# Timestamp Ordering

Lecture #19

Database Systems
15-445/15-645
Fall 2017

Andy Pavlo
Computer Science Dept.
Carnegie Mellon Univ.

# ADMINISTRIVIA

**Homework #4**: Monday November 13th @ 11:59pm

**Project #3**: Wednesday November 15th @ 11:59am

CARNEGIE MELLON
**DATABASE GROUP**

# CONCURRENCY CONTROL APPROACHES

## Two-Phase Locking (2PL)
→ Determine serializability order of conflicting operations at runtime while txns execute.

## Timestamp Ordering (T/O)
→ Determine serializability order of txns before they execute.

# CONCURRENCY CONTROL APPROACHES

**Two-Phase Locking (2PL)**
→ Determine serializability order of conflicting operations at runtime while txns execute.

**Pessimistic**

**Timestamp Ordering (T/O)**
→ Determine serializability order of txns before they execute.

**Optimistic**

# T/O CONCURRENCY CONTROL

Use timestamps to determine the serializability order of txns.

If $TS(T_i) < TS(T_j)$, then the DBMS must ensure that the execution schedule is equivalent to a serial schedule where $T_i$ appears before $T_j$.

CARNEGIE MELLON
DATABASE GROUP

# TIMESTAMP ALLOCATION

Each txn $T_i$ is assigned a unique fixed timestamp that is monotonically increasing.
→ Let $TS(T_i)$ be the timestamp allocated to txn $T_i$
→ Different schemes assign timestamps at different times during the txn.

Multiple implementation strategies:
→ System Clock.
→ Logical Counter.
→ Hybrid.

CARNEGIE MELLON
**DATABASE GROUP**

# TODAY'S AGENDA

Basic Timestamp Ordering

Optimistic Concurrency Control

Partition-based Timestamp Ordering

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O

Txns read and write objects without locks.

Every object **X** is tagged with timestamp of the last txn that successfully did read/write:
→ `W-TS(X)` – Write timestamp on **X**
→ `R-TS(X)` – Read timestamp on **X**

Check timestamps for every operation:
→ If txn tries to access an object "from the future", it aborts and restarts.

# BASIC T/O – READS

If $TS(T_i) < W-TS(X)$, this violates timestamp order of $T_i$ with regard to the writer of $X$.

→ Abort $T_i$ and restart it with same TS.

Else:

→ Allow $T_i$ to read $X$.

→ Update $R-TS(X)$ to $max(R-TS(X), TS(T_i))$

→ Have to make a local copy of $X$ to ensure repeatable reads for $T_i$.

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – WRITES

If $TS(T_i) < R\text{-}TS(X)$ or $TS(T_i) < W\text{-}TS(X)$
→ Abort and restart $T_i$.

Else:
→ Allow $T_i$ to write $X$ and update $W\text{-}TS(X)$
→ Also have to make a local copy of $X$ to
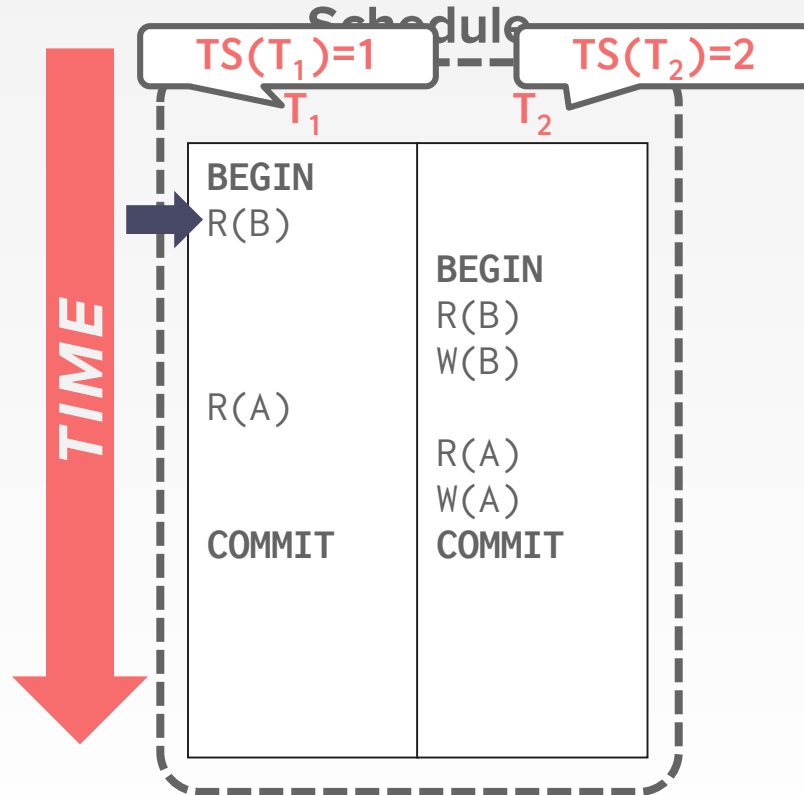ensure repeatable reads for $T_i$.

# BASIC T/O – EXAMPLE #1

**Schedule**

**Database**

TS(T₁)=1   TS(T₂)=2

T₁   T₂

**TIME**

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

| Object | R-TS | W-TS |
|--------|------|------|
| A | 0 | 0 |
| B | 0 | 0 |

# BASIC T/O – EXAMPLE #1

**Schedule**

TS(T$_1$)=1    TS(T$_2$)=2

**Database**

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

*TIME*

| Object | R-TS | W-TS |
|--------|------|------|
| A | 0 | 0 |
| B | 0 | 0 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #1

Schedule

Database

TS(T₁)=1  TS(T₂)=2

T₁          T₂

```
BEGIN
R(B)
            BEGIN
            R(B)
            W(B)

R(A)
            R(A)
            W(A)
COMMIT      COMMIT
```

TIME

| Object | R-TS | W-TS |
|--------|------|------|
| A      | 0    | 0    |
| B      | 1    | 0    |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #1

**Schedule**

**Database**

TS(T$_1$)=1

TS(T$_2$)=2

T$_1$          T$_2$

| Object | R-TS | W-TS |
|--------|------|------|
| A | 0 | 0 |
| B | 2 | 0 |

*TIME*

```
BEGIN
R(B)

                BEGIN
                R(B)
                W(B)

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

COMMIT          COMMIT
```

CARNEGIE MELLON
**DATABASE GROUP**

# BASIC T/O – EXAMPLE #1

Schedule

Database

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

**TIME**

$T_1$          $T_2$

```
BEGIN
R(B)

            BEGIN
            R(B)
      →     W(B)

R(A)

            R(A)
            W(A)

COMMIT      COMMIT
```

| Object | R-TS | W-TS |
|--------|------|------|
| A | 0 | 0 |
| B | 2 | 2 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #1

**Schedule**

TS(T₁)=1    TS(T₂)=2

**Database**

| Object | R-TS | W-TS |
|--------|------|------|
| A | 1 | 0 |
| B | 2 | 2 |

TIME

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

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #1

Schedule

Database

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

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

TIME

| Object | R-TS | W-TS |
|---|---|---|
| A | 2 | 0 |
| B | 2 | 2 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #1

Schedule

Database

$TS(T_1)=1$

$TS(T_2)=2$

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

TIME

| Object | R-TS | W-TS |
|--------|------|------|
| A      | 2    | 2    |
| B      | 2    | 2    |

No violations so both txns are safe to commit.

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule

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

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A | 0 | 0 |
| B | 0 | 0 |

*TIME*

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule

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

**TIME**

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A | 1 | 0 |
| B | 0 | 0 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule

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

**TIME**

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A | 1 | 2 |
| B | 0 | 0 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

Schedule

Database

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

| Object | R-TS | W-TS |
|---|---|---|
| A | 1 | 2 |
| B | 0 | 0 |

TIME

Violation:
TS(T₁) < W-TS(A)

T₁ cannot overwrite update by T₂, so the DBMS has to abort it!

CARNEGIE MELLON
DATABASE GROUP

# THOMAS WRITE RULE

If $TS(T_i) < R\text{-}TS(X)$:
→ Abort and restart $T_i$.

