# 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"*)

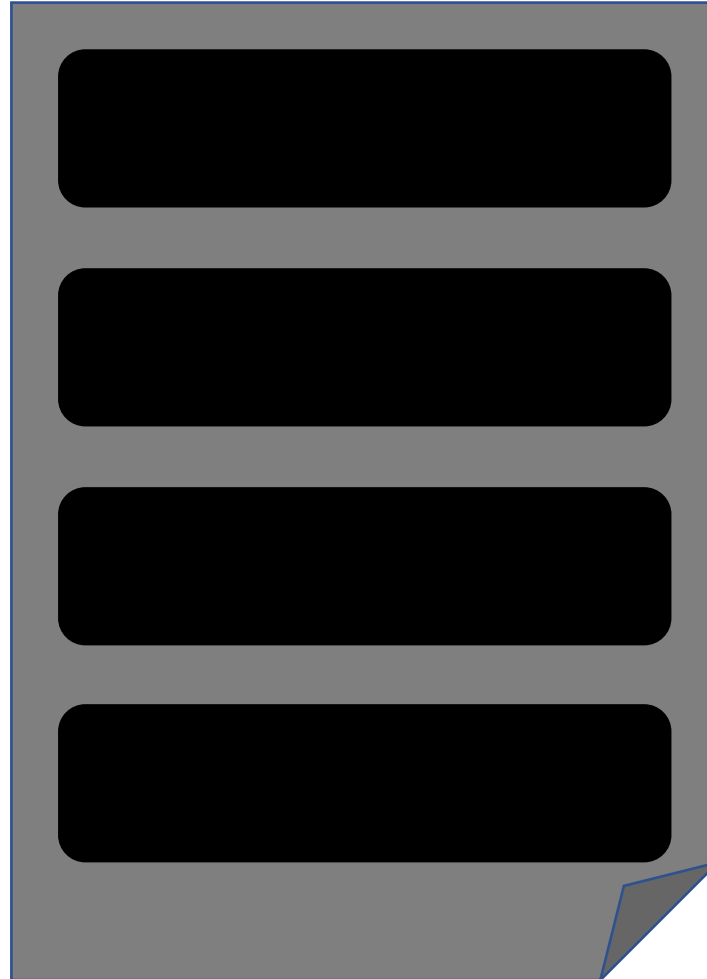# Let's revisit Zonemaps

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
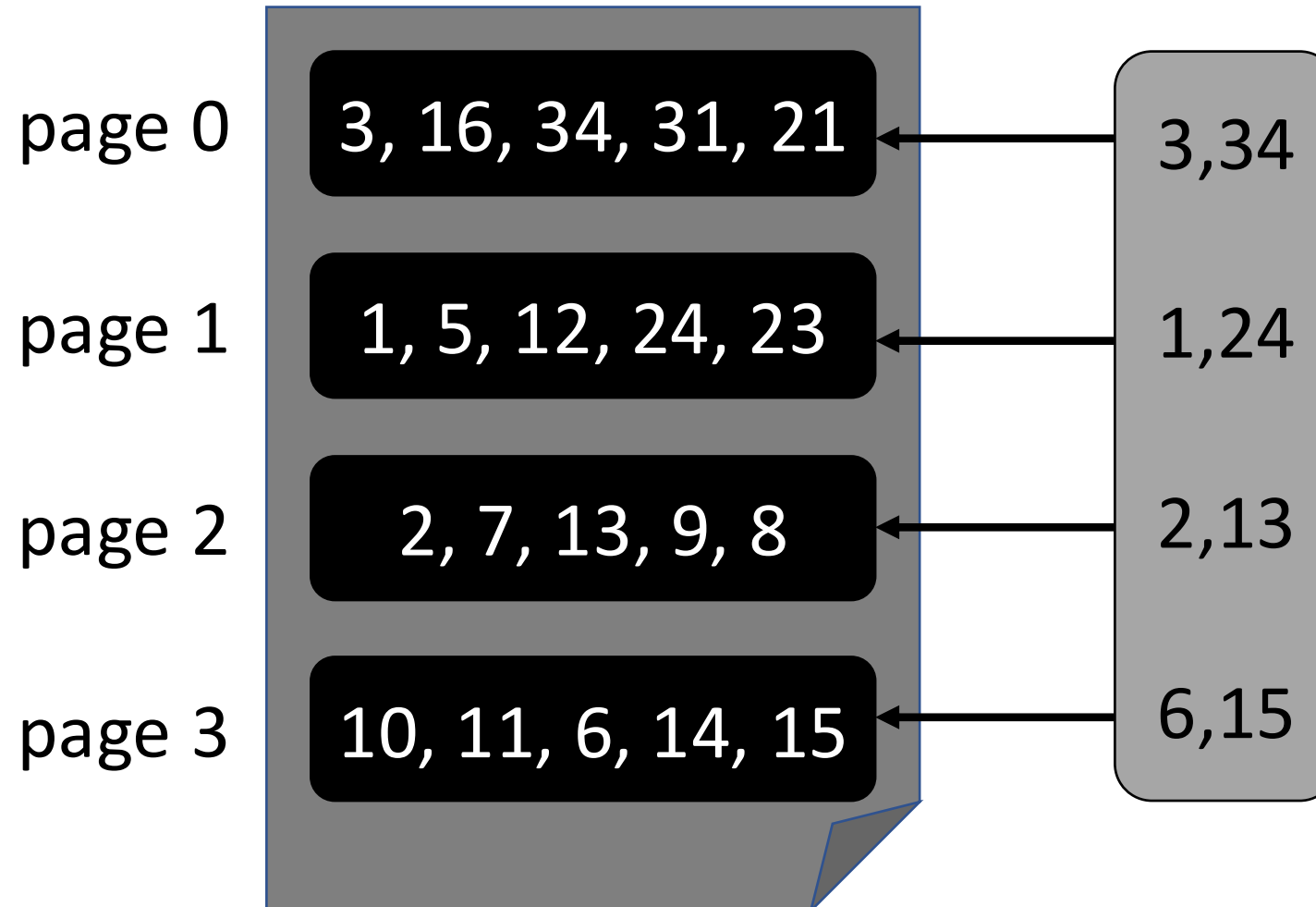
page 1 — 1, 5, 12, 24, 23

