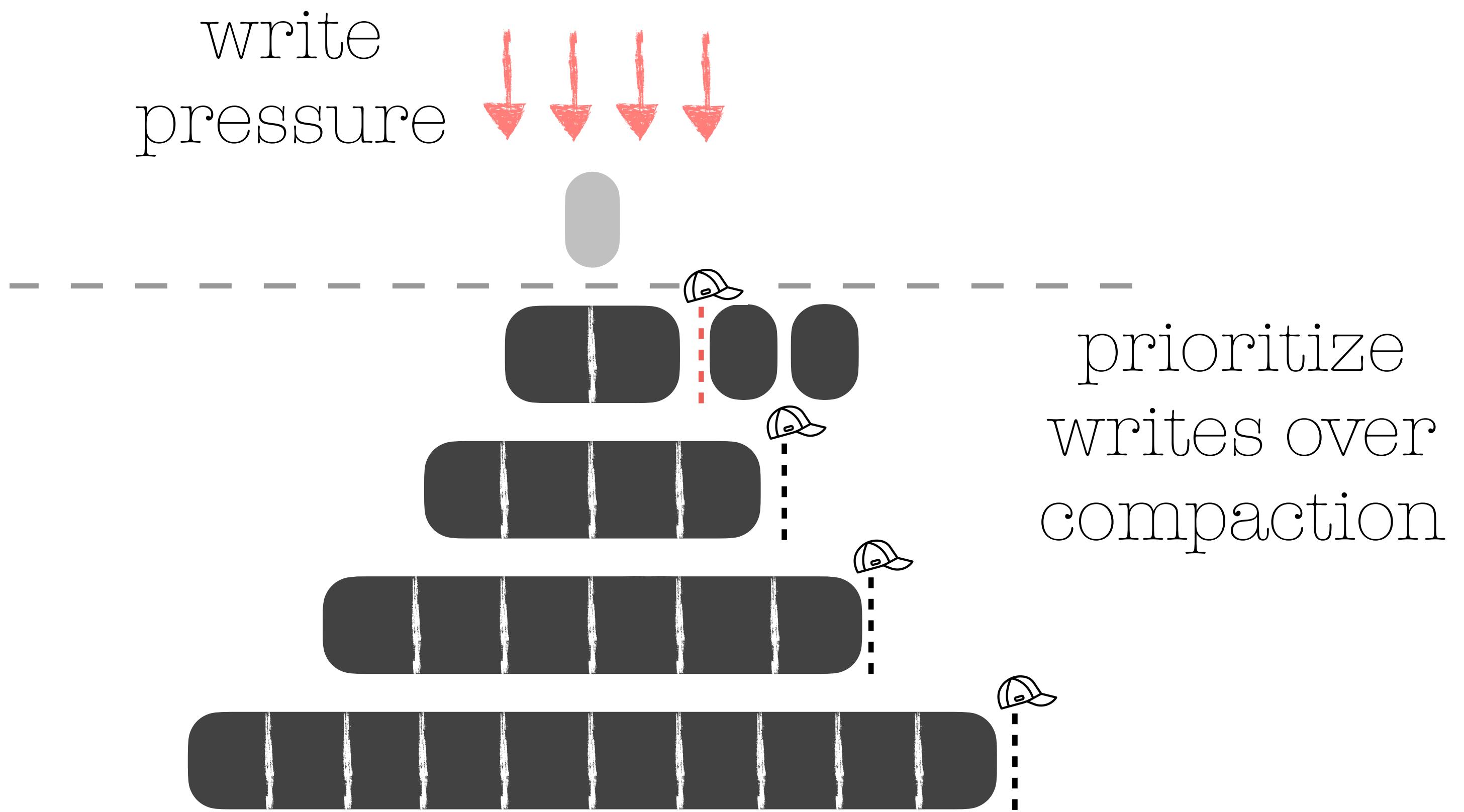
Optimizing Compactions



Background
Compactions

Compaction
Priority

Optimizing Compactions

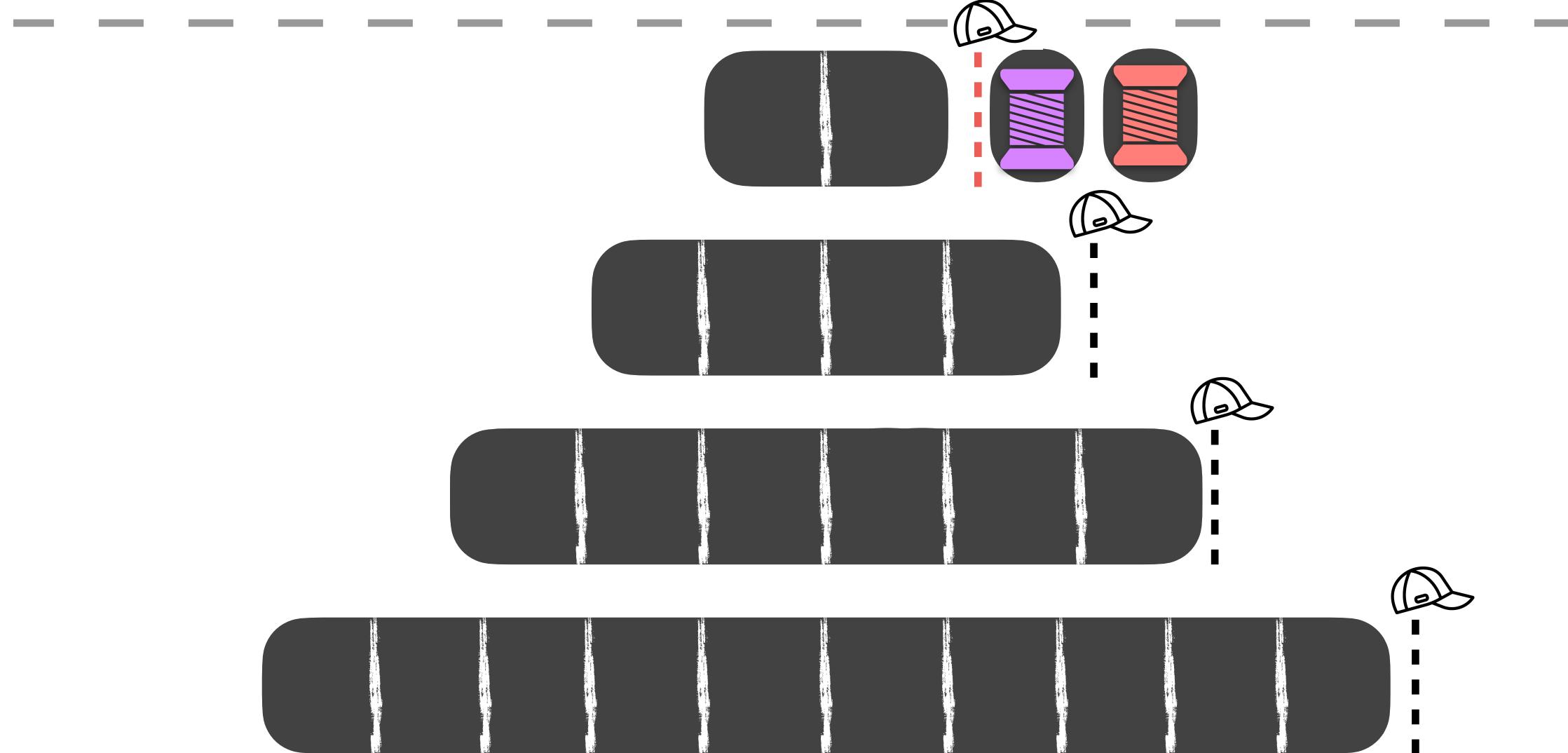
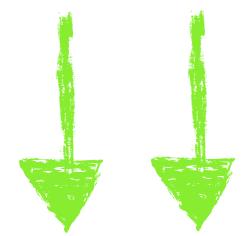


Background
Compactions

Compaction
Priority

Optimizing Compactions

write
pressure

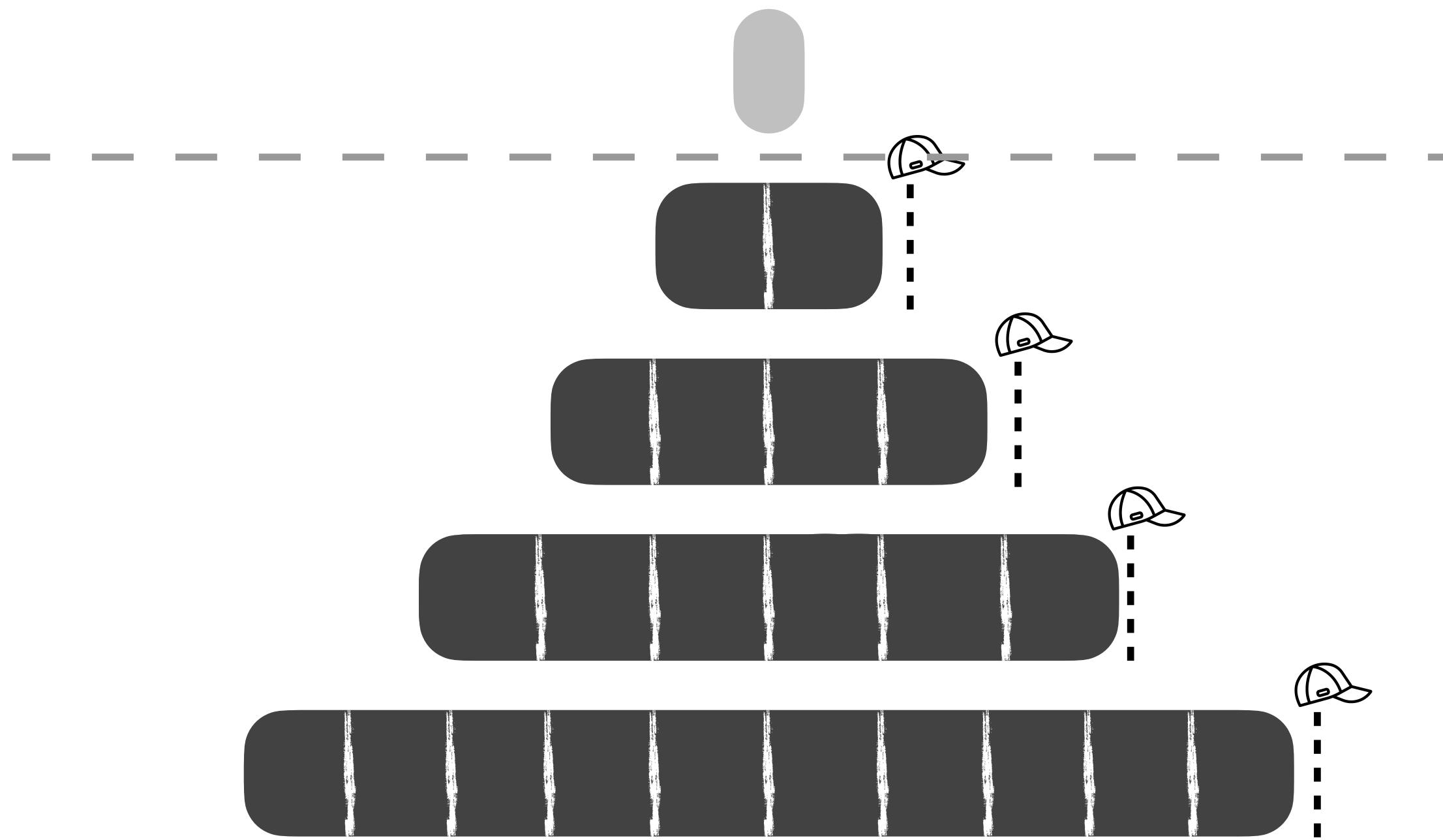


Background
Compactions

Compaction
Priority

- sustain heavy write bursts
- tree becomes out of shape

Optimizing Compactions

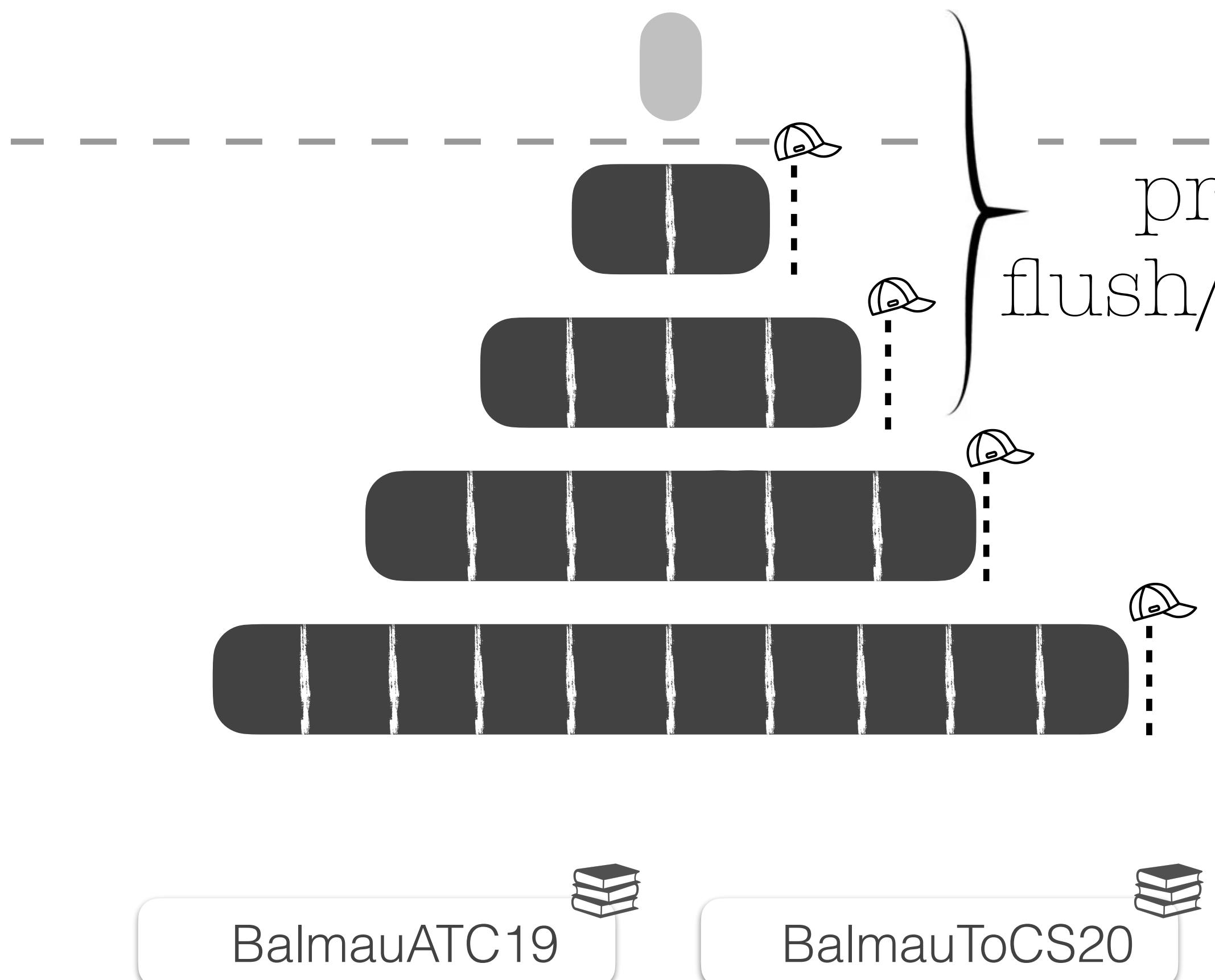


Background
Compactions

Compaction
Priority

I/O Scheduler

Optimizing Compactions



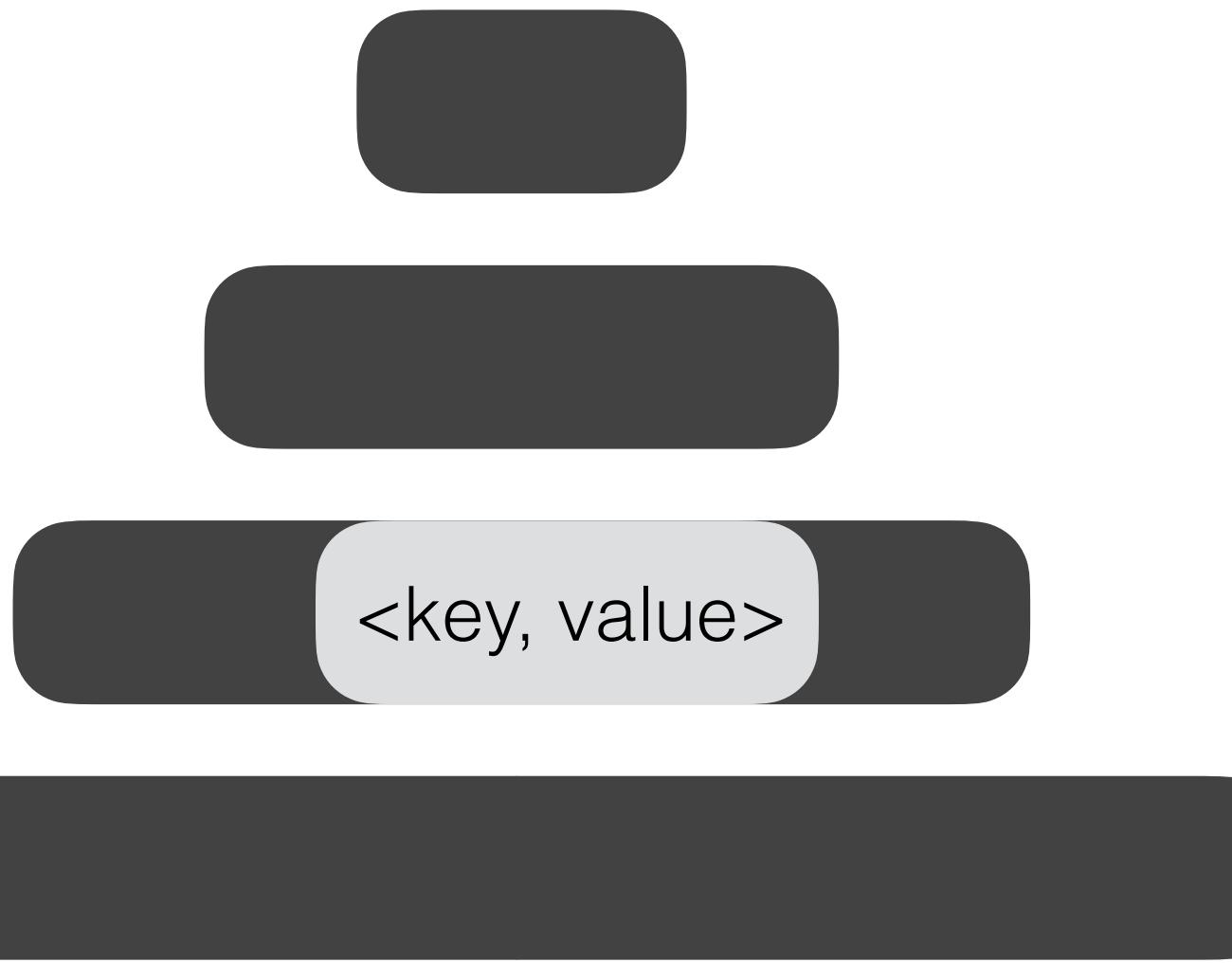
Background
Compactions

Compaction
Priority

I/O Scheduler

- eliminates write stalls
- no unnecessary high-priority compactions in lower levels

Data Placement Variations

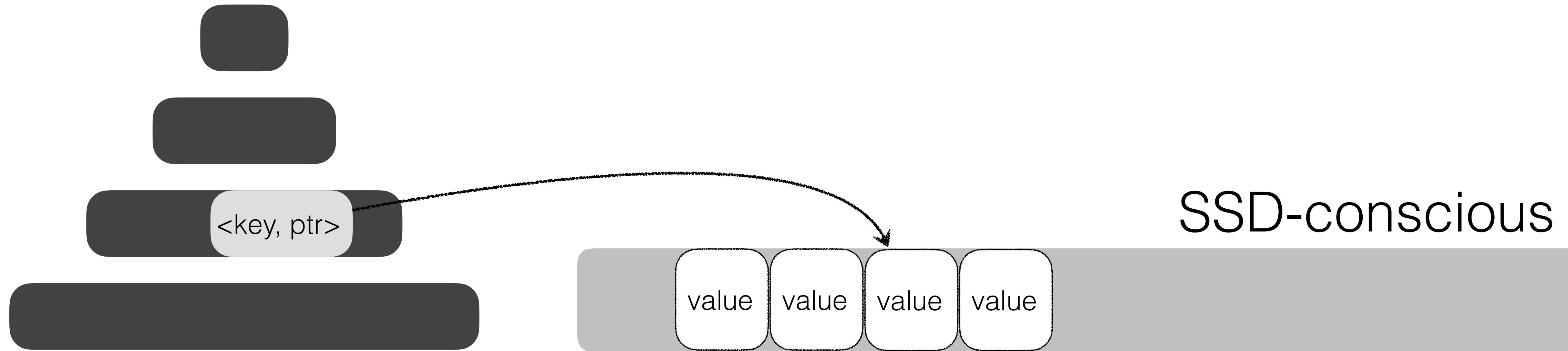


key-value separation



LuFAST16

Data Placement Variations



key-value separation

LuFAST16

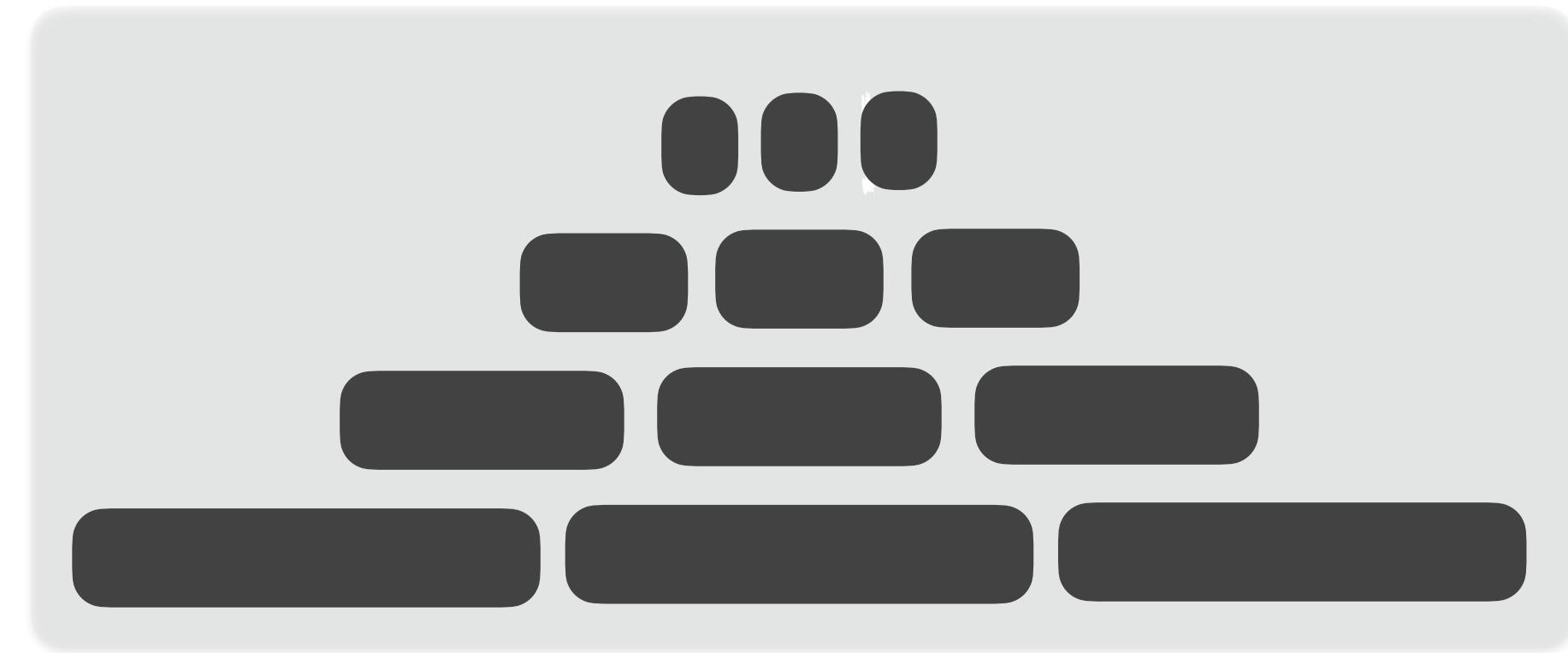
- reduced write amplification
- better read performance

Data Placement Variations



partitioning / sharding

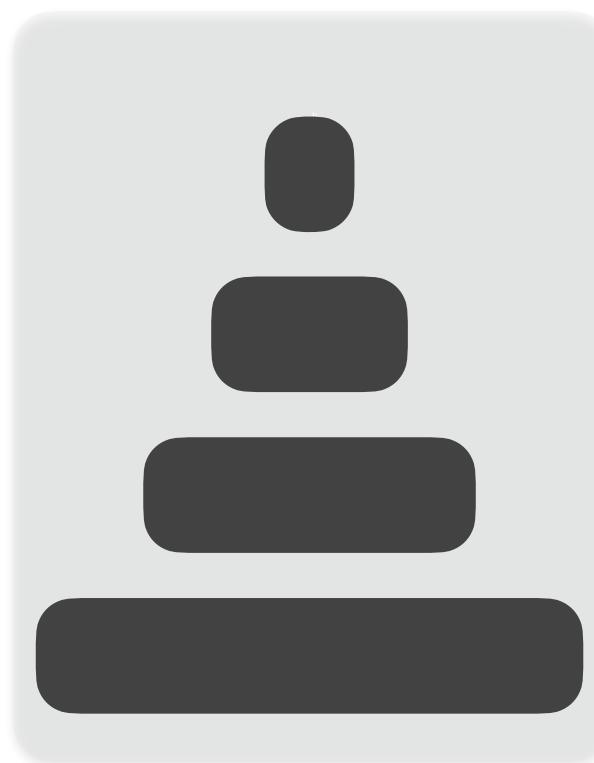
Data Placement Variations



storage

partitioning

RajuSOSP17



storage-1

storage-2

storage-3

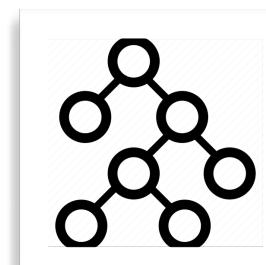
sharding

HuangSIGMOD21

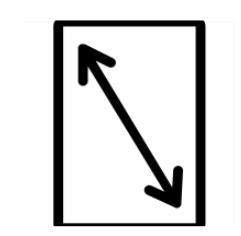
- improved ingestion throughput
- reduced write amplification

