*COSI 167A*
Advanced Data Systems

Class 3

**Data Systems Architecture**

Prof. Subhadeep Sarkar

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

# Today in COSI 167A

What's on the cards?

fundamentals of data **storage**

introduction to **row-stores** and **column-stores**
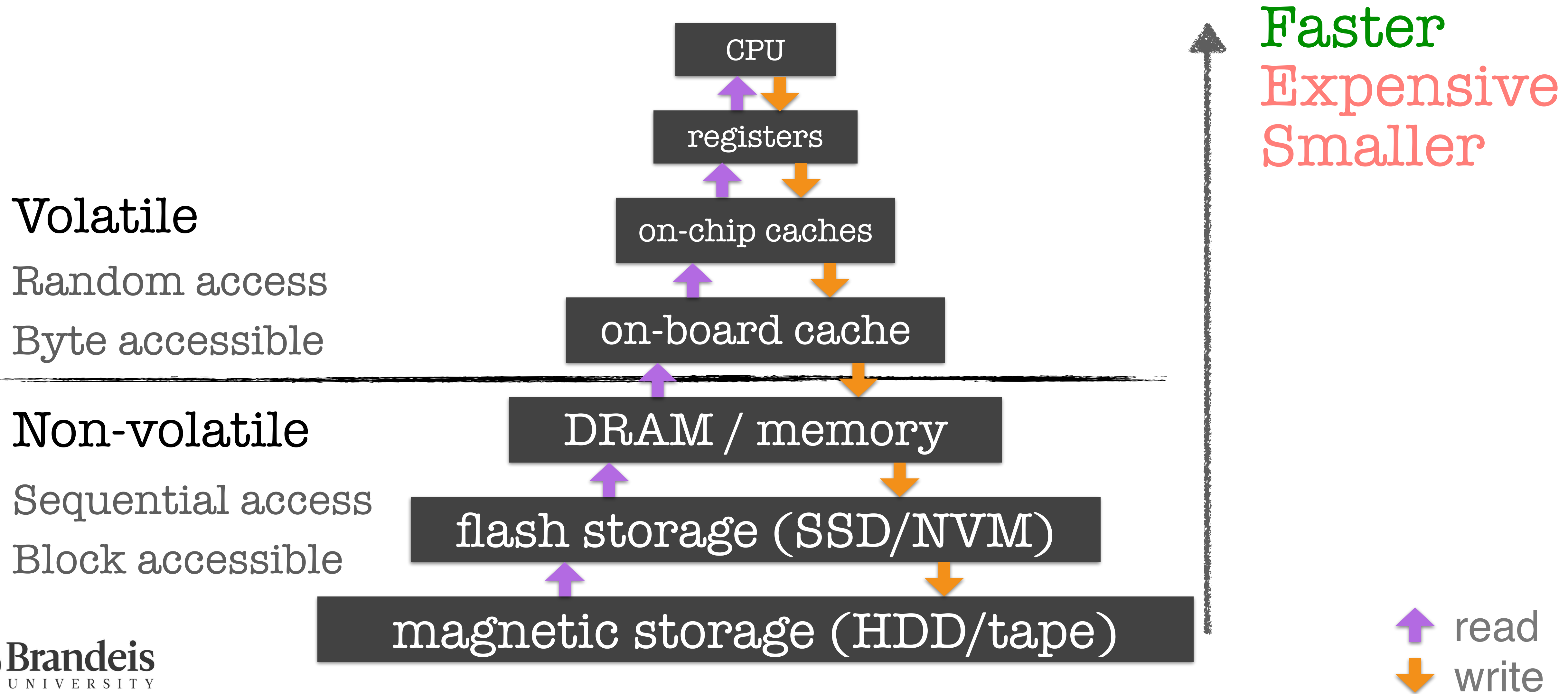
Brandeis
UNIVERSITY

# Class **logistics**
## and administrivia

**Project 1** (C++/Java) has been released (due on <span style="color:salmon">**Sep 20**</span>).

C/C++ learning resources at: https://ssd-brandeis.github.io/COSI-167A/assignments/

The **first technical question** is now available on the class website (due **before the class** on <span style="color:salmon">**Sep 10**</span>).
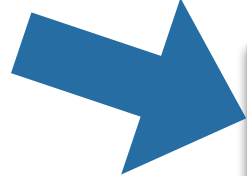
Brandeis
UNIVERSITY

# Recap: **Storage hierarchy**

How data moves!

**Volatile**

Random access

Byte accessible

**Non-volatile**

Sequential access

Block accessible

CPU

registers

on-chip caches

on-board cache

DRAM / memory

flash storage (SSD/NVM)

magnetic storage (HDD/tape)

**Faster**
**Expensive**
**Smaller**

read
write

Brandeis
UNIVERSITY

# Memory wall

Try not to jump the wall

computations
happen here

CPU

1 ns register

4 ns on-chip cache

be careful when you go below the green line

10 ns on-board cache

100 ns memory

flash storage

magnetic storage

Brandeis
UNIVERSITY

# Memory wall

Try not to jump the wall

computations
happen here

CPU

be careful when you go below the green line

1 ns      register

4 ns      on-chip cache

10 ns     on-board cache

100 ns    memory

          flash storage

          magnetic storage



performance

CPU ~20-25% perf increase annually

DRAM ~2-11% perf increase annually

time

# Memory wall

Try not to jump the wall

computations happen here

| CPU |
|---|

**1 ns** | register |

**4 ns** | on-chip cache |

**10 ns** | on-board cache |

**100 ns** | memory |

| flash storage |

| magnetic storage |

be careful when you go below the green line

doesn't matter how much faster CPUs become!

memory wall

performance

CPU perf increase

DRAM perf increase

time

old times

Brandeis
UNIVERSITY

# Memory wall

Try not to jump the wall

computations
happen here

| | |
|---|---|
| | CPU |
| **1 ns** | register |
| **4 ns** | on-chip cache |
| **10 ns** | on-board cache |
| **100 ns** | memory |
| **16,000 ns** | flash storage |
| **2 ms** | magnetic storage |

be careful when you go below the green line

be VERY careful when you go below the red line

Brandeis
U N I V E R S I T Y

# Recap: **Storing** data

Things to keep in mind

Disk is 6 orders of magnitude slower than CPU

SSDs are 4 orders of magnitude slower

Memory is 3 orders of magnitude slower

Brandeis
UNIVERSITY

# Recap: **Random** vs. **Sequential** access

So, be VERY careful!

**Avoid disk accesses** (reads/writes) whenever possible

**I/Os** to secondary storage is *always* **slow**!

Sequential access

read **each block exactly once**; process it; discard it; read next block

modern hardware can predict and **prefetch**; maximize performance

Random access

read a block; process it **partially**; discard it; may **read** the block **again**

often leads to **read amplification**

Brandeis UNIVERSITY

# Project 1

Testing the waters!

## Implementing a Simple Zone Map

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

**heap** file

zone map

w/o ZM: queries take 4 I/Os

with ZM:    query: x<4: 3 I/Os

x<12: 4 I/O

x=1: 1 I/O

x=20: 4 I/O

Thought Experiment

Are **zone maps** more or less useful if data is **sorted**?

# Project 1

Testing the waters!

## Implementing a Simple Zone Map

page 0 → 1 2 3 5 6 ← 1, 6

page 1 → 7 8 9 10 11 ← 7, 11

page 2 → 12 13 14 15 16 ← 12, 16

page 3 → 21 23 24 31 34 ← 21, 34

**sorted** file     zone map

w/o ZM: queries take 4 I/Os

with ZM:     query: x<4: 1 I/O

x<12: 2 I/Os

x=1: 1 I/O

x=20: 0 I/Os

Sorting is inherent indexing!

Zone map is a second index!

BRANDEIS
UNIVERSITY

How we store (write) data heavily determines the performance of the system

# The **design goals**

Biulding "efficient" data systems

build databases that can **write** data **fast** …

… and **process/analyze** that data **quickly**

Well, just get to work then!

# Performance **tradeoff**

The tug of war

hashmap

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

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

sorted array

R: $\mathcal{O}(logN)$

U: $\mathcal{O}(N)$

log

R: $\mathcal{O}(N)$

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

Query / **R**ead

Insert / **U**pdate

# Performance **tradeoff**

The tug of war

Storage / **M**emory

hashmap

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

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

M

sorted array

R: $\mathcal{O}(logN)$

U: $\mathcal{O}(N)$

M

Query / **R**ead

Insert / **U**pdate

log

R: $\mathcal{O}(N)$

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

M

# Performance **tradeoff**

The tug of war

Storage / **M**emory

**hashmap**

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

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

**M**

# There is **NO** perfect data structure!

**sorted array**

R: $\mathcal{O}(logN)$

U: $\mathcal{O}(N)$

**M**

Query / **R**ead

Insert / **U**pdate

**log**

R: $\mathcal{O}(N)$

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

**M**

# RUM conjecture

A three-way tradeoff

Storage / **M**emory

Query / **R**ead        Insert / **U**pdate

# RUM conjecture

A three-way tradeoff

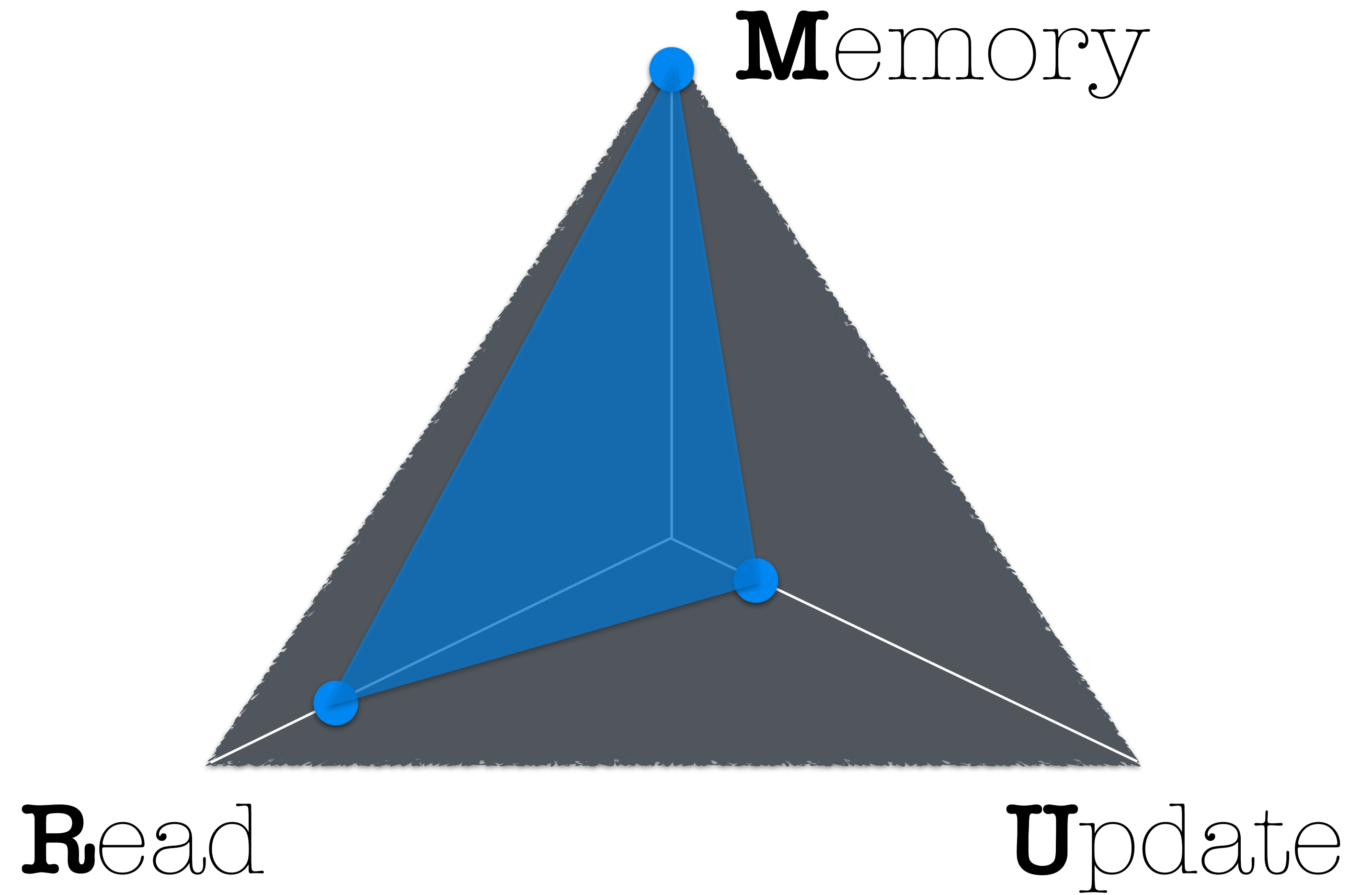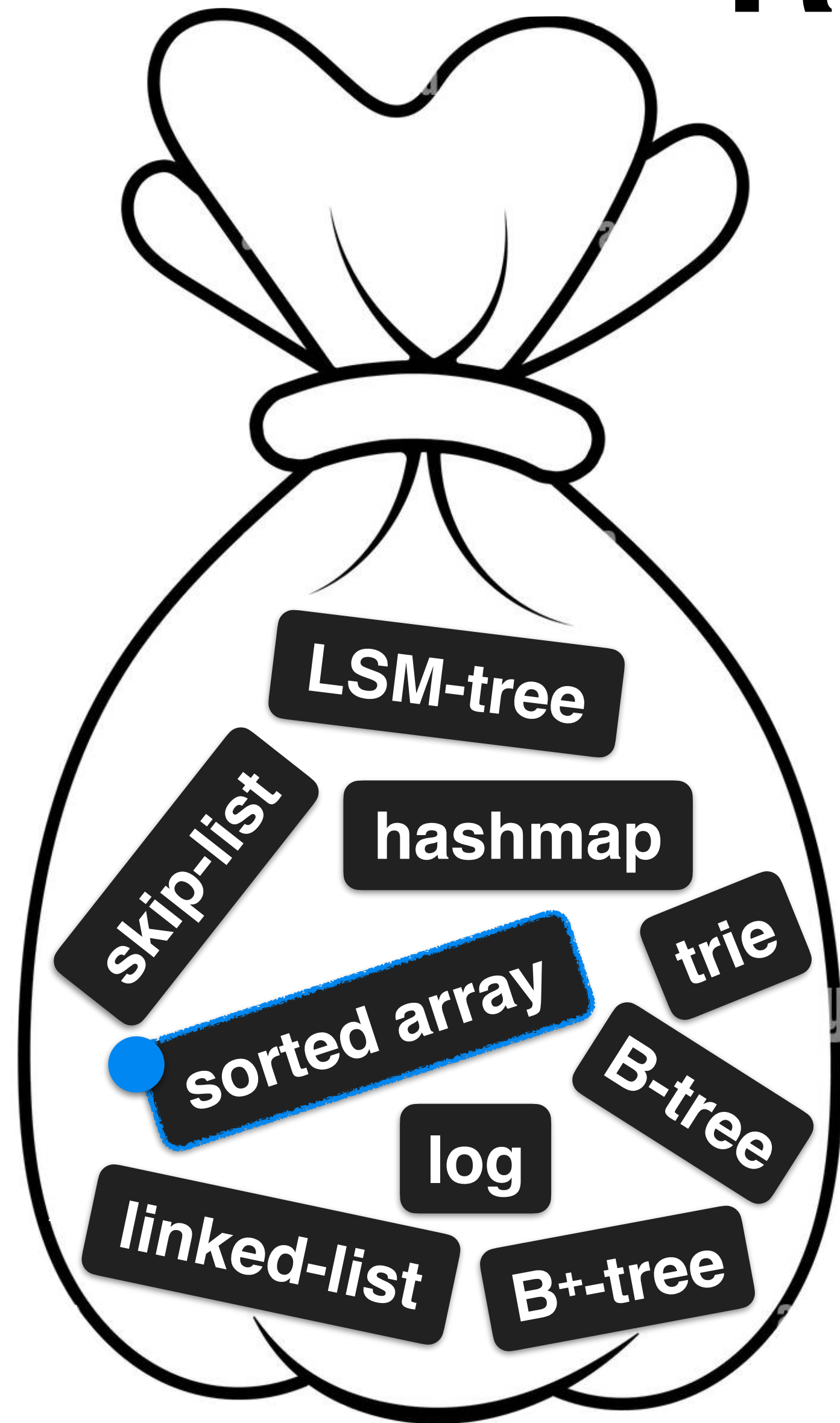**M**emory

