

class 6

Log-structured Merge Trees

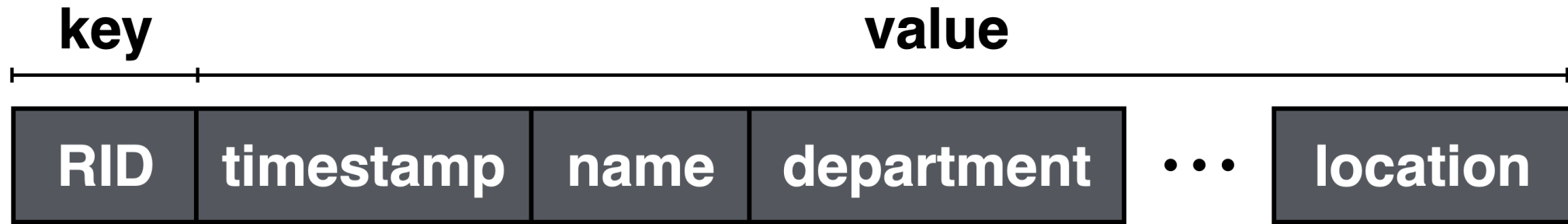
Prof. Manos Athanassoulis

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

Do we have a quiz today ... ?

Yes!

Key-Value Stores

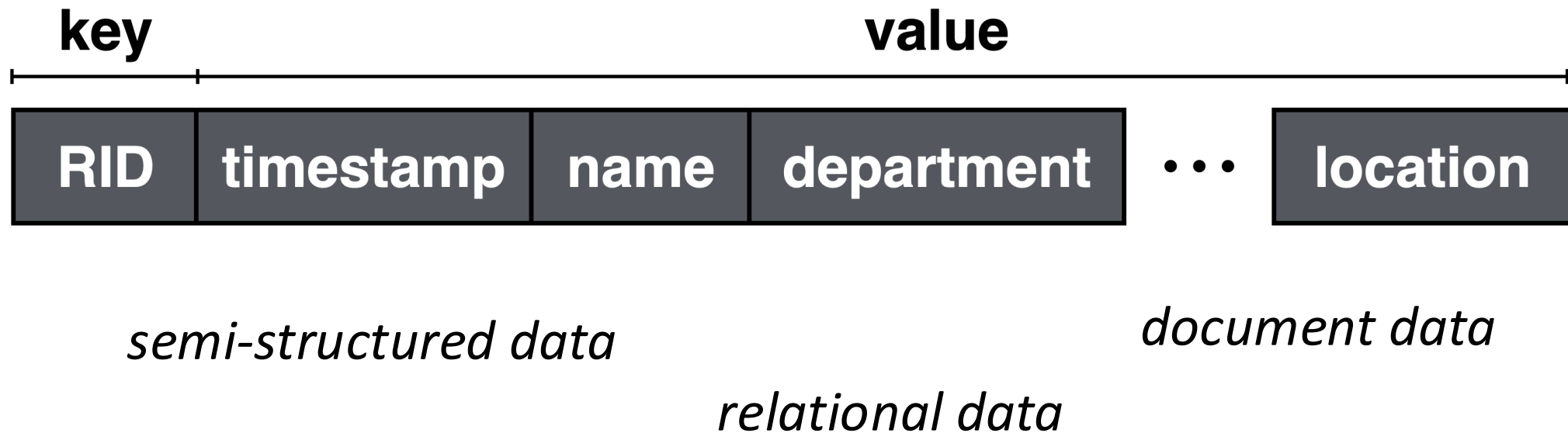


semi-structured data

relational data

document data

LSM-tree based Key-Value Stores



Log-**S**tructured **M**erge-tree

LSM-tree

LSM-tree

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

1996

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

LSM-tree
O'Neil *et al.*

1996

Patrick O'Neil
UMass Boston

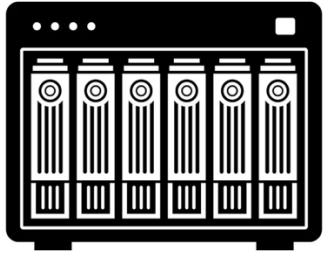


LSM-tree
O'Neil *et al.*

1996

● good random writes

● good reads



array
of discs

why?

RAID, striping ← ?

✗ LSM not explicitly needed

LSM-tree
O'Neil *et al.*

so, arrays of disks were enough!



how many IOPS?

10KRPM

max seek time 1.5ms

100 disks

10KRPM: 10K rev in 60s

$60/10000 = 6\text{ms}$ per rev

avg. rot. delay: 3ms (6ms/2)

avg. seek time: 0.75ms (1.5ms/2)

1 I/O / 3.75ms: 267 IOPS

100 disks: 26,700 IOPS



Bigtable

1980s

1996

2006

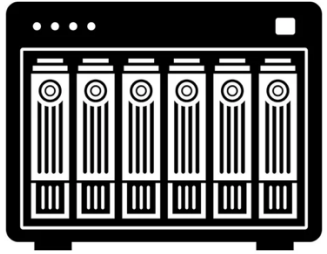
a decade

● good random writes

● good reads

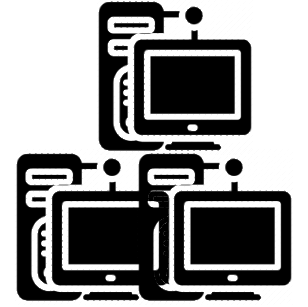
● poor ingestion perf.

● poor query perf.



array
of discs

what happened in 2006?



commodity
hardware

LSM-tree
O'Neil *et al.*

We set up a Bigtable cluster with N tablet servers to measure the performance and scalability of Bigtable as N is varied. The tablet servers were configured to use 1 GB of memory and to write to a GFS cell consisting of 1786 machines with two 400 GB IDE hard drives each.



Bigtable

1980s

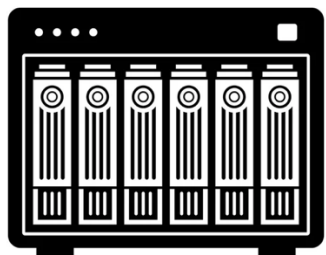
1996

2006

a decade

● good random writes

● good reads



array
of discs

SSD wear-friendly

competitive rand. reads

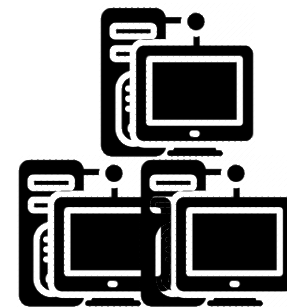
fast ingestion (sequential)

● poor ingestion perf.

● poor query perf.



why this
is important?



commodity
hardware



LSM-tree
O'Neil *et al.*

SSDs have dominated
storage since!

1980s

1996

2006

a decade



Bigtable

LSM-tree
O'Neil *et al.*

1996



Bigtable

2006

APACHE
HBASE 

2007

LSM-tree
O'Neil *et al.*

1996


Bigtable

2006

APACHE
HBASE 

2007


cassandra

2010

LSM-tree
O'Neil *et al.*

1996


Bigtable

2006

APACHE
HBASE 

2007


cassandra

2010


levelDB

2011

LSM-tree
O'Neil *et al.*

1996



2006



2007



2010

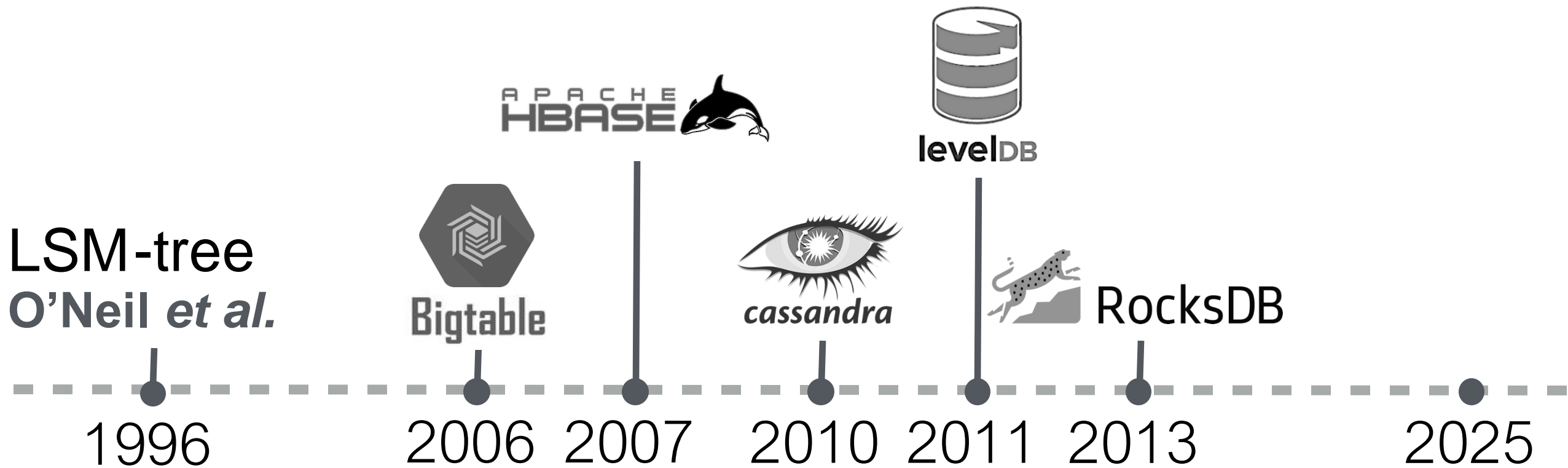


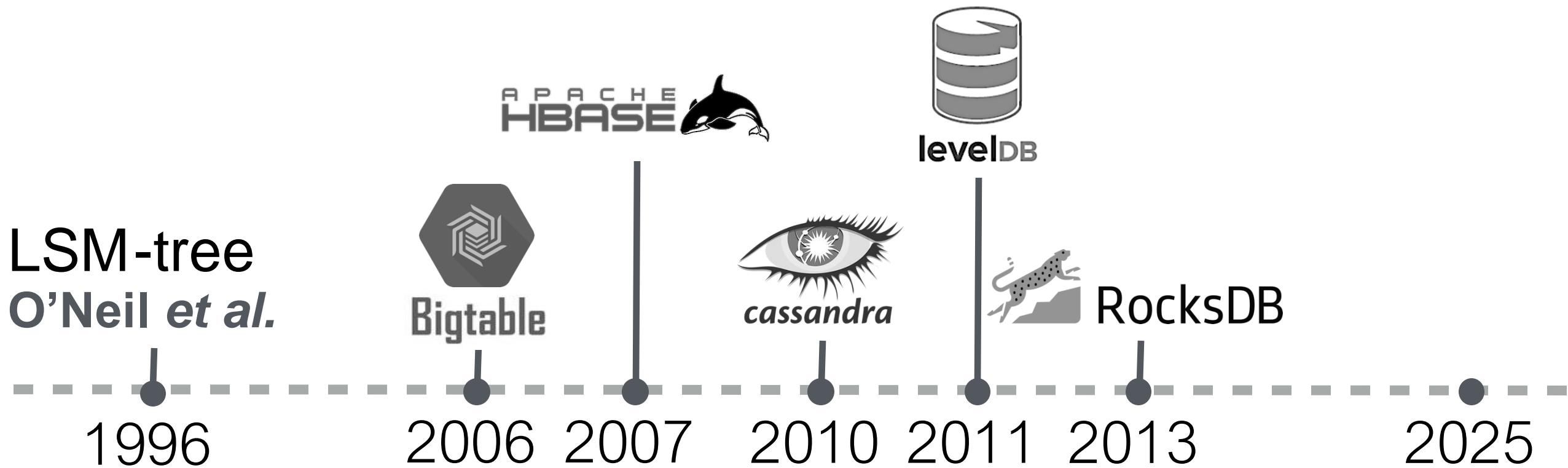
2011



RocksDB

2013





LSM-tree

NoSQL



RocksDB



levelDB



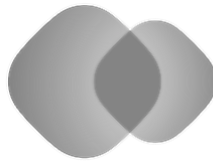
SCYLLA



DynamoDB



cassandra



tarantool



Bigtable

APACHE
HBASE



riak

accumulo™



SQLite



relational



influxdb



QuasarDB

time-series

2025

LSM-tree

NoSQL



RocksDB



levelDB



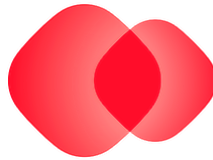
SCYLLA



DynamoDB



cassandra



tarantool



Bigtable

APACHE
HBASE



accumulo™

riak



SQLite



relational



influxdb

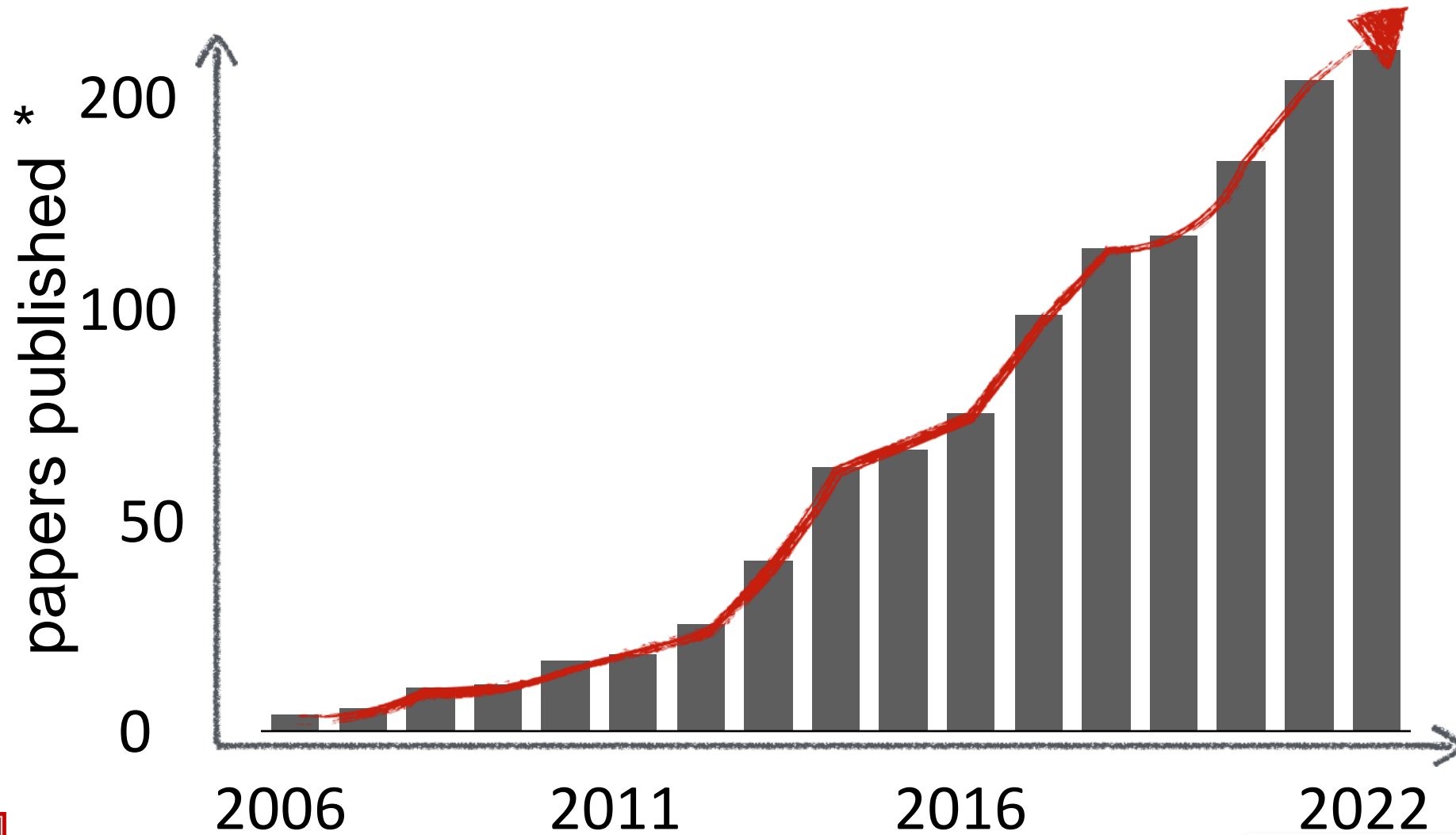


QuasarDB

time-series

2025

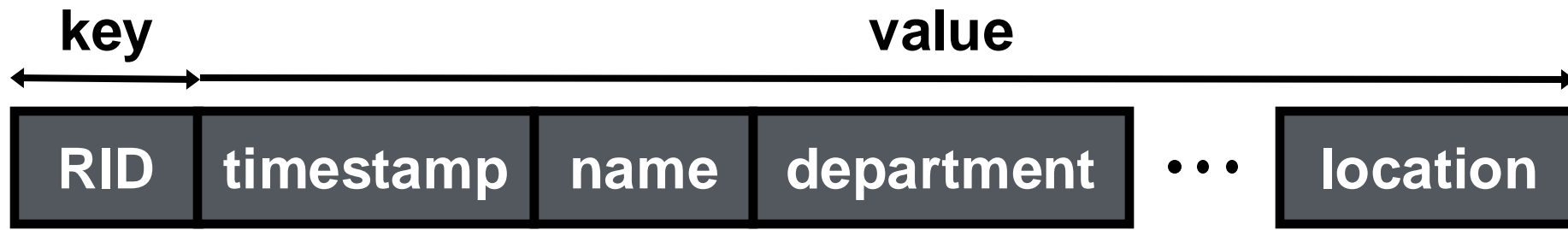
Research Trend



* data from Google scholar

LSM Basics

key-value
pairs



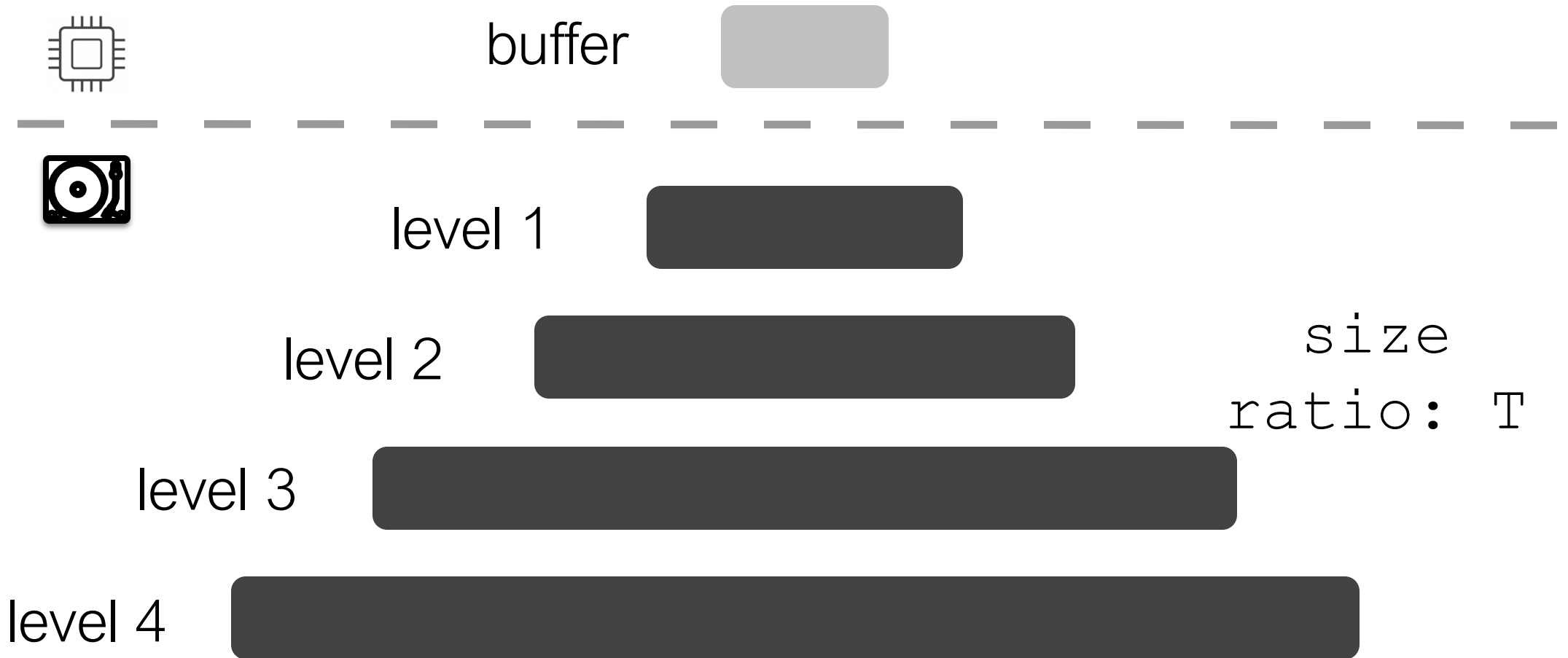
LSM Basics

key-value
pairs



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

LSM Basics



LSM Operating Principles

Buffering
ingestion

Immutable files on
storage

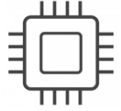
Out-of-place updates &
deletes

Periodic data layout
reorganization

Buffering ingestion

put(6)

put(2)

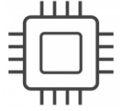


buffer

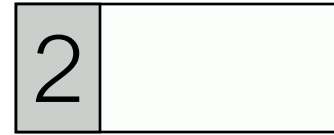


Buffering ingestion

put(1)
put(6)



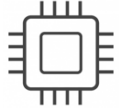
buffer



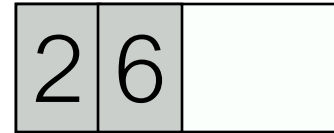
Buffering ingestion

put(4)

put(1)

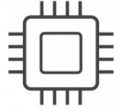


buffer



Buffering ingestion

put(4)

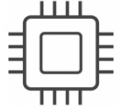


buffer

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



Buffering ingestion

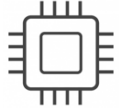


buffer

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



Buffering ingestion

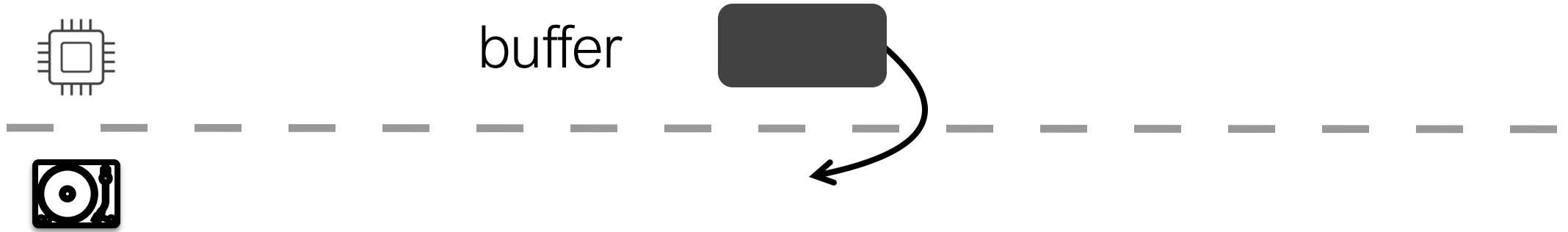


buffer

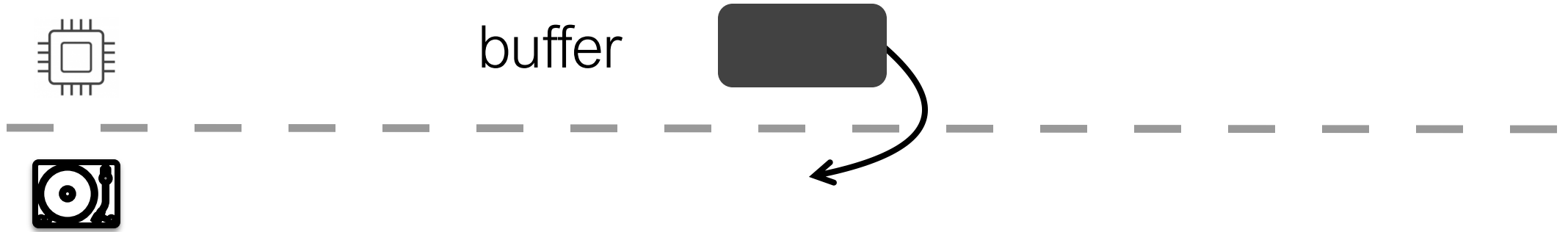
1	2	4	6
---	---	---	---



Buffering ingestion

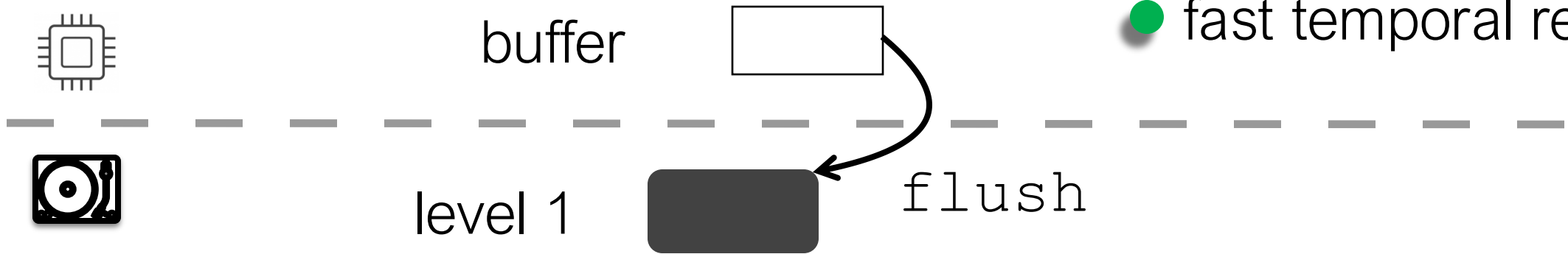


Buffering ingestion



Buffering ingestion

- low ingestion cost
- fast temporal reads



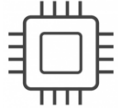
Immutable files on storage

- compact storage
- good ingestion throughput

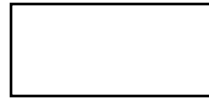


How do we update data?

put(6)

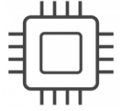


buffer



level 1

6



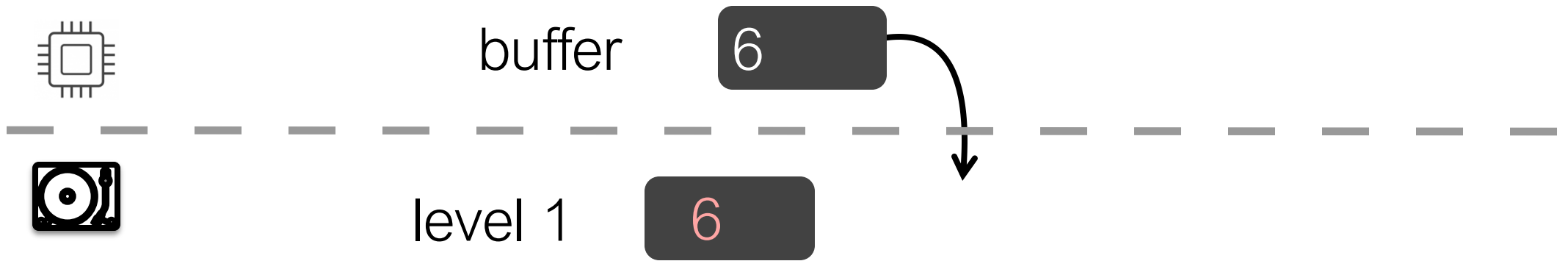
buffer

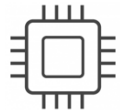


level 1



logically
invalidated





buffer

6



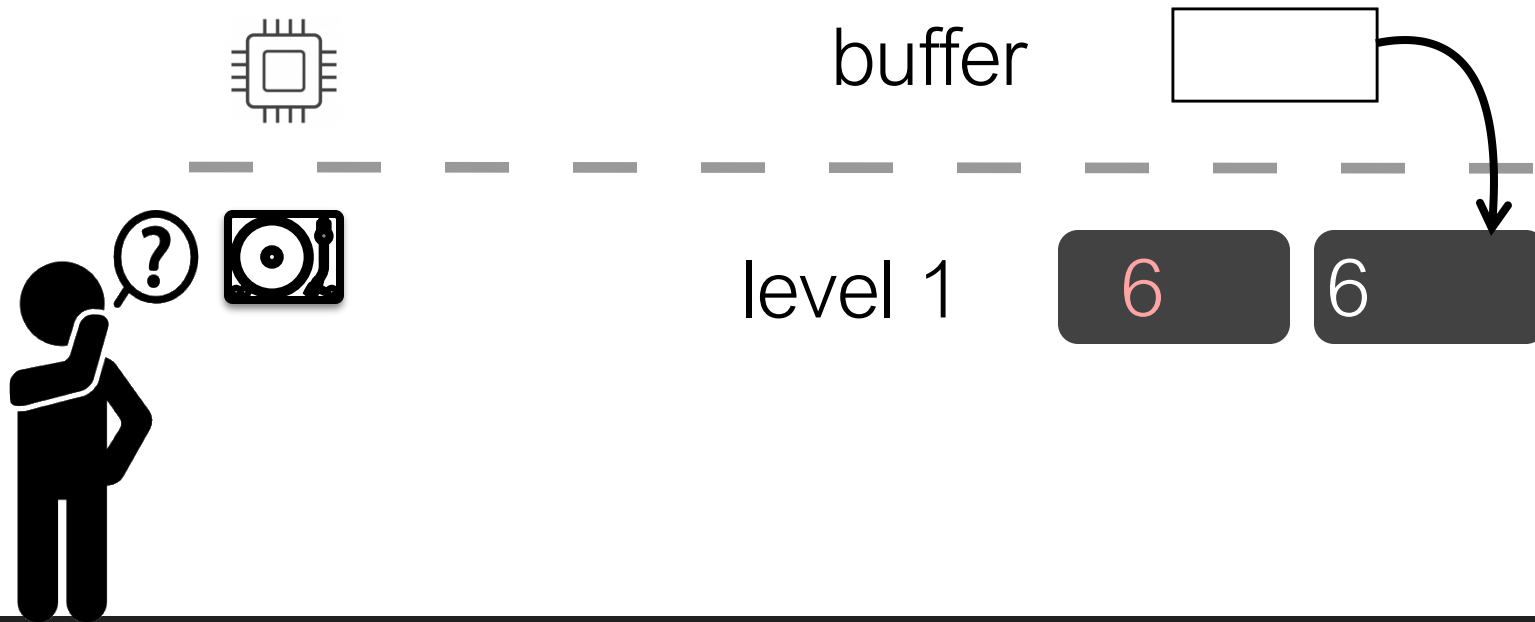
level 1

6

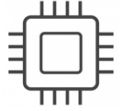


Out-of-place
updates

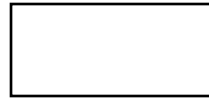
- fast ingestion
- space amplification
- slow reads



How do we reduce this space
amplification?



