

CS 561: Data Systems Architectures

class 6

Log-structured Merge Trees

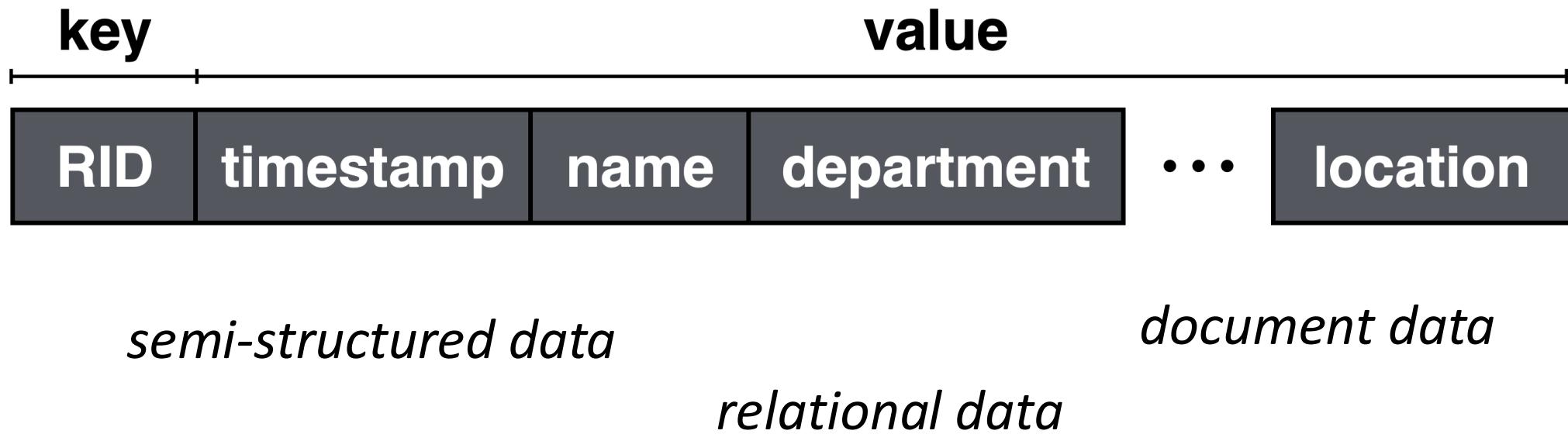
Prof. Manos Athanassoulis

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

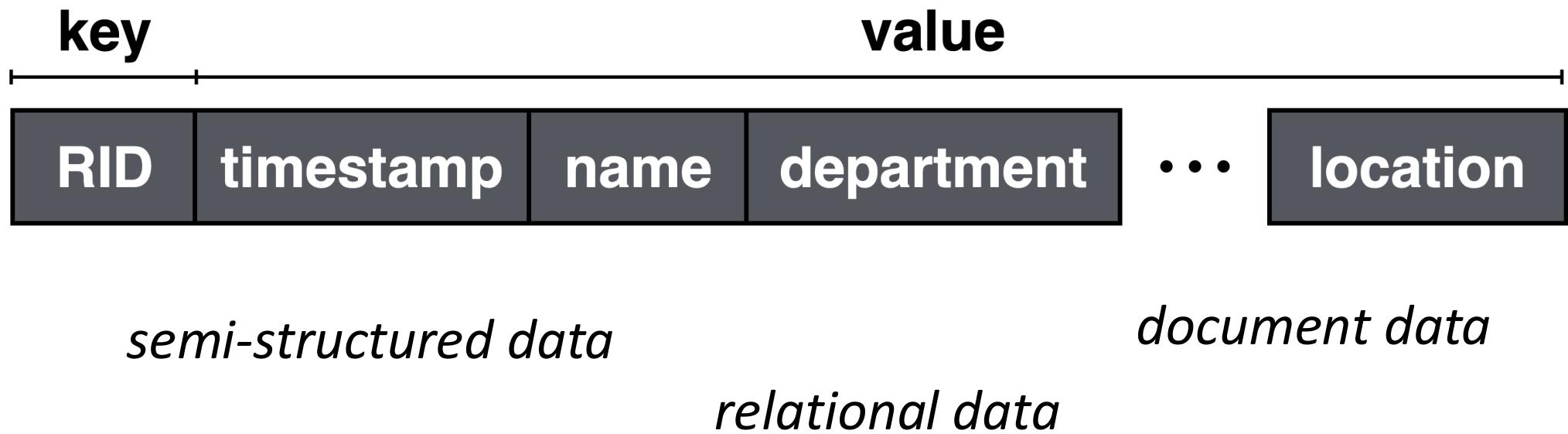
Do we have a quiz today ... ?

Yes!

Key-Value Stores



LSM-tree based Key-Value Stores



Log-Structured Merge-tree

LSM-tree

LSM-tree

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

1996

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

LSM-tree
O'Neil *et al.*

1996

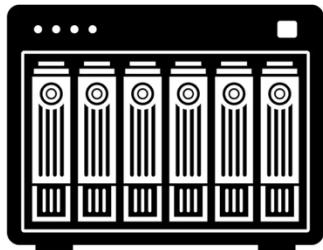
Patrick O'Neil
UMass Boston



LSM-tree
O'Neil et al.

1996

- good random writes
- good reads



array
of discs

1980s

why?
RAID, striping ← ?

✗ LSM not explicitly needed

LSM-tree
O'Neil et al.

1996

a decade



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 ($6\text{ms}/2$)

avg. seek time: 0.75ms ($1.5\text{ms}/2$)

1 I/O / 3.75ms: 267 IOPS

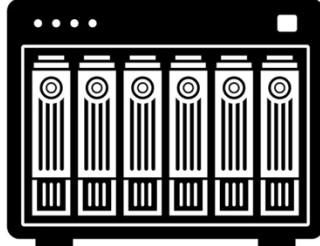
100 disks: 26,700 IOPS



Bigtable

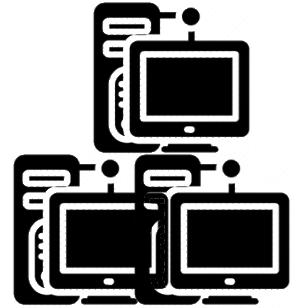
2006

- good random writes
- good reads
- poor ingestion perf.
- poor query perf.



array
of discs

what happened in 2006?



commodity
hardware

1980s

LSM-tree
O'Neil et al.

1996

a decade

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.

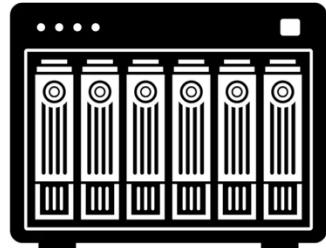


Bigtable

2006

- good random writes

- good reads



array
of discs

1980s

SSD wear-friendly

competitive rand. reads

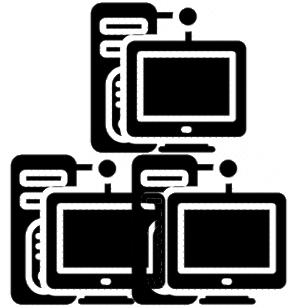
fast ingestion (sequential)

- poor ingestion perf.

- poor query perf.



why this
is important?



commodity
hardware



Bigtable

2006

SSDs have dominated
storage since!



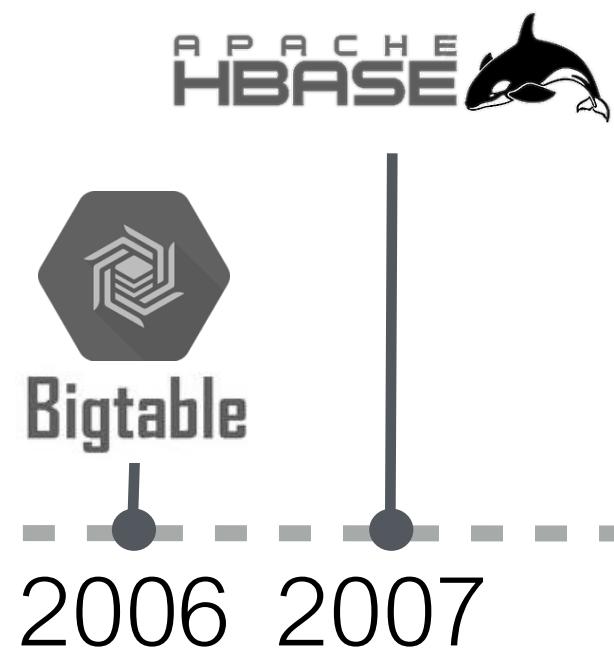
LSM-tree
O'Neil et al.

1996

a decade

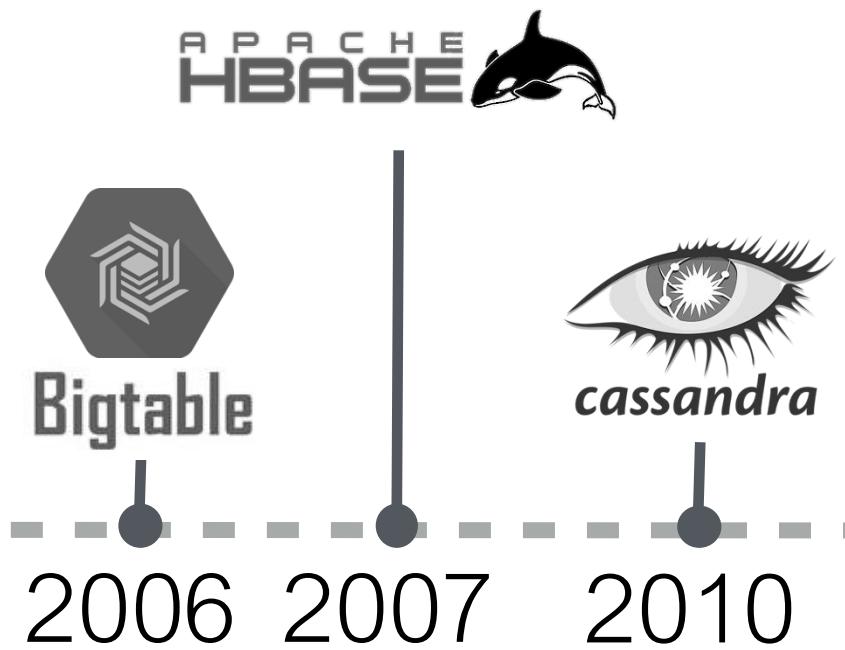
LSM-tree
O'Neil *et al.*

1996



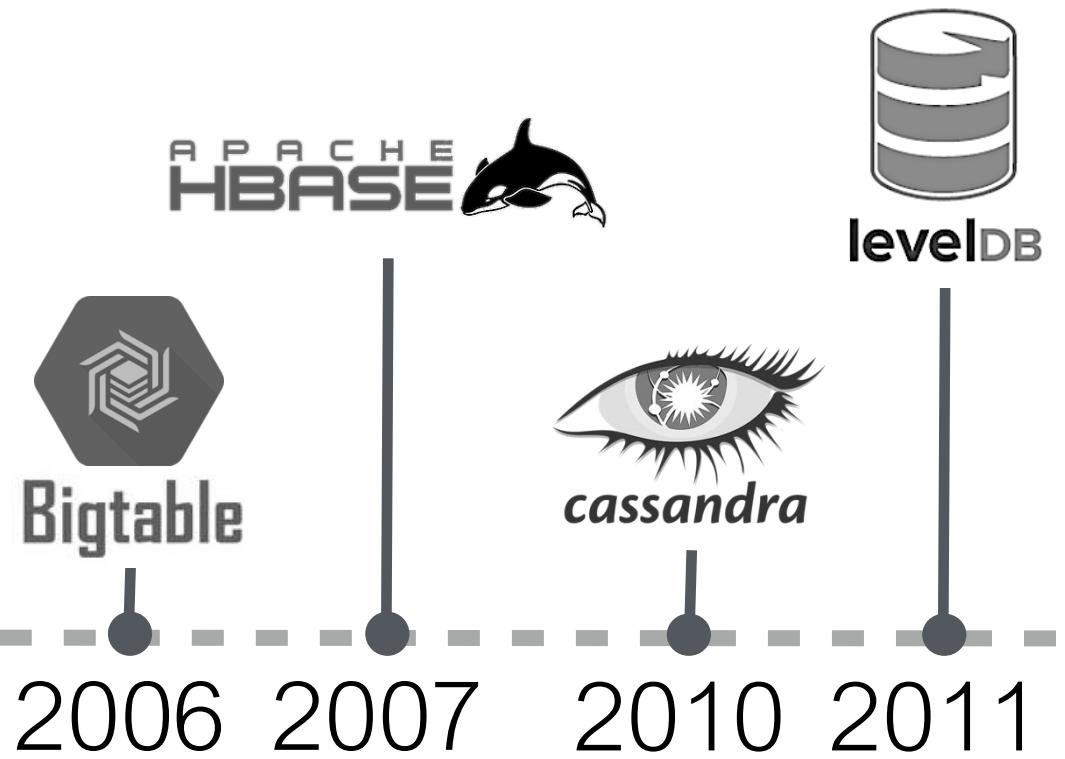
LSM-tree
O'Neil *et al.*

1996



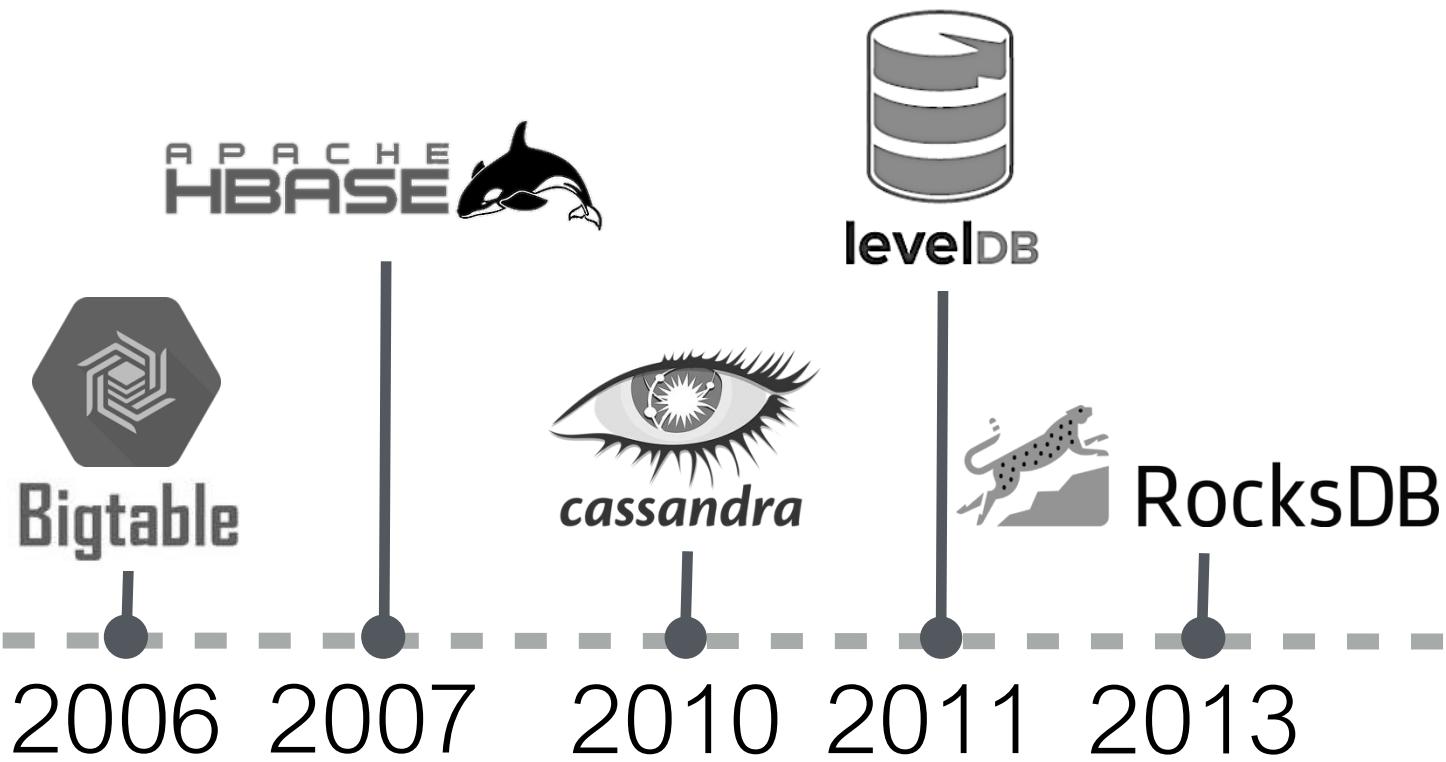
LSM-tree
O'Neil *et al.*

1996



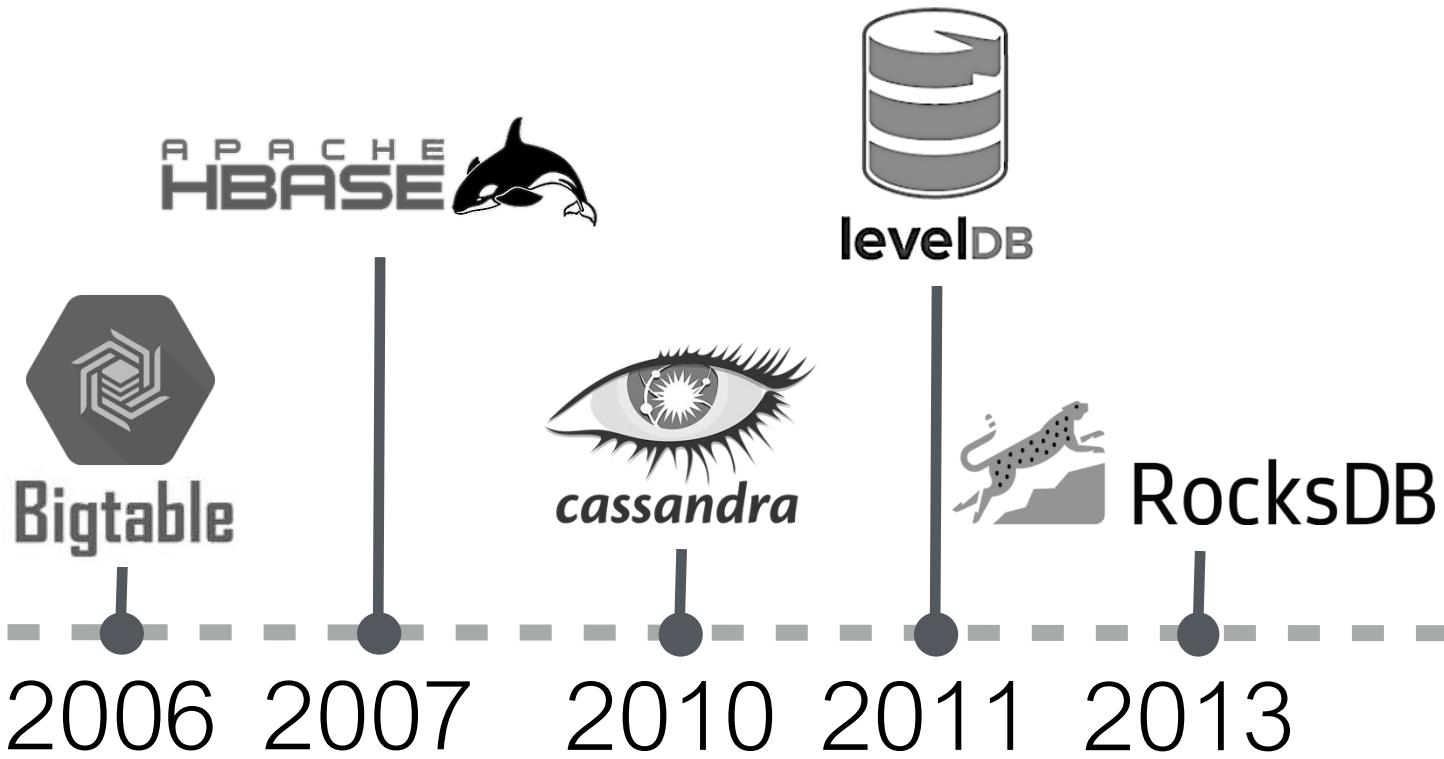
LSM-tree
O'Neil *et al.*

1996



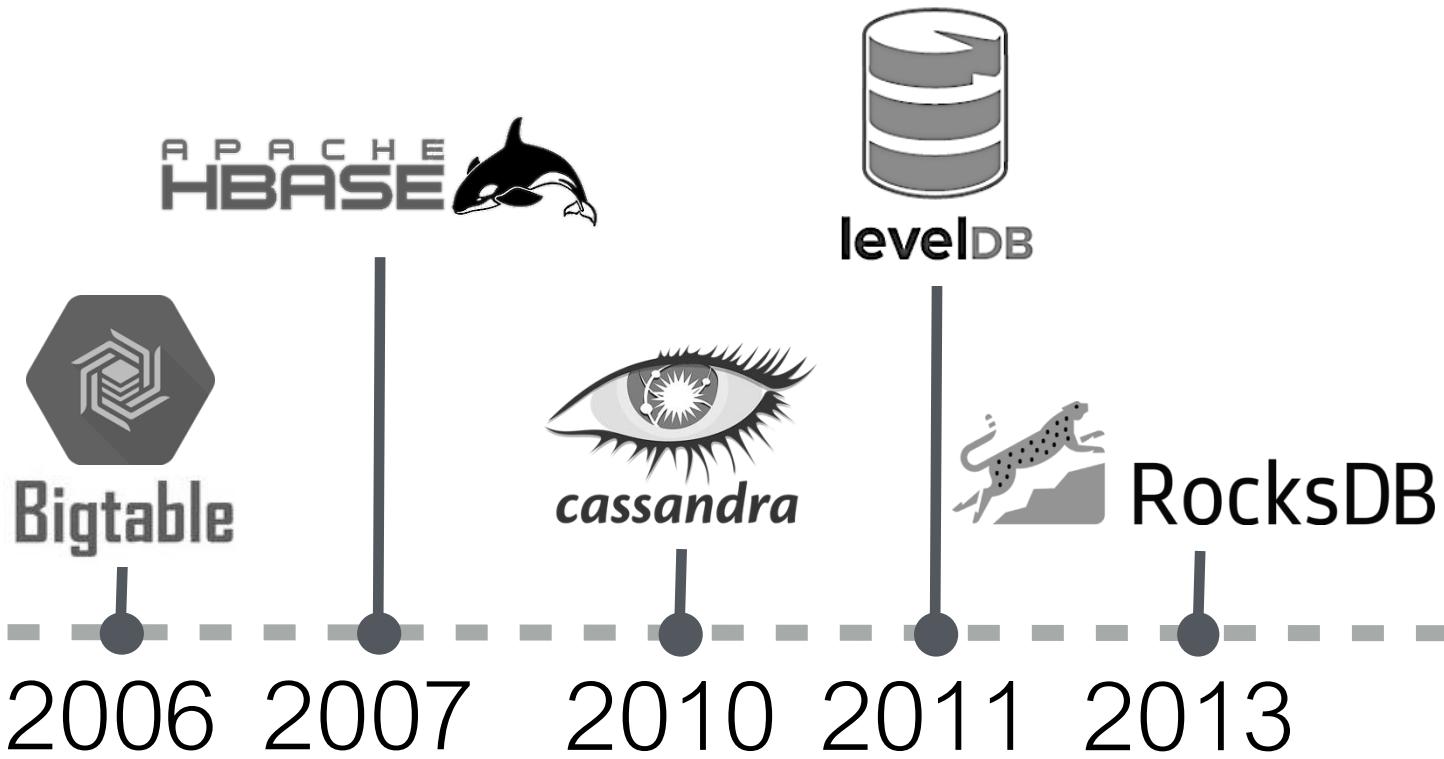
LSM-tree
O'Neil *et al.*

1996



LSM-tree
O'Neil *et al.*

1996



LSM-tree

NoSQL



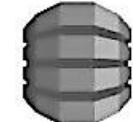
RocksDB



levelDB



SCYLLA



DynamoDB



cassandra



tarantool



Bigtable



ACCUMULO™



SQLite



influxdb



QuasarDB

relational

time-series

2025

LSM-tree

NoSQL



RocksDB



levelDB



SCYLLA



DynamoDB



cassandra



tarantool



Bigtable



accumulo



SQLite



influxdb



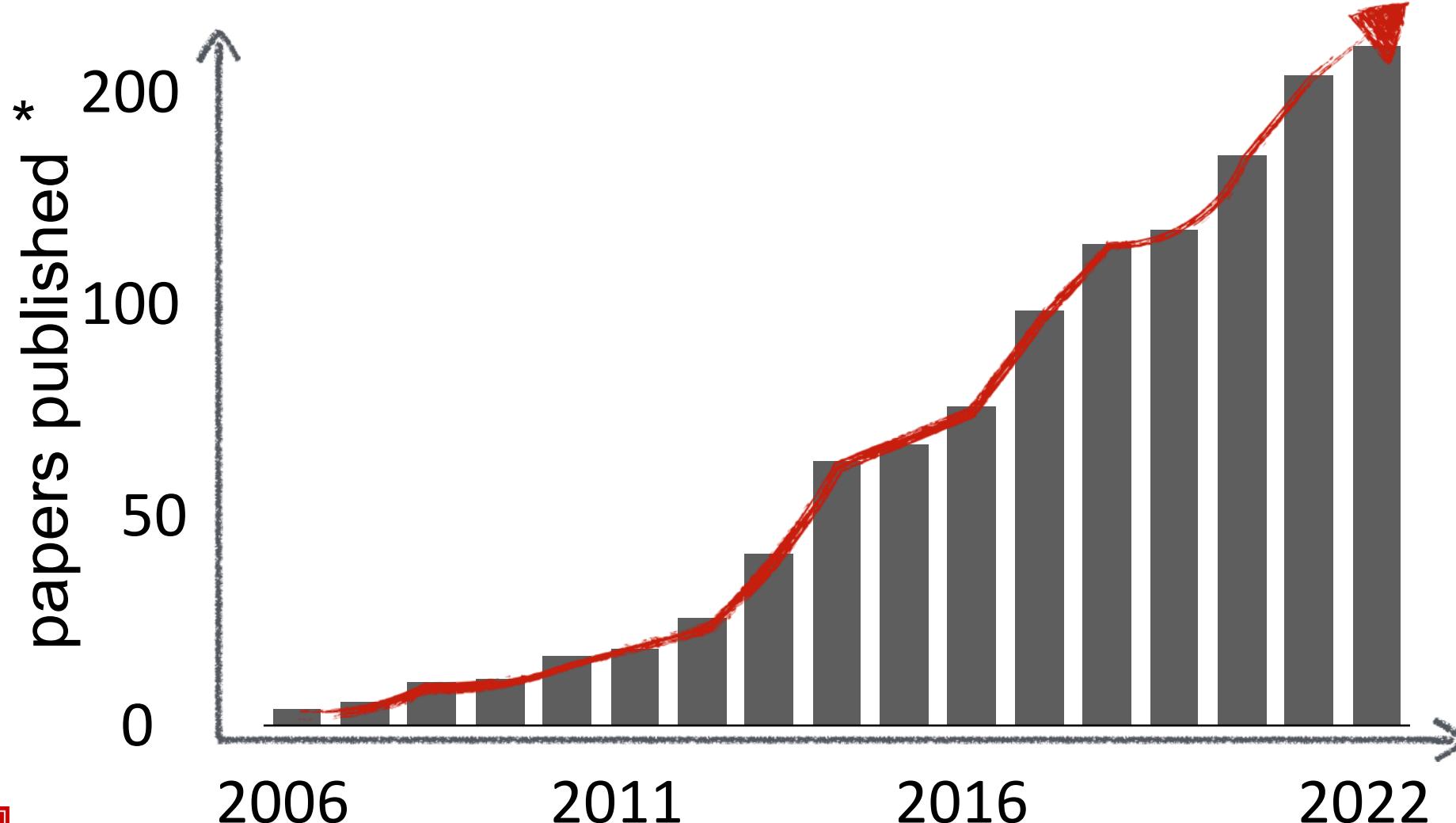
QuasarDB

relational

time-series

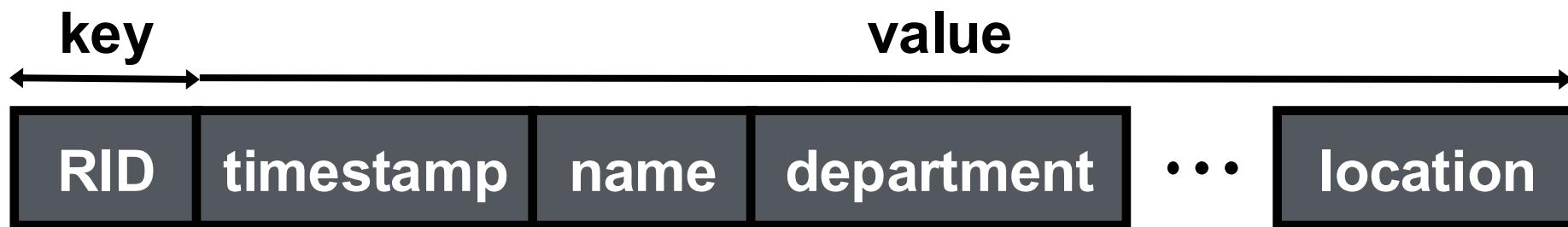
2025

Research Trend



LSM Basics

key-value
pairs



LSM Basics

key-value
pairs



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

LSM Basics



LSM Operating Principles

Buffering
ingestion

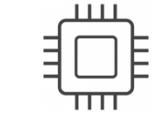
Immutable files on
storage

Out-of-place updates &
deletes

Periodic data layout
reorganization

Buffering ingestion

put(6)
put(2)

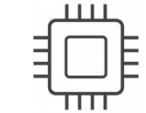


buffer



Buffering ingestion

put(1)
put(6)

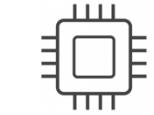


buffer

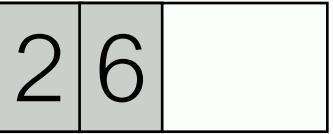


Buffering ingestion

put(4)
put(1)



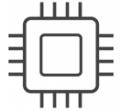
buffer



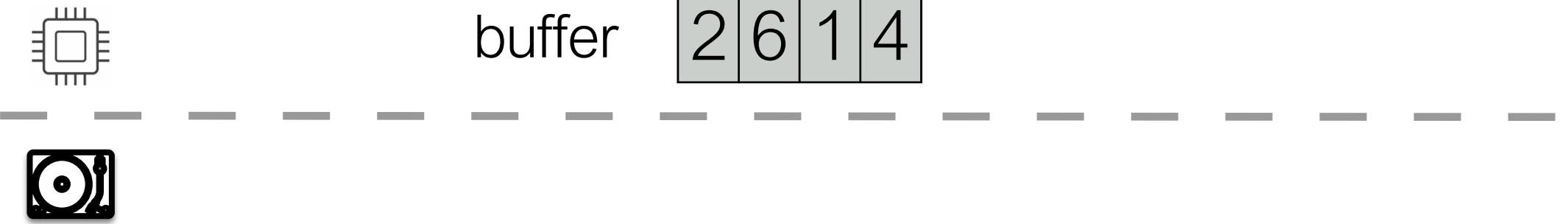
Buffering ingestion

put(4)