page 2 — 2, 7, 13, 9, 8

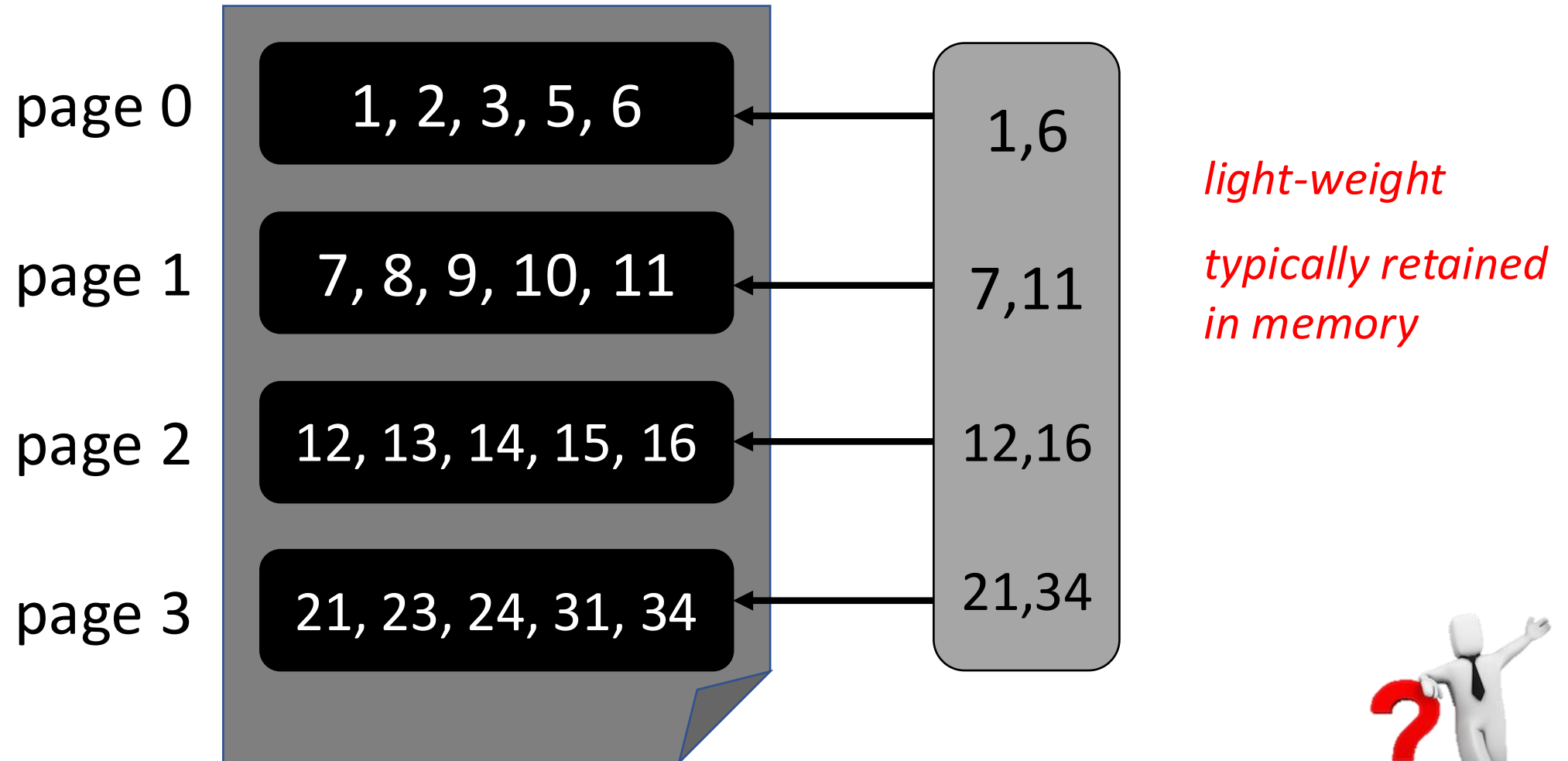page 3 — 10, 11, 6, 14, 15

# zonemaps

file = collection of pages

page 0 — **3, 16, 34, 31, 21** ← 3,34

*light-weight*

page 1 — **1, 5, 12, 24, 23** ← 1,24

*typically retained in memory*

page 2 — **2, 7, 13, 9, 8** ← 2,13

page 3 — **10, 11, 6, 14, 15** ← 6,15

But what if the data is sorted?

# zonemaps

file = collection of pages

page 0    **1, 2, 3, 5, 6**   ←   1,6

*light-weight*

page 1    **7, 8, 9, 10, 11**   ←   7,11

*typically retained in memory*

page 2    **12, 13, 14, 15, 16**   ←   12,16

page 3    **21, 23, 24, 31, 34**   ←   21,34

But what if the data is sorted?

# zonemaps

file

page 0    **3, 16, 34, 31, 21**   ←──── 3,34

page 1    **1, 5, 12, 24, 23**   ←──── 1,24

page 2    **2, 7, 13, 9, 8**   ←──── 2,13

page 3    **10, 11, 6, 14, 15**   ←──── 6,15

Range Queries

[25, 100]    1 page

[20, 25]    2 pages

[10, 13]    4 pages

# zonemaps

But what if the data is sorted?

file

| | | Range Queries |
|---|---|---|
| page 0 | 1, 2, 3, 5, 6 | 1,6 |
| | | [25, 100]  1 page |
| page 1 | 7, 8, 9, 10, 11 | 7,11 |
| | | [20, 25]  1 page |
| page 2 | 12, 13, 14, 15, 16 | 12,16 |
| | | [10, 13]  2 pages |
| page 3 | 21, 23, 24, 31, 34 | 21,34 |

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]

complex analytics

simple queries

access data

store, maintain, update

**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** } }

**how to index these attributes?**

index: { **name**, { **student_id** } }

index: { **yob**, { **student_id$_1$, student_id$_2$, ...** } }

what is the caveat?

other problems?

# 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$_1$, student_id$_2$, …** } }

# how to use a key-value store?

**basic interface**

put(k,v)

{v} = get(k)        {$v_1$, $v_2$, …} = get(k)

{$v_1$, $v_2$, …} = get_range($k_{min}$, $k_{max}$)        {$v_1$, $v_2$, …} = full_scan()

c = count($k_{min}$, $k_{max}$)

*deletes: delete(k)*

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

get set: {$v_1$, $v_2$, …} = get_set($k_1$, $k_2$, …)

more?

# how to build a key-value store?

**append**

if we have only *put* operations

if we mostly have *get* operations

*sort*

what about full scan?