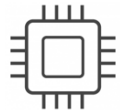
buffer



level 1

6

6



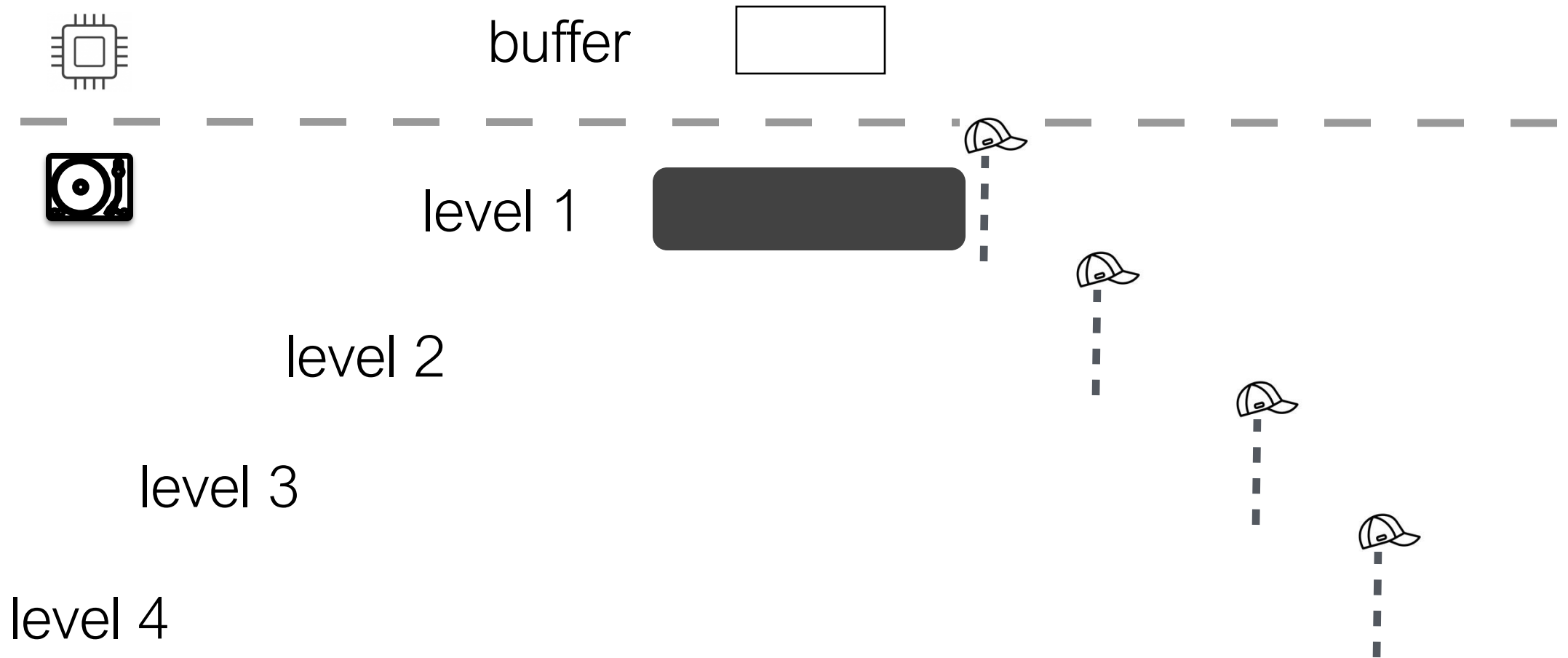
buffer

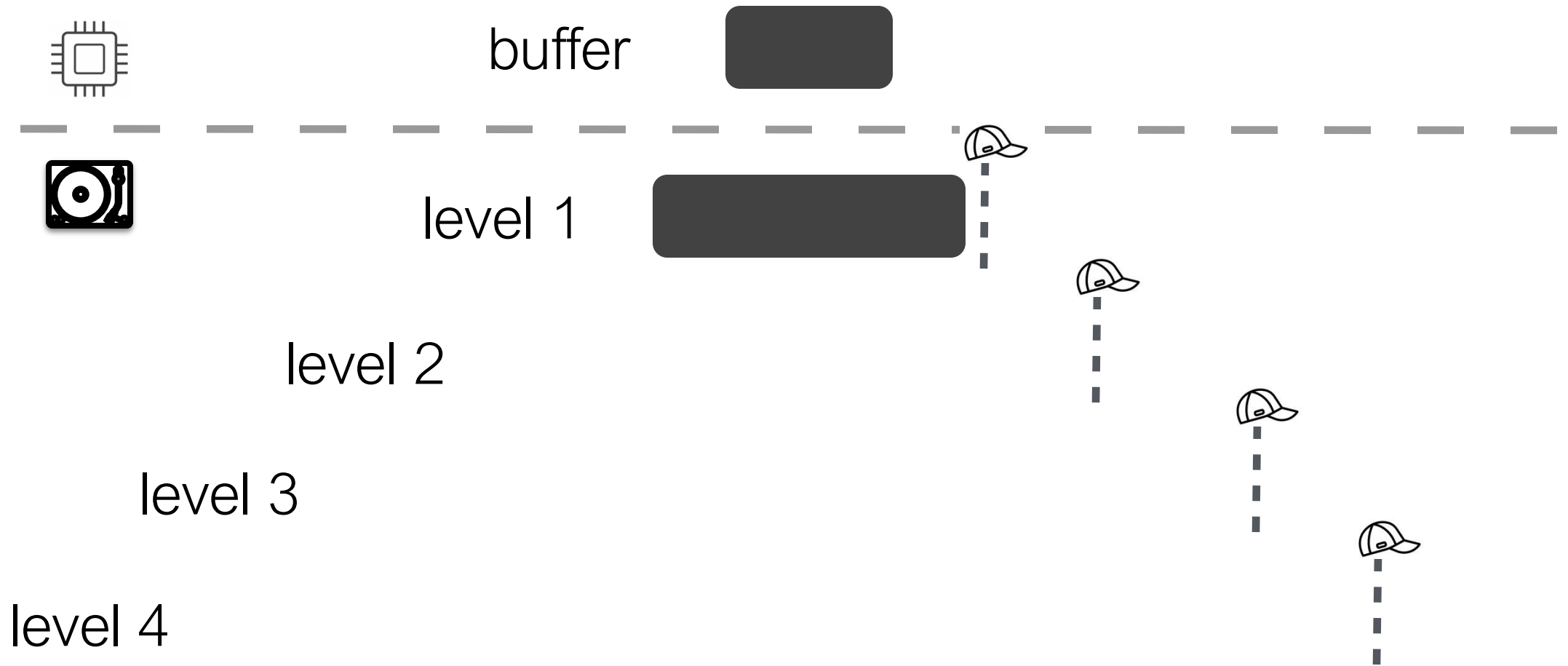


level 1

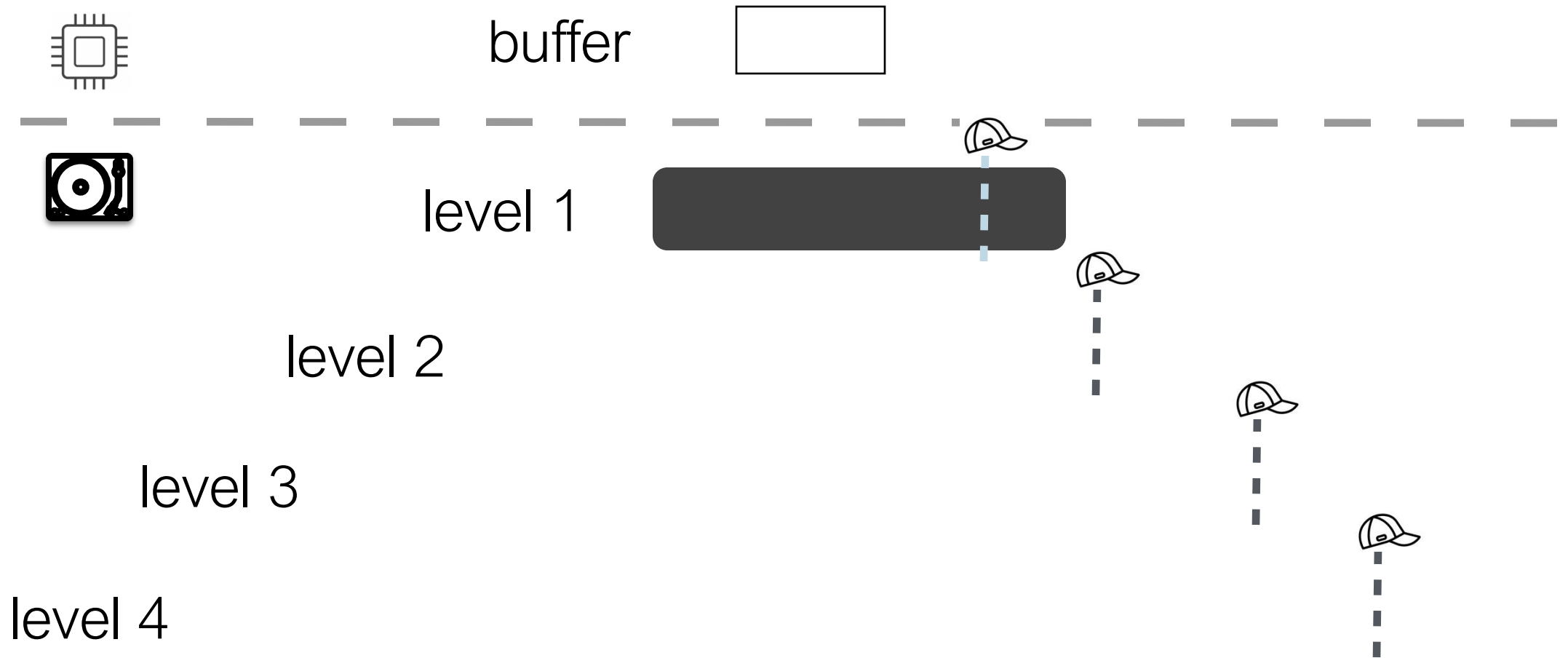
6

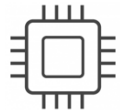
compaction



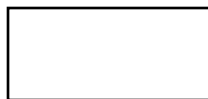








buffer



level 1

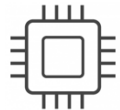
level 2



level 3

level 4





buffer



level 1

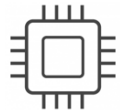
level 2



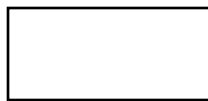
level 3

level 4





buffer



level 1



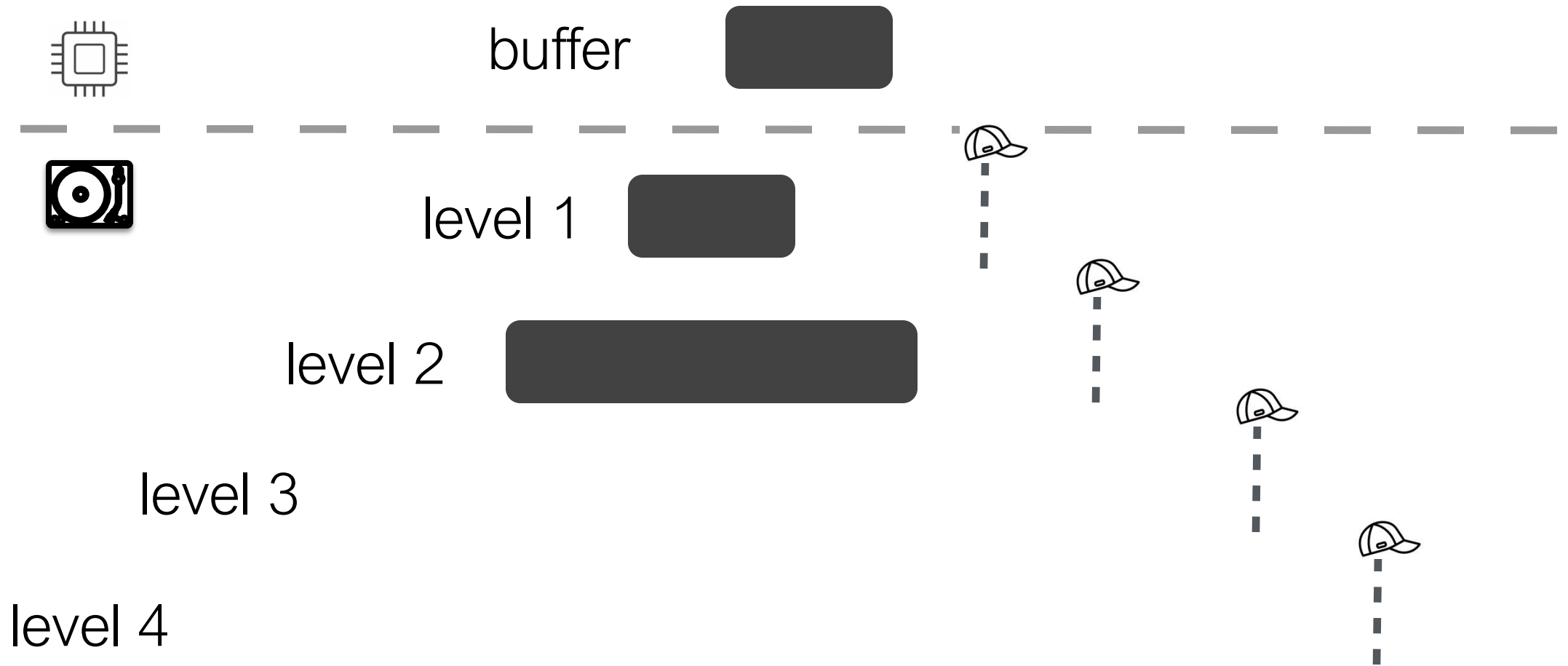
level 2

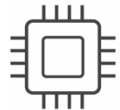


level 3

level 4







buffer



level 1



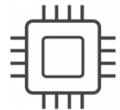
level 2



level 3



level 4



buffer



level 1



level 2

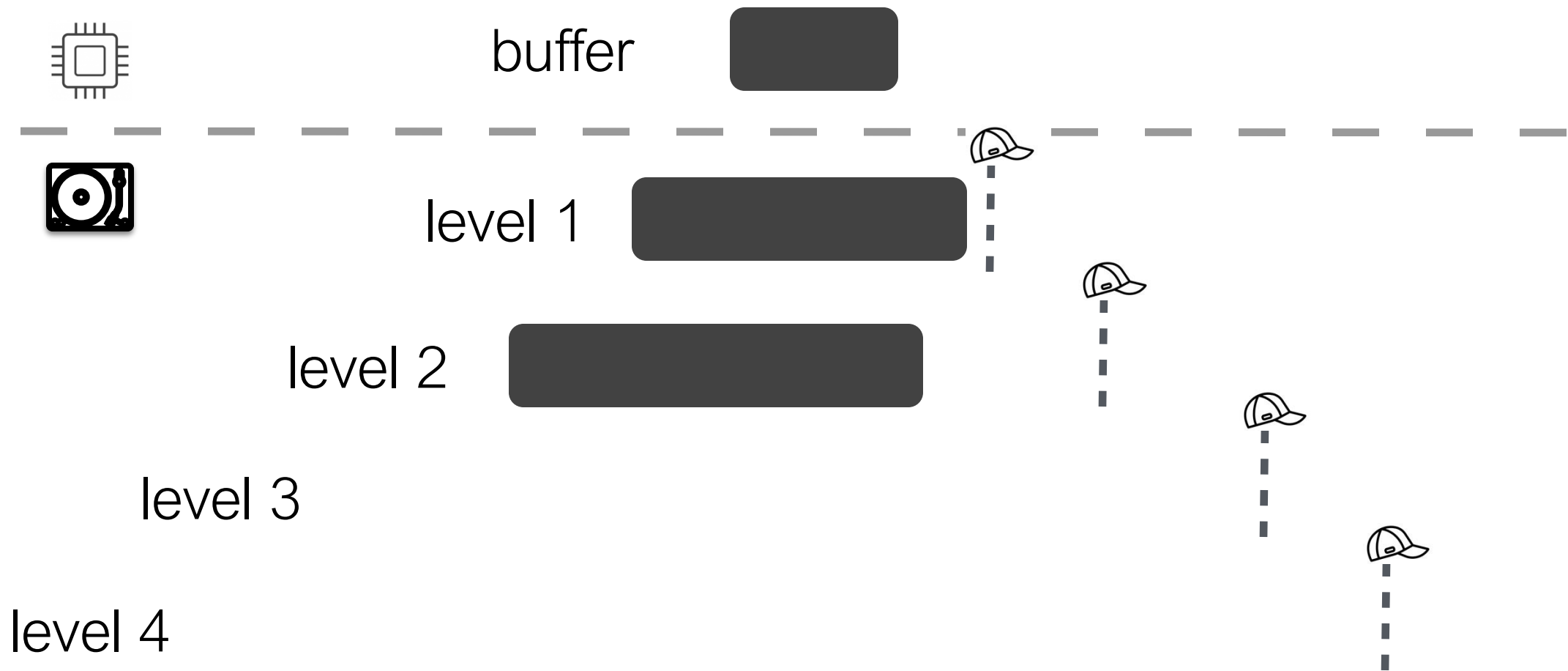


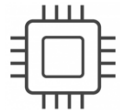
level 3



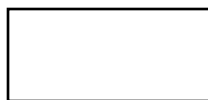
level 4







buffer



level 1



level 2

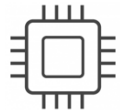


level 3

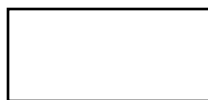


level 4





buffer



level 1



level 2

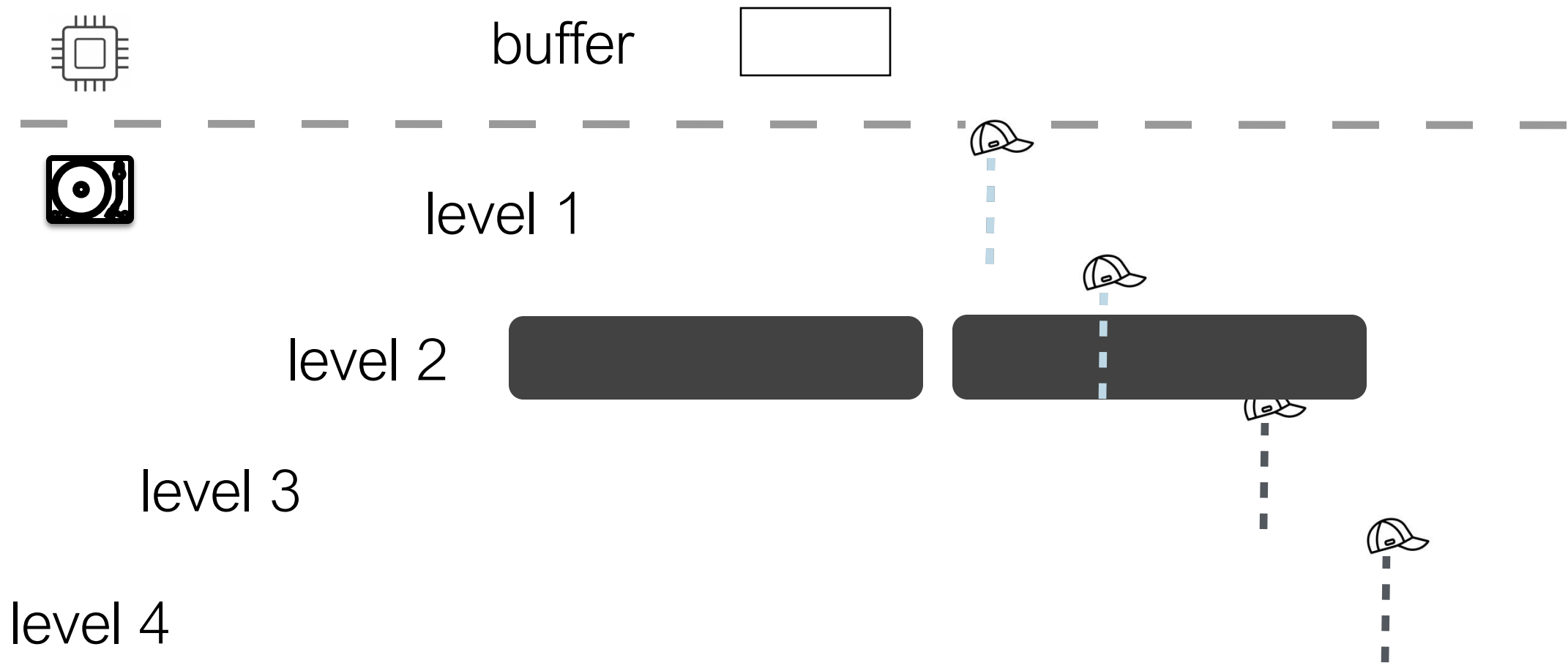


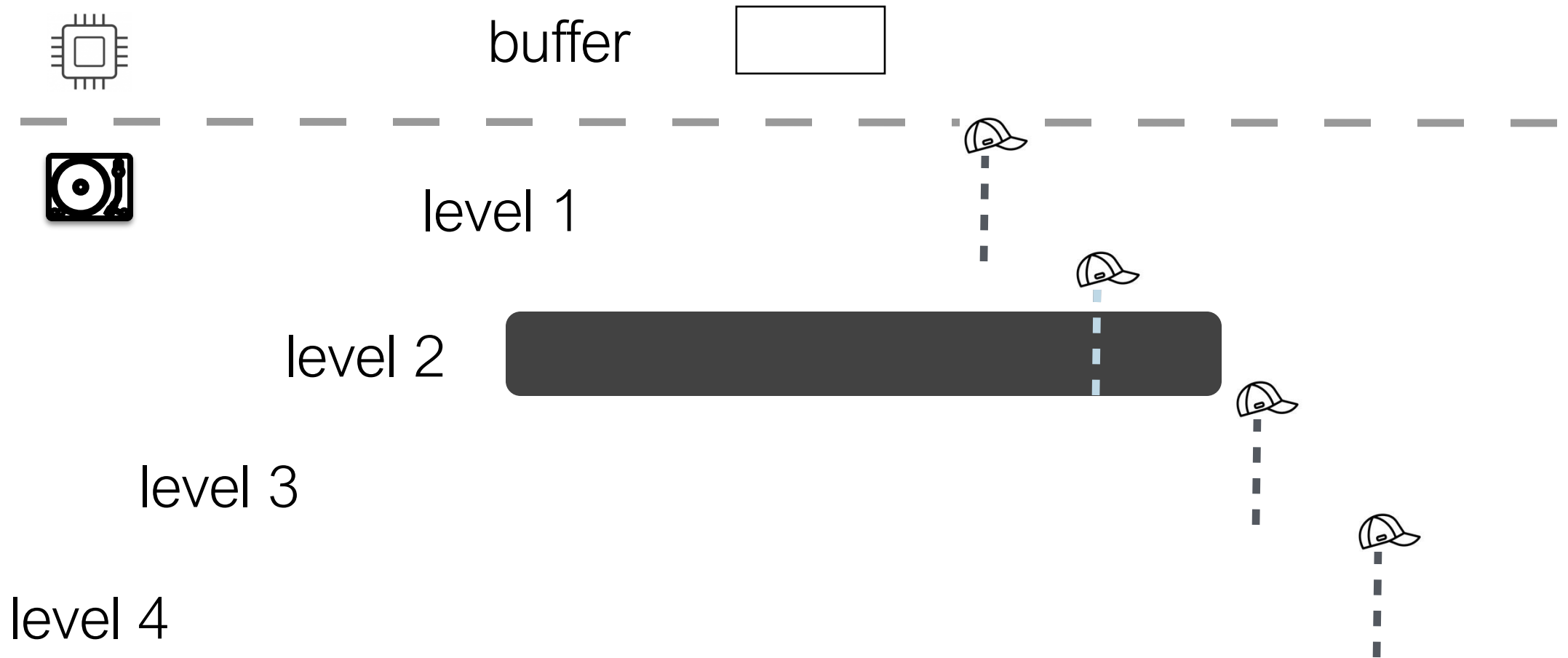
level 3

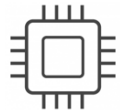


level 4









buffer



level 1



level 2



level 3



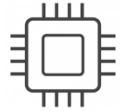
level 4



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

Periodic
compactions

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



buffer

$P \cdot B$

level 1



$P \cdot B \cdot T$

level 2



$P \cdot B \cdot T^2$

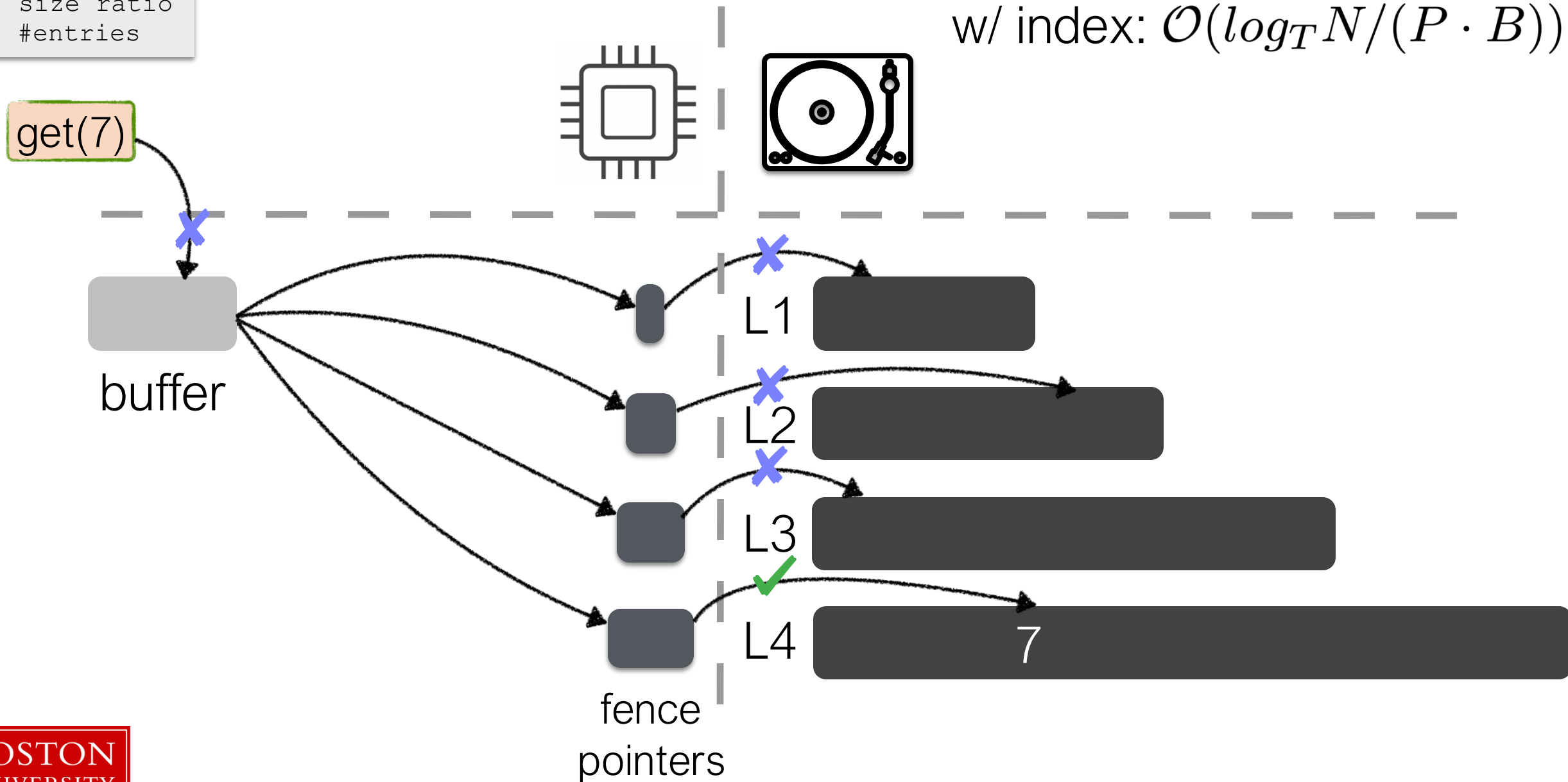
How about queries?

