# Carnegie Mellon University

# 18

# Multi-Version Concurrency Control

Intro to Database Systems
15-445/15-645
Fall 2021

**LM** Lin Ma
Computer Science Carnegie
Mellon University

# ADMINISTRIVIA

**Project #3** is due Sun Nov 14[nd] @ 11:59pm.

**Homework #4** is due Wed Nov 10[th] @ 11:59pm.

# UPCOMING DATABASE TALK

## Vertica – High Performance Over Varying Terrain

→ Mon Nov 8th @ 4:30pm ET
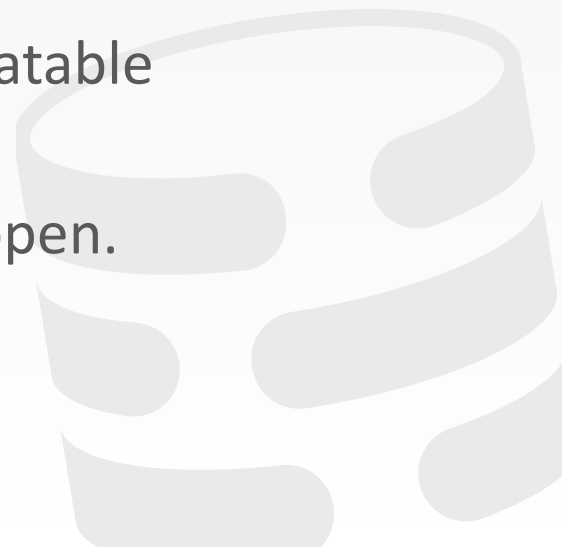
# ISOLATION LEVELS

*Isolation (High→Low)*

**SERIALIZABLE**: No phantoms, all reads repeatable, no dirty reads.

**REPEATABLE READS**: Phantoms may happen.

**READ COMMITTED**: Phantoms and unrepeatable reads may happen.

**READ UNCOMMITTED**: All of them may happen.

# ISOLATION LEVELS

| | Dirty Read | Unrepeatable Read | Phantom |
|---|---|---|---|
| SERIALIZABLE | No | No | No |
| REPEATABLE READ | No | No | Maybe |
| READ COMMITTED | No | Maybe | Maybe |
| READ UNCOMMITTED | Maybe | Maybe | Maybe |

# ISOLATION LEVELS

**SERIALIZABLE**: Obtain all locks first; plus index locks, plus strict 2PL.

**REPEATABLE READS**: Same as above, but no index locks.

**READ COMMITTED**: Same as above, but S locks are released immediately.

**READ UNCOMMITTED**: Same as above but allows dirty reads (no S locks).

# SQL-92 ISOLATION LEVELS

You set a txn's isolation level <u>before</u> you execute any queries in that txn.

Not all DBMS support all isolation levels in all execution scenarios
→ Replicated Environments

The default depends on implementation…

```
SET TRANSACTION ISOLATION LEVEL
   <isolation-level>;
```

```
BEGIN TRANSACTION ISOLATION LEVEL
   <isolation-level>;
```

# ISOLATION LEVELS (2013)

| | Default | Maximum |
|---|---|---|
| Actian Ingres 10.0/10S | SERIALIZABLE | SERIALIZABLE |
| Aerospike | READ COMMITTED | READ COMMITTED |
| Greenplum 4.1 | READ COMMITTED | SERIALIZABLE |
| MySQL 5.6 | REPEATABLE READS | SERIALIZABLE |
| MemSQL 1b | READ COMMITTED | READ COMMITTED |
| MS SQL Server 2012 | READ COMMITTED | SERIALIZABLE |
| Oracle 11g | READ COMMITTED | SNAPSHOT ISOLATION |
| Postgres 9.2.2 | READ COMMITTED | SERIALIZABLE |
| SAP HANA | READ COMMITTED | SERIALIZABLE |
| ScaleDB 1.02 | READ COMMITTED | READ COMMITTED |
| VoltDB | SERIALIZABLE | SERIALIZABLE |

Source: Peter Bailis

**CMU·DB**

**15-445/645 (Fall 2021)**

# ISOLATION LEVELS (2013)

| | *Default* | *Maximum* |
|---|---|---|
| Actian Ingres 10.0/10S | SERIALIZABLE | SERIALIZABLE |
| Aerospike | READ COMMITTED | READ COMMITTED |
| Greenplum 4.1 | READ COMMITTED | SERIALIZABLE |
| MySQL 5.6 | REPEATABLE READS | SERIALIZABLE |
| MemSQL 1b | READ COMMITTED | READ COMMITTED |
| MS SQL Server 2012 | READ COMMITTED | SERIALIZABLE |
| Oracle 11g | READ COMMITTED | SNAPSHOT ISOLATION |
| Postgres 9.2.2 | READ COMMITTED | SERIALIZABLE |
| SAP HANA | READ COMMITTED | SERIALIZABLE |
| ScaleDB 1.02 | READ COMMITTED | READ COMMITTED |
| VoltDB | SERIALIZABLE | SERIALIZABLE |

Source: Peter Bailis

# ISOLATION LEVELS (2013)

| | *Default* | *Maximum* |
|---|---|---|
| Actian Ingres 10.0/10S | SERIALIZABLE | SERIALIZABLE |
| Aerospike | READ COMMITTED | READ COMMITTED |
| Greenplum 4.1 | READ COMMITTED | SERIALIZABLE |
| MySQL 5.6 | REPEATABLE READS | SERIALIZABLE |
| MemSQL 1b | READ COMMITTED | READ COMMITTED |
| MS SQL Server 2012 | READ COMMITTED | SERIALIZABLE |
| Oracle 11g | READ COMMITTED | SNAPSHOT ISOLATION |
| Postgres 9.2.2 | READ COMMITTED | SERIALIZABLE |
| SAP HANA | READ COMMITTED | SERIALIZABLE |
| ScaleDB 1.02 | READ COMMITTED | READ COMMITTED |
| VoltDB | SERIALIZABLE | SERIALIZABLE |

