

Concurrency Control and Recovery Techniques, NOSQL Management

Dr. Jyotismita Chaki

Why concurrency control

- Assume that two people who go to electronic kiosks at the same time to buy a movie ticket for the same movie and the same show time.
- However, there is only one seat left in for the movie show in that particular theatre.
- Without concurrency control in DBMS, it is possible that both moviegoers will end up purchasing a ticket.
- However, concurrency control method does not allow this to happen.
- Both moviegoers can still access information written in the movie seating database.
- But concurrency control only provides a ticket to the buyer who has completed the transaction process first.

Why concurrency control

- Several problems can occur when concurrent transactions execute in an uncontrolled manner.
- Referring to a much simplified airline reservations database in which a record is stored for each airline flight.
- Each record includes the number of reserved seats on that flight as a named (uniquely identifiable) data item, among other information.
- Figure (a) shows a transaction T1 that transfers N reservations from one flight whose number of reserved seats is stored in the database item named X to another flight whose number of reserved seats is stored in the database item named Y.
- Figure (b) shows a simpler transaction T2 that just reserves M seats on the first flight (X) referenced in transaction T1.

(a)

T_1
read_item(X);
$X := X - N;$
write_item(X);
read_item(Y);
$Y := Y + N;$
write_item(Y);

(b)

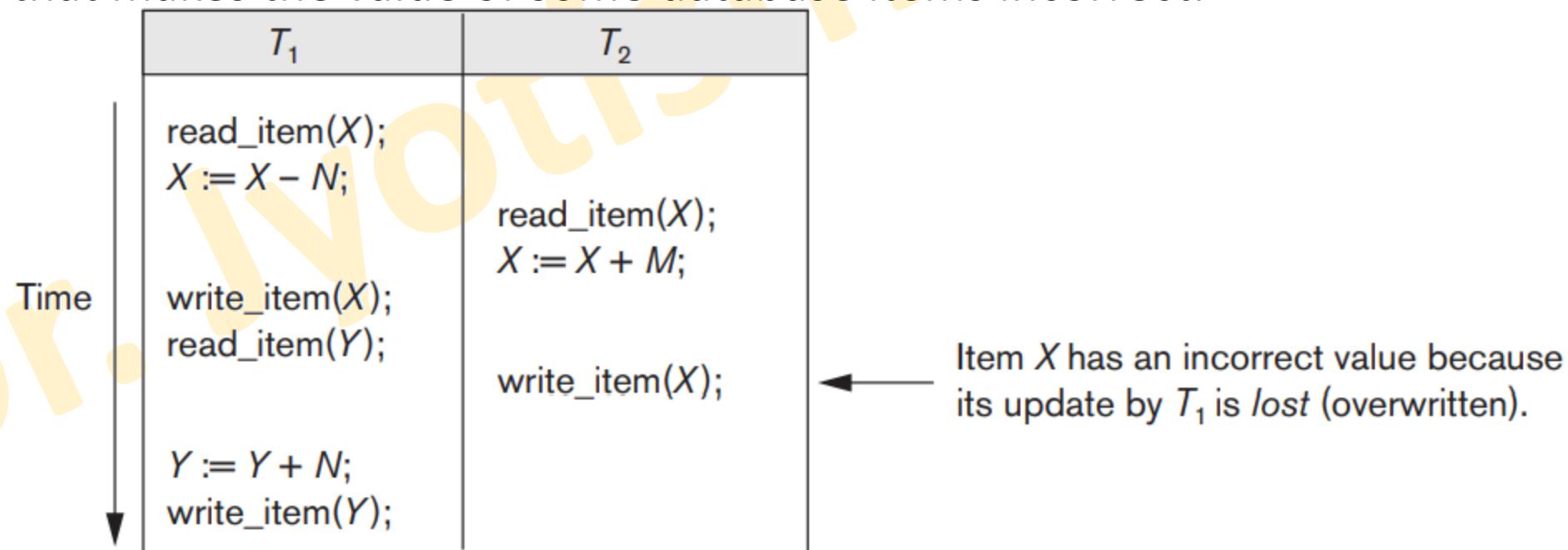
T_2
read_item(X);
$X := X + M;$
write_item(X);

Why concurrency control

- When a database access program is written, it has the flight number, the flight date, and the number of seats to be booked as parameters; hence, the same program can be used to execute many different transactions, each with a different flight number, date, and number of seats to be booked.
- For concurrency control purposes, a transaction is a particular execution of a program on a specific date, flight, and number of seats. In Figure (a) and (b), the transactions T1 and T2 are specific executions of the programs that refer to the specific flights whose numbers of seats are stored in data items X and Y in the database.

Why Concurrency control: Problems

- The types of problems we may encounter with these two simple transactions if they run concurrently.
 - **The Lost Update Problem:** This problem occurs when two transactions that access the same database items have their operations interleaved in a way that makes the value of some database items incorrect.



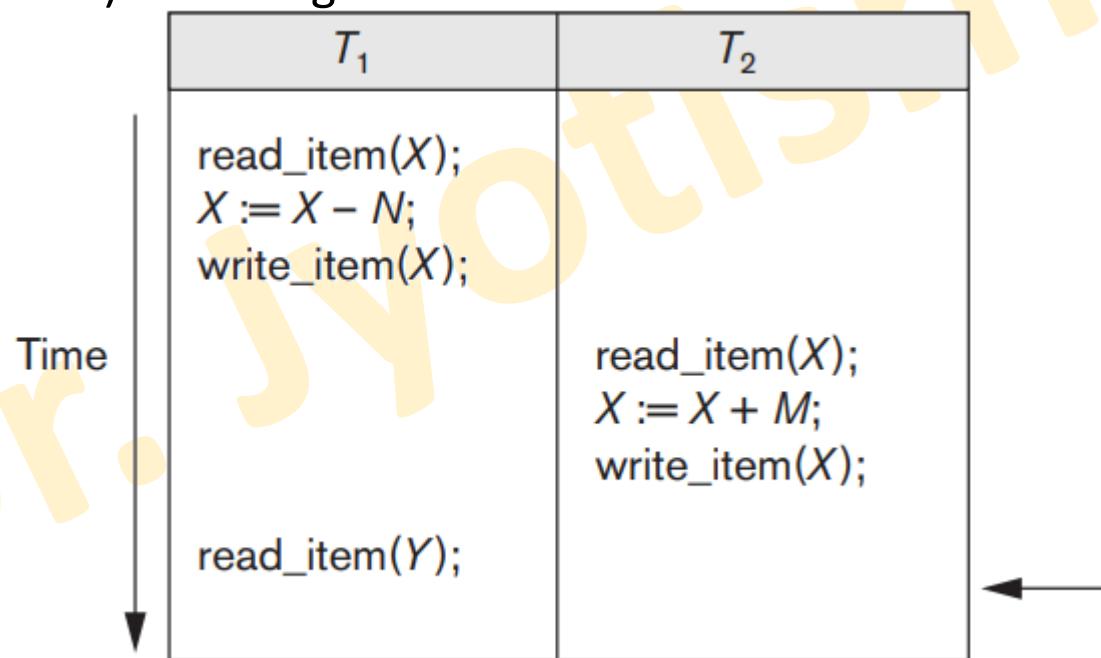
Why Concurrency control: Problems: The Lost Update Problem Example

Time	T ₁	T ₂	Comments
1	read_item(X);		Read operation is performed in T ₁ at time step 1.
2	X:=X-N;		Value of data item X is modified.
3		read_item(X);	Read in value of X (Which value is read in?)
4		X:=X+M	Value of data item X is modified.
5	write_item(X);		Write operation is performed in T ₁ .
6	read_item(Y);		Read operation is performed in T ₁ .
7		write_item(X);	Write data item X to database (What value is written in?)
8	Y:=Y+N;		Value of data item is modified.
9	write_item(Y);		Write data item Y to database.

Time	T ₁	T ₂	Value
1	read_item(X);		X = 80
2	X:=X-N;		X = 80 - 5 = 75 (which is not written into database)
3		read_item(X);	X = 80 (T ₂ still reads in the original value of X, the updated value of X is lost.)
4		X:=X+M	X = 80 + 4 = 84
5	write_item(X);		X = 75 is written into database
6	read_item(Y);		
7		write_item(X);	X = 84 over writes X = 75, a wrong record is written in database
8	Y:=Y+N;		
9	write_item(Y);		

Why Concurrency control: Problems

- The types of problems we may encounter with these two simple transactions if they run concurrently.
 - The Temporary Update (or Dirty Read) Problem:** This problem occurs when one transaction updates a database item and then the transaction fails for some reason. Meanwhile, the updated item is accessed (read) by another transaction before it is changed back (or rolled back) to its original value.



Transaction T_1 fails and must change the value of X back to its old value; meanwhile T_2 has read the *temporary* incorrect value of X .

Why Concurrency control: Problems: The Temporary Update (or Dirty Read) Problem: Example

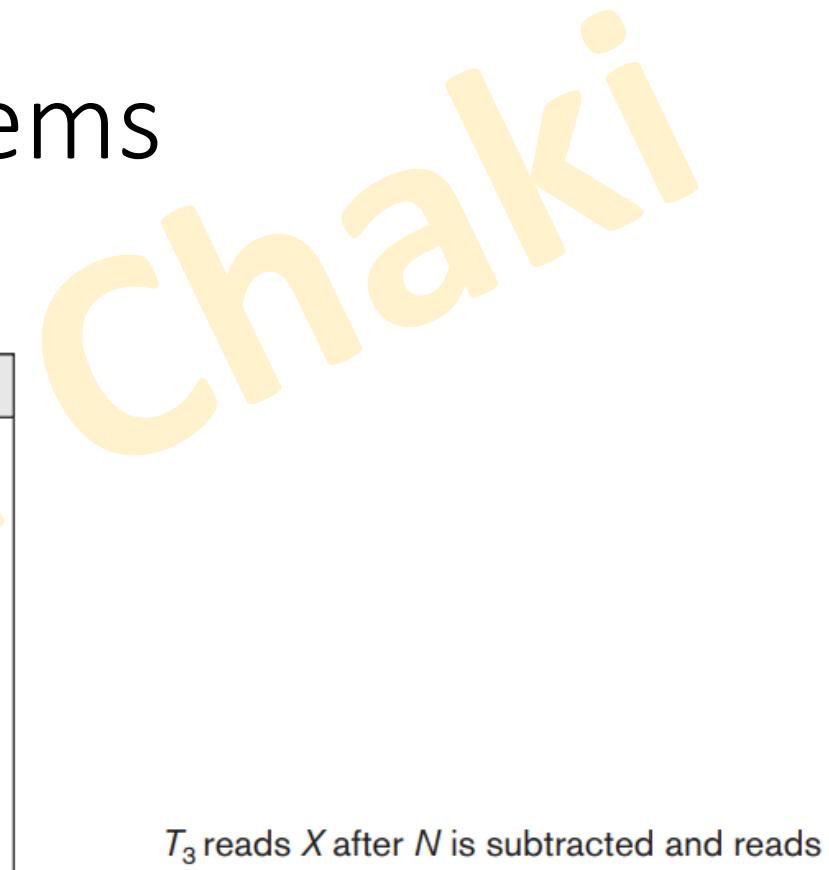
- The figure below shows an example where T1 updates item X and then fails before completion, so the system must change X back to its original value. Before it can do so, however, transaction T2 reads the 'temporary' value of X, which will not be recorded permanently in the database because of the failure of T1.
- The value of item X that is read by T2 is called dirty data, because it has been created by a transaction that has not been completed and committed yet; hence this problem is also known as the dirty read problem.
- Since the dirty data read in by T2 is only a temporary value of X, the problem is sometimes called temporary update too.

Time	T ₁	T ₂	Comment
1	read_item(X);		
2	X:=X-N;		
3	write_item(X);		X is temporarily updated
4		read_item(X);	
5		X:=X+M	
6		write_item(X);	
7	read_item(Y);		
...	
	ROLLBACK		T ₁ fails and must change the value of X back to its old value; meanwhile T ₂ has read the temporary incorrect value of X

Why Concurrency control: Problems

- The types of problems we may encounter with these two simple transactions if they run concurrently.
 - **The Incorrect Summary Problem:** If one transaction is calculating an aggregate summary function on a number of database items while other transactions are updating some of these items, the aggregate function may calculate some values before they are updated and others after they are updated.

T_1	T_3
read_item(X); $X := X - N$; write_item(X);	$sum := 0$; read_item(A); $sum := sum + A$; ⋮
read_item(Y); $Y := Y + N$; write_item(Y);	read_item(X); $sum := sum + X$; read_item(Y); $sum := sum + Y$;

T_3 reads X after N is subtracted and reads Y before N is added; a wrong summary is the result (off by N). 

Why Concurrency control: Problems: The Incorrect Summary Problem Example

- Consider the schedule S1 given below, in which, transaction T1 transfers money from account A to account B and in the mean time, transaction T2 calculates the sum of 3 accounts namely, A, B, and C. The third column shows the account balances and calculated values after every instruction is executed.

Transaction T2 reads the value of account A after A is updated and reads B before B is updated. [The portion that violates in T2 is highlighted in green color]. Hence, the aggregate operation ends up with an inconsistent result.

If all the instructions in T1 are executed before T2 starts, then A will be 950, B will be 1050 and average value will be 1000.

If all the instructions in T1 are executed after T2 finishes, then A will be 950, B will be 1050 and average value will be 1000.

But, due to this interleaved execution, the final value of A is 950, B is 1050, and average is 983.33 which is wrong.

Transaction T1	Transaction T2	A = 1000, B = 1000, C = 1000
<code>read(A); A := A - 50; write(A);</code>	<code>sum = 0; avg = 0; read(C); sum := sum + C;</code> <code>read(A); sum := sum + A; read(B); sum := sum + B; avg := sum/3; commit;</code>	<code>sum = 0 avg = 0 T2 read: C = 1000 sum = 1000 T1 read: A = 1000</code> <code>T1 write: A = 950 T2 read: A = 950 sum = 1950 t2 read: B = 1000 sum = 2950 avg = 983.33</code>
<code>read(B); B := B + 50; write(B); commit;</code>		<code>T2 read: B = 1000 T2 write: B = 1050</code>

Need for locking

- Some of the main techniques used to control concurrent execution of transactions are based on the concept of locking data items.
- A **lock** is a variable associated with a data item that describes the status of the item with respect to possible operations that can be applied to it.
- Generally, there is one lock for each data item in the database.
- Locks are used as a means of synchronizing the access by concurrent transactions to the database items.

Lock Based Protocols: Binary Locks

- Several types of locks are used in concurrency control.
 - **Binary Locks:**
 - A **binary lock** can have two **states** or **values**: **locked** and **unlocked** (or 1 and 0, for simplicity). A distinct lock is associated with each database item X.
 - If the value of the lock on X is 1, item X cannot be accessed by a database operation that requests the item. If the value of the lock on X is 0, the item can be accessed when requested, and the lock value is changed to 1.
 - We refer to the current value (or state) of the lock associated with item X as **lock(X)**.
 - Two operations, **lock_item** and **unlock_item**, are used with binary locking.
 - A transaction requests access to an item X by first issuing a **lock_item(X)** operation.
 - If $\text{LOCK}(X) = 1$, the transaction is forced to wait. If $\text{LOCK}(X) = 0$, it is set to 1 (the transaction locks the item) and the transaction is allowed to access item X.
 - When the transaction is through using the item, it issues an **unlock_item(X)** operation, which sets $\text{LOCK}(X)$ back to 0 (unlocks the item) so that X may be accessed by other transactions.
 - Hence, a binary lock enforces mutual exclusion on the data item.

Lock Based Protocols: Binary Locks

- If the simple binary locking scheme described here is used, every transaction must obey the following rules:
 - A transaction T must issue the operation `lock_item(X)` before any `read_item(X)` or `write_item(X)` operations are performed in T.
 - A transaction T must issue the operation `unlock_item(X)` after all `read_item(X)` and `write_item(X)` operations are completed in T.
 - A transaction T will not issue a `lock_item(X)` operation if it already holds the lock on item X.
 - A transaction T will not issue an `unlock_item(X)` operation unless it already holds the lock on item X.

Lock Based Protocols: Shared/Exclusive (or Read/Write) Locks

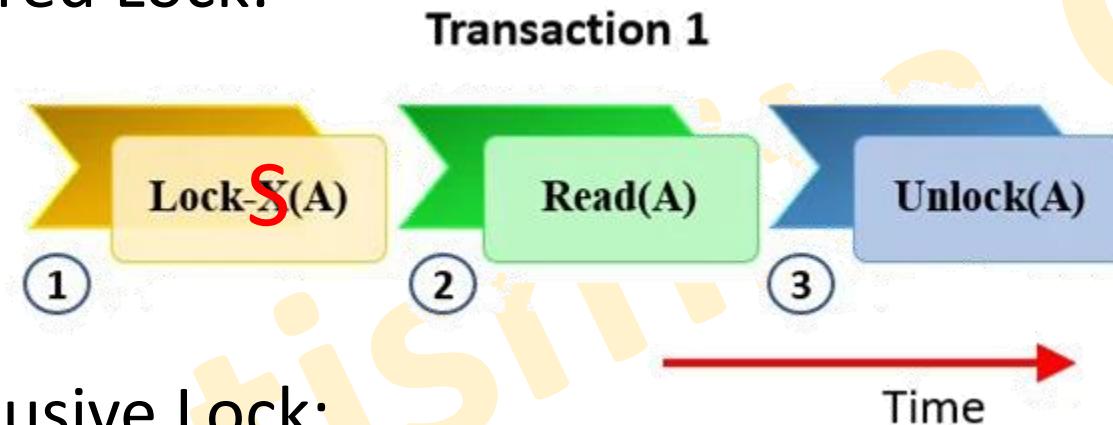
- The preceding binary locking scheme is too restrictive for database items because at most one transaction can hold a lock on a given item.
- We should allow several transactions to access the same item X if they all access X for reading purposes only.
- This is because read operations on the same item by different transactions are not conflicting.
- However, if a transaction is to write an item X, it must have exclusive access to X.
- For this purpose, a different type of lock, called a **multiple-mode lock**, is used.
- In this scheme—called **shared/exclusive** or **read/write locks**—there are three locking operations: `read_lock(X)`, `write_lock(X)`, and `unlock(X)`.

Lock Based Protocols: Shared/Exclusive (or Read/Write) Locks

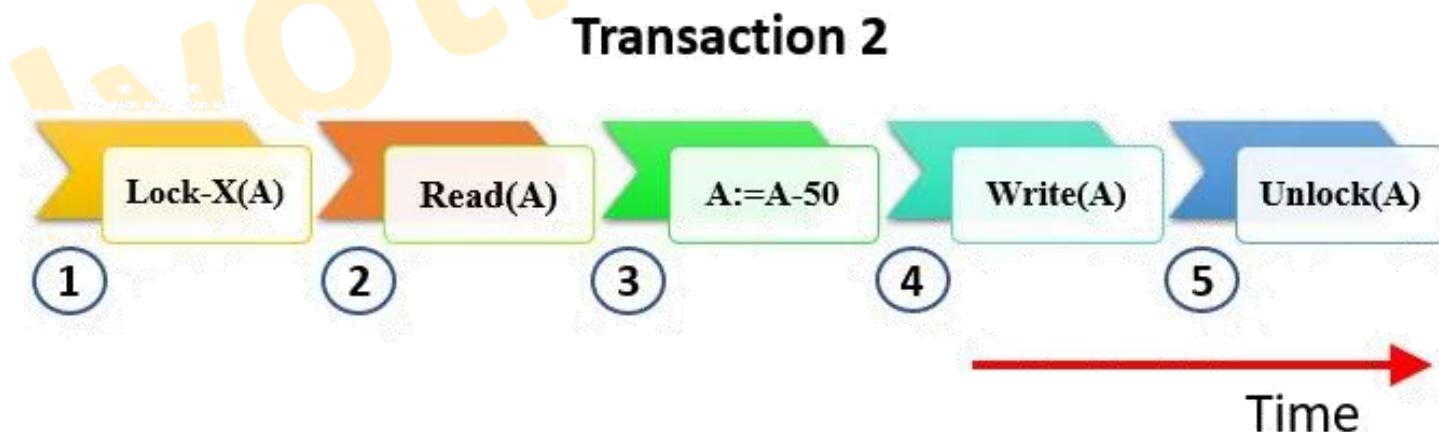
- A lock associated with an item X, $\text{LOCK}(X)$, now has three possible states: *read-locked*, *write-locked*, or *unlocked*.
- A **read-locked item** is also called **share-locked** (Often represented as **lock-S()**) because other transactions are allowed to read the item, whereas a **write-locked item** is called **exclusive-locked** (Often represented as **lock-X()**) because a single transaction exclusively holds the lock on the item.
- Multiple read locks can exist at the same time.
- The write transaction should wait for read locks to finish reading.
- If you have the write lock before the read lock, the write lock will block other transactions to read or write the same table. If you have the read lock before the write lock, the read lock will block the write transactions until the reading transaction finishes.

