

# COMP 421: Files & Databases

## L19: **M**ulti-**V**ersion **C**oncurrency **C**ontrol

# Multi-Version Concurrency Control

The DBMS maintains multiple **physical** versions of a single **logical** object in the database:

- When a txn writes to an object, the DBMS creates a new version of that object.
- When a txn reads an object, it reads the newest version that existed when the txn started.

# MVCC History

Protocol was first proposed in 1978 MIT PhD dissertation.

First implementations was Rdb/VMS and InterBase at DEC in early 1980s.

- Both were by Jim Starkey, co-founder of NuoDB.
- DEC Rdb/VMS is now “Oracle Rdb”.
- InterBase was open-sourced as Firebird.



**Rdb/VMS**



# Multi-Version Concurrency Control

**Writers do not block readers.**

**Readers do not block writers.**

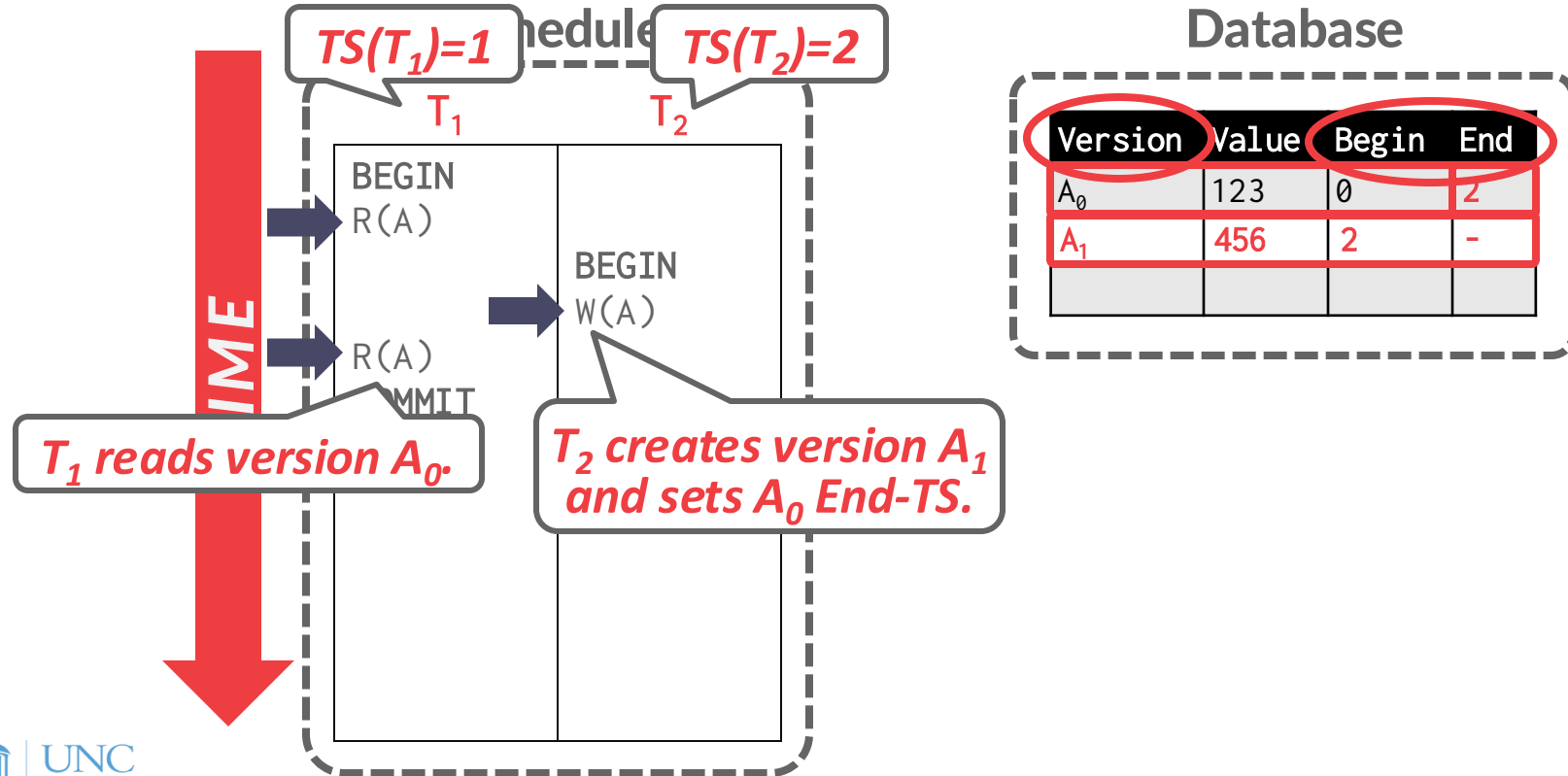
Read-only txns can read a consistent snapshot without acquiring locks.

