

# CS 561: Data Systems Architectures

class 4

Systems & Research Project

Prof. Manos Athanassoulis

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

# Let's revisit Zonemaps

- Light-weight auxiliary data structure (*“scan accelerator”*)

# Let's revisit Zonemaps

zonemaps

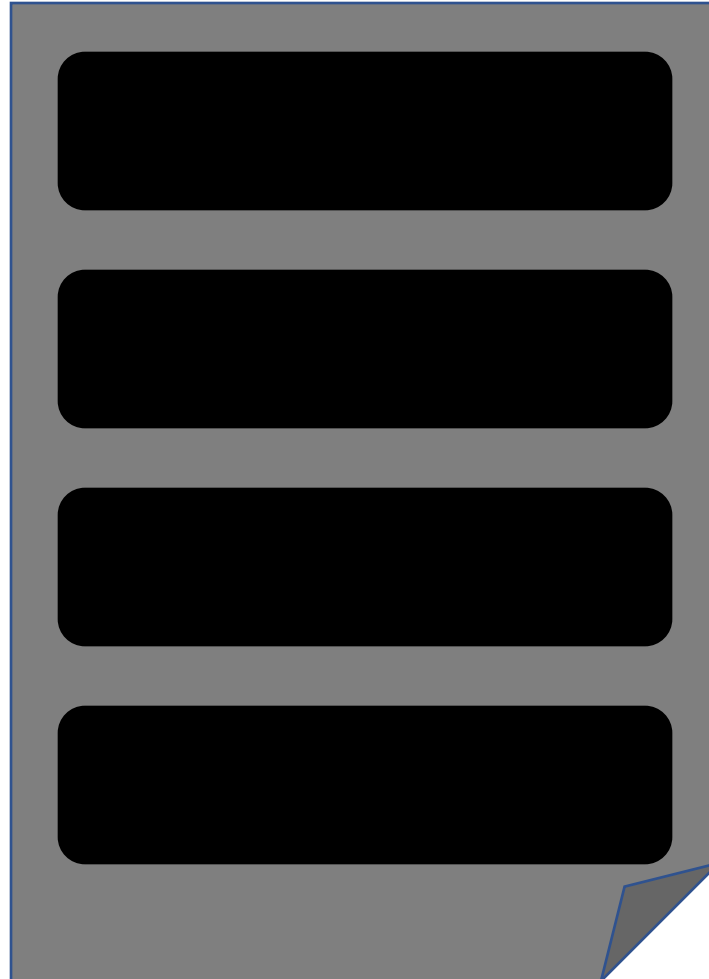
file = collection of pages

page 0

page 1

page 2

page 3



# zonemaps

file = collection of pages

page 0

3, 16, 34, 31, 21

page 1

1, 5, 12, 24, 23

page 2

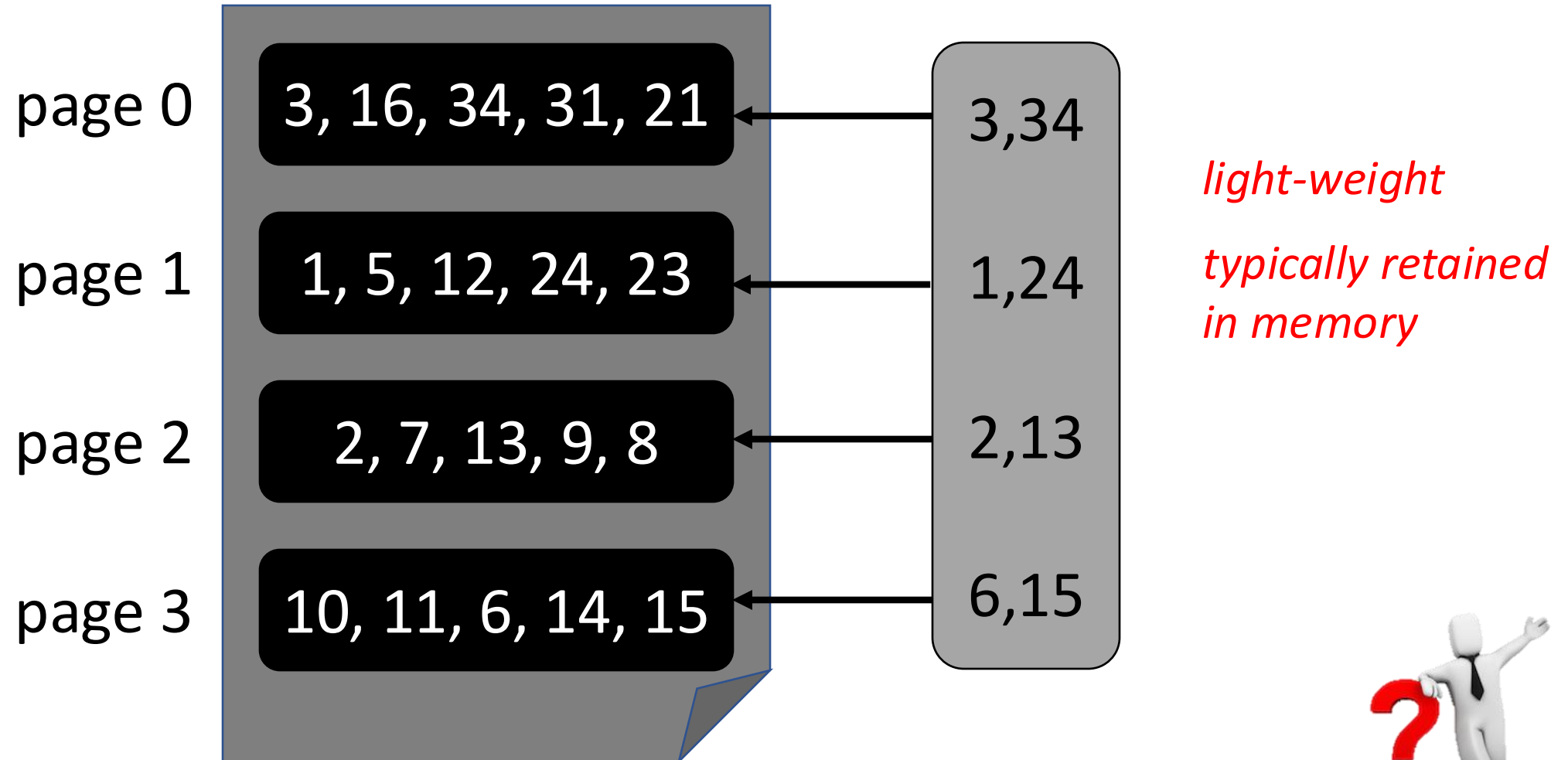
2, 7, 13, 9, 8

page 3

10, 11, 6, 14, 15

# zonemaps

file = collection of pages

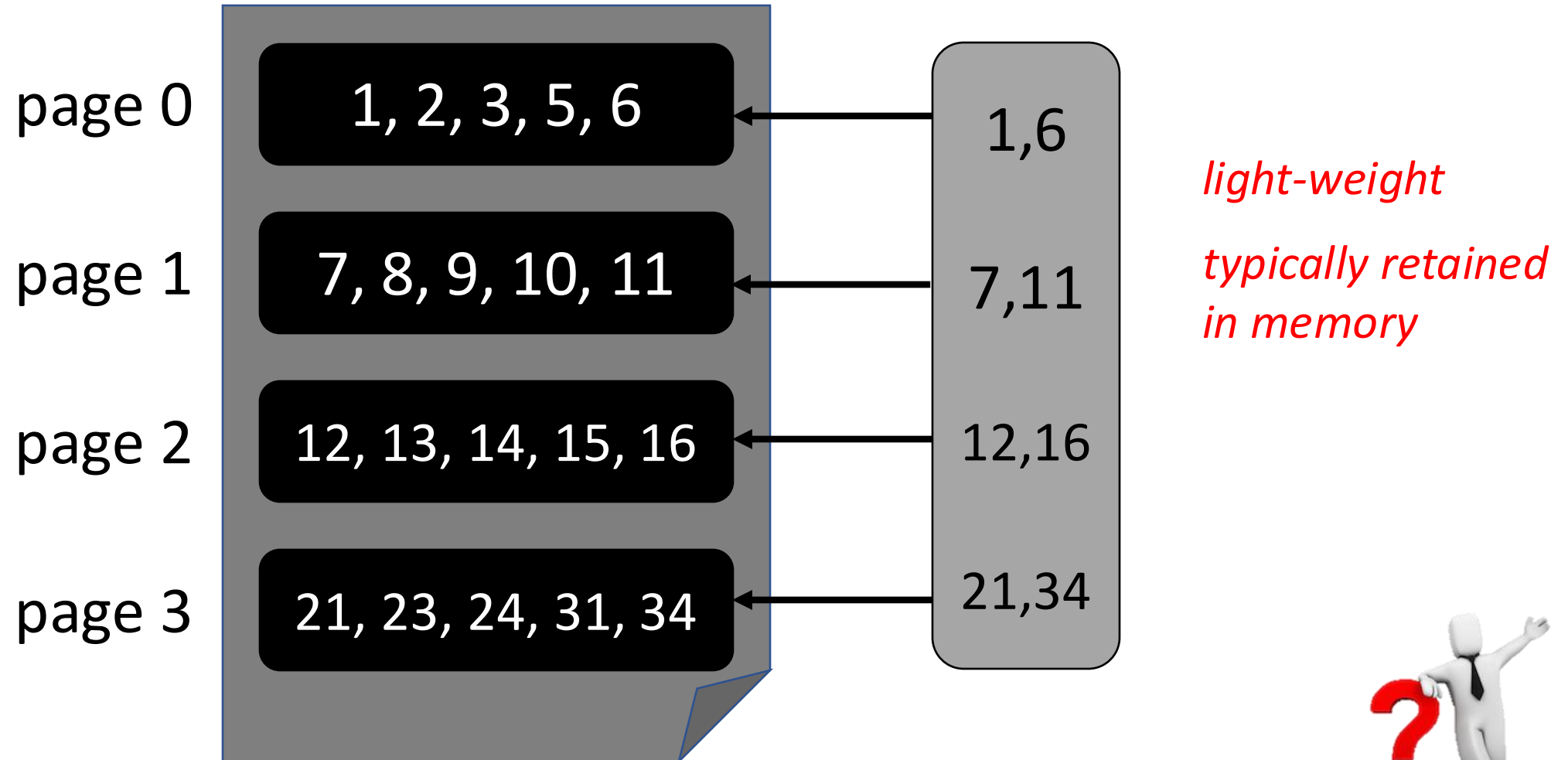


But what if the data is sorted?