Lock Based Protocols: Shared/Exclusive (or Read/Write) Locks

- Example of Shared Lock:



- Example of Exclusive Lock:



Lock Based Protocols: Shared/Exclusive (or Read/Write) Locks

- When we use the shared/exclusive locking scheme, the system must enforce the following rules:
 1. A transaction T must issue the operation `read_lock(X)` or `write_lock(X)` before any `read_item(X)` operation is performed in T.
 2. A transaction T must issue the operation `write_lock(X)` before any `write_item(X)` operation is performed in T.
 3. A transaction T must issue the operation `unlock(X)` after all `read_item(X)` and `write_item(X)` operations are completed in T.
 4. A transaction T will not issue a `read_lock(X)` operation if it already holds a read (shared) lock or a write (exclusive) lock on item X.
 5. A transaction T will not issue a `write_lock(X)` operation if it already holds a read (shared) lock or write (exclusive) lock on item X.
 6. A transaction T will not issue an `unlock(X)` operation unless it already holds a read (shared) lock or a write (exclusive) lock on item X.

Lock Based Protocols: Shared/Exclusive (or Read/Write) Locks: Conversion (Upgrading, Downgrading) of Locks

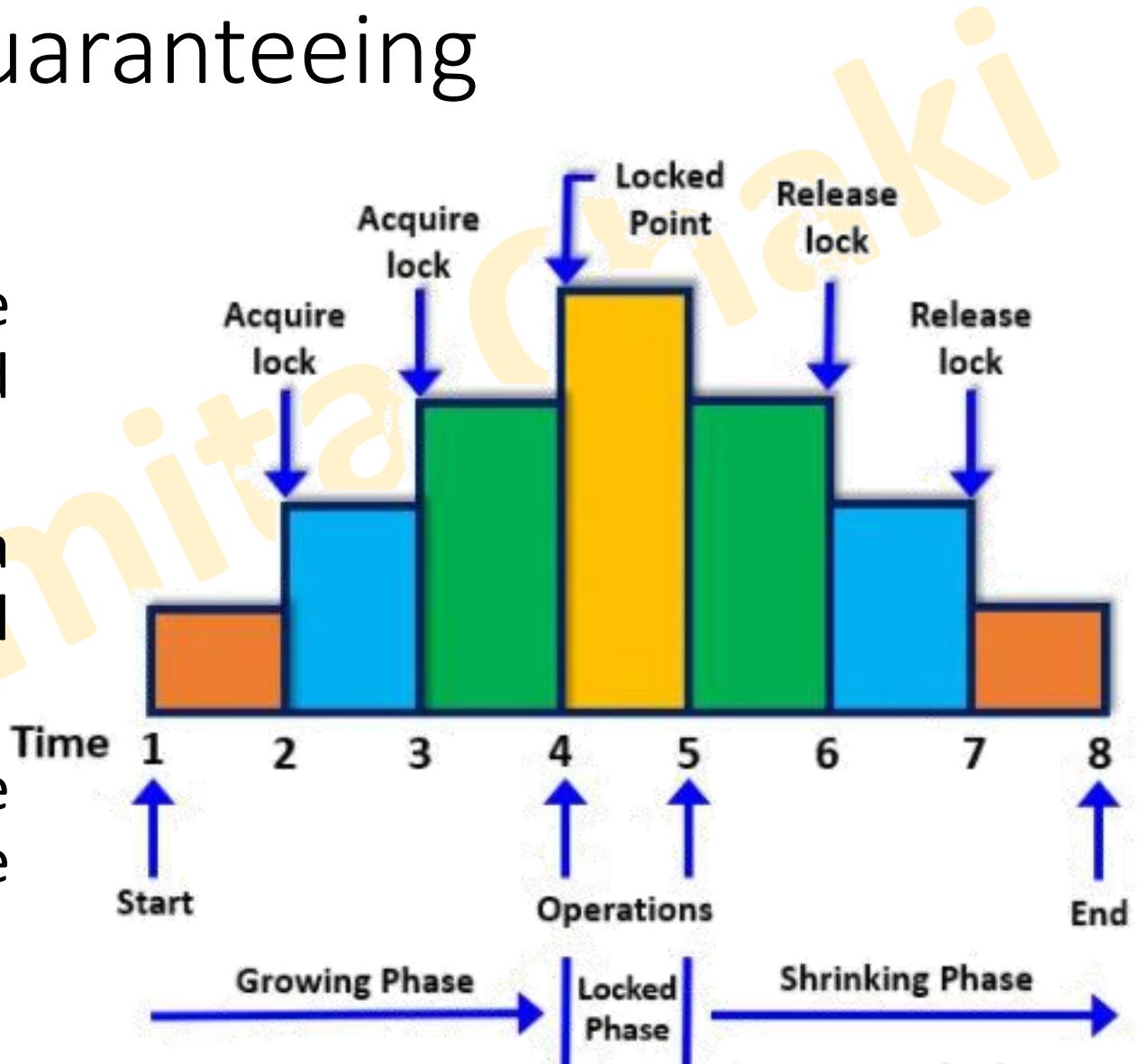
- It is desirable to relax conditions 4 and 5 in the preceding list in order to allow **lock conversion**; that is, a transaction that already holds a lock on item X is allowed under certain conditions to **convert** the lock from one locked state to another.
- For example, it is possible for a transaction T to issue a `read_lock(X)` and then later to **upgrade** the lock by issuing a `write_lock(X)` operation.
- If T is the only transaction holding a read lock on X at the time it issues the `write_lock(X)` operation, the lock can be upgraded; otherwise, the transaction must wait.
- It is also possible for a transaction T to issue a `write_lock(X)` and then later to **downgrade** the lock by issuing a `read_lock(X)` operation.

Two-Phase Locking: Guaranteeing Serializability

- Using binary locks or read/write locks in transactions, as described earlier, does not guarantee serializability of schedules on its own.
- A transaction is said to follow the **two-phase locking protocol** if all locking operations (read_lock, write_lock) precede the first unlock operation in the transaction.
- Such a transaction can be divided into two phases: an **expanding** or **growing (first) phase**, during which new locks on items can be acquired but none can be released; and a **shrinking (second) phase**, during which existing locks can be released but no new locks can be acquired.
- If lock conversion is allowed, then upgrading of locks (from read-locked to write-locked) must be done during the expanding phase, and downgrading of locks (from write-locked to read-locked) must be done in the shrinking phase.

Two-Phase Locking: Guaranteeing Serializability

- The above phases in a DBMS are determined by something called a '**Lock Point**'.
- Lock point** is the point where a transaction has achieved its final lock.
- It is also the point where the growing phase ends and the shrinking phase begins.



Two-Phase Locking: Guaranteeing Serializability

- Transactions T1 and T2 in top Figure do not follow the two-phase locking protocol because the `write_lock(X)` operation follows the `unlock(Y)` operation in T1, and similarly the `write_lock(Y)` operation follows the `unlock(X)` operation in T2. [Initial values: $X=20$, $Y=30$, *Result serial schedule T1 followed by T2: $X=50$, $Y=80$, Result of serial schedule T2 followed by T1: $X=70$, $Y=50$*]
- If we enforce two-phase locking, the transactions can be rewritten as T1' and T2', as shown in bottom Figure

T_1	T_2
<code>read_lock(Y);</code> <code>read_item(Y);</code> <code>unlock(Y);</code> <code>write_lock(X);</code> <code>read_item(X);</code> $X := X + Y;$ <code>write_item(X);</code> <code>unlock(X);</code>	<code>read_lock(X);</code> <code>read_item(X);</code> <code>unlock(X);</code> <code>write_lock(Y);</code> <code>read_item(Y);</code> $Y := X + Y;$ <code>write_item(Y);</code> <code>unlock(Y);</code>

T_1'	T_2'
<code>read_lock(Y);</code> <code>read_item(Y);</code> <code>write_lock(X);</code> <code>unlock(Y)</code> <code>read_item(X);</code> $X := X + Y;$ <code>write_item(X);</code> <code>unlock(X);</code>	<code>read_lock(X);</code> <code>read_item(X);</code> <code>write_lock(Y);</code> <code>unlock(X)</code> <code>read_item(Y);</code> $Y := X + Y;$ <code>write_item(Y);</code> <code>unlock(Y);</code>

Two-Phase Locking: Guaranteeing Serializability

- The schedule shown in this Figure is not permitted for T1' and T2' (with their modified order of locking and unlocking operations) under the rules of locking described because T1' will issue its `write_lock(X)` before it unlocks item Y; consequently, when T2' issues its `read_lock(X)`, it is forced to wait until T1' releases the lock by issuing an `unlock(X)` in the schedule.

T_1	T_2
<code>read_lock(Y);</code> <code>read_item(Y);</code> <code>unlock(Y);</code>	<code>read_lock(X);</code> <code>read_item(X);</code> <code>unlock(X);</code> <code>write_lock(Y);</code> <code>read_item(Y);</code> $Y := X + Y;$ <code>write_item(Y);</code> <code>unlock(Y);</code>
	<code>write_lock(X);</code> <code>read_item(X);</code> $X := X + Y;$ <code>write_item(X);</code> <code>unlock(X);</code>

Two-Phase Locking: Guaranteeing Serializability

- It can be proved that, if every transaction in a schedule follows the two-phase locking protocol, the schedule is guaranteed to be serializable, obviating the need to test for serializability of schedules.
- The locking protocol, by enforcing two-phase locking rules, also enforces serializability.
- Although the two-phase locking protocol guarantees serializability (that is, every schedule that is permitted is serializable), it does not permit all possible serializable schedules (that is, some serializable schedules will be prohibited by the protocol).

Two-Phase Locking: Solution of lost update problem



Time	T1	T2	A
t1		Write_lock(A)	100
t2	Write_lock(A)	Read_item(A)	100
t3	Wait	$A := A + 100$	100
t4	Wait	Write_item(A)	200
t5	Wait	Unlock(A)	200
t6	Wait	Commit	200
t7	Read_item(A)		200
t8	$A := A - 10$		200
t9	Write_item(A)		190
t10	Unlock(A)		190
t11	commit		190

Dr. J.

Two-Phase Locking: Solution of incorrect summary problem

Dr. JYC

Time	T1	T2	A = 100	B = 50	C = 25
t1	Write_lock(A)	1/ Sum=0	100	50	25
t2	Write_lock(C)		100	50	25
t3	Read_item(A)	Read_lock(A)	100	50	25
t4	A:= A-10	Wait	100	50	25
t5	Write_item(A)	Wait	90	50	25
t6	Read_item(C)	Wait	90	50	25
t7	C:= C+10	Wait	90	50	25
t8	Write_item(C)	Wait	90	50	35
t9	Unlock(A, C)	Wait	90	50	35
t10	Commit	Wait	90	50	35
t11		Read_lock(B)	90	50	35
t12		Read_lock(C)	90	50	35
t13		Read_item(A)	90	50	35
t14		Sum:= Sum+A	90	50	35
t15		Read_lock(B)	90	50	35
t16		Sum:= Sum+B	90	50	35
t17		Read_item(C)	90	50	35
t18		Sum:= Sum+C	90	50	35
t19		Unlock(A, B, C)	90	50	35
t20		Commit	90	50	35

Two-Phase Locking: Variations

- The technique just described is known as **basic 2PL**.
- A variation known as **conservative 2PL (or static 2PL)** requires a transaction to lock all the items it accesses before the transaction begins execution, by predeclaring its read-set and write-set.
- The most popular variation of 2PL is **strict 2PL**, which guarantees strict schedules. In this variation, a transaction T does not release any of its exclusive (write) locks until after it commits or aborts. Hence, no other transaction can read or write an item that is written by T unless T has committed, leading to a strict schedule for recoverability.
- A more restrictive variation of strict 2PL is **rigorous 2PL**, which also guarantees strict schedules. In this variation, a transaction T does not release any of its locks (exclusive or shared) until after it commits or aborts, and so it is easier to implement than strict 2PL.

Two-Phase Locking: Example

- **2 PL:** There is growing and shrinking phase.
- **Strict 2 PL:** There is Lock-x(B) and it is unlocked before commit so no strict 2 PL.
- **Rigorous:** If it is not strict 2 PL then it can't be Rigorous.
- **Conservative:** If it is not strict 2 PL then it can't be conservative.

T1
Lock-s(A)
Read(A)
Lock-x(B)
Read(B)
Unlock(A)
Write(B)
Unlock(B)

Two-Phase Locking: Example

- **2 PL:** There is growing and shrinking phase so it is 2 PL.
- **Strict 2 PL:** Exclusive locks are unlocked after commit.
So yes it is.
- **Rigorous:** We have taken Lock-s(A) and we have unlocked it before commit. So no rigorous.
- **Conservative:** We have not taken all the locks at first then start the transaction so no conservative.

T1
Lock-s(A)
Read(A)
Lock-x(B)
Unlock(A)
Read(B)
Write(B)
commit
Unlock(B)

Two-Phase Locking: Example

- **2 PL:** There is growing and shrinking phase so it is 2 PL.
- **Strict 2 PL:** Exclusive locks are unlocked after commit. So yes it is.
- **Rigorous:** We have unlocked all the locks after commit so it is rigorous.
- **Conservative:** We have not taken all the locks at first then start the transaction so no conservative.

T1
Lock-s(A)
Read(A)
Lock-x(B)
Read(B)
Write(B)
commit
Unlock(B)
Unlock(A)

Two-Phase Locking: Example

- **2 PL:** There is growing and shrinking phase so it is 2 PL.
- **Strict 2 PL:** Exclusive locks are unlocked after commit. So yes it is.
- **Rigorous:** We have unlocked all the locks after commit so it is rigorous.
- **Conservative:** We have taken all the locks at first then start the transaction so yes it is conservative.

T1
Lock-s(A)
Lock-x(B)
Read(B)
Write(B)
Read(A)
commit
Unlock(A)
Unlock(B)

Two-Phase Locking: Example

- **2 PL:** There is no growing and shrinking phase so it is not 2 PL.
- **Strict 2 PL:** Because it is not 2 PL so not either of it.
- **Rigorous:** Because it is not 2 PL so not either of it.
- **Conservative:** Because it is not 2 PL so not either of it.

T1
Lock-s(A)
Read(A)
Unlock(A)
Lock-x(B)
Read(B)
Write(B)
Unlock(B)
Unlock(A)
Commit

Two-Phase Locking: Variations

- The difference between strict and rigorous 2PL:
 - the former holds write-locks until it commits, whereas the latter holds all locks (read and write).
 - Also, the difference between conservative and rigorous 2PL is that the former must lock all its items before it starts, so once the transaction starts it is in its shrinking phase; the latter does not unlock any of its items until after it terminates (by committing or aborting), so the transaction is in its expanding phase until it ends.

Two-Phase Locking: Example

- Consider the following two transactions:
 - T31: read(A); read(B); if $A = 0$ then $B := B + 1$; write(B).
 - T32: read(B); read(A); if $B = 0$ then $A := A + 1$; write(A).
- Add lock and unlock instructions to transactions T31 and T32, so that they observe the two-phase locking protocol.
- Lock and unlock instructions:
 - T31: lock-S(A) read(A) lock-X(B) read(B) if $A = 0$ then $B := B + 1$ write(B)
unlock(A) unlock(B)
 - T32: lock-S(B) read(B) lock-X(A) read(A) if $B = 0$ then $A := A + 1$ write(A)
unlock(B) unlock(A)

Concurrency control based on timestamp

- The use of locking, combined with the 2PL protocol, guarantees serializability of schedules.
- The serializable schedules produced by 2PL have their equivalent serial schedules based on the order in which executing transactions lock the items they acquire.
- If a transaction needs an item that is already locked, it may be forced to wait until the item is released.
- Some transactions may be aborted and restarted because of the deadlock problem.
- A different approach to concurrency control involves using transaction timestamps to order transaction execution for an equivalent serial schedule.

Concurrency control based on timestamp: Timestamps

- **Timestamp** values are assigned in the order in which the transactions are submitted to the system, so a timestamp can be thought of as the transaction start time.
- We will refer to the timestamp of transaction T as **TS(T)**. Concurrency control techniques based on timestamp ordering do not use locks; hence, **deadlocks cannot occur**.
- Timestamps can be generated in several ways.
 - One possibility is to use a counter that is incremented each time its value is assigned to a transaction. The transaction timestamps are numbered 1, 2, 3, ... in this scheme.
 - A computer counter has a finite maximum value, so the system must periodically reset the counter to zero when no transactions are executing for some short period of time.
 - Another way to implement timestamps is to use the current date/time value of the system clock and ensure that no two timestamp values are generated during the same tick of the clock.

Concurrency control based on timestamp: Timestamp Ordering Algorithm

- The idea for this scheme is to enforce the equivalent serial order on the transactions based on their timestamps.
- A schedule in which the transactions participate is then serializable, and the only equivalent serial schedule permitted has the transactions in order of their timestamp values.
- This is called **timestamp ordering (TO)**.
- This differs from 2PL, where a schedule is serializable by being equivalent to some serial schedule allowed by the locking protocols.
- In timestamp ordering, however, the schedule is equivalent to the particular serial order corresponding to the order of the transaction timestamps.
- The algorithm allows interleaving of transaction operations, but it must ensure that for each pair of conflicting operations in the schedule, the order in which the item is accessed must follow the timestamp order.

Time	T1	T2	T3
1:00	Begin Transaction		
2:00		Begin Transaction	
3:00			Begin Transaction

$$TS(T1) = 10, TS(T2) = 20, TS (T3) = 30$$

Concurrency control based on timestamp: Timestamp Ordering Algorithm

- To do this, the algorithm associates with each database item X two timestamp (TS) values:
 1. **read_TS(X)**. The read timestamp of item X is the largest timestamp among all the timestamps of transactions that have successfully read item X—that is, $\text{read_TS}(X) = \text{TS}(T)$, where T is the youngest transaction that has read X successfully.

Time	T1	T2	T3
1:00	Begin Transaction		
2:00	R(X)		
2:15		Begin Transaction	
3:00		R(X)	
3:15			Begin Transaction
3:30			R (X)

$$\text{read_TS}(X) = 30$$

Concurrency control based on timestamp: Timestamp Ordering Algorithm

- To do this, the algorithm associates with each database item X two timestamp (TS) values:
 2. **write_TS(X)**. The write timestamp of item X is the largest of all the timestamps of transactions that have successfully written item X—that is, $\text{write_TS}(X) = \text{TS}(T)$, where T is the youngest transaction that has written X successfully. Based on the algorithm, T will also be the last transaction to write item X

Time	T1 (TS = 10)	T2 (TS = 20)	T3 (TS = 15)
1:00	Begin Transaction		
2:00	W(X)		
2:15			Begin Transaction
3:00			W (X)
3:15		Begin Transaction	
3:30		W(X)	

$$\text{write_TS}(X) = 20$$

Concurrency control based on timestamp: Basic Timestamp Ordering

- Whenever a transaction T issues a `write_item(X)` operation, the following check is performed:
 - If $TS(T_i) < \text{Read_TS}(X)$ then the operation is rejected, T_i abort.
 - If $TS(T_i) < \text{Write_TS}(X)$ then the operation is rejected and T_i is rolled back otherwise the operation is executed.
 - If the condition in part (a) and (b) does not occur, then execute the `write_item(X)` operation of T and set `write_TS(X)` to $TS(T)$.

Time	T1 (TS = 10)	T2 (TS = 20)
2:00		R(X)
2:15	W(X)	

a) T1 abort

Time	T1 (TS = 10)	T2 (TS = 20)
2:00		W(X)
2:15	W(X)	

b) T1 abort

Concurrency control based on timestamp: Basic Timestamp Ordering

- Whenever a transaction T issues a `read_item(X)` operation, the following check is performed:
 - a) If $\text{Write_TS}(X) > \text{TS}(T_i)$ then the operation is rejected.
 - b) If $\text{Write_TS}(X) \leq \text{TS}(T_i)$ then the operation is executed.

Time	T1 (TS = 10)	T2 (TS = 20)
2:00		W(X)
2:15	R(X)	

a) T1 abort

- Whenever the basic TO algorithm detects two conflicting operations that occur in the incorrect order, it rejects the later of the two operations by aborting the transaction that issued it.
- The schedules produced by basic TO are hence guaranteed to be conflict serializable.
- Deadlock does not occur with timestamp ordering.