Summary: Ingestion Optimization

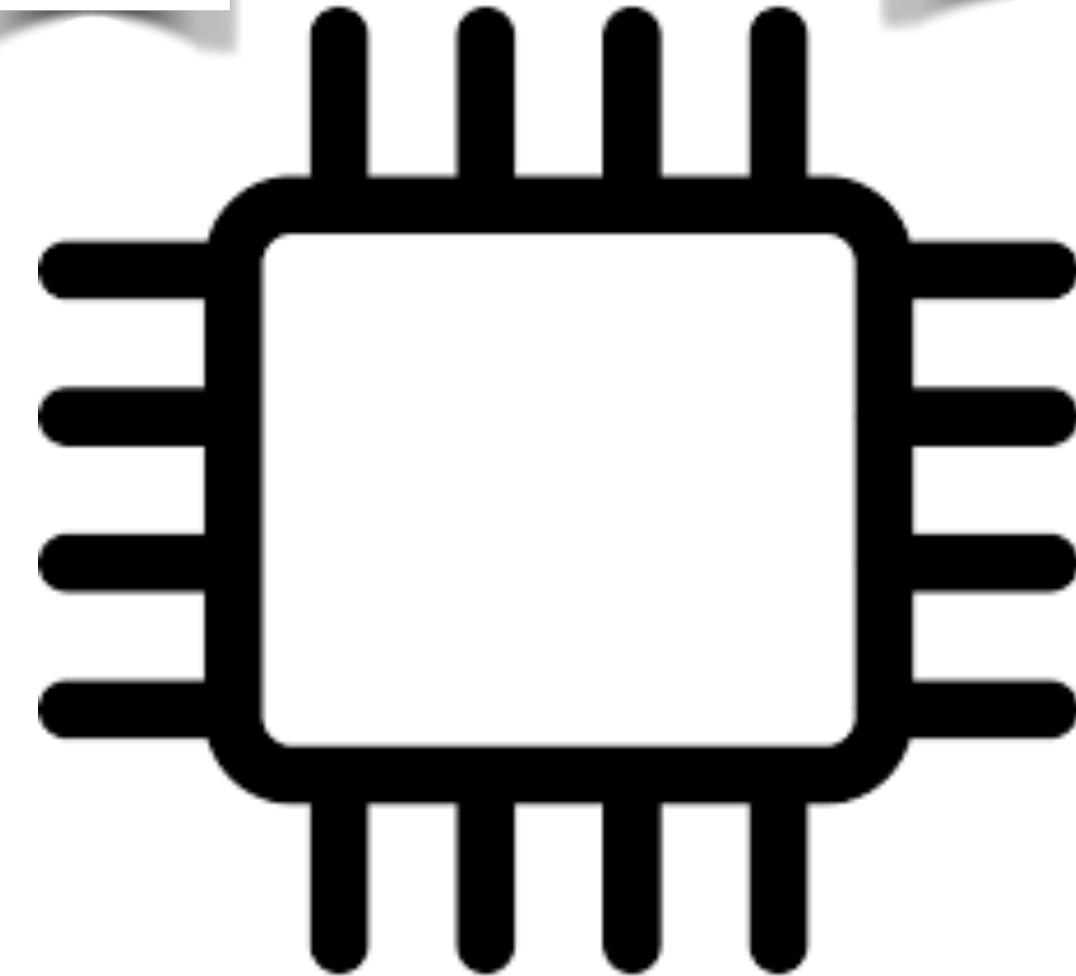
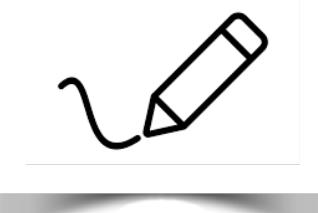
data
structure



size

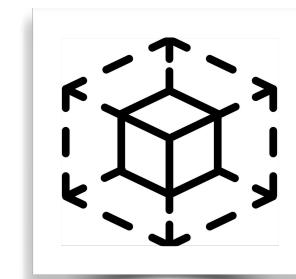


flush
strategy

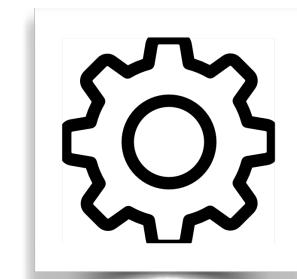


main memory
design elements

compaction
design space

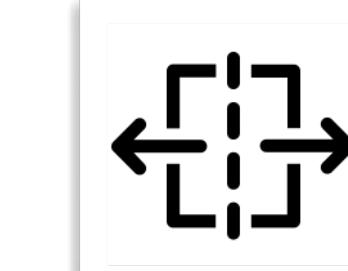


tuning
compactions



data layout
on storage

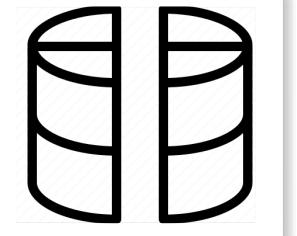
key-value
separation



data
partitioning



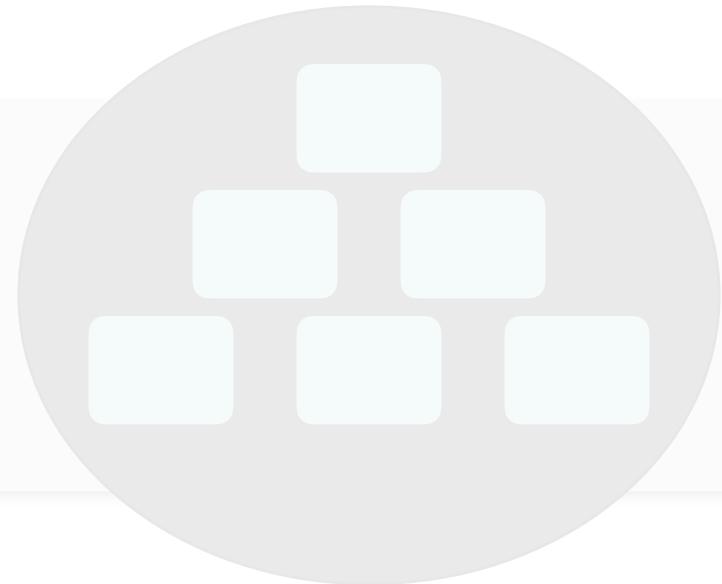
data
sharding



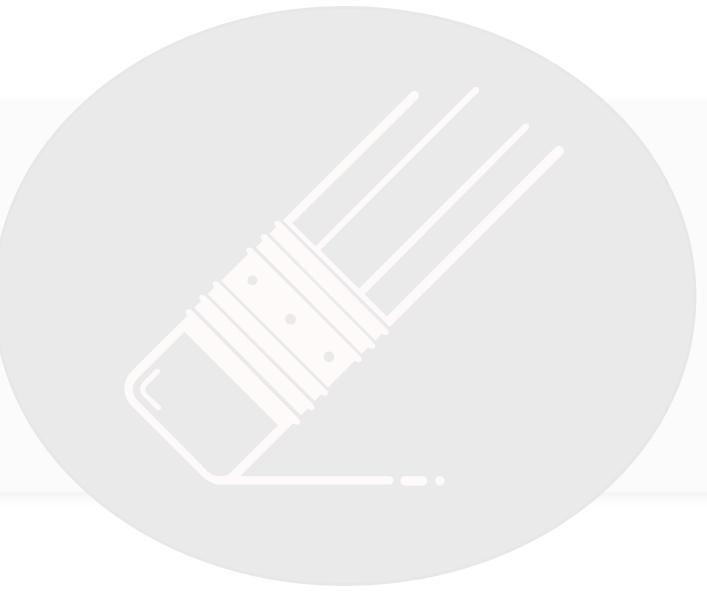
hardware-conscious
designs

Outline

Part 1: **LSM Basics**



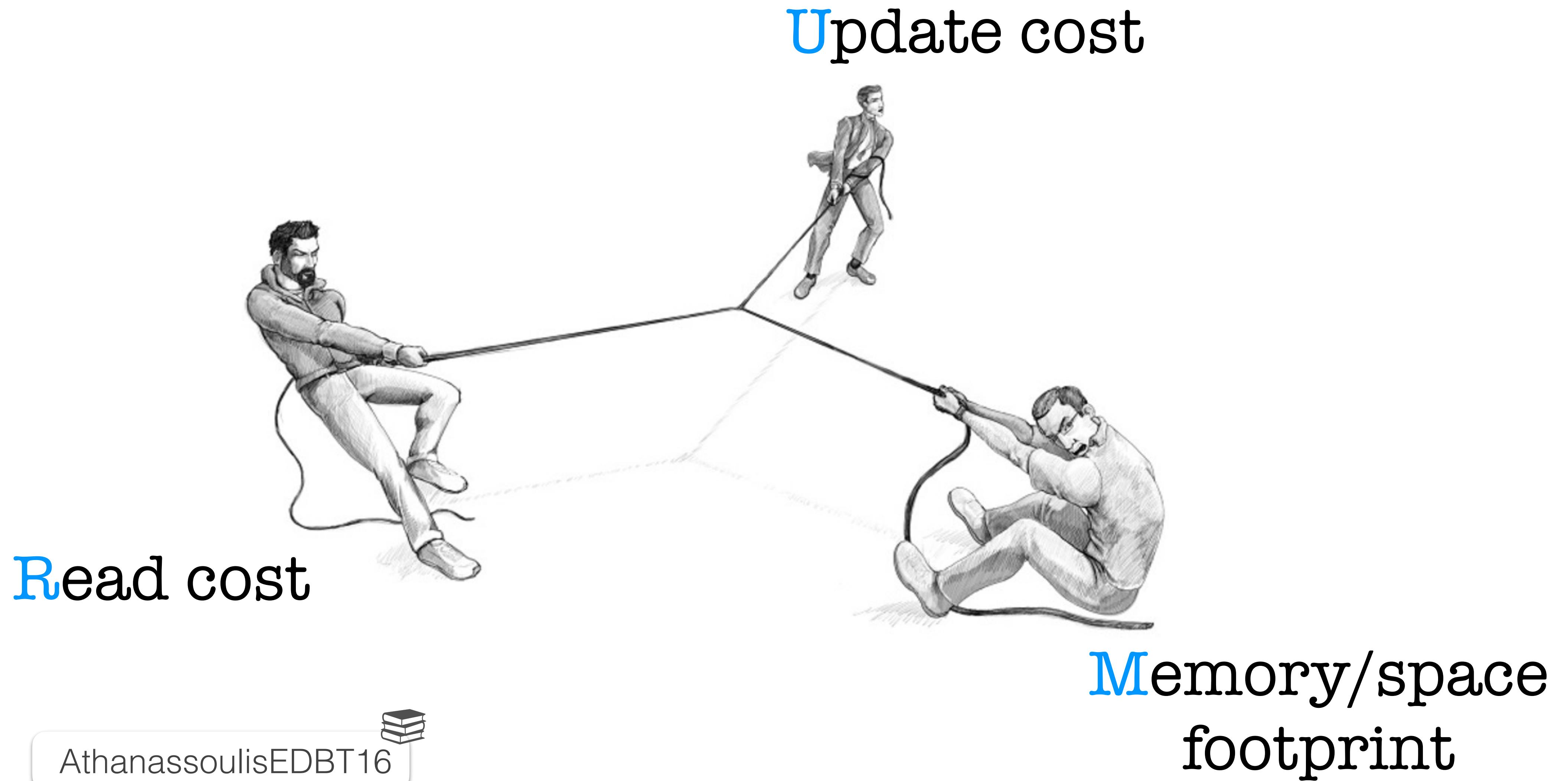
Part 2: Optimizing Ingestion in LSMS



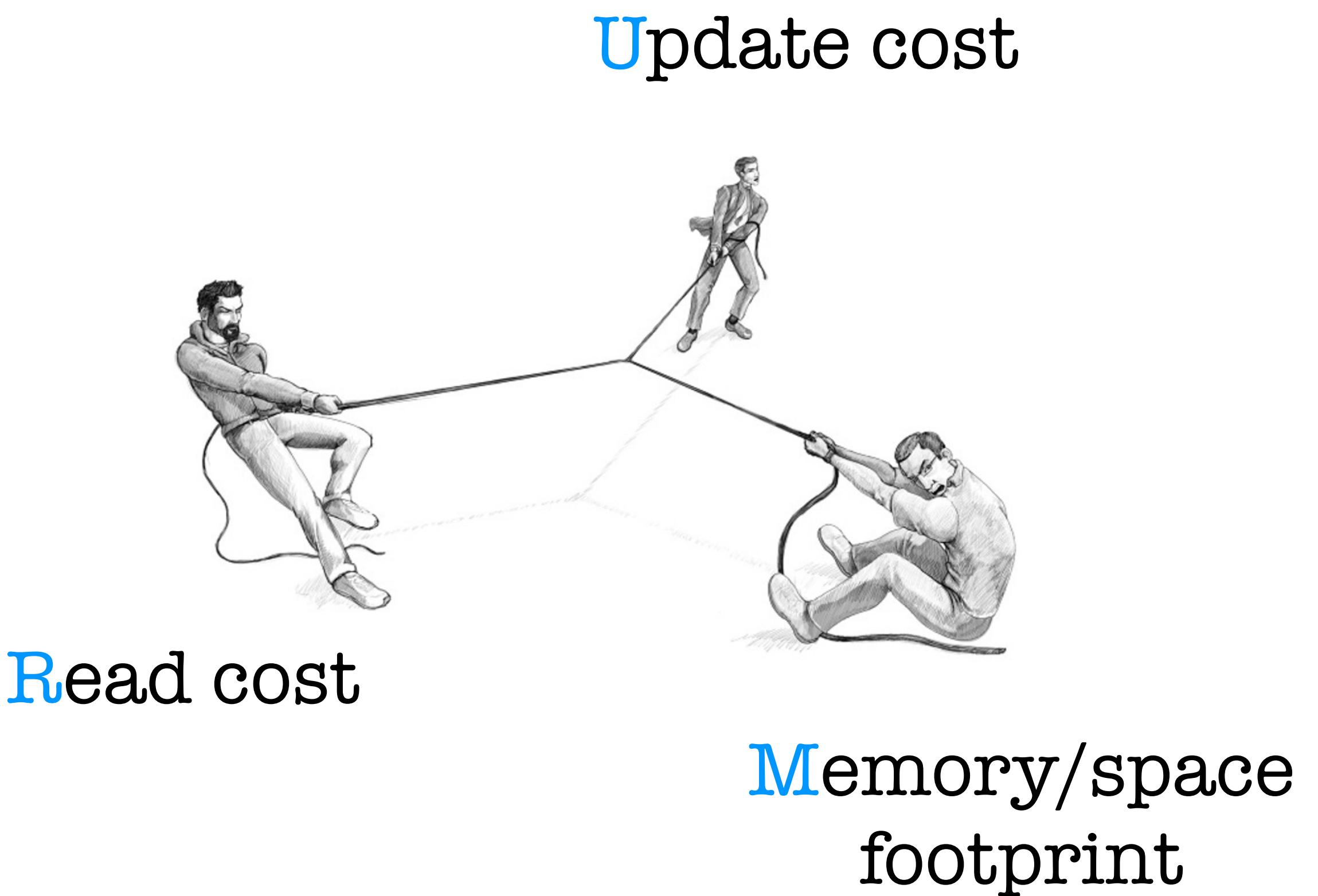
Part 3: Navigating the **LSM Design Space**



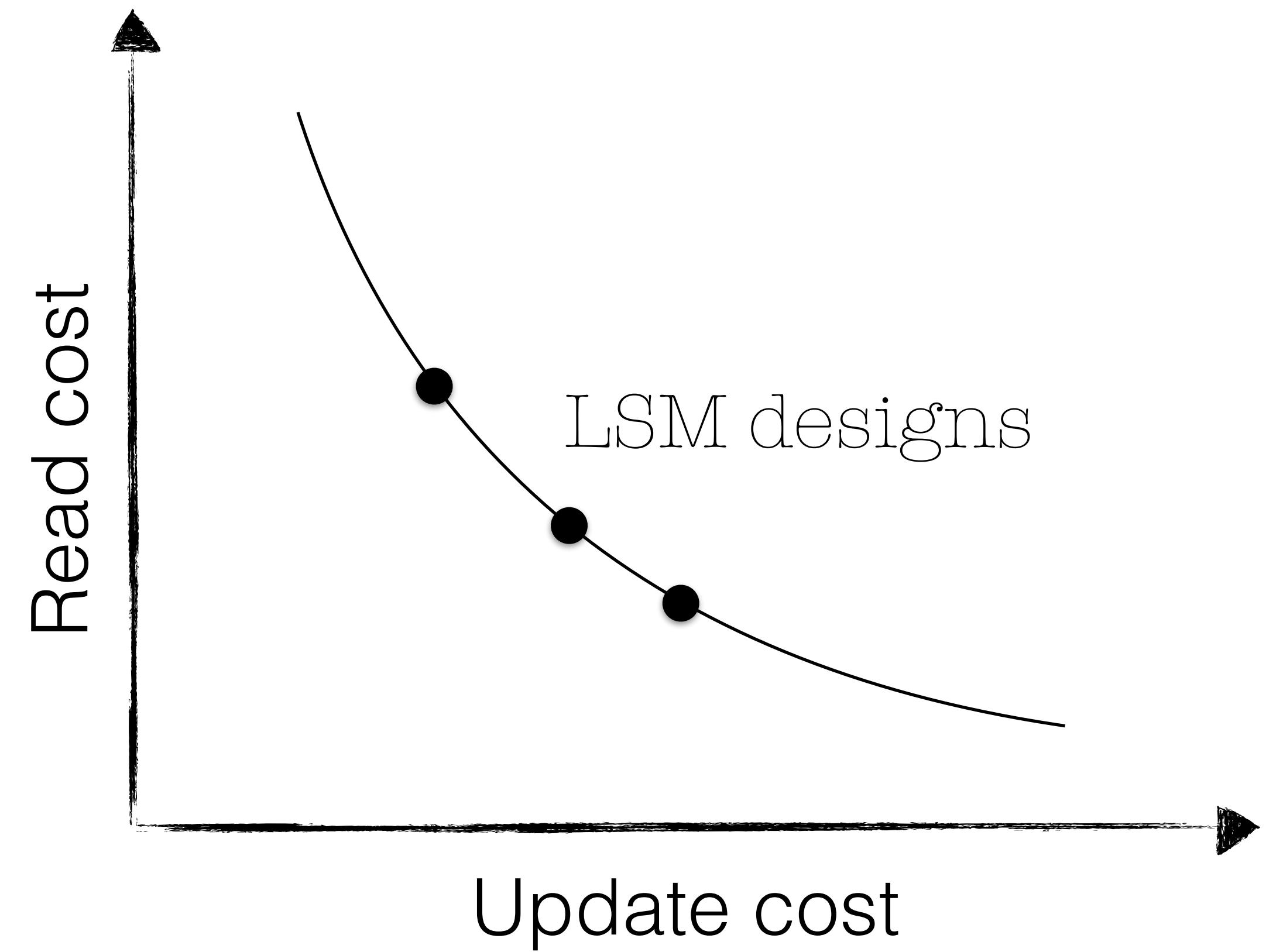
LSM Design Space



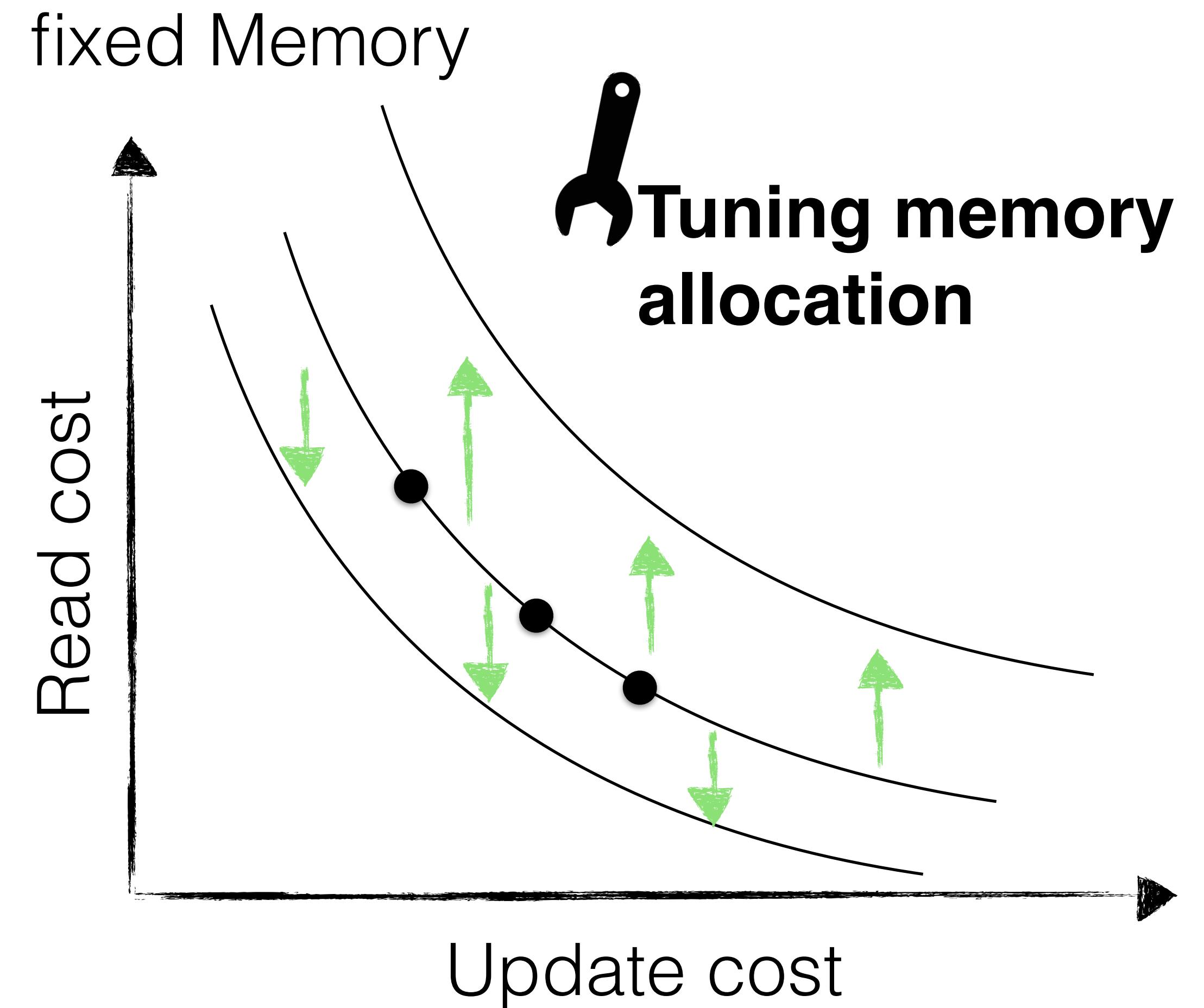
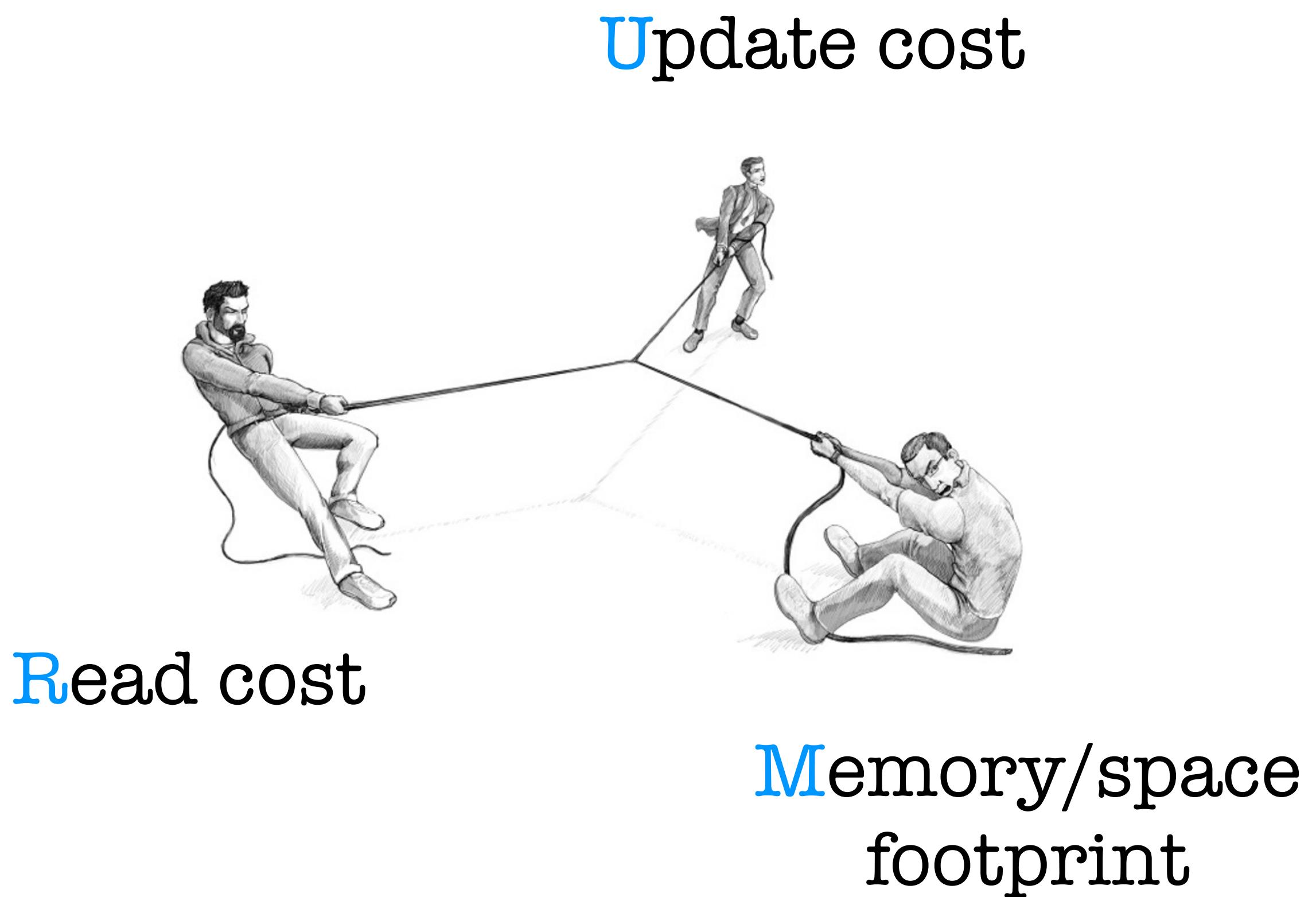
LSM Design Space