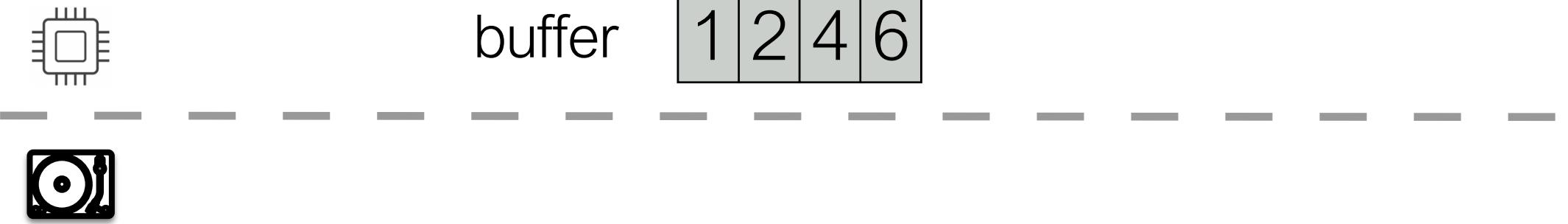
buffer



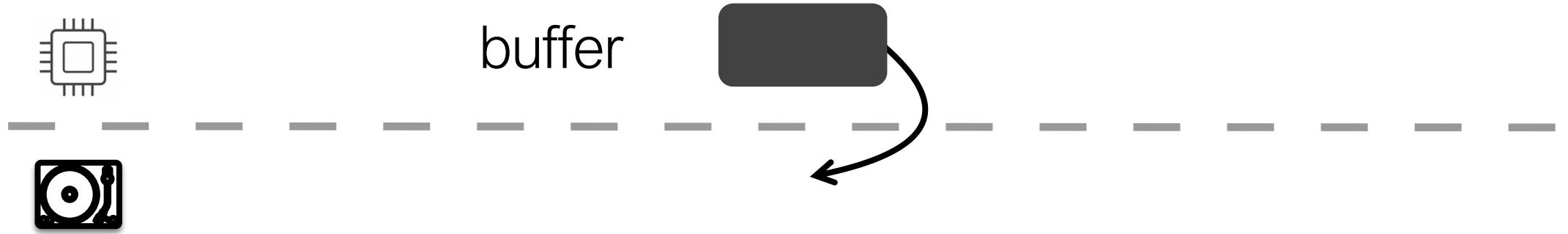
Buffering ingestion



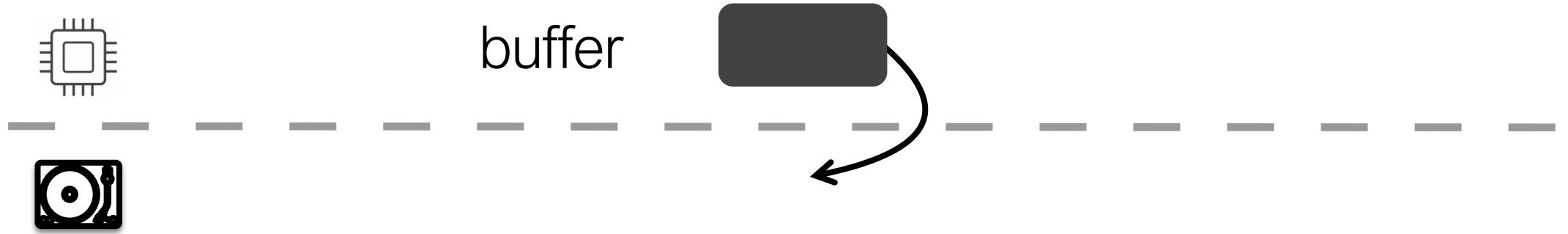
Buffering ingestion



Buffering ingestion

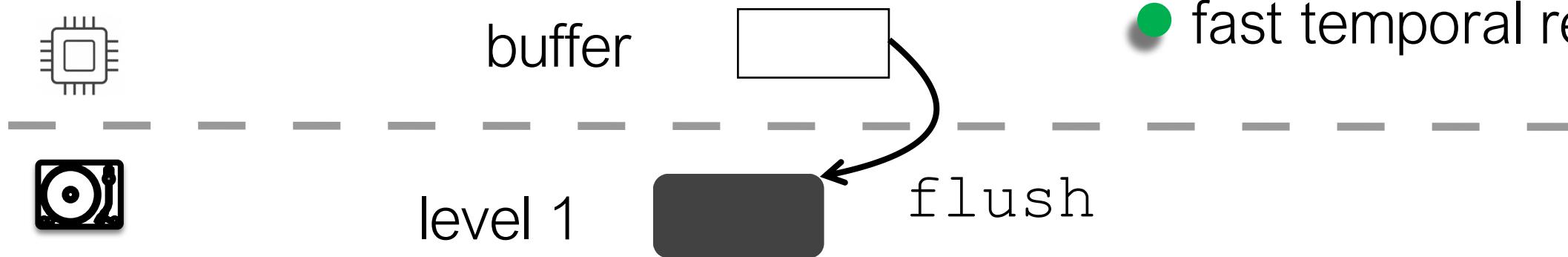


Buffering ingestion



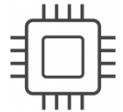
Buffering ingestion

- low ingestion cost
- fast temporal reads

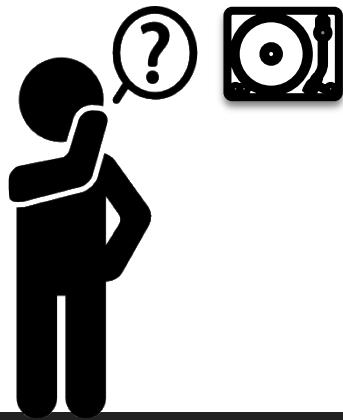


Immutable files on storage

- compact storage
- good ingestion throughput



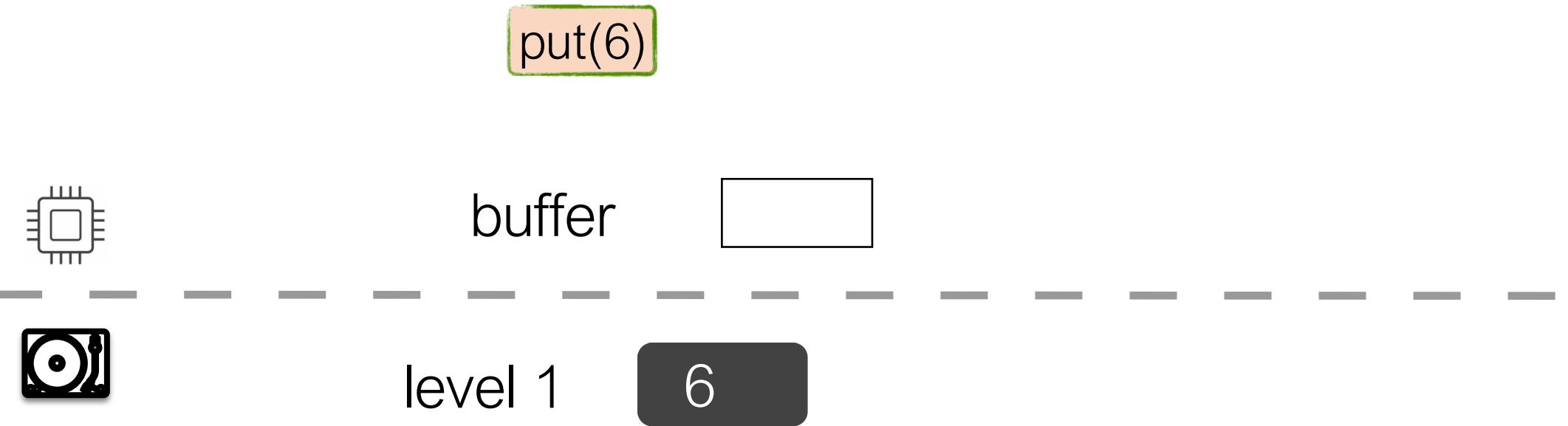
buffer

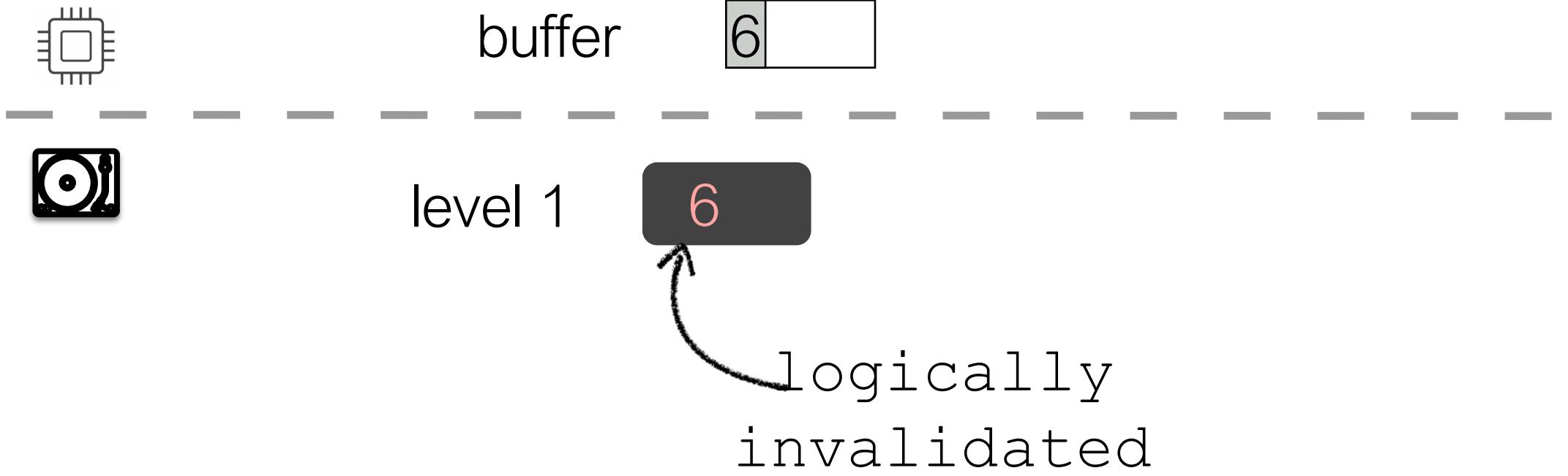


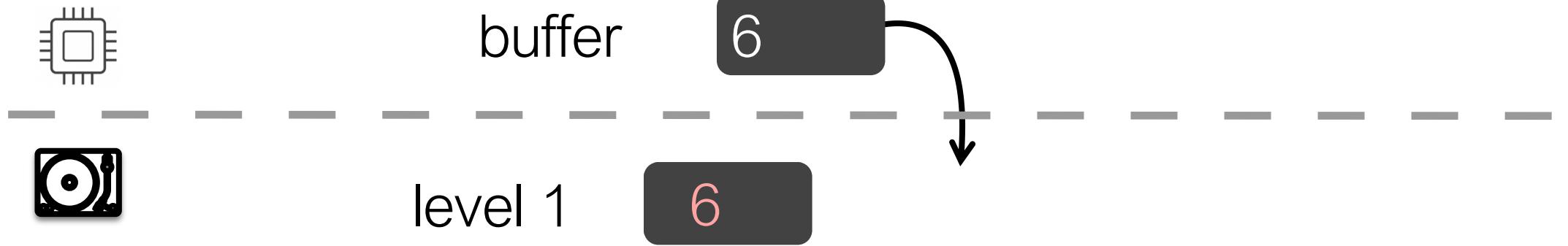
level 1

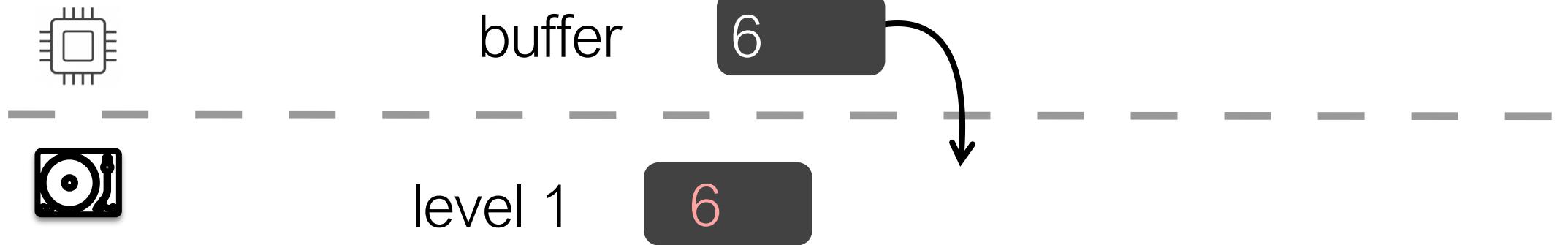


How do we update data?

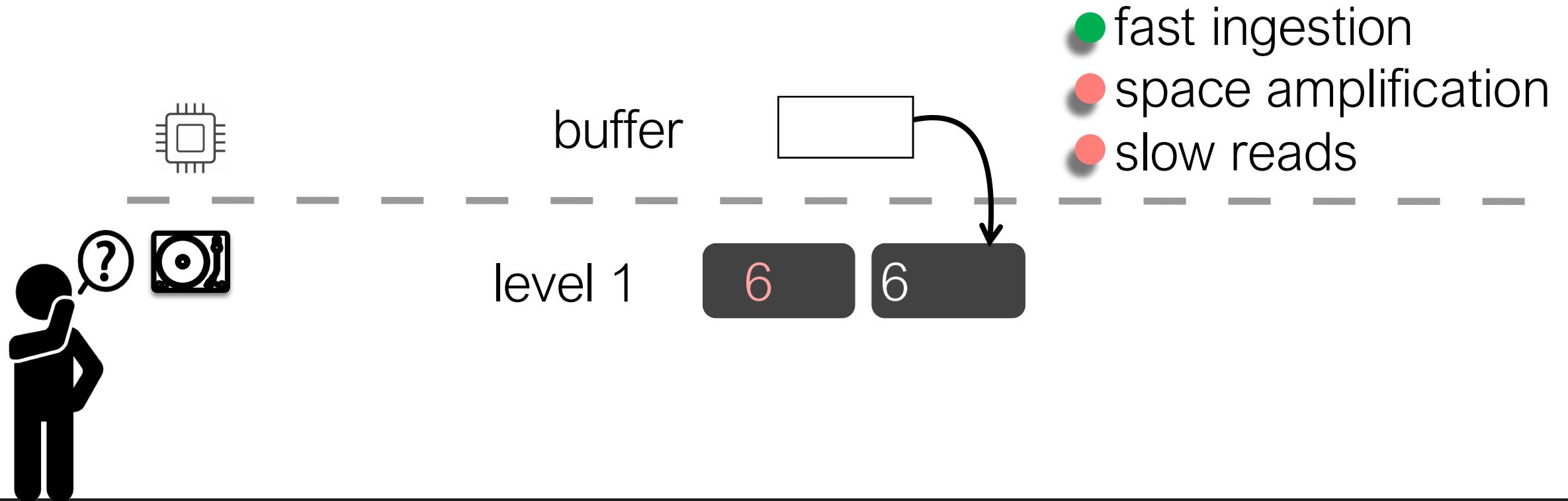




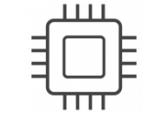




Out-of-place updates



How do we reduce this space amplification?

