

Formalising the Foundation Modelling and Verifying a Database Backend

Joint Work With:

Marc Shapiro (LIP6 & INRIA, Paris)
Annette Bieniusa (RPTU, Kaiserslautern)
Carla Ferreira (NOVA LINCS, Lisbon)

With Support From:

Saalik Hatia (Scaleway, Paris)
Emilie Ma (UBC, Canada)
Jaurel Fosset (Polytech, Dijon)
Gustavo Petri (Amazon Research)

Outline

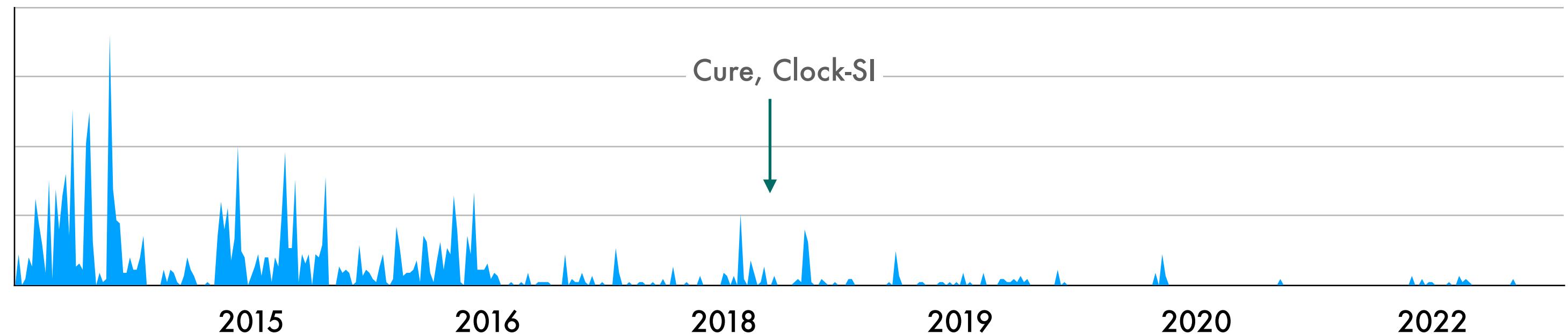
1. Background of this work and Motivation
2. Introduction
 1. Rich Datatypes
 2. Transactions
 3. Traces and their correctness (TCC)
 4. Database Backend
3. Stores
 1. Semantics
 2. Journal
 3. Map
4. Transaction Semantics
5. Correctness of transaction semantics
6. Improved definition of TCC
7. Proof Outline for correctness of the system specification
8. Specification to Implementation
9. Future

Learnings from AntidoteDB

Began ca. 2015

Clock-SI protocol in 2018

Bug fixing since then



What Happened?

Databases are large, complex

AntidoteDB: Planet-scale, highly available, replicated db

- ≈ 6 PhD Students => CRDTs, Clock-SI, Sharding, Lazy Replication, Lazy CC
- Heavily Engineered, Optimised. 60 KLoC ([Erlang](#))
- Too much code + tired PhD students = Unknown Bugs = Data Loss!

Redis ≈ 100 KLoc (C)

RocksDB ≈ 300 KLoC (C++)

Overarching Research Theme

We know, formal methods help avoid bugs in systems.

Can we specify and verify a set of building blocks for a database which, when wrapped and composed correctly, can be used to build a correct database?

Aside: Rich Data Types

Conflict Free Replicated Datatypes (CRDT)

Concurrent updates without synchronisation + Guaranteed Convergence

Key Properties:

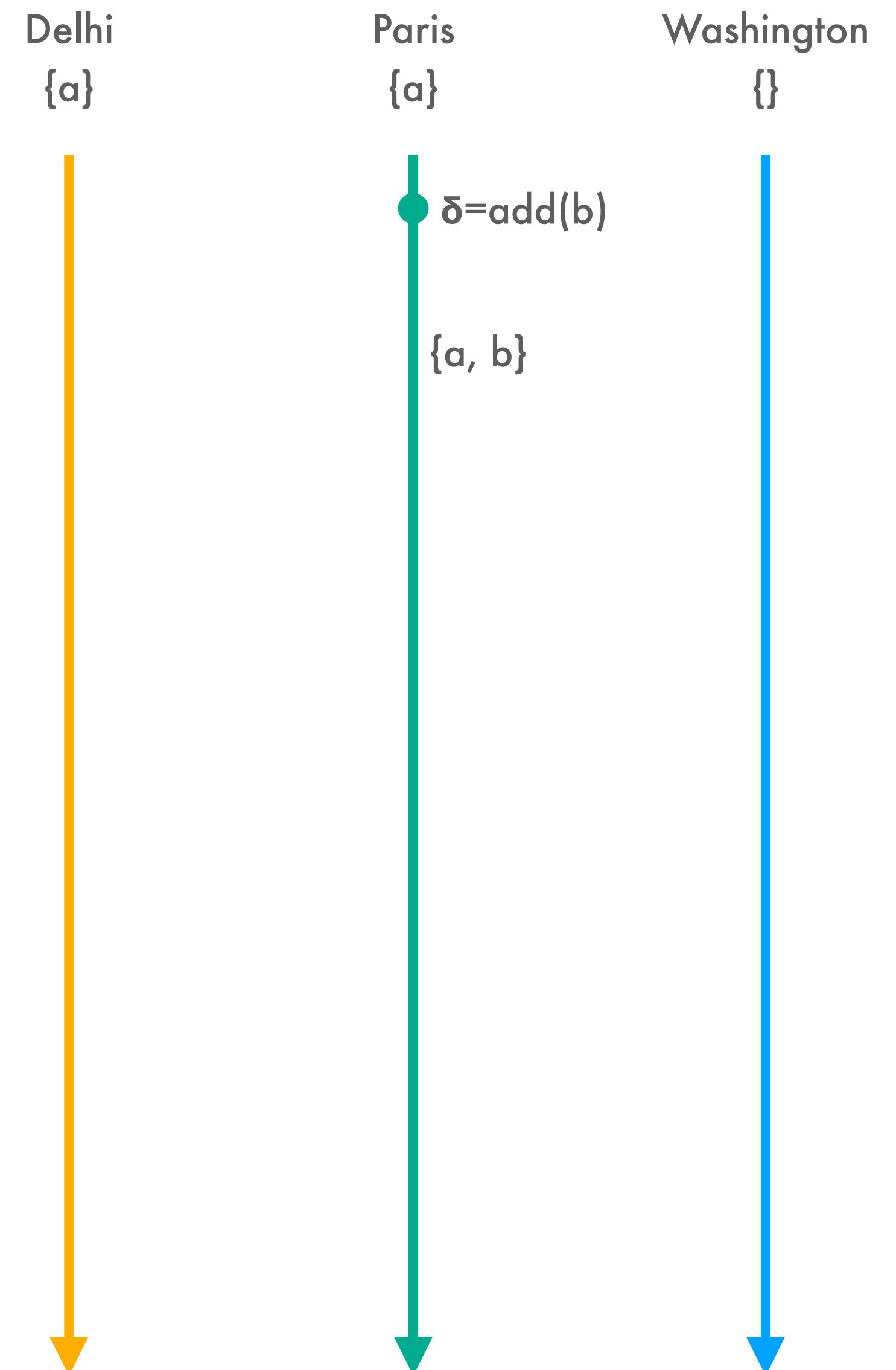
- Commutative, Associative and Idempotent Operations
- Mergeable States
- No Locks or consensus for updates

Implementation:

- **Effect (δ):** Takes the state of the CRDT and produces a new state
- **Merge:** Takes two CRDT states and **deterministically** produces a new state

Example: Add-Wins Set

- Effects: Add, Remove
- Merge: Between concurrent add and remove, add wins.



Aside: Rich Data Types

Conflict Free Replicated Datatypes (CRDT)

Concurrent updates without synchronisation + Guaranteed Convergence

Key Properties:

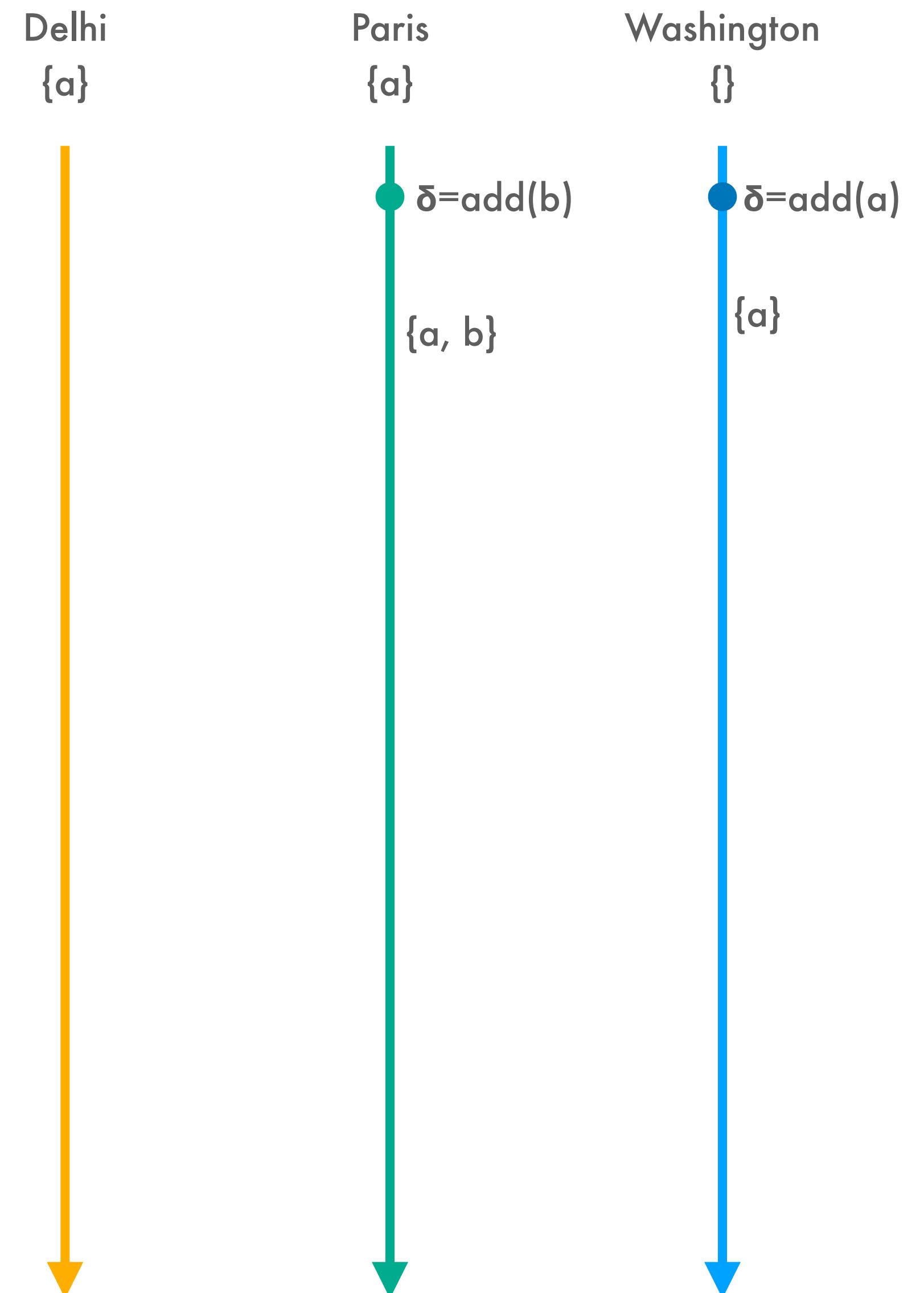
- Commutative, Associative and Idempotent Operations
- Mergeable States
- No Locks or consensus for updates

Implementation:

- **Effect (δ):** Takes the state of the CRDT and produces a new state
- **Merge:** Takes two CRDT states and **deterministically** produces a new state

Example: Add-Wins Set

- Effects: Add, Remove
- Merge: Between concurrent add and remove, add wins.



Aside: Rich Data Types

Conflict Free Replicated Datatypes (CRDT)

Concurrent updates without synchronisation + Guaranteed Convergence

Key Properties:

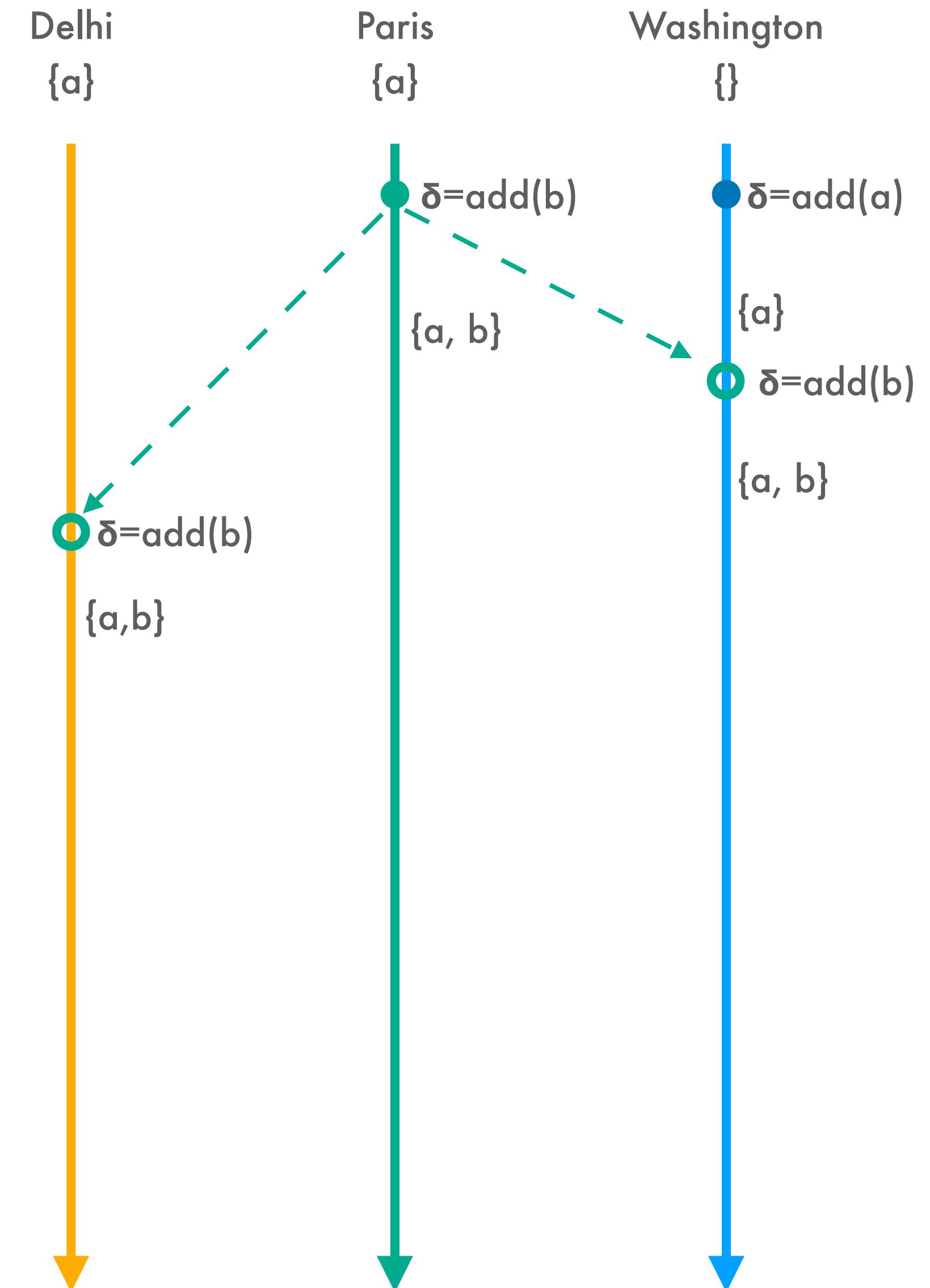
- Commutative, Associative and Idempotent Operations
- Mergeable States
- No Locks or consensus for updates

Implementation:

- **Effect (δ):** Takes the state of the CRDT and produces a new state
- **Merge:** Takes two CRDT states and **deterministically** produces a new state

Example: Add-Wins Set

- Effects: Add, Remove
- Merge: Between concurrent add and remove, add wins.



