# CS 561: Data Systems Architectures

# Systems & Research Project
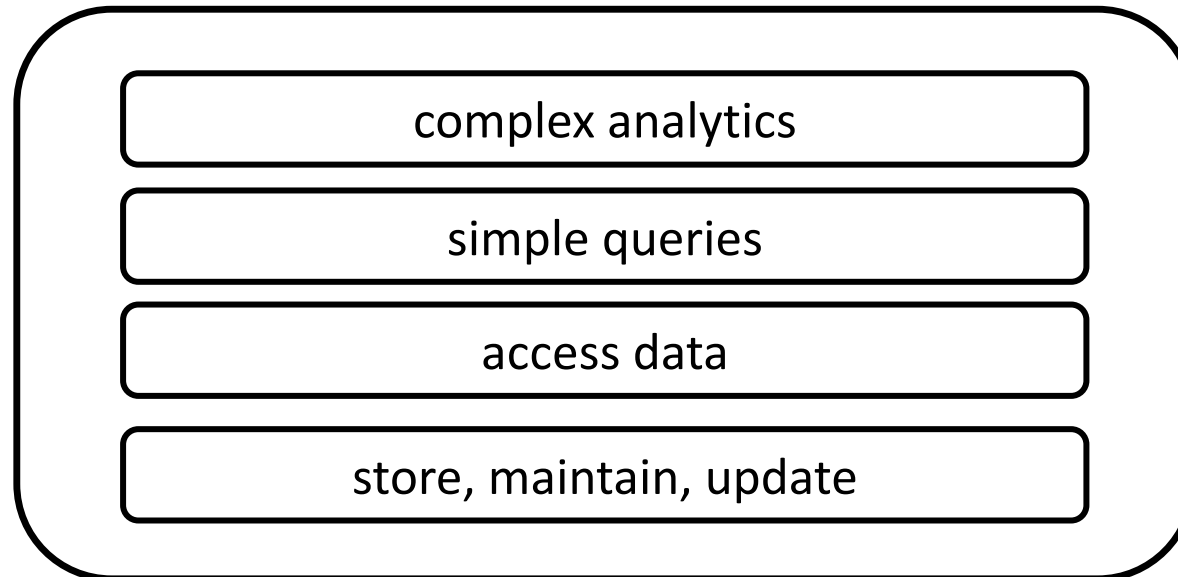
Prof. Manos Athanassoulis

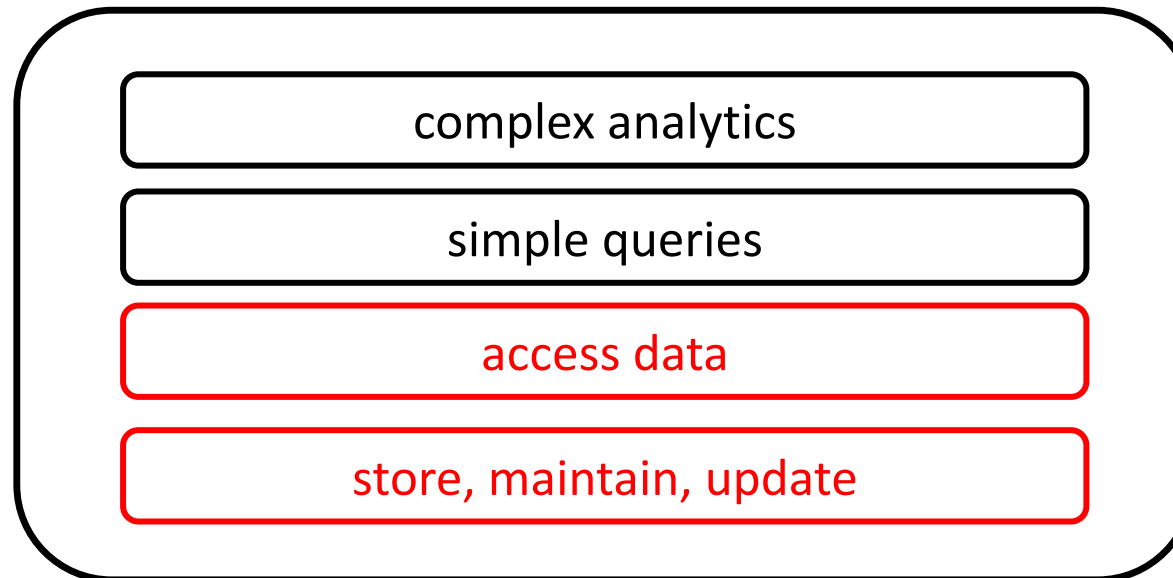https://midas.bu.edu/classes/CS591A1

# data systems

>$200B by 2020, growing at 11.7% every year

[The Forbes, 2016]

