Transactions & Concurrency Control

CSC443H1 Database System Technology

Please do the course evaluation!



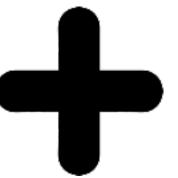
Sum of money in a system should stay constant



Sum of money in a system should stay constant

An account balance cannot drop below zero



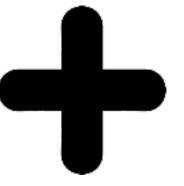


Sum of money in a system should stay constant

An account balance cannot drop below zero

A user must have a valid SIN







Sum of money in a system should stay constant

An account balance cannot drop below zero

A user must have a valid SIN

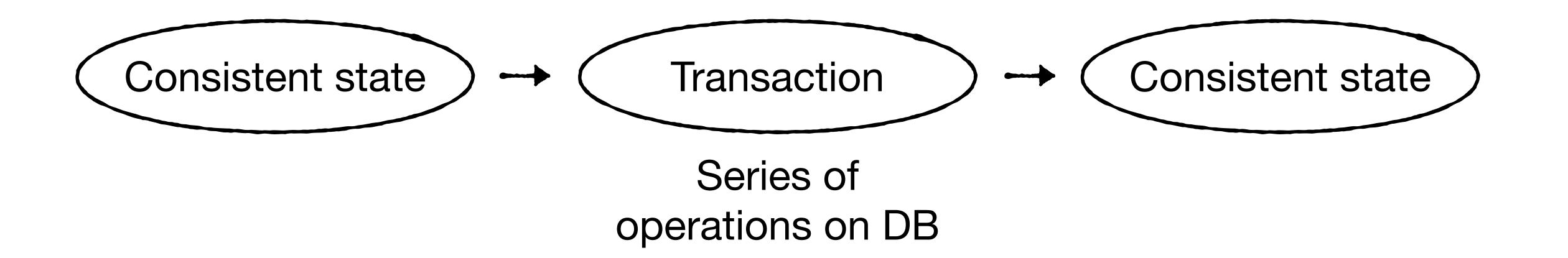
sum(balance) = X

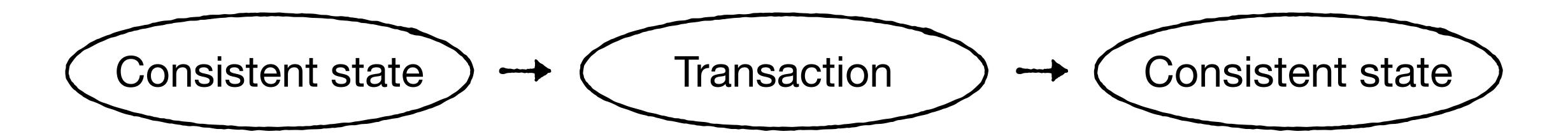
balance ≥ 0

SIN ≠ null

These can be set as integrity constraints

To maintain consistency, database APIs expose the concept of a transaction

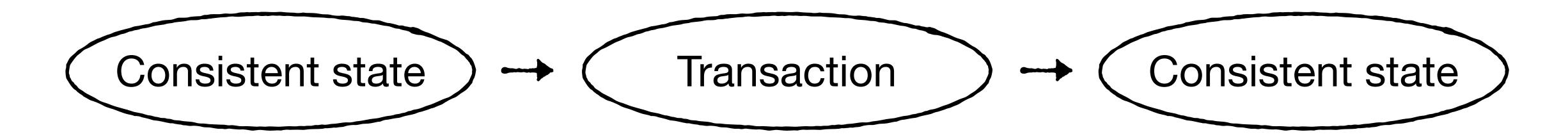




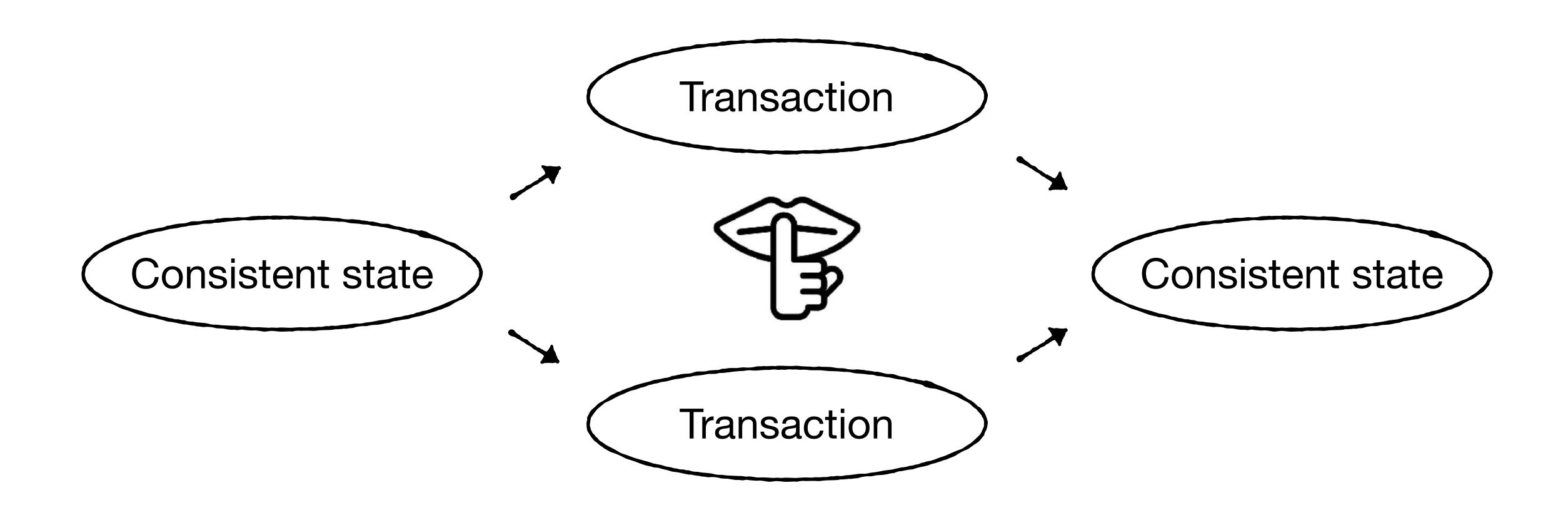
Example: Bank transfer

$$A = A - 100$$

$$B = B + 100$$



Example: Bank transfer



Transactions do not interact with each other (i.e., by exchanging messages). A transaction is oblivious to all other ongoing transactions.

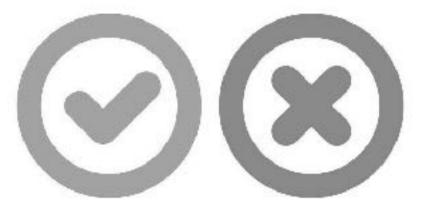
Transaction API

Begin Transaction

Some commands (e.g., in SQL)



Commit or Abort



Example: Money withdrawal

Begin Transaction

b = Select balance from accounts where id = x

If $b \leq 100$

Abort

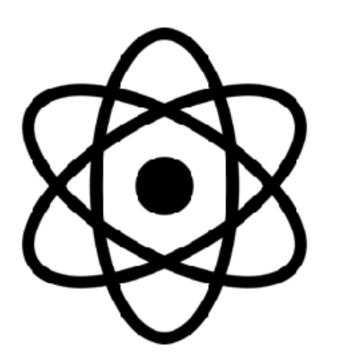
Else

Update accounts set balance = b - 100 where id = x

Commit

ACID semantics for Transactions

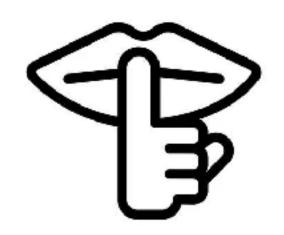
Atomicity



Consistency



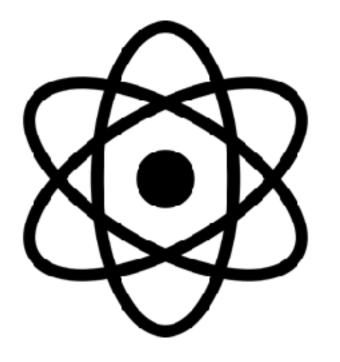
Isolation



Durability



Atomicity



All or nothing

Consistency



Isolation



Durability

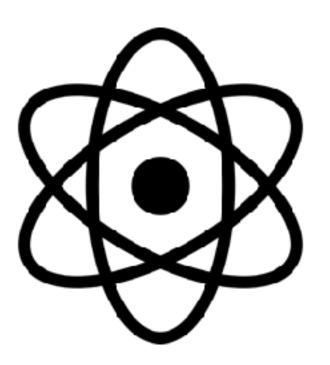


Atomicity

Consistency

Isolation

Durability









All or nothing

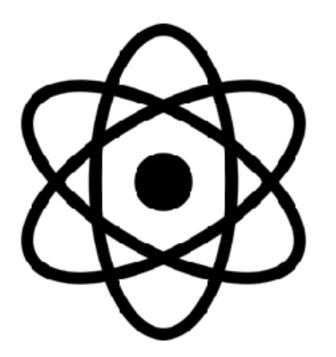
Transition across consistent states

Atomicity

Consistency

Isolation

Durability



All or nothing



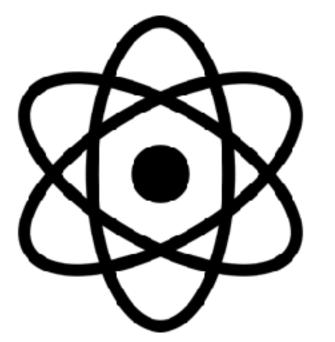
Transition across consistent states



Not corrupted by concurrency



Atomicity



All or nothing

Consistency



Transition across consistent states

Isolation



Not corrupted by concurrency

Durability



Recover from failure

Who is responsible for what?

Atomicity Consistency Isolation Durability

DB User DB DB

Durability

DB DB

What Endangers Consistency?

System Failure

(4)

Concurrency



Power failure

(4)

Hardware failure



Data center failure



Power failure

Example: Bank transfer



A = A - 100Power fails B = B + 100

Power failure

Example: Bank transfer



A = A - 100Power fails B = B + 100

Power fails before transaction finishes, leaving the system in an inconsistent state

(4)

Concurrency



Concurrent transactions can result in an inconsistent state

Concurrent transactions can result in an inconsistent state

Can you think of examples?

Concurrent transactions can result in an inconsistent state

Example 1 Interest & Payments



Example 2
Updates & statistics



T1: Interest payment

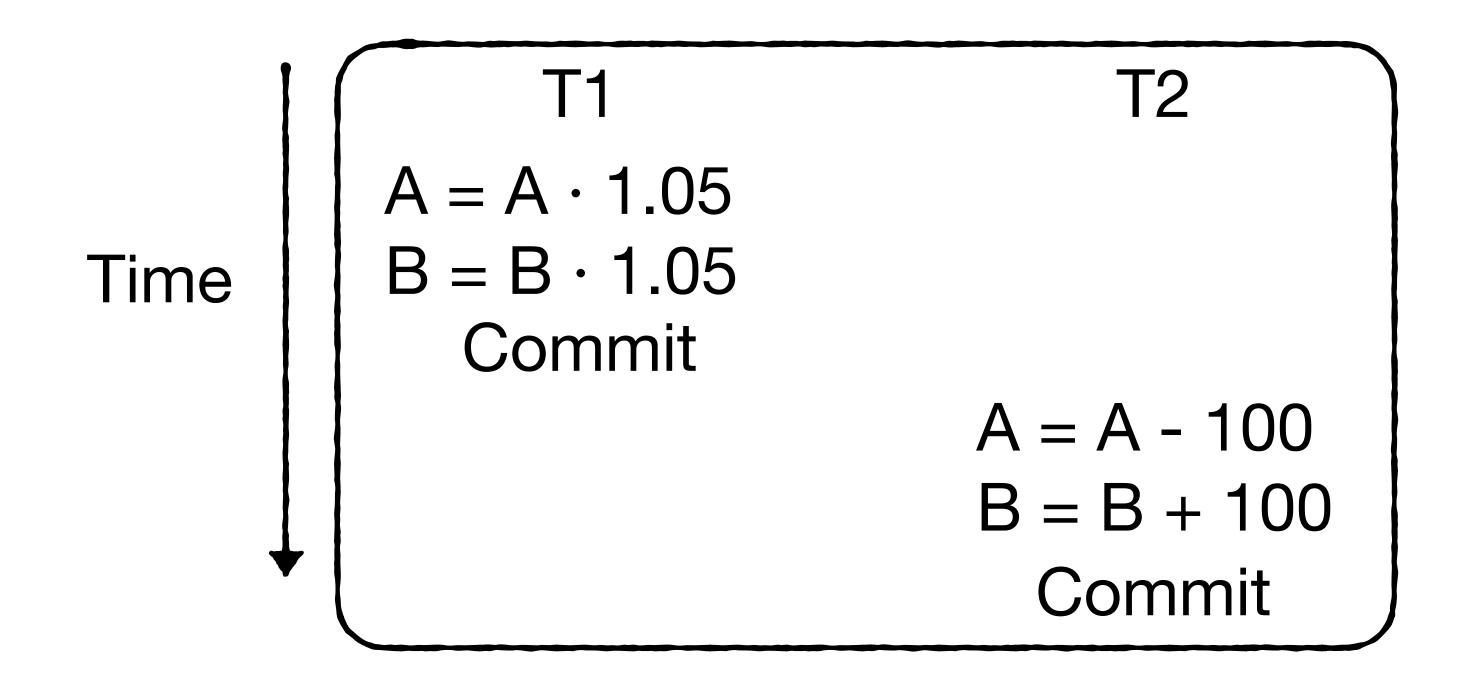
$$A = A \cdot 1.05$$

$$B = B \cdot 1.05$$

T2: Payments

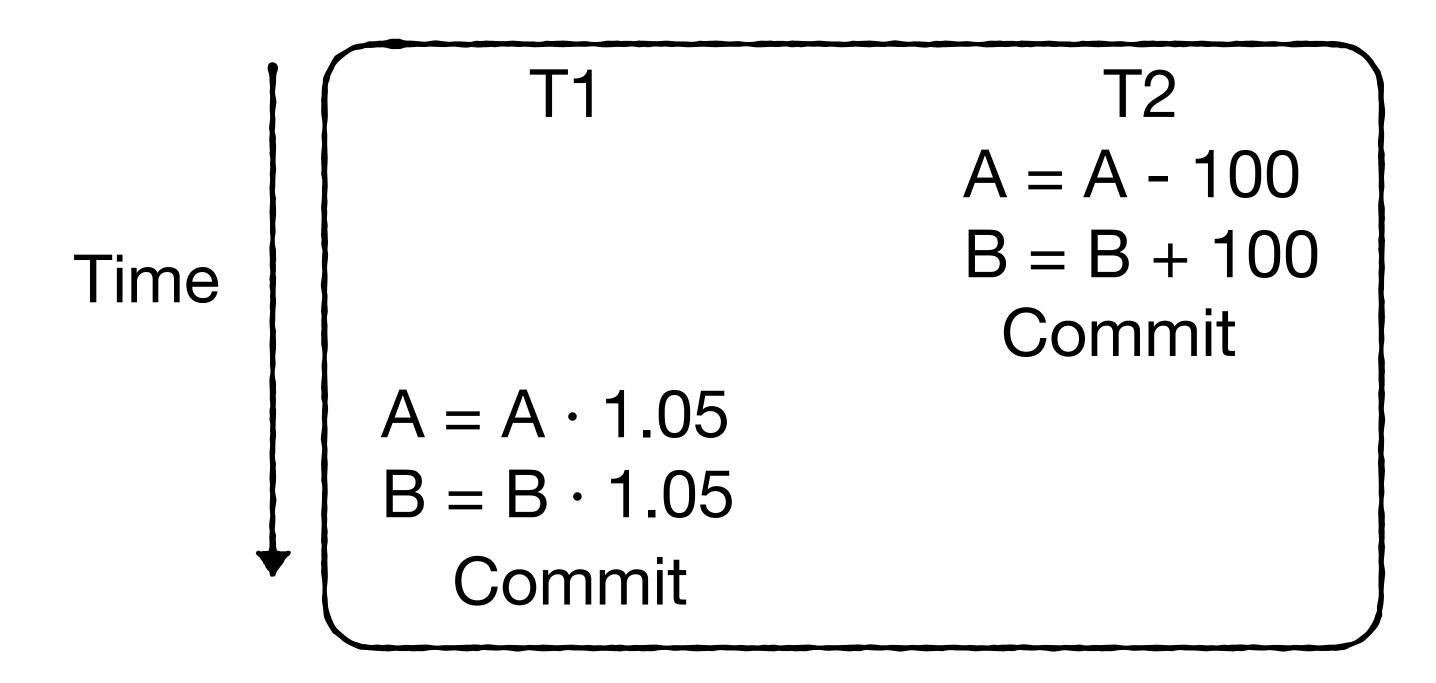
$$A = A - 100$$

$$B = B + 100$$



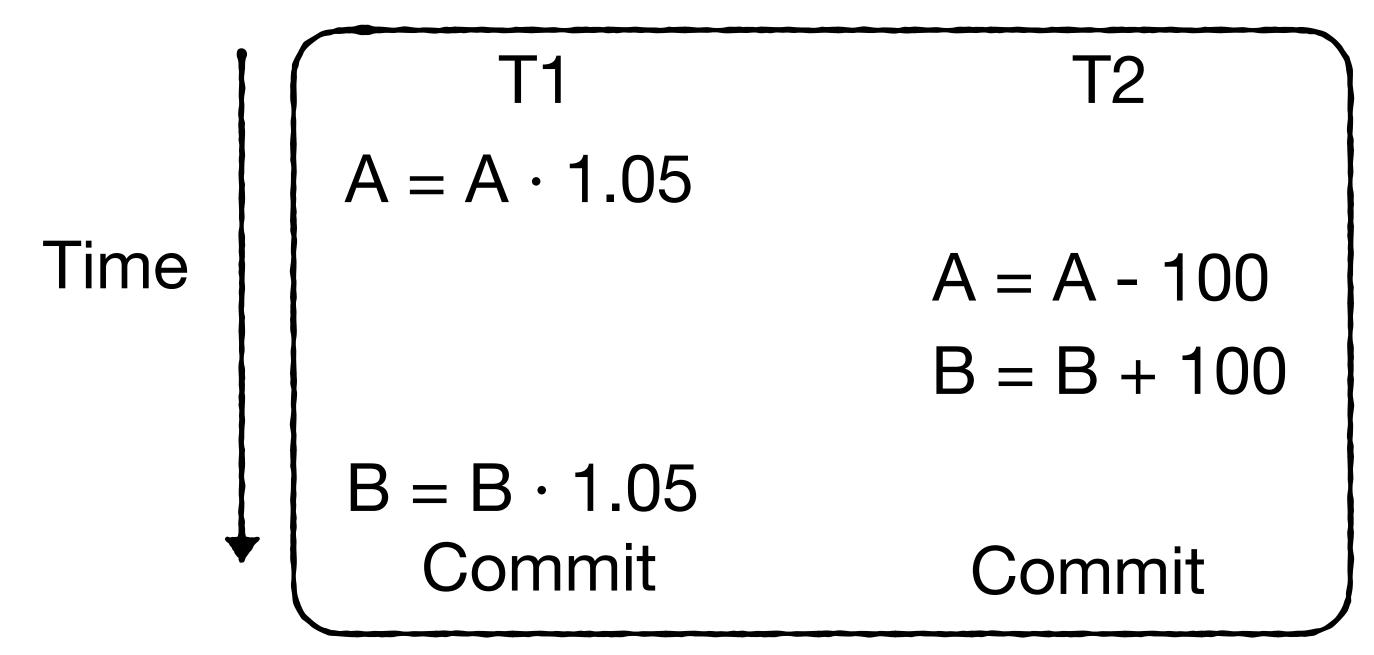
Initialization: A=1000, B=1000

Initialization: A=1000, B=1000



Valid Outcome 2: A = 945, B = 1155

Initialization: A=1000, B=1000



Invalid Outcome: A = 950, B = 1155

Valid Outcome 2: A = 945, B = 1155

Initialization: A=1000, B=1000

Dirty read: T2 reads a modified but uncommitted data item from T1

Invalid Outcome: A = 950, B = 1155

Valid Outcome 2: A = 945, B = 1155

Concurrent transactions can result in an inconsistent state

Example 1
Payments & Interest



Example 2
Updates & Statistics



T1: Account Updates

Balance += X

T2: Statistics Reporting

count(# accounts)
sum(account balances)
avg(account balances)

T1: Account Updates

bob_checking += X

T2: Statistics Reporting

count(# accounts) where user = "bob" sum(account balances) where user = "bob" avg(account balances) where user = "bob"

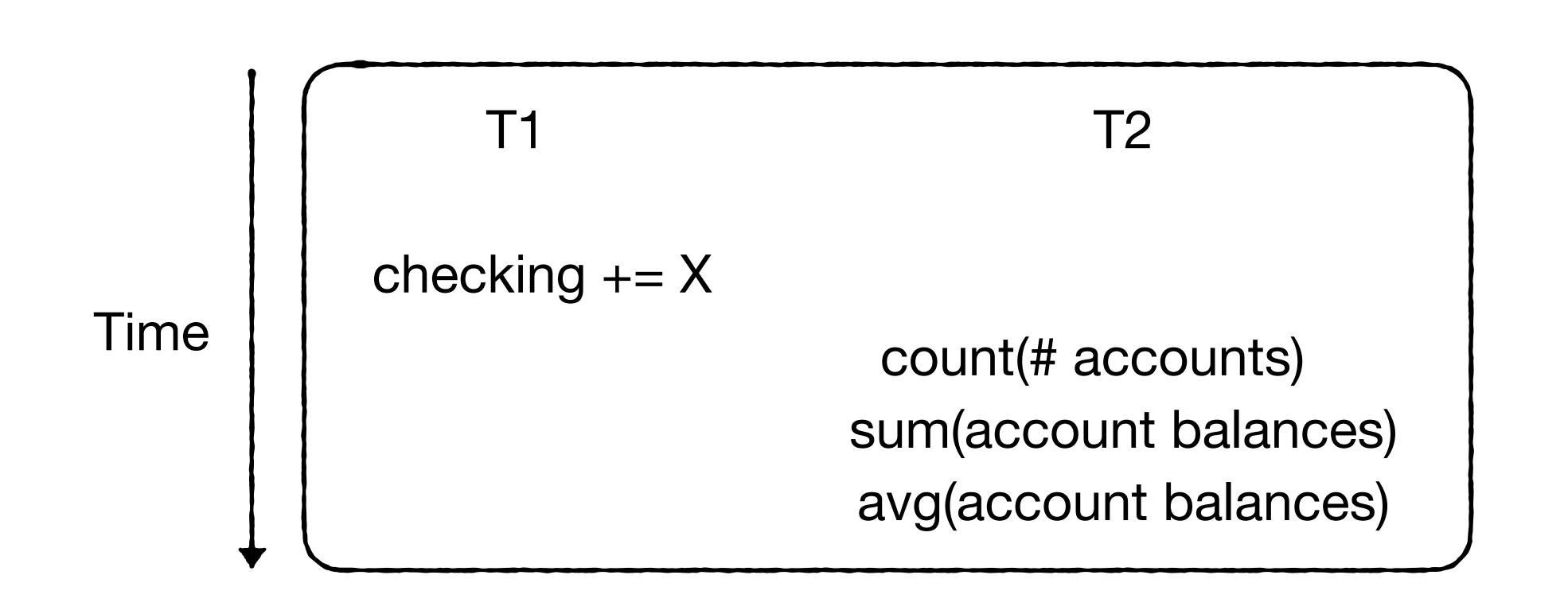
T1: Account Updates

bob_checking += X

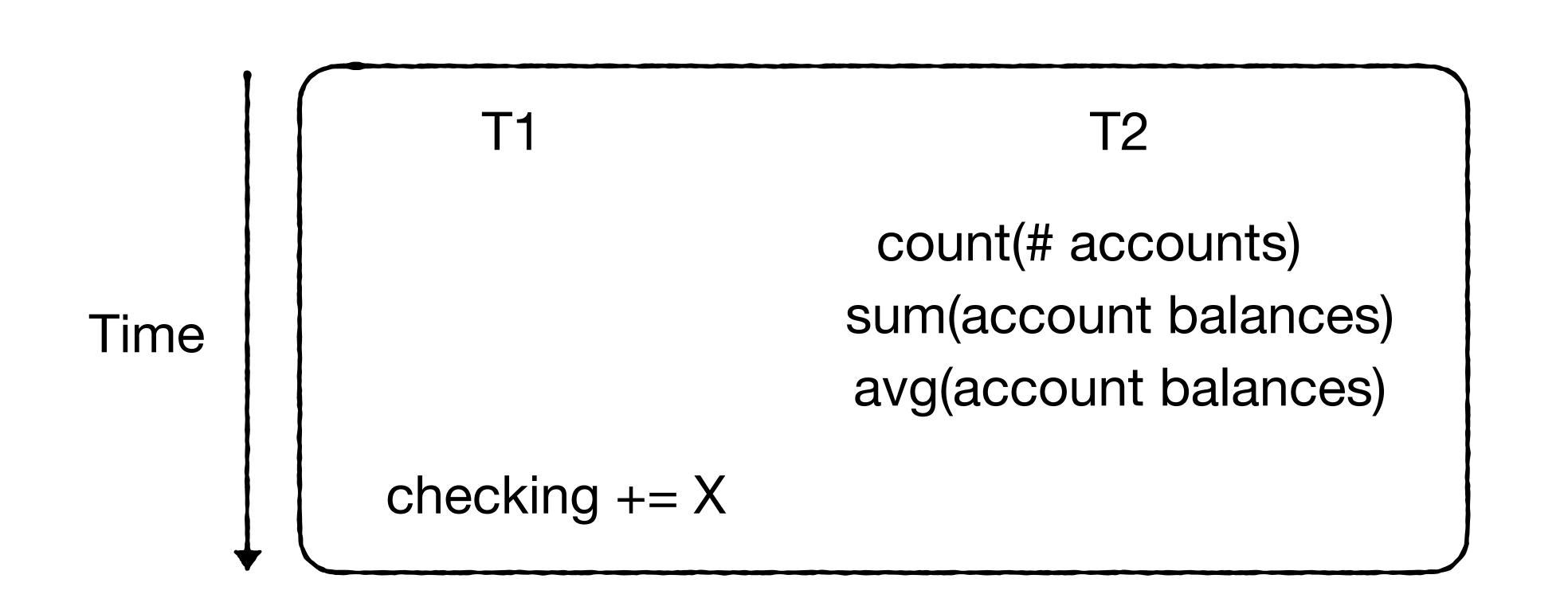
T2: Statistics Reporting

count(# accounts) where user = "bob" sum(account balances) where user = "bob" avg(account balances) where user = "bob"

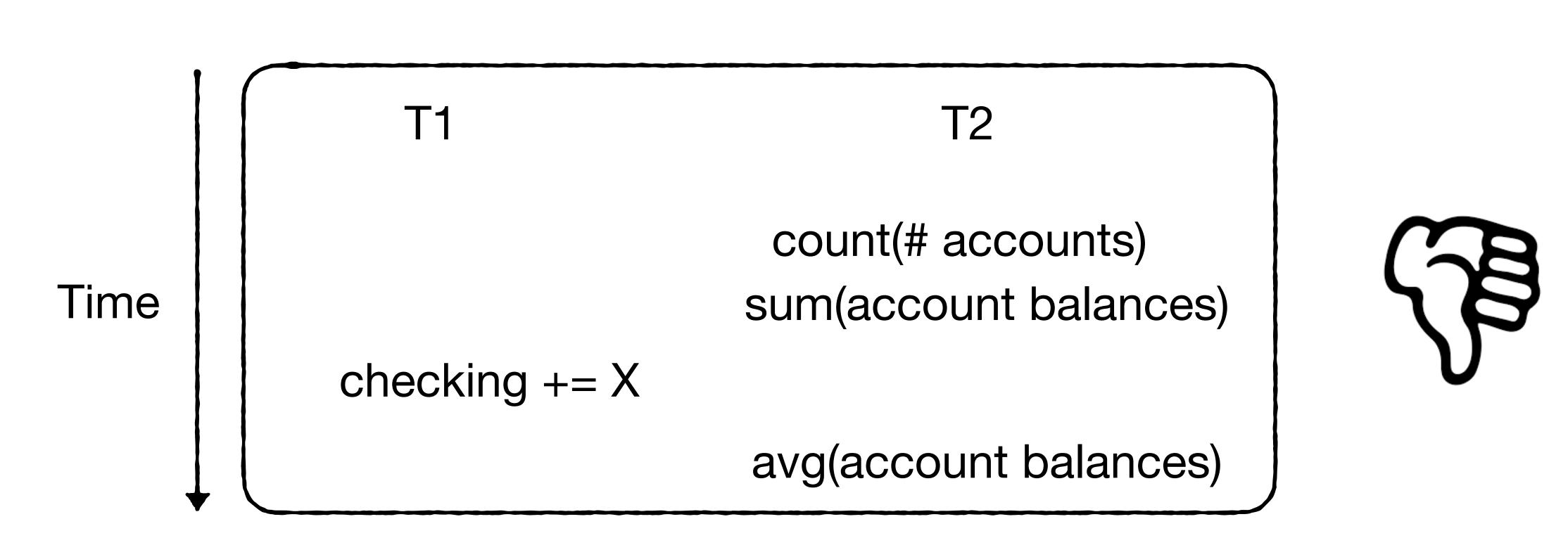
What can go wrong?



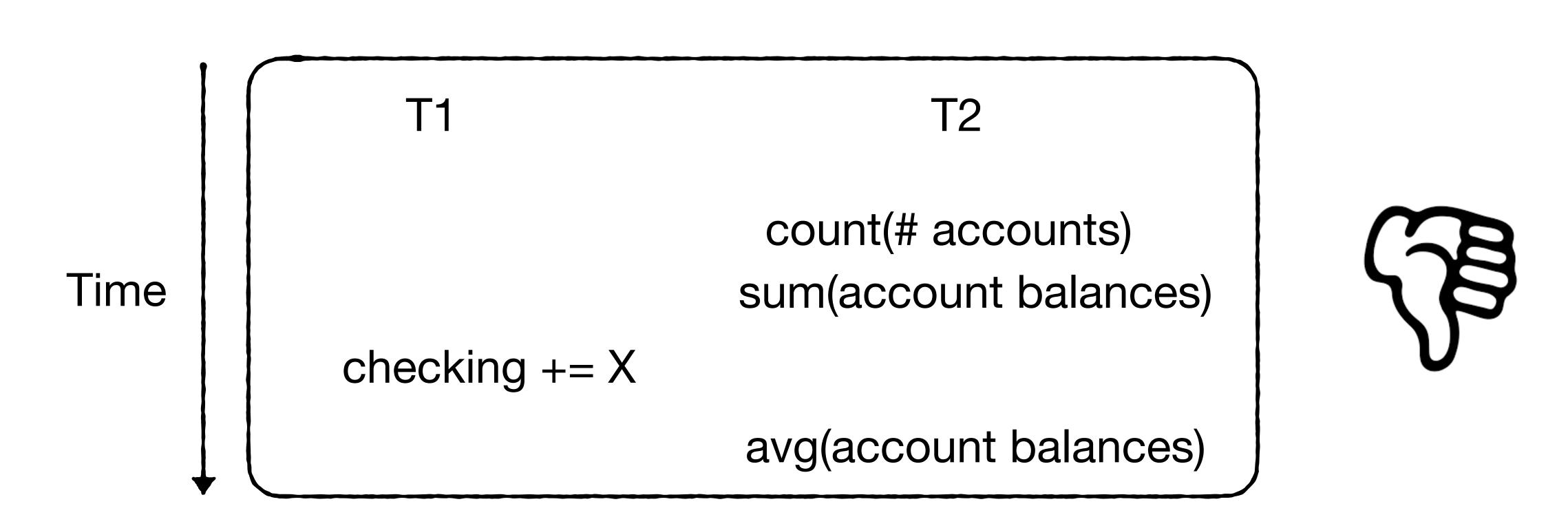




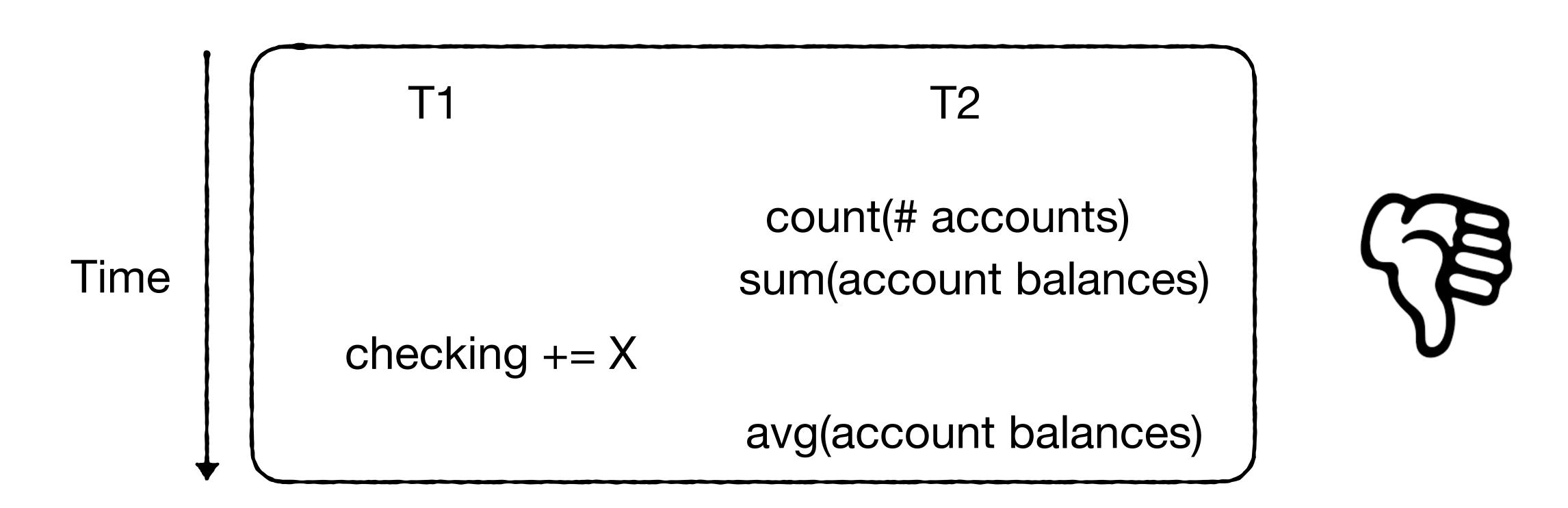








Problem? Inconsistent reporting: avg ≠ sum / count



More broadly, this is an unrepeatable read anomaly: subsequent reads of the same data are inconsistent as the data was changed in-between

Concurrency

Concurrent transactions can result in an inconsistent state

Example 1

Payments & Interest

Example 2

Updates & Statistics



Dirty reads

Problem:



Unrepeatable reads



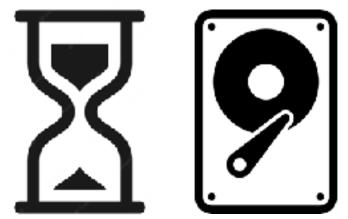
Problems?

