



Data Management

Maurizio Lenzerini

***Dipartimento di Informatica e Sistemistica “Antonio Ruberti”
Università di Roma “La Sapienza”***

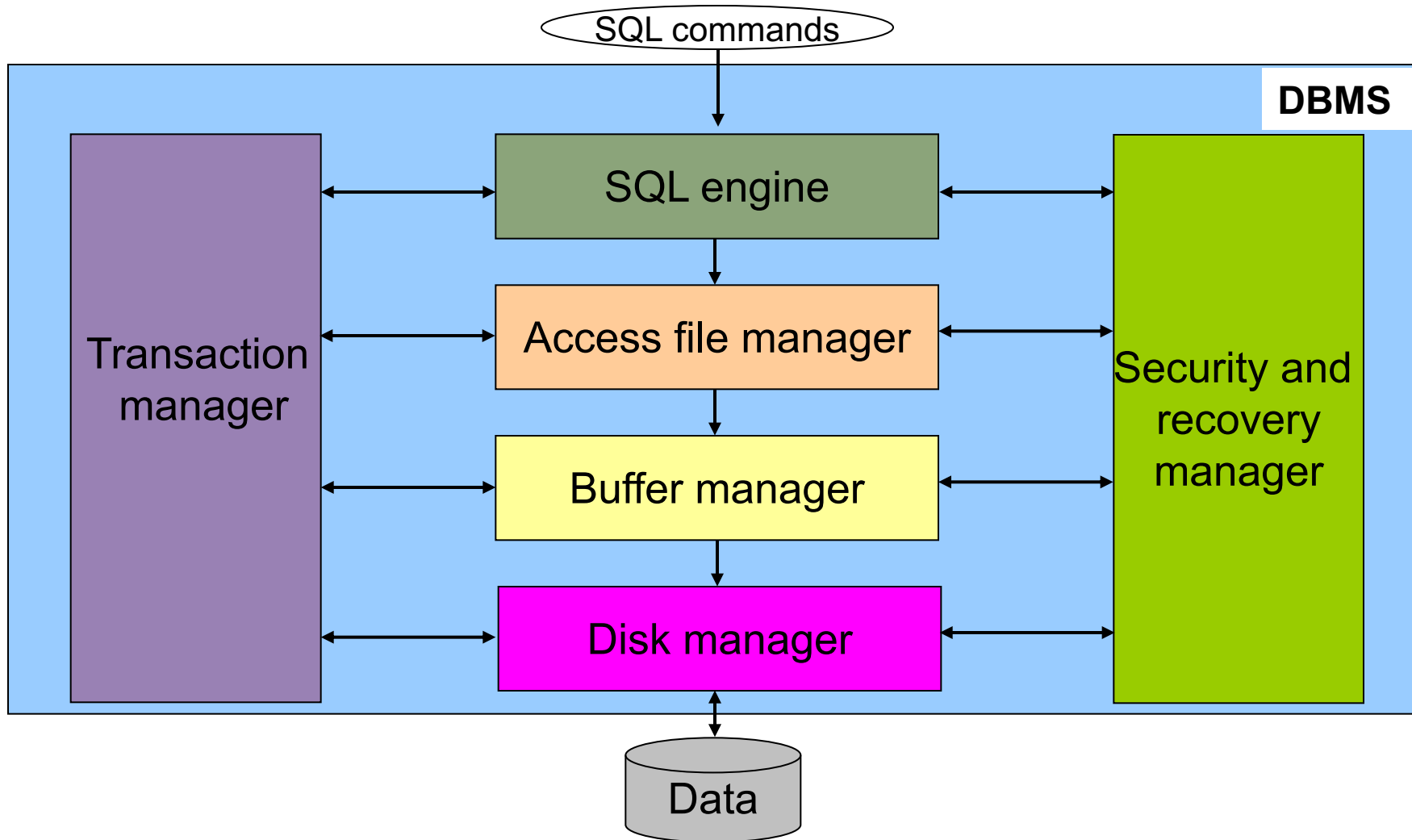
Academic Year 2018/2019

*Part 3
Transaction management and concurrency*

<http://www.dis.uniroma1.it/~lenzerin/index.html/?q=node/53>



Architecture of a DBMS





5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 Concurrency control through locks

5.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 Concurrency control through locks

5.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



Transactions

A **transaction** models the execution of a software procedure constituted by a set of instructions that may "read from" and "write on" a database, and **that form a single logical unit**.

Syntactically, we will assume that every transaction contains:

- one "begin" instruction
- one "end" instruction
- one among "commit" (confirm what you have done on the database so far) and "rollback" (undo what you have done on the database so far)

As we will see, each transaction should enjoy a set of properties (called ACID)



Concurrency

The **throughput** of a system is the number of transactions per second (tps) accepted by the system

In a DBMS, we want the throughput to be approximately **100-1000tps**

This means that the system should support a high degree of concurrency among the transactions that are executed

- **Example:** If each transaction needs 0.1 seconds in the average for its execution, then to get a throughput of 100tps, we must ensure that 10 transactions are executed concurrently in the average

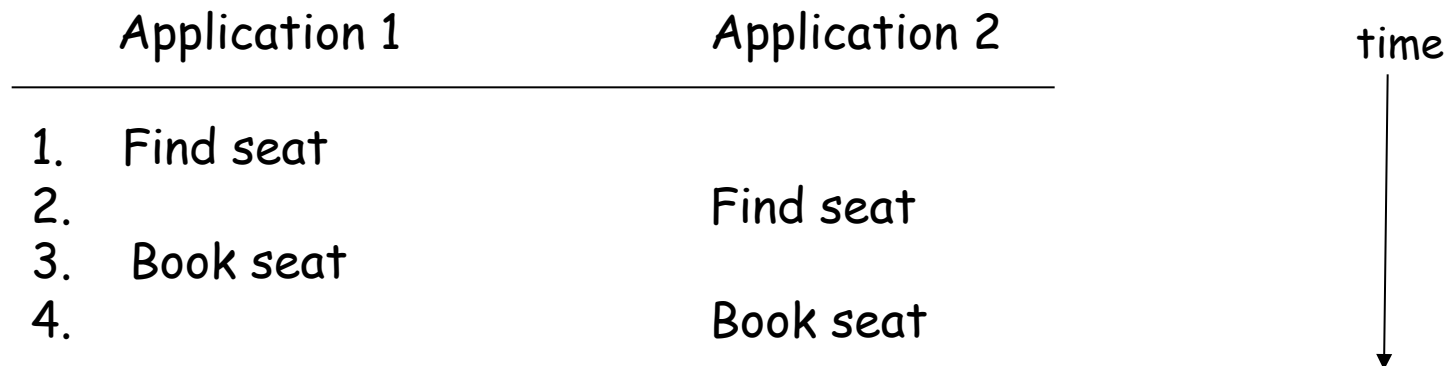
Typical applications: banks, flight reservations, web commerce...



Concurrency: example

Suppose that the same program is executed concurrently by two applications aiming at reserving a seat in the same flight

The following temporal evolution is possible:



The result is that we have two reservations for the same seat!



Isolation of transactions

The DBMS deals with this problem by ensuring the so-called “isolation” property for the transactions

This property for a transaction essentially means that it is executed like it was the only one in the system, i.e., without concurrent transactions

While isolation is essential, other properties are important as well



Desirable properties of transactions

The desirable properties in transaction management are called the **ACID** properties. They are:

1. **Atomicity**: for each transaction execution, either all or none of its actions have their effect
2. **Consistency**: each transaction execution brings the database to a correct state
3. **Isolation**: each transaction execution is independent of any other concurrent transaction executions
4. **Durability**: if a transaction execution succeeds, then its effects are registered permanently in the database



Schedules and serial schedules

Given a set of transactions $\{T_1, T_2, \dots, T_n\}$, a sequence S of executions of actions of such transactions respecting the order within each transaction (i.e., such that if action a is before action b in a transaction T_i , then a is before b also in S) is called **schedule on $\{T_1, T_2, \dots, T_n\}$** , or simply **schedule**.

A schedule on $\{T_1, T_2, \dots, T_n\}$ that does not contain all the actions of all transactions T_1, T_2, \dots, T_n is called **partial**

A schedule S is called **serial** if the actions of each transaction in S come before every action of a different transaction in S , i.e., if in S the actions of different transactions do not interleave.



Serializability

Example of serial schedules:

Given T1 ($x = x + x$; $x = x + 3$) and T2 ($x = x^2$; $x = x + 2$), possible serial schedules on them are:

Sequence 1: $x = x + x$; $x = x + 3$; $x = x^2$; $x = x + 2$

Sequence 2: $x = x^2$; $x = x + 2$; $x = x + x$; $x = x + 3$

A serial schedule is obviously “correct” with respect to concurrency, because it does not have interleaving. What about a schedule having interleaving?

Intuitively, we would like to say that a schedule that is not serial is “correct” with respect to concurrency if it is **serializable**, i.e., if for any initial state IS, its outcome is the same as the outcome of serial schedule constituted by the same transactions of S starting from IS.



Serializability

Definition of serializable schedule A schedule S on $\{T_1, T_2, \dots, T_n\}$ is serializable if there exists a serial schedule on $\{T_1, T_2, \dots, T_n\}$ that is “equivalent” to S .

But what does “equivalent” mean?

Definition of equivalent schedules: Two schedules S_1 and S_2 are said to be **equivalent** if, for each database state D , the execution of S_1 starting from the database state D produces the same outcome as the execution of S_2 starting from the same database state D .

Notice that when we talk about the outcome we talk about the final state of the process represented by the schedule, which incorporates both the state of the database, and the state of the local store.



Notation

A successful execution of transaction can be represented as a sequence of

- Commands of type **begin/commit**
- Actions that **read** and **write** an element (attribute, record, table) in the database
- Actions that **read** and **write** an element in the **local store**

T_1	T_2
begin	begin
READ(A,t)	READ(A,s)
$t := t+100$	$s := s*2$
WRITE(A,t)	WRITE(A,s)
READ(B,t)	READ(B,s)
$t := t+100$	$s := s*2$
WRITE(B,t)	WRITE(B,s)
commit	commit



A serial schedule

T ₁	T ₂	A	B
		25	25
begin			
READ(A,t)			
t := t+100			
WRITE(A,t)		125	
READ(B,t)			
t := t+100			
WRITE(B,t)			125
commit	begin		
	READ(A,s)		
	s := s*2		
	WRITE(A,s)	250	
	READ(B,s)		
	s := s*2		
	WRITE(B,s)		250
	commit		



A serializable schedule

T ₁	T ₂	A	B
		25	25
begin	begin		
READ(A,t)			
t := t+100			
WRITE(A,t)		125	
	READ(A,s)		
	s := s*2		
	WRITE(A,s)	250	
READ(B,t)			
t := t+100			
WRITE(B,t)			125
commit			
	READ(B,s)		
	s := s*2		
	WRITE(B,s)		250
	commit		

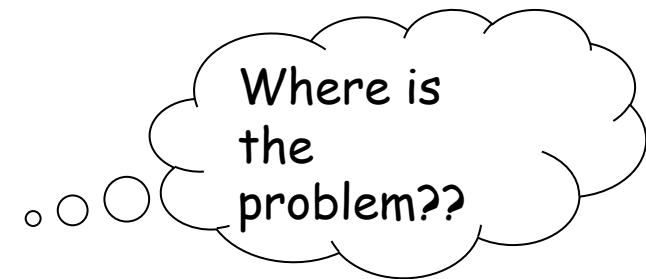
The final values of A and B are the same as the serial schedule T1, T2, no matter what the initial (identical) value of A and B.

We can indeed show that, if initially $A=B=c$ (c is a constant), then at the end of the execution of the schedule we have: $A=B=2(c+100)$



A non-serializable schedule

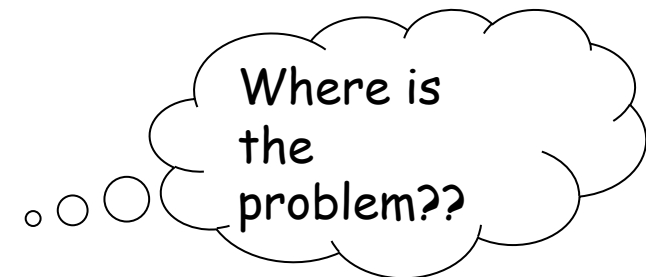
T ₁	T ₂	A	B
		25	25
begin	begin		
READ(A,t)			
t := t+100			
WRITE(A,t)		125	
	READ(A,s)		
	s := s*2		
	WRITE(A,s)	250	
	READ(B,s)		
	s := s*2		
	WRITE(B,s)		50
READ(B,t)			
t := t+100			
WRITE(B,t)			150
commit	commit		





A non-serializable schedule

T_1	T_2	A	B
		25	25
begin	begin		
READ(A,t)			
$t := t + 100$			
WRITE(A,t)		125	
	READ(A,s)		
	$s := s * 2$		
	WRITE(A,s)	250	
	READ(B,s)		
	$s := s * 2$		
	WRITE(B,s)		50
READ(B,t)			
$t := t + 100$			
WRITE(B,t)			150
commit	commit		





Anomaly 1: reading temporary data (WR anomaly)

T_1	T_2
begin	begin
READ(A,x)	
$x := x-1$	
WRITE(A,x)	
	READ(A,x)
	$x := x*2$
	WRITE(A,x)
	READ(B,x)
	$x := x*2$
	WRITE(B,x)
	commit
READ(B,x)	
$x := x+1$	
WRITE(B,x)	
commit	

Note that the interleaved execution is different from any serial execution. The problem comes from the fact that the value of A is read by T2 after T1 has written on A, whereas the value of B is read by T2 before T1 has written on B.

This is a **WR (write-read) anomaly**, because it shows up when a transaction writes an element, and another transaction reads such element.



Anomaly 2a: update loss (RW anomaly)

- Let T_1 , T_2 be two transactions, each of the form:
 $\text{READ}(A, x), x := x + 1, \text{WRITE}(A, x)$
- The serial execution with initial value $A=2$ produces $A=4$, which is the result of two subsequent updates
- Now, consider the following schedule:

T_1	T_2
begin	begin
READ(A,x)	
$x := x+1$	
	READ(A,x)
	$x := x+1$
WRITE(A,x)	
commit	
	WRITE(A,x)
	commit

Note that the interleaved execution is different from any serial execution. The final result is $A=3$, and the first update is lost: T_2 reads the initial value of A , and writes the final value. In this case, the update executed by T_1 is lost!



Anomaly 2a: update loss (RW anomaly)

- This kind of anomaly is called **RW anomaly** (read-write anomaly), because it shows up when a transaction reads an element, and another transaction writes the same element.
- Indeed, this anomaly comes from the fact that a transaction T2 could change the value of an object A that has been read by a transaction T1, while T1 is still in progress. The fact that T1 is still in progress means that the risk is that T1 works on A without taking into account the changes that T2 makes on A. Therefore, the update of T1 or T2 are lost.



Anomaly 2b: unrepeateable read (RW anomaly)

T_1 executes two consecutive reads of the same data (assume the initial vale of A is 20):

T_1	T_2
begin READ(A,x)	begin
	$x := 100$ WRITE(A,x) commit
READ(A,x) commit	

However, due to the concurrent update of T_2 , T_1 reads two different values.

Note that the interleaved execution is different from any serial execution. This is another kind of **RW (read-write) anomaly**.



Anomaly 3: ghost update (WW anomaly)

Assume the following integrity constraint $A = B$

T_1	T_2
begin WRITE(A,1)	begin WRITE(B,2)
WRITE(B,1) commit	WRITE(A,2) commit

Note that neither T_1 nor T_2 in isolation violate the integrity constraints. However, the interleaved execution is different from any serial execution. Transaction T_1 will see the update of A to 2 as a surprise, and transaction T_2 will see the update of B to 1 as a surprise.

Note that the interleaved execution is different from any serial execution. This is a **WW (write-write) anomaly**



Scheduler

The **scheduler** is part of the transaction manager, and works as follows:

- It deals with new transactions entered into the system, assigning them an identifier
- It instructs the buffer manager so as to read and write on the DB according to a particular sequence
- It is NOT concerned with specific operations on the local store of transactions, nor with constraints on the order of executions of transactions. The last conditions means that **every order by which transactions are entered into the system is acceptable to the schedule.**

It follows that we can simply characterize each transaction T_i (where i is a nonnegative integer identifying the transaction) in terms of its actions, where each action of transaction T_i is denoted by a letter (read, write, or commit) and the subscript i . **In other words, we ignore the operations on main memory.**

The transactions of the previous examples are written as:

$T_1: r_1(A) r_1(B) w_1(A) w_1(B) c_1$

$T_2: r_2(A) r_2(B) w_2(A) w_2(B) c_2$

An example of (complete) schedule on these transactions is:

$r_1(A) r_1(B) w_1(A) r_2(A) r_2(B) w_2(A) w_1(B) c_1 w_2(B) c_2$

T1 reads A

T2 writes A

T1 commit



Serializability and equivalence of schedules

As we saw before, the definition of serializability relies on the notion of equivalence between schedules.

Depending on the level of abstraction used to characterize the effects of transactions, we get **different notions of equivalence**, which in turn suggest **different definitions of serializability**.

