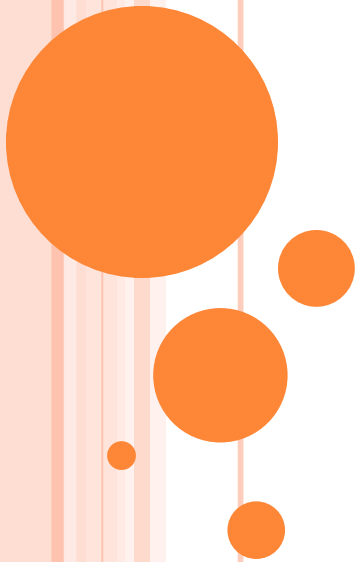# Transaction Management

# TRANSACTION CONCEPT

- A **transaction** is a *unit* of program execution that accesses and possibly updates various data items.
- A transaction must see a consistent database.
- During transaction execution the database may be temporarily inconsistent.
- When the transaction completes successfully (is committed), the database must be consistent.
- After a transaction commits, the changes it has made to the database persist, even if there are system failures.
- Multiple transactions can execute in parallel.
- Two main issues to deal with:
  - Failures of various kinds, such as hardware failures and system crashes
  - Concurrent execution of multiple transactions

# ACID PROPERTIES

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.
- **Consistency.** Execution of a transaction run by itself 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.

# EXAMPLE OF FUND TRANSFER

- 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, the system should ensure that its updates are not reflected in the database, else an inconsistency will result.

- **Consistency requirement** – the sum of A and B is unchanged by the execution of the transaction.

# EXAMPLE OF FUND TRANSFER (CONT.)

- **Isolation requirement** — if between steps 3 and 6, another transaction 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).

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

  - However, executing multiple transactions concurrently has significant benefits, as we will see later.

- **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 despite failures.
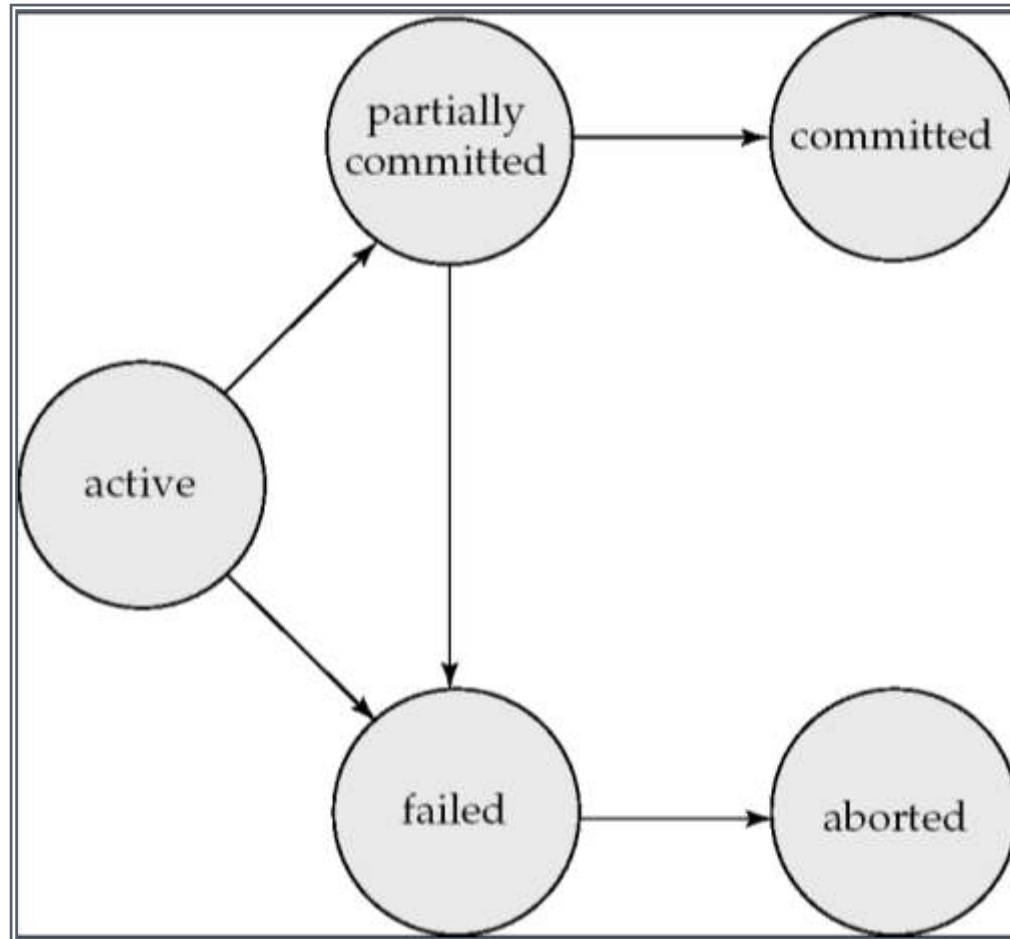
# TRANSACTION STATE

