

CS 561: Data Systems Architecture

Class 4

Systems & Research Projects

BOSTON
UNIVERSITY

Dr. Subhadeep Sarkar

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



Updates: Logistics

Some students in the waitlist have been added (deadline: **tomorrow**).

Order for Classes 4 through 7 has been changed.

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

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

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

Reminder: Presentations

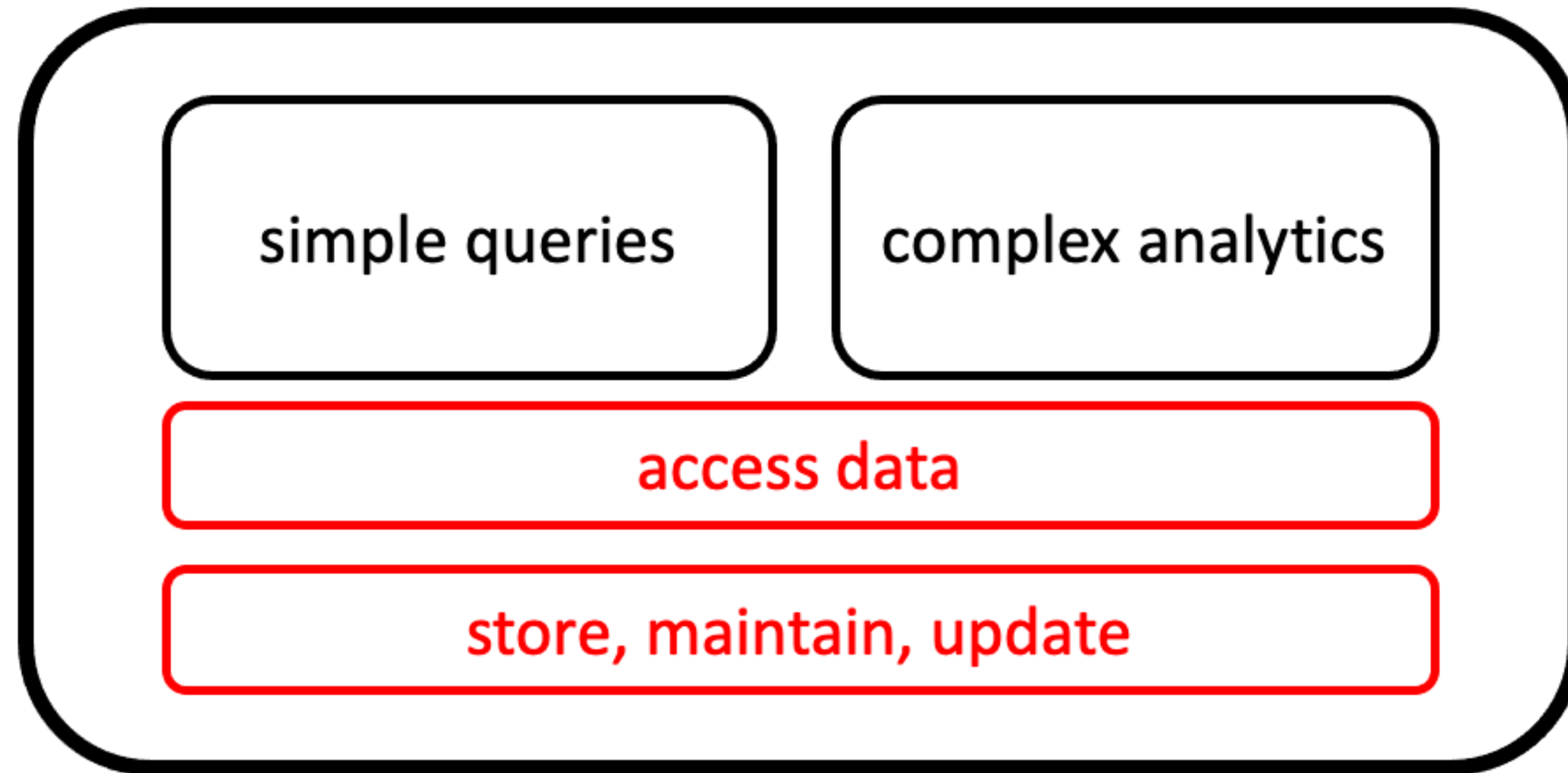
The first **student presentation** is in two weeks (on **02/14**)!

If you haven't done already, **select the paper** you will work on for your **presentation** (groups of 3-4 students)

<https://docs.google.com/document/d/1Lr-pwN7YYjL9Qk7riVvABwZyOhvVmAZsCjToghIYs8E/>

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

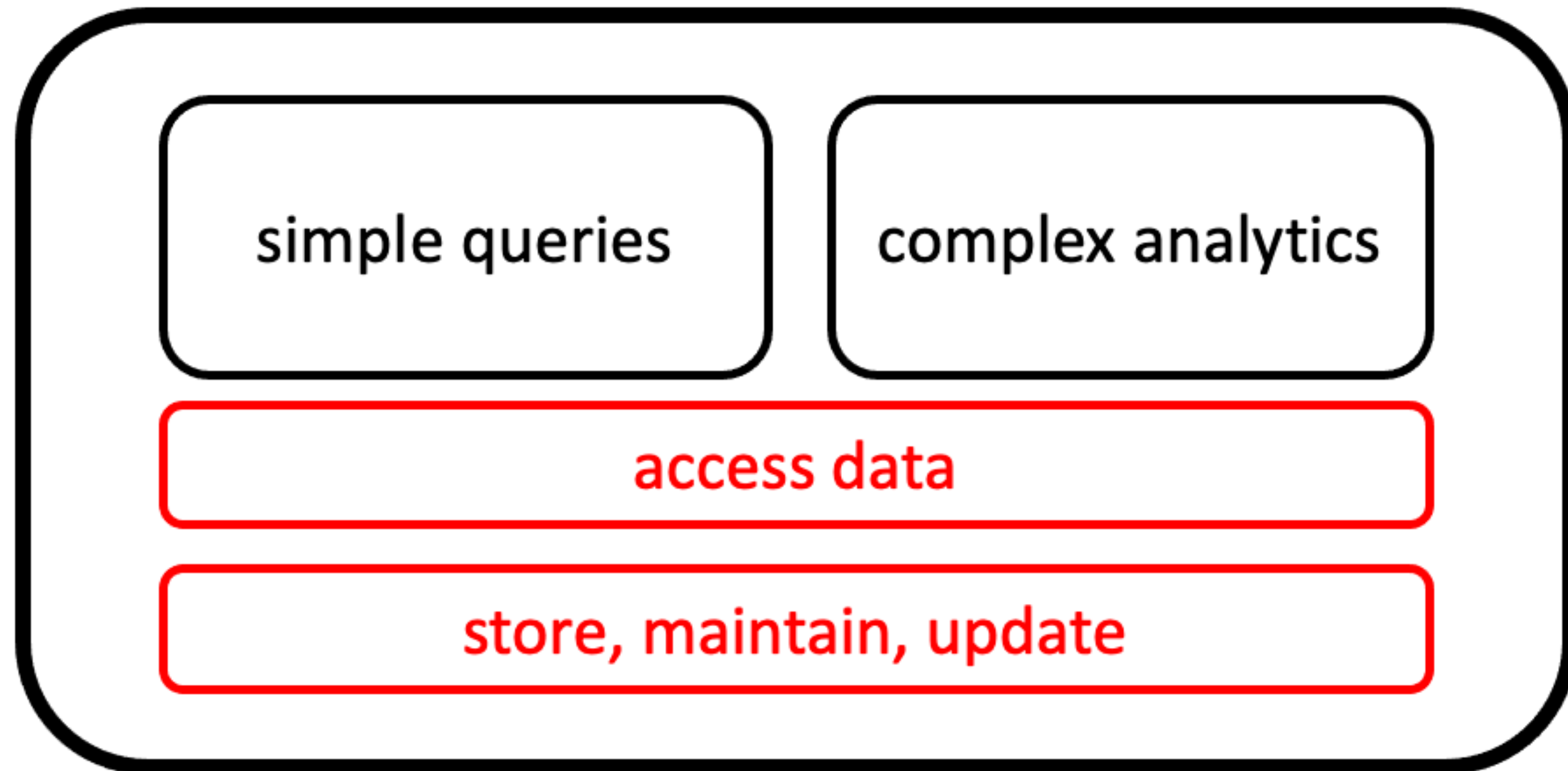
Data Systems



access methods

algorithms and data structures
for organizing and accessing data

Data Systems



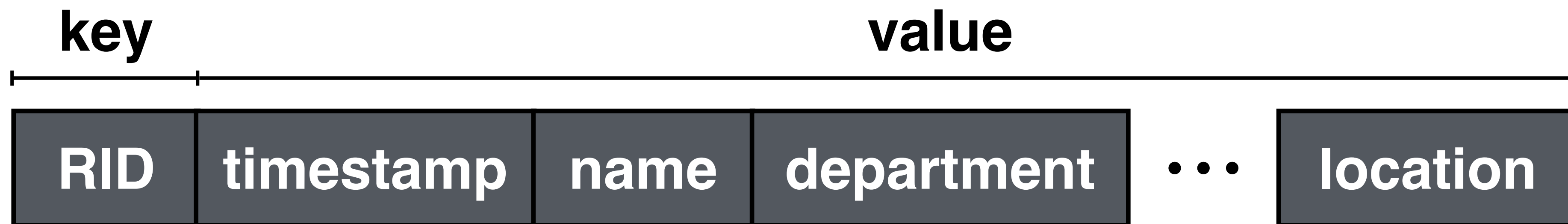
- how to **store** data?
- how to **access** data?
- how to **update** data?

At the core of our Research/Systems Projects!

Projects

Theme: **NoSQL key-value stores**

key-value pairs



Projects

Theme: **NoSQL key-value stores**

key-value pairs



How general is a key-value store?

can we store relational data?



yes! {<primary_key>, <rest_of_the_row>}

example: { student_id, { name, loginID, yob, gpa } }



what is the caveat?

how to index these attributes?



index: { loginID, { student_id } }



How to efficiently code if entry size is variable?

How to use a key-value store?

The key-value API

insert: `put(k,v)`

PQ: `{v} = get(k)` `{v1,v2,...} = get(k)` `{v1,v2,...} = get_set(k1,k2,...)`

RQ: `{v1,v2,...} = get_range(kmin,kmax)` `{v1,v2,...} = full_scan()`

count: `c = count(kmin,kmax)`

delete: `delete(k)` `delete(kmin,kmax)`

update: `update(k,vnew)` **not very different from put**

anything else?



How to build a key-value store?

append

if we have only *put* operations



if we mostly have *point get* operations



sort



if we mostly have *range get* operations

sort, find the min, scan

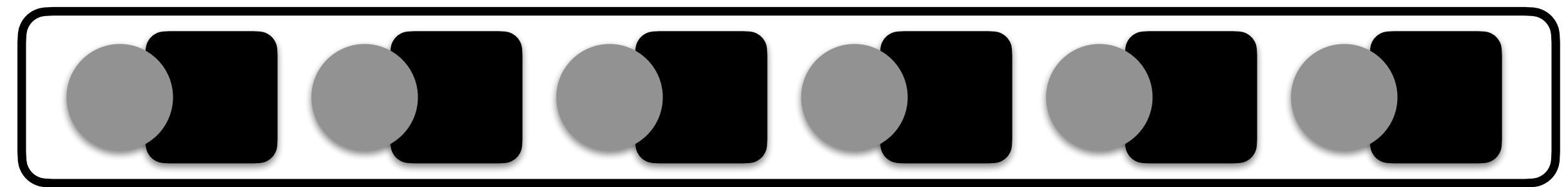


if we mostly have *full scan* operations

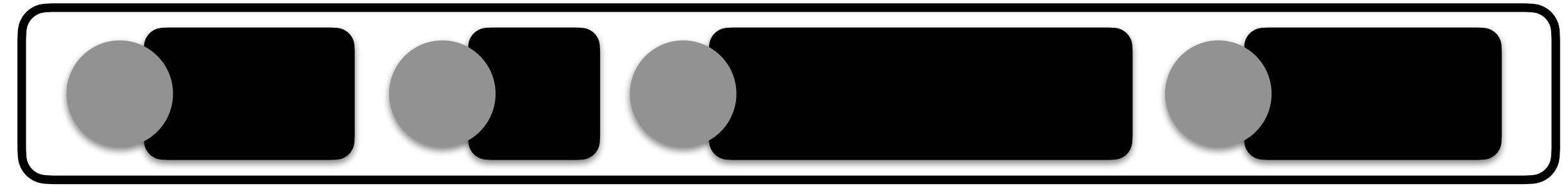
don't sort

What about **variable entry size**?

fixed entry size

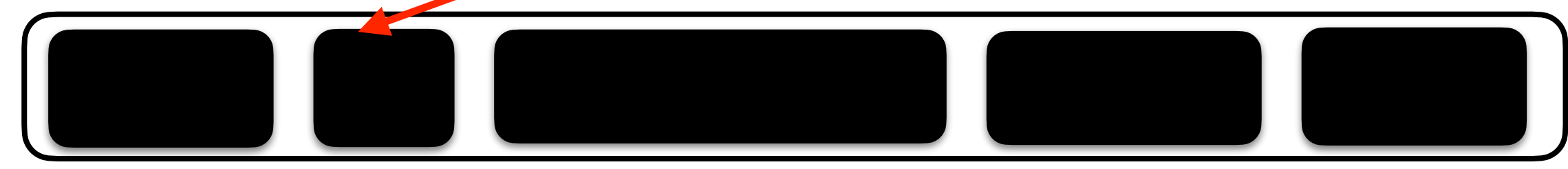
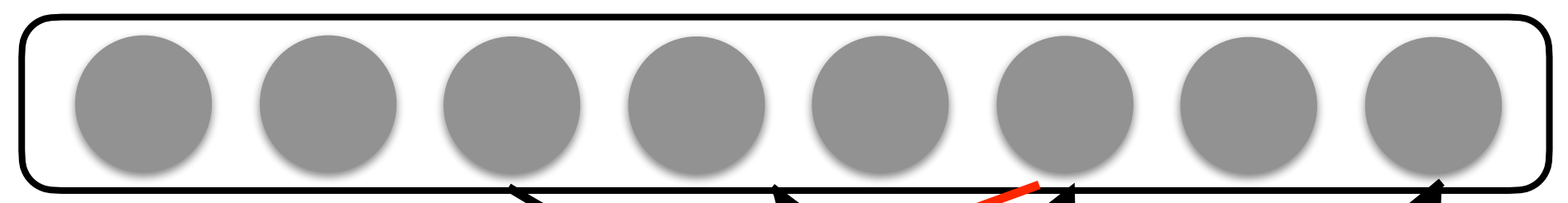


variable entry size



what can we do?