→ Use timestamps to determine visibility.

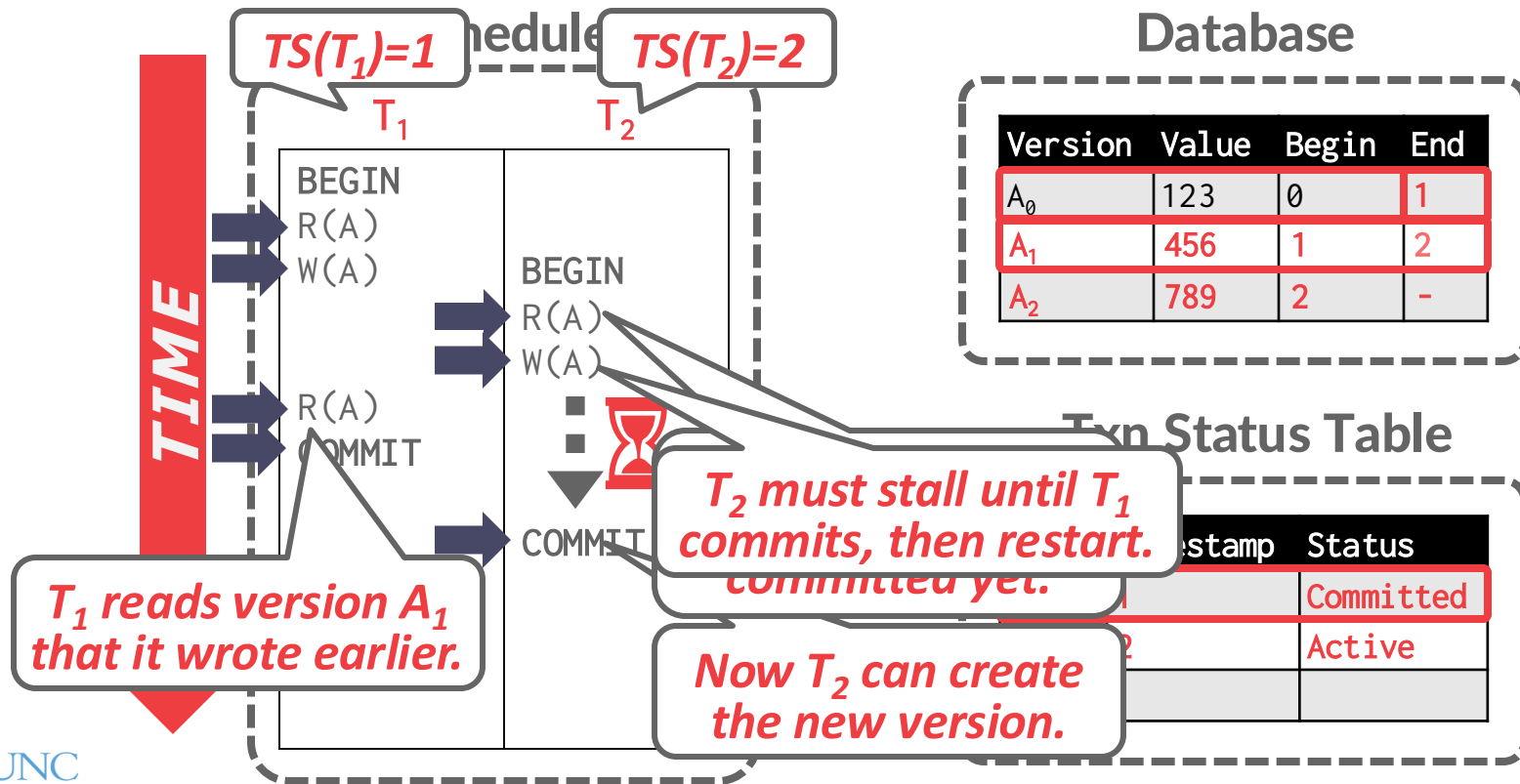
→ MVCC naturally supports Snapshot Isolation (SI).