- **Active** – the initial state; the transaction stays in this state while it is executing
- **Partially committed** – after the final statement has been executed.
- **Failed** -- after the discovery that normal execution can no longer proceed.
- **Aborted** – 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** – after successful completion.
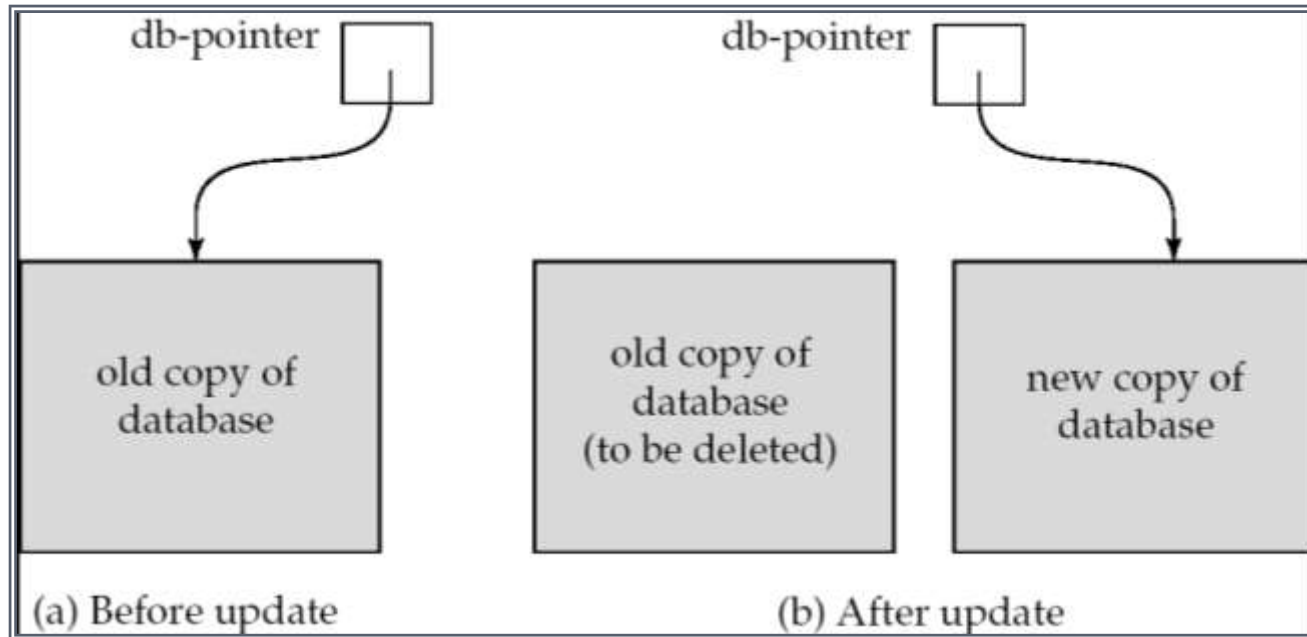
# Transaction State (Cont.)

# Implementation of Atomicity and Durability

- The **recovery-management** component of a database system implements the support for atomicity and durability.
- The *shadow-database* scheme:
  - assume that only one transaction is active at a time.
  - a pointer called db_pointer always points to the current consistent copy of the database.
  - all updates are made on a *shadow copy* of the database, and **db_pointer** is made to point to the updated shadow copy only after the transaction reaches partial commit and all updated pages have been flushed to disk.
  - in case transaction fails, old consistent copy pointed to by **db_pointer** can be used, and the shadow copy can be deleted.

The shadow-database scheme:



(a) Before update    (b) After update

- Assumes disks do not fail
- Useful for text editors, but
  - extremely inefficient for large databases
  - Does not handle concurrent transactions

# CONCURRENT EXECUTIONS

- Multiple transactions are allowed to run concurrently in the system. Advantages are:
  - **increased processor and disk utilization**, leading to better transaction *throughput:* 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; that is, to control the interaction among the concurrent transactions in order to prevent them from destroying the consistency of the database

# SCHEDULES

- **Schedule** – a sequences of instructions that specify the chronological 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 in which the instructions appear in each individual transaction.
- A transaction that successfully completes its execution will have a commit instructions as the last statement (will be omitted if it is obvious)
- A transaction that fails to successfully complete its execution will have an abort instructions as the last statement (will be omitted if it is obvious)

# SCHEDULE 1

- 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 ($B$) | |
| | read($A$) |
| | $temp := A * 0.1$ |
| | $A := A - temp$ |
| | write($A$) |
| | read($B$) |
| | $B := B + temp$ |
| | write($B$) |