- efficient key storage
- compiler-friendly
- need not sort values



is my lunch free ?

● range queries?

● locality?

● code complexity?

Log-Structured Merge-tree

LSM-tree

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

1996

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

LSM-tree
O'Neil *et al.*

1996


Bigtable

2006

APACHE
HBASE 

2007


cassandra

2010


levelDB

2011

 RocksDB

2013

2023

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 (black and orange letters), 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.

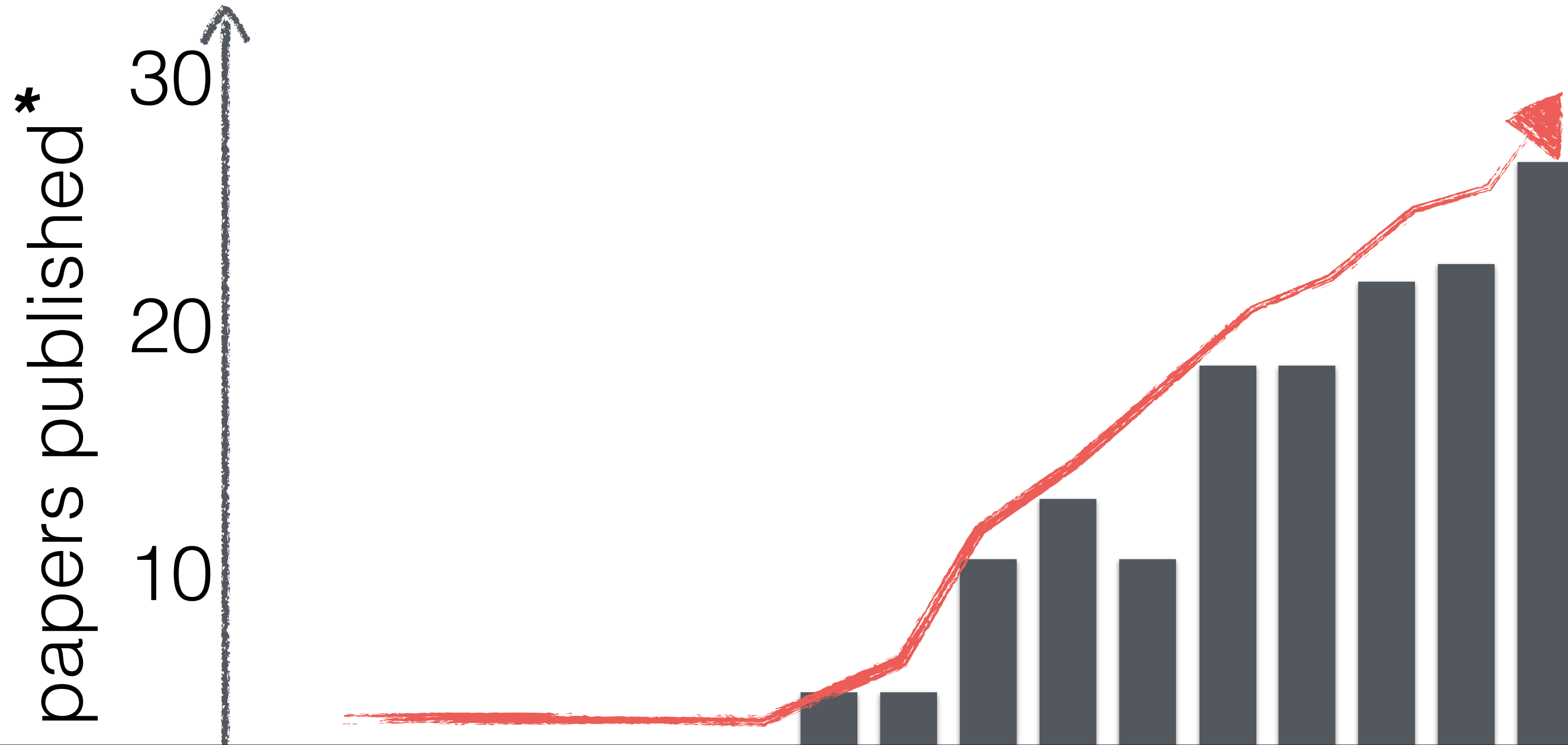
relational

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 that curves upwards.

time-series

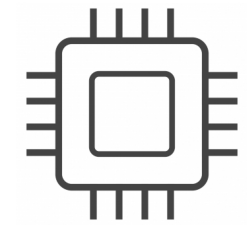
2023

Research Trend

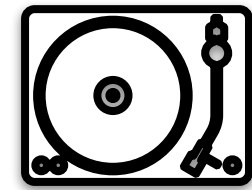


So, how does an LSM-tree look?

LSM Basics



buffer



level 1



level 2



size ratio: T

level 3



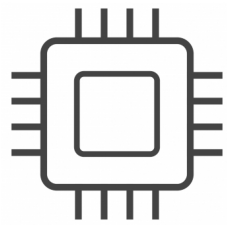
Great! But, how does it work?



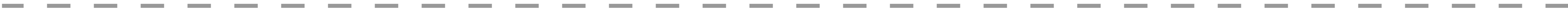
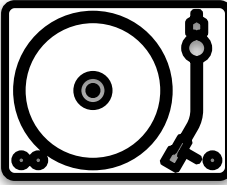
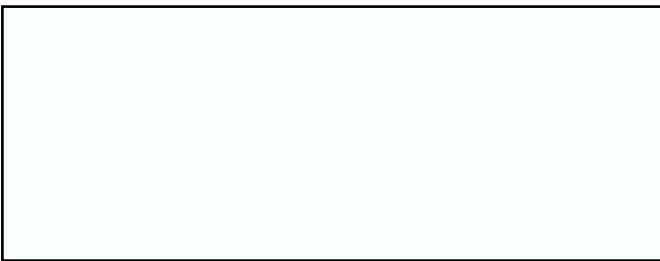
Buffering ingestion

put(6)

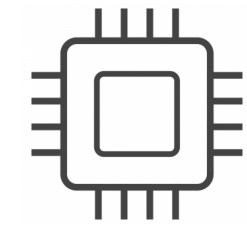
put(2)



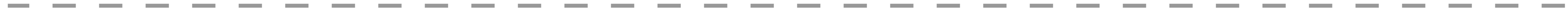
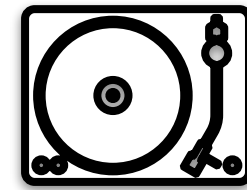
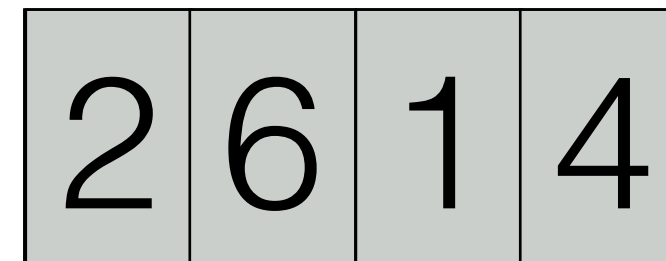
buffer



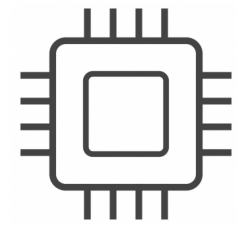
Buffering ingestion



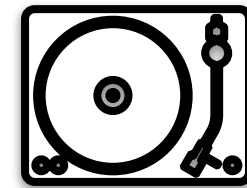
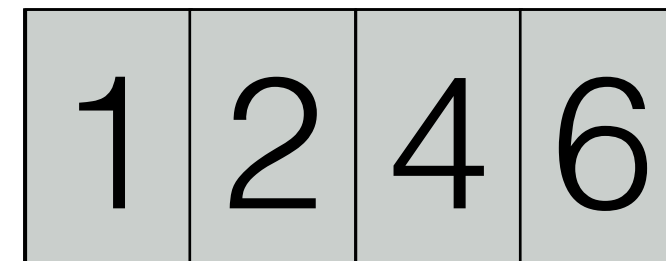
buffer



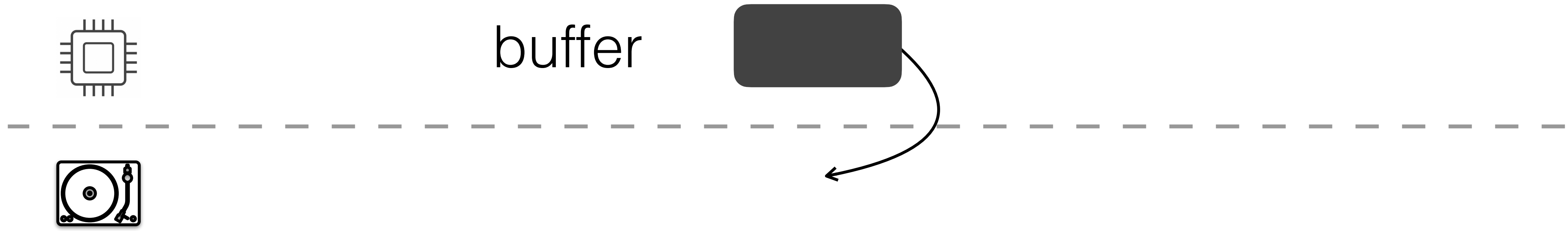
Buffering ingestion



buffer

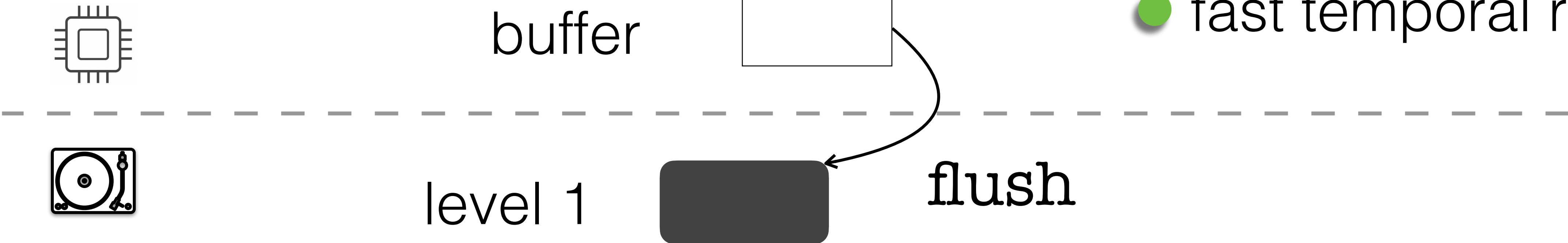


Buffering ingestion



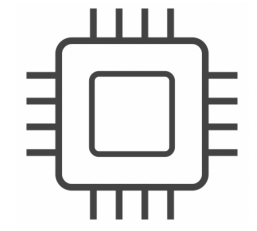
Buffering ingestion

- low ingestion cost
- fast temporal reads

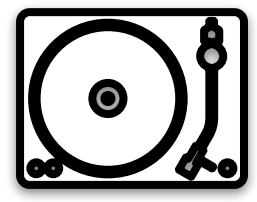
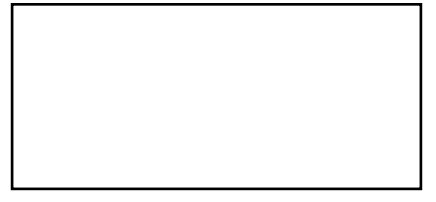


Immutable files on storage

- compact storage
- good ingestion throughput



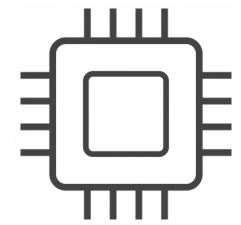
buffer



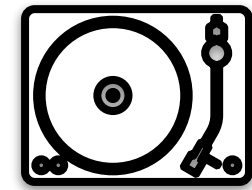
level 1



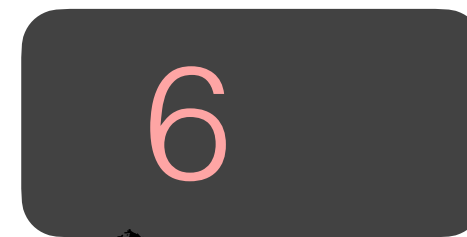
How do we update data?



