



# Chapter 14: Transactions

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# Transaction Concept

- 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



# Required Properties of a Transaction

- Consider a 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
    - 4 Failure could be due to software or hardware
  - 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.



# Required Properties of a Transaction (Cont.)

- **Consistency requirement** in above fund transfer example:
  - The sum of A and B is unchanged by the execution of the transaction
- In general, consistency requirements include
  - 4 Explicitly specified integrity constraints such as primary keys and foreign keys
  - 4 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, when starting to execute, must see a consistent database.
- During transaction execution the database may be temporarily inconsistent.
- When the transaction completes successfully the database must be consistent
  - Erroneous transaction logic can lead to inconsistency



# Required Properties of a Transaction (Cont.)

- **Isolation requirement** — if between steps 3 and 6 (of the fund transfer transaction) , 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**

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

**T2**

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**
  - That is, one after the other.
- However, executing multiple transactions concurrently has significant benefits, as we will see later.



# ACID Properties

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

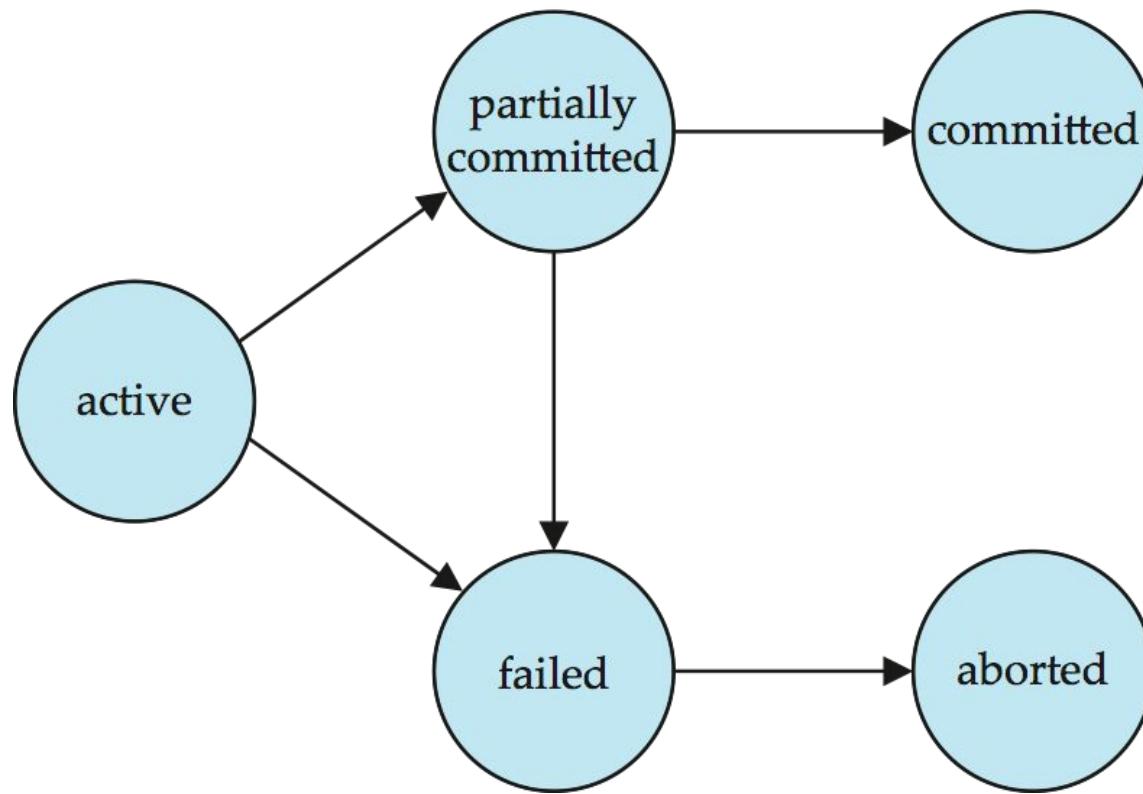


# Transaction State

- **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
    - 4 can be done only if no internal logical error
  - Kill the transaction
- **Committed** – after successful completion.



# Transaction State (Cont.)





# Concurrent Executions

- Multiple transactions are allowed to run concurrently in the system.  
Advantages are:
  - **Increased processor and disk utilization**, leading to better transaction *throughput*
    - 4 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
  - That is, to control the interaction among the concurrent transactions in order to prevent them from destroying the consistency of the database
    - 4 Will study in Chapter 15, after studying notion of correctness of concurrent executions.