... setting an <span style="color:green">upper bound</span> for <span style="color:green">two of the RUM axes</span>, implies a <span style="color:red">hard lower bound</span> for the <span style="color:red">third axis</span>

**R**ead

**U**pdate

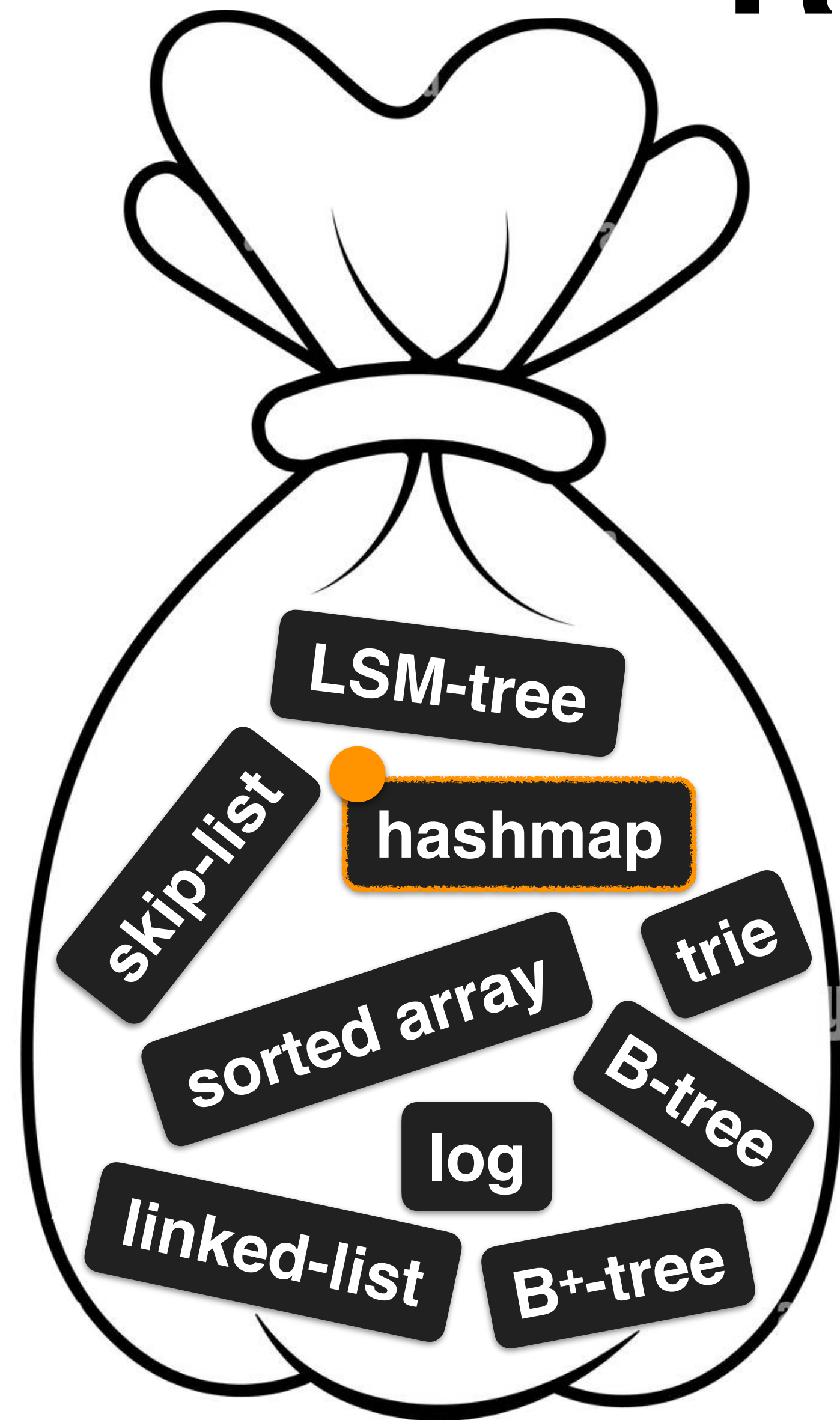# RUM conjecture

A three-way tradeoff

# RUM conjecture

A three-way tradeoff

**M**emory

**R**ead

**U**pdate

LSM-tree

skip-list

hashmap

trie

sorted array

B-tree

log

linked-list

B+-tree

# RUM conjecture

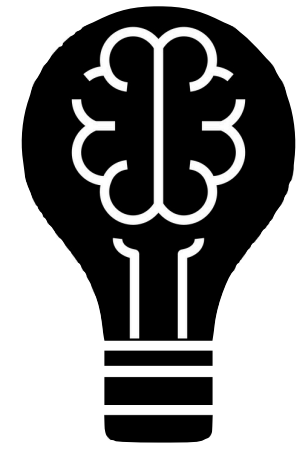A three-way tradeoff

# RUM conjecture

A three-way tradeoff



Memory

Read

Update

LSM-tree

skip-list

hashmap

sorted array

trie

log

B-tree

linked-list

B+-tree

# RUM conjecture

A three-way tradeoff

Thought Experiment 1

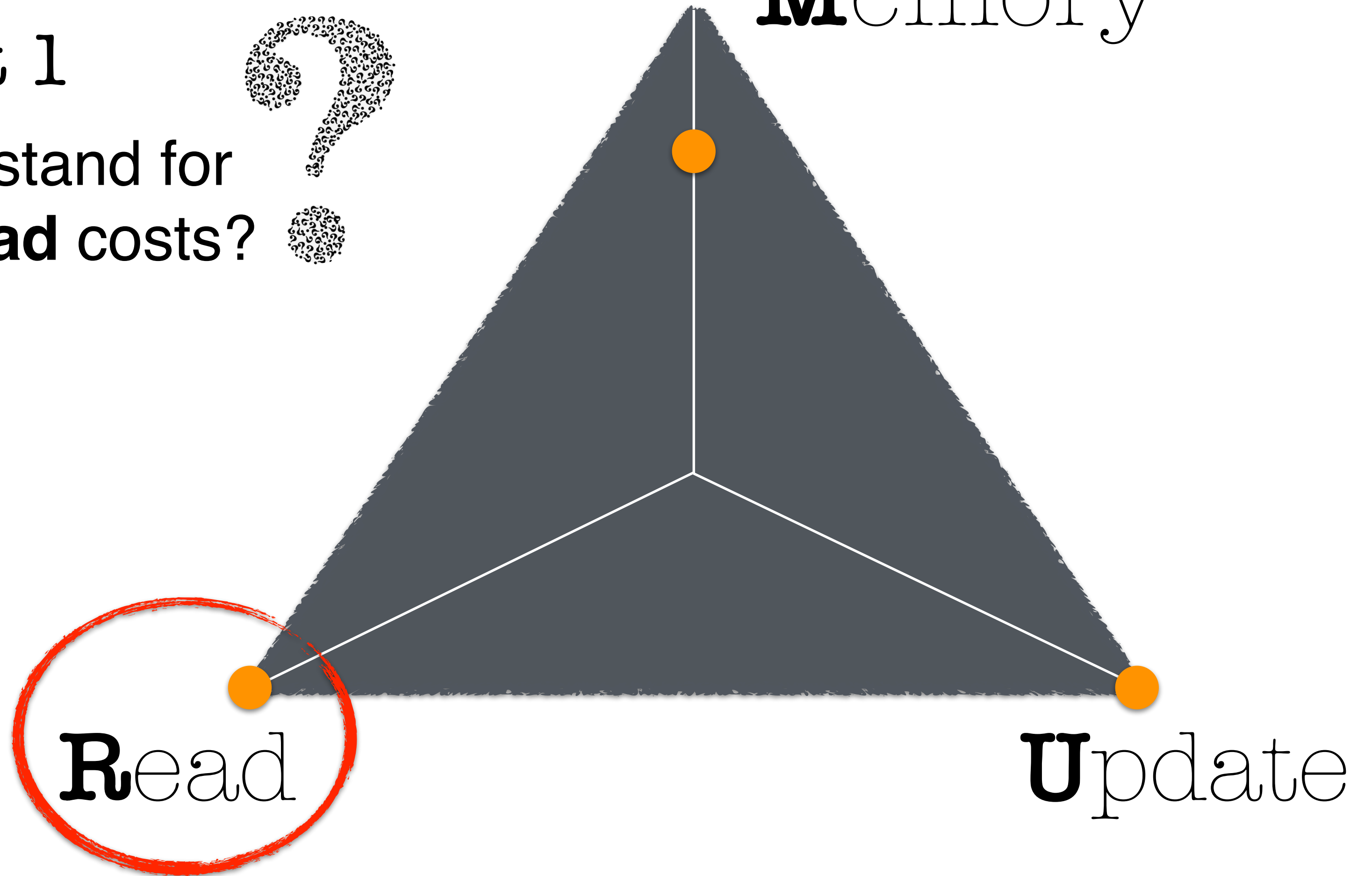Where does a **hashmap** stand for **memory**, **update**, and **read** costs?
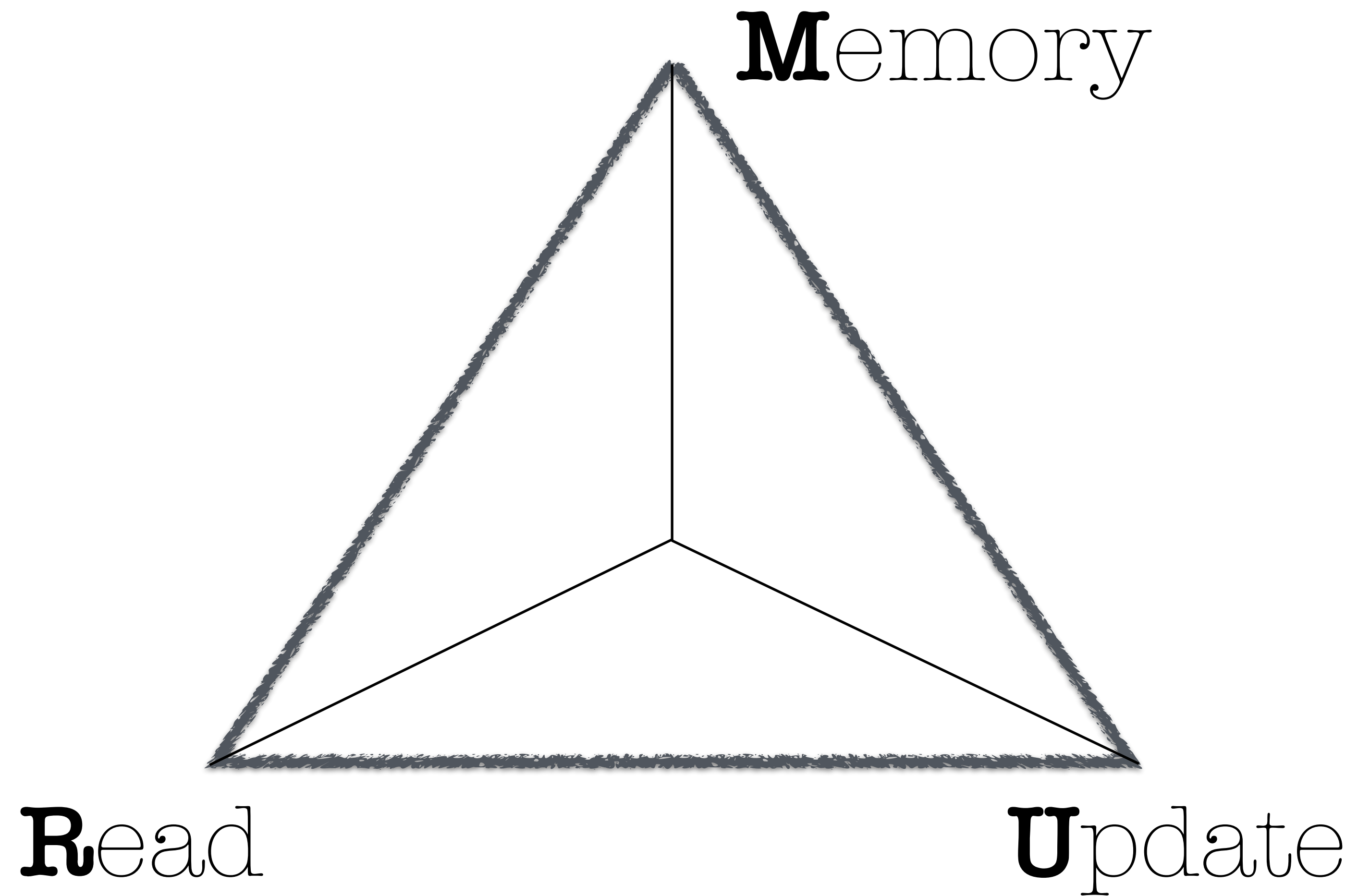
**M**emory