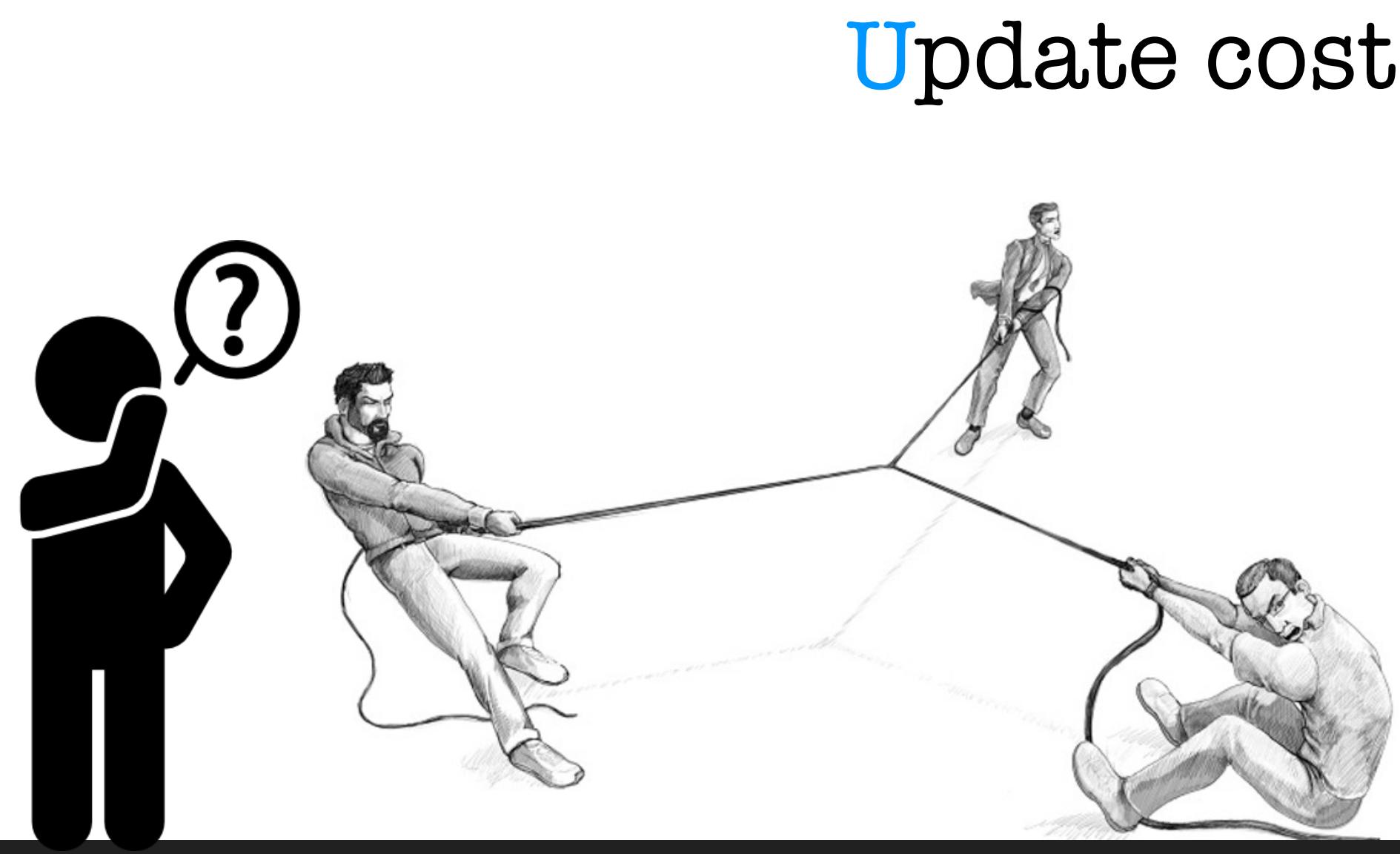
fixed Memory



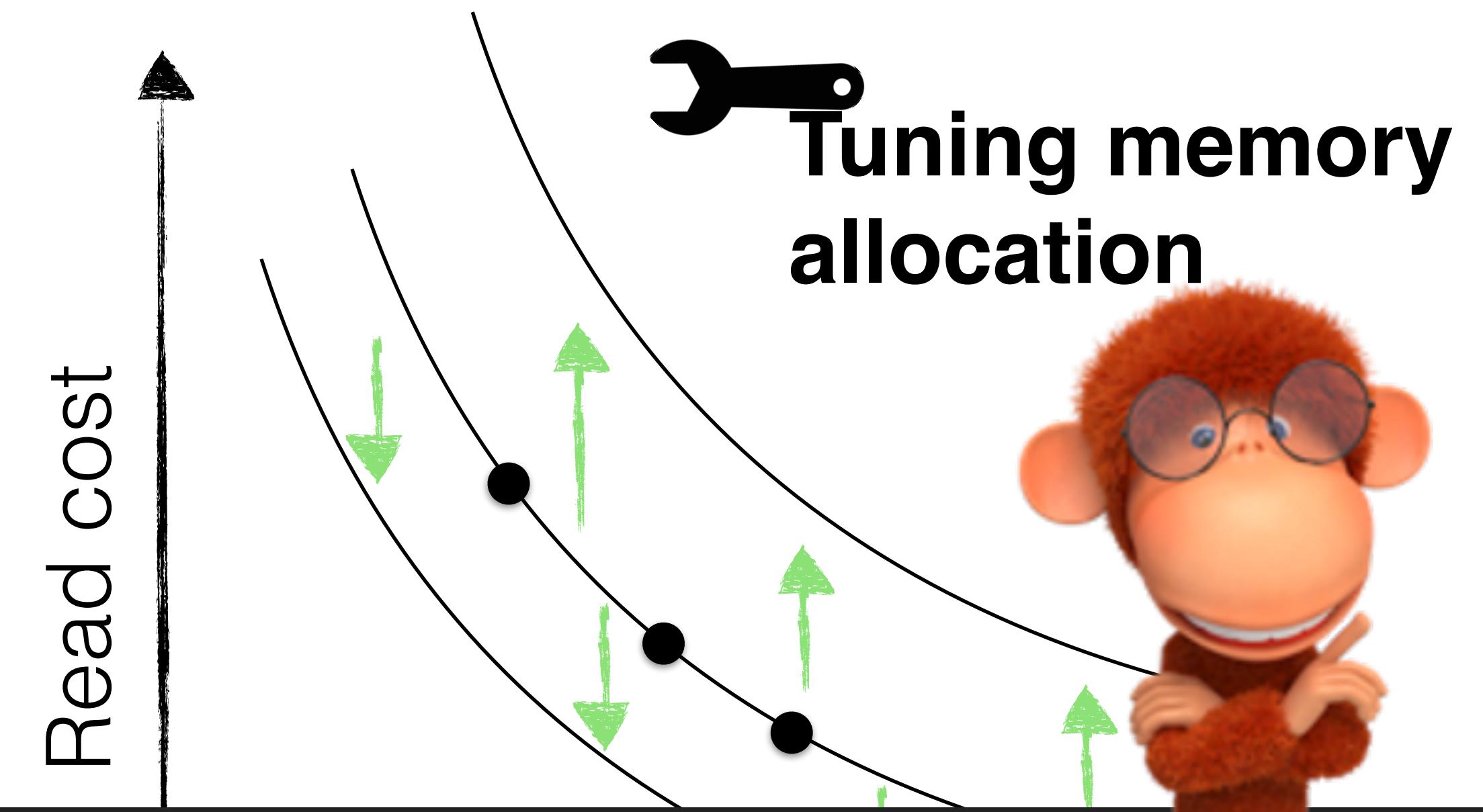
LSM Design Space



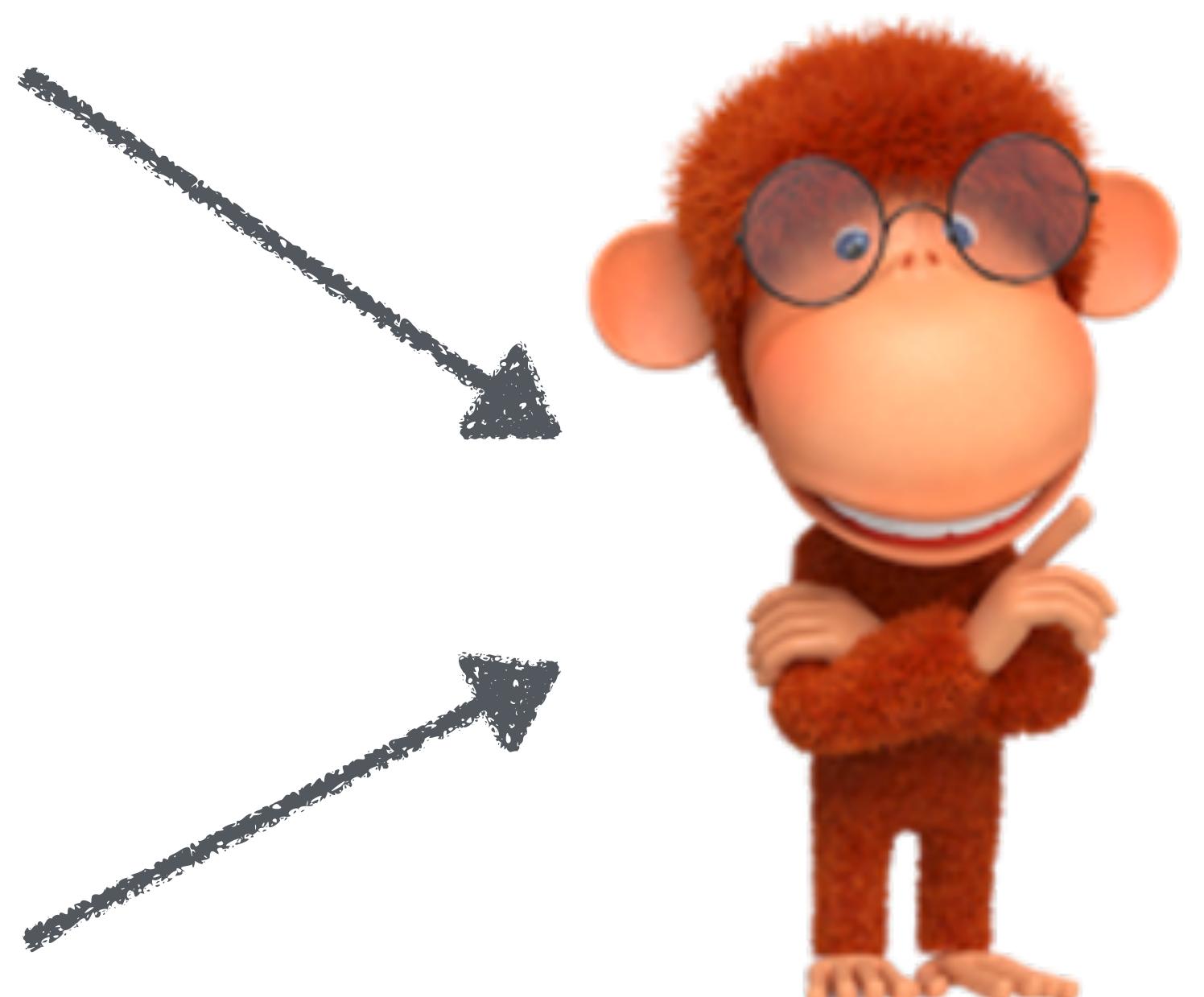
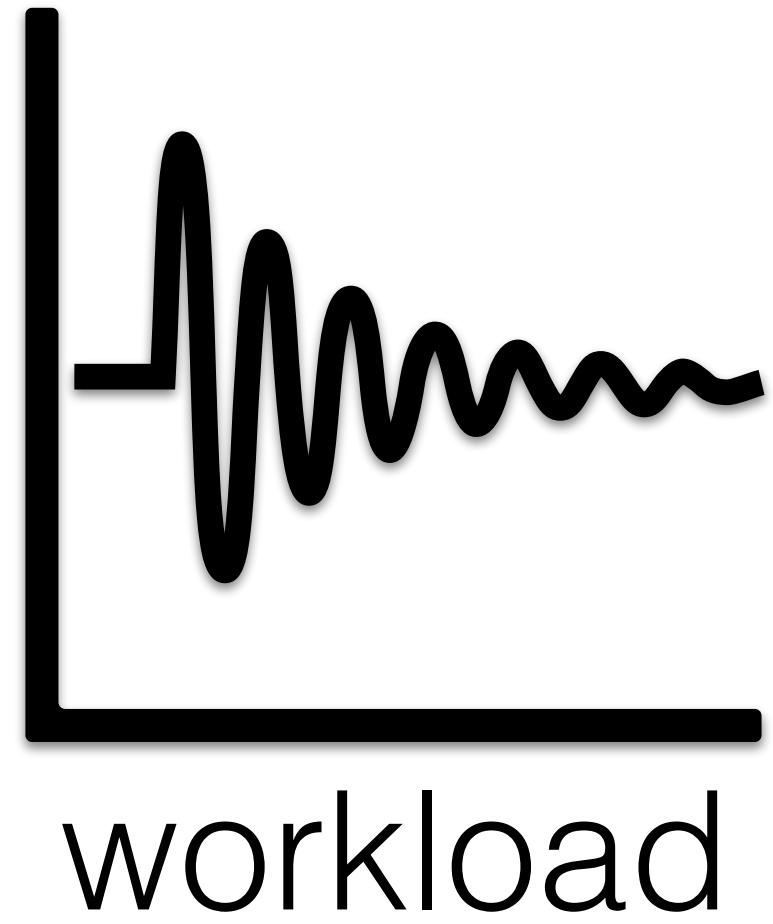
LSM Design Space



fixed Memory

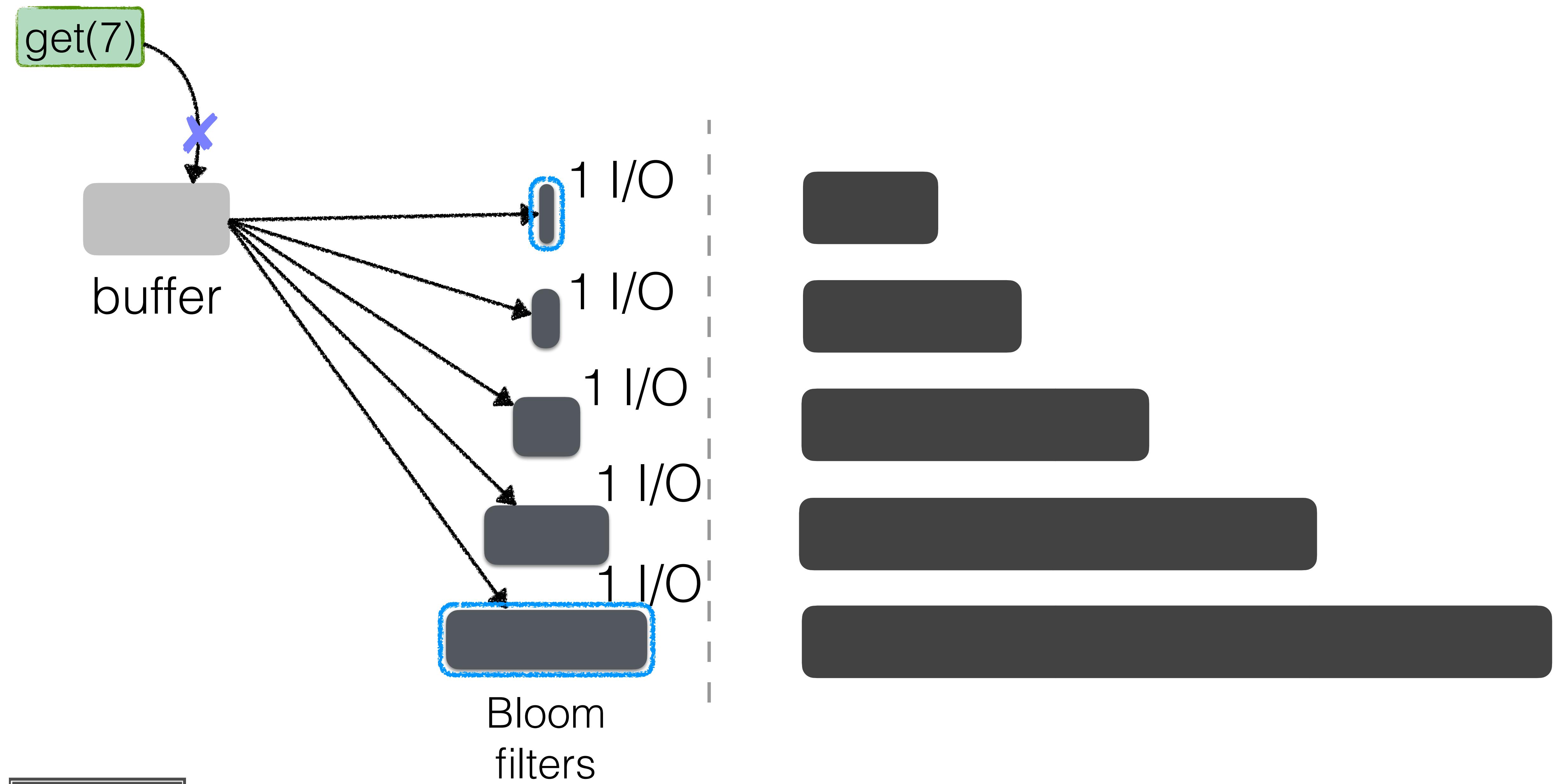


How to optimally allocate the available memory?

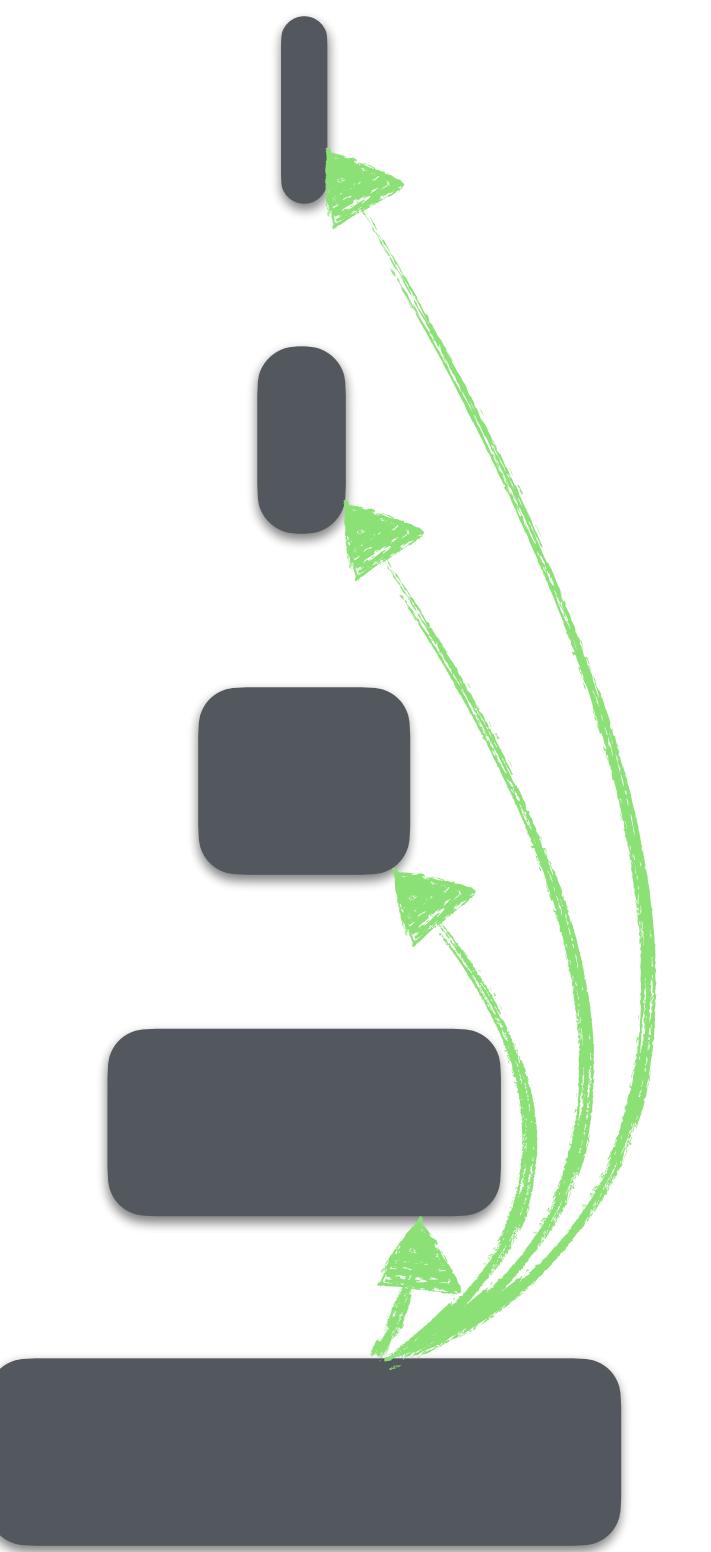


How to allocate memory
between buffer and BF

How to allocate memory
among BFs in LSM



Bloom filters



FPR

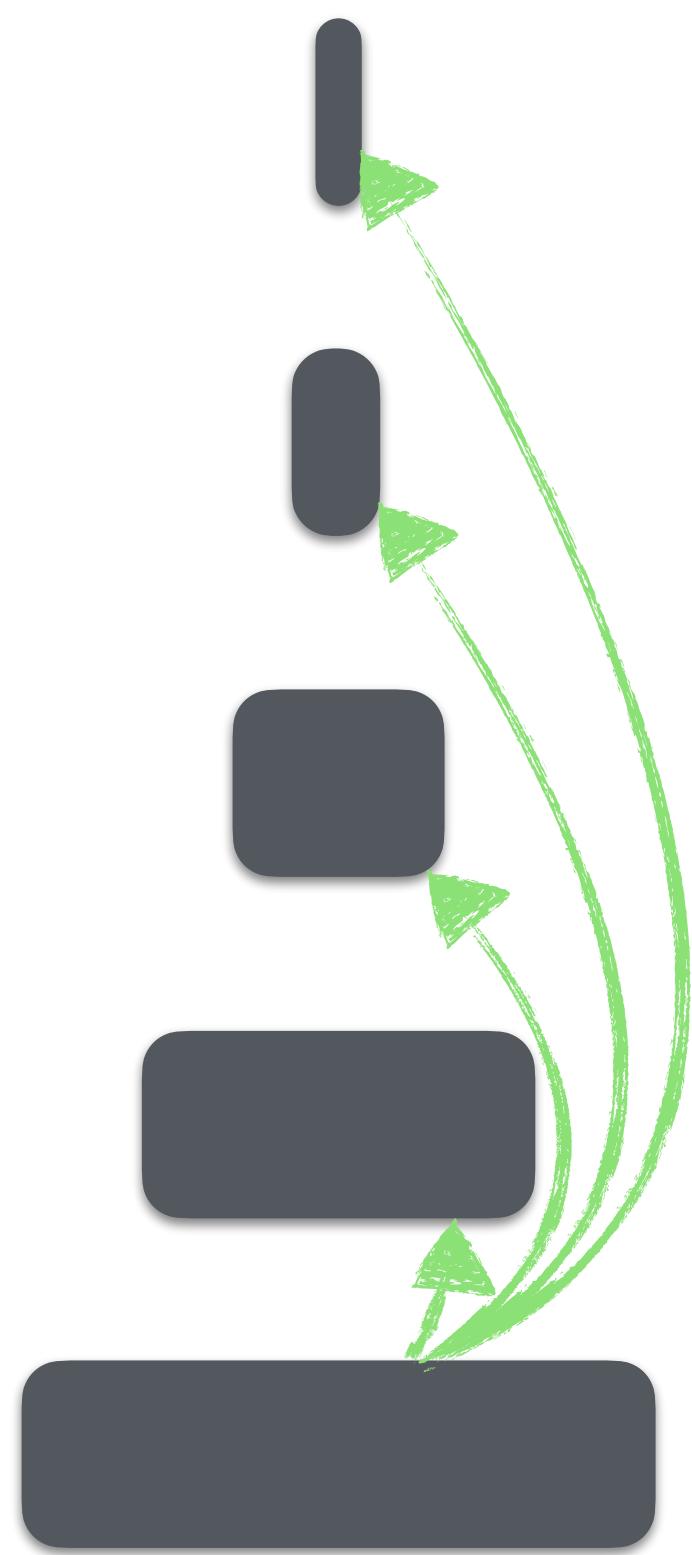
Monkey FPR

$$\begin{array}{ccl} \phi & > & \phi_1 = \phi_0/T^4 \\ \phi & > & \phi_2 = \phi_0/T^3 \\ \phi & > & \phi_3 = \phi_0/T^2 \\ \phi & > & \phi_4 = \phi_0/T \\ \phi & < & \phi_5 = \phi_0 \end{array}$$

↑
exponentially
decreasing



Bloom filters



FPR

Monkey FPR

$$\begin{array}{ccl} \phi & \textcolor{green}{>} & \phi_1 = \phi_0/T^4 \\ \phi & \textcolor{green}{>} & \phi_2 = \phi_0/T^3 \\ \phi & \textcolor{green}{>} & \phi_3 = \phi_0/T^2 \\ \phi & \textcolor{green}{>} & \phi_4 = \phi_0/T \\ \phi & \textcolor{red}{<} & \phi_5 = \phi_0 \end{array}$$