# 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** instruction as the last statement
  - By default transaction assumed to execute commit instruction as its last step
- A transaction that fails to successfully complete its execution will have an **abort** instruction as the last statement



# Schedule 1

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



# Schedule 2

- A **serial** schedule in which  $T_2$  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



# 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 ( $B$ ) commit	read ( $B$ ) $B := B + temp$ write ( $B$ ) commit

Note -- In schedules 1, 2 and 3, the sum “ $A + B$ ” is preserved.



# Schedule 4

- The following concurrent schedule does not preserve the sum 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 ( $B$ ) commit	$B := B + temp$ write ( $B$ ) commit



# 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**



# Conflicting Instructions

- Let  $l_i$  and  $l_j$  be two Instructions of transactions  $T_i$  and  $T_j$  respectively. Instructions  $l_i$  and  $l_j$  **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  $Q$ .
  - $l_i = \text{read}(Q)$ ,  $l_j = \text{read}(Q)$ .  $l_i$  and  $l_j$  don't conflict.
  - $l_i = \text{read}(Q)$ ,  $l_j = \text{write}(Q)$ . They conflict.
  - $l_i = \text{write}(Q)$ ,  $l_j = \text{read}(Q)$ . They conflict
  - $l_i = \text{write}(Q)$ ,  $l_j = \text{write}(Q)$ . 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 a 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

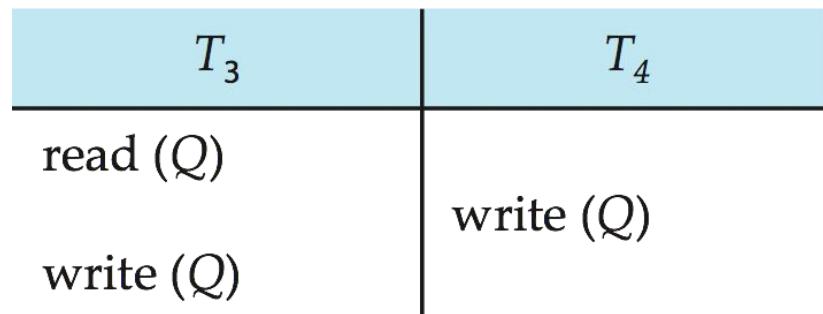
$T_1$	$T_2$
read ( $A$ ) write ( $A$ ) read ( $B$ ) write ( $B$ )	read ( $A$ ) write ( $A$ ) read ( $B$ ) write ( $B$ )

Schedule 6



# Conflict Serializability (Cont.)

- Example of a schedule that is not conflict serializable:

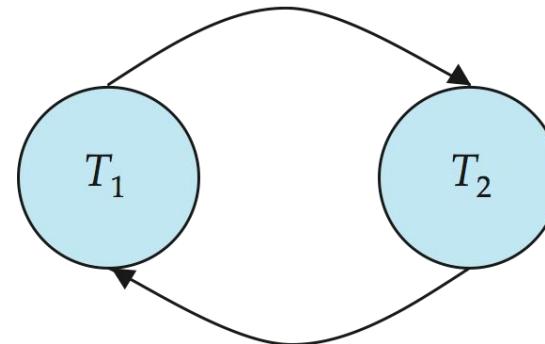


- We are unable to swap instructions in the above schedule to obtain either the serial schedule  $\langle T_3, T_4 \rangle$ , or the serial schedule  $\langle T_4, T_3 \rangle$ .



# Precedence Graph

- Consider some schedule of a set of transactions  $T_1, T_2, \dots, T_n$
- **Precedence graph** — a directed 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**





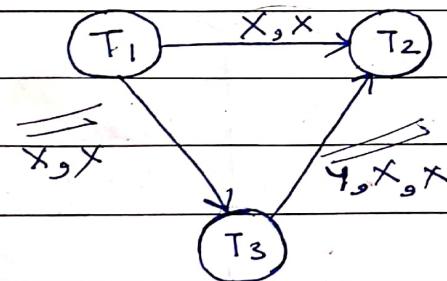
Q. consider the following schedule for transactions  
T<sub>1</sub>, T<sub>2</sub> and T<sub>3</sub>:

[Gate exam]

T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
Read(x)		
	Read(y)	
		Read(y)
	Write(y)	
Write(x)		
		Write(x)
	Read(x)	
	Write(x)	

Time

[gate 2010]



Since there are no cycles  
in the precedence graph  
hence it is conflict  
serializable.

