

## class 7

# Design Tradeoffs in Key-Value Stores

Prof. Manos Athanassoulis

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

Do we have a quiz today ... ?

**No!**

# what to do now?

- A) read the syllabus and the website
- B) register to Piazza + Gradescope
- C) finish project 0
- D) **finish project 1 (due 2/14)**
- E) **register for the student-presentations (by 2/17)**
- F) **start working on project proposal (due 2/23)**

# Fast Scans on Key-Value Stores (KVS)

Key-Value Stores are designed for *transactional* workloads (put and get operations)

APACHE  
HBASE



**Analytical** workloads require efficient scans and aggregations  
(typically offered by column-store systems)



**Can we do both in one system?**

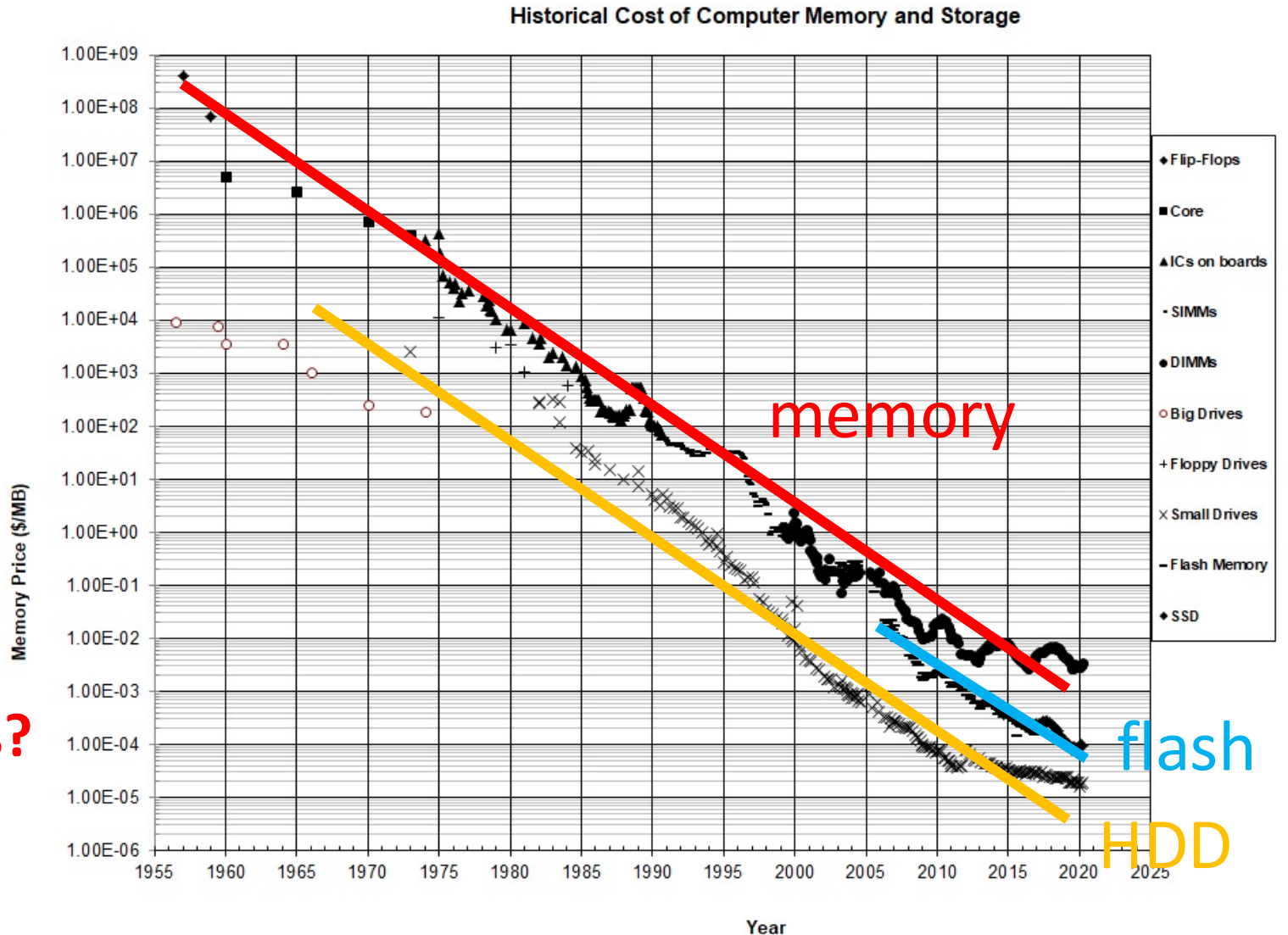
# Why combine KVS and analytical systems?

cheaper and cheaper storage

more data ingestion

need for write-optimized  
data structures

what about analytical queries?



# Both **transactional** and **analytical** systems

Most organizations maintain both

- ***transactional*** systems (often as key-value stores)
- ***analytical*** systems (often as column-stores)

problems?



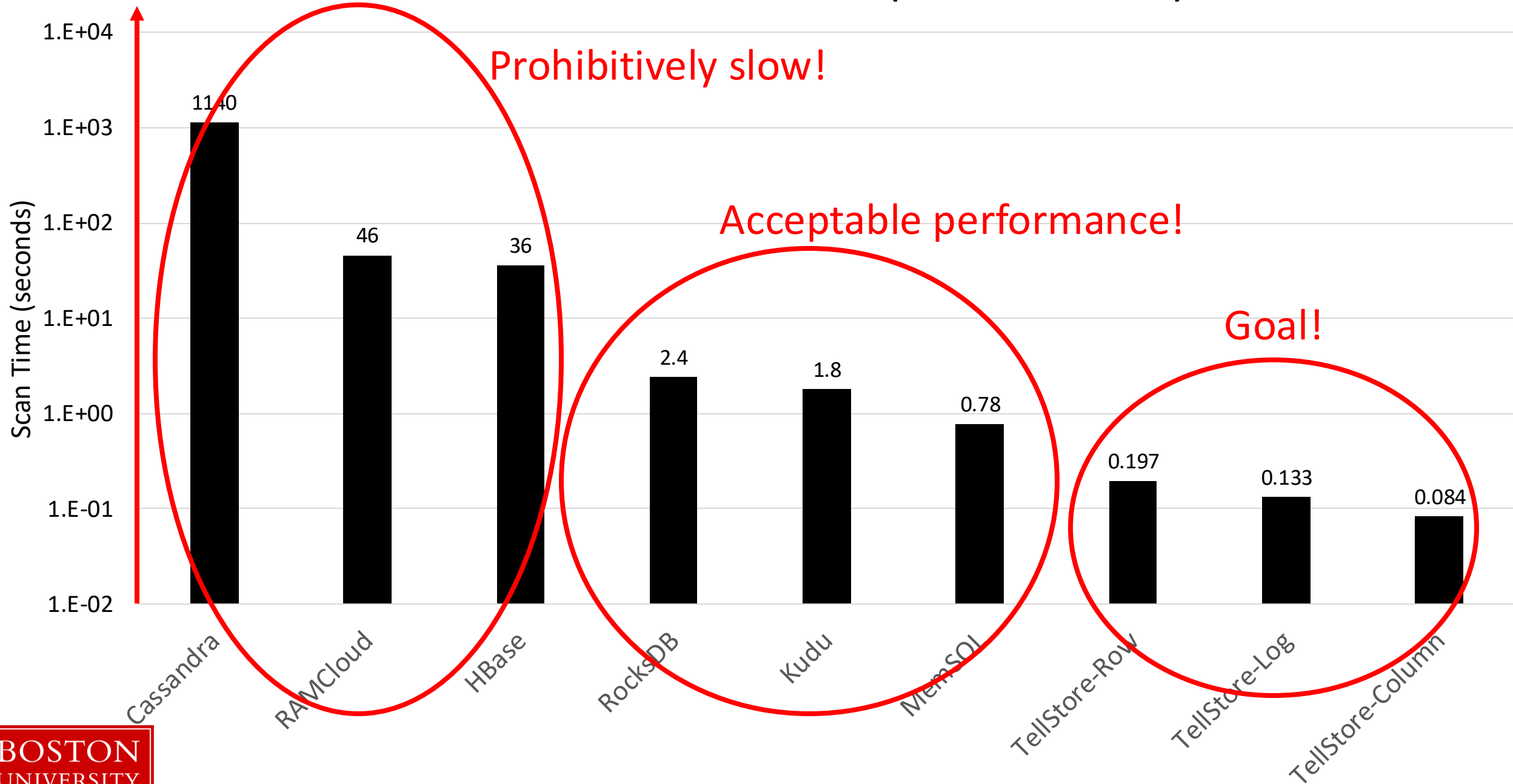
requires additional expertise and management (e.g., two DBAs)

harder to maintain (more systems, more code)

time consuming data integration/transfer

Log scale!

# Scan Time of 50M records (~4GB of data)



# Goals of this paper

Bridge the conflicting goals of *get/put* and *scan* operations

**get/put** operations need **sparse data structures** → *locality is not required, access one object*  
**scans** require **locality** (relevant data to be packed together)

we will discuss how to compromise, via the design of *Tellstore*

how to amend the *SQL-over-NoSQL* architecture for mixed workloads



# SQL over NoSQL

Elasticity

Snapshot Isolation

Processing  
Layer

e.g., MVCC, no locking, timestamp based

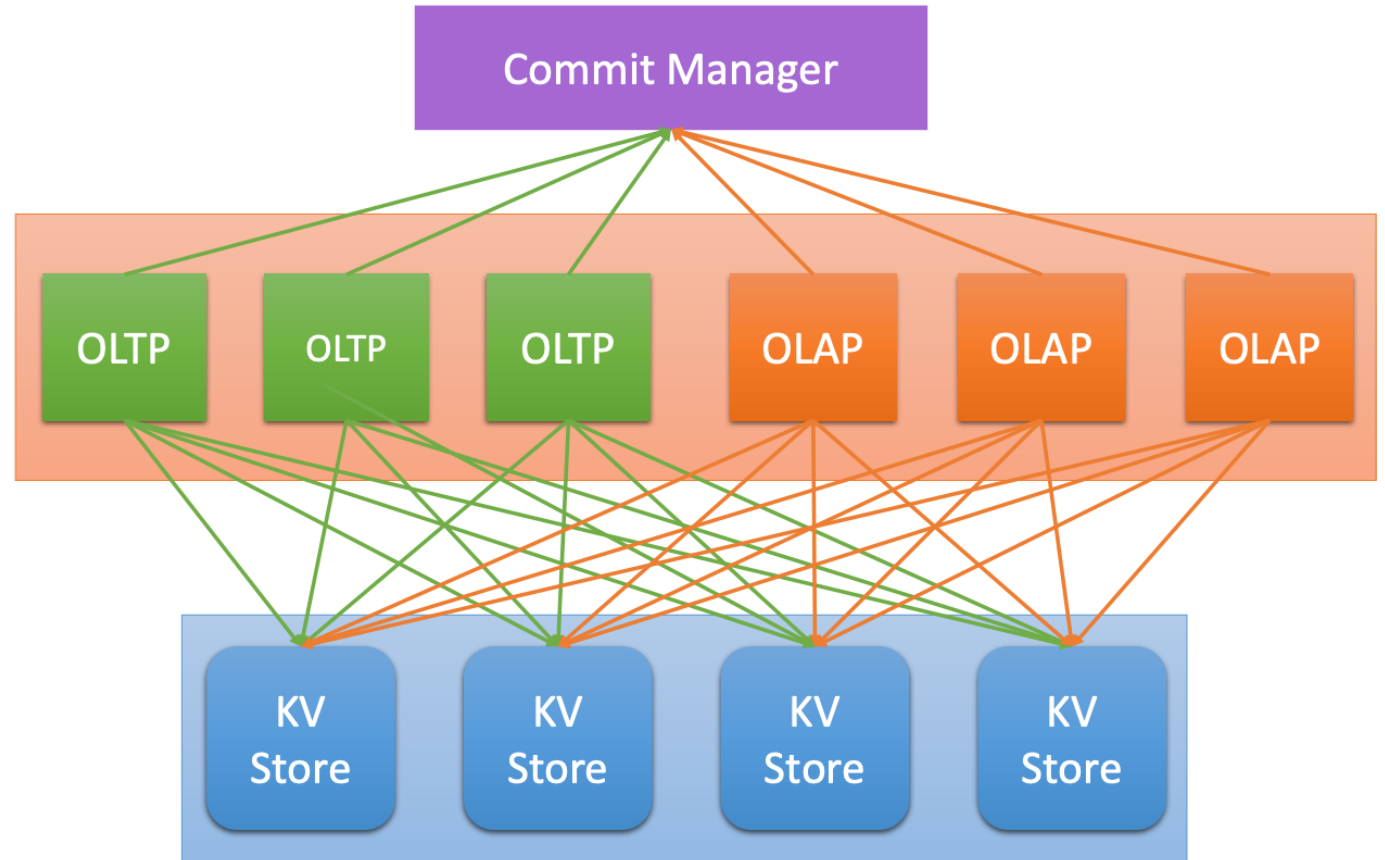
Support for:

Scans

Versioning

Batching

Storage  
Layer



# Scans

selection

projection

(simple) aggregates

shared scans

remember them?

# Versioning

multiple versions  
through timestamps

garbage collection

discarding old versions  
during scans might be costly

# Batching

batch several  
requests to the  
storage layer

amortize the  
network time



# Challenges

scans vs. get/put

Scans need columnar locality

#1, John, 2/4/88, Boston

get/put need row-wise locality

2/4/88  
2/1/87  
7/7/93  
4/1/92  
3/9/91  
9/3/96

why?



# Challenges

scans vs. get/put

scans vs. versioning

#1, John, 2/4/88, Boston, v1

#1, John, 2/4/88, Cambridge, v2



versioning reduces locality in scans

checking for the latest version in scans needs CPU time

# Challenges

scans vs. get/put

scans vs. versioning

scans vs. batching

batching **multiple scans** or **multiple put/get** requests is ok

but ...

**batching scans and puts/gets is a bad idea!**

why?