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*



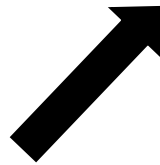*state-of-the-art* design

# 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** } }

**what is the caveat?**

**how to index these attributes?**

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

**other problems?**

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

# 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
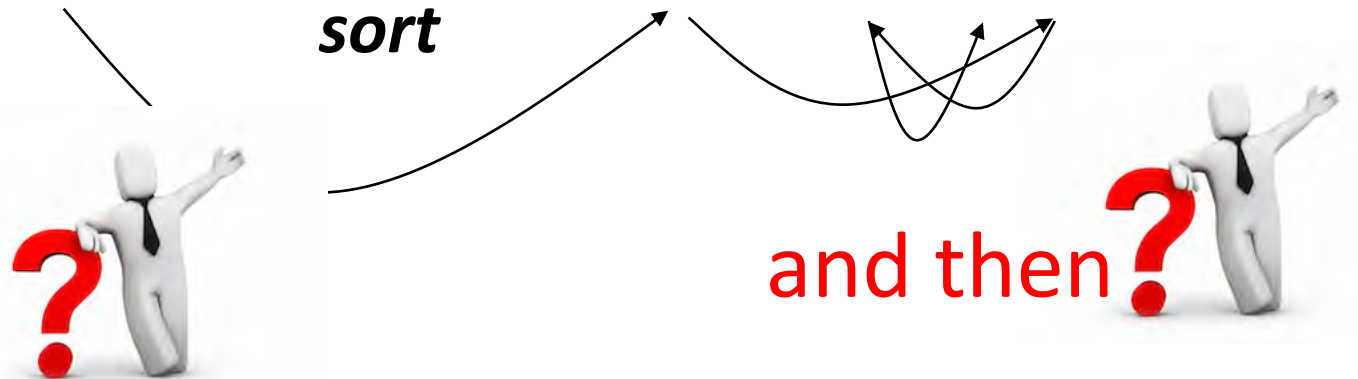
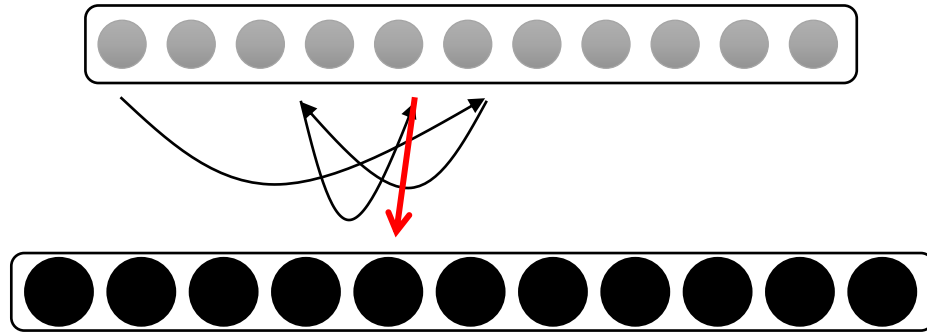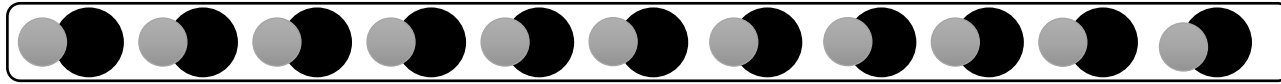if we mostly have *get* operations