capacity
(entries)

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

Cost analysis

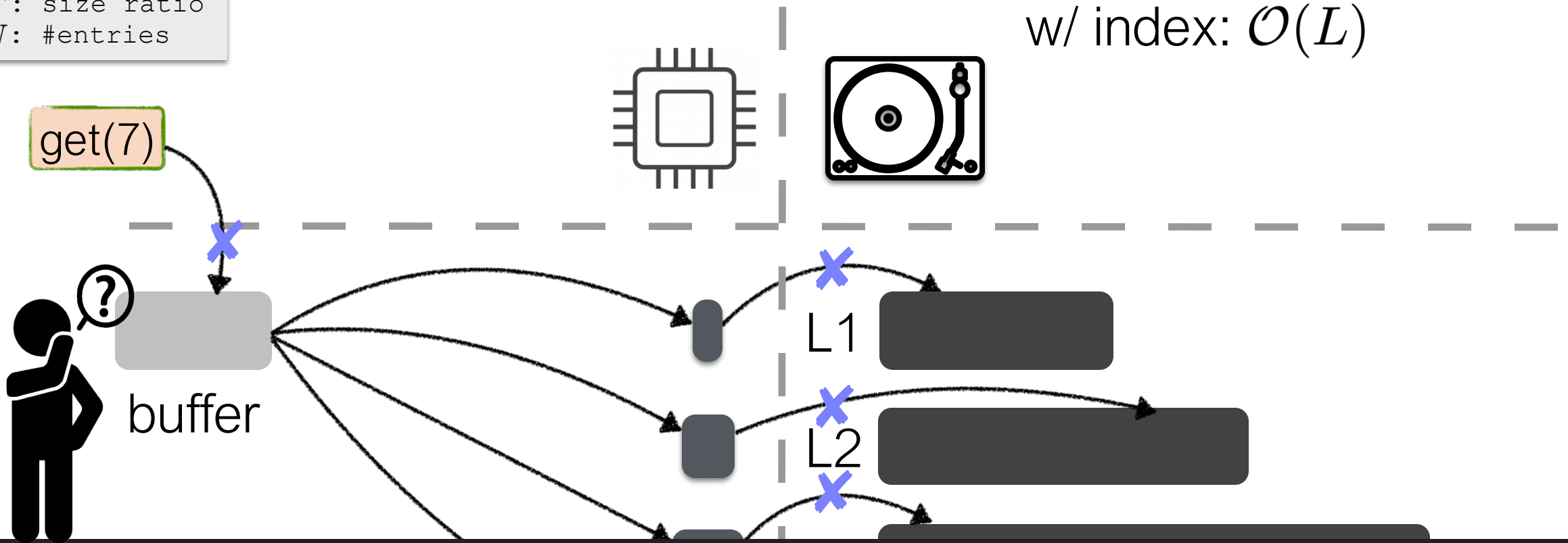
w/ index: $\mathcal{O}(\log_T N / (P \cdot B))$



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

Cost analysis

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



How to avoid unnecessary I/Os?

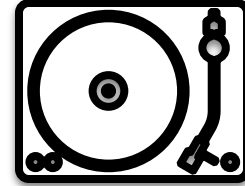
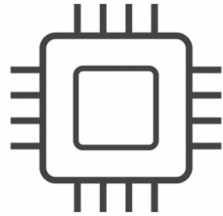
pointers

P : pages in
 B^{buffer}
 $B^{\text{entries/page}}$
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

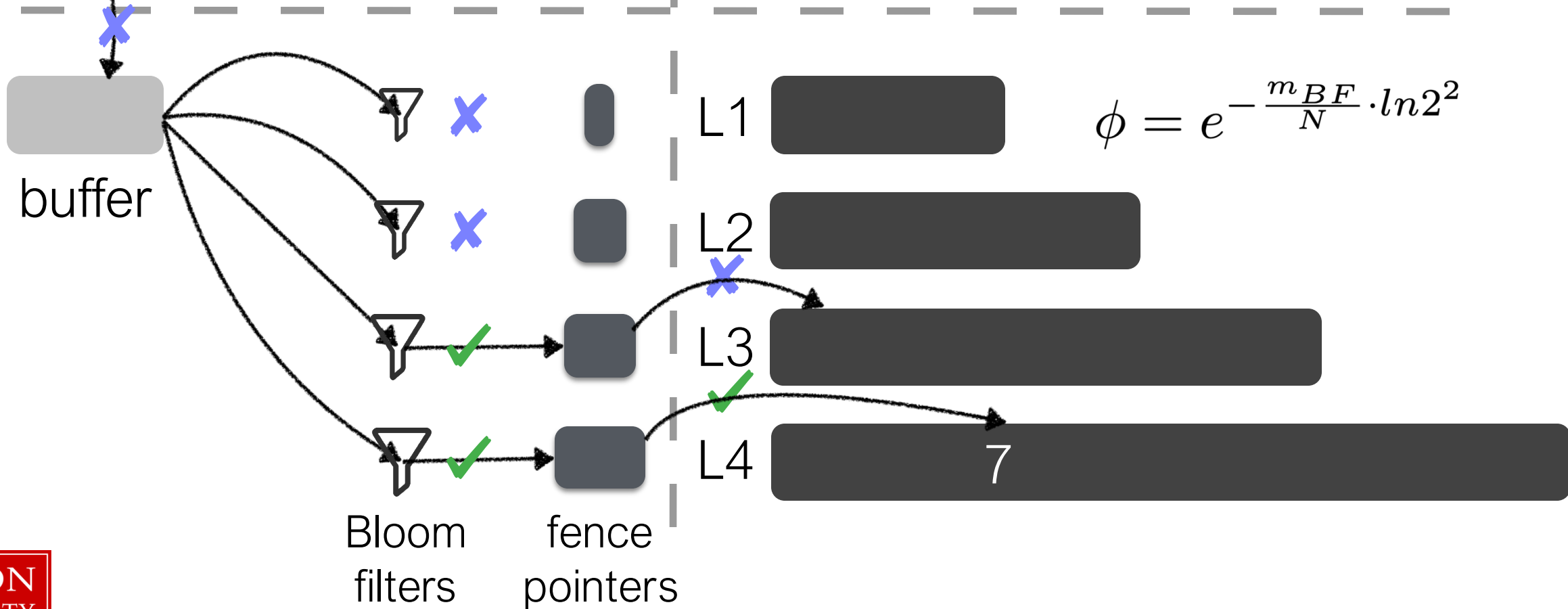
Cost analysis

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

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



get(7)

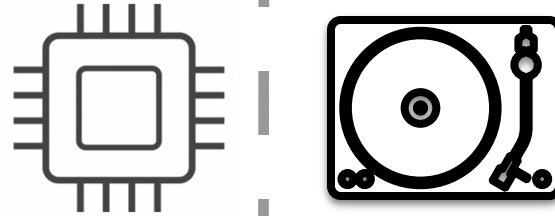


P : pages in
 B buffer
 B entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

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

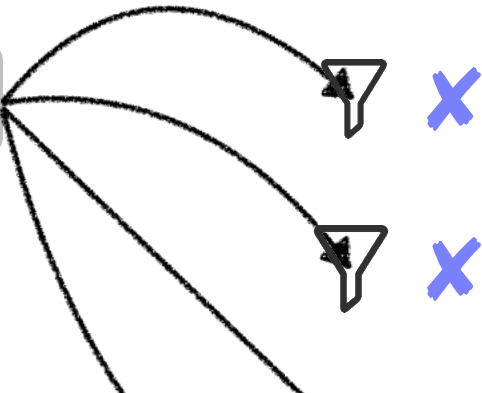
w F&I: $\mathcal{O}(\phi \cdot L)$



get(7)



buffer



L1



L2



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

How to manage memory?

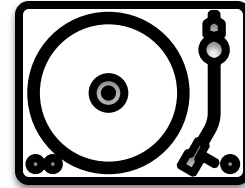
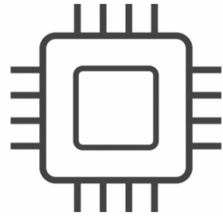
Bloom filters
 pointers

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

Cost analysis

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

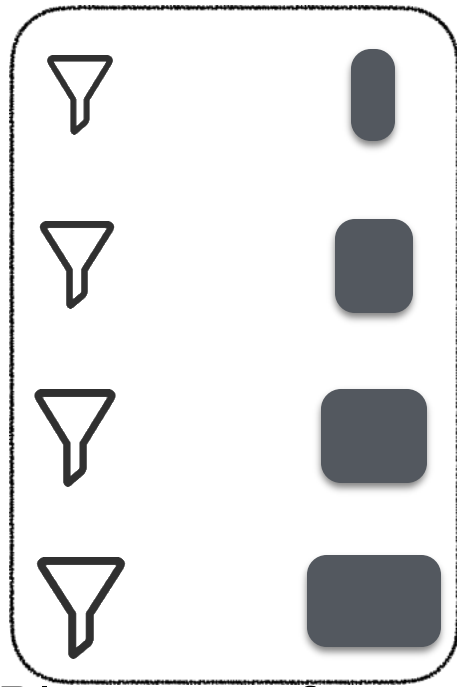
w F&I: $\mathcal{O}(\phi \cdot L)$



buffer



block cache



Bloom
filters

fence
pointers

L1



L2



L3



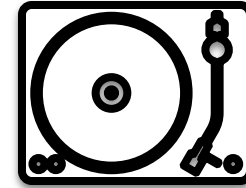
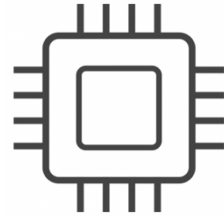
L4



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

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

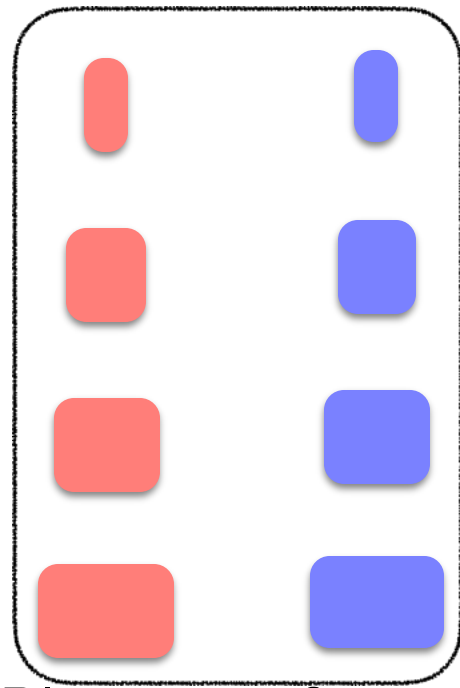
Block Cache



buffer



block cache



Bloom
filters

fence
pointers

L1



L2



L3

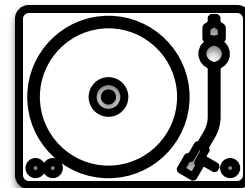



L4



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

Block Cache



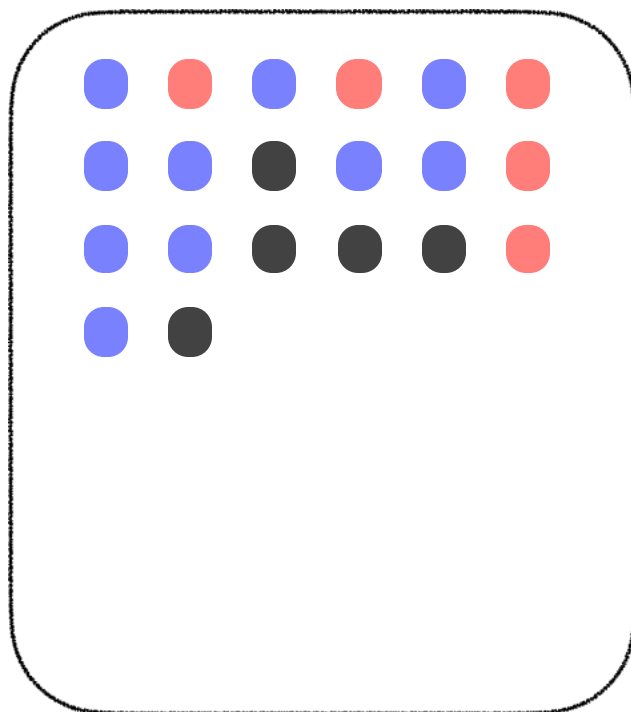

 Bloom
 filters


 fence
 pointers

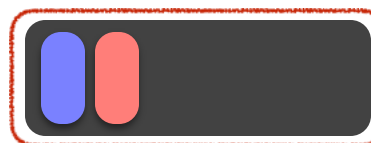
get(7)



buffer



L1



L2



L3



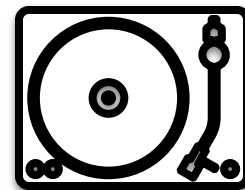
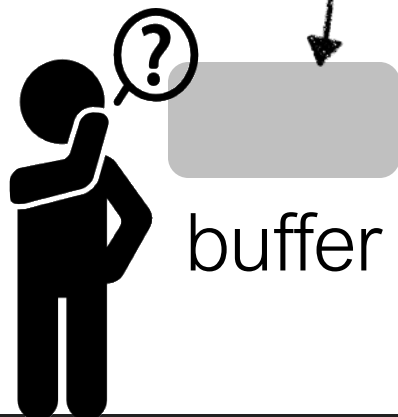
L4




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

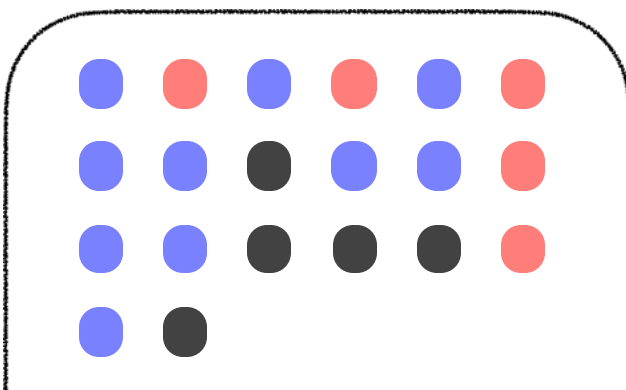
Block Cache

get(7)




 Bloom
 filters


 fence
 pointers



L1



L2

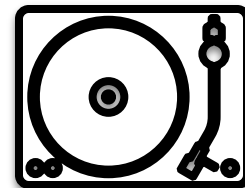
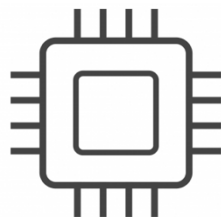


What about range queries?

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

s selectivity
 LRQ

Range Queries



buffer



L1



L2



L3



L4



Bloom
filters

fence
pointers

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

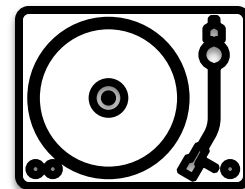
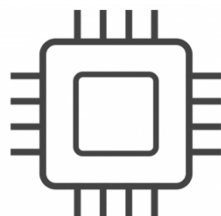
s selectivity
 LRQ

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

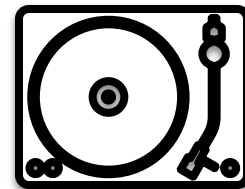
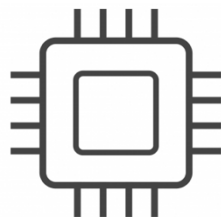
s selectivity
 SRQ

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



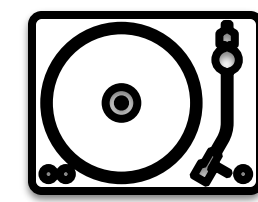
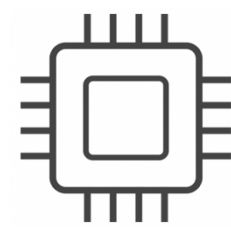
L4



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

Going back to point queries

Cost analysis

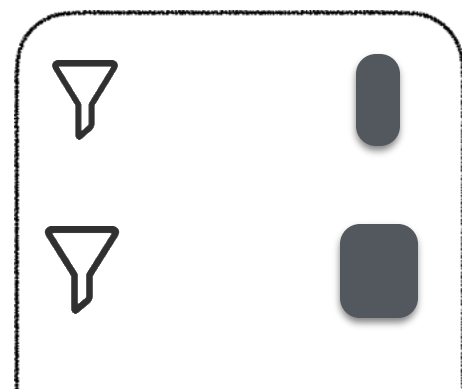


w/ index: $\mathcal{O}(L)$

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



buffer



L1



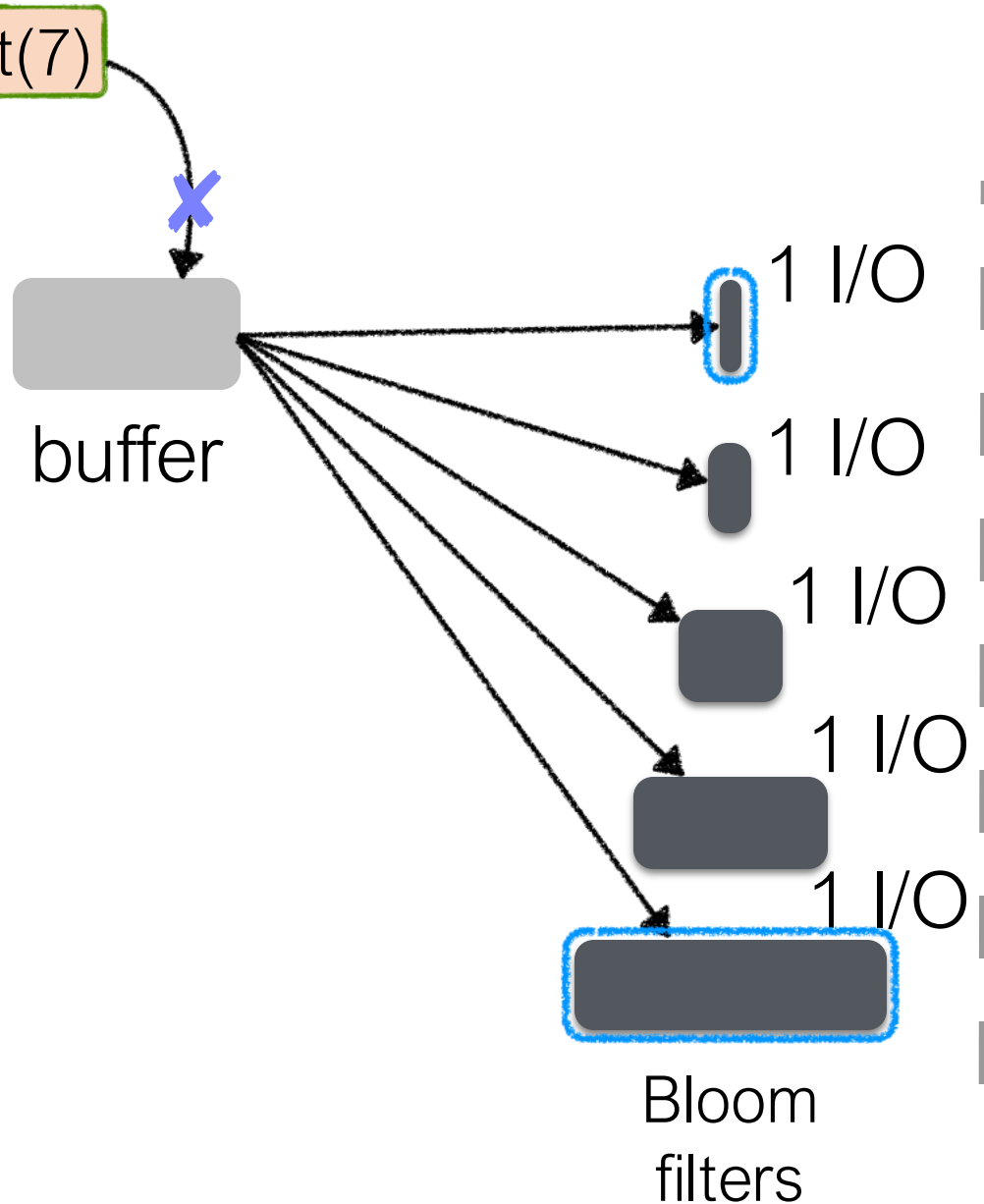
L2



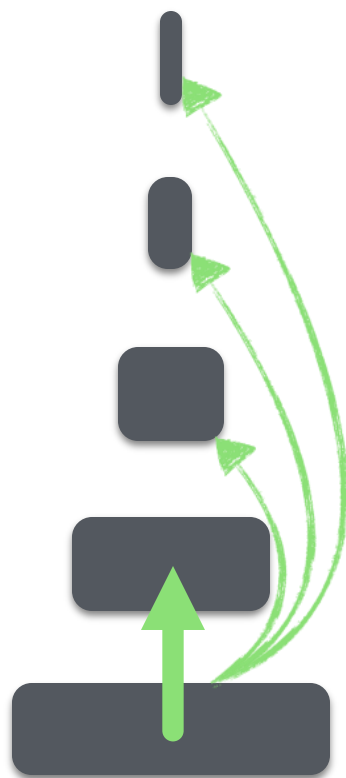
$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2} \ll 1$$

Should all BFs be equally accurate?

