Intro to Database Systems (15-445/645)

# 18 Multi-Version Concurrency Control

Carnegie Mellon University

FALL 2022

Andy Pavlo

# 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.

# MULTI-VERSION CONCURRENCY CONTROL

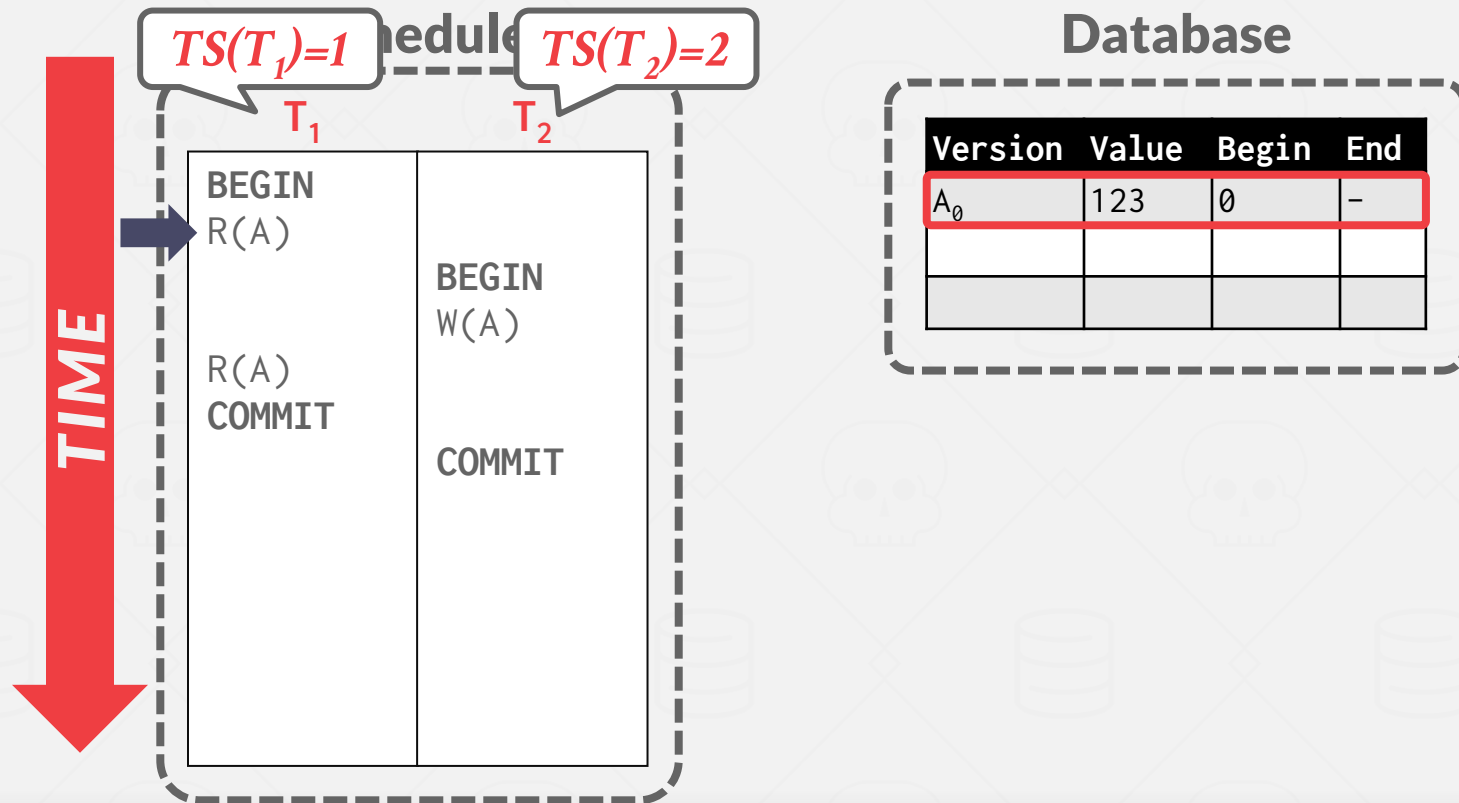**Writers do <u>not</u> block readers.**
**Readers do <u>not</u> block writers.**

