

# COSI 167A

# Advanced Data Systems

Class 7

# Class Projects

Prof. Subhadeep Sarkar

<https://ssd-brandeis.github.io/COSI-167A/>



# Class **logistics**

and administrivia

The **third technical question** is now available on the class website (due **before the class** on **Tue, Sep 24**).

**Project 1** is due in **~12 hours**.

**First guest lecture:** next **Tuesday** (**Sep 24**).

Register for **paper presentation!** (<https://shorturl.at/4P0IT>)



# Paper presentation

and discussion

Register for paper presentation! (<https://shorturl.at/4POIT>)

Brandeis

HOME

SCHEDULE

PROJECTS

PRESENTATION

RESOURCES

POLICIES



## Advanced Data Systems COSI 167A

SYLLABUS

GRADESCOPE

MOODLE



# Paper presentation

and discussion

Register for paper presentation! (<https://shorturl.at/4POIT>)

HOME

SCHEDULE

PROJECTS

PRESENTATION

RESOURCES

POLICIES



anced Data Systems

COSI 167A



# Paper presentation

and discussion

Register for paper presentation! (<https://shorturl.at/4POIT>)

**2 students** will be responsible for presenting the paper

learn the art of **technical presentation**, think as: **visualizing a review**

prepare slides **at least a week before** your presentation

get your **slides reviewed** by me **twice** before your final presentation

Pro tip: **Start EARLY!**



# Today in COSI 167A

What's on the cards?

**Cost analysis**

**Class projects!**



# Cost analysis

Counting all I/Os

## Ingestion cost

**expected #I/Os** performed to write a single entry to **disk**

note: if an entry is written **multiple times**, **count all I/Os**

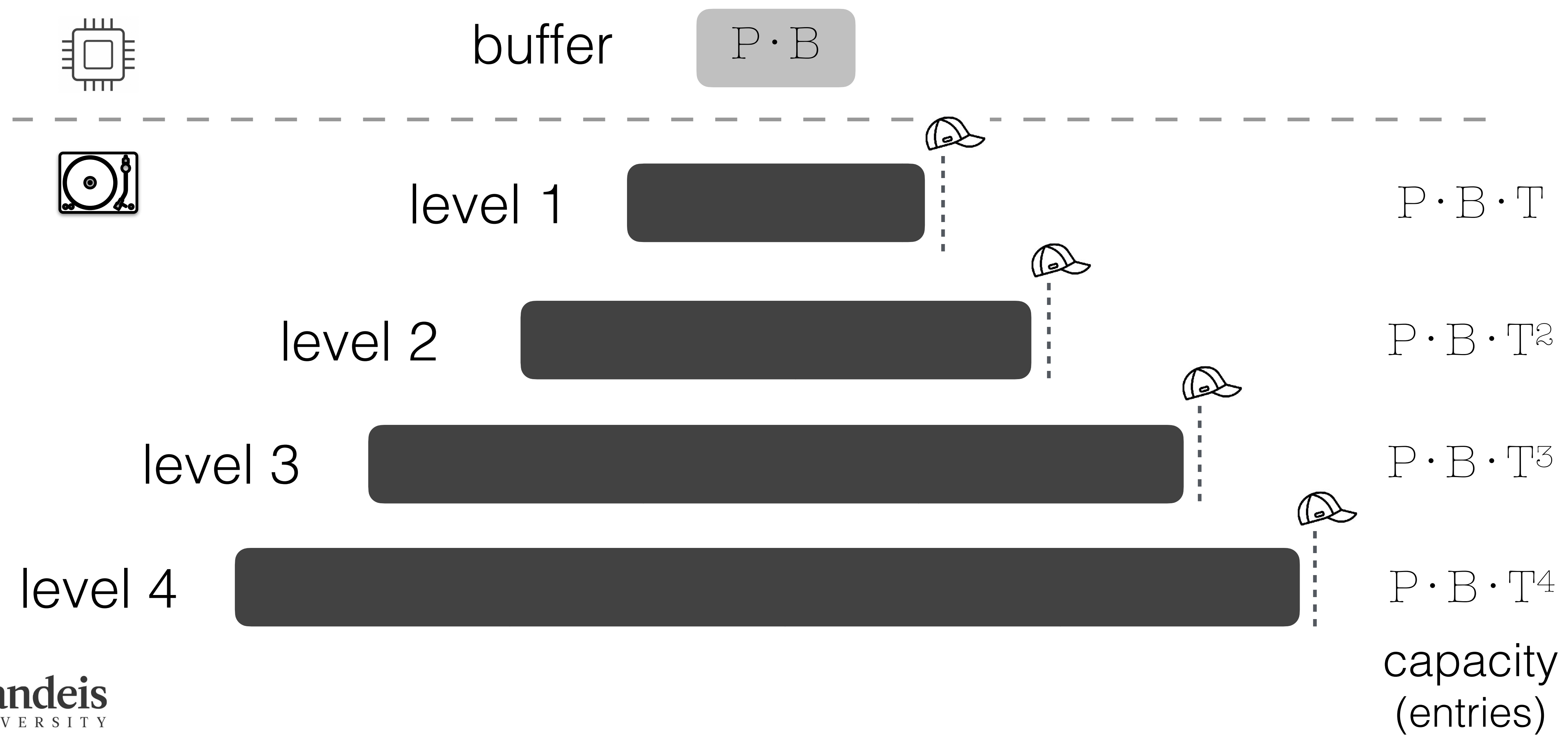
## Query cost

**expected #I/Os** performed to perform a single query

compare costs with and without **auxiliary** (helper) **data structures**



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



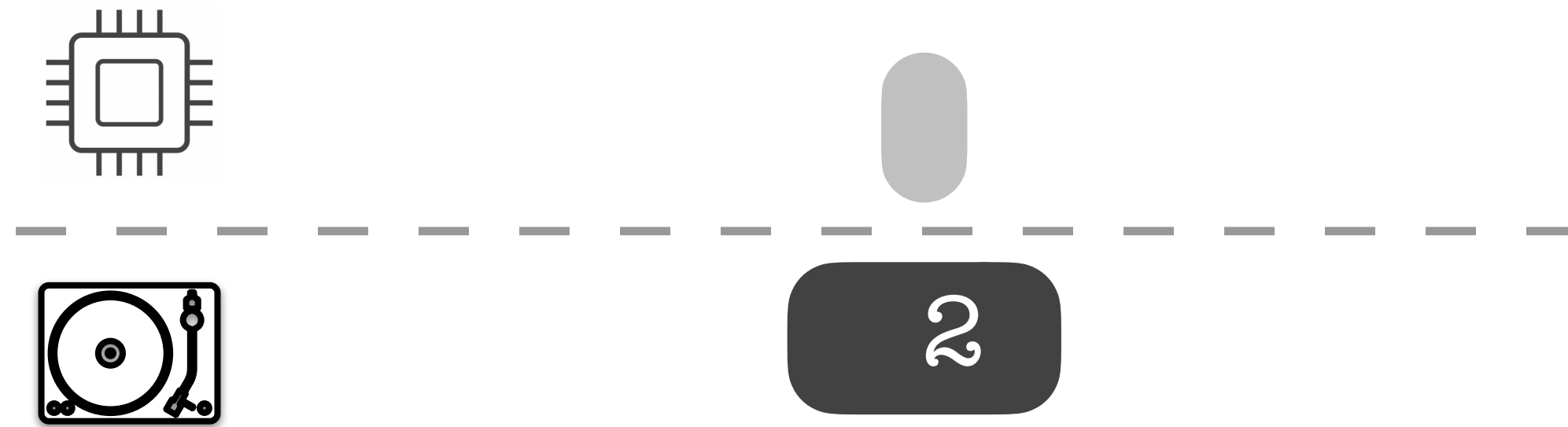


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Ingestion cost

Inserts and updates

leveled LSM-tree



in general, #times an entry is written to a level =  $T$

happens for all  $L$  levels on disk

Total #times an entry is written in the tree =  $L \cdot T$

$B$  entries are written to disk per I/O

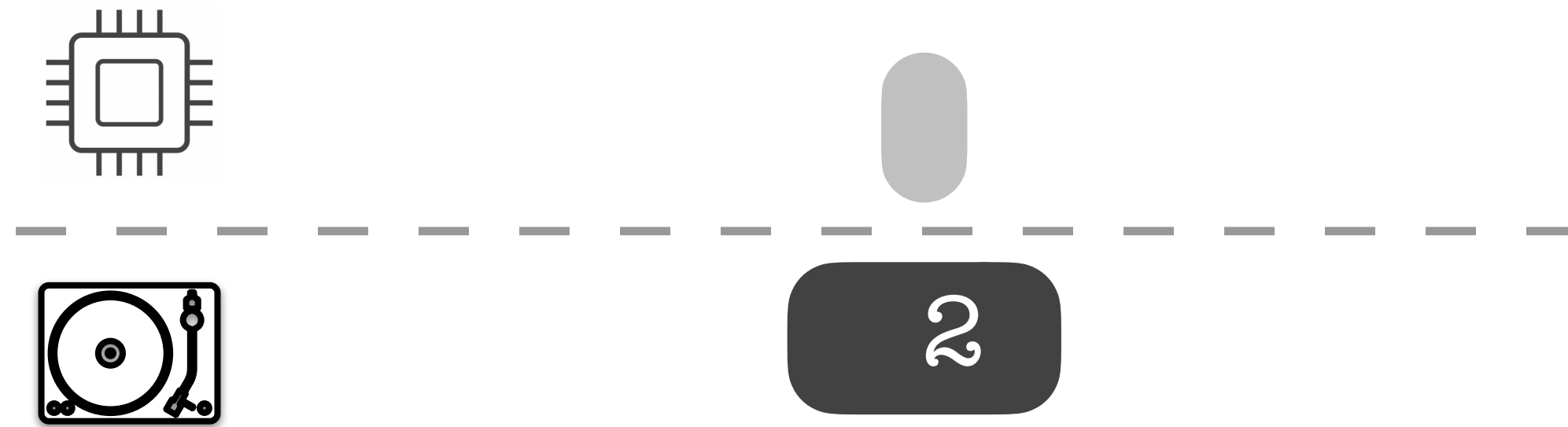


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Ingestion cost

Inserts and updates

leveled LSM-tree



in general, #times an entry is written to a level =  $T$

happens for all  $L$  levels on disk

Total #times an entry is written in the tree =  $L \cdot T$

$B$  entries are written to disk per I/O

Average number of I/Os to write a single entry =  $L \cdot T / B$



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$		
<b>Tiered LSM-tree</b>			
<b>B<sup>+</sup>-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			



# Cost analysis

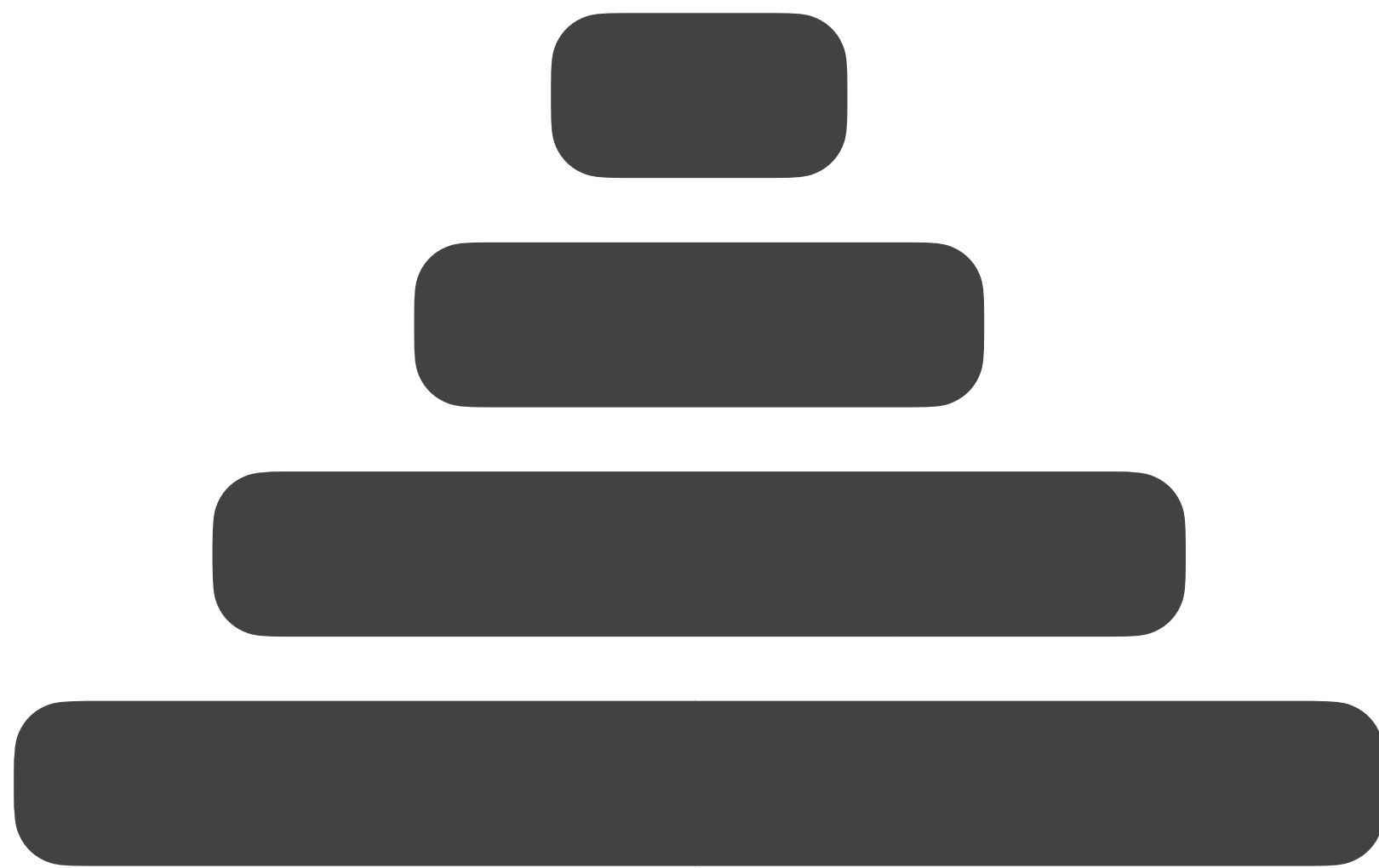
Counting all I/Os

data structure	ingestion cost	point lookup cost	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$		
<b>Tiered LSM-tree</b>	$O(L / B)$		
<b>B<sup>+</sup>-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

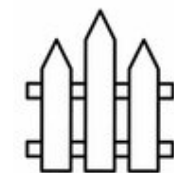


# Point lookup cost

Looking for a specific key

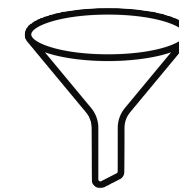


**leveled** LSM-tree



**fence  
pointers**

(page-wise zone map)



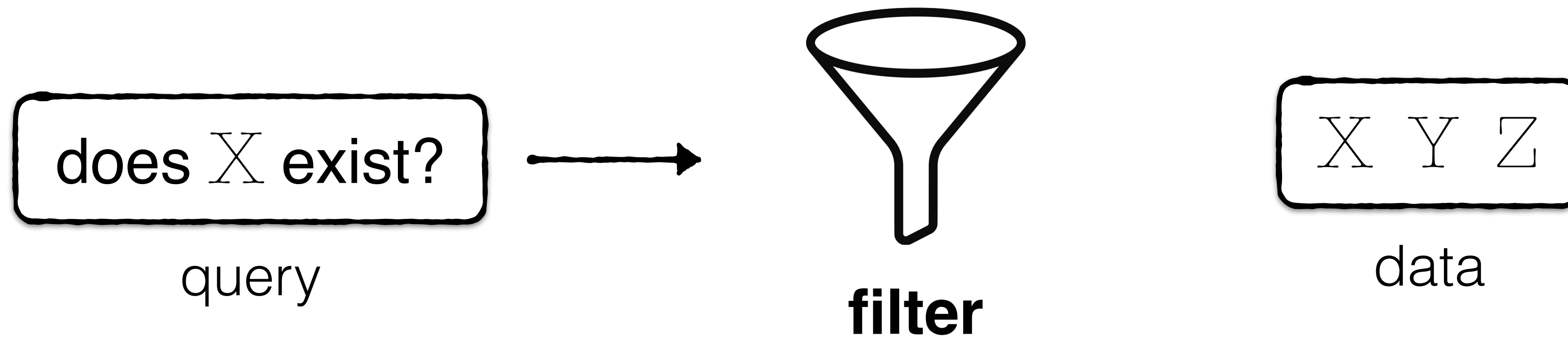
**filter**



# What is a **filter**?

Answers membership queries

answers  
**membership queries**



**no false negatives**



**false positives with tunable probability**

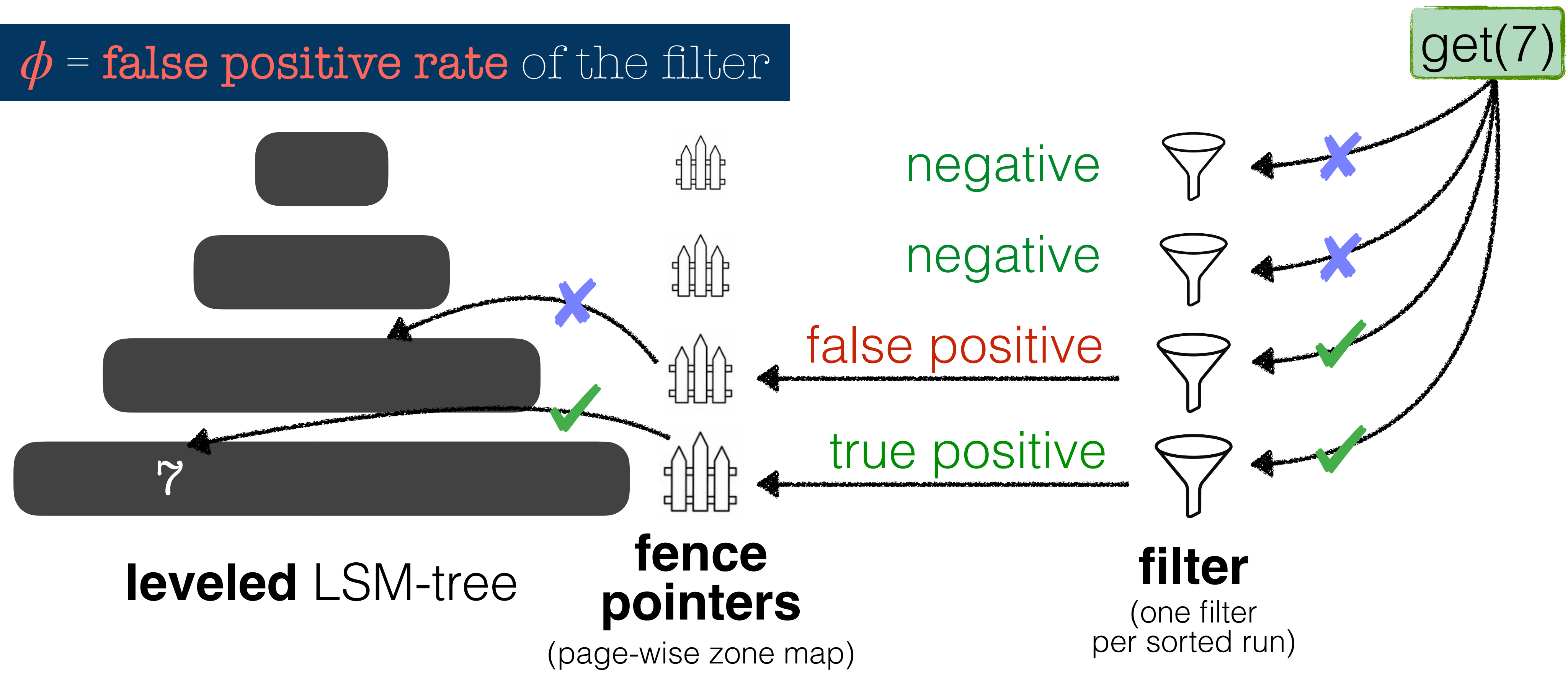


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Point lookup cost

Looking for a specific key

$\phi$  = false positive rate of the filter



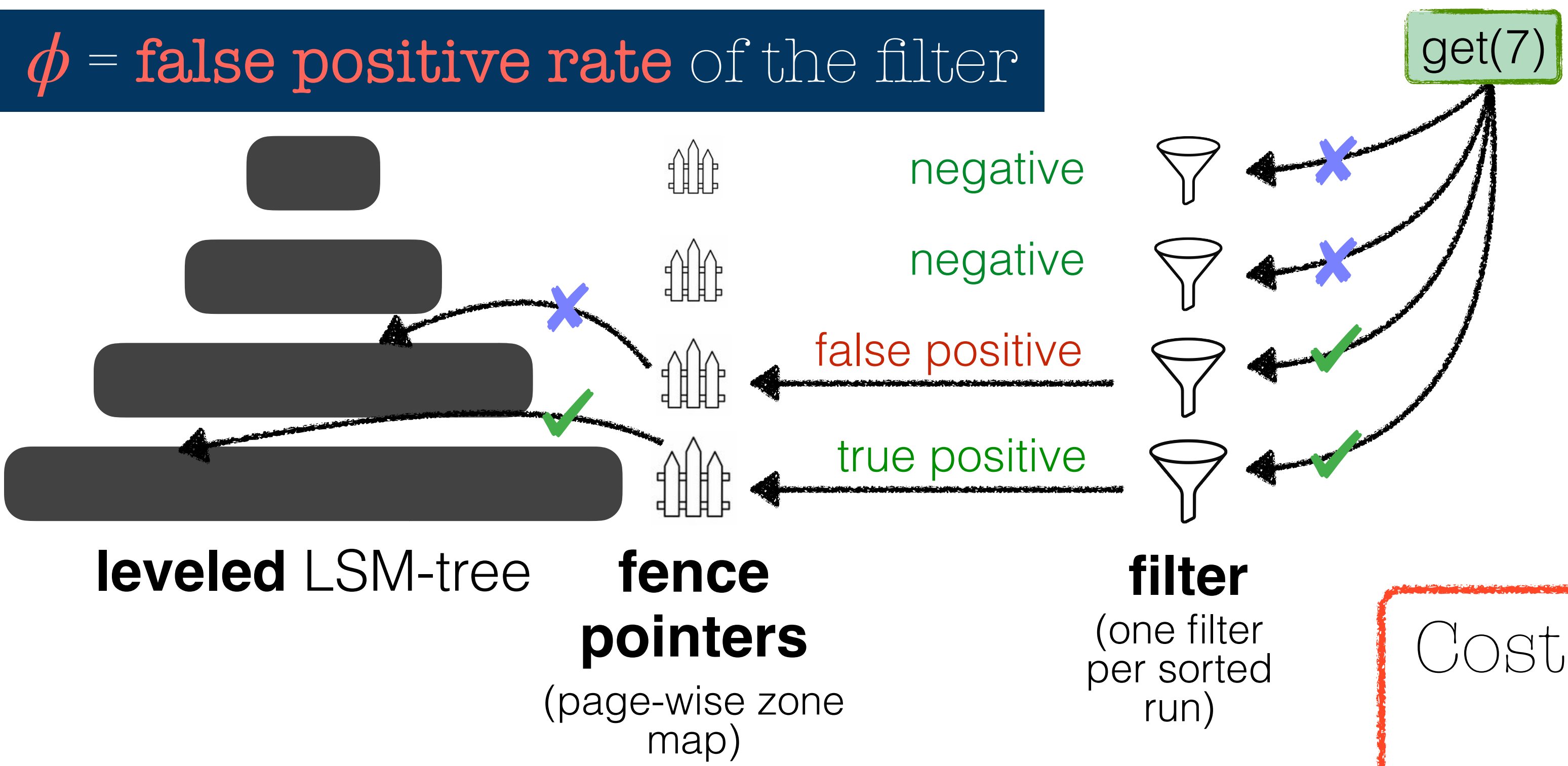


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Point lookup cost

Looking for a specific key

$\phi$  = false positive rate of the filter



1 I/O for the sorted run (level) containing the data

+  
1 I/O with probability  $\phi$  for all other sorted runs

1 sorted runs per level

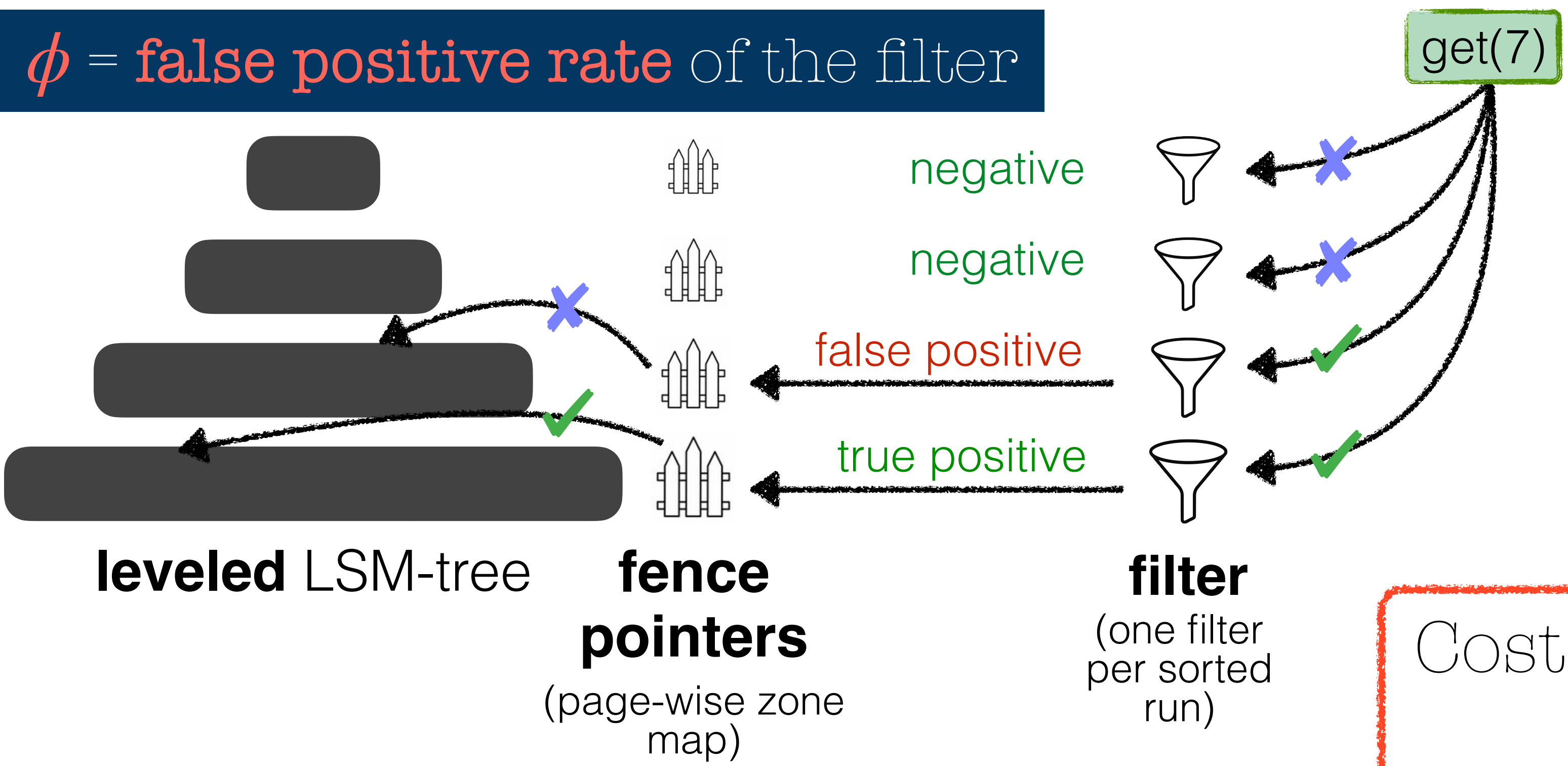
Cost of non-empty point lookup =  $1 + \phi \cdot (L - 1)$