If $TS(T_i) < W\text{-}TS(X)$:
→ **Thomas Write Rule**: Ignore the write and allow the txn to continue.
→ This violates timestamp order of $T_i$.

Else:
→ Allow $T_i$ to write $X$ and update $W\text{-}TS(X)$

# BASIC T/O – EXAMPLE #2

## Schedule

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

TIME

## Database

| Object | R-TS | W-TS |
|---|---|---|
| A | 0 | 0 |
| B | 0 | 0 |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule

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

**TIME**

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A      | 1    | 0    |
| B      | 0    | 0    |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule

|  | T₁ | T₂ |

**T₁**            **T₂**

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

W(A)
COMMIT
```

**TIME**

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A      | 1    | 2    |
| B      | 0    | 0    |

CARNEGIE MELLON
DATABASE GROUP

# BASIC T/O – EXAMPLE #2

## Schedule



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

## Database

| Object | R-TS | W-TS |
|--------|------|------|
| A      | 1    | 2    |
| B      | 0    | 0    |

**TIME**

CARNEGIE MELLON
**DATABASE GROUP**

# BASIC T/O – EXAMPLE #2

# BASIC T/O

Generates a schedule that is conflict serializable if you do **not** use the Thomas Write Rule.
→ No deadlocks because no txn ever waits.
→ Possibility of starvation for long txns if short txns keep causing conflicts.

Permits schedules that are not **recoverable**.

CARNEGIE MELLON
**DATABASE GROUP**

# RECOVERABLE SCHEDULES

A schedule is **recoverable** if txns commit only after all txns whose changes they read, commit.

Otherwise, the DBMS cannot guarantee that txns read data that will be restored after recovering from a crash.

# RECOVERABLE SCHEDULES

**Schedule**



$T_2$ is allowed to read the writes of $T_1$.

# RECOVERABLE SCHEDULES

# BASIC T/O – PERFORMANCE ISSUES

High overhead from copying data to txn's workspace and from updating timestamps.

Long running txns can get starved.
→ The likelihood that a txn will read something from a newer txn increases.

Suffers from timestamp bottleneck.

CARNEGIE MELLON
DATABASE GROUP

# OBSERVATION

If you assume that conflicts between txns are **rare** and that most txns are **short-lived**, then forcing txns to wait to acquire locks adds a lot of overhead.

A better approach is to optimize for the no-conflict case.

CARNEGIE MELLON
DATABASE GROUP

# OPTIMISTIC CONCURRENCY CONTROL

The DBMS creates a private workspace for each txn.
→ All modifications are applied to workspace.
→ Any object read is copied into workspace.

When a txn commits, the DBMS compares its workspace write set to see whether it conflicts with other txns.

If there are no conflicts, the write set is installed into the "global" database.

# OCC PHASES

**#1 – Read Phase**:
→ Track the read/write sets of txns and store their writes in a private workspace.

**#2 – Validation Phase**:
→ When a txn commits, check whether it conflicts with other txns.

**#3 – Write Phase:**
→ If validation succeeds, apply private changes to database. Otherwise abort and restart the txn.
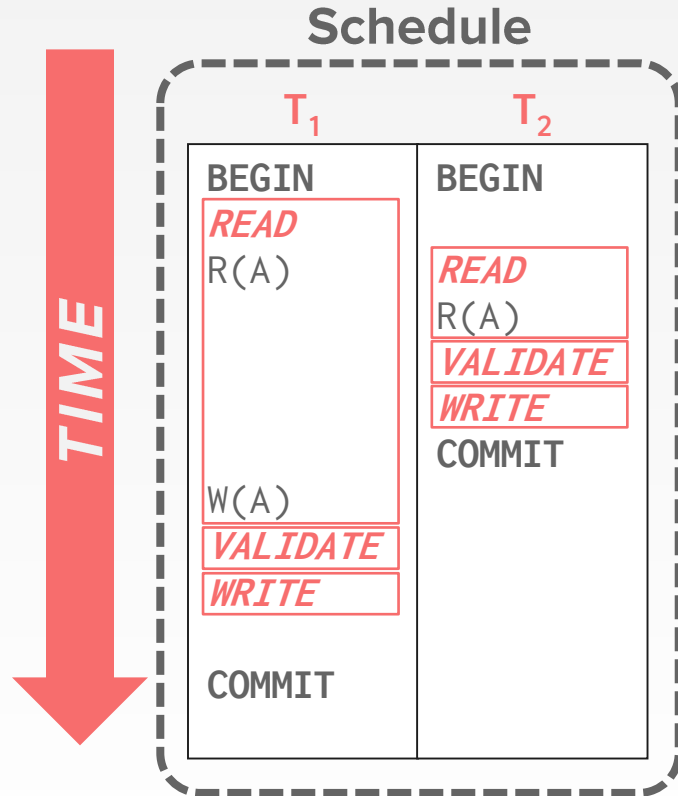
# OCC – EXAMPLE

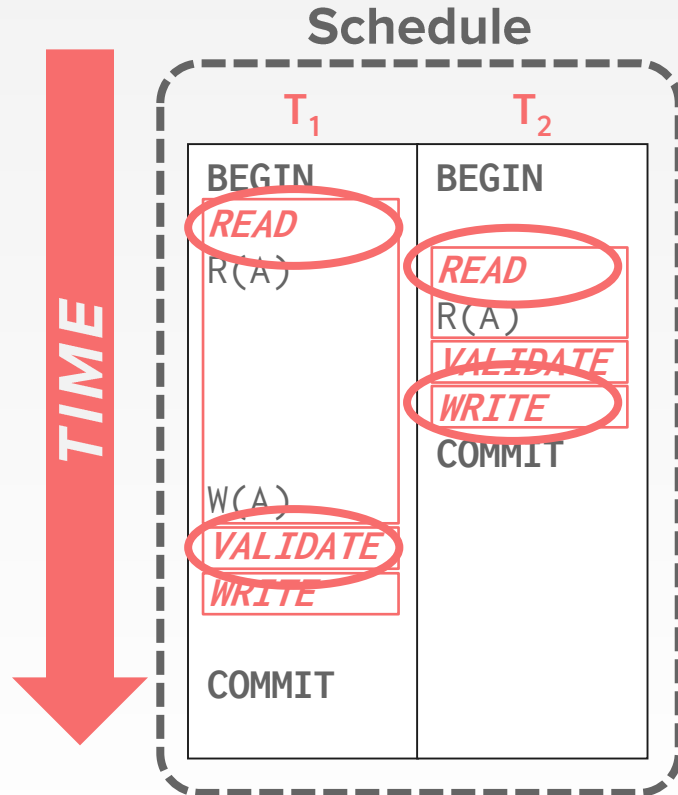## Schedule



## Database

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

# OCC – EXAMPLE

# OCC – EXAMPLE

# OCC – EXAMPLE

## Schedule



| T₁ | T₂ |
|---|---|
| BEGIN | BEGIN |
| *READb* | |
| R(A) | *READ* |
| | R(A) |
| | *VALIDATE* |
| | *WRITE* |
| | COMMIT |
| W(A) | |
| *VALIDATE* | |
| *WRITE* | |
| COMMIT | |

## Database

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

## T₁ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| –      | –     | –    |
| –      | –     | –    |

# OCC – EXAMPLE

**Schedule**

**Database**

TIME

| T<sub>1</sub> | T<sub>2</sub> |
|---|---|
| BEGIN | BEGIN |
| *READ* | |
| R(A) | *READ* |
| | R(A) |
| | *VALIDATE* |
| | *WRITE* |
| | COMMIT |
| W(A) | |
| *VALIDATE* | |
| *WRITE* | |
| COMMIT | |

| Object | Value | W-TS |
|---|---|---|
| A | 123 | 0 |
| – | – | – |

**T<sub>1</sub> Workspace**

| Object | Value | W-TS |
|---|---|---|
| A | 123 | 0 |
| – | – | – |

CARNEGIE MELLON
DATABASE GROUP

# OCC – EXAMPLE

## Schedule



**TIME**

### T₁
```
BEGIN
READ
R(A)



W(A)
VALIDATE
WRITE


COMMIT
```

### T₂
```
BEGIN

READ
R(A)
VALIDATE
WRITE
COMMIT
```

## Database

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

## T₁ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

## T₂ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| –      | –     | –    |
| –      | –     | –    |

CARNEGIE MELLON
DATABASE GROUP

# OCC – EXAMPLE

## Schedule



**TIME**

**T₁**