Insight: most of the memory allocated to BFs is helping us save only 1 I/O (90% for T=10)



Bloom
filters



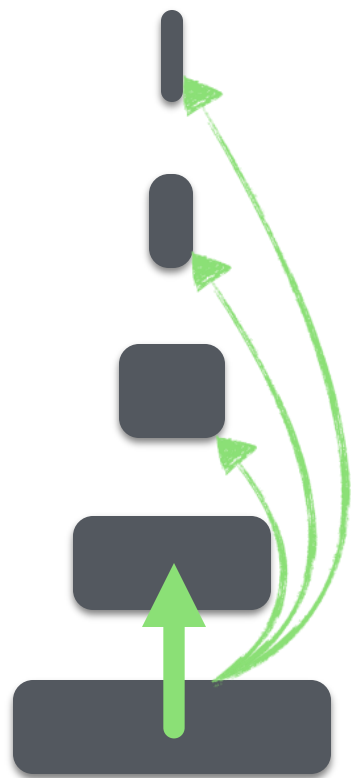
FPR Monkey
FPR

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

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

point lookup
cost

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

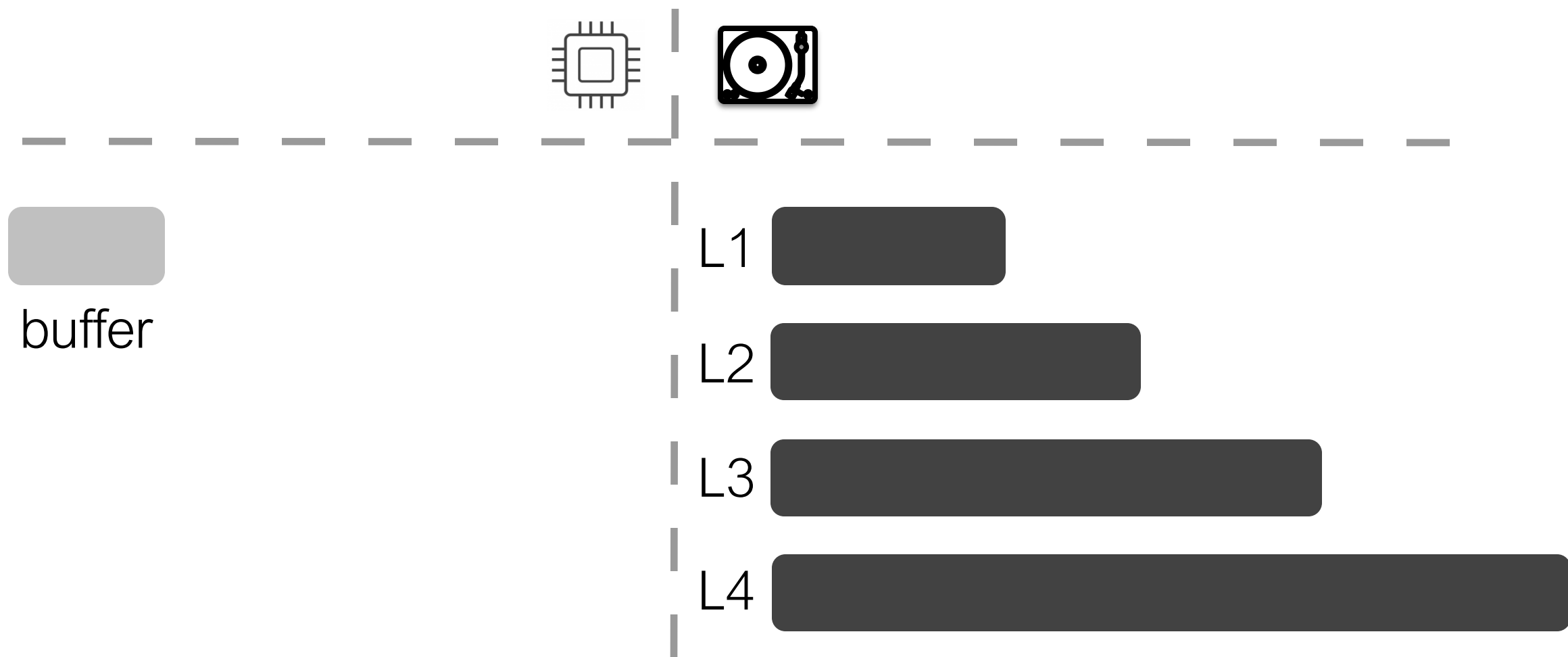
\checkmark

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



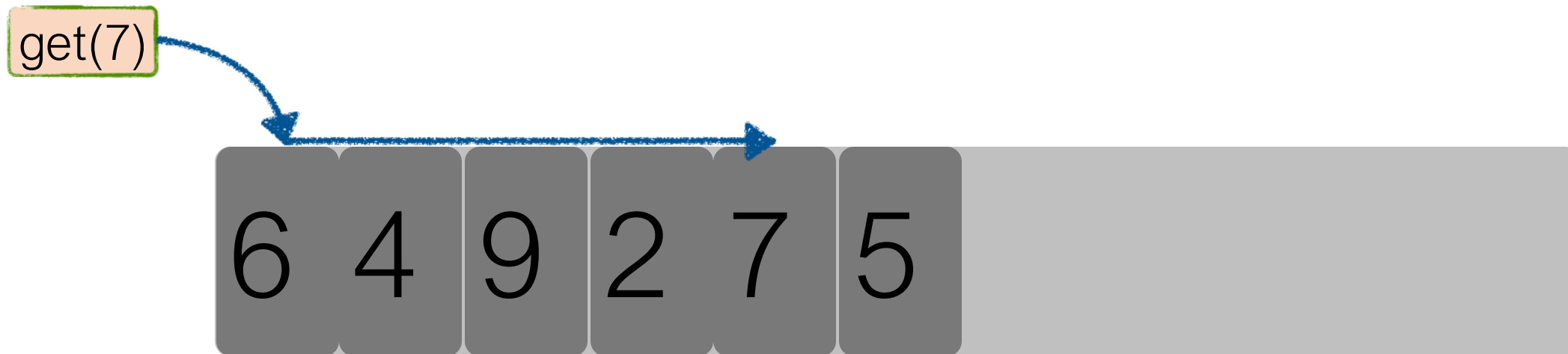
Buffer Optimizations

Buffer Optimizations



P : pages in
 B buffer
entries/page

Buffer Implementation: **vector**



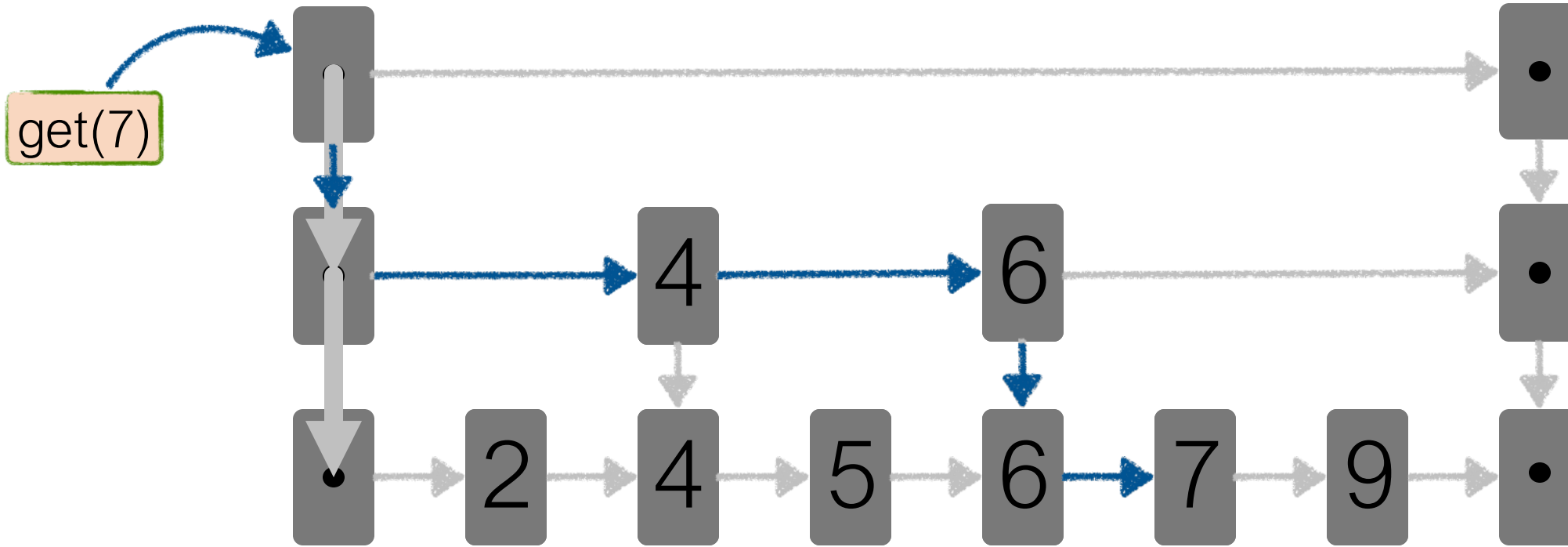
- great for ingestion-heavy w/l
- no extra space needed
- expensive points queries

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
 B buffer
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

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

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

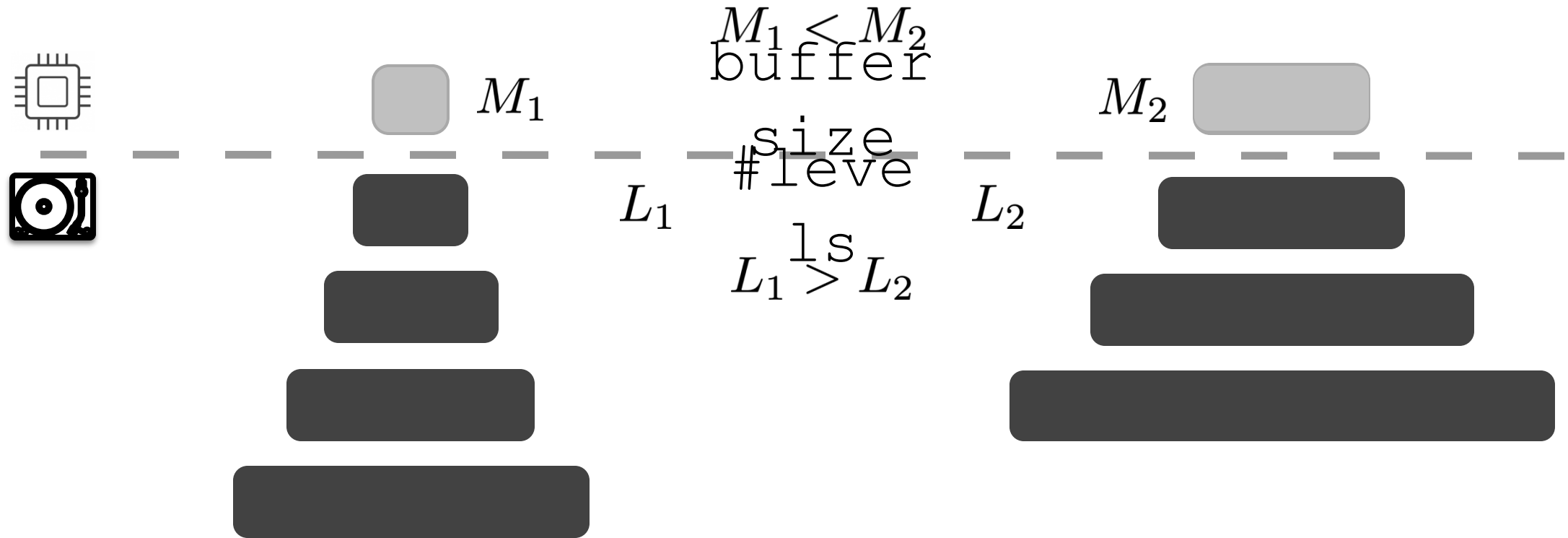
$$\mathcal{O}(1)$$

Ingestion-
only
workloads

Mixed
workloads

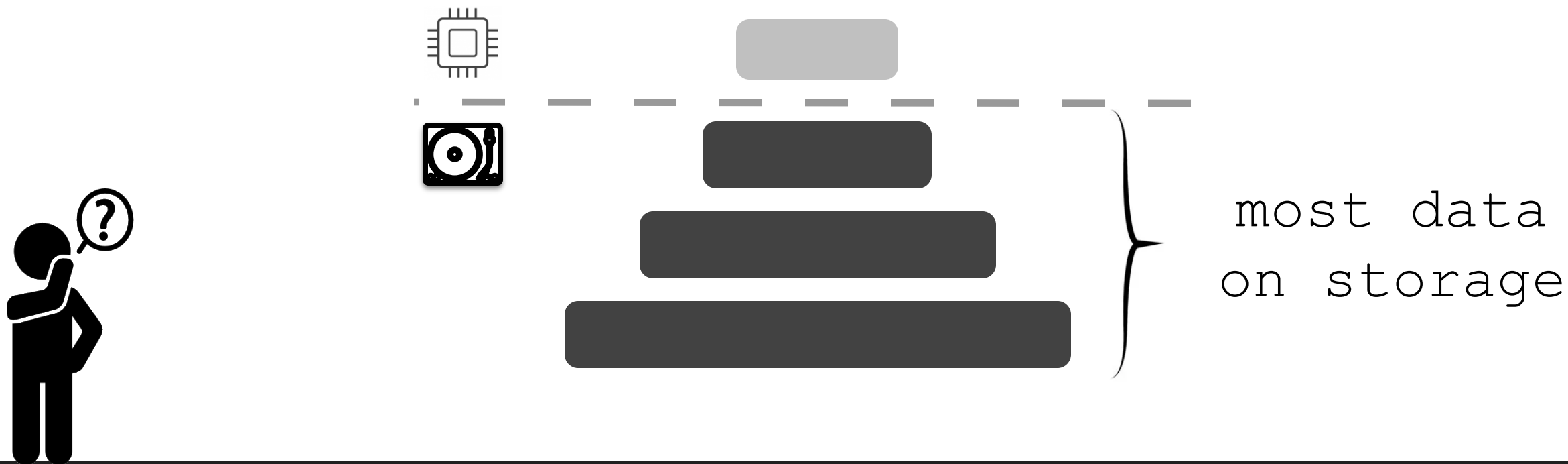
I/O-bound
workloads

Size of the Buffer



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

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



How does the storage layer affect ingestion?

L : #levels
 T : size ratio

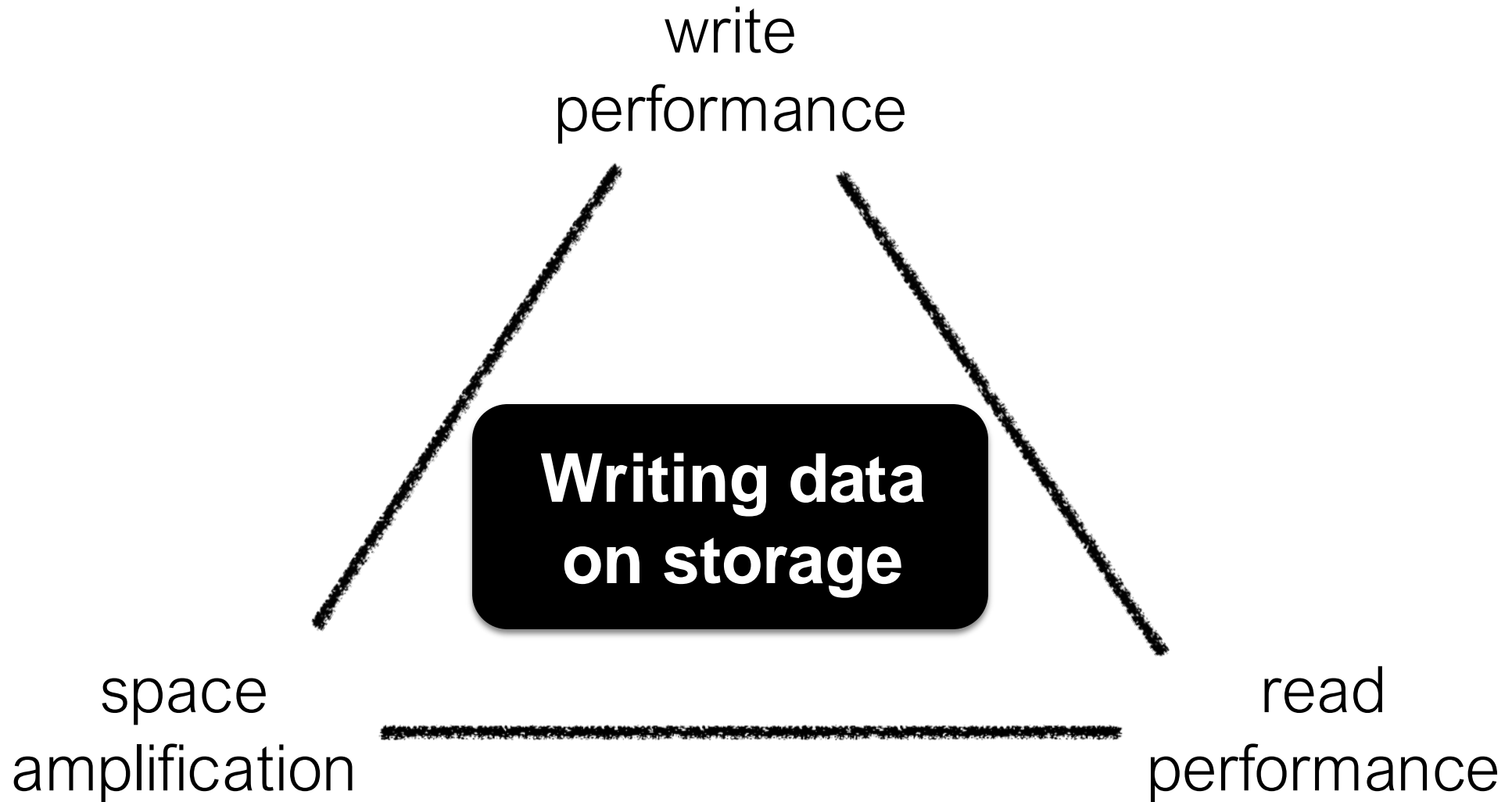


most data
on storage
if $T = 10$ & $L = 4$

99.9% on storage

How does the storage layer affect ingestion?

Storage Optimizations



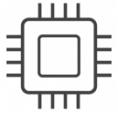
Data **Layout**

Classical LSM design: `leveling`
`g`
[eager merging]



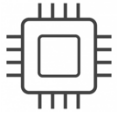
Data Layout

leveling [eager]



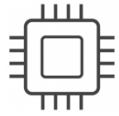
Data Layout

leveling [eager]



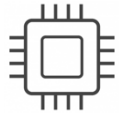
Data Layout

leveling [eager]



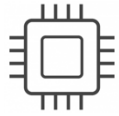
Data Layout

leveling [eager]



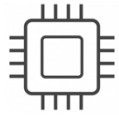
Data Layout

leveling [eager]

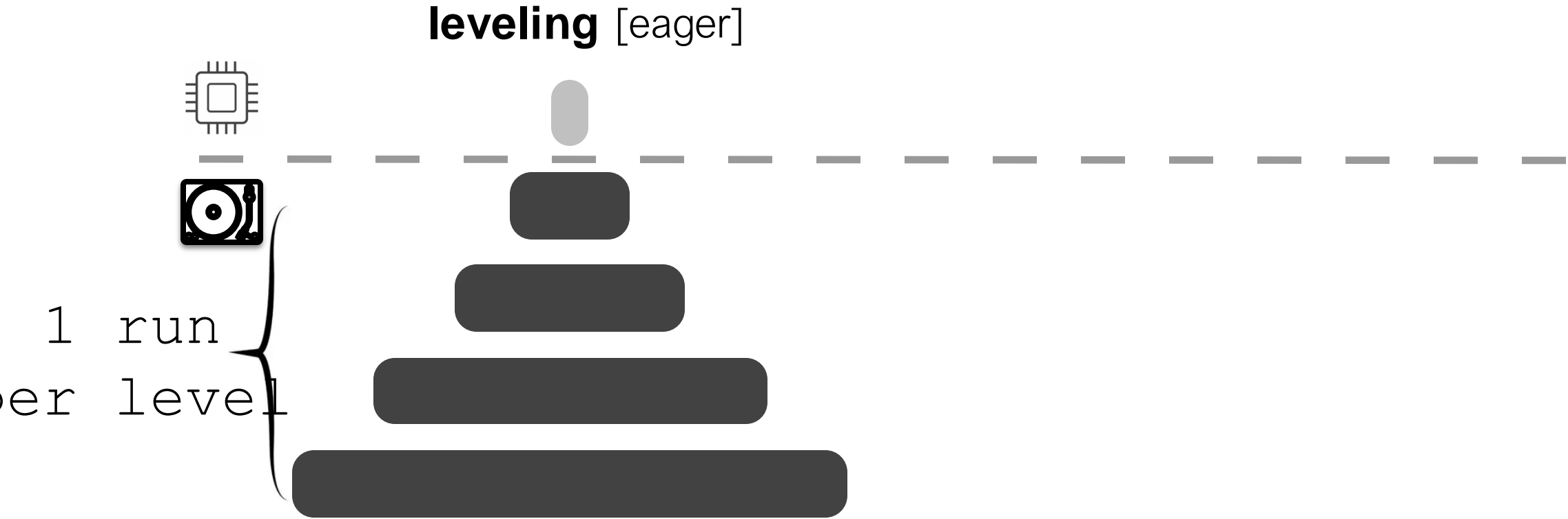


Data Layout

leveling [eager]

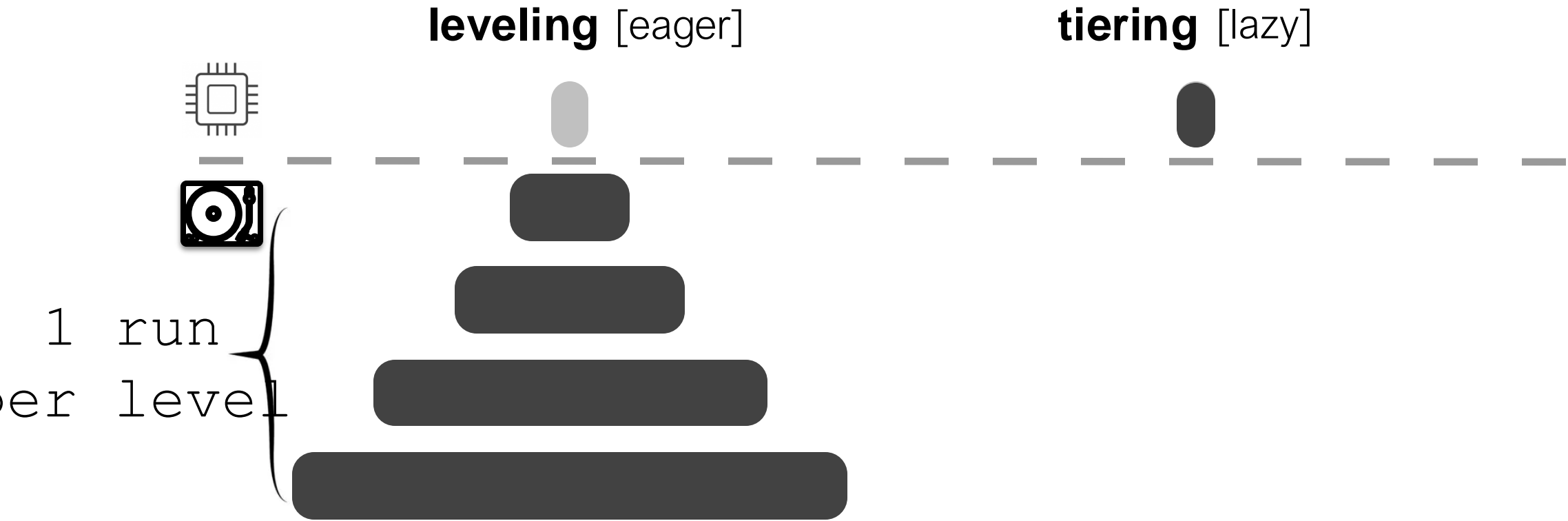


Data Layout



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

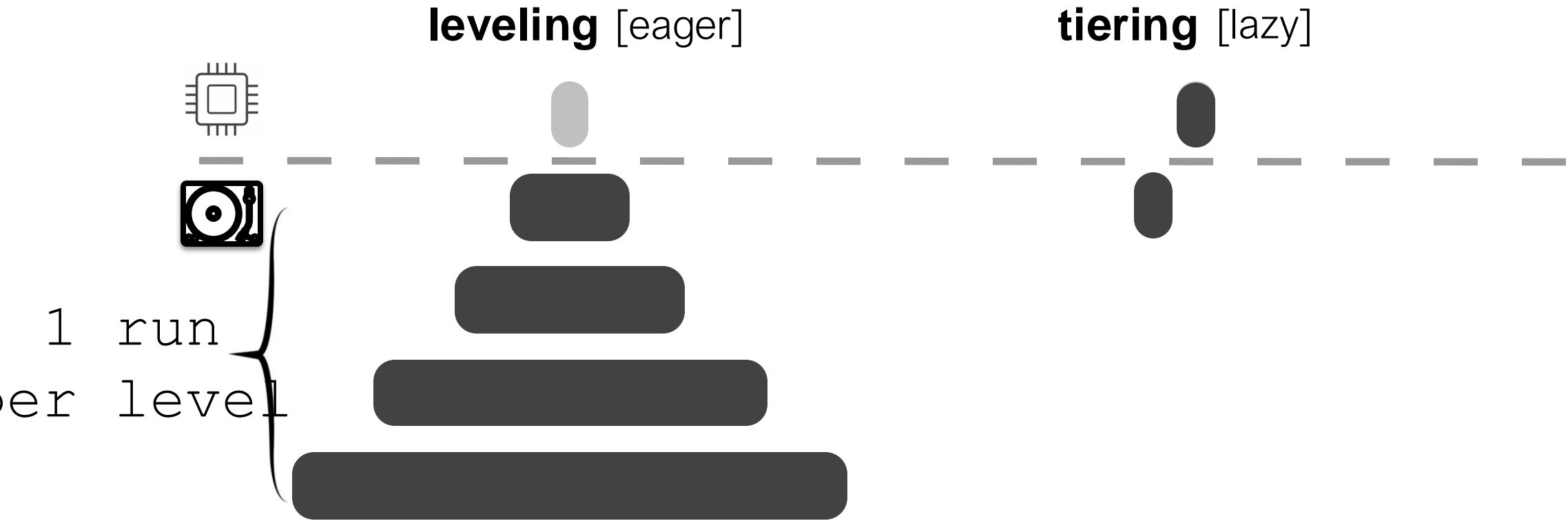
Data Layout



1 run
per level

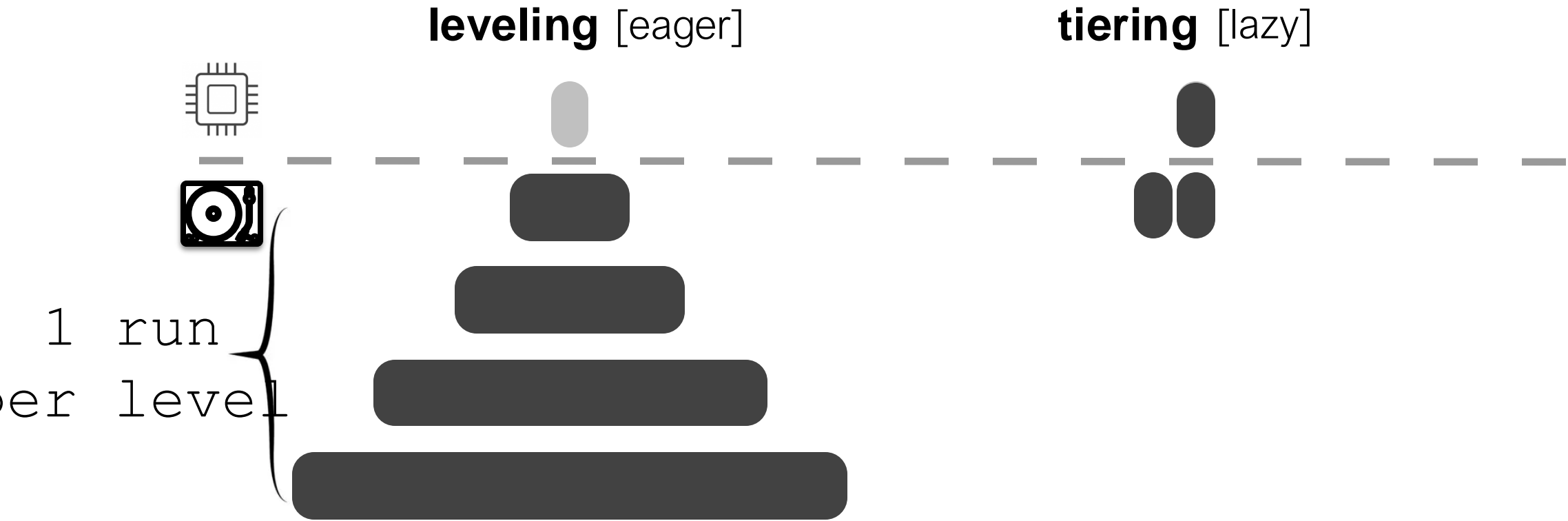
- good read performance
- good space amplification
- high write amplification

Data Layout



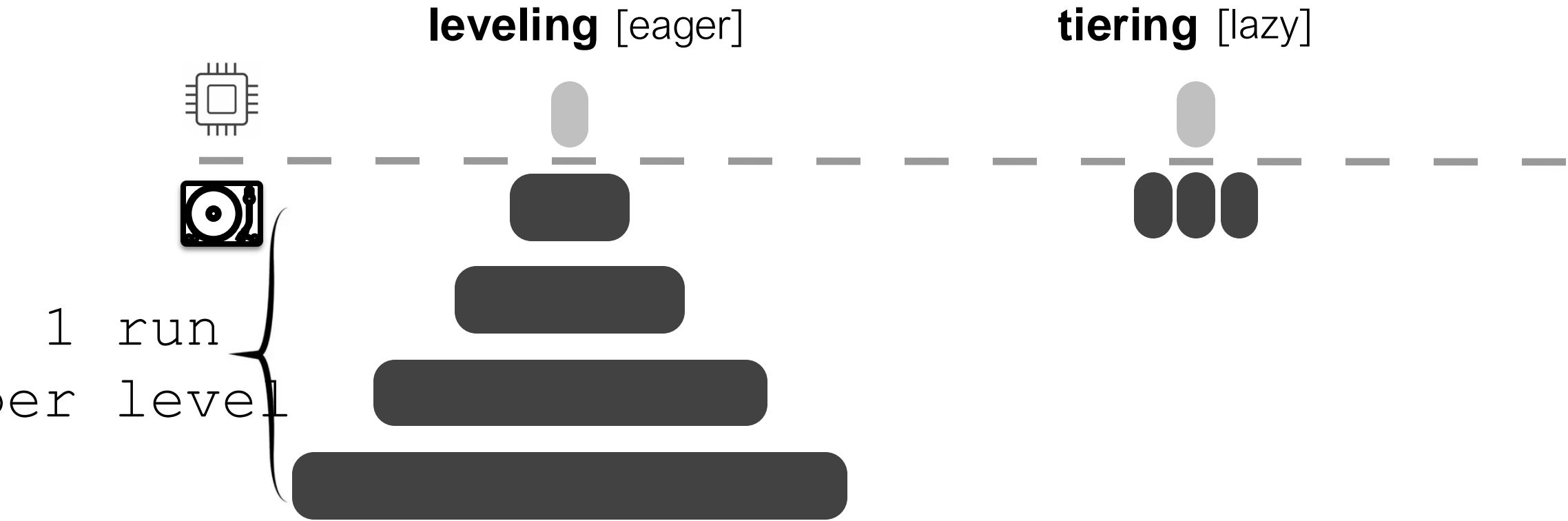
- good read performance
- good space amplification
- high write amplification