∴ equivalent serial schedule is

$$T_1 \rightarrow T_3 \rightarrow T_2$$



Q.3 consider the following schedules involving two transactions. Which one of the following statement is TRUE? [Gate]

$S_1 : r_1(x); r_1(y); r_2(x); r_2(y); w_2(y); w_1(x)$   
 $S_2 : r_1(x); r_2(x); r_2(y); w_2(y); r_1(y); w_1(x)$

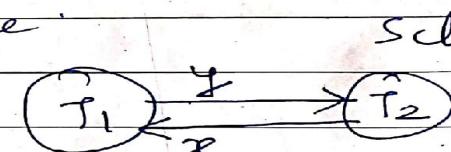
- (A) Both  $S_1$  and  $S_2$  are conflict serializable  
 (B)  $S_1$  is conflict serializable &  $S_2$  is not conflict serializable.

(C)  $S_1$  is not conflict serializable and  $S_2$  is conflict serializable

(D) Both  $S_1$  and  $S_2$  are not conflict serializable.

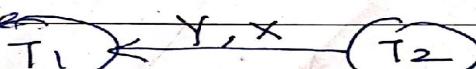
→ Schedule  $S_1$  ans is [C]

$T_1$	$T_2$	Time
$r_1(x)$		
$r_1(y)$	*	
	$r_2(x)$	
	$r_2(y)$	
	$w_2(y)$	
$w_1(x)$		



cycle is present in the above precedence graph, hence it is  $(S_1)$  not conflict serializable.

$T_1$	$T_2$	Time
$r_1(x)$		
	$r_2(x)$	
	$r_2(y)$	
	$w_2(y)$	
$r_1(y)$		
$w_1(x)$		



No cycle is present in also precedence graph

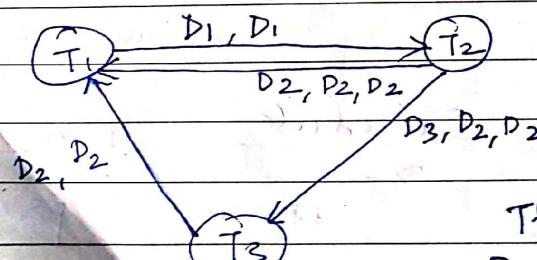
$S_2$  hence it is conflict serializable



Q. Consider three data items  $D_1$ ,  $D_2$  and  $D_3$  and the following execution schedule of transaction  $T_1$ ,  $T_2$  and  $T_3$ . In the diagram,  $R(D)$  and  $W(D)$  denote the actions ~~regarding~~ reading & writing the data item  $D$ , respectively.

M.0

$T_1$	$T_2$	$T_3$
	$R(D_3)$	
	$R(D_2)$	
	$W(D_2)$	
$R(D_1)$		
$W(D_1)$		
<del>Read D2</del>	<del>Read D1</del>	
		$R(D_2)$
		$R(D_3)$
		$W(D_2)$
		$W(D_3)$
	$R(D_1)$	
		$W(D_1)$



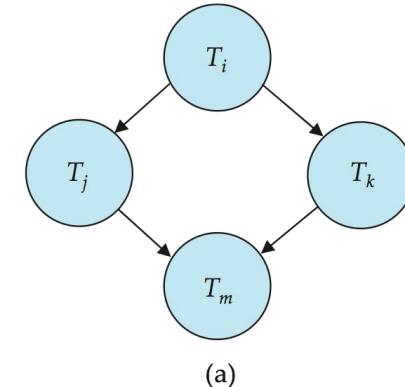
Serializable.

There are cycles in the Precedence Graph, hence it is not -

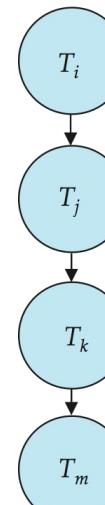


# Testing for Conflict Serializability

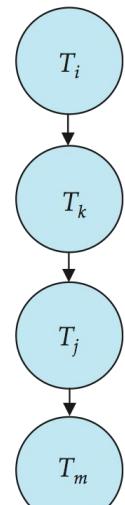
- A schedule is conflict serializable if and only if its precedence graph is acyclic.
- Cycle-detection algorithms exist which take order  $n^2$  time, where  $n$  is the number of vertices in the graph.
  - (Better algorithms take order  $n + e$  where  $e$  is the number of edges.)
- If precedence graph is acyclic, the serializability order can be obtained by a *topological sorting* of the graph.
  - That is, a linear order consistent with the partial order of the graph.
  - For example, a serializability order for the schedule (a) would be one of either (b) or (c)



(a)



(b)



(c)



# Recoverable Schedules

- **Recoverable schedule** — if a transaction  $T_j$  reads a data item previously written by a transaction  $T_i$ , then the commit operation of  $T_i$  **must** appear before the commit operation of  $T_j$ .
- The following schedule is not recoverable if  $T_9$  commits immediately after the read(A) operation.

