

# **6.5830 Lecture 13**

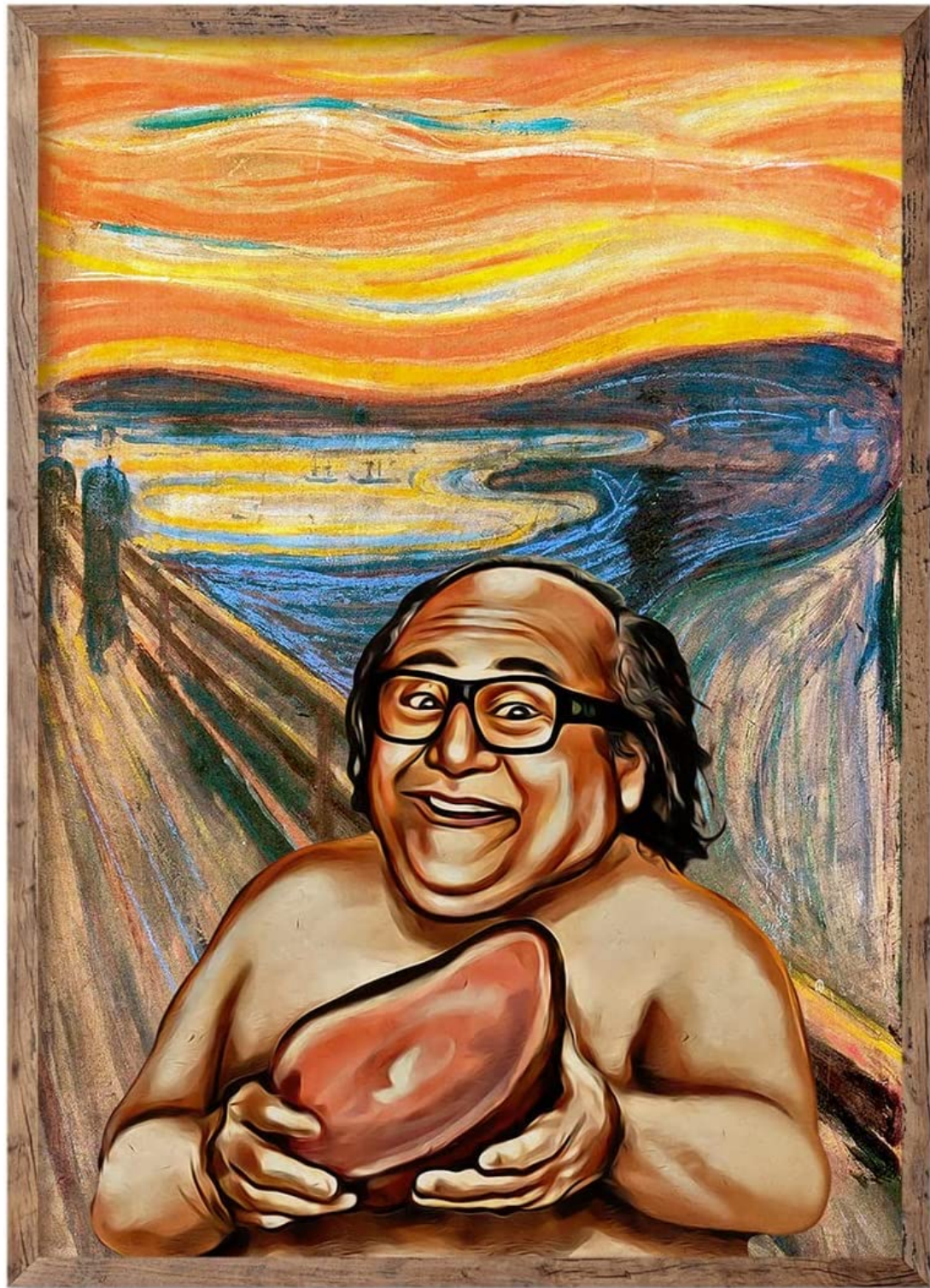
## **Two-phase Locking Recap**

### **Optimistic Concurrency Control**

October 26, 2022

Quiz 1 Back; Mean: 81, Median: 83, Std. Dev: 10

Lab 4 Out



# Transactions

- Group related sequence of actions so they are “all or nothing”
  - If the system crashes, partial effects are not seen
  - Other transactions do not see partial effects
- A set of implementation techniques that provides this abstraction with good performance

# ACID Properties of Transactions

- **A** tomicity – many actions look like one; “all or nothing”
- **C** onsistency – database preserves invariants
- **I** solation – concurrent actions don’t see each other’s results
- **D** urability – completed actions in effect after crash (“recoverable”)