Concurrency control based on timestamp:

Basic Timestamp Ordering

T1 (TS = 10)	T2 (TS = 20)	T3 (TS = 30)
R(X)		
	R(Y)	
W(Z)		
	R(Y)	
R(Z)		
	W(Y)	
		W(X)

W_{T2}(Y): T2 abort/roll back. After T3 again T2 will start with a new timestamp

W_{T3}(X)

	X	Y	Z
Read_TS	10	30	10
Write_TS	30	0	10

	X	Y	Z
Read_TS	0	0	0
Write_TS	0	0	0

	X	Y	Z
Read_TS	10	20	0
Write_TS	0	0	0

	X	Y	Z
Read_TS	10	30	0
Write_TS	0	0	10

R_{T1}(X)

	X	Y	Z
Read_TS	10	0	0
Write_TS	0	0	0

W_{T1}(Z)

	X	Y	Z
Read_TS	10	20	0
Write_TS	0	0	10

R_{T1}(Y)

	X	Y	Z
Read_TS	10	30	10
Write_TS	0	0	10

Concurrency control based on timestamp: Basic Timestamp Ordering

7. Assume basic timestamp ordering protocol and that time starts from 1; each operation takes unit amount of time and start of transaction T_i is denoted as S_i . The table of timestamp is given below:

Time	OP
1	S_1
2	$r_1(a)$
3	S_2
4	$r_2(b)$
5	$w_2(b)$
6	$w_1(a)$
7	S_3
8	$w_3(a)$
9	$w_3(b)$

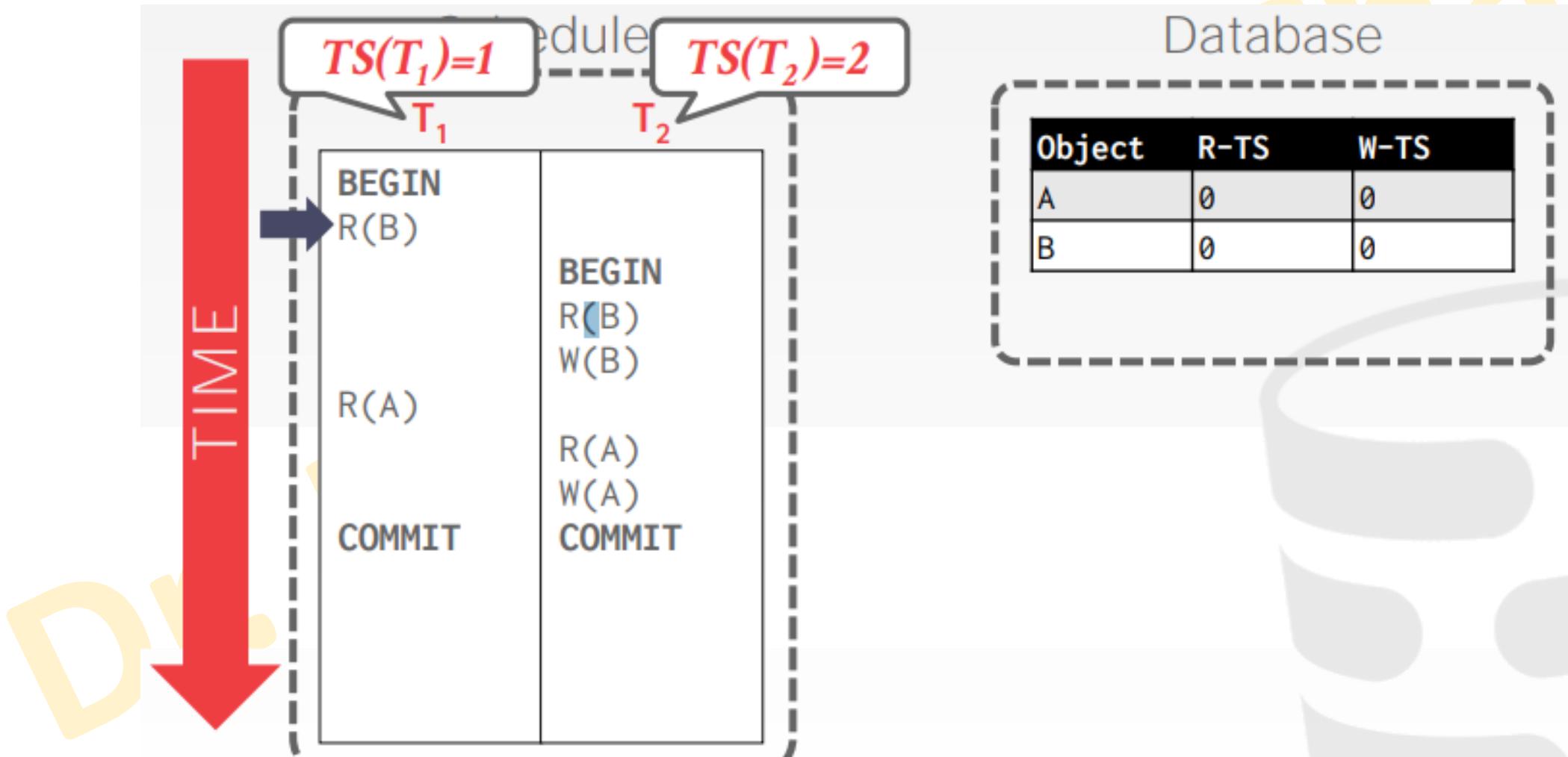
Find rts(a), wts(a), rts(b) and wts(b) at the end

(Marks: 2.00)

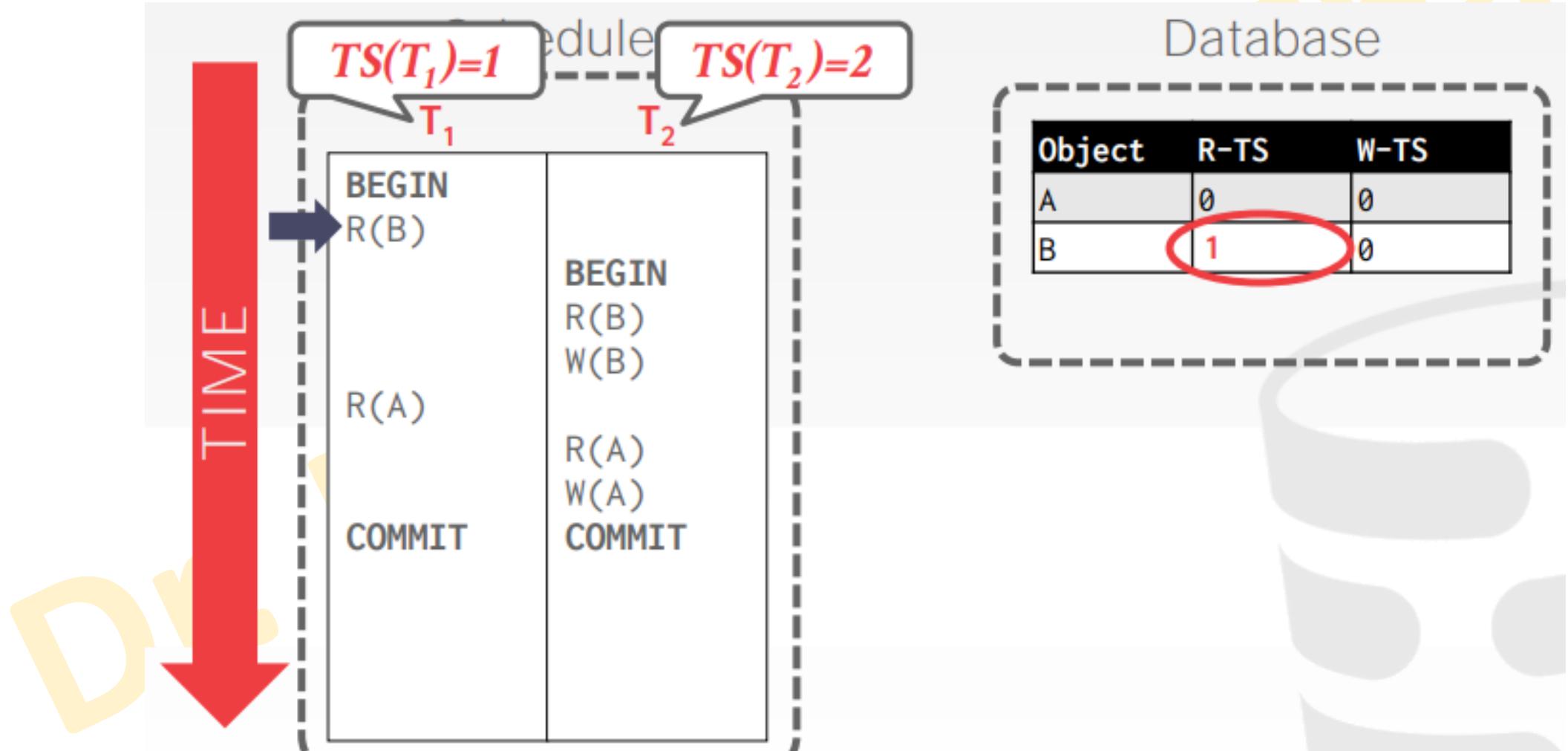
Concurrency control based on timestamp: Basic Timestamp Ordering

- T1 starts at TS =1
- T2 starts at TS = 3
- T3 starts at TS =7.
- While giving the TS for any read or write always look for youngest.
- **RTS(a)** = a is first read by T1 hence RTS(a) =1. (Read(a) is never done anywhere again hence it is youngest)
- **WTS(a)** = a is first written by T1 hence WTS(a) = 1. But again written by T3 which has higher TS (youngest) Hence final TS of WTS(a) = 7
- **RTS(b)** = b is first read by T2 hence RTS(b) =3. (Read(b) is never done anywhere again hence it is youngest)
- **WTS(b)** = b is first written by T2 hence WTS(b) = 3. But again written by T3 which has higher TS (youngest) Hence final TS of W(b) = 7

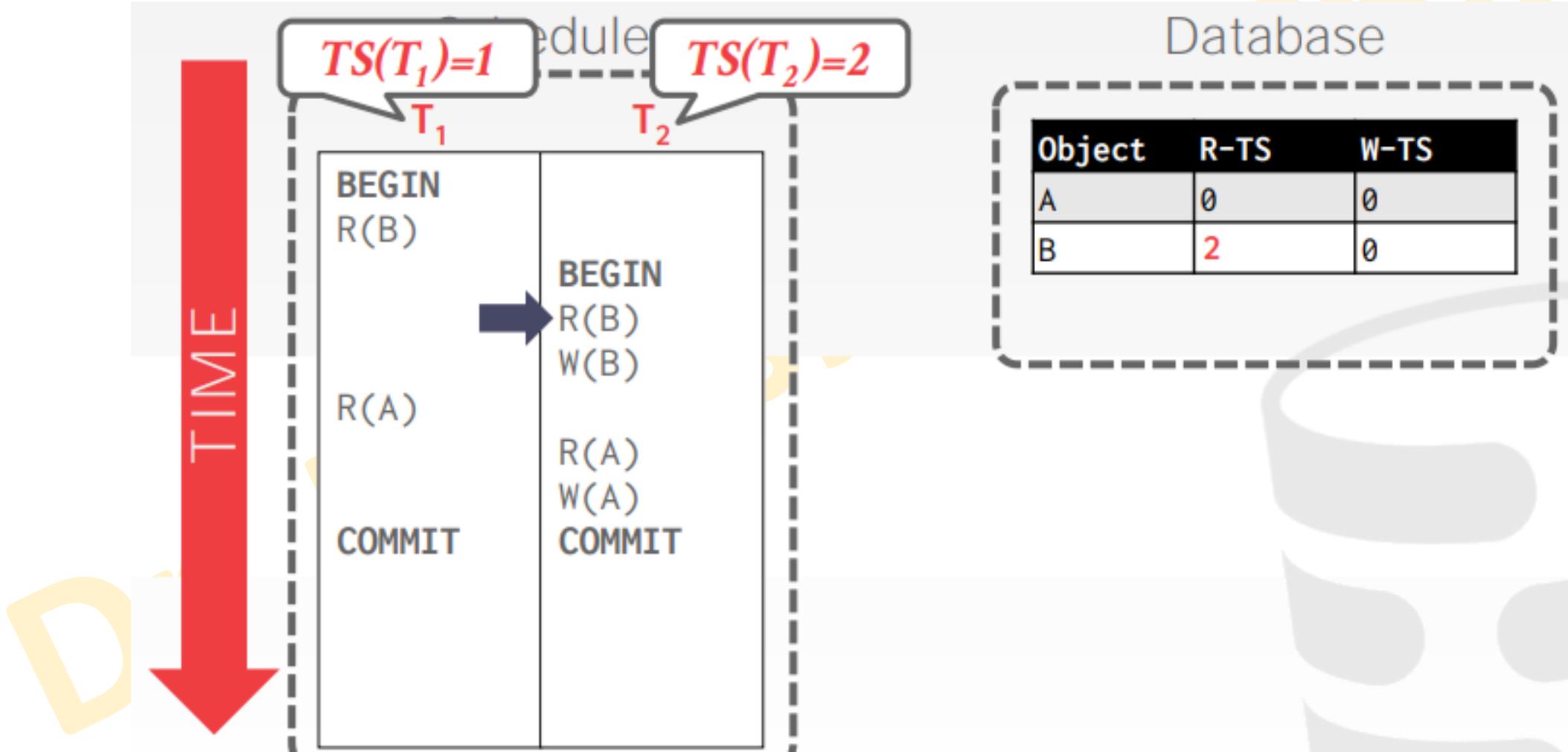
Concurrency control based on timestamp: Basic Timestamp Ordering



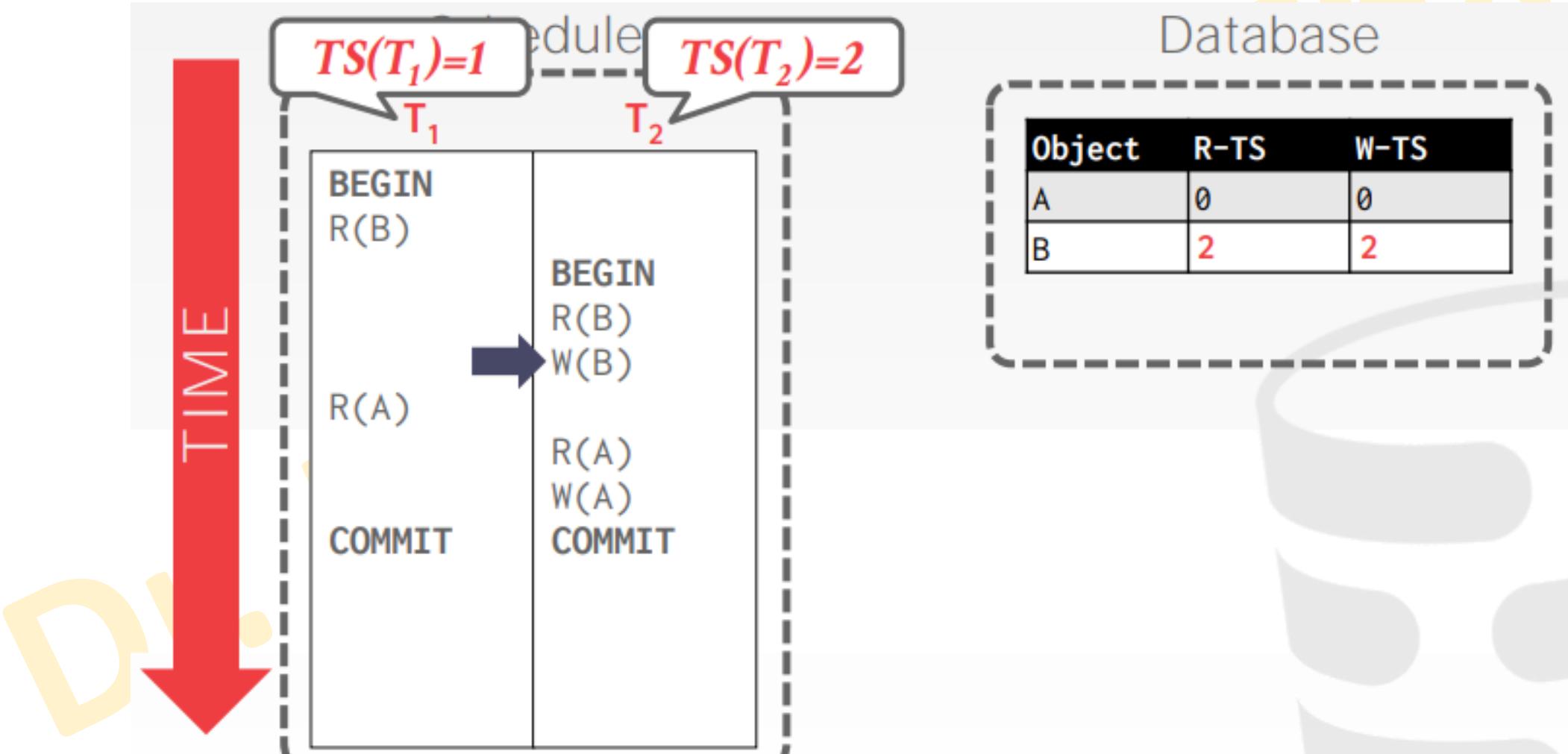
Concurrency control based on timestamp: Basic Timestamp Ordering



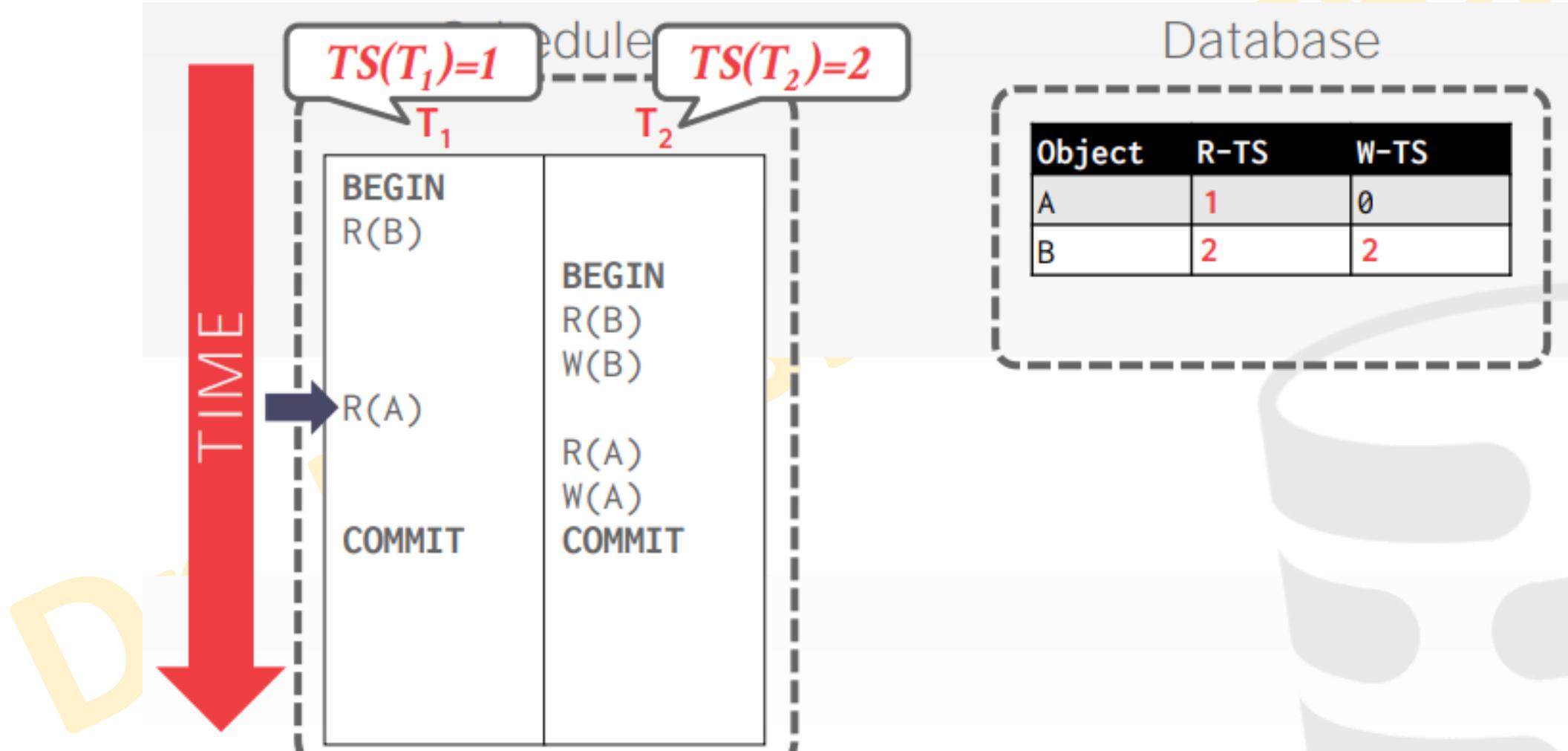
Concurrency control based on timestamp: Basic Timestamp Ordering