# SCHEDULE 2

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

# SCHEDULE 3

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

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

# SCHEDULE 4

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

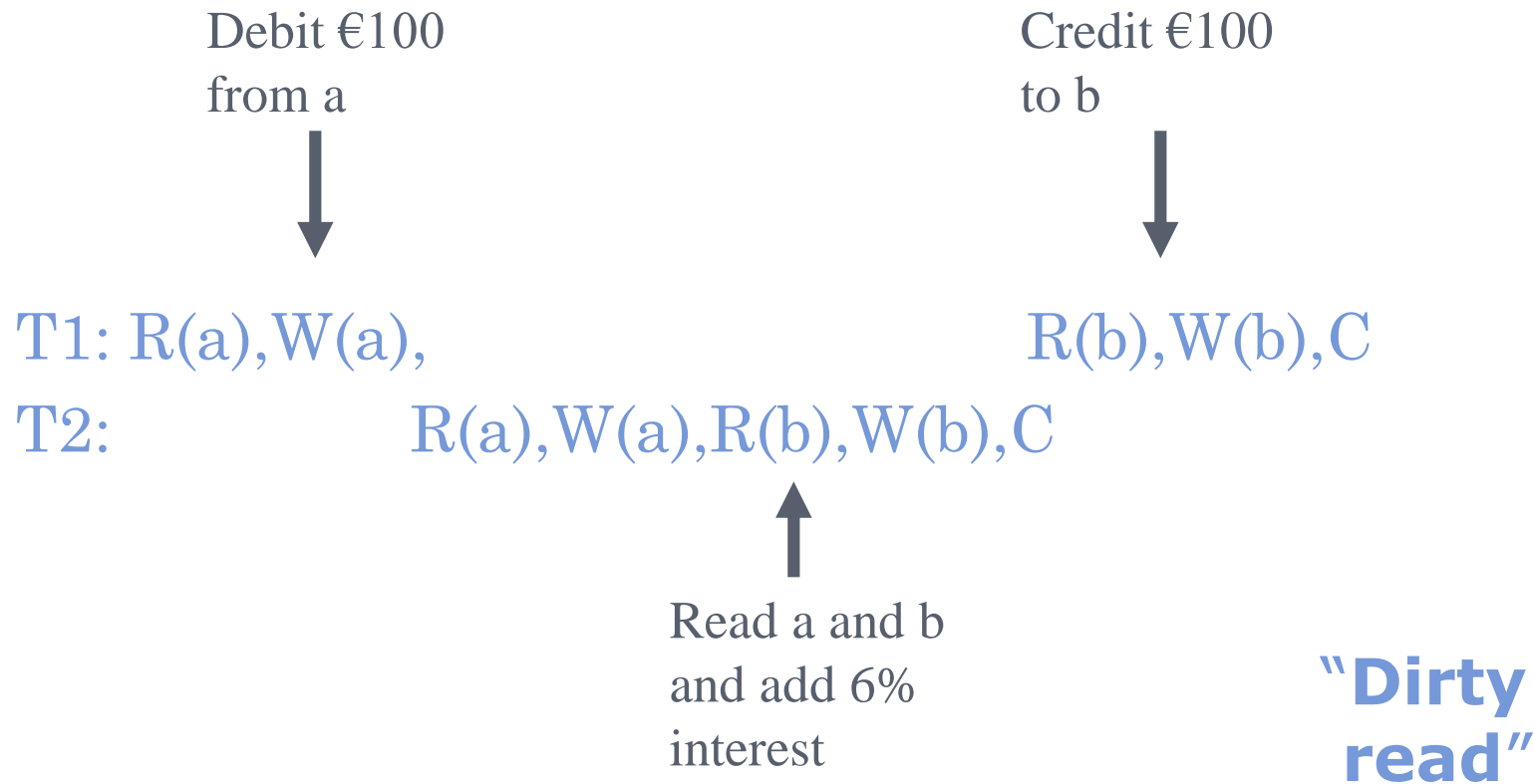| $T_1$ | $T_2$ |
|---|---|
| read($A$) | |
| $A := A - 50$ | |
| | read($A$) |
| | $temp := A * 0.1$ |
| | $A := A - temp$ |
| | write($A$) |
| | read($B$) |
| write($A$) | |
| read($B$) | |
| $B := B + 50$ | |
| write($B$) | |
| | $B := B + temp$ |
| | write($B$) |

# ANOMALIES WITH INTERLEAVED EXECUTION

- Two actions on the same data object **conflict** if at least one of them is a write.
- We'll now consider three ways in which a schedule involving two consistency-preserving transactions can leave a consistent database inconsistent

# WR CONFLICTS

- Transaction T2 reads a database object that has been modified by T1 which has not committed

Debit €100
from a

Credit €100
to b

T1: R(a),W(a),                               R(b),W(b),C

T2:                     R(a),W(a),R(b),W(b),C

Read a and b
and add 6%
interest

**"Dirty read"**

# RW CONFLICTS

- Transaction T2 could change the value of an object that has been read by a transaction T1, while T1 is still in progress

T1: R(a),                    R(a), W(a), C

T2:          R(a),W(a),C

**"Unrepeatable Read"**

Read A (5)    Write 5+1=6

T1: R(a),        W(a),C

T2:          R(a),            W(a),C

A is 4 ☹

Read A (5)    Write 5-1=4

# WW CONFLICTS

- Transaction T2 could overwrite the value of an object which has already been modified by T1, while T1 is still in progress

   T1: [W(Britney), W(gmb)]    "Set both salaries at $1000"

   T2: [W(gmb), W(Britney)]    "Set both salaries at $2000"