**puts/gets** need fast predictable performance

**scans** inherently have high and variable latency

# How to design KVS for efficient scans?

## **Key design decisions**

(A) Updates

(B) Layout

(C) Versioning

# How to design KVS for efficient scans?

## Key design decisions

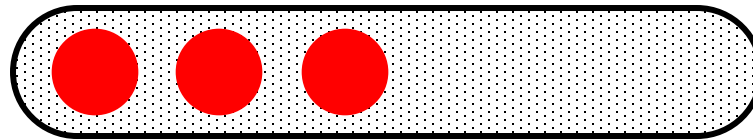
(A) Updates *in-place*



# How to design KVS for efficient scans?

## Key design decisions

(A) Updates      *in-place*      *log-structured*

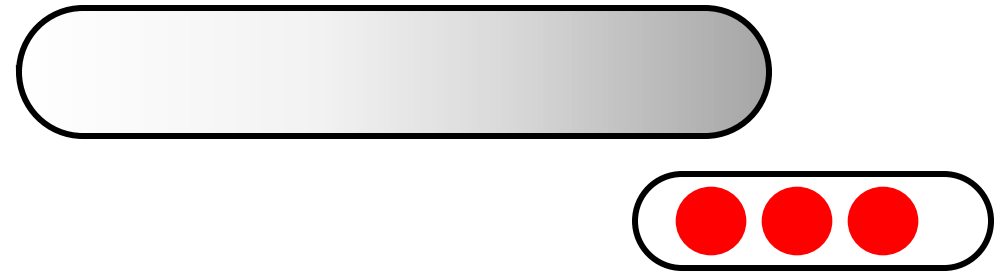




# How to design KVS for efficient scans?

## Key design decisions

(A) Updates      *in-place*      *log-structured*      *delta-main*

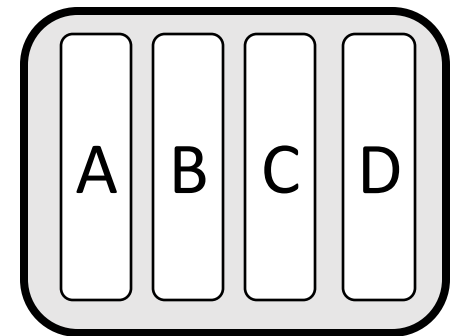
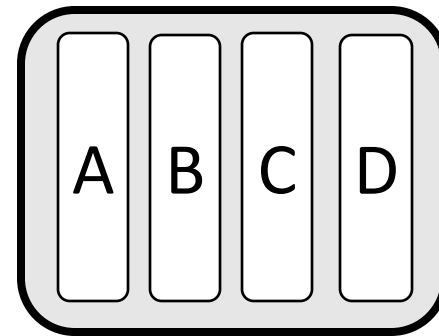
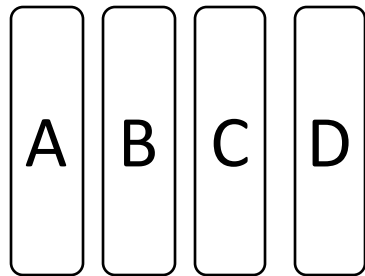


# How to design KVS for efficient scans?

## Key design decisions

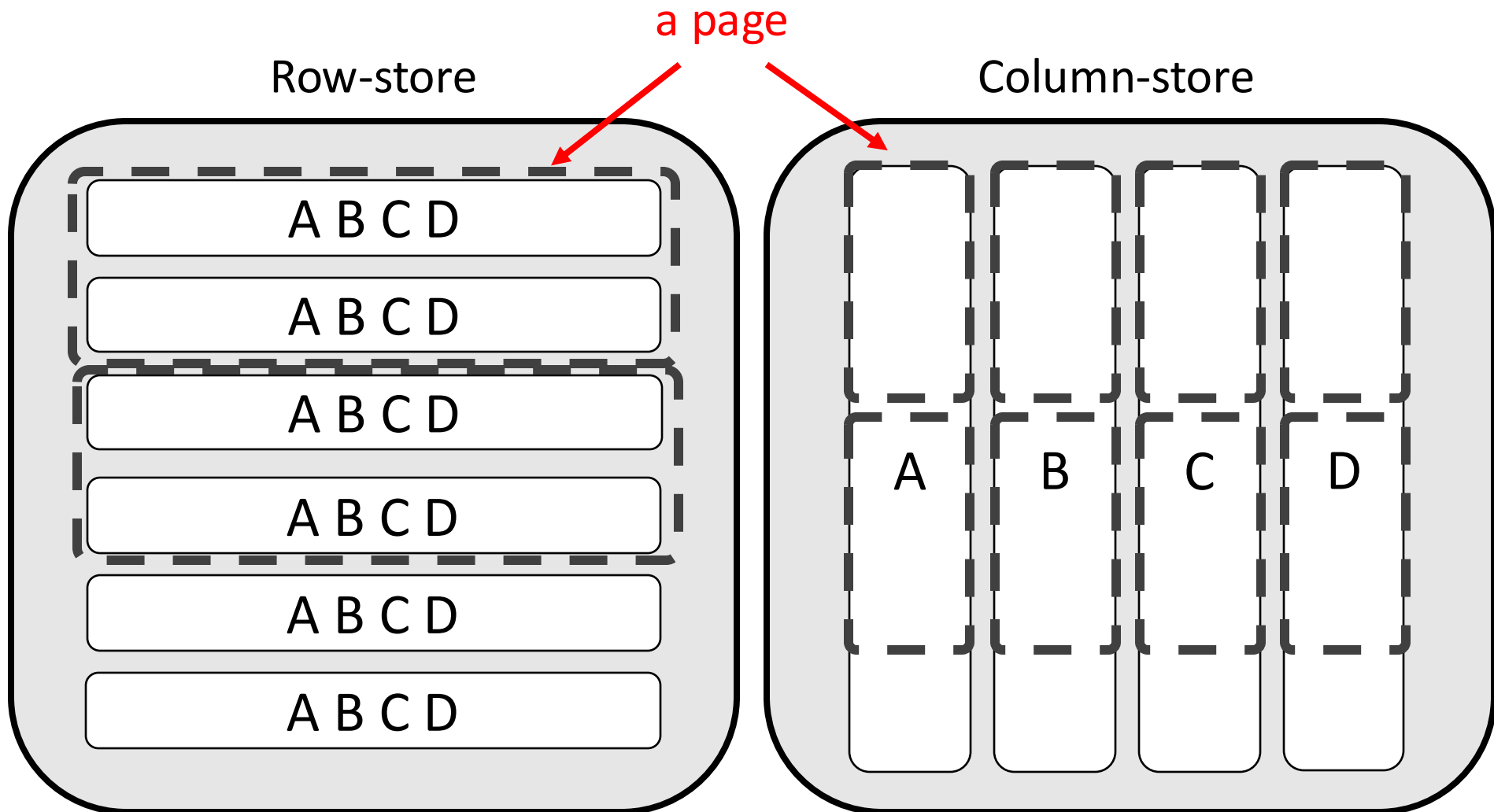
(A) Updates      *in-place*      *log-structured*      *delta-main*

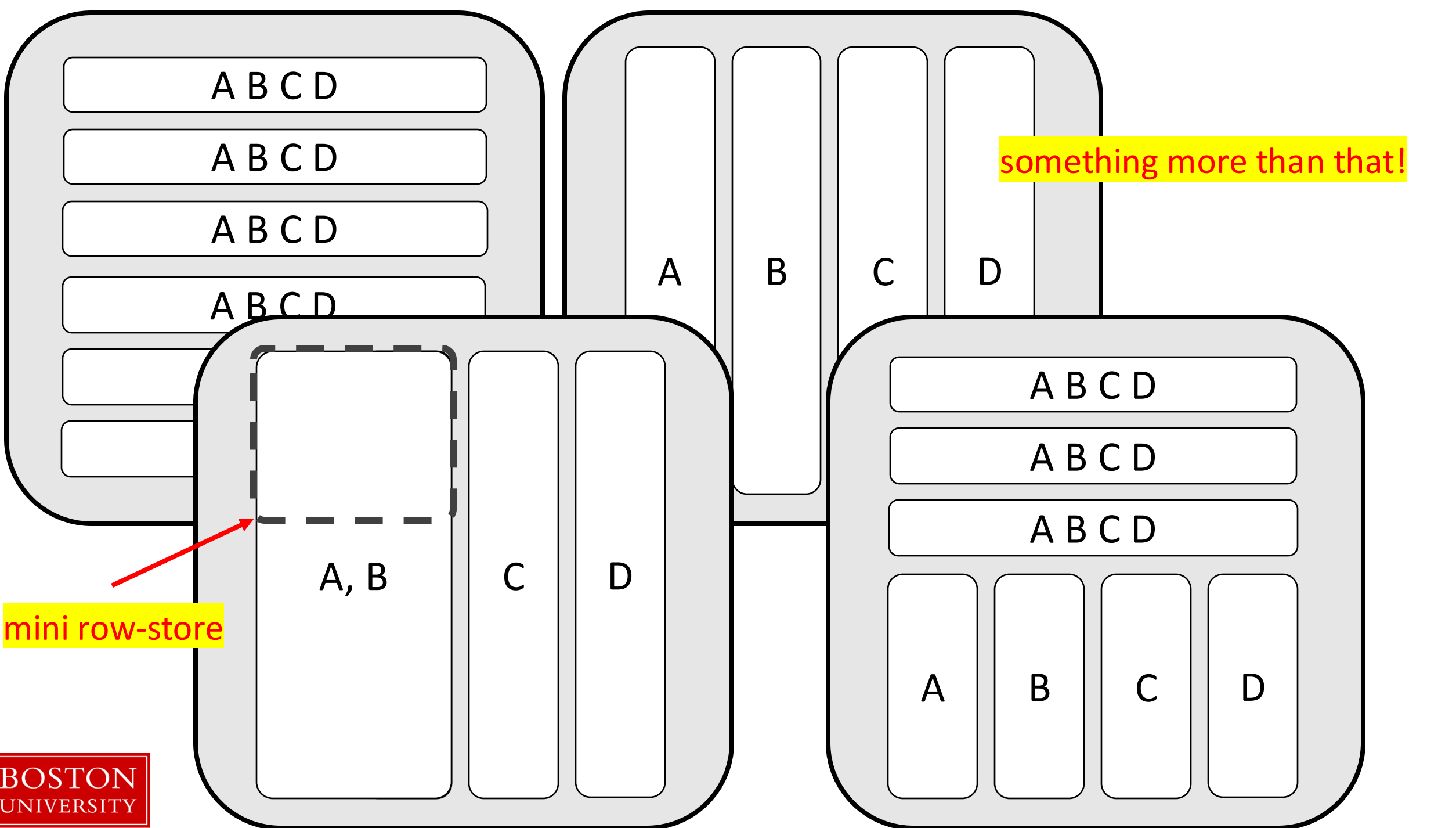
(B) Layout      *column*      *PAX (columnar per page)*



# Small detour: page layouts

middle ground?





# Partition Attributes Across (PAX)

Middle ground?



Decompose a slotted-page internally  
in mini-pages per attribute

✓ Cache-friendly

➤ Brings only relevant attributes to cache



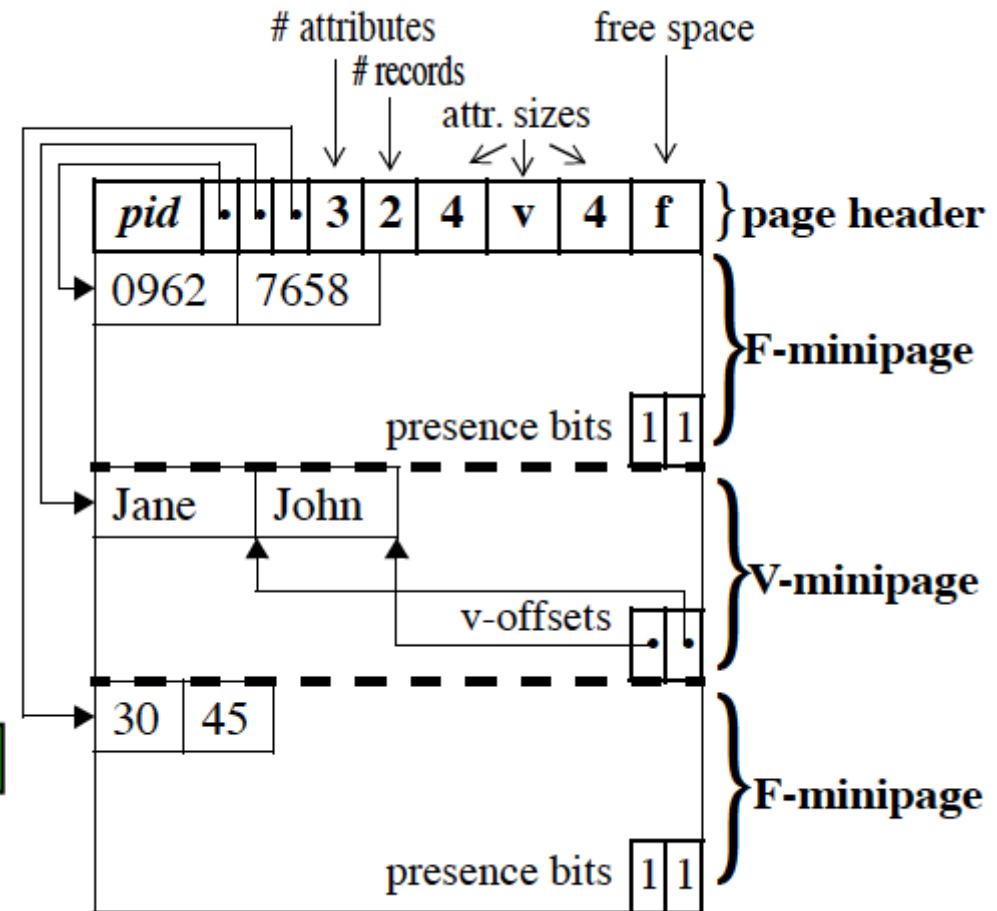
• Compatible with slotted-pages?