Aside: Rich Data Types

Conflict Free Replicated Datatypes (CRDT)

Concurrent updates without synchronisation + Guaranteed Convergence

Key Properties:

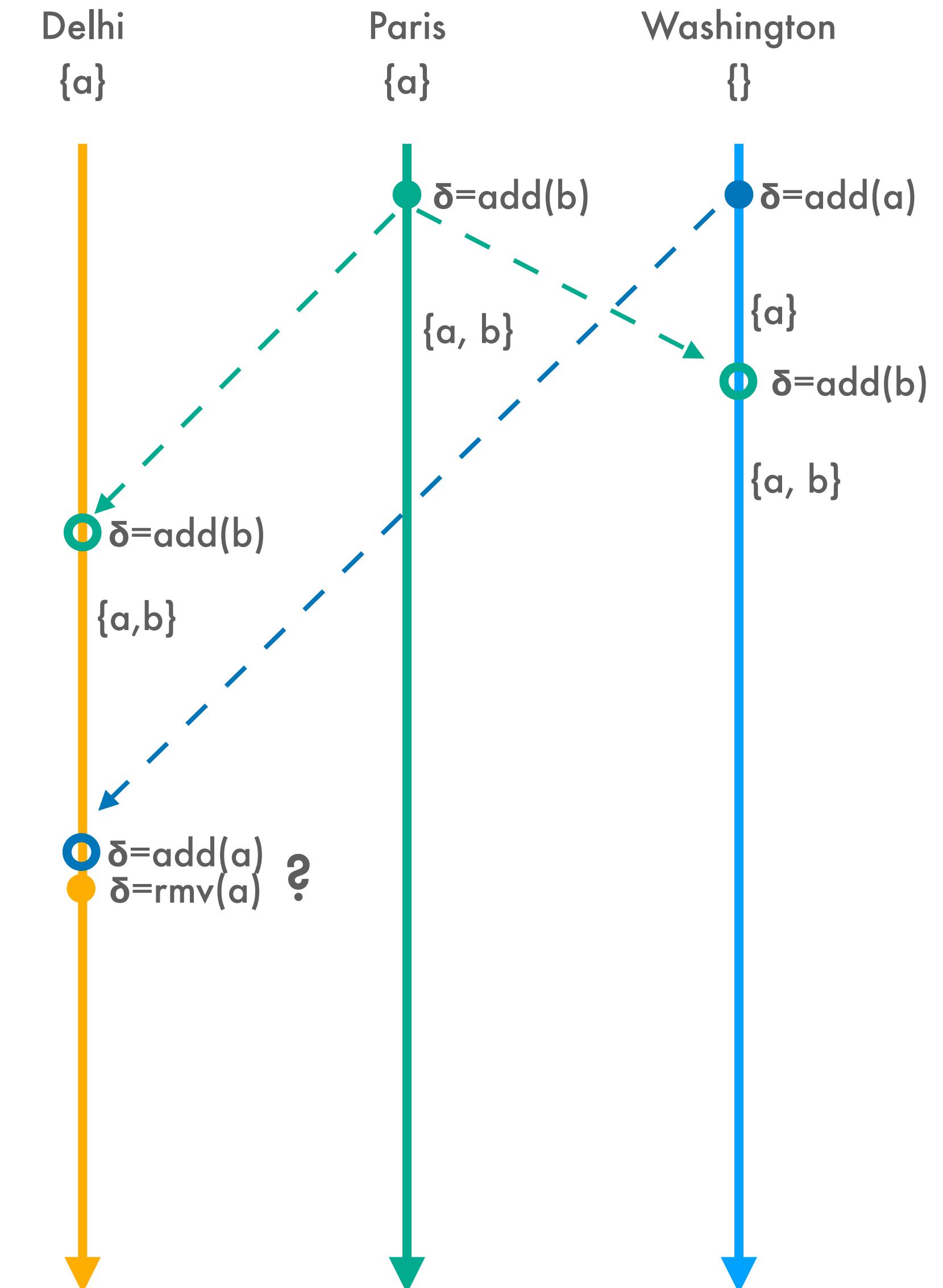
- Commutative, Associative and Idempotent Operations
- Mergeable States
- No Locks or consensus for updates

Implementation:

- **Effect (δ):** Takes the state of the CRDT and produces a new state
- **Merge:** Takes two CRDT states and **deterministically** produces a new state

Example: Add-Wins Set

- Effects: Add, Remove
- Merge: Between concurrent add and remove, add wins.



Aside: Rich Data Types

Conflict Free Replicated Datatypes (CRDT)

Concurrent updates without synchronisation + Guaranteed Convergence

Key Properties:

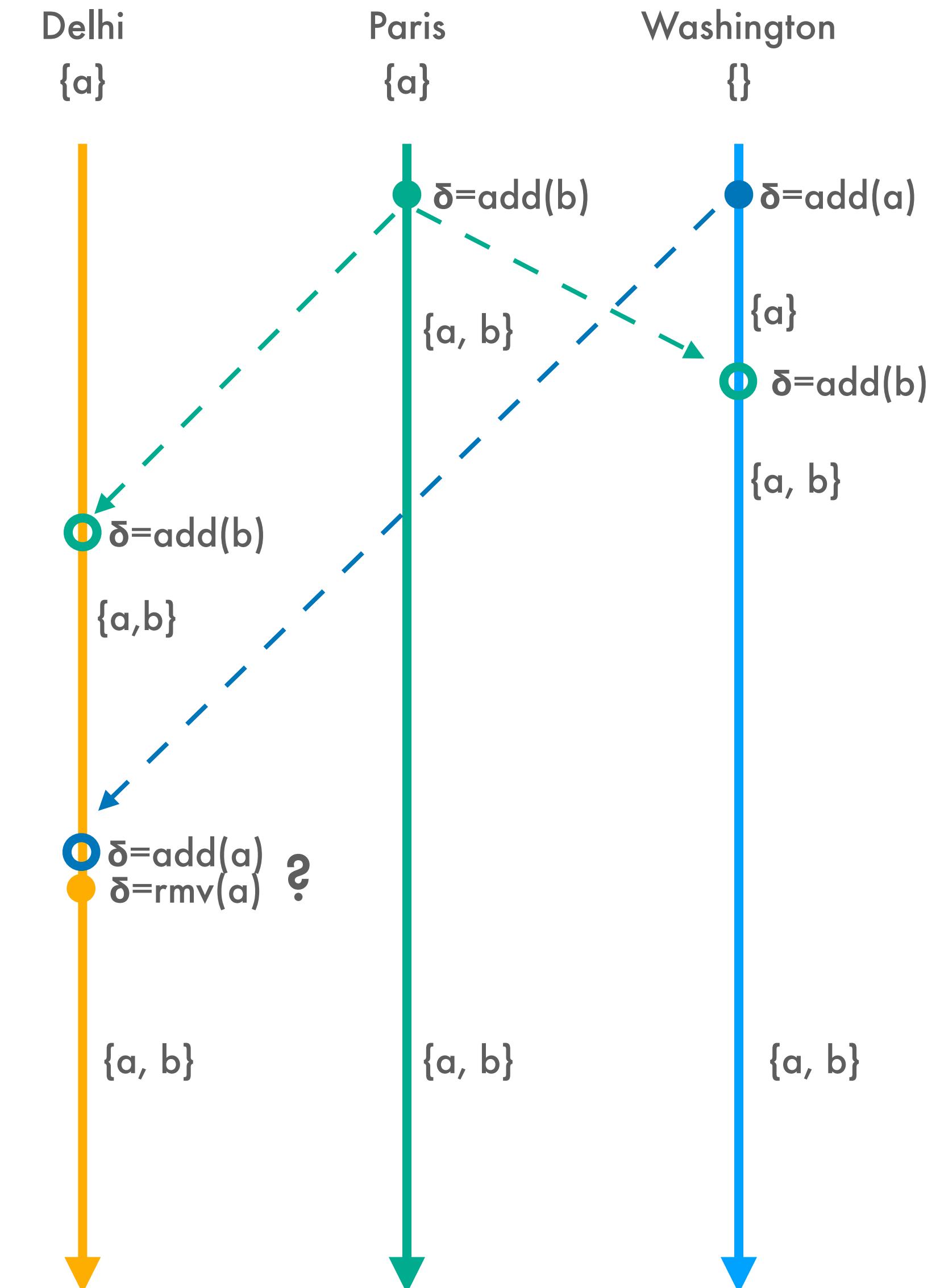
- Commutative, Associative and Idempotent Operations
- Mergeable States
- No Locks or consensus for updates

Implementation:

- **Effect (δ):** Takes the state of the CRDT and produces a new state
- **Merge:** Takes two CRDT states and **deterministically** produces a new state

Example: Add-Wins Set

- Effects: Add, Remove
- Merge: Between concurrent add and remove, add wins.



Transactions

Transaction is a block of operations \approx Session Order

- Begin $\circlearrowleft B$, Update $\circlearrowleft U$, Read, Commit $\circlearrowleft C$, Abort $\circlearrowleft A$

Begins with a snapshot timestamp (st) and commits with a commit timestamp (ct)

Sees updates from other transactions \approx Visibility

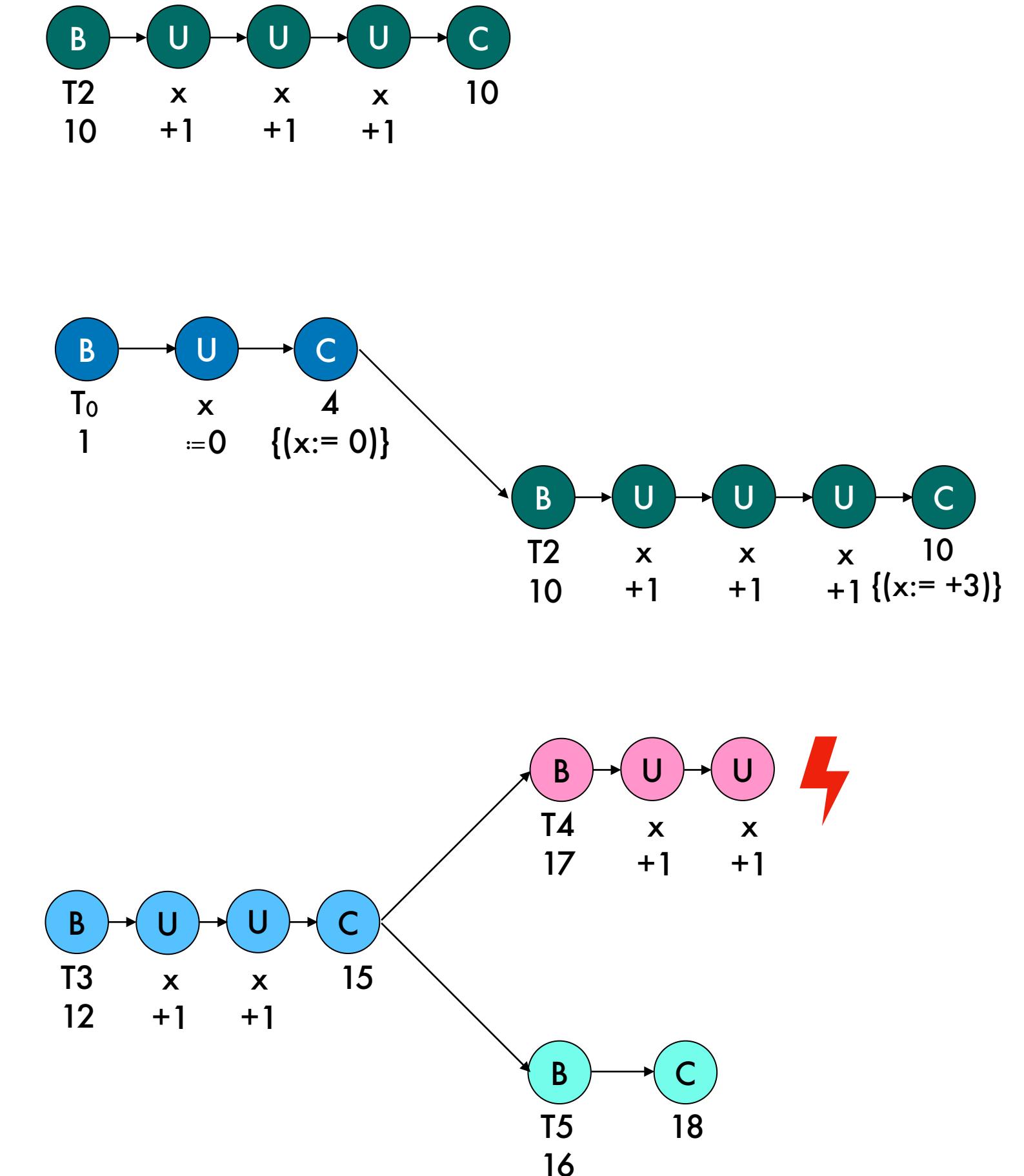
- $T_1 <_{vis} T_2 \Leftrightarrow T_1.ct < T_2.st$

Reads in a transaction:

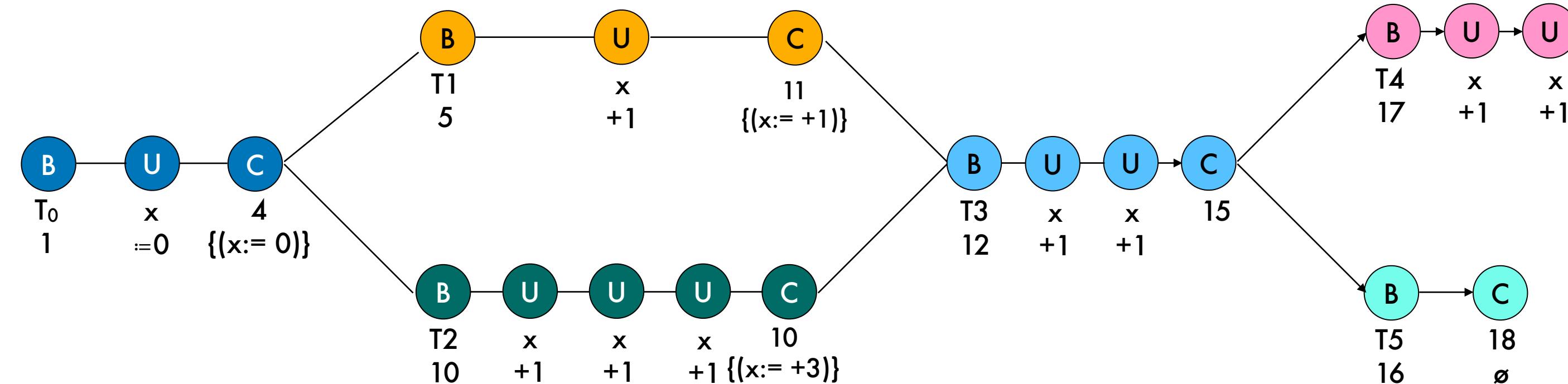
- read = snapshot \odot local updates

Crashes?

- Durable once committed, aborted otherwise



Traces - Histories of Transactions



Directed Acyclic Graph (DAG)

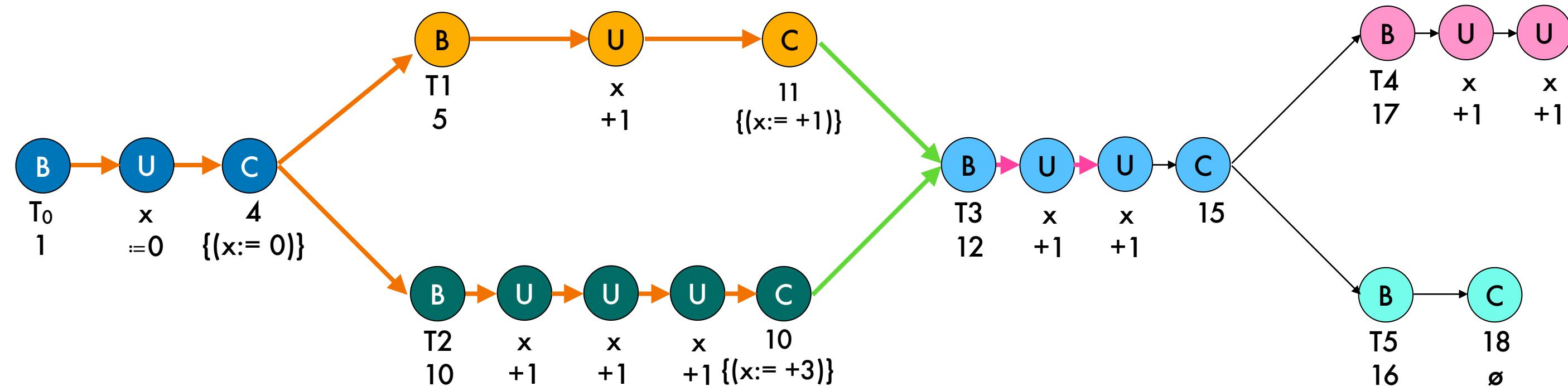
Vertices = Transaction events

Edges = Execution order

Some Metadata

Transaction IDs, Snapshot Timestamps, Datatypes, Effects etc.

Traces - Correctness



Correctness Criterion [Cerone et. al.] :

- Trace is Transactionally Causally Consistent (TCC)
- $TCC = Ext \wedge Int \wedge TransVis$

→ **Ext:** Snapshot Reads: First reads contains all previous transactions.