point read
vs
range read

**R**ead

**U**pdate

# RUM conjecture

A three-way tradeoff
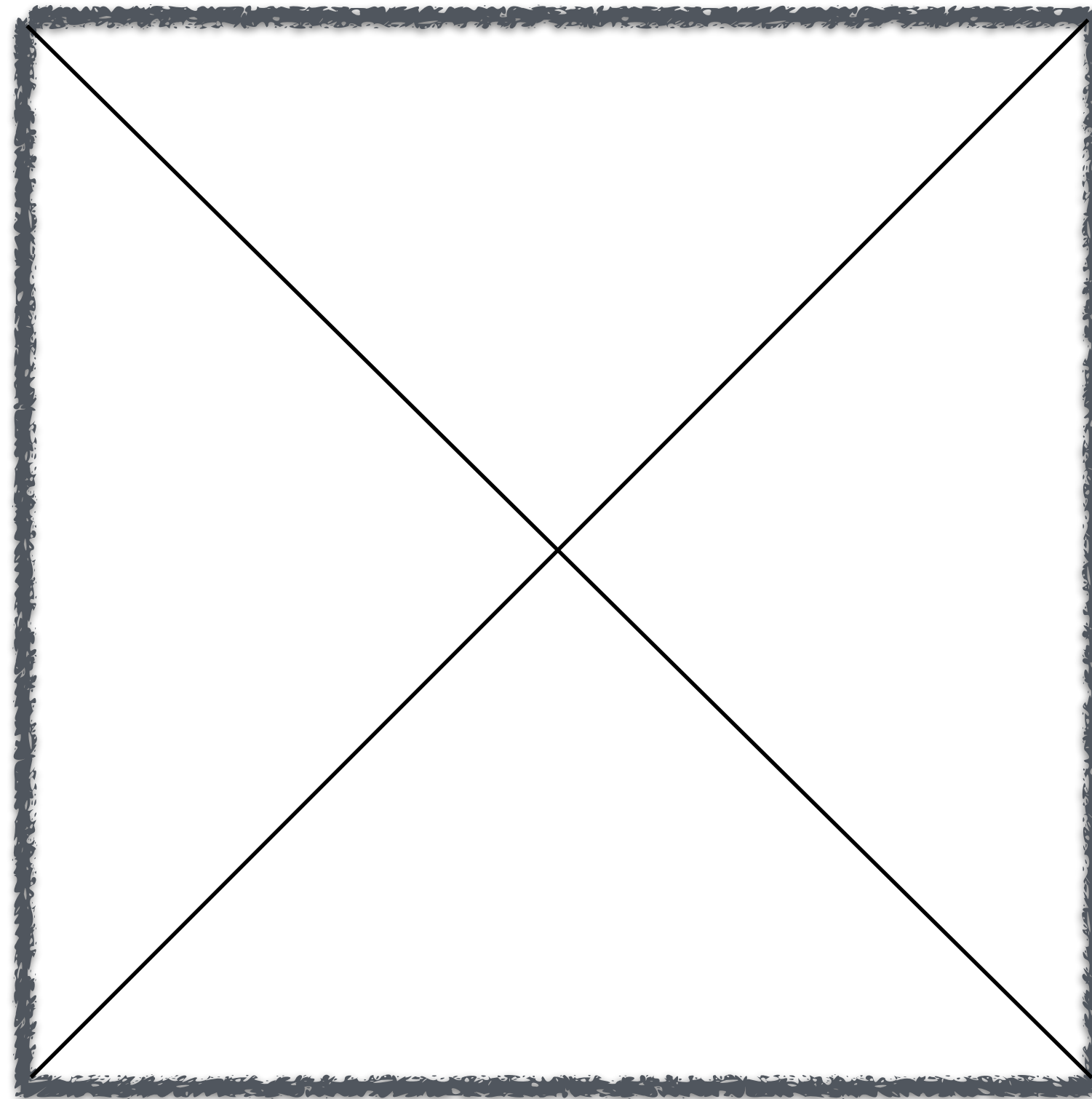


**M**emory

**R**ead **U**pdate

# Extending the **RUM tradeoff**

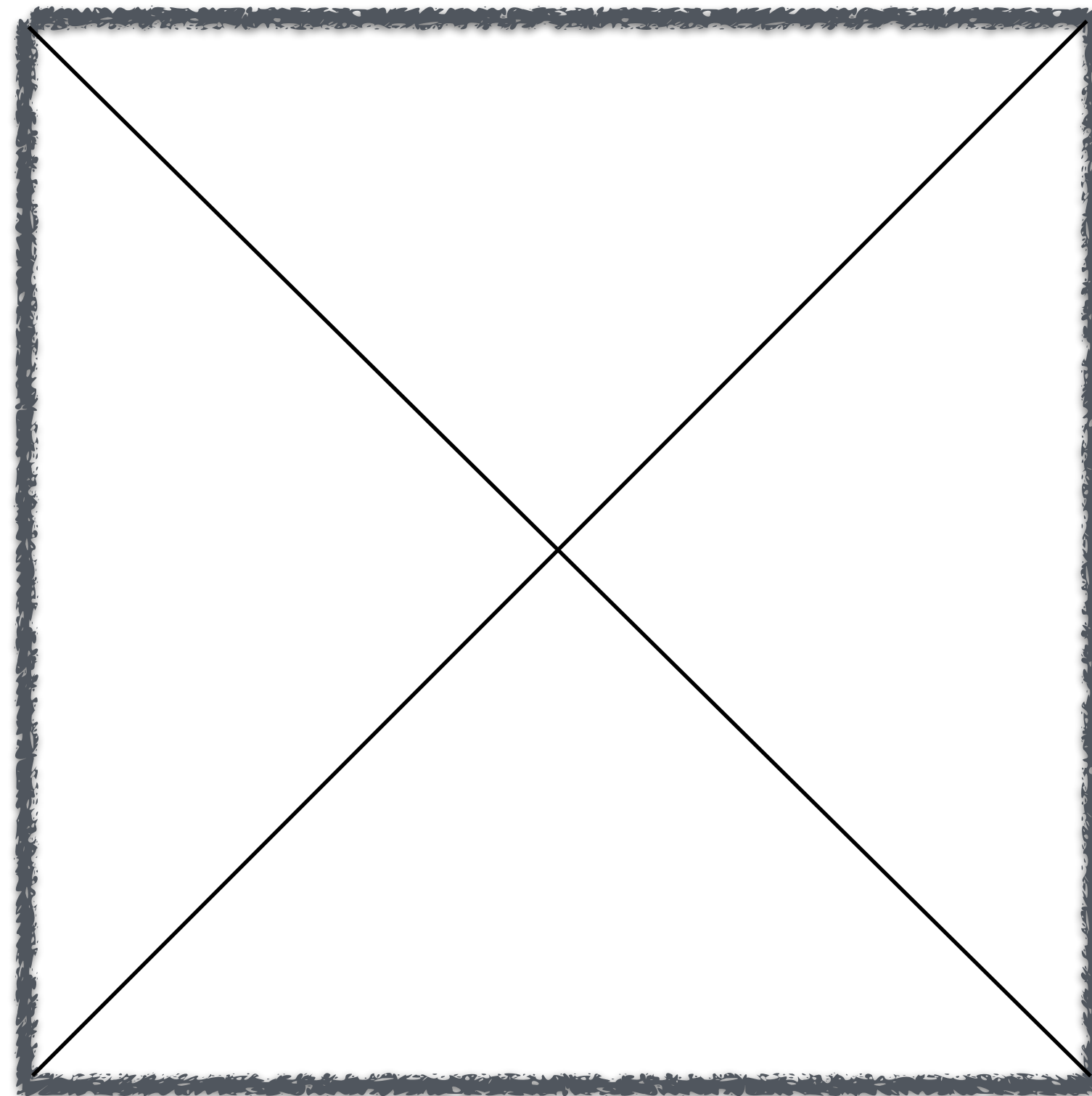A multi-way tradeoff

Point read

Memory

Range read

Update

# Extending the **RUM tradeoff**

### A multi-way tradeoff

Point read

Memory

But, what about **deletes**?