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

# Point lookup cost

Looking for a specific key

$\phi$  = false positive rate of the filter



1 I/O for the sorted run (level) containing the data

+

1 I/O with probability  $\phi$  for all other sorted runs

1 sorted runs per level

Cost of non-empty point lookup =  $1 + \phi \cdot (L - 1)$

Cost of **empty** point lookup =  $\phi \cdot L$



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost*	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)$	
<b>Tiered LSM-tree</b>	$O(L / B)$		
<b>B+-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

\* with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost*	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)$	
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)$	
<b>B+-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

\* with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)



# Cost analysis

Counting all I/Os

$\phi = 0.008$  with 10 BPK

data structure	ingestion cost	point lookup cost*	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)$	
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)$	
<b>B+-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

**Monkey** takes this **another step further!**

# Cost analysis

Counting all I/Os

$\phi = 0.008$  with 10 BPK

data structure	ingestion cost	point lookup cost*	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	
<b>B<sup>+</sup>-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

\* with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)





# Range lookup cost

Looking for keys in a range

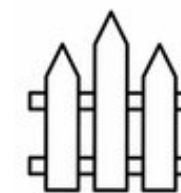
$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range

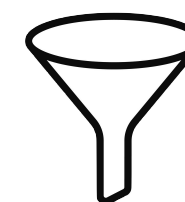
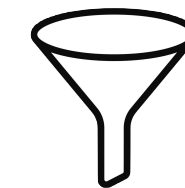


**leveled** LSM-tree



**fence  
pointers**

(page-wise zone map)



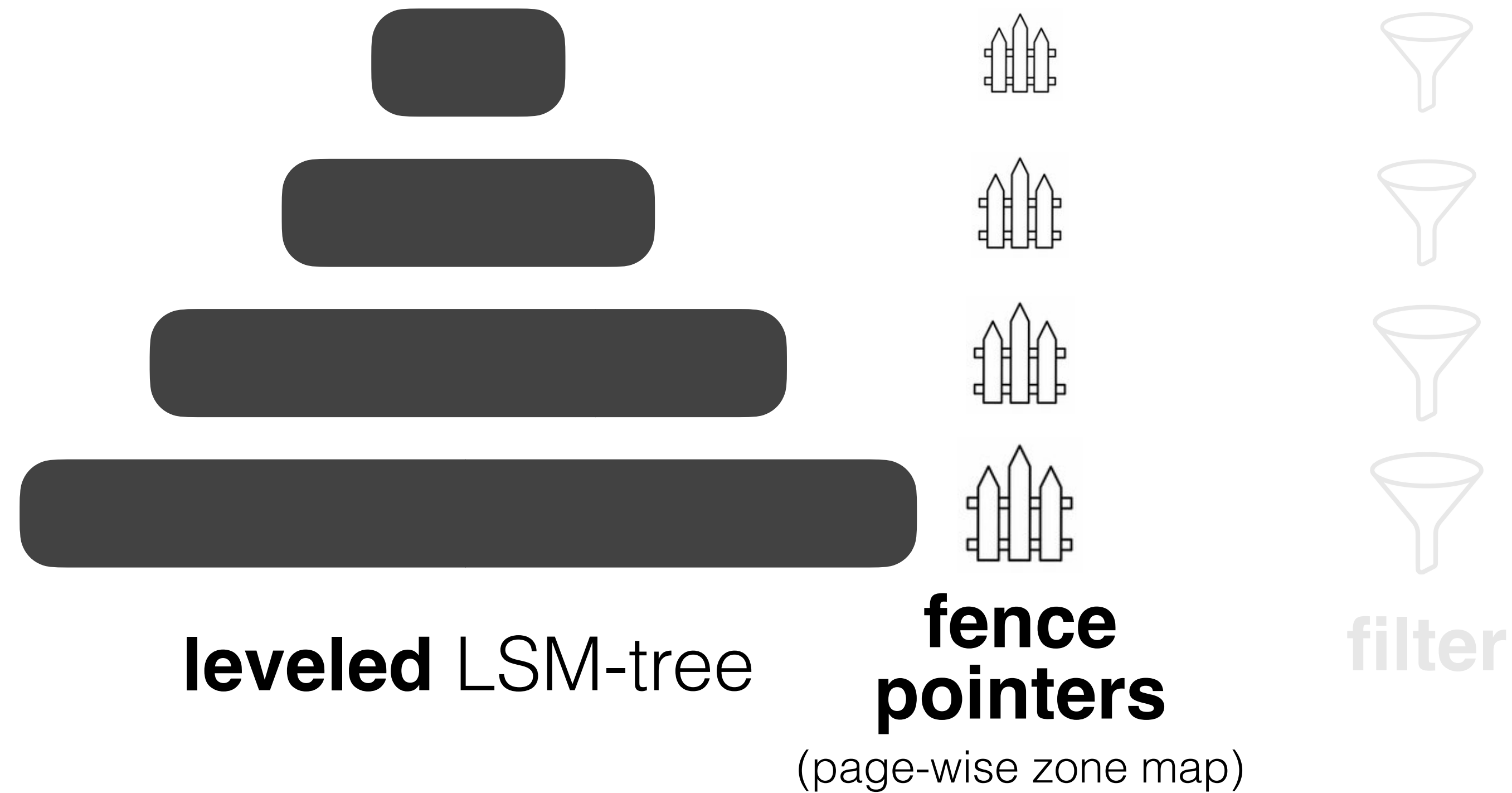
**filter**



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

# Range lookup cost

Looking for keys in a range

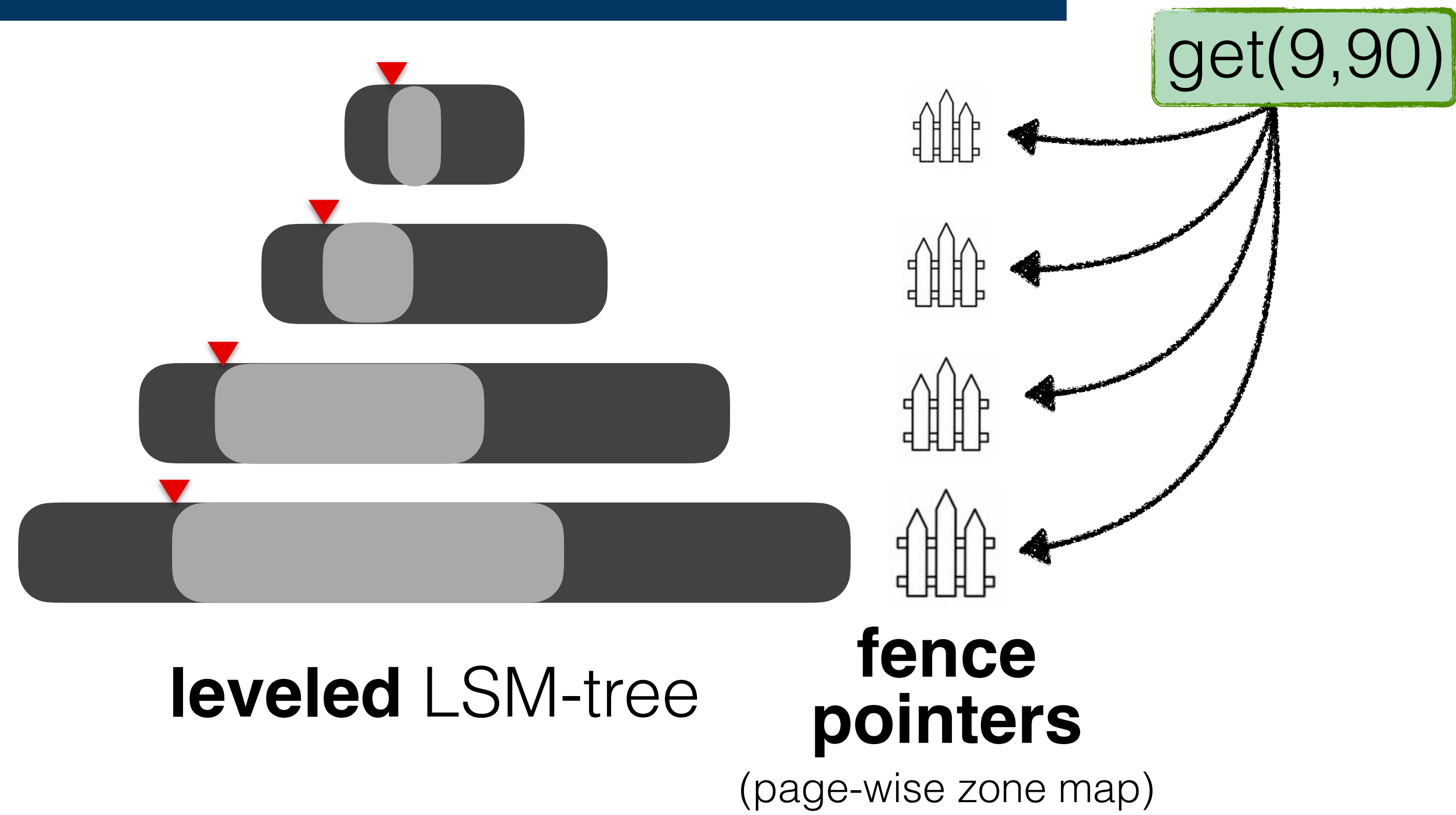


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range

$s$  = **selectivity** of the range query



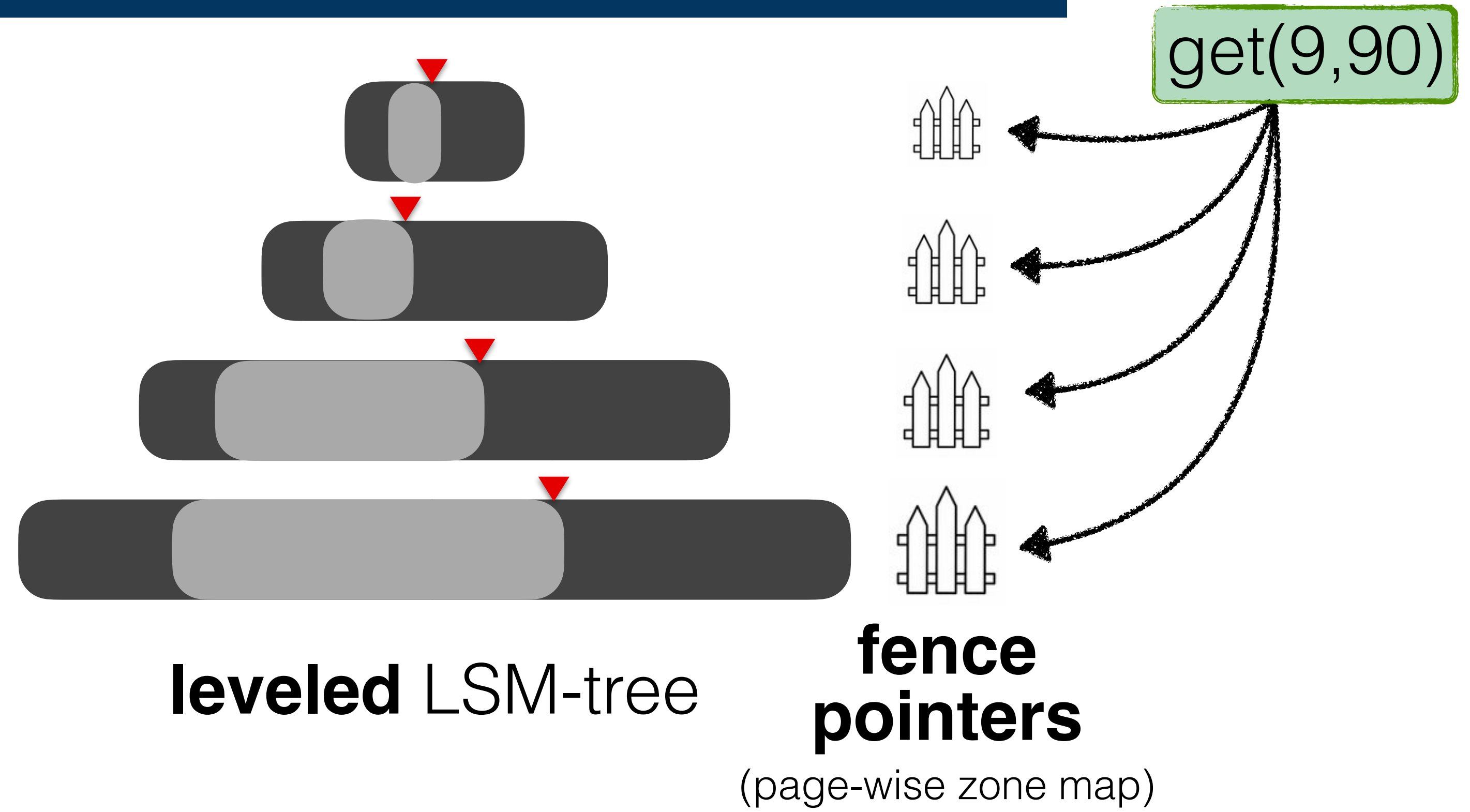


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range

$s$  = selectivity of the range query



---

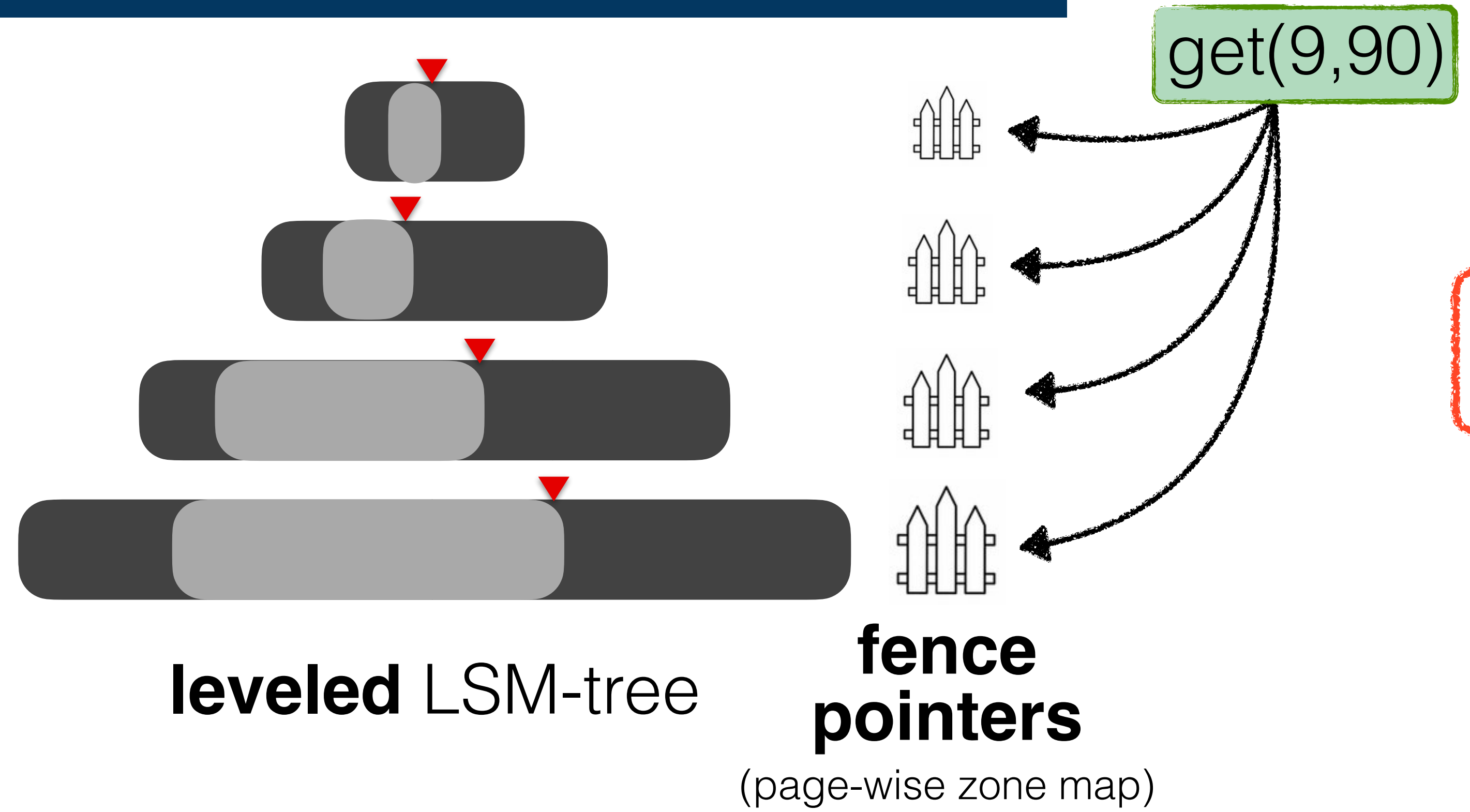
$$\begin{aligned} \text{total entries in tree} &= \mathbf{N} \\ \text{\#entries per page} &= \mathbf{B} \\ \hline \text{\#pages in tree} &= \mathbf{N/B} \end{aligned}$$

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

# Range lookup cost

Looking for keys in a range

$s$  = selectivity of the range query



total entries in tree =  $N$   
 #entries per page =  $B$   


---

 #pages in tree =  $N/B$

Cost of range lookup =  $s \cdot N/B$

What about the **range query cost** in a **tiered LSM-tree**?

same cost for  
 leveled & tiered LSM

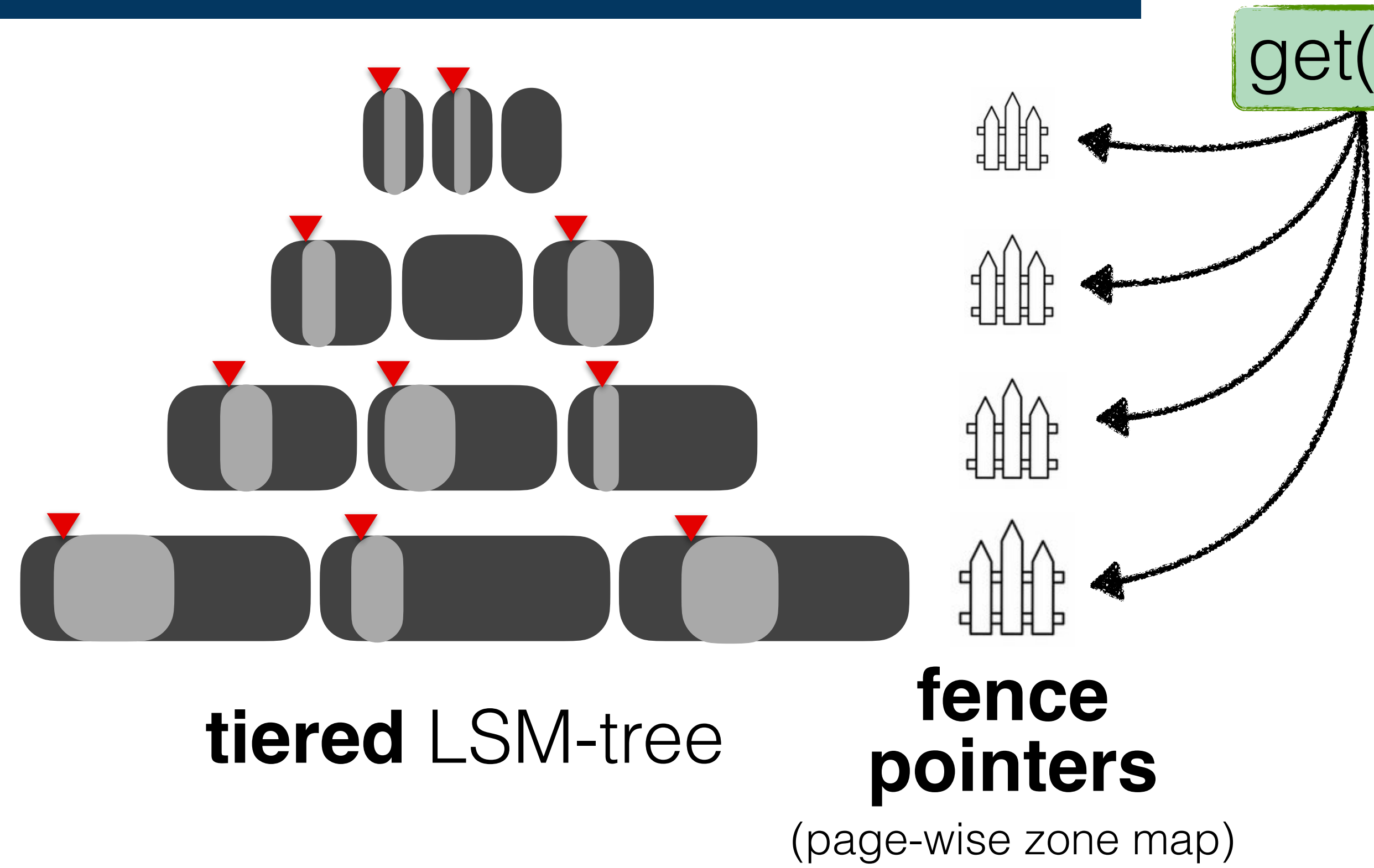


$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range

$s$  = selectivity of the range query



total entries in tree =  $N$   
 #entries per page =  $B$   


---

 #pages in tree =  $N/B$

Cost of range lookup =  $s \cdot N/B$

What about the **range query cost** in a **tiered LSM-tree**?



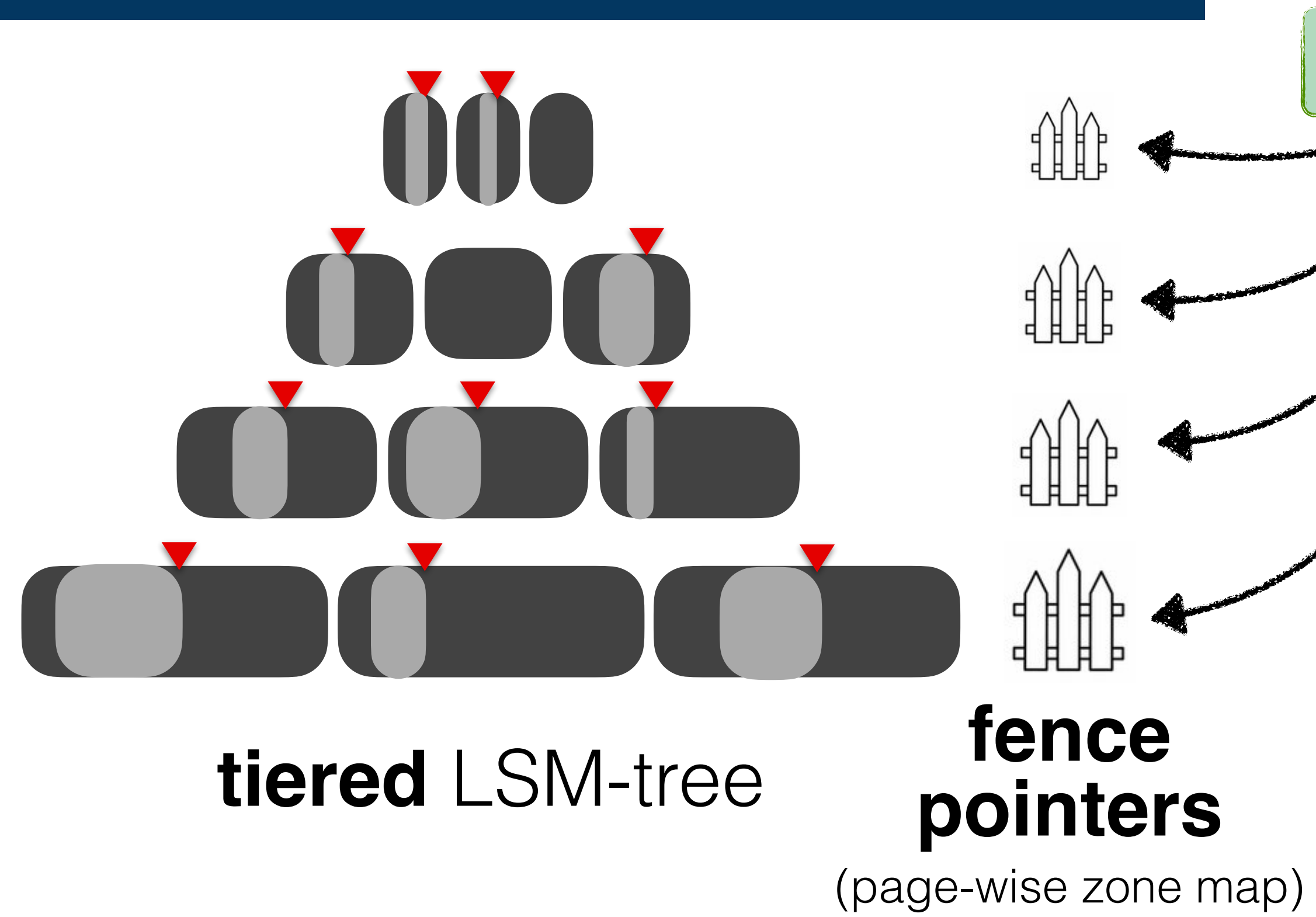
same cost for  
 leveled & tiered LSM

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

# Range lookup cost

Looking for keys in a range

$s$  = **selectivity** of the range query



total entries in tree =  $N$   
 #entries per page =  $B$   


---

 #pages in tree =  $N/B$

Cost of range lookup =  $s \cdot N/B$

What about the **range query cost** in a **tiered LSM-tree**?

same cost for  
 tiered & partitioned LSM

Very different story when workload has **update/deletes**.

