# Database Systems
# Lecture18 – Chapter 17: Transactions

SUNG KYUN KWAN
UNIVERSITY

Beomseok Nam (남범석)

bnam@skku.edu

- Transaction Concept
- Transaction State
- Concurrent Executions
- Serializability
- Recoverability
- Implementation of Isolation
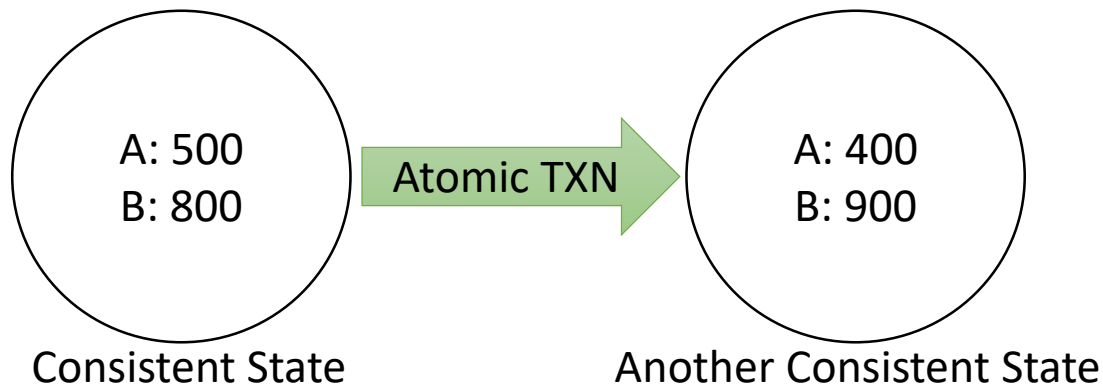- Transaction Definition in SQL
- Testing for Serializability.

- A **transaction** is a *unit* of program execution that accesses and possibly updates various data items.

- E.g., transaction to transfer $50 from account A to account B:
  1. **read**(*A*)
  2. *A := A* – 50
  3. **write**(*A*)
  4. **read**(*B*)
  5. *B := B* + 50
  6. **write**(*B)*

- Two main issues to deal with:
  - Failures of various kinds, such as hardware failures and system crashes
  - Concurrent execution of multiple transactions

- Transaction to transfer $50 from account A to account B:
    1. **read**(*A*)
    2. *A* := *A* − 50
    3. **write**(*A*)
    4. **read**(*B*)
    5. *B* := *B* + 50
    6. **write**(*B*)

- **Atomicity requirement**
    - If the transaction fails after step 3 and before step 6, money will be "lost" leading to an inconsistent database state
    - The system should ensure that updates of a partially executed transaction are not reflected in the database

- **Durability requirement** — once the user has been notified that the transaction has completed (i.e., the transfer of the $50 has taken place), the updates to the database by the transaction must persist even if there are software or hardware failures.

- **Consistency requirement** in above example:
  - The sum of A and B is unchanged by the execution of the transaction



A: 500
B: 800

Atomic TXN

A: 400
B: 900

Consistent State        Another Consistent State

- In general, consistency requirements include
  - Explicit integrity constraints such as primary keys and foreign keys
  - Implicit integrity constraints
    - e.g., sum of balances of all accounts, minus sum of loan amounts must equal value of cash-in-hand
  - A transaction must see a consistent database.

- **Isolation requirement:** if between steps 3 and 6, another transaction T2 is allowed to access the partially updated database, it will see an inconsistent database (the sum $A + B$ will be less than it should be).

    **T1**                                    **T2**

    1. **read**($A$)
    2. $A := A - 50$
    3. **write**($A$)

                        read(A), read(B), print(A+B)

    4. **read**($B$)
    5. $B := B + 50$
    6. **write**($B$)

- Isolation can be ensured trivially by running transactions **serially**, i.e., one after the other.

