# ADMINISTRIVIA

Project 2 due tonight!

Project 3 released today

Final exam Monday, May 1st, 8:30 – 11:30 a.m.

# LAST TIME: CONCURRENCY CONTROL

**Atomicity**

**Consistency**

**Isolation**
→ Serial execution schedules
→ Serializable

       Conflict serializable

       View serializable

**Durability** Linearizable

       Strict serializable

# CONCURRENCY CONTROL CONCLUSIONS

Concurrency control and recovery are among the most important functions provided by a DBMS.

# CONCURRENCY CONTRO

Concurrency control and reco
most important functions pro

ability problems that it brings [9, 10, 19]. We believe it is better to have application programmers deal with performance problems due to overuse of transactions as bottlenecks arise, rather than always coding around the lack of transactions. Running two-phase commit over Paxos

**Spanner: Google's Globally-Distributed Database**

James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, JJ Furman,
Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh,
Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura,
David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak,
Christopher Taylor, Ruth Wang, Dale Woodford

Google, Inc.

**Abstract**

Spanner is Google's scalable, multi-version, globally-distributed, and synchronously-replicated database. It is the first system to distribute data at global scale and support externally-consistent distributed transactions. This paper describes how Spanner is structured, its feature set, the rationale underlying various design decisions, and a novel time API that exposes clock uncertainty. This API and its implementation are critical to supporting external consistency and a variety of powerful features: non-blocking reads in the past, lock-free read-only transactions, and atomic schema changes, across all of Spanner.

**1 Introduction**

tency over higher availability, as long as they can survive 1 or 2 datacenter failures.

Spanner's main focus is managing cross-datacenter replicated data, but we have also spent a great deal of time in designing and implementing important database features on top of our distributed-systems infrastructure. Even though many projects happily use Bigtable [9], we have also consistently received complaints from users that Bigtable can be difficult to use for some kinds of applications: those that have complex, evolving schemas, or those that want strong consistency in the presence of wide-area replication. (Similar claims have been made by other authors [37].) Many applications at Google have chosen to use Megastore [5] because of its semi-relational data model and support for synchronous replication, despite its relatively poor write throughput. As a consequence, Spanner has evolved from a Bigtable-like key-value store into a temporal multi-version database. Data is stored in schematized semi-relational tables; data is versioned, and each version is automatically timestamped with its commit time; old versions of data are subject to configurable garbage-collection policies, and applications can read data at old timestamps. Spanner supports general-purpose transactions, and provides a SQL-based query language.

As a globally-distributed database, Spanner provides several interesting features. First, the replication configurations for data can be dynamically controlled at a fine grain by applications. Applications can specify constraints to control which datacenters contain which data, how far data is from its users (to control read latency), how far replicas are from each other (to control write latency), and how many replicas are maintained (to control durability, availability, and read performance). Data can also be dynamically and transparently moved between datacenters by the system to balance resource usage across datacenters. Second, Spanner has two features that are difficult to implement in a distributed database: it
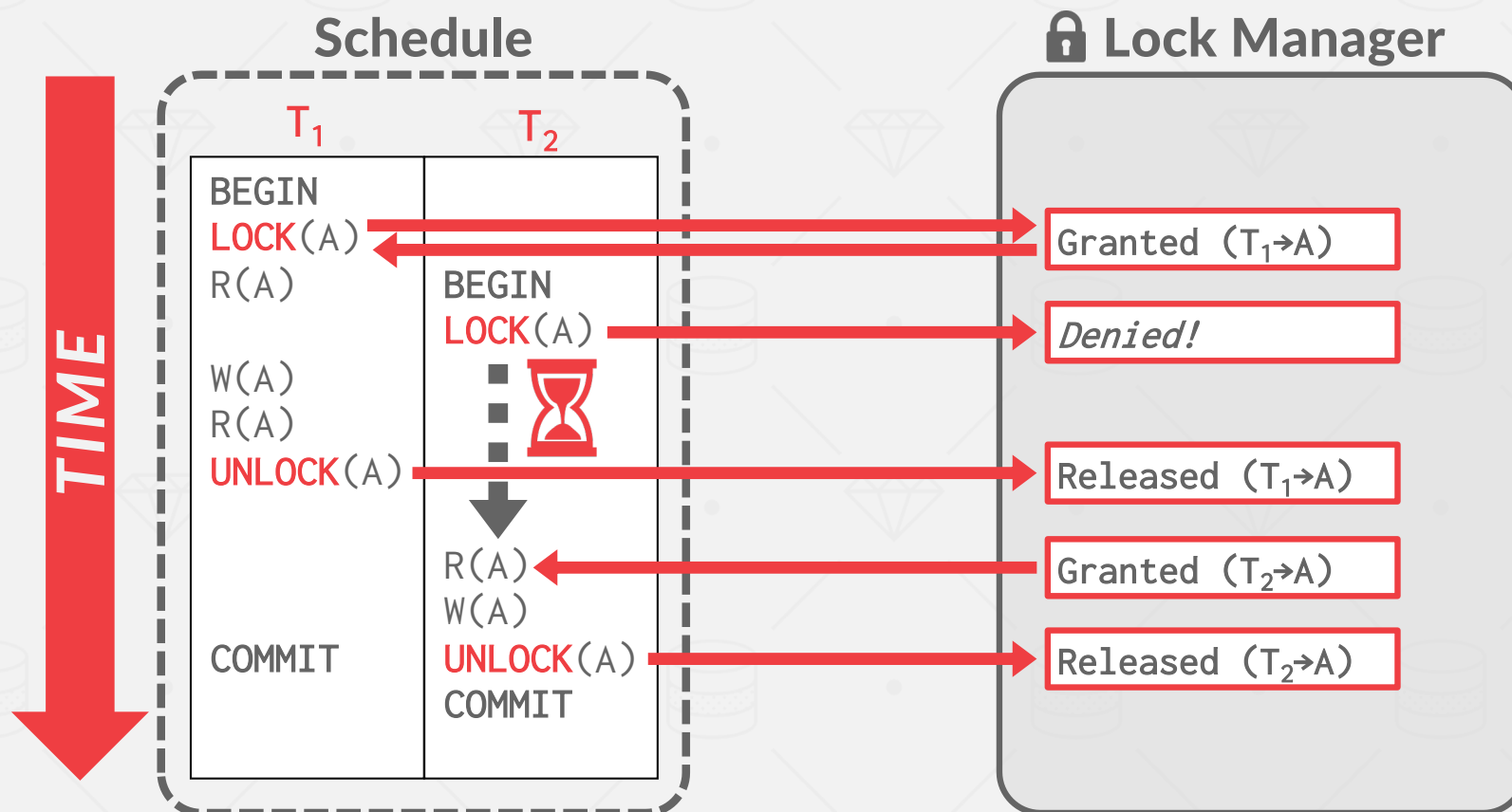
# OBSERVATION

We need a way to guarantee that all execution schedules are correct (i.e., serializable) without knowing the entire schedule ahead of time.

One solution: Use locks to protect database objects.
→ System automatically locks & unlocks objects as needed
→ Ensures that resulting execution is equivalent to some serial execution order

# EXECUTING WITH LOCKS

**Schedule**

🔒 **Lock Manager**

**TIME**

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

Granted ($T_1 \rightarrow$ A)

Denied!

Released ($T_1 \rightarrow$ A)

Granted ($T_2 \rightarrow$ A)

Released ($T_2 \rightarrow$ A)

# LOCKS VS. LATCHES

|  | *Locks* | *Latches* |
|---|---|---|
| **Separate…** | User transactions | Threads |
| **Protect…** | Database Contents | In-Memory Data Structures |
| **During…** | Entire Transactions | Critical Sections |
| **Modes…** | Shared, Exclusive, Update, Intention | Read, Write |
| **Deadlock** | Detection & Resolution | Avoidance |
| **…by…** | Waits-for, Timeout, Aborts | Coding Discipline |
| **Kept in…** | Lock Manager | Protected Data Structure |

Source: Goetz Graefe

# TODAY'S AGENDA

Lock Types

Two-Phase Locking

Deadlock Detection + Prevention

Hierarchical Locking

# BASIC LOCK TYPES

**S-LOCK**: Shared locks for reads.

**X-LOCK**: Exclusive locks for writes.

### Compatibility Matrix

|           | Shared | Exclusive |
|-----------|:------:|:---------:|
| Shared    | ✔      | X         |
| Exclusive | X      | X         |

## Compatibility of lock modes

The following table shows the compatibility of any two modes for page and row locks. No question of compatibility arises between page and row locks, because a partition or table space cannot use both page and row locks.