Problems: terrible for performance



Problems: terrible for performance

While one transaction waits for a storage I/O, another should be able to use the CPU

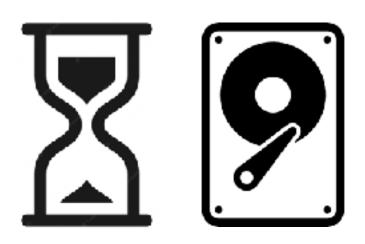




Problems: terrible for performance

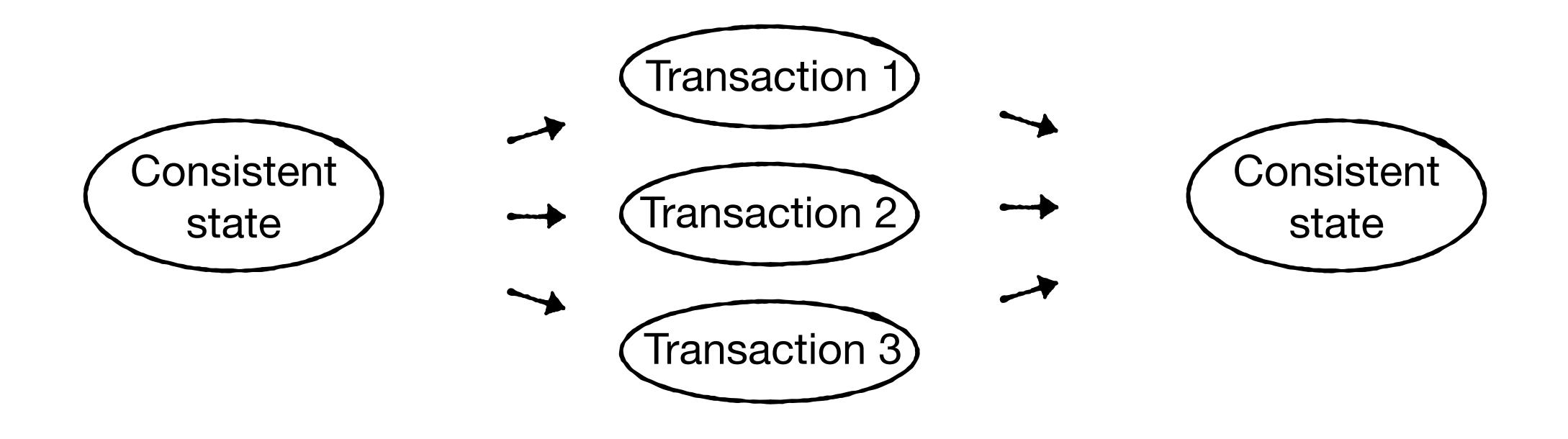
While one transaction waits for a storage I/O, another should be able to use the CPU

A short transaction should not have to wait until a long transaction completes

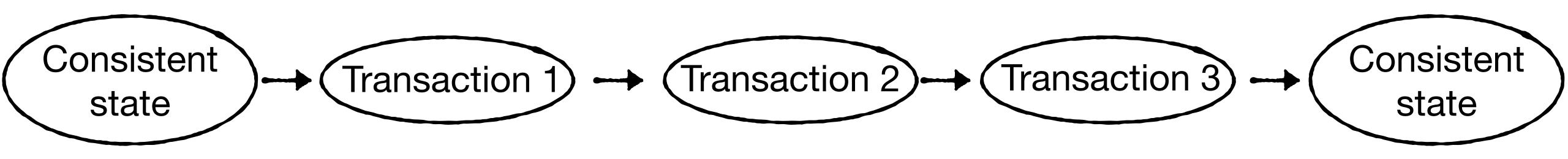




Transactions can be concurrent but must result in a consistent state

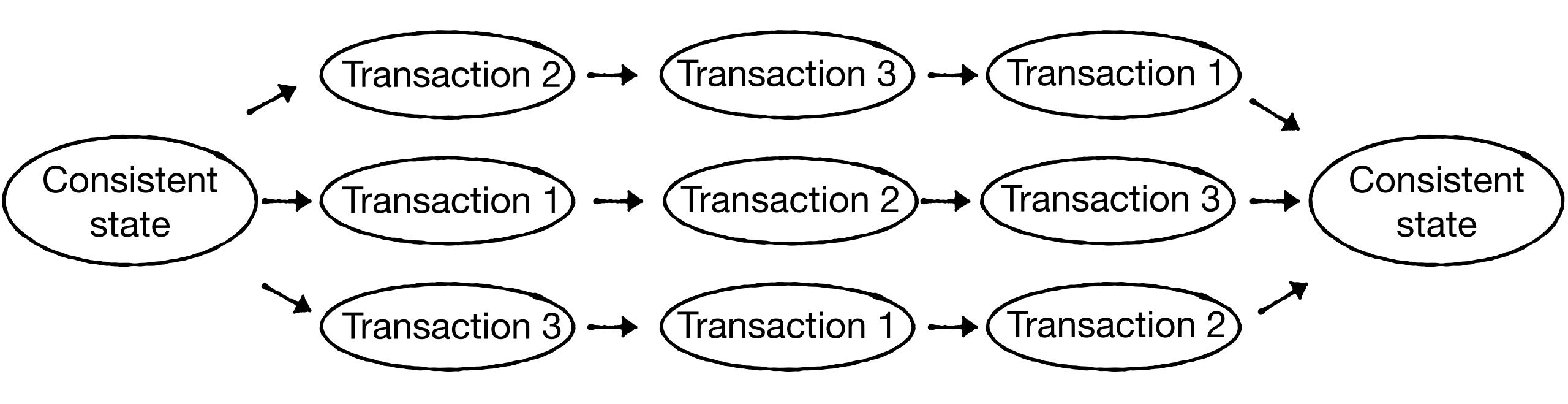


Transactions can be concurrent but must result in a consistent state



Any given state once a transaction commits should have been achievable through a serial execution of all committed transactions thus far

Transactions can be concurrent but must result in a consistent state



Equivalence to any serial execution is considered correct

How to achieve concurrency & serializability?

A transaction is divided into two phases

Phase 1 Phase 2

A transaction is divided into two phases

Phase 1 Phase 2

Lock a data item (e.g., row) when it is accessed for the first time

A transaction is divided into two phases

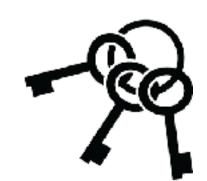
Phase 1

Phase 2

Lock a data item (e.g., row) when it is accessed for the first time

Reads take shared locks

Writes take exclusive locks





A transaction is divided into two phases

Phase 1

Lock a data item (e.g., row) when it is accessed for the first time

Reads take shared locks

Writes take exclusive locks





Phase 2

Release all locks during commit



A transaction is divided into two phases

Phase 1

Lock a data item (e.g., row) when it is accessed for the first time

Phase 2

Release all locks during commit

Invariant: a data item that has been accessed by this transaction cannot be modified by another transaction until this transaction commits

Let's fix our running examples

Example 1

Interest & Payments

Example 2

Updates & Statistics



Problems:

Dirty reads

Unrepeatable reads

T1: Interest payment

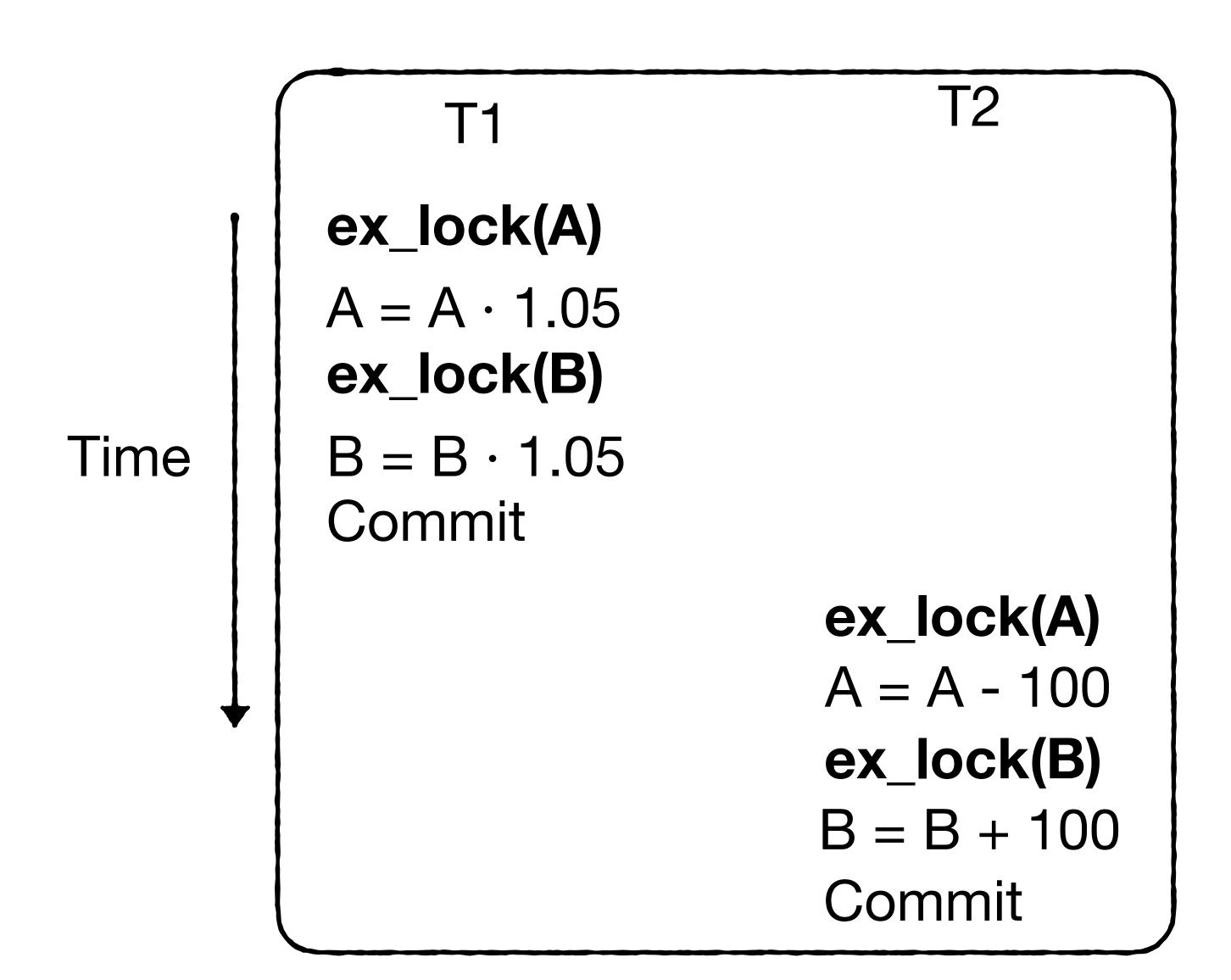
$$A = A \cdot 1.05$$

$$B = B \cdot 1.05$$

T2: Payments

$$A = A - 100$$

$$B = B + 100$$



Exclusive locks ensure T2 cannot modify A until T1 commits

Let's fix our running examples

Example 1

Interest & Payments

Example 2

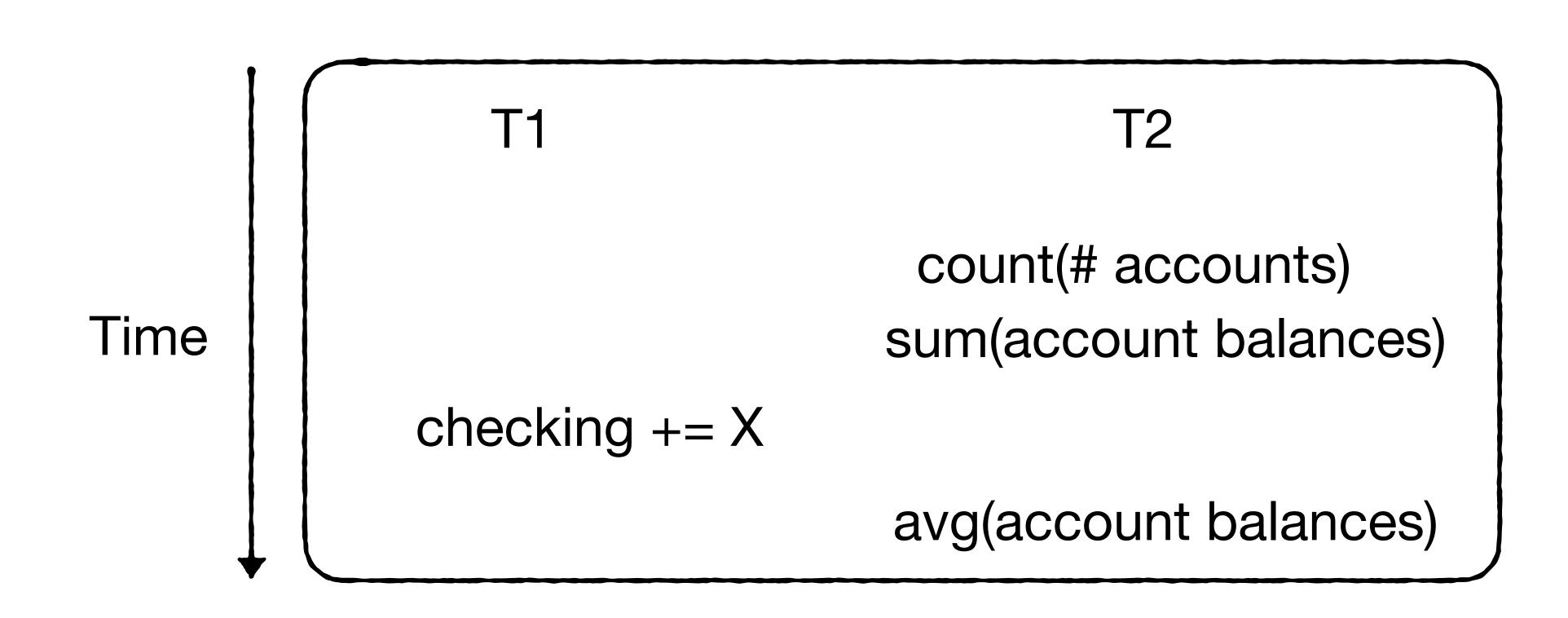
Updates & Statistics

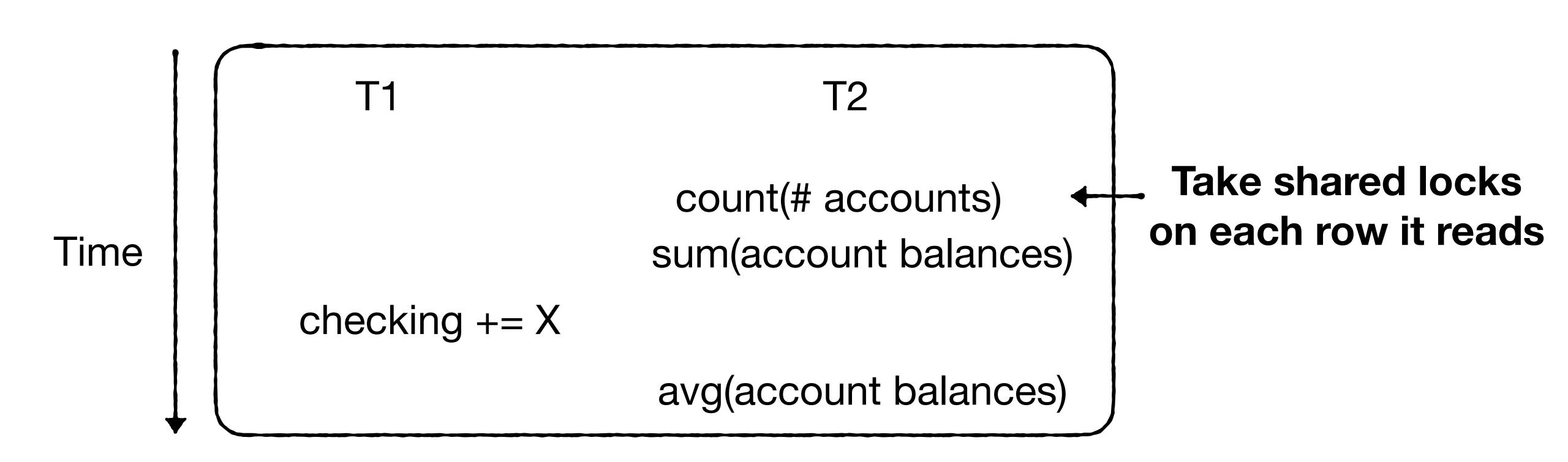


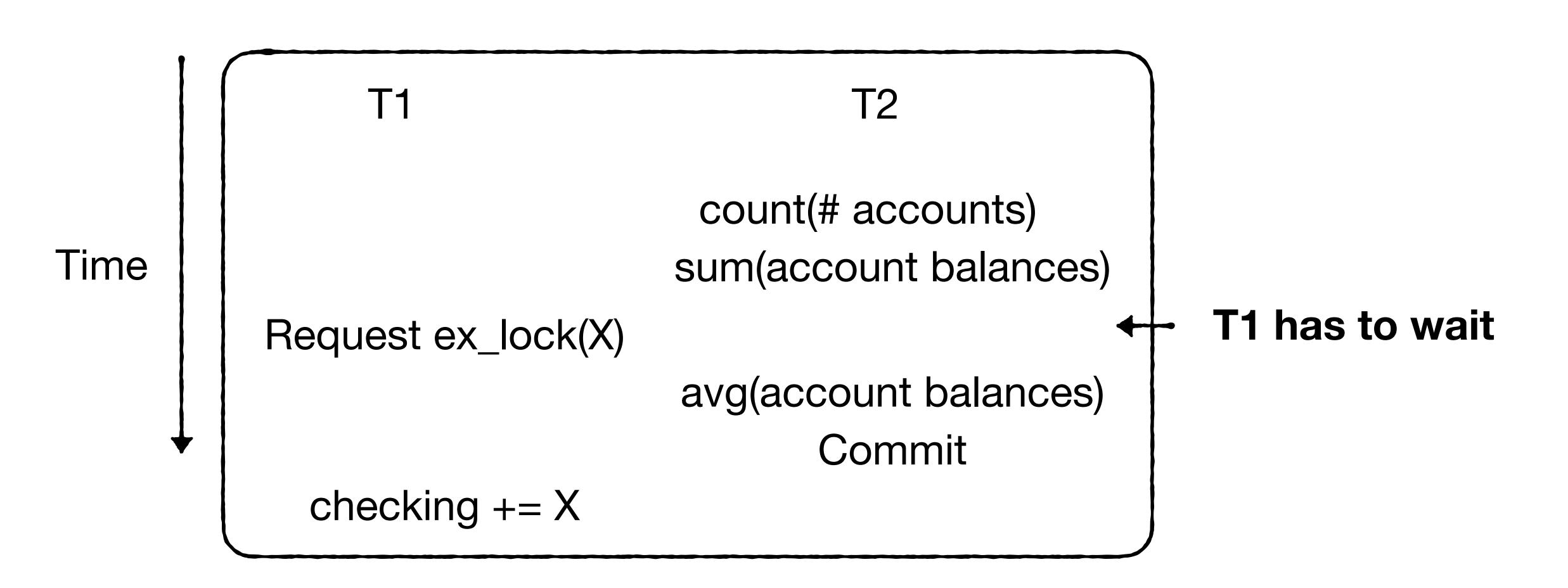
Problems:

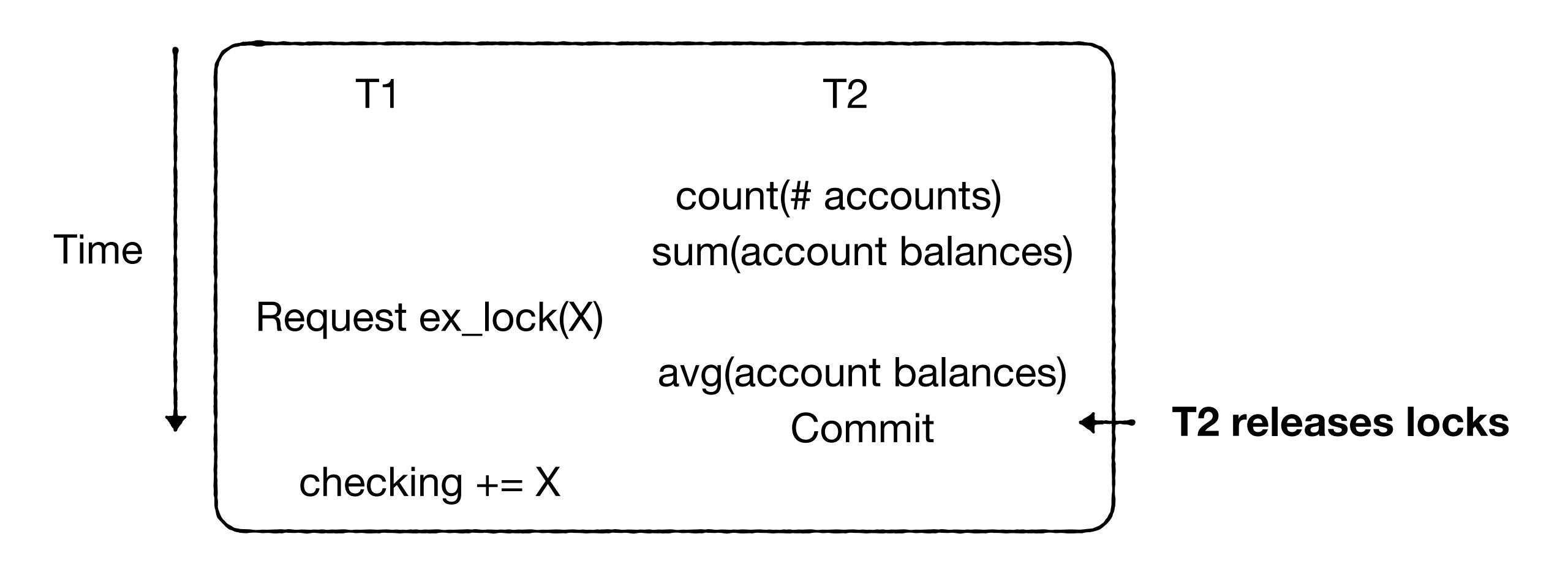
Dirty reads

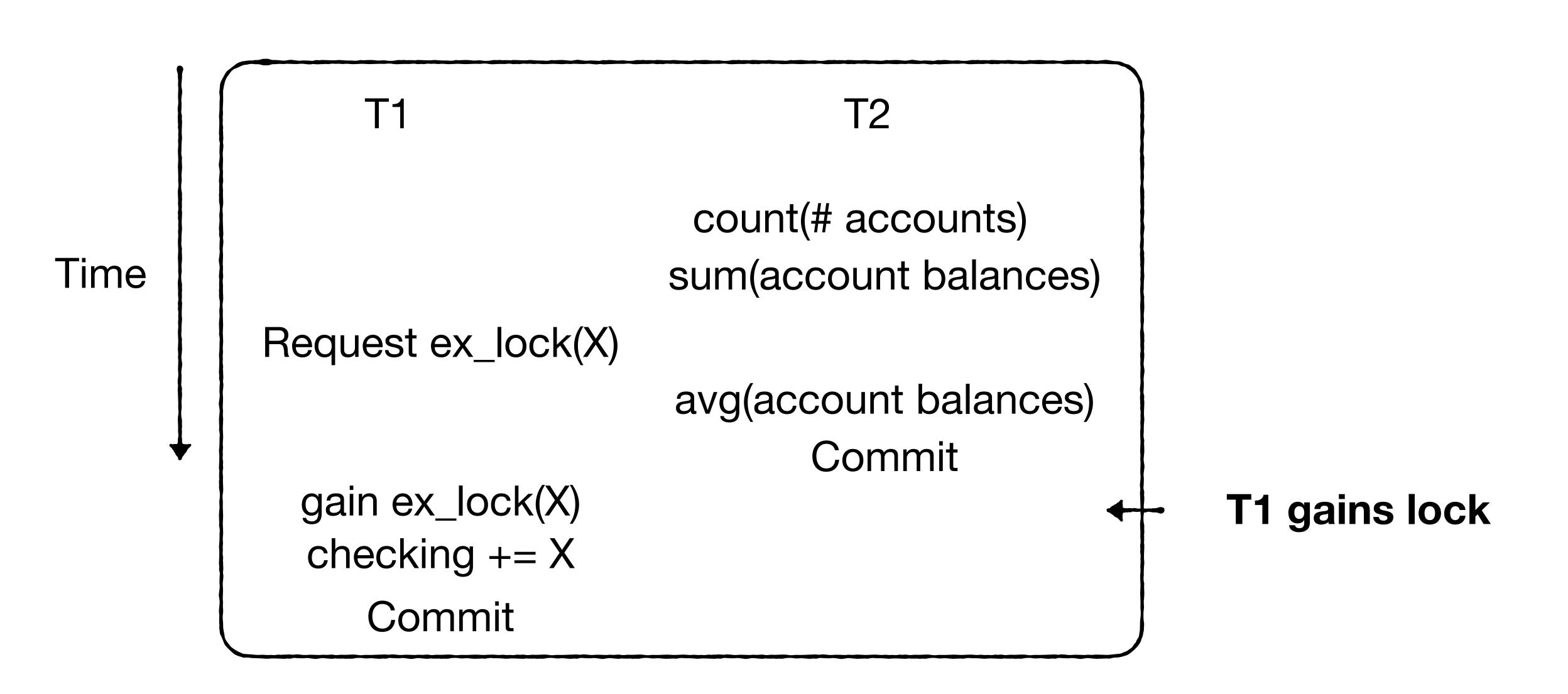
Unrepeatable reads

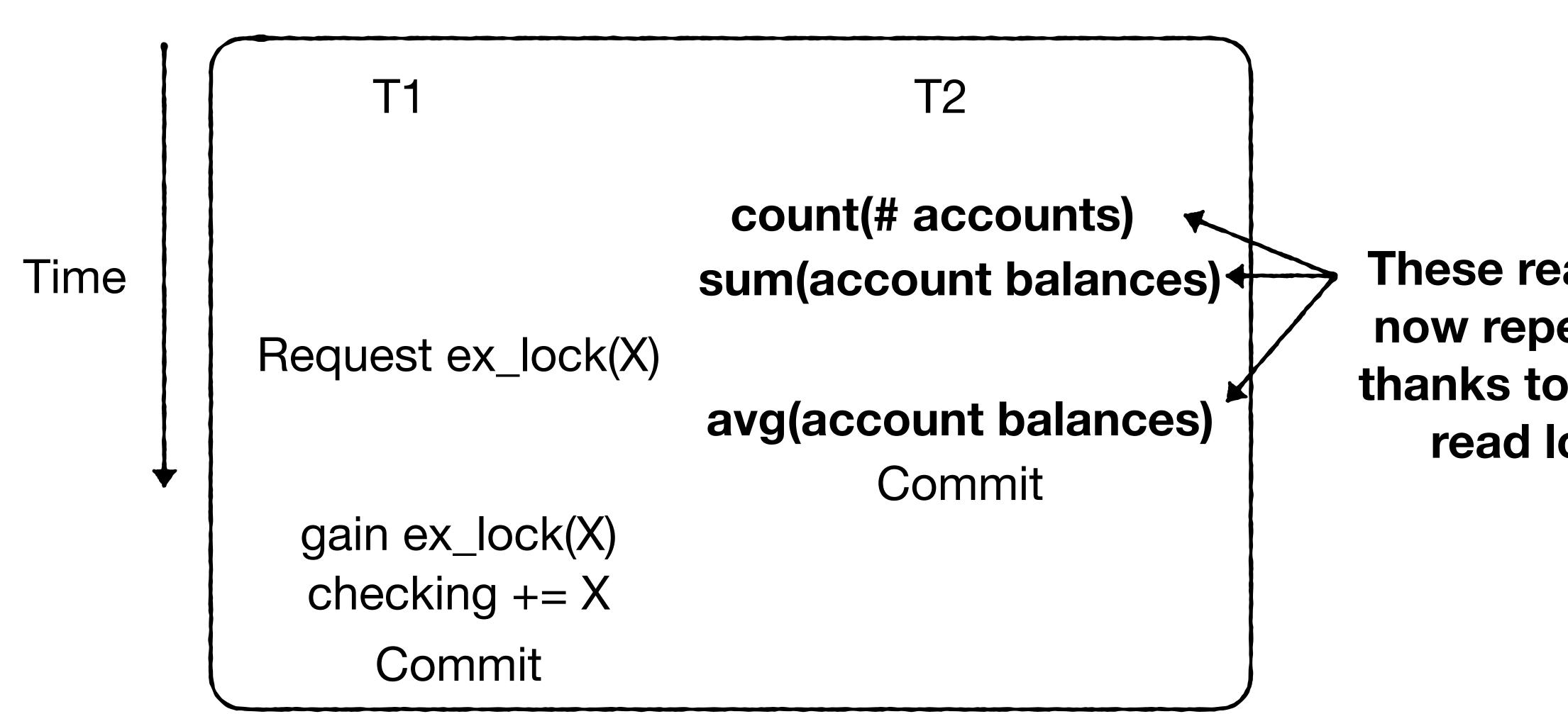












These reads are now repeatable thanks to shared read locks

Both examples now work correctly

Example 1

Interest & Payments

Example 2

Updates & Statistics





Problems:

Dirty reads

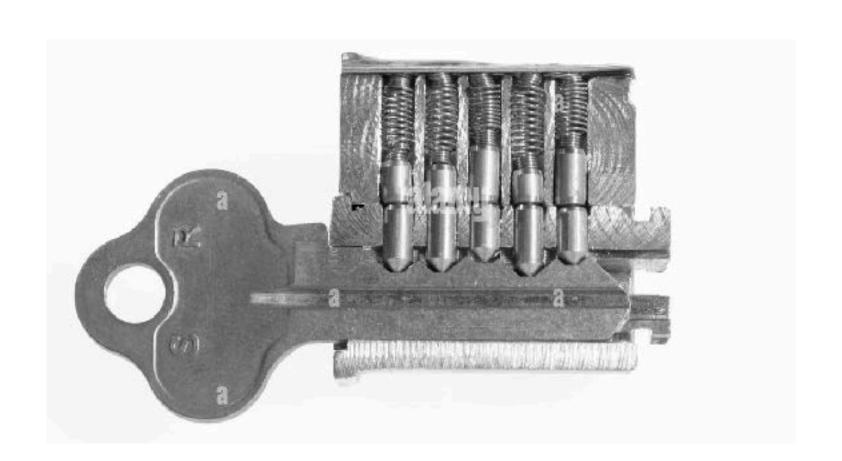
Unrepeatable reads

Solution:

Exclusive write locks

Shared read locks

How are locks implemented?