Given a certain definition of equivalence, we will be interested in

- two types of algorithms:
 - algorithms for **checking equivalence**: given two schedules, determine if they are equivalent
 - algorithms for **checking serializability**: given one schedule, check whether it is equivalent to any of the serial schedules on the same transactions
- rules that ensures serializability



Two important assumptions

In the next slides, we will work under two assumptions:

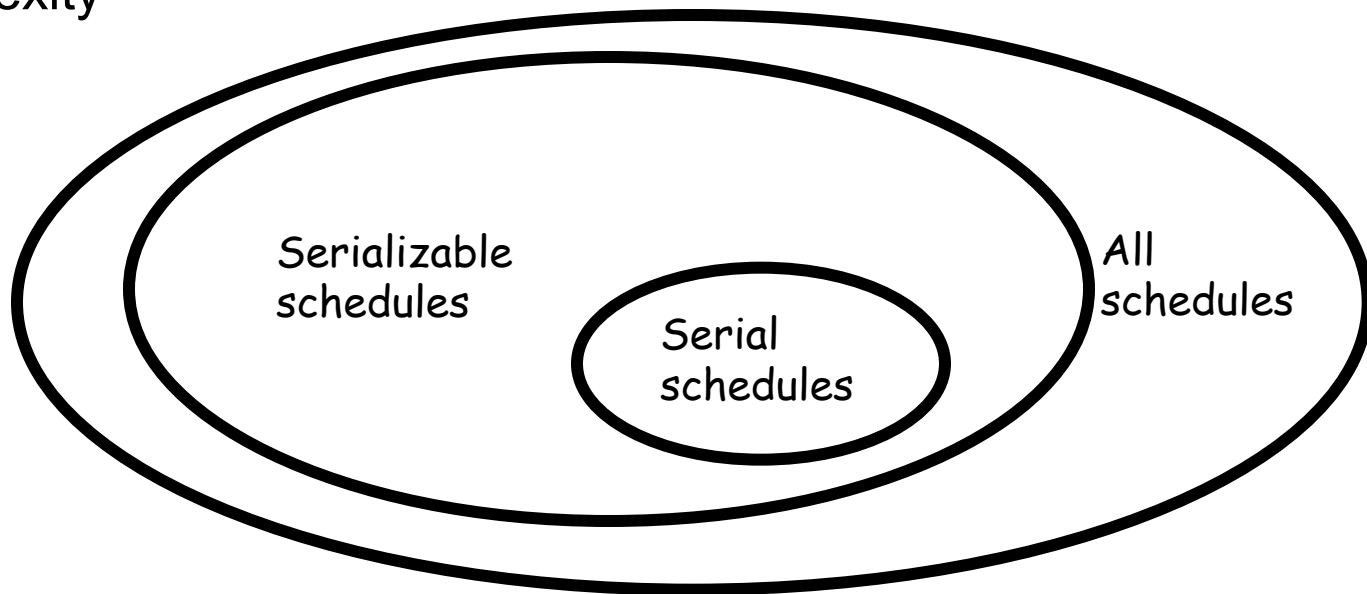
1. No transaction reads or writes the same element **twice** (sometimes, we will relax this assumption in examples), and **no transaction reads an element that it has written**
2. **No transaction executes the “rollback” command** (i.,e. all executions of transactions are successful)

Later on, we will remove the second assumption



Classes of schedules

Basic idea of our investigation: single out classes of schedules that are serializable, and such that the serializability check can be done (i.e., the problem is decidable), and can be done with reasonable computational complexity



We will define several notions of serializability, starting with

- view-serializability
- conflict-serializability



2. Transaction management

2.1 Transactions, concurrency, serializability

2.2 **View-serializability**

2.3 Conflict-serializability

2.4 Concurrency control through locks

2.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



View-equivalence and view-serializability

Preliminary definitions:

- In a schedule S , we say that $r_i(x)$ **READS-FROM** $w_j(x)$ if $w_j(x)$ precedes $r_i(x)$ in S , and there is no action of type $w_k(x)$ between $w_j(x)$ and $r_i(x)$
- In a schedule S , we say that $w_i(x)$ is a **FINAL-WRITE** if $w_i(x)$ is the last write action on x in S

Definition of view-equivalence: let $S1$ and $S2$ be two (complete) schedules on the same transactions. Then $S1$ is view-equivalent to $S2$ if $S1$ and $S2$ have the same READS-FROM relation, and the same FINAL-WRITE set.

Definition of view-serializability: a (complete) schedule S on $\{T1, \dots, Tn\}$ is view-serializable if there exists a serial schedule S' on $\{T1, \dots, Tn\}$ that is view-equivalent to S



View-serializability

read1(A,t) read2(A,s) s:=100 write2(A,s) t:=100 write1(A,t)

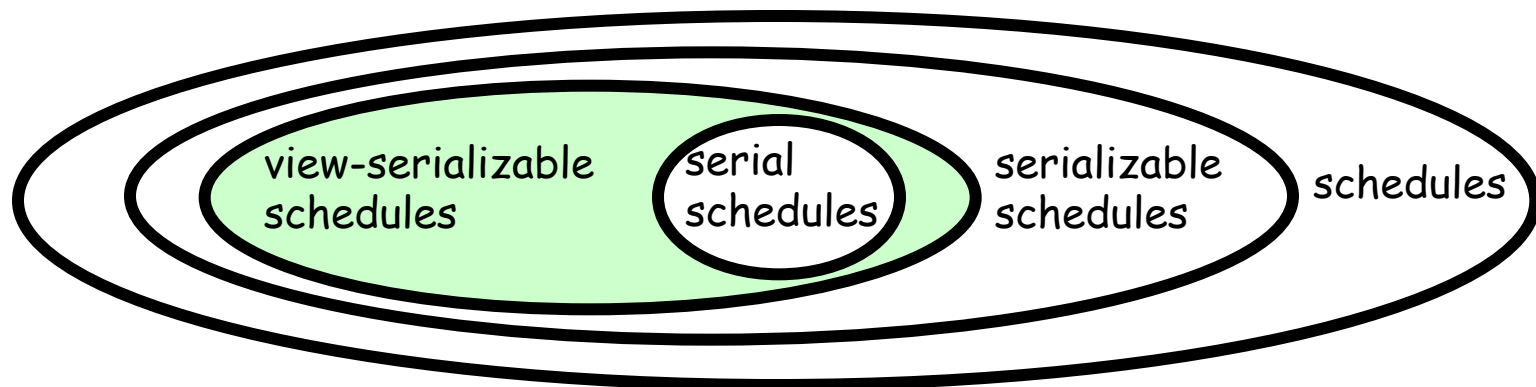
Serializable?

View-serializable?



View-serializability

- There are serializable schedules that are not view-serializable. For example,
 `read1(A,t) read2(A,s) s:=100 write2(A,s) t:=100 write1(A,t)`
is serializable, but not view-serializable
- Note however, that in order to realize that the above schedule is serializable, we need to take into account the operations performed on the local store
- If we limit our attention to our abstract model of transaction (where only read and write operations count), and we consider as outcome of a schedule the database state, then view-equivalence and view-serializability are the most general notions





Properties of view-equivalence

- Given two schedules, checking whether they are view-equivalent can be done in polynomial time
- Given one schedule, checking whether it is view-serializable is an NP-complete problem
 - It is easy to verify that the problem is in NP; this is a nondeterministic polynomial time algorithm for checking whether S is view-serializable or not: non deterministically guess a serial schedule S' on the transactions of S, and then check in polynomial time if S' is view-equivalent to S
 - Proving that the problem is NP-hard is much more difficult
- The above is one reason why view-serializability is not used in practice



Monotone classes of schedules

Notation:

- for a schedule S , $\text{Tran}(S)$ denotes the set of transactions present in S ;
- for $T \subseteq \text{Tran}(S)$, $\Pi_T(S)$ denotes the projection of S onto T , i.e., the schedule S' obtained from S by deleting all operations of the transactions that are not in T .

For example, if

$S = w_1(x) \ r_2(x) \ w_2(y) \ r_1(y) \ r_3(x) \ w_1(y) \ w_3(x) \ w_3(y) \ c_1 \ a_2$

and $T = \{t_1, t_2\}$, then

$\Pi_T(S) = w_1(x) \ r_2(x) \ w_2(y) \ r_1(y) \ w_1(y) \ c_1 \ a_2$

Definition of monotone class of schedule: A class E of schedules is called **monotone** if the fact that S is in E implies that for all $T \subseteq \text{Tran}(S)$, $\Pi_T(S)$ is in E too (i.e., E is closed under projection)



Monotone classes of schedules

- From the definition of monotonicity it follows that, if E is a monotone class of schedules, and a partial schedule constituted by a projection of a schedule S is not in E , then S is not in E too. A scheduler based on E can disregard a partial schedule s if it is not in E (on the basis of the fact that no extension of s can be in E)

- Unfortunately, the class of view-serializable schedules is not monotone, as this example shows:

$S = w_1(x) w_2(x) w_2(y) c_2 w_1(y) c_1 w_3(x) w_3(y) c_3$

It is easy to see that S is view-equivalent to $(t_1 t_2 t_3)$ and to $(t_2 t_1 t_3)$, and therefore it is view-serializable. However, $\Pi_{\{t_1, t_2\}}(S)$ is not view-serializable.

- Nonmonotonicity is another reason why view-serializability is **not** used in practice



Exercise 1a

- Consider the schedules:
 1. $w_0(x) \ r_2(x) \ r_1(x) \ w_2(x) \ w_2(z)$
 2. $w_0(x) \ r_1(x) \ r_2(x) \ w_2(x) \ w_2(z)$
 3. $w_0(x) \ r_1(x) \ w_1(x) \ r_2(x) \ w_1(z)$
 4. $w_1(x) \ r_2(x) \ w_2(y) \ r_1(y)$and tell which of them are view-serializable
- Consider the following schedules, verify that they are not view-serializable, and tell which anomalies they suffer from
 1. $r_1(x) \ r_2(x) \ w_1(x) \ w_2(x)$
 2. $r_1(x) \ w_2(x) \ r_1(x)$
 3. $w_1(x) \ w_2(y) \ w_1(y) \ w_2(x)$



Exercise 1b

Consider the three transactions T1, T2, T3 defined as follows:

- $T1 = r1(A), w1(A)$
- $T2 = r2(A), w2(A)$
- $T3 = r3(A), w3(A)$

and tell if there is at least one non-serial schedule on T1, T2, T3 that is view-serializable, motivating the answer.



5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 **Conflict-serializability**

5.4 Concurrency control through locks

5.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



Conflict-serializability: the notion of conflict

Definition of conflicting actions: Two actions are **conflicting** in a schedule if they belong to different transactions, they operate on the same element, and at least one of them is a write.

It is easy to see that:

- Two consecutive nonconflicting actions belonging to different transactions can be swapped without changing the effects of the schedule. Indeed,
 - Two consecutive reads of the same elements in different transactions can be swapped
 - One read of X in T1 and a consecutive read of Y in T2 (with $Y \neq X$) can be swapped
- The swap of two consecutive actions of the same transaction can change the effect of the transaction
- Two conflicting consecutive actions cannot be swapped without changing the effects of the schedule, because:
 - Swapping two write operations $w1(A) w2(A)$ on the same elements may result in a different final value for A
 - Swapping two consecutive operations such as $r1(A) w2(A)$ may cause T1 read different values of A (before and after the write of T2, respectively)



Conflict-equivalence

Definition of conflict-equivalence: Two schedules S1 and S2 on the same transactions are **conflict-equivalent** if S1 can be transformed into S2 through a sequence of swaps of consecutive nonconflicting actions

Exemple:

$S = r1(A) w1(A) r2(A) w2(A) r1(B) w1(B) r2(B) w2(B)$

is conflict-equivalent to:

$S' = r1(A) w1(A) r1(B) w1(B) r2(A) w2(A) r2(B) w2(B)$

because it can be transformed into S' through the following sequence of swaps:

$r1(A) w1(A) r2(A) \underline{w2(A) r1(B)} w1(B) r2(B) w2(B)$

$r1(A) w1(A) \underline{r2(A) r1(B)} w2(A) w1(B) r2(B) w2(B)$

$r1(A) w1(A) r1(B) r2(A) \underline{w2(A) w1(B)} r2(B) w2(B)$

$r1(A) w1(A) r1(B) \underline{r2(A) w1(B)} w2(A) r2(B) w2(B)$

$r1(A) w1(A) r1(B) w1(B) r2(A) w2(A) r2(B) w2(B)$



Exercise 2

Prove the following property:

Two schedules $S1$ and $S2$ on the same transactions $T1, \dots, Tn$ are **conflict-equivalent** **if and only if** there are no actions a_i of T_i and b_j of T_j (with T_i and T_j belonging to $T1, \dots, Tn$) such that

- a_i and b_j are conflicting, and
- the mutual position of the two actions in $S1$ is different from their mutual position in $S2$

This property is extremely important, because it allows us to check conflict-equivalence in a very direct way, without resorting to trying sequences of swaps.



Conflict-serializability

Definition of conflict-serializability: A schedule S is **conflict-serializable** if there exists a serial schedule S' that is conflict-equivalent to S

How can conflict-serializability be checked?

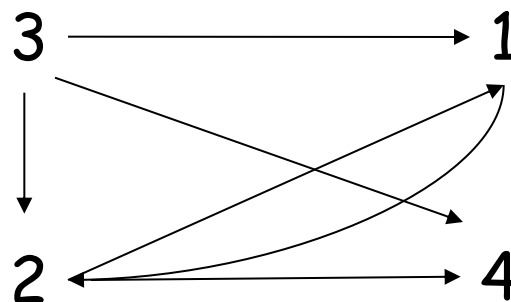
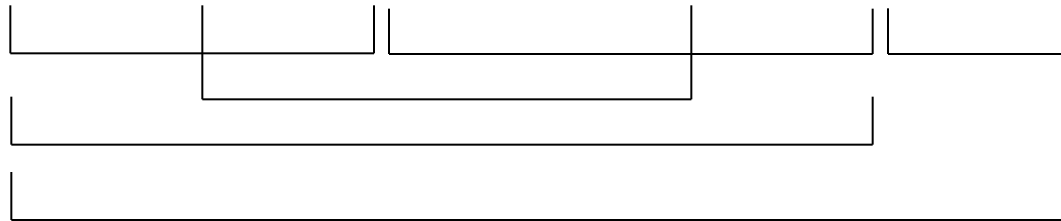
We can do it by analyzing the **precedence graph** associated to a schedule. Given a schedule S on T_1, \dots, T_n , the precedence graph $P(S)$ associated to S is defined as follows:

- the nodes of $P(S)$ are the transactions $\{T_1, \dots, T_n\}$ of S
- the edges E of $P(S)$ are as follows: the edge $T_i \rightarrow T_j$ is in E if and only if there exists two actions $P_i(A)$, $Q_j(A)$ of different transactions T_i and T_j in S operating on the same object A such that:
 - $P_i(A) <_S Q_j(A)$ (i.e., $P_i(A)$ appears before $Q_j(A)$ in S)
 - at least one between $P_i(A)$ and $Q_j(A)$ is a write operation



Example of precedence graph

S: $w_3(A)$ $w_2(C)$ $r_1(A)$ $w_1(B)$ $r_1(C)$ $w_2(A)$ $r_4(A)$ $w_4(D)$





How the precedence graph is used

Theorem (conflict-serializability) A schedule S is conflict-serializable if and only if the precedence graph $P(S)$ associated to S is acyclic.

To prove the theorem:

- we observe that if S is a serial schedule, then the precedence graph $P(S)$ is acyclic (easy to prove)
- we prove a preliminary lemma

Exercise 2': Prove that, if S is a serial schedule, then the precedence graph $P(S)$ is acyclic.



Preliminary lemma

Lemma If two schedules $S1$ and $S2$ are conflict-equivalent, then $P(S1) = P(S2)$

Proof Let $S1$ and $S2$ be two conflict-equivalent schedules, and assume that $P(S1) \neq P(S2)$. Then, $P(S1)$ and $P(S2)$ have the same nodes and different edges, i.e., there exists one edge $T_i \rightarrow T_j$ in $P(S1)$ that is not in $P(S2)$. But $T_i \rightarrow T_j$ in $P(S1)$ means that $S1$ has the form

... $p_i(A)$... $q_j(A)$...

with conflicting p_i, q_j . In other words, $p_i(A) <_{S1} q_j(A)$. Since $P(S2)$ has the same nodes as $P(S1)$, $S2$ contains $q_j(A)$ and $p_i(A)$, and since $P(S2)$ does not contain the edge $T_i \rightarrow T_j$, we can conclude that $q_j(A) <_{S2} p_i(A)$. But then, $S1$ and $S2$ differ in the order of a conflicting pair of actions, and therefore they cannot be transformed one into the other through the swap of two non-conflicting actions. This means that they are not conflict-equivalent, and we get a contradiction. Hence, we conclude that $P(S1)=P(S2)$.



The converse does not hold

If the converse of the previous lemma held, then the conflict-serializability theorem would already be proved. However, the converse does not hold. In fact, we can prove that $P(S1)=P(S2)$ does not imply that $S1$ and $S2$ are conflict-equivalent.

Indeed:

$S1 = w1(A) \ r2(A) \ w2(B) \ r1(B)$

$S2 = r2(A) \ w1(A) \ r1(B) \ w2(B)$

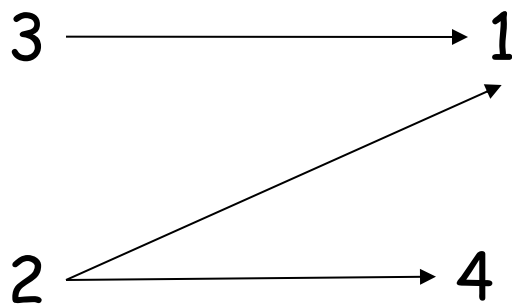
have the same precedence graph, but they are not conflict-equivalent, since to transform one of them into the other requires swapping two conflicting actions (or, equivalently, since they contain a pair of conflicting actions appearing in different order in the two schedules).



