

# TRANSACTIONS: OPTIMISTIC

CSE 4/562: Database Systems | Lecture 20

---

**DB. Sys.: T.C.B.:** Ch. 18.8-18.9

## **Schedule**

An order over the reads and writes of a transaction

## **Serial Schedule**

A schedule with no interleaving

## **Serializable Schedule**

A schedule guaranteed to produce the same output as a serial schedule

## **Conflict Serializable Schedule**

A schedule with the same “happens before” constraints as a serial schedule

## Acquire Phase

A transaction in the acquire phase may...

- Acquire Shared or eXclusive locks for any object
- Upgrade any shared lock into an exclusive one
- Read any object for which it holds a shared or exclusive lock
- Write any object for which it holds an exclusive lock
- Move to the release phase.

## Release Phase

A transaction in the release phase may...

- Read any object for which it holds a shared or exclusive lock
- Write any object for which it holds an exclusive lock
- Release any lock it holds
- End (Commit or Abort)

Any schedule created by 2-phase locking must be conflict serializable.  
(but the reverse is not true)

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓	W(A)		
↓		W(A)	
↓			W(A)
↓	W(B)		
↓		W(B)	
↓			W(B)
↓	W(C)		
↓		W(C)	
↓			W(C)

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓	W(A)		
↓		W(A)	
↓			W(A)
↓	W(B)		
↓		W(B)	
↓			W(B)
↓	W(C)		
↓		W(C)	
↓			W(C)

3 × T1 “happens before” T2  
“happens before” T3

Time	T1	T2	T3
↓	W(A)		
↓		W(A)	
↓			W(A)
↓	W(B)		
↓		W(B)	
↓			W(B)
↓	W(C)		
↓		W(C)	
↓			W(C)

3 × T1 “happens before” T2  
 “happens before” T3

The schedule is conflict-serializable.

Time	T1	T2	T3
↓	W(A)		
↓		W(A)	
↓			W(A)
↓	W(B)		
↓		W(B)	
↓			W(B)
↓	W(C)		
↓		W(C)	
↓			W(C)

3 × T1 “happens before” T2  
 “happens before” T3

The schedule is conflict-serializable.

... but there is no way to create this schedule with 2-Phase Locking.

## **Locking is...**

### **... expensive**

Costs are still incurred even if there are no problematic conflicts.

### **... restrictive**

We don't know what the transaction will do ahead of time, so we sometimes need to block unnecessarily.



We don't know what a transaction will do...



We don't know what a transaction will do... until it does.

**Idea:** Let the transaction do it!

**Idea:** Let the transaction do it!

(then fix anything it breaks, later)

## **Locking**

Pessimistically block transactions before they can act out of order.

## **Snapshot Isolation**

Run transactions on an image of the data, and revert changes if they accidentally ran out-of-order.

## **Timestamp Concurrency Control**

## **Locking**

Pessimistically block transactions before they can act out of order.

## **Snapshot Isolation**

Run transactions on an image of the data, and revert changes if they accidentally ran out-of-order.

## **Timestamp Concurrency Control**

# Snapshot Isolation



## **Phase 1: Read**

- The transaction executes on a private copy of all accessed objects
- The system records each of the transaction's reads and writes

## **Phase 1: Read**

- The transaction executes on a private copy of all accessed objects
- The system records each of the transaction's reads and writes

## **Phase 2: Validate**

- The system checks whether the transaction's writes can be applied without violating isolation.

## **Phase 1: Read**

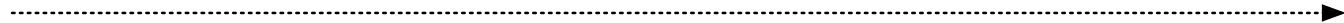
- The transaction executes on a private copy of all accessed objects
- The system records each of the transaction's reads and writes

## **Phase 2: Validate**

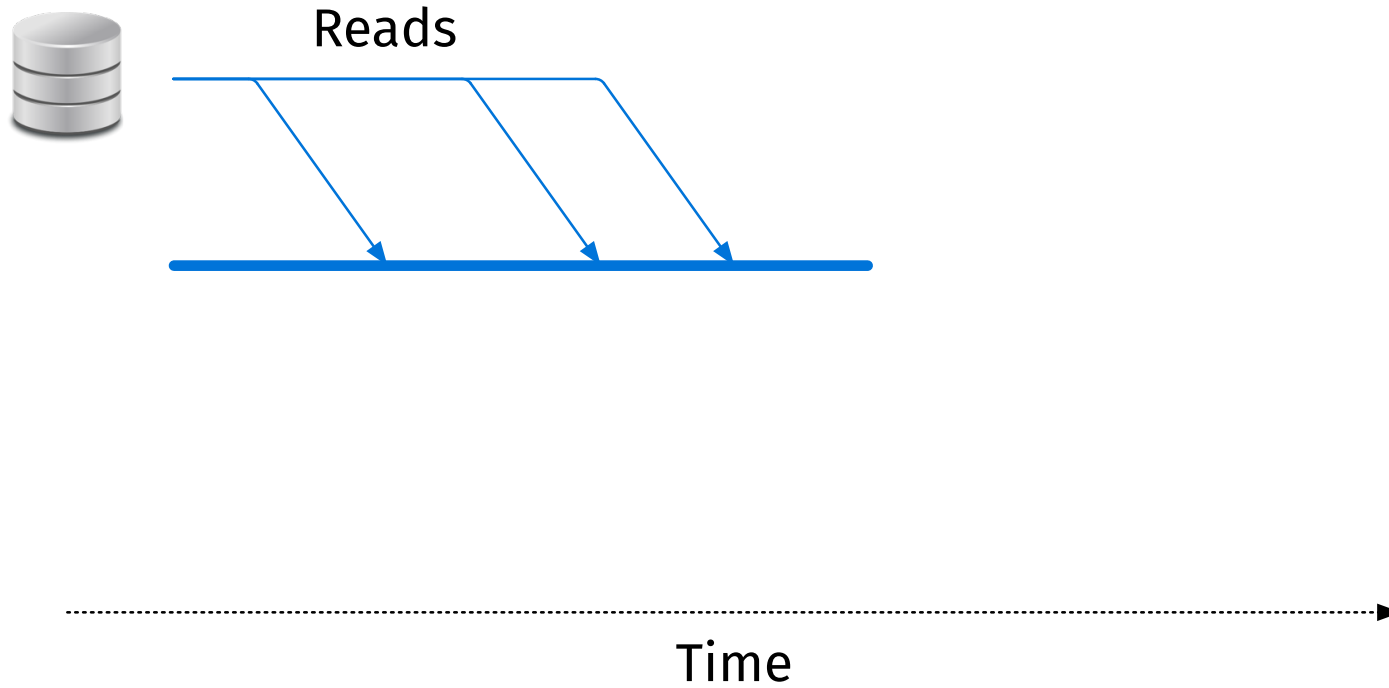
- The system checks whether the transaction's writes can be applied without violating isolation.

