



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG

FAKULTÄT  
FÜR MATHEMATIK, INFORMATIK  
UND NATURWISSENSCHAFTEN

# Databases and Information Systems (DIS)

Dr. Annett Ungethüm  
Universität Hamburg

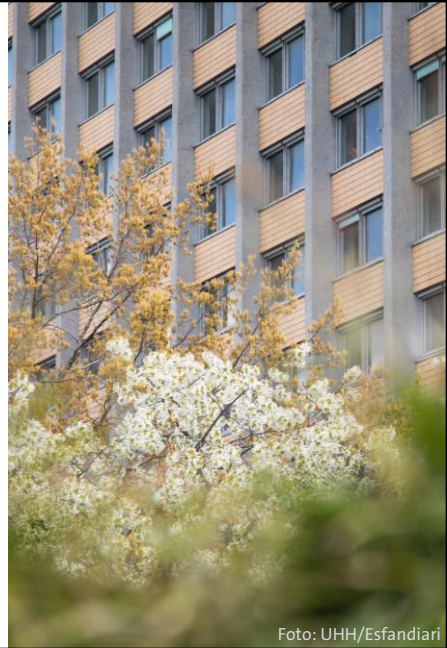
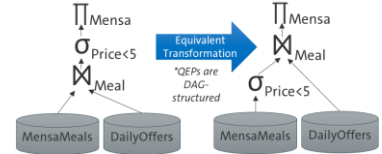
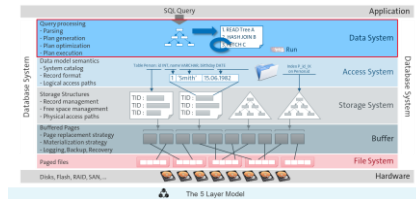


Foto: UHH/Esfandiari

# Summary

- Architecture of Database Systems
- Transaction Management
- Modern Database Technology
- Data Warehouses and OLAP
- Data Mining
- Big Data Analytics



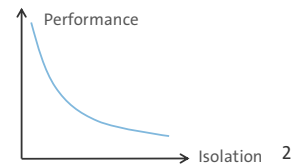
- Architectures of Database Systems
- Transaction Management
- Modern Database Technology
- Data Warehouses and OLAP
- Data Mining
- Big Data Analytics

Atomicity  
Consistency and semantic integrity  
Isolated execution  
Durability

Things that can go wrong, will go wrong: Anomalies

- Dirty Read
- Non-Repeatable Read
- Lost Update
- Phantom Problem

Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
1	r(A)			8		w(C)	
2		r(B)		9		w(A)	
3		r(C)		10			r(A)
4		w(B)		11		w(B)	
5	r(B)			12			w(C)
6	w(B)			13			w(A)
7		r(A)		14			



## ANSI-SQL isolation levels

Isolation Level	Lost Update possible?	Dirty Read possible?	Non-Repeatable Read possible?	Phantom Read possible?
Read Uncommitted	No	Yes	Yes	Yes
Read Committed	No	No	Yes	Yes
Repeatable Read	No	No	No	Yes
Serializable	No	No	No	No

Set isolation level with SQL:

SET TRANSACTION ISOLATION LEVEL

{ READ UNCOMMITTED | READ COMMITTED |  
REPEATABLE READ | SERIALIZABLE }

Special case: read uncommitted → Should avoid Lost Update, but does not always avoid them in reality

A more in-depth discussion and critique of the isolation levels can be found in the following paper:

Berenson, Hal, et al. "A critique of ANSI SQL isolation levels." *ACM SIGMOD Record* 24.2 (1995): 1-10. (<https://dl.acm.org/doi/pdf/10.1145/568271.223785>)

→ The list of authors might make it look heavily biased, which it might be. But the authors still have a point, and the paper is peer-reviewed

# Serializability

The concurrent execution of a set of transactions is considered to be correct, if there is a serial execution of the same set of transactions, leading to **the same resulting DB state** as well as **the same output values** as the original sequential execution.