**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
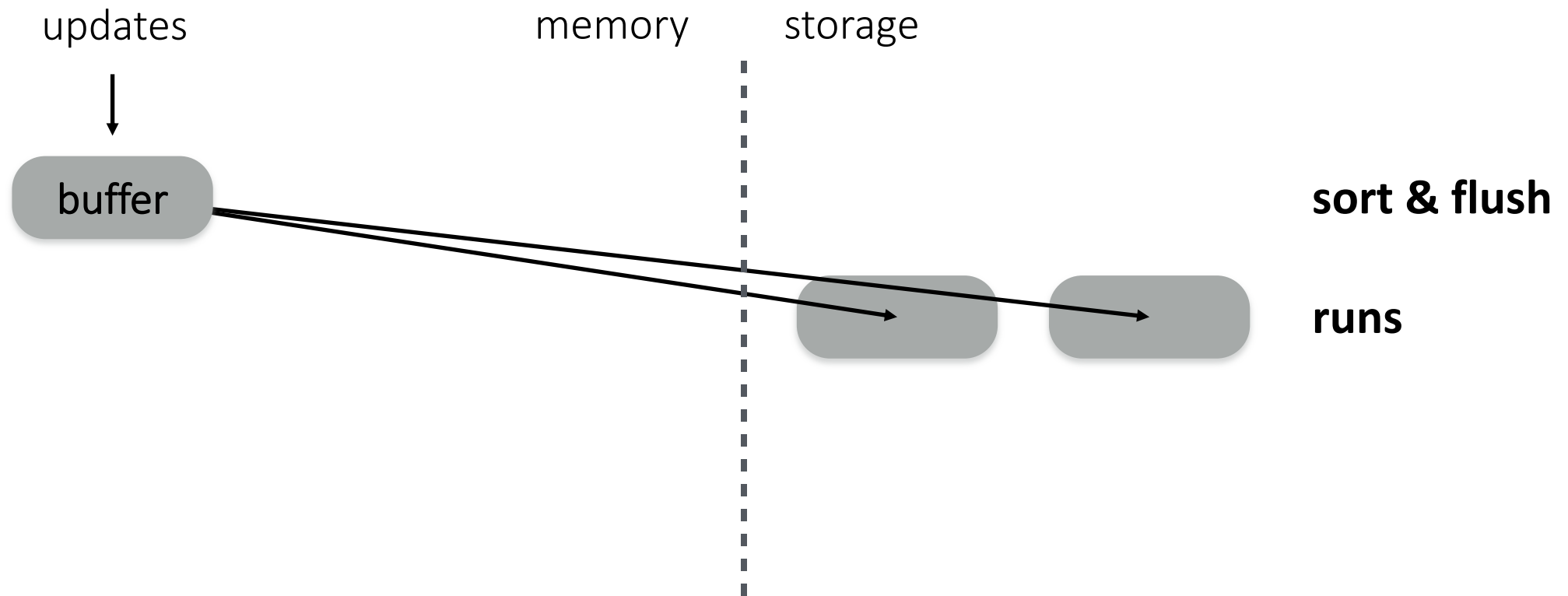# Key-Value Stores