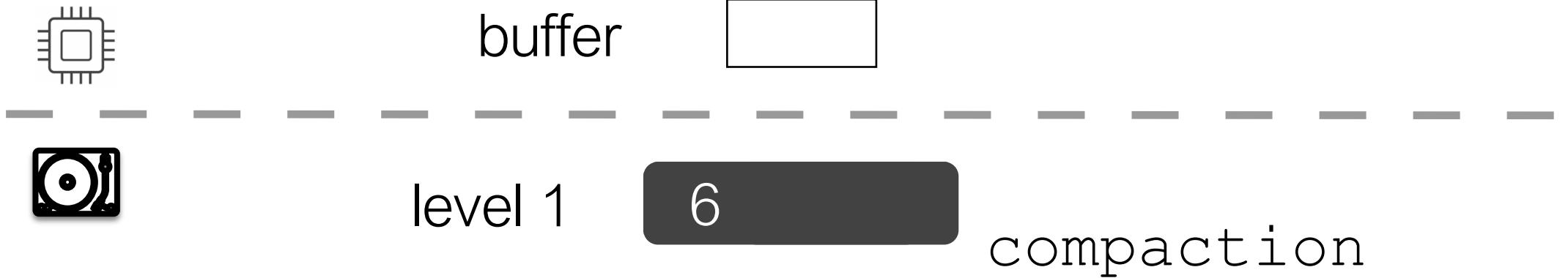


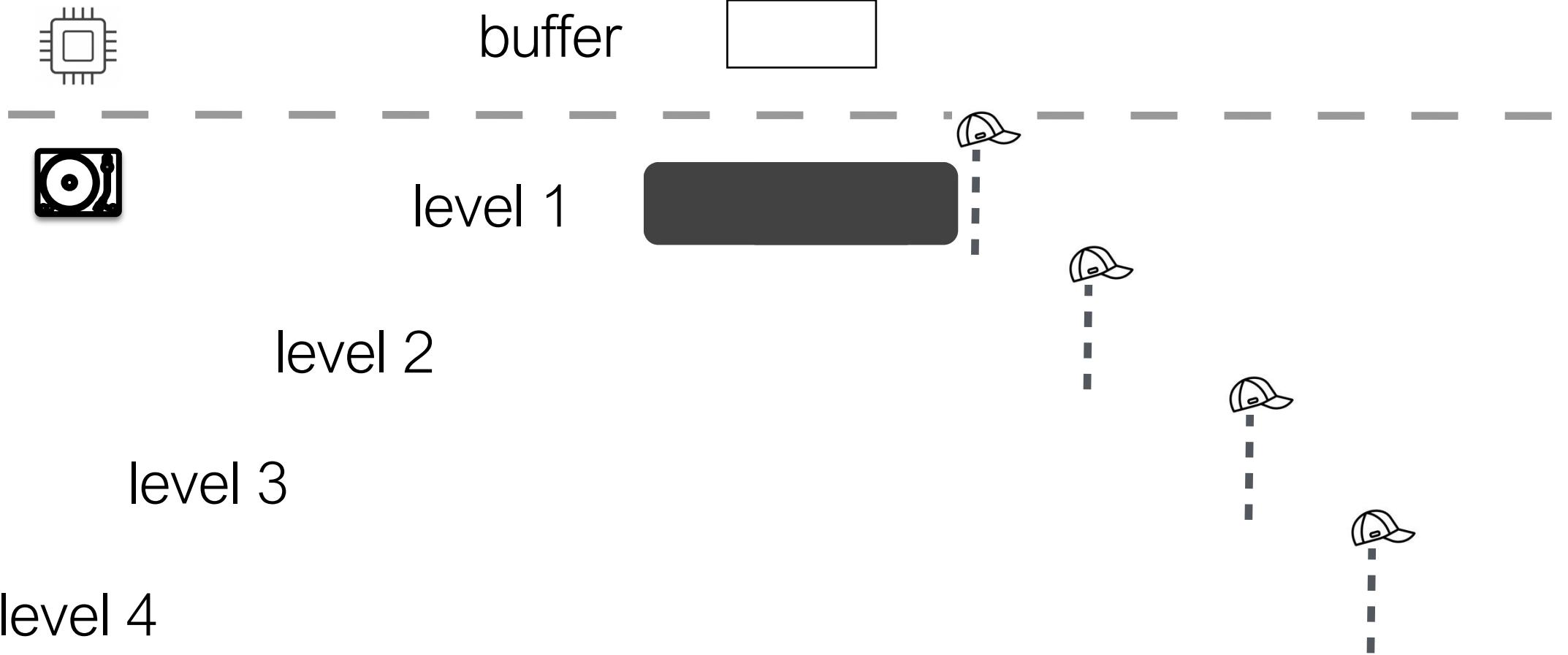
buffer

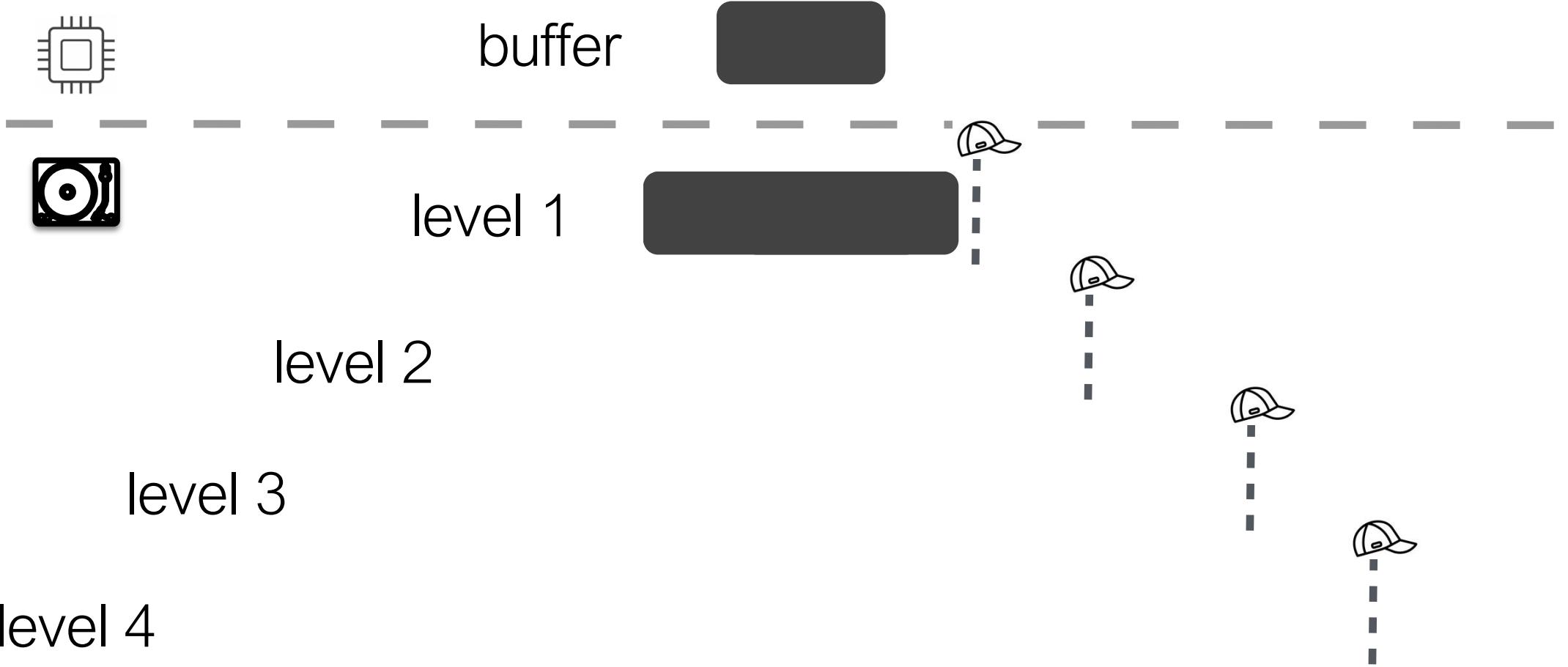


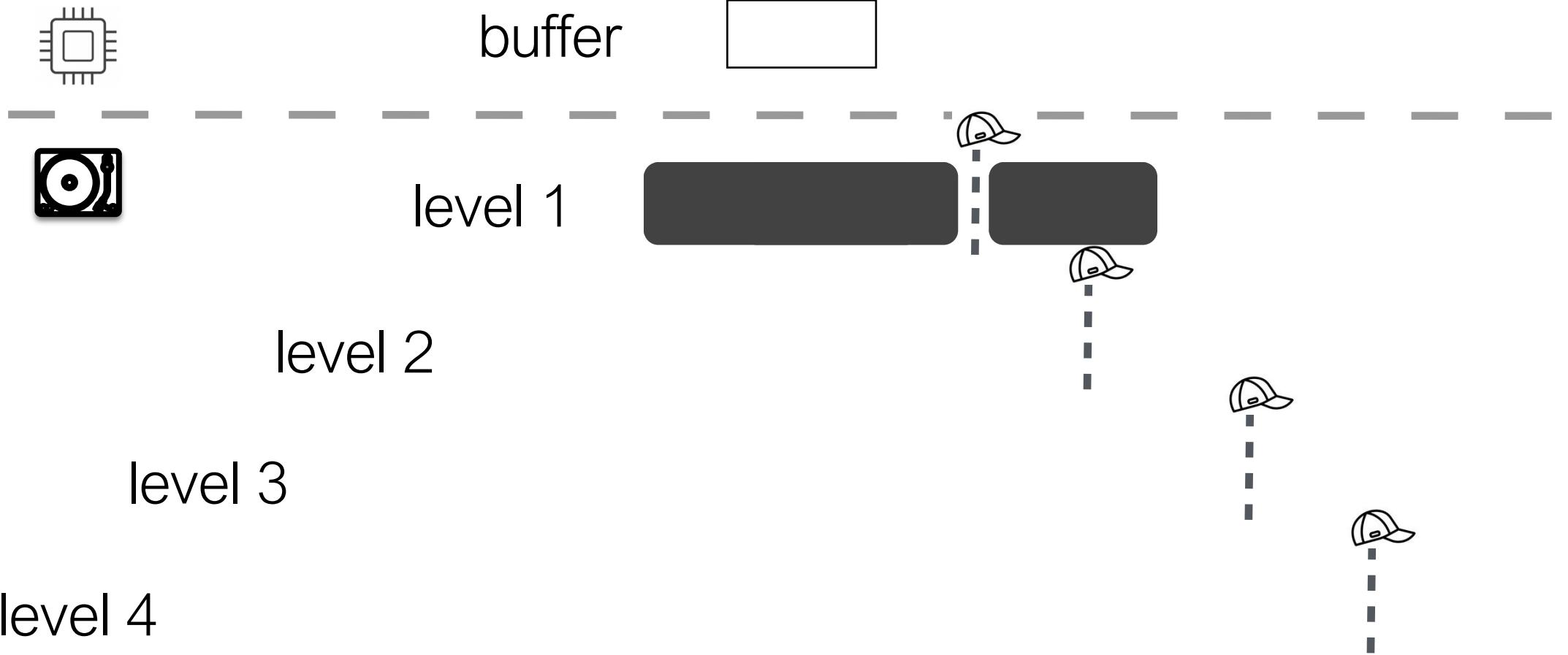
level 1

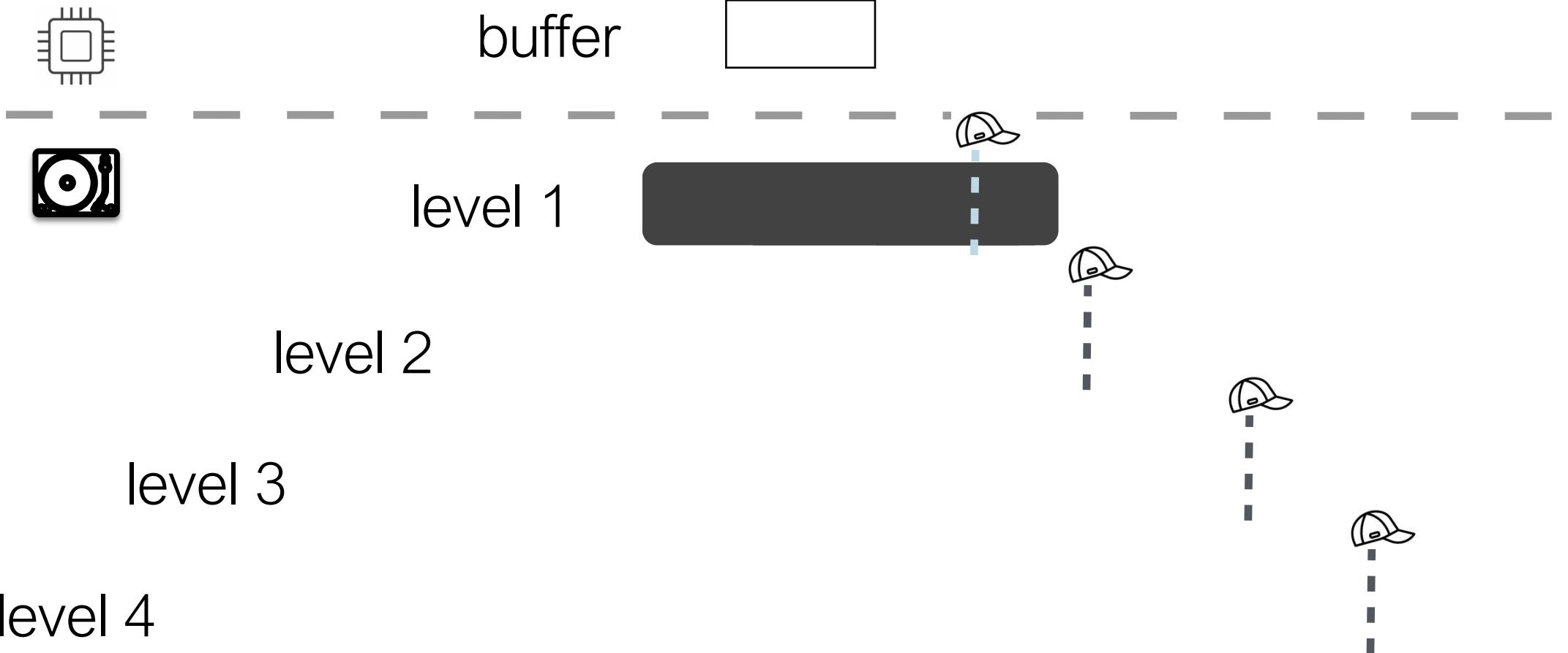


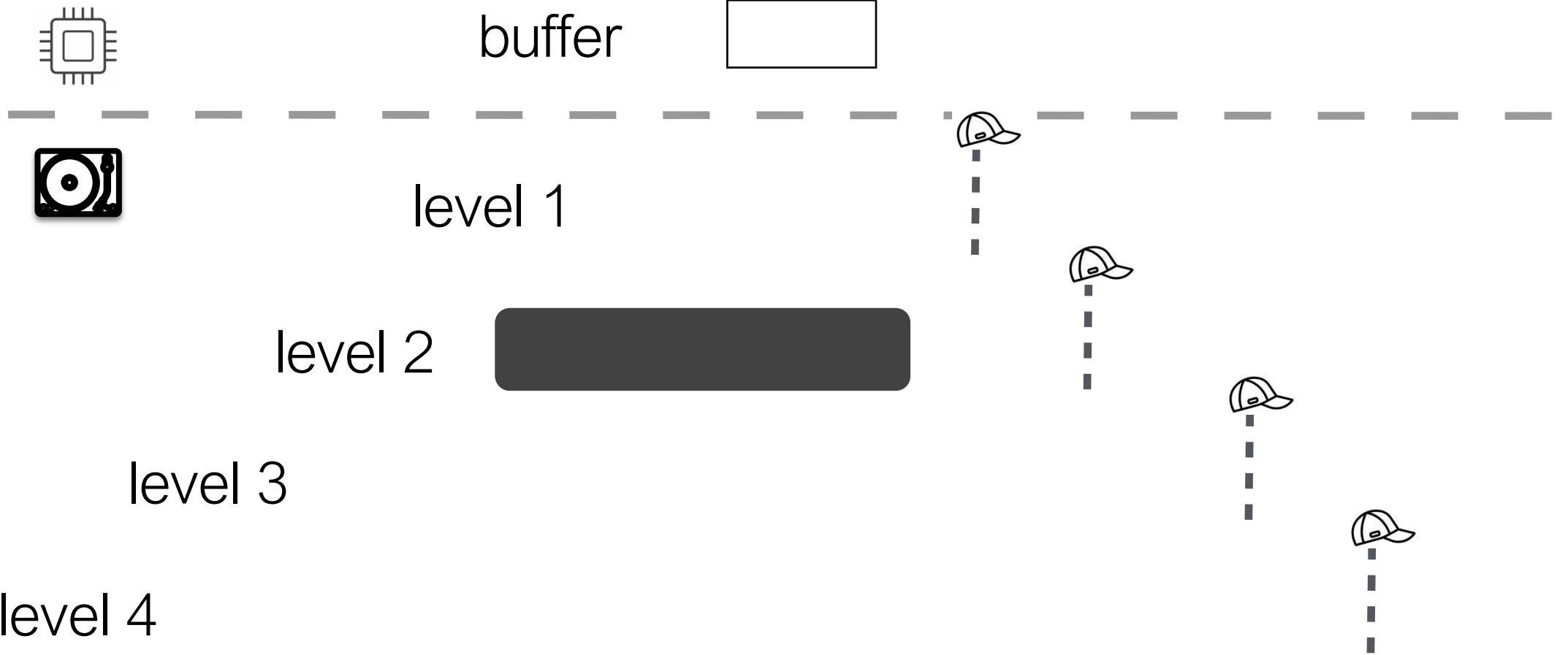


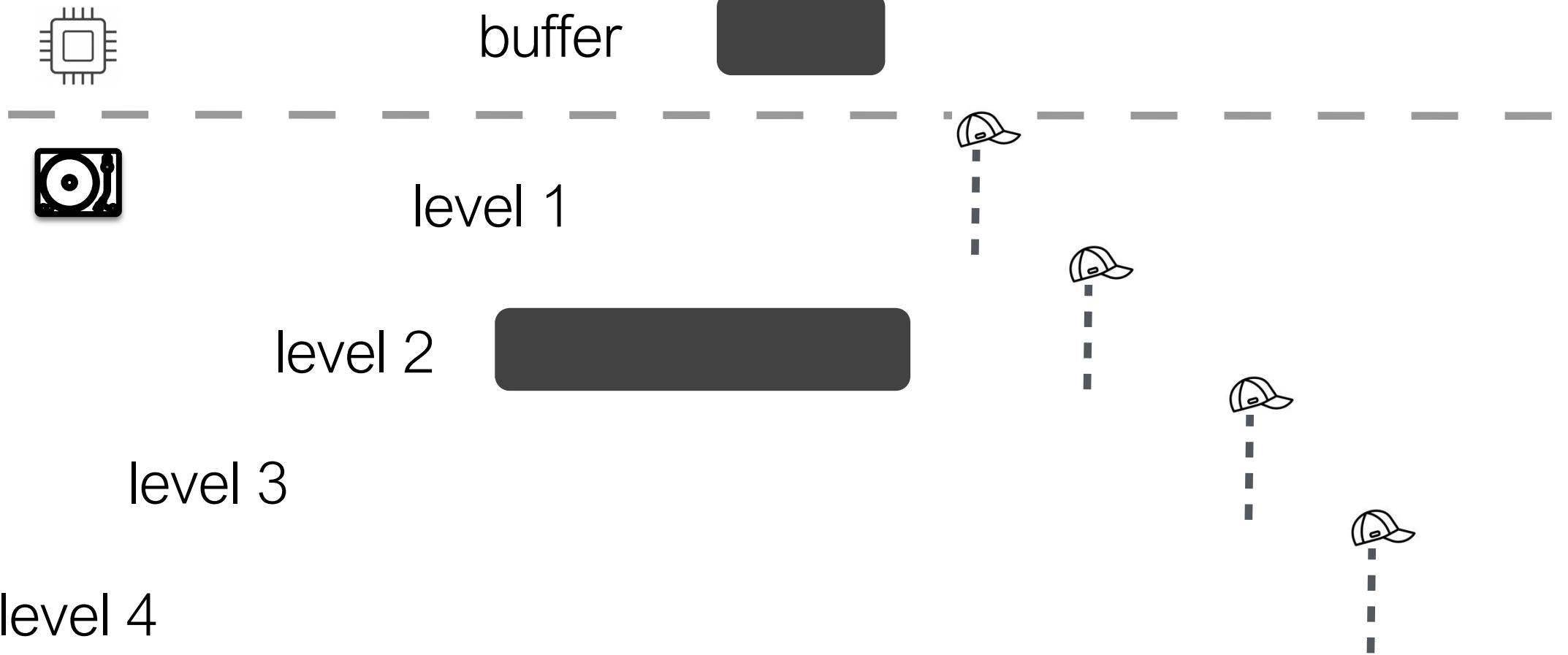


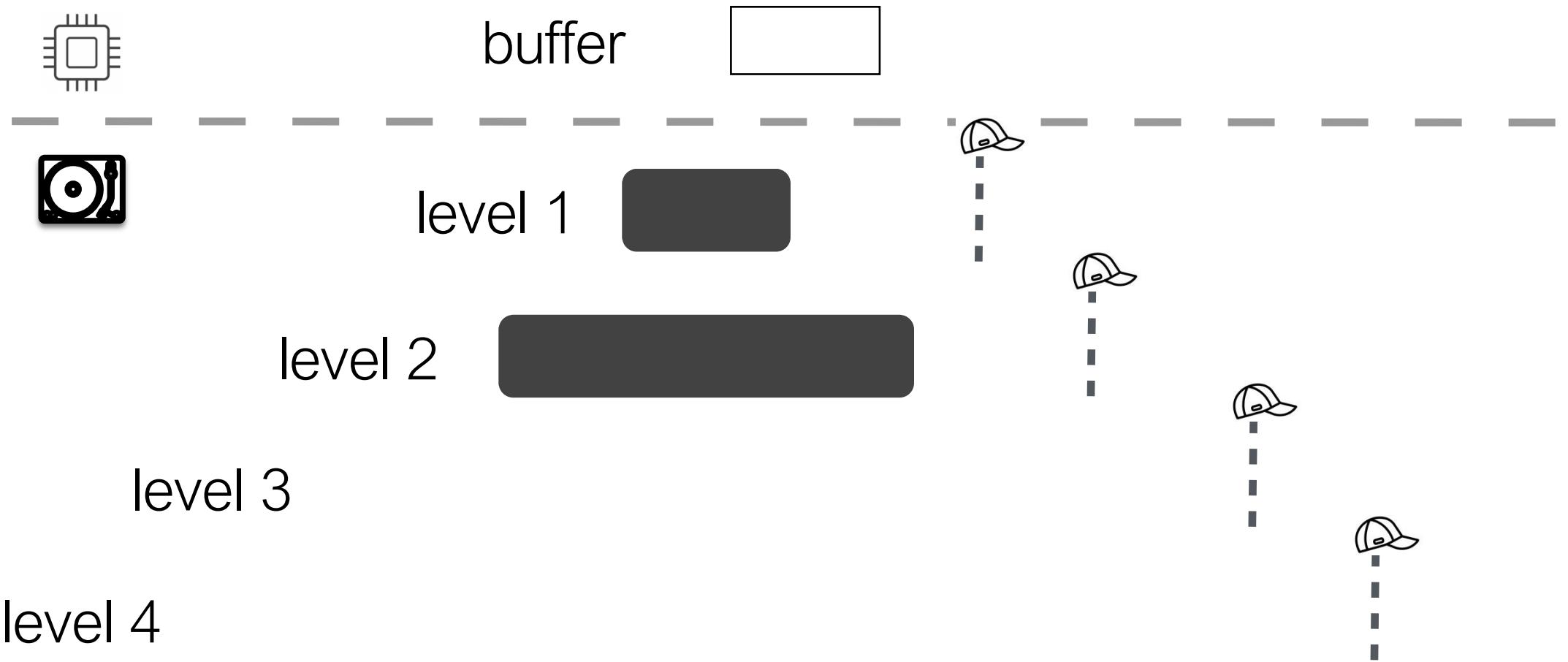


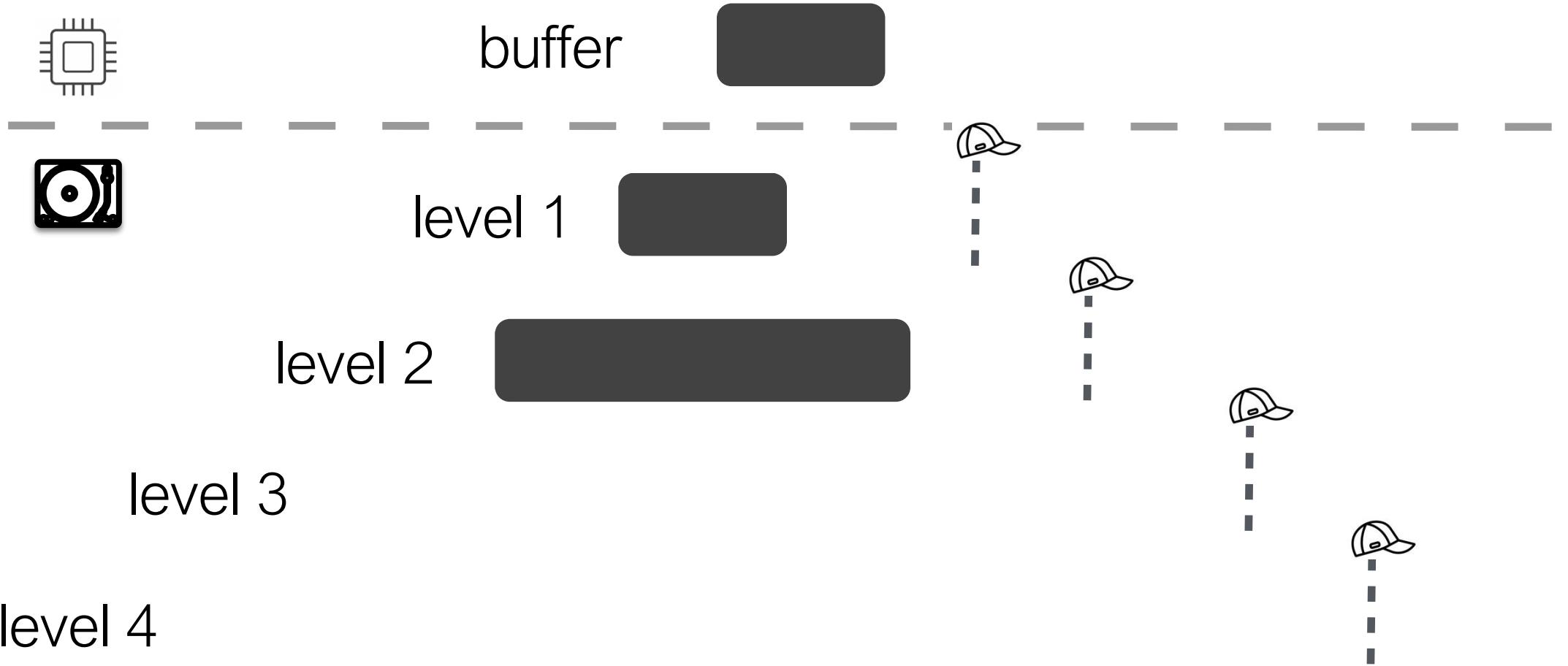


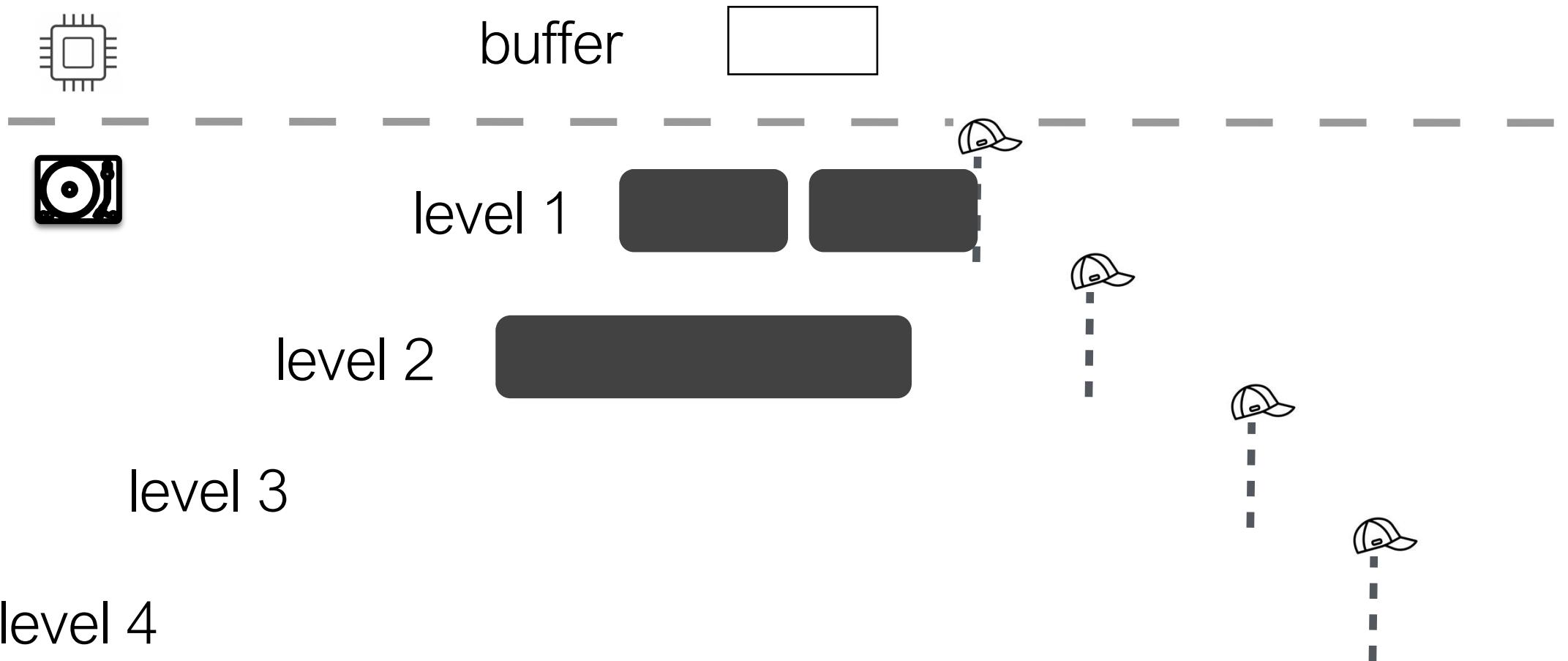


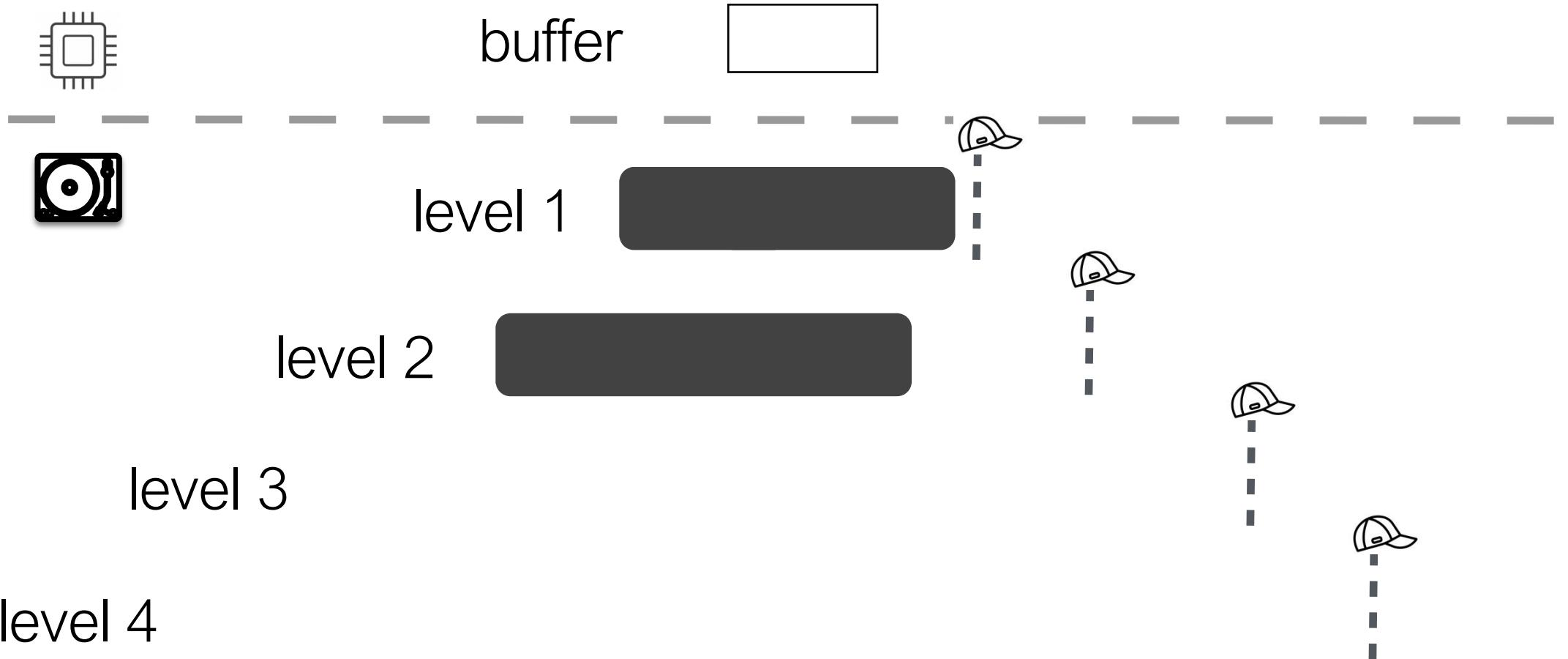


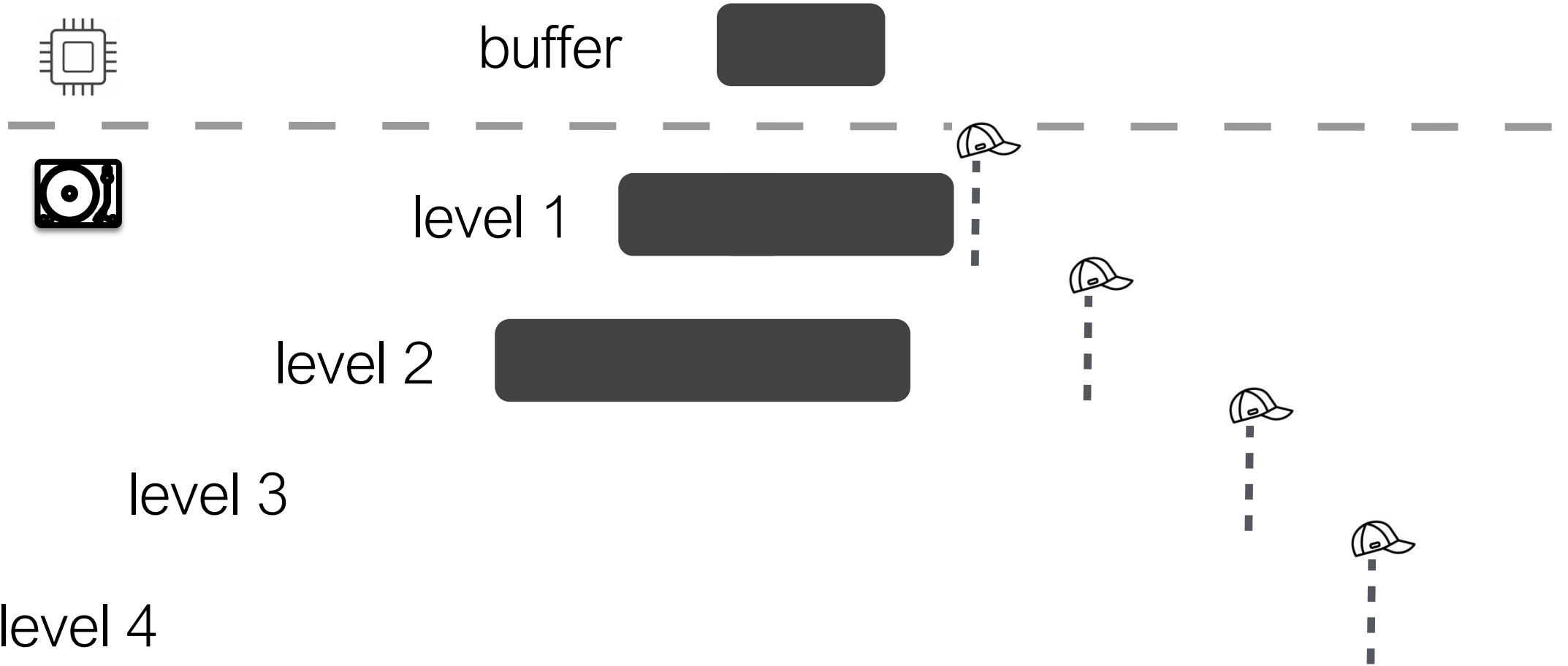


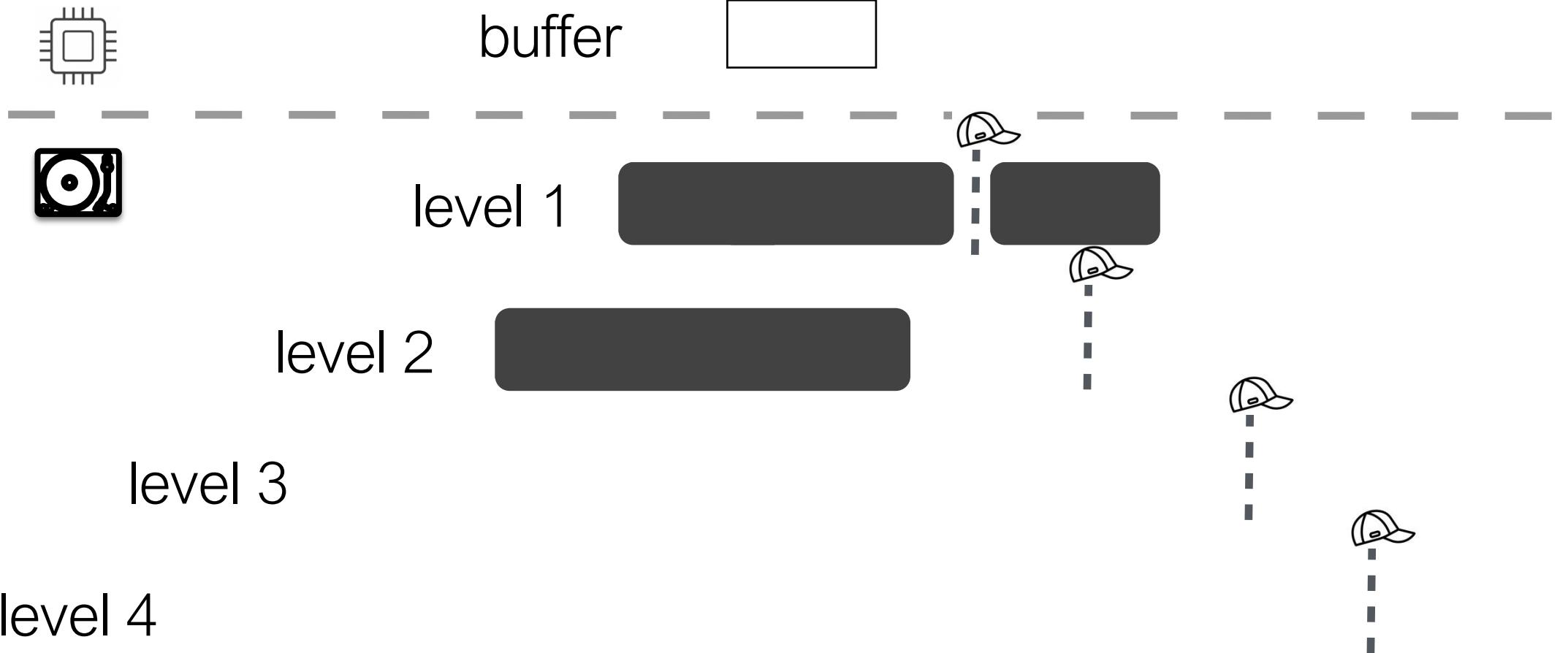


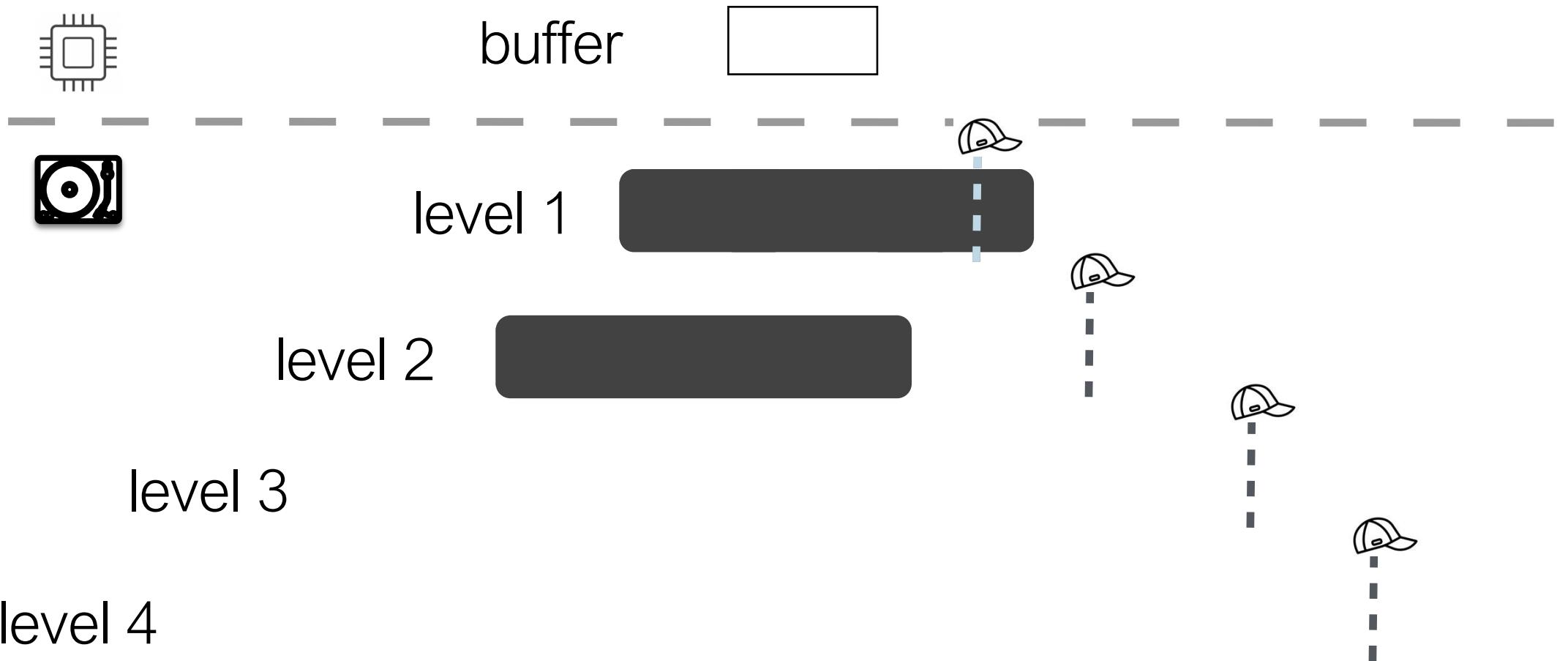


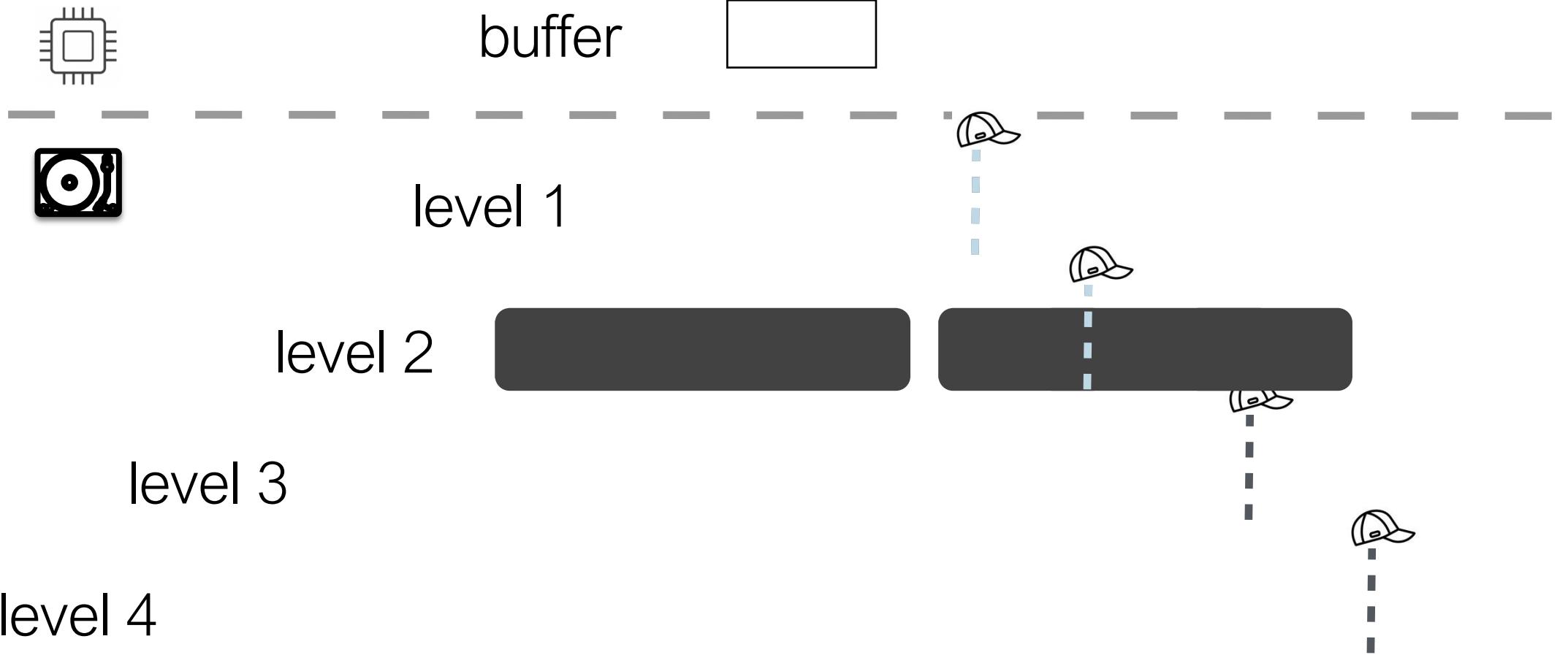




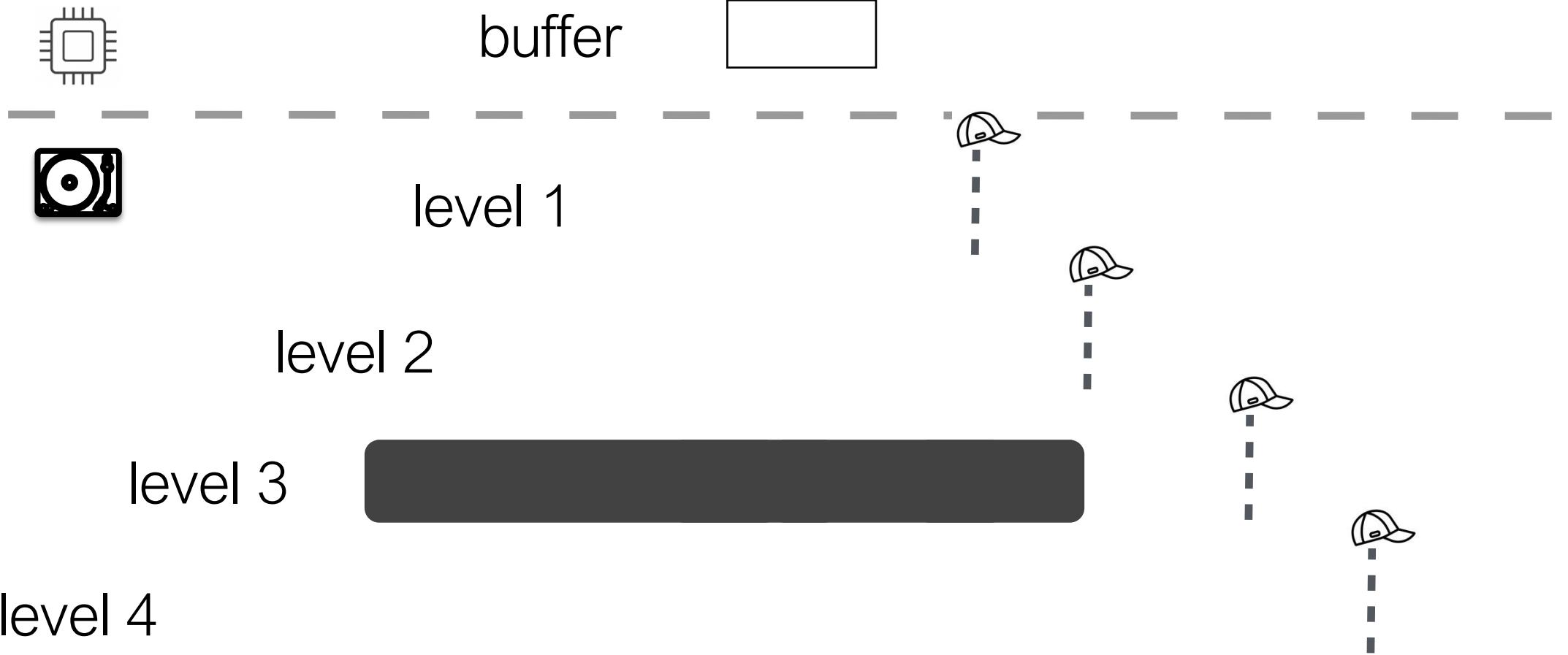








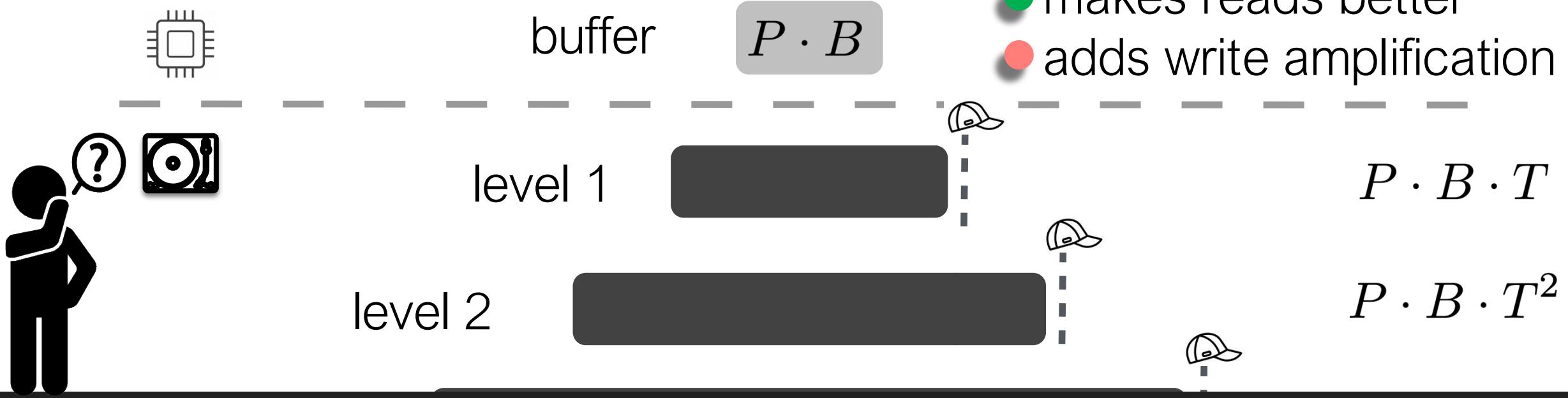




P : pages in
 B : buffer
 L : #levels
 T : size ratio

Periodic compactions

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



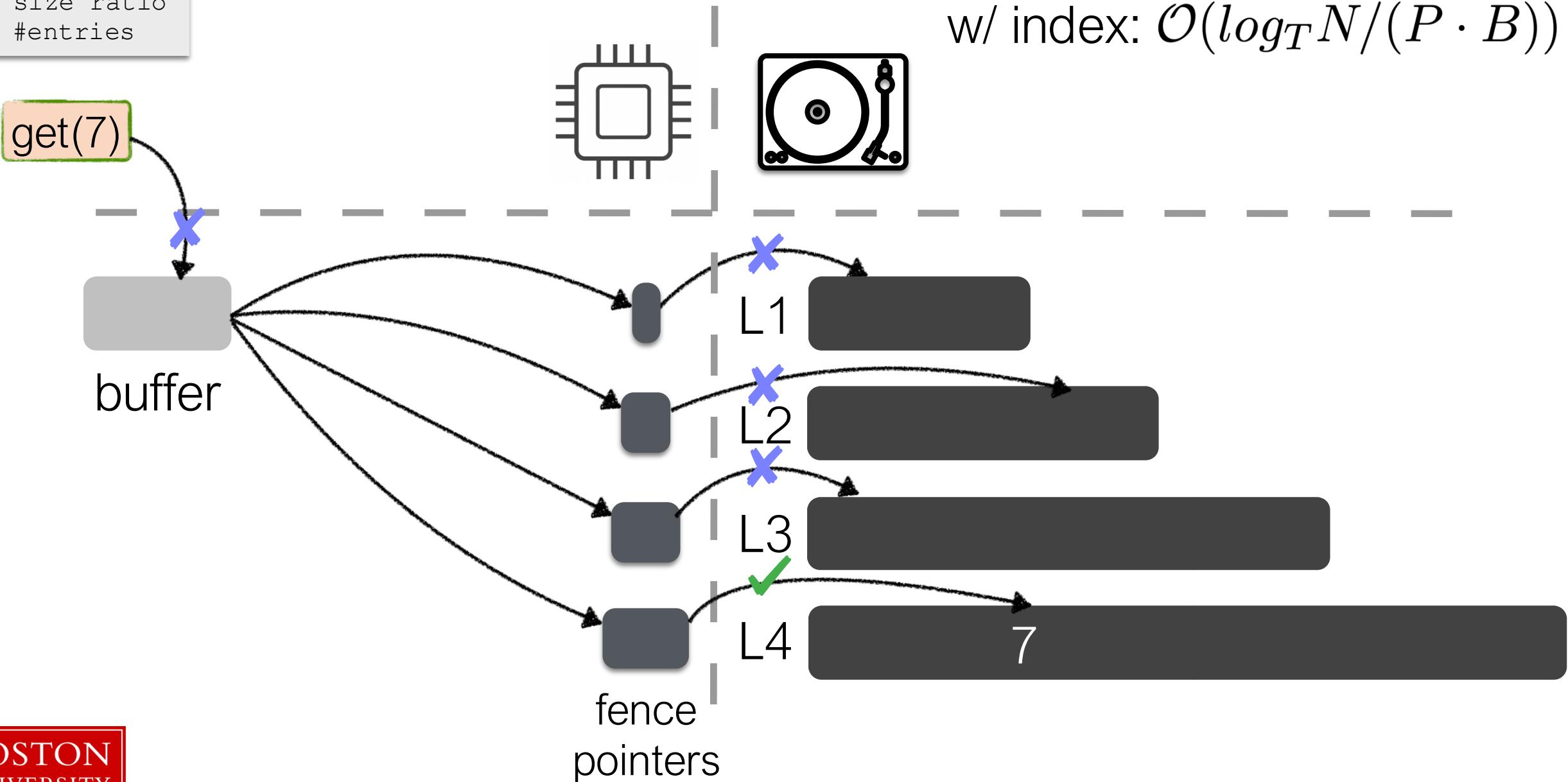
How about queries?

capacity
(entries)

Cost analysis

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

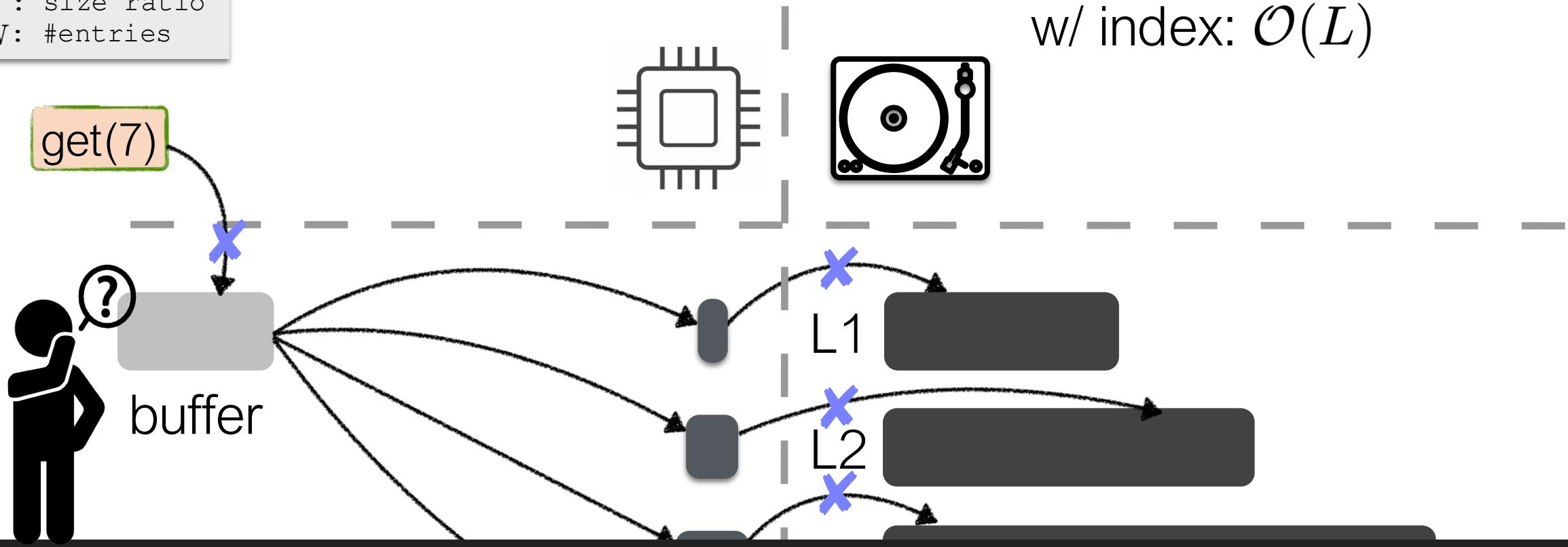
P : pages in
 B : buffer
 L : #levels
 T : size ratio
 N : #entries



Cost analysis

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

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



How to avoid unnecessary I/Os?

pointers

Cost analysis

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

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

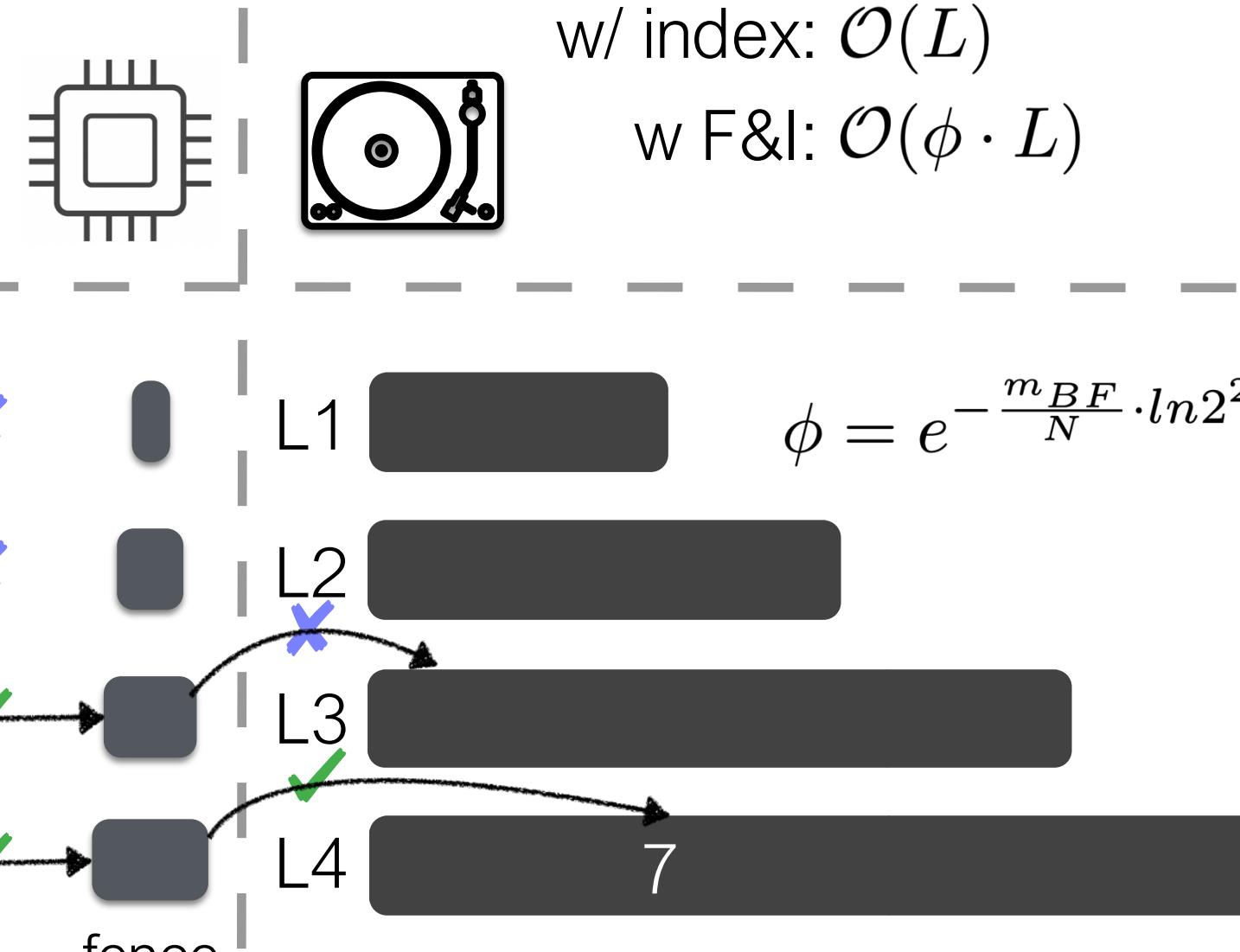
$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2}$$

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

get(7)

buffer

Bloom filters fence pointers



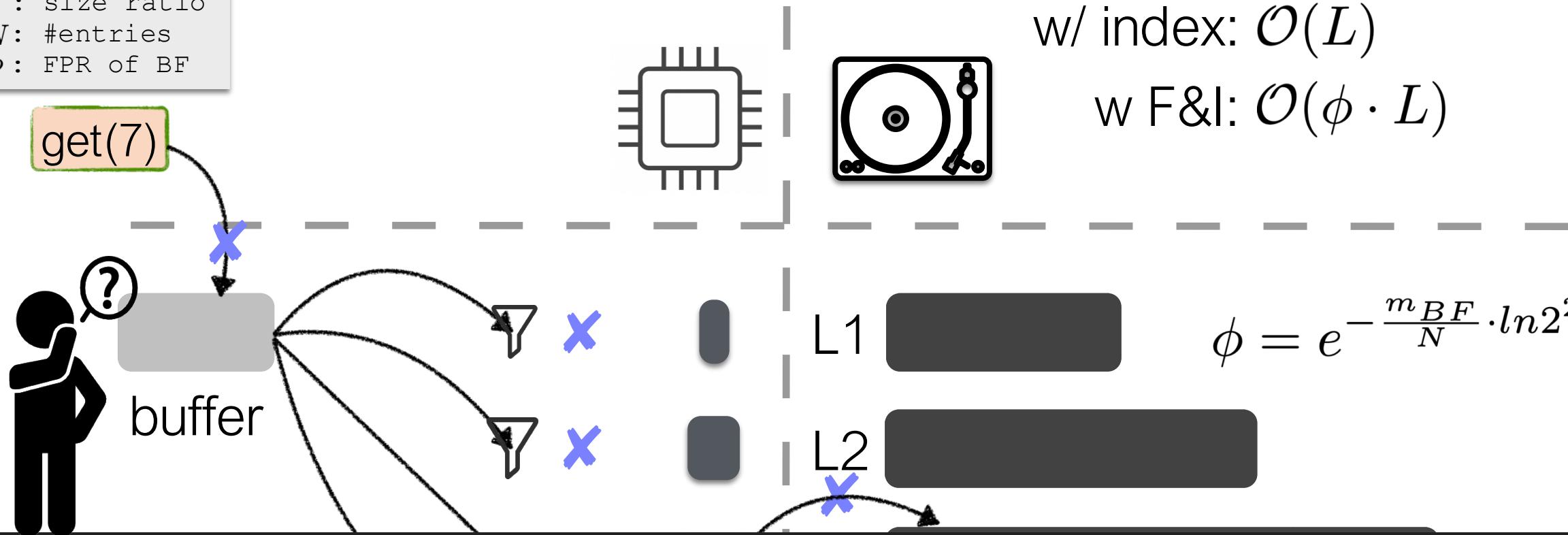
P : pages in
 B : buffer
 L : entries/page
 T : #levels
 N : size ratio
 ϕ : FPR of BF

Cost analysis

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

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

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



How to manage memory?