**Table 1. Compatibility matrix of page lock and row lock modes**

| Lock mode | Share (S-lock) | Update (U-lock) |
|---|---|---|
| Share (S-lock) | Yes | Yes |
| Update (U-lock) | Yes | No |
| Exclusive (X-lock) | | |

Compatibility for table space locks ...
modes for partition, table space, or ...

**Table 2. Compatibility of table and ...**

| Lock Mode | IS | IX | S |
|---|---|---|---|
| IS | Yes | Yes | Yes |
| IX | Yes | Yes | No |
| S | Yes | No | Yes |
| U | Yes | No | Yes |
| SIX | Yes | No | No |
| X | No | No | No |

**Existing granted mode**

| Requested mode | IS | S | U |
|---|---|---|---|
| Intent shared (IS) | Yes | Yes | |
| Shared (S) | Yes | Yes | |
| Update (U) | Yes | No | |
| Intent exclusive (IX) | Yes | No | |
| Shared with intent exclusive (SIX) | No | No | |

**Table 13-3 Summary of Table Locks**

| SQL Statement | Mode of Table Lock | Lock Modes Permitted? | | | | |
|---|---|---|---|---|---|---|
| | | RS | RX | S | SRX | X |
| SELECT...FROM table... | none | Y | Y | Y | Y | Y |
| INSERT INTO table ... | RX | Y | Y | N | N | N |
| UPDATE table ... | RX | Y* | Y* | N | N | N |
| DELETE FROM table ... | RX | Y* | Y* | N | N | N |
| SELECT ... FROM table FOR UPDATE OF ... | RS | Y* | Y* | Y | Y | N |
| LOCK TABLE table IN ROW SHARE MODE | RS | Y | Y | Y | Y | N |
| LOCK TABLE table IN ROW EXCLUSIVE MODE | RX | Y | Y | N | N | N |
| LOCK TABLE table IN SHARE MODE | S | Y | N | Y | N | N |
| LOCK TABLE table IN SHARE ROW EXCLUSIVE MODE | SRX | Y | N | N | N | N |
| LOCK TABLE table IN EXCLUSIVE MODE | X | N | N | N | N | N |

**Table 13.2. Conflicting Lock Modes**

| Requested Lock Mode | Existing Lock Mode | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | ACCESS SHARE | ROW SHARE | ROW EXCL. | SHARE UPDATE EXCL. | SHARE | SHARE ROW EXCL. | EXCL. | ACCESS EXCL. |
| ACCESS SHARE | | | | | | | | X |
| ROW SHARE | | | | | | | X | X |
| ROW EXCL. | | | | | X | X | X | X |
| SHARE UPDATE EXCL. | | | | X | X | X | X | X |
| SHARE | | | X | X | | X | X | X |
| SHARE ROW EXCL. | | | X | X | X | X | X | X |
| EXCL. | | X | X | X | X | X | X | X |
| ACCESS EXCL. | X | X | X | X | X | X | X | X |

## Table-level lock type compatibility is summarized in the following matrix

| | X | IX | S | IS |
|---|---|---|---|---|
| X | Conflict | Conflict | Conflict | Conflict |
| IX | Conflict | Compatible | Conflict | Compatible |
| S | Conflict | Conflict | Compatible | Compatible |
| IS | Conflict | Compatible | Compatible | Compatible |

# EXECUTING WITH LOCKS

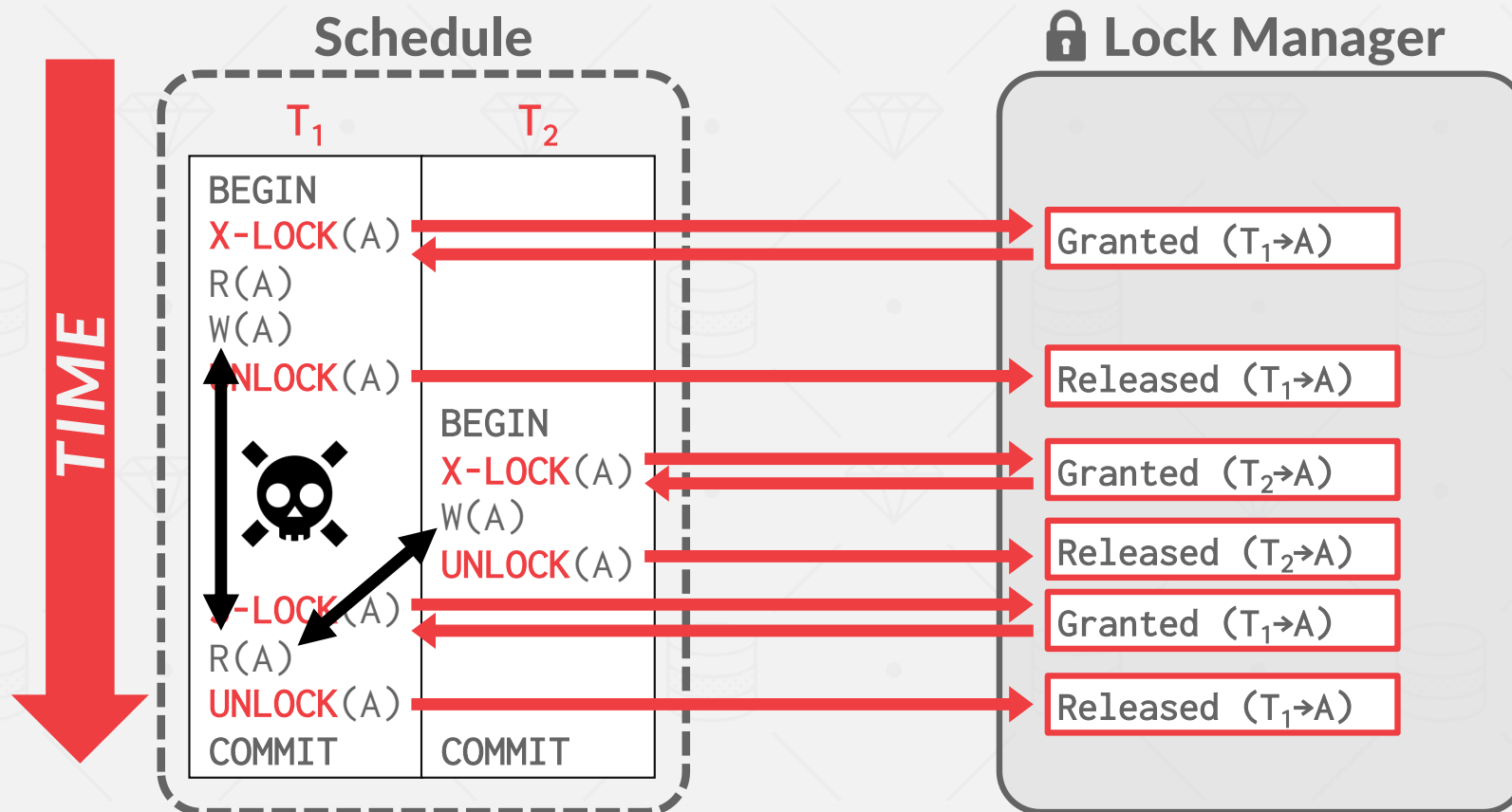Transactions request locks (or upgrades).

Lock manager grants or blocks requests.

Transactions release locks.

Lock manager updates its internal lock-table.
→ It keeps track of what transactions hold what locks and what transactions are waiting to acquire any locks.

# EXECUTING WITH LOCKS

**Schedule**                    🔒 **Lock Manager**

|  | $T_1$ | $T_2$ |
|---|---|---|
| BEGIN | | |
| X-LOCK(A) | | → Granted ($T_1$→A) |
| R(A) | | |
| W(A) | | |
| UNLOCK(A) | | → Released ($T_1$→A) |
| | BEGIN | |
| | X-LOCK(A) | → Granted ($T_2$→A) |
| | W(A) | |
| | UNLOCK(A) | → Released ($T_2$→A) |
| X-LOCK(A) | | → Granted ($T_1$→A) |
| R(A) | | |
| UNLOCK(A) | | → Released ($T_1$→A) |
| COMMIT | COMMIT | |

**TIME**

# CONCURRENCY CONTROL PROTOCOL

Two-phase locking (2PL) is a concurrency control protocol that determines whether a txn can access an object in the database at runtime.