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$
<b>B<sup>+</sup>-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

\* with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)





# Cost analysis

Counting all I/Os

\* **long** range lookups

data structure	ingestion cost	point lookup cost*	range lookup cost*
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$
<b>B<sup>+</sup>-tree</b>			
<b>Sorted array</b>			
<b>Log</b>			

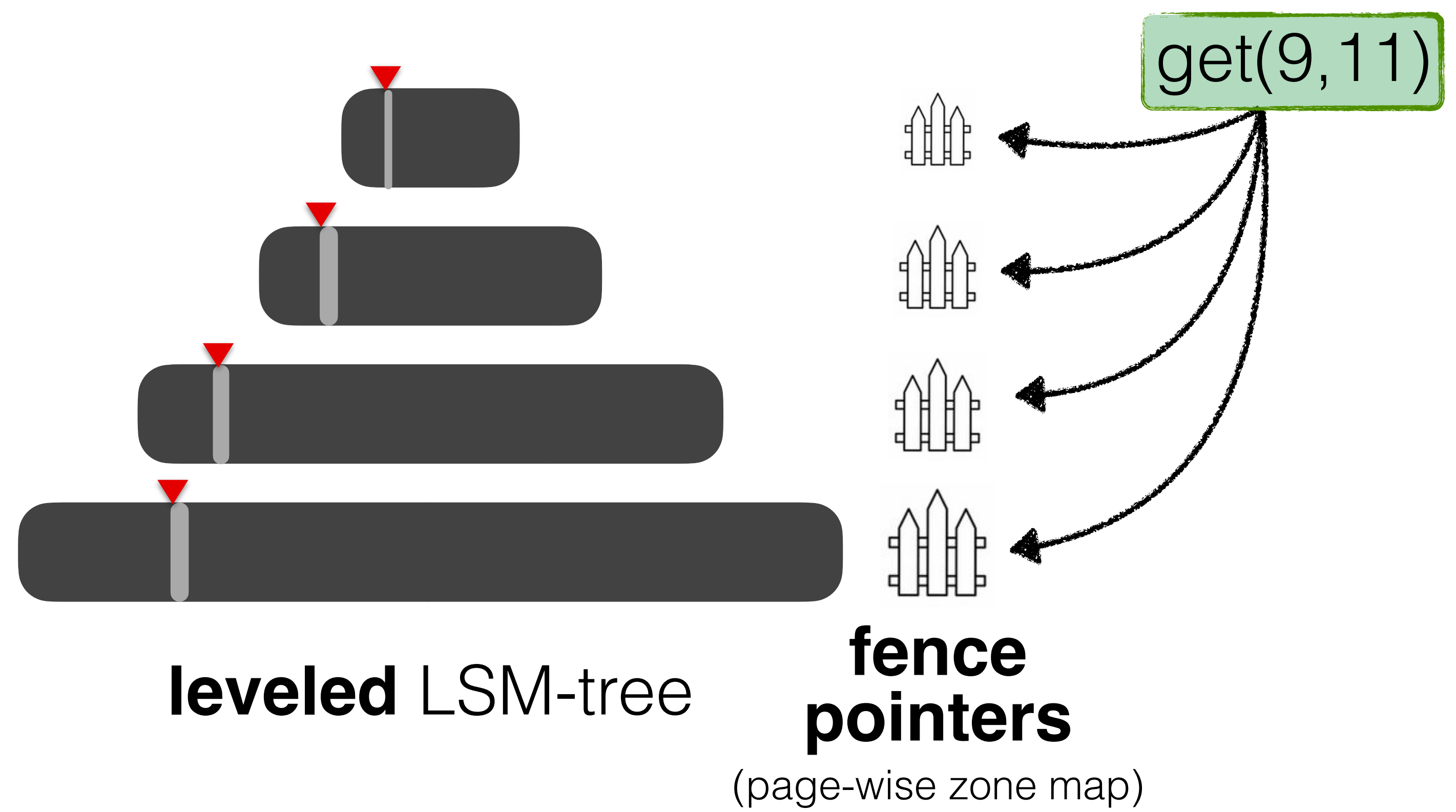
\* with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)



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

# Range lookup cost

Looking for keys in a range



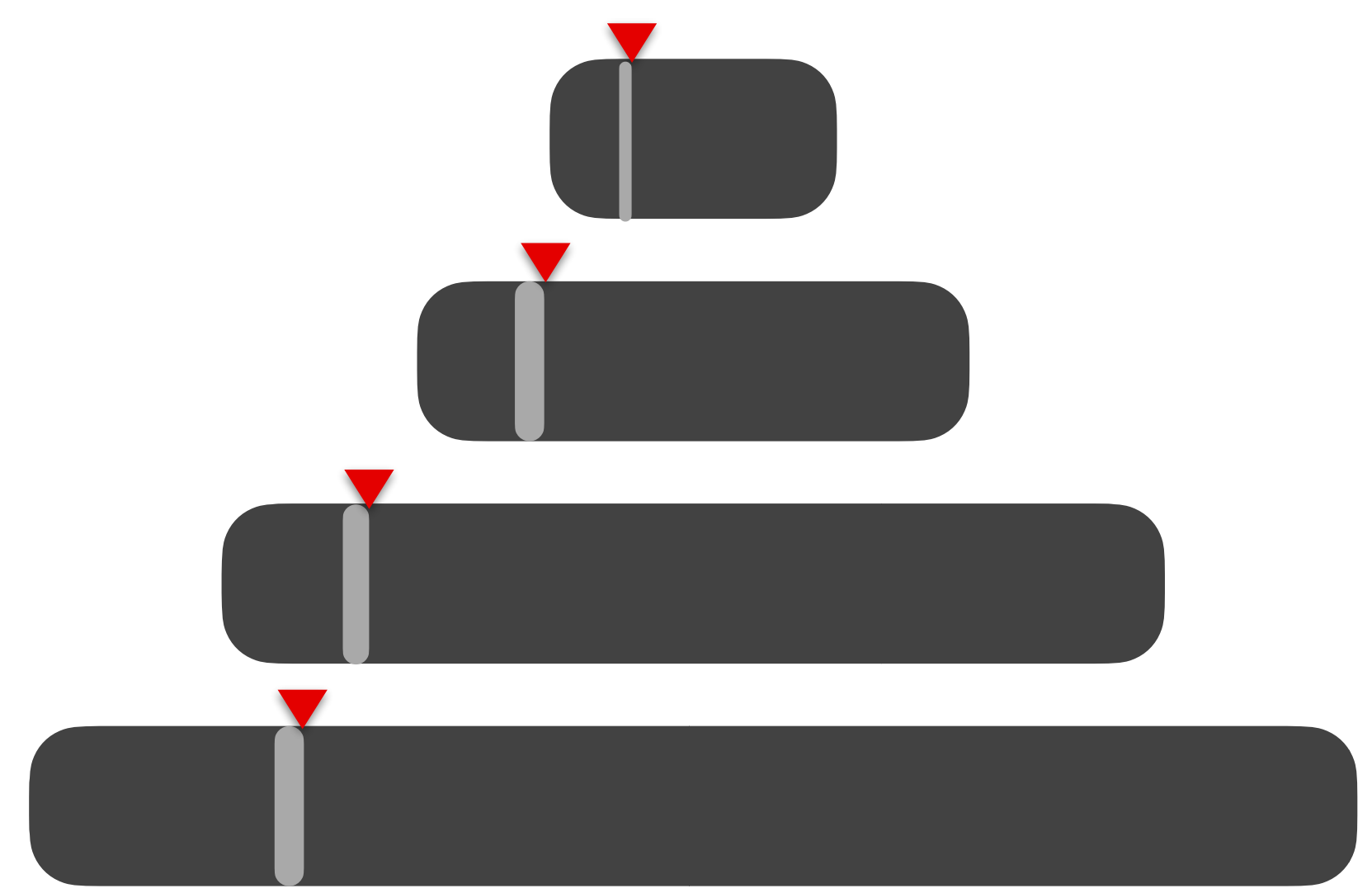
What if the **range query** has a **very low selectivity**?



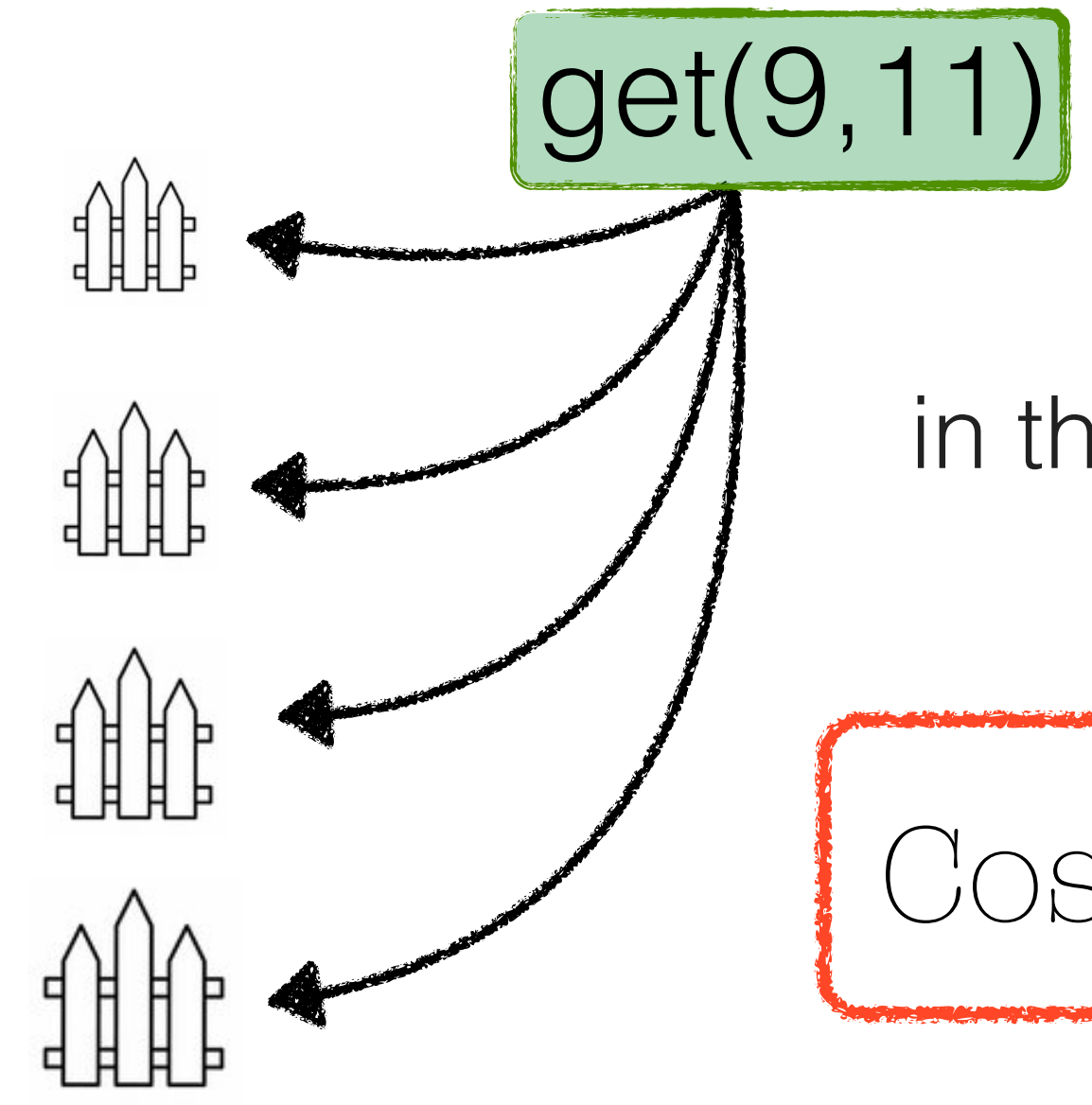
$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range



leveled LSM-tree



fence pointers  
(page-wise zone map)

What if the **range query** has a **very low selectivity**?



in the **worst case**,  
need to read **2 pages per sorted run**

Cost of short range lookup =  $2 \cdot L$

What about the **tiering**?

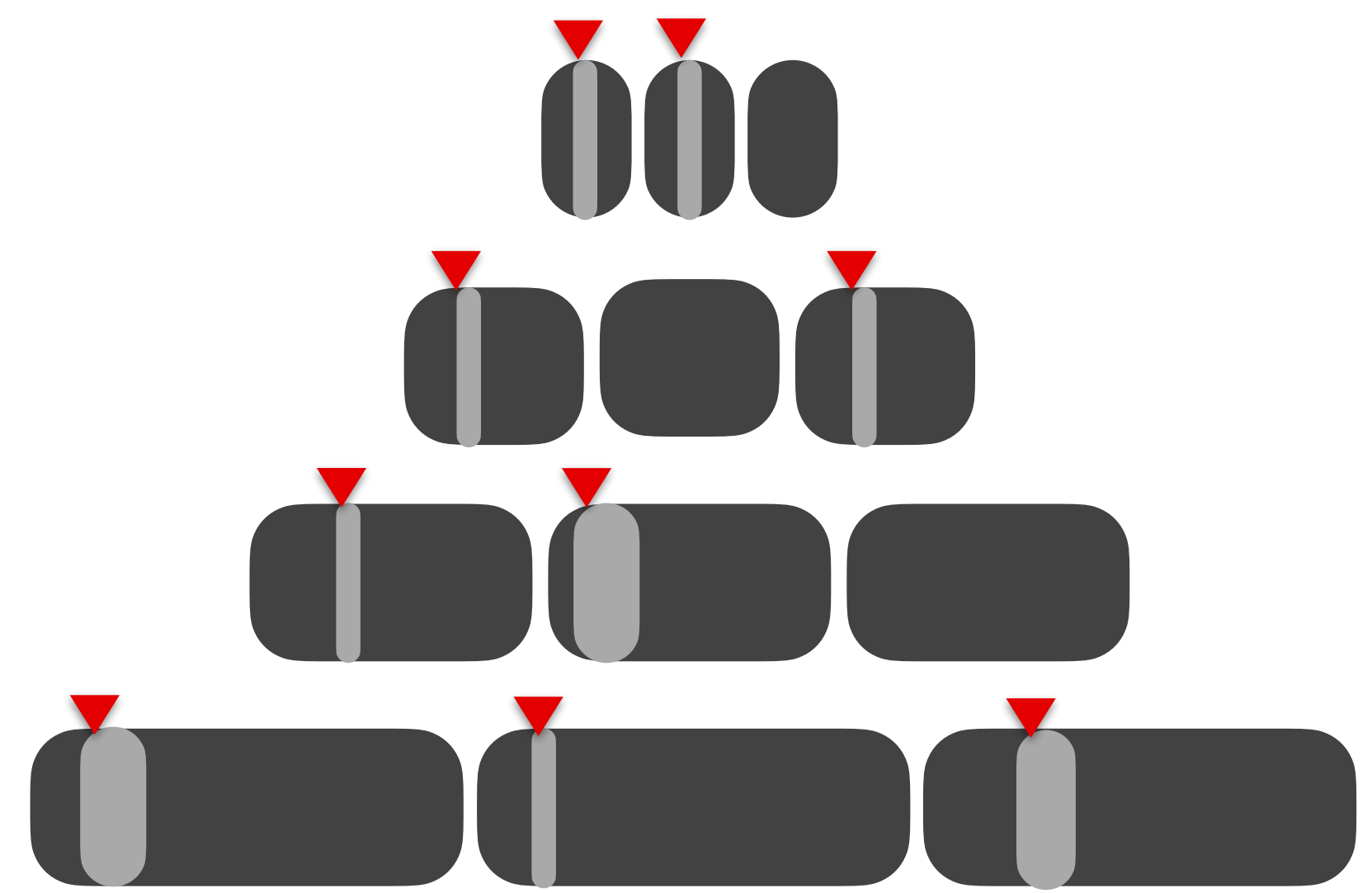




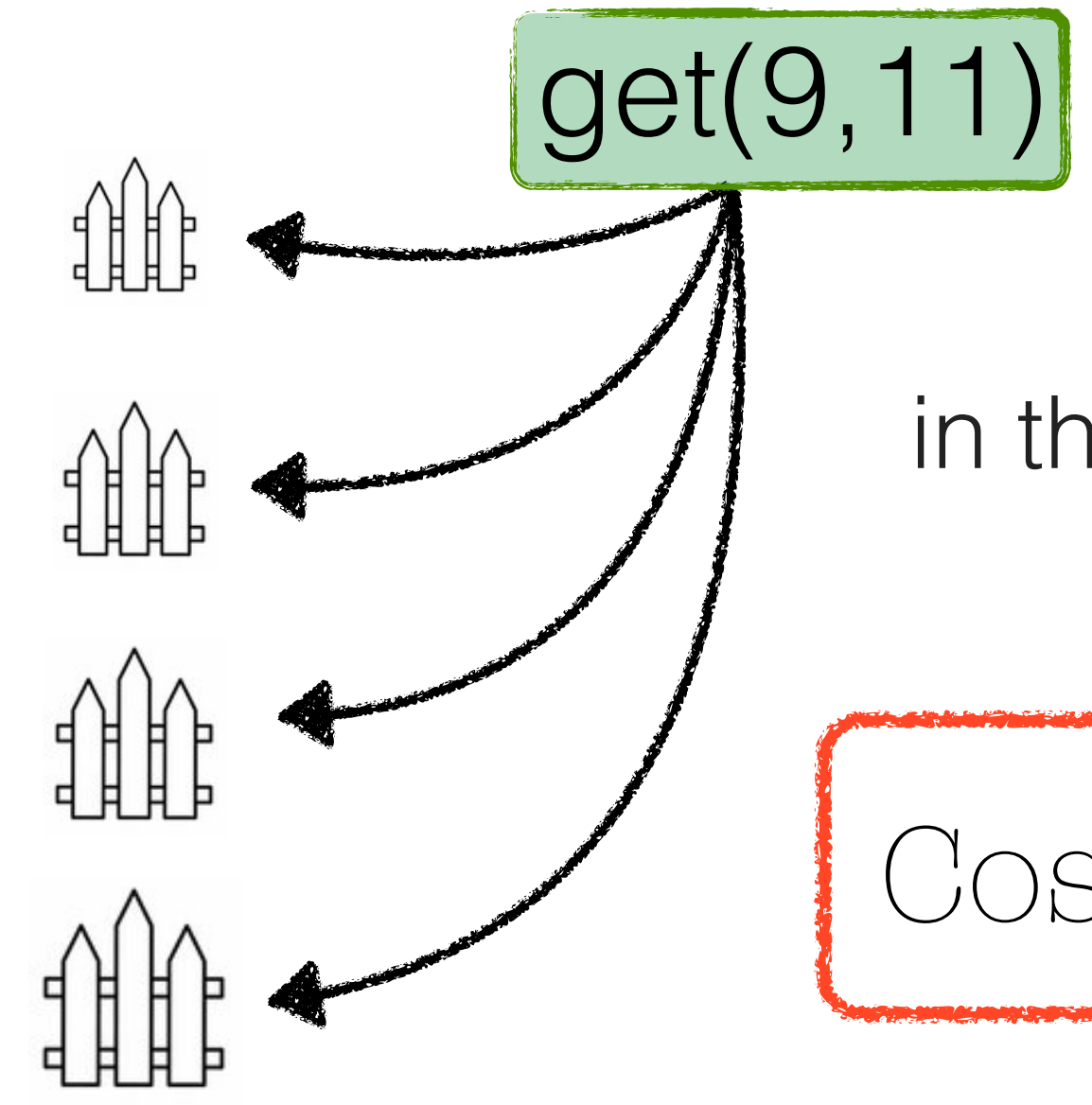
$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range



tiered LSM-tree



fence pointers

(page-wise zone map)

What if the **range query** has a **very low selectivity**?



in the **worst case**, need to read **2 pages per sorted run**

Cost of short range lookup =  $2 \cdot L$

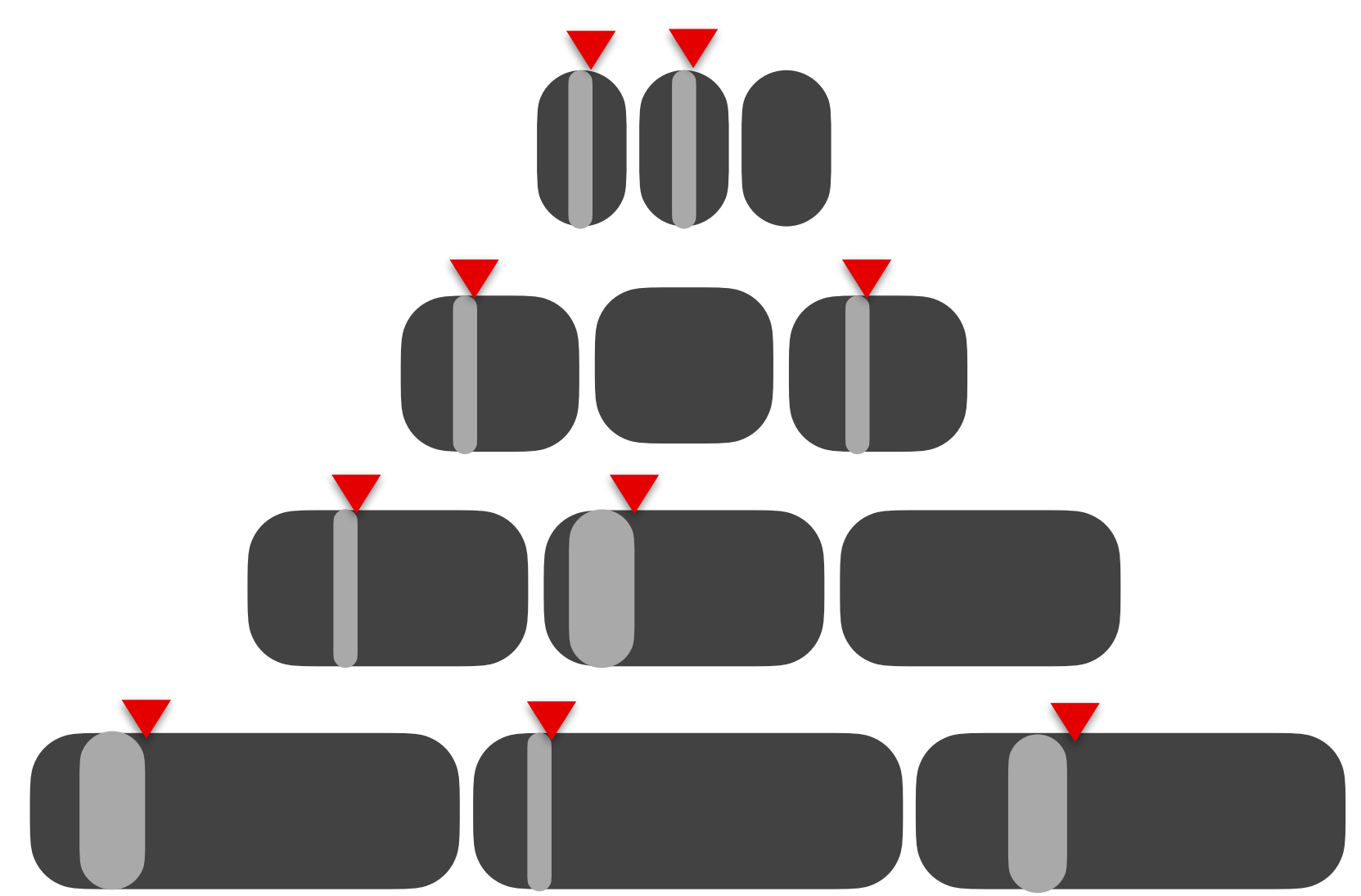
What about the **tiering**?



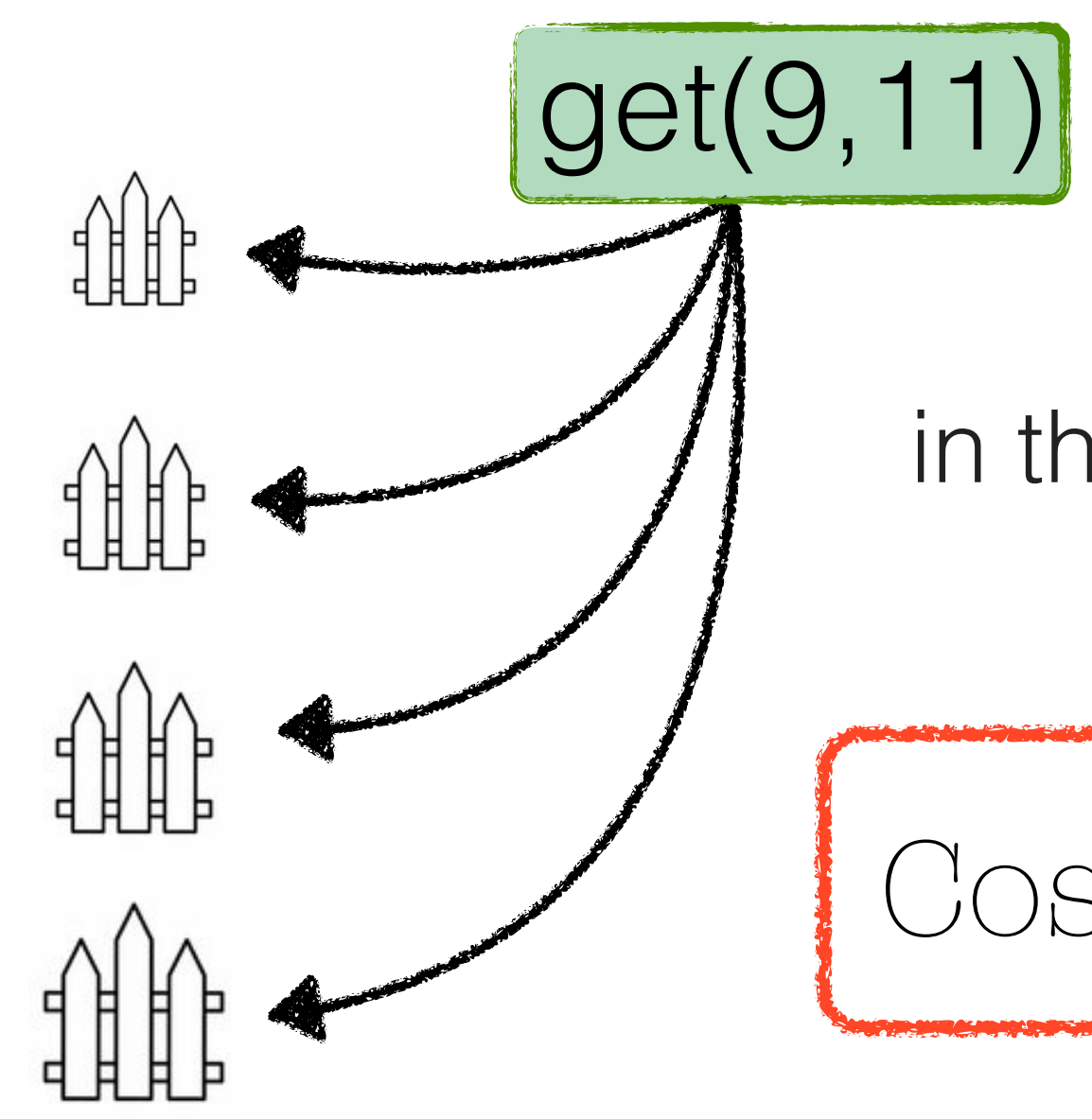
$P$ : pages in buffer  
 $B$ : entries/page  
 $L$ : #levels  
 $T$ : size ratio  
 $N$ : #entries in tree

# Range lookup cost

Looking for keys in a range



tiered LSM-tree



fence pointers

(page-wise zone map)

What if the **range query** has a **very low selectivity**?



in the **worst case**, need to read **2 pages per sorted run**

Cost of short range lookup =  $2 \cdot L$

What about the **tiering**?



in the **worst case**, need to read **2 pages per sorted run**

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

# Range lookup cost

Looking for keys in a range



What if the **range query** has a **very low selectivity**?

get(9, 11)

