# CS 561: **Data Systems Architectures**

## class 7

# Design Tradeoffs in Key-Value Stores

Prof. Manos Athanassoulis

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

Do we have a quiz today … ?

# Yes!

# what to do now?

A) read the syllabus and the website

B) register to Piazza + Gradescope

C) finish project 0

D) finish project 1

E) **finish project proposal (due 2/22)**

F) **[register](#) for the student-presentations (by 3/6)**

# Fast Scans on Key-Value Stores (KVS)

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



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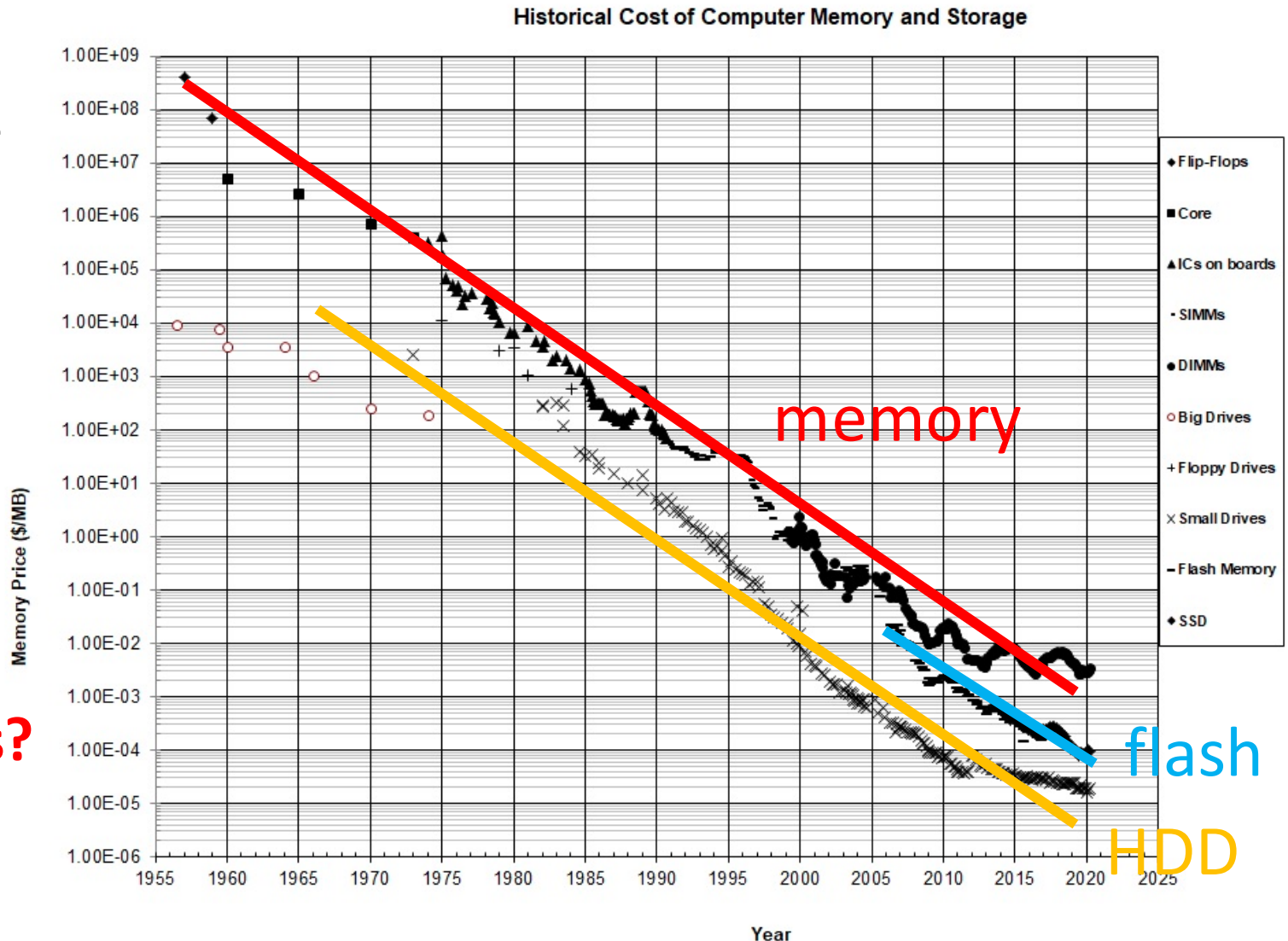


**Can we do both in one system?**

# Why combine KVS and analytical systems?

cheaper and cheaper storage

more data ingestion

**need for write-optimized data structures**

**what about analytical queries?**



Historical Cost of Computer Memory and Storage

memory

flash

HDD

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

Most organizations maintain both

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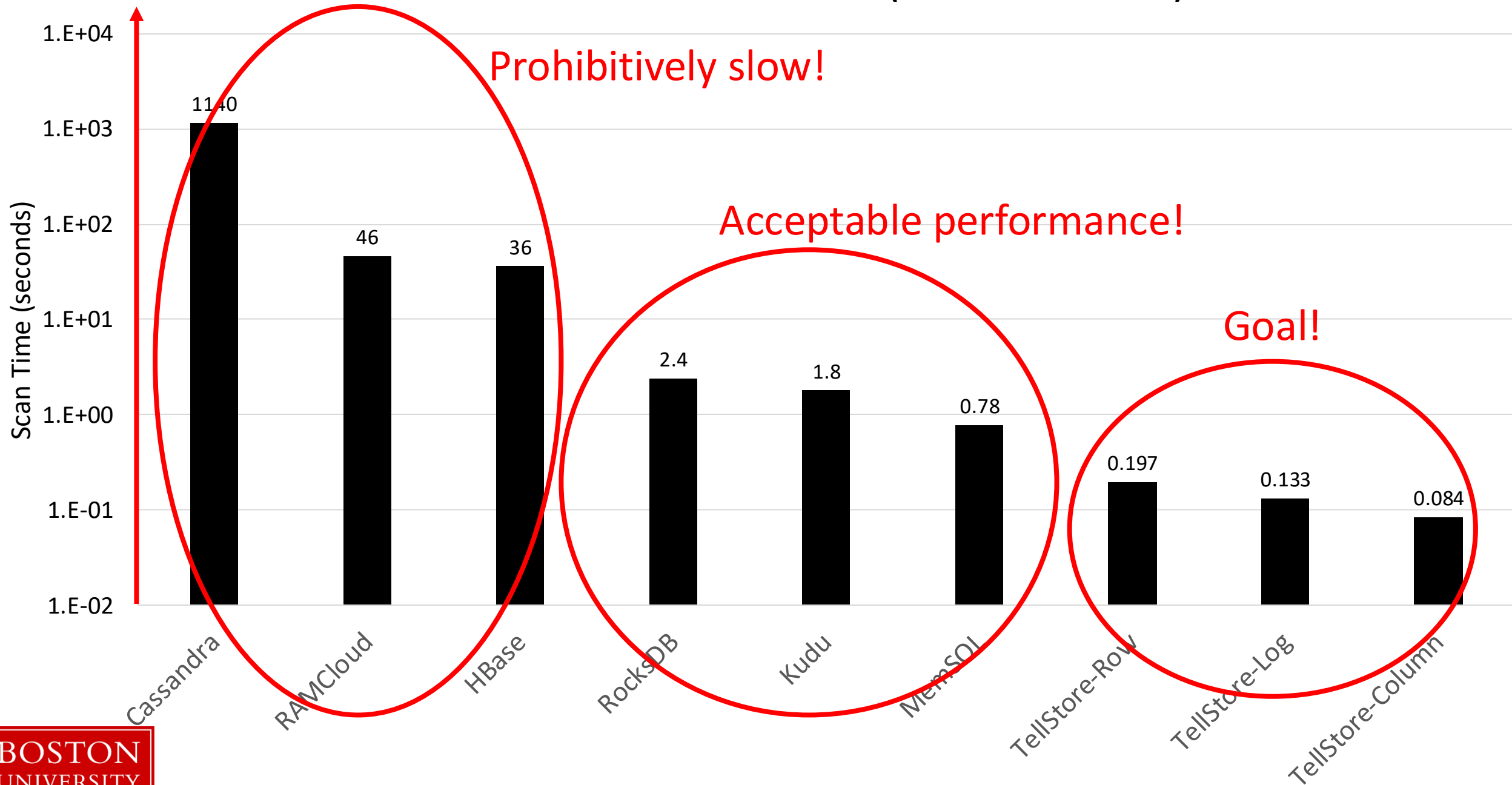- *analytical* systems (often as column-stores)

**problems?**

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

harder to maintain (more systems, more code)

time consuming data integration/transfer

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

Support for:
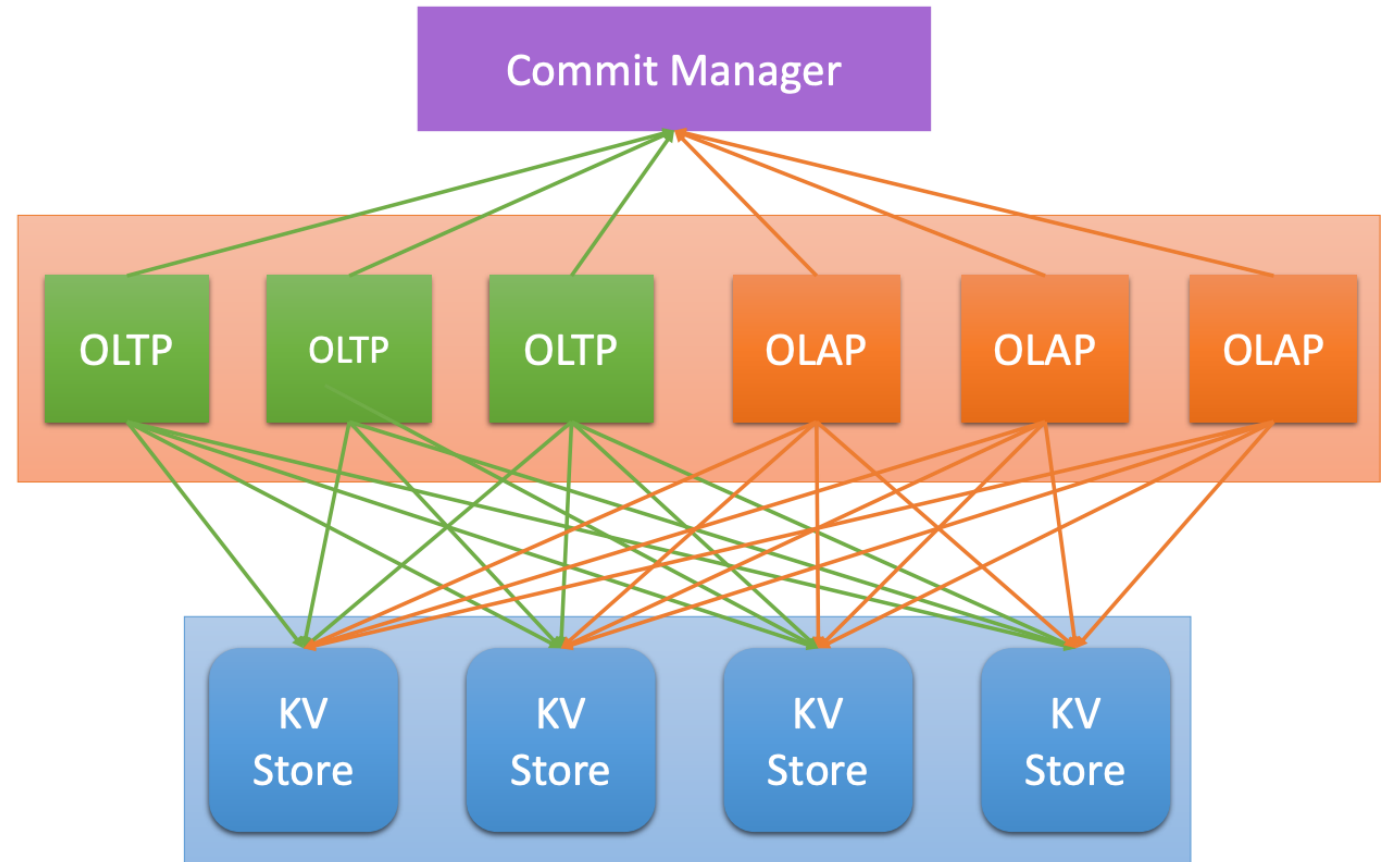 Scans
 Versioning
 Batching

| Scans | Versioning | Batching |
|---|---|---|
| selection | multiple versions through timestamps | batch several requests to the storage layer |
| projection | garbage collection | |
| (simple) aggregates | | amortize the network time |
| shared scans | discarding old versions during scans might be costly | |
| remember them? | | |

# Challenges

**Scans need columnar locality**

## scans vs. get/put

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

get/put need row-wise locality

why?