- Each schedule that has the same effect as a serial one is considered to be correct
- A schedule is serializable if there are no loops in its dependency graph

Complexity of topological sorting to get the possible serializable schedules is ( $O(N^2)$ )

More about the correctness of schedules:

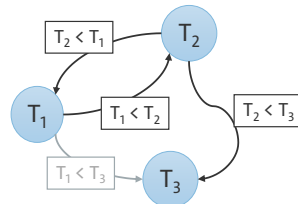
Vossen, Gottfried. "Database transaction models." *Computer Science Today: Recent Trends and Developments* (2005): 560-574.

## Example: Dependency Graph

Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
1	r(A)			8		w(C)	
2		r(B)		9		w(A)	
3		r(C)		10			r(A)
4		w(B)		11			r(C)
5	r(B)			12	w(B)		
6	w(A)			13			w(C)
7		r(A)		14			w(A)

### Conflicts

- $w_2(B) < r_1(B)$
- $w_1(A) < w_2(A)$
- $w_1(A) < r_2(A)$
- $w_2(A) < w_3(A)$
- $w_2(C) < r_3(C)$
- $w_1(A) < w_3(A)$
- $w_2(A) < r_3(A)$
- $w_2(B) < w_1(B)$
- $w_2(C) < w_3(C)$

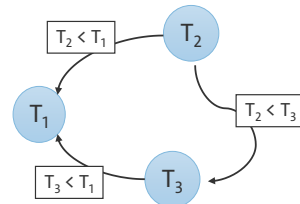


→ Not serializable!

- Cycle:  $T_1 \leftrightarrow T_2$
- $w_x(Z) < r_y(Z) \rightarrow$  Dirty Read
- $w_x(Z) < w_y(Z) \rightarrow$  Lost Update

## Example Schedule (II)

Step	$T_1$	$T_2$	$T_3$
1	$r(A)$		
2		$r(C)$	
3		$r(B)$	
4		$w(B)$	
5		$w(C)$	
6			$r(C)$
7			$r(B)$
8			$w(B)$
9	$r(B)$		



→ Serializable! 👍

- Conflicts:
  - $w_2(B) < r_1(B) \rightarrow DR$
  - $w_3(B) < r_1(B) \rightarrow DR$
  - $w_2(C) < r_3(C) \rightarrow DR$
  - $w_2(B) < r_3(B) \rightarrow DR$
  - $w_2(B) < w_3(B) \rightarrow LU$

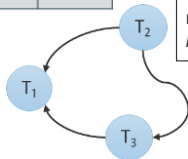
## Correct Schedules

Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
1	r(A)		
2		r(C)	
3		r(B)	
4		w(B)	
5		w(C)	
6			r(C)
7			r(B)
8			w(B)
9	r(B)		

$A:=1, B:=2, C:=3$   
 T<sub>1</sub>      T<sub>2</sub>      T<sub>3</sub>  
 read(A)      read(A)  
              read(B)  
              B:=B+10  
              C:=5  
              write(B)  
              write(C)  
  
 read(B)  
 Result: A=1, B=60, C=5

Schedule: T<sub>2</sub>, T<sub>3</sub>, T<sub>1</sub>

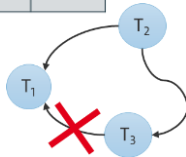
$A:=1, B:=2, C:=3$   
 T<sub>1</sub>      T<sub>2</sub>      T<sub>3</sub>  
          read(C)  
          read(B)  
          B:=B+10  
          C:=5  
          write(B)  
          write(C)  
  
          read(C)  
          read(B)  
          B:=B\*C  
          write(B)  
  
 read(A)  
 read(B)  
 Result: A=1, B=60, C=5



- Example for incorrect schedule: T<sub>1</sub>, T<sub>3</sub>, T<sub>2</sub>  
 → Results: A=1, B=2\*3+10=16≠60, C=5