• Same update abstraction?



– (insert in a page)

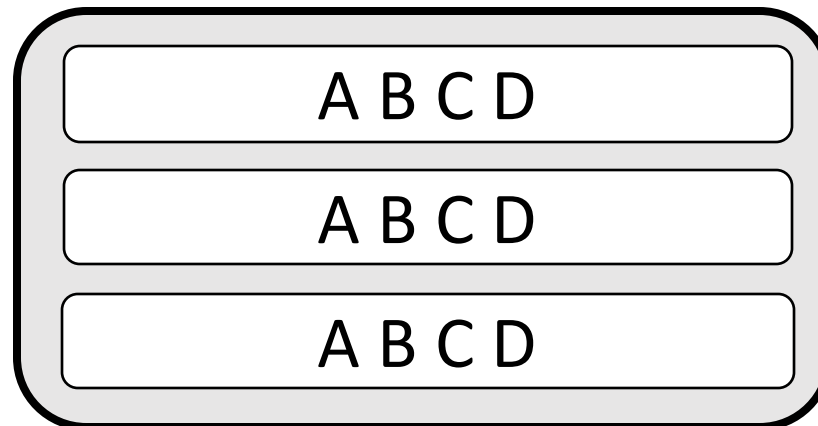


# How to design KVS for efficient scans?

## Key design decisions

(A) Updates     *in-place*     *log-structured*     *delta-main*

(B) Layout     *column (PAX)*     *row*



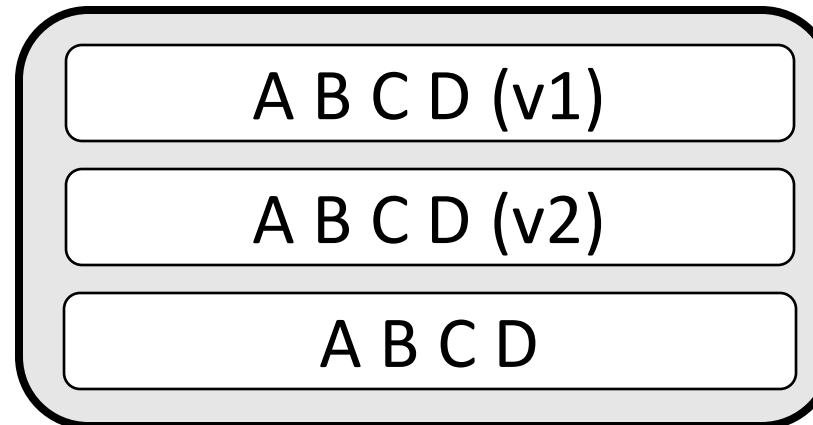
# How to design KVS for efficient scans?

## Key design decisions

(A) Updates     *in-place*     *log-structured*     *delta-main*

(B) Layout     *column (PAX)*     *row*

(C) Versioning     *clustered*



any other options?

# How to design KVS for efficient scans?

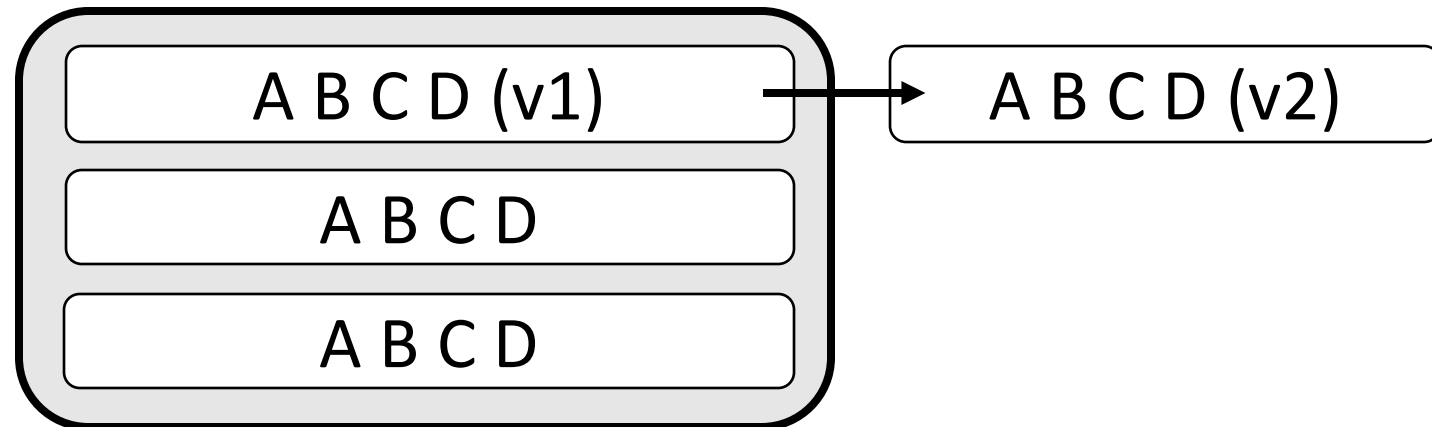
## Key design decisions

(A) Updates     *in-place*     *log-structured*     *delta-main*

(B) Layout     *column (PAX)*     *row*

(C) Versioning     *clustered*

*chained*





# How to design KVS for efficient scans?

## Key design decisions

(A) Updates     *in-place*     *log-structured*     *delta-main*

(B) Layout     *column (PAX)*     *row*

(C) Versioning     *clustered*     *chained*

what comes as a result of versioning?



# Garbage Collection (GC)

(A) Periodic      *separate dedicated thread(s)*

(B) Piggy-backed GC during scans

*increases scan time      but frequently read tables benefit*

*avoids re-reading for GC (since data is already accessed)*

# Design Space

Updates



Layout



Versioning



GC

*in-place*

*column (PAX)*

*clustered*

*periodic*

*log-structured*

*row*

*chained*

*piggy-backed*

*delta-main*



hybrid designs are also valid!  
should we consider all possible designs?

# Design Space

Updates

×

Layout

×

Versioning

×

GC

*in-place*

*column (PAX)*

*clustered*

*periodic*

*log-structured*

*row*

*chained*

*piggy-backed*

*delta-main*

some combinations can be discarded:

log-structured & column **worse than** delta-main & column

log-structured & clustered **worse than** log-structured & chained

note that each combination here represents multiple options

<i>Dimension</i>	<i>Approach</i>	<i>Advantages</i>	<i>Disadvantages</i>
<i>Update</i>	<b>update-in-place</b>	storage	versioning, concurrency
	<b>log-structured</b>	storage, concurrency	GC
	<b>delta-main</b>	compromise	
<i>Layout</i>	<b>column (PAX)</b>	scan	get/put
	<b>row</b>	get/put	scan
<i>Versions</i>	<b>clustered</b>	get/put	GC
	<b>chained</b>	GC	scan

**Table 2: Design Tradeoffs**

# Design Space

Updates

×

Layout

×

Versioning

×

GC

*in-place*

*column (PAX)*

*clustered*

*periodic*

*log-structured*

*row*

*chained*

*piggy-backed*

*delta-main*

focus on two extremes:

(1) log-structured & row & chained

***TellStore-Log***

# Design Space

Updates

×

Layout

×

Versioning

×

GC

*in-place*

*column (PAX)*

*clustered*

*periodic*

*log-structured*

*row*

*chained*

*piggy-backed*

*delta-main*

focus on two extremes:

(1) log-structured & row & chained

(2) delta-main & column & clustered

***TellStore-Log***

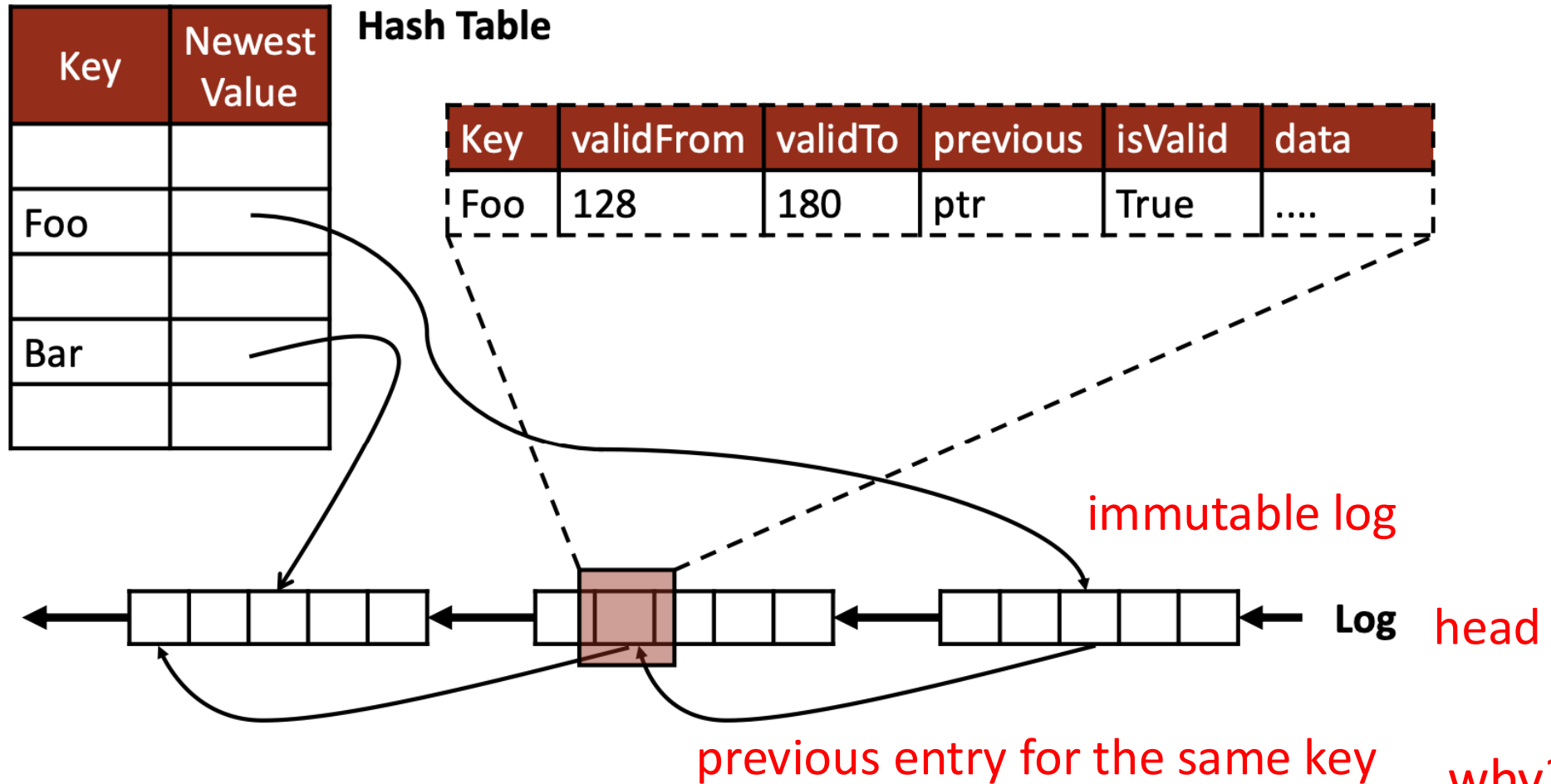
***TellStore-Col***

# TellStore-Log

one log per table (locality for scans)

**inserts, updates, and deletes** are all **logged**

lock-free hash table

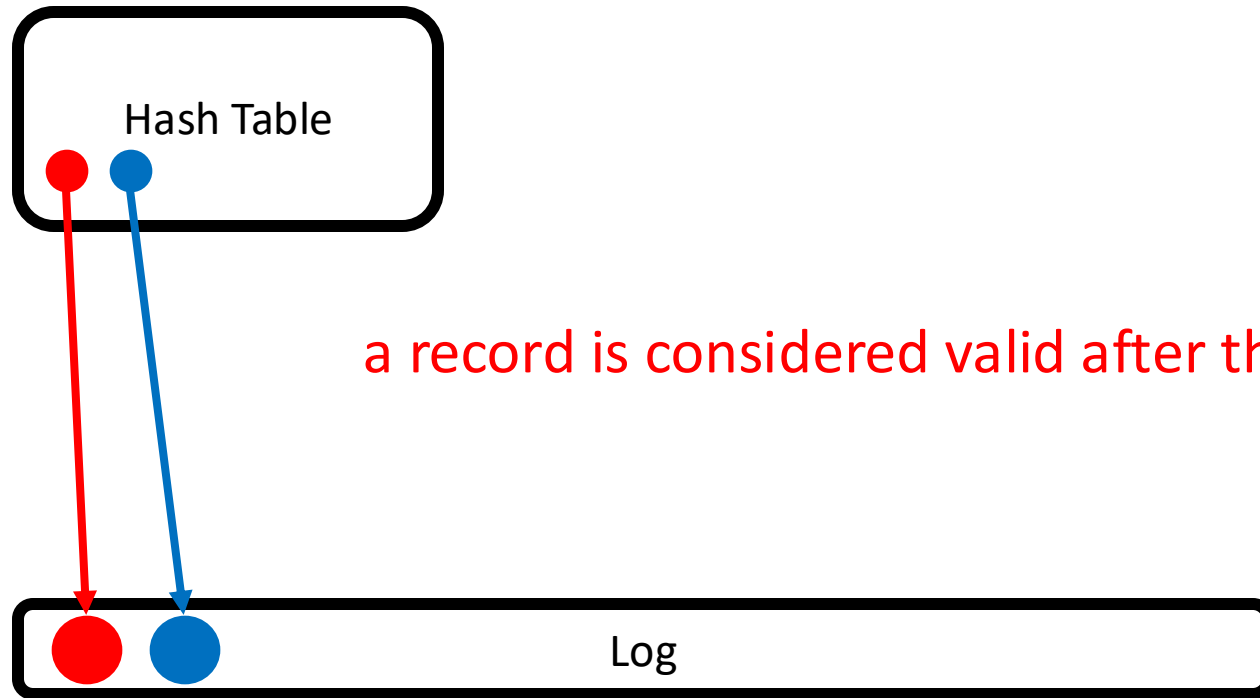


why?