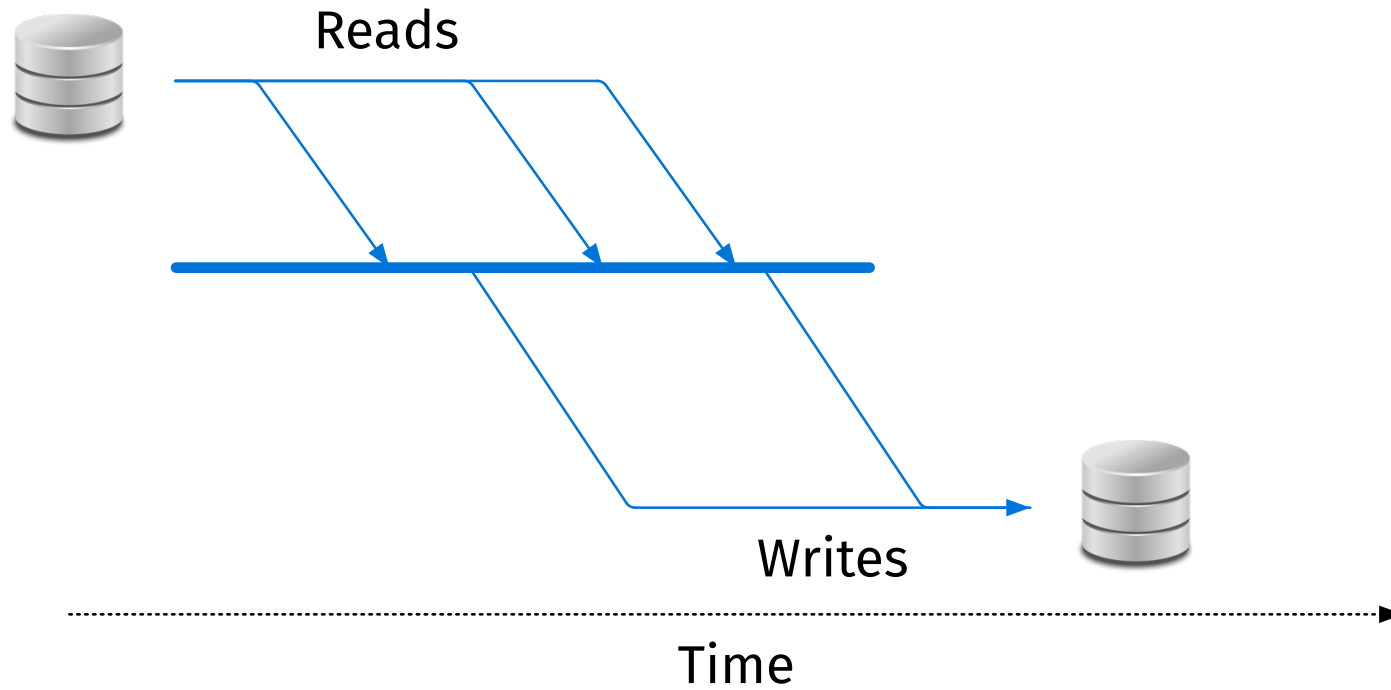
## **Phase 3: Write (Merge)**

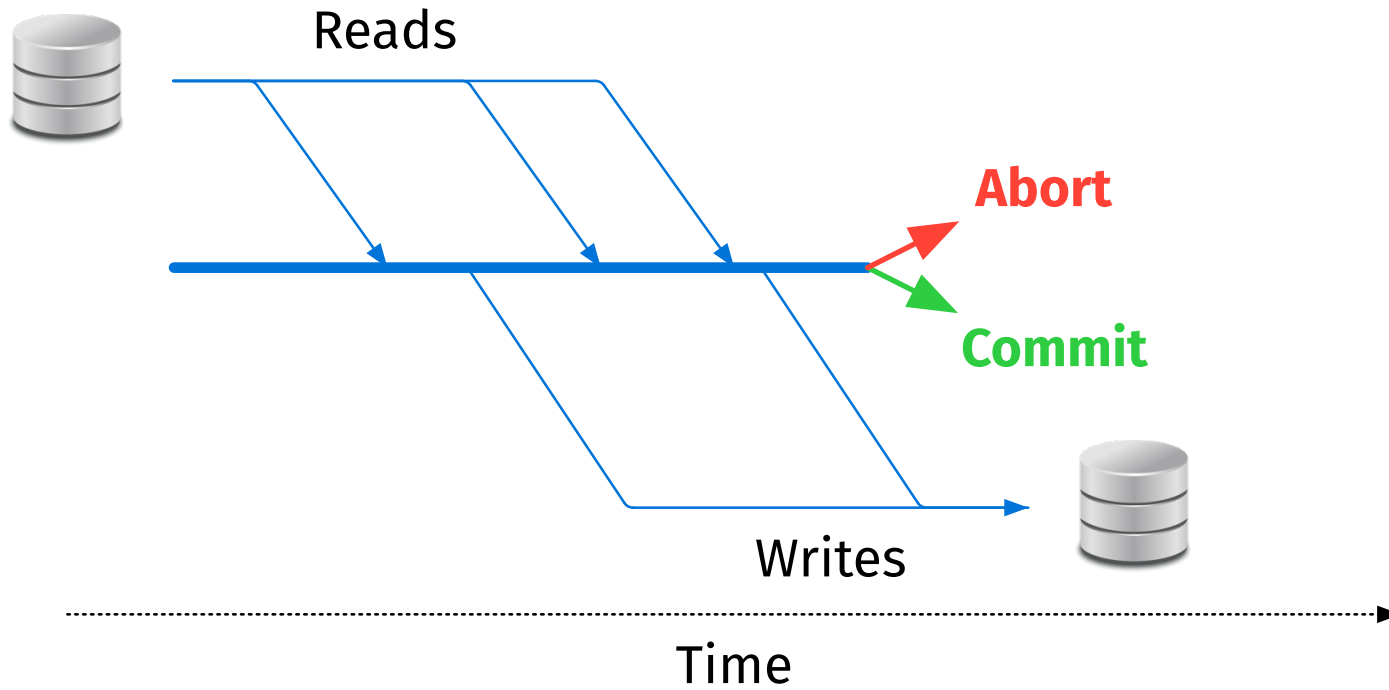
- The system replaces all written objects with the transaction's private copy.