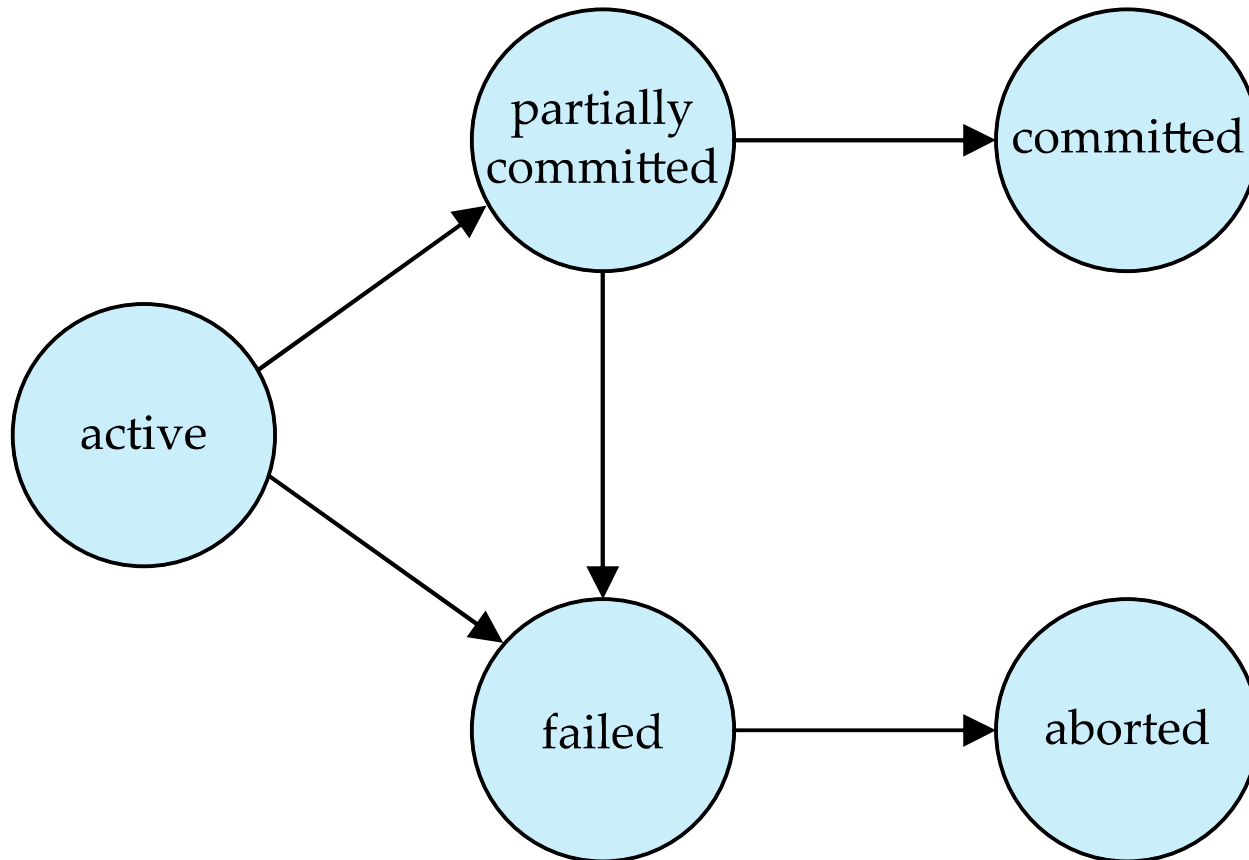- However, executing multiple transactions sequentially has performance problems.

- To preserve the integrity of data the database system must ensure:

- **Atomicity.** Either all operations of the transaction are properly reflected in the database or none are. All or Nothing.

- **Consistency.** Atomic execution of a transaction in isolation preserves the consistency of the database.

- **Isolation.** Although multiple transactions may execute concurrently, each transaction must be unaware of other concurrently executing transactions. Intermediate transaction results must be hidden from other concurrently executed transactions.
  - That is, for every pair of transactions $T_i$ and $T_j$, it appears to $T_i$ that either $T_j$ finished execution before $T_i$ started, or $T_j$ started execution after $T_i$ finished.

