



Dissecting, Designing, and Optimizing
LSM-Based Data Stores

Subhadeep Sarkar

Manos Athanassoulis

BOSTON
UNIVERSITY

lab
DiSC

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

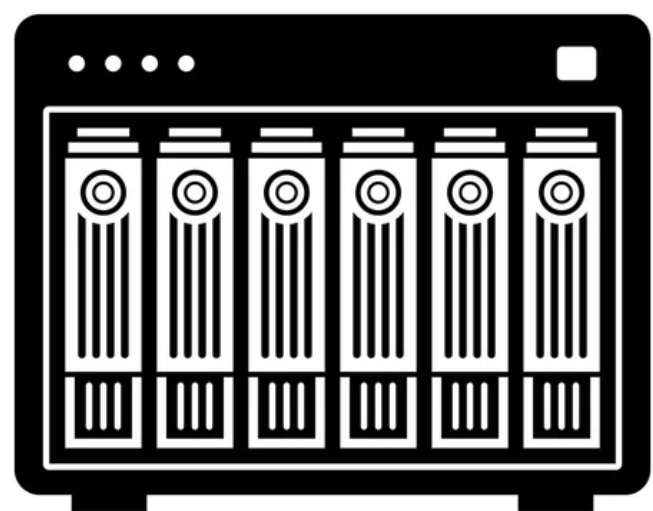
LSM-tree
O'Neil *et al.*



1996

● good random writes

● good reads



array of discs

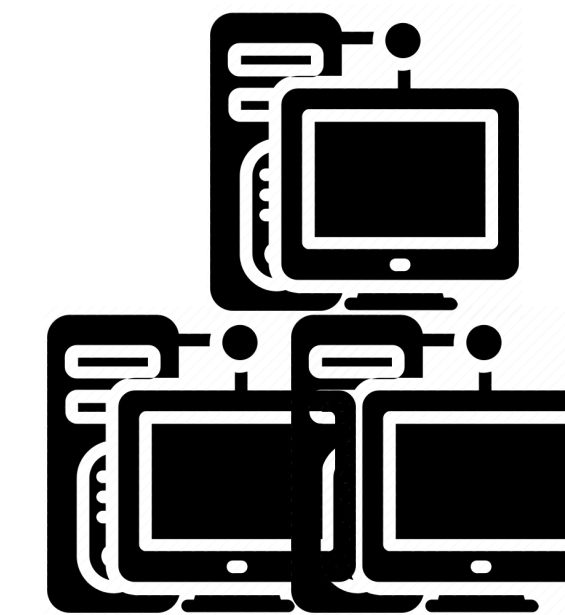
SSD wear-friendly

competitive rand. reads

fast ingestion

● poor ingestion perf.

● poor query perf.



commodity hardware



LSM-tree
O'Neil *et al.*

1980s

1996

2006

a decade



Bigtable

LSM-tree
O'Neil *et al.*

1996


Bigtable

2006

APACHE
HBASE 

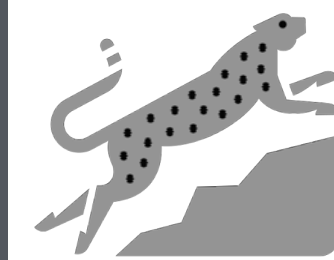
2007


cassandra

2010


levelDB

2011



RocksDB

2013

LSM-tree

NoSQL

This block contains logos for several NoSQL databases: RocksDB (yellow cheetah), WT (black and orange letters), levelDB (green cylinder), SCYLLA (blue alien head), DynamoDB (blue cylinder), cassandra (blue eye), tarantool (red Venn diagram), Bigtable (blue hexagon), APACHE HBASE (red and black text with orca), and riak (grey text with dots).

This block contains the SQLite logo (blue square with feather) and a dark blue square logo featuring a grey dolphin, representing a relational database.

relational

This block contains the influxdb logo (blue cube) and the QuasarDB logo (blue grid pattern).

time-series

2022

Why **LSM** ?

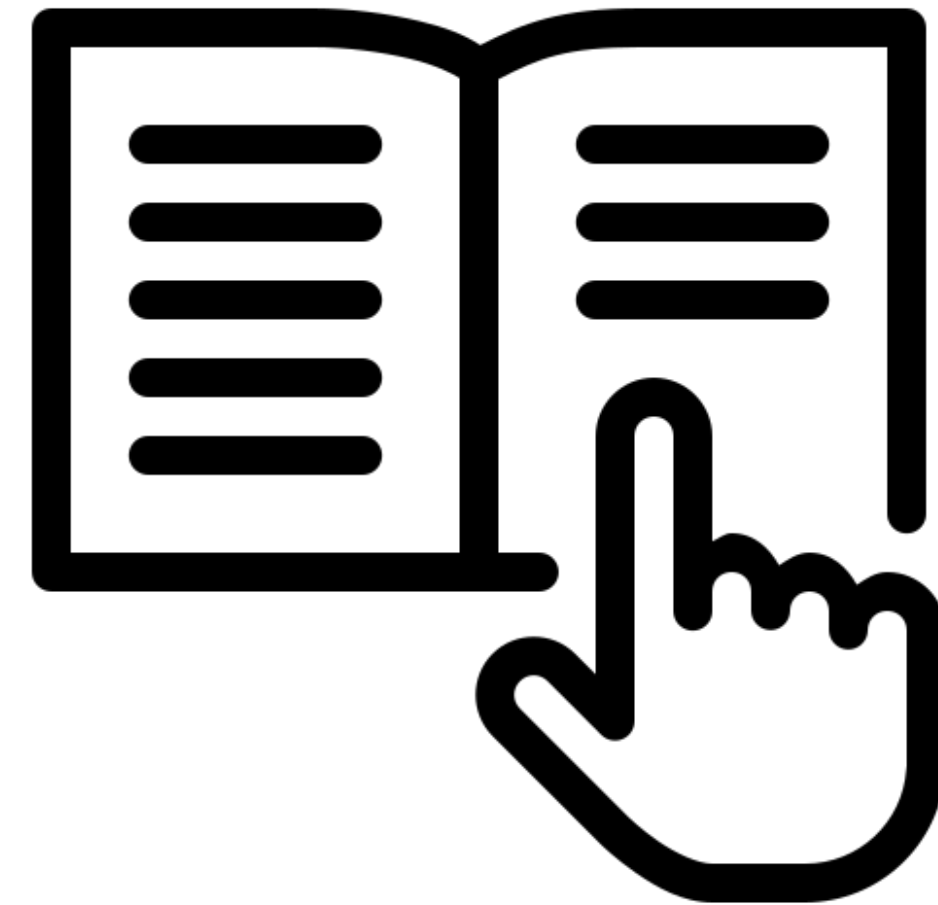


fast ingestion

Why **LSM** ?

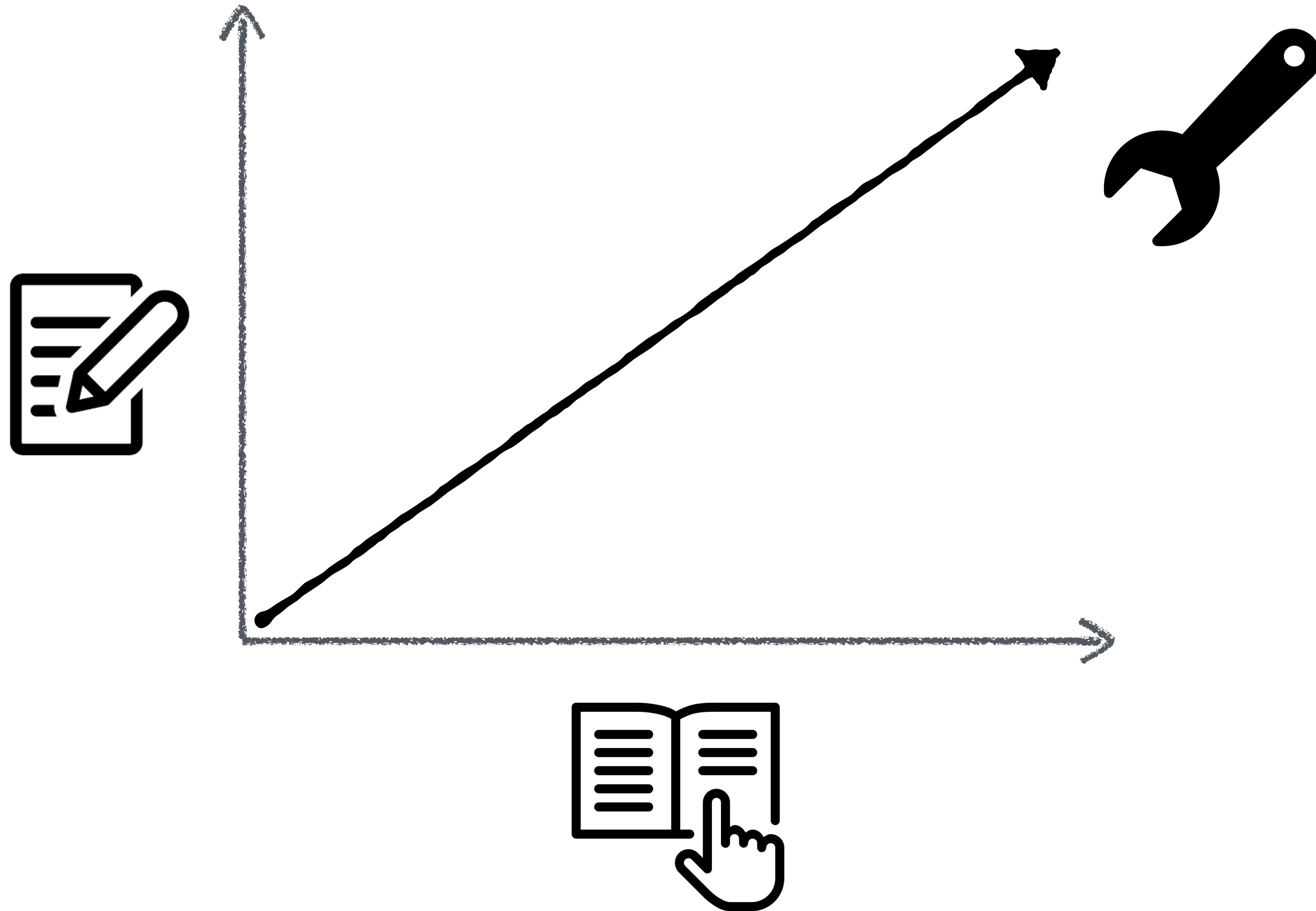


fast ingestion

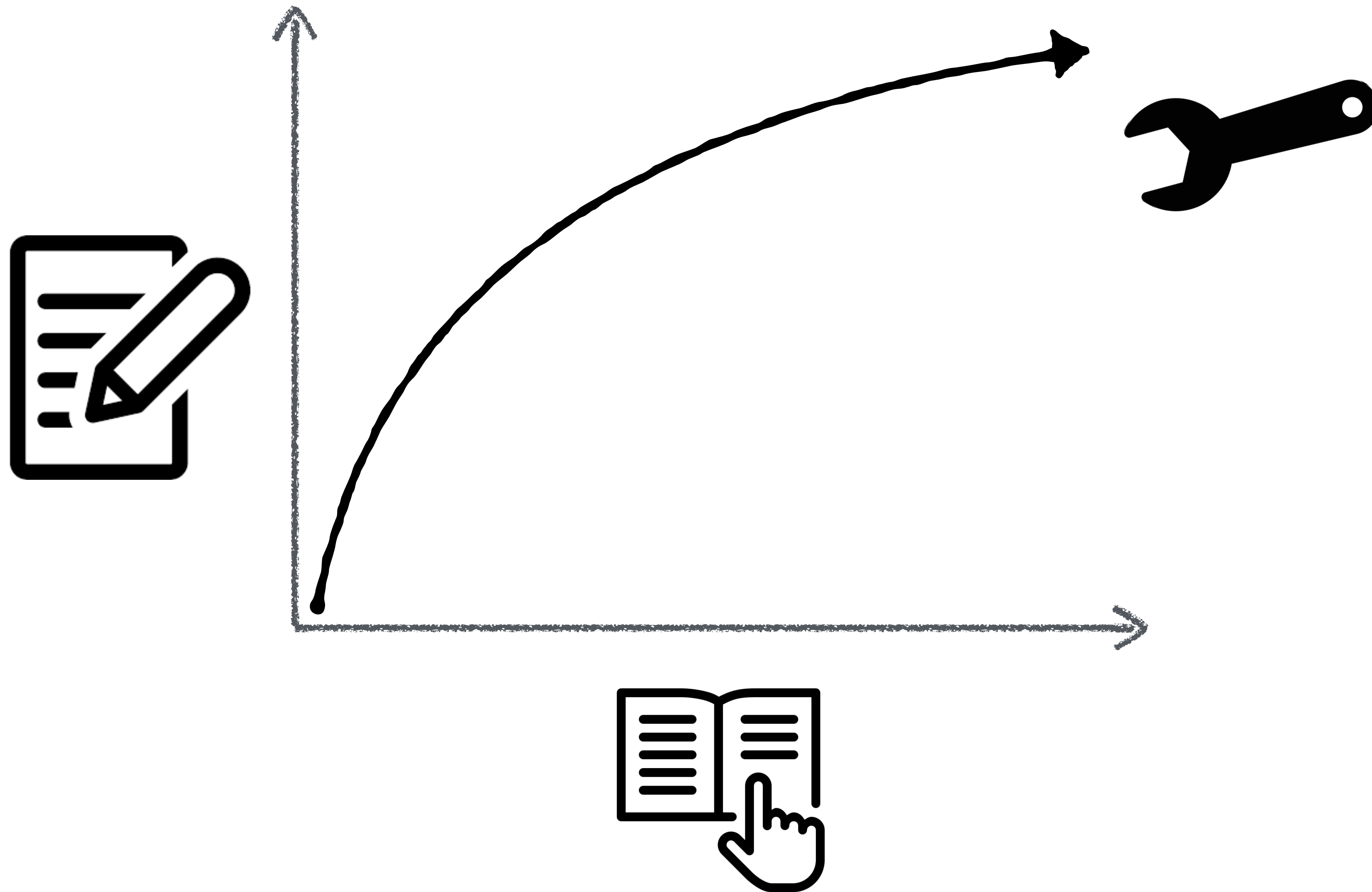


competitive reads

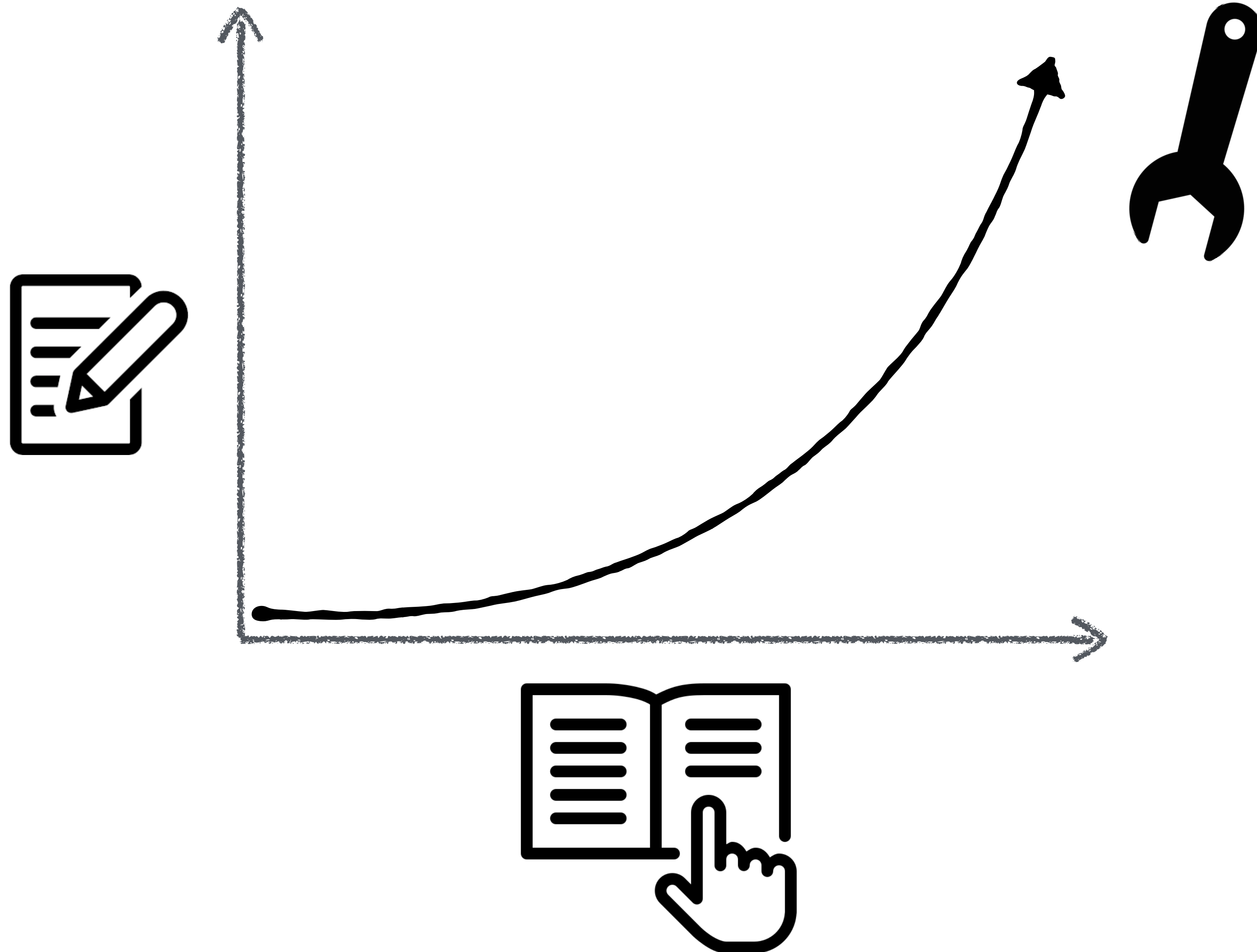
Why **LSM** ?



Why **LSM** ?



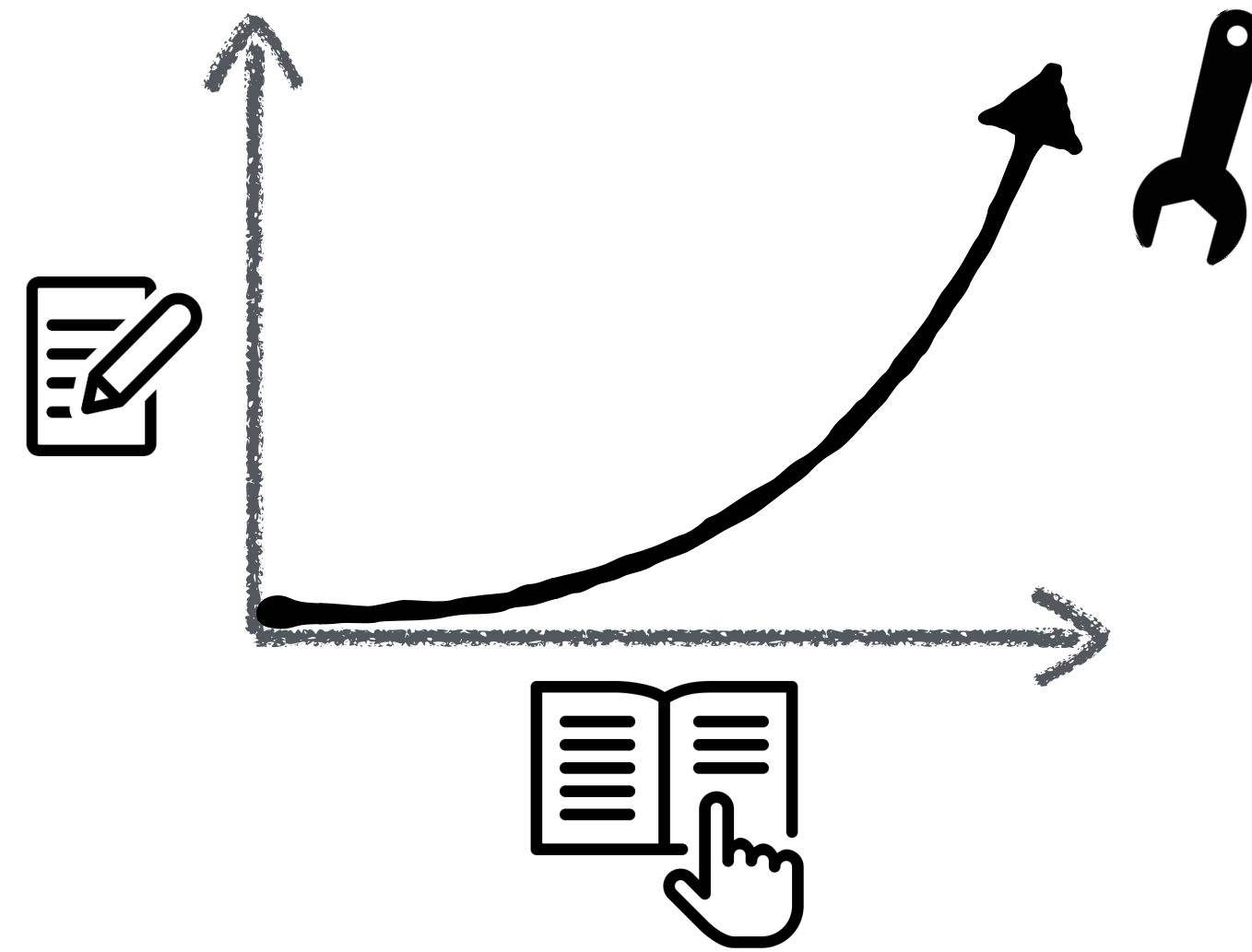
Why **LSM** ?



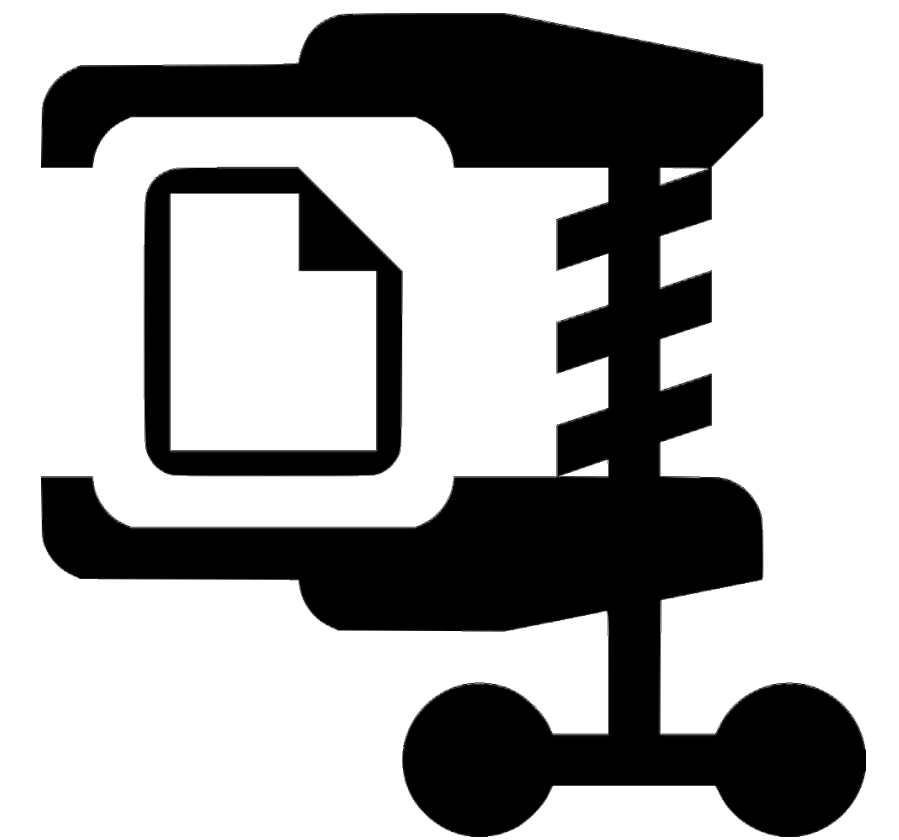
Why **LSM** ?



fast writes

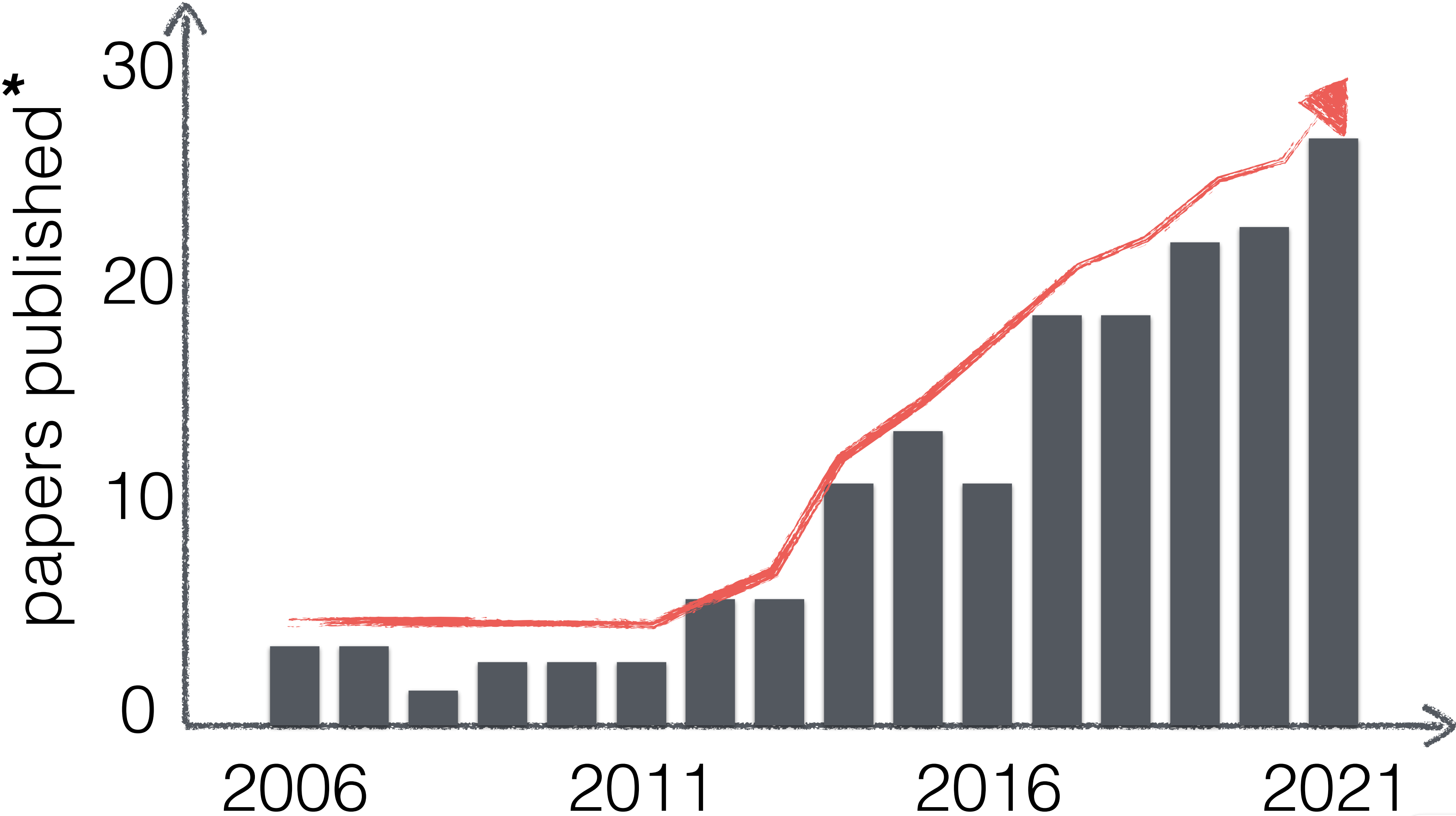


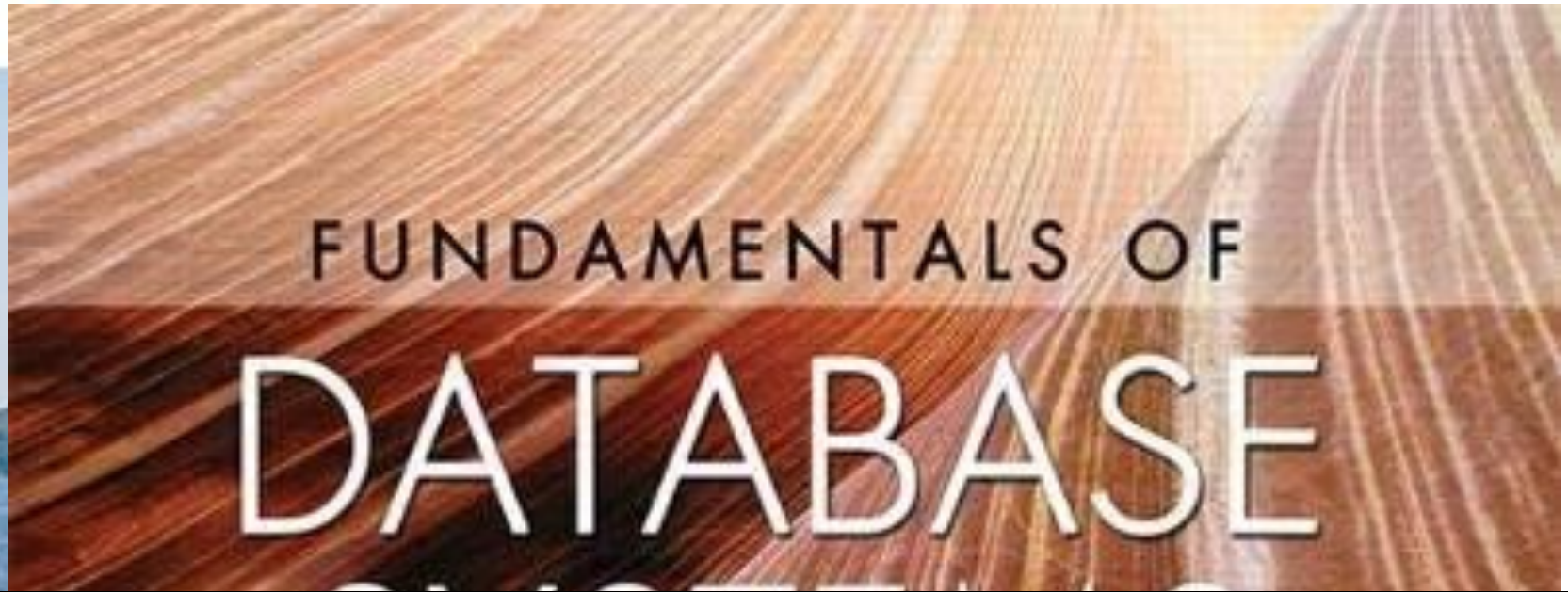
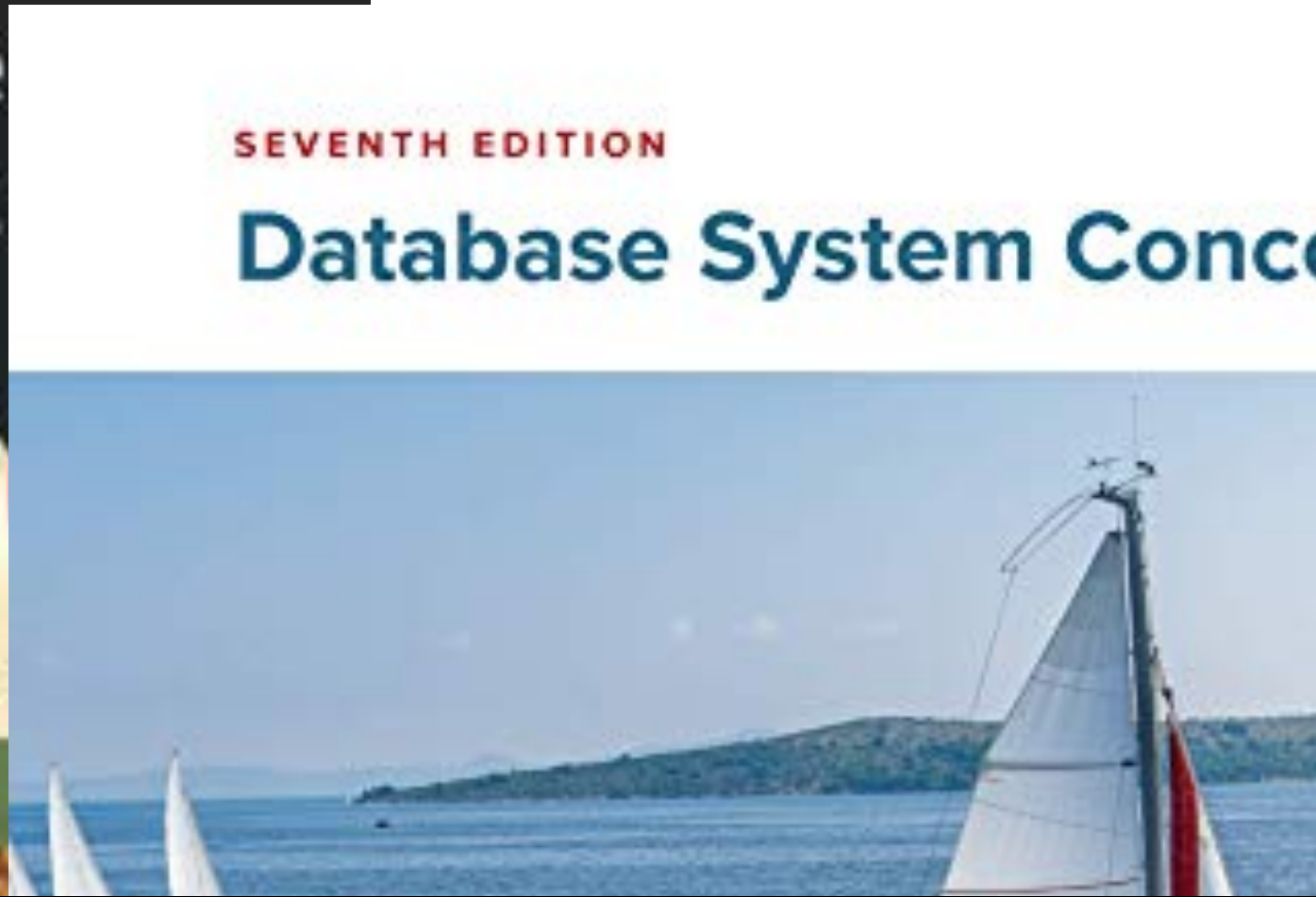
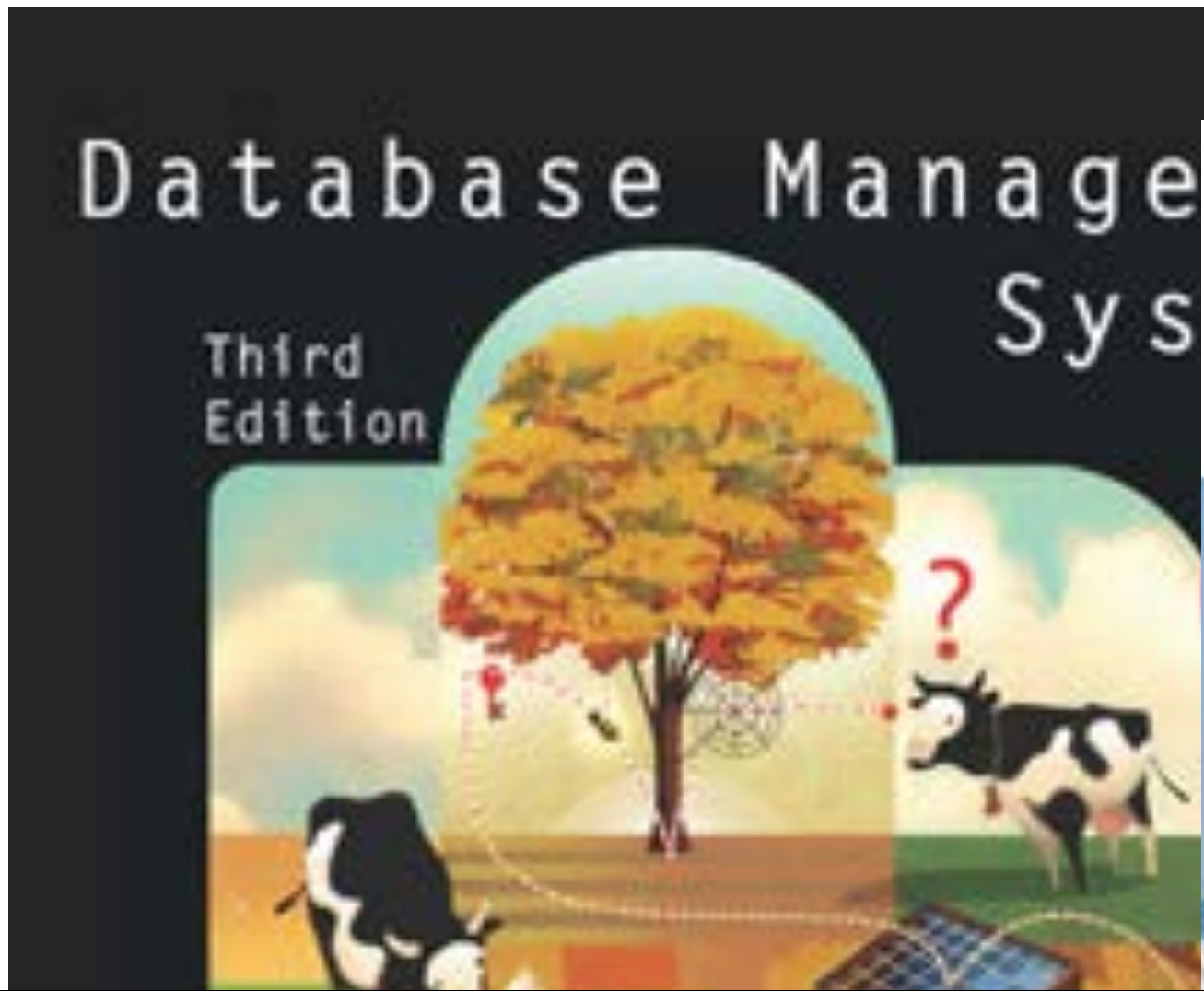
tunable read-write
performance



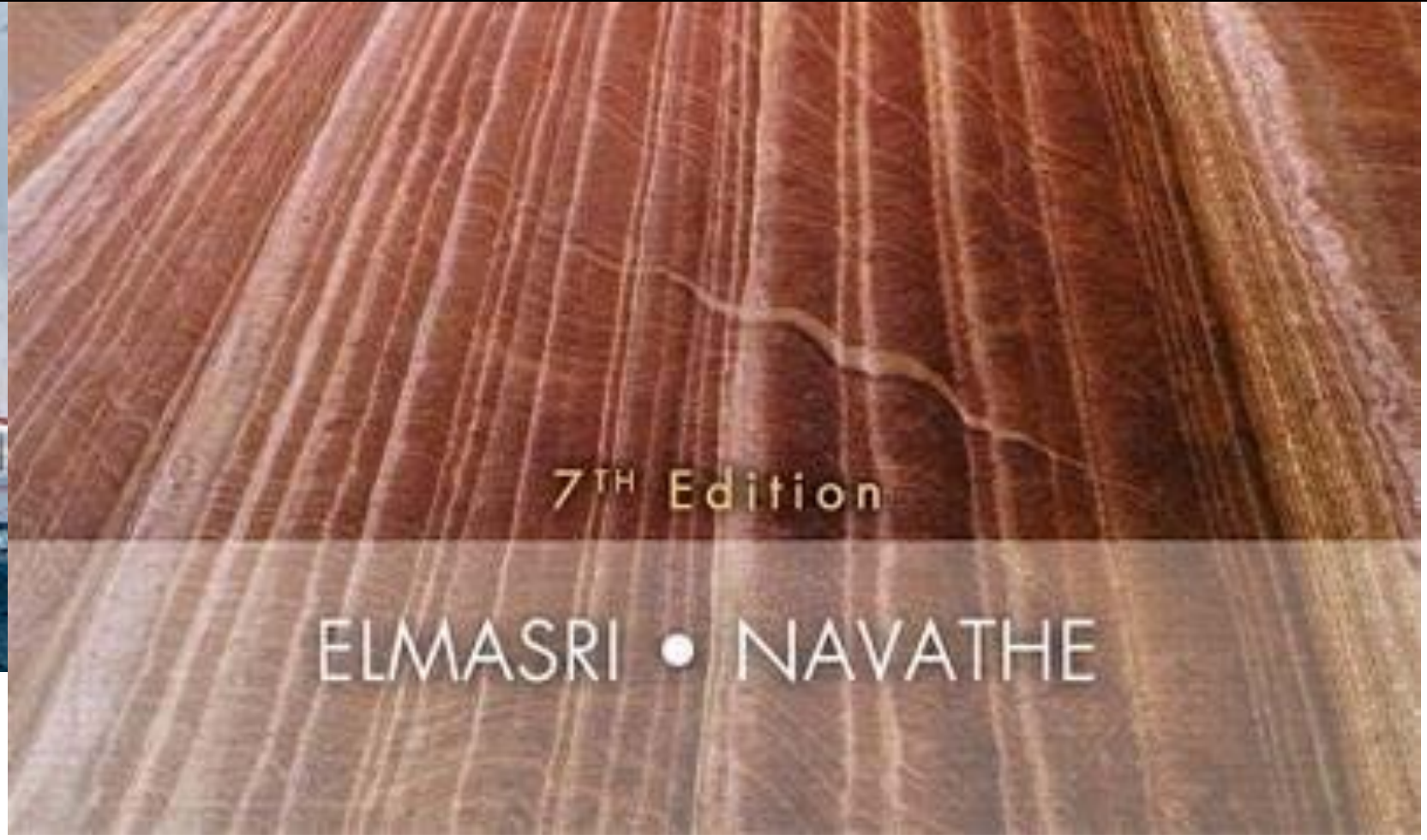
good space
utilization

Research Trend





No Textbook on LSMs !!



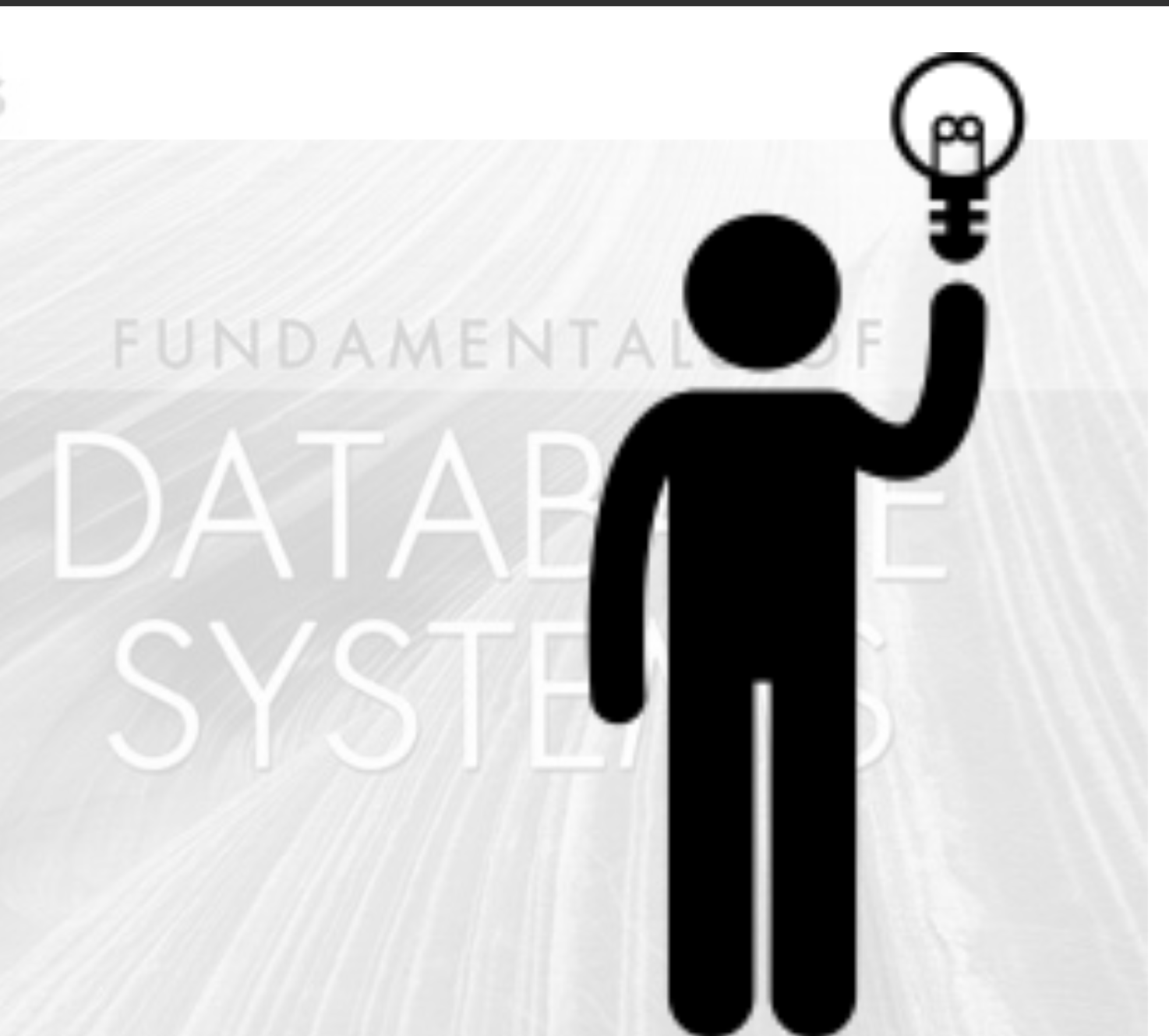
No Textbook on LSMs !!



explore the
LSM paradigm



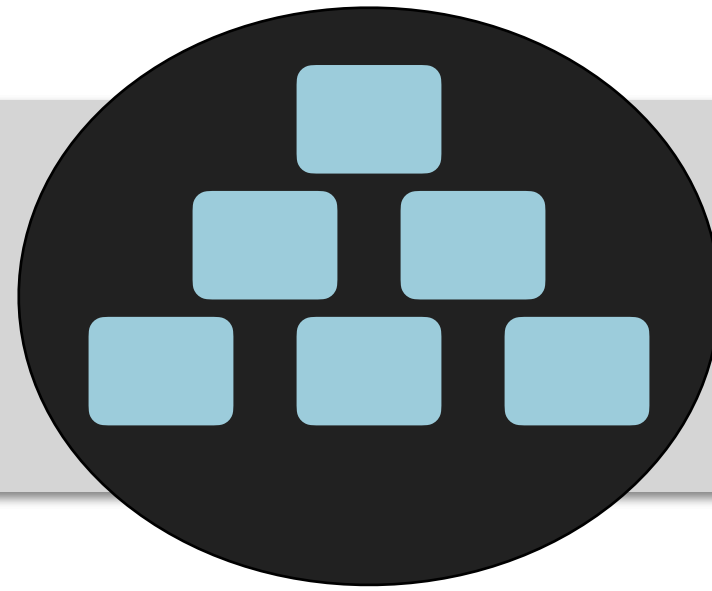
construct LSM
design space



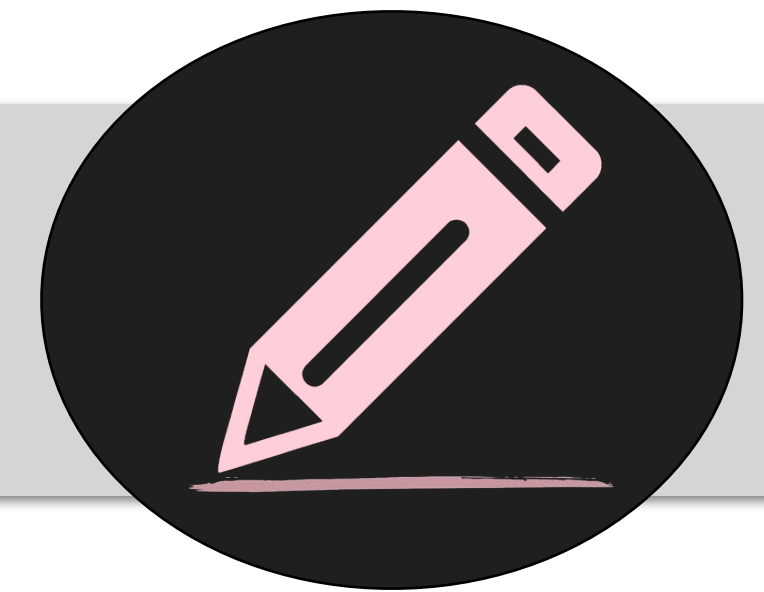
provide useful
insights

Outline

Part 1: **LSM Basics**



Part 2: **Optimizing Ingestion in LSMs**

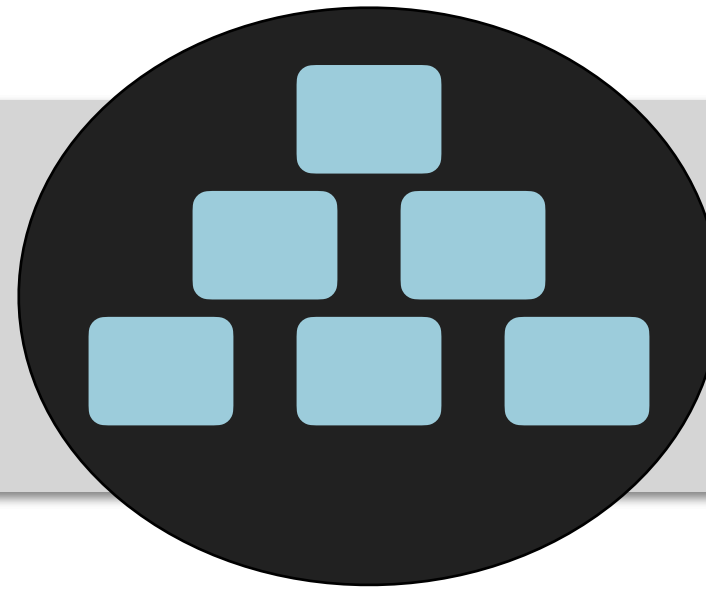


Part 3: **Navigating the LSM Design Space**

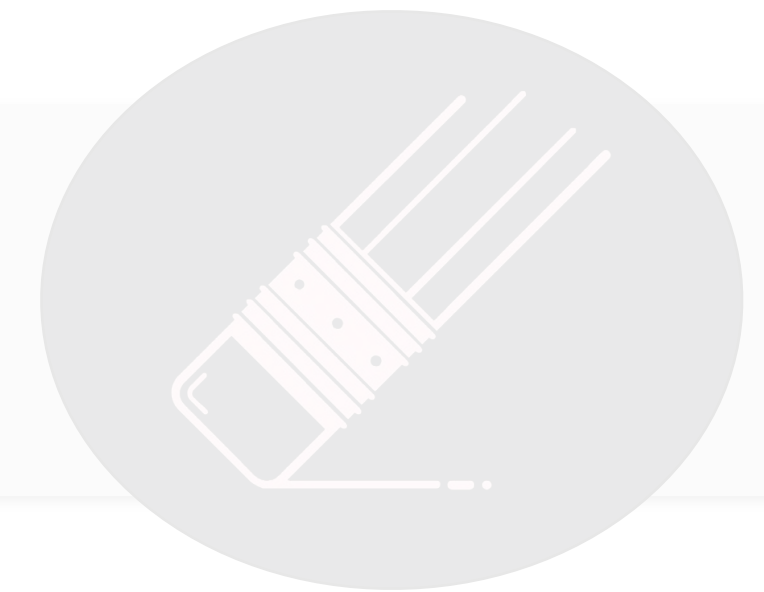


Outline

Part 1: **LSM Basics**



Part 2: Optimizing Ingestion in LSMs

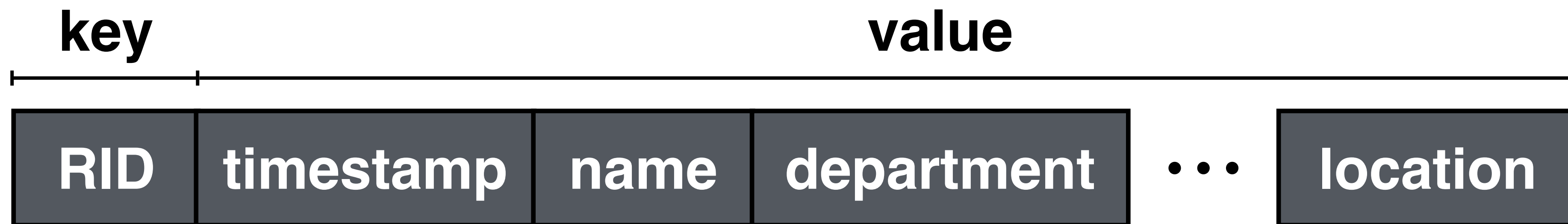


Part 3: Navigating the LSM Design Space



LSM Basics

key-value pairs



LSM Basics

key-value pairs



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

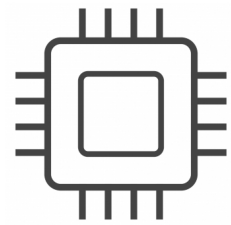
LSM Basics



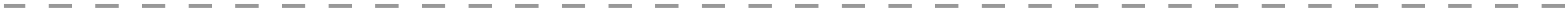
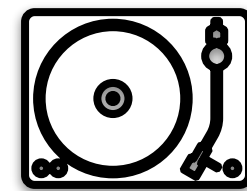
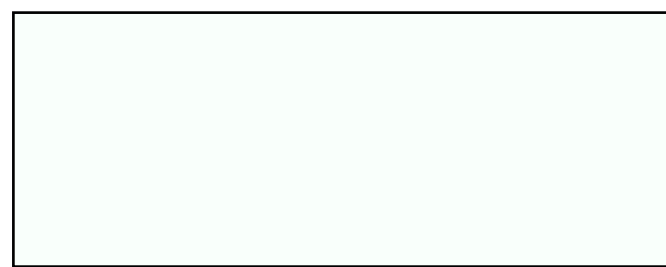
Buffering ingestion

put(6)

put(2)



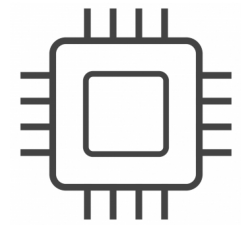
buffer



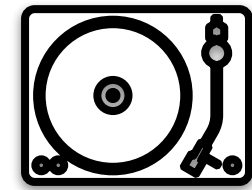
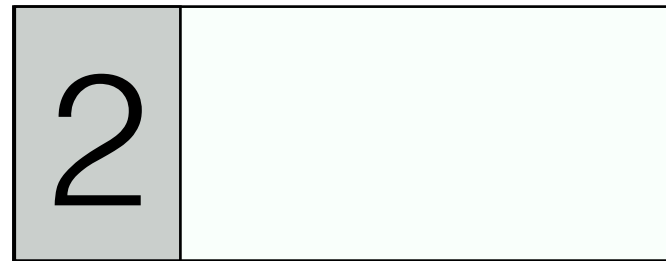
Buffering ingestion

put(1)

put(6)



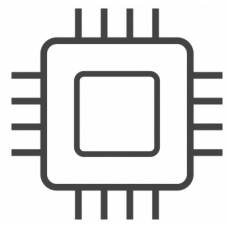
buffer



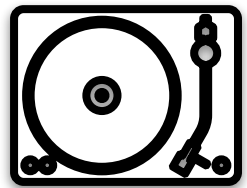
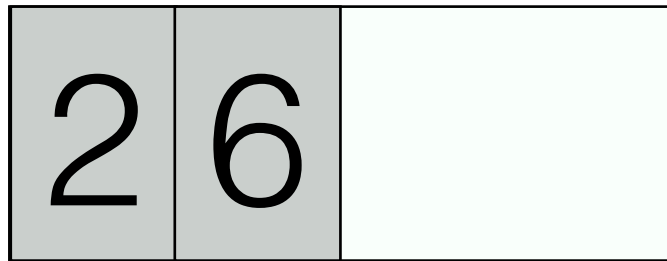
Buffering ingestion

put(4)

put(1)

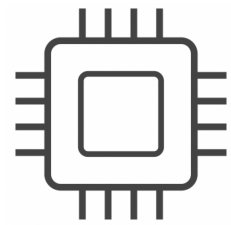


buffer

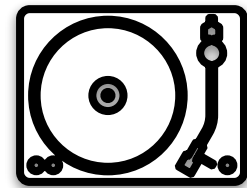
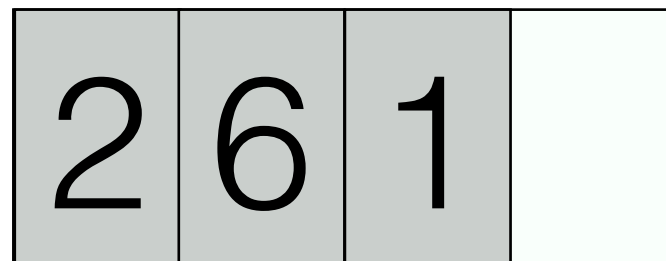


Buffering ingestion

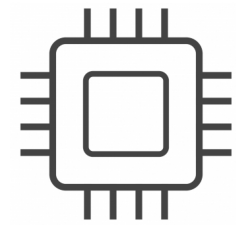
put(4)



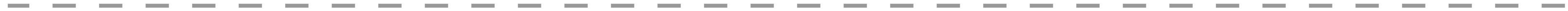
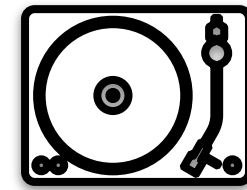
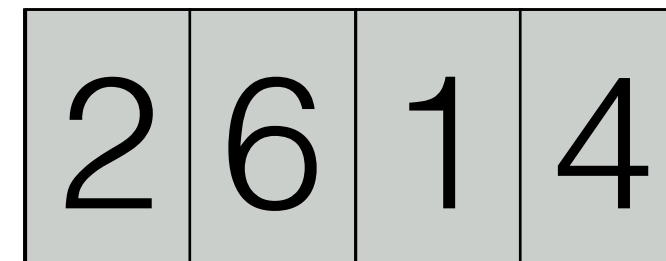
buffer



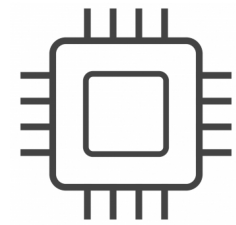
Buffering ingestion



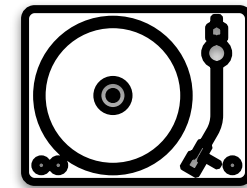
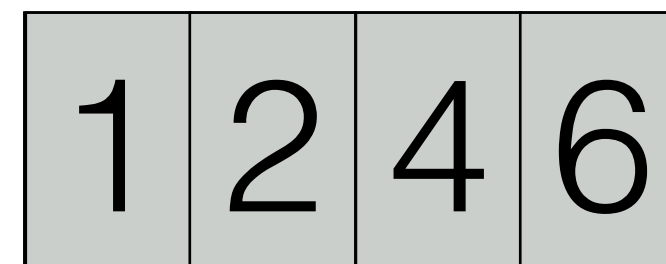
buffer



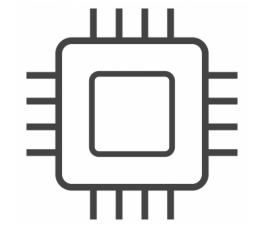
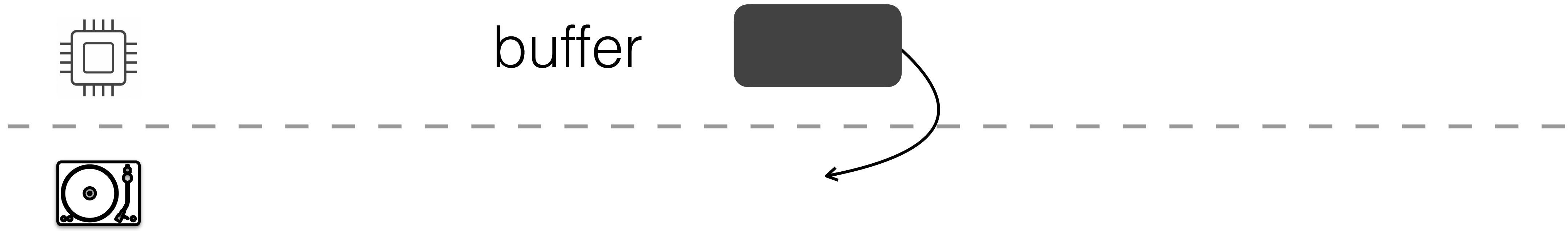
Buffering ingestion



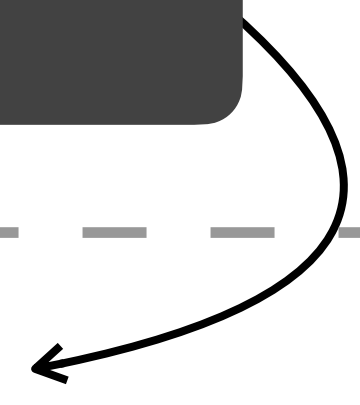
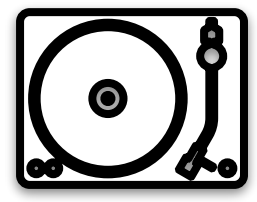
buffer



Buffering ingestion

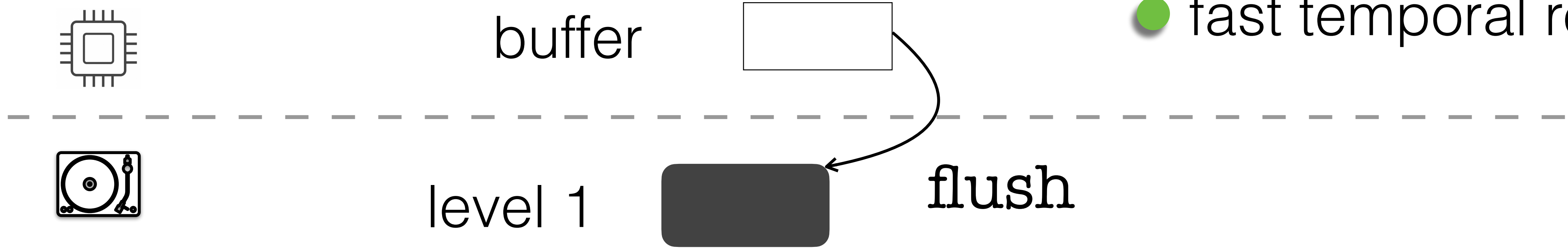


buffer



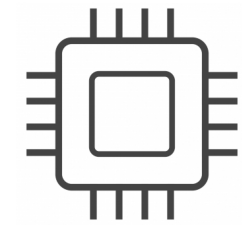
Buffering ingestion

- low ingestion cost
- fast temporal reads

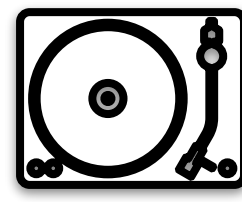


Immutable files on storage

- compact storage
- good ingestion throughput



buffer 

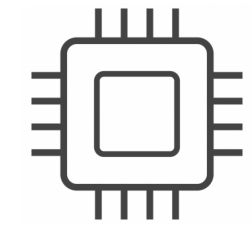


level 1 

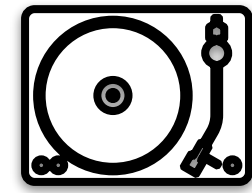
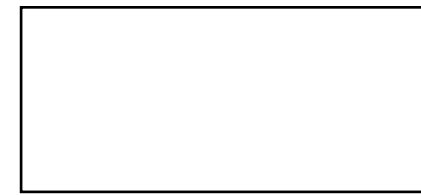


How do we update data?

put(6)

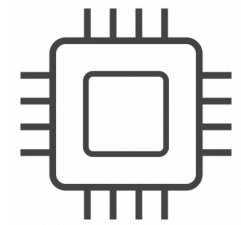


buffer

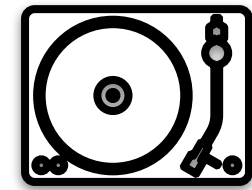


level 1

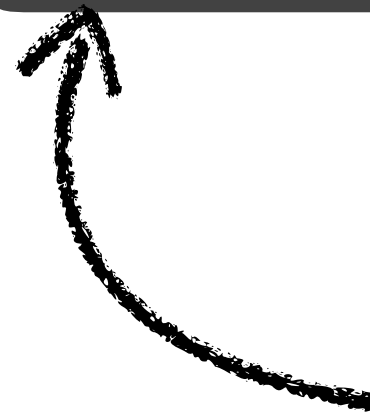
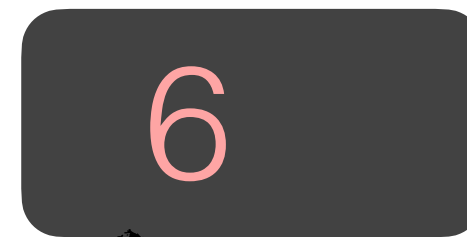




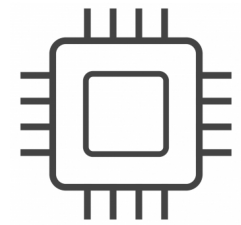
buffer



level 1

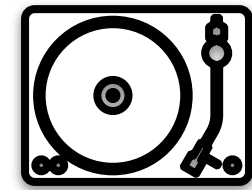
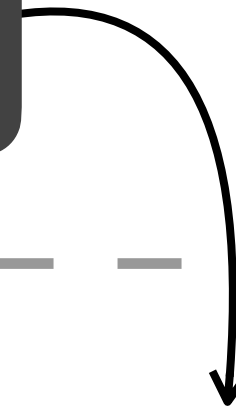


logically
invalidated



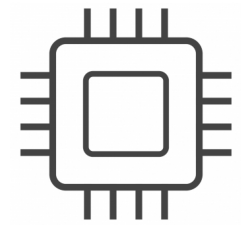
buffer

6



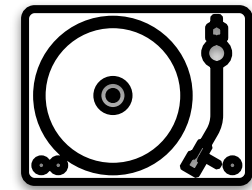
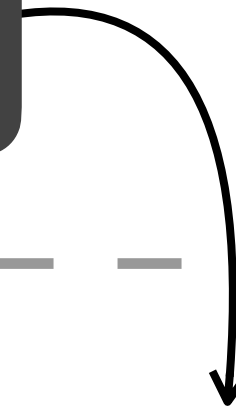
level 1

6



buffer

6

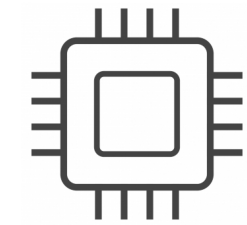


level 1

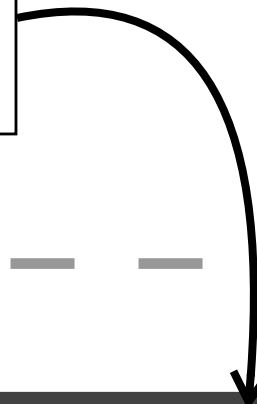
6

Out-of-place updates

- fast ingestion
- space amplification
- slow reads



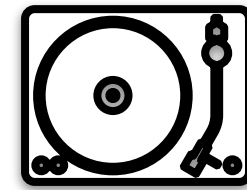
buffer



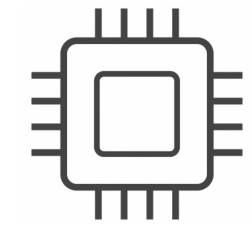
level 1

6

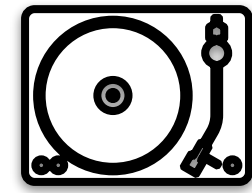
6



How do we reduce this space amplification?