$$L \cdot \phi$$

$$\sum_i \phi_i = c \cdot \phi_0$$

exponentially decreasing

point lookup cost

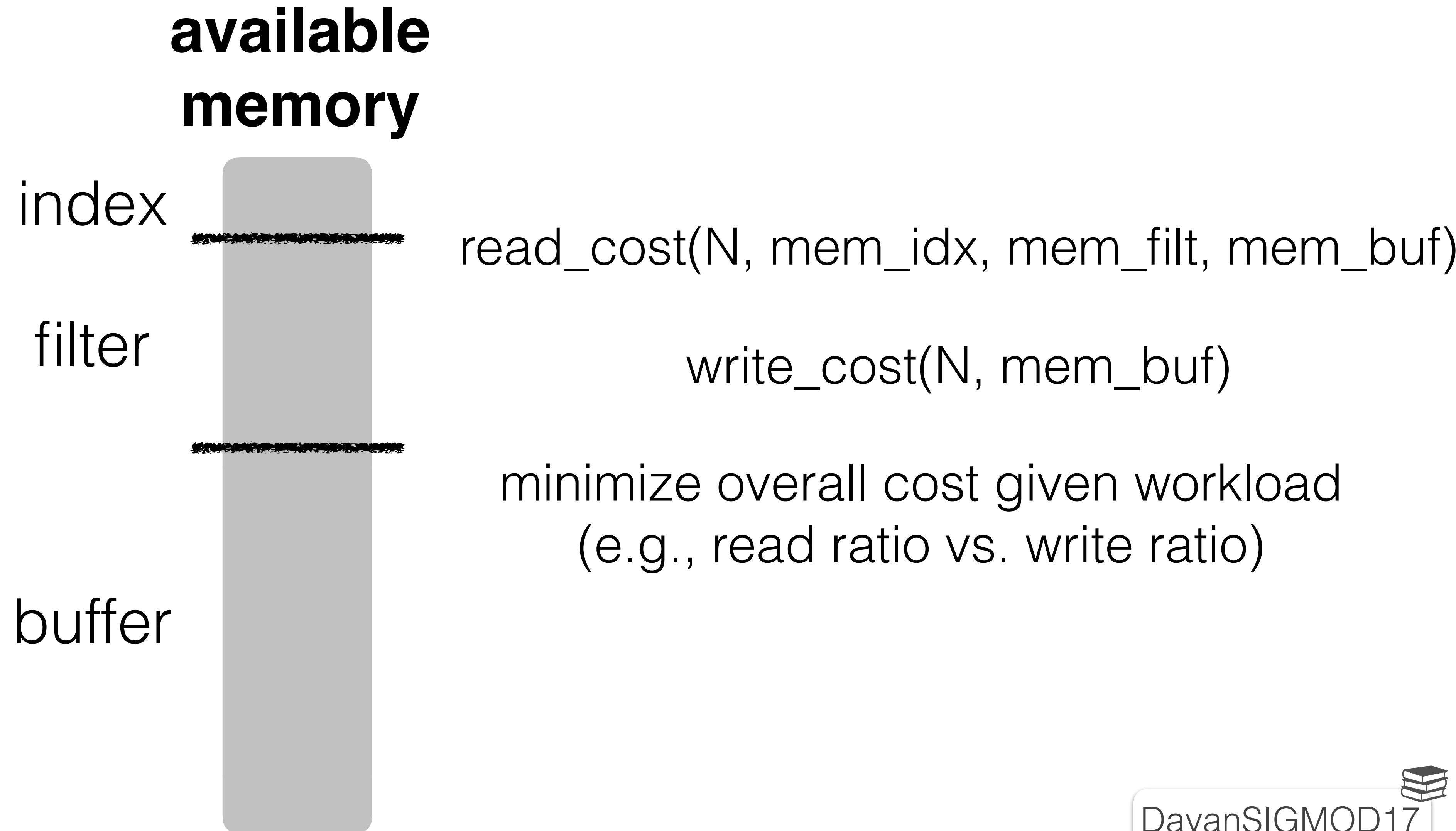
$$\mathcal{O}(L \cdot \phi)$$



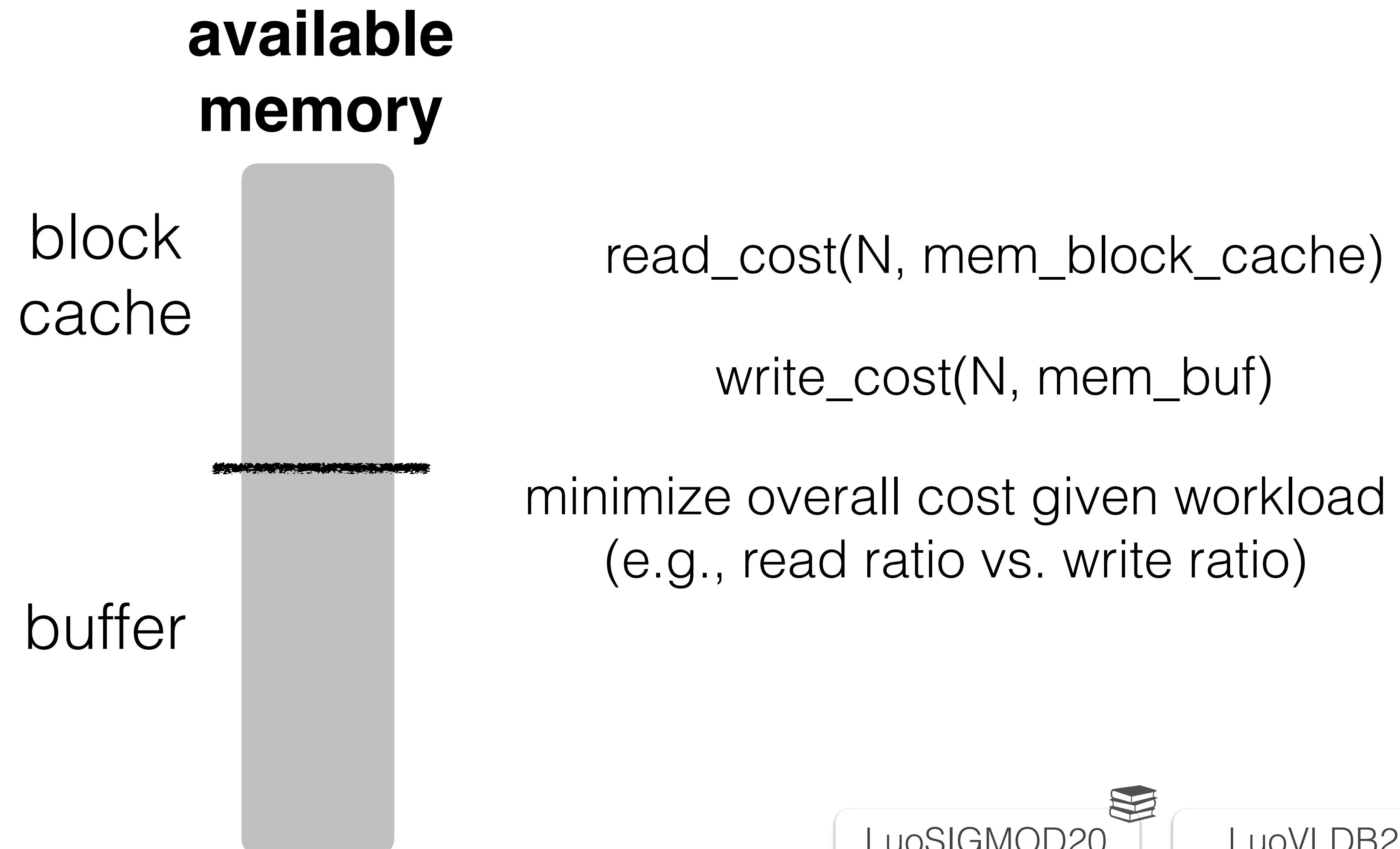
$$\mathcal{O}(c \cdot \phi_0)$$



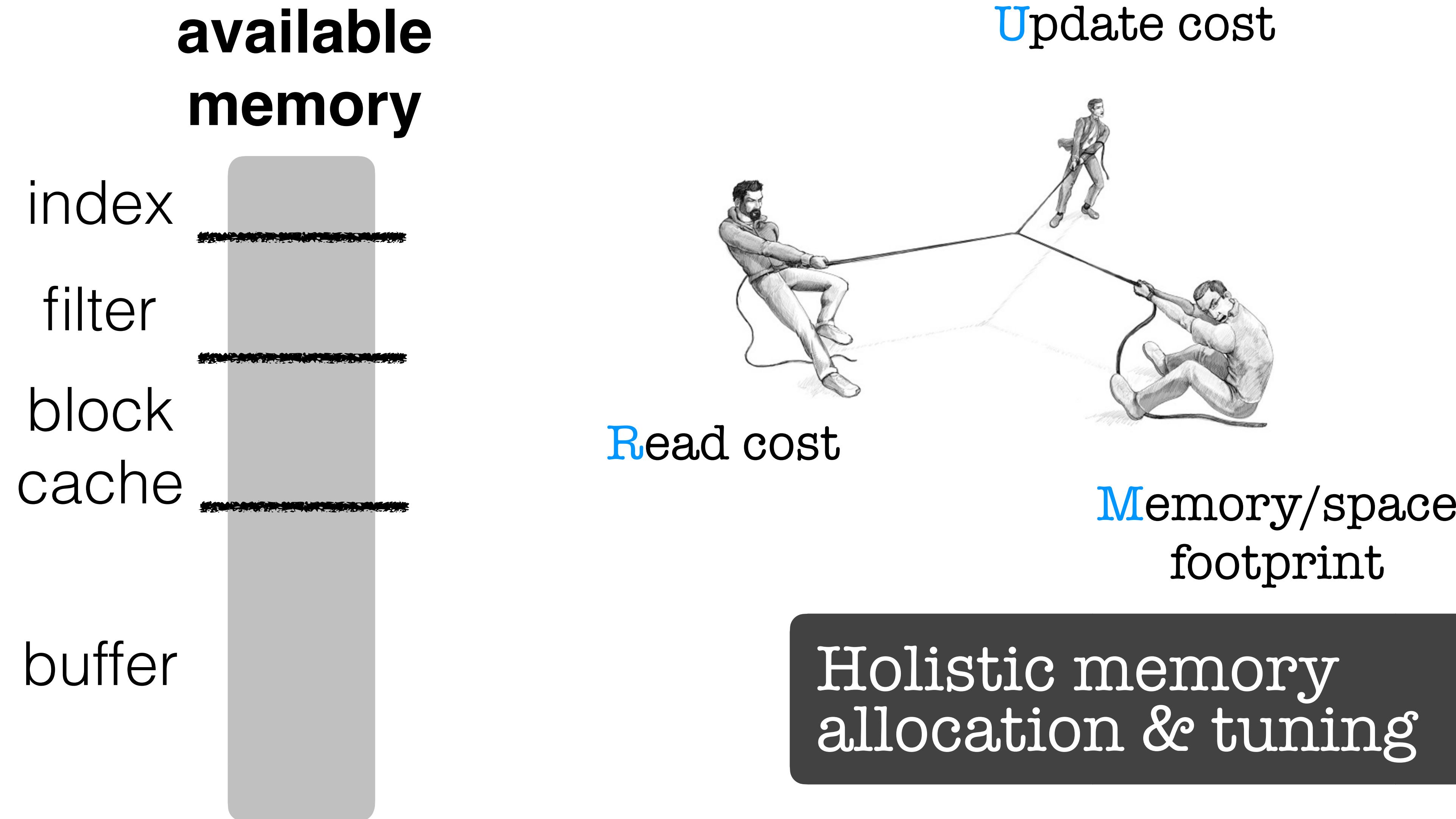
The Optimal Memory Allocation



The Optimal Memory Allocation

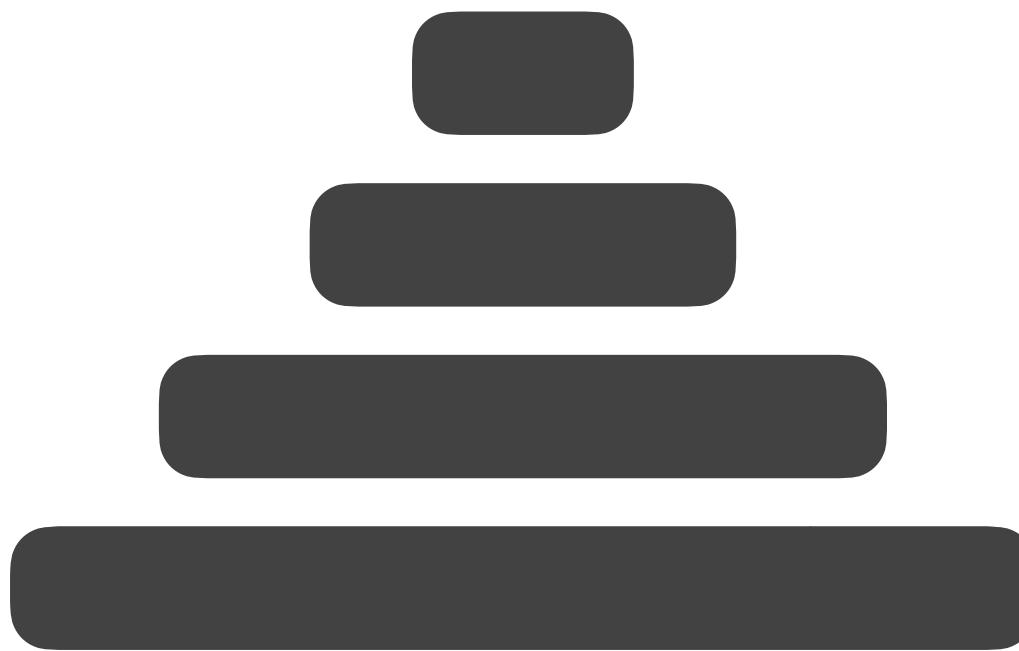


The Optimal Memory Allocation

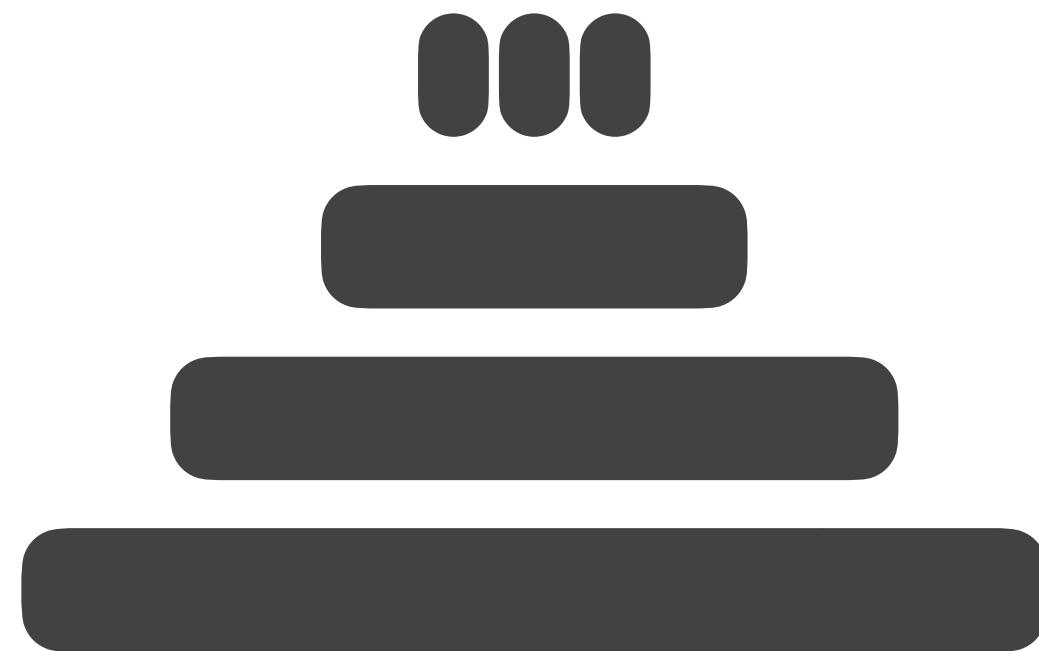


Storage Layer Design Continuum

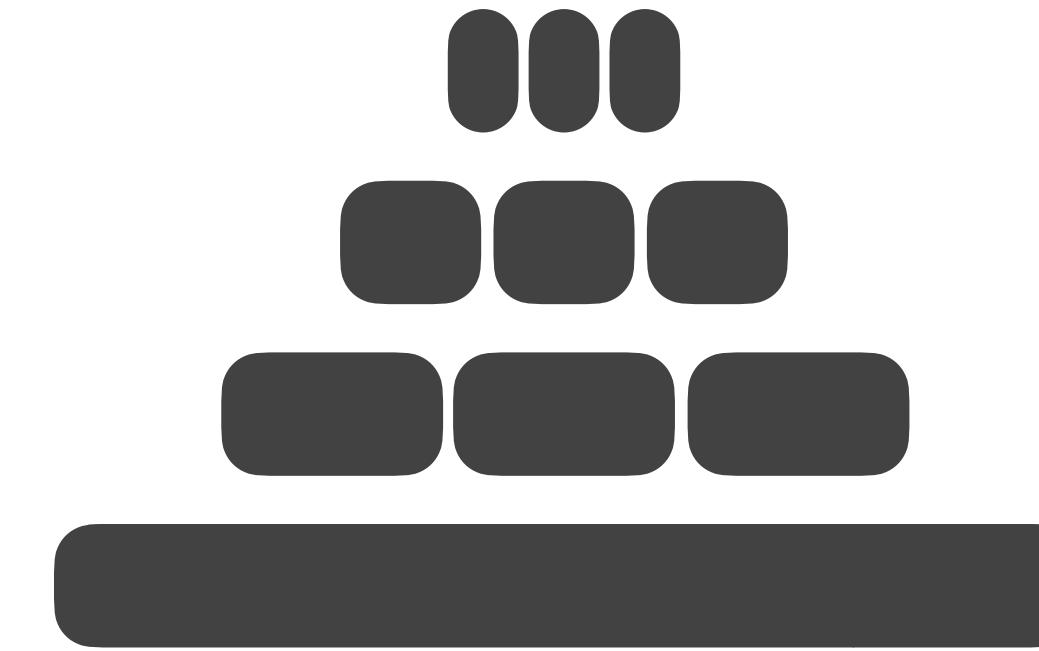
leveling



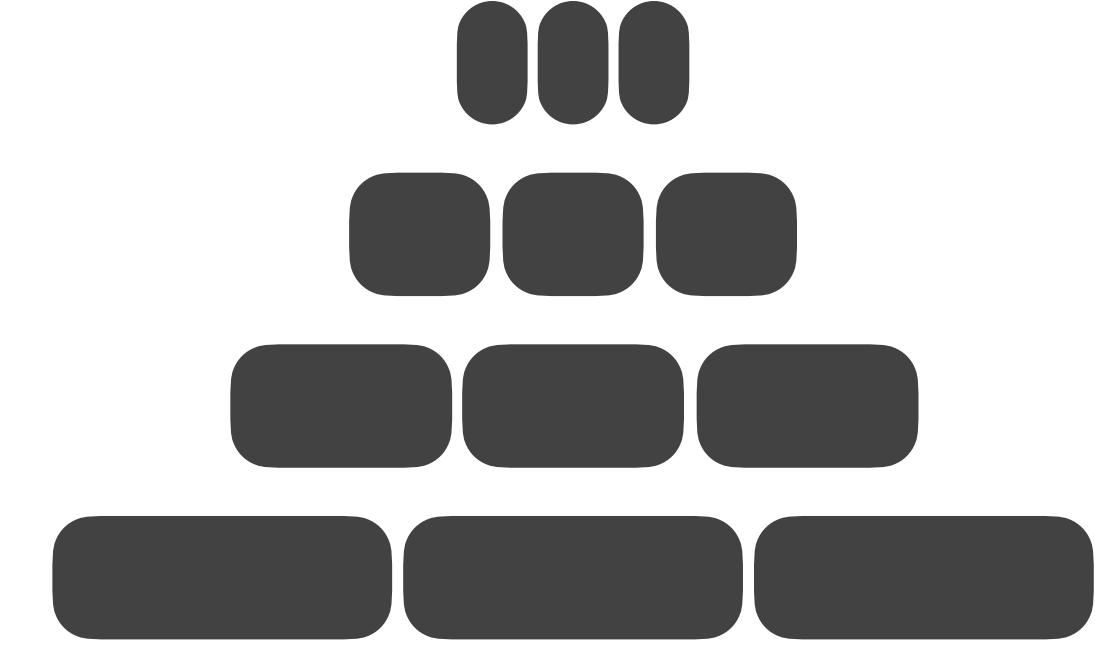
1-leveling



L-leveling



tiering



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

Storage Layer Design Continuum

leveling

1-leveling

L-leveling

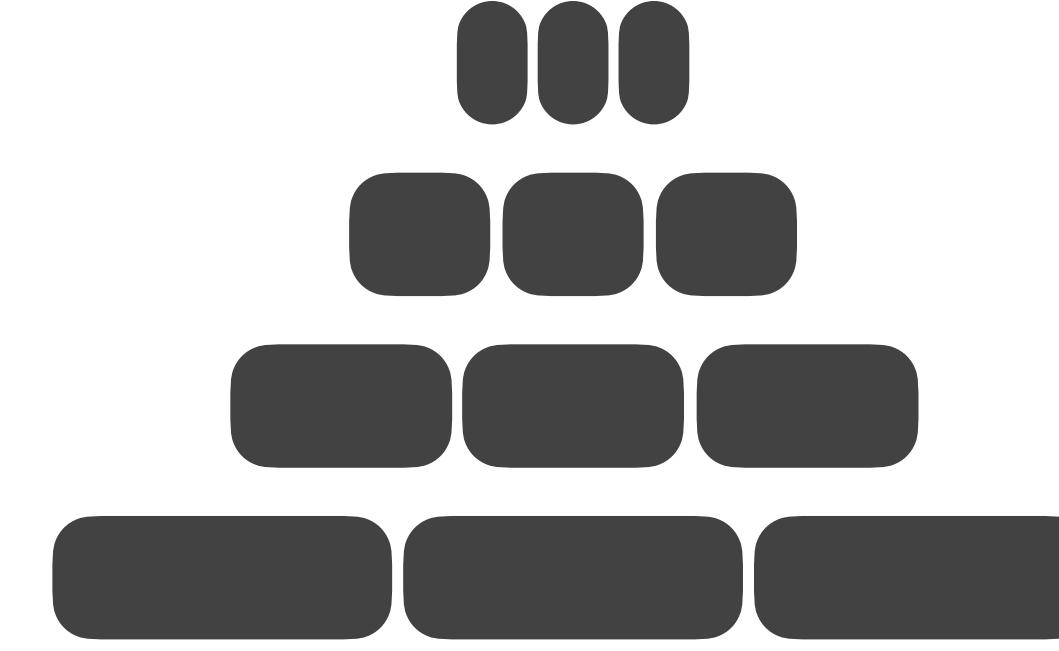
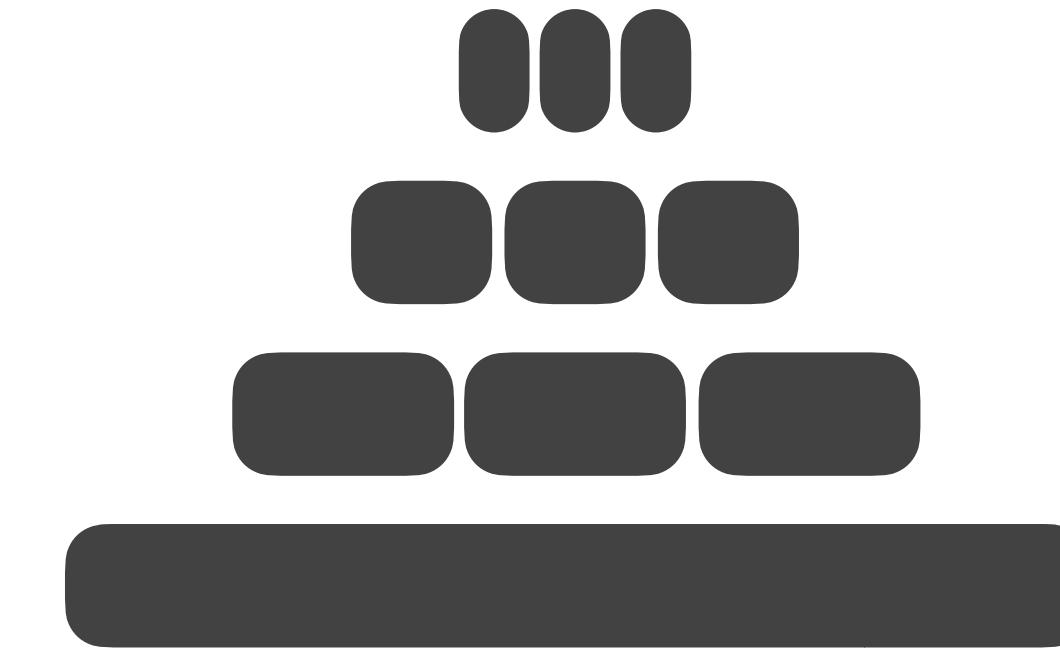
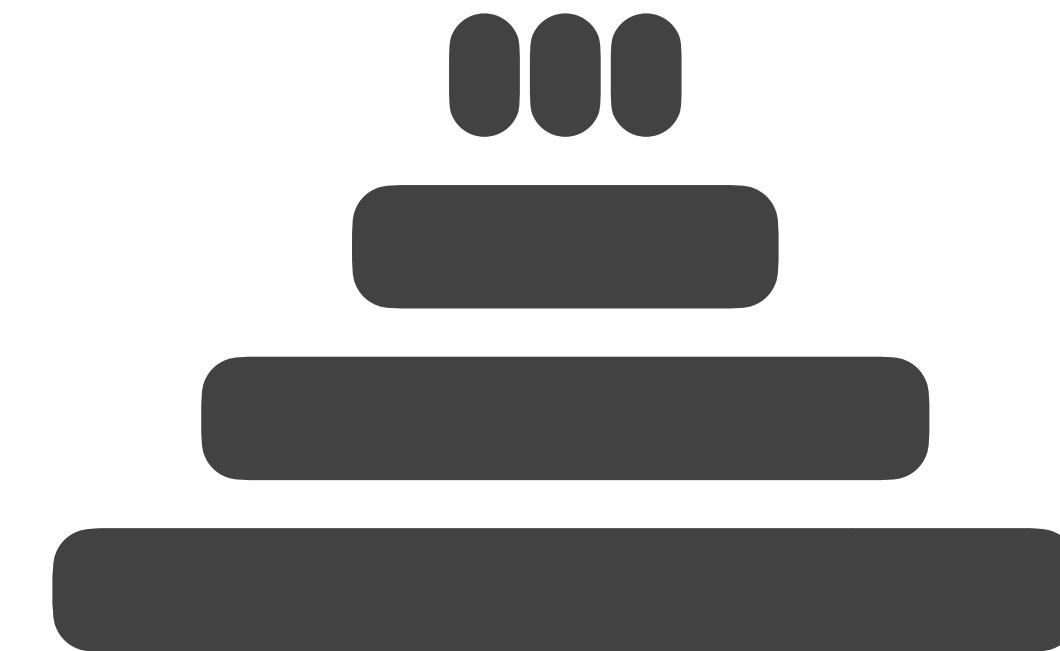
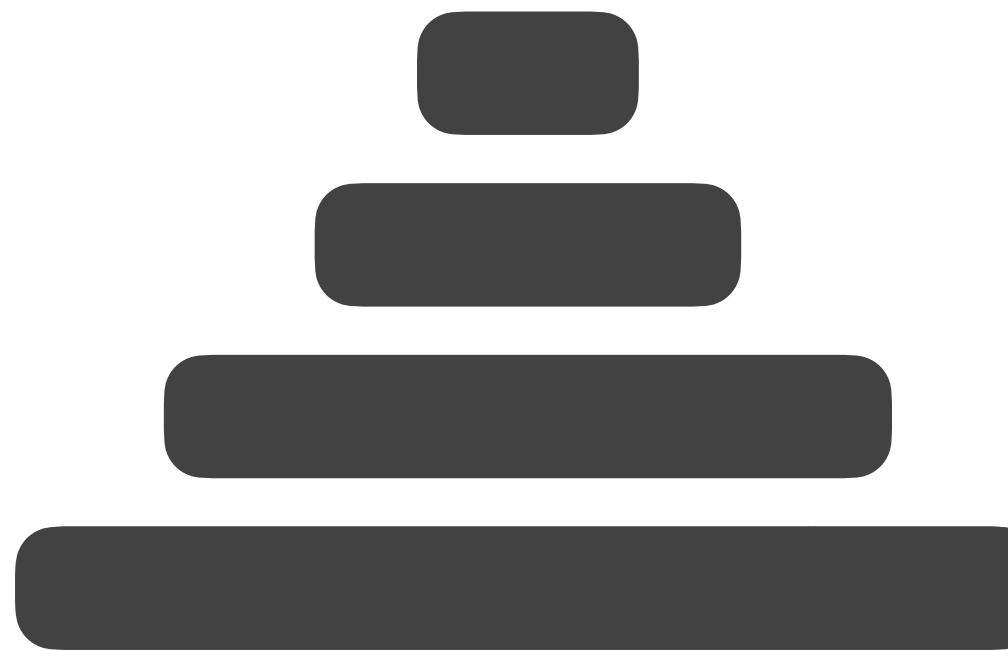
tiering

$(T, 0)$

$(T, 1)$

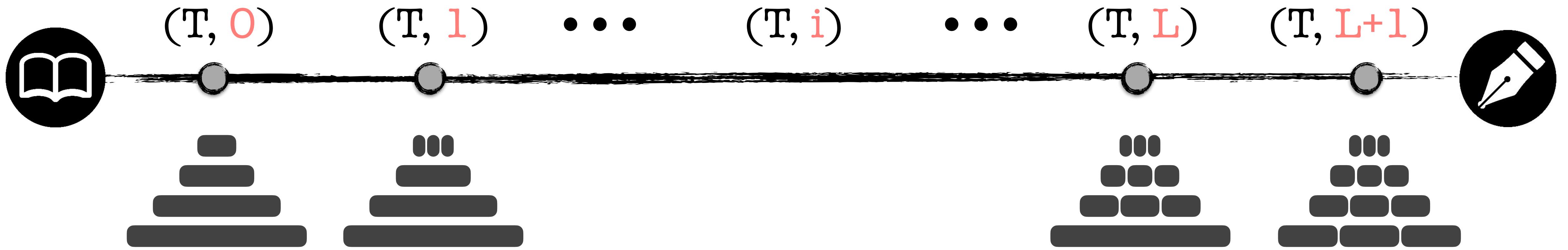
(T, L)

$(T, L+1)$

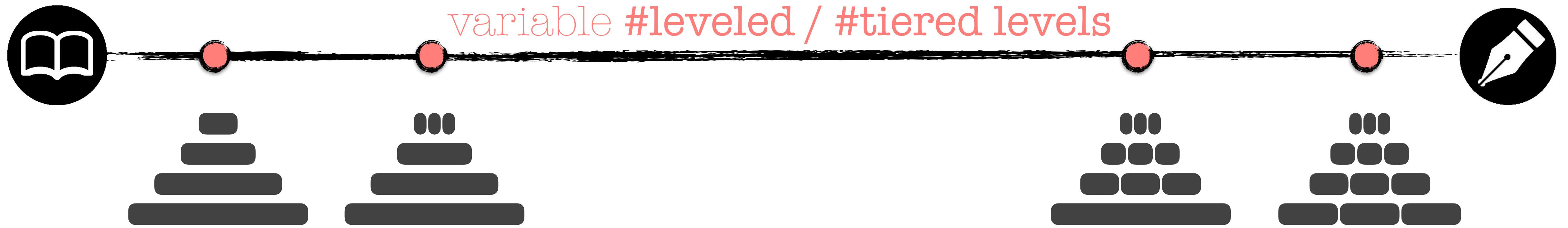


Any design can be defined by the tuple-set: (T, i)