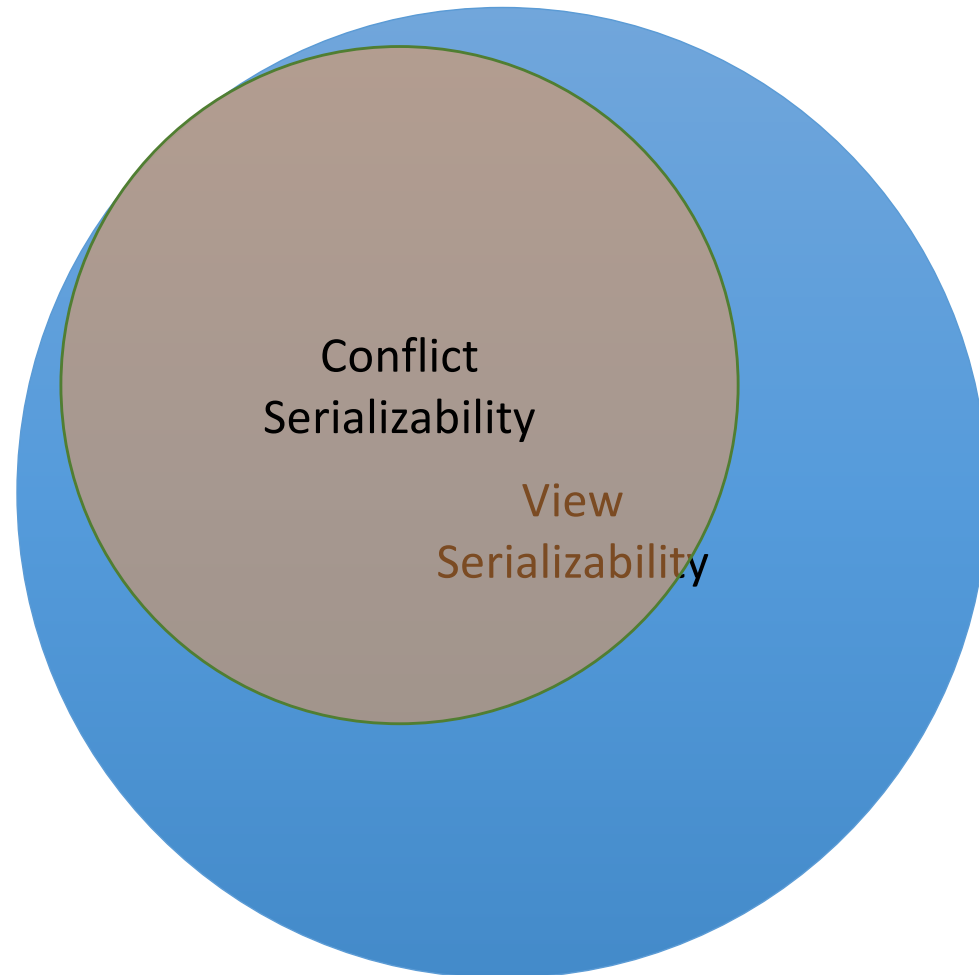
# Atomicity/2PL Recap

- Goal: given a set of transactions, consisting of multiple operations, schedule them such that they are equivalent to some serial execution of those transactions.

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

*Serially equivalent* to T1 then T2

# Testing for Serializability



*Any schedule that is conflict serializable is view serializable, but not vice-versa*

# Precedence Graph for Testing Conflict Serializability

Given transactions  $T_i$  and  $T_j$ ,  
Create an edge from  $T_i \rightarrow T_j$  if:

- $T_i$  reads/writes some  $A$  before  $T_j$  writes  $A$

$$RA_{T_i} < WA_{T_j} \text{ or } WA_{T_i} < WA_{T_j}$$

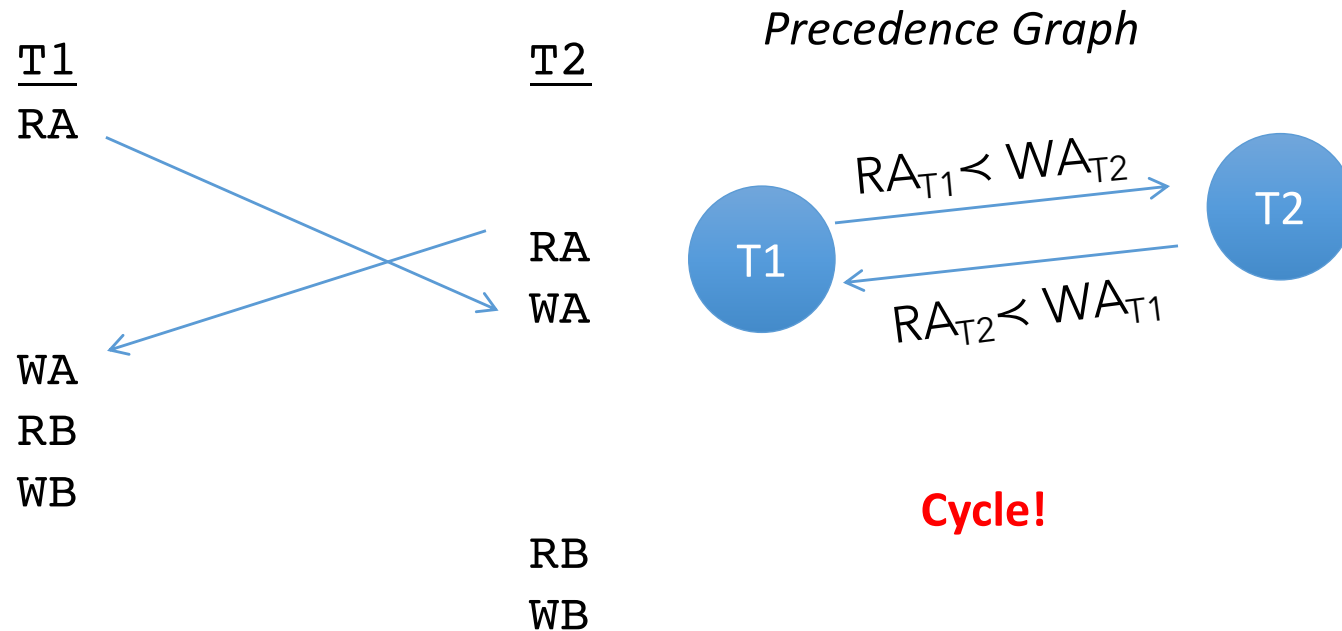
*or*

- $T_i$  writes some  $A$  before  $T_j$  reads  $A$

$$WA_{T_i} < RA_{T_j}$$

If there are cycles in this graph, schedule is not conflict serializable

# Non-Serializable Example



Create an edge from  $T_i \rightarrow T_j$  if:

$T_i$  reads/writes some A before  $T_j$  writes A, or

$$RA_{T_i} < WA_{T_j} \text{ or } WA_{T_i} < WA_{T_j}$$

$T_i$  writes some A before  $T_j$  reads A

$$WA_{T_i} < RA_{T_j}$$



# Two Phase Locking (2PL) Protocol

- Before every read, acquire a shared lock
- Before every write, acquire an exclusive lock (or "upgrade") a shared to an exclusive lock
- Release locks only after last lock has been acquired, and ops on that object are finished

# Refining 2PL

- Problems:
  - Deadlocks
  - Cascading Aborts
  - How do we know when we are done with all operations on an object?

# Rigorous Two-Phase Locking Protocol

- Before every read, acquire a shared lock
- Before every write, acquire an exclusive lock (or "upgrade") a shared to an exclusive lock
- Release locks only after the transaction commits
- Ensures cascadeless-ness, and
- *Commit order = serialization order*

# 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 * FROM emp WHERE SAL > 100
  ...
  SELECT * FROM emp WHERE SAL > 200
END