Easily support time-travel queries.

# MVCC – Example #1



# MVCC – Example #2



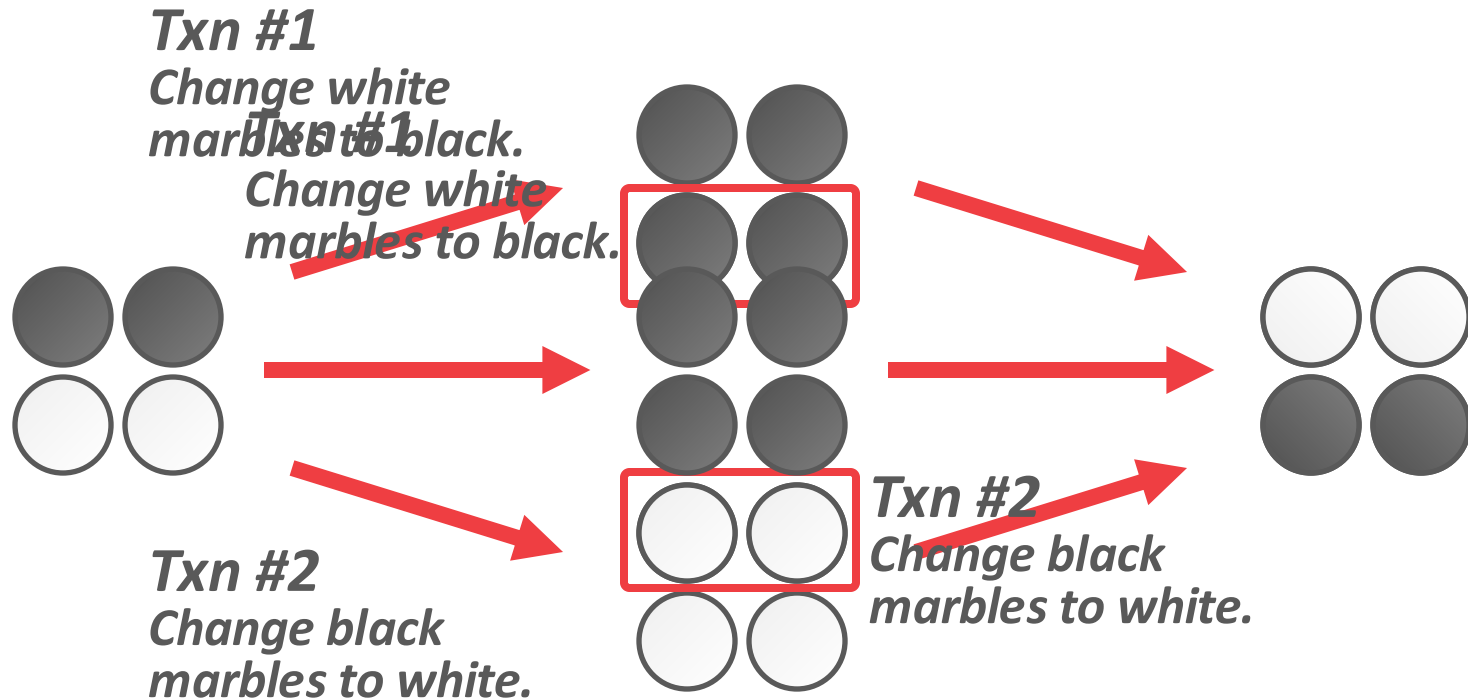
# Snapshot Isolation (SI)

When a txn starts, it sees a consistent snapshot of the database that existed when that the txn started.

