

6.5830 Lecture 14

Optimistic Concurrency Control

10/28/2024

The image features the iconic 'Jeopardy!' logo in a bold, metallic, 3D font. The letters are silver with a blue and yellow gradient, set against a vibrant space background of red, orange, and blue with stars and a planet's horizon. Below the logo, the text 'The Database Tour' is written in a bold, black, sans-serif font.

JEOPARDY!

The Database Tour

Rules

- Reverse Jeopardy: The answers are the answers
- Create groups of 3-5 people. You must answer as a group
- You receive candy for correct answers.
- Groups willing to answer the question must raise their hand
- Among all hand-raised groups the moderator randomly picks a group.
- At the discretion of the moderator the answer is deemed correct or wrong. If correct, the group can choose the next question. If wrong, another group can answer (but no candy)
- **You must pick questions from the top down**

Terminology

**Optimistic
Concurrency
Control**

2 Phase Locking

**Other
Concurrency
Techniques**

4 Candy

8 Candy

5 Candy

4 Candy

5 Candy

8 Candy

1 Candy

**What is the connection
between Atomicity,
Isolation, and Durability
with Concurrency Control,
Logging?**

2 Candy

A few days ago, Sam (MIT's CS area chair) and I almost dressed alike. Let's assume the process of dressing is done as part of a transaction, which ensures that we always dress differently.

What ACID properties were violated?



3 Candy

If a system ensures
Atomicity, **D**urability, and
Isolation, is it by definition
not also **C**onsistent?
Explain by means of an
example why/why not?

4 Candy

From the assigned reading from lecture 11:

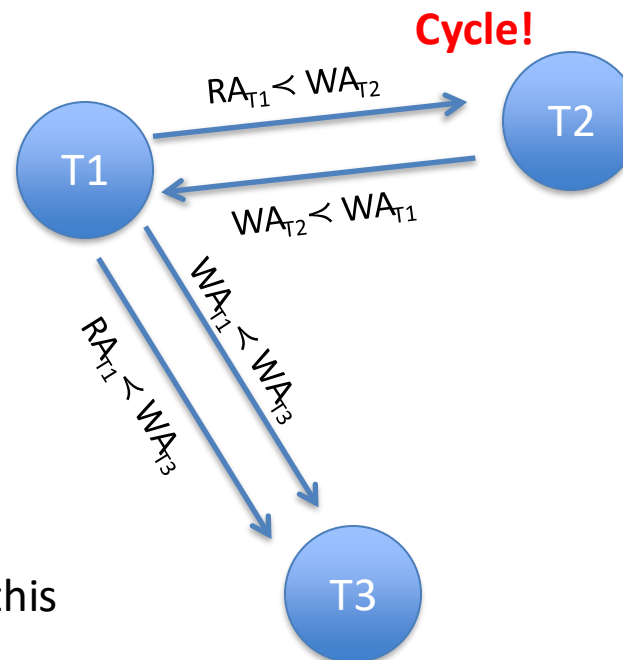
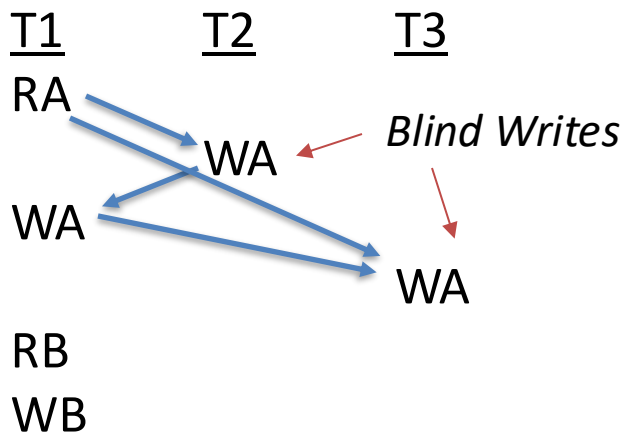
WAL, LSN, Undo, Redo, physical, logical, physiological are all important terms for describing logging. What do they stand for?

5 Candy

Name 2 types of serializability and give an example when one is valid under the definition but not the other?

View vs Conflict Serializable

- Testing for view serializability is NP-Hard
 - Have to consider all possible orderings
- Conflict serializability used in practice
 - Not because of NP-Hardness
 - Because we have a way to enforce it as transactions run
- Example of schedule that is view serializable but not conflict serializable:



Equivalent to T1, T2, T3

Conflict serializability does not permit this

Only happens with *blind writes*

4 Candy

What is the core idea of optimistic concurrency control?

What are the pros and cons?