```
BEGIN
READ
R(A)



W(A)
VALIDATE
WRITE

COMMIT
```

**T₂**

```
BEGIN

READ
R(A)
VALIDATE
WRITE
COMMIT
```

## Database

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

## T₁ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

## T₂ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| A      | 123   | 0    |
| –      | –     | –    |

CARNEGIE MELLON
DATABASE GROUP

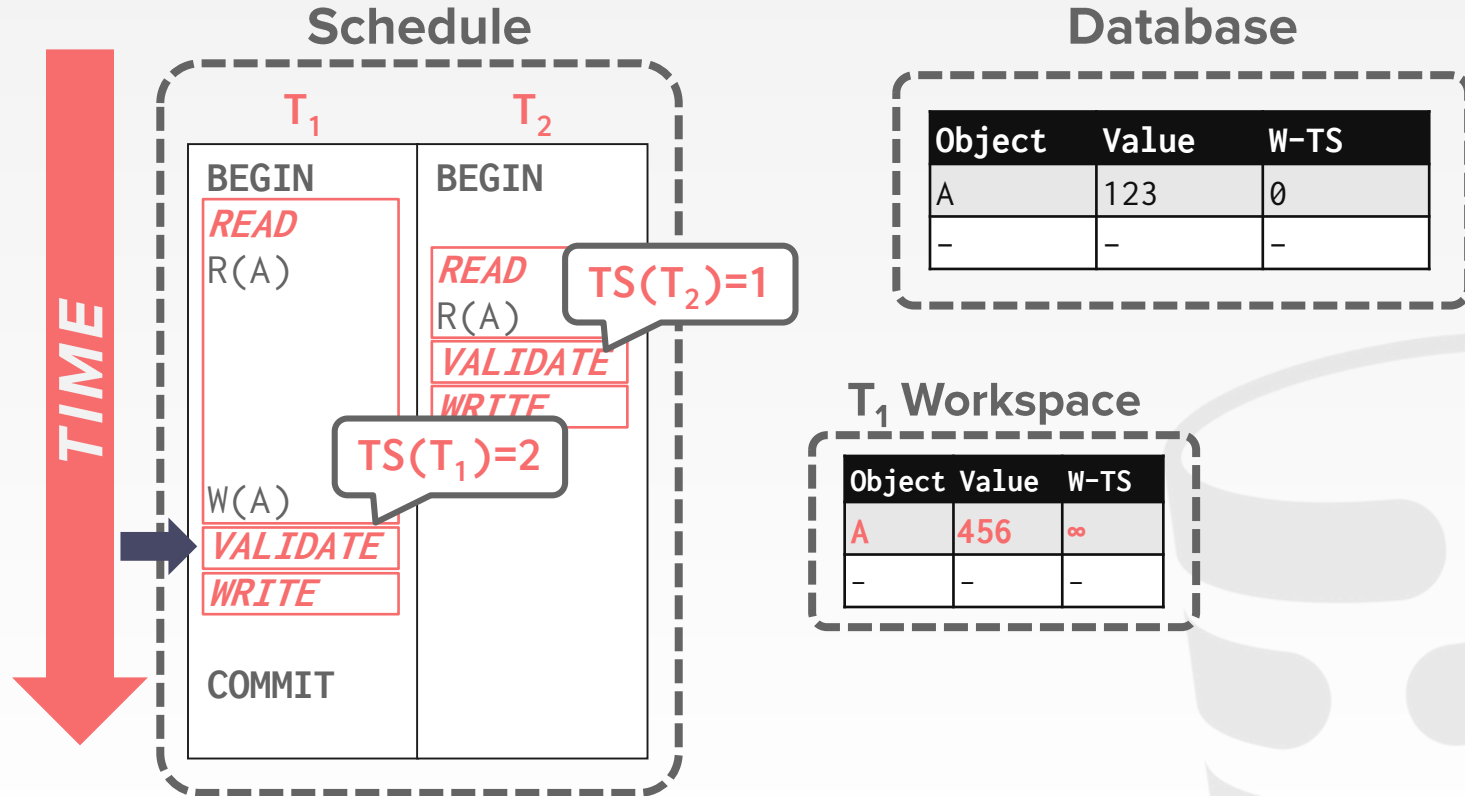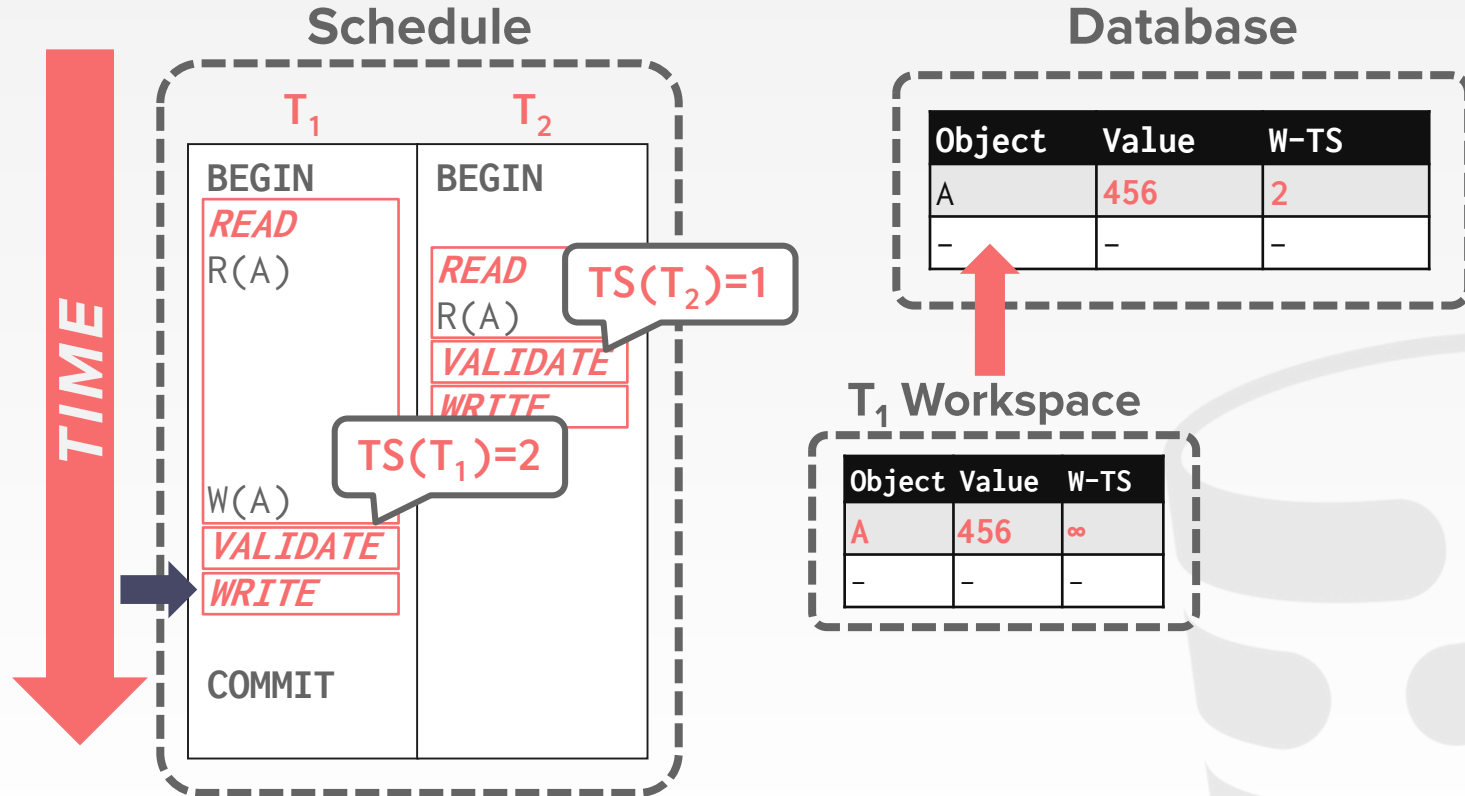# OCC – EXAMPLE

# OCC – EXAMPLE

# OCC – VALIDATION PHASE

The DBMS needs to guarantee only serializable schedules are permitted.

$T_i$ checks other txns for RW and WW conflicts and makes sure that all conflicts go one way (from older txns to younger txns).

# OCC – SERIAL VALIDATION

Maintain global view of all active txns.

Record read set and write set while txns are running and write into private workspace.

Execute **Validation** and **Write** phase inside a protected critical section.

# OCC – VALIDATION PHASE

Each txn's timestamp is assigned at the beginning of the validation phase.

Check the timestamp ordering of the committing txn with all other running txns.

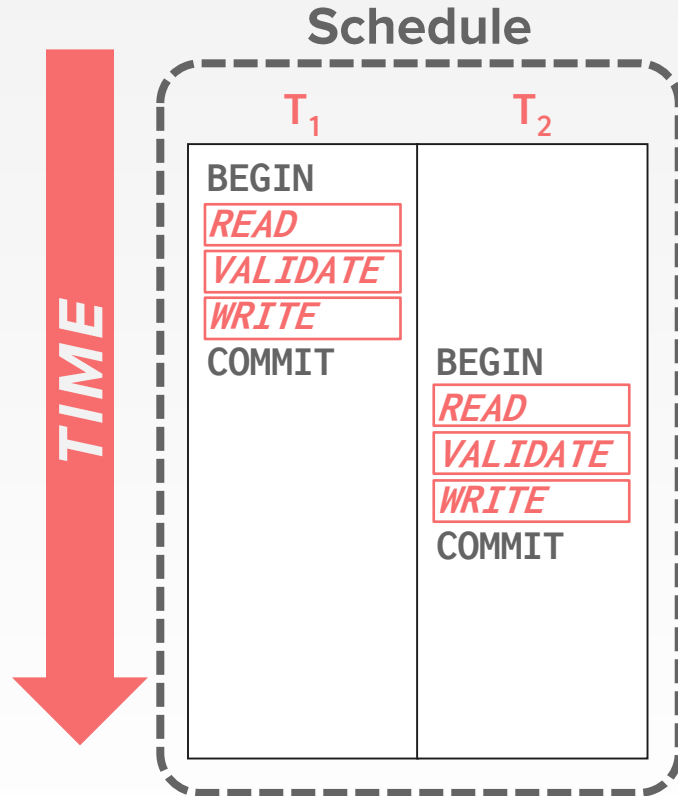If **TS(T$_i$)** < **TS(T$_j$)**, then one of the following three conditions must hold...

# OCC – VALIDATION STEP #1

$T_i$ completes all three phases before $T_j$ begins.

CARNEGIE MELLON
DATABASE GROUP

# OCC – VALIDATION STEP #1



Schedule

**T₁**