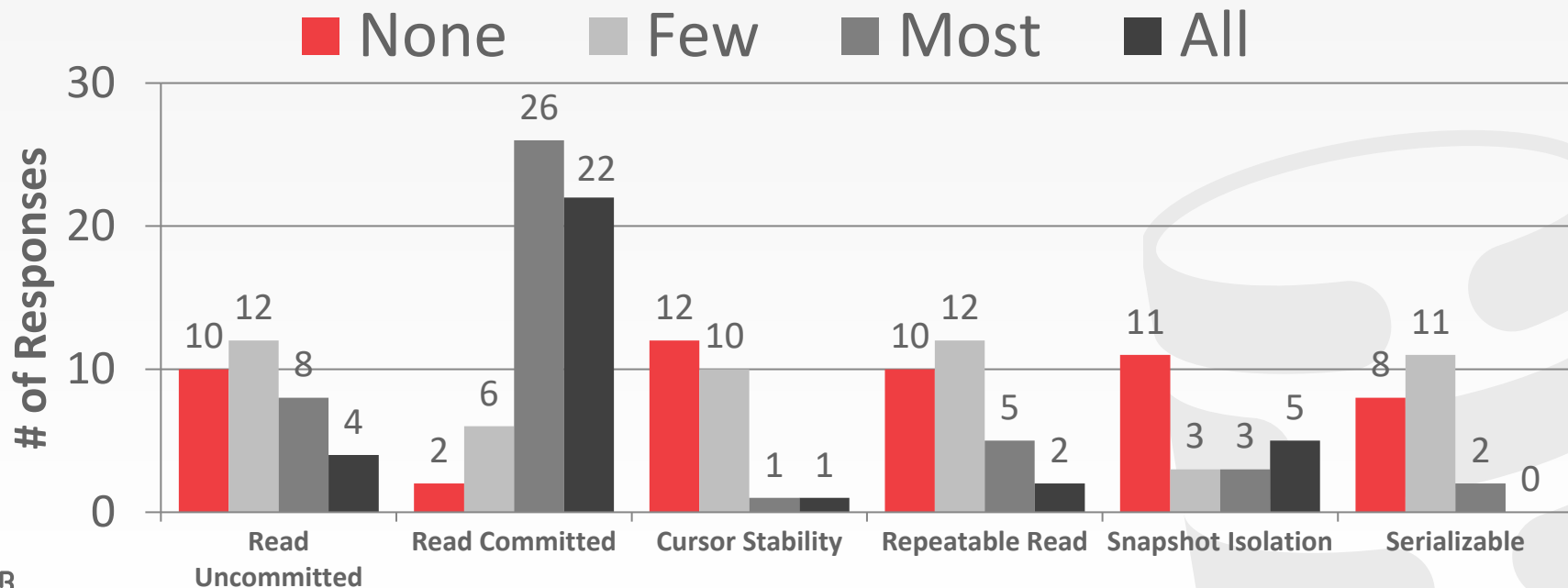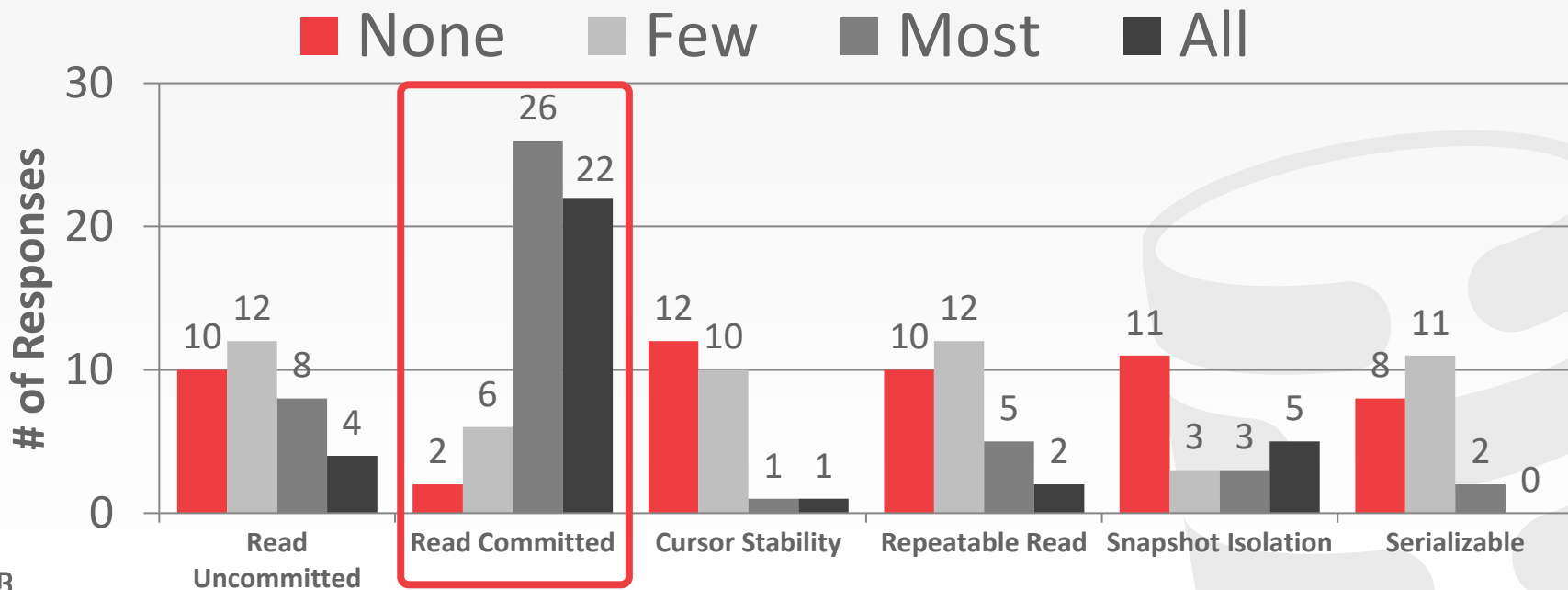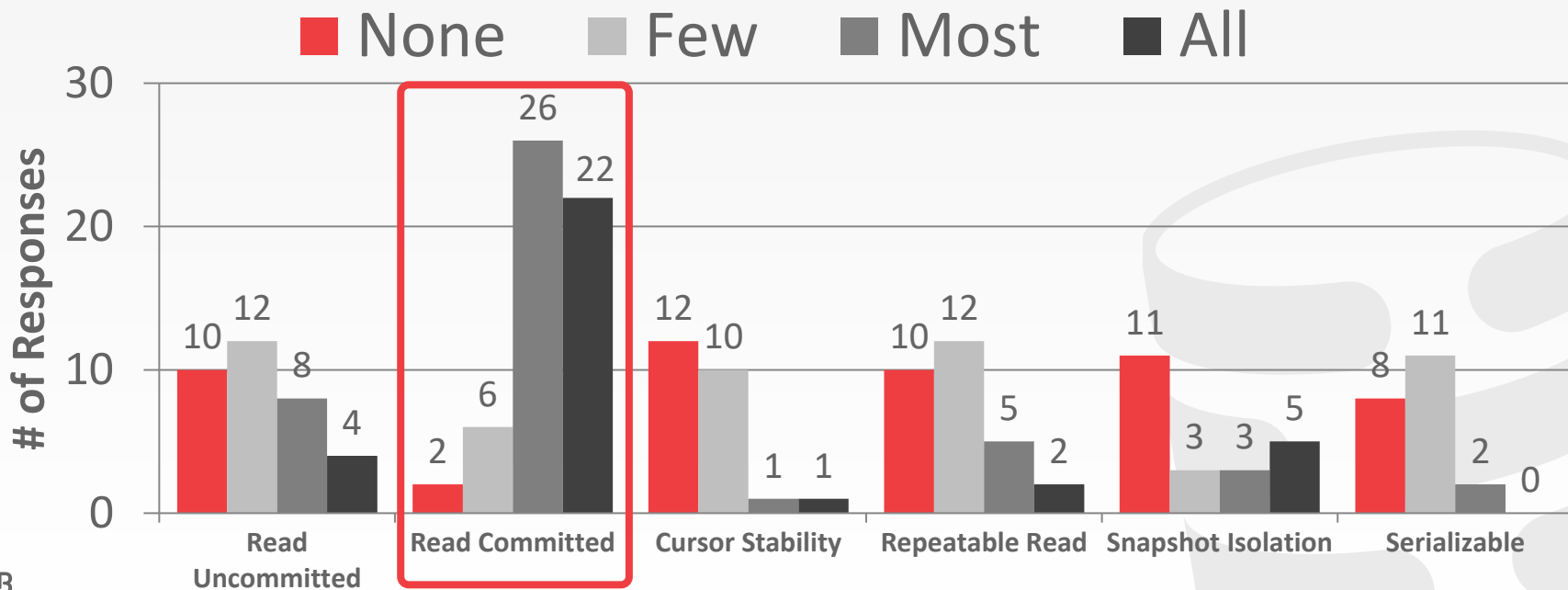Source: Peter Bailis

CMU·DB

15-445/645 (Fall 2021)

# DATABASE ADMIN SURVEY

What isolation level do transactions execute at on this DBMS?