Data Layout



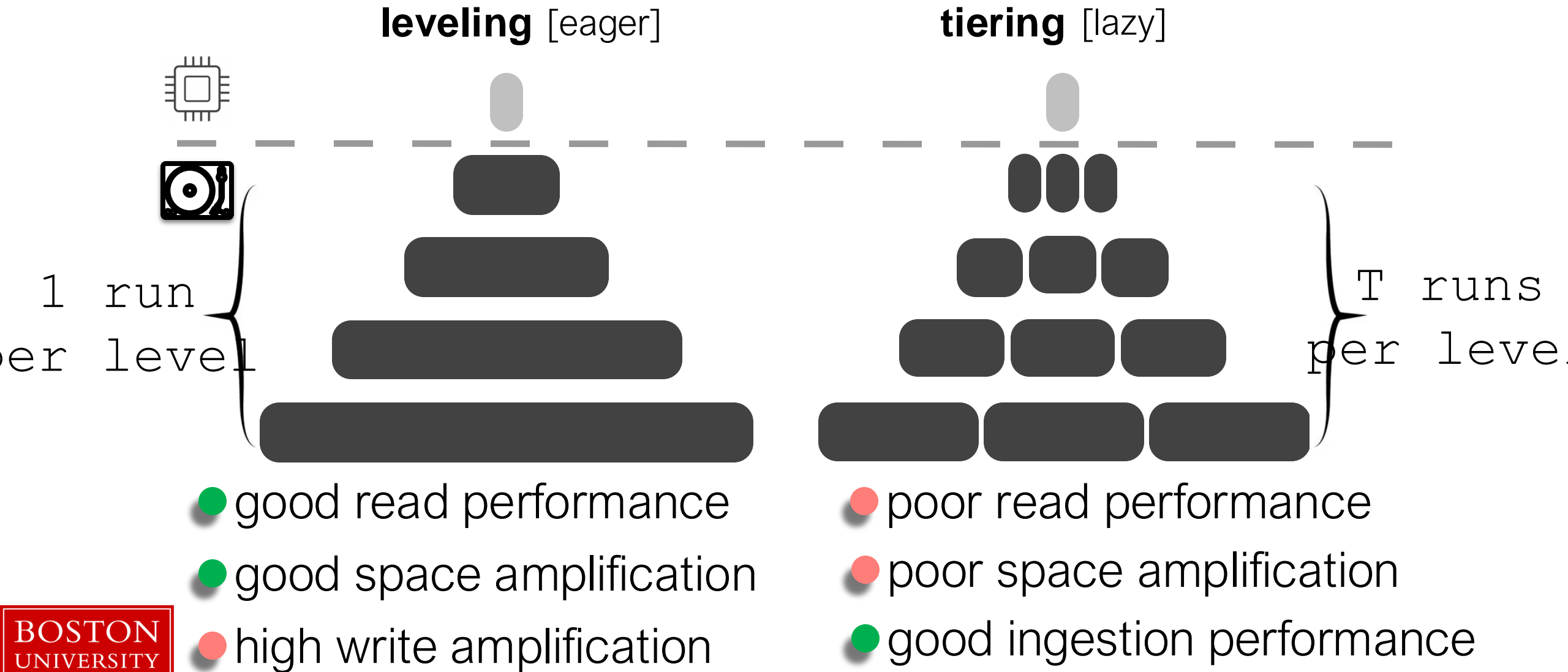
- good read performance
- good space amplification
- high write amplification

Data Layout



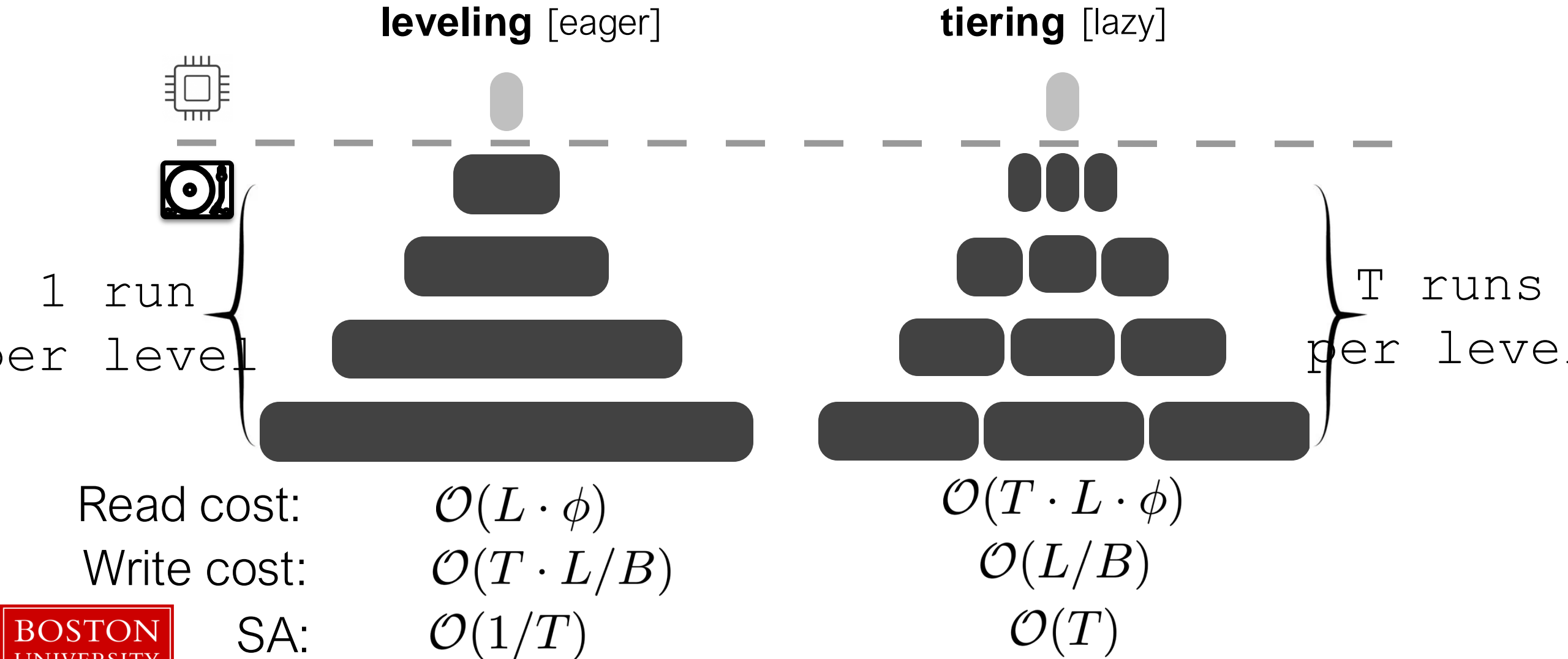
- good read performance
- good space amplification
- high write amplification

Data Layout



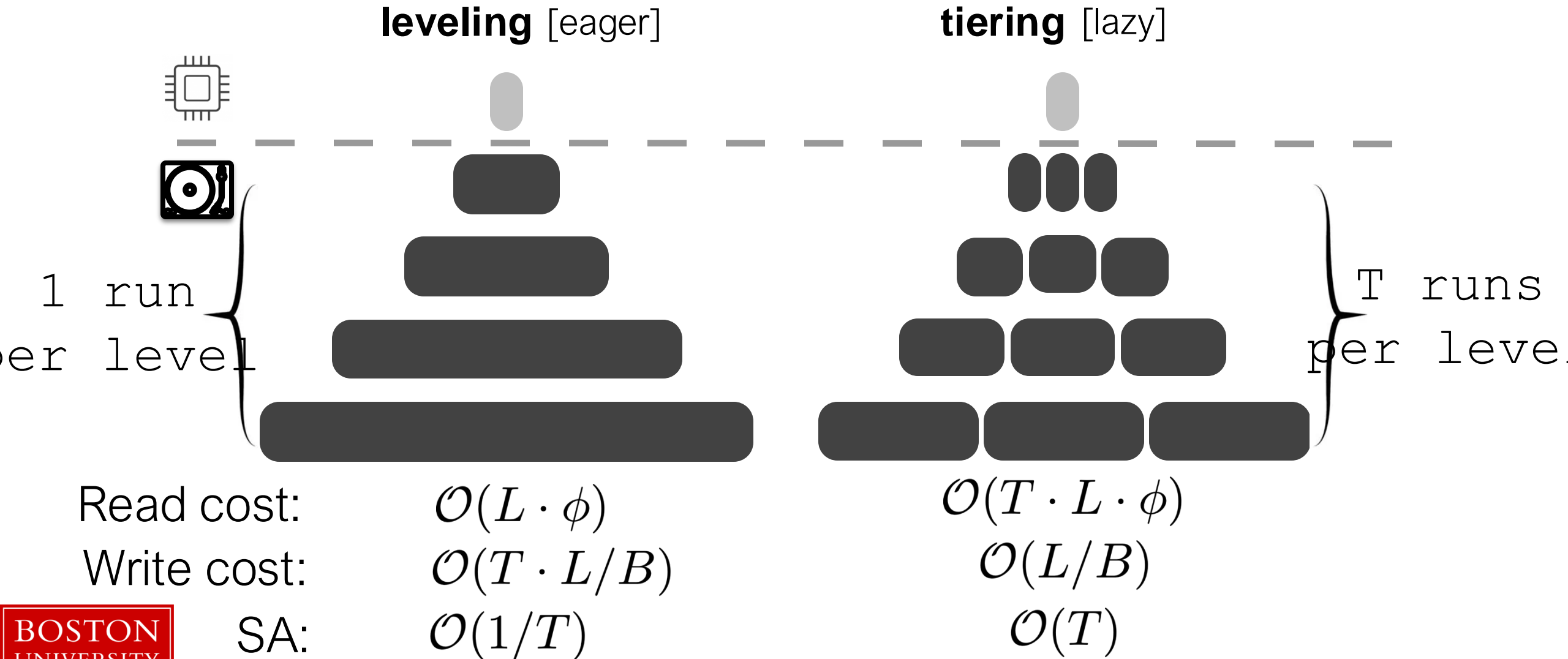
P : pages in
 buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Data Layout



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

Data Layout



Data Layout

hybrid
designs

leveling



tiering



read

optimized



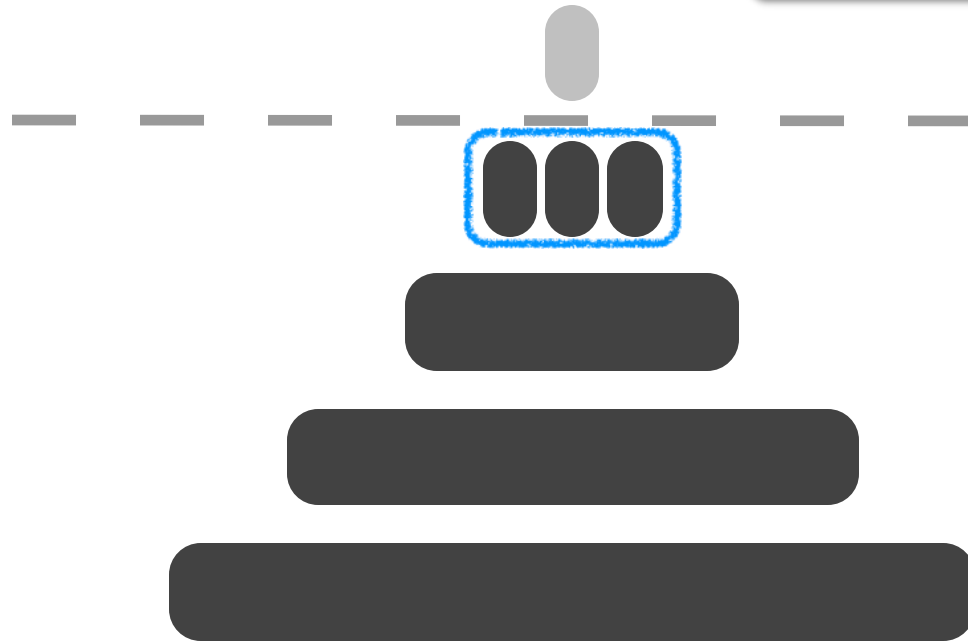
write

optimized

Data Layout

1-leveling

RocksDB20 

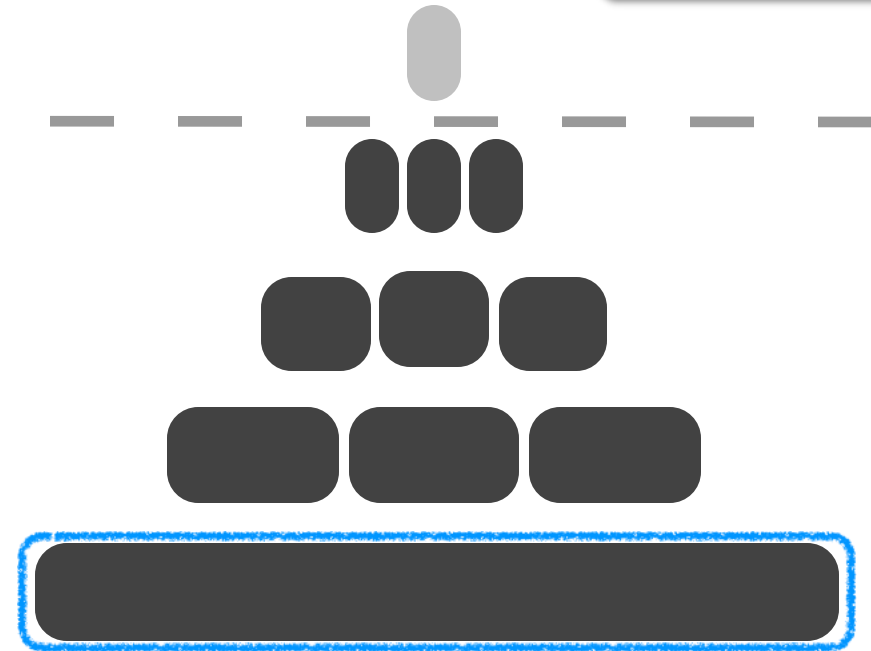


- fewer write stalls
- increased block cache hits

compared to leveling!

L-leveling

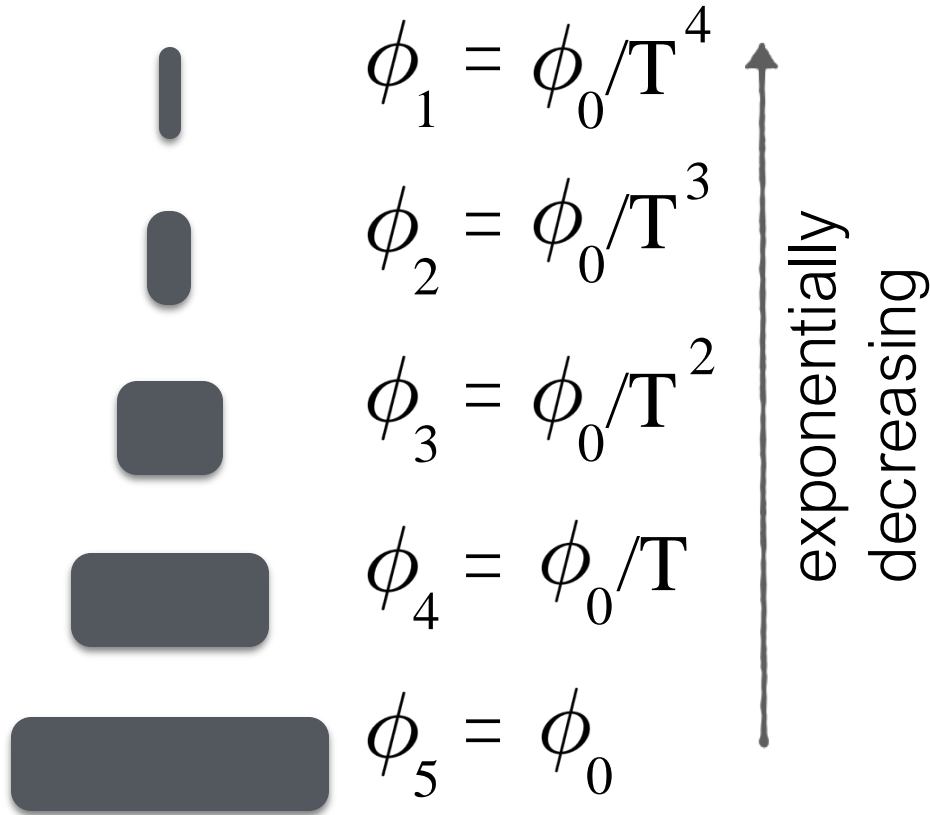
DayanSIGMOD18 



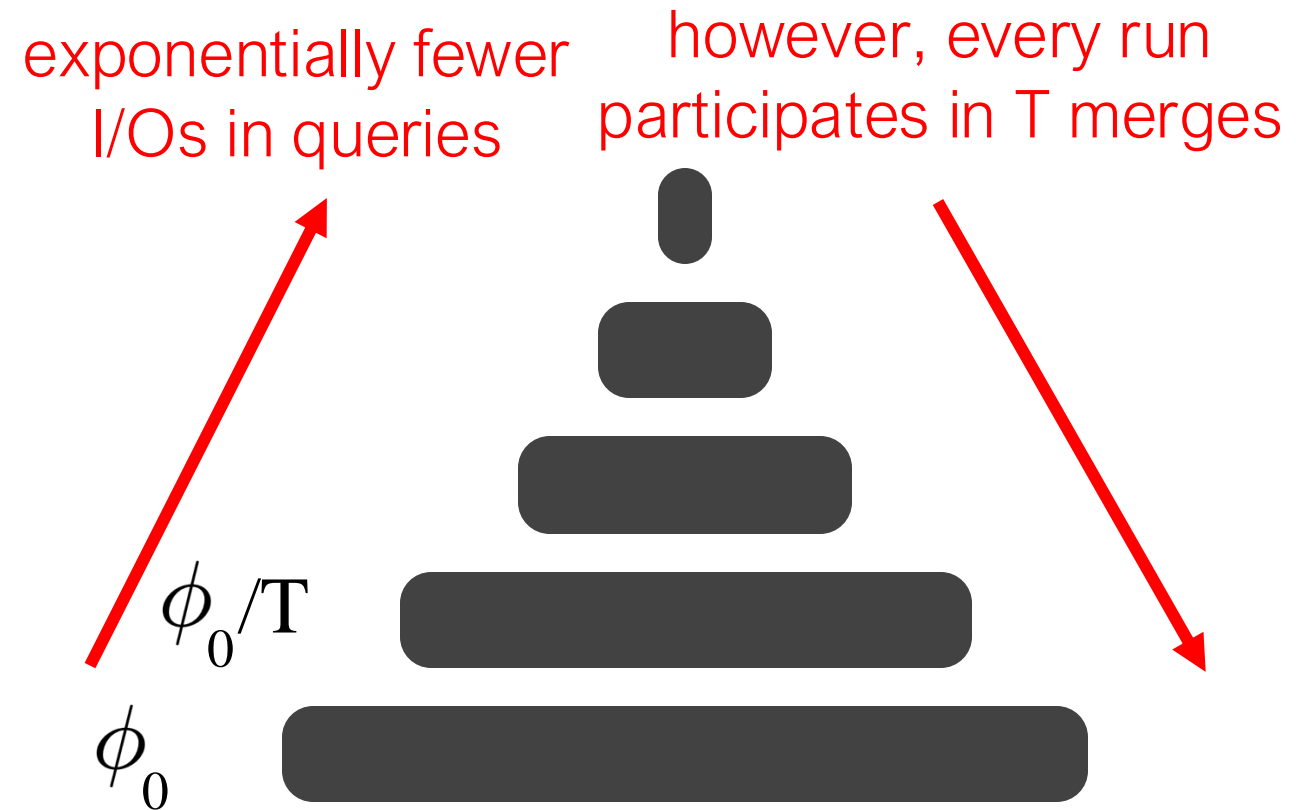
- low write amplification
- better read performance

compared to tiering!

Bloom
filters

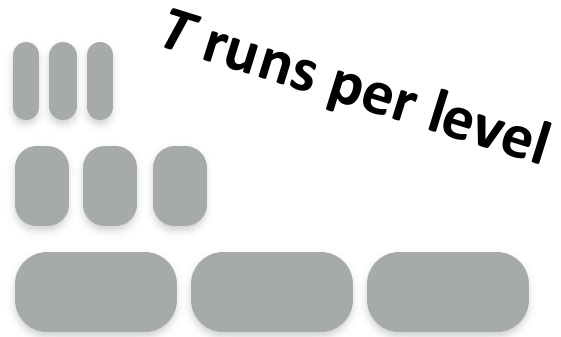


LSM-tree

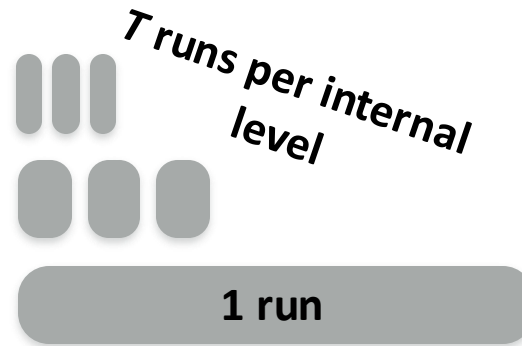


Insight: merging in top levels can
become lazy (tiering)!

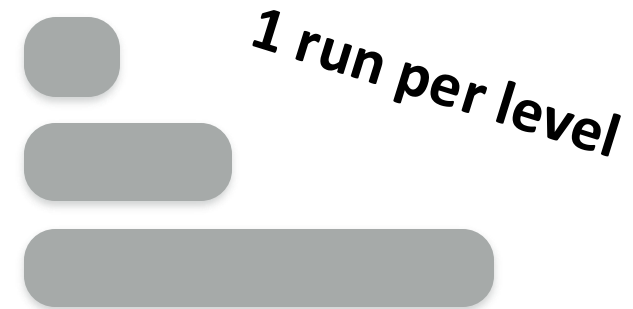
Tiering write-optimized



Lazy Leveling hybrid



Leveling read-optimized



lookup:

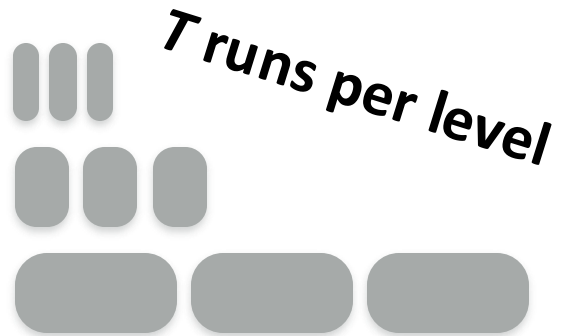
$$O(T \cdot \log_T(N) \cdot e^{-M/N})$$

runs per level levels false positive rate

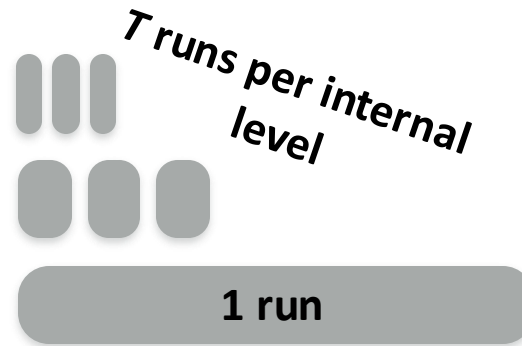
$$O(\log_T(N) \cdot e^{-M/N})$$

levels false positive rate

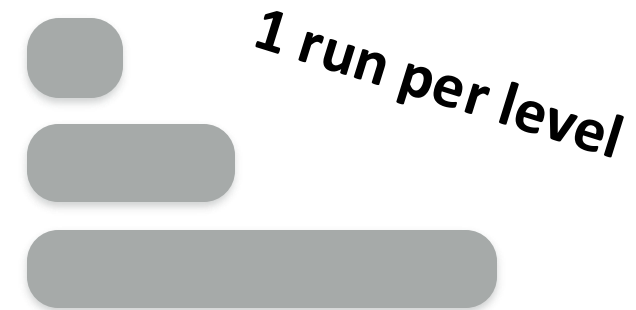
Tiering write-optimized



Lazy Leveling hybrid



Leveling read-optimized



lookup:

$$O(T \cdot \log_T(N) \cdot e^{-M/N})$$

runs per level levels false positive rate

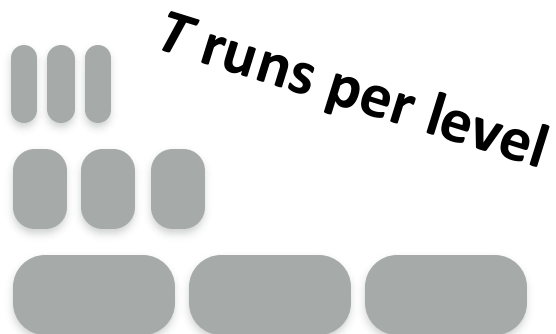
$$O(\log_T(N) \cdot e^{-M/N})$$

levels false positive rate

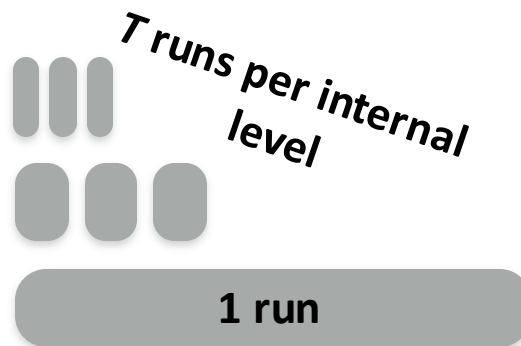


DayanSIGMOD17

Tiering
write-optimized



Lazy Leveling
hybrid



Leveling
read-optimized



lookup:

$$O(T \cdot e^{-M/N})$$

runs per level false positive rate

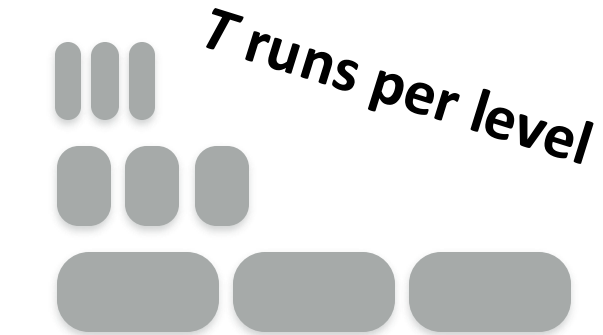
$$O(e^{-M/N})$$

 false positive rate

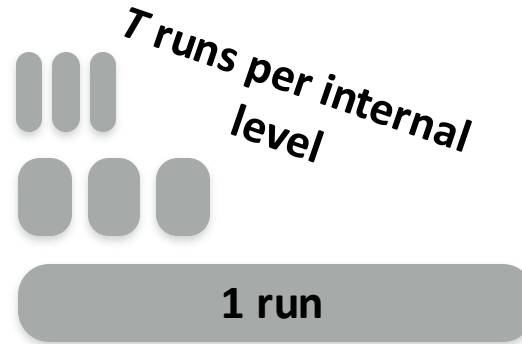
Tiering write-optimized



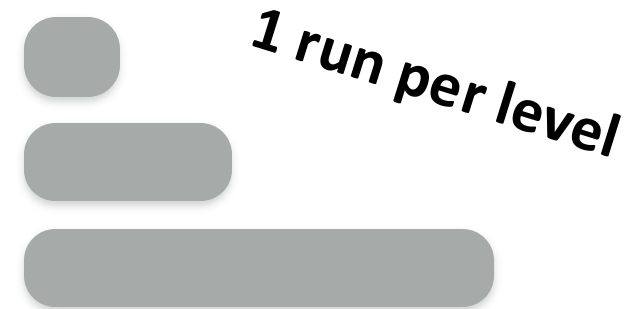
DayanSIGMOD18



Lazy Leveling hybrid



Leveling read-optimized



lookup:

$$O(T \cdot e^{-M/N})$$

runs
per level

false
positive rate

$$O(e^{-M/N})$$

there is a slightly higher
constant factor than leveling

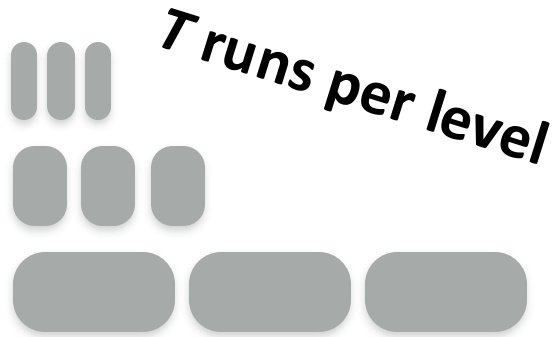
$$O(e^{-M/N})$$

false
positive rate

Tiering write-optimized



DayanSIGMOD18



lookup:

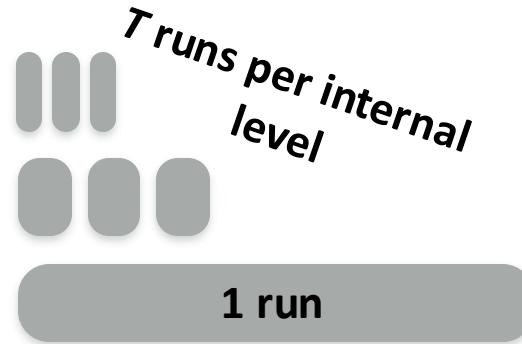
$$O(T \cdot e^{-M/N})$$

update cost:

$$O(\log_T(N))$$

↑
levels

Lazy Leveling hybrid



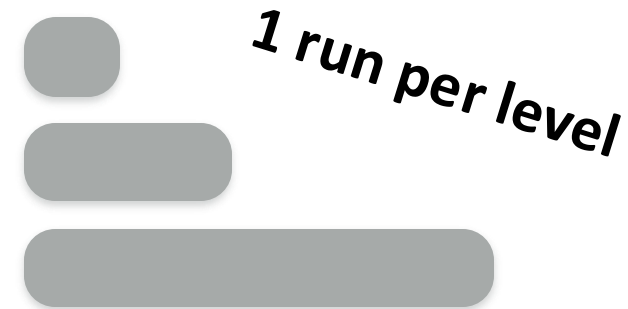
$$O(e^{-M/N})$$

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

↑
T merges for
the last level

↑
1 merge for all
internal levels

Leveling read-optimized



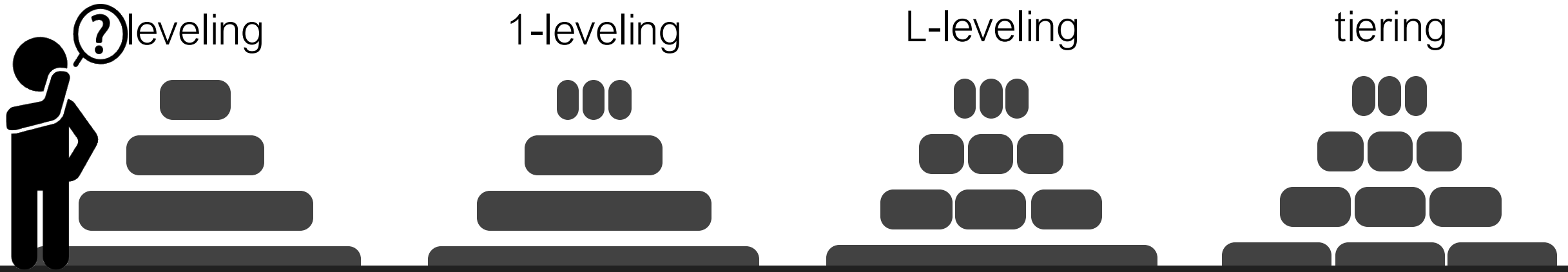
$$O(e^{-M/N})$$

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

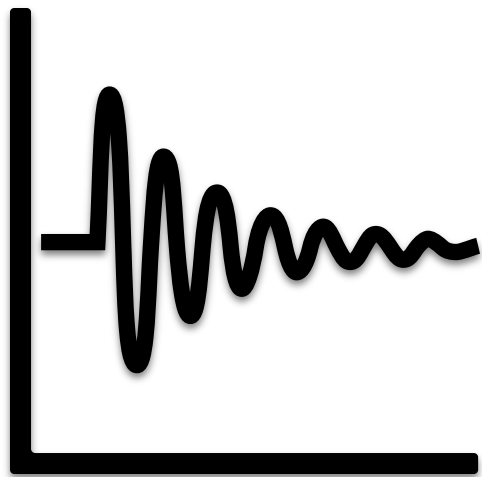
↑
merges per level

↑
levels

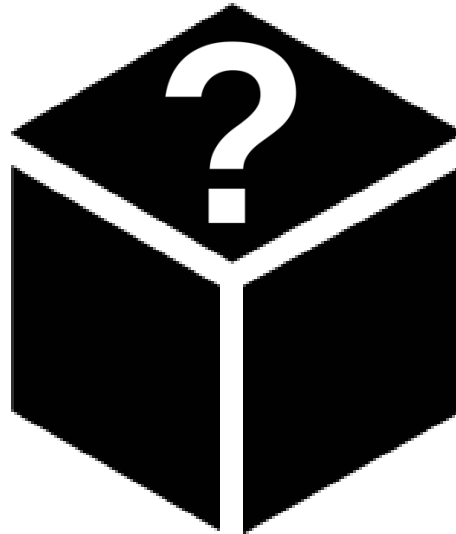
Data Layout



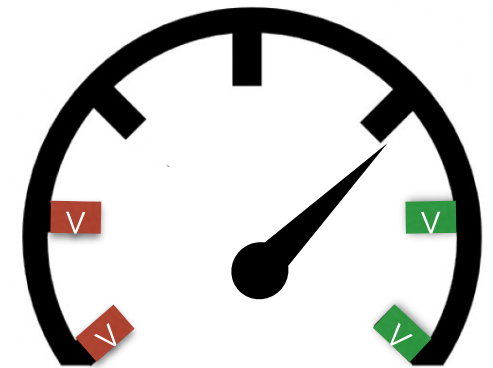
So, how do we reason about the data layout?



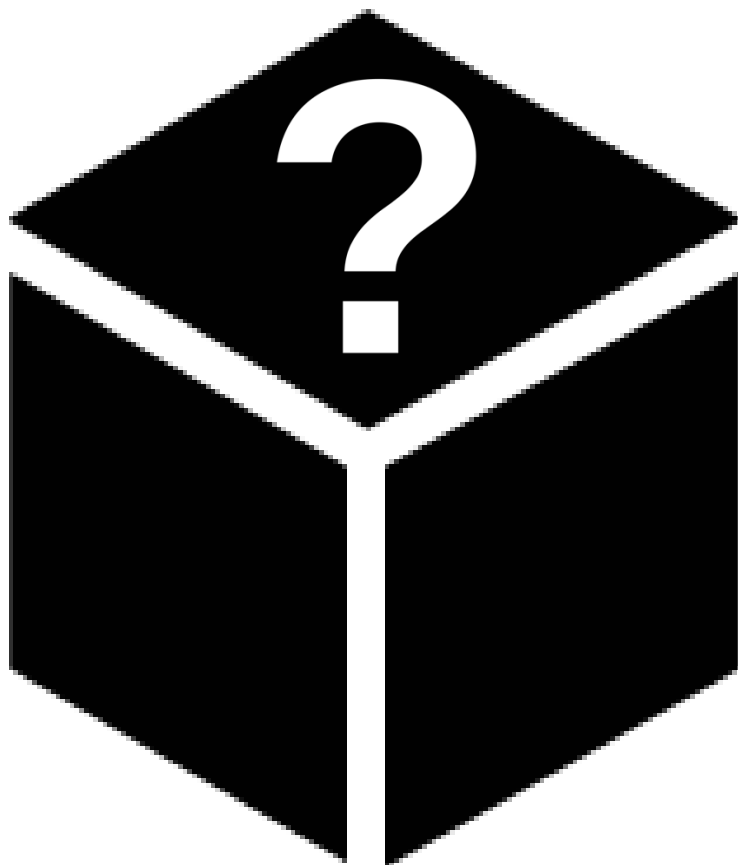
workload



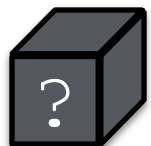
data layout



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

Granulari

Data
Movement
Policy

Compaction
Trigger

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?