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

# 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?

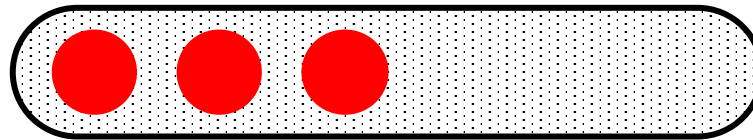**Key design decisions**

(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**

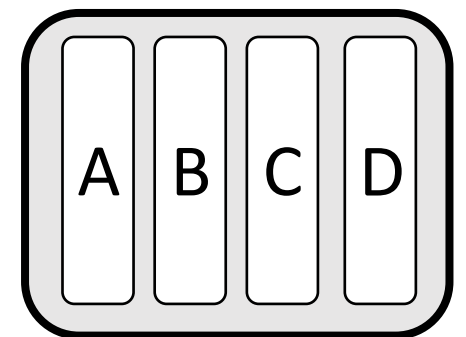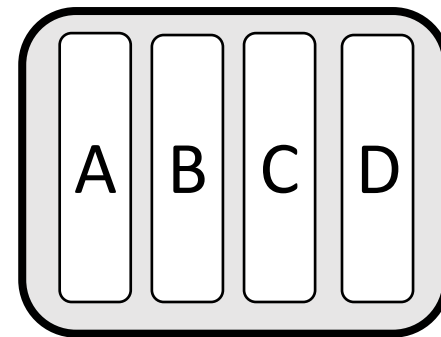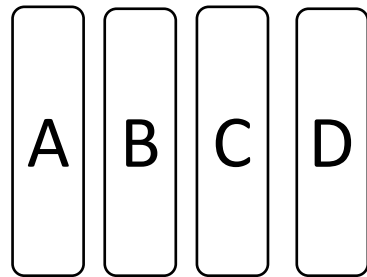(A) Updates      *in-place*      *log-structured*      *delta-main*

# How to design KVS for efficient scans?
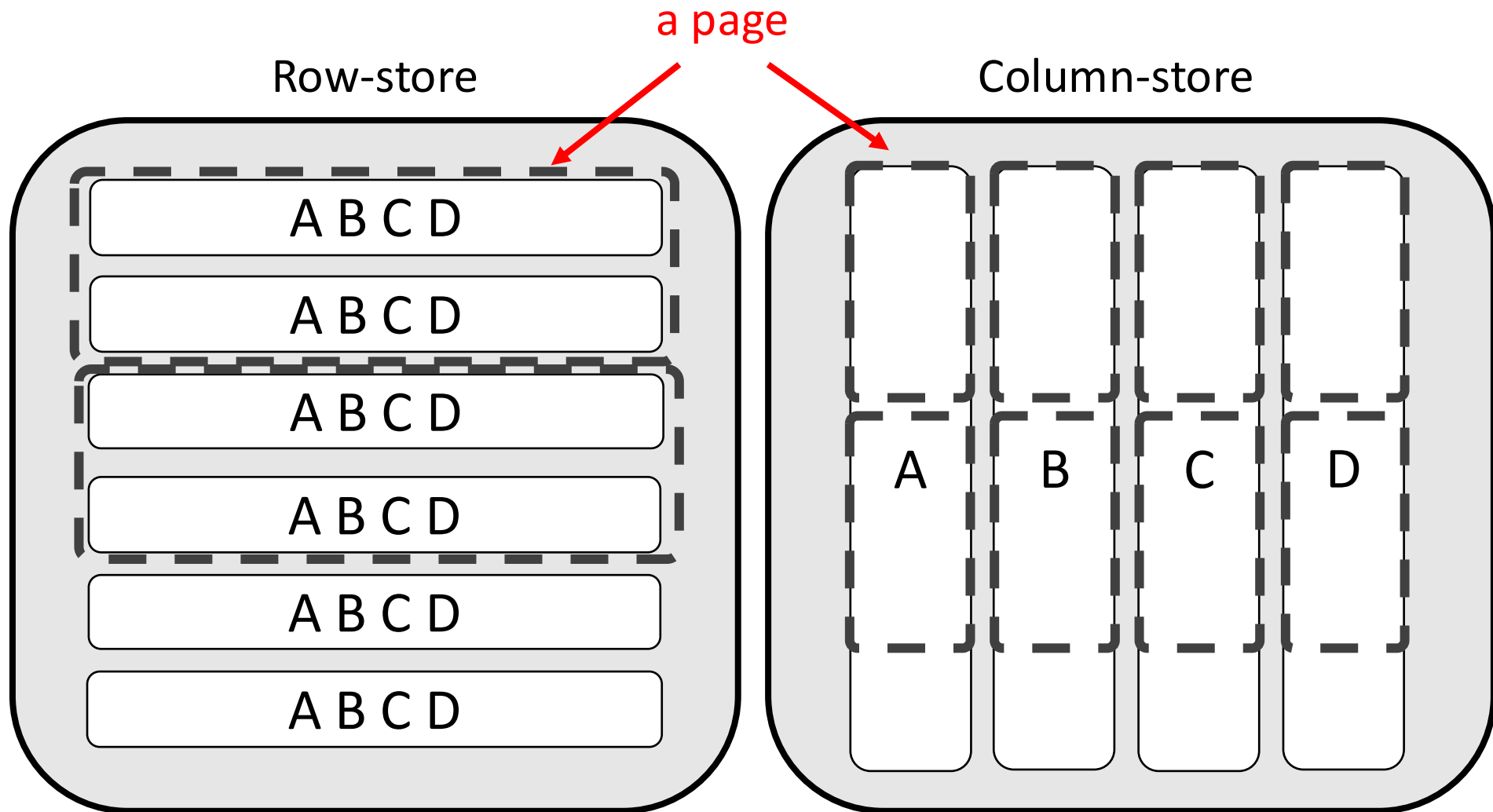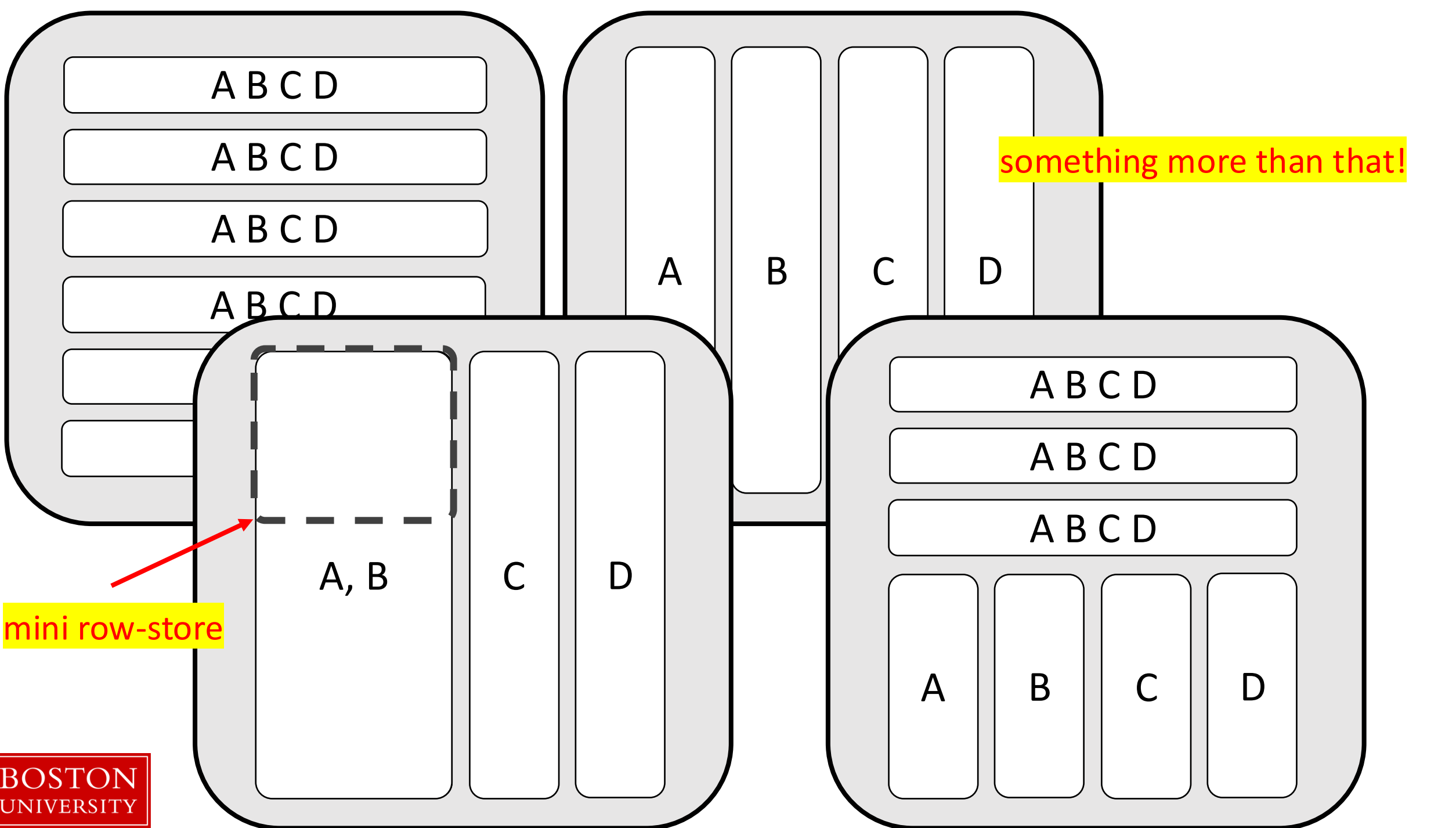
**Key design decisions**

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

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

# Small detour: page layouts

middle ground?

a page

Row-store

Column-store

| 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

BOSTON
UNIVERSITY

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

A B C D

something more than that!

A

B

C

D

mini row-store

A, B

C

D

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

A

B

C

D

BOSTON UNIVERSITY

# Partition Attributes Across (PAX)

Middle ground?