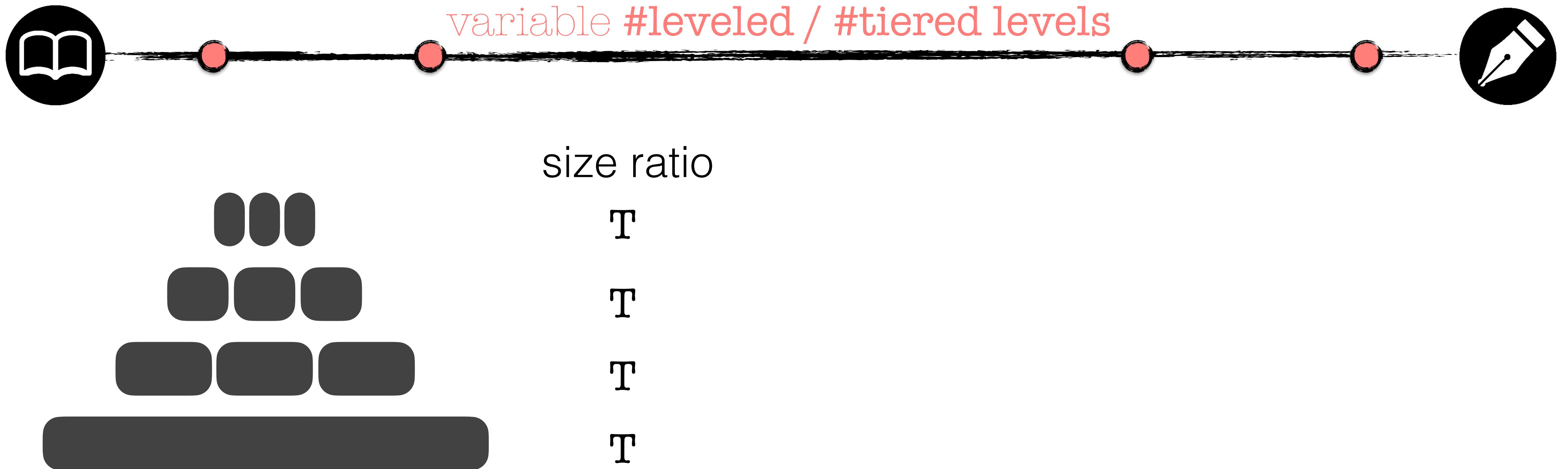
Storage Layer Design Continuum



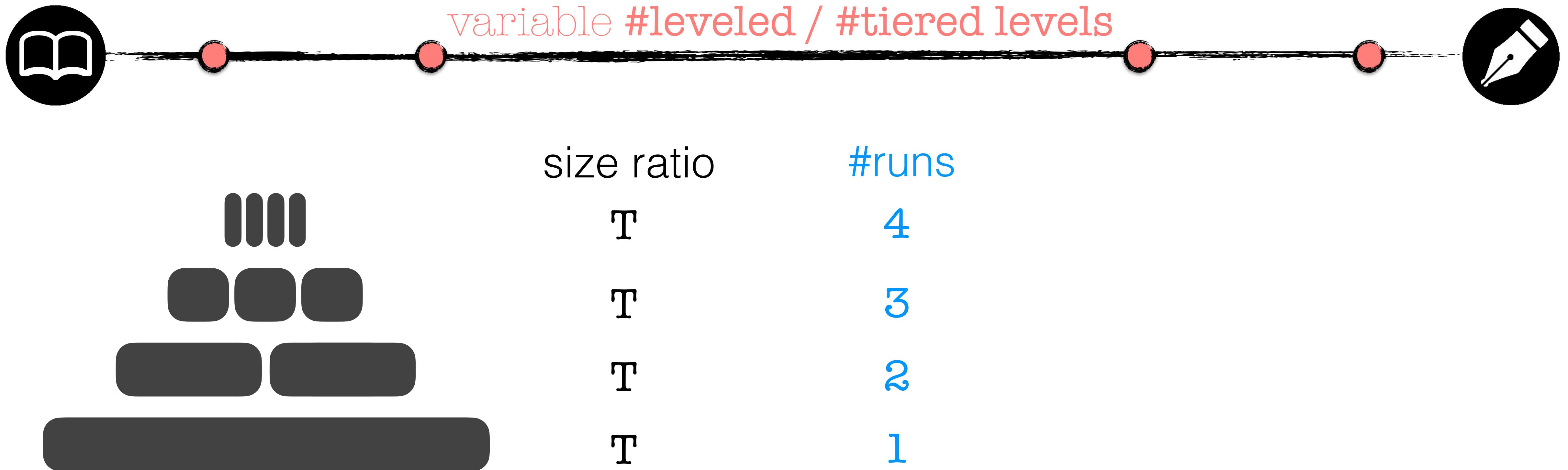
Storage Layer Design Continuum



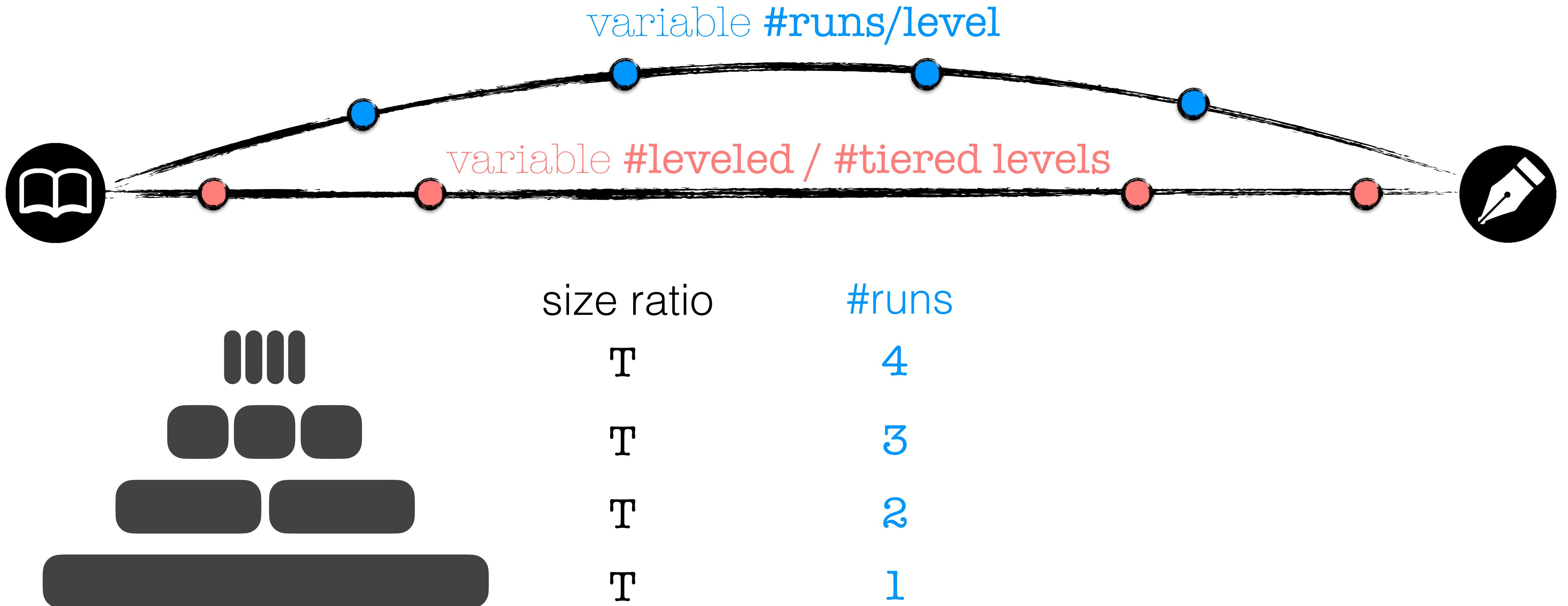
Storage Layer Design Continuum



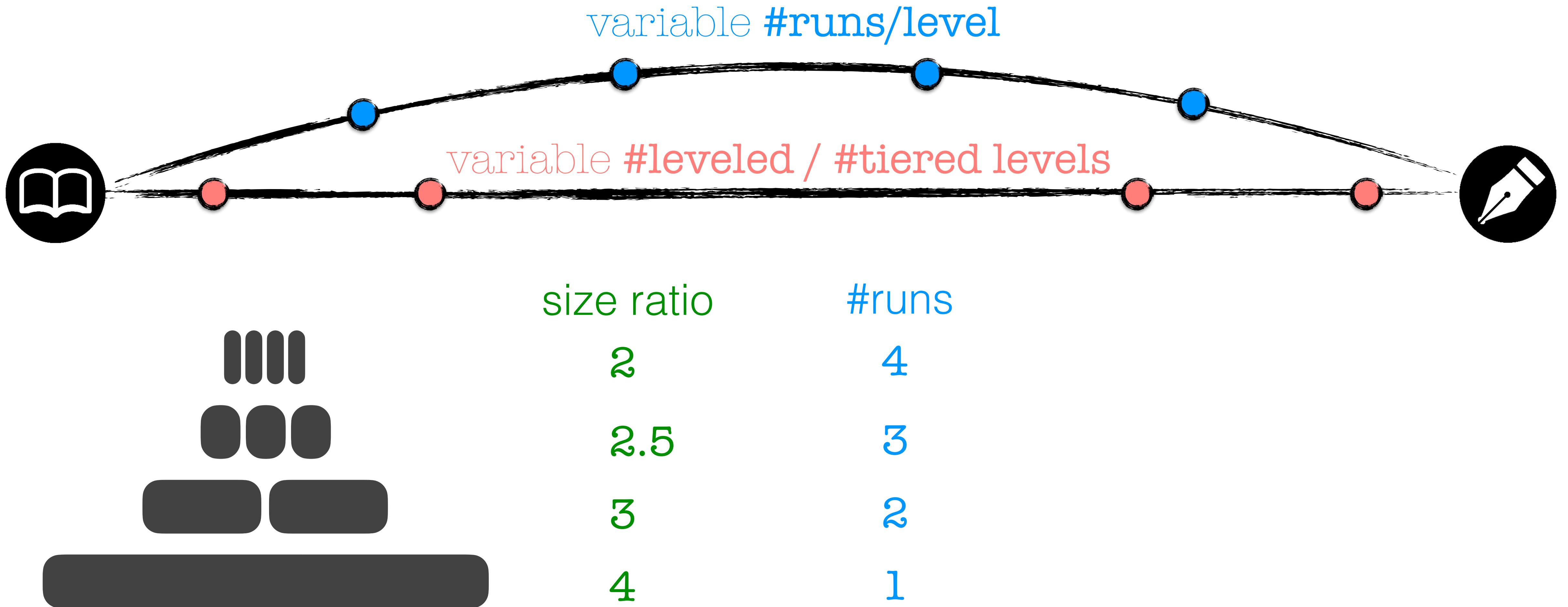
Storage Layer Design Continuum



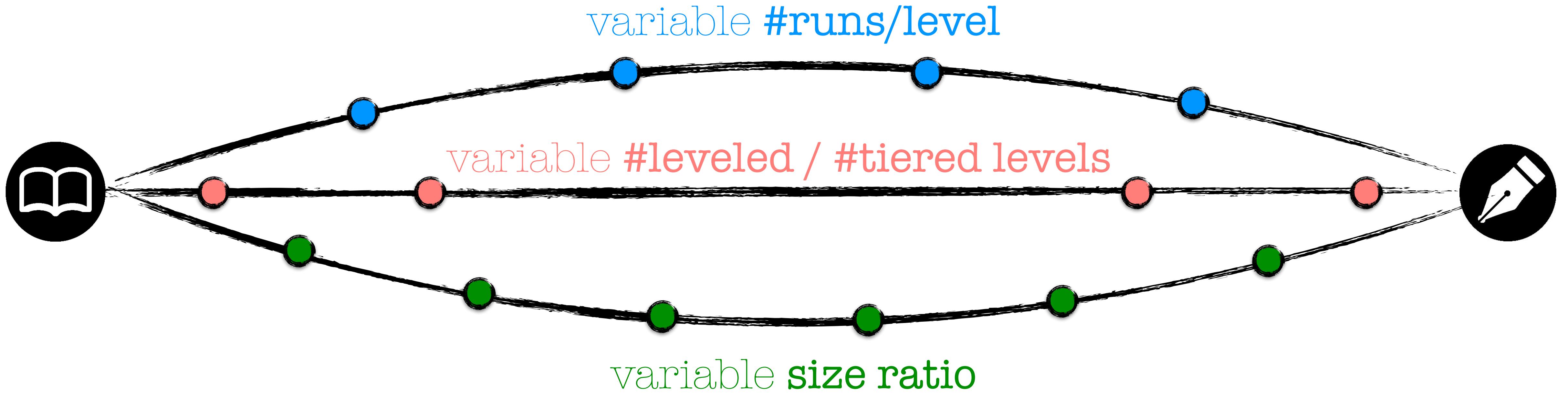
Storage Layer Design Continuum



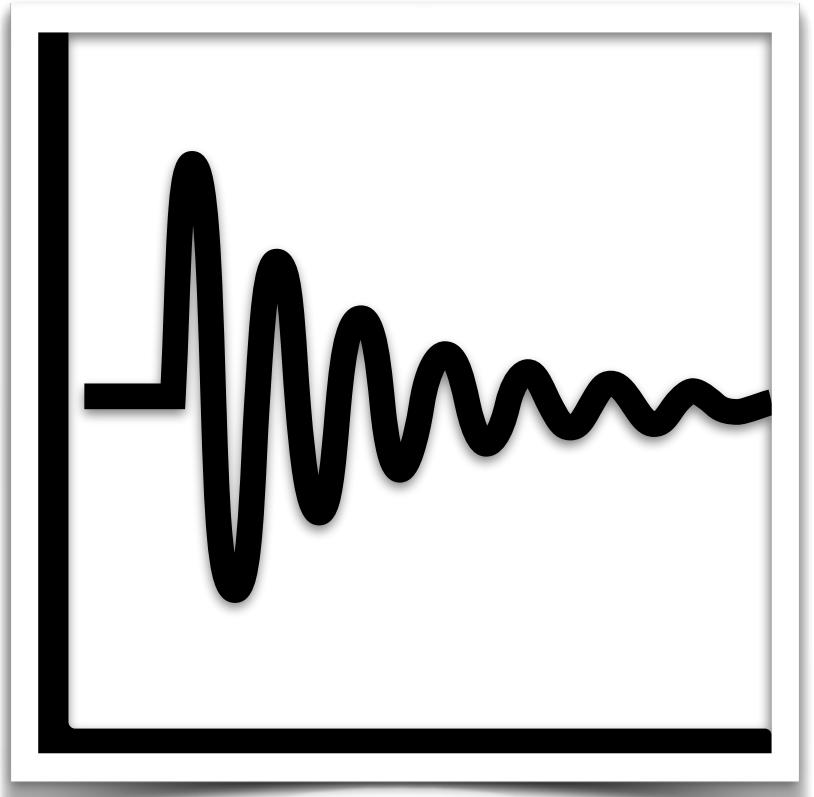
Storage Layer Design Continuum



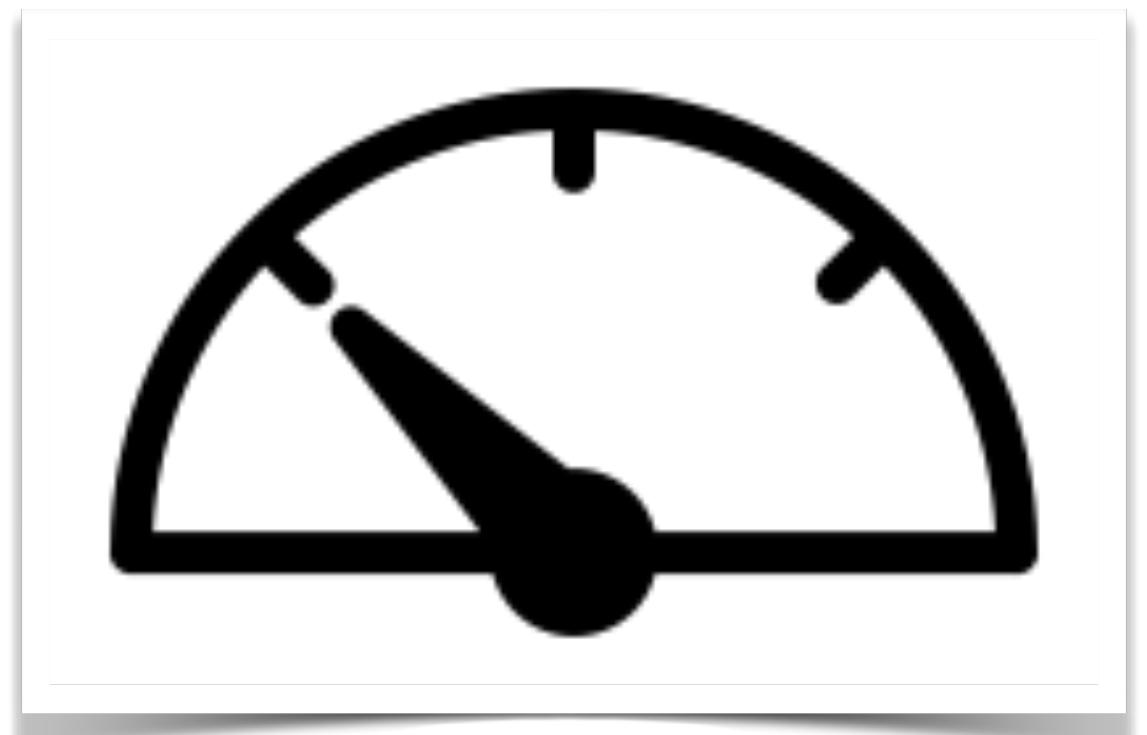
Storage Layer Design Continuum



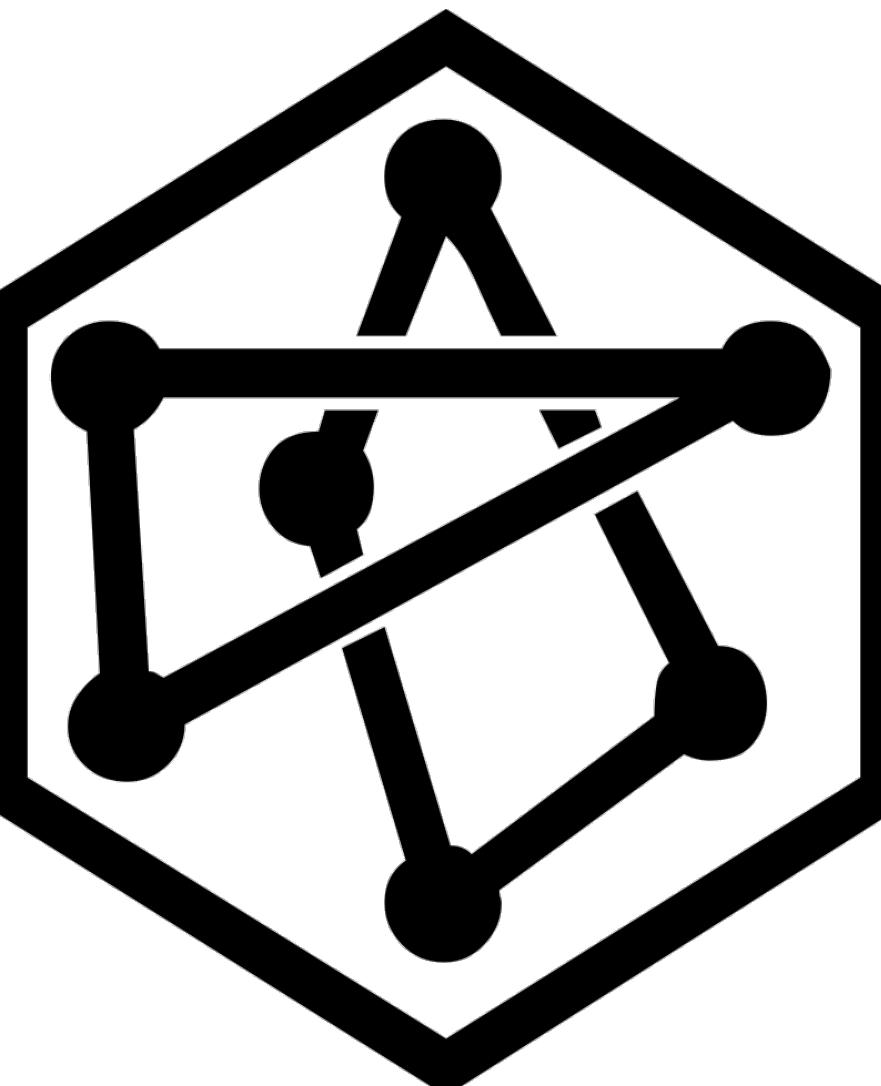
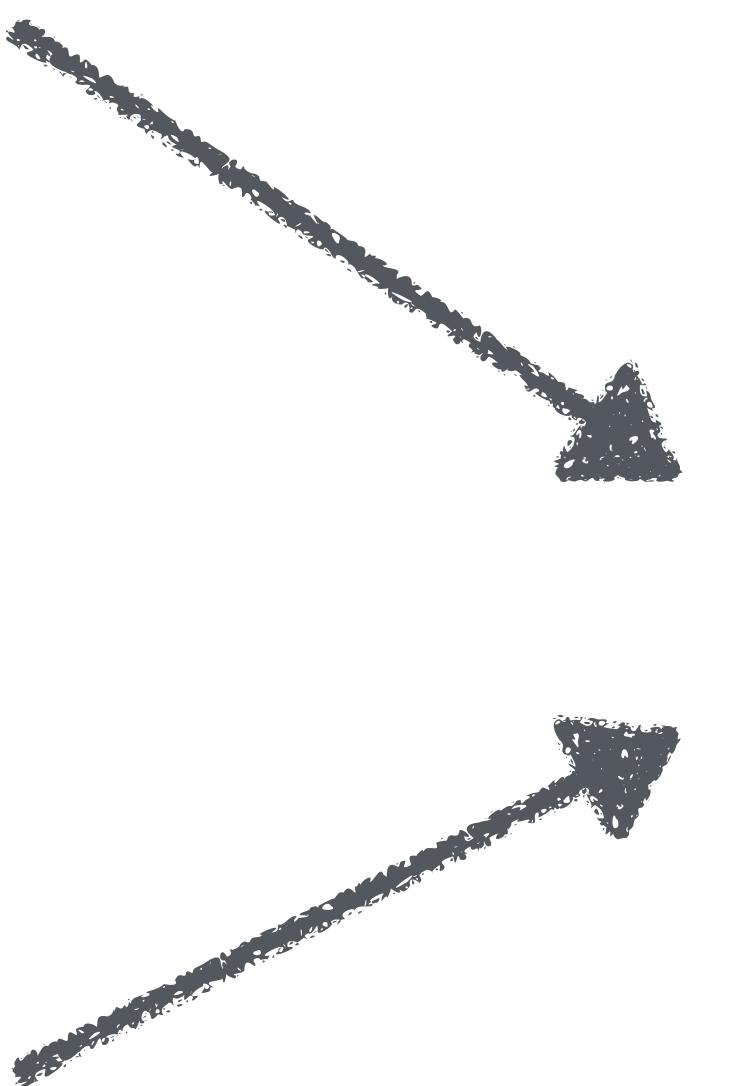
The LSM storage layer
design continuum



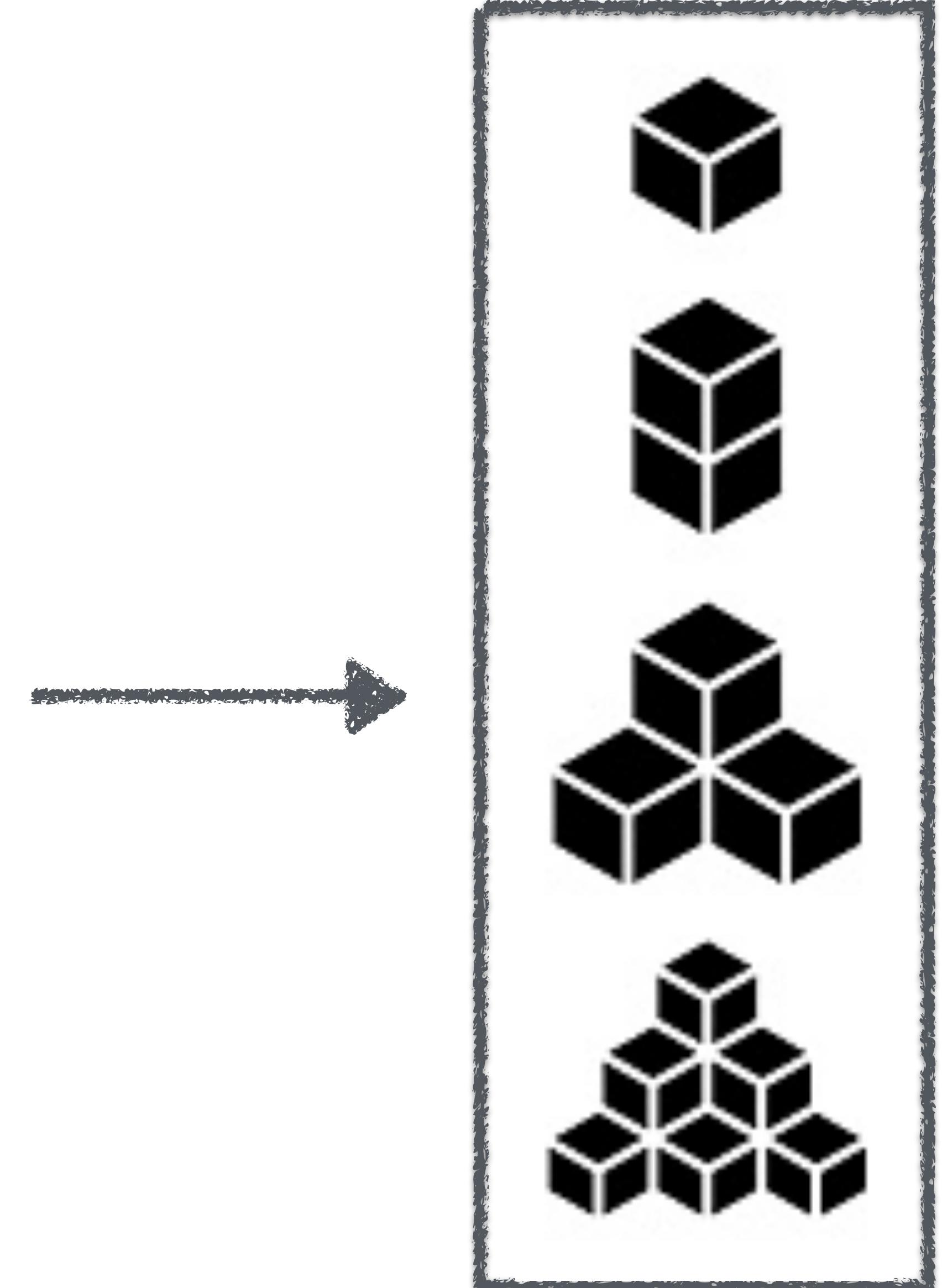
workload



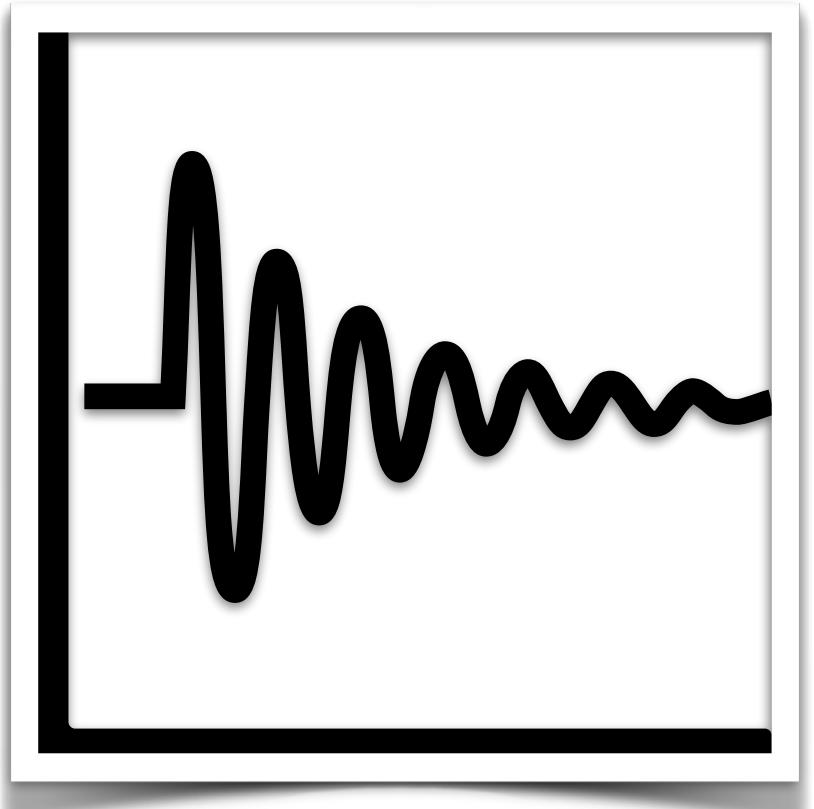
performance
target



performance
modeling



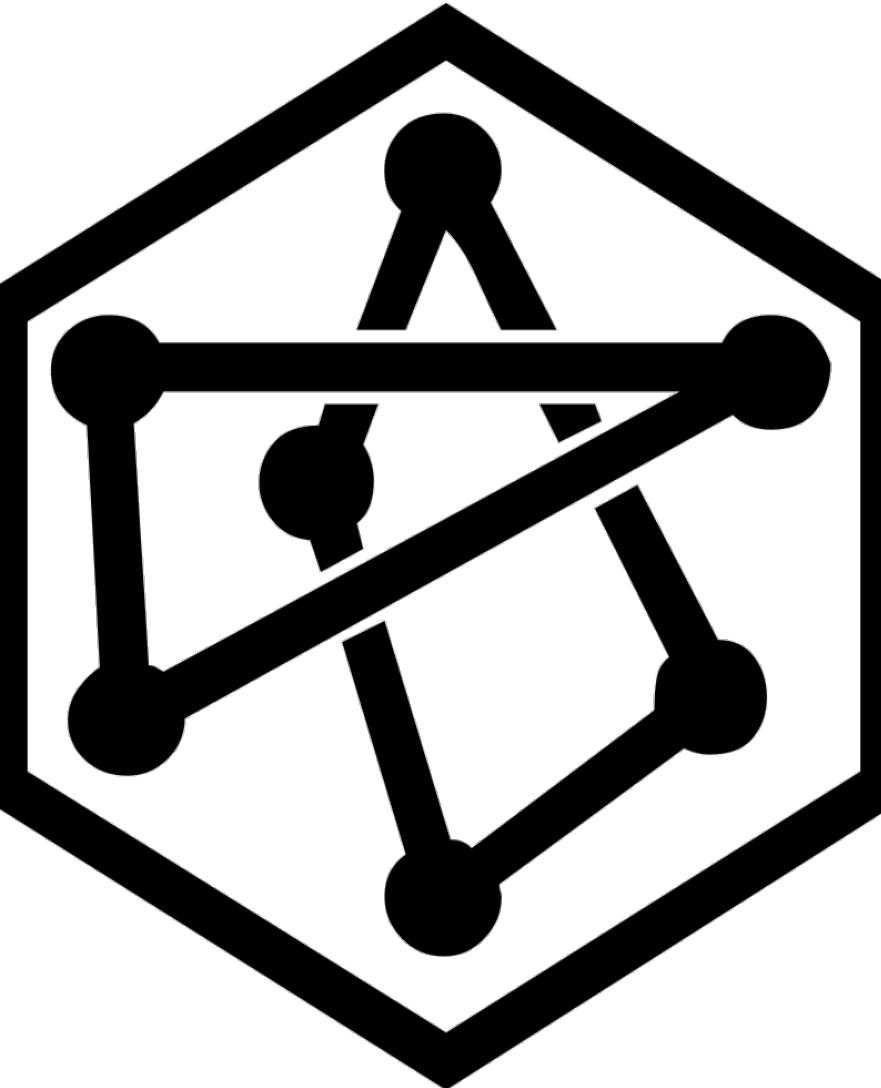
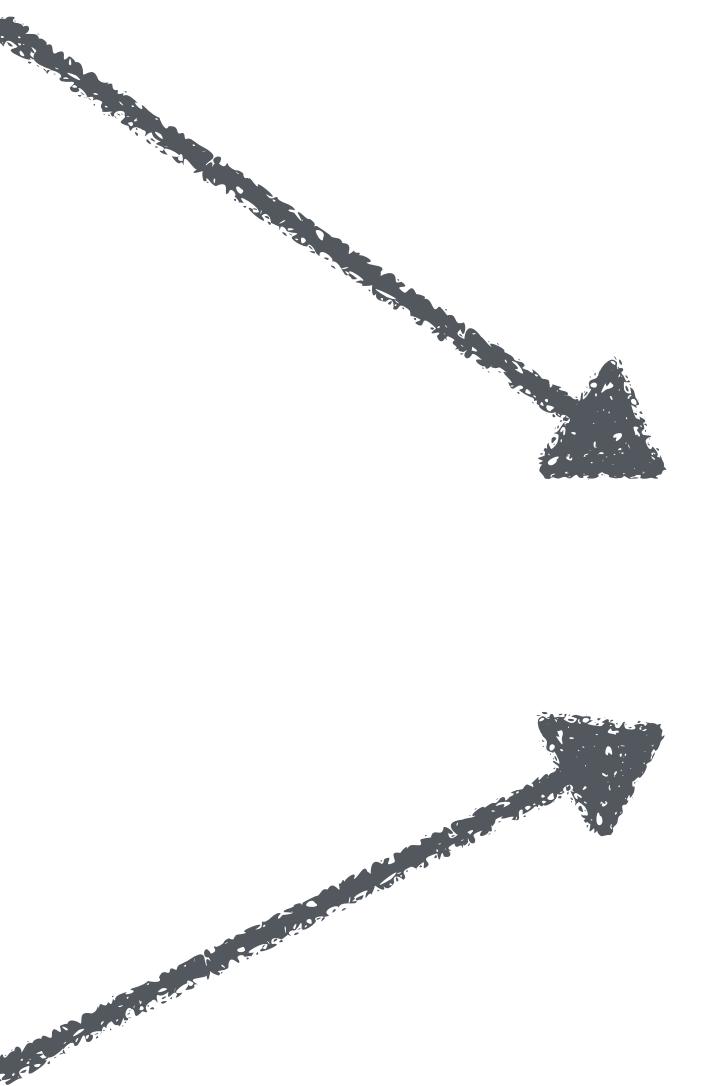
LSM designs