- Read rule: First read is a snapshot

→ **Int:** Read-My-Writes

- Read rule: later reads apply local effects

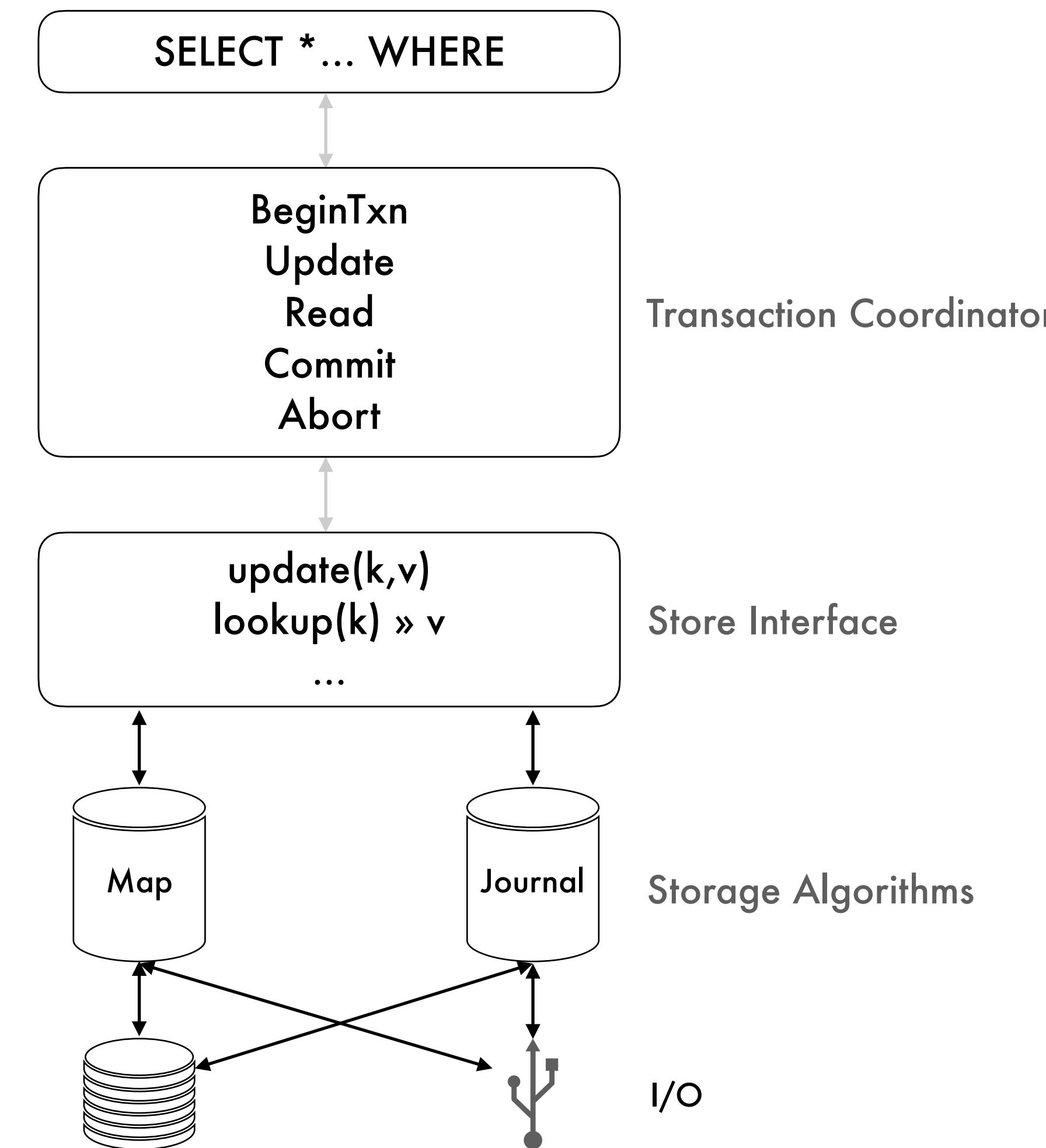
→ **TransVis:** Visibility is transitive

- By construction of a well formed trace

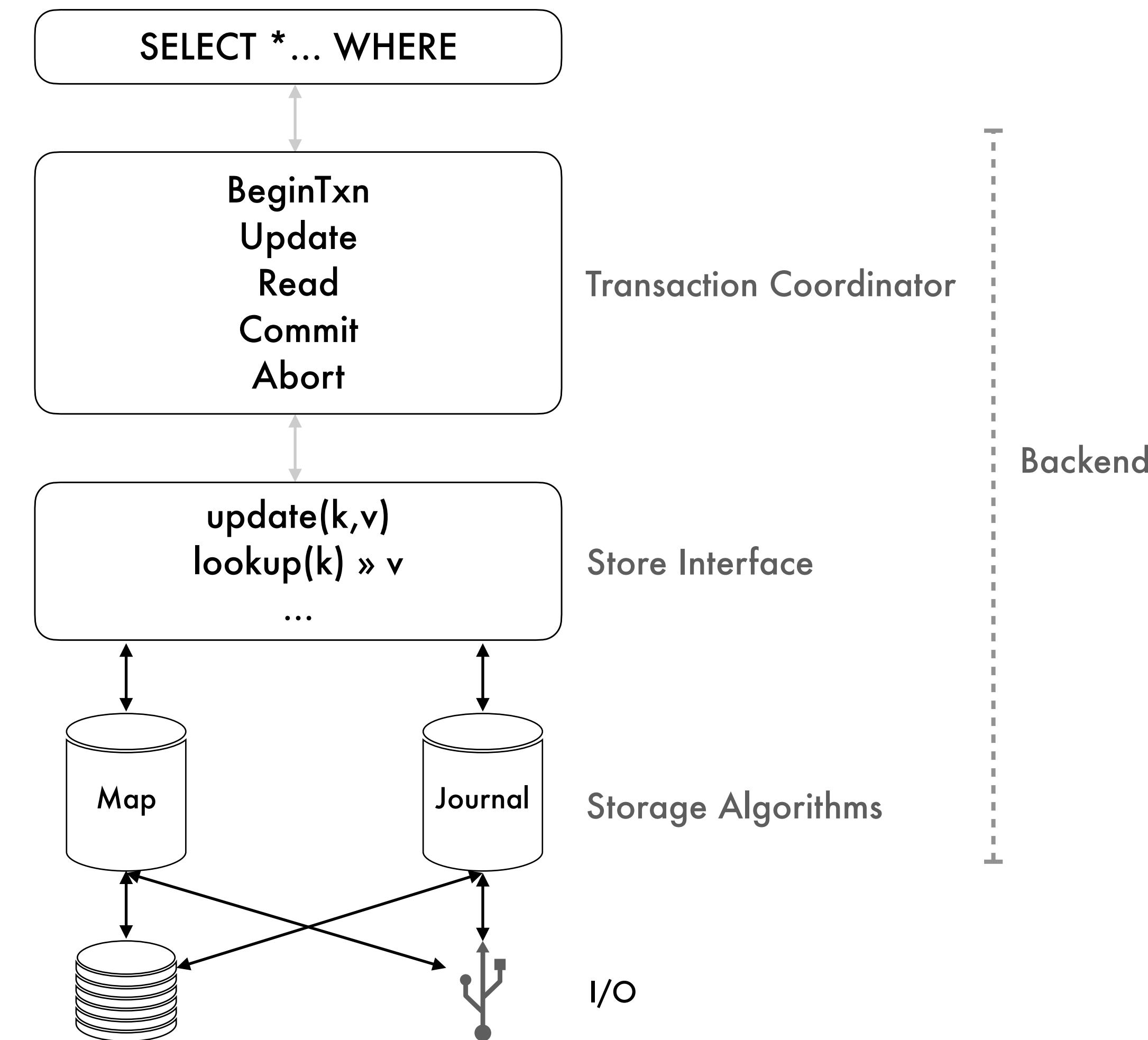
Why TCC?

TCC is the strongest consistency level which guarantees availability under partitions [CAP]

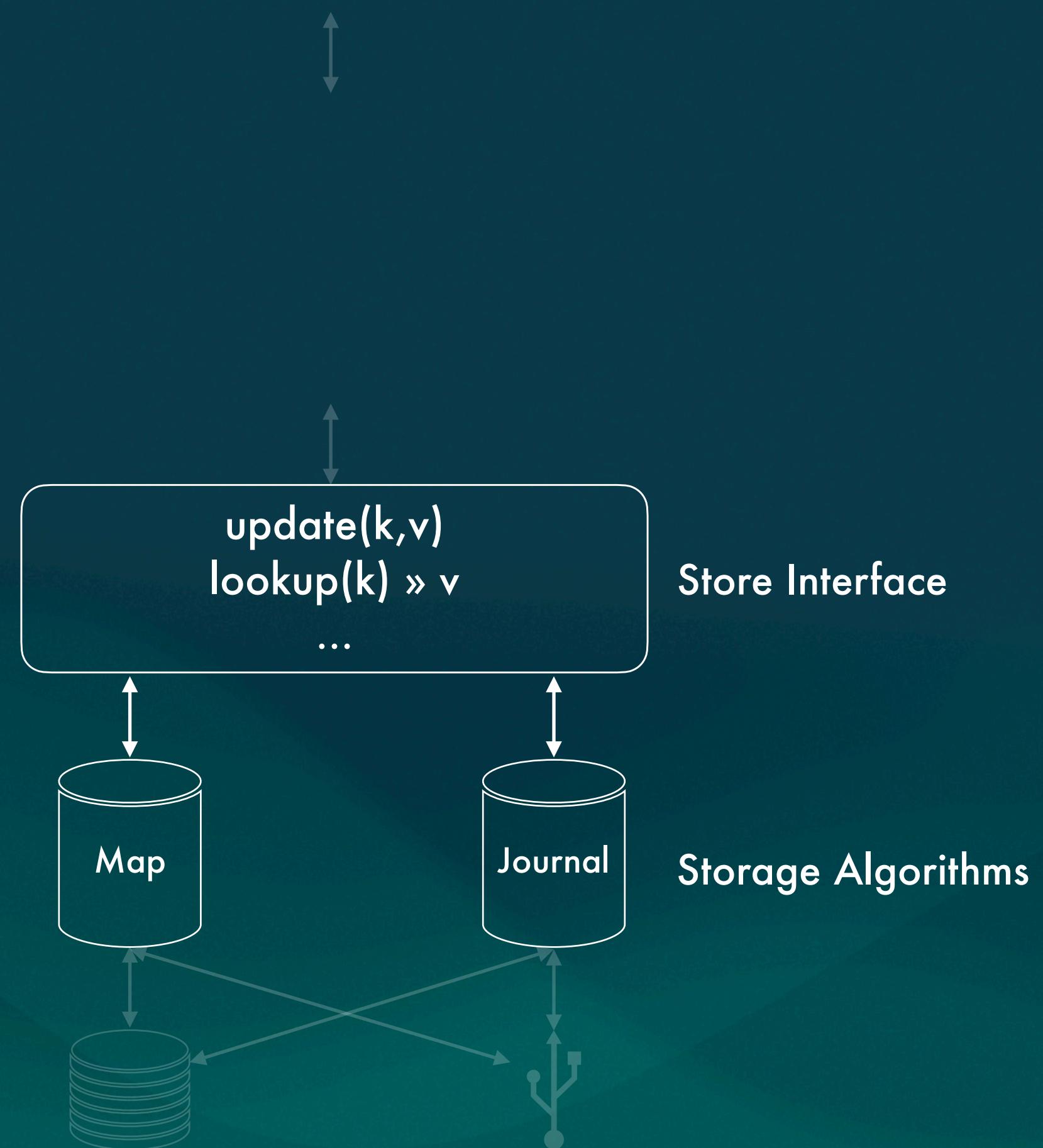
Database Backend?



Database Backend?



Stores



Store

Abstraction: Shared memory object (σ) that returns whatever is written to it.

Approach:

- $[\text{Key}, \text{Version}] \rightarrow \text{Value}_{\perp}$
- $\text{Lookup}(\sigma, k, t) \rightarrow \delta_{\perp}$

Variants:

Trivial:

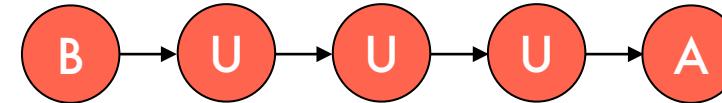
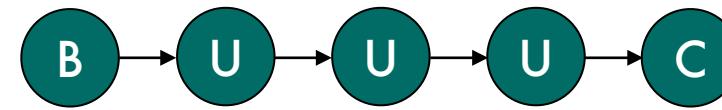
- Journal
- Map

Non-trivial compositions:

- Write Ahead Log (Map + Journal)
- LSM-Tree (Tree of Maps + Journal Blocks = WAL Memtables)

key	version	0	1	2	3	4	5	6	7	8	9	10	11	12	13	...
x		\perp														
y		\perp														
foobar		\perp		=	A	=	AB				=	0		=	\perp	
salary		\perp						+	100		+	200				

Store Interface



$doBegin(\sigma, \tau, st) \longmapsto \sigma'$

Start a transaction, add the transaction record to the store

$lookup(\sigma, k, st) \longmapsto \sigma$

Fetch the value of key 'k' at timestamp 'st'

$doUpdate(\sigma, \tau, k, st, \delta) \longmapsto \sigma'$

Update the value of key 'k' by applying the effect 'delta', within transaction 'tau' at timestamp 'st'.

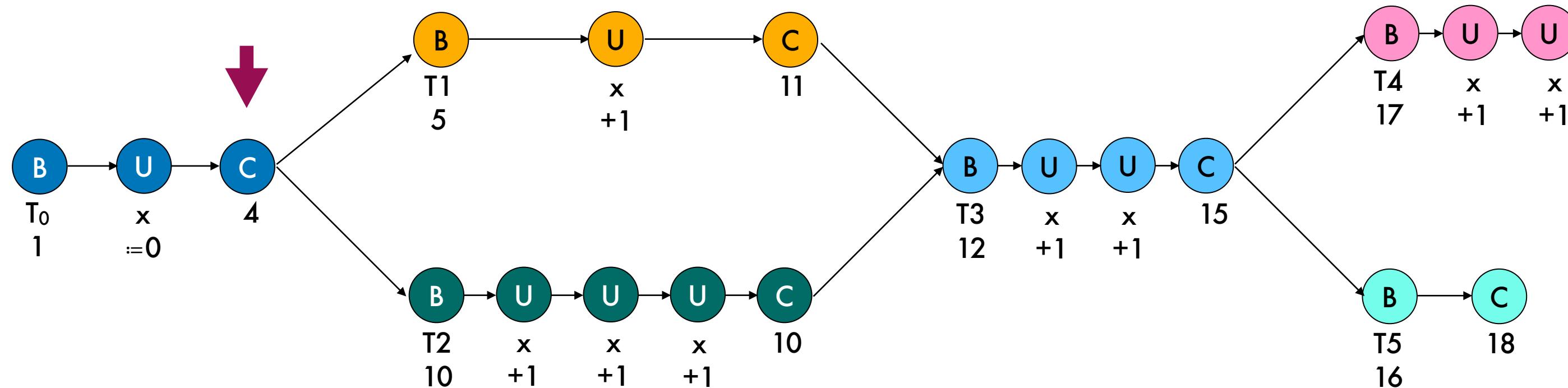
$doAbort(\sigma, \tau, st) \longmapsto \sigma'$

Abort transaction 'tau' which began at timestamp 'st'

$doCommit(\sigma, \tau, st, \mathcal{B}, ct) \longmapsto \sigma'$

Commit transaction 'tau', which began at 'st', with timestamp 'ct' and make effects from 'B' visible

Journal Store

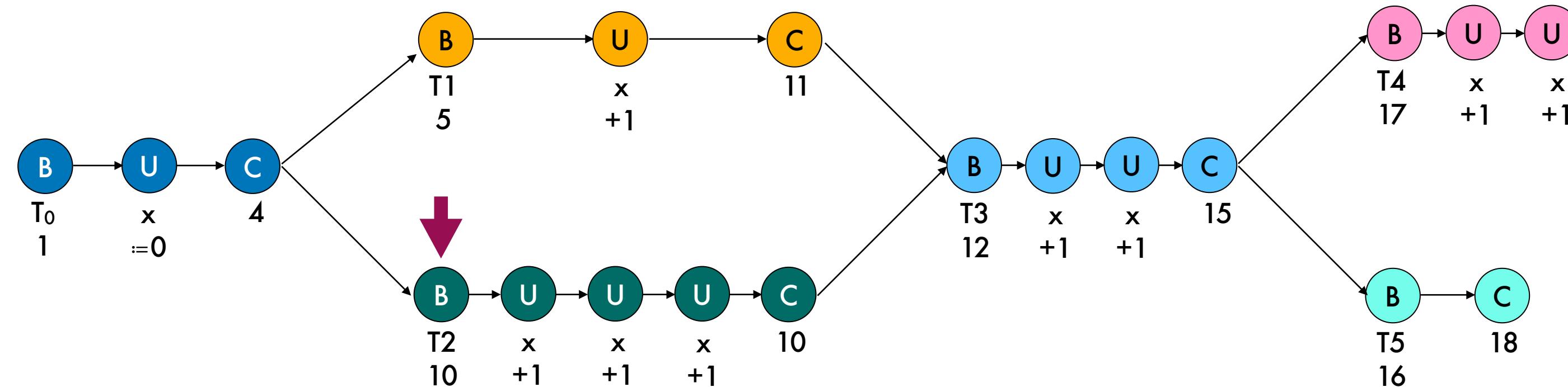


LSN	1	2	3
type	B	U	C
Txn ID	T ₀	T ₀	T ₀
st	1		
key	x		
effect	= 0		
ct			4
poststate(x)	⊥	0	0

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

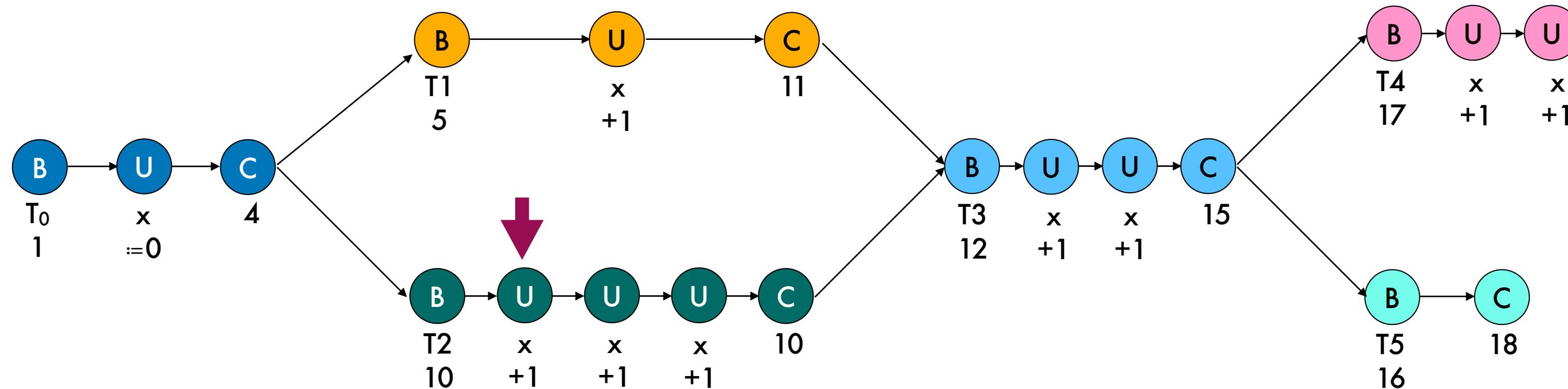


<i>LSN</i>	1	2	3	4
<i>type</i>	B	U	C	B
<i>Txn ID</i>	T0	T0	T0	T2
<i>st</i>	1			10
<i>key</i>	x			
<i>effect</i>	$\text{:= } 0$			
<i>ct</i>		4		
<i>poststate(x)</i>	⊥	0	0	0