(You should know that based on the assigned reading :)

Optimistic Concurrency Control (OCC)

- Alternative to locking for isolation
- Approach:
 - Store writes in a per-transaction buffer
 - Track read and write sets
 - At commit, check if transaction conflicted with earlier (concurrent) transactions
 - Abort transactions that conflict
 - Install writes at end of transaction
- “Optimistic” in that it does not block, hopes to “get lucky” arrive in serial interleaving



Tradeoff

- In OCC:
 - Never have to wait for locks
 - no deadlocks
- But...
 - Transactions that conflict often have to be restarted
 - **Transactions can "starve" -- e.g., be repeatedly restarted, never making progress**
- OCC will do better when the restart rate is low
 - (Less contention)
- Recent work on high performance transaction processing has focused on OCC because
 - OCC checks can be done between individual transactions
 - Unlike global shared lock table
 - Modern OCC systems obtain *insane* throughput (> 10M xactions / sec)

E.g., <https://people.eecs.berkeley.edu/~wzheng/silo.pdf>

4 Candy

What are the 3 phases of optimistic concurrency control?

OCC Implementation

- Divide transaction execution in 3 phases
 - **Read:** transaction executes on DB, stores local state
 - **Validate:** transaction checks if it can commit
 - **Write:** transaction writes state to DB

Read Phase

- Transactions execute, with updates affecting local copies of the data
- Build a list of data items that were read/written
 - *Read and write sets*
- Modify functions to read and write data from DB

OCC Write

twrite(object,value):

if object not in write_set: // never written, make copy

m = read(object)

copies[object] = m

write_set = write_set U {object}

write(copies[object], value)

By writing to local copies, we ensure dirty results aren't visible to other concurrent transactions

OCC Read

```
tread(object):  
    read_set = read_set U {object};  
    if object in write_set:  
        return read(copies[object]);  
    else:  
        return read(object);
```


Allows us to read our own writes!

Validation Phase (next)

- How do we know whether a transaction can commit?
 - *Validation Rules*
 - Check concurrent transactions for conflicts
- How do we implement validation efficiently?
 - *Validation Algorithm*

12 Candy

Under OCC there are 4 conditions. When T_j completes its read phase, require that for all $T_i < T_j$, **one of the following** conditions must be true for validation to succeed (T_j to commit) to ensure serializability:

- 1) T_i completes its write phase before T_j starts its read phase
- 2) $W(T_i)$ does not intersect with $R(T_j)$ or $W(t_j)$ 
- 3) $W(T_i)$ does not intersect $R(T_j)$, and T_i completes its write phase before T_j starts its write phase.
- 4) $W(T_i)$ does not intersect $R(T_j)$ or $W(T_j)$, and T_i completes its read phase before T_j completes its read phase.
- 5) $W(T_i)$ does not intersect $R(T_j)$ or $W(T_j)$, and $W(T_j)$ does not intersect $R(T_i)$ [no conflicts]

One of the conditions is not sufficient! Which one?

Validation Rules

XXX confusion about W and R as
Read / write sets vs read / write
phases

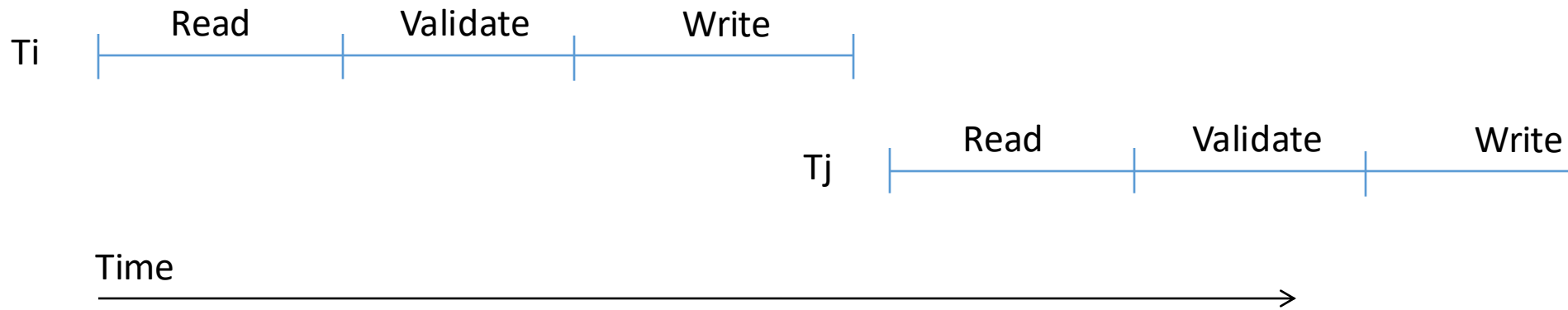
When T_j completes its read phase, require that for all $T_i < T_j$, one of the following conditions must be true for validation to succeed (T_j to commit):

