

# *CS 561: Data Systems Architecture*

## Class 5

# **Log-Structured Merge (LSM) Trees**

BOSTON  
UNIVERSITY

*Dr. Subhadeep Sarkar*

<https://bu-disc.github.io/CS561/>



# Updates: Logistics

First **technical question** is due on **02/07**.

First **review** is due on **02/14**.

Deadline for **Project 0** extended till **02/05**.

The first **student presentation** is in two weeks (on **02/14**)!

**Select the paper** you will work on for your **presentation** (groups of 3-4)

**A week before the presentation**, discuss the slides with me in OH.

# LSM-tree

# LSM-tree

NoSQL

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

SQLite

relational

influxdb QuasarDB

time-series

2023

# LSM-tree

NoSQL

This block contains logos for various NoSQL databases that utilize the LSM-tree architecture. The logos are arranged in two rows. The top row includes RocksDB (a yellow cheetah), WT (the letters 'WT' in black and orange), levelDB (a green cylinder), SCYLLA (a blue alien head), and DynamoDB (a blue cylinder). The bottom row includes cassandra (an eye), tarantool (two red circles), Bigtable (a blue hexagon), APACHE HBASE (a black orca), riak (a grey starburst), and accumulo (a grid of squares).

This block contains two logos. On the left is SQLite, featuring a blue square with a white feather and the text 'SQLite'. On the right is a dark blue square containing a white silhouette of a dolphin.

relational

This block contains two logos. On the left is influxdb, featuring a blue cube and the text 'influxdb'. On the right is QuasarDB, featuring a blue grid pattern.

time-series

2023

# Why **LSM** ?

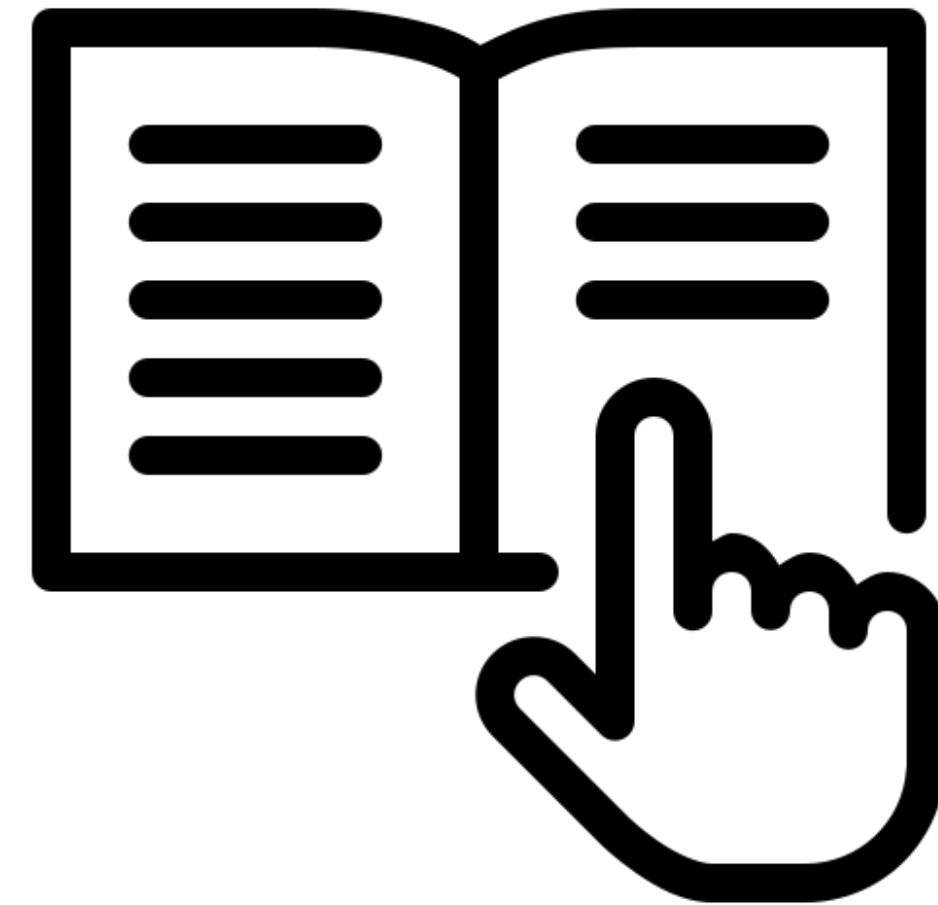


fast ingestion

# Why **LSM** ?

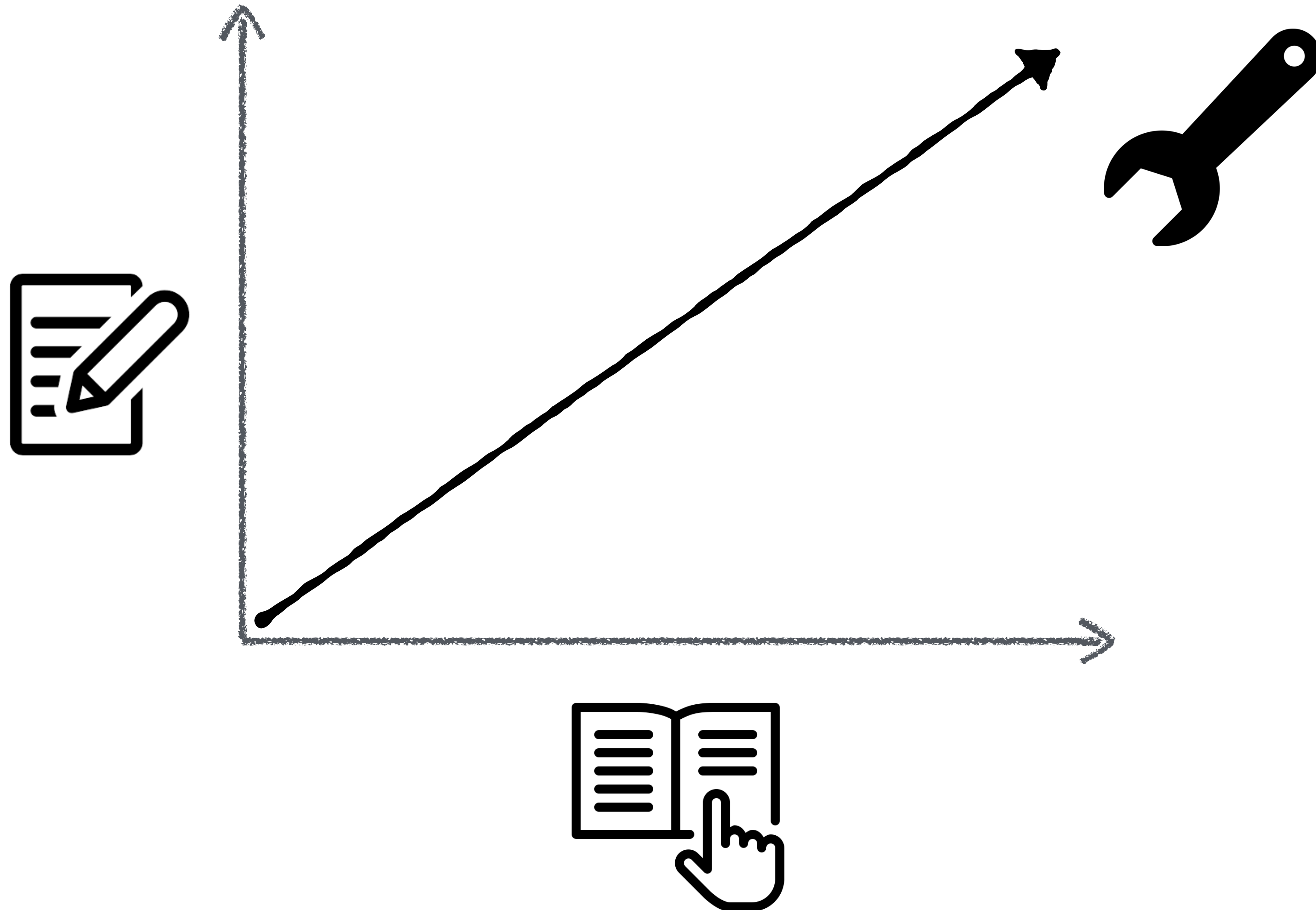


fast ingestion



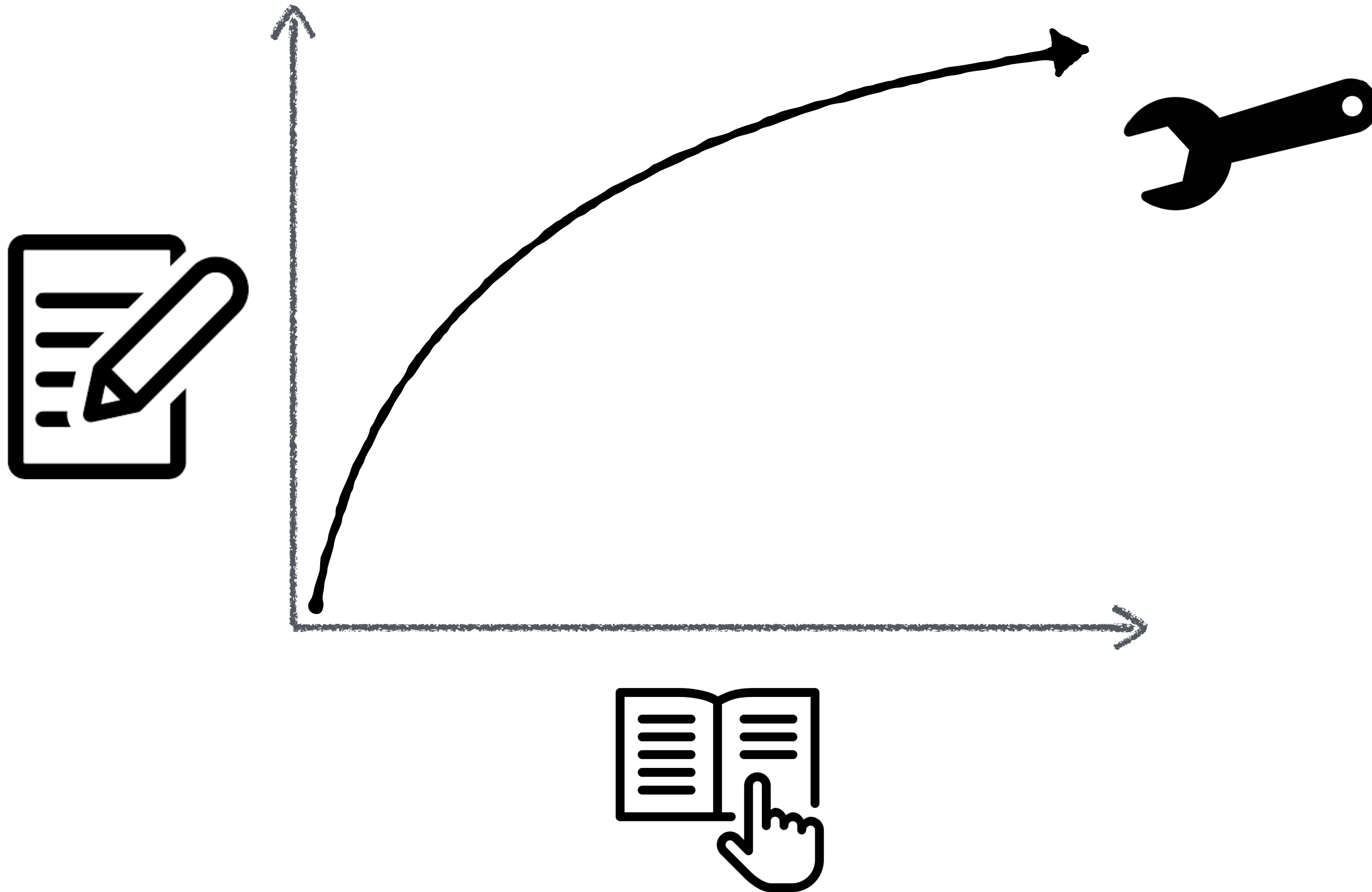
competitive reads

# Why **LSM** ?

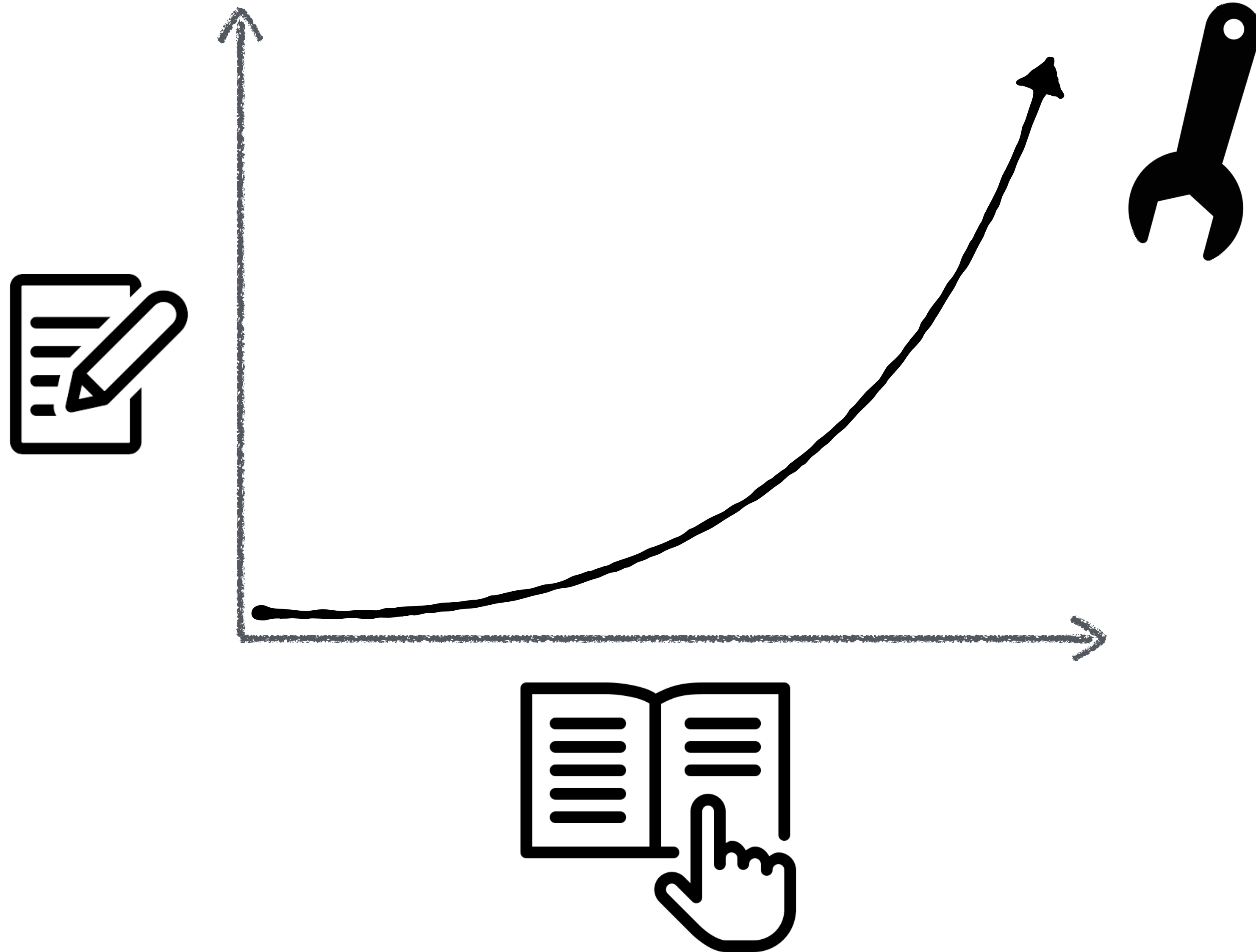




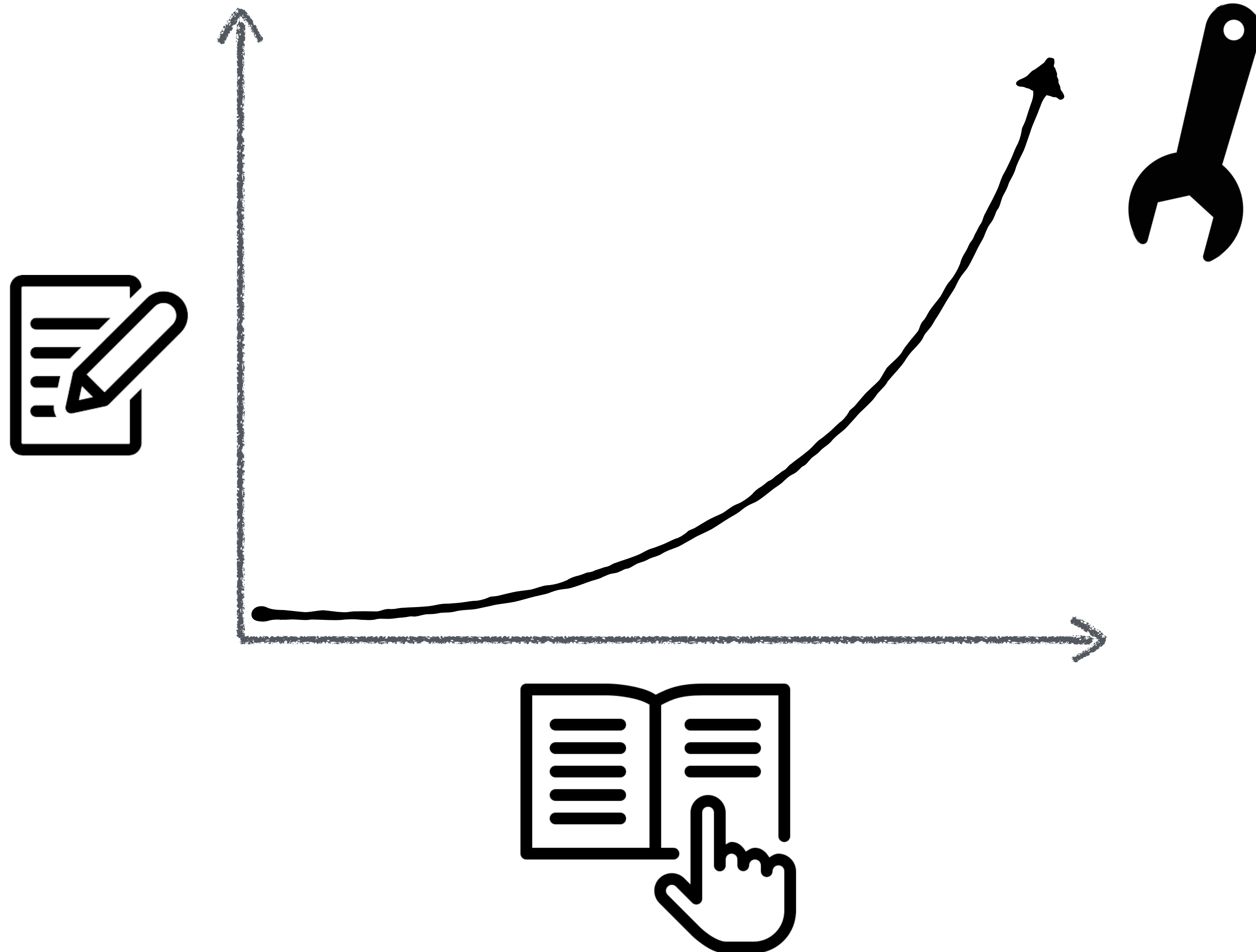
# Why **LSM** ?



# Why **LSM** ?



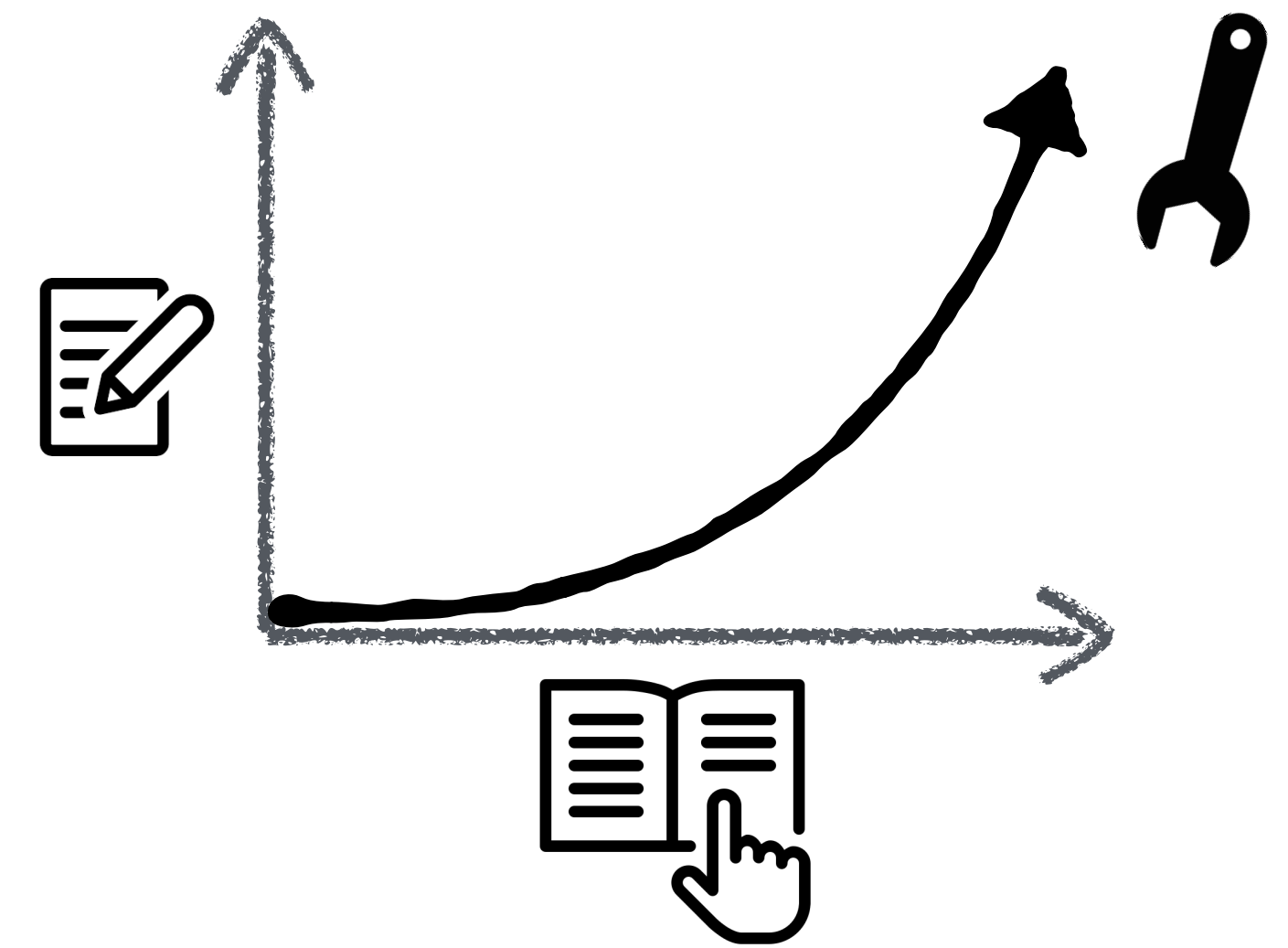
# Why **LSM** ?



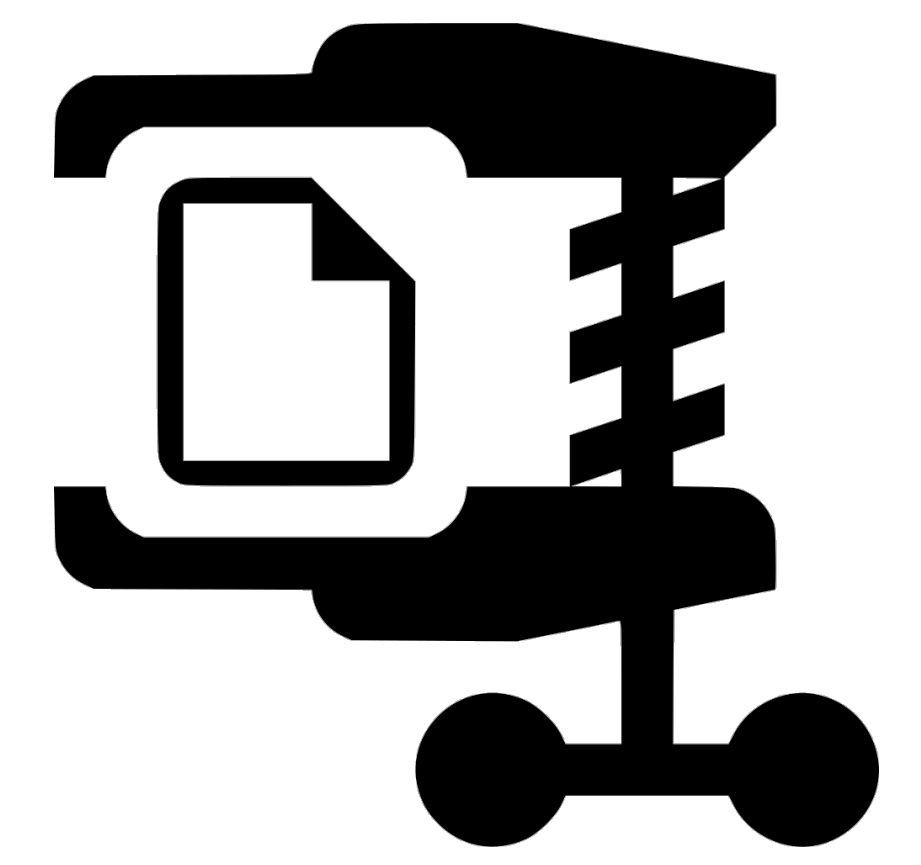
# Why **LSM** ?



fast writes

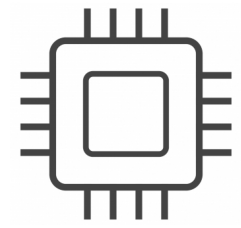


tunable read-write performance

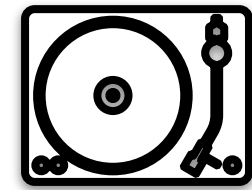


good space utilization

# LSM Basics



buffer



level 1



level 2



size ratio:  $T$

level 3

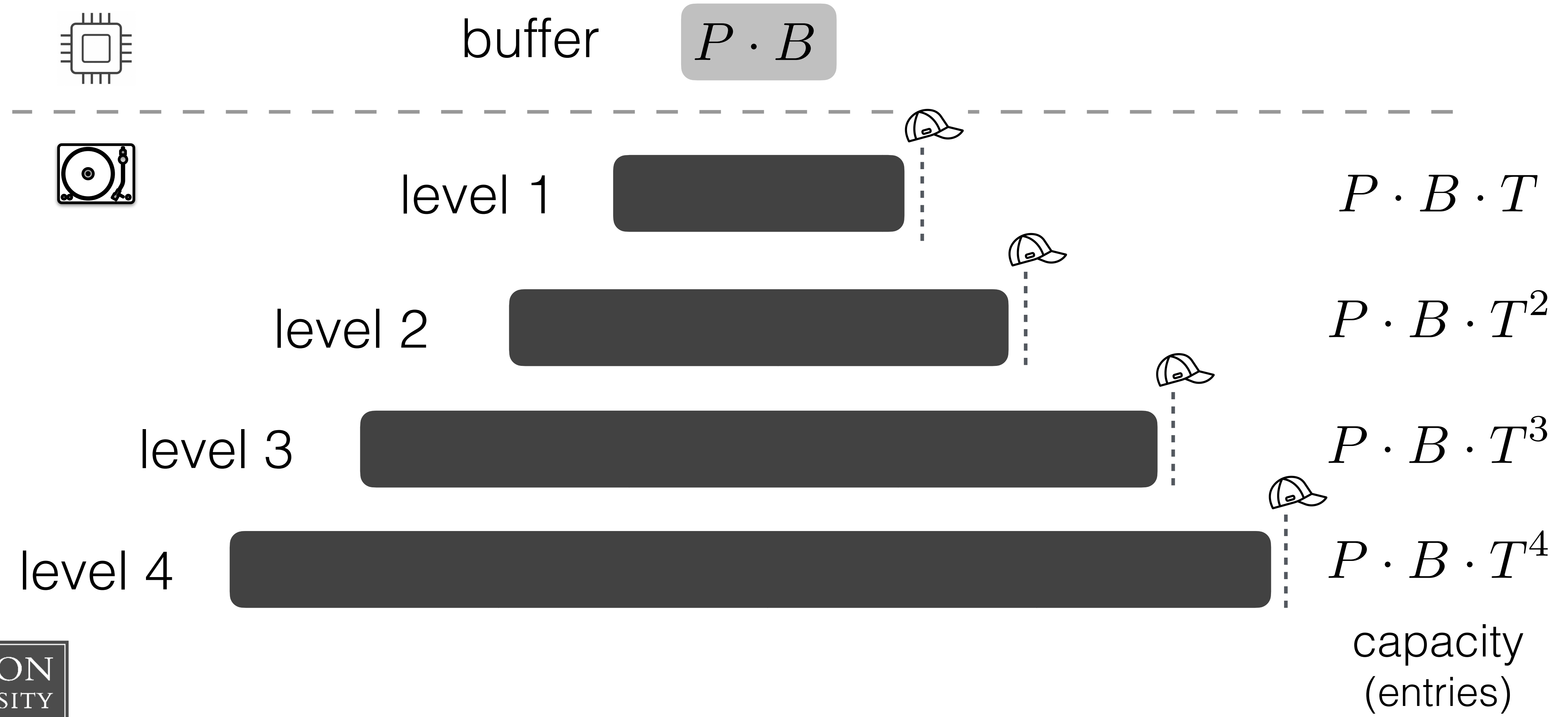


level 4



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

# LSM Basics



# LSM Operating Principles

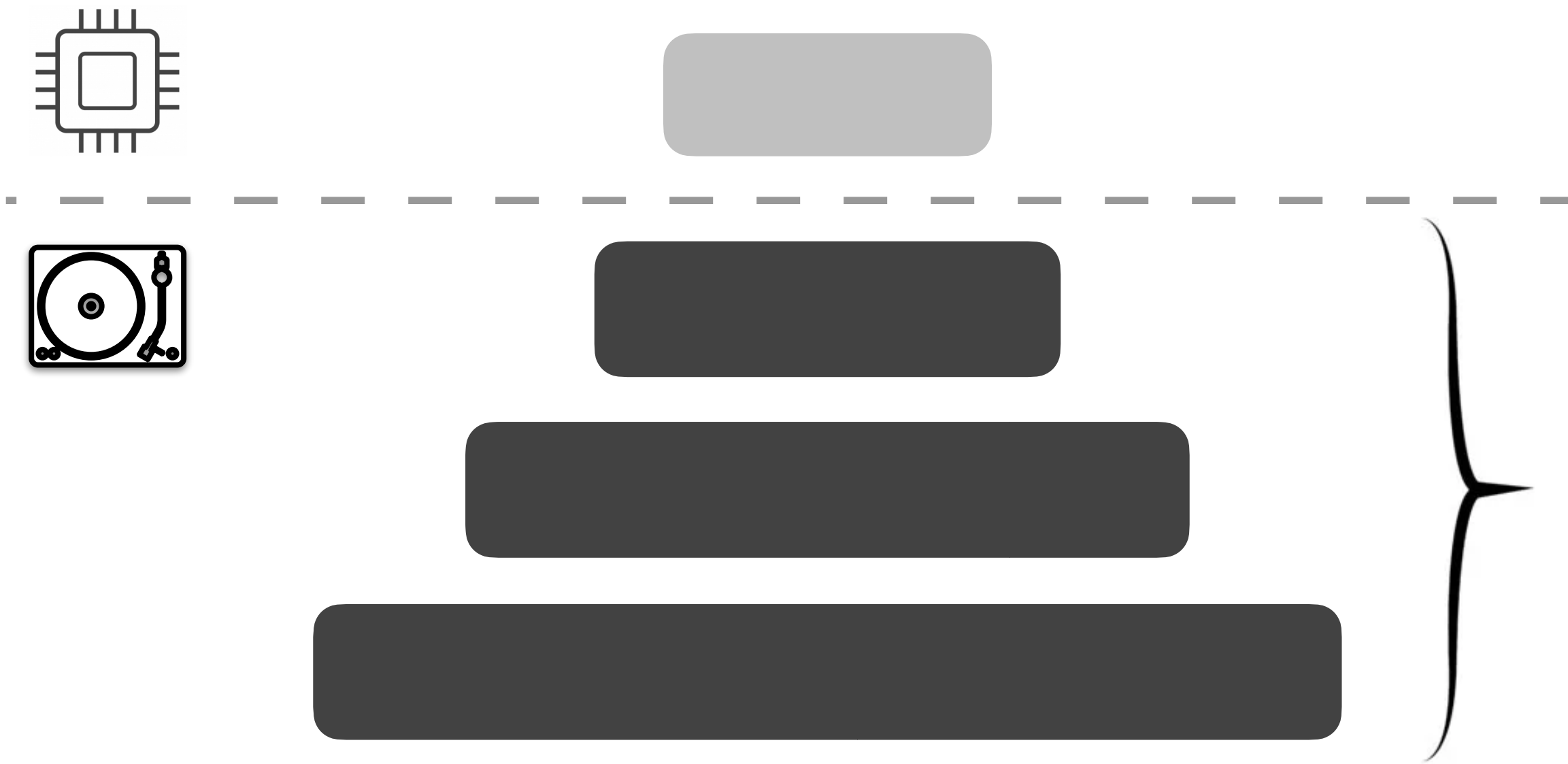
Buffering ingestion

Immutable files on storage

Out-of-place updates & deletes

Periodic data layout reorganization

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



most data  
on storage

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

99.9% on storage



How does the storage layer affect ingestion?