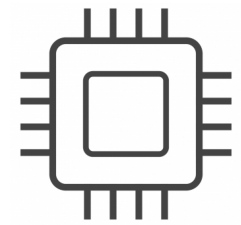
buffer



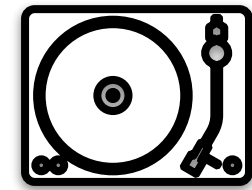
level 1

6

6



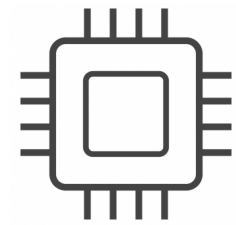
buffer



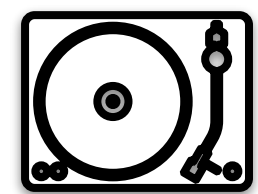
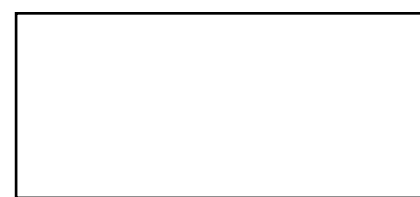
level 1

6

compaction



buffer



level 1



level 2

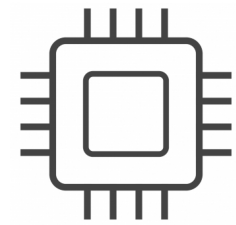


level 3

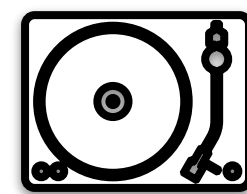


level 4





buffer



level 1



level 2

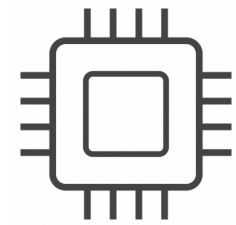


level 3

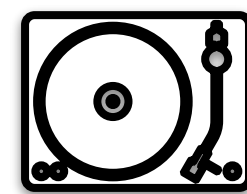
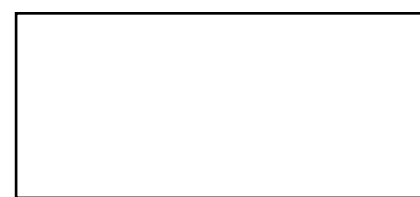


level 4





buffer



level 1



level 2

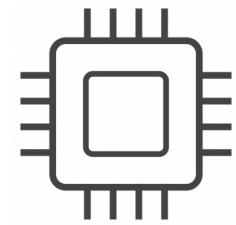


level 3

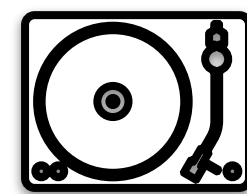
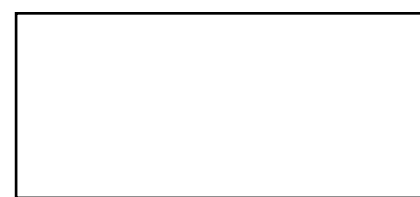


level 4





buffer



level 1



level 2

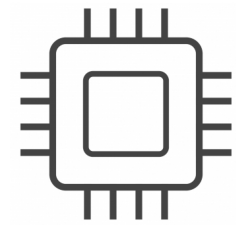


level 3

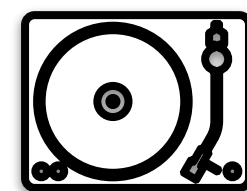
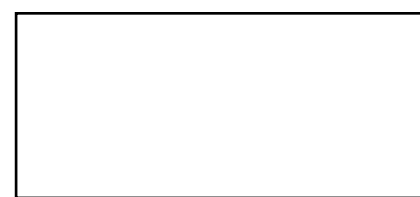


level 4





buffer



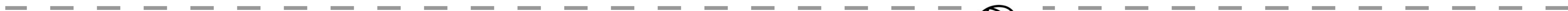
level 1

level 2



level 3

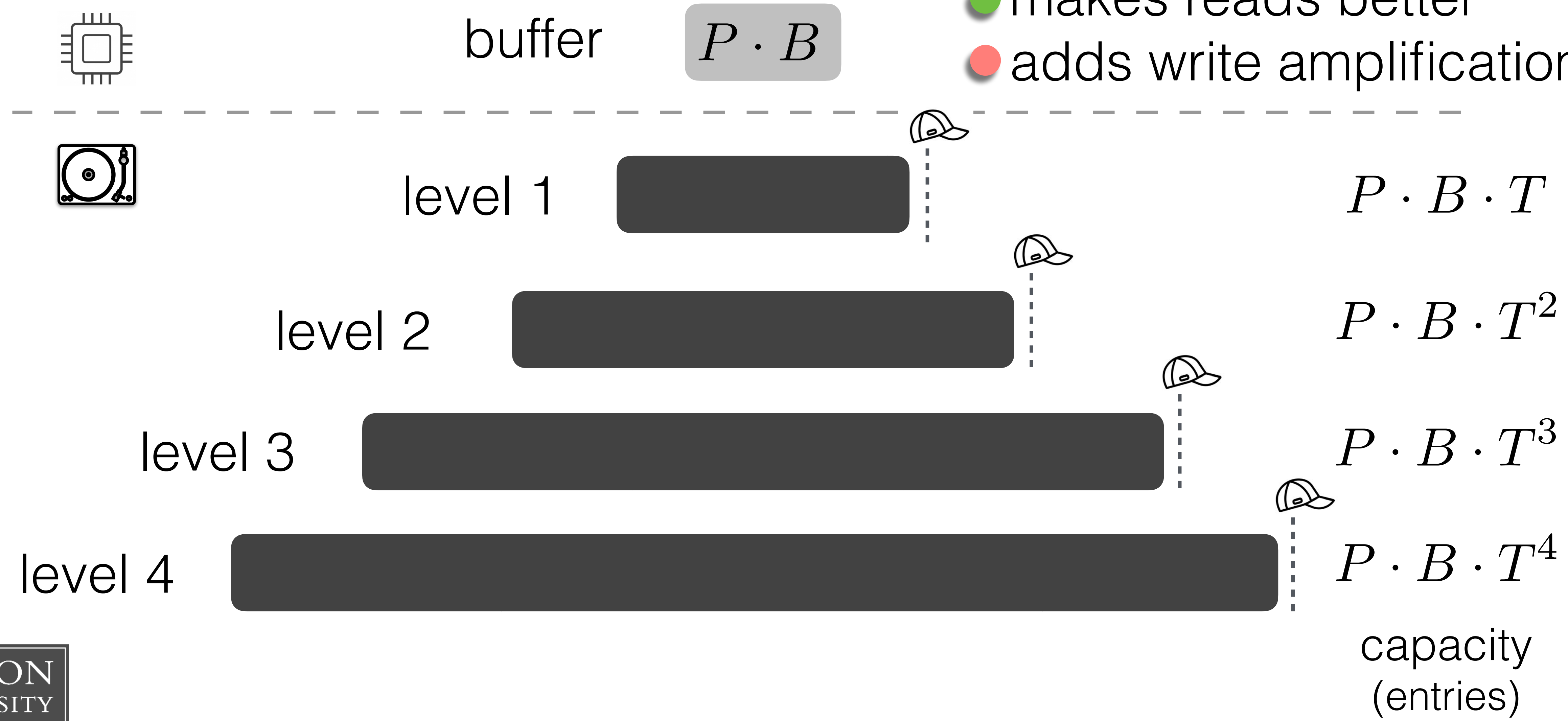
level 4



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

Periodic compactions

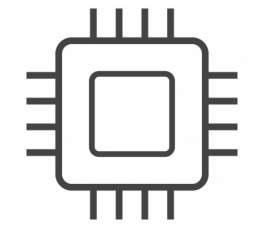
- low space amplification
- makes reads better
- adds write amplification



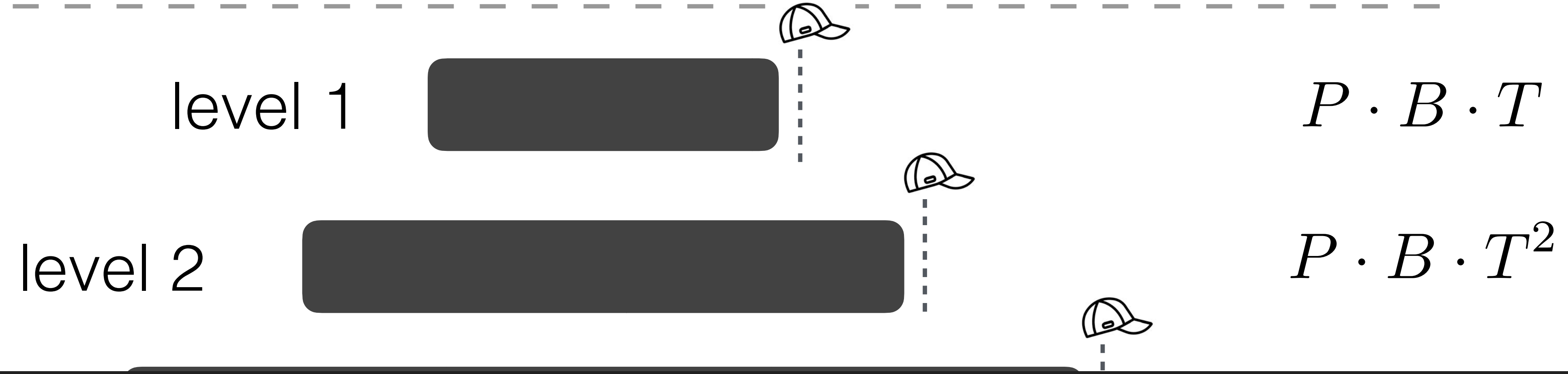
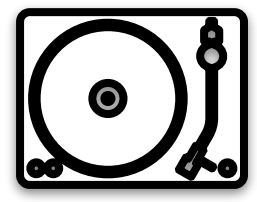
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio

Periodic compactions

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



buffer $P \cdot B$



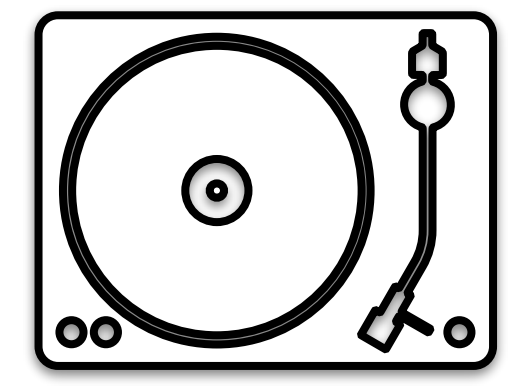
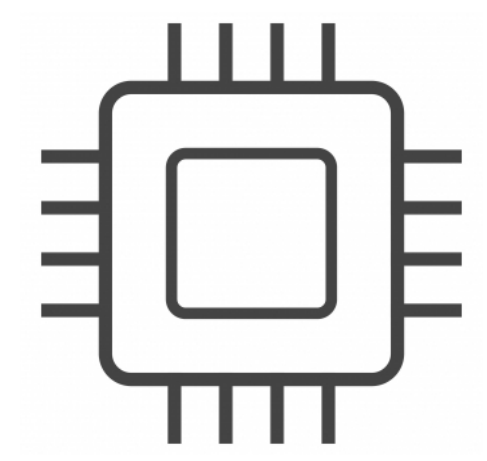
How about queries?

capacity
(entries)

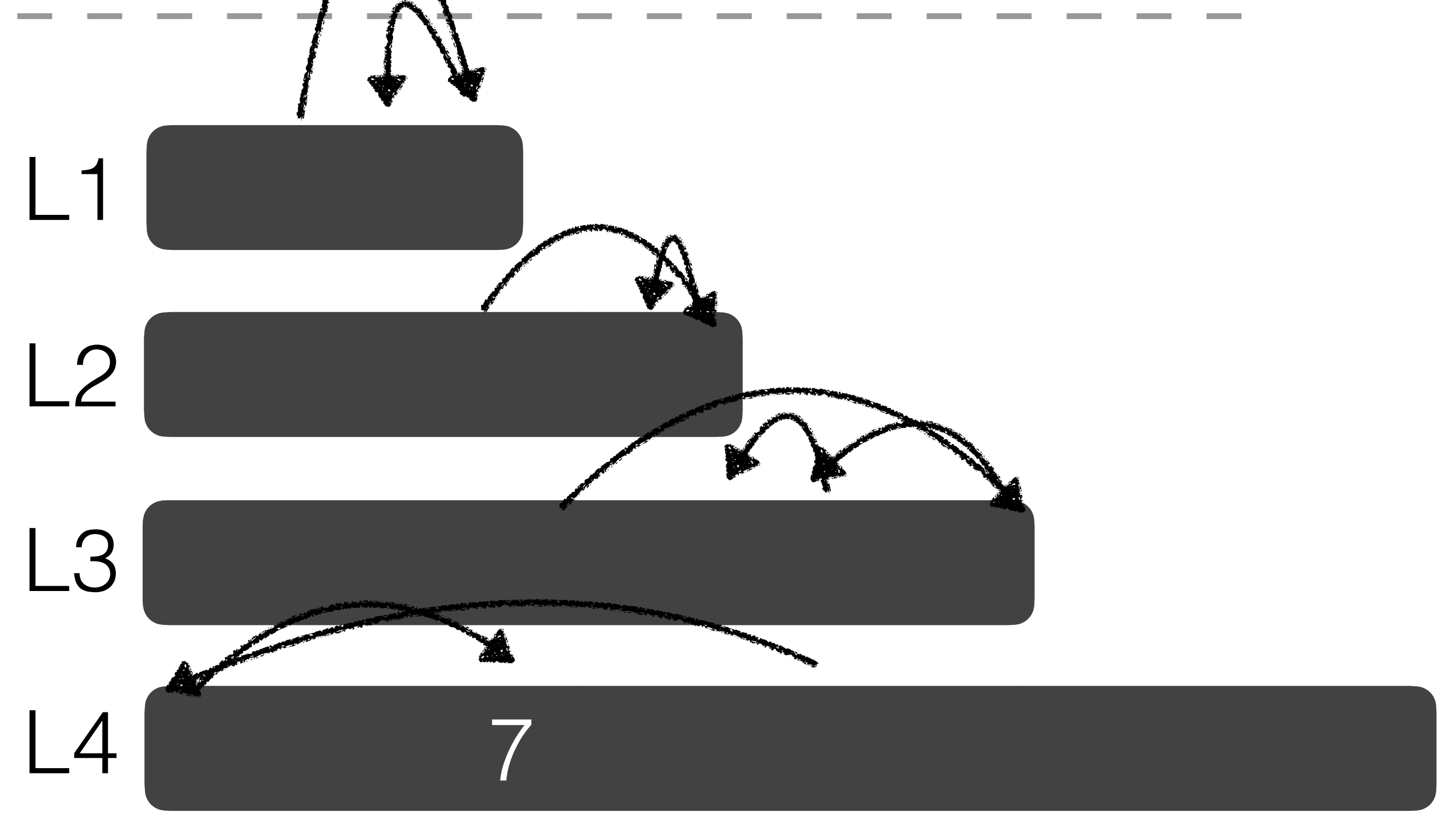
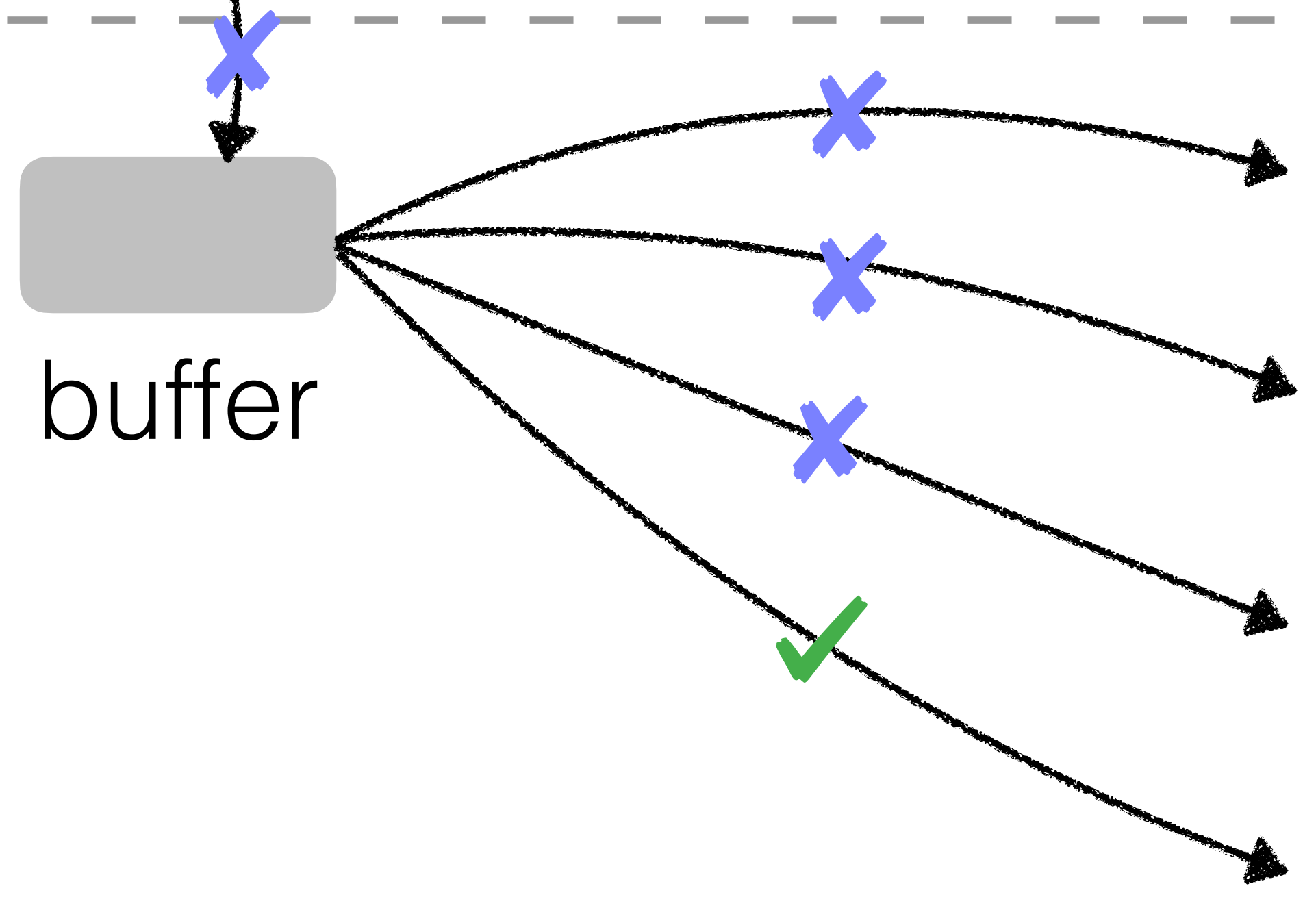
P: pages in buffer
B: entries/page
L: #levels
T: size ratio
N: #entries

Cost analysis

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



get(7)

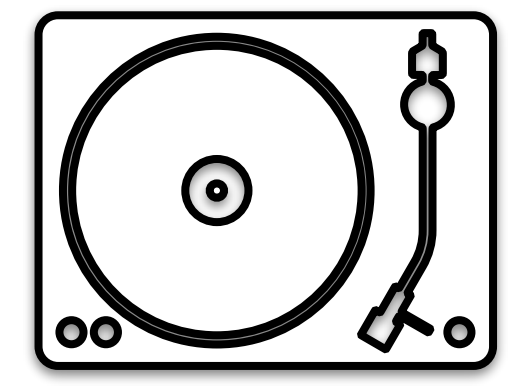
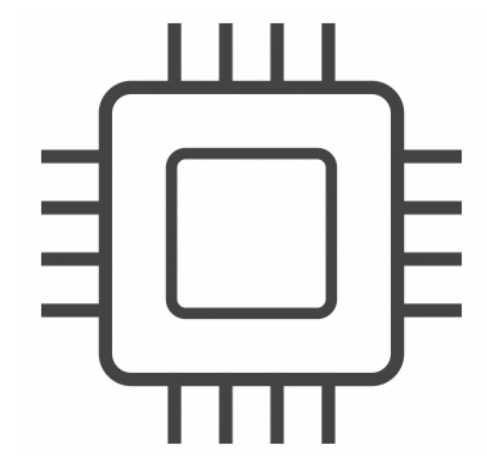


P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

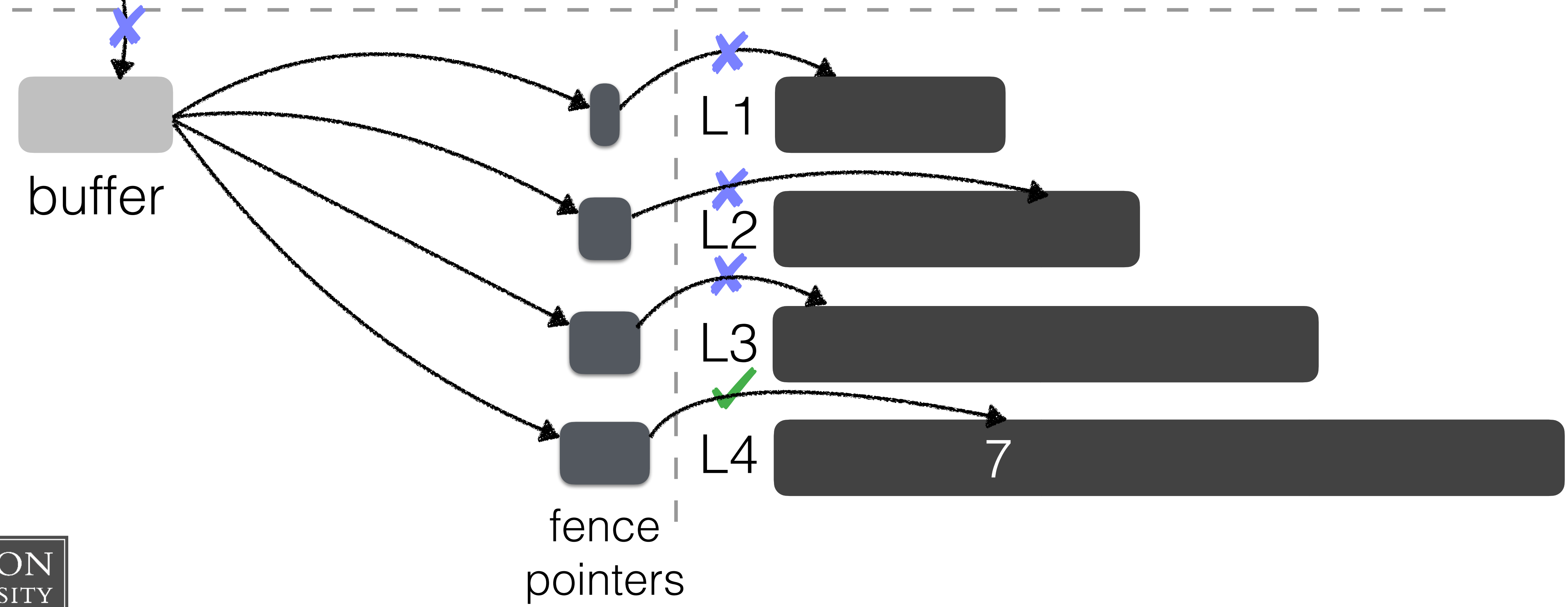
Cost analysis

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

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



get(7)

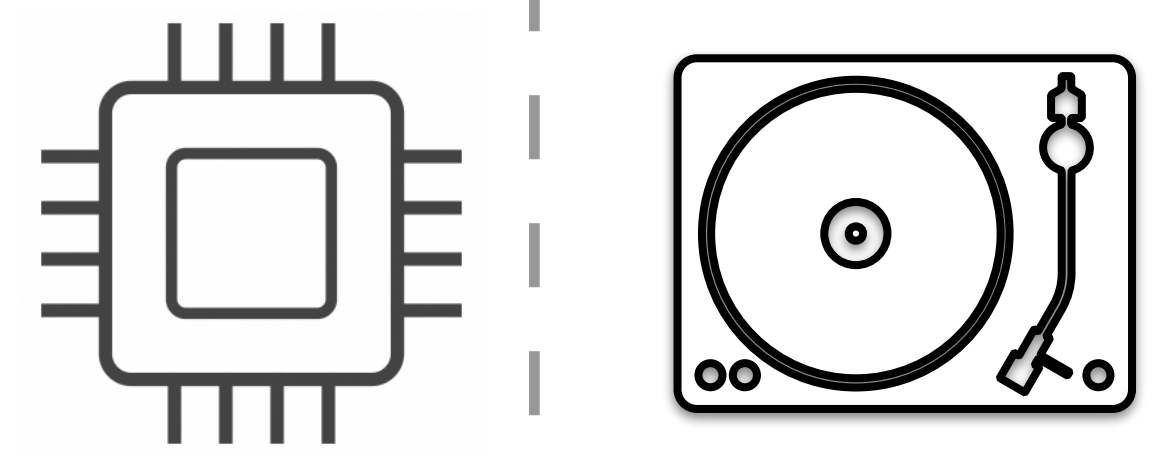


P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

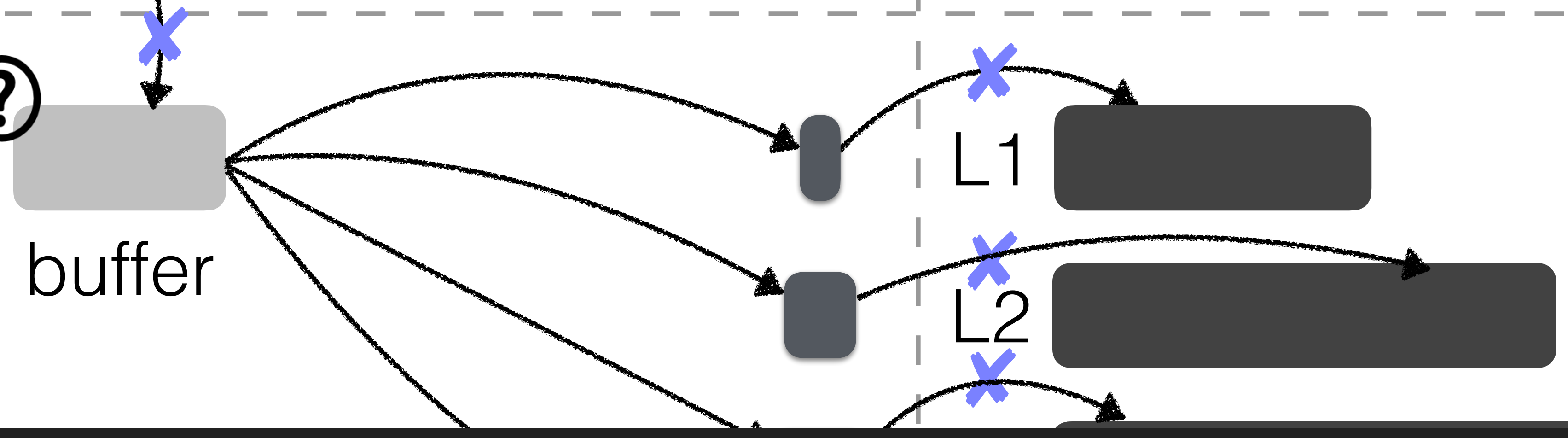
Cost analysis

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

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



get(7)



How to avoid unnecessary I/Os?

pointers

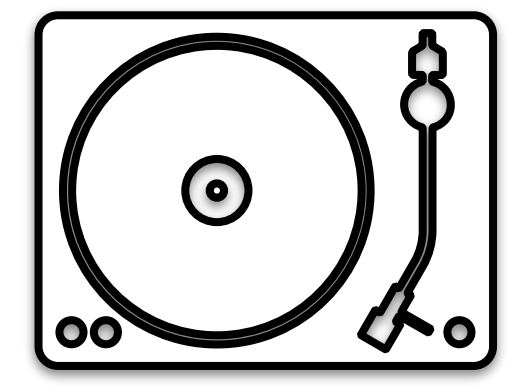
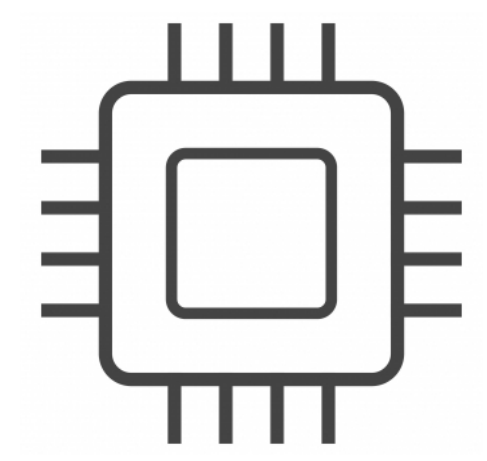
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

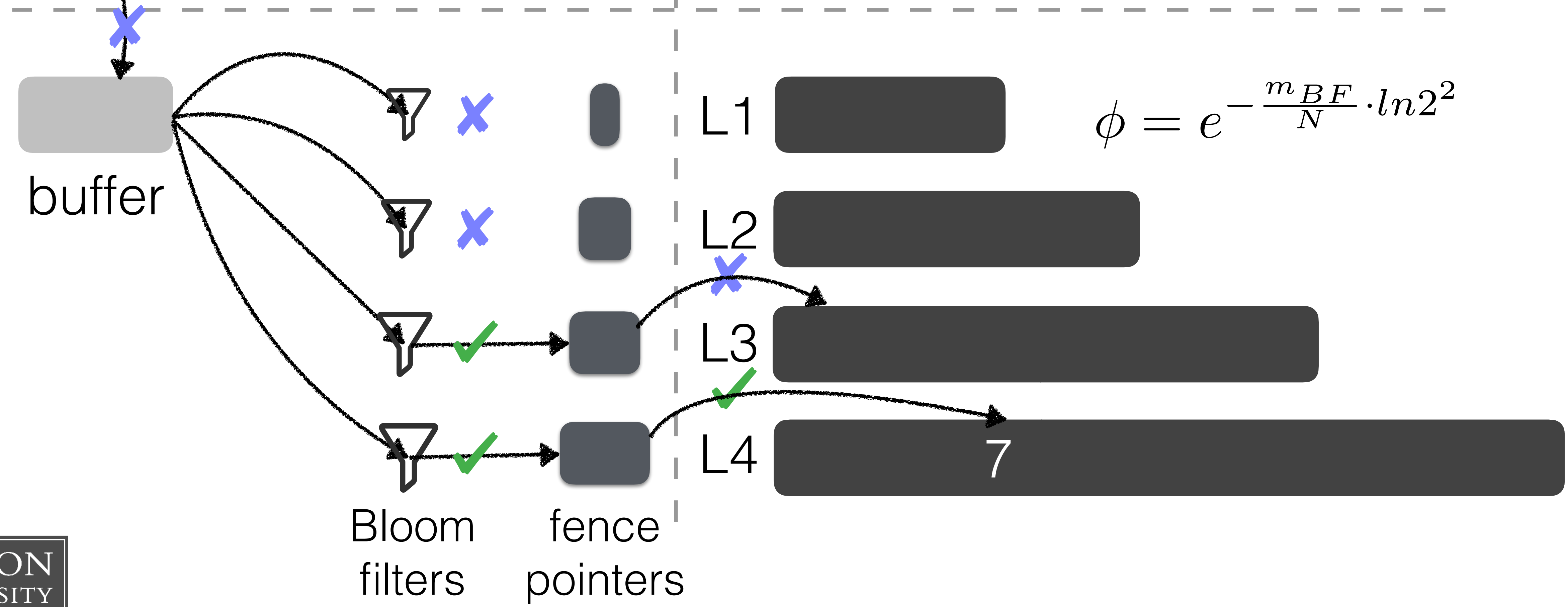
w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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

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



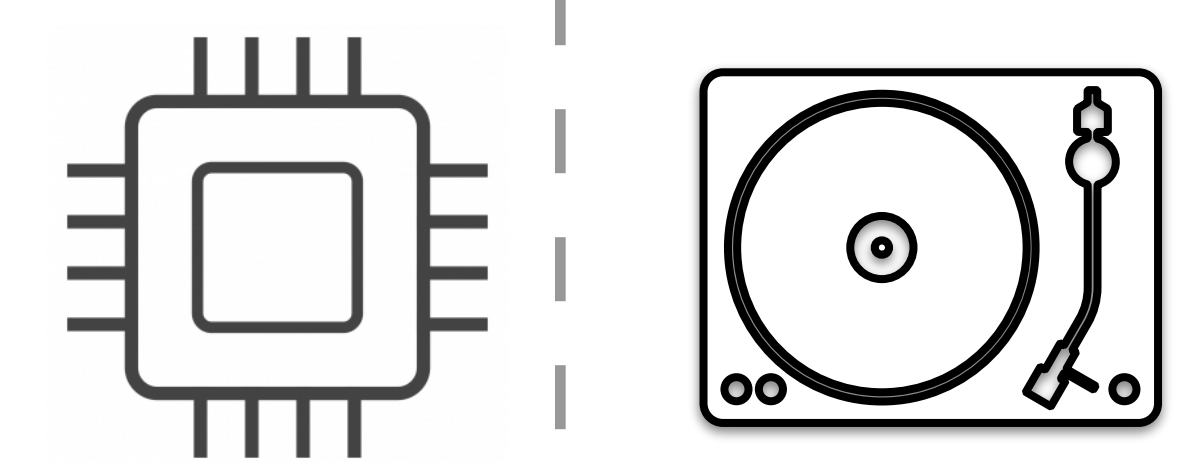
get(7)



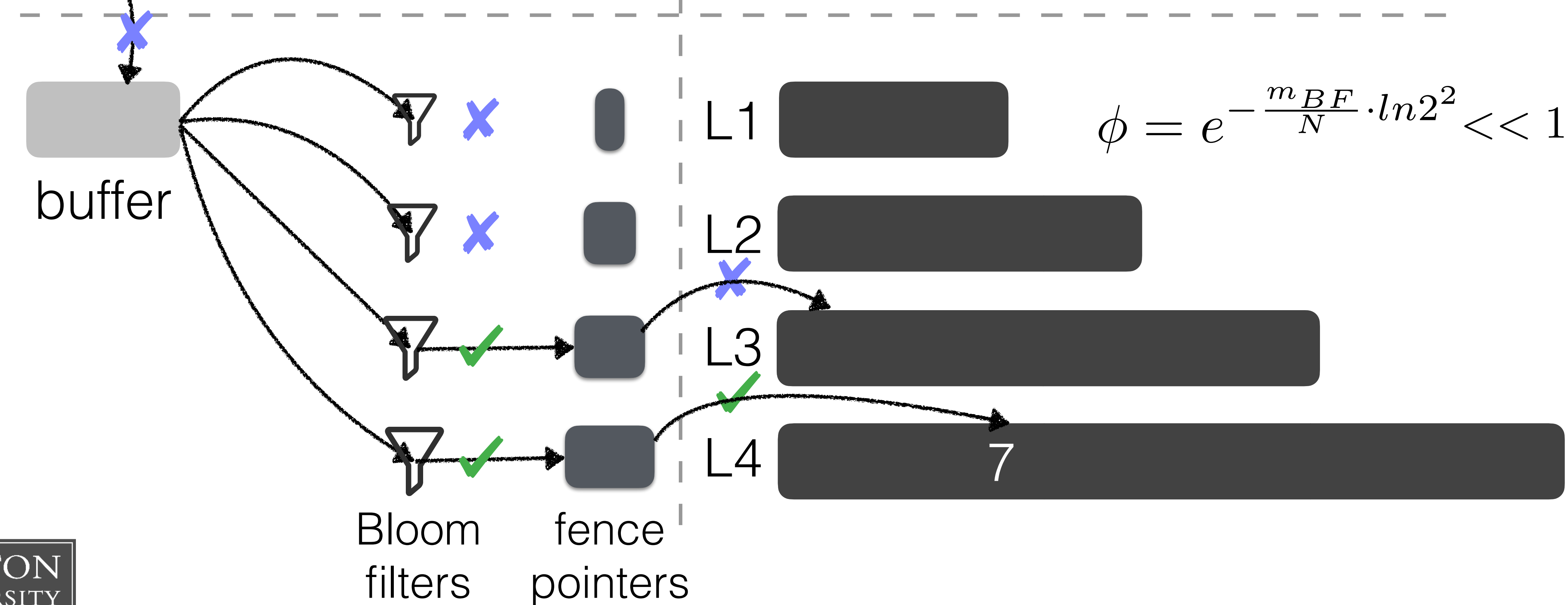
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

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



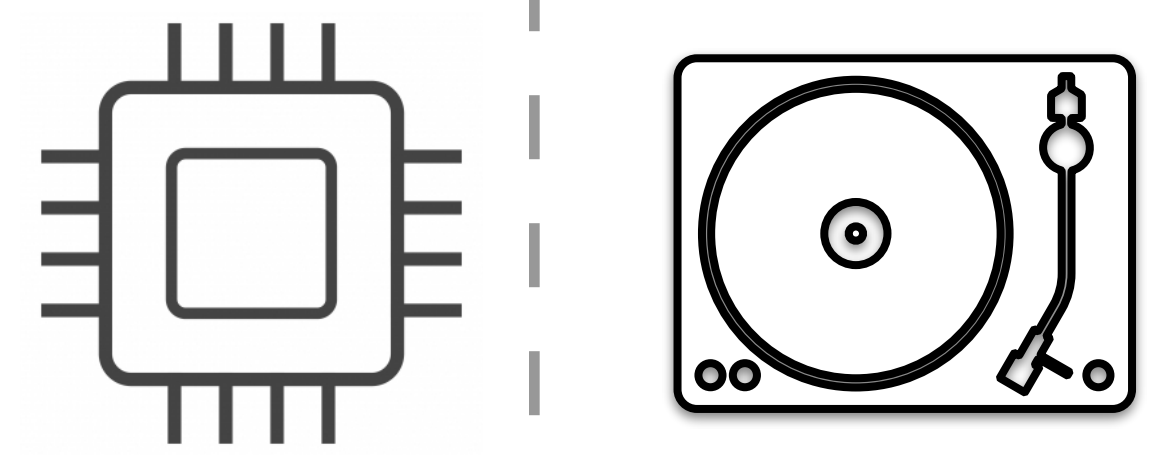
get(7)



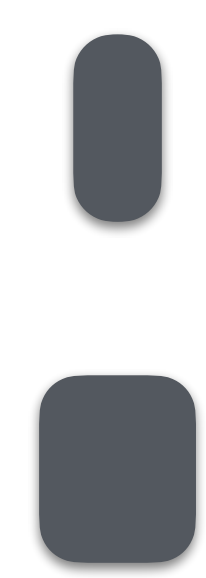
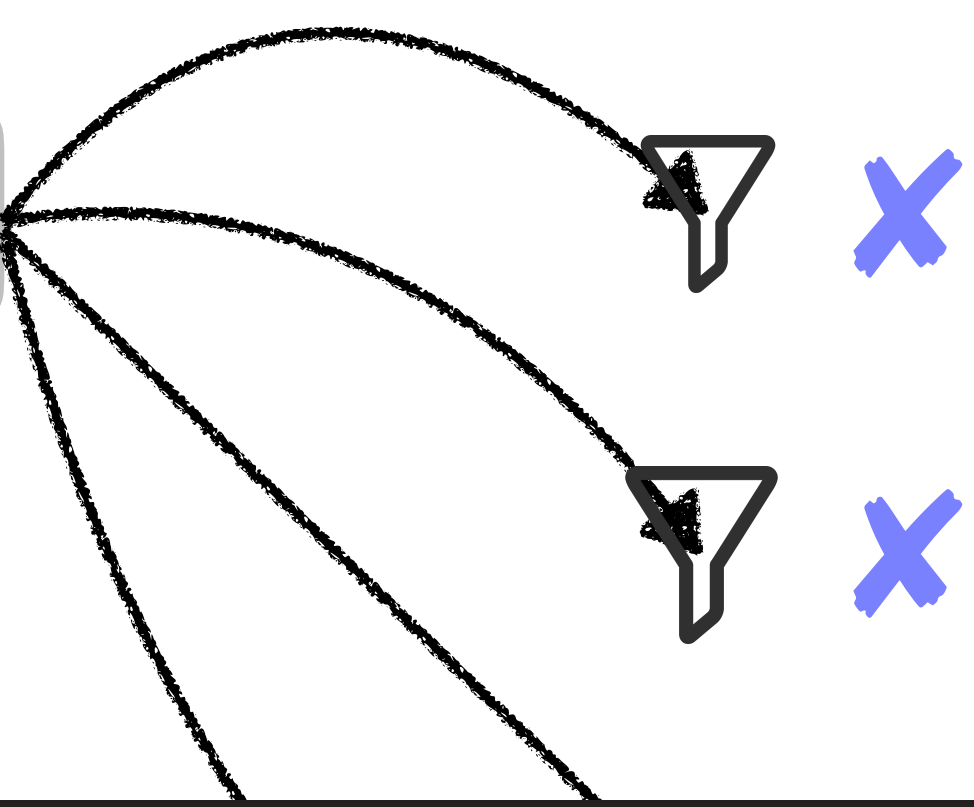
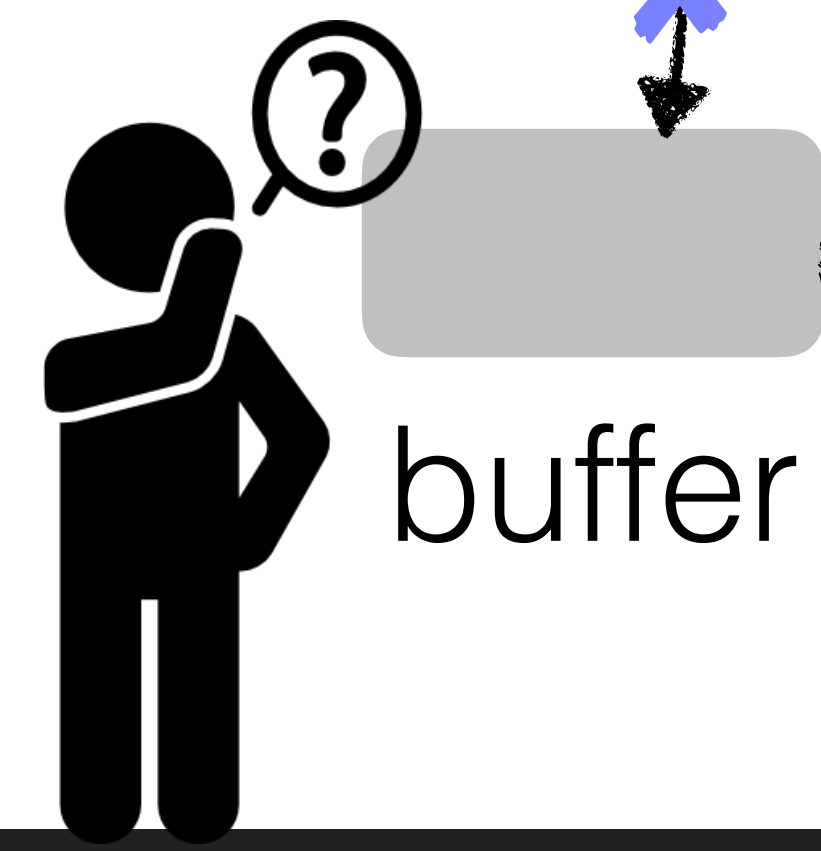
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

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



get(7)



L1
L2



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

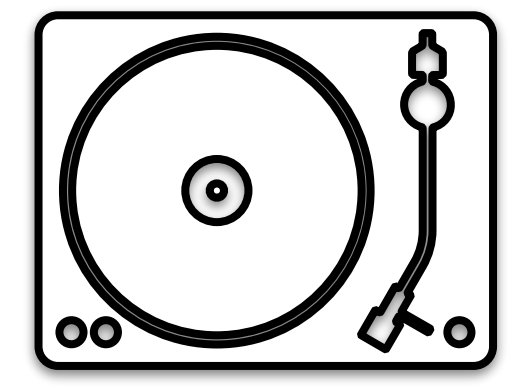
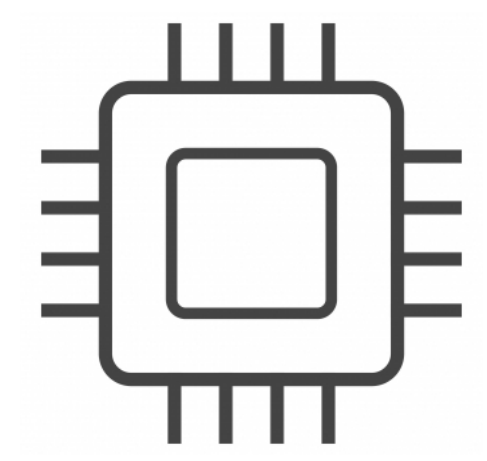
How to manage memory?

Bloom
 filters pointers

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

Cost analysis

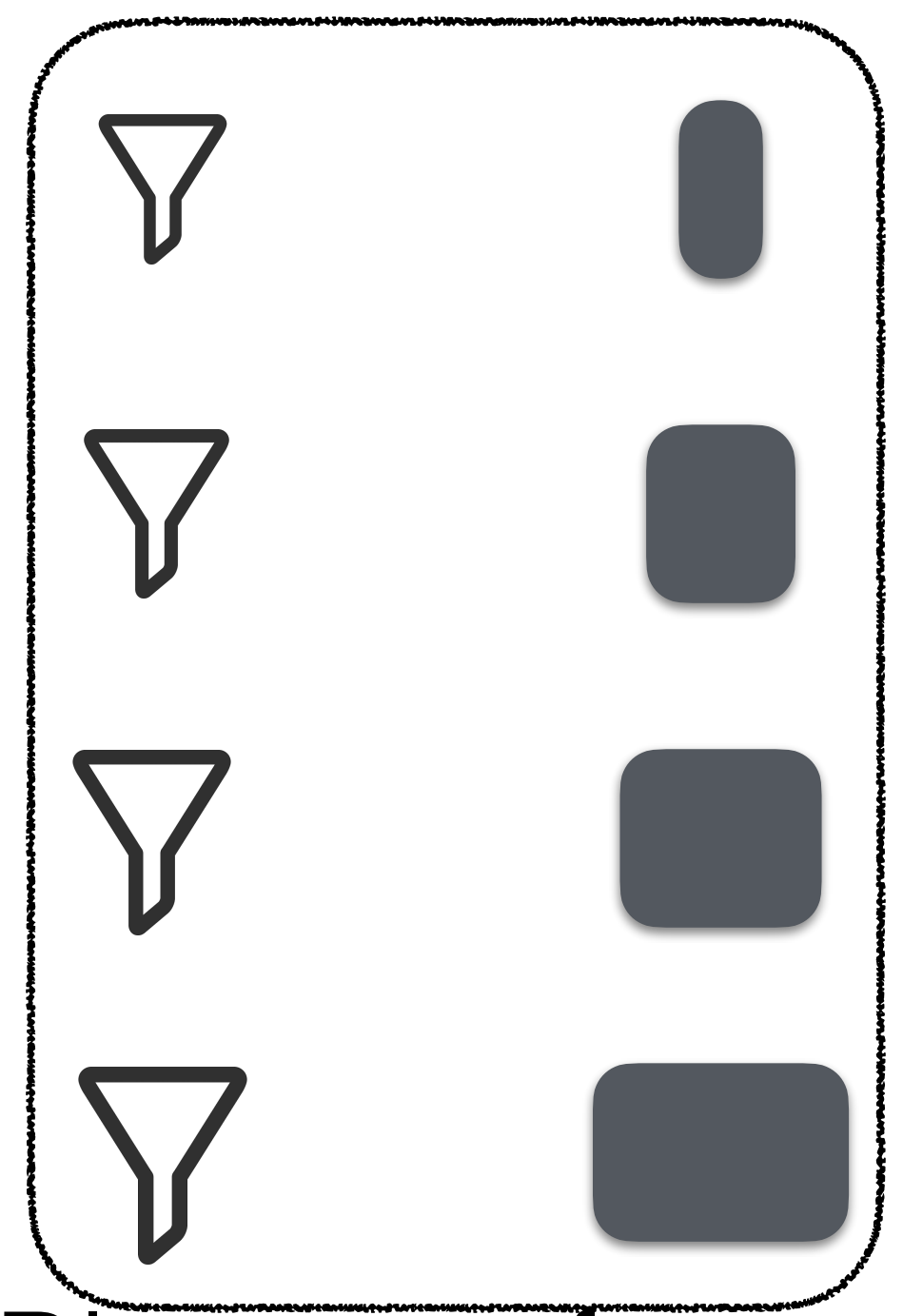
w/o F&I: $\mathcal{O}(\sum \log_2(P \cdot T^i))$
 w/ index: $\mathcal{O}(L)$
 w F&I: $\mathcal{O}(\phi \cdot L)$



buffer



block cache



Bloom filters fence pointers

L1



L2



L3



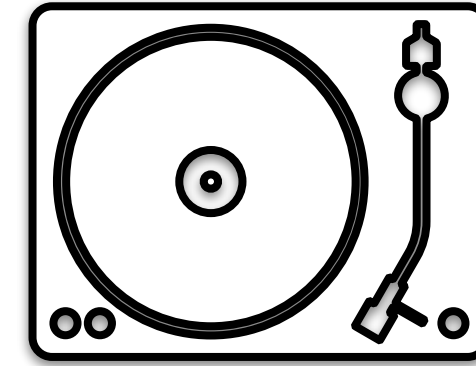
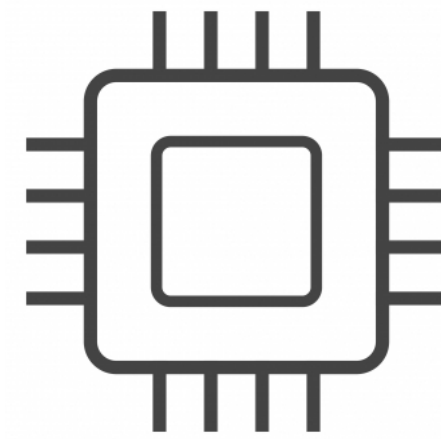
L4



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

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

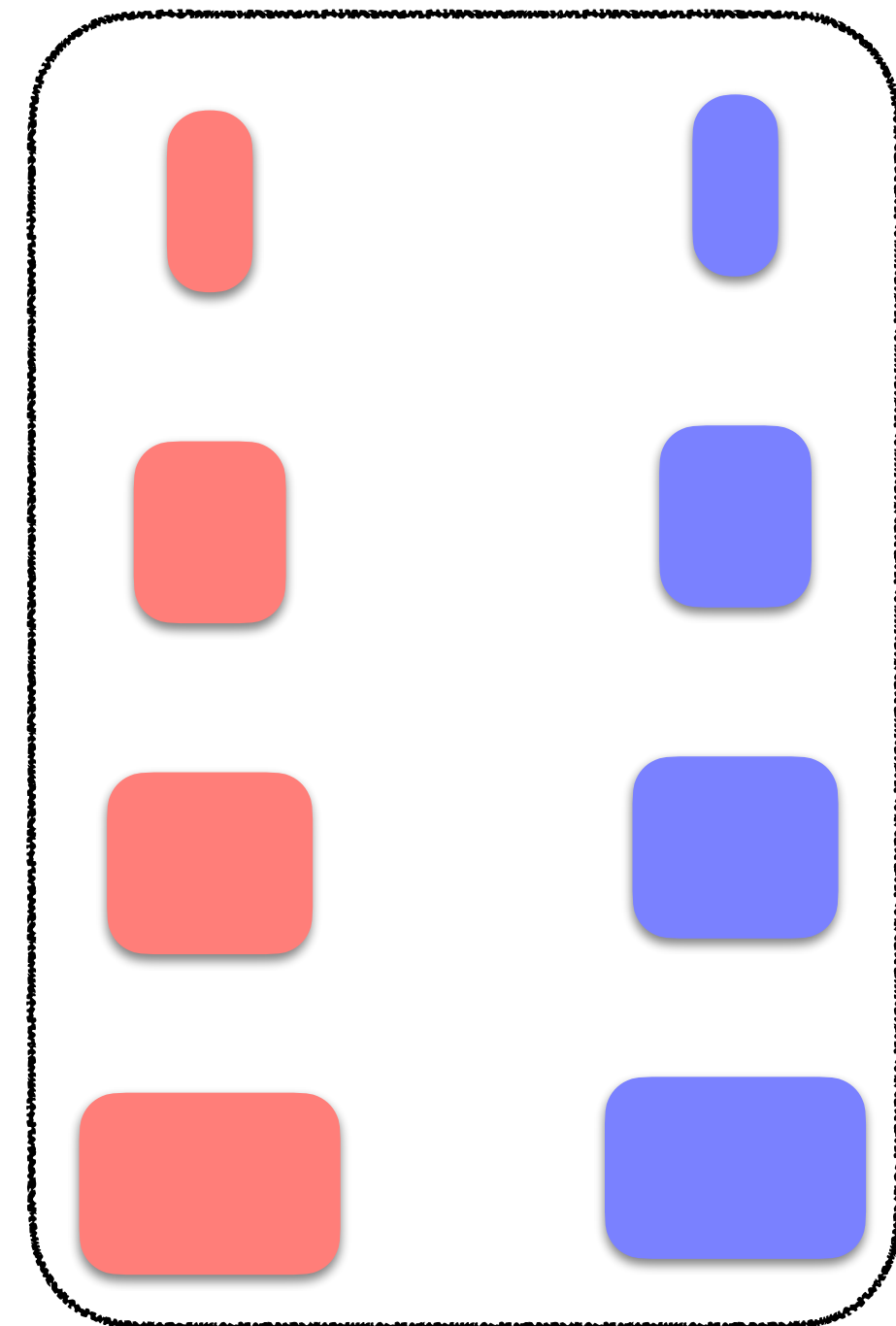
Block Cache



buffer



block cache



Bloom filters

fence pointers

L1



L2



L3

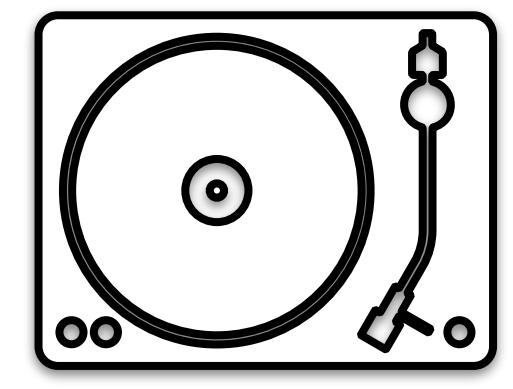


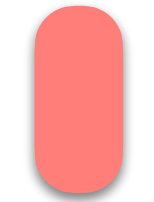

L4



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

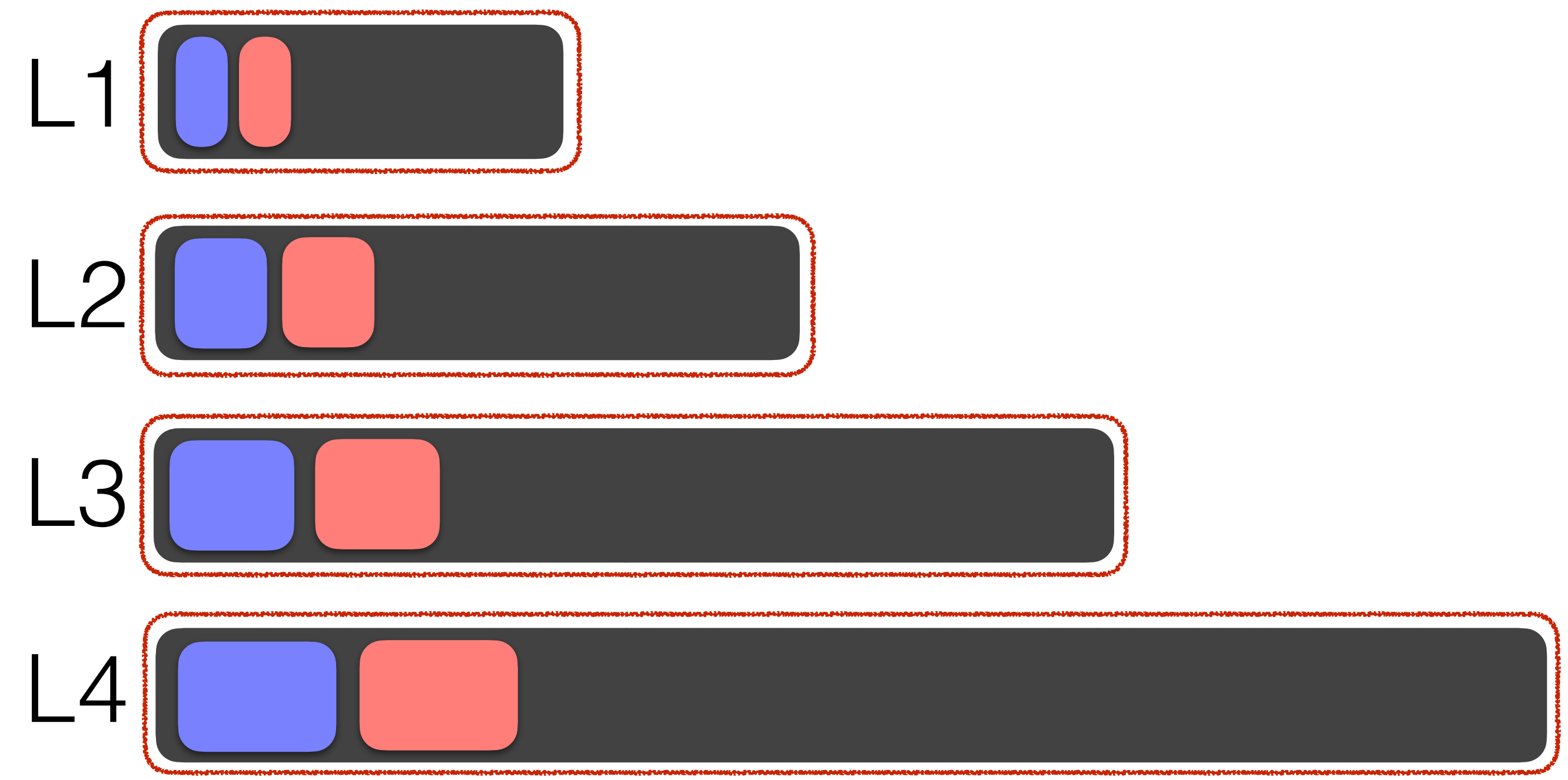
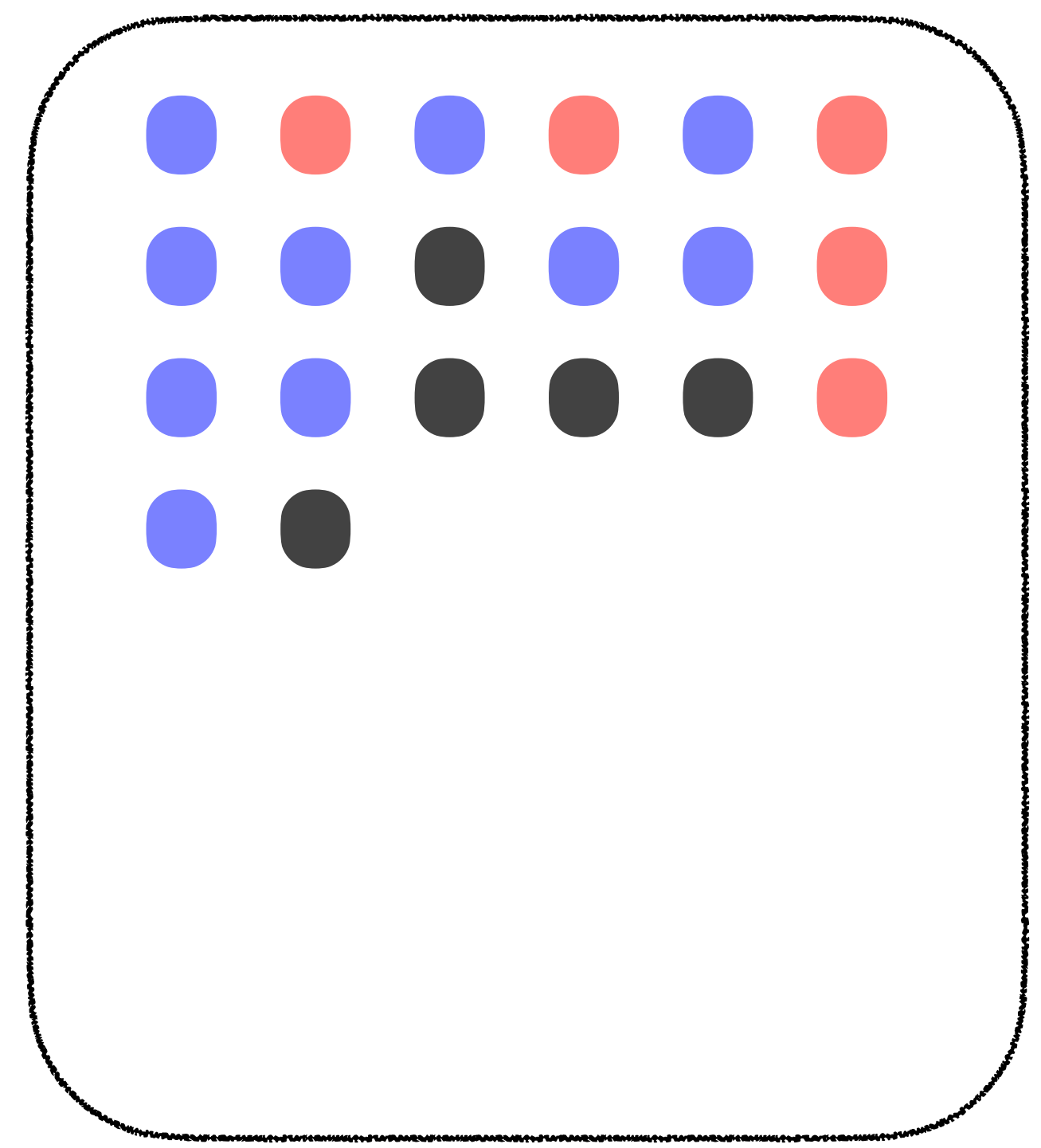
Block Cache



 Bloom filters
 fence pointers

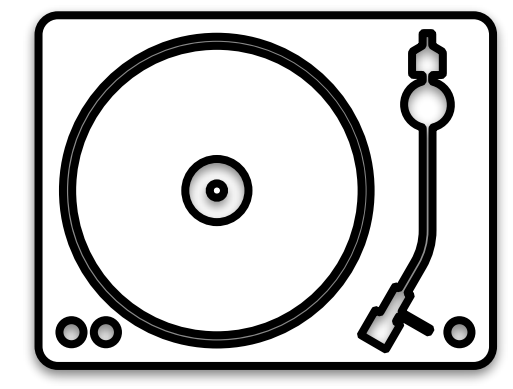
get(7)