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

✓ Cache-friendly
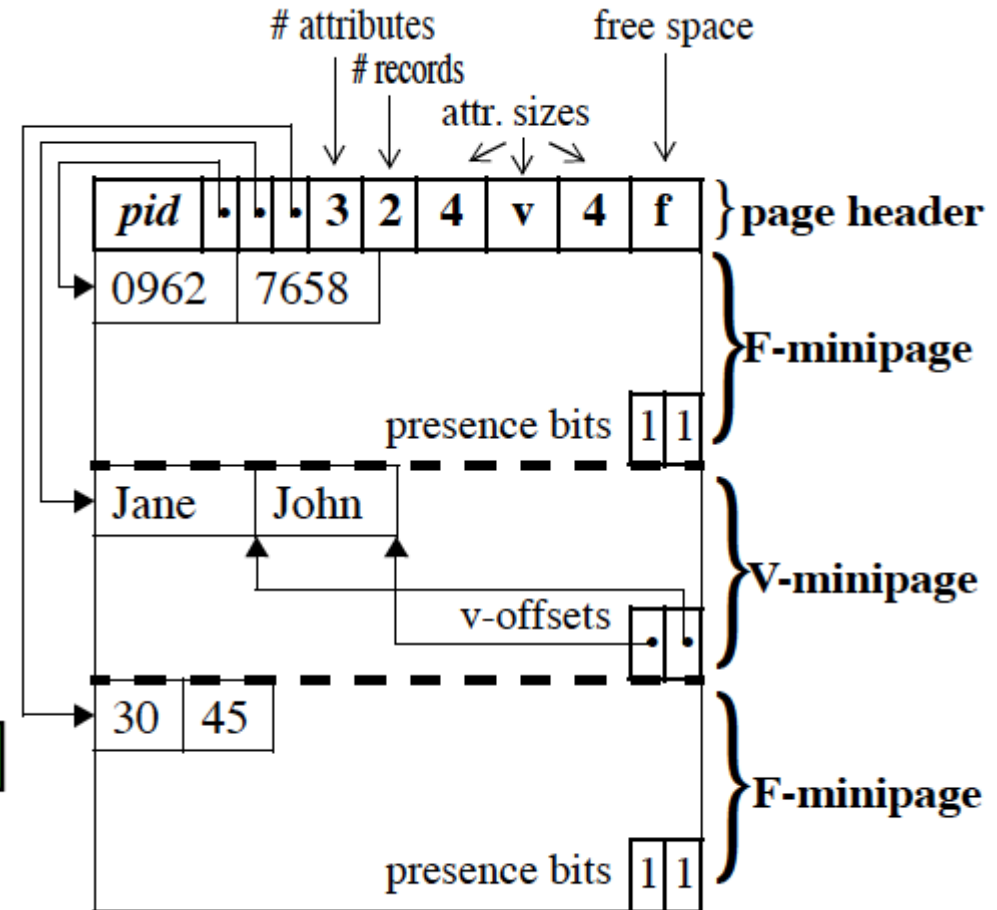  ➤ Brings only relevant attributes to cache

- Compatible with slotted-pages? ✅
- Same update abstraction? ✅
  – (insert in a page)

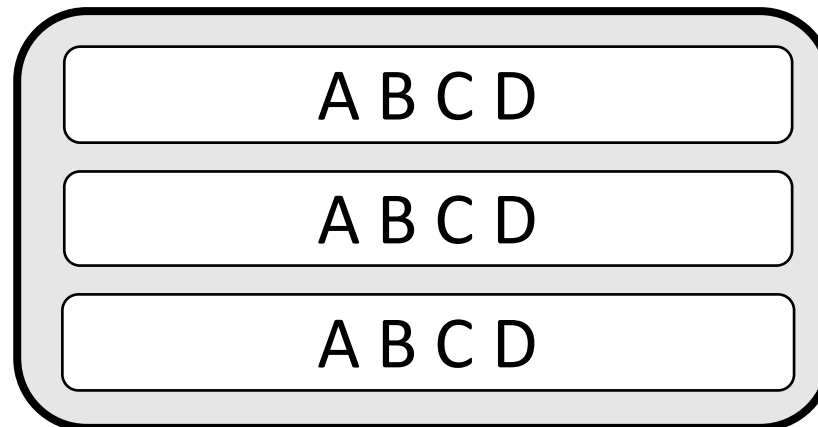# How to design KVS for efficient scans?

**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*

| A B C D (v1) |
| A B C D (v2) |
| A B C D |

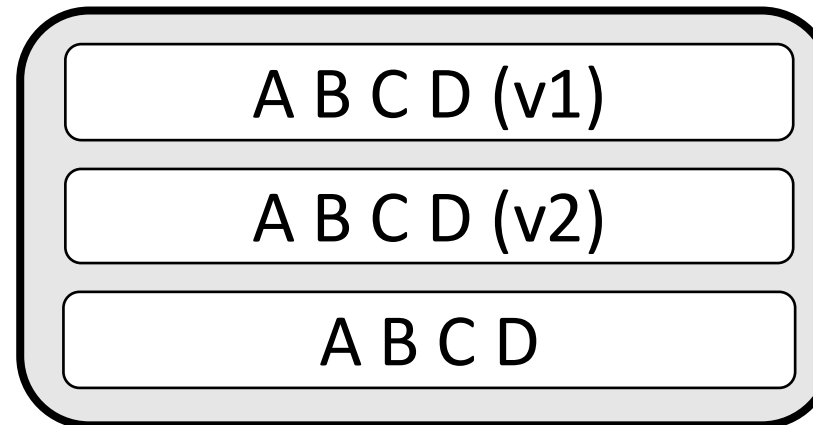any other options?

# 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*

A B C D (v1)   →   A B C D (v2)

A B C D

A B C D

# How to design KVS for efficient scans?
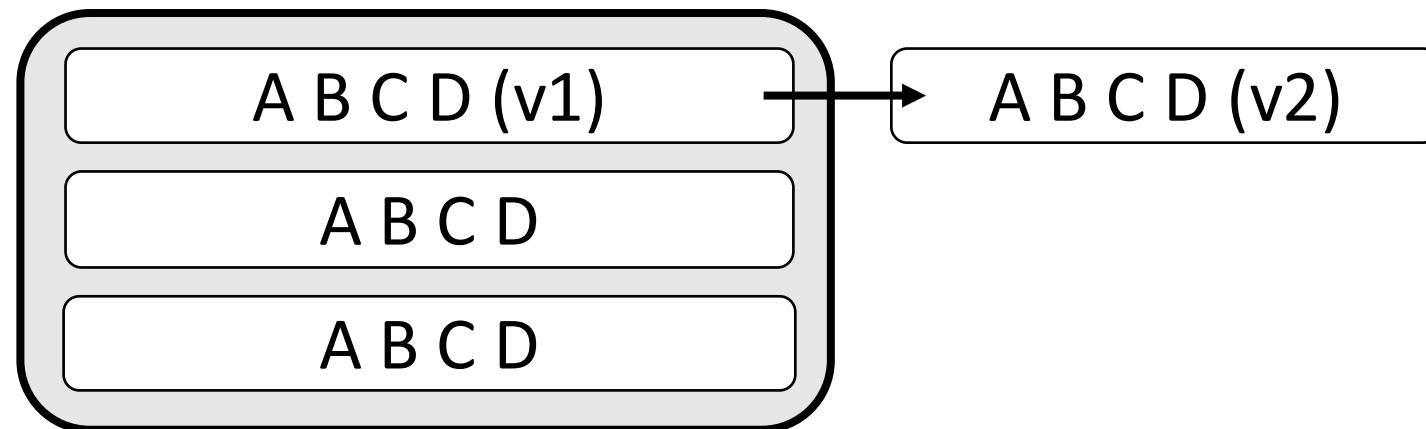
**Key design decisions**

(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**

<span style="color:red">hybrid designs are also valid!
should we consider all possible designs?</span>

BOSTON
UNIVERSITY

# 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

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

*TellStore-Log*

focus on two extremes:
     (1) log-structured & row & chained

# 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**

(2) delta-main & column & clustered

**TellStore-Col**

# TellStore-Log

lock-free hash table

**Hash Table**

| Key | Newest Value |
|-----|--------------|
|     |              |
| Foo |              |
|     |              |
| Bar |              |
|     |              |

| Key | validFrom | validTo | previous | isValid | data |
|-----|-----------|---------|----------|---------|------|
| Foo | 128       | 180     | ptr      | True    | .... |

immutable log

**Log**    head

previous entry for the same key    why?

BOSTON UNIVERSITY

# TellStore-Log Insertion

Hash Table

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

Log

the log contains *rows*

# TellStore-Log Update

Hash Table

Log

previous pointer
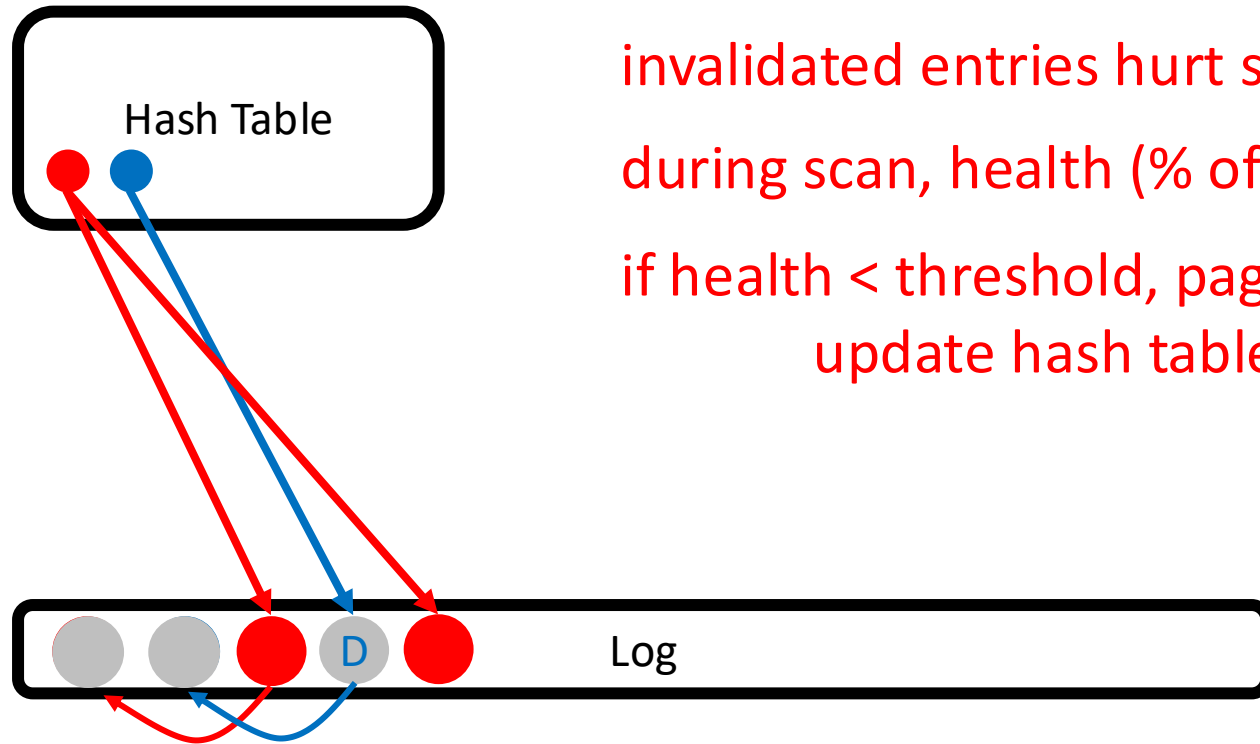
the log contains *rows*