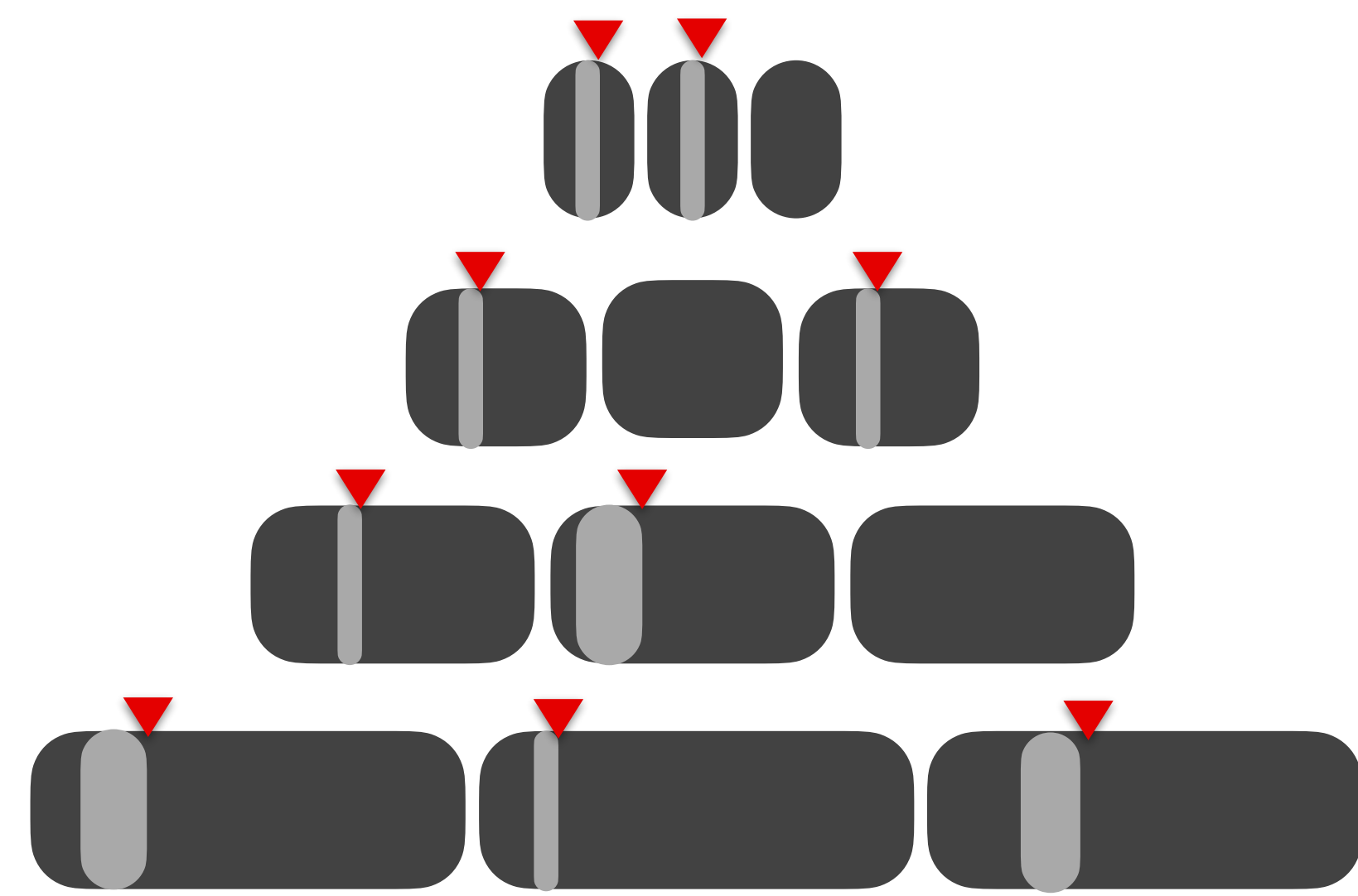
in the **worst case**, need to read **2 pages per sorted run**

Cost of short range lookup =  $2 \cdot L$

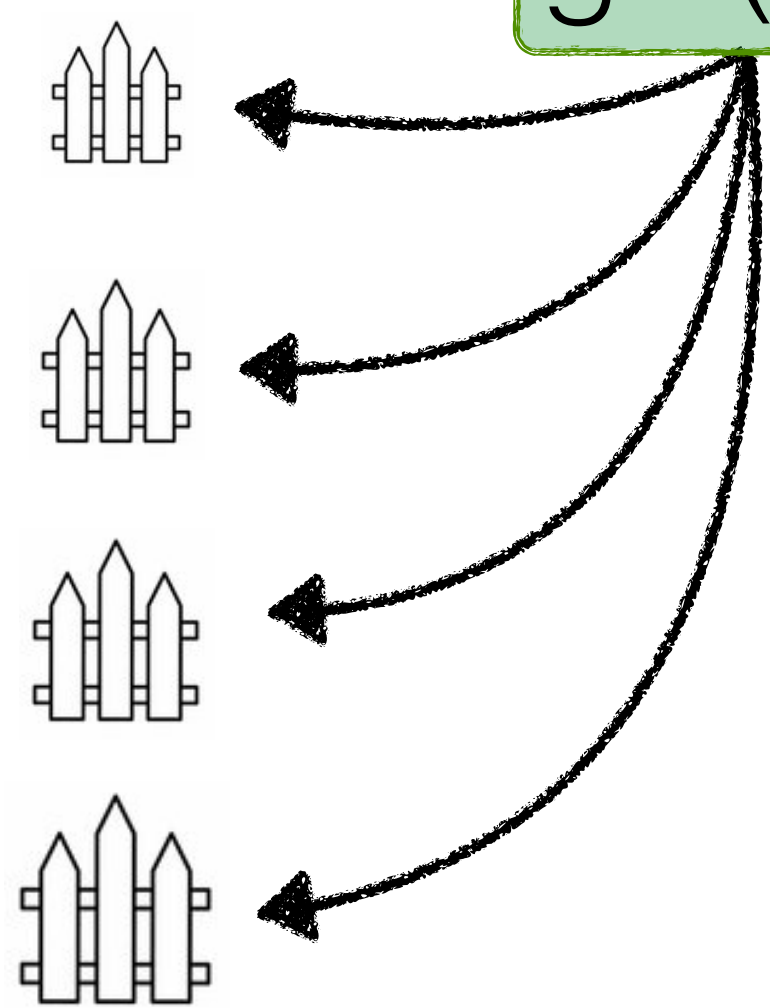
leveling

tiering

Cost of short range lookup =  $2 \cdot L \cdot T$



tiered LSM-tree



fence pointers

(page-wise zone map)



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree				
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	?	?	?	?
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# B<sup>+</sup>-trees

The go to index data structures for relational databases



*It could be said that the world's information is at our fingertips because of B-trees.*



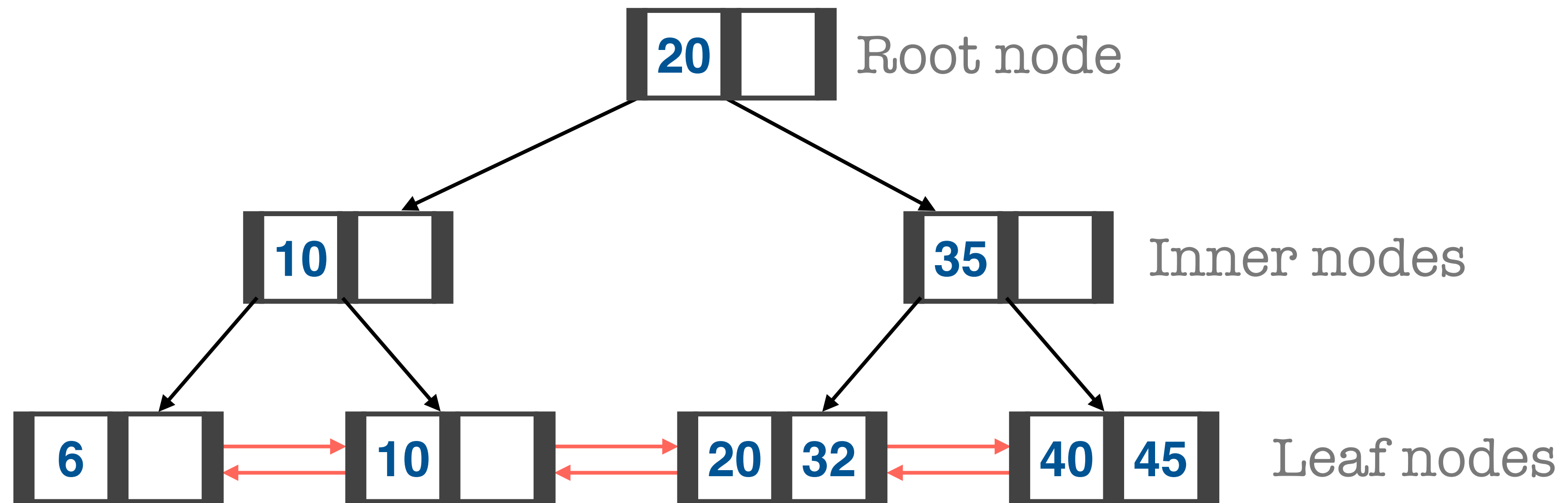
**Goetz Graefe**, Google, Ex-Microsoft, HP Fellow

Google ACM Software System Award



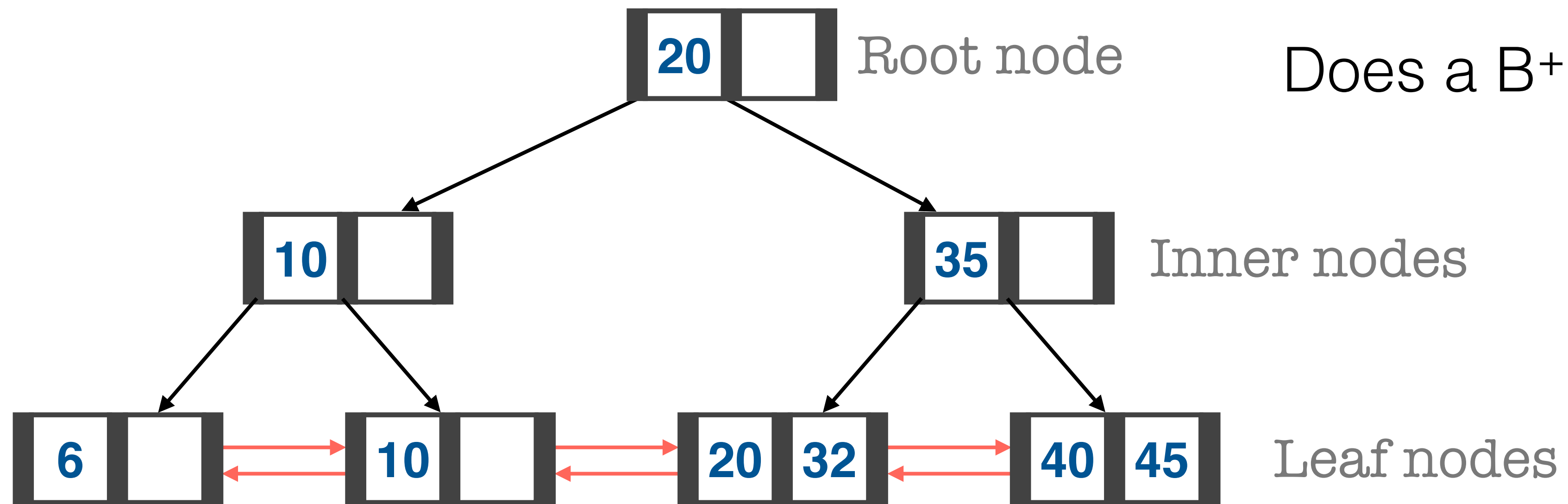
# B<sup>+</sup>-trees

The go to index data structures for relational databases

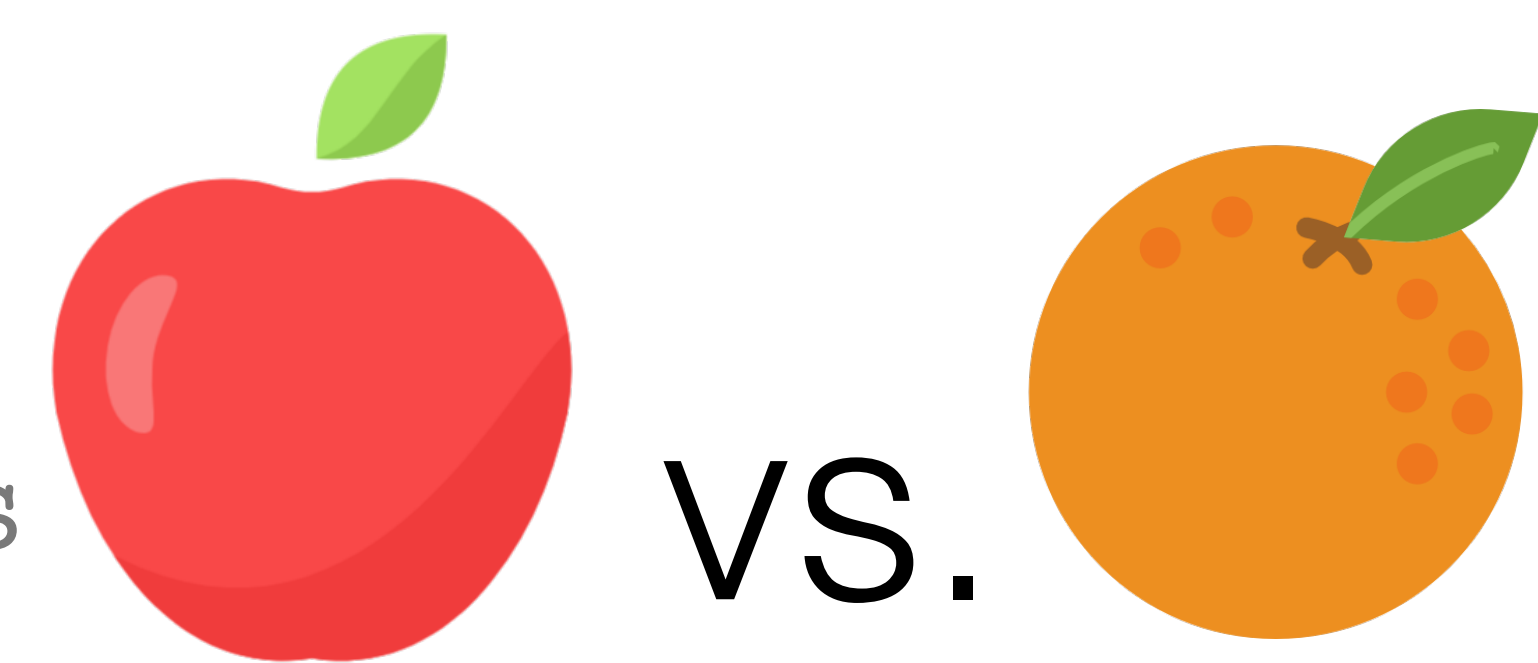


# B<sup>+</sup>-trees

The go to index data structures for relational databases

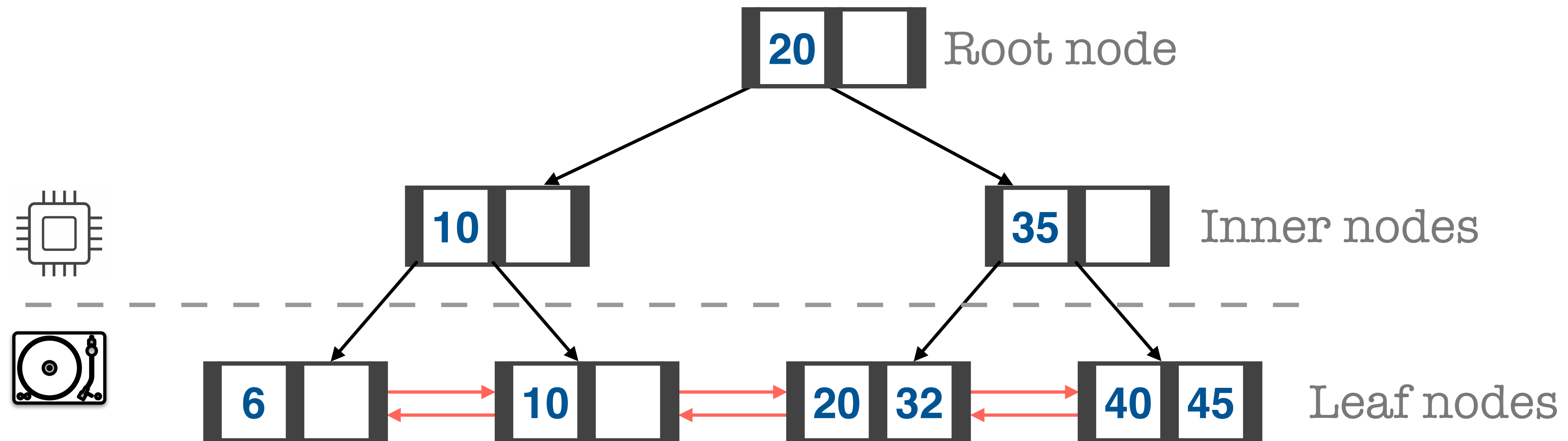


Does a B<sup>+</sup>-tree have a **buffer**?  
No!



# B<sup>+</sup>-trees

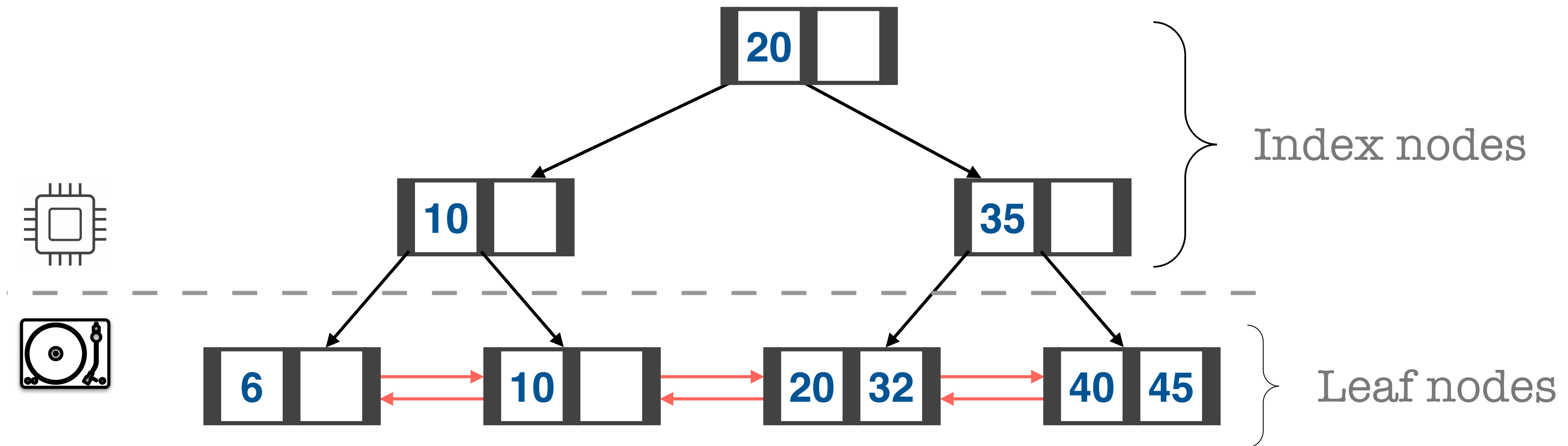
The go to index data structures for relational databases





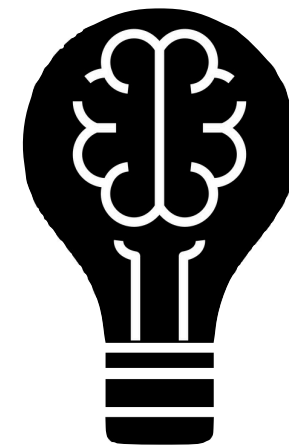
# B<sup>+</sup>-trees

The go to index data structures for relational databases



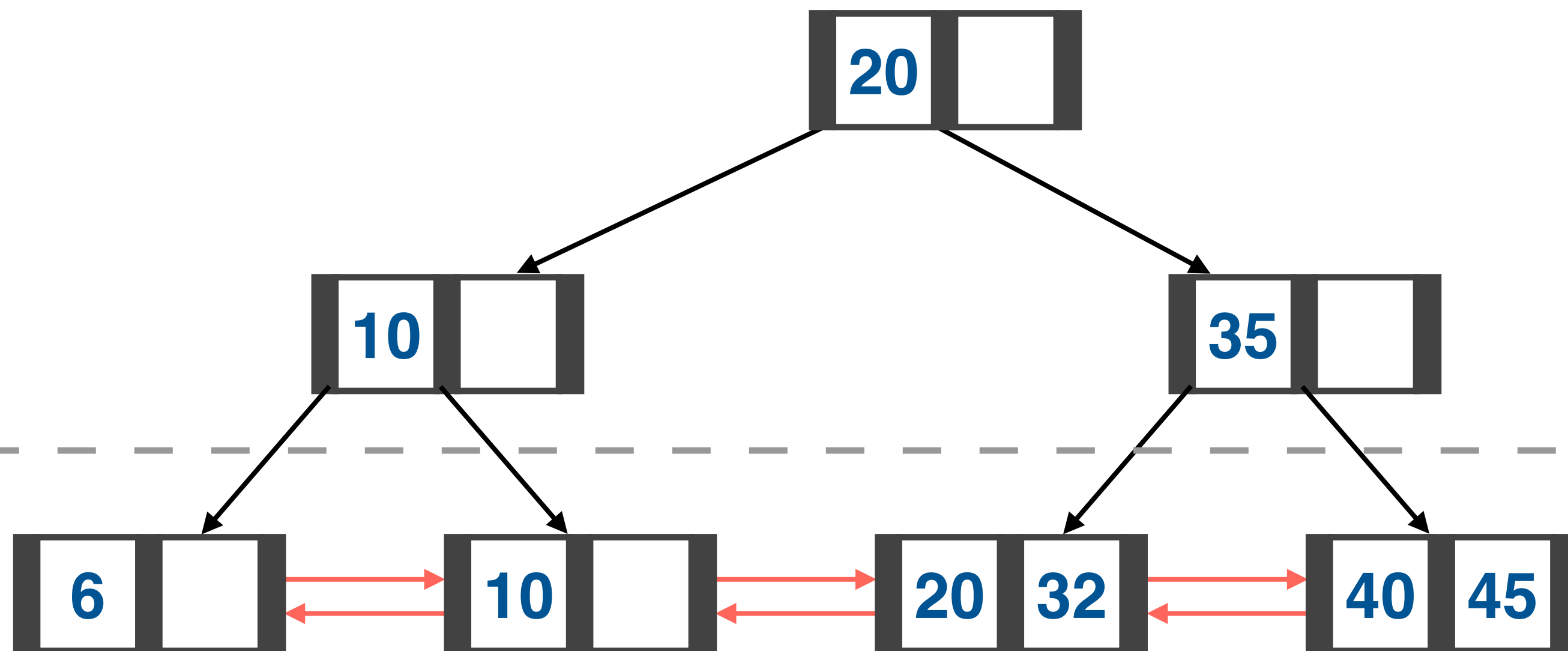
# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 1

What is the **insert cost** in a **buffered B<sup>+</sup>-tree**?

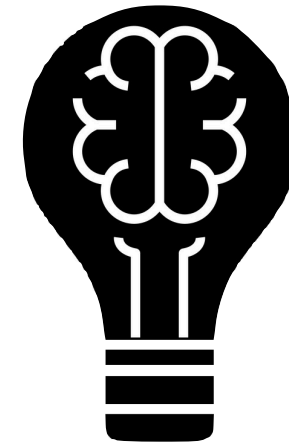


2 I/Os per insert

1 to read, 1 to write

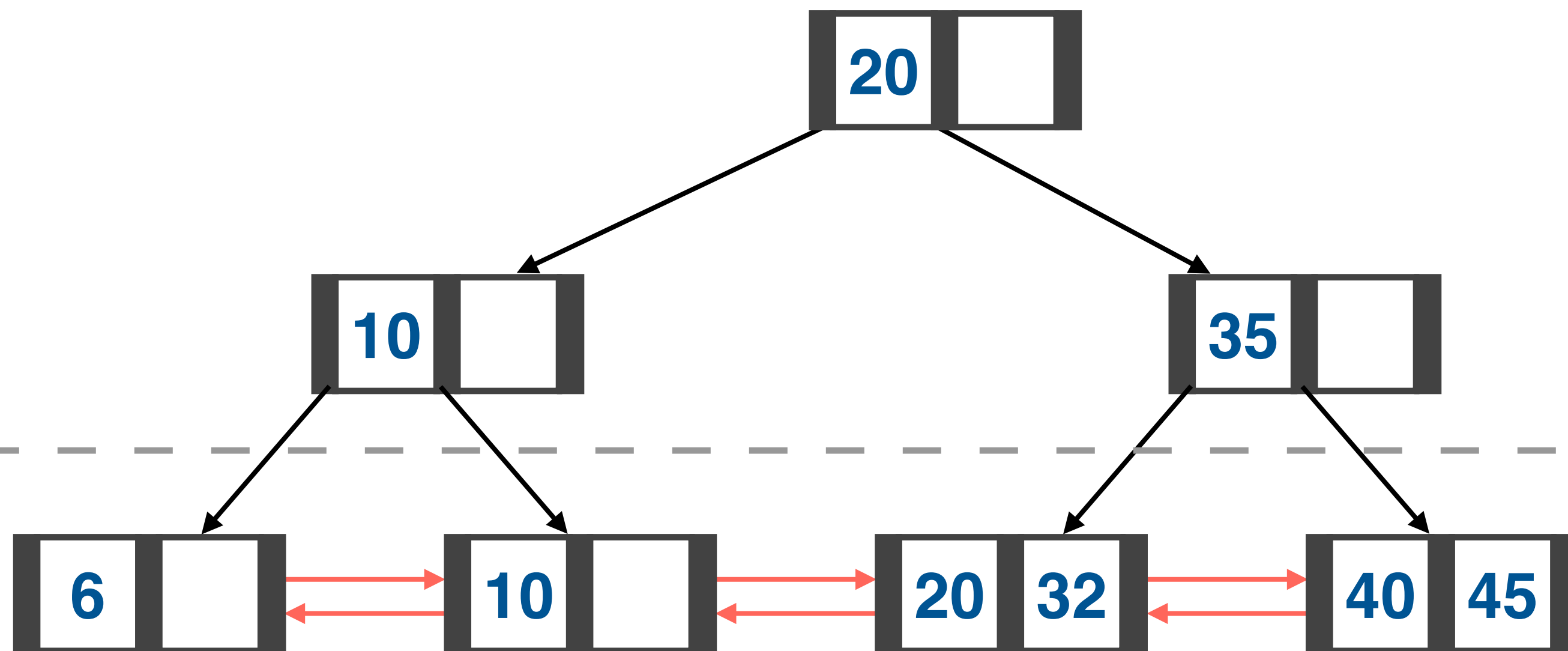
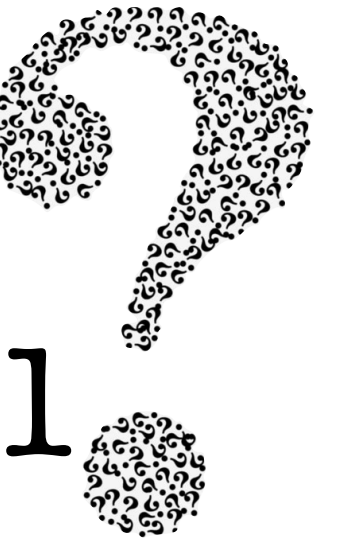
# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 1

What is the **insert cost** in a **buffered B<sup>+</sup>-tree**?



2 I/Os per insert

1 to read, 1 to write

Cost of ingestion = 2



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	?	?	?
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	?	?	?
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	?	?
Sorted array				
Log				

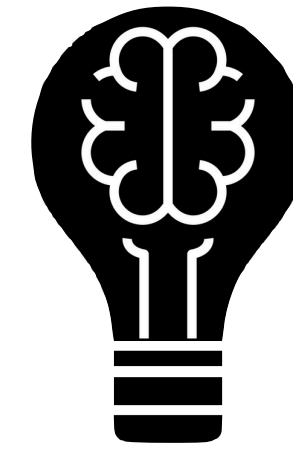
<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**   
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 2

What is the **range lookup cost** in a **buffered B<sup>+</sup>-tree**?