# DATABASE ADMIN SURVEY

What isolation level do transactions execute at on this DBMS?

# DATABASE ADMIN SURVEY

What isolation level do transactions execute at on this DBMS?

# DATABASE ADMIN SURVEY

What isolation level do transactions execute at on this DBMS?

# SQL-92 ACCESS MODES

You can provide hints to the DBMS about whether a txn will modify the database during its lifetime.

Only two possible modes:
→ READ WRITE  (Default)
→ READ ONLY

Not all DBMSs will optimize execution if you set a txn to in READ ONLY mode.

```
SET TRANSACTION <access-mode>;
```

```
BEGIN TRANSACTION <access-mode>;
```

# 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.**
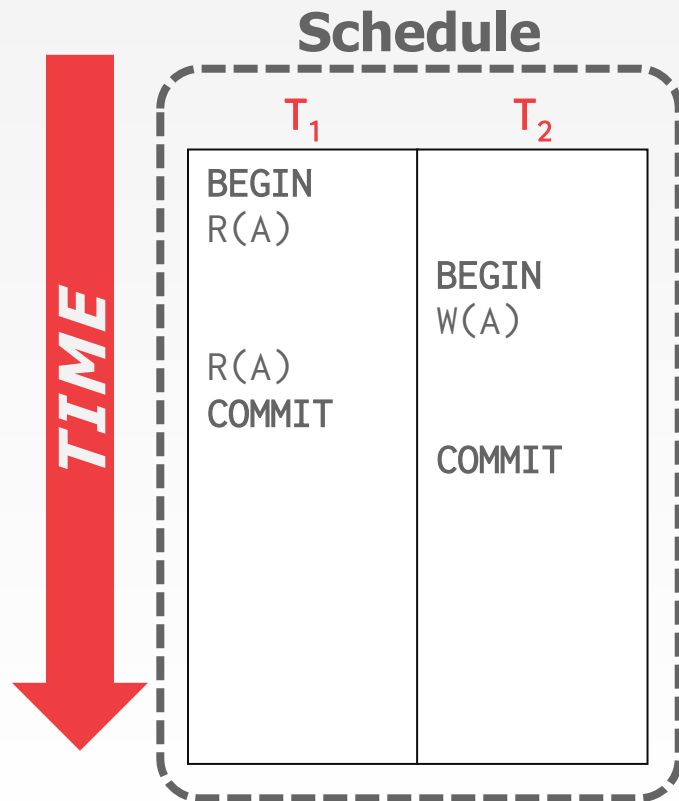**Readers do <u>not</u> block writers.**

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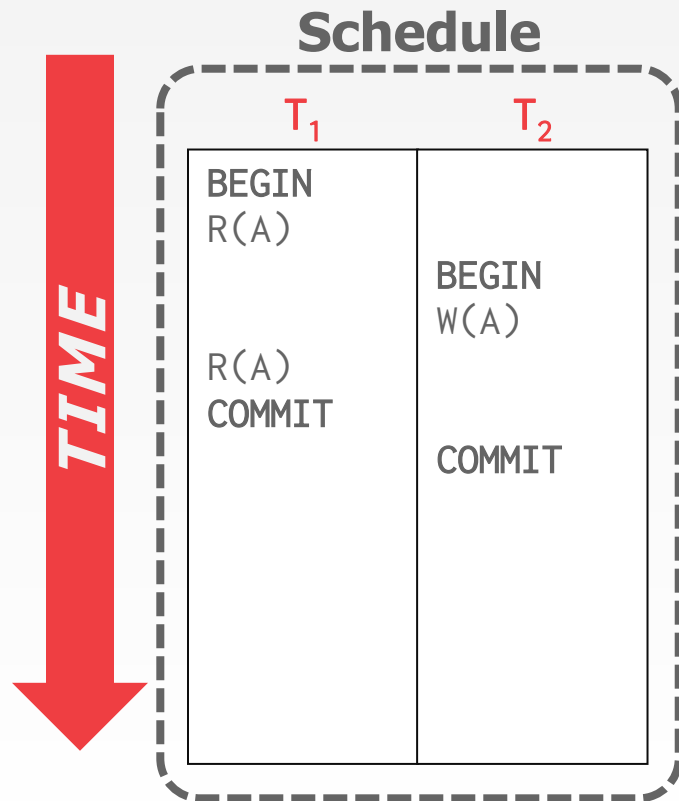
Easily support <u>time-travel</u> queries.

# MVCC – EXAMPLE #1

## Schedule

**TIME**

| T$_1$ | T$_2$ |
|---|---|
| BEGIN | |
| R(A) | |
| | BEGIN |
| | W(A) |
| R(A) | |
| COMMIT | |
| | COMMIT |

## Database

| Version | Value | Begin | End |
|---|---|---|---|
| A$_0$ | 123 | 0 | – |
| | | | |
| | | | |

# MVCC – EXAMPLE #1

## Schedule

**TIME**

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

## Database

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

# MVCC – EXAMPLE #1

## Schedule

**TIME**

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

## Database

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

# MVCC – EXAMPLE #1

**Schedule**

**TIME**

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

**Database**

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

# MVCC — EXAMPLE #1

**Schedule**

TS(T₁)=1   TS(T₂)=2

**TIME**

| | T₁ | T₂ |
|---|---|---|
| | BEGIN | |
| | R(A) | |
| | | BEGIN |
| | | W(A) |
| | R(A) | |
| | COMMIT | |
| | | COMMIT |

**Database**

| Version | Value | Begin | End |
|---|---|---|---|
| A₀ | 123 | 0 | – |
| | | | |
| | | | |

# MVCC – EXAMPLE #1



**TS(T₁)=1**

**TS(T₂)=2**

Schedule

T₁

T₂

```
BEGIN
R(A)


R(A)
COMMIT
```

```
BEGIN
W(A)


COMMIT
```

**TIME**

**Database**

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀ | 123 | 0 | – |
| | | | |
| | | | |