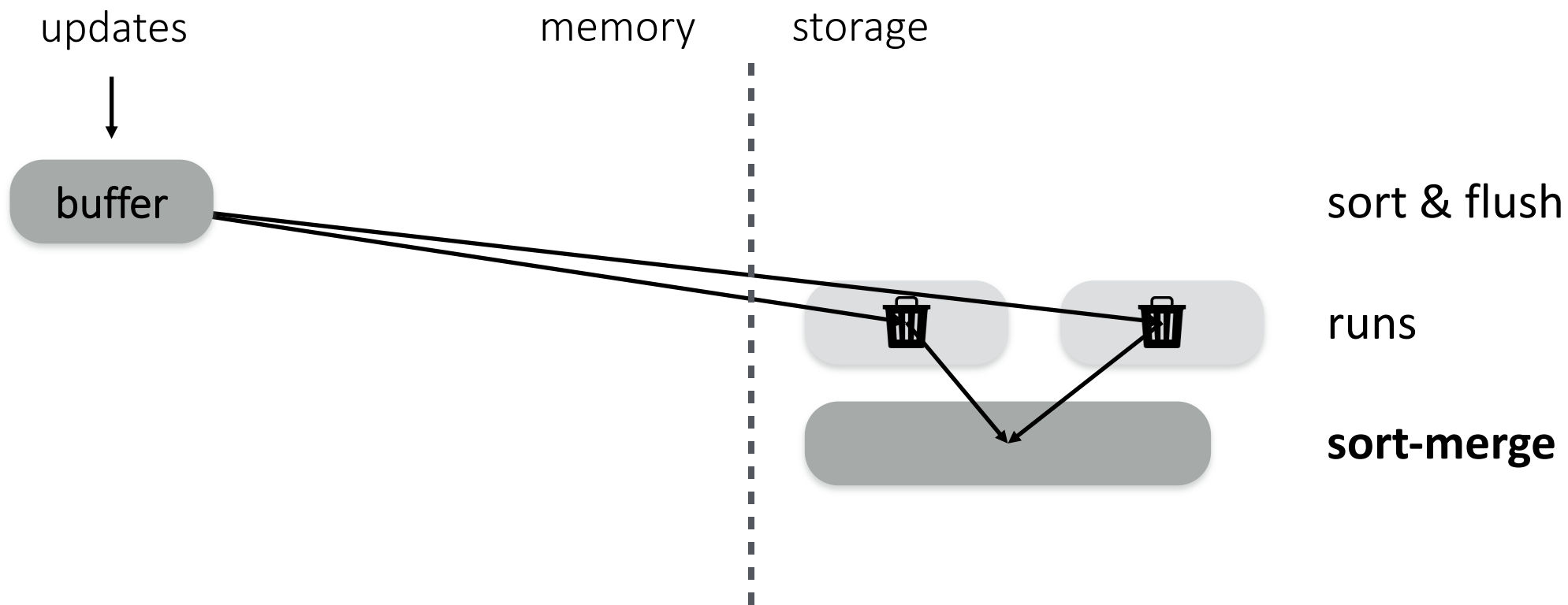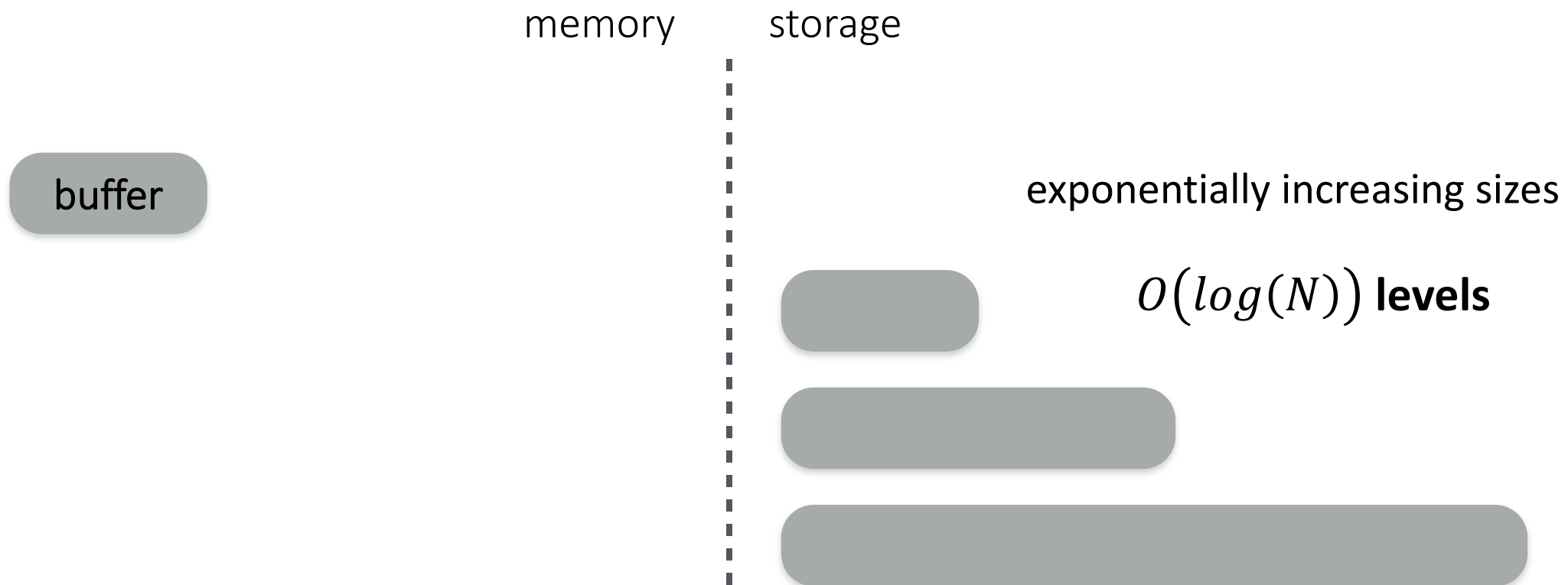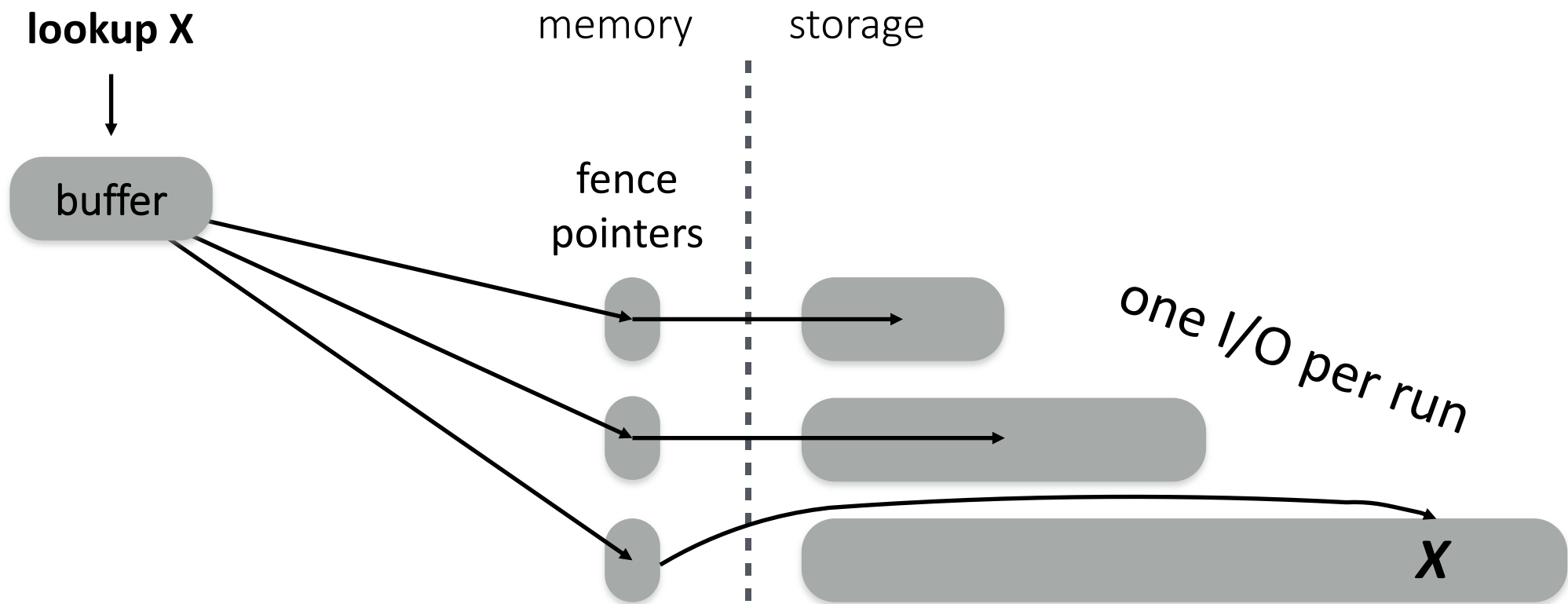
What are they really?