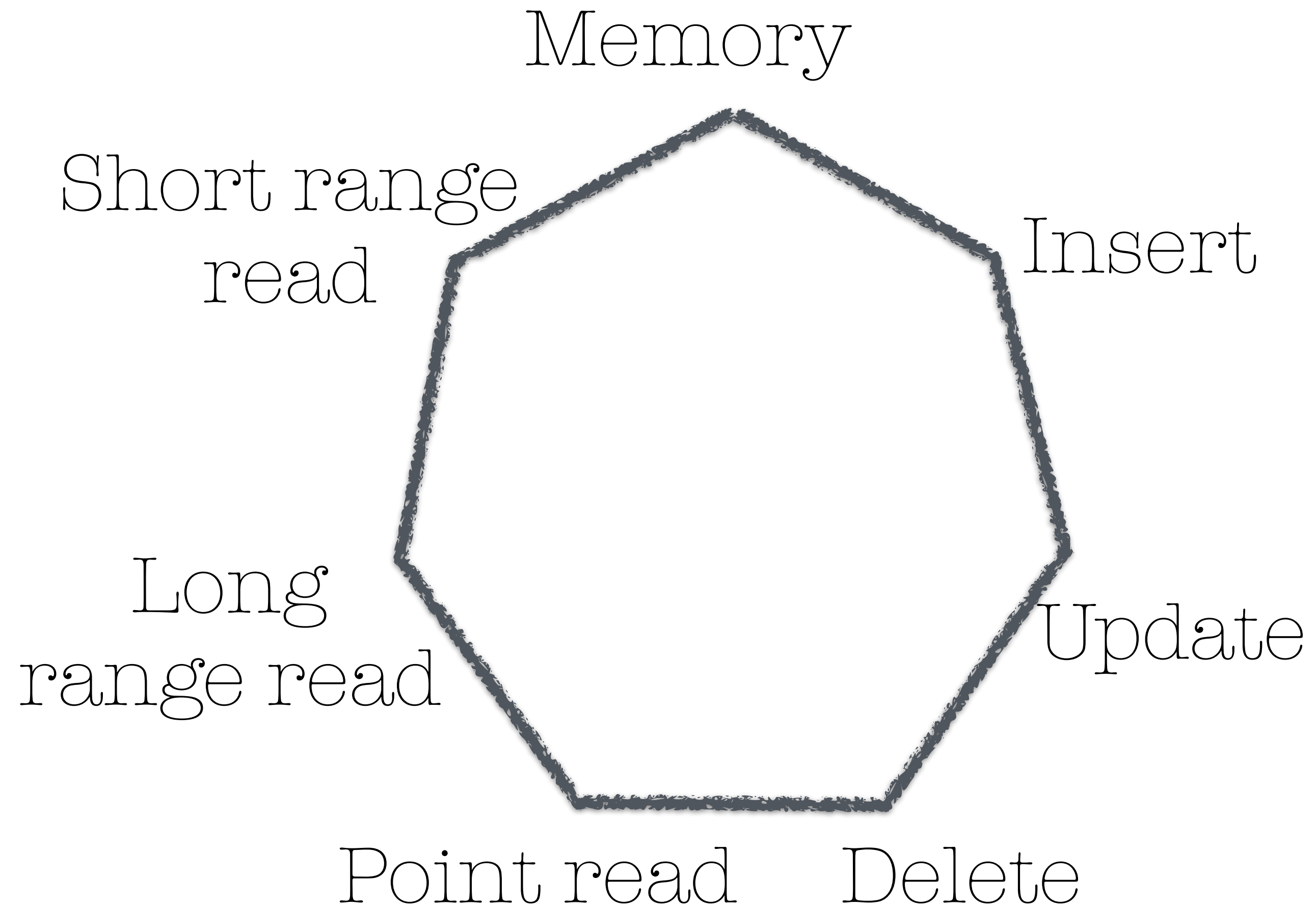Are **updates** the same as **inserts**?

What about **selectivity** of range reads?

Range read

Update

# Extending the **RUM tradeoff**

A multi-way tradeoff



But, what about **deletes**?

Are **updates** the same as **inserts**?

What about **selectivity** of range reads?

cloud cost

hardware

performance tradeoffs

Designing data systems = HARD PROBLEM

index design

access method

application requirements

# Designing data systems

Solving a hard problem



Ask for the HiPPO

Highest Paid Person's Opinion

# Designing data systems

Solving a hard problem

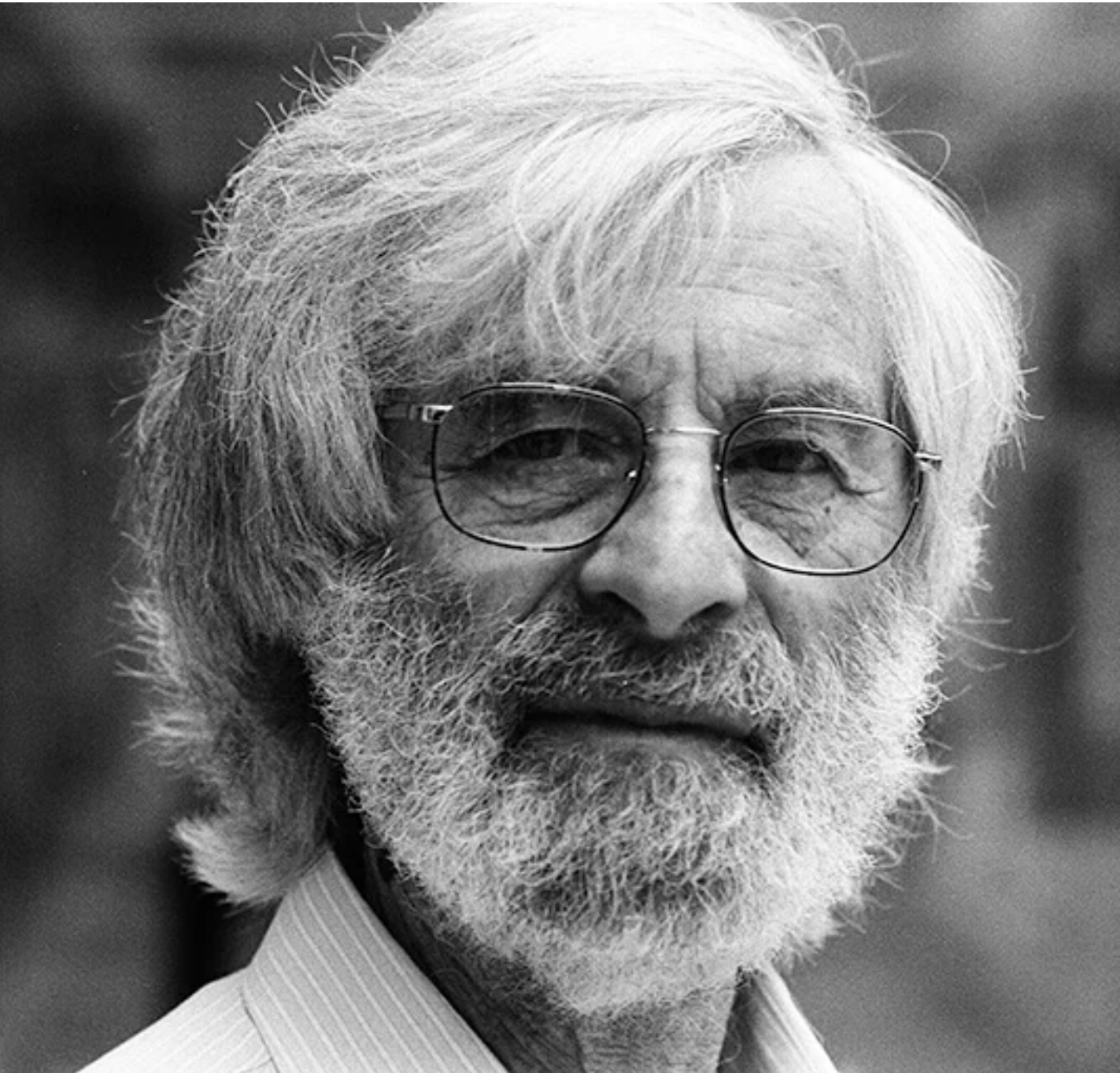No one knows everything!

Ask for the HiPPO

Highest Paid Person's Opinion

Brandeis
UNIVERSITY

# Designing data systems

Solving a hard problem

No one knows everything!

The Part-Time Parliament

LESLIE
Digital E

Recent ar
despite th
consistent
and the f
way of im

Categorie
uted Sys

**Paxos Made Simple**

Leslie Lamport

01 Nov 2001

**Abstract**

The Paxos algorithm, when presented in plain English, is very simple.

**Lesley Lamport**, Microsoft Research
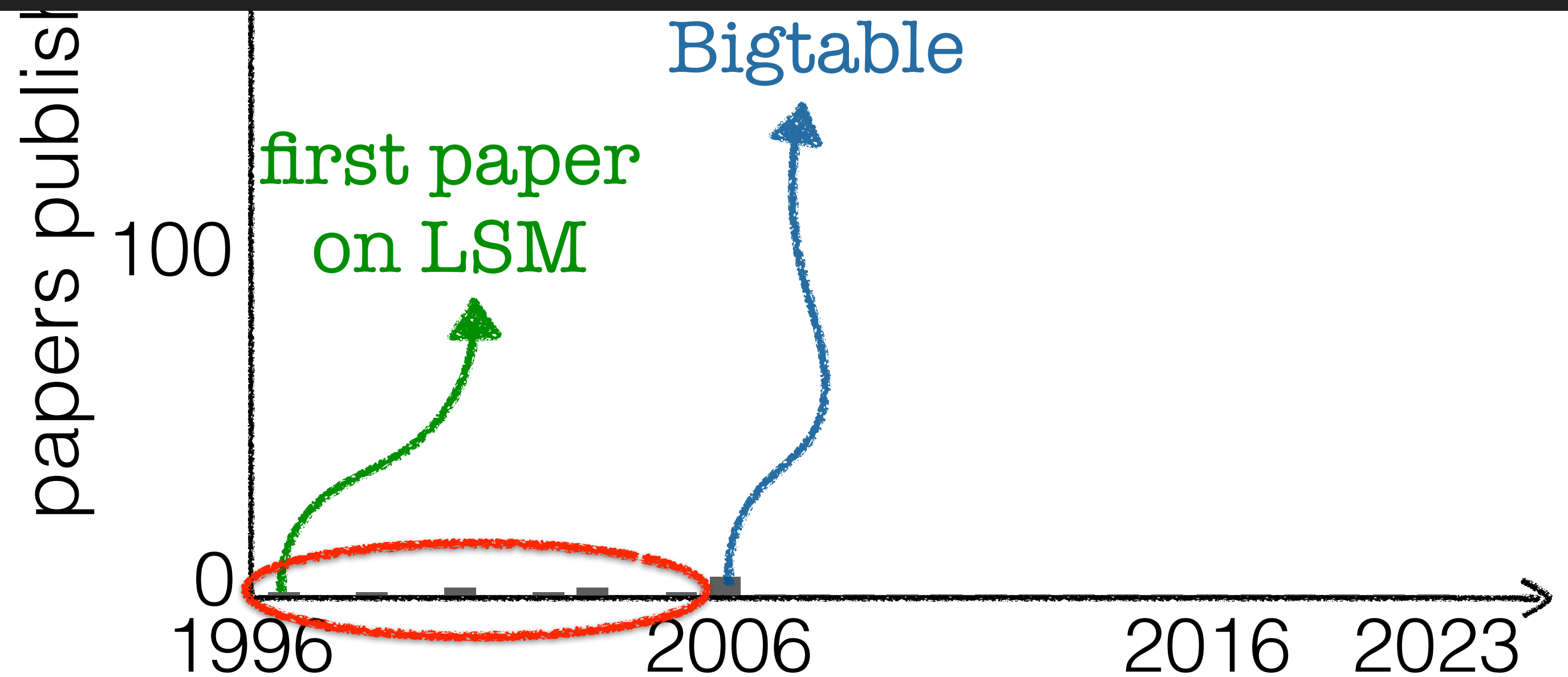MA '63, PhD '72 Brandeis University
ACM Turing Award 2013

Brandeis
UNIVERSITY

# Designing data systems

Solving a hard problem

No one knows everything!