- No torn writes from active txns.
- If two txns update the same object, then first writer wins.

SI is susceptible to the **Write Skew Anomaly**.

# Write Skew Anomaly





# Multi-Version Concurrency Control

MVCC is more than just a concurrency control protocol. It completely affects how the DBMS manages transactions and the database.



# MVCC Design Decisions

Concurrency Control Protocol

Version Storage

Garbage Collection

Index Management

Deletes

# Concurrency Control Protocol

## Approach #1: Timestamp Ordering

- Assign txns timestamps that determine serial order.
- Gives Snapshot Isolation

## Approach #2: Optimistic Concurrency Control

- Three-phase protocol from last class.
- Use private workspace for new versions.

## Approach #3: Two-Phase Locking

- Txns acquire appropriate lock on physical version before they can read/write a logical tuple.

# Version Storage

The DBMS uses the tuples' pointer field to create a **version chain** per logical tuple.

- This allows the DBMS to find the version that is visible to a particular txn at runtime.
- Indexes always point to the “head” of the chain.

Different storage schemes determine where/what to store for each version.

# Version Storage

## **Approach #1: Append-Only Storage**

→ New versions are appended to the same table space.

## **Approach #2: Time-Travel Storage**

→ Old versions are copied to separate table space.

## **Approach #3: Delta Storage**

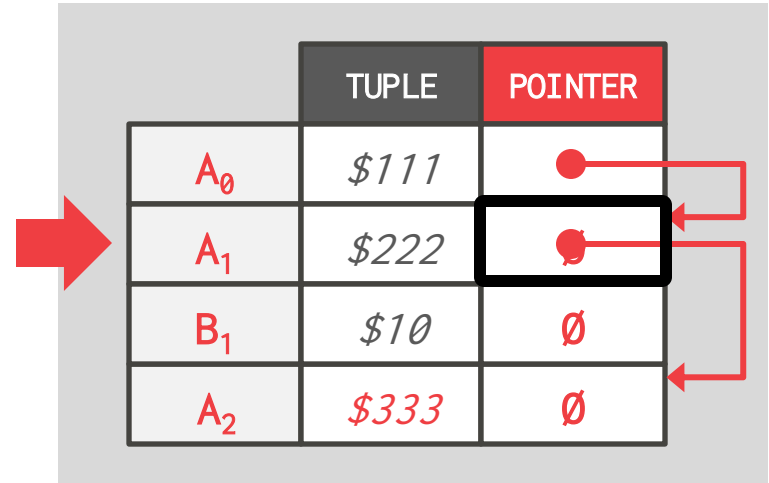
→ The original values of the modified attributes are copied into a separate delta record space.

# Append-Only Storage

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

*Main Table*



The diagram illustrates the Main Table structure for Append-Only Storage. It consists of a table with two columns: TUPLE and POINTER. The table contains four rows of data. A large red arrow points to the table from the left. Red arrows indicate the flow of updates: from A<sub>0</sub> to A<sub>1</sub>, and from A<sub>1</sub> to A<sub>2</sub>. The A<sub>1</sub> row's pointer cell is highlighted with a thick black border.

	TUPLE	POINTER
A <sub>0</sub>	\$111	●
A <sub>1</sub>	\$222	●
B <sub>1</sub>	\$10	∅
A <sub>2</sub>	\$333	∅

# Version Chain Ordering

## **Approach #1: Oldest-to-Newest (O2N)**


- Append new version to end of the chain.
- Must traverse chain on look-ups.

## **Approach #2: Newest-to-Oldest (N2O)**

- Must update index pointers for every new version.
- Do not have to traverse chain on look-ups.

# Time-Travel Storage


*Main Table*



TUPLE	POINTER
A <sub>3</sub>	\$333
B <sub>1</sub>	\$10

Overwrite primary version in the main table and update pointers.

*Time-Travel Table*




TUPLE	POINTER
A <sub>1</sub>	\$111
A <sub>2</sub>	\$222

On every update, copy the current version to the time-travel table. Update pointers.