# Data **Layout**

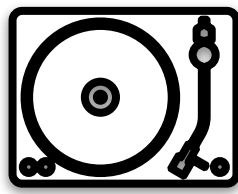
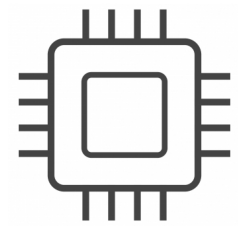
Classical LSM design: **leveling**

[eager merging]

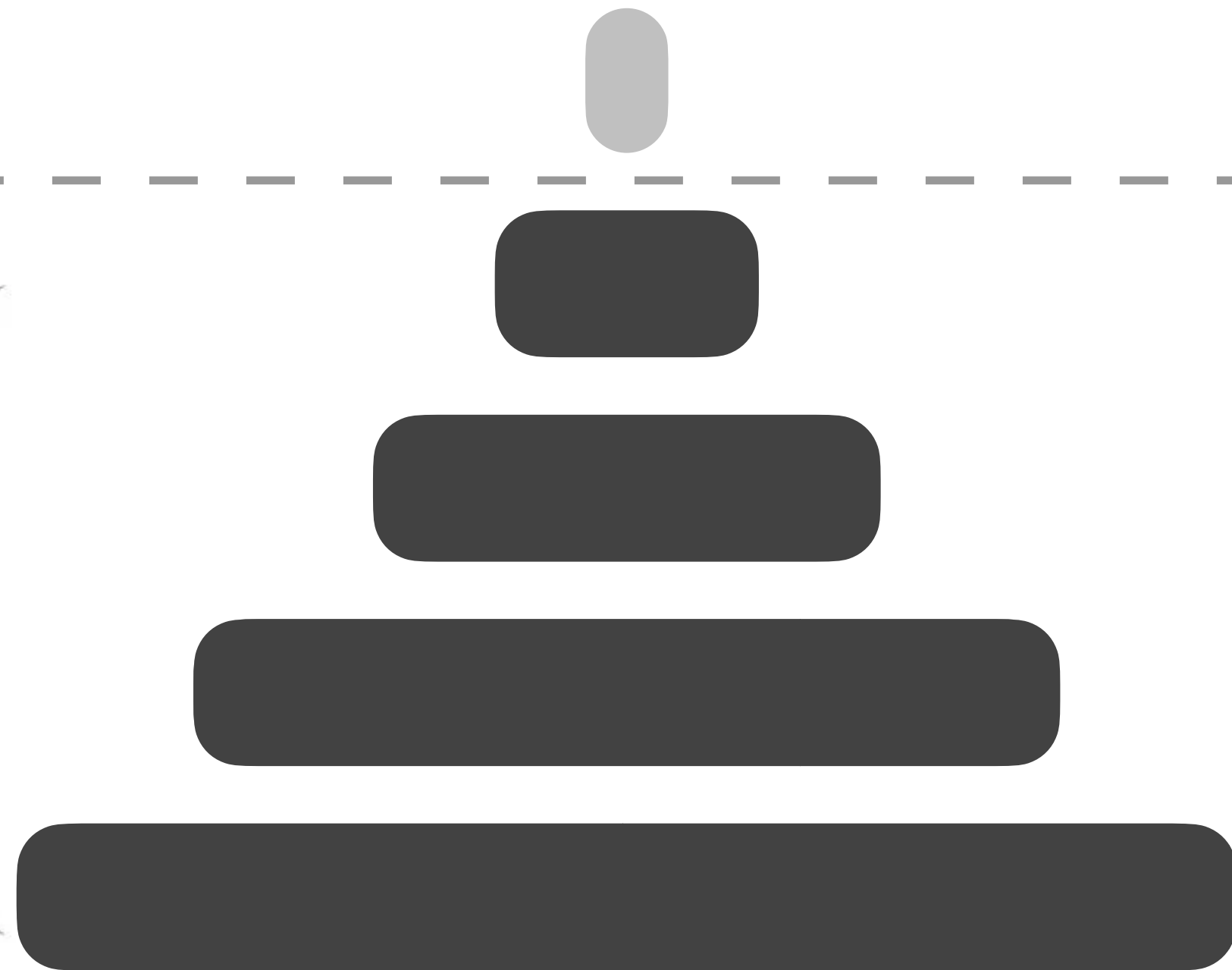


# Data Layout

**leveling** [eager]



1 run  
per level



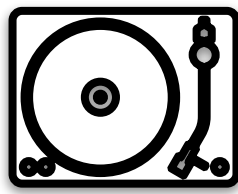
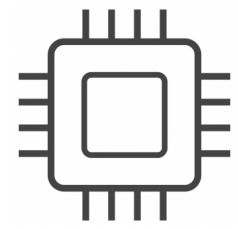
- good read performance
- good space amplification
- high write amplification

Seems like heaven!

# Data Layout

**leveling** [eager]

**tiering** [lazy]



1 run  
per level

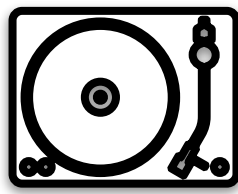
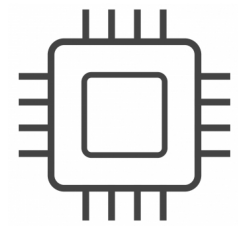


- good read performance
- good space amplification
- high write amplification

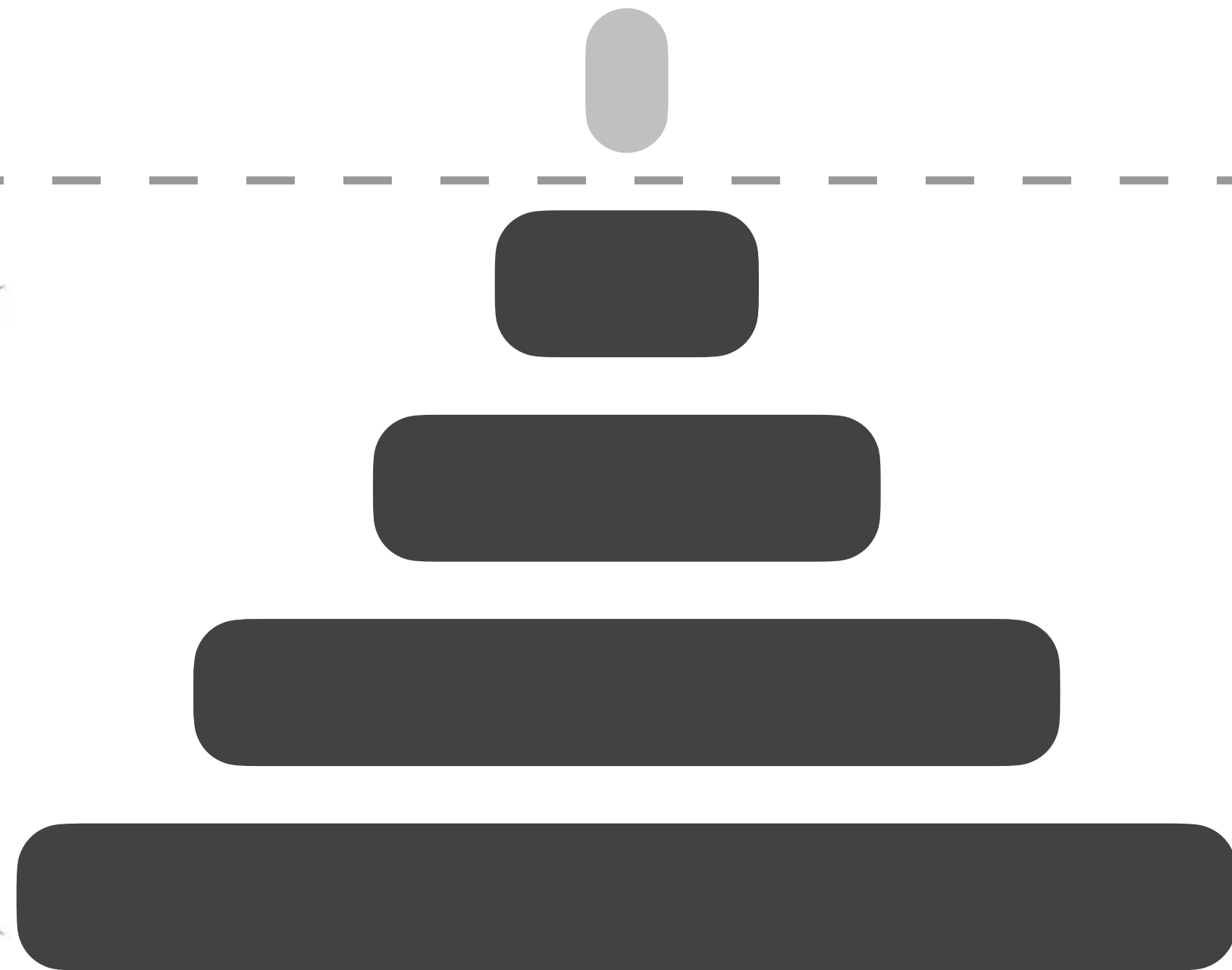
# Data Layout

**leveling** [eager]

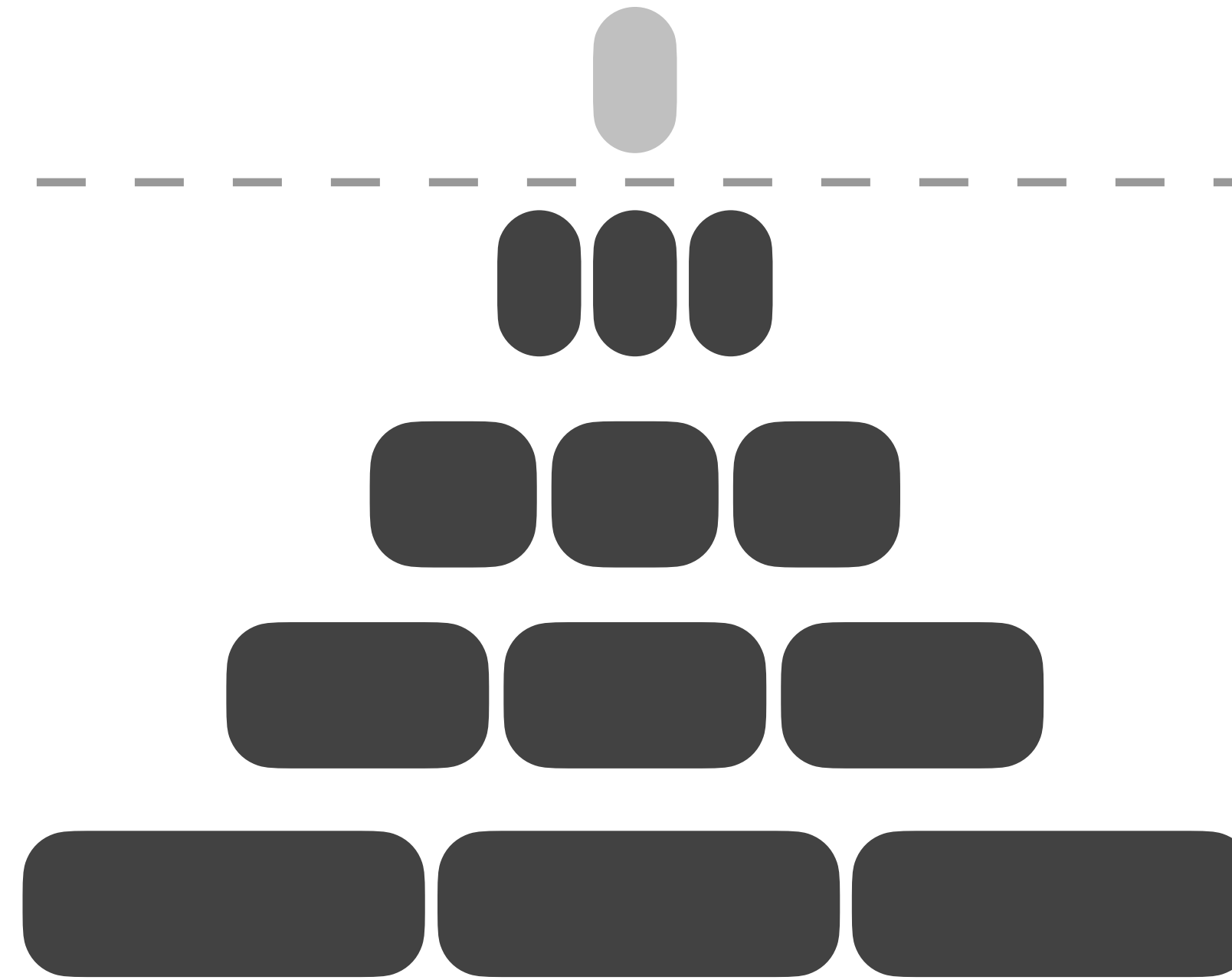
**tiering** [lazy]



1 run  
per level



T runs  
per level



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

- good ingestion performance

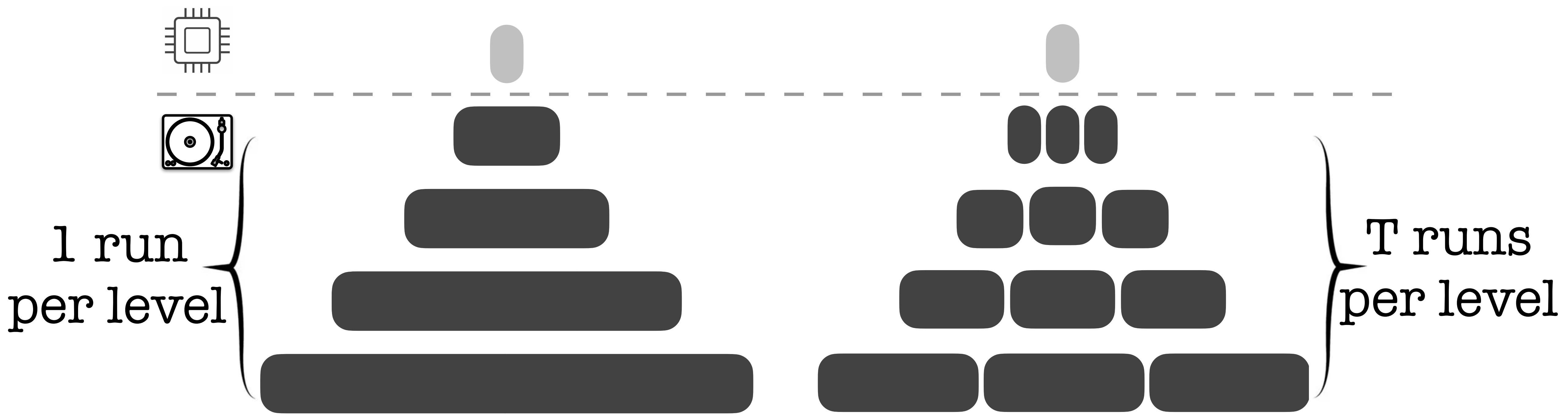
Limitations?

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

# Data Layout

**leveling** [eager]

**tiering** [lazy]



1 run  
per level

$T$  runs  
per level

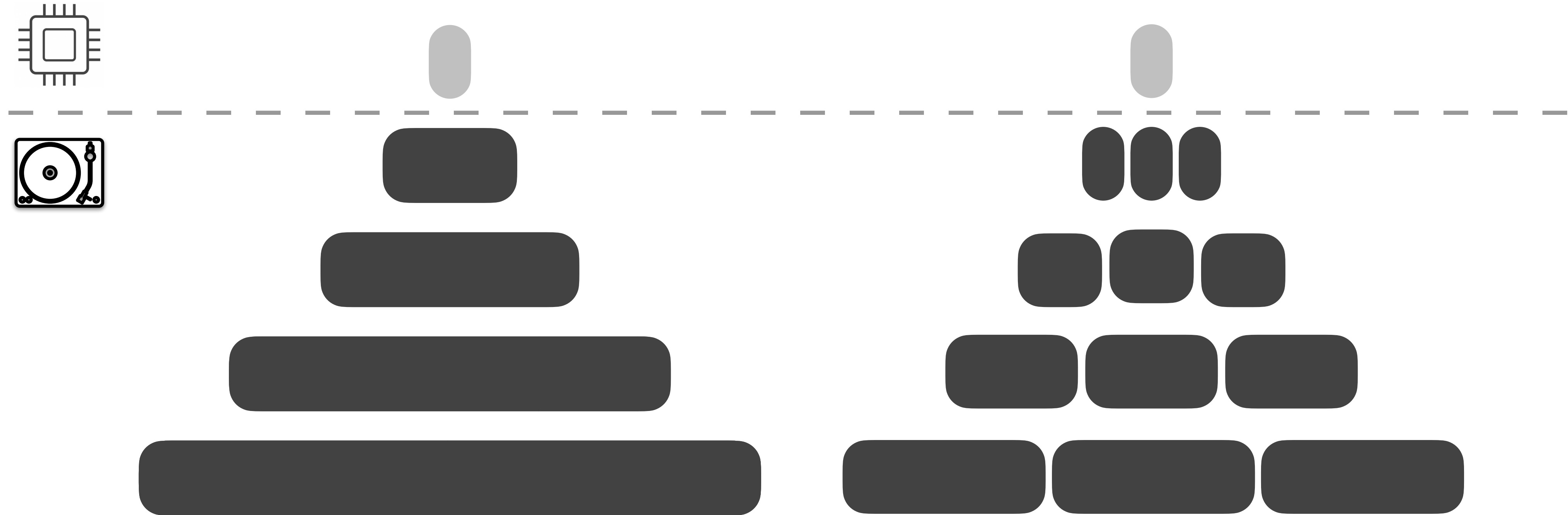
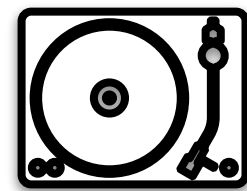
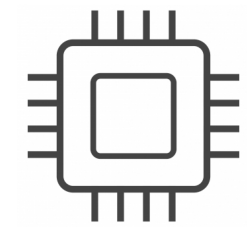
Read cost:  $\mathcal{O}(L \cdot \phi)$   
 Write cost:  $\mathcal{O}(T \cdot L/B)$   
 SA:  $\mathcal{O}(1/T)$

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

# Data Layout

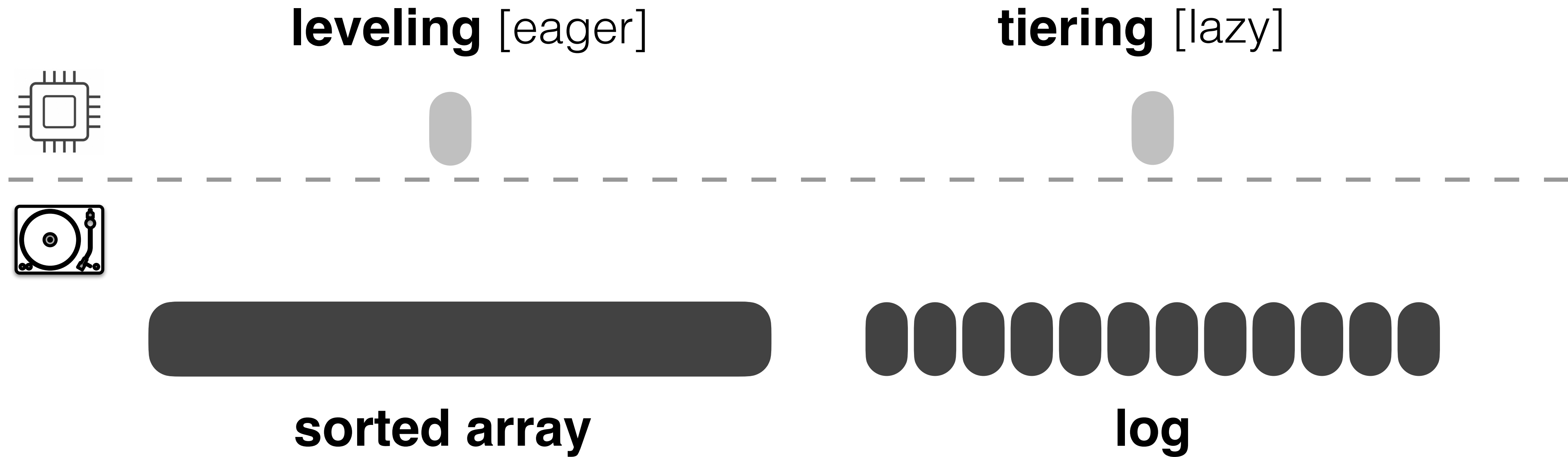
**leveling** [eager]

**tiering** [lazy]

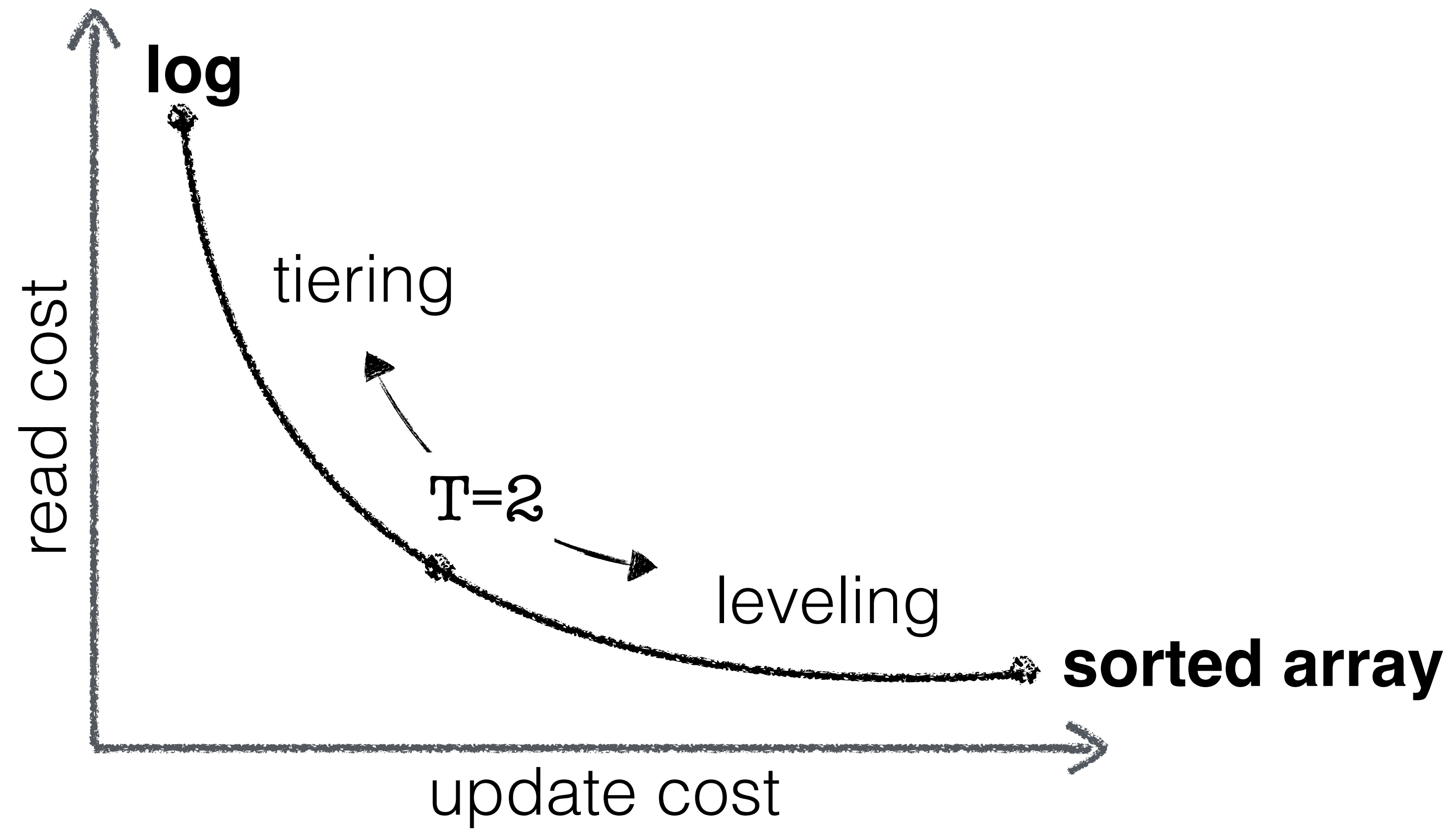


What happens if T becomes too large?

# Data Layout



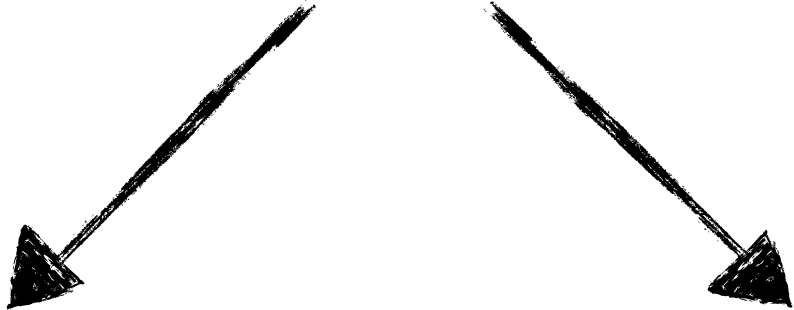
# Data Layout



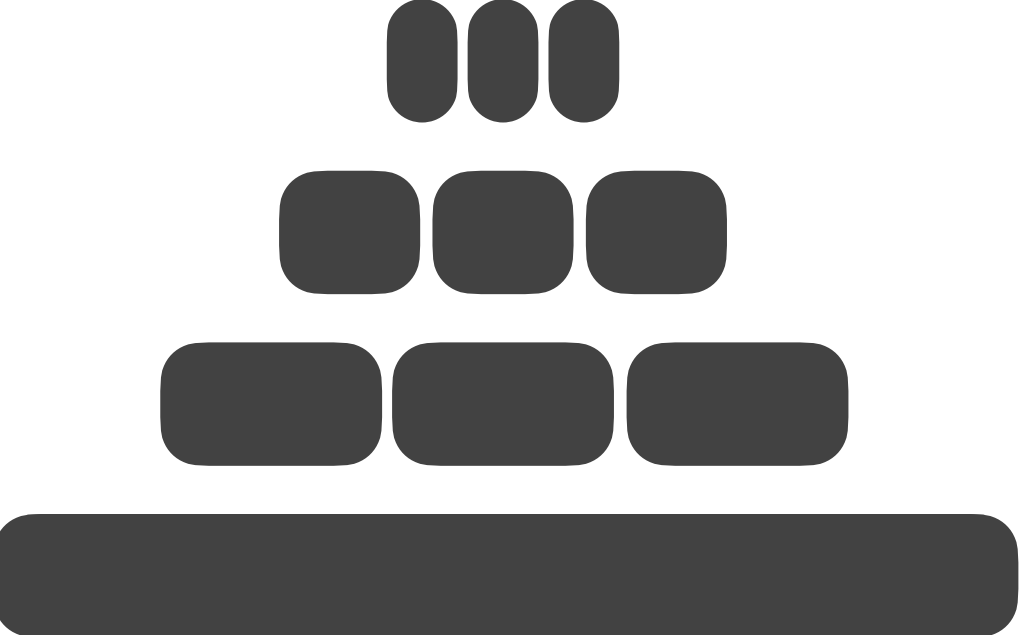
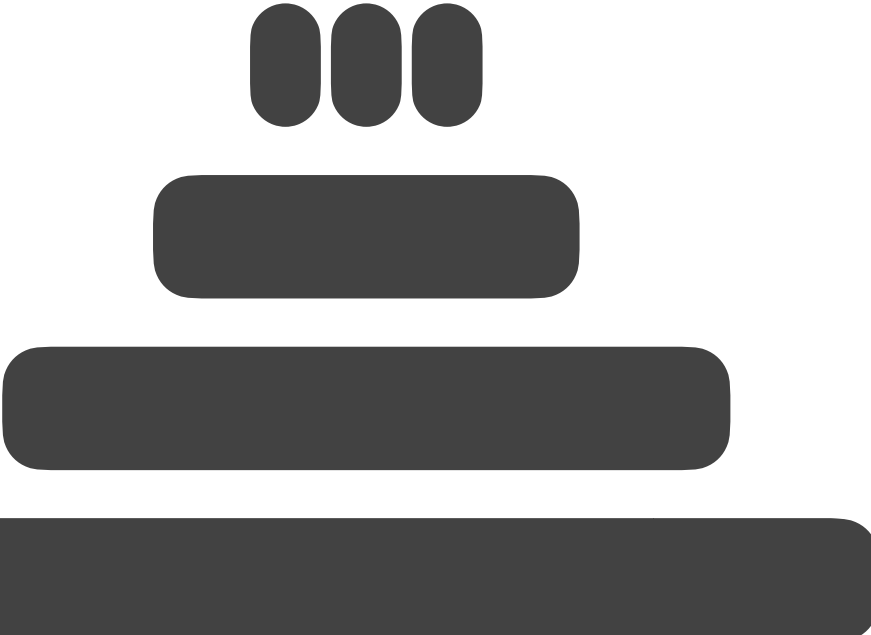
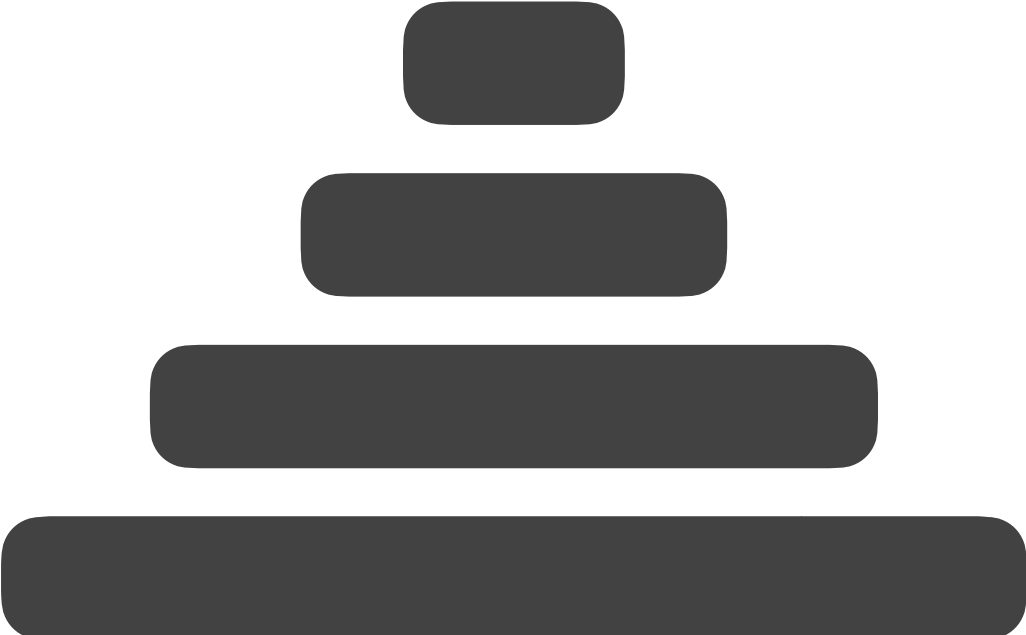