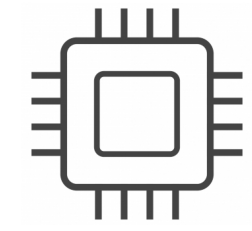
buffer



level 1

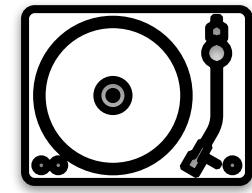
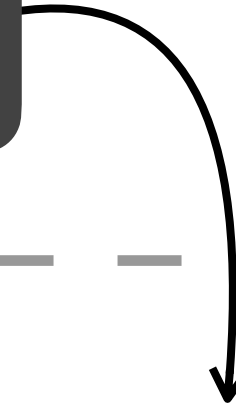


logically
invalidated



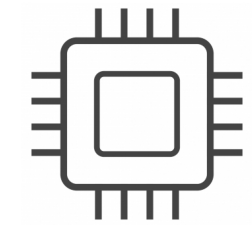
buffer

6



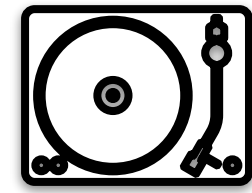
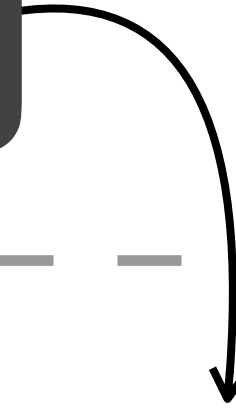
level 1

6



buffer

6

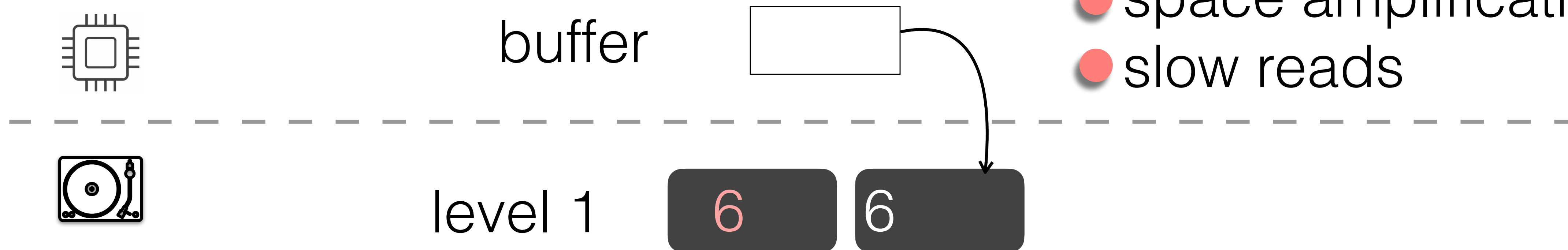


level 1

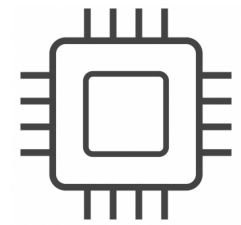
6

Out-of-place updates

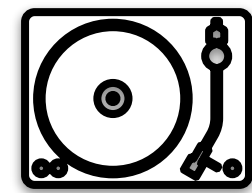
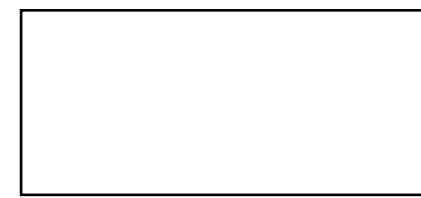
- fast ingestion
- space amplification
- slow reads



How do we reduce this space amplification?



buffer

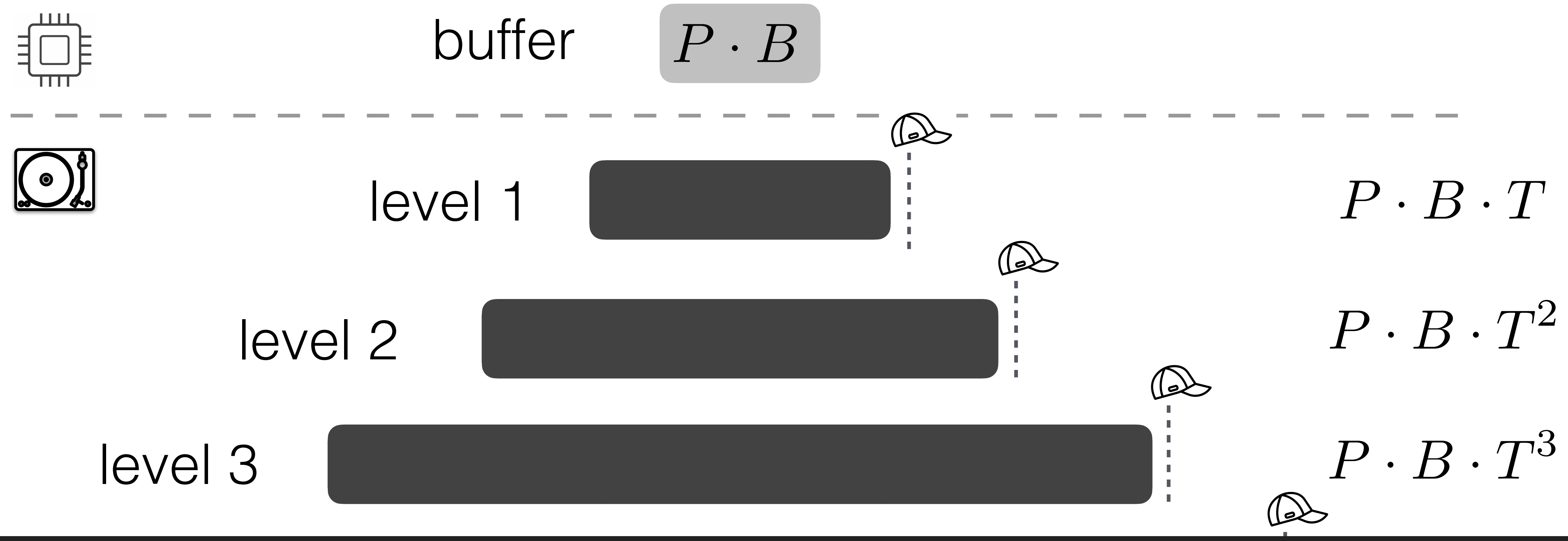


level 1

6

compaction

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

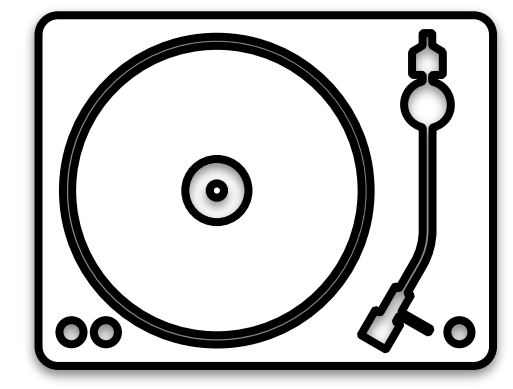
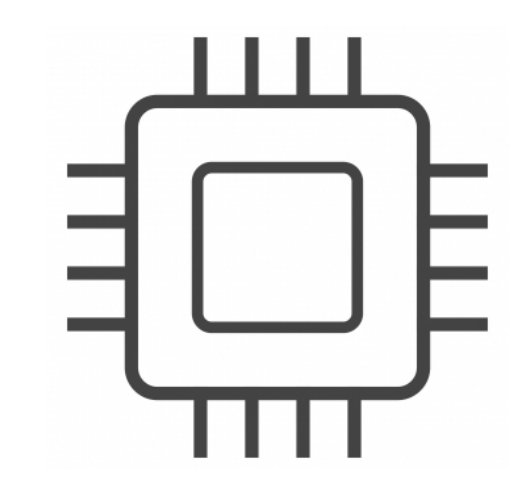


How about queries?

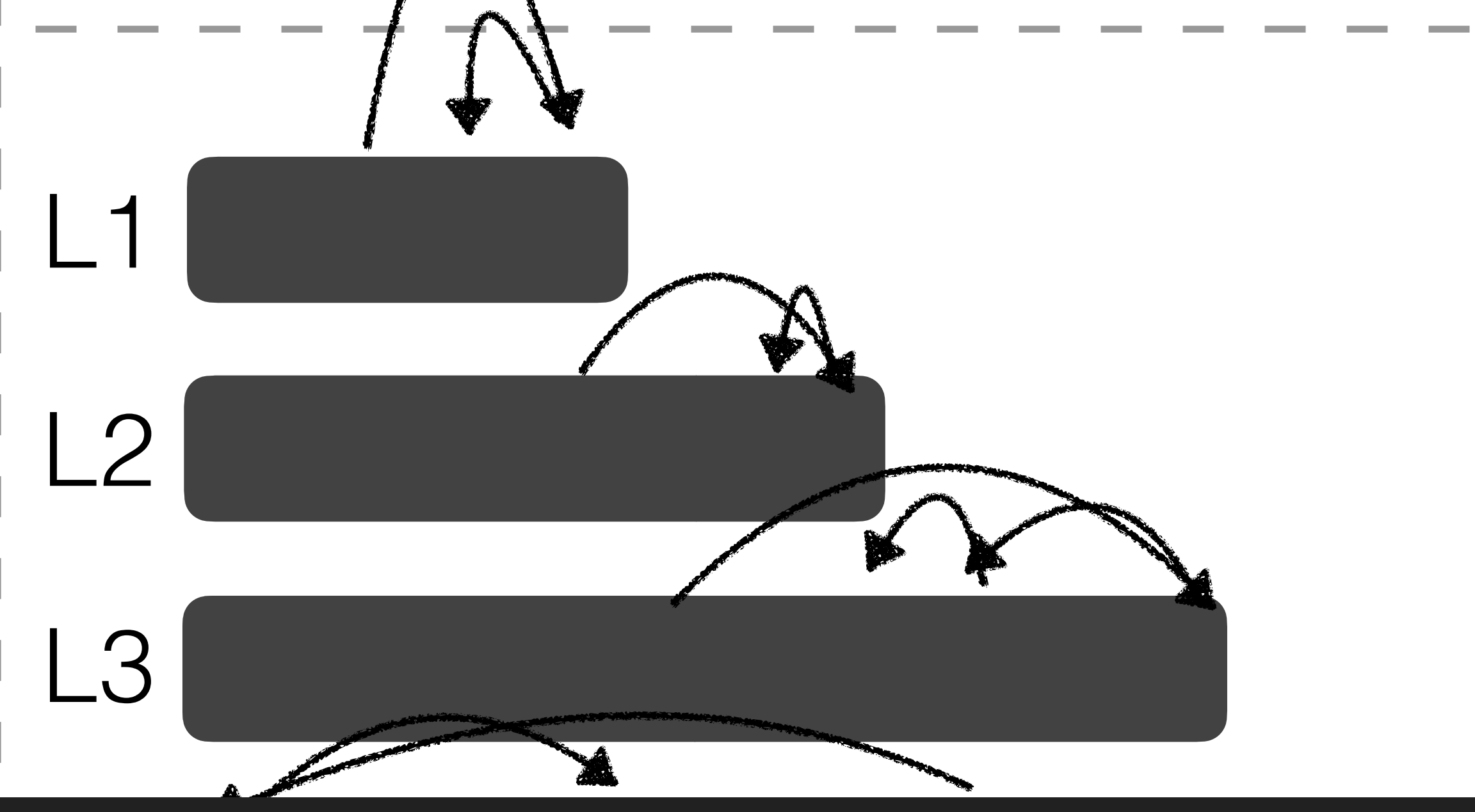
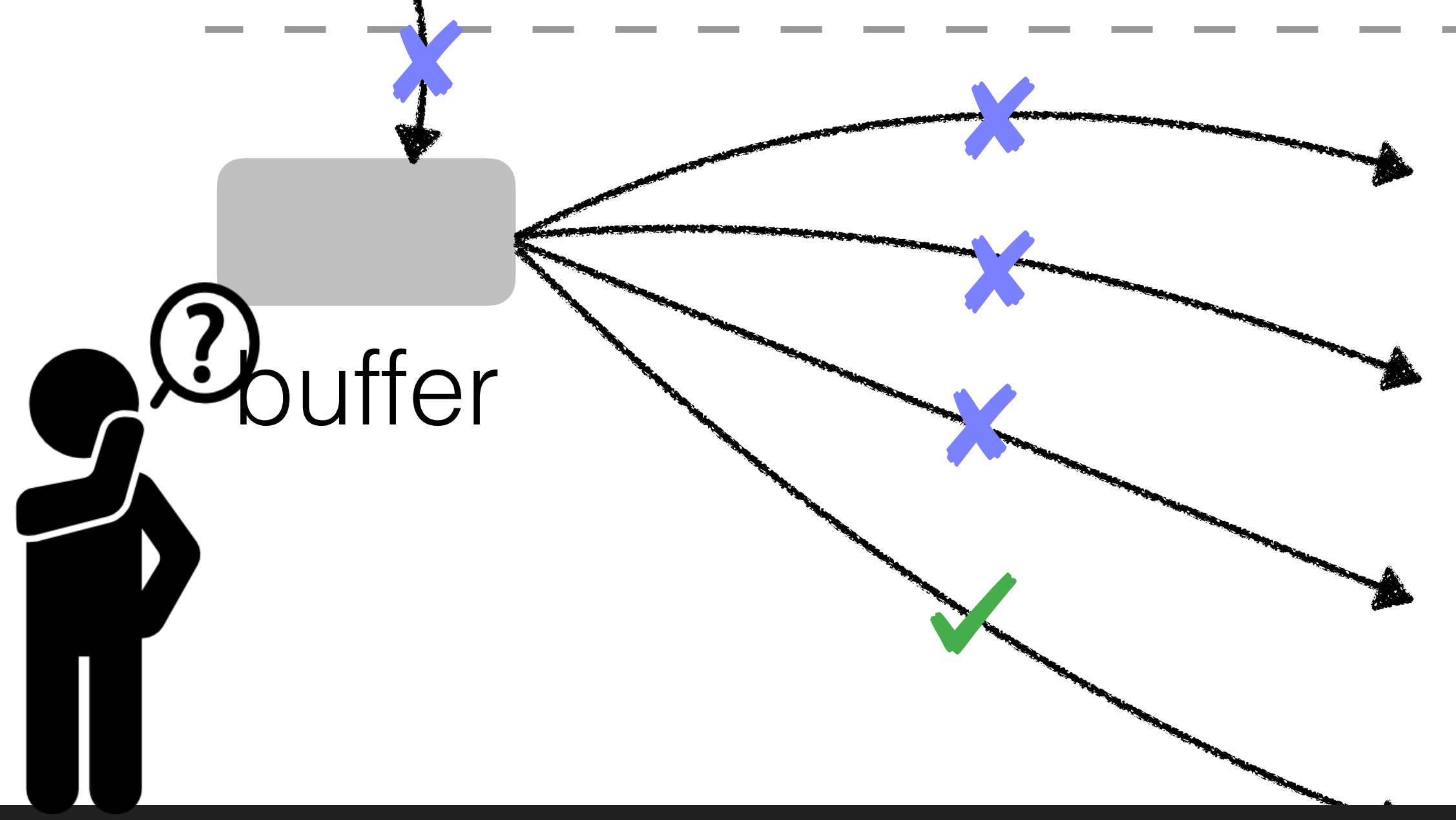
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