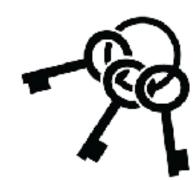
Lock Manager - a hash table

Object ID

Object ID → (type, lock count, queue of waiting requests)

Object ID --- (type, lock count, queue of waiting requests)

Shared/exclusive





Object ID → (type, lock count, queue of waiting requests)

In case of shared lock



Object ID → (type, lock count, queue of waiting requests)

What to invoke next when this lock is released



Object ID → (type, lock count, queue of waiting requests)





Object ID --> (type, lock count, queue of waiting requests)

We lock an object the first time it is accessed by a transaction



Object ID - (type, lock count, queue of waiting requests)

We lock an object the first time it is accessed by a transaction



Locking table implemented via OS locking primitives (e.g., mutexes)



Locks Latches



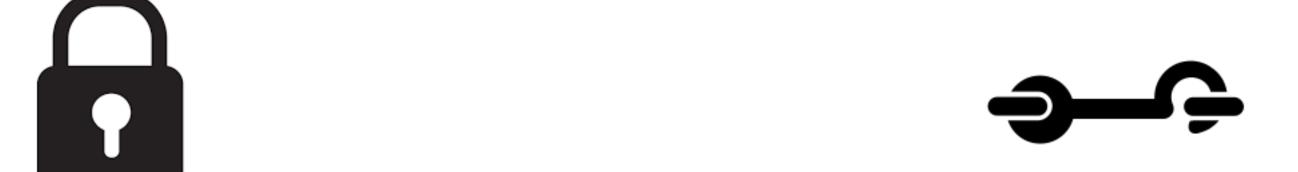


Locks Latches



Separate: Transactions Threads

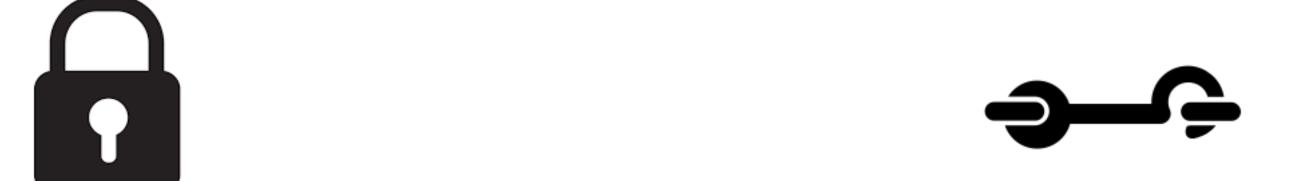
Locks Latches



Separate: Transactions Threads

Protect: DB content In-memory data structures (LRU queue)

Locks Latches



Separate: Transactions Threads

Protect: DB content In-memory data structures (LRU queue)

Duration: Transaction Critical section

Locks Latches



Separate: Transactions Threads

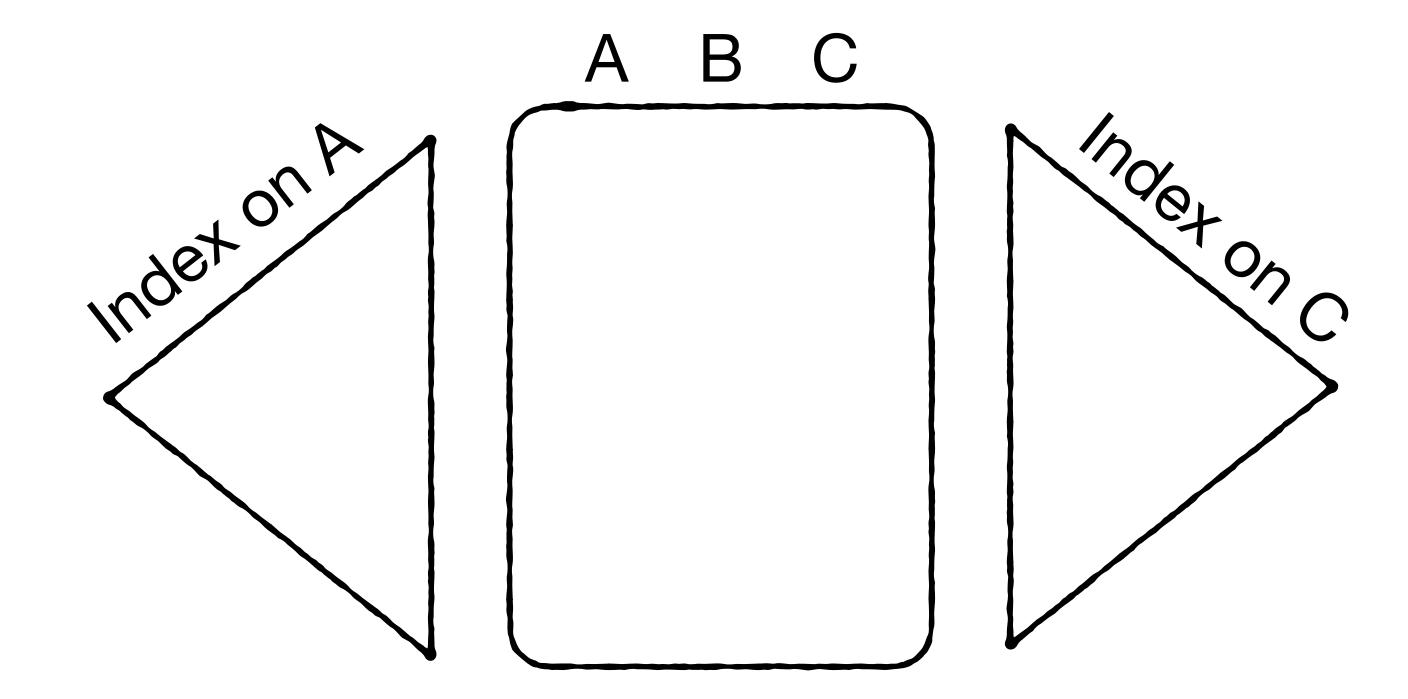
Protect: DB content In-memory data structures (LRU queue)

Duration: Transaction Critical section

Implementation: Lock manager Spin-locks

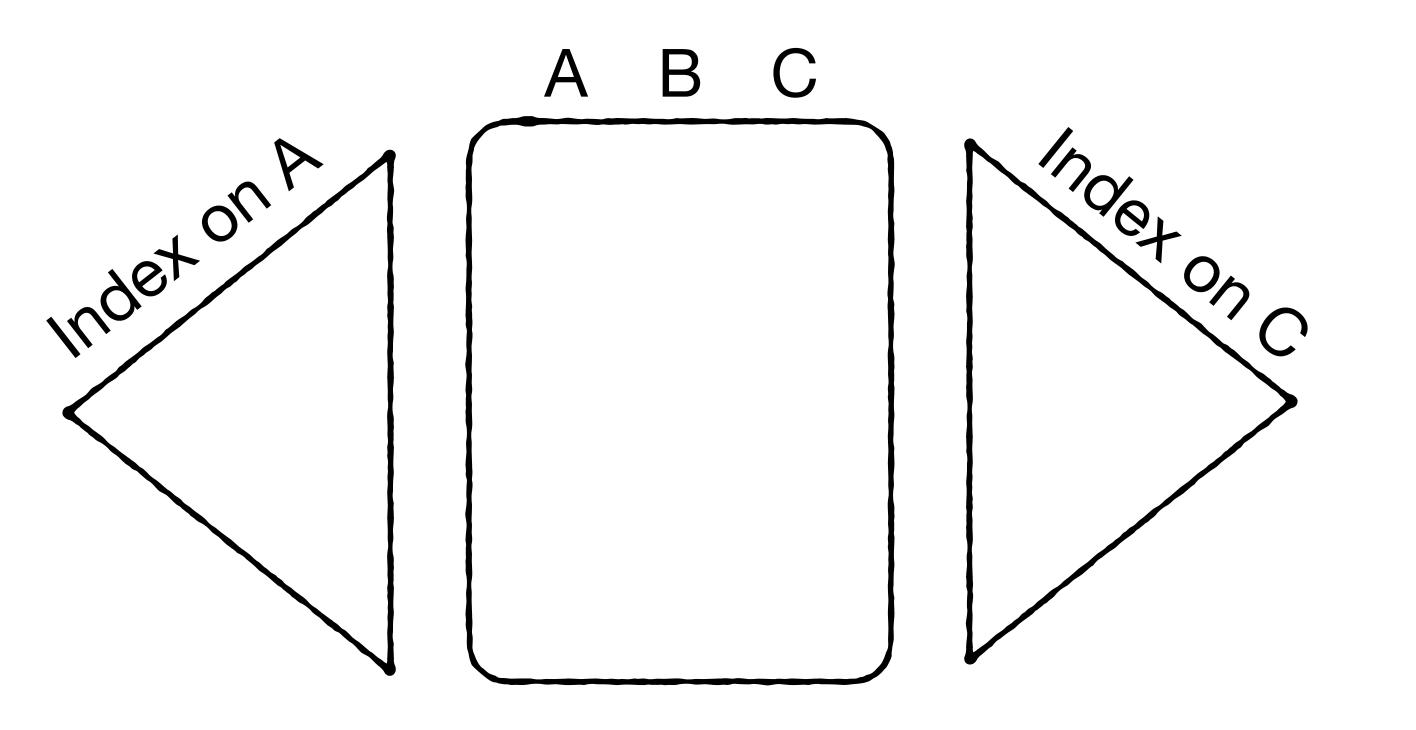
The DB modifies relevant indexes as a part of a transaction

e.g., update table set $C = \dots$ where $A \dots$



The DB modifies relevant indexes as a part of a transaction

e.g., update table set $C = \dots$ where $A \dots$



Schedule

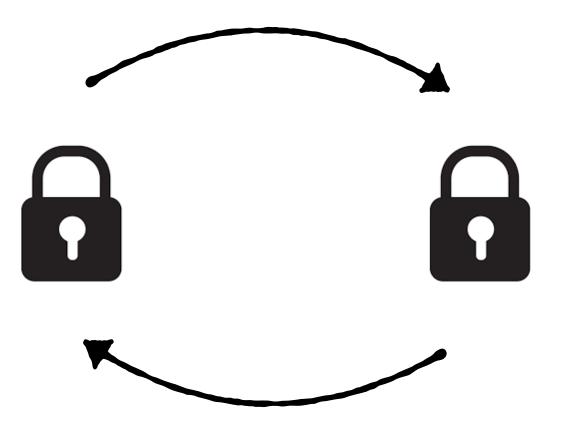
- (1) Begin transaction
- (2) Search index A
- (3) Access row and update C
- (4) Update entry in C's index
- (5) Commit

Anything to make you uneasy about two-phase locking so far?



Anything to make you uneasy about two-phase locking so far?

Deadlocks: transactions waiting on each other's locks forever



T1: A to B

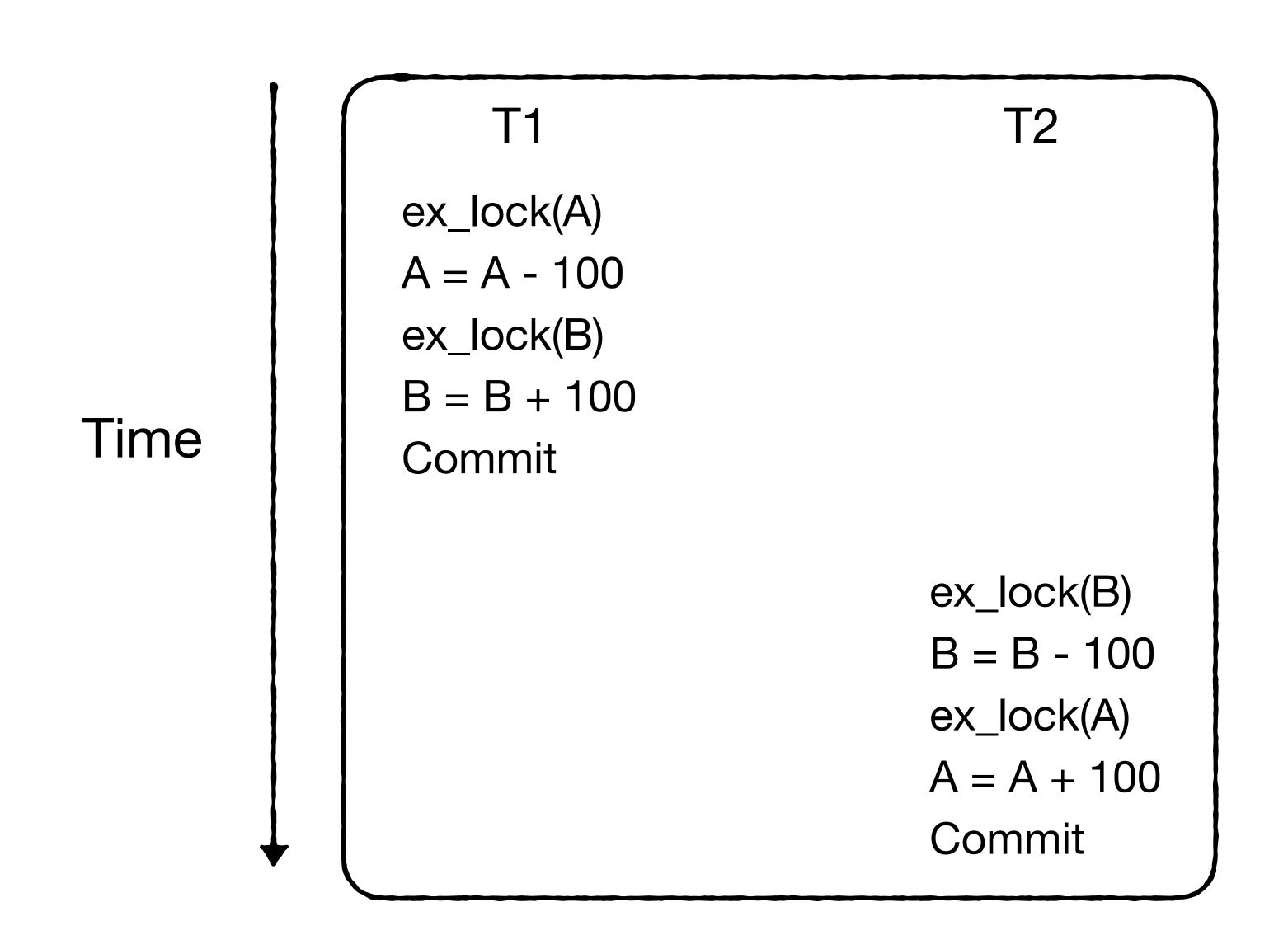
A = A - 100

B = B + 100

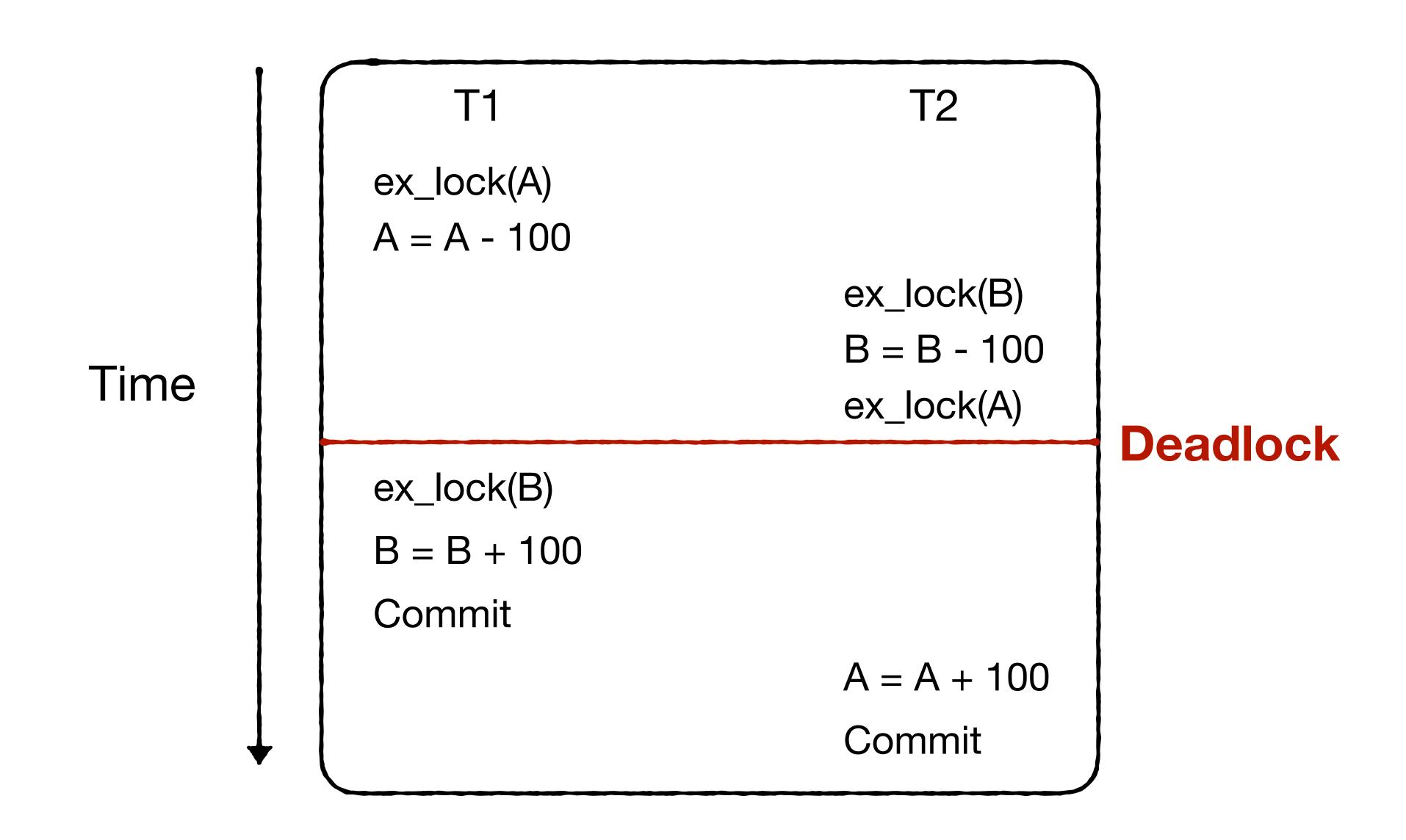
T2: B to A

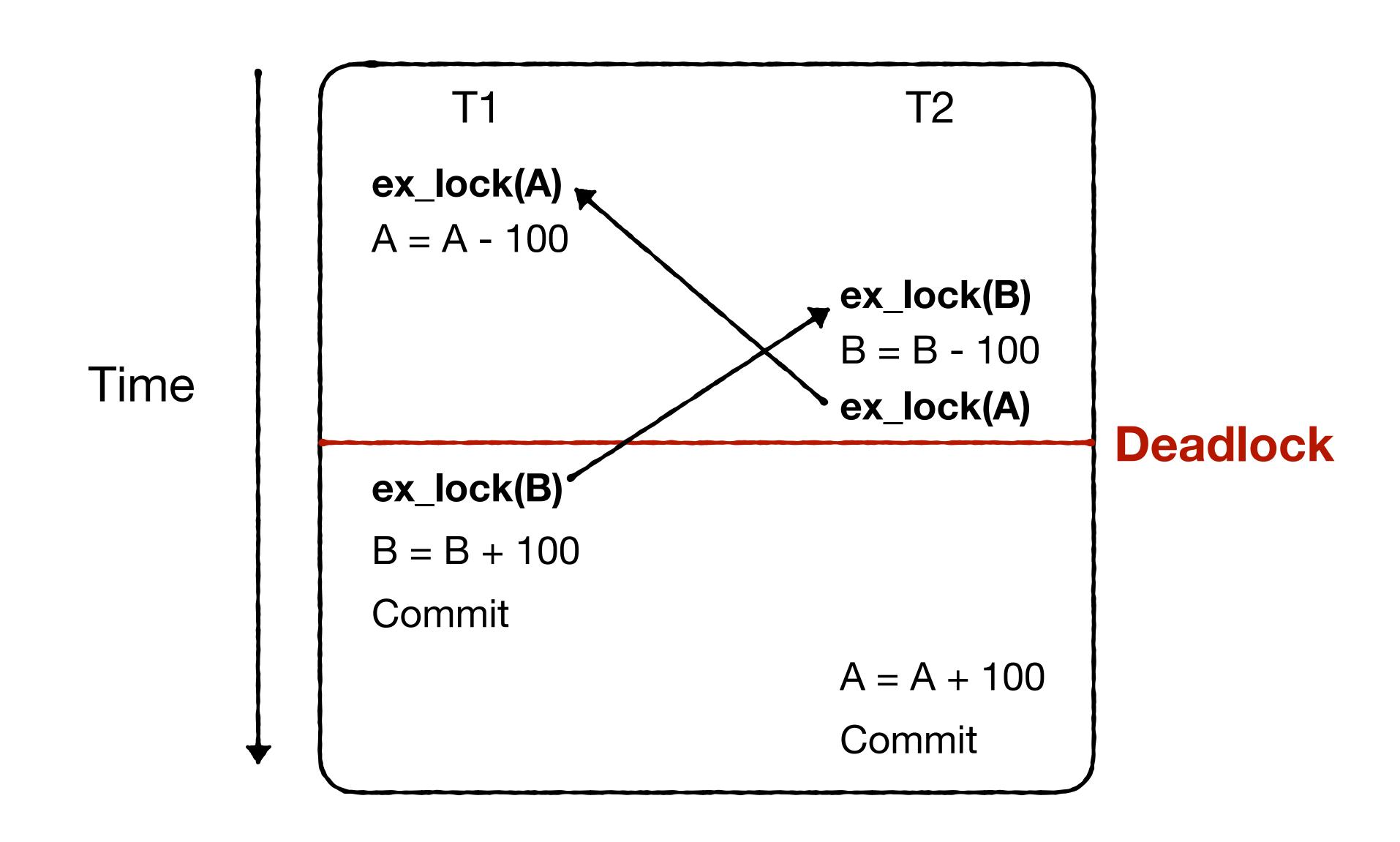
B = B - 100

A = A + 100



This is ok

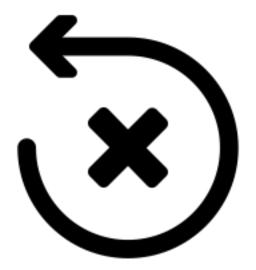




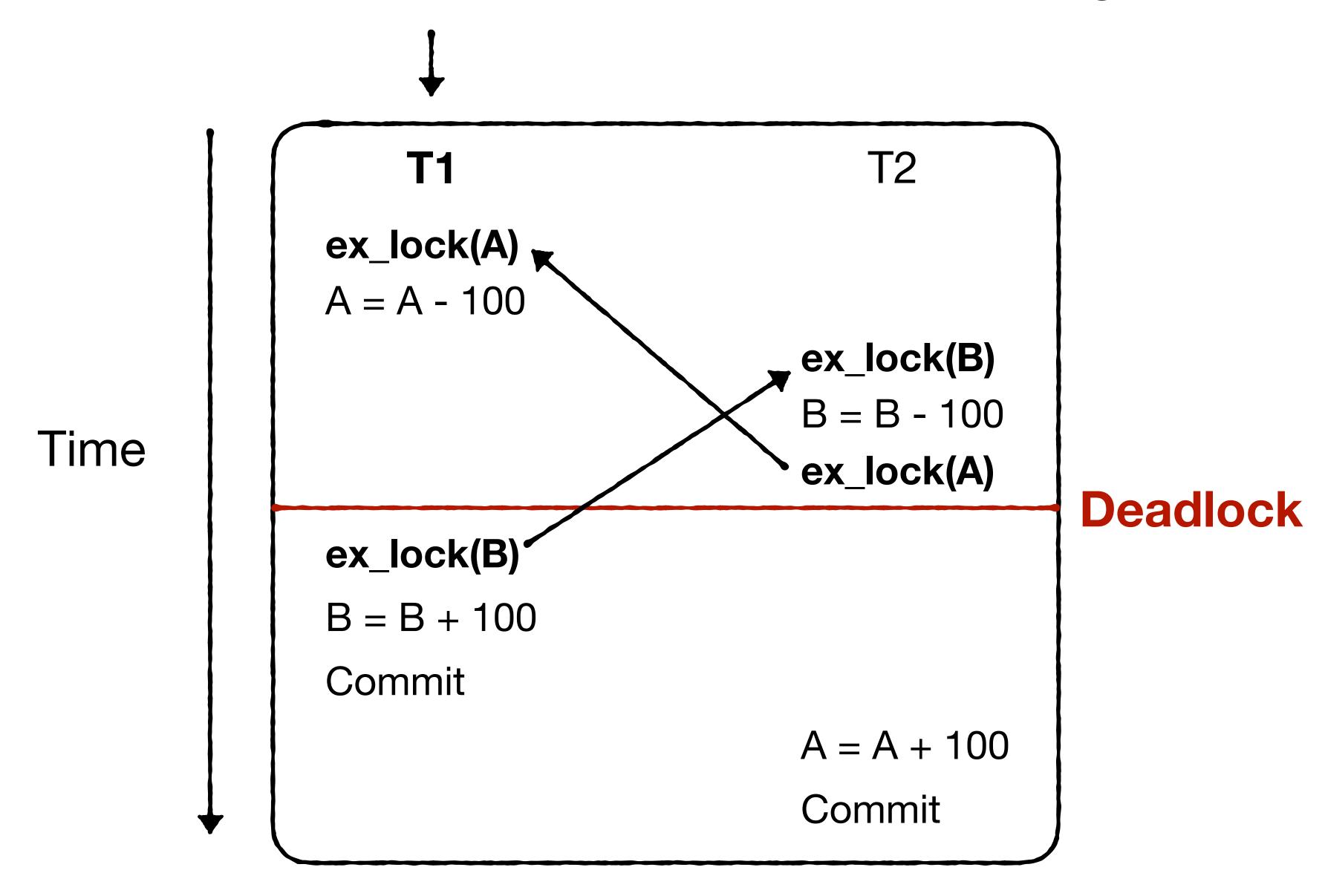
What should to do when we detect a deadlock?

What should to do when we detect a deadlock?

Abort transaction & undo all its changes



Abort transaction & undo all its changes



Abort transaction & undo all its changes



Abort T1

T2

ex_lock(A)

How to undo?

A = A - 100

ex_lock(B)

B = B - 100

ex_lock(A)

ex_lock(B)

B = B + 100

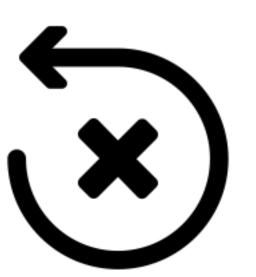
Commit

A = A + 100

Commit

Deadlock

How to undo?



How to undo?

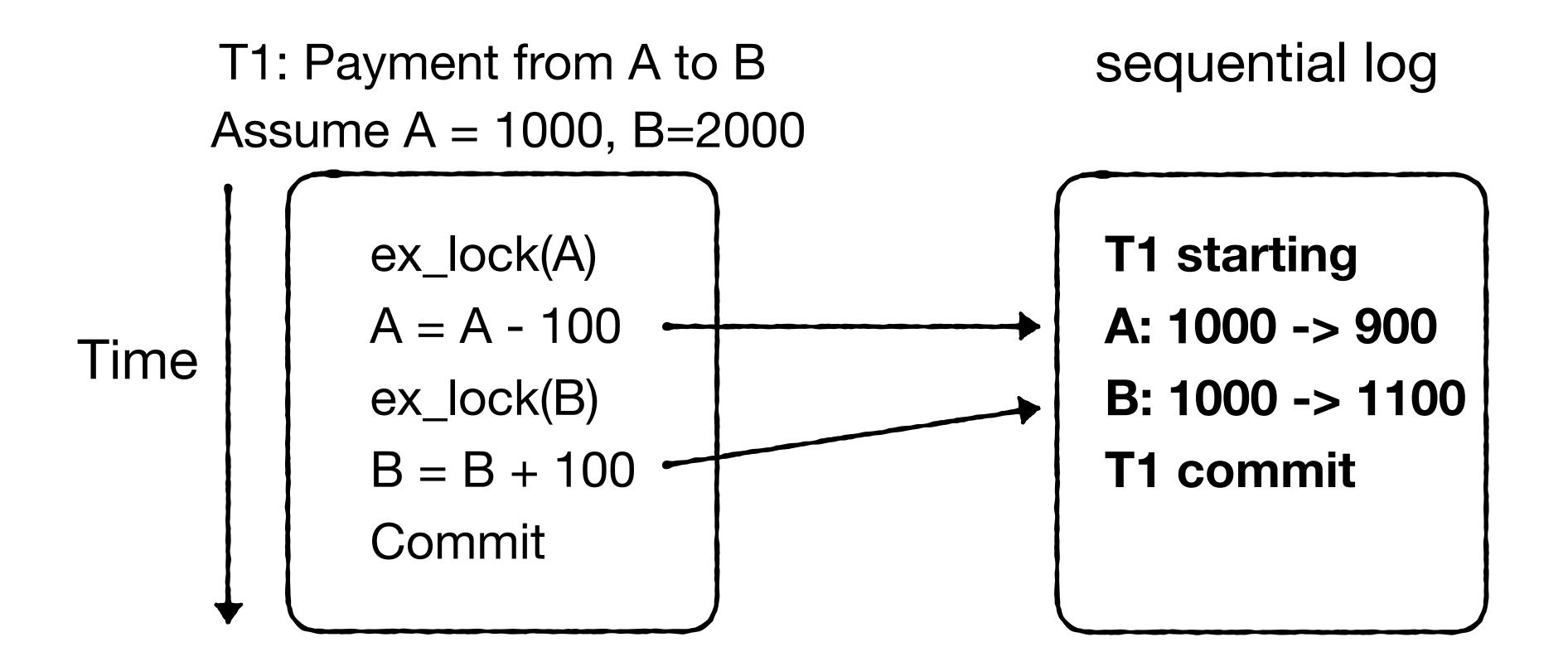
record before-image for all changes to the DB in a sequential log.

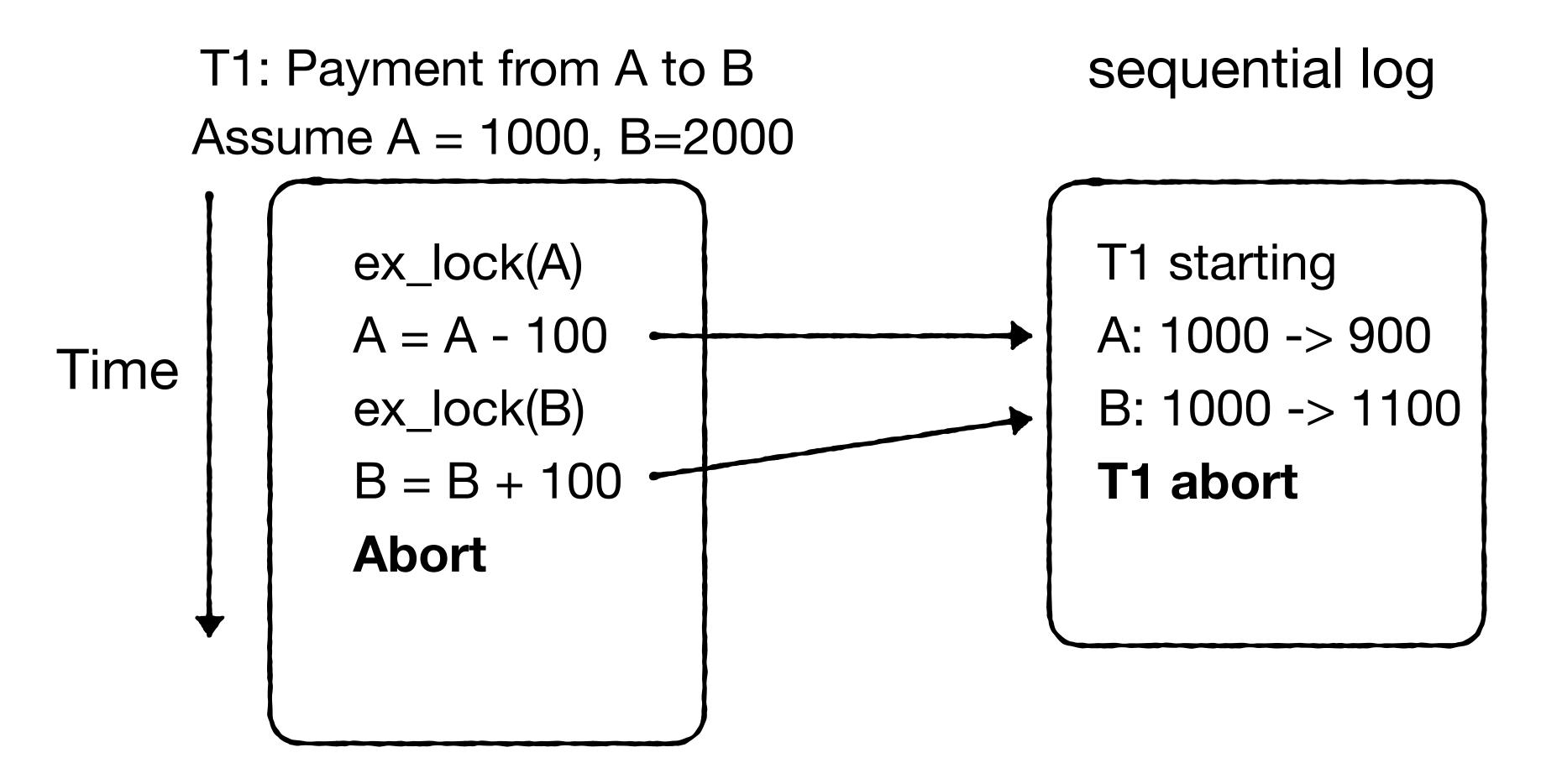


T1: Payment from A to B

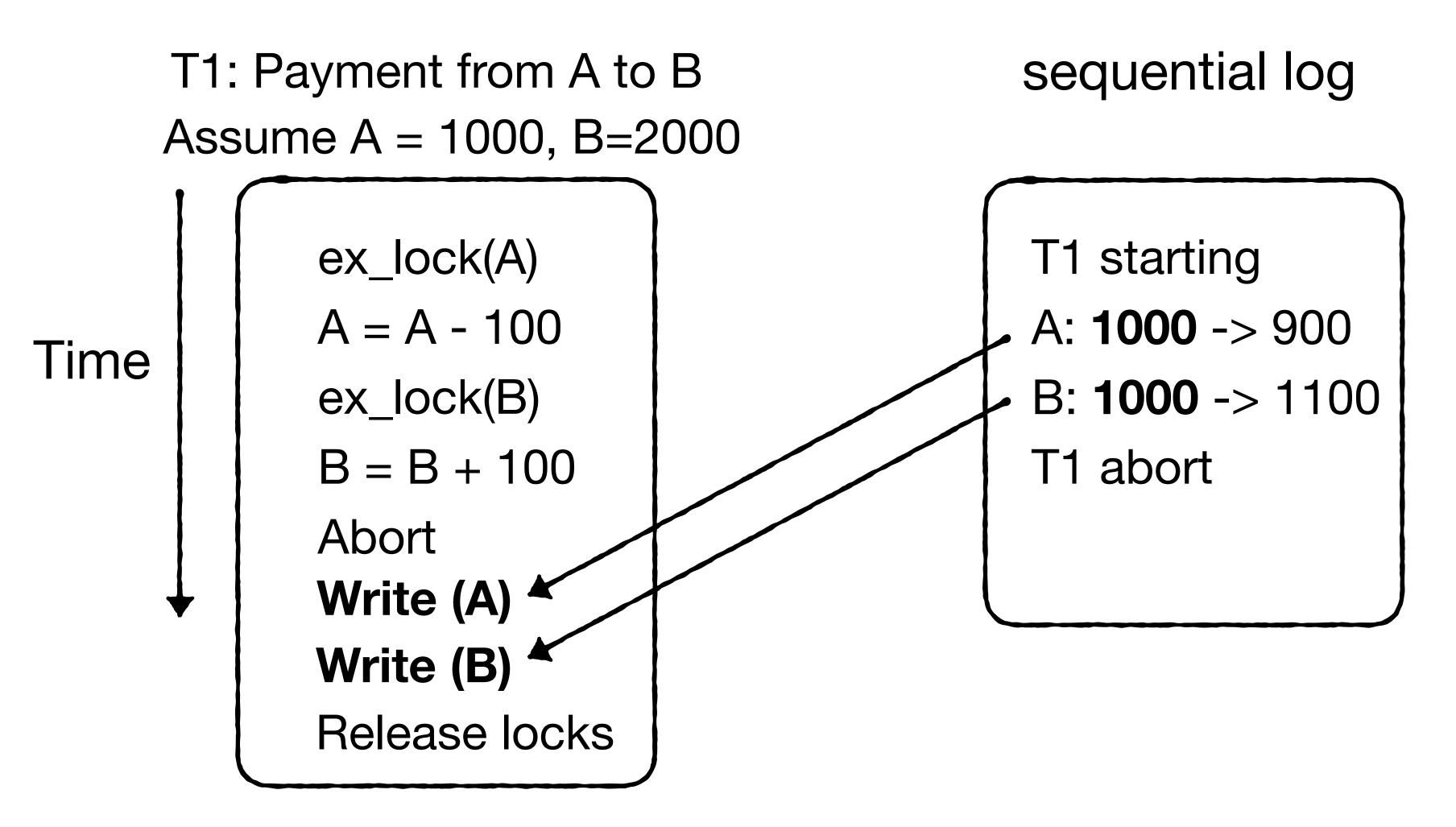
Assume A = 1000, B=2000

record before-image for all changes to the DB in a sequential log.





If transaction aborts before completing, we undo its changes via its before-images in the log



If transaction aborts before completing, we undo its changes via its before-images in the log

How can we detect & prevent deadlocks?

How can we detect & prevent deadlocks?

Timeouts

Conservative Two Phase Locking

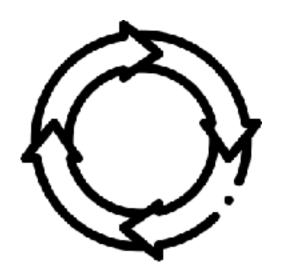
Abort on wait

Cycle Detection









Timeouts



Abort a transaction after waiting a certain time for a lock

Pros?

Cons?