Compaction **Granularity**

data moved per compaction



Compaction **Granularity**

data moved per compaction



consecutive
levels

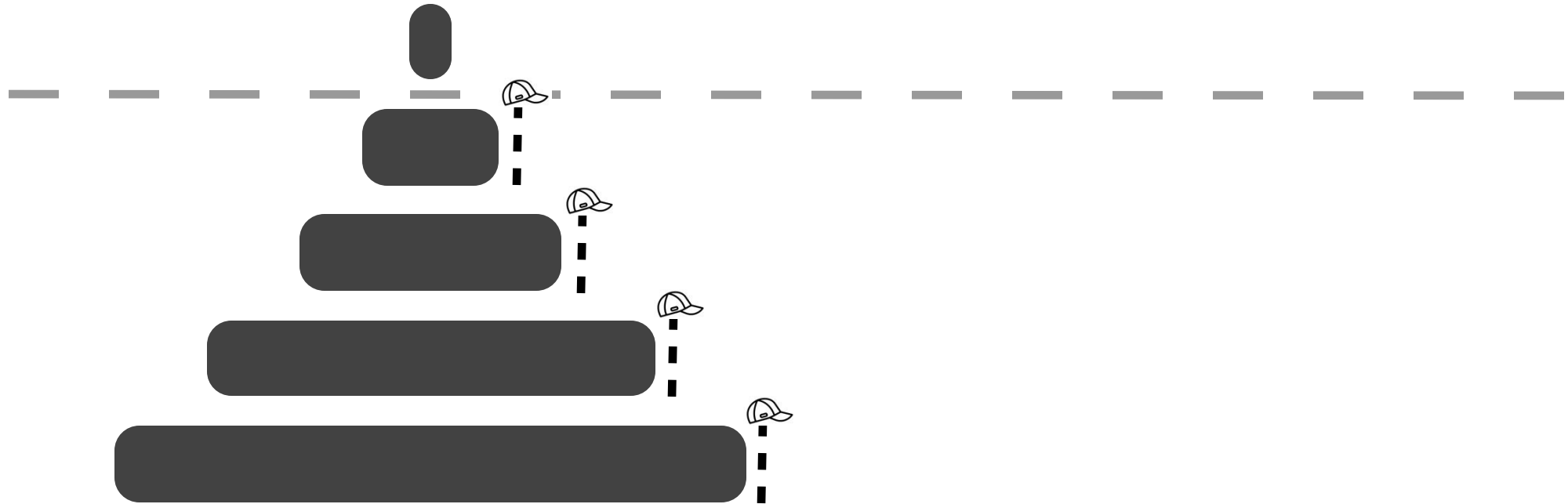
AsterixDB





Compaction **Granularity**

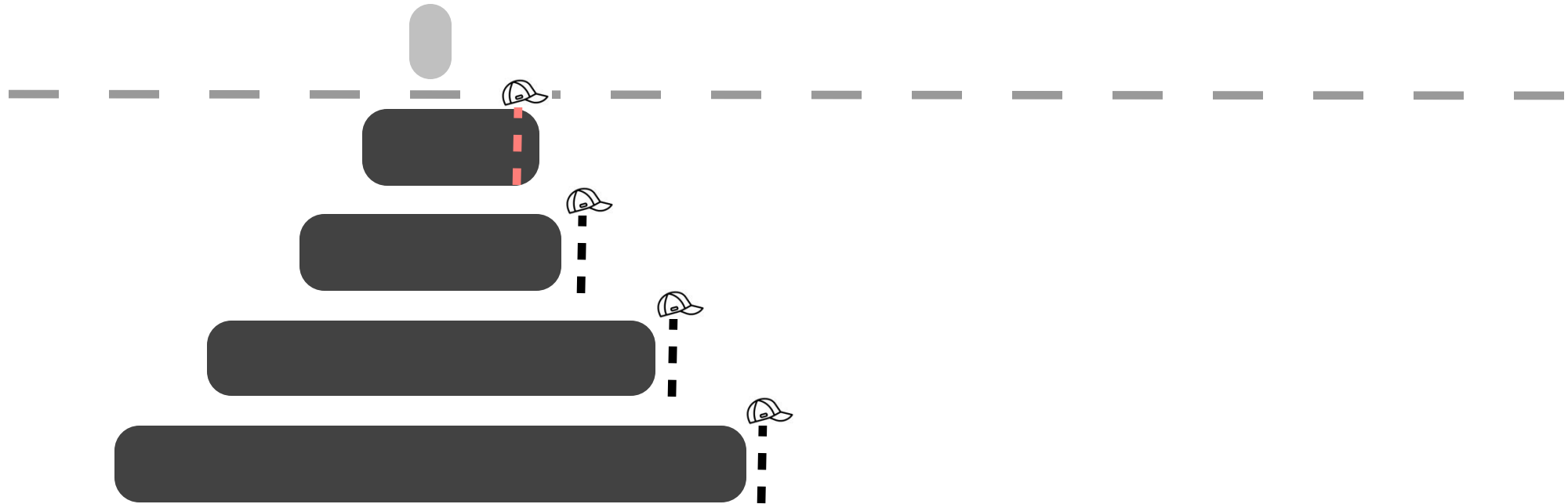
data moved per compaction





Compaction **Granularity**

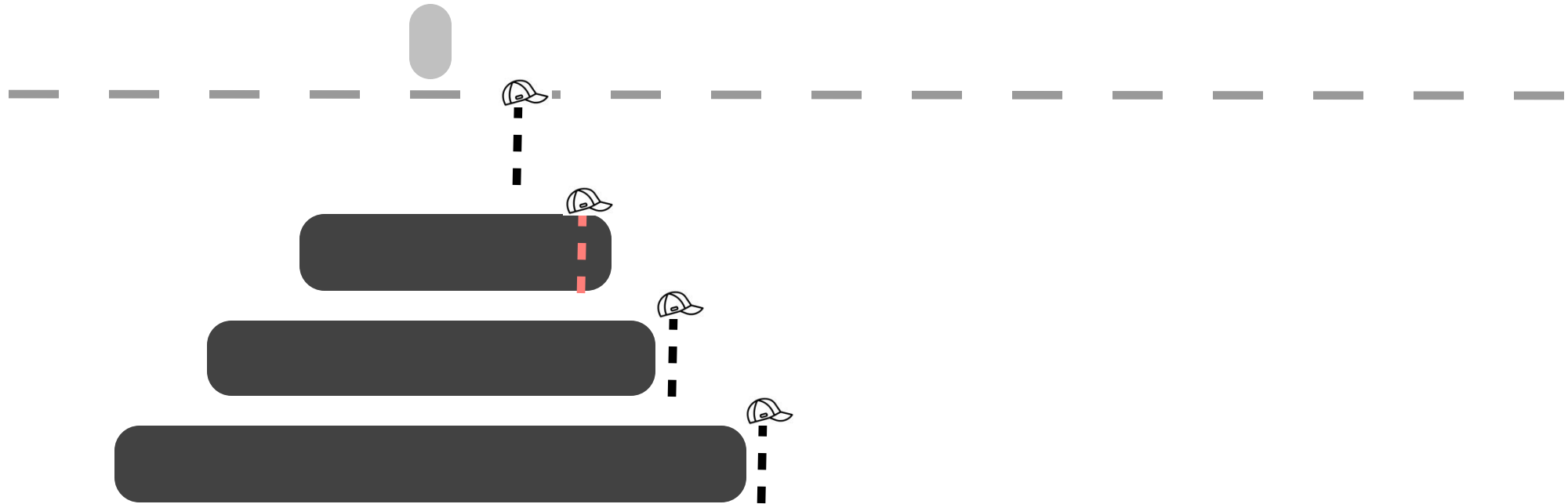
data moved per compaction





Compaction **Granularity**

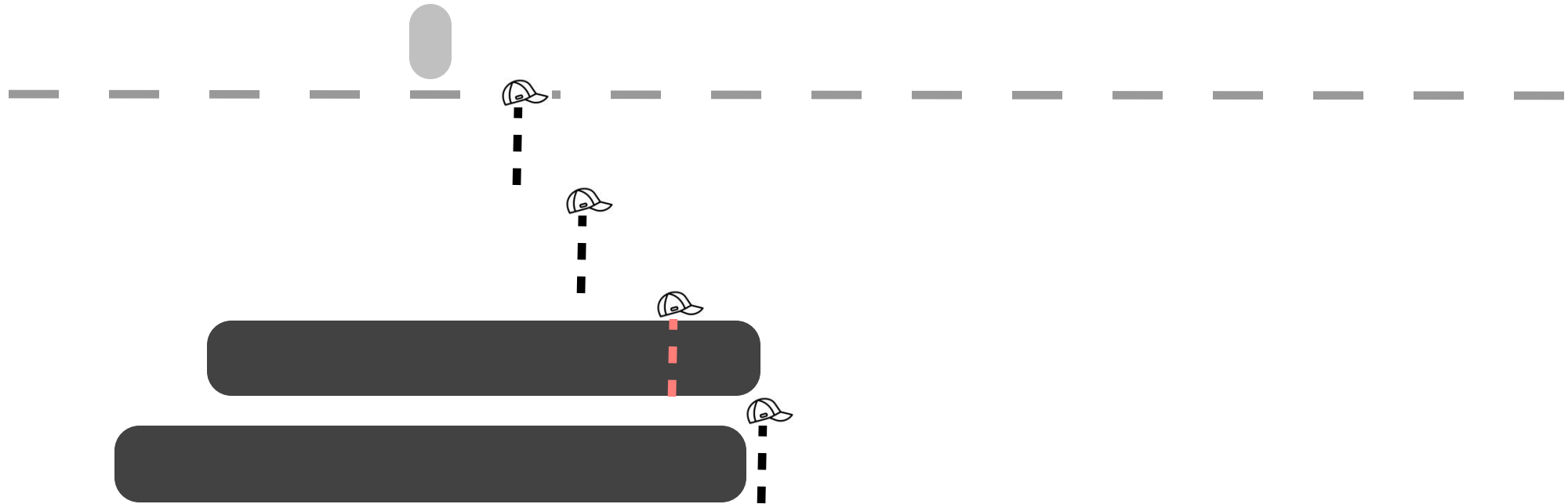
data moved per compaction





Compaction **Granularity**

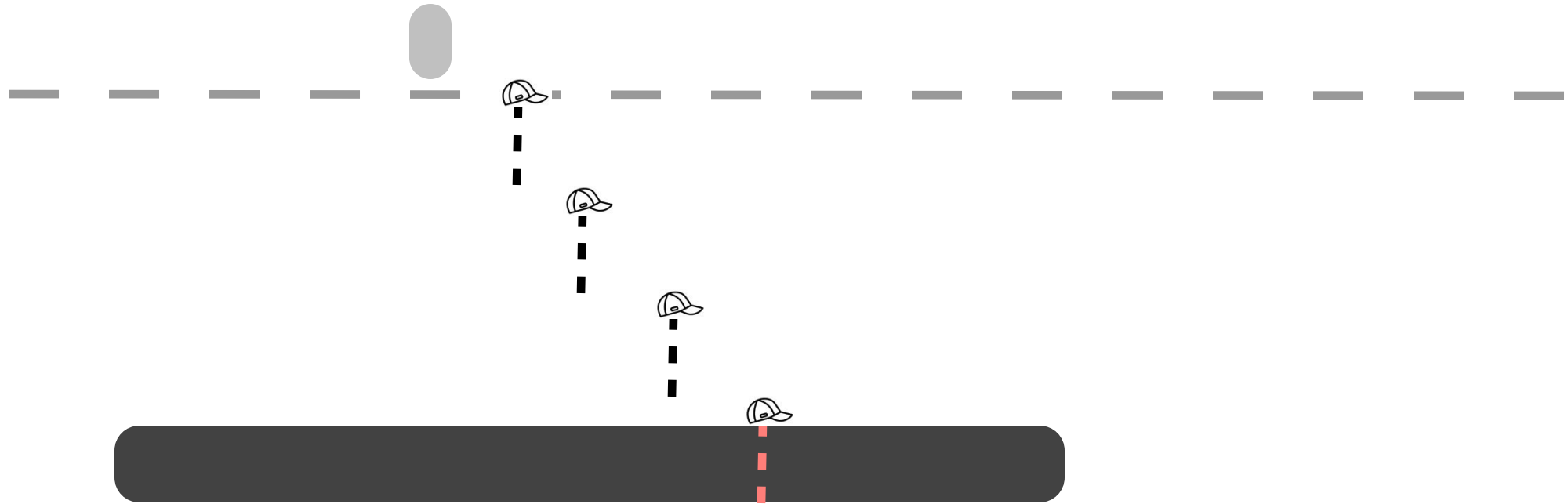
data moved per compaction





Compaction **Granularity**

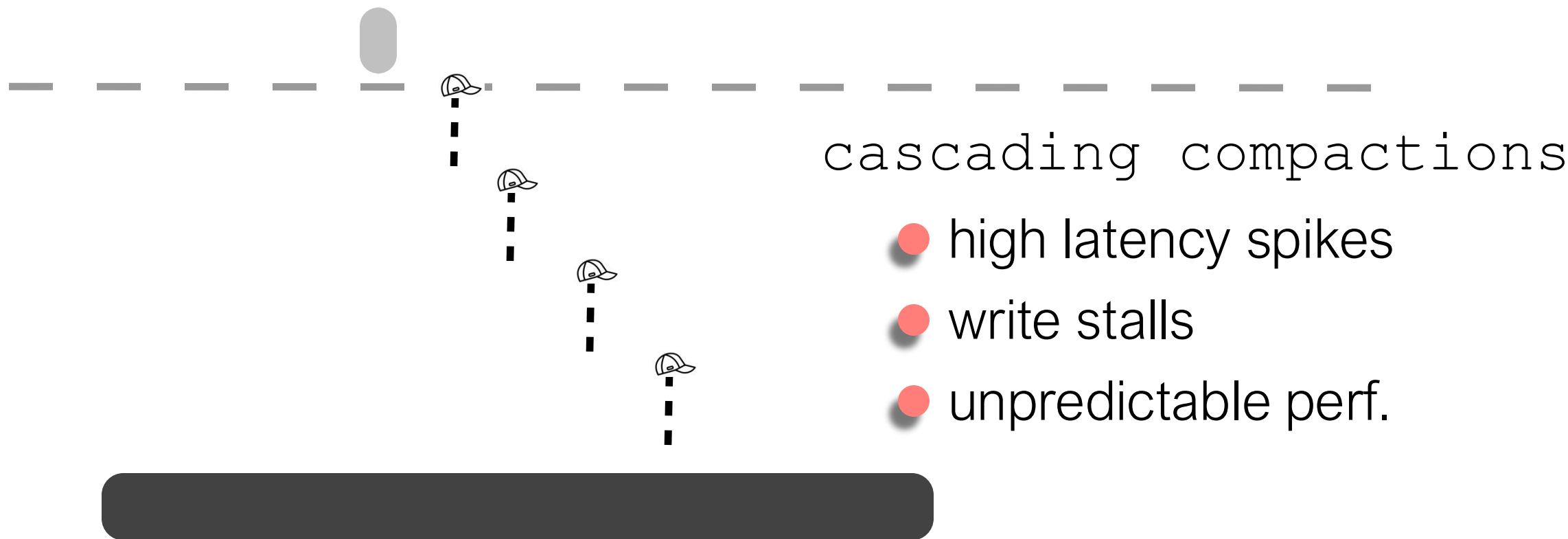
data moved per compaction





Compaction **Granularity**

data moved per compaction





Compaction **Granularity**

data moved per compaction

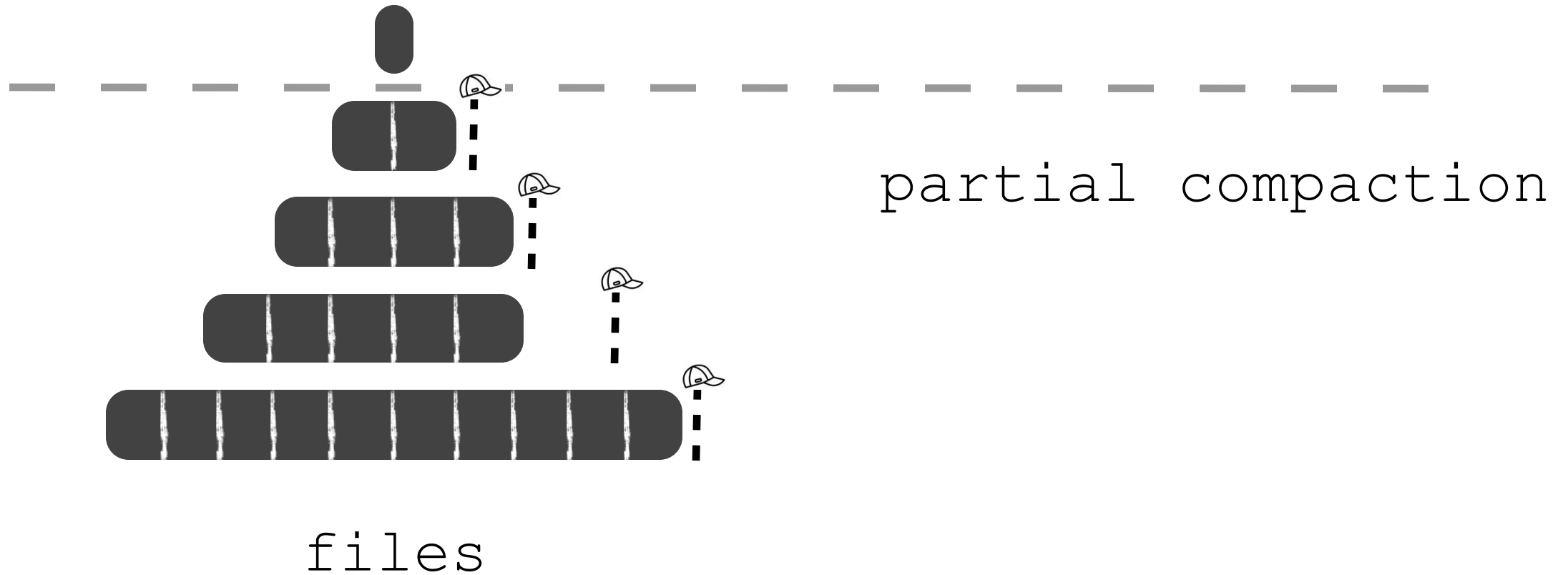
partial compaction

granularity:files



Compaction **Granularity**

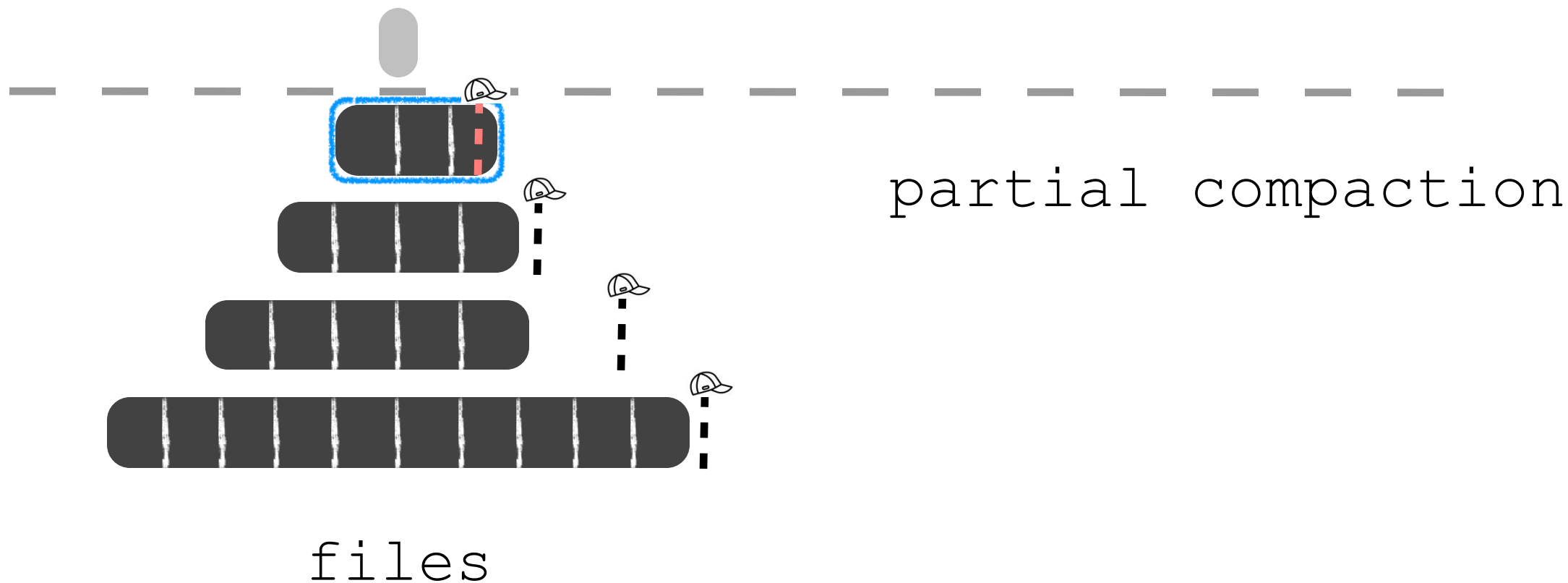
data moved per compaction





Compaction **Granularity**

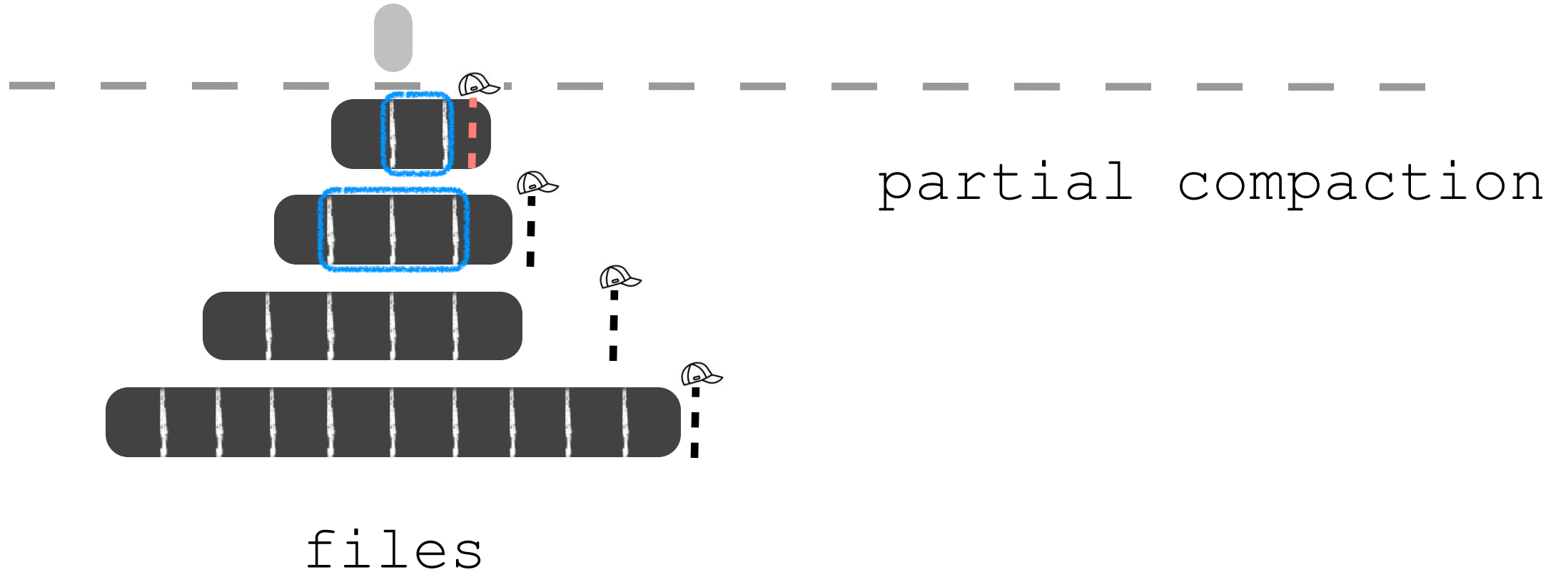
data moved per compaction





Compaction **Granularity**

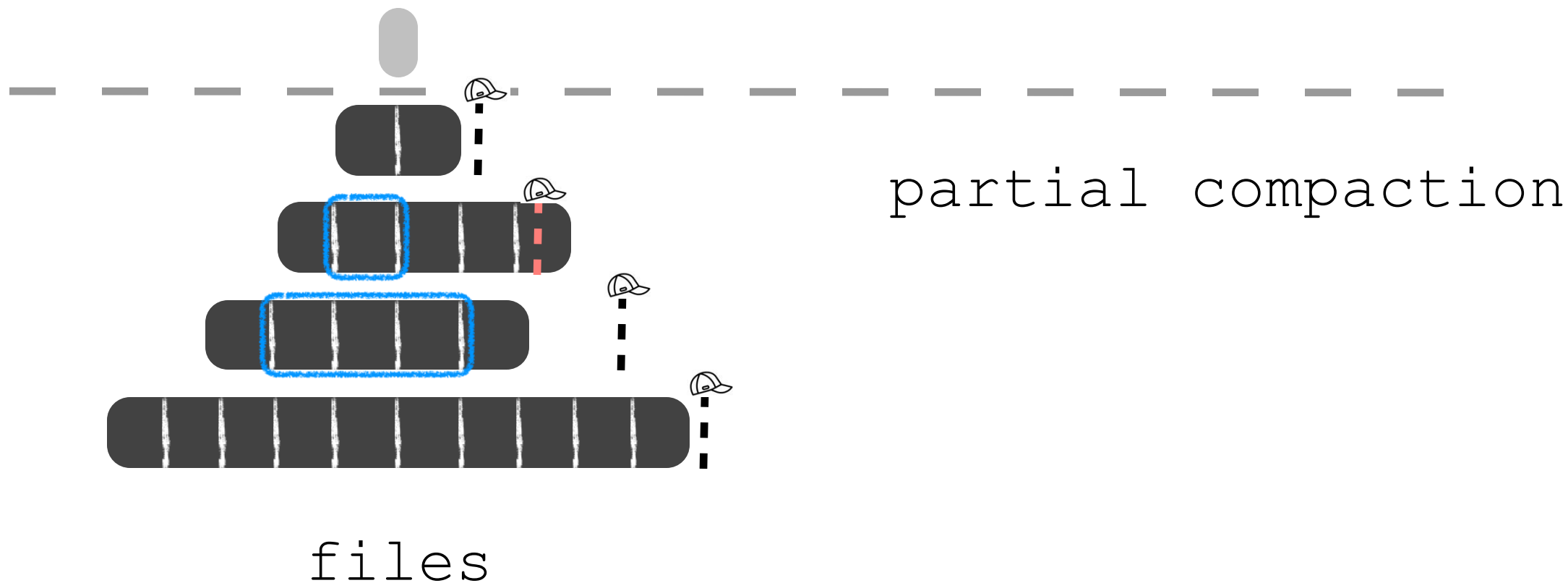
data moved per compaction





Compaction **Granularity**

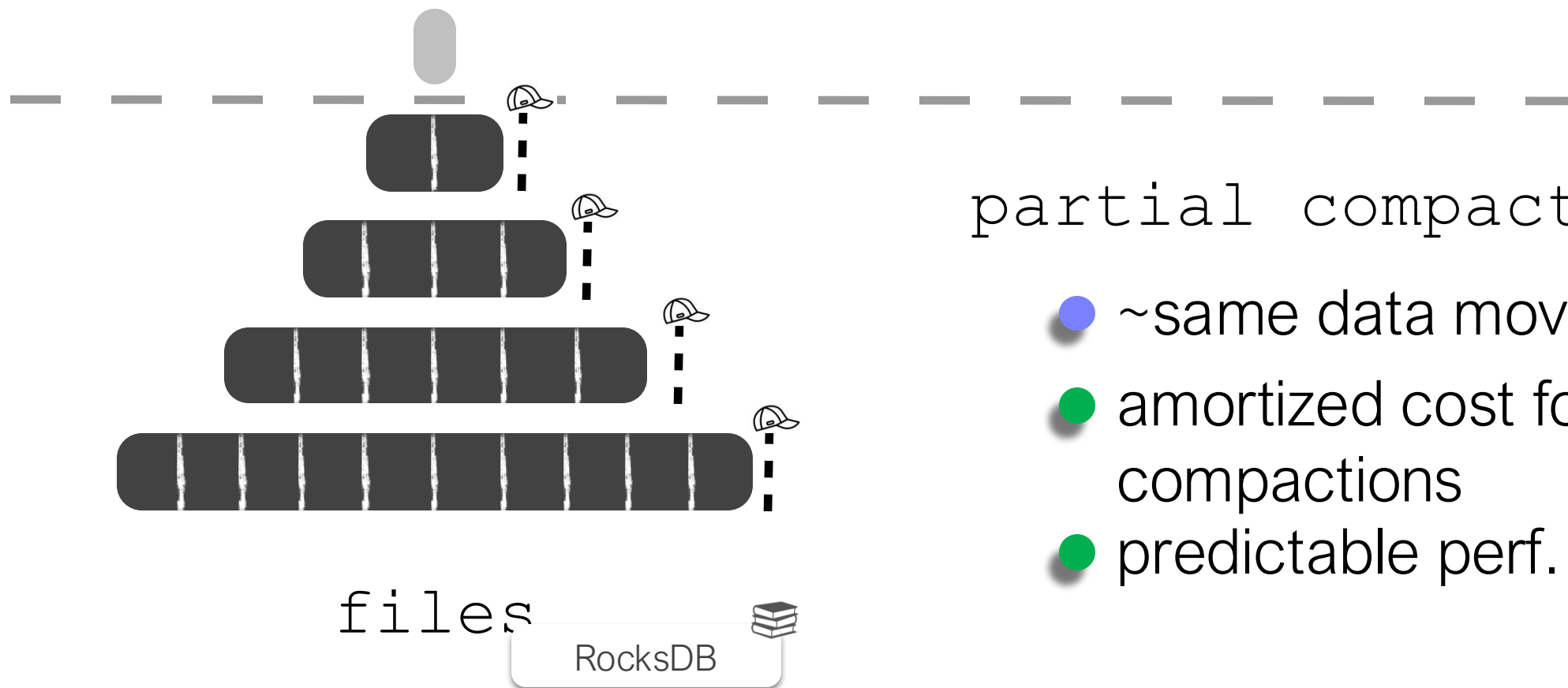
data moved per compaction





Compaction **Granularity**

data moved per compaction



partial compaction

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



Compaction **Granularity**

data moved per compaction

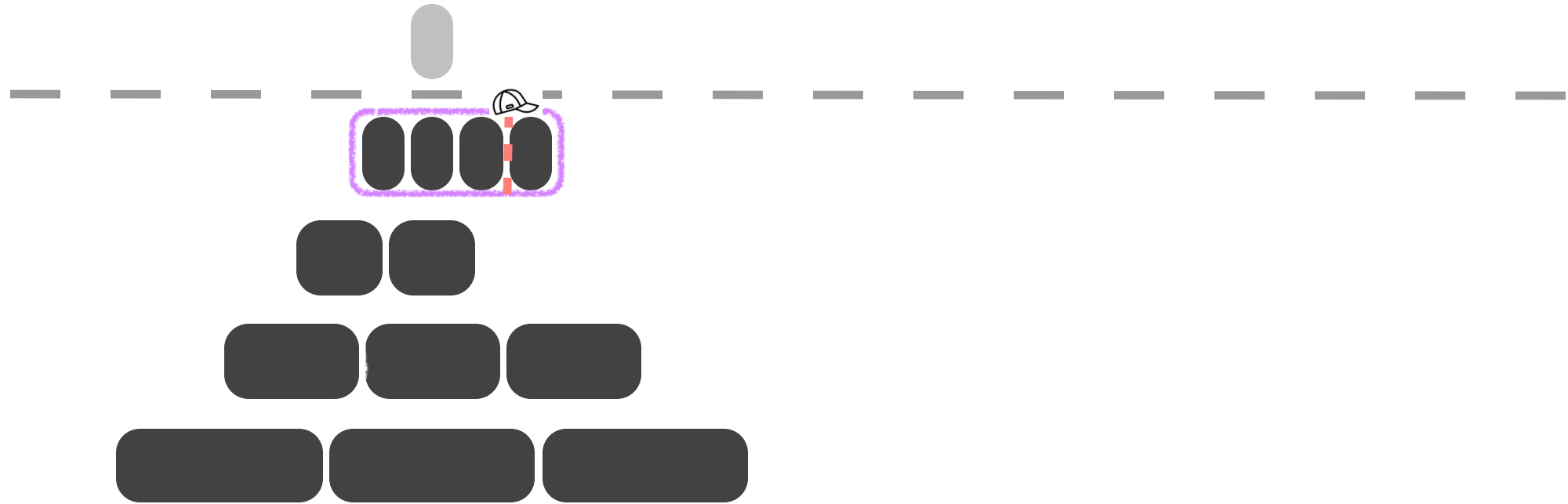


sorted runs in a
level



Compaction **Granularity**

data moved per compaction

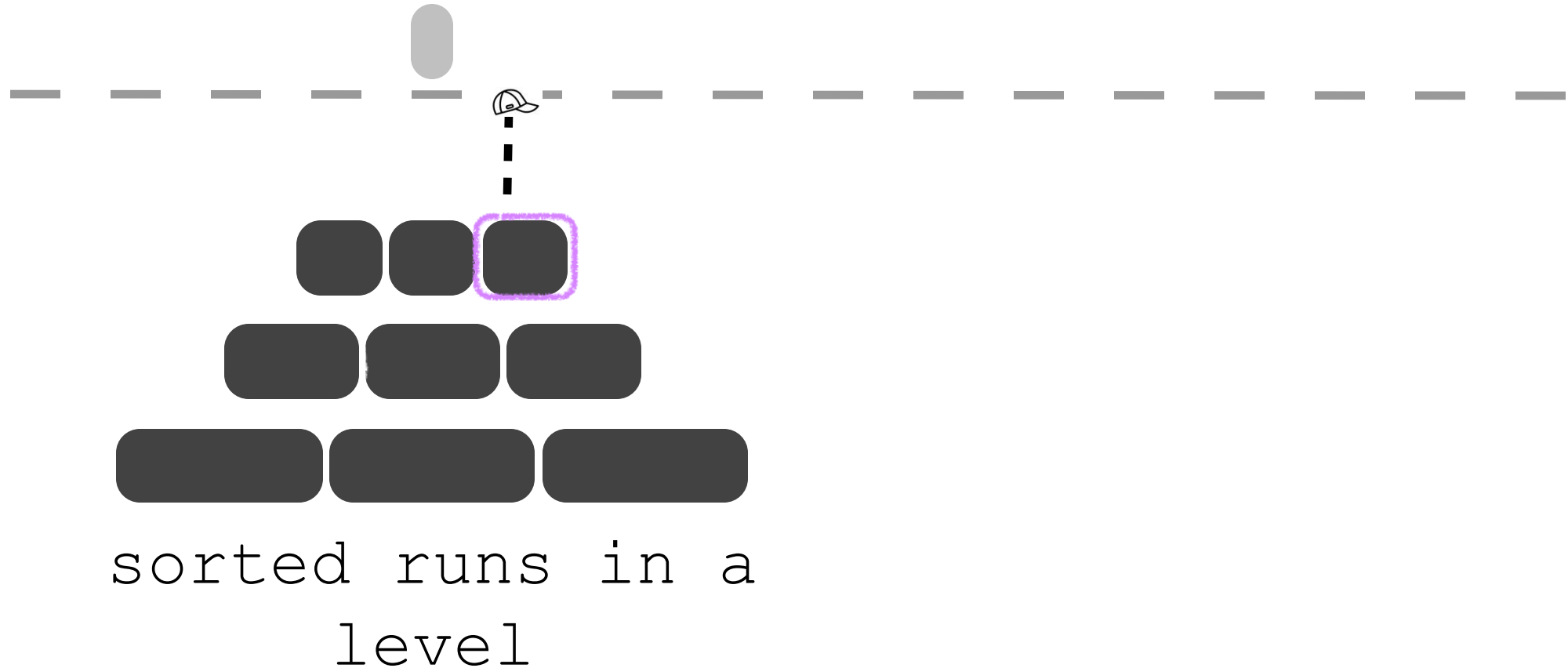


sorted runs in a
level



Compaction **Granularity**

data moved per compaction



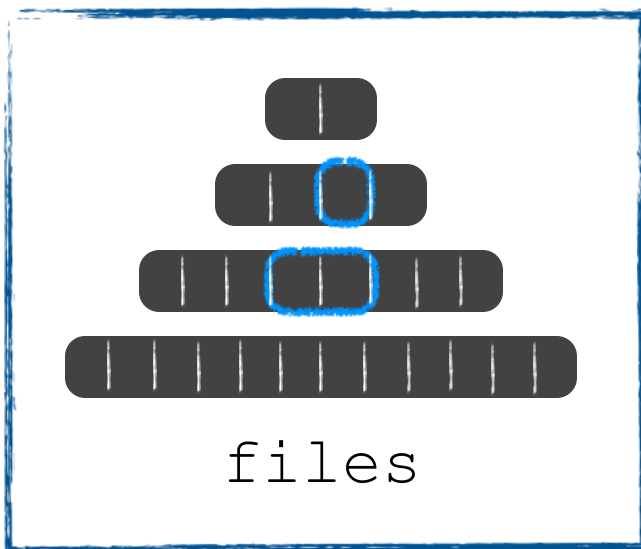


Compaction **Granularity**

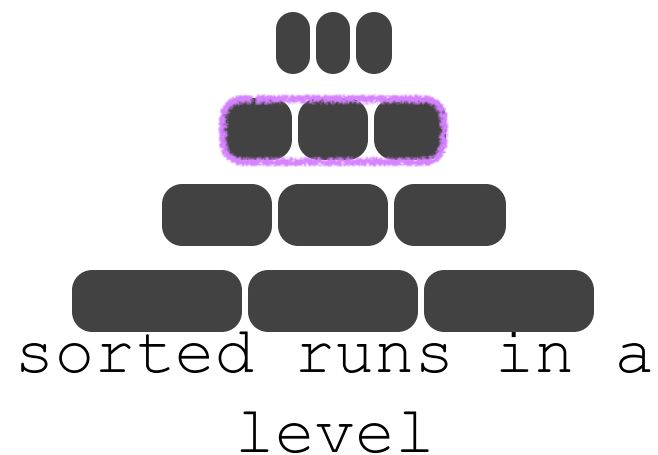
data moved per compaction



levels



files

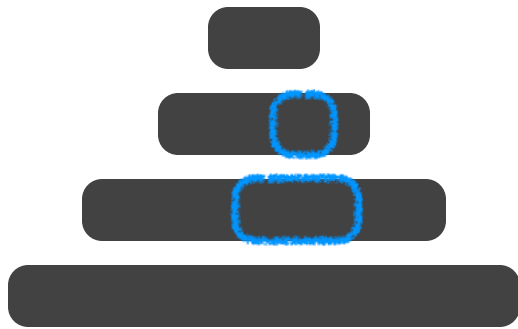


sorted runs in a
level



Data Movement Policy

which data to compact

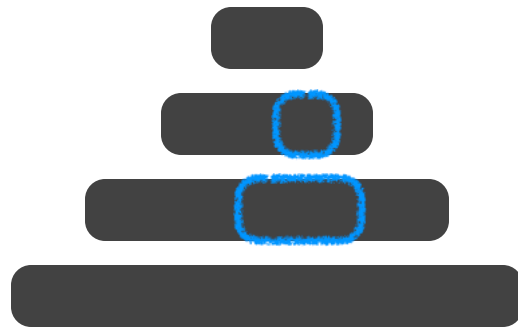


files



Data Movement Policy

which data to compact



files

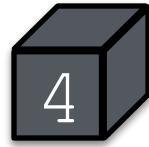
round-robin

minimum **overlap with parent** level

file with most **tombstones**

coldest file

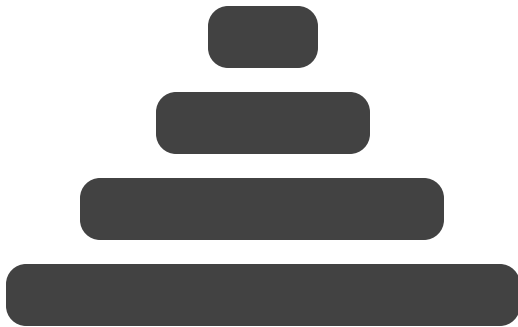


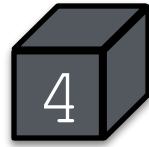


Compaction **Trigger**

invoking the compaction routine

level **saturation**