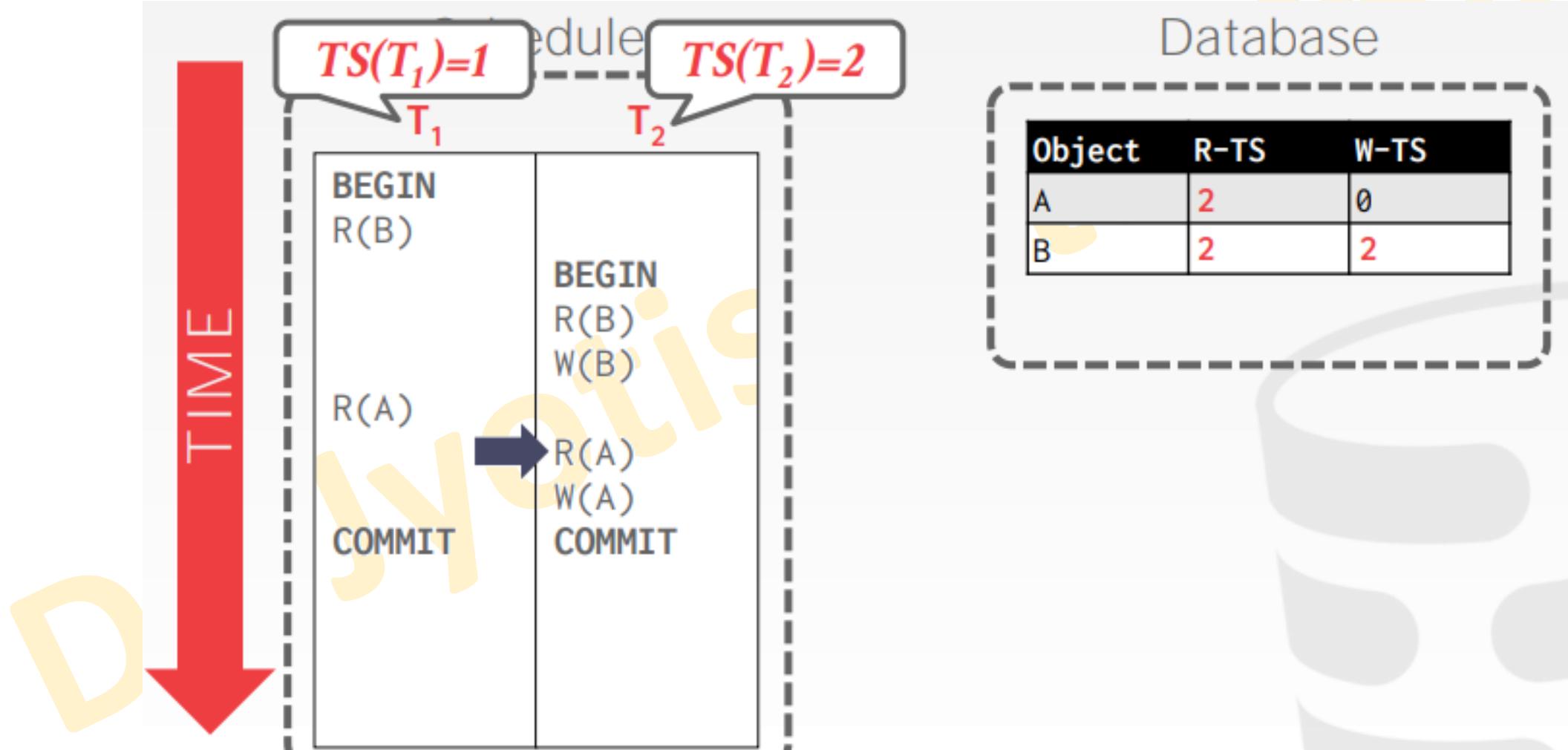
Concurrency control based on timestamp: Basic Timestamp Ordering



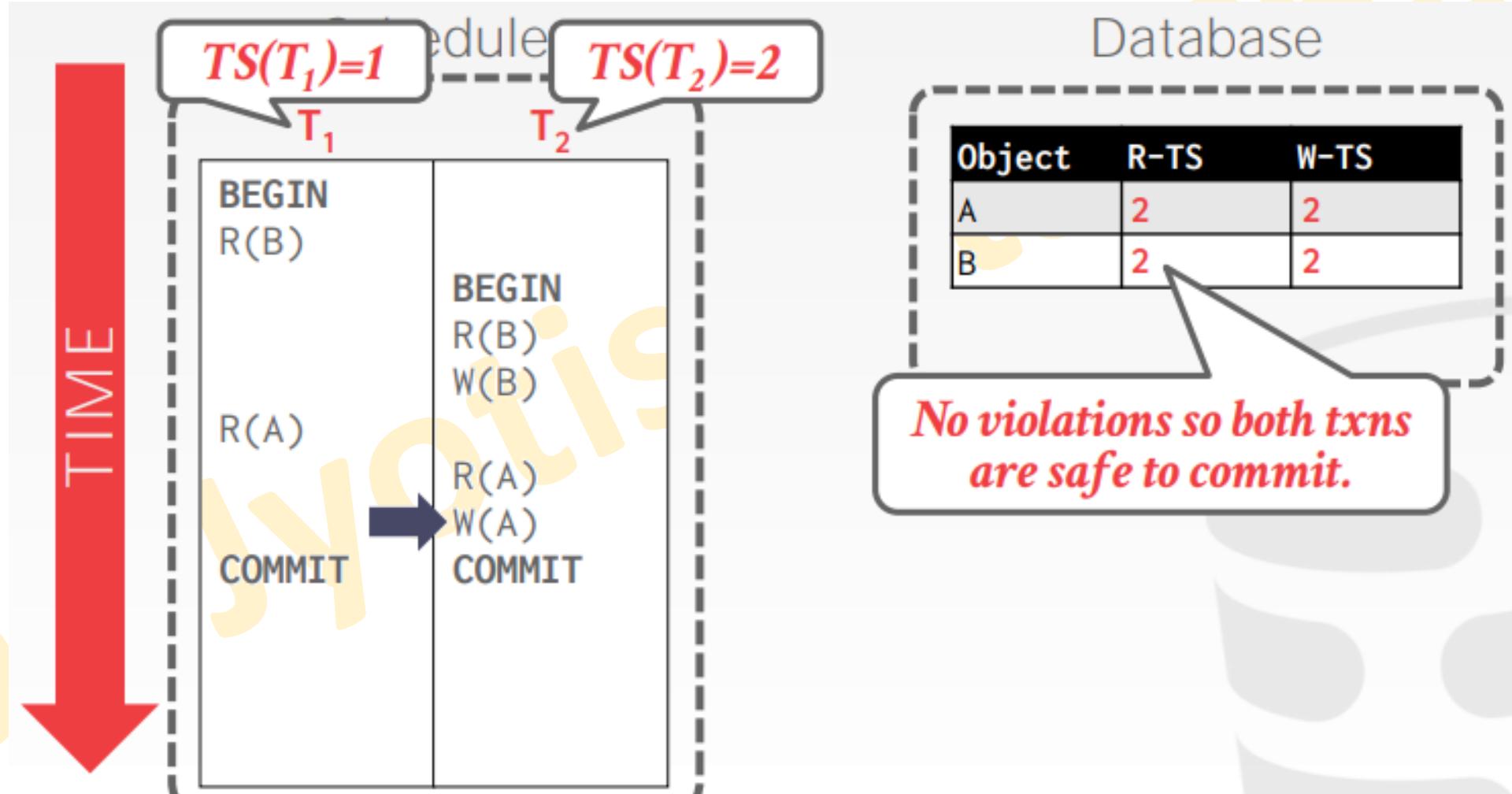
Concurrency control based on timestamp: Basic Timestamp Ordering



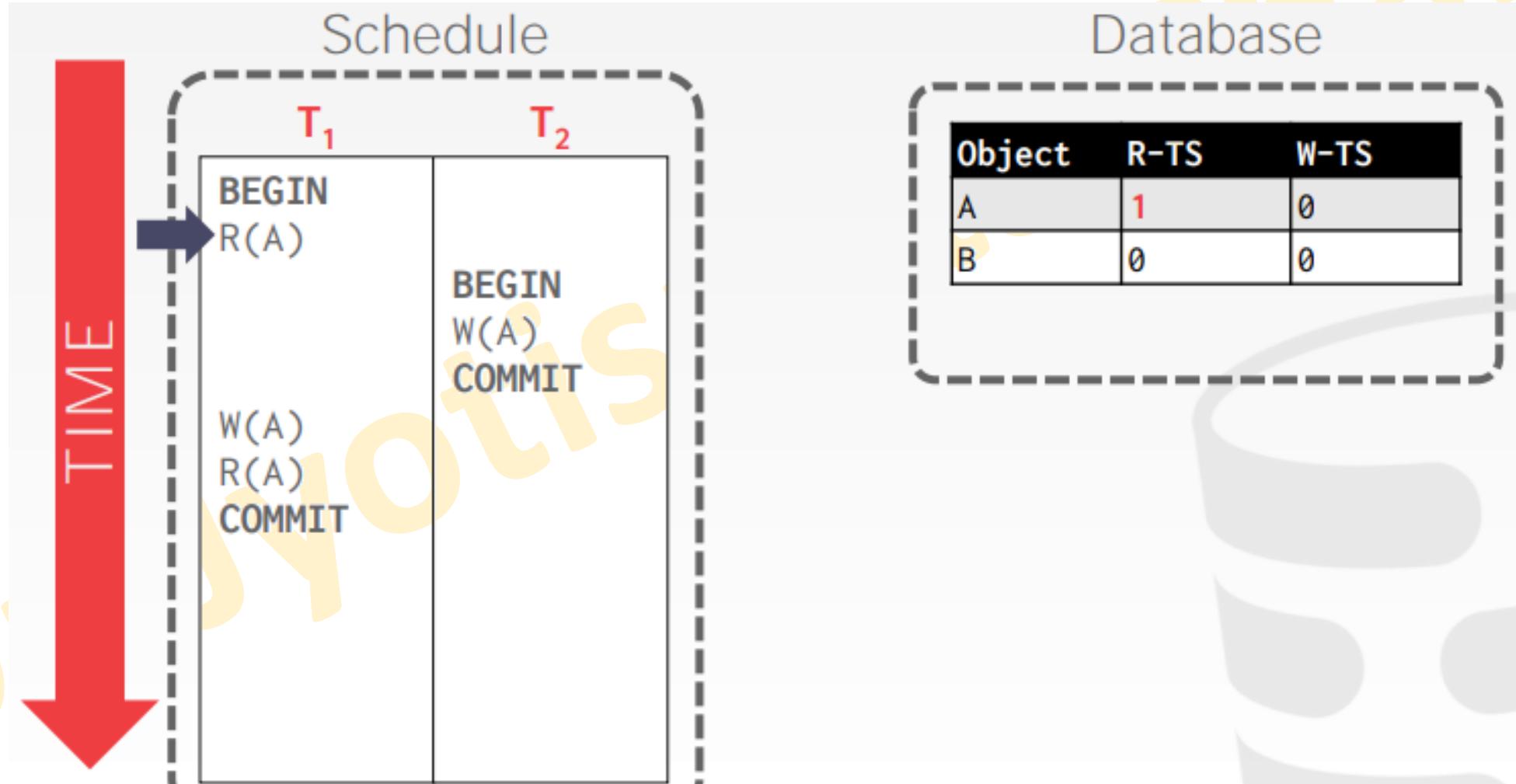
Concurrency control based on timestamp: Basic Timestamp Ordering



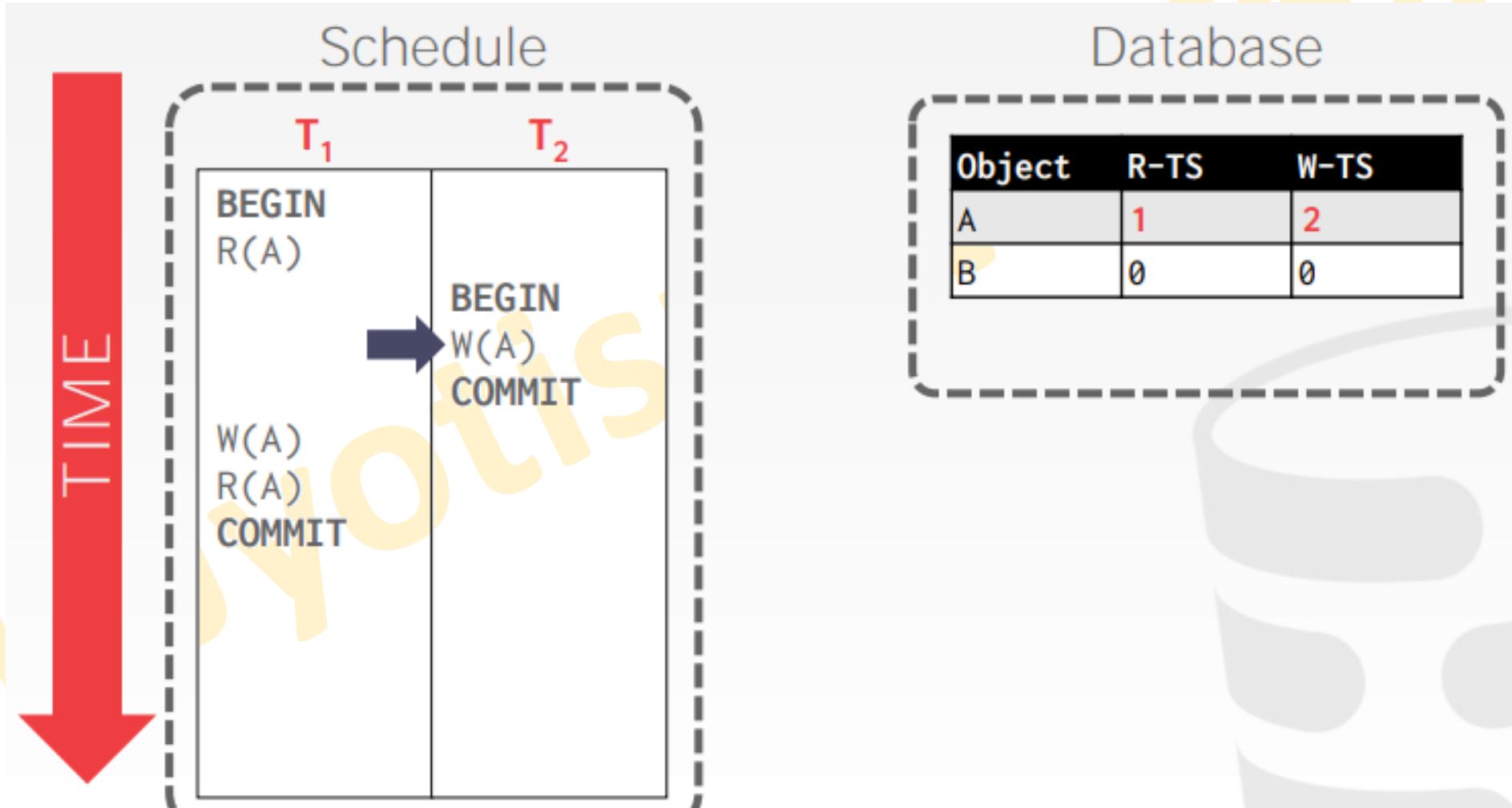
Concurrency control based on timestamp: Basic Timestamp Ordering



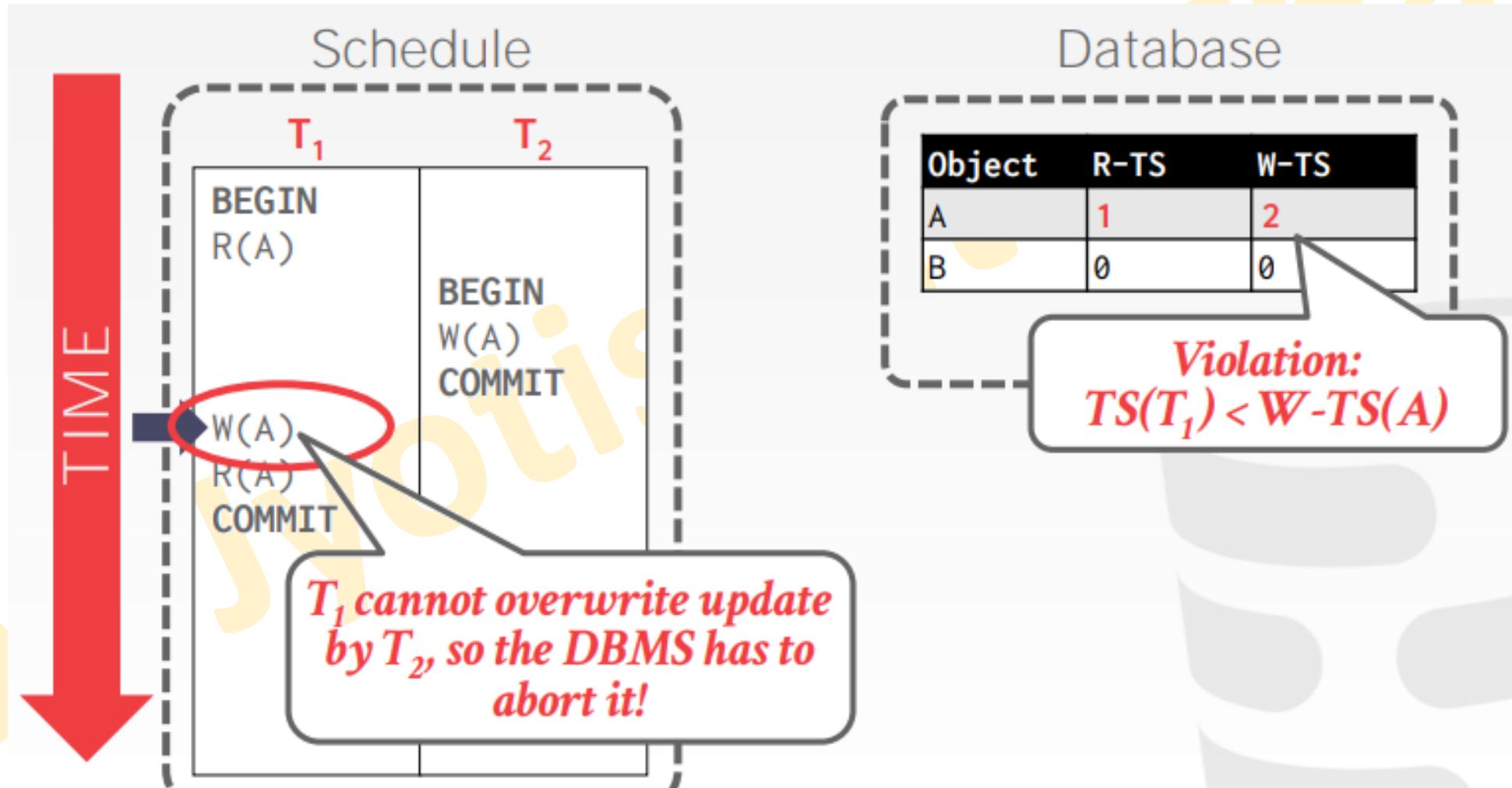
Concurrency control based on timestamp: Basic Timestamp Ordering



Concurrency control based on timestamp: Basic Timestamp Ordering



Concurrency control based on timestamp: Basic Timestamp Ordering



Concurrency control based on timestamp: Strict Timestamp Ordering

- A variation of basic TO called strict TO ensures that the schedules are both strict (for easy recoverability) and (conflict) serializable.
- In this variation, a transaction T issues a $\text{read_item}(X)$ or $\text{write_item}(X)$ such that $\text{TS}(T) > \text{write_TS}(X)$ has its read or write operation delayed until the transaction T' that wrote the value of X (hence $\text{TS}(T') = \text{write_TS}(X)$) has committed or aborted.
- To implement this algorithm, it is necessary to simulate the locking of an item X that has been written by transaction T' until T' is either committed or aborted.
- This algorithm does not cause deadlock, since T waits for T' only if $\text{TS}(T) > \text{TS}(T')$.

Concurrency control based on timestamp: Thomas's Write Rule

- A modification of the basic TO algorithm, known as Thomas's write rule, does not enforce conflict serializability, but it rejects fewer write operations by modifying the checks for the `write_item(X)` operation as follows:
 1. If $\text{read_TS}(X) > \text{TS}(T)$, then abort and roll back T and reject the operation.
 2. If $\text{write_TS}(X) > \text{TS}(T)$, then do not execute the write operation but continue processing. This is because some transaction with timestamp greater than $\text{TS}(T)$ —and hence after T in the timestamp ordering—has already written the value of X. Thus, we must ignore the `write_item(X)` operation of T because it is already outdated and obsolete. Notice that any conflict arising from this situation would be detected by case (1).
 3. If neither the condition in part (1) nor the condition in part (2) occurs, then execute the `write_item(X)` operation of T and set $\text{write_TS}(X)$ to $\text{TS}(T)$.

Concurrency control based on timestamp: Thomas's Write Rule

Time	T1	T2
t1	Begin transaction	
t2	Read (A)	
t3	A:= A+20	Begin transaction
t4		Read (A)
t5	Write (A)	
t6		A:= A+30
t7		Write (A)
t8		B:= 100
t9		Write (B)
t10		Commit
t11	Begin transaction	
t12	Read (A)	
t13	A:= A+20	
t14	Write (A)	
t15	Commit	

- T2 (t4) → dirty read problem.
- To get rid from the problem first T2 will execute, then T1 will restart (t11) [according to Thomas's write rule 1]

Concurrency control based on timestamp: Thomas's Write Rule

Time	T1	T2
t1	Begin transaction	
t2	Read (A)	
t3	A:= A+10	
t4	Write (A)	
t5		Begin transaction
t6		Read (B)
t7		B:= B+100
t8		Write (B)
t9		Read (C)
t10		C:= C+200
t11		Write (C)
t12		Commit
t13	C:= 50	
t14	Write (C)	
t15	Commit	

- T2 is younger than T1
- T2 first updates the value of C.
- Thus according to Thomas's write rule (2) the updated value of C done by T1 (at t13) will be cancelled out

Recovery concepts

- Recovery from transaction failures usually means that the database is restored to the most recent consistent state before the time of failure.
- To do this, the system must keep information about the changes that were applied to data items by the various transactions.
- This information is typically kept in the **system log**.
- A typical strategy for recovery may be summarized informally as follows:
 - If there is extensive damage to a wide portion of the database due to catastrophic failure, such as a disk crash, the recovery method restores a past copy of the database that was backed up to archival storage (typically tape or other large capacity offline storage media) and reconstructs a more current state by reapplying or redoing the operations of committed transactions from the backed-up log, up to the time of failure.
 - When the database on disk is not physically damaged, and a noncatastrophic failure has occurred, the recovery strategy is to identify any changes that may cause an inconsistency in the database.