and then

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*

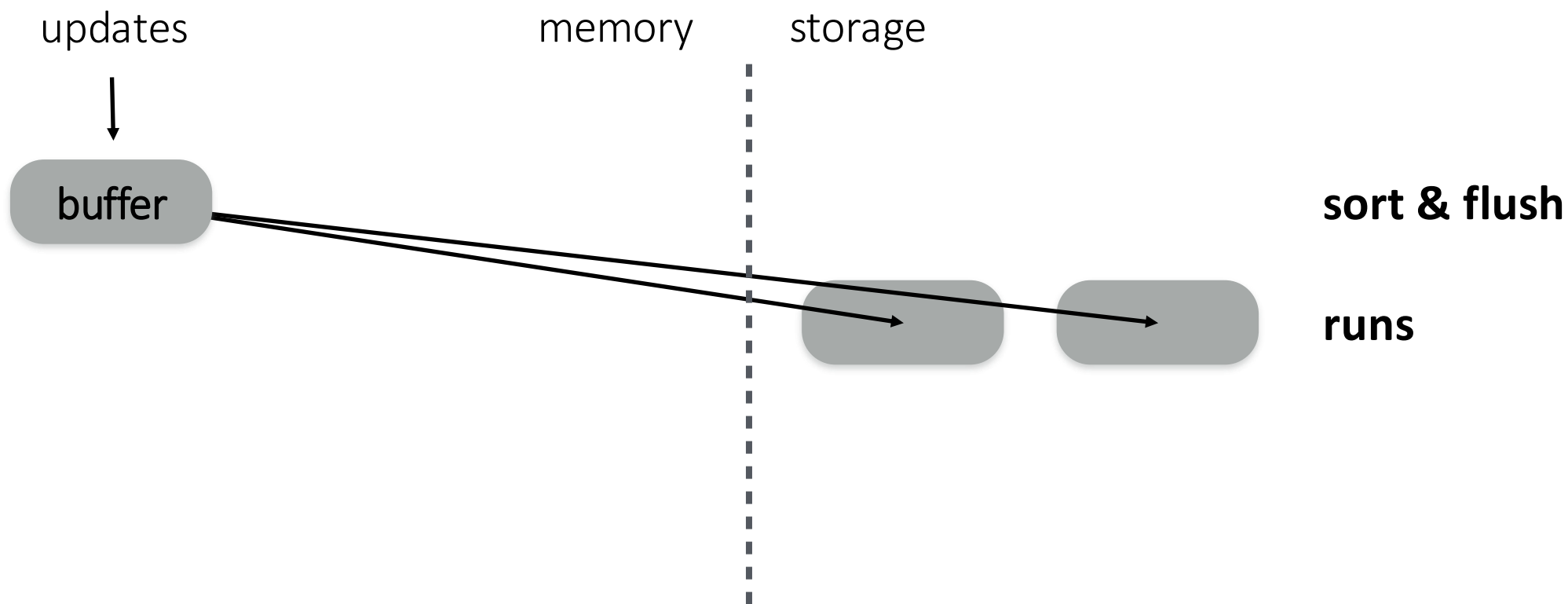