updates

memory    storage



buffer

BOSTON UNIVERSITY

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

buffer

**fence pointers**

one I/O per run

X

BOSTON UNIVERSITY

lookup X

memory     storage

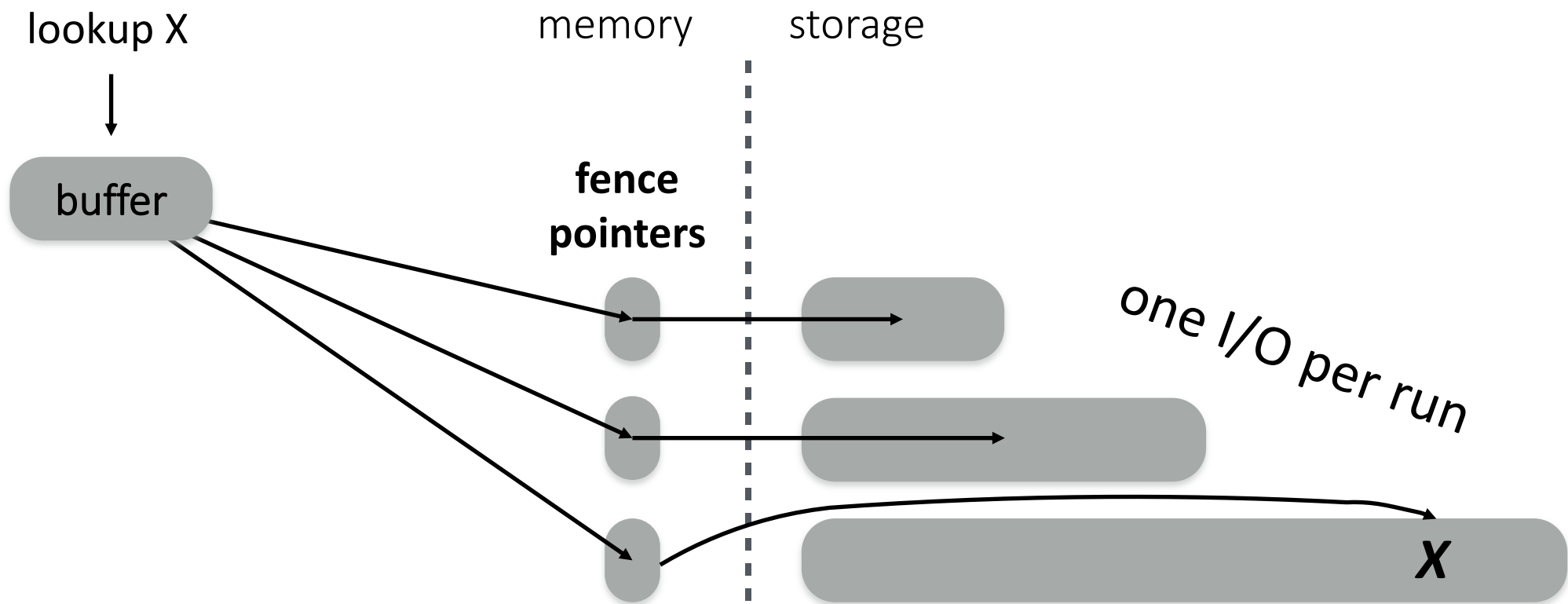**Bloom filters**     fence pointers

true negative

false positive

true positive
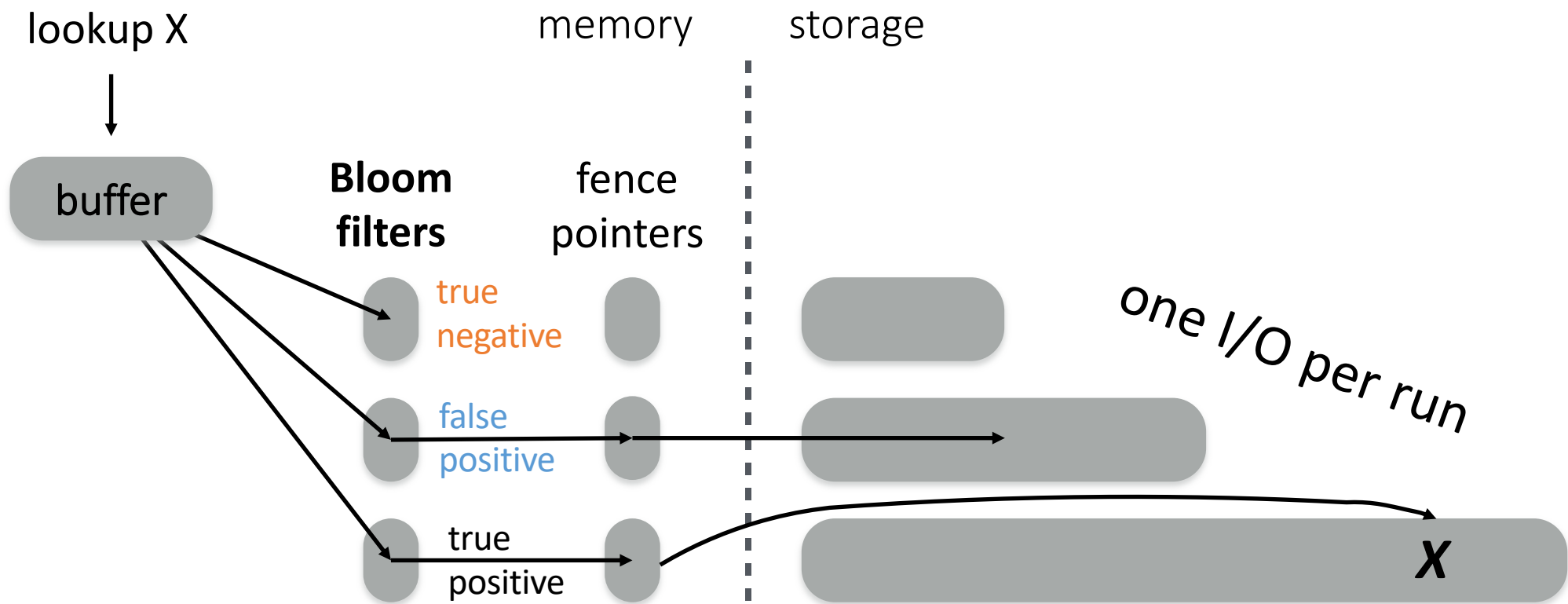
one I/O per run

X

BOSTON UNIVERSITY

# performance & cost trade-offs

# other operations

# remember merging?

updates

memory    storage



what strategies?