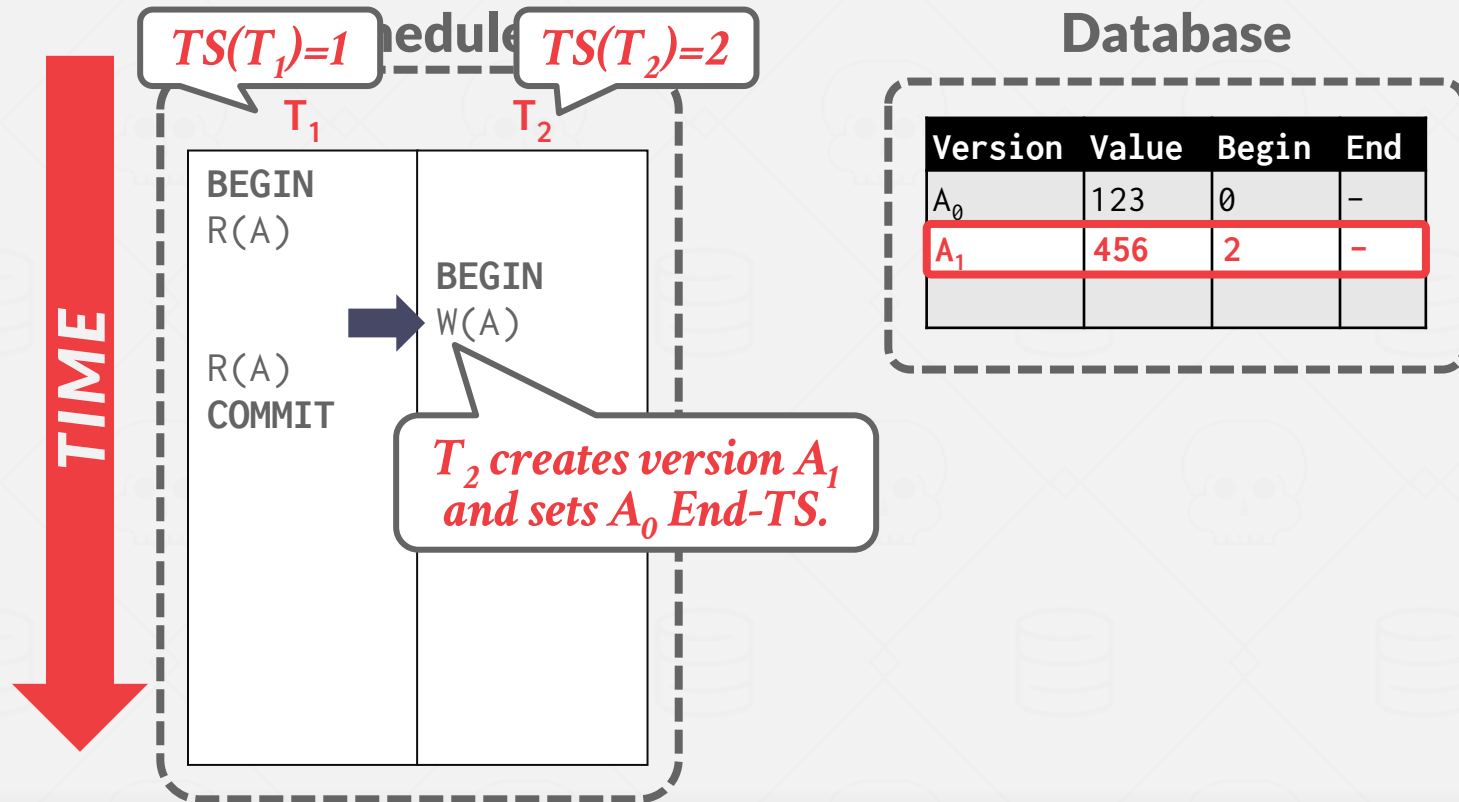Read-only txns can read a consistent <u>snapshot</u> without acquiring locks.
→ Use timestamps to determine visibility.
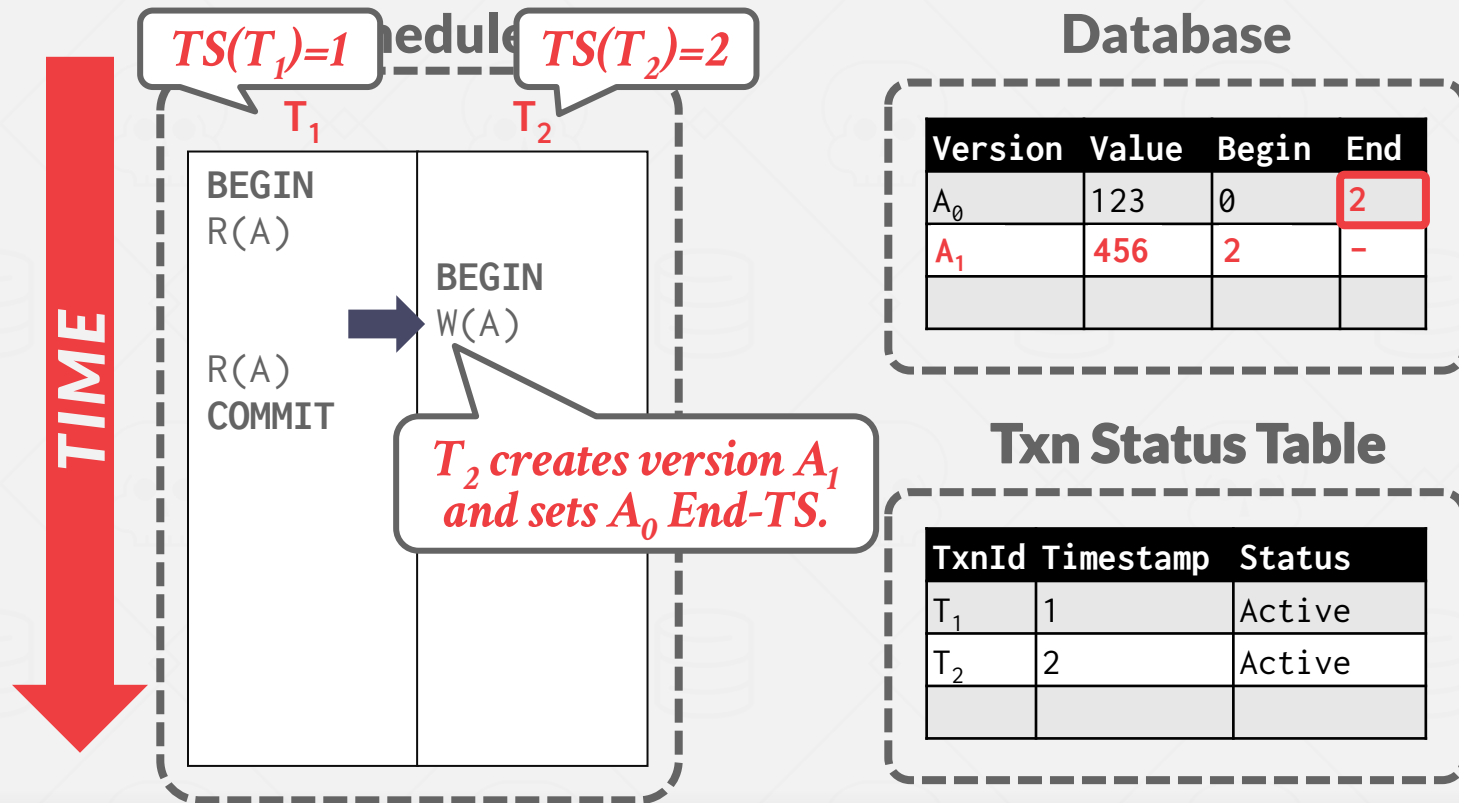
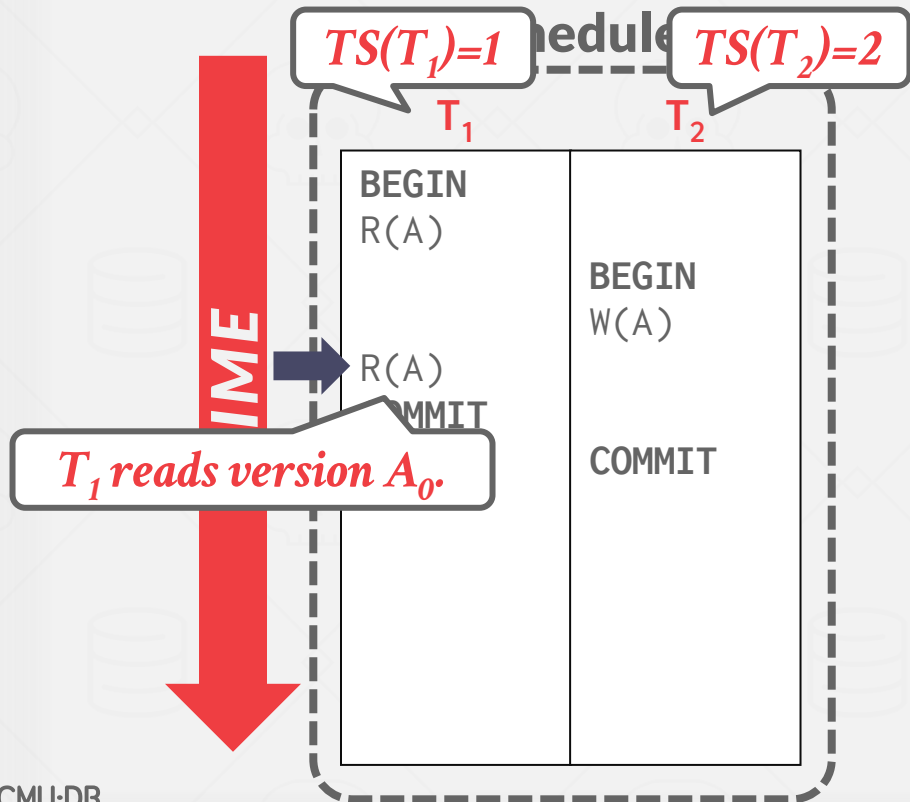Easily support <u>time-travel</u> queries.

# MVCC – EXAMPLE #1

**Schedule**

$TS(T_1)=1$

$TS(T_2)=2$

**Database**

**TIME**

| | $T_1$ | $T_2$ |
|---|---|---|
| | BEGIN | |
| →| R(A) | |
| | | BEGIN |
| | | W(A) |
| | R(A) | |
| | COMMIT | |
| | | COMMIT |

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| $A_0$ | 123 | 0 | – |
| | | | |
| | | | |

# MVCC – EXAMPLE #1

# MVCC – EXAMPLE #1



$TS(T_1)=1$   Schedule   $TS(T_2)=2$

**TIME**

| | $T_1$ | $T_2$ |
|---|---|---|
| | BEGIN | |
| | R(A) | |
| | | BEGIN |
| | | W(A) |
| | R(A) | |
| | COMMIT | |

**Database**

| Version | Value | Begin | End |
|---|---|---|---|
| $A_0$ | 123 | 0 | 2 |
| $A_1$ | 456 | 2 | – |
| | | | |

*$T_2$ creates version $A_1$ and sets $A_0$ End-TS.*

**Txn Status Table**

| TxnId | Timestamp | Status |
|---|---|---|
| $T_1$ | 1 | Active |
| $T_2$ | 2 | Active |
| | | |

# MVCC – EXAMPLE #1

Schedule

$TS(T_1)=1$

$TS(T_2)=2$

| | $T_1$ | | $T_2$ |
|---|---|---|---|
| | BEGIN | | |
| | R(A) | | |
| | | | BEGIN |
| | | | W(A) |
| | R(A) | | |
| | COMMIT | | |
| | | | COMMIT |

**TIME**

$T_1$ reads version $A_0$.