- **Durability.** After a transaction completes successfully, the changes it has made to the database persist, even if there are system failures.

- **Active** – the initial state; the transaction stays in this state while it is executing

- **Partially committed** – state after the final statement has been executed.

- **Failed** – state after the discovery that normal execution can no longer proceed.

- **Aborted** – state after the transaction has been rolled back and the database restored to its state prior to the start of the transaction.
    - Two options after it has been aborted:
        - Restart the transaction
            - Can be done only if no internal logical error
        - Kill the transaction

- **Committed** – state after successful completion.

# Concurrent Executions

- Multiple transactions are allowed to run concurrently in the system. Advantages are:
  - **Increased processor and disk utilization**, leading to better transaction *throughput*
    - E.g., one transaction can be using the CPU while another is reading from or writing to the disk
  - **Reduced average response time** for transactions: short transactions need not wait behind long ones.

- **Concurrency control schemes** – mechanisms  to achieve isolation
  - control the interaction among the concurrent transactions
  - prevent from destroying the consistency of the database
    - Will study in Chapter 18

# Schedules

- **Schedule** –specify the order in which instructions of concurrent transactions are executed
  - A schedule for a set of transactions must consist of all instructions of those transactions
  - Must preserve the order of instructions that appear in each individual transaction.

- A transaction that successfully completes its execution will have a commit instruction as the last statement

- A transaction that fails to complete its execution will have an abort instruction as the last statement

- Let $T_1$ transfer $50 from *A* to *B*, and $T_2$ transfer 10% of the balance from *A* to *B.*

- A serial schedule in which $T_1$ is followed by $T_2$ :

| $T_1$ | $T_2$ |
|---|---|
| read (*A*) | |
| *A* := *A* − 50 | |
| write (*A*) | |
| read (*B*) | |
| *B* := *B* + 50 | |
| write (*B*) | |
| commit | |
| | read (*A*) |
| | *temp* := *A* * 0.1 |
| | *A* := *A* - *temp* |
| | write (*A*) |
| | read (*B*) |
| | *B* := *B* + *temp* |
| | write (*B*) |
| | commit |

- A serial schedule where $T_2$ is followed by $T_1$

| $T_1$ | $T_2$ |
|---|---|
| | read ($A$) |
| | $temp := A * 0.1$ |
| | $A := A - temp$ |
| | write ($A$) |
| | read ($B$) |
| | $B := B + temp$ |
| | write ($B$) |
| | commit |
| read ($A$) | |
| $A := A - 50$ | |
| write ($A$) | |
| read ($B$) | |
| $B := B + 50$ | |
| write ($B$) | |
| commit | |

- Let $T_1$ and $T_2$ be the transactions defined previously. The following schedule is not a serial schedule, but it is *equivalent* to Schedule 1

| $T_1$ | $T_2$ |
|---|---|
| read $(A)$ <br> $A := A - 50$ <br> write $(A)$ | |
| | read $(A)$ <br> $temp := A * 0.1$ <br> $A := A - temp$ <br> write $(A)$ |
| read $(B)$ <br> $B := B + 50$ <br> write $(B)$ <br> commit | |
| | read $(B)$ <br> $B := B + temp$ <br> write $(B)$ <br> commit |

- In Schedules 1, 2 and 3, the sum A + B is preserved.

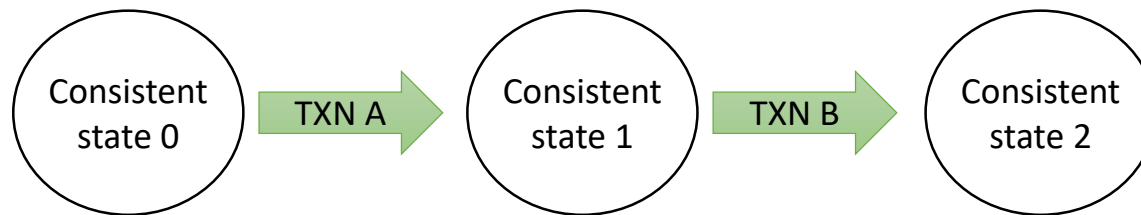▪ The following schedule does not preserve the value of ($A + B$).