Topological order of a graph

Definition of topological order: Given a graph G , the **topological order** of G is a total order S (i.e., a sequence) of the nodes of G such that if the edge $T_i \rightarrow T_j$ is in the graph G , then T_i appears before T_j in the sequence S .

Example



Possible topological order:

3 2 1 4

3 2 4 1

2 3 4 1

2 3 1 4

The following propositions are easy to prove:

- if the graph G is acyclic, then there exists at least one topological order of G
- if S is a topological order of G , and there exists a path from node n_1 to node n_2 in G , then n_1 is before n_2 in S



Exercise 3

Prove the above propositions, i.e.,

1. If the graph G is acyclic, then there exists at least one topological order of G
2. If S is a topological order of G , and there exists a path from node n_1 to node n_2 in G , then n_1 is before n_2 in S



Proof of the conflict-serializability theorem

(\Leftarrow) We have to show that if S is conflict-serializable, then the precedence graph $P(S)$ is acyclic. If S is conflict-serializable, then there exists a serial schedule S' on the same transactions that is conflict-equivalent to S . Since S' is serial, the precedence graph $P(S')$ associated to S' is acyclic. But for the preliminary lemma, since S is conflict-equivalent to S' , we have that $P(S)=P(S')$, and therefore $P(S)$ is acyclic.

(\Rightarrow) Let S be defined on the transactions T_1, \dots, T_n , and suppose that $P(S)$ is acyclic. Then there exists at least one topological order of $P(S)$, i.e., a sequence of its nodes such that if $T_i \rightarrow T_j$ is in $P(S)$, then T_i appears before T_j in the sequence. To such a topological order of $P(S)$, it corresponds the serial schedule S' on T_1, \dots, T_n such that, if $T_i \rightarrow T_j$ is in the graph, then all actions of T_i appear immediately before T_j in S' . It is easy to see that such a schedule S' is conflict-equivalent to S . Indeed, if S' is not conflict-equivalent to S , then there is a pair of conflicting actions a_h e b_k such that $(a_h <_{S'} b_k)$ and $(b_k <_S a_h)$. But $(b_k <_S a_h)$ means that the path $T_k \rightarrow T_h$ is in the graph $P(S)$, and therefore (see Exercise 3.2) T_k appears before T_h in every topological order of $P(S)$. However, $(a_h <_{S'} b_k)$ means that T_h appears before T_k in S' , and this contradicts the fact that S' corresponds to a topological order of $P(S)$.



Algorithm for conflict-serializability

The above theorem allows us to derive the following algorithm for checking whether a given schedule S is conflict-serializable:

- build the precedence graph $P(S)$ corresponding to S
- check whether $P(S)$ is acyclic or not
- return true if $P(S)$ is acyclic, false otherwise

It is immediate to verify that the time complexity of the algorithm is polynomial with respect to the size of the schedule S



Exercise 4

Check whether the following schedule is conflict-serializable

$w_1(x) \ r_2(x) \ w_1(z) \ r_2(z) \ r_3(x) \ r_4(z) \ w_4(z) \ w_2(x)$



Comparison with view-serializability

The main property to understand for comparing conflict-serializability and view-serializability is the following:

Theorem Let $S1$ and $S2$ be two schedules on the same transactions. If $S1$ and $S2$ are conflict-equivalent, then they are view-equivalent.

On the basis of this theorem, one can easily show the following:

Theorem If S is conflict-serializable, then it is view-serializable too.



Exercise 5

Prove the two theorems above.



Comparison with view-serializability

We have observed that every conflict-serializable schedule is also view-serializable.

It is important to note, however, that the converse does not hold. Indeed, there are schedules that are view-serializable and **not** conflict-serializable.

For example,

$$r1(x) \ w2(x) \ w1(x) \ w3(x)$$

is **view-serializable**, but not conflict-serializable



Comparison with view-serializability

Contrary to view-serializability, the class of conflict-serializable schedules is monotone.

Theorem The following two properties hold:

1. The class of conflict-serializable schedules is monotone.
2. S is in the class of conflict-serializable schedules if and only if for all $T \subseteq \text{Trans}(S)$, $\Pi_T(S)$ is in the class of view-serializable schedules.

What the above theorem says is that the class of conflict-serializable schedules is the largest monotone subclass of the class of view-serializable schedules.



Exercise 6

Consider the following schedule

$$w_1(y) \ w_2(y) \ w_2(x) \ w_1(x) \ w_3(x)$$

and

- check whether it is view-serializable or not,
- check whether it is conflict-serializable or not.



Order preserving conflict serializability

Let CSR denote the class of conflict serializable schedules.

Definition (Order Preservation)

A schedule S is **order preserving conflict serializable** if it is conflict equivalent to a serial schedule S' and for all $t, t' \in \text{trans}(S)$: if t completely precedes t' in S , then the same holds in S' . OCSR denotes the class of all schedules with this property.

Theorem

$\text{OCSR} \subset \text{CSR}$.

Example

$s = w_1(x) \ r_2(x) \ c_2 \ w_3(y) \ c_3 \ w_1(y) \ c_1 \quad \rightarrow \in \text{CSR}$
 $\rightarrow \notin \text{OCSR}$



Commit-order preserving conflict serializability

Definition (Commit Order Preservation)

A schedule S is **commit order preserving conflict serializable** if for all $t_i, t_j \in \text{tran}(S)$: if there are conflicting actions $p \in t_i, q \in t_j$ in S such that p precedes q in S , then c_i precedes c_j in S . COCSR denotes the class of schedules with this property.

Theorem

$\text{COCSR} \subset \text{CSR}$.

Theorem

A schedule S is in COCSR iff there is a serial schedule S' conflict equivalent to S such that for all $t_i, t_j \in \text{tran}(S)$: t_i precedes t_j in S' if and only if c_i precedes c_j in S .

Theorem

$\text{COCSR} \subset \text{OCSR}$.

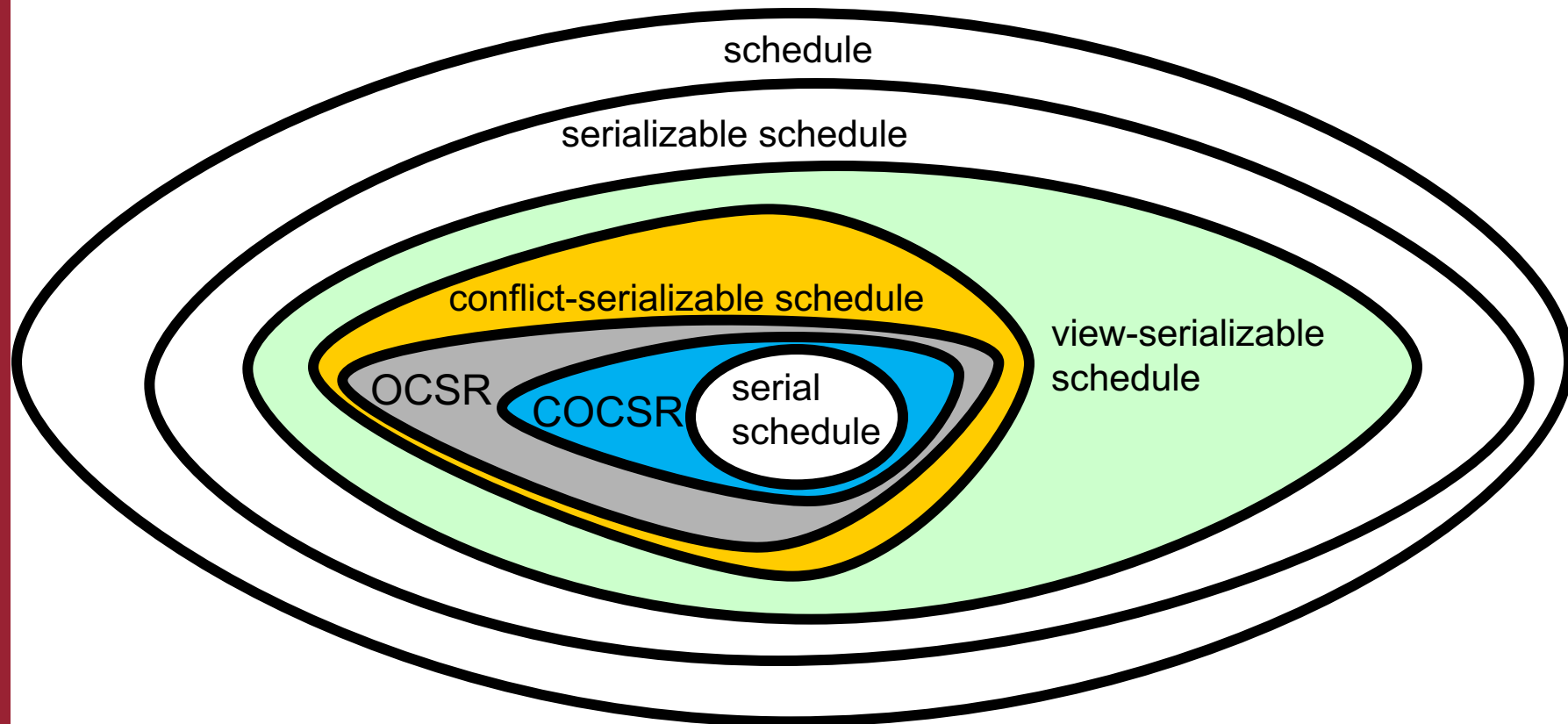
Example:

$s = w_3(y) \ c_3 \ w_1(x) \ r_2(x) \ c_2 \ w_1(y) \ c_1 \quad \rightarrow \in \text{OCSR}$
 $\quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \quad \rightarrow \notin \text{COCSR}$



View-serializability and conflict-serializability

The relationship between view-serializability and conflict-serializability (and its variants) can be visualized as follows:





Scheduler based on conflict-serializability

A scheduler based on conflict-serializability

- receives the sequence S of actions of the active transactions, in an interleaved order (such order depends on factors which are independent on the scheduler)
- manages the precedence graph associated to the sequence S
- once a new action is added to S , it updates the precedence graph of the current schedule (that is not necessarily complete), and
 - if a cycle appears in the graph, it kills the transaction where the action that has introduced the cycle appears (killing a transaction is a complex process)
 - otherwise, it accepts the action, and continues

Since maintaining the precedence graph can be very costly (the size of the graph can have thousands of nodes), **the notion of conflict-serializability is not used** in commercial systems.

However, contrary to view-serializability, conflict-serializability is used in some sophisticated applications where concurrency control has to be taken care of by a specialized module



5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 **Concurrency control through locks**

5.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



Concurrency control through locks

- We observed that view-serializability and conflict-serializability are not used in commercial systems
- We will now study a method for concurrency control that is used in commercial systems. Such method is based on the use of lock
- In the methods based on locks, a transaction must ask and get a permission in order to operate on an element of the database. The lock is a mechanism for a transaction to ask and get such a permission



Primitives for exclusive lock

- For the moment, we will consider exclusive locks. Later on, we will take into account more general types of locks
- We introduce two new operations (besides read and write) that can appear in schedules. Such operations are used to request and release the exclusive use of a resource (element A in the database):
 - **Lock** (exclusive): $l_i(A)$
 - **Unlock**: $u_i(A)$
- The lock operation $l_i(A)$ means that transaction T_i requests the exclusive use of element A of the database
- The unlock operation $u_i(A)$ means that transaction T_i releases the lock on A , i.e., it renounces the use of A



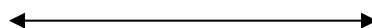
Well-formed transactions and legal schedules

When using exclusive locks, transactions and schedules should obey two rules:

- **Rule 1:** Every transaction is well-formed. A **transaction T_i is well-formed** if every action $p_i(A)$ (a read or a write on A) of T_i is contained in a “critical section”, i.e., in a sequence of actions delimited by a pair of lock-unlock on A :

$$T_i: \dots l_i(A) \dots p_i(A) \dots u_i(A) \dots$$

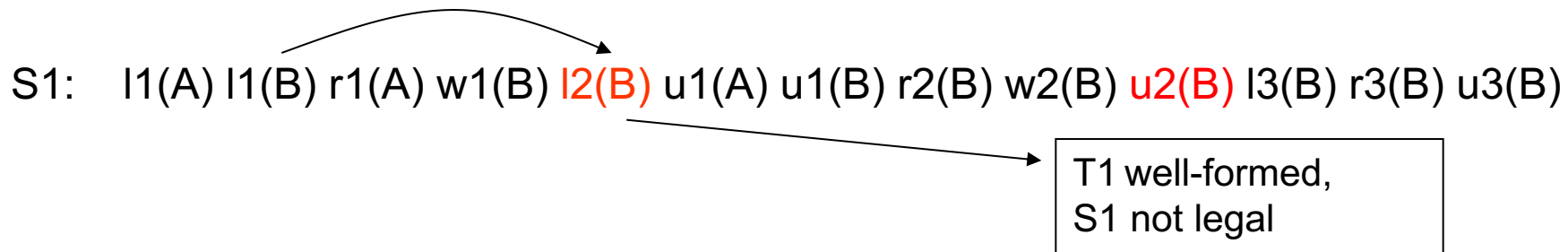
- **Rule 2:** The schedule is legal. A **schedule S with locks is legal** if no transaction in it locks an element A when a different transaction has granted the lock on A and has not yet unlocked A

$$S: \dots l_i(A) \dots u_i(A) \dots$$


no $l_j(A)$



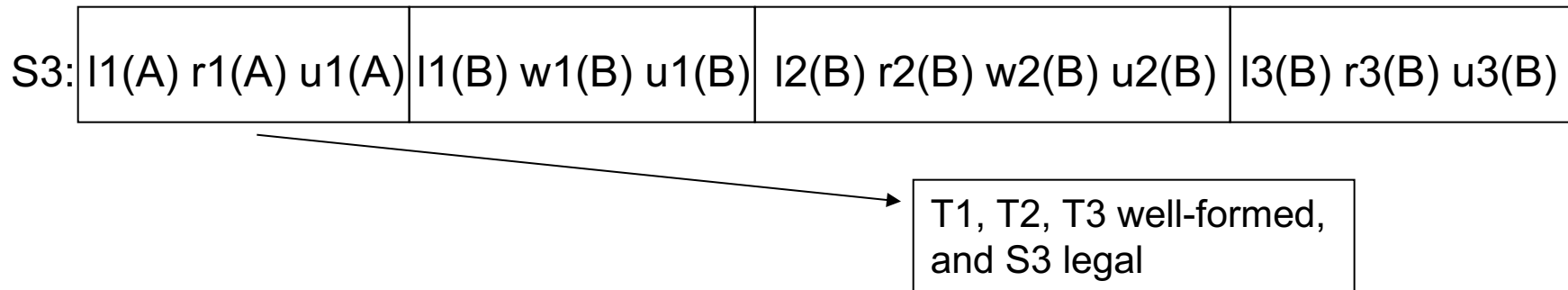
Schedule with exclusive locks: examples



S2: l1(A) r1(A) w1(B) u1(A) u1(B) l2(B) r2(B) w2(B) l3(B) r3(B) u3(B)

T1 ill-formed:
write without lock.
T2 ill-formed:
lock without unlock.

S2 not legal





Scheduler based on exclusive locks

A scheduler based on exclusive locks behaves as follows:

1. When an action request is issued by a transaction, the scheduler checks whether this request makes the transaction ill-formed, in which case the transaction is aborted by the scheduler.
2. When a lock request on A is issued by transaction T_i , while another transaction T_j has a lock on A, the scheduler does not grant the request (otherwise the schedule would become illegal), and T_i is blocked until T_j releases the lock on A.
3. To trace all the locks granted, the scheduler manages a table of locks, called **lock table**

In other words, the scheduler ensures that the current schedule is **legal** and all its transactions are **well-formed**.



Example of scheduler behaviour

T1	T2
l1(A); r1(A)	
A:=A+100; w1(A);	l2(A) - blocked!
l1(B); r1(B); u1(A);	l2(A) - re-started!
	r2(A)
	A:=Ax2; w2(A); u2(A)
B:=B+100; w1(B); u1(B)	l2(B); r2(B)
	B:=Bx2; w2(B); u2(B)



Is this sufficient for serializability?