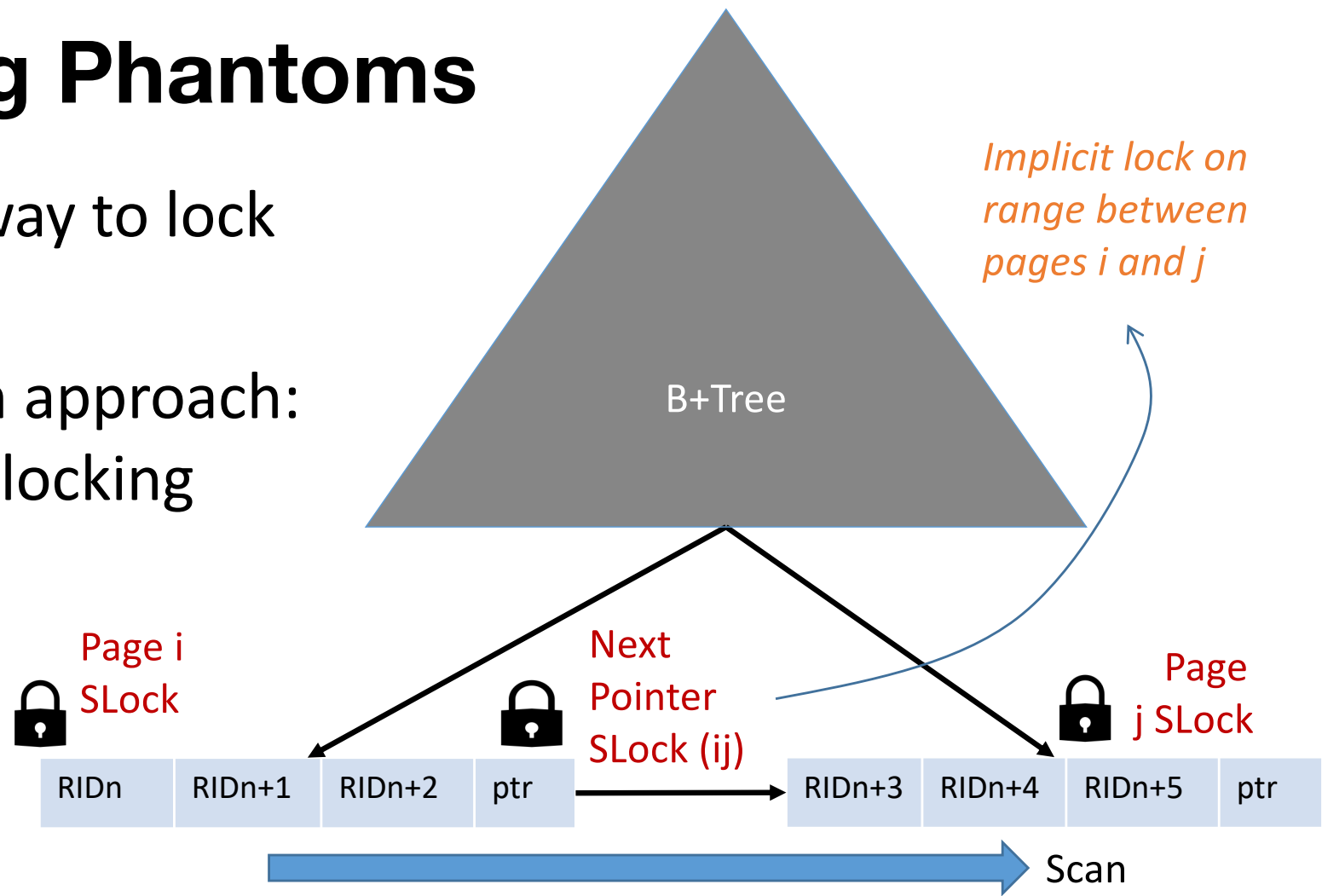
T2
BEGIN
  INSERT INTO EMP VALUES(...,sal=225)
END
```

A blue arrow points from the 'INSERT INTO EMP VALUES(...,sal=225)' statement in the T2 transaction to the second 'SELECT \* FROM emp WHERE SAL > 200' statement in the T1 transaction, illustrating the phantom read scenario.

*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!*

# Implementing 2PL

- **SimpleDB: 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, but not using special “lock” objects to represent the locks.

# Study Break

- Which of the following schedules would rigorous 2PL permit?

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

*T1 does not release lock on A until COMMIT*

# 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>

# 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

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

# 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

# Validation Rules

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):

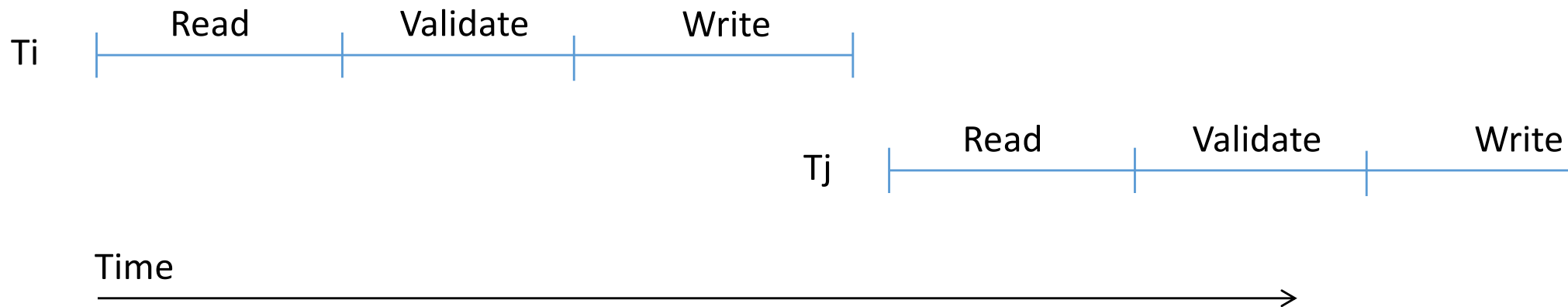
- 1)  $T_i$  completes its write phase before  $T_j$  starts its read phase
- 2)  $W(T_i)$  does not intersect  $R(T_j)$ , and  $T_i$  completes its write phase before  $T_j$  starts its write phase.