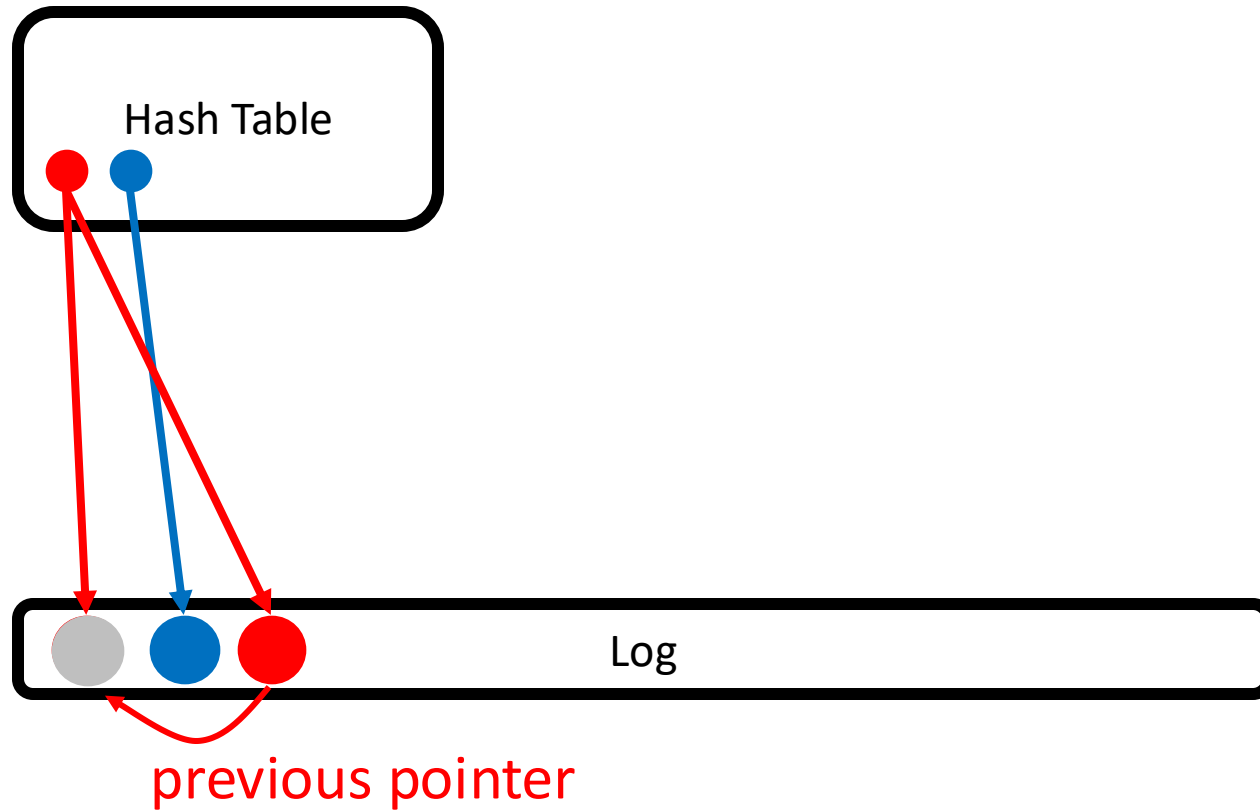
# TellStore-Log Insertion



a record is considered valid after the hash table pointer is updated

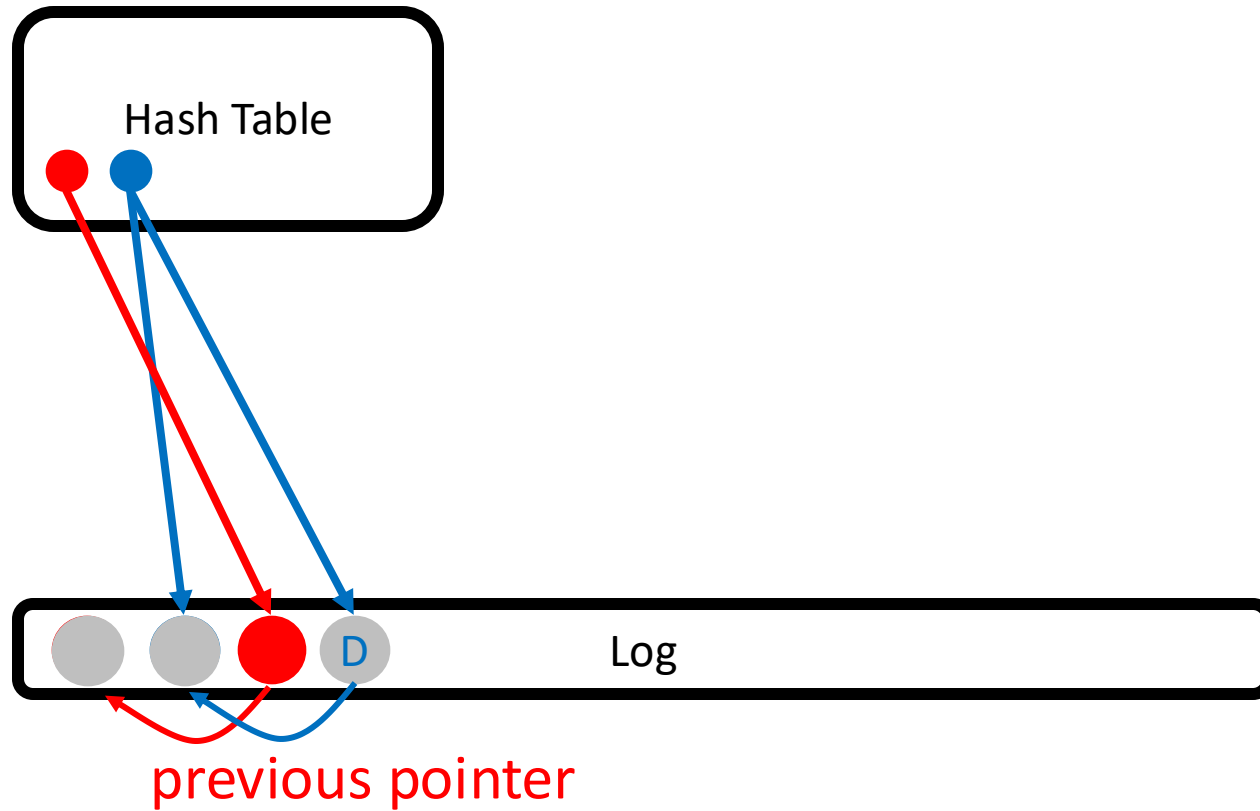
the log contains *rows*

# TellStore-Log Update



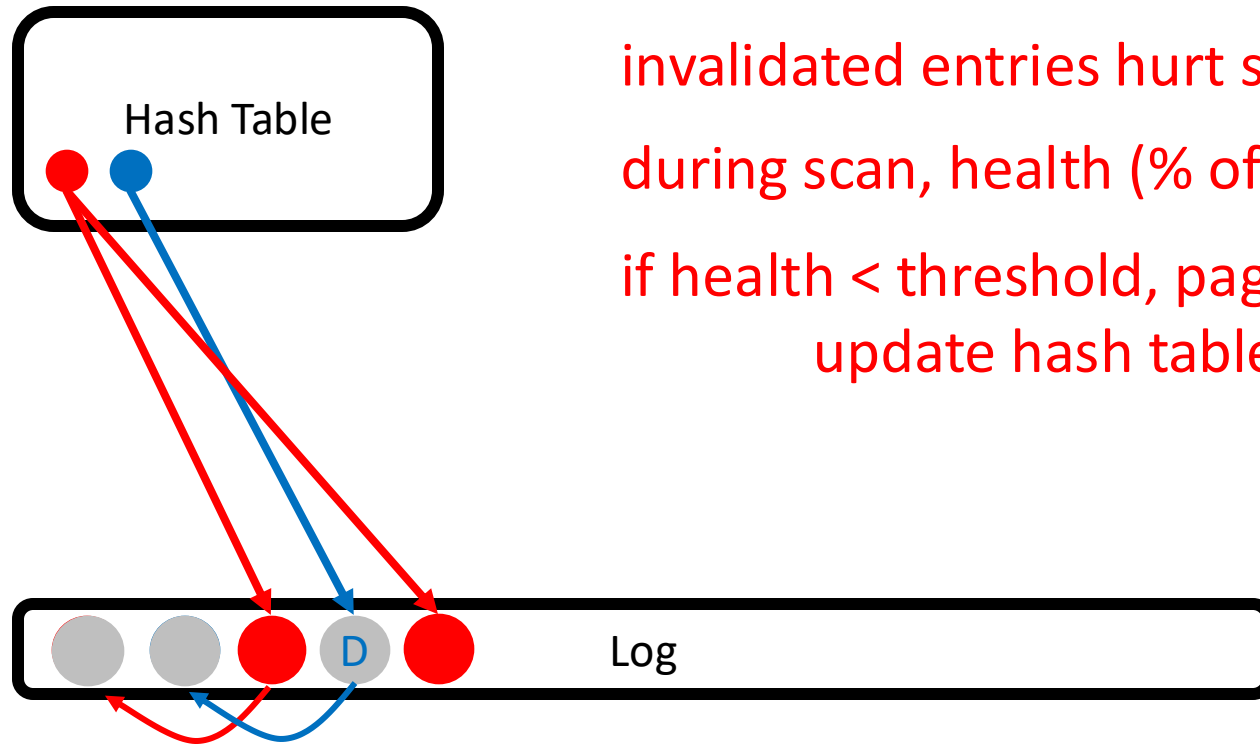
the log contains *rows*

# TellStore-Log Delete



the log contains **rows**

# TellStore-Log Garbage Collection



invalidated entries hurt scan performance

during scan, health (% of invalid entries) per page is calculated

if health < threshold, page is re-written in the head of the log & update hash table & old page is reclaimed

the log contains *rows*

# TellStore-Log in a nutshell

***log-structure:*** efficient puts

***hash-table:*** efficient gets (always points to the latest entry)

***snapshot Isolation:*** high throughput, no locks needed

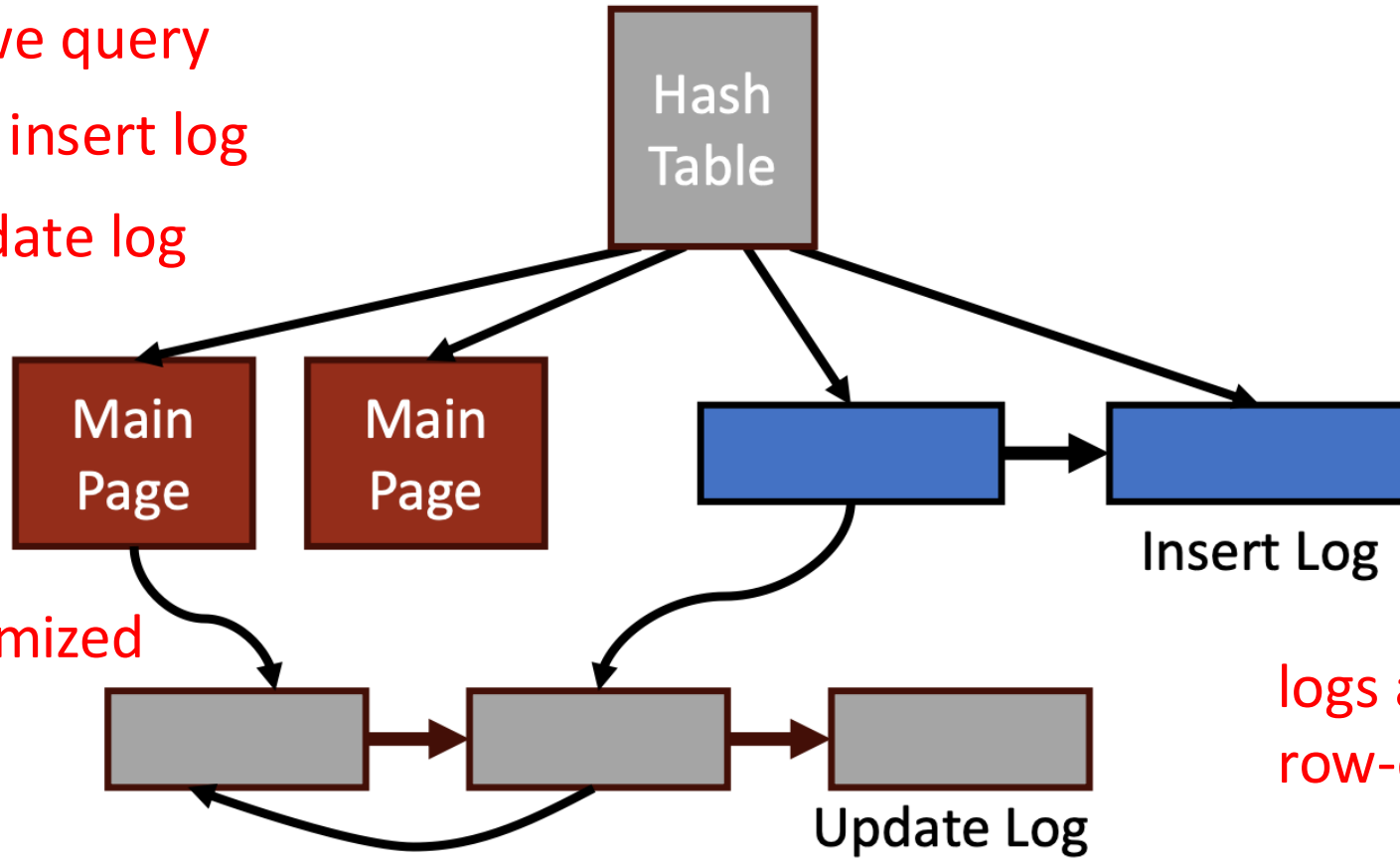
***self-contained log:*** efficient scans (valid from/to needed)