The protocol does <u>not</u> need to know all the queries that a txn will execute ahead of time.
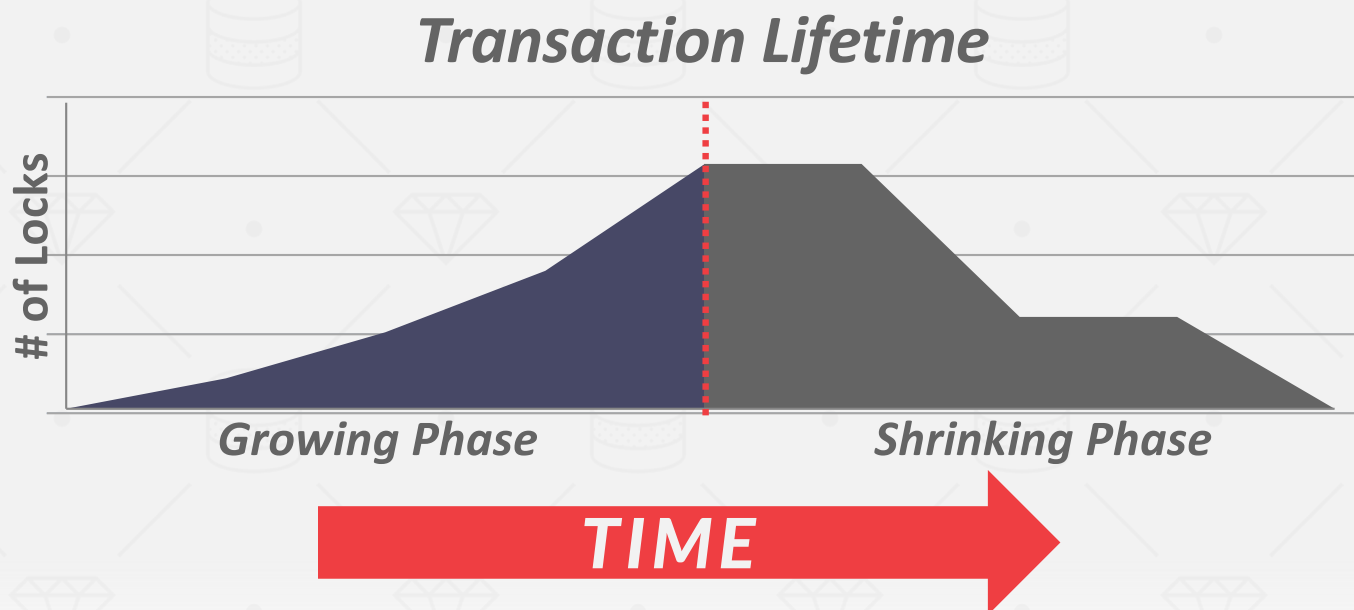
# TWO-PHASE LOCKING

**Phase #1: Growing**
→ Each txn requests the locks that it needs from the DBMS's lock manager.
→ The lock manager grants/denies lock requests.

**Phase #2: Shrinking**
→ The txn is allowed to only release/downgrade locks that it previously acquired. It cannot acquire new locks.
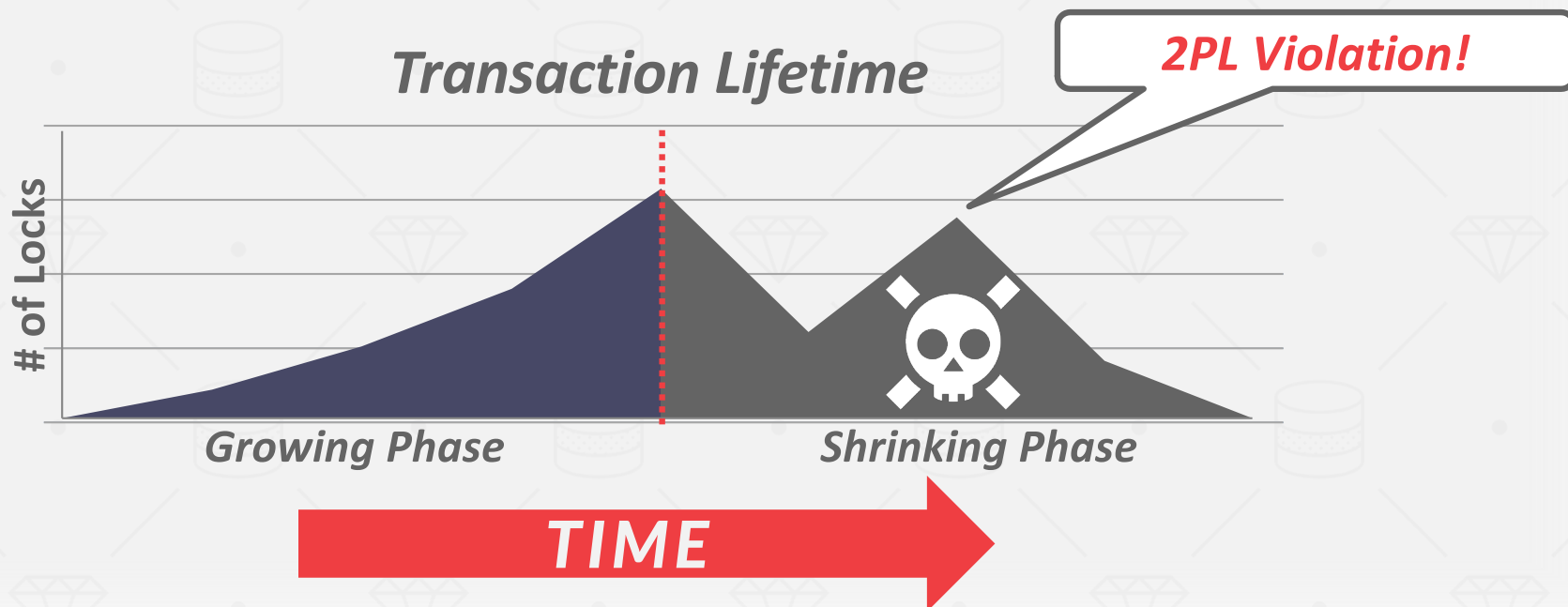
# TWO-PHASE LOCKING

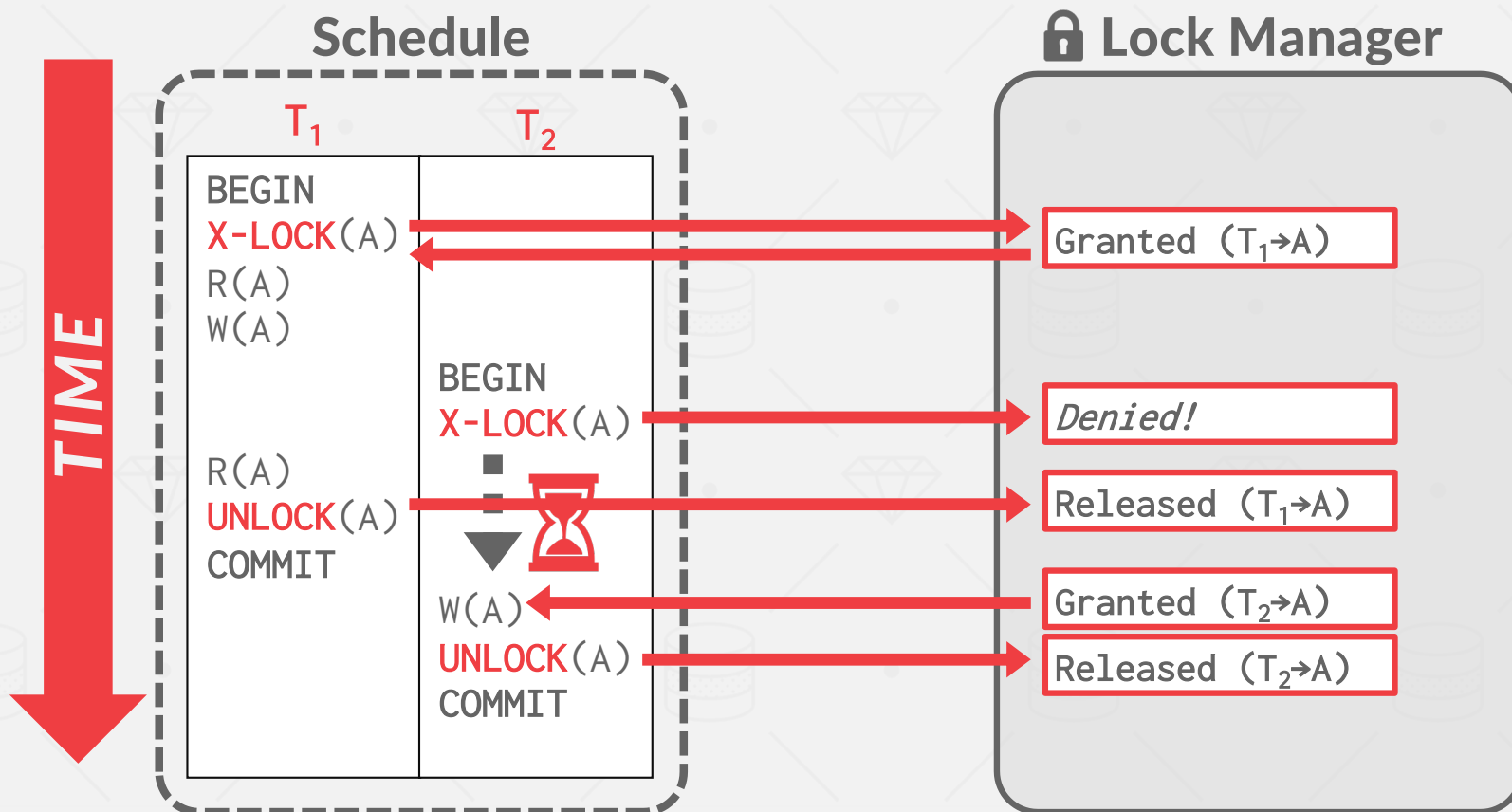The txn is not allowed to acquire/upgrade locks after the growing phase finishes.