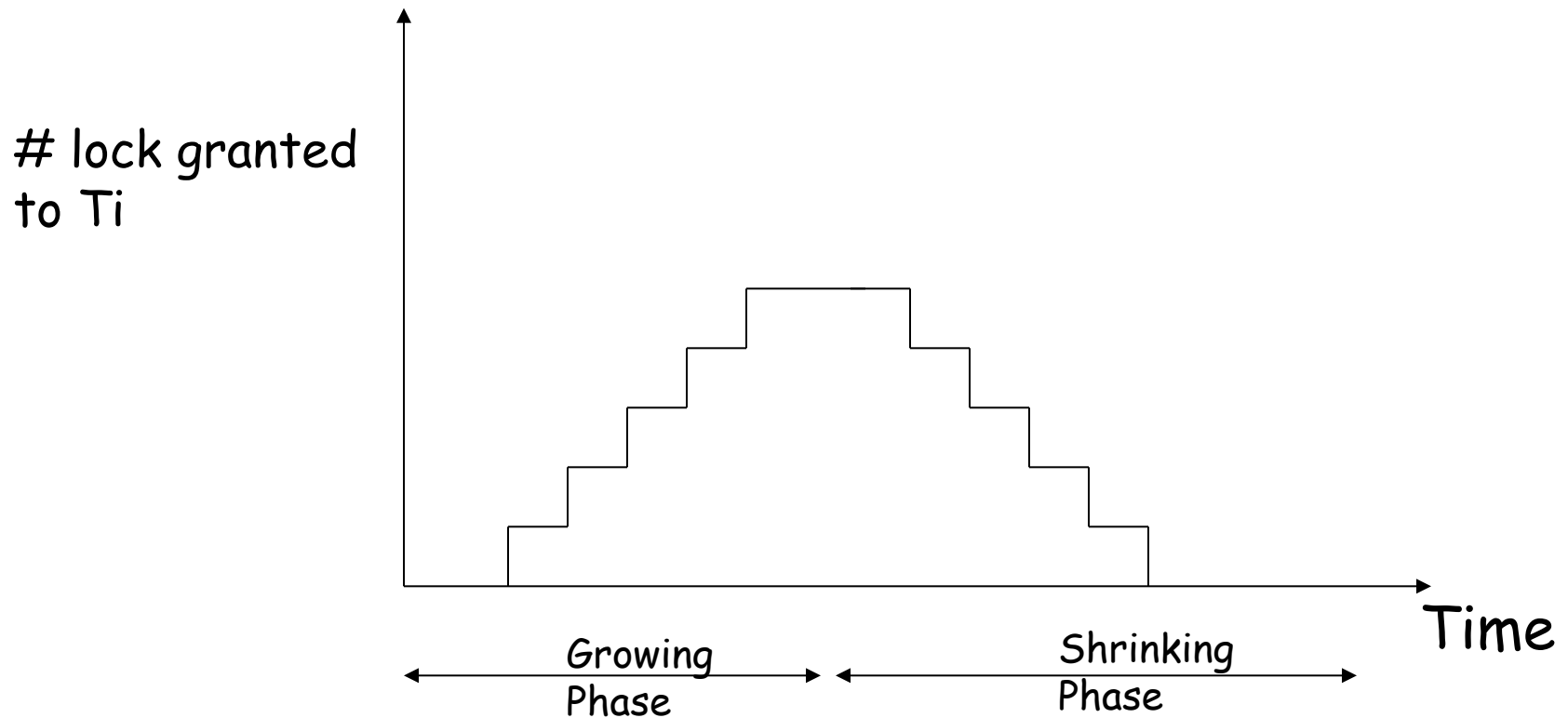
T1	T2	A 25	B 25
l1(A); r1(A) A:=A+100; w1(A); u1(A)	l2(A); r2(A) A:=A×2; w2(A); u2(A)	125	
	l2(B); r2(B) B:=B×2; w2(B); u2(B)	250	50
l1(B); r1(B) B:=B+100; w1(B); u1(B)			150
		250	150

Ghost update: isolation is not ensured by the use of locks



The two phases of Two-Phase Locking

Locking and unlocking scheme in a transaction following the 2PL protocol



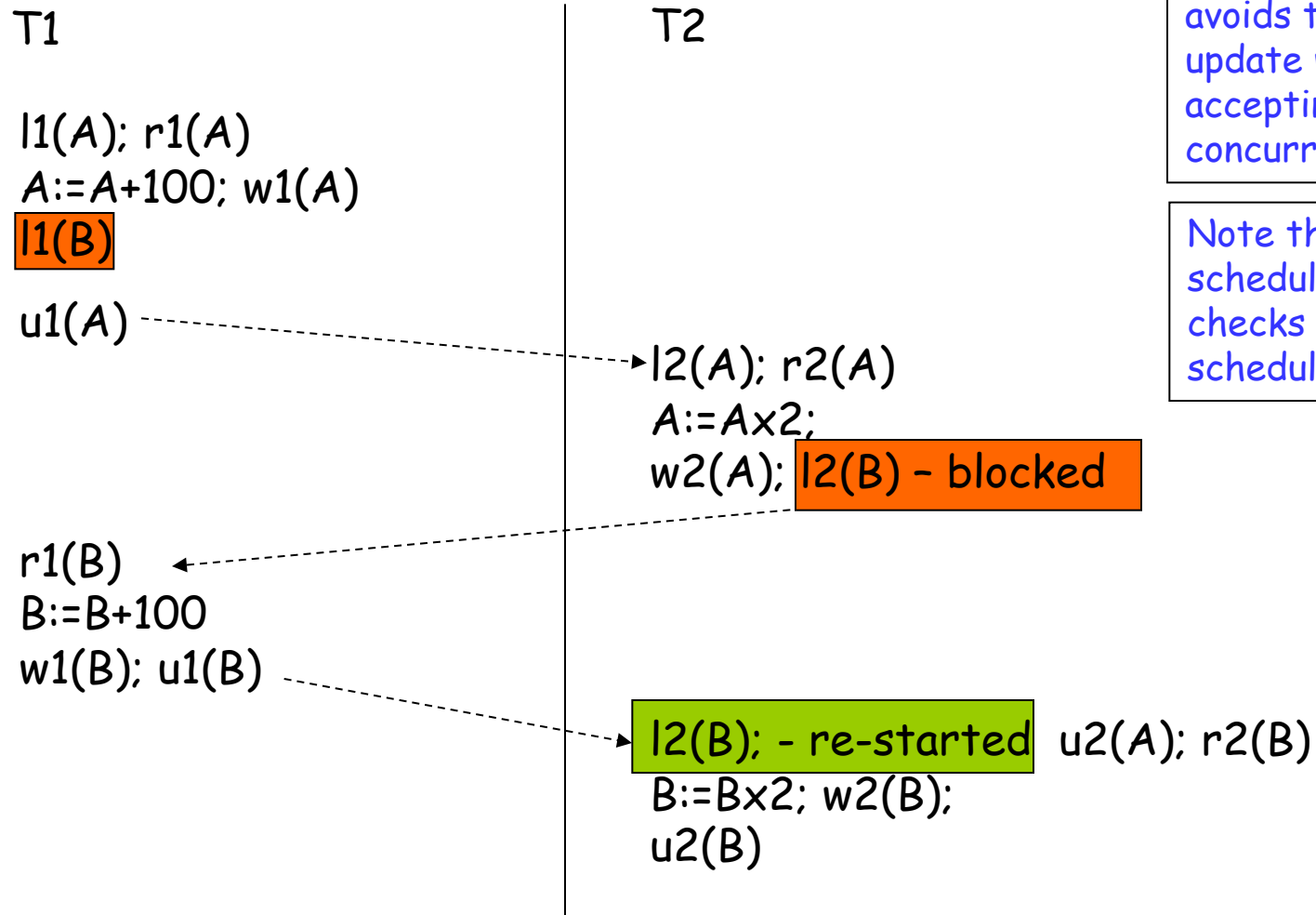


Example of a 2PL schedule

T1	T2	A 25	B 25
$l_1(A); r_1(A)$ $A := A + 100; w_1(A); u_1(A)$		125	
	$l_2(A); r_2(A)$ $A := A \times 2; w_2(A); u_2(A)$	250	
	$l_2(B); r_2(B)$ $B := B \times 2; w_2(B); u_2(B)$		50
$l_1(B); r_1(B)$ $B := B + 100; w_1(B); u_1(B)$			150
		250	150



How the scheduler works in the 2PL protocol



The 2PL protocol avoids the ghost update while accepting concurrency

Note that the scheduler still checks that the schedule is legal



The risk of deadlock

T1	T2
$l1(A); r1(A)$	
$A := A + 100;$	$l2(B); r2(B)$
	$B := B \times 2$
$w1(A)$	$w2(B)$
$l1(B)$ - blocked	$l2(A)$ - blocked

S: $l1(A) \ r1(A) \ l2(B) \ r2(B) \ w1(A) \ w2(B) \ l1(B) \ l2(A)$

To ensure that the schedule is legal, the scheduler blocks both T1 and T2, and none of the two transactions can proceed. This is a **deadlock** (we will come back to the methods for deadlock management).



Who issues the lock/unlock commands?

So far, we have assumed that transactions issue the lock/unlock commands. However, this is not necessary.

Indeed, we can design a scheduler in such a way that **it inserts the lock/unlock commands** while respecting the following conditions:

- Every transaction is well-formed
- The schedule is legal (if at all possible)
- Each transaction, extended with the inserted lock/unlock commands, follows the 2PL protocol

For this reason, even in the presence of locks, we will continue to denote a schedule by means of a sequence of read/write/commit commands. For example, the schedule

$I_1(A) \ r_1(A) \ I_1(B) \ u_1(A) \ I_2(A) \ w_2(A) \ r_1(B) \ w_1(B) \ u_1(B) \ I_2(B) \ u_2(A) \ r_2(B) \ w_2(B) \ u_2(B)$

can be denoted as:

$r_1(A) \ w_2(A) \ r_1(B) \ w_1(B) \ r_2(B) \ w_2(B)$



Scheduler based on exclusive locks and 2PL

We study how a scheduler based on exclusive locks and 2PL behaves during the analysis of the current schedule (obviously, not necessarily complete):

1. If a request by transaction T_i shows that T_i is not well-formed, then T_i is aborted by the scheduler
2. If a lock request by transaction T_i shows that T_i does not follow the 2PL protocol, then T_i is aborted by the scheduler
3. If a lock is requested for A by transaction T_i while A is used by a different transaction T_j , then the scheduler blocks T_i , until T_j releases the lock on A . If the scheduler figures out that a deadlock has occurred (or will occur), then the scheduler adopts a method for deadlock management
4. To trace all the locks granted, the scheduler manages a table of locks, called **lock table**

Note that (1) and (2) do not occur if the lock/unlock commands are automatically inserted by the scheduler.

Simply put, the above behaviour means that the scheduler ensures that

1. the current schedule is **legal**
2. all its transactions are **well-formed**
3. all its transactions follow the **2PL protocol**



2PL and conflict-serializability

To compare 2PL and conflict-serializability, we make use of the above observation, and note that every schedule that includes lock/unlock operations can be seen as a “traditional” schedule (by simply ignoring such operations)

Theorem Every legal schedule constituted by well-formed transactions following the 2PL protocol (with exclusive locks) is conflict-serializable.

Proof Let S be a legal schedule constituted by well-formed transactions following the 2PL protocol (with exclusive locks). To show that S is conflict-serializable, we proceed by induction on the number N of transactions in S .

Base step: If $N=1$, S is serial, and therefore is trivially conflict-serializable.



Proof continued

Inductive step: Suppose that S is defined on T_1, \dots, T_N ($N > 1$), and let T_i the first transaction that executes an unlock operation, say $ui(X)$, in S . We now show that we can move all operations of T_i in front of S , without swapping any pair of conflicting actions. We consider an action $w_i(Y)$ in T_i (analogous observation holds if we considered $ri(Y)$ instead of $w_i(Y)$), and we show that it cannot be preceded by any conflicting action in S . Indeed, suppose that there is a conflicting action $w_j(Y)$ in S preceding $w_i(Y)$ with j different from i :

... $w_j(Y)$... $uj(Y)$... $li(Y)$... $w_i(Y)$...

Since T_i is the first transaction that executes an unlock operation $ui(X)$ in S , we either have

... $ui(X)$... $w_j(Y)$... $uj(Y)$... $li(Y)$... $w_i(Y)$...

or

... $w_j(Y)$... $ui(X)$... $uj(Y)$... $li(Y)$... $w_i(Y)$...

In both cases, $ui(X)$ would appear before $li(Y)$ in S , and therefore T_i would not follow the 2PL protocol. We can then conclude that, by moving all actions of T_i in front of S , we get a schedule S'' that is conflict-equivalent to S , of the form

(actions of T_i) (remaining actions of S)

The part denoted by $S' =$ (remaining actions of S) is a legal schedule on $(N-1)$ transactions constituted by well-formed transactions following the 2PL protocol (with exclusive locks). For the inductive hypothesis, S' is conflict-serializable, which means that there is a serial schedule S''' on the $(N-1)$ transactions that is conflict equivalent to S' . Now, consider the schedule constituted by T_i followed by S''' : such a schedule is obviously conflict equivalent to S , which implies that S is conflict-serializable.



What does the theorem intuitively say

The theorem says that any legal schedule constituted by N well-formed transactions following the 2PL protocol (with exclusive locks) is conflict-equivalent to the serial schedule that orders the transactions according to the following rule:

1. Take as first transaction the one that executes the first unlock operation in S
2. Take as second transaction the one that executes the first unlock operation among the remaining $(N-1)$ transactions in S
3.
- $N-1$. Take as $(N-1)$ -th transaction the one that executes the first unlock operation among the remaining 2 transactions in S
- N . Take the last transaction as the N -th transaction



Another intuition

Suppose that S follows the 2PL protocol. If $T1$ and $T2$ in S contain two conflicting actions on element A , then there is an action $a1$ on A which is conflicting with an action $b2$ on the same A . Suppose that $a1$ precedes $b2$ in S . In order for S to be well-formed, $T1$ must get the lock on A , and execute the unlock on A before $b2$:

$l1(A) \dots a1(A) \dots u1(A) \dots l2(A) \dots b2(A)$

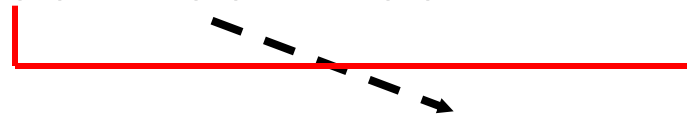
Now suppose that S is not conflict-serializable, because there is an action $c2$ on B before the conflicting action $d1$ on B in S . This would imply that $T2$ has got the lock on B and has unlocked B before $d1$:

$l2(B) \dots c2(B) \dots u2(B) \dots l1(B) \dots d1(B)$

But, since S follows the 2PL protocol, $u2(B)$ cannot appear before $l2(A)$, and this implies that $u1(A)$ appears before $l1(B)$, and therefore S does not follow the 2PL protocol:

$l1(A) \dots a1(A) \dots u1(A) \dots l2(A) \dots b2(A)$

$l2(B) \dots c2(B) \dots u2(B) \dots l1(B) \dots d1(B)$



In other words, 2PL ensures that, if $T1$ wins against $T2$ (i.e., it gets the lock on a competing element before $T2$), then $T1$ wins against $T2$ for any conflict on any other element. This is sufficient for conflict-serializability, because I can always assume that S is equivalent to any serial order compatible with the “win” relation.



Comment and exercise

We denote by “2PL schedule with exclusive locks” the class of legal schedules with exclusive locks constituted by well-formed transactions following the 2PL protocol.

1. Consider the following schedule S1:

$r1(x) \ w2(x) \ w3(y) \ w1(z) \ w3(z)$

and tell whether it is in the class of “2PL schedule with exclusive locks”, i.e., whether we can insert lock and unlock commands into S1 in such a way that the resulting sequence of actions is a 2PL schedule with exclusive locks.

2. Consider the following schedule S2:

$r1(x) \ w2(x) \ w3(y) \ w3(z) \ w1(z)$

and tell whether it is in the class of “2PL schedule with exclusive locks”, i.e., whether we can insert lock and unlock commands into S2 in such a way that the resulting sequence of actions is a 2PL schedule with exclusive locks.



2PL and conflict-serializability

Theorem There exists a conflict-serializable schedule that does not follow the 2PL protocol (with exclusive locks).