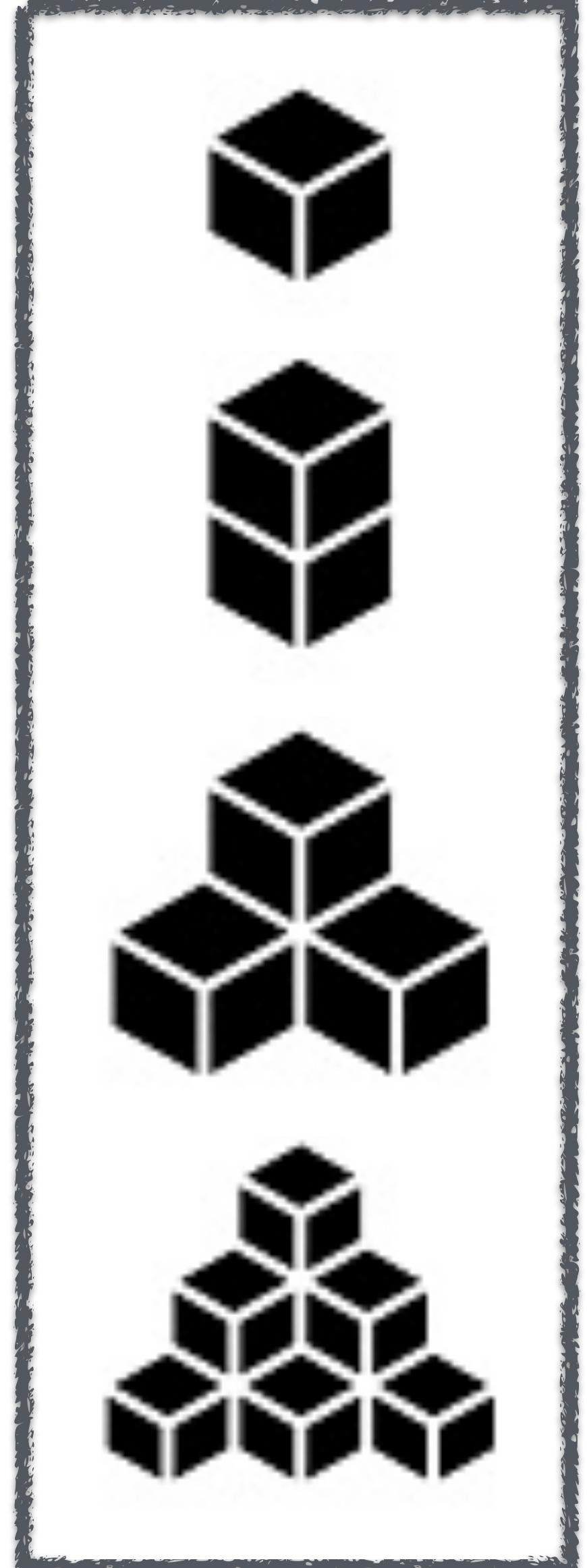
workload



performance
target



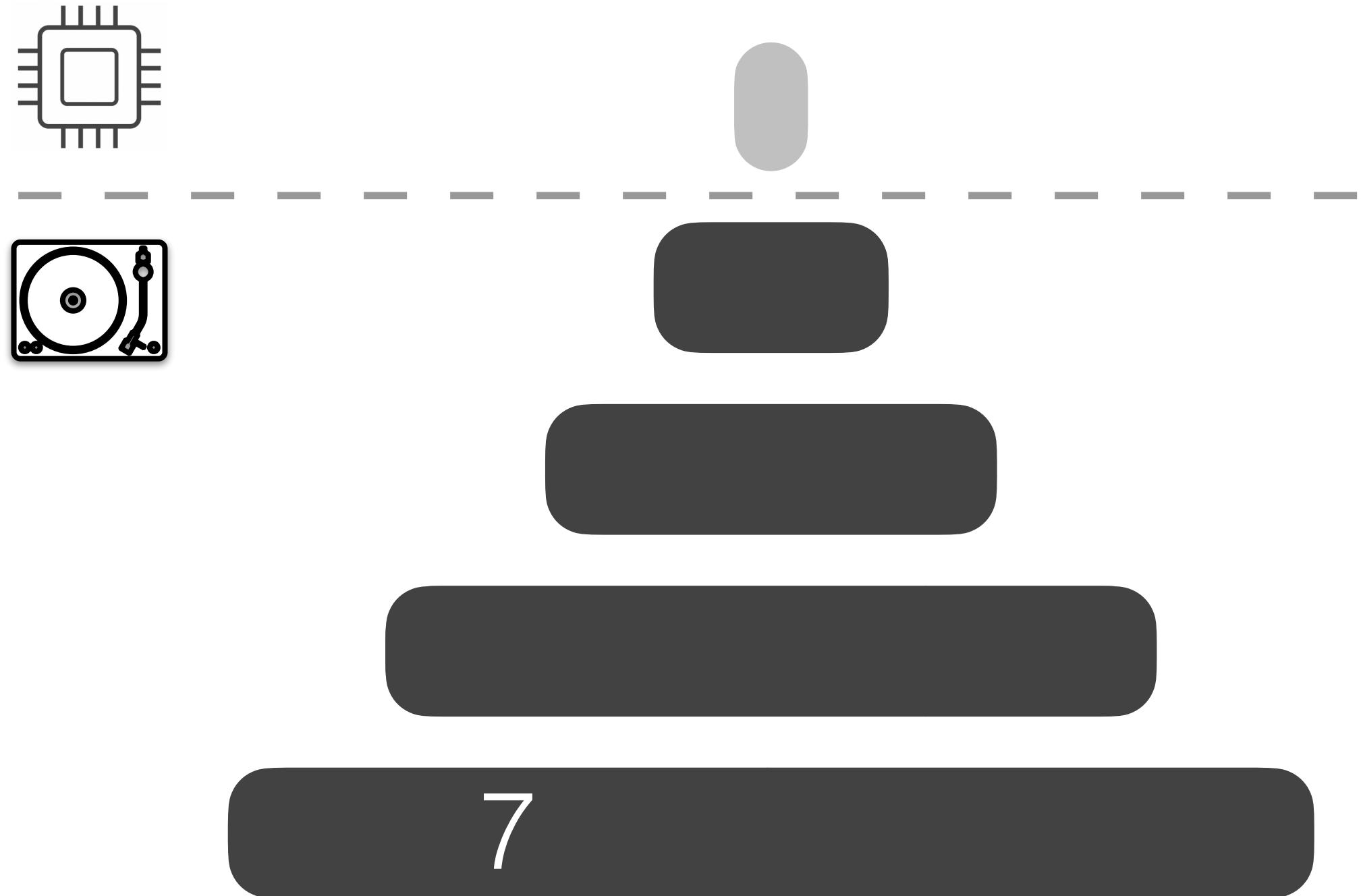
worst-case
performance
modeling



LSM designs

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

worst-case performance modeling



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

P : pages in buffer

B : entries/page

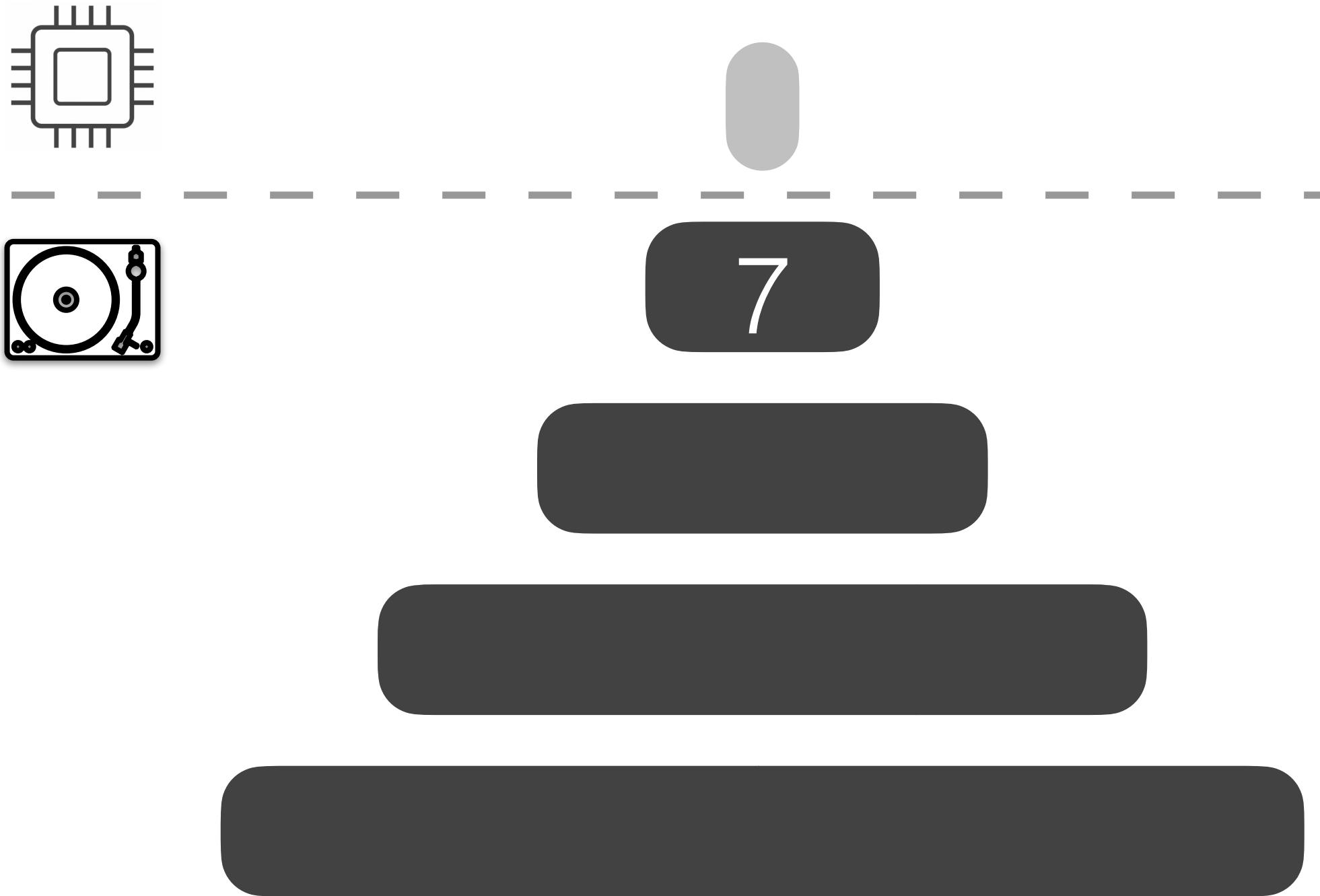
L : #levels

T : size ratio

N : #entries

ϕ : FPR of BF

worst-case performance modeling

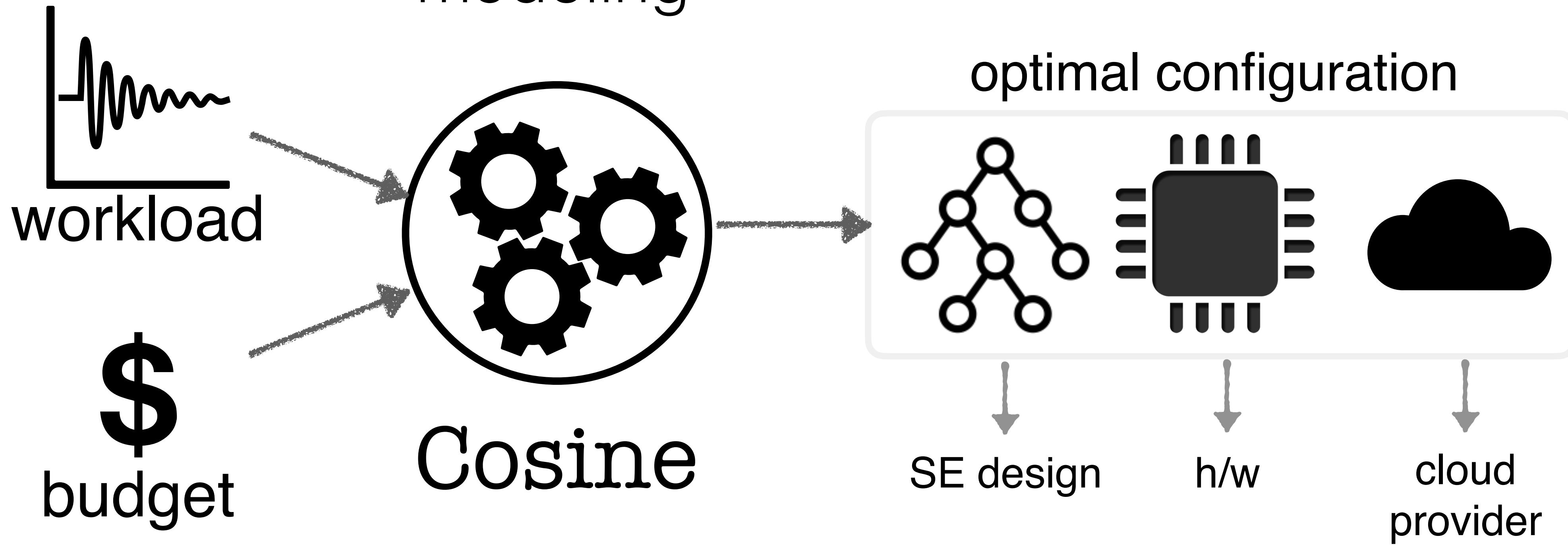


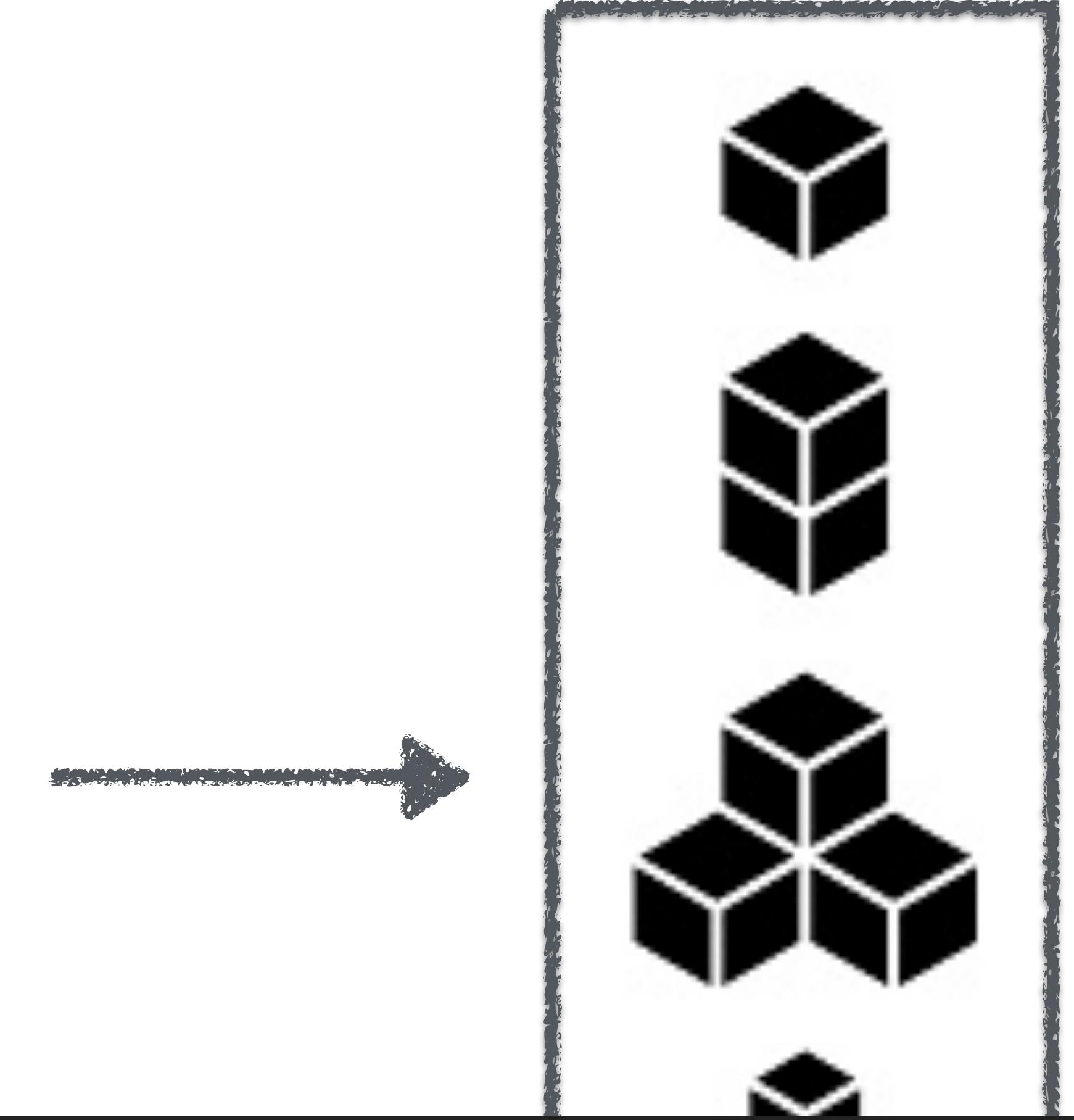
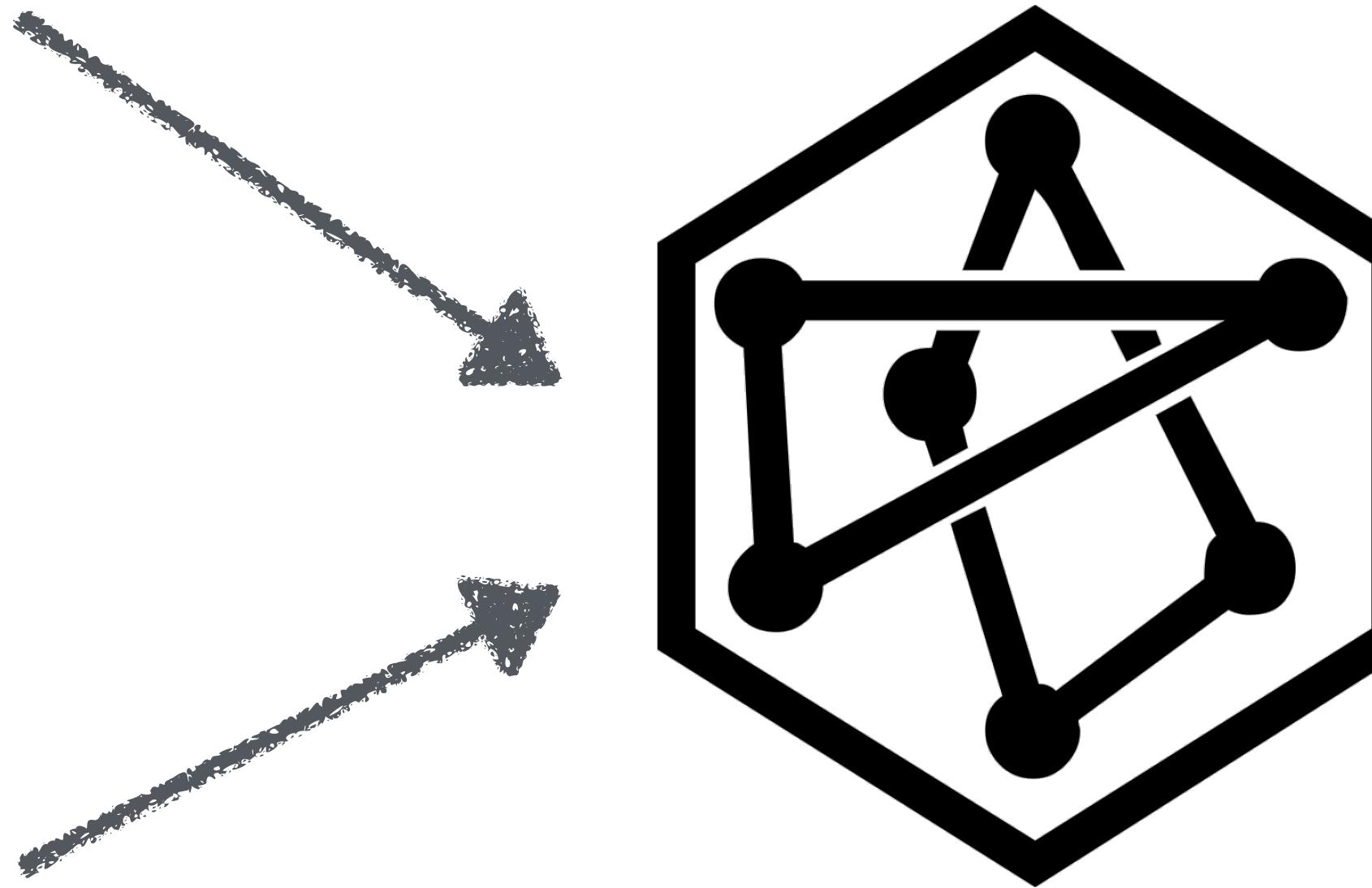
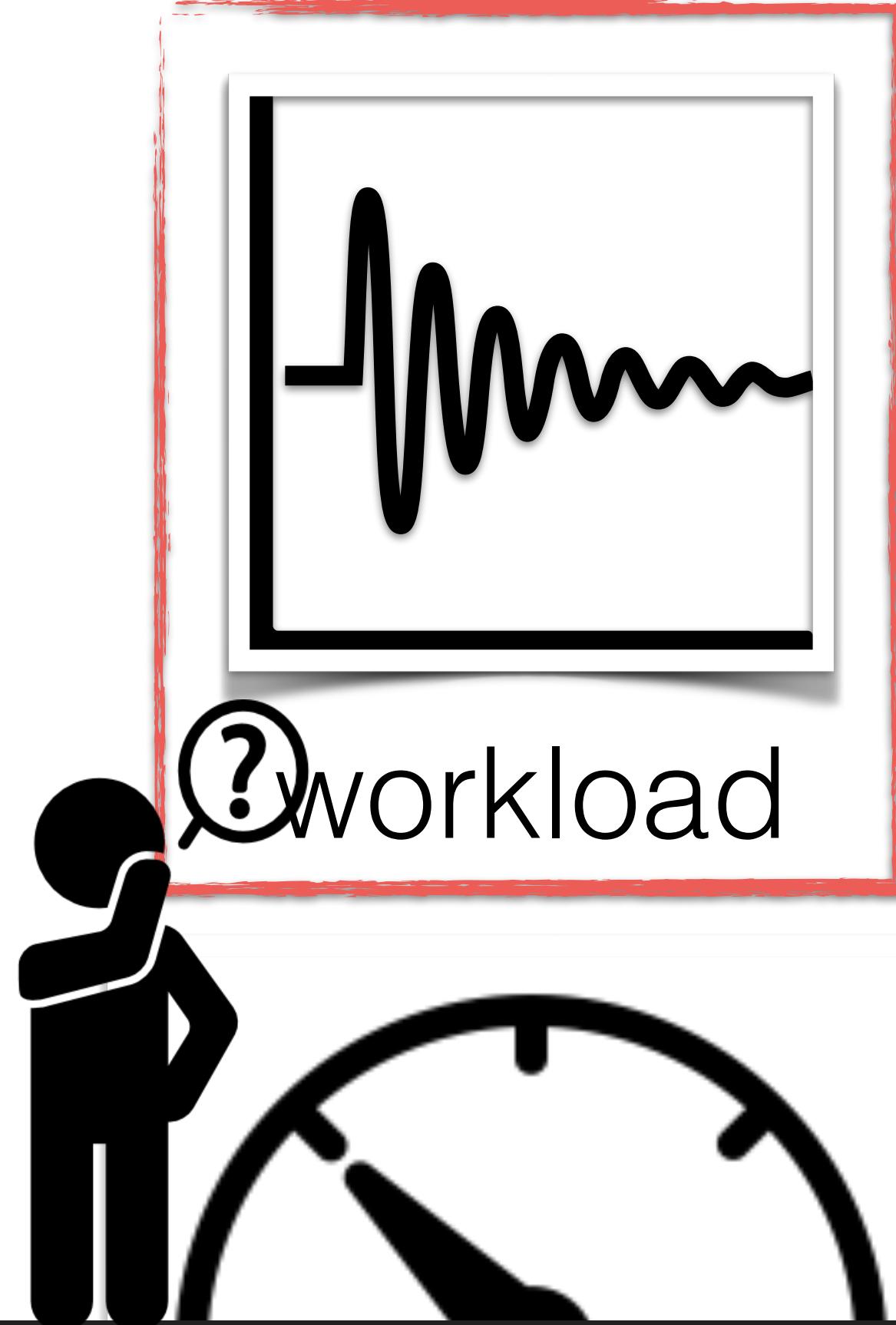
worst-case read cost: $1 + \sum_{i=1}^{L-1} \phi_i$

average-case performance modeling

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

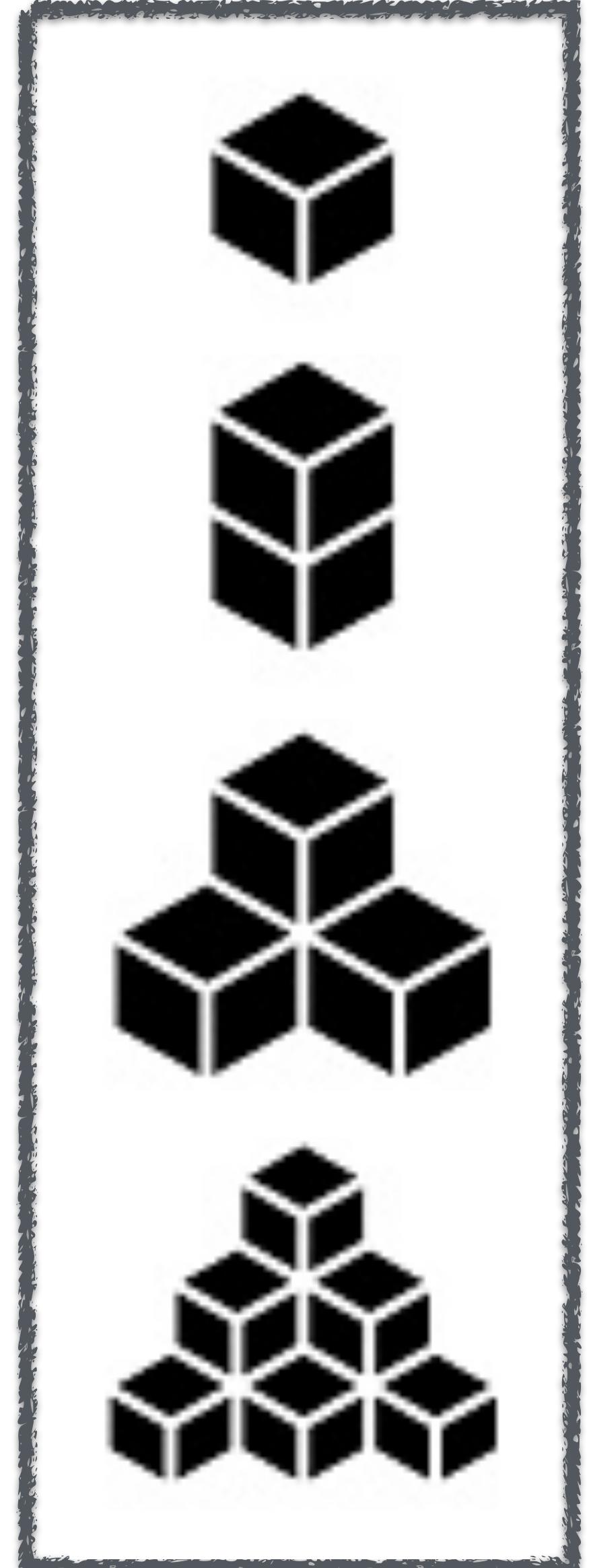
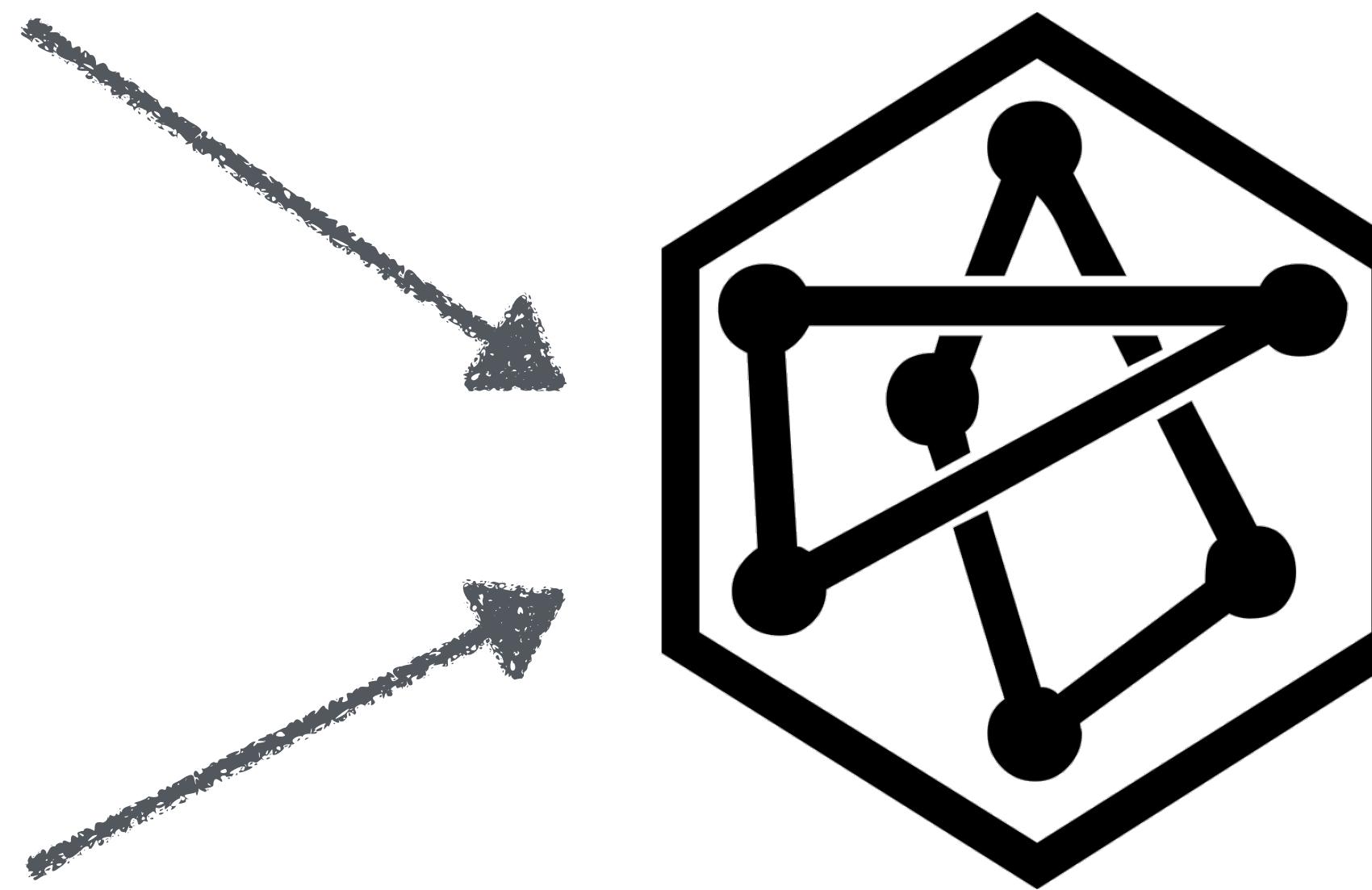
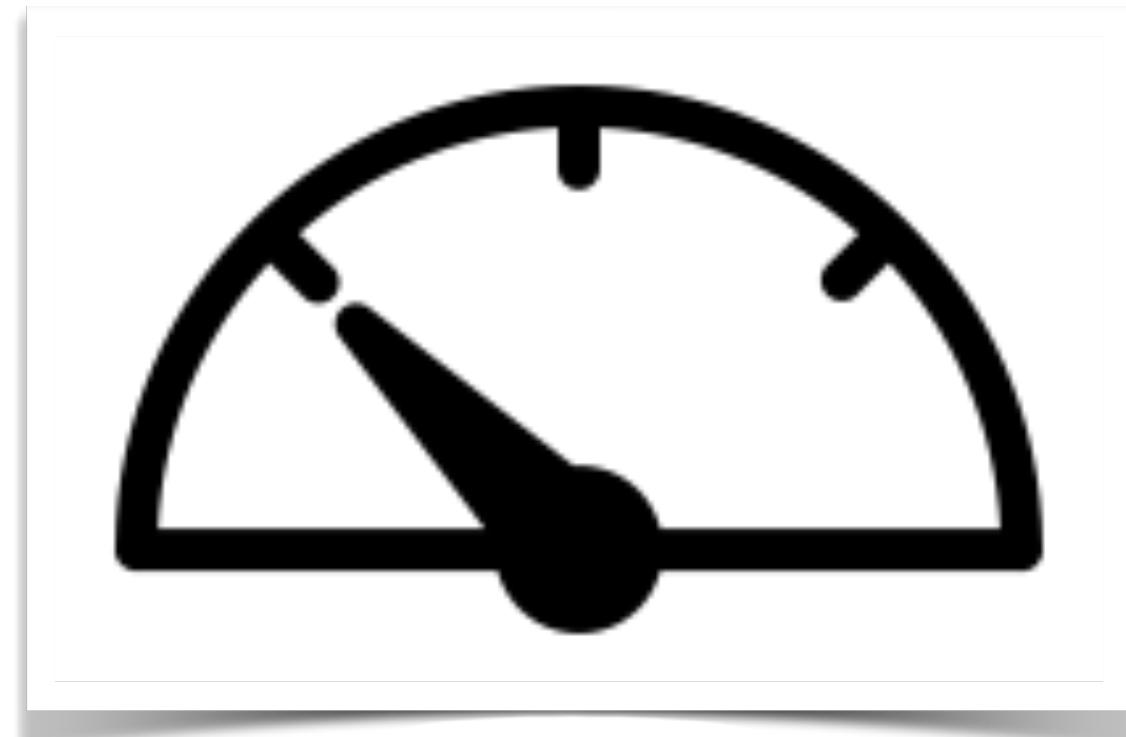
average-case performance modeling