## Database

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| $A_0$ | 123 | 0 | 2 |
| $A_1$ | 456 | 2 | – |
| | | | |

## Txn Status Table

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| $T_1$ | 1 | Active |
| $T_2$ | 2 | Active |
| | | |

# MVCC – EXAMPLE #2

**TS(T₁)=1**  edule  **TS(T₂)=2**

**Database**

| | T₁ | T₂ |
|---|---|---|
| **TIME** | BEGIN<br>R(A)<br>W(A)<br><br><br>R(A)<br>COMMIT | <br><br>BEGIN<br>R(A)<br>W(A)<br><br><br><br>COMMIT |

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| $A_0$   | 123   | 0     |     |
|         |       |       |     |
|         |       |       |     |

**Txn Status Table**

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| $T_1$ | 1         | Active |
|       |           |        |
|       |           |        |

# MVCC – EXAMPLE #2

TS(T₁)=1    chedule    TS(T₂)=2

**Database**

T₁                T₂

BEGIN
R(A)
W(A)         BEGIN
             R(A)
             W(A)

R(A)
COMMIT

             COMMIT

TIME

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀ | 123 | 0 | |
| A₁ | 456 | 1 | – |
| | | | |

**Txn Status Table**

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Active |
| | | |
| | | |

# MVCC – EXAMPLE #2

# MVCC – EXAMPLE #2



**Database**

**T_n Status Table**

TIME

TS(T_1)=1
TS(T_2)=2

Schedule

| | T_1 | | T_2 |
|---|---|---|---|
| | BEGIN | | |
| | R(A) | | |
| | W(A) | | BEGIN |
| | | | R(A) |
| | | | W(A) |
| | R(A) | | |
| | COMMIT | | |
| | | | COMMIT |

**Database**

| Version | Value | Begin | End |
|---|---|---|---|
| A_0 | 123 | 0 | 1 |
| A_1 | 456 | 1 | – |
| | | | |

**Tₓₙ Status Table**

| | Timestamp | Status |
|---|---|---|
| | | Active |
| T_2 | 2 | Active |
| | | |

*T_2 reads version A_0 because T_1 has not committed yet.*

# MVCC – EXAMPLE #2

# MVCC – EXAMPLE #2

# MVCC – EXAMPLE #2

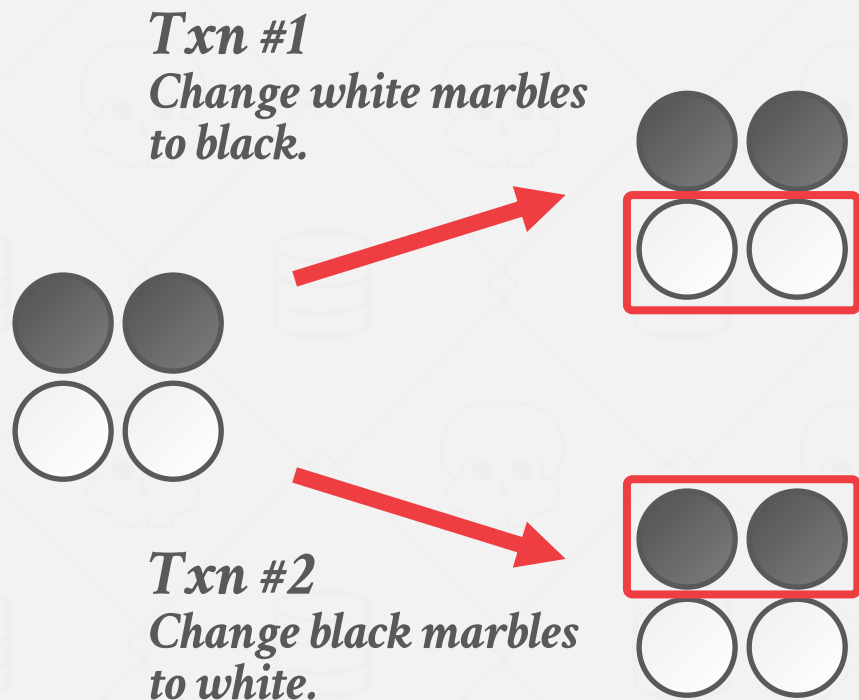# MVCC – EXAMPLE #2

# SNAPSHOT ISOLATION (SI)

When a txn starts, it sees a <u>consistent</u> snapshot of the database that existed when that the txn started.
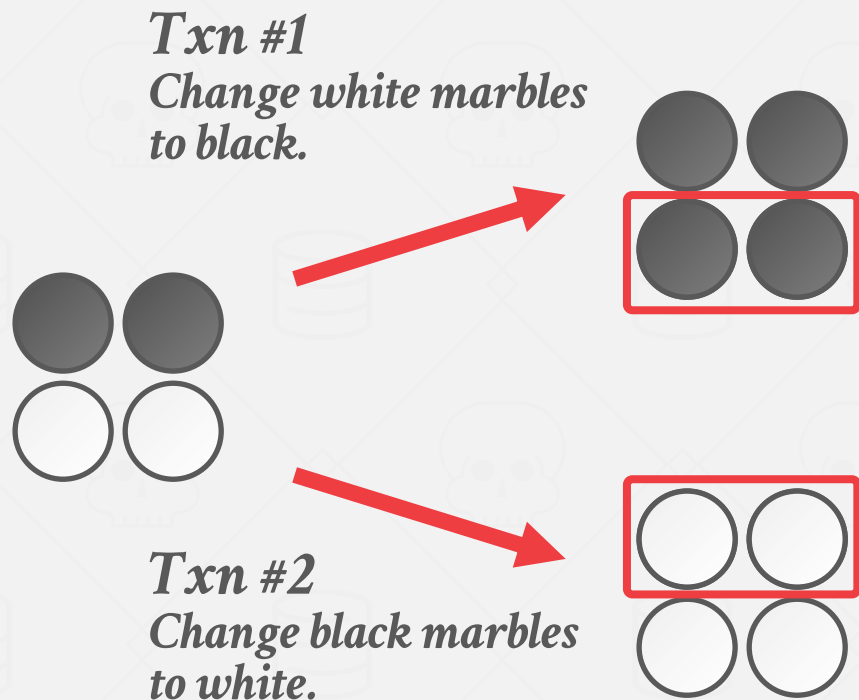→ No torn writes from active txns.
→ If two txns update the same object, then first writer wins.
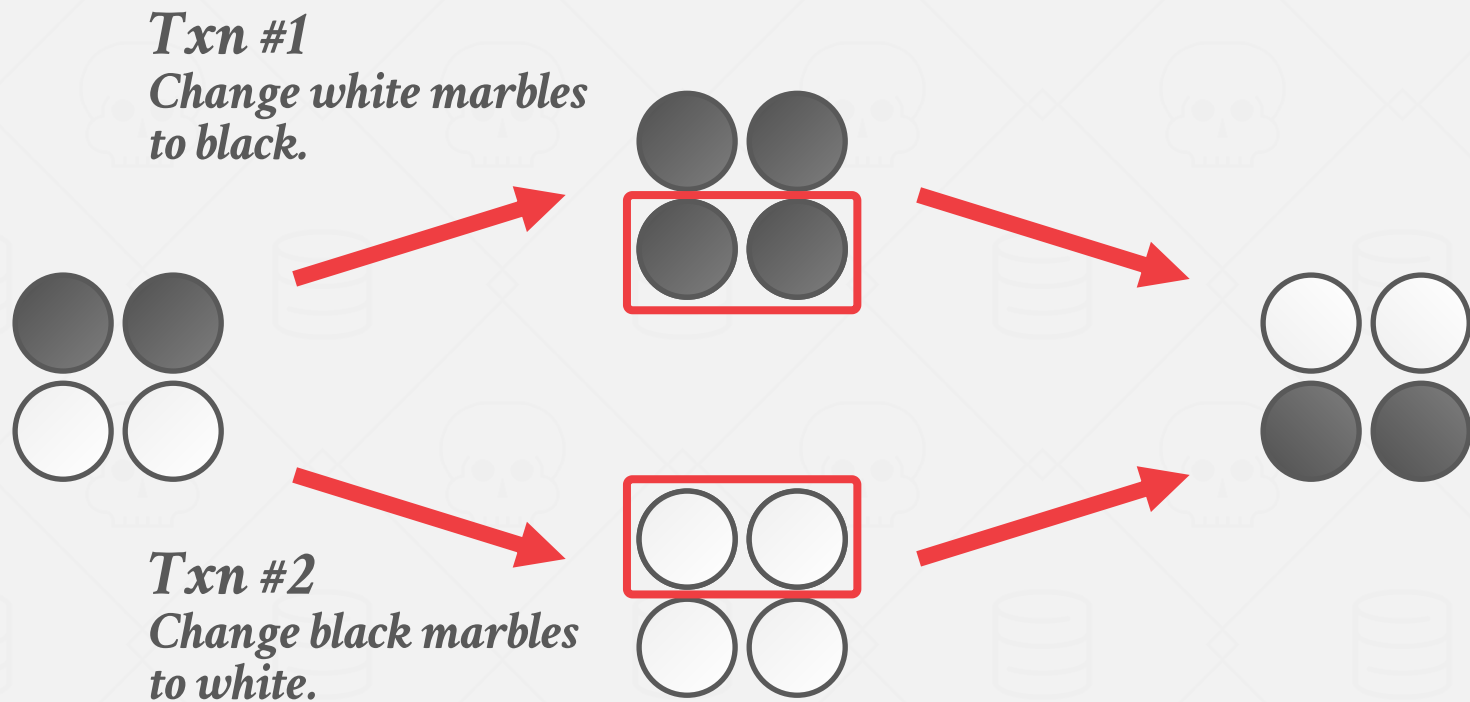
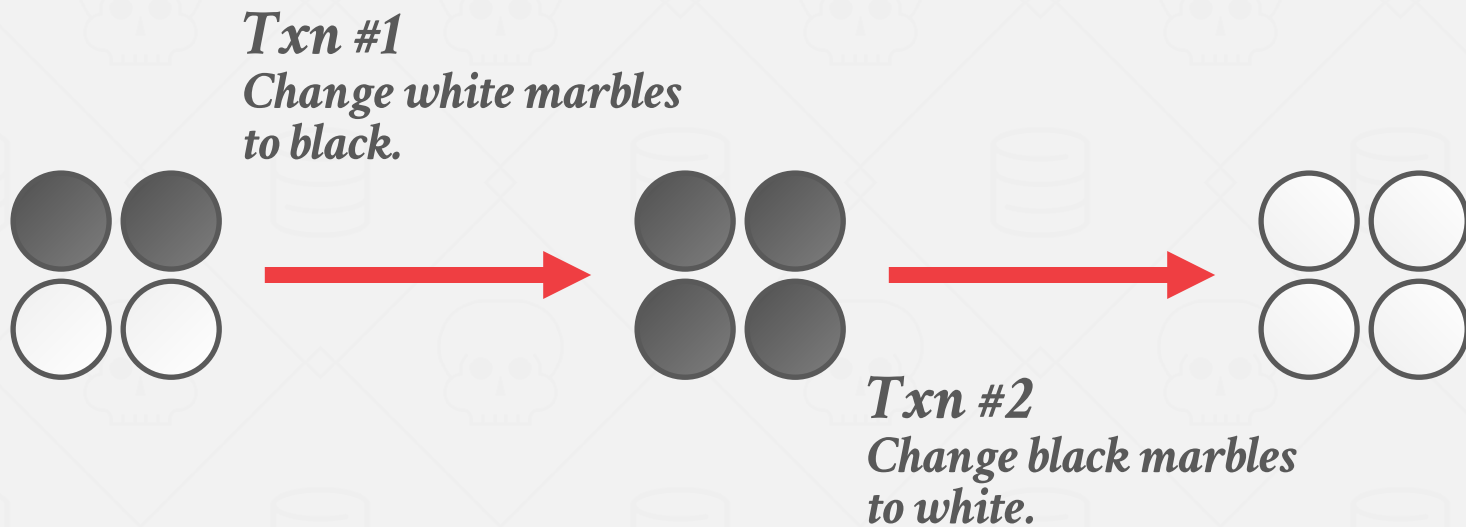SI is susceptible to the **Write Skew Anomaly**.