- 3)  $W(T_i)$  does not intersect  $R(T_j)$  or  $W(T_j)$ , and  $T_i$  completes its read phase before  $T_j$  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



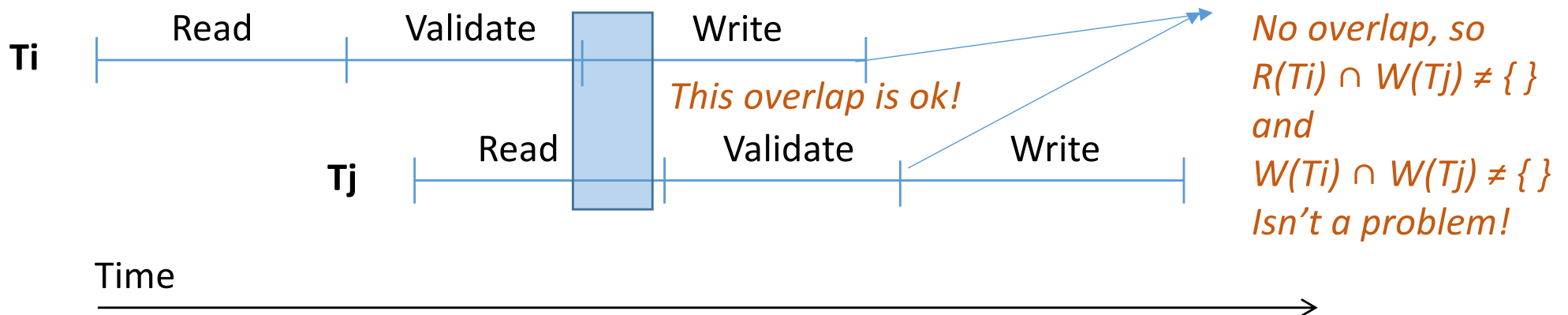
**Don't overlap at all.**

# 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.

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

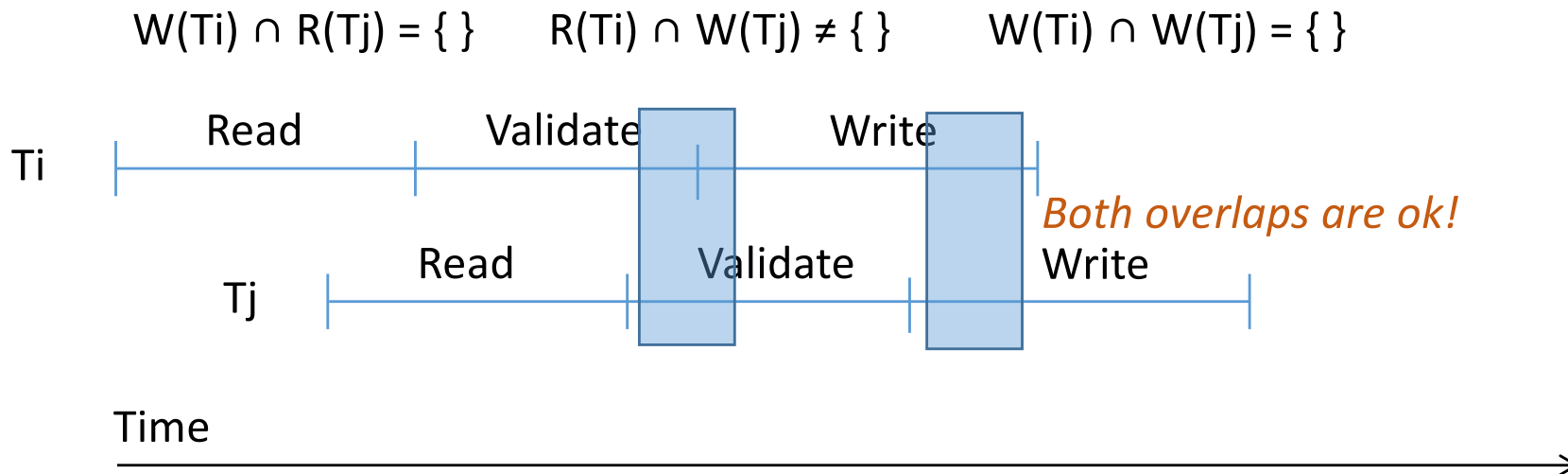


**$T_j$  doesn't read anything  $T_i$  wrote.  
Anything  $T_j$  wrote that  $T_i$  also wrote will be installed afterwards.  
Anything  $T_i$  read will not reflect  $T_j$ 's writes**

# Condition 3

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

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



**$T_j$  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

# 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

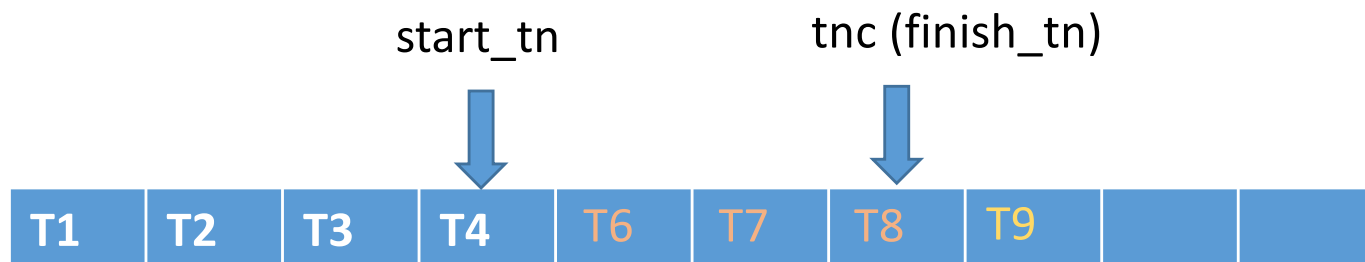
```
validateAndWrite(pastT[], start_tn, my_read_set, my_write_set)
{
    lock();
    int finish_tn = tnc; //prior transaction
    bool valid = true;
    for(int t = start_tn + 1; t <= finish_tn; t++)
        if(pastT[t].write_set intersects with my_read_set)
            valid = false;
    if (valid)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