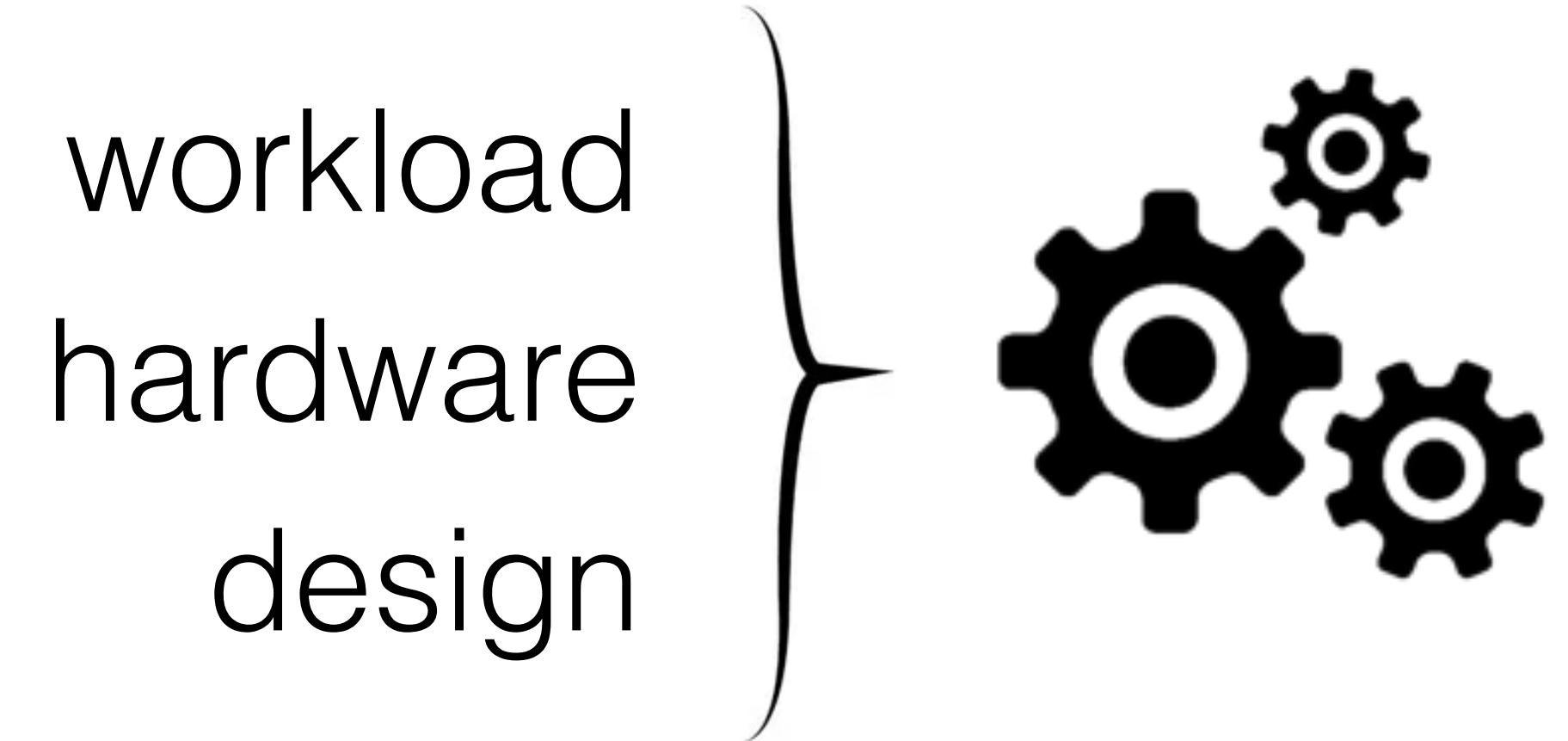


What if the workload changes?

LSM designs



Tuning



optimal
tuning

given workload \mathcal{W} , dataset \mathcal{D} , find a tuning \mathcal{T}^{opt} :

$$\mathcal{T}^{\text{opt}} = \operatorname{argmin}_{\mathcal{T}} (\text{cost}(\mathcal{W}, \mathcal{D}, \mathcal{T}))$$

Robust Tuning

unpredictability in

workload
hardware
design



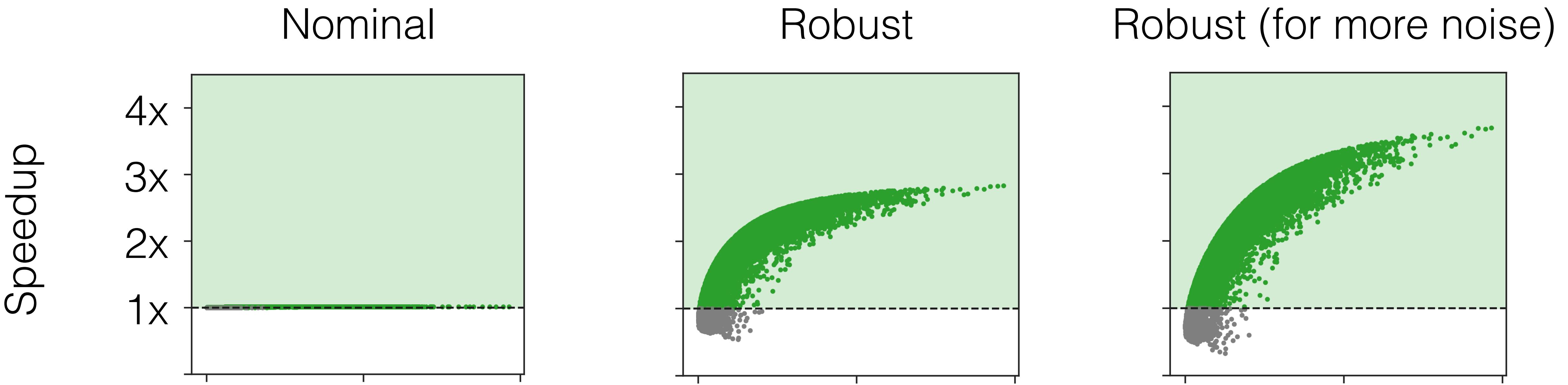
optimal
tuning

given workload \mathcal{W} , dataset \mathcal{D} , and a region $\mathcal{R}(\mathcal{W})$, find a tuning \mathcal{T}^{rob} :

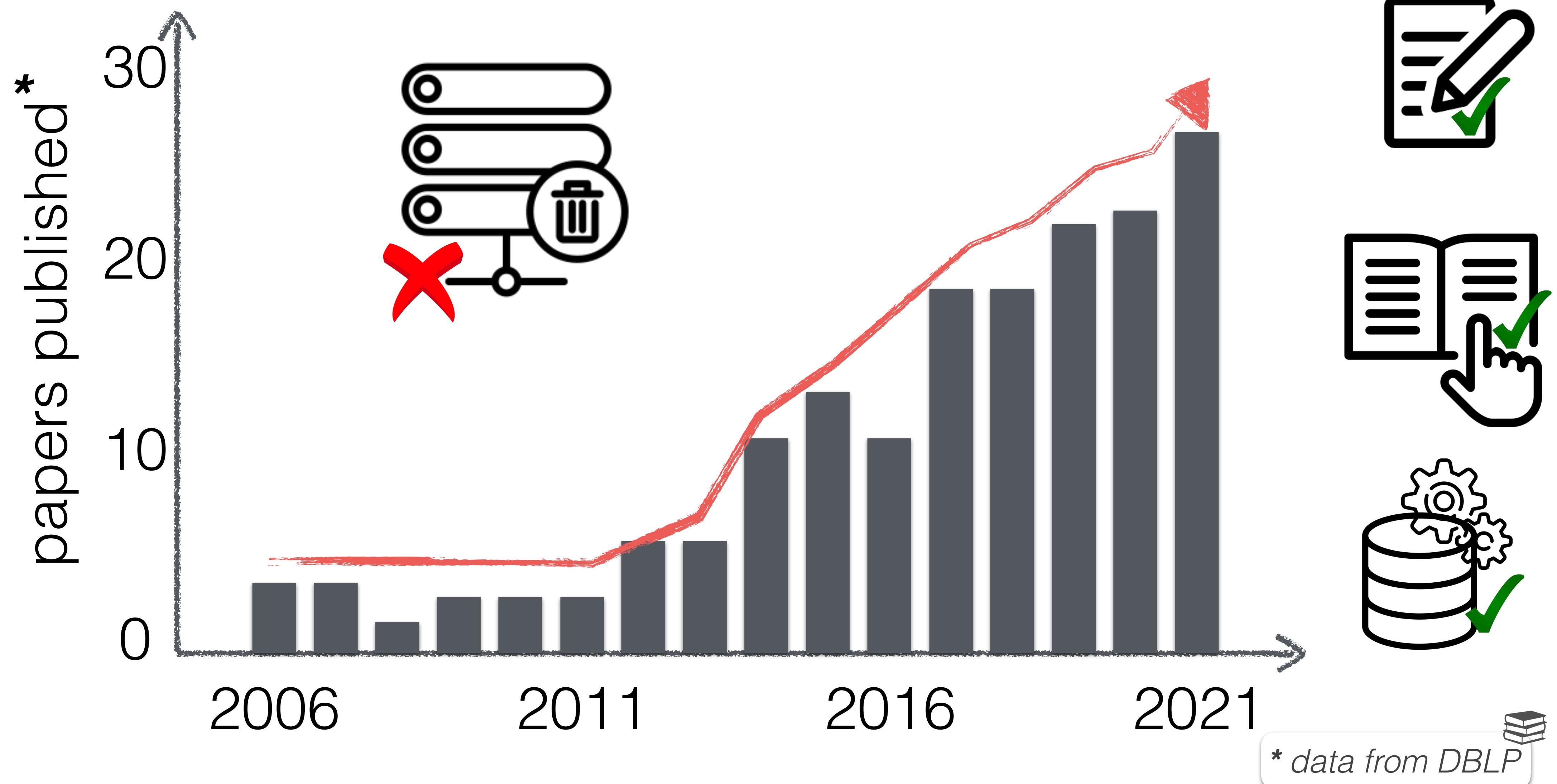
$$\mathcal{T}^{rob} = \operatorname{argmin}_{\mathcal{T}} \max_{\mathcal{W}' \in \mathcal{R}(\mathcal{W})} (\text{cost}(\mathcal{W}', \mathcal{D}, \mathcal{T}))$$



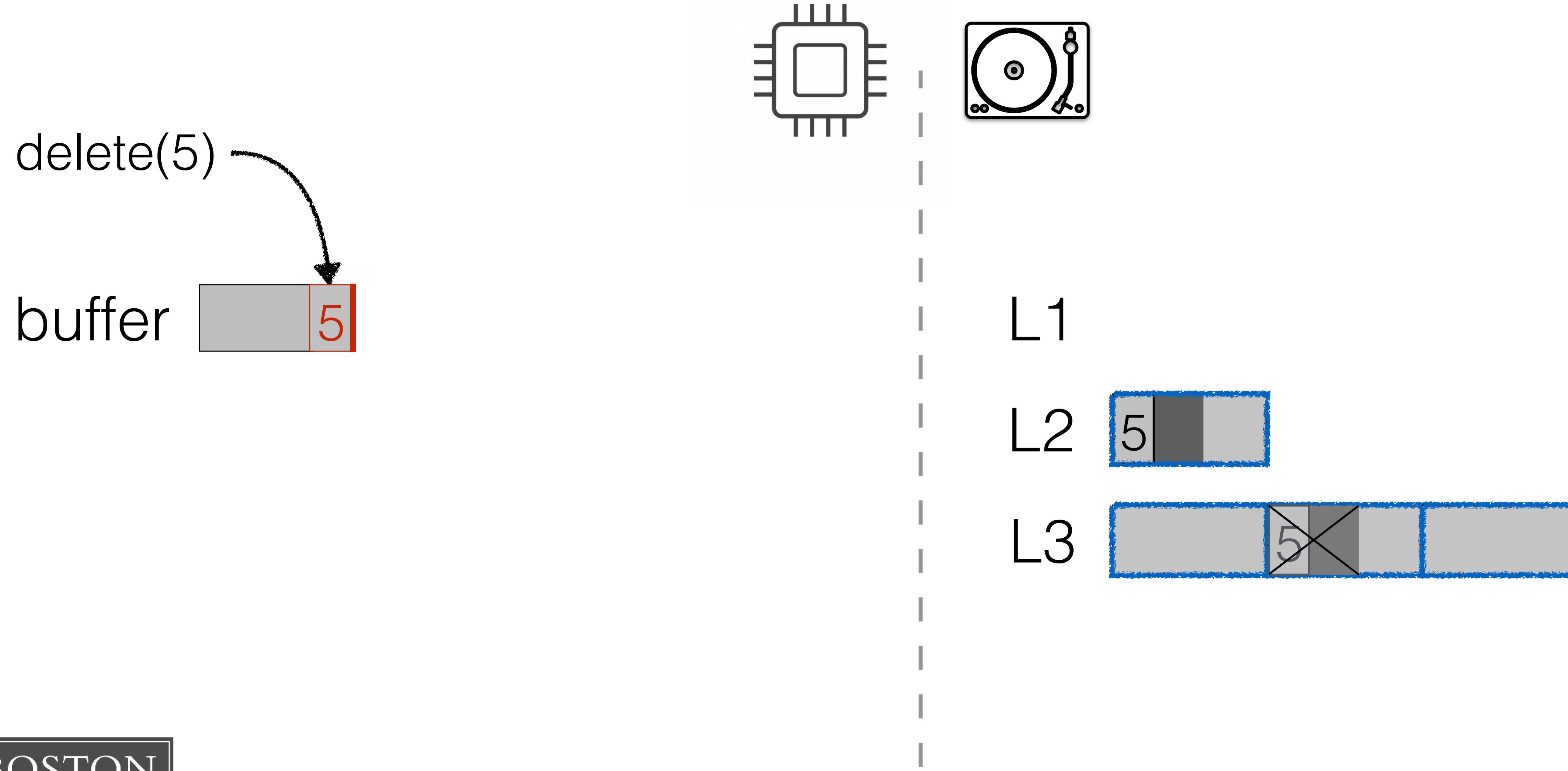
Robust Tuning



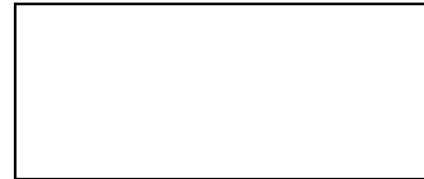
Research Trend: What about Deletes?

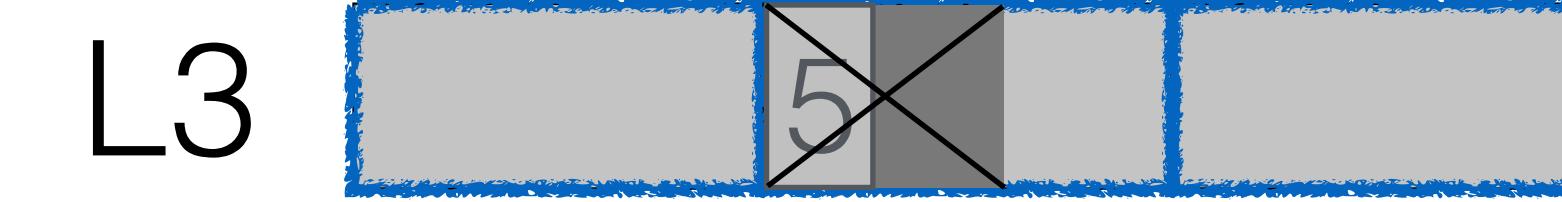
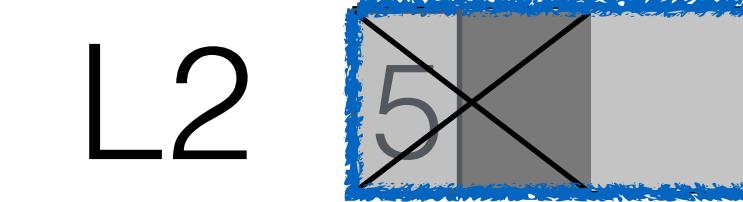
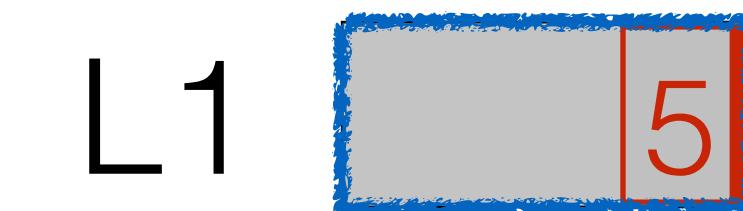
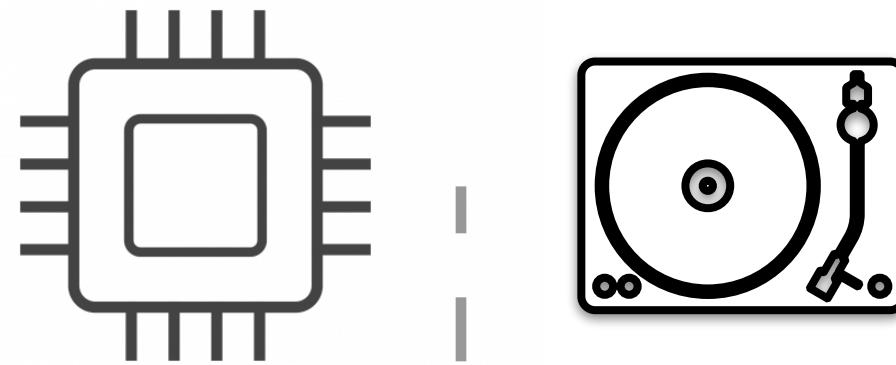


deletes in LSM-tree

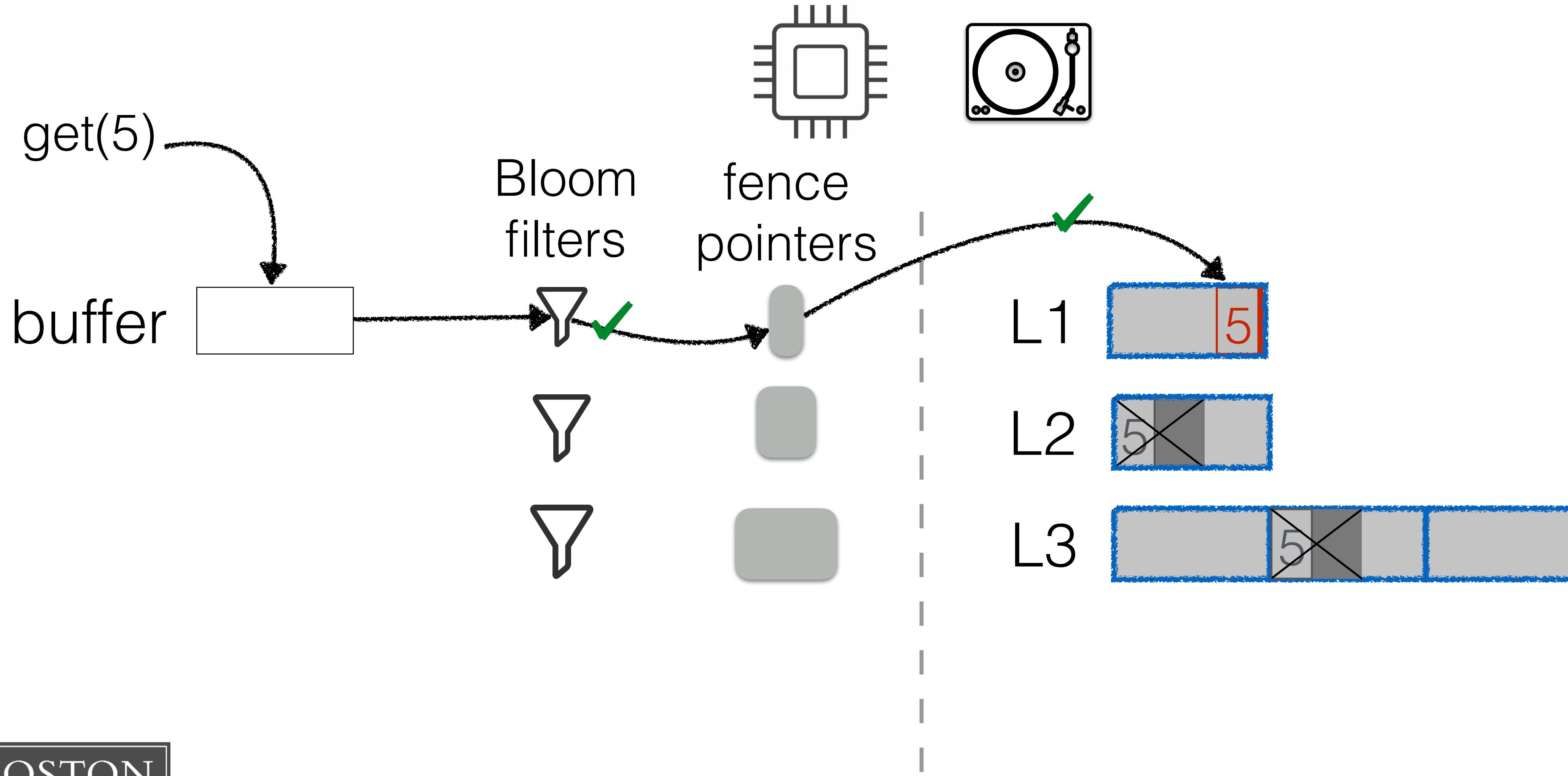


deletes in LSM-tree

buffer 



deletes in LSM-tree



Large-scale production

ZippyDB 

25.2M/day

UP2X 

92.5% merge through deletes

Internal DB ops

table drop 

data migration 

periodic cleanup 

Privacy



GDPR
(EU, UK) 



CCPA
(California) 



VCDPA
(Virginia) 

Facilitating Timely Deletes

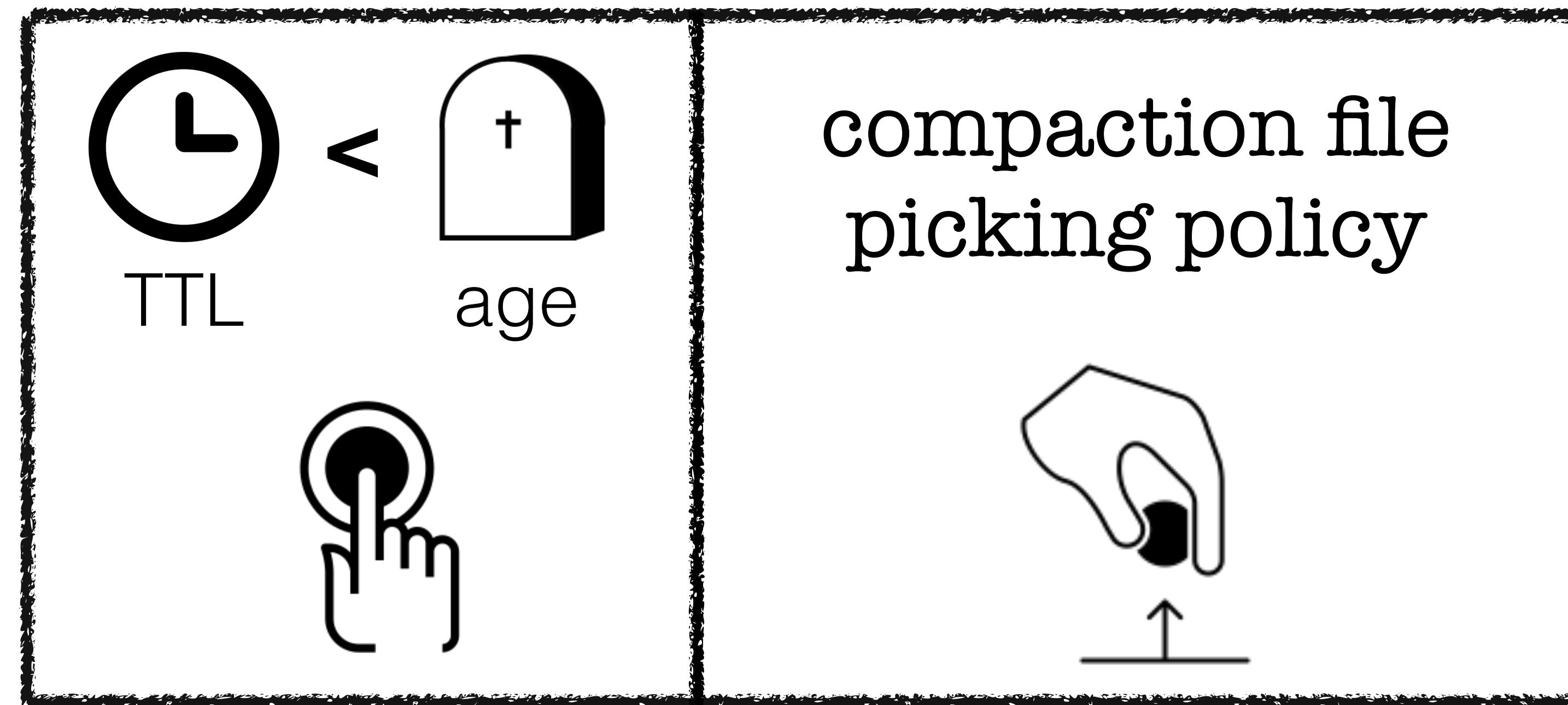
compaction
trigger



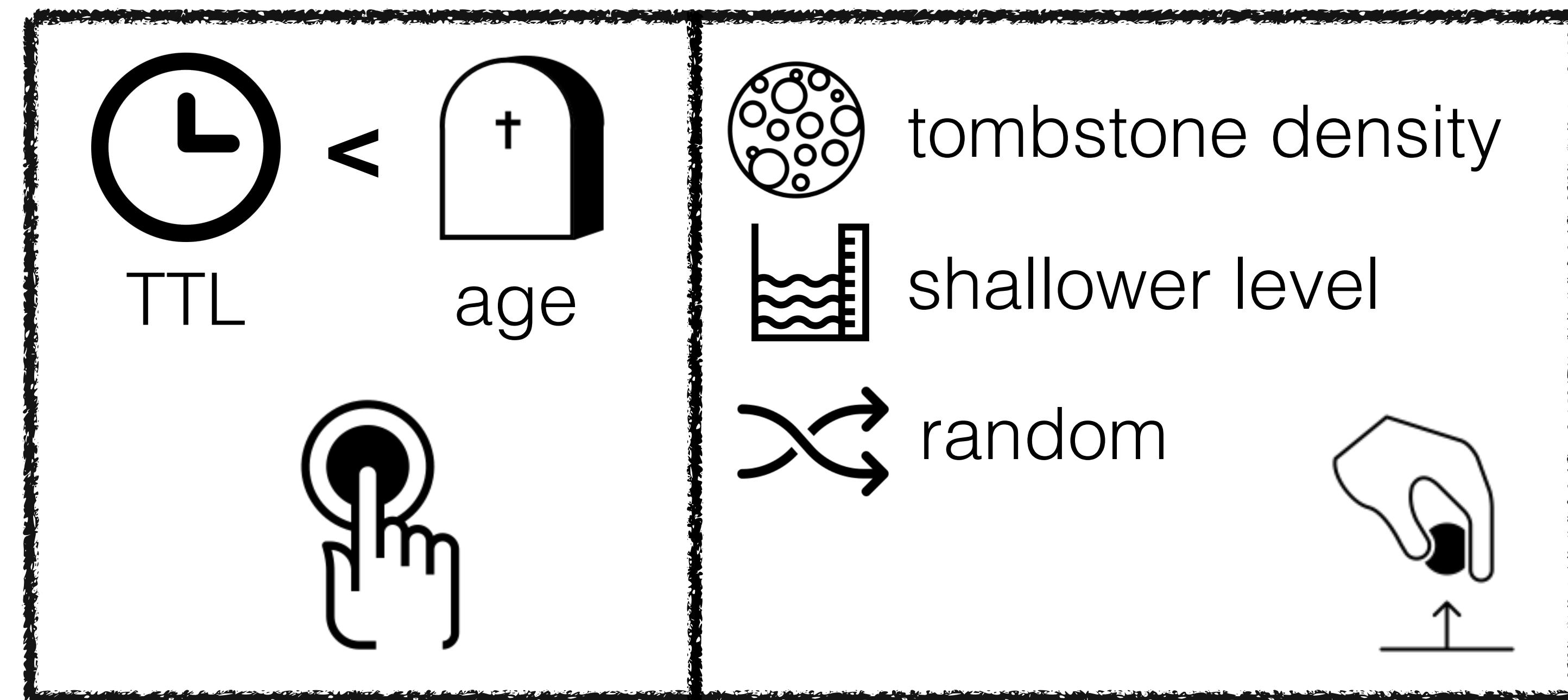
compaction file
picking policy



Facilitating Timely Deletes



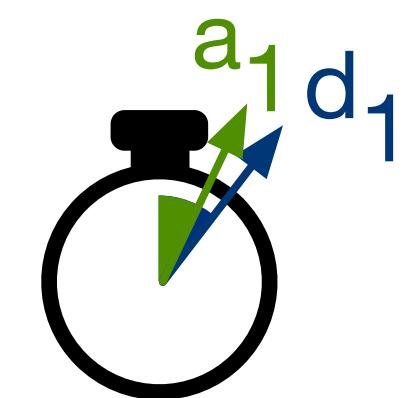
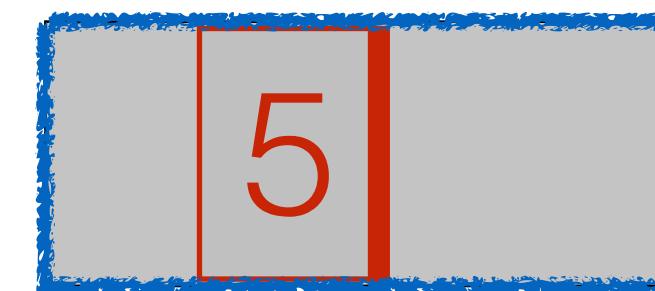
Facilitating Timely Deletes