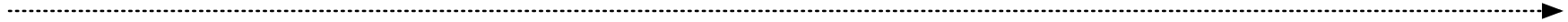


Time

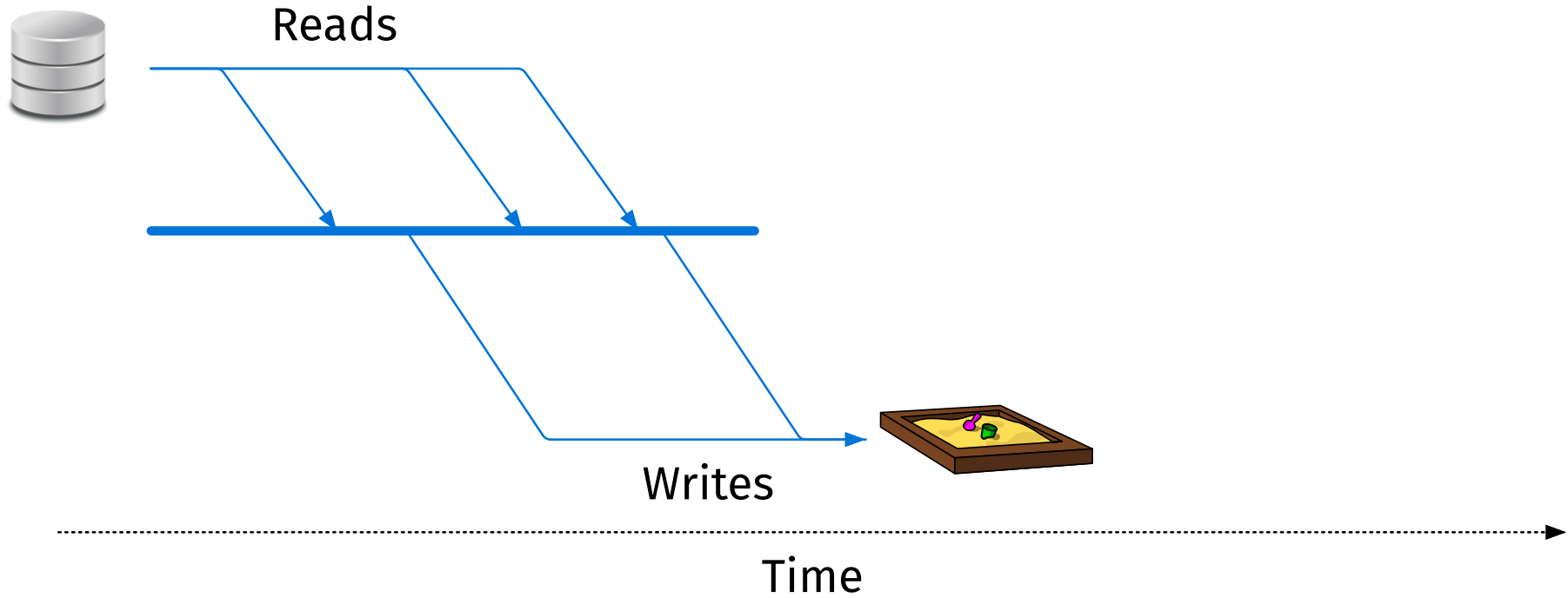


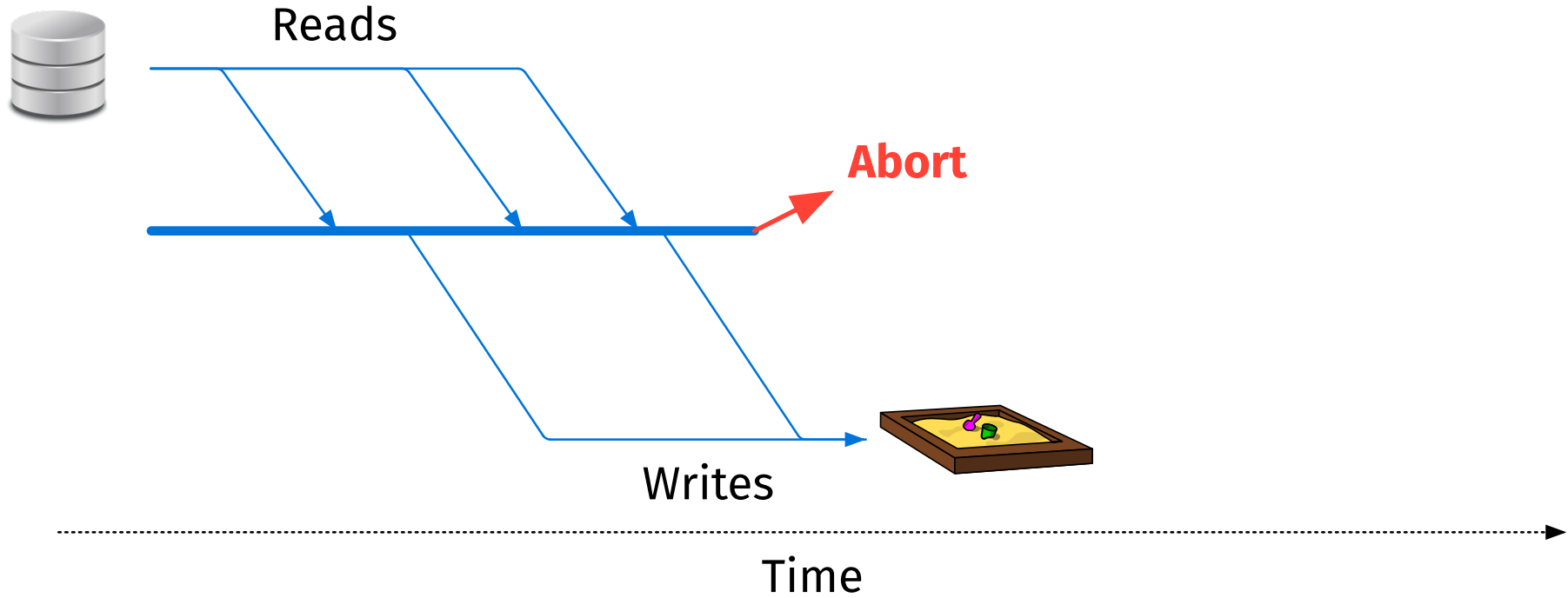


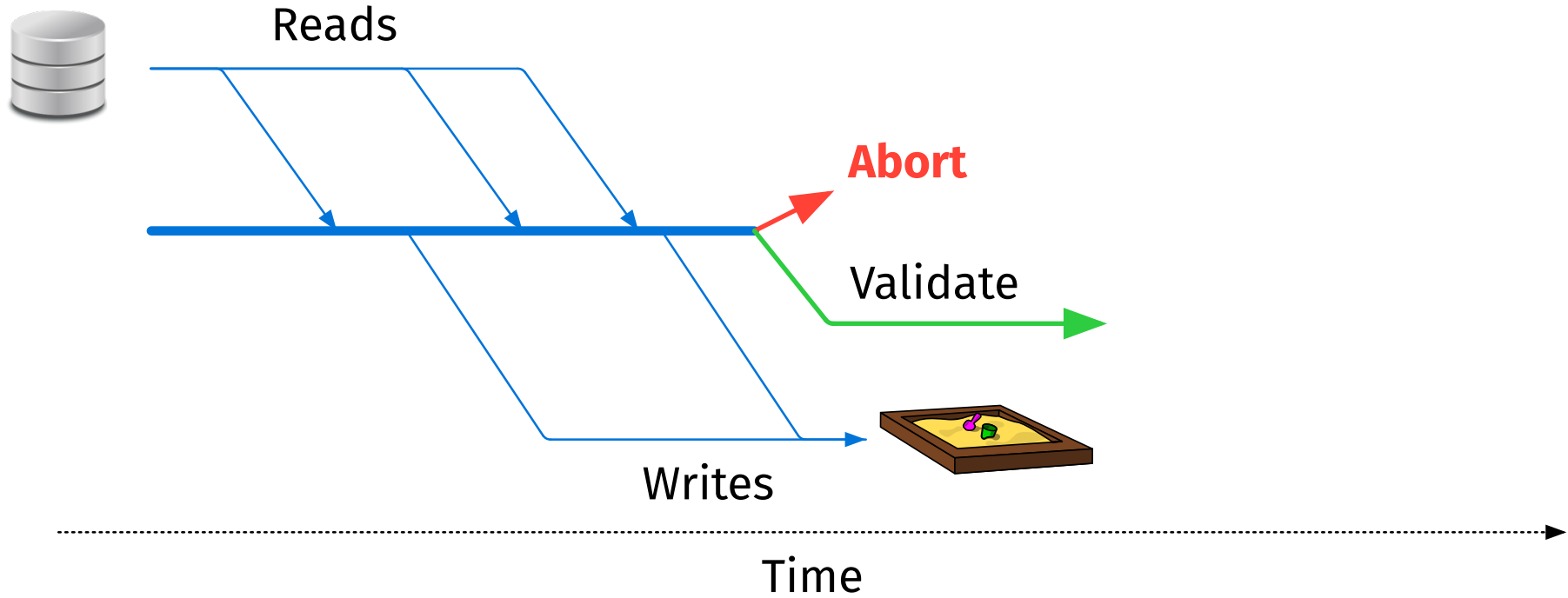


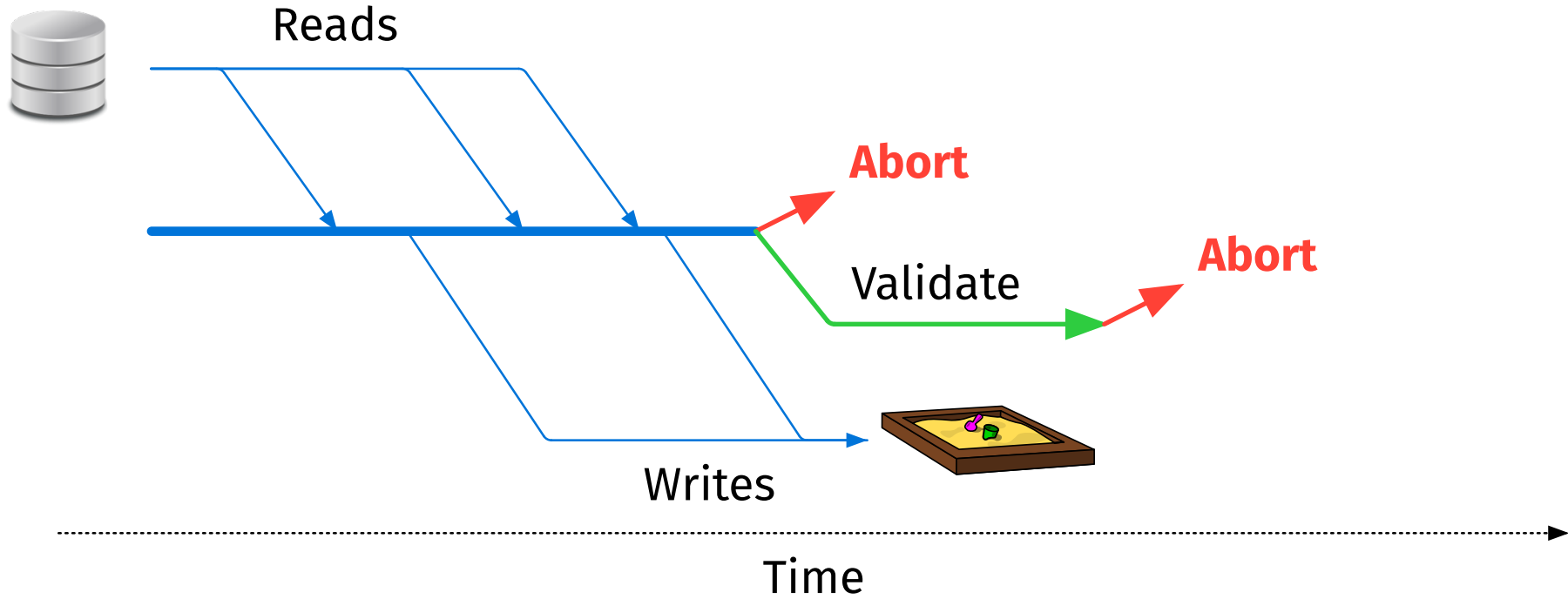


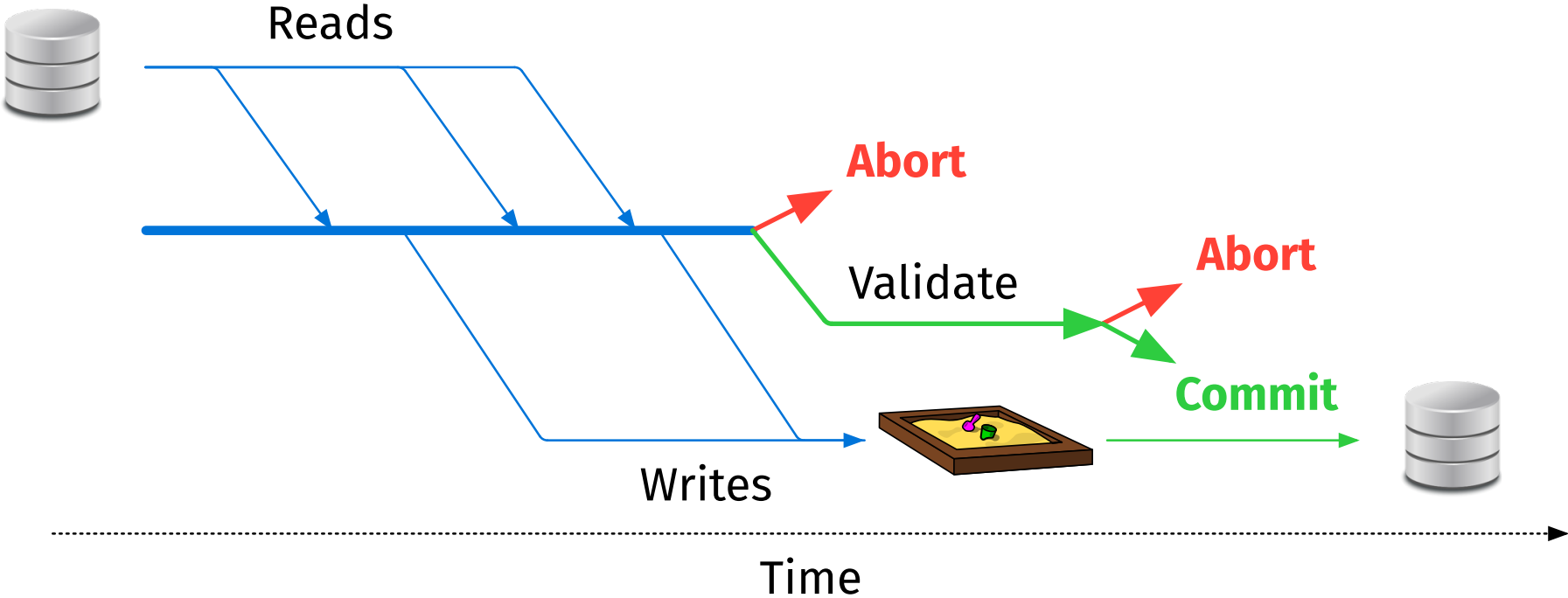
Time

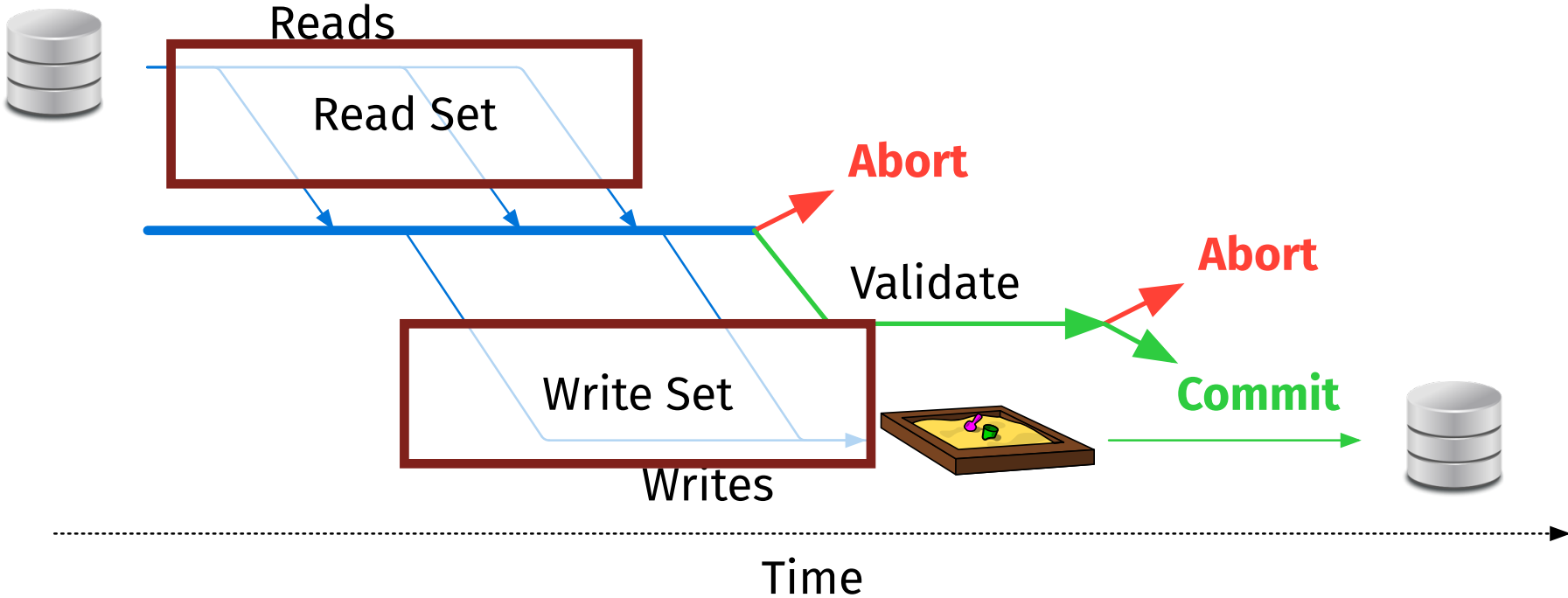












## Read Phase

Run the transaction on a snapshot of the data and collect:

ReadSet( $\mathcal{T}$ )

The set of objects read by  $\mathcal{T}$

WriteSet( $\mathcal{T}$ )

The set of objects written by  $\mathcal{T}$

## **Validation Phase**

1. Pick a serial order  
(e.g., the order in which transactions reach the validation phase)
2. Make sure that the transaction's operations are consistent with this order

T1 scheduled before T2 (So T1's validation step happens before T2's)