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

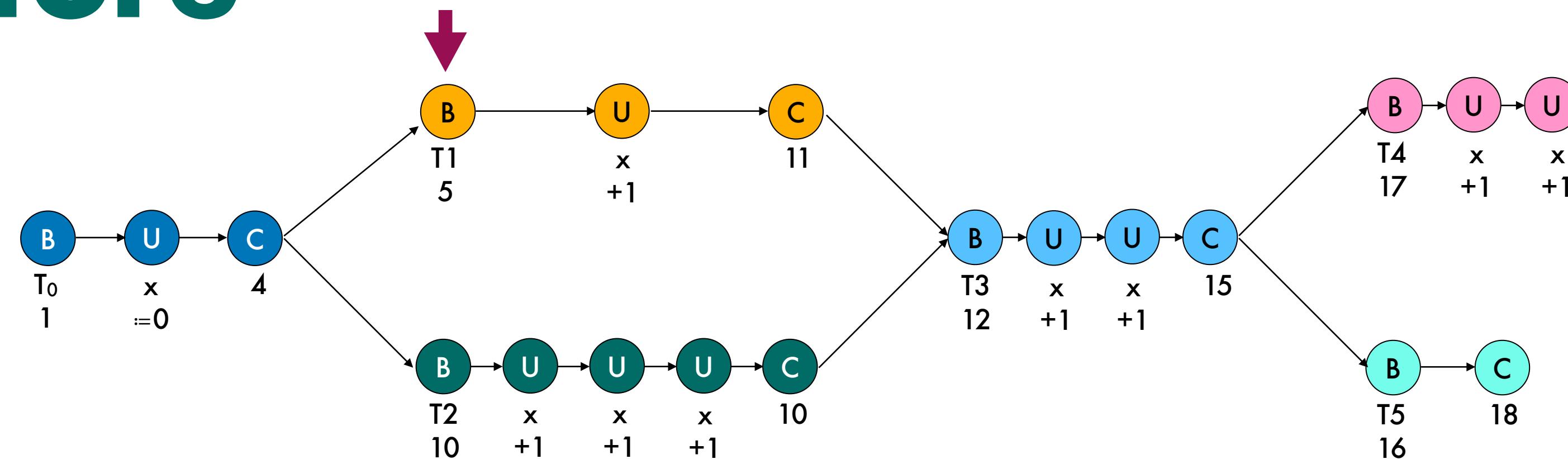


	1	2	3	4	5
LSN	1	2	3	4	5
type	B	U	C	B	U
Txn ID	T0	T0	T0	T2	T2
st	1			10	
key	x			x	
effect	$\text{:= } 0$			$+1$	
ct					
poststate(x)	\perp	0	0	4	1

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

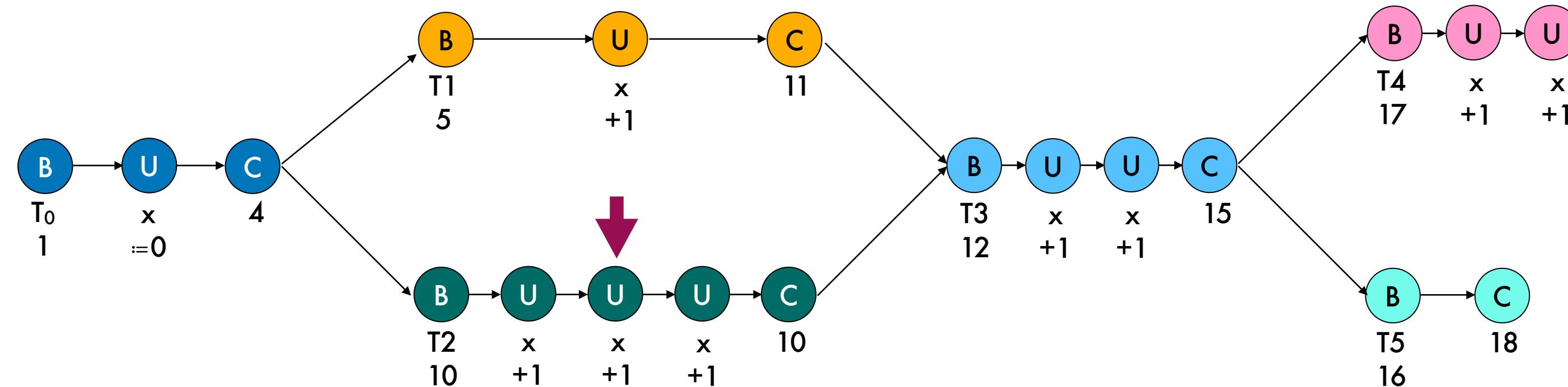


LSN	1	2	3	4	5	6
type	B	U	C	B	U	B
Txn ID	T0	T0	T0	T2	T2	T1
st	1			10		5
key	x			x		
effect	= 0			+1		
ct						
poststate(x)	⊥	0	0	1	1	

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

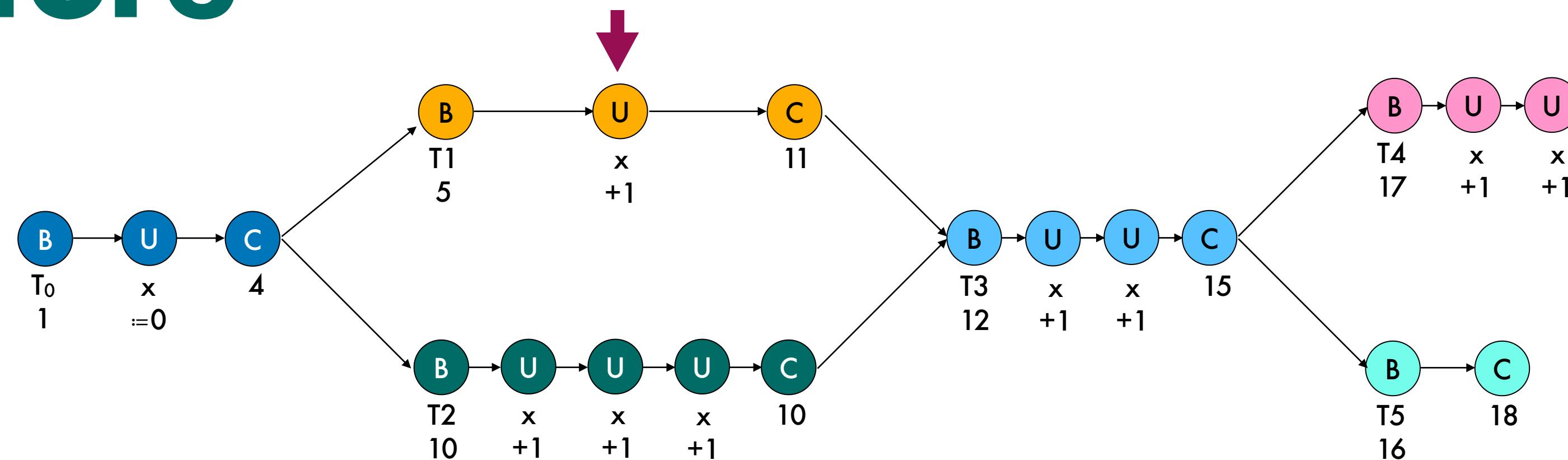


<i>LSN</i>	1	2	3	4	5	6	7
<i>type</i>	B	U	C	B	U	B	U
<i>Txn ID</i>	T0	T0	T0	T2	T2	T1	T2
<i>st</i>	1			10		5	
<i>key</i>	x			x		x	
<i>effect</i>	= 0			+1		+1	
<i>ct</i>							
<i>poststate(x)</i>	⊥	0	0	0	1	1	2

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

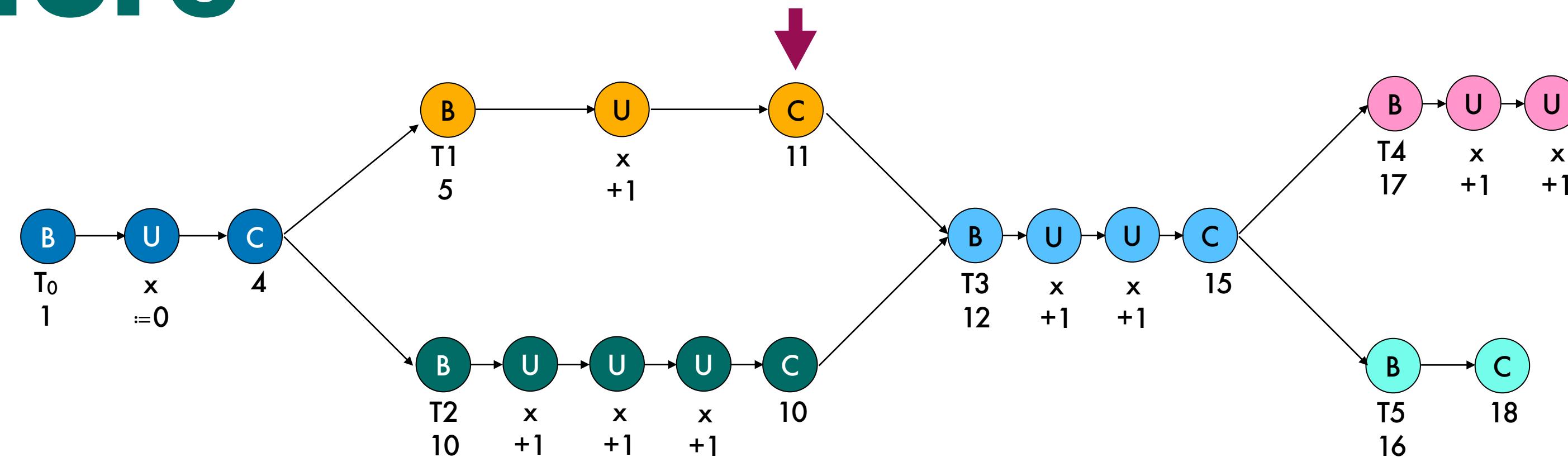


LSN	1	2	3	4	5	6	7	8
type	B	U	C	B	U	B	U	U
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1
st	1			10		5		
key		x			x		x	
effect		$x := 0$			+1		+1	+1
ct								
poststate(x)	⊥	0	0	0	1	1	2	1

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

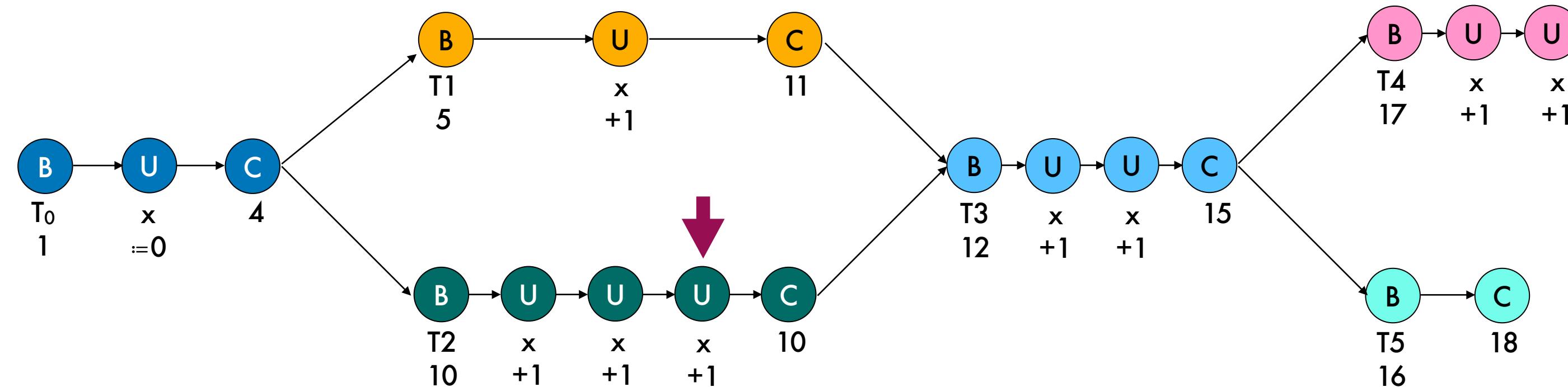
Journal Store



Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

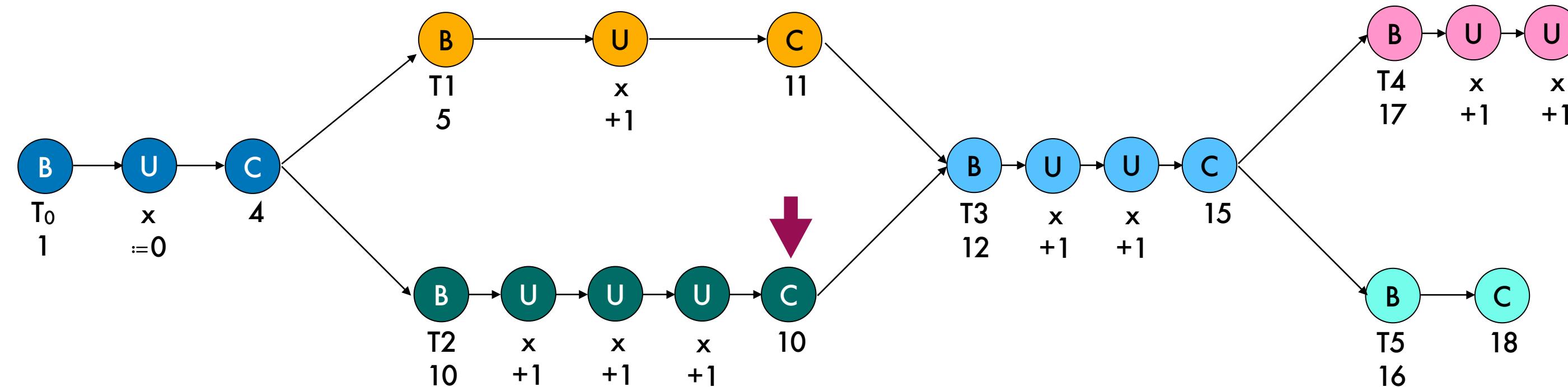


LSN	1	2	3	4	5	6	7	8	9	10
type	B	U	C	B	U	B	U	U	C	U
Txn ID	T ₀	T ₀	T ₀	T ₂	T ₂	T ₁	T ₂	T ₁	T ₁	T ₂
st	1			10		5				
key		x			x		x	x		x
effect		$x := 0$			+1		+1	+1		+1
ct								11		
poststate(x)	⊥	0	0	0	1	1	2	1	3	

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

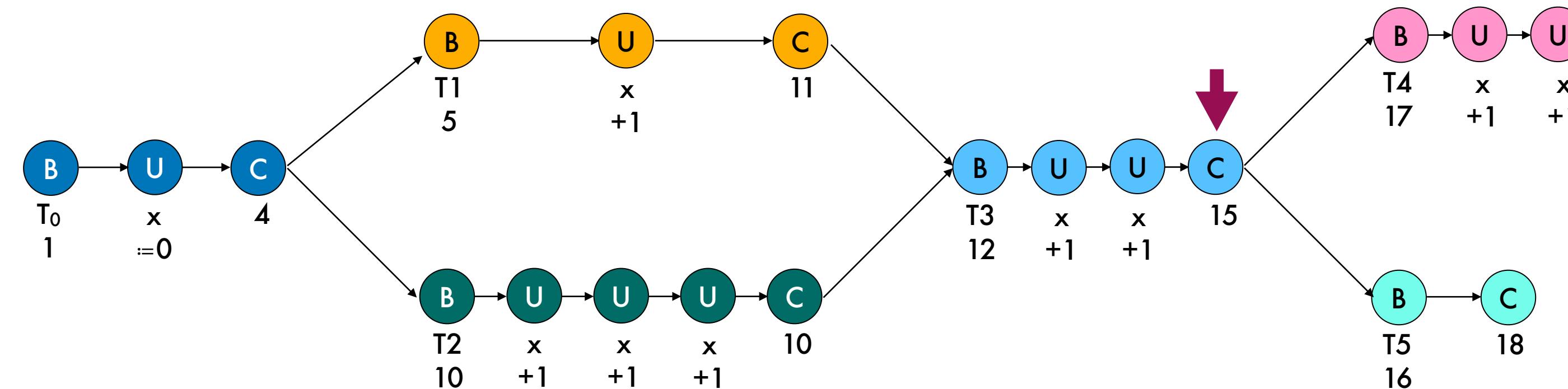


LSN	1	2	3	4	5	6	7	8	9	10	11
type	B	U	C	B	U	B	U	U	C	U	C
Txn ID	T ₀	T ₀	T ₀	T ₂	T ₂	T ₁	T ₂	T ₁	T ₁	T ₂	T ₂
st	1			10		5					
key	x			x		x	x	x	x		
effect	$x := 0$			+1		+1	+1	+1	+1		
ct							11		10		
poststate(x)	⊥	0	0	0	1	1	2	1	3	3	10

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

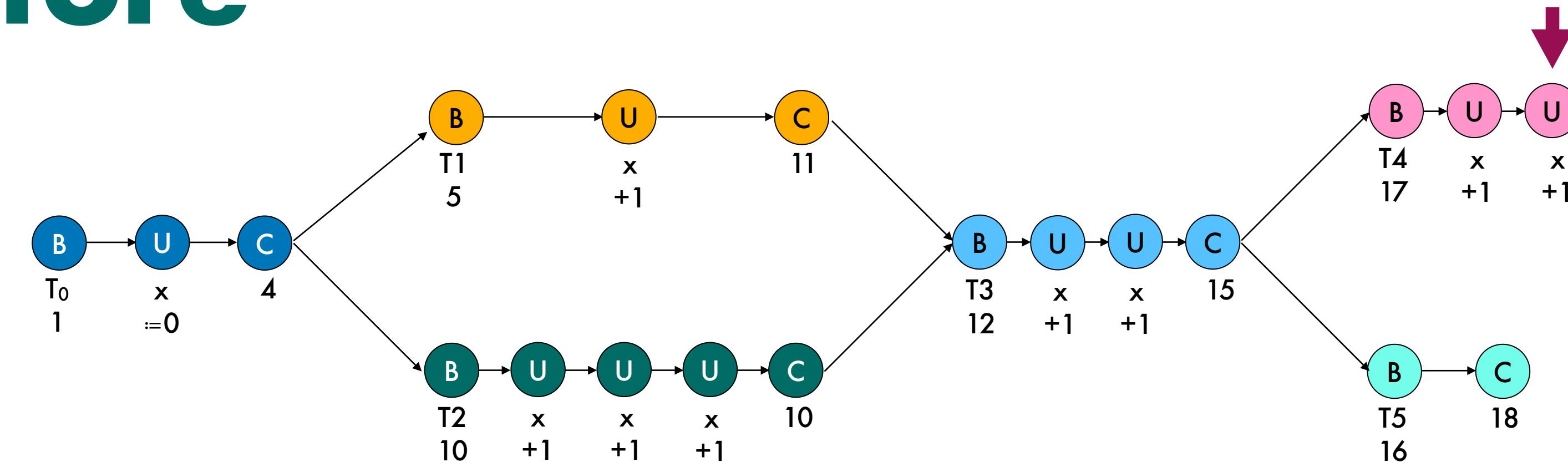


LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
type	B	U	C	B	U	B	U	U	C	U	C	B	U	U	C
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1	T1	T2	T2	T3	T3	T3	T3
st	1			10		5						12			
key	x			x			x	x		x		x	x		
effect	$x := 0$			+1			+1	+1		+1		+1	+1		
ct			4							10					15
poststate(x)	⊥	0	0	0	1	1	2	1	1	3	3	4	5	6	6

Every store operation is recorded in the journal as they occur.