**Transaction Lifetime**



# of Locks

*Growing Phase*          *Shrinking Phase*

**TIME**

# TWO-PHASE LOCKING

The txn is not allowed to acquire/upgrade locks after the growing phase finishes.
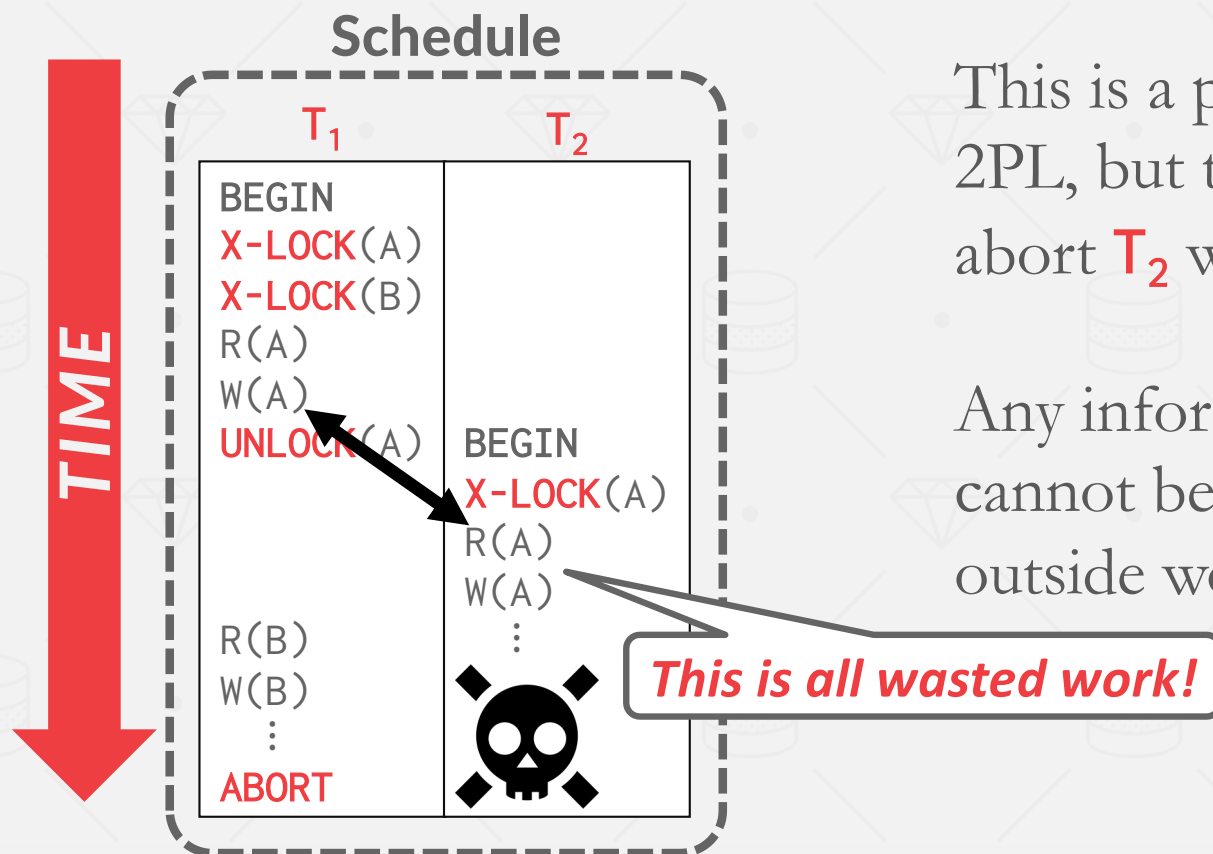
# EXECUTING WITH 2PL

**Schedule**

🔒 **Lock Manager**

|  | T₁ | T₂ |
| --- | --- | --- |

$T_1$     $T_2$

```
BEGIN
X-LOCK(A)
R(A)
W(A)

        BEGIN
        X-LOCK(A)
        ■
R(A)
UNLOCK(A)
COMMIT
        W(A)
        UNLOCK(A)
        COMMIT
```

Granted (T₁→A)

*Denied!*

Released (T₁→A)

Granted (T₂→A)

Released (T₂→A)

**TIME**

# TWO-PHASE LOCKING

2PL on its own is sufficient to guarantee conflict serializability because it generates schedules whose precedence graph is acyclic.

But it is subject to **cascading aborts**.

# 2PL – CASCADING ABORTS

**Schedule**

**TIME**

| $T_1$ | $T_2$ |
|-------|-------|
| BEGIN | |
| X-LOCK(A) | |
| X-LOCK(B) | |
| R(A) | |
| W(A) | |
| UNLOCK(A) | BEGIN |
| | X-LOCK(A) |
| | R(A) |
| | W(A) |
| R(B) | ⋮ |
| W(B) | |
| ⋮ | |
| ABORT | |

*This is all wasted work!*

This is a permissible schedule in 2PL, but the DBMS has to also abort $T_2$ when $T_1$ aborts.

Any information about $T_1$ cannot be "leaked" to the outside world.

# 2PL OBSERVATIONS

There are potential schedules that are serializable but would not be allowed by 2PL because locking limits concurrency.
→ Most DBMSs prefer correctness before performance.

May still have "dirty reads".
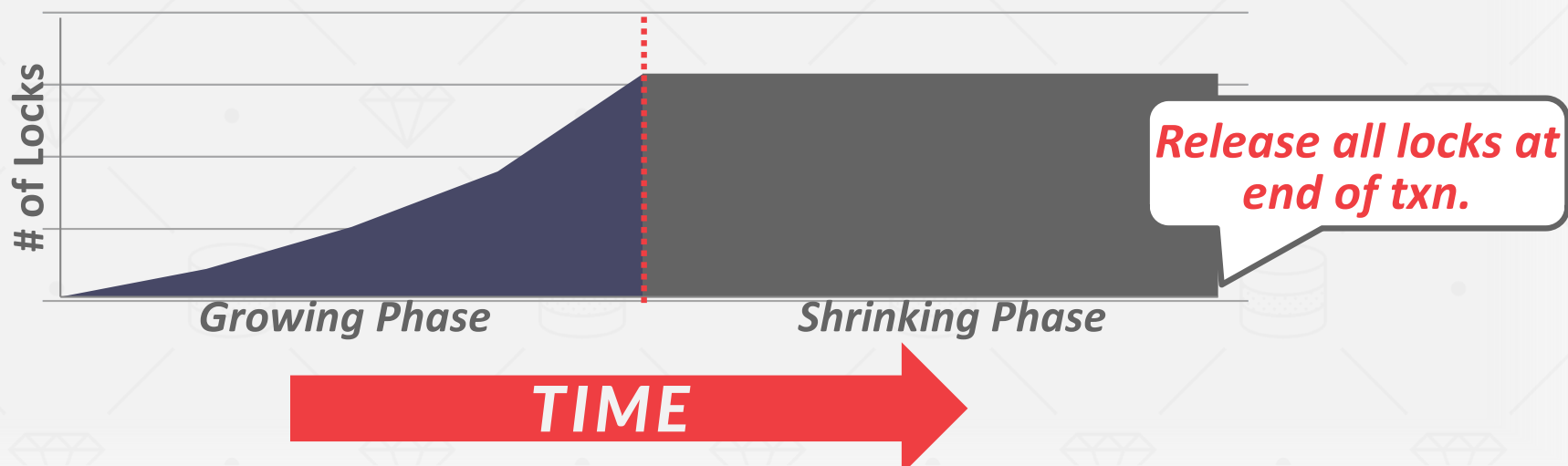→ Solution: **Strong Strict 2PL (aka Rigorous 2PL)**

May lead to deadlocks.
→ Solution: **Detection** or **Prevention**

# STRONG STRICT TWO-PHASE LOCKING

The txn is only allowed to release locks after it has ended (i.e., committed or aborted).