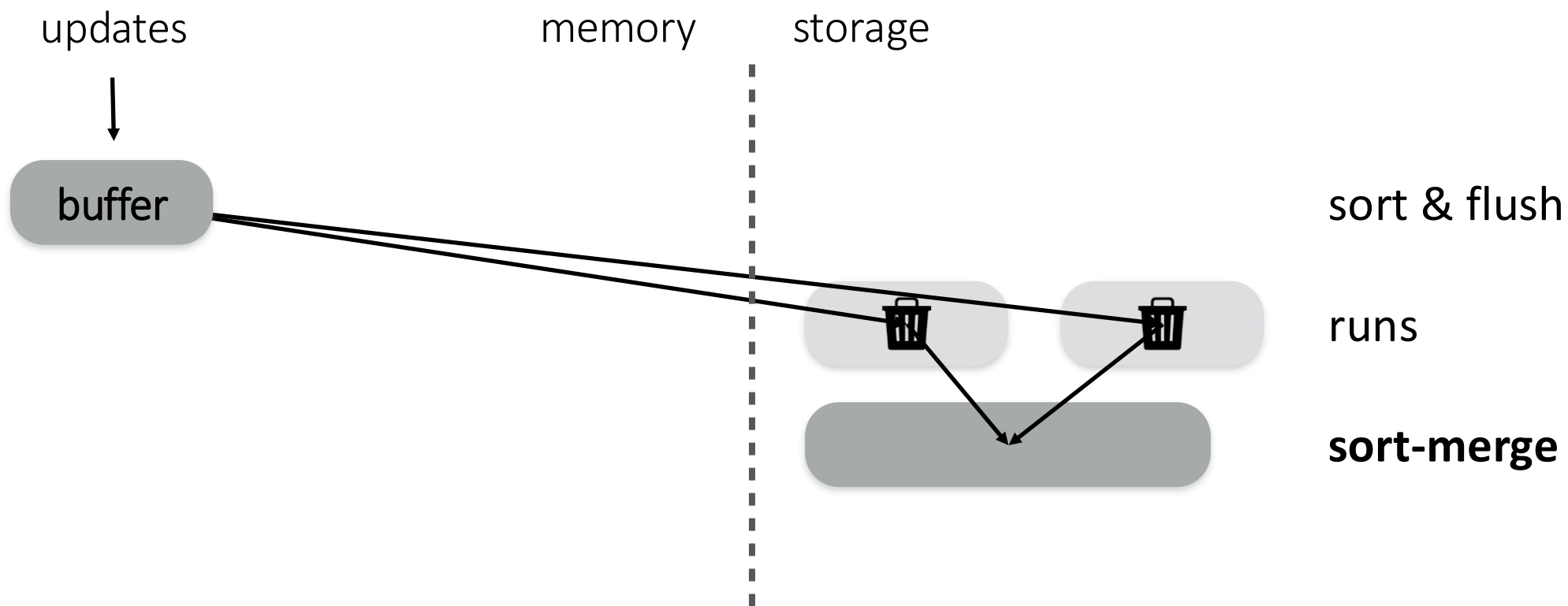how to bridge?
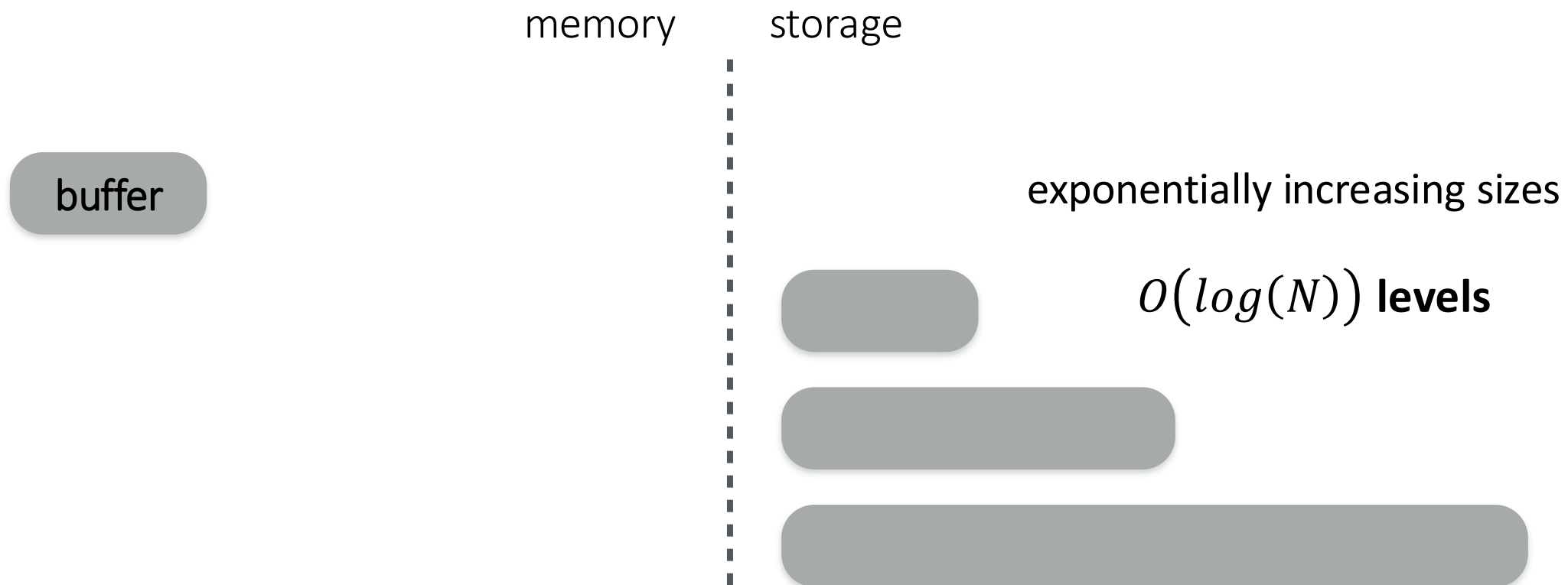