Writing is fast, Reading is slow.

Journal Store

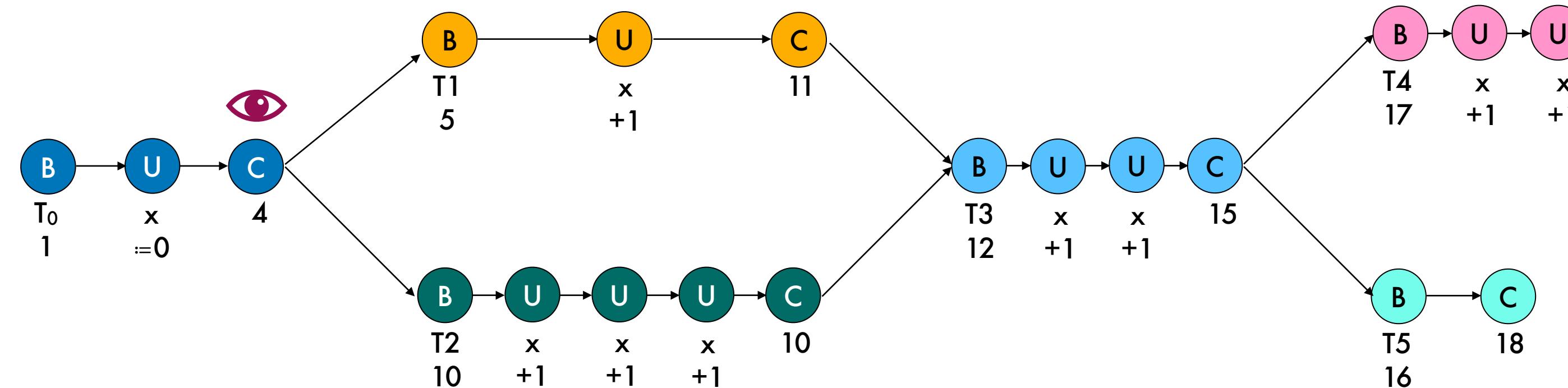


LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
type	B	U	C	B	U	B	U	U	C	U	C	B	U	U	C	B	U	B	C	U	
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1	T1	T2	T2	T3	T3	T3	T3	T4	T4	T5	T5	T4	
st	1			10		5						12				17		16		x	
key	x			x		x	x	x		x		x	x		x					x	
effect	= 0			+1		+1	+1	+1		+1		+1	+1		+1		+1		+1	+1	
ct				4						10					15			18		7	
poststate(x)	⊥	0	0	0	1	1	2	1	1	3	3	4	5	6	6	7	6	6	7	7	

Every store operation is recorded in the journal as they occur.

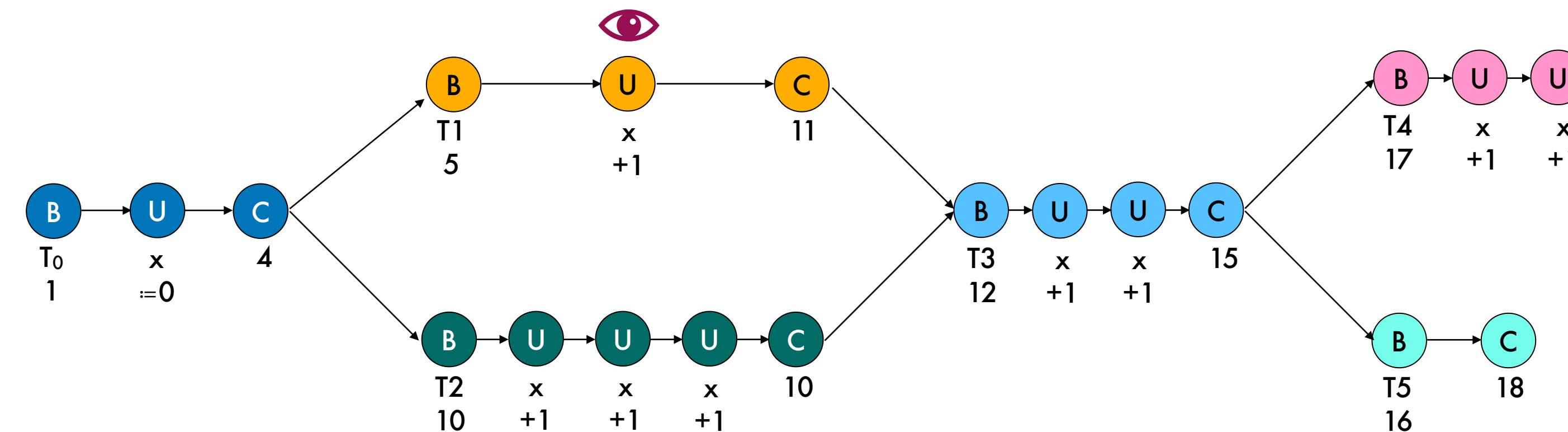
Writing is fast, Reading is slow.

Journal Store



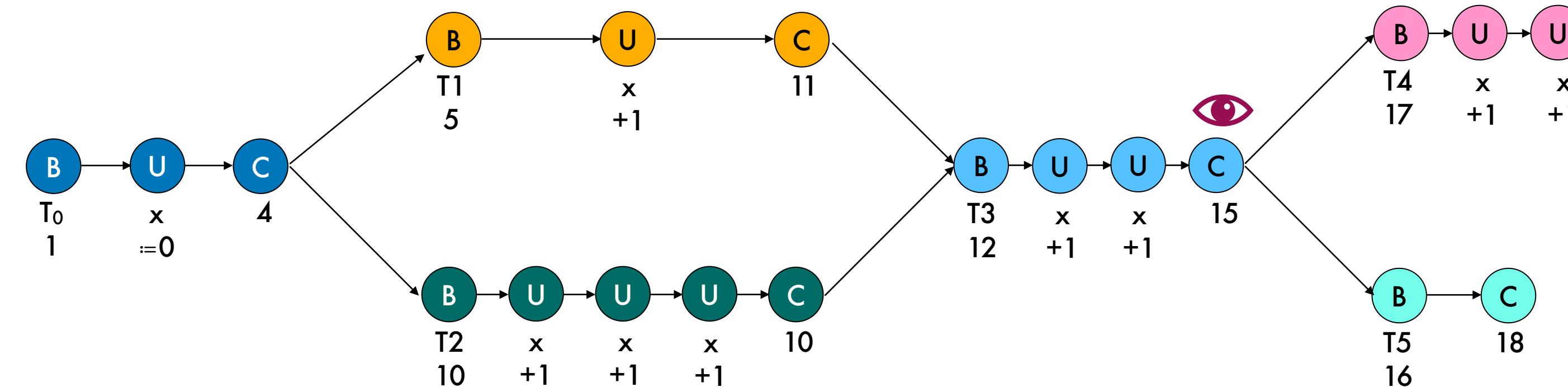
LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
type	B	U	C	B	U	B	U	C	U	C	B	U	U	C	B	U	B	C	B	U	
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1	T1	T2	T2	T3	T3	T3	T3	T4	T4	T5	T5	T4	
st	1			10		5					12				17		16				
key	x			x		x		x		x		x		x		x		x		x	
effect	= 0			+1		+1		+1		+1		+1		+1		+1		+1		+1	
ct			4																		
poststate(x)	⊥	0	0	0	1	1	2	1	3	3	10	4	5	6	15	6	18	6	7	7	

Journal Store



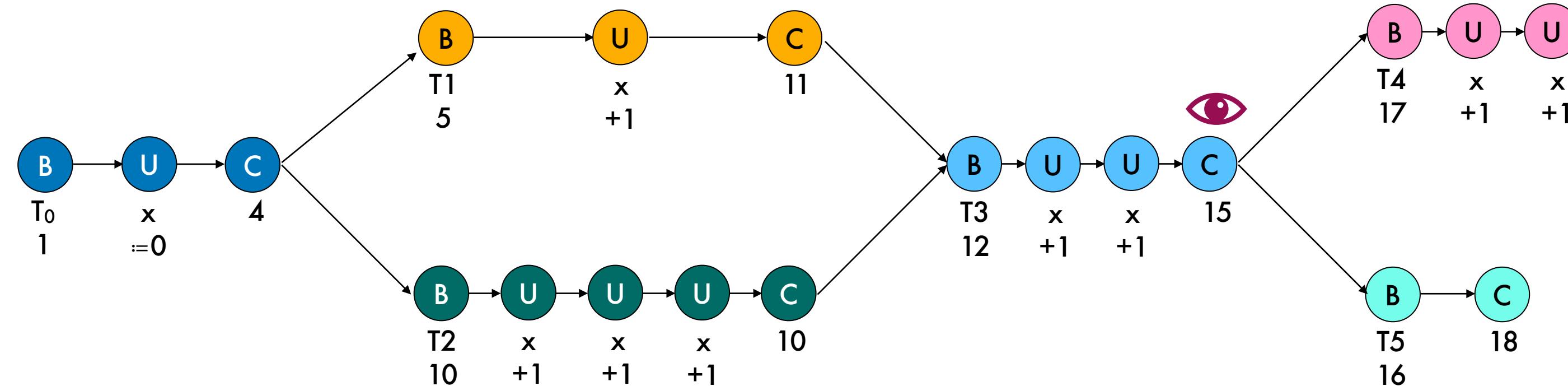
LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
type	B	U	C	B	U	B	U	U	C	U	C	B	U	U	C	B	U	B	C	U	
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1	T1	T2	T2	T3	T3	T3	T3	T4	T4	T5	T5	T4	
st	1			10		5					12					17		16			21
key	x				x		x	x		x		x	x			x			x		x
effect	$\text{:= } 0$				$+1$		$+1$	$+1$		$+1$		$+1$	$+1$		$+1$		$+1$		$+1$		$+1$
ct											11				15			18			
poststate(x)	\perp	0	0	0	1	1	2	1	3	3	10	4	5	6	6	6	7	6	7	7	

Journal Store



LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
type	B	U	C	B	U	B	U	U	C	U	C	B	U	U	C	B	U	B	C	U	
Txn ID	T ₀	T ₀	T ₀	T ₂	T ₂	T ₁	T ₂	T ₁	T ₁	T ₂	T ₂	T ₃	T ₃	T ₃	T ₃	T ₄	T ₄	T ₅	T ₅	T ₄	
st	1			10		5						12				17		16		x	
key		x			x		x	x		x			x	x			x			x	
effect		= 0			+1		+1	+1		+1			+1	+1			+1		+1	+1	+1
ct											10				15				18		
poststate(x)	⊥	0	0	0	1	1	2	1	1	3	3	4	5	6	6	6	7	6	7	7	

Journal Store

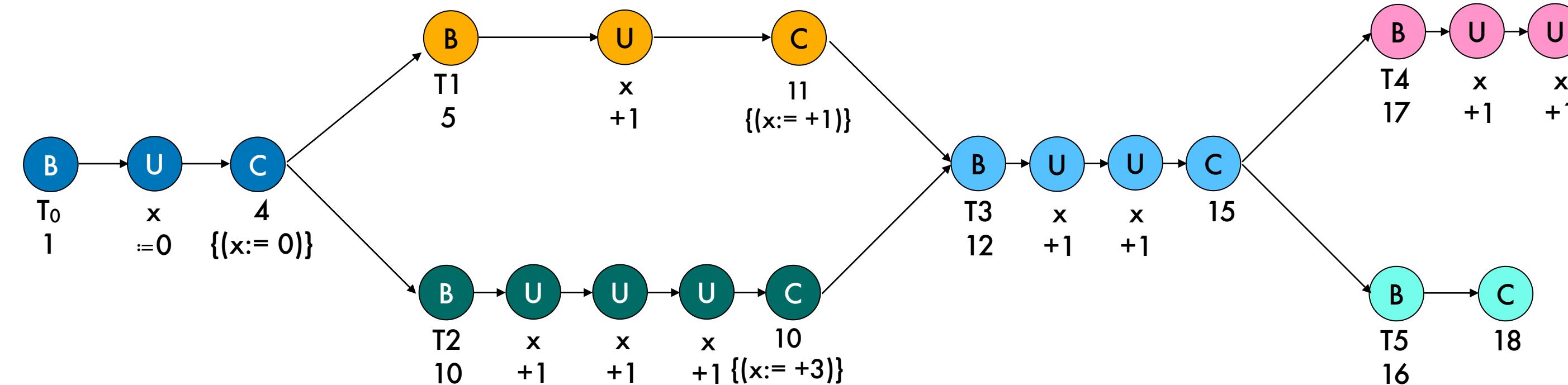


LSN	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
type	B	U	C	B	U	B	U	U	C	U	C	B	U	U	C	B	U	B	C	U	
Txn ID	T0	T0	T0	T2	T2	T1	T2	T1	T1	T2	T2	T3	T3	T3	T3	T4	T4	T5	T5	T4	
st	1			10		5					12				17		16				
key		x			x		x	x		x			x	x		x			x		
effect		≈ 0			+1		+1	+1		+1			+1	+1		+1		+1		+1	
ct				4						10				15				18			
poststate(x)	\perp	0	0	0	1	1	2	1	1	3	3	4	5	6	6	7	6	7	6	7	

$$\begin{aligned}
 \text{doBegin}(\sigma_J, \tau, \text{st}) &= \sigma_J \triangleright (\text{beginTxn}, \tau, \text{st}) \\
 \text{lookup}(\sigma_J, k, \text{st}) &= \text{poststate}_{\sigma_J}(r, k) \text{ where } r = (\text{beginTxn}, _, \text{st}) \\
 \text{doUpdate}(\sigma_J, \tau, k, \text{st}, \delta) &= \sigma_J \triangleright (\text{update}, \tau, k, \delta) \\
 \text{doAbort}(\sigma_J, \tau, \text{st}) &= \sigma_J \triangleright (\text{abort}, \tau) \\
 \text{doCommit}(\sigma, \tau, \text{st}, \mathcal{B}, \text{ct}) &= \sigma_J \triangleright (\text{commit}, \tau, \text{ct})
 \end{aligned}$$

$$\text{poststate}_{\sigma_J}(r, k) \stackrel{\text{def}}{=} \begin{cases} \text{merge} (\{ \text{poststate}_{\sigma_J}(r', k) : r' \in \text{parent}_J(r) \}) & \text{if } r = (\text{beginTxn}, _, _) \\ \text{poststate}_{\sigma_J}(\text{parent}_J(r), k) & \text{if } r = (\text{update}, _, k', _) \wedge k \neq k' \\ \text{poststate}_{\sigma_J}(\text{parent}_J(r), k) \odot \delta & \text{if } r = (\text{update}, _, k, \delta) \\ \text{poststate}_{\sigma_J}(\text{parent}_J(r), k) & \text{if } r = (\text{commit}, _, _) \vee r = (\text{abort}, _) \end{cases}$$

Map Store



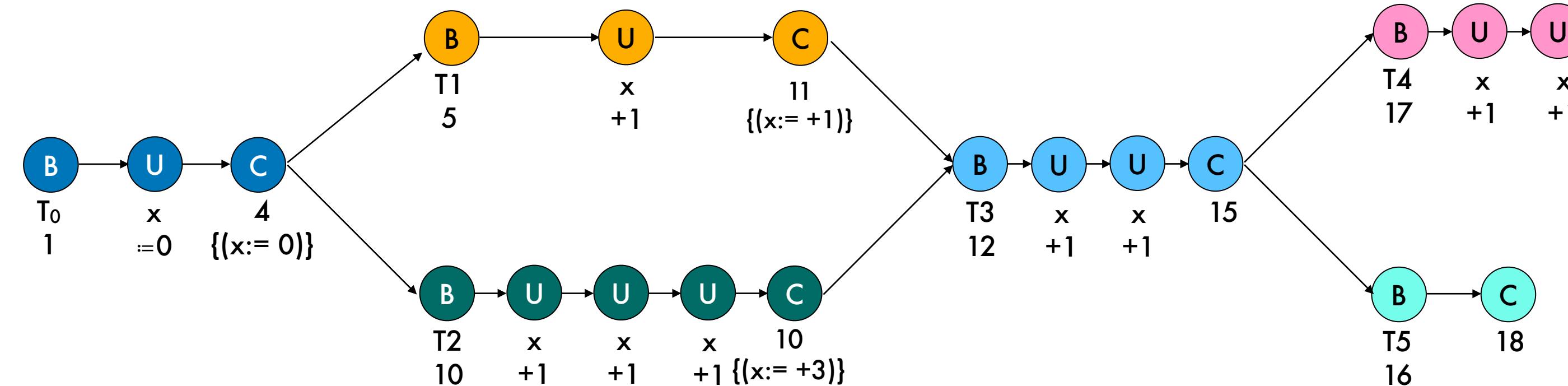
	key	x	x	x	x
version	4	11	10	15	
dependence	-	T0	T0	{T1, T2}	
Effect	:= 0	:= 1	:= 3	:= 6	

Only commits are recorded in the map, in no particular logical order.