# TellStore-Log Delete



Hash Table

Log

previous pointer

the log contains *rows*

# TellStore-Log Garbage Collection

Hash Table

Log

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)

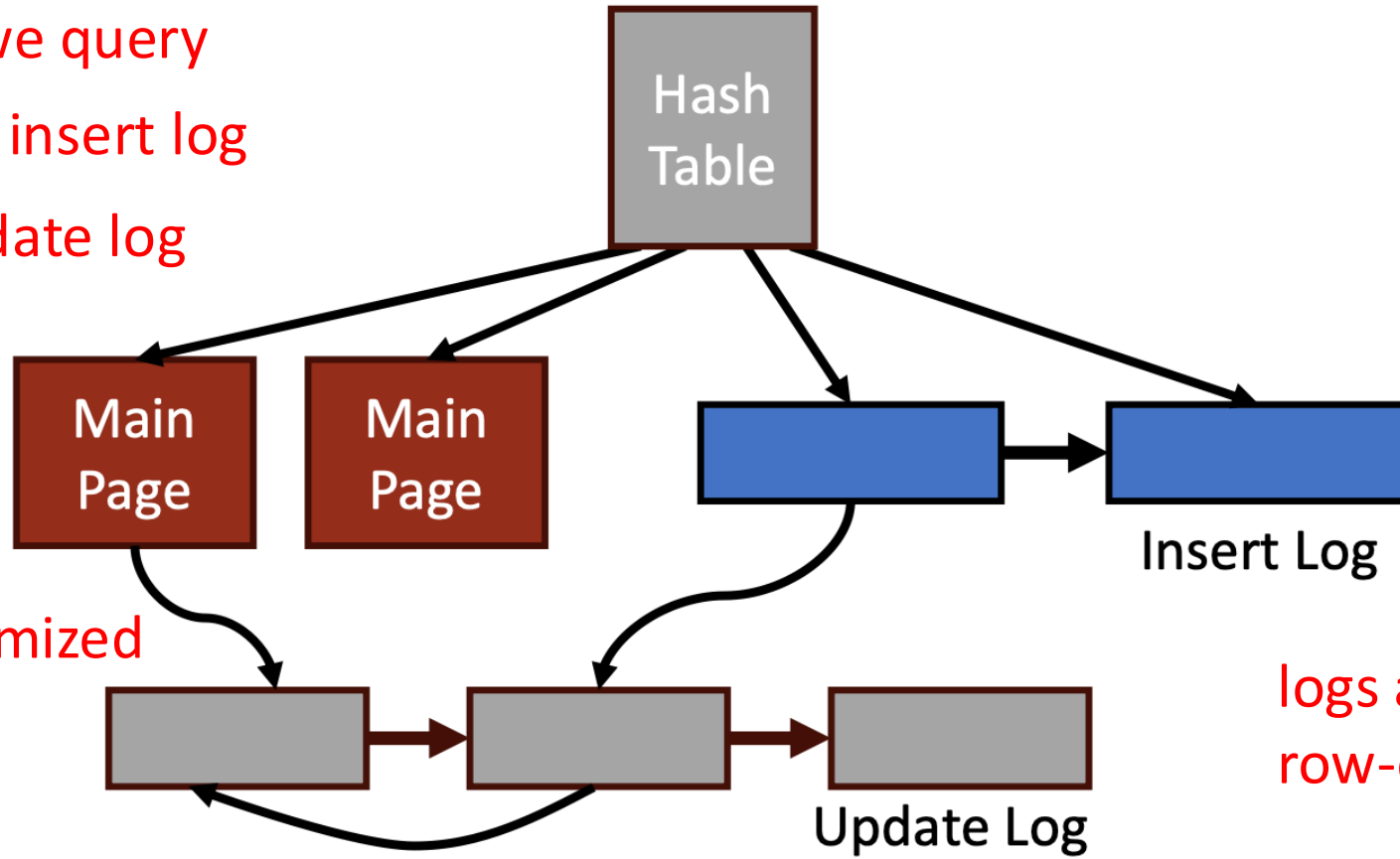**snapshot Isolation:** high throughput, no locks needed

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

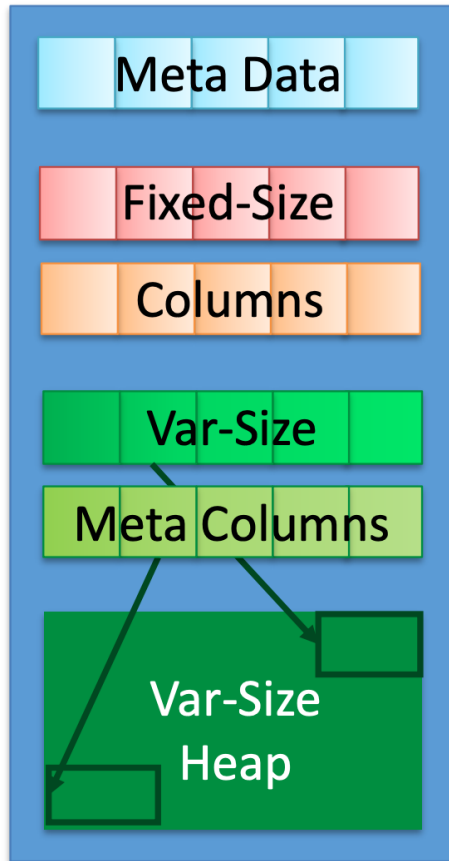**lazy GC:** Optimize tables that are scanned

# TellStore-Col

four data structures

before inserting we query

new entries go to insert log

updates go to update log



Hash Table

Main Page

Main Page

Insert Log

Update Log

main is read-only
columnar: read-optimized

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

# TellStore-Col Layout
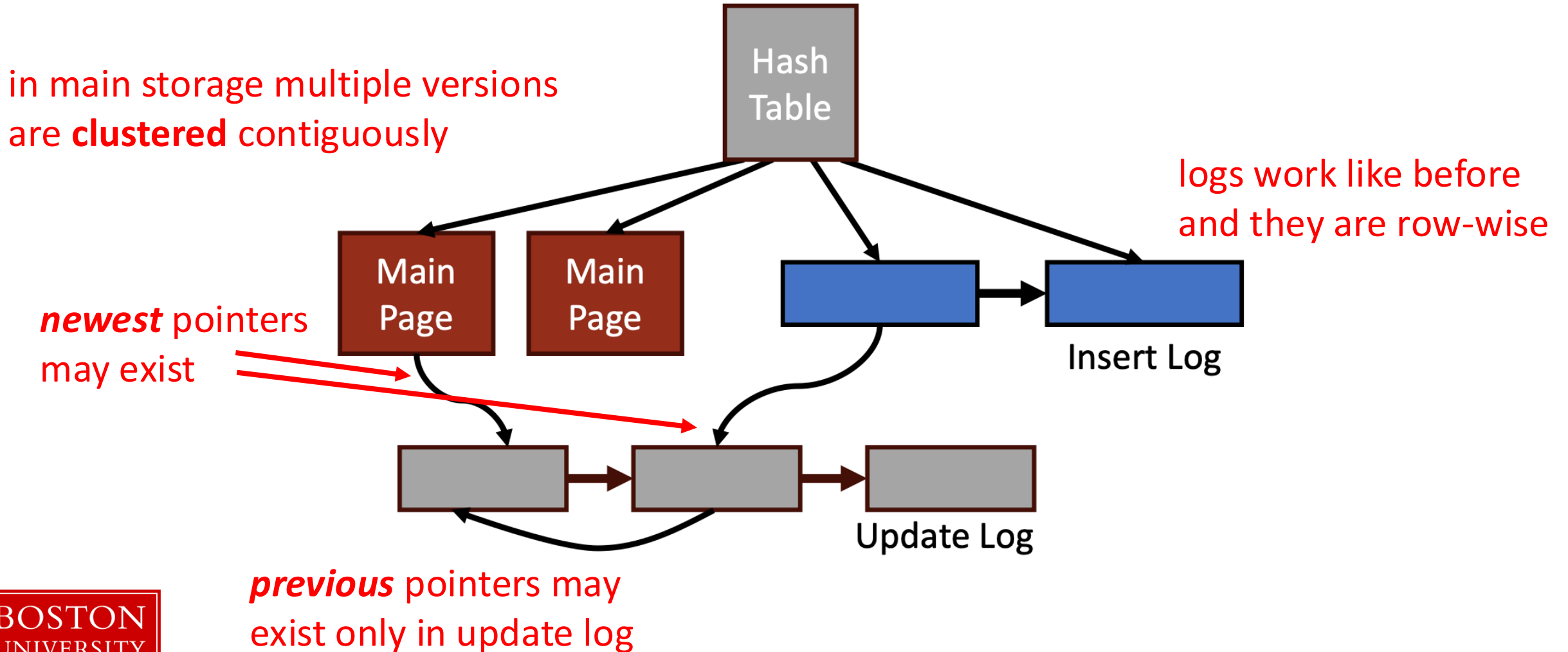


fixed-size data is stored in columnar format

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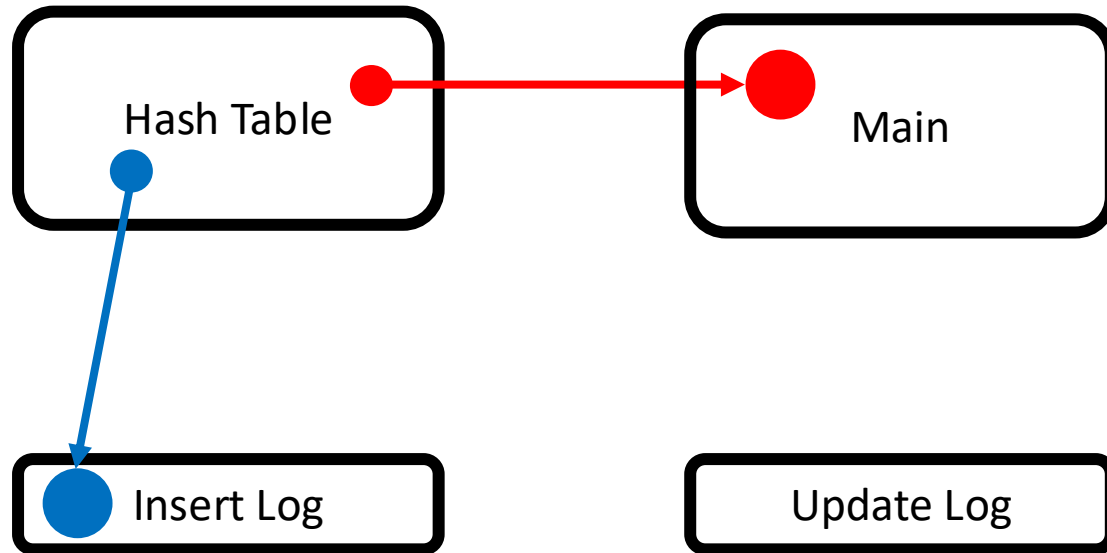but stored in row-wise format

why row-wise?