**T1 and T2 read the same object**

**T1 reads an object written by T2 (read-write conflict)**

**T2 reads an object written by T1 (write-read conflict)**

**T1 and T2 write the same object (write-write conflict)**

T1 scheduled before T2 (So T1's validation step happens before T2's)

**T1 and T2 read the same object**

Ok!

**T1 reads an object written by T2 (read-write conflict)**

**T2 reads an object written by T1 (write-read conflict)**

**T1 and T2 write the same object (write-write conflict)**

T1 scheduled before T2 (So T1's validation step happens before T2's)

**T1 and T2 read the same object**

Ok!

**T1 reads an object written by T2 (read-write conflict)**

Ok! (T2's write was sandboxed and invisible to T1)

**T2 reads an object written by T1 (write-read conflict)**

**T1 and T2 write the same object (write-write conflict)**

T1 scheduled before T2 (So T1's validation step happens before T2's)

## **T1 and T2 read the same object**

Ok!

## **T1 reads an object written by T2 (read-write conflict)**

Ok! (T2's write was sandboxed and invisible to T1)

## **T2 reads an object written by T1 (write-read conflict)**

Ok if T1's write phase finishes before T2's read phase starts

## **T1 and T2 write the same object (write-write conflict)**

T1 scheduled before T2 (So T1's validation step happens before T2's)

## **T1 and T2 read the same object**

Ok!

## **T1 reads an object written by T2 (read-write conflict)**

Ok! (T2's write was sandboxed and invisible to T1)

## **T2 reads an object written by T1 (write-read conflict)**

Ok if T1's write phase finishes before T2's read phase starts

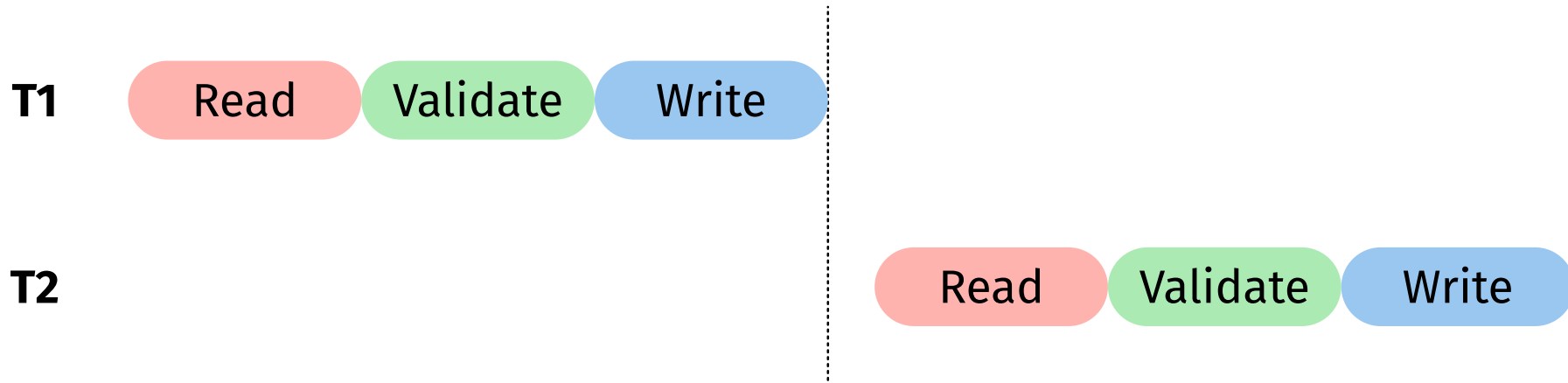
## **T1 and T2 write the same object (write-write conflict)**