Cost analysis

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



get(7)



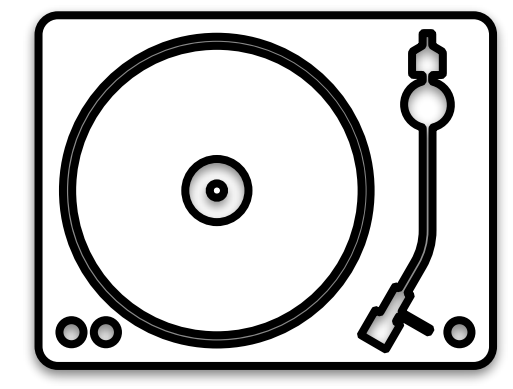
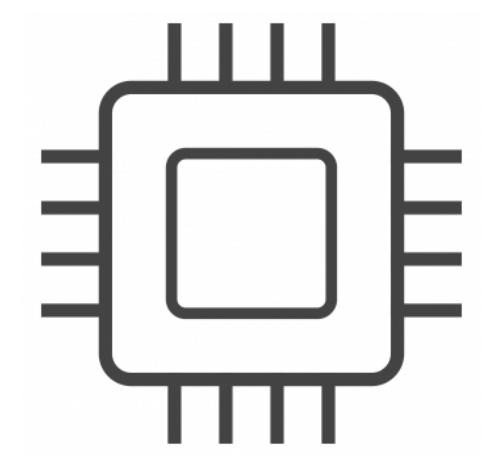
Can we do better?

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

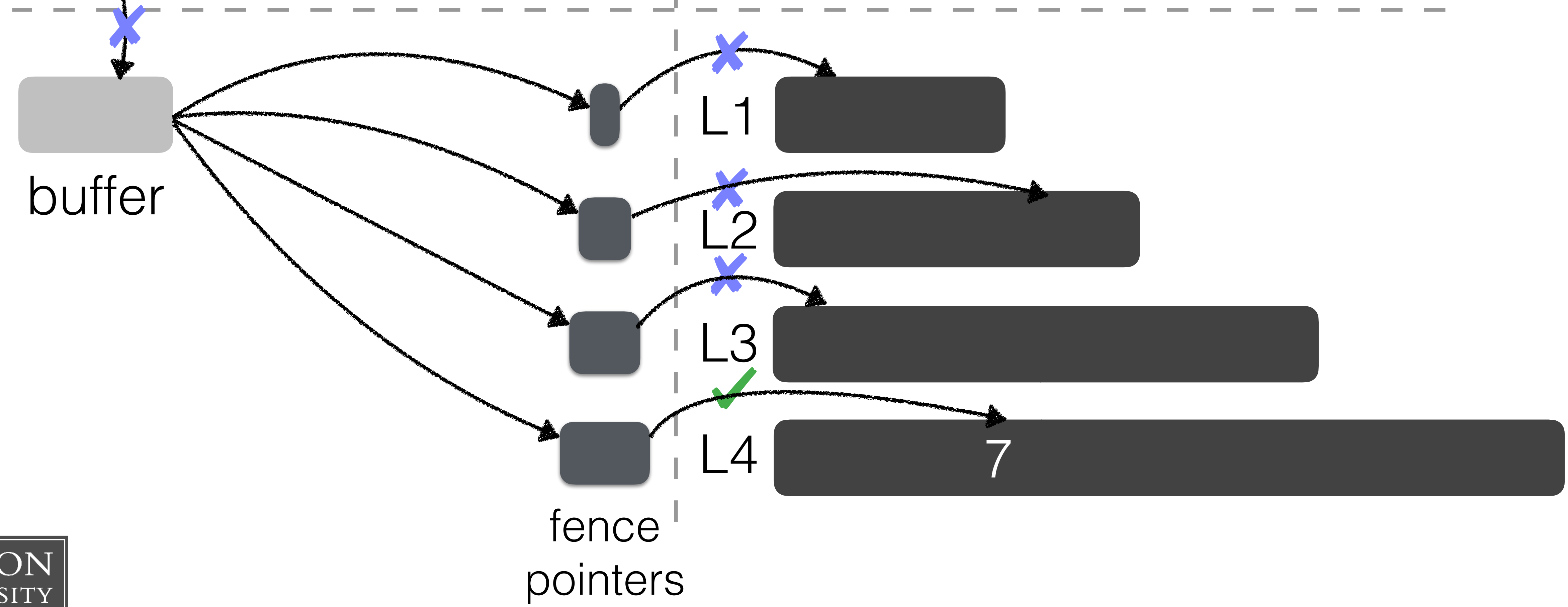
Cost analysis

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

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



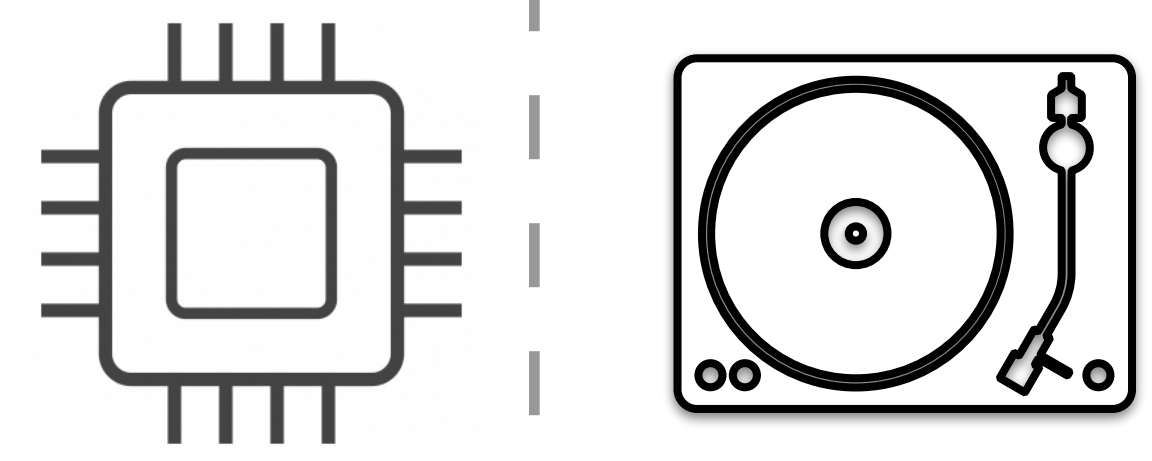
get(7)



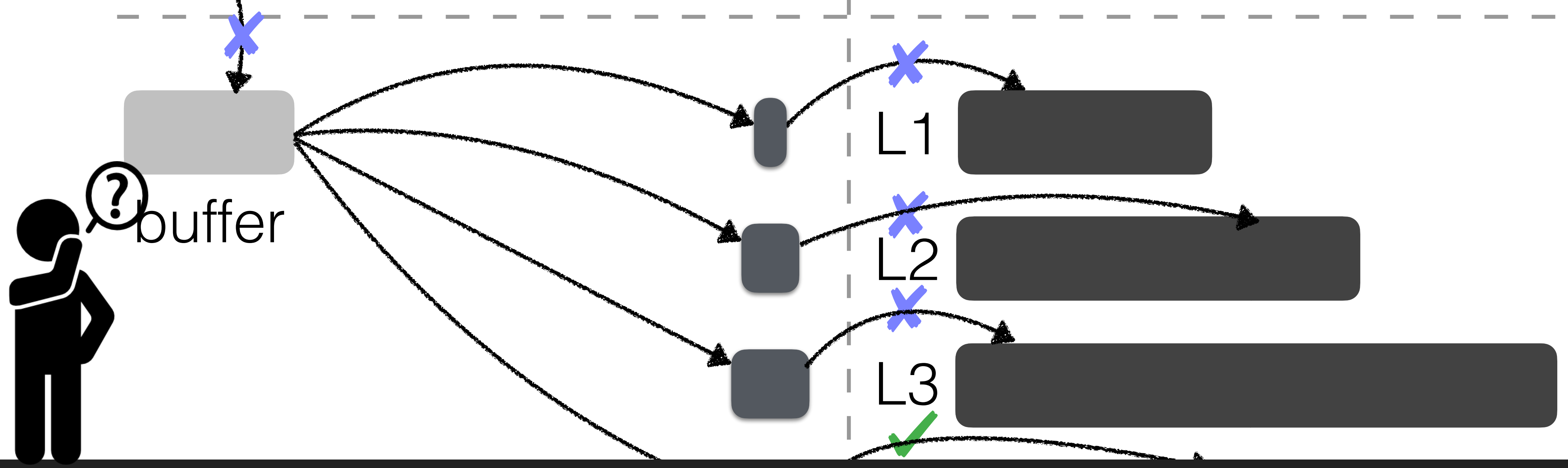
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries

Cost analysis

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



get(7)



Still expensive! Can't we do any better?

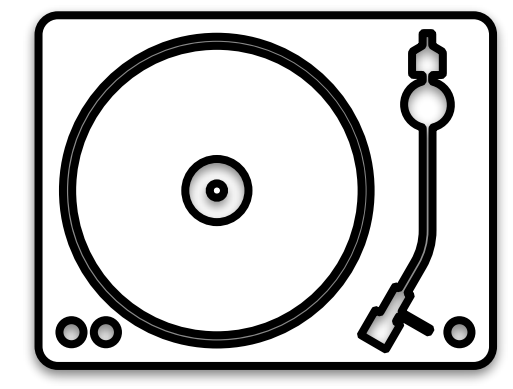
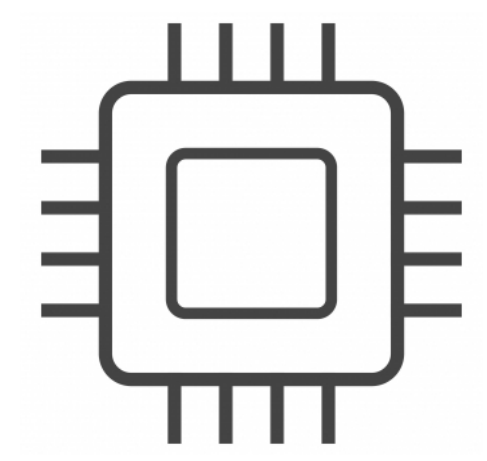
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

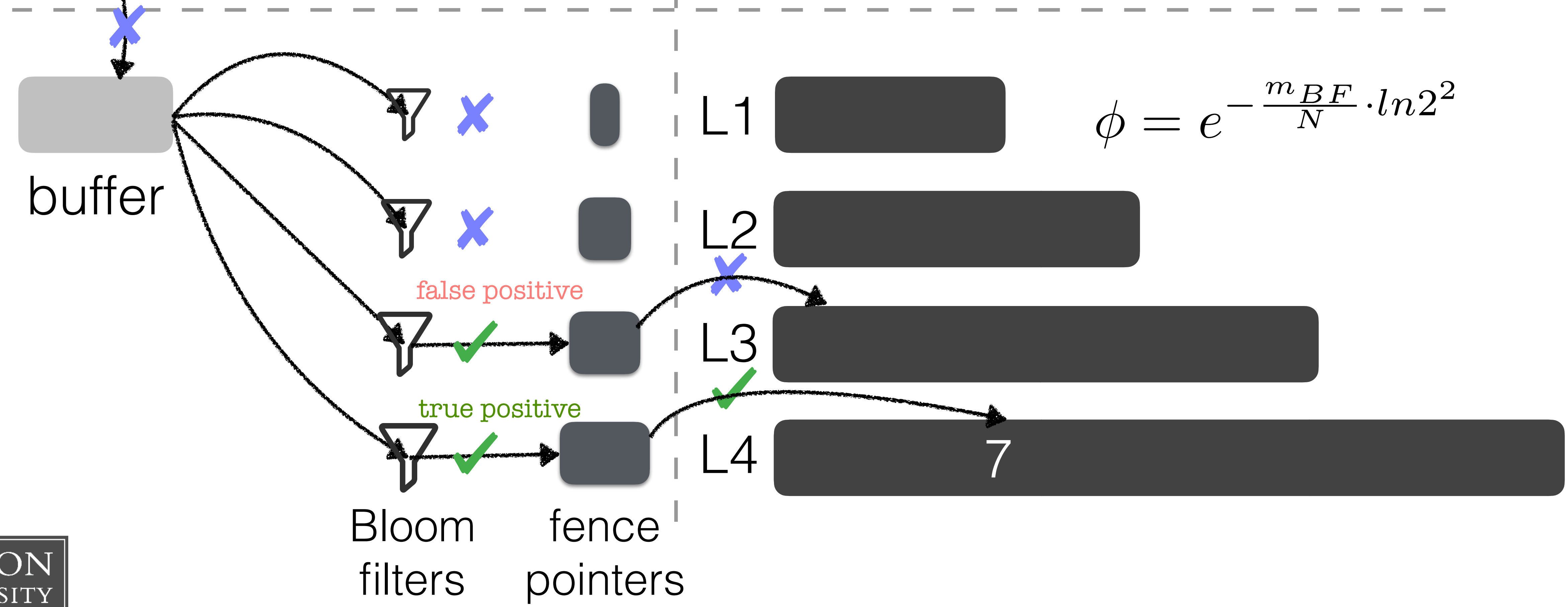
w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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