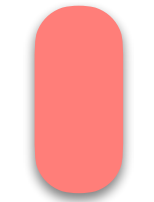
buffer



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

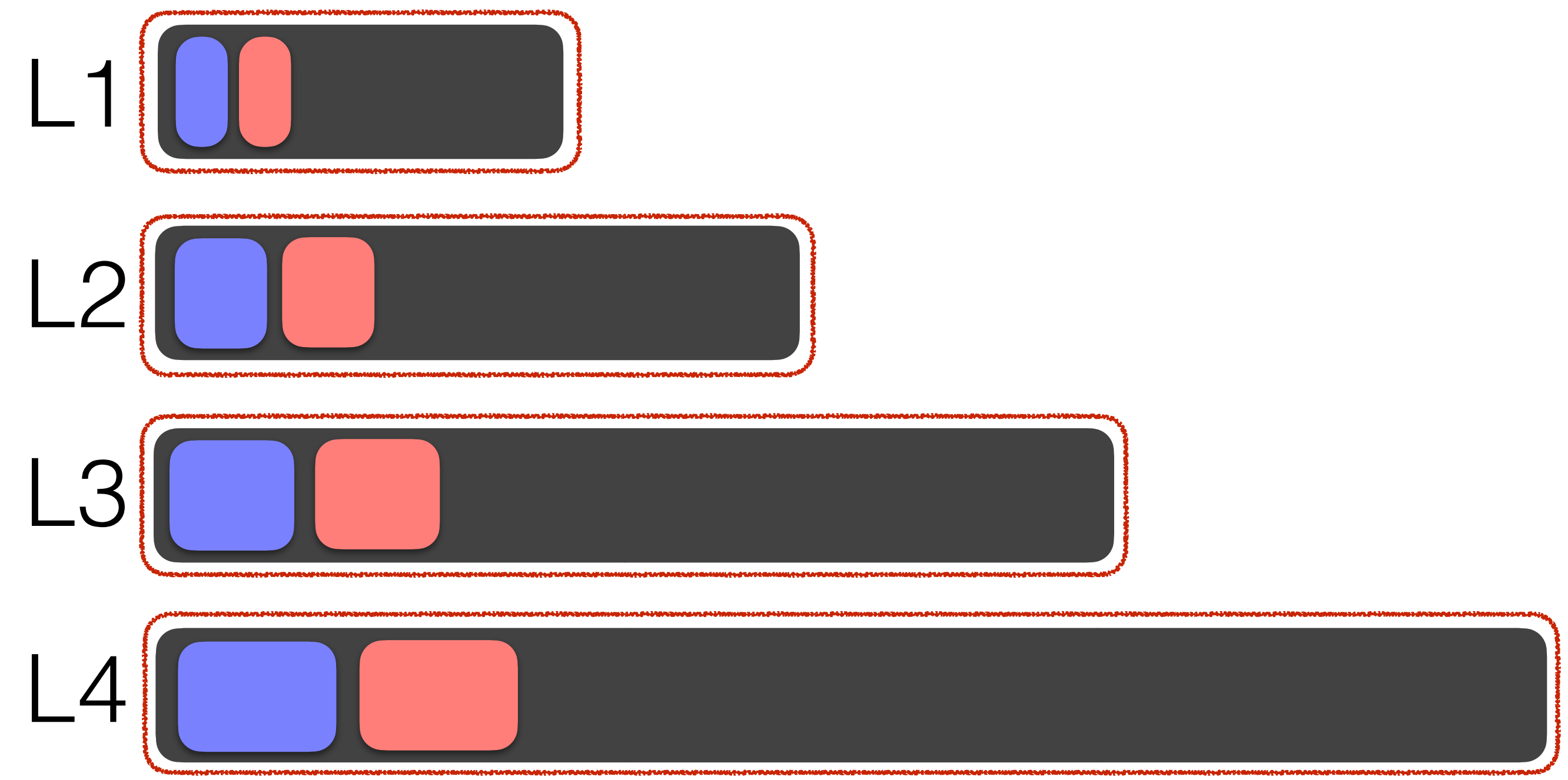
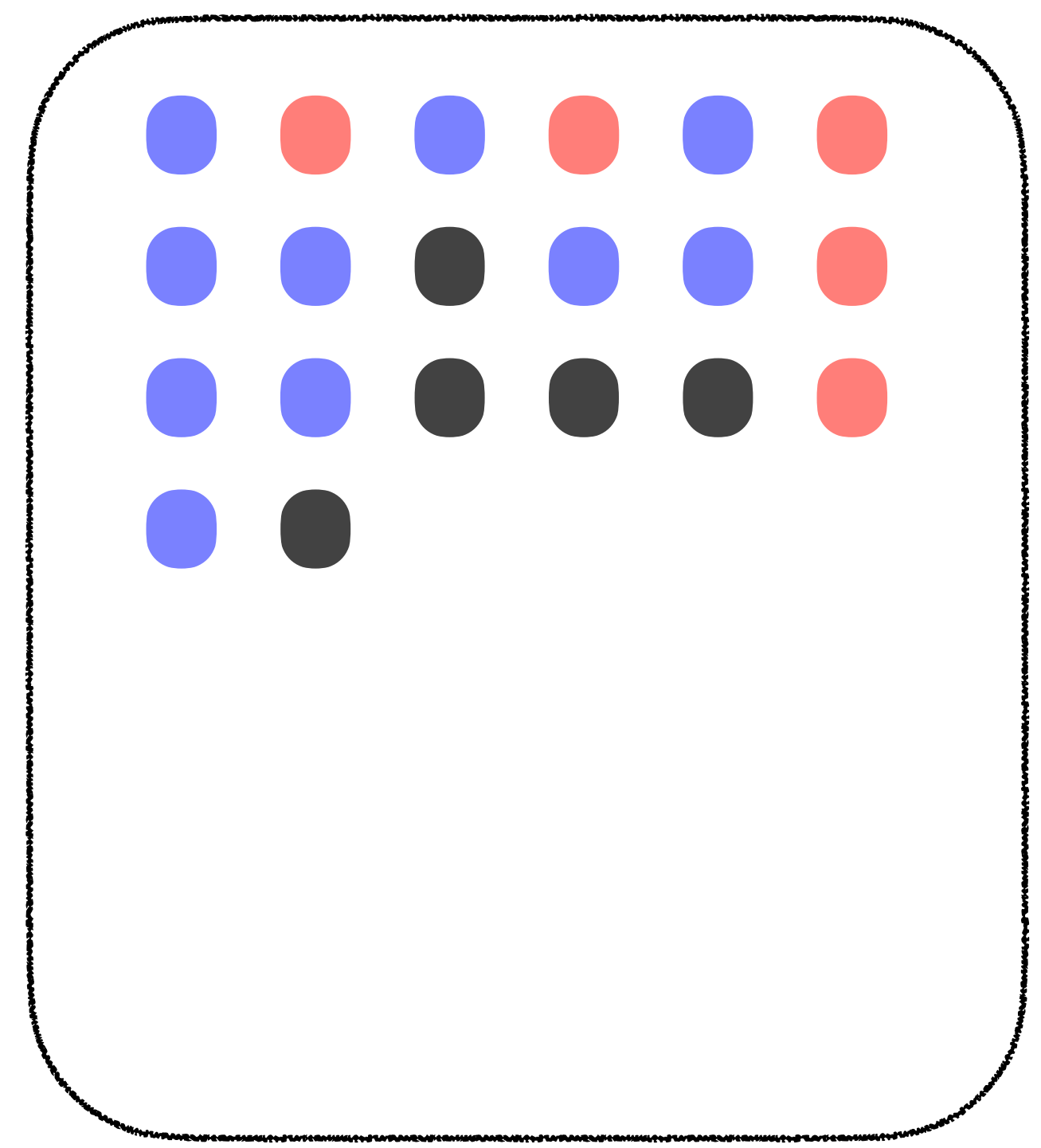
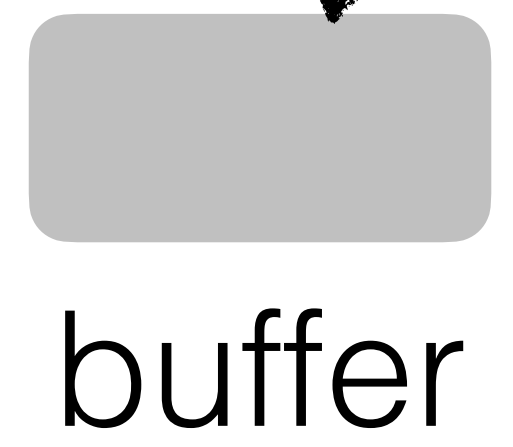
Block Cache




Bloom
filters

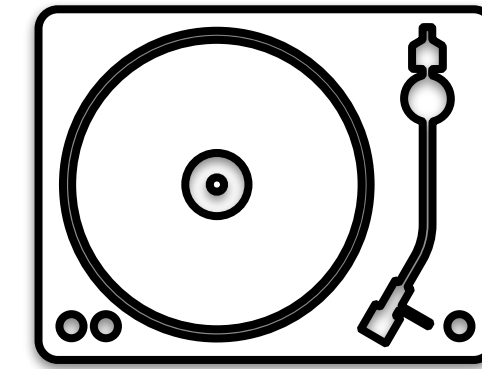

fence
pointers


get(7)




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

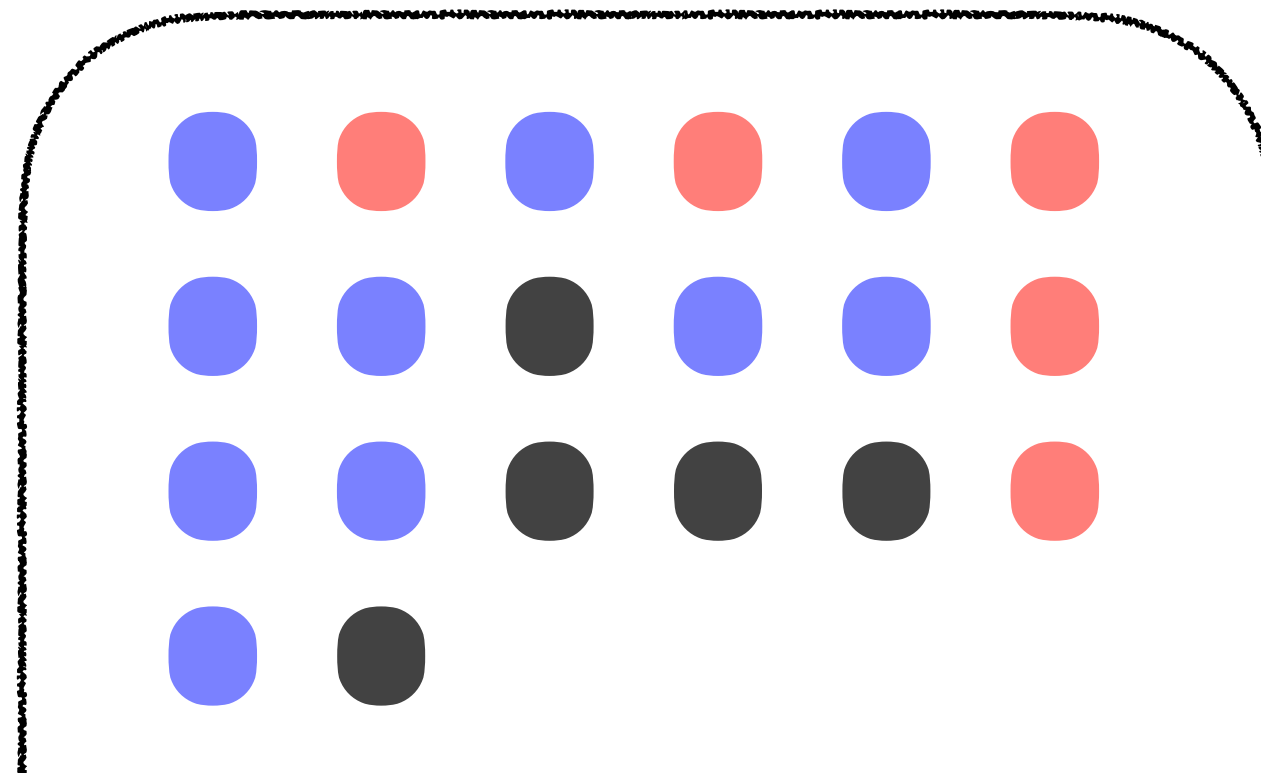
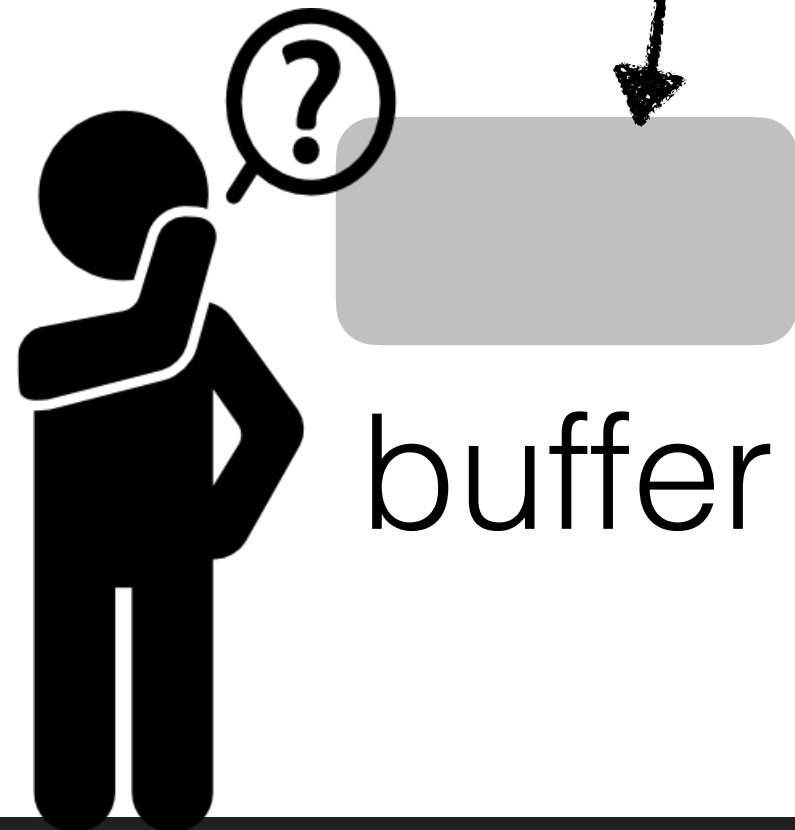
Block Cache




Bloom
filters


fence
pointers

get(7)



L1



L2

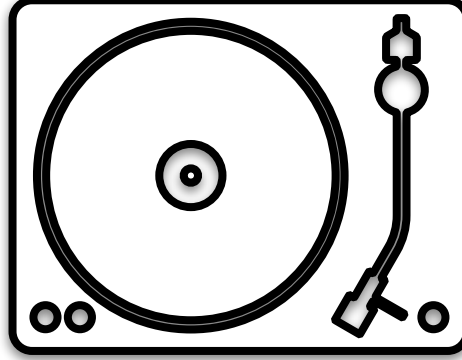
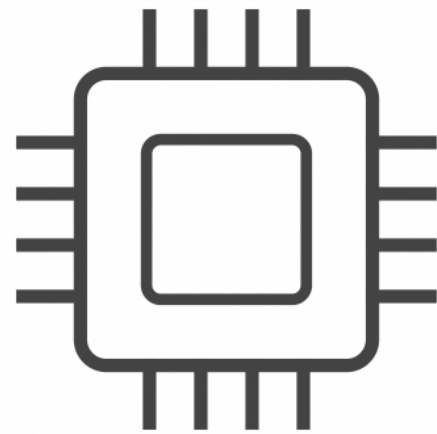


What about range queries?

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

s : selectivity LRQ

Range Queries



buffer



Bloom filters



fence pointers

L1

L2

L3

L4



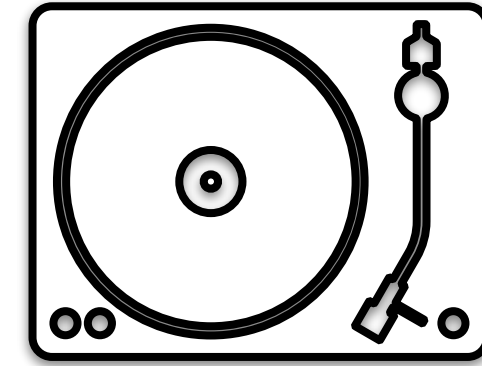
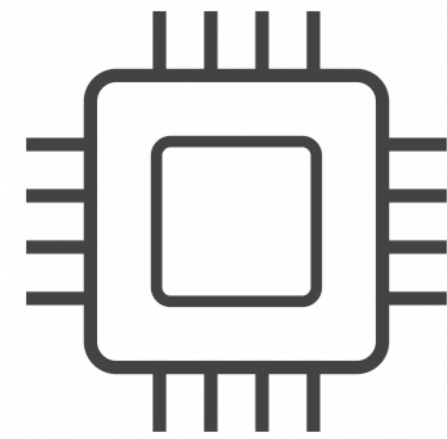
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

s : selectivity LRQ

Range Queries

Cost analysis

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



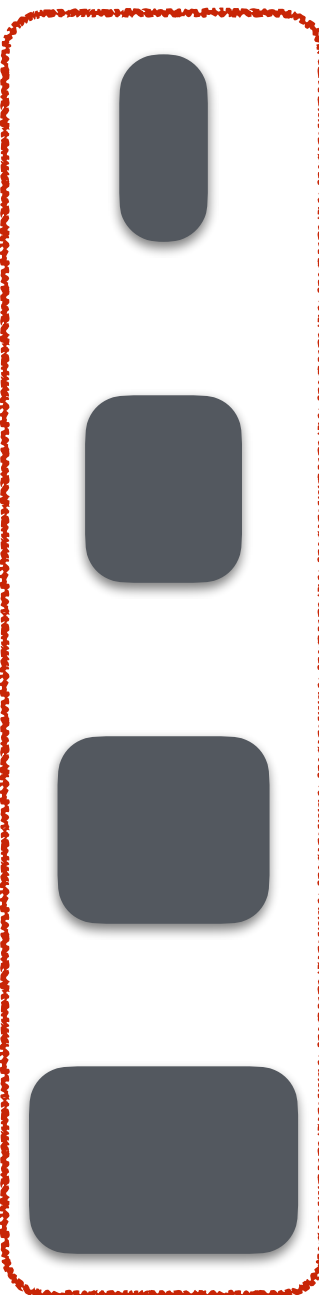
get(9,90)



buffer



Bloom
filters



fence
pointers

L1



L2



L3



L4



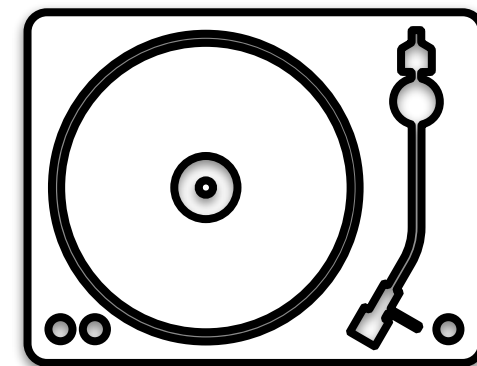
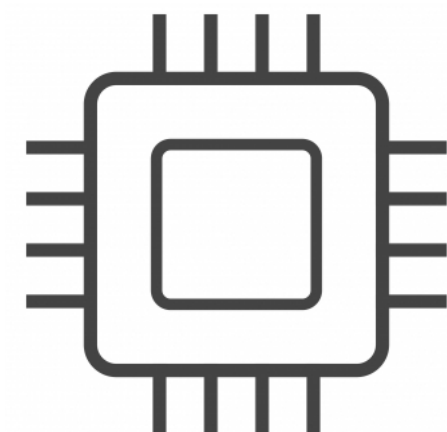
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

s : selectivity SRQ

Range Queries

Cost analysis

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



get(9, 15)



buffer



Bloom filters



fence pointers

L1

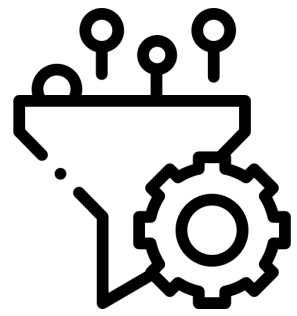
L2

L3

L4



More Read Optimizations



filter structures

Better performance

Chucky

DayanSIGMOD21



Less memory usage

Ribbon filter

DillingerArxiv21



Less CPU usage

Shared Hashing

ZhuDAMON21



Support range queries

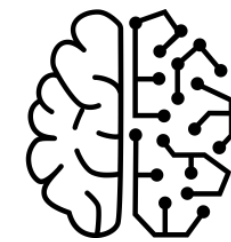
SuRF

ZhangSIGMOD18



Rosetta

LuoSIGMOD20



indexes

Better performance

Learned

DaiOSDI20



Less memory usage

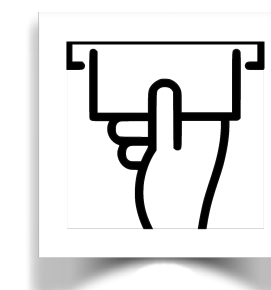
FerraginaVLDB20



block cache

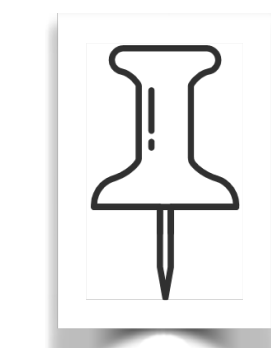
Better performance

Less memory usage



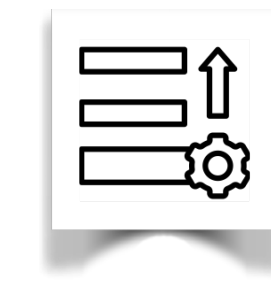
Prefetching

HuangVLDB20



Pinning

CallaghanSD18



Priority

CallaghanSD16

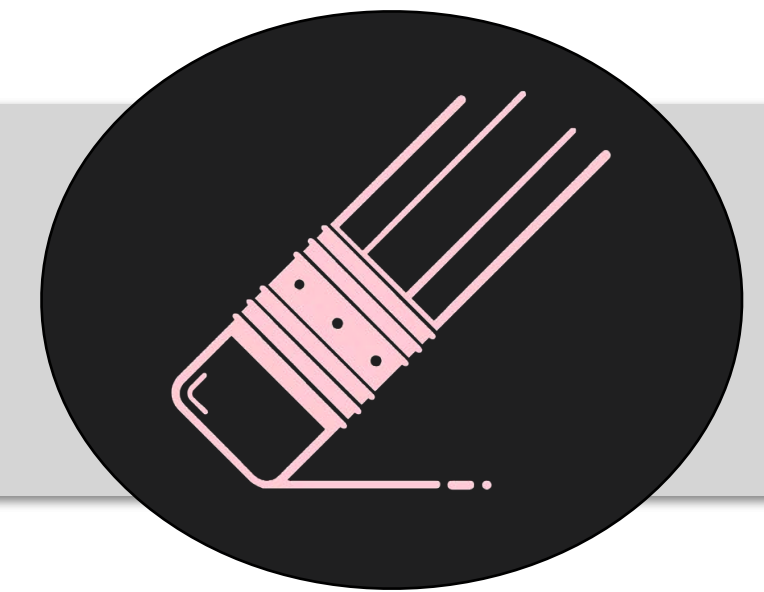


Outline

Part 1: LSM Basics



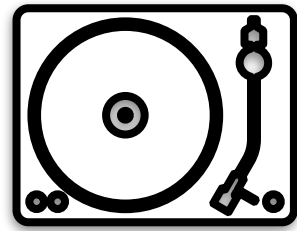
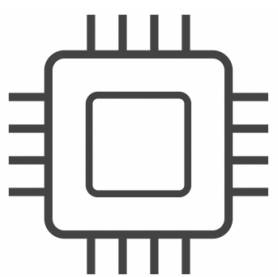
Part 2: **Optimizing Ingestion in LSMs**



Part 3: Navigating the LSM Design Space



Buffer Optimizations



buffer

L1



L2



L3

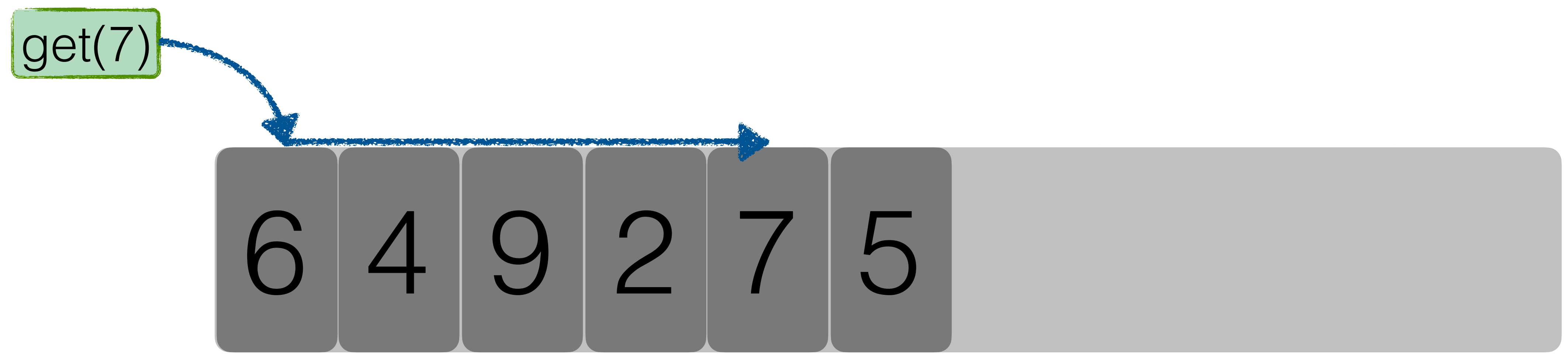


L4



P : pages in buffer
 B : entries/page

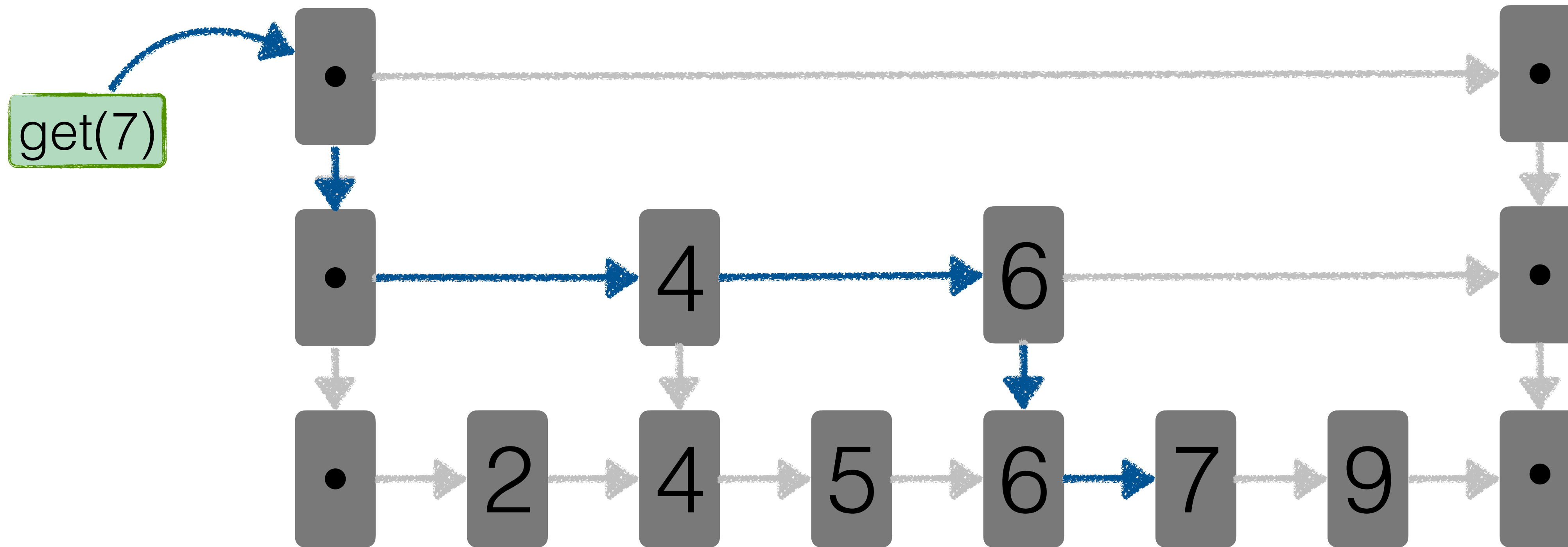
Buffer Implementation: **vector**



- 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

vector

skiplist

hashmap

ingestion
cost

$\mathcal{O}(1)$

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

$\mathcal{O}(1)$

space
complexity

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

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

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

point query
cost

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

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

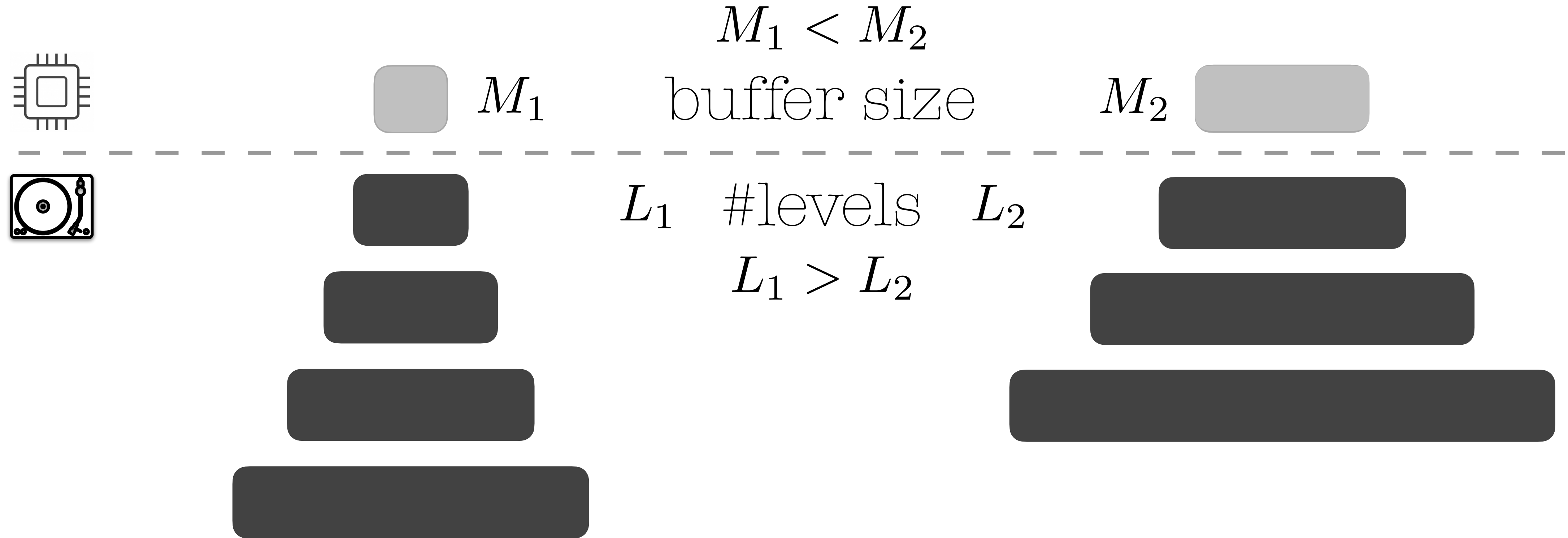
$\mathcal{O}(1)$

Ingestion-only
workloads

Mixed
workloads

I/O-bound
workloads

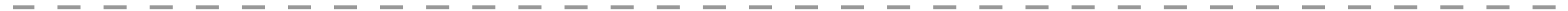
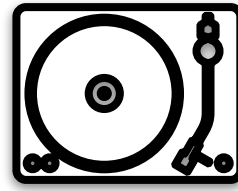
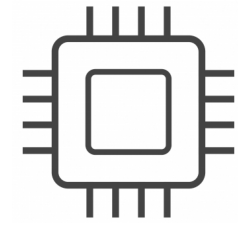
Size of the Buffer



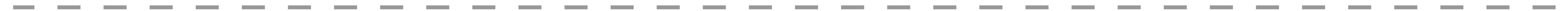
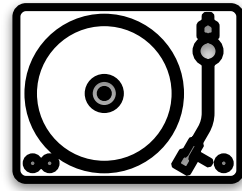
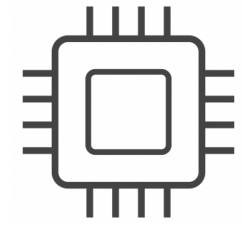
- frequent flushes
- smaller but more levels
- poor read performance

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

#Buffer Components

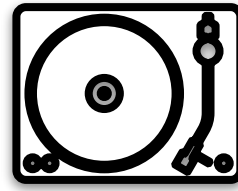
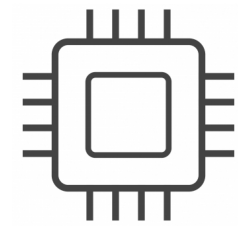


#Buffer Components



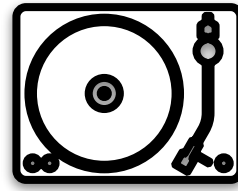
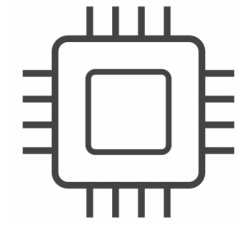
#Buffer Components

immutable
buffers



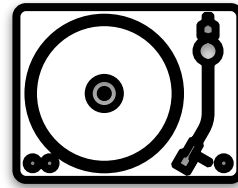
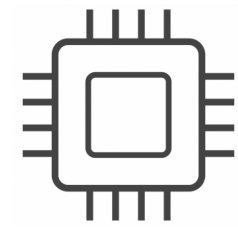
#Buffer Components

immutable
buffers

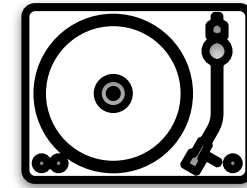
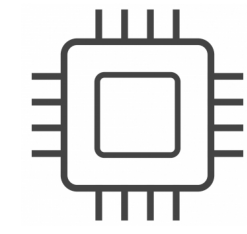


#Buffer Components

immutable
buffers

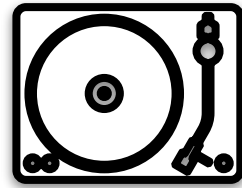
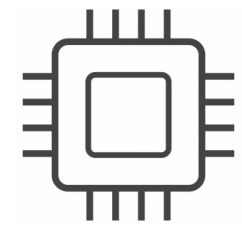


#Buffer Components

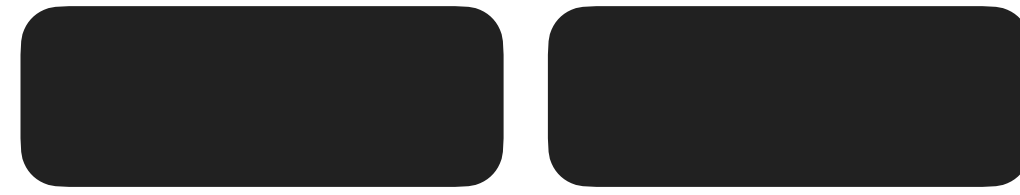


immutable
buffers

#Buffer Components

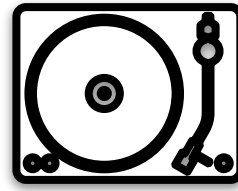
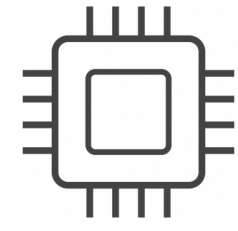


immutable
buffers

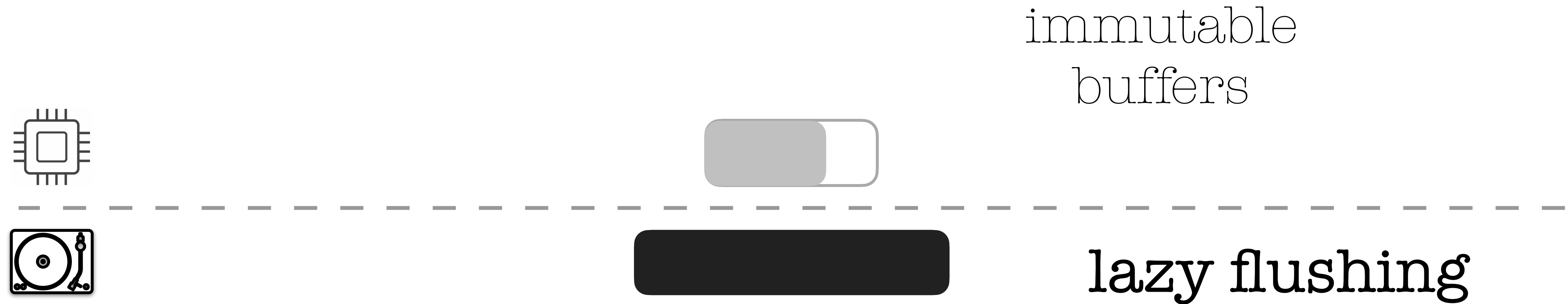


#Buffer Components

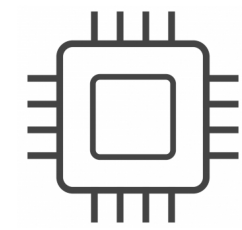
immutable
buffers



#Buffer Components



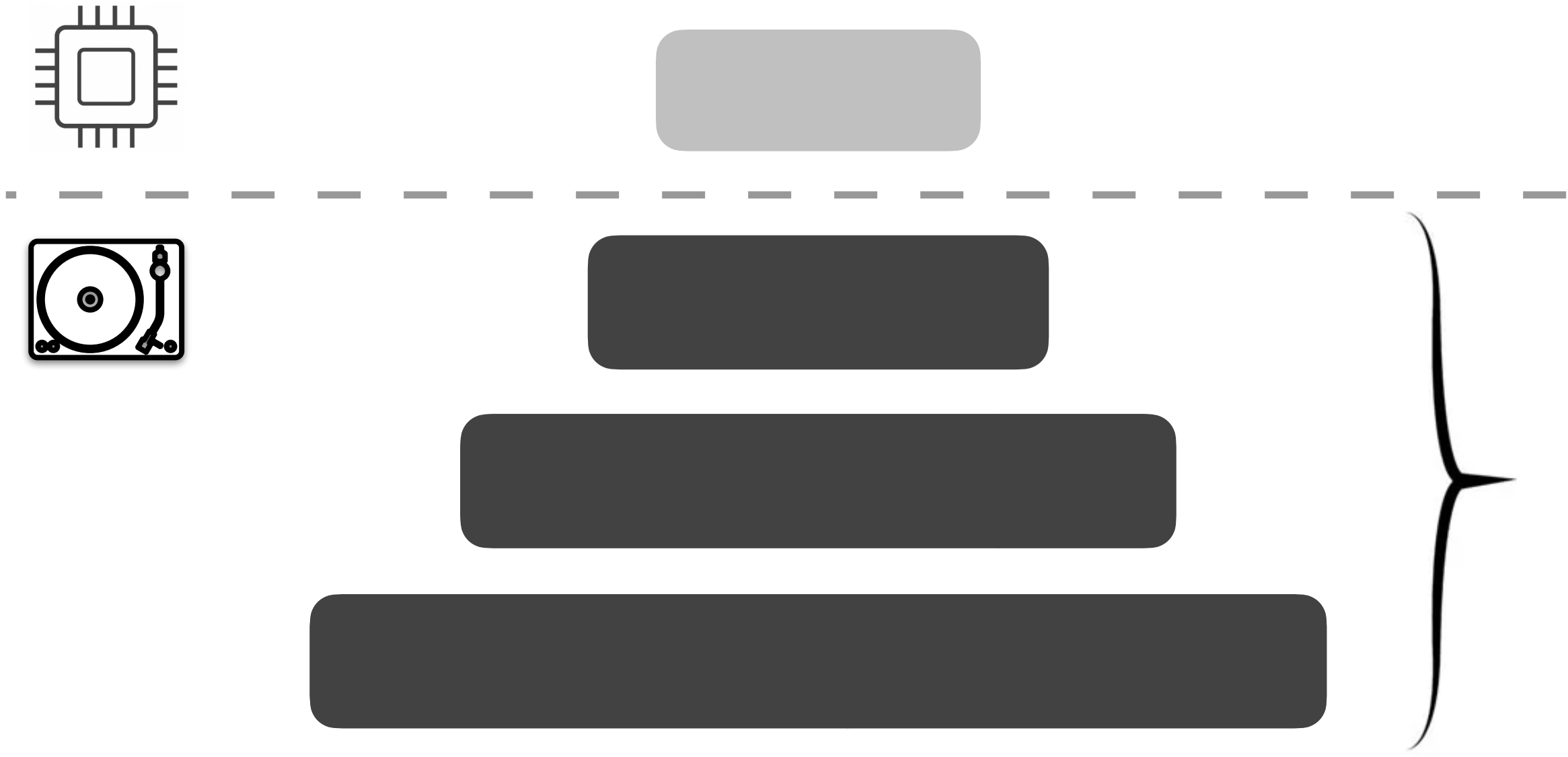
- avoid write stalls
- improved ingestion throughput
- better bandwidth utilization
- requires more memory



most data
on storage

How does the storage layer affect ingestion?

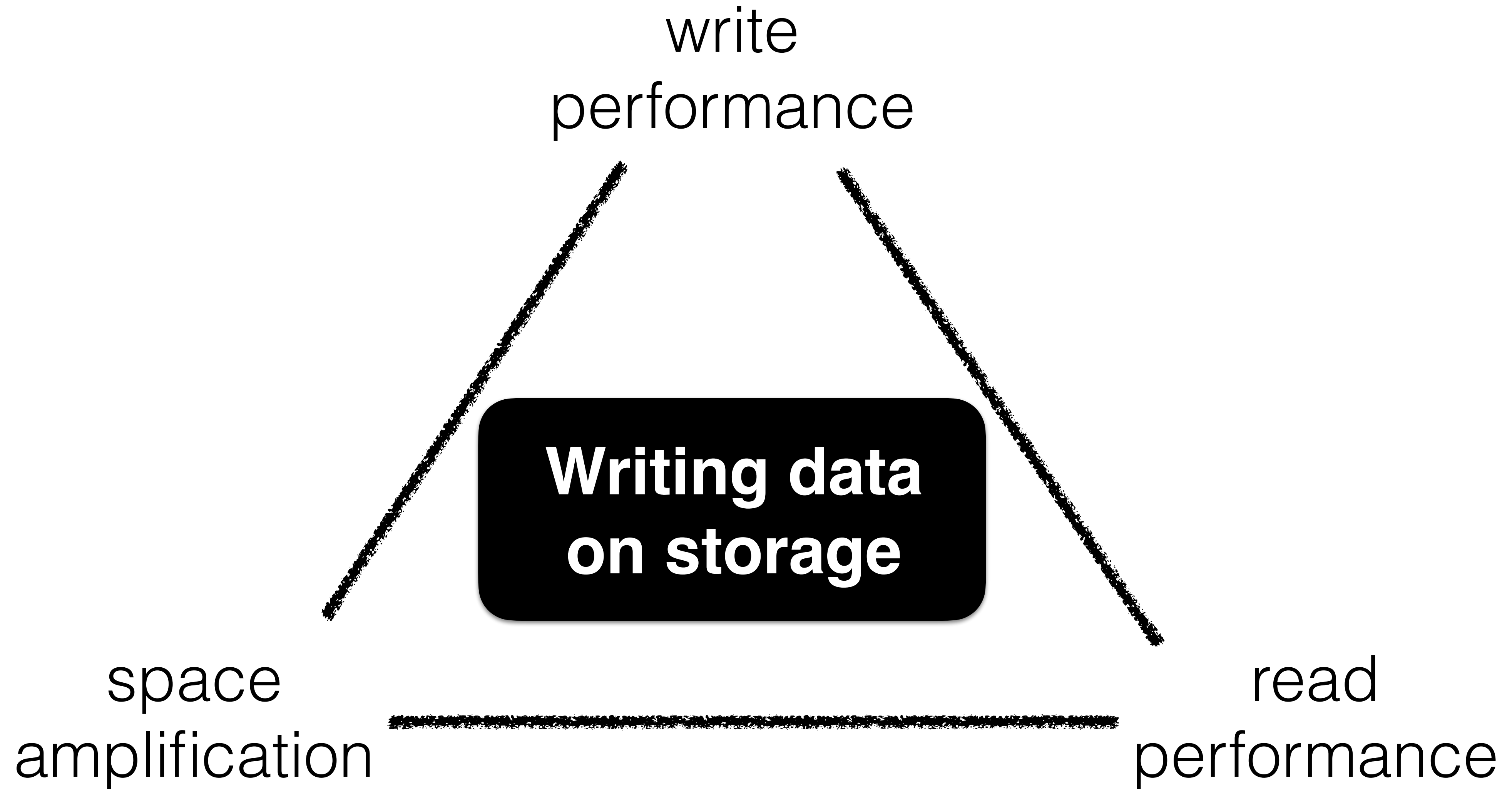
L : #levels
 T : size ratio



most data
on storage
if $T = 10$ & $L = 4$
99.9% on storage

How does the storage layer affect ingestion?

Storage Optimizations



Data **Layout**

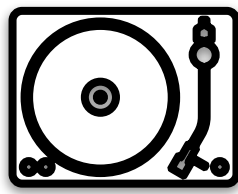
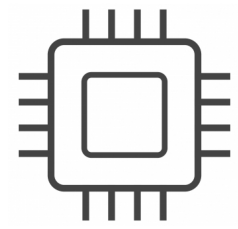
Classical LSM design: **leveling**

[eager merging]

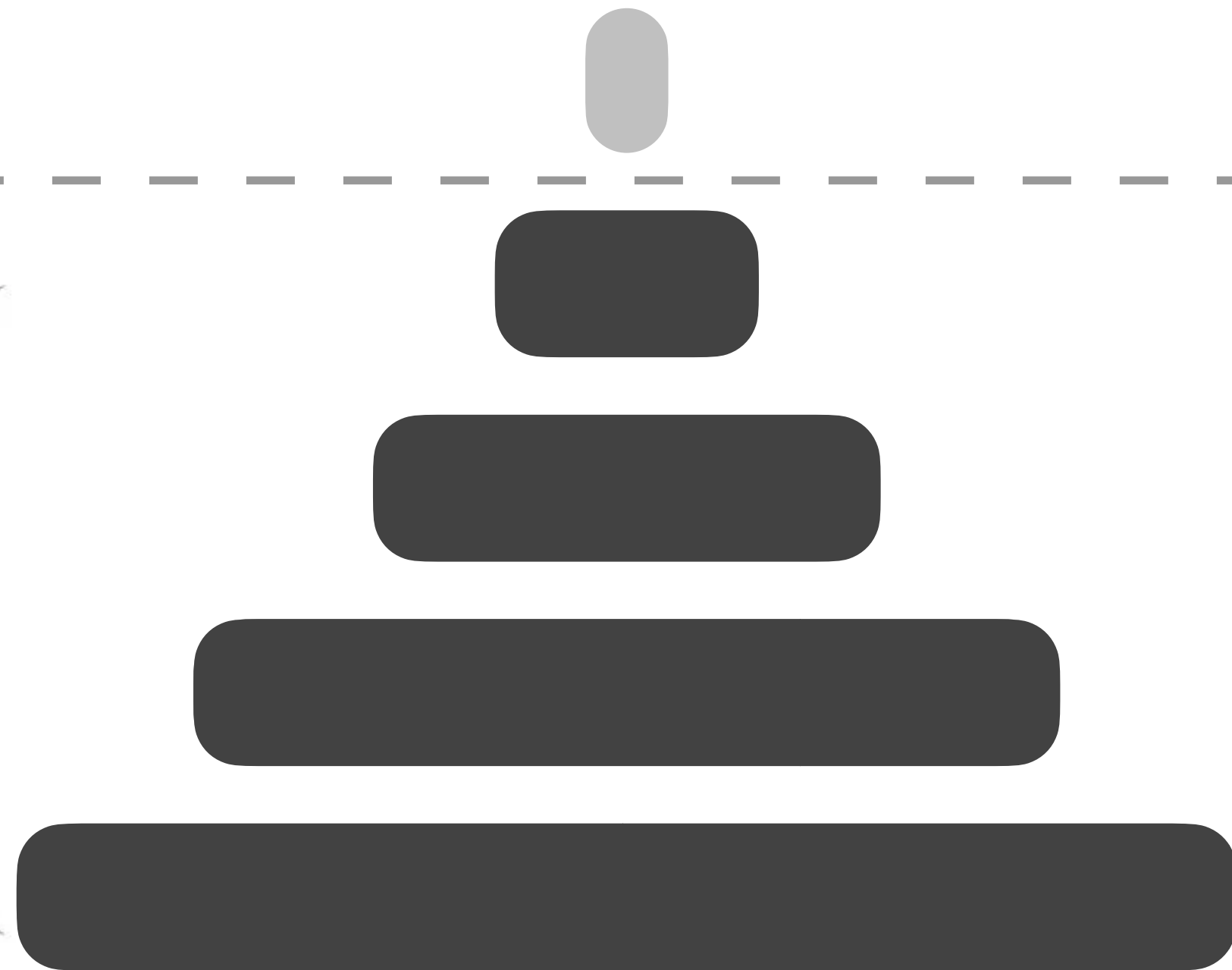


Data Layout

leveling [eager]



1 run
per level

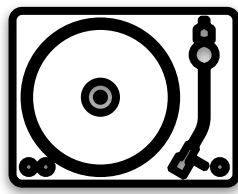
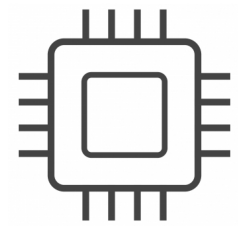


- good read performance
- good space amplification
- high write amplification

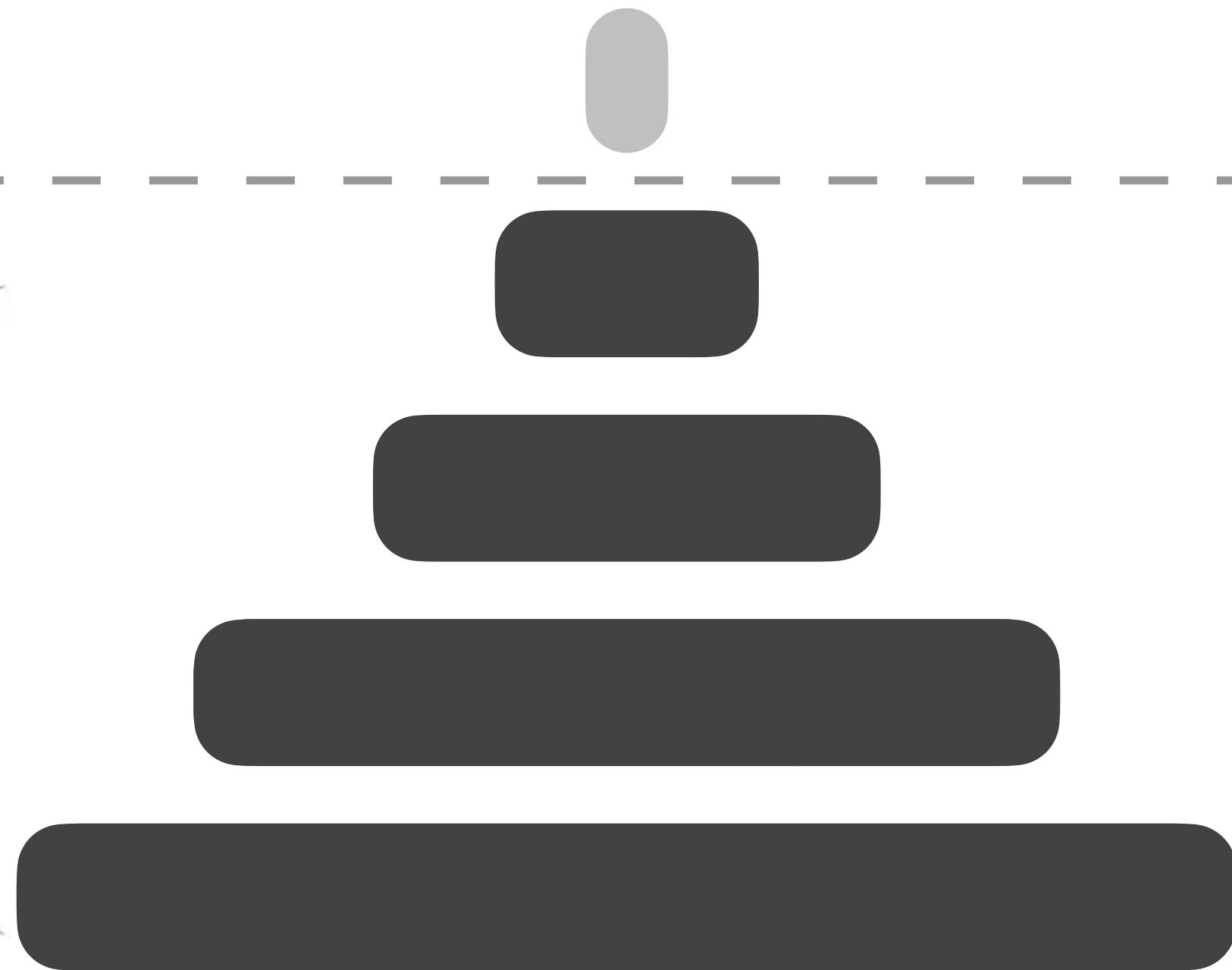
Data Layout

leveling [eager]

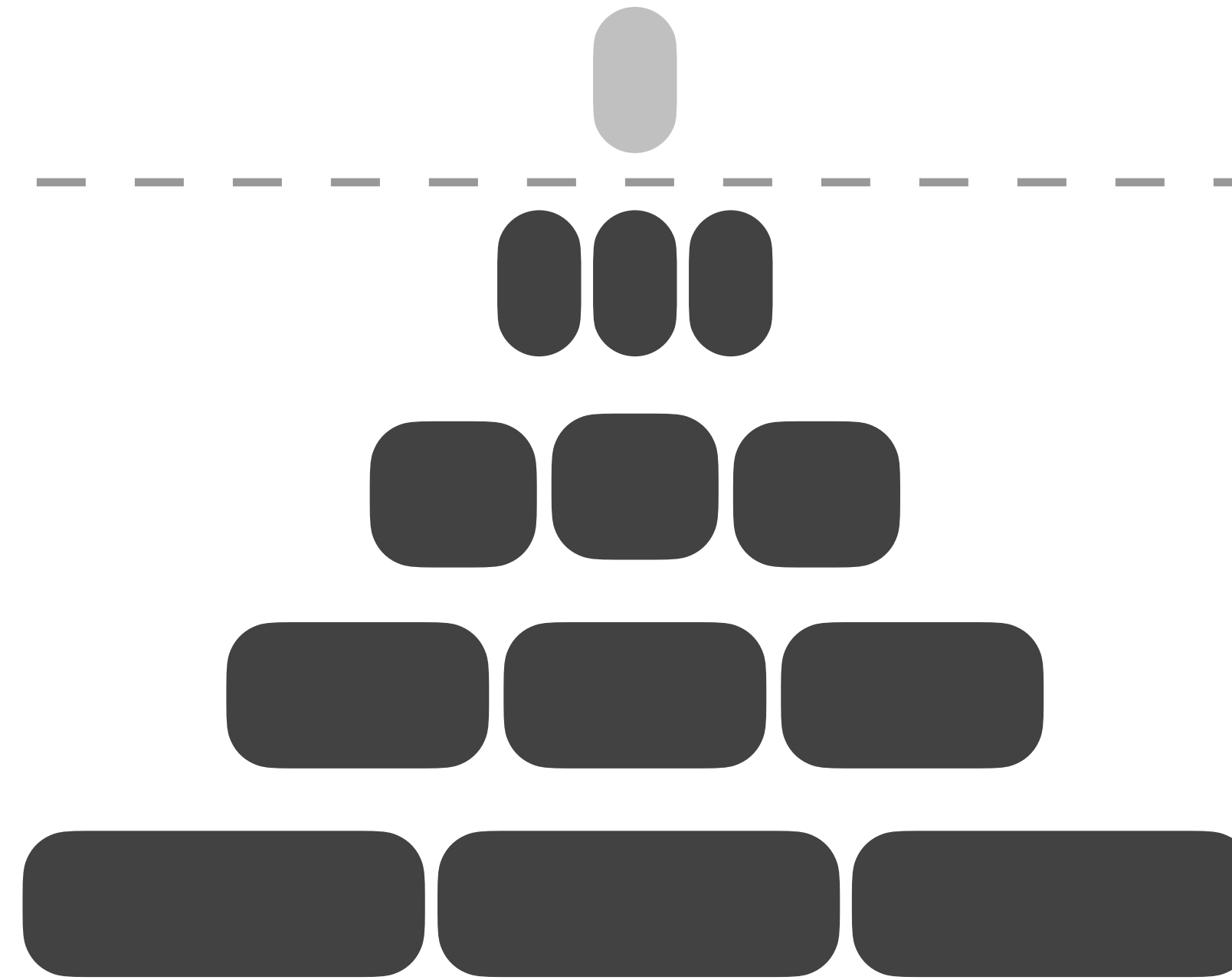
tiering [lazy]



1 run
per level



T runs
per level



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

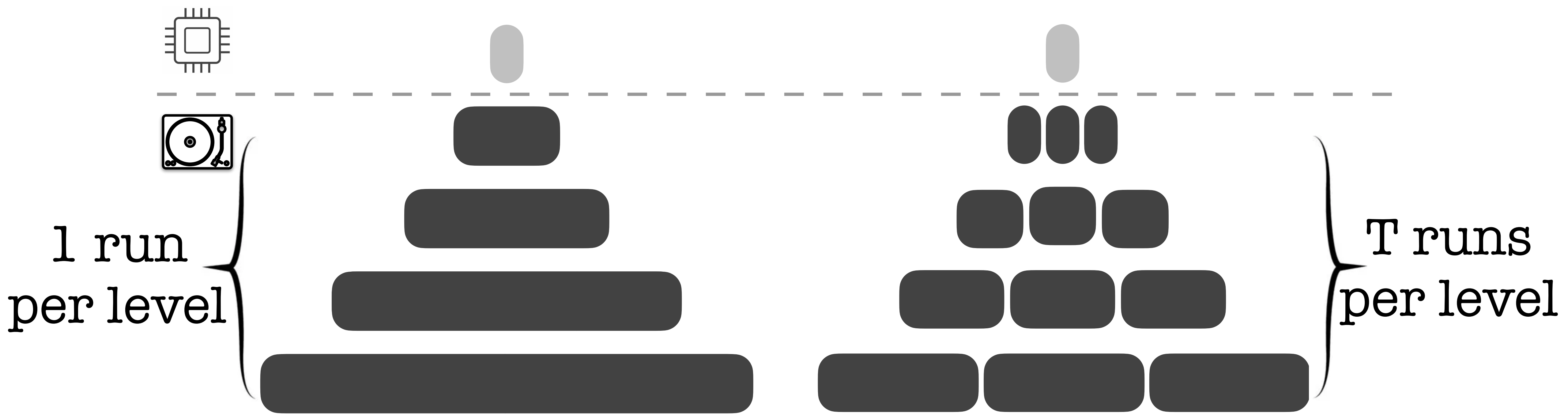
- poor read performance
- poor space amplification
- good ingestion performance

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

Data Layout

leveling [eager]

tiering [lazy]



1 run
per level

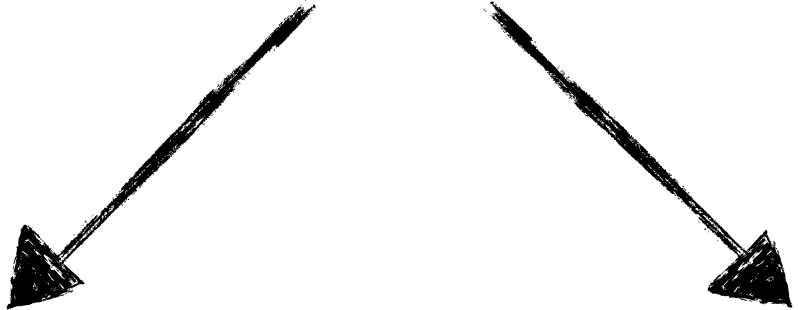
T runs
per level

Read cost: $\mathcal{O}(L \cdot \phi)$
 Write cost: $\mathcal{O}(T \cdot L/B)$
 SA: $\mathcal{O}(1/T)$

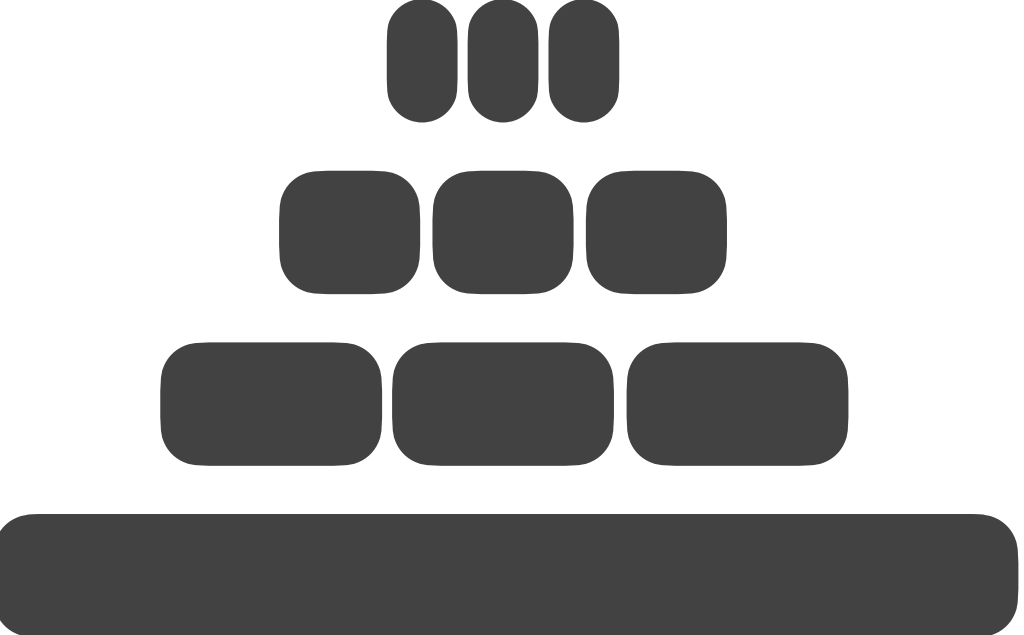
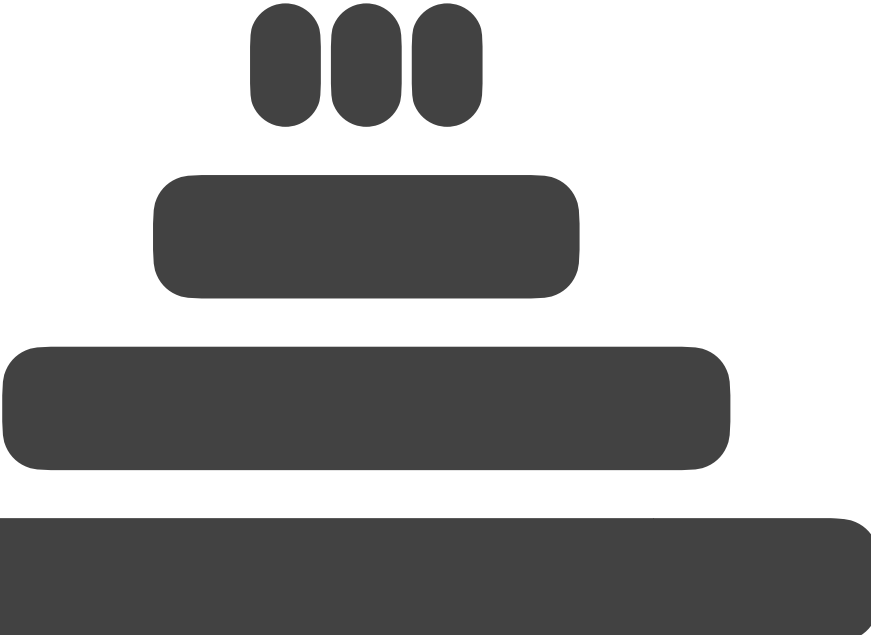
$\mathcal{O}(T \cdot L \cdot \phi)$
 $\mathcal{O}(L/B)$
 $\mathcal{O}(T)$

Data Layout

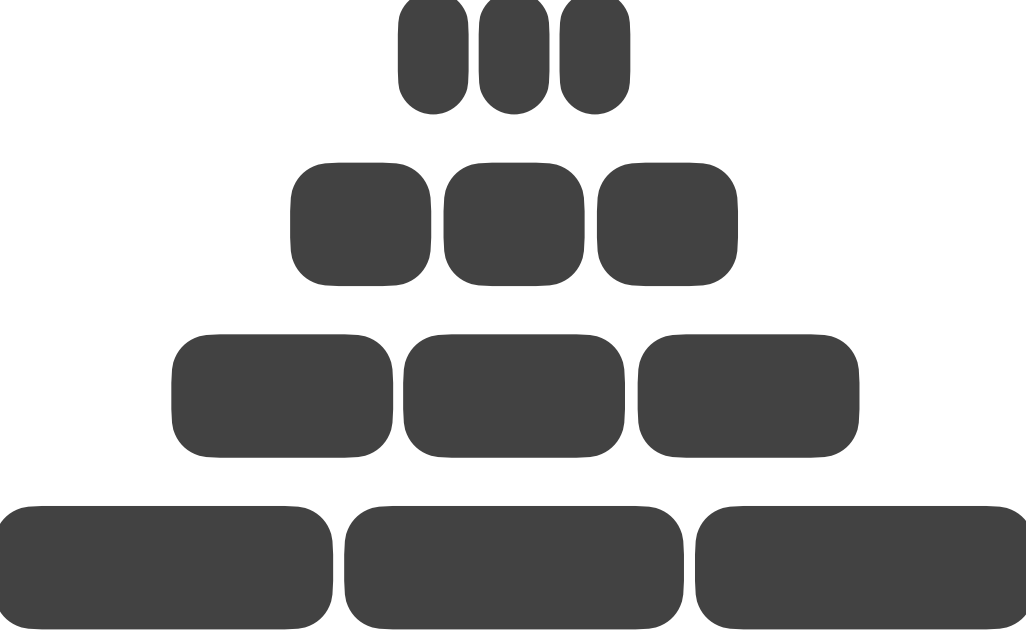
hybrid designs



leveling



tiering



read
optimized

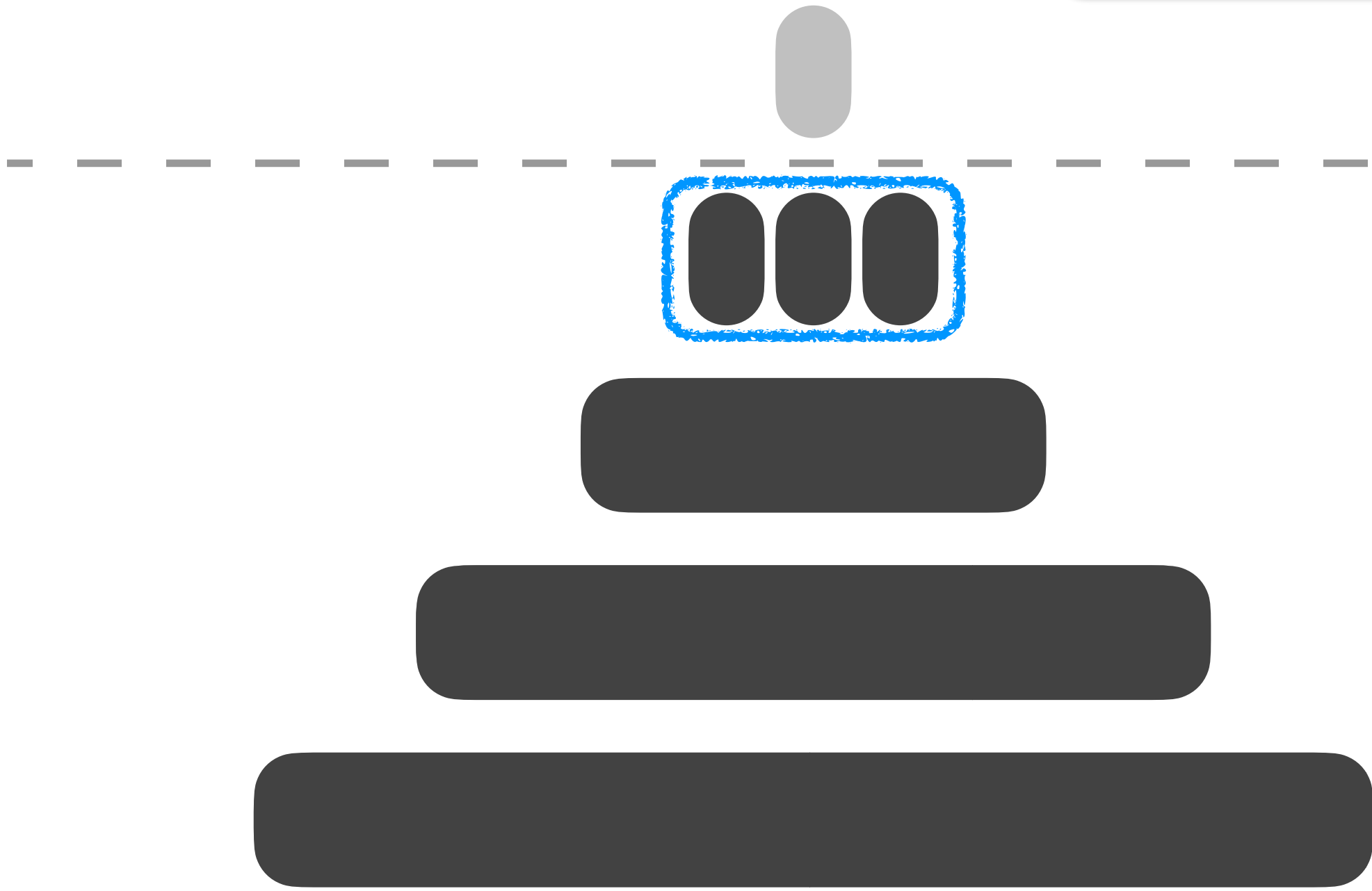


write
optimized

Data Layout

1-leveling

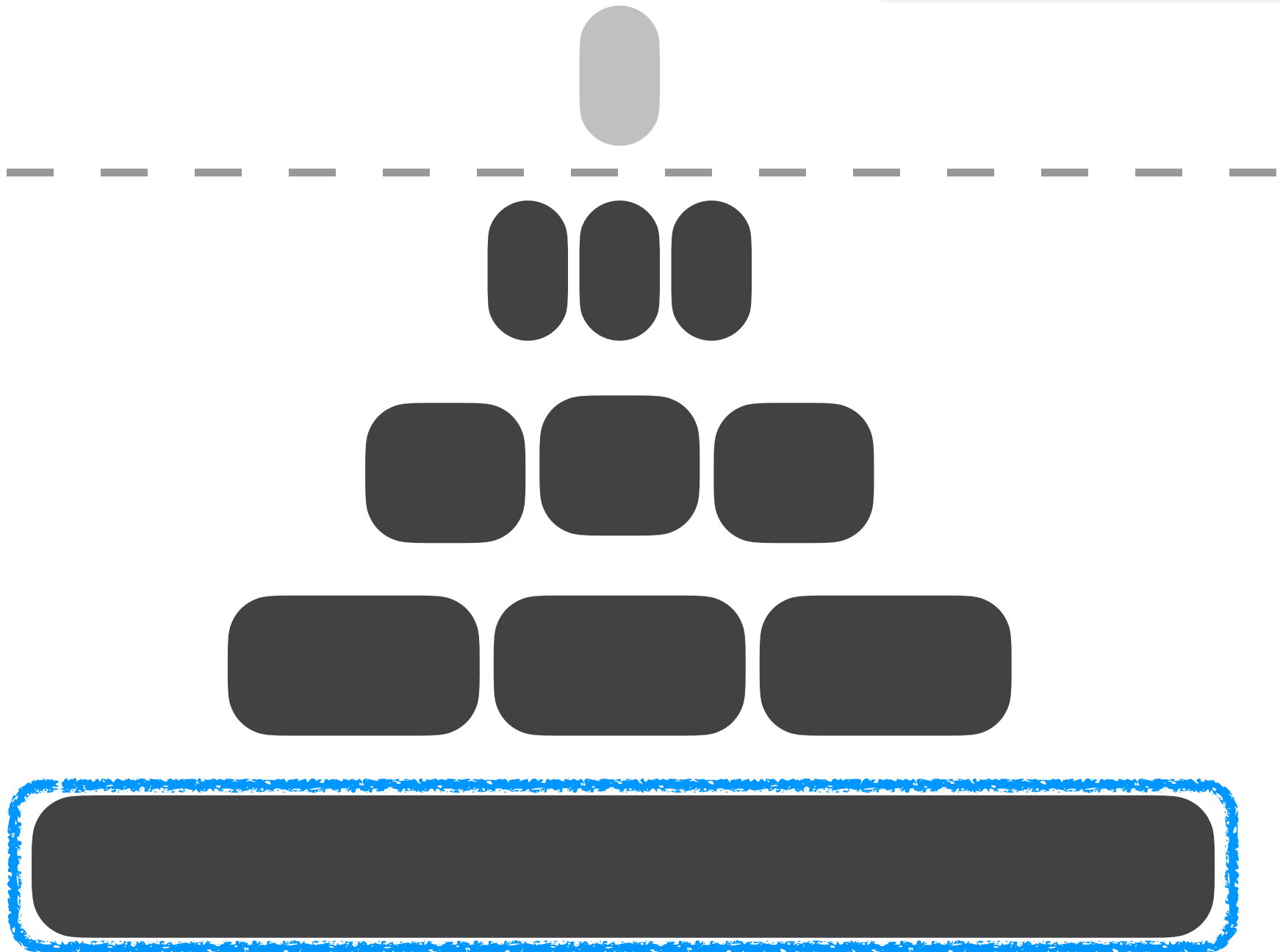
RocksDB20 



- fewer write stalls
- increased block cache hits

L-leveling

DayanSIGMOD18 



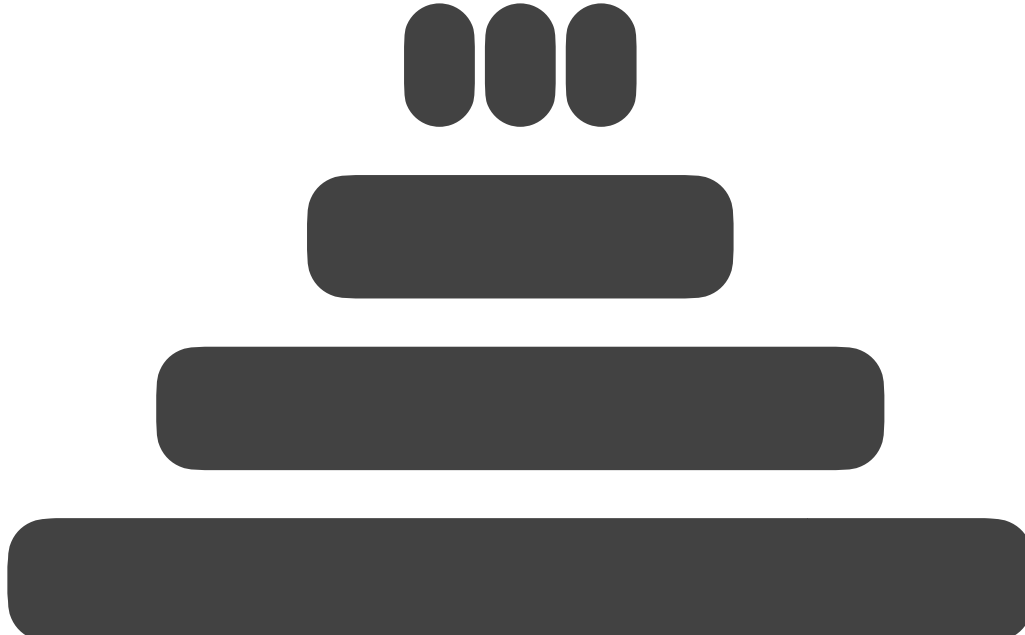
- low write amplification
- better read performance