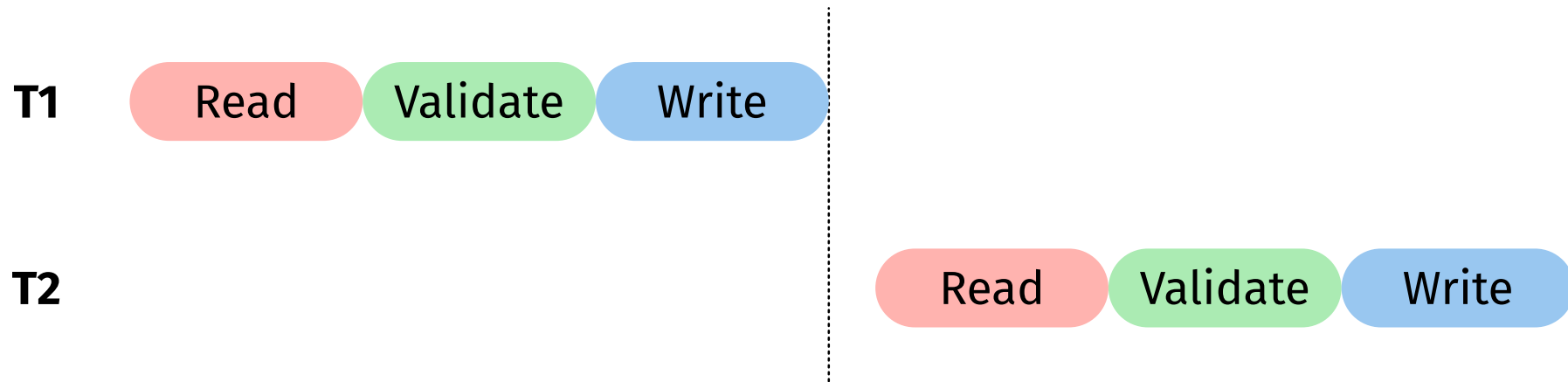
Ok if T1's write phase finishes before T2's write phase starts

Depending on how transaction phases overlap, we need to check for different conflicts.

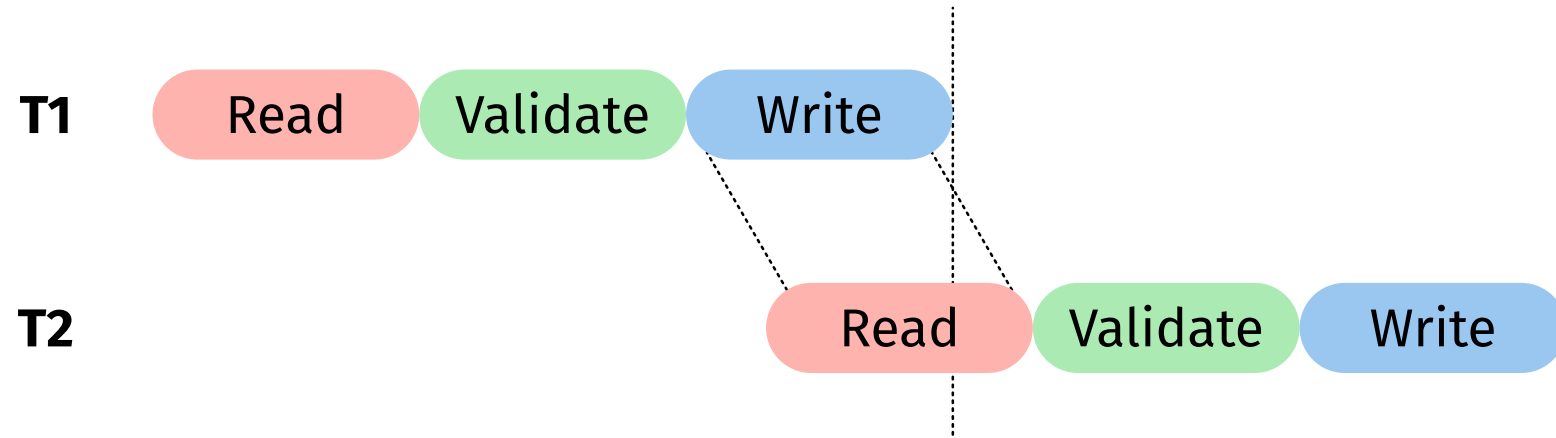
**Note:** The validation and write phases are not instantaneous.

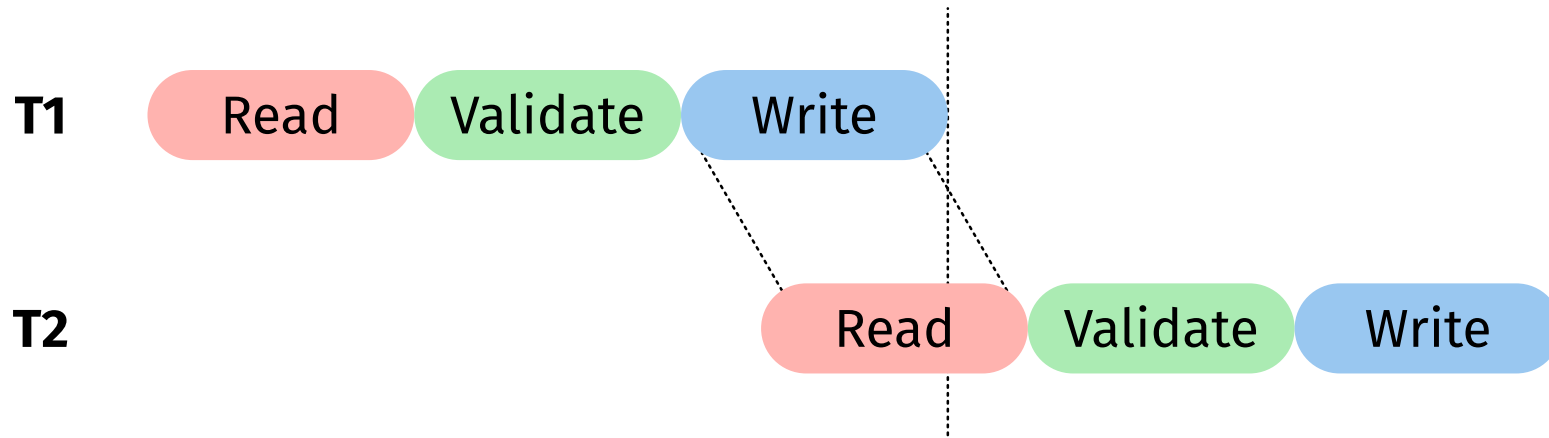
- Validation (and serial order selection) are a critical section, only one transaction is permitted at a time.
- (Compatible) write phases can proceed concurrently.