Allows only conflict serializable schedules, but it is often stronger than needed for some apps.



*Release all locks at end of txn.*

**# of Locks**

*Growing Phase*    *Shrinking Phase*

**TIME**

# STRONG STRICT TWO-PHASE LOCKING

A schedule is **strict** if a value written by a txn is not read or overwritten by other txns until that txn finishes.

Advantages:
→ Does not incur cascading aborts.
→ Aborted txns can be undone by just restoring original values of modified tuples.

# EXAMPLES

$T_1$ – Move $100 from Andy's account ($A$) to his bookie's account ($B$).

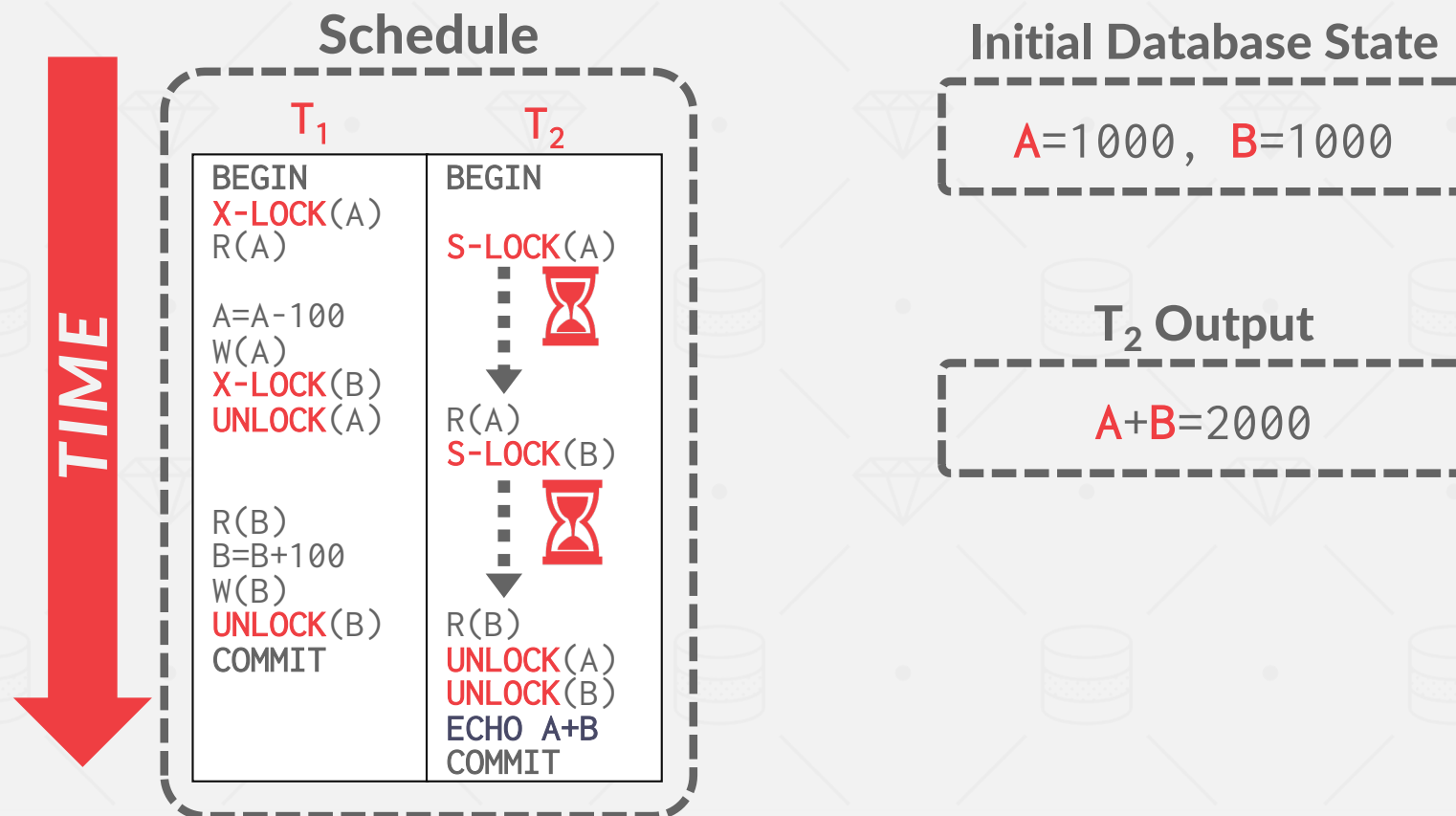$T_2$ – Compute the total amount in all accounts and return it to the application.
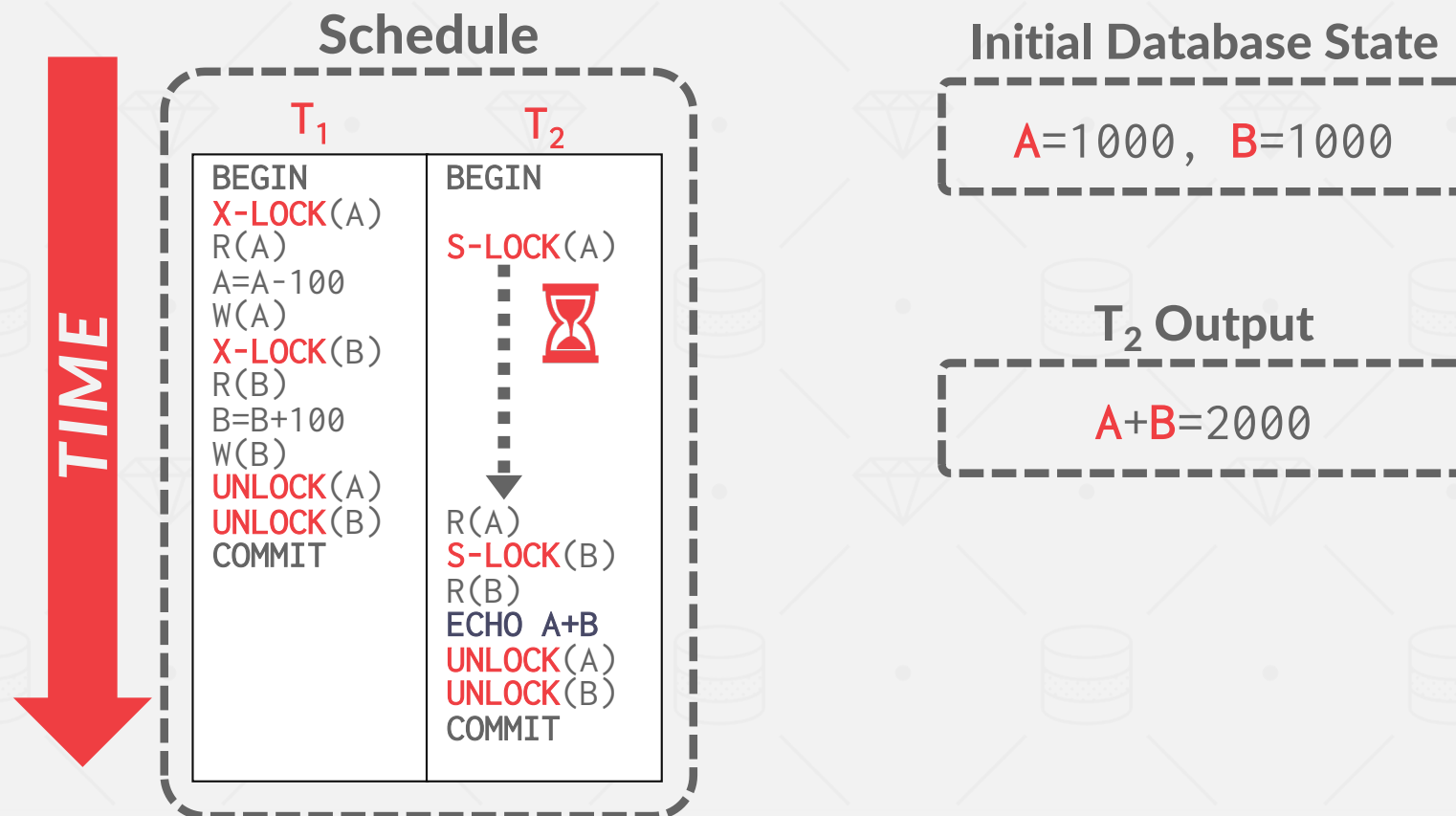
$T_1$

```
BEGIN
A=A-100
B=B+100
COMMIT
```
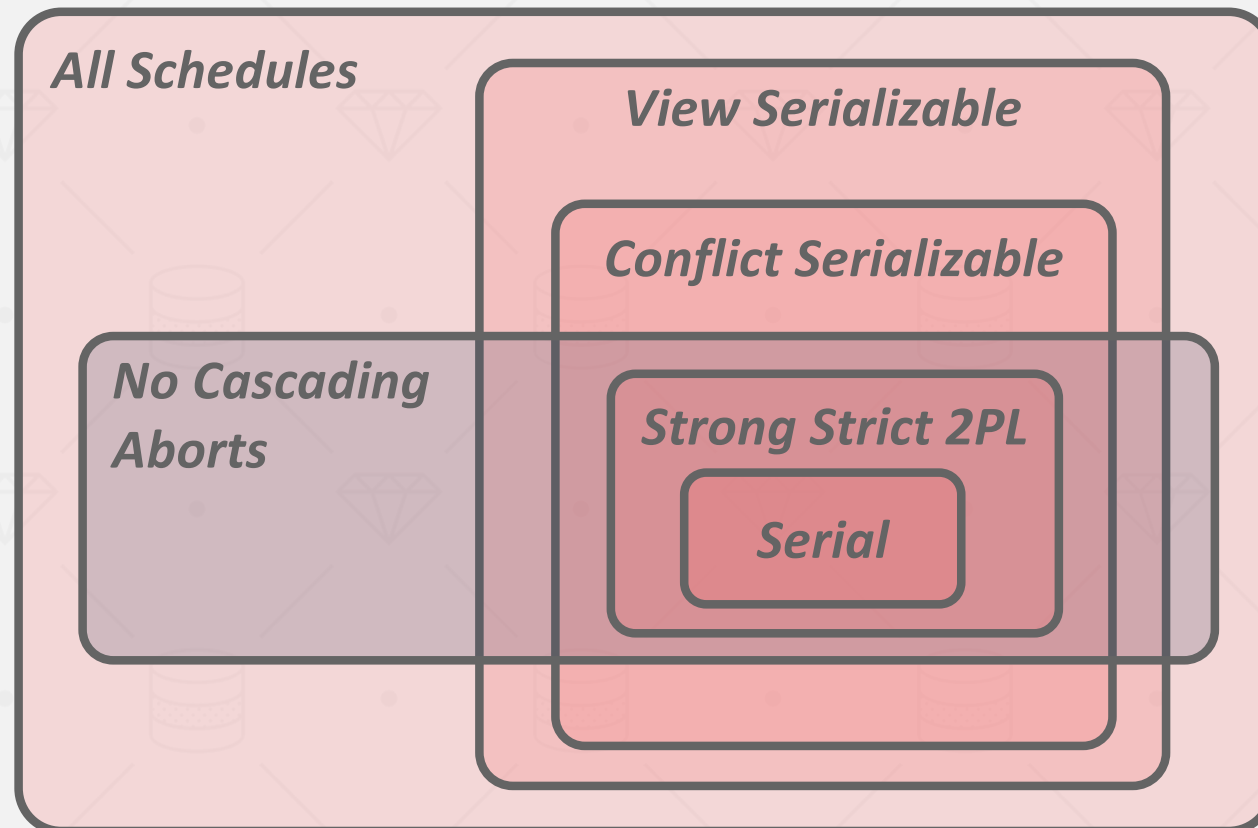
$T_2$

```
BEGIN
ECHO A+B
COMMIT
```