# zonemaps

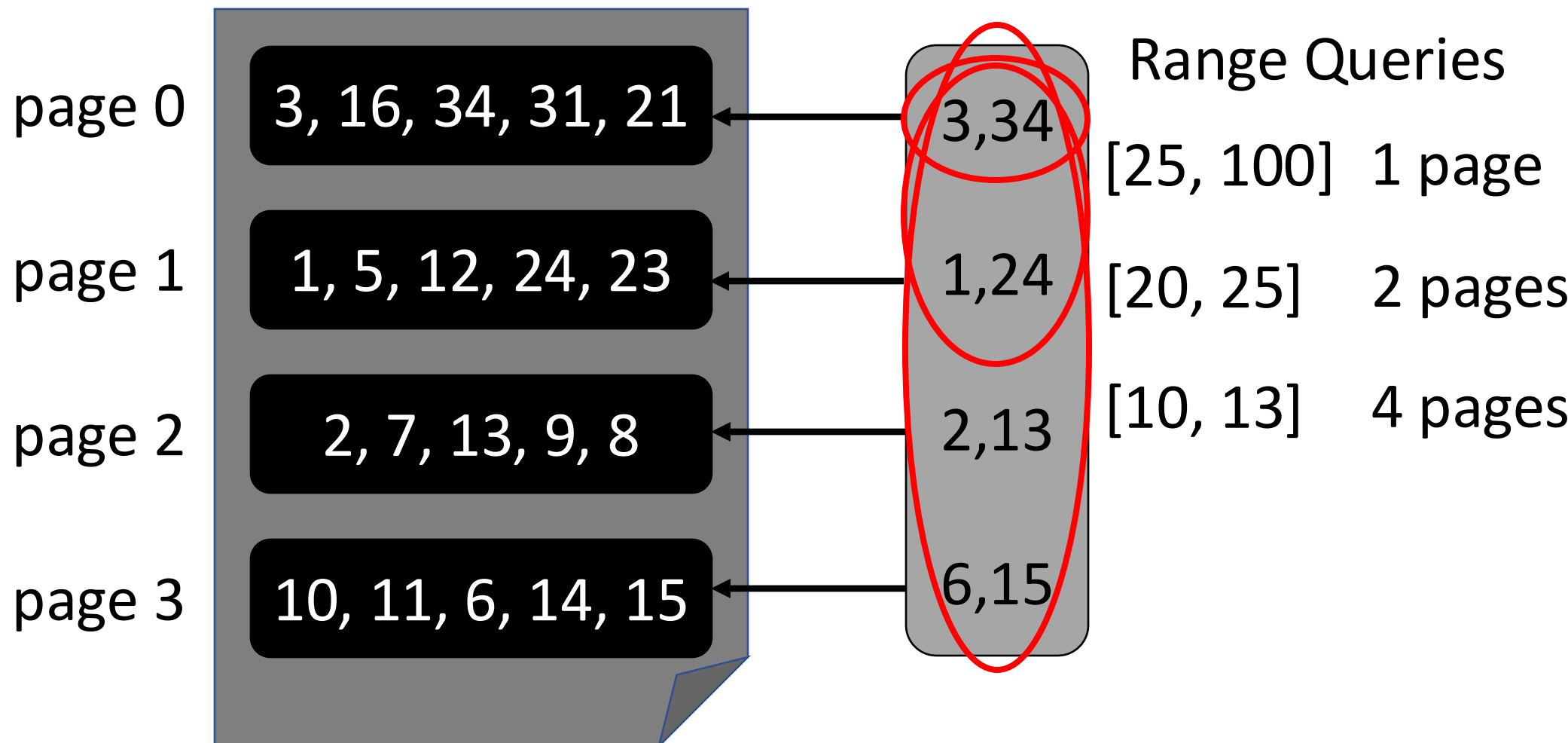
file = collection of pages



But what if the data is sorted?

# zonemaps

file

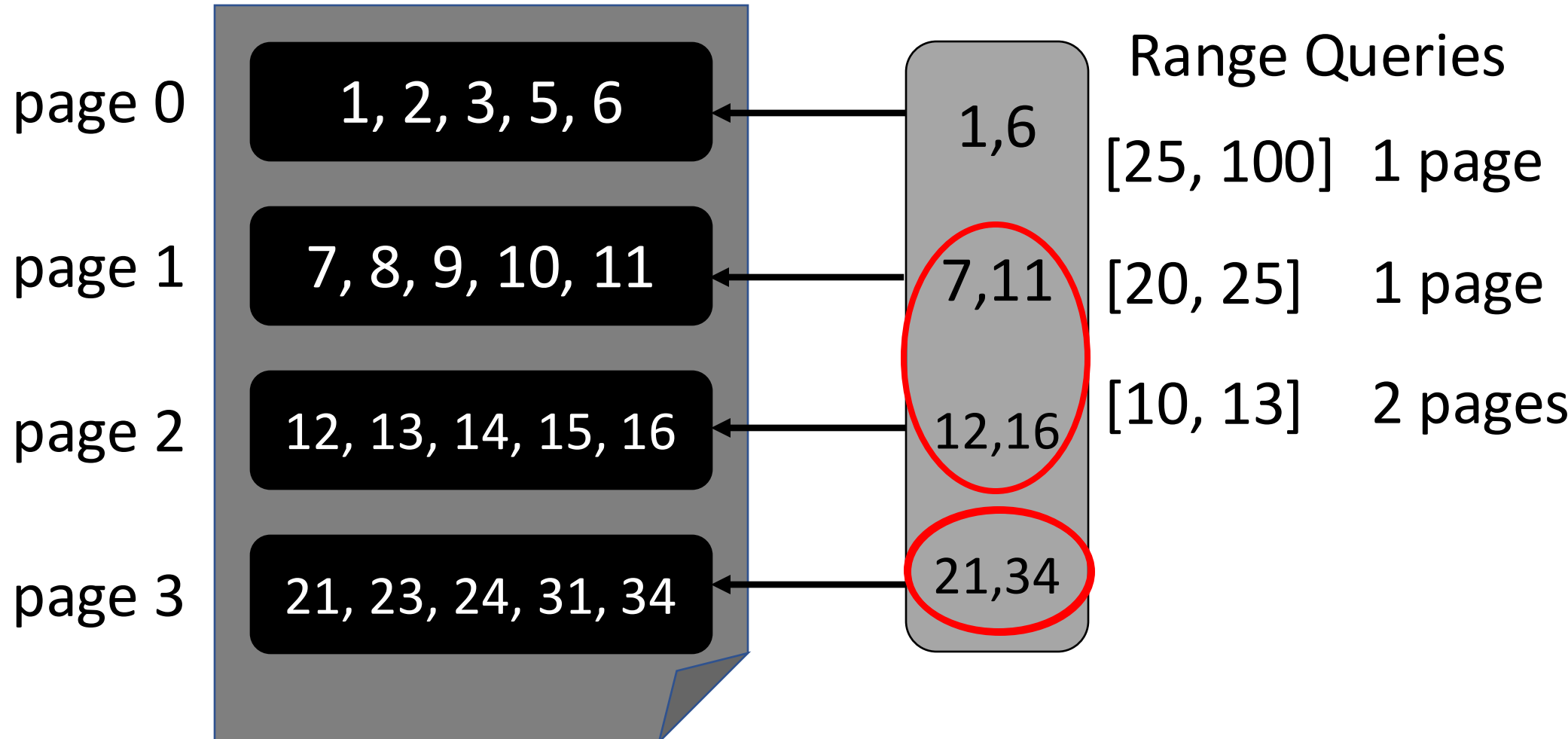


But what if the data is sorted?



zonemaps

file



Zonemaps efficiency depends on data & queries!