Timeouts



Abort a transaction after waiting a certain time for a lock

Pros? Simple

Cons? Speculative - wastes work aborting on non-deadlocks

How can we detect & prevent deadlocks?

Timeouts

Conservative Two Phase Locking

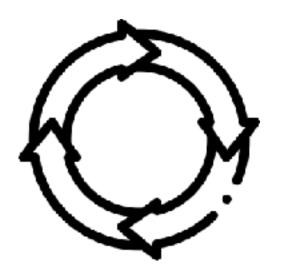
Abort on wait

Cycle Detection









Conservative Two Phase Locking

A transaction is divided into two phases

Phase 1

Take all locks the transaction could possibly need

Phase 2

Release all locks during commit

Conservative Two Phase Locking

A transaction is divided into two phases

Phase 1

Phase 2

Take all locks the transaction could possibly need

Release all locks during commit

Pros: no deadlocks

Con:

Conservative Two Phase Locking

A transaction is divided into two phases

Phase 1

Phase 2

Take all locks the transaction could possibly need

Release all locks during commit

Pros: no deadlocks

Con: takes more locks and holds them for longer

How can we detect & prevent deadlocks?

Timeouts

Conservative Two
Phase Locking

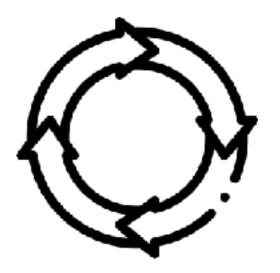
Abort on Wait

Cycle Detection



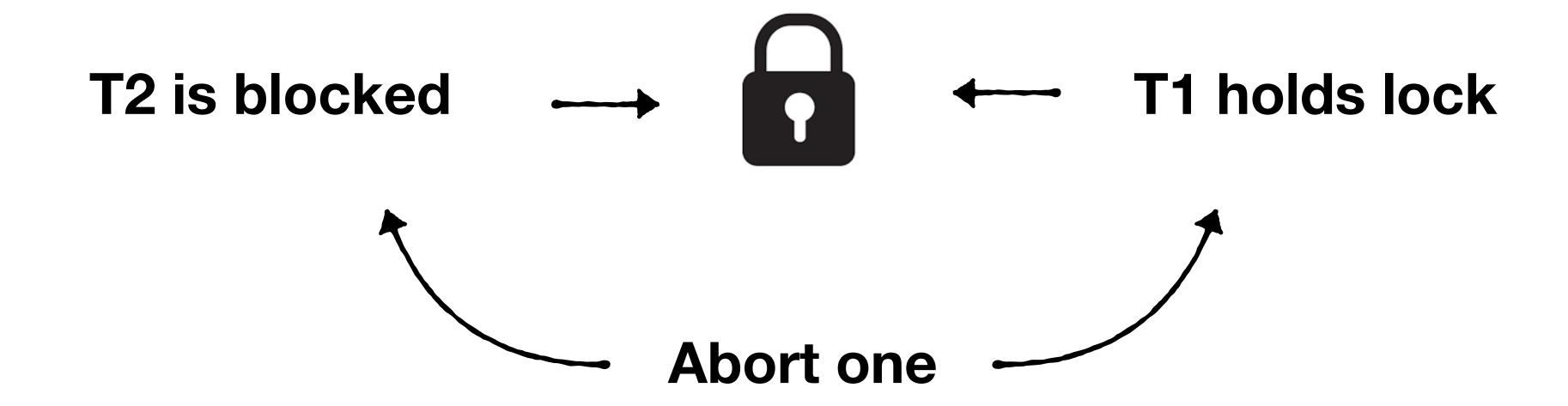






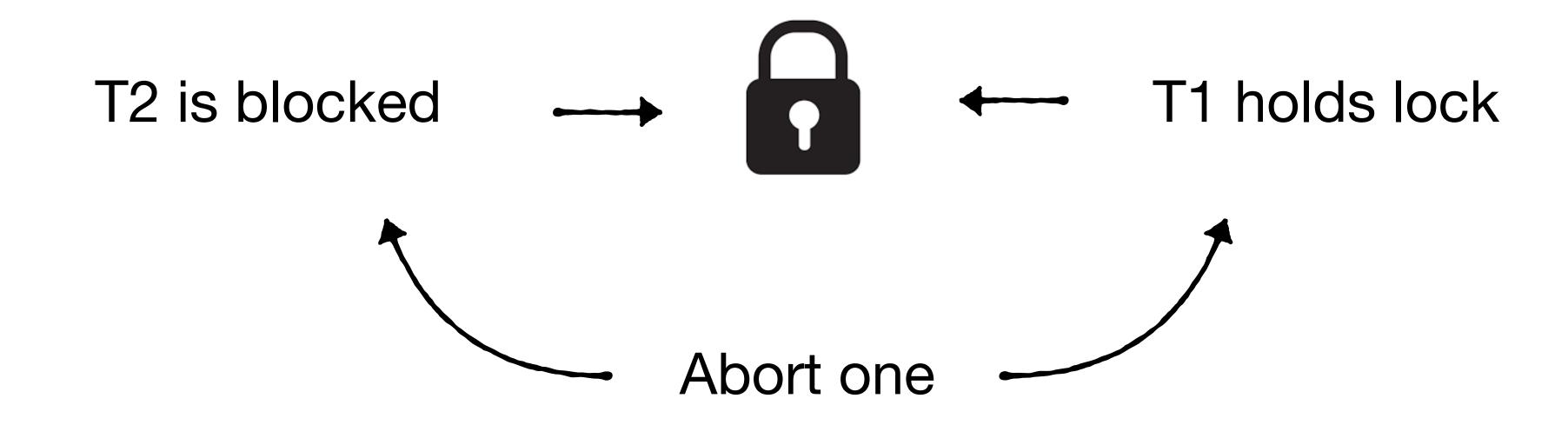
Abort on Wait

When a transaction is blocked, abort it or the transaction holding the lock.



Abort on Wait

When a transaction is blocked, abort it or the transaction holding the lock.



Pros: no deadlocks

Con: defeatist (wastes work, or aborts many transactions)

How can we detect & prevent deadlocks?

Timeouts

Conservative Two
Phase Locking

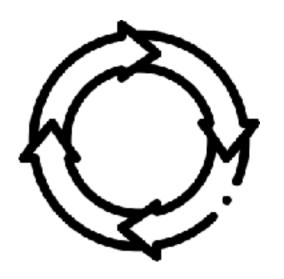
Abort on wait

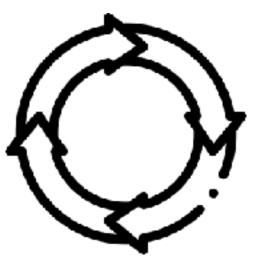
Cycle Detection

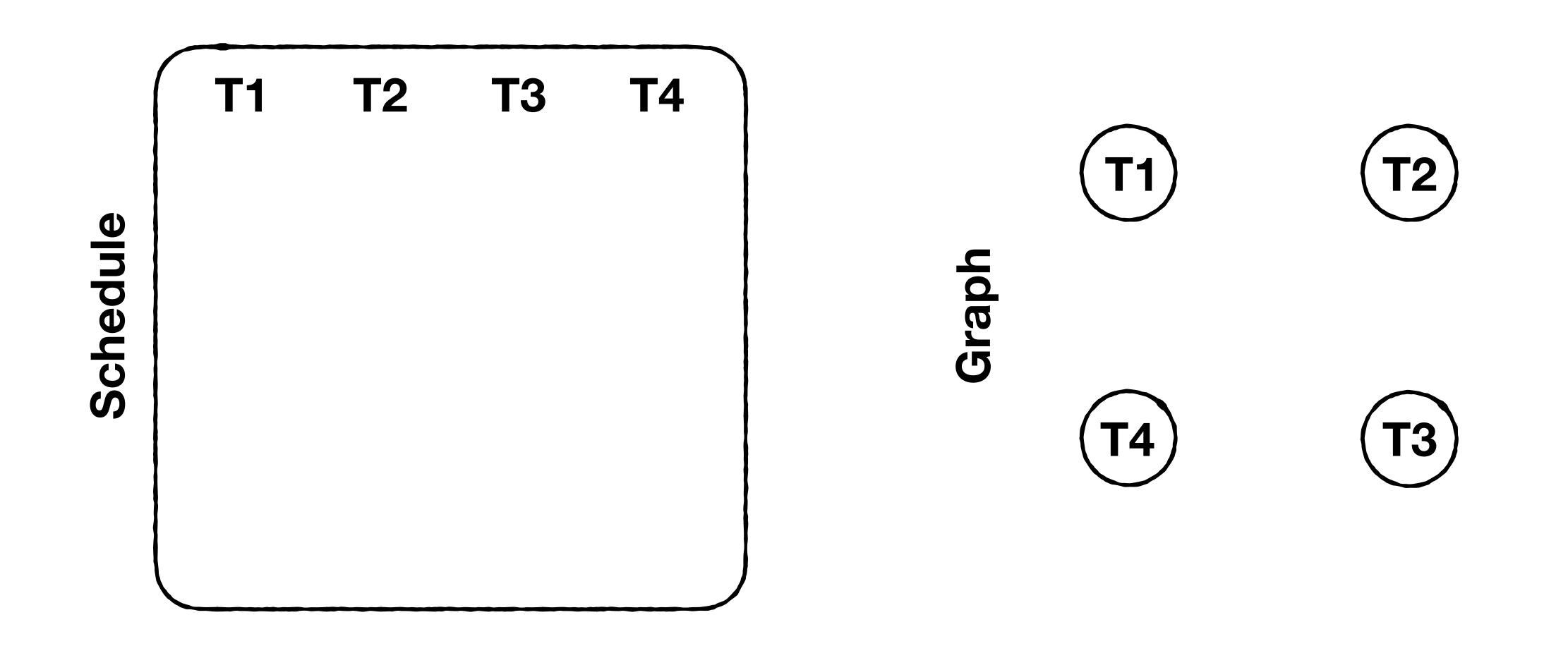


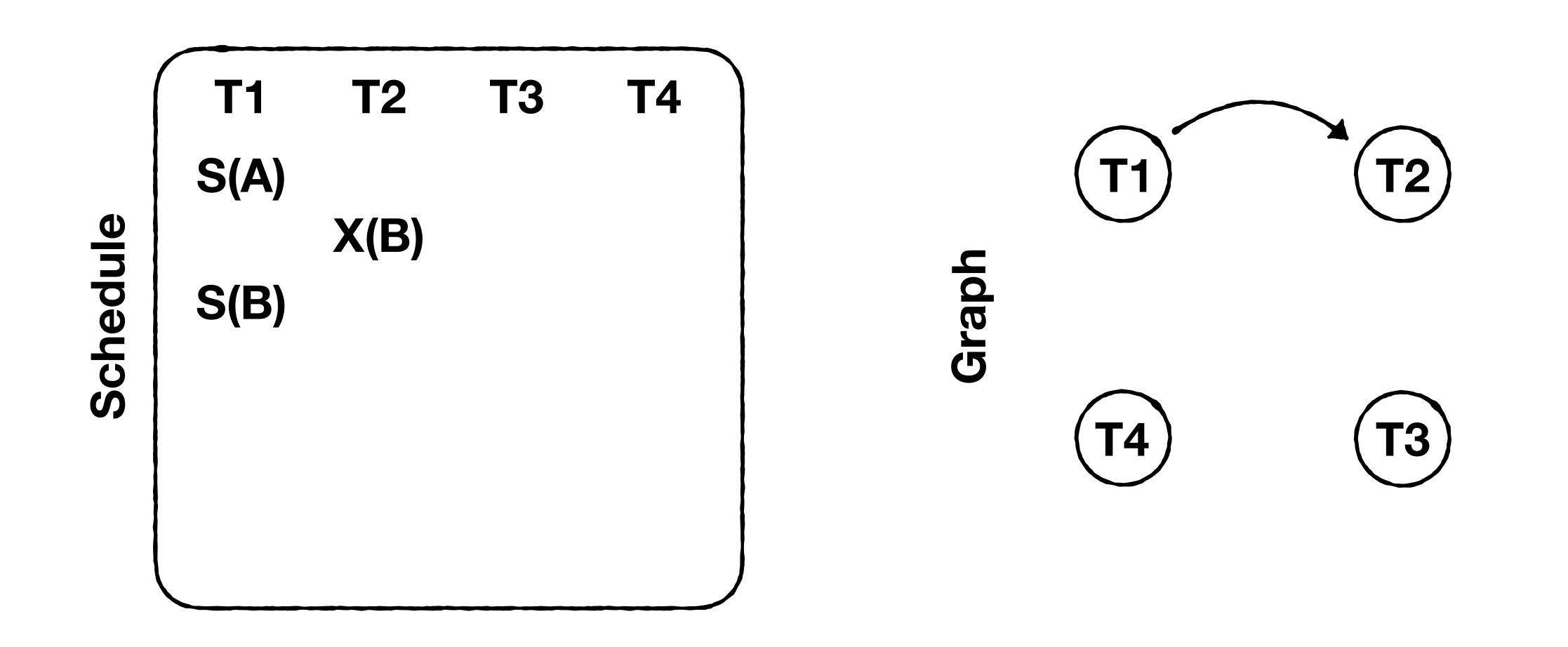


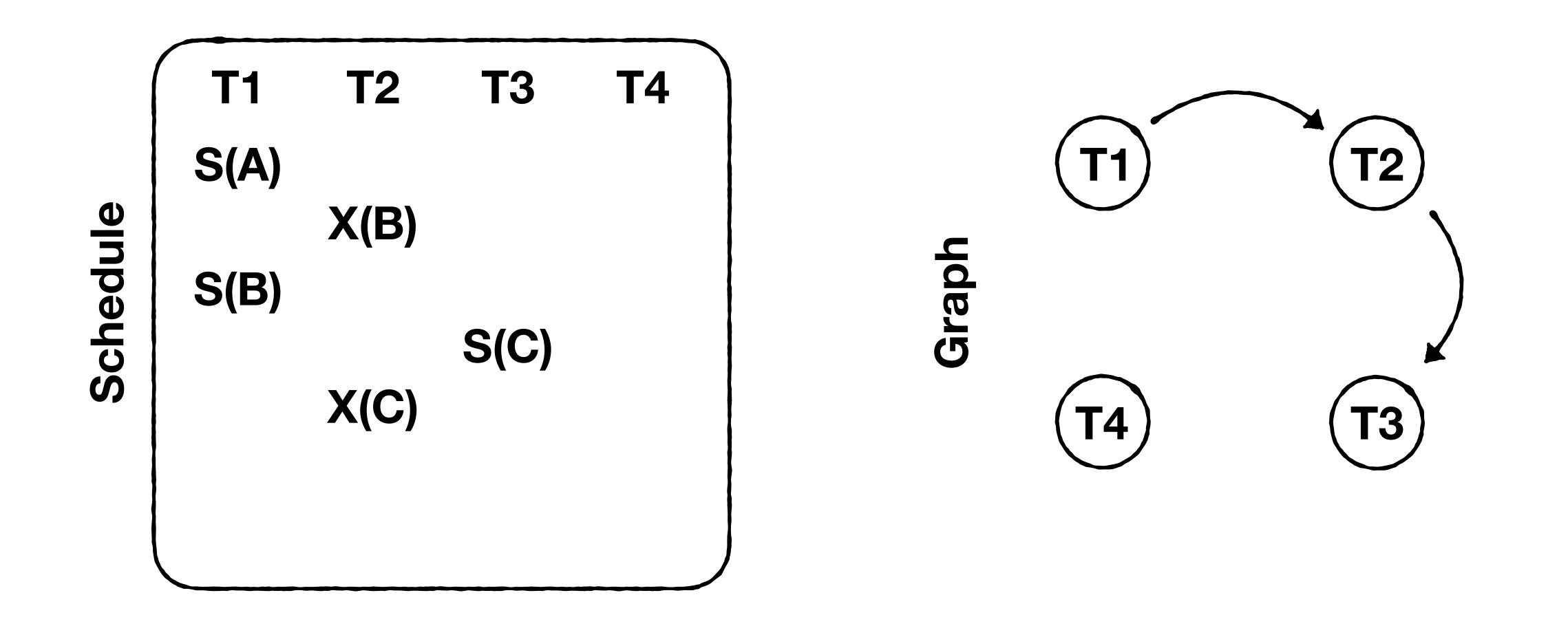


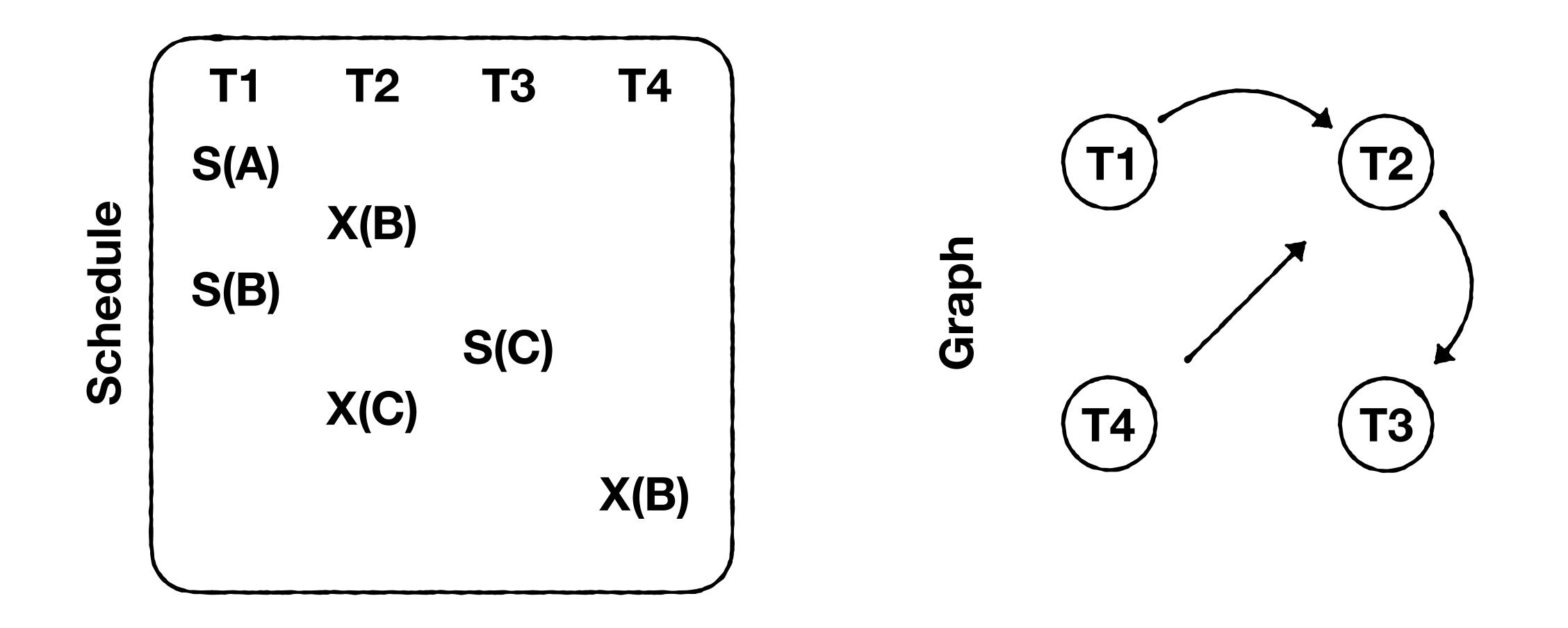


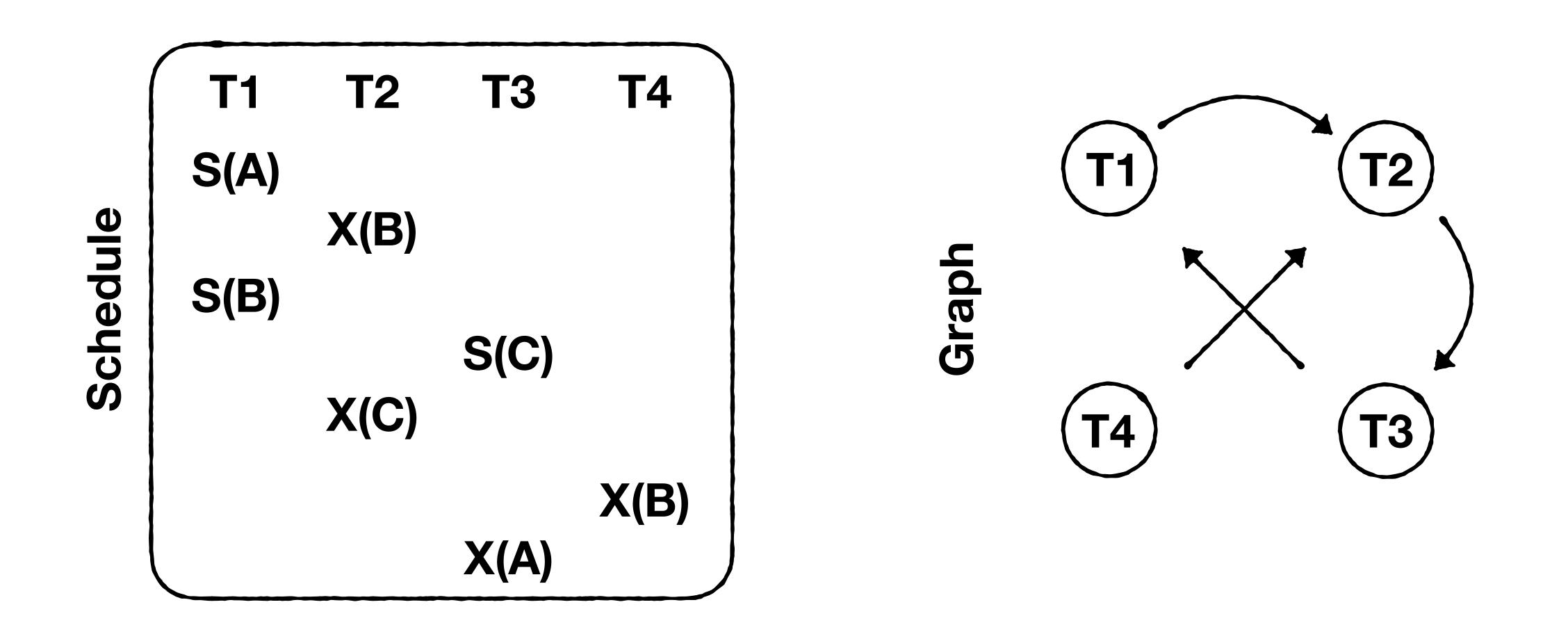








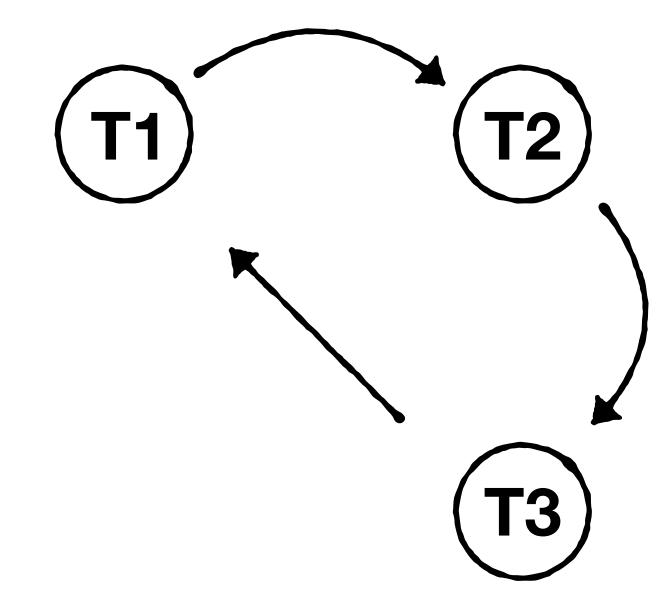




The DB can maintain graph of transactions waiting on each other.

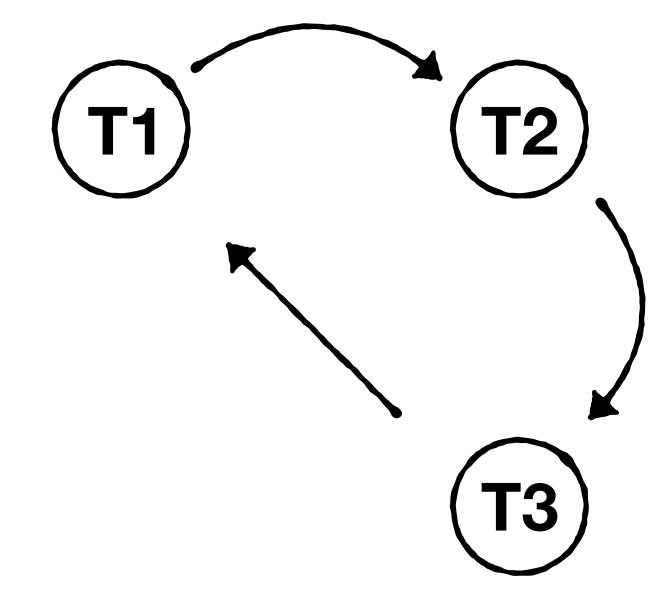
Other approaches are speculative

Here we detect a deadlock for sure



The DB can maintain graph of transactions waiting on each other.

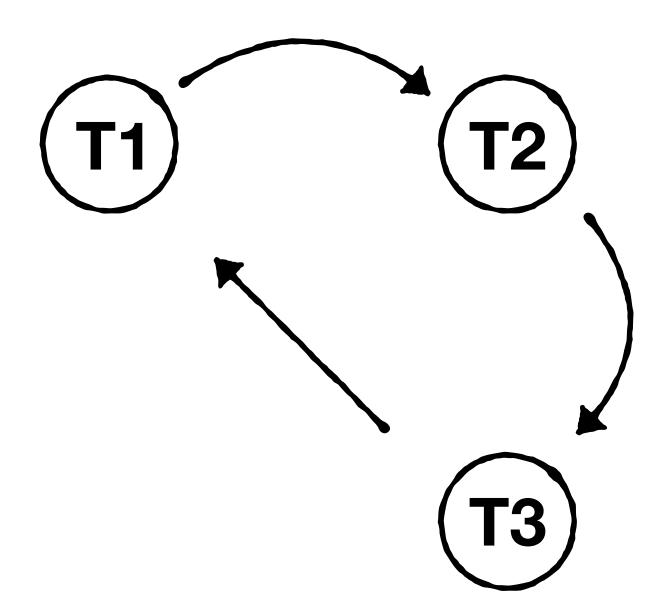
Should we abort T1, T2 or T3?



The DB can maintain graph of transactions waiting on each other.

Abort the transaction that:

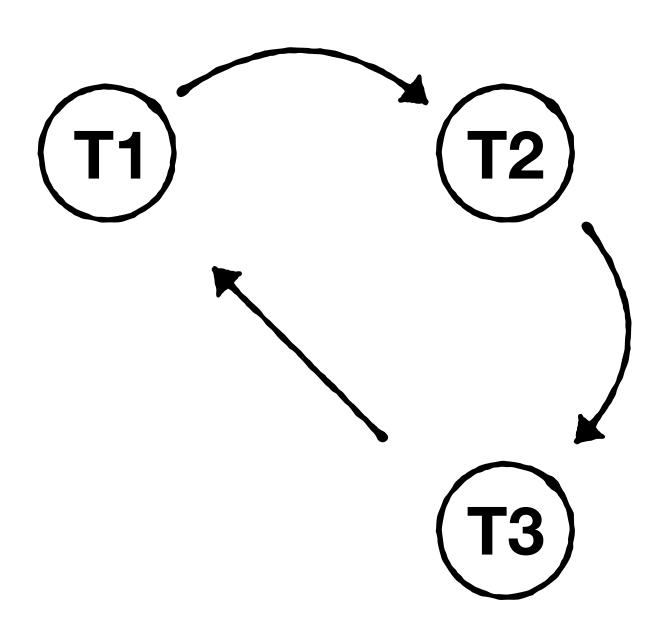
(1) has done the least work



The DB can maintain graph of transactions waiting on each other.

Abort the transaction that:

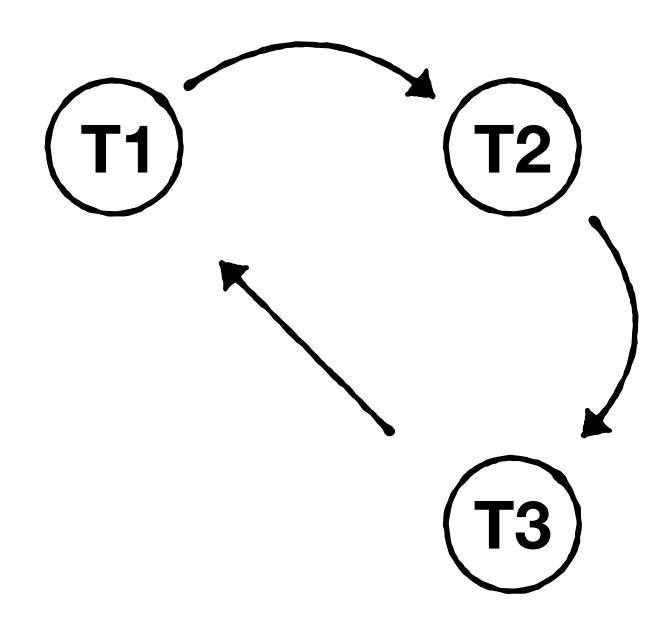
- (1) has done the least work
- (2) Farthest from completion



The DB can maintain graph of transactions waiting on each other.

Abort the transaction that:

- (1) has done the least work
- (2) Farthest from completion
- (3) Been aborted the least times



Timeouts

Conservative Two Phase Locking

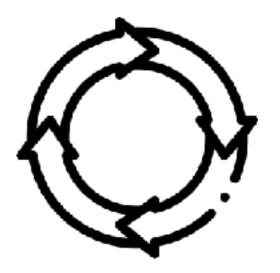
Abort on wait

Cycle Detection









There is still a problem

Let's return to our running examples

Example 1

Interest & Payments

Example 2

Updates & Statistics





Problems:

Dirty reads

Unrepeatable reads

Solution:

Exclusive write locks

Shared read locks

Let's return to our running examples

Example 1
Interest & Payments

Example 2
Updates & Statistics

Example 3 Inserts & Statistics





Dirty reads

Exclusive write locks Shared rea

Unrepeatable reads

Shared read locks

?

?