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



get(7)



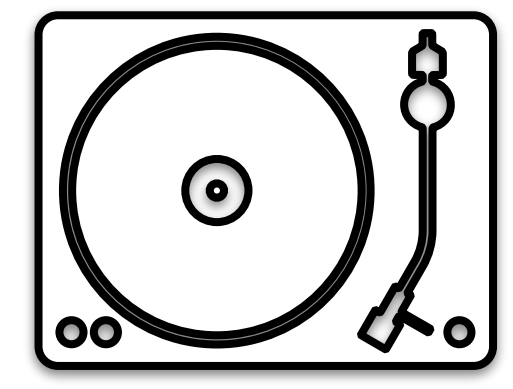
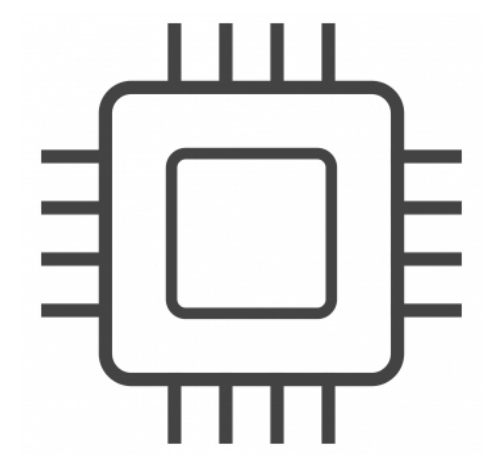
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

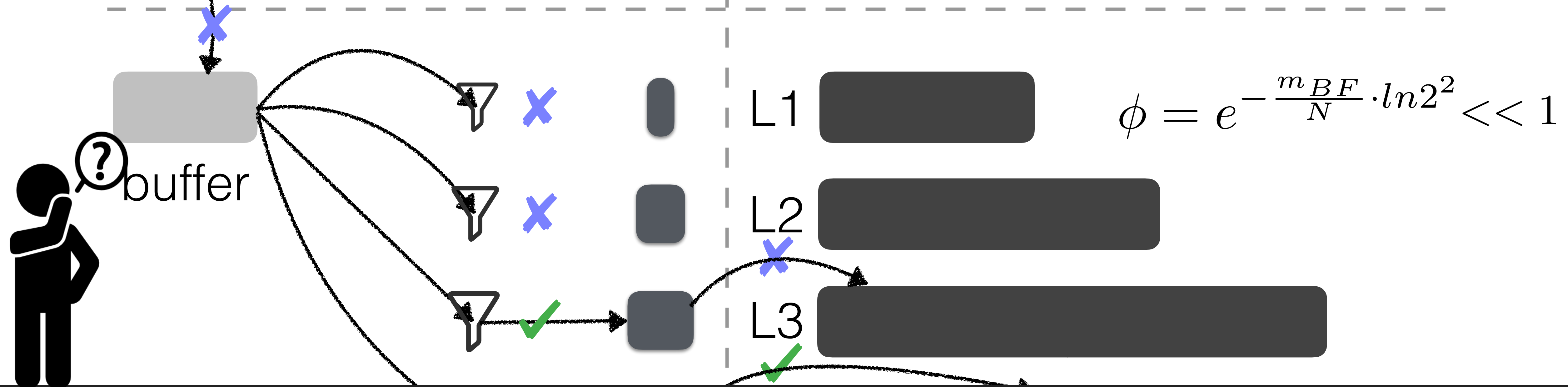
w/o F&I: $\mathcal{O}(\sum_i \log_2(P \cdot T^i))$

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

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



get(7)



How to manage memory?

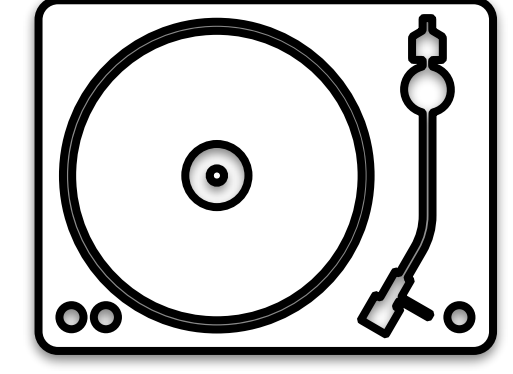
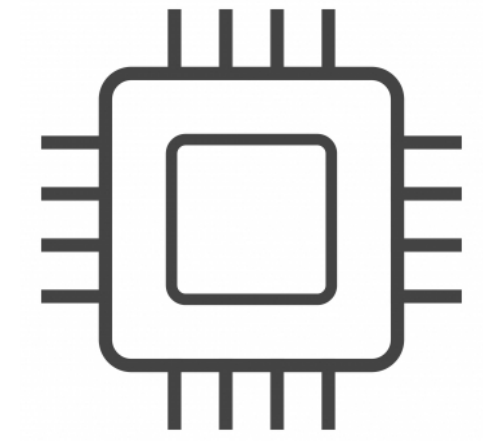
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