# Data Layout

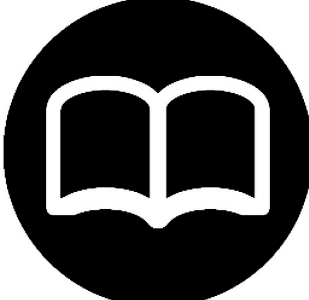
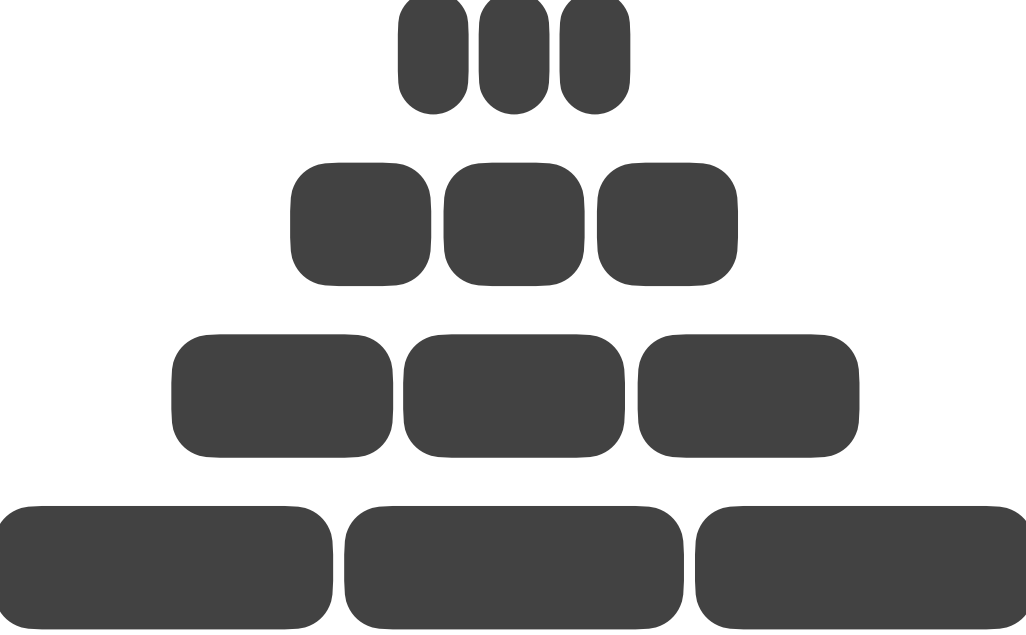
## hybrid designs



leveling



tiering



read  
optimized




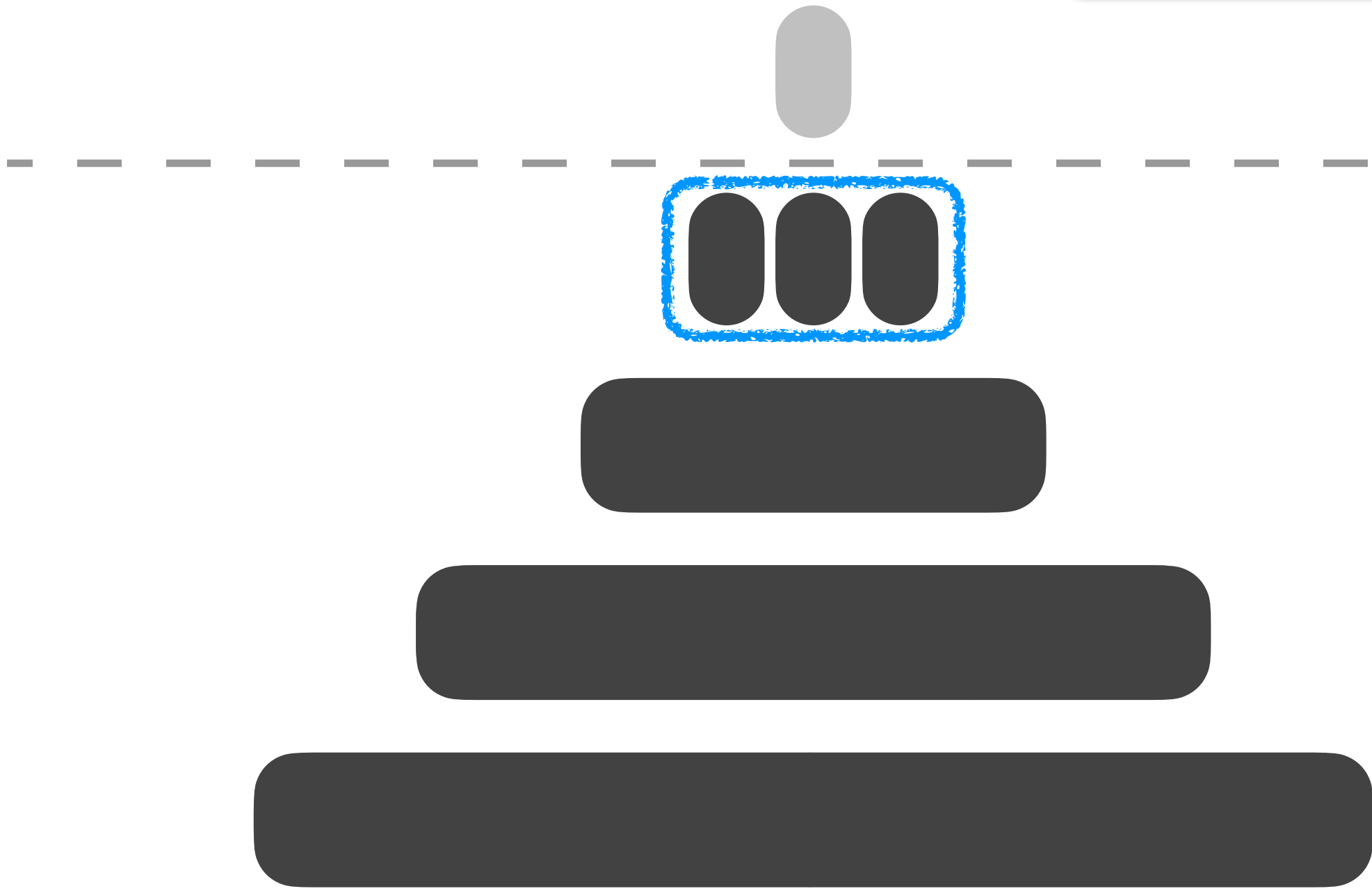
write  
optimized



# Data Layout

**1-leveling**

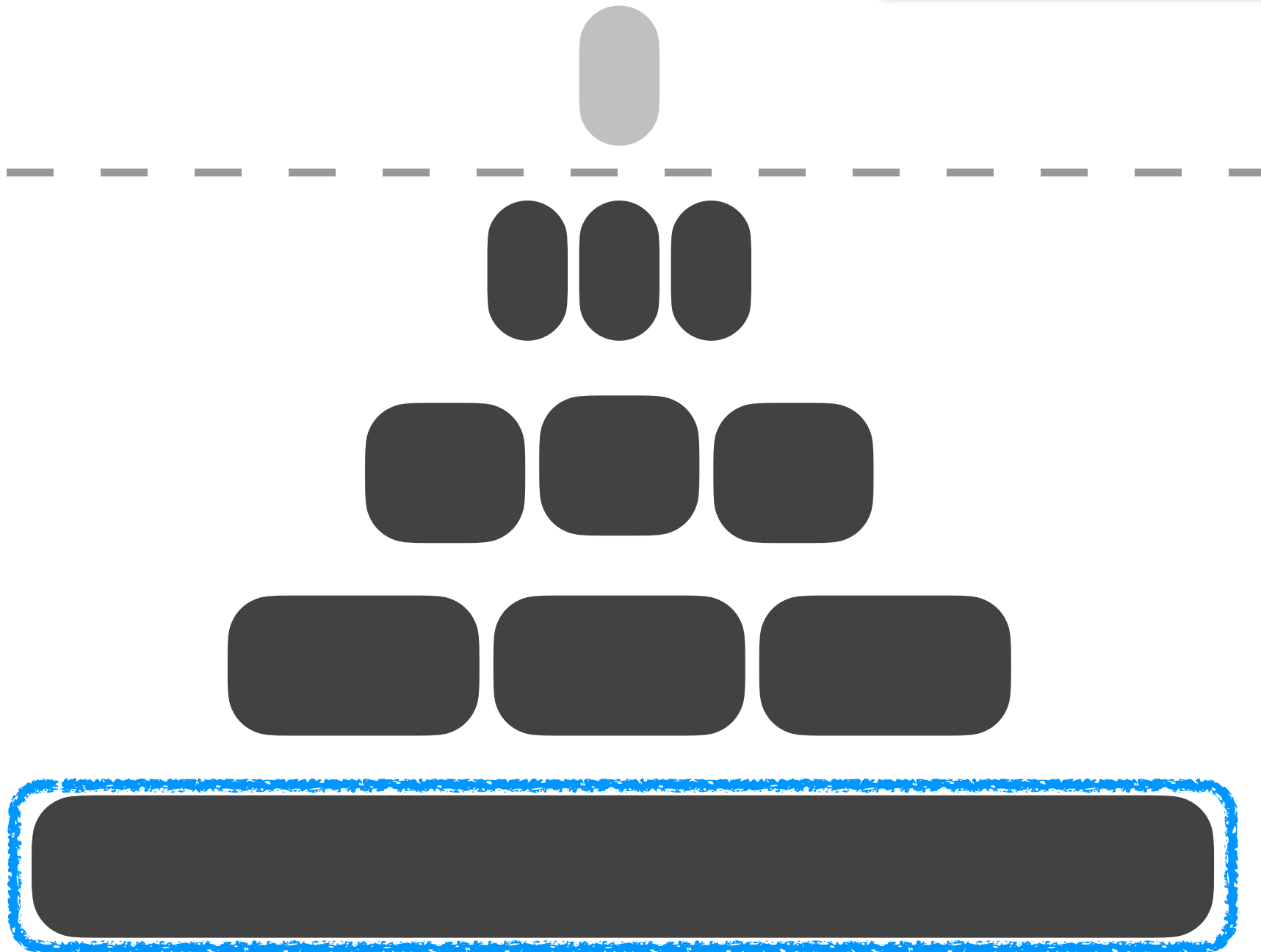
RocksDB20 



- fewer write stalls
- increased block cache hits

**L-leveling**

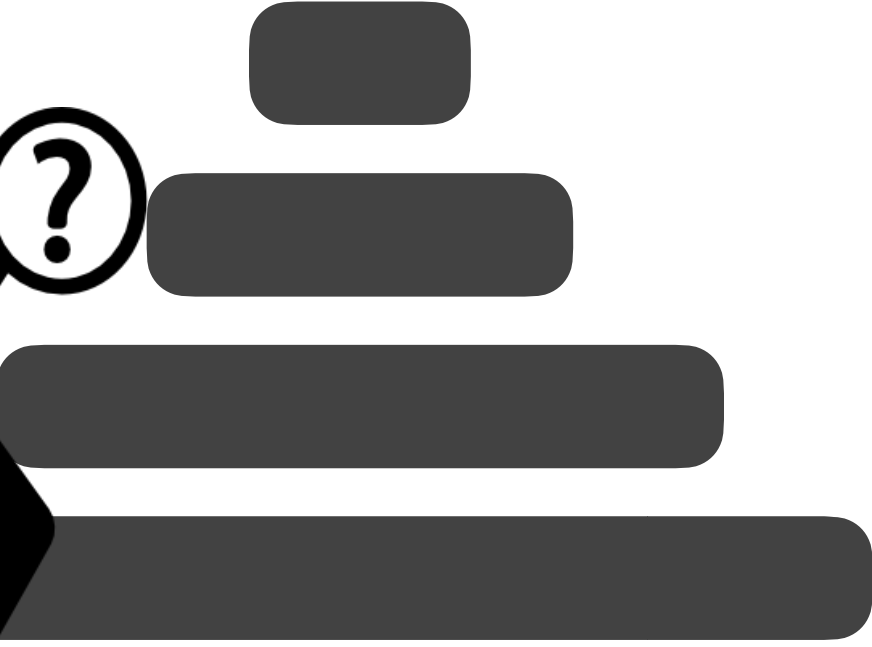
DayanSIGMOD18 



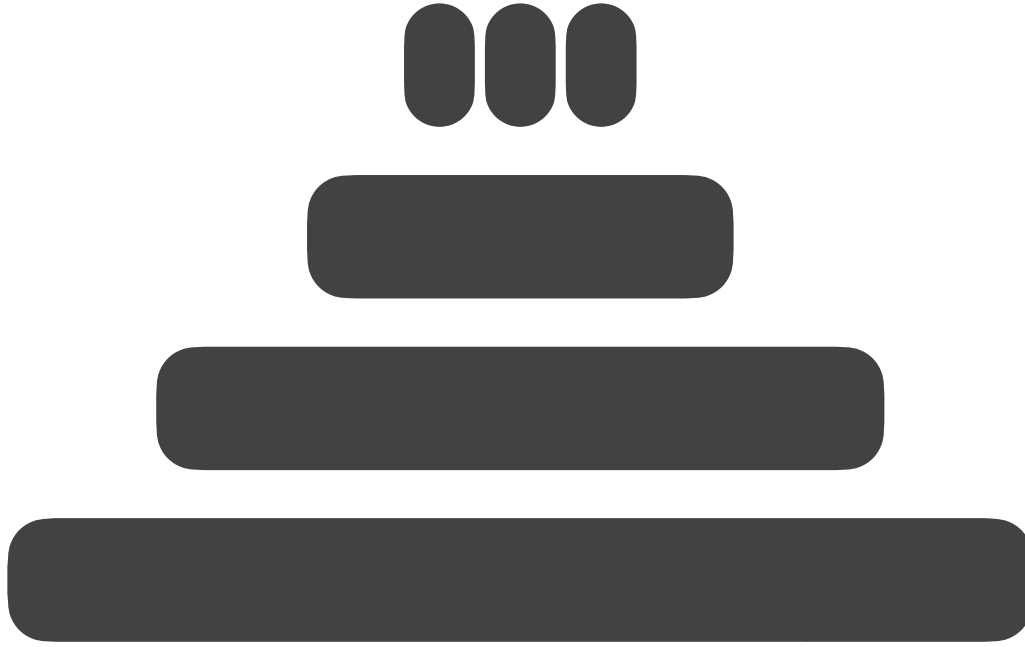
- low write amplification
- better read performance

# Data **L**ayout

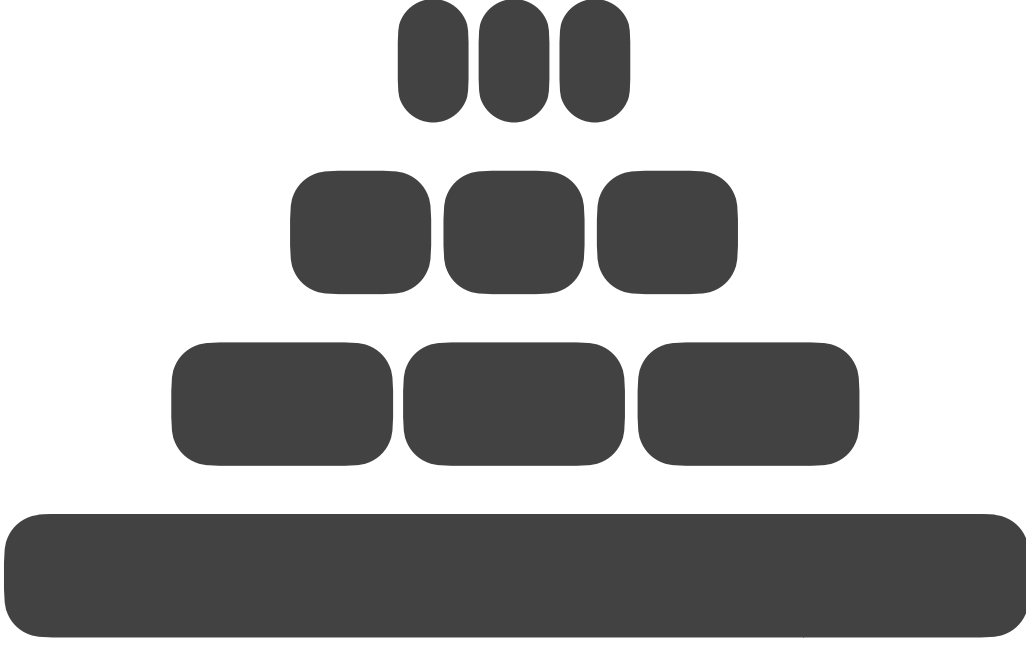
leveling



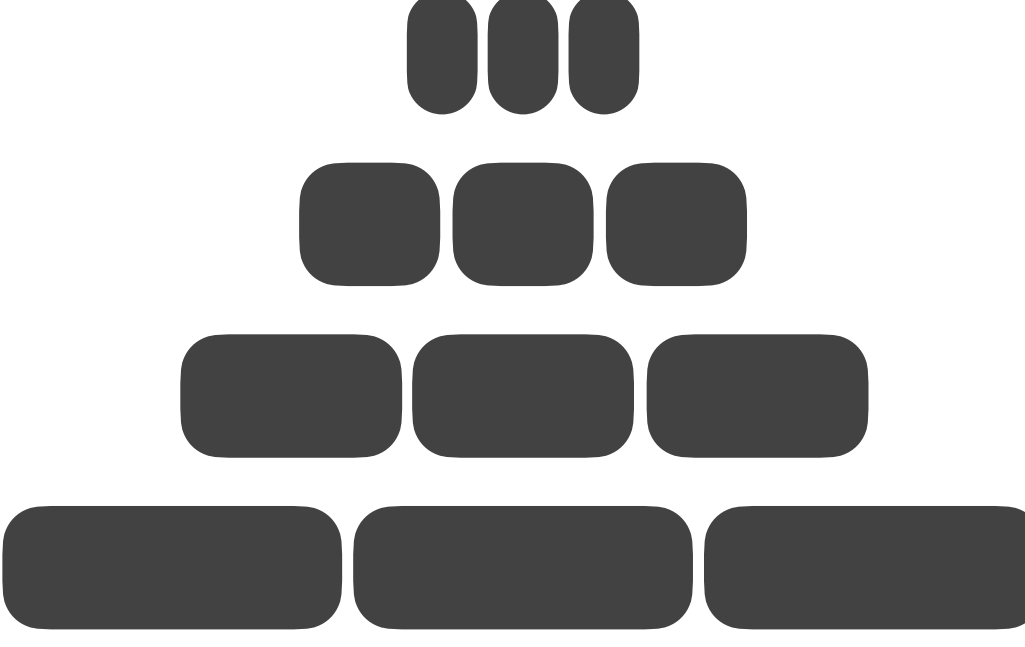
1-leveling



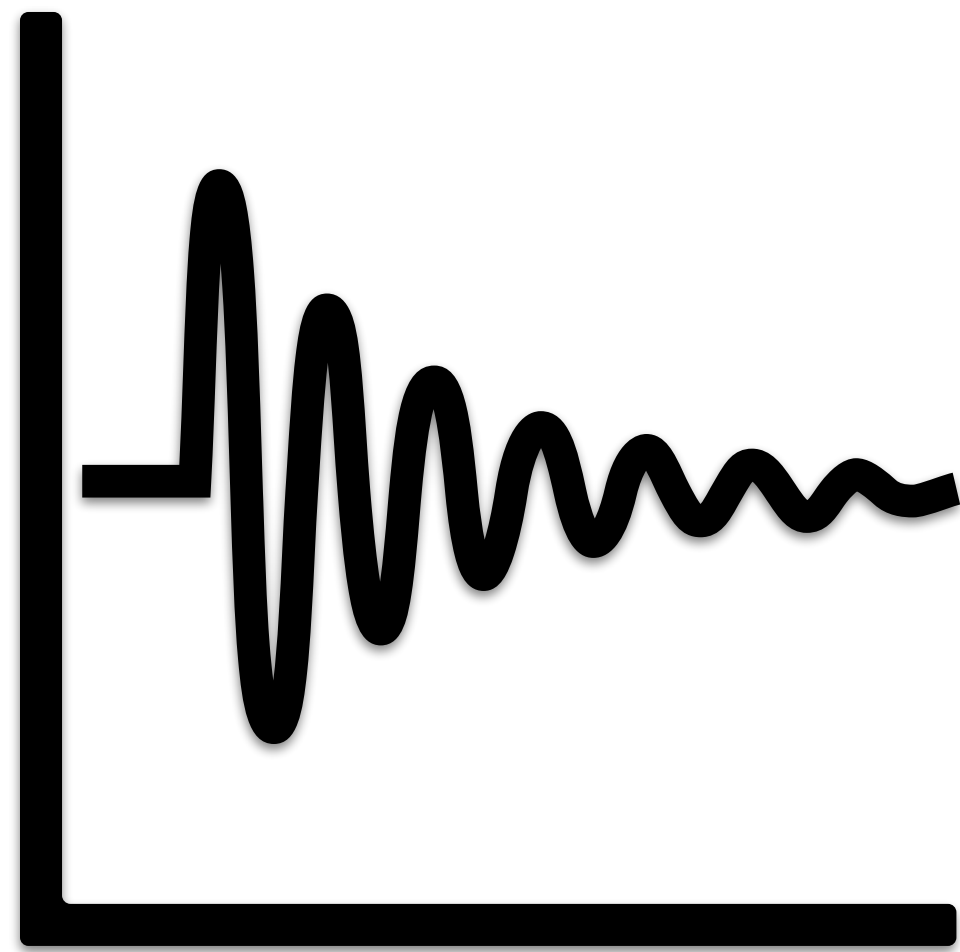
L-leveling



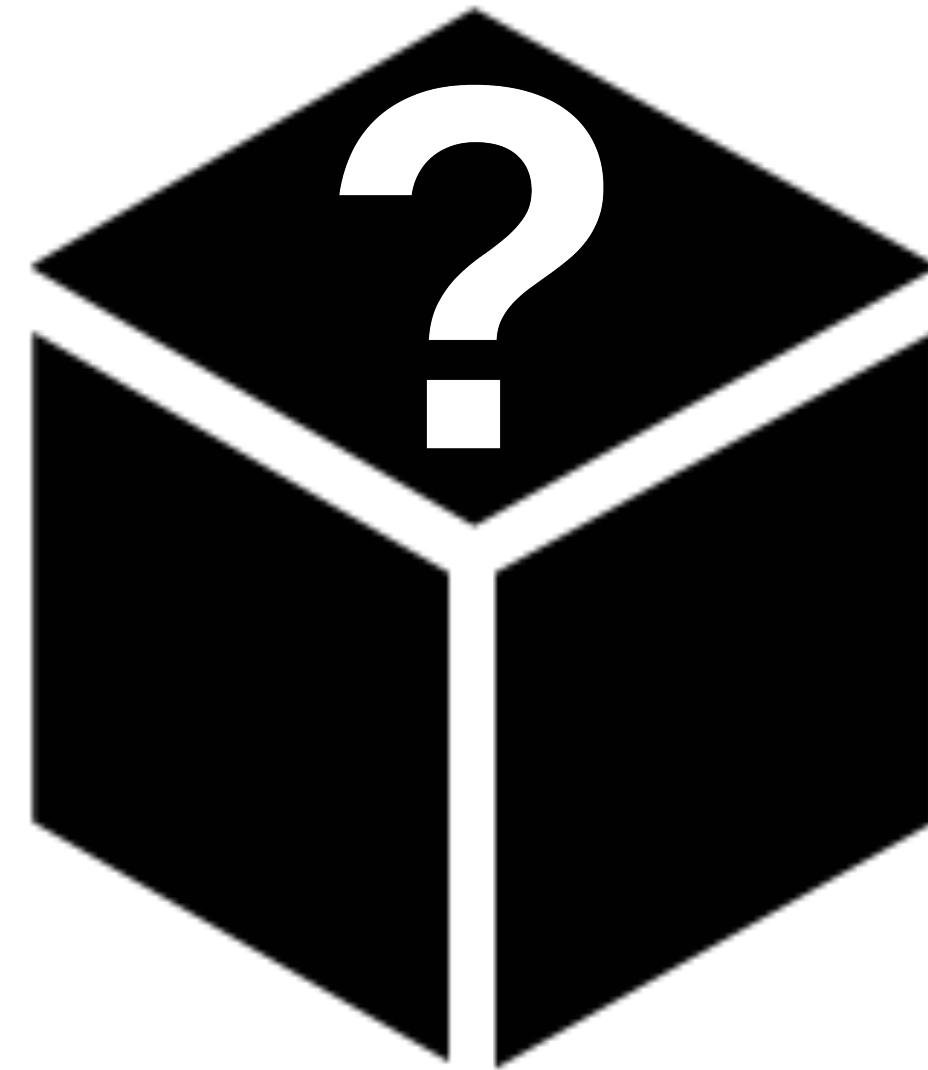
tiering



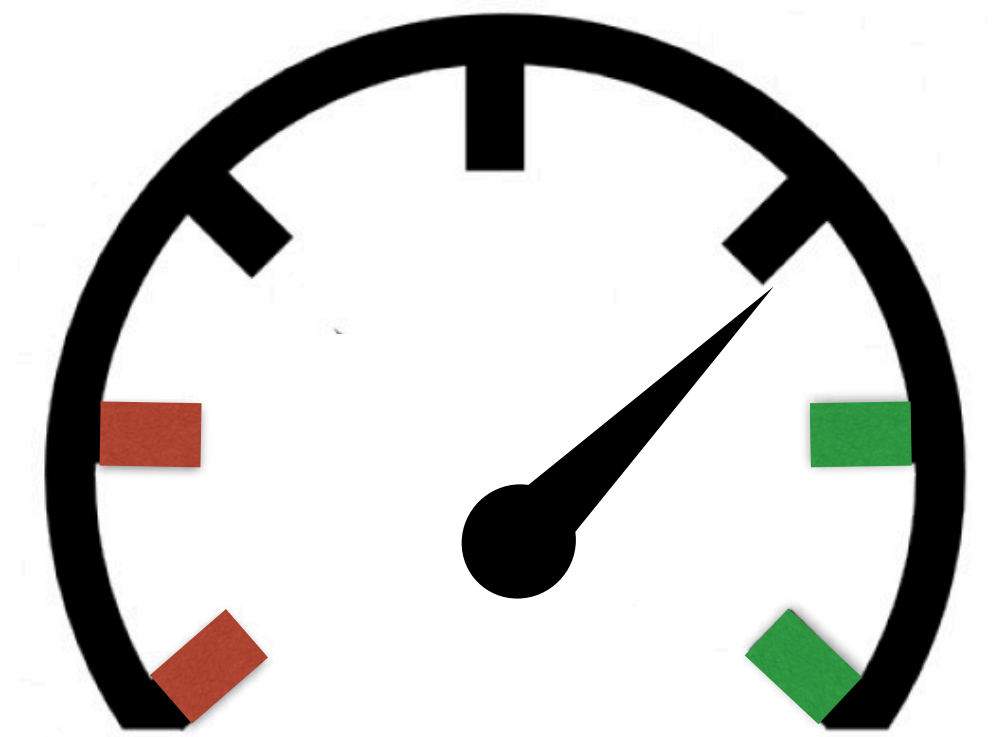
So, how do we reason about the data layout?



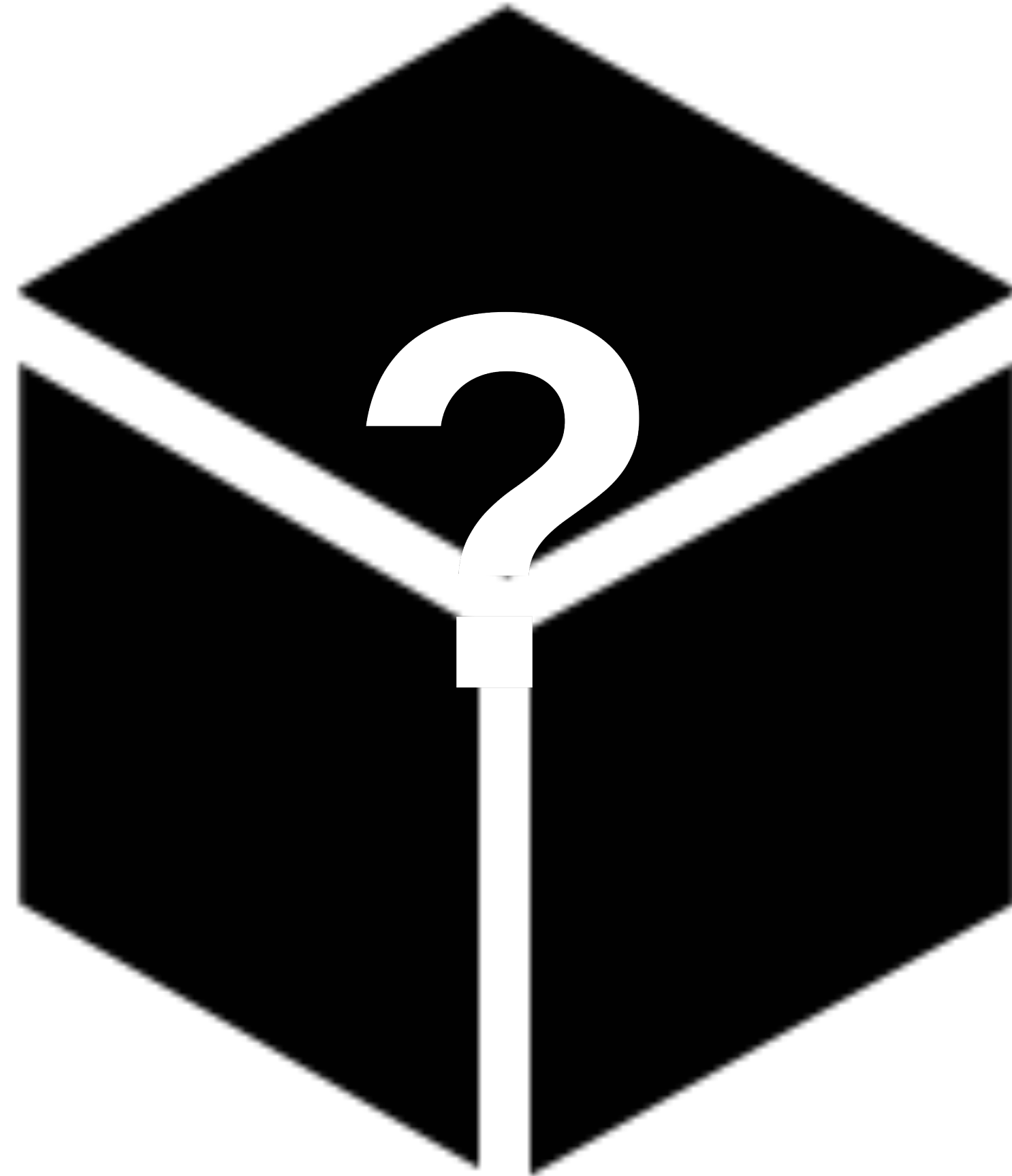
workload



**data layout**



performance



Compaction  
black box

1

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

2

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

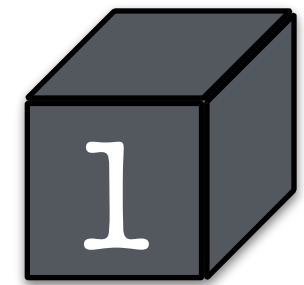
3

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

4

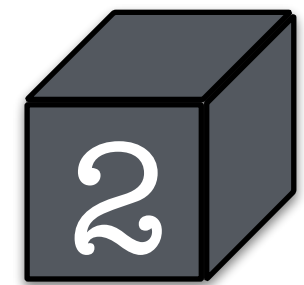
**When** to re-organize the data layout?