Data Layout

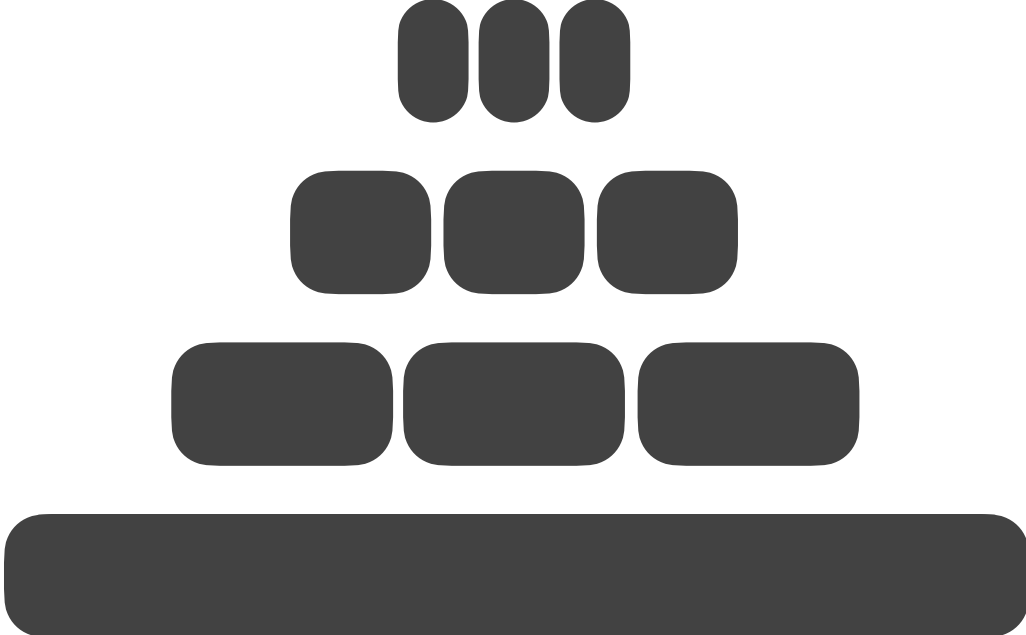
leveling



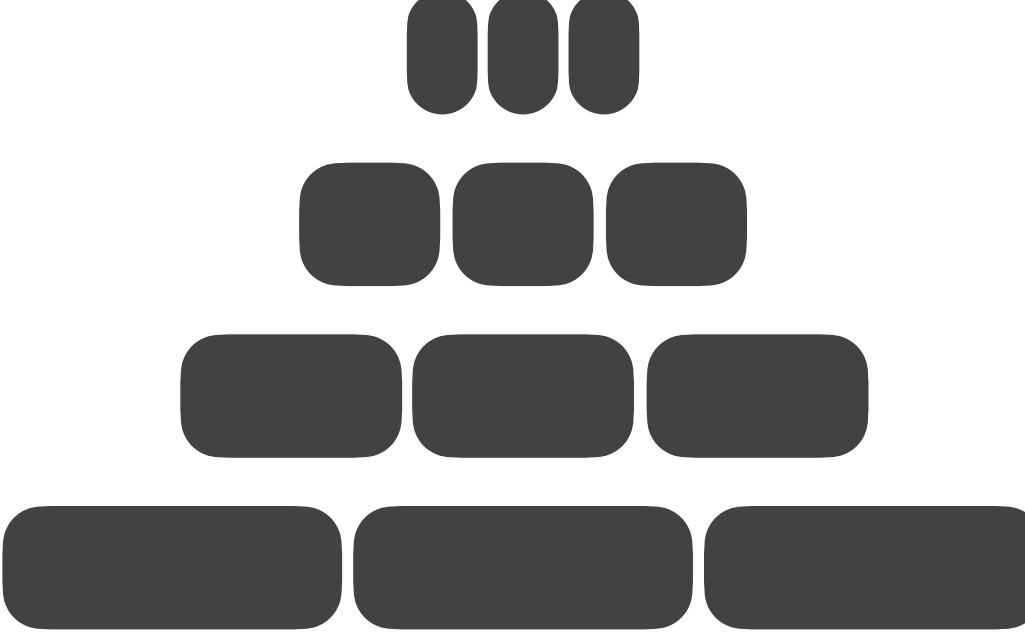
1-leveling



L-leveling



tiering



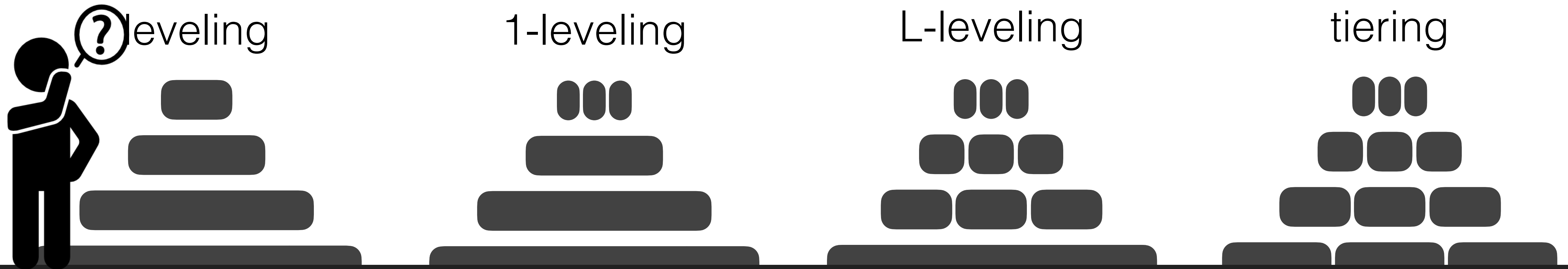
read
optimized



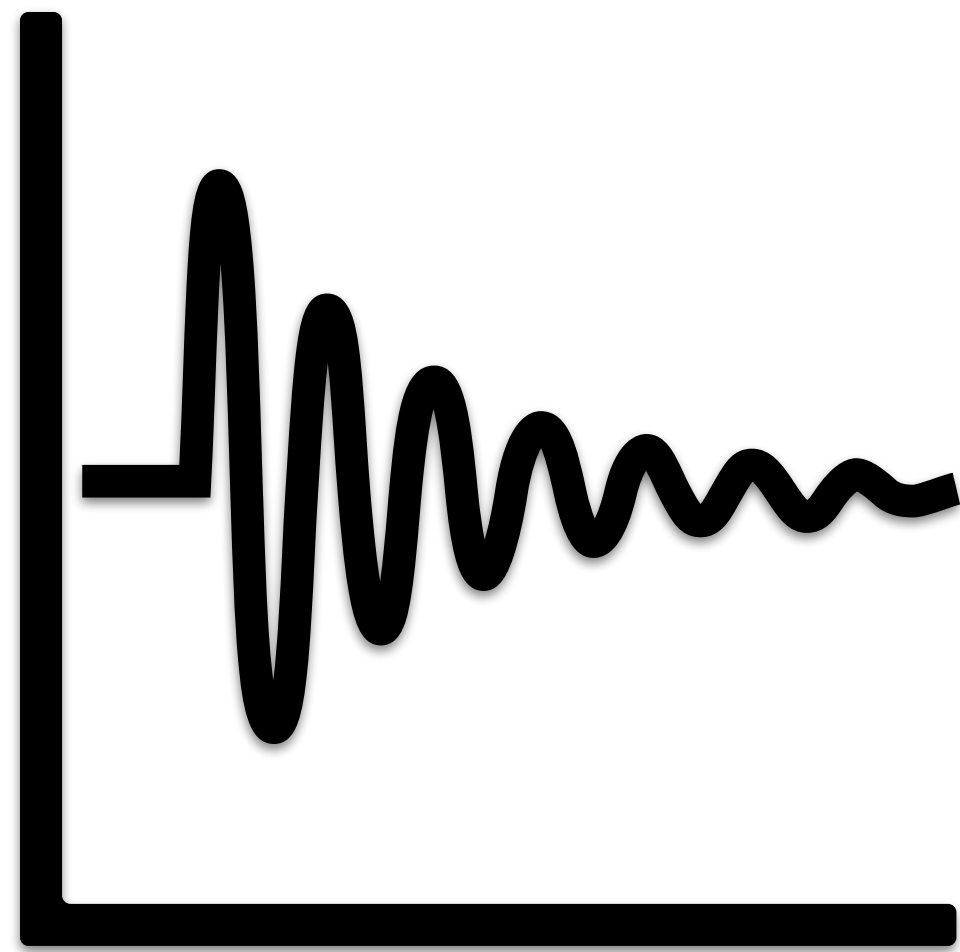
write
optimized



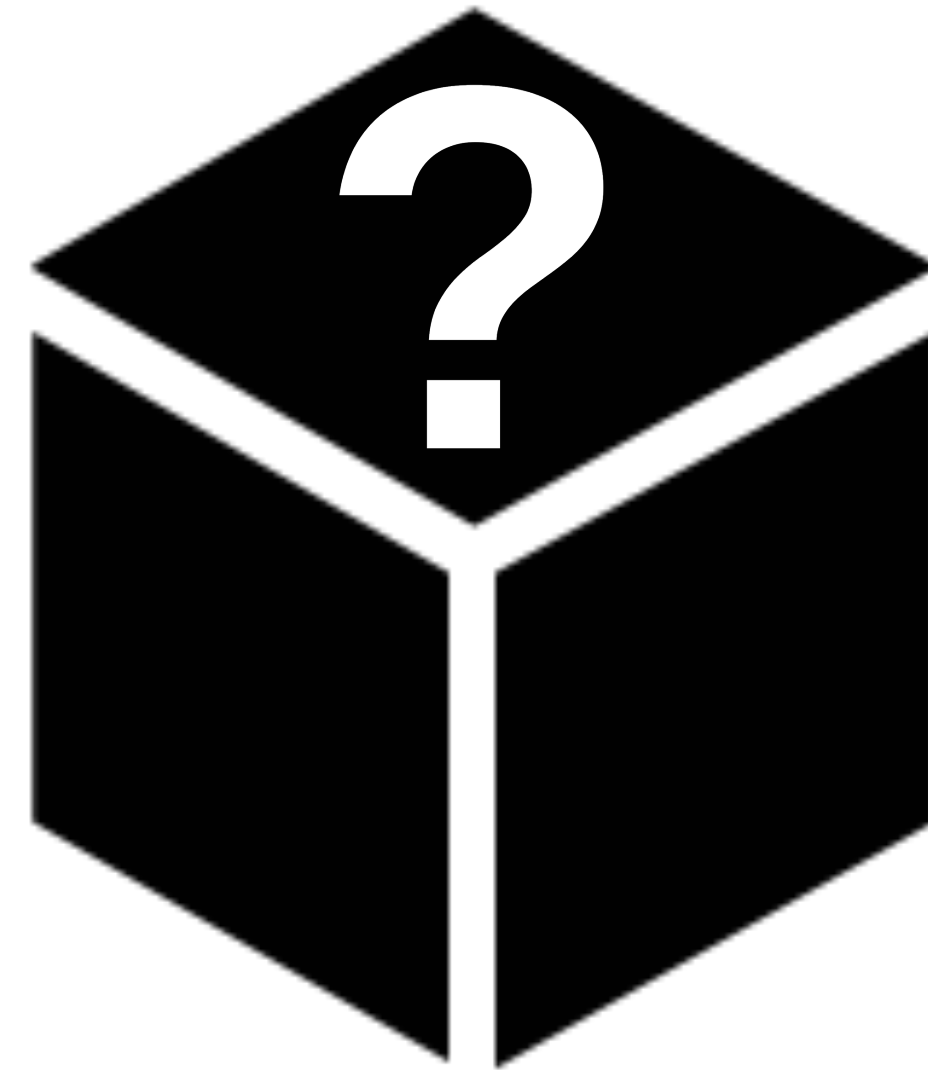
Data **L**ayout



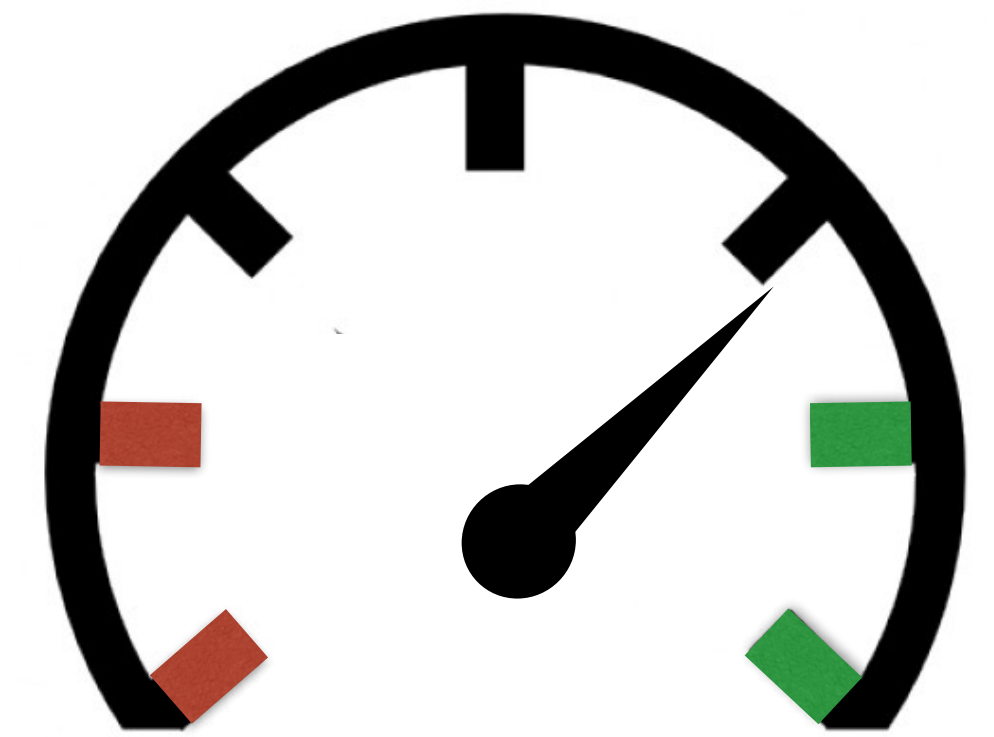
So, how do we reason about the data layout?



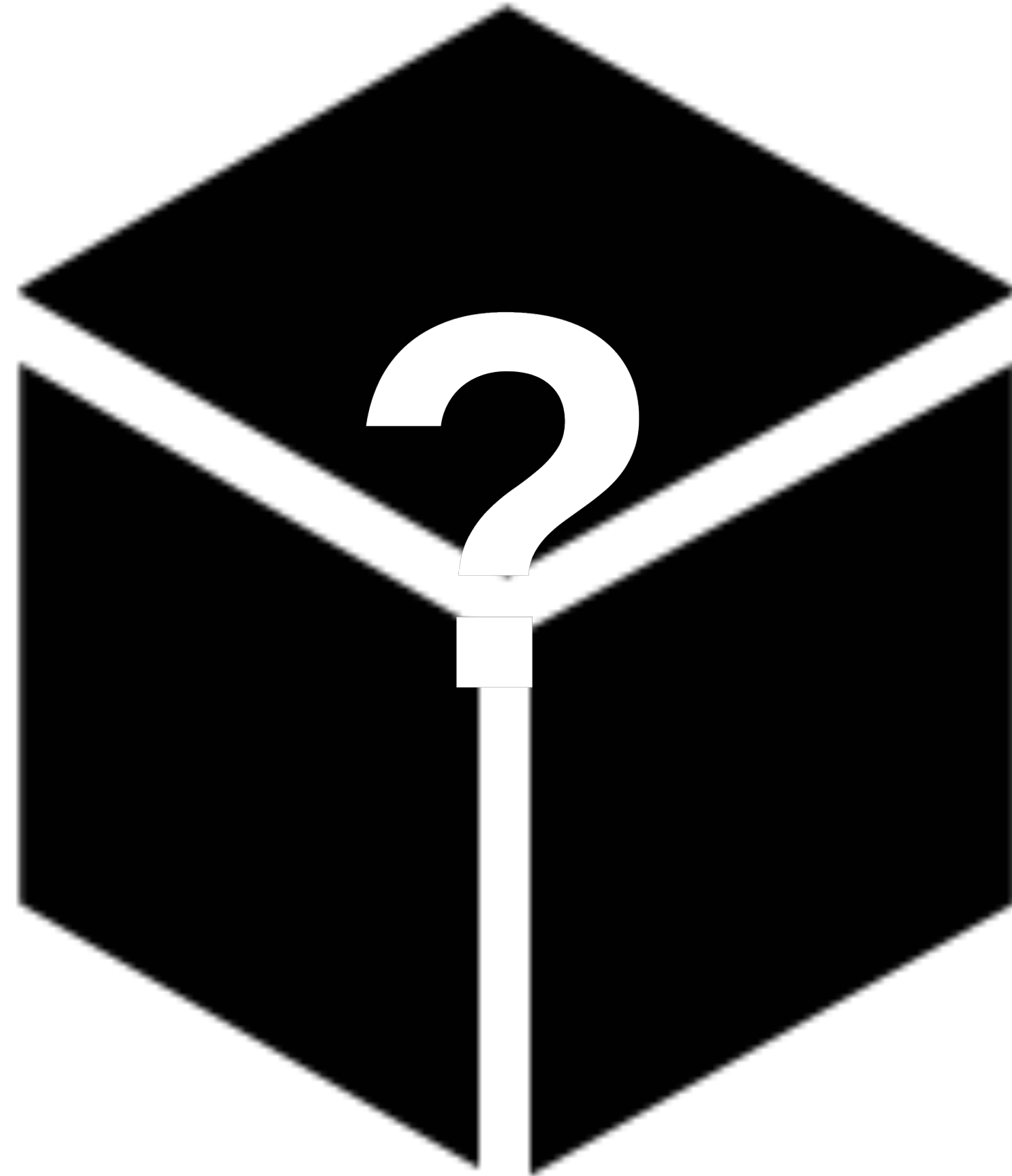
workload



data layout



performance



Compaction
black box

1

How to organize the data on device?

2

How much data to move at-a-time?

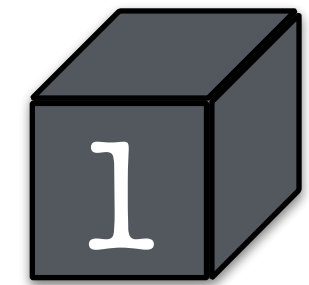
3

Which block of data to be moved?

4

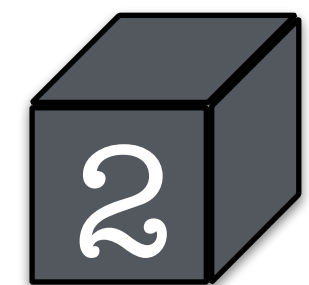
When to re-organize the data layout?

Data Layout



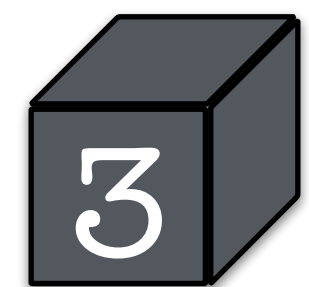
How to organize the data on device? ✓

Compaction
Granularity



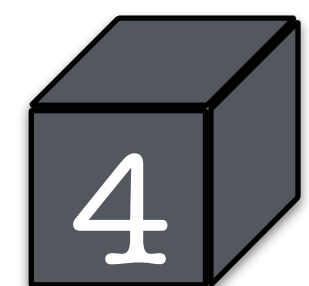
How much data to move at-a-time?

Data Movement
Policy

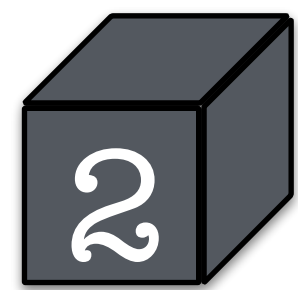


Which block of data to be moved?

Compaction
Trigger



When to re-organize the data layout?



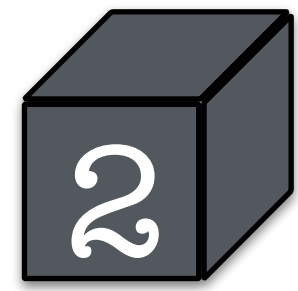
Compaction **Granularity**

data moved per compaction



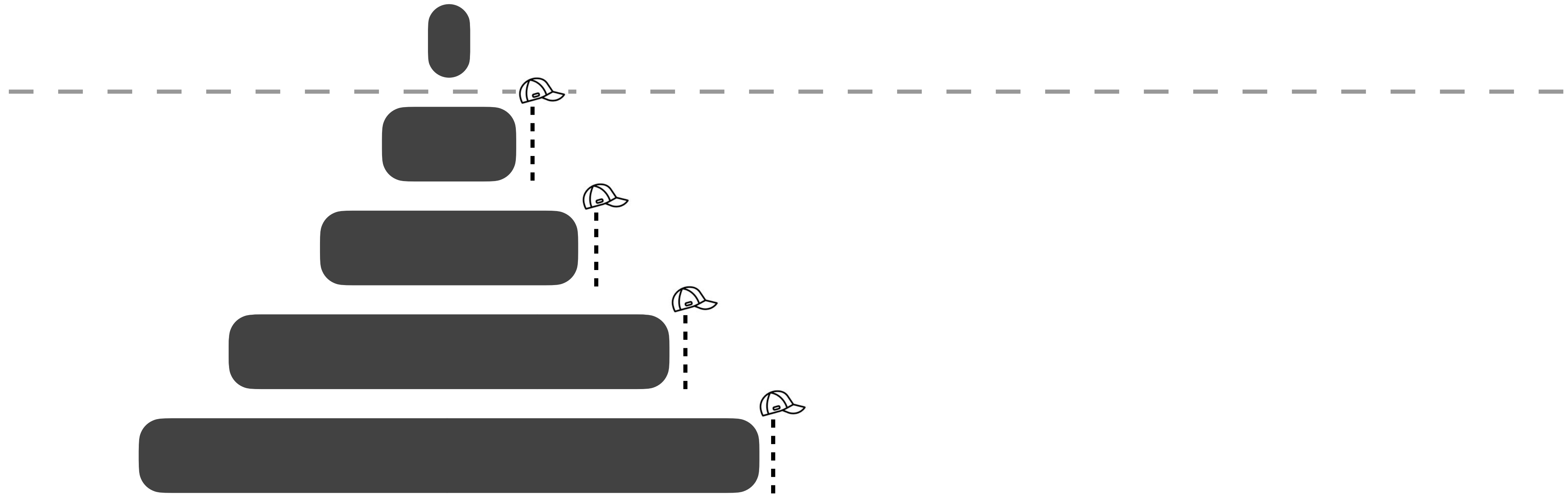
consecutive
levels

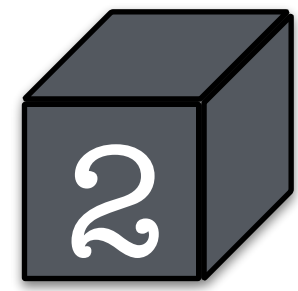
AsterixDB



Compaction **Granularity**

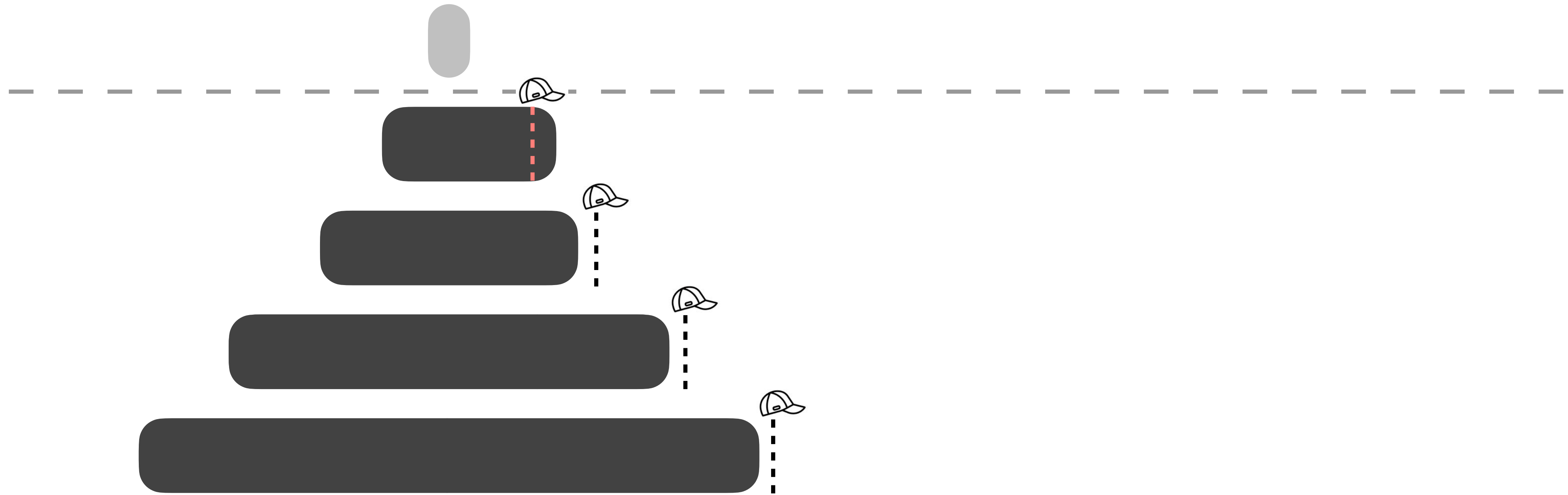
data moved per compaction

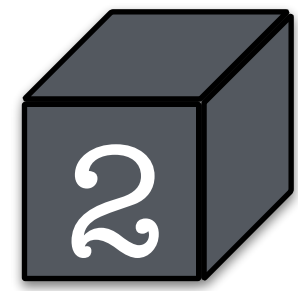




Compaction **Granularity**

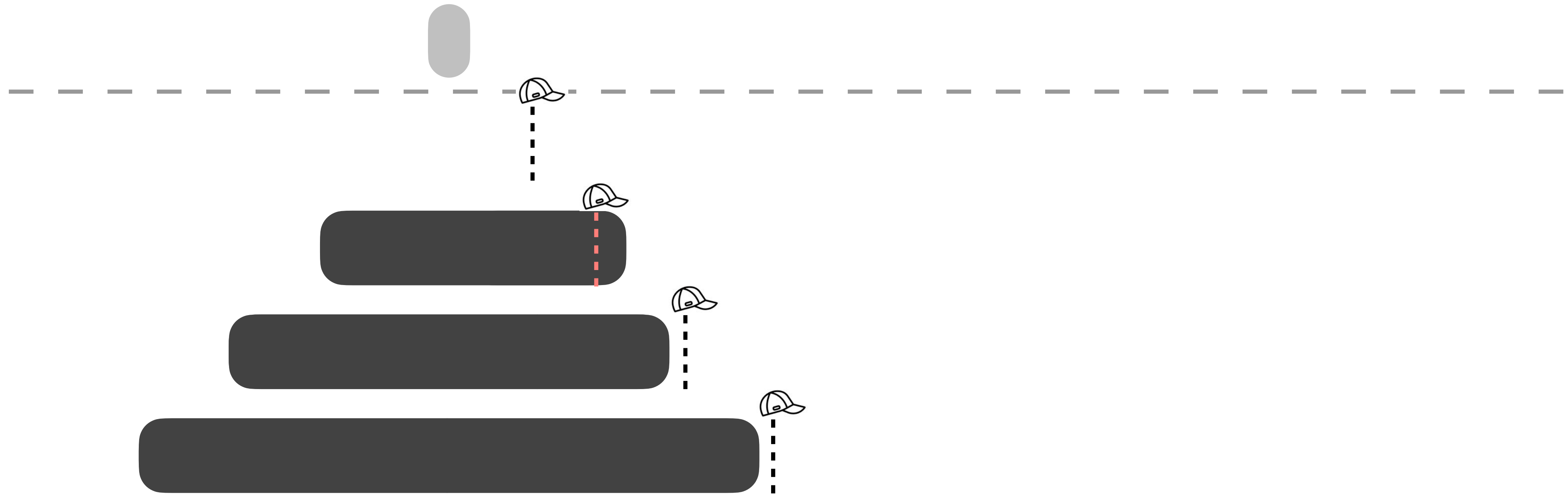
data moved per compaction

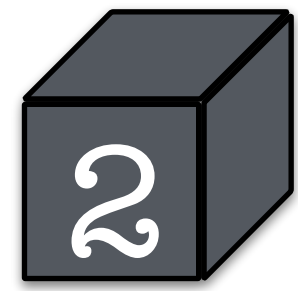




Compaction **Granularity**

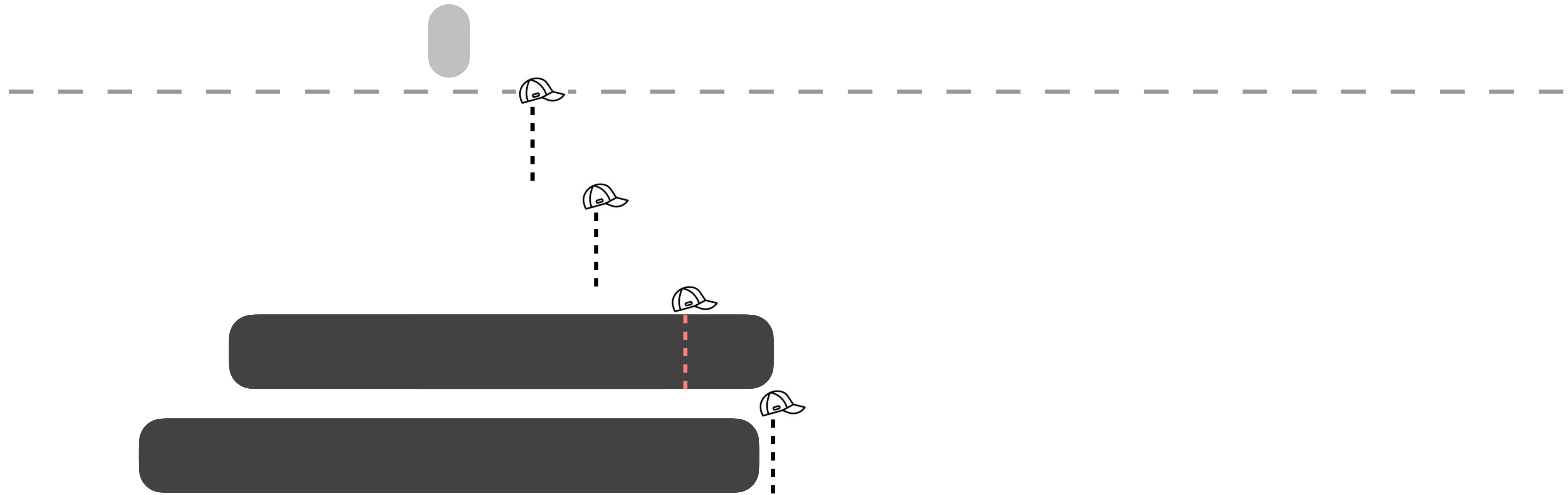
data moved per compaction

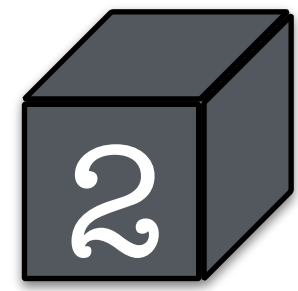




Compaction **Granularity**

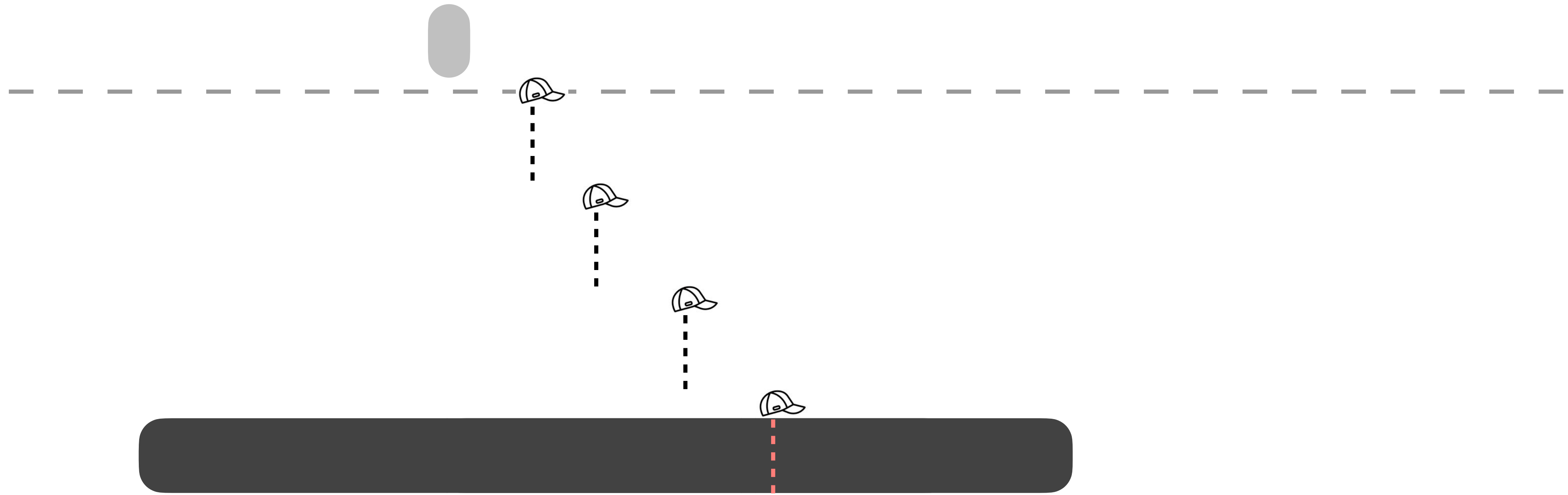
data moved per compaction

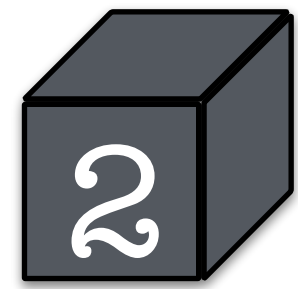




Compaction **Granularity**

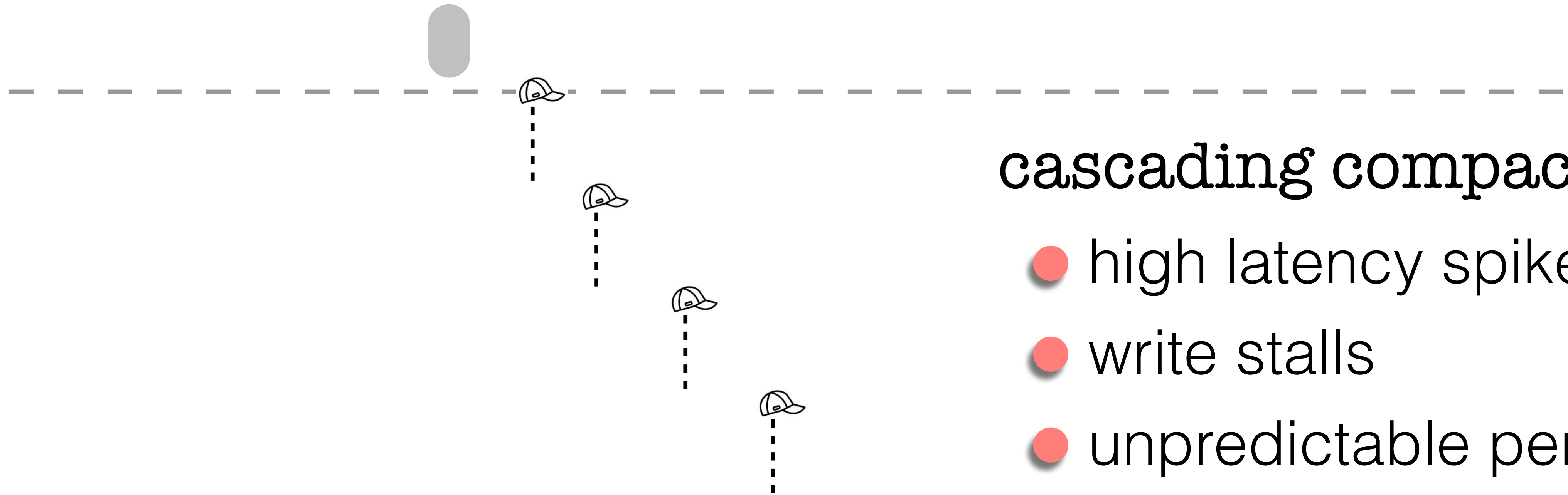
data moved per compaction





Compaction **Granularity**

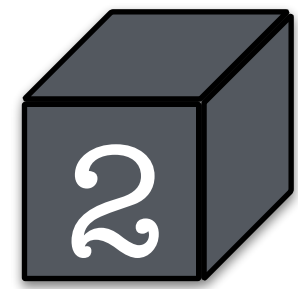
data moved per compaction



cascading compactions

- high latency spikes
- write stalls
- unpredictable perf.