- 1) T_i completes its write phase before T_j starts its read phase
- 2) $W(T_i)$ does not intersect $R(T_j)$, and T_i completes its write phase before T_j starts its write phase.
- 3) $W(T_i)$ does not intersect $R(T_j)$ or $W(T_j)$, and T_i completes its read phase before T_j completes its read phase.
- 4) $W(T_i)$ does not intersect $R(T_j)$ or $W(T_j)$, and $W(T_j)$ does not intersect $R(T_i)$
[no conflicts]

These rules will ensure serializability, with T_j being ordered after T_i with respect to conflicts

Condition 1

T_i completes its write phase before T_j starts its read phase



Don't overlap at all.

Condition 2

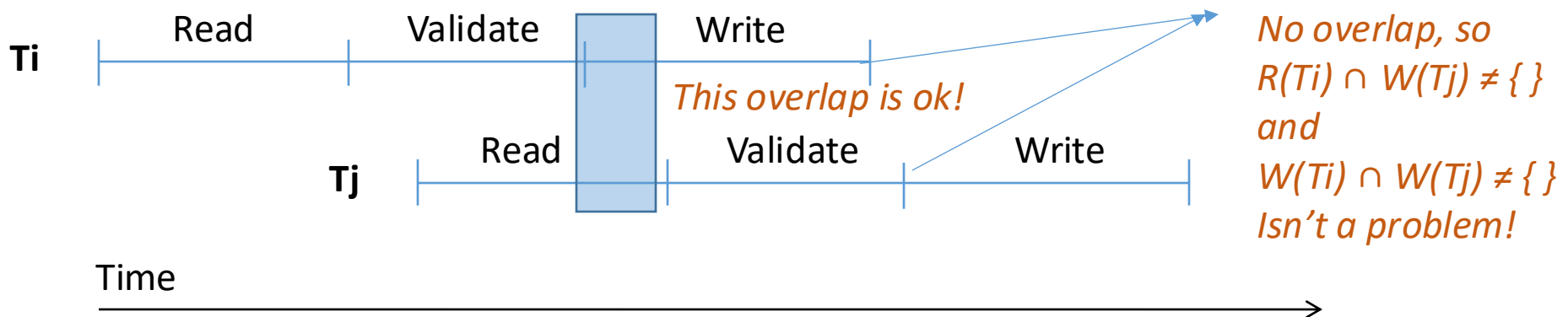
$W(T_i)$ intersects $W(T_j)$, i.e., T_j wrote something T_i wrote,

or

$R(T_i)$ intersects $W(T_j)$, i.e., T_i read something T_j wrote

$W(T_i)$ does not intersect $R(T_j)$, and T_i completes its write phase before T_j starts its write phase.

$$W(T_i) \cap R(T_j) = \{\} \quad R(T_i) \cap W(T_j) \neq \{\} \quad W(T_i) \cap W(T_j) \neq \{\}$$



T_j doesn't read anything T_i wrote.

Anything T_j wrote that T_i also wrote will be installed afterwards.

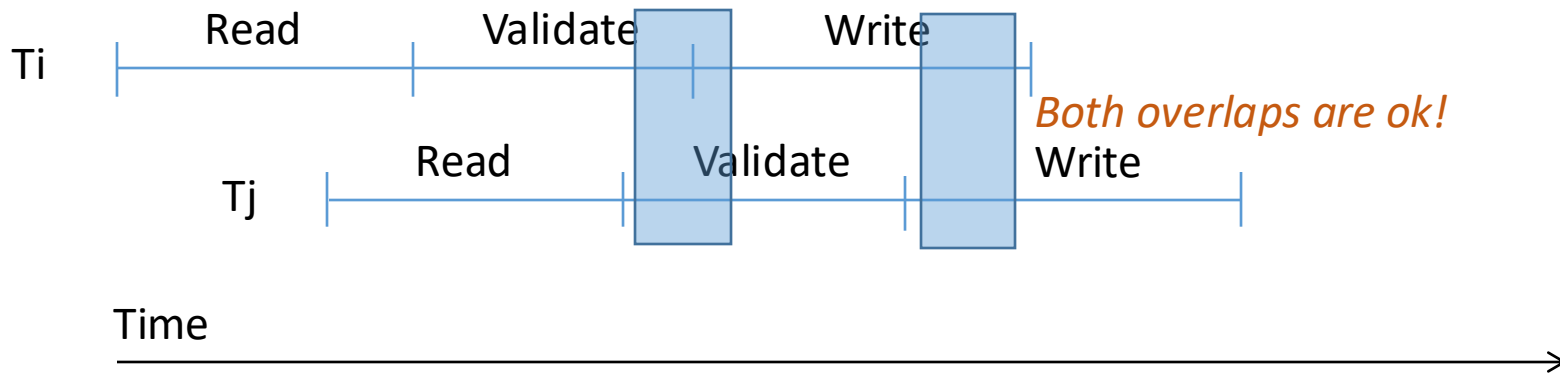
Anything T_i read will not reflect T_j 's writes