# 2PL EXAMPLE

## Schedule

**TIME**

| $T_1$ | $T_2$ |
|---|---|
| BEGIN | BEGIN |
| X-LOCK(A) | |
| R(A) | S-LOCK(A) |
| | ⧖ |
| A=A-100 | |
| W(A) | |
| X-LOCK(B) | |
| UNLOCK(A) | R(A) |
| | S-LOCK(B) |
| | ⧖ |
| R(B) | |
| B=B+100 | |
| W(B) | |
| UNLOCK(B) | R(B) |
| COMMIT | UNLOCK(A) |
| | UNLOCK(B) |
| | ECHO A+B |
| | COMMIT |

## Initial Database State

A=1000, B=1000

## $T_2$ Output

A+B=2000

# STRONG STRICT 2PL EXAMPLE

**Schedule**

| $T_1$ | $T_2$ |
|---|---|
| BEGIN | BEGIN |
| X-LOCK(A) | |
| R(A) | S-LOCK(A) |
| A=A-100 | |
| W(A) | |
| X-LOCK(B) | |
| R(B) | |
| B=B+100 | |
| W(B) | |
| UNLOCK(A) | |
| UNLOCK(B) | R(A) |
| COMMIT | S-LOCK(B) |
| | R(B) |
| | ECHO A+B |
| | UNLOCK(A) |
| | UNLOCK(B) |
| | COMMIT |

**TIME**

**Initial Database State**

A=1000, B=1000

**$T_2$ Output**

A+B=2000

# UNIVERSE OF SCHEDULES



All Schedules

View Serializable

Conflict Serializable

No Cascading Aborts

Strong Strict 2PL

Serial

# 2PL OBSERVATIONS

There are potential schedules that are serializable but would not be allowed by 2PL because locking limits concurrency.
→ Most DBMSs prefer correctness before performance.

May still have "dirty reads".
→ Solution: **Strong Strict 2PL (aka Rigorous 2PL)**

May lead to deadlocks.
→ Solution: **Detection** or **Prevention**

# IT JUST GOT REAL, SON

**Schedule**

🔒 **Lock Manager**

**TIME**

| $T_1$ | $T_2$ |
|-------|-------|
| BEGIN | BEGIN |
| X-LOCK(A) | |
| | S-LOCK(B) |
| | R(B) |
| | S-LOCK(A) |
| R(A) | |
| X-LOCK(B) | |

Granted ($T_1$→A)

Granted ($T_2$→B)

*Denied!*

*Denied!*

# 2PL DEADLOCKS

A **deadlock** is a cycle of transactions waiting for locks to be released by each other.

Two ways of dealing with deadlocks:
→ **Approach #1: Deadlock Detection**
→ **Approach #2: Deadlock Prevention**

# DEADLOCK DETECTION

The DBMS creates a **<u>waits-for</u>** graph to keep track of what locks each txn is waiting to acquire:
→ Nodes are transactions
→ Edge from $T_i$ to $T_j$ if $T_i$ is waiting for $T_j$ to release a lock.

The system periodically checks for cycles in *waits-for* graph and then decides how to break it.

# DEADLOCK DETECTION

## Schedule

|  | $T_1$ | $T_2$ | $T_3$ |
|---|---|---|---|
|  | BEGIN | BEGIN | BEGIN |
|  | S-LOCK(A) | | |
|  | | X-LOCK(B) | |
|  | | | S-LOCK(C) |
|  | S-LOCK(B) | | |
|  | | X-LOCK(C) | |
|  | | | X-LOCK(A) |

## *Waits-For* Graph

# DEADLOCK HANDLING

When the DBMS detects a deadlock, it will select a "victim" txn to rollback to break the cycle.

The victim txn will either restart or abort (more common) depending on how it was invoked.

There is a trade-off between the frequency of checking for deadlocks and how long txns wait before deadlocks are broken.

# DEADLOCK HANDLING: VICTIM SELECTION

Selecting the proper victim depends on a lot of different variables….
→ By age (lowest timestamp)
→ By progress (least/most queries executed)
→ By the # of items already locked
→ By the # of txns that we have to rollback with it

We also should consider the # of times a txn has been restarted in the past to prevent starvation.

# DEADLOCK HANDLING: ROLLBACK LENGTH

After selecting a victim txn to abort, the DBMS can also decide on how far to rollback the txn's changes.

**Approach #1: Completely**
→ Rollback entire txn and tell the application it was aborted.

**Approach #2: Partial (Savepoints)**
→ DBMS rolls back a portion of a txn (to break deadlock) and then attempts to re-execute the undone queries.

# DEADLOCK PREVENTION

When a txn tries to acquire a lock that is held by another txn, the DBMS kills one of them to prevent a deadlock.

This approach does <u>not</u> require a *waits-for* graph or detection algorithm.

# DEADLOCK PREVENTION

Assign priorities based on timestamps:
→ Older Timestamp = Higher Priority (e.g., $T_1 > T_2$)

**Wait-Die ("Old Waits for Young")**
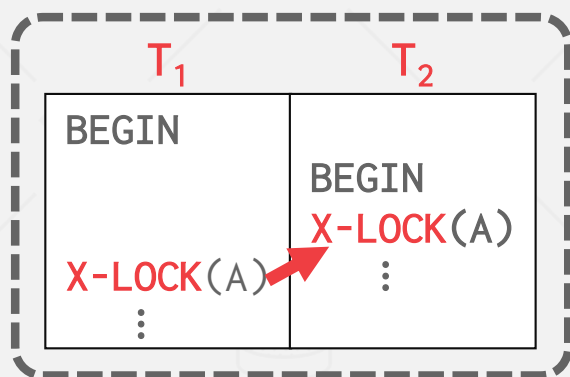→ If *requesting txn* has higher priority than *holding txn*, then *requesting txn* waits for *holding txn*.
→ Otherwise *requesting txn* aborts.