Writing is slow, Reading is fast.

Effects are buffered by the transactions until commit i.e. Correctness only when transactional.

Map Store



key	x	x	x	x
version	4	11	10	15
dependence	-	T_0 $5 > 4$	T_0 $10 > 4$	$\{T_1, T_2\}$ $12 > 11, 10$
Effect	$x := 0$	$x := 1$	$x := 3$	$x := 6$

$\sigma_M \in (Key \times Timestamp \rightarrow Timestamp \times Assign) \subseteq \Sigma$, initially, $\sigma_M[k, t] = \perp$ for all k, t

$\text{lookup}(\sigma_M, k, st) = \text{merge}(\max_{\prec_M}(\{\sigma_M[k, w] \neq \perp : w < st\}).\delta)$

$\text{doCommit}(\sigma_M, __, st, \mathcal{B}, ct)[k, t] = \begin{cases} (st, \mathcal{B}[k]) & \text{if } \mathcal{B}[k] \neq \perp \wedge t = ct \\ \sigma_M[k, t] & \text{otherwise} \end{cases}$

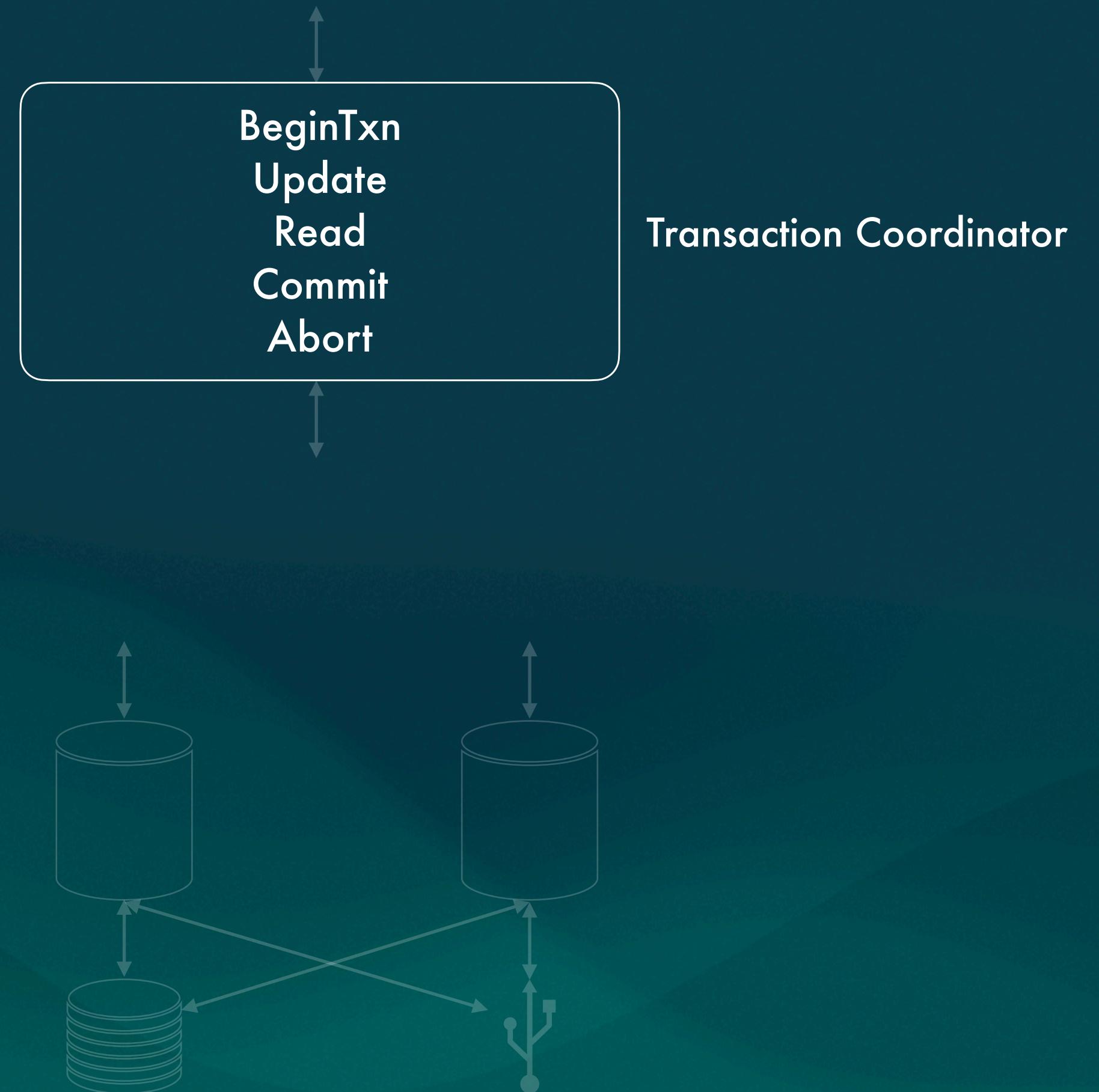
doBegin, doUpdate, doAbort: no-ops

Only commits are recorded in the map, in no particular logical order.

Writing is slow, Reading is fast.

Effects are buffered by the transactions until commit.

Transactions



Recap

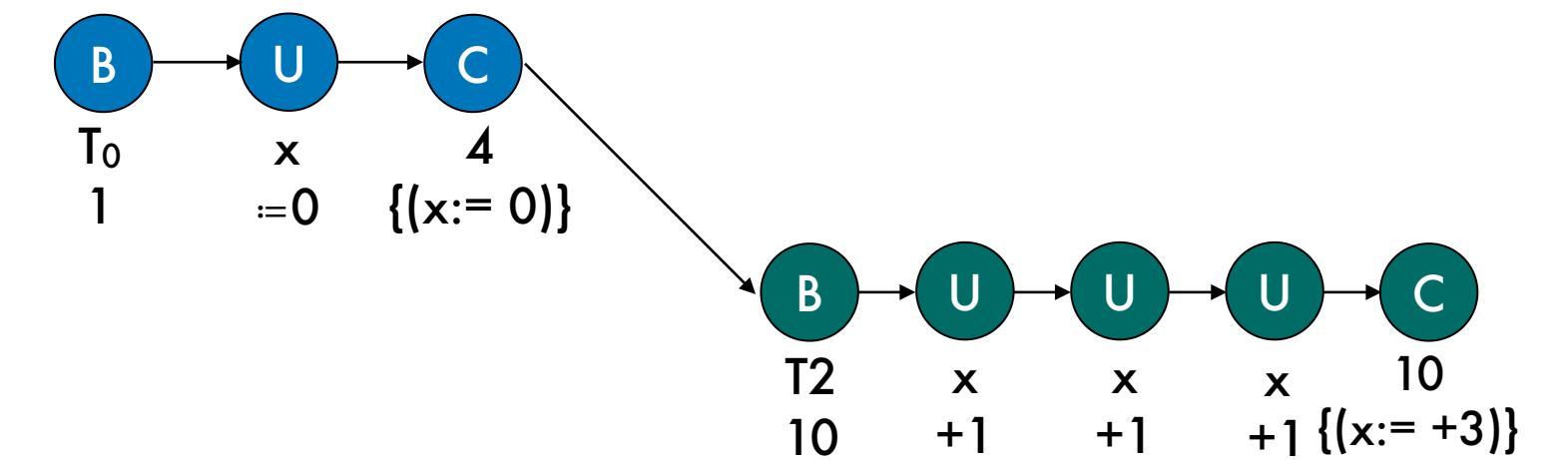
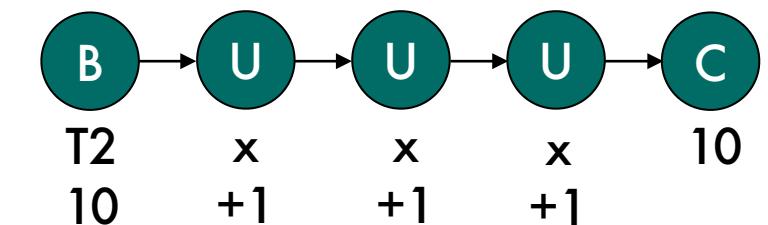
Block of Operations \approx Session Order

Has a st and a ct (if commits)

Sees updates from other transactions \approx Visibility

Reads values from other transactions and self

Durable once committed, aborted otherwise



API:

beginTxn(st)



read(k) \rightarrow value

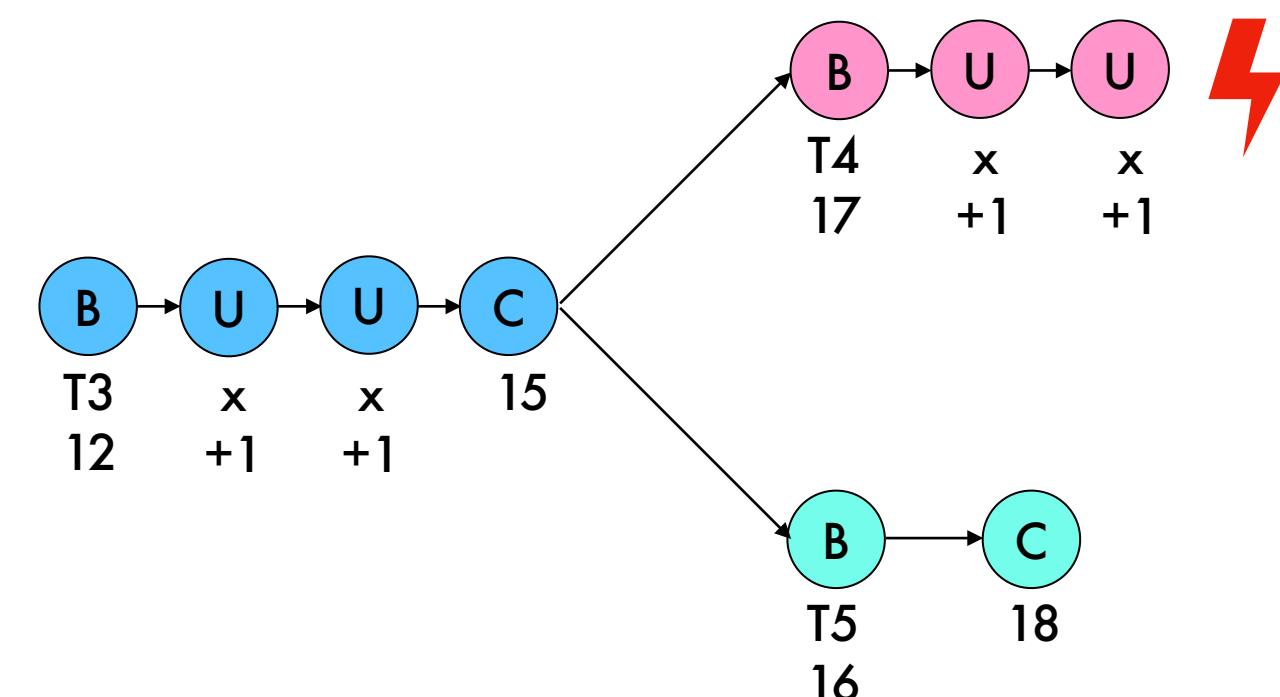


update(k, st, δ)

abort()



commit(ct)

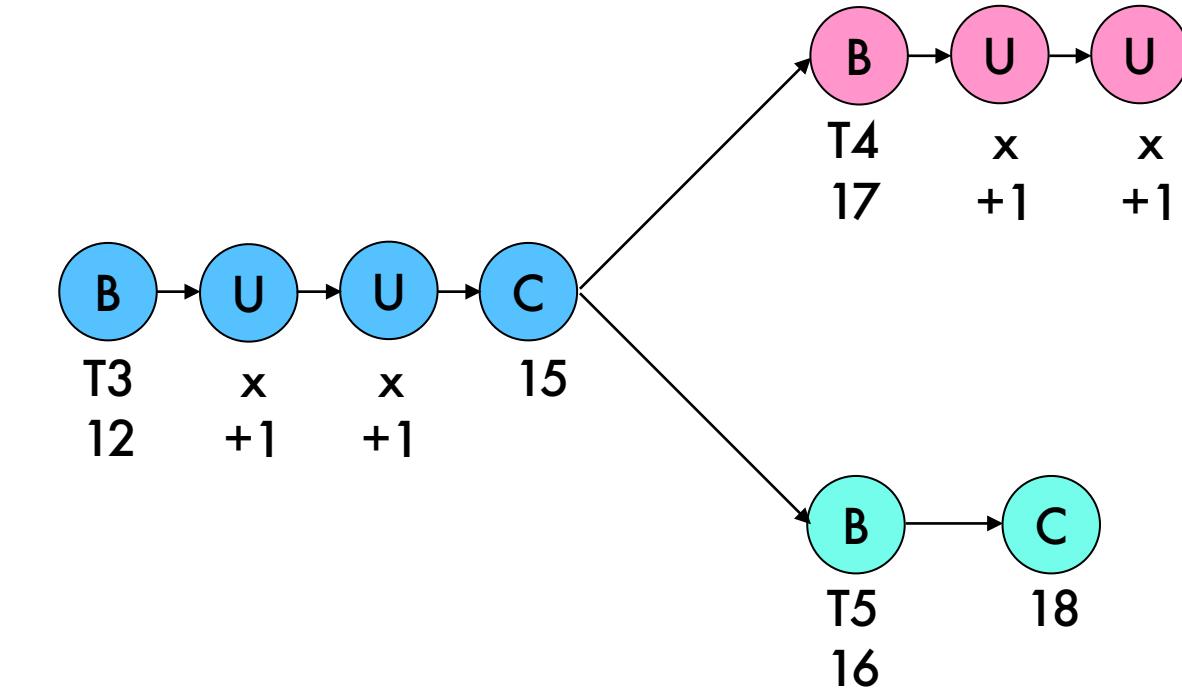


Transactions \leftrightarrow Stores

Transaction event...