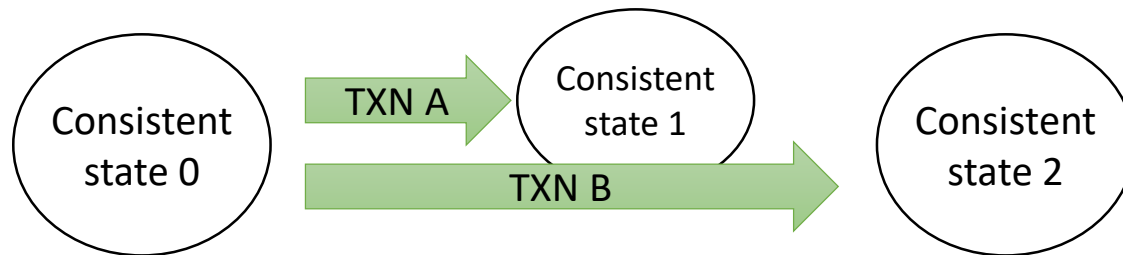| $T_1$ | $T_2$ |
|---|---|
| read ($A$)<br>$A := A - 50$ | |
| | read ($A$)<br>$temp := A * 0.1$<br>$A := A - temp$<br>write ($A$)<br>read ($B$) |
| write ($A$)<br>read ($B$)<br>$B := B + 50$<br>write ($B$)<br>commit | |
| | $B := B + temp$<br>write ($B$)<br>commit |

- **Basic Assumption**:
  - We assume each transaction preserves database consistency.
- Thus, serial execution of a set of transactions preserves database consistency.



- A schedule is *serializable* if it is equivalent to a serial schedule.



- Different forms of serializability:
  1. **Conflict serializability**
  2. **View serializability**

# *Simplified view of transactions*

- We ignore operations other than **read** and **write** instructions

- We assume that transactions may perform arbitrary computations on data in local buffers in between reads and writes.

- Our simplified schedules consist of only **read** and **write** instructions.

- Instructions $I_i$ and $I_j$ of transactions $T_i$ and $T_j$ **conflict** if and only if at least one of these instructions write a shared data item $Q.$

  1. $I_i = \textbf{read}(Q), I_j = \textbf{read}(Q).$  $I_i$ and $I_j$ don't conflict.
  2. $I_i = \textbf{read}(Q),$  $I_j = \textbf{write}(Q).$  They conflict.
  3. $I_i = \textbf{write}(Q), I_j = \textbf{read}(Q).$  They conflict
  4. $I_i = \textbf{write}(Q), I_j = \textbf{write}(Q).$  They conflict

- Note: a conflict forces an order between them.

- If *instructions* are consecutive in a schedule and they do not conflict, their results would remain the same even if we interchange them in the schedule.

## Conflict Serializability

# Conflict Serializability

- If a schedule *S* can be transformed into a schedule *S'* by a series of swaps of non-conflicting instructions, we say that *S* and *S'* are **conflict equivalent**.

- We say that a schedule *S* is **conflict serializable** if it is conflict equivalent to a serial schedule

- Schedule 3 can be transformed into Schedule 6, a serial schedule where $T_2$ follows $T_1$, by series of swaps of non-conflicting instructions. Therefore Schedule 3 is <span style="color:red">conflict serializable</span>.

| $T_1$ | $T_2$ |
|---|---|
| read ($A$) | |
| write ($A$) | |
| | read ($A$) |
| | write ($A$) |
| read ($B$) | |
| write ($B$) | |
| | read ($B$) |
| | write ($B$) |

Schedule 3

| $T_1$ | $T_2$ |
|---|---|
| read ($A$) | |
| write ($A$) | |
| read ($B$) | |
| write ($B$) | |
| | read ($A$) |
| | write ($A$) |
| | read ($B$) |
| | write ($B$) |

Schedule 6

- Example of a schedule that is not conflict serializable:

| $T_3$ | $T_4$ |
|---|---|
| read ($Q$) | |
| | write ($Q$) |
| write ($Q$) | |

- We are unable to swap instructions in the above schedule to obtain either the serial schedule < $T_3$, $T_4$ >, or the serial schedule < $T_4$, $T_3$ >.