Data Layout



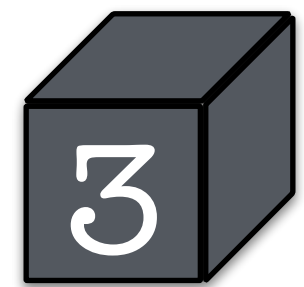
**How** to organize the data on device? ✓

Compaction  
Granularity



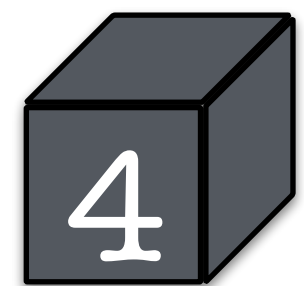
**How much** data to move at-a-time?

Data Movement  
Policy

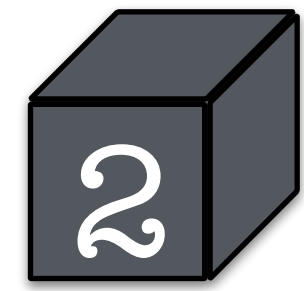


**Which** block of data to be moved?

Compaction  
Trigger



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

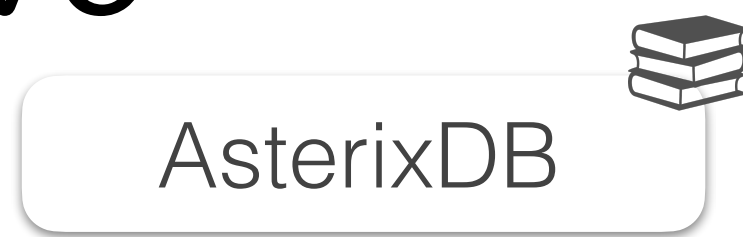


# Compaction **Granularity**

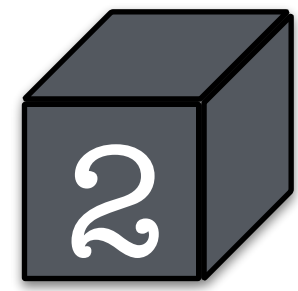
*data moved per compaction*



consecutive  
levels

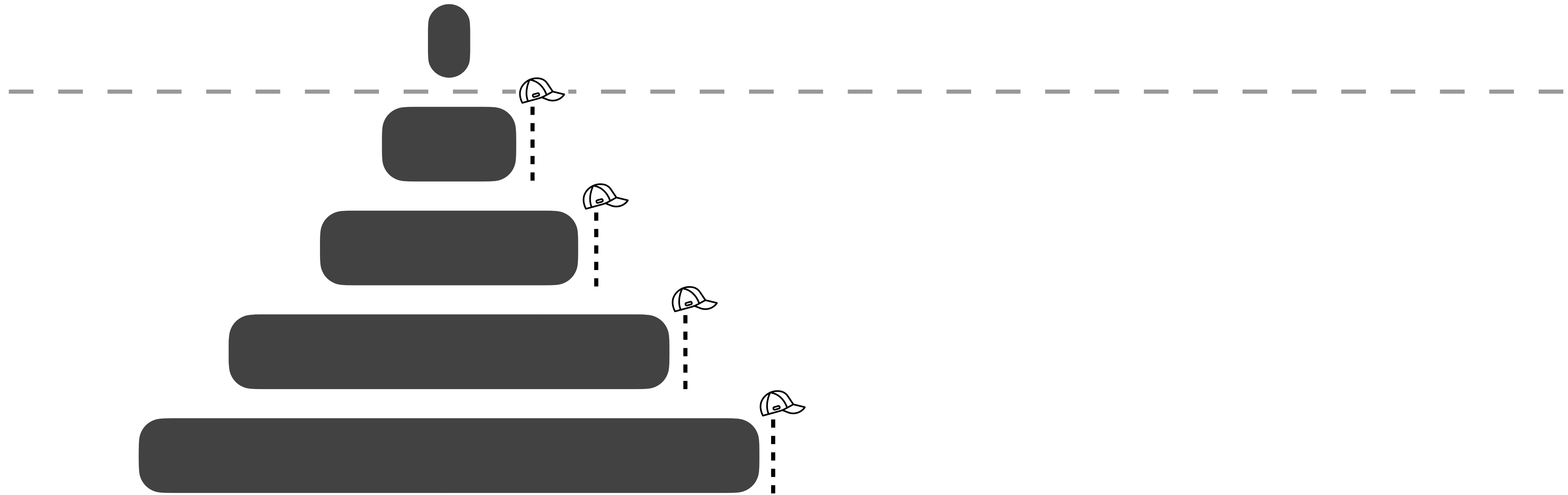


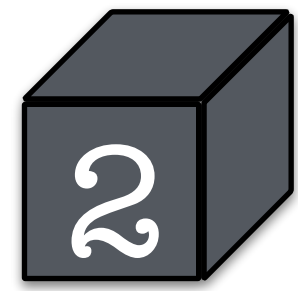




# Compaction **Granularity**

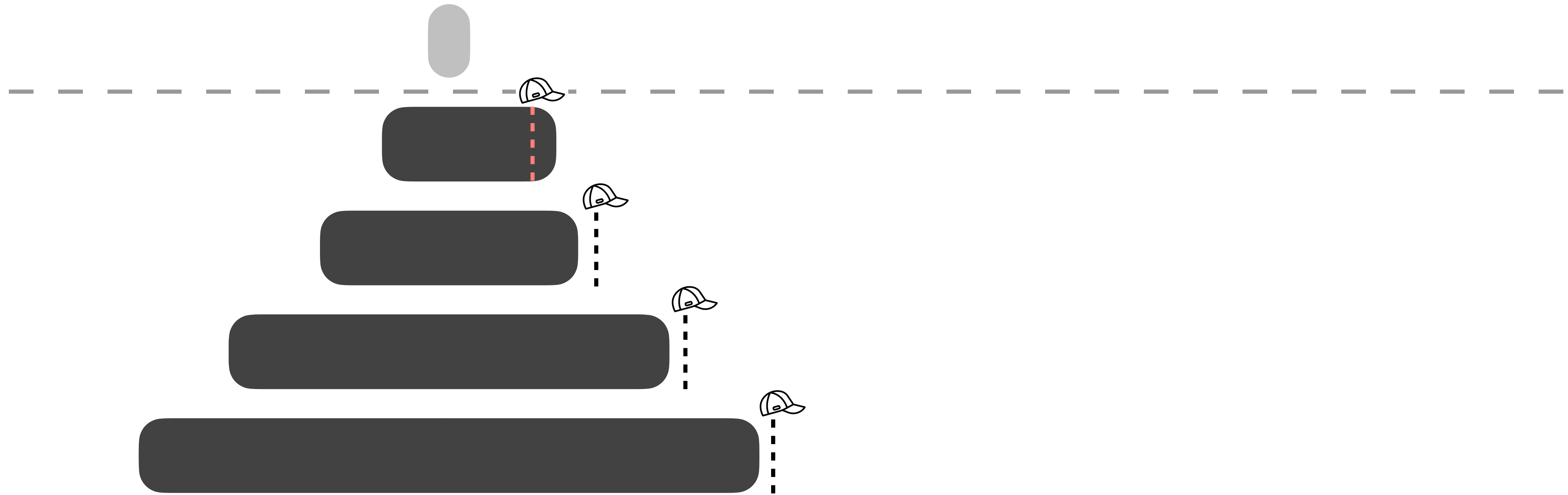
*data moved per compaction*

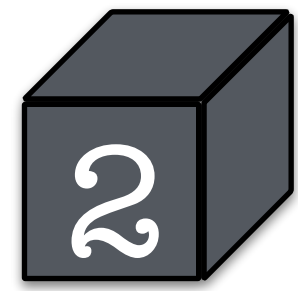




# Compaction **Granularity**

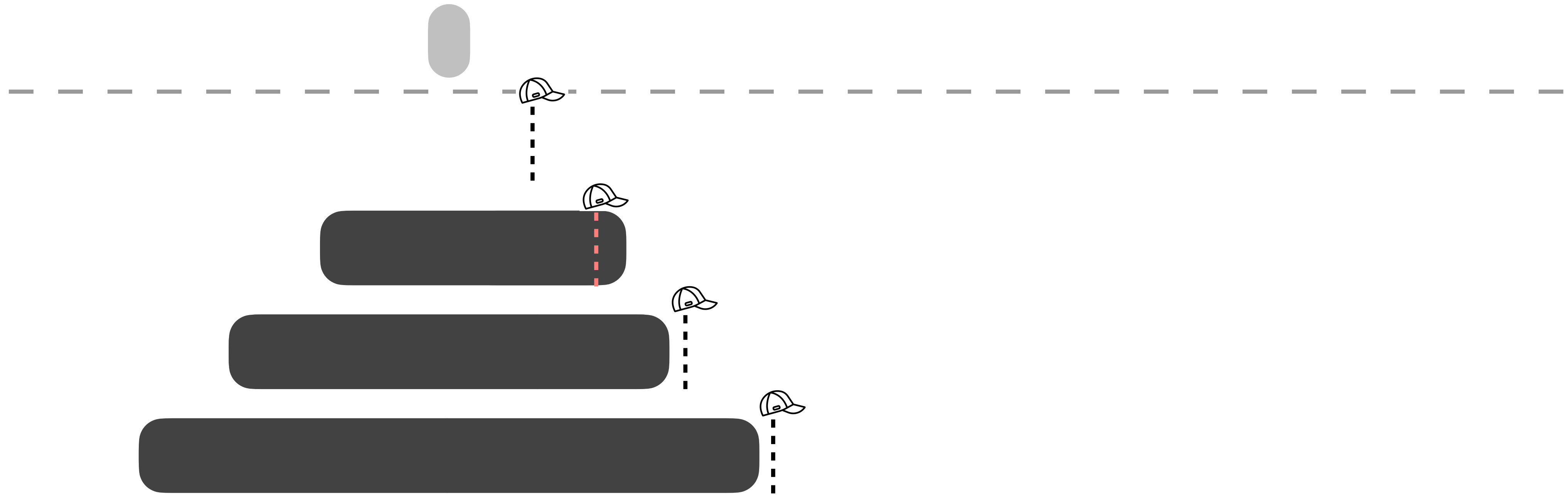
*data moved per compaction*

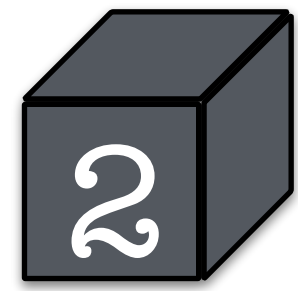




# Compaction **Granularity**

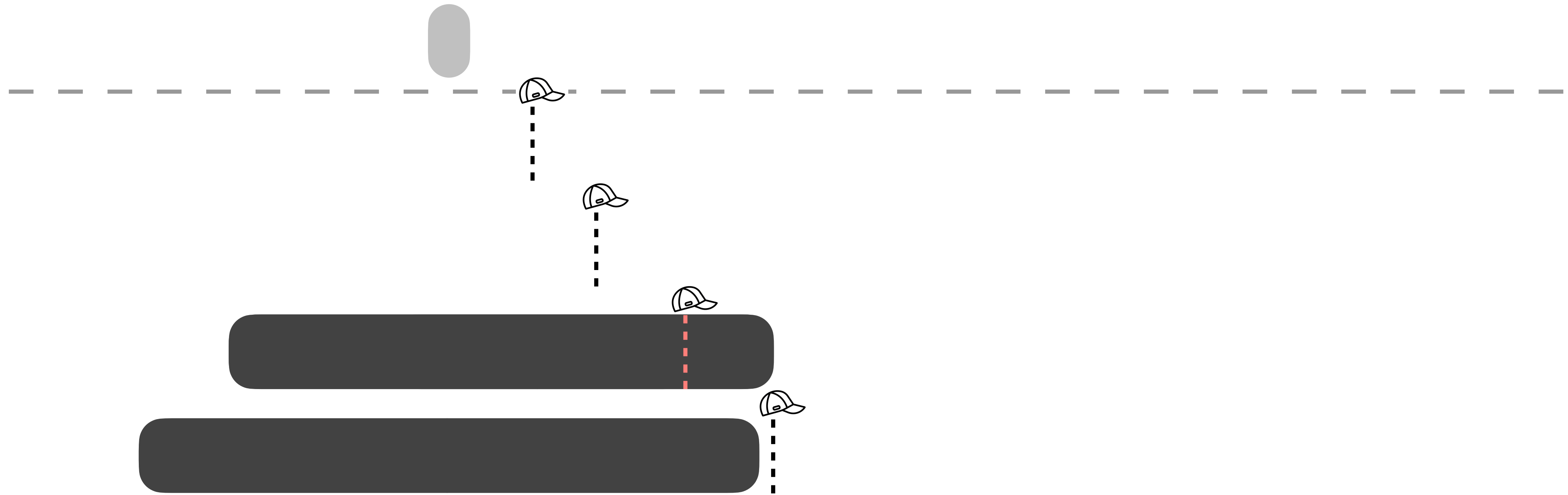
*data moved per compaction*

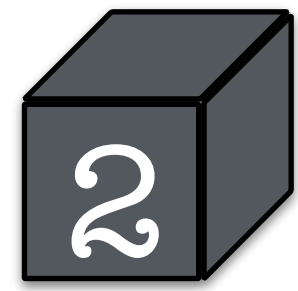




# Compaction **Granularity**

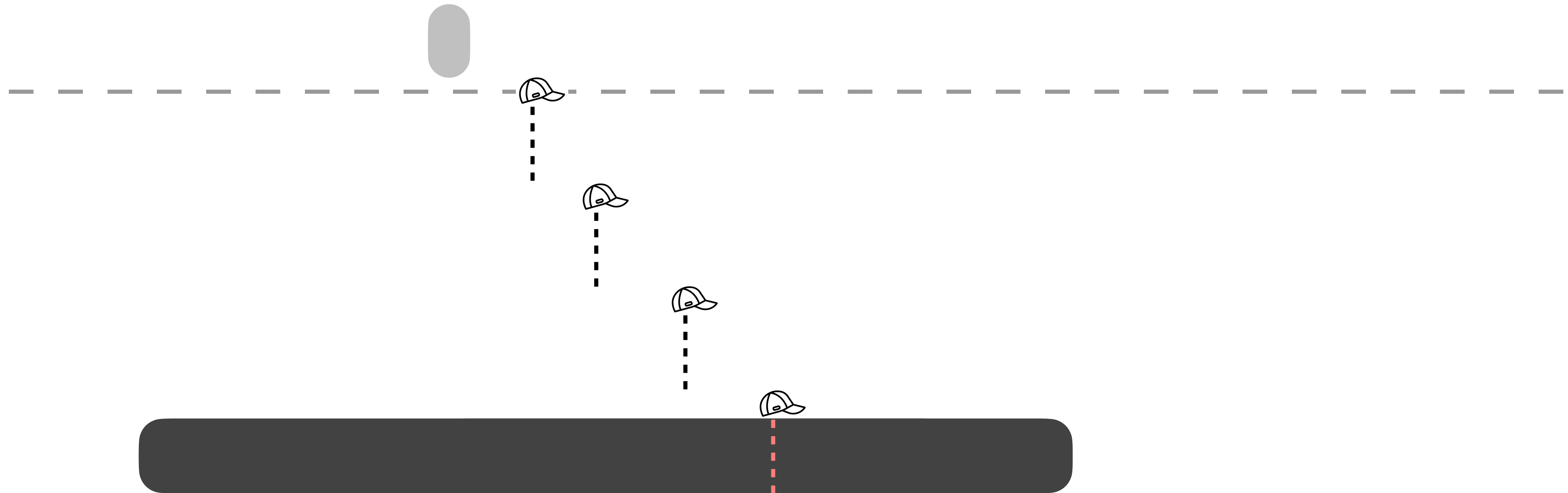
*data moved per compaction*

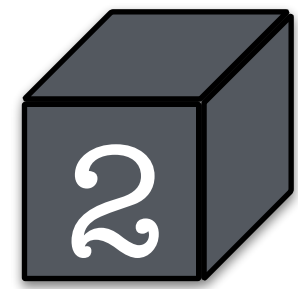




# Compaction **Granularity**

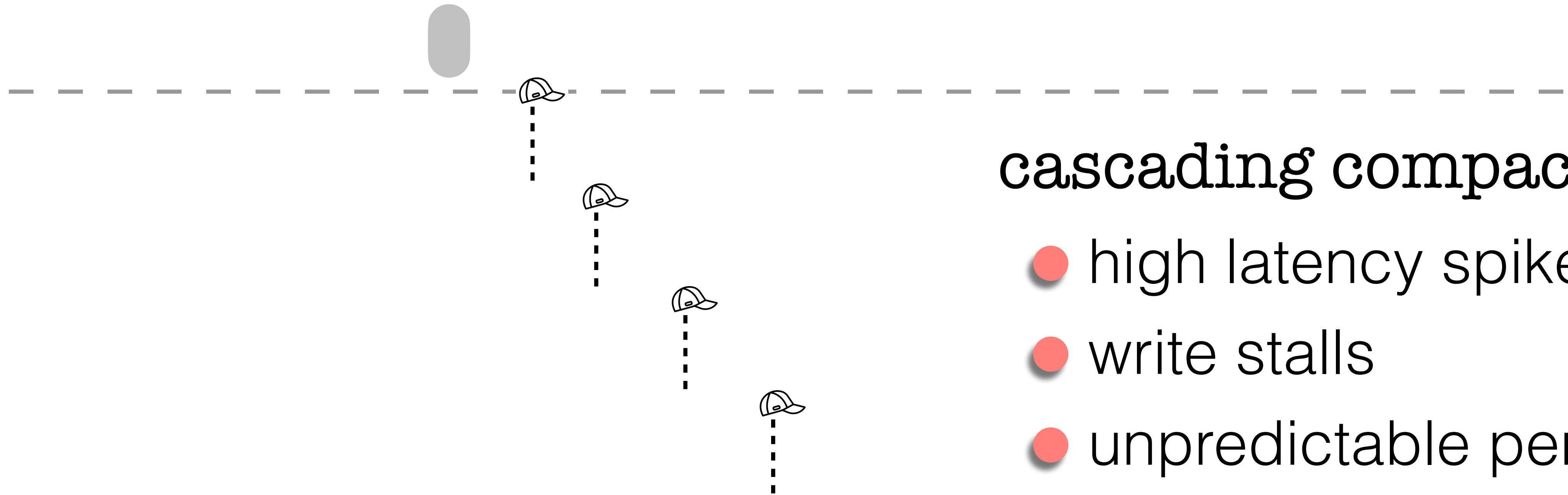
*data moved per compaction*





# Compaction **Granularity**

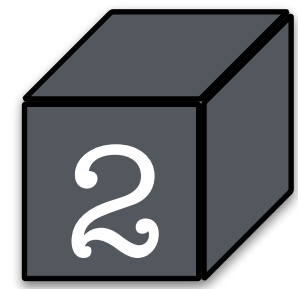
*data moved per compaction*



**cascading compactions**

- high latency spikes
- write stalls
- unpredictable perf.