(1) faster materialization (contiguous copying)
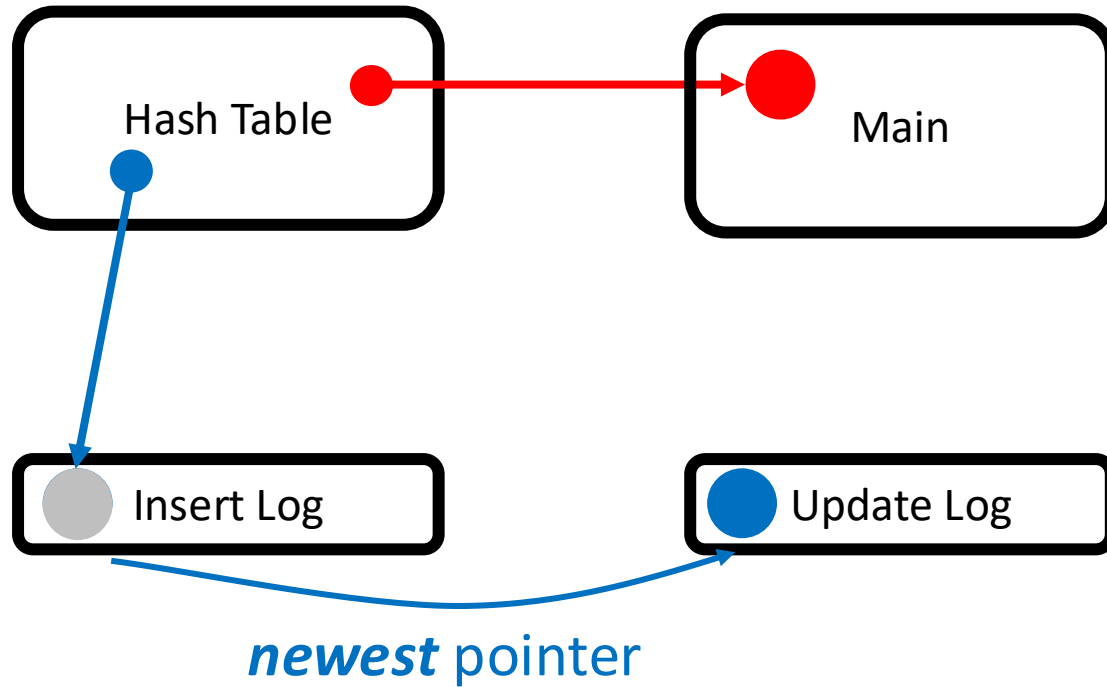
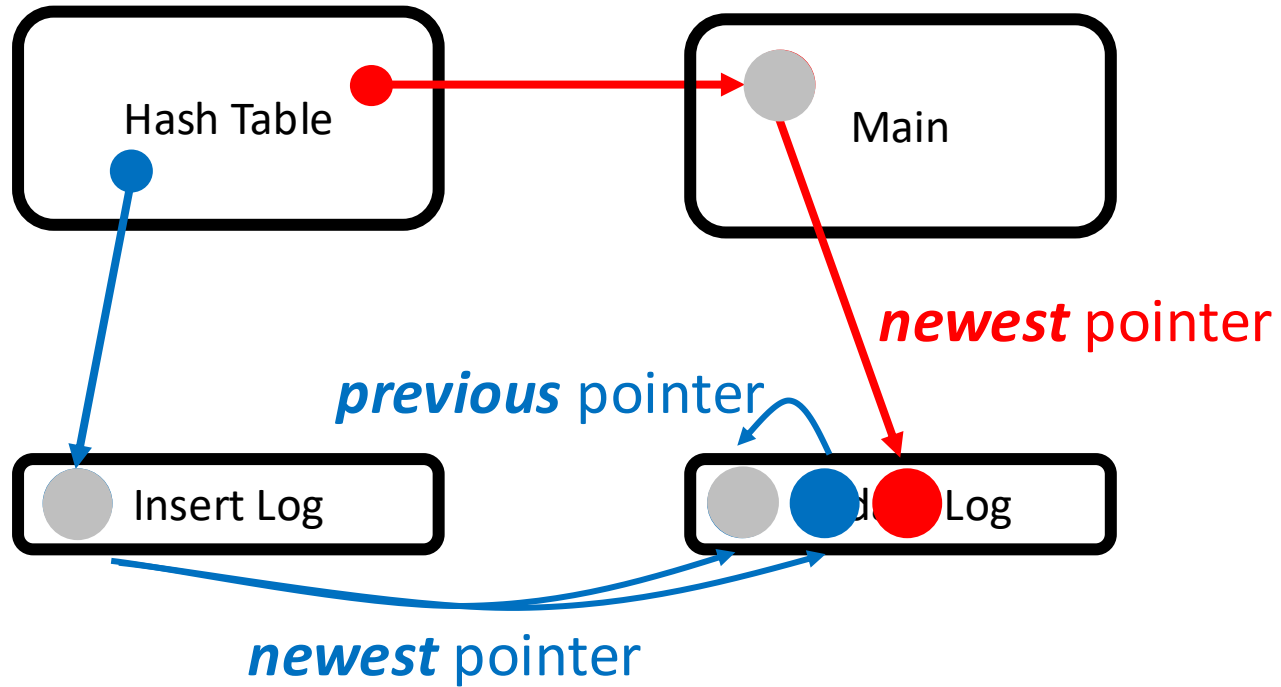(2) less metadata (one offset for many columns)

# TellStore-Col Versioning

in main storage multiple versions
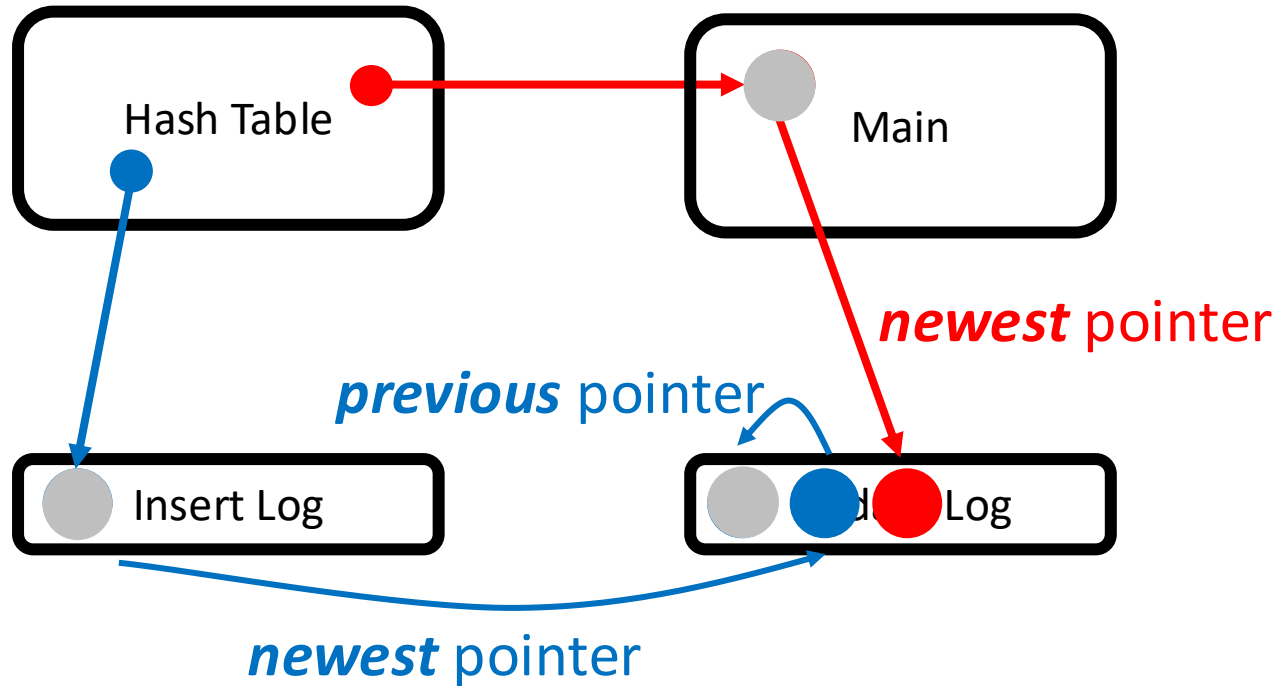are **clustered** contiguously

*newest* pointers
may exist

*previous* pointers may
exist only in update log

logs work like before
and they are row-wise

Hash
Table

Main
Page

Main
Page

Insert Log

Update Log

# TellStore-Col Insertion

# TellStore-Col Update



*newest* pointer

# TellStore-Col Update



Hash Table

Main

***newest*** pointer

***previous*** pointer

Insert Log

Log

***newest*** pointer

BOSTON UNIVERSITY

# TellStore-Col Garbage Collection



Hash Table

Main

**newest** pointer

**previous** pointer

Insert Log

Log

**newest** pointer

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

BOSTON
UNIVERSITY

# 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)

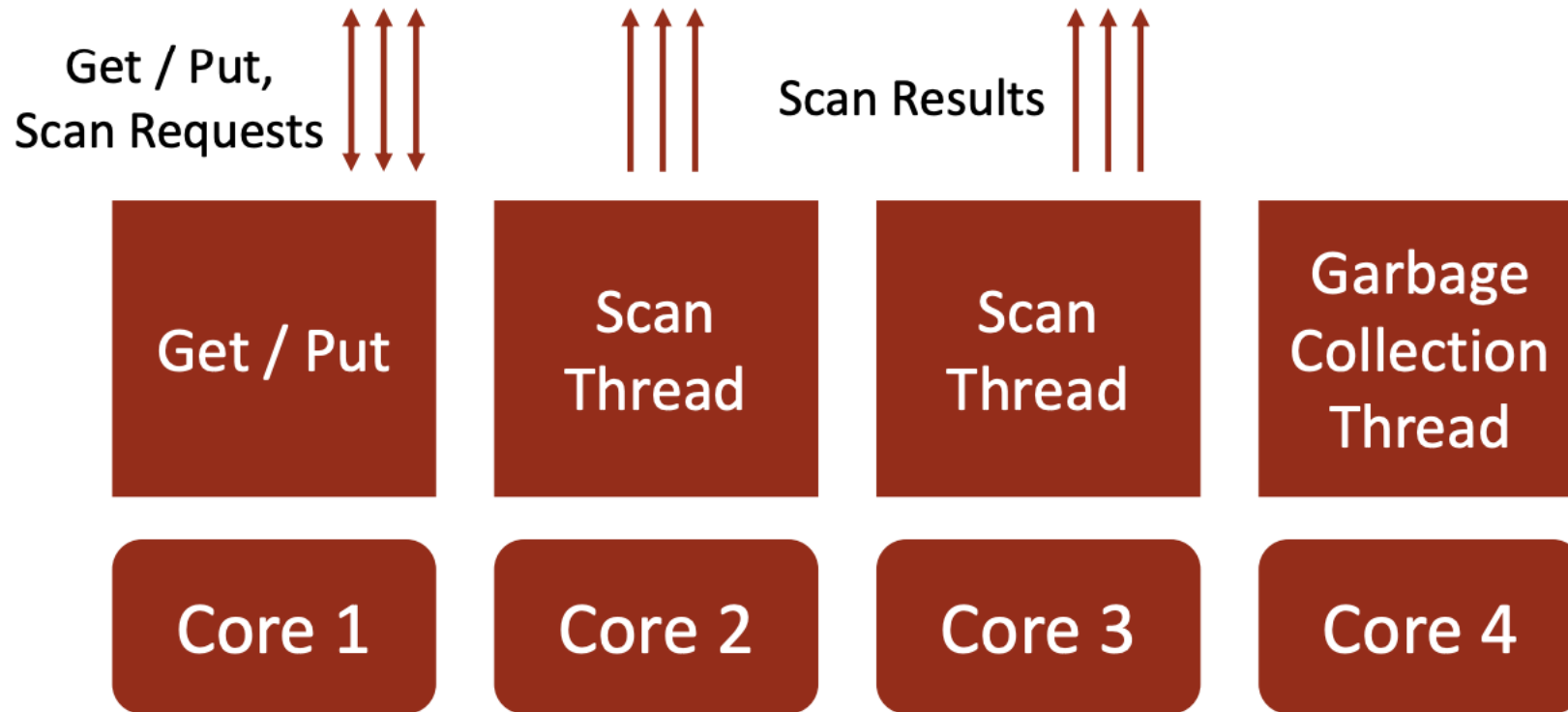*PAX layout:* minimize disk I/O, maintain locality for scans