- But:                    **"Blind Write or lost update"**
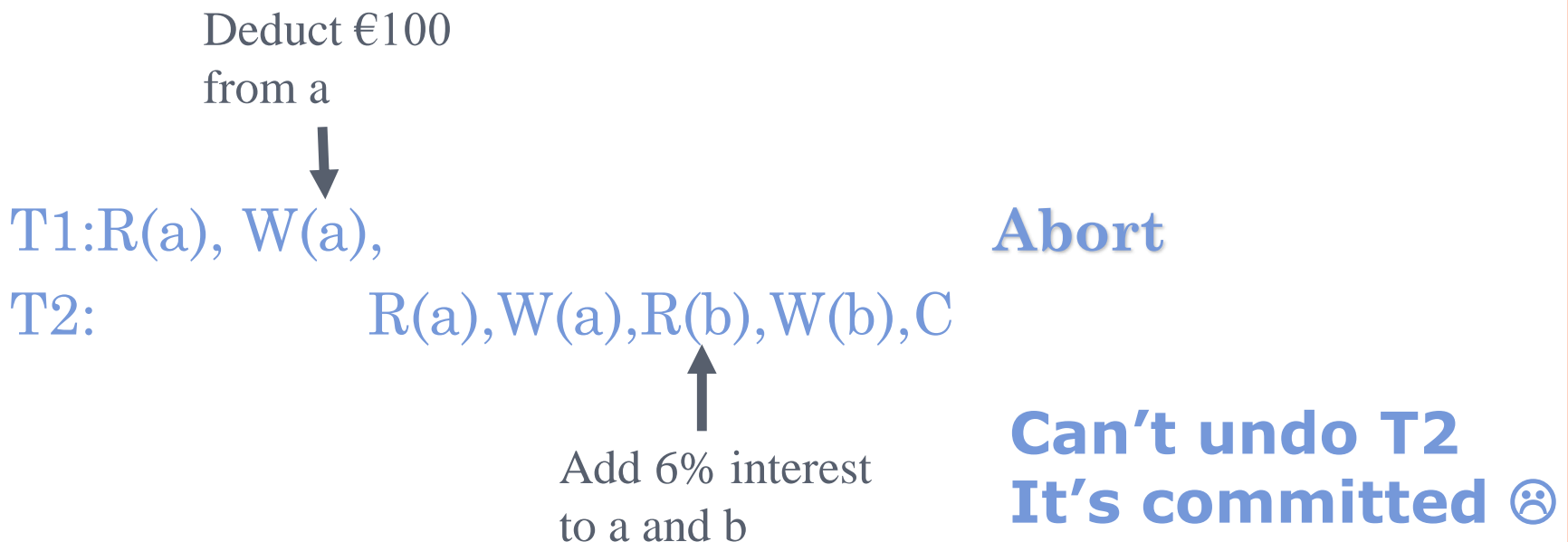
                                                    **gmb gets $1000**
T1: W(Britney),                          W(gmb)     **Britney gets $2000**
T2:                W(gmb), W(Britney)               ☺

# SERIALISABILITY AND ABORTS

- Things are more complicated when transactions can abort

Deduct €100
from a

T1:R(a), W(a),                    **Abort**
T2:            R(a),W(a),R(b),W(b),C

Add 6% interest
to a and b

**Can't undo T2**
**It's committed** ☹

# SERIALIZABILITY

- **Basic Assumption** – Each transaction preserves database consistency.

- Thus serial execution of a set of transactions preserves database consistency.

- A (possibly concurrent) schedule is serializable if it is equivalent to a serial schedule. Different forms of schedule equivalence give rise to the notions of:
  1. **conflict serializability**
  2. **view serializability**

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

# CONFLICTING INSTRUCTIONS

- Instructions $l_i$ and $l_j$ of transactions $T_i$ and $T_j$ respectively, **conflict** if and only if there exists some item $Q$ accessed by both $l_i$ and $l_j$, and at least one of these instructions wrote $Q$.

    1. $l_i = \mathbf{read}(Q)$, $l_j = \mathbf{read}(Q)$.   $l_i$ and $l_j$ don't conflict.
    2. $l_i = \mathbf{read}(Q)$,  $l_j = \mathbf{write}(Q)$.  They conflict.
    3. $l_i = \mathbf{write}(Q)$, $l_j = \mathbf{read}(Q)$.   They conflict
    4. $l_i = \mathbf{write}(Q)$, $l_j = \mathbf{write}(Q)$.  They conflict

- Intuitively, a conflict between $l_i$ and $l_j$ forces a (logical) temporal order between them.

  - If $l_i$ and $l_j$ are consecutive in a schedule and they do not conflict, their results would remain the same even if they had been interchanged in the schedule.

# 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

# CONFLICT SERIALIZABILITY (CONT.)

- 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 conflict serializable.

| $T_1$ | $T_2$ |
|---|---|
| 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

# CONFLICT SERIALIZABILITY (CONT.)

- 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 $< T_4, T_3 >$.

# VIEW SERIALIZABILITY

- 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. For each data item $Q$, if transaction $T_i$ reads the initial value of $Q$ in schedule $S$, then transaction $T_i$ must, in schedule $S'$, also read the initial value of $Q$.

  2. For each data item $Q$ if transaction $T_i$ executes **read**$(Q)$ in schedule $S$, and that value was produced by transaction $T_j$ (if any), then transaction $T_i$ must in schedule $S'$ also read the value of $Q$ that was produced by transaction $T_j$.

  3. For each data item $Q$, the transaction (if any) that performs the final **write**$(Q)$ operation in schedule $S$ must perform the final **write**$(Q)$ operation in schedule $S'$.

  As can be seen, view equivalence is also based purely on **reads** and **writes** alone.

# VIEW SERIALIZABILITY (CONT.)

- A schedule $S$ is **view serializable** 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.
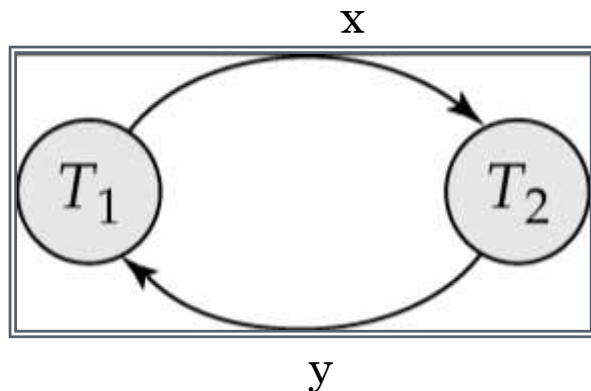
| $T_3$ | $T_4$ | $T_6$ |
|---|---|---|
| read($Q$) | | |
| | write($Q$) | |
| write($Q$) | | |
| | | write($Q$) |