Example 3

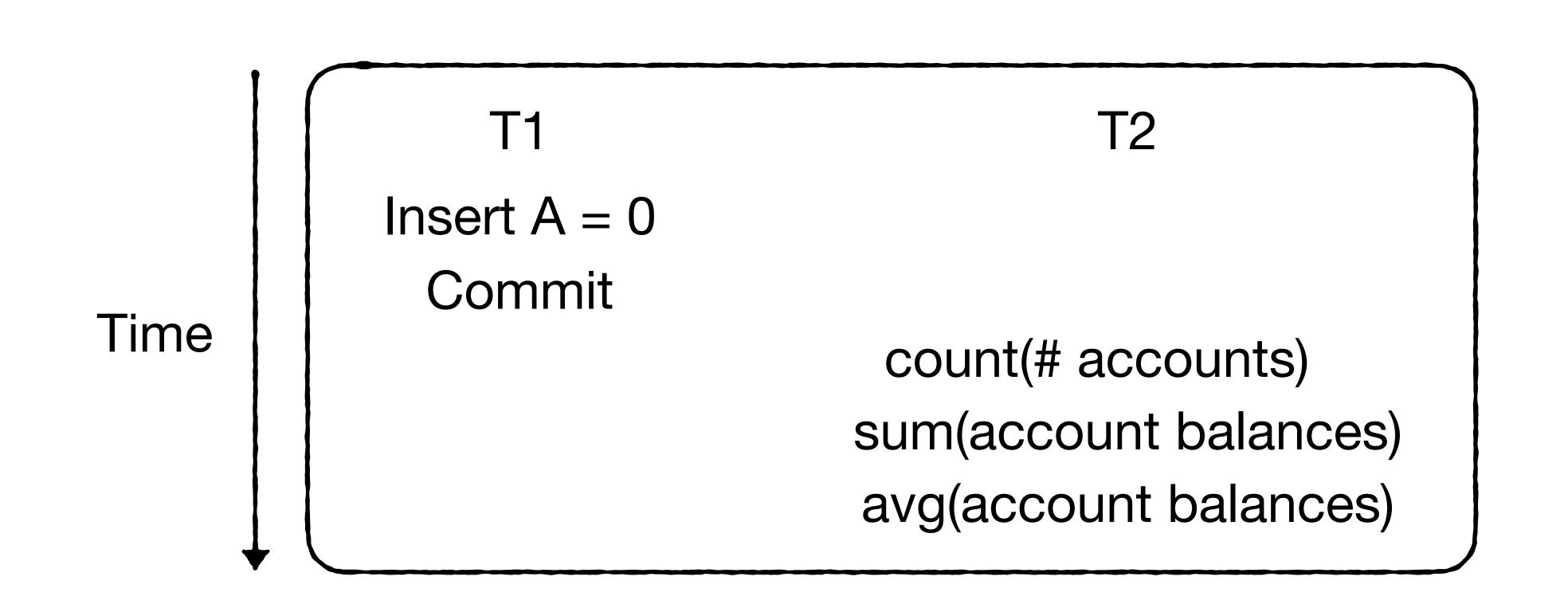
T1: insert account

Insert A = 0

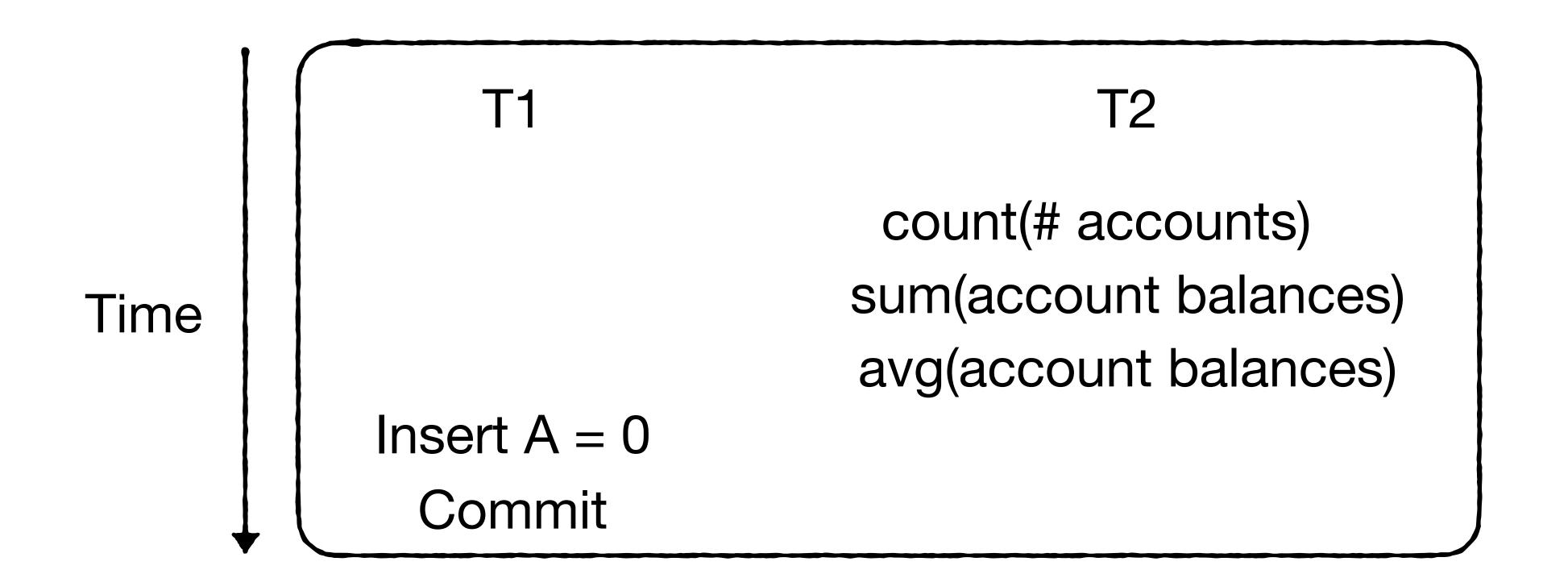
T2: statistics reporting

count(# accounts)
sum(account balances)
avg(account balances)

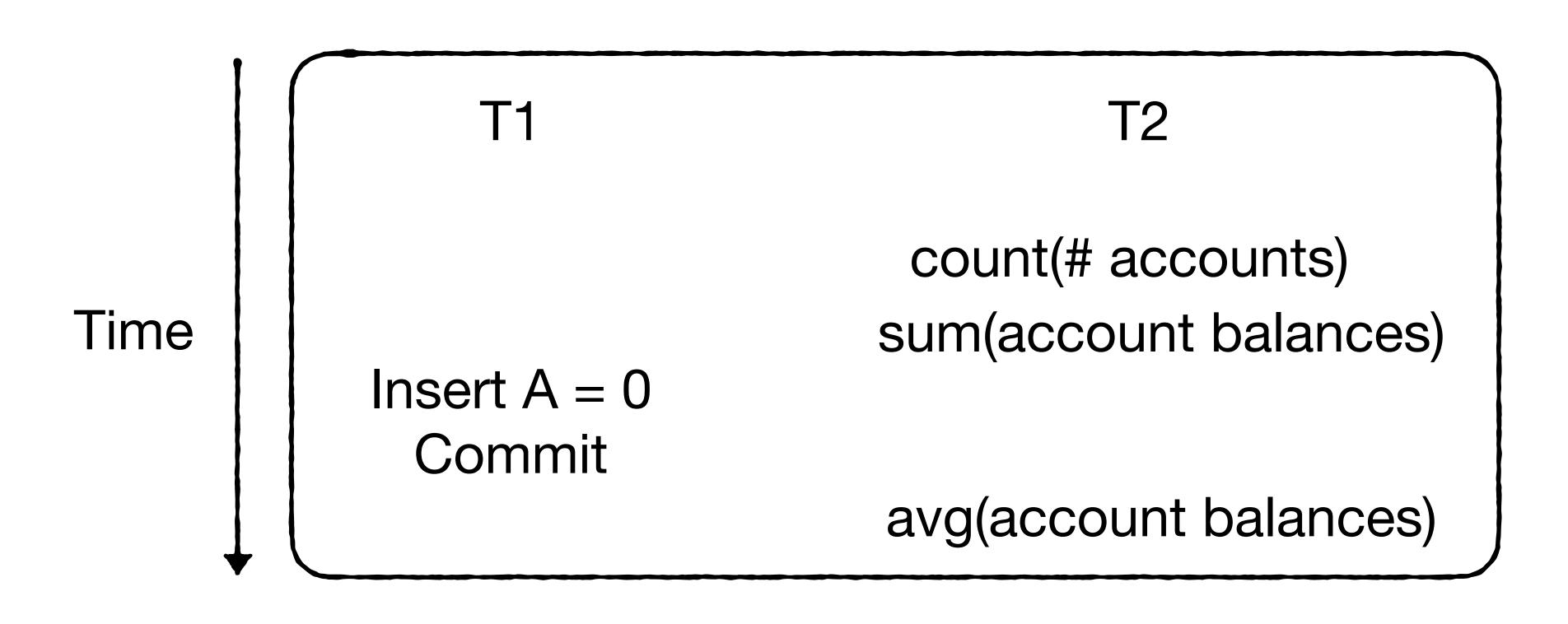
Example 3





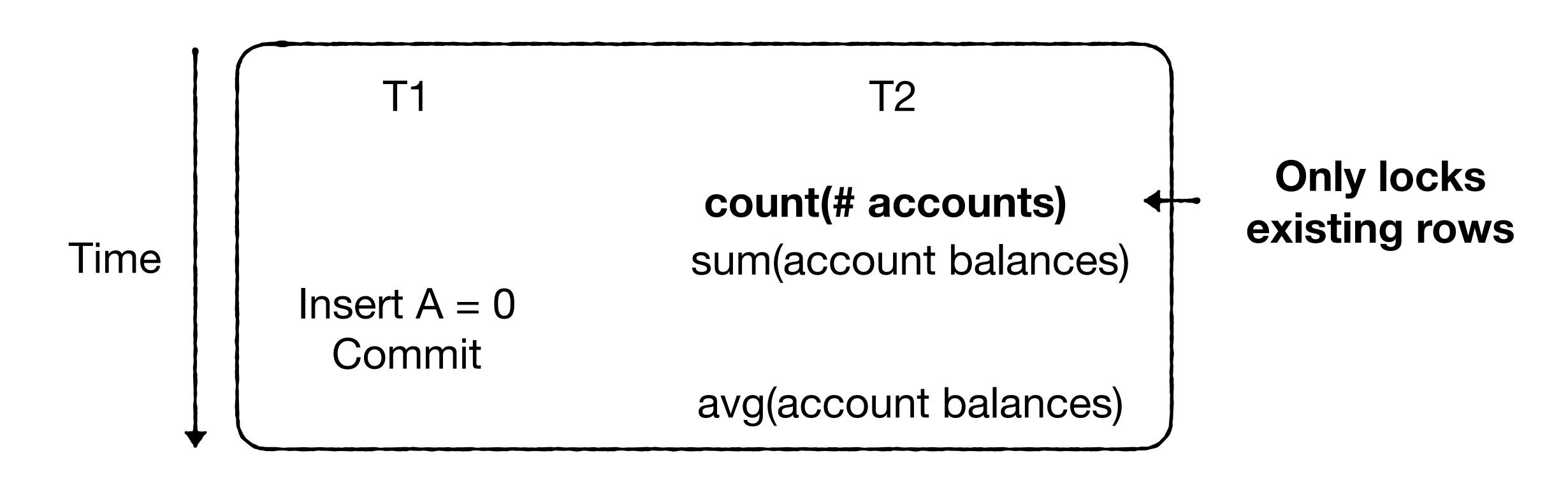




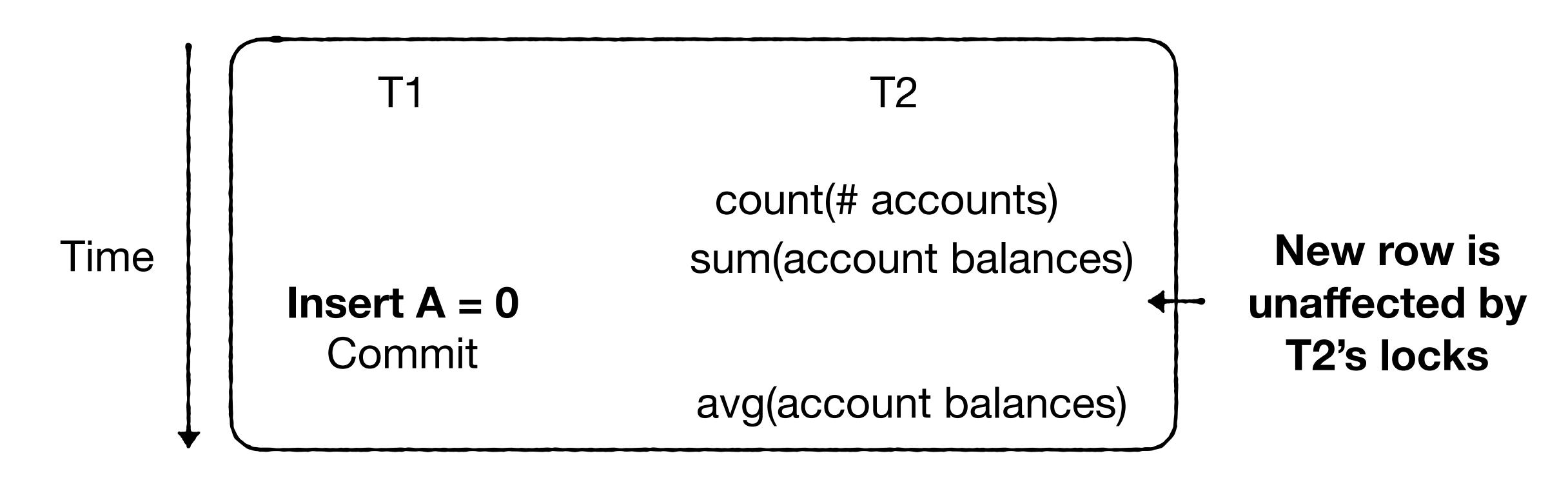




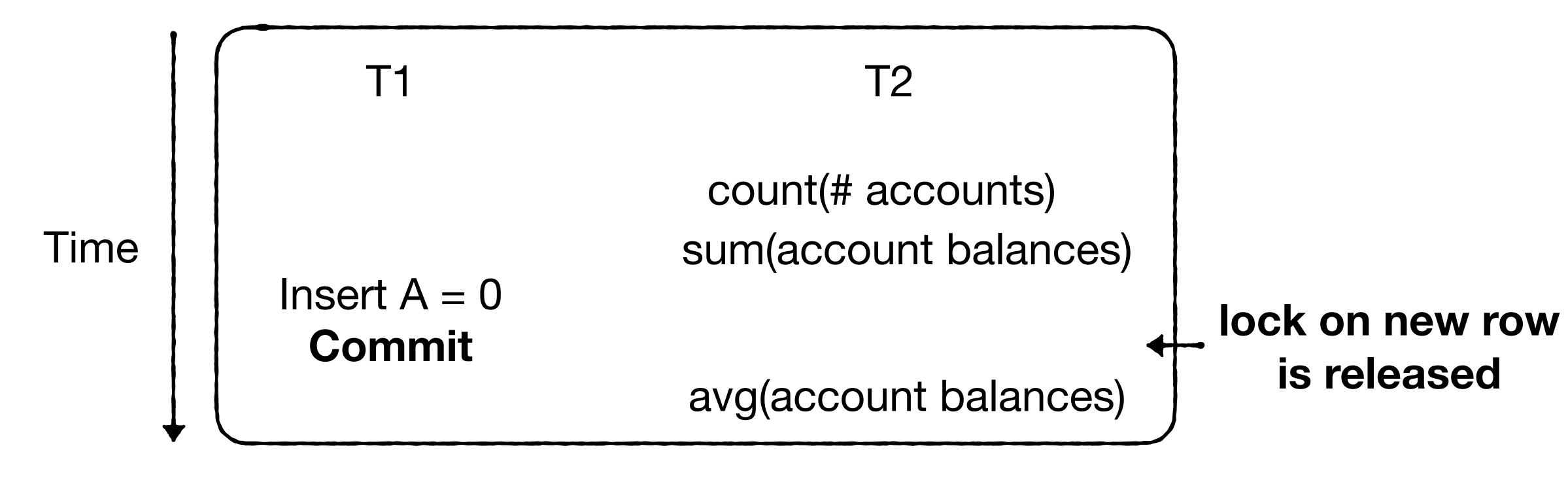
Problem?



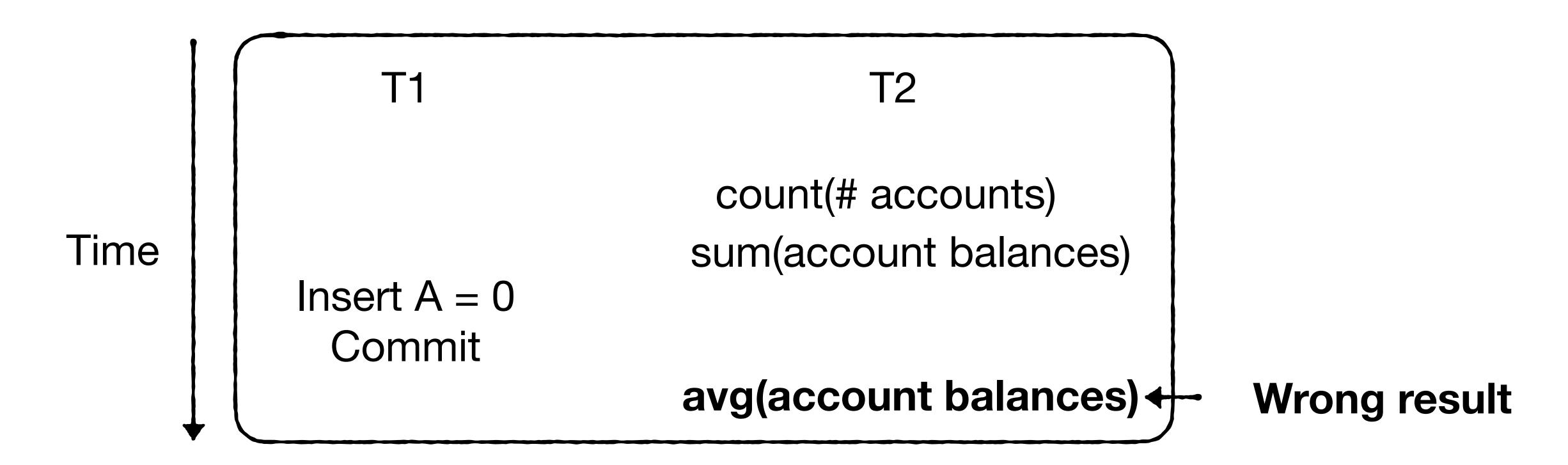
Problem?



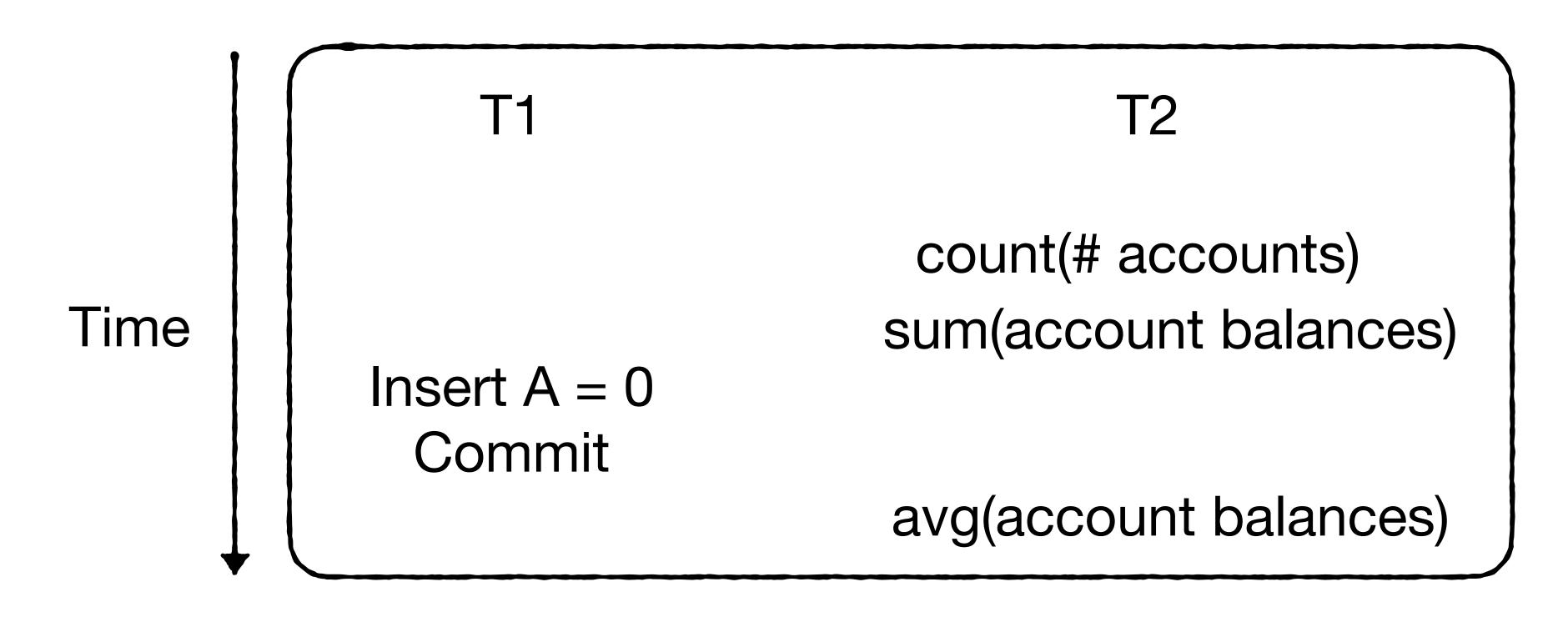
Problem?



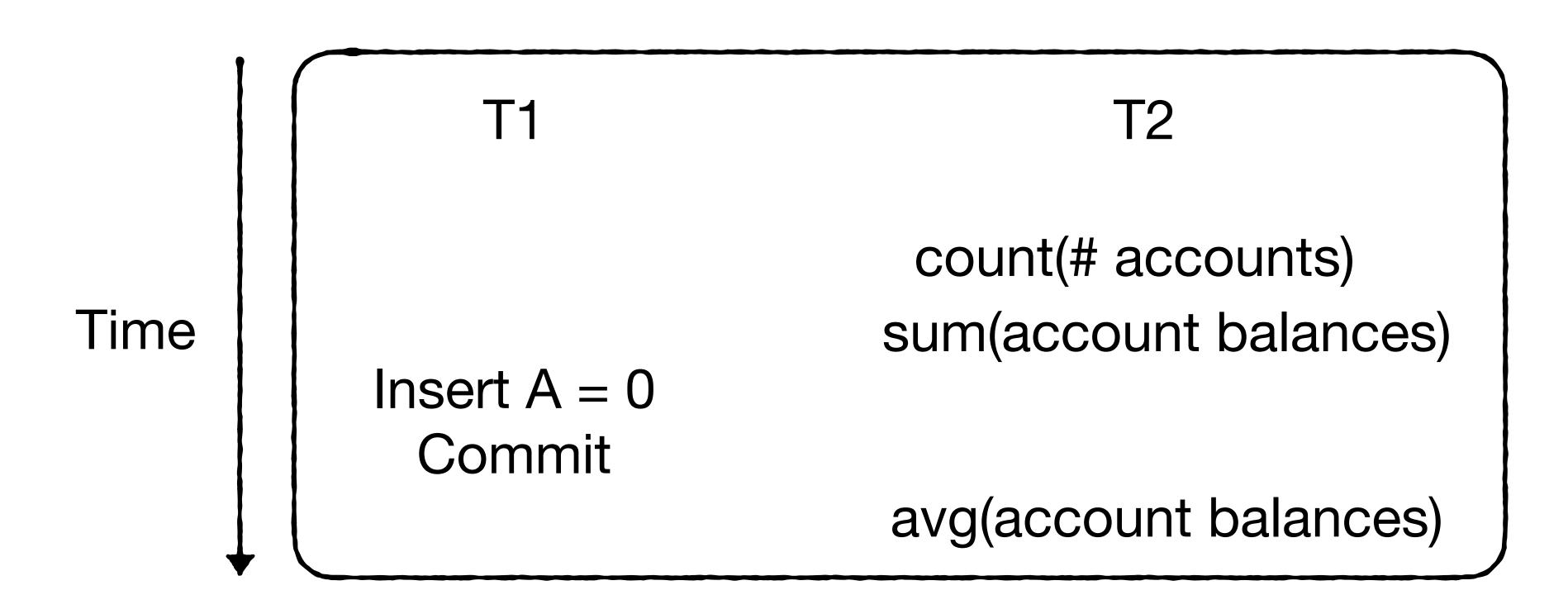
Problem?



Problem? Inconsistent reporting: avg ≠ sum / count



More broadly, this is a **phantom read**: a transaction accesses a set of rows twice, and qualifying rows are added in-between



More broadly, this is a **phantom read**: a transaction accesses a set of rows twice, and qualifying rows are added in-between

lock table



Aggressive

lock table

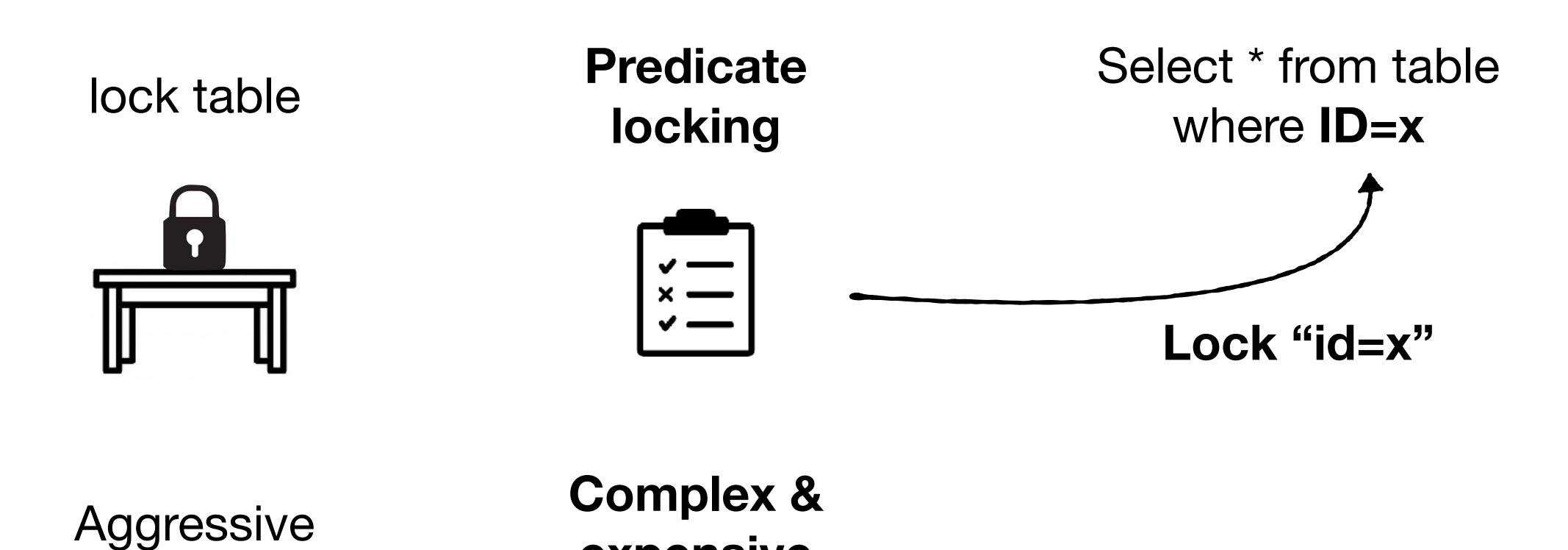


Aggressive

Predicate locking



Complex & expensive



expensive

lock table



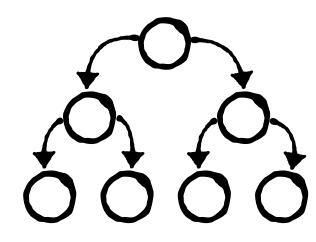
Aggressive

Predicate locking



Complex & expensive

Index locking



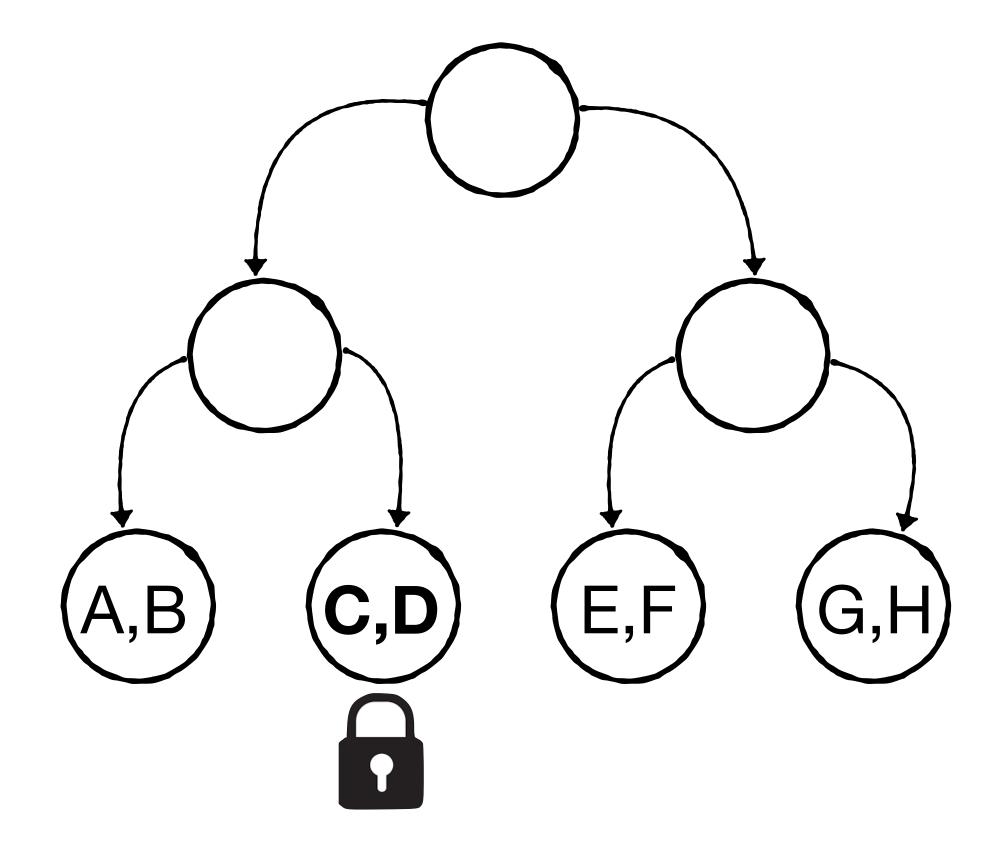
Good but requires an index

Index Locking

Lock B-tree leaf storing relevant range

Example: Select * from accounts

where ID="Cindy"



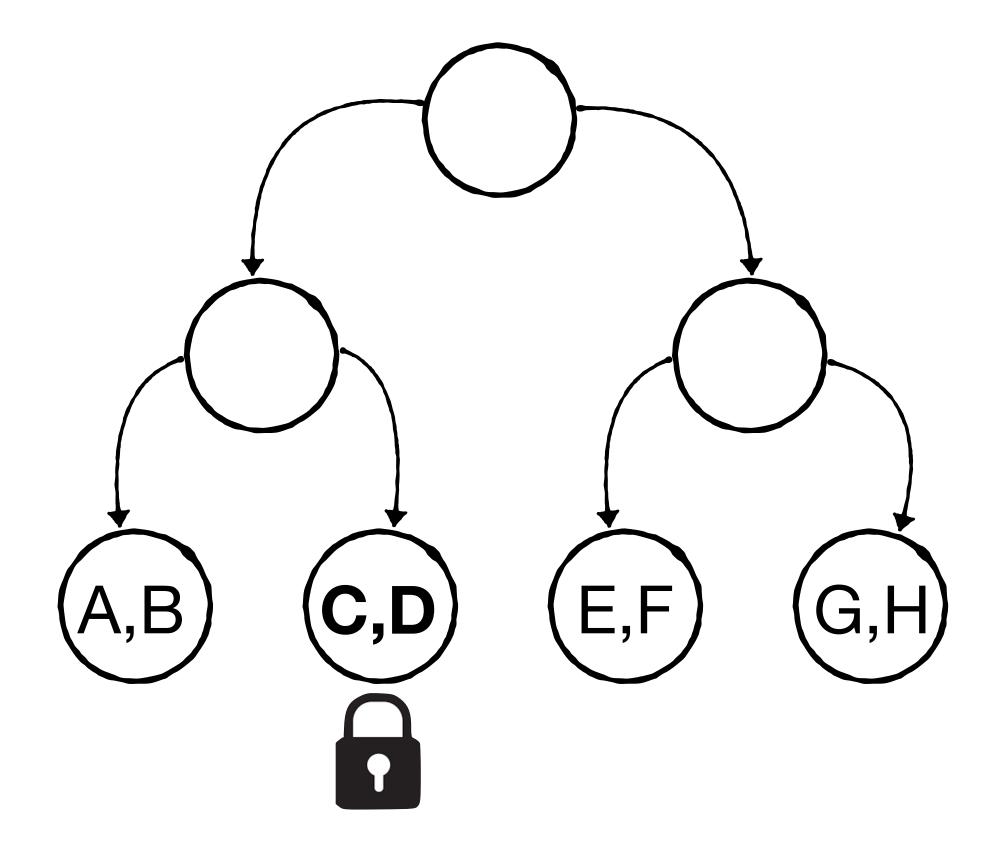
Index Locking

Lock B-tree leaf storing relevant range

Example: Select * from accounts

where ID="Cindy"

So we do not have to lock the whole table as long as we have an index:)



Example 1
Interest & Payments



Dirty reads

Exclusive write locks

Example 2
Updates & Statistics



Unrepeatable reads
Shared read locks

Example 3 Inserts & Statistics



Phantom Read
Range locks

Example 1
Interest & Payments

Example 2
Updates & Statistics

Example 3
Inserts & Statistics







Dirty reads

Exclusive write locks

Unrepeatable reads
Shared read locks

Phantom Read
Range locks

More correct but slower

Example 1
Interest & Payments

Example 2
Updates & Statistics

Example 3
Inserts & Statistics







Dirty reads

Exclusive write locks

Unrepeatable reads
Shared read locks

Phantom Read
Range locks

Compromising correctness is fine for some applications

Dirty read Unrepeatable read Phantom reads

Read uncommitted

Read committed

Repeatable reads

Dirty read Unrepeatable read Phantom reads

Read uncommitted Possible Possible Possible

Read committed

Repeatable reads

	Dirty read	Unrepeatable read	Phantom reads
Read uncommitted	Possible	Possible	Possible

Read committed Not Possible Possible Possible

Repeatable reads

Repeatable reads	Not Possible	Not Possible	Possible
Read committed	Not Possible	Possible	Possible
Read uncommitted	Possible	Possible	Possible
	Dirty read	Unrepeatable read	Phantom reads

Serializable	Not Possible	Not Possible	Not Possible
Repeatable reads	Not Possible	Not Possible	Possible
Read committed	Not Possible	Possible	Possible
Read uncommitted	Possible	Possible	Possible
	Dirty read	Unrepeatable read	Phantom reads

Read uncommitted Exclusive write locks not held until a transaction ends

Read committed

Repeatable reads

Read uncommitted

Exclusive write locks not held until a transaction ends

Read committed

Shared read locks not held until a transaction ends

Repeatable reads

Read uncommitted

Read committed

Repeatable reads

Serializable

Exclusive write locks not held until a transaction ends

Shared read locks not held until a transaction ends

Range locks not employed (e.g., no tree/table locking)

Read uncommitted Exclusive write locks not held until a transaction ends

Read committed Shared read locks not held until a transaction ends

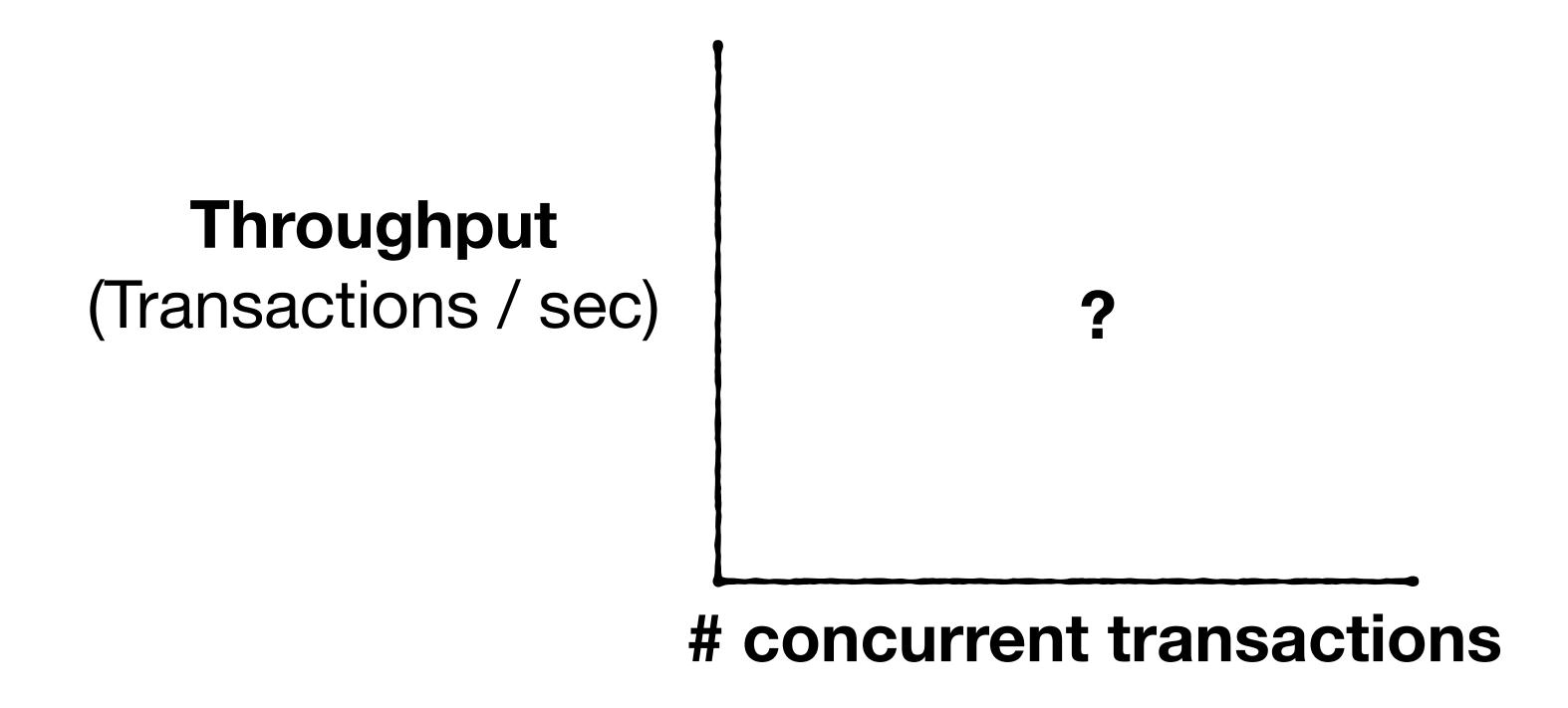
Repeatable reads Range locks not employed (e.g., no tree/table locking)