Recovery concepts

- For recovery from any type of failure data values prior to modification (BFIM - BeFore Image) and the new value after modification (AFIM – AFter Image) are required.
- These values and other information is stored in a sequential file called system log.

T ID	Operation	Data item	BFIM	AFIM
T1	Begin			
T1	Write	X	X = 100	X = 200
T2	Begin			
T1	W	Y	Y = 50	Y = 100
T1	R	M	M = 200	M = 200
T3	R	N	N = 400	N = 400
T1	End			

Recovery concepts

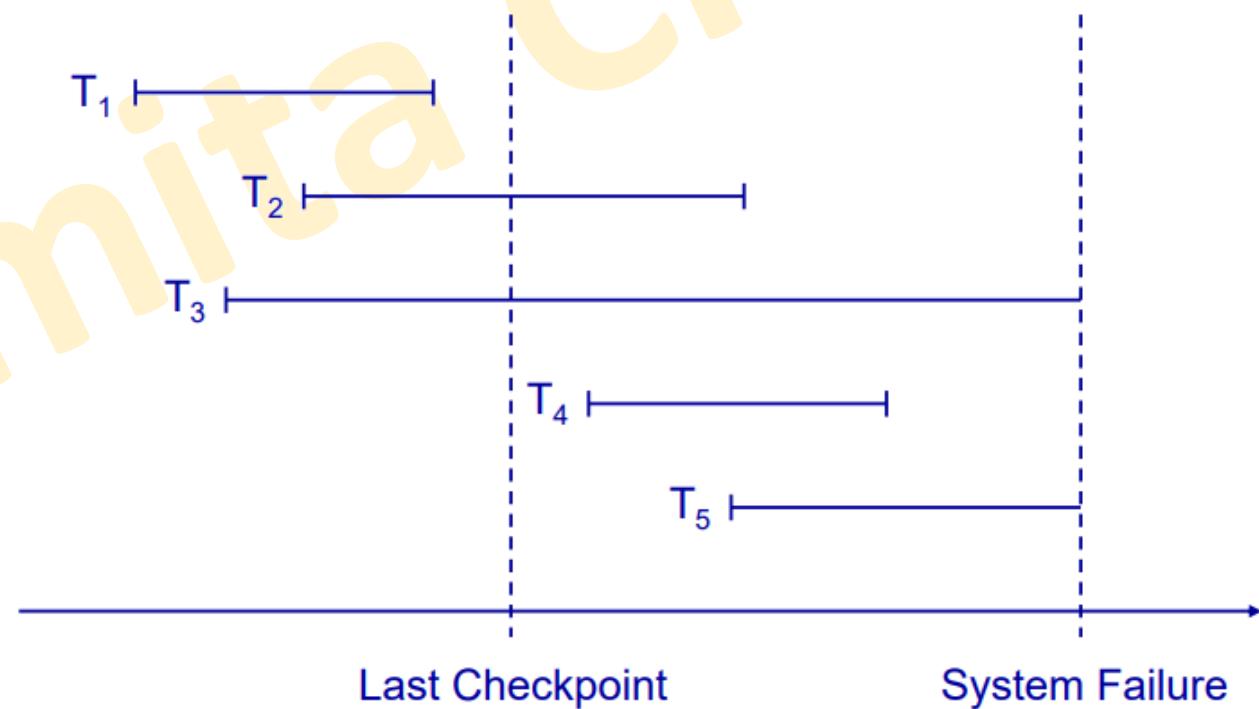
- To maintain atomicity, a transaction's operations are redone or undone.
 - **Undo:** Restore all BFIMs on to disk (Remove all AFIMs).
 - **Redo:** Restore all AFIMs on to disk.
- Database recovery is achieved either by performing only **Undos** or only **Redos** or by a combination of the two.

Recovery concepts: Checkpointing

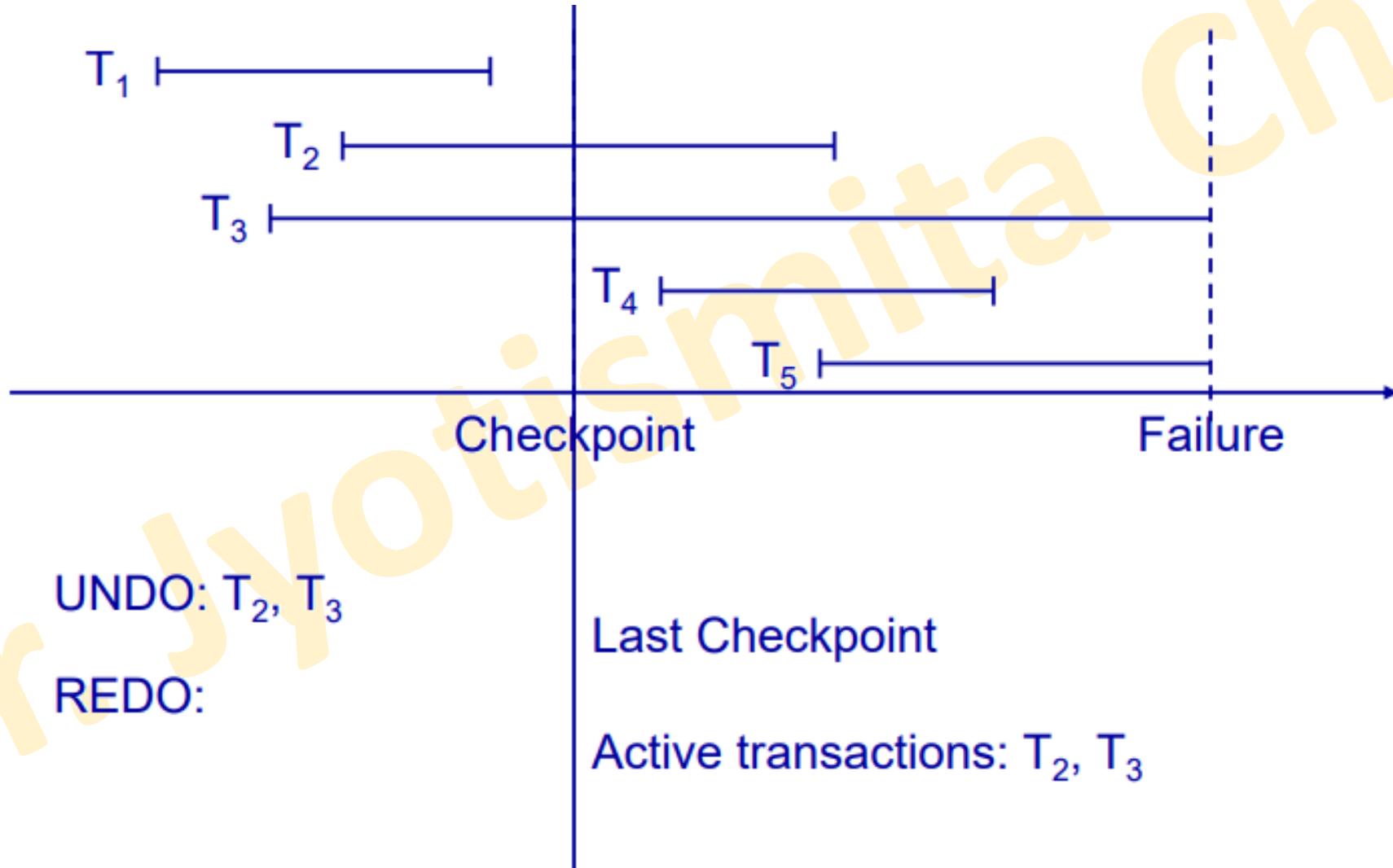
- From time to time the database flushes its buffer to database disk to minimize the task of recovery.
- Following steps defines a **checkpoint** operation:
 1. Suspend execution of transactions temporarily.
 2. Force write modified buffer data to disk (unless a no-UNDO).
 3. Write a [checkpoint] record to the log, save the log to disk.
 4. Resume normal transaction execution.
- During recovery redo or undo is required to transactions appearing after [checkpoint] record.

Recovery concepts: Checkpointing

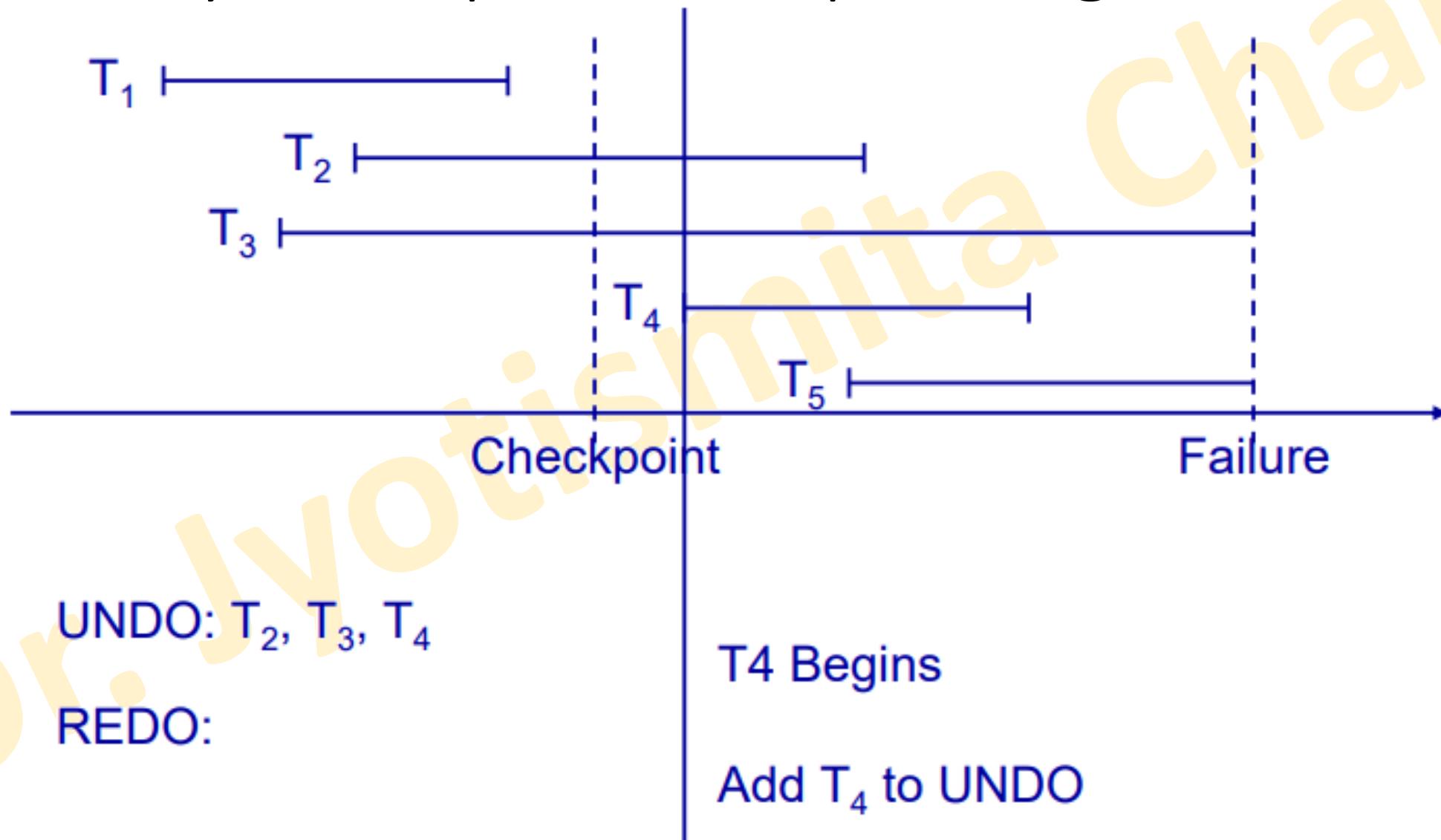
- Any transaction that was running at the time of failure needs to be undone and restarted.
- Any transactions that committed since the last checkpoint need to be redone.
- Transactions of type T1 need no recovery.
- Transactions of type T3 or T5 need to be undone and restarted.
- Transactions of type T2 or T4 need to be redone.