- What serial schedule is above equivalent to?
- Every view serializable schedule that is not conflict serializable has **blind writes.**

# TESTING FOR SERIALIZABILITY

- 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 transaction conflict, and $T_i$ accessed the data item on which the conflict arose earlier.
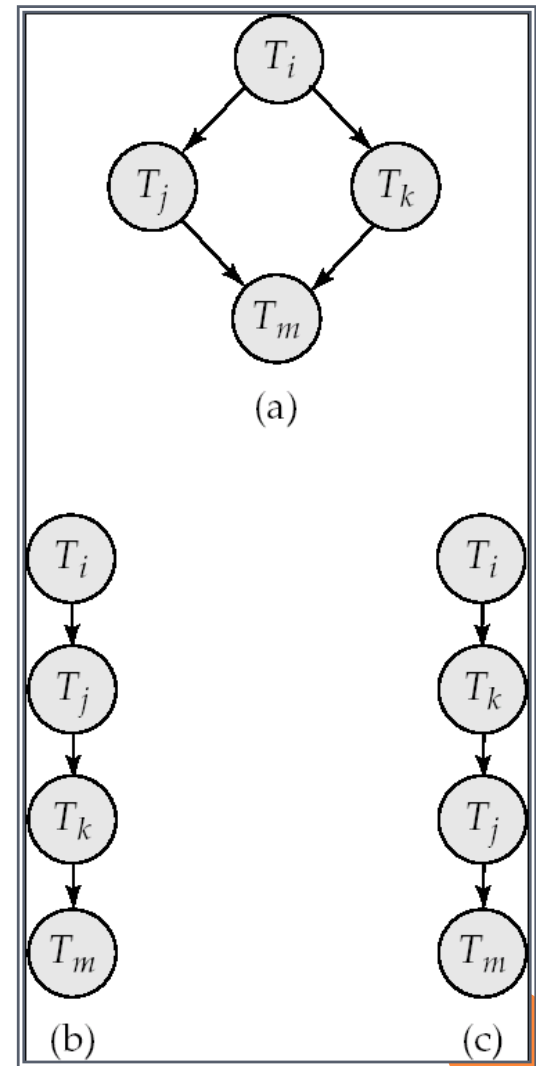- We may label the arc by the item that was accessed.
- **Example 1**

x



y

The set of edges consists of all edges
*Ti →Tj for which one of three conditions holds:*
**1. Ti executes write(Q) before Tj executes read(Q).**
**2. Ti executes read(Q) before Tj executes write(Q).**
**3. Ti executes write(Q) before Tj executes write(Q).**

# TEST FOR CONFLICT SERIALIZABILITY

- 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.
- 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.
  - For example, a serializability order for Schedule A would be $Ti \rightarrow Tj \rightarrow Tk \rightarrow T_m$
    - Are there others?

# TEST FOR VIEW SERIALIZABILITY

- 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.

# RECOVERABLE SCHEDULES

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 (Schedule 11) is not recoverable if $T_9$ commits immediately after the read

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

- 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 ROLLBACKS

- **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($A$) | | |
| read($B$) | | |
| write($A$) | | |
| | read($A$) | |
| | write($A$) | |
| | | read($A$) |

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

- **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 $T_j$.

- Every cascadeless schedule is also recoverable

- It is desirable to restrict the schedules to those that are cascadeless

# Concurrency Control

- A database must provide a mechanism that will ensure that all possible schedules are
  - either conflict or view serializable, and
  - are recoverable and preferably cascadeless
- A policy in which only one transaction can execute at a time generates serial schedules, but provides a poor degree of concurrency
  - Are serial schedules recoverable/cascadeless?
- Testing a schedule for serializability *after* it has executed is a little too late!
- **Goal** – to develop concurrency control protocols that will assure serializability.

# Concurrency Control vs. Serializability Tests

- Concurrency-control protocols allow concurrent schedules, but ensure that the schedules are conflict/view serializable, and are recoverable and cascadeless.
- Concurrency control protocols generally do not examine the precedence graph as it is being created
  - Instead a protocol imposes a discipline that avoids nonseralizable schedules.
- Different concurrency control protocols provide different tradeoffs between the amount of concurrency they allow and the amount of overhead that they incur.
- Tests for serializability help us understand why a concurrency control protocol is correct.

# LOCK-BASED PROTOCOLS

- A lock is a mechanism to control concurrent access to a data item

- Data items can be locked in two modes :

  1.  *exclusive* *(X) mode*. Data item can be both read as well as written. X-lock is requested using **lock-X** instruction.

  2.  *shared* *(S) mode*. Data item can only be read. S-lock is requested using **lock-S** instruction.

- Lock requests are made to concurrency-control manager. Transaction can proceed only after request is granted.

# LOCK-BASED PROTOCOLS (CONT.)

- Lock-compatibility matrix

|   | S | X |
|---|---|---|
| S | true | false |
| X | false | false |

- A transaction may be granted a lock on an item if the requested lock is compatible with locks already held on the item by other transactions

- Any number of transactions can hold shared locks on an item,
  - but if any transaction holds an exclusive on the item no other transaction may hold any lock on the item.

- If a lock cannot be granted, the requesting transaction is made to wait till all incompatible locks held by other transactions have been released. The lock is then granted.

# LOCK-BASED PROTOCOLS (CONT.)

- Example of a transaction performing locking:

$T_1$: lock-X($B$);
    read($B$);
    $B := B - 50$;
    write($B$);
    unlock($B$);
    lock-X($A$);
    read($A$);
    $A := A + 50$;
    write($A$);
    unlock($A$).