***lazy GC:*** Optimize tables that are scanned

four data structures

# TellStore-Col

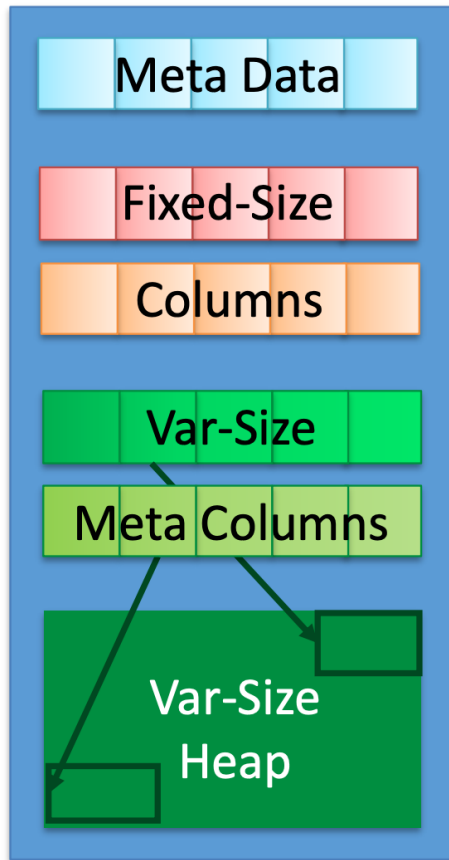
before inserting we query  
new entries go to insert log  
updates go to update log



main is read-only  
columnar: read-optimized

logs are write-optimized  
row-oriented: append-only

# TellStore-Col Layout



fixed-size data is stored in columnar format

variable-size data is indexed in columnar format  
but stored in row-wise format

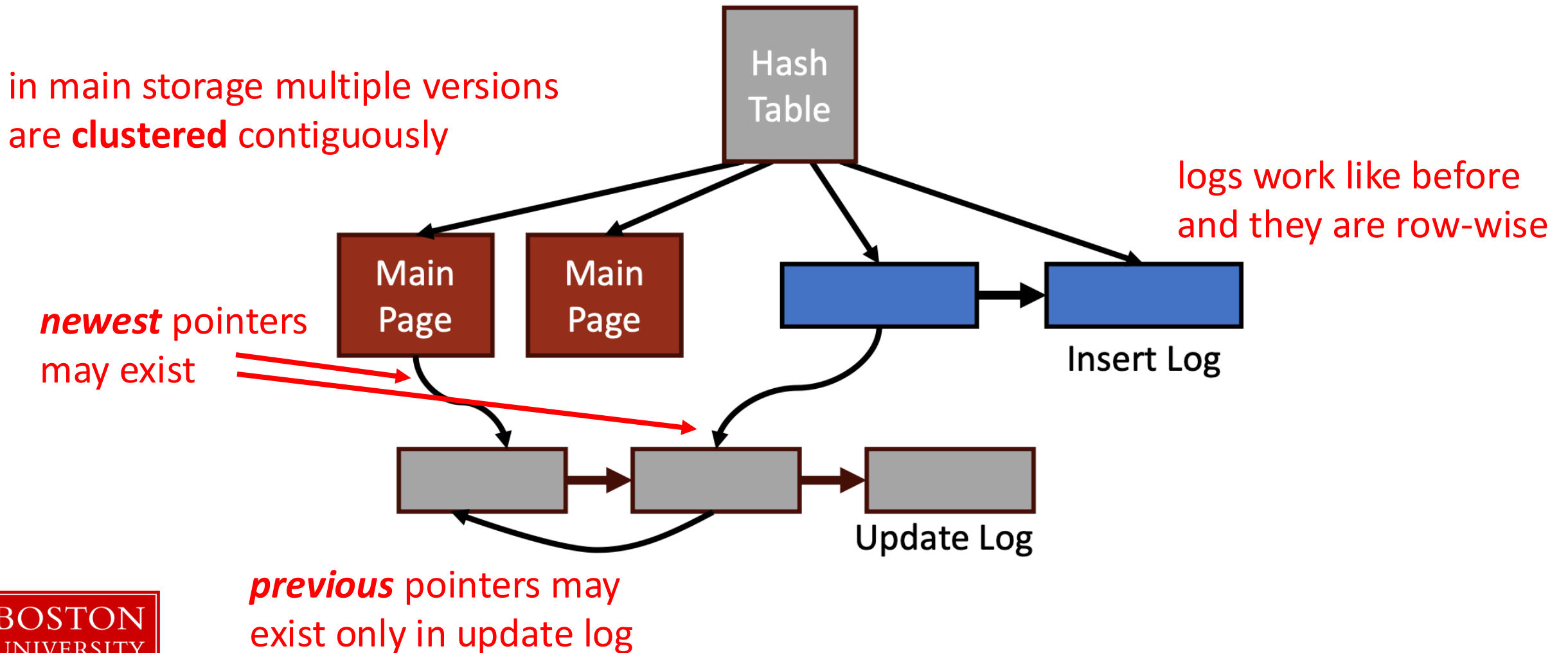
why row-wise?



(1) faster materialization (contiguous copying)

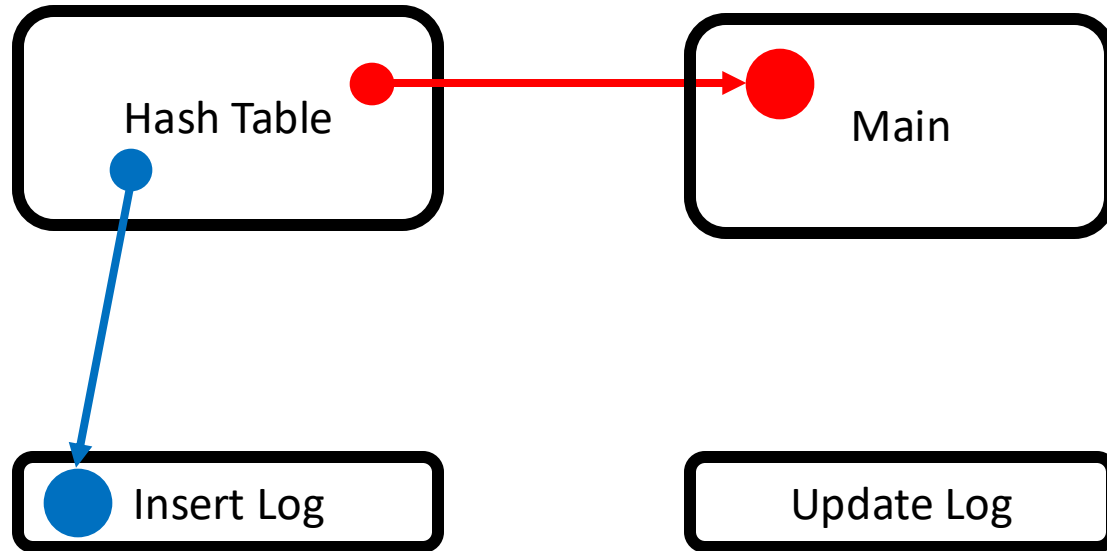
(2) less metadata (one offset for many columns)