# WRITE SKEW ANOMALY



Txn #1
Change white marbles
to black.

Txn #2
Change black marbles
to white.

# WRITE SKEW ANOMALY

**Txn #1**
*Change white marbles to black.*

**Txn #2**
*Change black marbles to white.*

# WRITE SKEW ANOMALY

**Txn #1**
*Change white marbles to black.*

**Txn #2**
*Change black marbles to white.*

# WRITE SKEW ANOMALY

*Txn #1*
*Change white marbles to black.*

*Txn #2*
*Change black marbles to white.*

# 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.

**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.

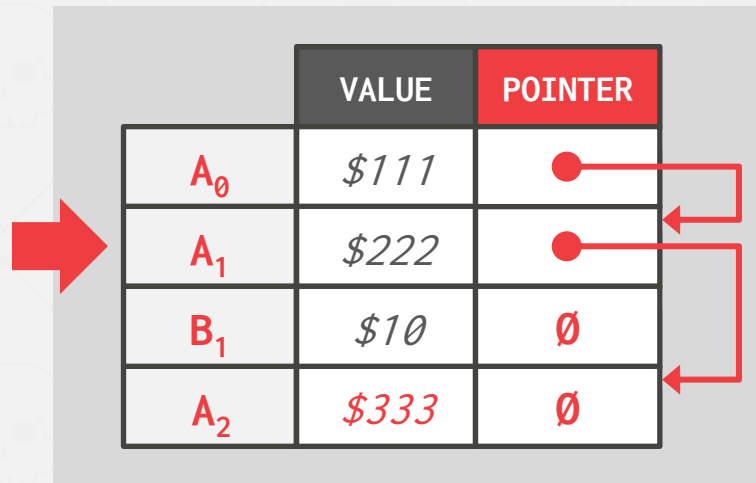On every update, append a new version of the tuple into an empty space in the table.

*Main Table*

| | VALUE | POINTER |
|---|---|---|
| $A_0$ | $111 | ● |
| $A_1$ | $222 | Ø |
| $B_1$ | $10 | Ø |
| | | |

CMU·DB

# 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*

| | VALUE | POINTER |
|---|---|---|
| $A_0$ | $111 | ● |
| $A_1$ | $222 | Ø |
| $B_1$ | $10 | Ø |
| $A_2$ | $333 | Ø |

# APPEND-ONLY STORAGE

*Main Table*

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.



| | VALUE | POINTER |
|---|---|---|
| $A_0$ | $111 | ● |
| $A_1$ | $222 | ● |
| $B_1$ | $10 | Ø |
| $A_2$ | $333 | Ø |

# 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.
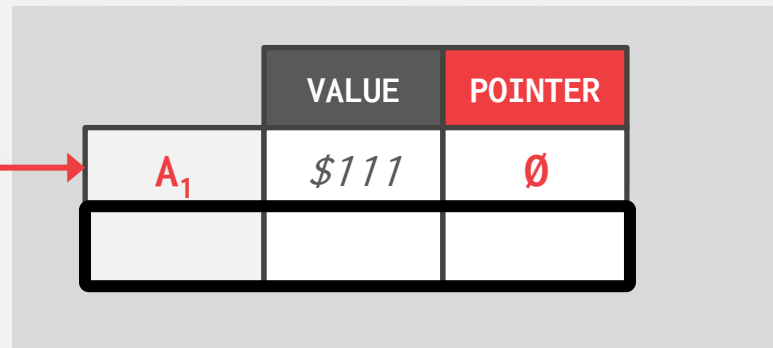→ Do not have to traverse chain on look-ups.

# TIME-TRAVEL STORAGE

*Main Table*

*Time-Travel Table*

| | VALUE | POINTER |
|---|---|---|
| A$_2$ | $222 | ● |
| B$_1$ | $10 | |

| | VALUE | POINTER |
|---|---|---|
| A$_1$ | $111 | Ø |
| | | |

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

# TIME-TRAVEL STORAGE

*Main Table*

*Time-Travel Table*

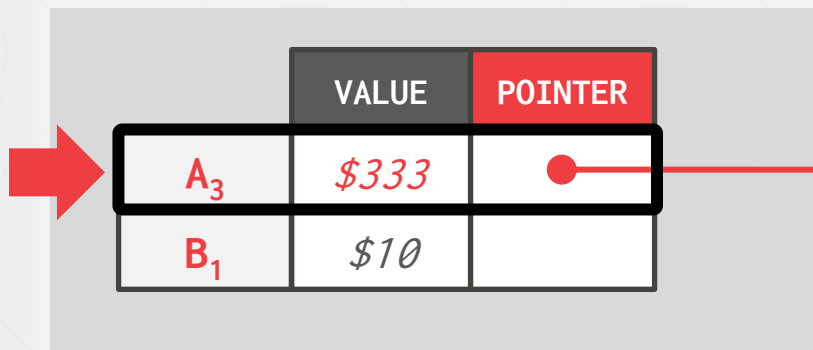| | VALUE | POINTER |
|---|---|---|
| A₂ | $222 | ● |
| B₁ | $10 | |

| | VALUE | POINTER |
|---|---|---|
| A₁ | $111 | Ø |
| A₂ | $222 | ● |

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

# TIME-TRAVEL STORAGE

*Main Table*



| | VALUE | POINTER |
|---|---|---|
| A₃ | $333 | ● |
| B₁ | $10 | |

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

*Time-Travel Table*

| | VALUE | POINTER |
|---|---|---|
| A₁ | $111 | Ø |
| A₂ | $222 | ● |

Overwrite master version in the main table and update pointers.