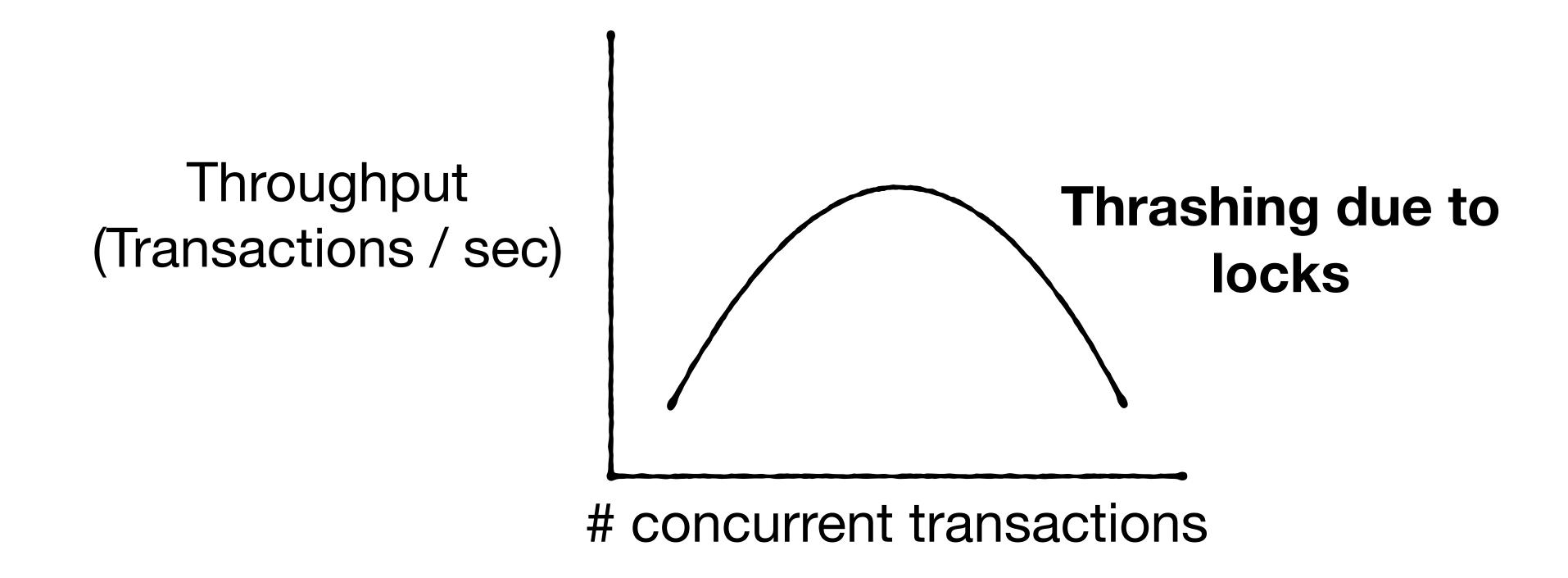
Serializable Everything is fully correct

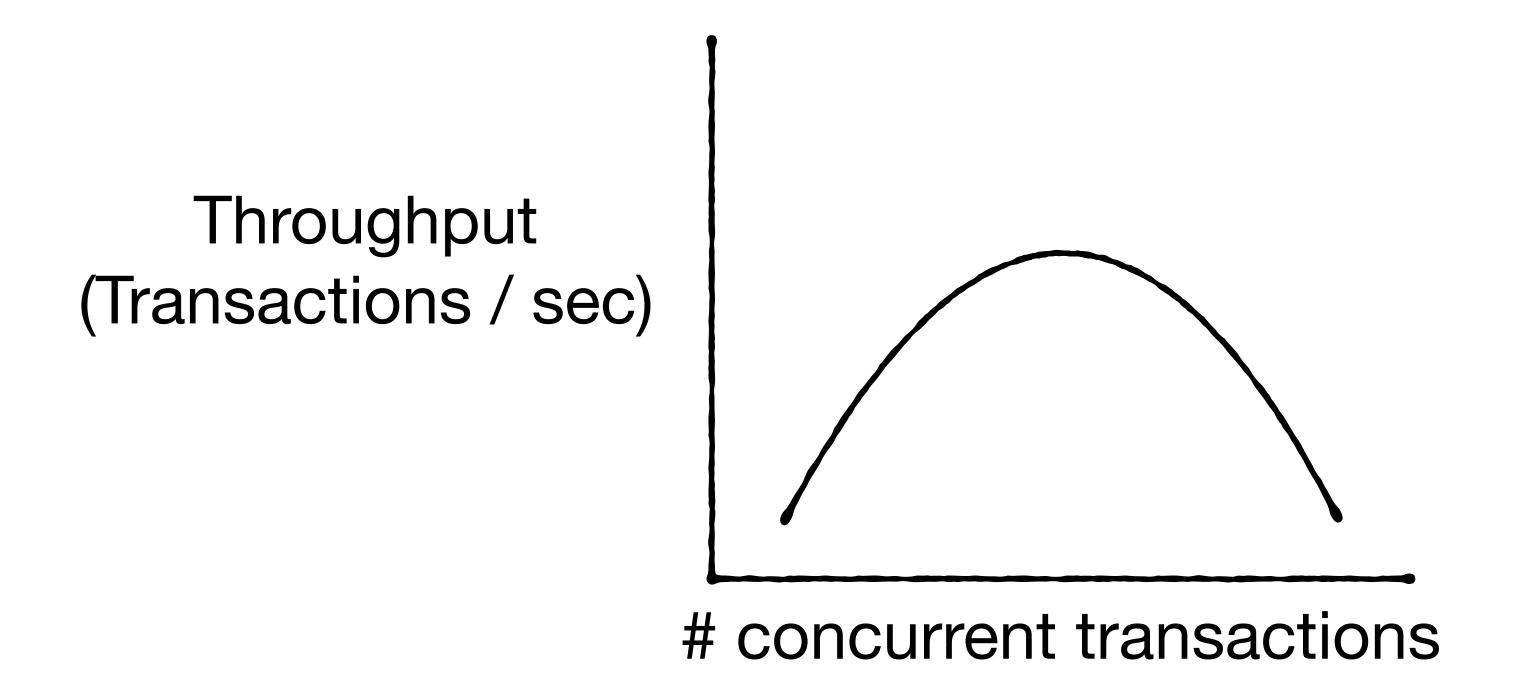
The default in many relational DB systems



Repeatable reads Range locks not employed (e.g., no tree/table locking)



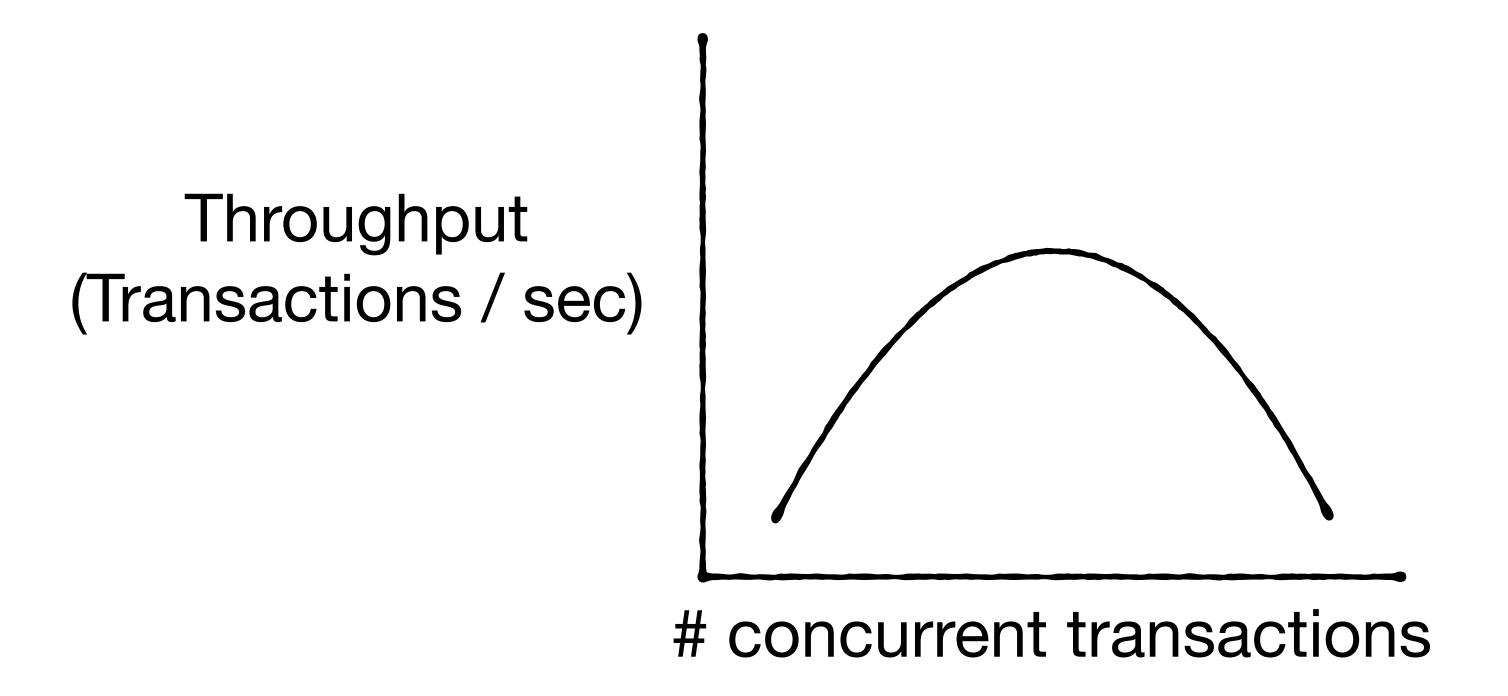




How can we address this?

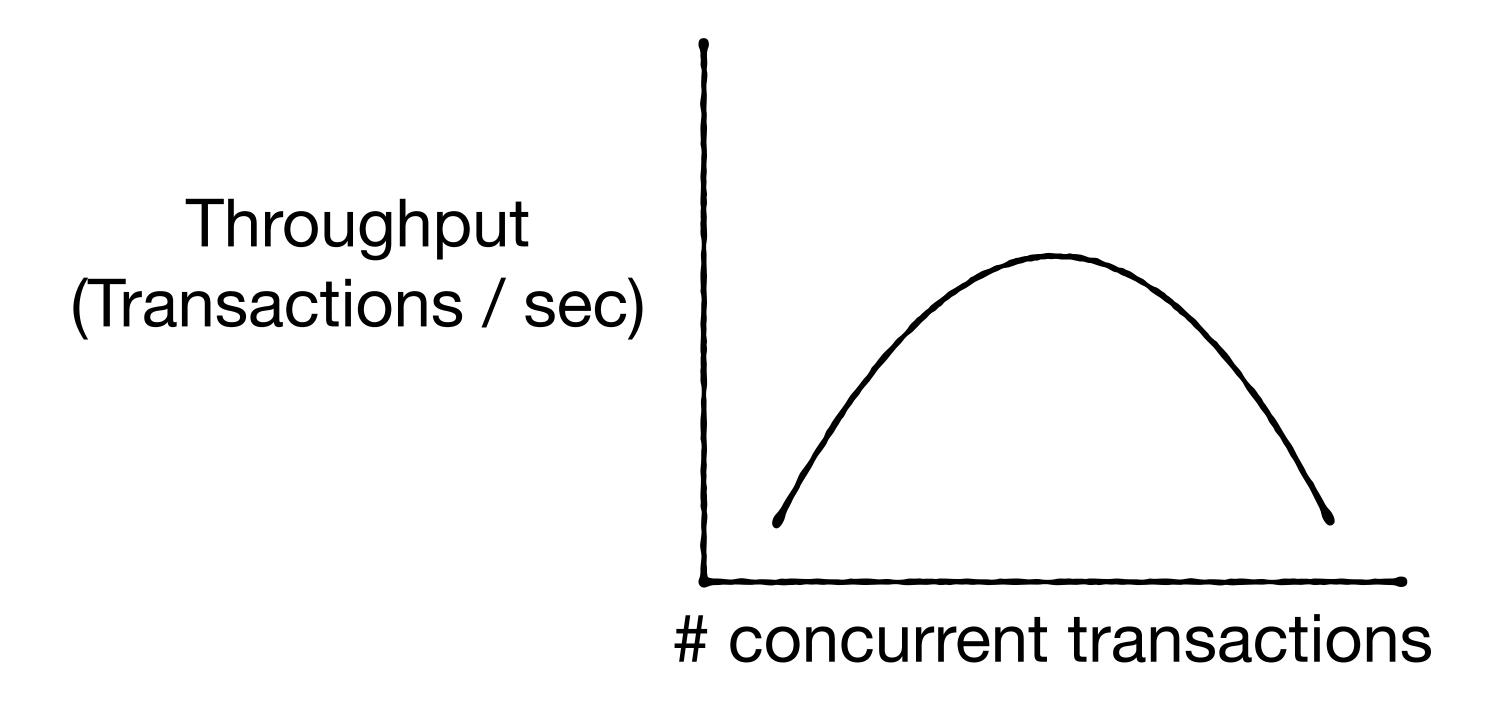
(1)

(2)



How can we address this?

- (1) Design the application such that transactions are short
- **(2)**

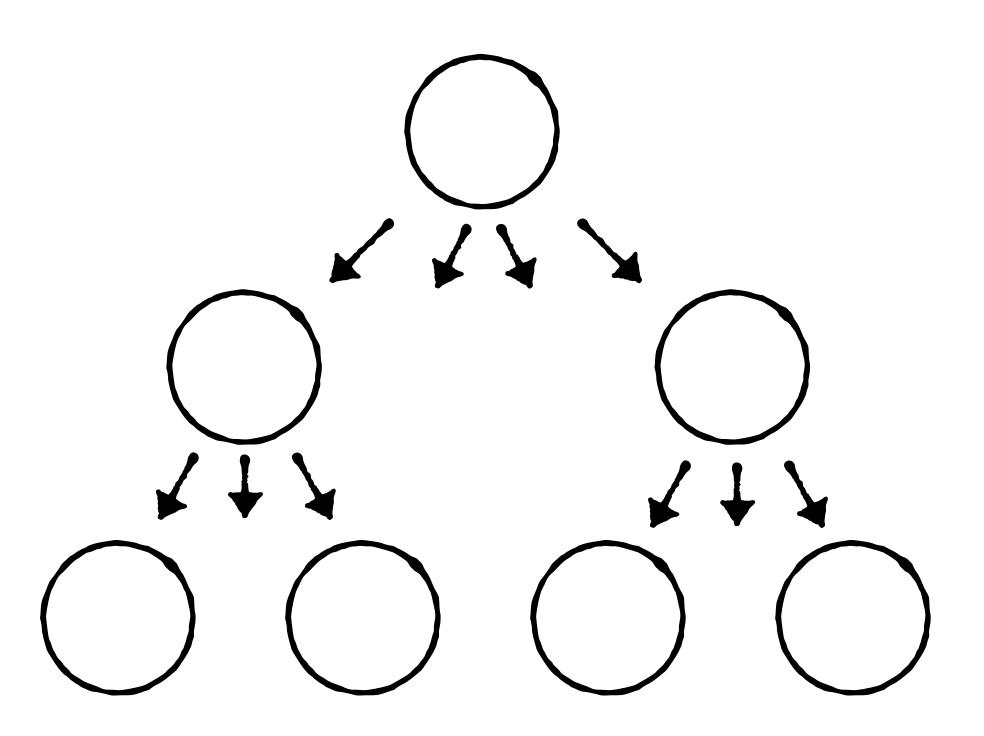


How can we address this?

- (1) Design the application such that transactions are short
- (2) Restrict the maximum number of concurrent transactions

With Strict Two-Phase Locking, a transaction locks all objects on its path until it commits

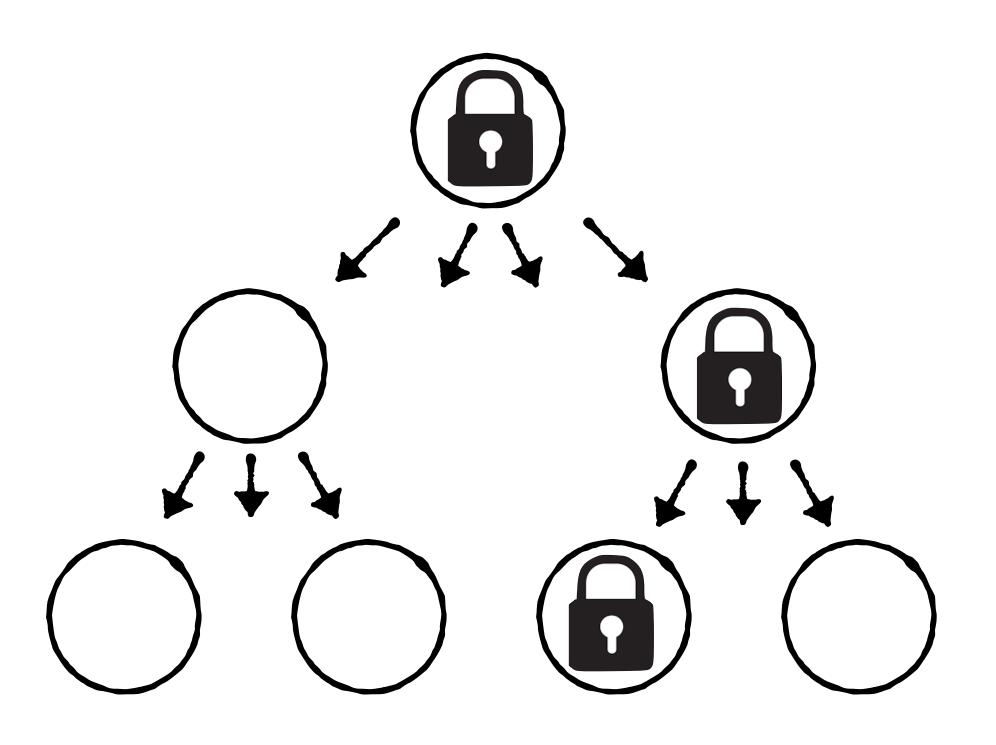
Why is this a problem for B-tree update?



With Strict Two-Phase Locking, a transaction locks all objects on its path until it commits

Why is this a problem for B-tree update?

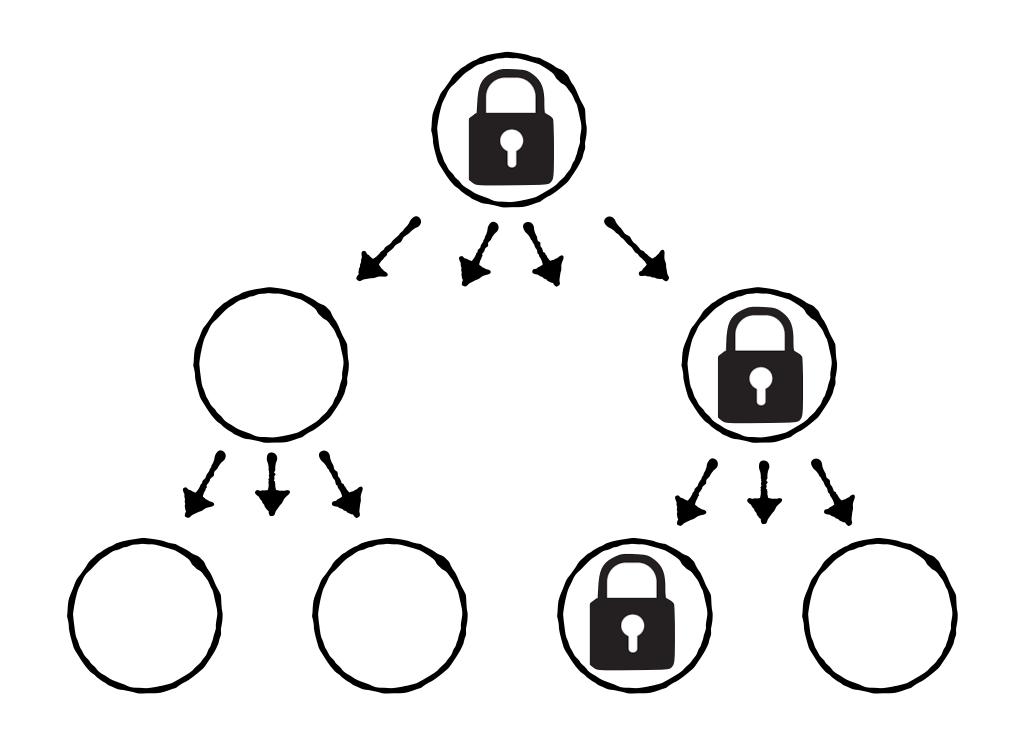
Lots of locking contention at root and upper levels



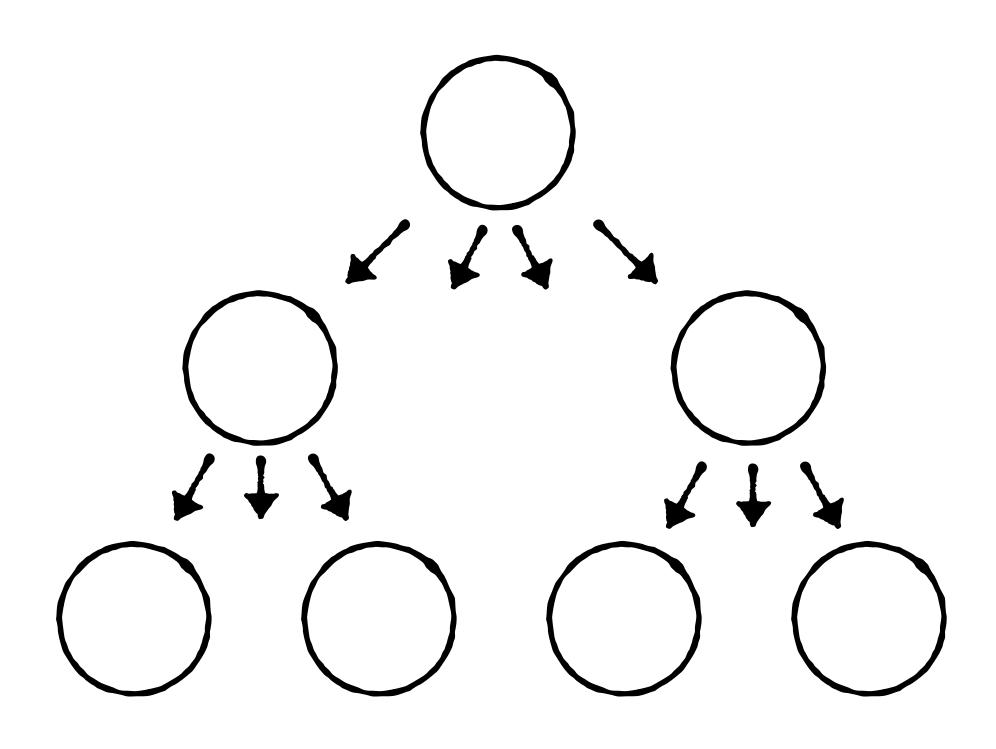
With Strict Two-Phase Locking, a transaction locks all objects on its path until it commits

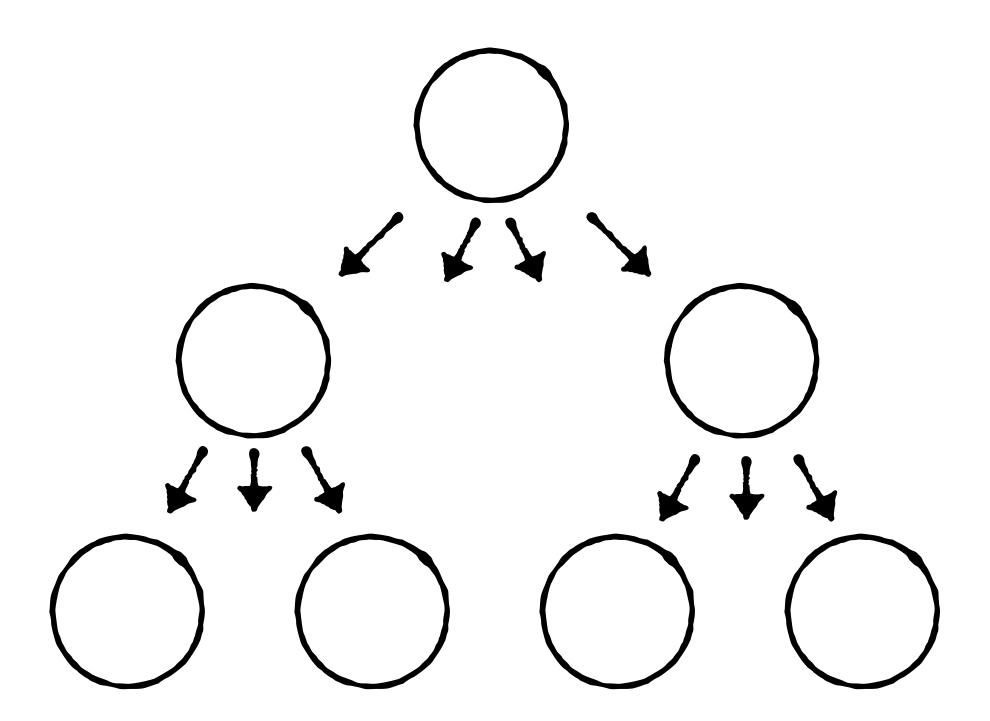
Why is this a problem for B-tree update?

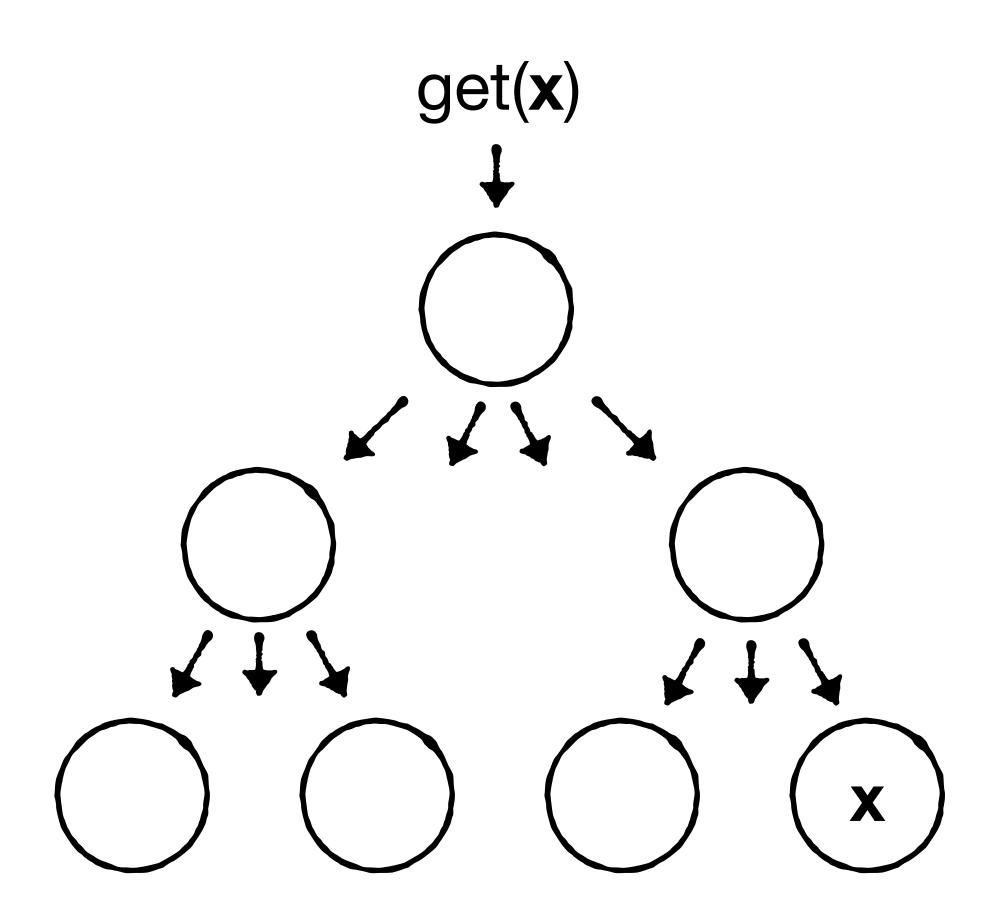
Lots of locking contention at root and upper levels

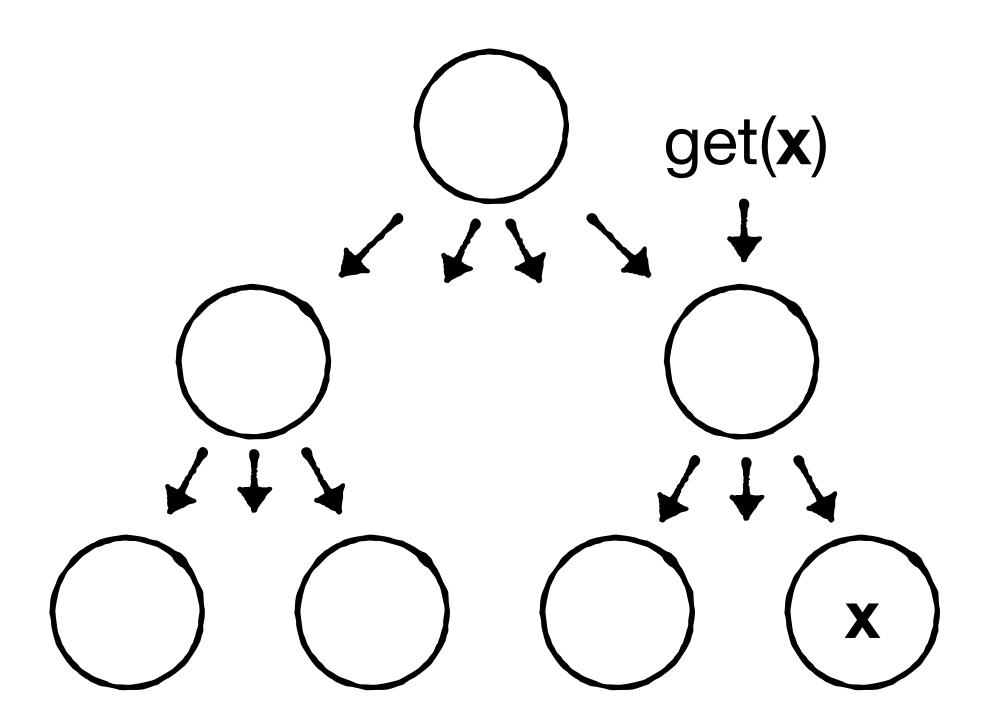


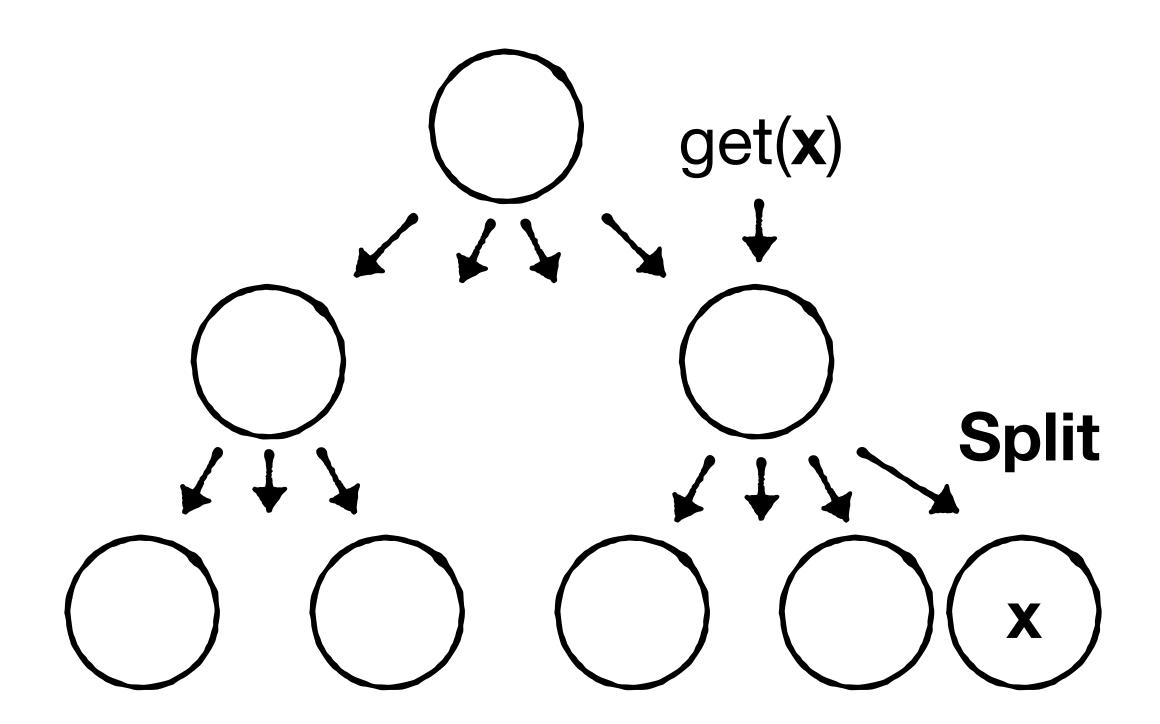
Any solutions?