Condition 3

R(Ti) intersects W(Tj), i.e., Ti read something Tj wrote

W(Ti) does not intersect R(Tj) or W(Tj), and Ti completes its read phase before Tj completes its read phase.

$$W(T_i) \cap R(T_j) = \{ \} \quad R(T_i) \cap W(T_j) \neq \{ \} \quad W(T_i) \cap W(T_j) = \{ \}$$



T_j doesn't read or write anything T_i wrote (but T_i may read something T_j writes).

T_i definitely won't see any of T_j 's writes, because it finishes reading before T_j starts validation, so T_i ordered before T_j .

T_i will always complete its read phase before T_j b/c xaction IDs assigned after read phase

If no conditions apply, abort!

Restating previous rules, aborts required if:

1) $W(T_i) \cap R(T_j) \neq \{ \}$, and T_i does not finish writing before T_j starts, T_j must abort, because T_j may have only seen some of what T_i wrote

or

2) $W(T_i) \cap (W(T_j) \cup R(T_j)) \neq \{ \}$, and T_j overlaps with T_i validation or write phase, T_j must abort because it needs its writes to all appear after T_i 's writes

12 Candy

Validate Implementation

Validate Implementation

- Several different implementations designed to provide different levels of concurrency during writeback
- Transaction initialization:

```
tnc = 0; // current transaction id

void tbegin {

    read_set = new Set();

    write_set = new Set();
    start_tn = tnc; //the transaction that
                  //finished just before this
                  //one started

}
```

Serial Validation

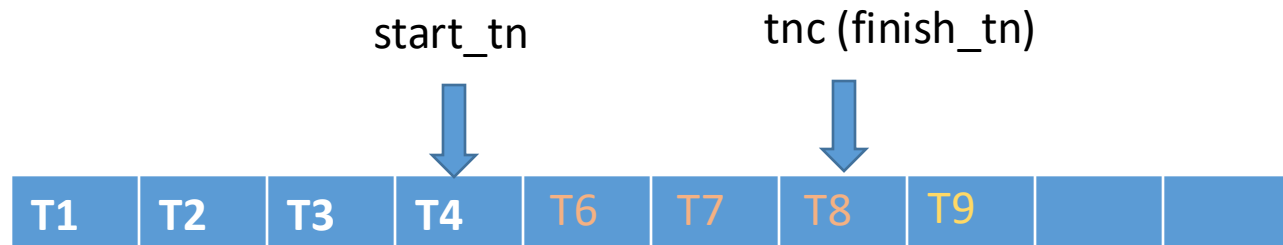
```
validateAndWrite(pastT[], start_tn, my_read_set, my_write_set)
{
    lock();
    int finish_tn = tnc; //prior transaction
    bool valid = true;
    for(int t = start_tn + 1; t <= finish_tn; t++)
        if(pastT[t].write_set intersects with my_read_set)
            valid = false;
    if (valid) {
        write_phase();
        tnc = tnc+1;
        tn = tnc;
    }
    unlock();
}
```

1. $W(T_i) \cap R(T_j) \neq \{ \}$, and T_i does not finish writing before T_j starts, T_j must abort
- ~~2. $W(T_i) \cap (W(T_j) \cup R(T_j)) \neq \{ \}$, and T_j overlaps with T_i validation or write phase, T_j must abort~~

2nd condition doesn't occur because if T_i completes its read phase before T_j , it will also complete its write phase before T_j .