sort & flush
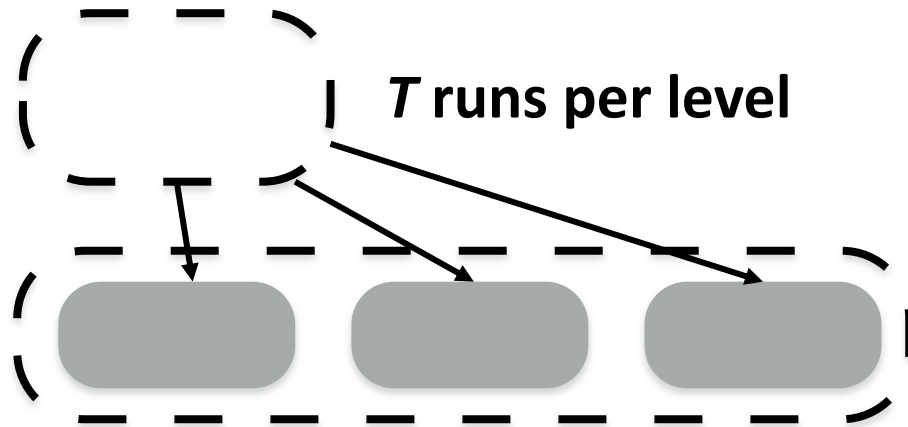
runs

**sort-merge**

# Merge Policies

**Tiering**
write-optimized

**Leveling**
read-optimized

# Tiering
## write-optimized

# Leveling
## read-optimized

*T* **runs per level**

Tiering
write-optimized
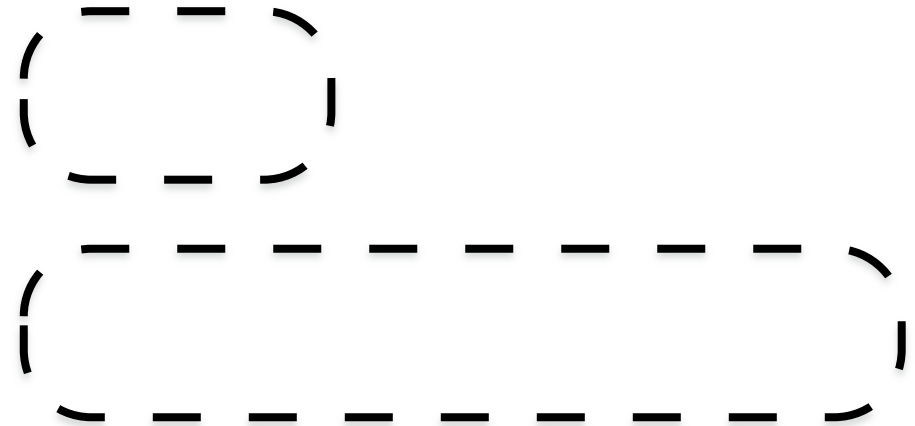
*T* runs per level

**merge & flush** ↓

Leveling
read-optimized

# 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

memory          storage

buffer

**Bloom filters**    fence pointers

true negative

false positive

one I/O per run

true positive
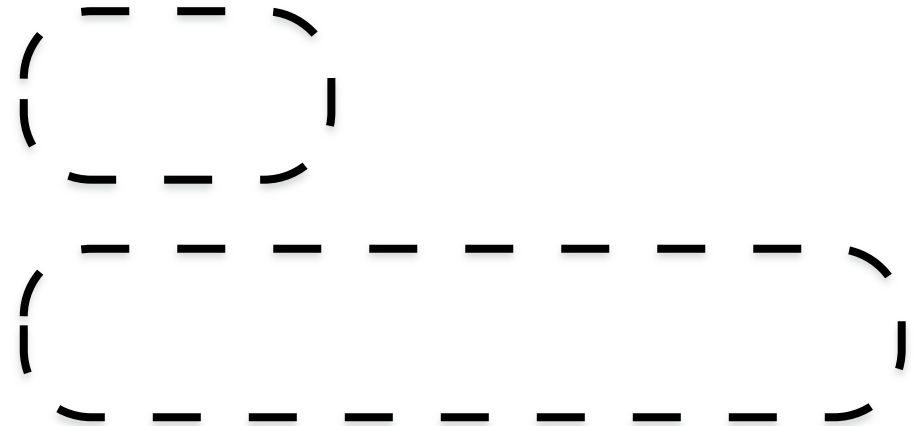
*X*

BOSTON UNIVERSITY

more on LSM-Tree performance

# Tiering
## write-optimized

# Leveling
## read-optimized

**T runs per level**

**1 run per level**

**lookup cost:**

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

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

runs
per level

levels

false
positive rate

levels

false
positive rate

# Tiering
## write-optimized



*T runs per level*

# Leveling
## read-optimized



*1 run per level*

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:** $\boldsymbol{O(log_T(N))}$

$\boldsymbol{O(T \cdot log_T(N))}$

levels

merges per level

levels

# 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

*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 | Leveling |
| :---: | :---: |
| write-optimized | read-optimized |

$O(N)$ runs per level

1 run per level

**log**

**sorted array**

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

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

**N**

**for size ratio T** ⌃⌃

read
cost

log

Tiering

T=2

Leveling

sorted array

update cost

T : size ratio

# Research Question on LSM-Trees

how to do range scans?

how to delete?        how to delete *quickly*?

how to allocate memory between buffer/Bloom filters/fence pointers?

what is the CPU overhead of Bloom filters?

buffer        **Bloom filters**    **fence pointers**

what if data items come ordered?

what if data items come *almost ordered*?