# TellStore-Col Versioning

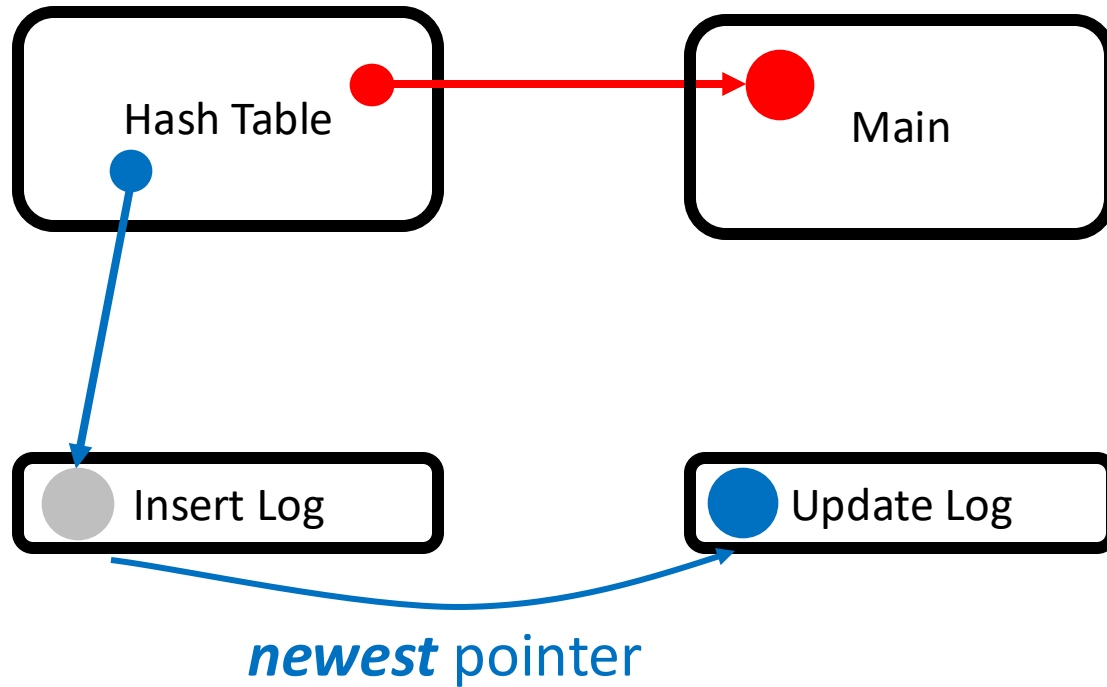




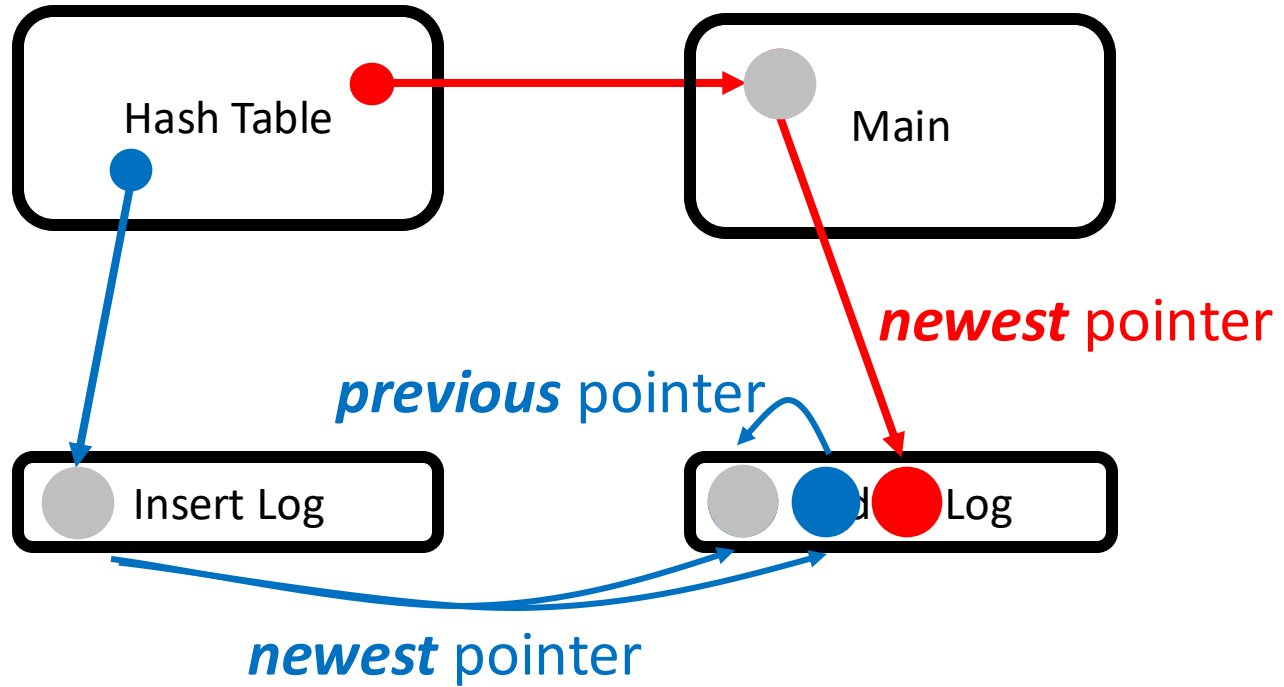
# TellStore-Col Insertion



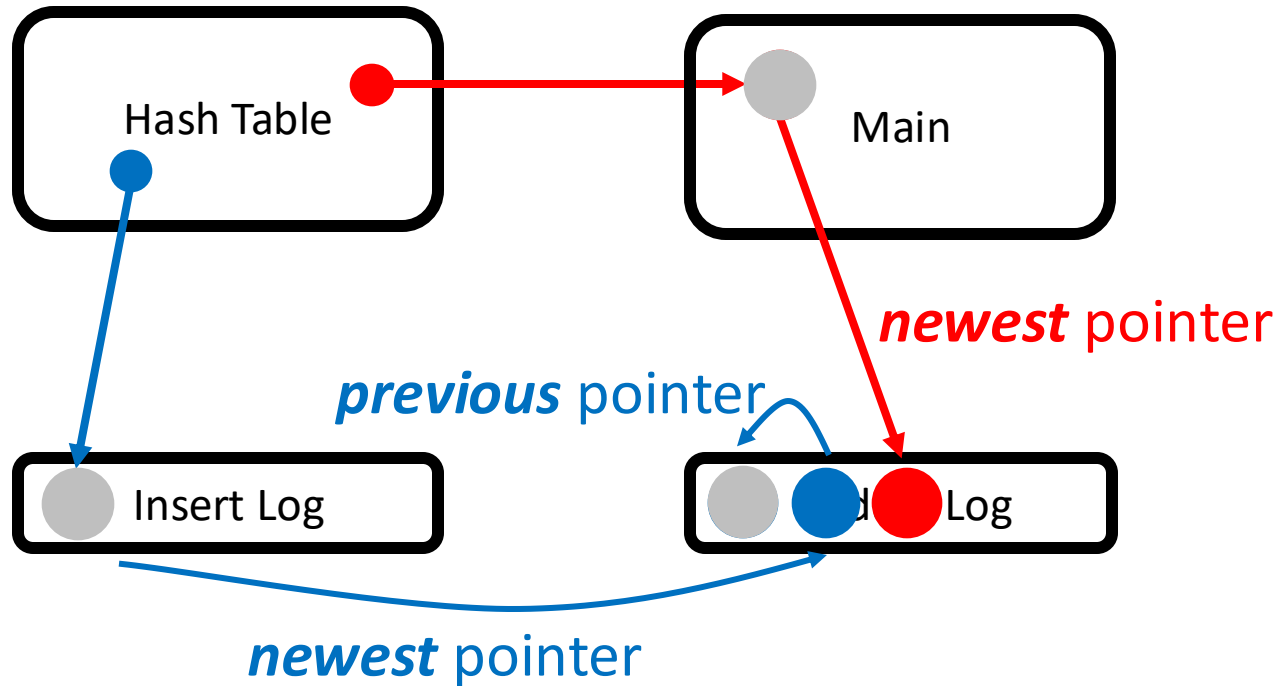
# TellStore-Col Update



# TellStore-Col Update



# TellStore-Col Garbage Collection



dedicated thread  
*(conversion from row to column)*

all main pages with invalid entries

all pages from insert log + update  
to main

run GC frequently + truncate logs

# TellStore-Col in a nutshell

***delta-main:*** compromise between puts and scans

***hash-table:*** efficient gets (always points to the latest entry, may need one more pointer to follow)

***PAX layout:*** minimize disk I/O, maintain locality for scans

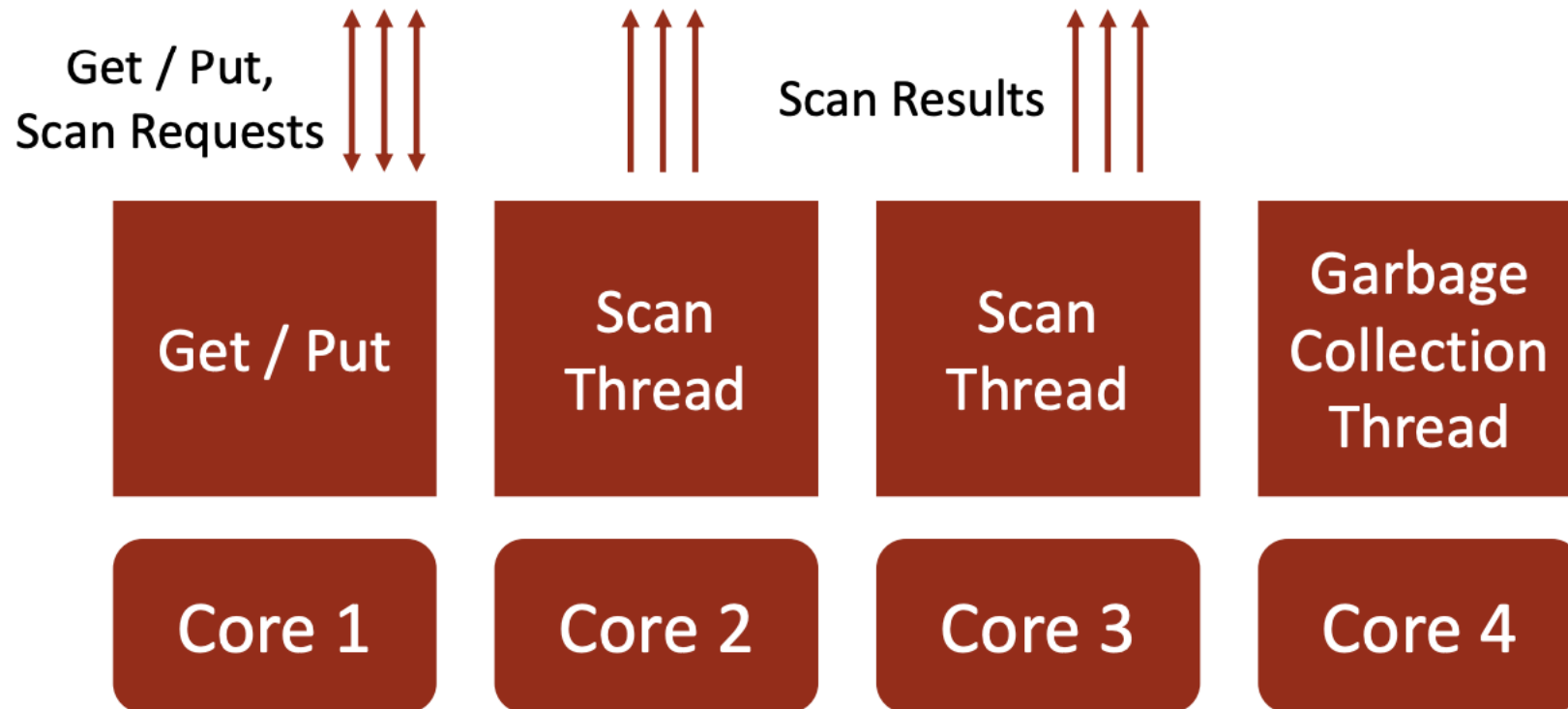
***separate insert/update logs:*** efficient GC

***eager GC:*** improve scans

# Implementation Details

scans are assigned to dedicated threads

scan coordinator for shared scans

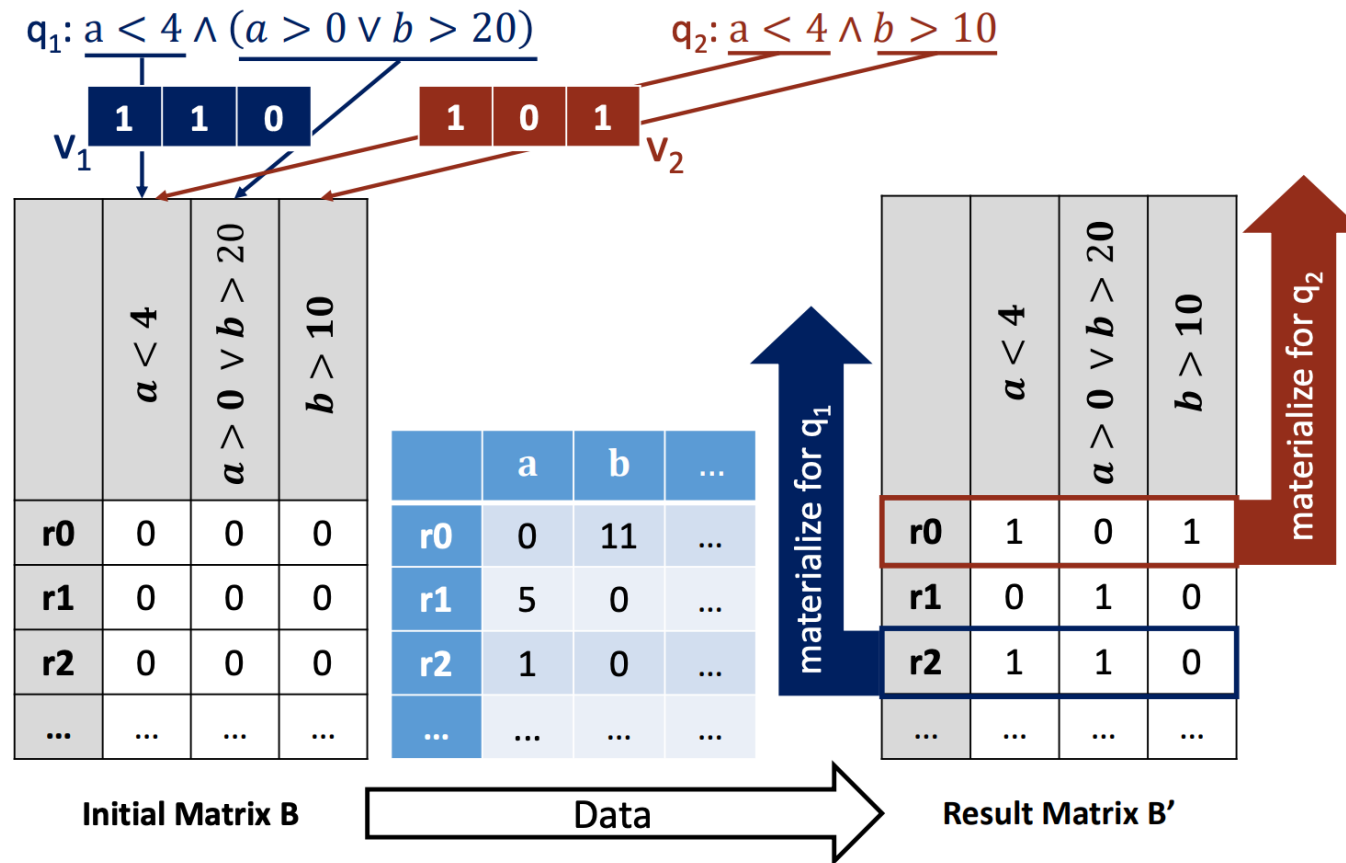


# Implementation Details

efficient predicate evaluation via code generation and predicate pushdown

all queries in CNF

reuse work



# Yahoo! Cloud Serving Benchmark# (YCSB#)

*based on YSCB, a put/get benchmark*

**main\_table** (P, A, B, C, D, E, F, G, H, I, J)     **P: 8-byte key** | A-H: 2-bytes, 4-bytes, 8-bytes | I-J: strings 12-16 bytes

- *Query 1:* A simple aggregation on the first floating point column to calculate the maximum value:

```
SELECT max(B) FROM main_table
```

- *Query 2:* The same aggregation as Query 1, but with an additional selection on a second floating point column and selectivity of about 50%:

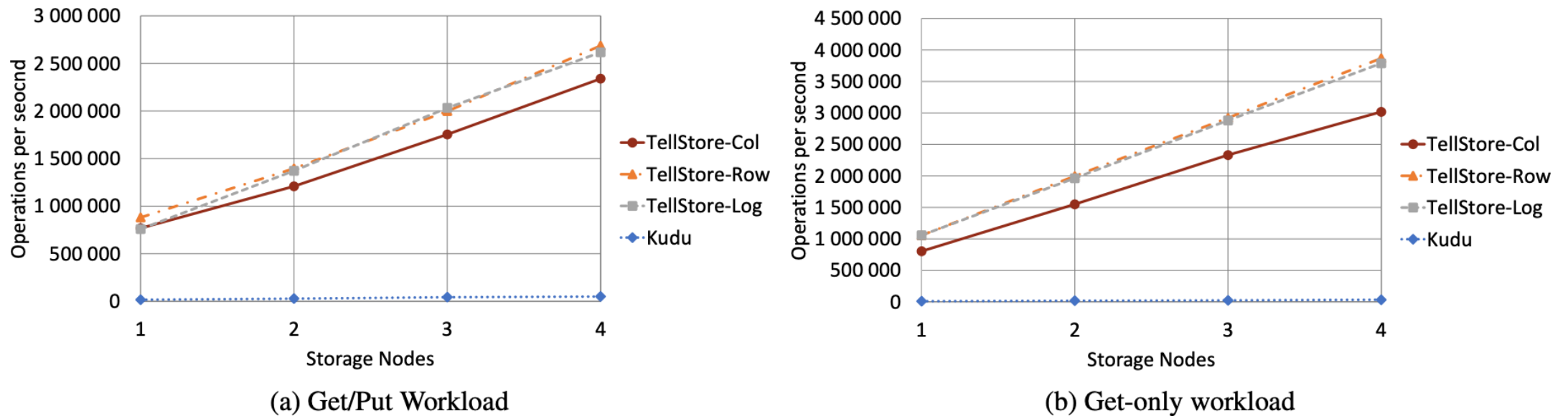
```
SELECT max(B) FROM main_table  
WHERE H > 0 and H < 0.5
```

- *Query 3:* A selection with approximately 10% selectivity:

```
SELECT * FROM main_table  
WHERE F > 0 and F < 26
```



# Experiments: Transactional Workload



**Figure 8: Exp 1, Throughput: YCSB, TellStore Variants and Kudu, Vary Storage Nodes**

Kudu is used as it was the most competitive to begin with

All TellStore approaches are not that far!