**Patrick O'Neil**, UMass Boston
Co-inventor of the LSM-tree

papers publish... 100 0

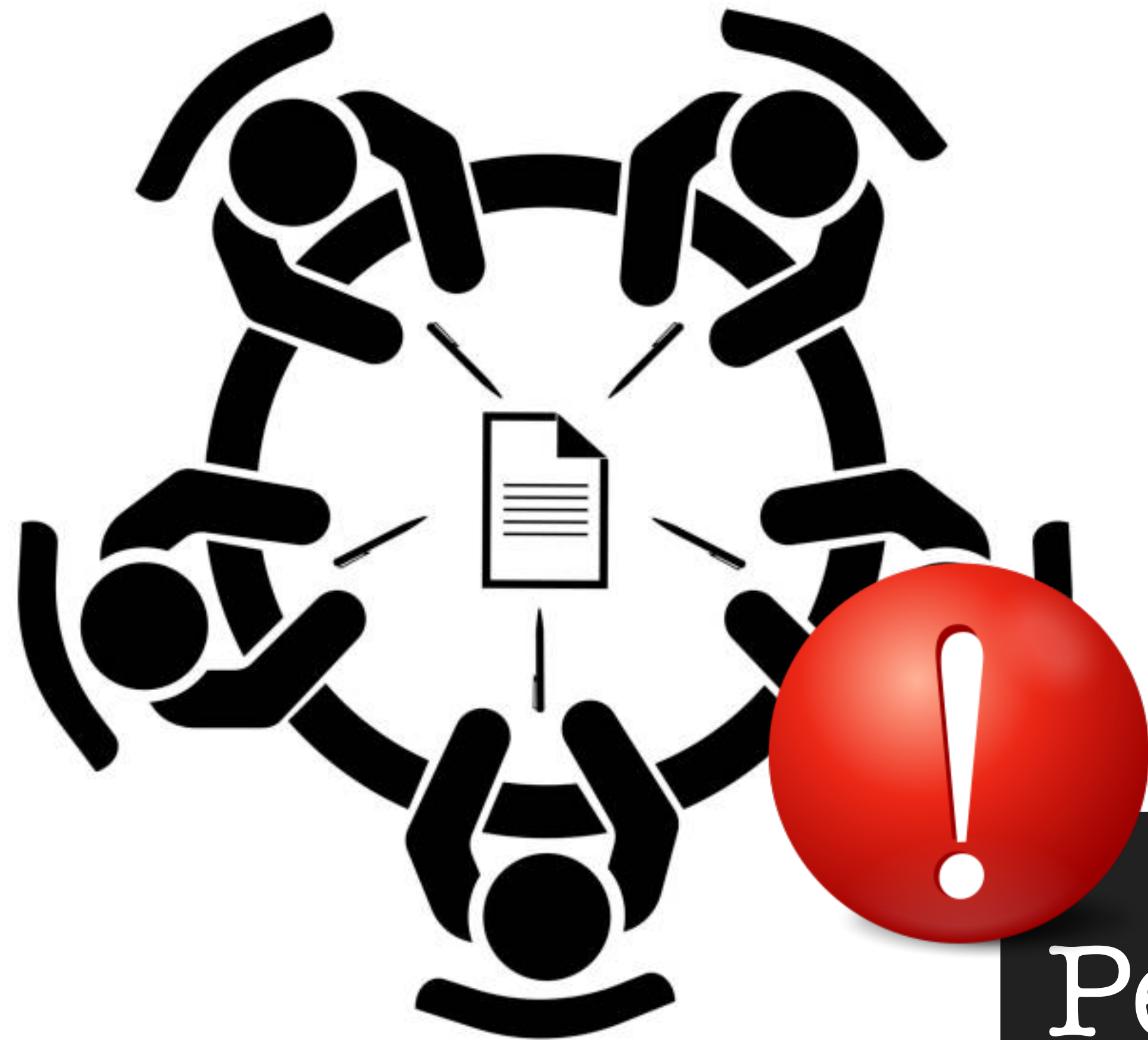first paper on LSM

Bigtable

1996    2006    2016  2023

Brandeis
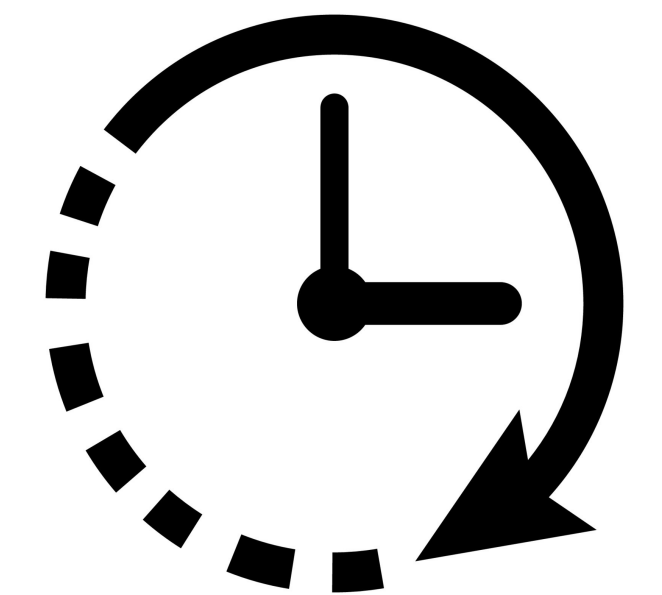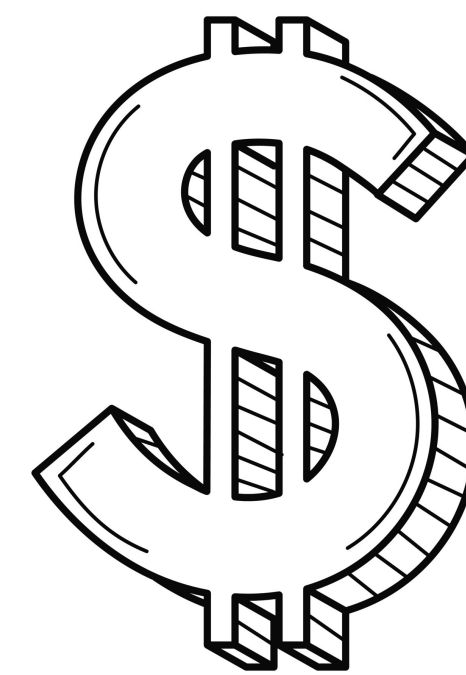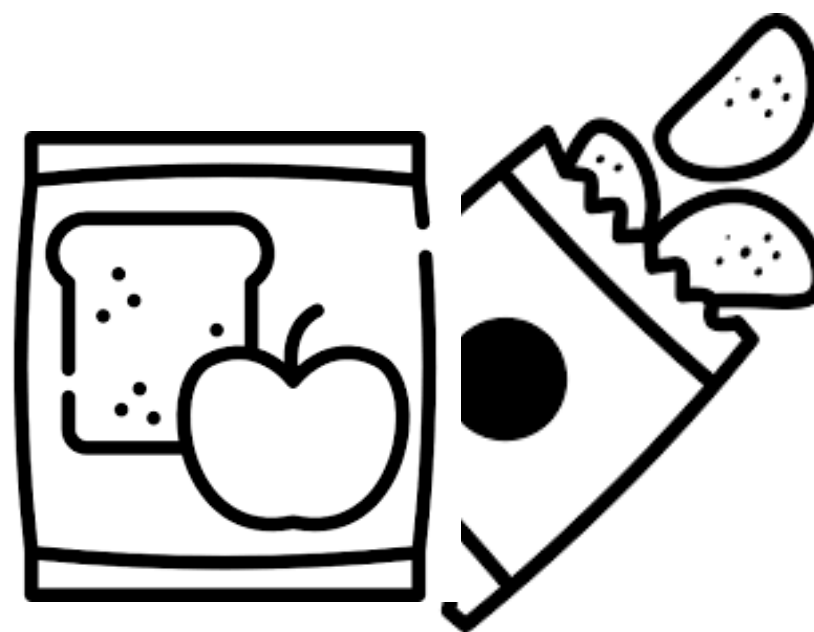UNIVERSITY

* data from Google scholar

# Designing data systems

Solving a hard problem

Get the COSI 167 students

Perfectly tuned data system!