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

# What "almost ordered" even mean?

Research question on **sortedness**

How to quantify it?          Need a metric!

How does the sortedness of the data affect
the behavior of LSM-Trees, B-Trees, Zonemaps?

similar question to:

how does the order of the values in an array affect a sorting algorithm

# How to tune our system?

if we know the workload …

LSM-Trees: memory (Buffer/BF/FP) – what about caching?

**Back to column-stores**: do we need to sort?

*partition* the data?

add *empty slots* in the column for future inserts?

# Workload-based tuning

**find** *Tuning*, s.t.
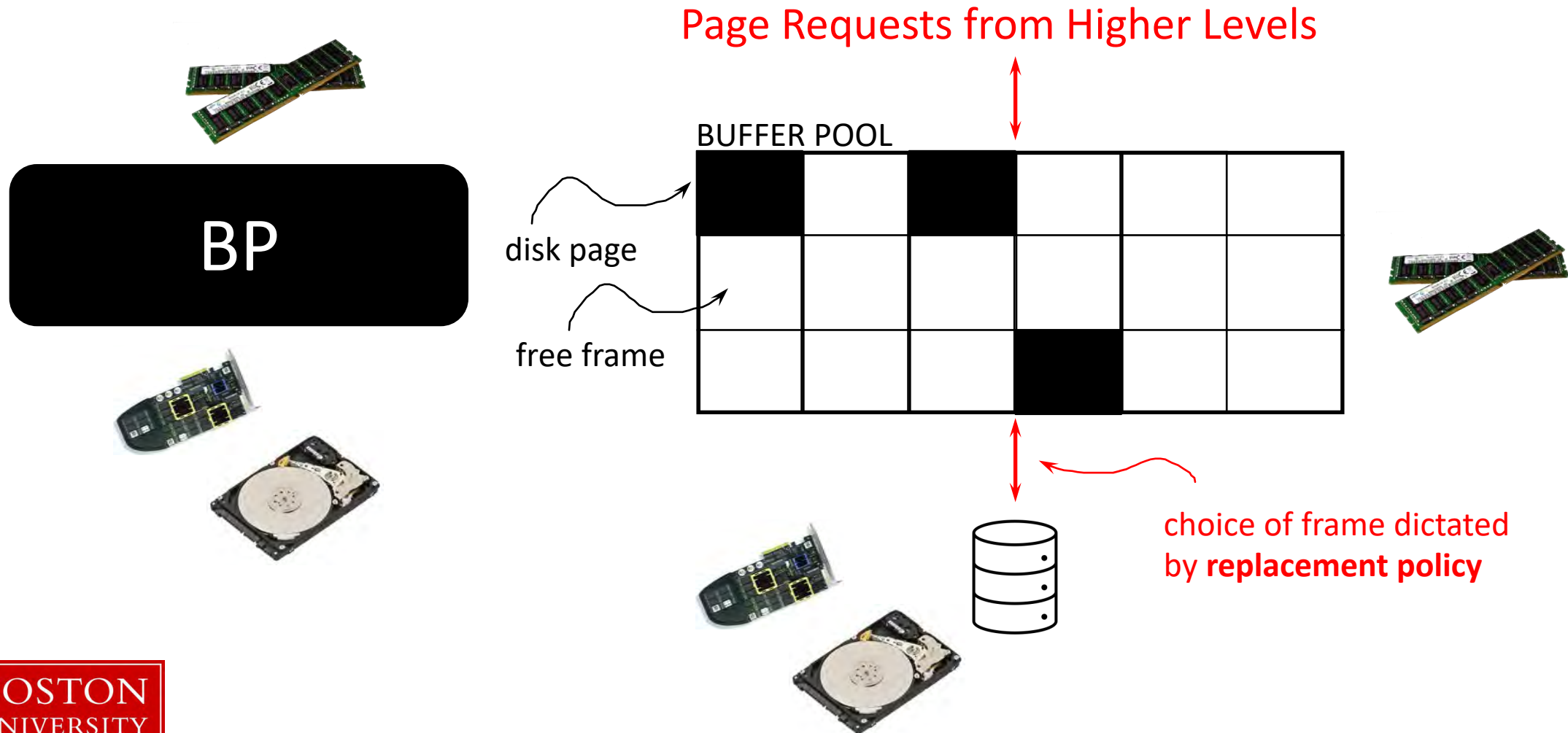**min** ***cost****(Workload, Data, Tuning)*
given *Workload* and *Data*


what if workload information is a bit wrong?