# data systems



®



complex analytics

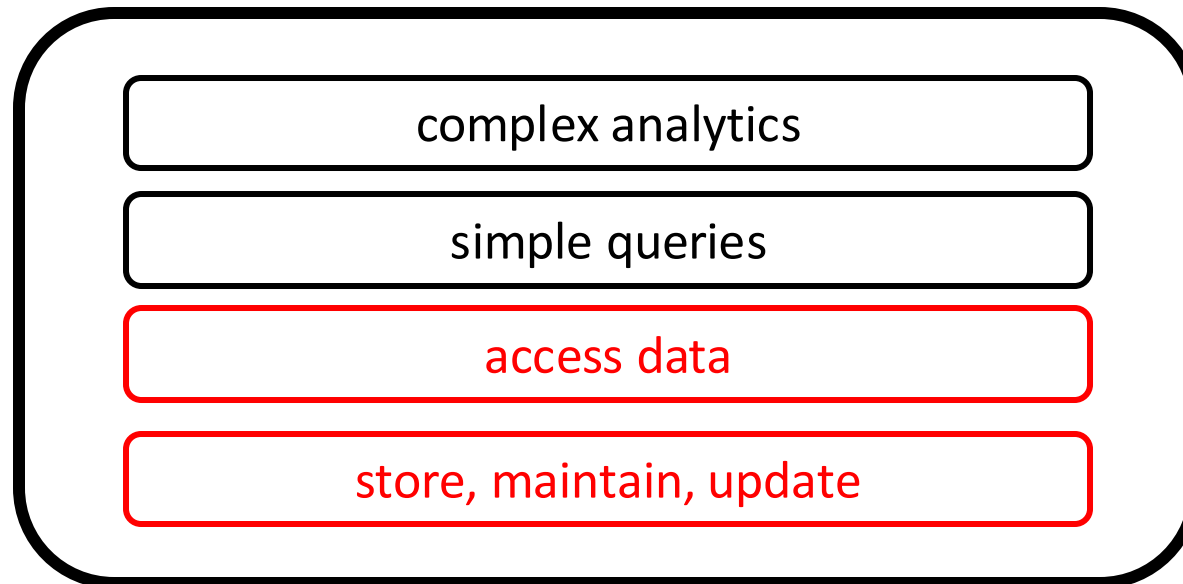
simple queries

access data

store, maintain, update

# data systems

>\$200B by 2020, growing at 11.7% every year  
[The Forbes, 2016]



***access methods\****

**\*algorithms and data structures**  
for organizing and accessing data

# data systems core: storage engines

## main decisions

how to ***store*** data?

how to ***access*** data?

how to ***update*** data?

let's simplify: **key-value** storage engines

collection of keys-value pairs

query on the key, return both key and value

*remember*



*state-of-the-art* design

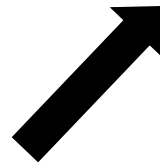
# how general is a key value store?

can we store relational data?



yes! {<primary\_key>, <rest\_of\_the\_row>}

example: { **student\_id**, { **name**, **login**, **yob**, **gpa** } }



what is the caveat?

how to index these attributes?

other problems?

index: { **name**, { **student\_id** } }

index: { **yob**, { **student\_id<sub>1</sub>**, **student\_id<sub>2</sub>**, ... } }

# how general is a key value store?

can we store relational data?



yes!    {<primary\_key>,<rest\_of\_the\_row>}

how to efficiently code if we do not know  
the structure of the “*value*”

index: { **yob**, { **student\_id<sub>1</sub>**, **student\_id<sub>2</sub>**, ... } }

# how to use a key-value store?

## *basic interface*

`put(k,v)`

`{v} = get(k)`      `{v1, v2, ...} = get(k)`

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

`c = count(kmin, kmax)`

*deletes: delete(k)*

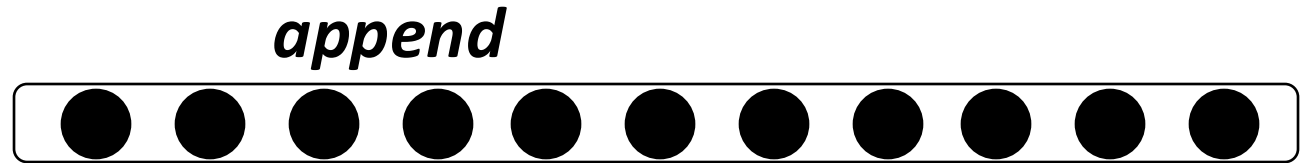
*updates: update(k,v)*    is it different than put?

get set: `{v1, v2, ...} = get_set(k1, k2, ...)`



# how to build a key-value store?

if we have only **put** operations



if we mostly have **get** operations

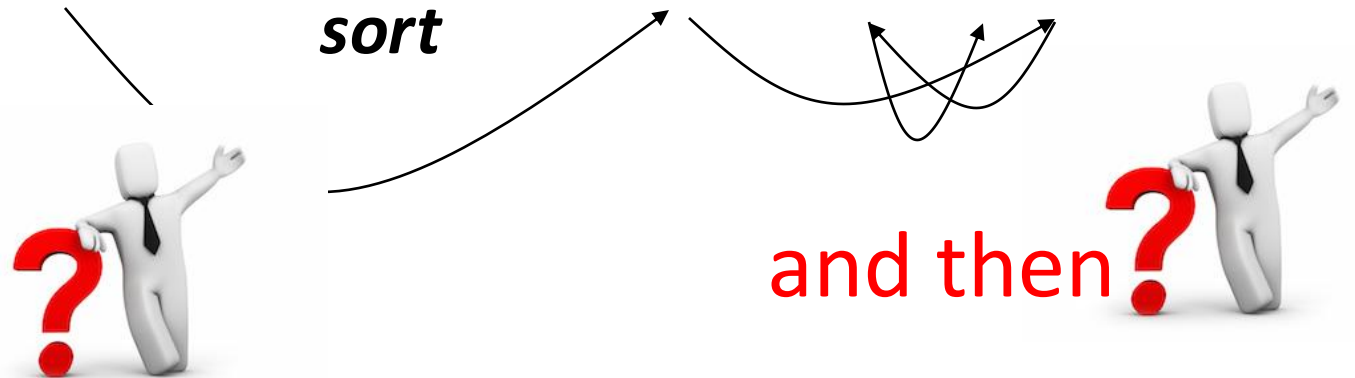


*sort*

and then?

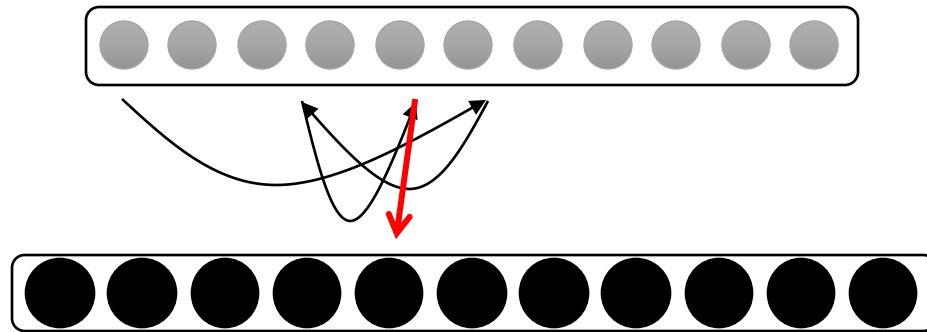
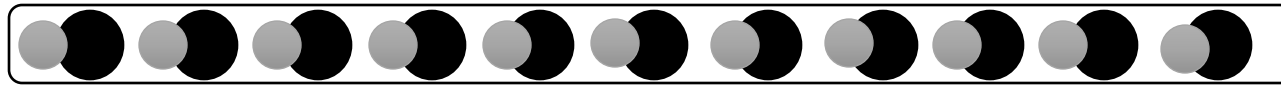
what about full scan?

range queries?





# can we separate keys and values?



at what price?



locality? code?

read queries  
(point or range)



inserts  
(or updates)

*sort data*

*simply append*

*amortize sorting cost*

*avoid resorting after every update*

how to bridge?



# LSM-tree

## Key-Value Stores

What are they really?