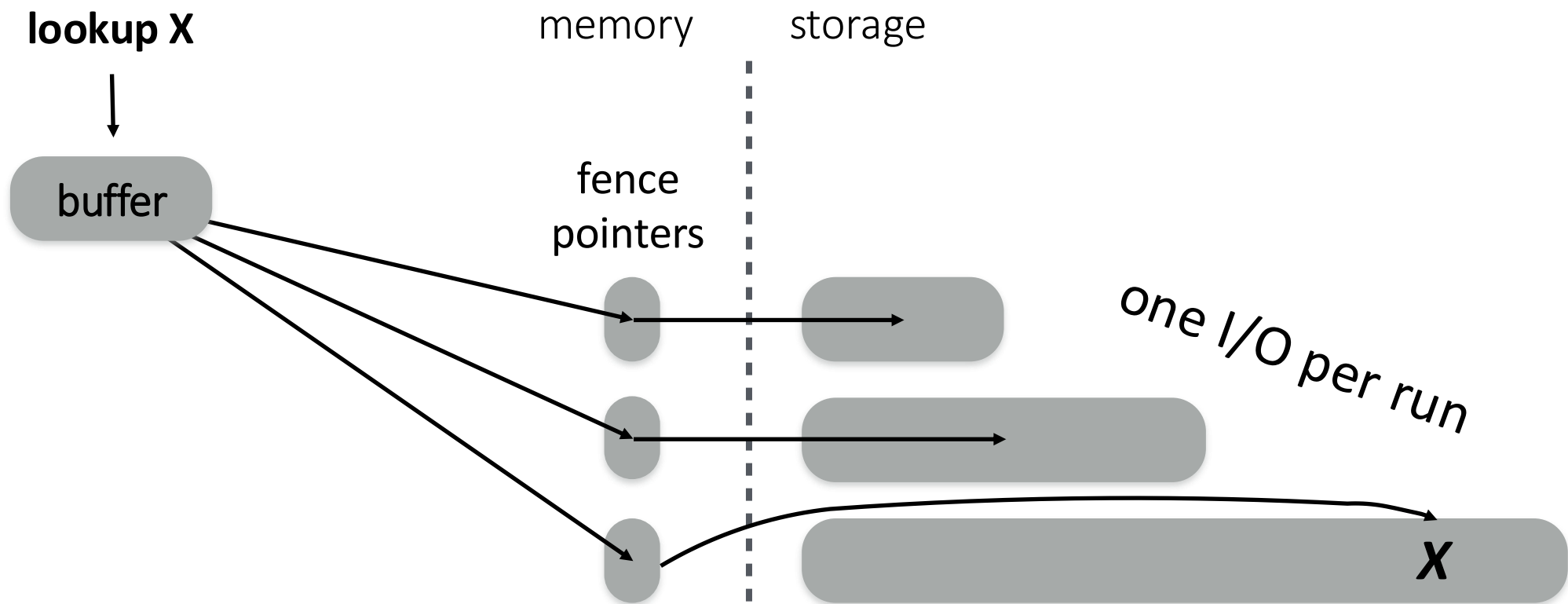# LSM-tree
# Key-Value Stores