```
BEGIN
READ
VALIDATE
WRITE
COMMIT
```

**T₂**

```
BEGIN
READ
VALIDATE
WRITE
COMMIT
```
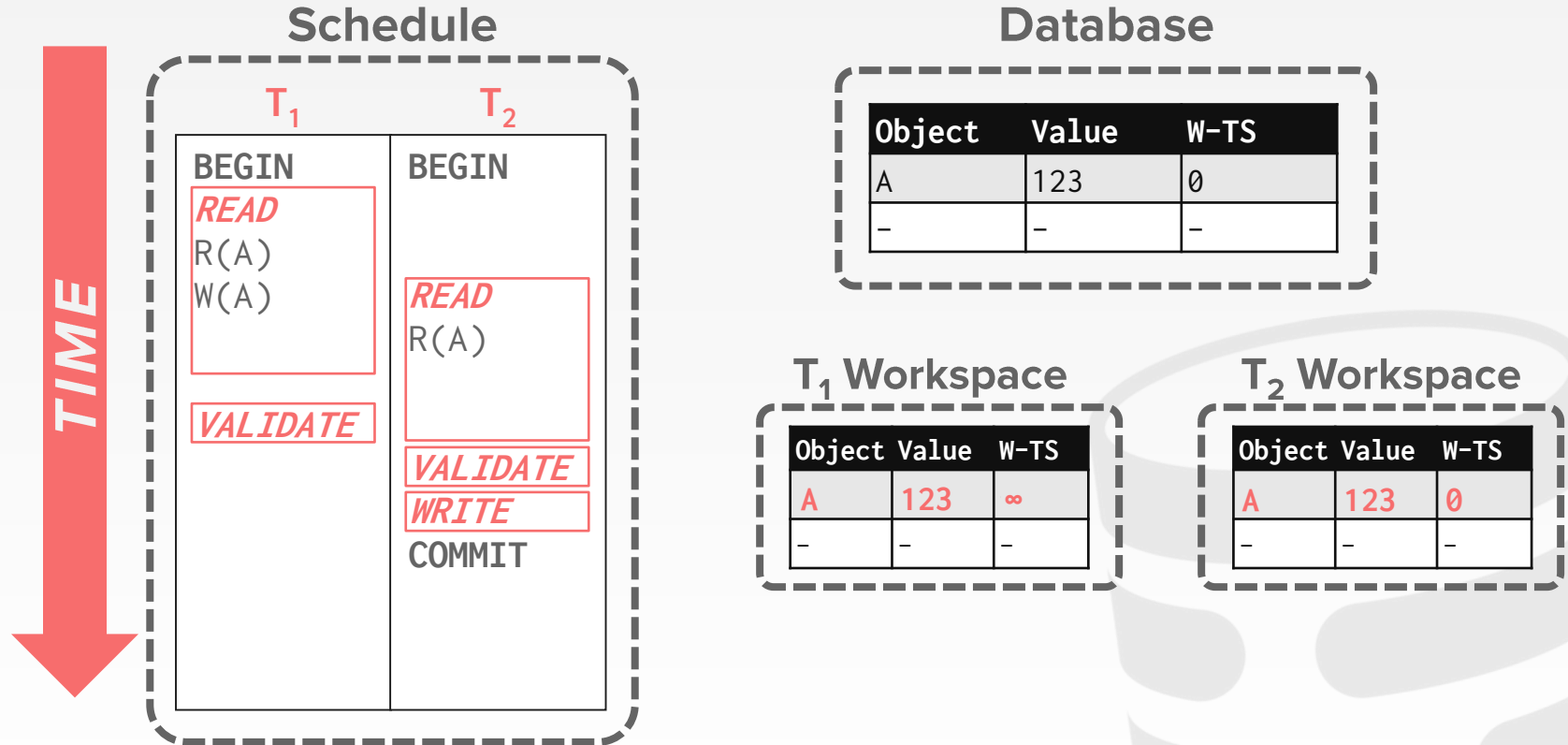
TIME

CARNEGIE MELLON
DATABASE GROUP

# OCC — VALIDATION STEP #2

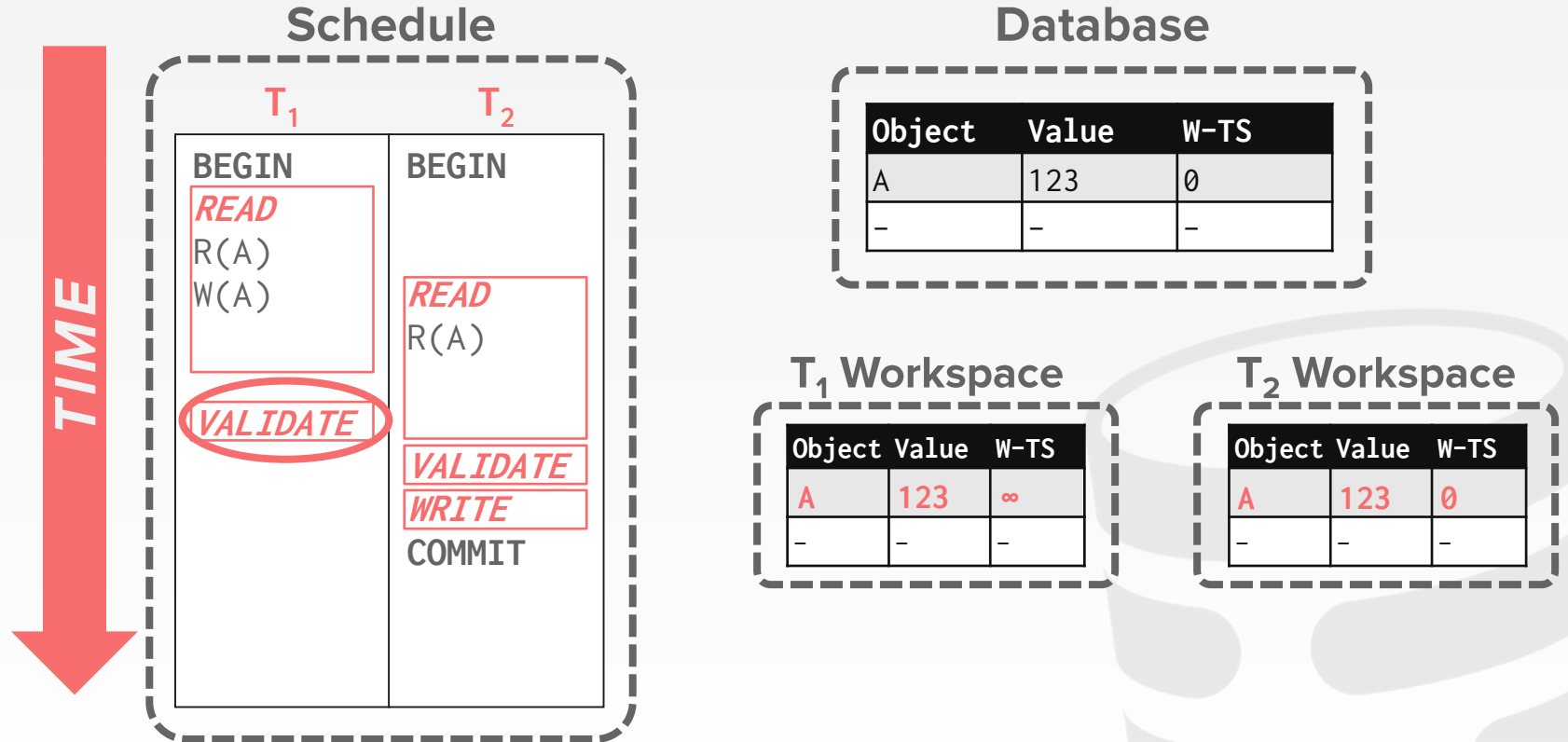$T_i$ completes before $T_j$ starts its **Write** phase, and $T_i$ does not write to any object read by $T_j$.
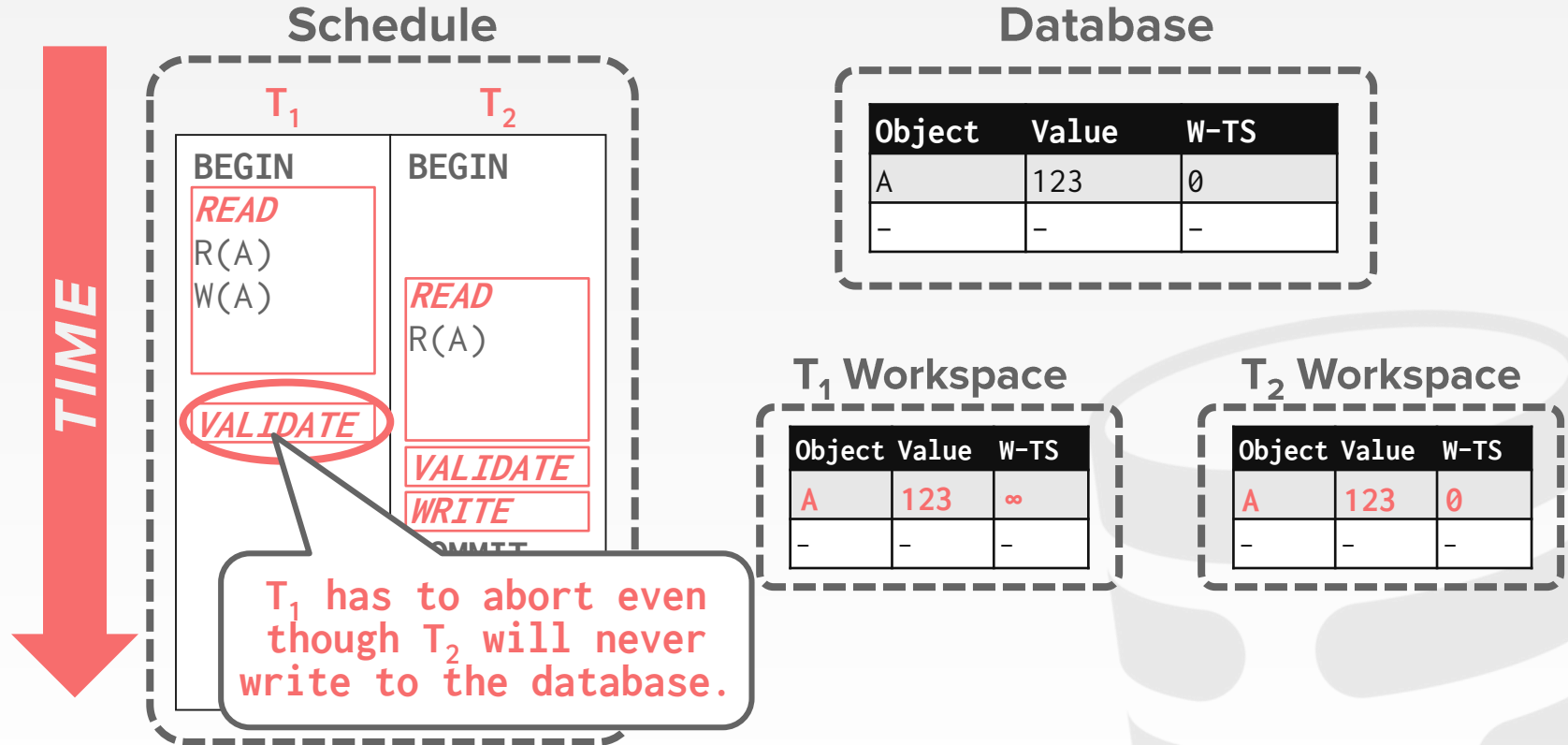→ $\text{WriteSet}(T_i) \cap \text{ReadSet}(T_j) = \varnothing$

CARNEGIE MELLON
DATABASE GROUP