Proof It is sufficient to consider the following schedule S:

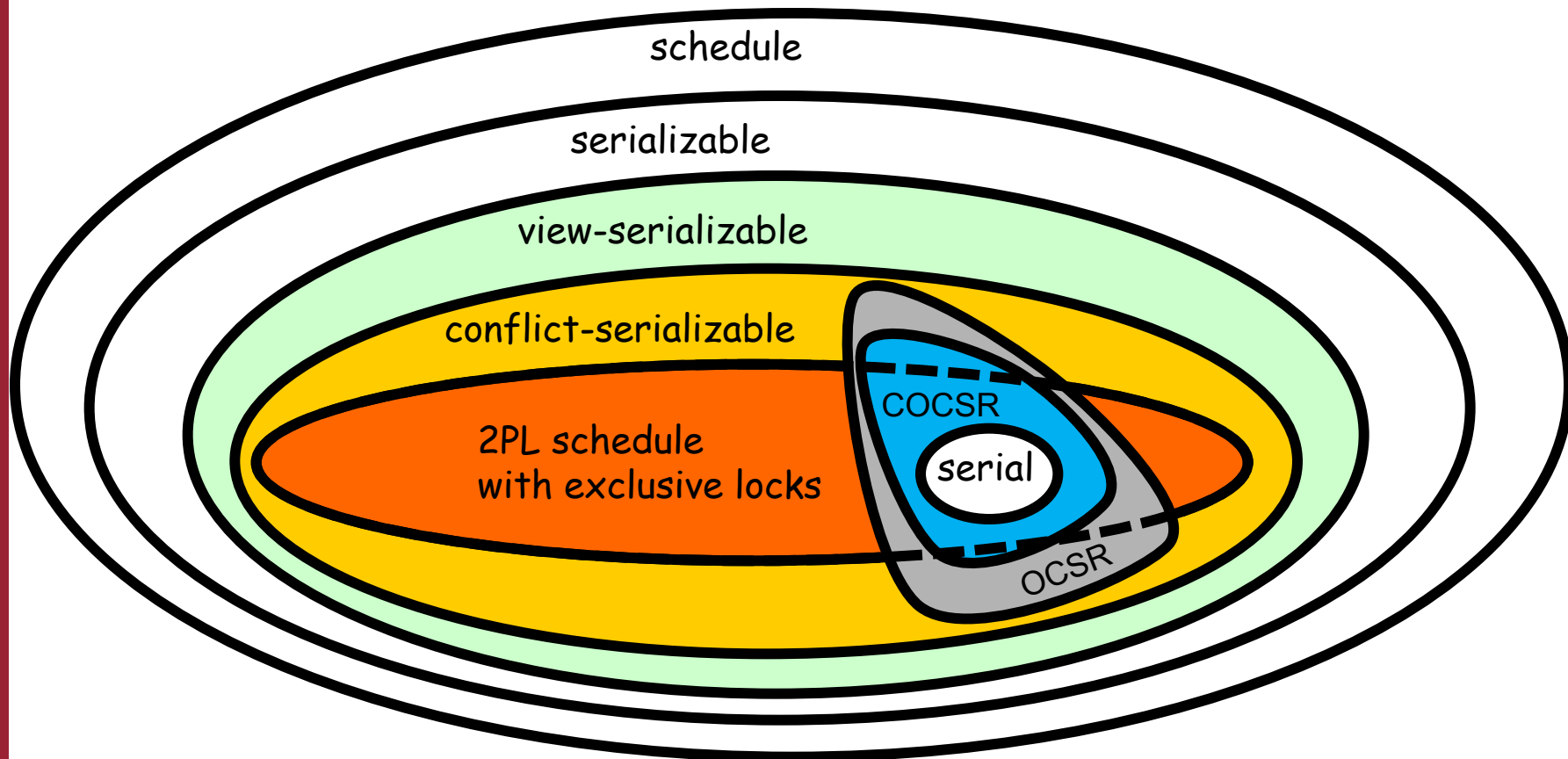
$r1(x) \ w1(x) \ r2(x) \ r3(y) \ w1(y)$

S is obviously conflict-serializable (the serial schedule $T3, T1, T2$ is conflict-equivalent to S), but it is easy to show that we cannot insert in S the lock/unlock commands in such a way that all transactions are well-formed and follow the 2PL protocol, and the resulting schedule is legal. Indeed, it suffices to notice that we should insert in S the command $u1(x)$ before $r2(x)$, because in order for T2 to read x it must hold the exclusive lock on x, and we should insert in S the command $l1(y)$ after $r3(y)$, because in order for T3 to read y it must hold the exclusive lock on y, and therefore, the command $l1(y)$, which is necessary for executing $w3(y)$, cannot be issued before $r3(y)$. It follows that we cannot insert into S the lock/unlock commands in such a way that the 2PL protocol is respected.



2PL and conflict-serializability

We denote by “2PL schedule with exclusive locks” the class of legal schedules with exclusive locks constituted by well-formed transactions following the 2PL protocol. Graphically, the relationship between conflict-serializability and 2PL can be represented as follows:





Shared locks

With exclusive locks, a transaction reading A must unlock A before another transaction can read the same element A:

S: ... l1(A) r1(A) u1(A) ... l2(A) r2(A) u2(A) ...

Actually, this looks too restrictive, because the two read operations do not create any conflict. To remedy this situation, we introduce a new type of lock: the **shared lock**. We denote by sli(A) the command for the transaction Ti to ask for a shared lock on A.

With the use of shared locks, the above example changes as follows:

S: ... sl1(A) r1(A) sl2(A) r2(A) u1(A) u2(A)

The primitive for locks are now as follows:

xli(A): exclusive lock (also called write lock)

sli(A): shared lock (also called read lock)

ui(A): unlock



Well-formed transactions with shared locks

With shared and exclusive locks, the following rule must be respected.

Rule 1: We say that a **transaction T_i is well-formed** if

- every read $ri(A)$ is preceded either by $sli(A)$ or by $xli(A)$, with no $ui(A)$ in between,
- every $wi(A)$ is preceded by $xli(A)$ with no $ui(A)$ in between,
- every lock (sl or xl) on A by T_i is followed by an unlock on A by T_i .

Note that we allow T_i to first execute $sli(A)$, probably for reading A , and then to execute $xli(A)$, probably for writing A without the unlock of A by means of T . The transition from a shared lock on A by T to an exclusive lock on the same element A by T (without an unlock on A by T) is called “**lock upgrade**”.



Legal schedule with shared locks

With shared and exclusive locks, the following rule must also be respected.

Rule 2: We say that a **schedule S is legal** if

- an $xli(A)$ is not followed by any $xlj(A)$ or by any $slj(A)$ (with j different from i) without an $ui(A)$ in between
- an $sli(A)$ is not followed by any $xlj(A)$ (with j different from i) without an $ui(A)$ in between



Two-phase locking (with shared locks)

With shared locks, the two-phase locking rule becomes:

Definition of two-phase locking (with exclusive and shared locks): A schedule S (with shared and exclusive locks) follows the **2PL protocol** if in every transaction T_i of S , all lock operations (either for exclusive or for shared locks) precede all unlocking operations of T_i .

In other words, no action $sli(X)$ or $xli(X)$ can be preceded by an operation of type $ui(Y)$ in the schedule.



How locks are managed

- The scheduler uses the so-called “compatibility matrix” (see below) for deciding whether a lock request should be granted or not.
- In the matrix, “S” stands for shared lock, “X” stands for exclusive lock, “yes” stands for “requested granted” and “no” stands for “requested not granted”

		New lock requested by $T_j \neq T_i$ on A	
		S	X
Lock already granted to T_i on A	S	yes	no
	X	no	no



How locks are managed

- The problem for the scheduler of automatically inserting the lock/unlock commands becomes more complex in the presence of shared locks.
- Also, the execution of the unlock commands requires more work. Indeed, when an unlock command on A is issued by T_i , there may be several transactions waiting for a lock (either shared or exclusive) on A, and the scheduler must decide to which transaction to grant the lock. Several methods are possible:
 - First-come-first-served
 - Give priorities to the transactions asking for a shared lock
 - Give priorities to the transactions asking for a lock upgrade

The first method is the most used one, because it avoids “**starvation**”, i.e., the situation where a request of a transaction is never granted.



Exercise 7

Consider the following schedule S:

r1(A) r2(A) r2(B) w1(A) w2(D) r3(C) r1(C) w3(B) c2 r4(A) c1 c4 c3

and tell whether S is in the class of 2PL schedules with shared and exclusive locks



Exercise 7: solution

The schedule S:

$r1(A) \ r2(A) \ r2(B) \ w1(A) \ w2(D) \ r3(C) \ r1(C) \ w3(B) \ c2 \ r4(A) \ c1 \ c4 \ c3$

is in the class of 2PL schedules with shared and exclusive locks. This can be shown as follows:

$sl1(A) \ r1(A) \ sl2(A) \ r2(A) \ sl2(B) \ r2(B) \ xl2(D) \ u2(A) \ xl1(A) \ w1(A) \ w2(D) \ sl3(C) \ r3(C) \ sl1(C) \ r1(C) \ u1(C) \ u1(A) \ u2(B) \ u2(D) \ xl3(B) \ w3(B) \ u3(B) \ u3(C) \ c2 \ sl4(A) \ r4(A) \ u4(A) \ c1 \ c4 \ c3$



Properties of two-phase locking (with shared locks)

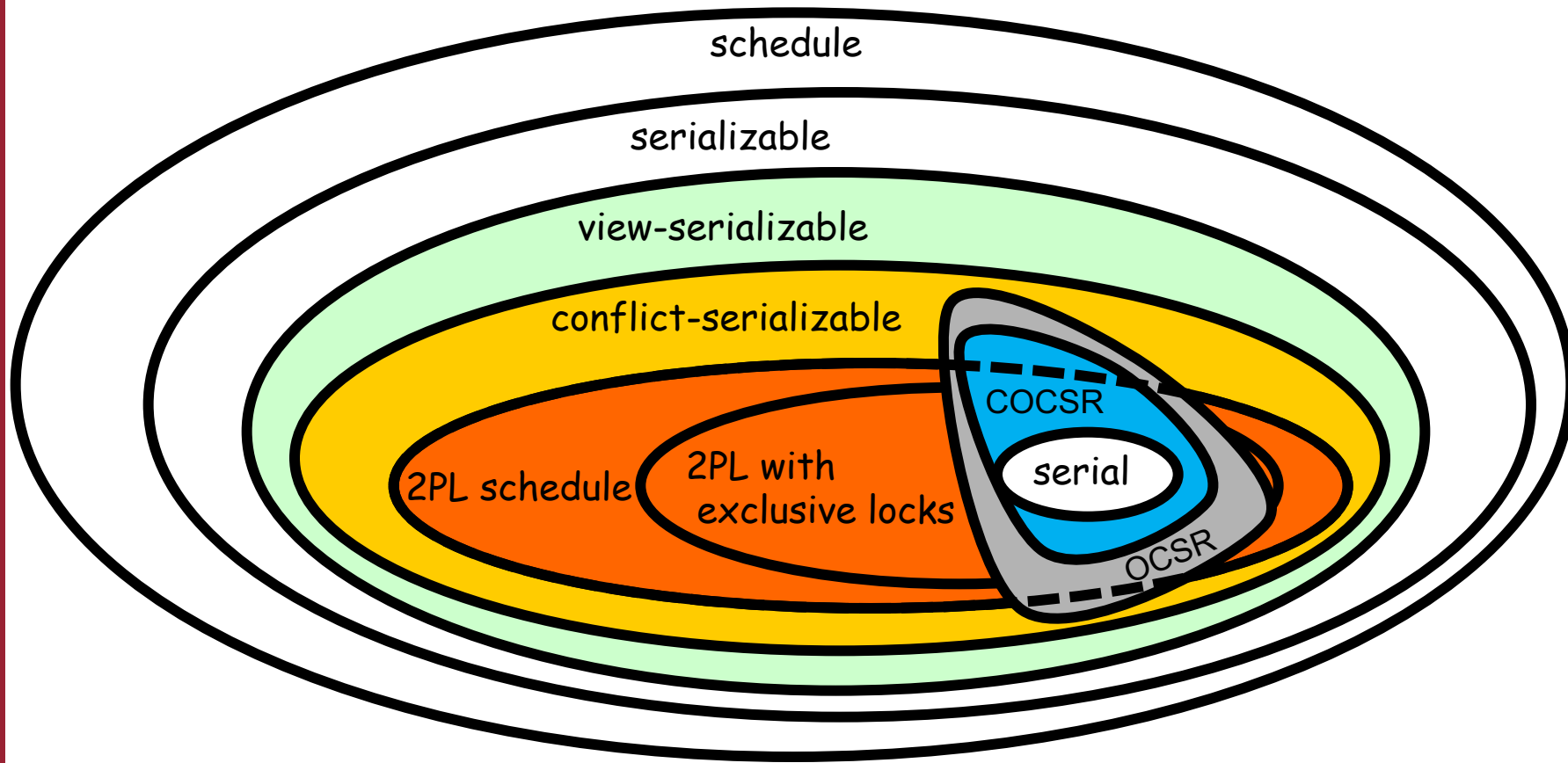
The properties of two-phase locking with shared and exclusive locks are similar to the case of exclusive locks only:

- **Theorem** Every legal schedule with well-formed transactions following the two-phase locking protocol (with exclusive and shared locks) is conflict-serializable.
- **Theorem** There exists a conflict-serializable schedule that does not follow the 2PL protocol (with exclusive and shared locks).
- With shared locks, the risk of deadlock is still present, like in:
$$sl1(A) \quad sl2(A) \quad xl1(A) \quad xl2(A)$$



2PL and conflict-serializability

We denote by “2PL schedule” the class of legal schedules with shared and exclusive locks constituted by well-formed transactions following the 2PL protocol.





Deadlock management

- We recall that the **deadlock** occurs when two transactions T1 and T2 have the use of two elements A and B, and each of them asks for an exclusive lock on the element of the other one, and therefore no one can proceed
- The probability of deadlock **grows linearly with the number of transactions** and **quadratically with the number of lock requests** in the transactions

T1	T2
$xl1(A); r1(A)$	
	$xl2(B); r2(B)$
$A := A + 100;$	$B := B \times 2$
$w1(A)$	$w2(B)$
$sl1(B)$ - blocked!	$sl2(A)$ - blocked!



Techniques for deadlock management

1. Timeout
2. Deadlock recognition and solution
3. Deadlock prevention



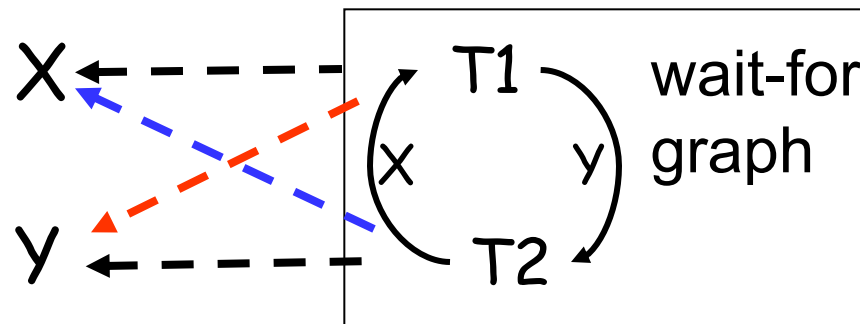
Timeout

- The system fixes a timeout t after which a transaction waiting for a lock is killed
- Advantages
 - very simple
- Disadvantages
 - if t is high, the risk is to be late in solving the problem
 - if t is low, too many transactions are killed
 - risk of **individual block** (same transactions killed several times)



Deadlock recognition

- A graph (**wait-for graph**) is incrementally maintained: the nodes are the transactions, and the edge from T_i to T_j means that T_i is waiting for T_j to release a lock
- When a cycle appears in the graph, the deadlock is solved by killing one of the involved transactions, for example the one that made the fewer operations (individual block is a risk)
- Example: $sl_1(X)$ $sl_2(Y)$ $r_1(X)$ $r_2(Y)$ $l_1(Y)$ $xl_2(X)$





Deadlock prevention: wait-die

To each transaction T_i a priority $pr(T_i)$ is assigned (for example, a number indicating how old is the transaction), in such a way that different transactions have different priorities

The following rule is applied: in case of conflict on a lock, T_i is allowed to wait for T_j only if T_i has greater priority, i.e., if $pr(T_i) > pr(T_j)$, otherwise T_i is killed.

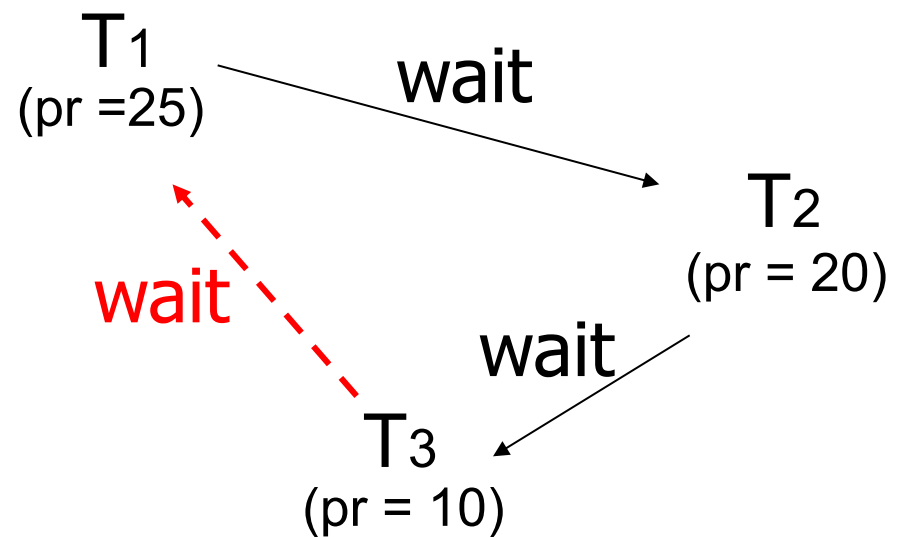
In practice, when a new edge $T_i \rightarrow T_j$ is added in the wait-for graph:

- if $pr(T_i) > pr(T_j)$: ok
- if $pr(T_i) \leq pr(T_j)$: T_i is killed



Example of wait-die

...
x|1(Y) T1 uses Y
x|3(X) T3 uses X
x|2(X) T2 waits for T3
x|1(X) T1 waits for T2
x|3(Y) T3 killed



T3 killed



Example of wait-die

...
x|3(X) T3 uses X
x|2(X) T2 waits for T3
x|1(X) T1 waits for T2 and for T3 ?

T₁
(pr =22)

T₂
(pr =25)

wait
T₃
(pr =20)

Note that $pr(T_1) > pr(T_3)$, and $pr(T_1) < pr(T_2)$. If we allow T₁ to wait for T₃, we have two options when T₃ releases the lock on X:

- (1) T₁ proceeds – in this case T₂ will wait for T₁, with the risk of starvation;
- (2) T₂ proceeds, and T₁ waits for T₂ – this violates the rule that only transactions with higher priorities wait.

So **the right choice is to kill T₁**.