# Experiments: Scans

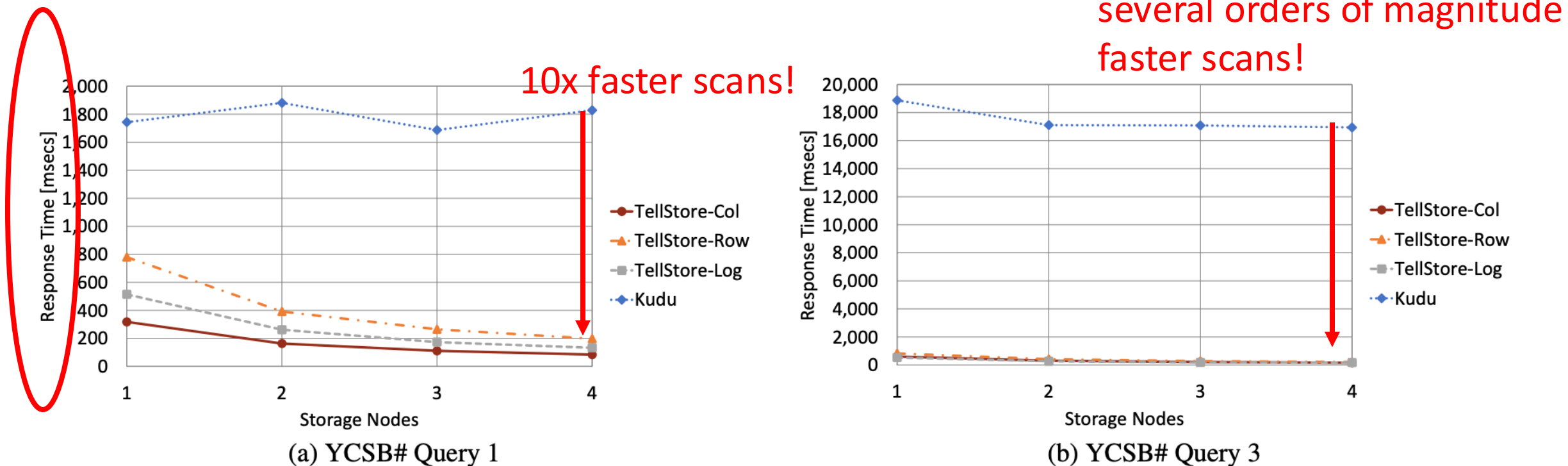
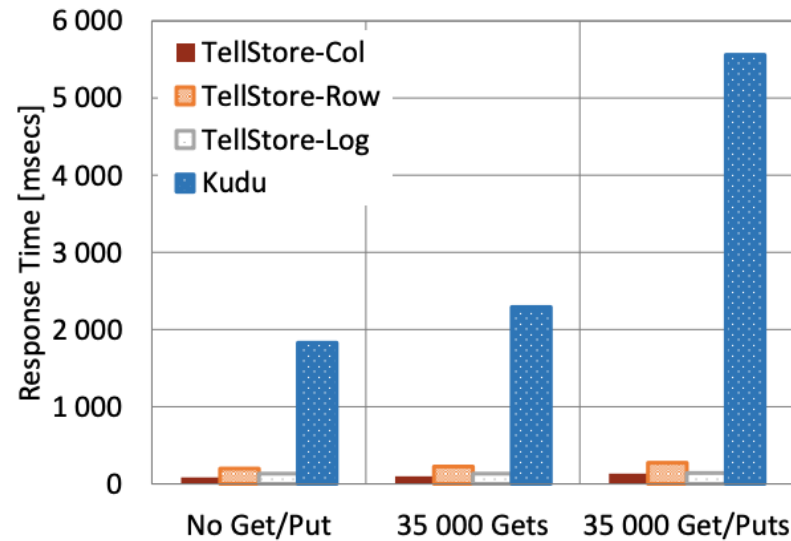


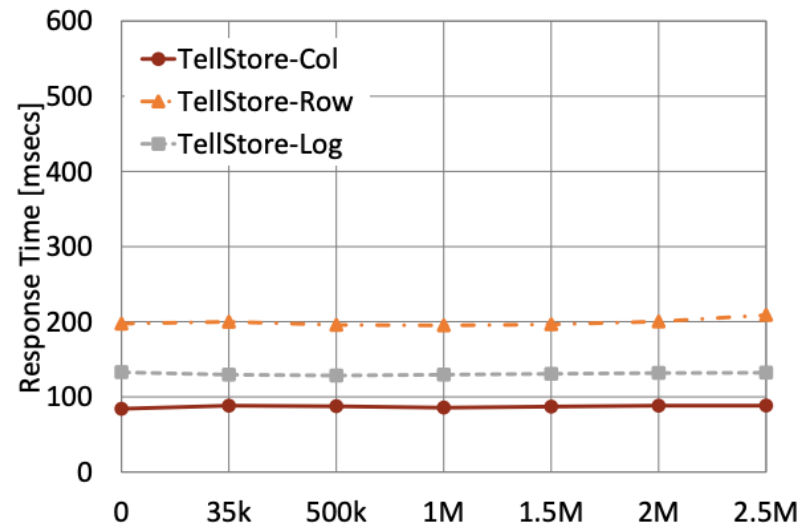
Figure 10: Exp 3, Response Time: YCSB#, Vary Storage Nodes

Q3 does not have projections,  
so no benefit from columnar

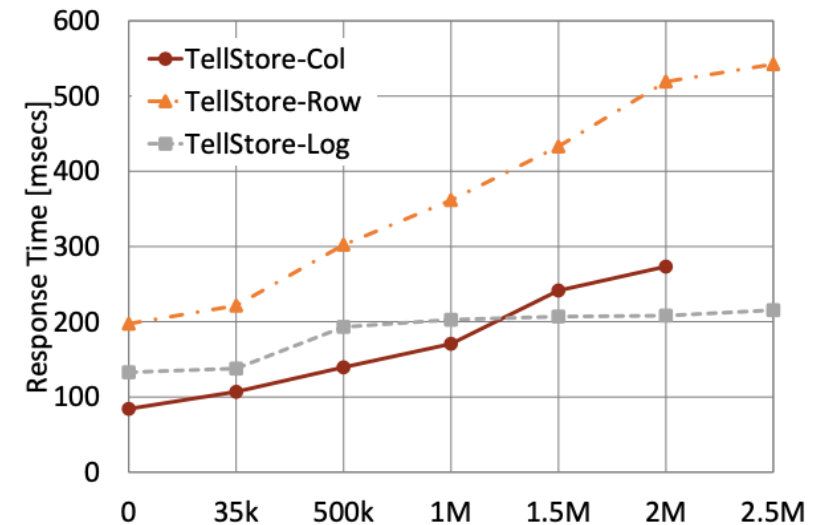
# Experiments: Mixed Workload



(a) TellStore vs Kudu



(b) TellStore, Scale Get Requests



(c) TellStore, Scale Get/Put Requests

**Figure 11: Exp 4, Response Time: YCSB# Query 1, 4 Storage Nodes**

Contrary to competition,  
scan perf. is stable with  
more gets/puts

In the absence of updates  
TellStore scales perfectly:  
scans+gets go to different  
cores

With 50% updates  
eventually logging wins

# Things to remember

KVS vs. Scans: how to compromise, navigate the design space

- ✓ delta-main vs. log-structure
- ✓ chained vs. clustered versions
- ✓ row-major vs. column-major
- ✓ lazy vs. eager GC

## F2: Designing a Key-Value Store for Large Skewed Workloads

Konstantinos Kanellis\*

University of Wisconsin-Madison

kkanellis@cs.wisc.edu

Badrish Chandramouli

Microsoft Research

badrishc@microsoft.com

Shivaram Venkataraman

University of Wisconsin-Madison

shivaram@cs.wisc.edu

A new set of requirements:

**point** put/get operations (need for high throughput, ideally full in-memory for hot data)

**larger-than-memory** working set (so full in-memory not possible)

**high skew** in access patterns (reads/writes)

The target application has **no scans**!

which design should they follow?



# Key-value stores

## Log-structured Merge (LSM) Trees

- Handle *larger-than-memory* workloads
- Organized in levels; first is in-memory
- Support both *point* & range queries
- Avoid I/O by employing (Bloom) *filters*
- Judicious use of main memory



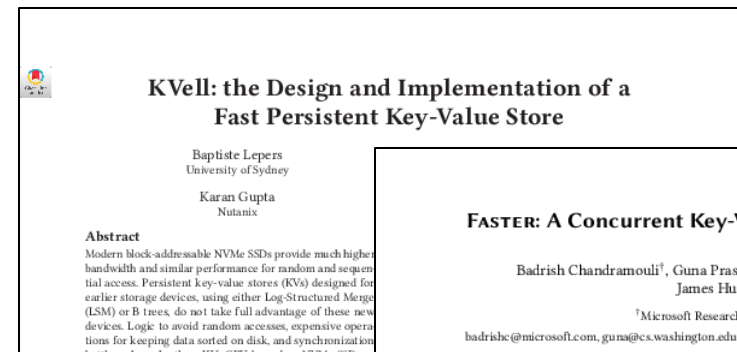
RocksDB



SPLINTERDB

## Point-optimized Stores

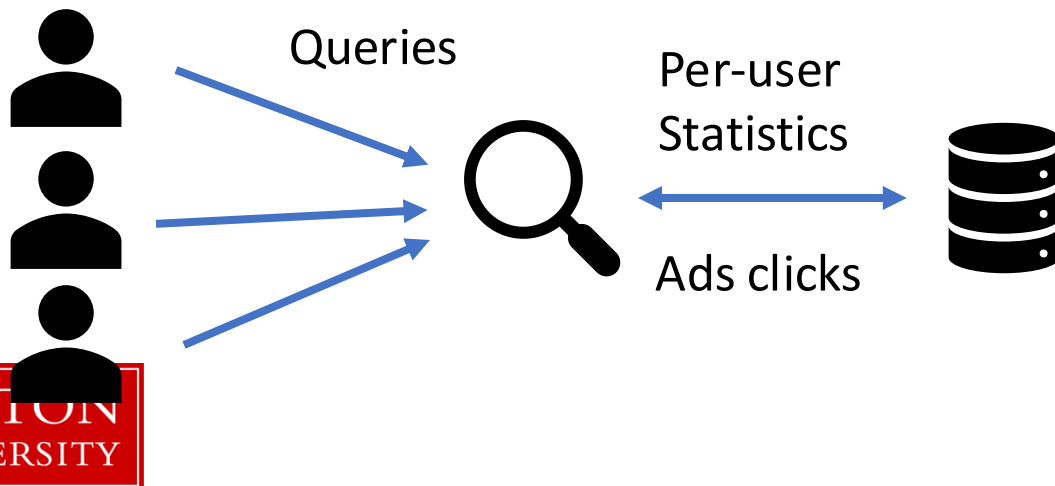
- Focus on use-cases like *web caching*
- Large *in-memory* index structures
- *Latch-free* concurrent designs
- Saturate I/O (even for NVMe SSDs)
- Very *high* throughput (>1M ops/sec)



# Real-world, large *skewed* workloads

- **Point** queries and **high** throughput *paramount*
- **Working sets** *larger* than main-memory – most data rarely accessed or updated
- Total **indexed** data order of magnitude *larger* than main-memory
- Natural **skew** in key access pattern – both for *reads* and *writes*
- Memory resources **scarce** – disk **wear** is a practical *concern*

## Search Engine Workload



Insert / Update **statistics** (e.g., clicks, statistics)

Millions of **active** users at any time – *critical* path

Many more *inactive*, but we still need to keep data!

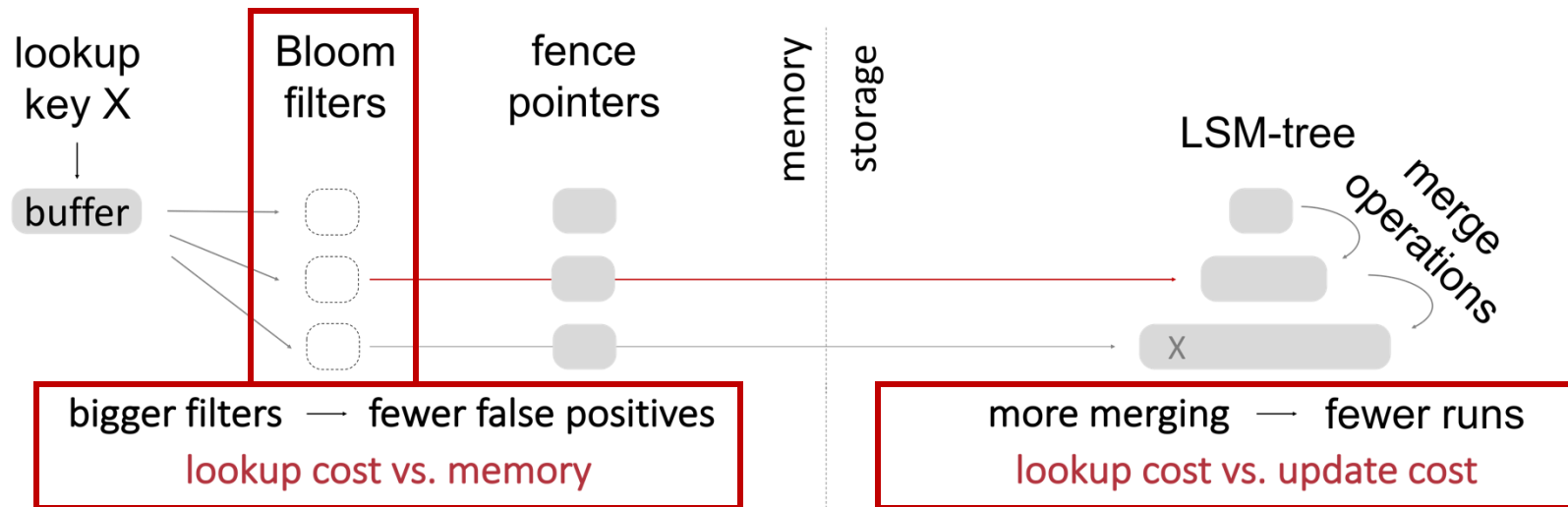
We fetch clicks for active users (during browsing)

We count viewed ads, from active users

# Limitations of Existing Systems (1)

## Log-structured Merge (LSM) Trees

- + Enable and tune (Bloom) *filters*
- + Use hash indices
- + Efficient compaction policies
- Filters may no longer fit in memory
- CPU overhead (10s of filters / query)
- Need tuning





# Limitations of Existing Systems (2)

## KVell

- Adjust page-cache size
  - Large in-memory B-tree (19B per key)
  - B-tree continuously paged out to disk

## FASTER

- Tune record log in-memory size / hash index
  - Reducing index size increases I/O ops
  - Log compaction “pollutes” in-memory log

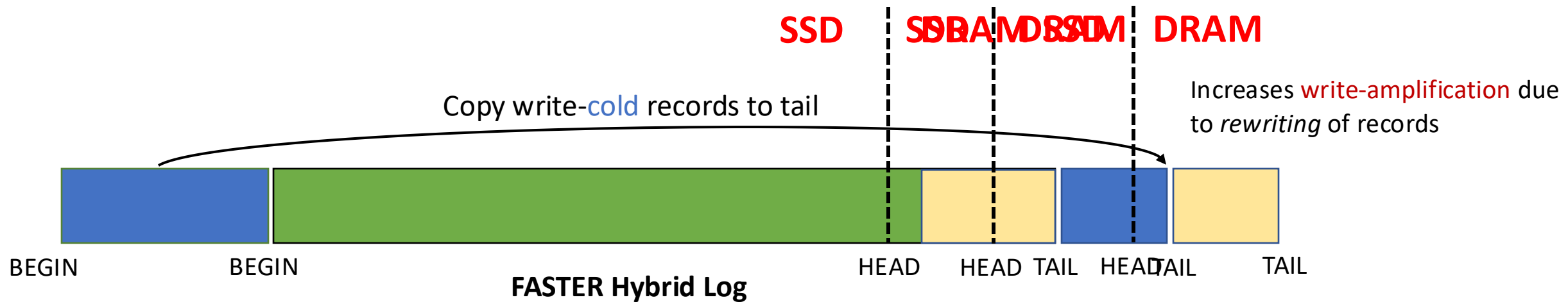
# Limitations of Existing Systems (2)