## Correct Schedules (II)

Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
1	r(A)		
2		r(A)	
3		r(B)	
4		w(B)	
5		w(A)	
6			r(C)
7			r(B)
8			w(B)
9	r(A)		



Schedule: T<sub>2</sub>, T<sub>3</sub>, T<sub>1</sub>

A:=1, B:=2, C:=3

T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
	read(A)	
	read(B)	
	B:=B+10	
	A:=5	
	write(B)	
	write(A)	
read(A)		read(C)
		read(B)
		B:=B*C
		write(B)

Result: A=5, B=36, C=3

Schedule: T<sub>2</sub>, T<sub>1</sub>, T<sub>3</sub>

A:=1, B:=2, C:=3

T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>
	read(A)	
	read(B)	
	B:=B+10	
	A:=5	
	write(B)	
	write(A)	
read(A)		
		read(C)
		read(B)
		B:=B*C
		write(B)

Result: A=5, B=36, C=3

Schedules are equivalent!

- If T<sub>1</sub> does not read B, B can be changed before or after T<sub>1</sub> has finished  
 → T<sub>1</sub> does not require results from T<sub>3</sub> anymore  
 → Multiple correct schedules
- Effect of a schedule that is not correct:
  - T<sub>3</sub>, T<sub>2</sub>, T<sub>1</sub> → result: 2\*3+10=16 ~~≠36~~



## Equivalence

Let  $<_H$  and  $<_G$  be the ordering relations of the schedules  $H$  and  $G$ , then the two execution plans  $H$  and  $G$  are considered equivalent if the following conditions are true for all objects  $A$ , and Transactions  $i$  and  $j$ :

$$\begin{aligned} \text{read}_i[A] <_H \text{write}_j[A] &\Leftrightarrow \text{read}_i[A] <_G \text{write}_j[A] \\ \text{write}_i[A] <_H \text{read}_j[A] &\Leftrightarrow \text{write}_i[A] <_G \text{read}_j[A] \\ \text{write}_i[A] <_H \text{write}_j[A] &\Leftrightarrow \text{write}_i[A] <_G \text{write}_j[A] \end{aligned}$$

→ Not all correct schedules are equivalent!

## Correctness and Equivalency

A:= 1000, B:=2000

T <sub>1</sub>	T <sub>2</sub>
read(A)	read(A)
A:=A-50	X := A*0.1
write(A)	A := A-X
read(B)	write(A)
B := B+50	read(B)
write(B)	B := B+X
	write(B)

Both schedules are correct, but they are not equivalent!

T<sub>1</sub>  
read(A)  
A:=A-50  
write(A)  
read(B)  
B := B+50  
write(B)

T<sub>2</sub>  
  
  
  
  
  
  
  
  
read(A)  
X := A\*0.1  
A := A-X  
write(A)  
read(B)  
B := B+X  
write(B)

Results  
A= (1000-50) - (1000-50)\*0.1 = 855  
B= 2000+50+ (1000-50)\*0.1 = 2145

T <sub>1</sub>	T <sub>2</sub>
	read(A)
	X := A*0.1
	A := A-X
	write(A)
	read(B)
	B := B+X
	write(B)
read(A)	
A:=A-50	
write(A)	
read(B)	
B := B+50	
write(B)	

Results  
A= 1000-(1000\*0.1)-50 = 850  
B= 2000+(1000\*0.1)+50 = 2150

For n Transactions, there are n! possible serial schedules

→ Reminder from slide 13: Correctness means, that a plan is equivalent to one of these serial schedules, i.e. it has the same effect as one of the serial schedules!