# OCC – VALIDATION STEP #2

# OCC – VALIDATION STEP #2



Schedule

| T₁ | T₂ |
| BEGIN | BEGIN |
| READ | |
| R(A) | READ |
| W(A) | R(A) |
| VALIDATE | |
| | VALIDATE |
| | WRITE |
| | COMMIT |

TIME

Database

| Object | Value | W-TS |
|--------|-------|------|
| A | 123 | 0 |
| – | – | – |

T₁ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| A | 123 | ∞ |
| – | – | – |

T₂ Workspace

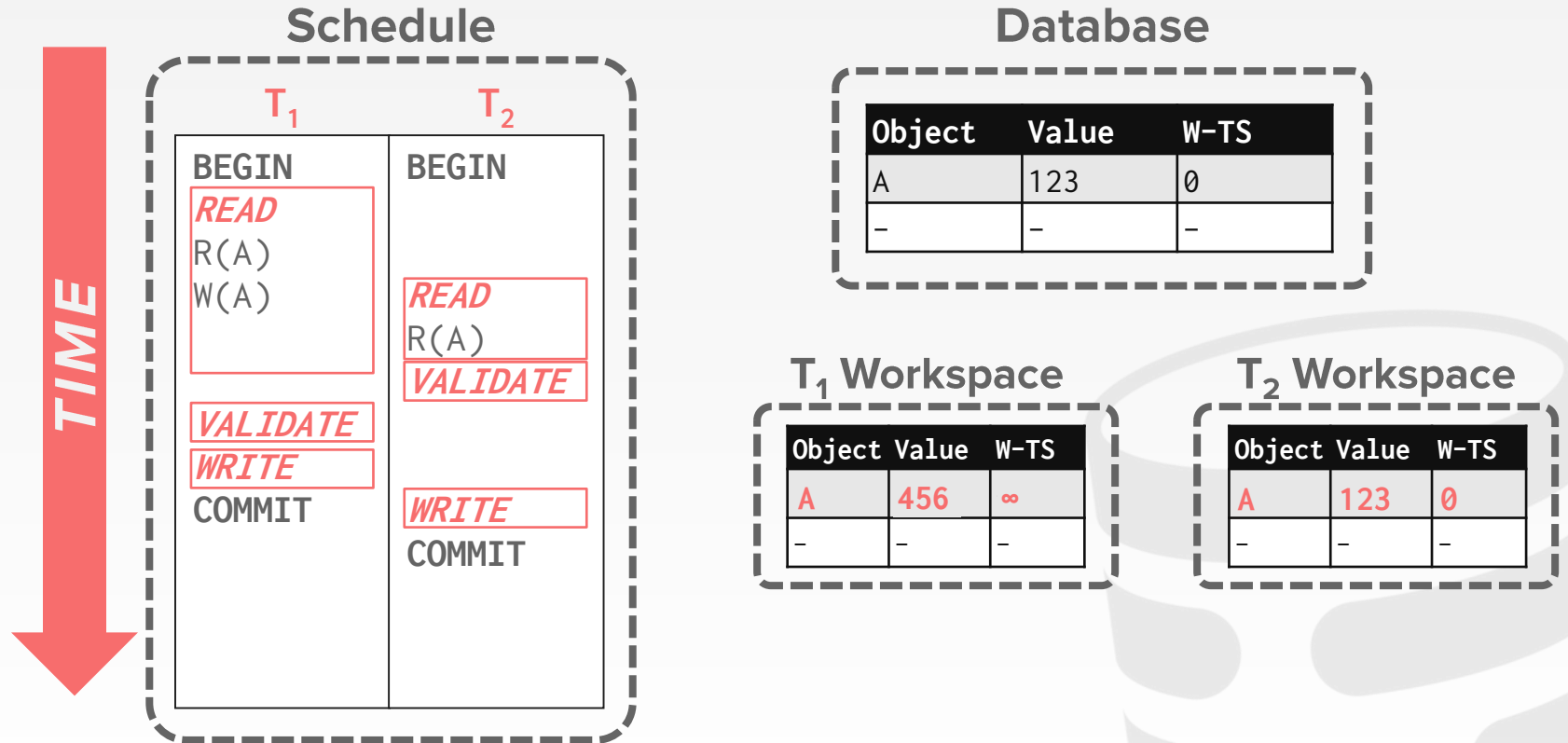| Object | Value | W-TS |
|--------|-------|------|
| A | 123 | 0 |
| – | – | – |