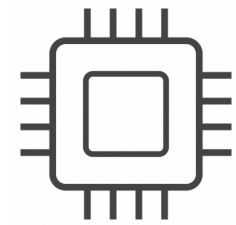


# Compaction **Granularity**

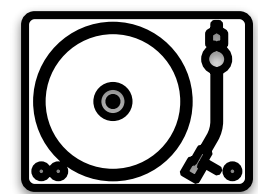
*data moved per compaction*

partial compaction

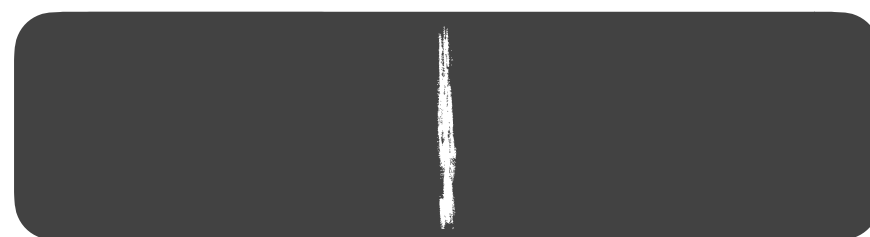
granularity: files



buffer

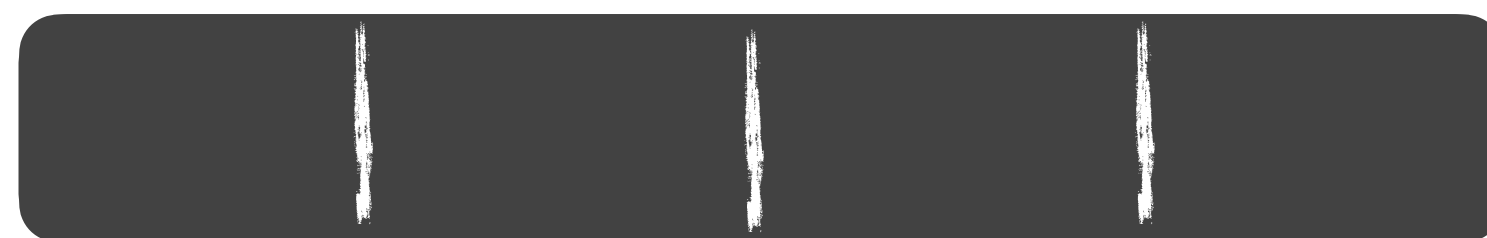


level 1

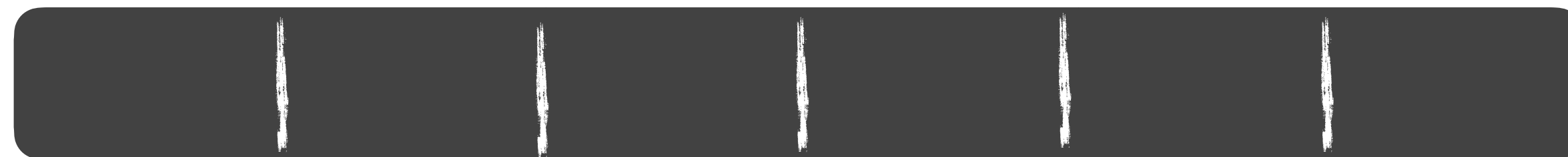


partial compaction

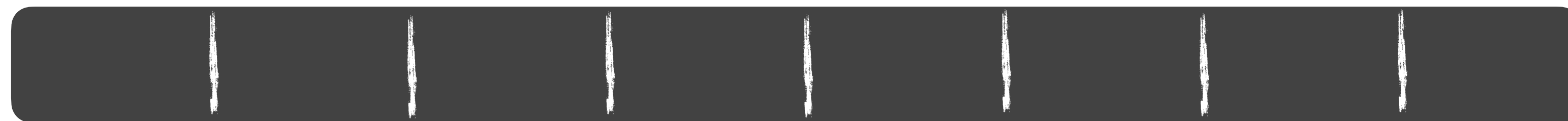
level 2



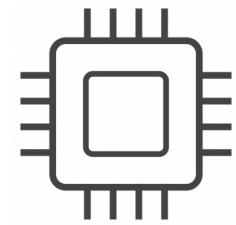
level 3



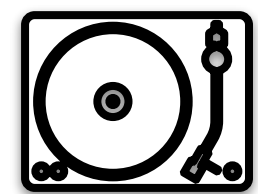
level 4



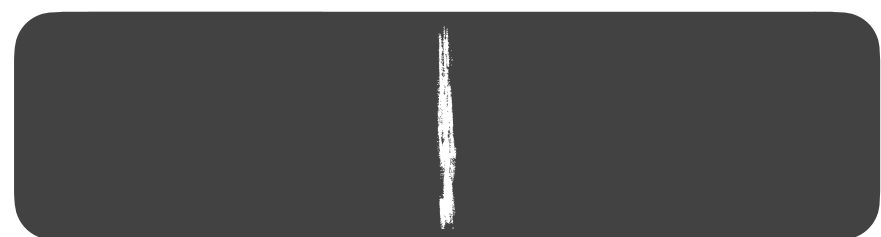




buffer

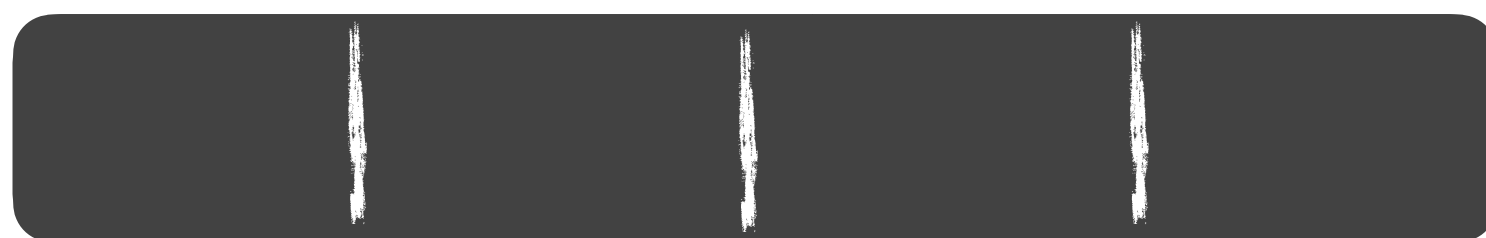


level 1



partial compaction

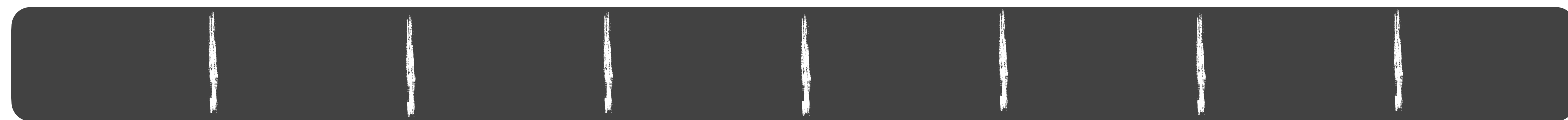
level 2

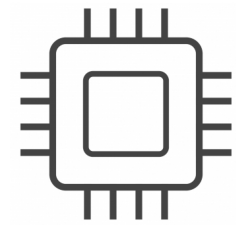


level 3

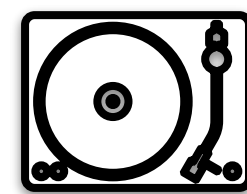


level 4

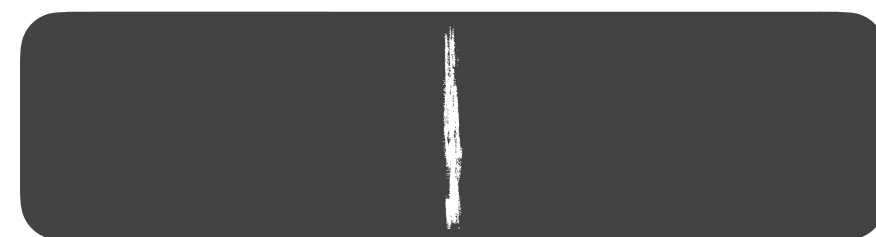




buffer

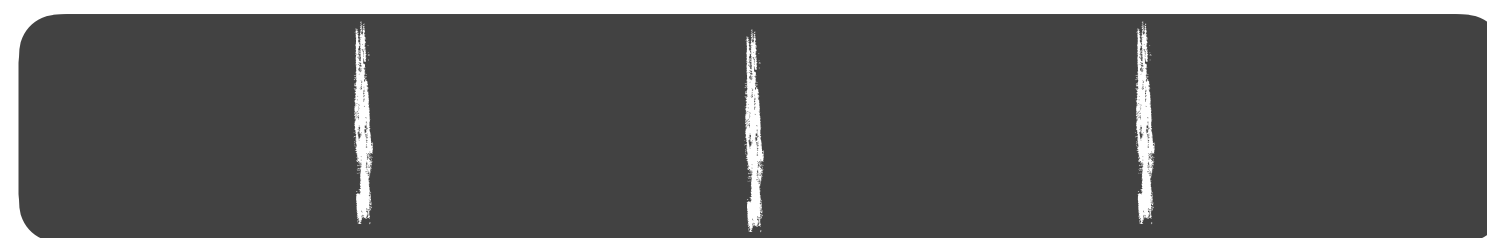


level 1



partial compaction

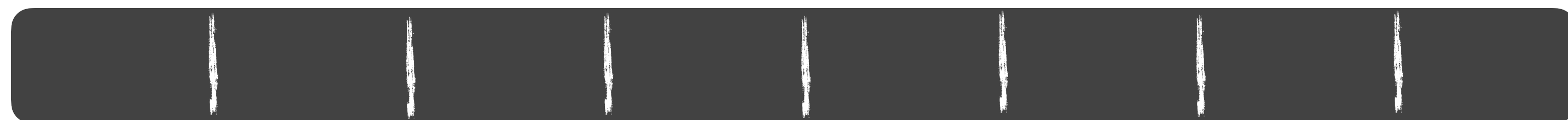
level 2

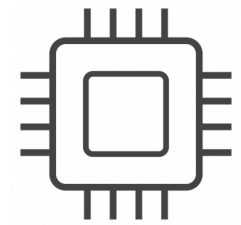


level 3

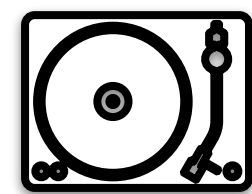


level 4





buffer



level 1



partial compaction

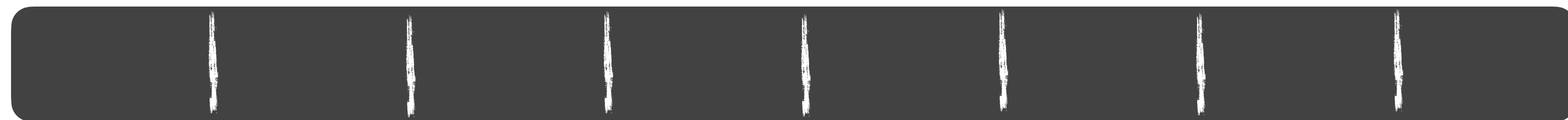
level 2

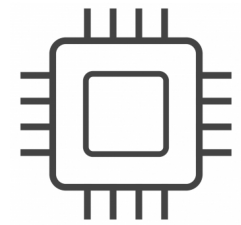


level 3

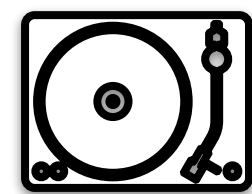


level 4

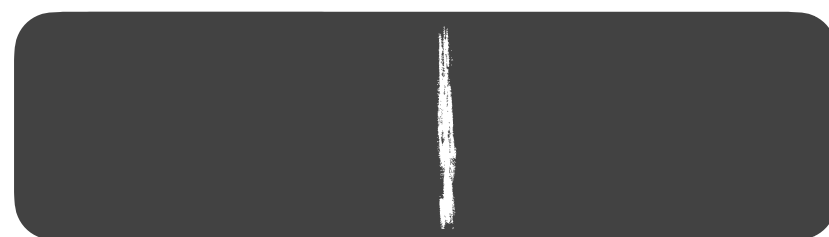




buffer



level 1



partial compaction

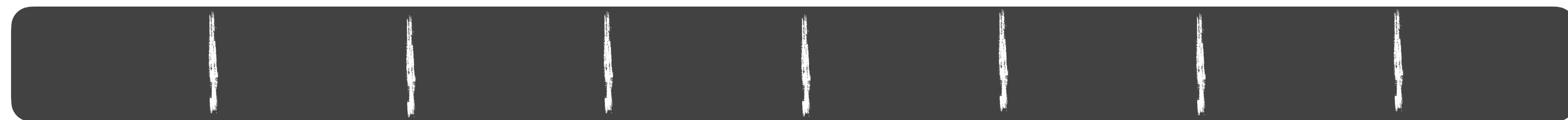
level 2

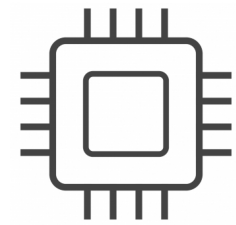


level 3

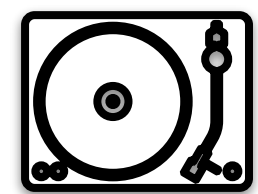


level 4

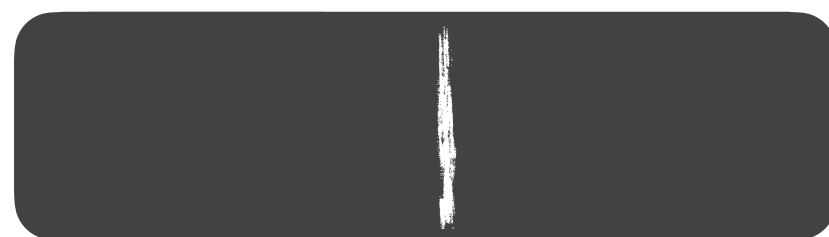




buffer



level 1



partial compaction

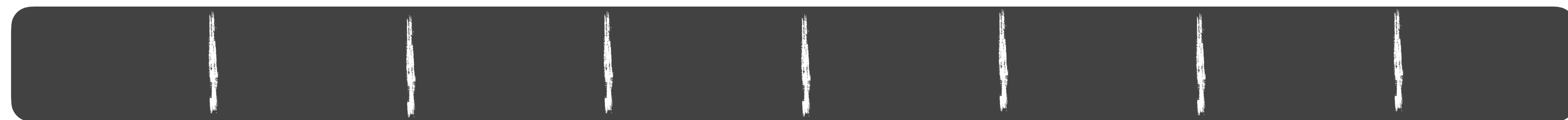
level 2

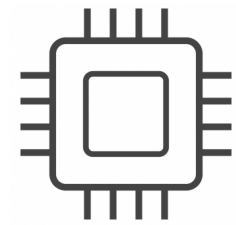


level 3

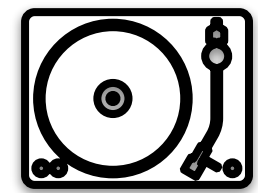


level 4

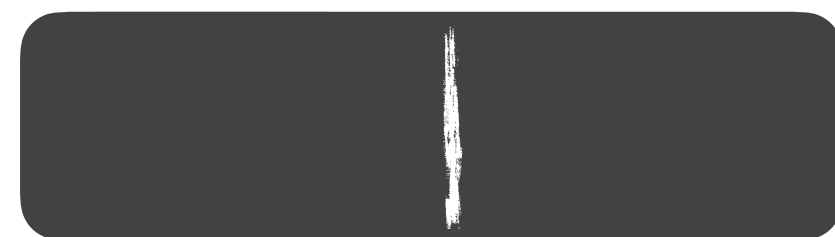




buffer

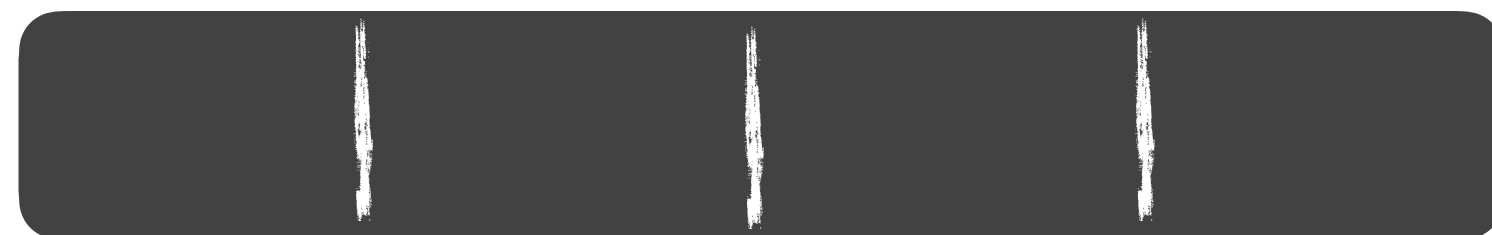


level 1

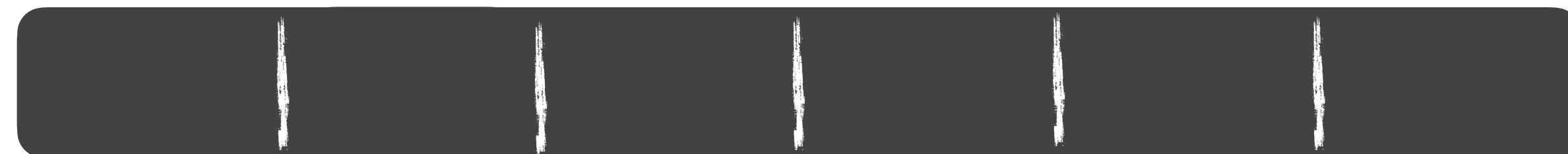


partial compaction

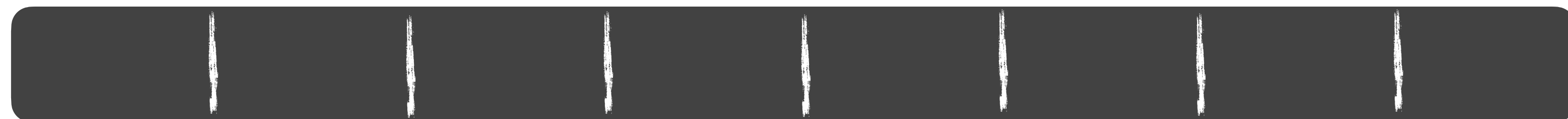
level 2

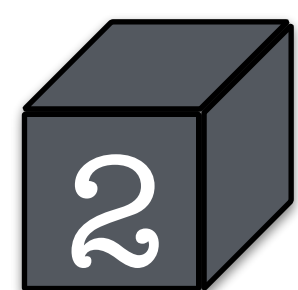


level 3



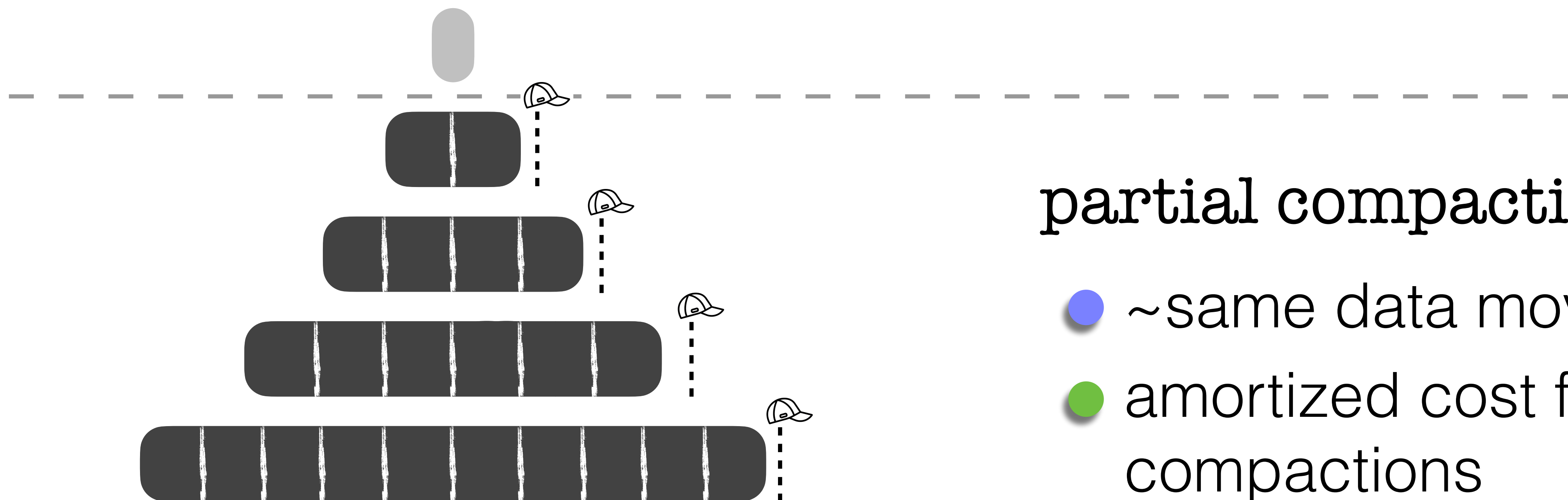
level 4





# Compaction Granularity

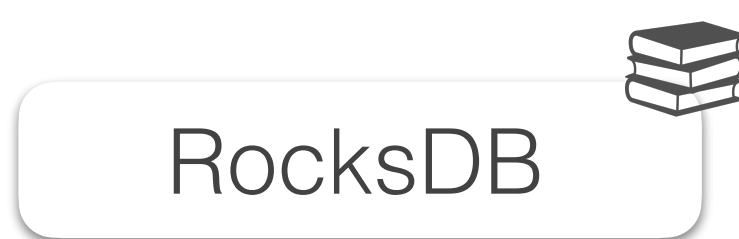
*data moved per compaction*

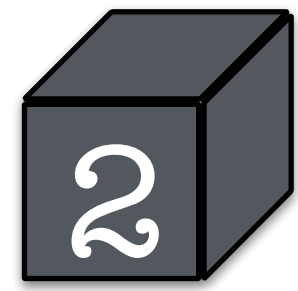


partial compaction

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

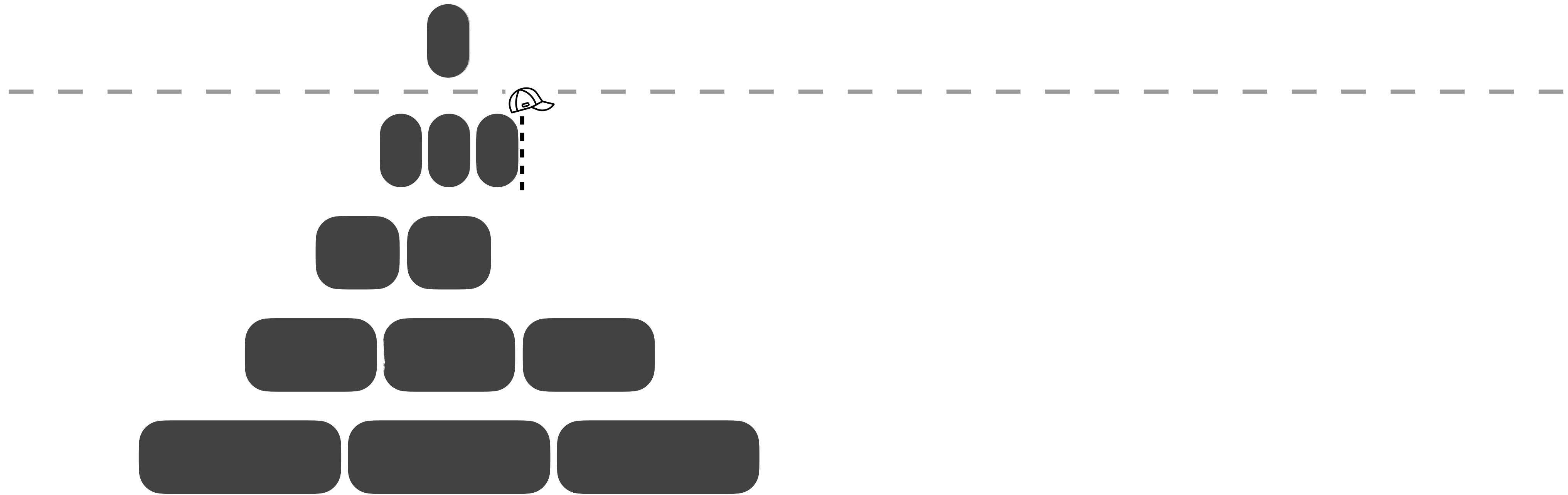
files





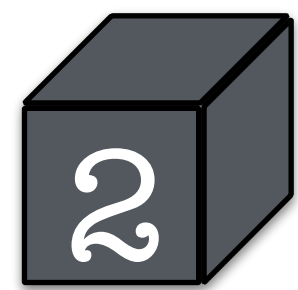
# Compaction **Granularity**

*data moved per compaction*



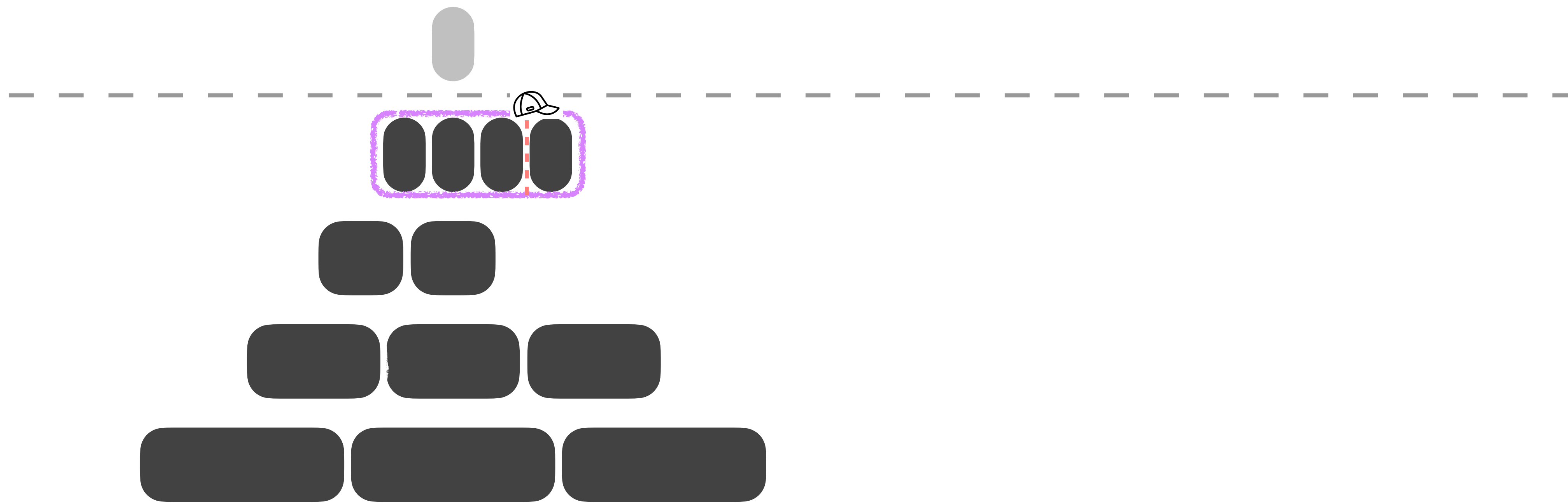
sorted runs in a level



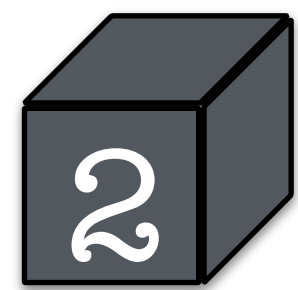


# Compaction **Granularity**

*data moved per compaction*

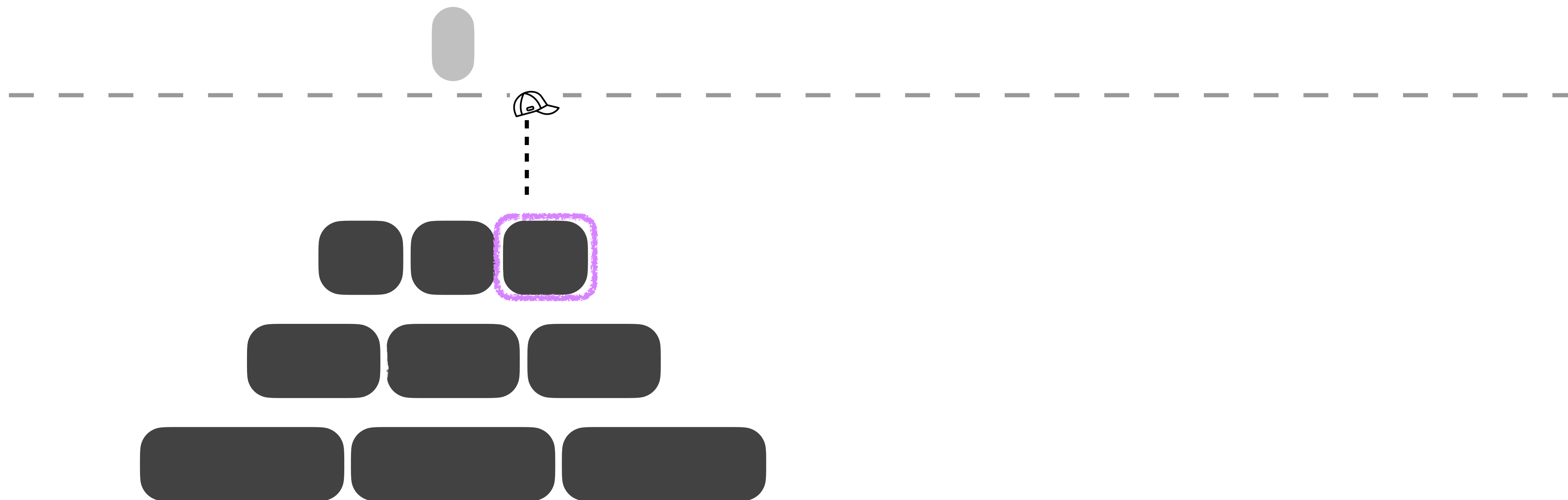


sorted runs in a level

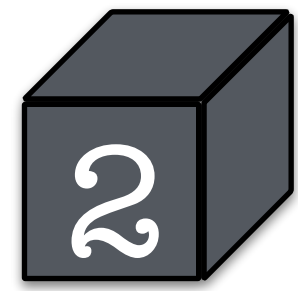


# Compaction **Granularity**

*data moved per compaction*

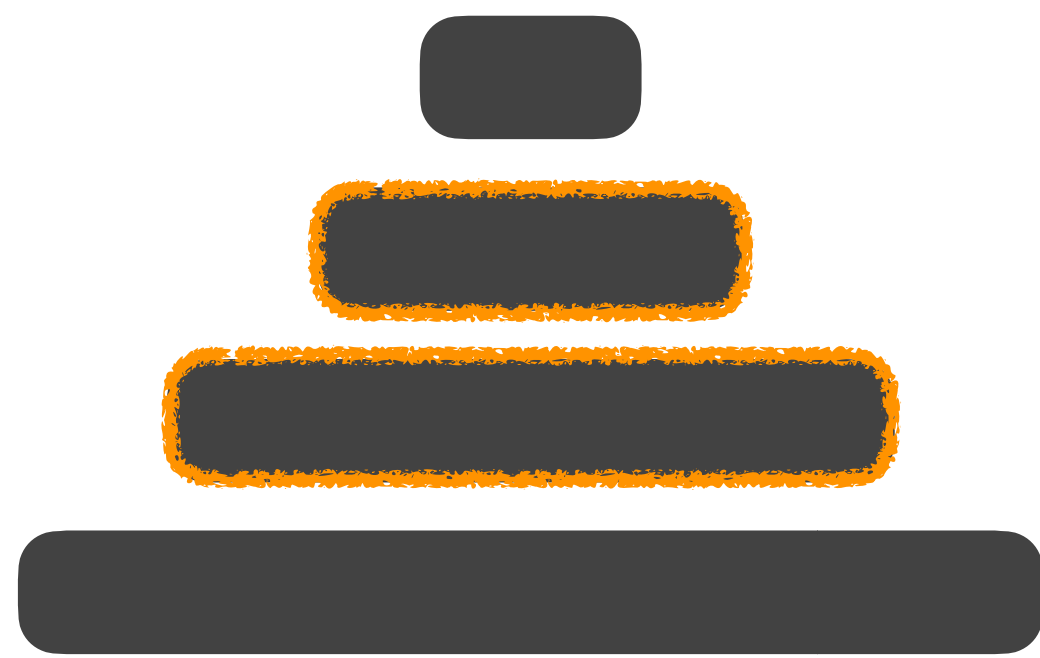


sorted runs in a level

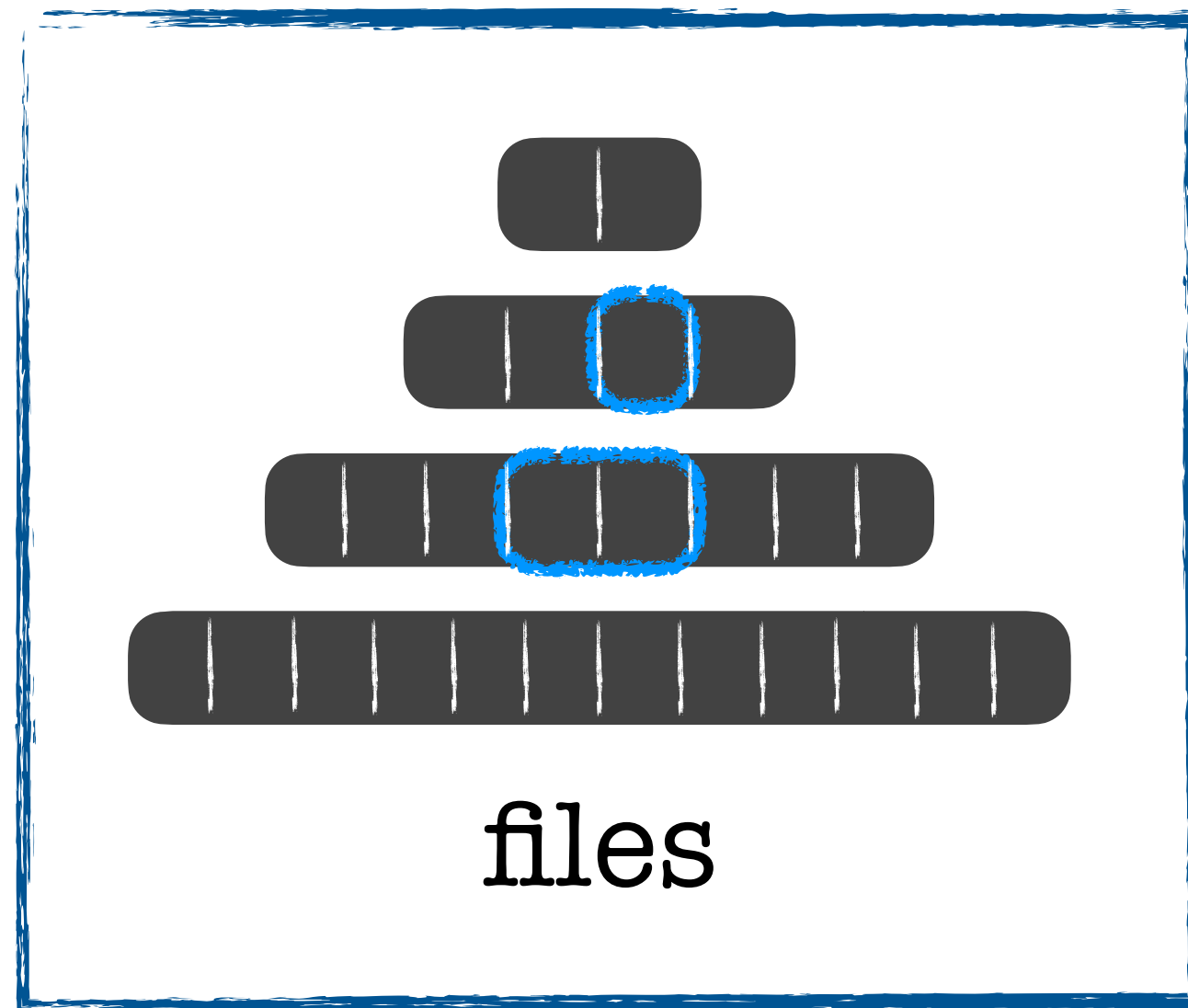


# Compaction Granularity

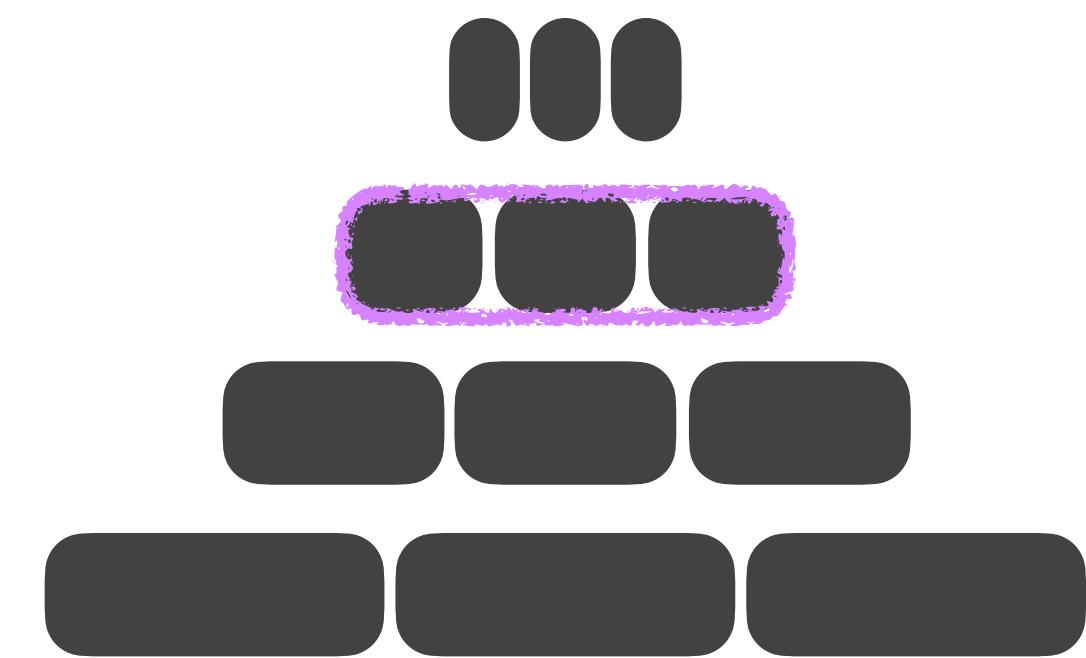
*data moved per compaction*



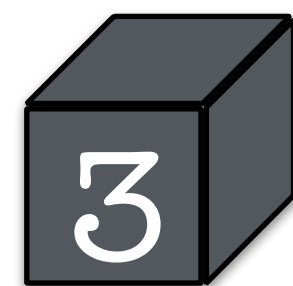
levels



files

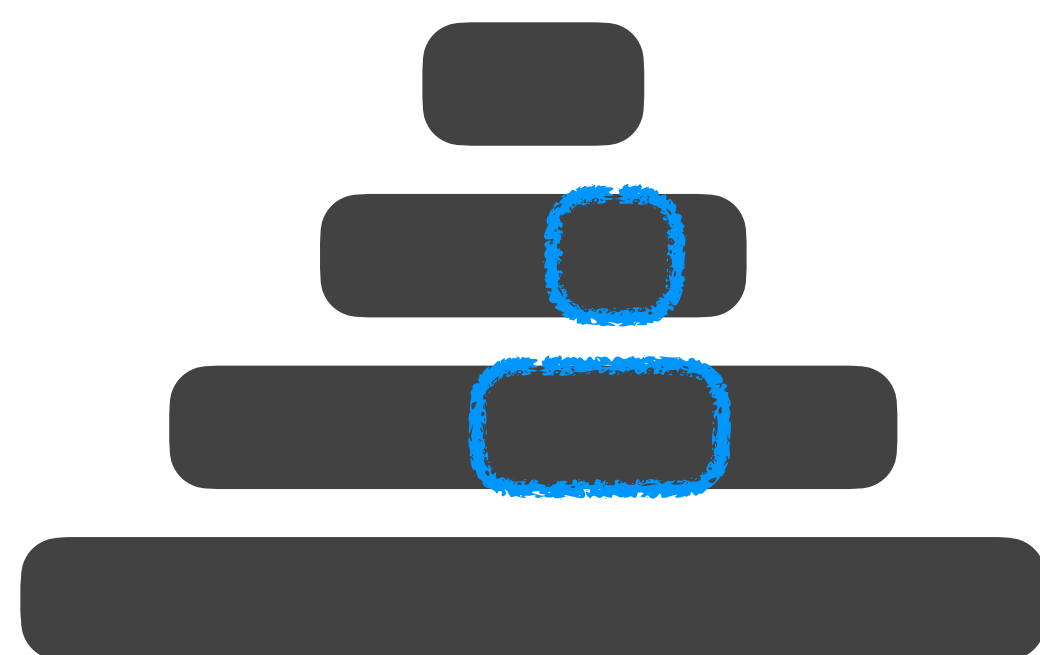