Cost analysis

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

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

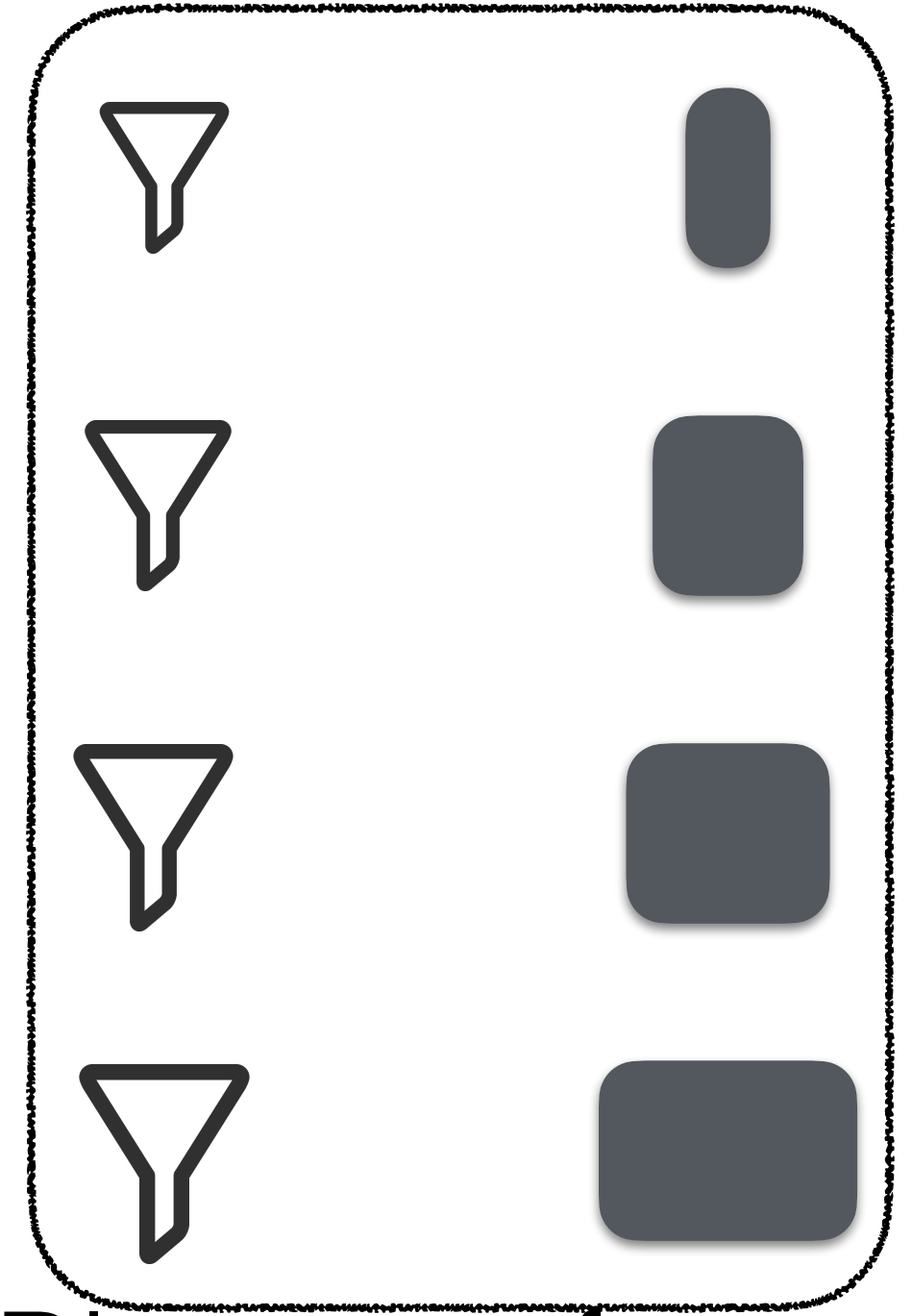
w F&I: $\mathcal{O}(\phi \cdot L)$



buffer



block cache



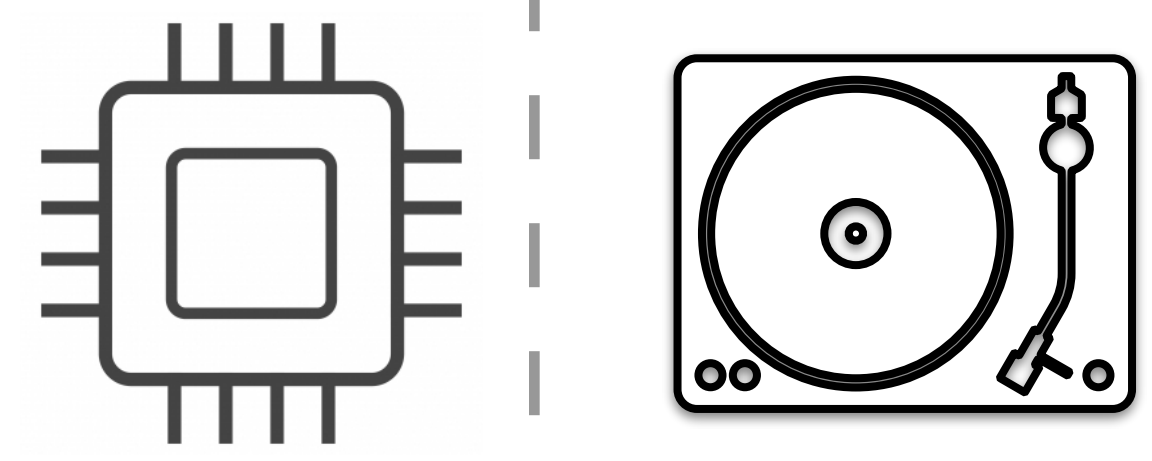
Bloom filters fence pointers



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

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

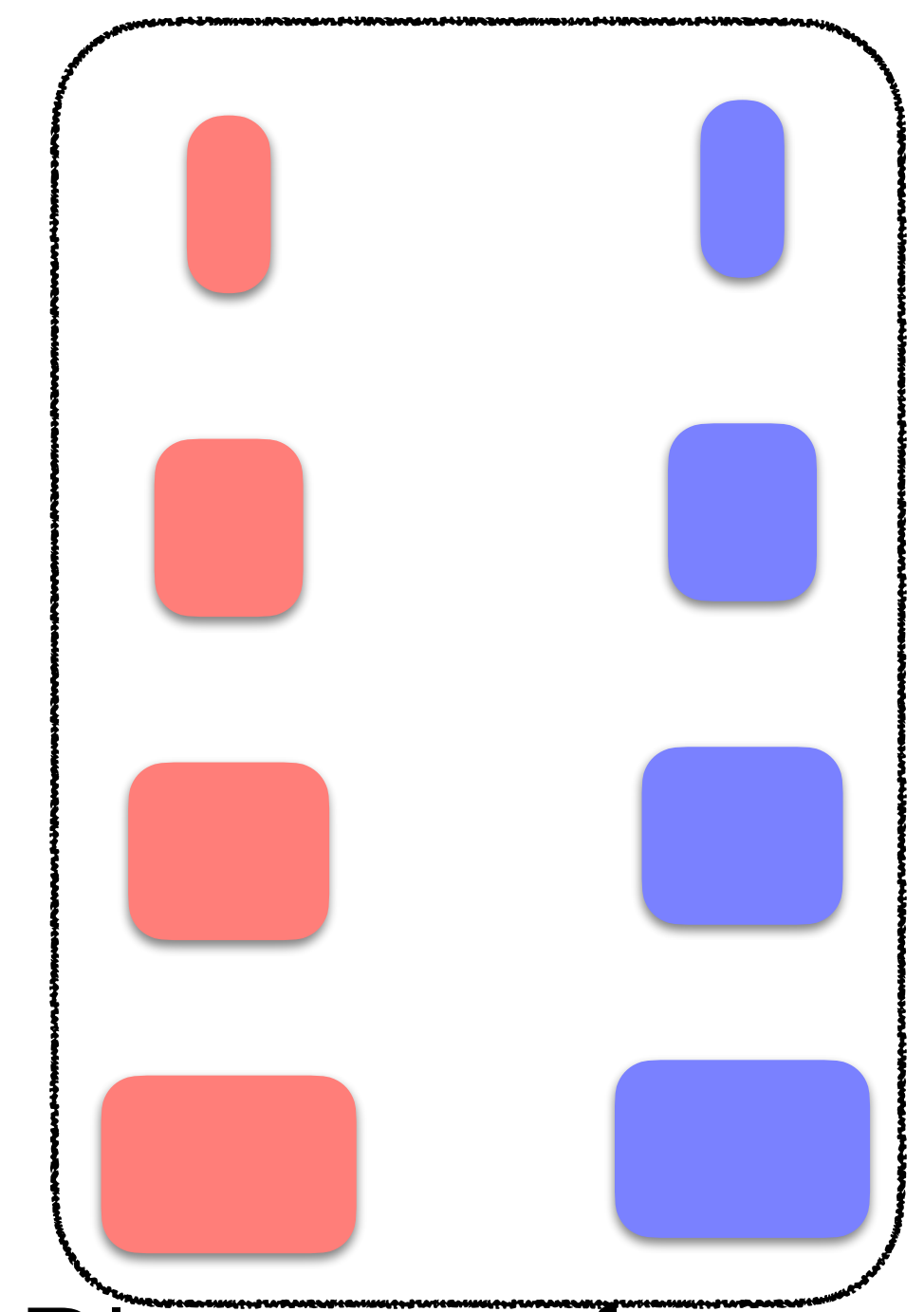
Block Cache



buffer



block cache



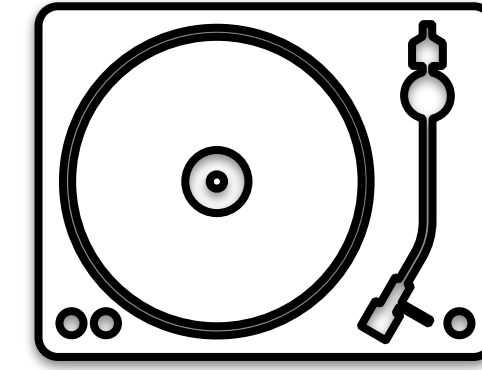
Bloom filters fence pointers


L1
L2
L3
L4




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

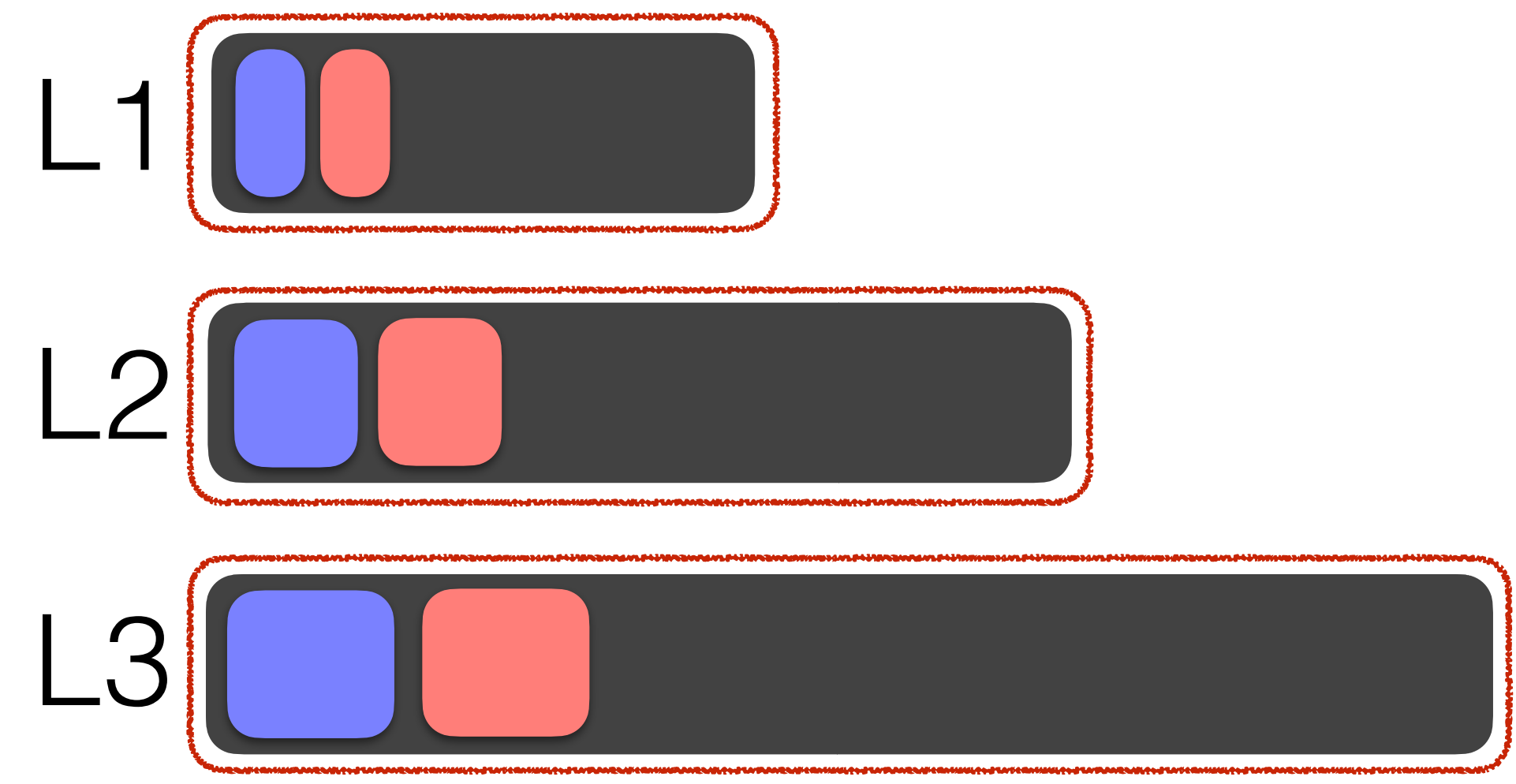
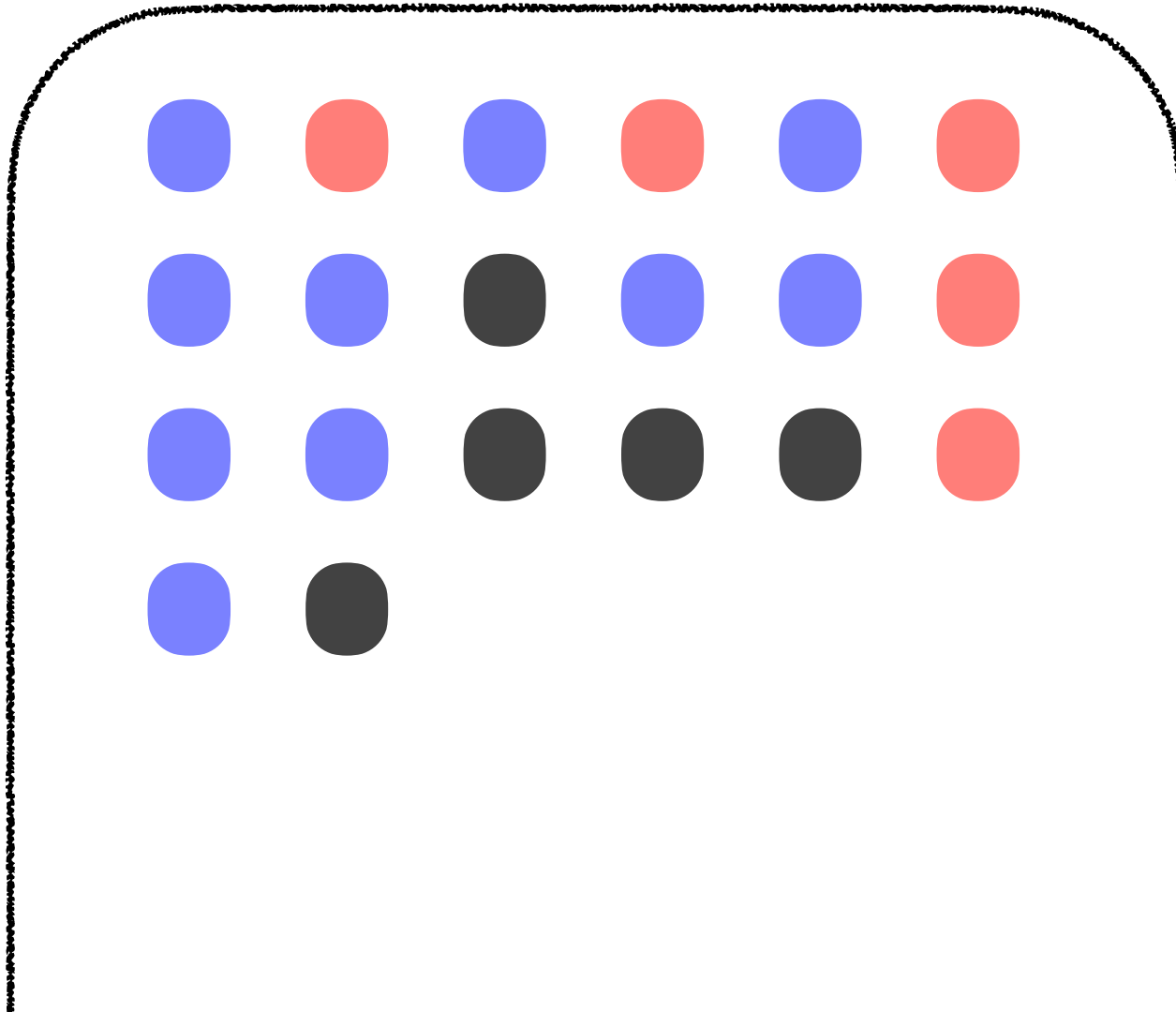
Block Cache




 Bloom filters


 fence pointers

get(7)



What about range queries?

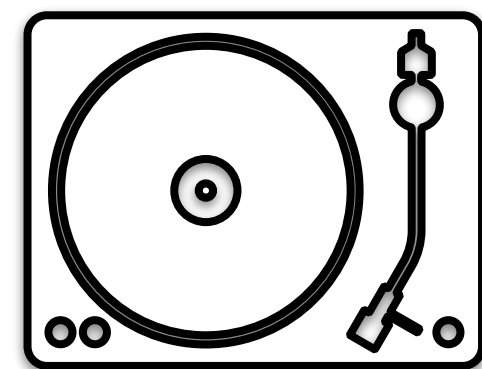
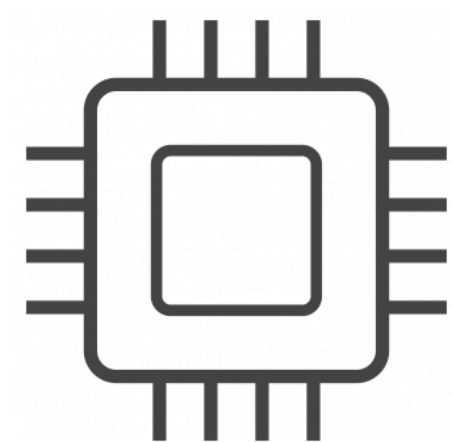
P : pages in buffer
 B : entries/page
 L : #levels
 T : size ratio
 N : #entries
 ϕ : FPR of BF

s : selectivity LRQ

Range Queries

Cost analysis

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



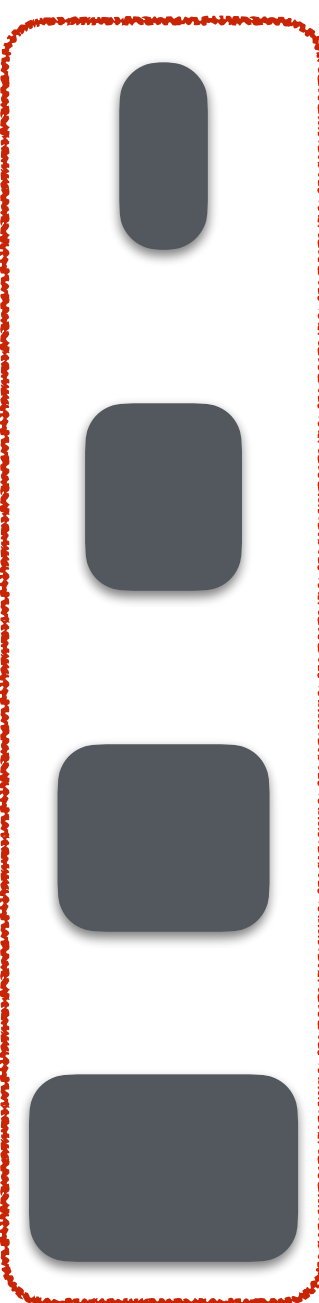
get(9,90)



buffer



Bloom filters



fence pointers

L1



L2



L3



L4



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

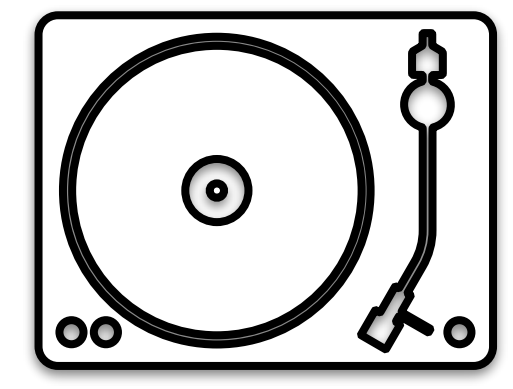
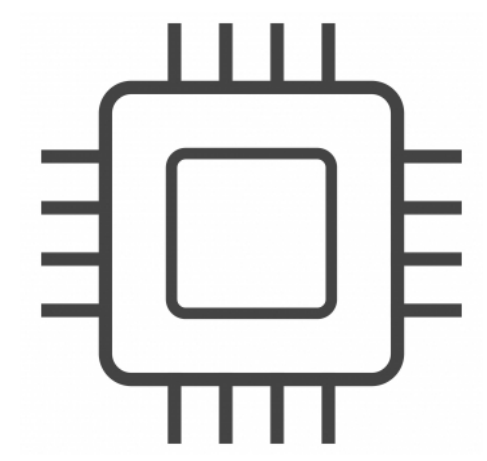
s : selectivity SRQ