## Properties of Transaction Schedules

### Recoverable (RC)

A transaction can only finish if all transactions it is reading from have already finished

→ Required for Durability

Example for not recoverable schedule  
 $H = w_1[x], r_2[x], w_2[y], c_2, a_1$

### Strict execution (ST)

Objects that have been changed by a transaction  $T_i$  cannot be read or changed by another transaction before transaction  $T_i$  has finished

→ Transactions can be reset using a before-image

### Avoid cascading aborts (ACA)

Each transaction  $T_a$  can only see results from other transactions  $T_b$  when the transactions  $T_b$  have finished

### RC

- A schedule is recoverable if the following condition is met:  
 If a transaction  $T_i$  reads from another transaction  $T_j$  ( $i \neq j$ ), then  $c_j < c_i$   
 → The transaction that was read from must commit before the reading transaction
- Example of not recoverable schedule  $H$   
 →  $T_2$  reads from  $T_1$ , but  $T_2$  commits first

### ST

- Physical logging using before images  
 → Physical logging stores values (instead of operations)

### ACA

- If reading results from another TA  $T$  is allowed before it has finished, an abort can lead to cascading resets  
 → All TAs that have read from the aborted TA  $T$  have to be reset + all TAs that used a result of  $T$  indirectly

Strict execution avoids cascades, but not necessarily the other way round

## (Counter) Examples ACA & ST

Schedule A

Step	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	T <sub>4</sub>	T <sub>5</sub>
0.	...				
1.	w <sub>1</sub> [x]				
2.		r <sub>2</sub> [x]			
3.		w <sub>2</sub> [y]			
4.			r <sub>3</sub> [y]		
5.			w <sub>3</sub> [z]		
6.				r <sub>4</sub> [z]	
7.				w <sub>4</sub> [v]	
8.					r <sub>5</sub> [v]
9.	a <sub>1</sub>				

Schedule B

B = r<sub>1</sub>[x] w<sub>1</sub>[x] w<sub>2</sub>[x] c<sub>1</sub> c<sub>2</sub>

Schedule A:

A  $\notin$  ACA  $\rightarrow$  When T<sub>1</sub> is aborted, T<sub>2</sub> – T<sub>5</sub> must also be rolled back

Schedule B:

B  $\notin$  ST  $\rightarrow$  w<sub>2</sub>(x) takes place before T<sub>1</sub> has finished

But

B  $\in$  ACA  $\rightarrow$  No transaction reads the result of the other transaction

# Database Scheduler

- Sorts the operations and ensures that the schedule/history is serializable and can be rolled back
- Operations can be
  - Executed immediately
  - Rejected
  - On hold
- Strategies:
  - Optimistic
  - Pessimistic



## Pessimistic Scheduling

- Operations are set back → Scheduler finds a sequence for operations
- Used when many transactions are expected to change the data
- Most common implementation is based on Locks → Operations are on hold until according lock is released

## Optimistic Scheduling

- Scheduler tries to execute operations as soon as possible
- Used When only few transactions are expected to change the data
  - Damage might happen and must be repaired afterwards
  - Popular implementations: Time stamp based scheduling, optimistic synchronization

## Pessimistic Synchronization

- Locks for exclusive access of data objects

		Active lock		
		NL	R	X
Requested lock	R	+	+	-
	X	+	-	-

NL – No Lock  
R – Read Lock  
X – eXclusive Lock

- Realizes a logical single-user setup
- For each data object, a central lock table stores the lock mode
- R lock: Object can only be read by other transactions, it cannot be changed
- X lock: Object can neither be read nor changed by another transaction
- Example: If an object is already locked with a read lock (R), an exclusive lock (X) cannot be acquired for that object, but another read lock (R) can be acquired.
- Read Lock is also sometimes referred to as “Shared lock (S)”

## Working with Locks

### Static locks

Claimed at the start of a transaction  
(Preclaiming)

### Dynamic locks

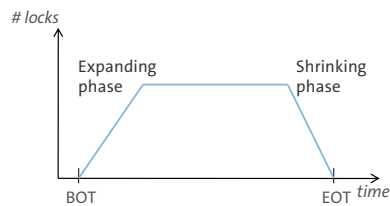
Claimed when needed  
→ Possibility of deadlocks

- Locks have to be kept until the transaction finishes to guarantee serializability
- This constraint can be relaxed for optimization reasons, e.g. early release of read locks