- Let $S$ and $S'$ be two schedules with the same set of transactions.
- $S$ and $S'$ are **view equivalent** if the following three conditions are met

    1. If in schedule S, transaction $T_i$ reads the initial value of $Q$, then in schedule $S'$ also transaction $T_i$ must read the initial value of $Q$.

    2. If in schedule S transaction $T_i$ executes **read**($Q$), and that value was produced by transaction $T_j$ (if any), then in schedule $S'$ also transaction $T_i$ must read the value of $Q$ that was produced by the same **write**(Q) operation of transaction $T_j$.

    3. The transaction (if any) that performs the final **write**($Q$) operation in schedule $S$ must also perform the final **write**($Q$) operation in schedule $S'$.

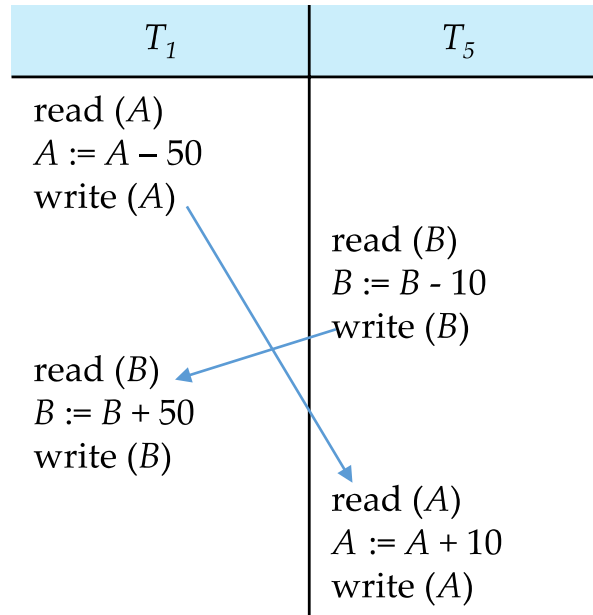- View equivalence is also based purely on **reads** and **writes** alone.

- A schedule *S* is **view serializable** if it is view equivalent to a serial schedule.

- Every conflict serializable schedule is also view serializable.

- Below is a schedule which is view-serializable but *not* conflict serializable.

| $T_{27}$ | $T_{28}$ | $T_{29}$ |
|----------|----------|----------|
| read ($Q$) | | |
| | write ($Q$) | |
| write ($Q$) | | |
| | | write ($Q$) |

- What serial schedule is above equivalent to?

- Every view serializable schedule that is not conflict serializable has **blind writes**.

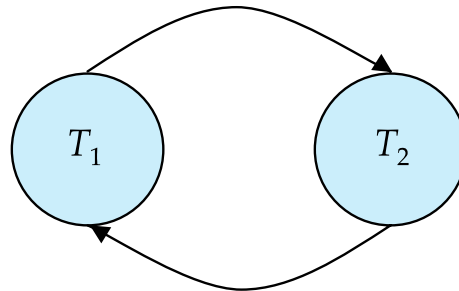- The schedule below produces same outcome as the serial schedule < $T_1$, $T_5$ >, yet is not conflict equivalent or view equivalent to it.

| $T_1$ | $T_5$ |
|---|---|
| read ($A$) | |
| $A := A - 50$ | |
| write ($A$) | |
| | read ($B$) |
| | $B := B - 10$ |
| | write ($B$) |
| read ($B$) | |
| $B := B + 50$ | |
| write ($B$) | |
| | read ($A$) |
| | $A := A + 10$ |
| | write ($A$) |