updates

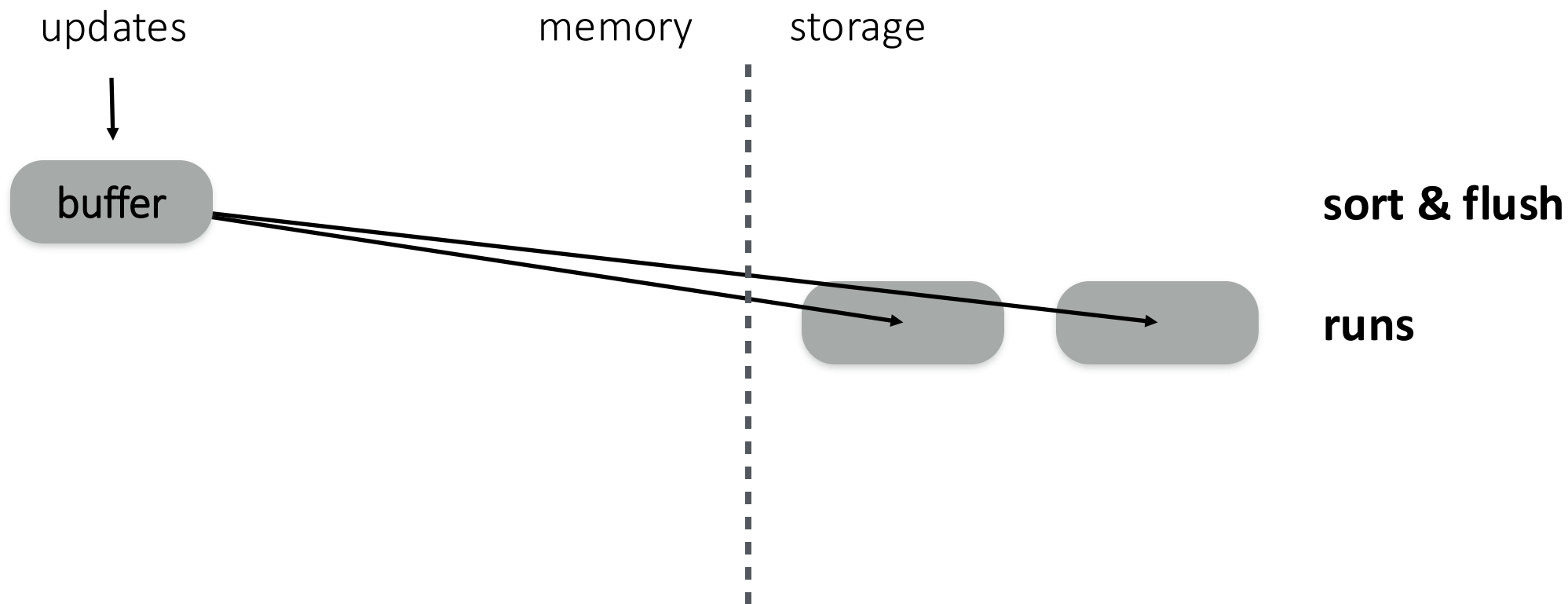


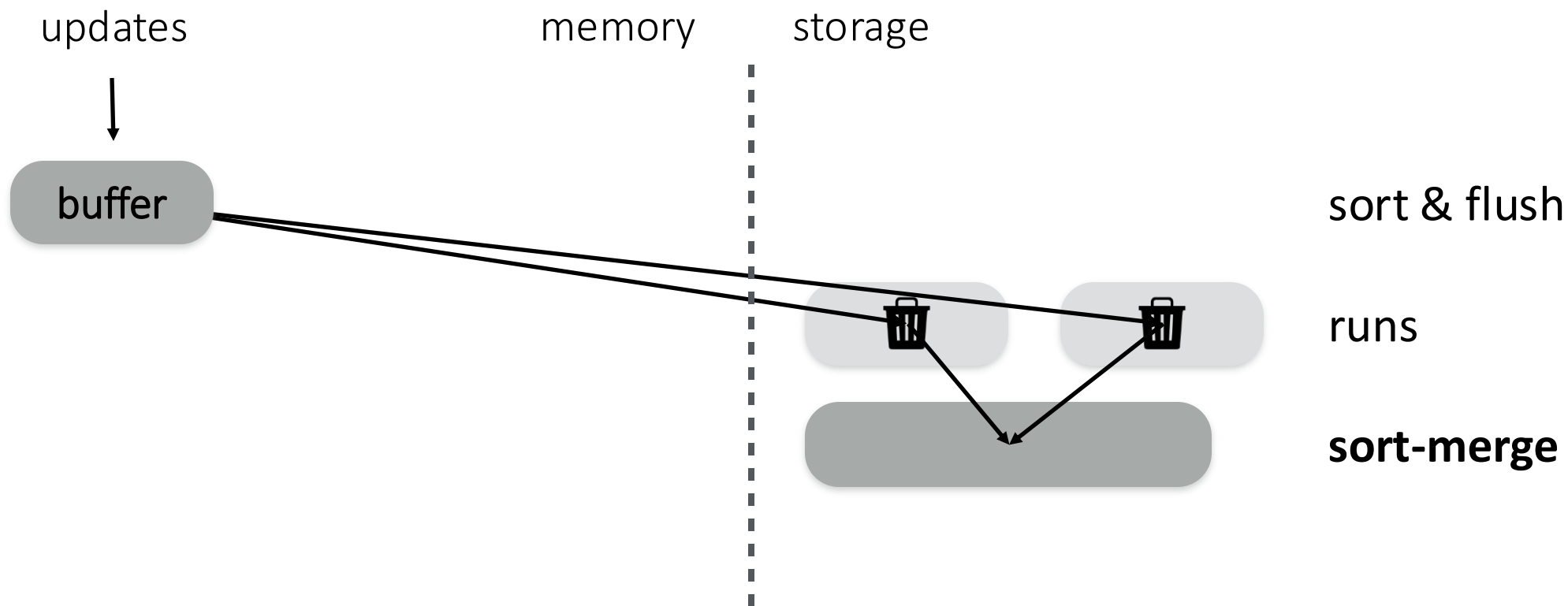
buffer

memory

storage







buffer

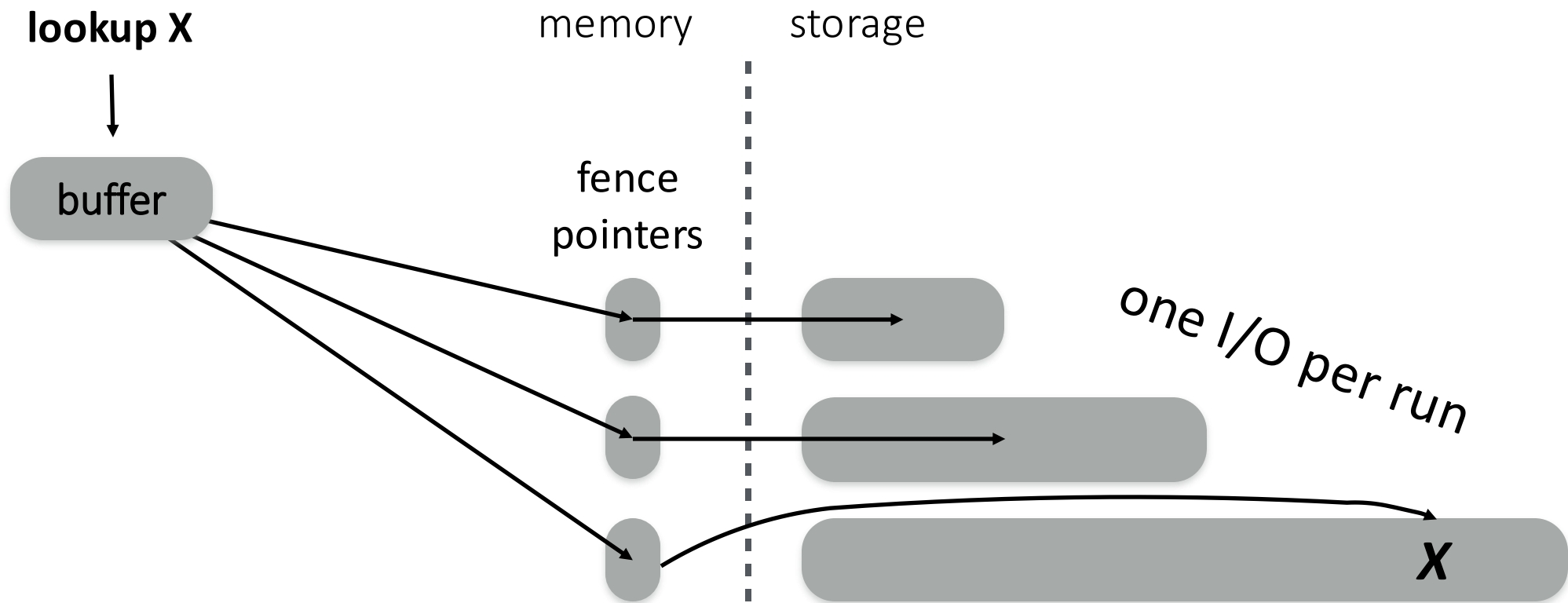
memory

storage

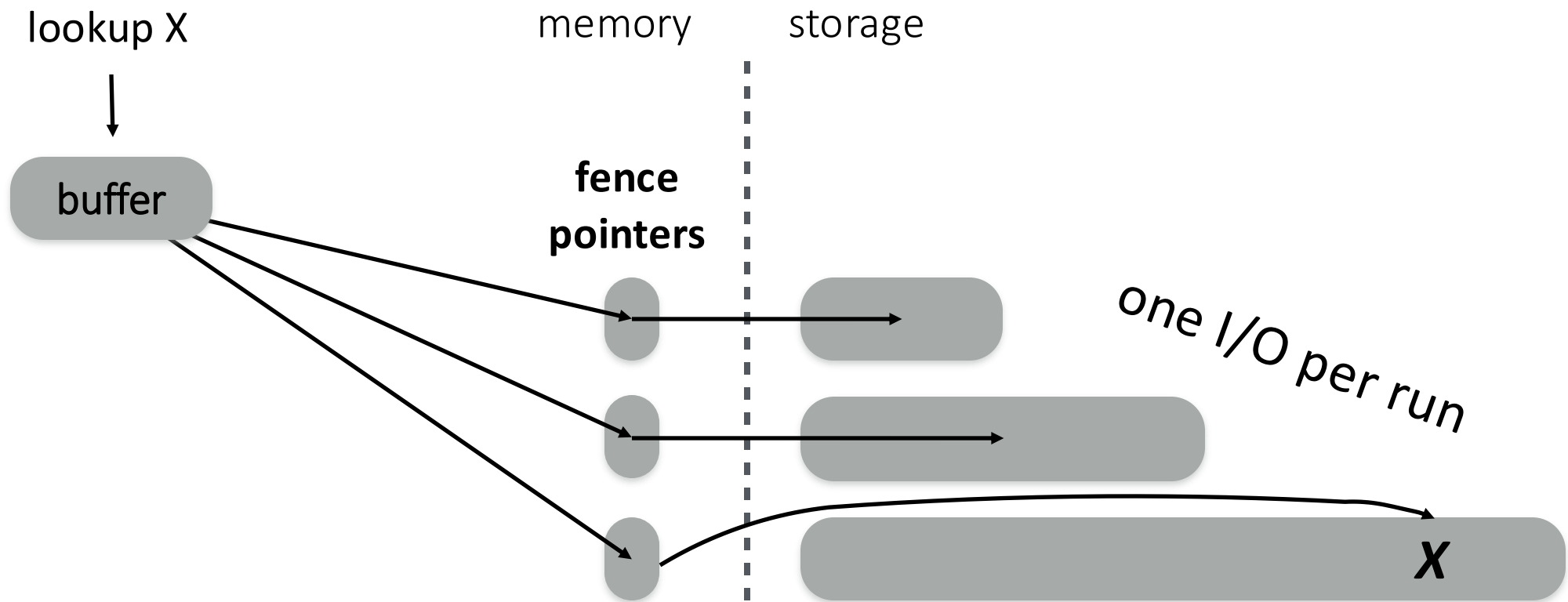
exponentially increasing sizes

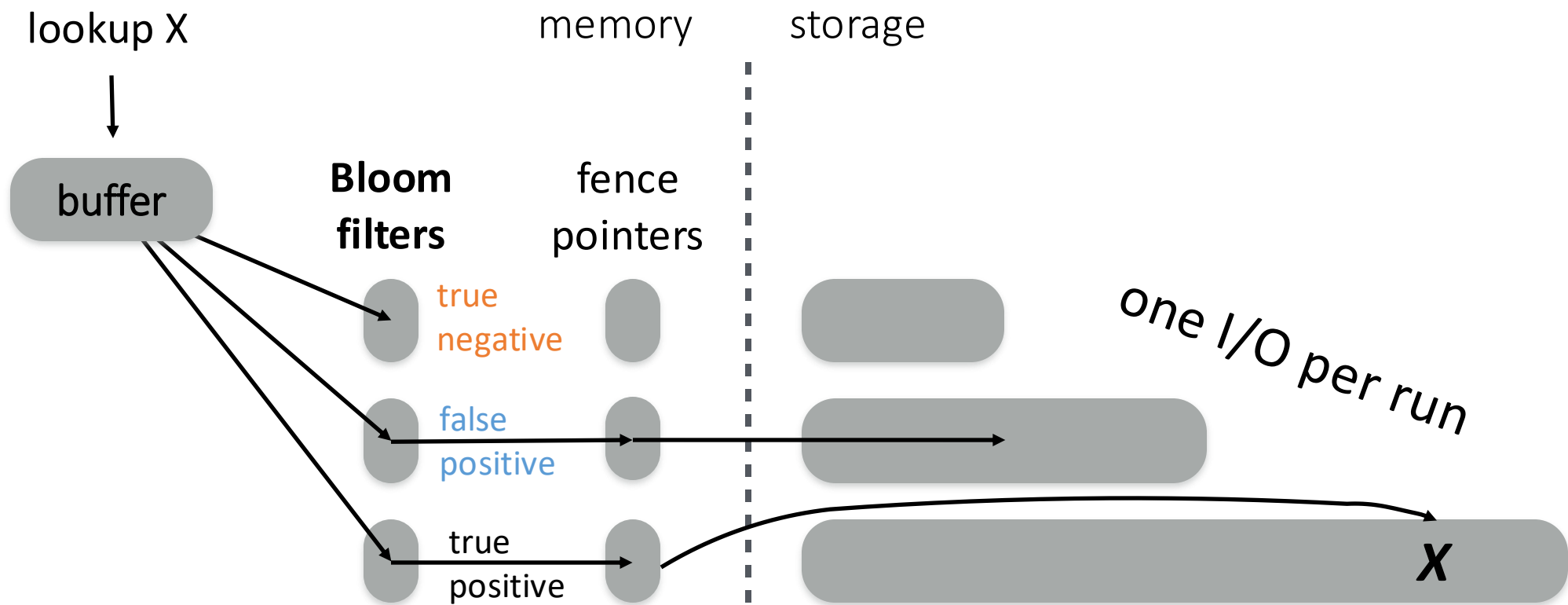
$O(\log(N))$  levels



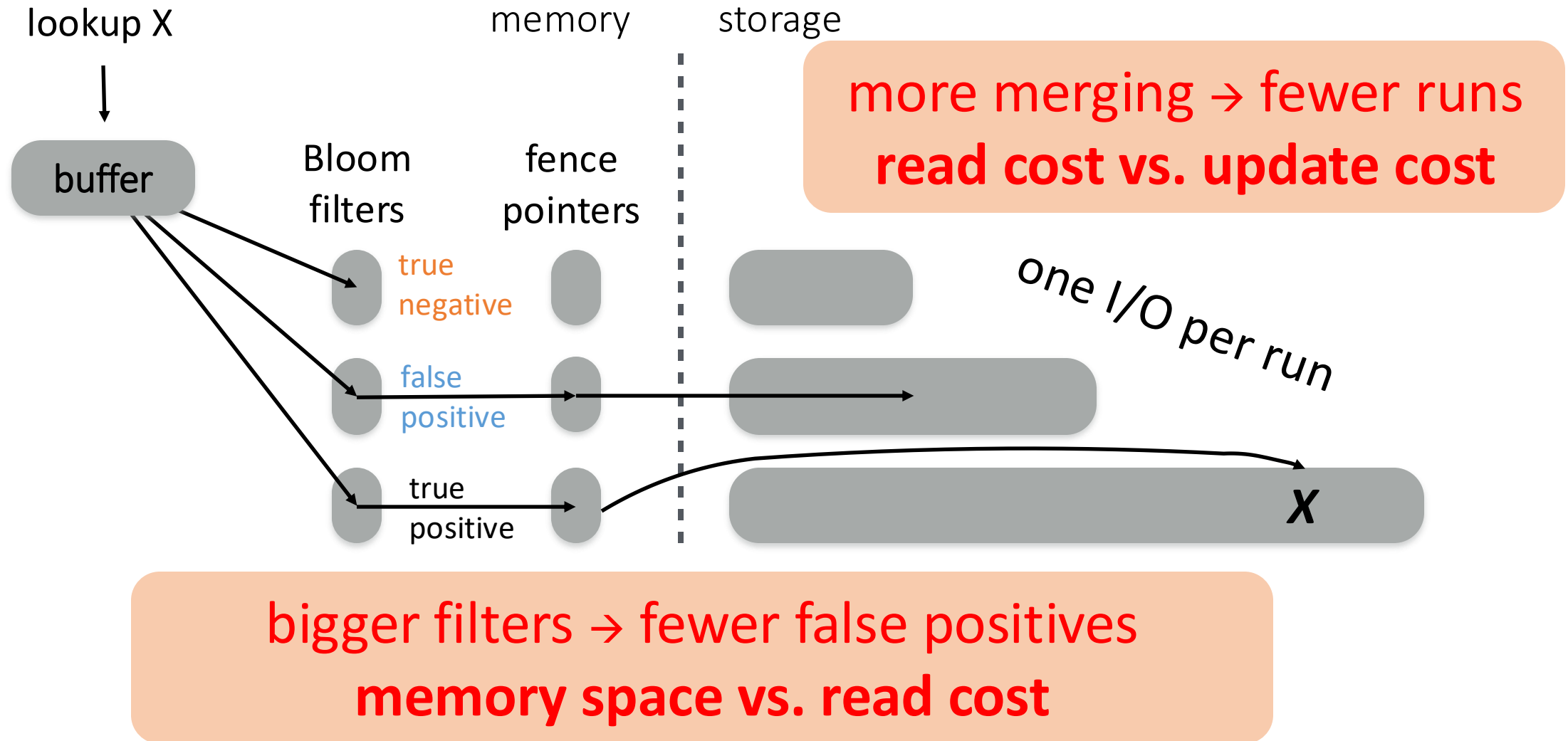




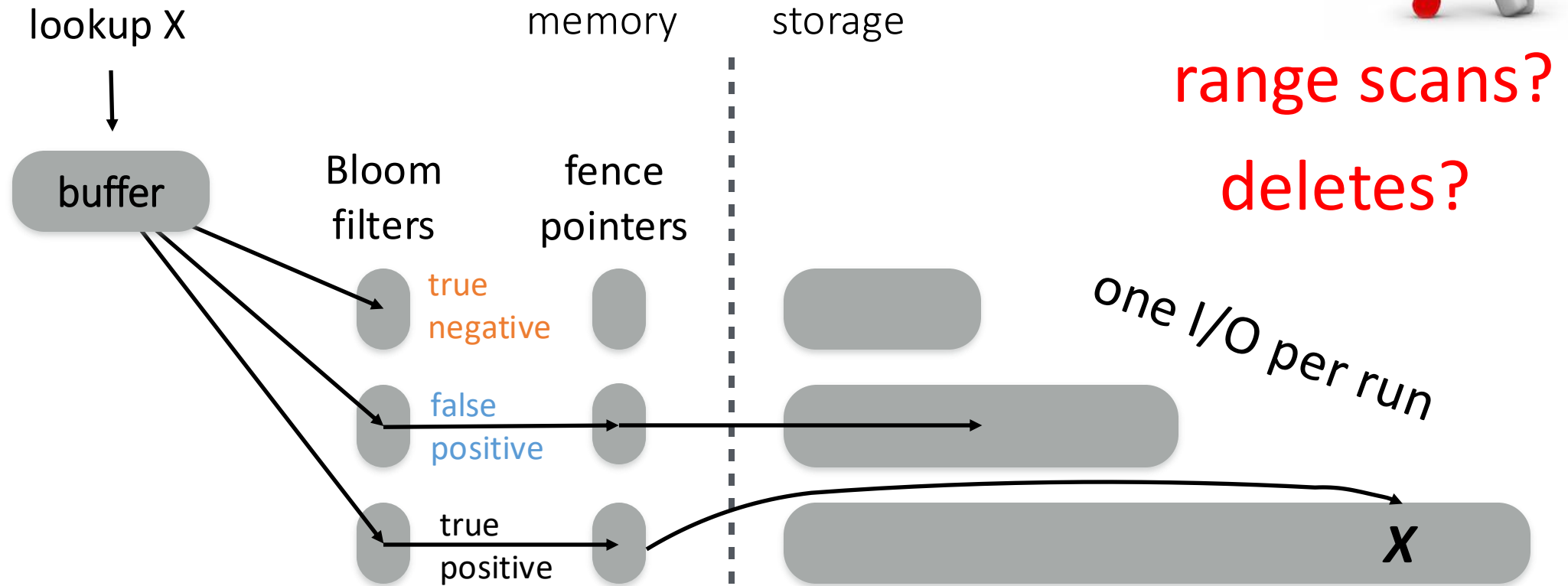




# performance & cost trade-offs



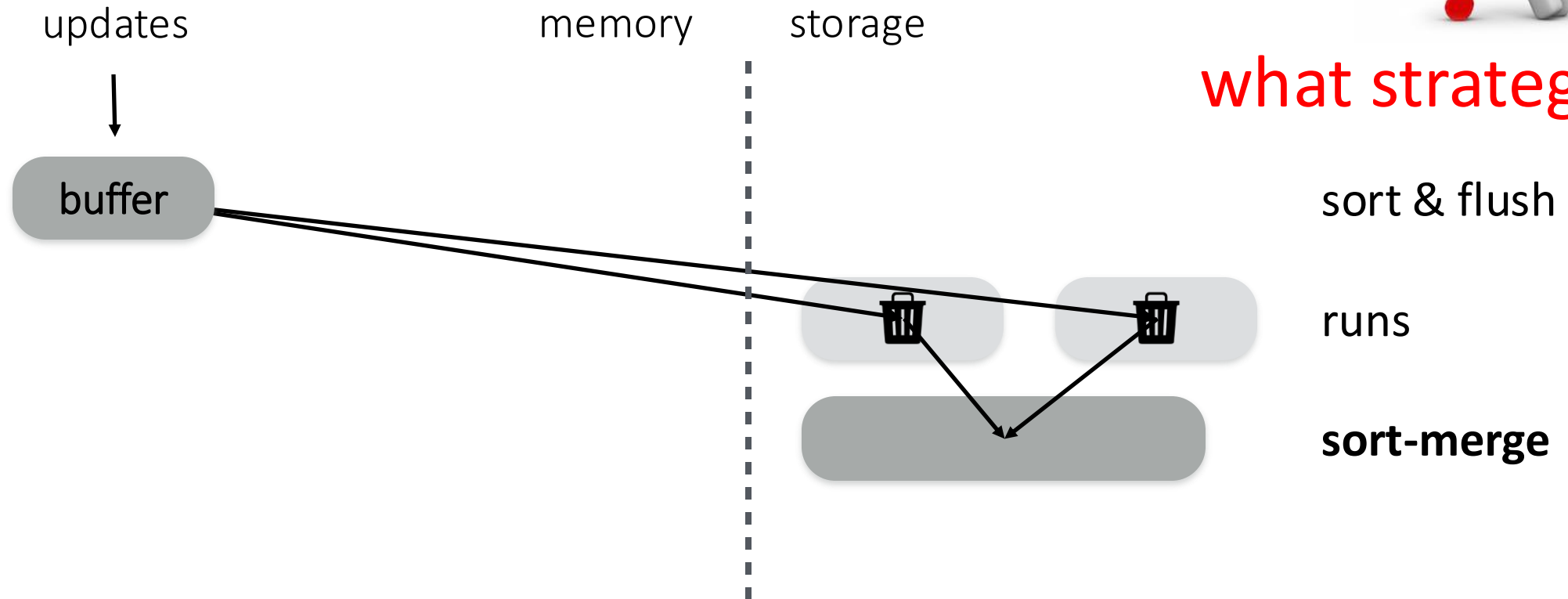
# other operations



# remember merging?



## what strategies?



# Merge Policies

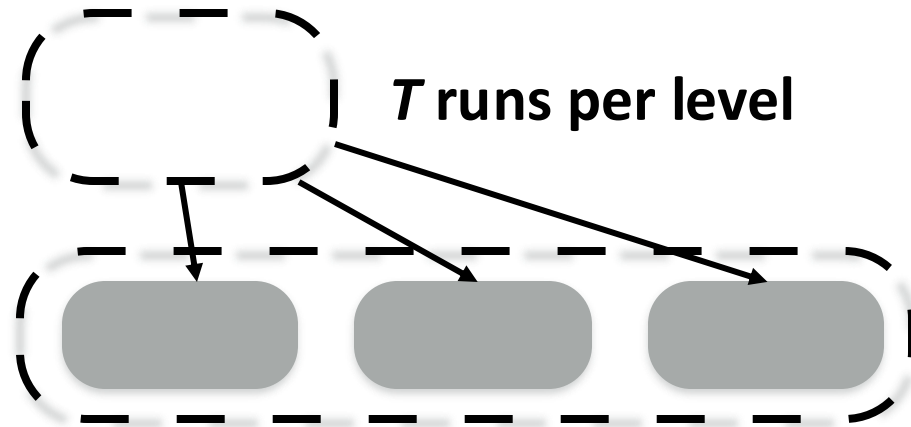