Example

*Last transaction that
validated before this
transaction started*



*Transactions that
validated and wrote while
this transaction was in
read phase*

Have to compare against T6, T7, and T8


Which of the following transactions would serial validation allow to commit, assuming the transactions are concurrent and T_i completes its write phase before T_j starts its write phase $\forall i < j$

Aborts required if:

1) $W(T_i) \cap R(T_j) \neq \{ \}$, and T_i does not finish writing before T_j starts, T_j must abort, because T_j may have only some of what T_i wrote

or

2) $W(T_i) \cap (W(T_j) \cup R(T_j)) \neq \{ \}$, and T_j overlaps with T_i validation or write phase, T_j must abort because it reads/writes to all appear after T_i 's writes

A.	<u>T1</u>		<u>T2</u>	C.	<u>T1</u>	<u>T2</u>
	RA		RC		RA	RA
	WA		WA		WA	RA
	RB					WA
	WB					

Blind write

B.	<u>T1</u>	<u>T2</u>	<u>T3</u>
	RA	RB	RA
	WA	WB	RA
			WC

4 Candy

Goal: assign transaction ids T_1, \dots, T_n , such that this is the serial equivalent order.

When do you assign the Transaction Identifier?

- a) At the beginning of the trx,**
- b) At the start of the validation phase**
- c) At the start of the write phase**
- d) At the end of the write phase**

Transaction Identifier Assignment

- Goal: assign transaction ids T_1, \dots, T_n , such that this is the serial equivalent order
- When should we assign transaction identifiers?
- At start of read phase?
 - No! Would be “pessimistic” – don’t want to pre-assign the transaction order before transactions finish running
 - Long running transactions would have to commit before later short transactions
- Assign at end of read phase, just before validation starts

1 Candy

Is the following schedule permitted by two-phase locking?

T1	T2	T3
READ A		
		READ A
	WRITE A	
		WRITE B
WRITE A		
	WRITE B	
COMMIT		
	COMMIT	
		COMMIT

Is this schedule view serializable? Is it conflict serializable?

2 Candy

Which of the following schedules would rigorous 2PL permit?

<u>T1</u>		<u>T2</u>
RA	X	
WA		
		RA
		WA
RB		
WB		
COMMIT		
		RB
		WB
		COMMIT

<u>T1</u>		<u>T2</u>
RA	✓	
WA		
		RC
		WC
RB		
WB		
COMMIT		
		RB
		WB
		COMMIT

<u>T1</u>		<u>T2</u>
RA	✓	
WA		
		RB
		WB
		COMMIT
RB		
WB		
COMMIT		

T1 does not release lock on A until COMMIT (note under rigorous 2PL T2 would wait so this schedule would not happen)

3 Candy

Consider the following banking schema (a typical design)

account_nb	owner	status	account
1	Tim	open	
2	Karin	open	
...	

id	account_nb	description	amount	bank_trx
1	1	Lunch Money	\$ +200	
3	2	Candy for Tim	\$ -40	
4	1	Ski pass	\$ -50	
...	

You have decided with your girl-/boyfriend to consolidate the accounts by closing ACCOUNT_NB=1 with the following trx:

```
acct_sum = SELECT SUM(amount) FROM bank_trx WHERE account_nb = 1
INSERT (100, 2, "Transfer from Acct 1", acct_sum) INTO bank_trx
INSERT (101, 1, "Transfer to Acct 2", -acct_sum) INTO bank_trx
UPDATE account SET status=closed WHERE account_nb = 1
```

Assuming transfers always check the account status, why is 2-Phase Locking from class not sufficient to ensure that the balance of account 1 is 0 after closing it? How would you prevent the problem?