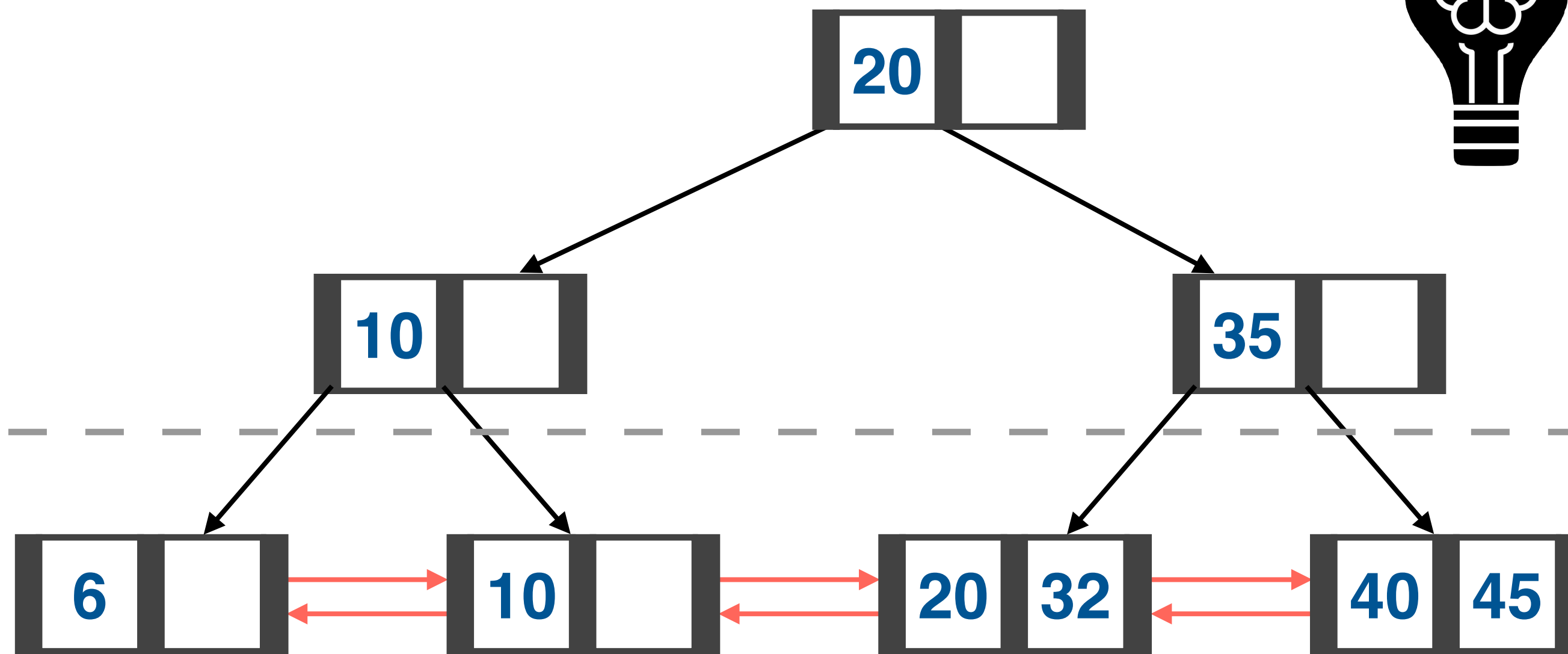
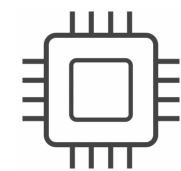
total entries in tree = **N**

#entries per page = **B**

---

#pages in tree = **N/B**

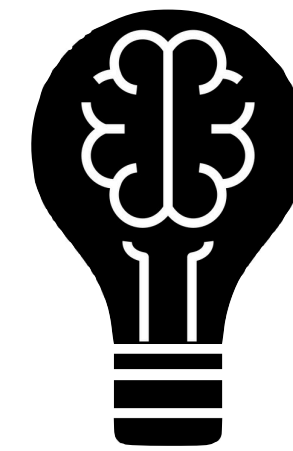
Cost of range lookup = **s · N/B**



**s** = **selectivity** of the range query

# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 2

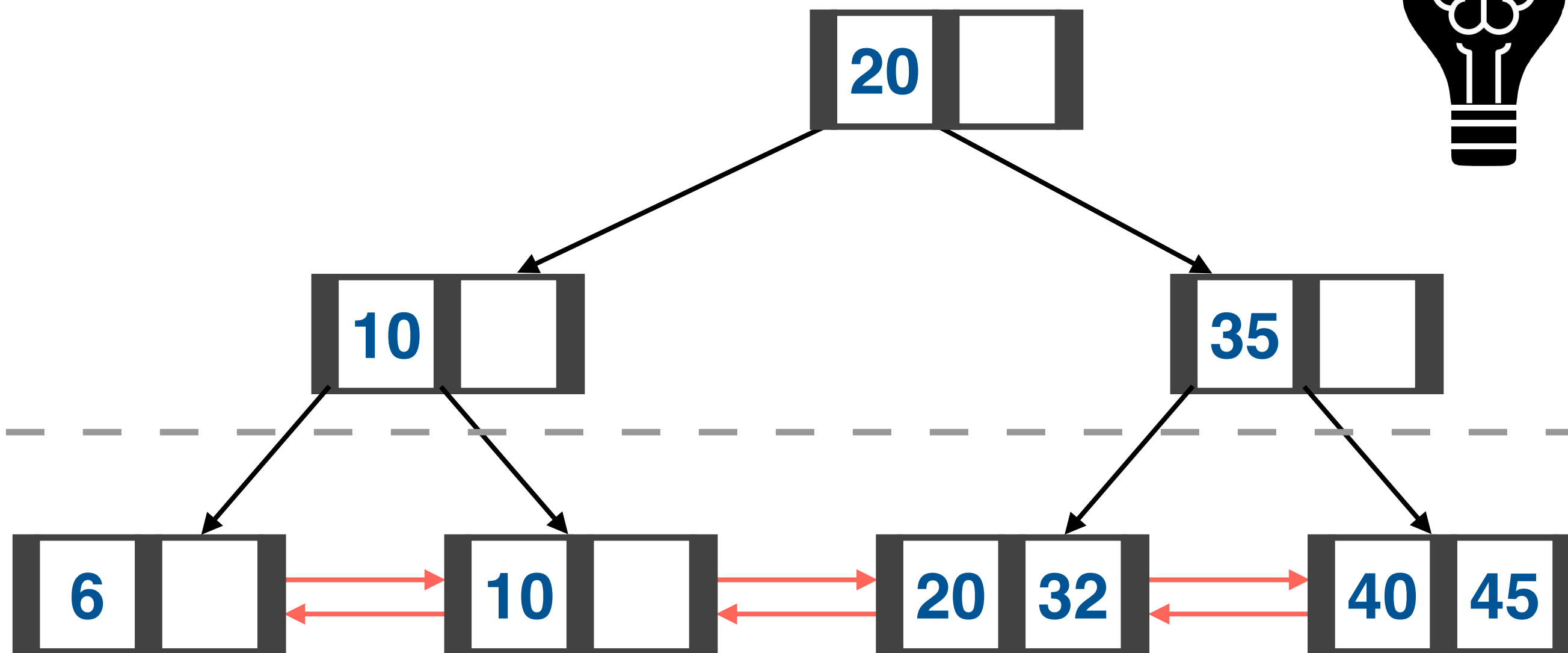
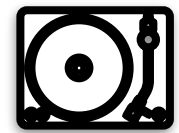
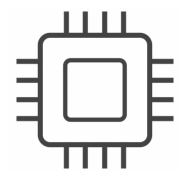
What is the **range lookup cost** in a **buffered B<sup>+</sup>-tree**?

total entries in tree = **N**

#entries per page = **B**

#pages in tree =  $N/B$

Cost of range lookup =  $s \cdot N/B$

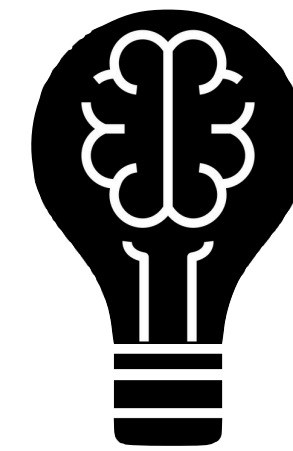


$s$  = **selectivity** of the range query

On average, nodes in a B<sup>+</sup>-tree are **67% full!**

# B<sup>+</sup>-trees

The go to index data structures for relational databases



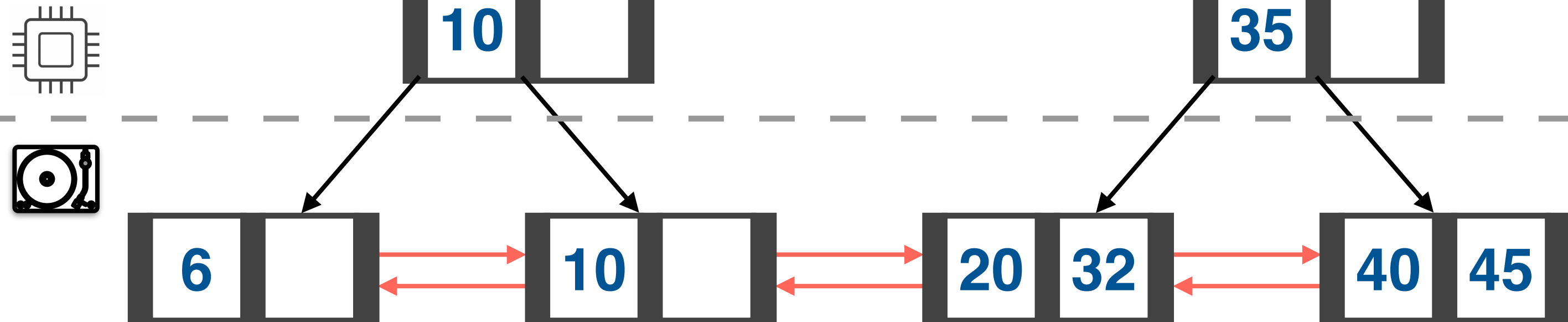
## Thought Experiment 2

What is the **range lookup cost** in a **buffered B<sup>+</sup>-tree**?

total entries in tree = **N**

#entries per page = **B**

#pages in tree = **N/B**



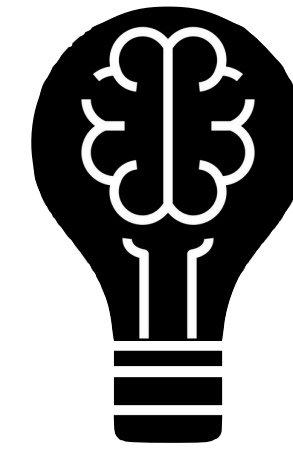
On average, nodes in a B<sup>+</sup>-tree are **67% full!**

On average, B<sup>+</sup>-tree have a **50% read amplification!**



# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 2

What is the **range lookup cost** in a **buffered B<sup>+</sup>-tree**?

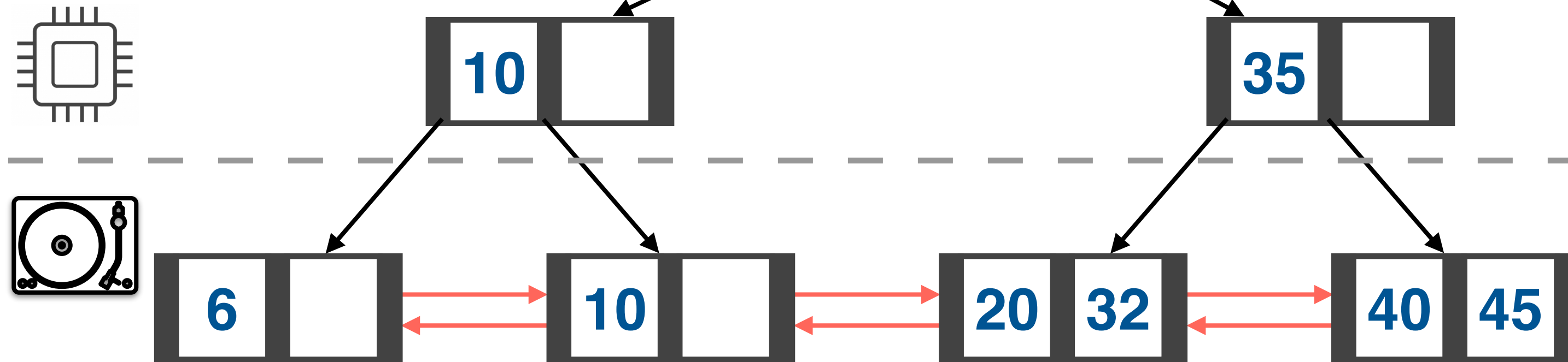
total entries in tree = **N**

#entries per page = **B**

#pages in tree = **N/B**

Cost of range lookup = **s · N/B**

**x1.5**

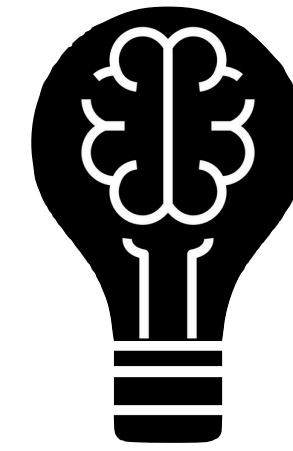


**s** = **selectivity** of the range query



# B<sup>+</sup>-trees

The go to index data structures for relational databases



## Thought Experiment 2

What is the **range lookup cost** in a **buffered B<sup>+</sup>-tree**?

total entries in tree = **N**

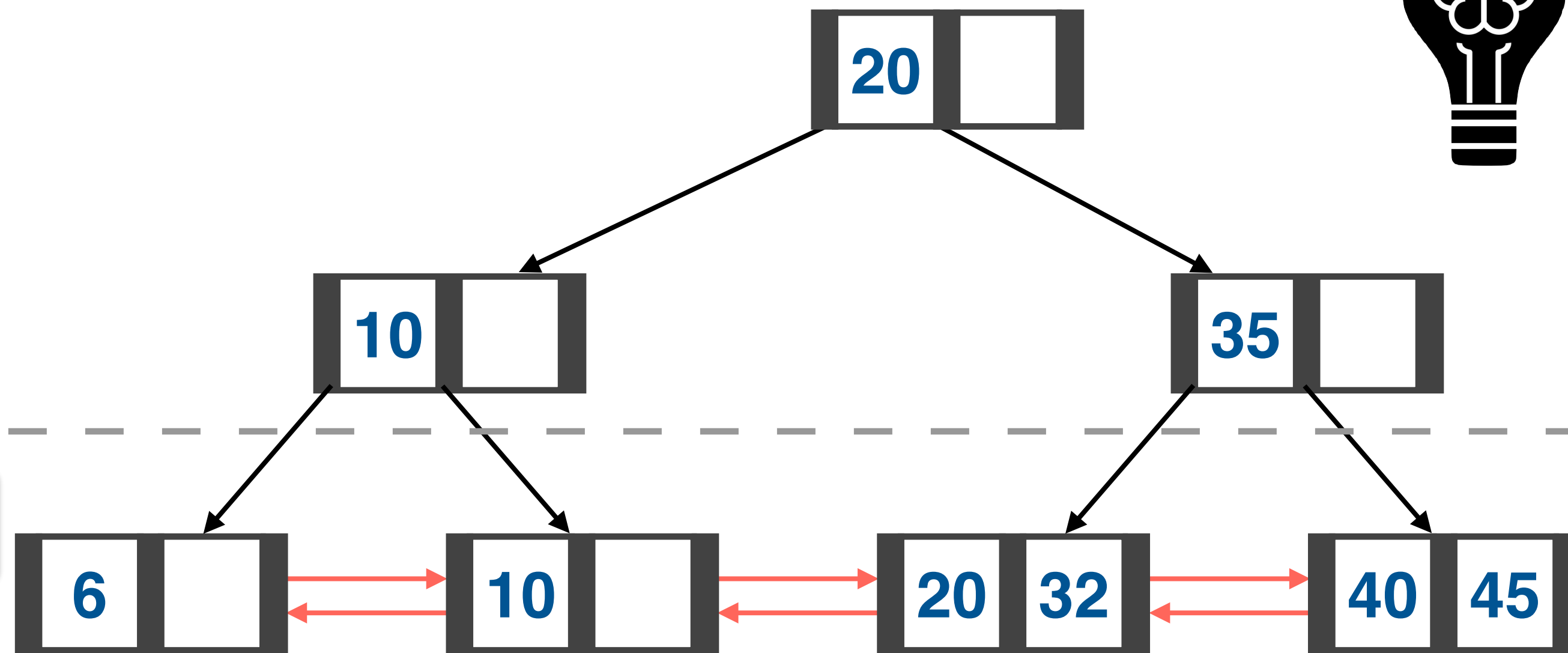
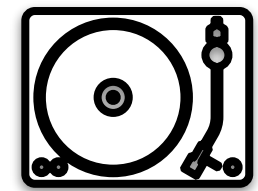
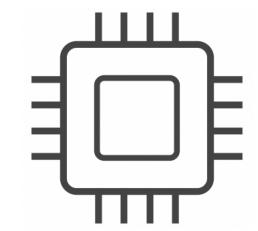
#entries per page = **B**

#pages in tree = **N/B**

$$\text{Cost of range lookup} = s \cdot N/B$$

x1.5

$$\text{Cost of range lookup} = 1.5 \cdot s \cdot N/B$$



**s** = **selectivity** of the range query

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	?
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**  <sup>\*</sup> 1.5x costlier than LSM-trees  
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	?
Sorted array				
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	?	?	?	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

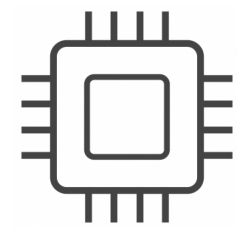
cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

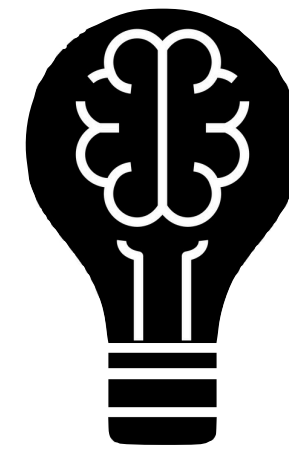
<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Sorted array

Optimizing for queries

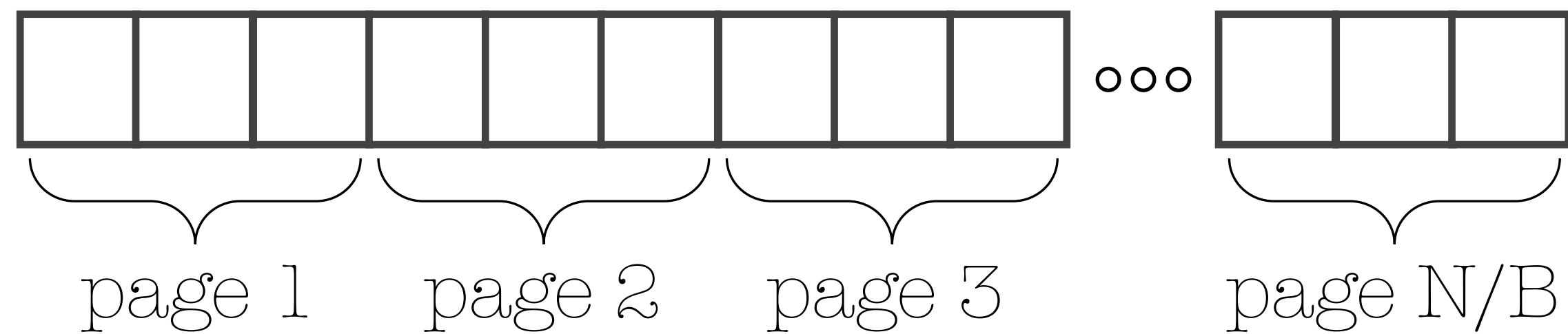


Buffer



## Thought Experiment 3

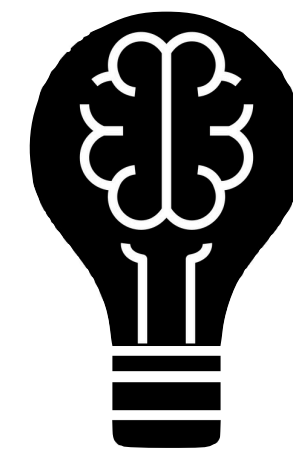
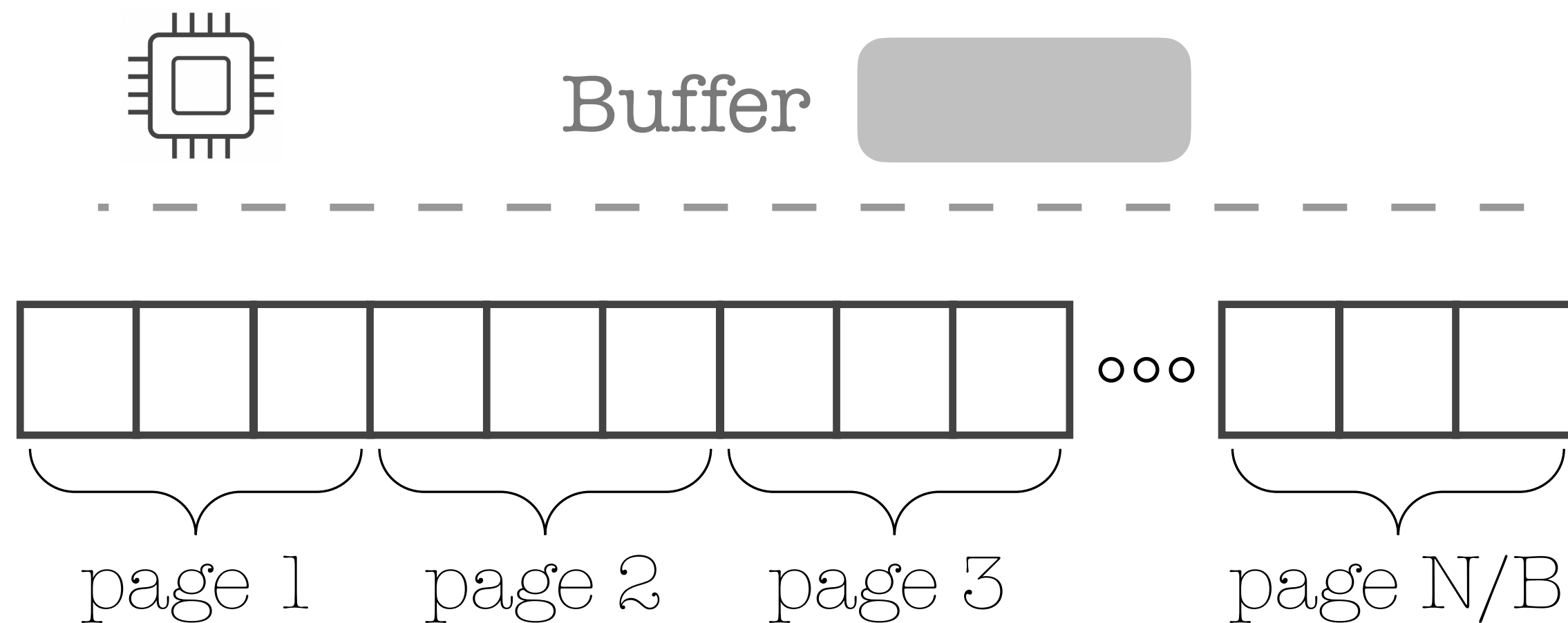
What is the **insert cost** in a **sorted array** with a **buffer**?





# Sorted array

Optimizing for queries



## Thought Experiment 3

What is the **insert cost** in a **sorted array** with a **buffer**?

**read + write all  $N/B$  pages**  
for uniform data distribution

**B entries** are written to disk per I/O

$$\text{Cost of ingestion} = N/B^2$$



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	?	?	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	?	?	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

<sup>\*</sup> 1.5x costlier than LSM-trees

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	?	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	?	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B+-tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	?
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log				

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B+-tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	?	?	?	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	?	?	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B+-tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	?	?	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove "1 +" to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the "L" factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	?	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	?	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**  <sup>\*</sup> 1.5x costlier than LSM-trees  
 cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	?

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
<b>Leveled LSM-tree</b>	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
<b>Tiered LSM-tree</b>	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
<b>B+-tree</b>	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
<b>Sorted array</b>	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
<b>Log</b>	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B <sup>+</sup> -tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B+-tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
B+-tree	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost

# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
<b>B<sup>+</sup>-tree</b>	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**  <sup>\*</sup> 1.5x costlier than LSM-trees  
 cost for **non-empty lookups** (remove "1 +" to get cost of empty lookups)

<sup>\*</sup> **Monkey** shaves off the "L" factor from the cost



# Cost analysis

Counting all I/Os

data structure	ingestion cost	point lookup cost <sup>*</sup>	long range lookup cost	short range lookup cost
Leveled LSM-tree	$O(L \cdot T / B)$	$O(1 + \phi \cdot L)^*$	$O(s \cdot N / B)$	$O(L)$
Tiered LSM-tree	$O(L / B)$	$O(1 + \phi \cdot L \cdot T)^*$	$O(s \cdot N / B)$	$O(L \cdot T)$
<b>B<sup>+</sup>-tree</b>	$O(1)$	$O(1)$	$O(s \cdot N / B)^*$	$O(1)$
Sorted array	$O(N/B^2)$	$O(1)$	$O(s \cdot N / B)$	$O(1)$
Log	$O(1/B)$	$O(N/B)$	$O(N/B)$	$O(N/B)$

<sup>\*</sup> with **fence pointers & Bloom filter** with **FPR =  $\phi$**

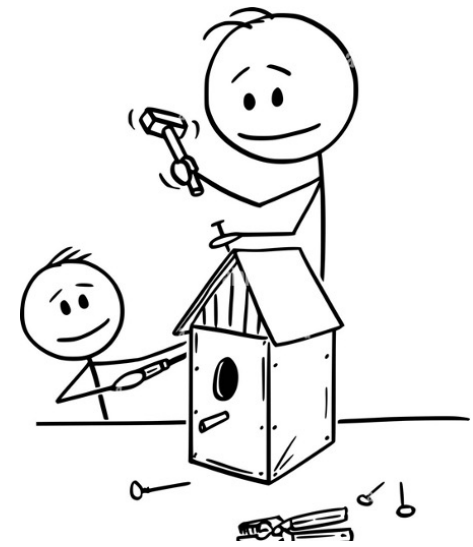
cost for **non-empty lookups** (remove “1 +” to get cost of empty lookups)

<sup>\*</sup> 1.5x costlier than LSM-trees

<sup>\*</sup> **Monkey** shaves off the “L” factor from the cost



# Projects



# Projects

Let's get our hands dirty

Class Project

**Groups of 2**

Out **tonight**; due in **December**

**Multiple milestones** in between

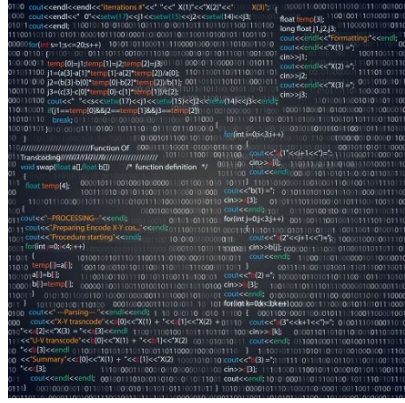
**Systems project**

**Research project**









# Class project: **Theme**

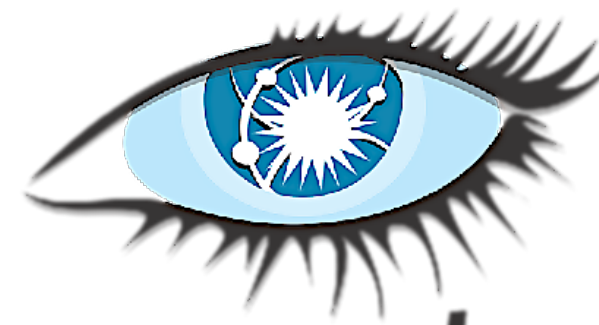


Working with state-of-the-art systems

## NoSQL key-value stores



RocksDB



cassandra

APACHE  
HBASE



MEMCACHED

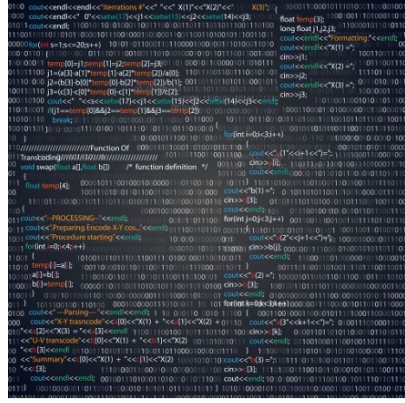
Asterix\*DB

Google™  
BigTable

amazon  
DynamoDB

redis





# Class project: **Theme**



Working with state-of-the-art systems

## NoSQL key-value stores



**RocksDB**

Asterix\*DB™



*cassandra*

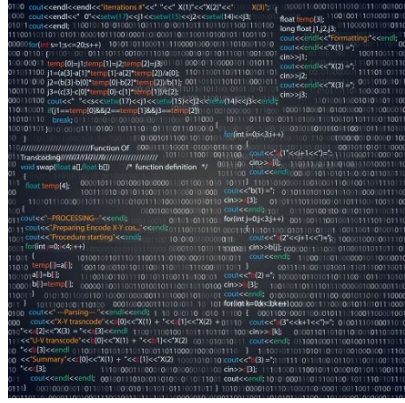
APACHE  
**HBASE**



MEMCACHED

Google  
BigTable





# Systems Project 1

Developing your way through!