What are they really?

memory    storage

buffer

exponentially increasing sizes

$O\big(log(N)\big)$ **levels**

**lookup X**

memory     storage

buffer

fence
pointers

one I/O per run
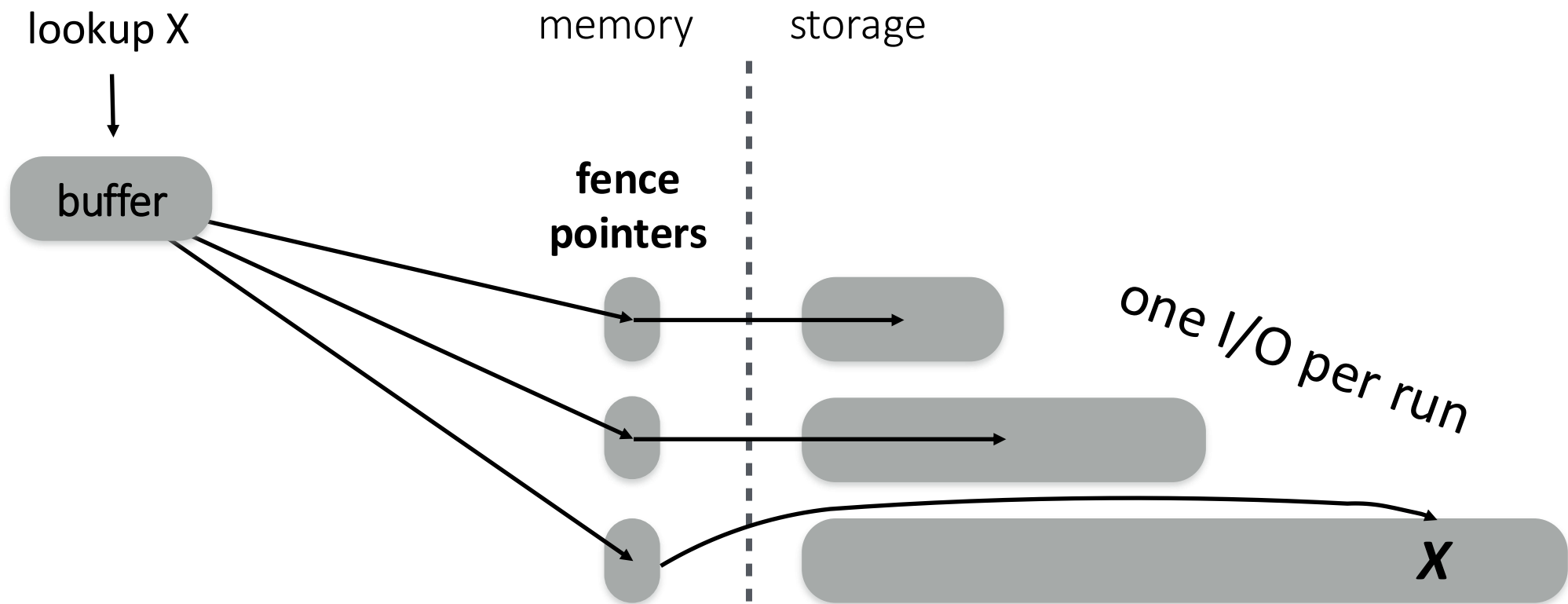
*X*

BOSTON
UNIVERSITY

lookup X

memory    storage

**fence pointers**

buffer

one I/O per run

X

# performance & cost trade-offs

# other operations

lookup X

memory     storage

buffer

Bloom filters

fence pointers

true negative

false positive

true positive

range scans?