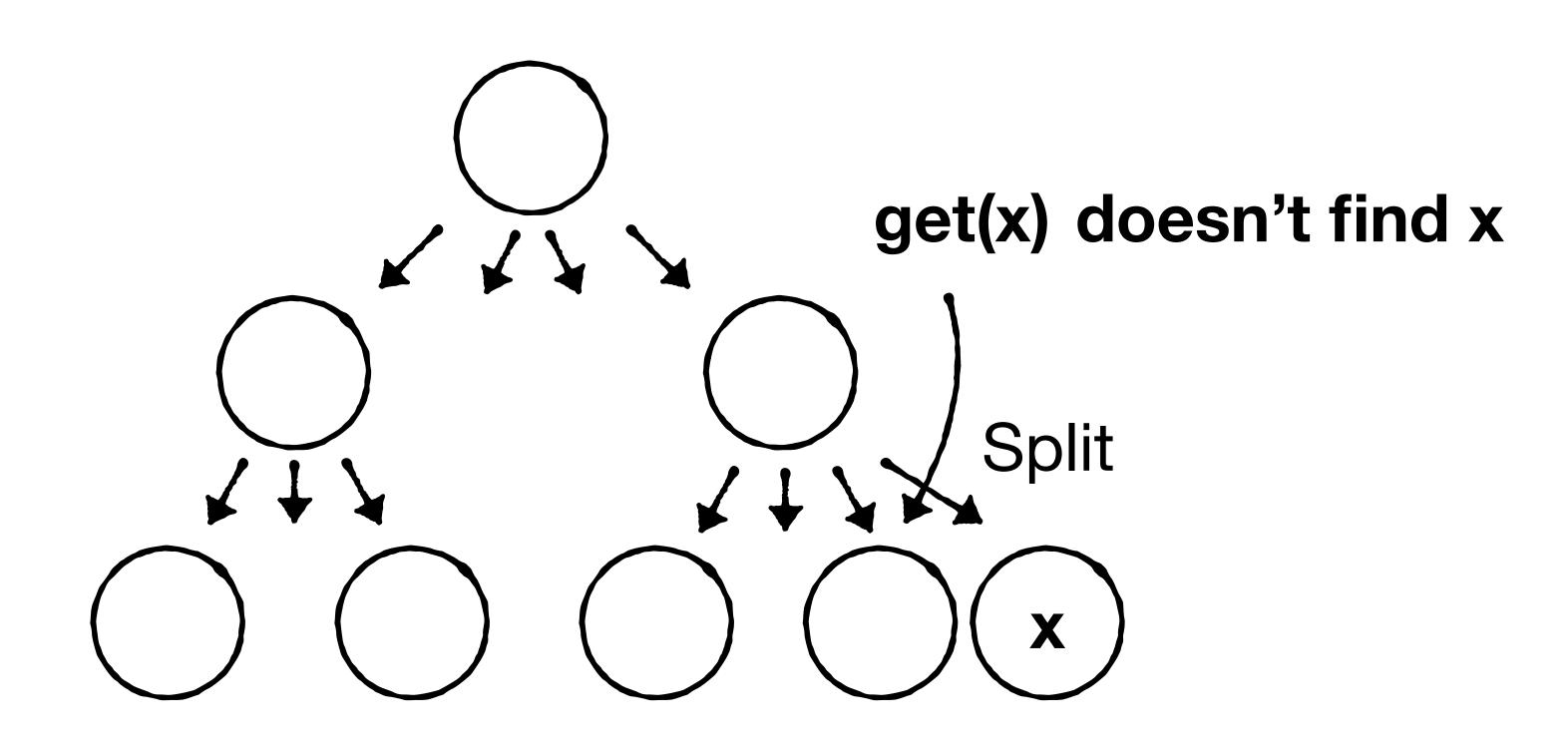


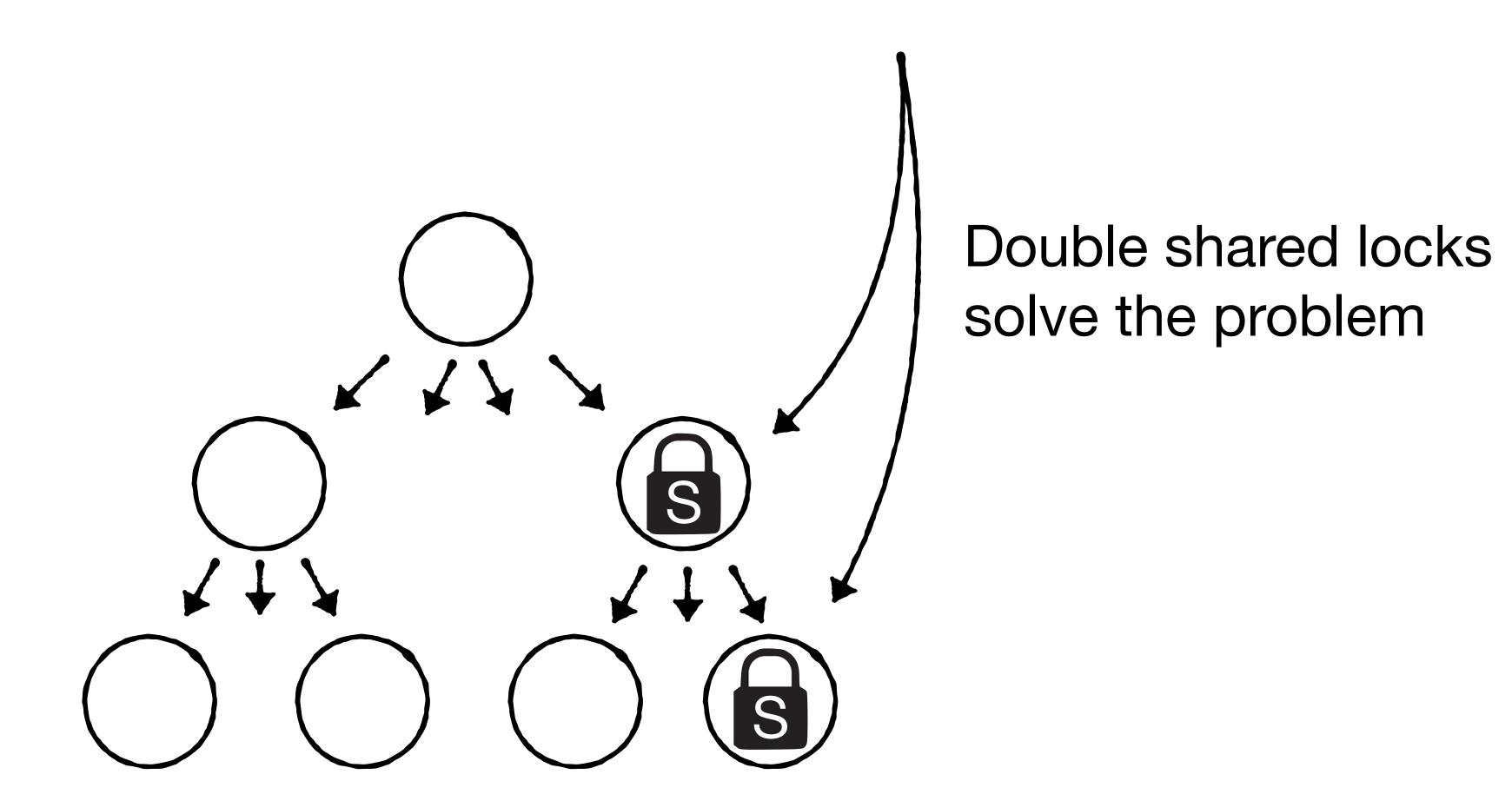


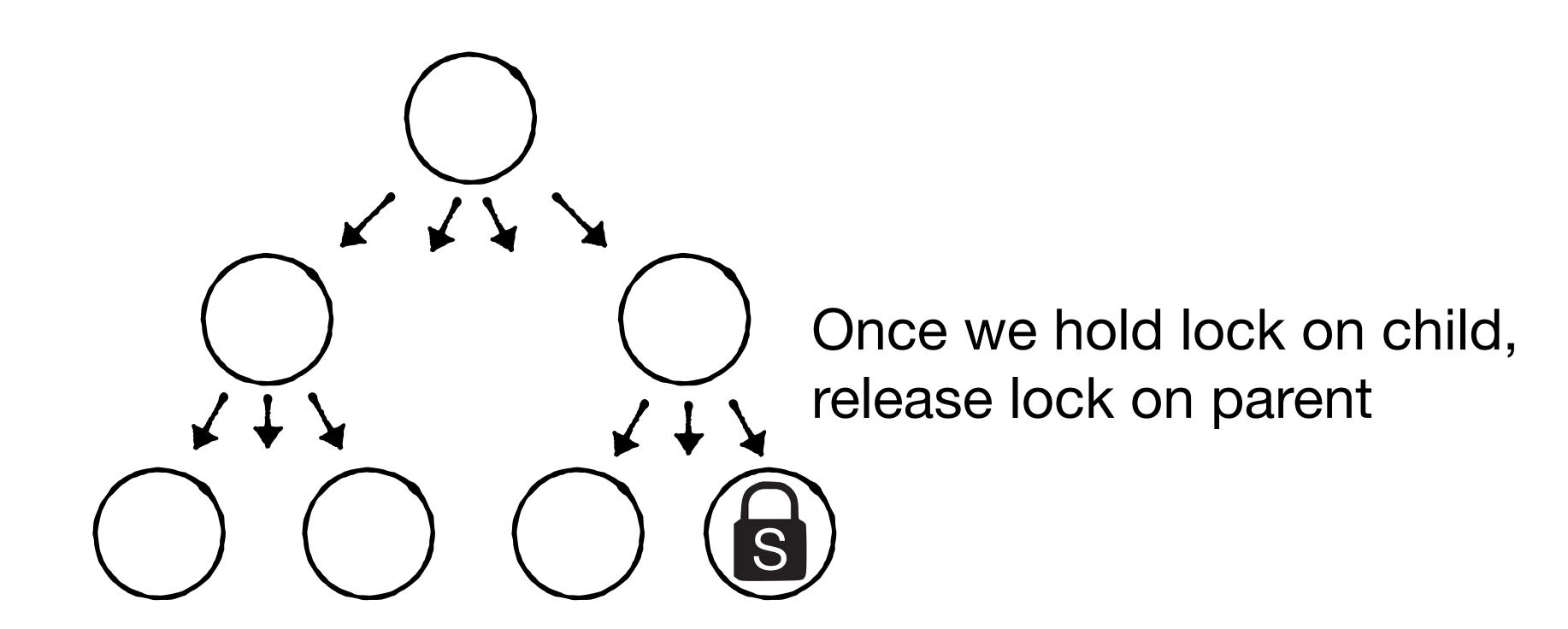


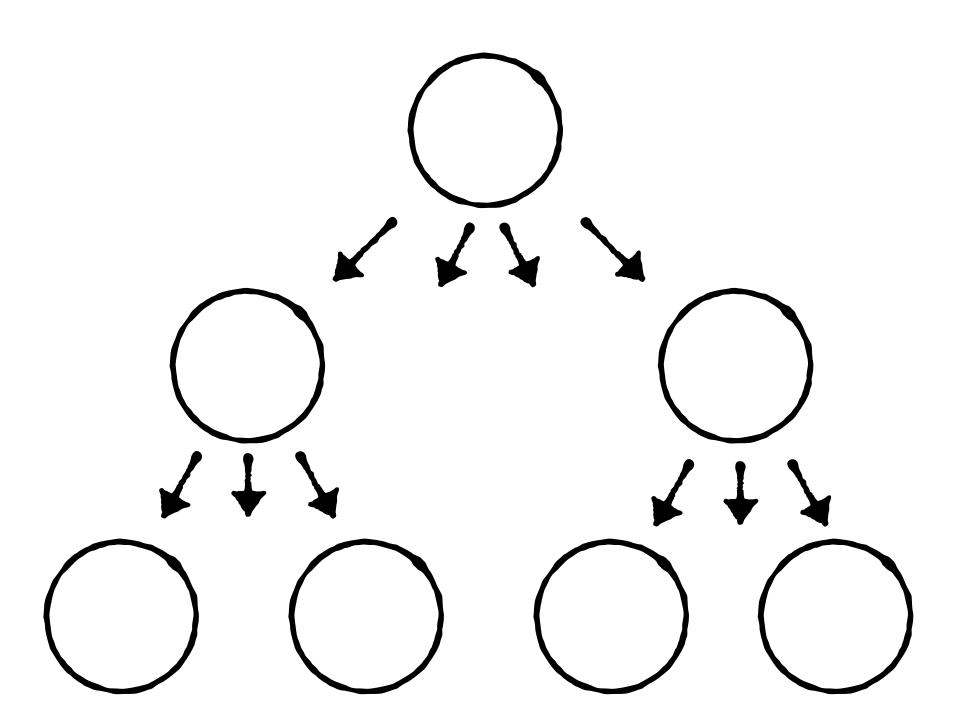


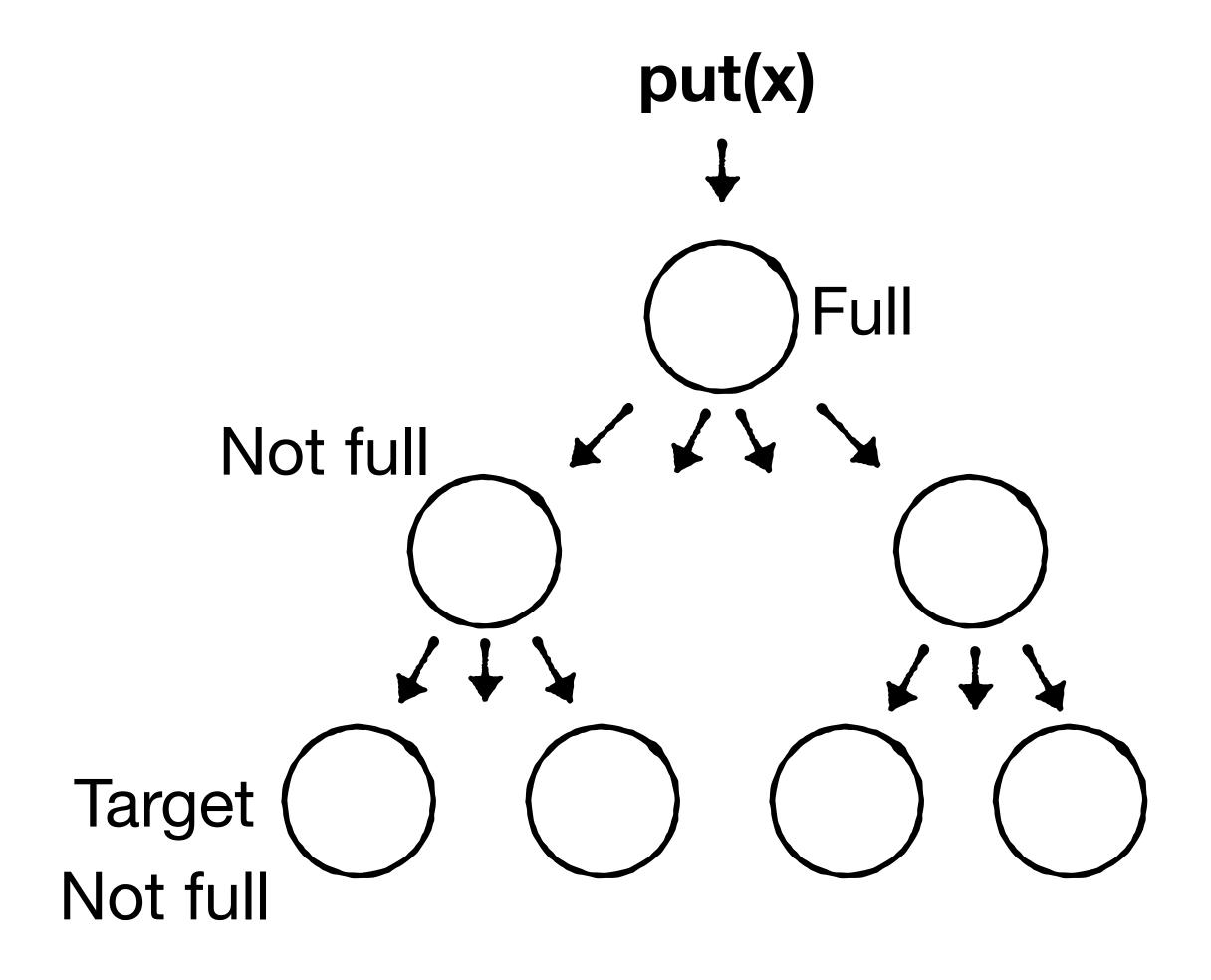


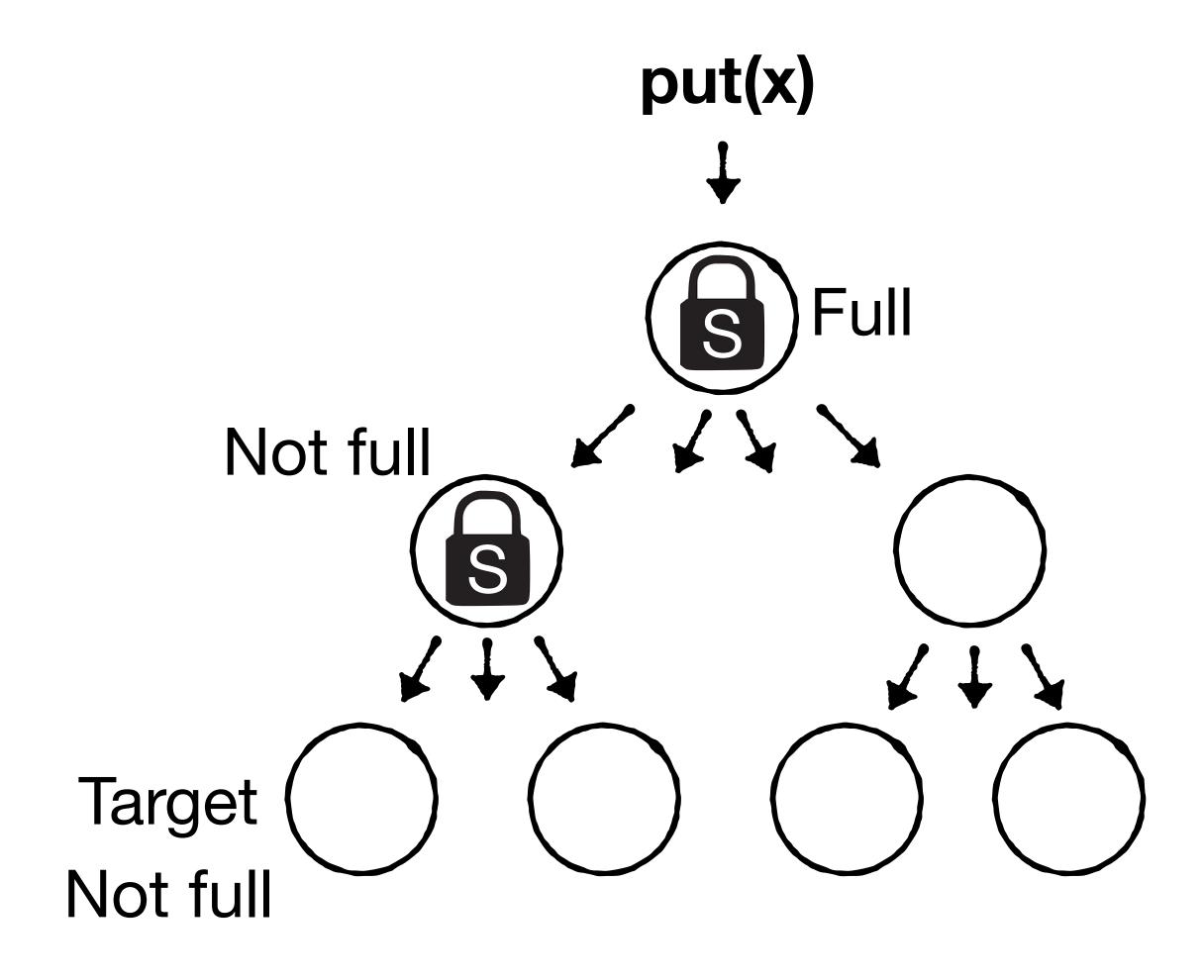


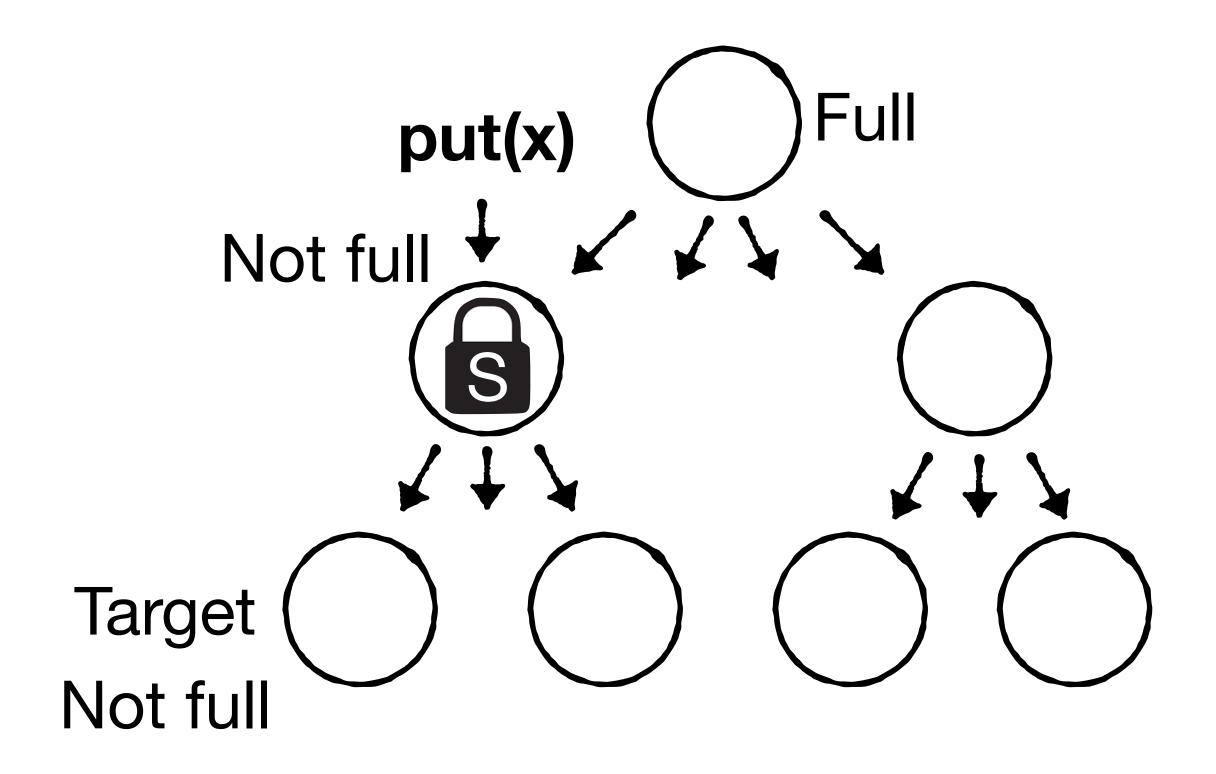


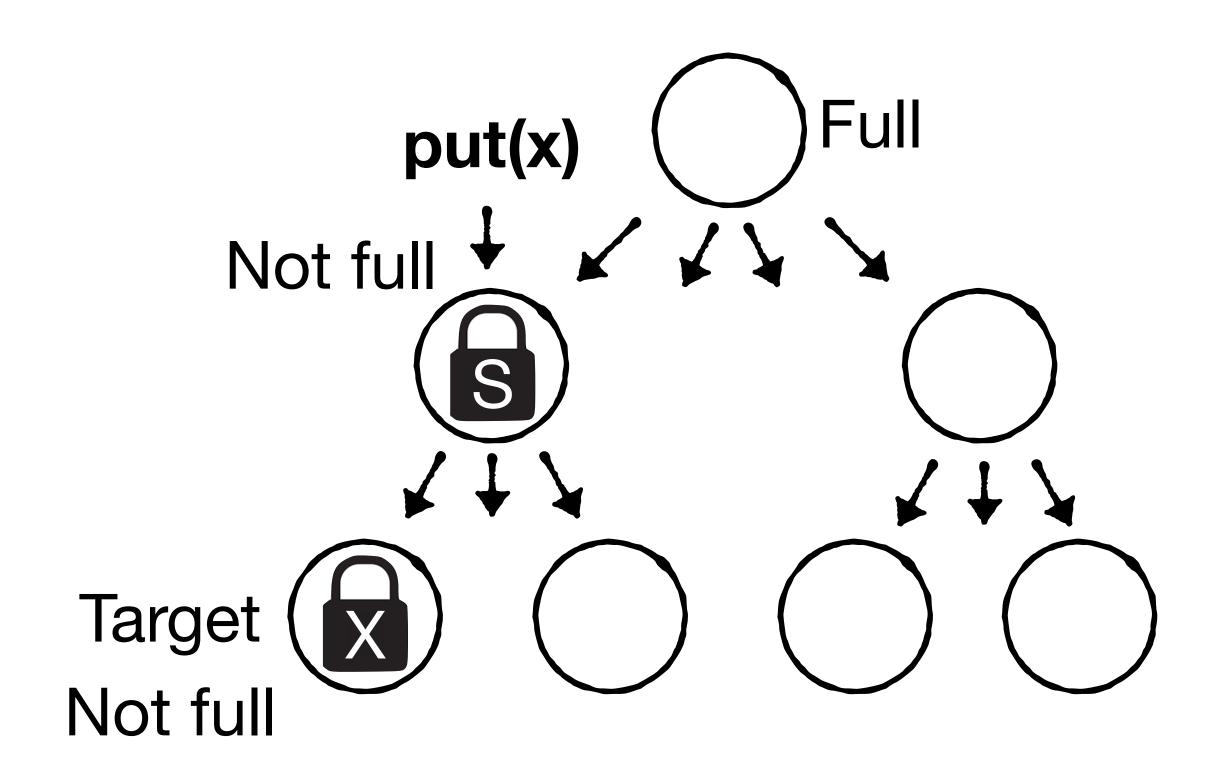


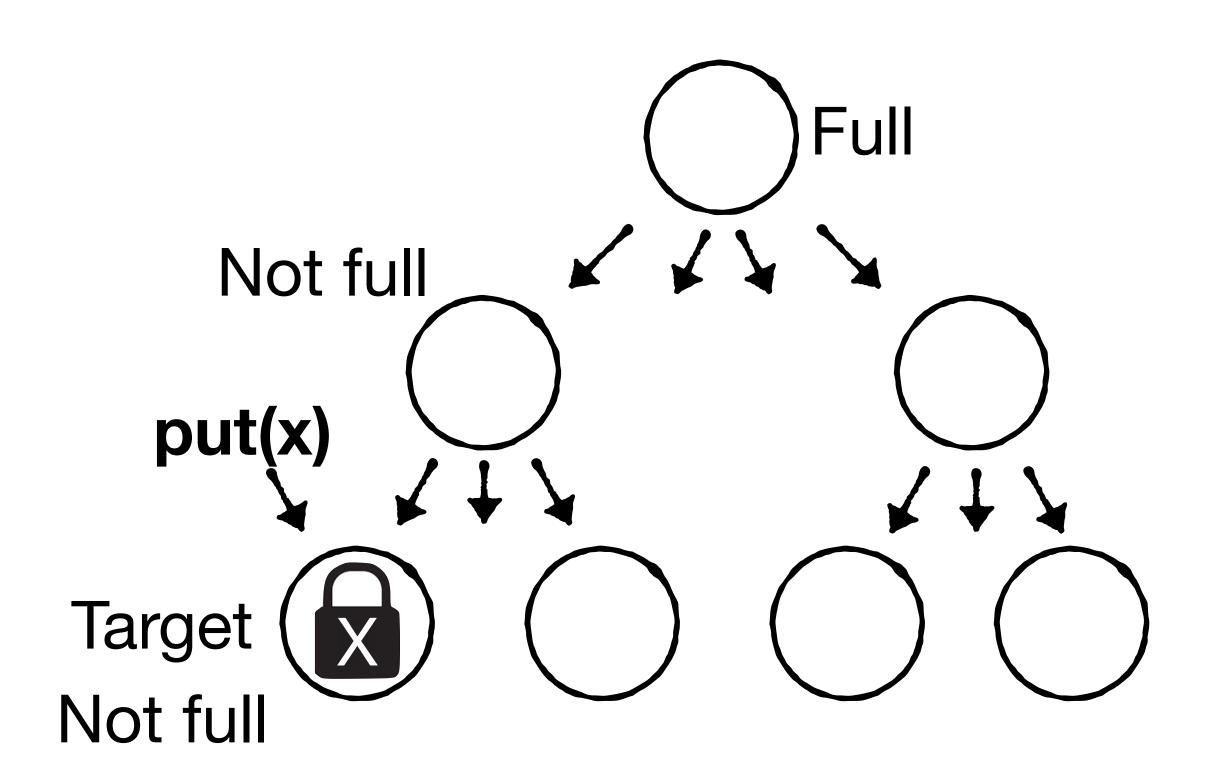


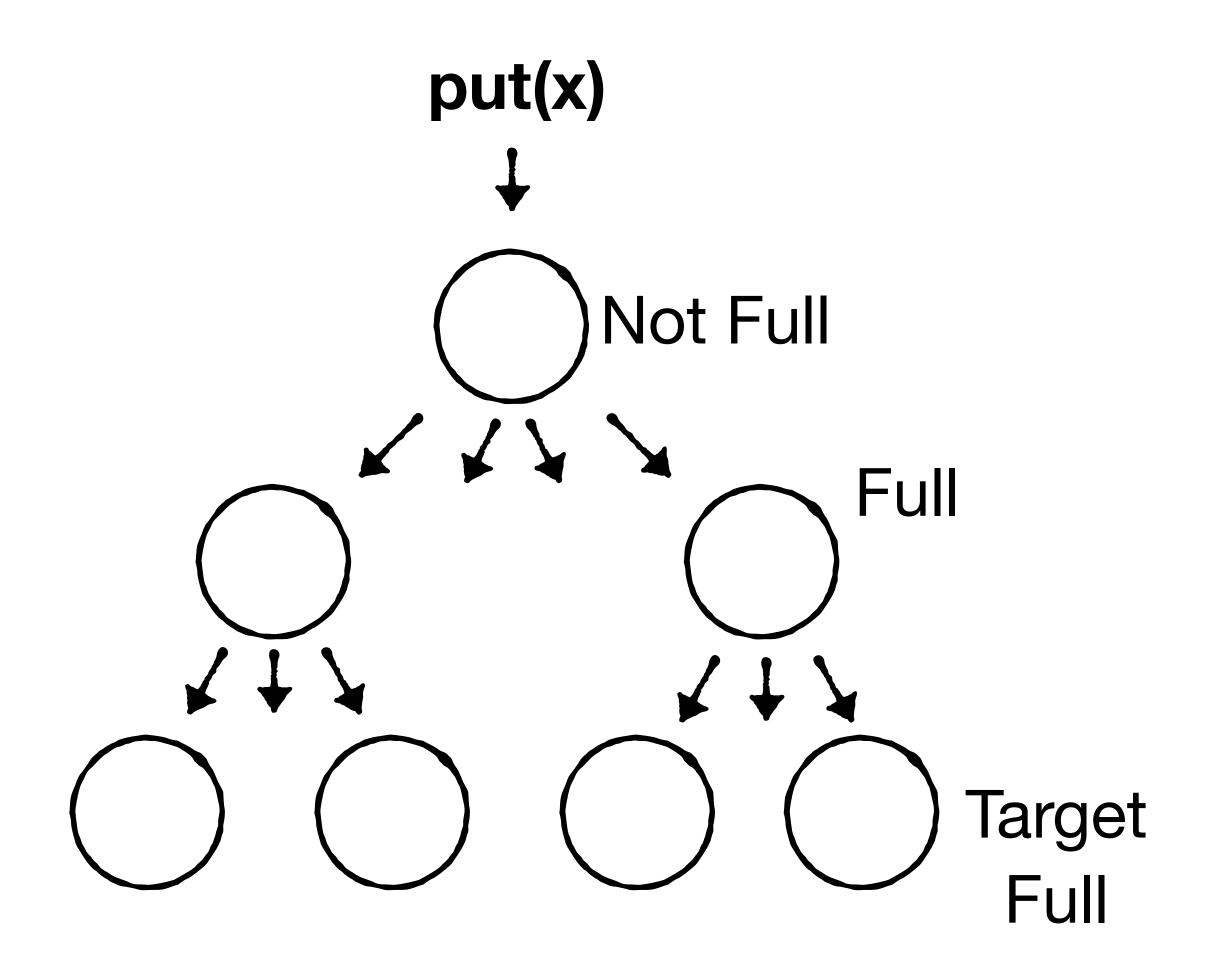


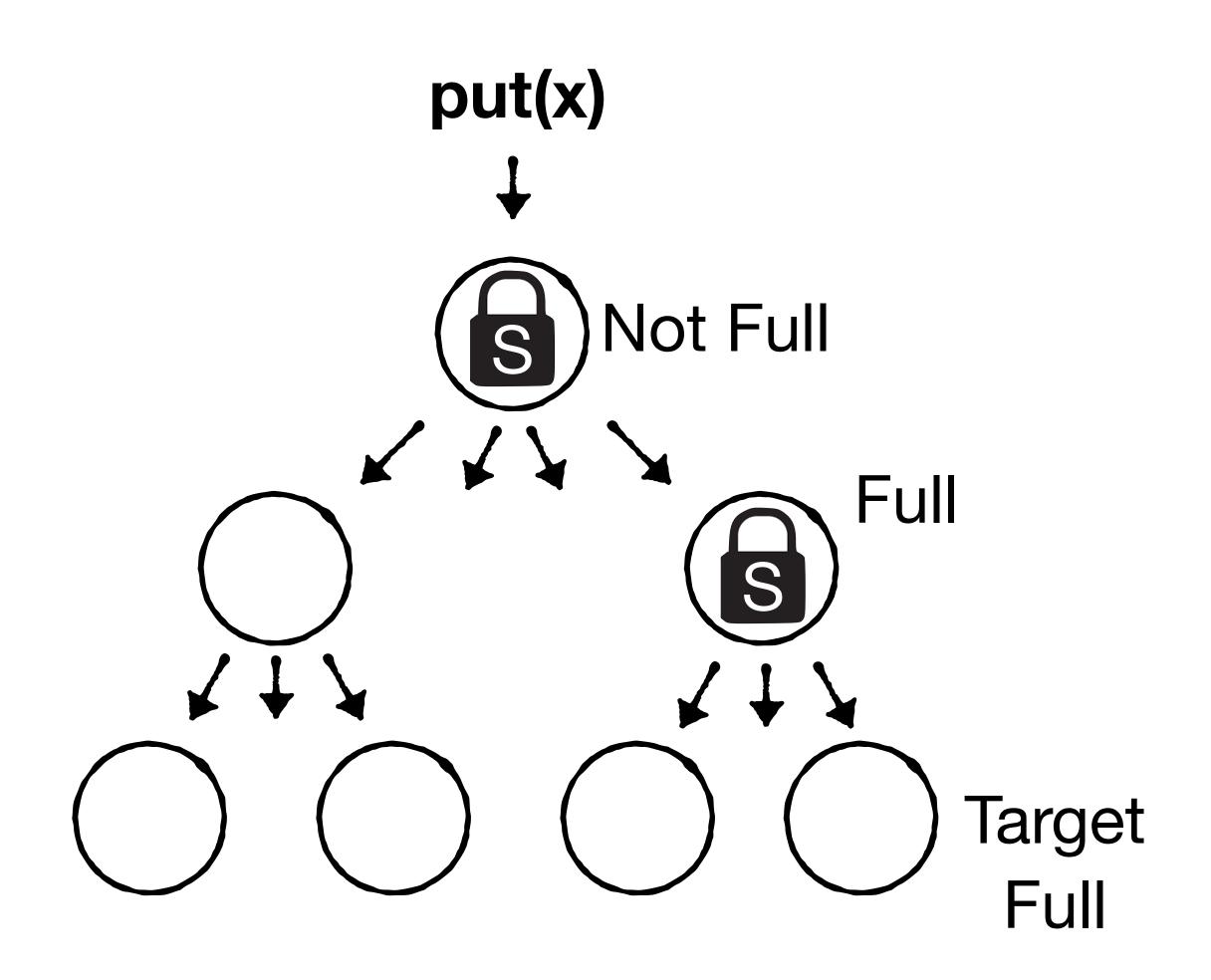


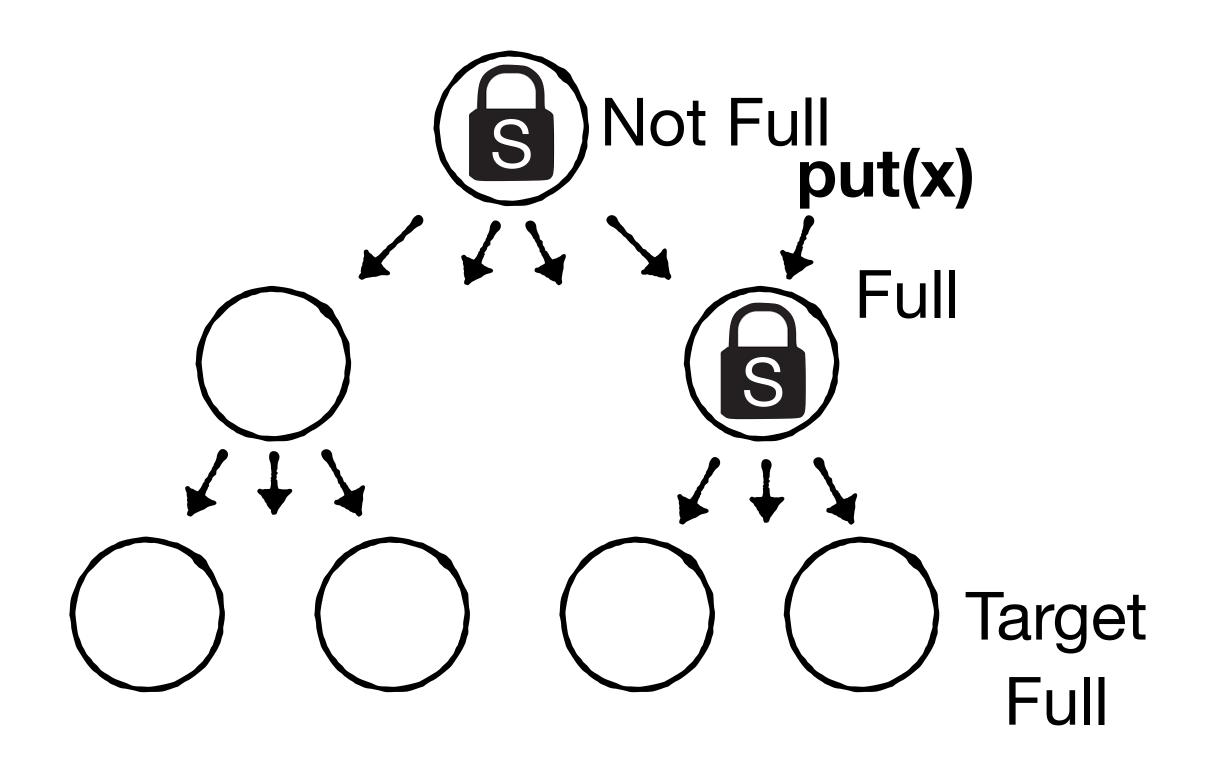


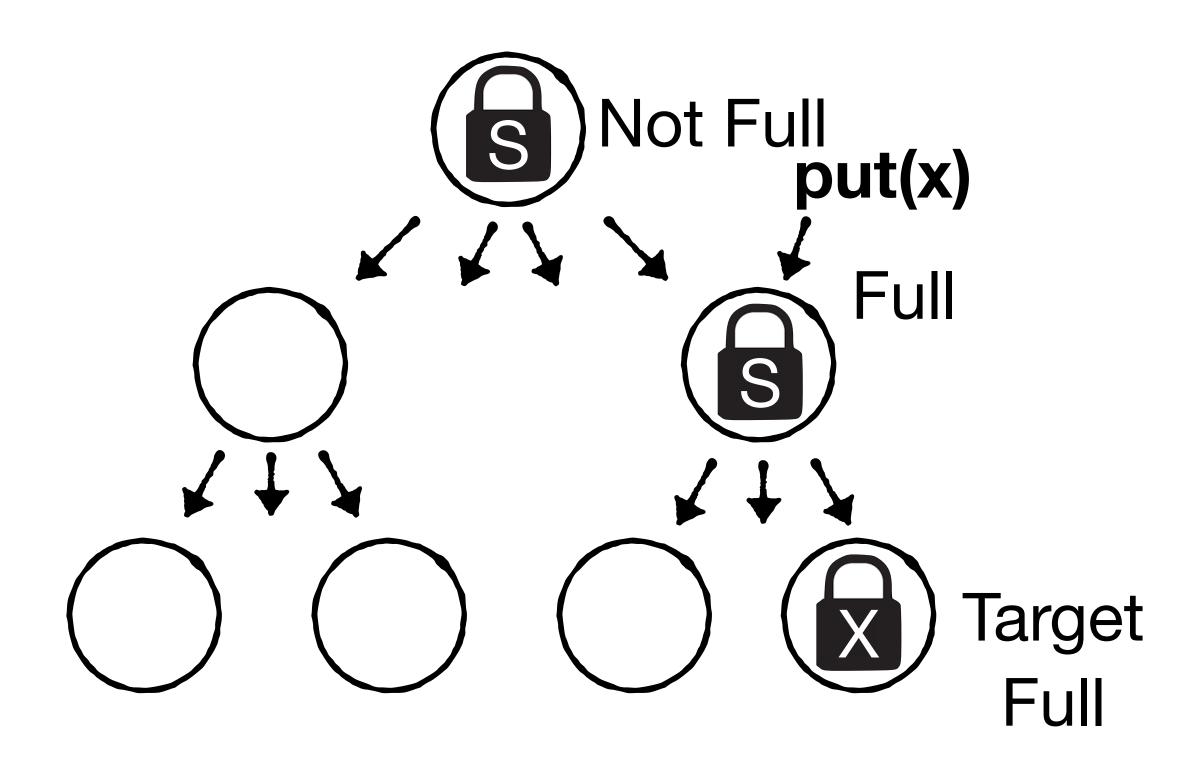