# Designing data systems

Solving a hard problem

Manual designing &
hand-tuning hundreds of knobs
do NOT SCALE

# Designing data systems

Solving a hard problem

Manual designing &
hand-tuning hundreds of knobs
do NOT SCALE

Self-designing, self-tuning, and
adaptive data systems

# Designing data systems

Solving a hard problem

Self-designing, self-tuning, and adaptive data systems

**Adaptive** data layouts: row stores vs. column stores vs. hybrids

**Self-tuning** LSM-engines

**Adaptive** (**adaptive**) indexing

**Workload-aware** data re-organization
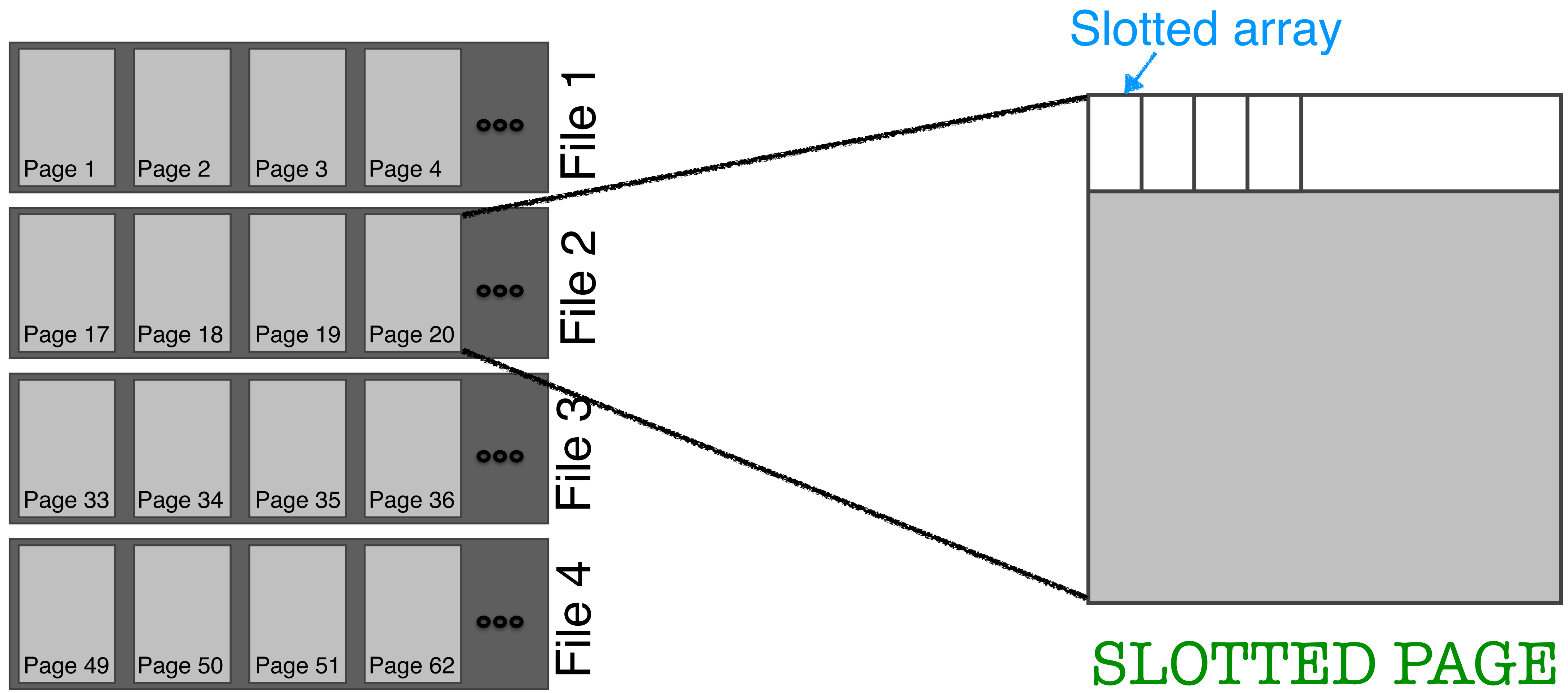
**Self-designing** storage engines

**Hardware-conscious** memory manager

**Cloud cost-optimized** data systems
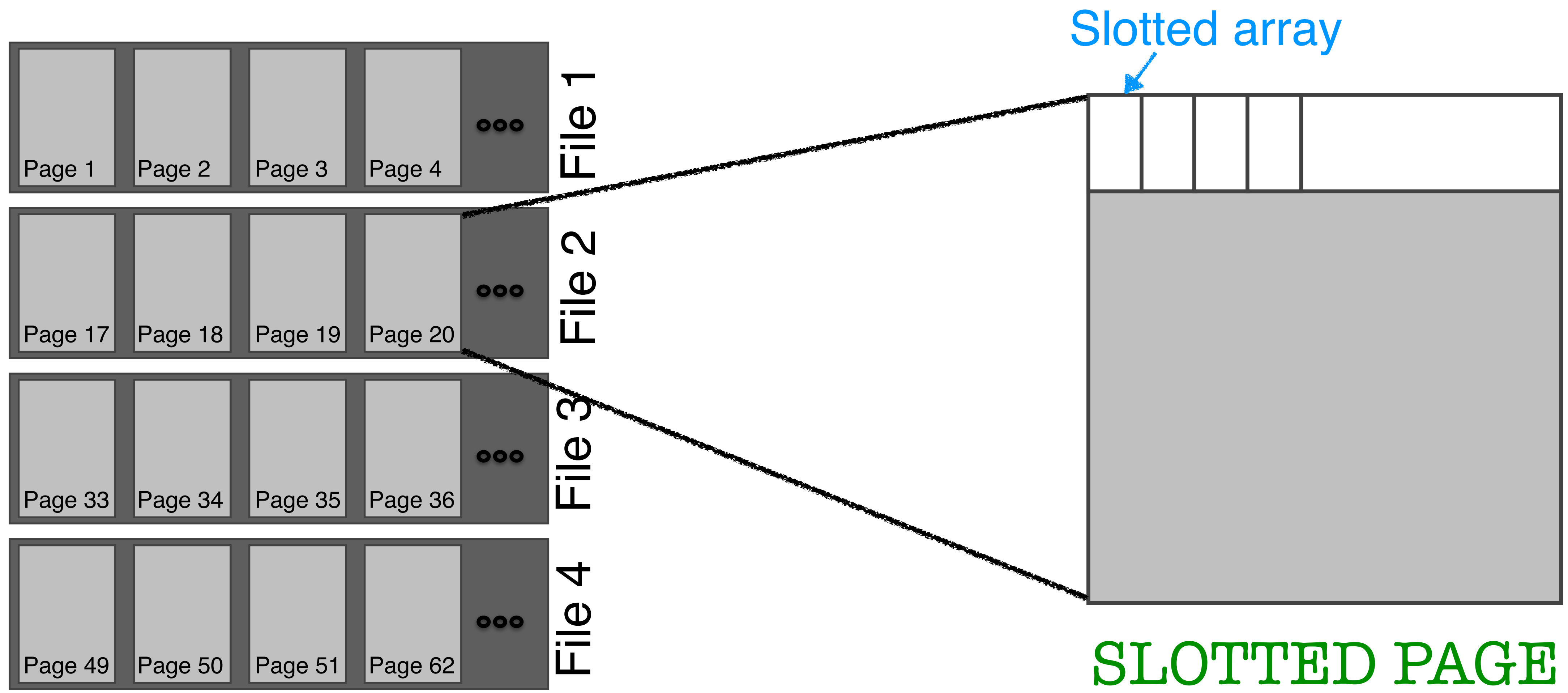
# Understanding data placement

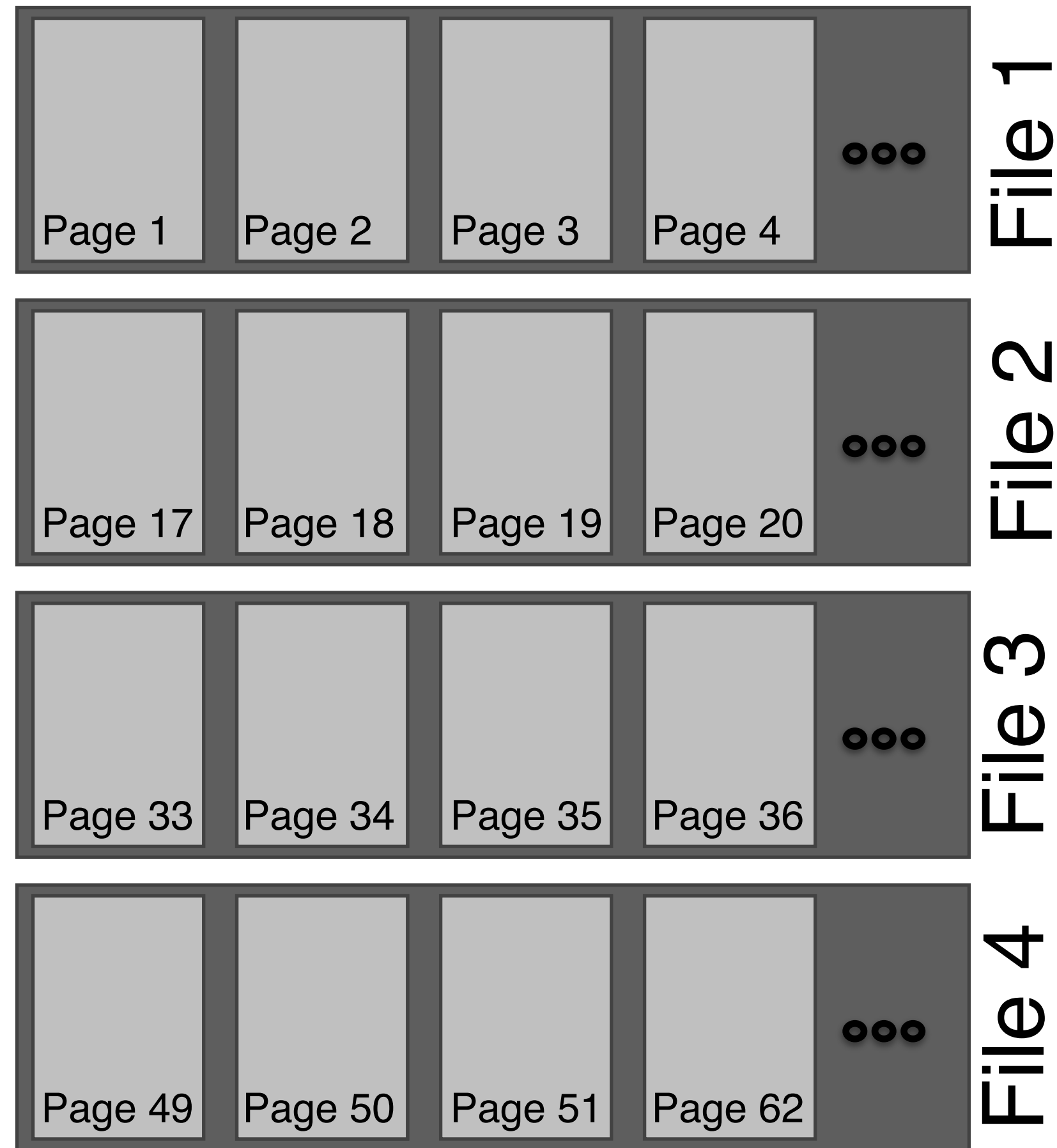# Understanding **data placement**

Files, pages, tuples

File 1

Page 1    Page 2    Page 3    Page 4    ●●●

File 2

Page 17    Page 18    Page 19    Page 20    ●●●

File 3

Page 33    Page 34    Page 35    Page 36    ●●●

File 4

Page 49    Page 50    Page 51    Page 62    ●●●

Slotted array

SLOTTED PAGE

# Understanding **data placement**

Files, pages, tuples



File 1

Page 1    Page 2    Page 3    Page 4    ●●●

File 2

Page 17    Page 18    Page 19    Page 20    ●●●

File 3

Page 33    Page 34    Page 35    Page 36    ●●●

File 4

Page 49    Page 50    Page 51    Page 62    ●●●

Slotted array

SLOTTED PAGE

Brandeis
UNIVERSITY

# Understanding **data placement**

Files, pages, tuples

| | | | | |
|---|---|---|---|---|
| Page 1 | Page 2 | Page 3 | Page 4 | ●●● | File 1
| Page 17 | Page 18 | Page 19 | Page 20 | ●●● | File 2
| Page 33 | Page 34 | Page 35 | Page 36 | ●●● | File 3
| Page 49 | Page 50 | Page 51 | Page 62 | ●●● | File 4

Slotted array

insert(Tuple 1)

SLOTTED PAGE

# Understanding **data placement**

Files, pages, tuples



File 1

Page 1　Page 2　Page 3　Page 4

File 2

Page 17　Page 18　Page 19　Page 20

File 3

Page 33　Page 34　Page 35　Page 36

File 4

Page 49　Page 50　Page 51　Page 62

Slotted array

insert(Tuple 1)

Tuple 1

SLOTTED PAGE

Brandeis
UNIVERSITY

# Understanding **data placement**

Files, pages, tuples

Page 1 | Page 2 | Page 3 | Page 4 | ••• | File 1

Page 17 | Page 18 | Page 19 | Page 20 | ••• | File 2

Page 33 | Page 34 | Page 35 | Page 36 | ••• | File 3

Page 49 | Page 50 | Page 51 | Page 62 | ••• | File 4

Slotted array

insert(Tuple 1)

insert(Tuple 2)

**Tuple 2** | **Tuple 1**

SLOTTED PAGE

Brandeis
UNIVERSITY

# Understanding **data placement**

Files, pages, tuples

Page 1 | Page 2 | Page 3 | Page 4 | File 1

Page 17 | Page 18 | Page 19 | Page 20 | File 2

Page 33 | Page 34 | Page 35 | Page 36 | File 3

Page 49 | Page 50 | Page 51 | Page 62 | File 4

insert(Tuple 1)
insert(Tuple 2)
insert(Tuple 3)

Slotted array

Tuple 3

Tuple 2 | Tuple 1

SLOTTED PAGE

Brandeis
UNIVERSITY

# Understanding **data placement**

Files, pages, tuples

Slotted array

File 1

Page 1 | Page 2 | Page 3 | Page 4 | ...

File 2

Page 17 | Page 18 | Page 19 | Page 20 | ...

File 3

Page 33 | Page 34 | Page 35 | Page 36 | ...

File 4

Page 49 | Page 50 | Page 51 | Page 62 | ...

insert(Tuple 1)

insert(Tuple 2)

insert(Tuple 3)

insert(Tuple 4)

Tuple 4      Tuple 3

Tuple 2      Tuple 1

SLOTTED PAGE

Brandeis
UNIVERSITY

# Understanding **data placement**

Files, pages, tuples

| | | | | |
|---|---|---|---|---|
| Page 1 | Page 2 | Page 3 | Page 4 | ••• | File 1 |
| Page 17 | Page 18 | Page 19 | Page 20 | ••• | File 2 |
| Page 33 | Page 34 | Page 35 | Page 36 | ••• | File 3 |
| Page 49 | Page 50 | Page 51 | Page 62 | ••• | File 4 |

insert(Tuple 1)

insert(Tuple 2)

insert(Tuple 3)

insert(Tuple 4)

Slotted array

Tuple 4

Tuple 3

Tuple 2

Tuple 1

SLOTTED PAGE

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

file

| A | B | C | D |
| A | B | C | D |

| A | B | C | D |
| A | B | C | D |

| A | B | C | D |
| A | B | C | D |

| A | B | C | D |
| A | B | C | D |

○○○

💡 Thought Experiment 2

select A,B,C,D from R

💡 Thought Experiment 3

select A from R

Problem?

Read amplification

Solution?

Column stores

Brandeis
UNIVERSITY

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

file

A   B   C   D

o o o

💡 Thought Experiment 3

select A from R

💡 Thought Experiment 2

select A,B,C,D from R

Problem?

Tuple reconstruction

Brandeis
UNIVERSITY

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

file

A    B    C    D

o o o

Thought Experiment 4

select A+B from R

Can we do something better?

Brandeis
UNIVERSITY

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

file

A    B    C    D

○○○

Thought Experiment 4

select A+B from R

Can we do something better?

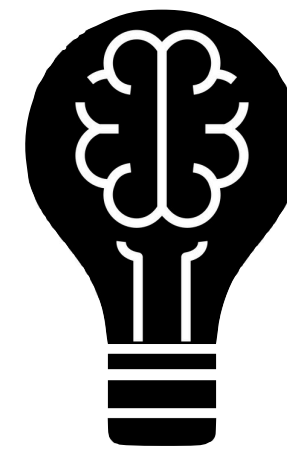# Querying over slotted pages

### Understanding the schema

schema: **R** (A, B, C, D)

file



**Thought Experiment 5**

select A+B from R
select A,B,C,D from R
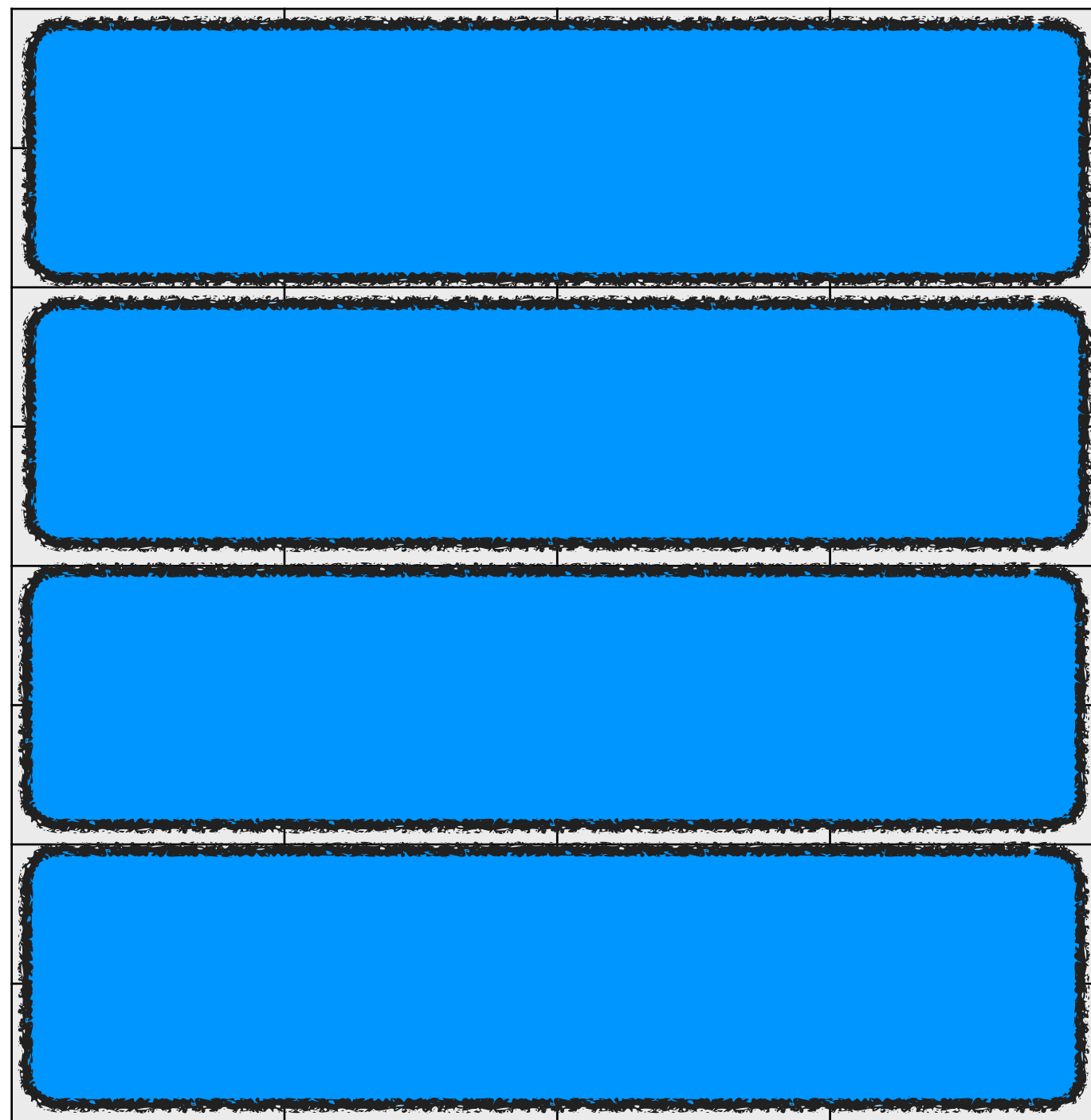select A from R

What if we have all three queries?