...calls store method

<i>beginTxn (st)</i>	<i>doBegin (σ, τ, st)</i>
<i>update (k, δ)</i>	<i>doUpdate ($\sigma, \tau, k, st, \delta$)</i>
<i>read (k) $\rightarrow v$</i>	<i>lookup (σ, k, st) $\rightarrow \delta$</i>
<i>abort ()</i>	<i>doAbort (σ, τ, st)</i>
<i>commit (ct)</i>	<i>doCommit ($\sigma, \tau, st, \Delta, ct$)</i>



σ	store
τ	Transaction ID
<i>st</i>	snapshot timestamp
<i>k</i>	key
δ	effect
Δ	buffered effects
<i>ct</i>	commit timestamp

Transaction Semantics

$$\text{BEGINTXN} \quad \frac{\begin{array}{c} \forall T \in \mathcal{X}_a \uplus \mathcal{X}_c \uplus \mathcal{X}_r, T.\tau \neq \tau \\ \sigma' = \text{doBegin}(\sigma, \tau, \text{st}) \\ \mathcal{X}'_r = \mathcal{X}_r \uplus \{(\tau, \text{st}, \emptyset, \emptyset, \emptyset, _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{beginTxn(st)}} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

$$\text{UPDATE} \quad \frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \\ \delta \neq \perp \\ \mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta] \\ \sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta) \\ \mathcal{R}' = (\text{if } k \notin \mathcal{I} \text{ then } \mathcal{R} \text{ else } \mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]) \\ \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}', _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{update}(k, \delta)} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

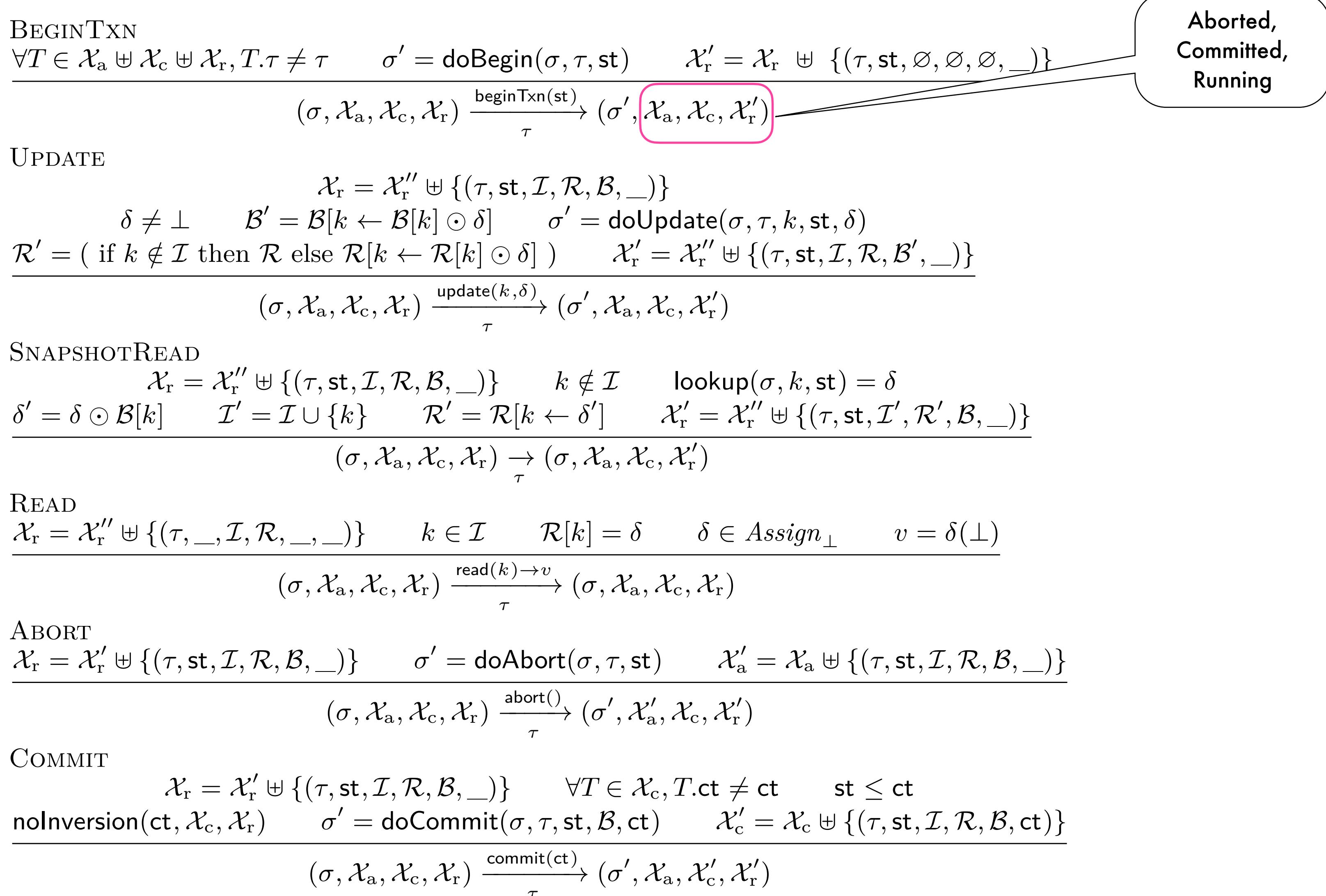
$$\text{SNAPSHOTREAD} \quad \frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \\ k \notin \mathcal{I} \\ \text{lookup}(\sigma, k, \text{st}) = \delta \\ \delta' = \delta \odot \mathcal{B}[k] \\ \mathcal{I}' = \mathcal{I} \cup \{k\} \\ \mathcal{R}' = \mathcal{R}[k \leftarrow \delta'] \\ \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}', \mathcal{R}', \mathcal{B}, _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

$$\text{READ} \quad \frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, _, \mathcal{I}, \mathcal{R}, _, _) \} \\ k \in \mathcal{I} \\ \mathcal{R}[k] = \delta \\ \delta \in \text{Assign}_{\perp} \\ v = \delta(\perp) \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{read}(k) \rightarrow v} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r)}$$

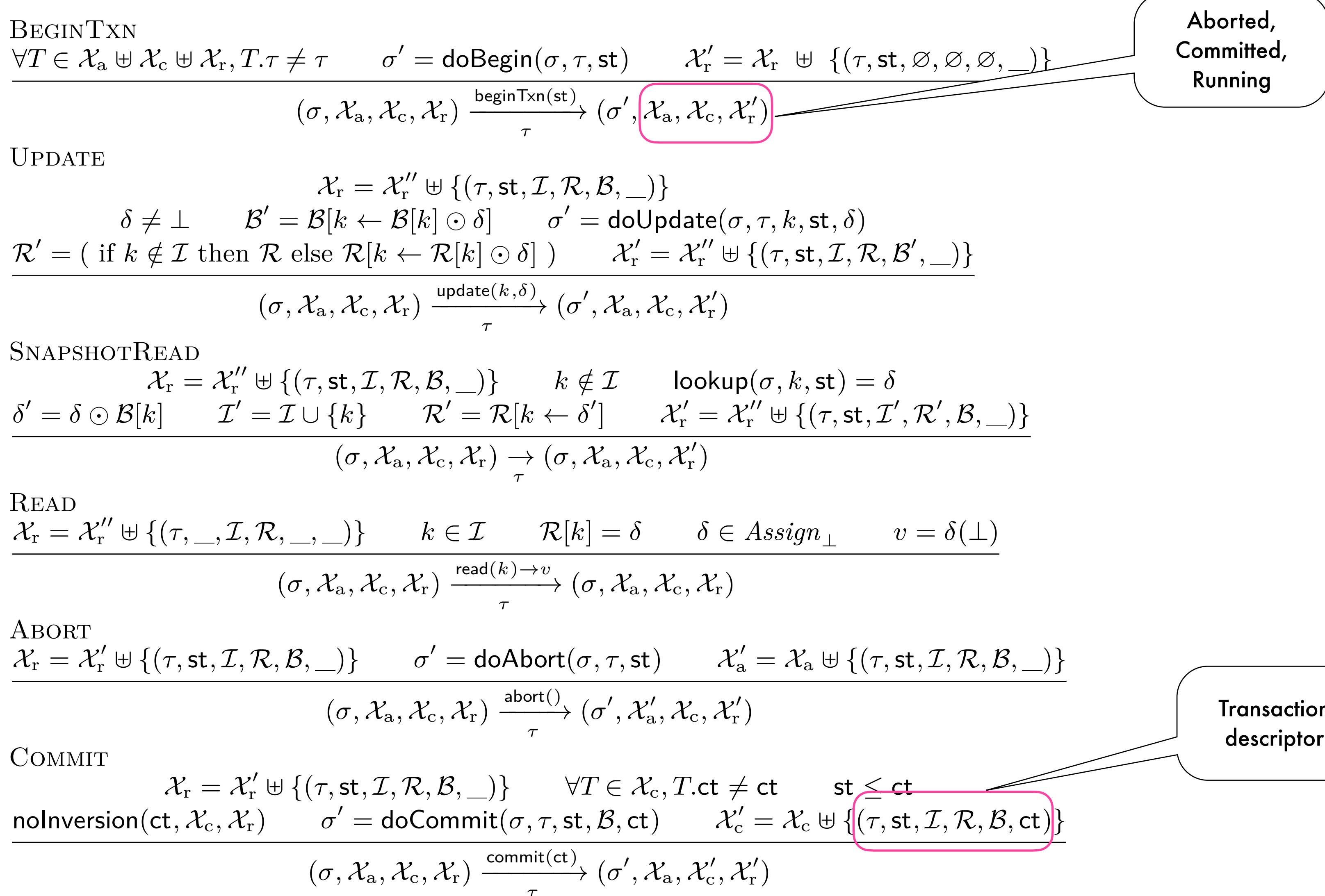
$$\text{ABORT} \quad \frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \\ \sigma' = \text{doAbort}(\sigma, \tau, \text{st}) \\ \mathcal{X}'_a = \mathcal{X}_a \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{abort}()} (\sigma', \mathcal{X}'_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

$$\text{COMMIT} \quad \frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \\ \forall T \in \mathcal{X}_c, T.\text{ct} \neq \text{ct} \\ \text{st} \leq \text{ct} \\ \text{noInversion}(\text{ct}, \mathcal{X}_c, \mathcal{X}_r) \\ \sigma' = \text{doCommit}(\sigma, \tau, \text{st}, \mathcal{B}, \text{ct}) \\ \mathcal{X}'_c = \mathcal{X}_c \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, \text{ct}) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{commit(ct)}} (\sigma', \mathcal{X}_a, \mathcal{X}'_c, \mathcal{X}'_r)}$$

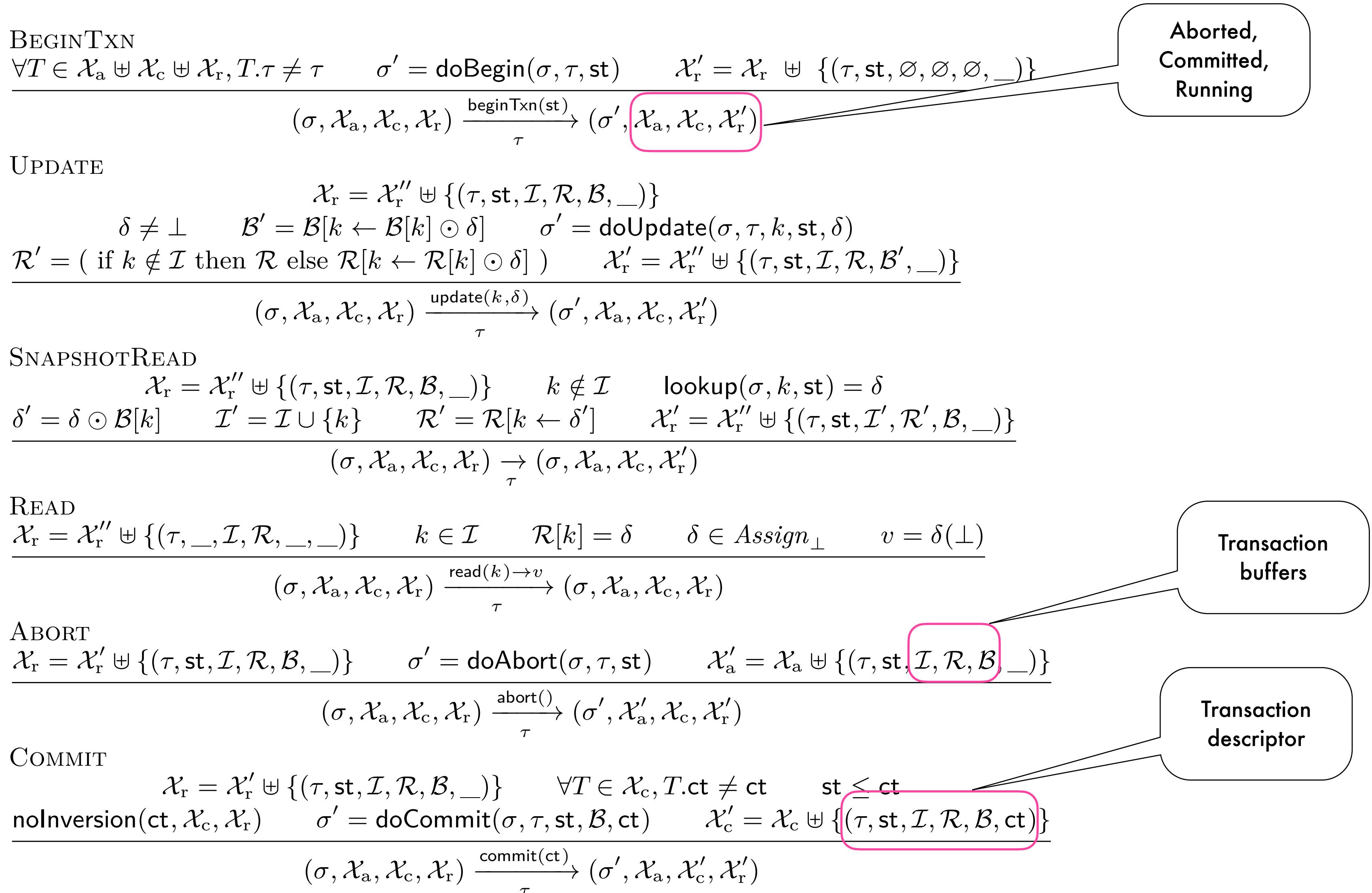
Transaction Semantics



Transaction Semantics



Transaction Semantics



Transaction Semantics - BeginTxn

$$\frac{\text{BEGINTXN} \quad \begin{array}{c} \forall T \in \mathcal{X}_a \uplus \mathcal{X}_c \uplus \mathcal{X}_r, T.\tau \neq \tau \\ \sigma' = \text{doBegin}(\sigma, \tau, \text{st}) \\ \mathcal{X}'_r = \mathcal{X}_r \uplus \{(\tau, \text{st}, \emptyset, \emptyset, \emptyset, _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{beginTxn(st)}} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

Pre Conditions:

- 1. Transaction not already running, aborted or committed

State Changes:

- 1. Add do begin record to the store
- 2. Add transaction to the set of running transactions

Transaction Semantics - Commit

COMMIT

$$\frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, st, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \forall T \in \mathcal{X}_c, T.ct \neq ct \quad st \leq ct \\ noInversion(ct, \mathcal{X}_c, \mathcal{X}_r) \quad \sigma' = doCommit(\sigma, \tau, st, \mathcal{B}, ct) \quad \mathcal{X}'_c = \mathcal{X}_c \uplus \{(\tau, st, \mathcal{I}, \mathcal{R}, \mathcal{B}, ct)\} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{commit(ct)} (\sigma', \mathcal{X}_a, \mathcal{X}'_c, \mathcal{X}'_r)}$$

Pre Conditions for Update:

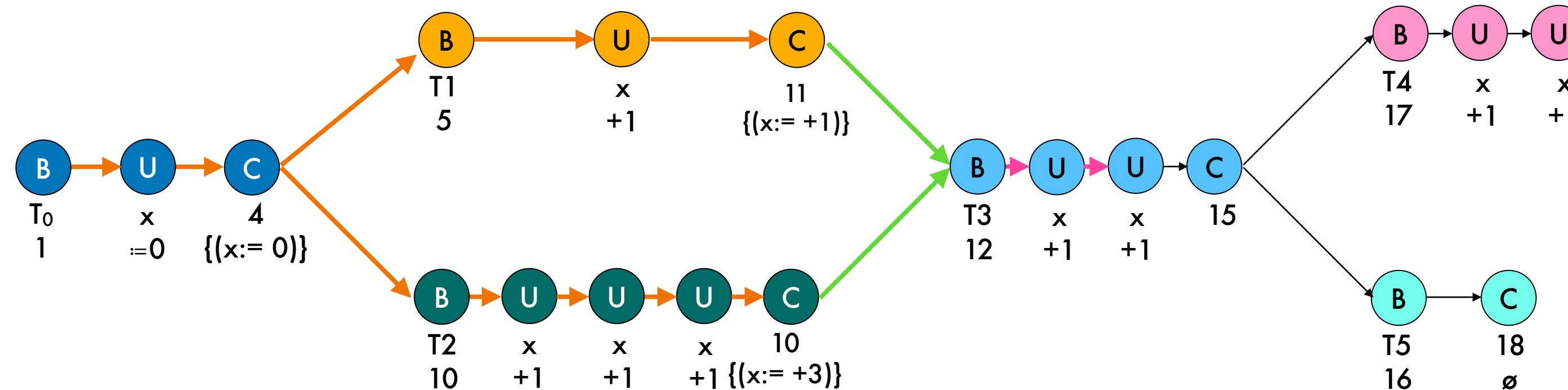
- 1. Transaction should be running
- 2. Commit timestamp should be unique
- 3. Snapshot timestamp \leq commit timestamp
- 4. Dependency inversion should not occur

State Changes:

- 1. Add commit record to the store
- 2. Update buffers in the Txn descriptor

Execution traces and Consistency

Traces - Correctness?



Correctness Criterion [Cerone et. al.] :

- Trace is Transactionally Causally Consistent (TCC)
- TCC = Ext \wedge Int \wedge TransVis

→ **Ext:** Snapshot Reads: First reads contains all previous transactions.