Range Queries

Cost analysis

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

```
get(9, 15)
```



buffer



Bloom filters



fence pointers

L1

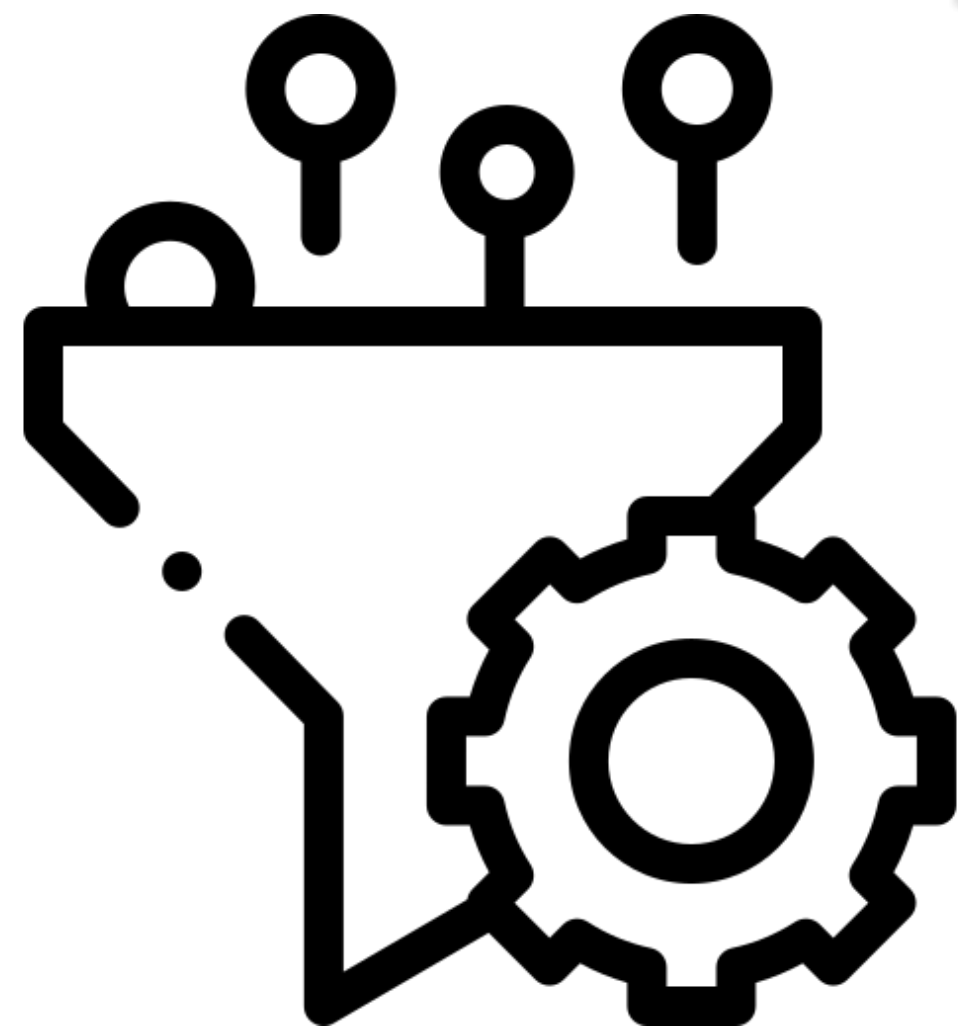
L2

L3

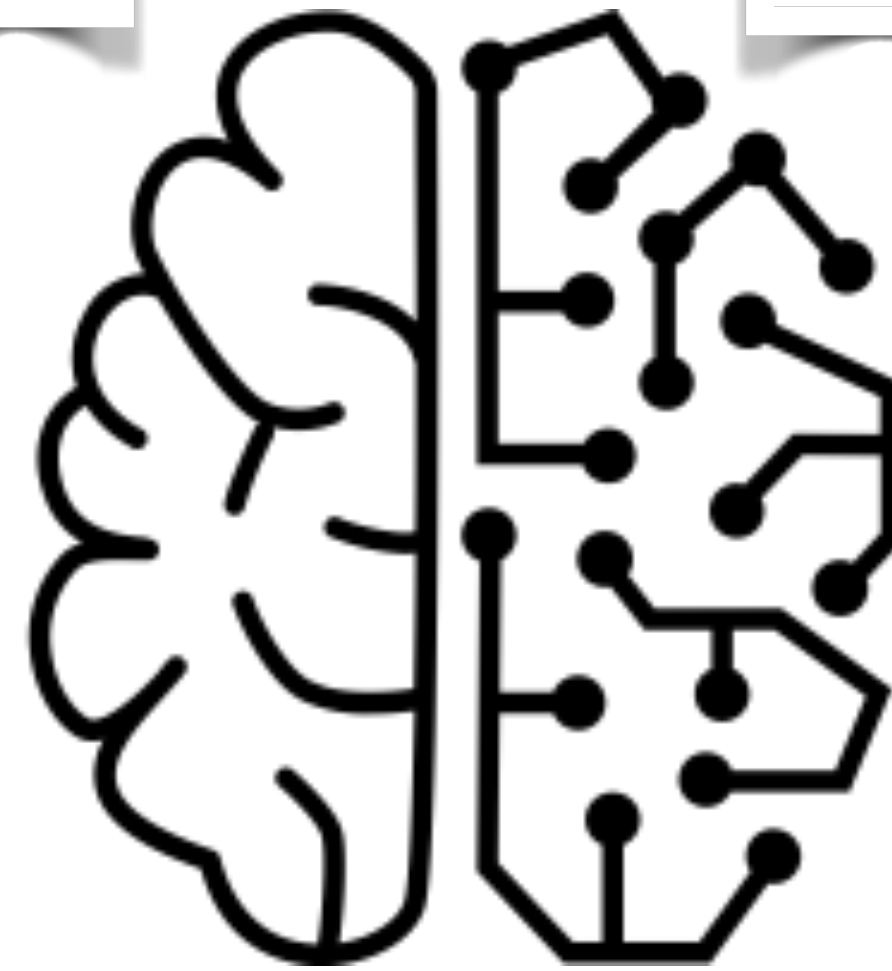
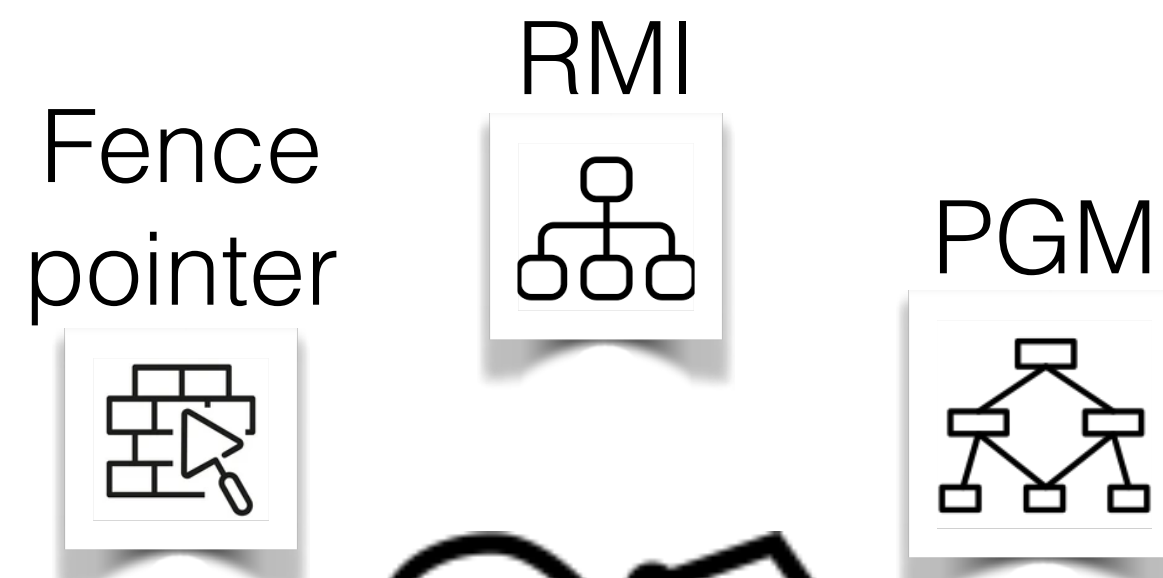
L4