Final Wrinkle: Phantoms



- T1 scans a range; T2 later inserts into that range
- If T1 scans the range again, it will see a new value

```
T1
BEGIN
  SELECT COUNT(*) FROM emp WHERE SAL > 100
  ...
  SELECT COUNT(*) FROM emp WHERE SAL > 200
END

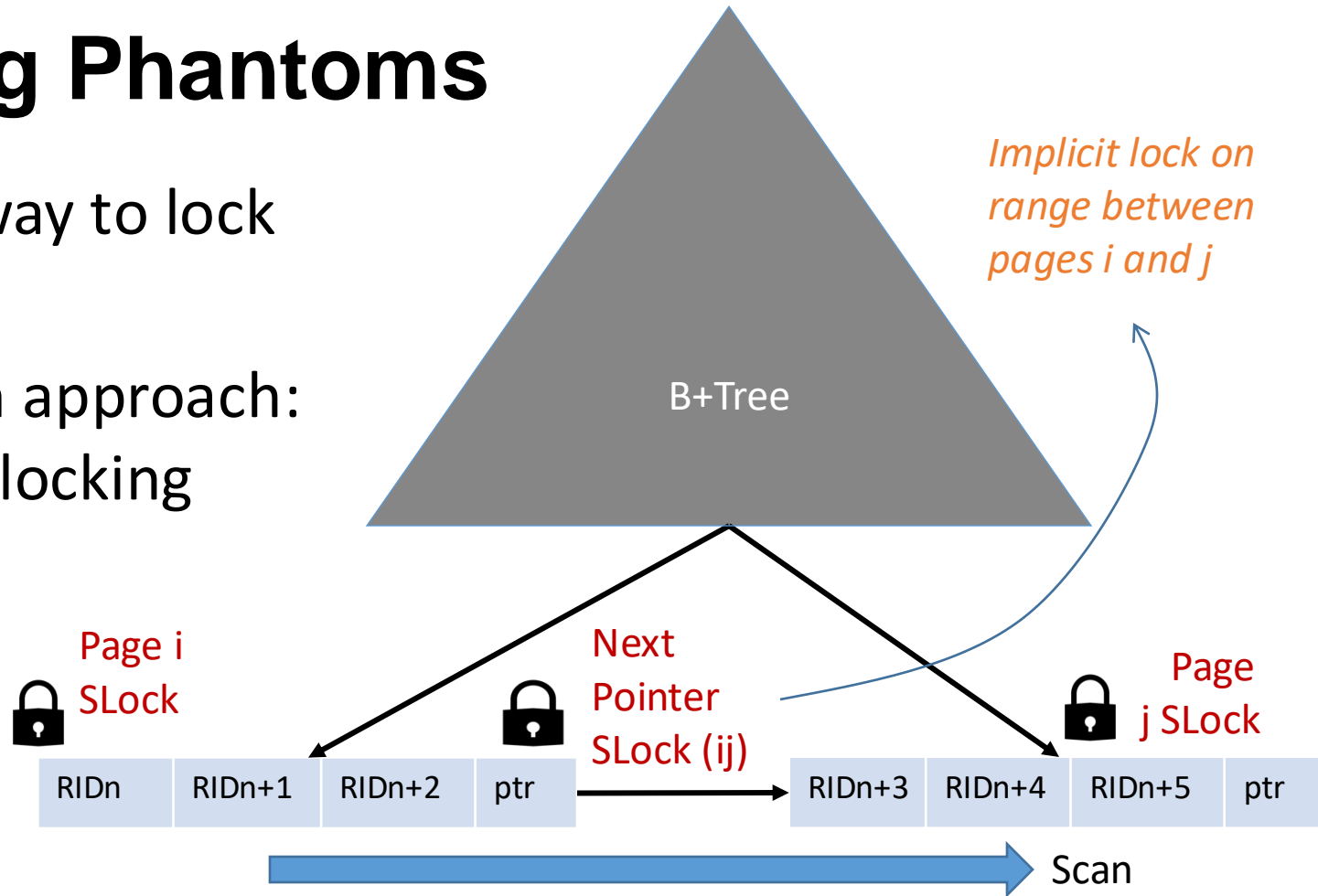
T2
BEGIN
  INSERT INTO EMP VALUES(...,sal=225)
  END
```

A blue arrow originates from the 'INSERT INTO EMP VALUES(...,sal=225)' statement in the T2 block and points to the second 'SELECT COUNT(*) FROM emp WHERE SAL > 200' statement in the T1 block, illustrating the phantom insertion.

If we are just locking, e.g., records, this insertion would be allowed in all 2PL algos we have studied, but is not serializable (since this couldn't happen in a serial execution).

Solving Phantoms

- Need a way to lock ranges
- Common approach: next key locking



On insert(val), Xlock ij next pointer if $val > \max(\text{page } i)$ and $< \min(\text{page}(j))$

Only works for ranges with indexes

For unindexed tables, must read the whole table, so just use a table lock

More details next lecture!

4 Candy

How would you implement 2PL and prevent Phantoms in GoDB without a Btree?

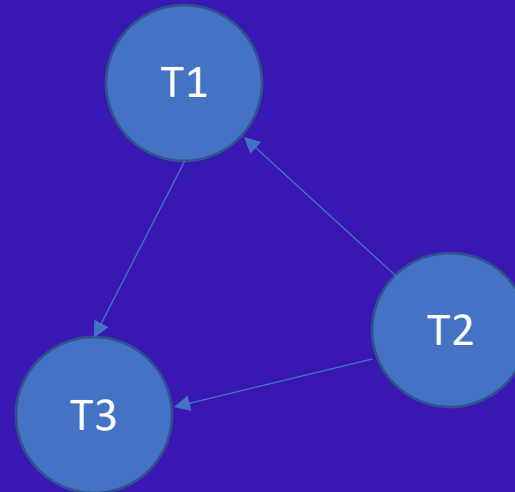
Implementing 2PL

- **GoDB: Lock Table**
 - Buffer pool maintains a table of locks per page
 - Transactions acquire locks on reads/writes of pages
 - Release locks at commit
- Access to lock table will need to be synchronized

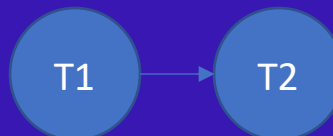
5 Candy

What is the conflict graph for this schedule, and is it serializable?