5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 Concurrency control through locks

5.5 **Recoverability of transactions**

5.6 Concurrency control through timestamps

5.7 Concurrency control in SQL



The rollback problem

So far, we have carried out our study under the assumption that no transaction are rolled back. Now, we relax this strong assumption, and we study the problem of rollback.

The first observation is that, with rollbacks, the notion of serializability that we have considered up to now is not sufficient for achieving the ACID properties.

This fact is testified by the existence of a new anomaly, called “dirty read”.



A new anomaly: dirty read (WR anomaly)

Consider two transactions T1 and T2, both with the commands:

READ(A,x), $x := x + 1$, WRITE(A,x)

Now consider the following schedule (where T1 executes the rollback):

T ₁	T ₂
begin	begin
READ(A,x)	
$x := x + 1$	
WRITE(A,x)	
	READ(A,x)
	$x := x + 1$
rollback	
	WRITE(A,x)
	commit

The problem is that T2 reads a value written by T1 before T1 commits or rollbacks.

Therefore, T2 reads a “dirty” value, that is shown to be incorrect when the rollback of T1 is executed. The behavior of T2 depends on an incorrect input value.

This is another form of **WR (write-read) anomaly**.



Commit o rollback?

Recall that, at the end of transaction T_i :

- If T_i has executed the commit operation:
 - the system should ensure that the effects of the transactions are recorded permanently in the database
- If T_i has executed the rollback operation:
 - the system should ensure that the transaction has no effect on the database



Cascading rollback

Note that the rollback of a transaction T_i can trigger the rollback of other transactions, in a cascading mode. In particular:

- If a transaction T_j different from T_i has read from T_i , we should kill T_j (in other words, T_j should rollback)
- If another transaction T_h has read from T_j , T_h should in turn rollback
- and so on...

This situation, called **cascading rollback**, should be avoided, since it causes several performance problems.



Recoverable schedules

If in a schedule S , a transaction T_i that has read from T_j commits before T_j , the risk is that T_j then rollbacks, so that T_i leaves an effect on the database that depends on an operation (of T_j) that never existed. To capture this concept, we say that T_i is not recoverable.

A schedule S is **recoverable** if no transaction in S commits before all other transactions it has “read from”, commit.

Example of recoverable schedule:

$S: w_1(A) w_1(B) w_2(A) r_2(B) c_1 c_2$

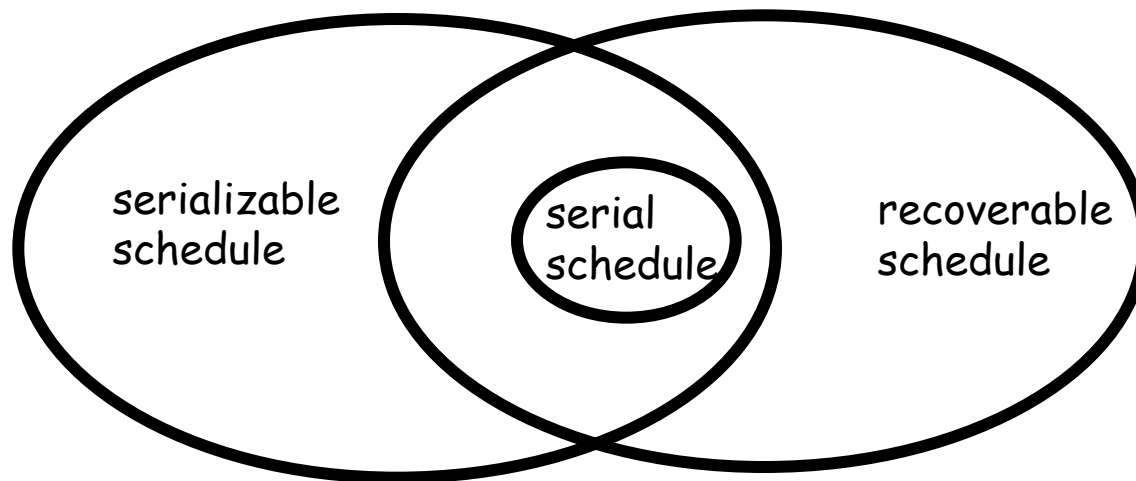
Example of non-recoverable schedule:

$S: w_1(A) w_1(B) w_2(A) r_2(B) r_3(A) c_1 c_3 c_2$



Serializability and recoverability

Serializability and recoverability are two orthogonal concepts: there are recoverable schedules that are non-serializable, and serializable schedules that are not recoverable. Obviously, every serial schedule is recoverable.



For example, the schedule

S1: w2(A) w1(B) w1(A) r2(B) c1 c2

is recoverable, but not serializable (it is not view-serializable), whereas the schedule

S2: w1(A) w1(B) w2(A) r2(B) c2 c1

is serializable (in particular, conflict-serializable), but not recoverable



Recoverability and cascading rollback

Recoverable schedules can still suffer from the cascading rollback problem (the correct situation can be recovered, but in order to recover it, we should kill several transactions).

For example, in this recoverable schedule

S: w2(A) w1(B) w1(A) r2(B)

if T1 rolls back, T2 must be killed.

To avoid cascading rollback, we need a stronger condition wrt recoverability: a schedule S **avoids cascading rollback** (i.e., the schedule is ACR, Avoid Cascading Rollback) if every transaction in S reads values that are written by transactions that have already committed.

For example, this schedule is ACR

S: w2(A) w1(B) w1(A) **c1** r2(B) c2

In other words, an ACR schedule blocks the dirty data anomaly.



Summing up

- S is **recoverable** if no transaction in S commits before the commit of all the transactions it has “read from”

Example:

w1(A) w1(B) w2(A) r2(B) c1 c2

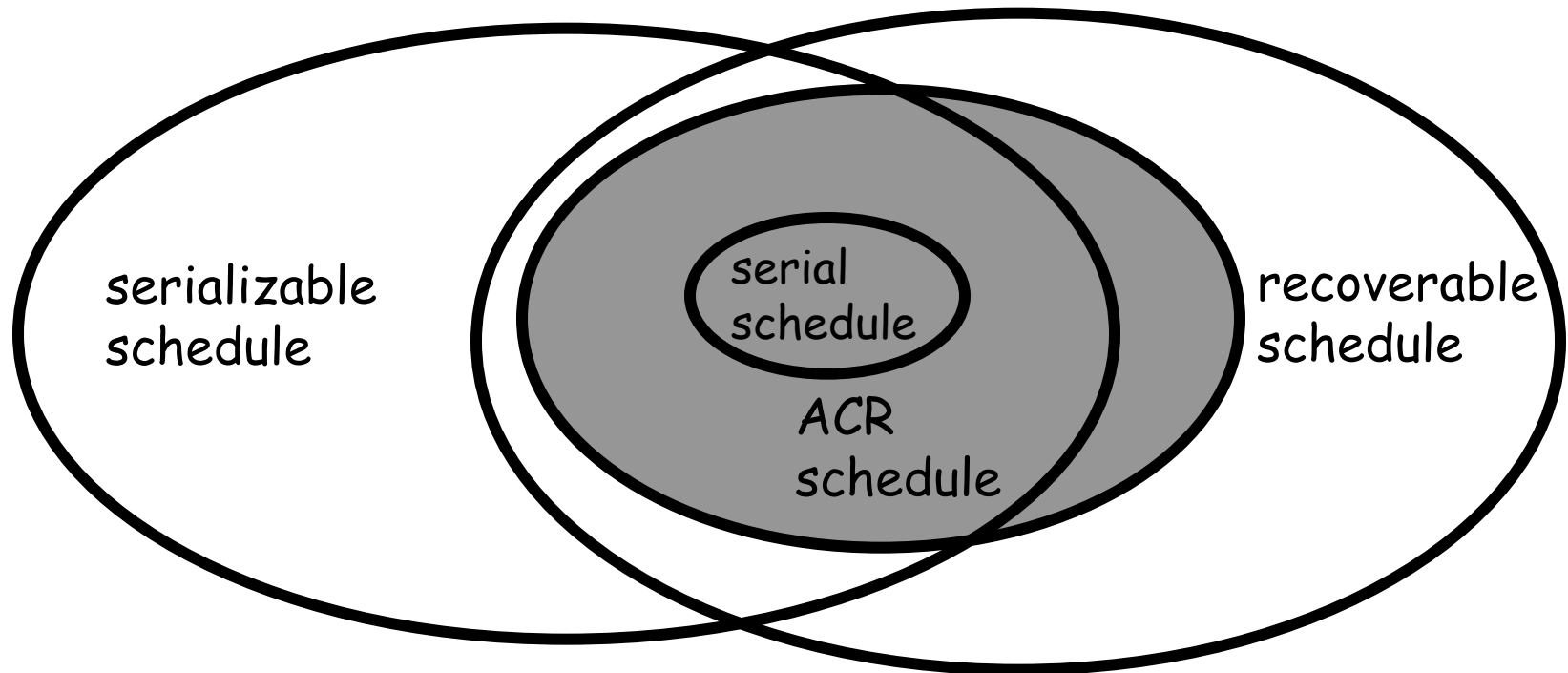
- S is ACR, i.e., **avoids cascading rollback**, if no transaction “reads from” a transaction that has not committed yet

Example:

w1(A) w1(B) w2(A) c1 r2(B) c2



Recoverability and ACR



Analogously to recoverable schedules, not all ACR schedules are serializable. Obviously, every ACR schedule is recoverable, and every serial schedule is ACR.

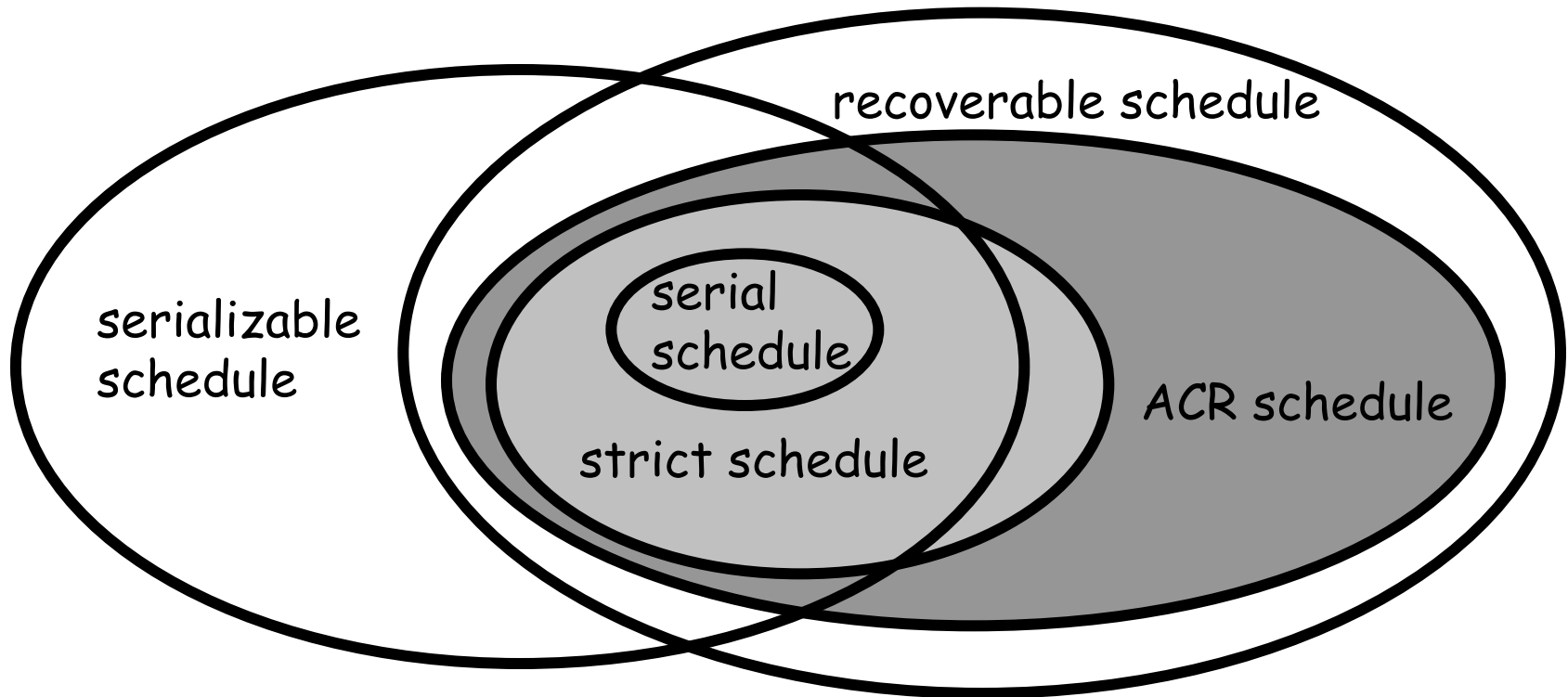


Strict schedules

- We say that, in a schedule S , a transaction T_i writes on T_j if there is a $w_j(A)$ in S followed by $w_i(A)$, and there is no write action on A in S between these two actions
- We say that a schedule S is **strict** if every transaction *reads* only values written by transactions that have already committed, and *writes* only on transactions that have already committed
- It is easy to verify that every strict schedule is ACR, and therefore recoverable
- Note that, for a strict schedule, when a transaction T_i rolls back, it is immediate to determine which are the values that have to be stored back in the database to reflect the rollback of T_i , because no transaction may have written on this values after T_i



Strict schedules and ACR



Obviously, every serial schedule is strict, and every strict schedule is ACR, and therefore recoverable. However, not all ACR schedules are strict.

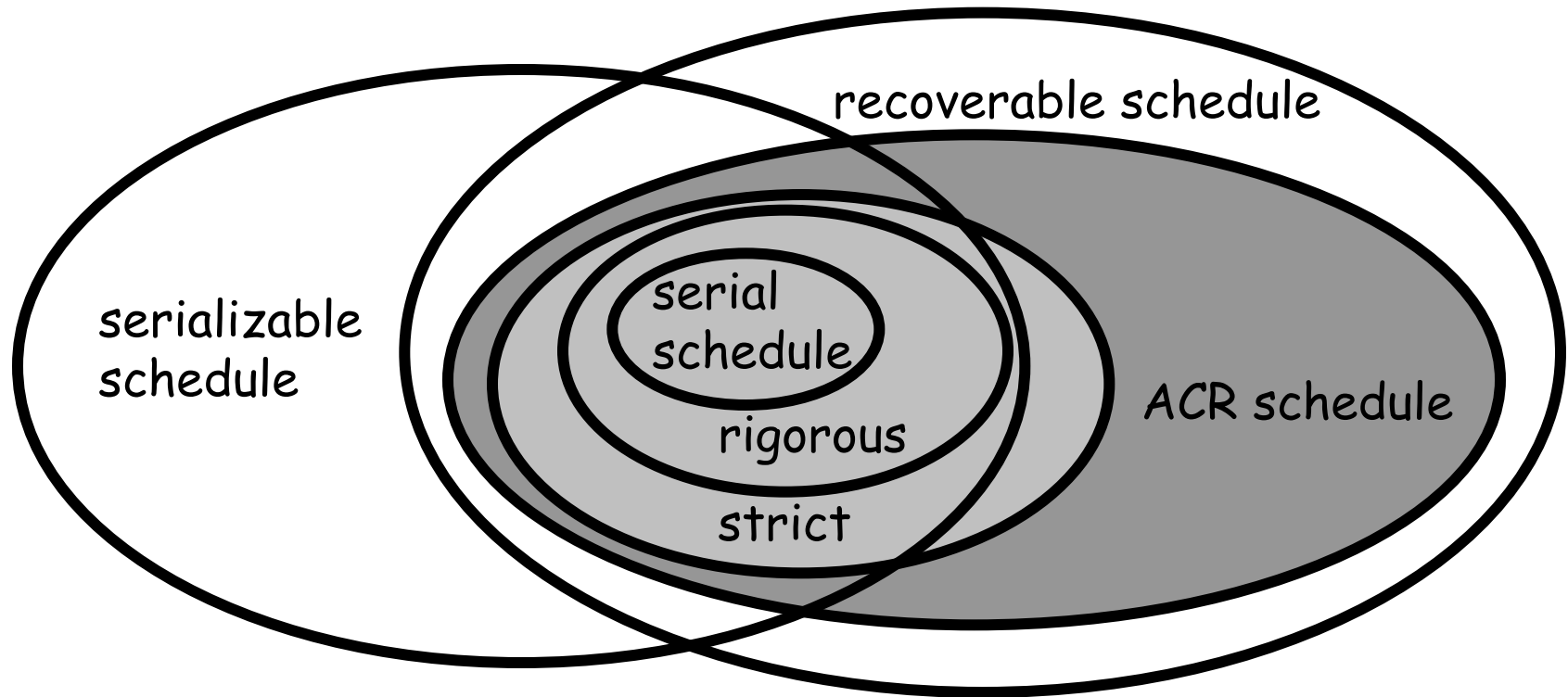


Rigorous schedules

- We say that a schedule S is **rigorous** if for each pair of conflicting actions a_i (belonging to transaction T_i) and b_j (belonging to transaction T_j) appearing in S , the commit command c_i of T_i appears in S between a_i and b_j .
- It is easy to verify that every rigorous schedule is strict.



Strict schedules and ACR



Obviously, every serial schedule is rigorous, and every rigorous schedule is strict, and therefore ACR, and recoverable. However, not all strict schedules are rigorous.