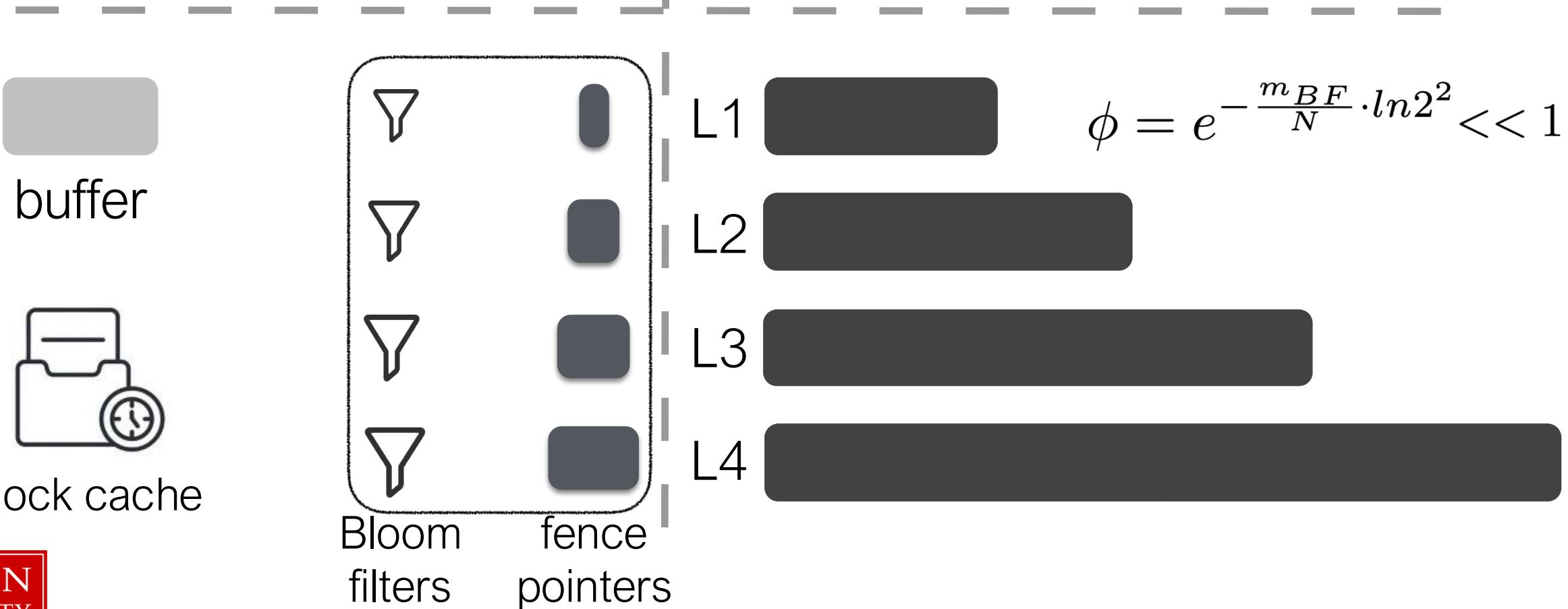
Bloom filters
pointers

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

Cost analysis

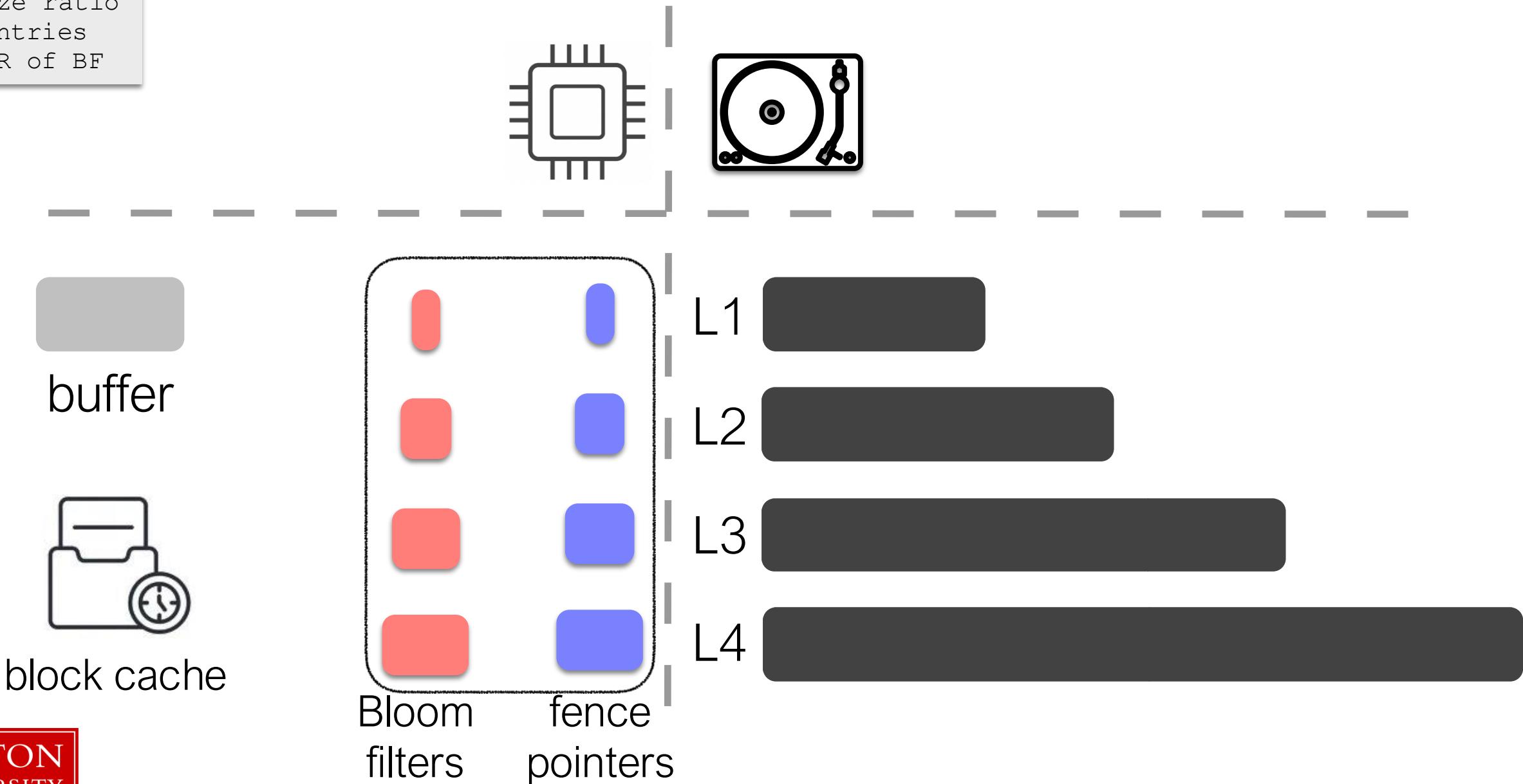
w/ index: $\mathcal{O}(L)$

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



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

Block Cache

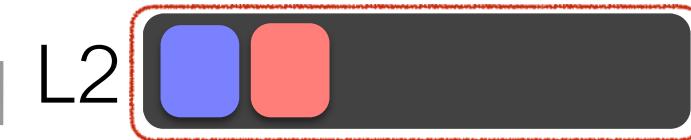
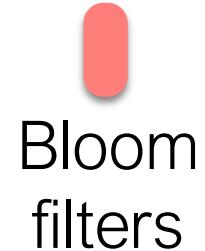
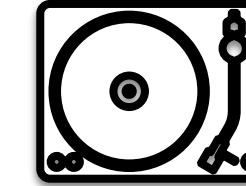
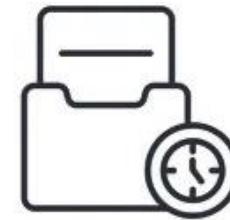
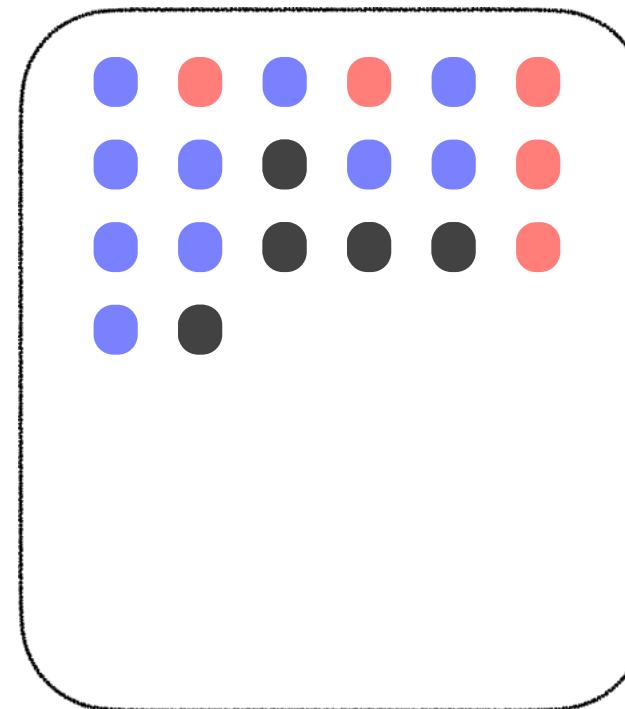


P : pages in
 B_{buffer}
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Block Cache

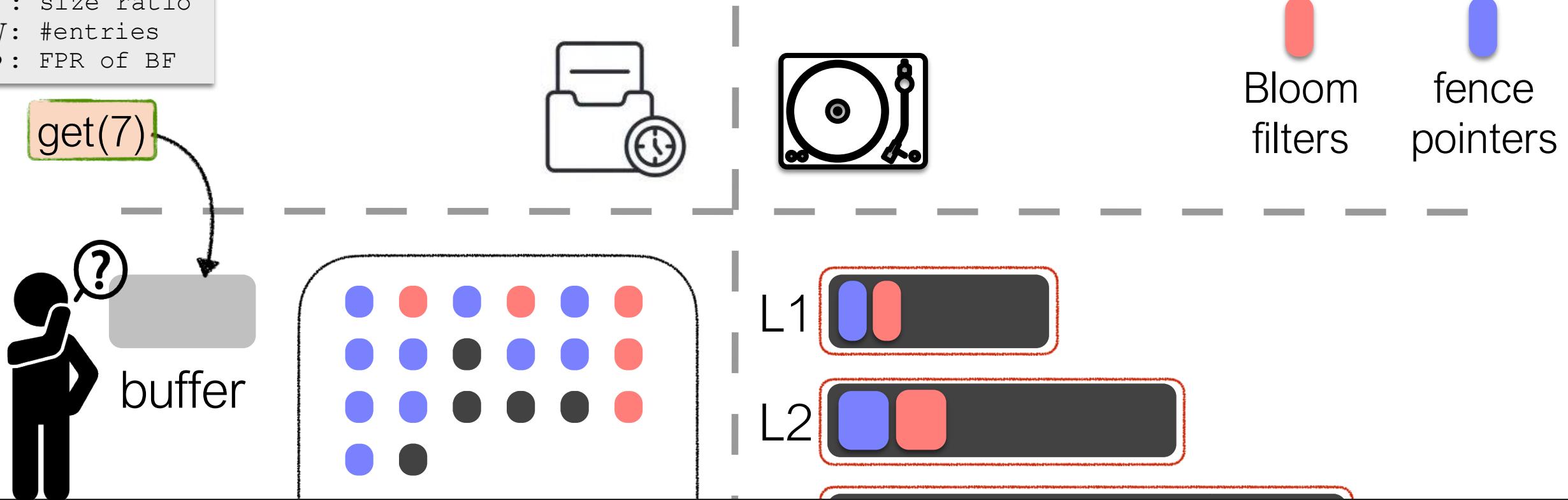
get(7)

buffer



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

Block Cache

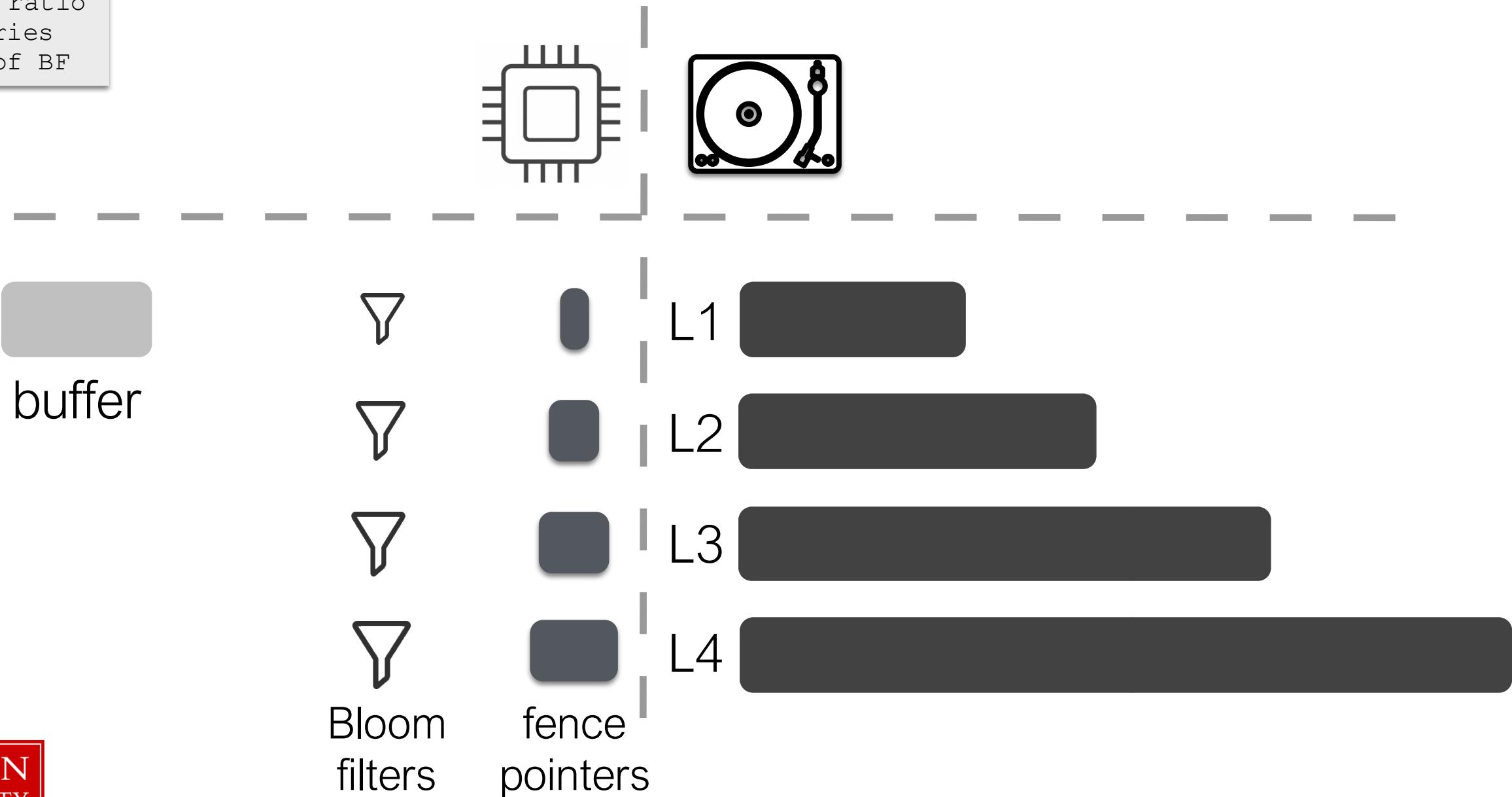


What about range queries?

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

s selectivity
LRQ

Range Queries

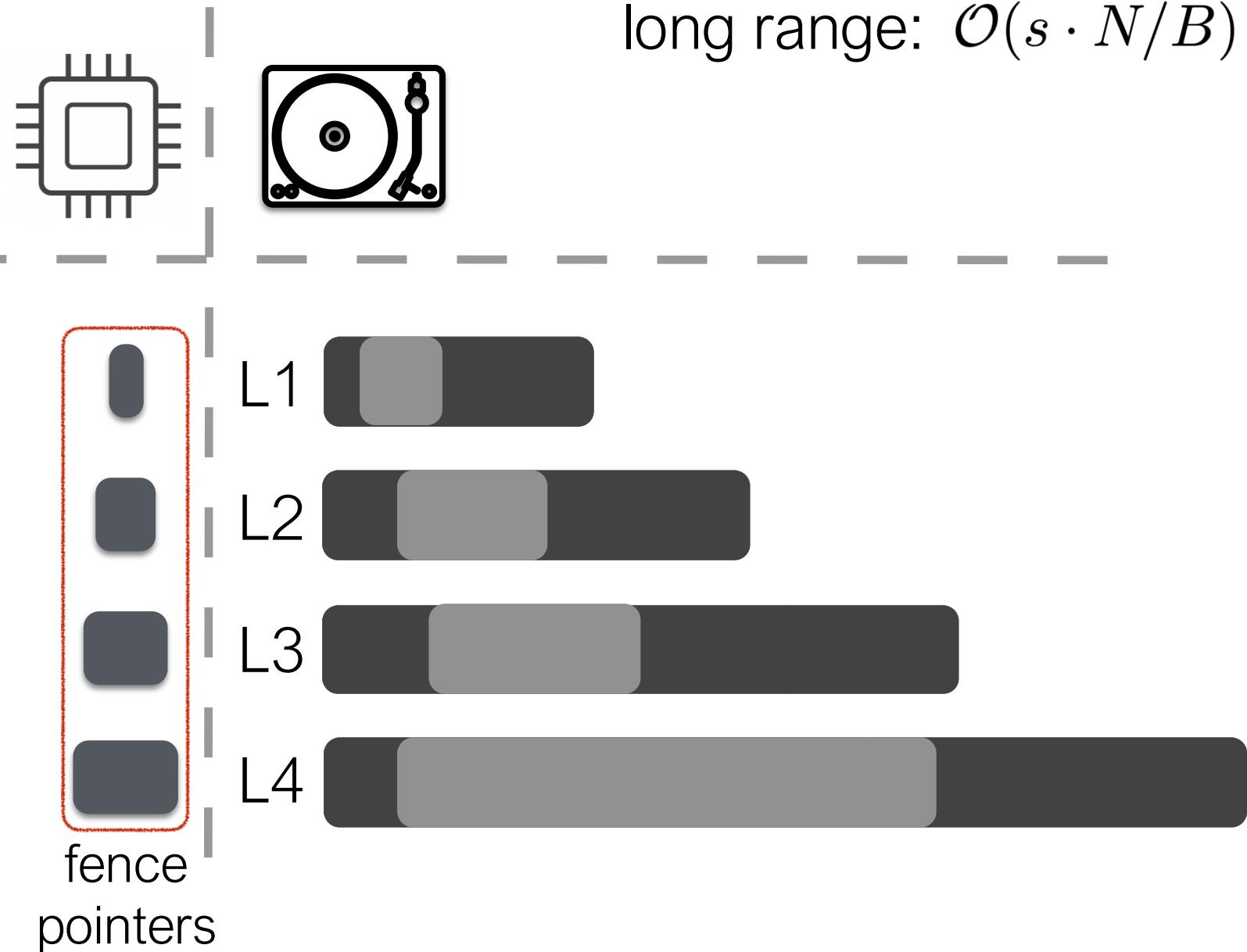


P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

s selectivity
LRQ

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

get(9,90)



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

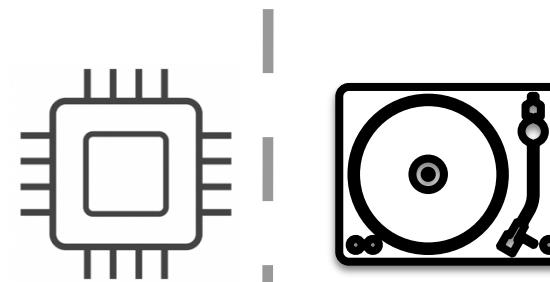
s selectivity
SRQ

Range Queries

Cost analysis

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

get(9,15)



buffer



Bloom filters



L1



L2



L3



L4

fence pointers

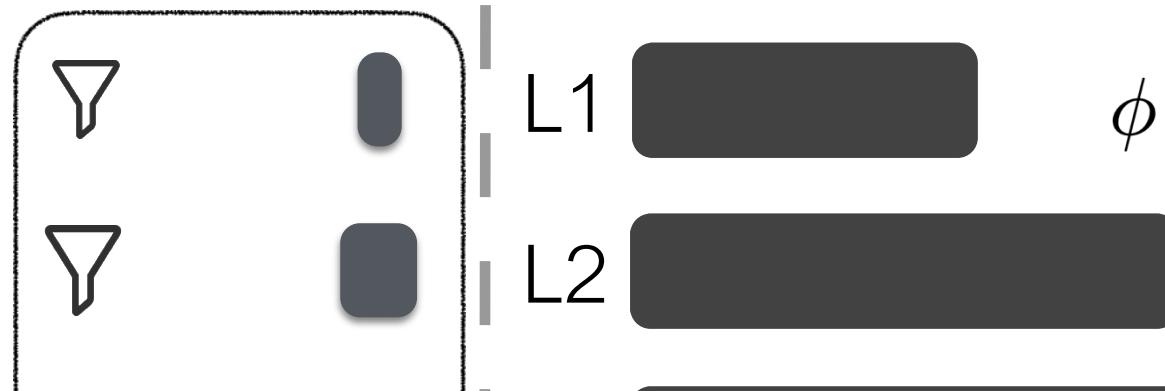
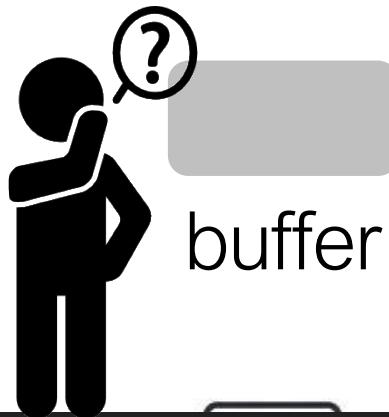
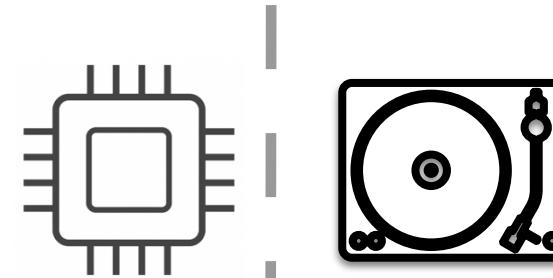
P : pages in
 B_{buffer}
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Going back to point queries

Cost analysis

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

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

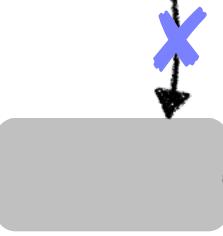


$$\phi = e^{-\frac{m_{BF}}{N} \cdot \ln 2^2} \ll 1$$

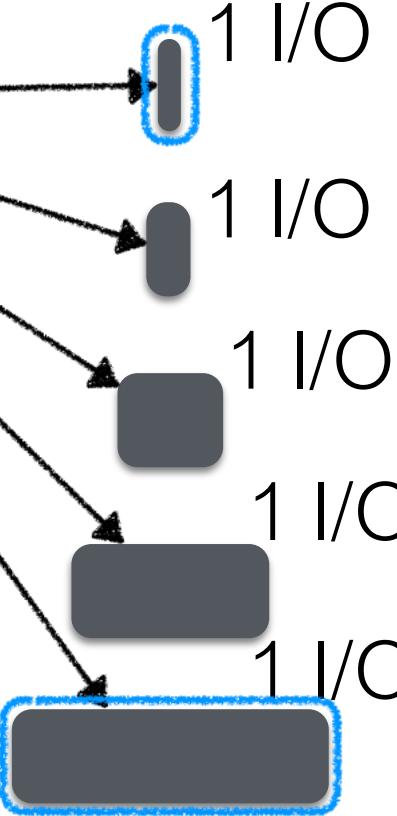
Should all BFs be equally accurate?

Bloom filters Index pointers

get(7)



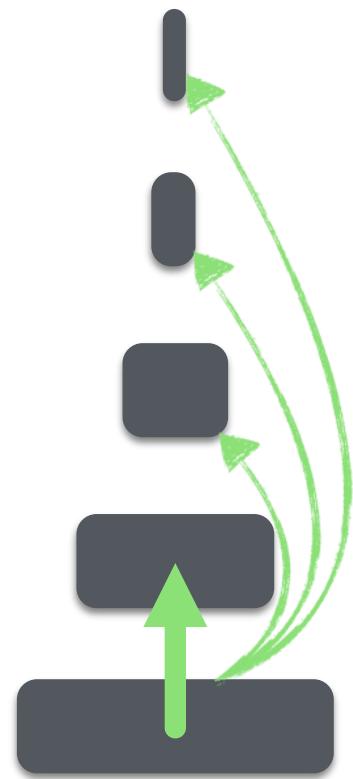
get(7)



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



Bloom filters

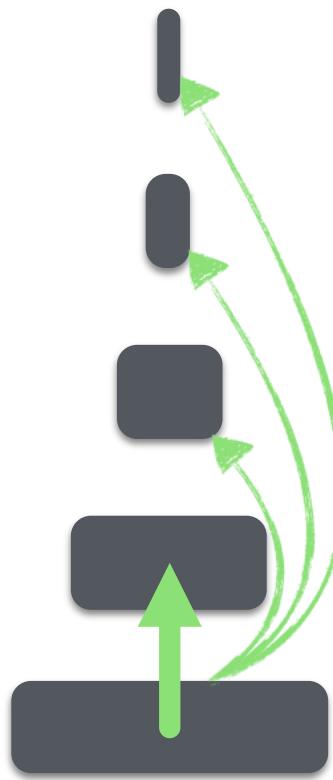


FPR Monkey FPR

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

↑
exponentially decreasing

Bloom
filters



FPR

Monkey
FPR

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

exponentially decreasing

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

point lookup
cost

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

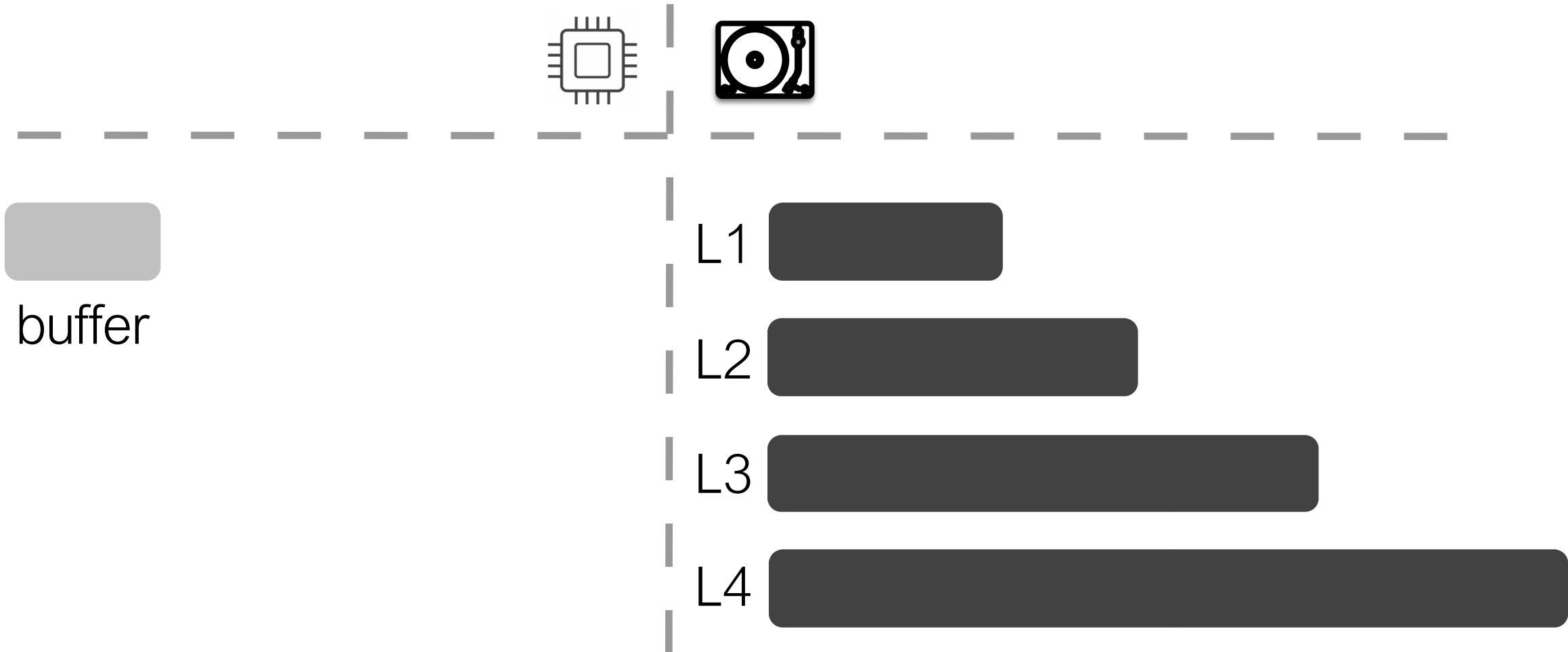


$$\mathcal{O}(c \cdot \phi_0)$$

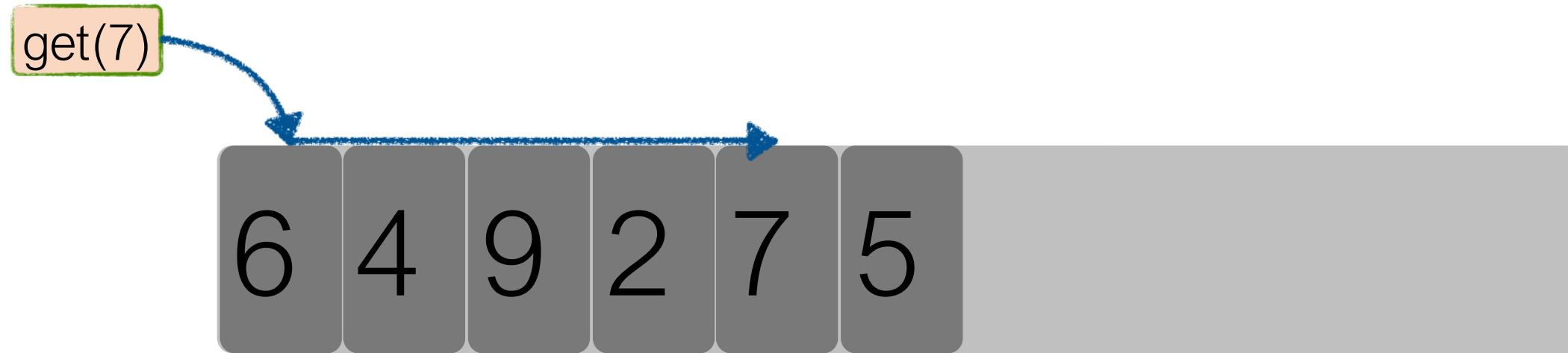


Buffer Optimizations

Buffer Optimizations



Buffer Implementation: **vector**



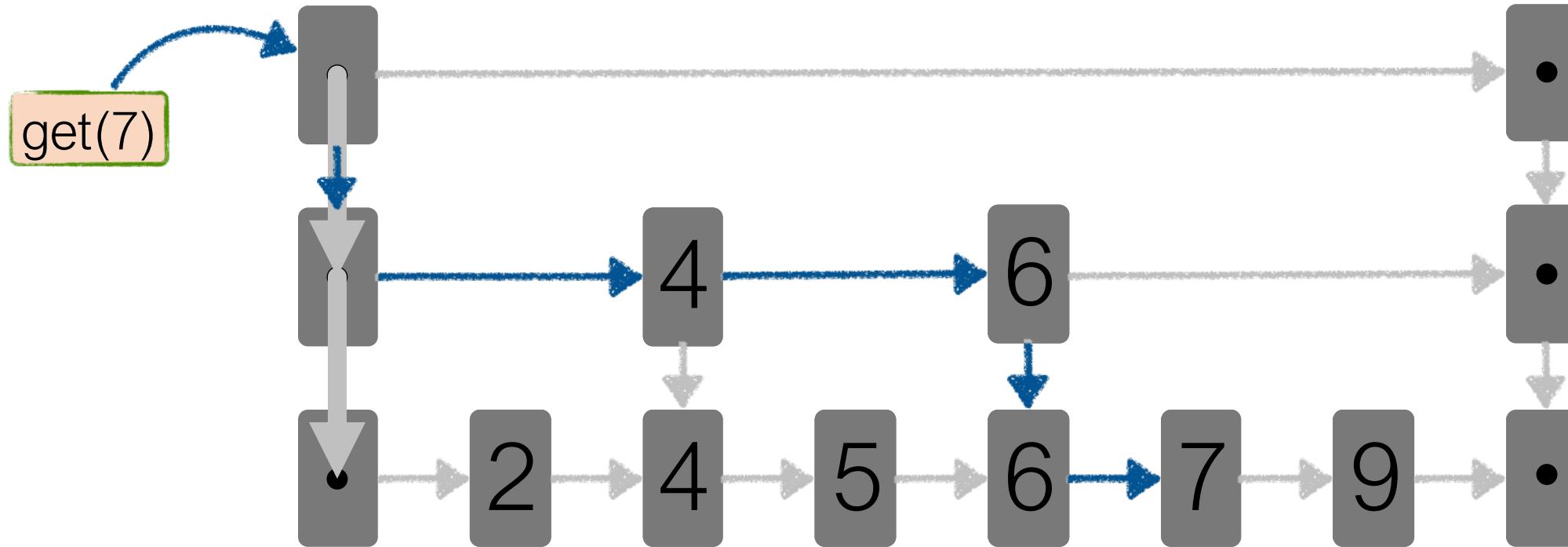
- great for ingestion-heavy w/l
- no extra space needed
- expensive points queries

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



Buffer Implementation

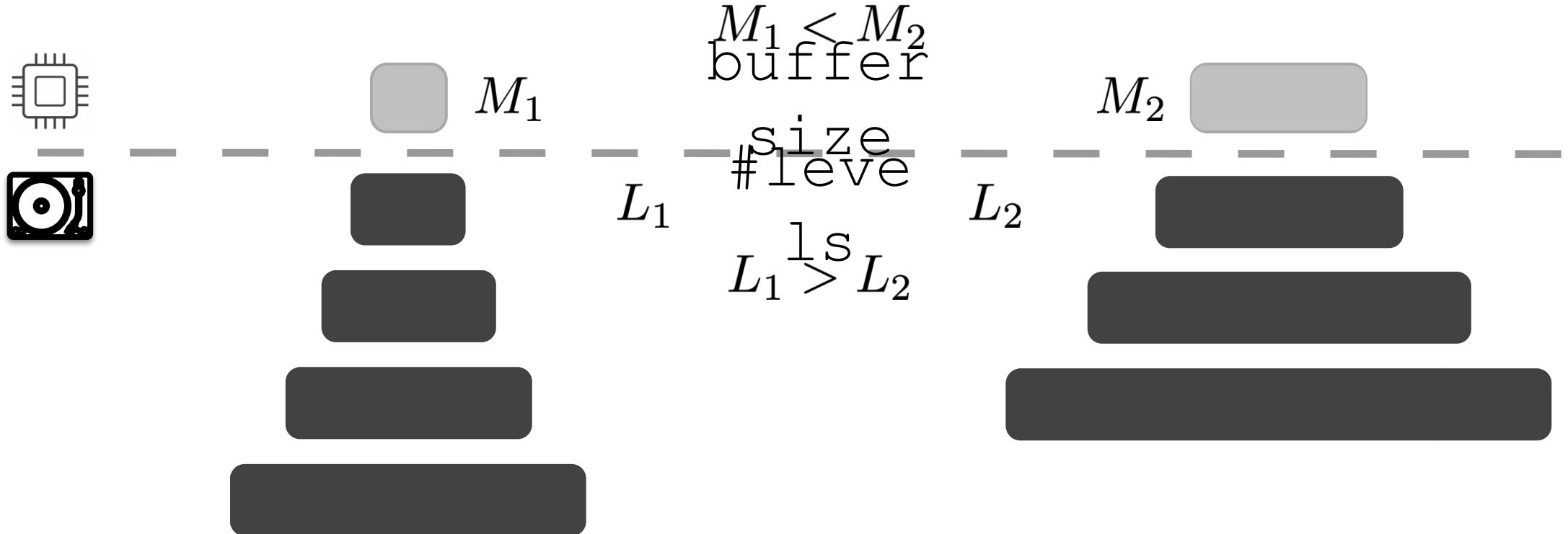
	vector	skiplist	hashmap
ingestion cost	$\mathcal{O}(1)$	$\mathcal{O}(\log(P \cdot B))$	$\mathcal{O}(1)$
space complexity	$\mathcal{O}(P \cdot B)$	$\mathcal{O}(P \cdot B)$	$\mathcal{O}(P \cdot B)$
point query cost	$\mathcal{O}(P \cdot B)$	$\mathcal{O}(\log(P \cdot B))$	$\mathcal{O}(1)$

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



most data
on storage

How does the storage layer affect ingestion?