- Determining such equivalence requires analysis of operations other than read and write.
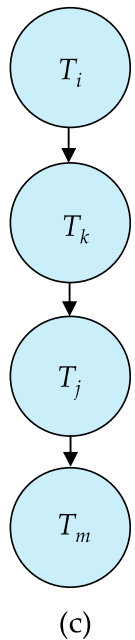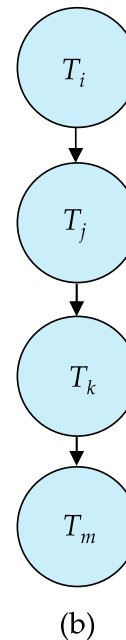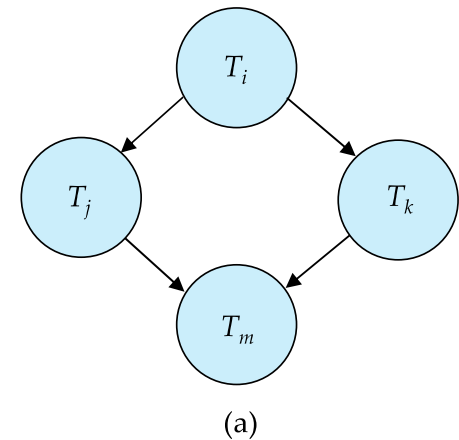
- Consider some schedule of a set of transactions $T_1$, $T_2$, ..., $T_n$

- **Precedence graph** — a direct graph where the vertices are the transactions (names).

- We draw an arc from $T_i$ to $T_j$ if the two transactions conflict, and $T_i$ accessed the conflicting data item earlier than $T_j$.

- We may label the arc by the item that was accessed.

- Example of a precedence graph

- A schedule is conflict serializable if and only if its precedence graph is acyclic.

- Cycle-detection algorithms exist which take order $n^2$ time, where $n$ is the number of vertices in the graph.
  - Better algorithms take order $n + e$ where $e$ is the number of edges.

- If precedence graph is acyclic, the serializability order can be obtained by a *topological sorting* of the graph.
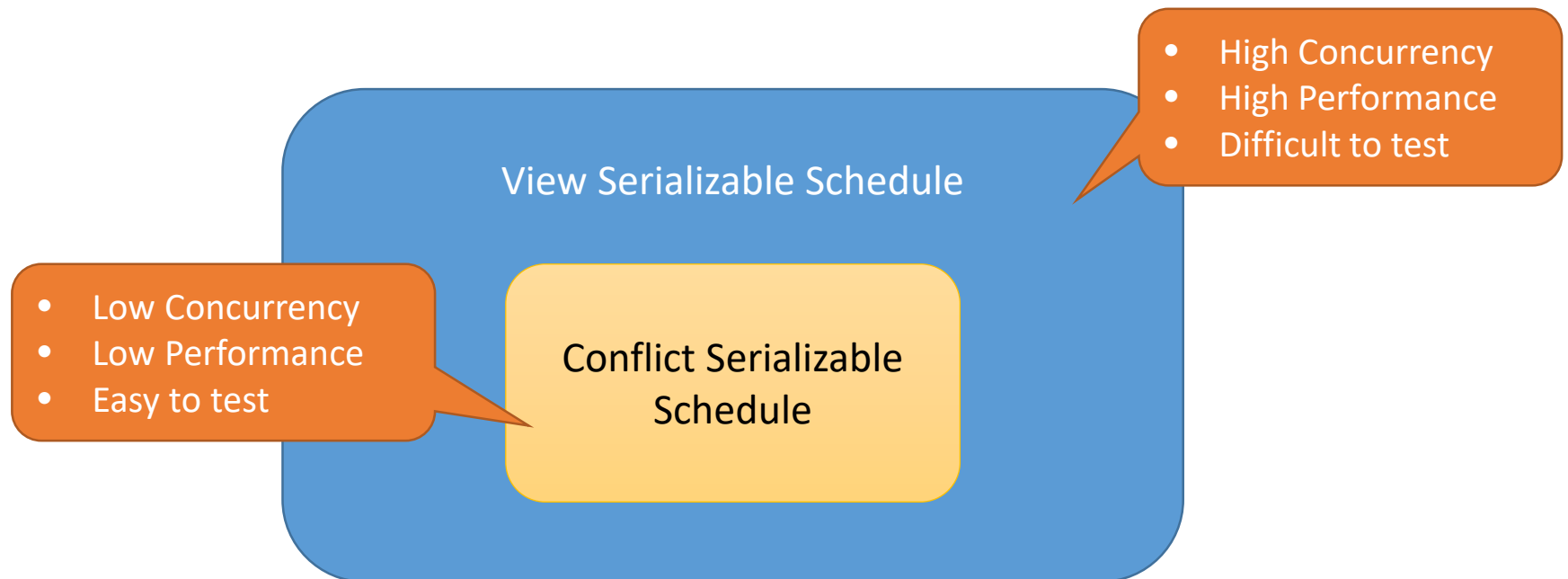  - This is a linear order consistent with the partial order of the graph.