# Study Break




- 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 needs to write to all appear after  $T_i$ 's writes

<p><b>A.</b></p> <table border="0"> <tr> <td><u>T1</u></td> <td></td> <td><u>T2</u></td> </tr> <tr> <td>RA</td> <td></td> <td>RC</td> </tr> <tr> <td>WA</td> <td></td> <td>WA</td> </tr> <tr> <td>RB</td> <td></td> <td></td> </tr> <tr> <td>WB</td> <td></td> <td></td> </tr> </table>	<u>T1</u>		<u>T2</u>	RA		RC	WA		WA	RB			WB			<p><i>Blind write</i></p>	<p><b>C.</b></p> <table border="0"> <tr> <td><u>T1</u></td> <td><u>T2</u></td> </tr> <tr> <td>RA</td> <td>RA</td> </tr> <tr> <td>WA</td> <td><del>RA</del></td> </tr> <tr> <td></td> <td>WA</td> </tr> </table>	<u>T1</u>	<u>T2</u>	RA	RA	WA	<del>RA</del>		WA
<u>T1</u>		<u>T2</u>																							
RA		RC																							
WA		WA																							
RB																									
WB																									
<u>T1</u>	<u>T2</u>																								
RA	RA																								
WA	<del>RA</del>																								
	WA																								

<p><b>B.</b></p> <table border="0"> <tr> <td><u>T1</u></td> <td><u>T2</u></td> <td><u>T3</u></td> </tr> <tr> <td>RA</td> <td>RB</td> <td>RA</td> </tr> <tr> <td>WA</td> <td>WB</td> <td><del>RA</del></td> </tr> <tr> <td></td> <td></td> <td>WC</td> </tr> </table>	<u>T1</u>	<u>T2</u>	<u>T3</u>	RA	RB	RA	WA	WB	<del>RA</del>			WC
<u>T1</u>	<u>T2</u>	<u>T3</u>										
RA	RB	RA										
WA	WB	<del>RA</del>										
		WC										



# Problems w/ Serial Validation

- Only one transaction can commit at a time
  - Severely limits transaction throughput
- Idea:
  - Reduce critical sections
  - Allow multiple transactions to write at a time
  - Check that concurrent writers didn't intersect committing transactions write set (condition #2)