L : #levels

T : size ratio

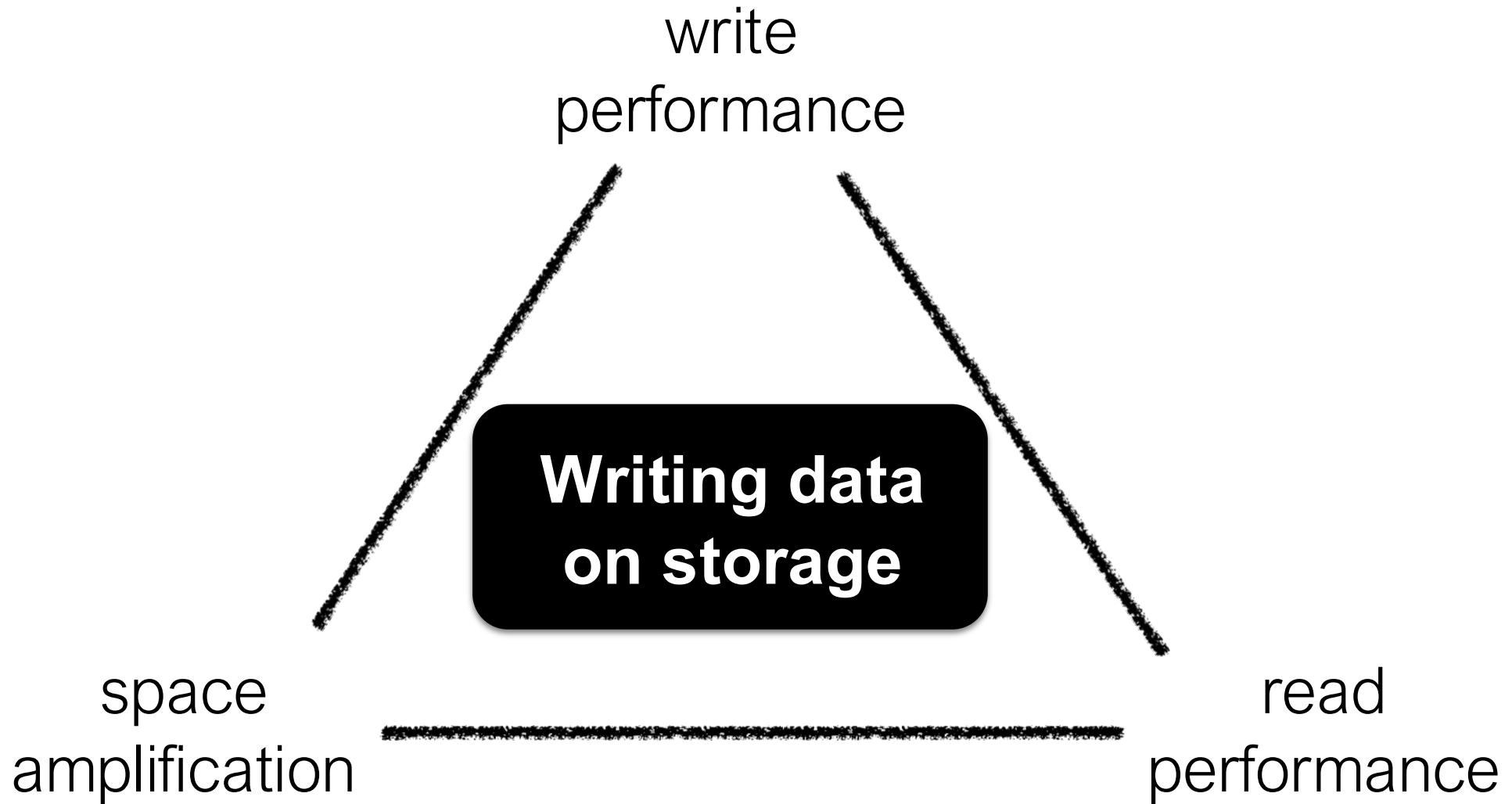


most data
on storage
if $T = 10$ & $L = 4$

99.9% on storage

How does the storage layer affect ingestion?

Storage Optimizations



Data Layout

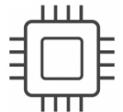
Classical LSM design:

level in
g
[eager merging]



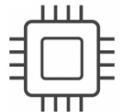
Data Layout

leveling [eager]



Data Layout

leveling [eager]



Data Layout

leveling [eager]



Data Layout

leveling [eager]



Data Layout

leveling [eager]



Data Layout

leveling [eager]



Data Layout

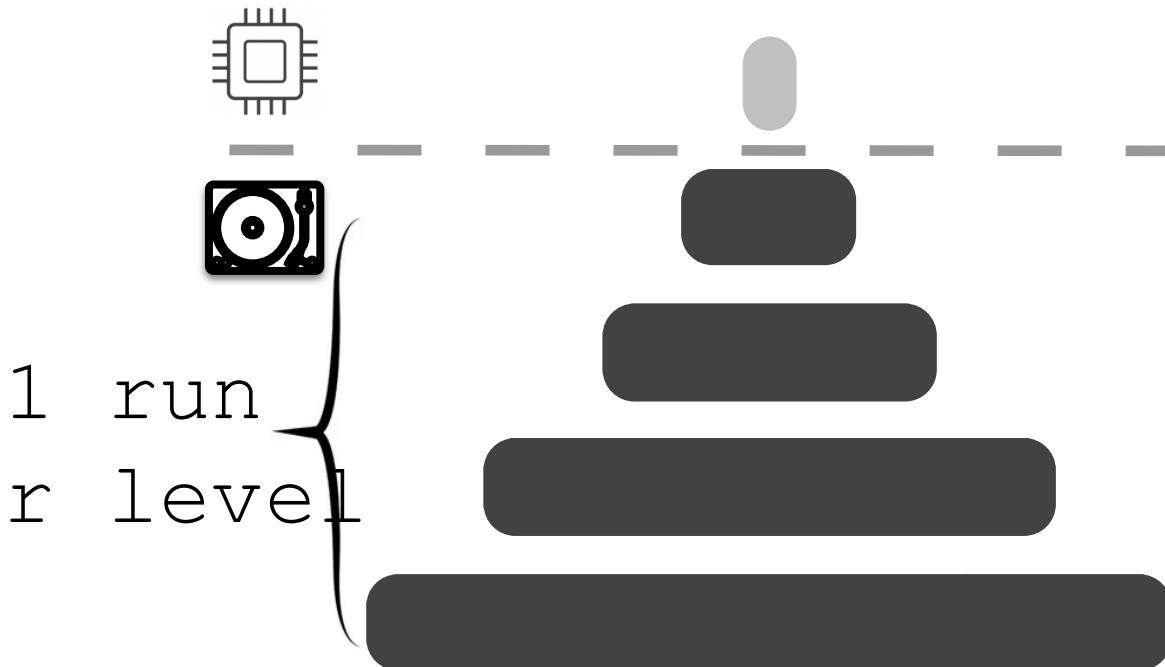
leveling [eager]



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

Data Layout

leveling [eager]



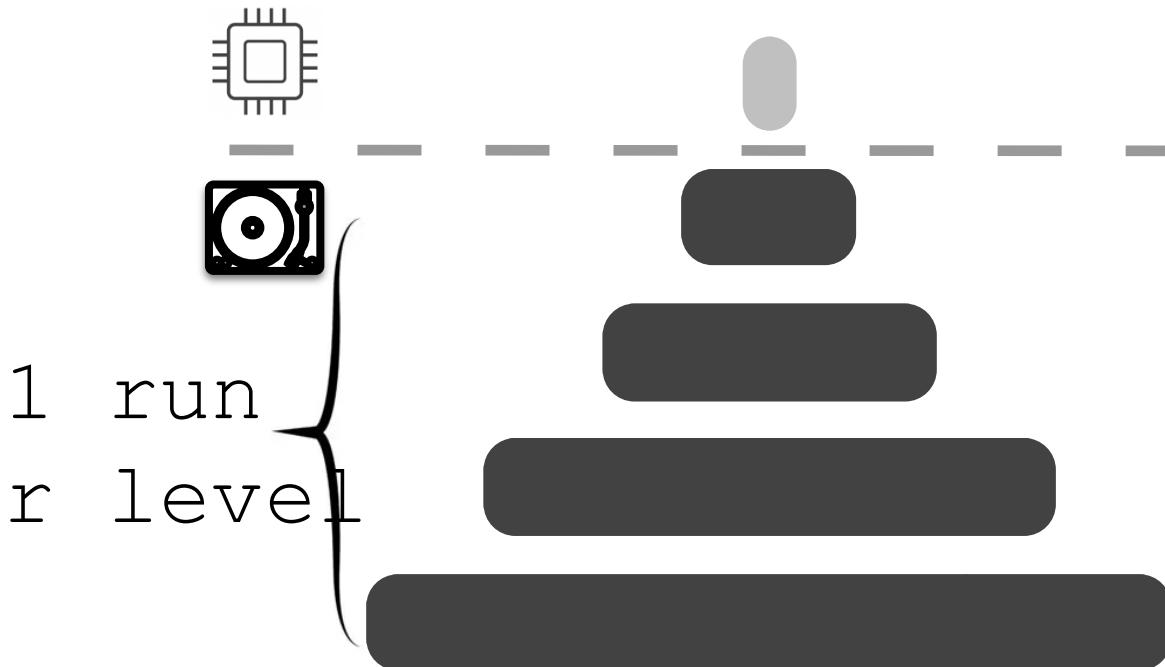
tiering [lazy]



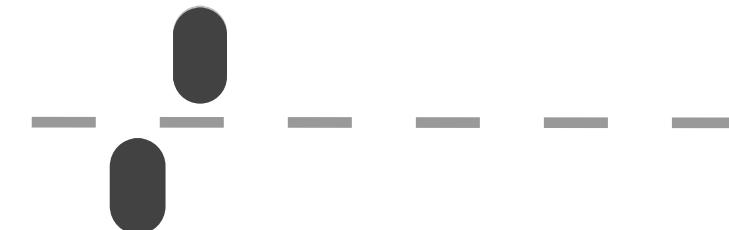
- good read performance
- good space amplification
- high write amplification

Data Layout

leveling [eager]



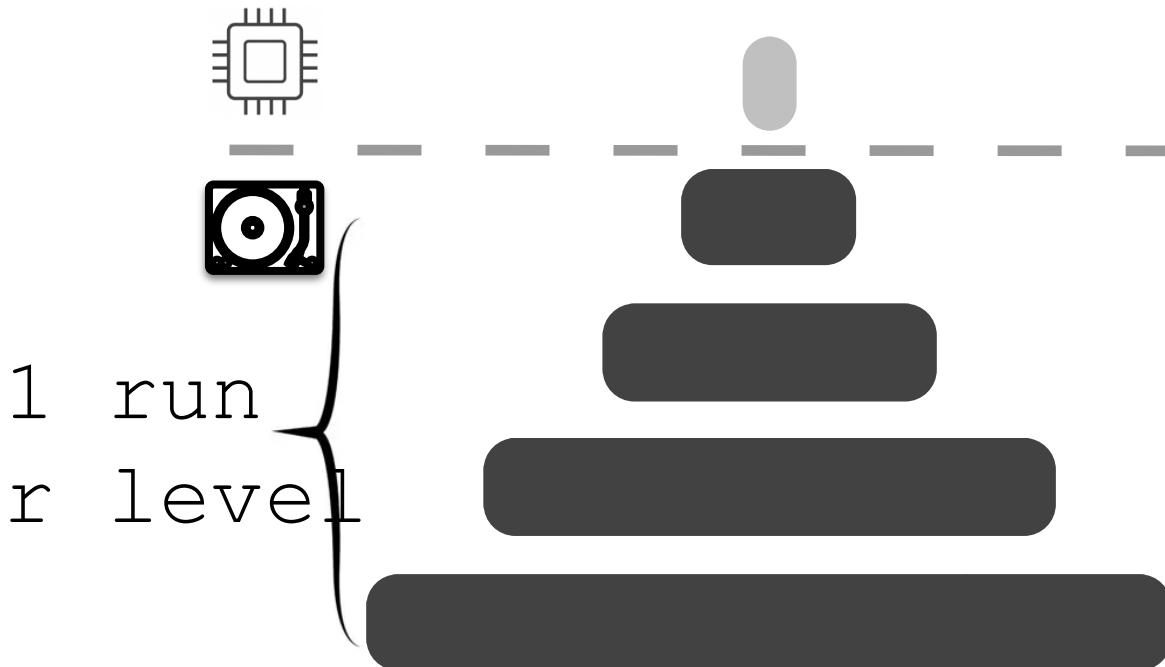
tiering [lazy]



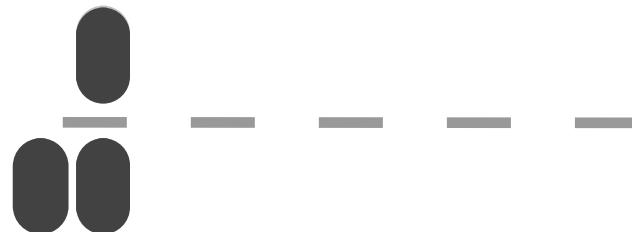
- good read performance
- good space amplification
- high write amplification

Data Layout

leveling [eager]



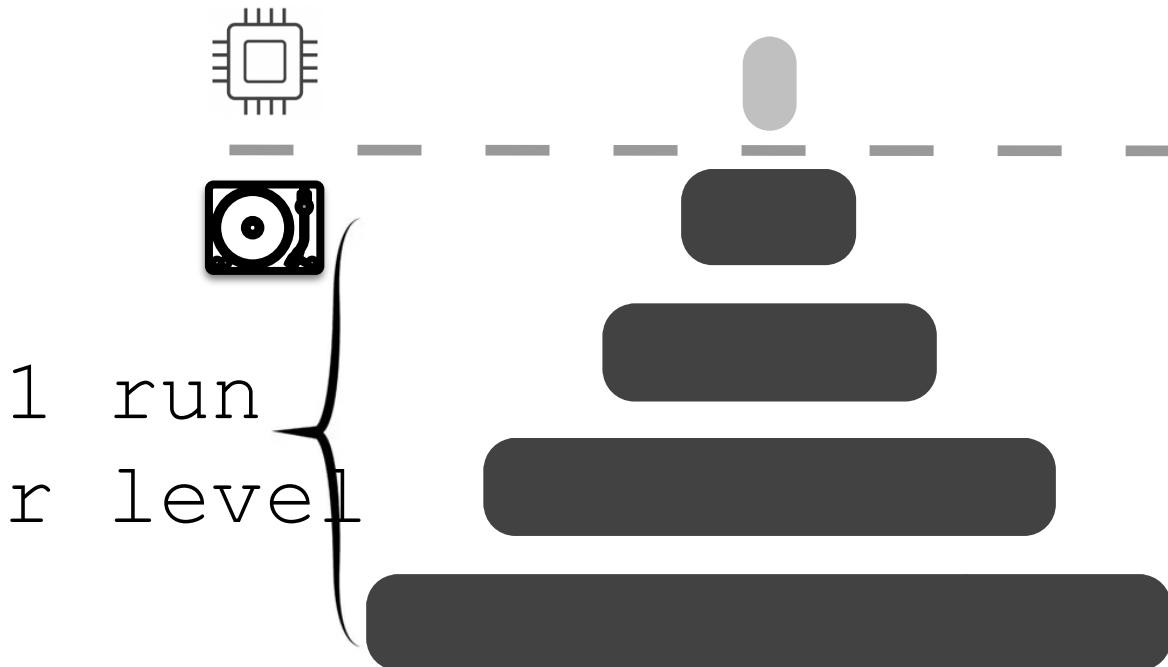
tiering [lazy]



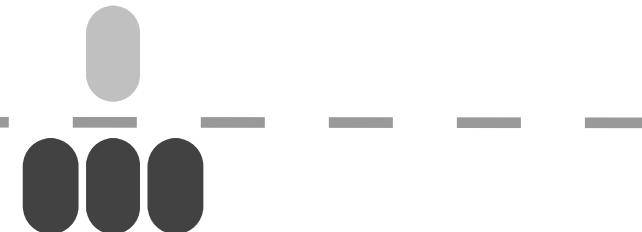
- good read performance
- good space amplification
- high write amplification

Data Layout

leveling [eager]



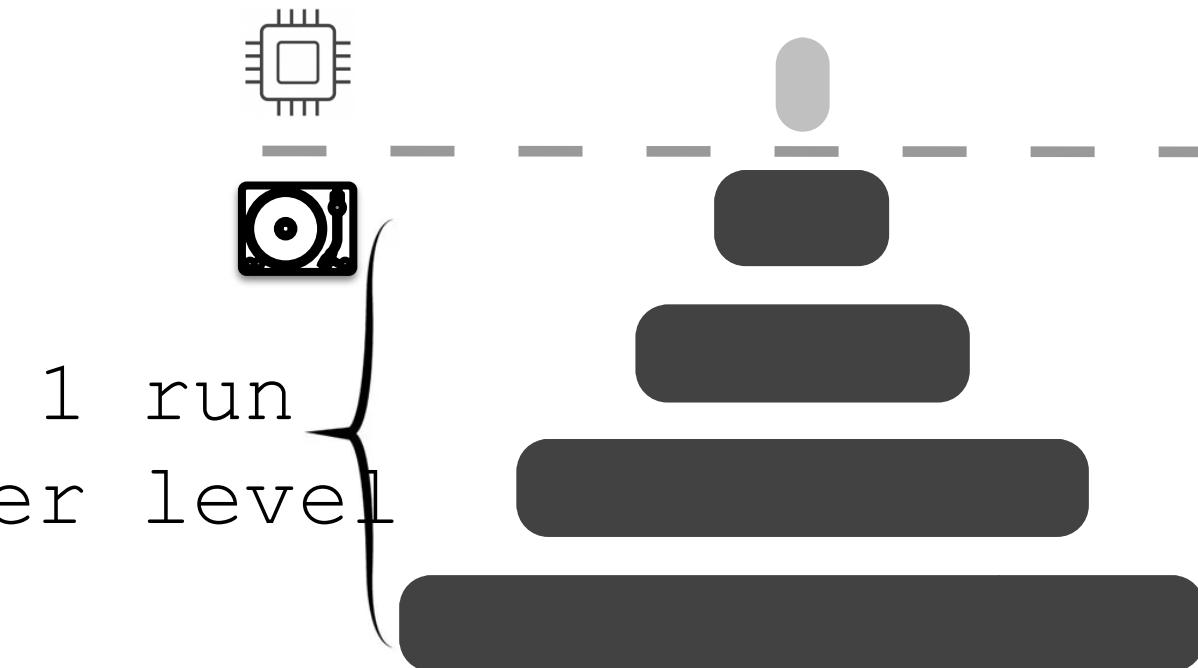
tiering [lazy]



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

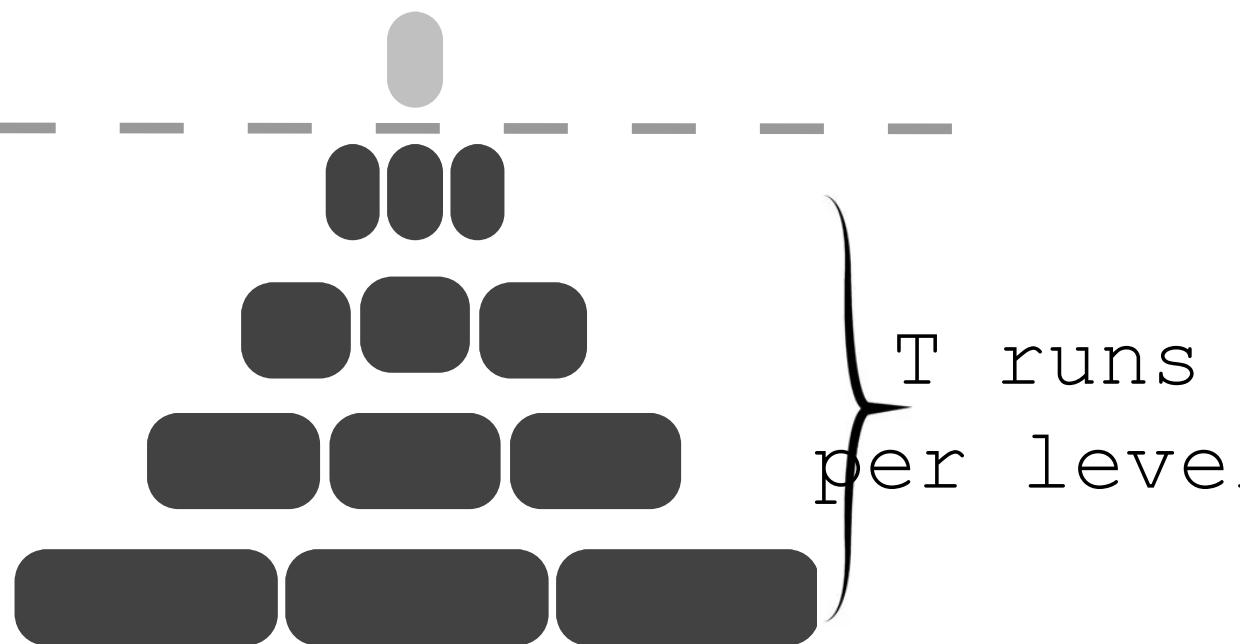
Data Layout

leveling [eager]



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

tiering [lazy]

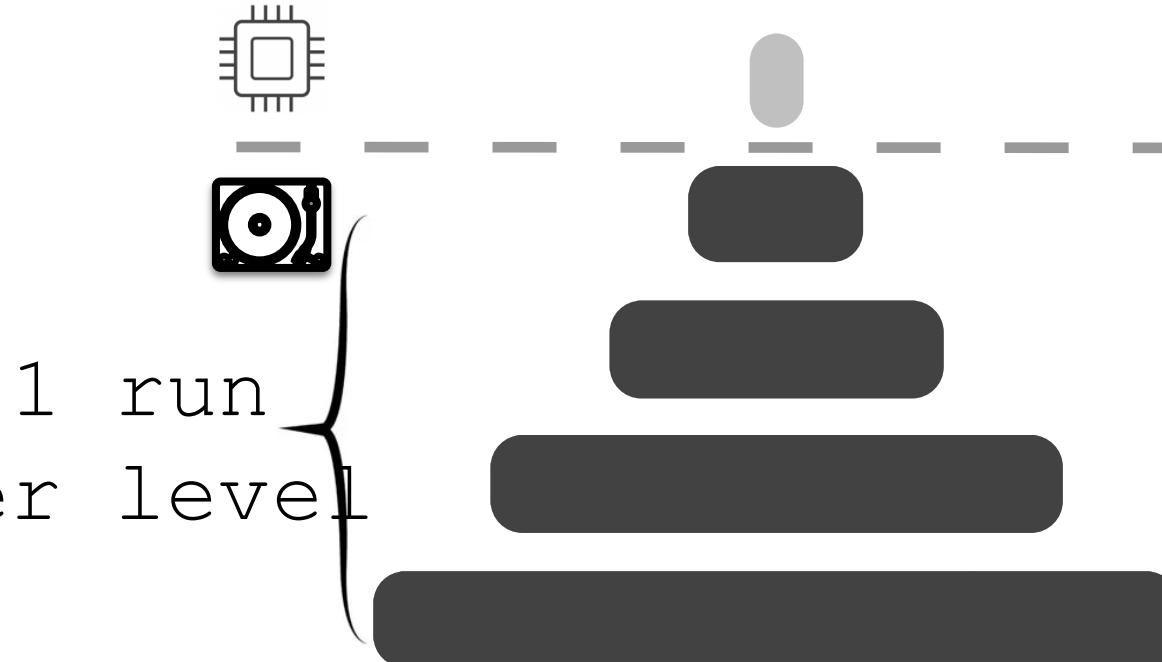


- poor read performance
- poor space amplification
- good ingestion performance

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

Data Layout

leveling [eager]



Read cost:

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

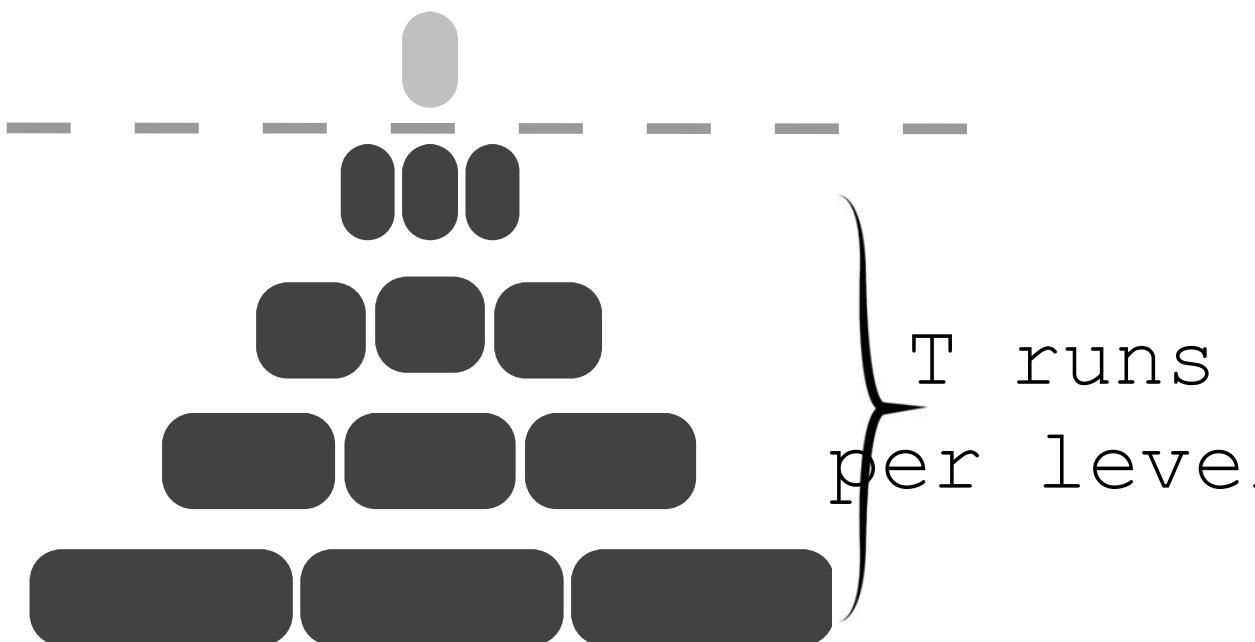
Write cost:

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

SA:

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

tiering [lazy]



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

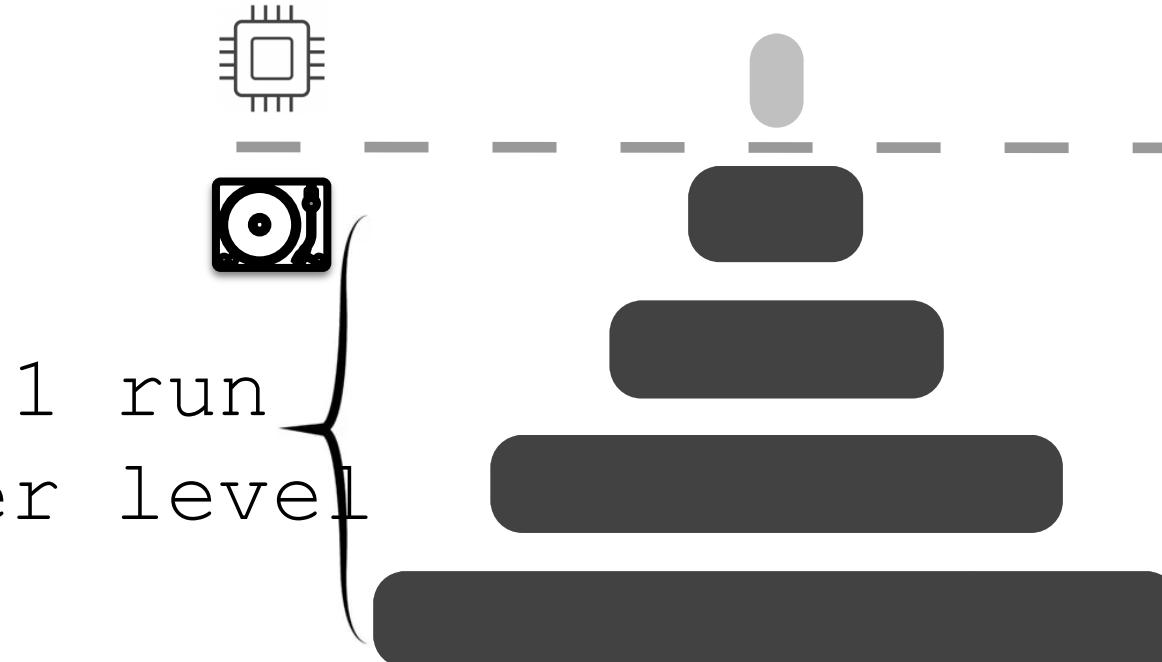
$$\mathcal{O}(L/B)$$

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

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

Data Layout

leveling [eager]



Read cost:

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

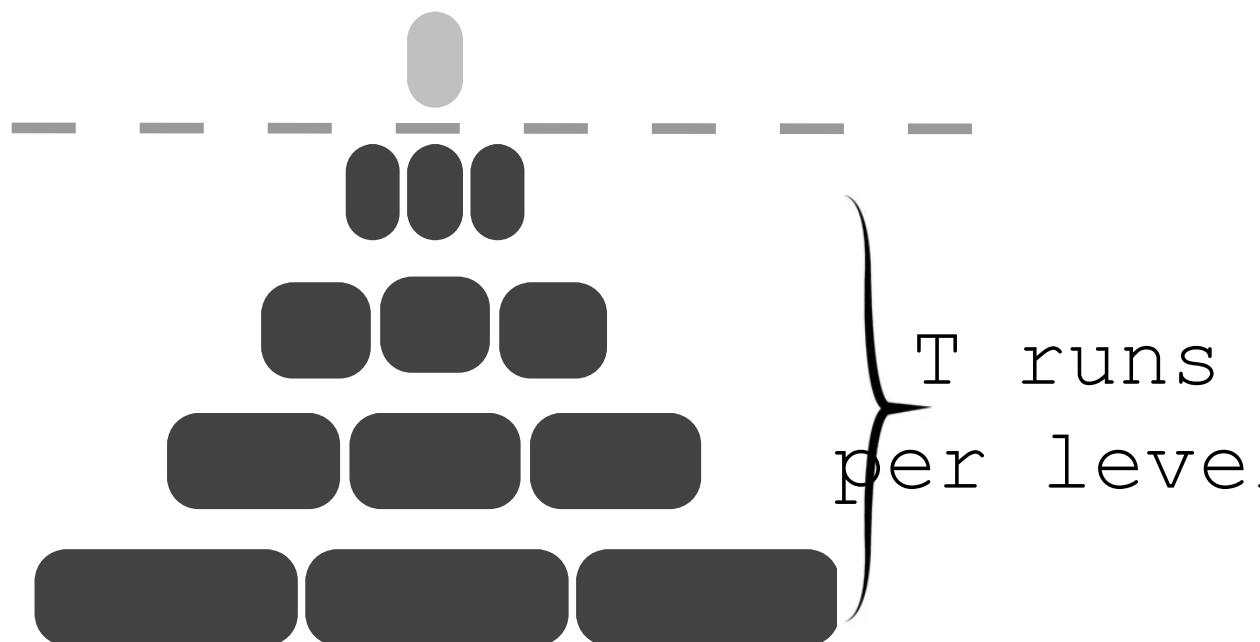
Write cost:

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

SA:

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

tiering [lazy]