# MVCC — EXAMPLE #1

**TS(T₁)=1** **TS(T₂)=2**

**Database**

Schedule

T₁          T₂

```
BEGIN
R(A)

         BEGIN
         W(A)

R(A)
COMMIT

         COMMIT
```

TIME

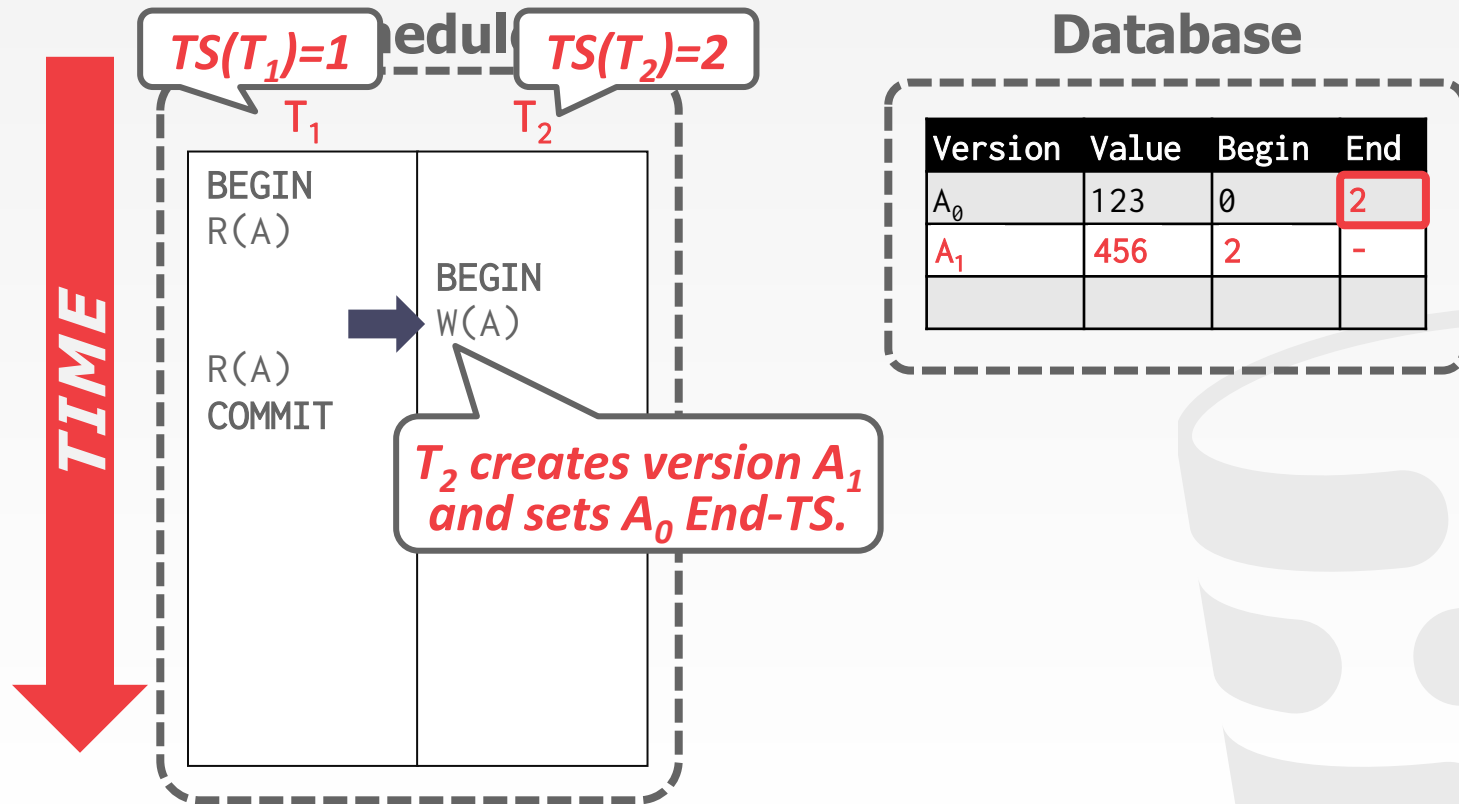| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀      | 123   | 0     | –   |
|         |       |       |     |
|         |       |       |     |