$T_8$	$T_9$
read (A) write (A)	
read (B)	read (A) commit

- 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 ( $A$ )  abort	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** — 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
- Example of a schedule that is NOT cascadeless

$T_{10}$	$T_{11}$	$T_{12}$
read ( $A$ ) read ( $B$ ) write ( $A$ )		
	read ( $A$ ) write ( $A$ )	read ( $A$ )
abort		



# Concurrency Control

- A database must provide a mechanism that will ensure that all possible schedules are both:
  - Conflict serializable.
  - 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
- Concurrency-control schemes tradeoff between the amount of concurrency they allow and the amount of overhead that they incur
- Testing a schedule for serializability *after* it has executed is a little too late!
  - Tests for serializability help us understand why a concurrency control protocol is correct
- **Goal – to develop concurrency control protocols that will assure serializability.**



# Weak Levels of Consistency

- Some applications are willing to live with weak levels of consistency, allowing schedules that are not serializable
  - E.g., a read-only transaction that wants to get an approximate total balance of all accounts
  - E.g., database statistics computed for query optimization can be approximate (why?)
  - Such transactions need not be serializable with respect to other transactions
- Tradeoff accuracy for performance



# Levels of Consistency in SQL-92

- **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, but successive reads of record may return different (but committed) values.
  - **Read uncommitted** — even uncommitted records may be read.
- 
- Lower degrees of consistency useful for gathering approximate information about the database
  - Warning: some database systems do not ensure serializable schedules by default
    - E.g., Oracle and PostgreSQL by default support a level of consistency called snapshot isolation (not part of the SQL standard)



# Transaction Definition in SQL

- Data manipulation language must include a construct for specifying the set of actions that comprise a transaction.
- 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
    - 4 E.g. in JDBC, `connection.setAutoCommit(false);`



# Other Notions of Serializability



# 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, for each data item  $Q$ ,
  - 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  $Q$ .
  - 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  $T_j$ .
  - The transaction (if any) that performs the final **write( $Q$ )** operation in schedule  $S$  must also perform the final **write( $Q$ )** operation in schedule  $S'$ .
- As can be seen, view equivalence is also based purely on **reads** and **writes** alone.



# View Serializability (Cont.)

- 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**.



# 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.



# Chapter 15 : Concurrency Control

Database System Concepts, 6<sup>th</sup> Ed.

©Silberschatz, Korth and Sudarshan

See [www.db-book.com](http://www.db-book.com) for conditions on re-use



# 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 the concurrency-control manager by the programmer. Transaction can proceed only after request is granted.
- 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.



# 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:

```
 $T_2$ : 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.



# The Two-Phase Locking Protocol

- This protocol 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).



# Lock Conversions

- Two-phase locking with lock conversions:
  - First Phase:
    - can acquire a lock-S on item
    - can acquire a lock-X on item
    - can convert a lock-S to a lock-X (upgrade)
  - Second Phase:
    - can release a lock-S
    - can release a lock-X
    - can convert a lock-X to a lock-S (downgrade)
- This protocol assures serializability. But still relies on the programmer to insert the various locking instructions.



# Automatic Acquisition of Locks

- A transaction  $T_i$  issues the standard read/write instruction, without explicit locking calls.
- The operation **read( $D$ )** is processed as:

```
if  $T_i$  has a lock on  $D$ 
  then
    read( $D$ )
  else begin
    if necessary wait until no other
      transaction has a lock-X on  $D$ 
    grant  $T_i$  a lock-S on  $D$ ;
    read( $D$ )
  end
```



# Automatic Acquisition of Locks (Cont.)

- **write( $D$ )** is processed as:  
**if**  $T_i$  has a **lock-X** on  $D$   
**then**  
    **write( $D$ )**  
**else begin**  
    if necessary wait until no other transaction has any lock on  $D$ ,  
    **if**  $T_i$  has a **lock-S** on  $D$   
        **then**  
            **upgrade lock on  $D$  to lock-X**  
        **else**  
            **grant  $T_i$  a lock-X on  $D$**   
            **write( $D$ )**  
**end;**
- All locks are released after commit or abort



# Deadlocks

- 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.



Q. 15.2

10/2/12

Color?

Ans:

(T34)  
lock-S(A)

read(A)

lock-X(B)

read(B)

if A=0

then B:=B+1

write(B)

unlock(A)

unlock(B)

T34:

read(A);

read(B);

if A=0 then B:=B+1;

write(B).

T35: read(B);

read(A);

if B=0 then A:=A+1;

write(A) to

(T35)  
lock-S(B)

read(B)

lock-X(A)

read(A)

if B=0

then A:=A+1

write(A)

unlock(B)

unlock(A)

add lock and unlock instructions T34 and T35, so  
that they observe the two phase locking protocol.  
Can the execution of these transactions result  
in a deadlock.

T34

lock-S(A)

read(A)

unlock-S(A)

lock-X(B)

read(B)

if A=0 then B:=B+1

write(B)

unlock-X(B)

T35

lock-S(B)

read(B)

unlock-S(B)

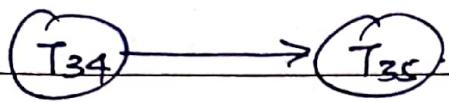
lock-X(A)

read(A)

if B=0 then A:=A+1;

write(A)

unlock-X(A)



Serial execution of Transactions T<sub>34</sub> & T<sub>35</sub> does not result in deadlock.

10/2/12

### 1) Recoverable schedules :-

T <sub>6</sub>	T <sub>7</sub>
read(A) write(A)	
read(B)	

T <sub>34</sub>	T <sub>35</sub>
lock-S(A)	
read(A) lock-X(B)	lock-S(B) read(B)

lock-X(A)  
partial schedule

partial schedule

lege of E

Ti

ST3

P32

+

W

O



# Deadlocks (Cont.)

- Two-phase locking *does not* ensure freedom from deadlocks.
- In addition to deadlocks, there is a possibility of **starvation**.
- **Starvation** occurs if the 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.



# Deadlocks (Cont.)

- The potential for deadlock exists in most locking protocols. Deadlocks are a necessary evil.