deletes?

one I/O per run

X

# remember merging?

updates    memory    storage

## what strategies?

buffer

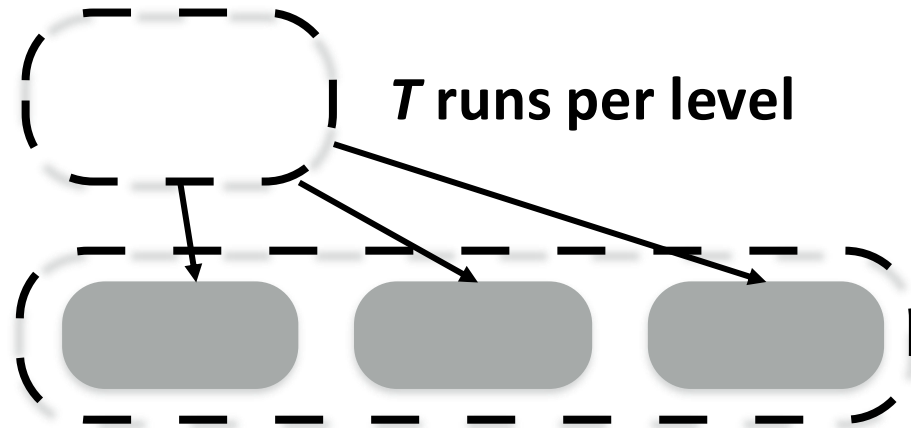sort & flush
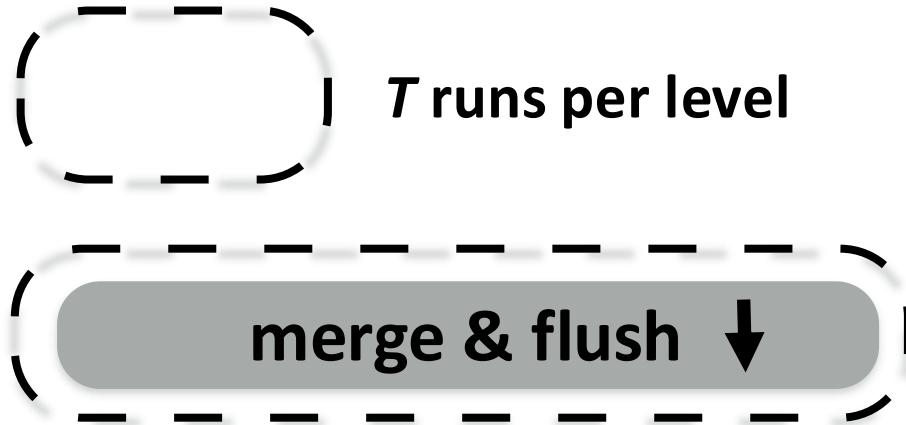
runs

**sort-merge**

# Merge Policies

**Tiering**
write-optimized

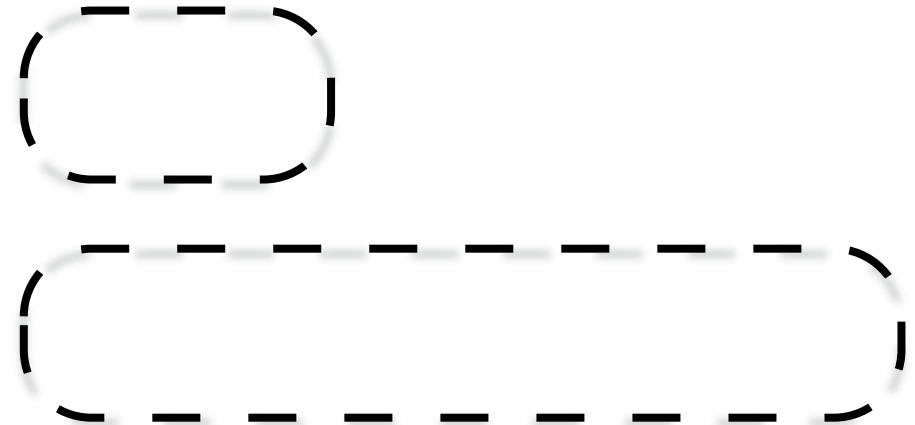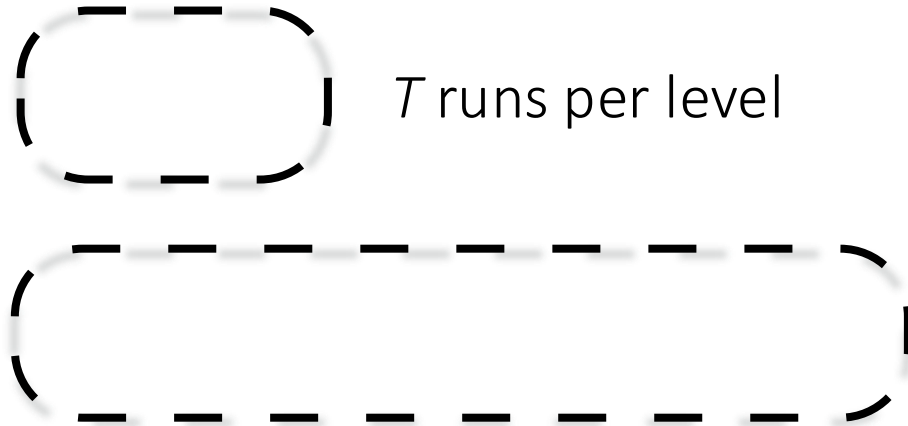**Leveling**
read-optimized