T1	T2	T3
READ A		
	READ A	
	WRITE B	
		READ B
WRITE A		
READ B		
		WRITE B



For an operation pair that conflicts in T1 and T2, the operation happens first in T1



4 Candy

**What if
Serializability isn't
needed?**

What If Serializability Isn't Needed?

- E.g., application only needs to read committed data
- Databases provide different isolation levels
 - READ UNCOMMITTED
 - Ok to read other transaction's dirty data
 - READ COMMITTED
 - Only read committed values
 - REPEATABLE READS
 - If R1 read $A=x$, R2 will read $A=x \forall A$
- Many database systems default to READ COMMITTED

READ UNCOMMITTED w/ Locking

- If OK reading uncommitted data, no need to check if records that are read are locked
- However, to prevent other transactions from overwriting data, need to hold write locks for the duration of the transaction
- May be OK if, e.g., just reporting some statistic, like number of users or views

READ COMMITTED w/ Locking

- To ensure that a transaction only reads committed values, need to acquire locks before reading
 - If all other transactions hold write locks (as in READ UNCOMMITTED), it will never read a dirty value
- Since we don't care about always reading the same value, OK to release locks after a value is read
- As in READ UNCOMMITTED, write locks still need to be held for the duration of the transaction

Question: Is this schedule permitted under (a) read uncommitted, (b) read committed, (c) repeatable read