- When a deadlock occurs there is a possibility of cascading roll-backs.
- Cascading roll-back is possible under two-phase locking. To avoid this, follow a modified protocol called **strict two-phase locking** -- 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.

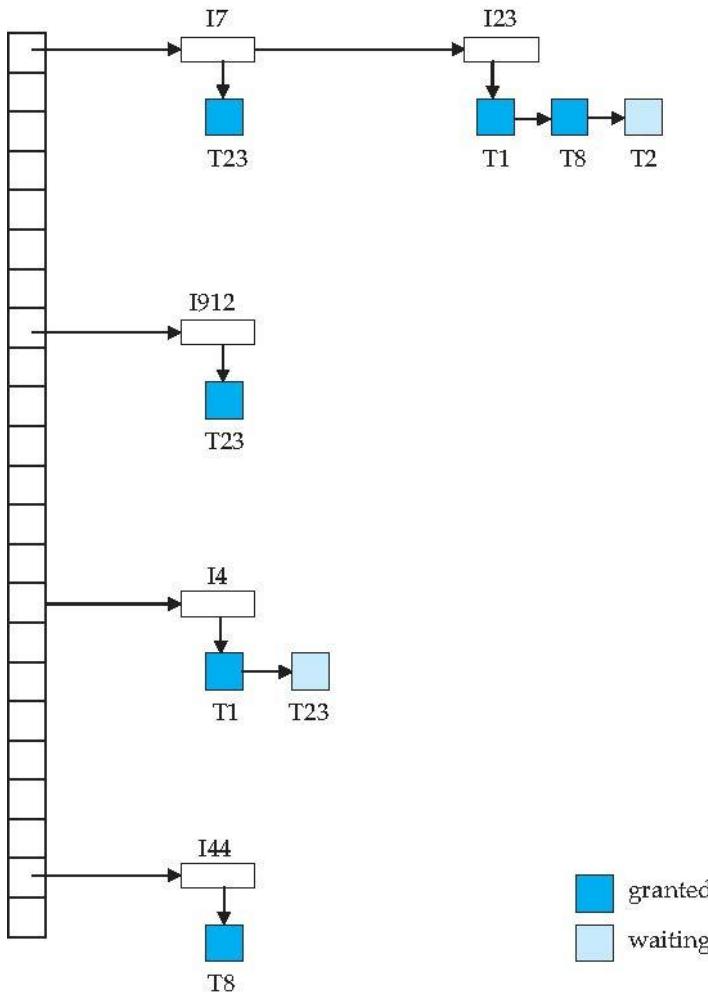


# Implementation of Locking

- A **lock manager** can be implemented as a separate process to which transactions send lock and unlock requests
- The lock manager replies to a lock request by sending a lock grant messages (or a message asking the transaction to roll back, in case of a deadlock)
- The requesting transaction waits until its request is answered
- The lock manager maintains a data-structure called a **lock table** to record granted locks and pending requests
- The lock table is usually implemented as an in-memory hash table indexed on the name of the data item being locked



# Lock Table



- Dark blue rectangles indicate granted locks; light blue indicate waiting requests
- Lock table also records the type of lock granted or requested
- New request is added to the end of the queue of requests for the data item, and granted if it is compatible with all earlier locks
- Unlock requests result in the request being deleted, and later requests are checked to see if they can now be granted
- If transaction aborts, all waiting or granted requests of the transaction are deleted
  - lock manager may keep a list of locks held by each transaction, to implement this efficiently



# Deadlock Handling

- System is deadlocked if there is a set of transactions such that every transaction in the set is waiting for another transaction in the set.
- **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 (predeclaration).
  - Impose partial ordering of all data items and require that a transaction can lock data items only in the order specified by the partial order.

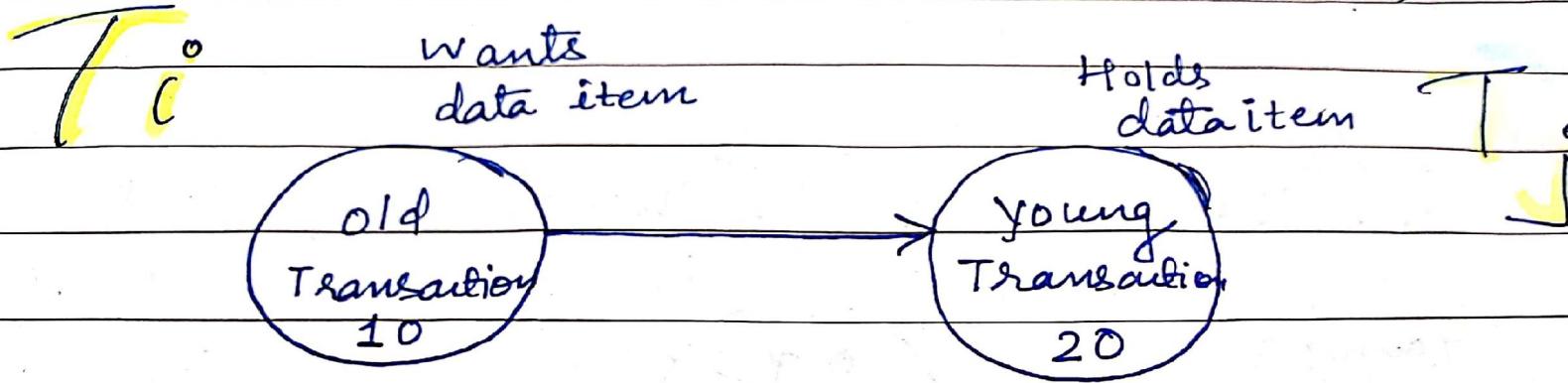


# More Deadlock Prevention Strategies

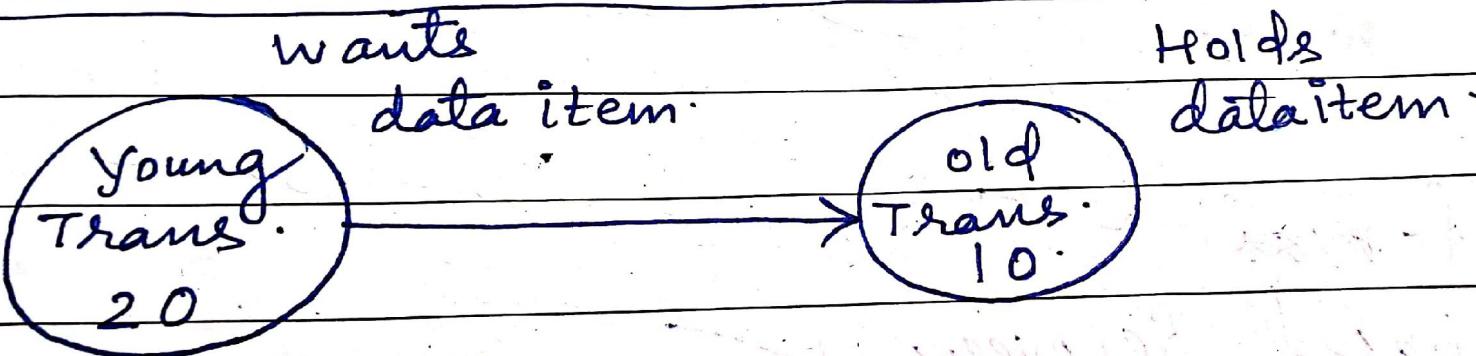
- Following schemes use transaction timestamps(TS) for the sake of deadlock prevention alone.
- **wait-die** scheme — non-preemptive(**Scheme number 1**)