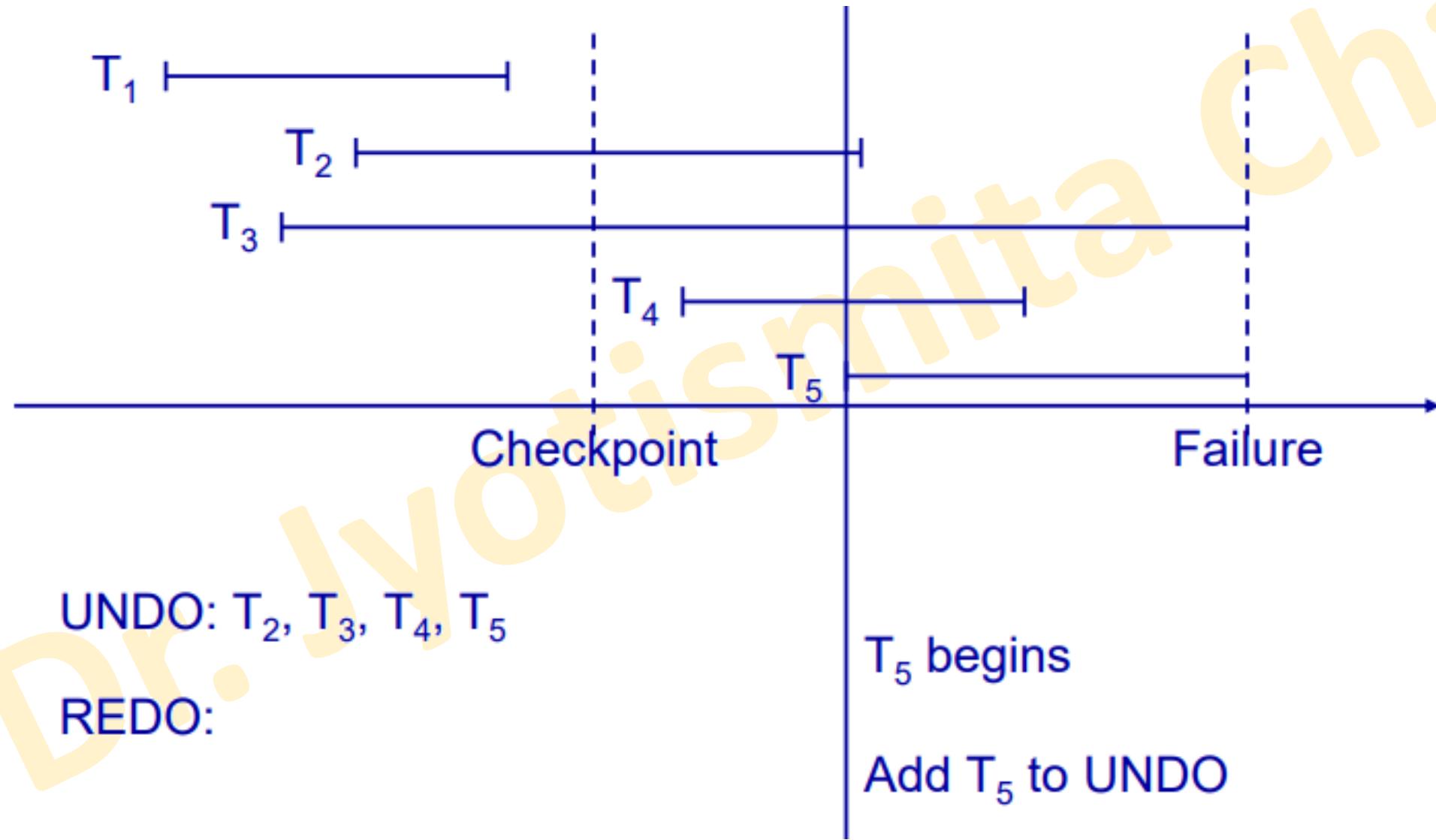
Recovery concepts: Checkpointing



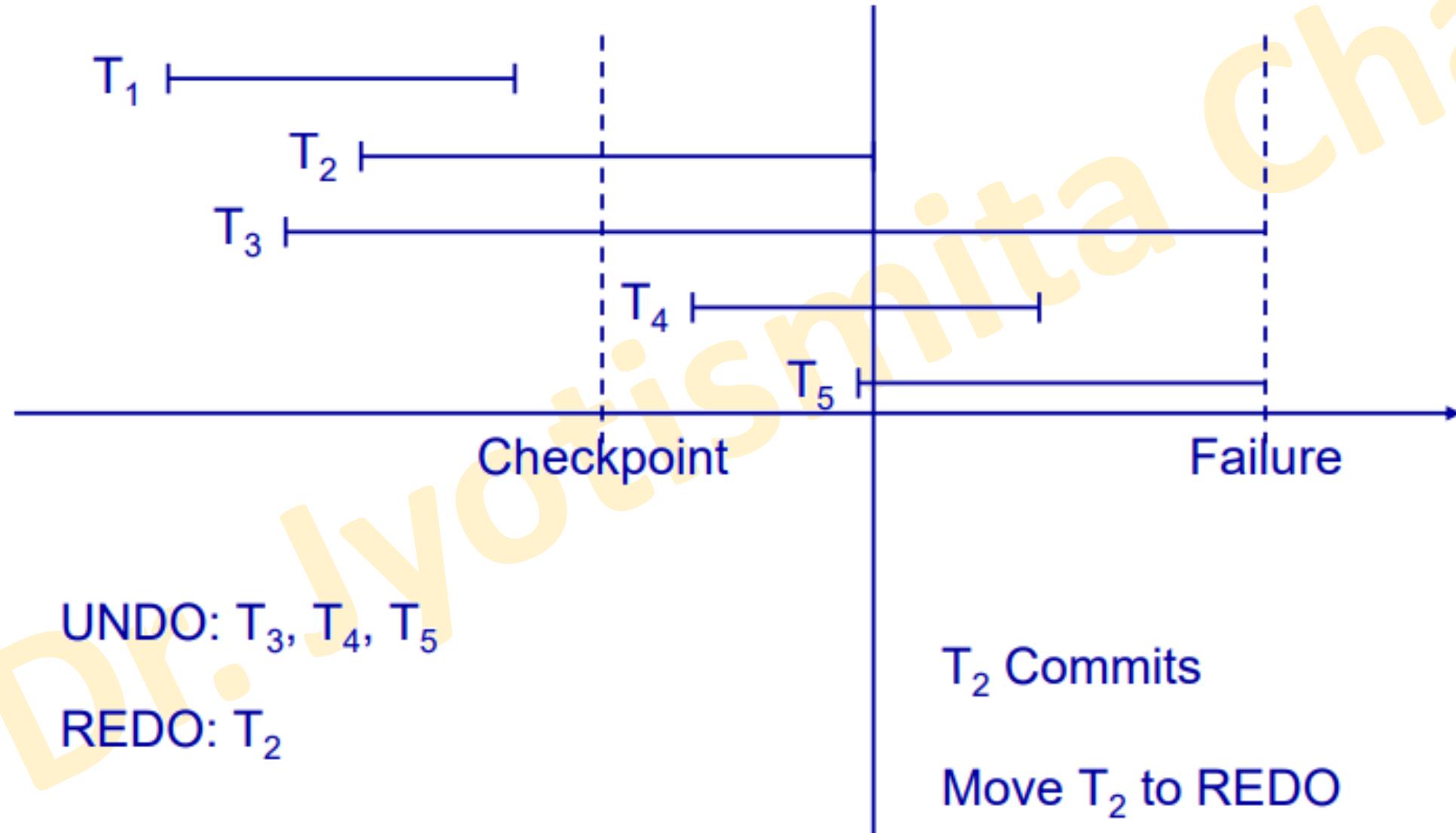
Recovery concepts: Checkpointing



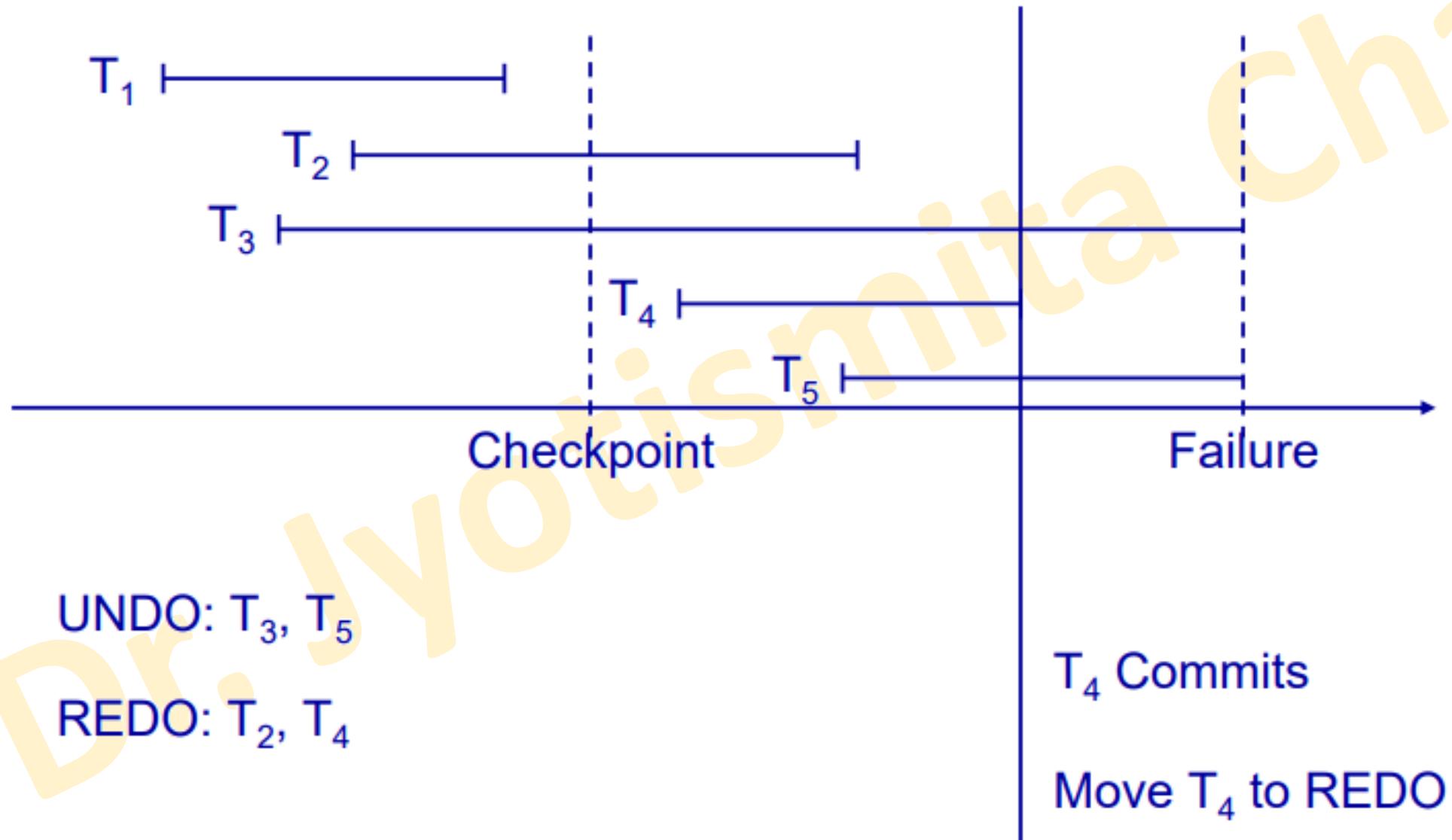
Recovery concepts: Checkpointing



Recovery concepts: Checkpointing



Recovery concepts: Checkpointing



Recovery based on deferred update

- Conceptually, we can distinguish two main policies for recovery from noncatastrophic transaction failures: **deferred update** and **immediate update**.
- The **deferred update** techniques do not physically update the database on disk until after a transaction commits; then the updates are recorded in the database.
- Before reaching commit, all transaction updates are recorded in the local transaction workspace or in the main memory buffers that the DBMS maintains.
- Before commit, the updates are recorded persistently in the log file on disk, and then after commit, the updates are written to the database from the main memory buffers.
- If a transaction fails before reaching its commit point, it will not have changed the database on disk in any way, so UNDO is not needed.
- It may be necessary to REDO the effect of the operations of a committed transaction from the log, because their effect may not yet have been recorded in the database on disk.
- Hence, deferred update is also known as the **NO-UNDO/REDO** algorithm.

Recovery based on deferred update

- Two tables are required for implementing this protocol.
 - **Active table:** All active transactions are entered in this table.
 - **Commit table:** Transactions to be committed are entered in this table.
- These tables are filled by scanning through the log from the last checkpoint during recovery.
- During recovery
 - All transactions of the **commit** table are redone in the order they were written to log, and
 - All transactions of **active** tables are ignored.

Recovery based on deferred update: Single Transaction

- $A = 100, B = 200$

T1
Read (A)
$A := A + 100$
Write (A)
Read (B)
$B := B + 200$
Write (B)
Commit

System Log
<start_transaction, T1>
<write_item, T1, A, 200>
<write_item, T1, B, 400>
<Commit, T1>

New value

T1
Read (A)
$A := A + 100$
Write (A)
Read (B)
$B := B + 200$
Write (B)

System Log
<start_transaction, T1>
<write_item, T1, A, 200>
<write_item, T1, B, 400>

- If T1 fails before Commit \rightarrow no REDO
- After recovery of T1: Value of A = 100, B = 200

- After commit, the data stored in HD will be updated

- After commit, if T1 fails \rightarrow REDO
- After Recovery of T1: Updated value A = 200, B = 400

Recovery based on deferred update: Multiple Transaction

T_1
read_item(A)
read_item(D)
write_item(D)

T_4
read_item(B)
write_item(B)
read_item(A)
write_item(A)

T_2
read_item(B)
write_item(B)
read_item(D)
write_item(D)

T_3
read_item(A)
write_item(A)
read_item(C)
write_item(C)

The READ and WRITE operations of four transactions.

[start_transaction, T_1]
[write_item, $T_1, D, 20$]
[commit, T_1]
[checkpoint]
[start_transaction, T_4]
[write_item, $T_4, B, 15$]
[write_item, $T_4, A, 20$]
[commit, T_4]
[start_transaction, T_2]
[write_item, $T_2, B, 12$]
[start_transaction, T_3]
[write_item, $T_3, A, 30$]
[write_item, $T_2, D, 25$]

System log at the point of crash

← System crash



T_2 and T_3 are ignored because they did not reach their commit points. T_4 is redone because its commit point is after the last system checkpoint.

Recovery techniques based on immediate update

- In the **immediate update** techniques, the database may be updated by some operations of a transaction before the transaction reaches its commit point.
- However, these operations must also be recorded in the log on disk by force-writing before they are applied to the database on disk, making recovery still possible.
- If a transaction fails after recording some changes in the database on disk but before reaching its commit point, the effect of its operations on the database must be undone; that is, the transaction must be rolled back.
- In the general case of immediate update, both undo and redo may be required during recovery. This technique, known as the **UNDO/REDO** algorithm, requires both operations during recovery and is used most often in practice.
- A variation of the algorithm where all updates are required to be recorded in the database on disk before a transaction commits requires undo only, so it is known as the **UNDO/NO-REDO** algorithm.

Recovery techniques based on immediate update

- Theoretically, we can distinguish two main categories of immediate update algorithms:
 1. If the recovery technique ensures that all updates of a transaction are recorded in the database on disk before the transaction commits, there is never a need to REDO any operations of committed transactions. This is called the UNDO/NO-REDO recovery algorithm. In this method, all updates by a transaction must be recorded on disk before the transaction commits, so that REDO is never needed.
 2. If the transaction is allowed to commit before all its changes are written to the database, we have the most general case, known as the UNDO/REDO recovery algorithm. This is also the most complex technique, but the most commonly used in practice.

Recovery techniques based on immediate update

- $A = 100, B = 200$

T1
Read (A)
$A := A + 100$
Write (A)
Read (B)
$B := B + 200$
Write (B)
Commit

System Log
<start_transaction, T1>
<write_item, T1, A, 100, 200>
<write_item, T1, B, 200, 400>
<Commit, T1>

New value

Old value

Before commit (immediate after writing to main memory), the data stored in HD will be updated

T1
Read (A)
$A := A + 100$
Write (A)
Read (B)
$B := B + 200$
Write (B)

System Log
<start_transaction, T1>
<write_item, T1, A, 100, 200>
<write_item, T1, B, 200, 400>

- After commit, if T1 fails → REDO
- After Recovery of T1: Updated value $A = 200, B = 400$

- If T1 fails before Commit → UNDO
- After recovery of T1 (old value is fetched from the log): Value of $A = 100, B = 200$

Shadow paging

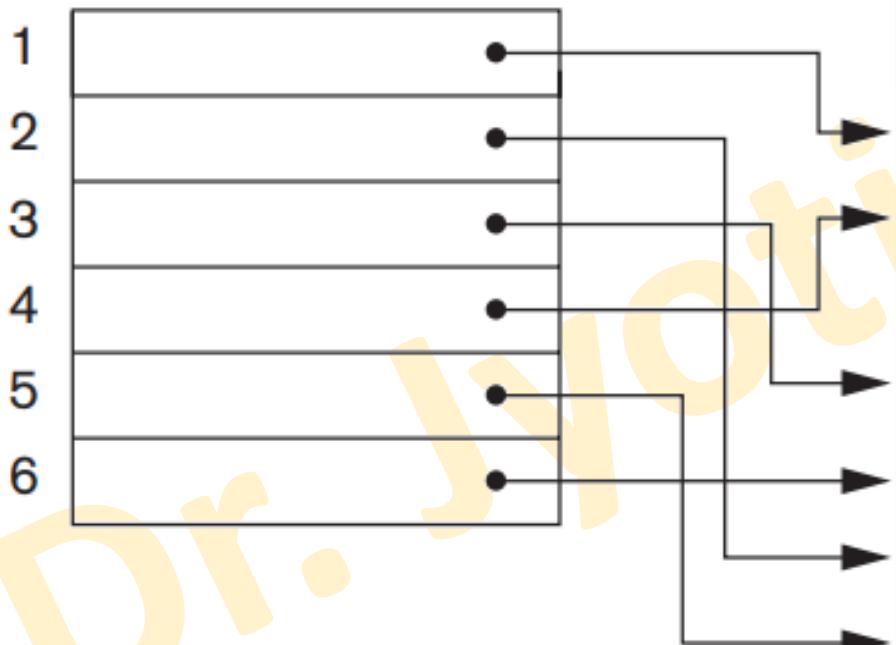
- This recovery scheme does not require the use of a log in a single-user environment.
- In a multiuser environment, a log may be needed for the concurrency control method.
- Shadow paging considers the database to be made up of a number of fixedsize disk pages (or disk blocks)—say, n —for recovery purposes.
- A directory with n entries is constructed, where the i -th entry points to the i -th database page on disk.
- The directory is kept in main memory if it is not too large, and all references—reads or writes—to database pages on disk go through it.
- When a transaction begins executing, the current directory—whose entries point to the most recent or current database pages on disk—is copied into a shadow directory.
- The shadow directory is then saved on disk while the current directory is used by the transaction.

Shadow paging

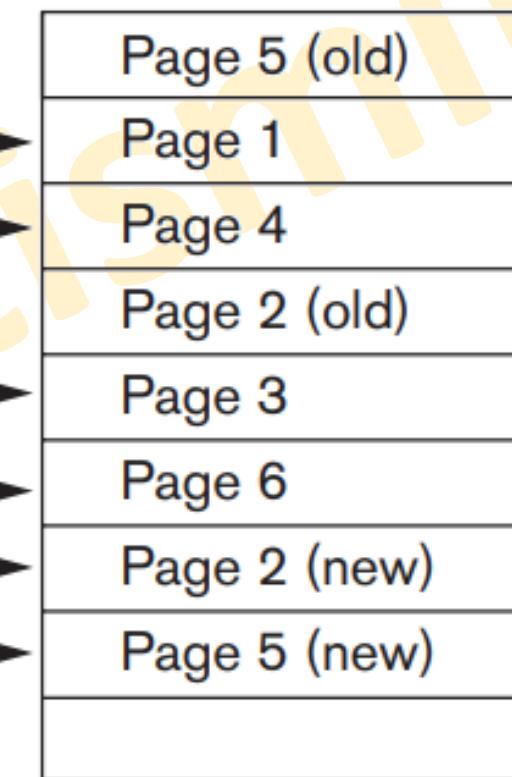
- During transaction execution, the shadow directory is never modified.
- When a `write_item` operation is performed, a new copy of the modified database page is created, but the old copy of that page is not overwritten.
- Instead, the new page is written elsewhere—on some previously unused disk block.
- The current directory entry is modified to point to the new disk block, whereas the shadow directory is not modified and continues to point to the old unmodified disk block.
- For pages updated by the transaction, two versions are kept.
- The old version is referenced by the shadow directory and the new version by the current directory.

Shadow paging

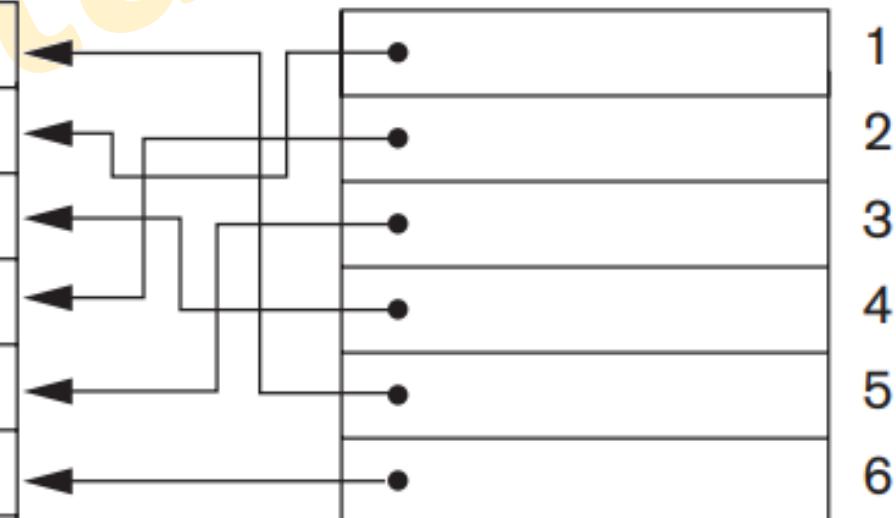
Current directory
(after updating
pages 2, 5)



Database disk
blocks (pages)



Shadow directory
(not updated)



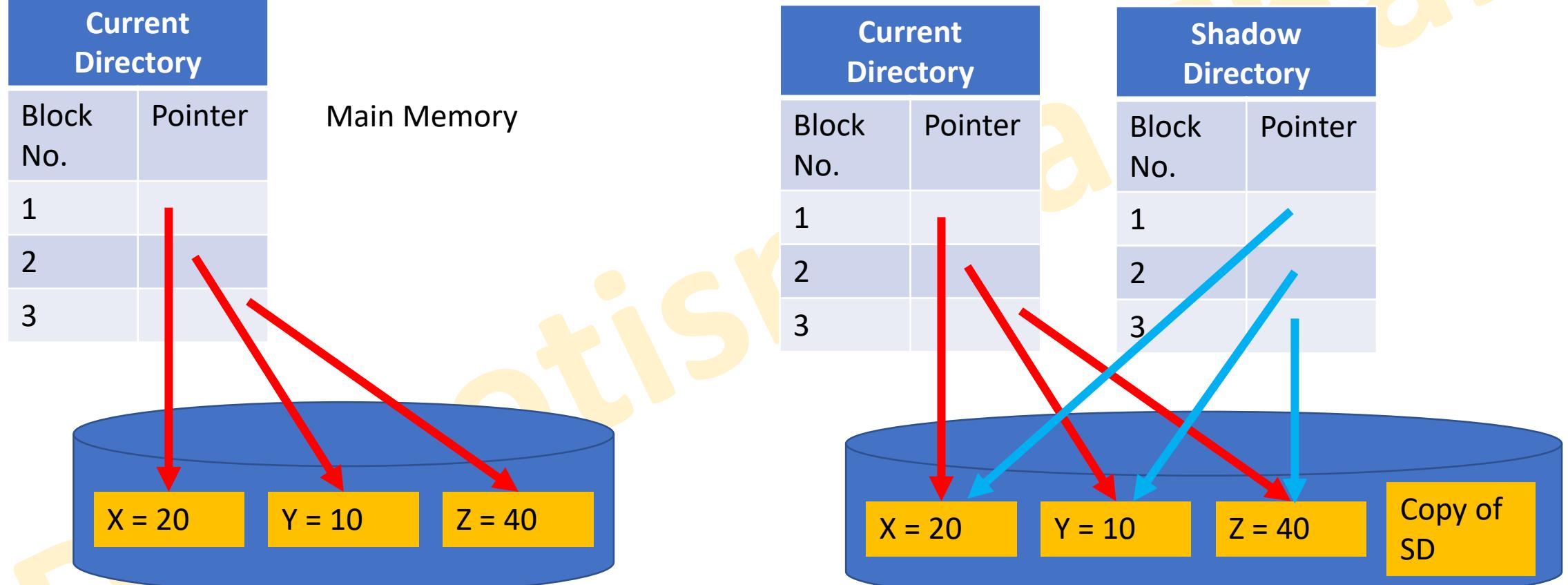
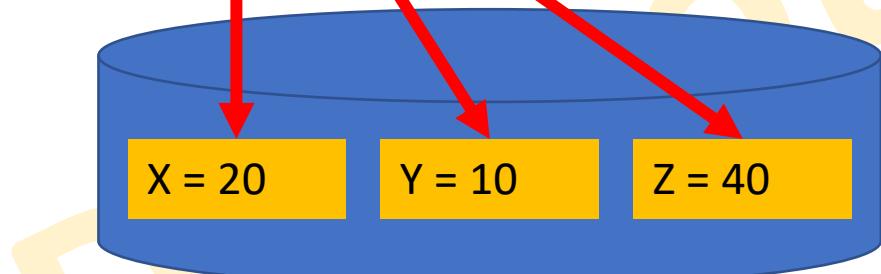
Shadow paging

- To recover from a failure during transaction execution, it is sufficient to free the modified database pages and to discard the current directory.
- The state of the database before transaction execution is available through the shadow directory, and that state is recovered by reinstating the shadow directory.
- The database thus is returned to its state prior to the transaction that was executing when the crash occurred, and any modified pages are discarded.
- Committing a transaction corresponds to discarding the previous shadow directory.
- Since recovery involves neither undoing nor redoing data items, this technique can be categorized as a NO-UNDO/NO-REDO technique for recovery.