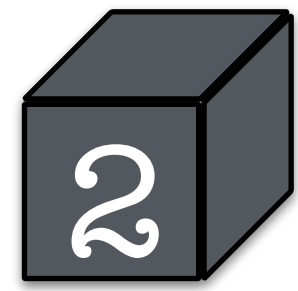


Compaction **Granularity**

data moved per compaction

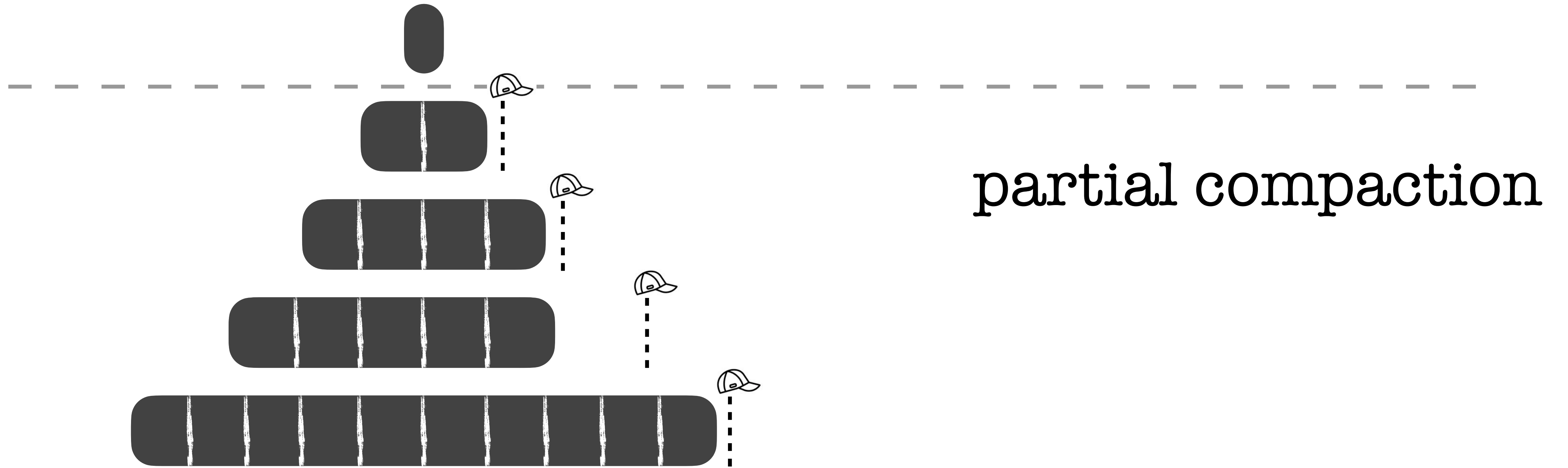
partial compaction

granularity: files

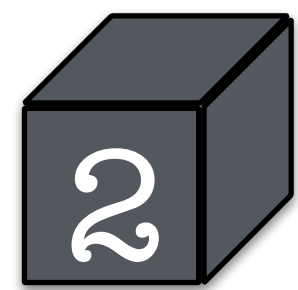


Compaction **Granularity**

data moved per compaction

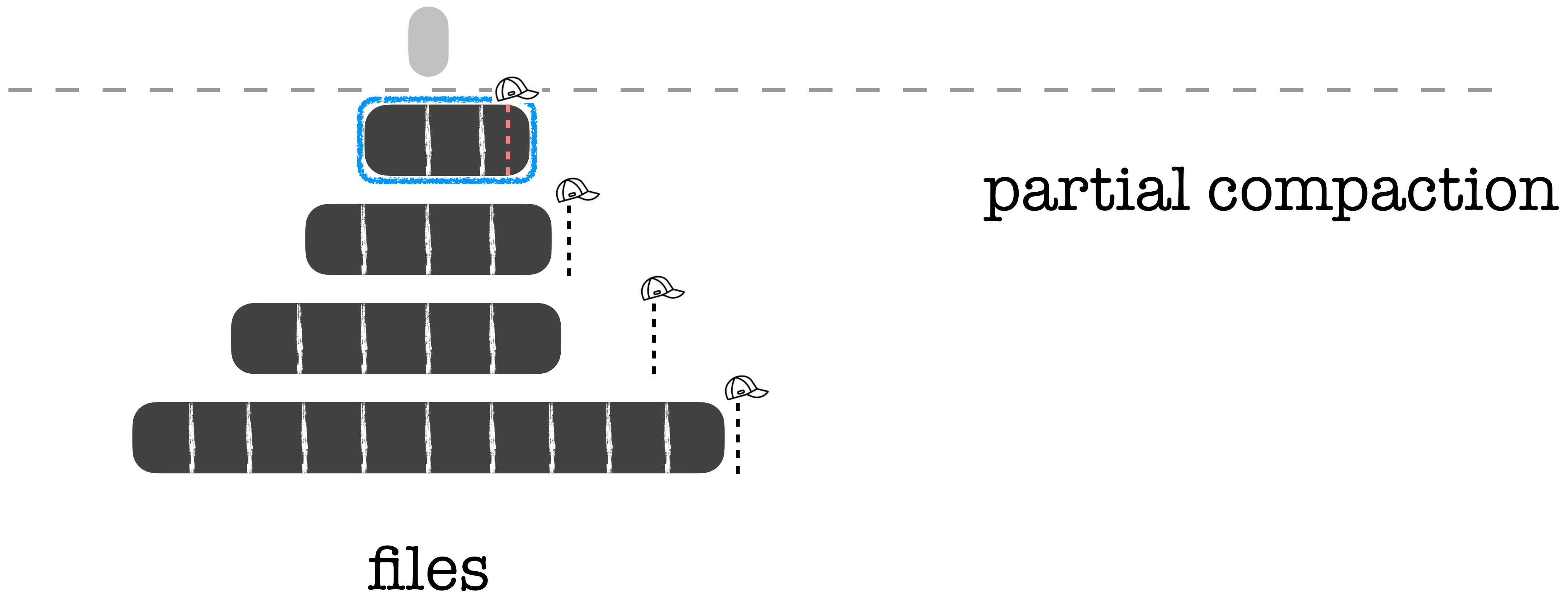


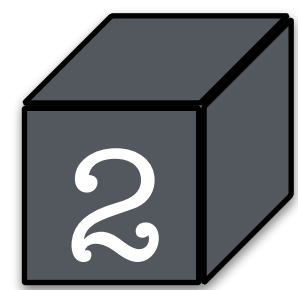
files



Compaction **Granularity**

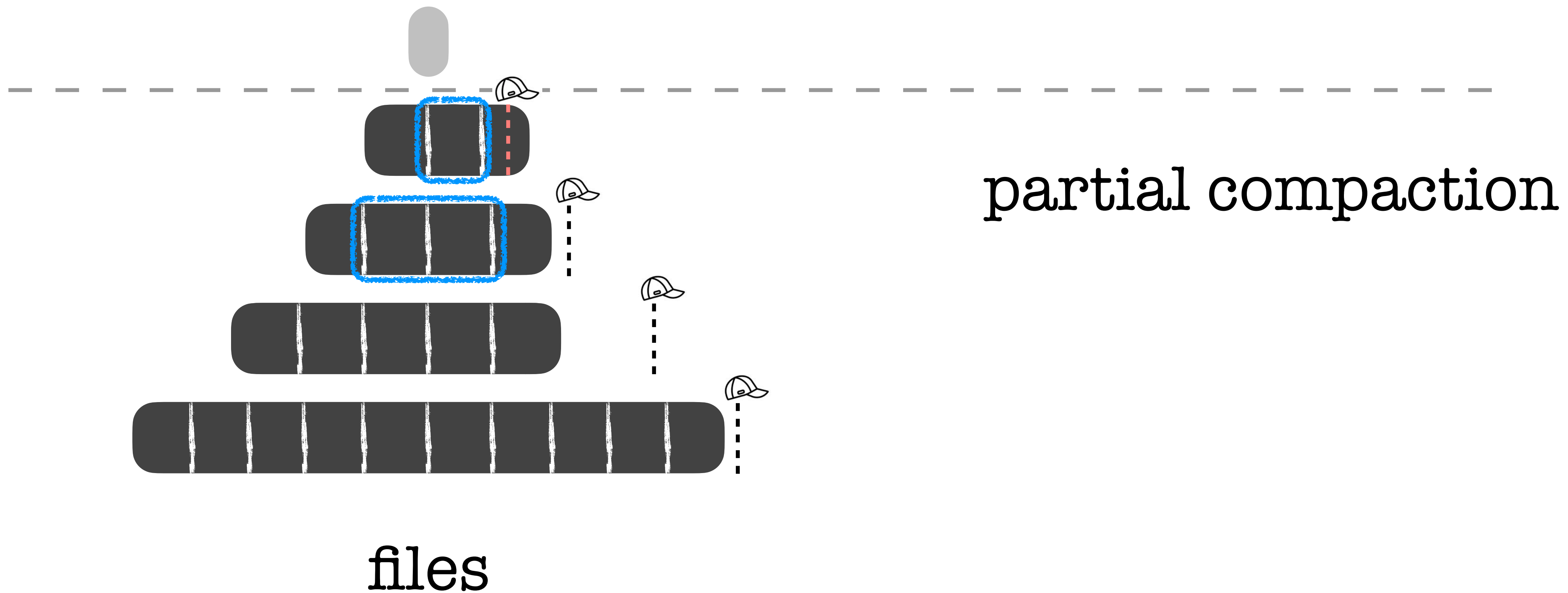
data moved per compaction

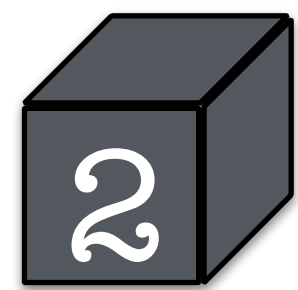




Compaction Granularity

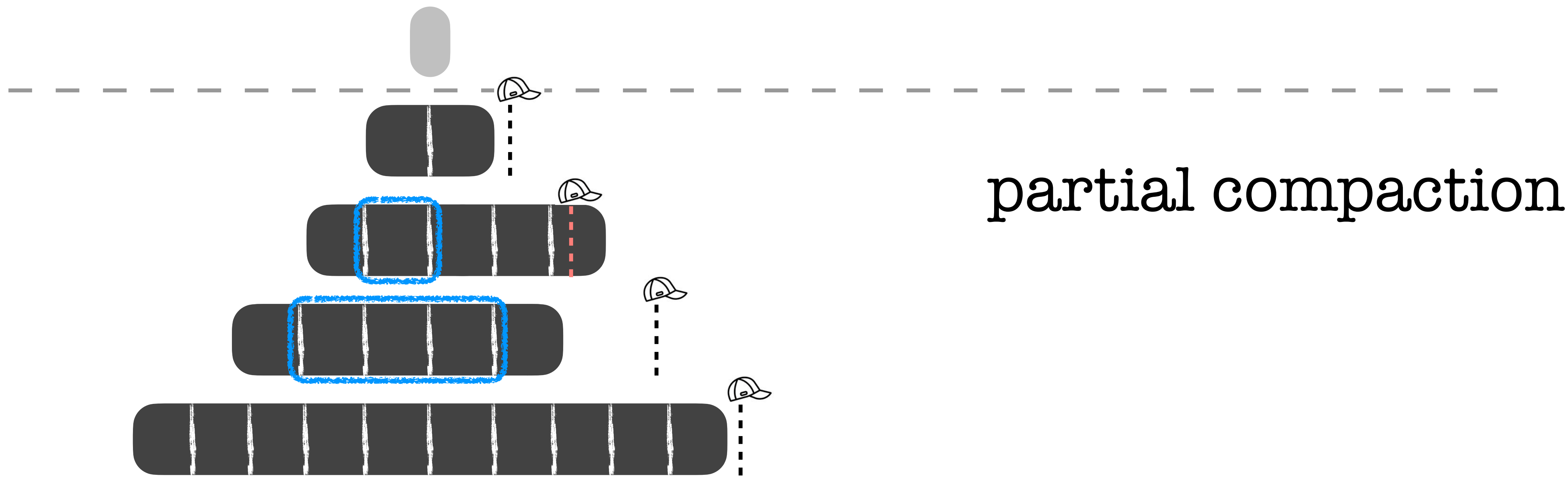
data moved per compaction





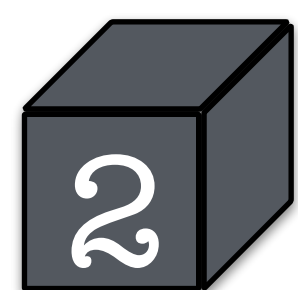
Compaction **Granularity**

data moved per compaction



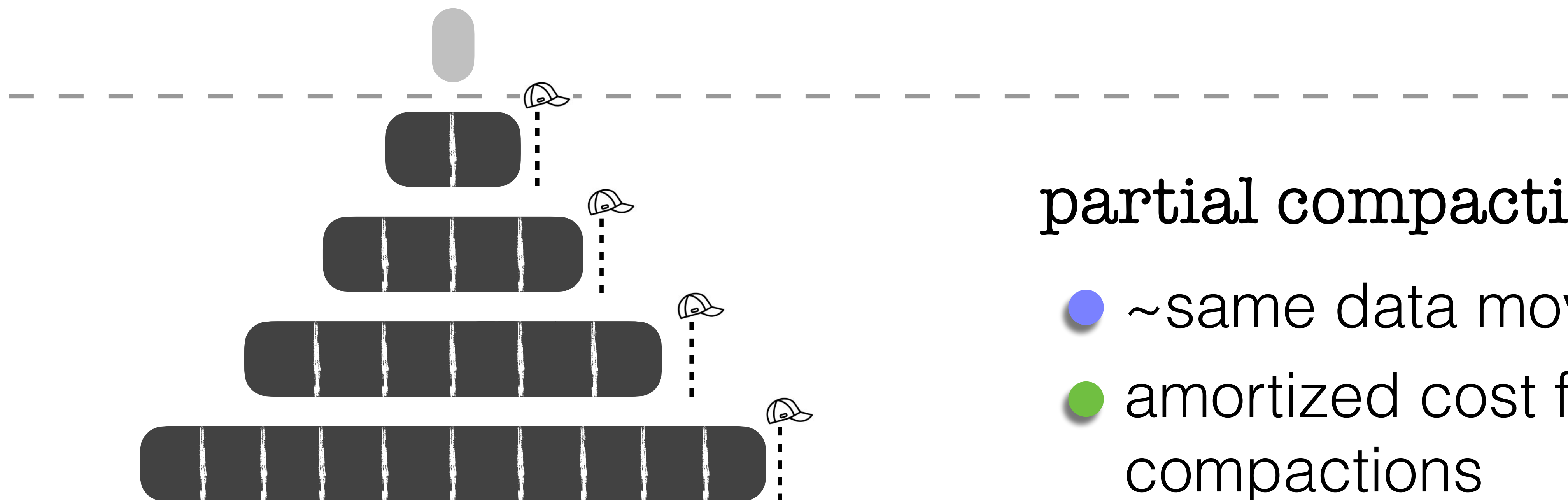
partial compaction

files



Compaction Granularity

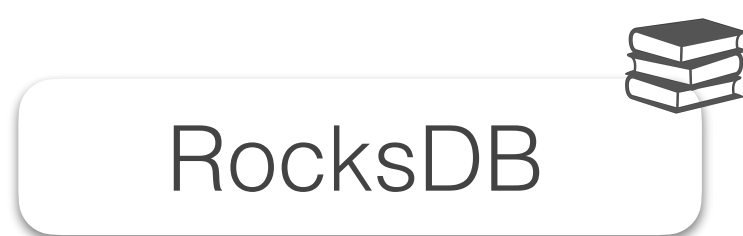
data moved per compaction

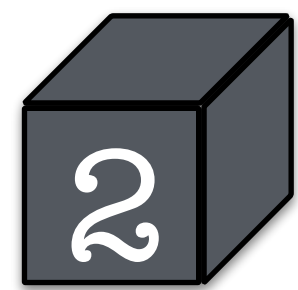


partial compaction

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

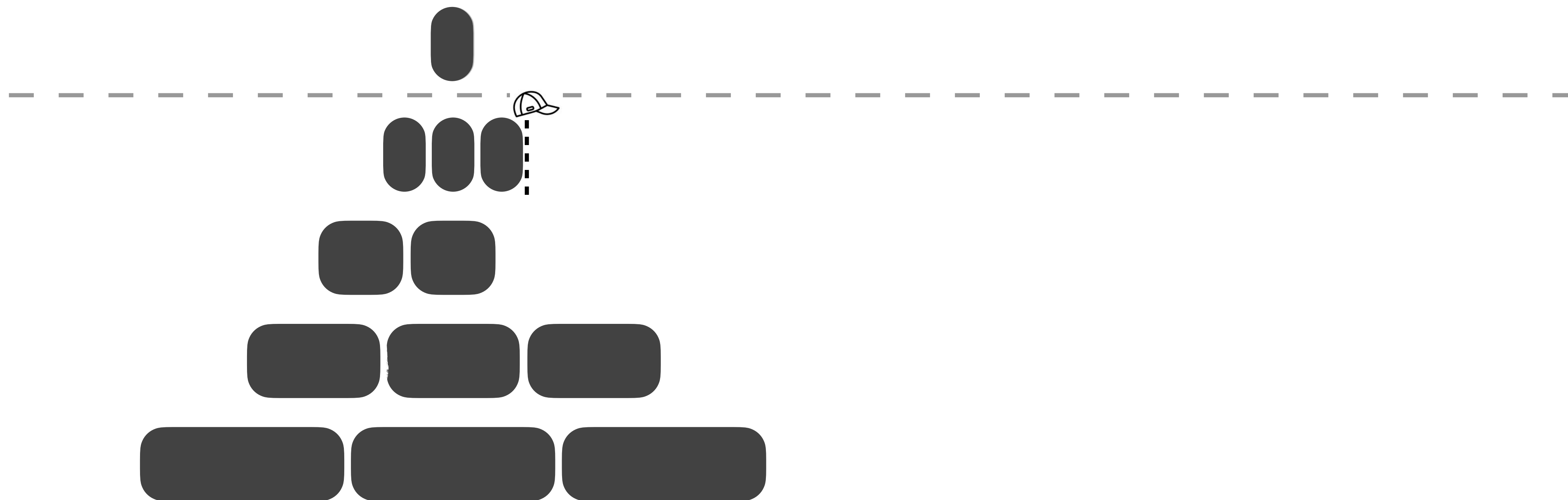
files



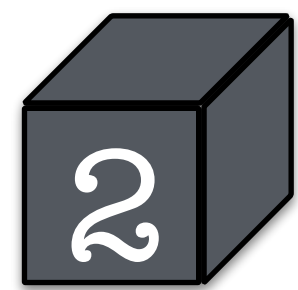


Compaction **Granularity**

data moved per compaction

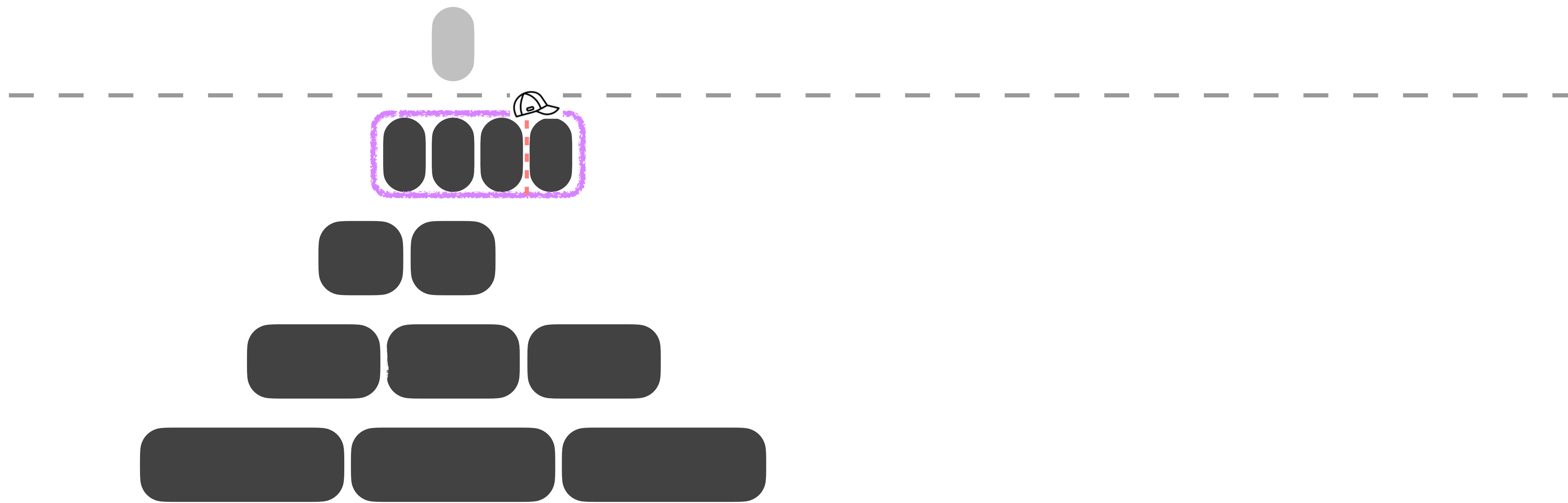


sorted runs in a level

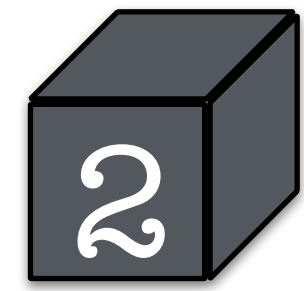


Compaction **Granularity**

data moved per compaction

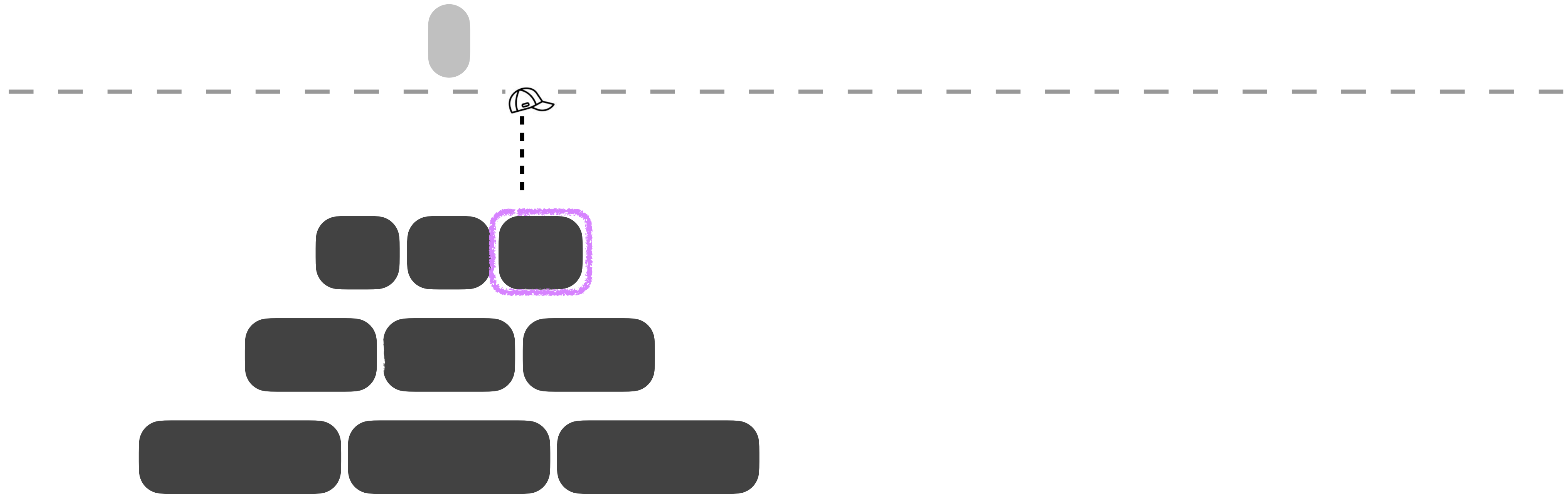


sorted runs in a level

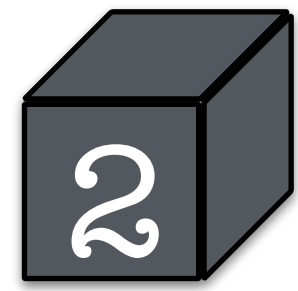


Compaction **Granularity**

data moved per compaction

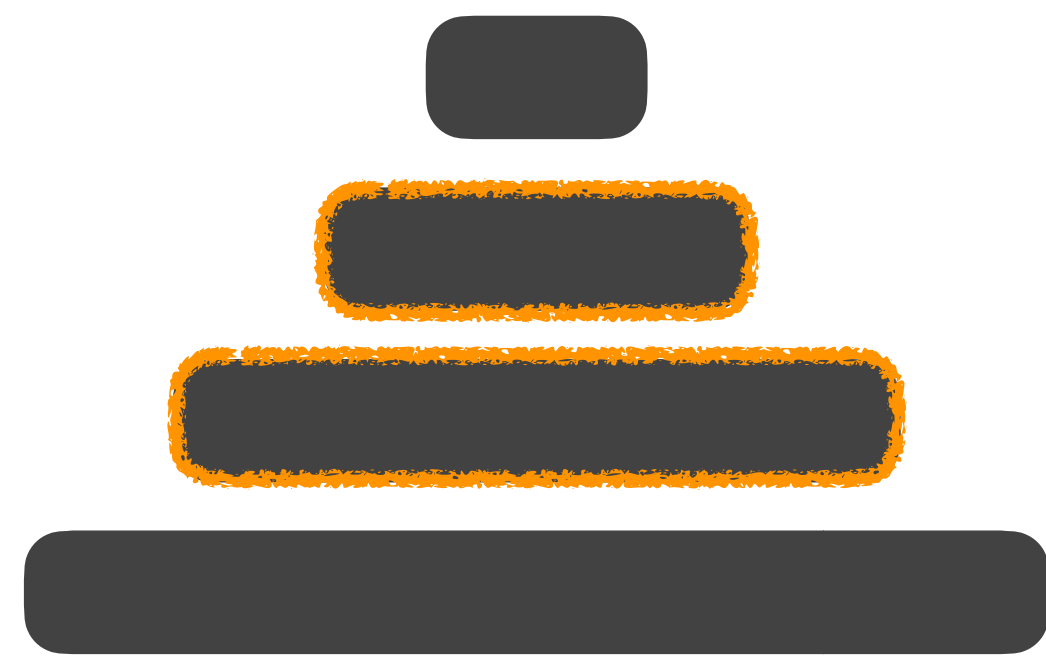


sorted runs in a level

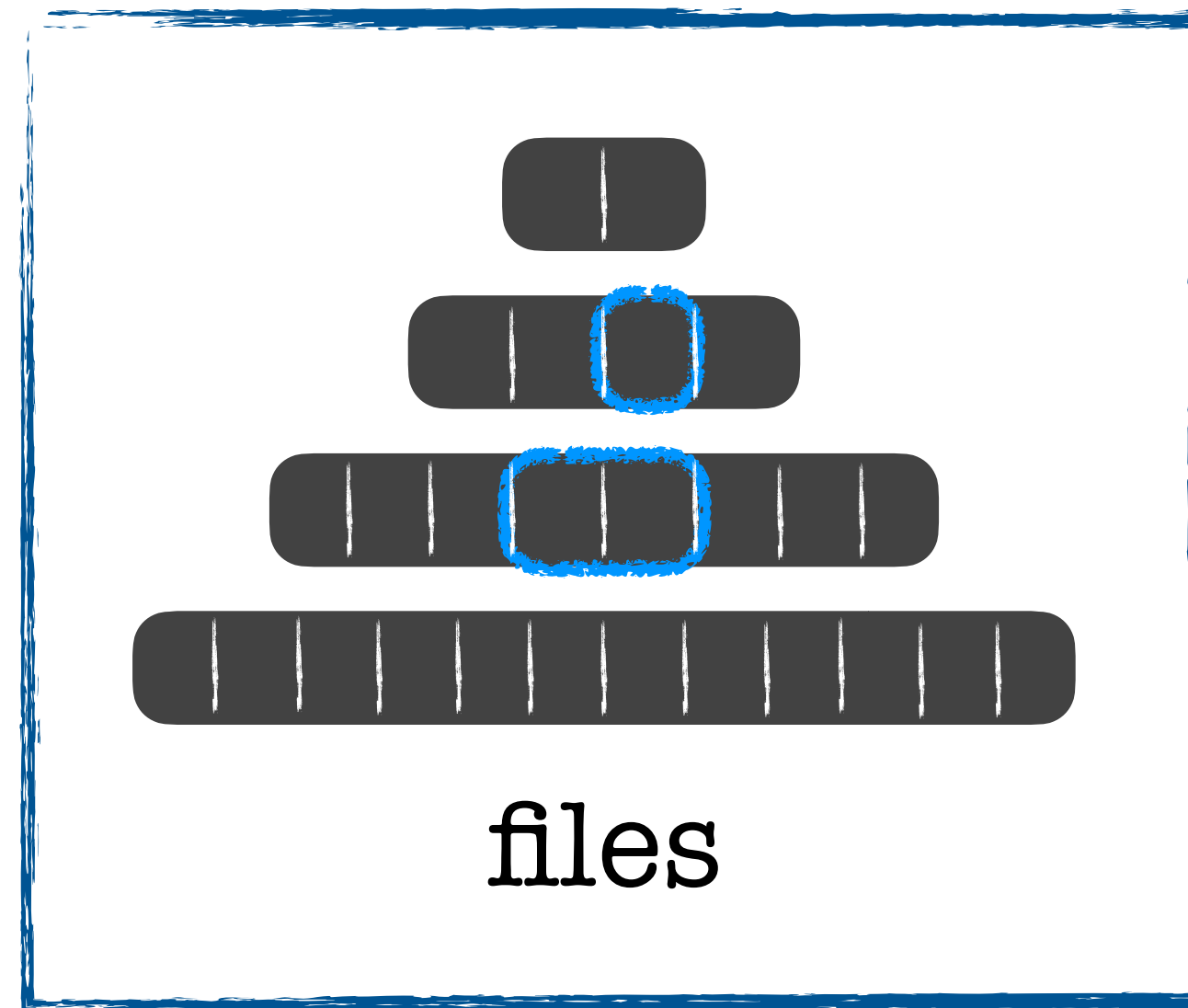


Compaction **Granularity**

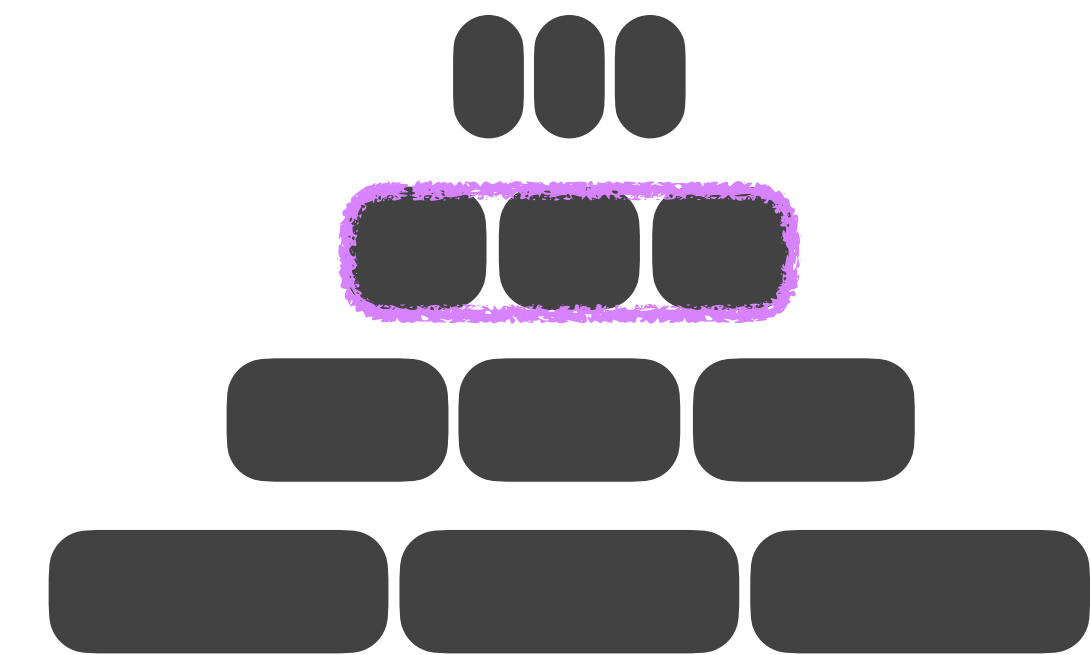
data moved per compaction



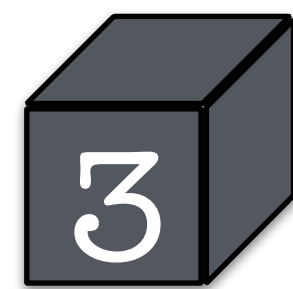
levels



files

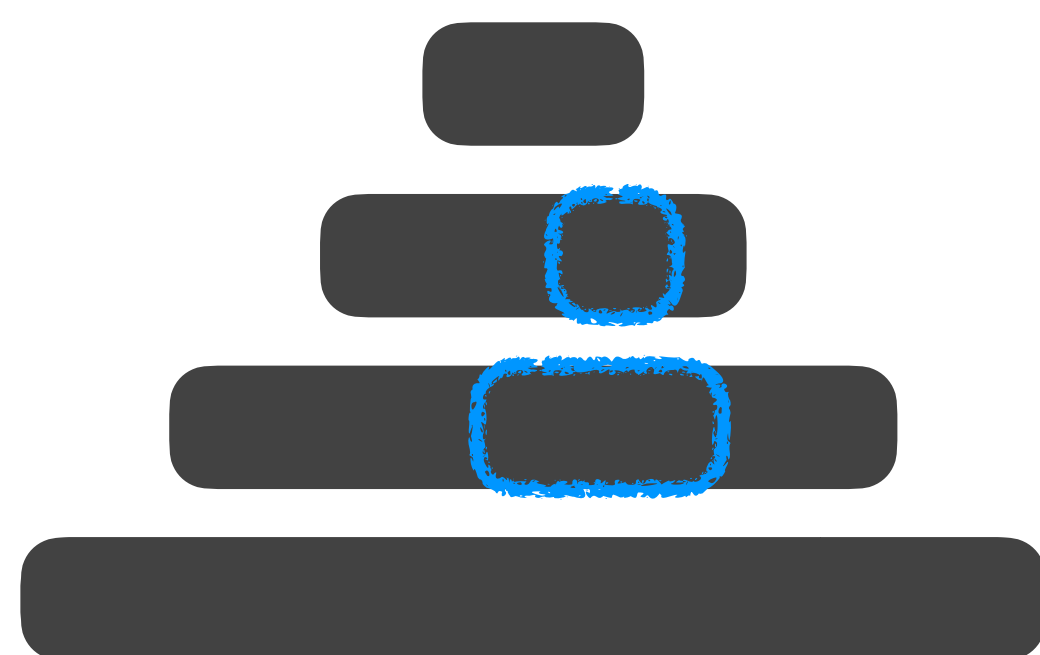


sorted runs in a level

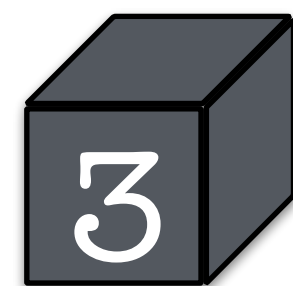


Data Movement Policy

which data to compact

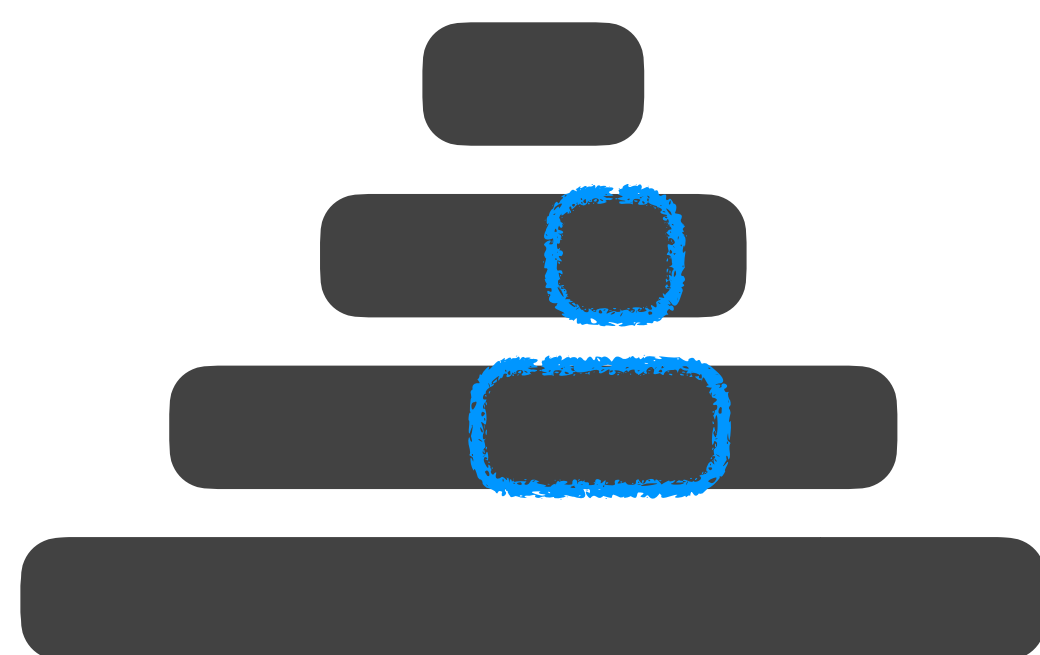


files



Data Movement Policy

which data to compact



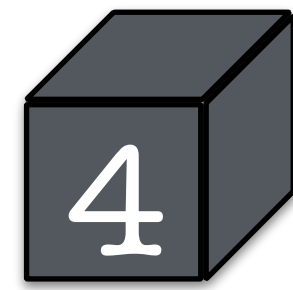
files

round-robin 

minimum **overlap with parent** level 

file with most **tombstones** 

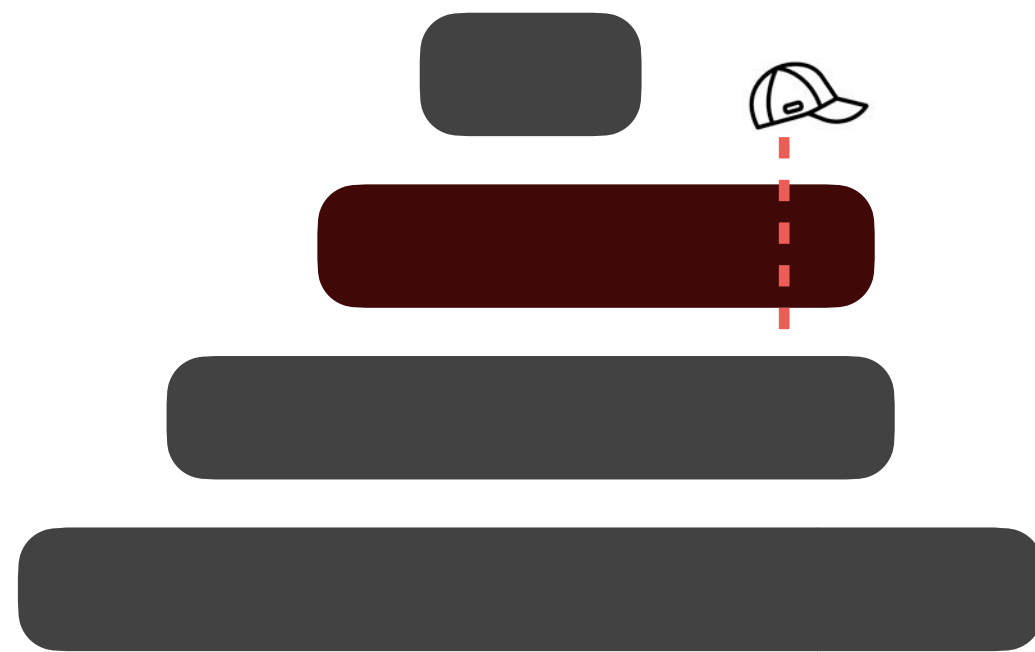
coldest file 

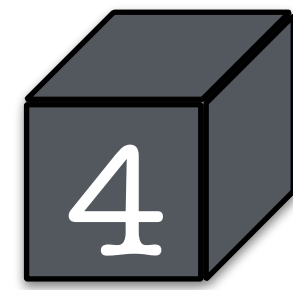


Compaction **Trigger**

invoking the compaction routine

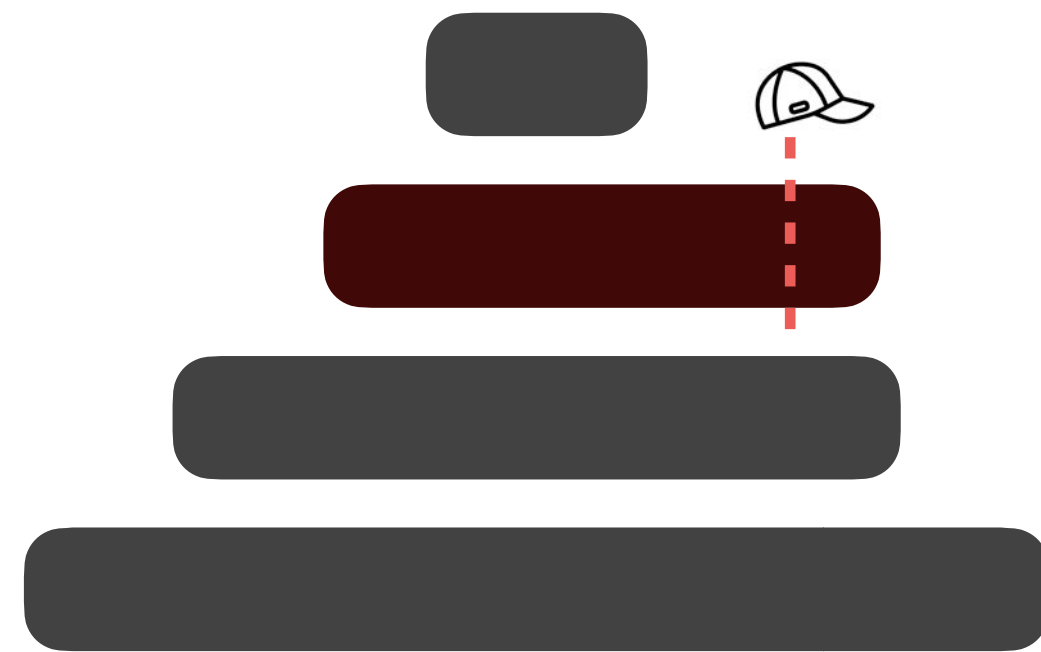
level **saturation**





4 Compaction **Trigger**

invoking the compaction routine



level **saturation**

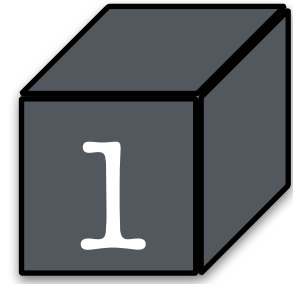
number of **sorted runs**

space amplification

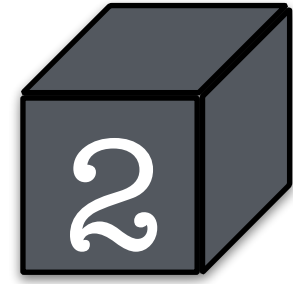


age of a file

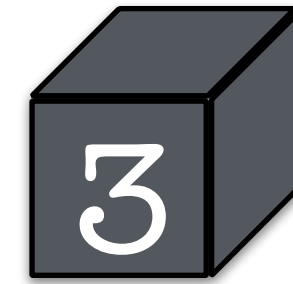




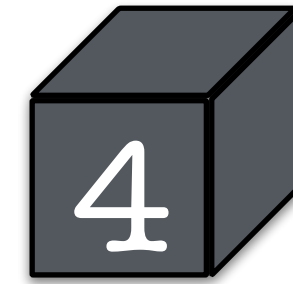
Data Layout



Compaction
Granularity



Data Movement
Policy



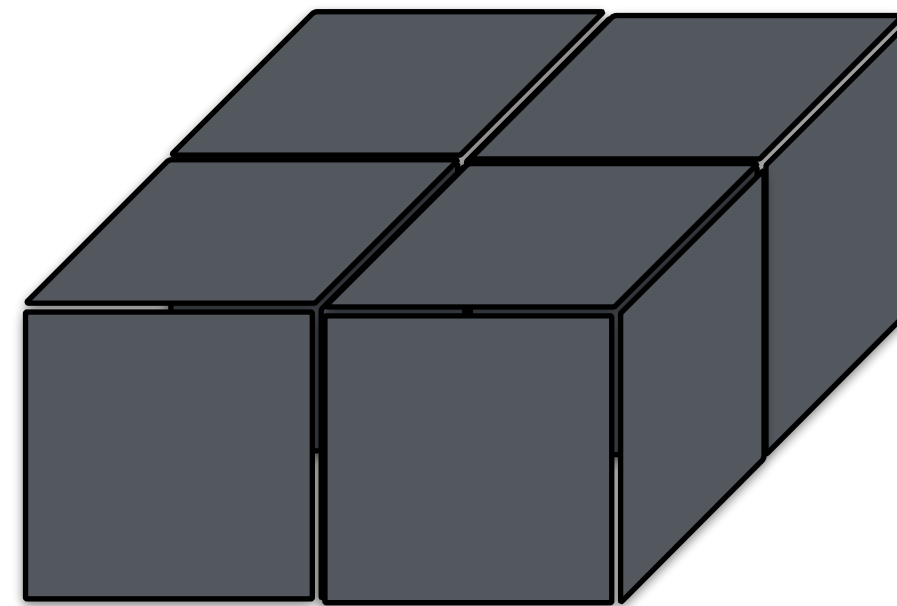
Compaction
Trigger

Data Layout

Compaction
Granularity

Data Movement
Policy

Compaction
Trigger



***Any* Compaction Algorithm**

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]	<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		✓					✓									✓

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]	<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		✓					✓									✓

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]	<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		✓					✓									✓

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]	<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

Comparative Analysis Individual Analysis

Workload

#Entries: 10000000

Entry size: 128 B

Main Memory Parameter

Buffer size: 16 MB

BF size / Entry: 10 MB

Disk Parameters

File size (in terms of buffer): 1

Page size: 4 KB

Data Layout

Size ratio: 4

Leveling Tiering

Step: 1

Advanced settings

Progress: 92% 9276935

Play Pause Finish

Vanilla-LSM

Leveling Tiering

▶

Partial Compaction

Leveled Compaction

▶

Hybrid Strategy

Lazy Leveling

▶

Build-Your-Own

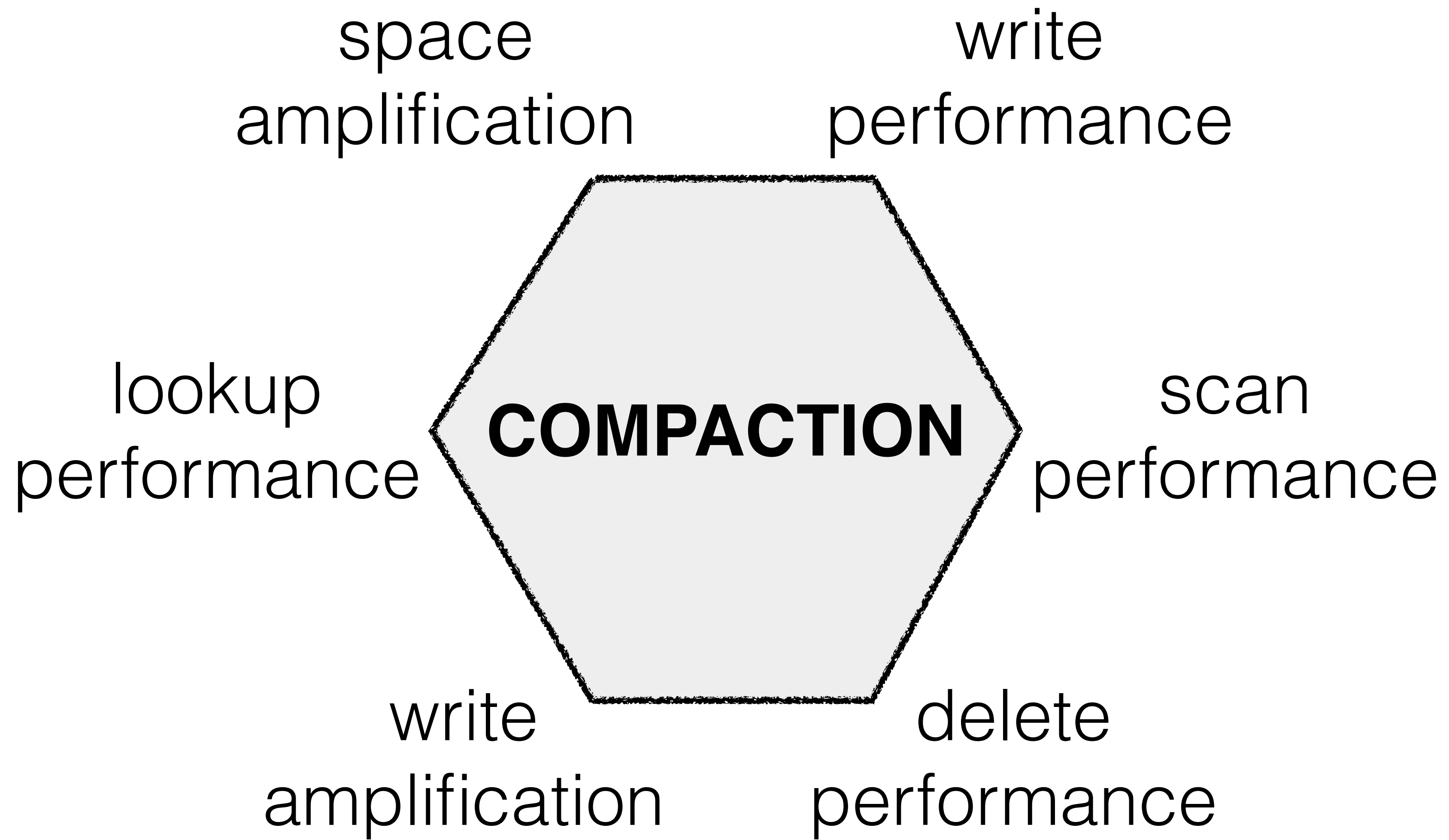
Tiered #Levels: 1

▶

Levels	4	3	4	4
#sorted runs	3	3	5	5
#compactions	67	186	31	31

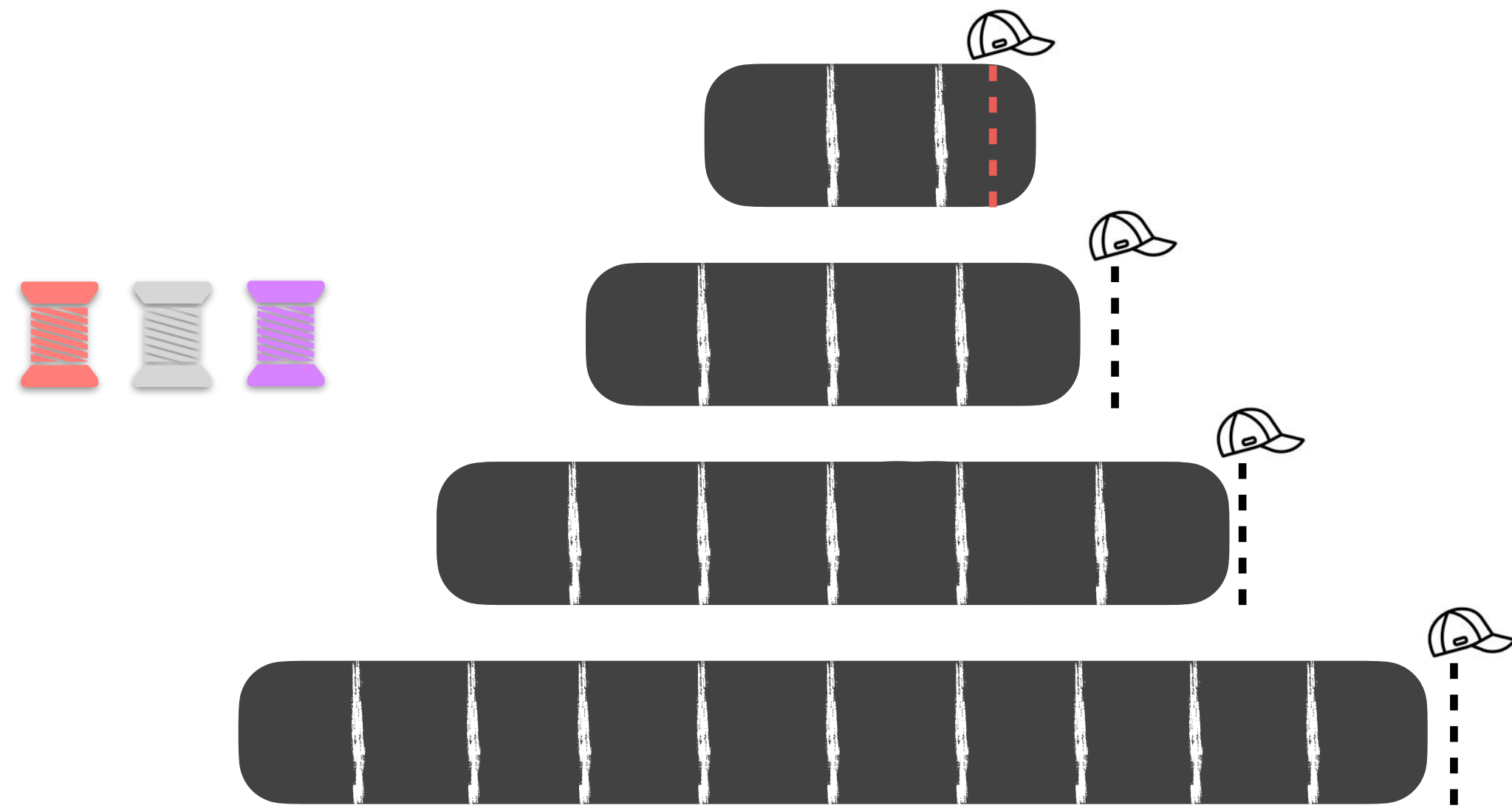
SarkarSIGMOD22

Tuesday
4-6PM



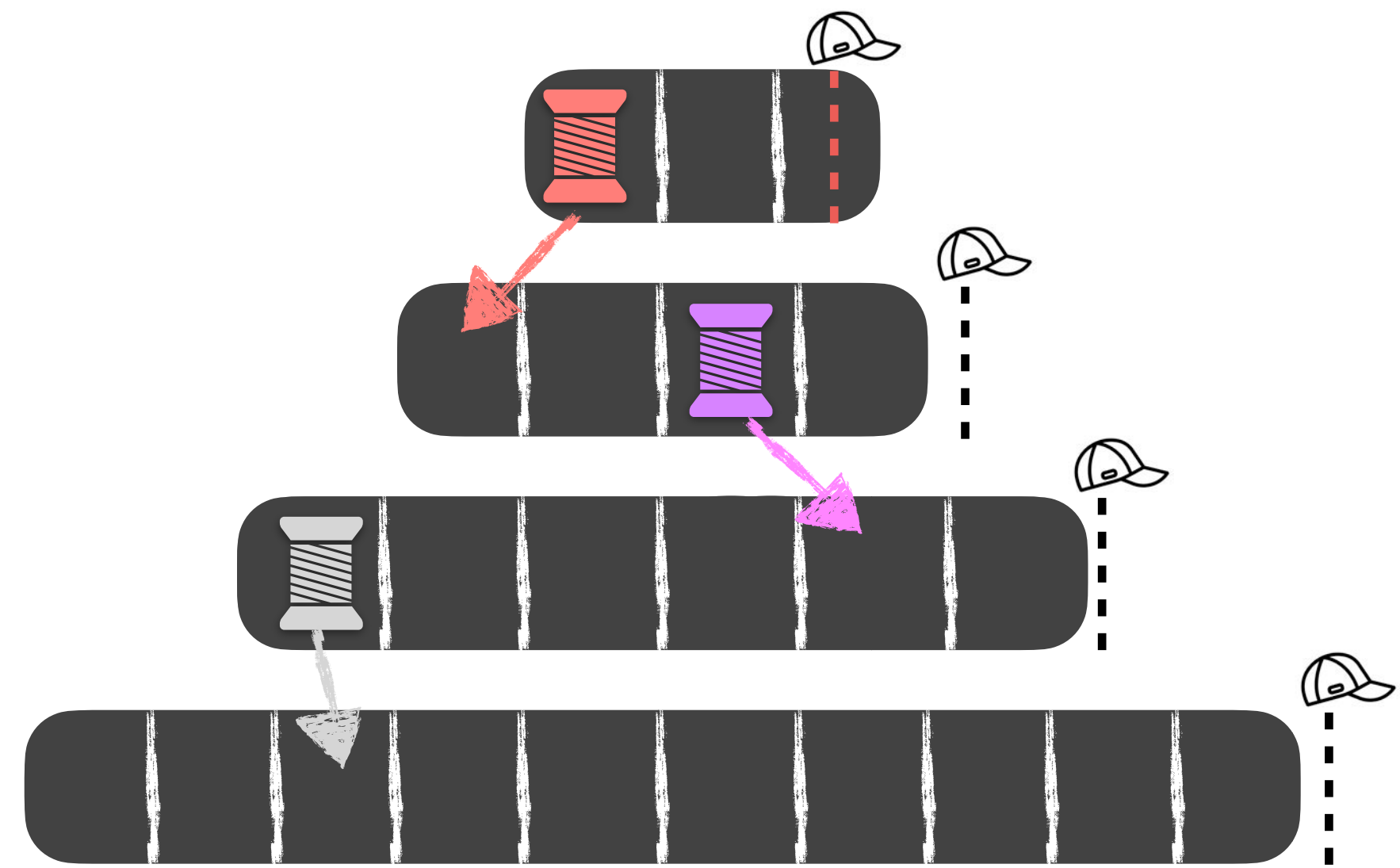
Optimizing Compactions

Background
Compactions



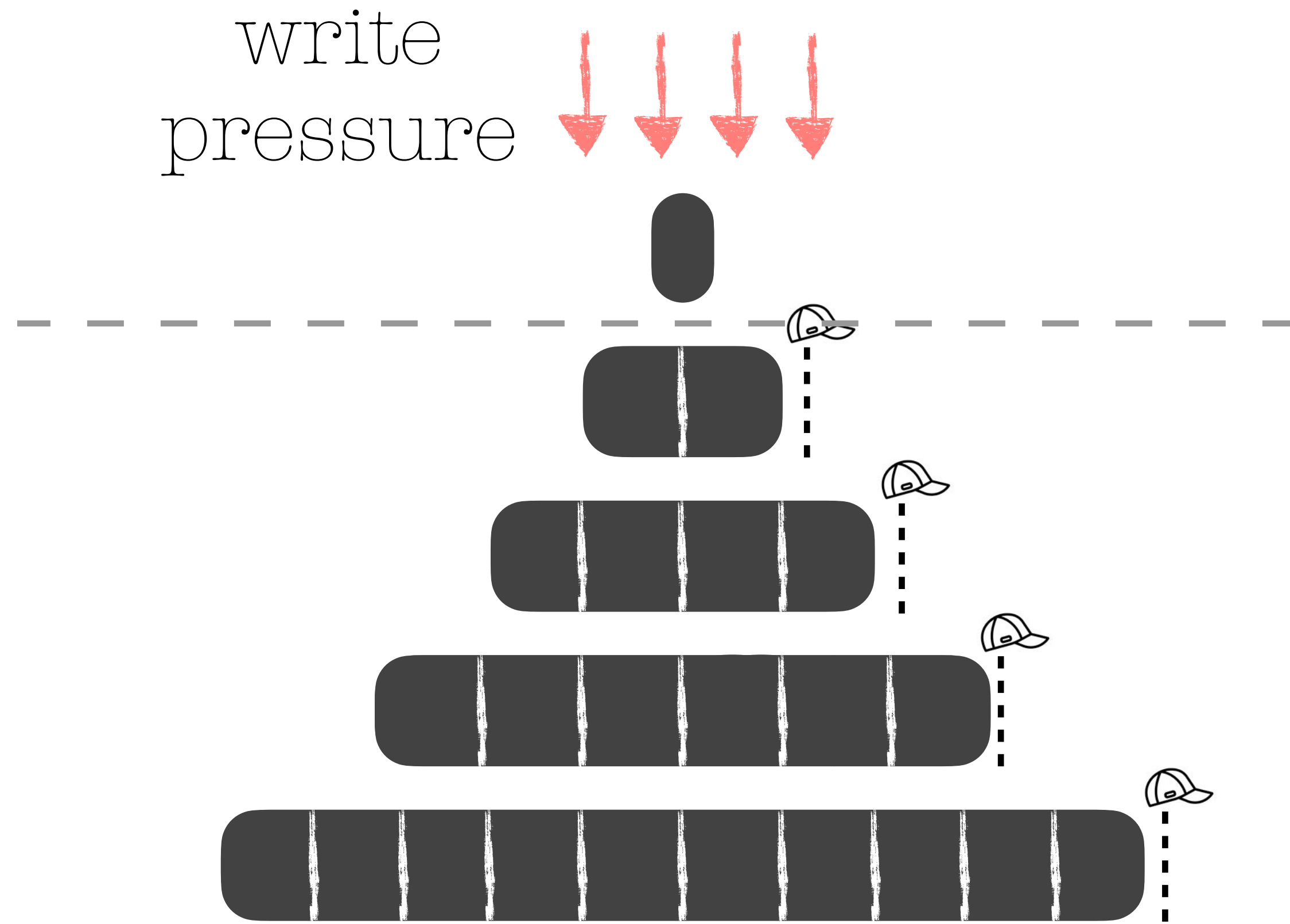
Optimizing **Compactions**

Background
Compactions



- non-blocking reads/writes
- improves write throughput

Optimizing Compactions

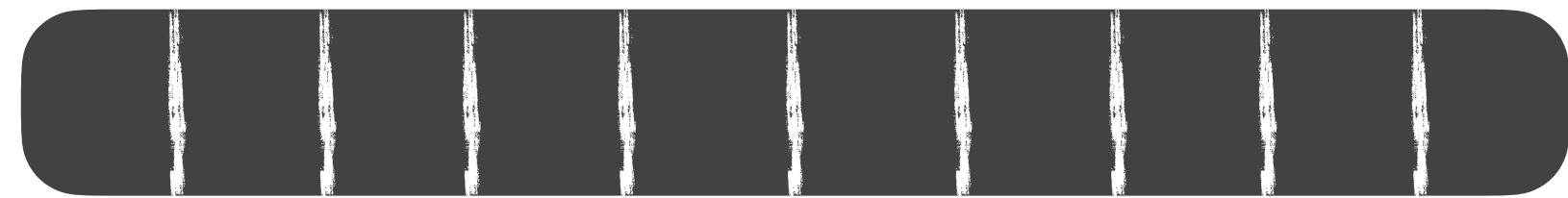
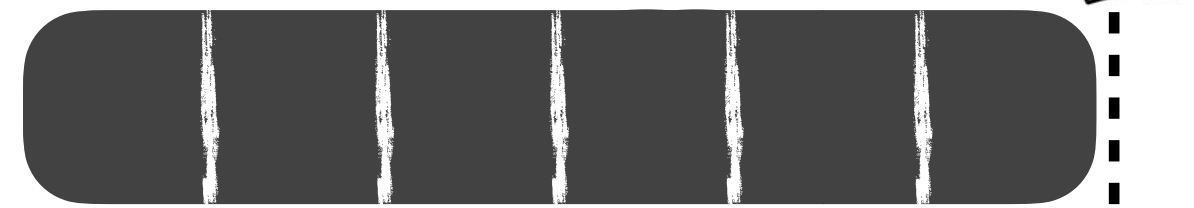
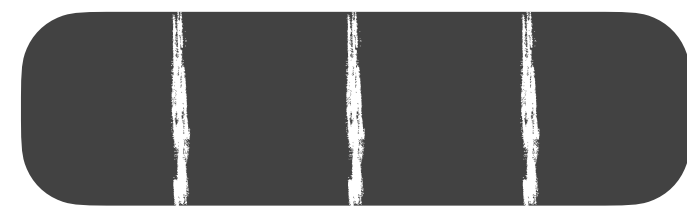
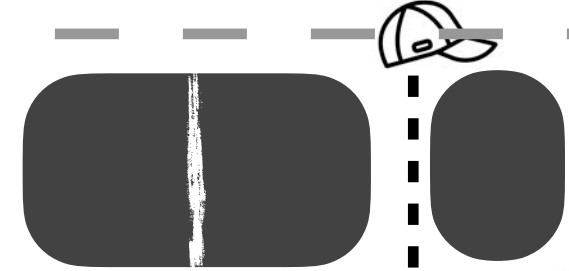
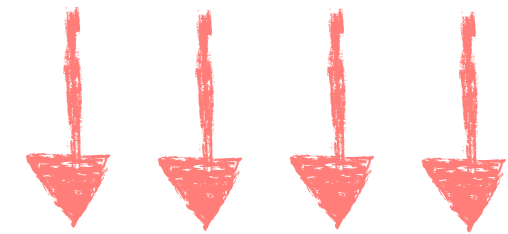


Background
Compactions

Compaction
Priority

Optimizing Compactions

write
pressure

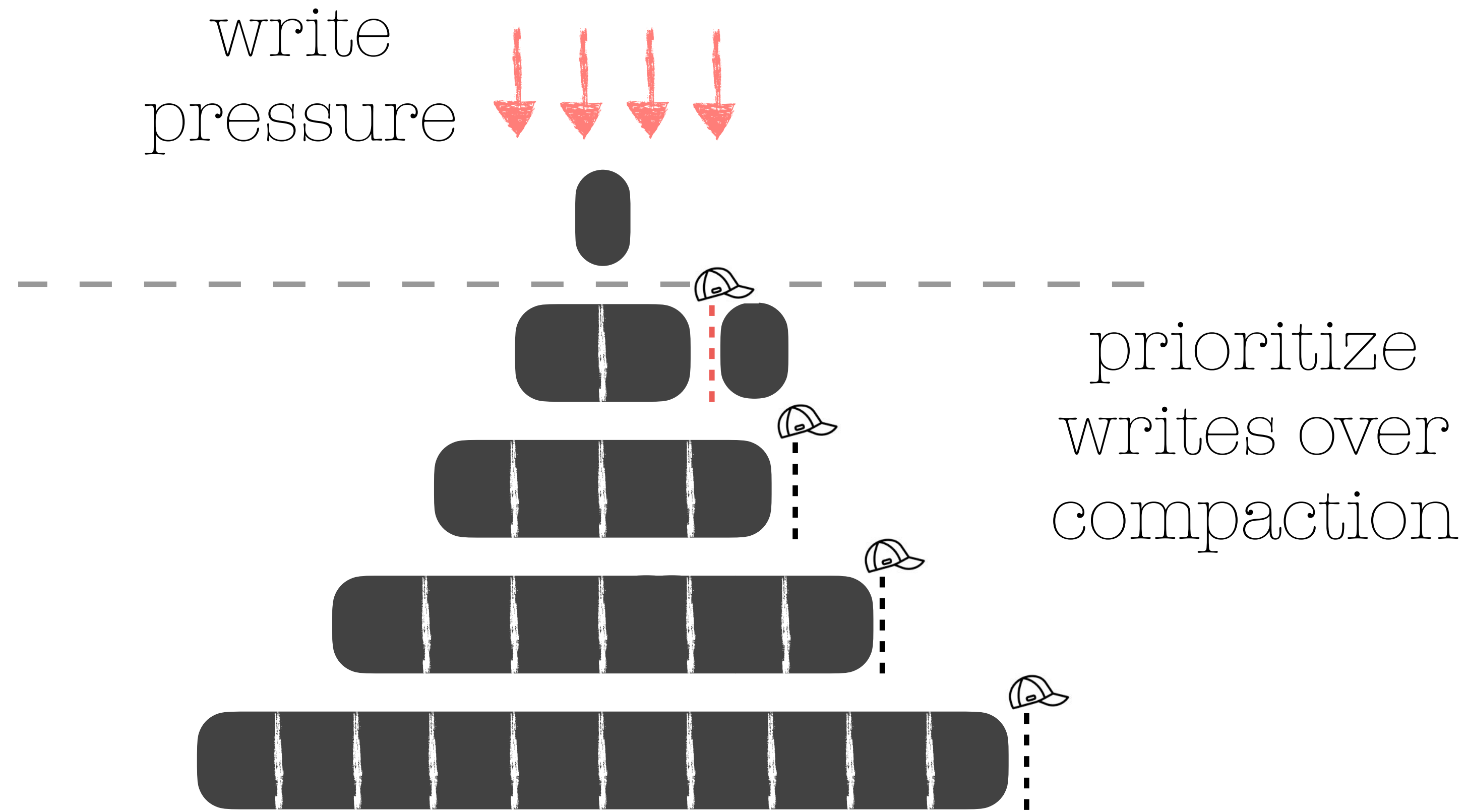


prioritize
writes over
compaction

Background
Compactions

Compaction
Priority

Optimizing Compactions

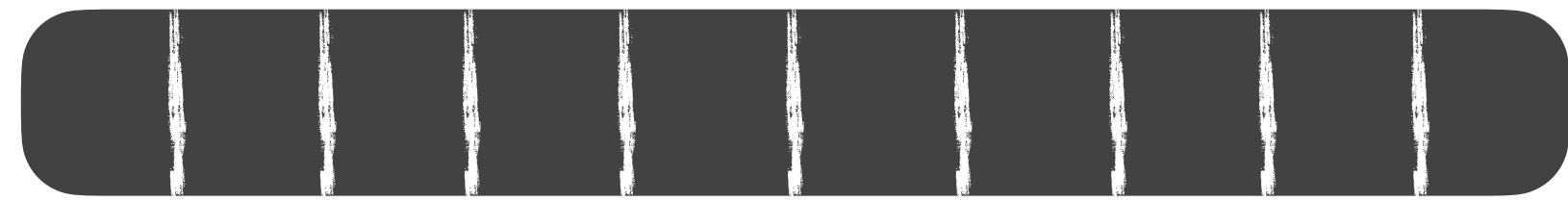
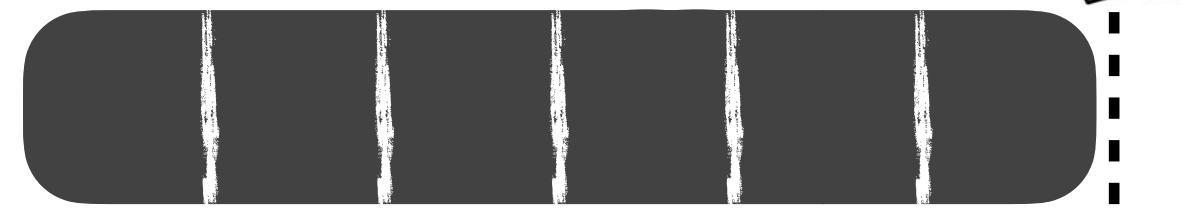
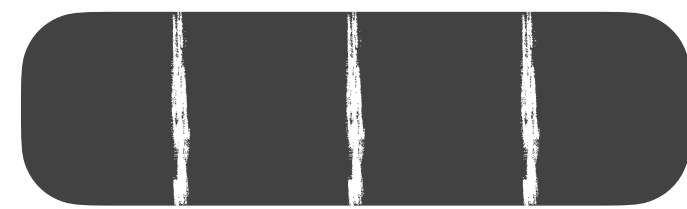
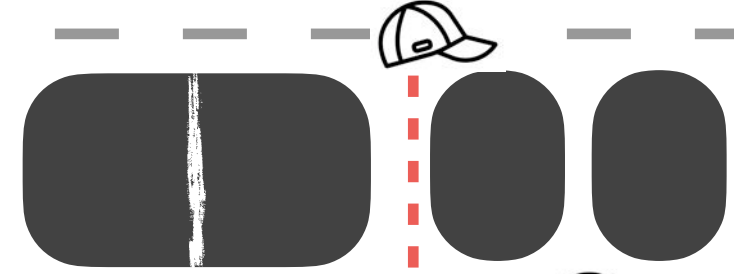
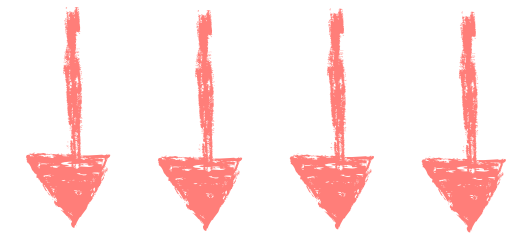


Background Compactions

Compaction Priority

Optimizing Compactions

write
pressure

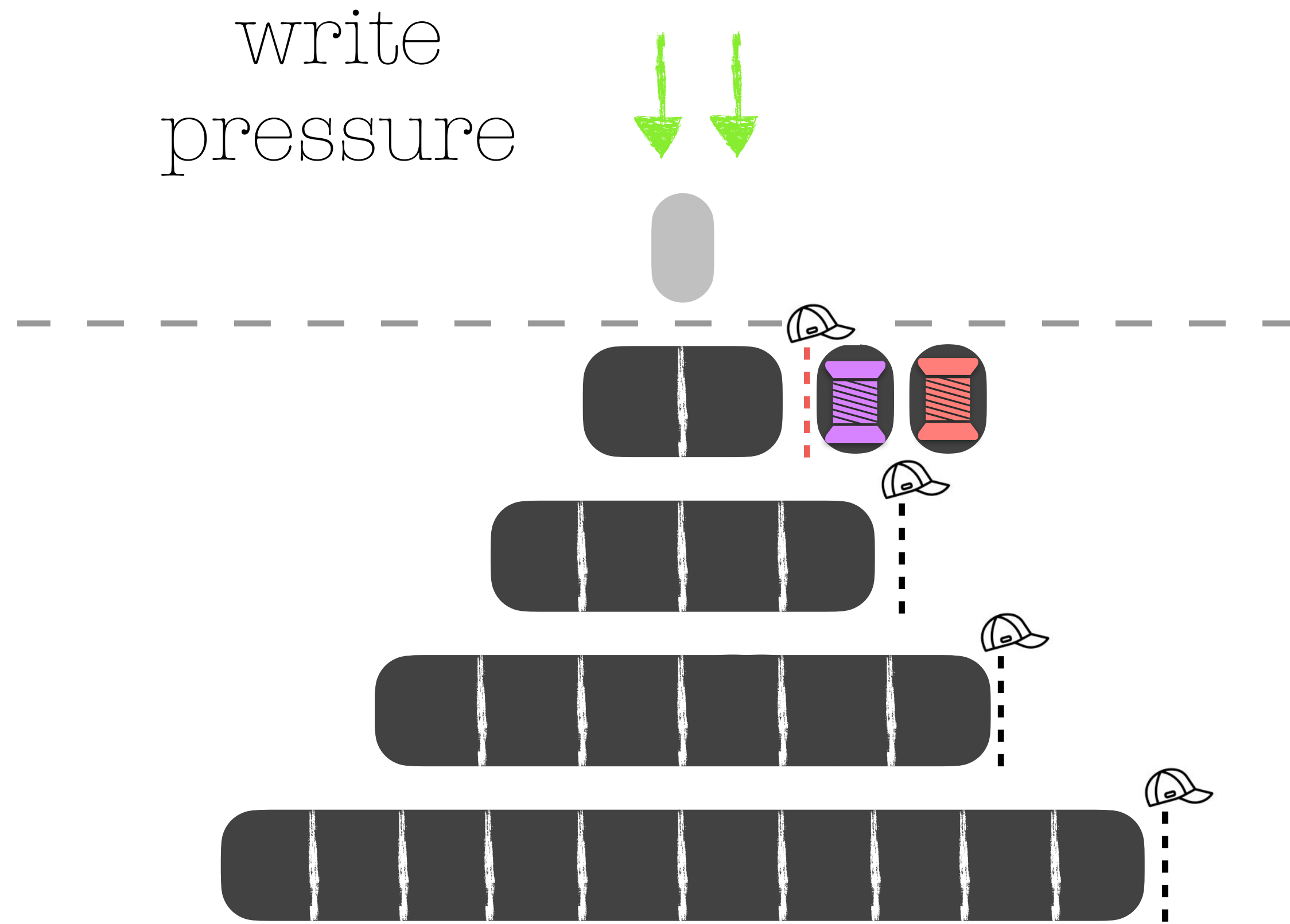


prioritize
writes over
compaction

Background
Compactions

Compaction
Priority

Optimizing Compactions

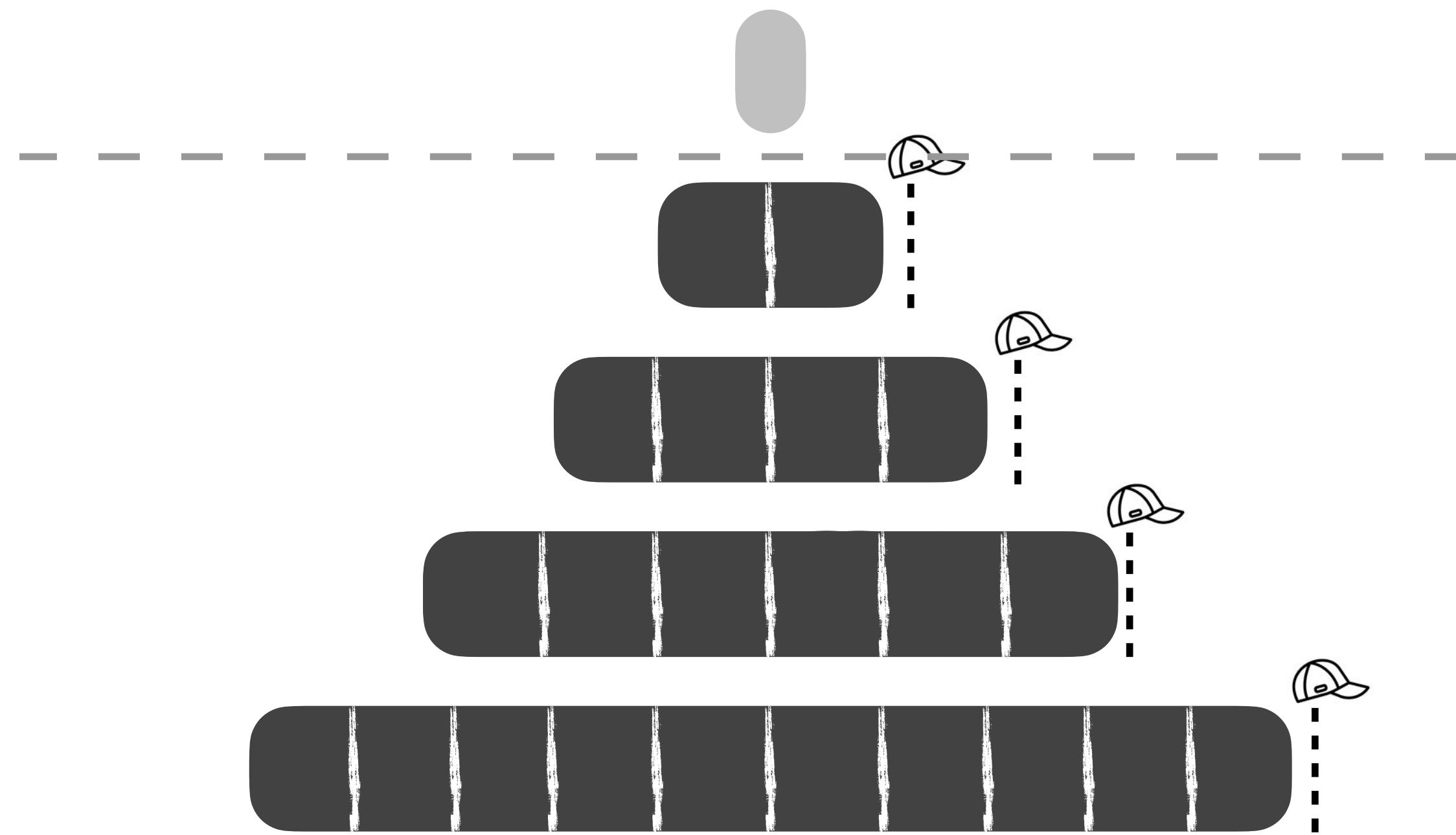


Background
Compactions

Compaction
Priority

- sustain heavy write bursts
- tree becomes out of shape

Optimizing **Compactions**

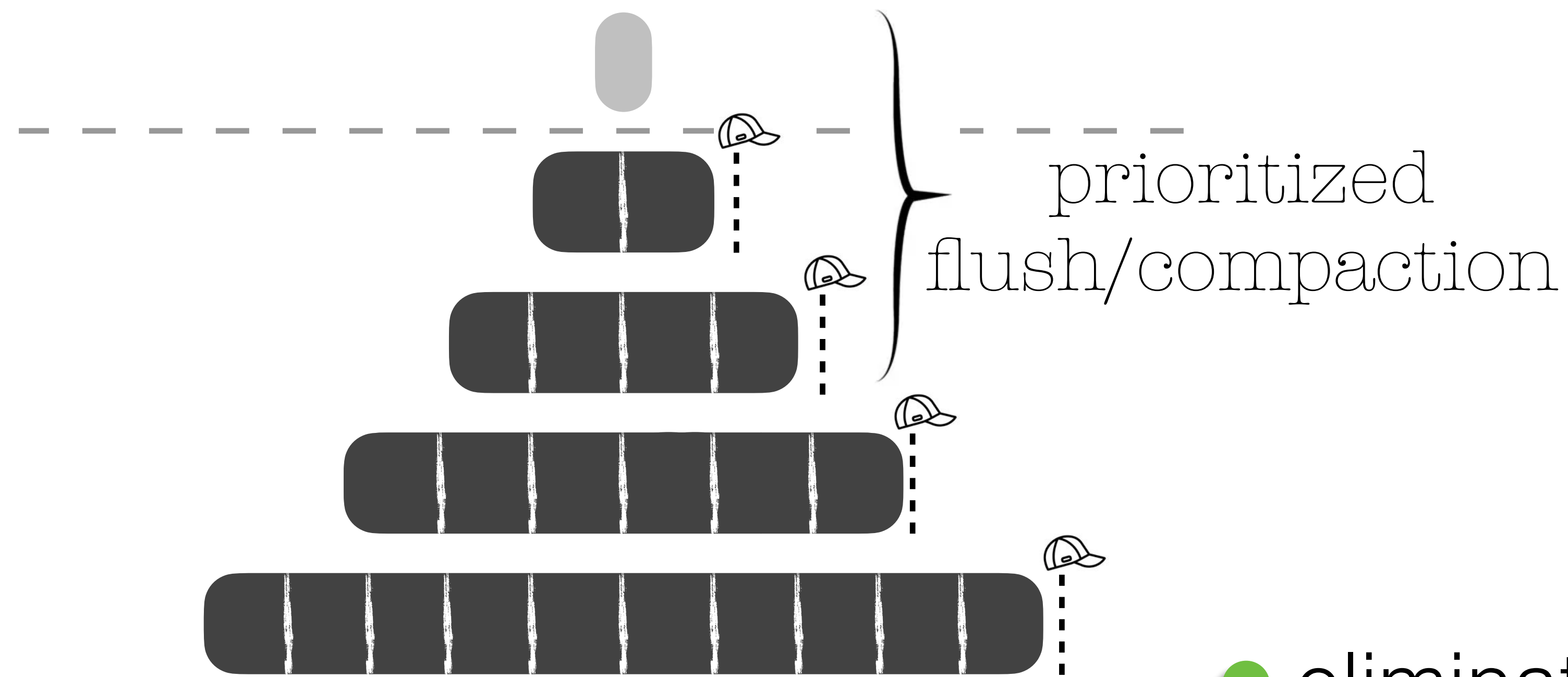


Background
Compactions

Compaction
Priority

I/O Scheduler

Optimizing **Compactions**



BalmauATC19

BalmauToCS20

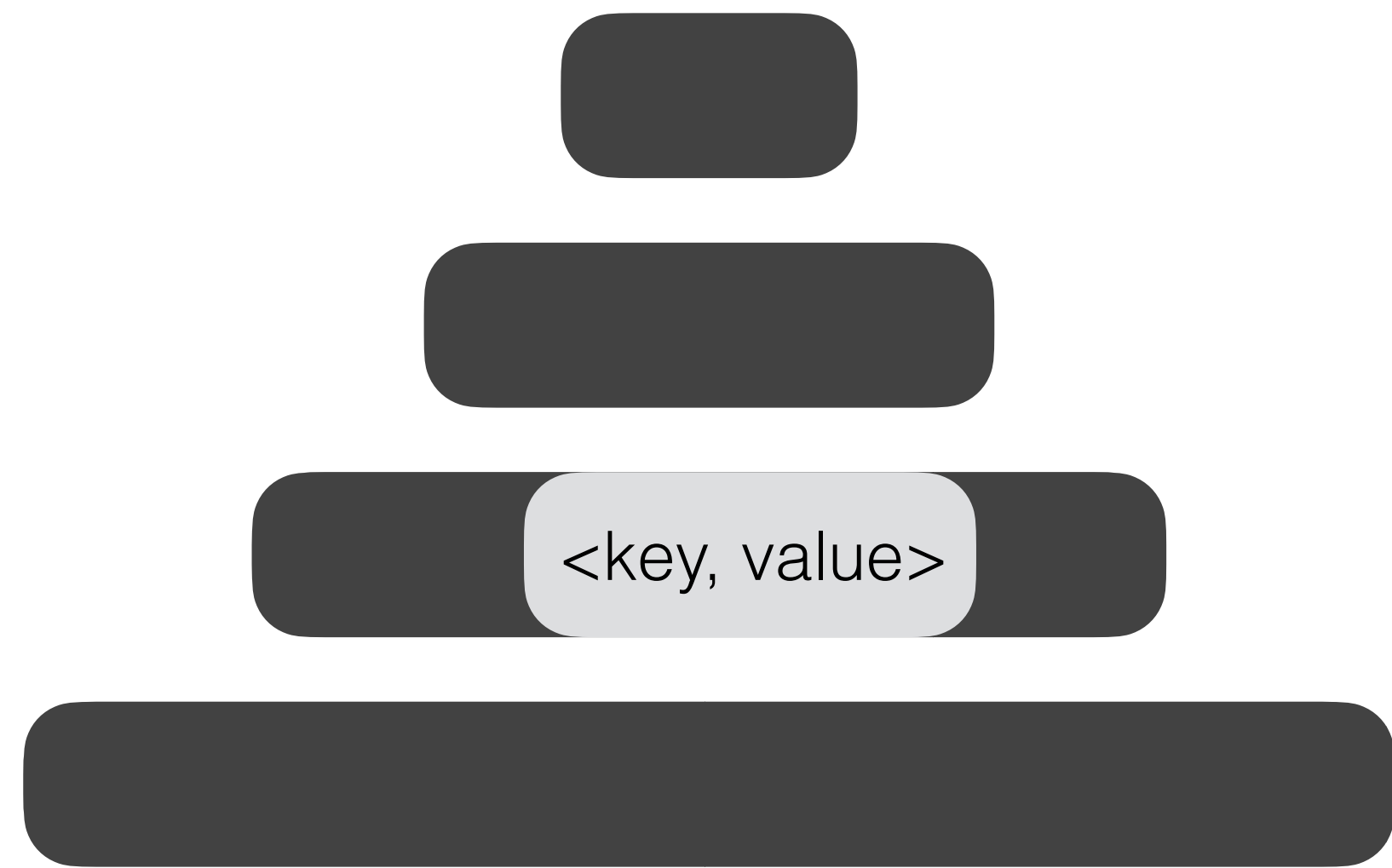
Background
Compactions

Compaction
Priority

I/O Scheduler

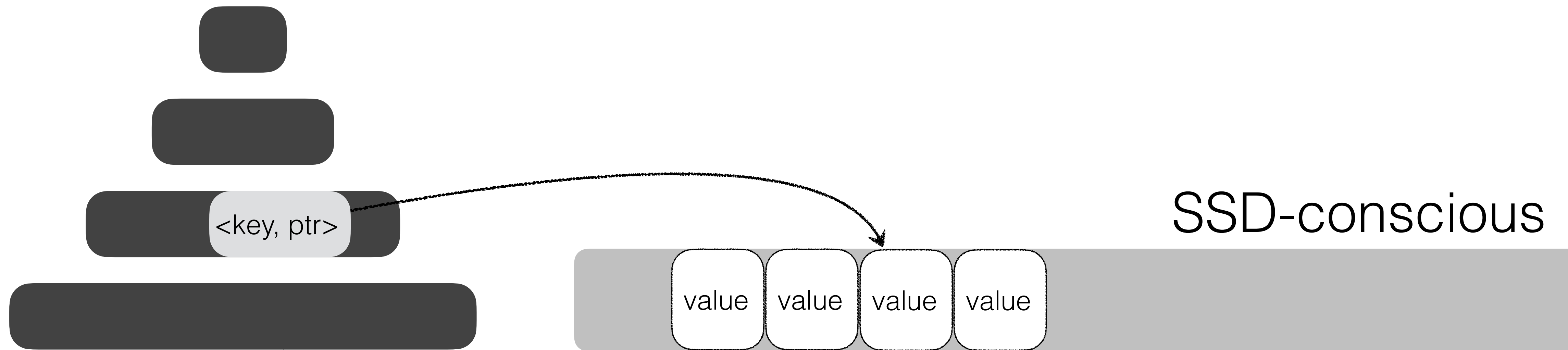
- eliminates write stalls
- no unnecessary high-priority compactions in lower levels