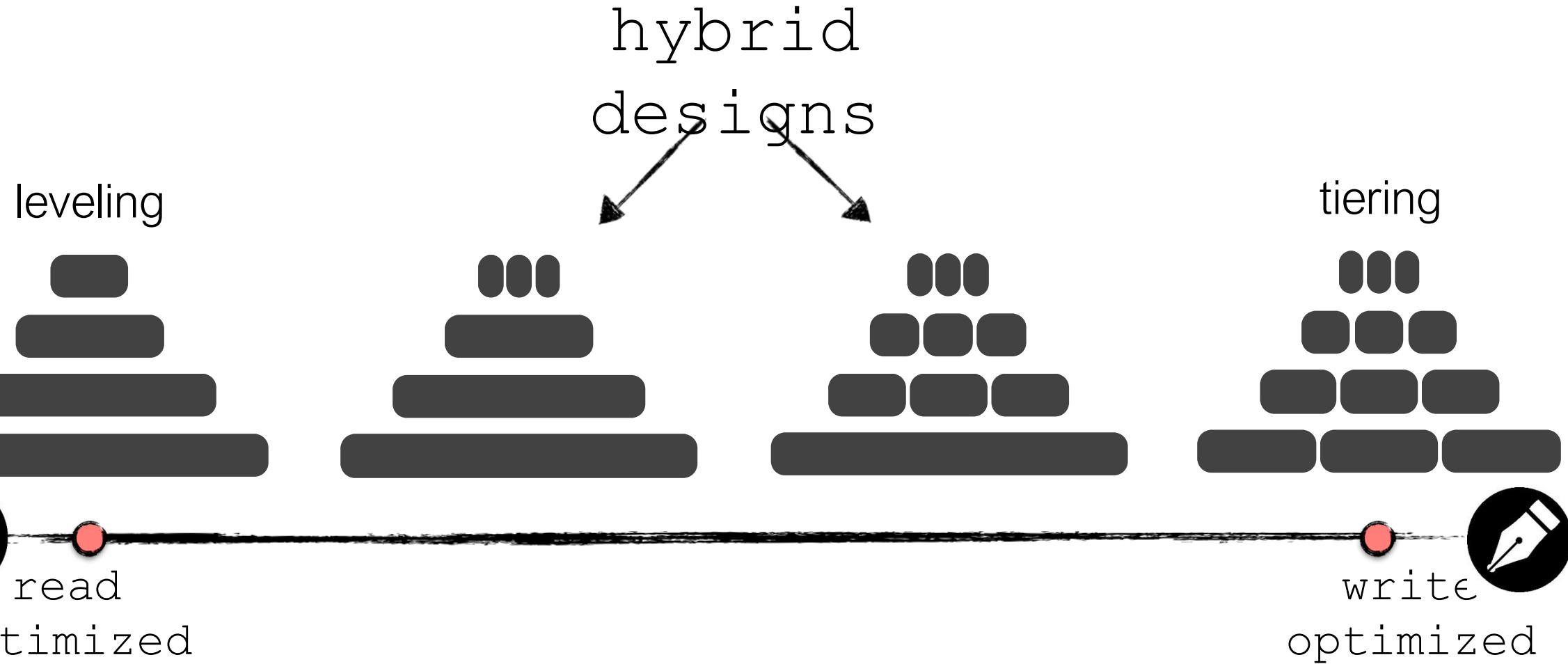


$$\mathcal{O}(T \cdot L \cdot \phi)$$

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

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

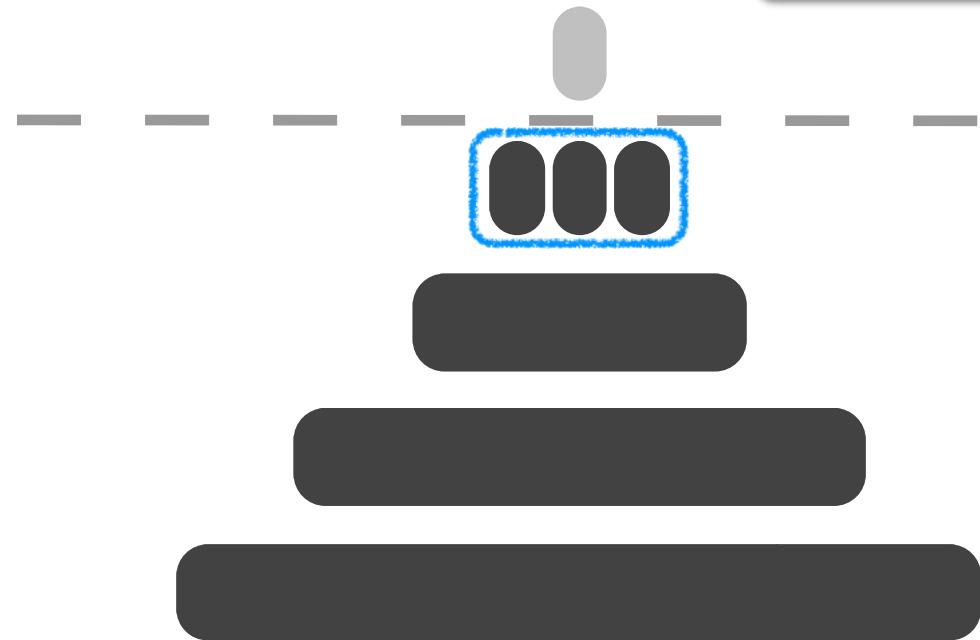
Data Layout



Data Layout

1-leveling

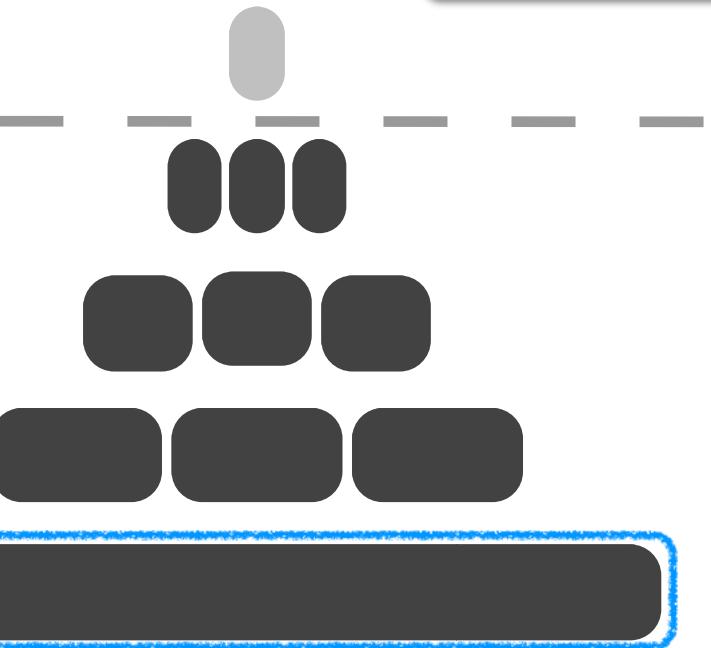
RocksDB20



- fewer write stalls
- increased block cache hits
compared to leveling!

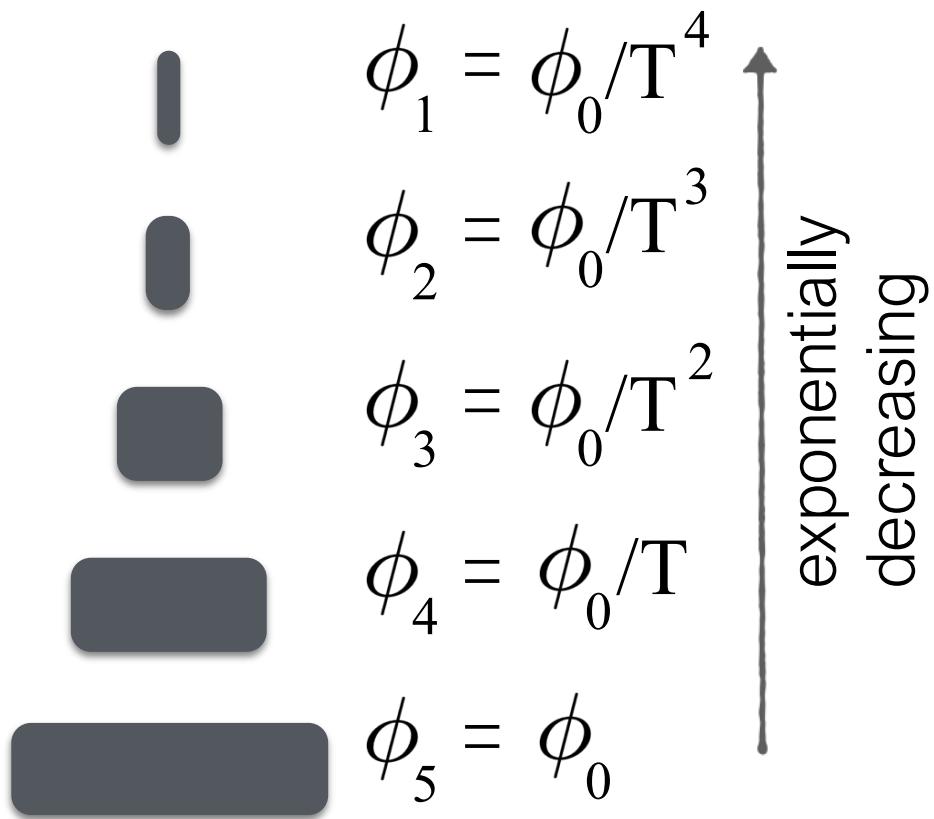
L-leveling

DayanSIGMOD18



- low write amplification
- better read performance
compared to tiering!

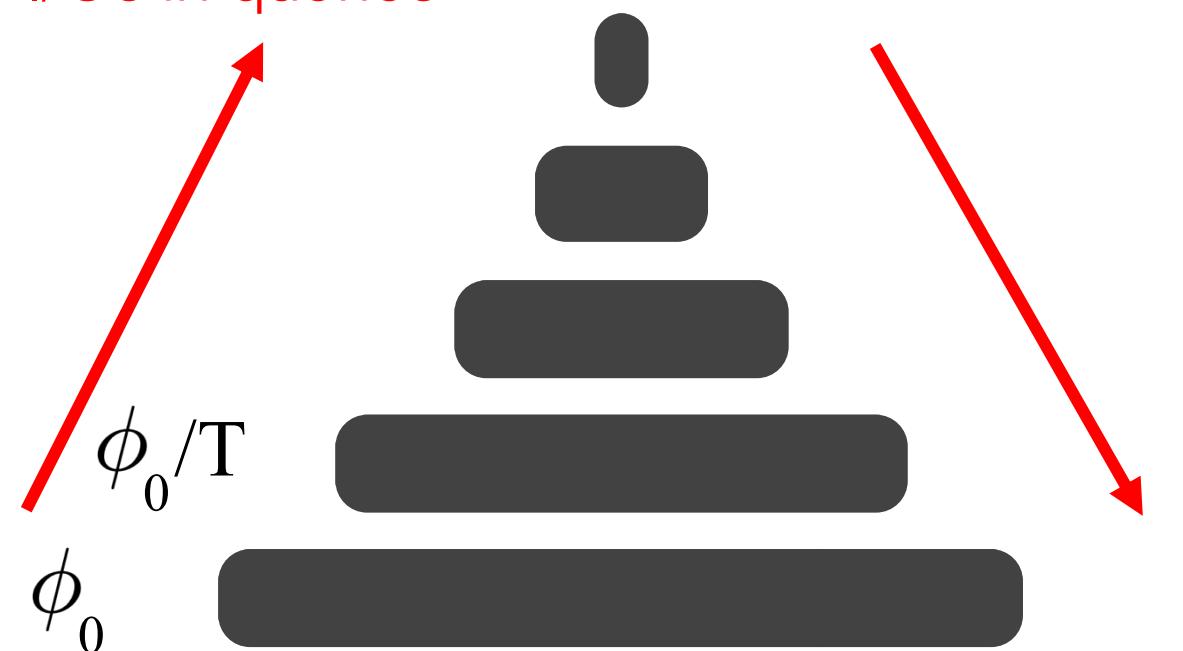
Bloom filters



LSM-tree

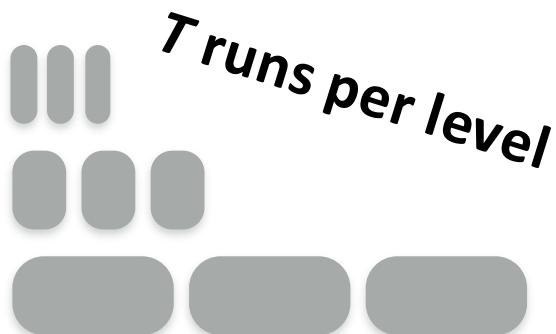
exponentially fewer I/Os in queries

however, every run participates in T merges

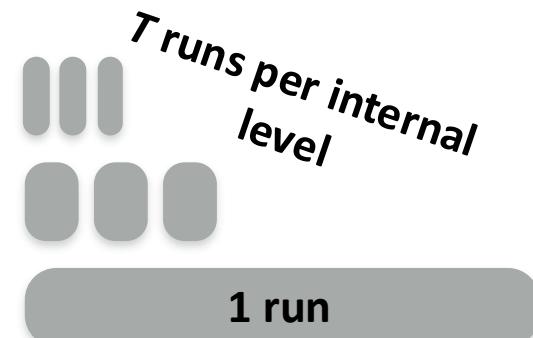


Insight: merging in top levels can become lazy (tiering)!

Tiering write-optimized



Lazy Leveling hybrid



Leveling read-optimized



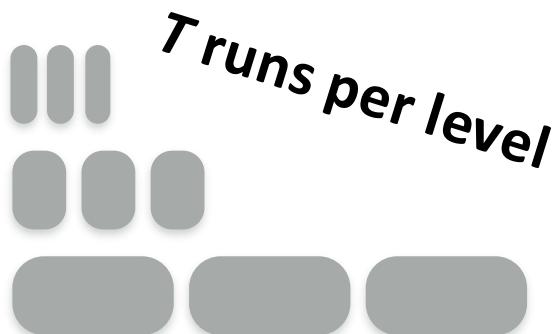
lookup: $O(T \cdot \log_T(N) \cdot e^{-M/N})$

runs per level levels false positive rate

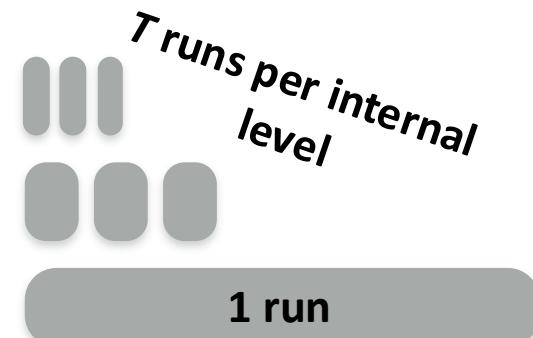
$O(\log_T(N) \cdot e^{-M/N})$

levels false positive rate

Tiering write-optimized



Lazy Leveling hybrid



Leveling read-optimized



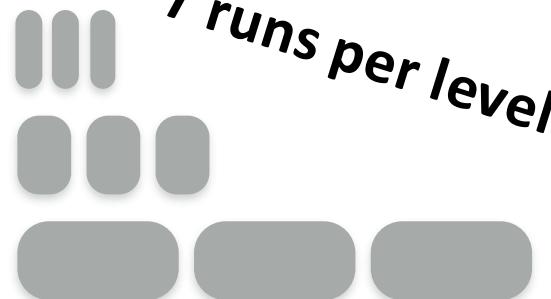
lookup: $O(T \cdot \log_T(N) \cdot e^{-M/N})$

runs per level levels false positive rate

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

levels false positive rate

Tiering write-optimized



DayanSIGMOD17

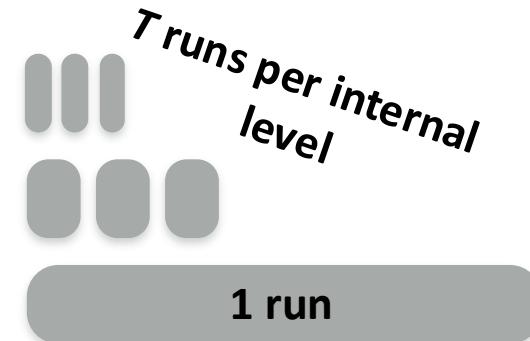
lookup:

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

runs per level false positive rate

Arrows point from the words "runs per level" and "false positive rate" to the variables T and N in the formula respectively.

Lazy Leveling hybrid



Leveling read-optimized

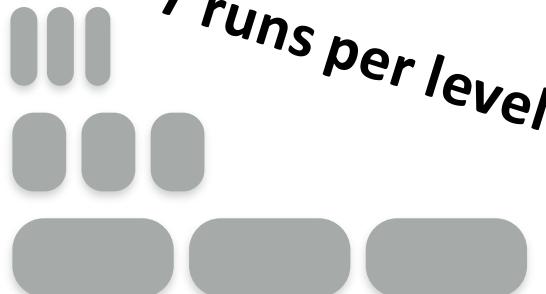


$$O(e^{-M/N})$$

false positive rate

An arrow points from the words "false positive rate" to the variable N in the formula.

Tiering write-optimized



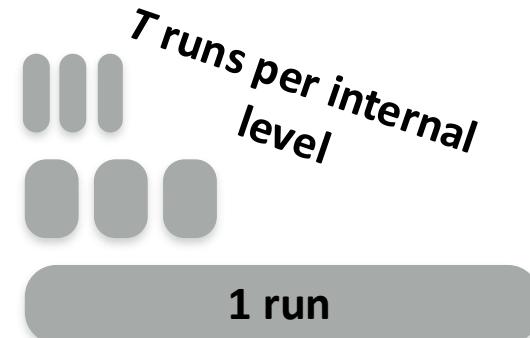
DayanSIGMOD18

lookup:

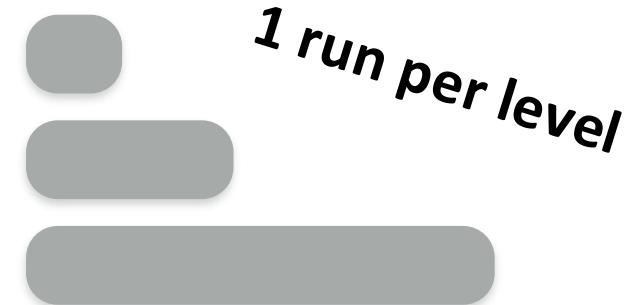
$$O(T \cdot e^{-M/N})$$

runs per level false positive rate

Lazy Leveling hybrid



Leveling read-optimized



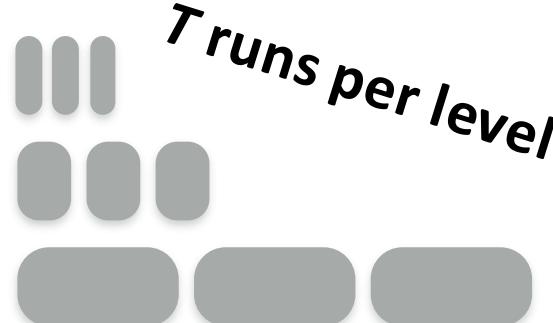
$$O(e^{-M/N})$$

there is a slightly higher constant factor than leveling

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

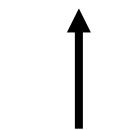
false positive rate

Tiering write-optimized



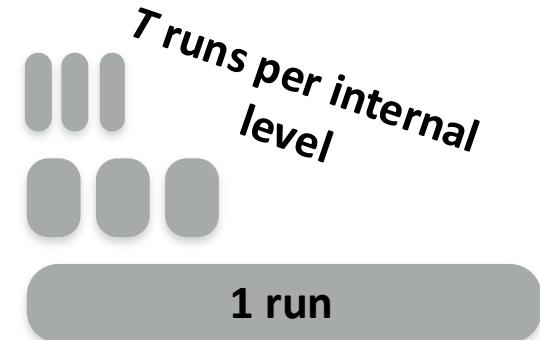
lookup: $O(T \cdot e^{-M/N})$

update cost: $O(\log_T(N))$



levels

Lazy Leveling hybrid



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

update cost: $O(T + \log_T(N))$



T merges for the last level 1 merge for all internal levels

Leveling read-optimized



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

update cost: $O(T \cdot \log_T(N))$

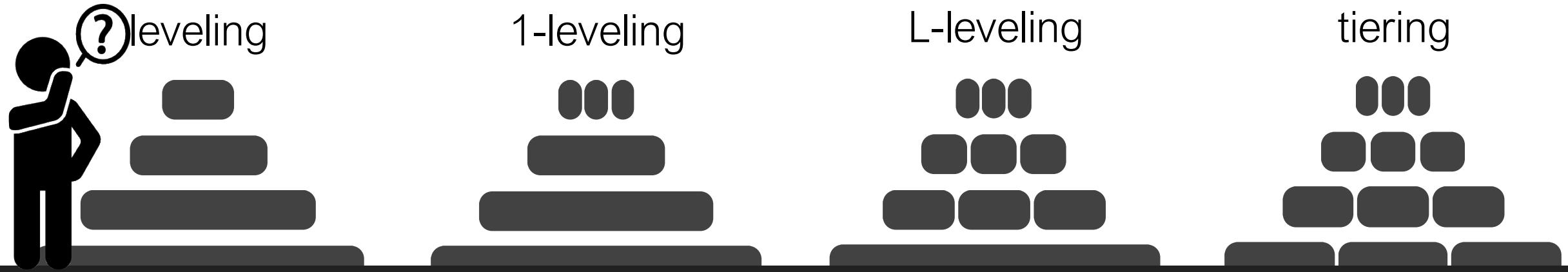


merges per level

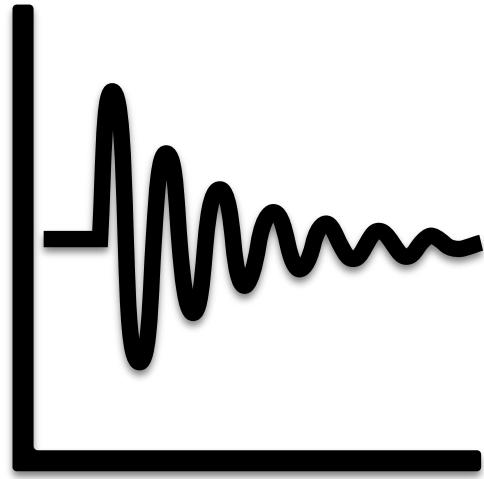
levels

121

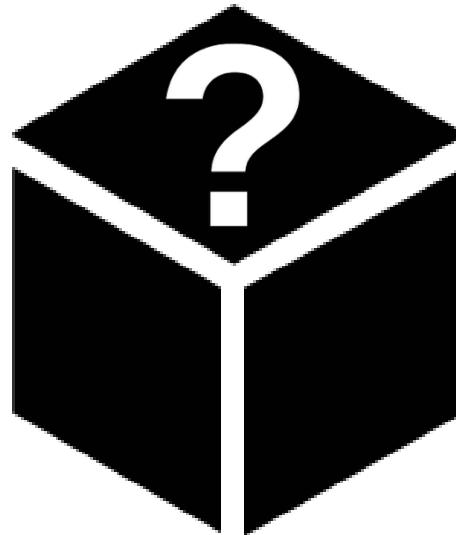
Data Layout



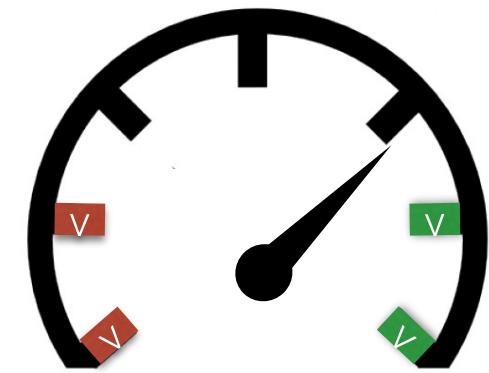
So, how do we reason about the data layout?



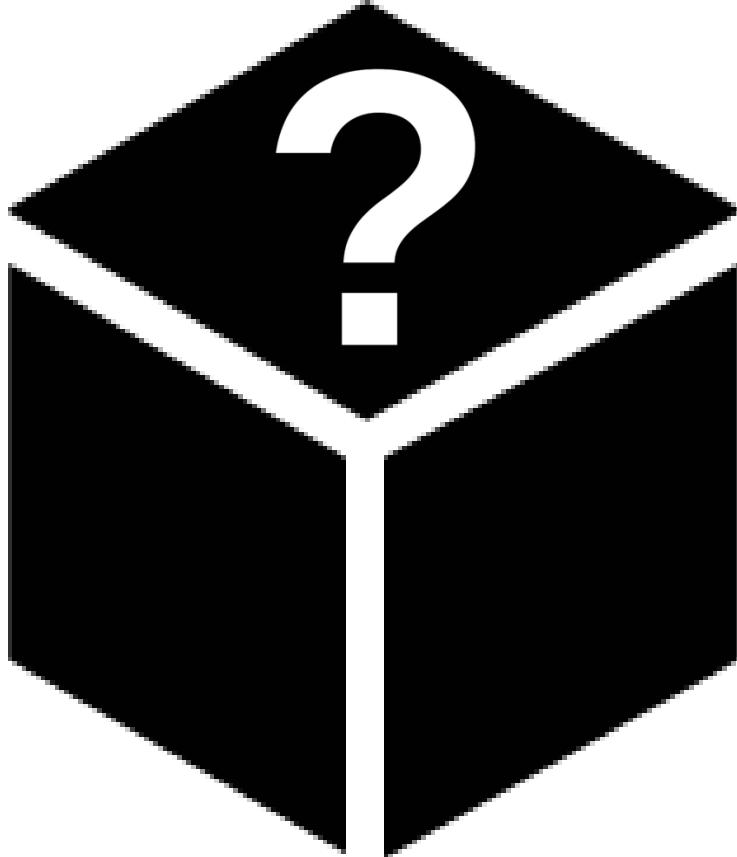
workload



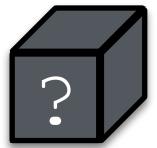
data layout



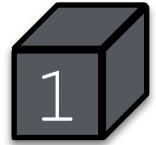
performance



Compaction
black box



SarkarVLDB21



How to organize the data on device?



How much data to move at-a-time?



Which block of data to be moved?



When to re-organize the data layout?



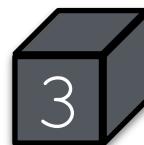
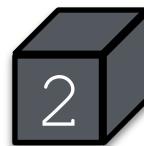
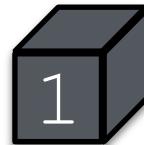
SarkarVLDB21

Data
Layout.

+
n
Granulari
-

Data
Movement
Policy

Compaction
Trigger



How to organize the data on device? 

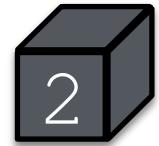
How much data to move at-a-time?

Which block of data to be moved?

When to re-organize the data layout?