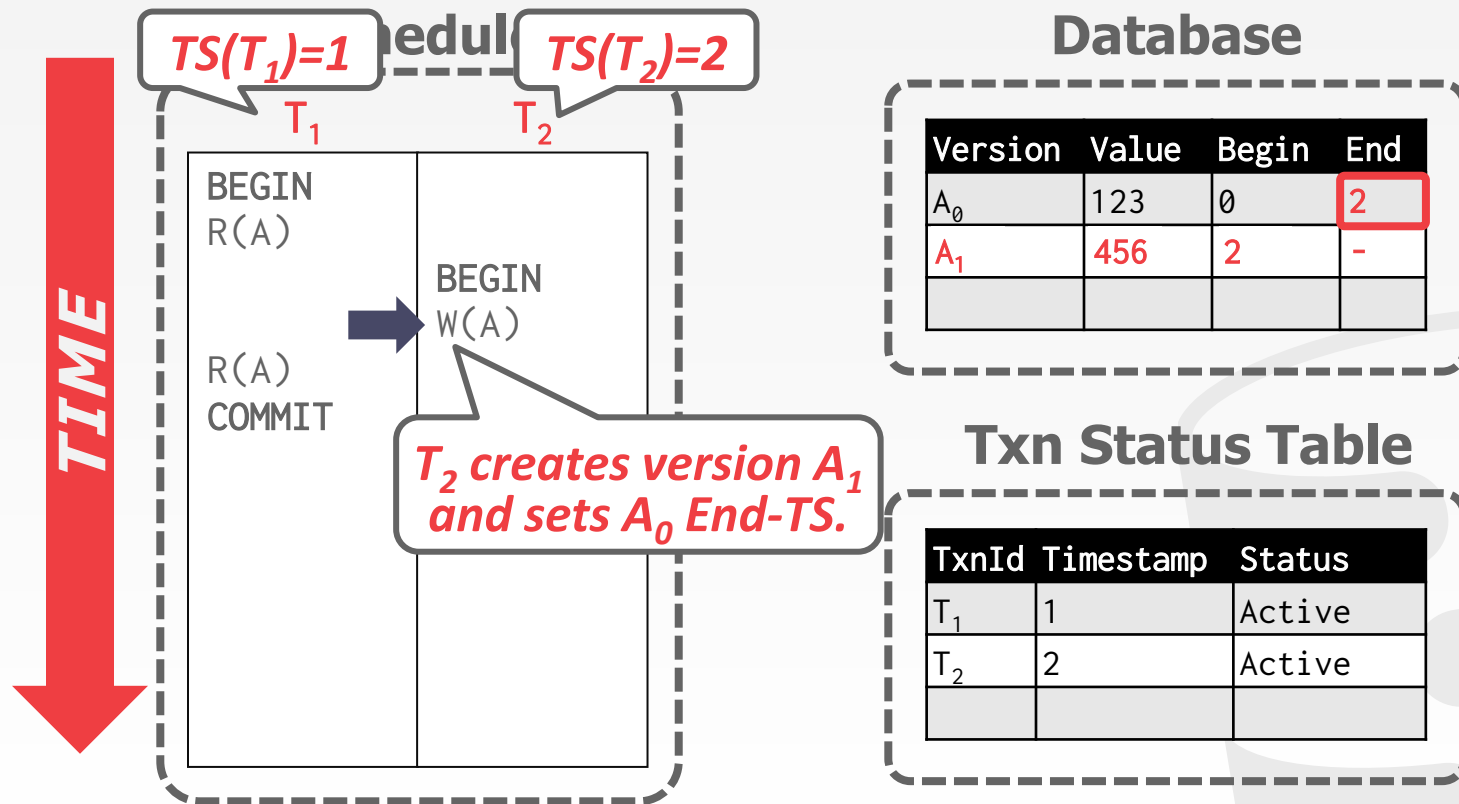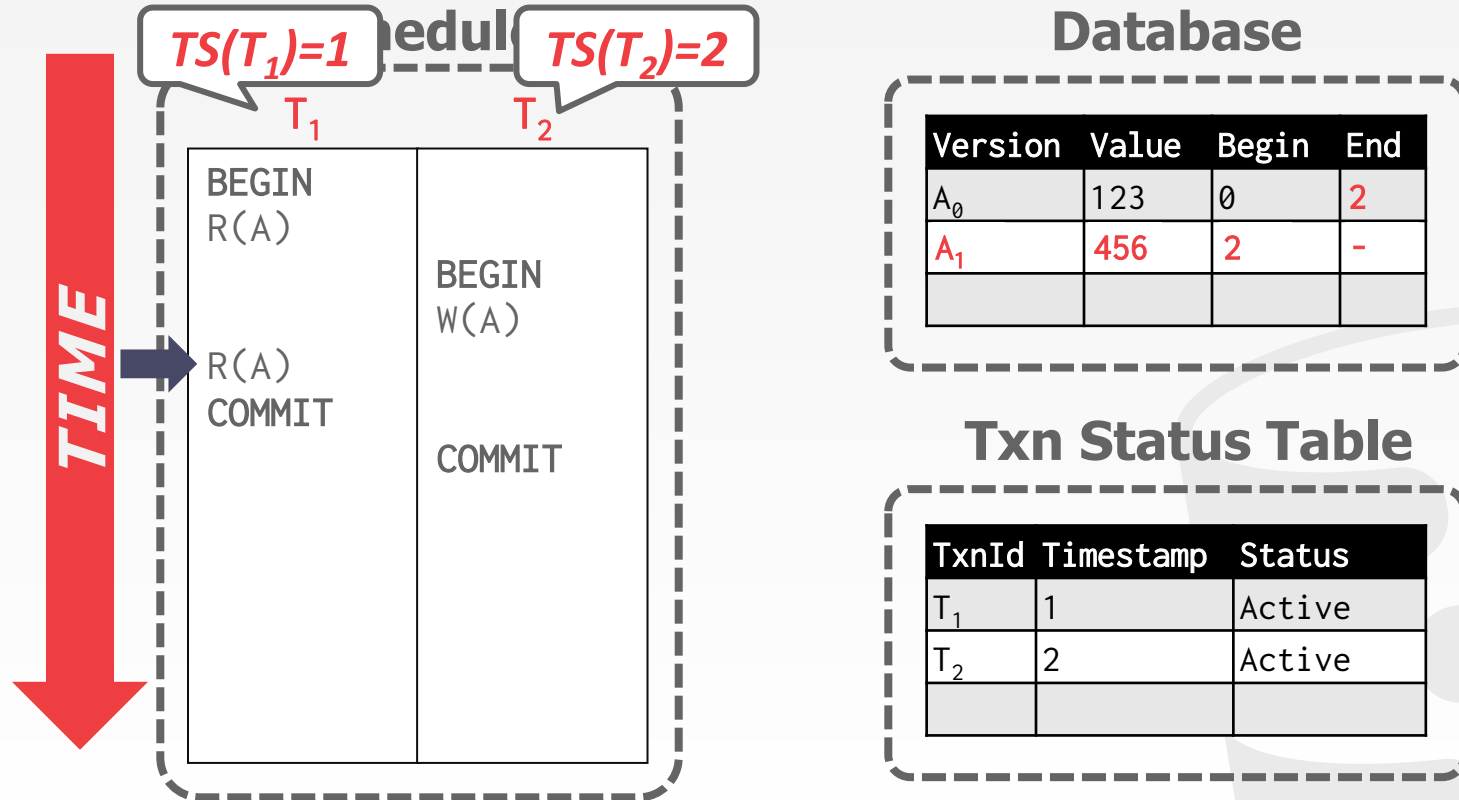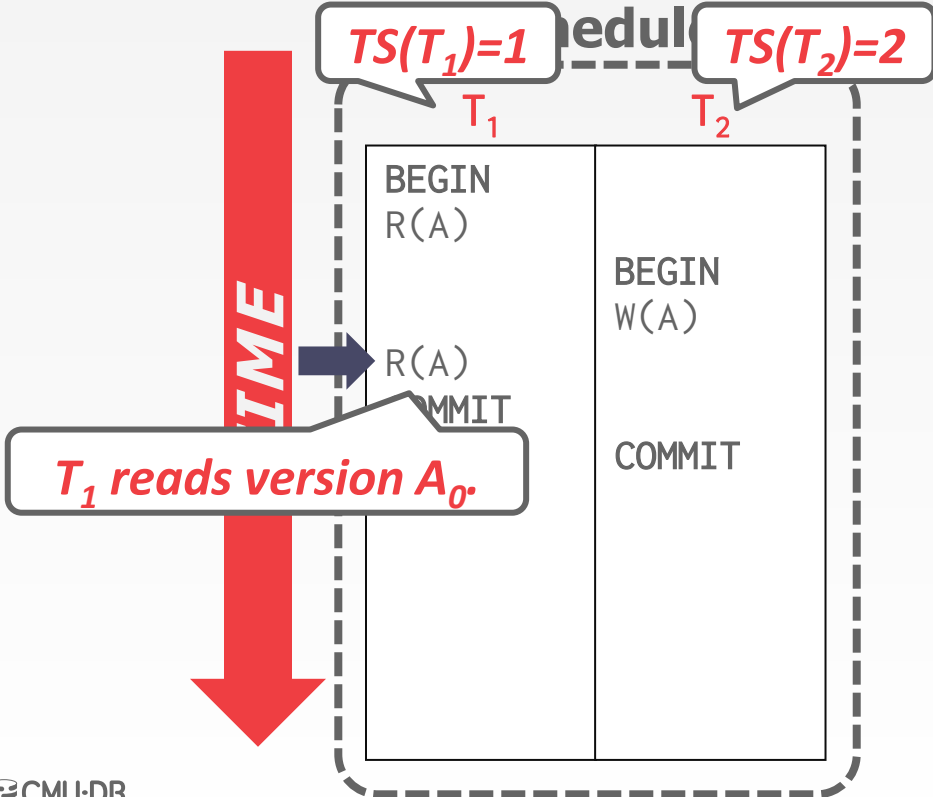# MVCC – EXAMPLE #1

# MVCC — EXAMPLE #1

# MVCC — EXAMPLE #1



**Schedule**

*TS(T₁)=1*    *TS(T₂)=2*

T₁          T₂

```
BEGIN
R(A)

          BEGIN
          W(A)

R(A)
COMMIT
```

TIME

*T₂ creates version A₁ and sets A₀ End-TS.*

**Database**

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀ | 123 | 0 | 2 |
| A₁ | 456 | 2 | - |
|  |  |  |  |

**Txn Status Table**

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Active |
| T₂ | 2 | Active |
|  |  |  |