# Tiering
## write-optimized

# Leveling
## read-optimized

**T runs per level**

# Tiering
## write-optimized

# Leveling
## read-optimized

*T* **runs per level**

**merge & flush** ⬇

# Tiering
## write-optimized

# Leveling
## read-optimized



*T* runs per level

**merge**

# Tiering
## write-optimized

# Leveling
## read-optimized

*T* runs per level

**merge**

# Tiering
## write-optimized

# Leveling
## read-optimized

*T* runs per level

***T* times bigger**

**flush**

# Systems Project: LSM-Trees

*tuning knobs*

**merge policy**

**size ratio**

lookup X

buffer

**Bloom filters**

fence pointers

memory

storage

true negative

false positive

true positive

one I/O per run

X

BOSTON UNIVERSITY

more on LSM-Tree performance

# Tiering
## write-optimized

*T runs per level*

# Leveling
## read-optimized

*1 run per level*

**lookup cost:**

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

runs per level → levels → false positive rate

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

levels → false positive rate

# Tiering
## write-optimized



*T runs per level*

# Leveling
## read-optimized

*1 run per level*

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

**update cost:** $\qquad$ $\boldsymbol{O(log_T(N))}$ $\qquad$ $\boldsymbol{O(T \cdot log_T(N))}$

levels

merges per level $\qquad$ levels

# Tiering
## write-optimized



T runs per level

# Leveling
## read-optimized



1 run per level

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

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

**for size ratio T**  ≫

BOSTON
UNIVERSITY

42

Tiering
write-optimized

Leveling
read-optimized

1 run per level

1 run per level

lookup cost: $O\left(log_T(N) \cdot e^{-M/N}\right) = O\left(log_T(N) \cdot e^{-M/N}\right)$

update cost: $O\left(log_T(N)\right) = O\left(log_T(N)\right)$

**for size ratio T** ⩒

# Tiering
## write-optimized

# Leveling
## read-optimized

*T runs per level*

*1 run per level*

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

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

**for size ratio T** ⟰

# Tiering
## write-optimized

# Leveling
## read-optimized

$O(N)$ runs per level

1 run per level

**log**

**sorted array**

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

update cost: $O\left(log_N(N)\right) = \boldsymbol{O(1)}$          $O\left(N \cdot log_N(N)\right) = \boldsymbol{O(N)}$

**N**

**for size ratio T** ⟪

log

Tiering

read cost

T=2

Leveling

sorted array

update cost

T : size ratio

46

# Research Question on LSM-Trees

how to reduce temporary space amplification at compaction time?

how to support variable deletion persistent thresholds?

how to address high number of contiguous tombstones?

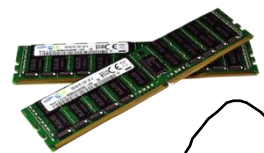buffer    **Bloom filters**    **fence pointers**    how to efficiently support range deletes?

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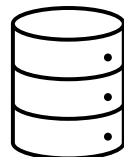design space using Facebook's RocksDB

# Systems Project: Bufferpool

*Implementation of a bufferpool*

**Page Requests from Higher Levels**

BUFFER POOL



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 on PostgreSQL



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

How can we implement a sortendess-aware join algorithm?

Integrate light-weight eviction policies in PostgreSQL bufferpool

# Research Topics on SSDs

how to design a flexible ZNS SSD interface?

understanding Garbage Collection in applications on ZNS SSDs

understanding the impact of Zone Sizes on ZNS SSDs

# Research Topics on Bitmap Indexing

how to do access path selection in the presence of bitmap indexes?

how to execute index-based joins with bitmap indexes?

# 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**
~~form groups of 2~~

*come to OH/Labs*
submit project 0 <u>this Sunday</u> on 2/1
start working on project 1 (due on 2/15)
submit semester project proposal on 2/22

# CS 561: **Data Systems Architectures**

# class 4

# Systems & Research Project

Prof. Manos Athanassoulis

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