*separate insert/update logs:* efficient GC

*eager GC:* improve scans

# Implementation Details

scans are assigned to dedicated threads          scan coordinator for shared scans

Get / Put,
Scan Requests

Scan Results

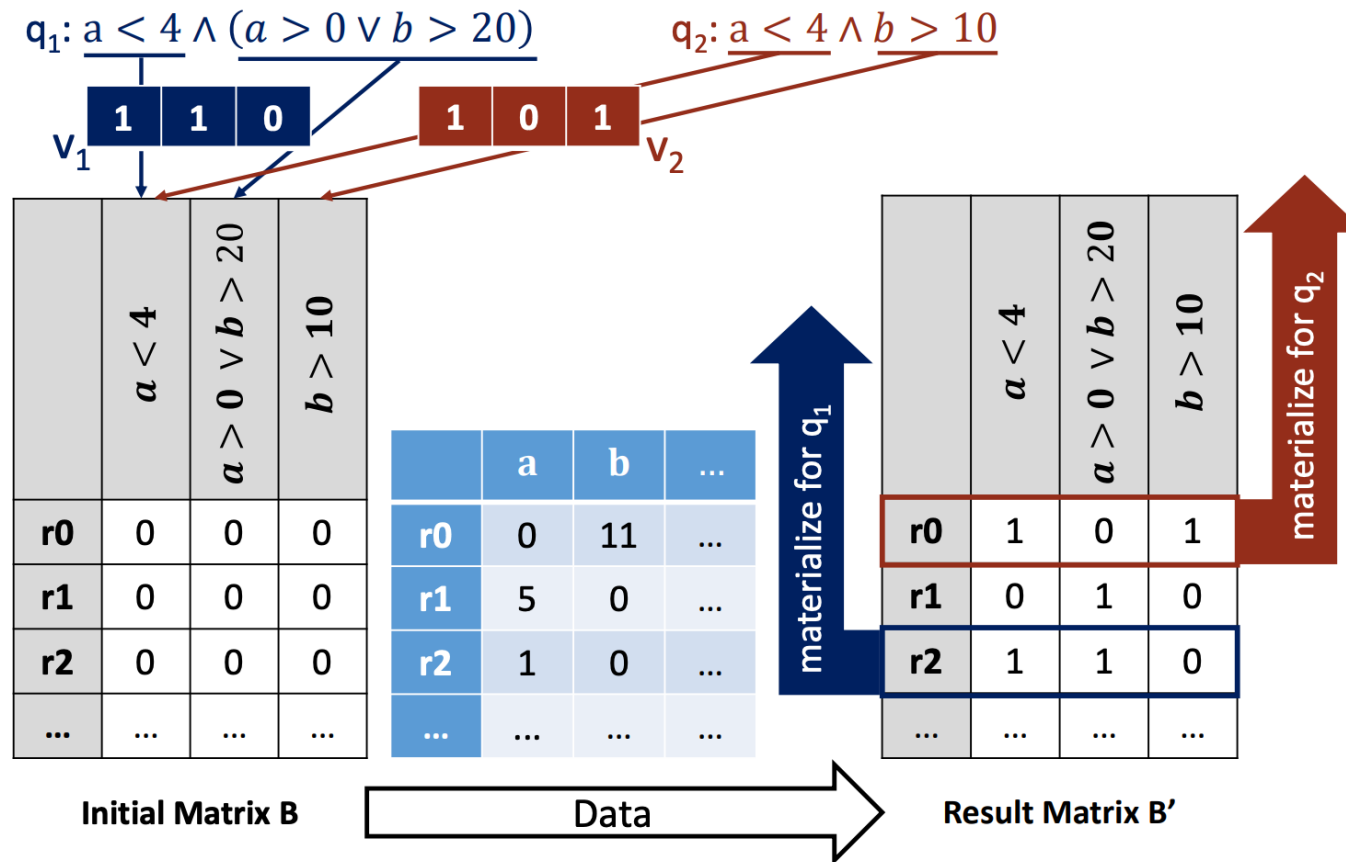| Get / Put | Scan Thread | Scan Thread | Garbage Collection Thread |
|-----------|-------------|-------------|---------------------------|
| Core 1 | Core 2 | Core 3 | Core 4 |