**Tiering**

write-optimized

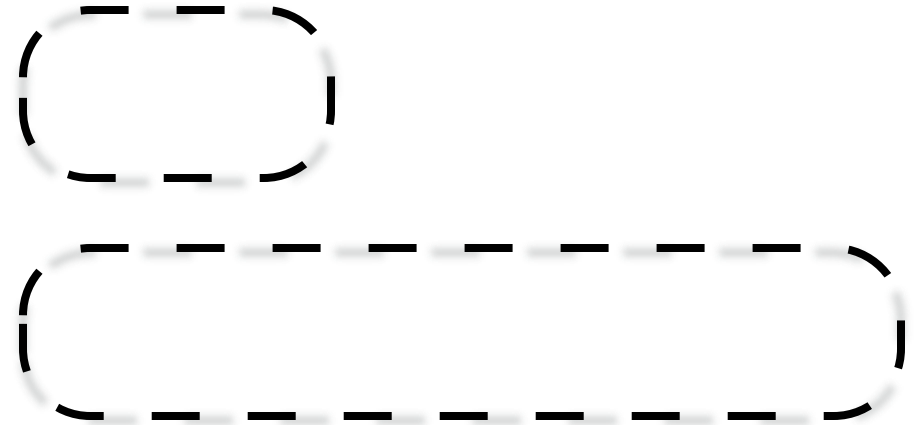
**Leveling**

read-optimized

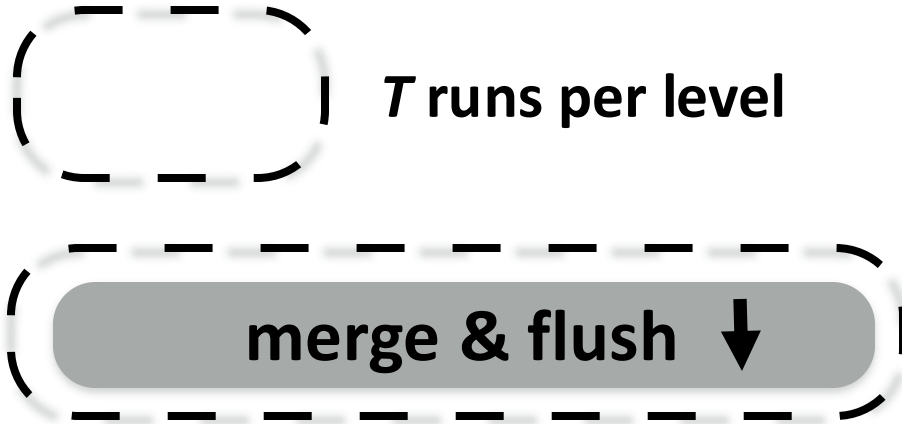
Tiering  
write-optimized



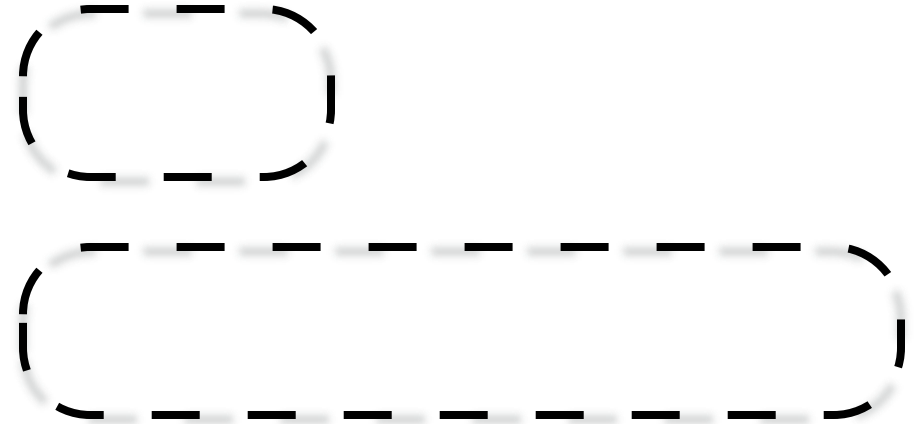
Leveling  
read-optimized



Tiering  
write-optimized

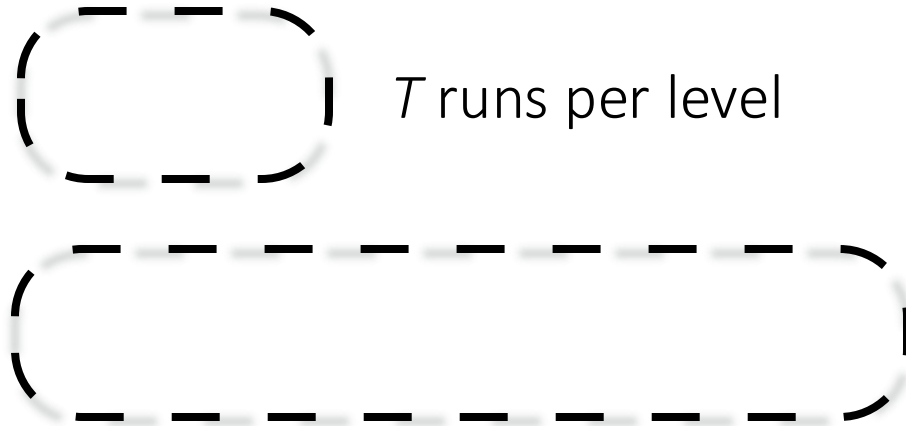


Leveling  
read-optimized

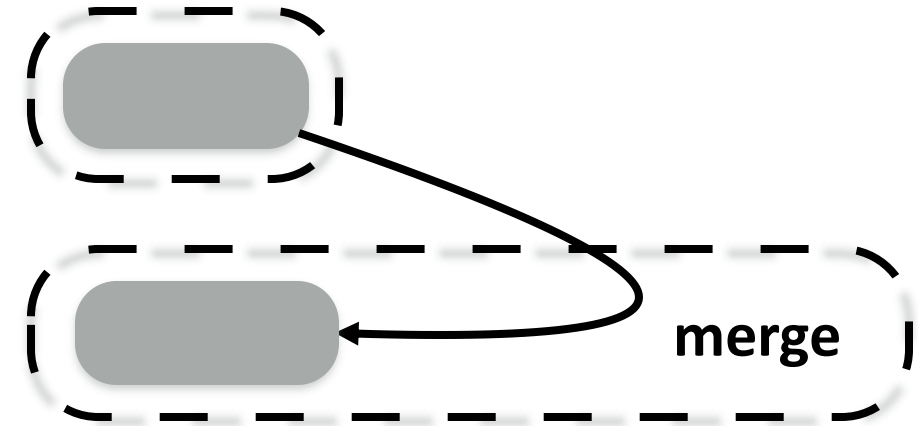




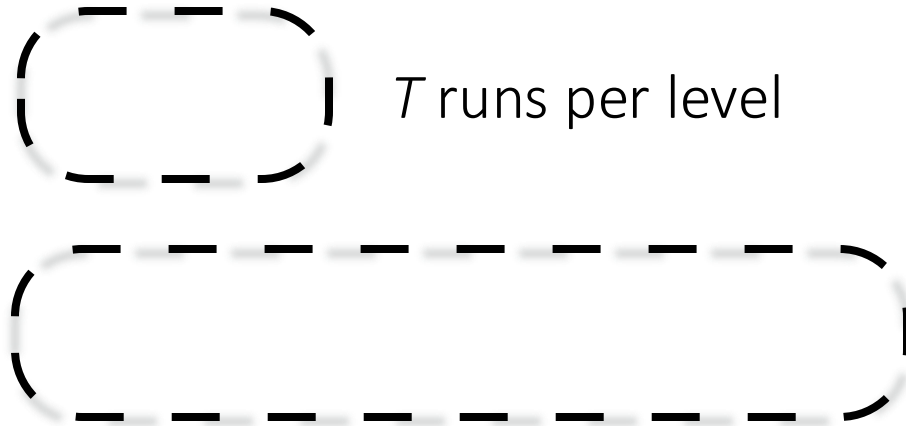
## Tiering write-optimized



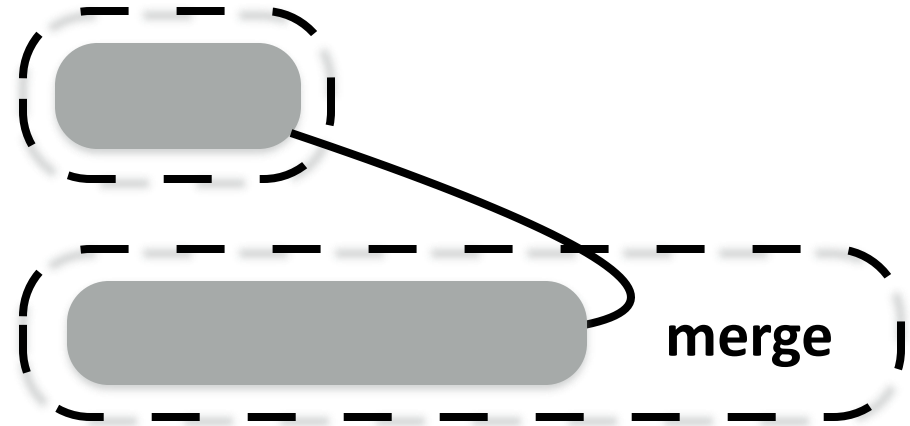
## Leveling read-optimized



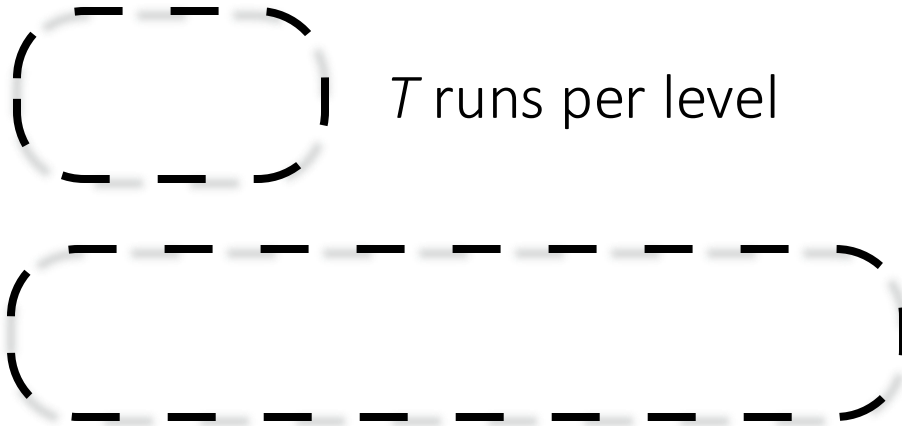
## Tiering write-optimized



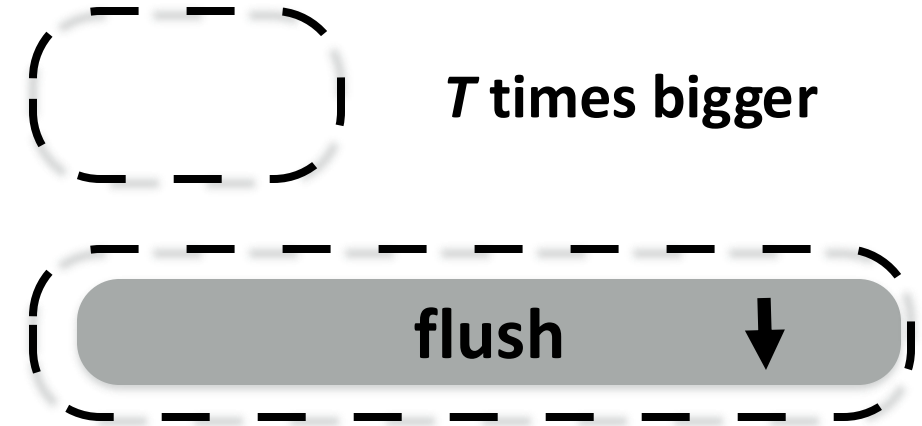