(a)

(b)

(c)

- The precedence graph test for conflict serializability cannot be used directly to test for view serializability.
  - Extension to test for view serializability has cost exponential in the size of the precedence graph.

- The problem of checking if a schedule is view serializable falls in the class of *NP*-complete problems.
  - Thus, existence of an efficient algorithm is *extremely* unlikely.

- However practical algorithms that just check some **sufficient conditions** for view serializability can still be used.
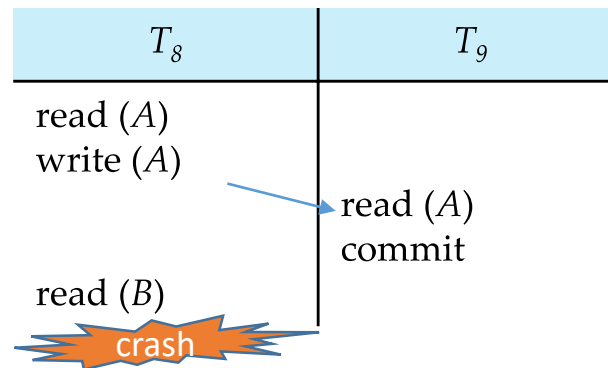
- Conflict serializability is more popular in DBMS in practice
- But, conflict serializability actually permit only a subset of serializable schedules that do not have consistency problems.
- The general form of view serializability is very expensive to test, and only a very restricted form of it is used for concurrency control.

View Serializable Schedule

- High Concurrency
- High Performance
- Difficult to test

Conflict Serializable Schedule

- Low Concurrency
- Low Performance
- Easy to test

- Need to address the effect of transaction failures on concurrently running transactions.

- **Recoverable schedule** — if a transaction $T_j$ reads a data item previously written by a transaction $T_i$, then the commit operation of $T_i$ appears before the commit operation of $T_j$.

- The following schedule is not recoverable

| $T_8$ | $T_9$ |
|---|---|
| read ($A$) | |
| write ($A$) | |
| | read ($A$) |
| | commit |
| read ($B$) | |
| crash | |

- If $T_8$ should abort, $T_9$ would have read (and possibly shown to the user) an inconsistent database state. Hence, database must ensure that schedules are recoverable.

- **Cascading rollback** – a single transaction failure leads to a series of transaction rollbacks.  Consider the following schedule where none of the transactions has yet committed (so the schedule is recoverable)

| $T_{10}$ | $T_{11}$ | $T_{12}$ |
|---|---|---|
| read $(A)$<br>read $(B)$<br>write $(A)$ | | |
| | read $(A)$<br>write $(A)$ | |
| | | read $(A)$ |
| abort | | |

If $T_{10}$ fails, $T_{11}$ and $T_{12}$ must also be rolled back.