Recoverability and 2PL

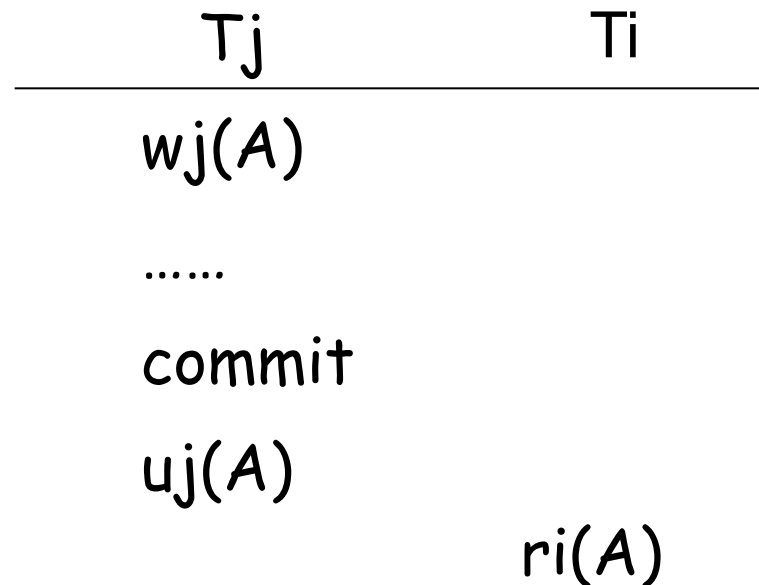
- So far, when discussing about recoverability, ACR, strictness and rigorousness we focused on:
 - read, write
 - rollback
 - commit
- We still have to study the impact of these notions on the locking mechanisms and the 2PL protocol



Strict two-phase locking (strict 2PL)

A schedule S is said to be in the **strict 2PL** class if

- S is in 2PL, and
- S is strict.





Strict two-phase locking (strict 2PL)

With the goal of capturing the class of strict 2PL the following protocol has been defined: A schedule S follows the **strict 2PL protocol** if it follows the 2PL protocol, and all exclusive locks of every transaction T are kept by T until either T commits or rollbacks.

T_j	T_i
$w_j(A)$	
$r_j(B)$	
.....	
$u_j(B)$	
commit	
$u_j(A)$	
	$ri(A)$



Properties of strict 2PL

- Every schedule following the strict 2PL protocol is strict:
(see exercise 8)
- Every schedule following the strict 2PL protocol is serializable
 - Obvious, since every 2PL schedule is conflict-serializable!



Exercise 8

- Prove or disprove the following statement:

Every schedule following the strict 2PL protocol is strict.

- Prove or disprove the following statement:

Every schedule that is strict and follows the 2PL protocol also follows the strict 2PL protocol.



Strong strict two-phase locking (SS2PL)

A schedule S follows the **strong strict 2PL** protocol if it follows the 2PL protocol, and all locks of every transaction T are kept by T until either T commits or rollbacks.

T_j	T_i
$w_j(A)$	
$r_j(B)$	
.....	
commit	
$u_j(A)$	
$u_j(B)$	
	$ri(A)$



Properties of strong strict 2PL

- Every schedule following the strong strict 2PL protocol is rigorous
(see exercise 9)
- Every schedule S following the strong strict 2PL protocol is obviously serializable, and the commit order of S is also a conflict-serializability order. Indeed, it can be shown that every strong strict 2PL schedule S is conflict-equivalent to the serial schedule S' obtained from S by ignoring the transactions that have rolled back, and by choosing the order of transactions determined by the order of commit (the first transaction in S' is the first that has committed, the second transaction in S' is the second that has committed, and so on)



Exercise 9

- Prove or disprove the following statement:

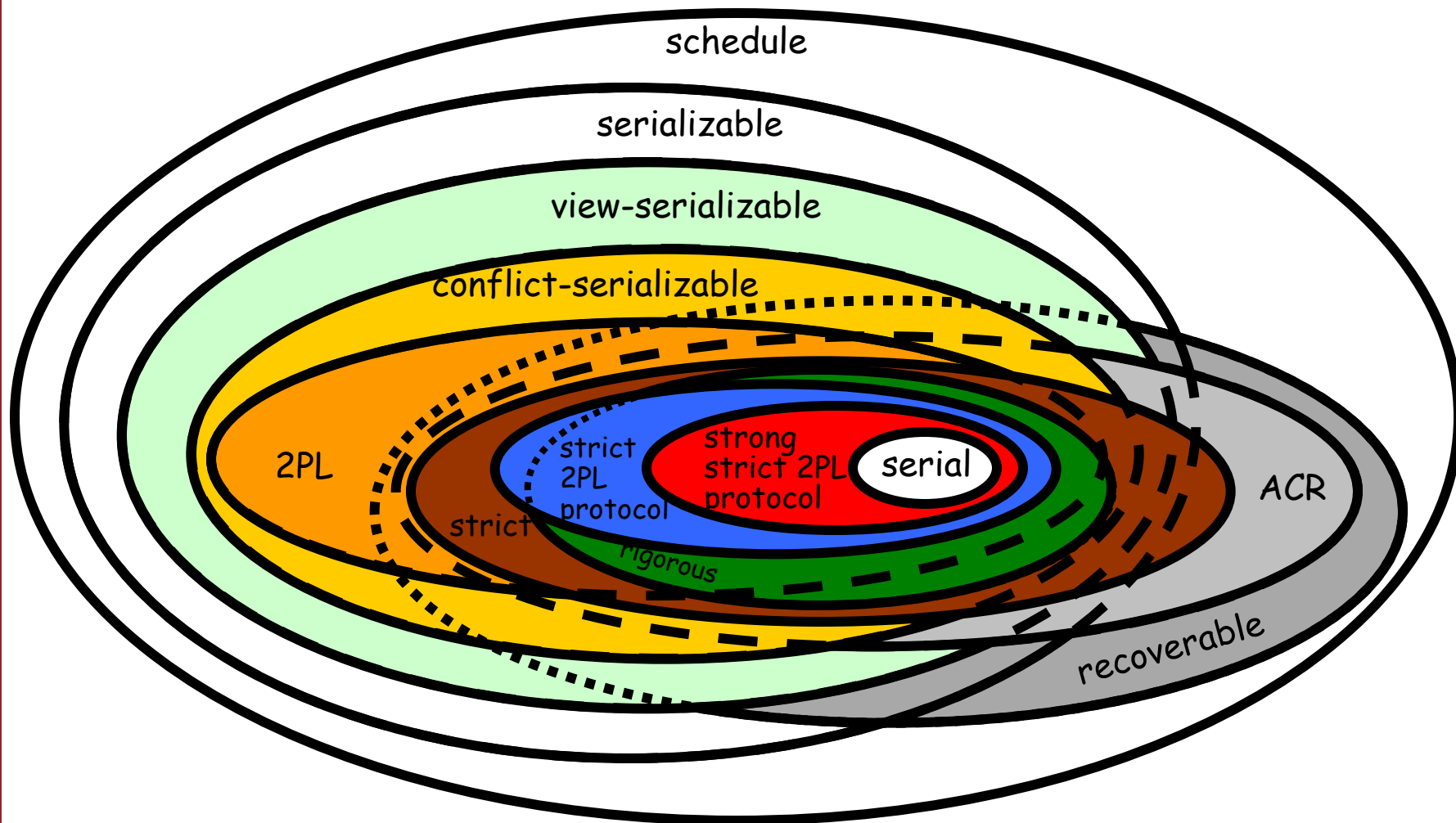
Every schedule following the strong strict 2PL protocol is rigorous.

- Prove or disprove the following statement:

Every schedule that is rigorous and follows the 2PL protocol also follows the strong strict 2PL protocol.



The complete picture





5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 Concurrency control through locks

5.5 Recoverability of transactions

5.6 **Concurrency control through timestamps**

5.7 Concurrency control in SQL



Concurrency based on timestamps

- Each transaction T has an associated **timestamp $ts(T)$** that is unique among the active transactions, and is such that $ts(T_j) < ts(T_h)$ whenever transaction T_j arrives at the scheduler before transaction T_h . In what follows, we assume that the timestamp of transaction T_i is simply i : $ts(T_i) = i$.
- Note that the timestamps actually define a total order on transactions, in the sense that they can be considered ordered according to the order in which they arrive at the scheduler.
- **Note also that every schedule respecting the timestamp order is conflict-serializable, because it is conflict-equivalent to the serial schedule corresponding to the timestamp order.**
- Obviously, the use of timestamp avoids the use of locks. Note, however, that deadlock may still occur.



The use of timestamp

- Transactions execute without any need of protocols.
- The basic idea is that, at each action execution, the scheduler checks whether the involved timestamps violates the serializability condition according to the order induced by the timestamps.
- In particular, we maintain the following data for each element X:
 - **rts(X)**: the highest timestamp among the active transactions that have read X
 - **wts(X)**: the highest timestamp among the active transactions that have written X (this coincides with the timestamp of the last transaction that wrote X)
 - **wts-c(X)**: the timestamp of the last committed transaction that has written X
 - **cb(X)**: a bit (called commit-bit), that is false if the last transaction that wrote X has not committed yet, and true otherwise.

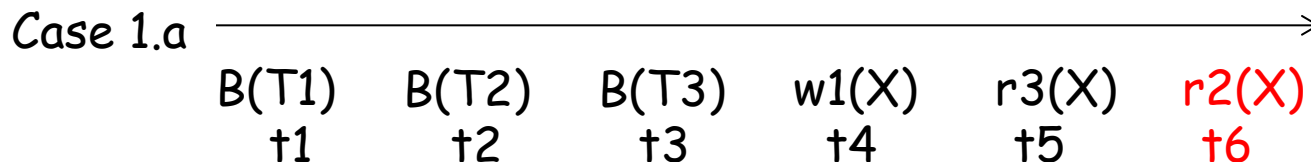


The rules for timestamps

- Basic idea:
 - the actions of transaction T in a schedule S must be considered as being logically executed in one spot
 - the logical time of an action of T is the timestamp of T , i.e., $ts(T)$
 - the commit-bit is used to avoid the dirty read anomaly
- The system manages two “temporal axes”, corresponding to the “physical” and to the “logical” time. The values $rts(X)$ and $wts(X)$ indicate the timestamp of the transaction that was the last to read and write X according to the logical time.
- An action of transaction T executed at the physical time t is accepted if its ordering according to the physical temporal order is compatible with respect to the logical time $ts(T)$
- This “compatibility principle” is checked by the scheduler.
- As we said before, we assume that the timestamp of each transaction T_i coincide with the subscript i : $ts(T_i)=i$. In what follows, t_1, \dots, t_n will denote physical times.



Rules – case 1a (read ok)



Consider $r2(X)$ with respect to the last write on X , namely $w1(X)$:

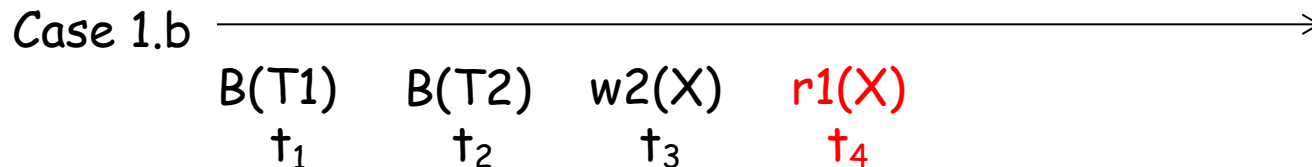
- the physical time of $r2(X)$ is $t6$, that is greater than the physical time of $w1$ ($t4$)
- the logical time of $r2(X)$ is $ts(T2)$, that is greater than the logical time of $w1(X)$, which is $wts(X) = ts(T1)$

We conclude that there is no incompatibility between the physical and the logical time, and therefore we proceed as follows:

1. if $cb(X)$ is true, then
 - generally speaking, after a read on X of T , $rts(X)$ should be set to the maximum between $rts(X)$ and $ts(T)$ – in the example, although according to the physical time $r2(X)$ appears after the last read $r3(X)$ on X , it logically precedes $r3(X)$, and therefore, if $cb(X)$ was true, $rts(X)$ would remain equal to $ts(T3)$
 - $r2(X)$ is executed, and the schedule goes on
2. if $cb(X)$ is false (as in the example), then $T2$ is put in a state waiting for the commit or the rollback of the transaction T' that was the last to write X (i.e., a state waiting for $cb(X)$ equal true -- indeed, $cb(X)$ is set to true both when T' commits, and when T' rollbacks, because the transactions T'' that was the last to write X before T' obviously committed, otherwise T' would be still blocked)



Rules – case 1b (read too late)



Consider $r1(X)$ with respect to the last write on X , namely $w2(X)$:

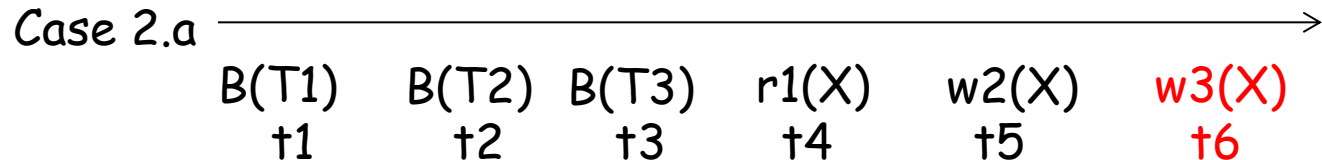
- the physical time of $r1(X)$ is t_4 , that is greater than the physical time of $w2(X)$, that is t_3
- the logical time of $r1(X)$ is $ts(T1)$, that is less than the logical time of $w2(X)$, i.e., $wt_s(X) = ts(T2)$

We conclude that $r1(X)$ and $w2(X)$ are incompatible.

Action $r1(X)$ of $T1$ cannot be executed, $T1$ rollbacks, and a new execution of $T1$ starts, with a new timestamp.



Rules – case 2a (write ok)



Consider $w3(X)$ with respect to the last read on X ($r1(X)$) and the last write on X ($w2(X)$):

- the physical time of $w3(X)$ is greater than that of $r1(X)$ and $w2(X)$
- the logical time of $w3(X)$ is greater than that of $r1(X)$ and $w2(X)$

We can conclude that there is no incompatibility. Therefore:

1. if $cb(X)$ is true or no active transaction wrote X , then
 - we set $wts(X)$ to $ts(T3)$
 - we set $cb(X)$ to false
 - action $w3(X)$ of $T3$ is executed, and the schedule goes on
2. else $T3$ is put in a state waiting for the commit or the rollback of the transaction T' that was the last to write X (i.e., a state waiting for $cb(X)$ equal true -- indeed, $cb(X)$ is set to true both when T' commits, and when T' rolls back, because the transactions T'' that was the last to write X before T' obviously committed, otherwise T' would be still blocked)



Rules – case 2b (Thomas rule)

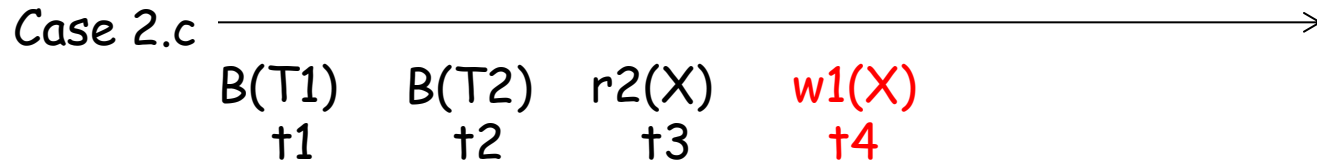
Case 2.b

B(T1)	B(T2)	B(T3)	r1(X)	w2(X)	w1(X)
t1	t2	t3	t4	t5	t6

- Consider w1(X) with respect to the last read r1(X) on X: the physical time of w1(X) is greater than the physical time of r1(X), and, since w1(X) and r1(X) belong to the same transaction, there is no incompatibility with respect to the logical time.
- However, on the logical time dimension, w2(X) comes after the write w1(X), and therefore, the execution of w1(X) would correspond to an update loss. Therefore:
 1. If cb(X) is true, we simply **ignore** w1(X) (i.e., w1(X) is not executed). In this way, the effect is to correctly overwrite the value written by T1 on X with the value written by T2 on X (it is like pretending that w1(X) came before w2(X))
 2. if cb(X) is false, we let T1 waiting either for the commit or for the rollback of the transaction that was the last to write X (i.e., a state waiting for cb(X) equal true -- indeed, cb(X) is set to true both when T' commits, and when T' rollbacks, because the transactions T'' that was the last to write X before T' obviously committed, otherwise T' would be still blocked)



Rules – case 2c (write too late)



Consider $w1(X)$ with respect to the last read $r2(X)$ on X :

- the physical time of $w1(X)$ is $t4$, that is greater than the physical time of $r2(X)$, i.e., $t3$
- the logical time of $w1(X)$ is $ts(T1)$, that is less than the logical time of $r2(X)$, that is $rts(X) = ts(T2)$

We conclude that $w1(X)$ and $r2(X)$ are incompatible.

Action $w1(X)$ is not executed, $T1$ is aborted, and is executed again with a new timestamp.