# TIME-TRAVEL STORAGE

**Main Table**

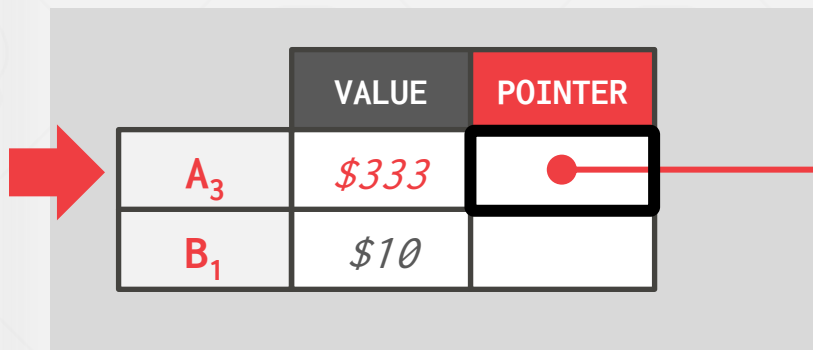| | VALUE | POINTER |
|---|---|---|
| A₃ | $333 | ● |
| B₁ | $10 | |

**Time-Travel Table**

| | VALUE | POINTER |
|---|---|---|
| A₁ | $111 | Ø |
| A₂ | $222 | ● |

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

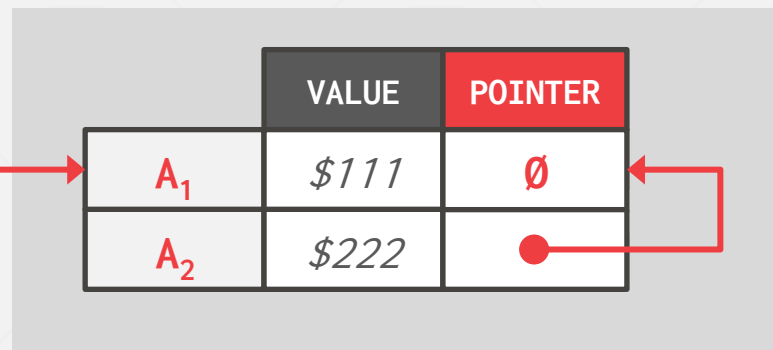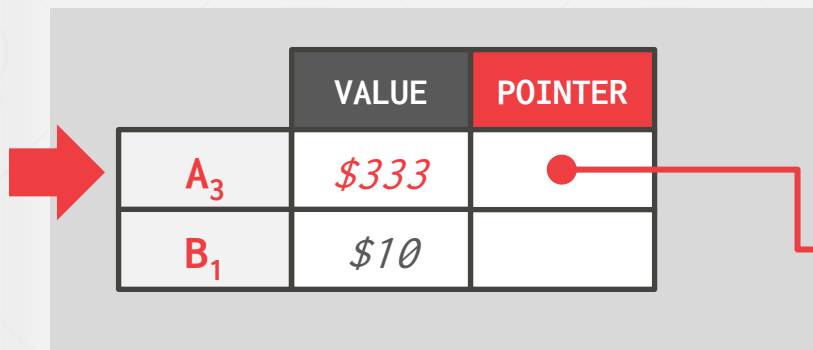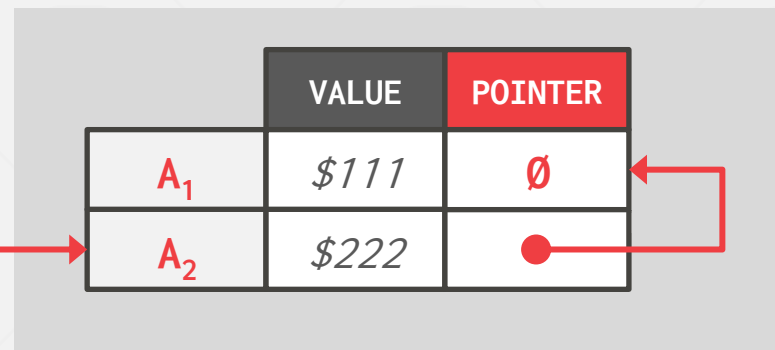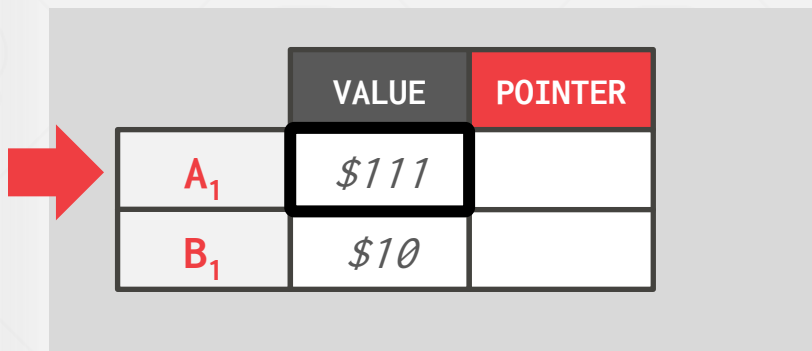Overwrite master version in the main table and update pointers.

# TIME-TRAVEL STORAGE

*Main Table*



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

*Time-Travel Table*



Overwrite master version in the main table and update pointers.

# DELTA STORAGE

*Main Table*

| | VALUE | POINTER |
|---|---|---|
| **A₁** | *$111* | |
| **B₁** | *$10* | |

*Delta Storage Segment*

| | DELTA | POINTER |
|---|---|---|
| **A₁** | *(VALUE→$111)* | *Ø* |

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

# DELTA STORAGE

*Main Table*

*Delta Storage Segment*

| | VALUE | POINTER |
|---|---|---|
| A₂ | *$222* | ● |
| B₁ | *$10* | |

| | DELTA | POINTER |
|---|---|---|
| A₁ | *(VALUE→$111)* | Ø |

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

# DELTA STORAGE

*Main Table*

*Delta Storage Segment*

| | VALUE | POINTER |
|---|---|---|
| A₂ | *$222* | ● |
| B₁ | *$10* | |

| | DELTA | POINTER |
|---|---|---|
| A₁ | *(VALUE→$111)* | Ø |
| A₂ | *(VALUE→$222)* | ● |

On every update, copy only the values that were modified to the delta storage and overwrite the master version.
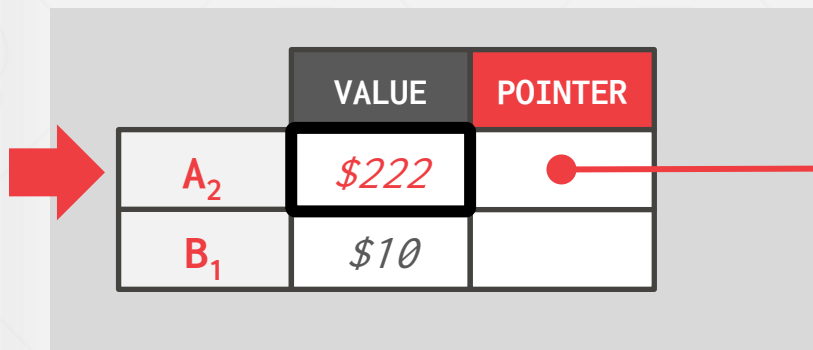
# DELTA STORAGE

*Main Table*



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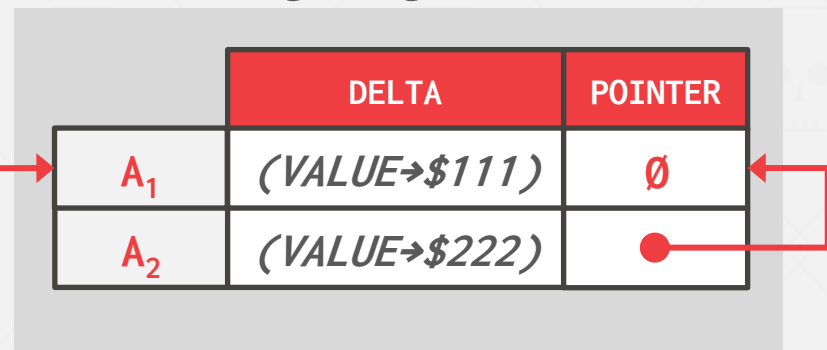to the delta storage and
overwrite the master version.