# 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 ley** | 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
  ```
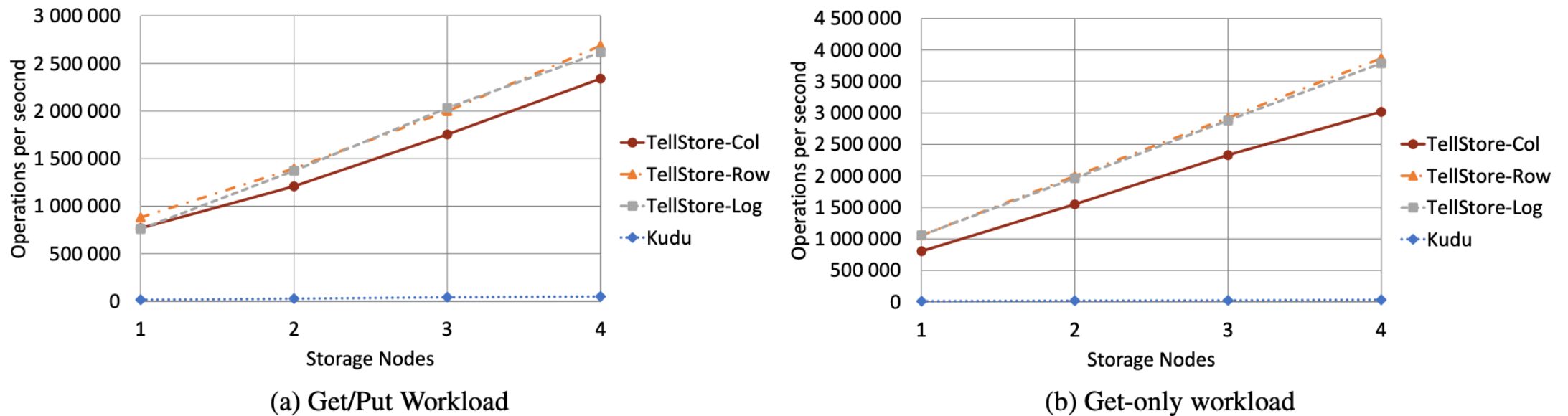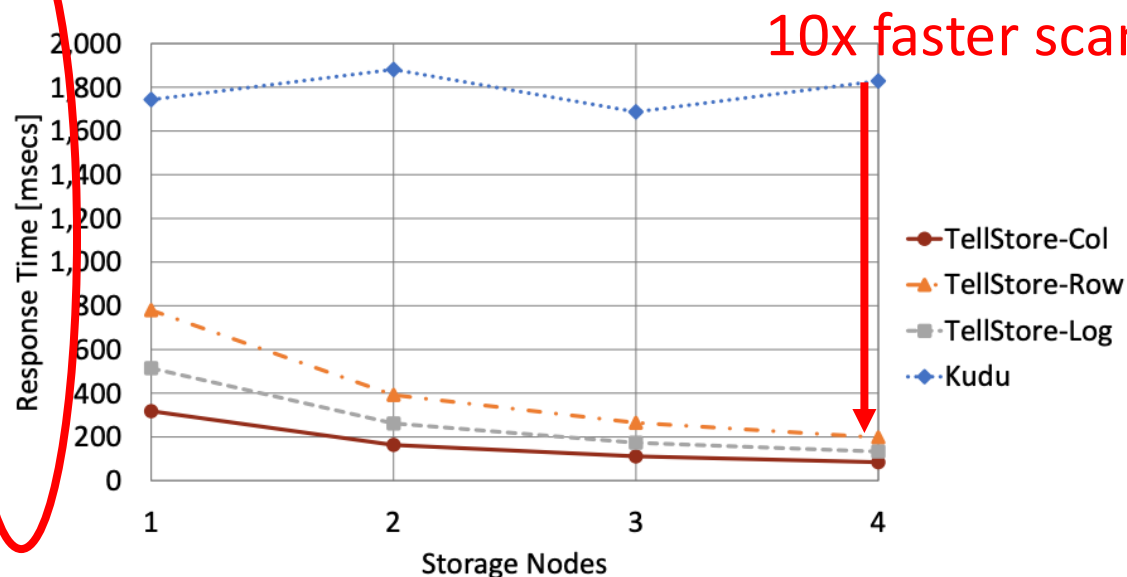
# Experiments: Transactional Workload



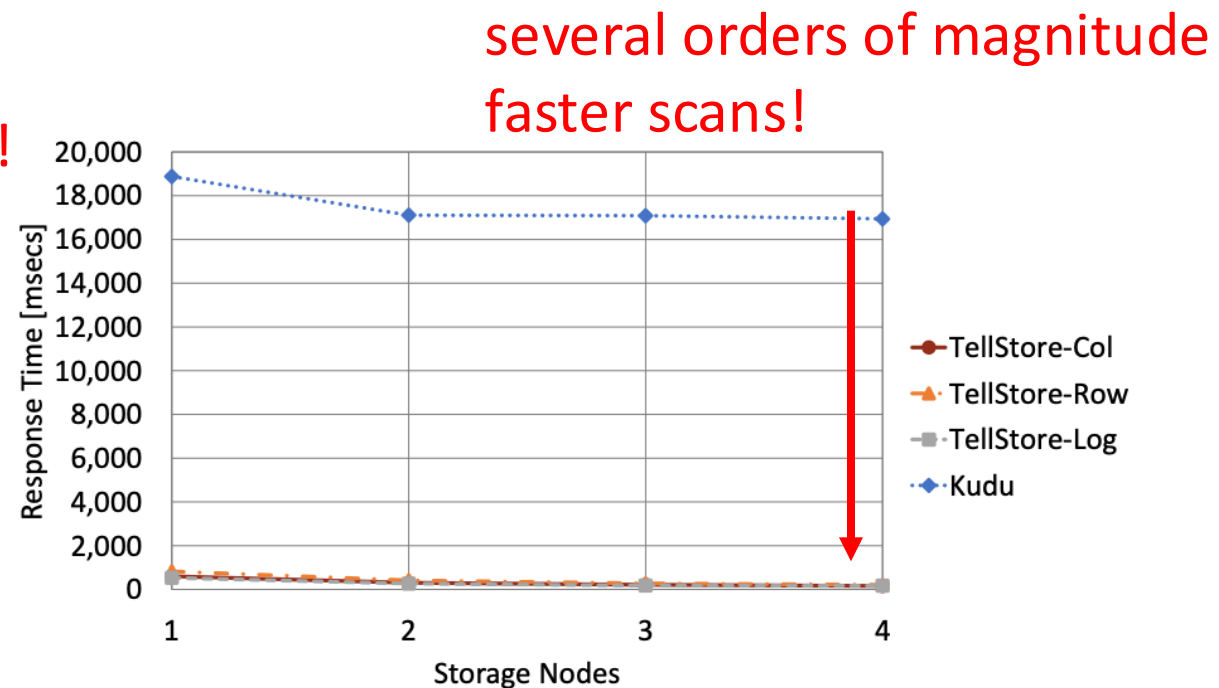Figure 8: Exp 1, Throughput: YCSB, TellStore Variants and Kudu, Vary Storage Nodes

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

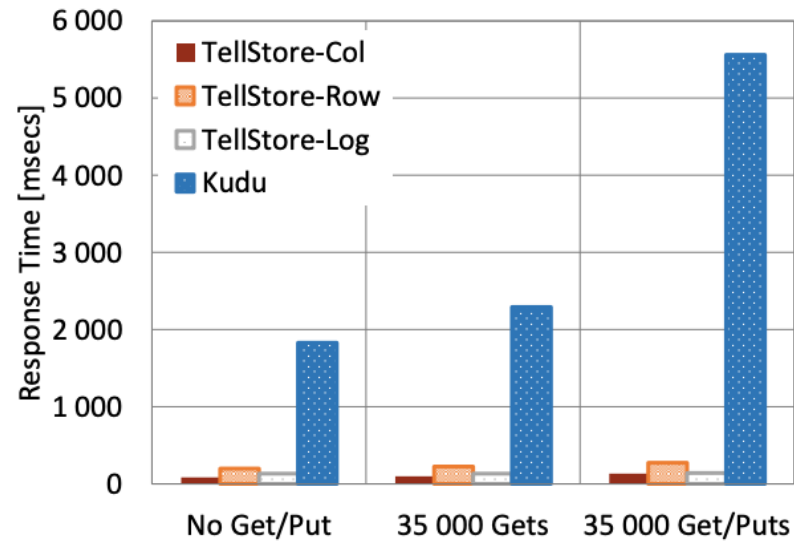All TellStore approaches are not that far!
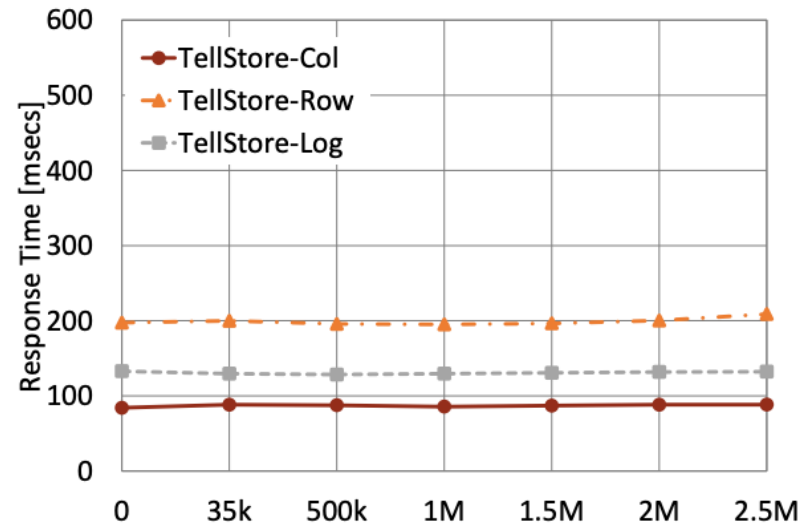
# Experiments: Scans



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

# Experiments: Mixed Workload



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

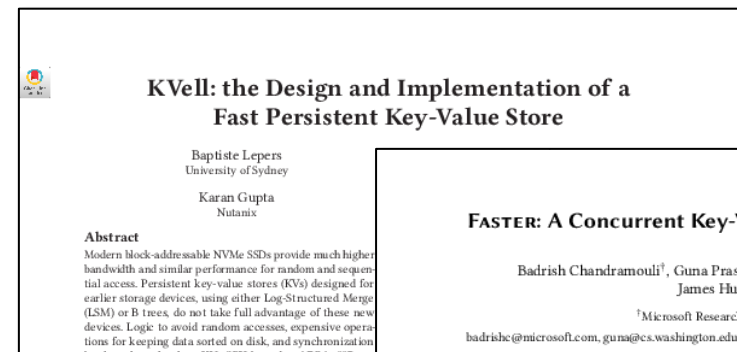- Avoid I/O by employing (Bloom) *filters*

- Judicious use of main memory
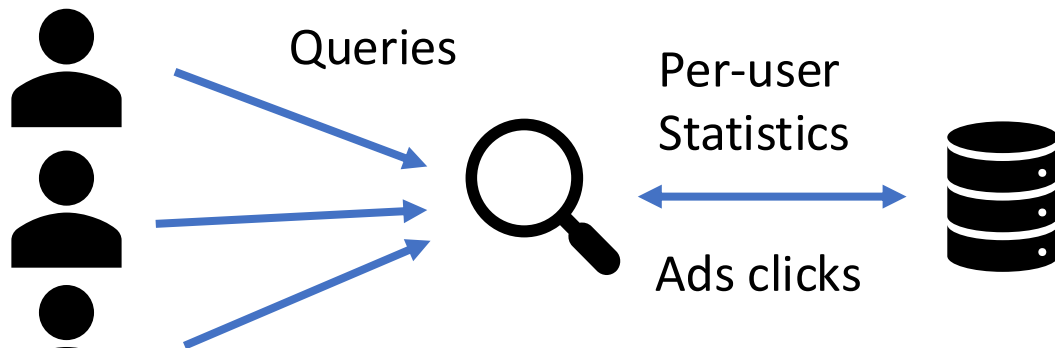
### Point-optimized Stores

- 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)



RocksDB

SPLINTERDB

KVell: the Design and Implementation of a Fast Persistent Key-Value Store

FASTER: A Concurrent Key-Value Store with In-Place Updates

BOSTON UNIVERSITY

# 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

Queries

Per-user
Statistics

Ads clicks

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

Millions of active users at any time – *critical* path
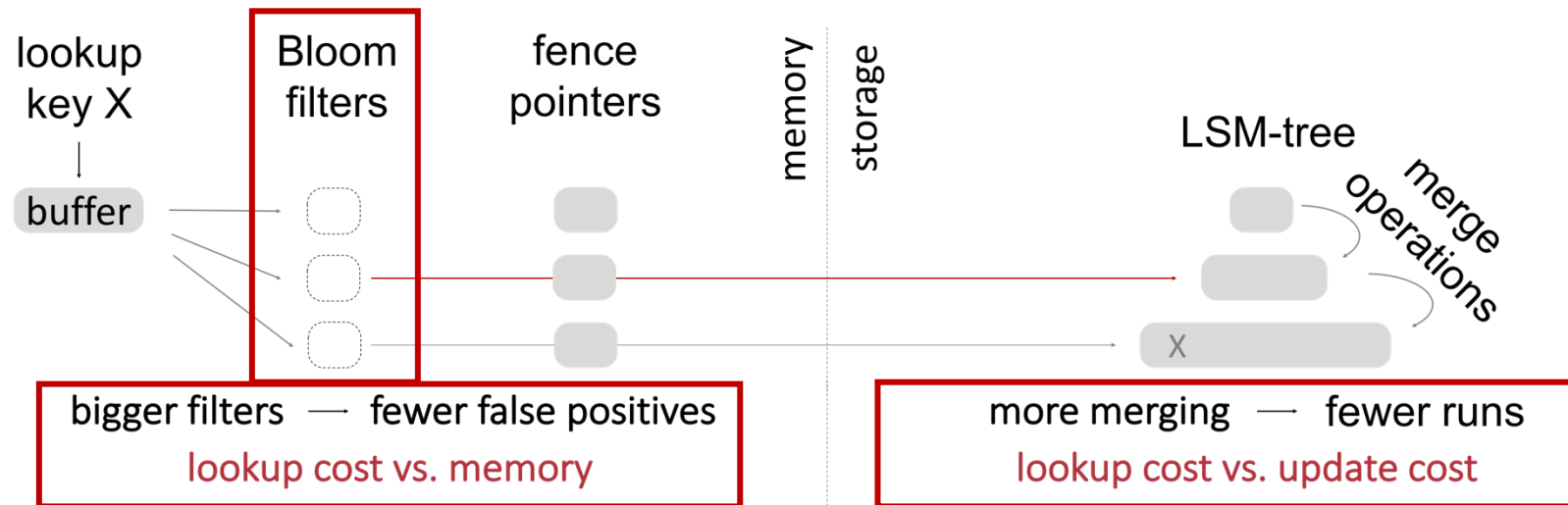
Many more *inactive,* but we still need to keep data!

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



[1] Dayan et. al. Monkey: Optimal Navigable Key-Value Store (SIGMOD'17)

# 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

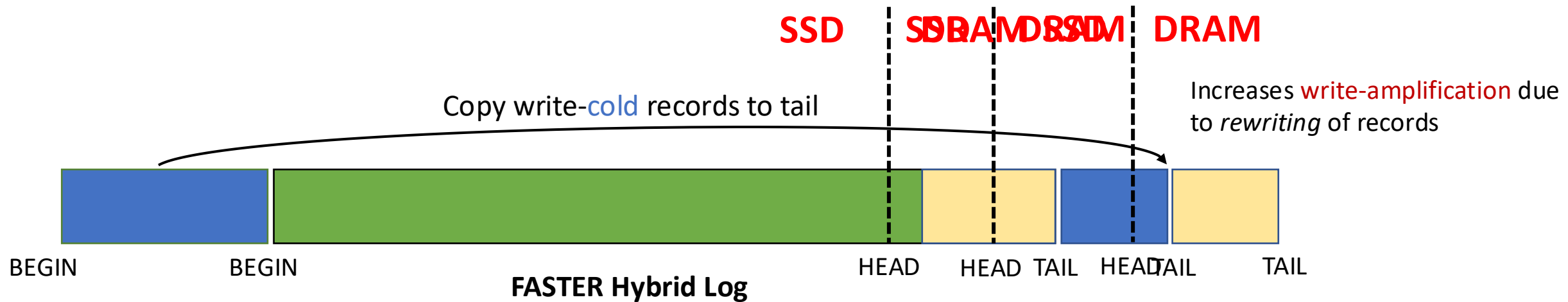– Log compaction "pollutes" in-memory log