What if we only have inserts/updates?

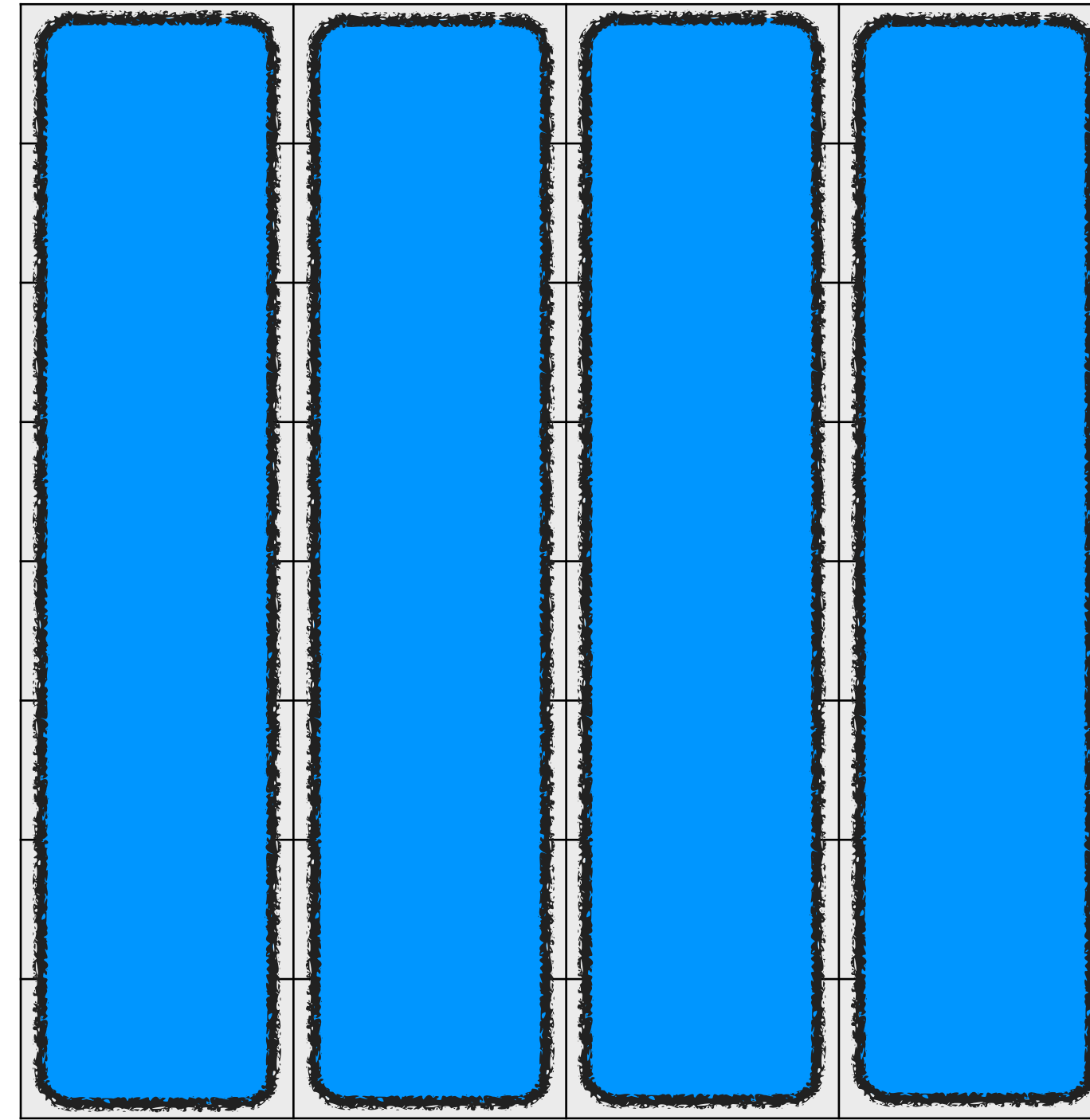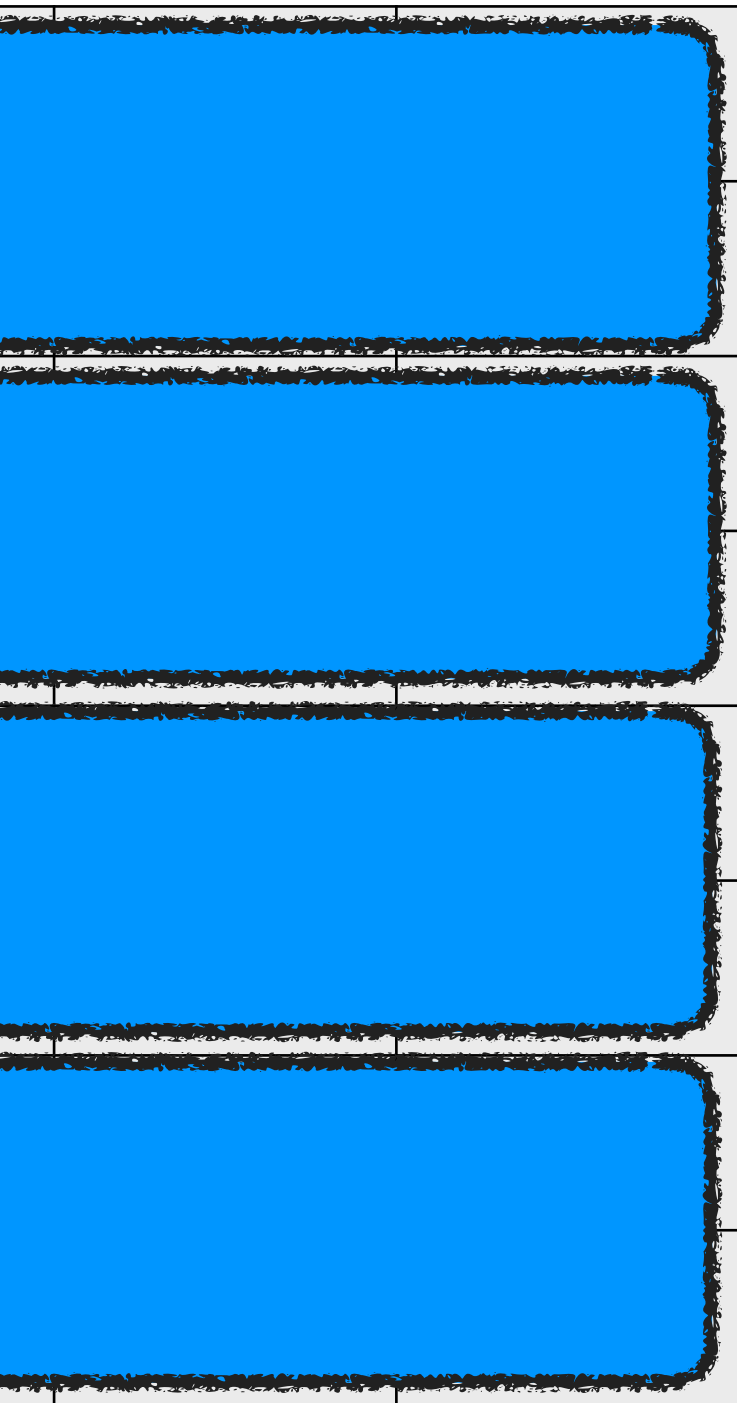Brandeis
UNIVERSITY