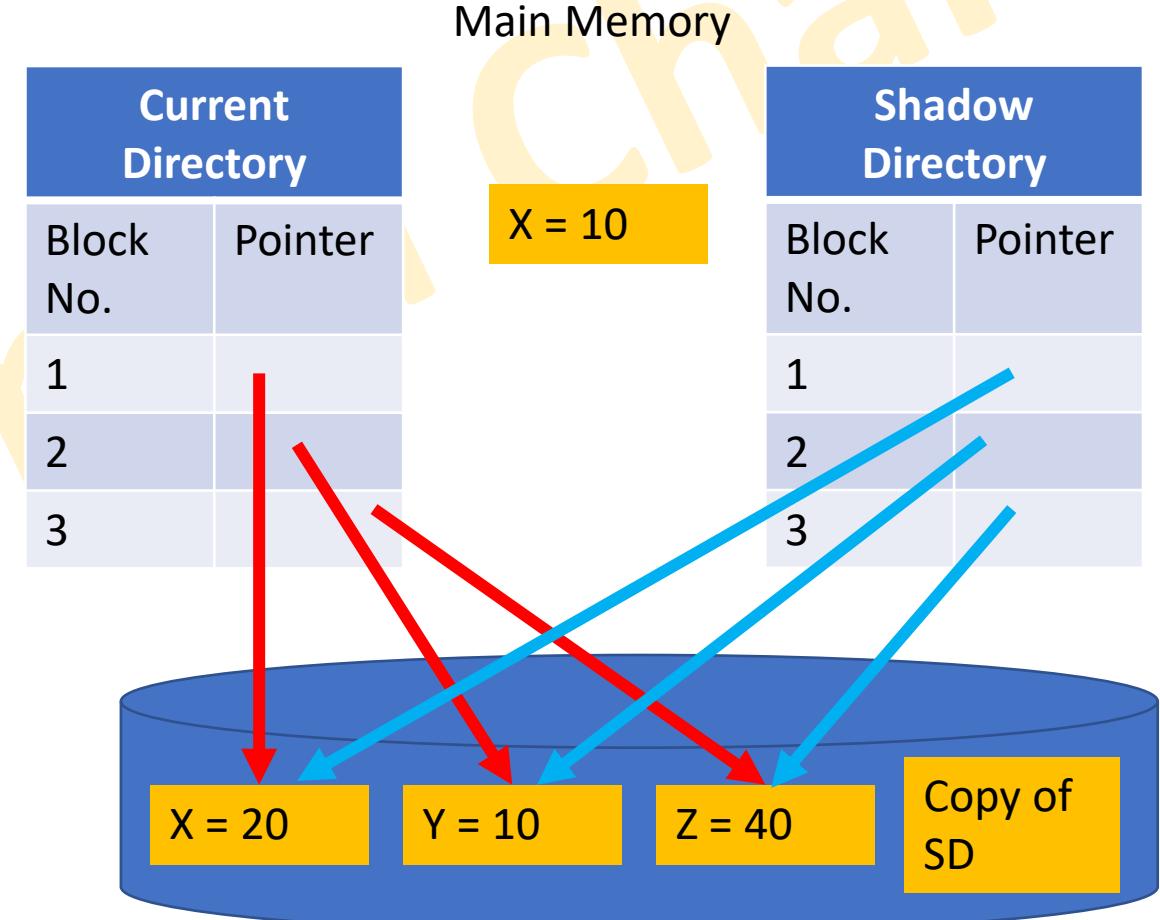
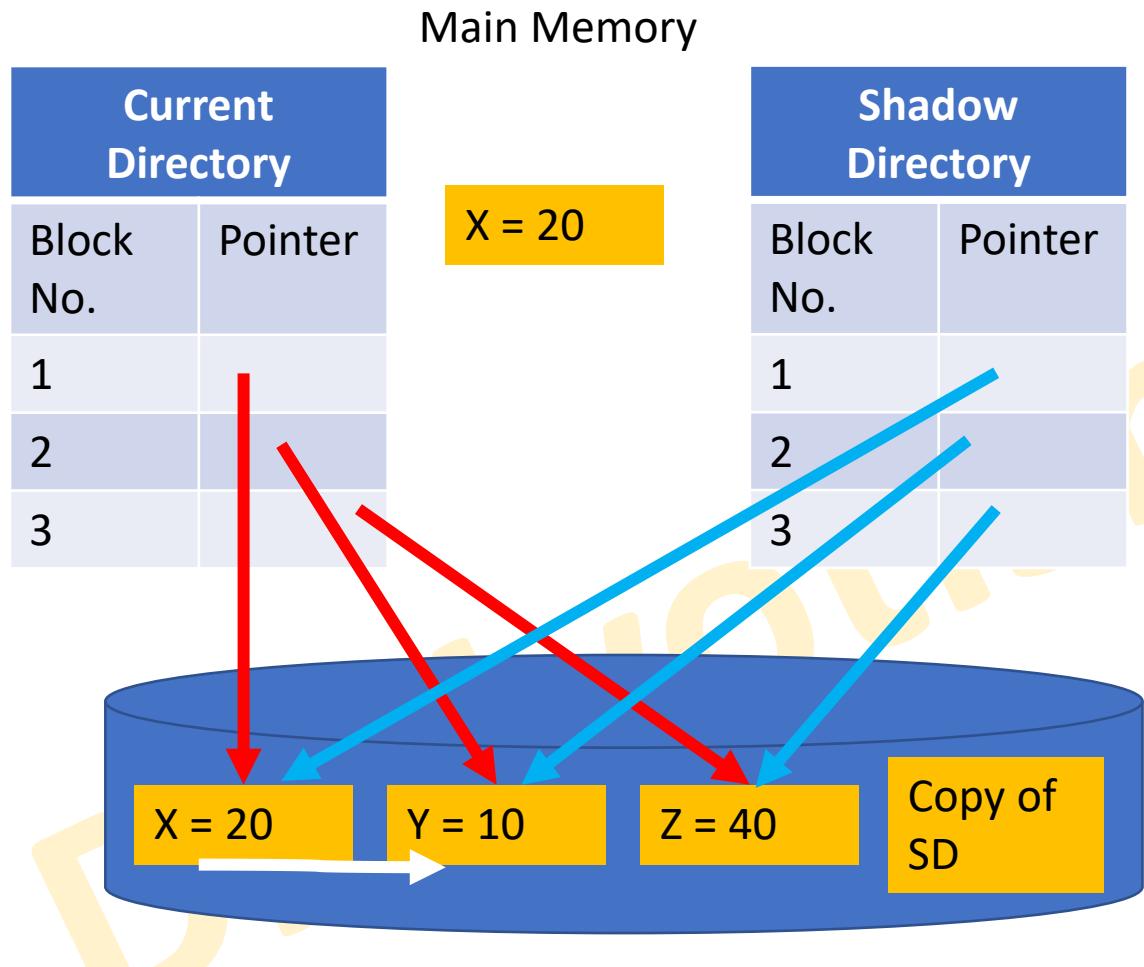
Shadow paging: Example

Current Directory	
Block No.	Pointer
1	
2	
3	

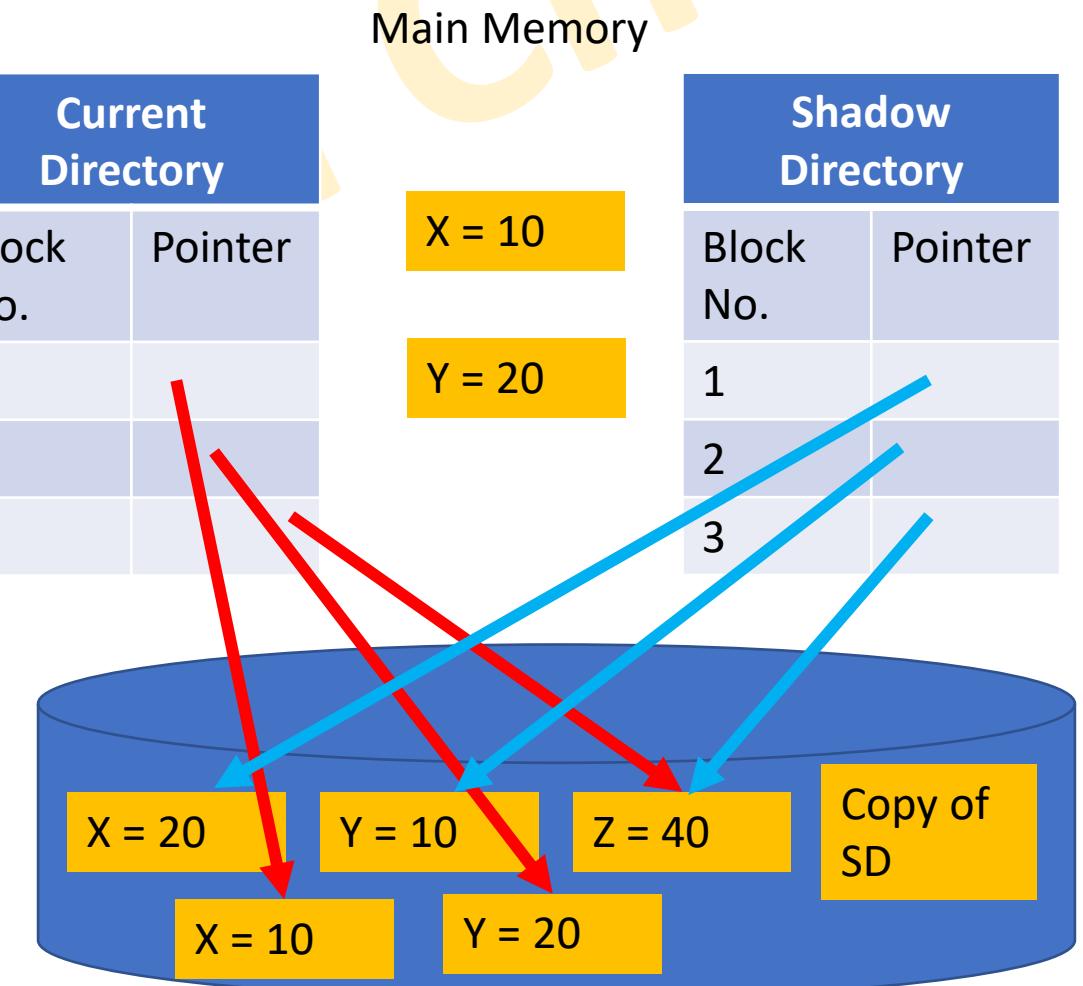
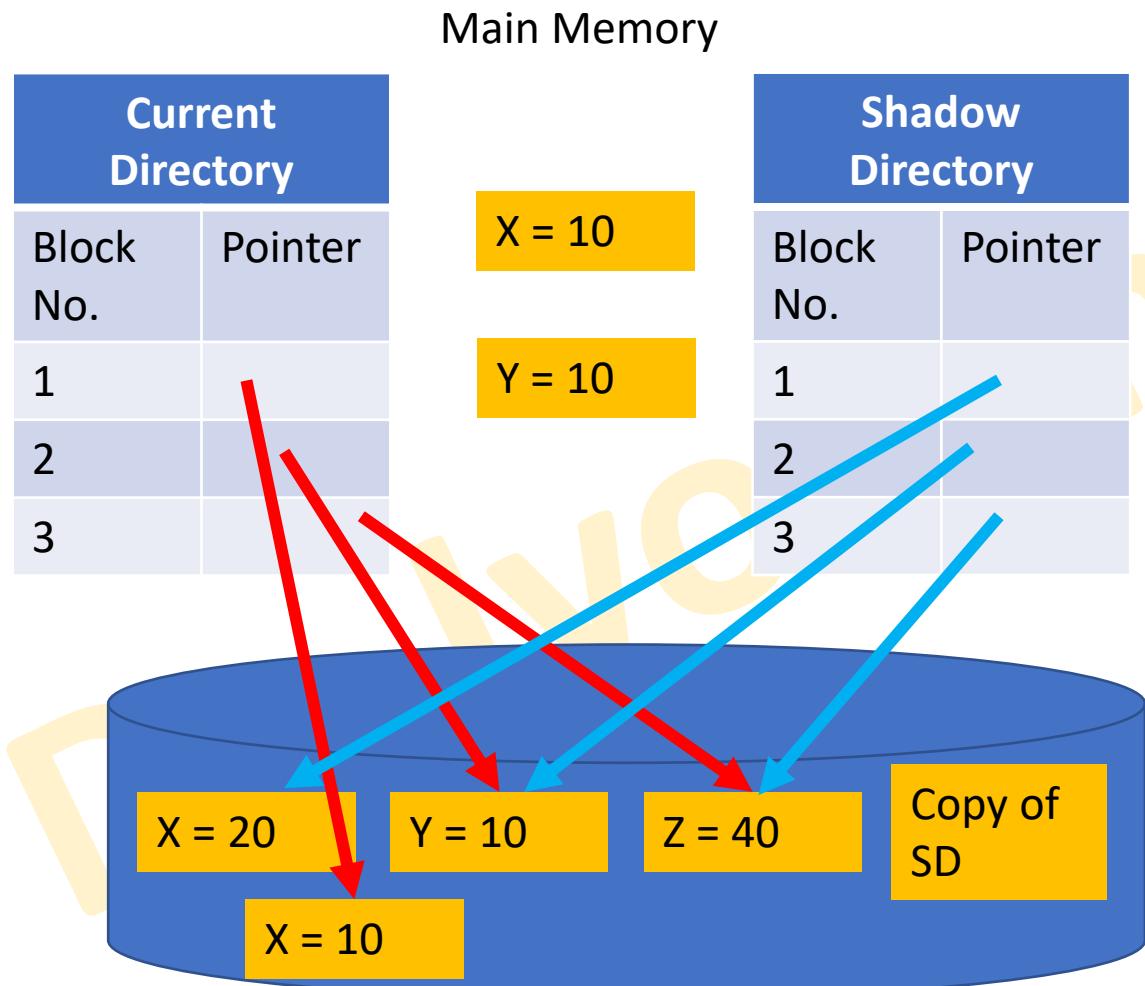
Main Memory



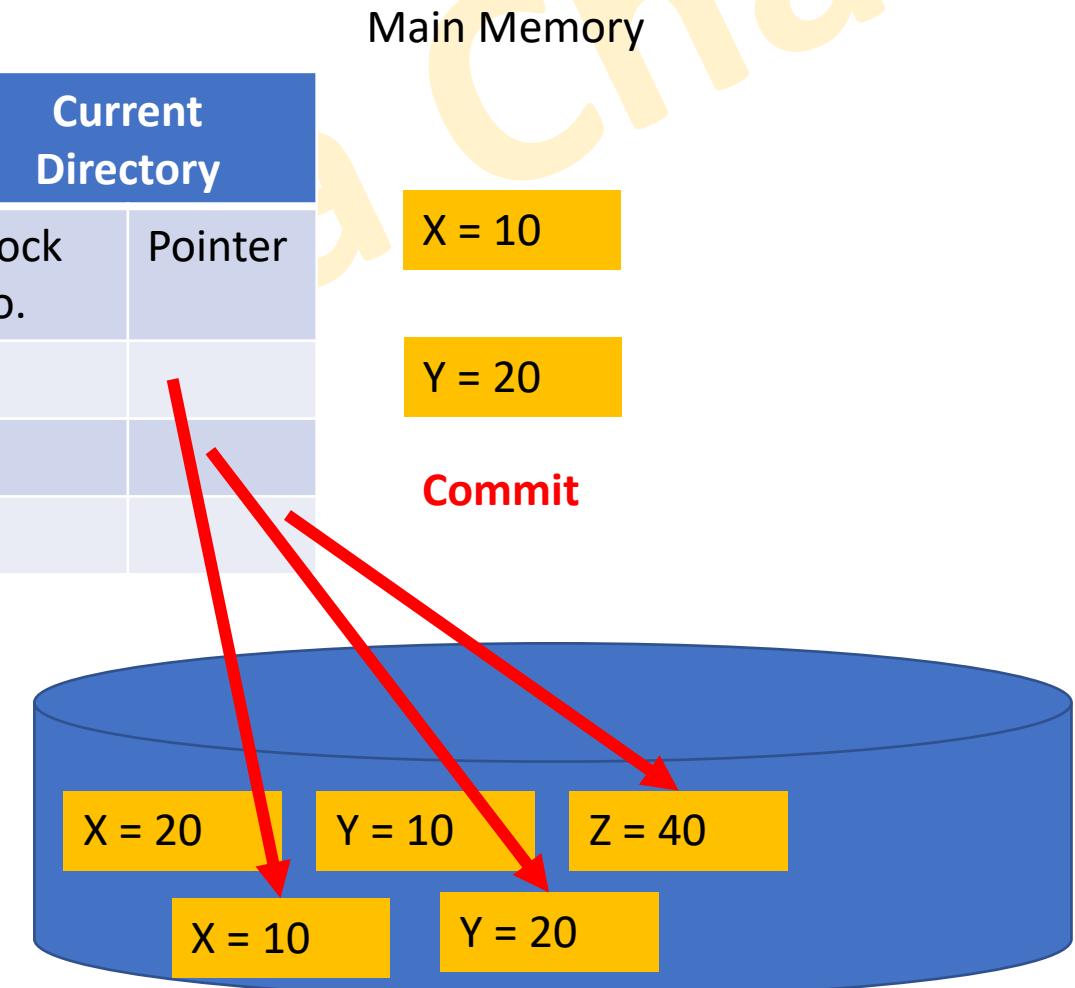
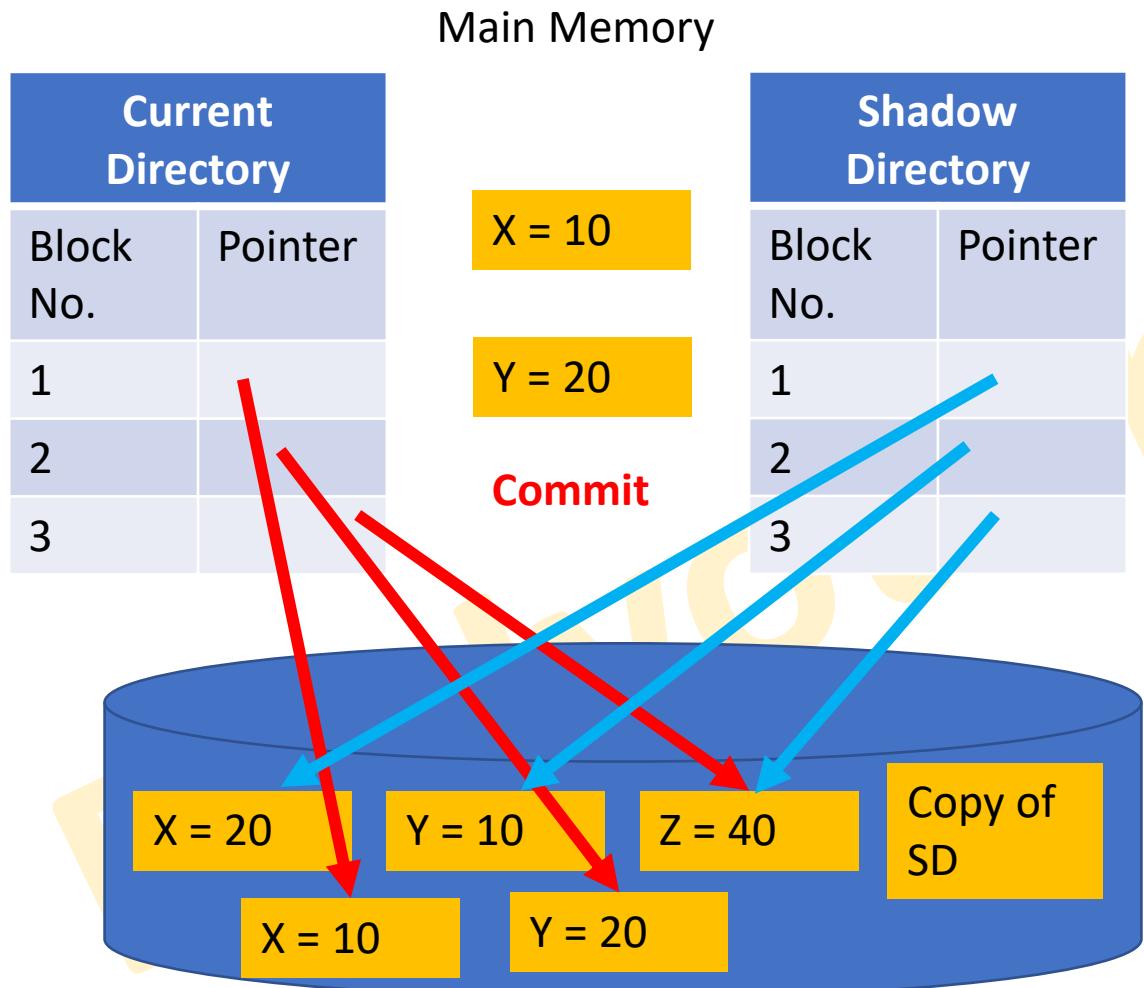
Shadow paging: Example