- Remember: We do not always need perfect serializability
- An early release of read locks can lead to different results
- More rules for working with locks:
  - Each object that is used by a transaction, has to be locked before its usage
  - A transaction that is currently holding a lock, cannot claim the same lock again
  - If a lock cannot be claimed (see last 14), the transaction is queued
  - A transaction must not acquire another lock after the first lock has been released (→ 2-phase-locking-protocol)
  - When a transaction ends, it must release all locks

## 2-Phase-Locking-Protocol (2PL)

- Expanding phase: Locks are acquired but not released
- Shrinking phase: Locks are released but not claimed



- Locks can only start to be released after all locks have been acquired
  - 2PL is almost(!) sufficient for a correct execution of transactions
- Deadlocks can still happen



## 2PL Example

Step	T <sub>1</sub>	T <sub>2</sub>	Comment
1.	BOT		
2.	lockX[x]		
3.	r[x]		
4.	w[x]		
5.		BOT	
6.		lockR[x]	T <sub>2</sub> must wait
7.	lockX[y]		
8.	r[y]		
9.	unlockX[x]		Wake up T <sub>2</sub>
10.		r[x]	
11.		lockR[y]	T <sub>2</sub> must wait
12.	w[y]		
13.	unlockX[y]		Wake up T <sub>2</sub>
14.		r[y]	
15.	commit		
16.		unlockR[x]	
17.		unlockR[y]	
18.		commit	

Expanding phase of T<sub>1</sub>: Step 1 to 7

Expanding phase of T<sub>2</sub>: Step 5 to 11

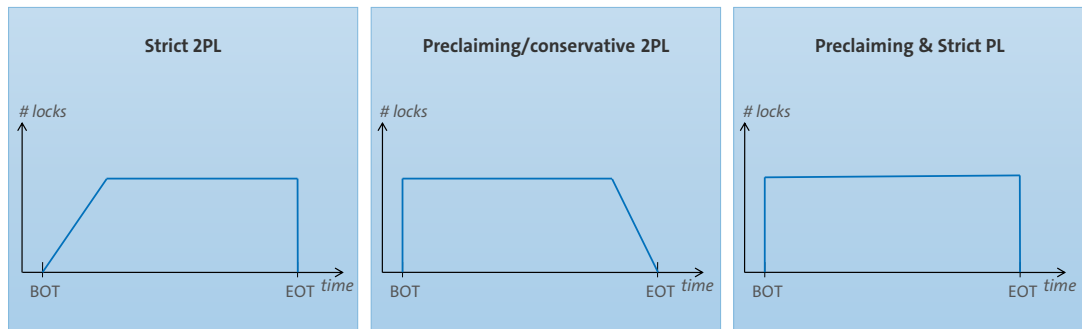
Shrinking phase of T<sub>1</sub>: Step 9 to 15

Shrinking phase of T<sub>2</sub>: Step 16 to 18

Step 6: T<sub>2</sub> cannot acquire lock for x because T<sub>1</sub> holds an exclusive lock on x → T<sub>2</sub> must wait

Step 9: The lock of x is released in T<sub>1</sub>, T<sub>2</sub> can continue

## 2PL Variations



### Strict 2PL

All locks are released after all other operations have finished

- Avoids cascading aborts of transactions that read the results of a rolled back transaction

### Preclaiming

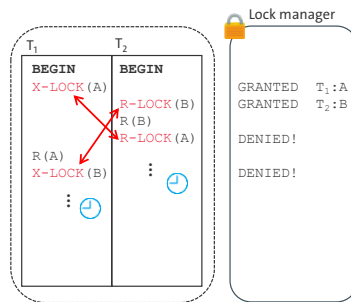
All locks are set before any other operations are executed

- Avoids deadlocks
- Decreases parallelism
- Requires to know all necessary locks at the BOT (not trivial in statements with branches)

### Strict 2PL and Preclaiming

Enforces sequential execution

# Deadlocks



➔ Next week