Data Placement Variations



key-value separation 
LuFAST16

Data Placement Variations

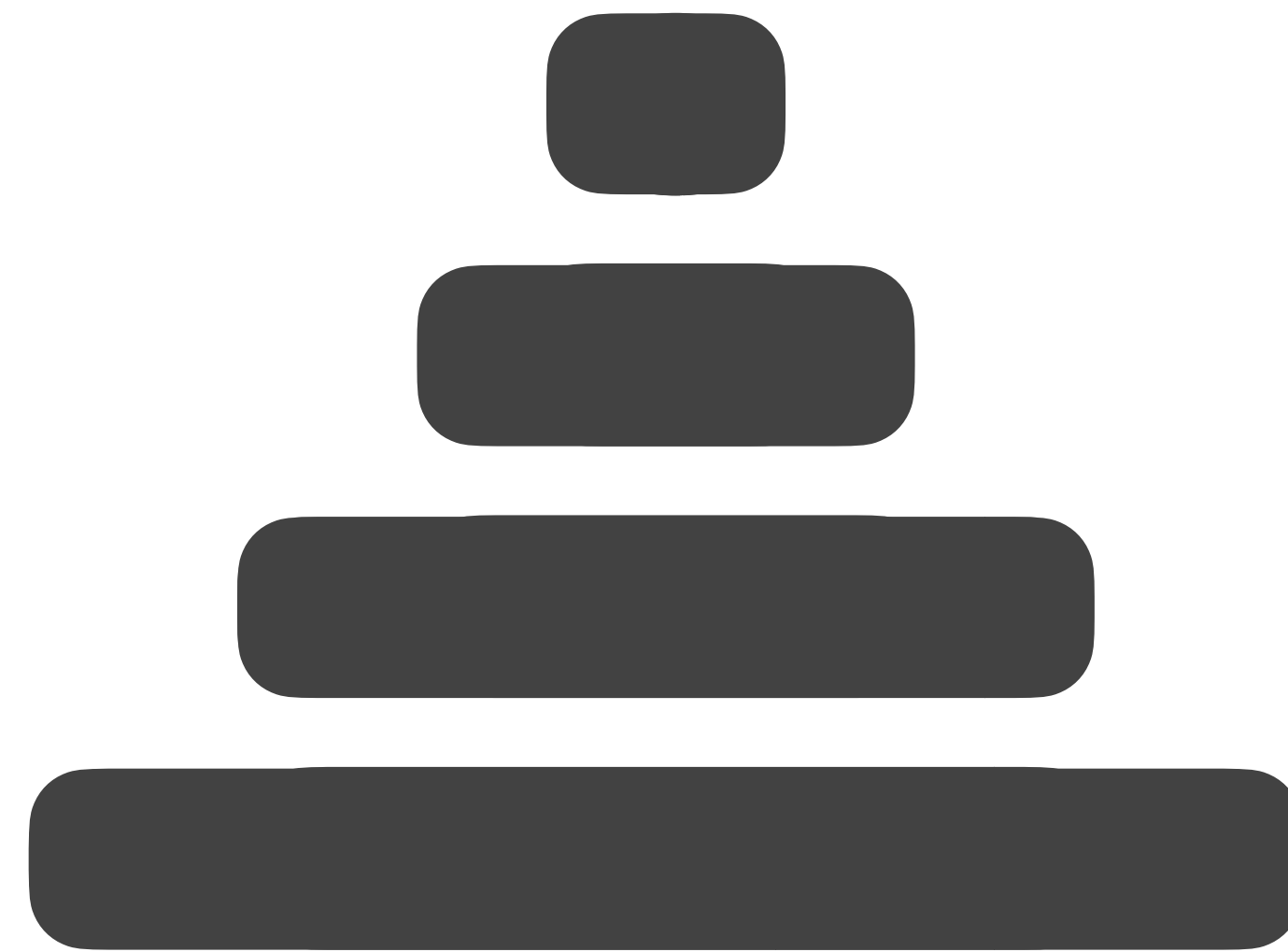


key-value separation 

LuFAST16

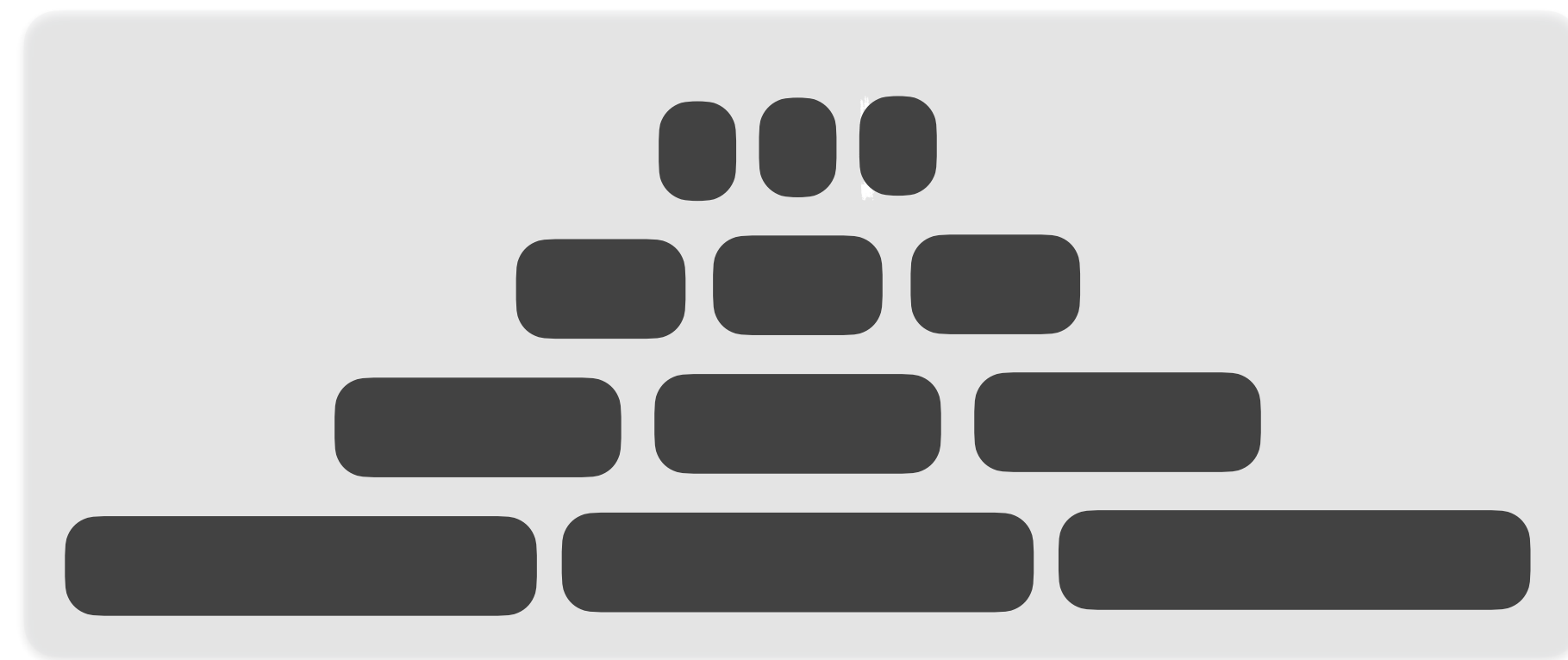
- reduced write amplification
- better read performance

Data Placement Variations



partitioning / sharding

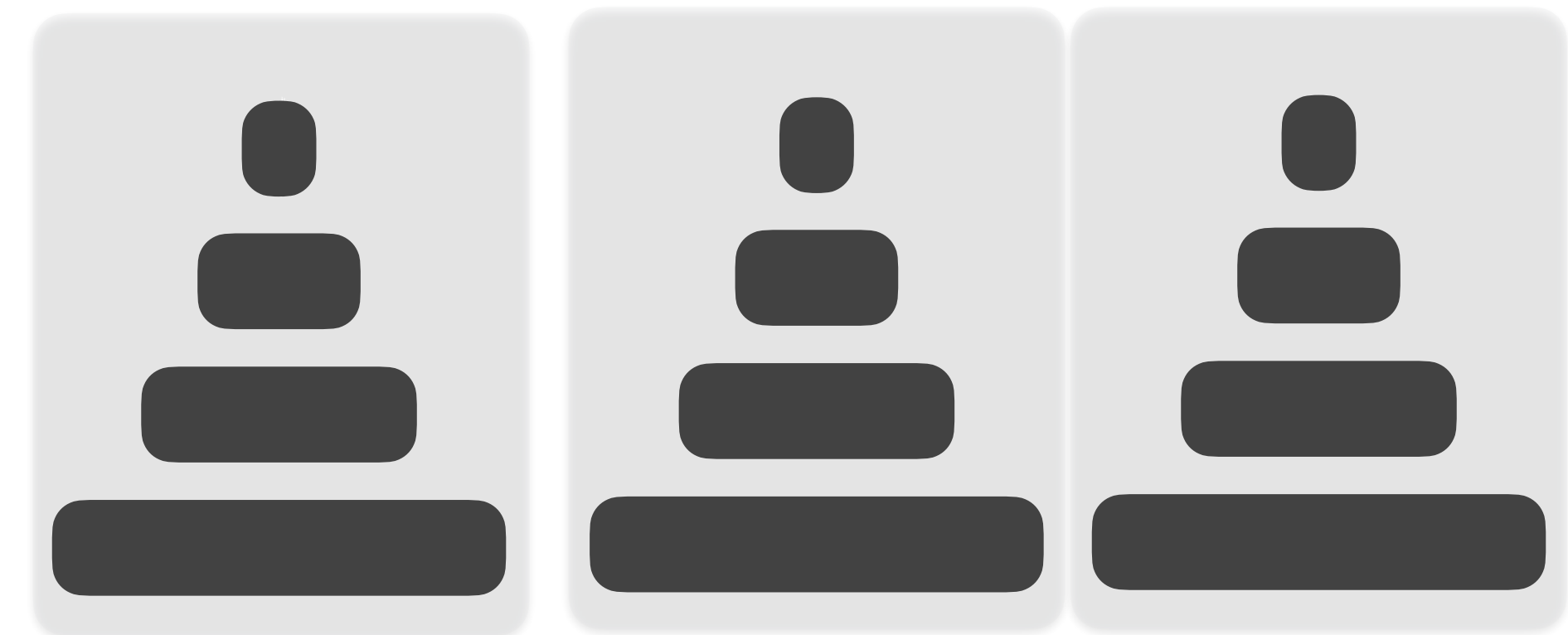
Data Placement Variations



storage

partitioning

RajuSOSP17



storage-1

storage-2

storage-3

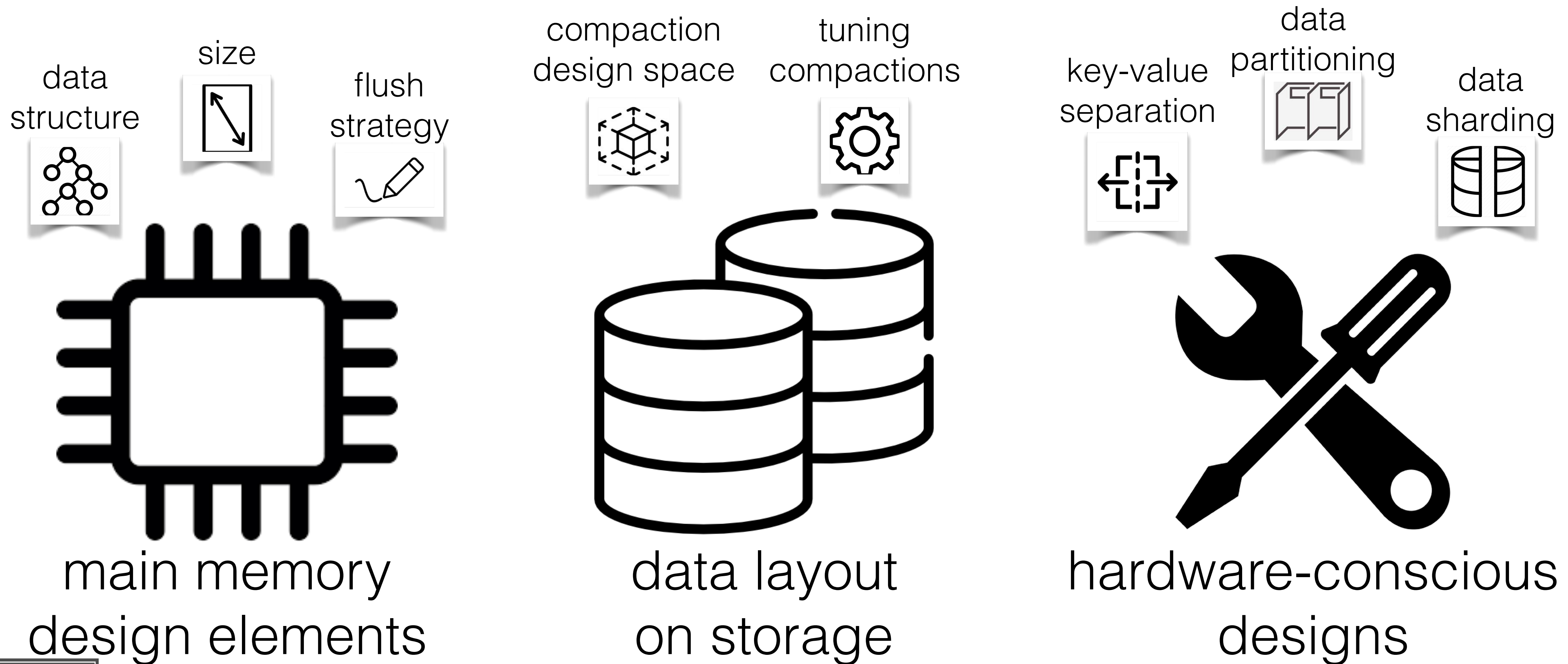
sharding

HuangSIGMOD21



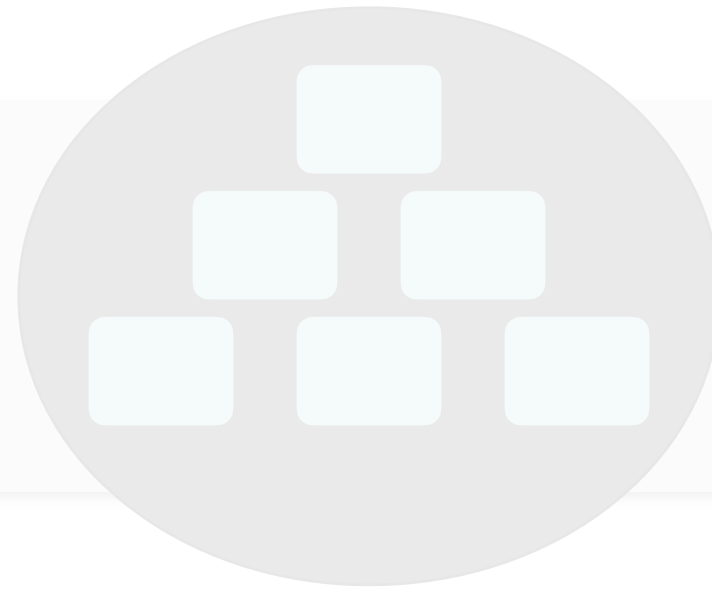
- improved ingestion throughput
- reduced write amplification

Summary: Ingestion Optimization

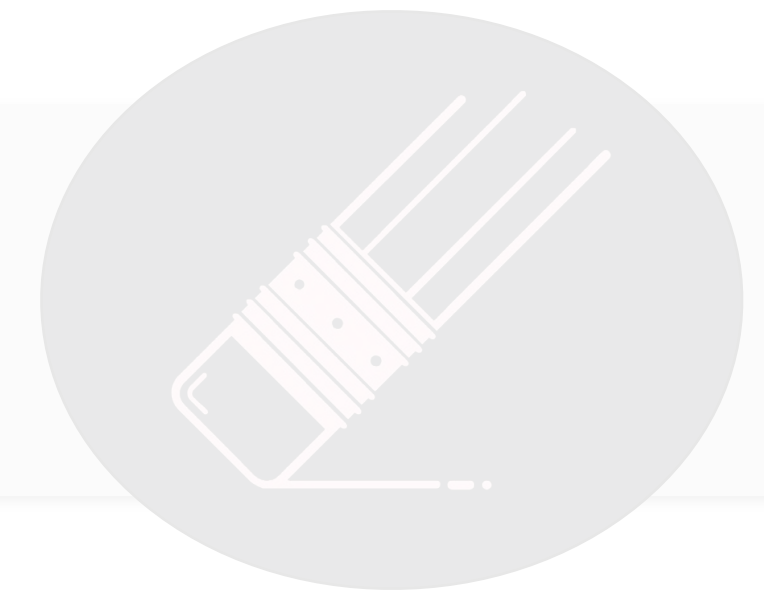


Outline

Part 1: **LSM Basics**



Part 2: **Optimizing Ingestion in LSMs**

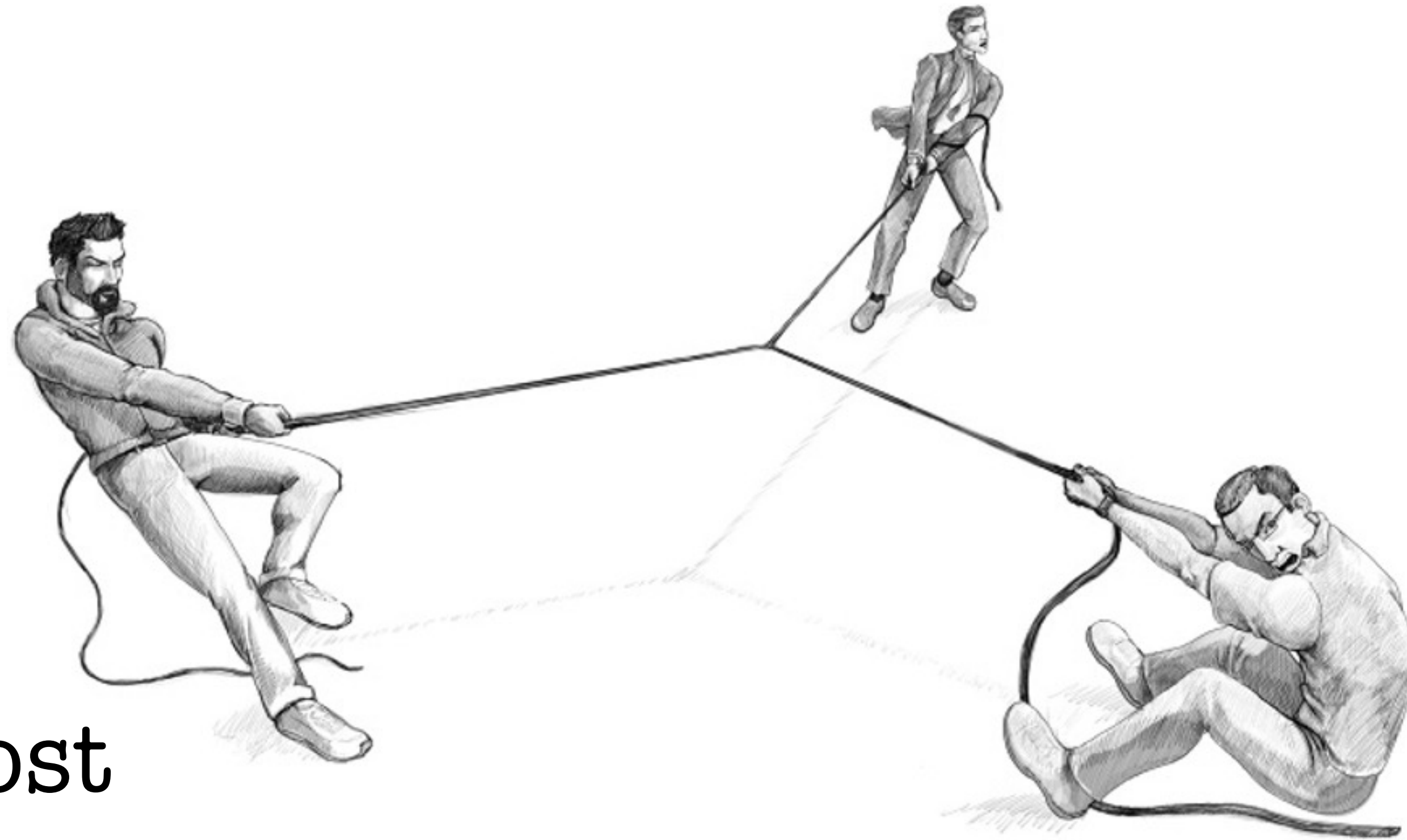


Part 3: **Navigating the LSM Design Space**



LSM Design Space

Update cost



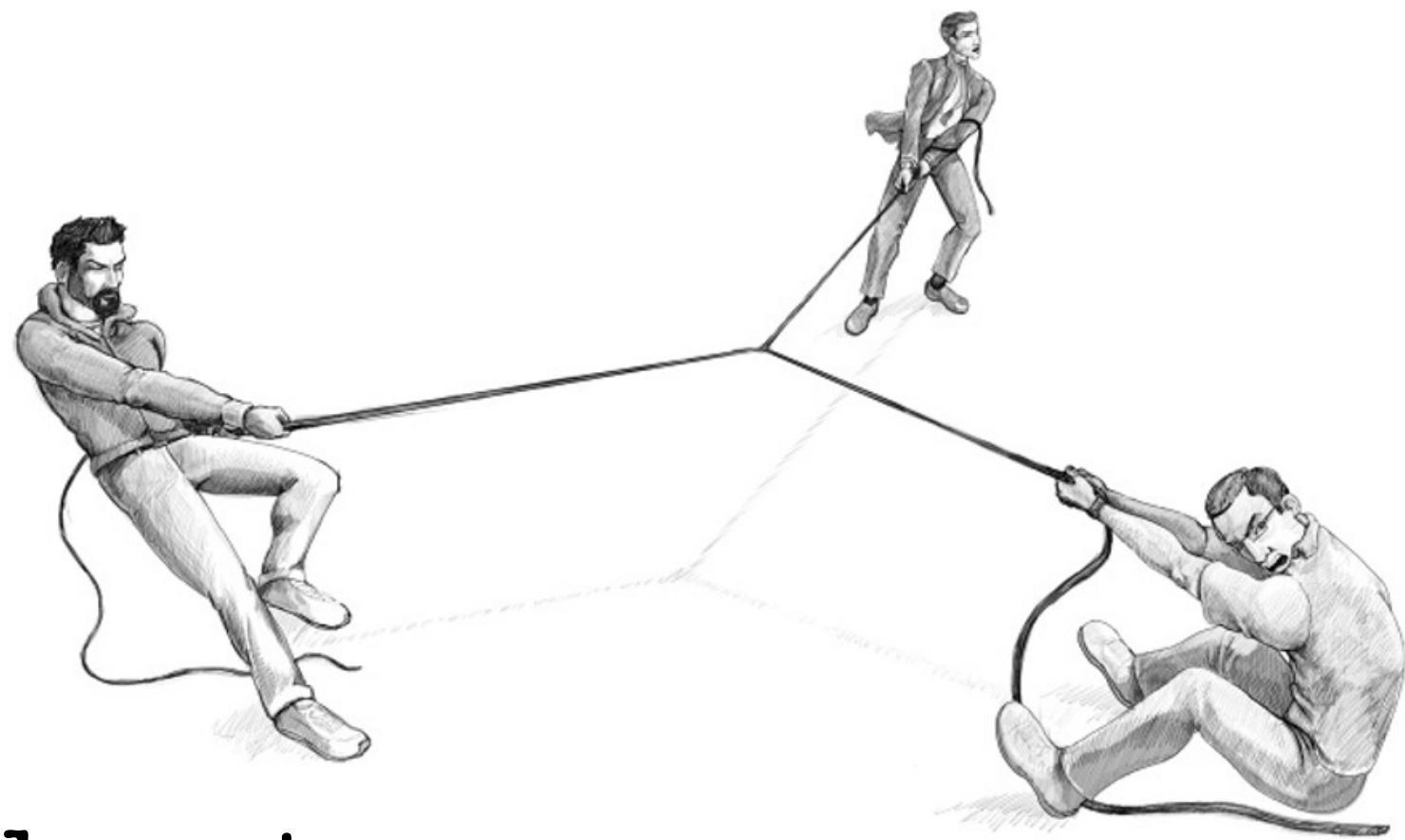
Read cost

Memory/space footprint

LSM Design Space

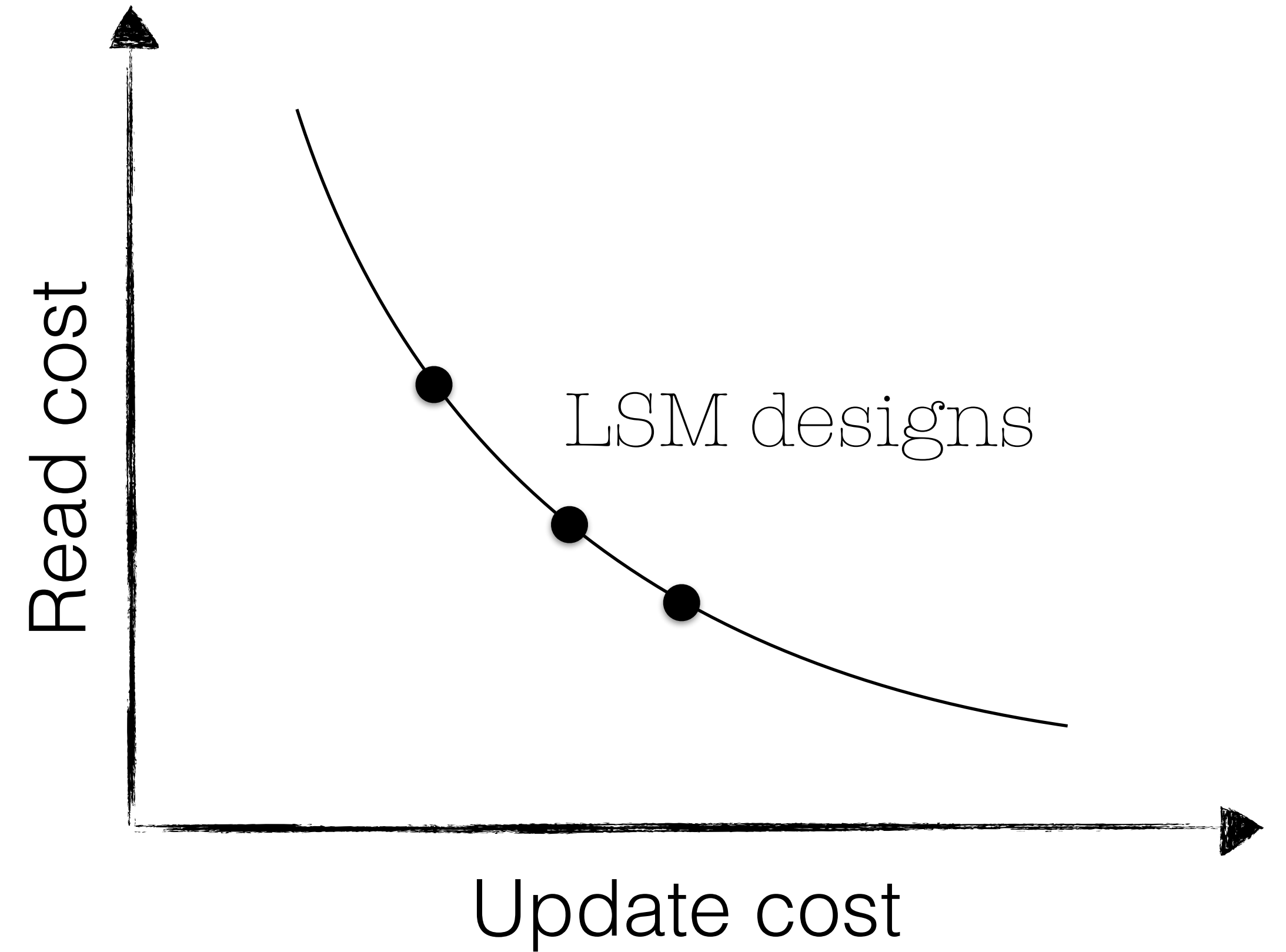
fixed Memory

Update cost

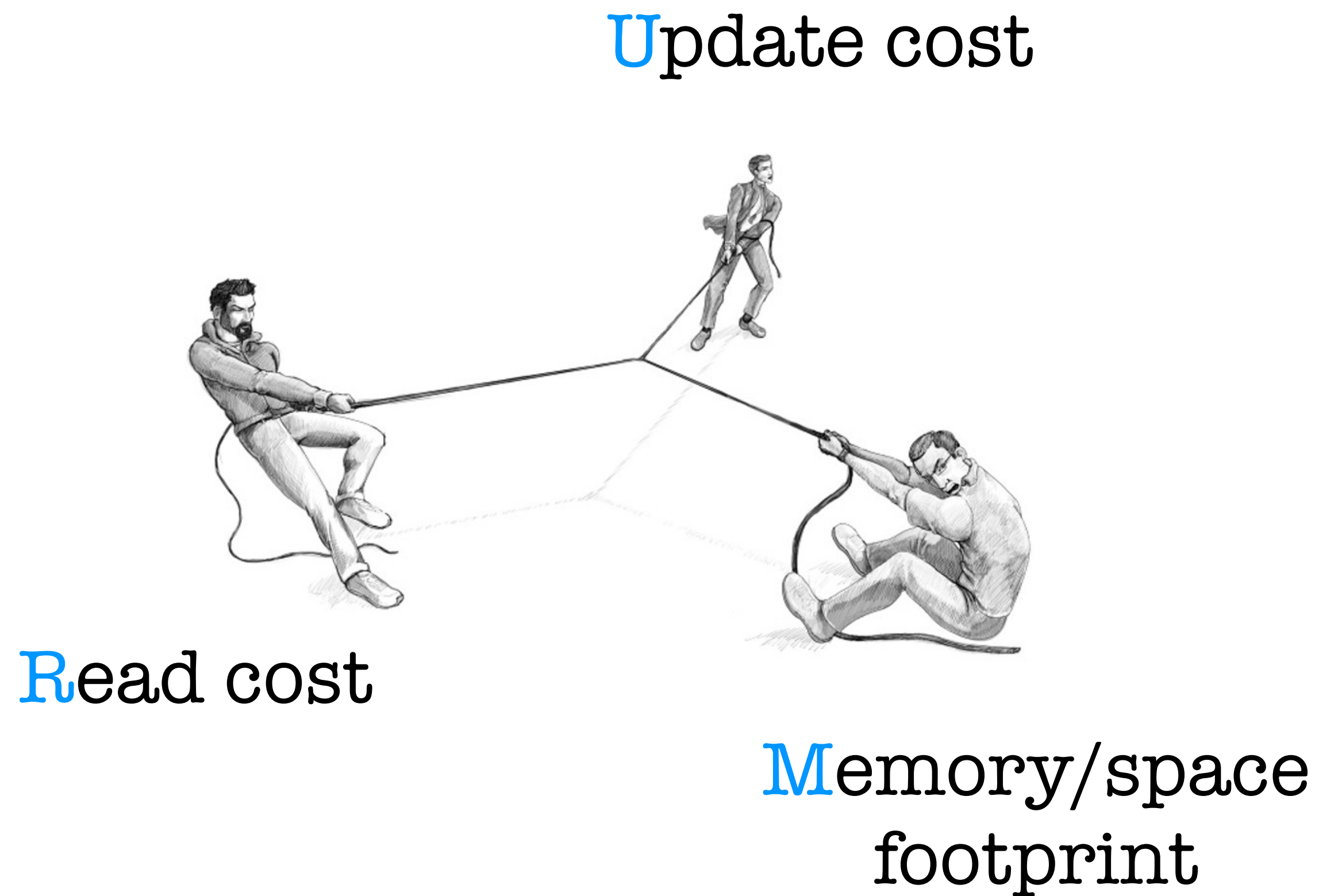


Read cost

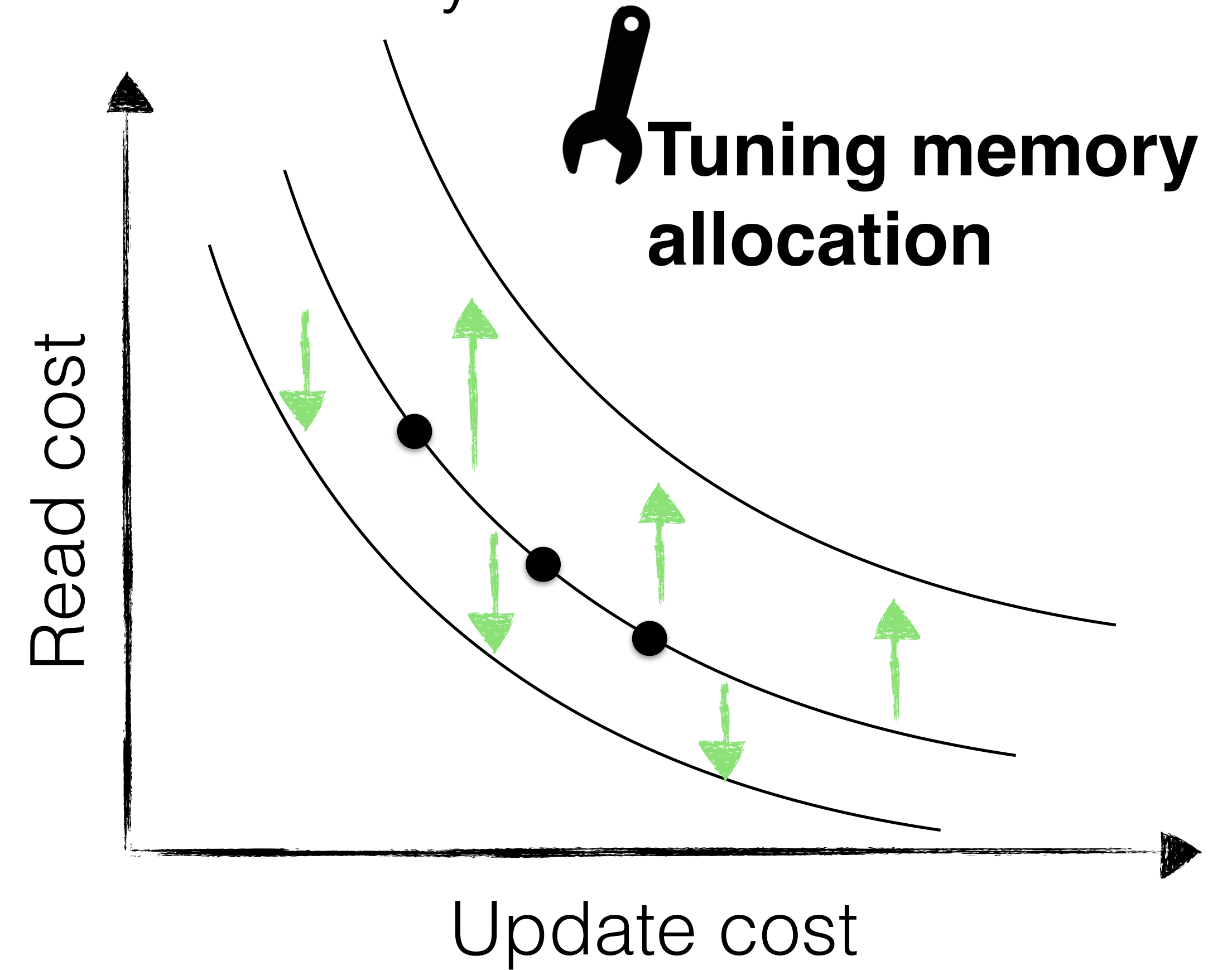
Memory/space footprint



LSM Design Space

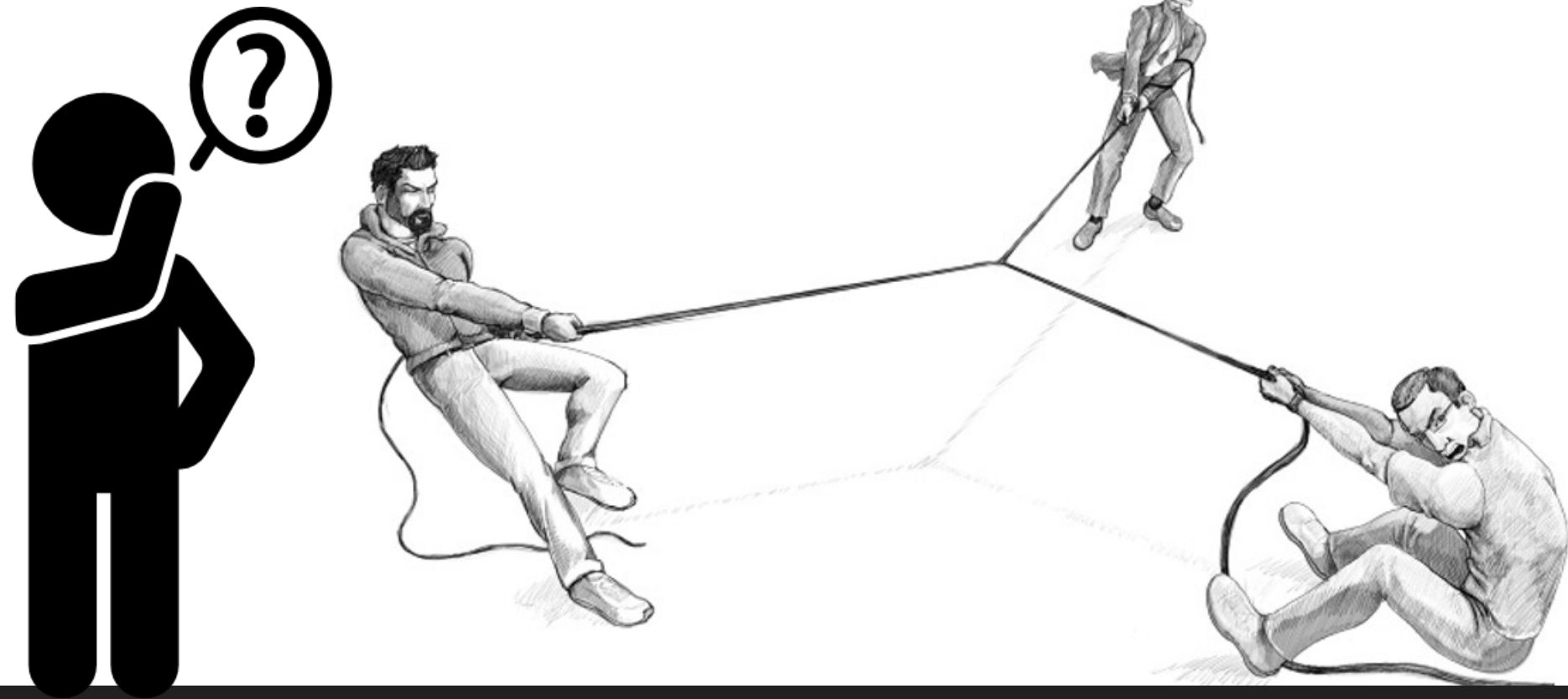


fixed Memory

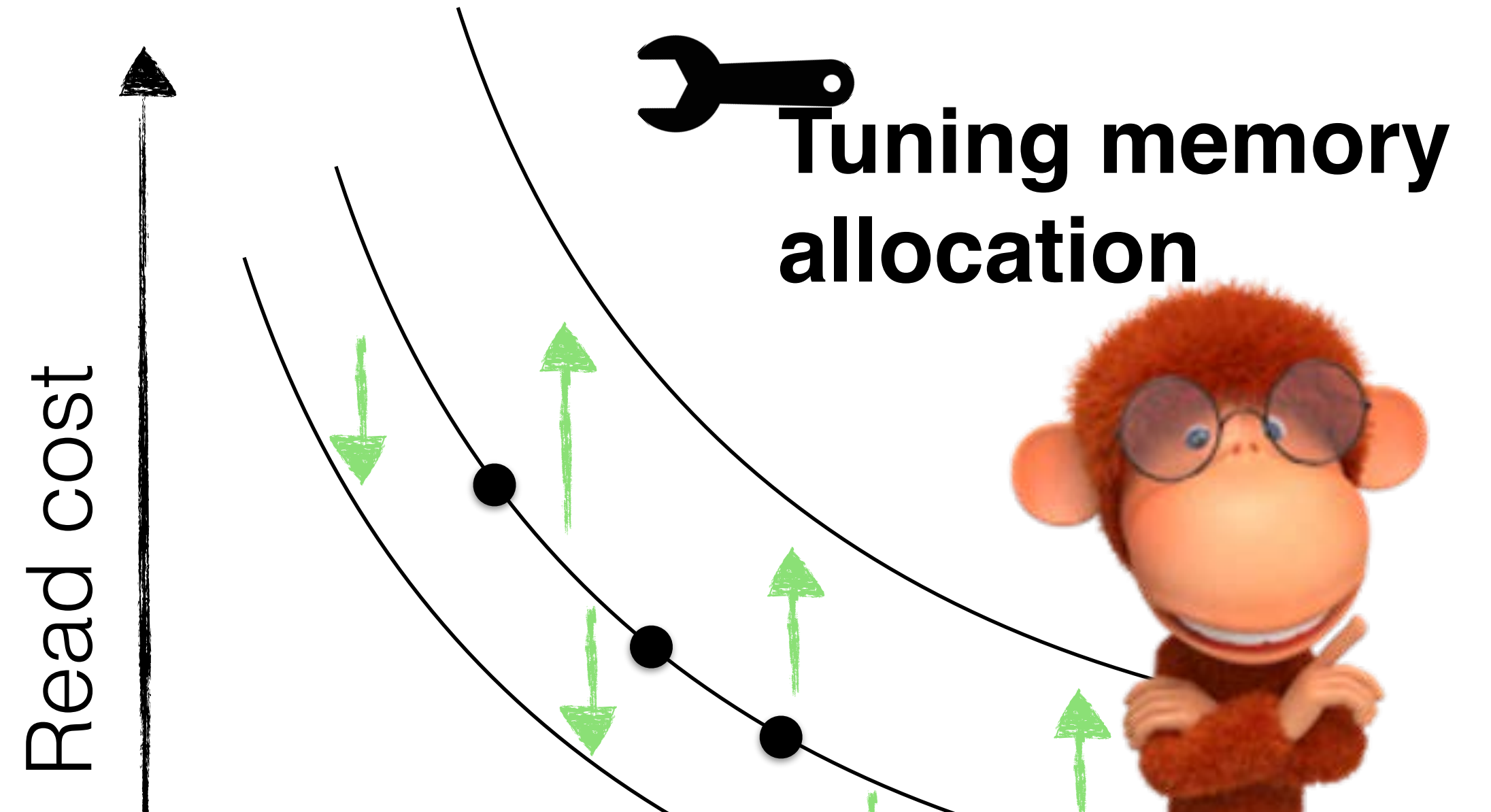


LSM Design Space

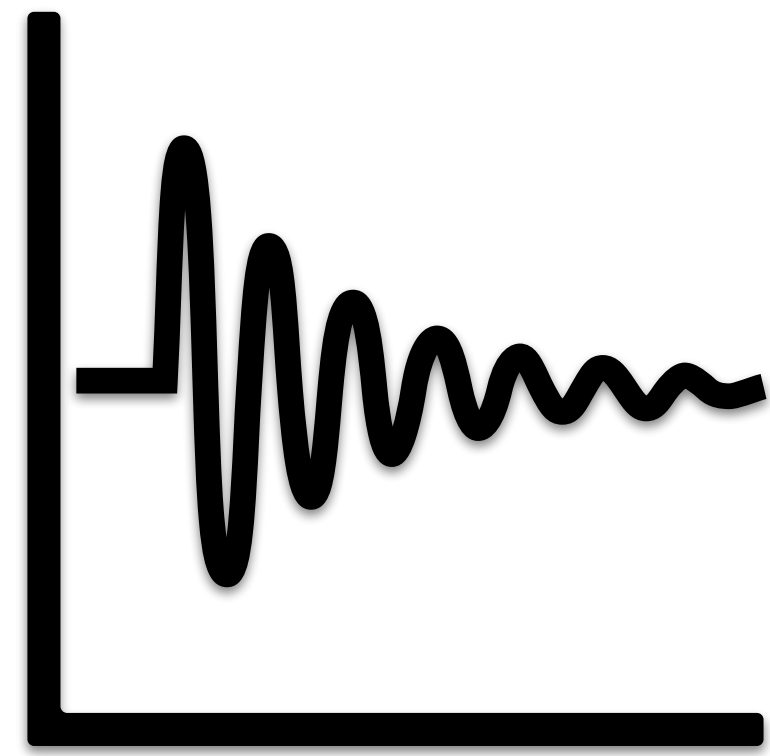
Update cost



fixed Memory



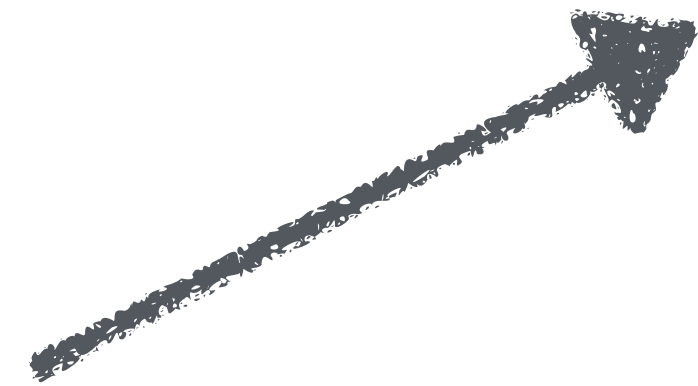
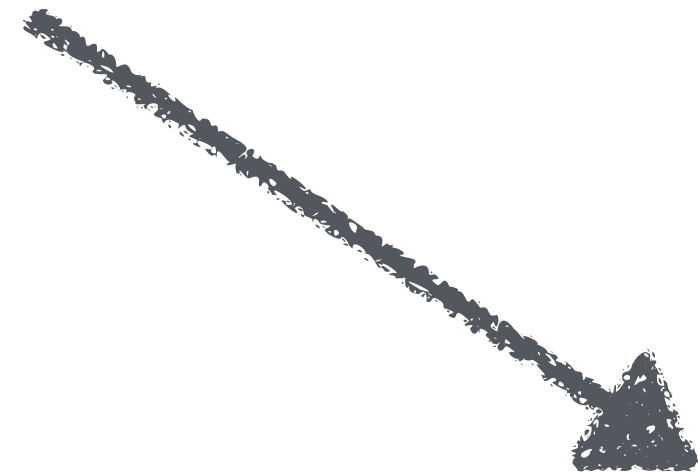
How to optimally allocate the available memory?



workload



memory
budget



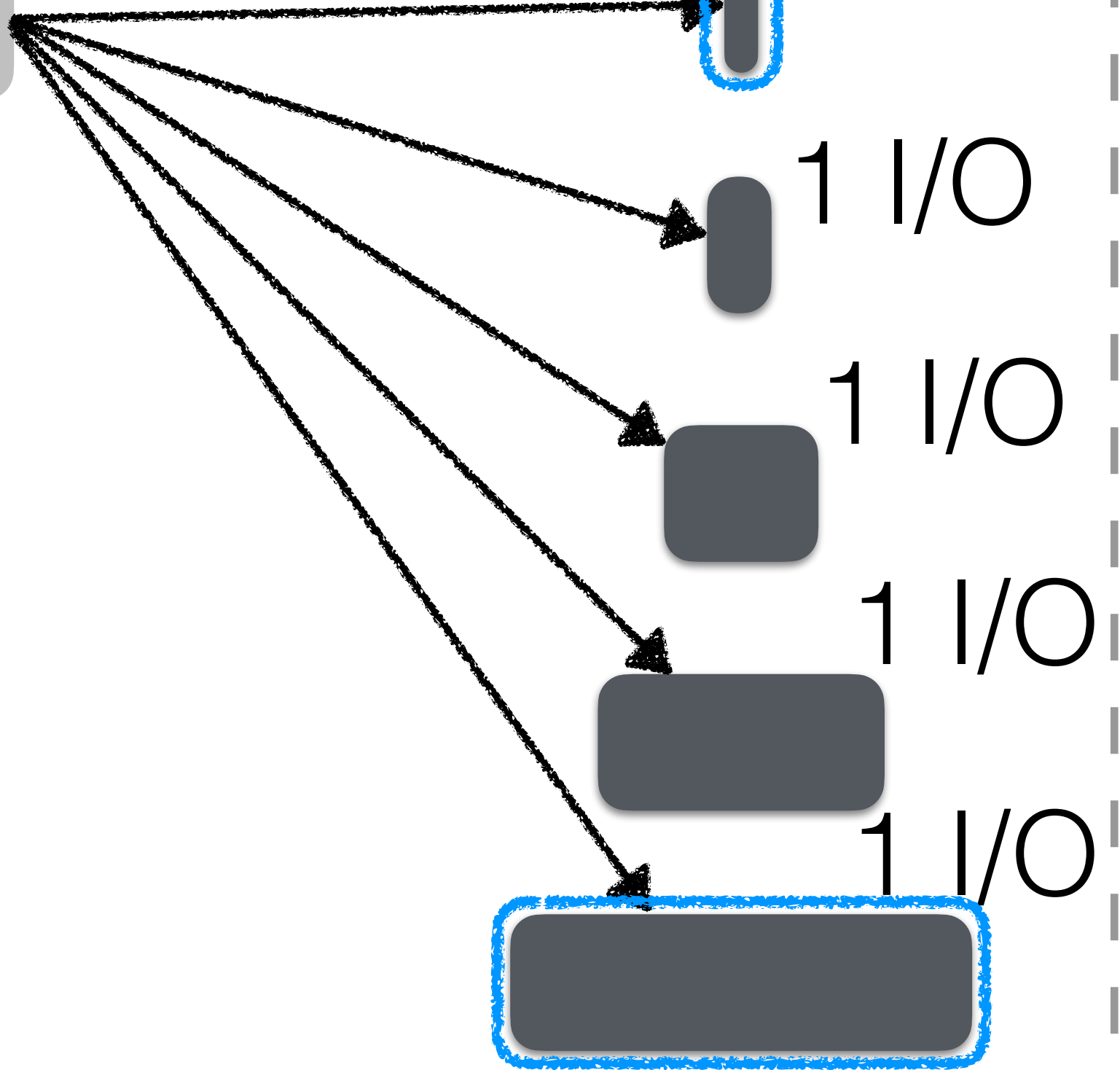
How to allocate memory
between buffer and BF

How to allocate memory
among BF's in LSM

get(7)



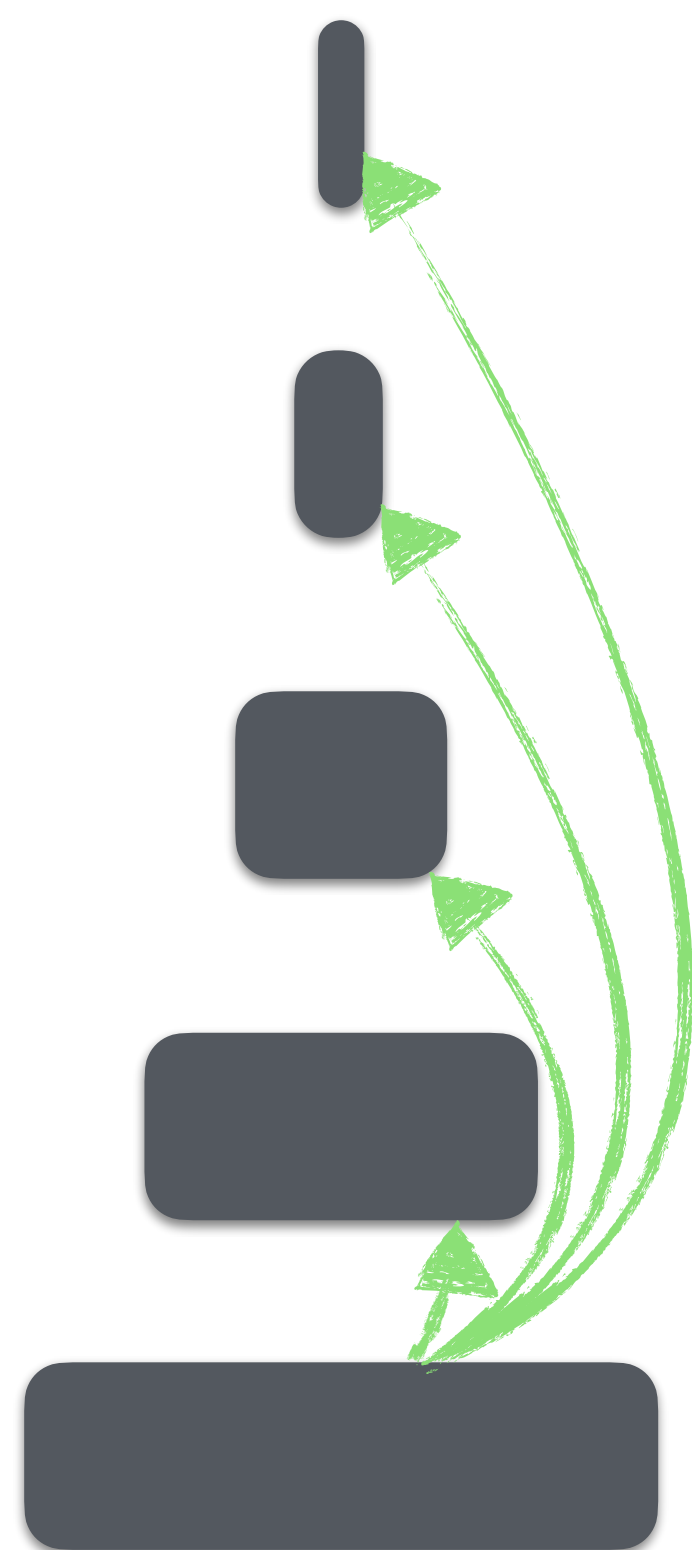
buffer



Bloom filters



Bloom
filters



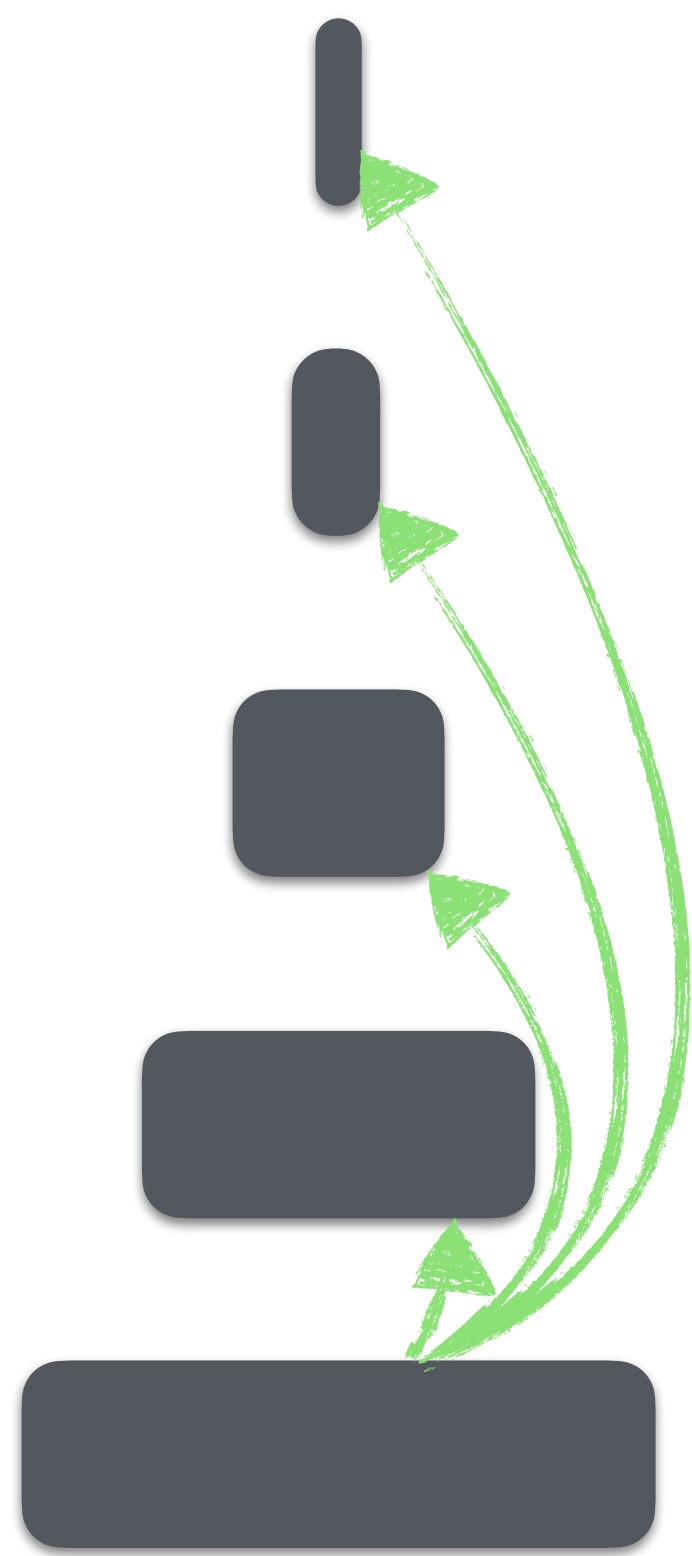
FPR Monkey
FPR

ϕ	$>$	$\phi_1 = \phi_0 / T^4$
ϕ	$>$	$\phi_2 = \phi_0 / T^3$
ϕ	$>$	$\phi_3 = \phi_0 / T^2$
ϕ	$>$	$\phi_4 = \phi_0 / T$
ϕ	$<$	$\phi_5 = \phi_0$

↑
exponentially
decreasing



Bloom filters



FPR		Monkey FPR
ϕ	>	$\phi_1 = \phi_0 / T^4$
ϕ	>	$\phi_2 = \phi_0 / T^3$
ϕ	>	$\phi_3 = \phi_0 / T^2$
ϕ	>	$\phi_4 = \phi_0 / T$
ϕ	<	$\phi_5 = \phi_0$

↑ exponentially decreasing

point lookup cost

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

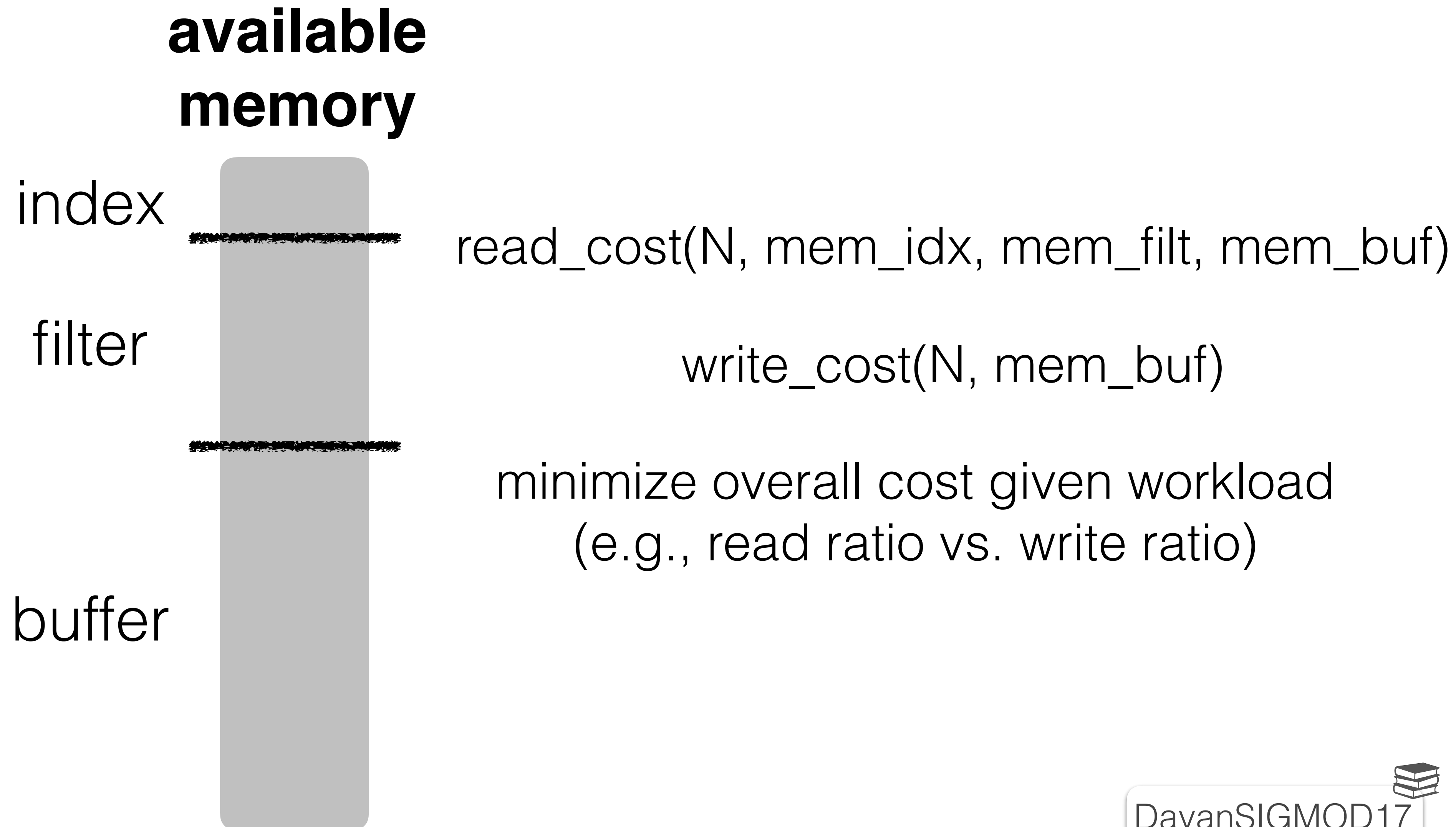
V

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

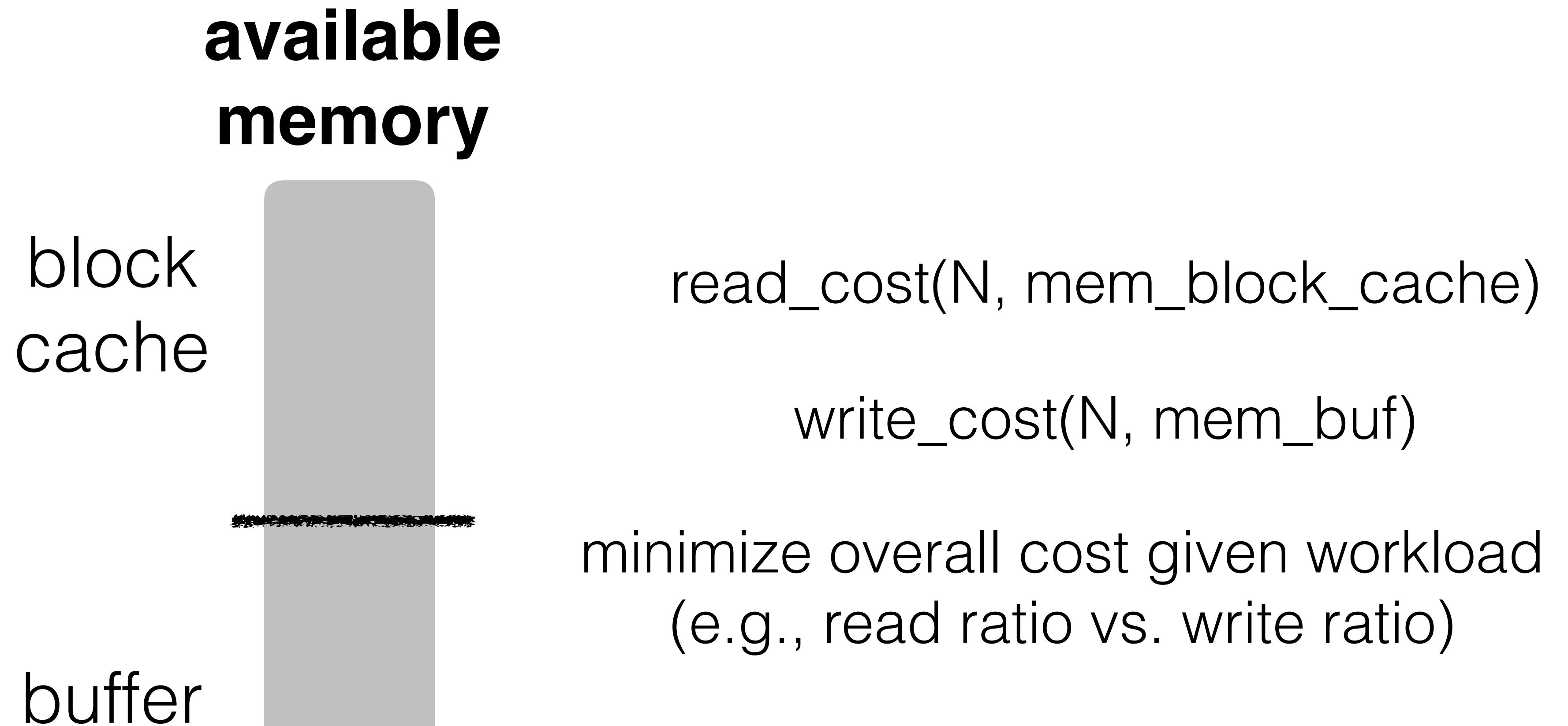


$$L \cdot \phi \quad \sum_i \phi_i = c \cdot \phi_0$$

The **Optimal** Memory Allocation

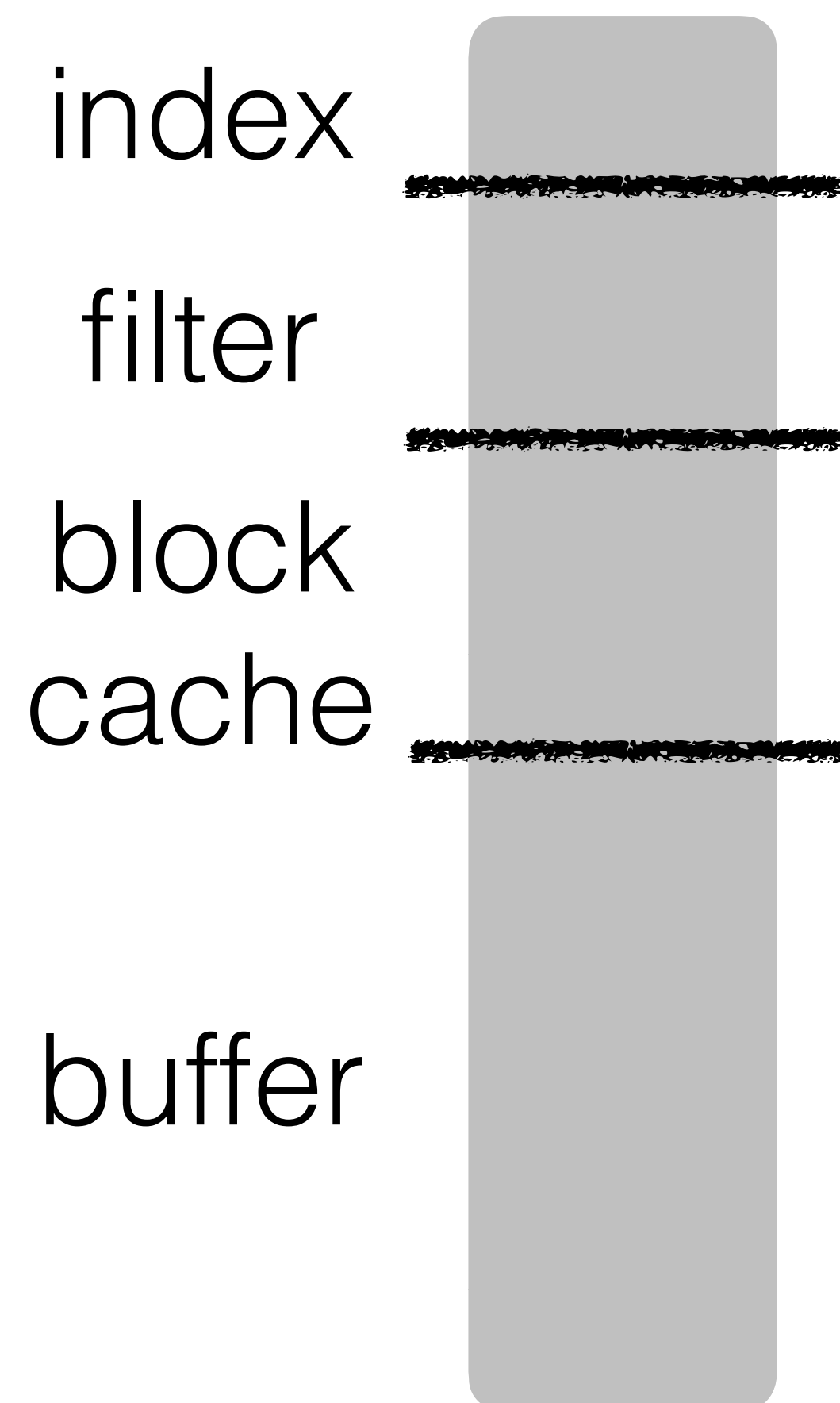


The **Optimal** Memory Allocation

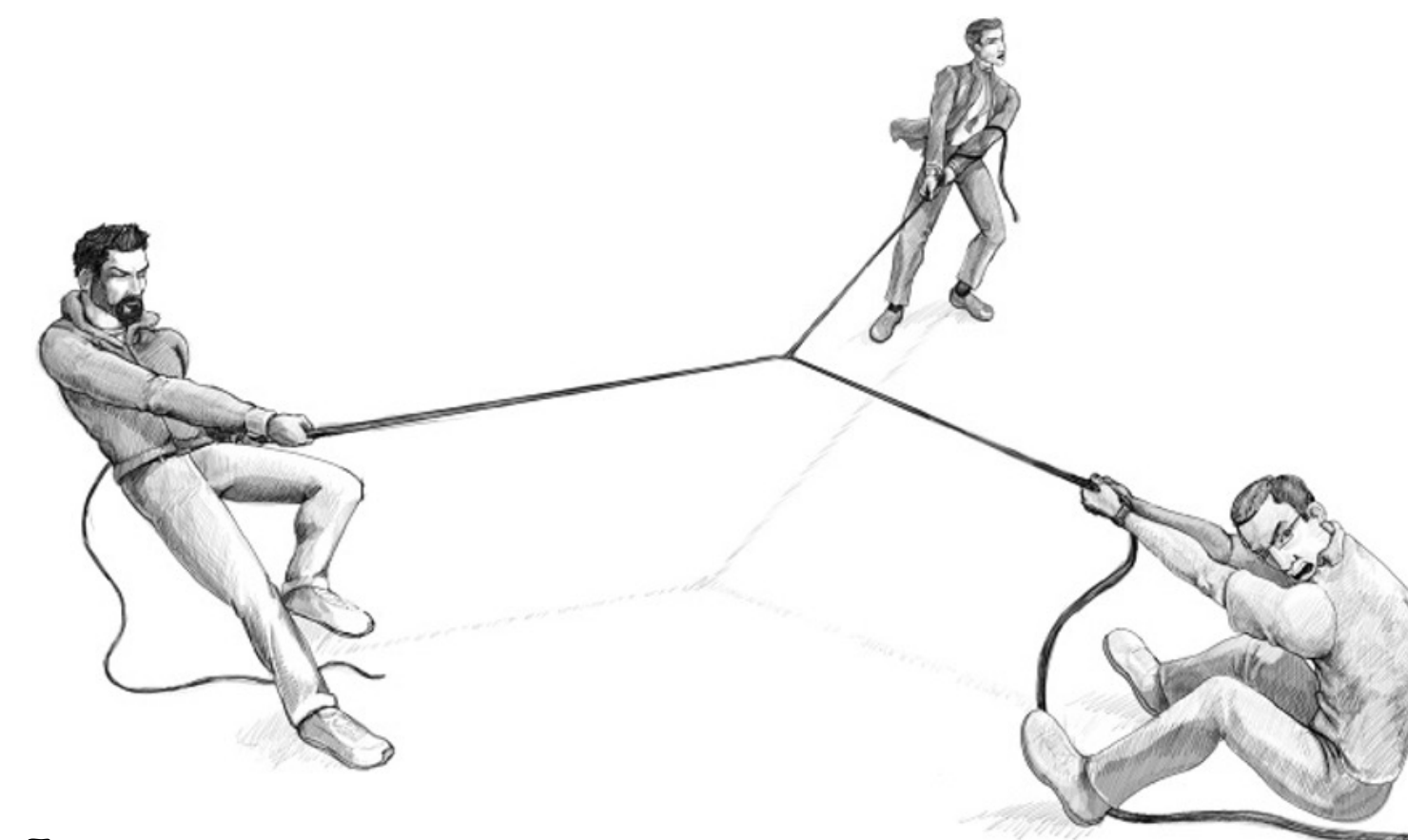


The **Optimal** Memory Allocation

**available
memory**



Uupdate cost

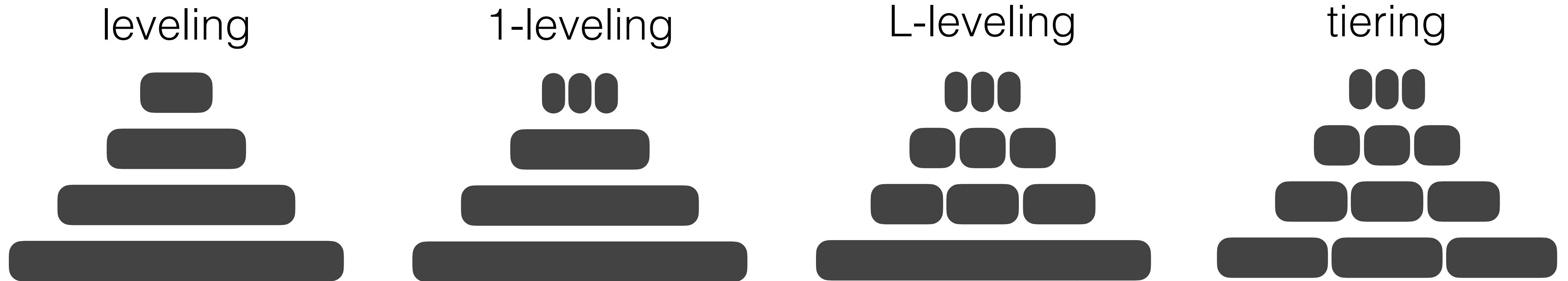


Read cost

Memory/space
footprint

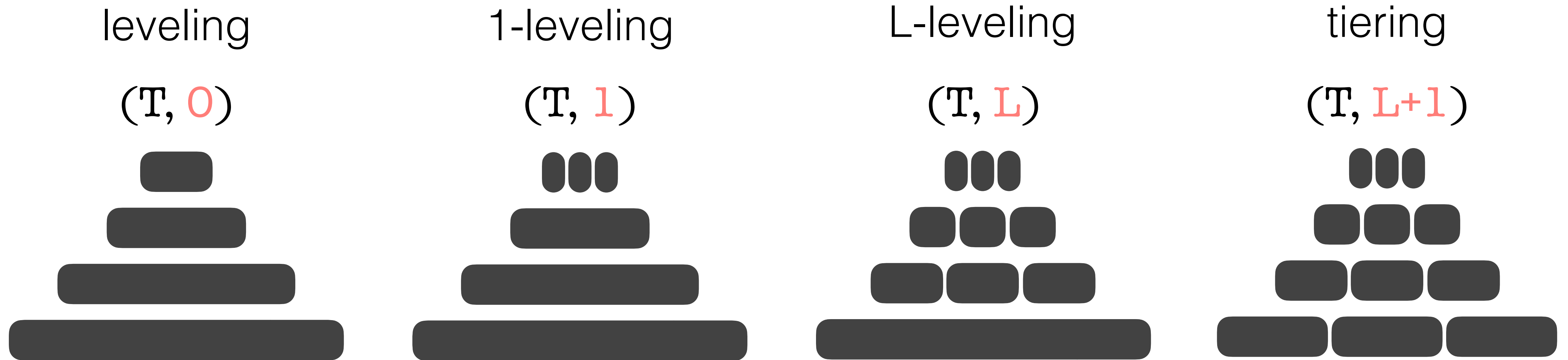
Holistic memory
allocation & tuning

Storage Layer **Design Continuum**



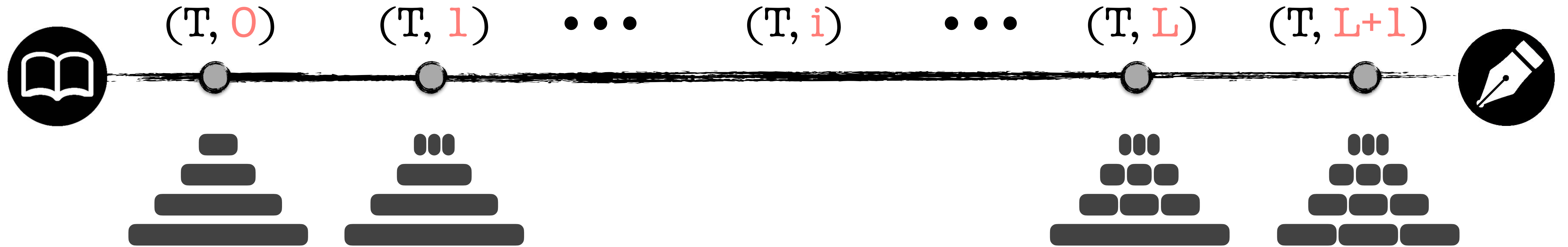
Any design can be defined by the tuple-set: (T, i)