Facilitating Timely Deletes

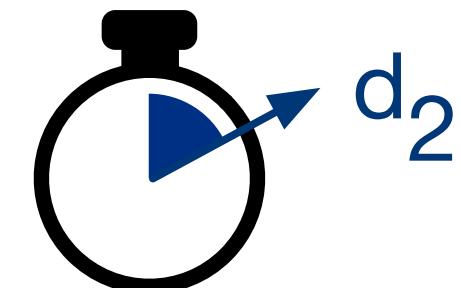
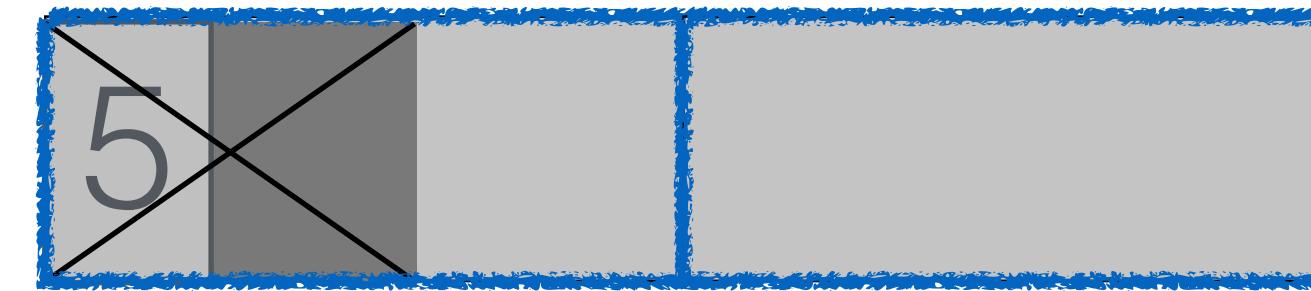
delete(5) within a threshold time: D_{th}

L1



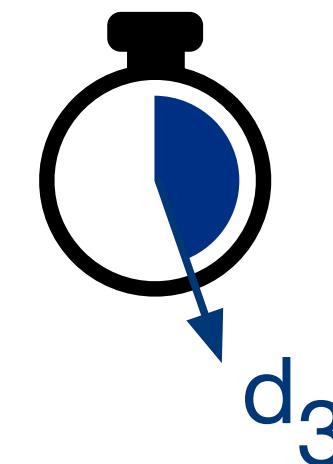
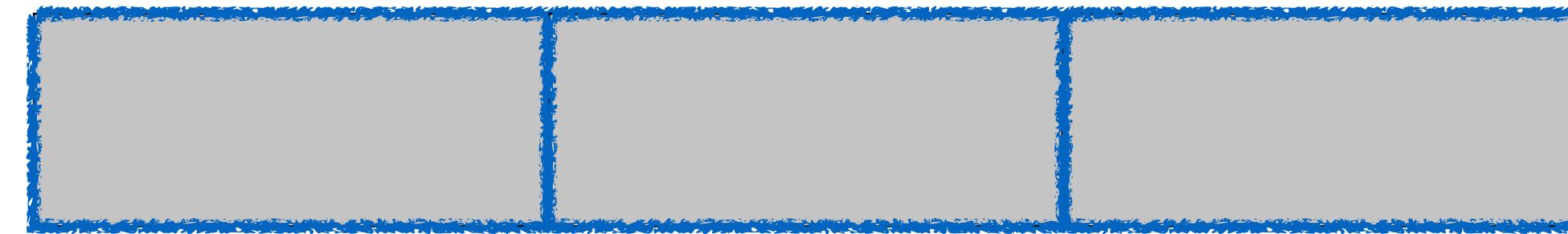
$$\sum_{i=1}^{L-1} d_i \leq D_{th}$$

L2



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

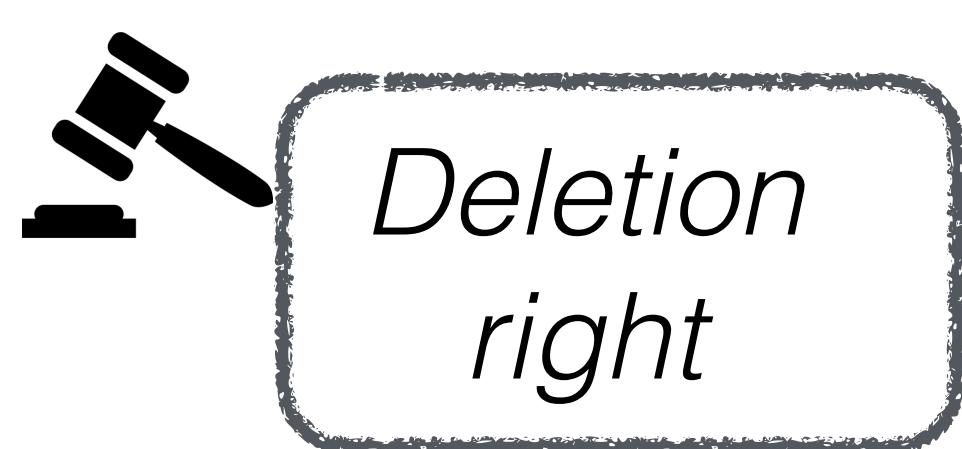
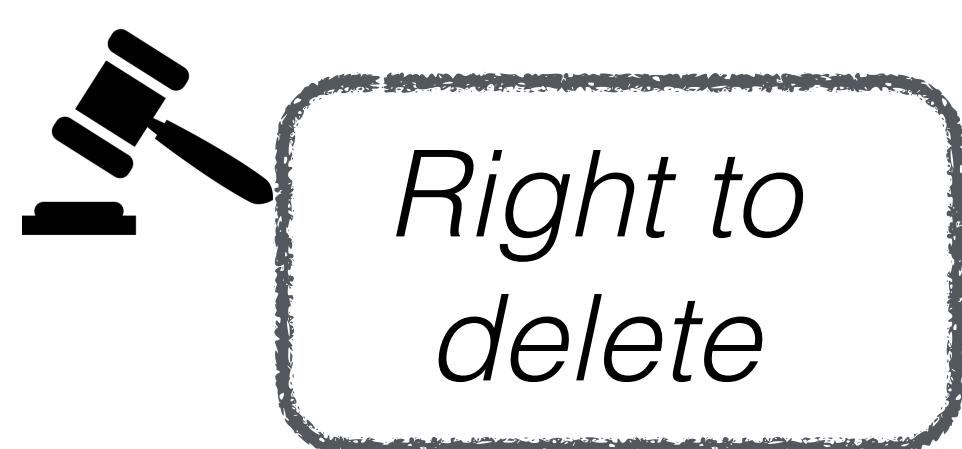
L3



L4



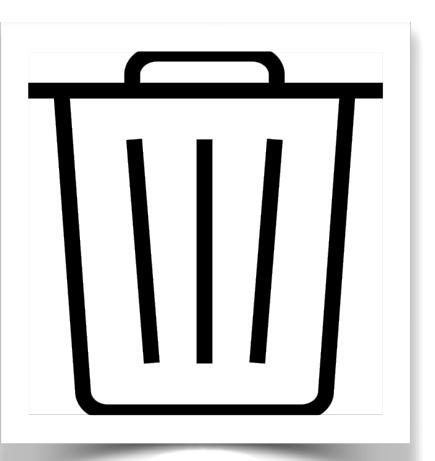
Policy layer



Requirements layer



Retention-based Deletes



On-demand Deletes

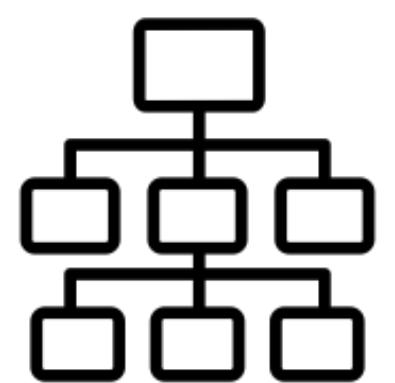
Application layer

```
CREATE TABLE R (...)  
WITH RET_DUR  
{ARBITRARY | FIXED(...)}  
WITH DPT  
{ARBITRARY | FIXED(...)};
```

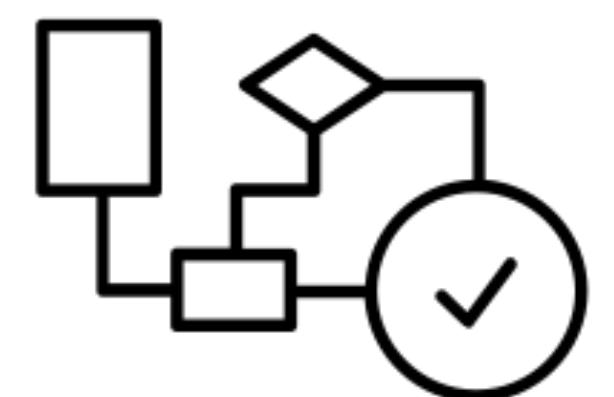
```
INSERT INTO R (...)  
WITH RET_DUR {<t>|t<i>};
```

```
DELETE FROM R  
WHERE (...)  
WITH DPT {<d>|d<i>};
```

System layer



Data layout re-organization



Data deletion algorithms

Policy layer



Right to be forgotten



Right to delete



Deletion right

Requirements layer



Retention-based Deletes



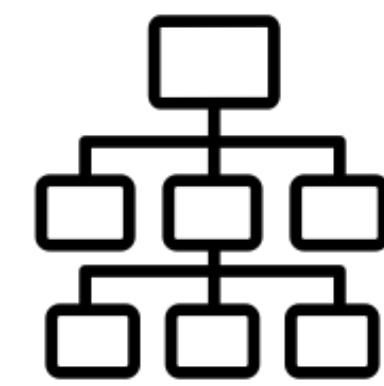
On-demand Deletes

Application layer

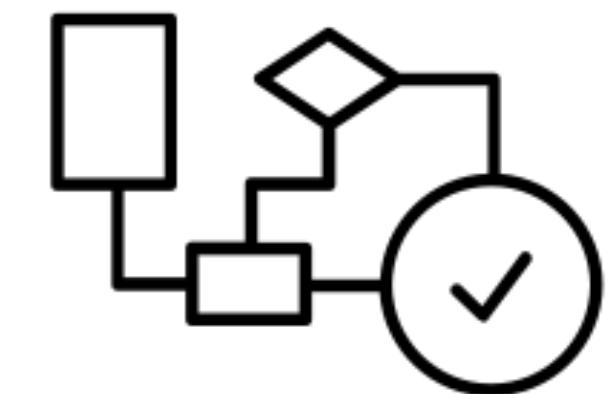
```
CREATE TABLE R (...)  
WITH RET_DUR  
{ARBITRARY | FIXED(...)}  
WITH DPT  
{ARBITRARY | FIXED(...)};
```

```
INSERT INTO R (...)  
WITH RET_DUR {<t>|t<i>};
```

```
DELETE FROM R  
WHERE (...)  
WITH DPT {<d>|d<i>};
```



Data layout
re-organization



Data deletion
algorithms



Policy layer



Right to be forgotten



Right to delete

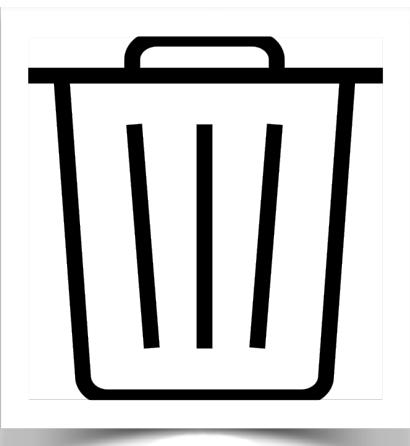


Deletion right

Requirements layer



Retention-based Deletes



On-demand Deletes

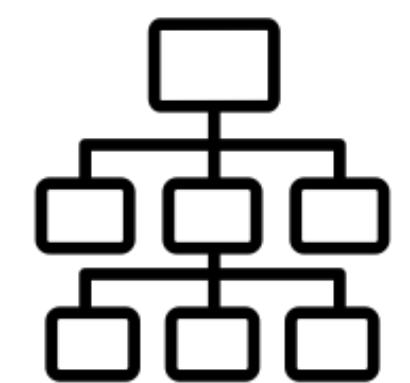
Application layer

```
CREATE TABLE R (...)  
WITH RET_DUR  
{ARBITRARY | FIXED(...)}  
WITH DPT  
{ARBITRARY | FIXED(...)};
```

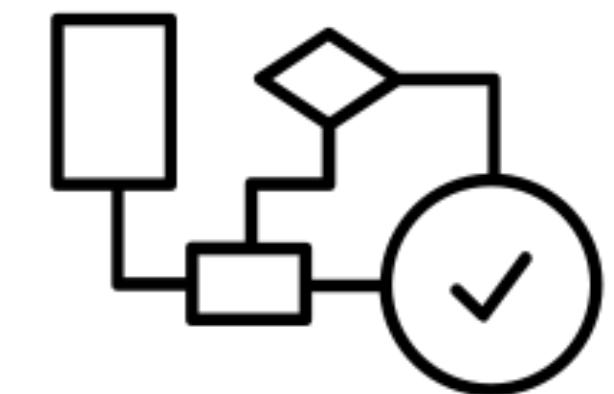
```
INSERT INTO R (...)  
WITH RET_DUR {<t>|t<i>};
```

```
DELETE FROM R  
WHERE (...)  
WITH DPT {<d>|d<i>};
```

System layer



Data layout re-organization



Data deletion algorithms



The Key Takeaways

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

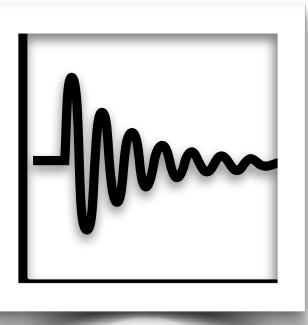
Compactions are key to building ingestion-friendly systems.

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

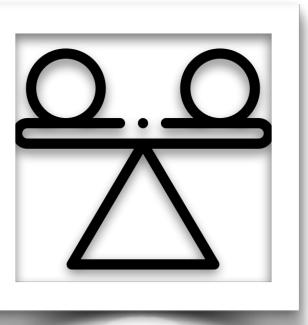
Open Research Challenges



Reduce write amplification



Workload-aware compactions & layout transformation



Performance Stability & Holistic Tuning



Automatic Tuning & Adaptive Behavior



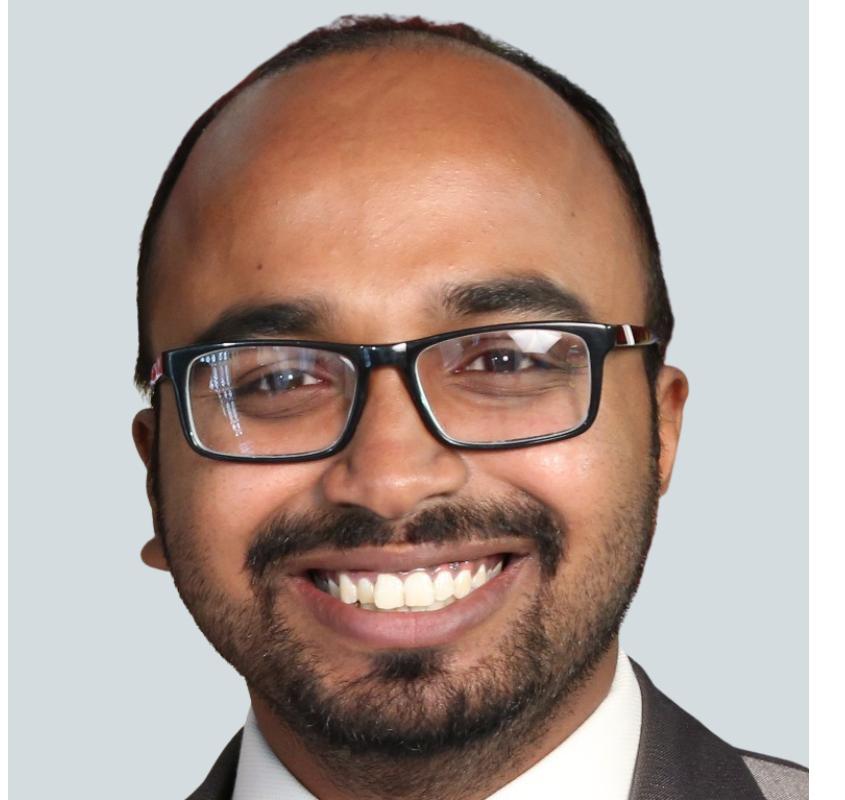
Privacy-aware LSM designs

Please see our manuscript for all references!

REFERENCES

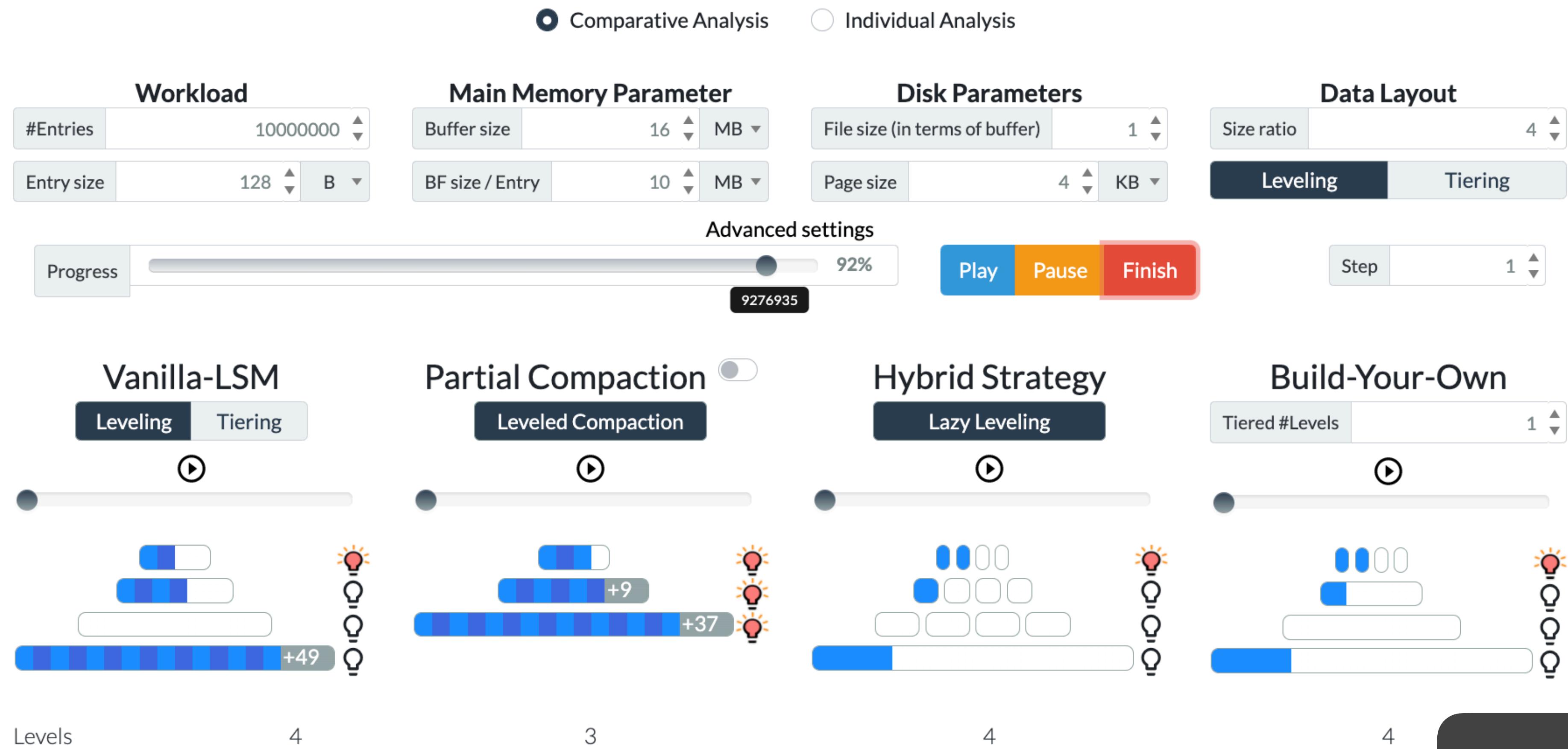
- [1] Regulation (EU) 2016/679 of the European Parliament and of the council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC. Official Journal of the European Union – L119/88, 2016.
- [2] California Consumer Privacy Act.
- [3] The California Privacy Protection Act.
- [4] Virginia Code of Publication Law.
- [5] H. Abu-Libdeh, S. Idreos, D. Li, Z. Zhu, and A. Ly. Designing Key-Value Stores for Self-Validating Data. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [6] W. Y. Alkov, S. Idreos, K. Kester, D. Guo, and D. Staratzis. Compaction Framework for Log-Structured Data Structures. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [7] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [8] Amazon. CloudWatch Metrics. <https://docs.aws.amazon.com/AmazonCloudWatch/latest/metrics-api/>.
- [9] Apache. Accumulo. <http://accumulo.apache.org>.
- [10] Apache. HDFS. <http://hadoop.apache.org/docs/r2.7.3/hdfs/HDFS.html>.
- [11] Apache. Cassandra. <http://cassandra.apache.org>.
- [12] M. Athanassoulis, S. Idreos, and D. Staratzis. Efficient Online Indexing for Log-Structured Key-Value Stores. In *Proceedings of the International Conference on Management of Data*, pages 217–228, 2012.
- [13] M. Athanassoulis, S. Idreos, and D. Staratzis. LSMT: Efficient Online Indexing for Log-Structured Key-Value Stores. In *Proceedings of the International Conference on Management of Data*, pages 217–228, 2012.
- [14] M. Athanassoulis, S. Idreos, and D. Staratzis. LSMT: Efficient Online Indexing for Log-Structured Key-Value Stores. In *Proceedings of the International Conference on Management of Data*, pages 217–228, 2012.
- [15] M. Athanassoulis, S. Idreos, and D. Staratzis. LSMT: Efficient Online Indexing for Log-Structured Key-Value Stores. In *Proceedings of the International Conference on Management of Data*, pages 217–228, 2012.
- [16] O. Balmau, S. Idreos, and M. Callaghan. LSMT: Efficient Online Indexing for Log-Structured Key-Value Stores. In *Proceedings of the International Conference on Management of Data*, pages 217–228, 2012.
- [17] S. Idreos and M. Callaghan. Key-Value Storage Engines. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [18] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [19] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [20] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [21] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [22] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [23] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [24] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [25] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [26] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [27] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuum and the Path Toward Self-Validating Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [28] Y. Chen, Y. Lu, F. Yang, Q. Wang, Y. Wang, and J. Shu. FlatStore: An Efficient Log-Structured Key-Value Storage Engine for Persistent Memory. In *ASPLOS '20: Architectural Support for Programming Languages and Operating Systems*, 2020.
- [29] C. Luo and M. J. Carey. LSM-based Storage Techniques: A Survey. *The VLDB Journal*, 29(1):393–418, 2020.
- [30] S. Luo, S. Chatterjee, R. Ketsetsidis, N. Dayan, W. Qin, and S. Idreos. Rosetta: A Robust Space-Time Optimized Range Filter for Key-Value Stores. In *Proceedings of the VLDB Endowment*, 13(9):1388–1399, 2020.
- [31] X. Wu, Y. Xu, Z. Shao, and S. Jiang. LSM-trie: An LSM-tree-based Ultra-Large Key-Value Store for Small Data Items. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 71–82, 2015.
- [32] L. Yang, H. Wu, T. Zhang, X. Cheng, F. Li, L. Zou, Y. Wang, R. Chen, J. Wang, and G. Huang. Leaper: A Learned Prefetcher for Cache Invalidations in LSM-tree based Storage Engines. *Proceedings of the VLDB Endowment*, 13(11):1976–1989, 2020.
- [33] T. Yao, J. Wan, P. Huang, X. He, Q. Gui, F. Wu, and C. Xie. A Light-weight Compaction Tree to Reduce I/O Amplification toward Efficient Key-Value Stores. In *Proceedings of the IEEE Symposium on Mass Storage Systems and Technologies (MSST)*, 2017.
- [34] T. Yao, J. Wan, P. Huang, X. He, Q. Gui, F. Wu, and C. Xie. Building Efficient Key-Value Stores via a Lightweight Compaction Tree. *ACM Transactions on Storage (TOS)*, 13(4):29:1–29:28, 2017.
- [35] H. Zhang, H. Lim, V. Leis, D. G. Andersen, M. Kaminsky, K. Keeton, and A. Pavlo. SuRF: Practical Range Query Filtering with Fast Succinct Tries. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 323–336, 2018.
- [36] H. Zhang, H. Lim, V. Leis, D. G. Andersen, M. Kaminsky, K. Keeton, and A. Pavlo. Succinct Range Filters. *ACM Transactions on Database Systems (TODS)*, 45(2):5:1–5:31, 2020.
- [37] W. Zhang, Y. Xu, Y. Li, and D. Li. Improving Write Performance of LSMT-Based Key-Value Store. In *22nd IEEE International Conference on Parallel and Distributed Systems, ICPADS 2016, Wuhan, China, December 13–16, 2016*, pages 553–560, 2016.
- [38] W. Zhang, X. Zhao, S. Jiang, and H. Jiang. ChameleonDB: a key-value store for optane persistent memory. In *EuroSys '21: Sixteenth European Conference on Computer Systems*, 2021.

The DiSC Lab



disc.bu.edu

Don't Miss: Compactionary



disc-projects.bu.edu/compactionary

Tuesday
4-6PM

Dissecting, Designing, and Optimizing **LSM-Based Data Stores**

Subhadeep Sarkar

Manos Athanassoulis

Thank You!

Questions?