# MVCC − EXAMPLE #1

TS(T₁)=1 edul TS(T₂)=2

T₁          T₂

| | |
|---|---|
| BEGIN<br>R(A)<br><br><br>R(A)<br>COMMIT | BEGIN<br>W(A)<br><br><br><br>COMMIT |

**TIME**

## Database

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀ | 123 | 0 | 2 |
| A₁ | 456 | 2 | - |
| | | | |

## Txn Status Table

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Active |
| T₂ | 2 | Active |
| | | |

# MVCC – EXAMPLE #1

**TS(T₁)=1**

**TS(T₂)=2**

Schedule

| T₁ | T₂ |
|---|---|
| BEGIN | |
| R(A) | |
| | BEGIN |
| | W(A) |
| R(A) | |
| COMMIT | |
| | COMMIT |

*TIME*

**T₁ reads version A₀.**

## Database

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀ | 123 | 0 | 2 |
| A₁ | 456 | 2 | - |
| | | | |

## Txn Status Table

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Active |
| T₂ | 2 | Active |
| | | |

# MVCC — EXAMPLE #1

# MVCC – EXAMPLE #2

**TS(T₁)=1**

**TS(T₂)=2**

Schedule

T₁

T₂

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

                    BEGIN
                    R(A)
                    W(A)

R(A)
COMMIT


                    COMMIT
```

## Database

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| A₀      | 123   | 0     |     |
|         |       |       |     |
|         |       |       |     |

## Txn Status Table

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

**TIME**

# MVCC – EXAMPLE #2



Schedule

$TS(T_1)=1$

$TS(T_2)=2$

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

**TIME**

## Database

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

## Txn Status Table

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

# MVCC – EXAMPLE #2



**Schedule**

*TS(T₁)=1*  $TS(T_1)=1$

*TS(T₂)=2*  $TS(T_2)=2$

$T_1$

$T_2$

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

R(A)
COMMIT


             COMMIT
```

**TIME**

## Database

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

## Txn Status Table

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

# MVCC – EXAMPLE #2

# MVCC — EXAMPLE #2



**TS(T₁)=1**

**TS(T₂)=2**

T₁ / T₂ Schedule

T₁:
```
BEGIN
R(A)
W(A)


R(A)
COMMIT
```

T₂:
```
BEGIN
R(A)
W(A)



COMMIT
```

**Database**

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

**Txn Status Table**

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

# MVCC – EXAMPLE #2

# MVCC – EXAMPLE #2



**Schedule**

*TS(T₁)=1* — $TS(T_1)=1$

*TS(T₂)=2* — $TS(T_2)=2$

$T_1$

$T_2$

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

R(A)
COMMIT

             COMMIT
```

**TIME**

**Database**

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

**Txn Status Table**

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

# MVCC – EXAMPLE #2



**TS(T₁)=1**  **TS(T₂)=2**

Schedule

|  | T₁ | T₂ |
|---|---|---|
|  | BEGIN |  |
|  | R(A) |  |
|  | W(A) |  |
|  |  | BEGIN |
|  |  | R(A) |
|  |  | W(A) |
|  | R(A) |  |
|  | COMMIT |  |
|  |  | COMMIT |

## Database

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

## Txn Status Table

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

# MVCC – EXAMPLE #2

# MVCC — EXAMPLE #2

# MVCC — EXAMPLE #2

# MVCC — EXAMPLE #2

Schedule

*TS(T₁)=1* — $TS(T_1)=1$

*TS(T₂)=2* — $TS(T_2)=2$

$T_1$

$T_2$

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


R(A)
COMMIT
```

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



COMMIT
```

**TIME**

## Database

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

## Txn Status Table

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

# MVCC – EXAMPLE #2



**Schedule**

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

$T_1$ | $T_2$

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

                BEGIN
                R(A)
                W(A)

R(A)
COMMIT
                COMMIT
```

**TIME**

*$T_1$ reads version $A_1$ that it wrote earlier.*

**Database**

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

**Txn Status Table**

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

# MVCC – EXAMPLE #2

**TS(T₁)=1**

**TS(T₂)=2**

Schedule

T₁

T₂

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


R(A)
COMMIT
```

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



COMMIT
```

## Database

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

## Txn Status Table

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Active |
| T₂ | 2 | Active |
| | | |

*TIME*

# MVCC — EXAMPLE #2



**Database**

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

**Txn Status Table**

| TxnId | Timestamp | Status |
|-------|-----------|--------|
| T₁ | 1 | Committed |
| T₂ | 2 | Active |
| | | |

# MVCC — EXAMPLE #2



**TS(T₁)=1**

**TS(T₂)=2**

Schedule

$T_1$

$T_2$

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


R(A)
COMMIT
```

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



COMMIT
```

**TIME**

**Database**

| Version | Value | Begin | End |
|---------|-------|-------|-----|
| $A_0$ | 123 | 0 | 1 |
| $A_1$ | 456 | 1 | 2 |
| $A_2$ | 789 | 2 | - |

**Txn Status Table**

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

*Now T₂ can create the new version.*

# MULTI-VERSION CONCURRENCY CONTROL

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

# 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 |
| | | |

# 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 |
| | | |

# 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*

# 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

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 | ∅ |
| $B_1$ | $10 | ∅ |
| $A_2$ | $333 | ∅ |

# 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 | ● |
| $B_1$ | $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**

| | VALUE | POINTER |
|---|---|---|
| $A_2$ | *$222* | ● |
| $B_1$ | *$10* | |

**Time-Travel Table**

| | VALUE | POINTER |
|---|---|---|
| $A_1$ | *$111* | Ø |
| | | |

# 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_**

| | VALUE | POINTER |
|---|---|---|
| $A_2$ | *$222* | ● |
| $B_1$ | *$10* | |

**_Time-Travel Table_**

| | 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 |
|---|---|---|
| A₂ | $222 | ● |
| 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**

| | 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**

| | 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**

| | 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.

# DELTA STORAGE

**Main Table**

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

**Delta Storage Segment**

# DELTA STORAGE

## Main Table

| | VALUE | POINTER |
|---|---|---|
| $A_1$ | *$111* | |
| $B_1$ | *$10* | |

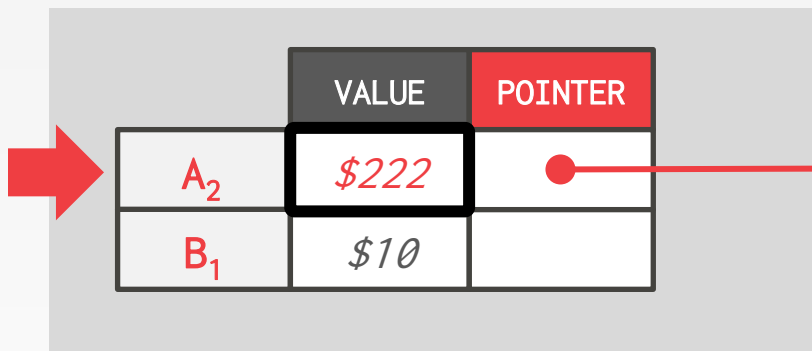## Delta Storage Segment

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 |
|---|---|---|
| A$_1$ | *$111* | |
| B$_1$ | *$10* | |