  - older transaction may wait for younger one to release data item. (older means smaller timestamp) Younger transactions never wait for older ones; they are rolled back instead.
  - a transaction may die several times before acquiring needed data item
- **Wait-die. (Navathe book)**

If  $TS(T_i) < TS(T_j)$ , then ( $T_i$  older than  $T_j$ )  $T_i$  is allowed to wait; otherwise ( $T_i$  younger than  $T_j$ ) abort  $T_i$  ( $T_i$  dies) and restart it later with the same timestamp.

otherwise  $T_i$  is rolled back (dies).



waits



Dies

No preemption technique  
in wait-die.



# More Deadlock Prevention Strategies

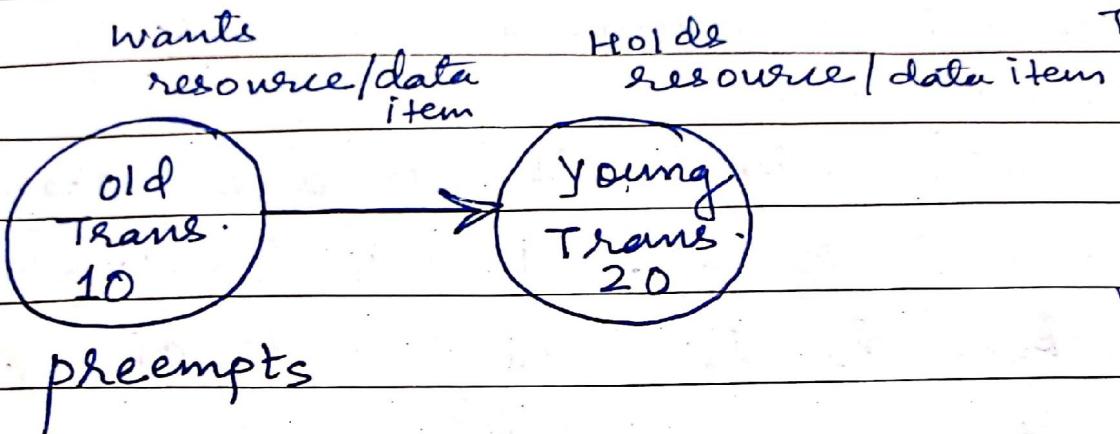
- **wound-wait** scheme — preemptive (**Scheme number 2**)

- older transaction *wounds* (forces rollback) of younger transaction instead of waiting for it.
- Younger transactions may wait for older ones. may be fewer rollbacks than *wait-die* scheme.

- **Wound-wait. (Navathe book)**

If  $TS(T_i) < TS(T_j)$ , then ( $T_i$  older than  $T_j$ ) abort  $T_j$  ( $T_i$  wounds  $T_j$ ) and restart it later with the same timestamp; otherwise ( $T_i$  younger than  $T_j$ )  $T_i$  is allowed to wait.

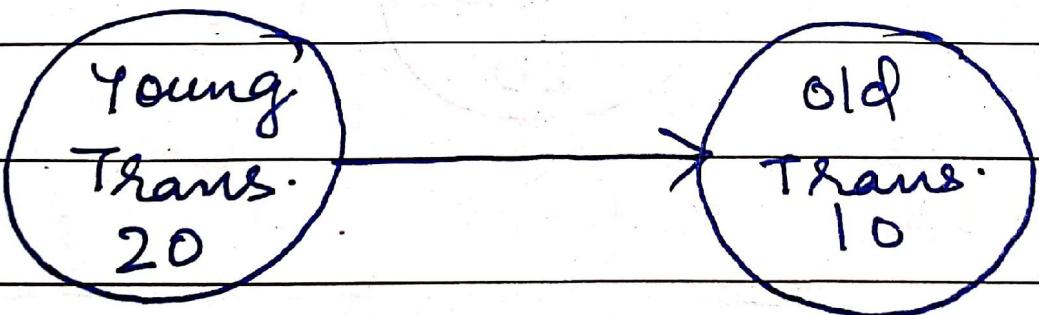
(That is,  $T_i$  is younger than  $T_j$ ) otherwise  
 $T_j$  is rolled back



preempt is  
abnormal is  
possible

wants  
data item

Holds  
data item



Waits.

wound-wait ✓



# Summary of Deadlock Prevention

- In wait-die, an older transaction is allowed to wait for a younger transaction, whereas a younger transaction requesting an item held by an older transaction is aborted and restarted.
- The wound-wait approach does the opposite: A younger transaction is allowed to wait for an older one, whereas an older transaction requesting an item held by a younger transaction preempts the younger transaction by aborting it. Both schemes end up aborting the younger of the two transactions (the transaction that started later) that may be involved in a deadlock, assuming that this will waste less processing. It can be shown that these two techniques are deadlock-free, since in wait-die, transactions only wait for younger transactions so no cycle is created. Similarly, in wound-wait, transactions only wait for older transactions so no cycle is created. However, both techniques may cause some transactions to be aborted and restarted needlessly, even though those transactions may never actually cause a deadlock.