$T_2$: LOCK-S($A$);
    READ ($A$);
    UNLOCK($A$);
    LOCK-S($B$);
    READ ($B$);
    UNLOCK($B$);
    DISPLAY($A+B$)

- Locking as above is not sufficient to guarantee serializability — if $A$ and $B$ get updated in-between the read of $A$ and $B$, the displayed sum would be wrong.
- A **locking protocol** is a set of rules followed by all transactions while requesting and releasing locks. Locking protocols restrict the set of possible schedules.

| $T_1$ | $T_2$ | concurreny-control manager |
|---|---|---|
| lock-X($B$) | | |
| | | grant-X($B$, $T_1$) |
| read($B$) | | |
| $B := B - 50$ | | |
| write($B$) | | |
| unlock($B$) | | |
| | lock-S($A$) | |
| | | grant-S($A$, $T_2$) |
| | read($A$) | |
| | unlock($A$) | |
| | lock-S($B$) | |
| | | grant-S($B$, $T_2$) |
| | read($B$) | |
| | unlock($B$) | |
| | display($A + B$) | |
| lock-X($A$) | | |
| | | grant-X($A$, $T_1$) |
| read($A$) | | |
| $A := A - 50$ | | |
| write($A$) | | |
| unlock($A$) | | |

$T_3$: lock-X($B$);
    read($B$);
    $B := B - 50$;
    write($B$);
    lock-X($A$);
    read($A$);
    $A := A + 50$;
    write($A$);
    unlock($B$);
    unlock($A$).

Transaction $T_3$ (transaction $T_1$ with unlocking delayed).

$T_4$: lock-S($A$);
    read($A$);
    lock-S($B$);
    read($B$);
    display($A + B$);
    unlock($A$);
    unlock($B$).

Transaction $T_4$ (transaction $T_2$ with unlocking delayed).

# PITFALLS OF LOCK-BASED PROTOCOLS

- Consider the partial schedule

| $T_3$ | $T_4$ |
|---|---|
| lock-x($B$) | |
| read($B$) | |
| $B := B - 50$ | |
| write($B$) | |
| | lock-s($A$) |
| | read($A$) |
| | lock-s($B$) |
| lock-x($A$) | |

- Neither $T_3$ nor $T_4$ can make progress — executing **lock-S(B)** causes $T_4$ to wait for $T_3$ to release its lock on $B$, while executing **lock-X(A)** causes $T_3$ to wait for $T_4$ to release its lock on $A$.
- Such a situation is called a **deadlock**.
  - To handle a deadlock one of $T_3$ or $T_4$ must be rolled back and its locks released.

- The potential for deadlock exists in most locking protocols. Deadlocks are a necessary evil.

- **Starvation** is also possible if concurrency control manager is badly designed. For example:
  - A transaction may be waiting for an X-lock on an item, while a sequence of other transactions request and are granted an S-lock on the same item.
  - The same transaction is repeatedly rolled back due to deadlocks.

- Concurrency control manager can be designed to prevent starvation.

# The Two-Phase Locking Protocol

- This is a protocol which ensures conflict-serializable schedules.
- Phase 1: Growing Phase
  - transaction may obtain locks
  - transaction may not release locks
- Phase 2: Shrinking Phase
  - transaction may release locks
  - transaction may not obtain locks
- The protocol assures serializability. It can be proved that the transactions can be serialized in the order of their **lock points**  (i.e. the point where a transaction acquired its final lock).

- Two-phase locking *does not* ensure freedom from deadlocks
- Cascading roll-back is possible under two-phase locking. To avoid this, follow a modified protocol called **strict two-phase locking**. Here a transaction must hold all its exclusive locks till it commits/aborts.
- **Rigorous two-phase locking** is even stricter: here *all* locks are held till commit/abort. In this protocol transactions can be serialized in the order in which they commit.

| $T_5$ | $T_6$ | $T_7$ |
|---|---|---|
| lock-x($A$) | | |
| read($A$) | | |
| lock-s($B$) | | |
| read($B$) | | |
| write($A$) | | |
| unlock($A$) | | |
| | lock-x($A$) | |
| | read($A$) | |
| | write($A$) | |
| | unlock($A$) | |
| | | lock-s($A$) |
| | | read($A$) |

| T1 | T2 |
|---|---|
| X(A) | |
| R(A) | |
| W(A) | |
| X(B) | |
| R(B) | |
| W(B) | |
| Commit | |
| | X(A) |
| | R(A) |
| | W(A) |
| | X(B) |
| | R(B) |
| | W(B) |
| | Commit |

Schedule Illustrating Strict 2PL with Serial Execution

| T1 | T2 |
|---|---|
| S(A) | |
| R(A) | |
| | S(A) |
| | R(A) |
| | X(B) |
| | R(B) |
| | W(B) |
| | Conllnit |
| X(C) | |
| R(C) | |
| W(C) | |
| Commit | |

Schedule Following Strict 2PL with Interleaved Actions

# ABORTING

- If a transaction $T_i$ is aborted, then all actions must be undone
  - Also, if $T_j$ reads object last written by $T_i$, then $T_j$ must be aborted!
- Most systems avoid **cascading aborts** by releasing locks only at commit time (strict protocols)
  - If $T_i$ writes an object, then $T_j$ can only read this after $T_i$ finishes
- In order to undo changes, the DBMS maintains a **log** which records every write