## Leveling read-optimized



## Tiering write-optimized



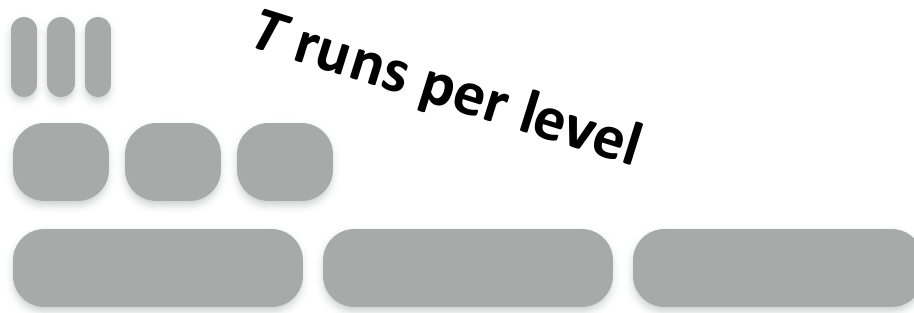
## Leveling read-optimized





more on LSM-Tree performance

## Tiering write-optimized



## Leveling read-optimized



**lookup cost:**

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

Diagram illustrating the lookup cost for Tiering (write-optimized):

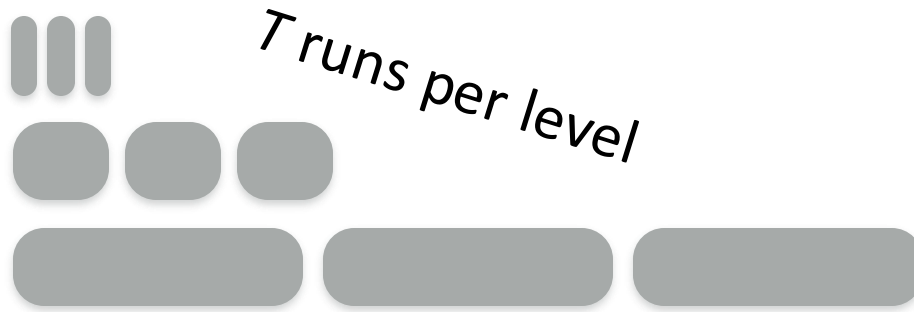
- $T$ : runs per level
- $\log_T(N)$ : levels
- $e^{-M/N}$ : false positive rate

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

Diagram illustrating the lookup cost for Leveling (read-optimized):