# Delta Storage


*Main Table*



	VALUE	ID	POINTER
A <sub>3</sub>	\$333		●
B <sub>1</sub>	\$10		

On every update, copy only the column values that were modified to the delta storage and overwrite the primary version.

*Delta Storage Segment*



	DELTA	POINTER
A <sub>1</sub>	(VALUE→\$111)	∅
A <sub>2</sub>	(VALUE→\$222)	●

Txns can recreate old versions by applying the delta in reverse order.

# Garbage Collection

The DBMS needs to remove reclaimable physical versions from the database over time.

- No active txn in the DBMS can “see” that version (SI).
- The version was created by an aborted txn.

Two additional design decisions:

- How to look for expired versions?
- How to decide when it is safe to reclaim memory?

# Garbage Collection

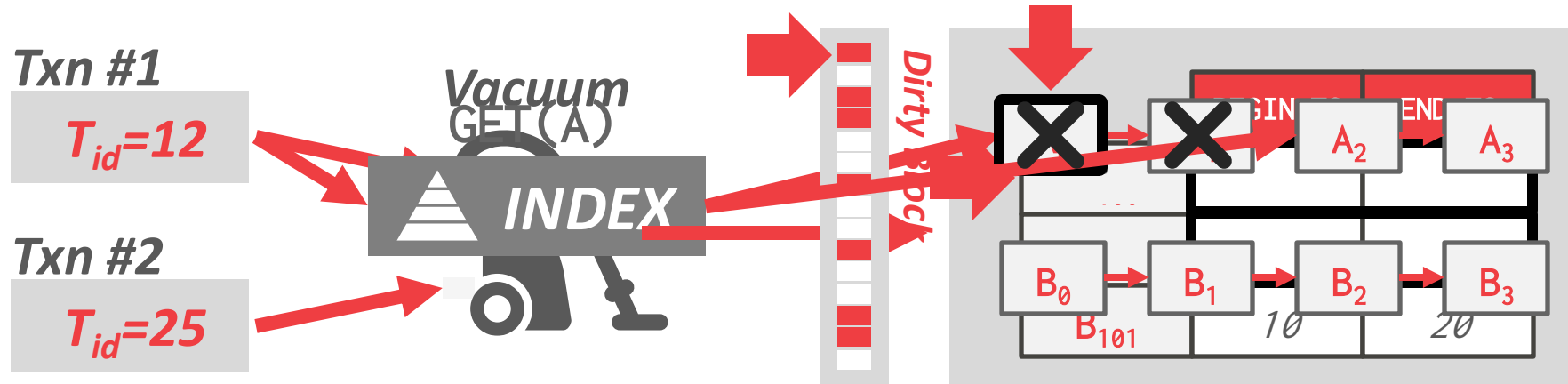
## Approach #1: Tuple-level

- Find old versions by examining tuples directly.
- Background Vacuuming vs. Cooperative Cleaning

## Approach #2: Transaction-level

- Txns keep track of their old versions so the DBMS does not have to scan tuples to determine visibility.

# Tuple-Level GC



**Background Vacuuming:**  
 Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

**Cooperative Cleaning:**  
 Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

# Transaction-Level GC

Each txn keeps track of its read/write set.

On commit/abort, the txn provides this information to a centralized vacuum worker.

The DBMS periodically determines when all versions created by a finished txn are no longer visible.

# Transaction-Level GC

**Txn #1**

BEGIN @ 10  
COMMIT @ 15

**Old Versions**



	BEGIN-TS	END-TS	DATA
$A_2$	1	10	-
$B_6$	8	10	-
$A_3$	10	$\infty$	-
$B_7$	10	$\infty$	-