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

# 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*



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**

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 |
|---|---|---|
| A₃ | *$333* | ● |
| 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*

## *Delta Storage Segment*



| | VALUE | POINTER |
|---|---|---|
| $A_3$ | *$333* | ● |
| $B_1$ | *$10* | |

| | DELTA | POINTER |
|---|---|---|
| $A_1$ | *(VALUE→$111)* | Ø |
| $A_2$ | *(VALUE→$222)* | ● |

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

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?
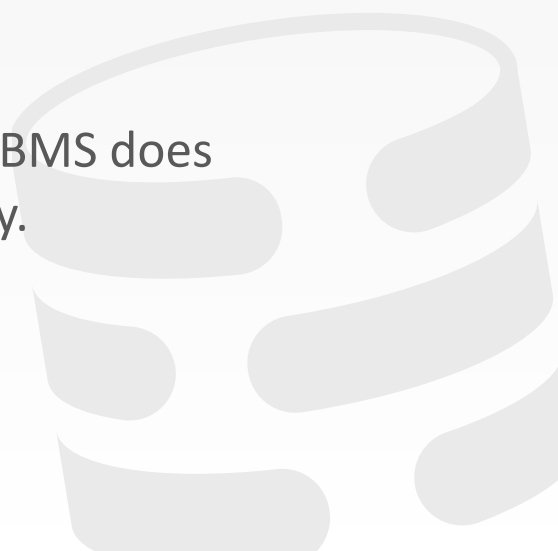→ How to decide when it is safe to reclaim memory?

# GARBAGE COLLECTION

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

**Thread #1**

$T_{id}=12$

**Thread #2**

$T_{id}=25$

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

# TUPLE-LEVEL GC

**Thread #1**

$T_{id}$=12

**Thread #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

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

**Thread #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

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

**Thread #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

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

**Vacuum**

**Thread #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.
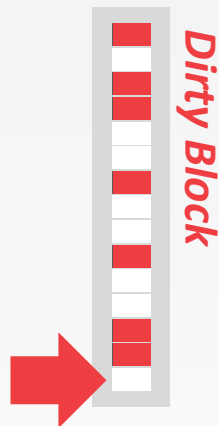
# TUPLE-LEVEL GC

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

**Thread #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

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

**Vacuum**

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

| | 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

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

**Vacuum**

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

*Dirty Block*

| | 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

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

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

*Vacuum*

*Dirty Block*

|  | 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

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

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

*Vacuum*

*Dirty Block*

| | 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.
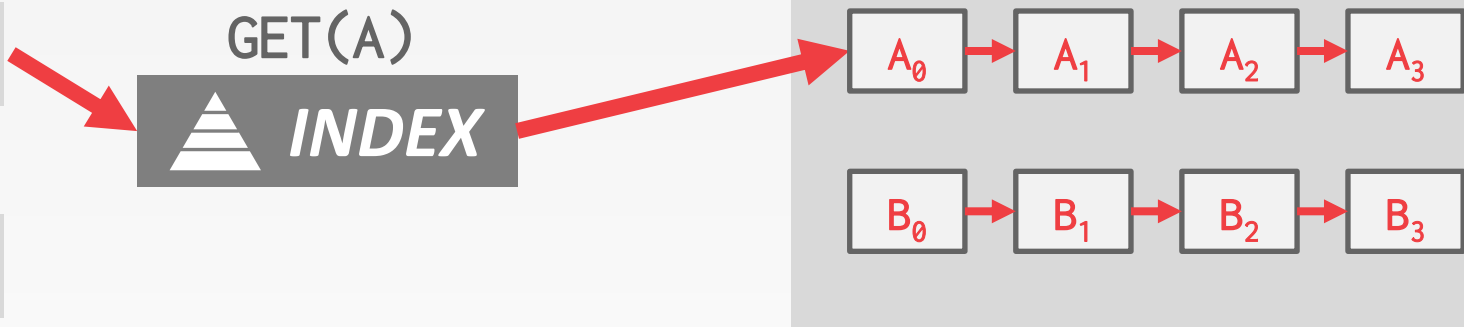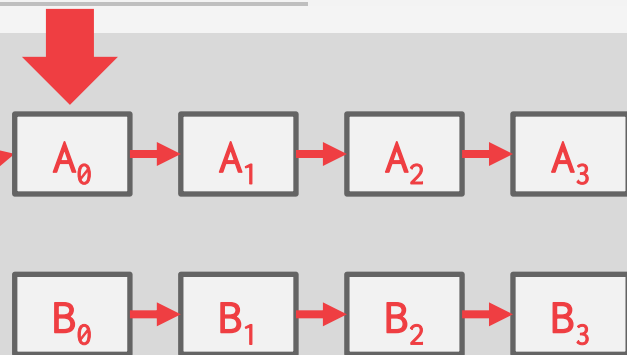
# TUPLE-LEVEL GC

**Thread #1**

$T_{id}=12$

**Thread #2**

$T_{id}=25$

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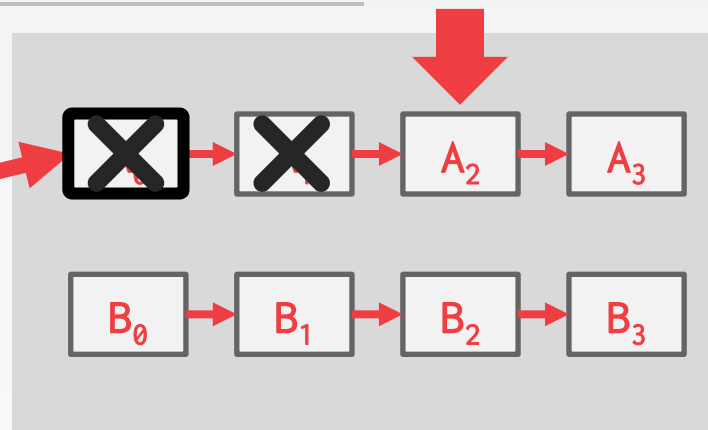
# TUPLE-LEVEL GC

**Thread #1**

$T_{id}=12$

**Thread #2**

$T_{id}=25$

**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

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

**INDEX**

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

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

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

**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.

# TUPLE-LEVEL GC



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

GET(A)

**INDEX**

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

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

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

**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.
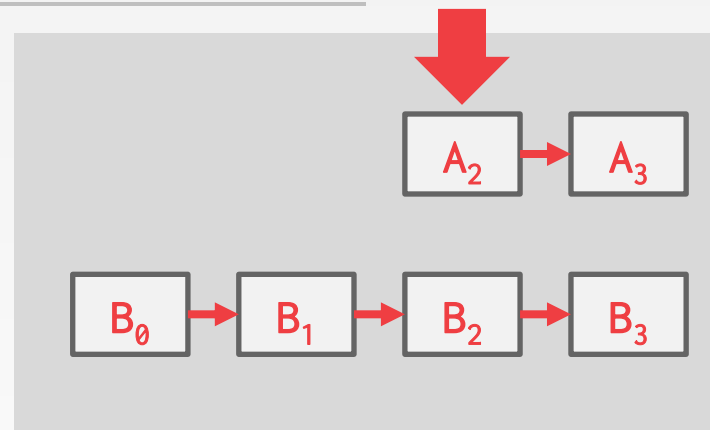
# TUPLE-LEVEL GC

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

GET(A)