*Delta Storage Segment*

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

# GARBAGE COLLECTION

The DBMS needs to remove **<u>reclaimable</u>** 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

**Txn #1**

$T_{id}=12$

**Txn #2**

$T_{id}=25$

*Vacuum*



|  | BEGIN-TS | END-TS |
|---|---|---|
| $A_{100}$ | 1 | 9 |
| $B_{100}$ | 1 | 9 |
| $B_{101}$ | 10 | 20 |

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

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

*Vacuum*

**Txn #2**

$T_{id}=25$



|  | BEGIN-TS | END-TS |
|---|---|---|
| $A_{100}$ | 1 | 9 |
| $B_{100}$ | 1 | 9 |
| $B_{101}$ | 10 | 20 |

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

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

**Txn #2**

$T_{id}=25$

*Vacuum*

|  | BEGIN-TS | END-TS |
|---|---|---|
|  |  |  |
|  |  |  |
| $B_{101}$ | *10* | *20* |

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

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

**Txn #2**

$T_{id}=25$

*Vacuum*

*Dirty Block BitMap*

| | BEGIN-TS | END-TS |
|---|---|---|
| | | |
| | | |
| $B_{101}$ | 10 | 20 |

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

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

*Vacuum*

**Txn #2**

$T_{id}=25$

*Dirty Block BitMap*

|  | BEGIN-TS | END-TS |
|---|---|---|
|  |  |  |
|  |  |  |
| $B_{101}$ | *10* | *20* |

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

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

GET(A)


△ INDEX

$A_0 \rightarrow A_1 \rightarrow A_2 \rightarrow A_3$

**Txn #2**

$T_{id}=25$

$B_0 \rightarrow B_1 \rightarrow B_2 \rightarrow B_3$

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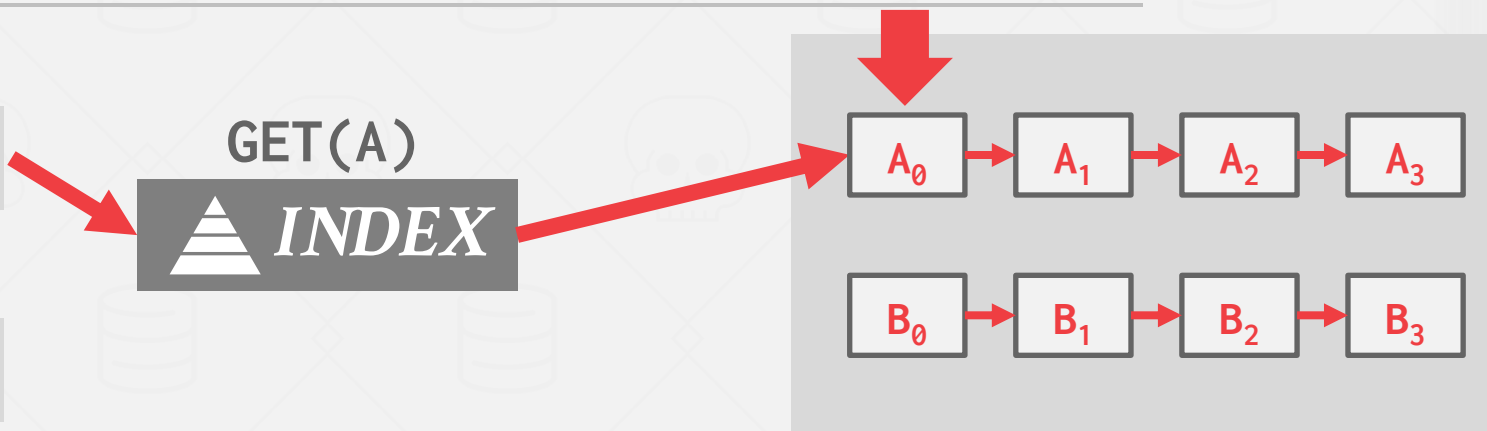Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

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

# TUPLE-LEVEL GC



**Txn #1**
$T_{id}=12$

GET(A)

**INDEX**

$A_1 \to A_2 \to A_3$

**Txn #2**
$T_{id}=25$

$B_0 \to B_1 \to B_2 \to B_3$

**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.
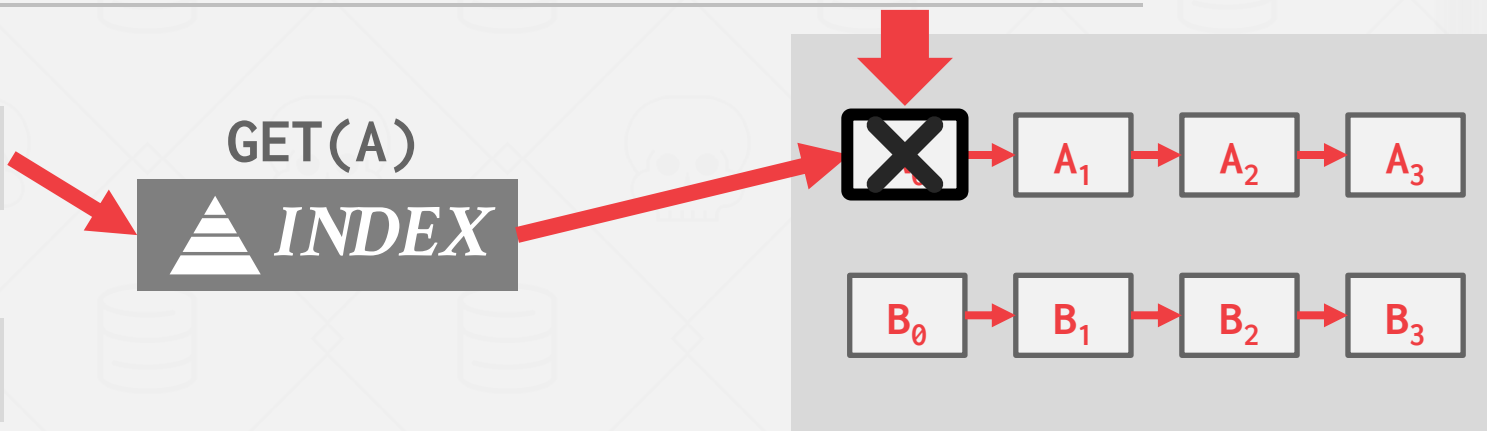
# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

GET(A)

**INDEX**

**Txn #2**

$T_{id}=25$

A₂   A₃

B₀ → B₁ → B₂ → B₃

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

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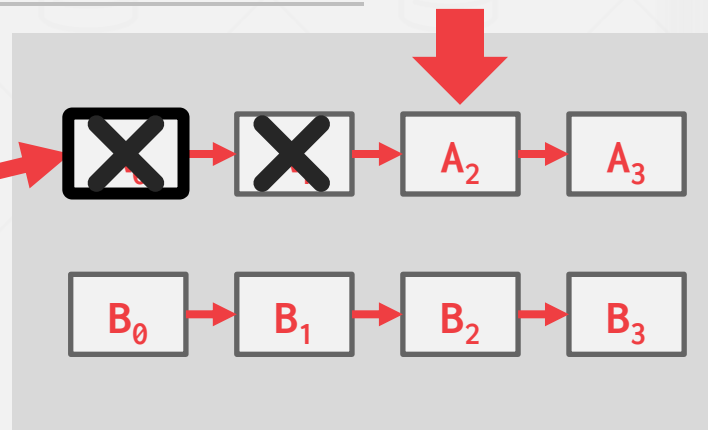Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

CMU·DB

# TUPLE-LEVEL GC

**Txn #1**

$T_{id}=12$

**Txn #2**

$T_{id}=25$

GET(A)