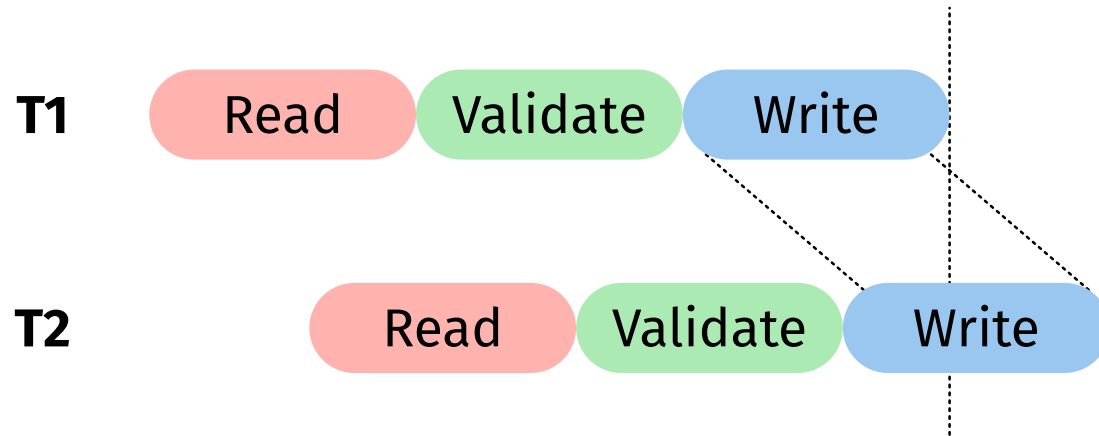
If **T1**'s write phase ends before **T2** starts, then **T2** sees all of **T1**'s writes

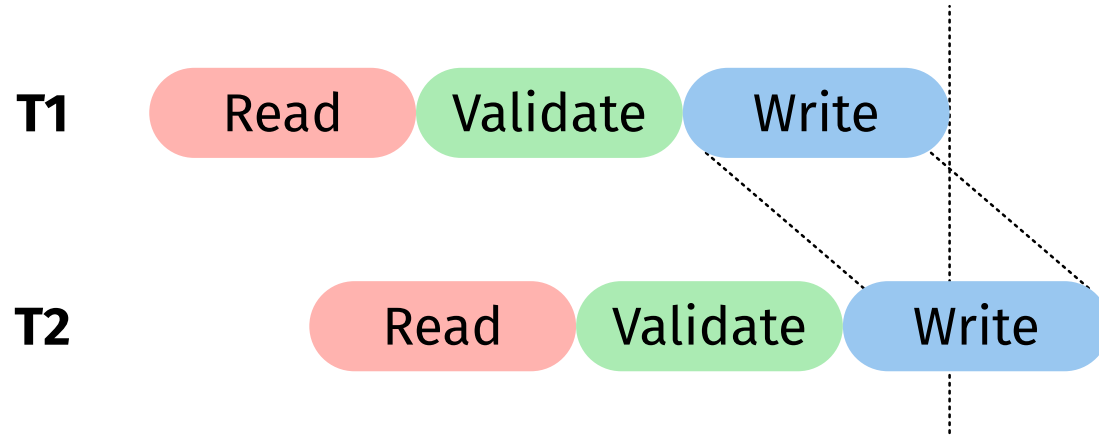




If **T1**'s write phase ends after **T2**'s read phase starts, there is a possibility of write-read conflicts. Abort (and restart) **T2** if:

$$\text{WriteSet}(\mathcal{T}_1) \cap \text{ReadSet}(\mathcal{T}_2) \neq \emptyset$$





If **T1**'s write phase overlaps with **T2**'s write phase, there is also a possibility of write-write conflicts. Abort (and restart) **T2** if either:

$$\text{WriteSet}(\mathcal{T}_1) \cap \text{ReadSet}(\mathcal{T}_2) \neq \emptyset$$

$$\text{WriteSet}(\mathcal{T}_1) \cap \text{WriteSet}(\mathcal{T}_2) \neq \emptyset$$

**Locking**

**Snapshot Isolation**

## Locking

- **Pro:** No need to restart failed transactions in general.
- **Pro:** Low overheads when there are no conflicts.

## Snapshot Isolation

## Locking

- **Pro:** No need to restart failed transactions in general.
- **Pro:** Low overheads when there are no conflicts.

## Snapshot Isolation

- **Pro:** Even lower overheads if transactions are “safe” by overhead.

## Locking

- **Pro:** No need to restart failed transactions in general.
- **Pro:** Low overheads when there are no conflicts.

## Snapshot Isolation

- **Pro:** Even lower overheads if transactions are “safe” by overhead.
- **Con:** Validation doesn't happen until the validation step.

**Idea:** Check for conflicts while the transaction is running  
(abort and restart immediately if a violation is detected)

# Timestamp Concurrency Control

## Timestamp Concurrency Control (Idealized)

Each object  $A$  gets:

- A read timestamp  $RTS(A)$
- A write timestamp  $WTS(A)$

Each transaction  $\mathcal{T}$  gets:

- A timestamp  $TS(\mathcal{T})$

(note that timestamps can be logical timestamps like sequence numbers)

We require transactions to follow the serial order defined by  $TS(\mathcal{T})$

Transactions that violate this order are restarted with a fresh timestamp

**When  $\mathcal{T}$  reads object  $A$**

**When  $\mathcal{T}$  reads object  $A$**

**if  $\text{WTS}(A) > \text{TS}(\mathcal{T})$  then**

**When  $\mathcal{T}$  reads object  $A$**

**If  $WTS(A) > TS(\mathcal{T})$  then**

- 1.** The value stored is a “later” version (write-read), so  $\mathcal{T}$  is aborted.

## When $\mathcal{T}$ reads object $A$

**If**  $WTS(A) > TS(\mathcal{T})$  **then**

1. The value stored is a “later” version (write-read), so  $\mathcal{T}$  is aborted.

**If**  $WTS(A) \leq TS(\mathcal{T})$  **then**

## When $\mathcal{T}$ reads object $A$

**If**  $WTS(A) > TS(\mathcal{T})$  **then**

1. The value stored is a “later” version (write-read), so  $\mathcal{T}$  is aborted.

**If**  $WTS(A) \leq TS(\mathcal{T})$  **then**

1. The value stored is an earlier version, so is safe to read.
2. Update  $RTS(A) \leftarrow \max(RTS(A), TS(\mathcal{T}))$

**When  $\mathcal{T}$  writes object  $A$**

**When  $\mathcal{T}$  writes object  $A$**

**If  $\text{RTS}(A) > \text{TS}(\mathcal{T})$  then**

## When $\mathcal{T}$ writes object $A$

**If**  $\text{RTS}(A) > \text{TS}(\mathcal{T})$  **then**

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

## When $\mathcal{T}$ writes object $A$

**If  $\text{RTS}(A) > \text{TS}(\mathcal{T})$  then**

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

**If  $\text{WTS}(A) > \text{TS}(\mathcal{T})$  then**

## When $\mathcal{T}$ writes object $A$

**If  $RTS(A) > TS(\mathcal{T})$  then**

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

**If  $WTS(A) > TS(\mathcal{T})$  then**

1. A later transaction has already written a new version (write-write), so  $\mathcal{T}$  is aborted.

## When $\mathcal{T}$ writes object $A$

**If  $\text{RTS}(A) > \text{TS}(\mathcal{T})$  then**

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

**If  $\text{WTS}(A) > \text{TS}(\mathcal{T})$  then**

1. A later transaction has already written a new version (write-write), so  $\mathcal{T}$  is aborted.

**Otherwise**

## When $\mathcal{T}$ writes object $A$

**If**  $\text{RTS}(A) > \text{TS}(\mathcal{T})$  **then**

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

**If**  $\text{WTS}(A) > \text{TS}(\mathcal{T})$  **then**

1. A later transaction has already written a new version (write-write), so  $\mathcal{T}$  is aborted.

**Otherwise**

1. The value is safe to write.
2. Update  $\text{WTS}(A) \leftarrow \max(\text{WTS}(A), \text{TS}(\mathcal{T}))$

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓		W(A)	
↓	W(A)		
↓			W(A)
↓	W(B)		
↓		W(B)	