T1
RA
WA
COMMIT

T2

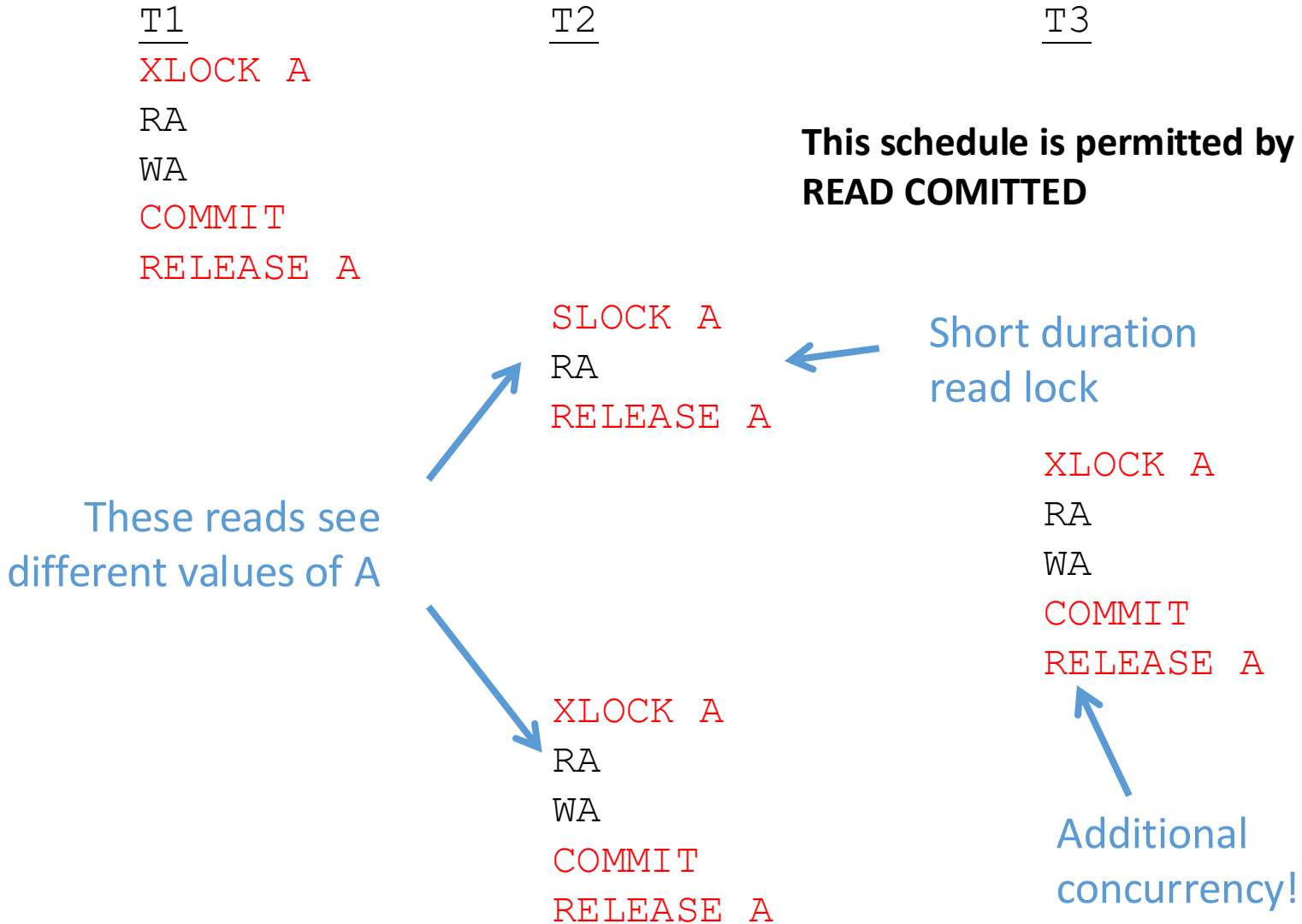
RA

T3

RA
WA
COMMIT

RA
WA
COMMIT

READ COMMITTED Example



REPEATABLE READ w/ Locking

- If we want to always read the same value, need to hold read locks for transaction duration
- So how is this different from SERIALIZABLE?
- SERIALIZABLE also needs to prevent phantoms



6 Candy

**Multi-version concurrency
control (MVCC) /
Snapshot Isolation (SI)**

REPEATABLE READ vs SERIALIZABLE

- Some systems, e.g., Postgres implement REPEATABLE READ through a different mechanism based on database snapshots taken at the start of transaction
 - Called “multiversion concurrency control” – yet another way of achieving isolation!
- This has other problems besides phantoms – so called “read skew anomalies”
 - See: <https://www.cockroachlabs.com/blog/what-write-skew-looks-like/>

Snapshot Isolation

- When a TA starts it receives a timestamp, T.
- All reads are carried out *as of* the DB version of T.
 - Need to keep historic versions of all objects!!!
- All writes are carried out in a separate buffer.
 - Writes only become visible after a commit.
- When TA commits, DBMS checks for conflicts
 - Abort TA1 with timestamp T1 if exists TA2 such that
 - TA2 committed after T1 and before TA1
 - TA1 and TA2 updated the same object

UNDER SI

A) Does T2 abort or commit?

B) Does T3 abort or commit?

T1	T2	T3
		R(Y)
W(Y) Commit		
	Start R(X) R(Y)	
		W(X) W(Z) Commit/Abort
	R(Z) R(Y) W(X) Abort?/Commit?	

Solution

- A transaction T1 executing with Snapshot Isolation
 - takes snapshot of committed data at start
 - always reads/modifies data in its own snapshot
 - updates of concurrent transactions are not visible to T1
 - writes of T1 complete when it commits
 - **First-committer-wins rule:**
 - Commits only if no other concurrent transaction has already written data that T1 intends to write.

Concurrent updates not visible
 Own updates are visible
 Not first-committer of X
 Serialization error, T2 is rolled back

T1	T2	T3
		R(Y) → v0
W(Y := v1) Commit		
	Start R(X) → v0 R(Y) → v1	
		W(X := v2) W(Z := v1) Commit
	R(Z) → v0 R(Y) → v1 W(X := v3) Abort	

8 Candy

For what type of workloads should you use 2PL for what type of workloads SI?

8 Candy

Create a real-world-inspired example, (e.g., from banking, hotel booking, ...) for which two transactions executed with Snapshot Isolation guarantees would violate serializability.

- Observation: Snapshot isolation does not prevent Write-Read conflicts, since it doesn't check whether it read something another transaction wrote
- This leads to so-called write-skew, i.e.:

<u>T1</u>	<u>T2</u>
RX	
	RY
WY	
	WX

- Neither transaction saw the other's write; this would not be permitted under serializability

Real world example:

Employee X and employee Y are signing up for shifts
They hate each other and don't want to work together
EY checks that EX isn't working (RX)
EX checks that EY isn't working (RY)
Both update their schedule to work on the same day

8 Candy

Does SI isolation have the phantom problem?

Under 2-Phase Locking, the following locks got requested from the Lock Manager. Is there a deadlock? Can you think of a way to avoid deadlocks?

T1	T2	T3
RA	RB	RC
WB	RC	RA

time	Type	TID
t1	Req. ReadLock A	T1
t2	Req. ReadLock B	T2
t3	Req. ReadLock C	T3
t4	Req. ReadLock A	T3
t5	Req. WriteLock B	T1
t6	Req. ReadLock C	T2