SarkarVLDB21



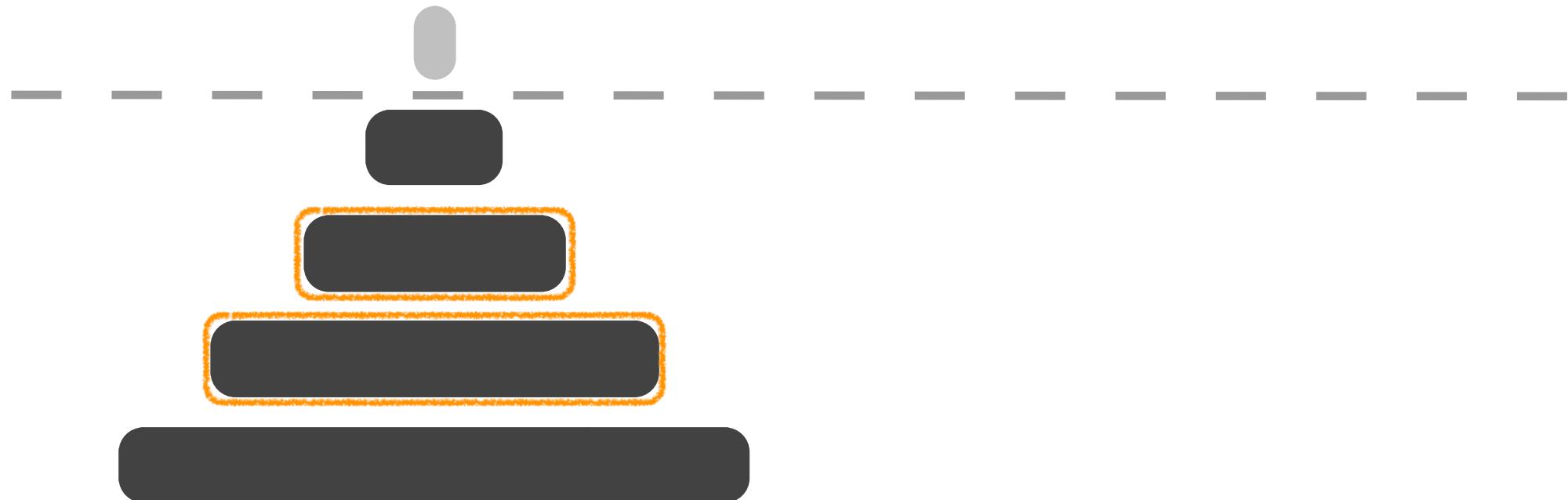
Compaction **Granularity**

data moved per compaction



Compaction **Granularity**

data moved per compaction



consecutiv

e levels

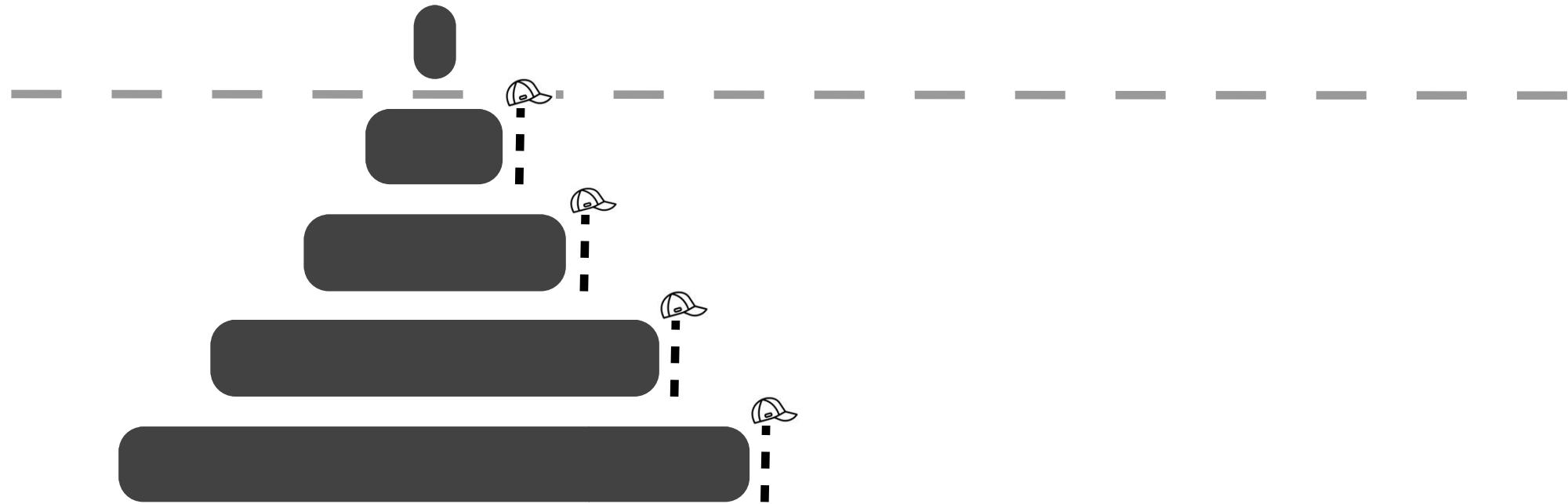
AsterixDB





Compaction **Granularity**

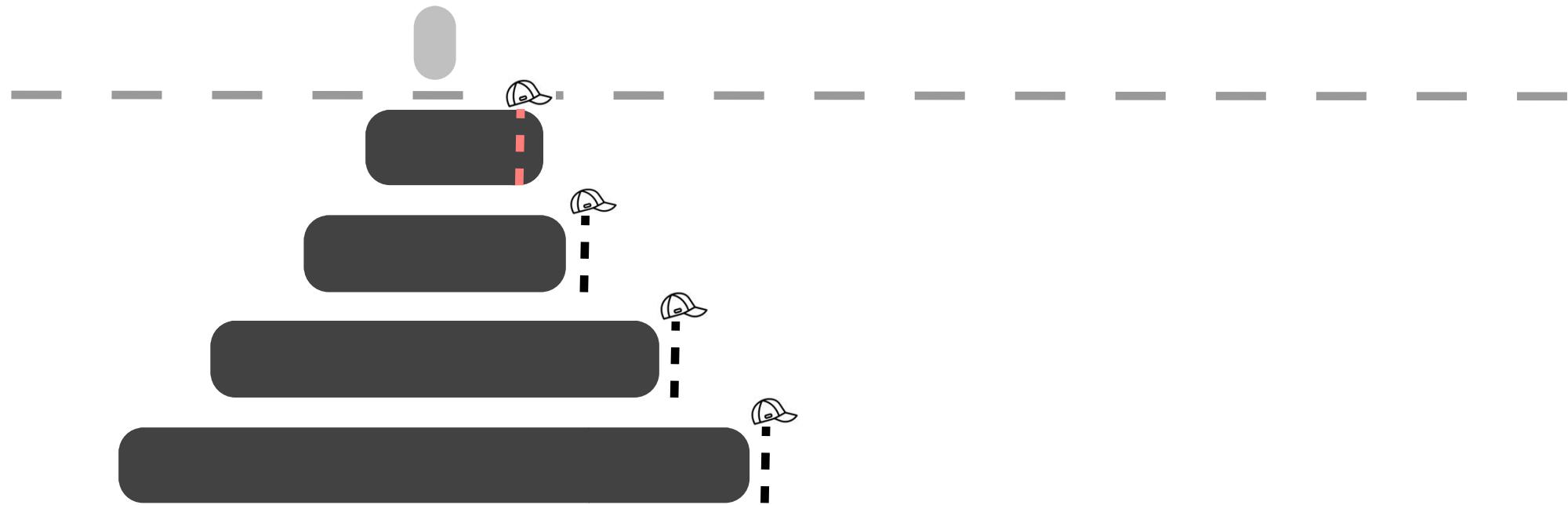
data moved per compaction





Compaction **Granularity**

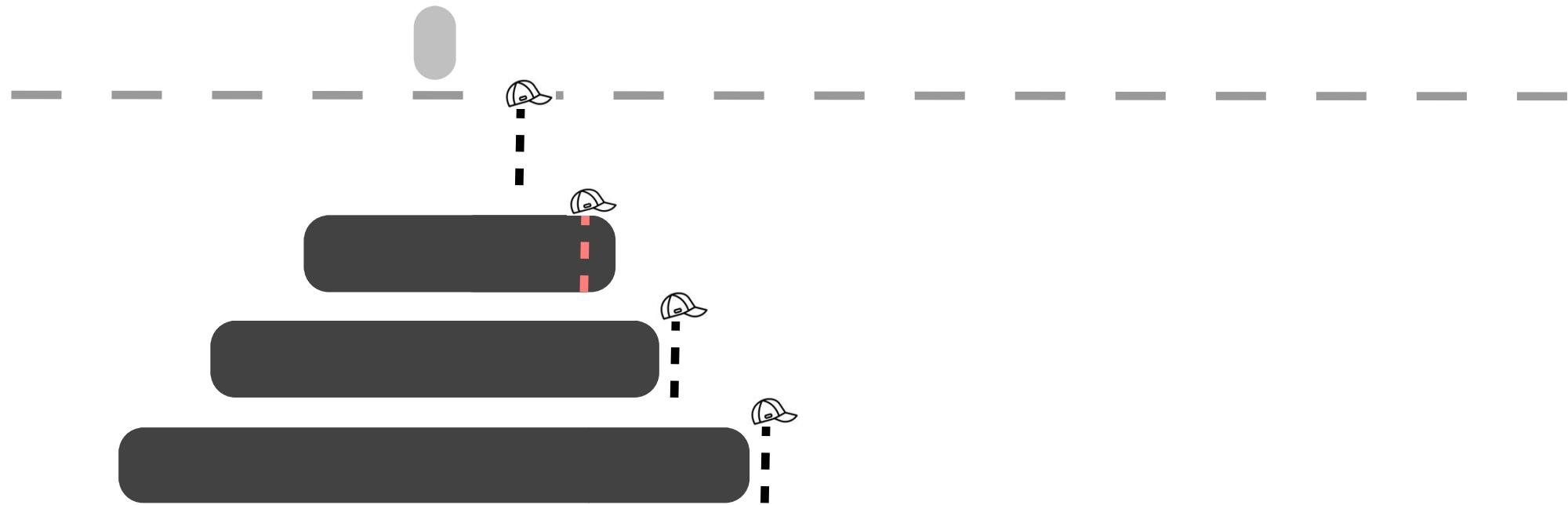
data moved per compaction





Compaction **Granularity**

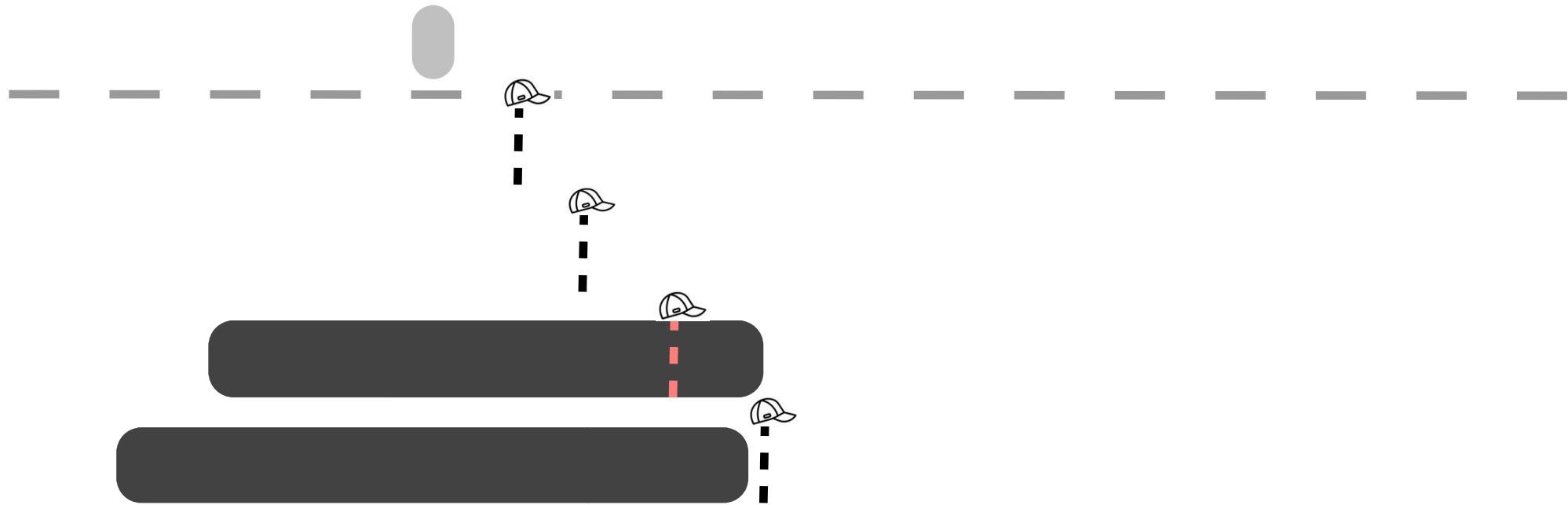
data moved per compaction





Compaction **Granularity**

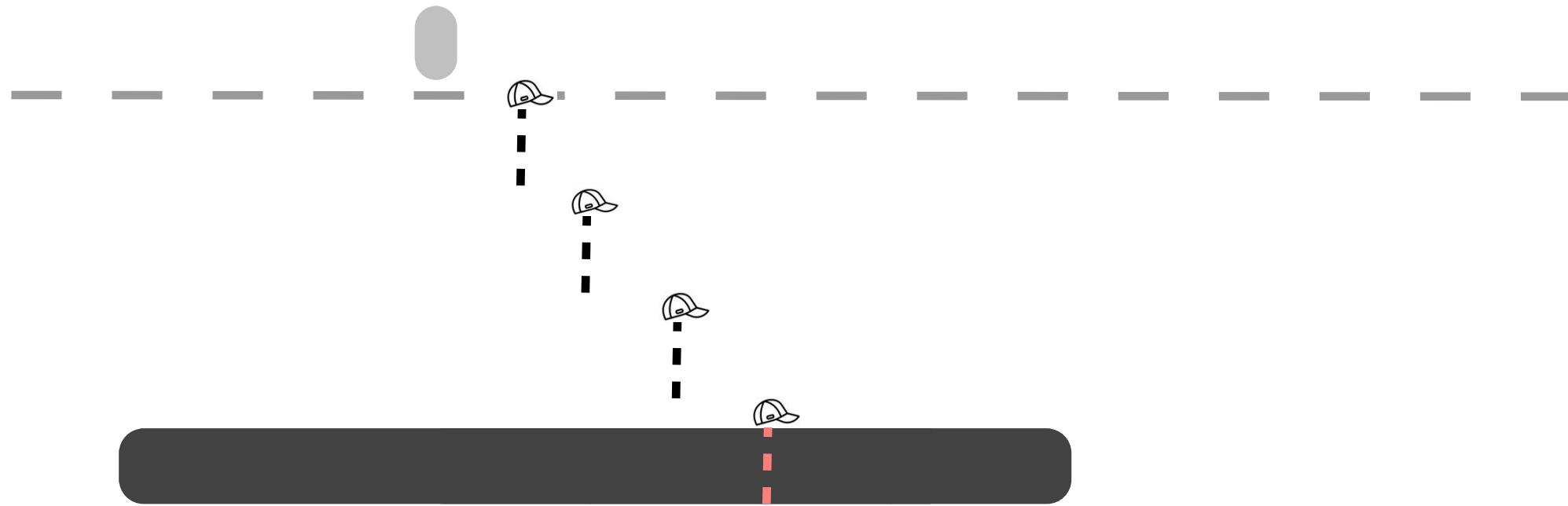
data moved per compaction





Compaction **Granularity**

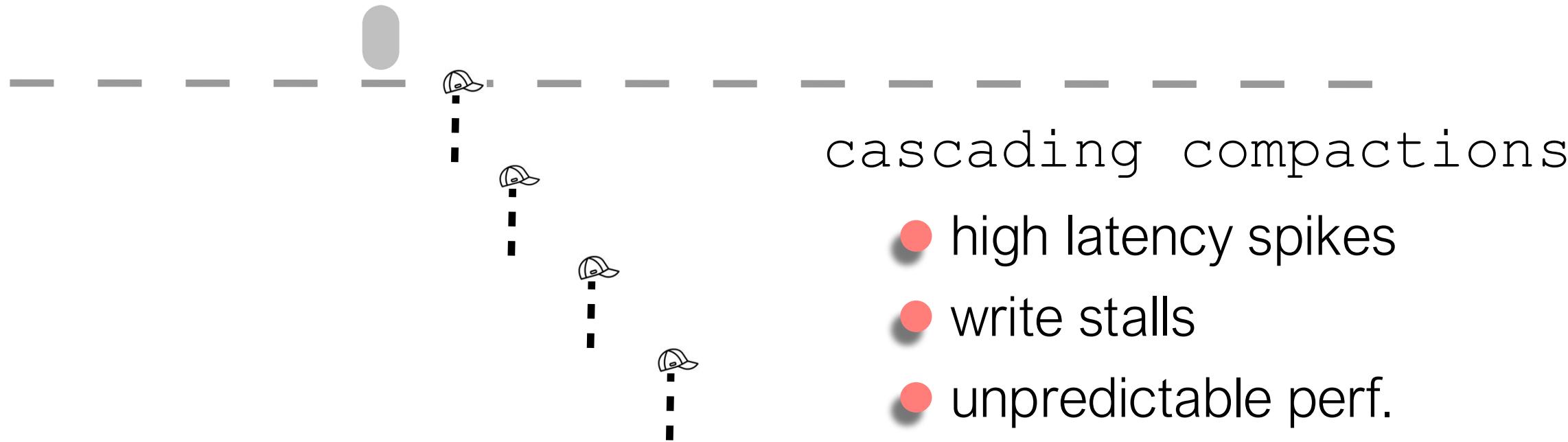
data moved per compaction





Compaction **Granularity**

data moved per compaction





Compaction **Granularity**

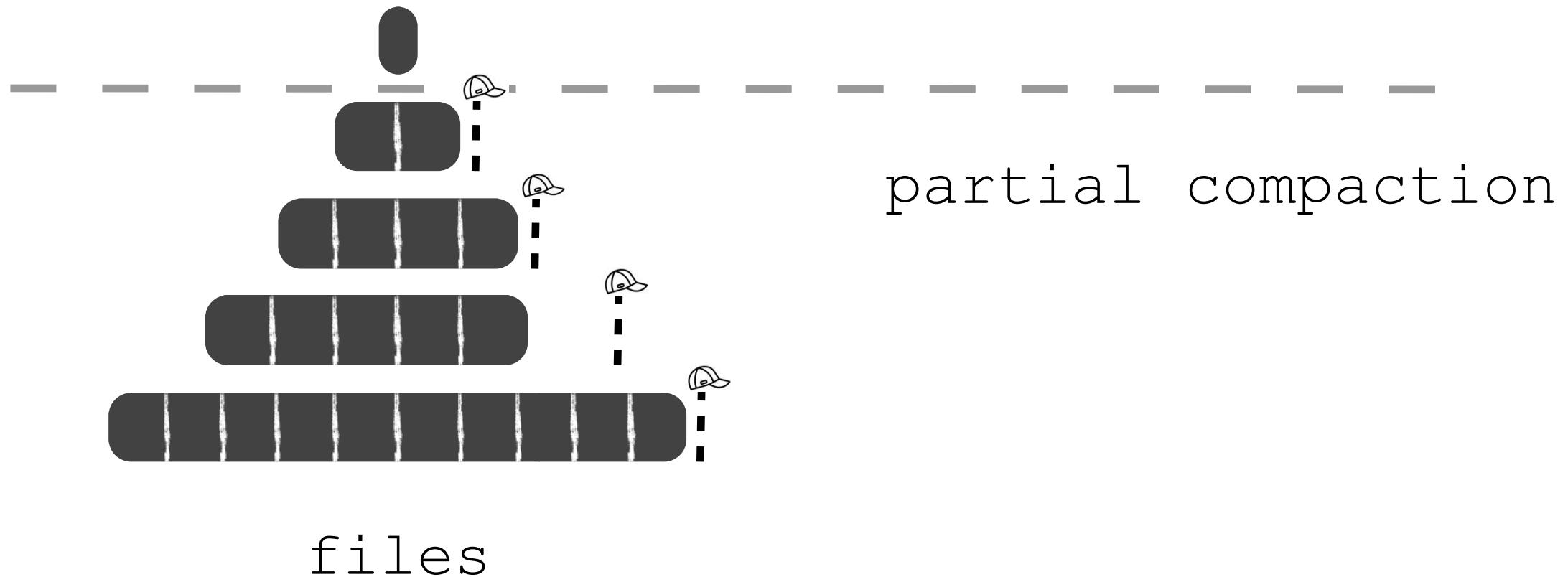
data moved per compaction

partial compaction
granularity: files



Compaction **Granularity**

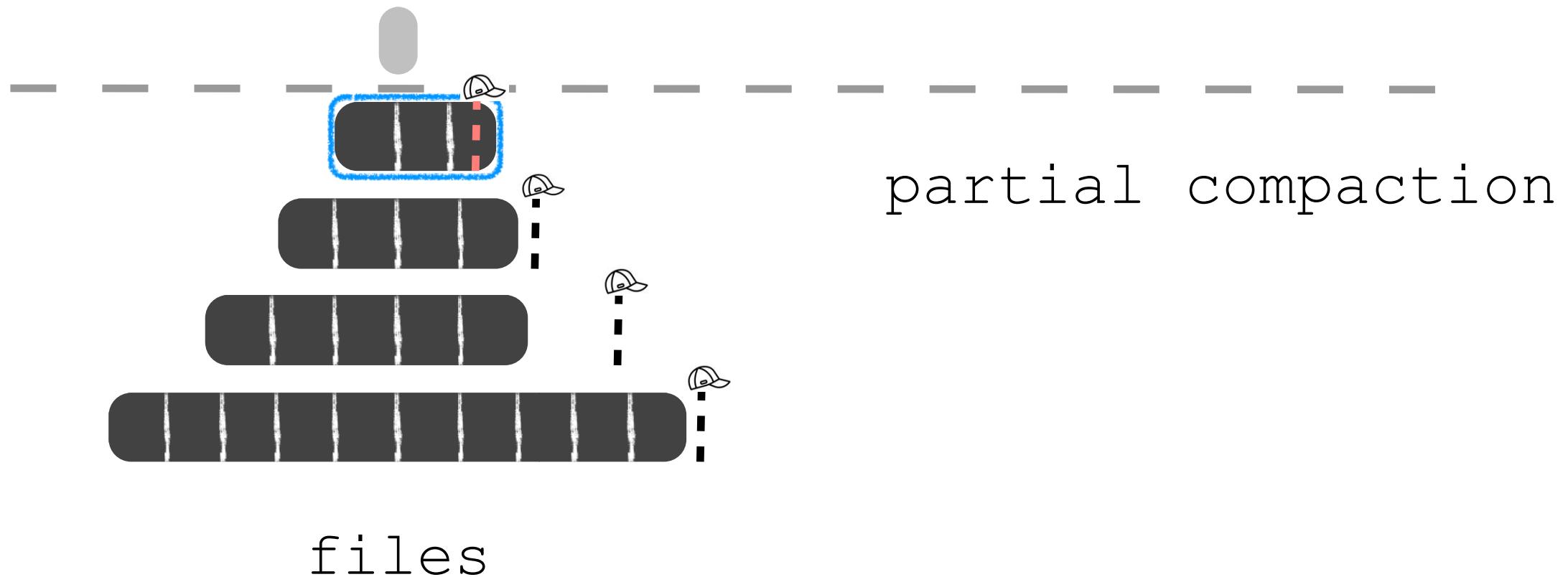
data moved per compaction





Compaction **Granularity**

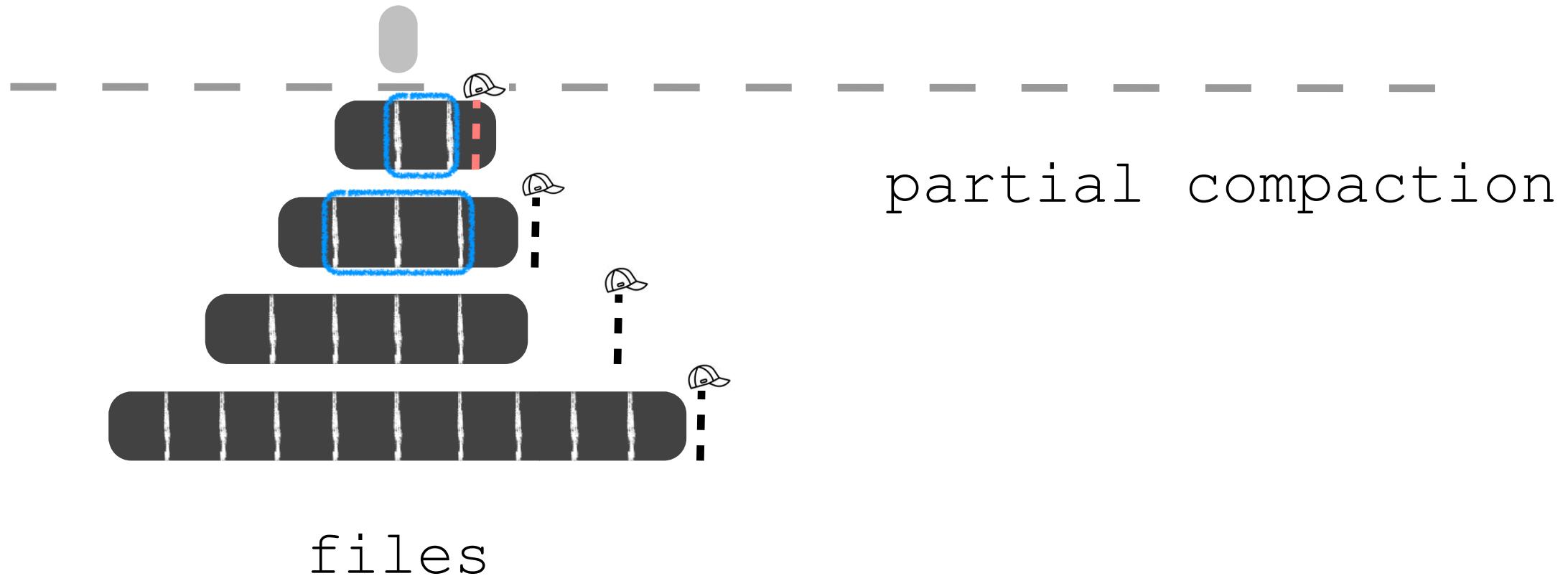
data moved per compaction





Compaction **Granularity**

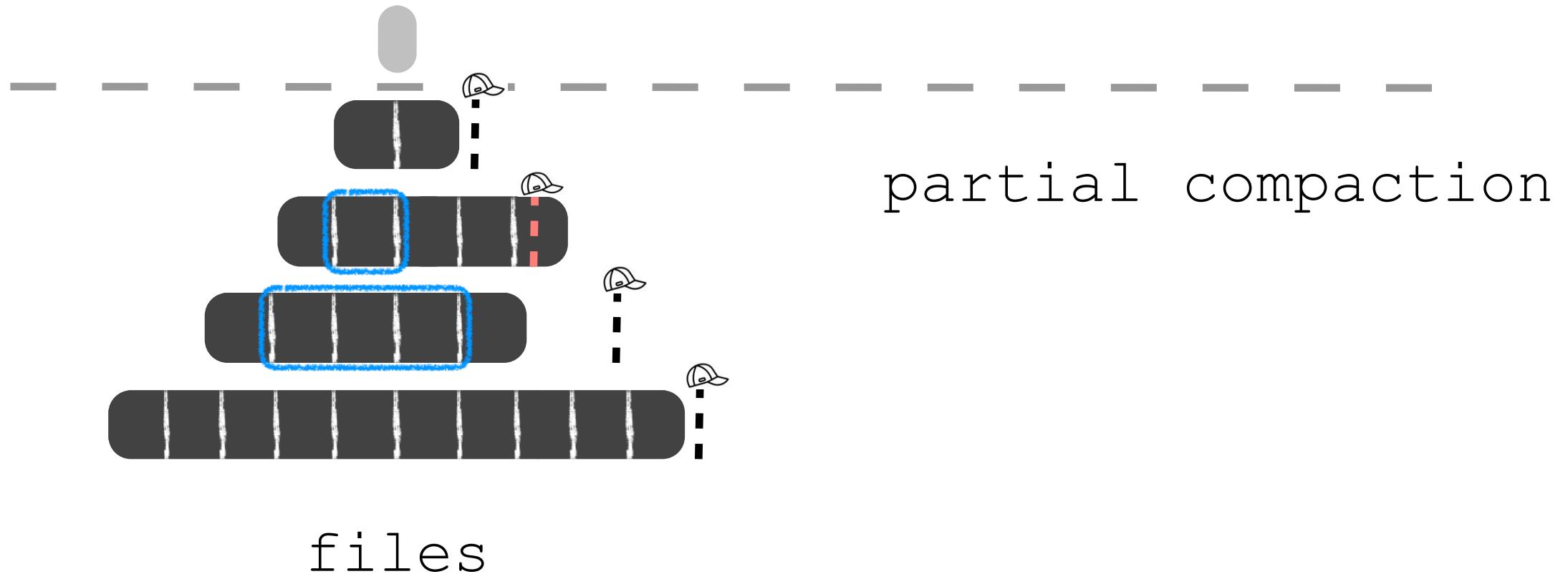
data moved per compaction





Compaction **Granularity**

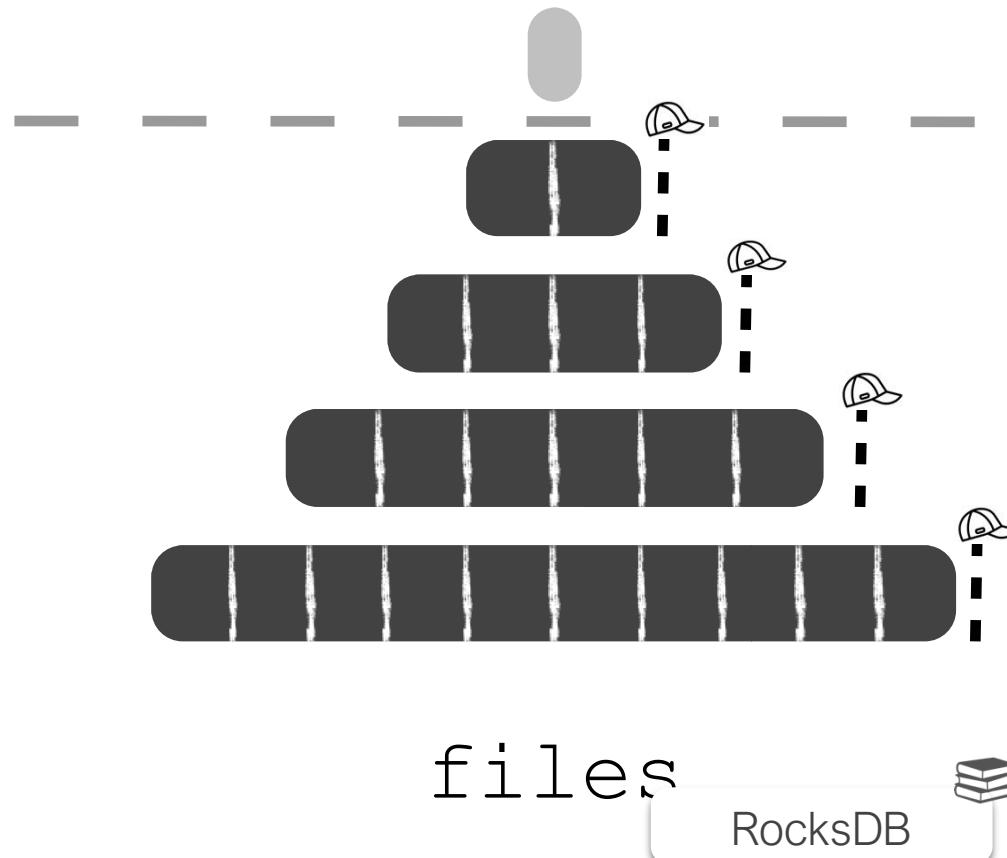
data moved per compaction





Compaction **Granularity**

data moved per compaction



partial compaction

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



Compaction **Granularity**

data moved per compaction

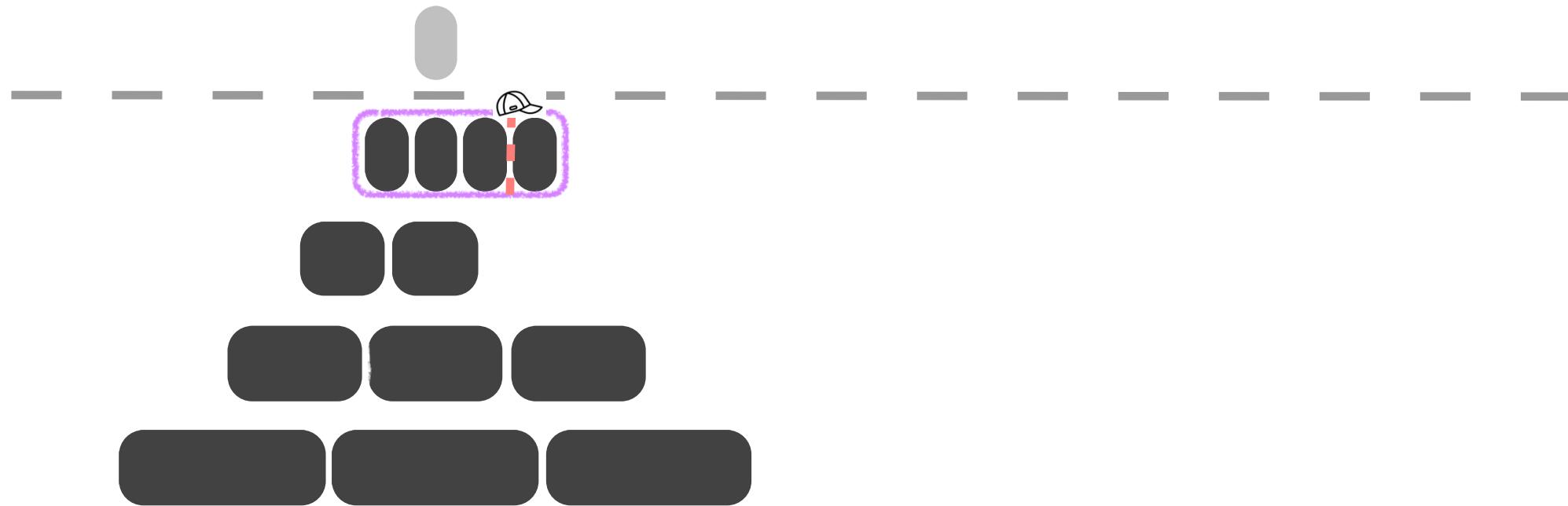


sorted runs in a
level



Compaction **Granularity**

data moved per compaction

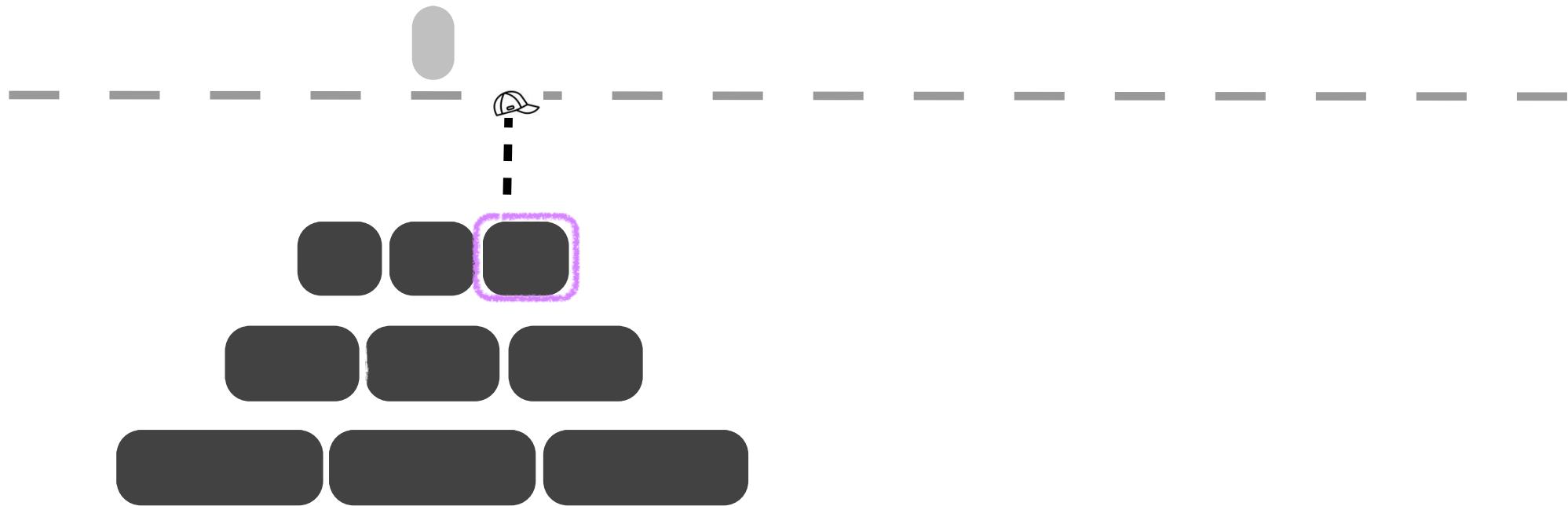


sorted runs in a
level



Compaction **Granularity**

data moved per compaction

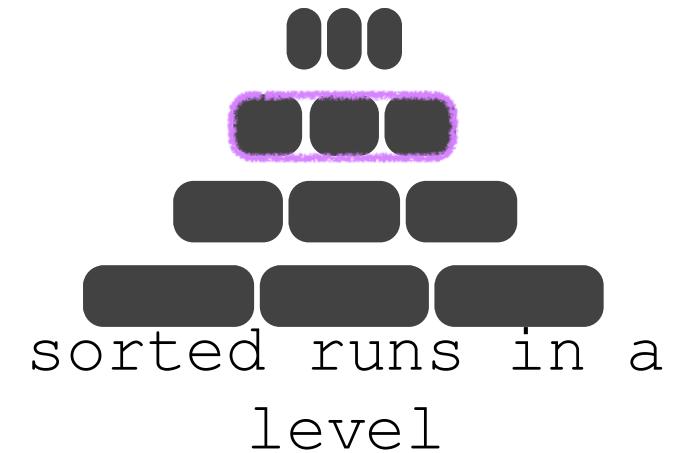
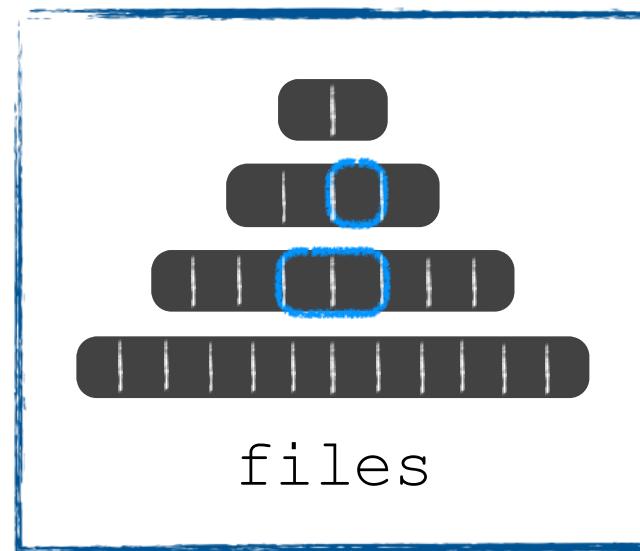