△ **INDEX**

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

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

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

**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.
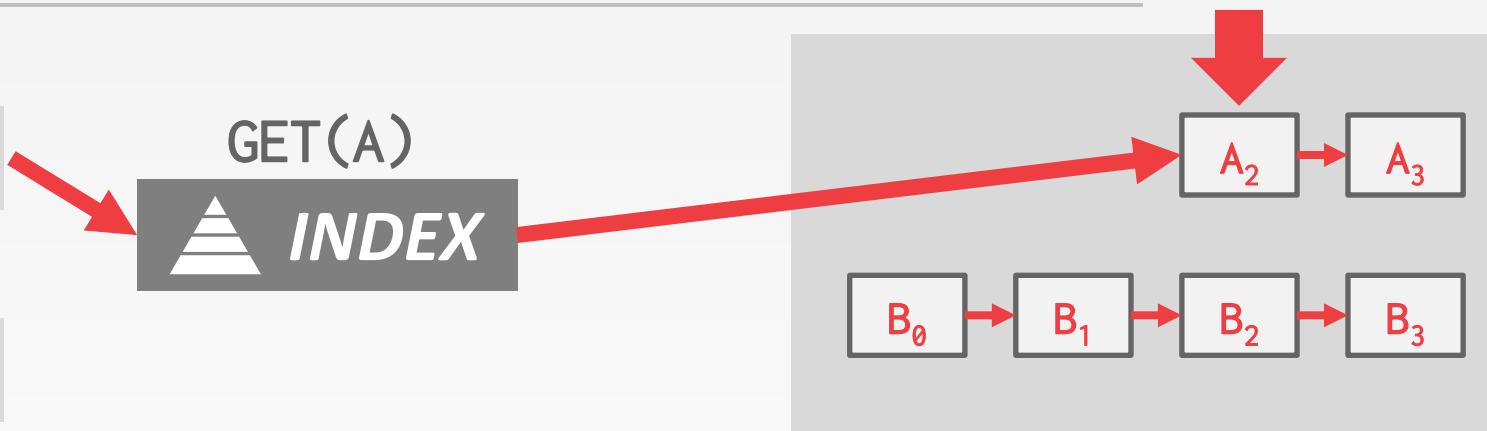
# TUPLE-LEVEL GC

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

GET(A)

**INDEX**

$A_0$ → $A_1$ → $A_2$ → $A_3$

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

$B_0$ → $B_1$ → $B_2$ → $B_3$

**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

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

GET(A)

**INDEX**

$A_1$ → $A_2$ → $A_3$

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

$B_0$ → $B_1$ → $B_2$ → $B_3$

**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

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

GET(A)

▲ **INDEX**

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



$A_2$ → $A_3$

$B_0$ → $B_1$ → $B_2$ → $B_3$

**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

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

GET(A)


**INDEX**

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



**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

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

GET(A)

**INDEX**

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

$A_2 \rightarrow A_3$

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

**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.

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

May still require multiple threads to reclaim the memory fast enough for the workload.

# TRANSACTION-LEVEL GC

**Thread #1**

**Begin @ 10**

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| A$_2$ | 1 | ∞ | - |
| B$_6$ | 8 | ∞ | - |
| | | | |
| | | | |

# TRANSACTION-LEVEL GC

**Thread #1**
*Begin @ 10*

UPDATE(A)

| | BEGIN-TS | END-TS | DATA |
|---|---|---|---|
| A$_2$ | *1* | ∞ | - |
| B$_6$ | *8* | ∞ | - |
| | | | |
| | | | |

# TRANSACTION-LEVEL GC

**Thread #1**
**Begin @ 10**

UPDATE(A)

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

# TRANSACTION-LEVEL GC

**Thread #1**
*Begin @ 10*

UPDATE(A)

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

# TRANSACTION-LEVEL GC

**Thread #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

**Thread #1**
*Begin @ 10*

UPDATE(A)

**Old Versions**

$A_2$

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

# TRANSACTION-LEVEL GC

**Thread #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



**Thread #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

**Thread #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

**Thread #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

**Thread #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

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4$ → $A_3$ → $A_2$ → $A_1$

# INDEX POINTERS

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

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

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4$ → $A_3$ → $A_2$ → $A_1$

*Append-Only Newest-to-Oldest*

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

*Physical Address*

$A_4$ → $A_3$ → $A_2$ → $A_1$

} *Append-Only Newest-to-Oldest*

# INDEX POINTERS

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

} *Append-Only Newest-to-Oldest*

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

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

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

*Physical Address*

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

*Append-Only Newest-to-Oldest*

CMU·DB

# 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

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**

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

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

# INDEX POINTERS

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$

*Append-Only Newest-to-Oldest*

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

**TupleId→Address**

$$A_4 \rightarrow A_3 \rightarrow A_2 \rightarrow A_1$$

} *Append-Only Newest-to-Oldest*

# INDEX POINTERS

GET(A)

**PRIMARY INDEX**

**SECONDARY INDEX**

*TupleId*

**TupleId→Address**

*Physical Address*

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

**} 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

*Index*



| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| A$_1$ | 1 | ∞ | Ø |
| | | | |
| | | | |

# MVCC DUPLICATE KEY PROBLEM

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**
**Begin @ 10**

READ(A)

**Thread #2**
**Begin @ 20**

UPDATE(A)

*Index*

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

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**

*Begin @ 10*

READ(A)

**Thread #2**

*Begin @ 20*

UPDATE(A)

*Index*

|  | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| $A_1$ | *1* | *20* | ● |
| $A_2$ | *20* | *∞* | *∅* |
|  |  |  |  |

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**
**Begin @ 10**

READ(A)

**Thread #2**
**Begin @ 20**

UPDATE(A)  DELETE(A)

*Index*

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

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**
*Begin @ 10*

READ(A)

**Thread #2**
*Begin @ 20*
*Commit @ 25*

UPDATE(A)  DELETE(A)

*Index*



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

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**
*Begin @ 10*

READ(A)

**Thread #2**
*Begin @ 20*
*Commit @ 25*

UPDATE(A)   DELETE(A)

*Index*

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

# MVCC DUPLICATE KEY PROBLEM

**Thread #1**
*Begin @ 10*

READ(A)

**Thread #2**
*Begin @ 20*
*Commit @ 25*

UPDATE(A)    DELETE(A)

**Thread #3**
*Begin @ 30*

INSERT(A)

*Index*

| | BEGIN-TS | END-TS | POINTER |
|---|---|---|---|
| A₁ | 1 | 20 | ● |
| ✕ | 20 | 20 | Ø |
| | | | |

# MVCC DUPLICATE KEY PROBLEM

# MVCC DUPLICATE KEY PROBLEM



**Thread #1**
*Begin @ 10*

**Thread #2**
*Begin @ 20*
*Commit @ 25*

**Thread #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* | Ø |
| $A_1$ | *30* | *∞* | Ø |

# 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 | 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 |
| NoisePage | MV-OCC | Delta | Txn-level | 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

No class on Wed November 11$^{th}$