**Wound-Wait ("Young Waits for Old")**
→ If *requesting txn* has higher priority than *holding txn*, then *holding txn* aborts and releases lock.
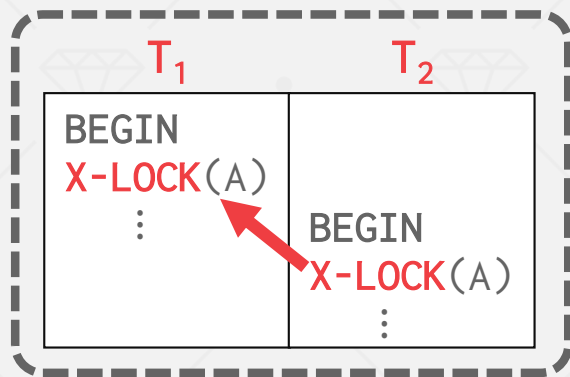→ Otherwise *requesting txn* waits.

# DEADLOCK PREVENTION



**Wait-Die**

$T_1$ waits

**Wound-Wait**

$T_2$ aborts

**Wait-Die**

$T_2$ aborts

**Wound-Wait**

$T_2$ waits

# DEADLOCK PREVENTION

*Why do these schemes guarantee no deadlocks?*

Only one "type" of direction allowed when waiting for a lock.

*When a txn restarts, what is its (new) priority?*

Its original timestamp to prevent it from getting starved for resources like an old man at a corrupt senior center.

# OBSERVATION

All these examples have a one-to-one mapping from database objects to locks.

If a txn wants to update one billion tuples, then it must acquire one billion locks.

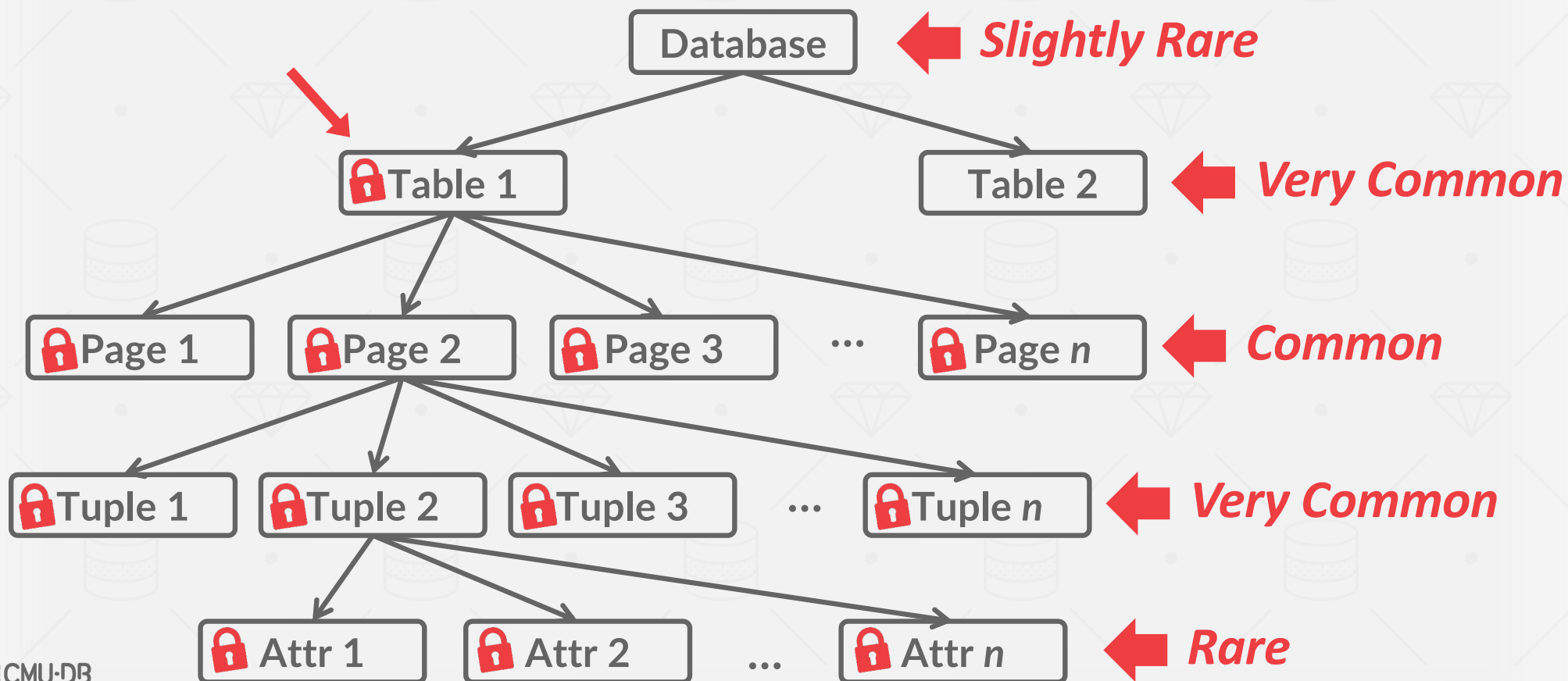Acquiring locks is a more expensive operation than acquiring a latch even if that lock is available.

# LOCK GRANULARITIES

When a txn wants to acquire a "lock", the DBMS can decide the granularity (i.e., scope) of that lock.
→ Attribute? Tuple? Page? Table?

The DBMS should ideally obtain fewest number of locks that a txn needs.

Trade-off between <u>parallelism</u> versus <u>overhead</u>.
→ Fewer Locks, Larger Granularity vs. More Locks, Smaller Granularity.

# DATABASE LOCK HIERARCHY



Database **← *Slightly Rare***

🔒Table 1 Table 2 **← *Very Common***

🔒Page 1 🔒Page 2 🔒Page 3 ⋯ 🔒Page *n* **← *Common***

🔒Tuple 1 🔒Tuple 2 🔒Tuple 3 ⋯ 🔒Tuple *n* **← *Very Common***

🔒 Attr 1 🔒 Attr 2 … 🔒 Attr *n* **← *Rare***

# INTENTION LOCKS

An **intention lock** allows a higher-level object to be locked in **shared** or **exclusive** mode without having to check all descendent objects.

If an object is locked in an intention mode, then some txn is doing explicit locking at a lower level.

# INTENTION LOCKS

## Intention-Shared (IS)

→ Indicates explicit locking at lower level with S locks.

→ Intent to get S lock(s) at finer granularity.

## Intention-Exclusive (IX)

→ Indicates explicit locking at lower level with X locks.

→ Intent to get X lock(s) at finer granularity.

## Shared+Intention-Exclusive (SIX)

→ The subtree rooted by that node is locked explicitly in S mode and explicit locking is being done at a lower level with X locks.

# LOCKING PROTOCOL

Each txn obtains appropriate lock at highest level of the database hierarchy.

To get **S** or **IS** lock, the txn must hold at least **IS** on parent.

To get **X**, **IX**, or **SIX** lock, must hold at least **IX** on parent.

# EXAMPLE

$T_1$ – Get the balance of Andy's shady off-shore bank account.