# Limitations of Existing Systems (2)

## KVell

- Adjust page-cache size

– Large in-memory B-tree (19B per key)
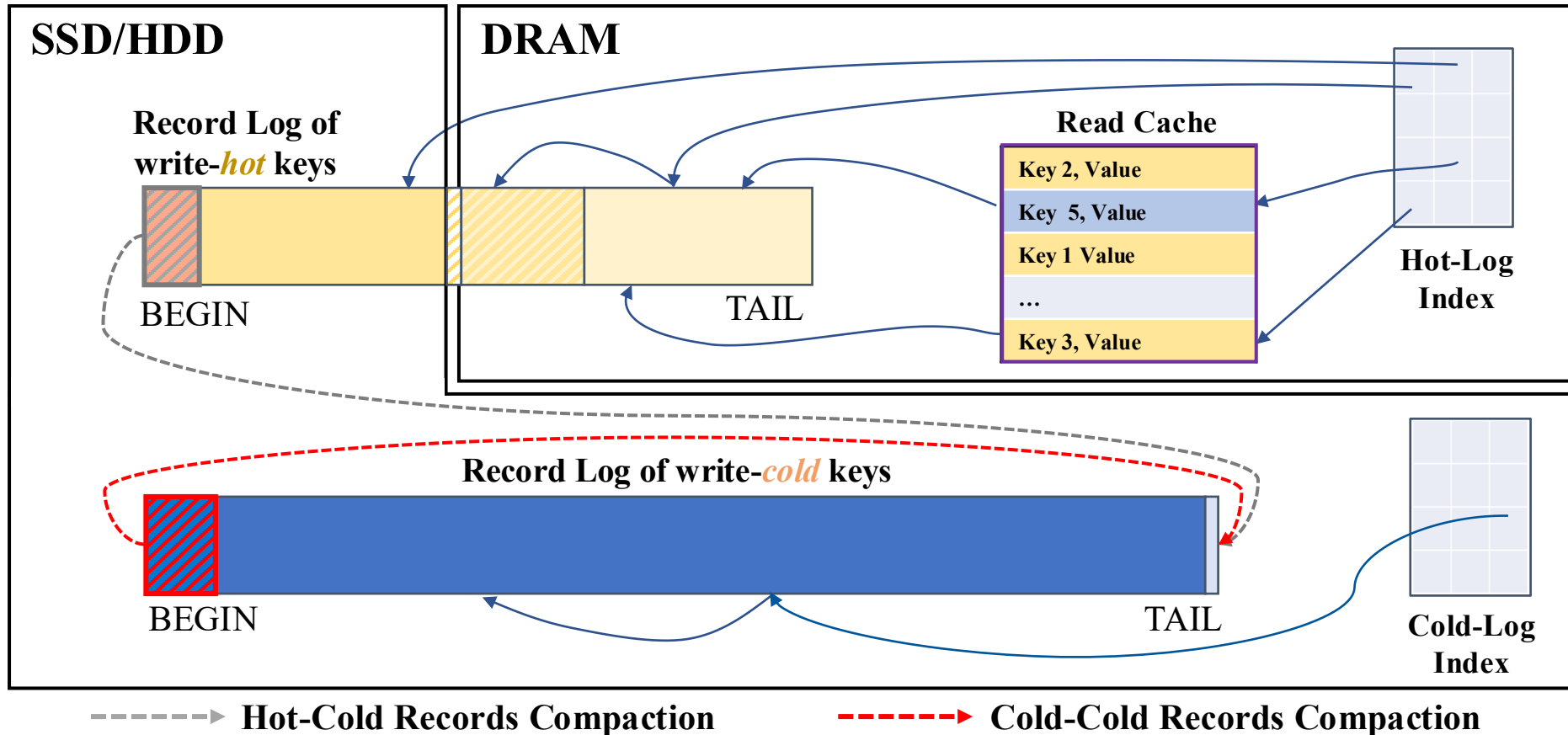
– B-tree continuously paged out to disk

## FASTER

- Tune record log in-memory size / hash index

– Reducing index size increases I/O ops

– Log compaction "pollutes" in-memory log

SSD    SSD DRAM DRAM    DRAM

Copy write-cold records to tail

Increases write-amplification due to *rewriting* of records

**FASTER Hybrid Log**

BEGIN    BEGIN    HEAD    HEAD    TAIL    HEAD TAIL    TAIL

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**              **column (PAX)**         **clustered**              **periodic**

**log-structured**        **row**                  **chained**                **piggy-backed**

**delta-main**

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)

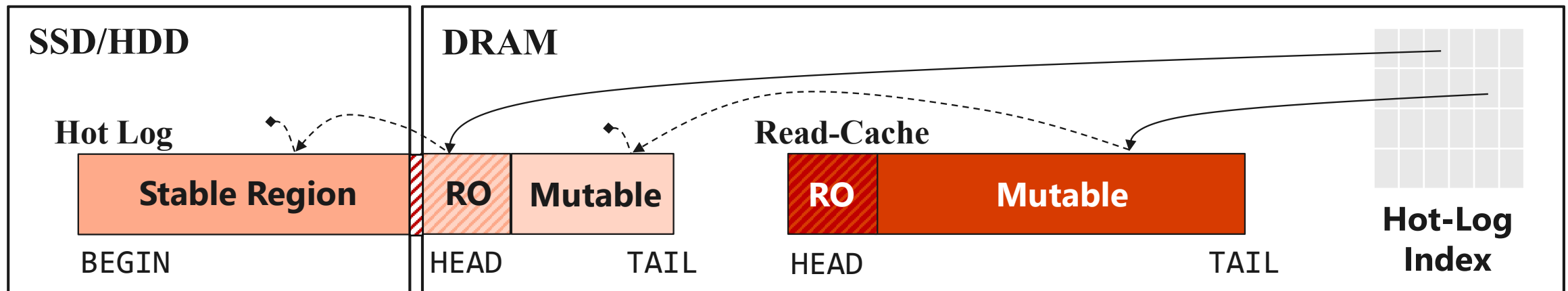**Reads** go from **hot**-log index → (optionally) read cache → **hot** log

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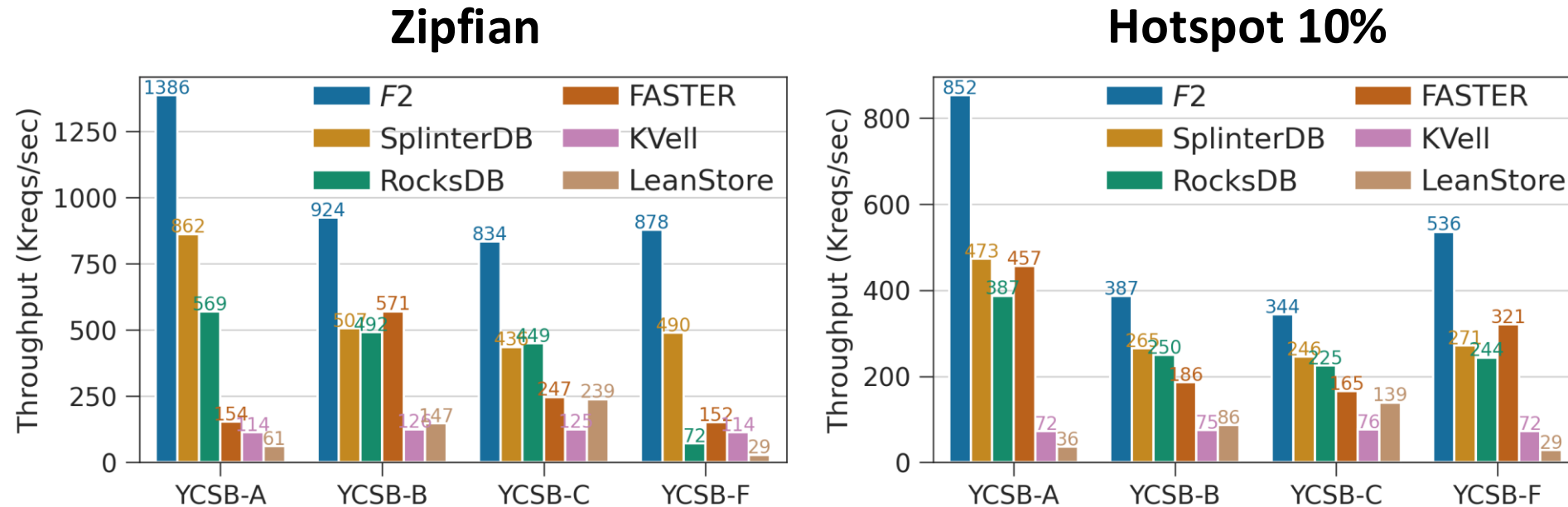  - If record with same key in read cache, it's **invalidated**!

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

# F2 – Performance Comparison

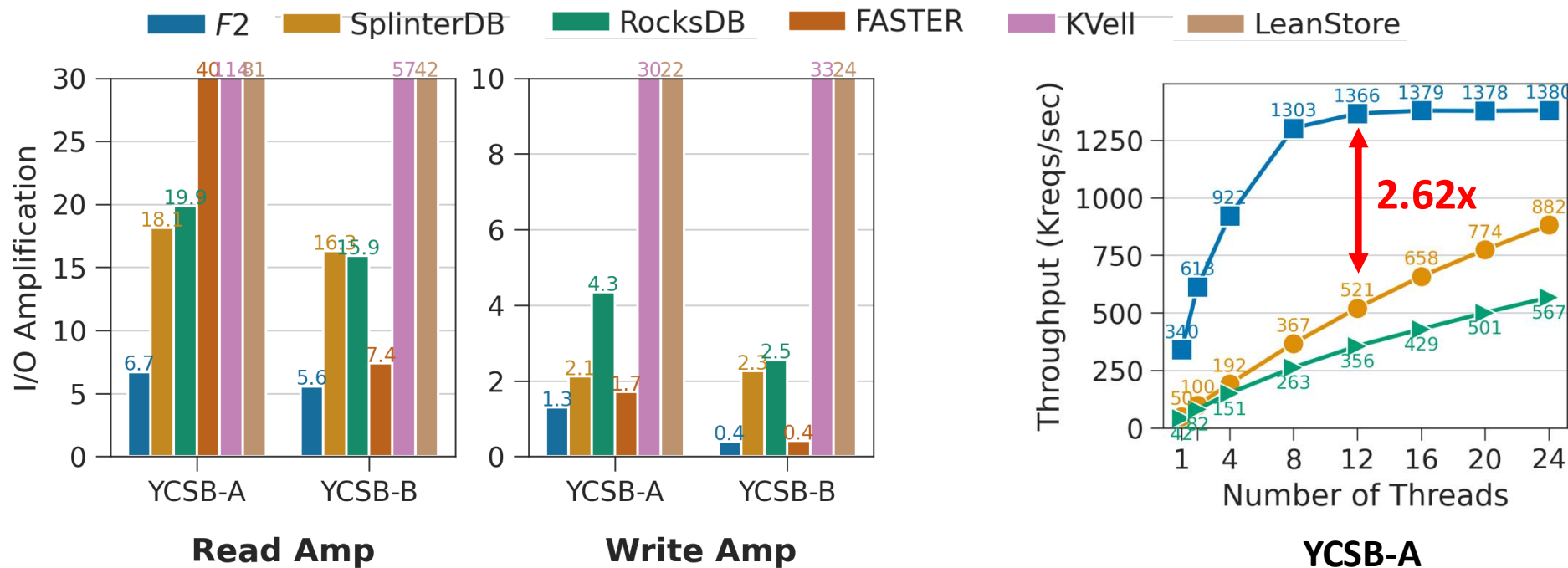YCSB: 250M keys, 8B keys, 100B values; **3GiB** mem budget (**10%** of dataset size); 24 threads, NVMe SSD



**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



**Read Amp**

**Write Amp**

**YCSB-A**

**F2 achieve less WA than Baselines**

**F2 saturates at 8-12 threads**

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

CS 561: **Data Systems Architectures**

class 7

Design Tradeoffs in Key-Value Stores

Prof. Manos Athanassoulis

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