**Vacuum**



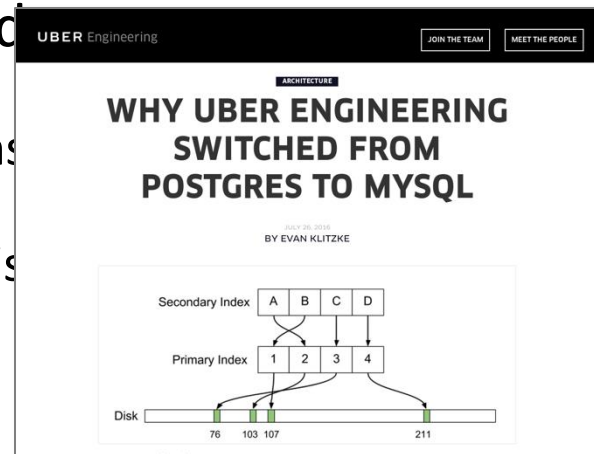
$TS < 10$

# Index Management

Primary key indexes point to version chain head

- How often the DBMS must update the pkey index depends on whether the system creates new versions when a tuple is updated.
- If a txn updates a tuple's pkey attribute(s), then this is treated as a **DELETE** followed by an **INSERT**.

Secondary indexes are more complicated...



# Secondary Indexes

## **Approach #1: Logical Pointers**

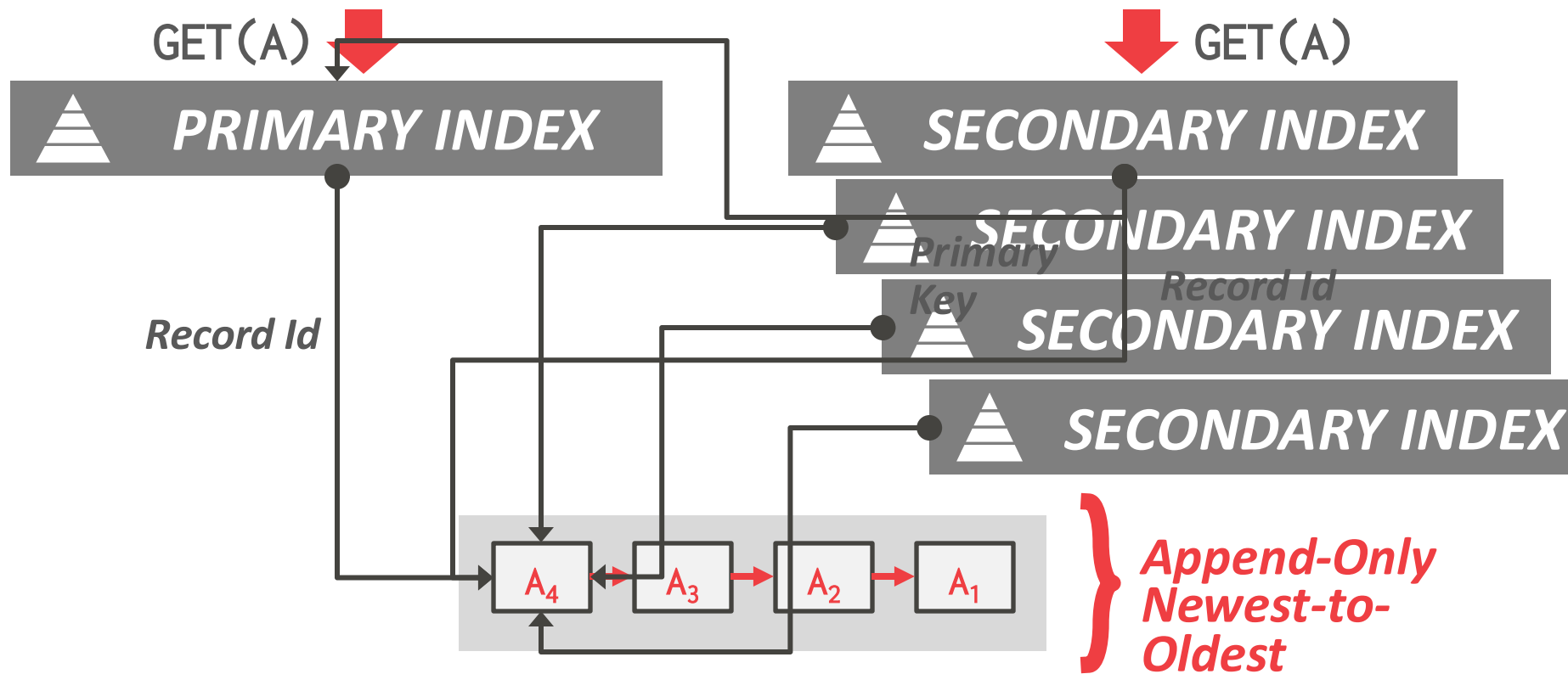
- Use a fixed identifier per tuple that does not change.
- Requires an extra indirection layer.
- Primary Key vs. Tuple Id