robust optimization (come and find me)
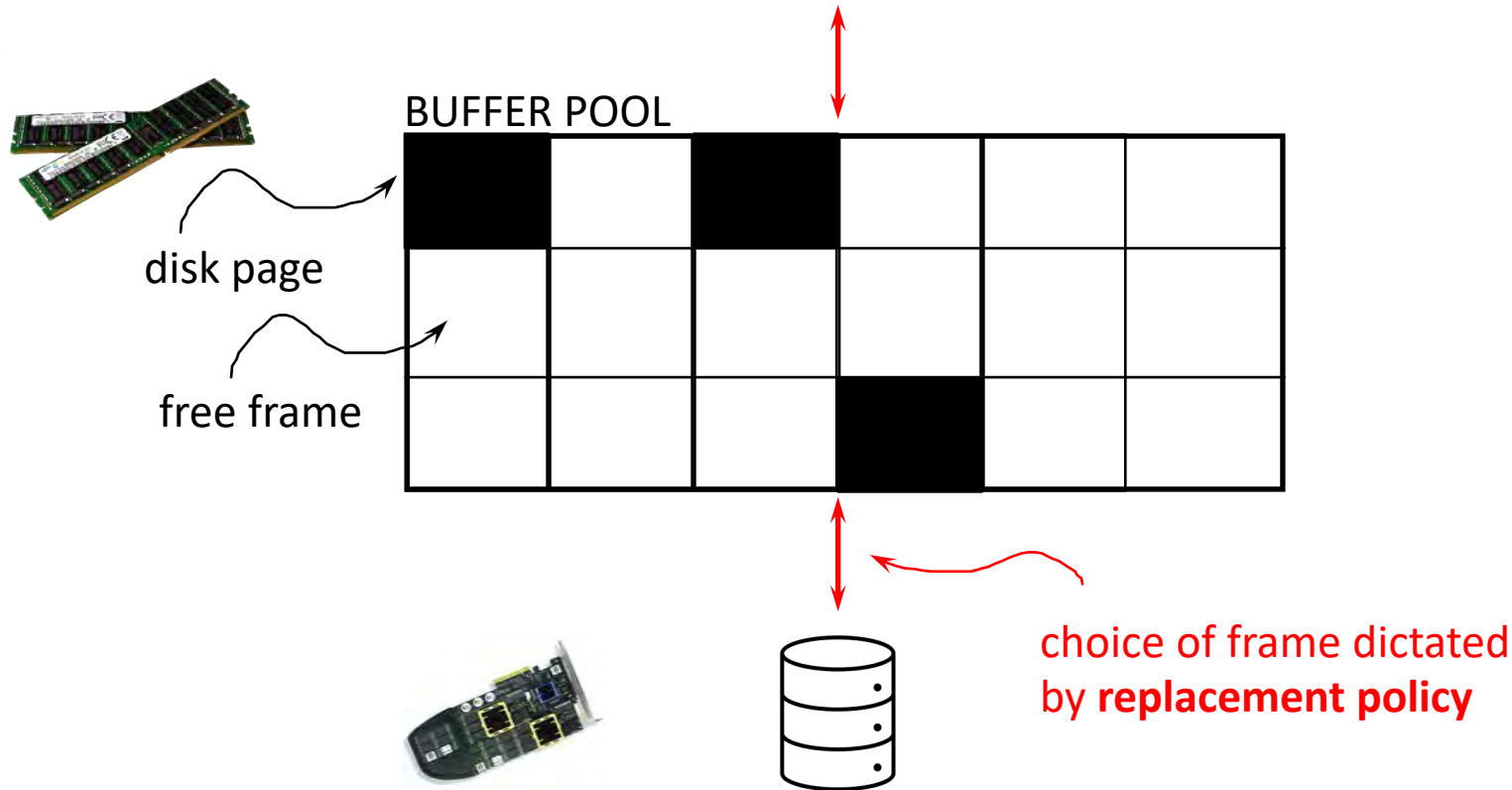
# Asynchronous Bufferpool

**what is the bufferpool?**

**BP**

**Page Requests from Higher Levels**

BUFFER POOL

disk page

free frame

choice of frame dictated by **replacement policy**

BOSTON UNIVERSITY

# Systems Project: Bufferpool

## Page Requests from Higher Levels

BUFFER POOL

disk page

free frame

choice of frame dictated by **replacement policy**

## *Implementation of a bufferpool*

- **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 …*

BOSTON UNIVERSITY

# Asynchronous Bufferpool

**what is the bufferpool?**

**BP**

manages available memory
reads/writes from/to disk

what happens when full?

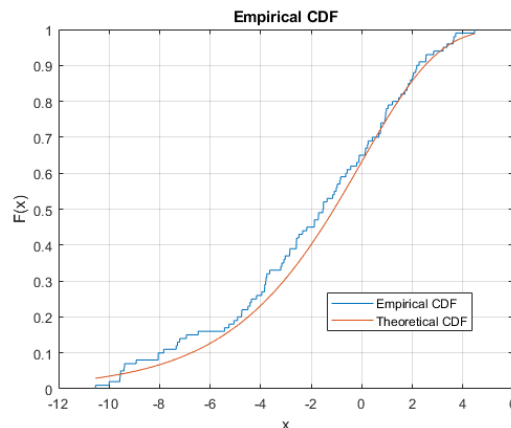writes one page back and reads on page

is this optimal?

# what is an index?



**sorted data**

1 1 1 2 3 5 10 11 12 13 18 19 20 50 54 58 62 98 101 102



$$position(val) = CDF(val) \cdot array\_size$$

can you learn the CDF?

what is the best way to do so?

how to update that?

# what to do now?

**systems project**
form groups of 1-2
(speak to me in OH if you want to work on your own)


**research project**
form groups of 2-3
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 1-2
(speak to me in OH if you want to work on your own)


**research project**
form groups of 2-3
pick one of the subjects & read background
material
define the behavior you will study and address
sketch approach and success metric

come to OH
finalize your project in 1-2 weeks (by Feb 14th)
submit proposal on February 21st

# CS 561: Data Systems Architectures

## Systems & Research Project

Prof. Manos Athanassoulis

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