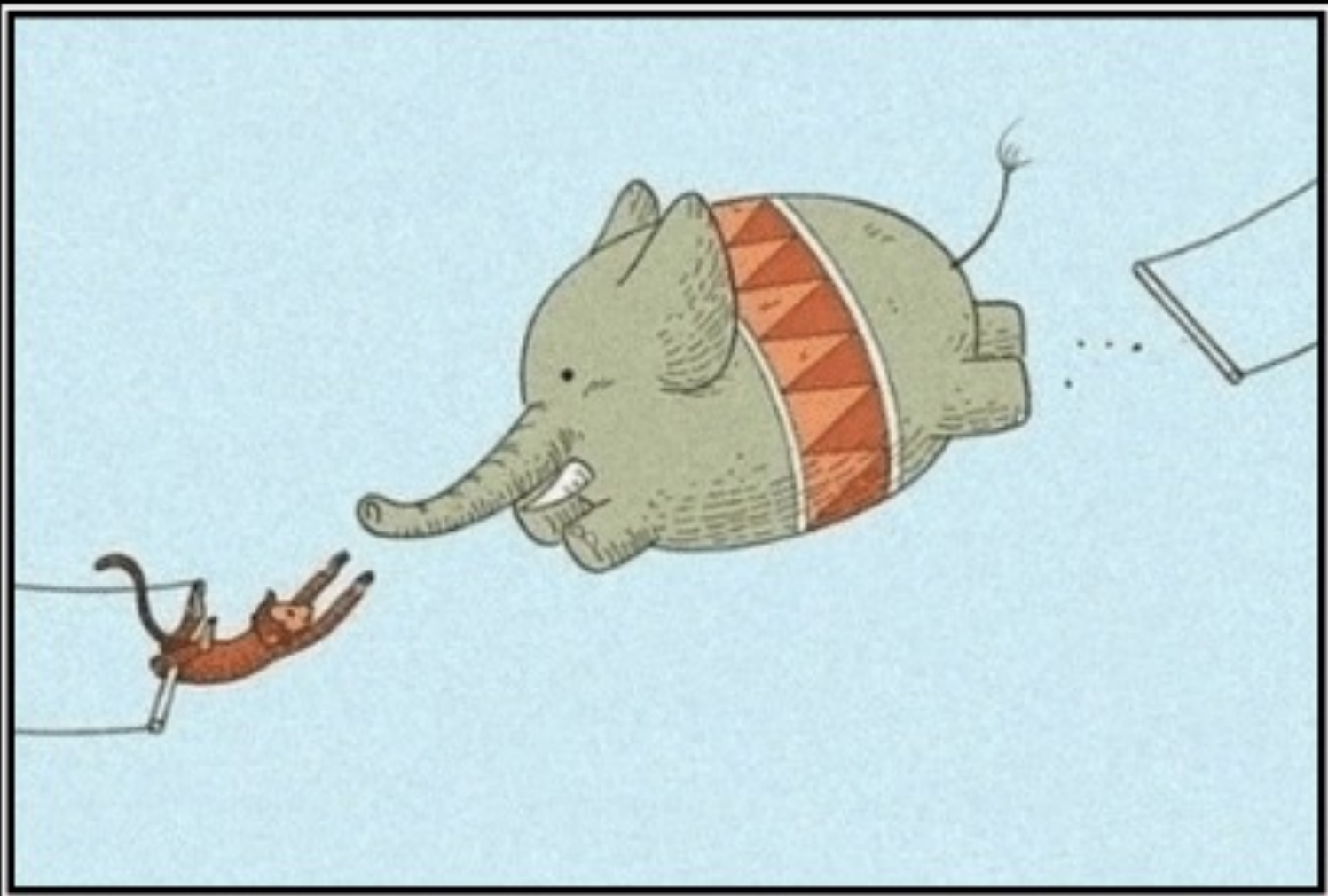
# Parallel Validation

```
List<Transaction> active = new List();
validateAndWrite(pastT[], start_tn,
                 my_read_set, my_write_set) {
    lock();
    int finish_tn = tnc;
    //transactions writing concurrently
    List<Transaction> finish_active = active.copy();
    active.append(me.id);
    unlock();
    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; Condition 1
    for (id t : finish_active)
        if(pastT[t].write_set intersects with Condition 1 &2
           (my_read_set U my_write_set))
            valid = false;
```

```
if(valid) {
    //note that concurrent writes
    //are all to different objects
    write_phase();
    lock();
    tnc = tnc+1;
    tn = tnc;
    active.remove(me.id);
    unlock();
} else {
    lock();
    active.remove(me.id);
    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

*We must check both!*



# OPTIMISM

Not always the best option.

# 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

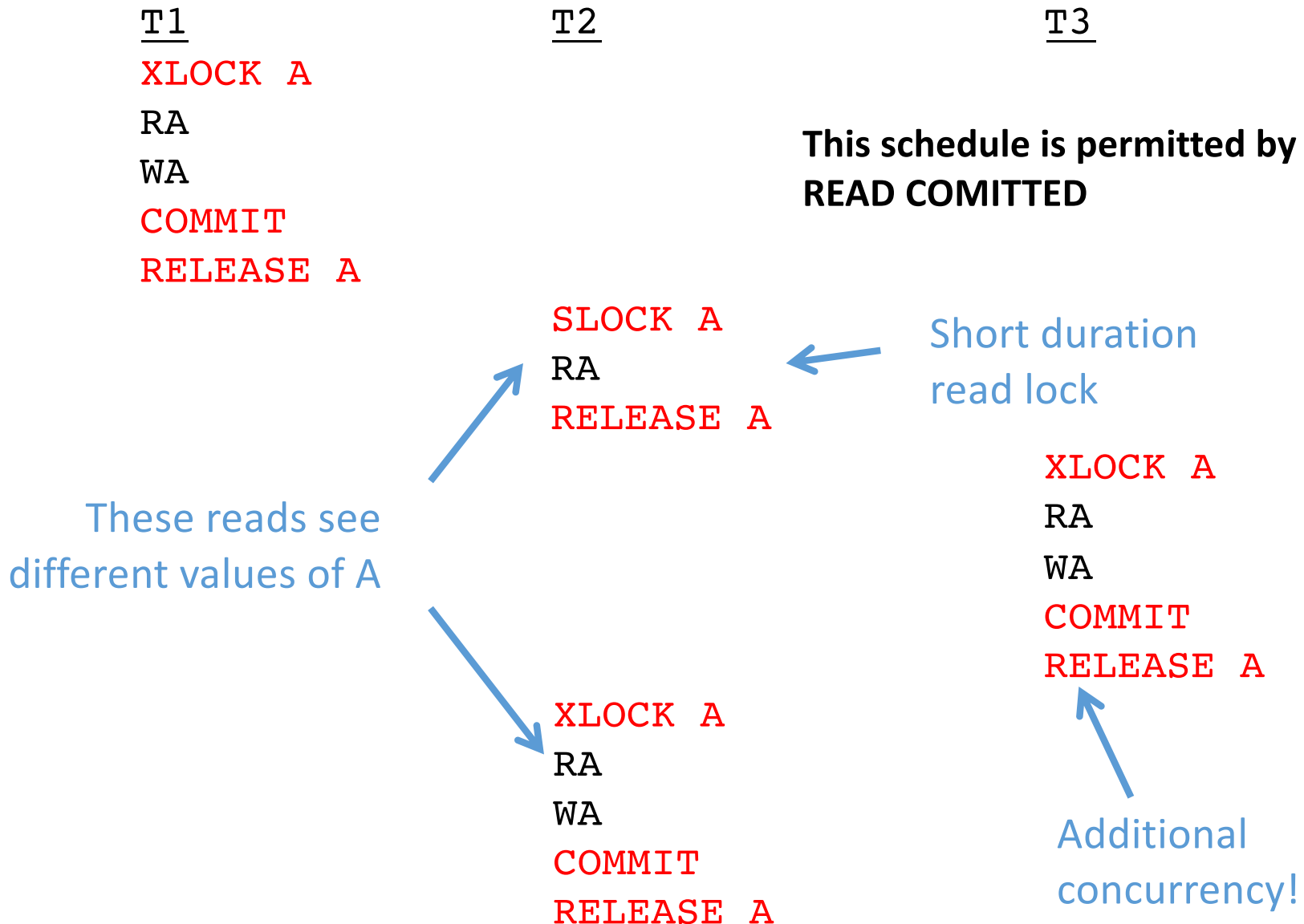
# READ UNCOMMITTED w/ Locking

- If OK reading uncommitted data, no need to check if records that are read are locked