- $\log_T(N)$ : levels
- $e^{-M/N}$ : false positive rate

## Tiering write-optimized



## Leveling read-optimized



lookup cost:

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

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

update cost:

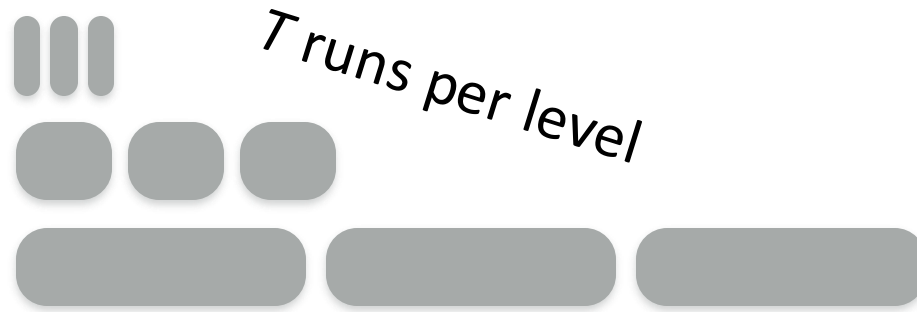
$$O(\log_T(N))$$

↑  
levels

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

↑      ↓  
merges per level      levels

## Tiering write-optimized



## Leveling read-optimized



lookup cost:

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

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

update cost:

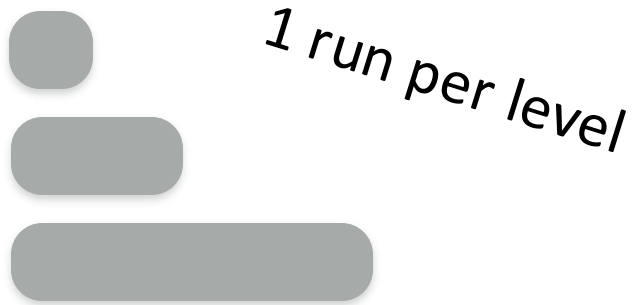
$$O(\log_T(N))$$

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

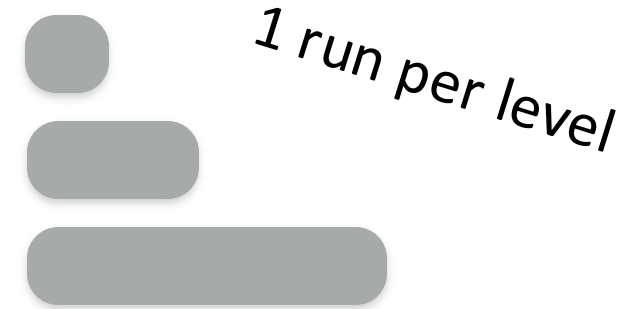
for size ratio  $T \gg$



Tiering  
write-optimized



Leveling  
read-optimized



lookup cost:

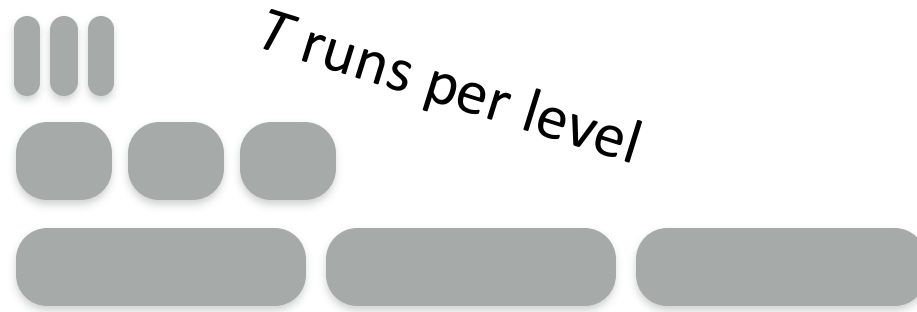
$$O(\log_T(N) \cdot e^{-M/N}) = O(\log_T(N) \cdot e^{-M/N})$$

update cost:

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

**for size ratio  $T \gg 1$**

## Tiering write-optimized



## Leveling read-optimized



lookup cost:

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

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

update cost:


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

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

**for size ratio T**  $\gg$

Tiering  
write-optimized

$O(N)$  runs per level



**log**

Leveling  
read-optimized

1 run per level



**sorted array**

lookup cost:

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

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

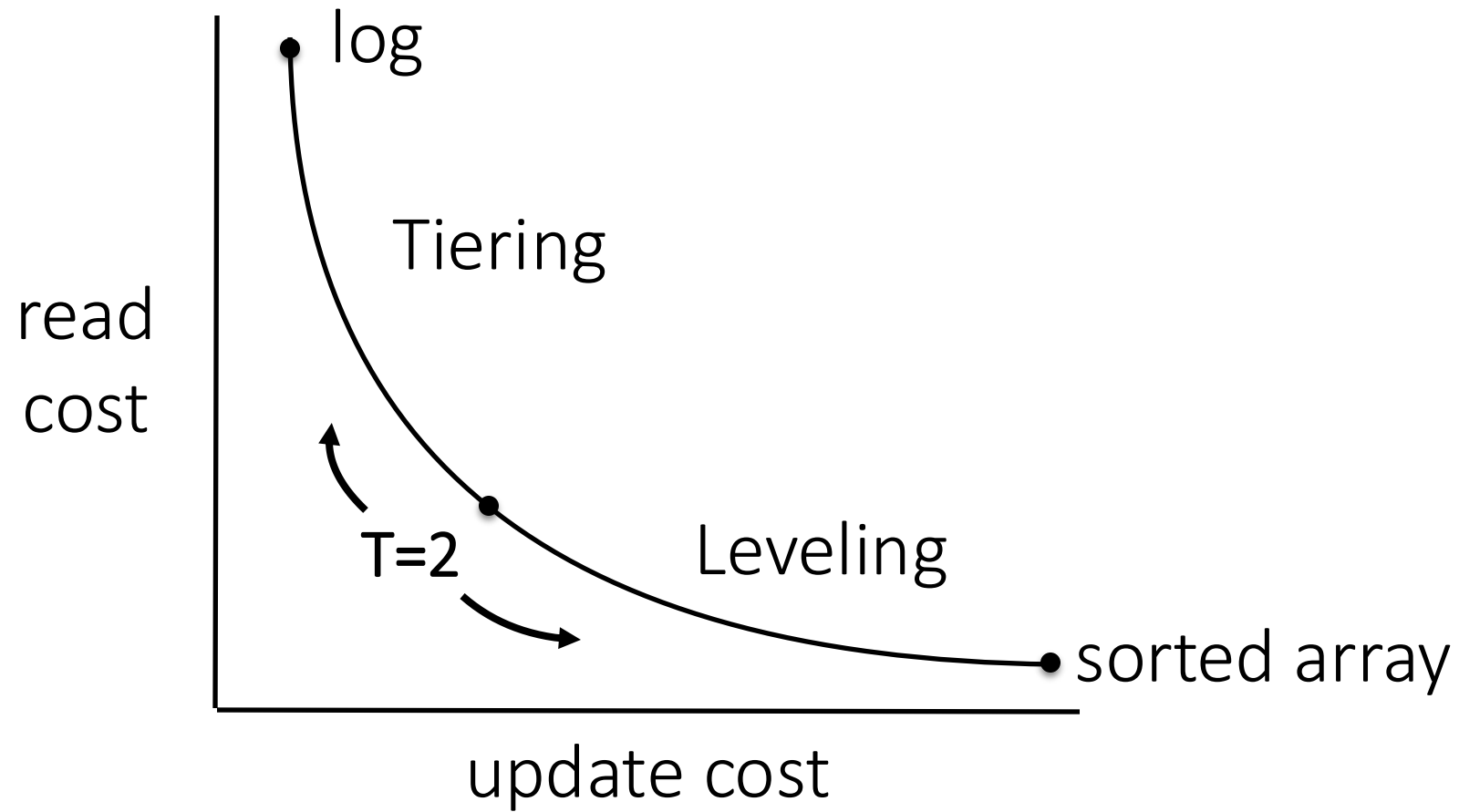
update cost:

$$O(\log_N(N)) = \mathbf{O(1)}$$

$$O(N \cdot \log_N(N)) = \mathbf{O(N)}$$

**for size ratio T**

**N**  
**⌋**



T : size ratio

# Research Question on LSM-Trees

how can we minimize the duplicate space during compaction?



how to ensure that we can tune without sharing workload details?

How much is the *real* write-amplification on SSDs?

buffer

Bloom  
filters

fence  
pointers



study these questions and navigate LSM  
design space using Facebook's RocksDB

# Research on PostgreSQL



***A state-of-the-art relational database***

How can we implement a skew-aware efficient join algorithm?

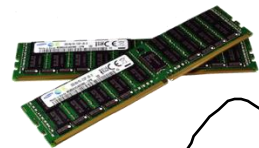
How much a noise in the existing cardinality estimation can impact the selected query plan, and the overall query performance



# Systems Project: Bufferpool

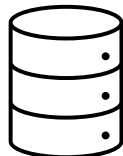
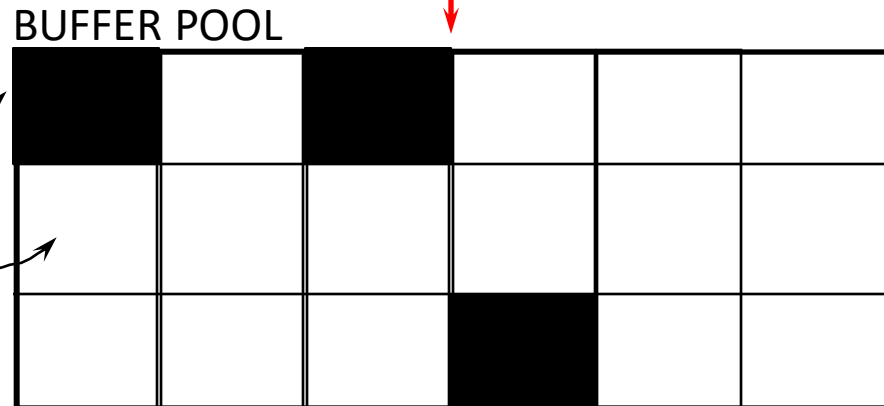
## Implementation of a bufferpool

### Page Requests from Higher Levels



disk page

free frame



choice of frame dictated by replacement policy

- Application requests a page
  - If in the **bufferpool** return it
  - If **not in the bufferpool** fetch it from the disk
    - If bufferpool is full select page to **evict**

### Core Idea: Eviction Policy

- *Least Recently Used*
- *First In First Out*
- *more ...*

# Research Topics on Buffer Management & SSDs

How to deploy ACE buffer management on emulated SSDs?

How can we achieve storage parallelism in ZNS SSDs?



# Other Research Topics on Indexing and beyond

How to apply sortedness-aware concepts on Adaptive Radix Tree?

How to build cache-friendly sortedness-aware indexing?

How to design a sortedness-aware join algorithm?

# what to do now?

## **systems project**

form groups of 3

(speak to me in OH if you want a different arrangement)

## **research project**

form groups of 3

pick one of the subjects & read background  
material

define the behavior you will study and address  
sketch approach and success metric  
(if LSM-related get familiar with RocksDB)

# what to do now?

**systems project**

form groups of 2

(speak to me in OH if you want to work on your own)

**research project**

6-8

*come to OH/Labs*

submit **project 0** this Friday on 1/31

start working on **project 1** (due on 2/14)

submit **semester project proposal** on 2/23

# CS 561: Data Systems Architectures

class 4

Systems & Research Project

Prof. Manos Athanassoulis

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