$T_2$ – Increase Chi's account balance by 6%.
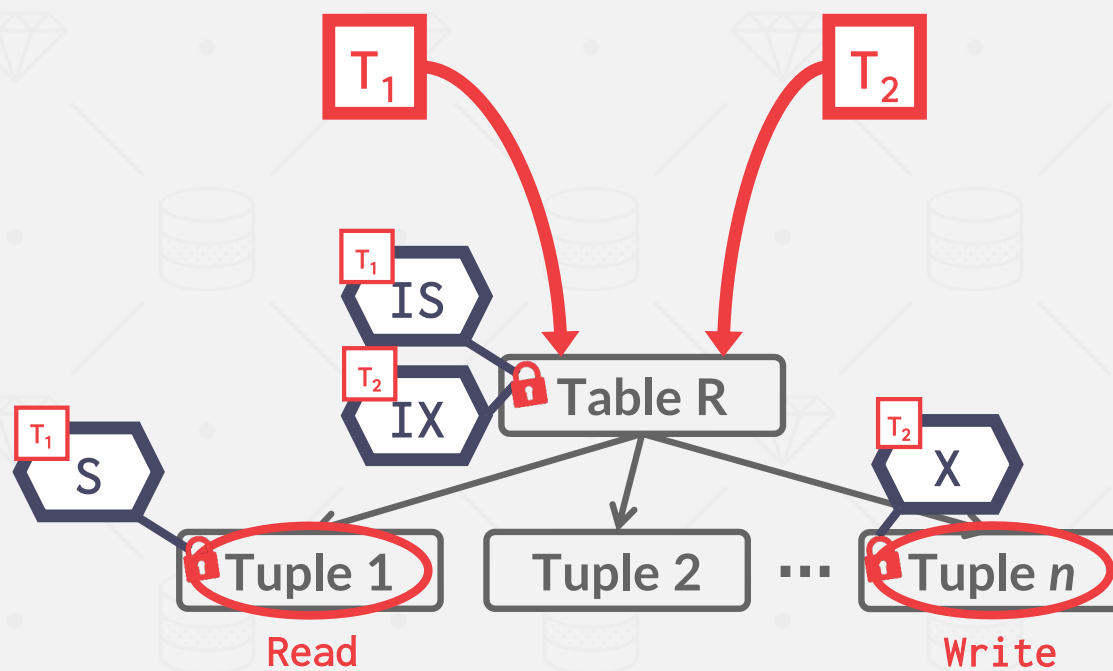
***What locks should these txns obtain?***
→ **Exclusive** + **Shared** for leaves of lock tree.
→ Special **Intention** locks for higher levels.

# EXAMPLE – TWO-LEVEL HIERARCHY
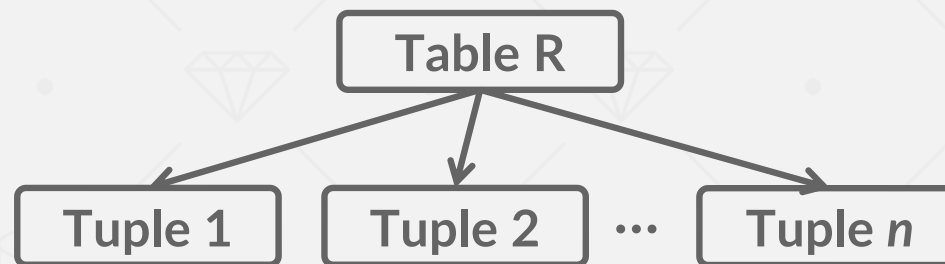


Read Andy's record in **R**.

Update Chi's record in **R**.

# EXAMPLE – THREE QUERIES

Assume three txns execute at same time:
→ $T_1$ – Scan all tuples in **R** and update one tuple.
→ $T_2$ – Read a single tuple in **R**.
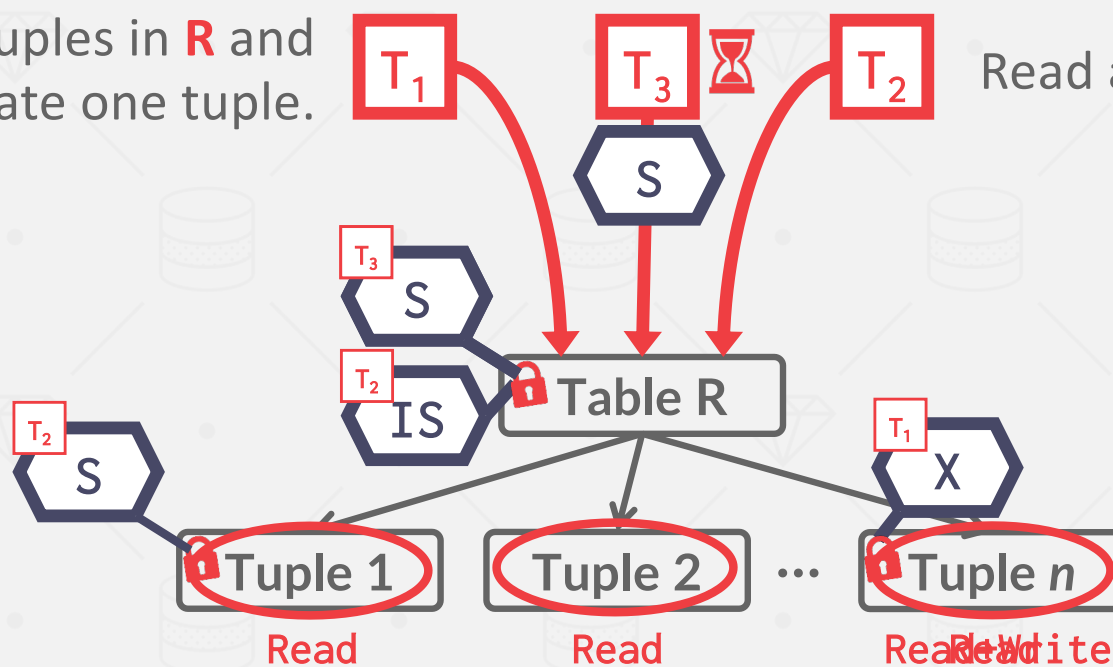→ $T_3$ – Scan all tuples in **R**.

# EXAMPLE – THREE QUERIES



Scan all tuples in **R**.

Scan all tuples in **R** and update one tuple.

Read a single tuple in **R**.

# COMPATIBILITY MATRIX

**T$_2$** Wants

**T$_1$** Holds

|     | IS | IX | S | SIX | X |
|-----|----|----|----|-----|----|
| IS  | ✓  | ✓  | ✓  | ✓   | ✗  |
| IX  | ✓  | ✓  | ✗  | ✗   | ✗  |
| S   | ✓  | ✗  | ✓  | ✗   | ✗  |
| SIX | ✓  | ✗  | ✗  | ✗   | ✗  |
| X   | ✗  | ✗  | ✗  | ✗   | ✗  |

# LOCK ESCALATION

The DBMS can automatically switch to coarser-grained locks when a txn acquires too many low-level locks.

This reduces the number of requests that the lock manager must process.

# LOCKING IN PRACTICE

Applications typically don't acquire a txn's locks manually (i.e., explicit SQL commands).

Sometimes you need to provide the DBMS with hints to help it to improve concurrency.
→ Update a tuple after reading it.

Explicit locks are also useful when doing major changes to the database.

# LOCK TABLE

Explicitly locks a table. Not part of the SQL standard.
→ Postgres/DB2/Oracle Modes: SHARE, EXCLUSIVE
→ MySQL Modes: READ, WRITE

```
LOCK TABLE <table> IN <mode> MODE;
```

```
SELECT 1 FROM <table> WITH (TABLOCK, <mode>);
```

```
LOCK TABLE <table> <mode>;
```

# SELECT...FOR UPDATE

Perform a select and then sets an exclusive lock on the matching tuples.

Can also set shared locks:
→ Postgres: FOR SHARE
→ MySQL: LOCK IN SHARE MODE

```
SELECT * FROM <table>
 WHERE <qualification> FOR UPDATE;
```

# CONCLUSION

2PL is used in almost every DBMS.

Automatically generates correct interleaving:
→ Locks + protocol (2PL, SS2PL ...)
→ Deadlock detection + handling
→ Deadlock prevention

# PROJECT #3 – QUERY EXECUTION

You will add support for executing queries in BusTub.

BusTub supports (basic) SQL with a rule-based optimizer for converting AST into physical plans.



*Prompt: A realistic photo of a bath tub with wheels and cartoon eyes driving down a city street.*

https://15445.courses.cs.cmu.edu/spring2023/project3/

# PROJECT #3 – TASKS

## Plan Node Executors
→ Access Methods: Sequential Scan, Index Scan
→ Modifications: Insert, Update, Delete
→ Joins: Nested Loop Join, Hash Join
→ Miscellaneous: Aggregation, Limit, Sort, Top-N

## Optimizer Rules:
→ Convert Nested Loop Join into a Hash Join
→ Convert ORDER BY + LIMIT into a Top-N

# PROJECT #3 - LEADERBOARD

The leaderboard requires you to add additional rules to the optimizer to generate query plans.
→ It will be impossible to get a top ranking by just having the fastest implementations in Project #1 + Project #2.

# DEVELOPMENT HINTS

Implement the **Insert** and **Sequential Scan** executors first so that you can populate tables and read from it.

You do **not** need to worry about transactions.

The aggregation and hash join hash tables do not need to be backed by the buffer pool (i.e., use STL)

Gradescope is meant for grading, not debugging. Write your own local tests.

# THINGS TO NOTE

Do **<u>not</u>** change any file other than the ones that you submit to Gradescope.

Make sure you pull in the latest changes from the BusTub main branch.

Post your questions on Piazza or come to TA office hours.

Compare against our <u>solution in your browser</u>!

# NEXT CLASS

Timestamp Ordering Concurrency Control