sorted runs in a
level



Compaction **Granularity**

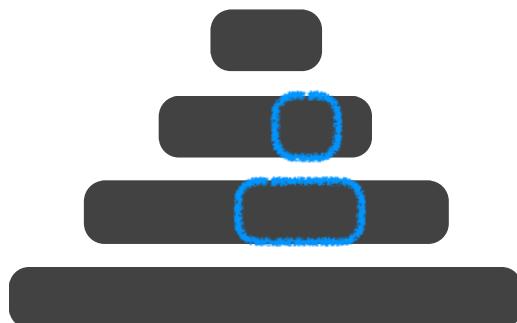
data moved per compaction





Data Movement Policy

which data to compact

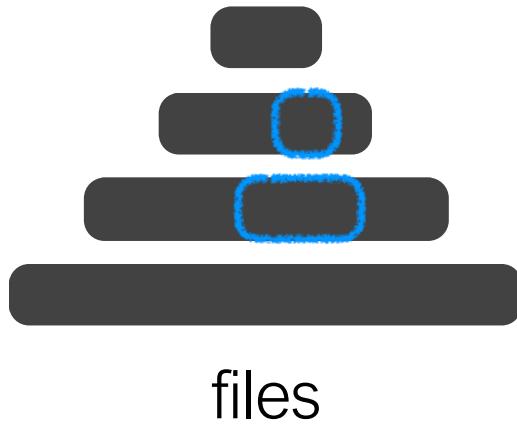


files



Data Movement Policy

which data to compact



round-robin

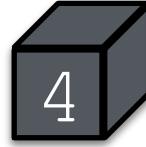


minimum **overlap with parent level**

file with most **tombstones**



coldest file

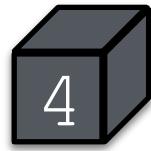


Compaction Trigger

invoking the compaction routine

level **saturation**

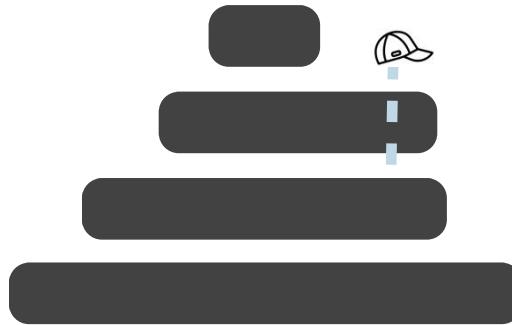


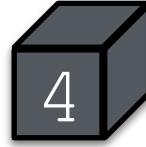


Compaction Trigger

invoking the compaction routine

level **saturation**

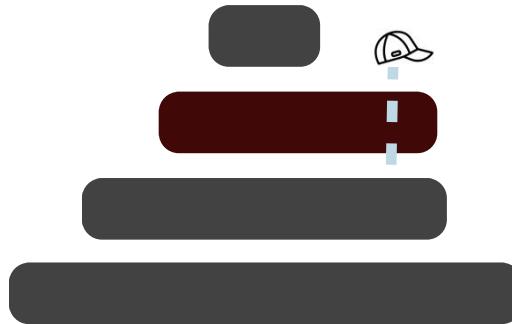




Compaction Trigger

invoking the compaction routine

level **saturation**

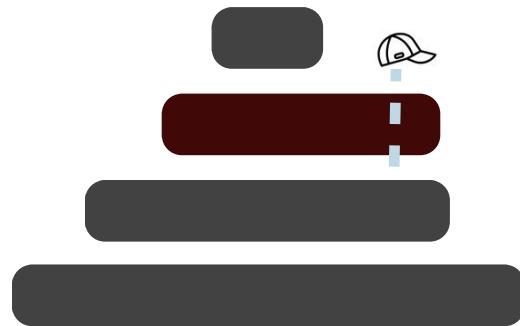




Compaction Trigger

invoking the compaction routine

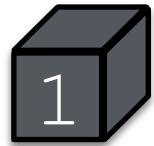
level **saturation**



number of **sorted runs**

space amplification SA

age of a file De



Data Layout



Compaction
Granularity



Data Movement
Policy



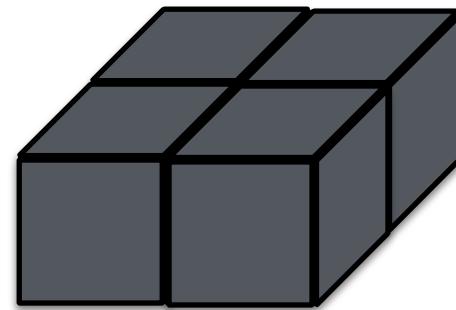
Compaction
Trigger

Data Layout

Compaction
Granularity

Data Movement
Policy

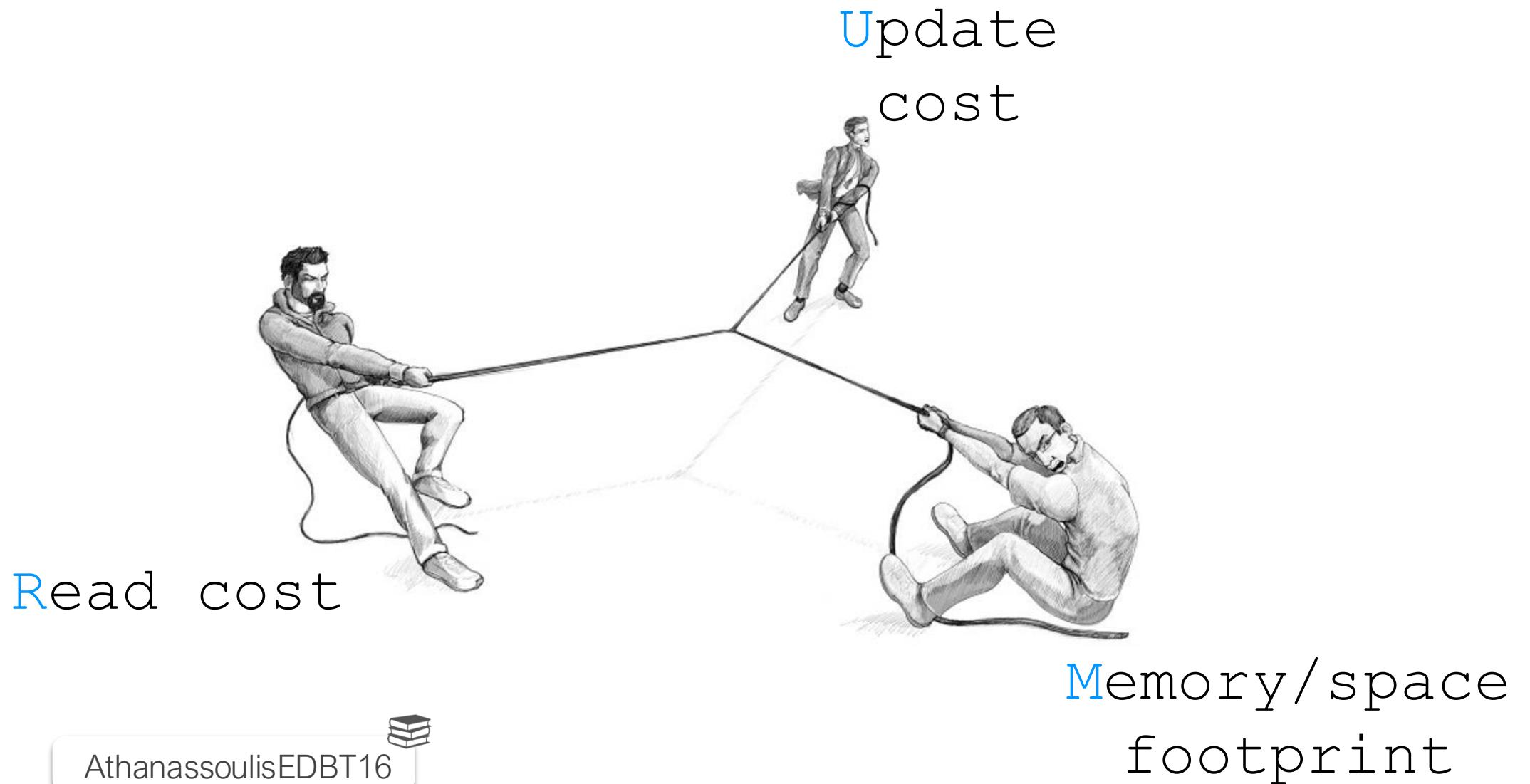
Compaction
Trigger



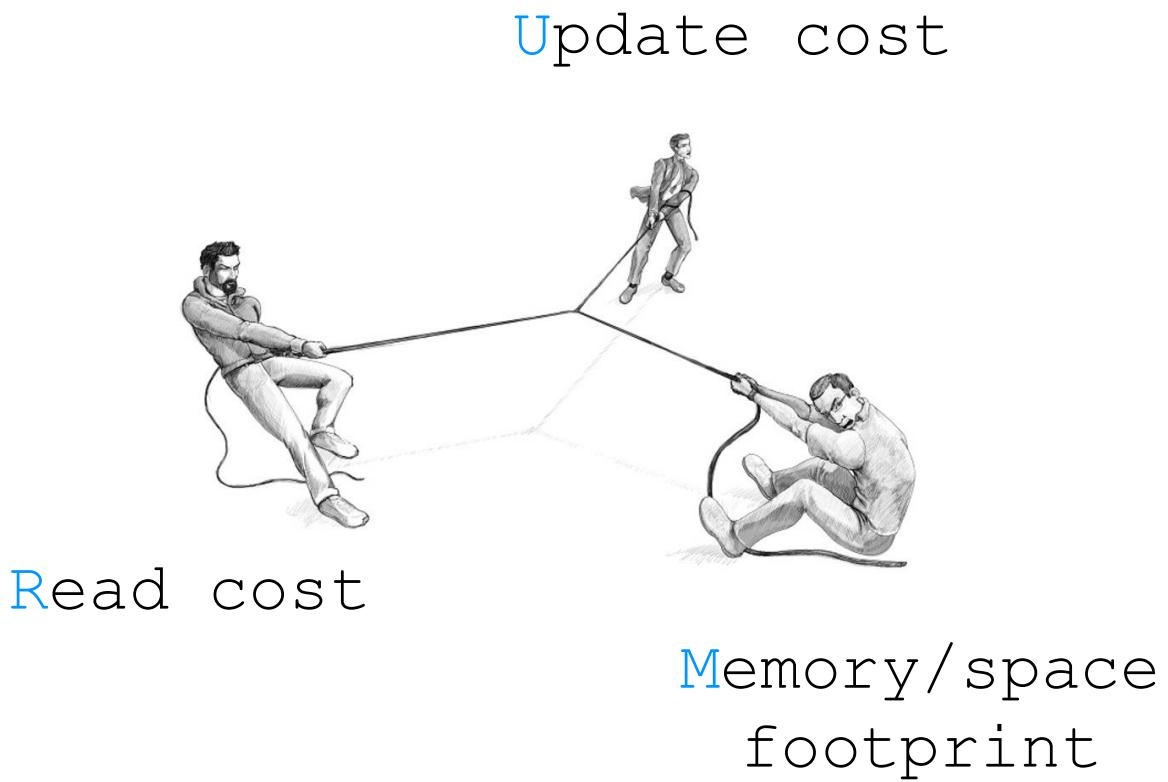
Any Compaction Algorithm

LSM Design Space

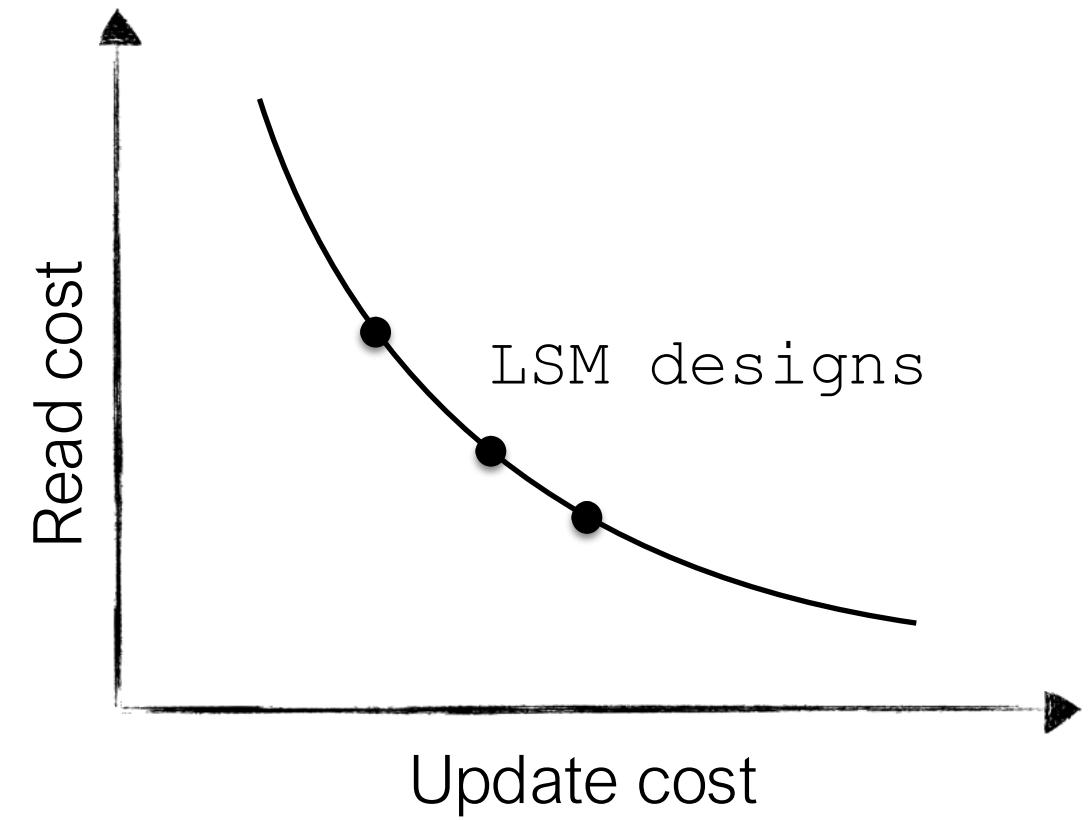
LSM Design Space



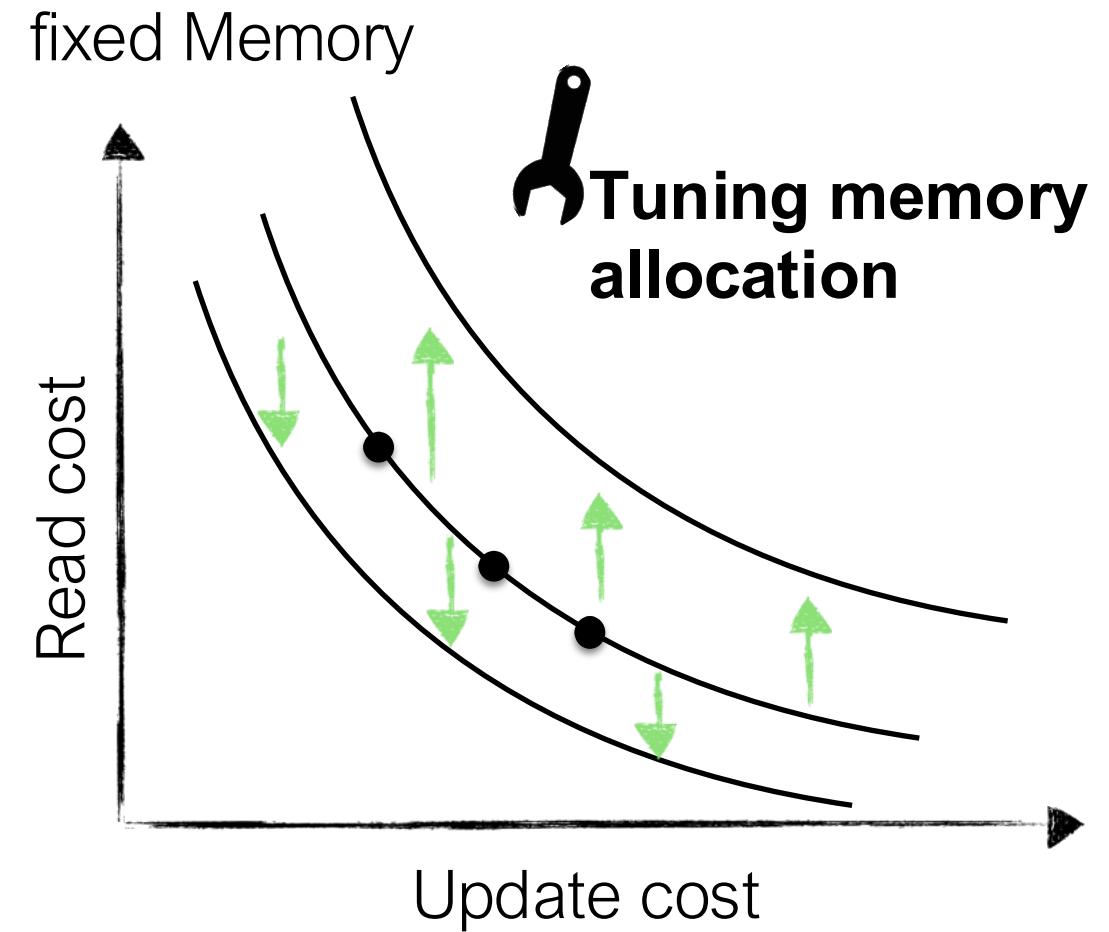
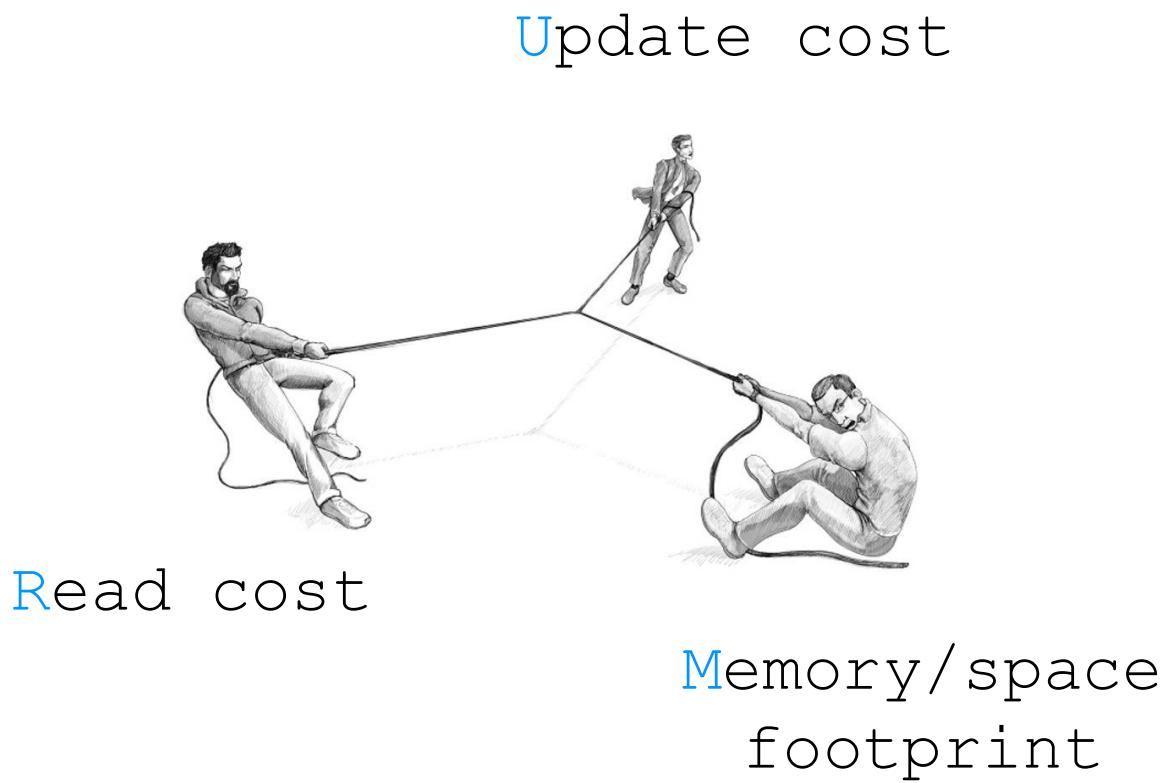
LSM Design Space



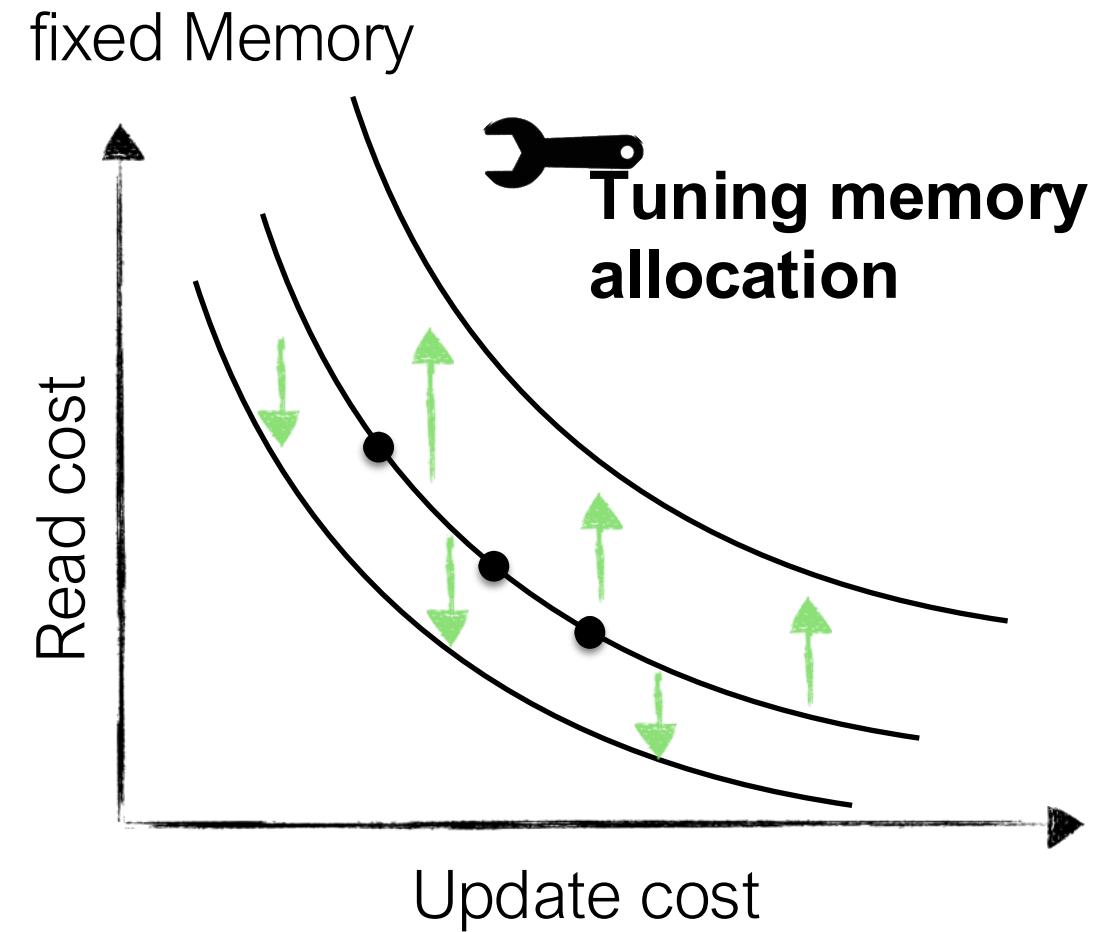
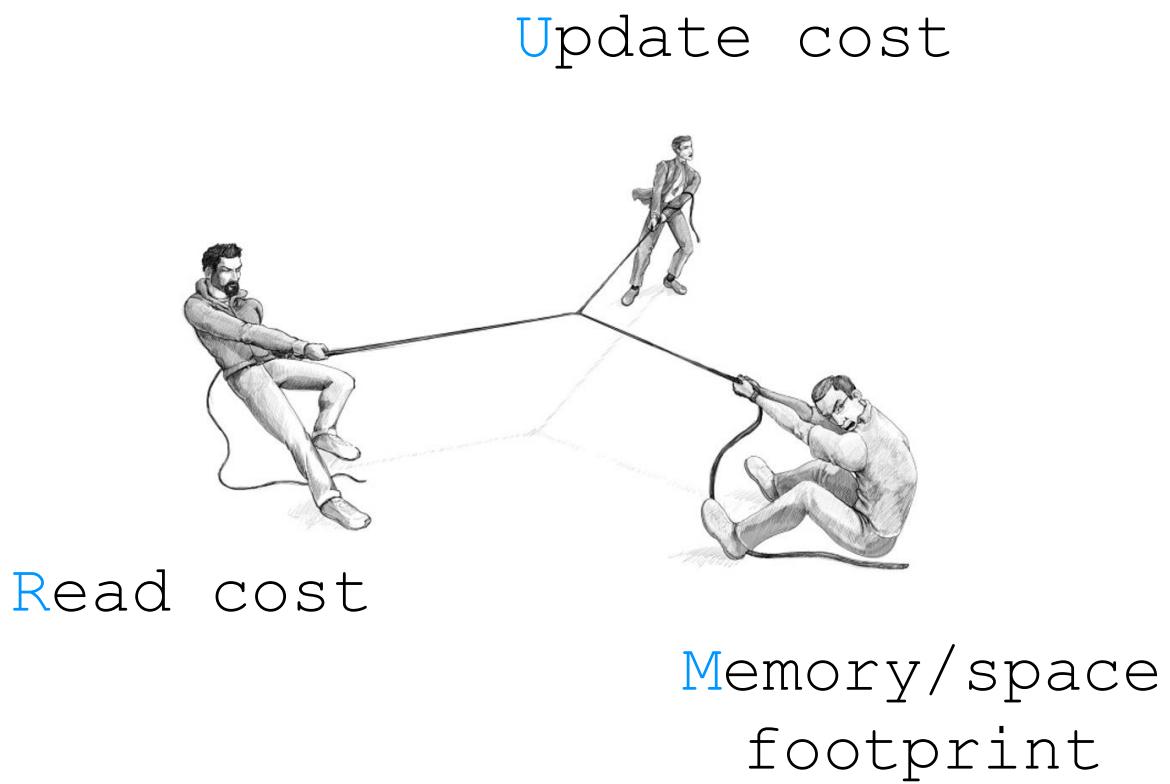
fixed Memory



LSM Design Space



LSM Design Space

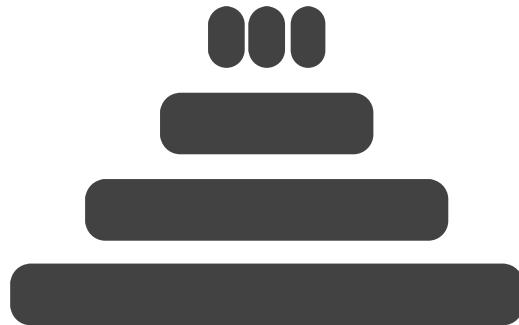


Storage Layer Design Continuum

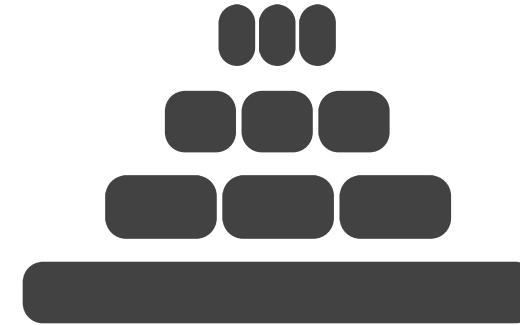
leveling



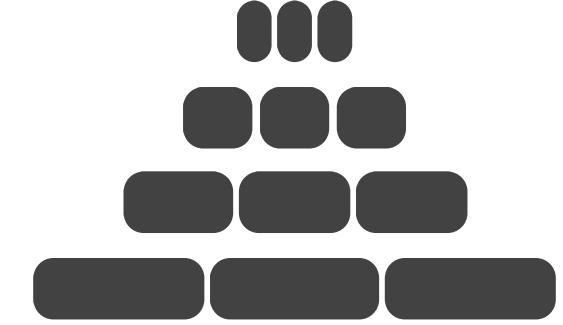
1-leveling



L-leveling



tiering

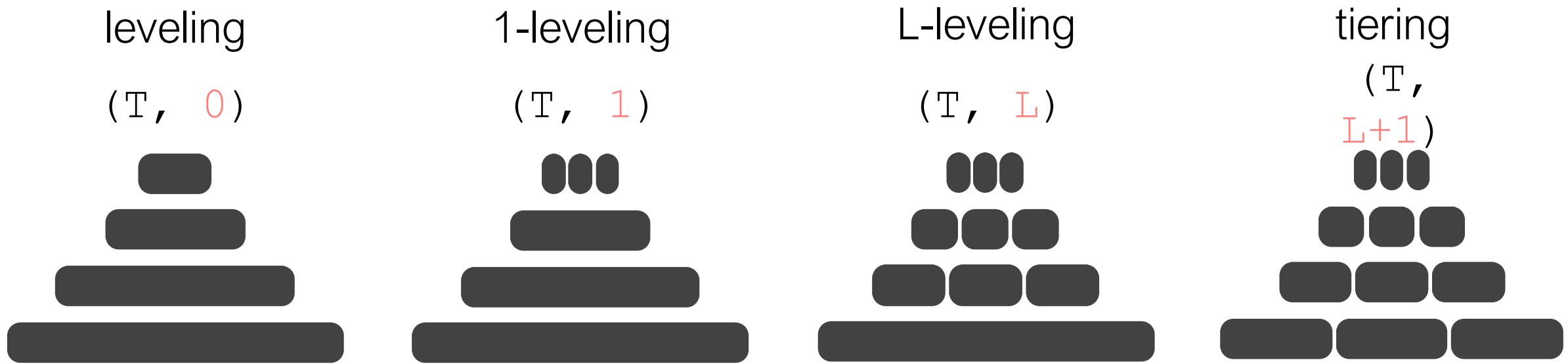


Any design can be defined by the tuple-set:

$(T,$
 $i)$



Storage Layer Design Continuum

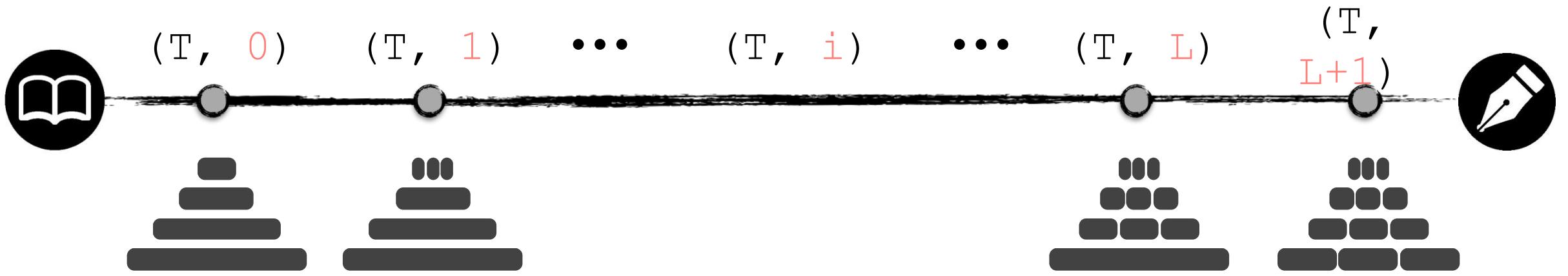


Any design can be defined by the tuple-set:

(T, i)



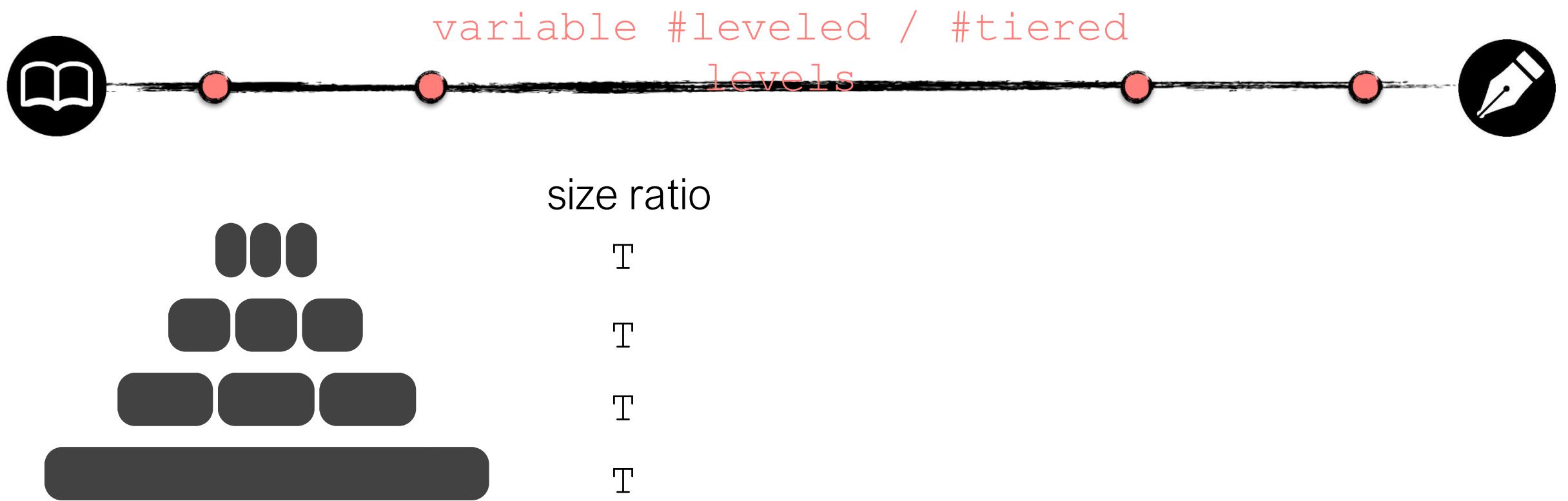
Storage Layer Design Continuum



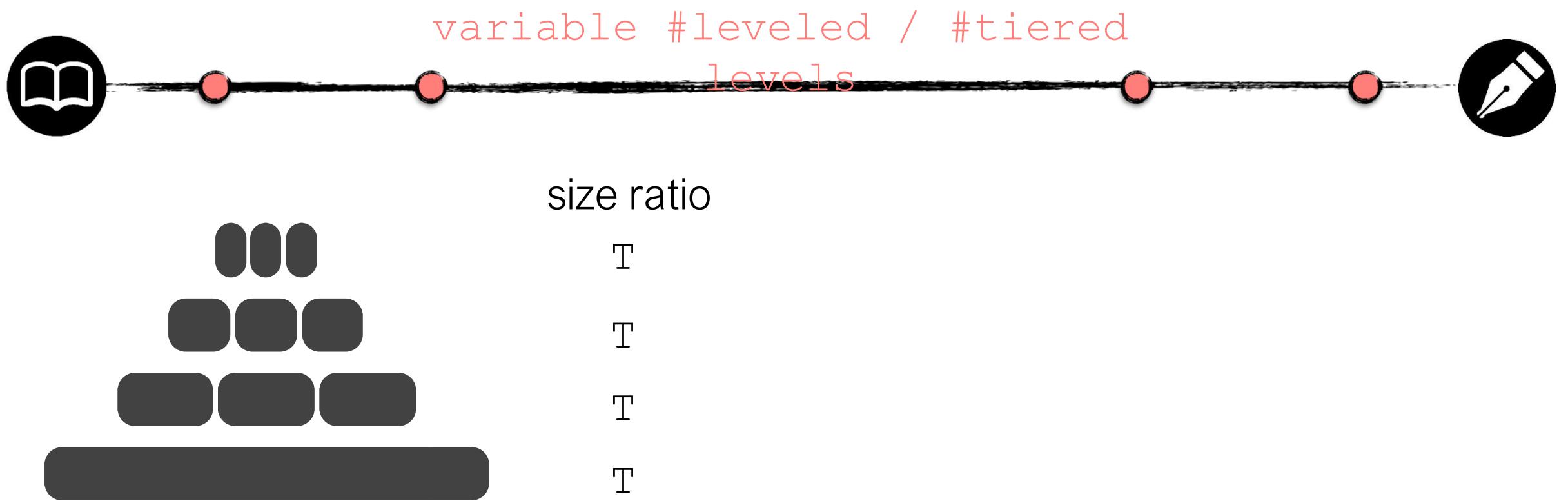
Storage Layer Design Continuum



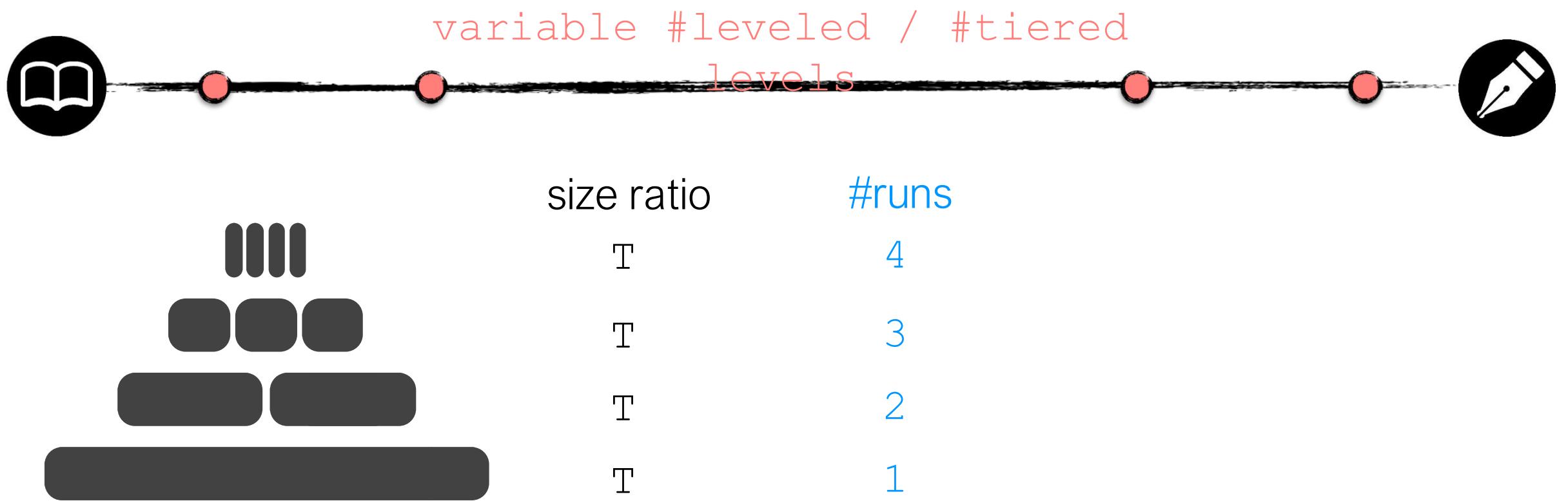
Storage Layer Design Continuum



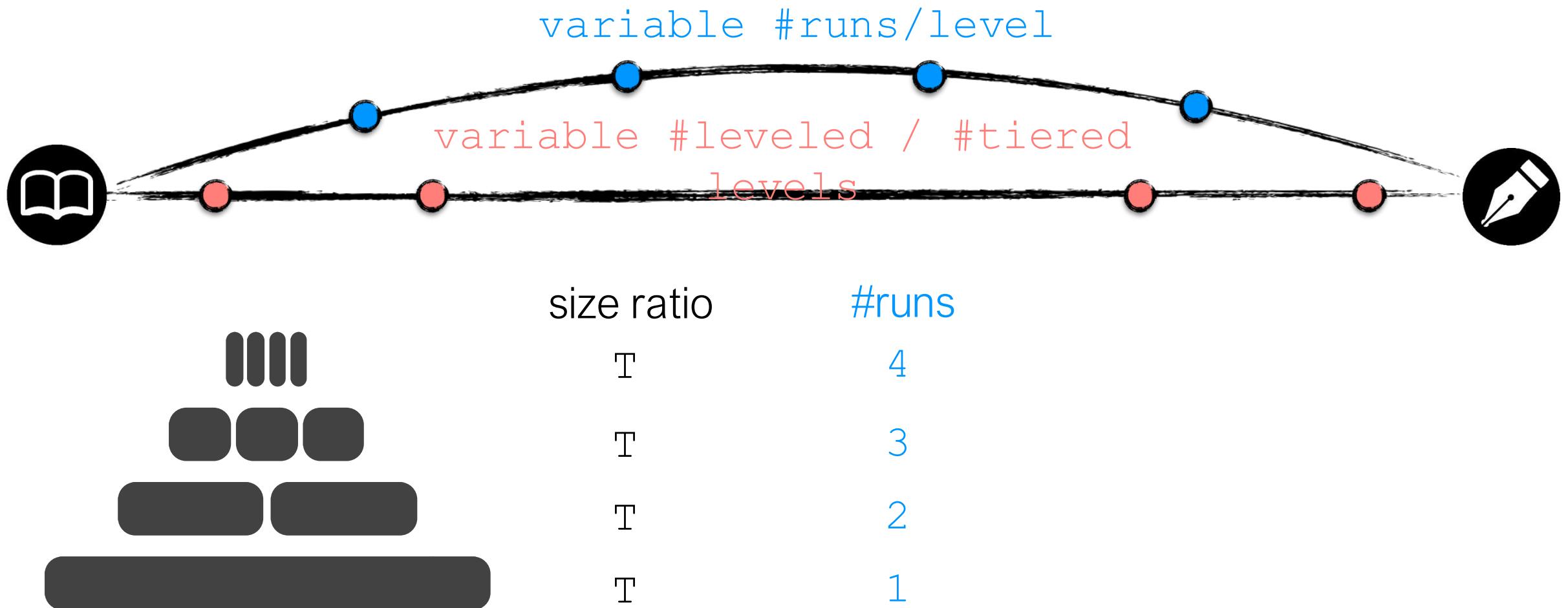
Storage Layer Design Continuum



Storage Layer Design Continuum



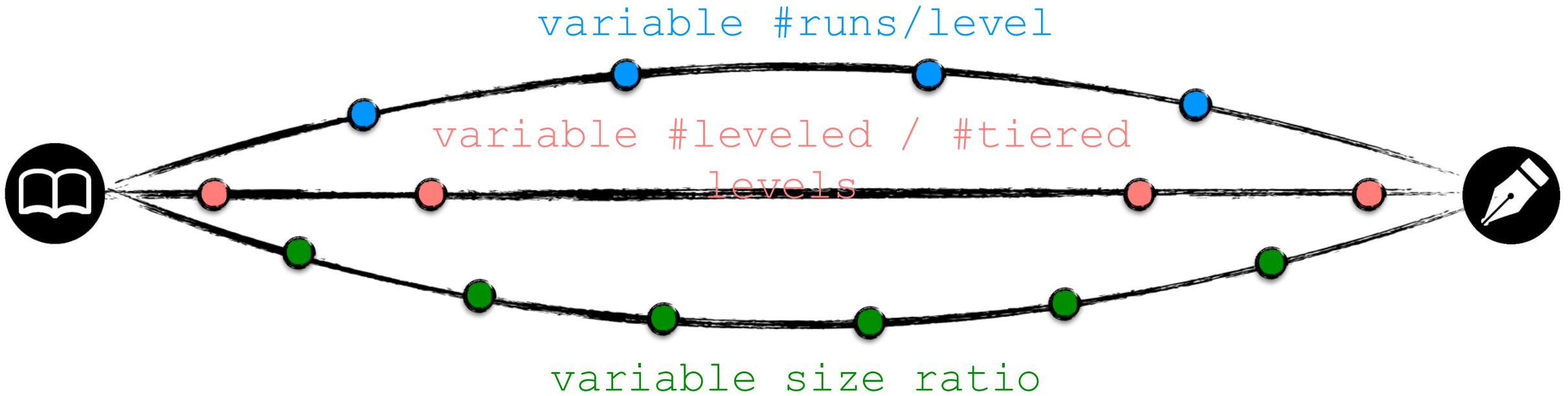
Storage Layer Design Continuum



Storage Layer Design Continuum



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/>