sorted runs in a level

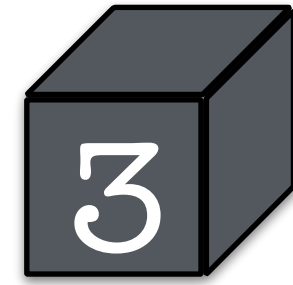


# Data Movement Policy

*which data to compact*

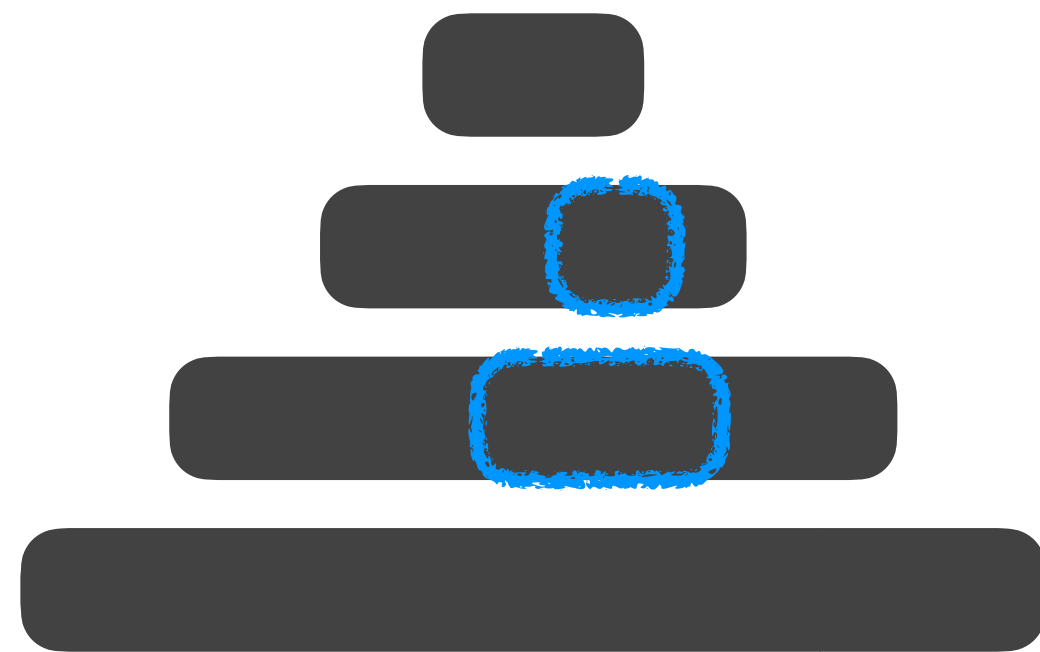


files



# Data Movement Policy

*which data to compact*



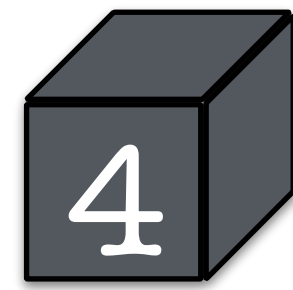
files

**round-robin** 

minimum **overlap with parent** level 

file with most **tombstones** 

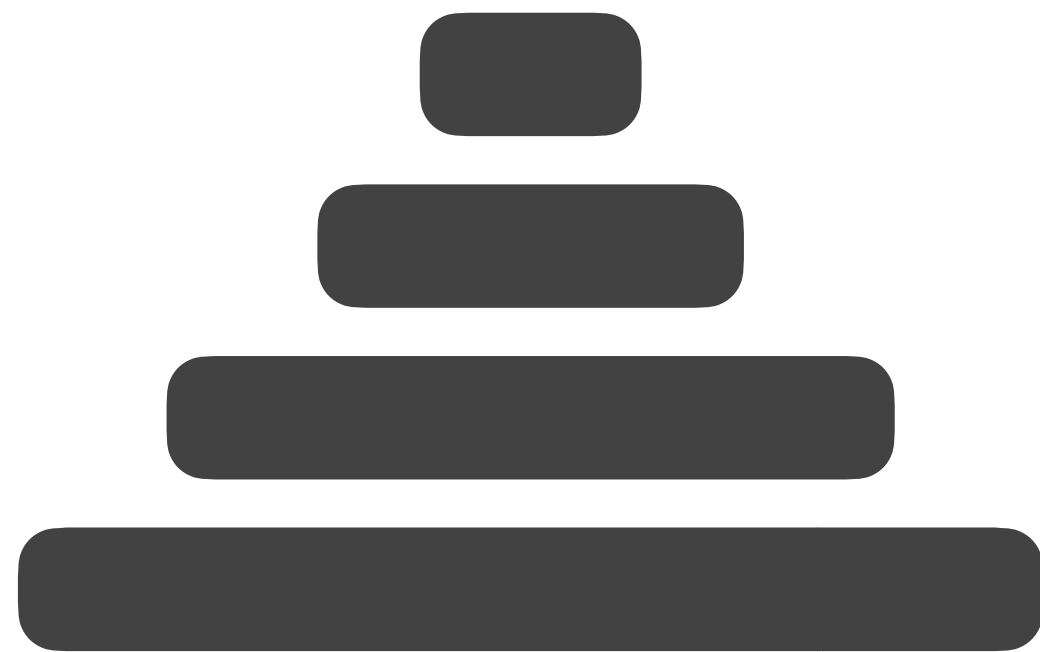
**coldest** file 

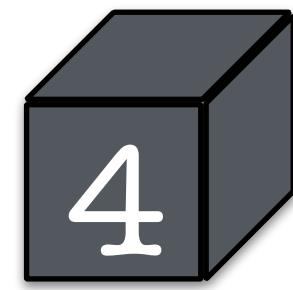


# Compaction **Trigger**

*invoking the compaction routine*

level **saturation**

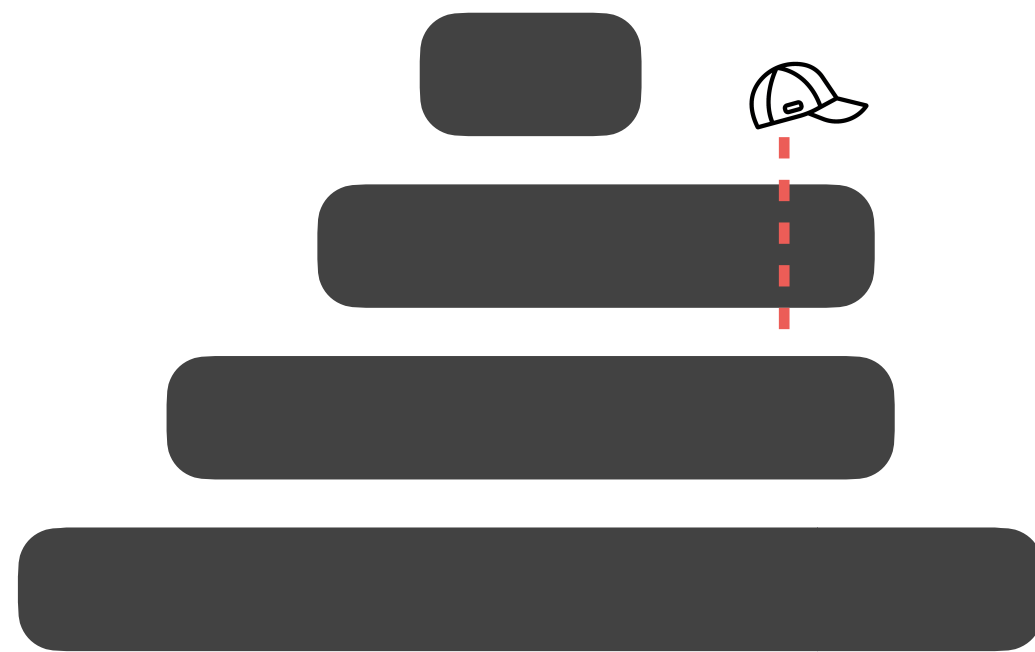


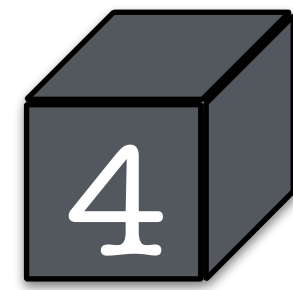


# Compaction **Trigger**

*invoking the compaction routine*

level **saturation**

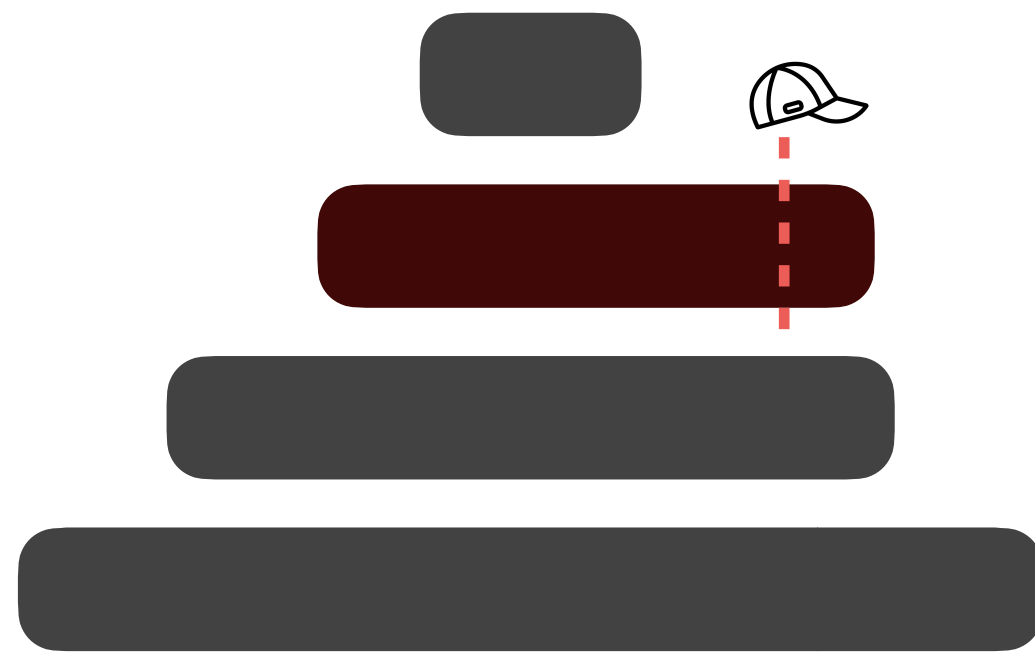




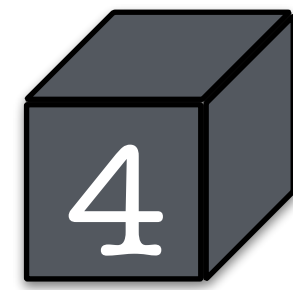
# Compaction **Trigger**

*invoking the compaction routine*

level **saturation**

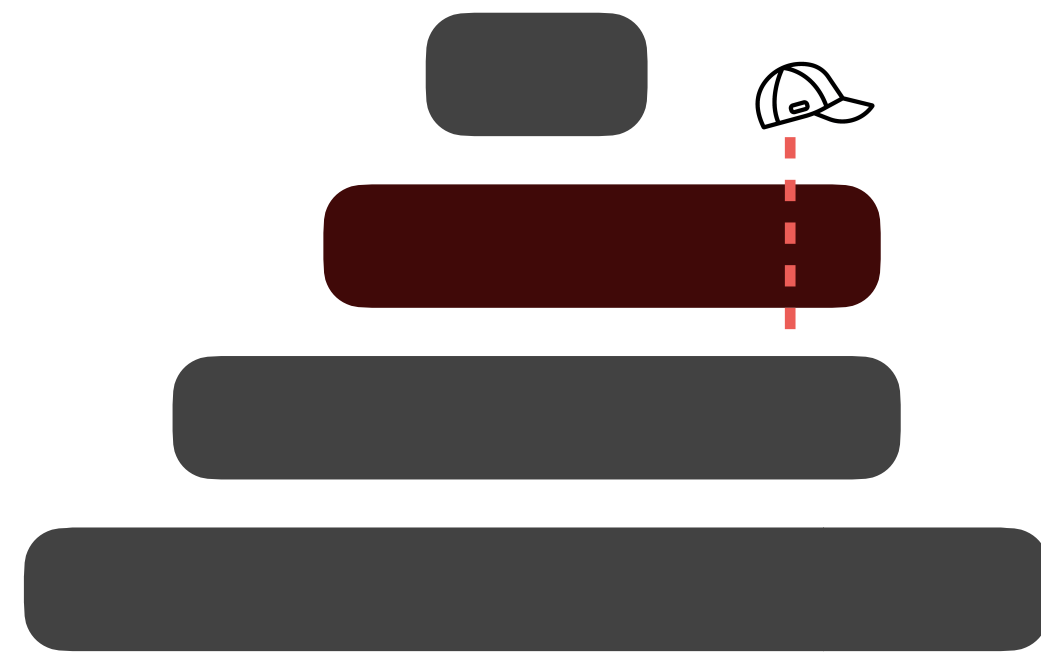






# Compaction **Trigger**

*invoking the compaction routine*



level **saturation**

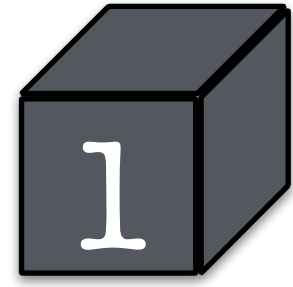
number of **sorted runs**

**space amplification**

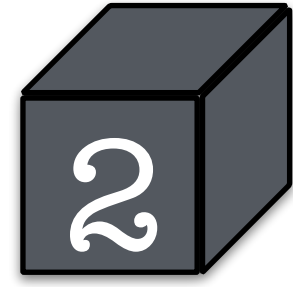


**age** of a file

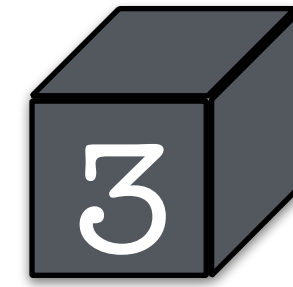




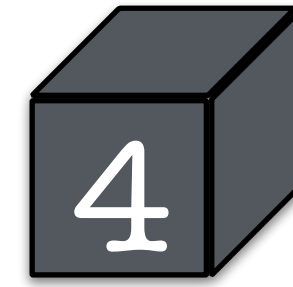
Data Layout



Compaction  
Granularity



Data Movement  
Policy



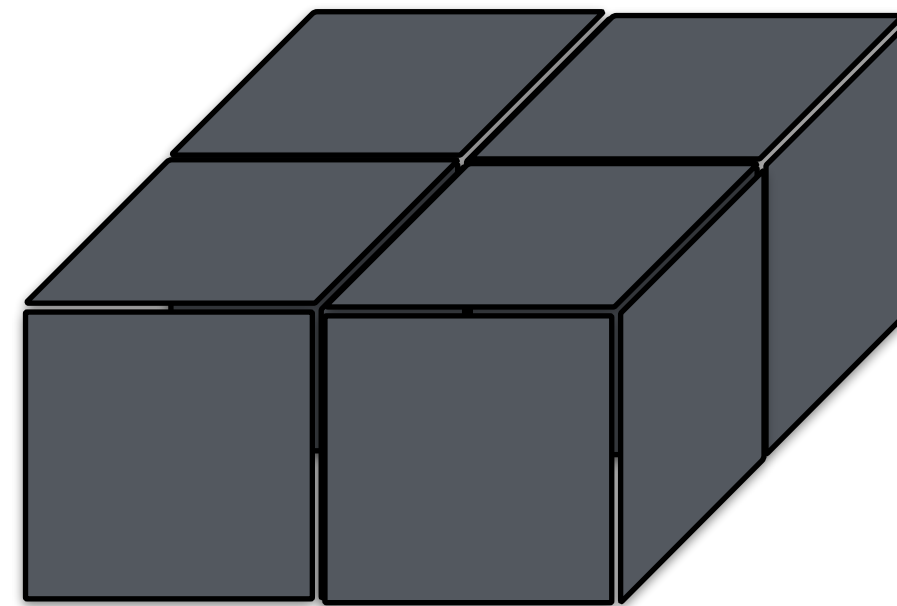
Compaction  
Trigger

Data Layout

Compaction  
Granularity

Data Movement  
Policy

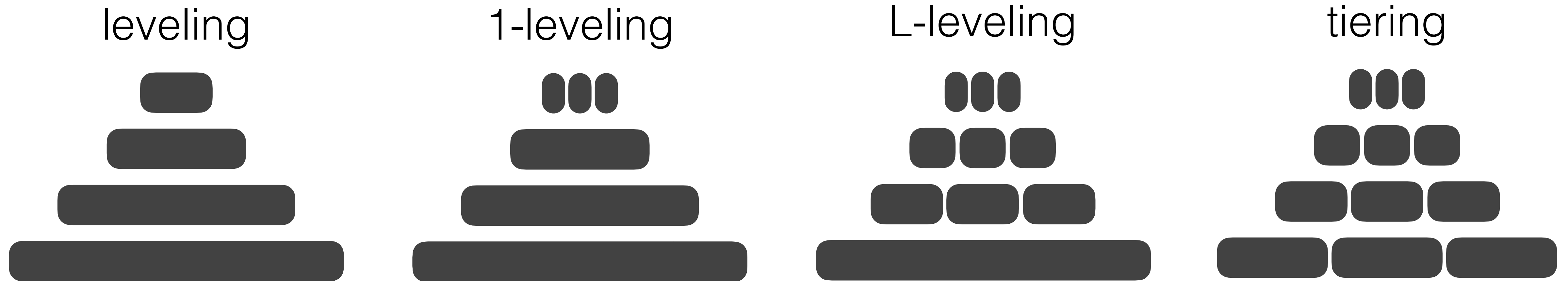
Compaction  
Trigger



***Any* Compaction Algorithm**

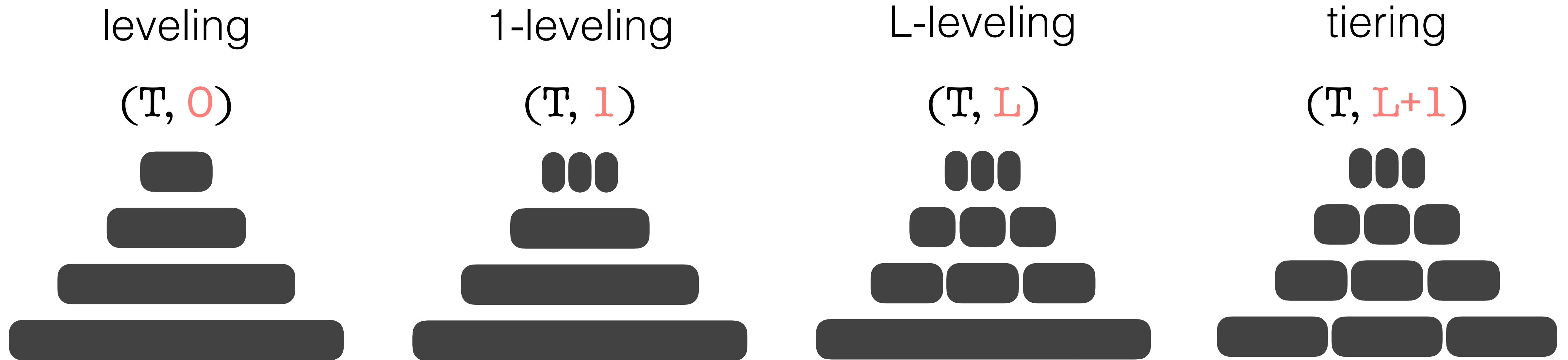
Database	Data layout	Compaction Trigger					Compaction Granularity				Data Movement Policy						
		Level saturation	#Sorted runs	File staleness	Space amp.	Tombstone-TTL	Level	Sorted run	File (single)	File (multiple)	Round-robin	Least overlap (+1)	Least overlap (+2)	Coldest file	Oldest file	Tombstone 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	✓ <sup>L</sup>	✓ <sup>T</sup>				✓ <sup>L</sup>	✓ <sup>T</sup>			✓ <sup>L</sup>						✓ <sup>T</sup>
LSM-Bush [24]	Hybrid leveling	✓ <sup>L</sup>	✓ <sup>T</sup>				✓ <sup>L</sup>	✓ <sup>T</sup>			✓ <sup>L</sup>						✓ <sup>T</sup>
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		✓					✓									✓

# Storage Layer **Design Continuum**



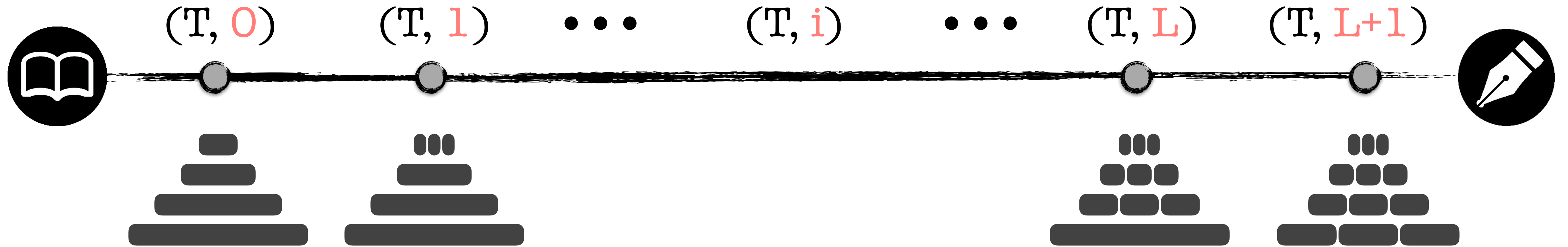
Any design can be defined by the tuple-set:  $(T, i)$

# Storage Layer **Design Continuum**

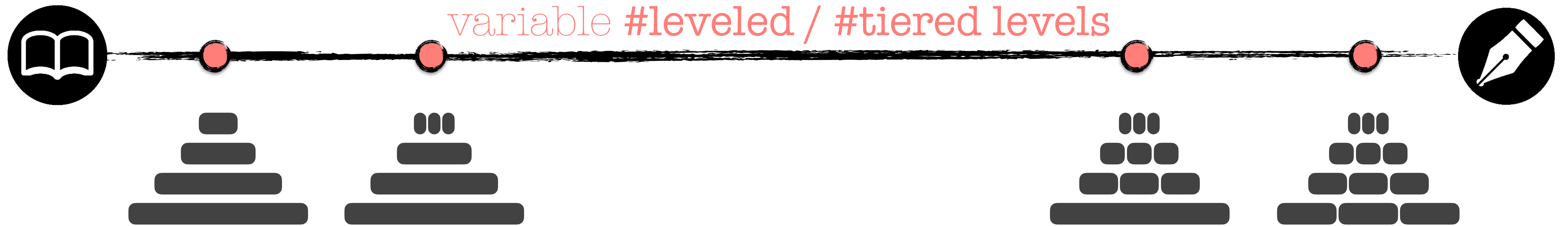


Any design can be defined by the tuple-set:  $(T, i)$

# Storage Layer Design Continuum

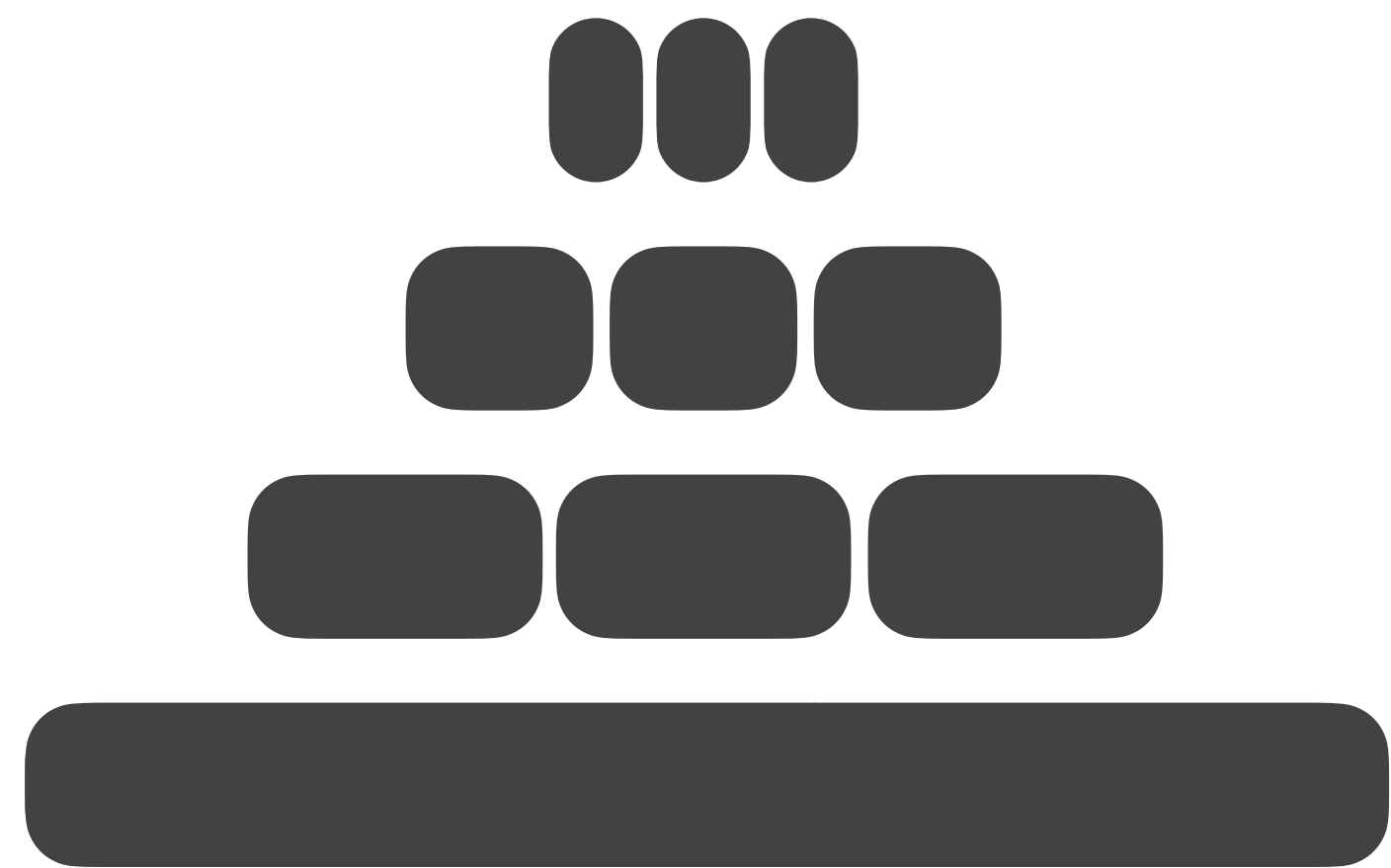
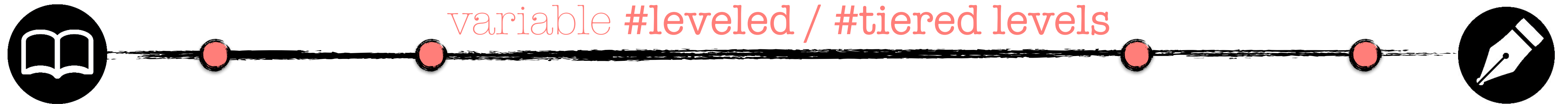


# Storage Layer Design Continuum





# Storage Layer Design Continuum



size ratio

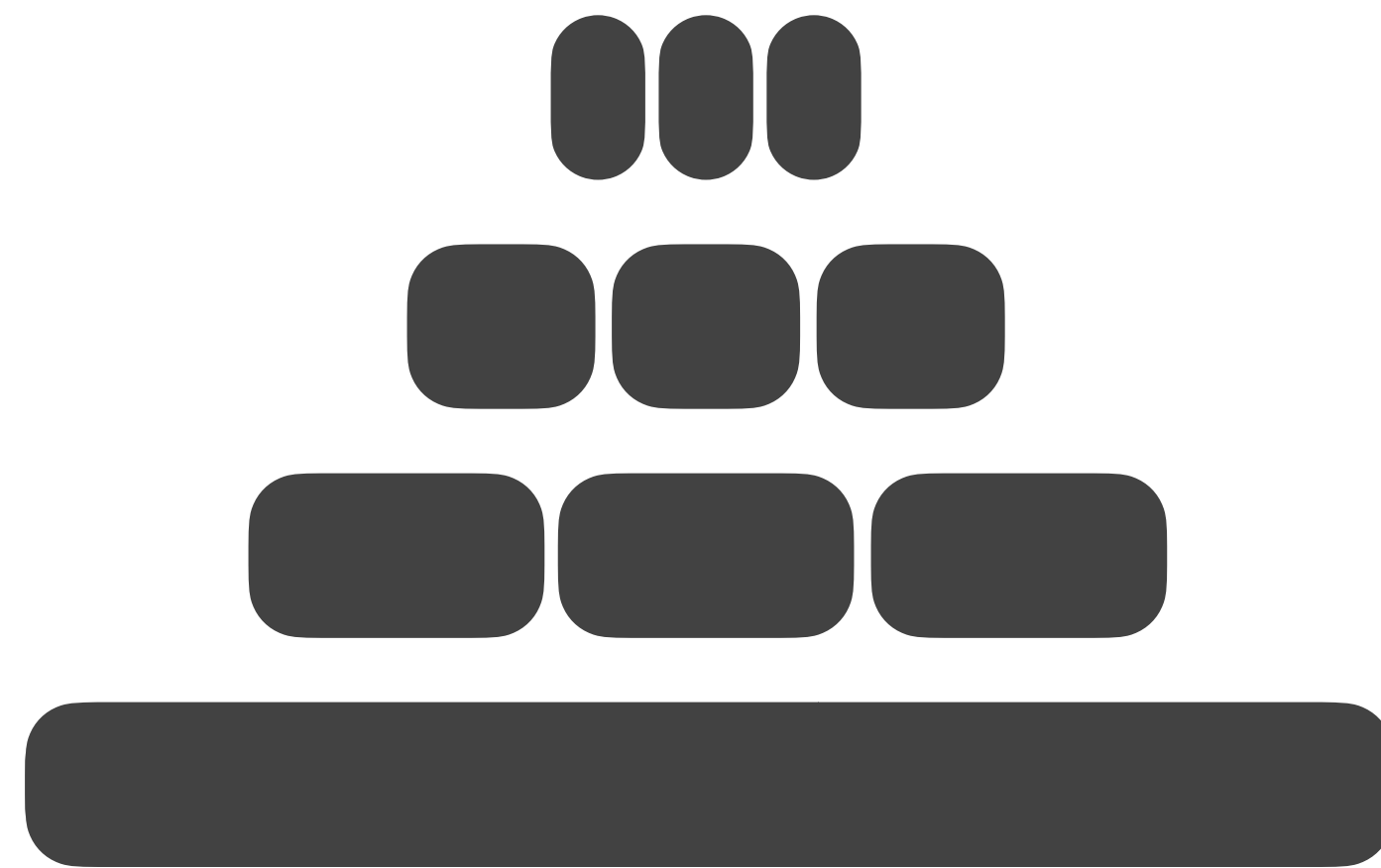
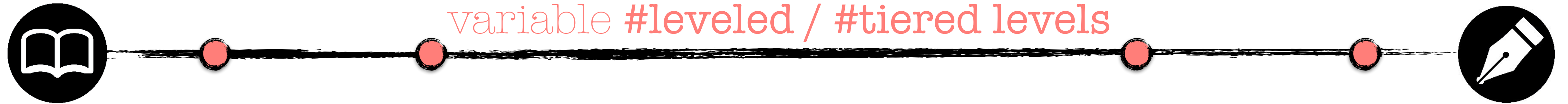
T

T

T

T

# Storage Layer Design Continuum



size ratio

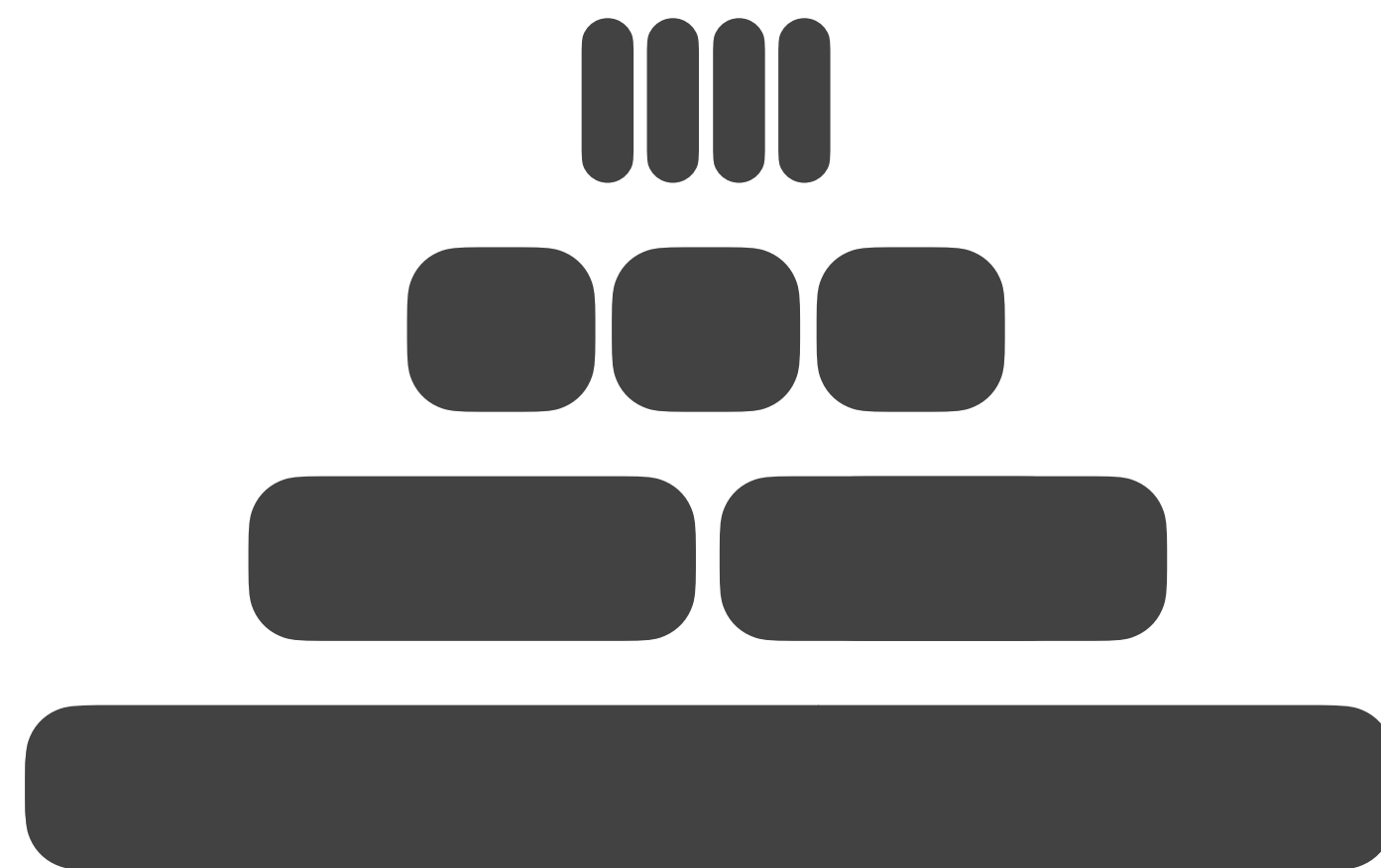
T

T

T

T

# Storage Layer Design Continuum



size ratio

#runs

T

4

T

3

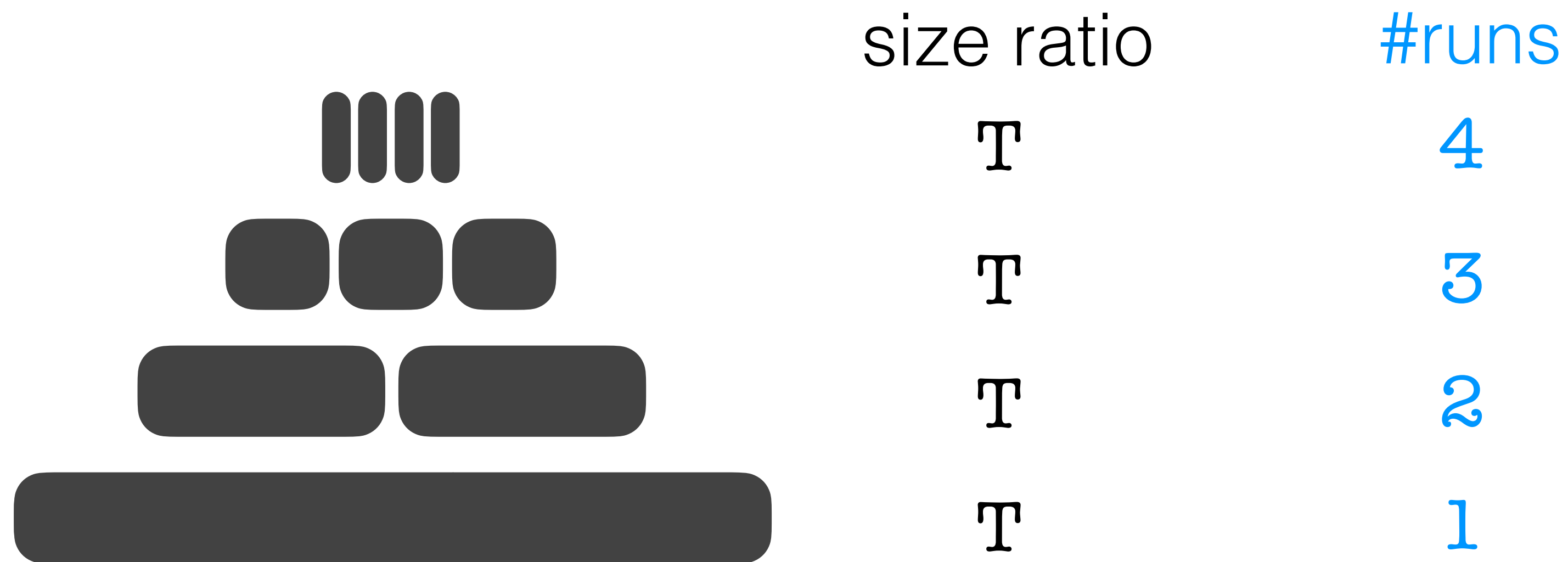
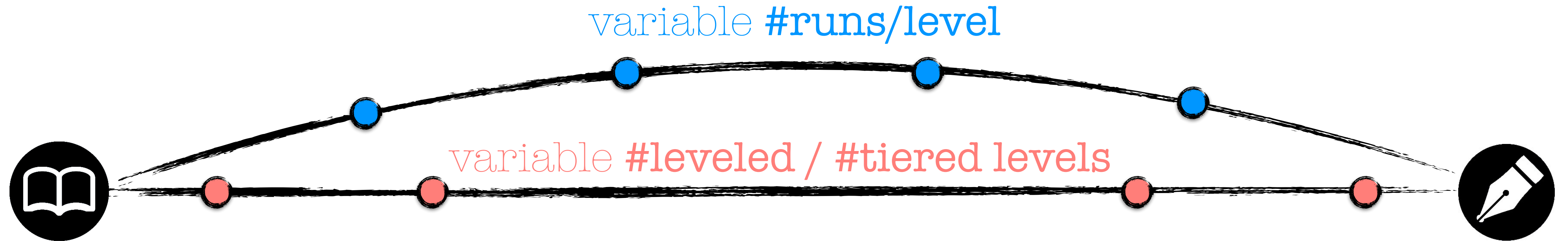
T