## **Approach #2: Physical Pointers**

- Use the physical address to the version chain head.



# Index Pointers



# MVCC Indexes

MVCC DBMS indexes (usually) do not store version information about tuples with their keys.  
→ Exception: Index-organized tables (e.g., MySQL)

Every index must support duplicate keys from different snapshots:  
→ The same key may point to different logical tuples in different snapshots.

# MVCC Duplicate Key Problem

**Txn #1**

BEGIN @ 10



READ(A)



READ(A)

**Txn #2**

BEGIN @ 20

COMMIT @ 25



UPDATE(A)



DELETE(A)

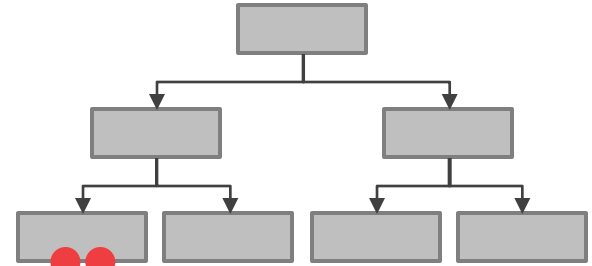
**Txn #3**

BEGIN @ 30



INSERT(A)

**Index**



	BEGIN-TS	END-TS	POINTER
$A_1$	1	20	
	20	20	$\emptyset$
$A_3$	30	$\infty$	$\emptyset$

# MVCC Indexes

Each index's underlying data structure must support the storage of non-unique keys.

Use additional execution logic to perform conditional inserts for pkey / unique indexes.

→ Atomically check whether the key exists and then insert.

Workers may get back multiple entries for a single fetch. They then must follow the pointers to find the proper physical version.

# MVCC Deletes

The DBMS physically deletes a tuple from the database only when all versions of a logically deleted tuple are not visible.

- If a tuple is deleted, then there cannot be a new version of that tuple after the newest version.
- No write-write conflicts / first-writer wins

We need a way to denote that tuple has been logically delete at some point in time.

# MVCC Deletes

## Approach #1: Deleted Flag

- Maintain a flag to indicate that the logical tuple has been deleted after the newest physical version.
- Can either be in tuple header or a separate column.

## Approach #2: Tombstone Tuple

- Create an empty physical version to indicate that a logical tuple is deleted.
- Use a separate pool for tombstone tuples with only a special bit pattern in version chain pointer to reduce the storage overhead.

# MVCC Implementations

	<i>Protocol</i>	<i>Version Storage</i>	<i>Garbage Collection</i>	<i>Indexes</i>
Oracle	MV2PL	Delta	Vacuum	Logical
Postgres	MV-2PL/MV-TO	Append-Only	Vacuum	Physical
MySQL-InnoDB	MV-2PL	Delta	Vacuum	Logical
HYRISE	MV-OCC	Append-Only	–	Physical
Hekaton	MV-OCC	Append-Only	Cooperative	Physical
MemSQL (2015)	MV-OCC	Append-Only	Vacuum	Physical
SAP HANA	MV-2PL	Time-travel	Hybrid	Logical
NuoDB	MV-2PL	Append-Only	Vacuum	Logical
HyPer	MV-OCC	Delta	Txn-level	Logical
CockroachDB	MV-2PL	Delta (LSM)	Compaction	Logical

# CONCLUSION

MVCC is the widely used scheme in DBMSs.  
Even systems that do not support multi-statement txns (e.g., NoSQL) use it.



# NEXT CLASS

Logging and recovery!