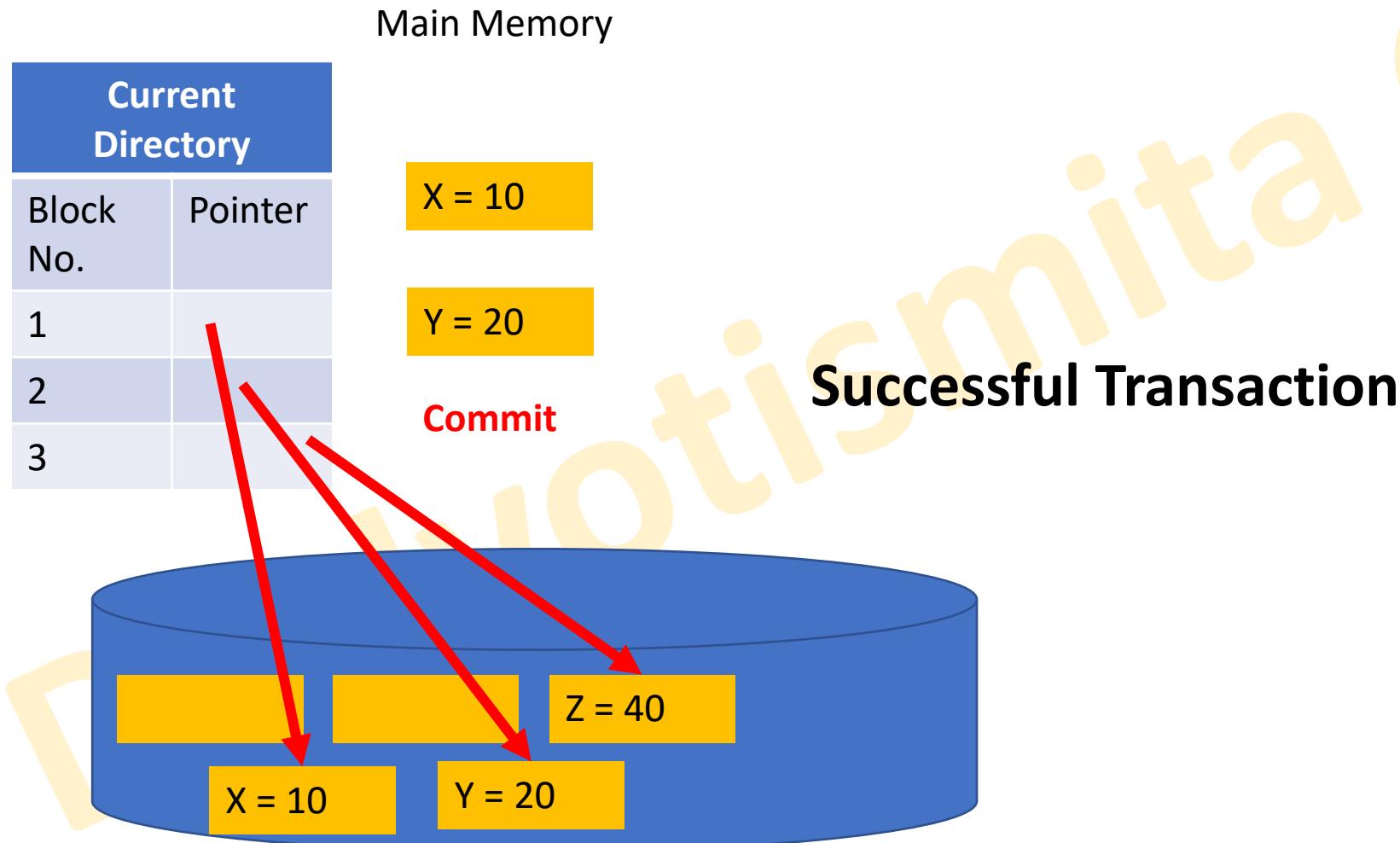
Shadow paging: Example



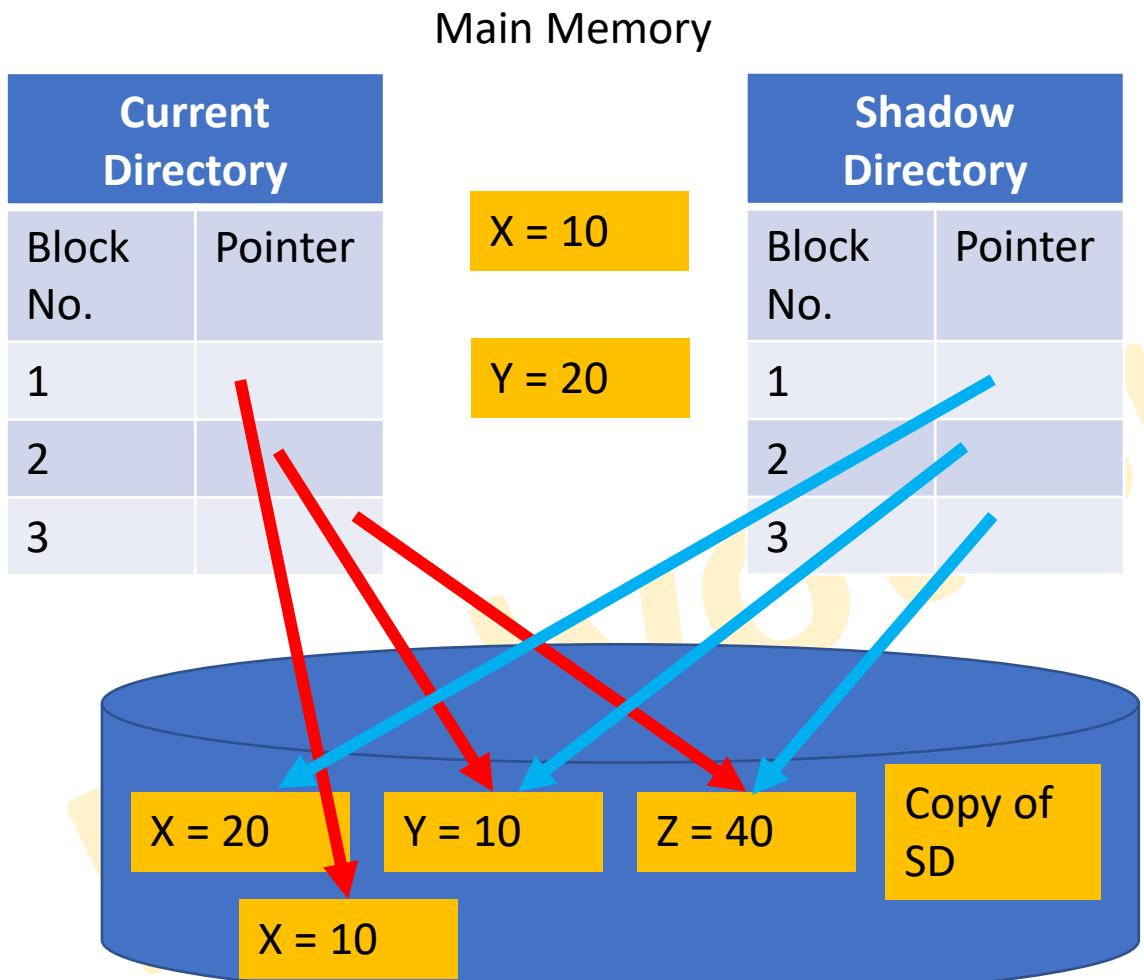
Shadow paging: Example



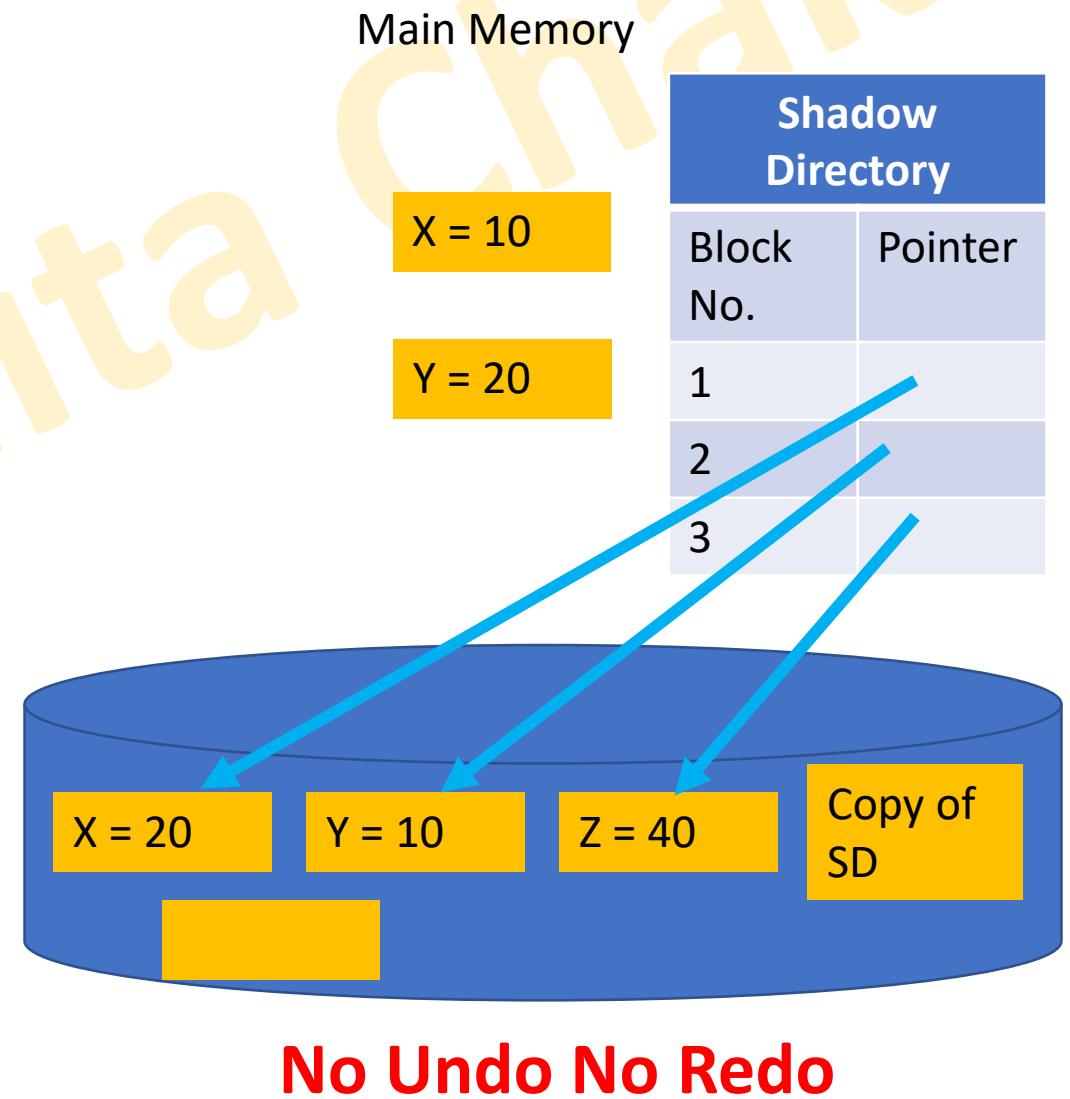
Shadow paging: Example



Shadow paging: Example



Transaction Fails



Introduction to NOSQL

- The term **NOSQL** is generally interpreted as **Not Only SQL**—rather than NO to SQL—and is meant to convey that many applications need systems other than traditional relational SQL systems to augment their data management needs.
- Most NOSQL systems are distributed databases or distributed storage systems, with a focus on semistructured data storage, high performance, availability, data replication, and scalability as opposed to an emphasis on immediate data consistency, powerful query languages, and structured data storage.
- developed to manage large amounts of data in organizations such as Google, Amazon, Facebook, and Twitter and in applications such as social media, Web links, user profiles, marketing and sales, posts and tweets, road maps and spatial data, and e-mail.

Need of NOSQL

- Many companies and organizations are faced with applications that store vast amounts of data. Consider a free e-mail application, such as Google Mail or Yahoo Mail or other similar service—this application can have millions of users, and each user can have thousands of e-mail messages. There is a need for a storage system that can manage all these e-mails; a structured relational SQL system may not be appropriate because
 - (1) SQL systems offer too many services (powerful query language, concurrency control, etc.), which this application may not need; and
 - (2) a structured data model such the traditional relational model may be too restrictive.

Need of NOSQL

- As another example, consider an application such as Facebook, with millions of users who submit posts, many with images and videos; then these posts must be displayed on pages of other users using the social media relationships among the users.
- User profiles, user relationships, and posts must all be stored in a huge collection of data stores, and the appropriate posts must be made available to the sets of users that have signed up to see these posts.
- Some of the data for this type of application is not suitable for a traditional relational system and typically needs multiple types of databases and data storage systems.

NOSQL: Few examples

- **Google** developed a proprietary NOSQL system known as **BigTable**, which is used in many of Google's applications that require vast amounts of data storage, such as Gmail, Google Maps, and Web site indexing. Apache Hbase is an open source NOSQL system based on similar concepts. Google's innovation led to the category of NOSQL systems known as **column-based** or **wide column** stores; they are also sometimes referred to as **column family** stores.
- **Amazon** developed a NOSQL system called **DynamoDB** that is available through Amazon's cloud services. This innovation led to the category known as **key-value** data stores or sometimes **key-tuple** or **key-object** data stores.

NOSQL: Few examples

- **Facebook** developed a NOSQL system called **Cassandra**, which is now open source and known as Apache Cassandra. This NOSQL system uses concepts from both key-value stores and column-based systems.
- Other software companies started developing their own solutions and making them available to users who need these capabilities—for example, **MongoDB** and **CouchDB**, which are classified as **document-based** NOSQL systems or **document stores**.
- Another category of NOSQL systems is the **graph-based** NOSQL systems, or **graph databases**; these include **Neo4J** and **GraphBase**, among others.

Characteristics of NOSQL Systems

- **Scalability:** There are two kinds of scalability in distributed systems: **horizontal** and **vertical**. In NOSQL systems, **horizontal scalability** is generally used, where the distributed system is expanded by adding more nodes for data storage and processing as the volume of data grows. **Vertical scalability**, on the other hand, refers to expanding the storage and computing power of existing nodes.
- **Availability, Replication and Eventual Consistency:** Many applications that use NOSQL systems require continuous system availability. To accomplish this, data is replicated over two or more nodes in a transparent manner, so that if one node fails, the data is still available on other nodes.

Characteristics of NOSQL Systems

- **Replication Models:** Two major replication models are used in NOSQL systems: master-slave and master-master replication. **Master-slave replication** requires one copy to be the master copy; all write operations must be applied to the master copy and then propagated to the slave copies. The **master-master replication** allows reads and writes at any of the replicas but may not guarantee that reads at nodes that store different copies see the same values. Different users may write the same data item concurrently at different nodes of the system, so the values of the item will be temporarily inconsistent.

Characteristics of NOSQL Systems

- **Sharding of Files:** In many NOSQL applications, files (or collections of data objects) can have many millions of records (or documents or objects), and these records can be accessed concurrently by thousands of users. So it is not practical to store the whole file in one node. **Sharding** (also known as **horizontal partitioning**) of the file records is often employed in NOSQL systems. This serves to distribute the load of accessing the file records to multiple nodes.

Characteristics of NOSQL Systems

- **High-Performance Data Access:** In many NOSQL applications, it is necessary to find individual records or objects (data items) from among the millions of data records or objects in a file. To achieve this, most systems use one of two techniques: **hashing** or **range partitioning** on object keys. In **hashing**, a **hash function** $h(K)$ is applied to the key K , and the location of the object with key K is determined by the value of $h(K)$. In **range partitioning**, the location is determined via a range of key values. In applications that require range queries, where multiple objects within a range of key values are retrieved, range partitioned is preferred.

Categories of NOSQL Systems

- NOSQL systems have been characterized into four major categories, with some additional categories that encompass other types of systems. The most common categorization lists the following four major categories:
 1. **Document-based NOSQL systems:** These systems store data in the form of documents using well-known formats, such as JSON (JavaScript Object Notation). Documents are accessible via their document id, but can also be accessed rapidly using other indexes.
 2. **NOSQL key-value stores:** These systems have a simple data model based on fast access by the key to the value associated with the key; the value can be a record or an object or a document or even have a more complex data structure.

Categories of NOSQL Systems

- NOSQL systems have been characterized into four major categories, with some additional categories that encompass other types of systems. The most common categorization lists the following four major categories:
 3. **Column-based or wide column NOSQL systems:** These systems partition a table by column into column families, where each column family is stored in its own files. They also allow versioning of data values.
 4. **Graph-based NOSQL systems:** Data is represented as graphs, and related nodes can be found by traversing the edges using path expressions.

CAP theorem

- The three letters in CAP refer to three desirable properties of distributed systems with replicated data: **consistency** (among replicated copies), **availability** (of the system for read and write operations) and **partition tolerance** (in the face of the nodes in the system being partitioned by a network fault).
- **Availability** means that each read or write request for a data item will either be processed successfully or will receive a message that the operation cannot be completed.
- **Partition tolerance** means that the system can continue operating if the network connecting the nodes has a fault that results in two or more partitions, where the nodes in each partition can only communicate among each other.
- **Consistency** means that the nodes will have the same copies of a replicated data item visible for various transactions.

CAP theorem

- The **CAP theorem** states that it is not possible to guarantee all three of the desirable properties—consistency, availability, and partition tolerance—at the same time in a distributed system with data replication.
- If this is the case, then the distributed system designer would have to choose two properties out of the three to guarantee.

NoSQL data models: Document-based NOSQL

- Document-based or document-oriented NOSQL systems typically store data as **collections** of similar **documents**.
- These types of systems are also sometimes known as **document stores**.
- A major difference between document-based systems versus object and object-relational systems and XML is that there is no requirement to specify a schema—rather, the documents are specified as **self-describing data**.
- Although the documents in a collection should be similar, they can have different data elements (attributes), and new documents can have new data elements that do not exist in any of the current documents in the collection.

NoSQL data models: Document-based NOSQL

- The system basically extracts the data element names from the self-describing documents in the collection, and the user can request that the system create indexes on some of the data elements.
- Documents can be specified in various formats, such as XML.
- A popular language to specify documents in NOSQL systems is JSON (JavaScript Object Notation).
- There are many document-based NOSQL systems, including MongoDB and CouchDB, among many others.

NoSQL data models: Document-based

NOSQL: MongoDB

- MongoDB documents are stored in BSON (Binary JSON) format, which is a variation of JSON with some additional data types and is more efficient for storage than JSON.
- Individual **documents** are stored in a **collection**.
- For example, the following command can be used to create a collection called **project** to hold PROJECT objects.
 - `db.createCollection("project", { capped : true, size : 1310720, max : 500 })`
 - The first parameter “project” is the name of the collection, which is followed by an optional document that specifies collection options.
 - In our example, the collection is capped; this means it has upper limits on its storage space (size) and number of documents (max).
 - The **capping parameters** help the system choose the storage options for each collection.

NoSQL data models: Document-based

NOSQL: MongoDB

- For our example, we will create another document collection called worker to hold information about the EMPLOYEEs who work on each project; for example:
 - db.createCollection("worker", { capped : true, size : 5242880, max : 2000 })
- In Figure (a), the workers information is embedded in the project document; so there is no need for the "worker" collection.
- This is known as the **denormalized** pattern.
- A list of values that is enclosed in square brackets [...] within a document represents a field whose value is an array.

NoSQL data models: Document-based

NOSQL: MongoDB

(a) project document with an array of embedded workers:

```
{  
    "_id": "P1",  
    "Pname": "ProductX",  
    "Plocation": "Bellaire",  
    "Workers": [  
        { "Ename": "John Smith",  
          "Hours": 32.5  
        },  
        { "Ename": "Joyce English",  
          "Hours": 20.0  
        }  
    ]  
};
```

NoSQL data models: Document-based

NOSQL: MongoDB

- Another option is to use the design in Figure (b), where worker references are embedded in the project document, but the worker documents themselves are stored in a separate “worker” collection.

(b) project document with an embedded array of worker ids:

```
{  
    _id: "P1",  
    Pname: "ProductX",  
    Plocation: "Bellaire",  
    WorkerIds: [ "W1", "W2" ]  
  
    {  
        _id: "W1",  
        Ename: "John Smith",  
        Hours: 32.5  
    }  
  
    {  
        _id: "W2",  
        Ename: "Joyce English",  
        Hours: 20.0  
    }  
}
```

Key-value stores

- Key-value stores focus on high performance, availability, and scalability by storing data in a distributed storage system.
- The data model used in key-value stores is relatively simple, and in many of these systems, there is no query language but rather a set of operations that can be used by the application programmers.
- The **key** is a unique identifier associated with a data item and is used to locate this data item rapidly.
- The **value** is the data item itself, and it can have very different formats for different key-value storage systems.

Key-value stores: DynamoDB

- The DynamoDB system is an Amazon product and is available as part of Amazon's **AWS/SDK** platforms (Amazon Web Services/Software Development Kit).
- It can be used as part of **Amazon's cloud computing services**, for the data storage component.
- The basic data model in DynamoDB uses the concepts of tables, items, and attributes.
- A **table** in DynamoDB does not have a **schema**; it holds a collection of **self-describing** items.
- Each **item** will consist of a number of (attribute, value) pairs, and attribute values can be single-valued or multivalued.
- When a table is created, it is required to specify a **table name** and a **primary key**; the primary key will be used to rapidly locate the items in the table.
- Thus, the primary key is the **key** and the item is the **value** for the DynamoDB key-value store.

Key-value stores: Voldemort

- Voldemort is an open source system available through Apache 2.0 open source licensing rules.
- It is based on Amazon's DynamoDB.
- The focus is on high performance and horizontal scalability, as well as on providing replication for high availability.
- A collection of **(key, value)** pairs is kept in a Voldemort store.
- We will assume the store is called `s`.
 - The operation `s.put(k, v)` inserts an item as a key-value pair with key `k` and value `v`.
 - The operation `s.delete(k)` deletes the item whose key is `k` from the store.
 - The operation `v = s.get(k)` retrieves the value `v` associated with key `k`.

Column families

- Another category of NOSQL systems is known as **column-based** or **wide column** systems.
- The Google distributed storage system for big data, known as **BigTable**, is a well-known example of this class of NOSQL systems, and it is used in many Google applications that require large amounts of data storage, such as Gmail.
- BigTable uses the **Google File System (GFS)** for data storage and distribution.

Column families

- A **table** is associated with one or more **column families**.
- Each column family will have a name, and the column families associated with a table must be specified when the table is created and cannot be changed later.
- Figure (a) shows how a table may be created; the table name is followed by the names of the column families associated with the table.
- When the data is loaded into a table, each column family can be associated with many **column qualifiers**, but the column qualifiers are not specified as part of creating a table.
- So the column qualifiers make the model a self-describing data model because the qualifiers can be dynamically specified as new rows are created and inserted into the table.

(a) **creating a table:**

```
create 'EMPLOYEE', 'Name', 'Address', 'Details'
```

Column families

- A column is specified by a combination of ColumnFamily:ColumnQualifier.
- Basically, column families are a way of grouping together related columns (attributes in relational terminology) for storage purposes, except that the column qualifier names are not specified during table creation.
- Rather, they are specified when the data is created and stored in rows, so the data is self-describing since any column qualifier name can be used in a new row of data.

(b) inserting some row data in the EMPLOYEE table:

```
put 'EMPLOYEE', 'row1', 'Name:Fname', 'John'
put 'EMPLOYEE', 'row1', 'Name:Lname', 'Smith'
put 'EMPLOYEE', 'row1', 'Name:Nickname', 'Johnny'
put 'EMPLOYEE', 'row1', 'Details:Job', 'Engineer'
put 'EMPLOYEE', 'row1', 'Details:Review', 'Good'
put 'EMPLOYEE', 'row2', 'Name:Fname', 'Alicia'
put 'EMPLOYEE', 'row2', 'Name:Lname', 'Zelaya'
put 'EMPLOYEE', 'row2', 'Name:MName', 'Jennifer'
put 'EMPLOYEE', 'row2', 'Details:Job', 'DBA'
put 'EMPLOYEE', 'row2', 'Details:Supervisor', 'James Borg'
put 'EMPLOYEE', 'row3', 'Name:Fname', 'James'
put 'EMPLOYEE', 'row3', 'Name:Minit', 'E'
put 'EMPLOYEE', 'row3', 'Name:Lname', 'Borg'
put 'EMPLOYEE', 'row3', 'Name:Suffix', 'Jr.'
put 'EMPLOYEE', 'row3', 'Details:Job', 'CEO'
put 'EMPLOYEE', 'row3', 'Details:Salary', '1,000,000'
```