2

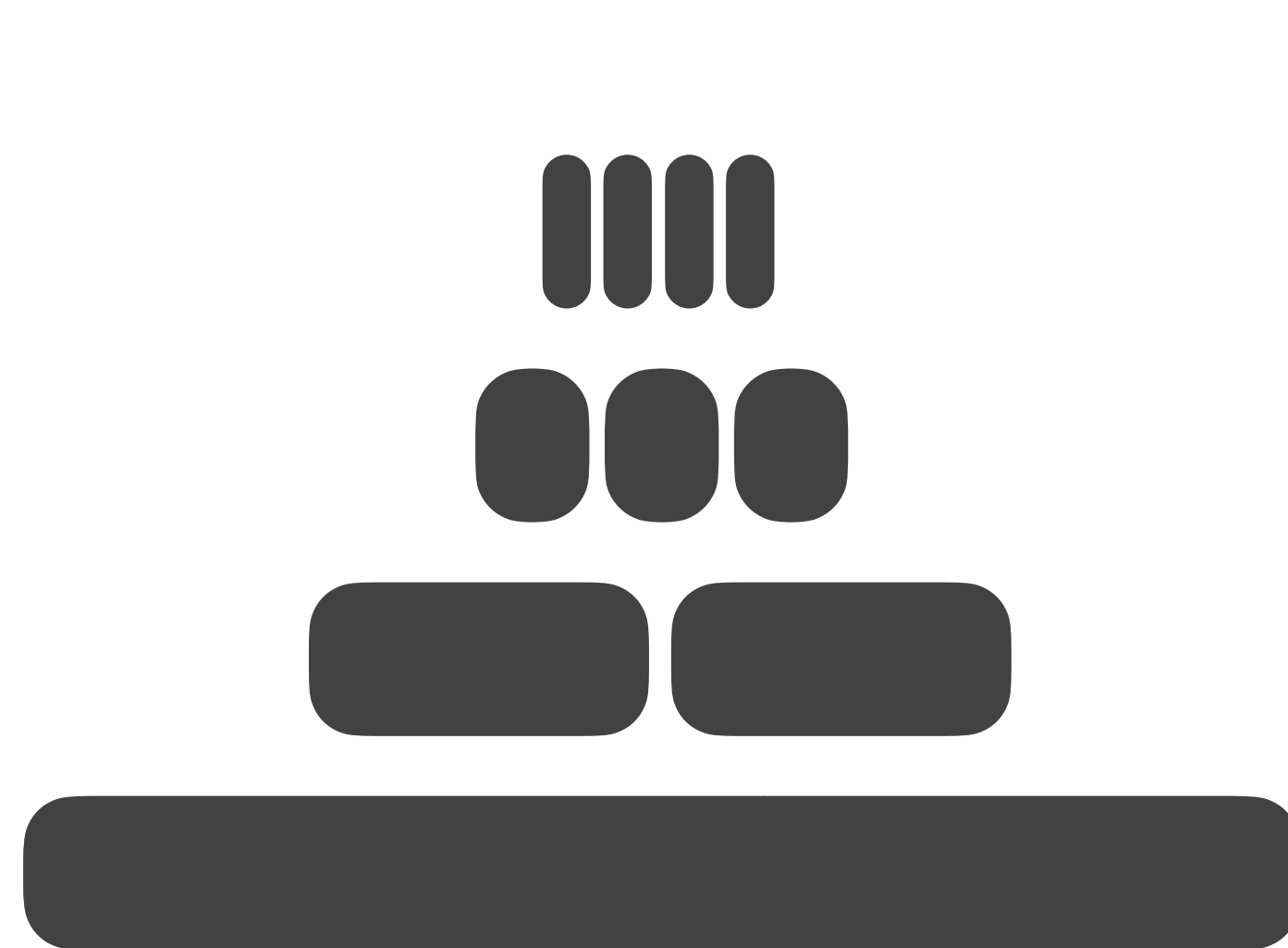
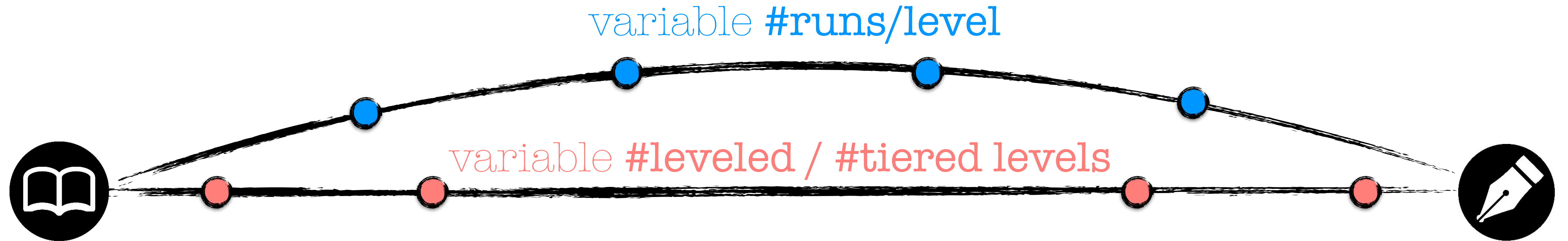
T

1

# Storage Layer Design Continuum



# Storage Layer Design Continuum



size ratio

2

2.5

3

4

#runs

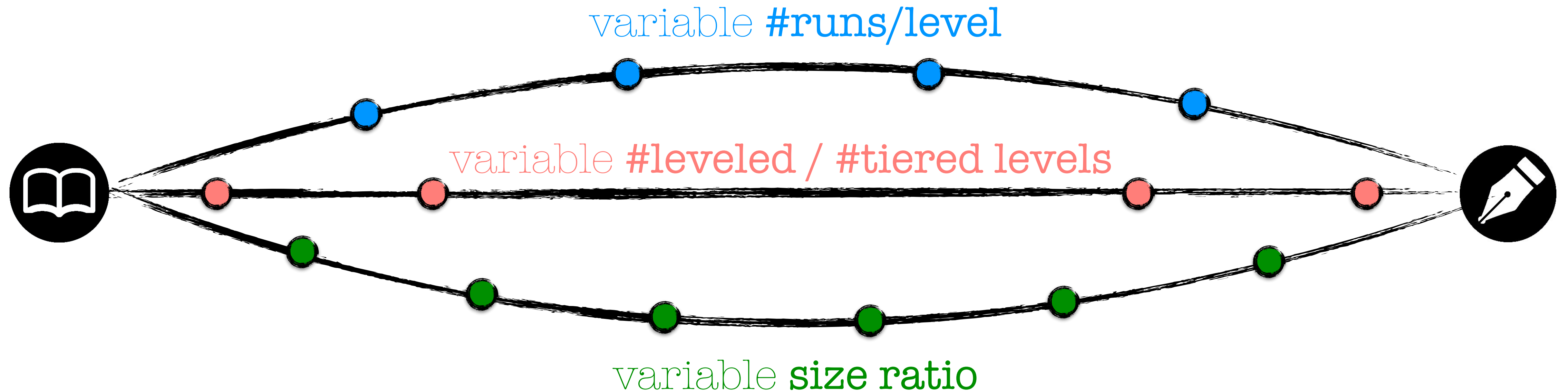
4

3

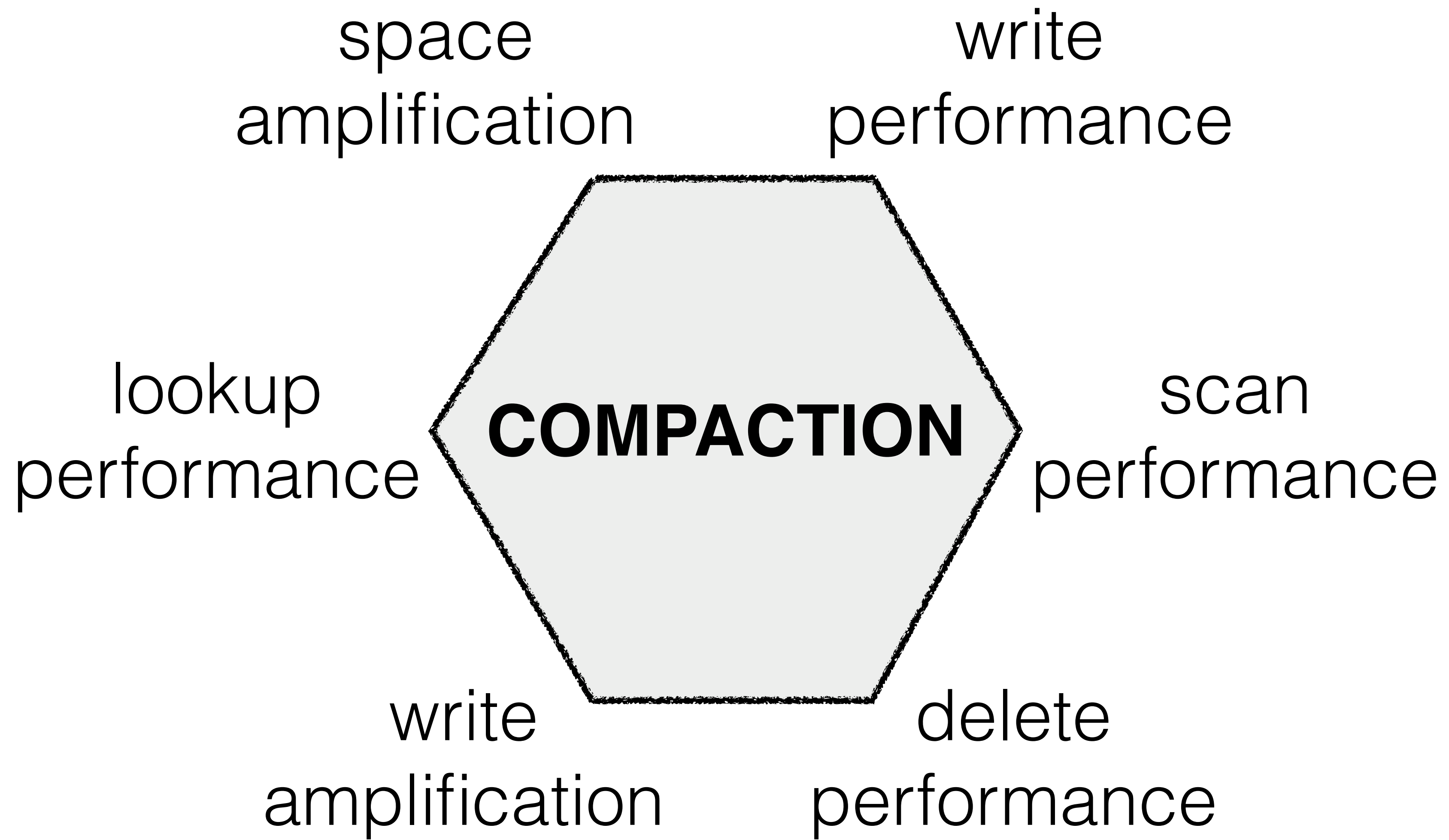
2

1

# Storage Layer Design Continuum

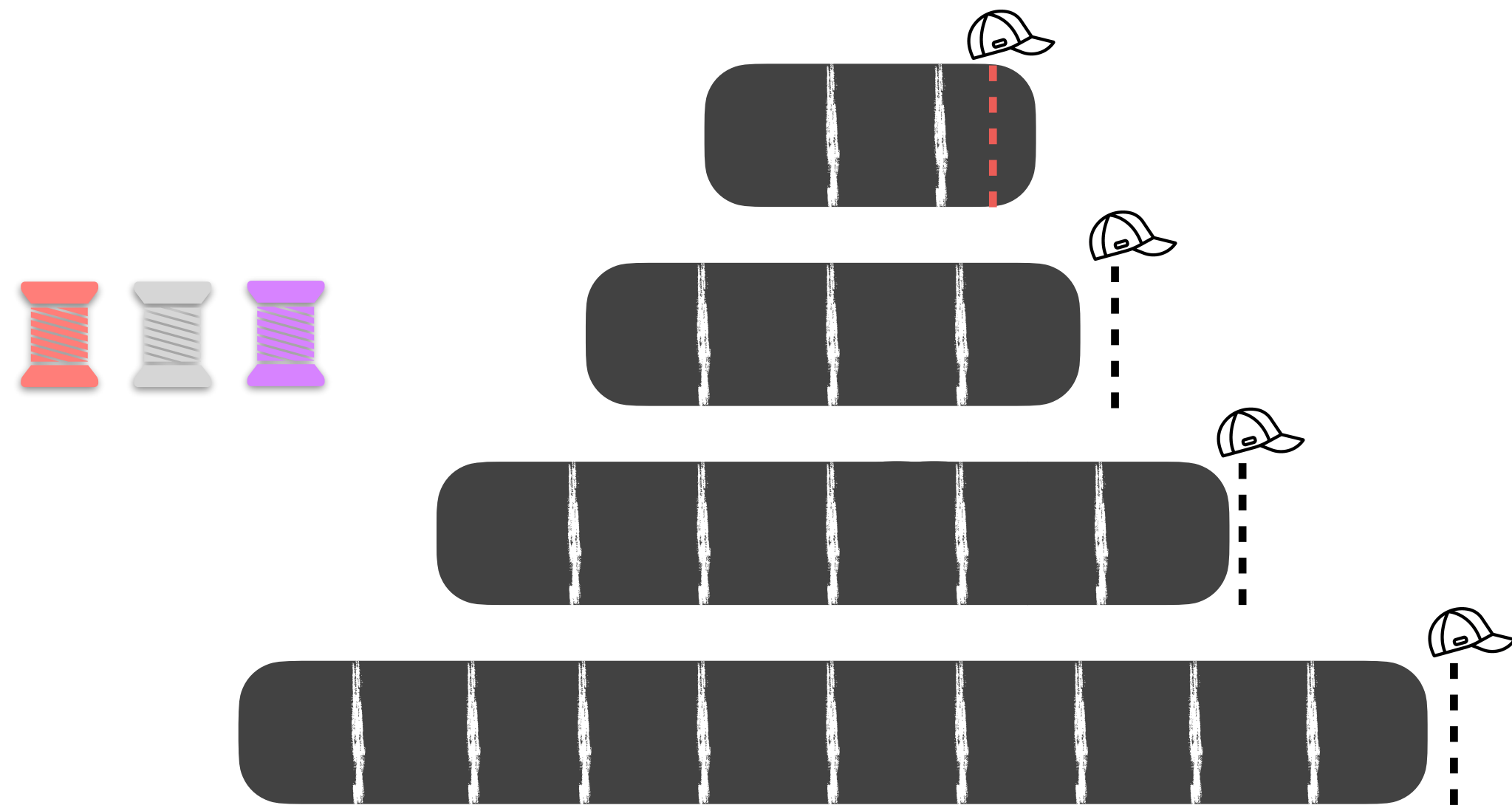


The LSM storage layer design continuum



# Optimizing Compactions

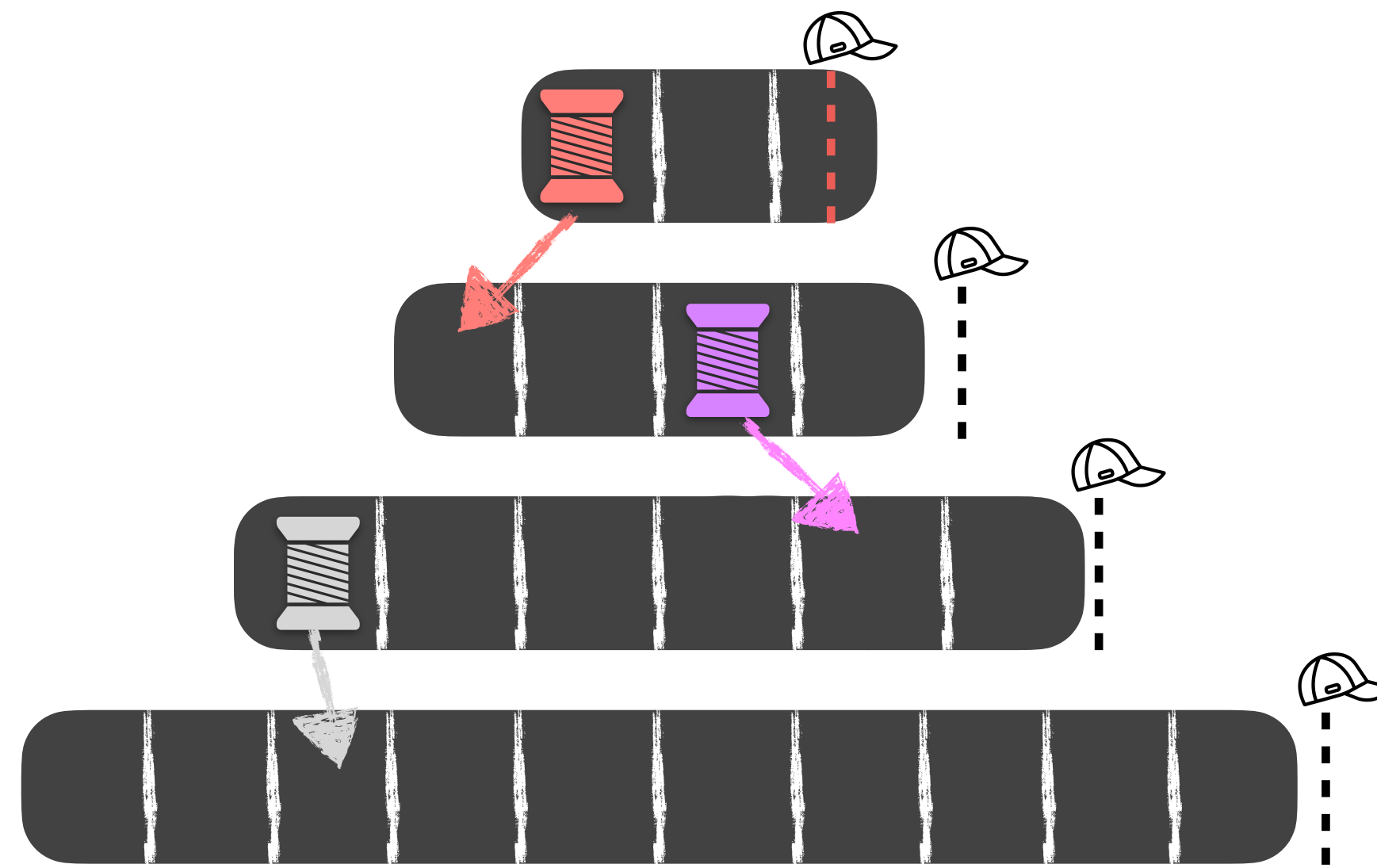
Background  
Compactions





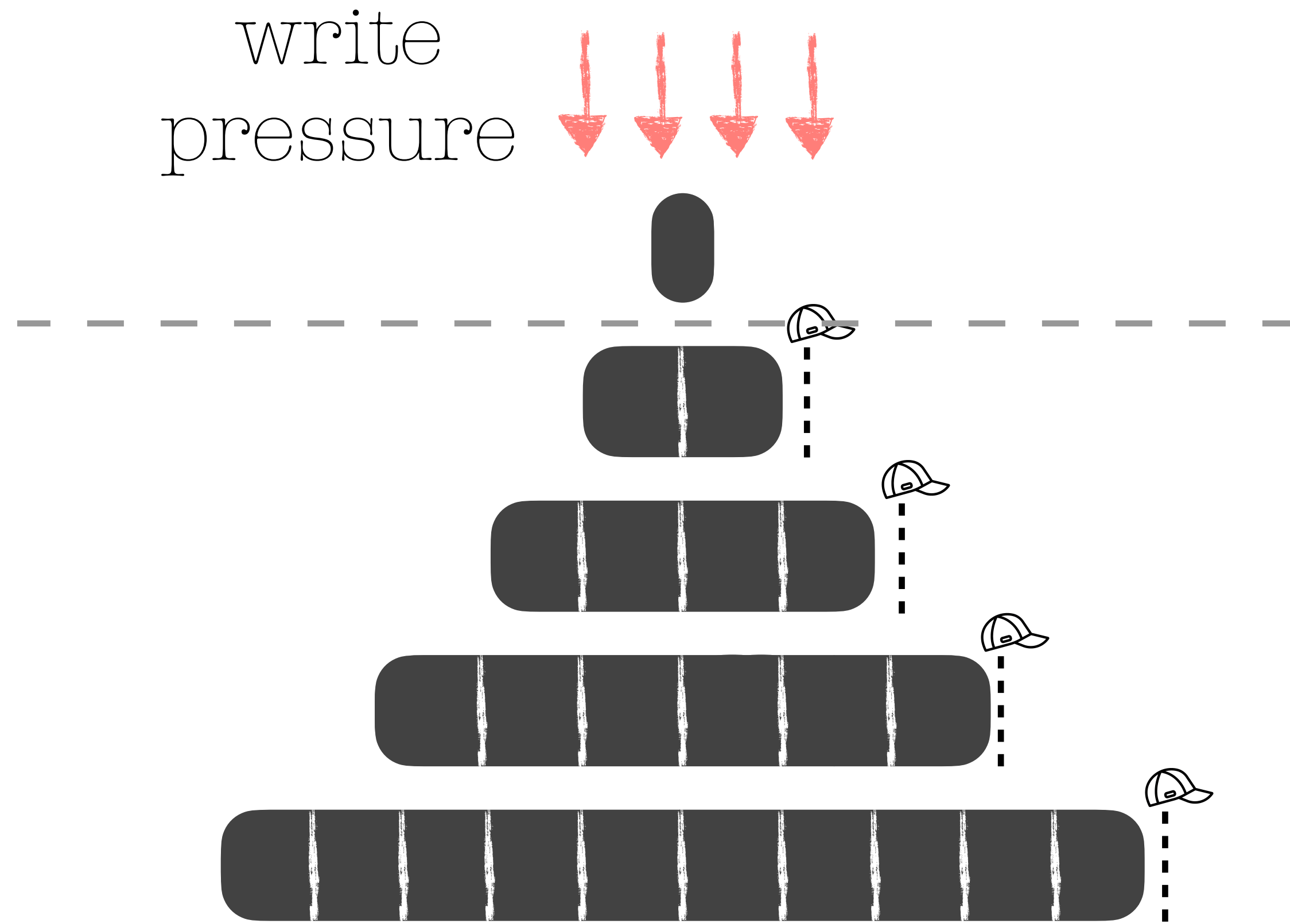
# Optimizing **Compactions**

## Background Compactions



- non-blocking reads/writes
- improves write throughput

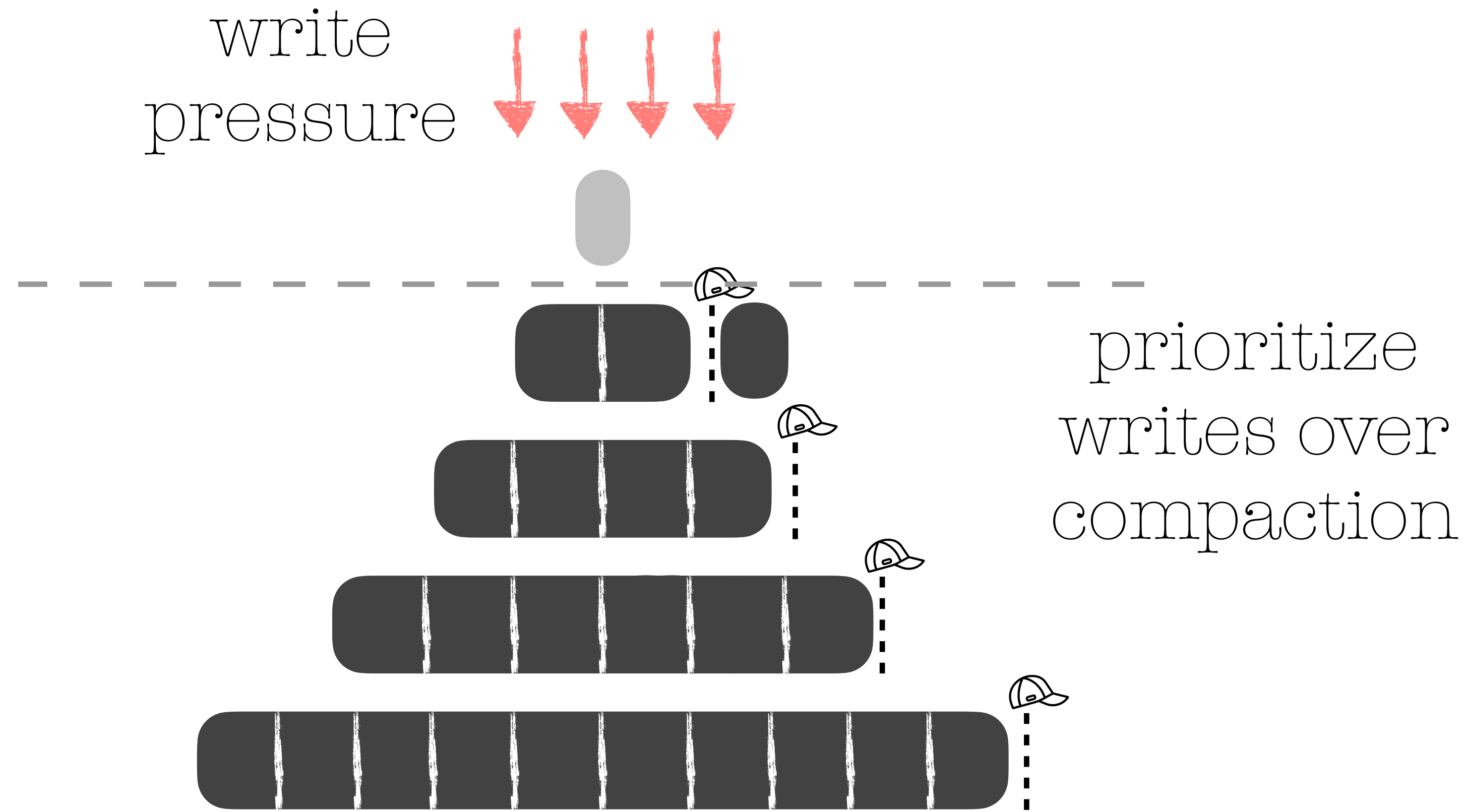
# Optimizing Compactions



Background  
Compactions

Compaction  
Priority

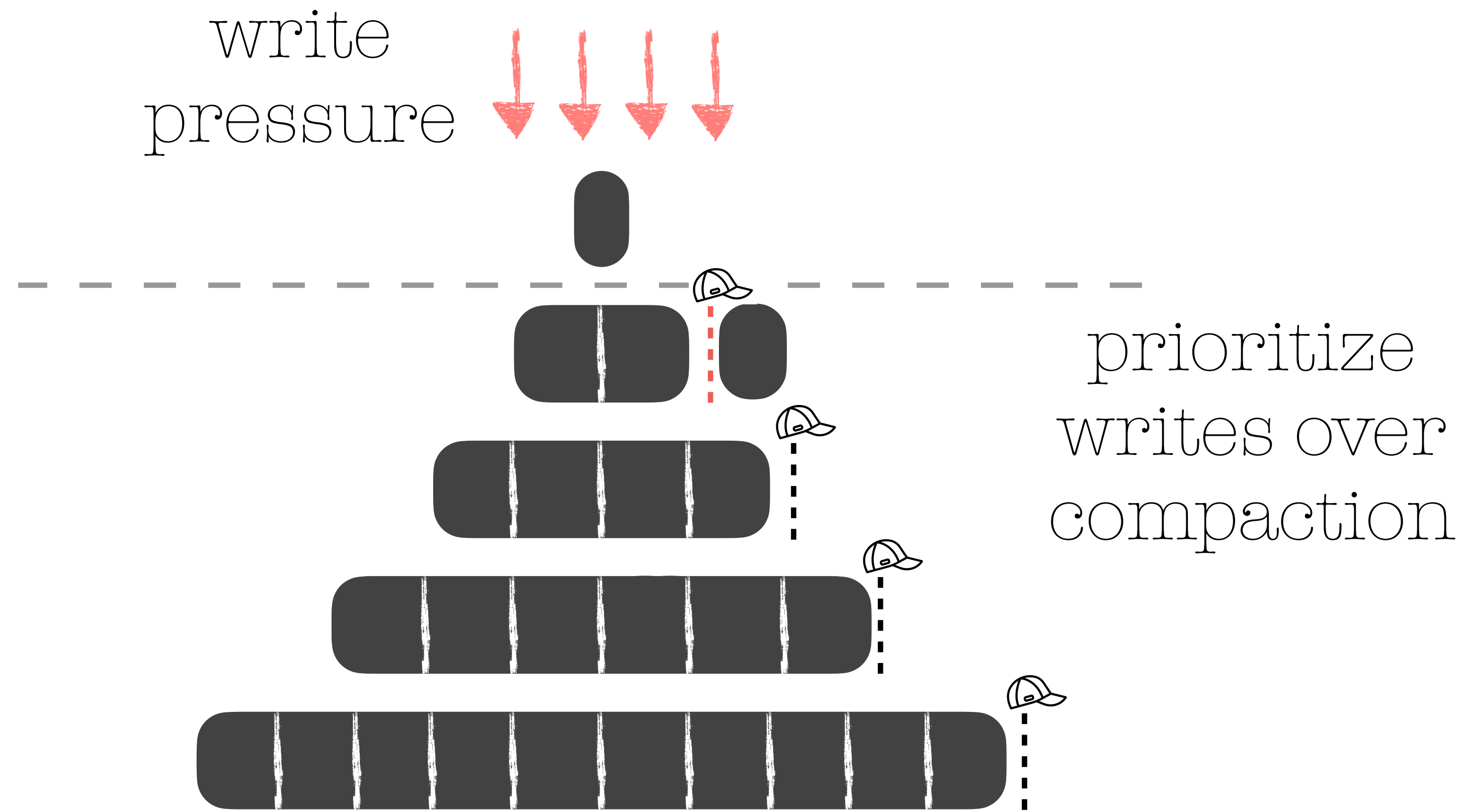
# Optimizing Compactions



Background Compactions

Compaction Priority

# Optimizing Compactions

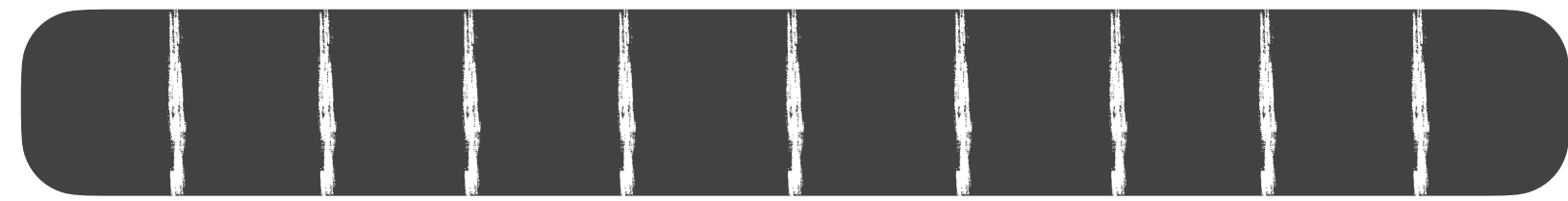
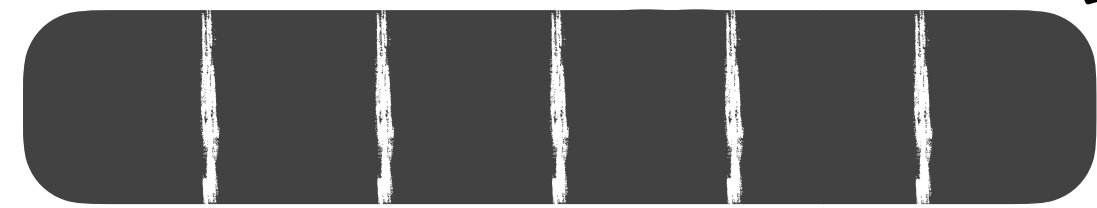
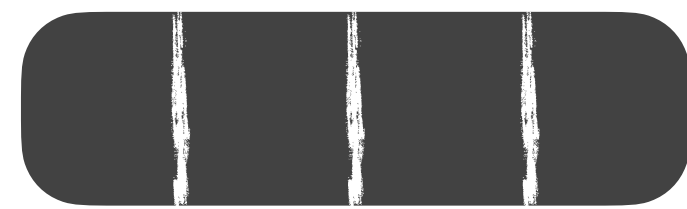
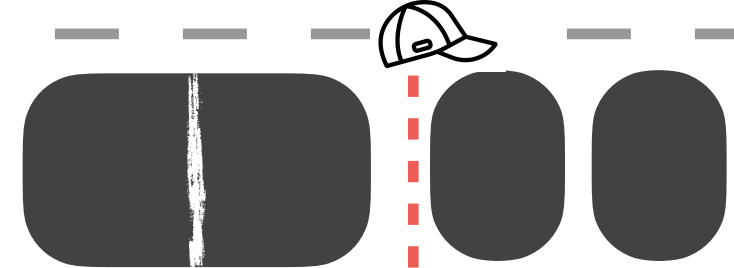
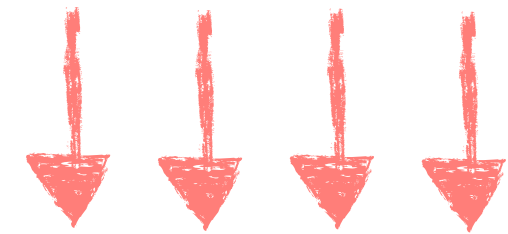