## Implementing an LSM-based KV store

Data layout: **leveling & tiering**

File structure: **sorted run as one or more files; files have multiples pages**

Auxiliary data structures: **Bloom filters & fence pointers**

Granularity: filters **per sorted run / file**; fence pointers **per page**

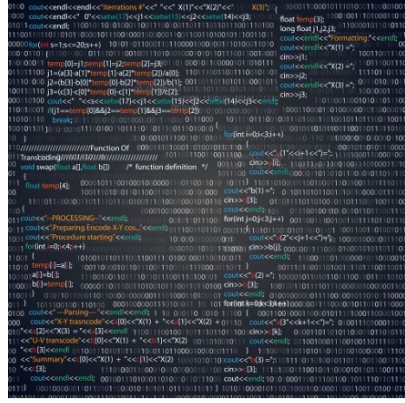
APIs: **insert, update, delete, point** and **range queries**

**Parameterize everything!**

**Report all relevant performance numbers!**







# Systems Project 2

Developing your way through!

## Implementing Compactions on RocksDB

Working with **RocksDB** — **C++**

Data layout: **leveling, tiering, hybrids**

Data movement policy: **partial-level, full-level, tiered**

File picking policy: **coldest, oldest, most tombstones**

**Benchmark** performance for different workloads





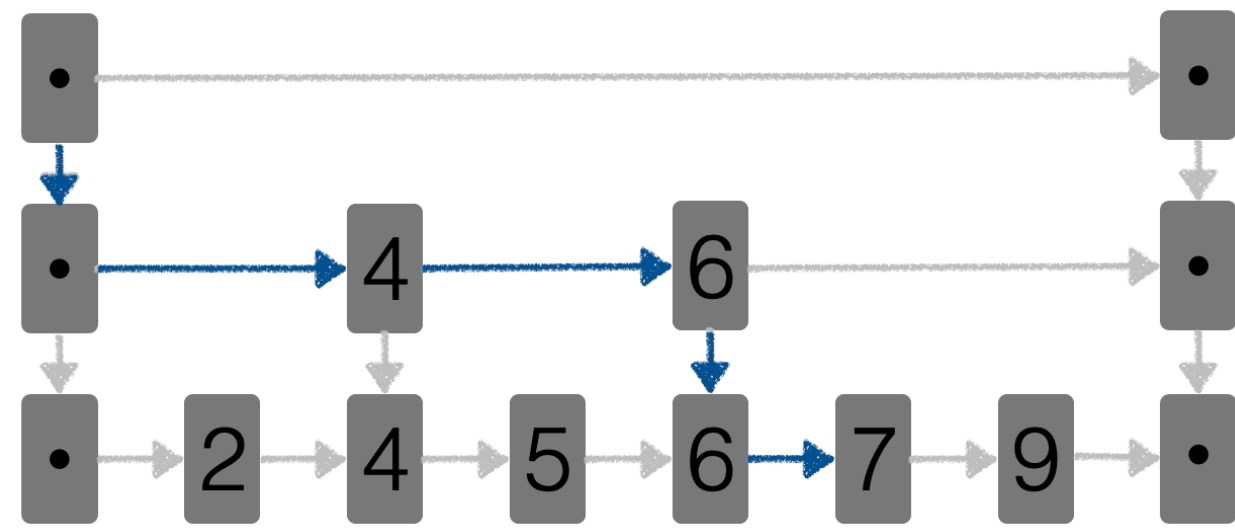
# Reserach projects

Solving problems on your way through!

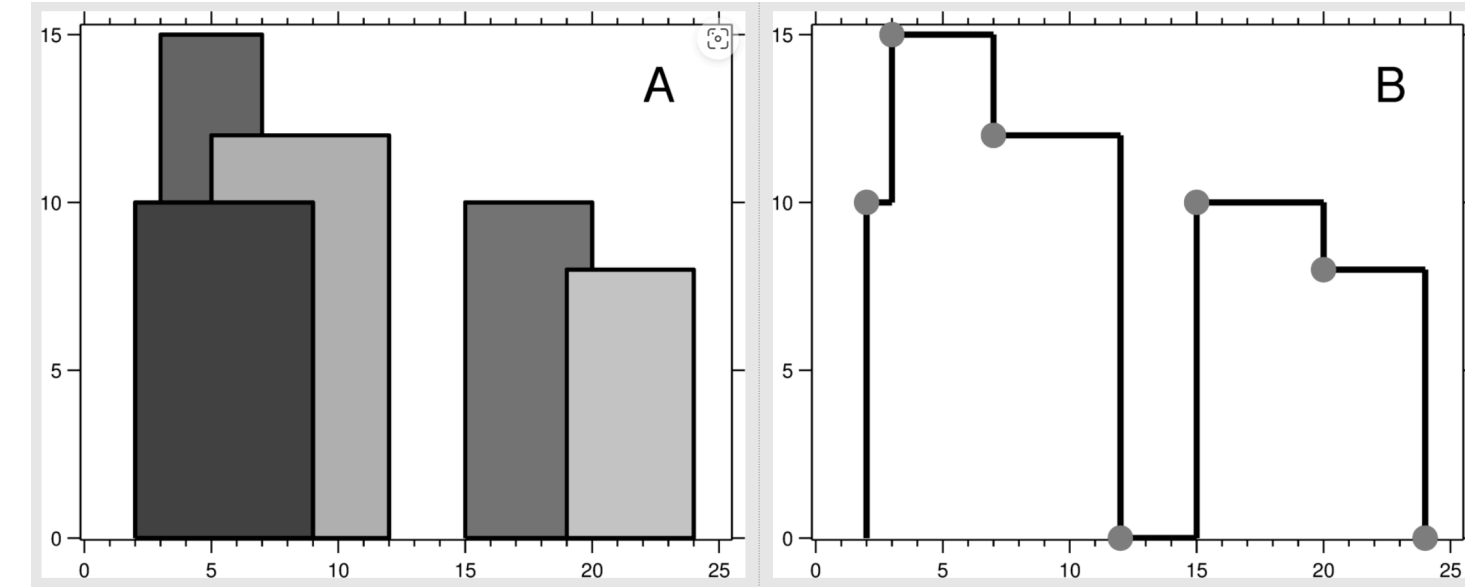


# Reserach projects

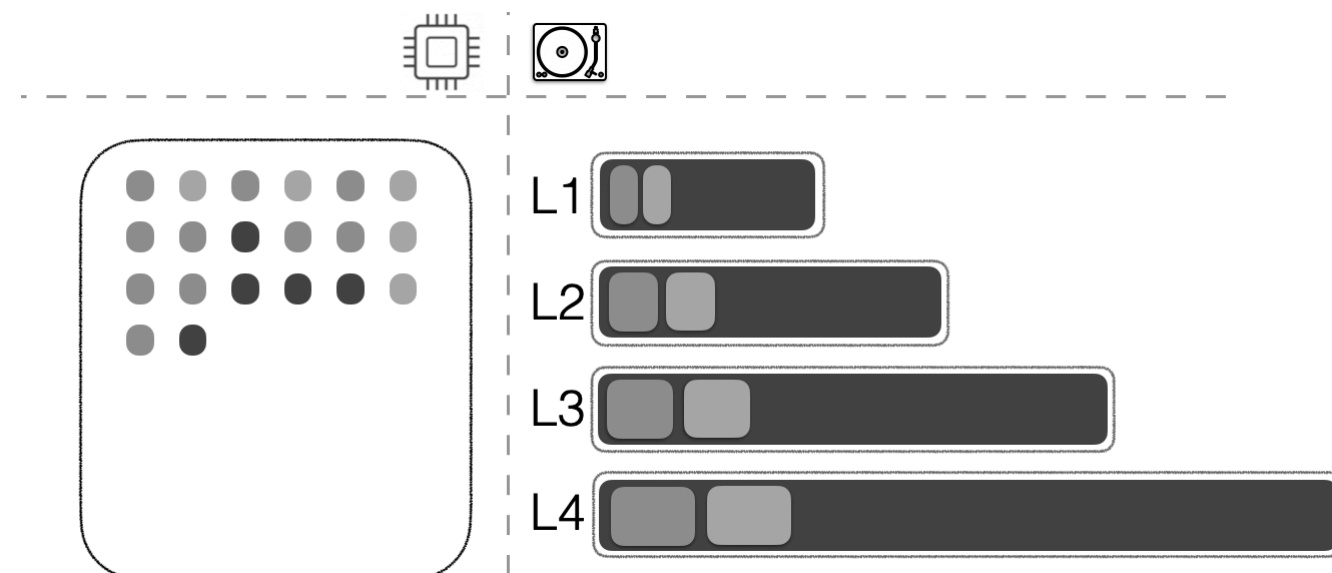
Solving problems on your way through!



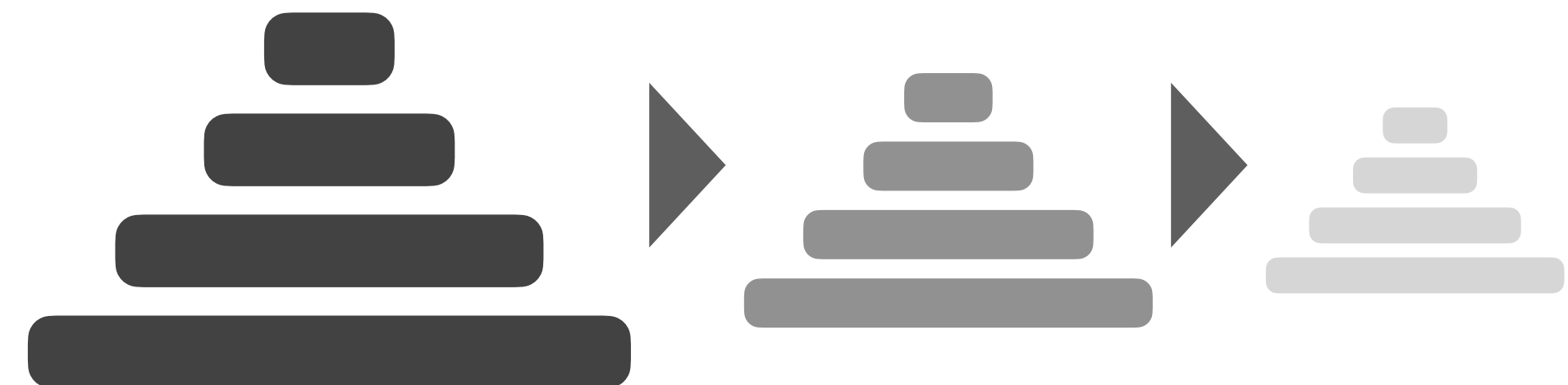
## Self-designing LSM Buffer



## Range Deletes in LSM



## Optimal Caching in LSM

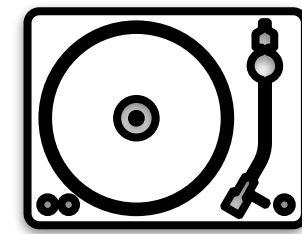
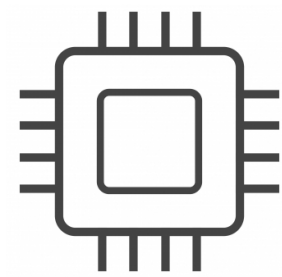


## Time-Traveling LSM

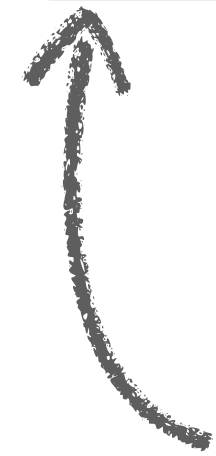
# Buffer Design



# Buffer Design



buffer



A **small** component  
acting as a **bottleneck**

Offers many **design choices**

L1



L2



L3



L4



$P$ : pages in buffer  
 $B$ : entries/page

# Buffer Implementation

**vector**

**skiplist**

**hashmap**

**ingestion**  
cost

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

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

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

**point query**  
cost

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

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

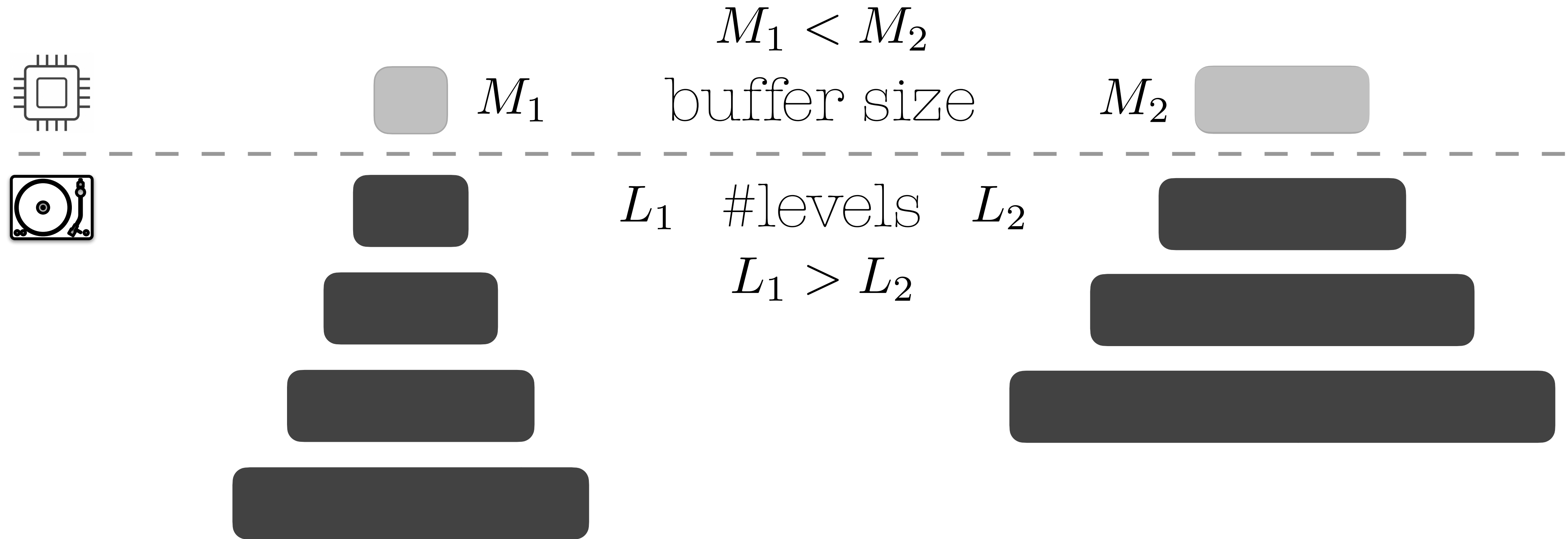
$$\mathcal{O}(1)$$

Ingestion-only  
workloads

Mixed  
workloads

I/O-bound  
workloads

# Size of the Buffer

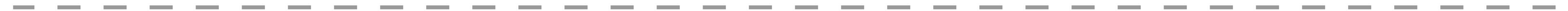
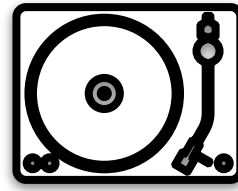
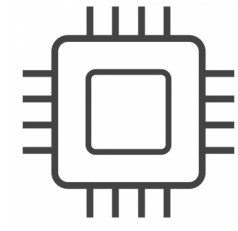


- frequent flushes
- smaller but more levels
- poor read performance

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

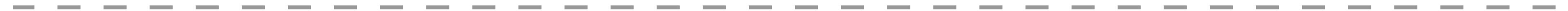
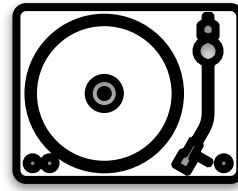
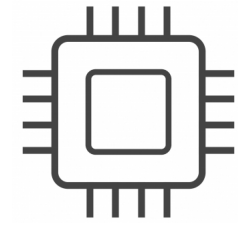


# #Buffer Components



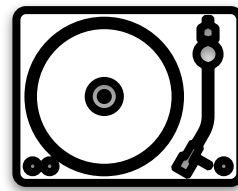
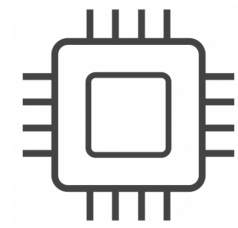


# #Buffer Components



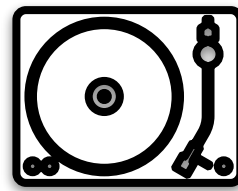
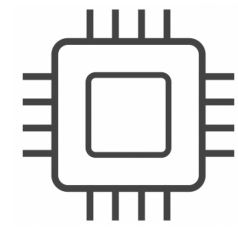
# #Buffer Components

immutable  
buffers



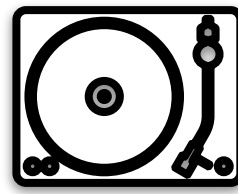
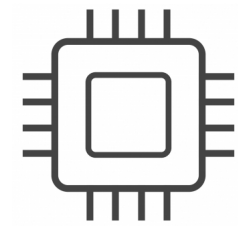
# #Buffer Components

immutable  
buffers



# #Buffer Components

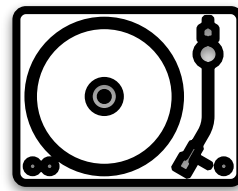
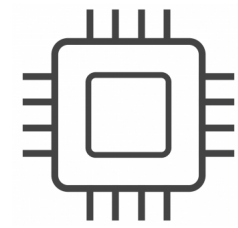
immutable  
buffers





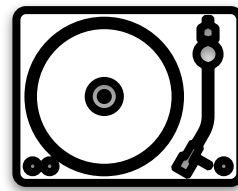
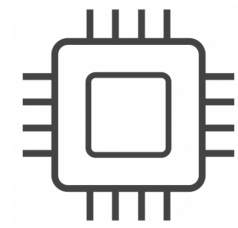
# #Buffer Components

immutable  
buffers



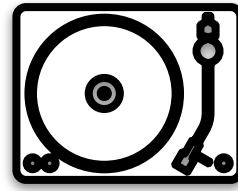
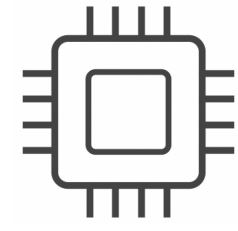
# #Buffer Components

immutable  
buffers

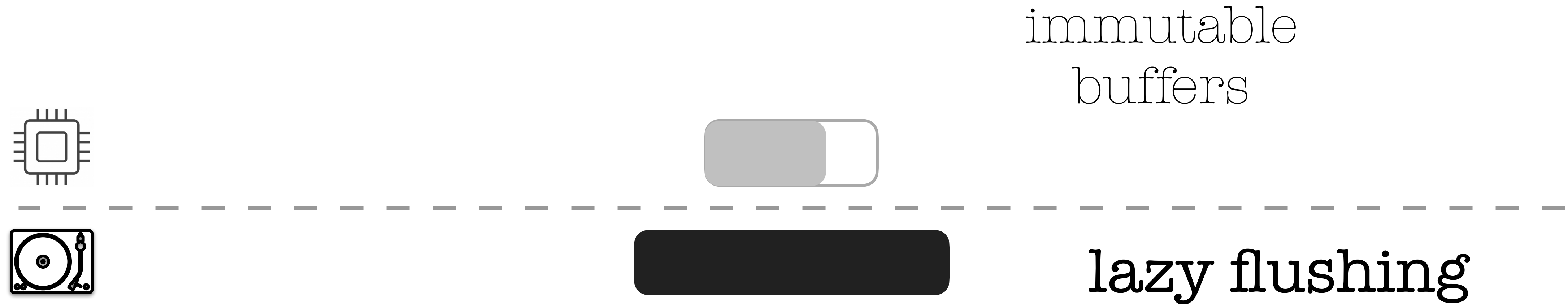


# #Buffer Components

immutable  
buffers



# #Buffer Components



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





# Research Project 1

Solving problems on your way through!

## Self-designing LSM Buffer

Which **buffer implementation** to choose?

How to **tune** the particular buffer implementation?

What if the **workload changes**?

What if the **application requirements** change?

**Benchmark performance** for different workloads & performance goals

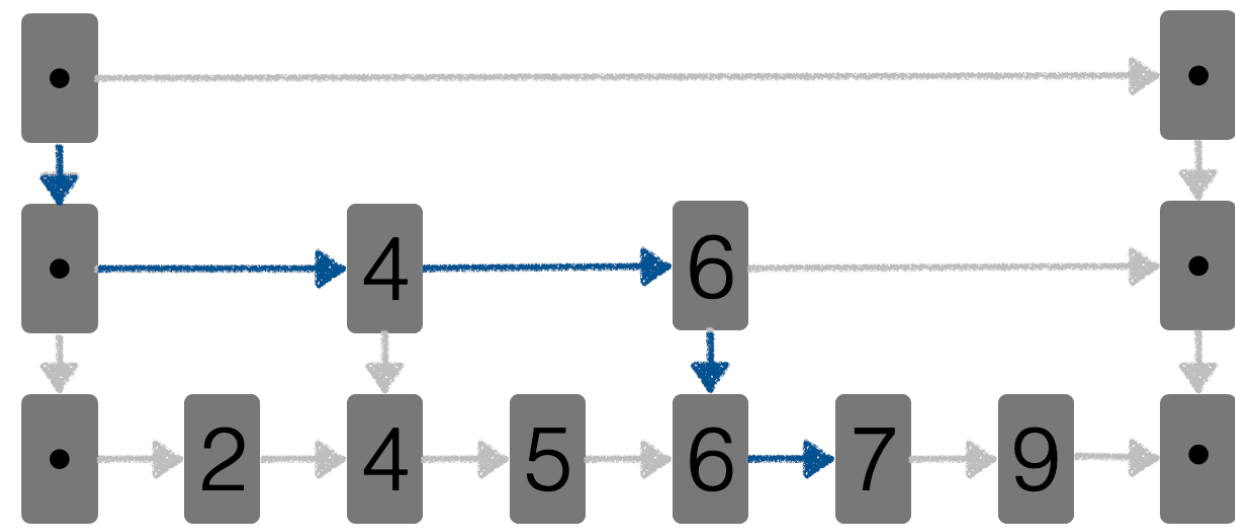
Design an **ML-based optimizer**

Integrate the **self-designing buffer** with RocksDB

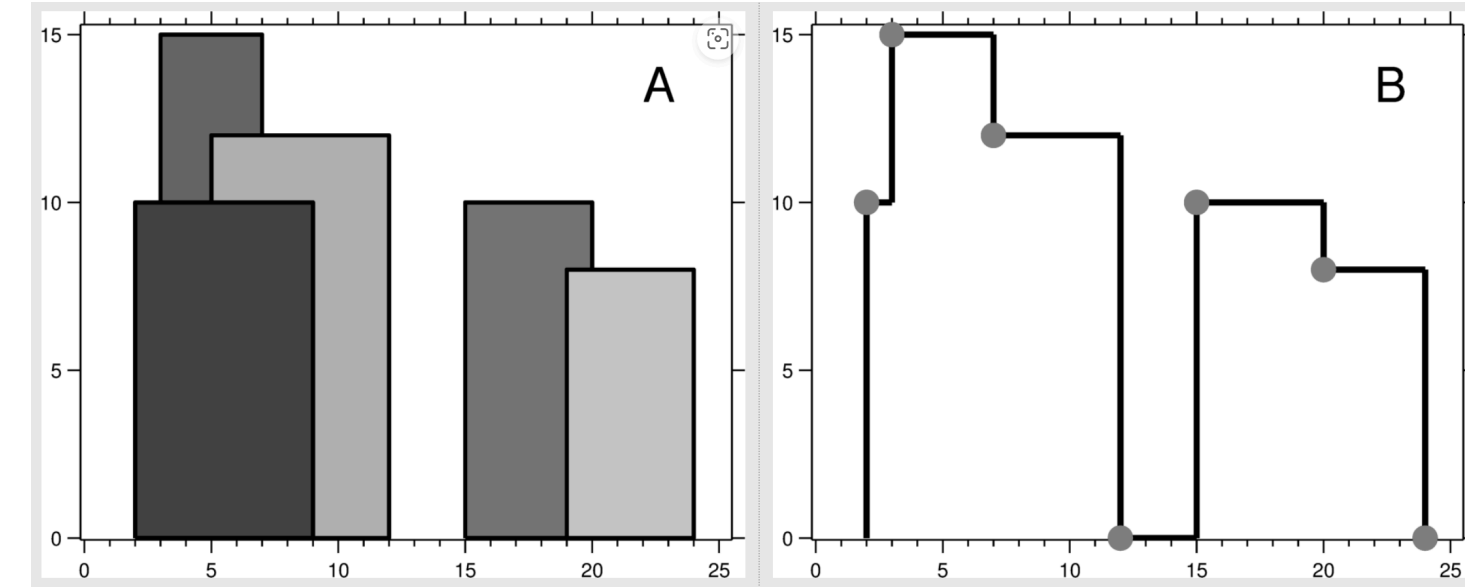


# Reserach projects

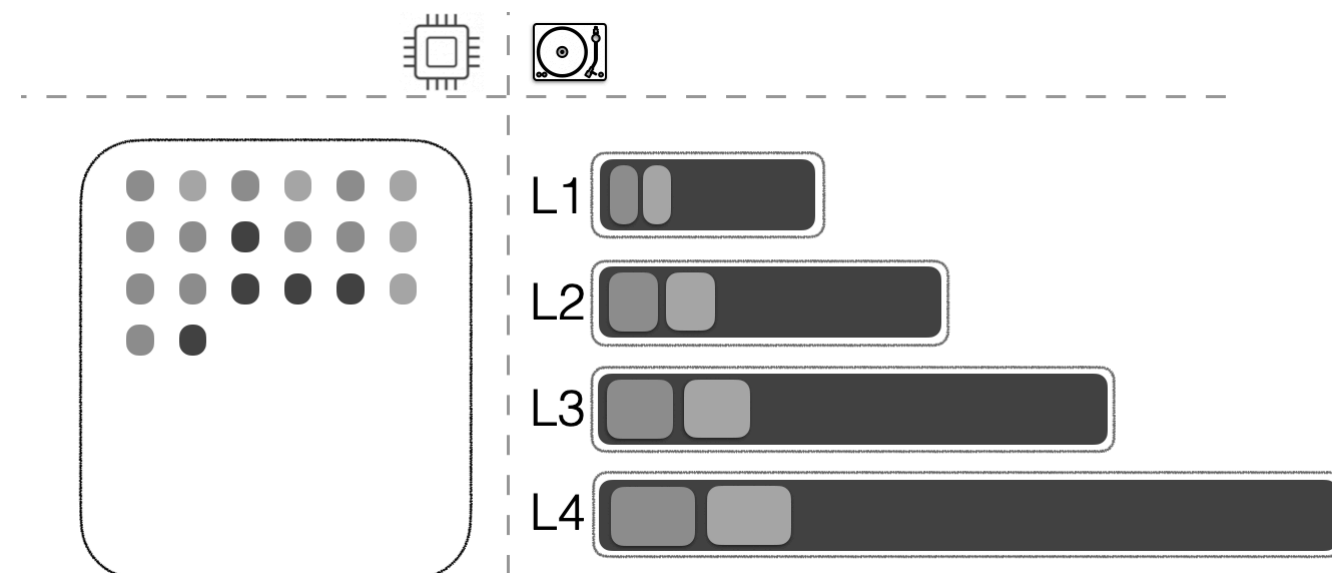
Solving problems on your way through!