# LOCK-BASED CONCURRENCY CONTROL

- *Strict Two-phase Locking (Strict 2PL) Protocol*:
  - Each Xact must obtain a S (*shared*) lock on object before reading, and an X (*exclusive*) lock on object before writing.
  - All locks held by a transaction are released when the transaction completes
    - (Non-strict) 2PL Variant: Release locks anytime, but cannot acquire locks after releasing any lock.
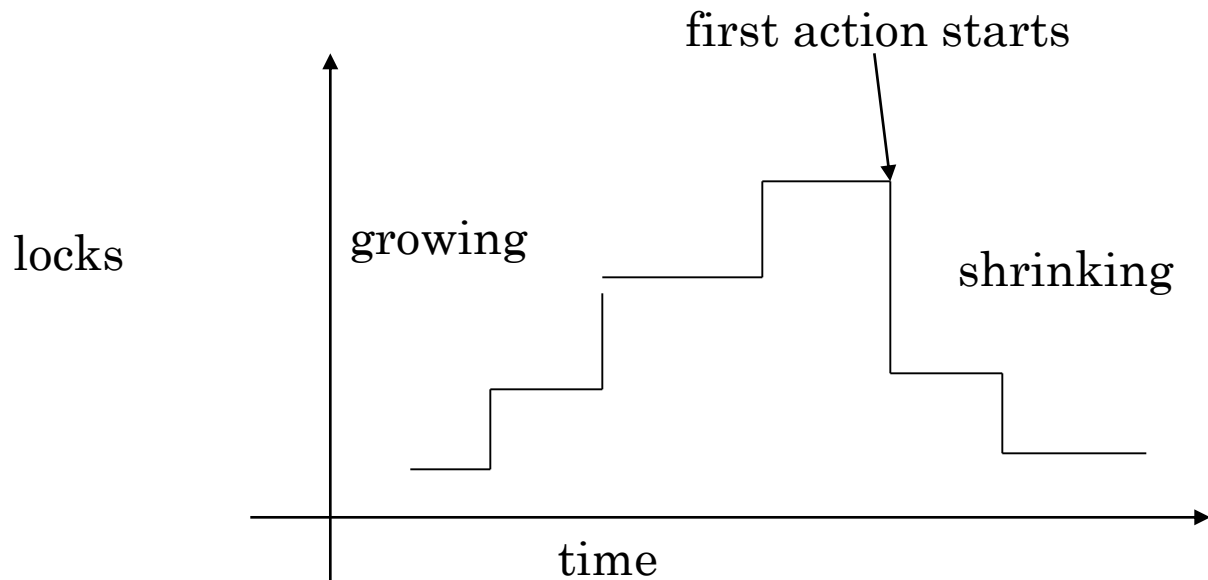  - If an Xact holds an X lock on an object, no other Xact can get a lock (S or X) on that object.
- Strict 2PL allows only serializable schedules.
  - Additionally, it simplifies transaction aborts
  - (Non-strict) 2PL also allows only serializable schedules, but involves more complex abort processing
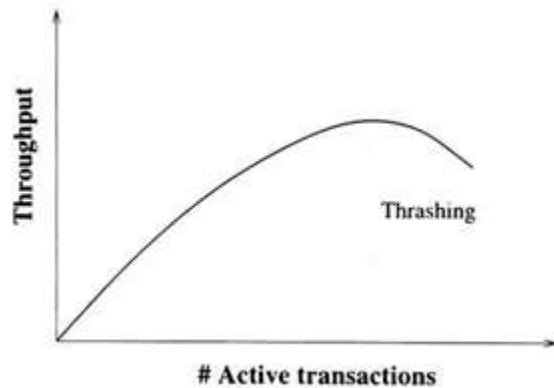
# CONSERVATIVE 2PL

- Lock all items it needs then transaction starts execution
  - If any locks can not be obtained, then do not lock anything
- Difficult but deadlock free

first action starts

growing

locks

shrinking

time

# PERFORMANCE OF LOCKING



Lock Thrashing

•Delays due to blocking increases with the number of active transactions and throughput increases more slowly than number of active transactions.

•Adding another transaction may reduce the throughput at some point of time

•Throughput can be increased in three ways

    •By locking the smallest sized objects possible

    •By reducing the time that transaction hold locks

    •By reducing hotspots

# DEADLOCKS

- Deadlock: Cycle of transactions waiting for locks to be released by each other.
- Two ways of dealing with deadlocks:
  - Deadlock prevention
  - Deadlock detection
  - Deadlock Recovery

| $T_3$ | $T_4$ |
|---|---|
| lock-X($B$) | |
| read($B$) | |
| $B := B - 50$ | |
| write($B$) | |
| | lock-S($A$) |
| | read($A$) |
| | lock-S($B$) |
| lock-X($A$) | |

# Deadlock Prevention

- ***Deadlock prevention*** protocols ensure that the system will *never* enter into a deadlock state. Some prevention strategies:
  - Require that each transaction locks all its data items before it begins execution (pre-declaration).