Storage Layer **Design Continuum**

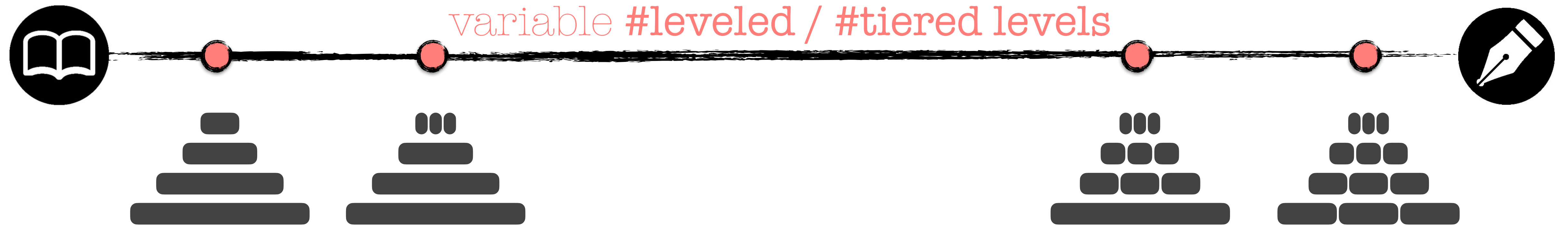


Any design can be defined by the tuple-set: (T, i)

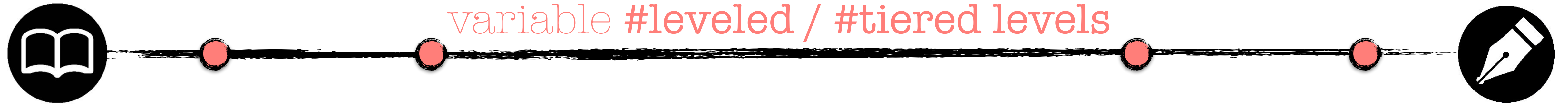
Storage Layer Design Continuum



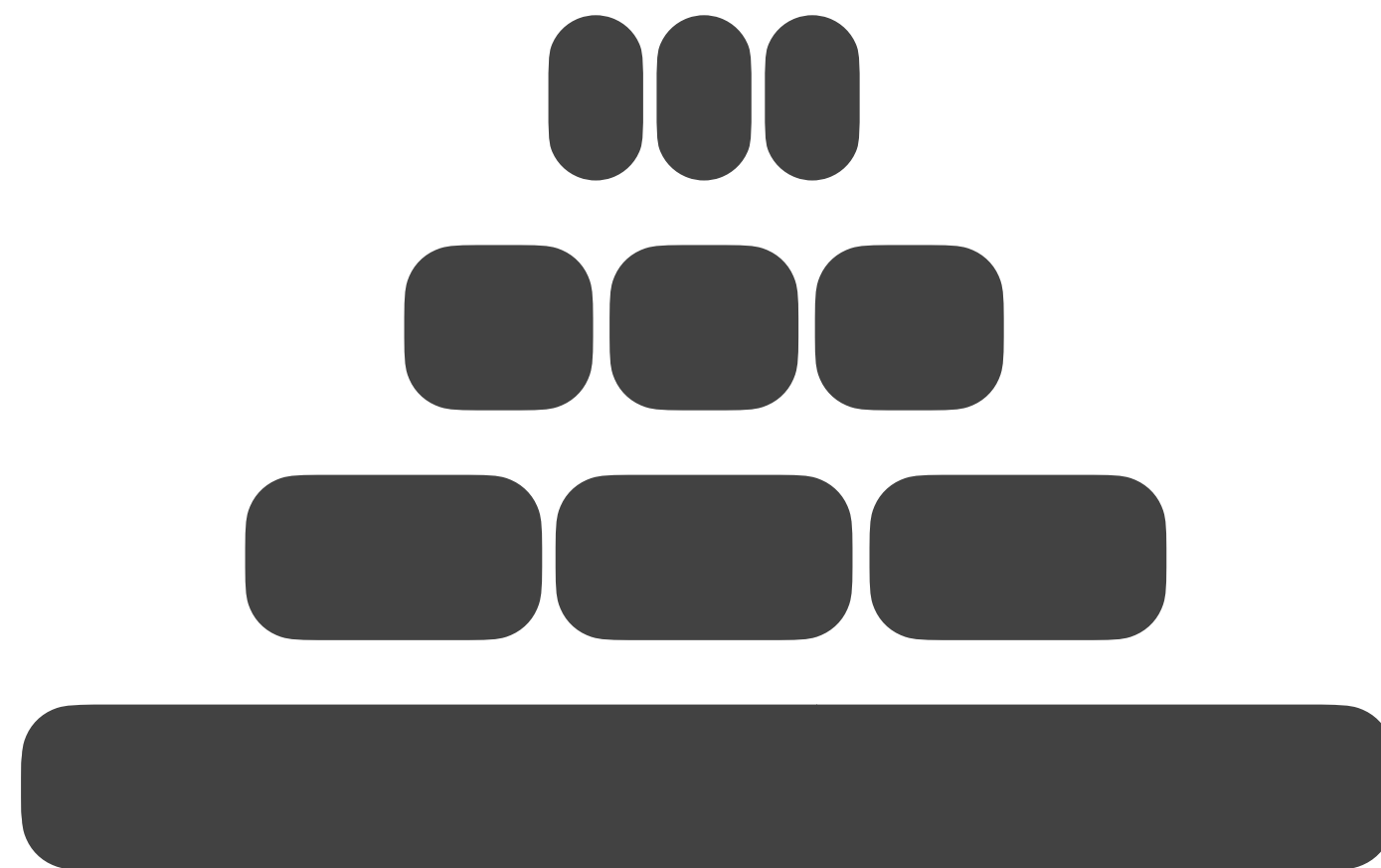
Storage Layer Design Continuum



Storage Layer Design Continuum



size ratio



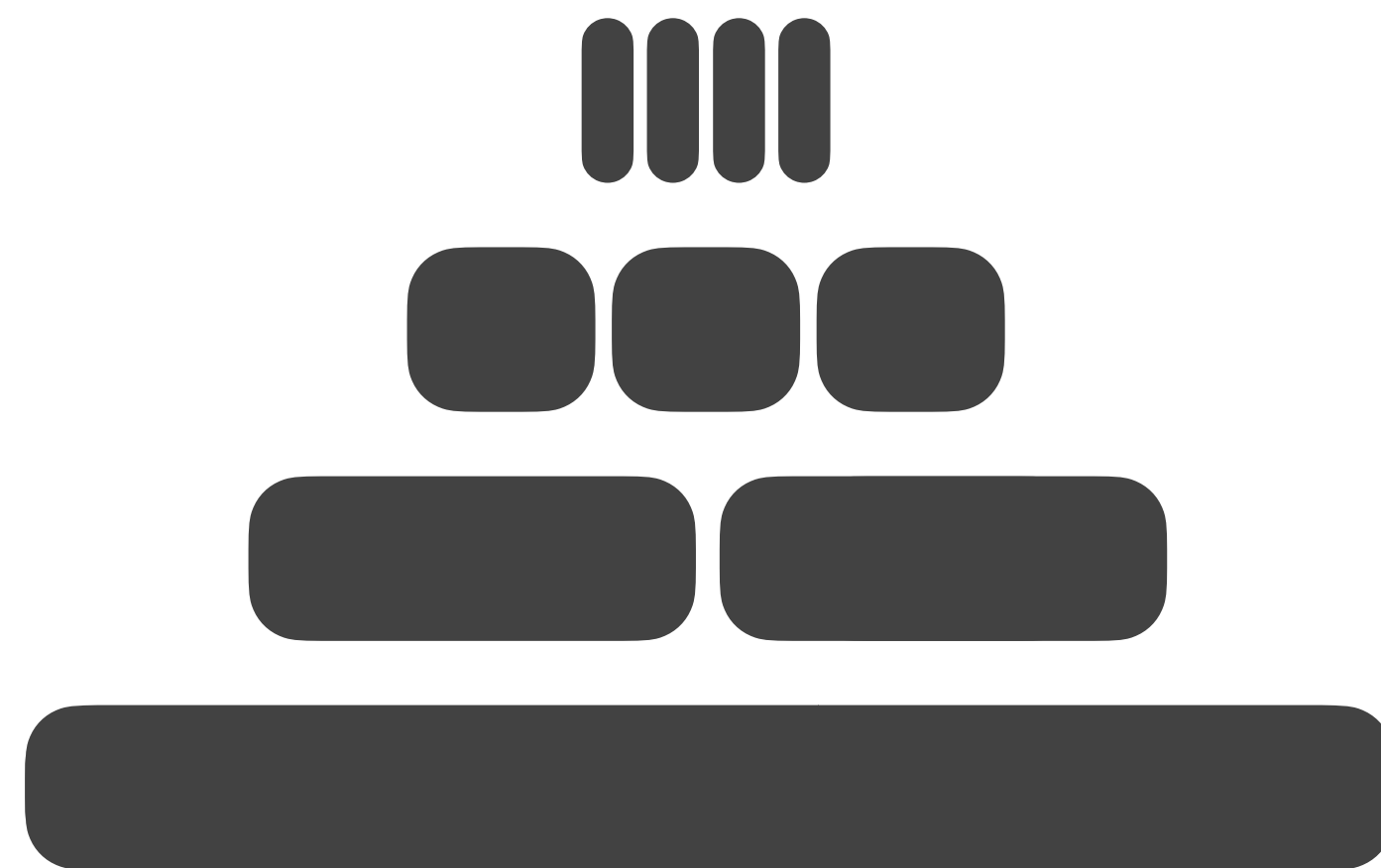
T

T

T

T

Storage Layer Design Continuum



size ratio

#runs

T

4

T

3

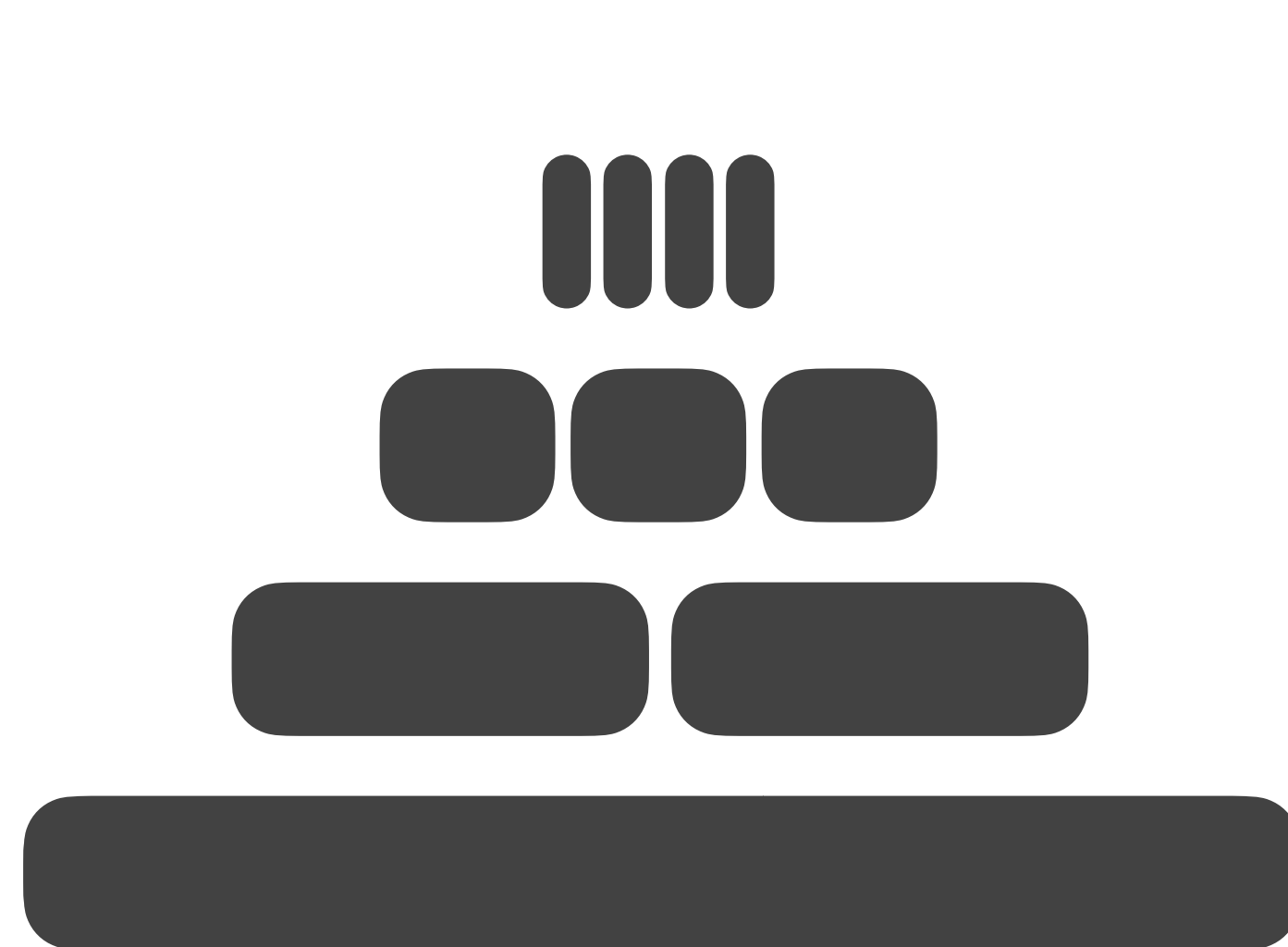
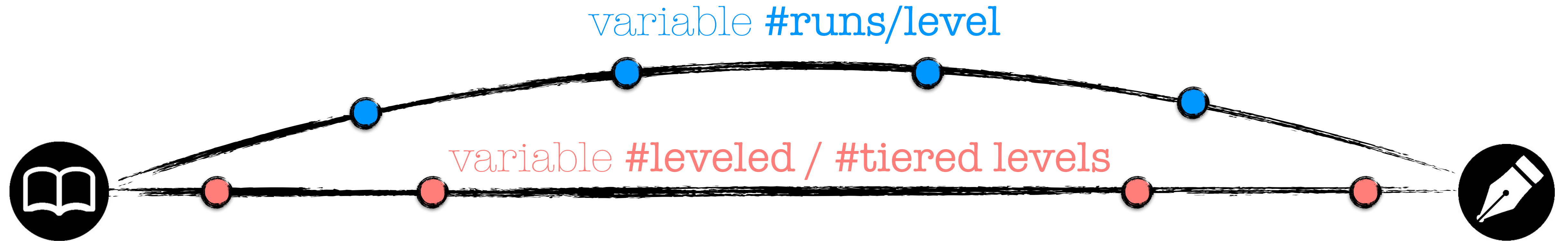
T

2

T

1

Storage Layer Design Continuum



size ratio

#runs

T

4

T

3

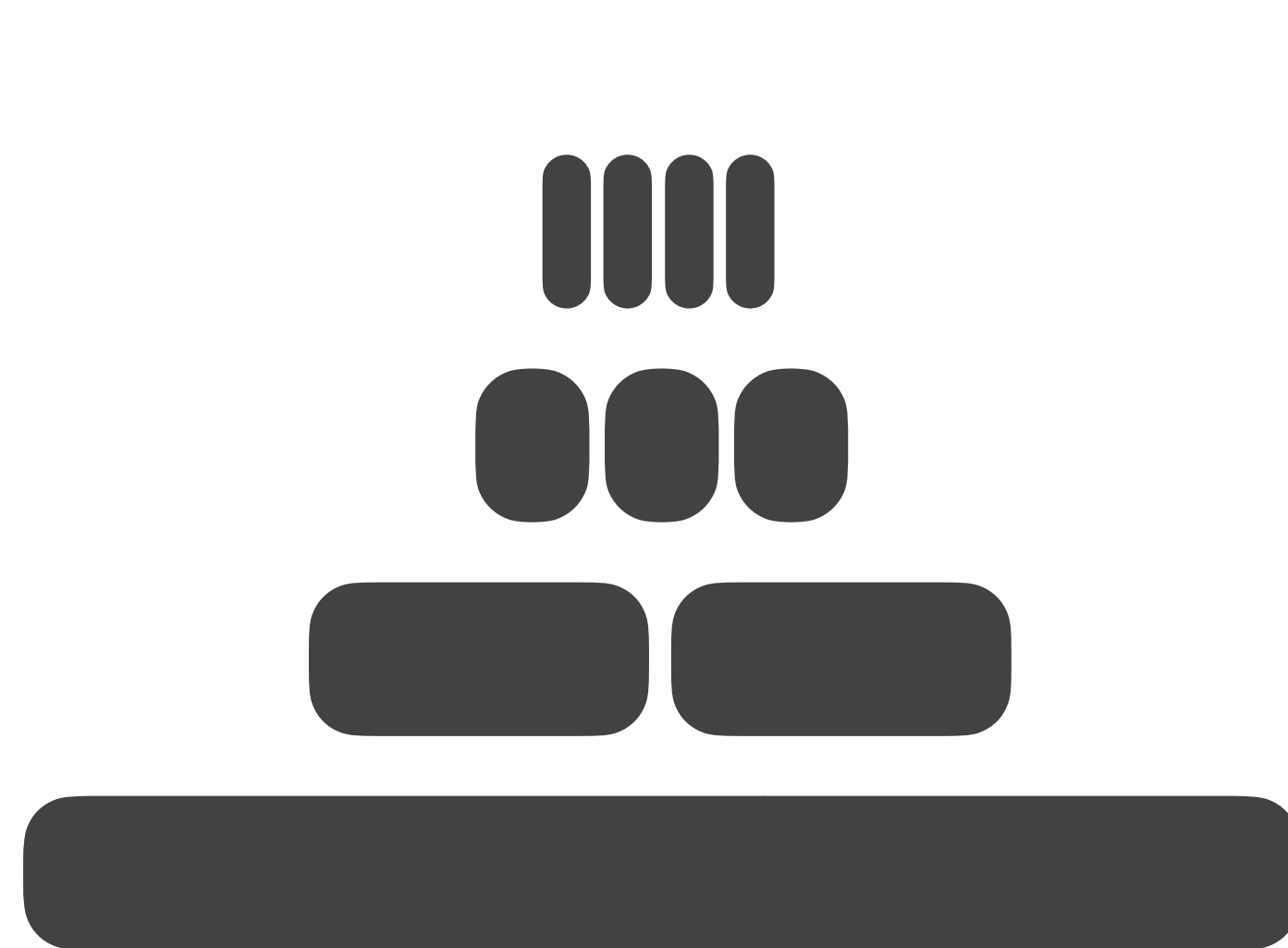
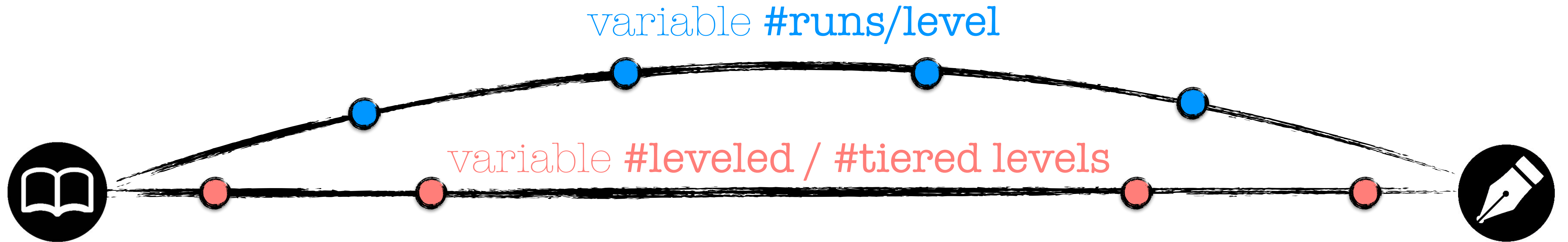
T

2

T

1

Storage Layer Design Continuum



size ratio

2

2.5

3

4

#runs

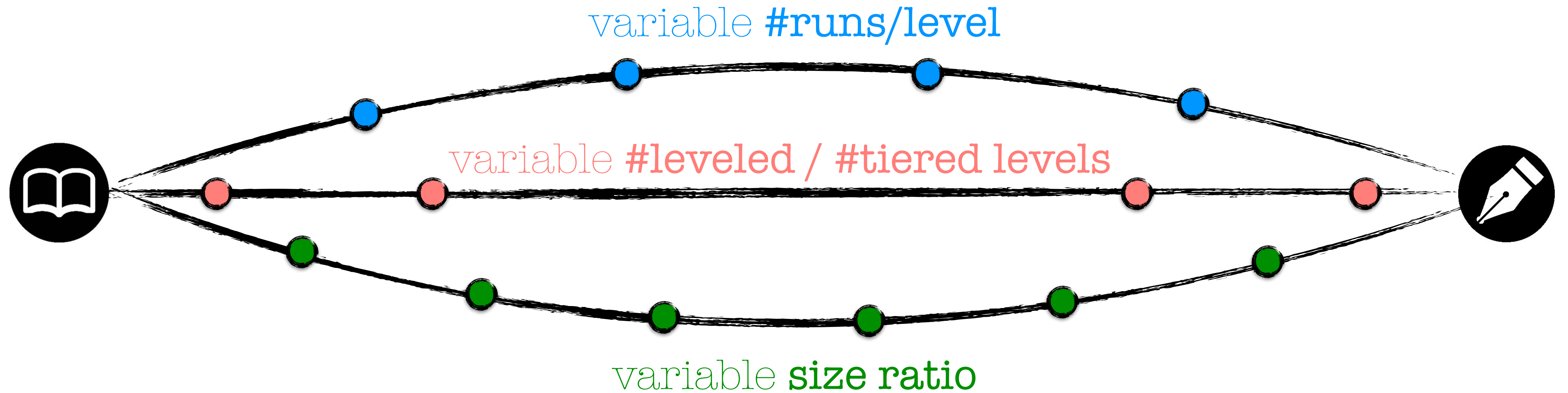
4

3

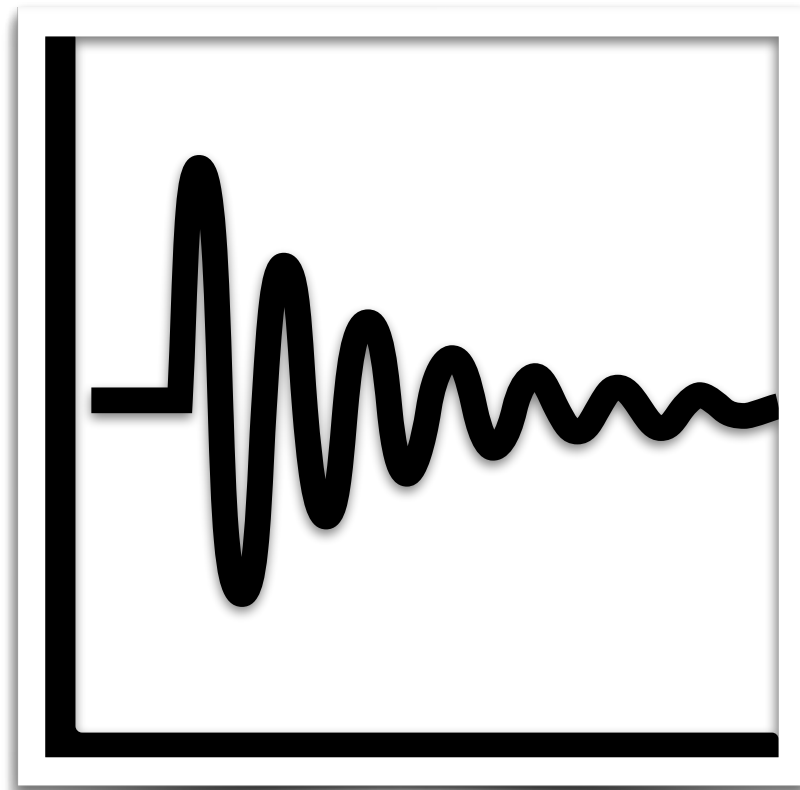
2

1

Storage Layer Design Continuum



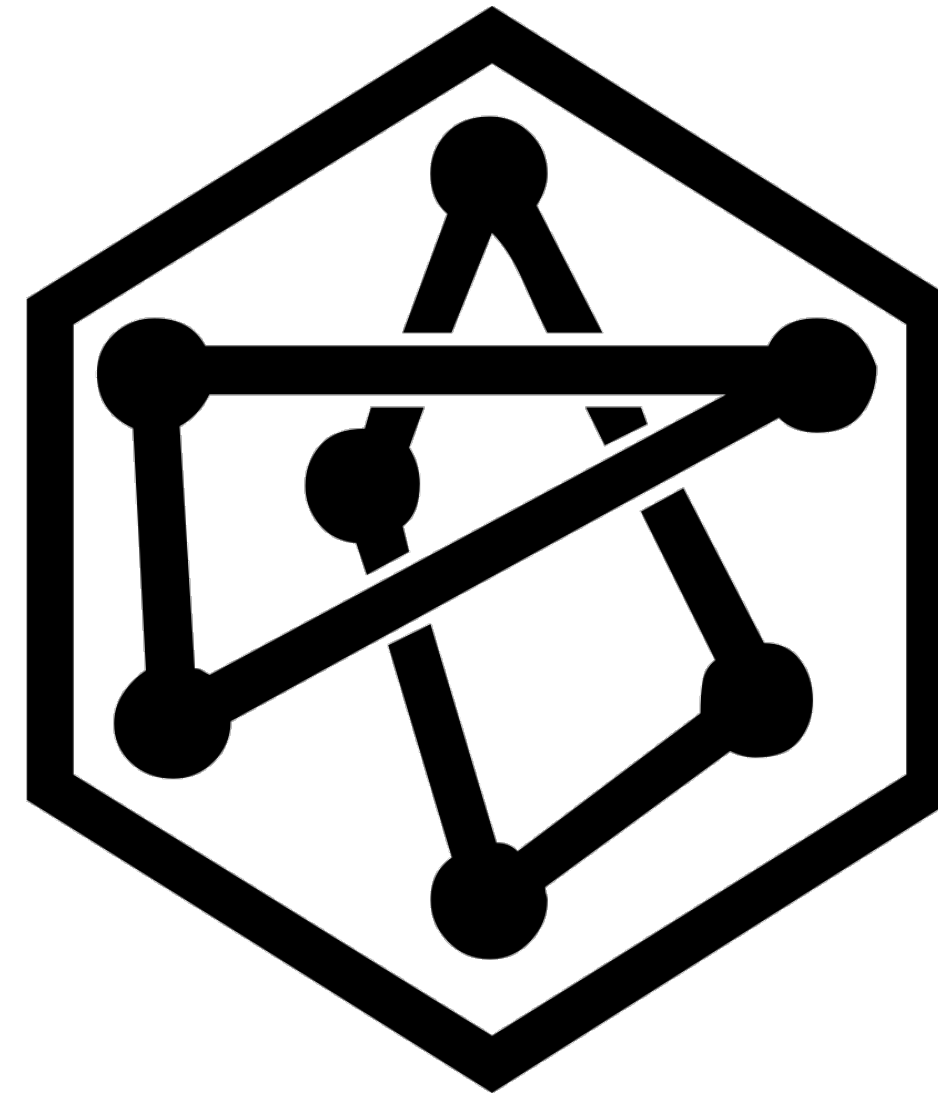
The LSM storage layer design continuum



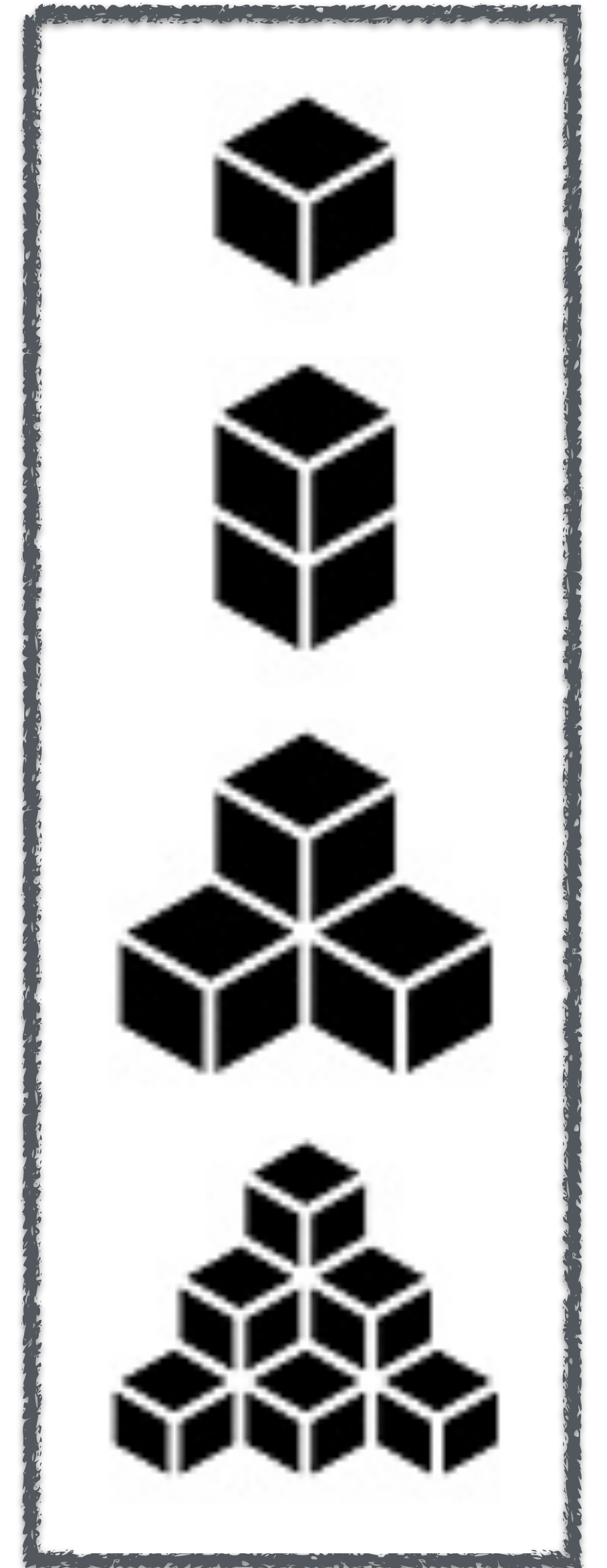
workload



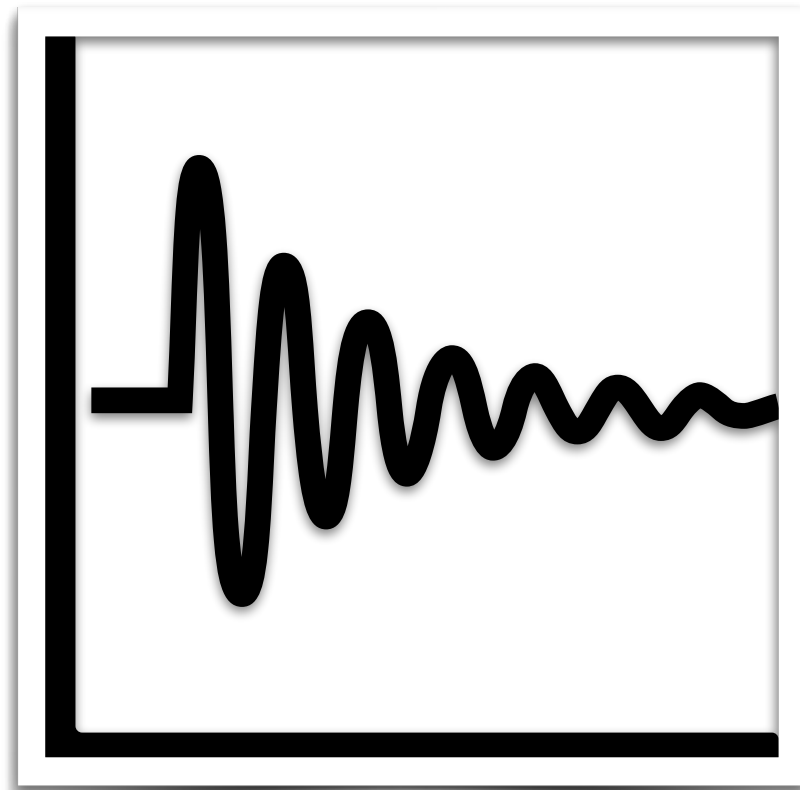
performance
target



performance
modeling



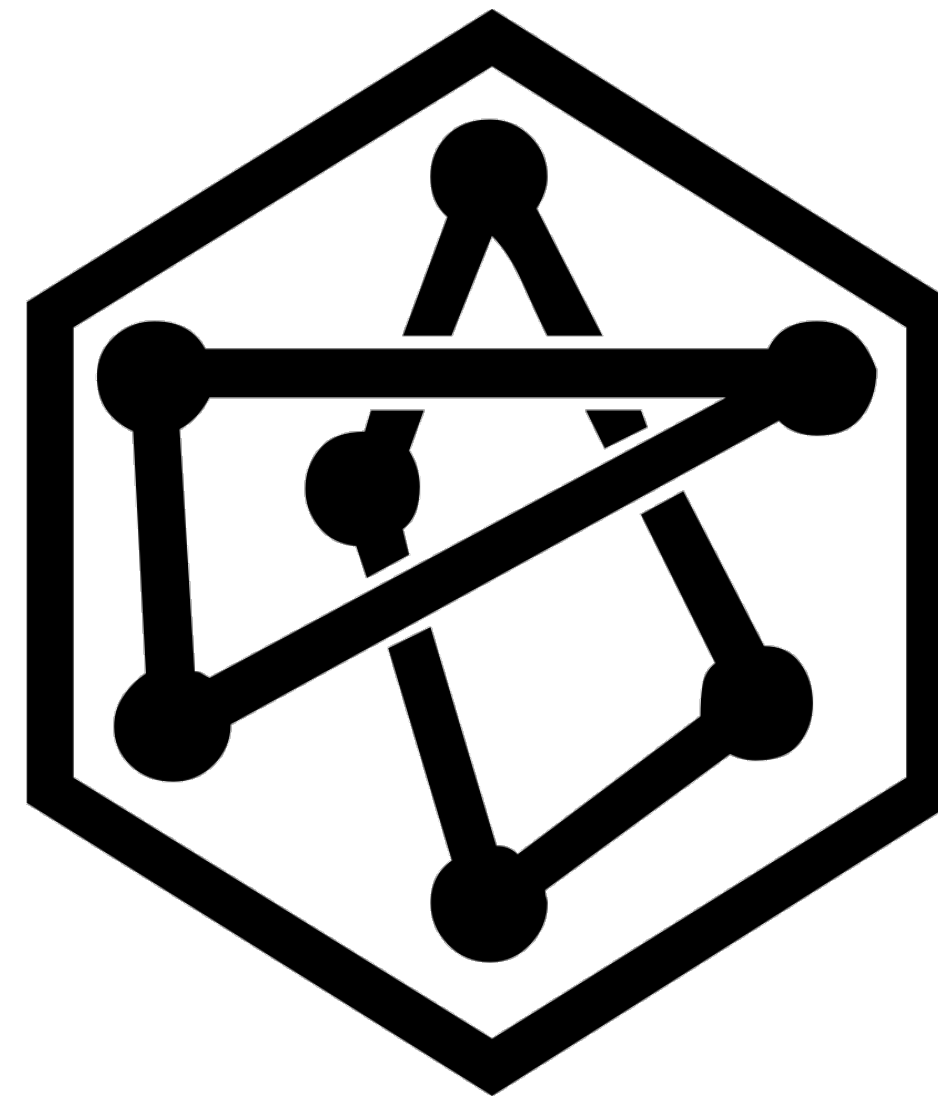
LSM designs



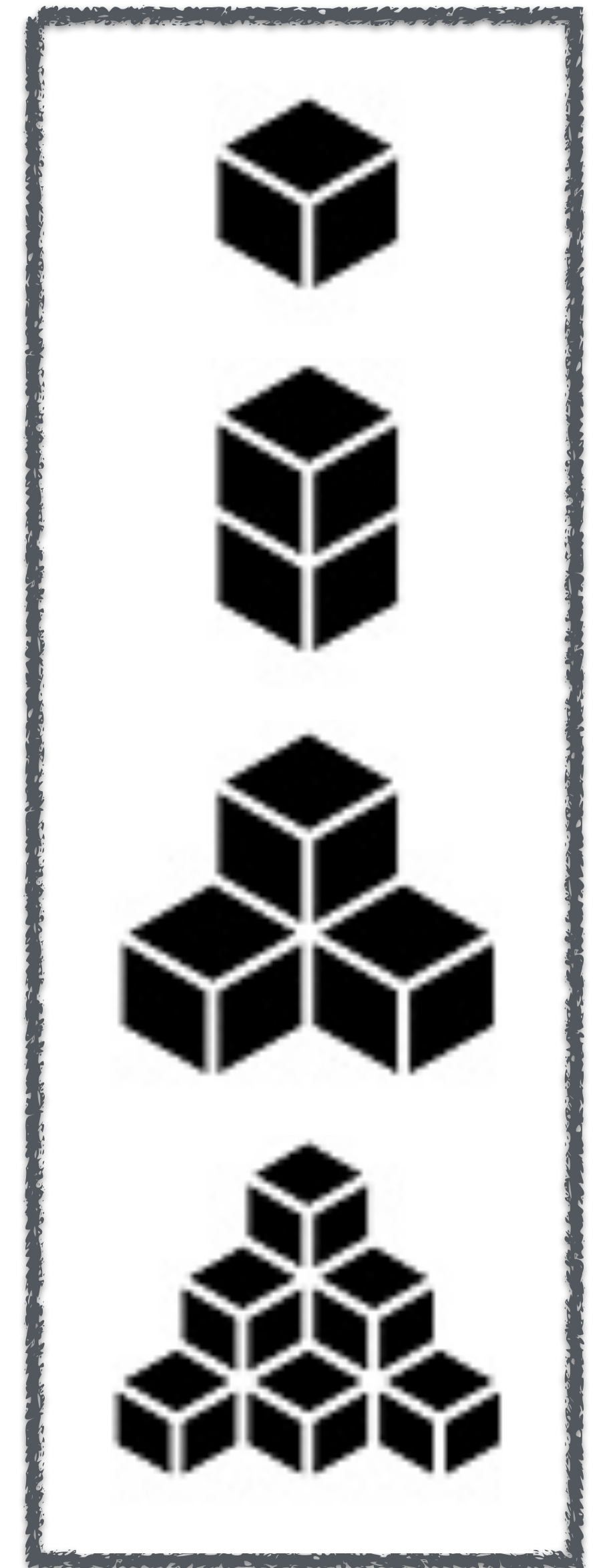
workload



performance
target



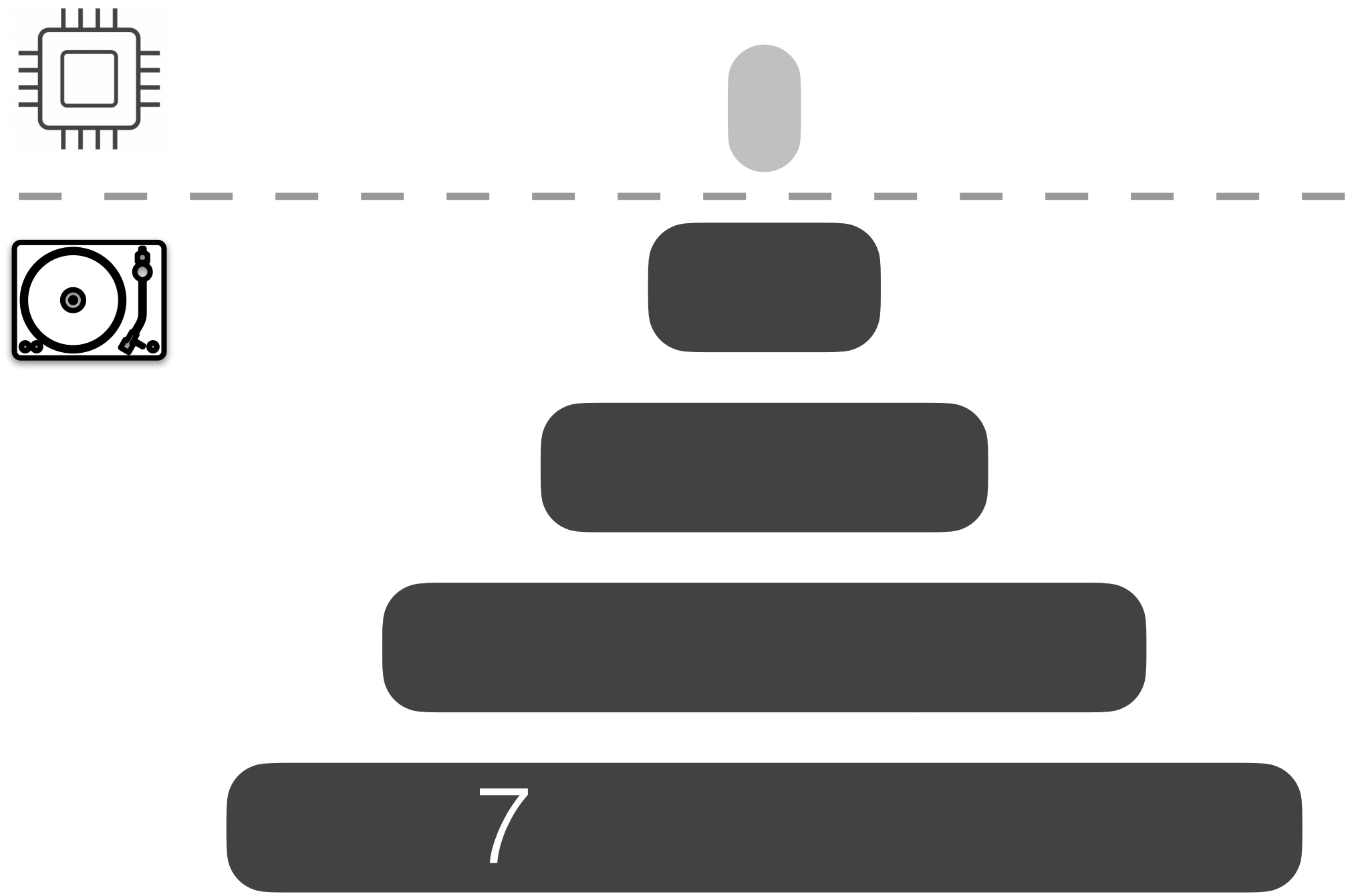
worst-case
performance
modeling



LSM designs

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

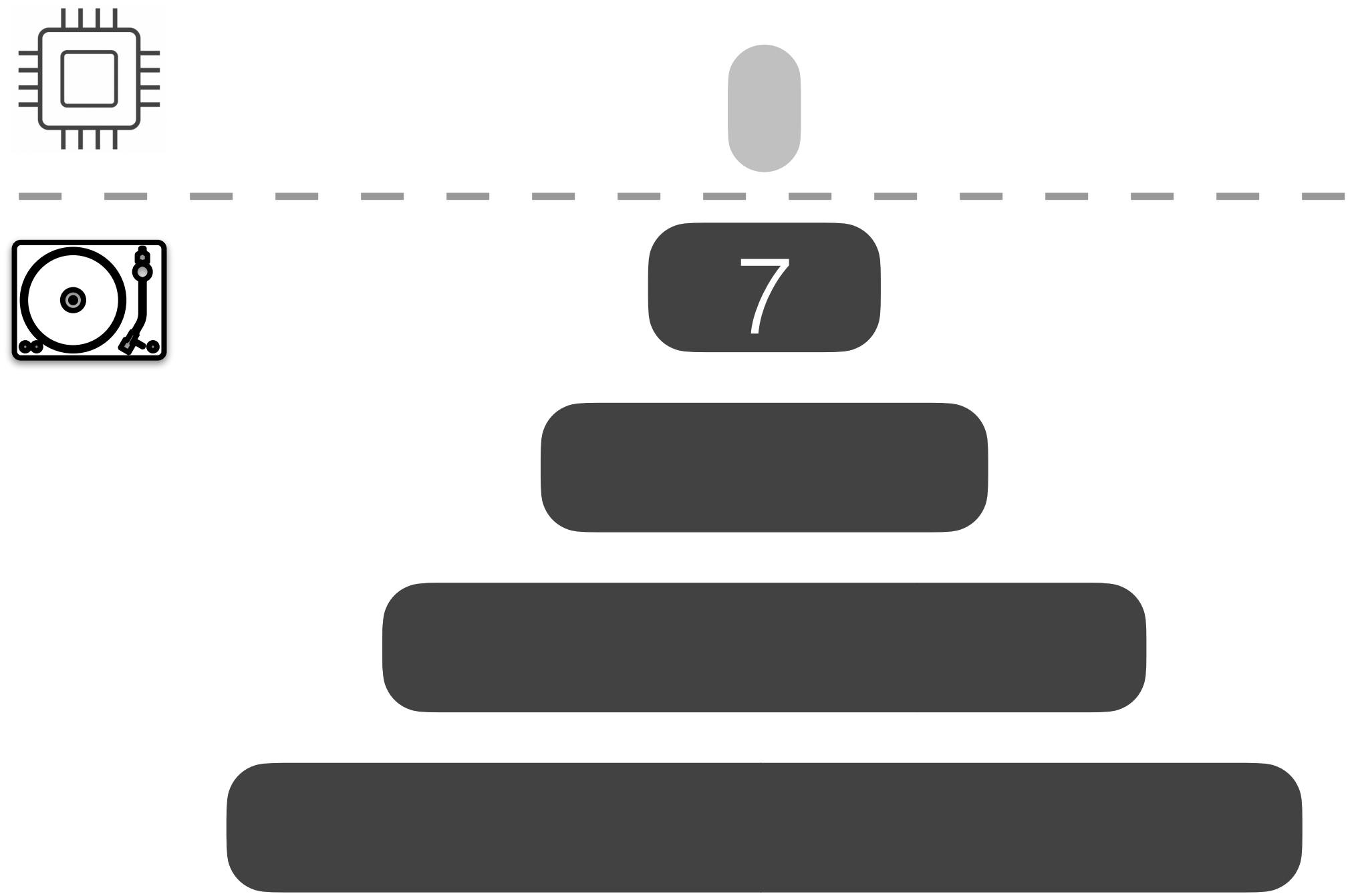
worst-case performance modeling



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

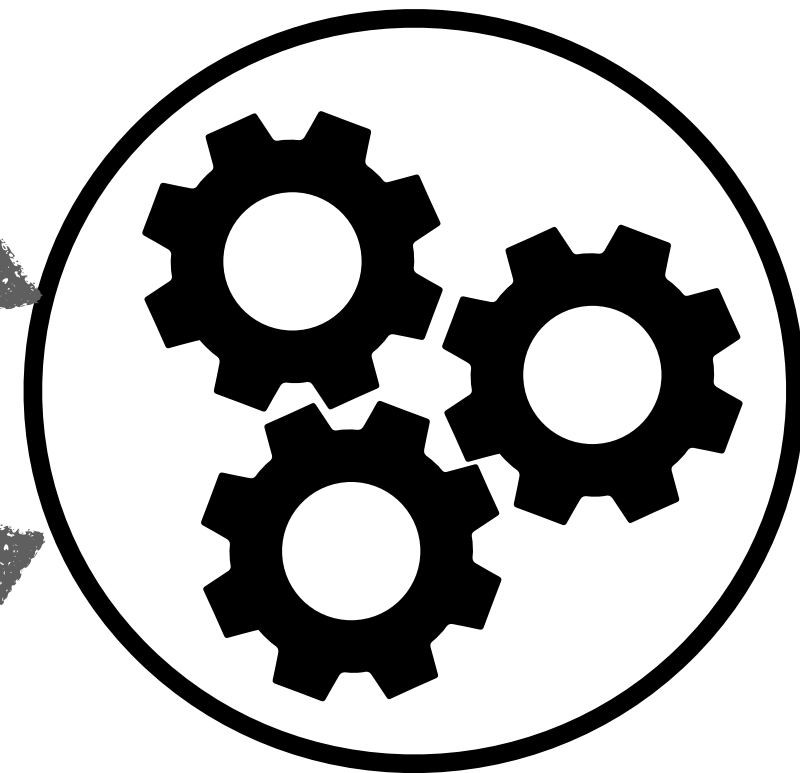
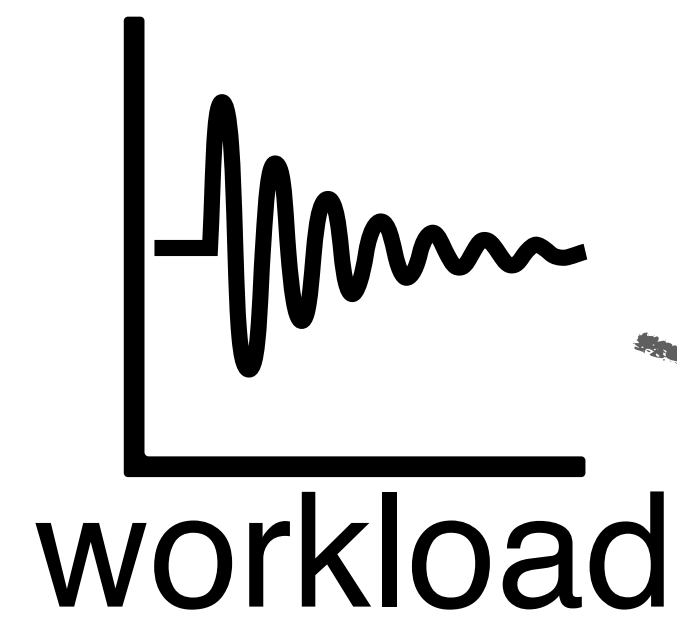


worst-case read cost: $1 + \sum_{i=1}^{L-1} \phi_i$

average-case performance modeling

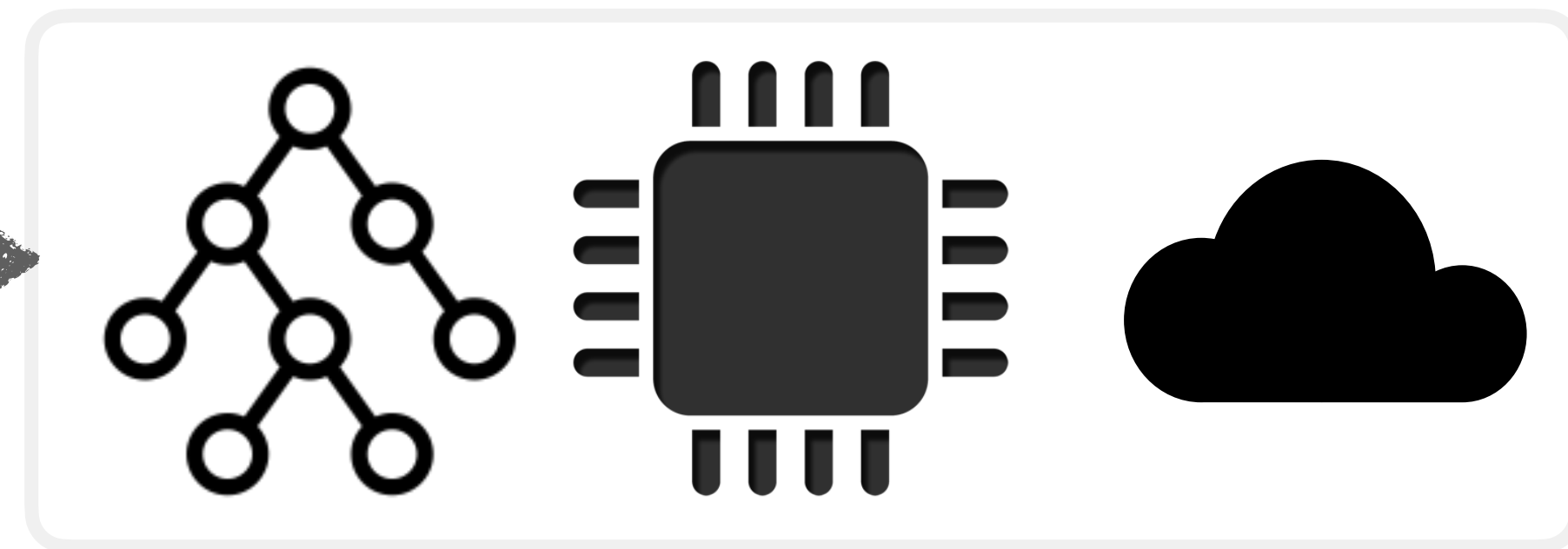
$$\sum_{i=1}^L (\mathbb{P}[\text{query in } L_i] \cdot (1 + \sum_{j=1}^{i-1} \phi_j))$$

average-case
performance
modeling



Cosine

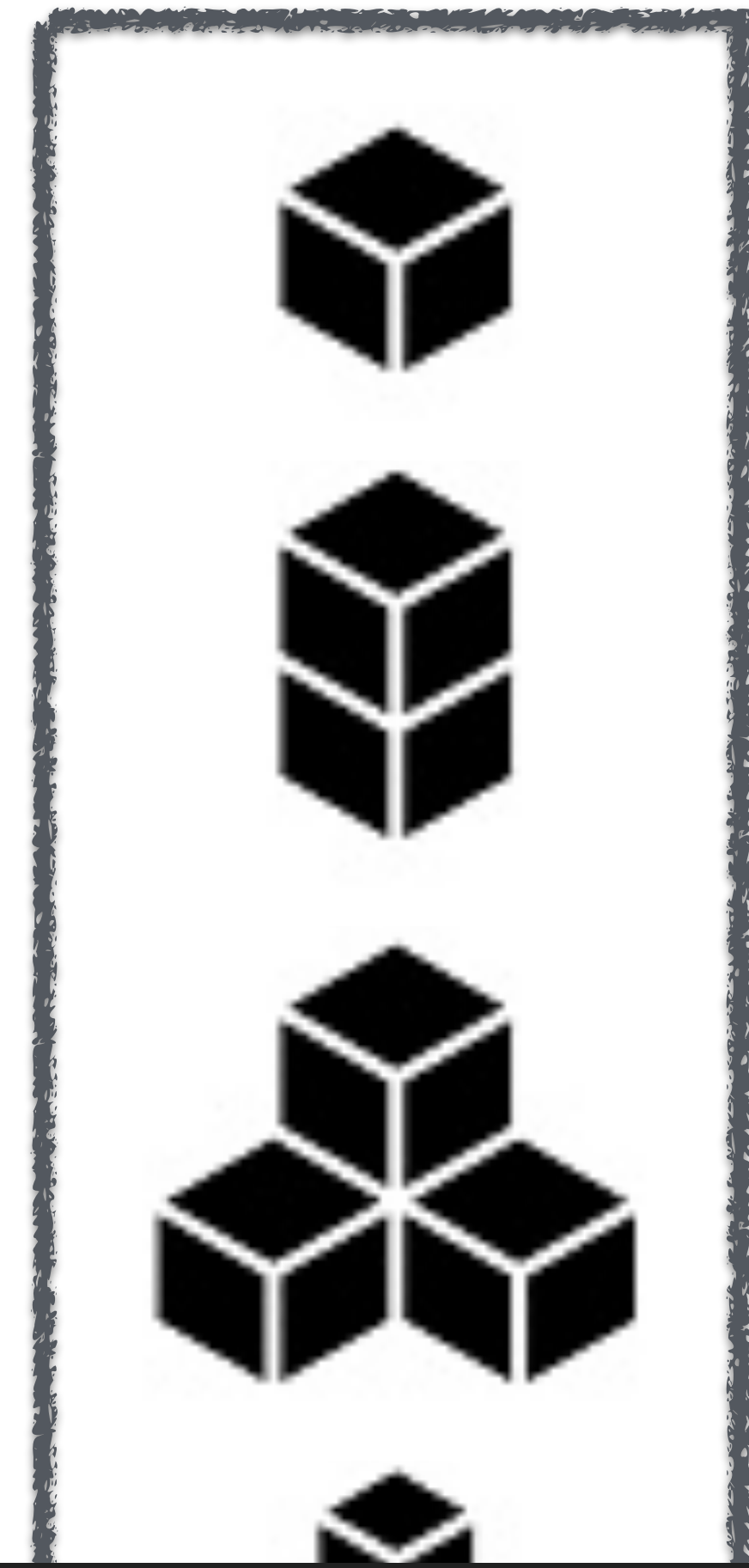
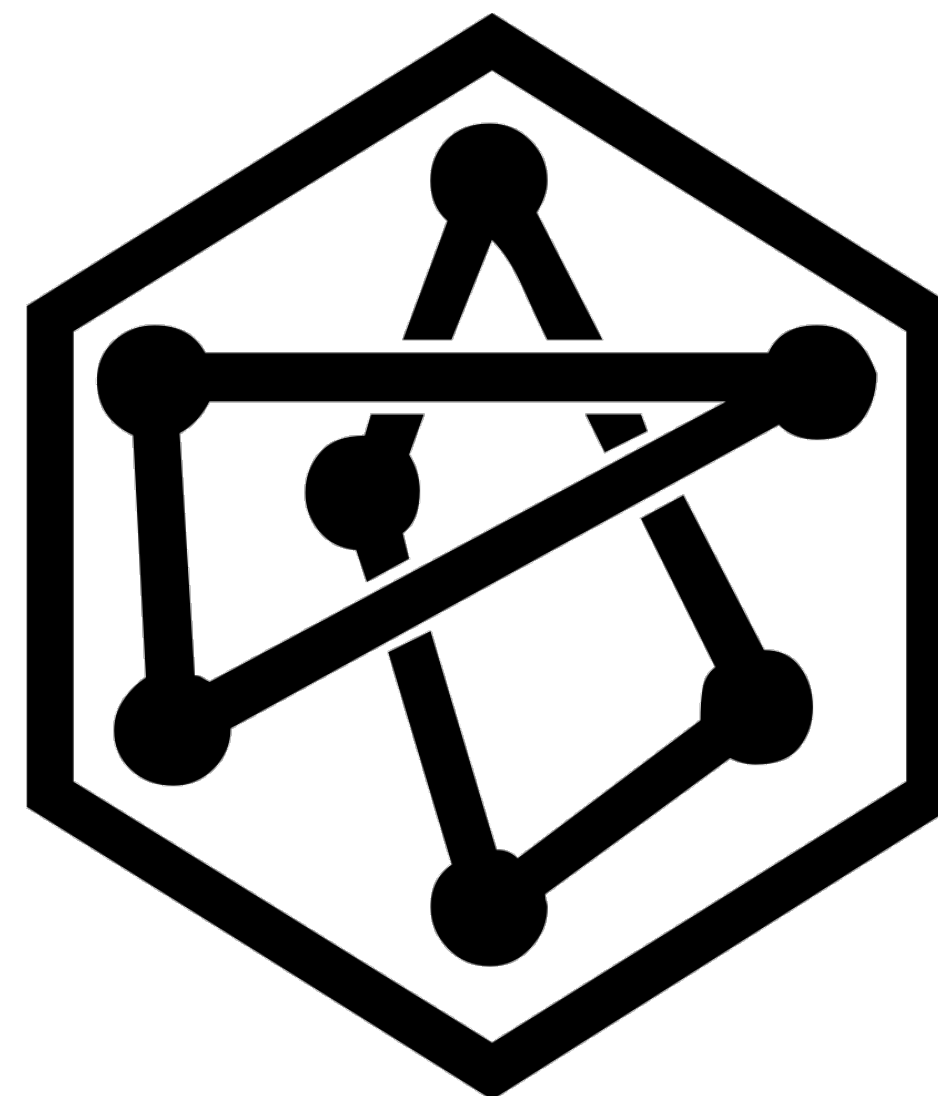
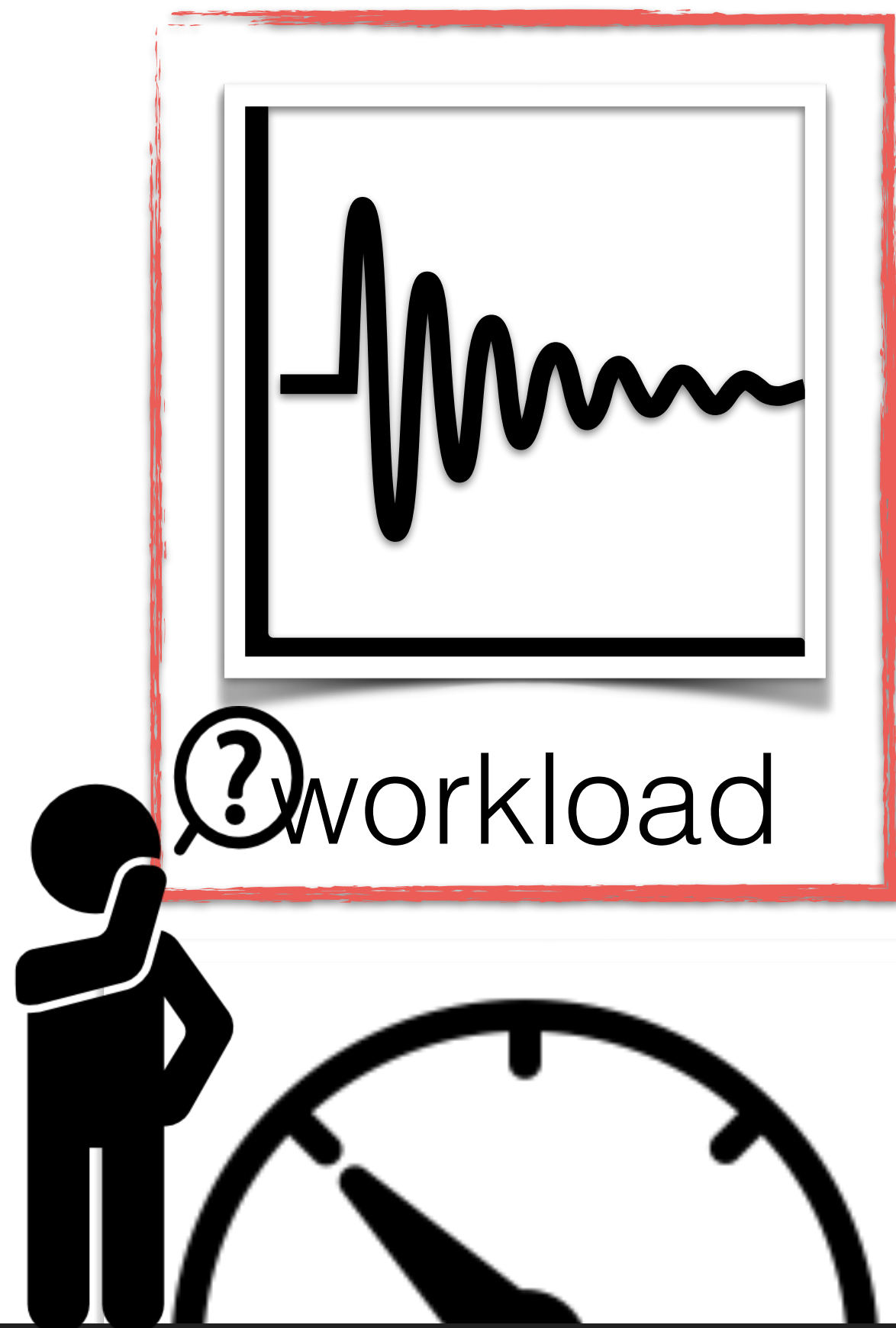
optimal configuration