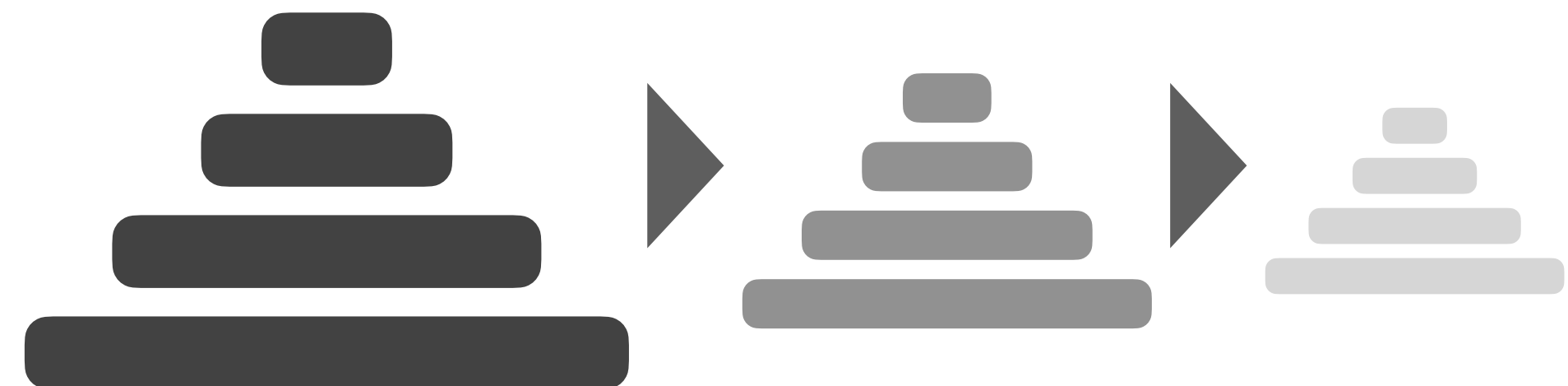
## Self-designing LSM Buffer



## Range Deletes in LSM



## Optimal Caching in LSM

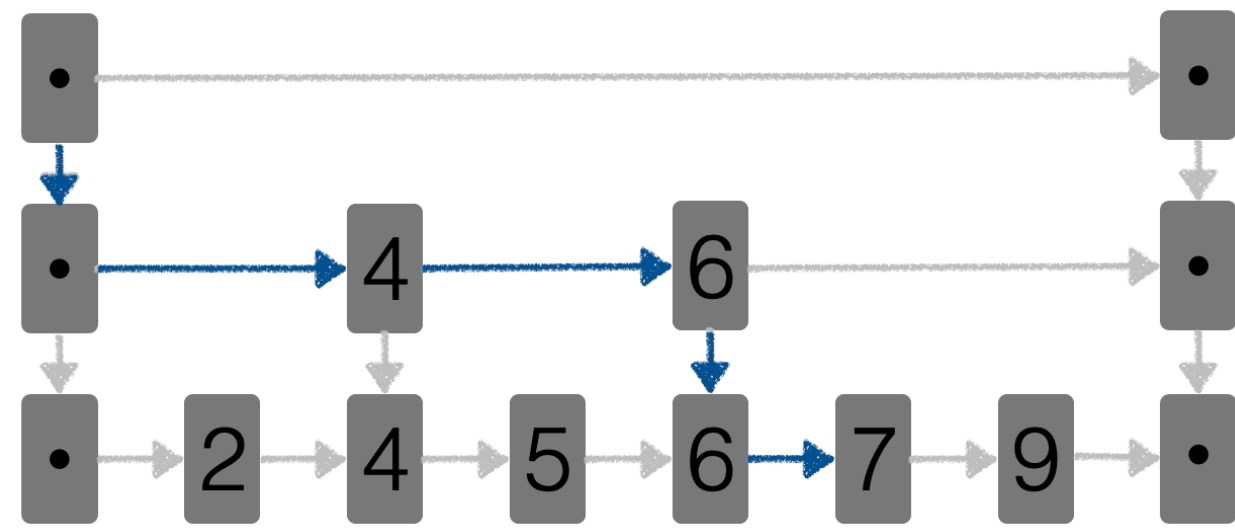


## Time-Traveling LSM

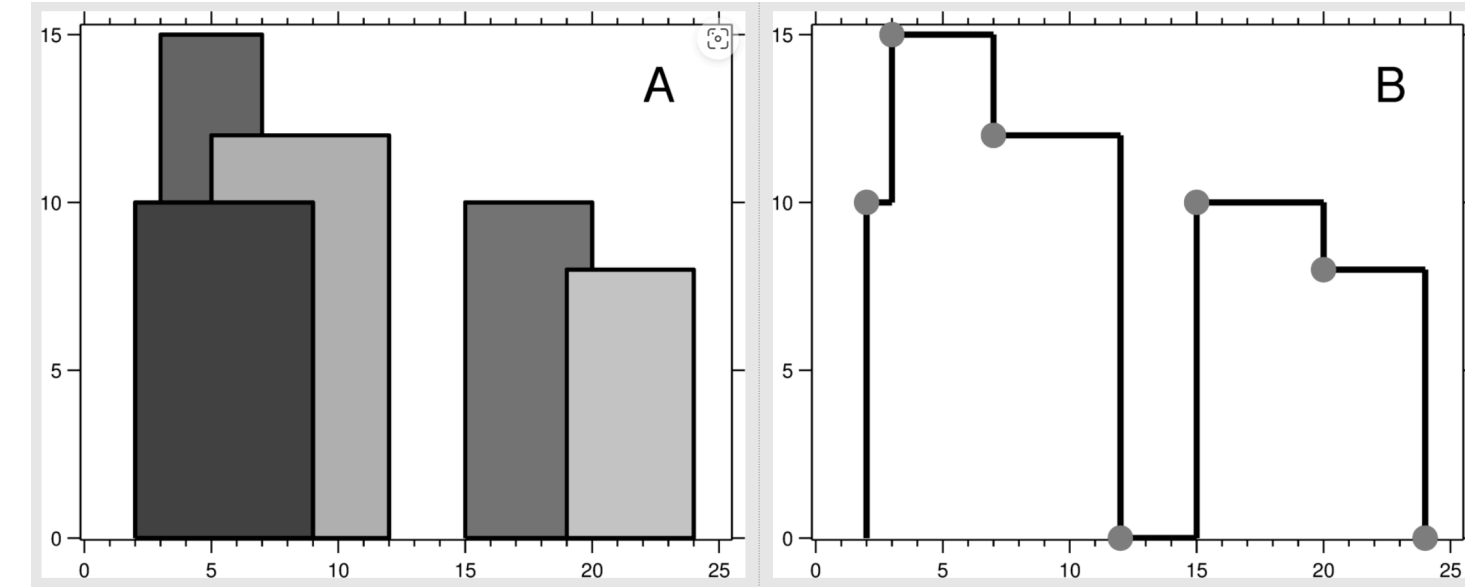


# Reserach projects

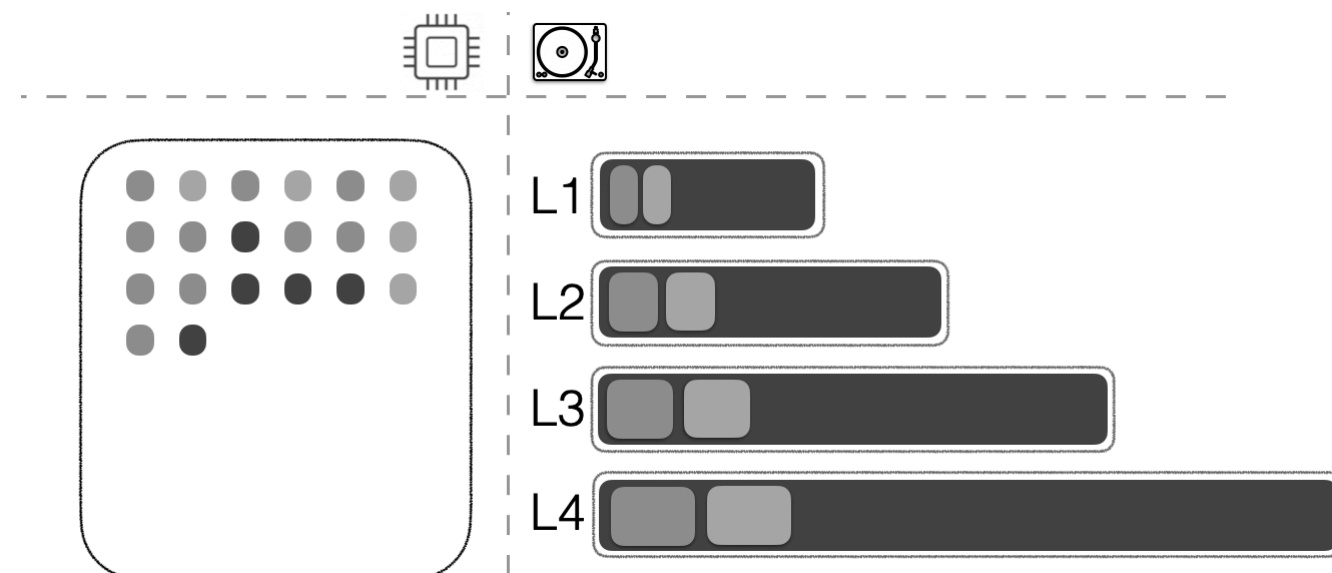
Solving problems on your way through!



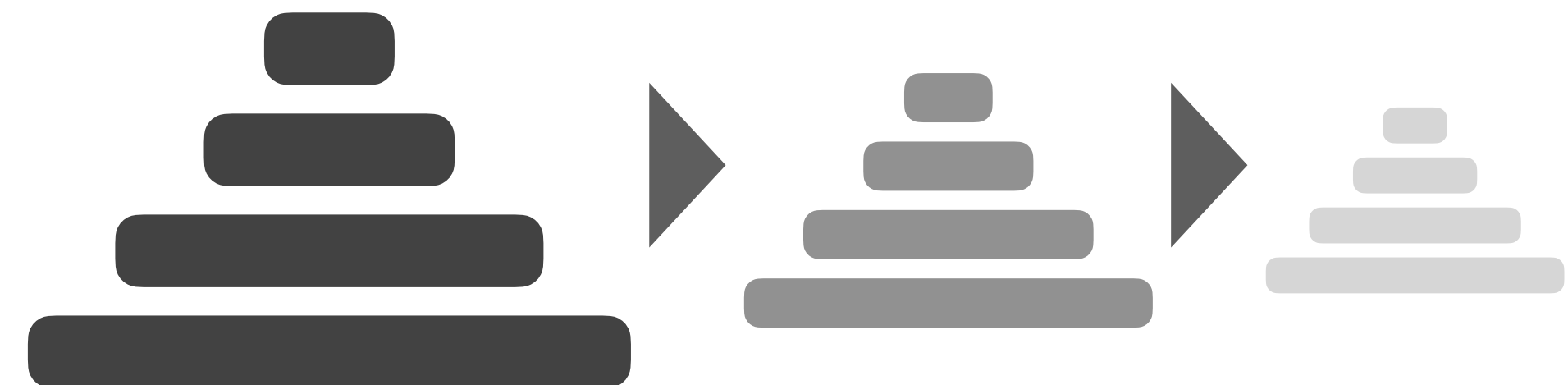
Self-designing LSM Buffer



Range Deletes in LSM



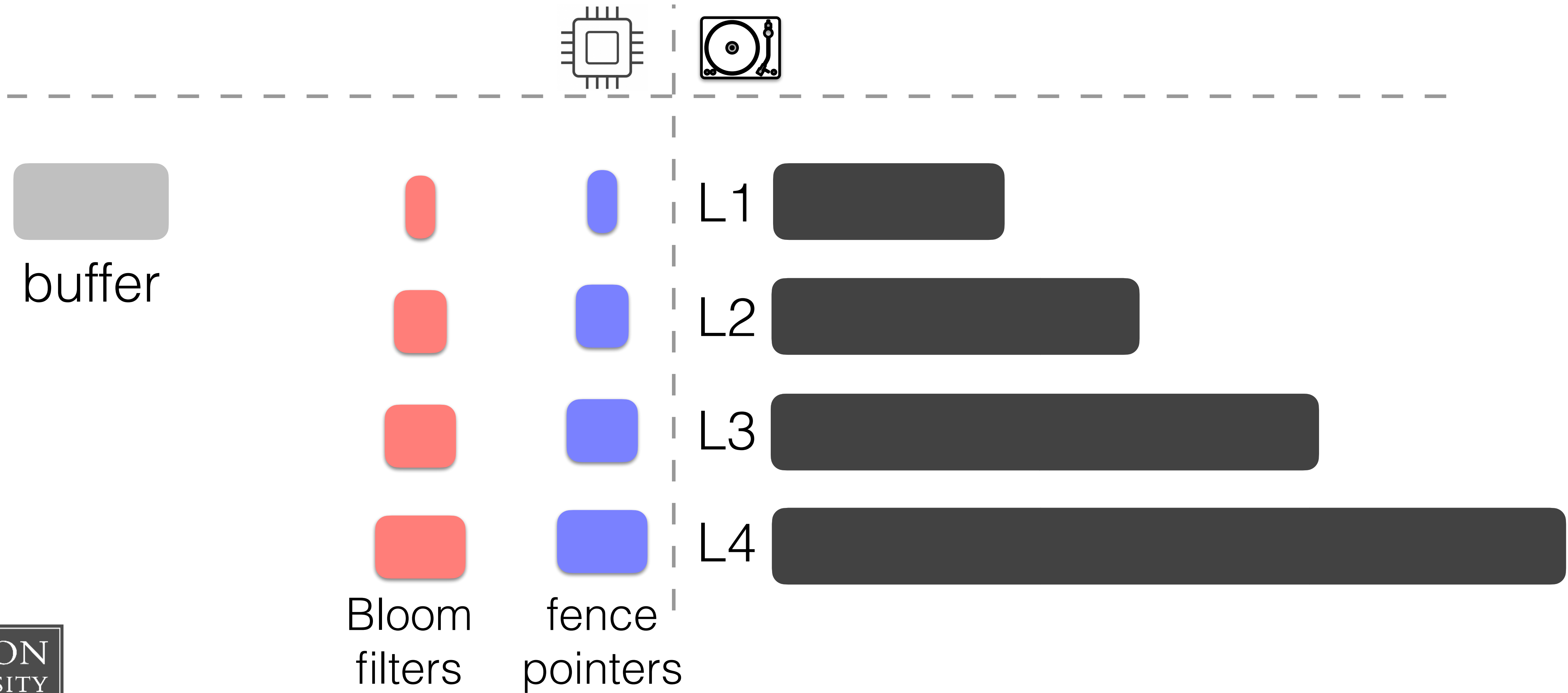
Optimal Caching in LSM



Time-Traveling LSM

# Block cache

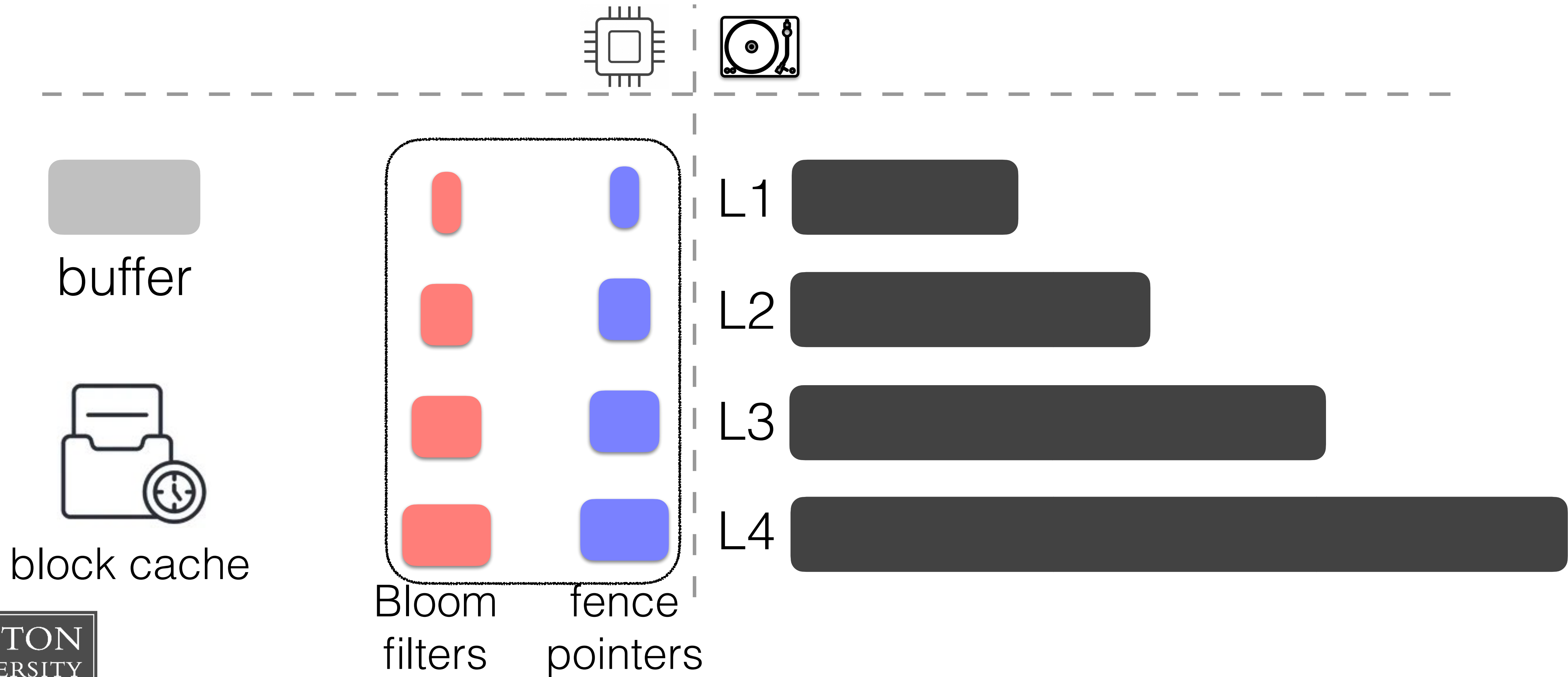
Saving trips to disk





# Block cache

Saving trips to disk

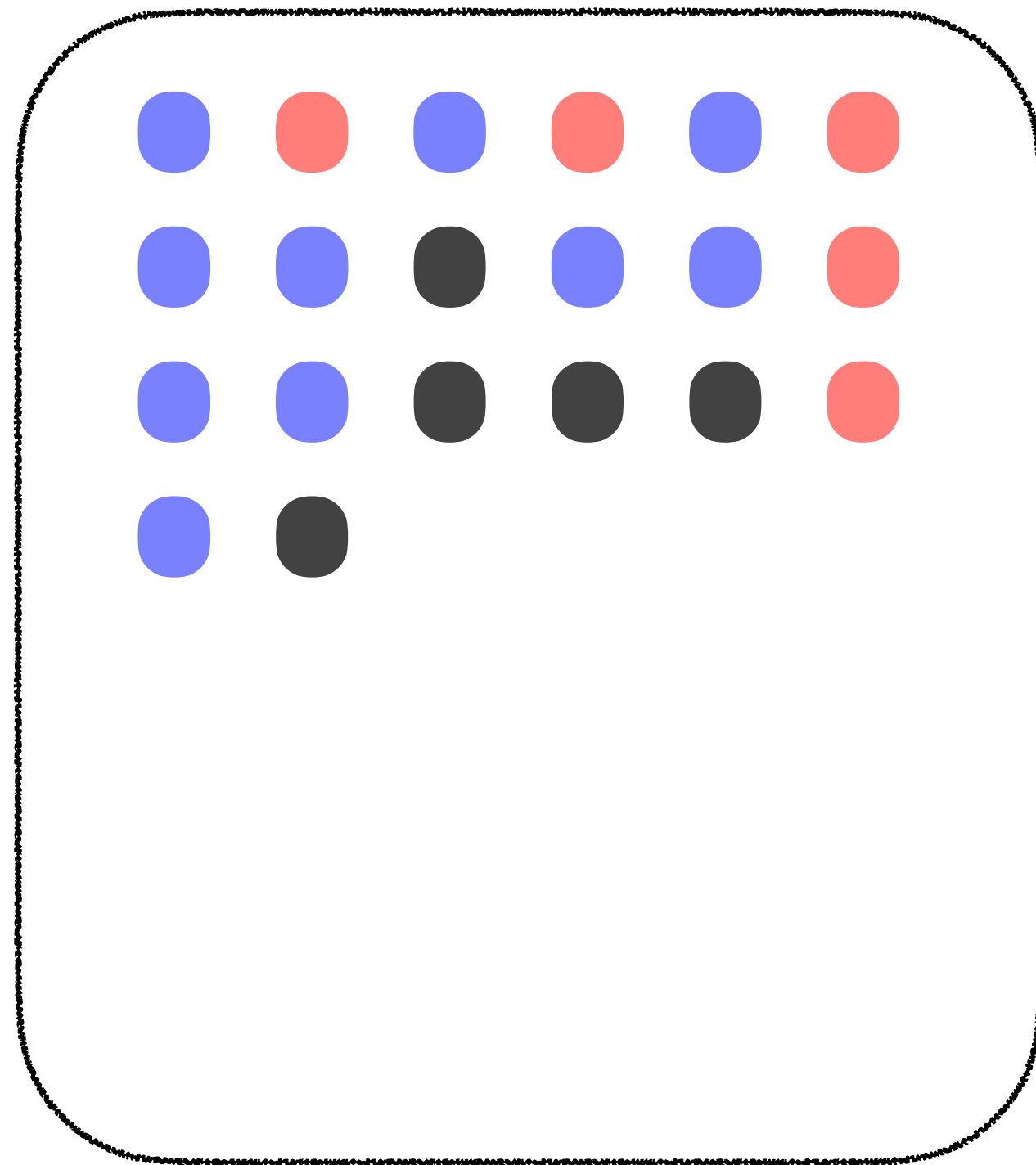


# Block cache

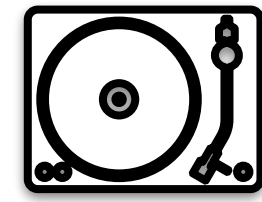
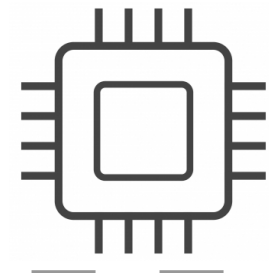
Saving trips to disk

get(7)

buffer

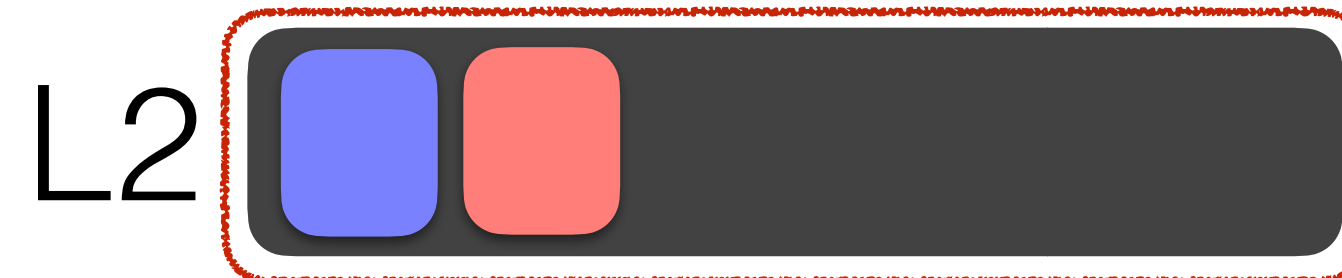
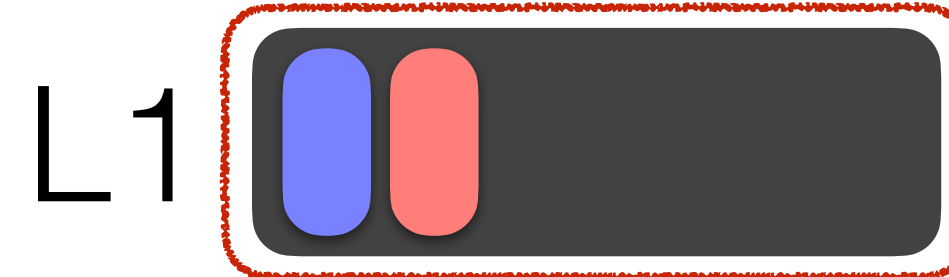


block cache



Bloom filters

fence pointers





# Research Project 2

Solving problems on your way through!

## Optimal Caching in LSM

**Memory is finite!**

How useful is **block-based caching**?

Which blocks to cache: **index, filter, data**?

What is the impact of **pinning blocks**?