△ *INDEX*

$A_2 \rightarrow A_3$

$B_0 \rightarrow B_1 \rightarrow B_2 \rightarrow B_3$

**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.
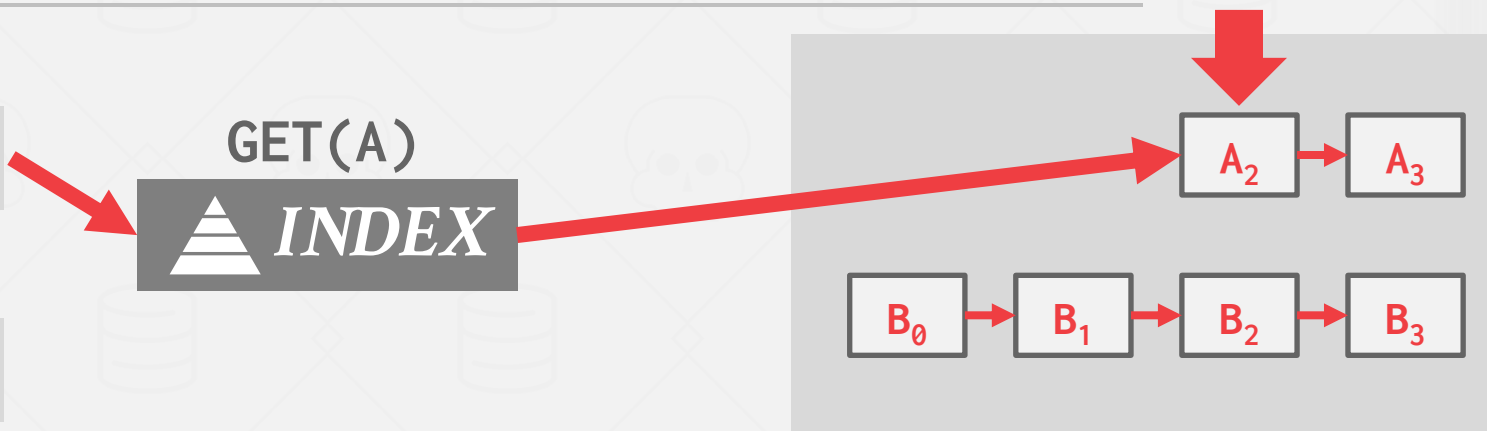
# TRANSACTION-LEVEL GC

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

UPDATE(A)

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | *1* | ∞ | – |
| $B_6$ | *8* | ∞ | – |
| | | | |
| | | | |

# TRANSACTION-LEVEL GC

*Txn #1*

BEGIN @ 10

UPDATE(A)

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | ∞ | – |
| $A_3$ | 10 | ∞ | – |
| | | | |

# TRANSACTION-LEVEL GC

*Txn #1*

**BEGIN @ 10**

*Old Versions*

$A_2$

UPDATE(A)

|  | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | *1* | *10* | – |
| $B_6$ | *8* | ∞ | – |
| $A_3$ | *10* | ∞ | – |
|  |  |  |  |

# TRANSACTION-LEVEL GC

*Txn #1*

**BEGIN @ 10**

*Old Versions*

$A_2$

UPDATE(A)

UPDATE(B)

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | ∞ | – |
| $A_3$ | 10 | ∞ | – |
| | | | |

# TRANSACTION-LEVEL GC

*Txn #1*

**BEGIN @ 10**

*Old Versions*

$A_2$

UPDATE(A)

UPDATE(B)

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | 10 | – |
| $A_3$ | 10 | ∞ | – |
| $B_7$ | 10 | ∞ | – |

# TRANSACTION-LEVEL GC

*Txn #1*

**BEGIN @ 10**

UPDATE(A)

UPDATE(B)

*Old Versions*

$A_2$

$B_6$

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | 10 | – |
| $A_3$ | 10 | ∞ | – |
| $B_7$ | 10 | ∞ | – |

# TRANSACTION-LEVEL GC

*Txn #1*

BEGIN @ 10
COMMIT @ 15

**Old Versions**

$A_2$

$B_6$

UPDATE(A)

UPDATE(B)

|  | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | 10 | – |
| $A_3$ | 10 | ∞ | – |
| $B_7$ | 10 | ∞ | – |

# TRANSACTION-LEVEL GC

*Txn #1*

BEGIN @ 10
COMMIT @ 15

*Old Versions*

UPDATE(A)

UPDATE(B)

|  | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| $A_2$ | 1 | 10 | – |
| $B_6$ | 8 | 10 | – |
| $A_3$ | 10 | ∞ | – |
| $B_7$ | 10 | ∞ | – |

*Vacuum*

$TS<10$ { $A_2$ $B_6$

# 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**.
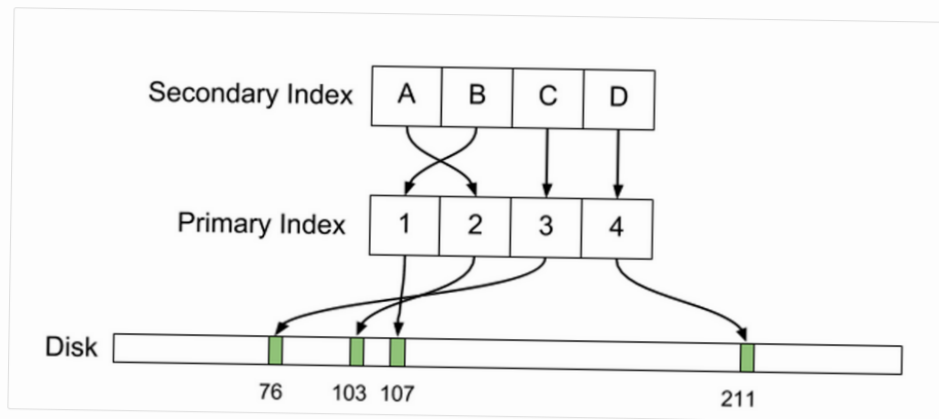
Secondary indexes are more complicated…