- Prevention Strategies
  - Assign priorities based on timestamps. Assume Ti wants a lock that Tj holds. Two policies are possible:
    - Wait-Die: If Ti has higher priority, Ti waits for Tj; otherwise Ti aborts(non-preemptive)
    - Wound-wait: If Ti has higher priority, Tj aborts; otherwise Ti waits(preemptive)
  - If a transaction re-starts, make sure it has its original timestamp
  - Timeout-Based Schemes:
    - A transaction waits for a lock only for a specified amount of time. After that, the wait times out and the transaction is rolled back.
    - Ensures that deadlocks get resolved by timeout if they occur

# DEADLOCK DETECTION

- Create a waits-for graph:
  - Nodes are transactions
  - There is an edge from Ti to Tj if Ti is waiting for Tj to release a lock
- Periodically check for cycles in the waits-for graph

| T1 | T2 | T3 | T4 |
|------|------|------|------|
| S(A) | | | |
| R(A) | | | |
| | X(B) | | |
| | W(B) | | |
| 8(B) | | | |
| | | 8(C) | |
| | | R(C) | |
| | X(C) | | |
| | | | X(B) |
| | | X(A) | |

Figure 17.3    Schedule Illustrating Deadlock
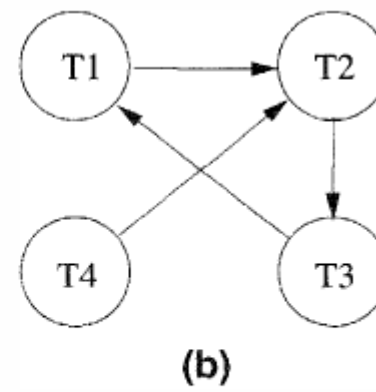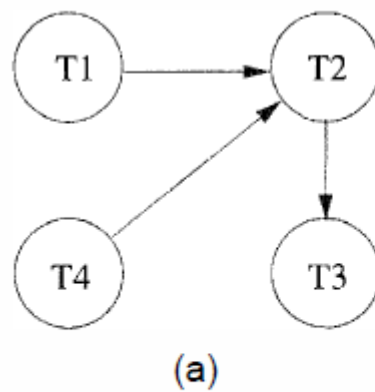


(a)                    (b)

Figure 17.4    Waits-for Graph Before and After Deadlock

# DEADLOCK RECOVERY

- When deadlock is detected :
  - Some transaction will have to rolled back (made a **victim**) to break deadlock cycle.
    - Select that transaction as victim that will incur minimum cost
  - Rollback -- determine how far to roll back transaction
    - **Total rollback**: Abort the transaction and then restart it.
    - **Partial rollback**: Roll back victim transaction only as far as necessary to release locks that another transaction in cycle is waiting for

- Starvation can happen
  - One solution: oldest transaction in the deadlock set is never chosen as victim

# Log-Based Recovery

- A **log** is a sequence of **log records**. The records keep information about update activities on the database.
  - The **log** is kept on stable storage
- When transaction $T_i$ starts, it registers itself by writing a

  **$<T_i$ start$>$** log record

- *Before* $T_i$ executes **write**($X$), a log record

  $<T_i, X, V_1, V_2>$

  is written, where $V_1$ is the value of $X$ before the write (the **old**

  **value**)**,** and $V_2$ is the value to be written to $X$ (the **new value**).

- When $T_i$ finishes it last statement, the log record $<T_i$ **commi**t$>$ is written.
- Two approaches using logs
  - Immediate database modification
  - Deferred database modification.

- The **immediate-modification** scheme allows updates of an uncommitted transaction to be made to the buffer, or the disk itself, before the transaction commits
- Update log record must be written *before* database item is written
  - We assume that the log record is output directly to stable storage
- Output of updated blocks to disk can take place at any time before or after transaction commit
- Order in which blocks are output can be different from the order in which they are written.
- The **deferred-modification** scheme performs updates to buffer/disk only at the time of transaction commit
  - Simplifies some aspects of recovery
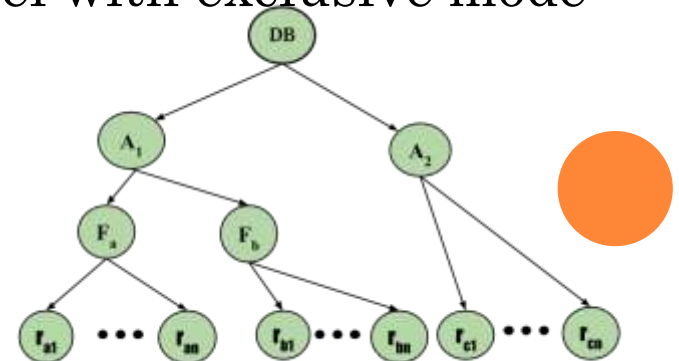  - But has overhead of storing local copy

# INTENT LOCK

- **Intention Mode Lock –**
  In addition to **S** and **X** lock modes, there are three additional lock modes with multiple granularities:
  - **Intention-Shared (IS):** explicit locking at a lower level of the tree but only with shared locks.
  - **Intention-Exclusive (IX):** explicit locking at a lower level with exclusive or shared locks.
  - **Shared & Intention-Exclusive (SIX):** the subtree rooted by that node is locked explicitly in shared mode and explicit locking is being done at a lower level with exclusive mode locks.

The compatibility matrix for these lock modes are described below:

| | IS | IX | S | SIX | X |
|---|---|---|---|---|---|
| **IS** | ✔ | ✔ | ✔ | ✔ | ✘ |
| **IX** | ✔ | ✔ | ✘ | ✘ | ✘ |
| **S** | ✔ | ✘ | ✔ | ✘ | ✘ |
| **SIX** | ✔ | ✘ | ✘ | ✘ | ✘ |
| **X** | ✘ | ✘ | ✘ | ✘ | ✘ |

**IS** : Intention Shared          **X** : Exclusive
**IX** : Intention Exclusive     **SIX** : Shared & Intention Exclusive
**S**   : Shared