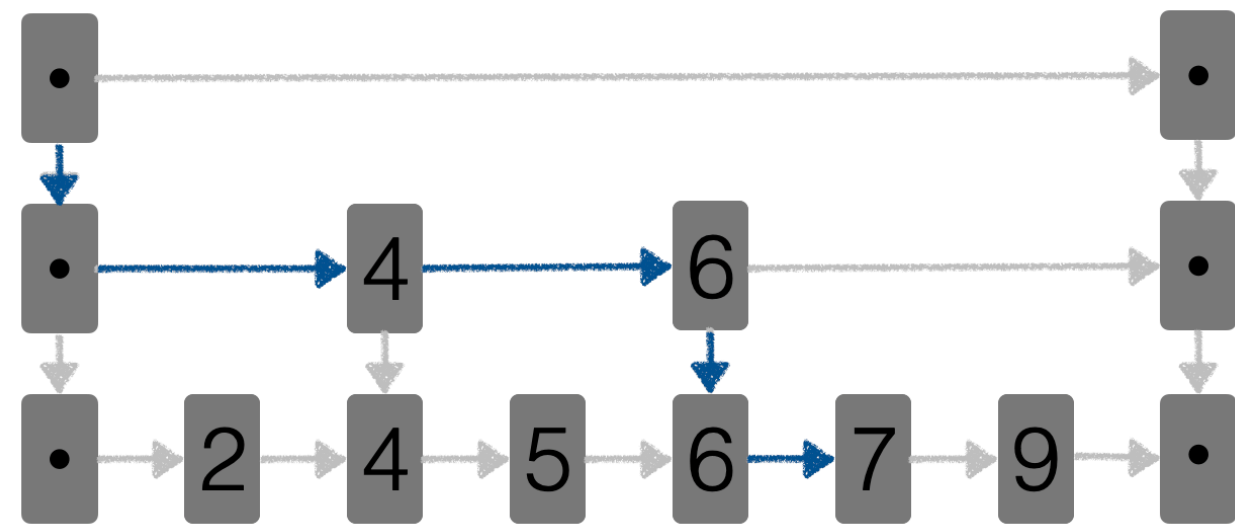
Which blocks to **evict**?

**Benchmark performance** for different caching schemes

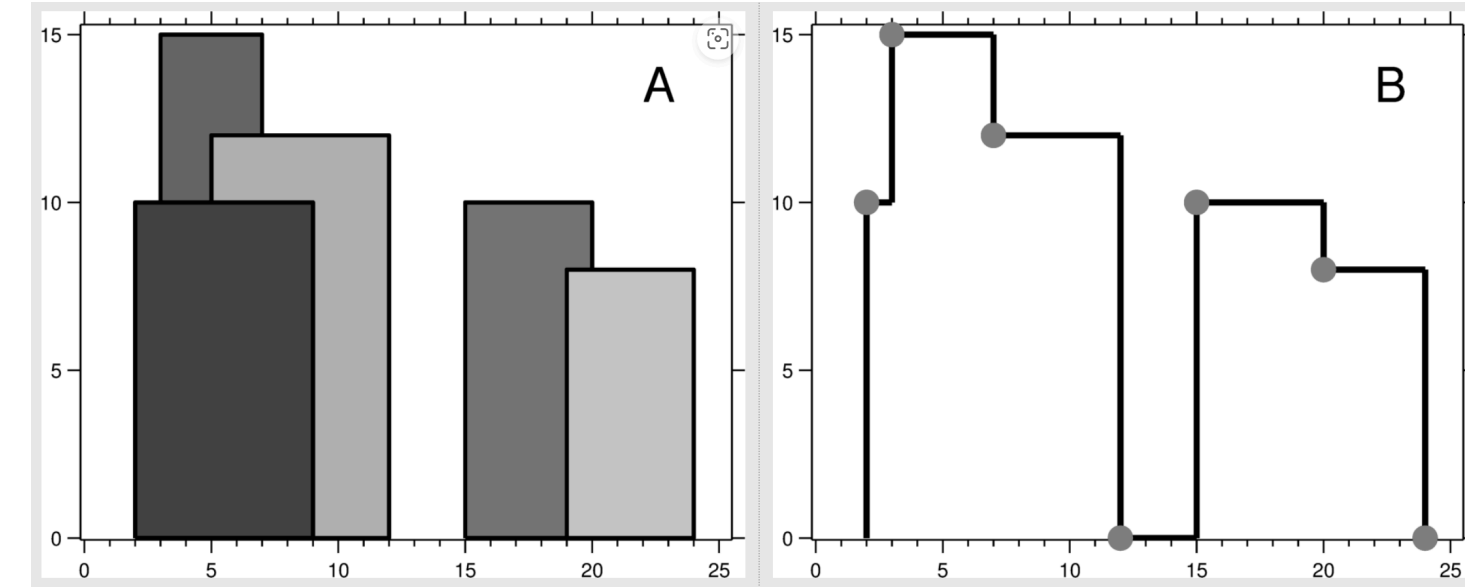


# Reserach projects

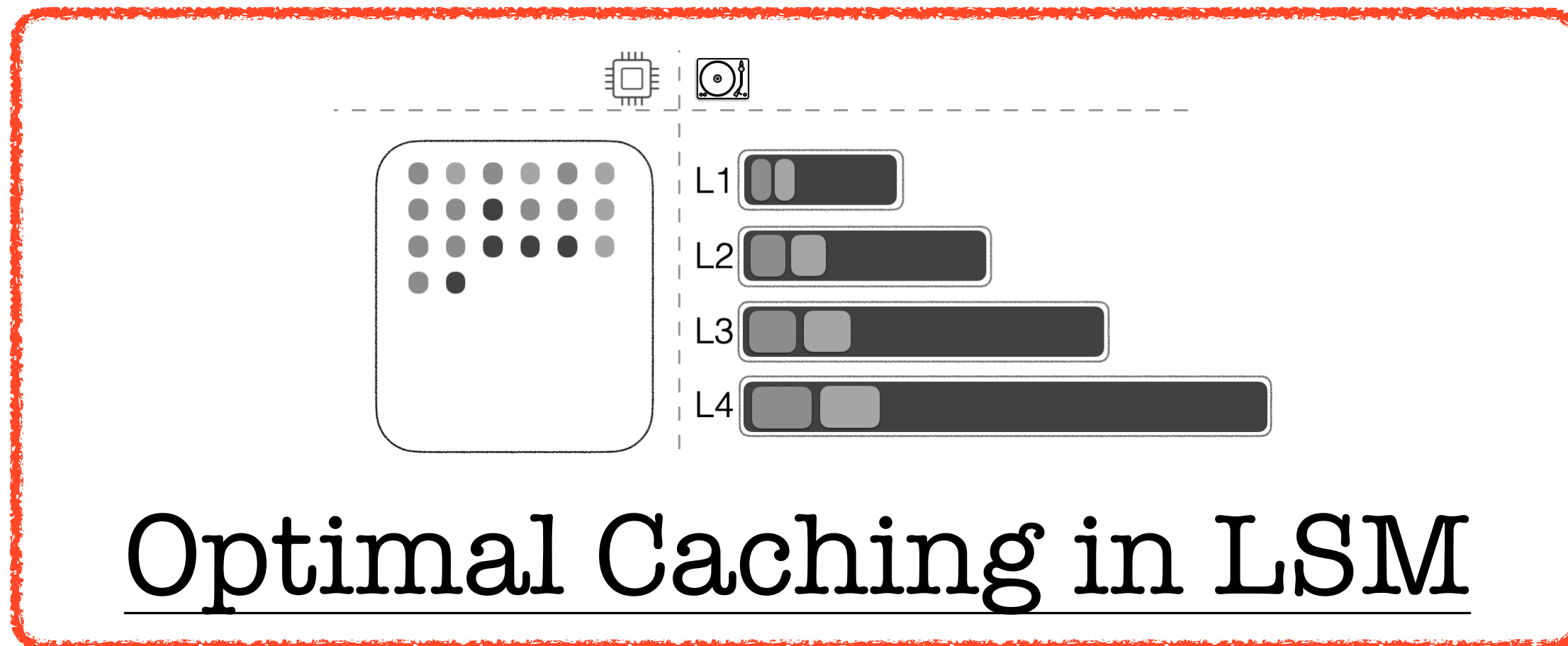
Solving problems on your way through!



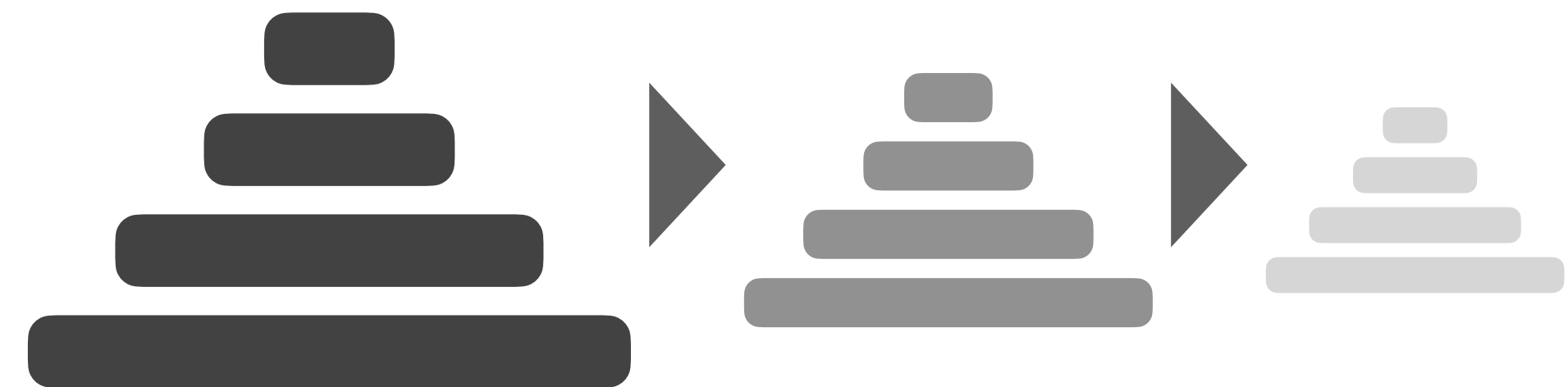
Self-designing LSM Buffer



Range Deletes in LSM



Optimal Caching in LSM



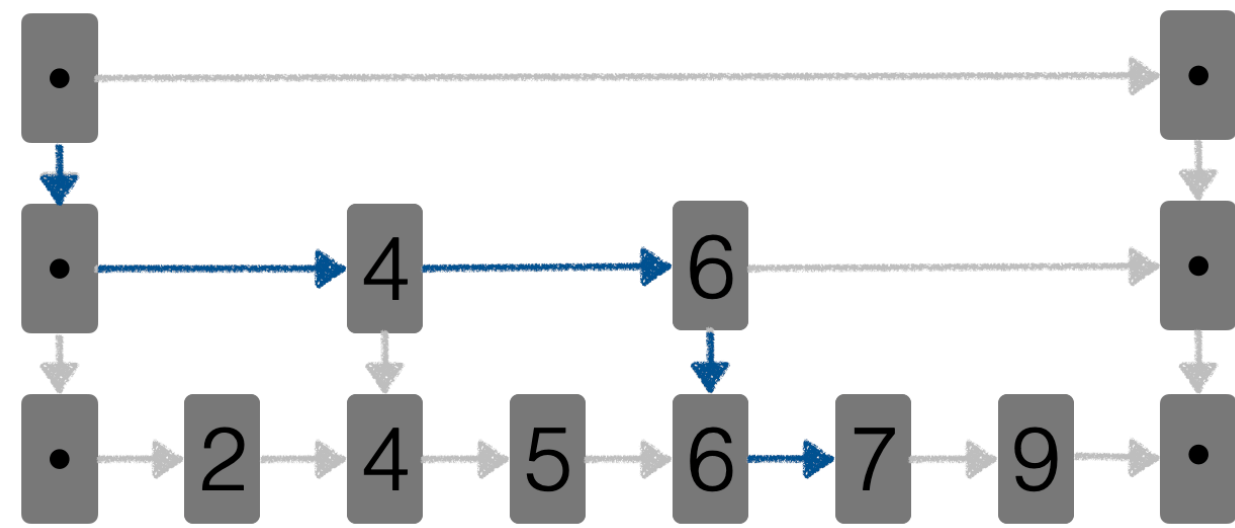
Time-Traveling LSM



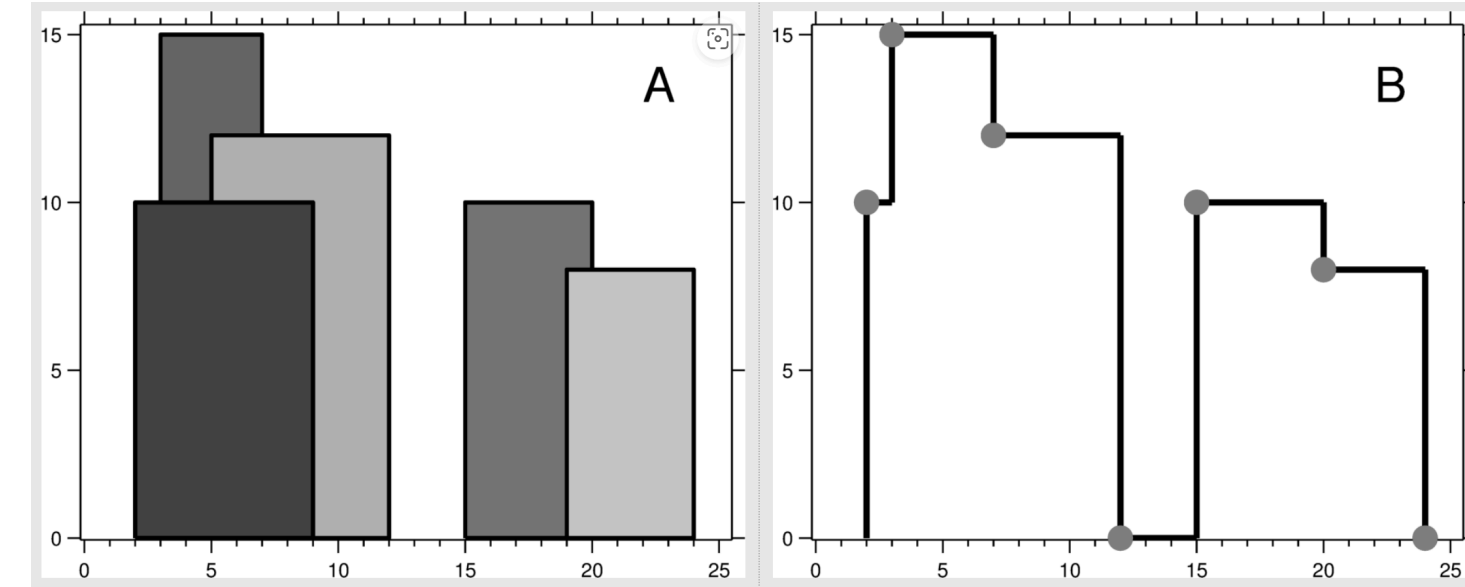


# Reserach projects

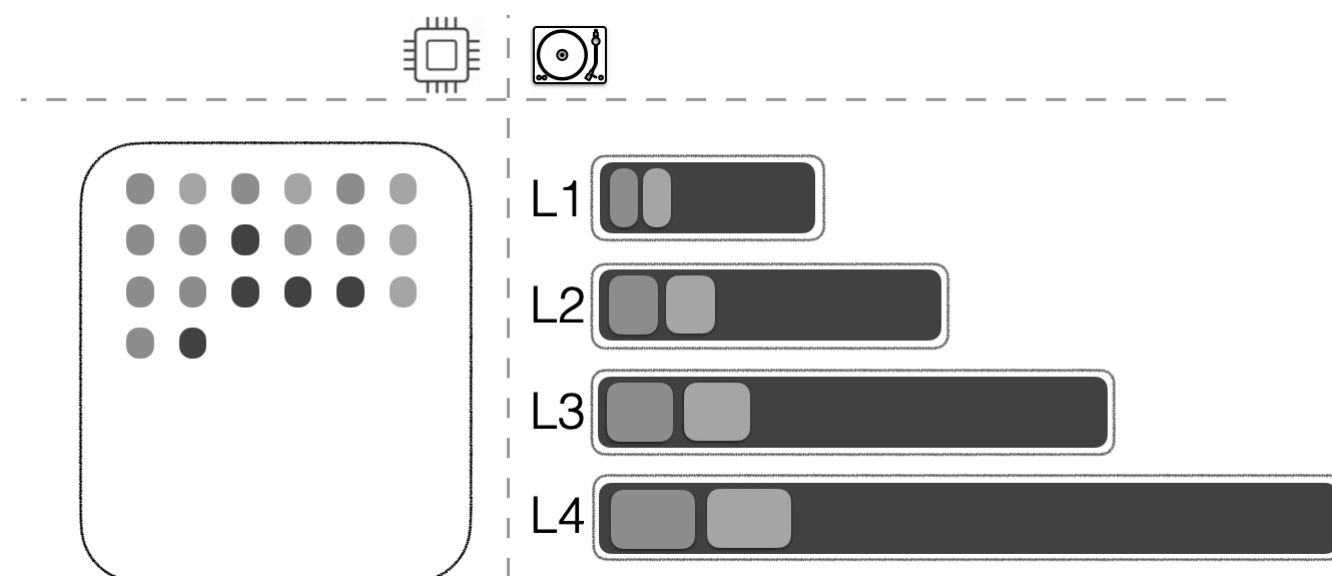
Solving problems on your way through!



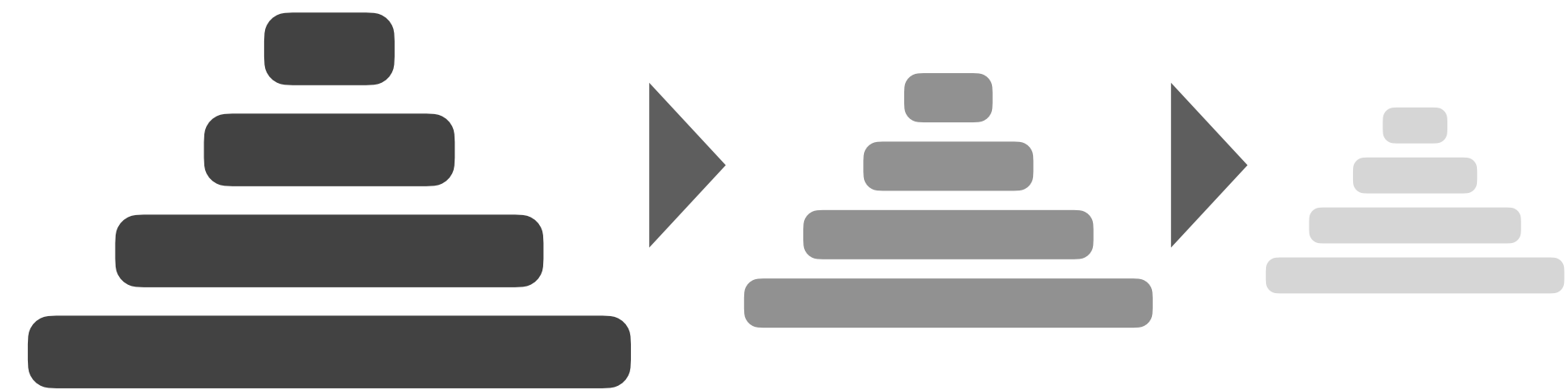
Self-designing LSM Buffer



Range Deletes in LSM



Optimal Caching in LSM

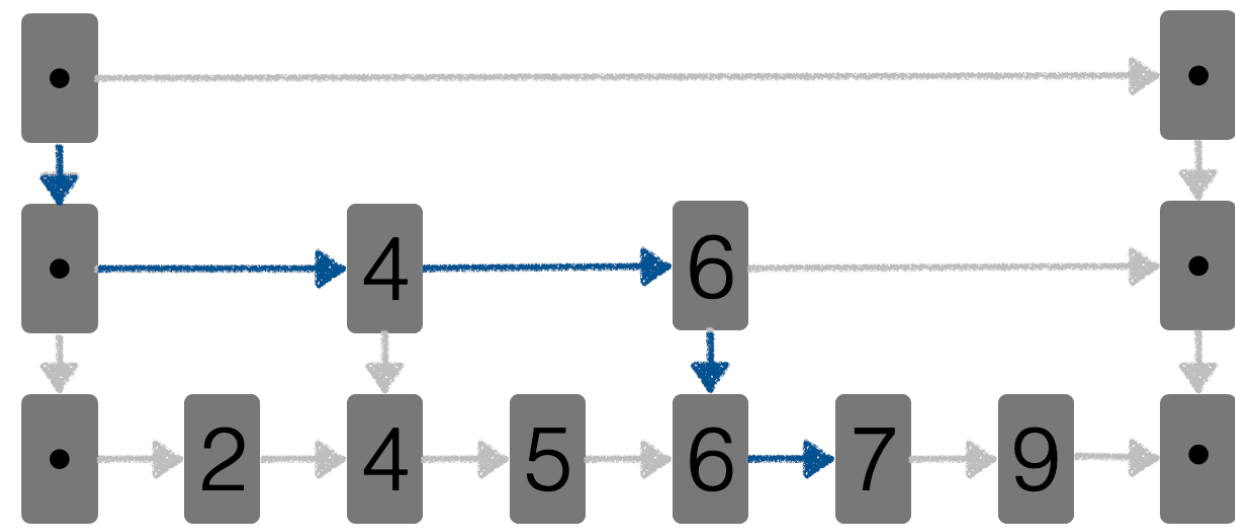


Time-Traveling LSM

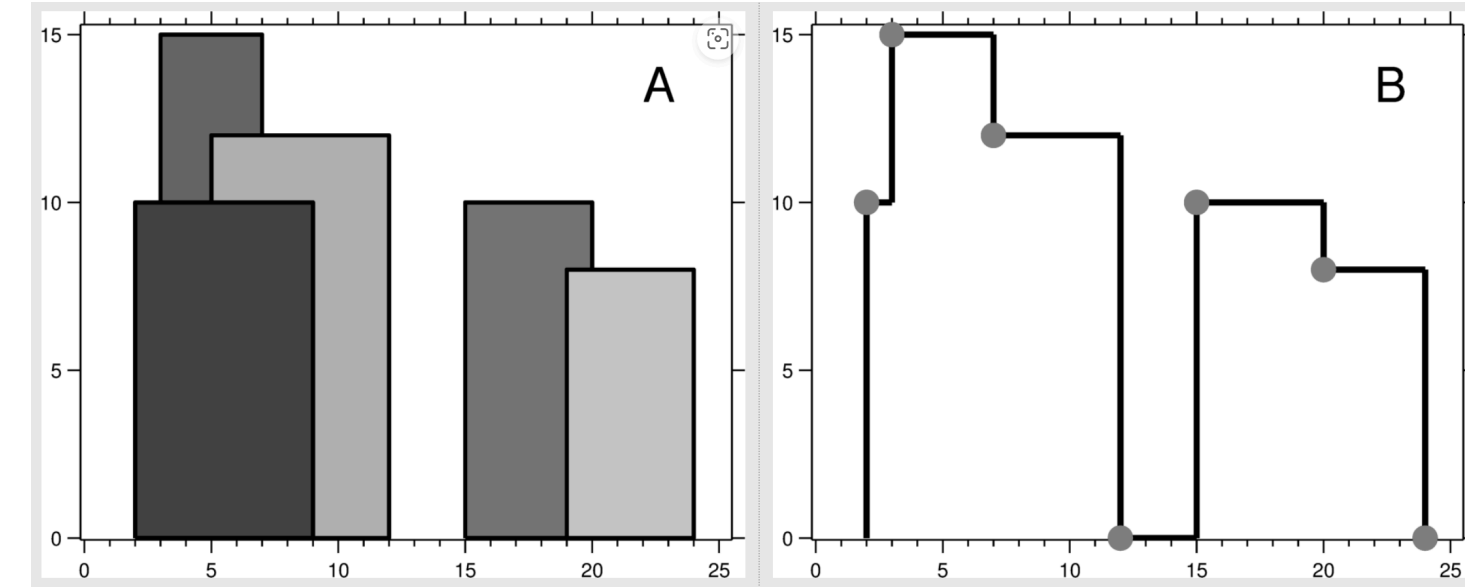


# Reserach projects

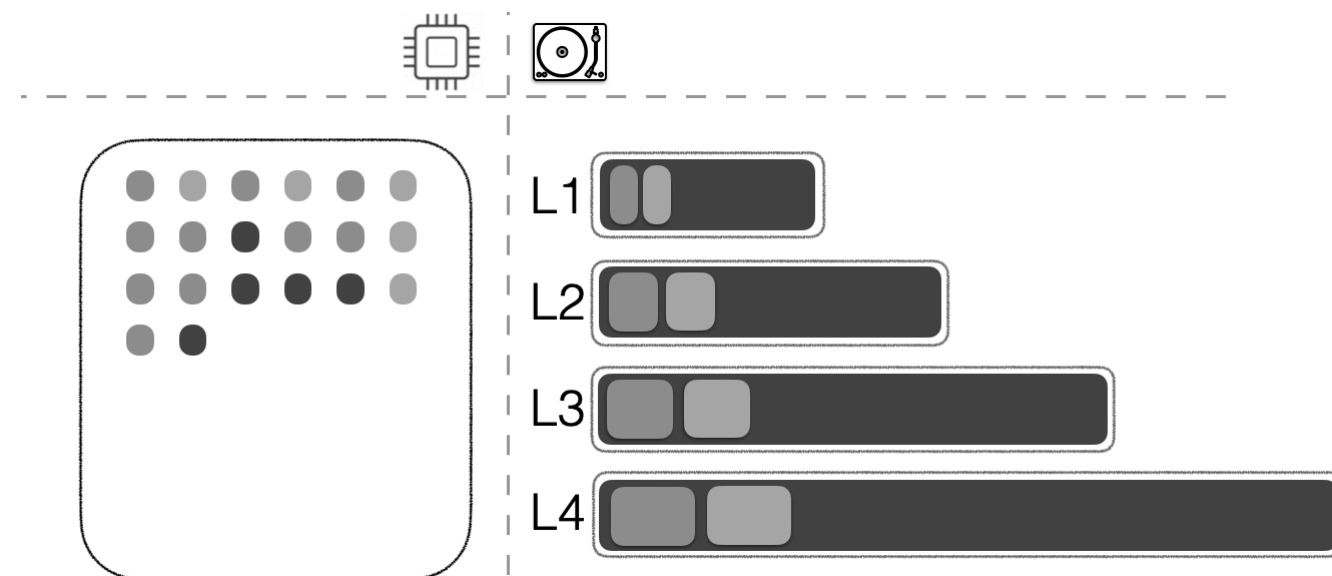
Solving problems on your way through!



Self-designing LSM Buffer



Range Deletes in LSM



Optimal Caching in LSM



Time-Traveling LSM

# Next time in COSI 167A

More on LSMs

## LSM buffer design space

[P] ["Anatomy of the LSM Memory Buffer: Insights & Implications"](#), *DBTest*, 2024

**TECHNICAL QUESTION 3** [How does the memory allocation strategy affect the query latency for a pre-allocated vector vs. a dynamically allocated one? Given a workload, how would you decide the prefix length for a hash-hybrid buffer implementation?](#)

[B] ["Breaking Down Memory Walls: Adaptive Memory Management in LSM-based Storage Systems"](#), *VLDB*, 2020



# COSI 167A

# Advanced Data Systems

Class 7

# Class Projects

Prof. Subhadeep Sarkar

<https://ssd-brandeis.github.io/COSI-167A/>