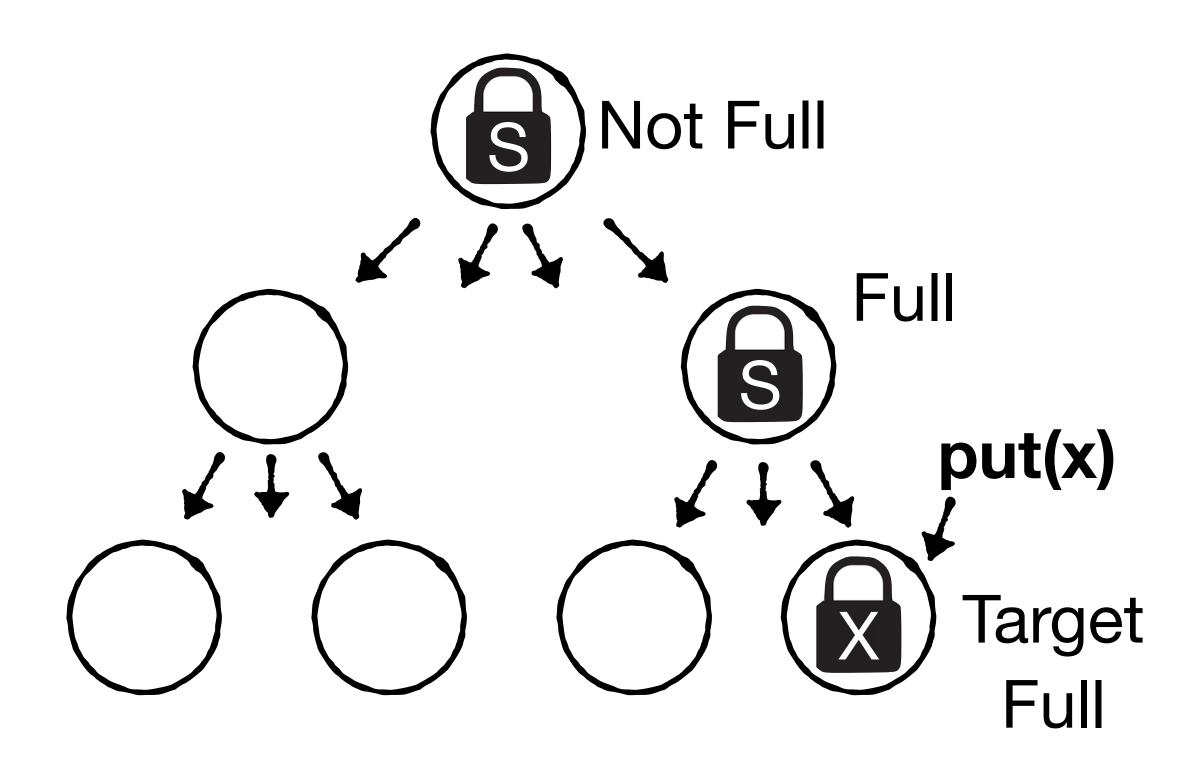


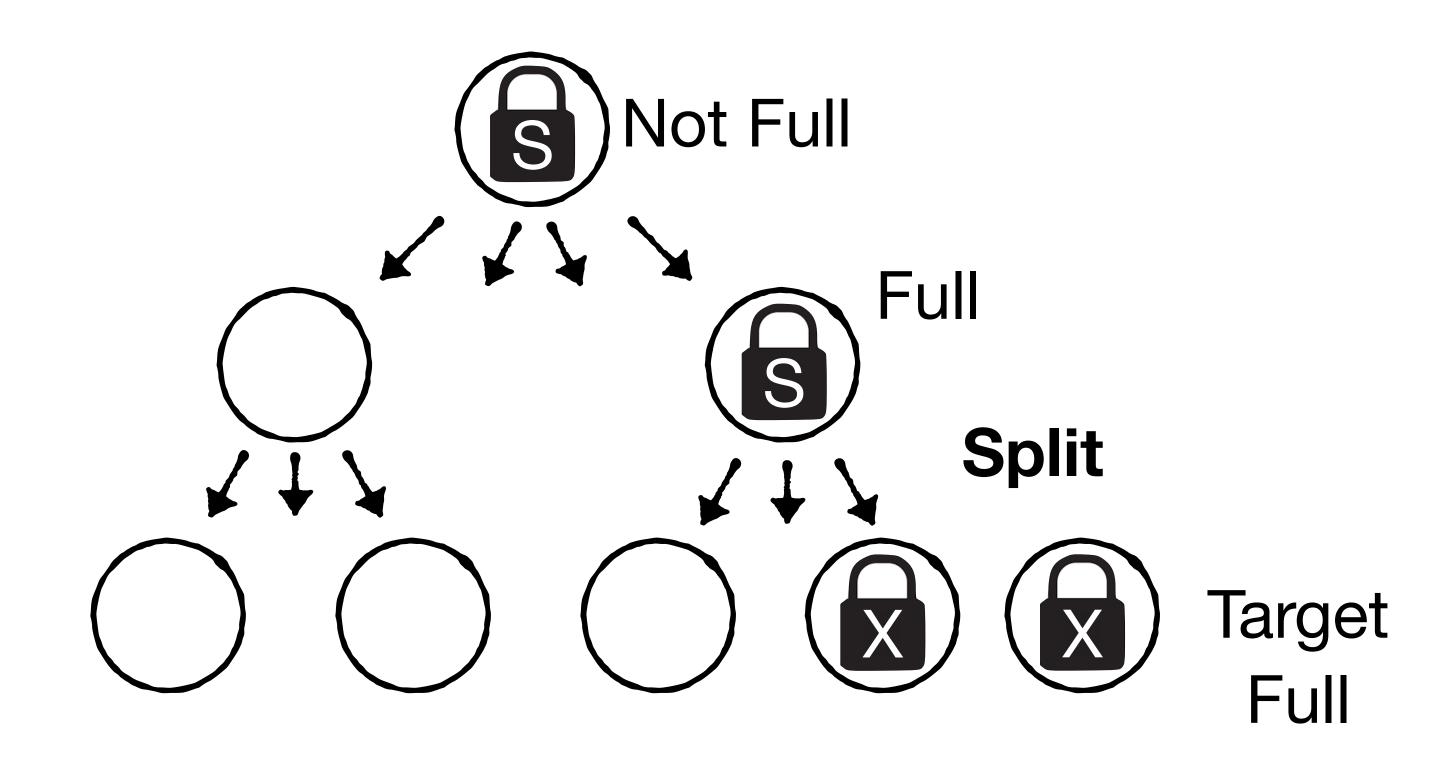


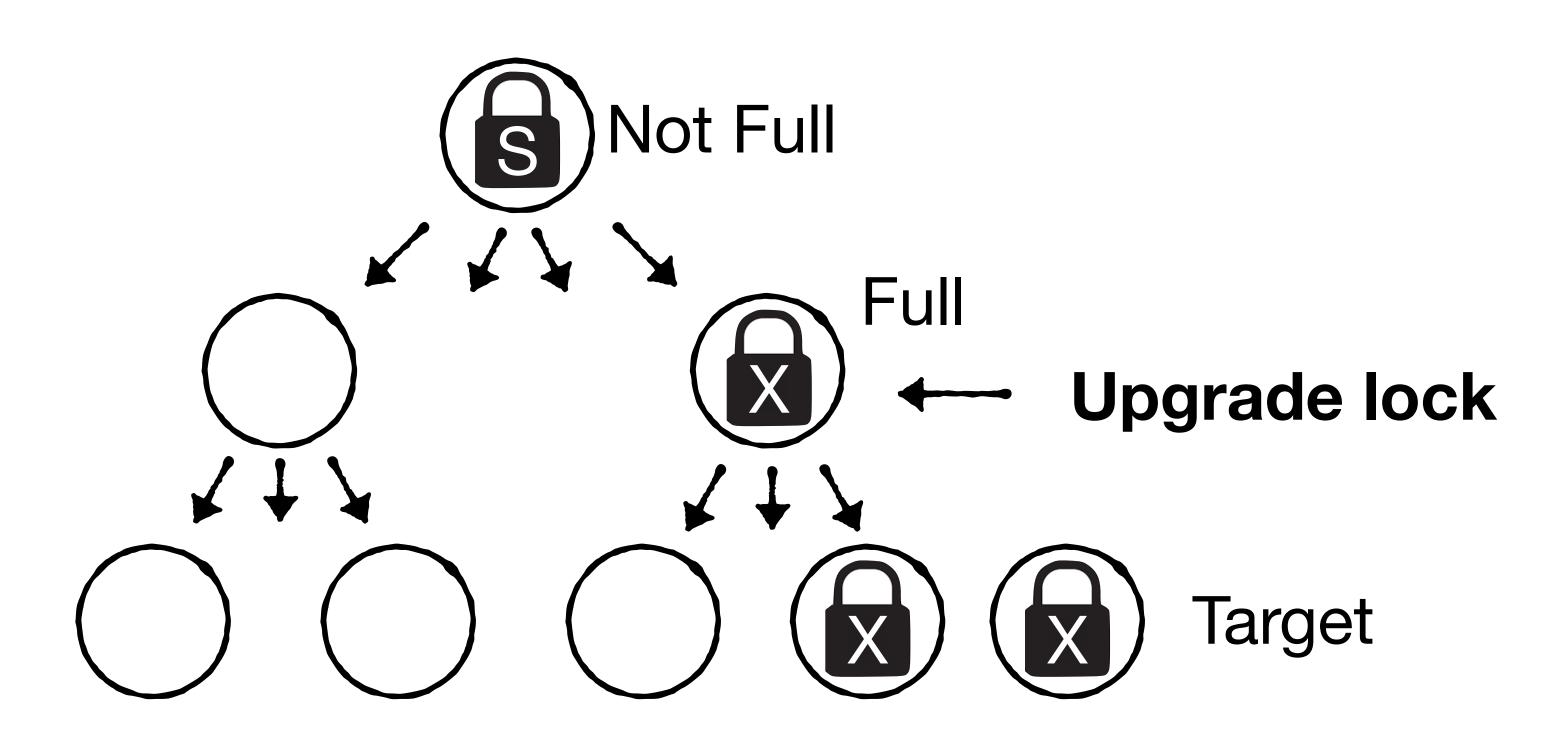


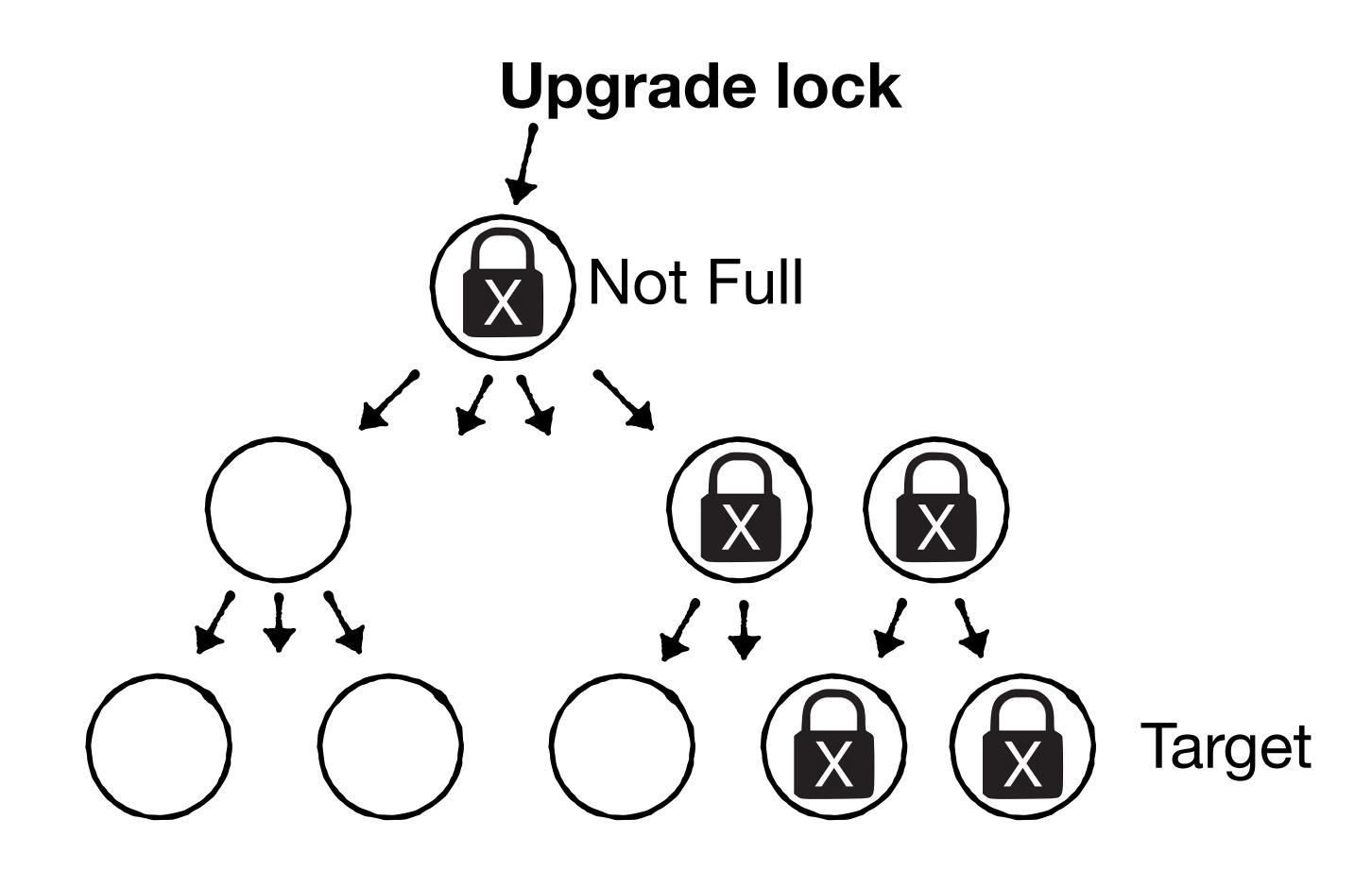


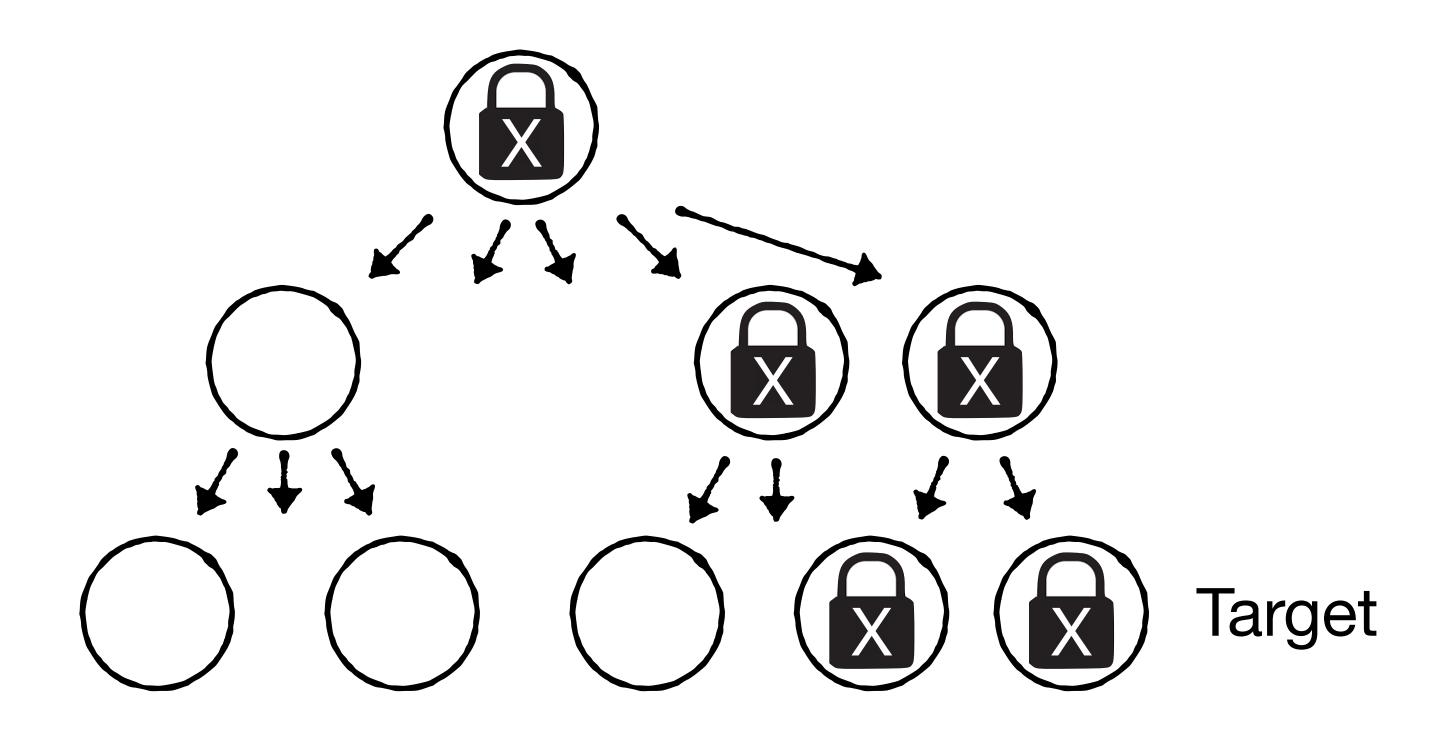




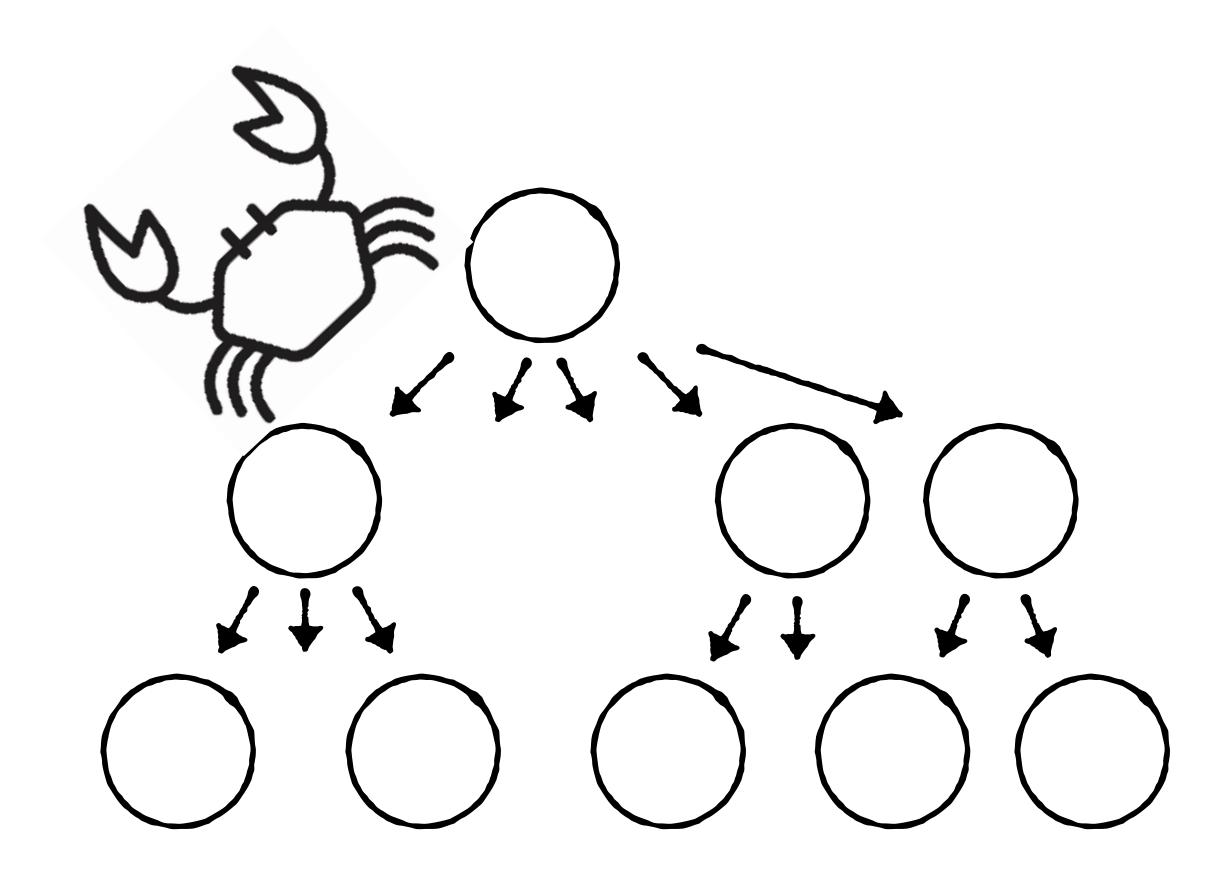




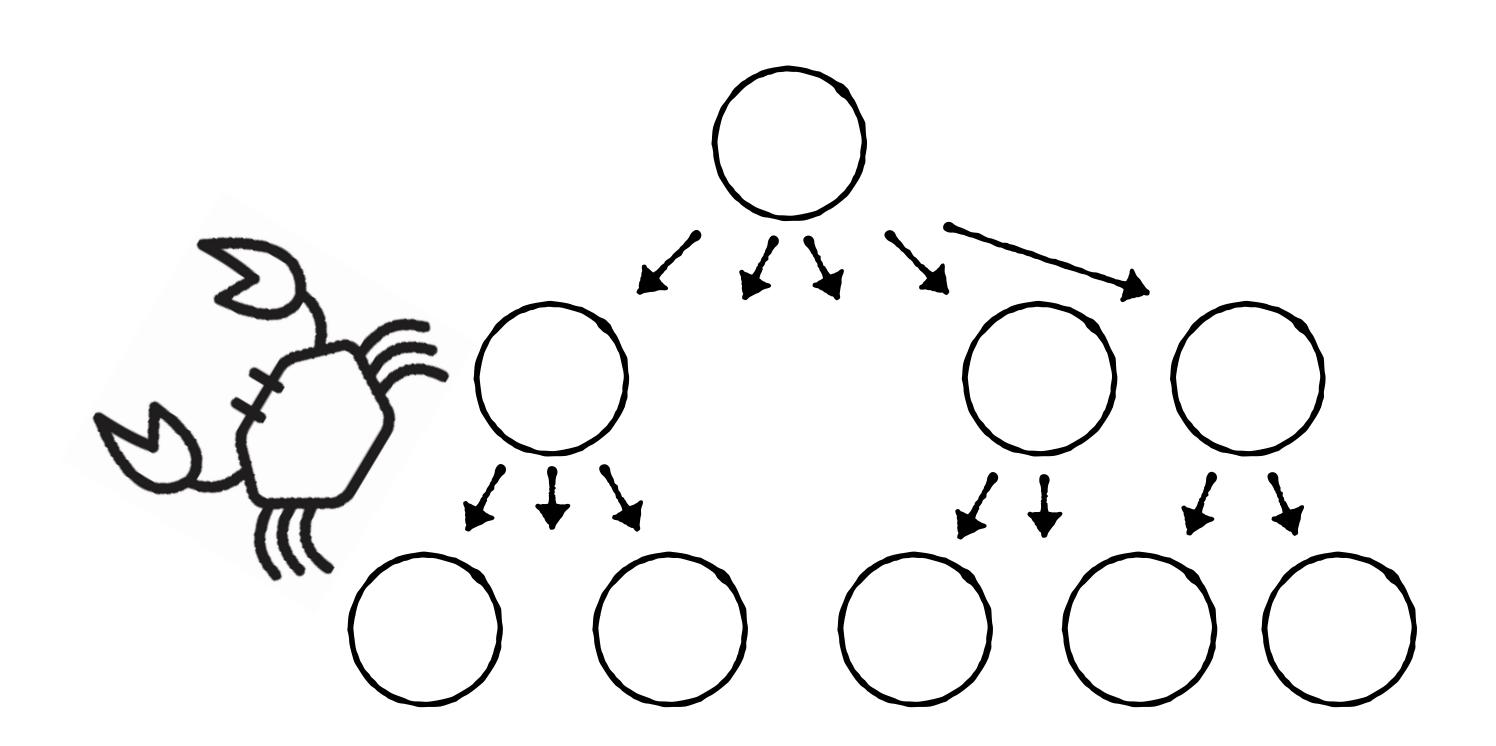




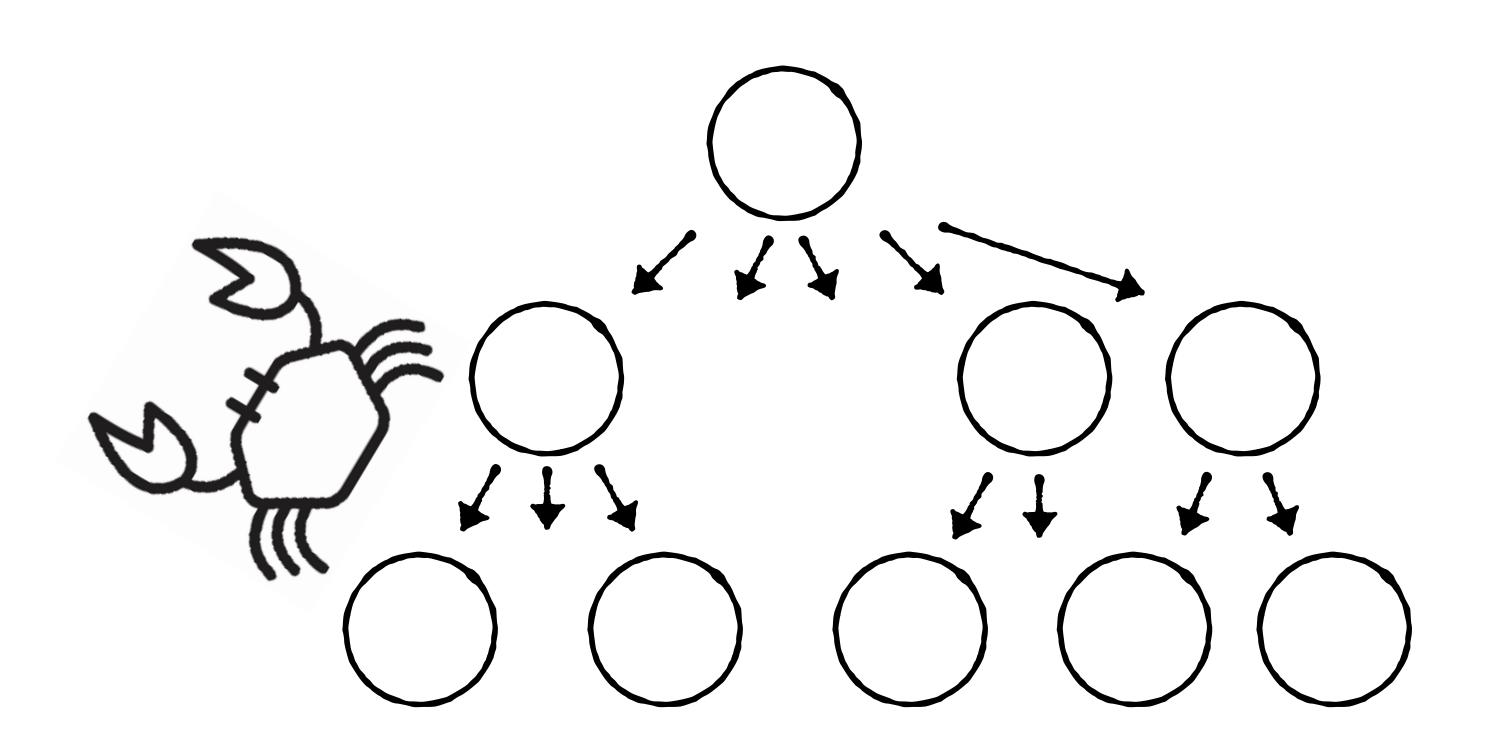
Protocol known as Lock-Coupling, or Crabbing



Protocol known as Lock-Coupling, or Crabbing



Shows that while strict two-phase locking is correct, it can sometimes be relaxed while maintaining correctness



And now, office hours