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓		W(A)	
↓	W(A)		
↓			W(A)
↓	W(B)		
↓		W(B)	

This schedule is not conflict serializable

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓		W(A)	
↓	W(A)		
↓			W(A)
↓	W(B)		
↓		W(B)	

This schedule is not conflict serializable ... but it could be

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓		IGNORE	
↓	W(A)		
↓			W(A)
↓	W(B)		
↓		W(B)	

This schedule is not conflict serializable ... but it could be

## When $\mathcal{T}$ reads object $A$

### If $\text{RTS}(A) > \text{TS}(\mathcal{T})$ then

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

### If $\text{WTS}(A) > \text{TS}(\mathcal{T})$ then

1. A later transaction has already written a new version (write-write), so  $\mathcal{T}$  is aborted.

### Otherwise

1. The value is safe to write.
2. Update  $\text{WTS}(A) \leftarrow \max(\text{WTS}(A), \text{TS}(\mathcal{T}))$

## When $\mathcal{T}$ reads object $A$

### If $\text{RTS}(A) > \text{TS}(\mathcal{T})$ then

1. A later transaction has already read the stored version (read-write), so  $\mathcal{T}$  is aborted.

### If $\text{WTS}(A) > \text{TS}(\mathcal{T})$ then

1. A later transaction has already written a new version (write-write), so **the write is ignored**

### Otherwise

1. The value is safe to write.
2. Update  $\text{WTS}(A) \leftarrow \max(\text{WTS}(A), \text{TS}(\mathcal{T}))$

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓			W(A)
↓		W(A)	
↓	W(A)		
↓	W(B)		
↓		W(B)	

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓			W(A)
↓		IGNORE	
↓	W(A)		
↓	W(B)		
↓		W(B)	

<b>Time</b>	<b>T1</b>	<b>T2</b>	<b>T3</b>
↓			W(A)
↓		IGNORE	
↓	IGNORE		
↓	W(B)		
↓		W(B)	

Time	T1	T2	T3
↓			W(A)
↓		IGNORE	
↓	IGNORE		
↓	W(B)		
↓		W(B)	

... but this schedule isn't conflict serializable

# View Serializability

## **Conflict Equivalence**

The order of two operations has no effect on the transaction's visible if the operations:

1. ... act on different objects
2. ... or are both reads

## View Equivalence

The order of two operations has no effect on the transaction's visible if the operations:

1. ... act on different objects
2. ... or are both reads
3. ... are both writes and their effects are hidden by a third write

## View Equivalence

The order of two operations has no effect on the transaction's visible if the operations:

1. ... act on different objects
2. ... or are both reads
3. ... are both writes and their effects are hidden by a third write

A schedule is **view serializable** if it is view equivalent to a serial schedule

Timestamp Concurrency Control guarantees view serializable schedules.

- Serializability is guaranteed by...
  - View serializability, which is guaranteed by...
    - Timestamp concurrency control
    - Conflict serializability, which is guaranteed by...
      - 2-Phase Locking
      - Optimistic Concurrency Control (aka Snapshot Isolation<sup>1</sup>)

---

<sup>1</sup>Database vendors, in their markety wisdom (incorrectly) refer to a weaker form of TSCC that we'll talk about next as "Snapshot Isolation". This term is, accordingly, ambiguous.

**Problem:** Adding timestamp checks on reads is **expensive**.

**Solution?:** Ignore read timestamps.

(Major database vendors (incorrectly) call this “snapshot isolation”)

**Write-Write conflicts**

**Write-Read conflicts**

**Read-Write conflicts**

## **Write-Write conflicts**

Detected!

## **Write-Read conflicts**

## **Read-Write conflicts**

## **Write-Write conflicts**

Detected!

## **Write-Read conflicts**

Detected!

## **Read-Write conflicts**

## **Write-Write conflicts**

Detected!

## **Write-Read conflicts**

Detected!

## **Read-Write conflicts**

Ignored!

```
A = 100      def T1():      def T2():  
B = 200      A = B+1      B = A+1
```

<b>Time</b>	<b>T1</b>	<b>T2</b>
↓		R(A)
↓	R(B)	
↓	W(A)	
↓		W(B)

A = 100  
B = 200

def T1():  
A = B+1

def T2():  
B = A+1

Time	T1	T2
↓		R(A) A → 100
↓	R(B)	
↓	W(A)	
↓		W(B)

A = 100  
B = 200

def T1():  
A = B+1

def T2():  
B = A+1

Time	T1	T2
↓		R(A) A → 100
↓	R(B) B → 200	
↓	W(A)	
↓		W(B)

```
A = 100      def T1():      def T2():  
B = 200      A = B+1      B = A+1
```

Time	T1	T2
↓		R(A) A → 100
↓	R(B) B → 200	
↓	W(A) A ← 201	
↓		W(B)

A = 100  
B = 200

def T1():  
A = B+1

def T2():  
B = A+1

Time	T1	T2
↓		R(A) A → 100
↓	R(B) B → 200	
↓	W(A) A ← 201	
↓		W(B) B ← 101

```

A = 100      def T1():      def T2():
B = 200      A = B+1        B = A+1

```

Time	T1	T2
↓		R(A) A → 100
↓	R(B) B → 200	
↓	W(A) A ← 201	
↓		W(B) B ← 101

This schedule would be permitted!

# Recoverability

<b>Time</b>	<b>T1</b>	<b>T2</b>
↓	W(A)	
↓		R(A)
↓		W(B)
↓		COMMIT

<b>Time</b>	<b>T1</b>	<b>T2</b>
↓	W(A)	
↓		R(A)
↓		W(B)
↓		COMMIT
↓	ABORT	

Oops...

**Until a transaction  $\mathcal{T}$  commits successfully...**

**Until a transaction  $\mathcal{T}$  commits successfully...**

1. Buffer all writes (but update WTS immediately)

## **Until a transaction $\mathcal{T}$ commits successfully...**

1. Buffer all writes (but update WTS immediately)
2. Block transactions reading from  $\mathcal{T}$ 's writes from committing

## **Until a transaction $\mathcal{T}$ commits successfully...**

1. Buffer all writes (but update WTS immediately)
2. Block transactions reading from  $\mathcal{T}$ 's writes from committing
3. Abort transactions reading from  $\mathcal{T}$ 's writes if  $\mathcal{T}$  aborts

# Multiversion Timestamp Concurrency Control

Write-Read conflicts are avoidable  
(the necessary data was available at one point)

**Idea:** Keep old versions of objects around

## **Each object version tracks...**

- The writing transaction's WTS
- The greatest TS of any transaction to read it as an RTS

## When $\mathcal{T}$ reads an object $A$

1. Find the newest version  $A_i$  such that  $WTS(A_i) < TS(\mathcal{T})$   
This is always safe!
2. Update  $RTS(A_i) \leftarrow \max(RTS(A_i), TS(\mathcal{T}))$

## When $\mathcal{T}$ writes an object $A$

1. Find the newest version  $A_i$  such that  $WTS(A_i) < TS(\mathcal{T})$

**If**  $RTS(A_i) < TS(\mathcal{T})$

2. The read is unrecoverable, as a later transaction should have read this version.  
Abort  $\mathcal{T}$

## Otherwise

2. Write the new version in between  $A_i$  and  $A_{i+1}$

Next time...

AC ⚡ ID