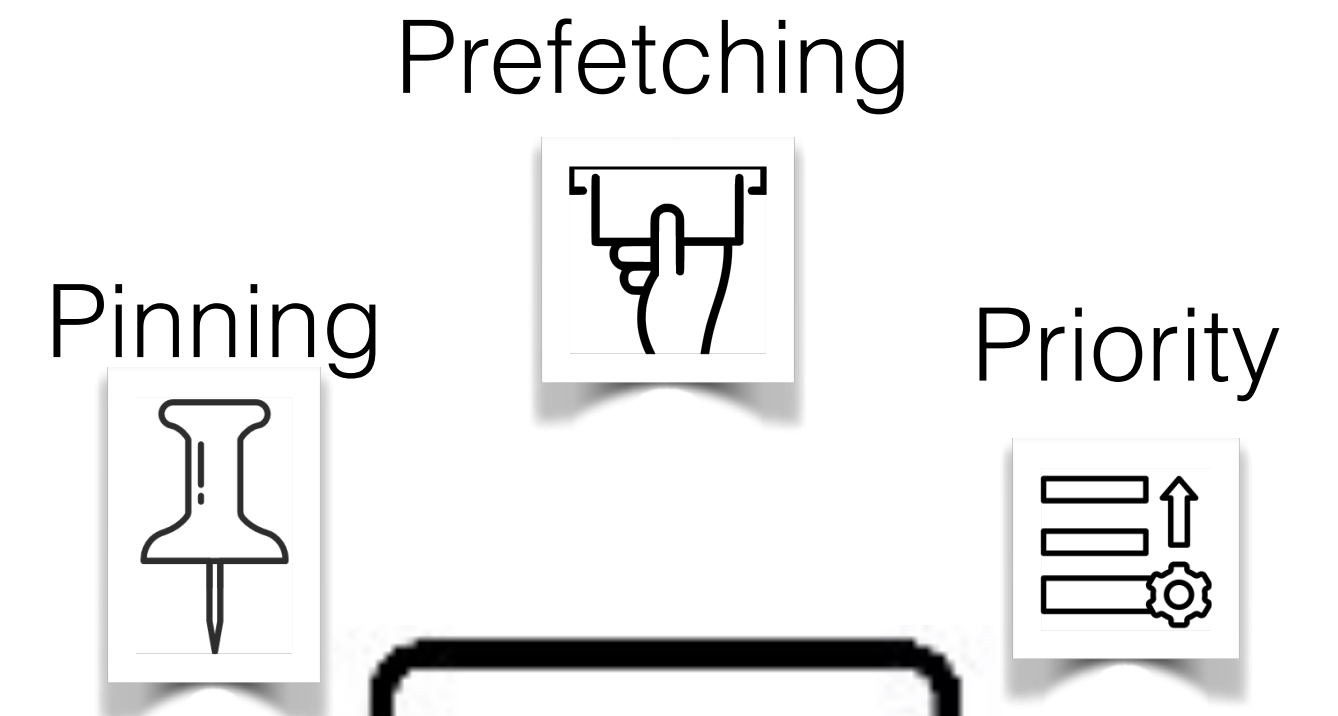
More Read Optimizations



filter structures



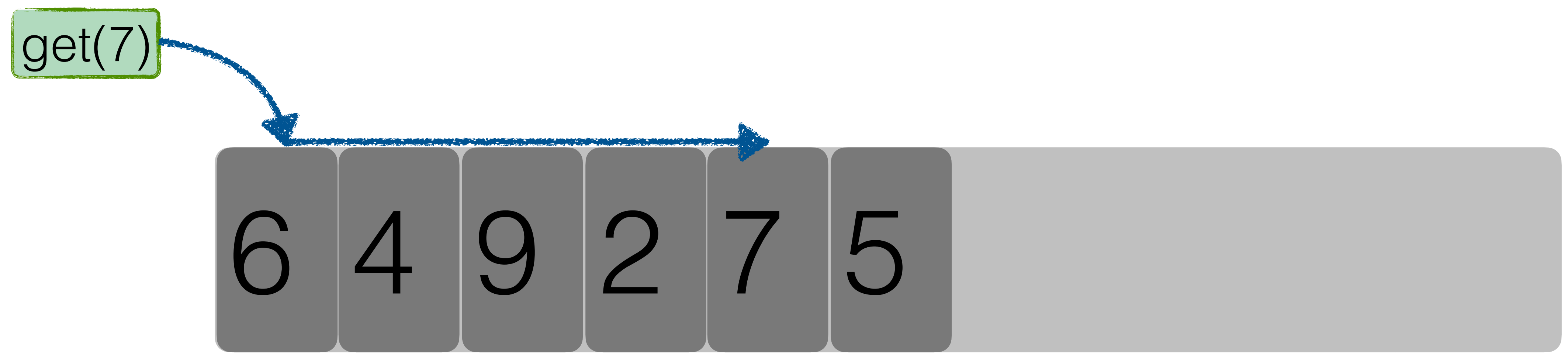
indexes



block cache

P : pages in buffer
 B : entries/page

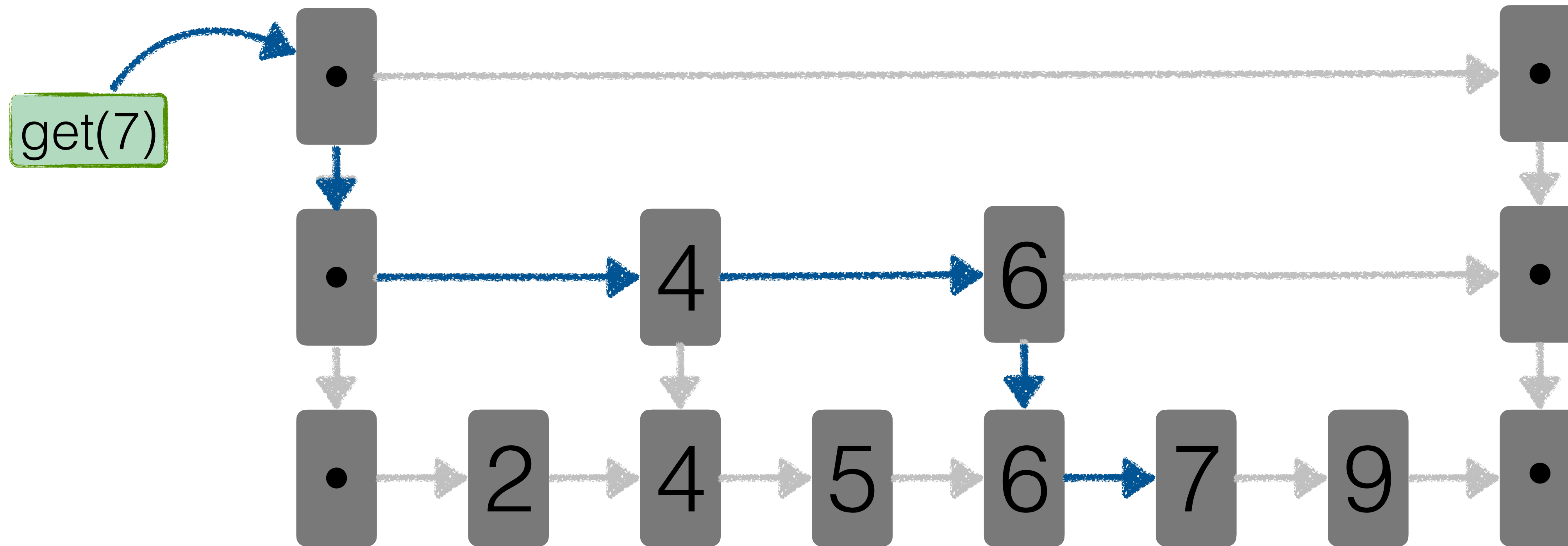
Buffer Implementation: **vector**



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

ingestion cost: $\mathcal{O}(1)$
space complexity: $\mathcal{O}(P \cdot B)$
point query cost: $\mathcal{O}(P \cdot B)$

Buffer Implementation: **skiplist**



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

P : pages in buffer
 B : entries/page

Buffer Implementation

vector

skiplist

hashmap

ingestion
cost

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

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

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

space
complexity

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

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

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

point query
cost

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

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

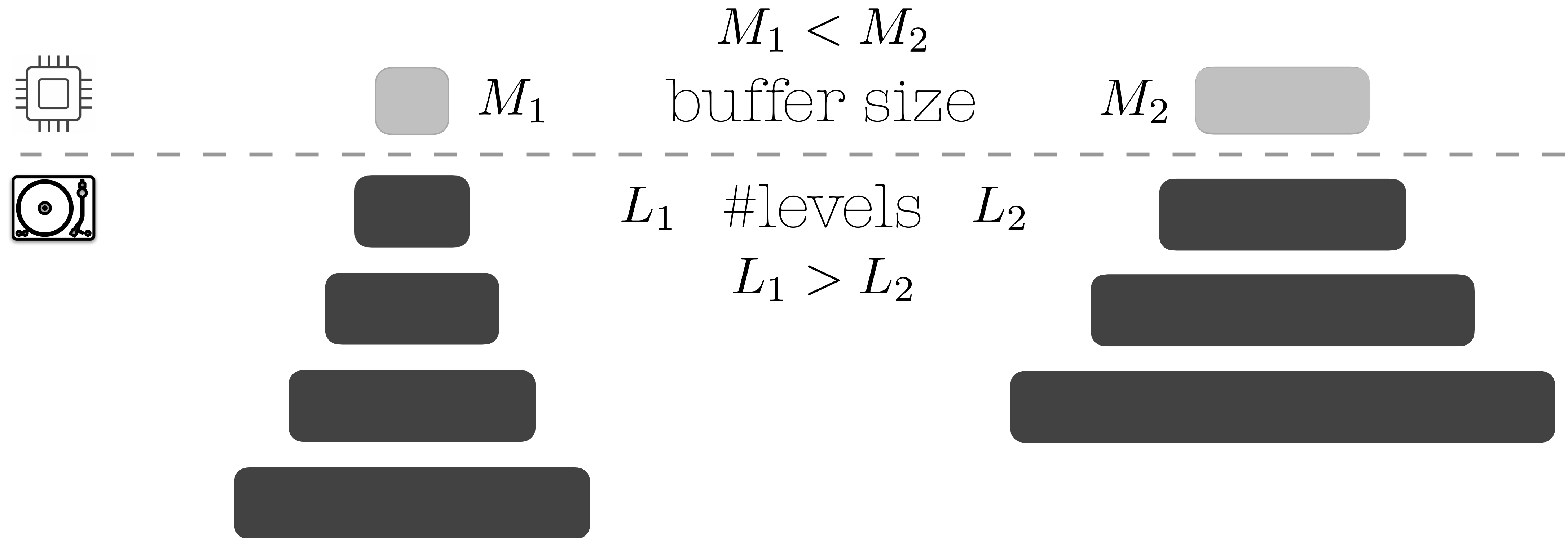
$$\mathcal{O}(1)$$

Ingestion-only
workloads

Mixed
workloads

I/O-bound
workloads

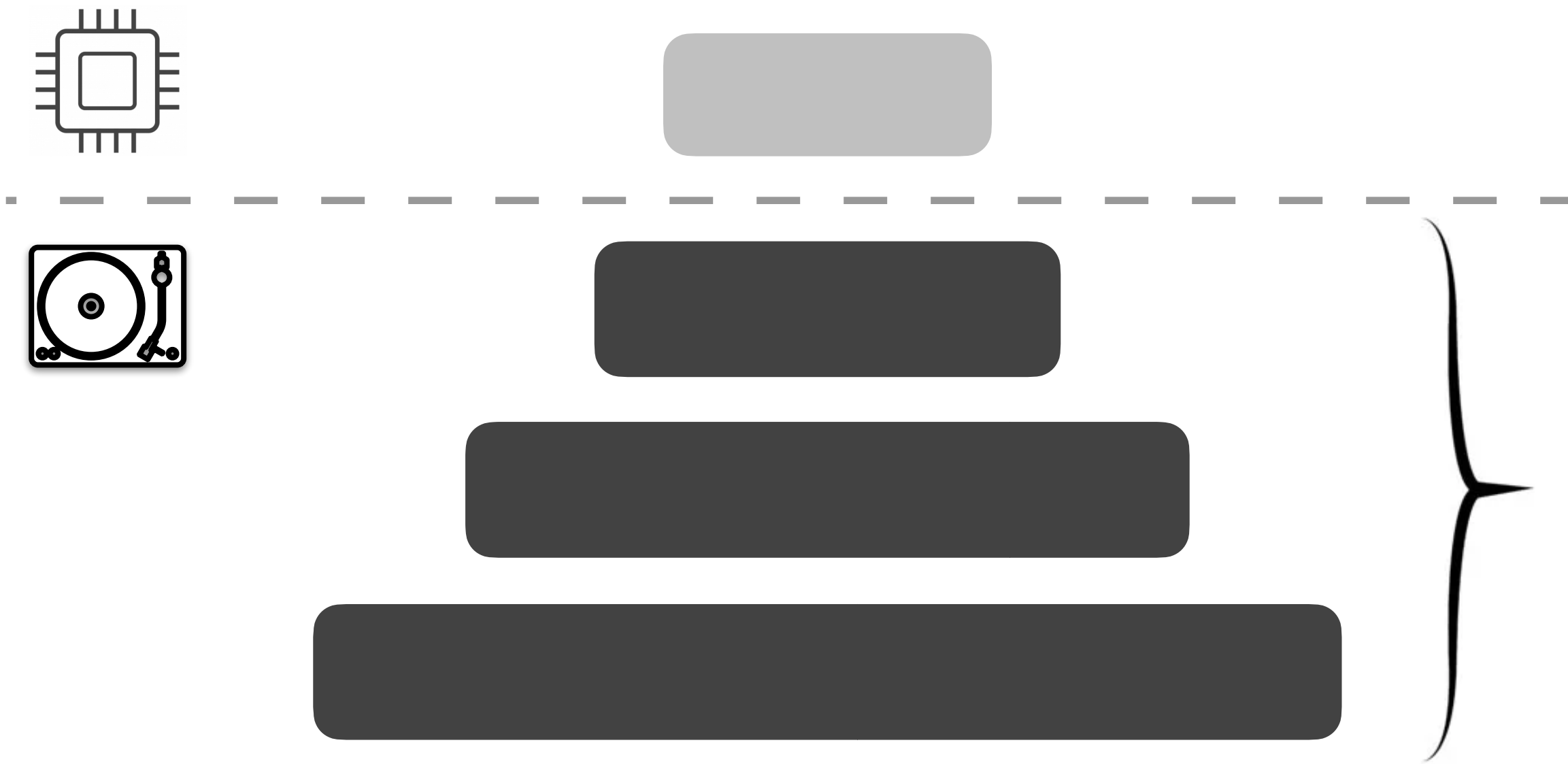
Size of the Buffer



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

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

L : #levels
 T : size ratio



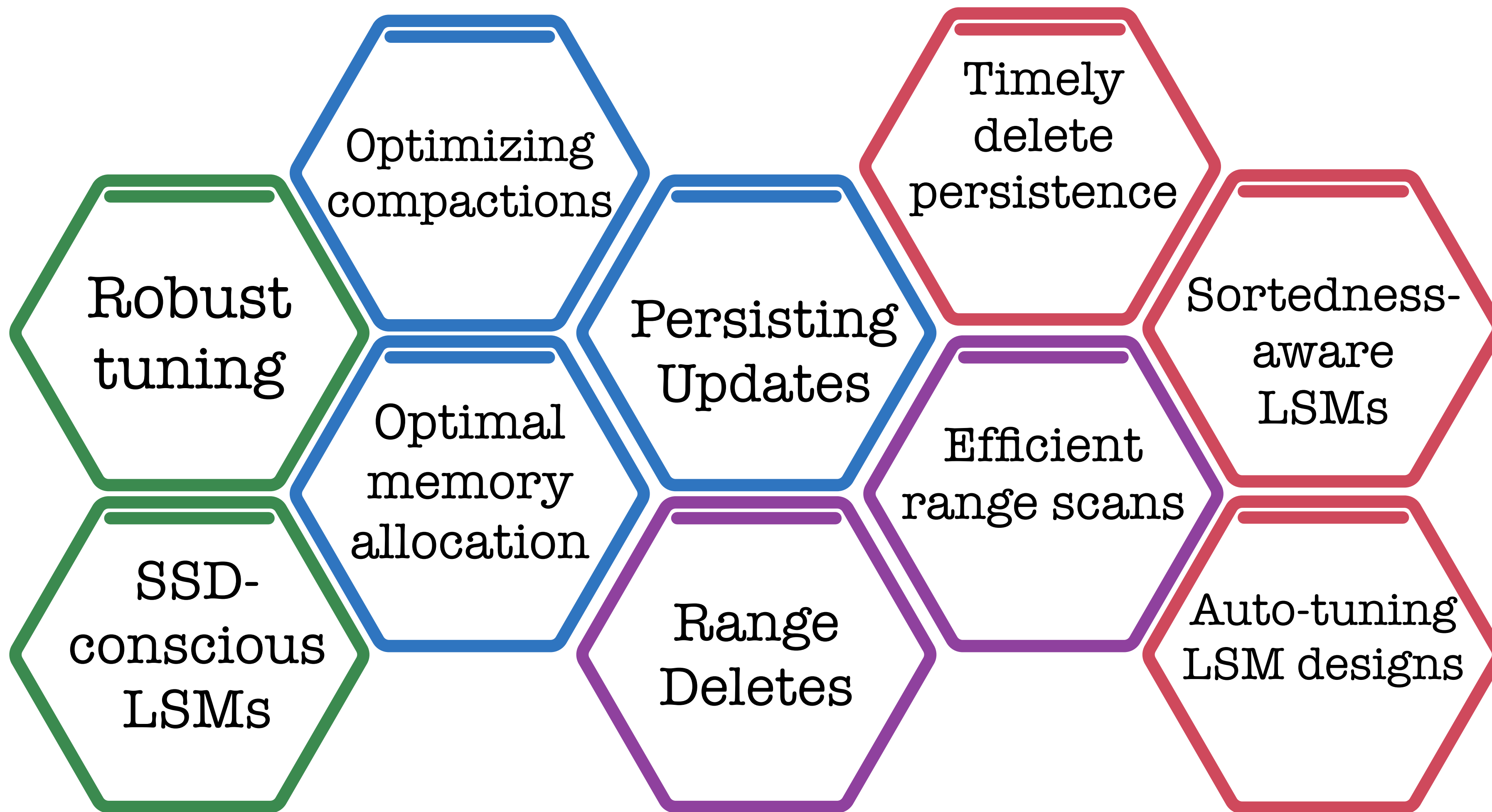
most data
on storage

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

99.9% on storage



How does the storage layer affect ingestion?



CS 561: Data Systems Architecture

Class 4

Systems & Research Projects

BOSTON
UNIVERSITY

Dr. Subhadeep Sarkar

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