CARNEGIE MELLON
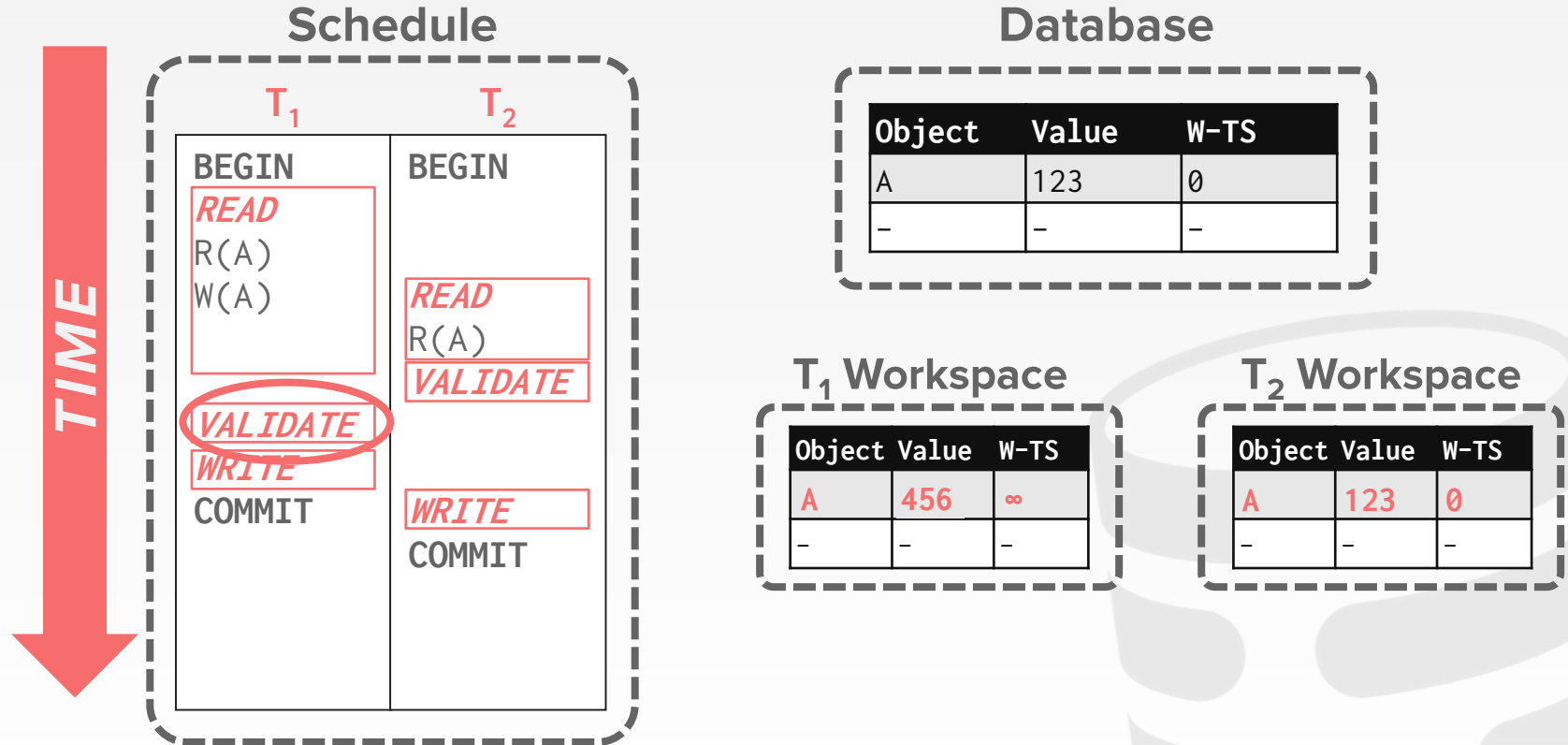DATABASE GROUP

# OCC – VALIDATION STEP #2
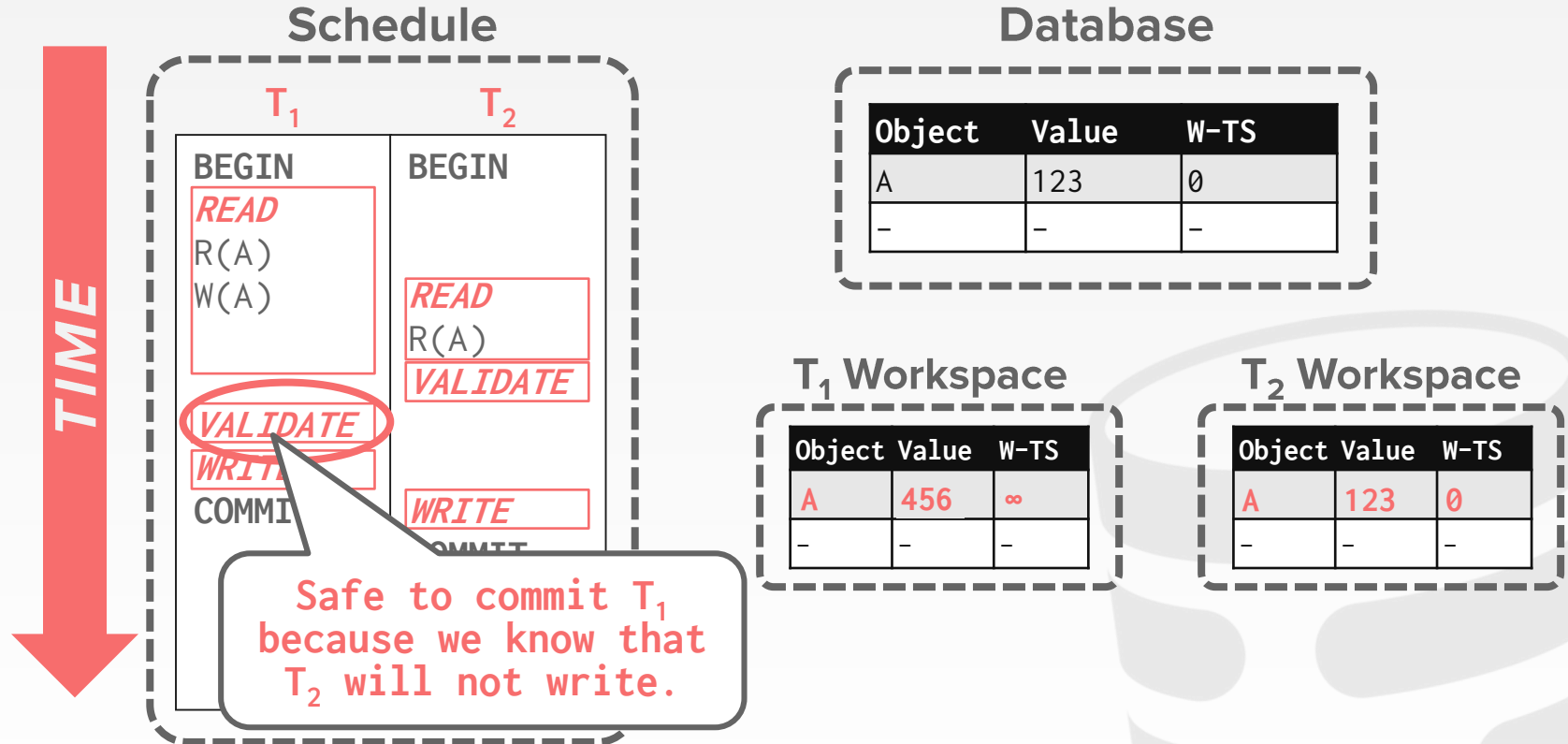
# OCC – VALIDATION STEP #2

# OCC – VALIDATION STEP #2

# OCC – VALIDATION STEP #2

# OCC – VALIDATION STEP #3

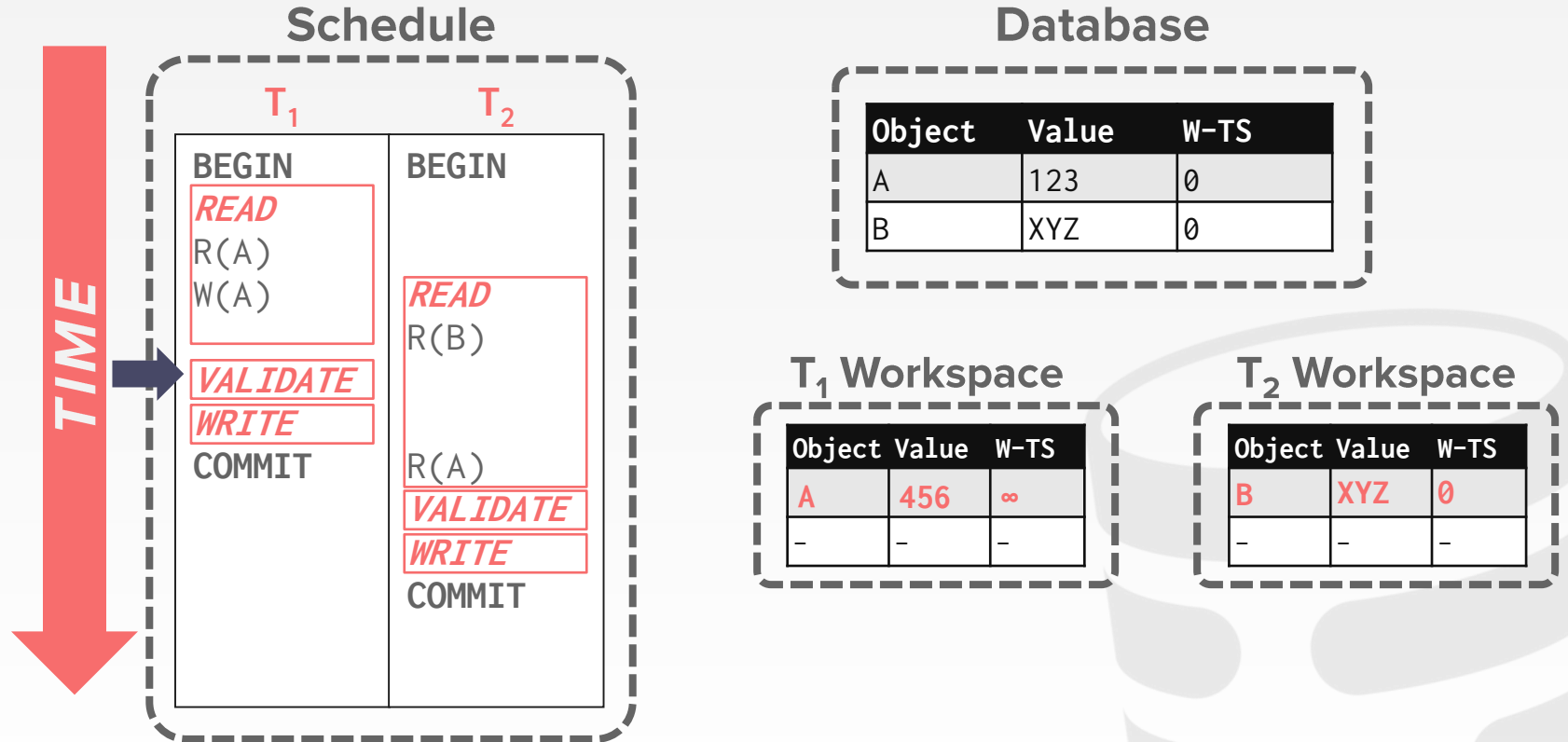$T_i$ completes its **Read** phase before $T_j$ completes its **Read** phase

And $T_i$ does not write to any object that is either read or written by $T_j$:

→ $WriteSet(T_i) \cap ReadSet(T_j) = \varnothing$
→ $WriteSet(T_i) \cap WriteSet(T_j) = \varnothing$

CARNEGIE MELLON
**DATABASE GROUP**

# OCC – VALIDATION STEP #3



**Schedule**

**Database**

| Object | Value | W-TS |
|--------|-------|------|
| A | 123 | 0 |
| B | XYZ | 0 |

**T₁ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| A | 456 | ∞ |
| – | – | – |

**T₂ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| B | XYZ | 0 |
| – | – | – |

CARNEGIE MELLON
DATABASE GROUP

# OCC – VALIDATION STEP #3

## Schedule



**T₁**
```
BEGIN
READ
R(A)
W(A)