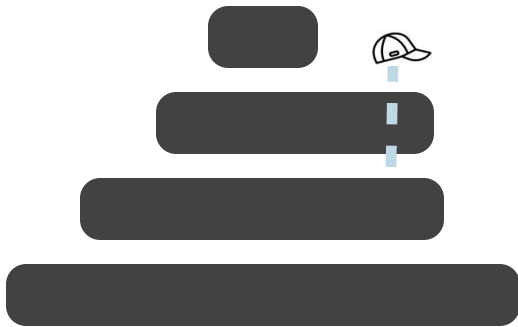


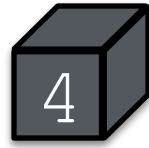


Compaction **Trigger**

invoking the compaction routine

level **saturation**

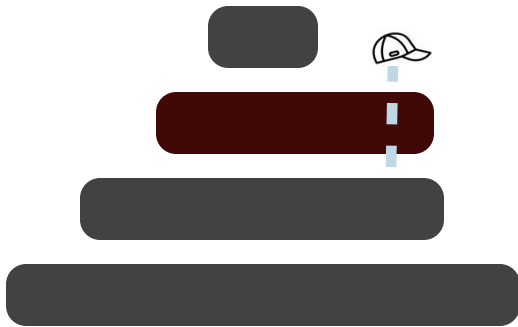


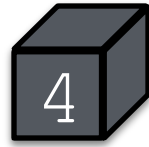


Compaction **Trigger**

invoking the compaction routine

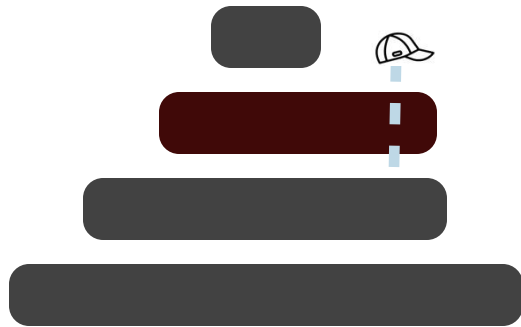
level **saturation**





Compaction **Trigger**

invoking the compaction routine



level **saturation**

number of **sorted runs**

space amplification

SA

age of a file

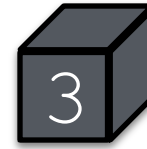
De



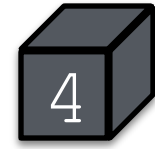
Data Layout



Compaction
Granularity



Data Movement
Policy



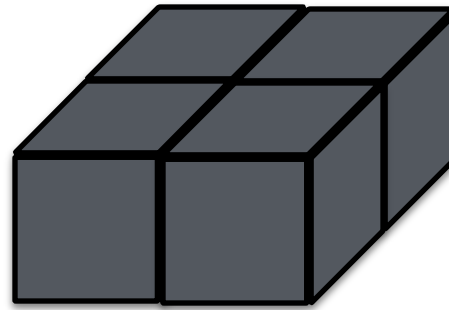
Compaction
Trigger

Data Layout

Compaction
Granularity

Data Movement
Policy

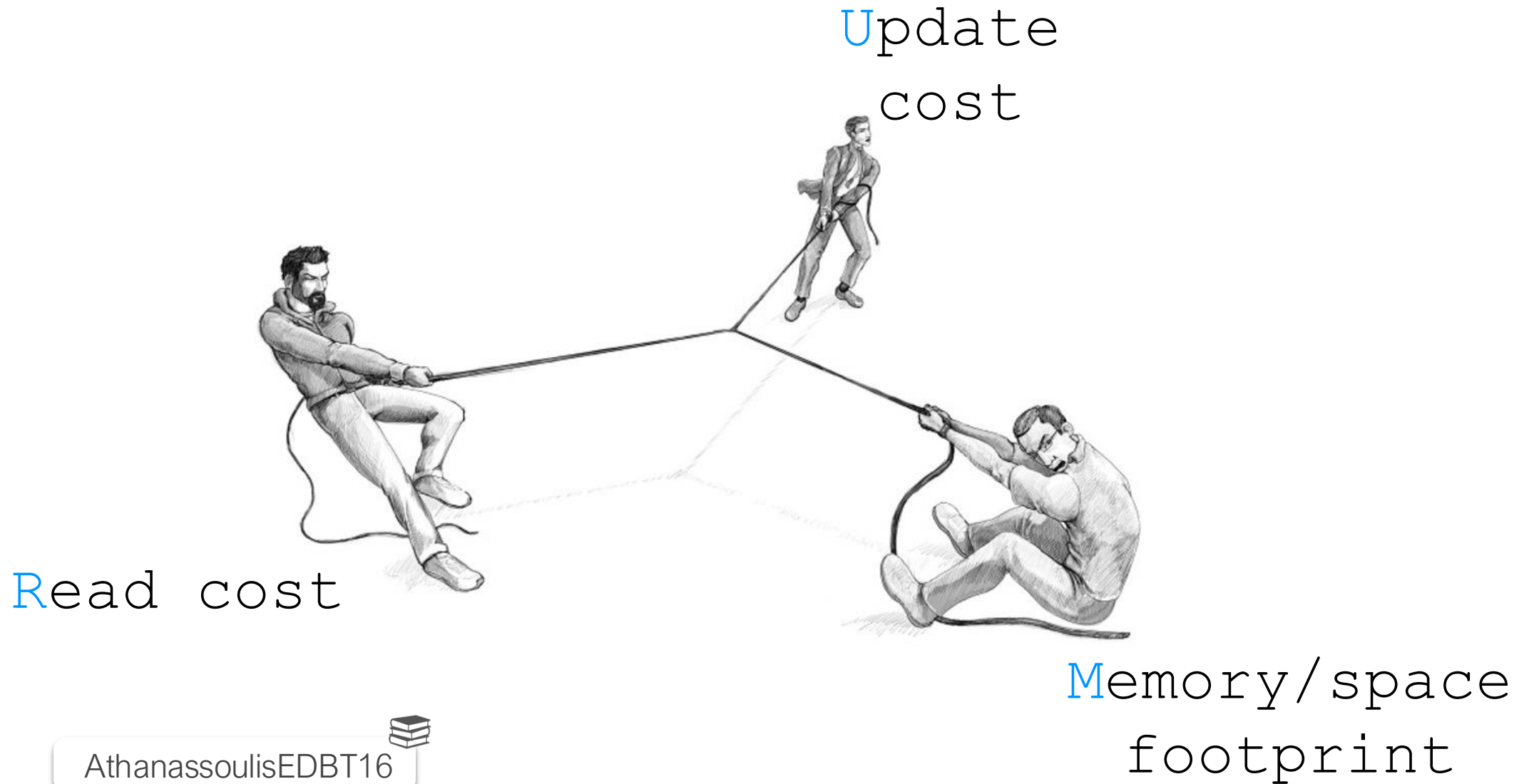
Compaction
Trigger



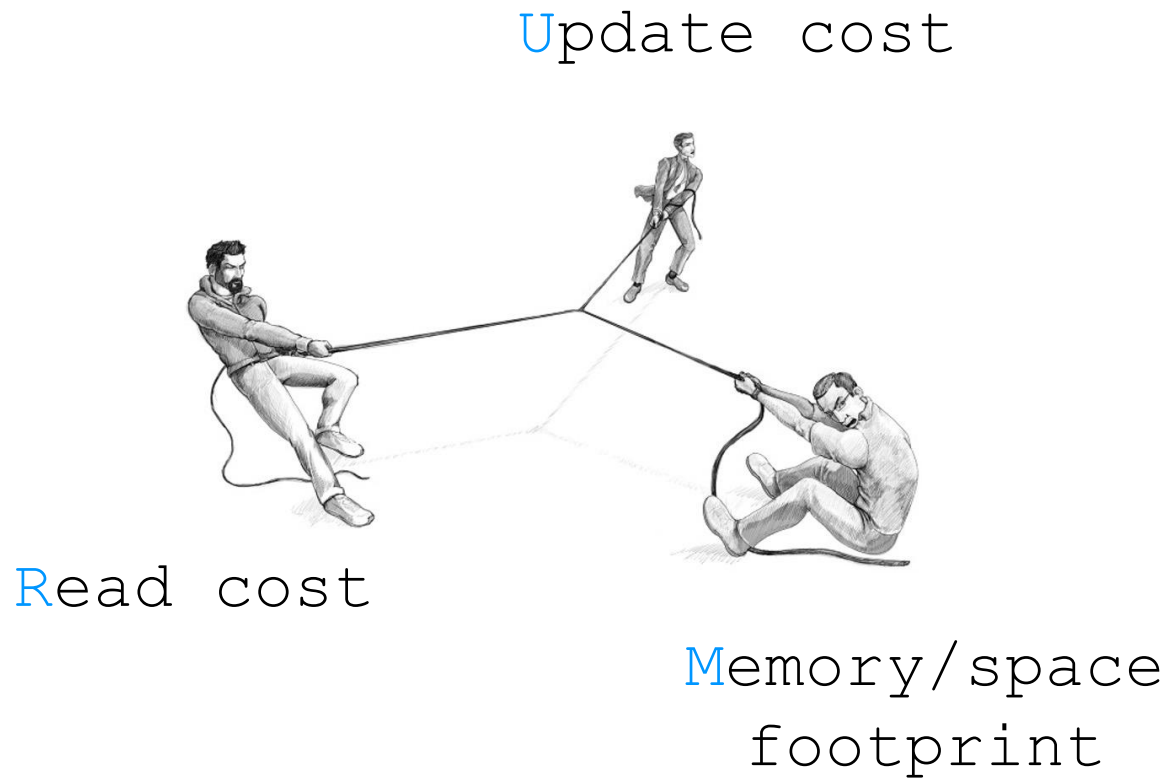
***Any* Compaction Algorithm**

LSM **Design Space**

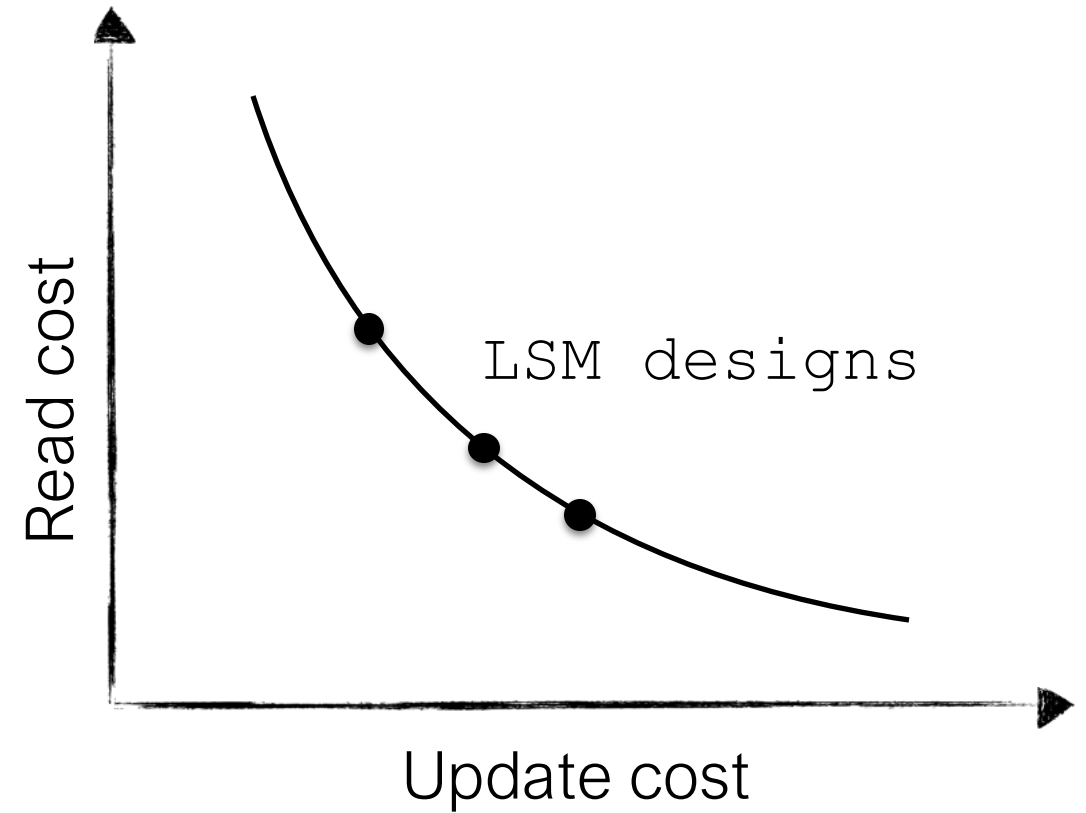
LSM Design Space



LSM Design Space



fixed Memory

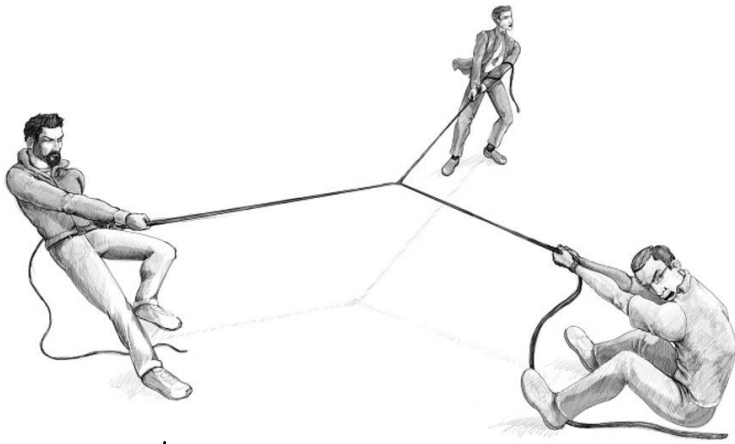


LSM Design Space

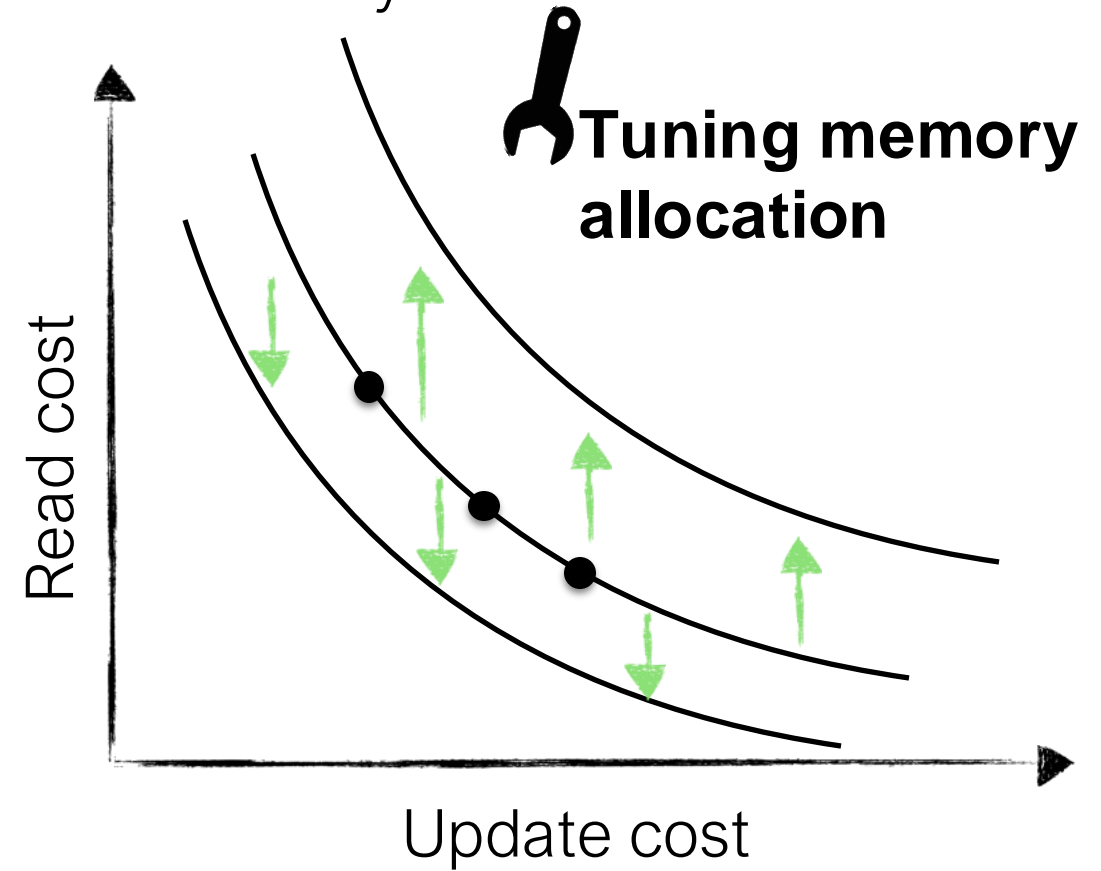
Update cost

Read cost

Memory/space
footprint



fixed Memory

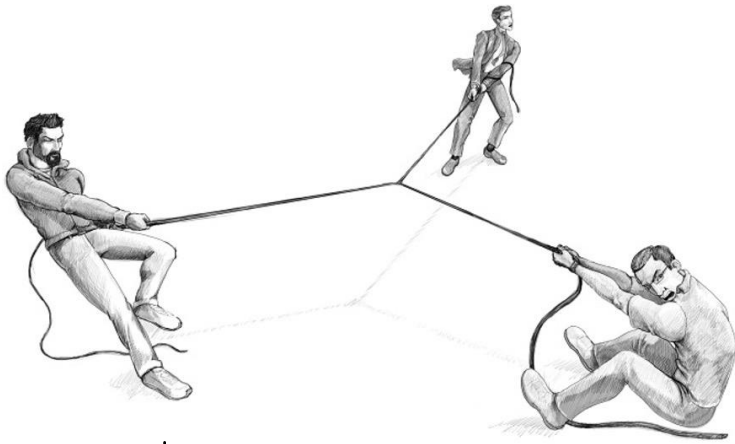


LSM Design Space

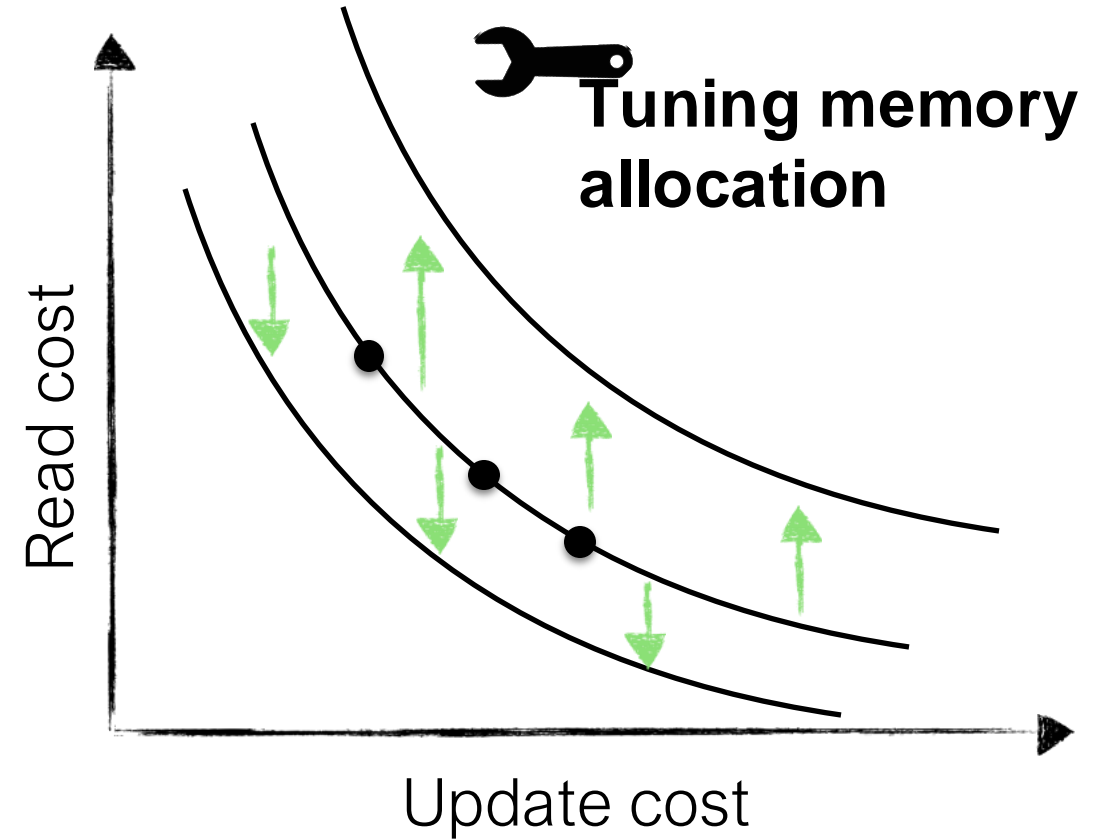
Update cost

Read cost

Memory/space
footprint



fixed Memory



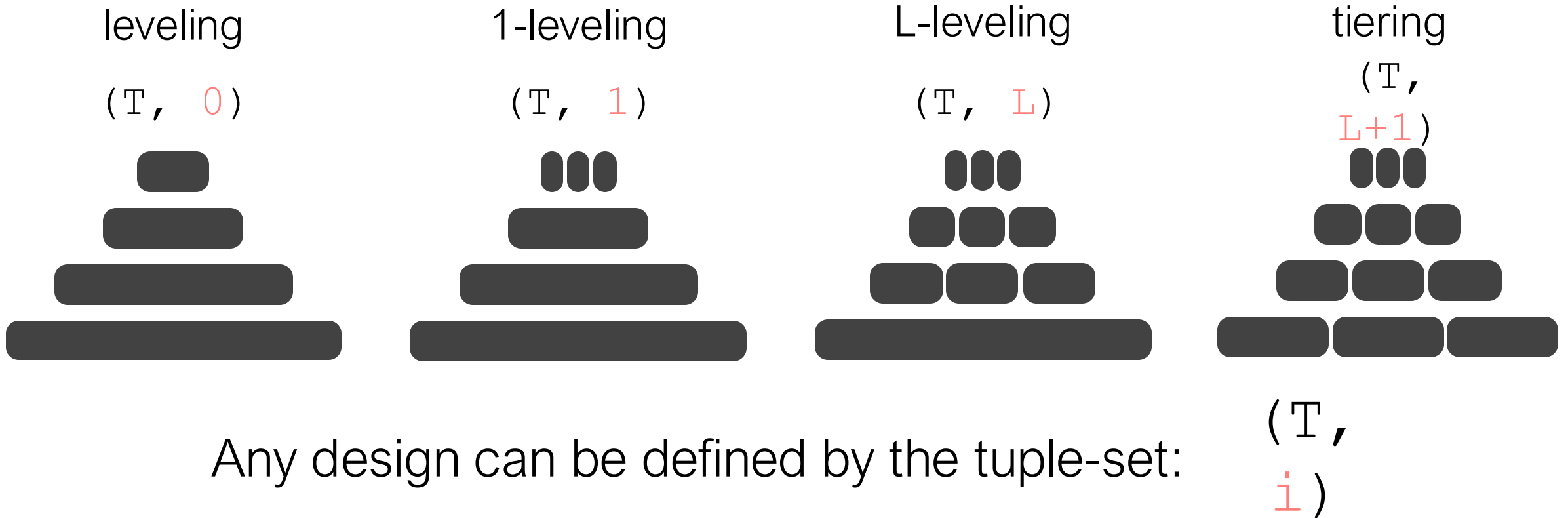
Storage Layer **Design Continuum**



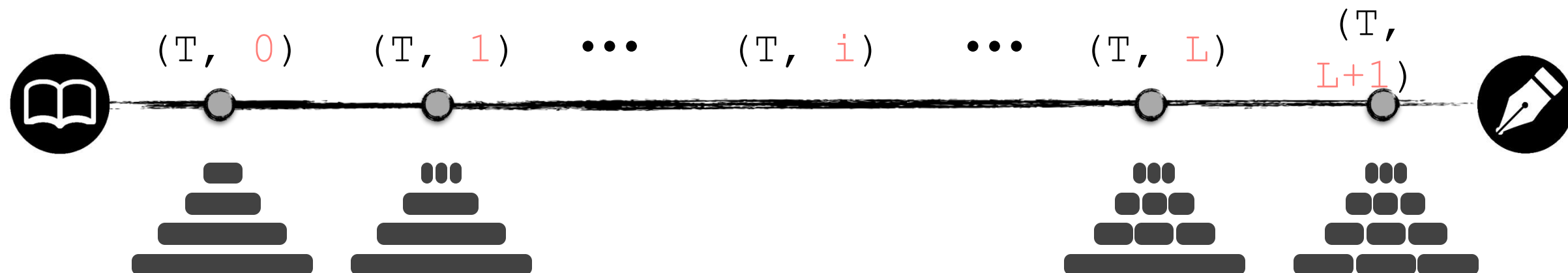
Any design can be defined by the tuple-set: $(T, \underset{\text{red}}{i})$



Storage Layer **Design Continuum**



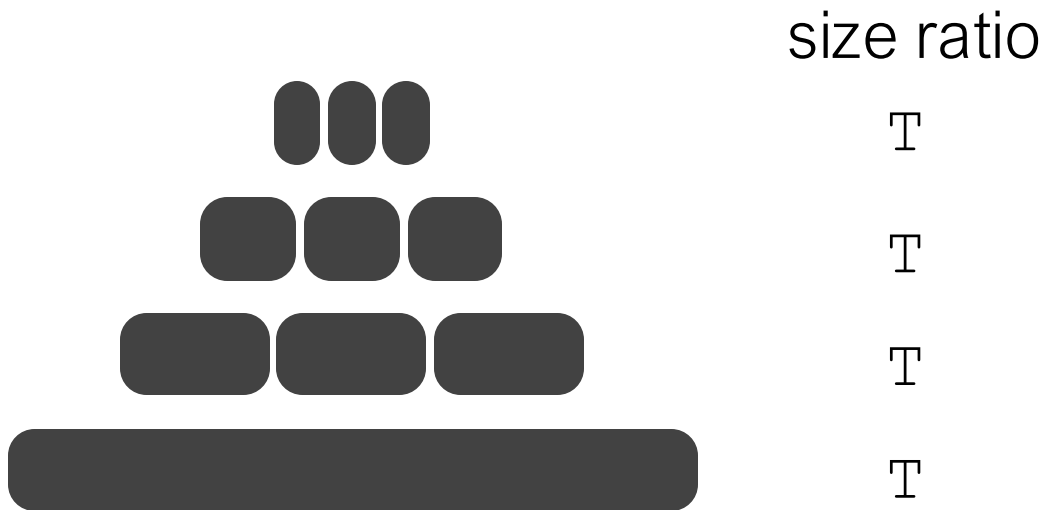
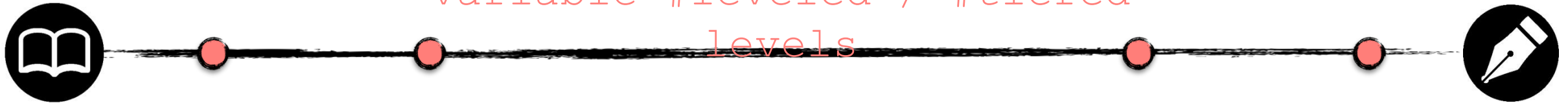
Storage Layer **Design Continuum**



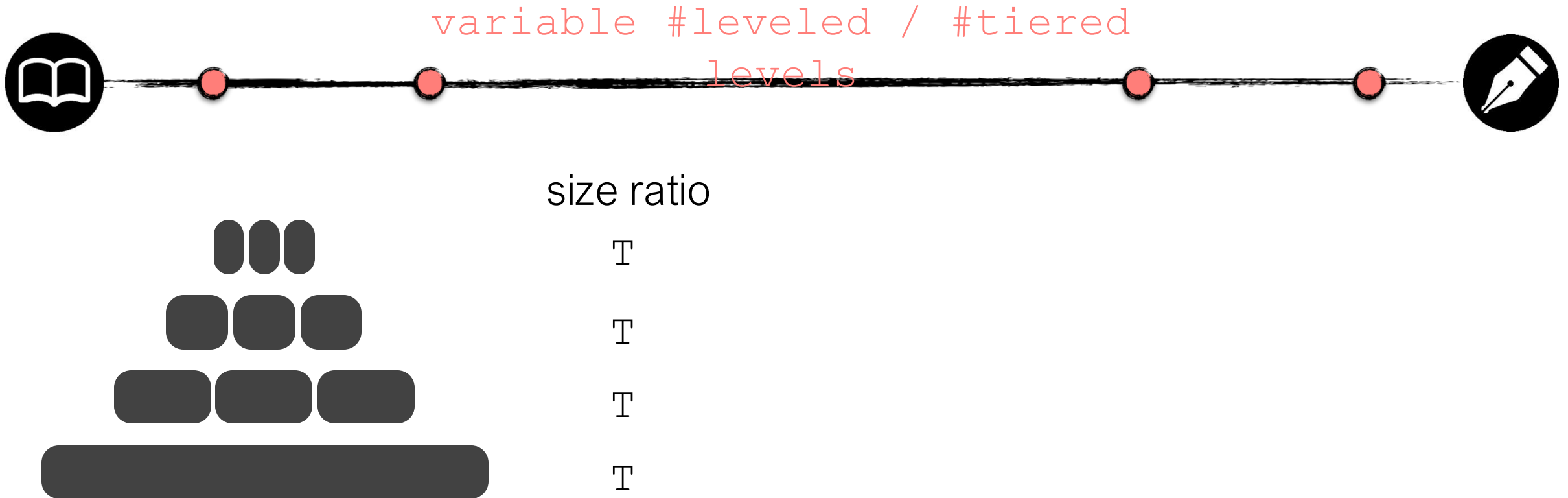
Storage Layer **Design Continuum**



Storage Layer **Design Continuum**

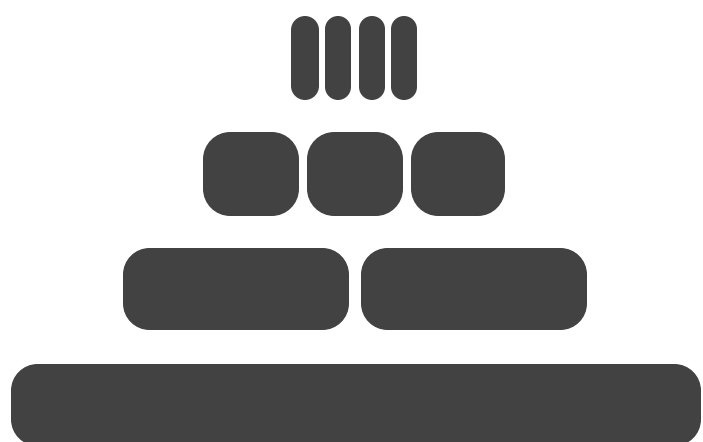


Storage Layer **Design Continuum**







Storage Layer **Design Continuum**



	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



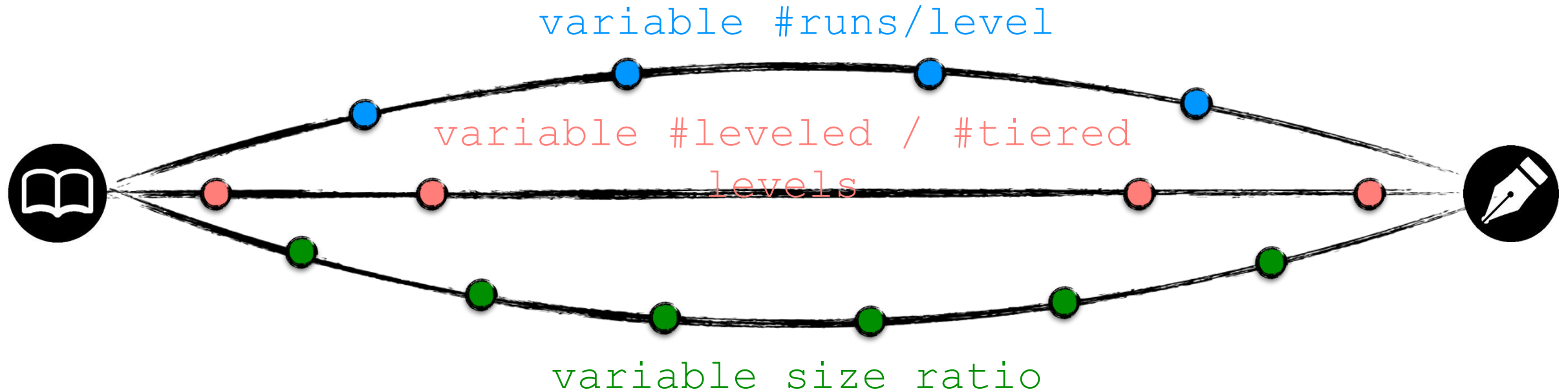
Storage Layer Design Continuum



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



Storage Layer **Design Continuum**



The LSM storage layer
design continuum

CS 561: Data Systems Architectures

class 6

Log-structured Merge Trees

Prof. Manos Athanassoulis

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