Background  
Compactions

Compaction  
Priority

# Optimizing Compactions

write  
pressure



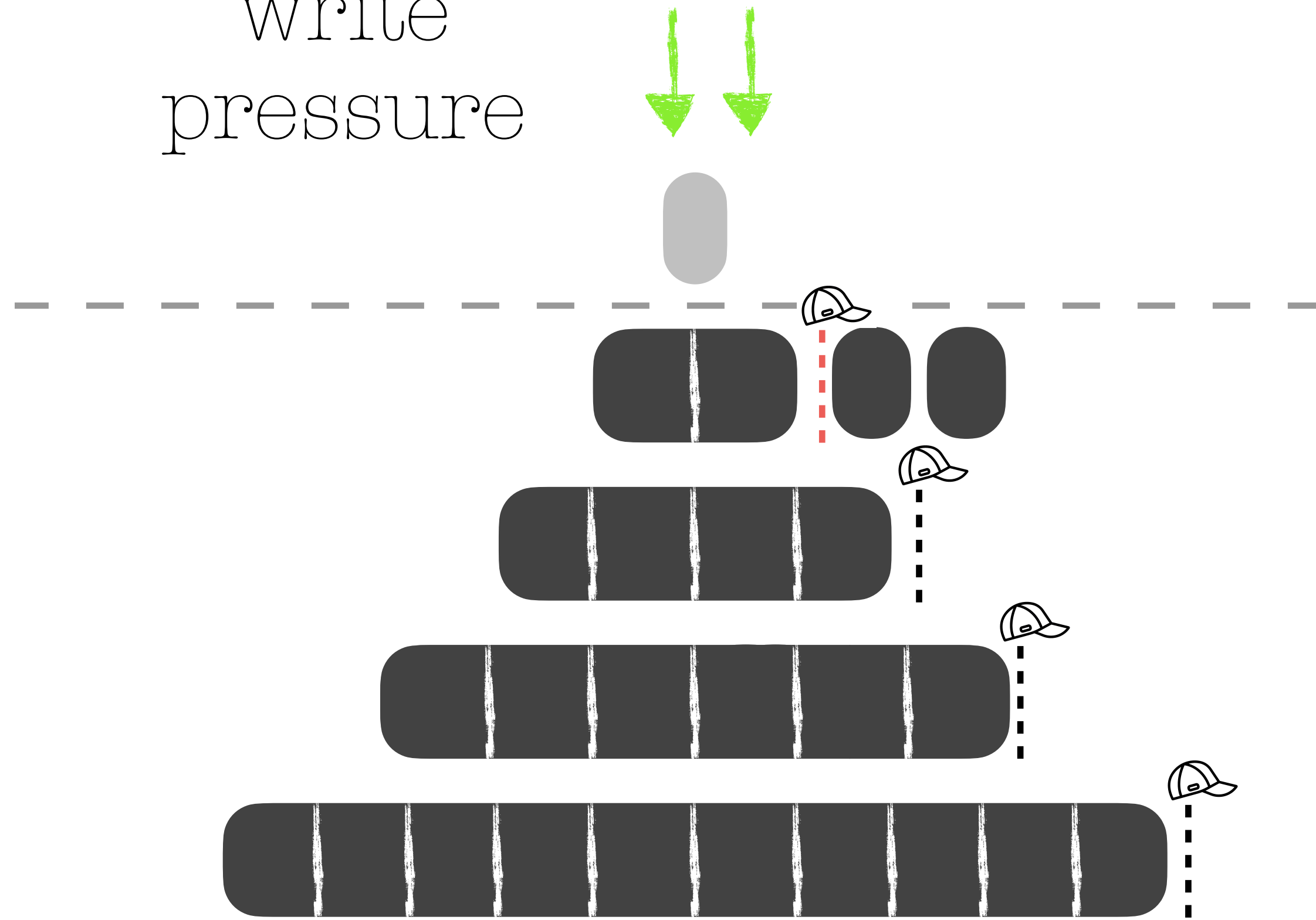
prioritize  
writes over  
compaction

Background  
Compactions

Compaction  
Priority

# Optimizing Compactions

write  
pressure

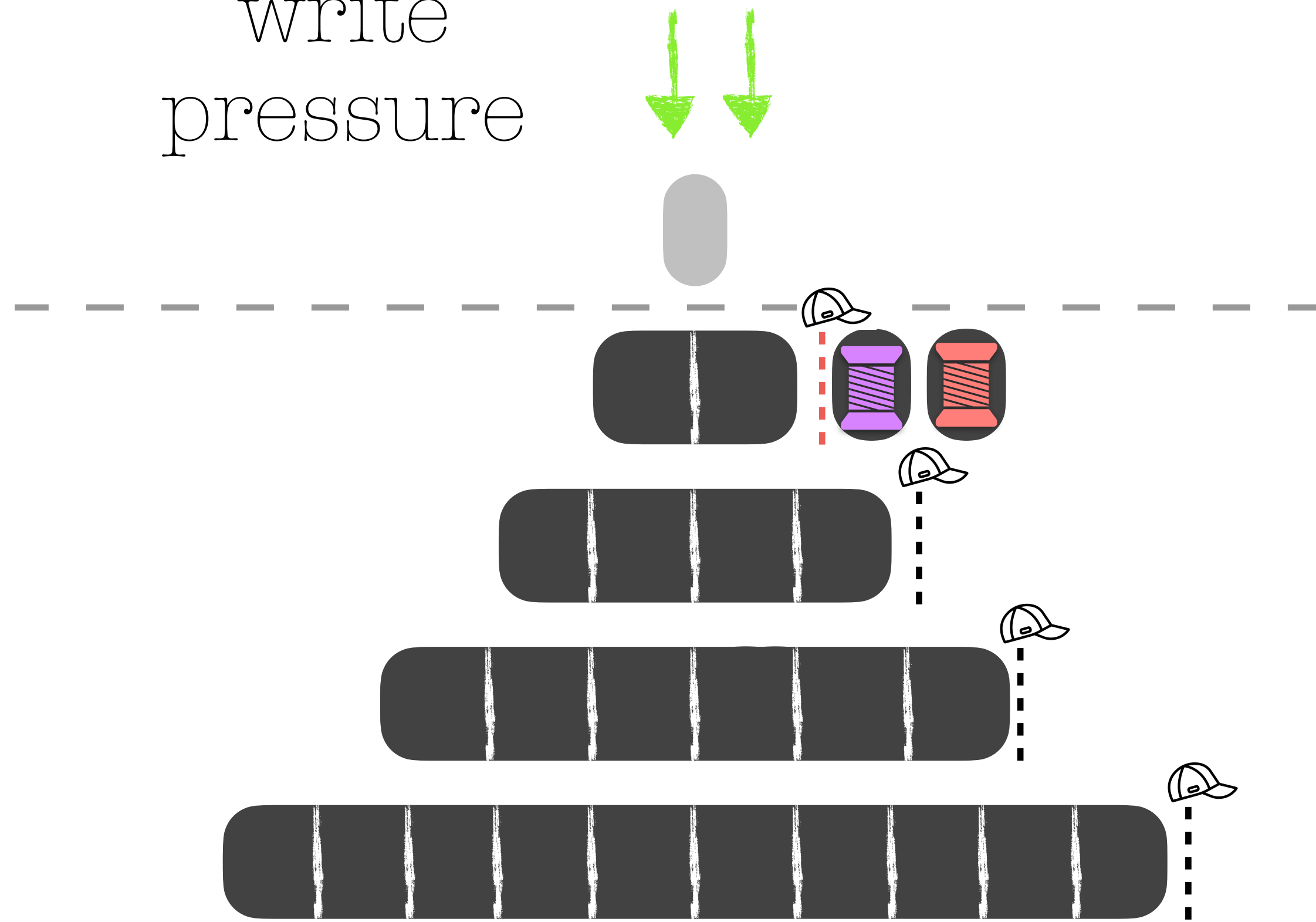


Background  
Compactions

Compaction  
Priority

# Optimizing Compactions

write  
pressure

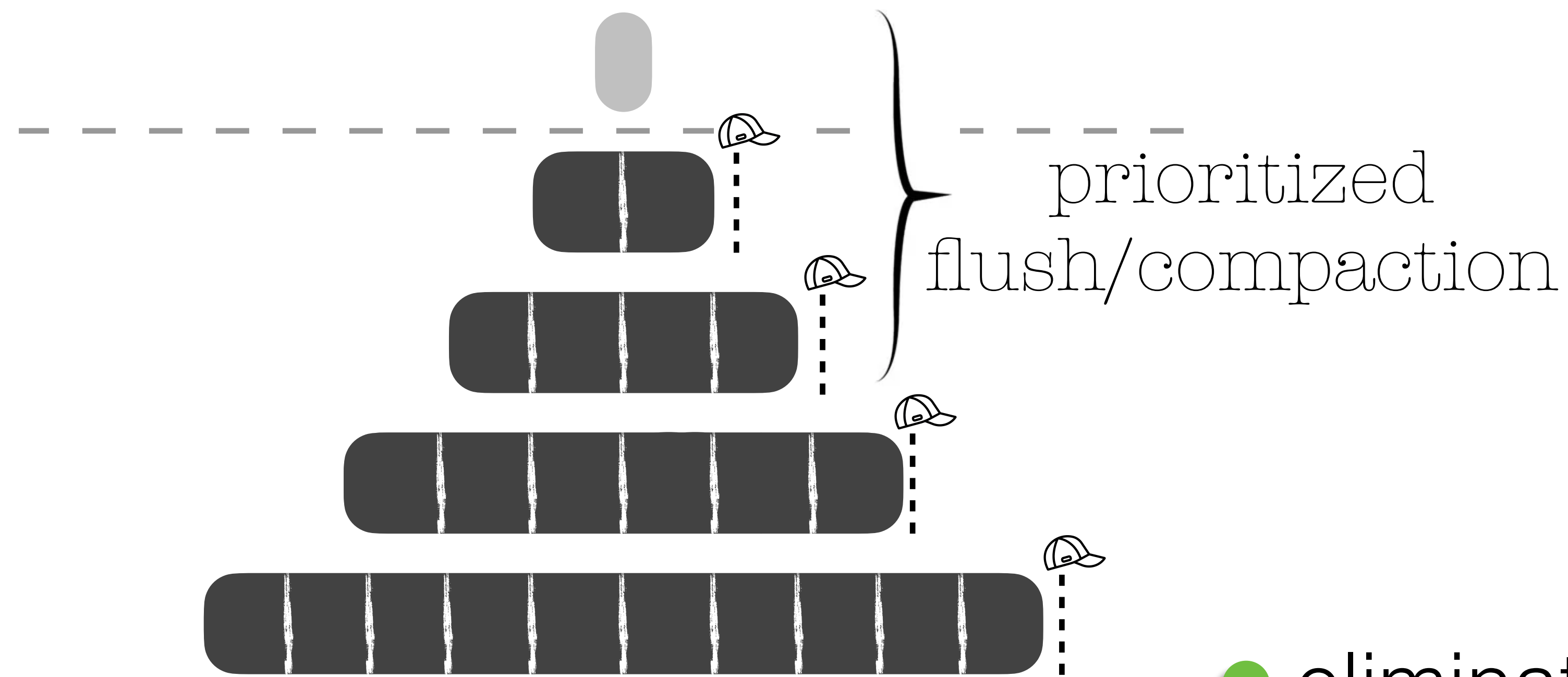


Background  
Compactions

Compaction  
Priority

- sustain heavy write bursts
- tree becomes out of shape

# Optimizing Compactions



BalmauATC19

BalmauToCS20

Background  
Compactions

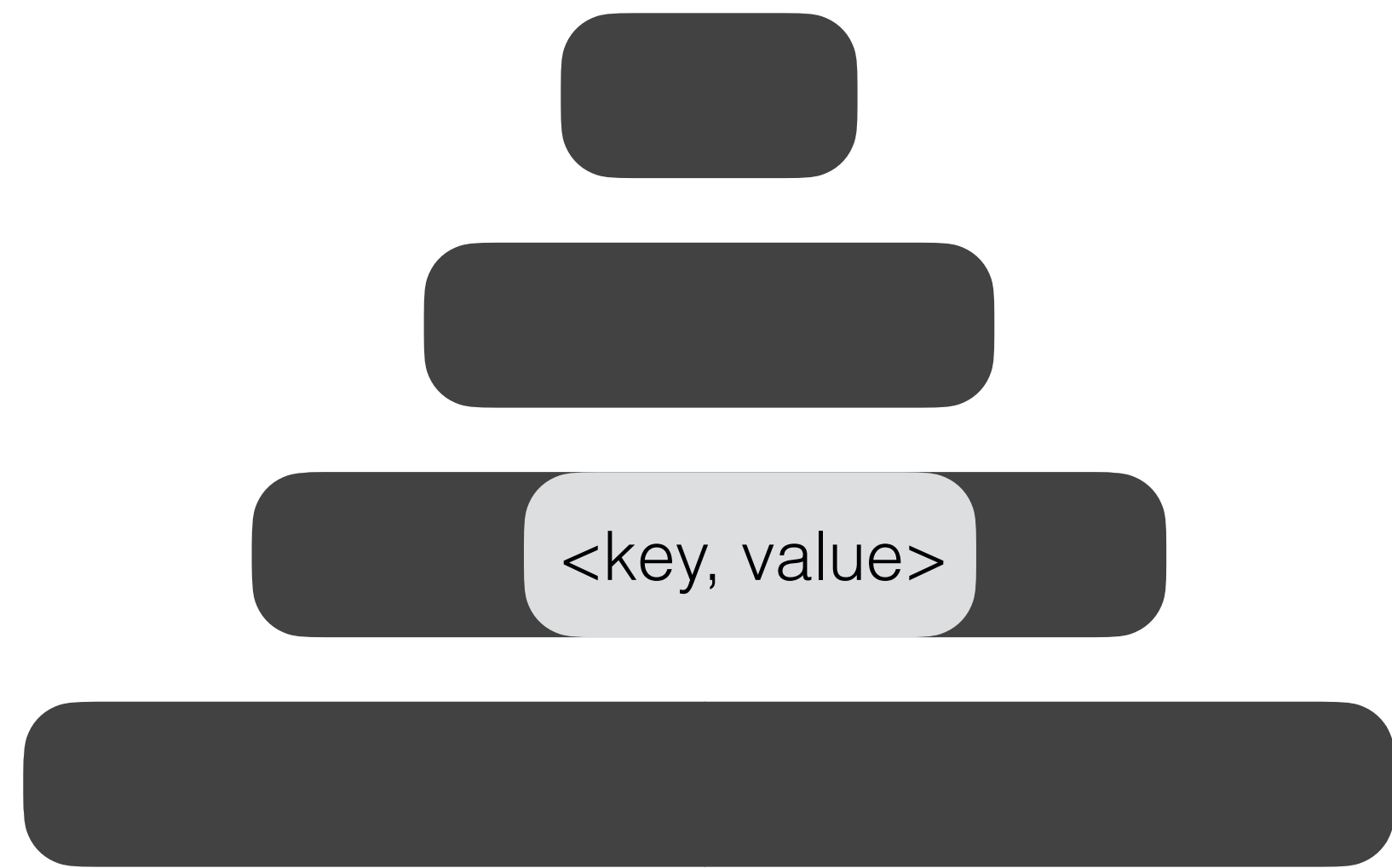
Compaction  
Priority

I/O Scheduler

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

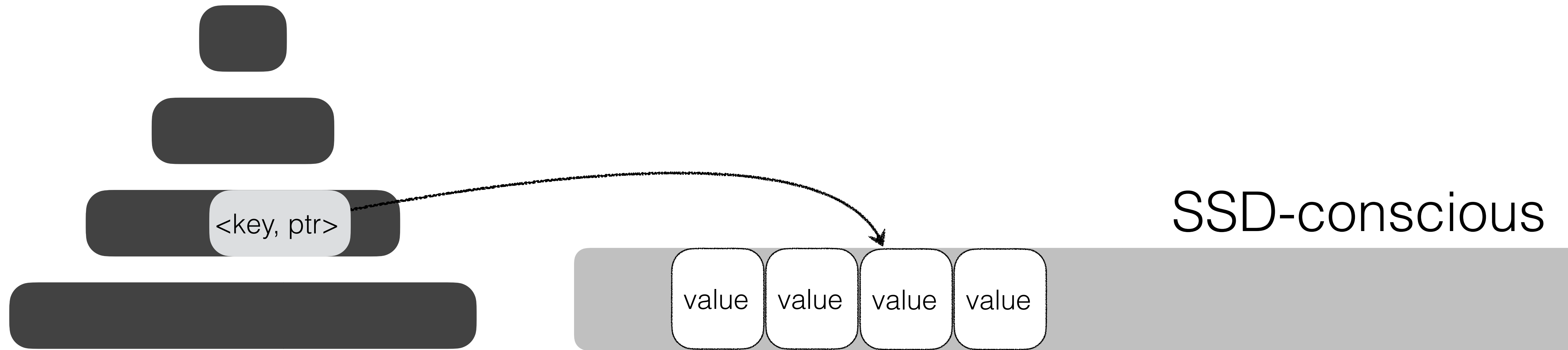


# Data Placement Variations



key-value separation   
LuFAST16

# Data Placement Variations



key-value separation 

LuFAST16

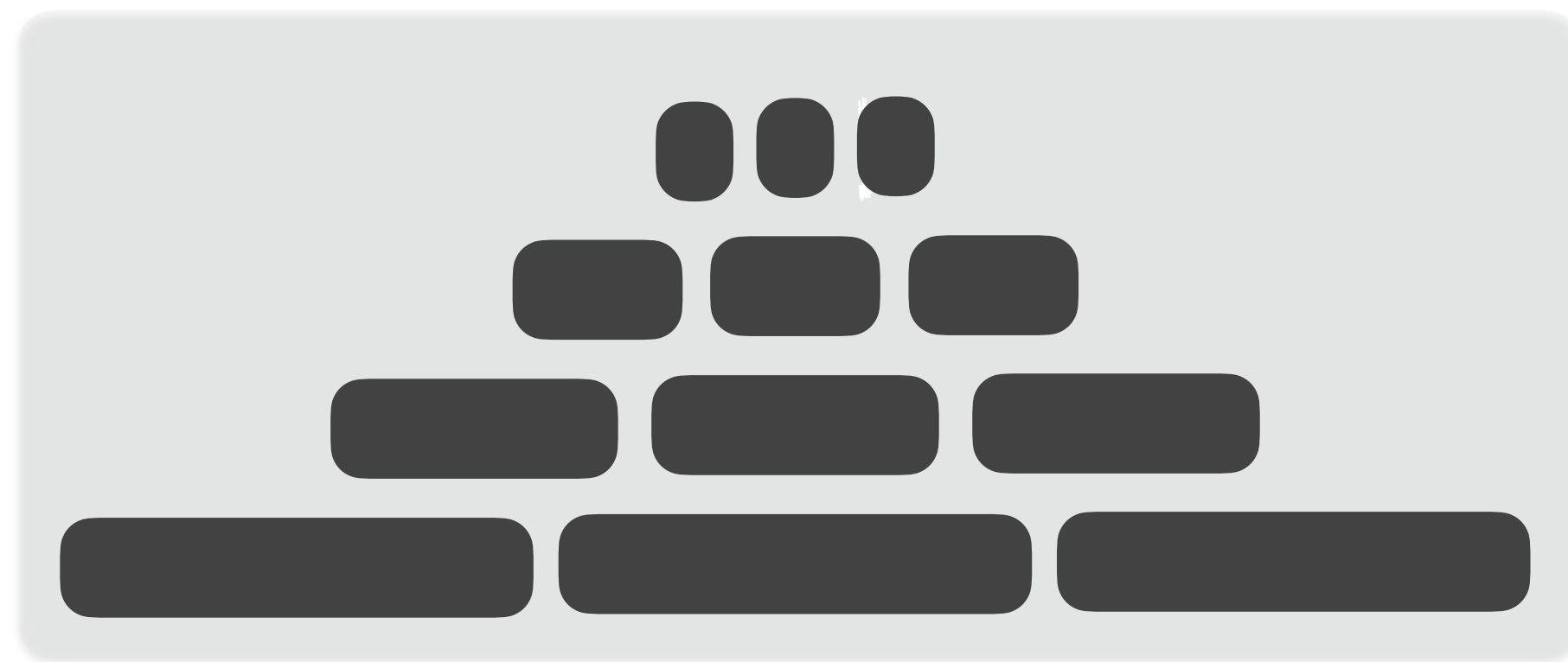
- reduced write amplification
- better read performance

# Data Placement Variations



partitioning / sharding

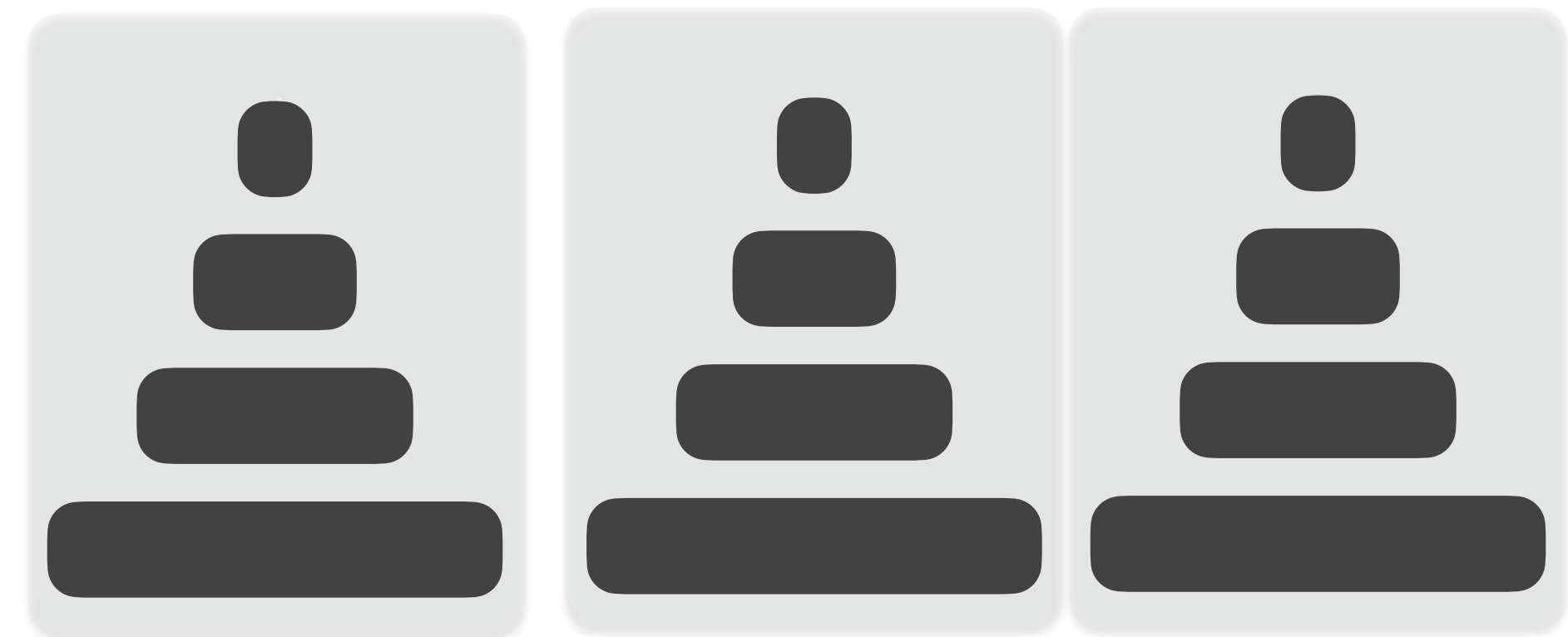
# Data Placement Variations



storage

partitioning

RajuSOSP17



storage-1

storage-2

storage-3

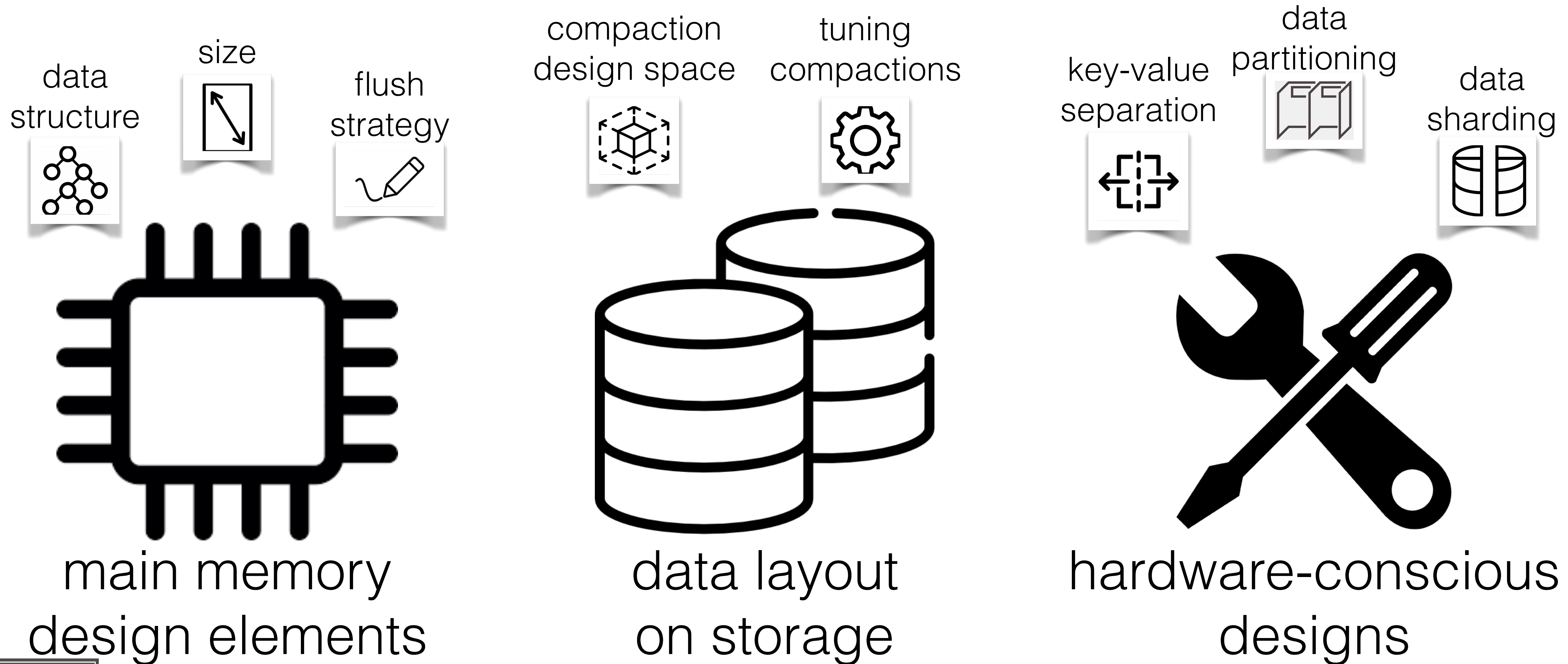
sharding

HuangSIGMOD21



- improved ingestion throughput
- reduced write amplification

# Summary: Ingestion Optimization



**No Labs on 02/03**

stay  
warm

# *CS 561: Data Systems Architecture*

## Class 5

# **Log-Structured Merge (LSM) Trees**

BOSTON  
UNIVERSITY

*Dr. Subhadeep Sarkar*

<https://bu-disc.github.io/CS561/>