- However, to prevent other transactions from seeing dirty 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 doesn'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

# 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





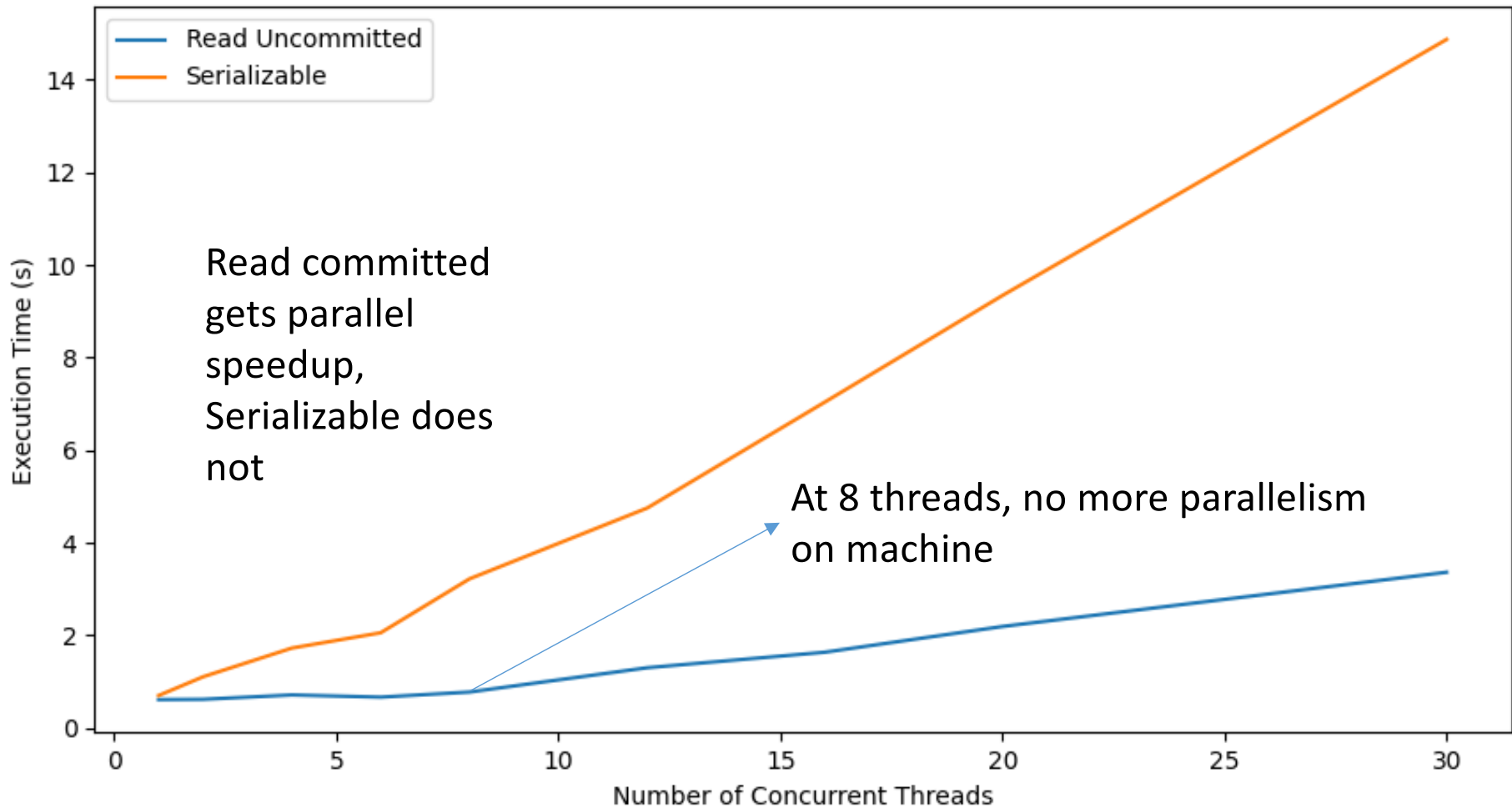
# REPEATABLE READ vs SERIALIZABLE

- 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/>

# Demo

- Table with 1000 rows t (a int, b int, c int)
- Transactions that:
  - Count all rows
  - Update a random row
- Threads that run this transaction 1000 times
- Vary number of threads (level of concurrency) from 1 to 20
  - READ UNCOMMITTED doesn't have conflicts between counts and updates
  - SERIALIZABLE xactions may have to block/abort on update due to concurrent readers
- Measure runtime w/ READ UNCOMMITTED and SERIALIZABLE on Postgres
- Postgres multi-version concurrency can abort transactions on conflict; retry aborts

# Results



# Summary

- Optimistic concurrency control provides another way to provide serializability
- Good for low contention workloads
- Basis for many modern high throughput transactional systems
  
- Reduced consistency levels that lock fewer records common in practice
  - Permit greater concurrency at price of serializability