# Deadlock prevention (Cont.)

- Both in *wait-die* and in *wound-wait* schemes, a rolled back transaction is restarted with its original timestamp. Older transactions thus have precedence over newer ones, and starvation is hence avoided.
- **Timeout-Based Schemes:**
  - a transaction waits for a lock only for a specified amount of time. If the lock has not been granted within that time, the transaction is rolled back and restarted,
  - Thus, deadlocks are not possible
  - simple to implement; but starvation is possible. Also difficult to determine good value of the timeout interval.

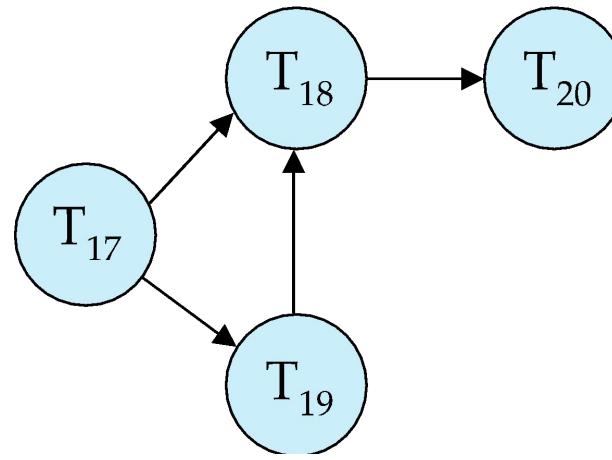


# Deadlock Detection

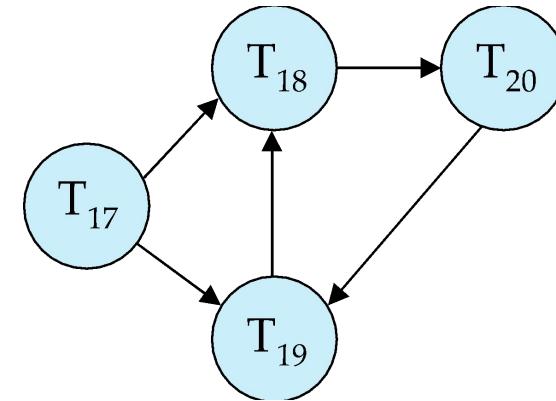
- Deadlocks can be described as a *wait-for graph*, which consists of a pair  $G = (V,E)$ ,
  - $V$  is a set of vertices (all the transactions in the system)
  - $E$  is a set of edges; each element is an ordered pair  $T_i \rightarrow T_j$ .
- If  $T_i \rightarrow T_j$  is in  $E$ , then there is a directed edge from  $T_i$  to  $T_j$ , implying that  $T_i$  is waiting for  $T_j$  to release a data item.
- When  $T_i$  requests a data item currently being held by  $T_j$ , then the edge  $T_i \rightarrow T_j$  is inserted in the wait-for graph. This edge is removed only when  $T_j$  is no longer holding a data item needed by  $T_i$ .
- The system is in a deadlock state if and only if the wait-for graph has a cycle. Must invoke a deadlock-detection algorithm periodically to look for cycles.



# Deadlock Detection (Cont.)



Wait-for graph without a cycle



Wait-for graph with a cycle



# Deadlock Recovery

- When deadlock is detected :
  - Some transaction will have to rolled back (made a victim) to break deadlock. Select that transaction as victim that will incur minimum cost.
  - Rollback -- determine how far to roll back transaction
    - 4 Total rollback: Abort the transaction and then restart it.
    - 4 More effective to roll back transaction only as far as necessary to break deadlock.
  - Starvation happens if same transaction is always chosen as victim. Include the number of rollbacks in the cost factor to avoid starvation