VALIDATE
WRITE
COMMIT
```

**T₂**
```
BEGIN

READ
R(B)

R(A)
VALIDATE
WRITE
COMMIT
```

TIME

## Database

| Object | Value | W-TS |
|--------|-------|------|
| A | 456 | 1 |
| B | XYZ | 0 |

## T₁ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| A | 456 | ∞ |
| – | – | – |

## T₂ Workspace

| Object | Value | W-TS |
|--------|-------|------|
| B | XYZ | 0 |
| – | – | – |

CARNEGIE MELLON
DATABASE GROUP

# OCC – VALIDATION STEP #3

# OCC – VALIDATION STEP #3

**Schedule**



**Database**

| Object | Value | W-TS |
|--------|-------|------|
| A | 456 | 1 |
| B | XYZ | 0 |

**T₁ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| – | – | – |
| – | – | – |

**T₂ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| B | XYZ | 0 |
| – | – | – |

CARNEGIE MELLON
DATABASE GROUP

# OCC – VALIDATION STEP #3

**Schedule**

**Database**

TIME

|       | $T_1$ | $T_2$ |
|-------|-------|-------|
|       | BEGIN | BEGIN |
|       | *READ* |       |
|       | R(A)  |       |
|       | W(A)  |       |
|       |       | *READ* |
|       |       | R(B)  |
|       | *VALIDATE* |   |
|       | *WRITE* |     |
|       | COMMIT |      |
|       |       | R(A)  |
|       |       | *VALIDATE* |
|       |       | *WRITE* |
|       |       | COMMIT |

| Object | Value | W-TS |
|--------|-------|------|
| A      | 1     | 1    |
| B      | XYZ   | 0    |

**$T_1$ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| –      | –     | –    |
| –      | –     | –    |

**$T_2$ Workspace**

| Object | Value | W-TS |
|--------|-------|------|
| B      | XYZ   | 0    |
| A      | 456   | 1    |

CARNEGIE MELLON
DATABASE GROUP

# OCC – OBSERVATIONS

**When does OCC work well?**

When # of conflicts is low:
→ All txns are read-only (ideal).
→ Txns access disjoint subsets of data.

If the database is large and the workload is not skewed, then there is a low probability of conflict, so again locking is wasteful.

CARNEGIE MELLON
DATABASE GROUP

# OCC – PERFORMANCE ISSUES

High overhead for copying data locally.

Validation/Write phase bottlenecks.

Aborts are more wasteful because they only occur after a txn has already executed.

Suffers from timestamp allocation bottleneck.

CARNEGIE MELLON
DATABASE GROUP

# OBSERVATION

When a txn commits, all previous T/O
schemes check to see whether there is
a conflict with concurrent txns.
→ This requires latches.

If you have a lot of concurrent txns,
then this is slow even if the conflict rate
is low.

CARNEGIE MELLON
DATABASE GROUP

# PARTITION-BASED T/O

Split the database up in disjoint subsets called **partitions** (aka shards).

Only check for conflicts between txns that are running in the same partition.

# DATABASE PARTITIONING

**Schema**

# DATABASE PARTITIONING

## Schema



```
CREATE TABLE WAREHOUSE (
  w_id INT PRIMARY KEY,
  w_name VARCHAR UNIQUE,
  ⋮
);
```

```
CREATE TABLE DISTRICT (
  d_id INT,
  d_w_id INT REFERENCES
         ↪ WAREHOUSE (w_id),
  ⋮
  PRIMARY KEY (d_w_id, d_id)
);
```

CARNEGIE MELLON
DATABASE GROUP

# DATABASE PARTITIONING



Schema

Schema Tree

CARNEGIE MELLON
DATABASE GROUP
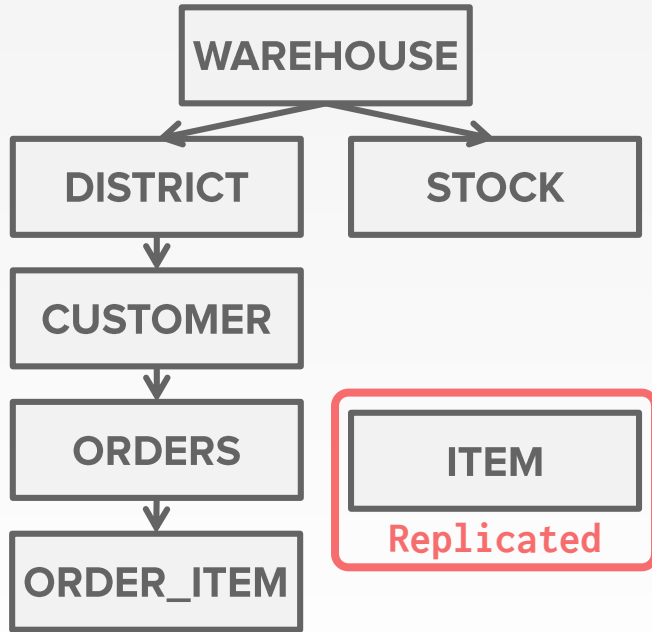
# DATABASE PARTITIONING



Schema

Schema Tree

CARNEGIE MELLON
DATABASE GROUP

# DATABASE PARTITIONING

# DATABASE PARTITIONING

Partitions

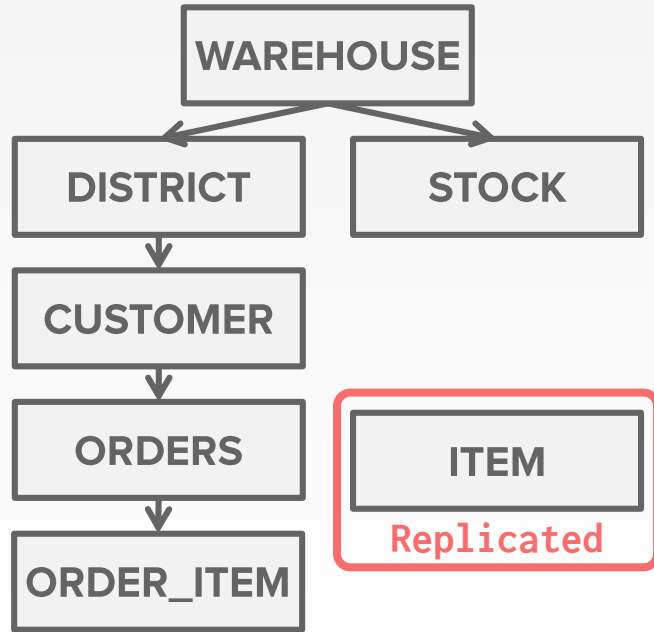| P1 | P2 | P3 | P4 |

| P1 | P2 | P3 | P4 | | P1 | P2 | P3 | P4 |

| P1 | P2 | P3 | P4 |

| P1 | P2 | P3 | P4 |

ITEM

Replicated

| P1 | P2 | P3 | P4 |

# DATABASE PARTITIONING



**Partitions**

WAREHOUSE

DISTRICT STOCK

CUSTOMER

ORDERS ITEM
Replicated

ORDER_ITEM

P1 P2
P4 P3

CARNEGIE MELLON
**DATABASE GROUP**
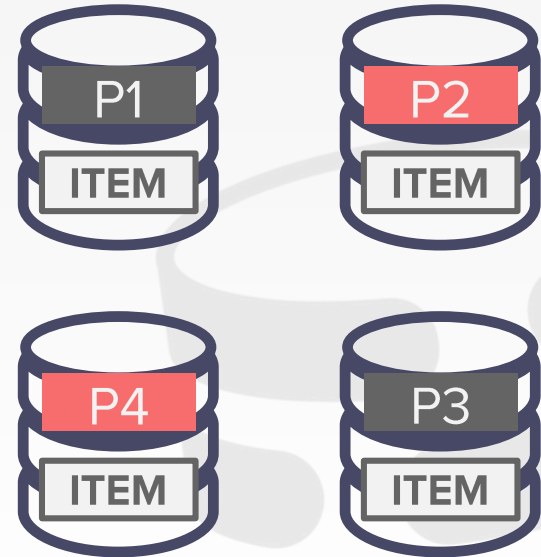
# DATABASE PARTITIONING

# PARTITION-BASED T/O

Txns are assigned timestamps based on when they arrive at the DBMS.

Partitions are protected by a single lock:
→ Each txn is queued at the partitions it needs.
→ The txn acquires a partition's lock if it has the lowest timestamp in that partition's queue.
→ The txn starts when it has all of the locks for all the partitions that it will read/write.

CARNEGIE MELLON
DATABASE GROUP

# PARTITION-BASED T/O – READS

Txns can read anything that they want at the partitions that they have locked.

If a txn tries to access a partition that it does not have the lock, it is aborted + restarted.

# PARTITION-BASED T/O — WRITES

All updates occur in place.
→ Maintain a separate in-memory buffer to undo changes if the txn aborts.

If a txn tries to write to a partition that it does not have the lock, it is aborted + restarted.

CARNEGIE MELLON
DATABASE GROUP

# PARTITION-BASED T/O – PERFORMANCE ISSUES

Partition-based T/O protocol is very fast if:
→ The DBMS knows what partitions the txn needs before it starts.
→ Most (if not all) txns only need to access a single partition.

Multi-partition txns causes partitions to be idle while txn executes.

# CONCLUSION

Every concurrency control can be broken down into the basic concepts that I've described in the last two lectures.

I'm not showing benchmark results because I don't want you to get the wrong idea.

CARNEGIE MELLON
**DATABASE GROUP**

# NEXT CLASS

Multi-Version Concurrency Control

CARNEGIE MELLON
**DATABASE GROUP**