- Can lead to the undoing of a significant amount of work

- **Cascadeless schedules** — cascading rollbacks cannot occur;
  - For each pair of transactions $T_i$ and $T_j$ such that $T_j$ reads a data item previously written by $T_i$, the commit operation of $T_i$ appears before the read operation of $T_j$.
  - Read committed data only!
- Every Cascadeless schedule is also recoverable
- It is desirable to restrict the schedules to those that are cascadeless

- **Goal:** A DBMS must make scheduling decisions that are
  - either conflict or view serializable, and
  - recoverable and preferably cascadeless

- DBMS scheduler should provide a high degree of concurrency
  - A policy in which only one transaction can execute at a time generates serial schedules, but provides a poor degree of concurrency.
  - Concurrency-control schemes tradeoff between the amount of concurrency they allow and the amount of overhead that they incur.
  - Some schemes allow only conflict-serializable schedules to be generated, while others allow view-serializable schedules that are not conflict-serializable.
  - We study concurrency control protocols in Chapter 18.

- Some applications (OLAP) are willing to live with weak levels of consistency, allowing schedules that are not serializable
  - Tradeoff accuracy for performance

- On Line Transaction Processing –OLTP
  - Maintains a database that is an accurate model of some real-world enterprise
    - Short simple transactions
    - Transactions access only a small fraction of the database

- On Line Analytic Processing –OLAP
  - Uses database to guide strategic decisions.
    - Complex queries
    - Transactions access a large fraction of the database
    - Data need not be up-to-date
    - E.g., A bank CEO wants to get an approximate total balance of all customers' accounts
  - OLAP transactions need not be serializable with respect to other OLTP transactions

# Phenomena caused by Concurrent Transactions

- **dirty read**
  - A transaction reads data written by a concurrent uncommitted transaction.

- **nonrepeatable read**
  - A transaction runs the same query twice and finds that a data has been modified by another transaction.

- **phantom read**
  - A transaction runs the same query twice and finds that the result set has new records due to another recently-committed transaction.

- **serialization anomaly**
  - The result of successfully committing a group of transactions is inconsistent with all possible orderings of running those transactions one at a time.

- E.g., Transaction 1:
  **select** *ID, name*  **from**  *instructor*  **where** *salary* > 90000

- E.g., Transaction 2:
  **insert into** *instructor* **values** ('11111', 'James', 'Marketing', 100000)

- Suppose
  - T1 starts, finds tuples salary > 90000 using index and locks them
  - And then T2 executes.
  - Do T1 and T2 conflict?  Does tuple level locking detect the conflict?
  - Instance of the **phantom phenomenon**

- **Serializable** — default

- **Repeatable read** — only committed records to be read.
  - Repeated reads of same record must return same value.
  - However, a transaction may not be serializable – it may find some records inserted by a transaction but not find others.

- **Read committed** — only committed records can be read.
  - Successive reads of record may return different (but committed) values.

- **Read uncommitted** — even uncommitted records may be read.

| Isolation Level | Dirty Read | Nonrepeatable Read | Phantom Read | Serialization Anomaly |
|---|---|---|---|---|
| Read uncommitted | Allowed, but not in PG | Possible | Possible | Possible |
| Read committed | Not possible | Possible | Possible | Possible |
| Repeatable read | Not possible | Not possible | Allowed, but not in PG | Possible |
| Serializable | Not possible | Not possible | Not possible | Not possible |

- In SQL, a transaction begins implicitly.

- A transaction in SQL ends by:
  - **Commit work** commits current transaction and begins a new one.
  - **Rollback work** causes current transaction to abort.

- In almost all database systems, by default, every SQL statement also commits implicitly if it executes successfully
  - Implicit commit can be turned off by a database directive
    - E.g., in JDBC -- connection.setAutoCommit(false);

- Isolation level can be set at database level

- Isolation level can be changed at start of transaction
  - E.g. In SQL **set transaction isolation level serializable**
  - E.g. in JDBC -- connection.setTransactionIsolation( Connection.TRANSACTION_SERIALIZABLE)

- Locking
  - Lock on whole database vs lock on items
  - How long to hold lock?
  - Shared vs exclusive locks

- Timestamps
  - Transaction timestamp assigned e.g. when a transaction begins
  - Data items store two timestamps
    - Read timestamp
    - Write timestamp
  - Timestamps are used to detect out of order accesses

- Multiple versions of each data item
  - Allow transactions to read from a "snapshot" of the database