Timestamp-based method: the scheduler

Action $ri(X)$:

```
if      ts(Ti) >= wts(X)
then    if cb(X)=true or ts(Ti) = wts(X)           // (case 1.a)
          then set rts(X) = max(ts(Ti), rts(X)) and execute ri(X) // (case 1.a.1)
          else put Ti in “waiting” for the commit or the
                rollback of the last transaction that wrote X // (case 1.a.2)
else    rollback(Ti)                             // (case 1.b)
```

Action $wi(X)$:

```
if      ts(Ti) >= rts(X) and ts(Ti) >= wts(X)
then    if cb(X) = true
          then set wts(X) = ts(Ti), cb(X) = false, and execute wi(X) // (case 2.a.1)
          else put Ti in “waiting” for the commit or the
                rollback of the last transaction that wrote X // (case 2.a.2)
else    if ts(Ti) >= rts(X) and ts(Ti) < wts(X) // (case 2.b)
          then if cb(X)=true
                then ignore wi(X) // (case 2.b.1)
                else put Ti in “waiting” for the commit or the
                      rollback of the last transaction that wrote X // (case 2.b.2)
          else rollback(Ti) // (case 2.c)
```



Timestamp-based method: the scheduler

When T_i executes ci :

for each element X written by T_i ,
 set $cb(X) = \text{true}$
 for each transaction T_j waiting for $cb(X)=\text{true}$ or for the
 rollback of the transaction that was the last to
 write X , allow T_j to proceed
choose the transaction that proceeds

When T_i executes the rollback bi :

for each element X written by T_i , set $wts(X)$ to be $wts-c(X)$, i.e., the
 timestamp of the transaction T_j that wrote X before T_i , and set
 $cb(X)$ to true (indeed, T_j has surely committed)
 for each transaction T_j waiting for $cb(X)=\text{true}$ or for the
 rollback of the transaction that was the last to
 write X allow T_j to proceed
choose the transaction that proceeds



Deadlock with the timestamps

Unfortunately, the method based on timestamps does not avoid the risk of deadlock (although the probability is lower than in the case of locks).

The deadlock is related to the use of the commit-bit. Consider the following example:

$w1(B), w2(A), w1(A), r2(B)$

When executing $w1(A)$, T1 is put in waiting for the commit or the rollback of T2. When executing $r2(B)$, T2 is put in waiting for the commit or the rollback of T1.

The deadlock problem in the method based on timestamps is handled with the same techniques used in the 2PL method.



The method based on timestamp: example

Action	Effect	New values
r6(A)	ok	rts(A) = 6
r8(A)	ok	rts(A) = 8
r9(A)	ok	rts(A) = 9
w8(A)	no	T8 aborted
w11(A)	ok	wts(A) = 11
r10(A)	no	T10 aborted
c11	ok	cb(A) = true



Timestamps and conflict-serializability

- There are conflict-serializable schedules that are not accepted by the timestamp-based scheduler, such as:

$r_1(Y) \ r_2(X) \ w_1(X)$

- If the schedule S is processed by the timestamp-based scheduler without using the Thomas rule, then the schedule obtained from S by removing all actions of rolled back transactions is conflict-serializable
- If the schedule S is accepted by the timestamp-based scheduler using the Thomas rule, then S may be not conflict-serializable, like for example:

$r_1(A) \ w_2(A) \ c_2 \ w_1(A) \ c_1$

However, if the schedule S is processed by the timestamp-based scheduler using the Thomas rule, then the schedule obtained from S by removing all actions ignored by the Thomas rules and all actions of rolled back transactions is conflict-serializable



Comparison between timestamps and 2PL

- There are schedules that are accepted by timestamp-based schedulers that are not in 2PL, such as

r1(A) w2(A) r3(A) r1(B) w2(B) r1(C) w3(C) r4(C) w4(B) w5(B)

(that is not 2PL because T2 must release the lock on A before asking for the lock on B)

- Obviously, there are schedules that are accepted by the timestamp-based schedulers and are also strict 2PL schedules, such as the serial schedule:

r1(A) w1(A) r2(A) w2(A)

- There are strong strict 2PL schedules that are not accepted by the timestamp-based scheduler, such as:

r1(B) r2(A) w2(A) r1(A) w1(A)



Comparison between timestamps and 2PL

- Waiting stage
 - 2PL: transactions waiting for locks are put in waiting stage
 - TS: transactions reading too late or writing too late are killed and re-started; the waiting stage is only for transactions waiting for other transaction to commit or rollback
- Serialization order
 - 2PL: determined by conflicts
 - TS: determined by timestamps
- Need to wait for commit by other transactions
 - 2PL: solved by the strong strict 2PL protocol
 - TS: buffering of write actions (waiting for $cb(X) = \text{true}$)
- Deadlock
 - 2PL: risk of deadlock
 - TS: deadlock is less probable

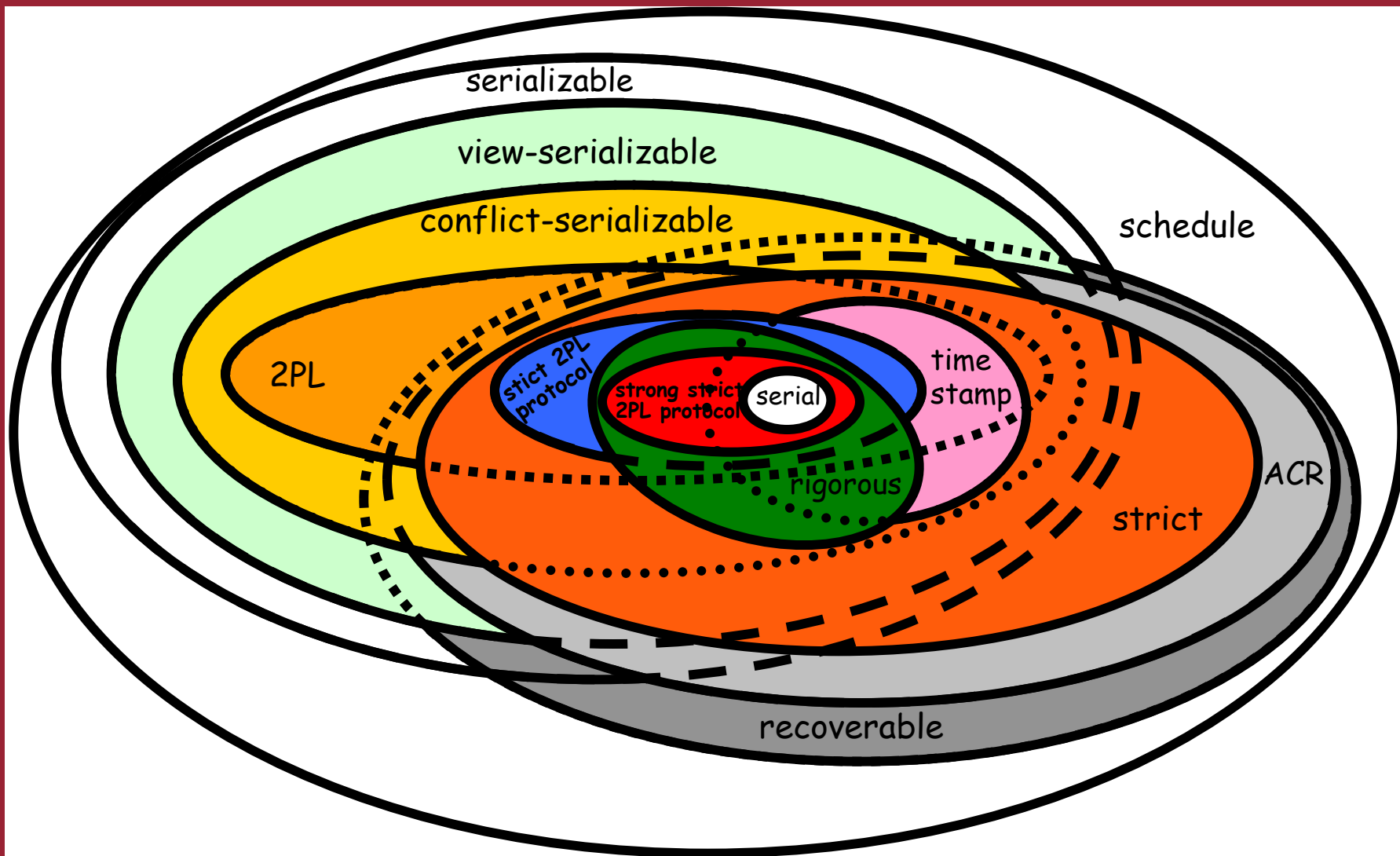


Comparison between timestamps and 2PL

- Timestamp-based method is superior when transactions are “read-only”, or when concurrent transactions rarely write the same elements
- 2PL is superior when the number of conflicts is high because:
 - although locking may delay transactions and may cause deadlock (and therefore rollback),
 - the probability of rollback is higher in the case of the timestamp-based method, and this causes a greater global delay of the system
- In the following picture, the set indicated by “timestamp” denotes the set of schedules generated by the timestamp-based scheduler, where all actions ignored by the Thomas rule and all actions of rolled back transactions are removed



The final picture





Multi-version concurrency control

Example 5.1:

$S = w_0(x) w_0(y) r_1(x) w_1(x) r_2(x) w_2(y) r_1(y) w_1(z) c_0 c_1 c_2$ Not conflict serializable

but: S would be tolerable if $r_1(y)$ could read the **old version** y_0 of y , to be "coherent" with $r_1(x)$ that read x_0

→ S would then be equivalent to serial $S' = t_0 t_1 t_2$

Approach:

- each write action creates a new version
- each read action can choose which version it wants/needs to read
- versions are transparent to application and transient (i.e., subject to garbage collection)



Snapshot isolation

- It is a kind of Multiversion Concurrency Control Mechanism (used in Postgres, for example)
- A transaction executing under snapshot isolation appears to operate on a personal snapshot of the database, taken at the start of the transaction
- When the transaction concludes, it will successfully commit only if the values updated by the transaction have not been changed externally since the snapshot was taken. If such a write-write conflict occur, it will cause the transaction to abort
- Readers never conflict with writers
- Does not guarantee serializability
- We will not study snapshot isolation; rather, we concentrate in the following on another multiversion concurrency control mechanism, i.e., multiversion timestamp-based method.



Multiversion timestamp

Idea: do not block the read actions! This is done by introducing **different versions $X_1 \dots X_n$ of element X** , so that **every read can be always executed**, provided that the “right” version (according to the logical time determined by the timestamp) is chosen

- Every “legal” write $w_i(X)$ generates a new version X_i (in our notation, the subscript corresponds to the timestamp of the transaction that generated X)
- To each version X_h of X , the timestamp $wts(X_h)=ts(Th)$ is associated, denoting the ts of the transaction that wrote that version
- To each version X_h of X , the timestamp $rts(X_h)=ts(T_i)$ is associated, denoting the highest ts among those of the transactions that read X_h

The properties of the multiversion timestamp are similar to those of the timestamp method.



New rules for the use of timestamps

The scheduler uses timestamps as follows:

- when executing $w_i(X)$: if a read $r_j(X_k)$ such that $wts(X_k) < ts(T_i) < ts(T_j)$ already occurred, then the write is refused (it is a “write too late” case, because transaction T_j , that is older than T_i , has already read a version of X that precedes version X_i), otherwise the write is executed on a new version X_i of X , and we set $wts(X_i) = ts(T_i)$.
- when executing $r_i(X)$: the read is executed on the version X_j such that $wts(X_j)$ is the highest write timestamp among the versions of X having a write timestamp less than or equal to $ts(T_i)$, i.e.: X_j is such that $wts(X_j) \leq ts(T_i)$, and there is no version X_h such that $wts(X_j) < wts(X_h) \leq ts(T_i)$. Note that such a version always exists, because it is impossible that all versions of X are greater than $ts(T_i)$. Obviously, $rts(X_j)$ is updated in the usual way.
- For X_j with $wts(X_j)$ such that no active transaction has timestamp less than j , the versions of X that precede X_j are deleted, from the oldest to the newest.
- To ensure recoverability, the commit of T_i is delayed until all commit of the transactions T_j that wrote versions read by T_i are executed.



New rules for the use of timestamps

The scheduler uses suitable data structures:

- For each version X_i the scheduler maintains a range $\text{range}(X_i) = [\text{wts}, \text{rts}]$, where wts is the timestamp of the transaction that wrote X_i , and rts is the highest timestamp among those of the transactions that read X_i (if no one read X_i , then $\text{rts} = \text{wts}$).
- We denote with $\text{ranges}(X)$ the set:
$$\{ \text{range}(X_i) \mid X_i \text{ is a version of } X \}$$
- When $\text{ri}(X)$ is processed, the scheduler uses $\text{ranges}(X)$ to find the version X_j such that $\text{range}(X_j) = [\text{wts}, \text{rts}]$ has the highest wts that is less than or equal to the timestamp $\text{ts}(T_i)$ of T_i . Moreover, if $\text{ts}(T_i) > \text{rts}$, then the rts of $\text{range}(X_j)$ is set to $\text{ts}(T_i)$.
- When $\text{wi}(x)$ is processed, the scheduler uses $\text{ranges}(X)$ to find the version X_j such that $\text{range}(X_j) = [\text{wts}, \text{rts}]$ has the highest wts that is less than or equal to the timestamp $\text{ts}(T_i)$ of T_i . Moreover, if $\text{rts} > \text{ts}(T_i)$, then $\text{wi}(X)$ is rejected, else $\text{wi}(X_i)$ is accepted, and the version X_i with $\text{range}(X_i) = [\text{wts}, \text{rts}]$, with $\text{wts} = \text{rts} = \text{ts}(T_i)$ is created.



Multiversion timestamp: example

Suppose that the current version of A is A0, with $rts(A0)=0$.

$T1_{(ts=1)}$ $T2_{(ts=2)}$ $T3_{(ts=3)}$ $T4_{(ts=4)}$ $T5_{(ts=5)}$

$r1(A)$

$w1(A)$

$r2(A)$

$w2(A)$

$r4(A)$

$r5(A)$

$w3(A)$

reads A0, and set $rts(A0)=1$

writes the new version A1

reads A1, and set $rts(A1)=2$

writes the new version A2

reads A2, and set $rts(A2)=4$

reads A2, and set $rts(A2)=5$

rollback T3



Typical strategy of commercial systems

The scheduler distinguishes between two classes of transactions:

- The transactions with read and write are executed under the 2PL protocol
- The transactions that are “read only” are executed under a “multiversion-based” method (such as the multiversion timestamp-based method)



5. Transaction management

5.1 Transactions, concurrency, serializability

5.2 View-serializability

5.3 Conflict-serializability

5.4 Concurrency control through locks

5.5 Recoverability of transactions

5.6 Concurrency control through timestamps

5.7 **Concurrency control in SQL**



Transactions in SQL

- The SQL engine is used by sessions.
- If a session does not explicitly define a transaction, then every SQL command (select, insert, update, etc.) is considered a transaction that ends when the execution finishes
- If a session does define a transaction, then it starts with the BEGIN command, and will end with the COMMIT command or, equivalently, with END
- Within a transaction we may have the ROLLBACK commands (undoing all actions of the transaction)
- Within a transaction we may also have LOCK/UNLOCK commands, if the DBMS uses locking



The anomalies considered in SQL

- **Dirty read**

as we have seen, this anomaly occurs when a transaction reads an element written by a transaction that has not committed or rolled back yet

- **Nonrepeatable read**

as we have seen, this anomaly occurs when a transaction reads the same element twice

- **Phantom read**

this is a new kind of anomaly, that occurs when a transaction T1 executes a “range” query (like “select * from person where age > 10 and age < 40”), then another transaction T2 inserts a new tuple satisfying the range, and then T1 executes again the the same range query, finding different results. It can be avoided by “range” locks.



Concurrency in SQL - standard

Isolation Level	<i>Dirty Read</i>	<i>Nonrepeatable Read</i>	<i>Phantom Read</i>
<i>Read uncommitted</i>	Possible	Possible	Possible
<i>Read committed</i>	Not possible	Possible	Possible
<i>Repeatable read</i>	Not possible	Not possible	Possible
<i>Serializable</i>	Not possible	Not possible	Not possible

- Except for the “Read uncommitted” level, each other level guarantees the absence of a specified set of anomalies (see table above). Serializability is the maximum level of correctness
- Every transaction decides its level of isolation (SET TRANSACTION ISOLATION LEVEL <level>); we can always know the current isolation level (command SHOW TRANSACTION ISOLATION LEVEL)
- The standard does not impose any constraints on the implementation of the concurrency control mechanisms



Possible implementation with locking

- “Read uncommitted”: no action is taken for controlling concurrency.
- “Read committed”: a lock-based concurrency control implementation keeps write locks until the end of the transaction, but read locks are released as soon as the SELECT operation is performed. Range-locks are not managed. It makes no promise whatsoever that if the transaction re-issues the read, it will find the same data
- “Repeatable read”: a lock-based concurrency control implementation keeps read and write locks until the end of the transaction. However, range-locks are not managed, so phantom reads can occur.
- “Serializable”: a lock-based concurrency control implementation keeps read and write locks to be released at the end of the transaction, as in strong strict 2PL. Also range-locks (locks on all possible elements satisfying the range) must be acquired when a SELECT query uses ranged WHERE clause, to avoid the phantom reads phenomenon.

NOTE: The SQL standard permits a DBMS to run a transaction at an isolation level stronger than that requested (e.g., a "Read committed" transaction may actually be performed at a "Repeatable read" isolation level).



Example: the Postgres DBMS

- The minimum isolation level is “read committed”, which is the default level
- The concurrency control strategy of Postgres is a sort of multiversion control
- Reads are never blocked
- The implementation of the strategy is based on locking
- Explicit LOCK/UNLOCK commands are allowed
- Deadlock management based on recognition