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

schema: **R** (A, B, C, D)
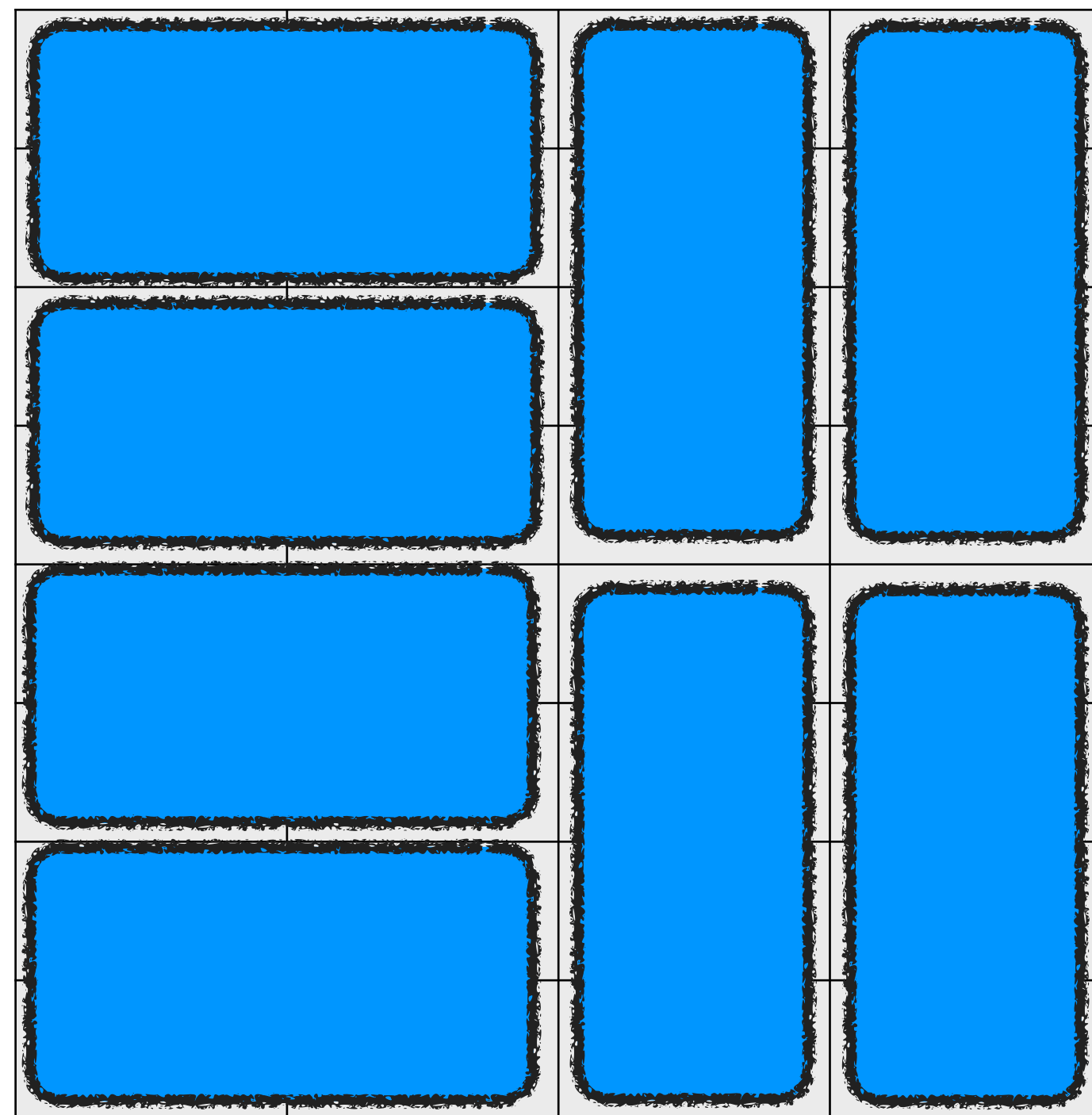
row store

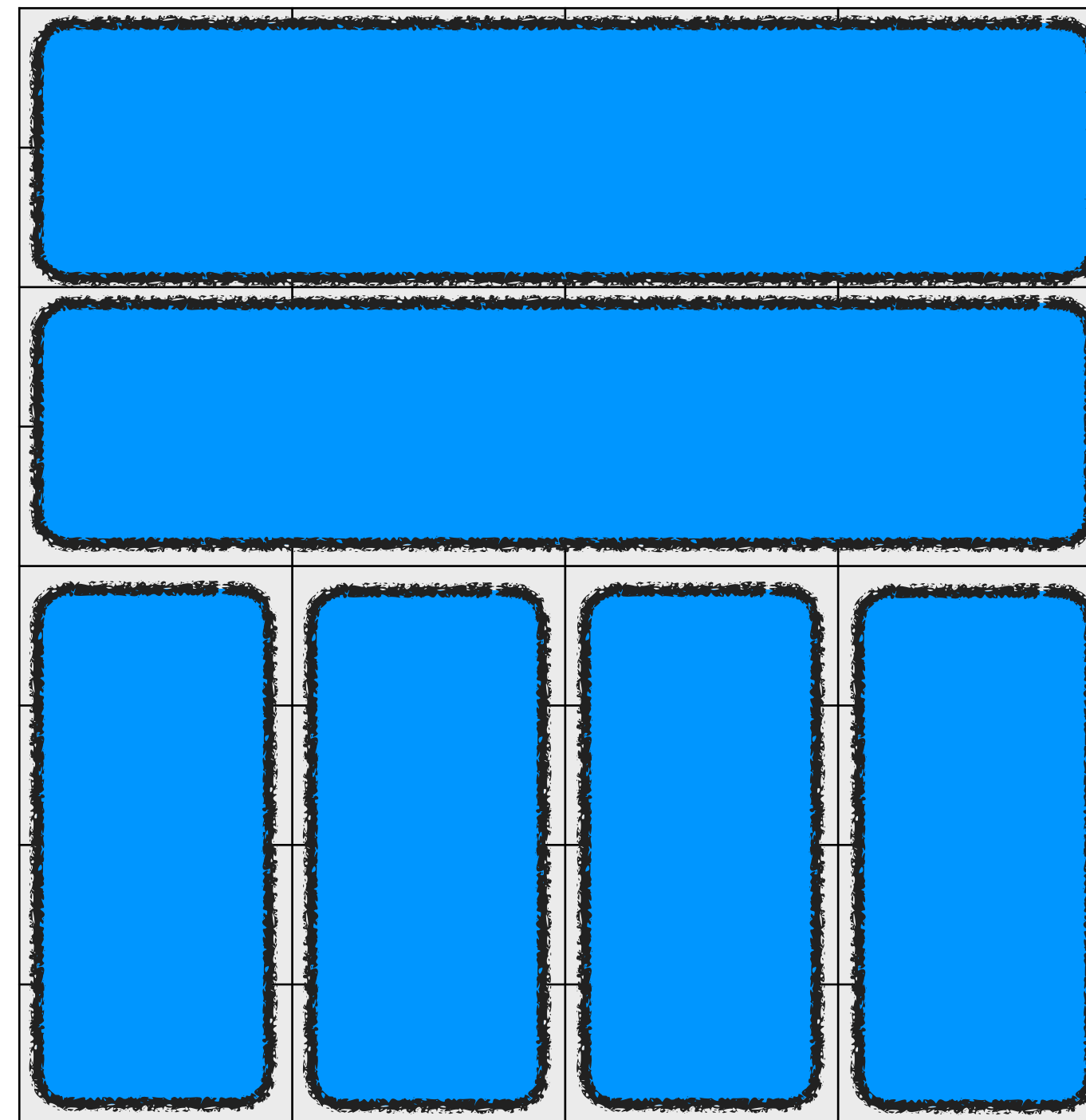column store

# Querying over slotted pages

Understanding the schema

**R** (A, B, C, D)

schema: **R** (A, B, C, D)

schema: **R** (A, B, C, D)

schema: **R** (A,

store

hybrid data layouts

Brandeis
UNIVERSITY
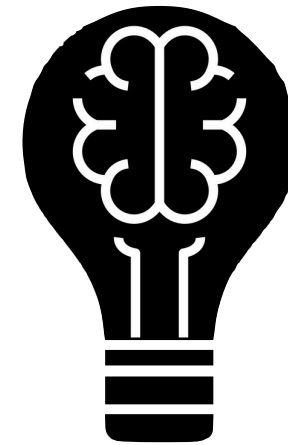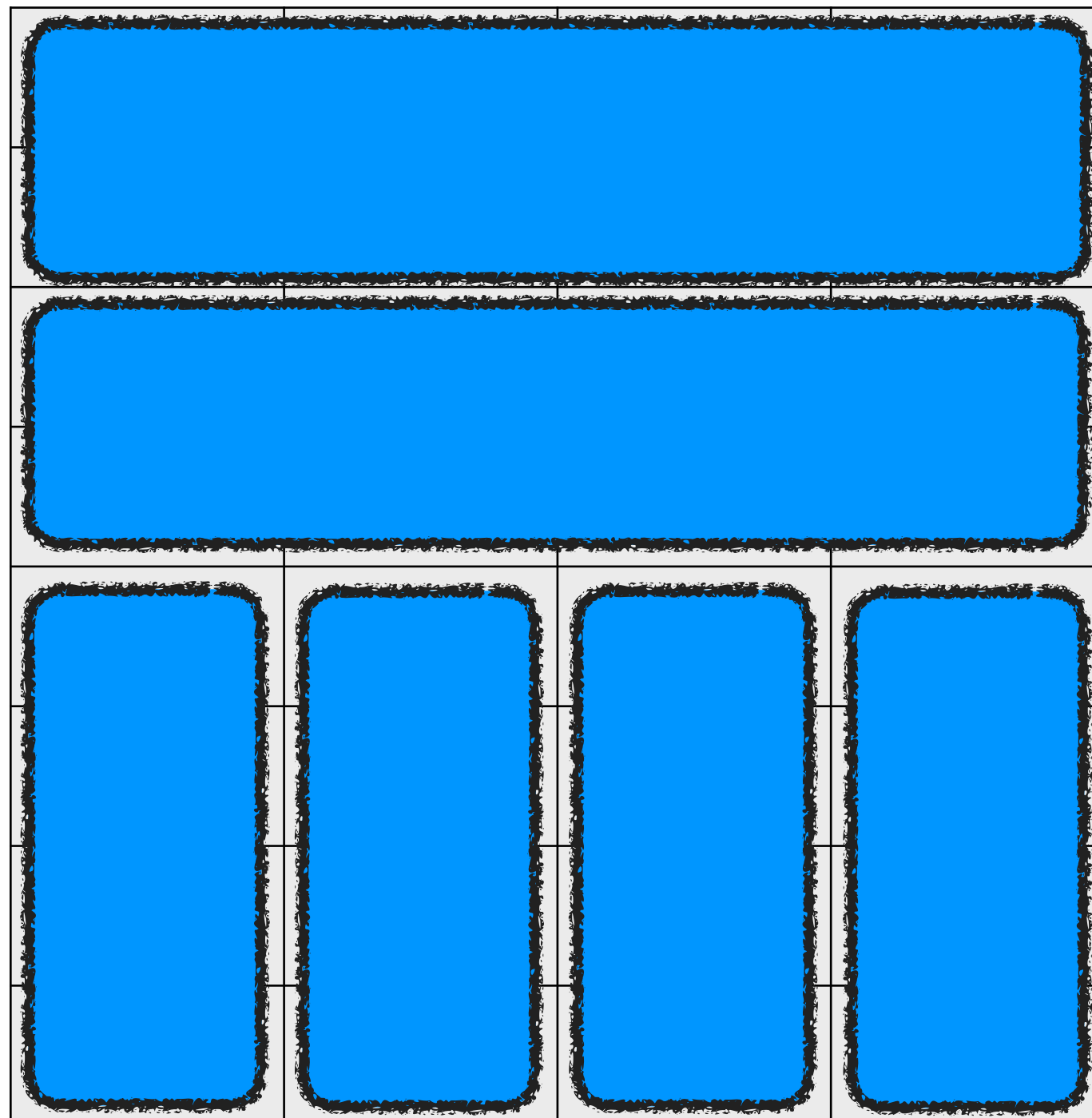
# **Querying** over **slotted pages**

Understanding the schema

schema: **R** (A, B, C, D)



Thought Experiment 6

When would you have this data layout?

Queries on a subset of data

Brandeis
UNIVERSITY

# **Evolution** of column store

From row sotres to column stores

1985: first complete
column-store model

2000: first complete
column-store system

2001: first idea for
hybrid layouts

2012: expanding on
hybrid layouts

60s     70s     80s     90s     00s     10s     20s

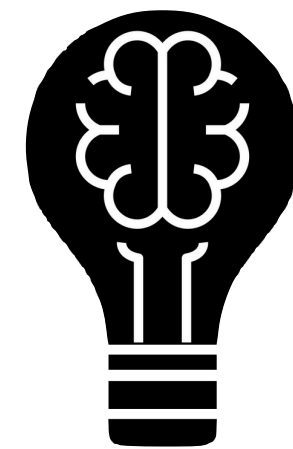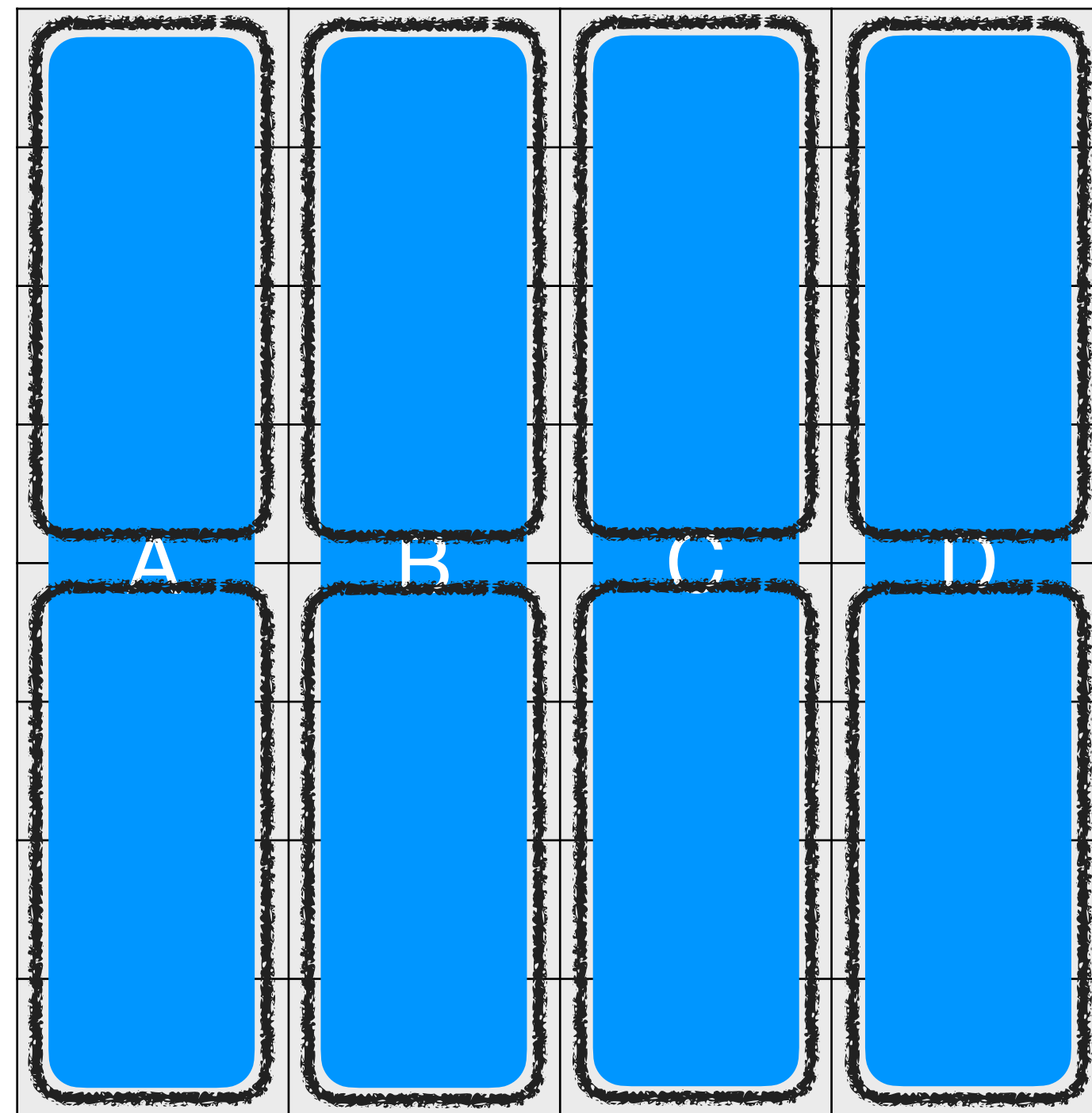rows    rows    rows    rows    rows    rows*

# Querying over slotted pages

Understanding the schema

schema: **R** (A, B, C, D)

file

A    B    C    D

○○○

Thought Experiment

select max(B) from R
where A>5 and C<10

Home work!

Brandeis
UNIVERSITY

# Next time in COSI 167A

Row stores vs. Column stores

**[P]** "Column-Stores vs. Row-Stores: How Different Are They Really?", *SIGMOD*, 2008

**TECHNICAL QUESTION 1**

[B] "C-Store: A Column-oriented DBMS", *VLDB*, 2005

Brandeis
UNIVERSITY

*COSI 167A*
Advanced Data Systems

Class 3

**Data Systems Architecture**

Prof. Subhadeep Sarkar

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