# 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

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

*Record Id*

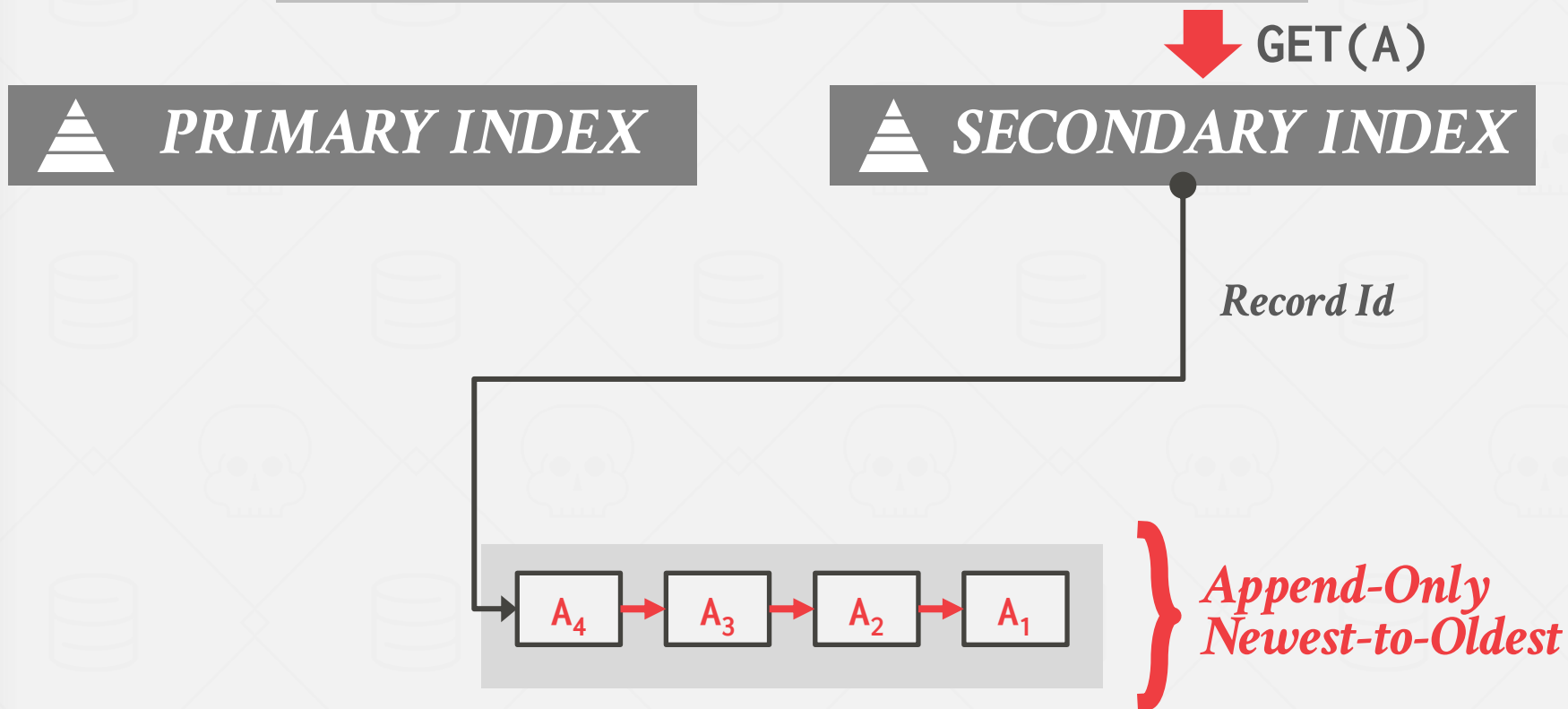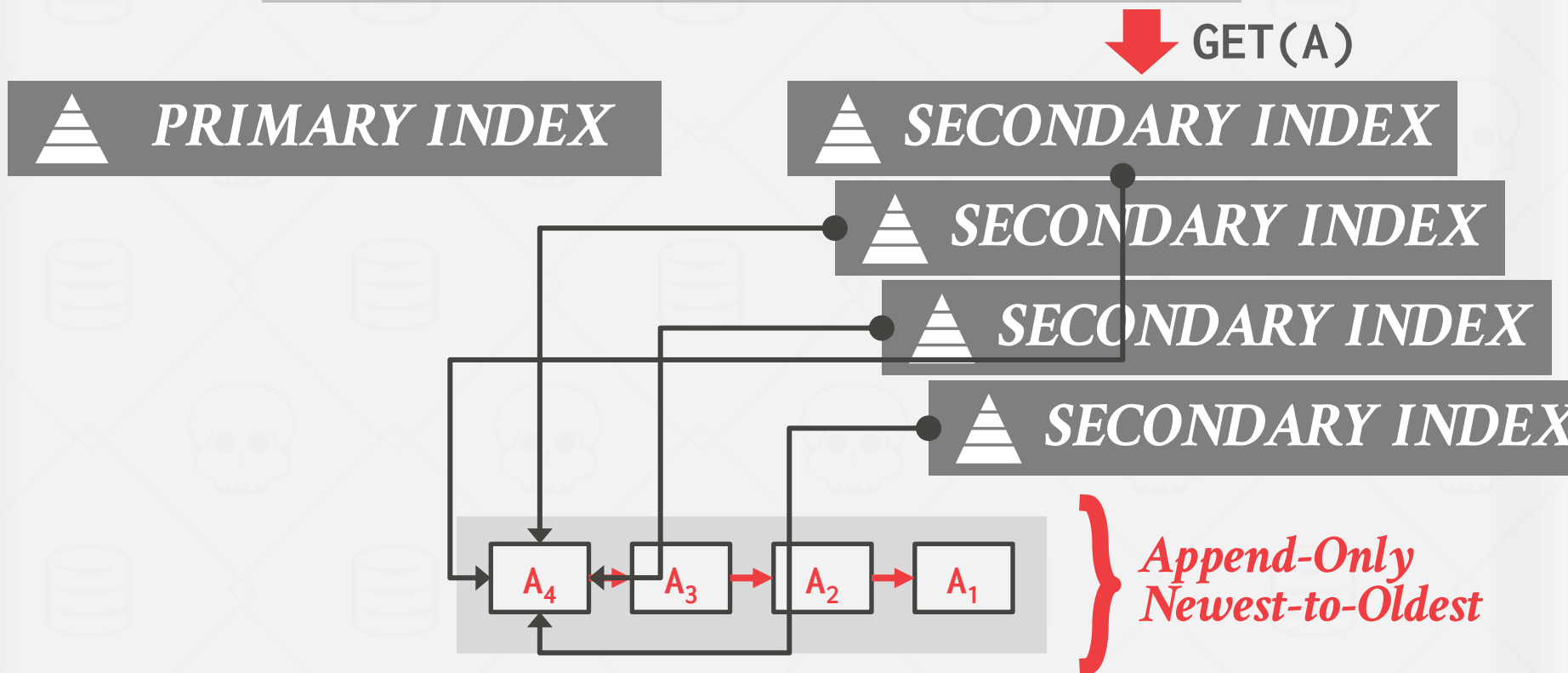| A₄ | → | A₃ | → | A₂ | → | A₁ |

$$\}$$ *Append-Only Newest-to-Oldest*

# INDEX POINTERS

GET(A)

PRIMARY INDEX

SECONDARY INDEX

Record Id

$A_4$ → $A_3$ → $A_2$ → $A_1$

} *Append-Only*
*Newest-to-Oldest*

# INDEX POINTERS

GET(A)

PRIMARY INDEX

SECONDARY INDEX

SECONDARY INDEX

SECONDARY INDEX

SECONDARY INDEX

$A_4$   $A_3$   $A_2$   $A_1$

} *Append-Only Newest-to-Oldest*

# INDEX POINTERS



GET(A)

PRIMARY INDEX

SECONDARY INDEX

Primary
Key

Record Id

A₄ → A₃ → A₂ → A₁

} *Append-Only*
*Newest-to-Oldest*

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

*TupleId*

*TupleId→Address*

*Record Id*

| $A_4$ | → | $A_3$ | → | $A_2$ | → | $A_1$ |

} *Append-Only Newest-to-Oldest*

# 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.
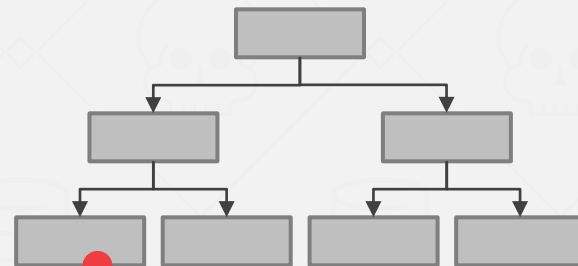
# MVCC DUPLICATE KEY PROBLEM

*Txn #1*

BEGIN @ 10

READ(A)

*Index*



| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| $A_1$ | *1* | ∞ | Ø |
| | | | |
| | | | |

# MVCC DUPLICATE KEY PROBLEM

*Txn #1*

BEGIN @ 10

READ(A)

*Txn #2*

BEGIN @ 20

UPDATE(A)

*Index*

| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| A₁ | 1 | 20 | ● |
| A₂ | 20 | ∞ | ∅ |
| | | | |

# MVCC DUPLICATE KEY PROBLEM

**Txn #1**

BEGIN @ 10

READ(A)

**Txn #2**

BEGIN @ 20

UPDATE(A)    DELETE(A)

*Index*

| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| A₁ | *1* | *20* | ● |
| ✕ | *20* | *∞* | *∅* |
| | | | |

# MVCC DUPLICATE KEY PROBLEM

*Txn #1*

BEGIN @ 10

READ(A)

*Txn #2*

BEGIN @ 20
COMMIT @ 25

UPDATE(A)

DELETE(A)

*Index*

| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| $A_1$ | 1 | 20 | ● |
| ✗ | 20 | 20 | ∅ |
| | | | |

CMU·DB

# MVCC DUPLICATE KEY PROBLEM

**Txn #1**

BEGIN @ 10

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 | ∅ |
| $A_1$ | 30 | ∞ | ∅ |

# MVCC DUPLICATE KEY PROBLEM

*Txn #1*

BEGIN @ 10

*Txn #2*

BEGIN @ 20
COMMIT @ 25

*Txn #3*

BEGIN @ 30

READ(A)  READ(A)

UPDATE(A)  DELETE(A)

INSERT(A)

*Index*



| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| $A_1$ | 1 | 20 | ● |
| ✖ | 20 | 20 | ∅ |
| $A_1$ | 30 | ∞ | ∅ |

# 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 <u>physically</u> deletes a tuple from the database only when all versions of a <u>logically</u> 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

| | *Protocol* | *Version Storage* | *Garbage Collection* | *Indexes* |
|---|---|---|---|---|
| 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!