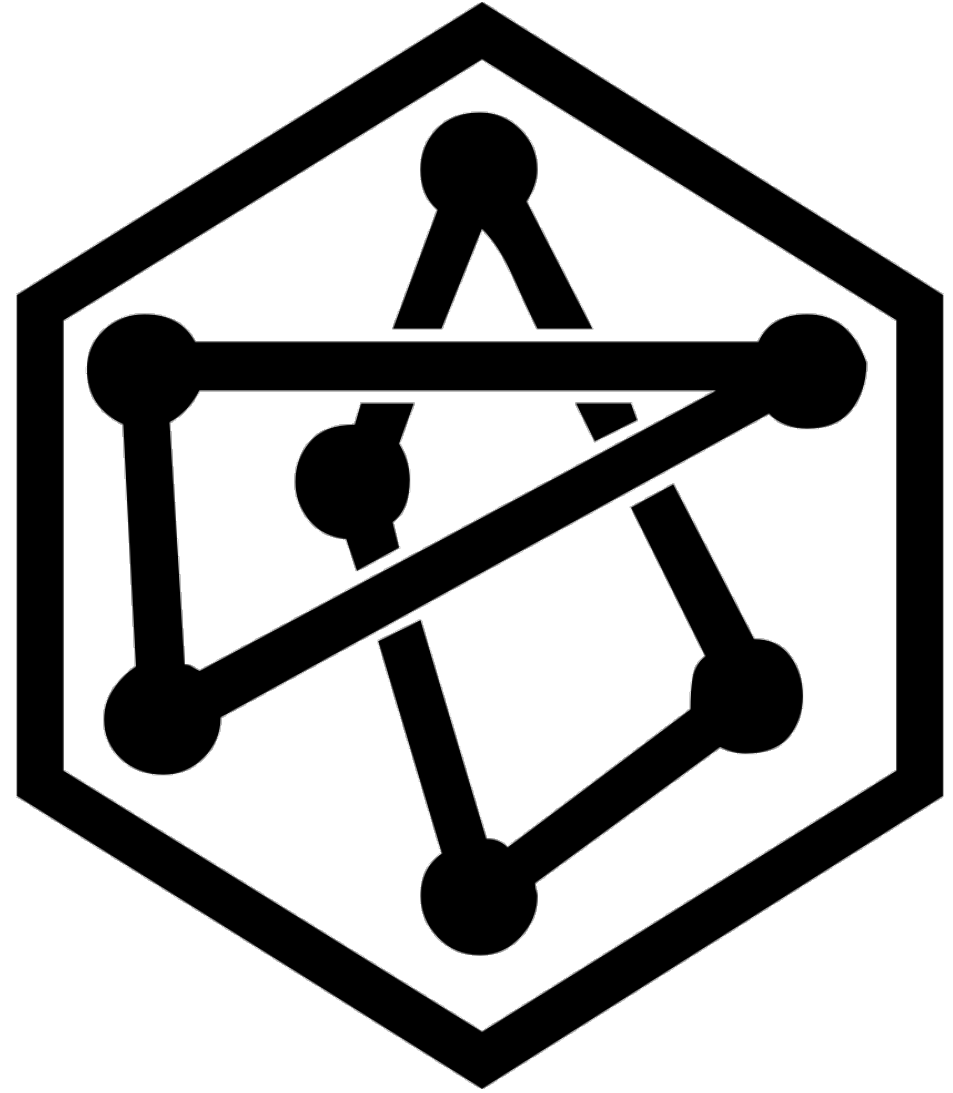
SE design

h/w

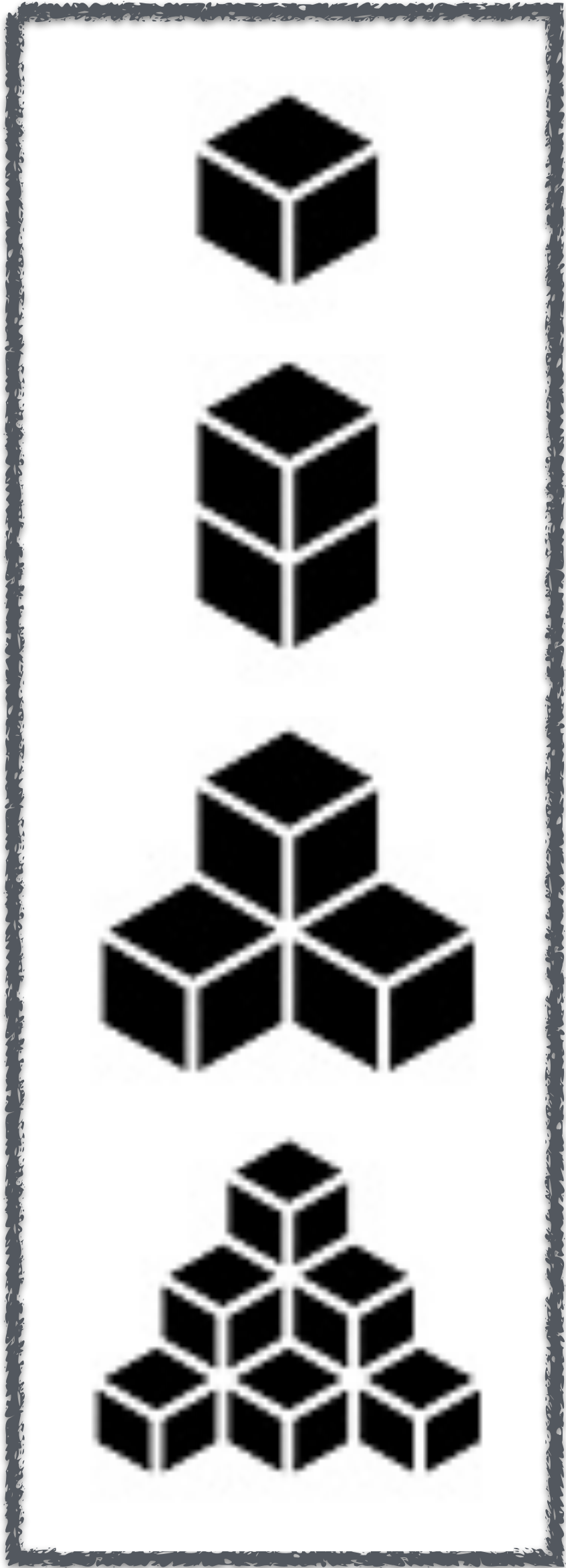
cloud
provider



What if the workload changes?



average-case
performance modeling



LSM designs

Tuning

workload
hardware
design



optimal
tuning

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

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

Robust Tuning

unpredictability in

workload
hardware
design



optimal
tuning

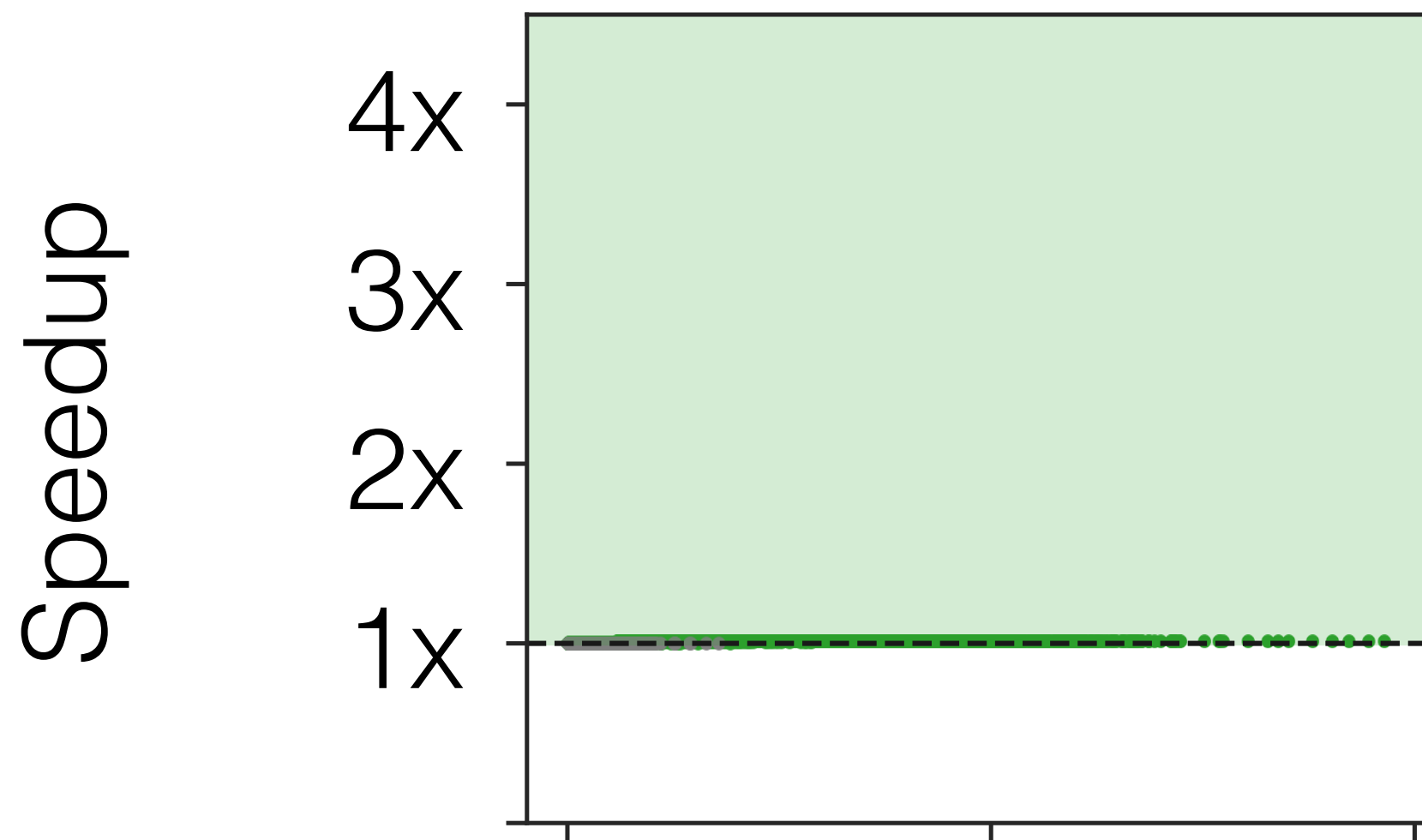
given workload \mathcal{W} , dataset \mathcal{D} , and a region $\mathcal{R}(\mathcal{W})$, find a tuning \mathcal{T}^{rob} :

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

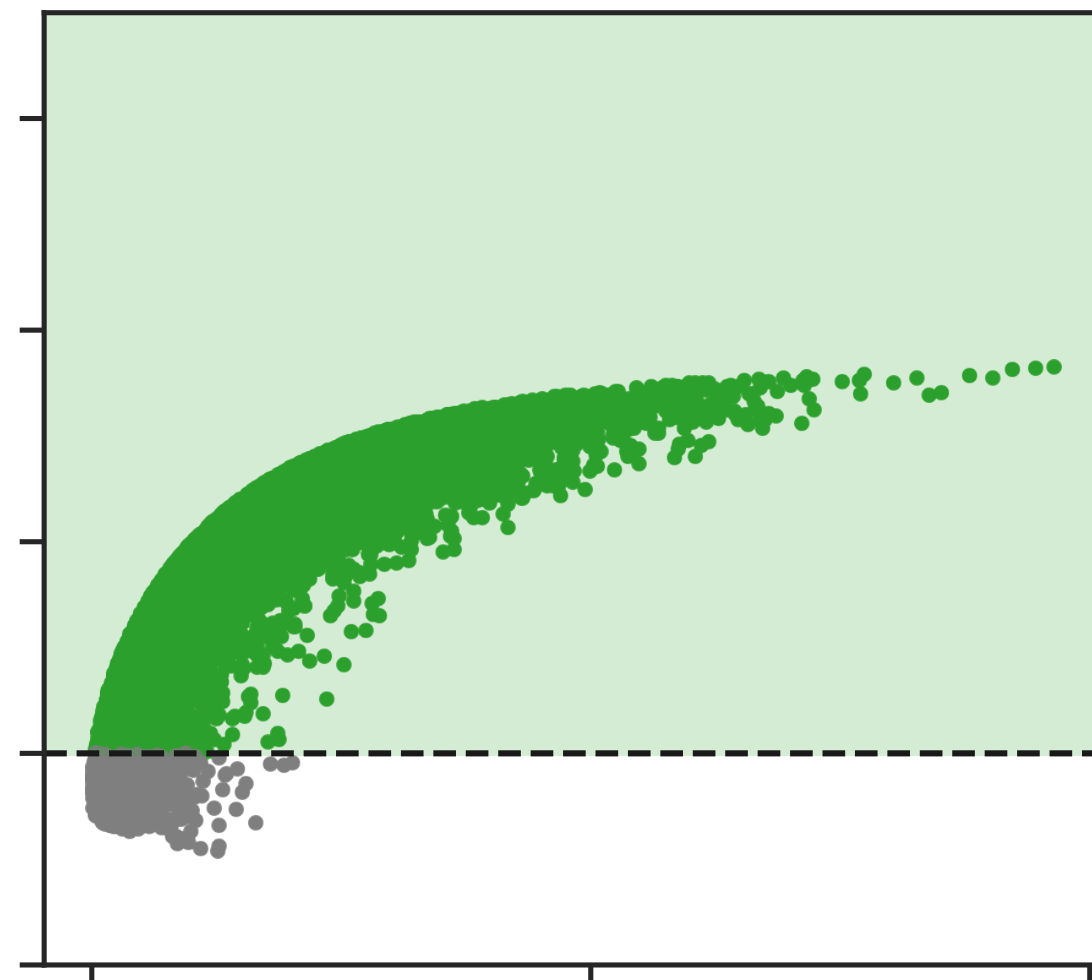


Robust Tuning

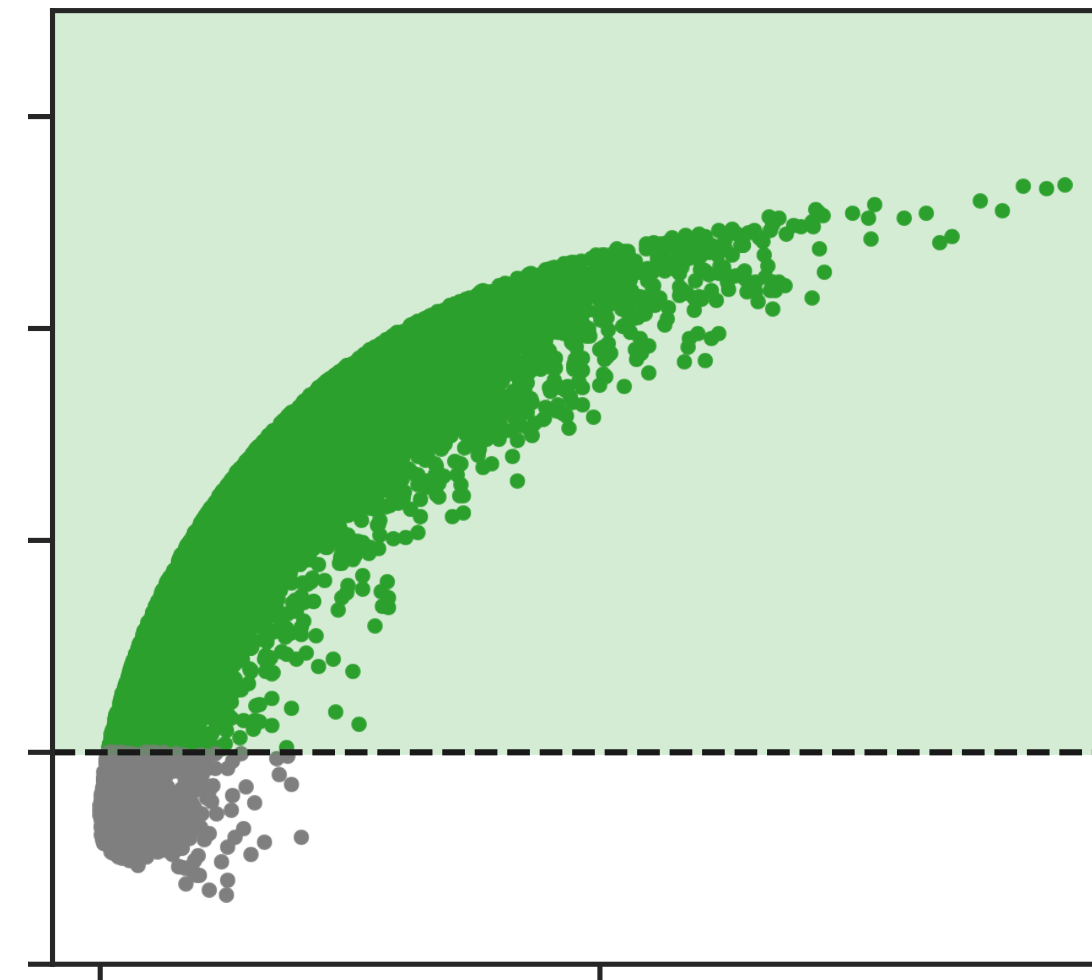
Nominal



Robust



Robust (for more noise)



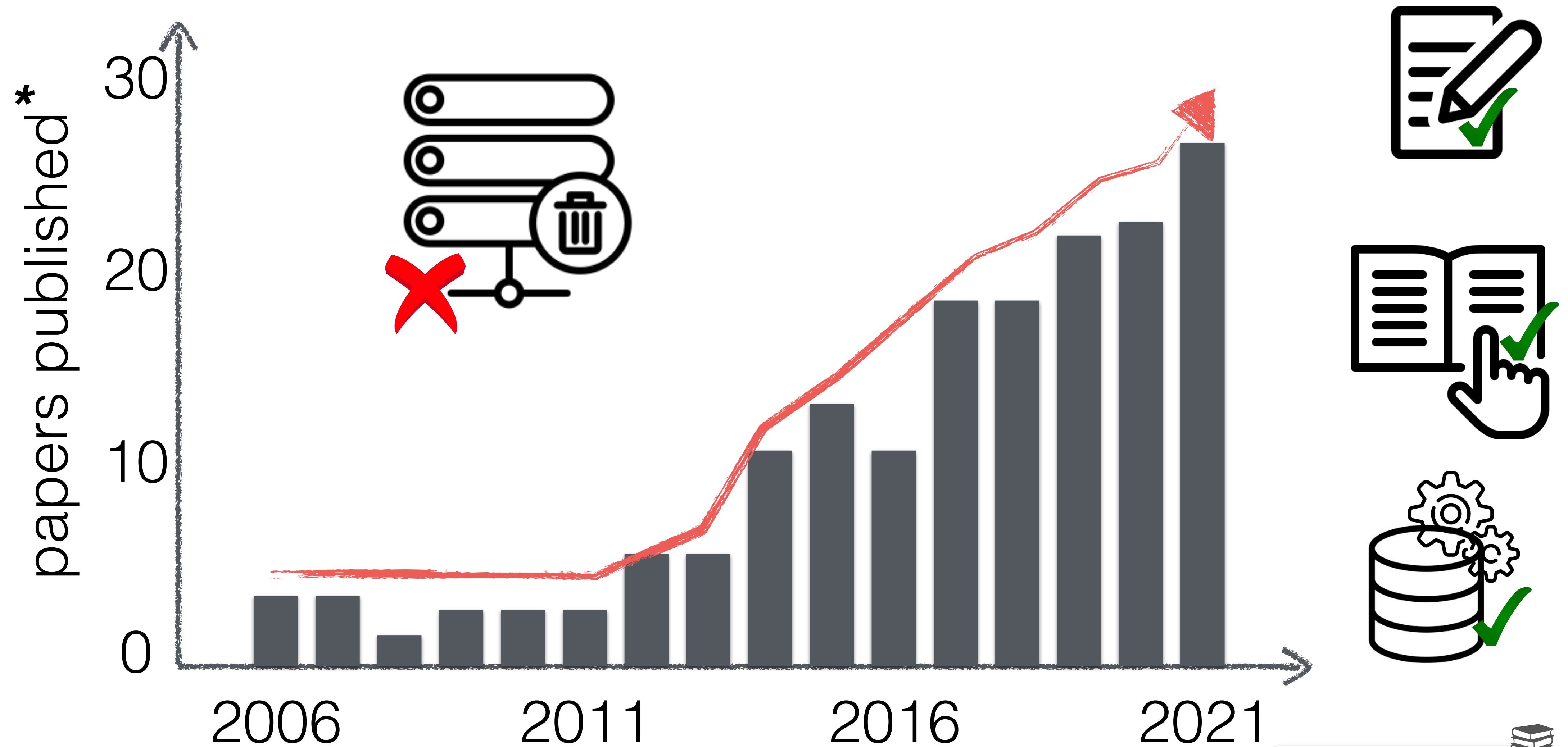
Speedup



Workload Noise

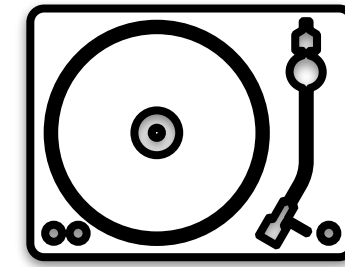
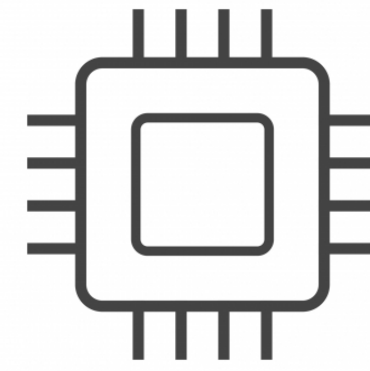
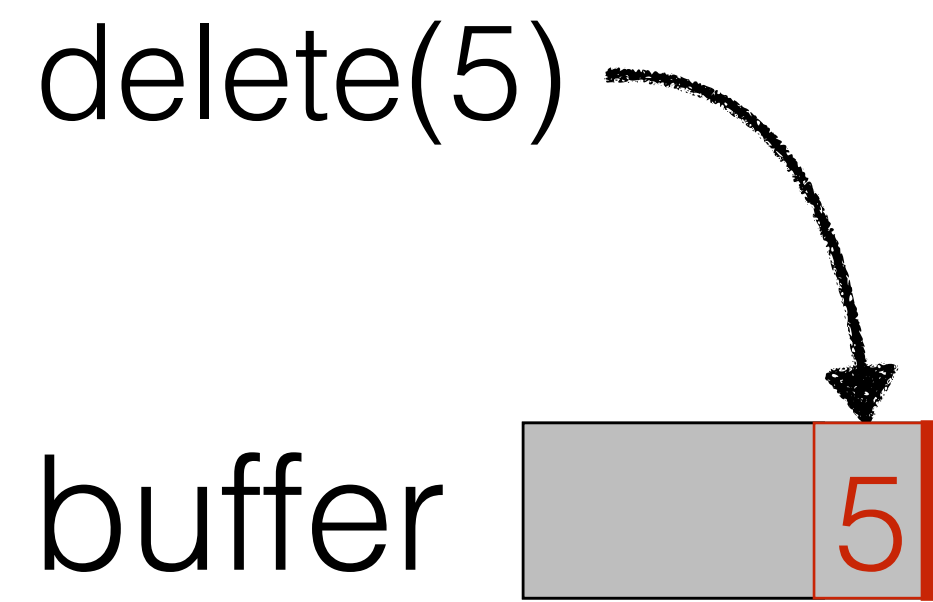


Research Trend: **What about Deletes?**



* data from DBLP

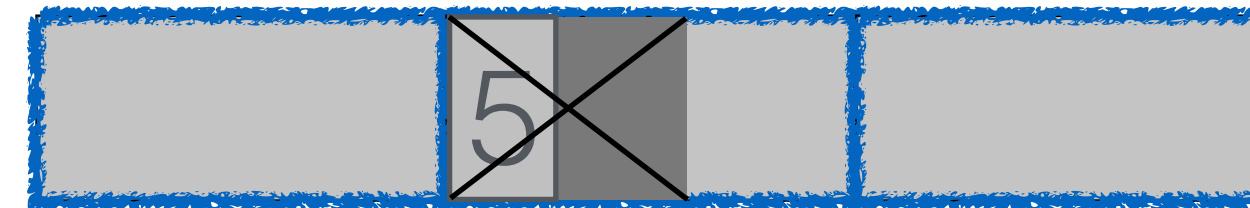
deletes in LSM-tree



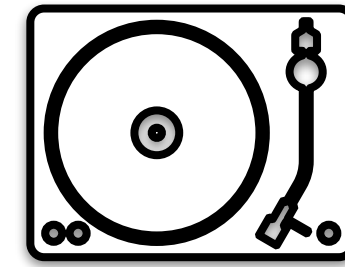
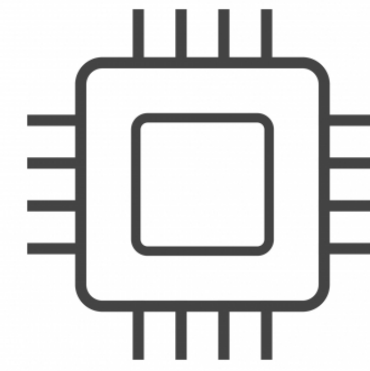
L1

L2

L3



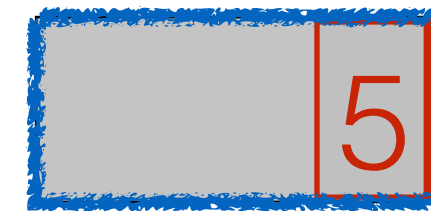
deletes in LSM-tree



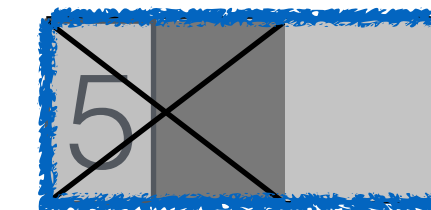
buffer



L1



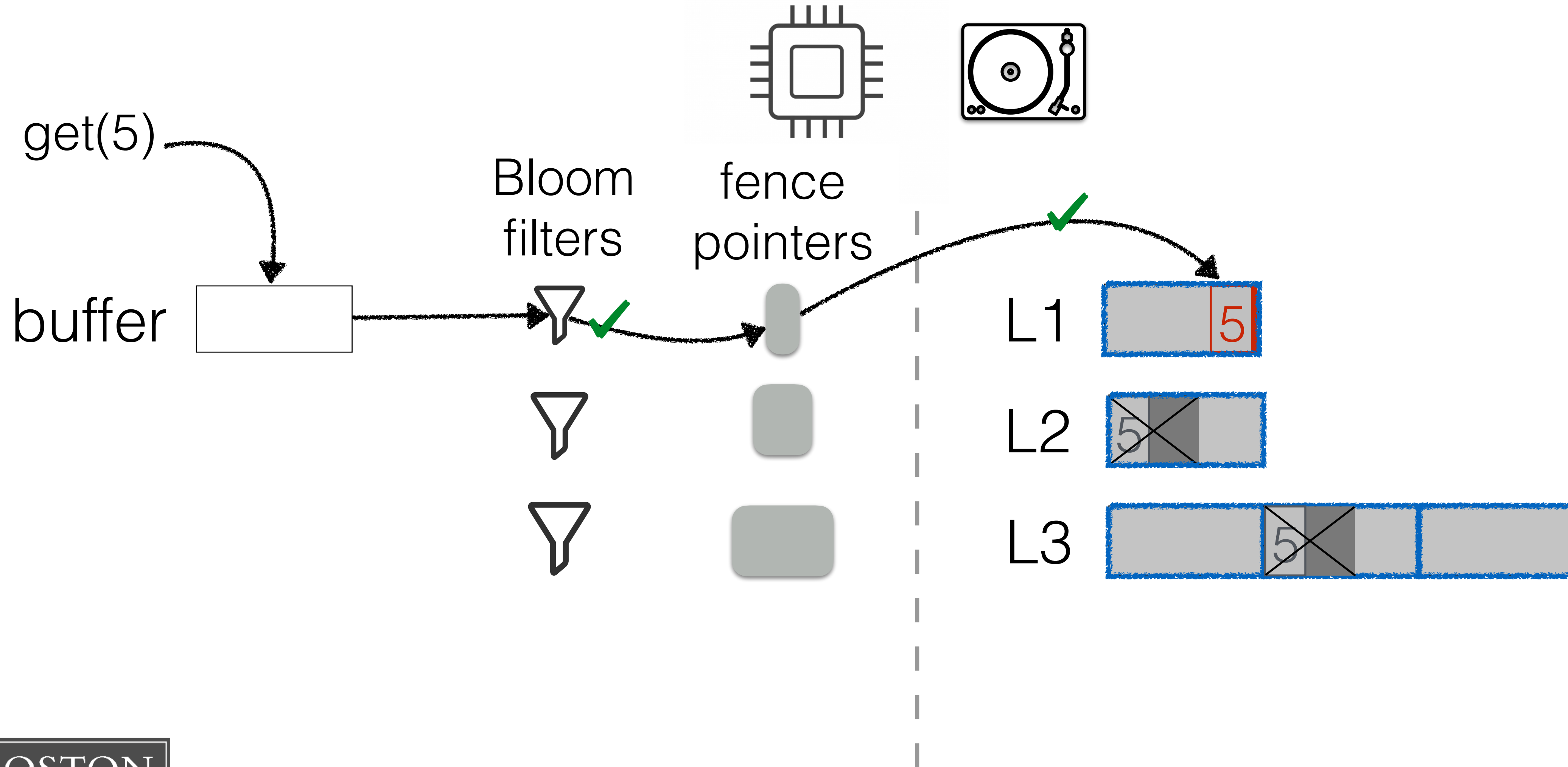
L2



L3



deletes in LSM-tree



Large-scale production

Internal DB ops

Privacy

ZippyDB 

25.2M/day

table drop 

data 

migration




GDPR 
(EU, UK)




CCPA 
(California)

UP2X 

**92.5% merge
through deletes**

periodic 
cleanup



VCDPA 
(Virginia)

Facilitating **Timely Deletes**

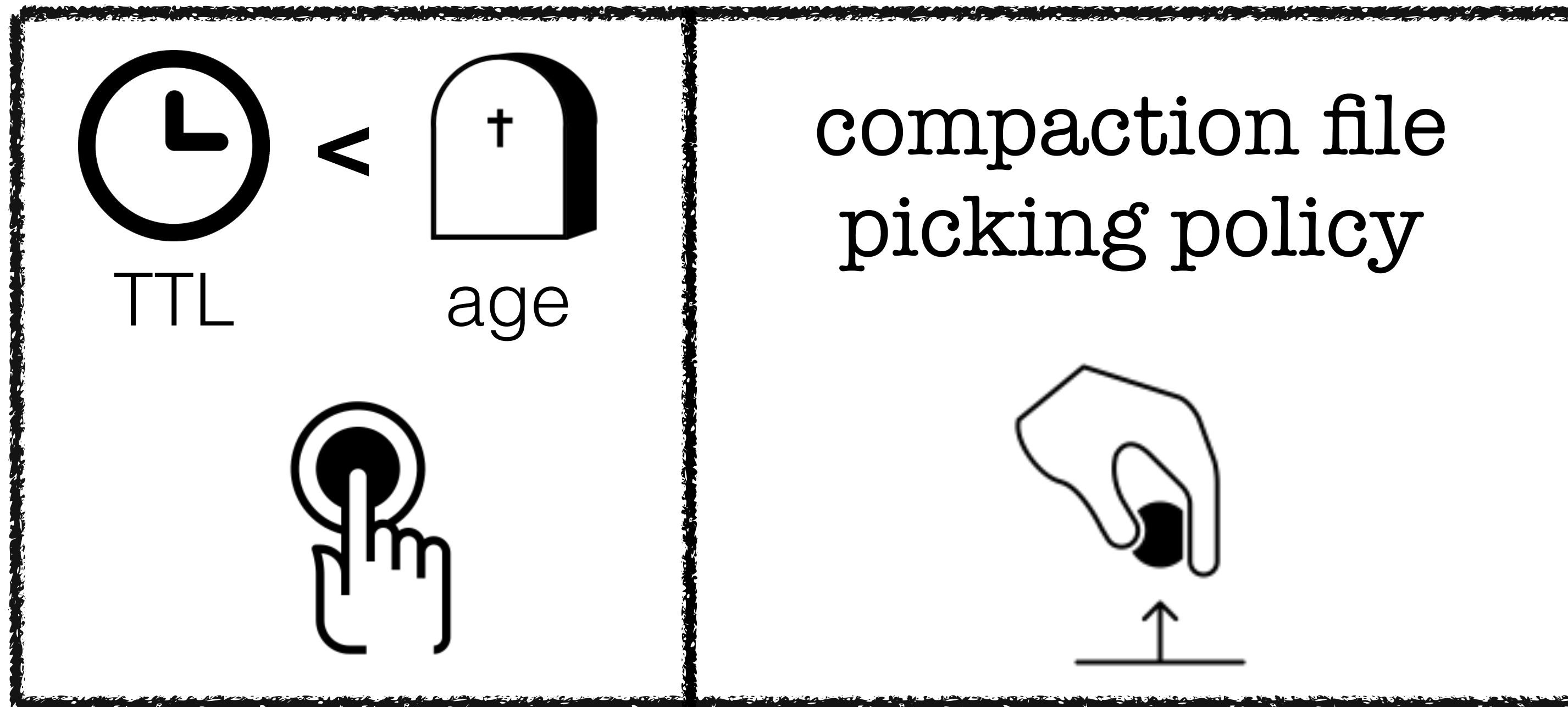
compaction
trigger



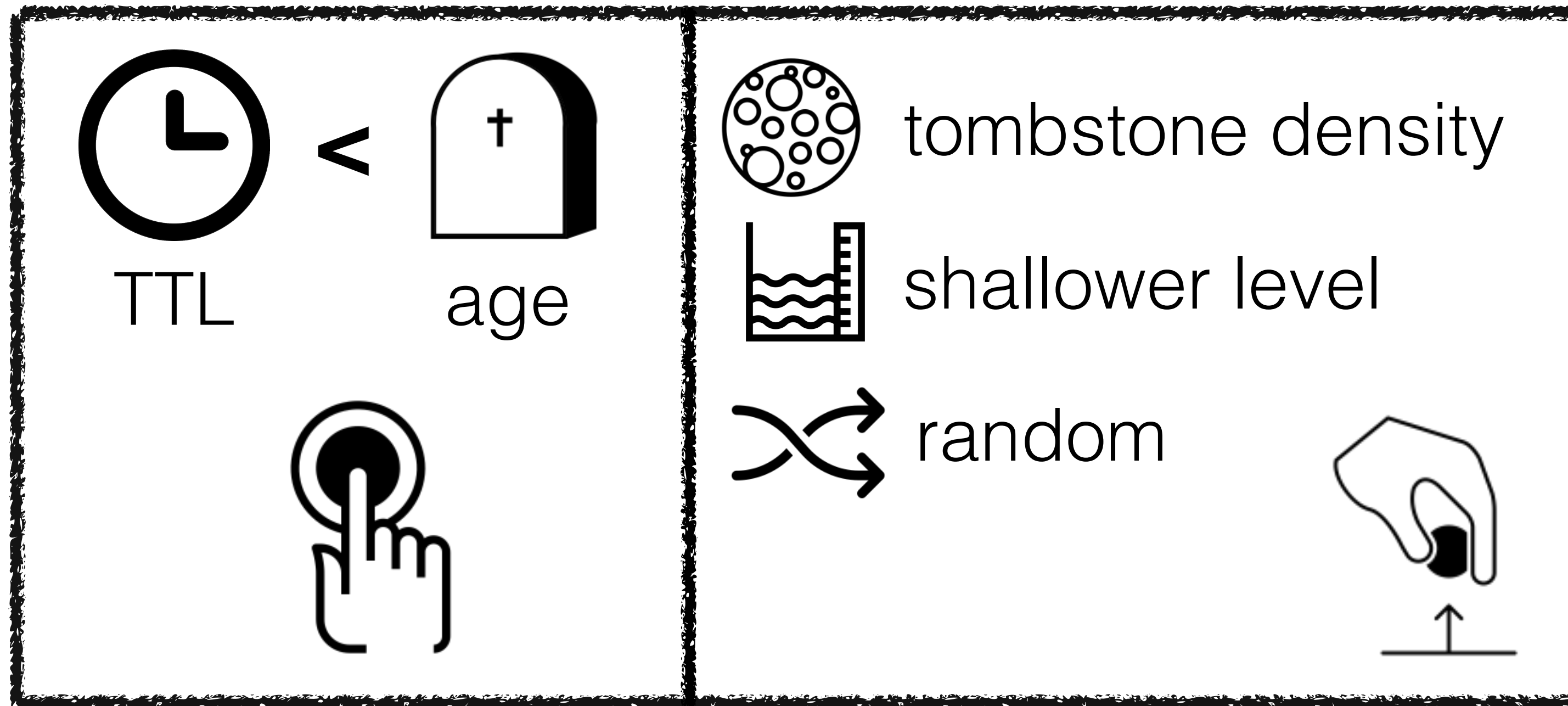
compaction file
picking policy



Facilitating **Timely Deletes**

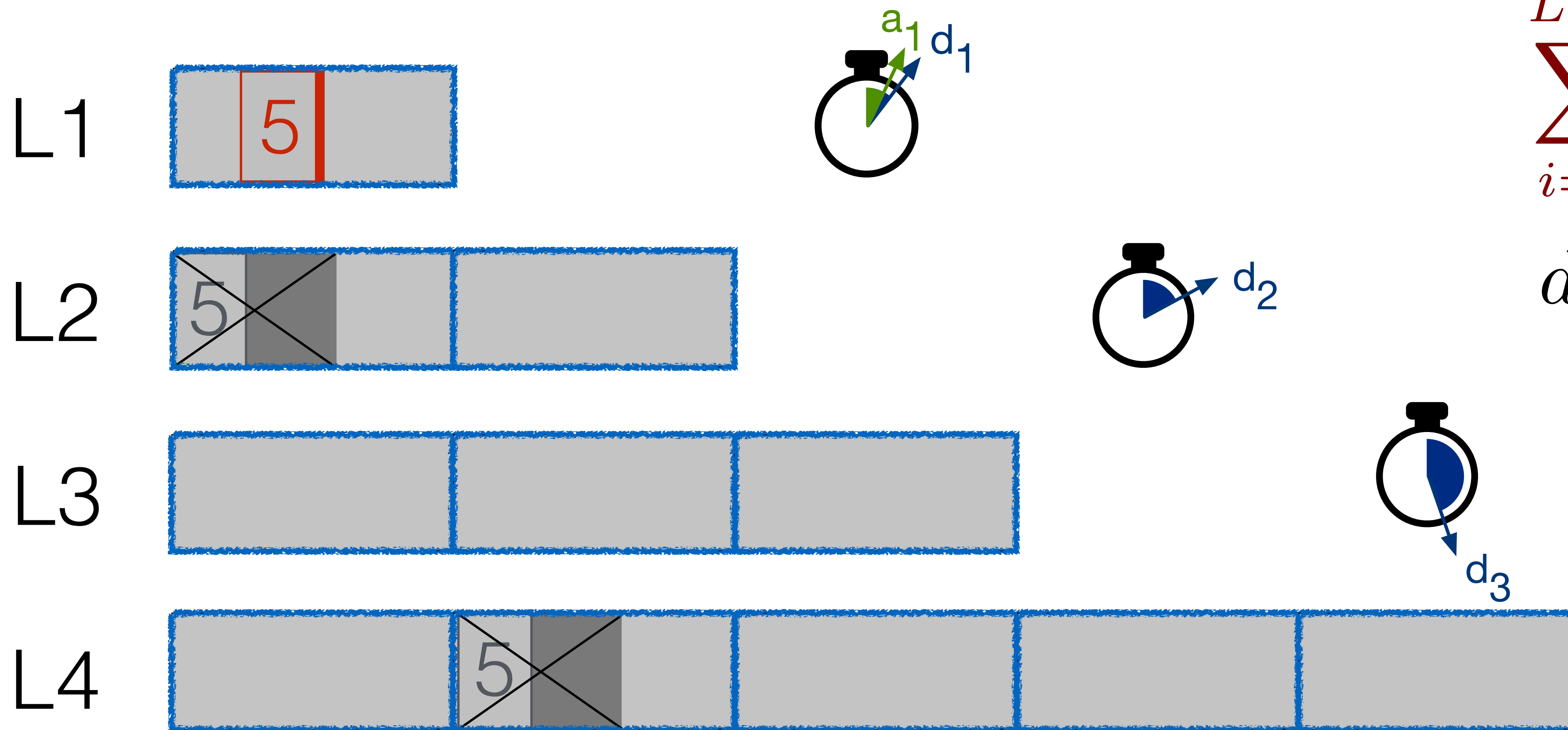


Facilitating **Timely Deletes**



Facilitating **Timely Deletes**

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



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

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

Policy layer

Right to be forgotten

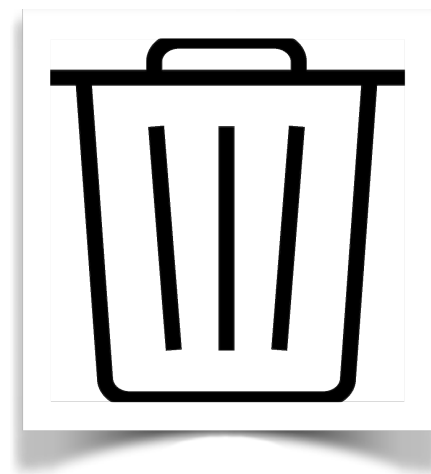
Right to delete

Deletion right

Requirements layer



Retention-based Deletes



On-demand Deletes

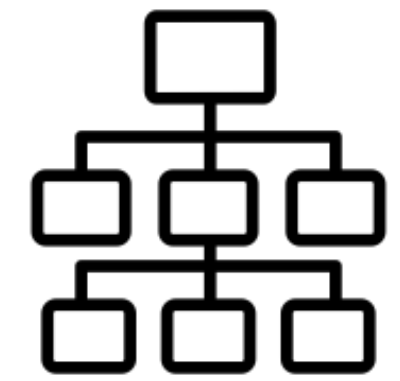
Application layer

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

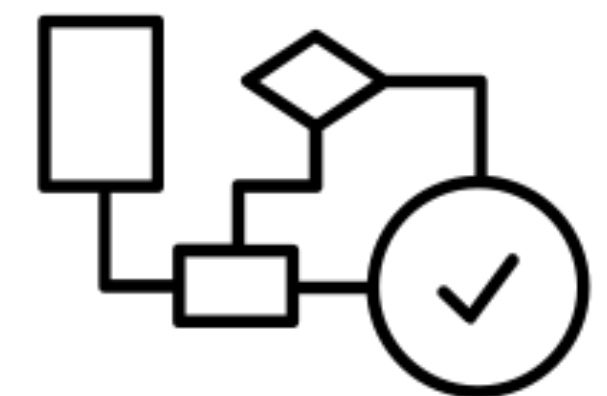
```
INSERT INTO R (...)
  WITH RET_DUR {<t>|t<i>};
```

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

System layer



Data layout re-organization



Data deletion algorithms

Policy layer

Right to be forgotten

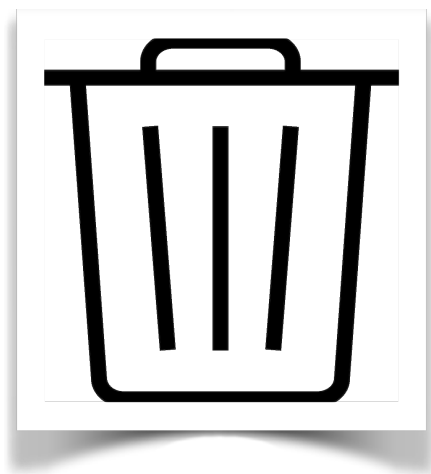
Right to delete

Deletion right

Requirements layer



Retention-based Deletes



On-demand Deletes

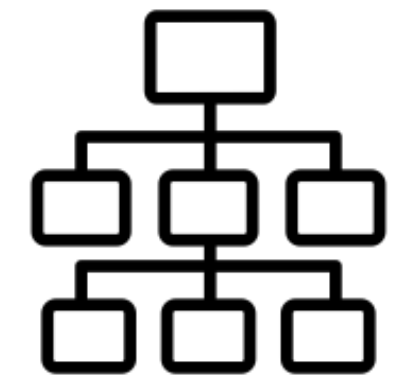
Application layer

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

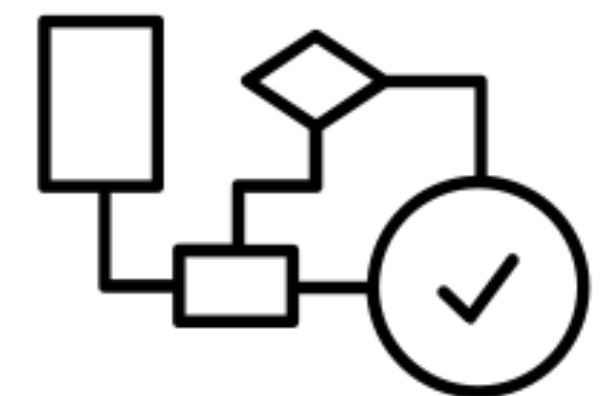
```
INSERT INTO R (...)
  WITH RET_DUR {<t>|t<i>};
```

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

System layer



Data layout re-organization



Data deletion algorithms

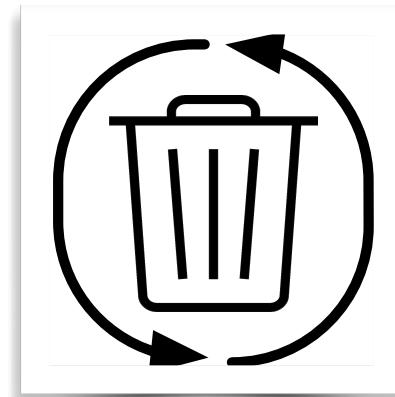
Policy layer

Right to be forgotten

Right to delete

Deletion right

Requirements layer



Retention-based Deletes



On-demand Deletes

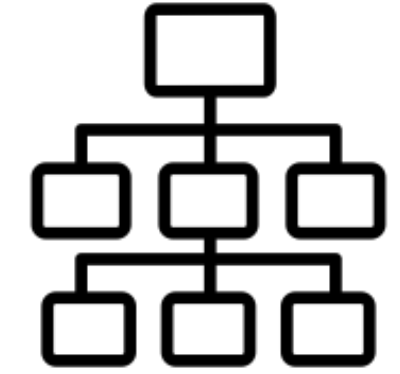
Application layer

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

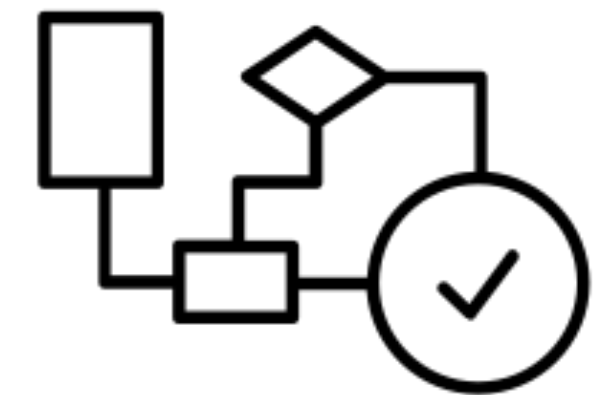
```
INSERT INTO R (...)
  WITH RET_DUR {<t> | t<i>};
```

```
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**.

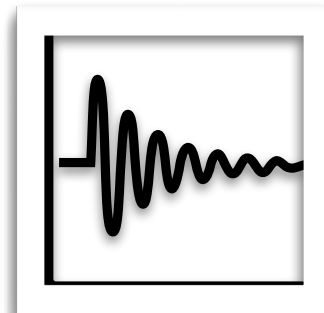
Compactions are key to building ingestion-friendly systems.

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

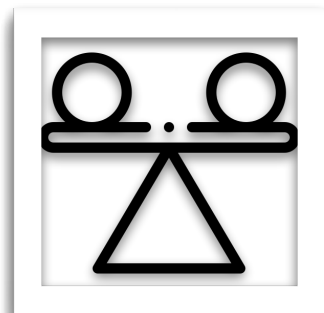
Open Research Challenges



Reduce write amplification



Workload-aware compactions & layout transformation



Performance Stability & Holistic Tuning



Automatic Tuning & Adaptive Behavior



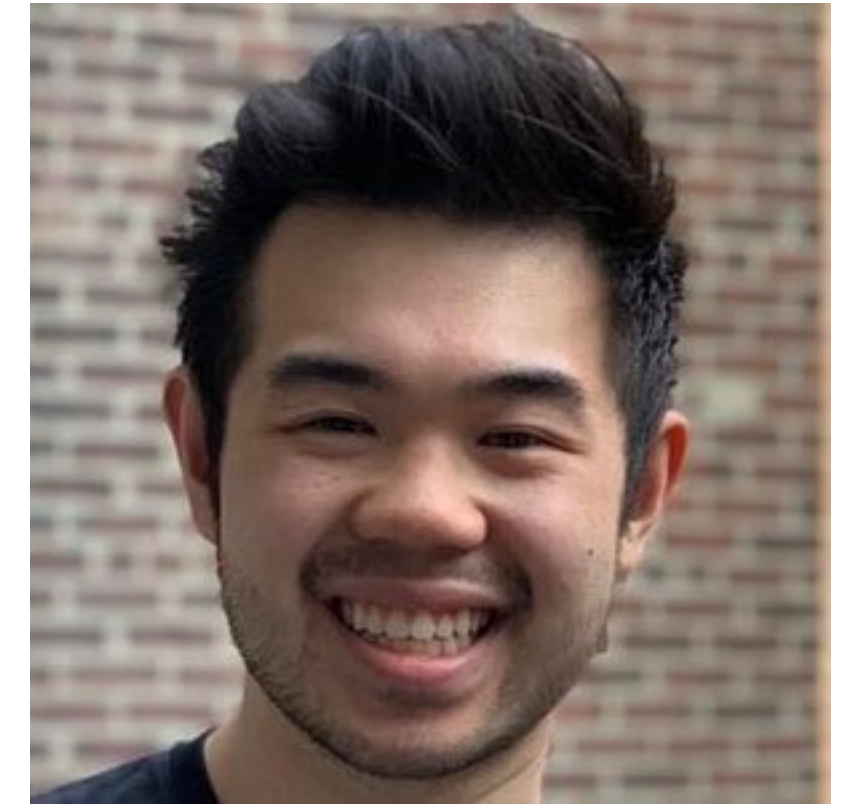
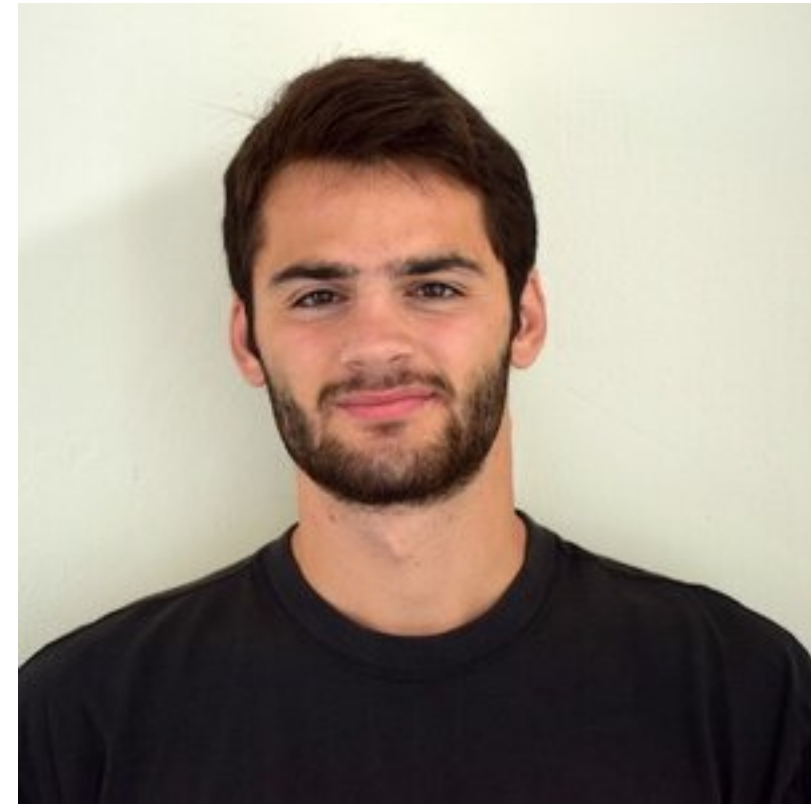
Privacy-aware LSM designs

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 of 2018. *California Civil Code*, 2018.
- [3] The California Consumer Privacy Act of 2018. *California Civil Code*, 2018.
- [4] Virginia Consumer Data Protection Act. *Virginia Code*, 2018.
- [5] H. Abu-Libdeh, A. Ly, and S. Idreos. Designing Key-Value Stores. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [6] W. Y. Alkover, S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [7] S. Alsubaiee, I. Cetindil, S. Kim, C. Li, C. Wang, and T. West. Designing Key-Value Stores. *the VLDB Endowment*, 13(9):1388–1400, 2019.
- [8] Amazon. *Amazon DynamoDB*. <https://docs.aws.amazon.com/dynamodb/latest/APIReference/>, 2017.
- [9] Apache. *Apache HBase*. <https://hbase.apache.org/>, 2017.
- [10] Apache. *Apache Cassandra*. <https://cassandra.apache.org/>, 2017.
- [11] M. Athanassoulis, S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Efficient On-Disk Key-Value Stores. *International Conference on Database Systems (ICDE)*, 2018.
- [12] M. Athanassoulis, S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Efficient On-Disk Key-Value Stores. *International Conference on Database Systems (ICDE)*, 2018.
- [13] M. Athanassoulis, S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Efficient On-Disk Key-Value Stores. *International Conference on Database Systems (ICDE)*, 2018.
- [14] M. Athanassoulis, S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Efficient On-Disk Key-Value Stores. *International Conference on Database Systems (ICDE)*, 2018.
- [15] M. Athanassoulis, S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Efficient On-Disk Key-Value Stores. *International Conference on Database Systems (ICDE)*, 2018.
- [16] O. Balmau, S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Designing Key-Value Stores. *the VLDB Endowment*, 13(9):1388–1400, 2019.
- [17] S. Idreos, N. Dayan, W. Qin, M. Akmanalp, S. Hilgard, A. Ross, J. Lennon, V. Jain, H. Gupta, D. Li, and Z. Zhu. Design Continuums and the Path Toward Self-Organizing Key-Value Stores. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [18] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [19] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [20] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [21] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [22] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [23] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [24] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [25] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [26] S. Idreos, K. Kester, D. G. Andersen, and M. Callaghan. Data Structure for Key-Value Stores. *VLDB Endowment*, 13(9):1388–1400, 2019.
- [27] 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.
- [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*, pages 241–254, 2020.
- [29] 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*, pages 241–254, 2020.
- [30] 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*, pages 241–254, 2020.
- [31] 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*, pages 241–254, 2020.
- [32] 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*, pages 241–254, 2020.
- [33] C. Luo and M. J. Carey. LSM-based Storage Techniques: A Survey. *The VLDB Journal*, 29(1):393–418, 2020.
- [34] 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 ACM SIGMOD International Conference on Management of Data*, pages 2667–2672, 2020.
- [35] S. Sarkar, D. Staratzis, Z. Zhu, and M. Athanassoulis. Constructing and Analyzing the LSM Compaction Design Space. *Proceedings of the VLDB Endowment*, 14(11):2216–2229, 2021.
- [36] ScyllaDB. Online reference. <https://www.scylladb.com/>.
- [37] R. Sears and R. Ramakrishnan. bLSM: A General Purpose Log Structured Merge Tree. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 217–228, 2012.
- [38] M. I. Seltzer. Berkeley DB: A Retrospective. *IEEE Data Engineering Bulletin*, 30(3):21–28, 2007.
- [39] W. Tan, S. Tata, Y. R. Tang, and L. L. Fong. Diff-Index: Differentiated Index in Distributed Log-Structured Data Stores. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pages 700–711, 2014.
- [40] Y. R. Tang, A. Iyengar, W. Tan, L. L. Fong, L. Liu, and B. Palanisamy. Deferred Lightweight Indexing for Log-Structured Key-Value Stores. In *15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGrid 2015, Shenzhen, China, May 4-7, 2015*, pages 11–20, 2015.
- [41] D. Teng, L. Guo, R. Lee, F. Chen, S. Ma, Y. Zhang, and X. Zhang. LSbM-tree: Re-Enabling Buffer Caching in Data Management for Mixed Reads and Writes. In *37th IEEE International Conference on Distributed Computing Systems, ICDCS 2017, Atlanta, GA, USA, June 5-8, 2017*, pages 68–79, 2017.
- [42] R. Thonangi, S. Babu, and J. Yang. A Practical Concurrent Index for Solid-State Drives. In *Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM)*, pages 1332–1341, 2012.
- [43] R. Thonangi and J. Yang. On Log-Structured Merge for Solid-State Drives. In *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, pages 683–694, 2017.
- [44] T. Vinçon, S. Hardock, C. Riegger, J. Oppermann, A. Koch, and I. Petrov. NoFTL-2017.
- [45] 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*, pages 241–254, 2020.
- [46] 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*, pages 241–254, 2020.
- [47] 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*, pages 241–254, 2020.
- [48] 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*, pages 241–254, 2020.
- [49] 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*, pages 241–254, 2020.
- [50] 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*, pages 241–254, 2020.
- [51] 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.
- [52] 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 Invalidation in LSM-tree based Storage Engines. *Proceedings of the VLDB Endowment*, 13(11):1976–1989, 2020.
- [53] 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.
- [54] T. Yao, J. Wan, P. Huang, X. He, 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.
- [55] 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.
- [56] 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.
- [57] 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.
- [58] W. Zhang, X. Zhao, S. Jiang, and H. Jiang. ChameleonDB: a key-value store for optane persistent memory. In *EuroSys '21: Sixteenth European Conference*

The DiSC Lab



disc.bu.edu

Don't Miss: Compactionary

Comparative Analysis Individual Analysis

Workload
#Entries: 10000000
Entry size: 128 B

Main Memory Parameter
Buffer size: 16 MB
BF size / Entry: 10 MB

Disk Parameters
File size (in terms of buffer): 1
Page size: 4 KB

Data Layout
Size ratio: 4
Leveling **Leveling** Tiering
Step: 1

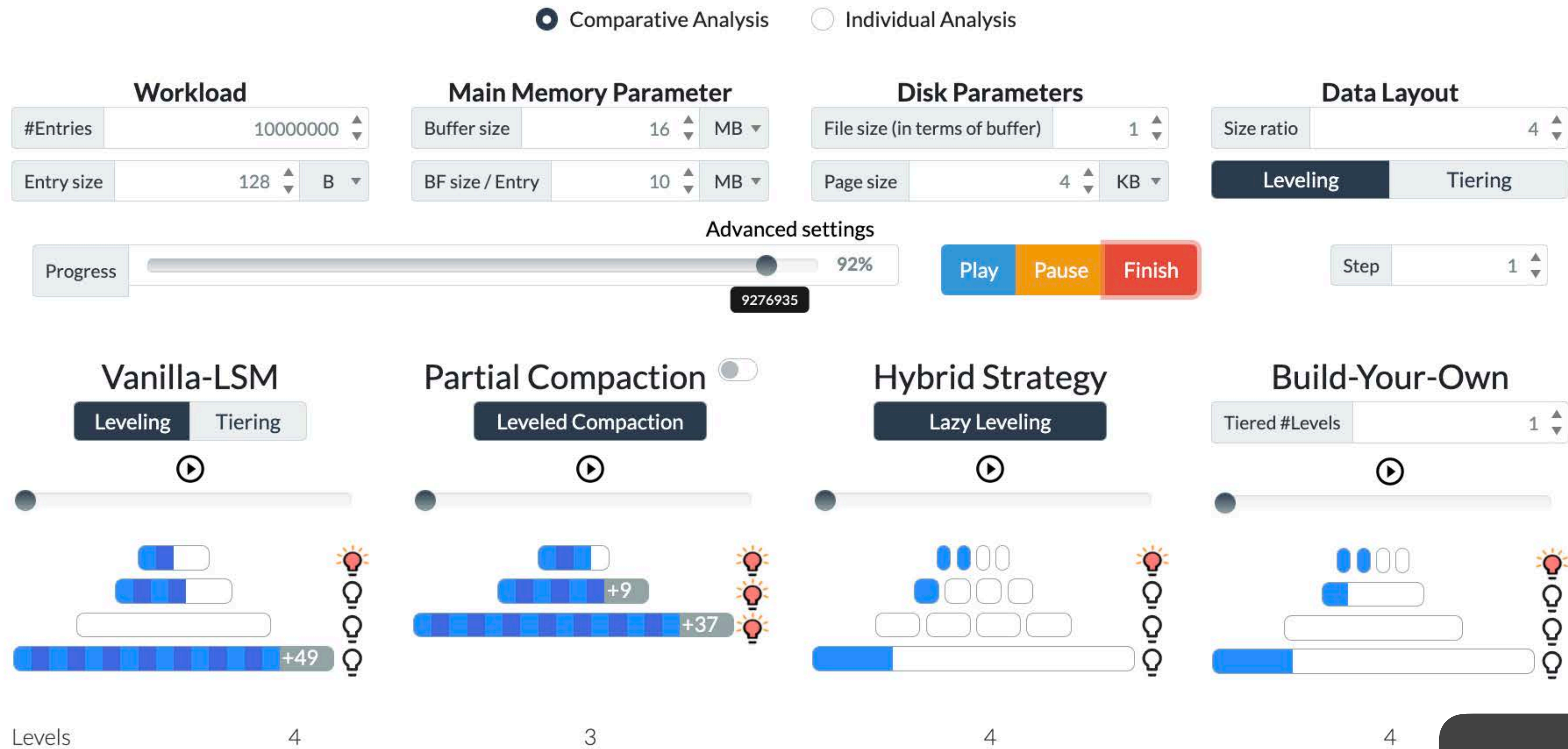
Advanced settings
Progress: 92% (9276935)
Play Pause Finish

Vanilla-LSM
Leveling **Leveling** Tiering
Levels: 4

Partial Compaction
Leveled Compaction
Levels: 3

Hybrid Strategy
Lazy Leveling
Levels: 4

Build-Your-Own
Tiered #Levels: 1
Levels: 4



disc-projects.bu.edu/compactionary



4
5
31

Tuesday
4-6PM

Dissecting, Designing, and Optimizing
LSM-Based Data Stores

Subhadeep Sarkar

Manos Athanassoulis

Thank You!

Questions?

BOSTON
UNIVERSITY

lab
DiSC