KVell

- Adjust page-cache size
- Large in-memory B-tree (19B per key)
- B-tree continuously paged out to disk

FASTER

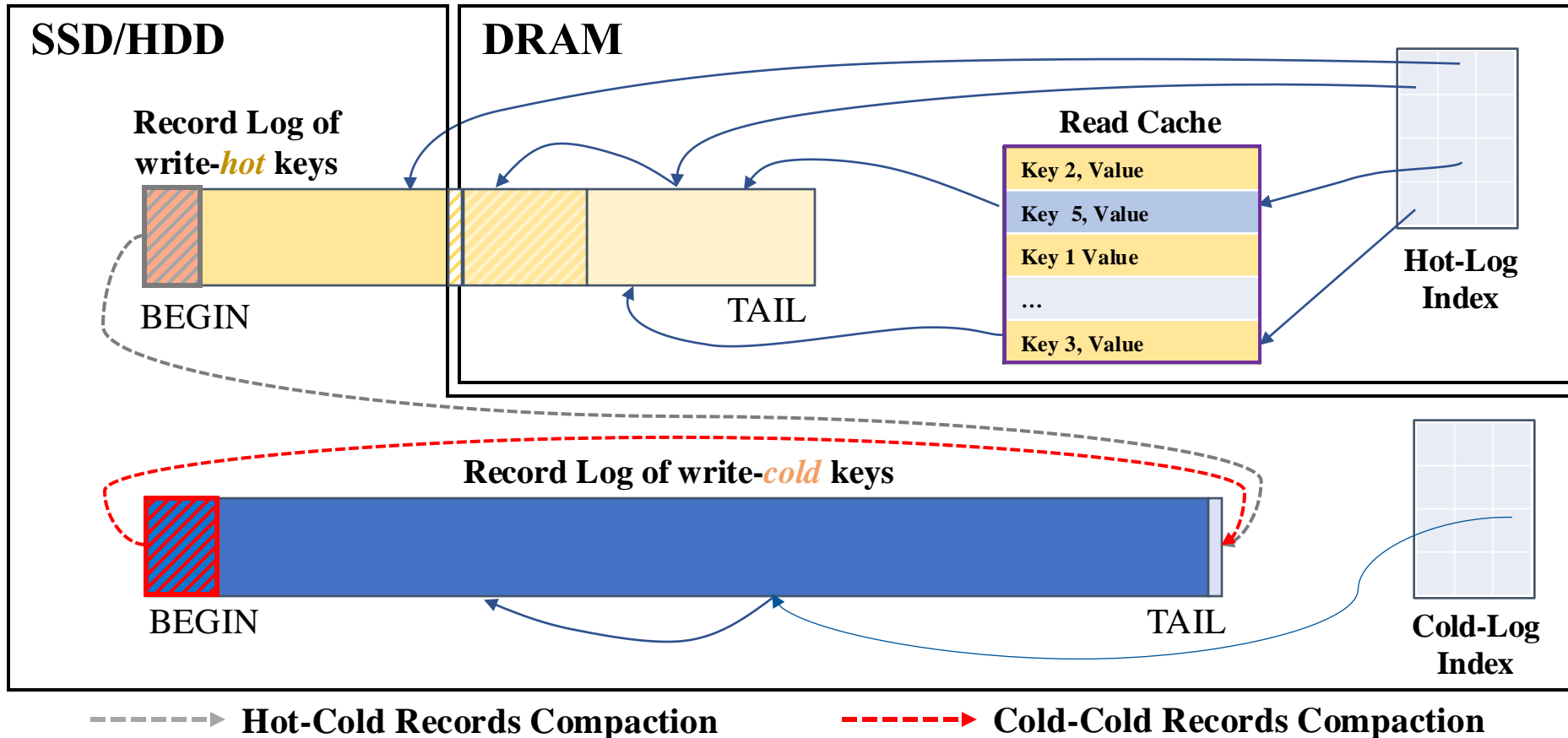
- Tune record log in-memory size / hash index
- Reducing index size increases I/O ops
- Log compaction “pollutes” in-memory log



New records **compete** with old *compacted* records for a spot at *in-memory* regions

# Introducing $F2$

Key Idea: separate **management of records** across both read/write and hot/cold domains



# Design Space

Updates

×

Layout

×

Versioning

×

GC

*in-place*

*log-structured*

*delta-main*

*column (PAX)*

*row*

*clustered*

*chained*

*periodic*

*piggy-backed*

similar to TellStore-Log, but with **periodic compaction**

**compaction aggressively optimized**

# F2 – Read Cache

Contains *read*-hot records of both **hot** log and **cold** log

- Hash index entries can point to *either* read cache entries or (**hot**) record log (single bit in address)

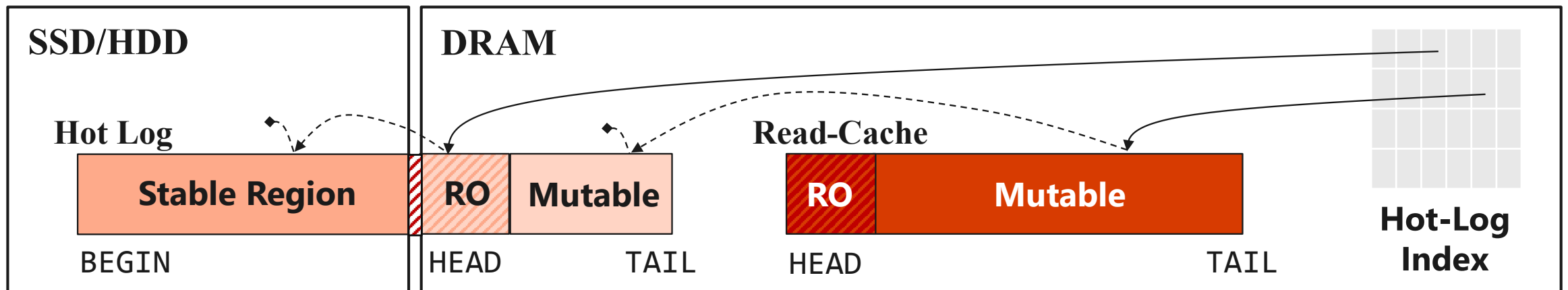
**Reads** go from **hot**-log index → (optionally) read cache → **hot** log

**Reads** from **cold**-log are *always* inserted at the tail of the read cache

**Upserts** and **RMWs** write *directly* to **hot** log tail, eliminating read cache chain

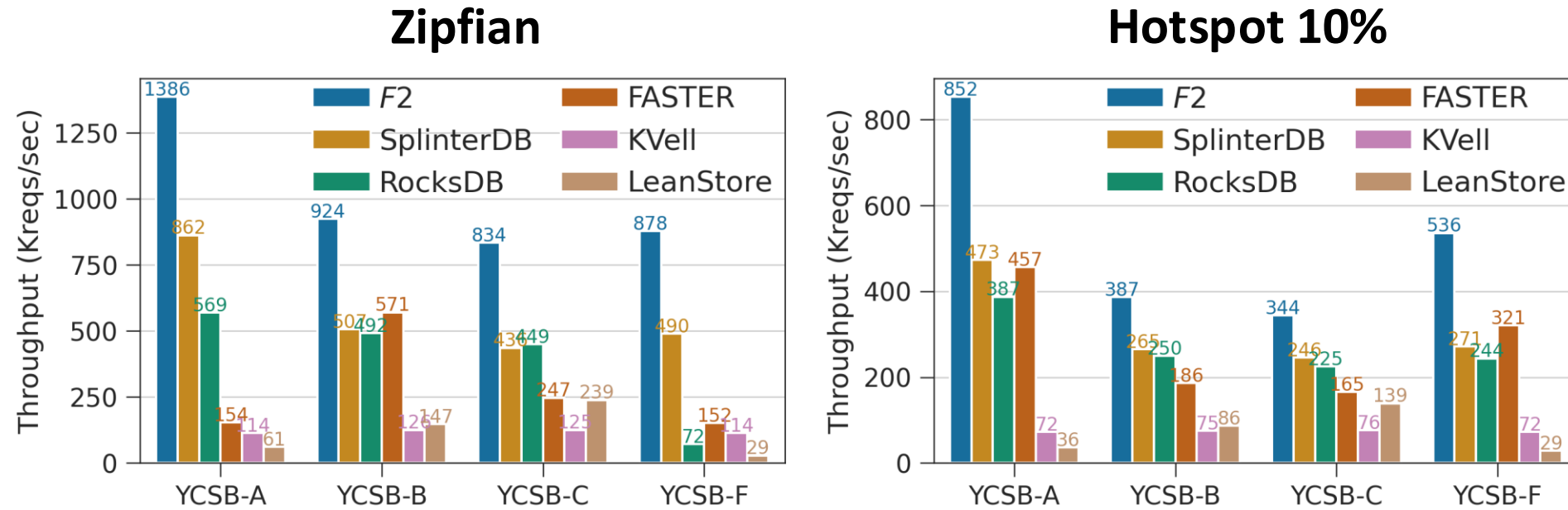
- If record with same key in read cache, it's **invalidated**!

Periodically, read cache is *evicting* in-memory records (HEAD), by altering the **hash chains**



# F2 – Performance Comparison

YCSB: 250M keys, 8B keys, 100B values; **3GiB** mem budget (**10%** of dataset size); 24 threads, NVMe SSD

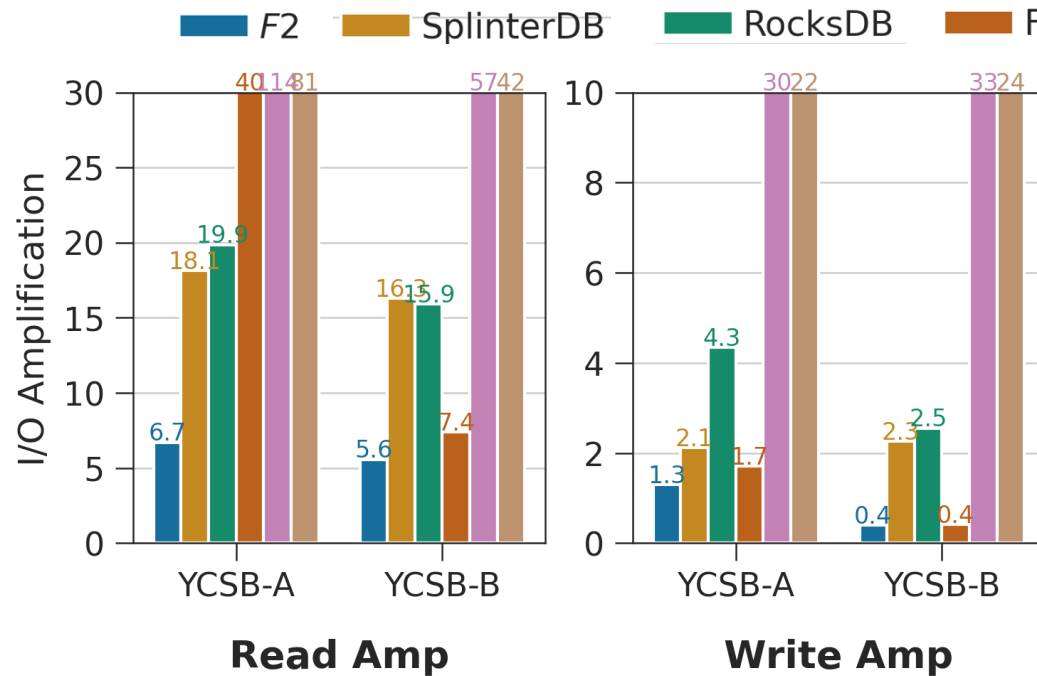


**Zipfian:** SplinterDB (**1.78x**), RocksDB (**4.61x**), FASTER (**4.94x**)

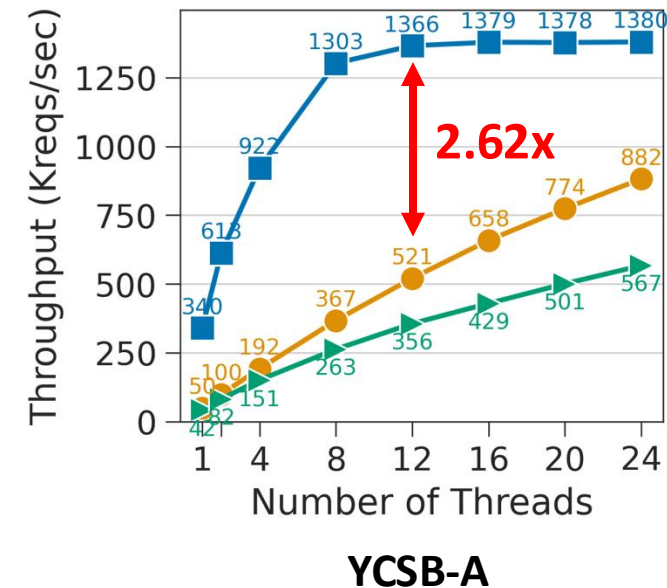
**Hotspot 10%:** SplinterDB (**1.66x**), RocksDB (**1.88x**), FASTER (**1.92x**)

# F2 – I/O Amplification & Scalability Comparison

YCSB: 250M keys, 8B keys, 100B values; **3GiB** mem budget (**10%** of dataset size); 24 threads, NVMe SSD



F2 achieve **less** WA than Baselines



F2 saturates at **8-12** threads

→ Varying memory-budget, YCSB skewness; F2 detailed evaluation ←

## class 7

# Design Tradeoffs in Key-Value Stores

Prof. Manos Athanassoulis

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