- Read rule: First read is a snapshot

→ **Int:** Read-My-Writes

- Read rule: later reads apply local effects

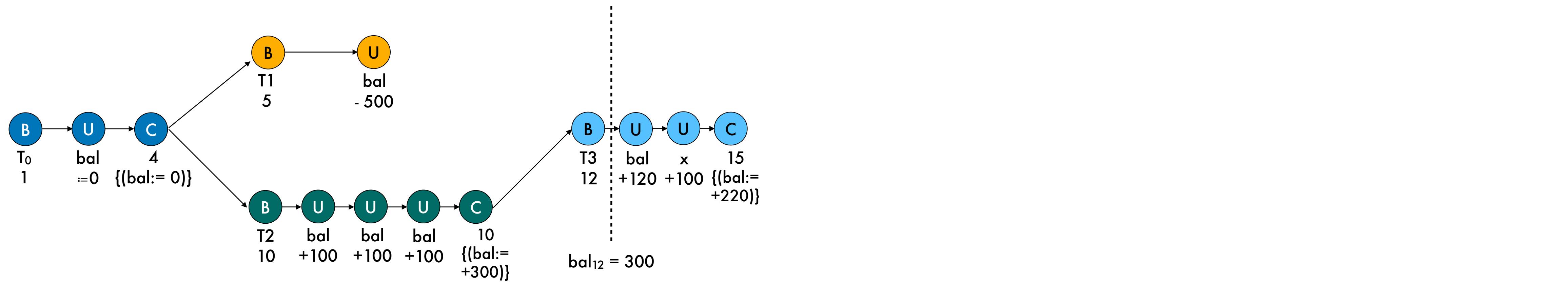
→ **TransVis:** Visibility is transitive

- By construction of a well formed trace

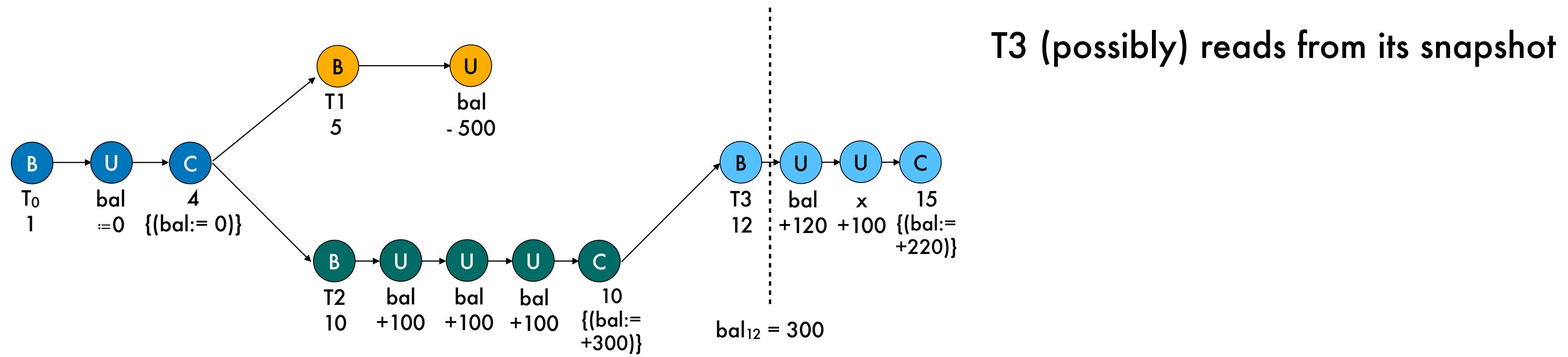
$$\frac{\text{SNAPSHOTREAD} \quad \mathcal{X}_r = \mathcal{X}_r'' \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad k \notin \mathcal{I} \quad \text{lookup}(\sigma, k, \text{st}) = \delta}{\delta' = \delta \odot \mathcal{B}[k] \quad \mathcal{I}' = \mathcal{I} \cup \{k\} \quad \mathcal{R}' = \mathcal{R}[k \leftarrow \delta'] \quad \mathcal{X}_r' = \mathcal{X}_r'' \uplus \{(\tau, \text{st}, \mathcal{I}', \mathcal{R}', \mathcal{B}, _) \}} \quad (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r')}$$

$$\frac{\text{READ} \quad \mathcal{X}_r = \mathcal{X}_r'' \uplus \{(\tau, _, \mathcal{I}, \mathcal{R}, _, _) \} \quad k \in \mathcal{I} \quad \mathcal{R}[k] = \delta \quad \delta \in \text{Assign}_{\perp} \quad v = \delta(\perp)}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{read}(k) \rightarrow v} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r)}$$

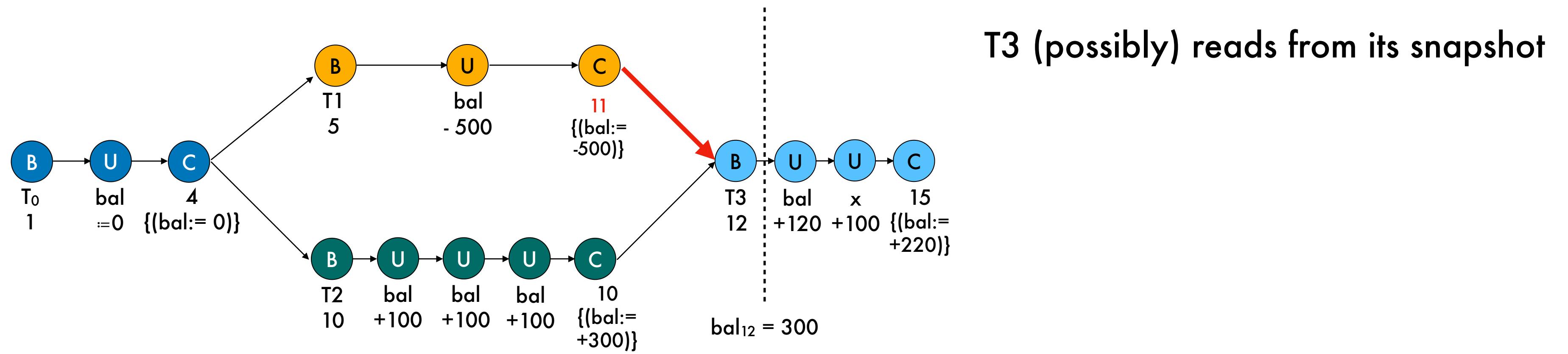
Ext + Int + TransViz \neq TCC



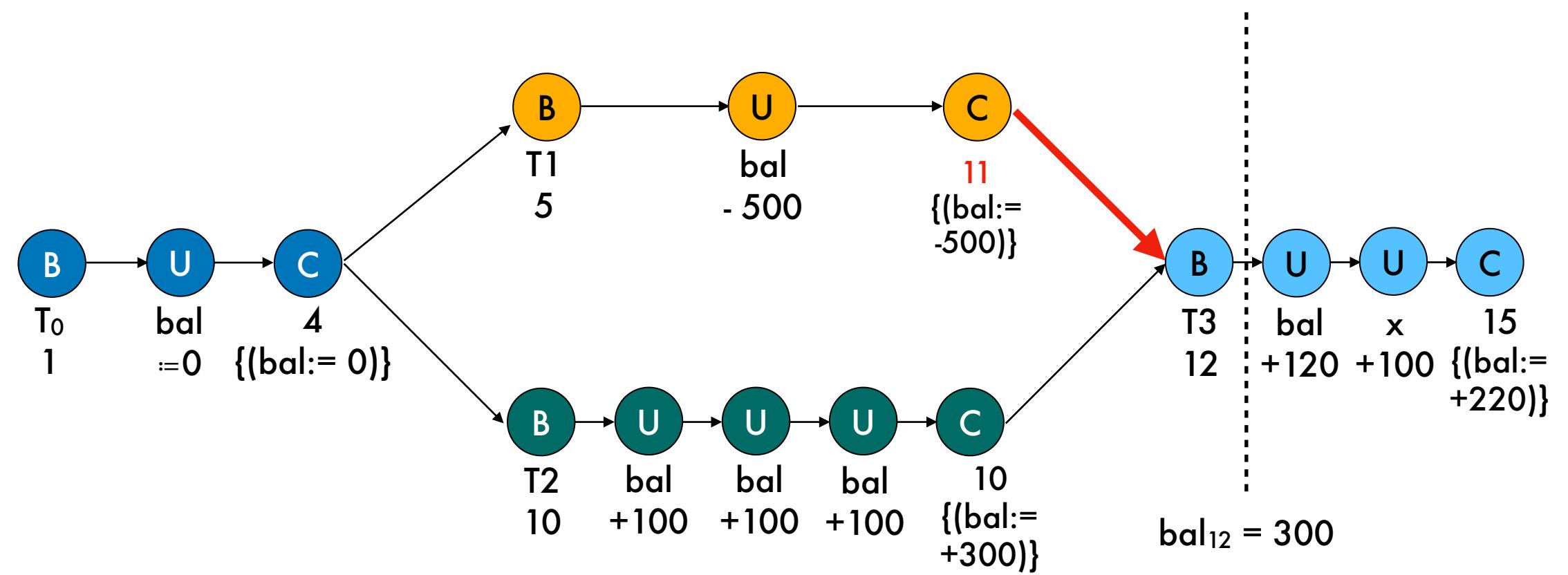
Ext + Int + TransViz \neq TCC



Ext + Int + TransViz \neq TCC



Ext + Int + TransViz \neq TCC

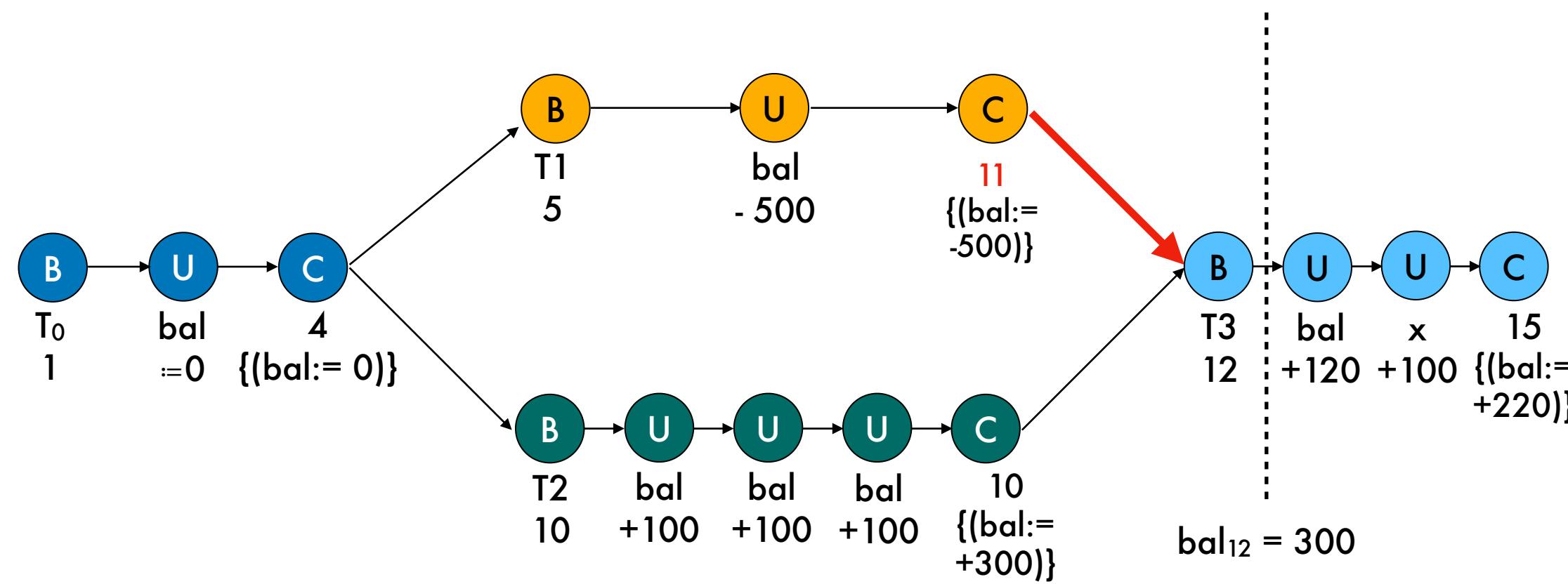


T3 (possibly) reads from its snapshot

T1 becomes visible to T3 and changes its causal past: **anomaly!**

The balance i read has changed, Money is lost!!

Ext + Int + TransViz \neq TCC



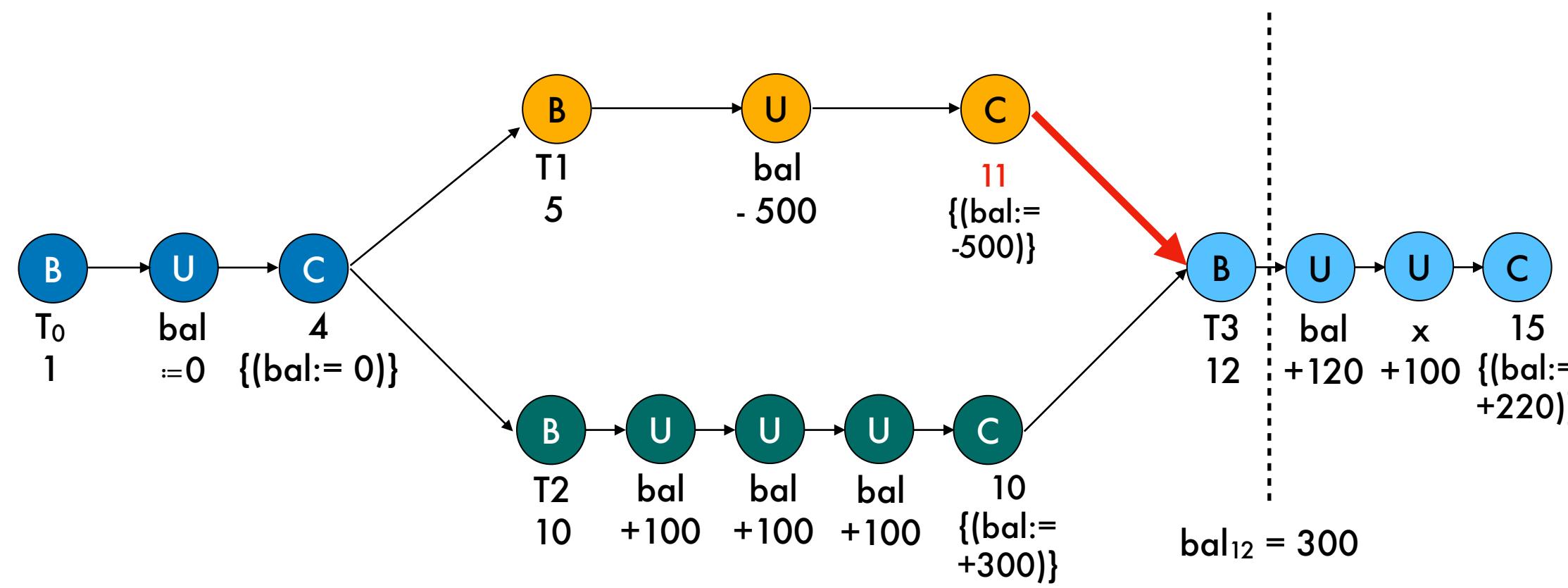
T₃ (possibly) reads from its snapshot

T₁ becomes visible to T₃ and changes its causal past: **anomaly!**

The balance i read has changed, Money is lost!!

Does not violate the wellformedness rules of a trace

Ext + Int + TransViz \neq TCC



T3 (possibly) reads from its snapshot

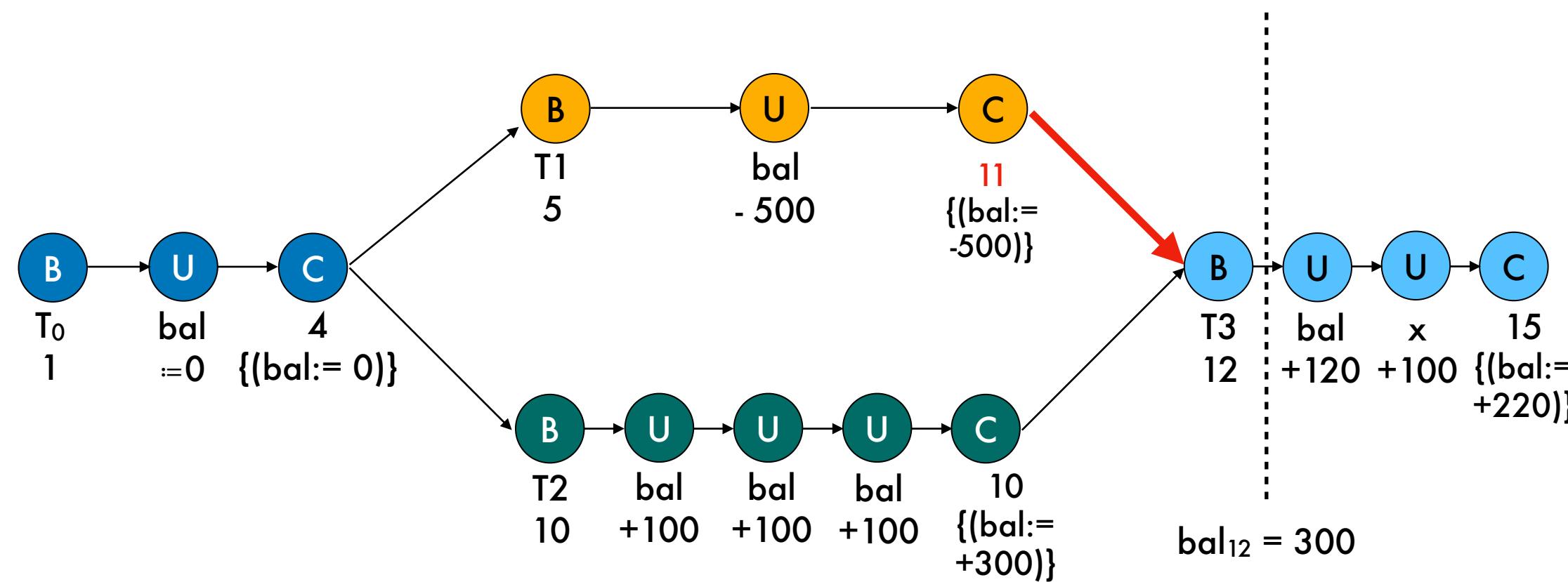
T1 becomes visible to T3 and changes its causal past: **anomaly!**

The balance i read has changed, Money is lost!!

Does not violate the wellformedness rules of a trace

Does not violate the definition of TCC (remains a DAG)

Ext + Int + TransViz \neq TCC



T3 (possibly) reads from its snapshot

T1 becomes visible to T3 and changes its causal past: **anomaly!**

The balance i read has changed, Money is lost!!

Does not violate the wellformedness rules of a trace

Does not violate the definition of TCC (remains a DAG)

Definition assumes correctness on a history

In practice, execution traces are dynamic

How can we guarantee that a visible prefix is immutable?

No Inversion Condition

Prevents transactions from choosing a commit timestamp which change causal past.

When adding a commit node, we ensure the following:

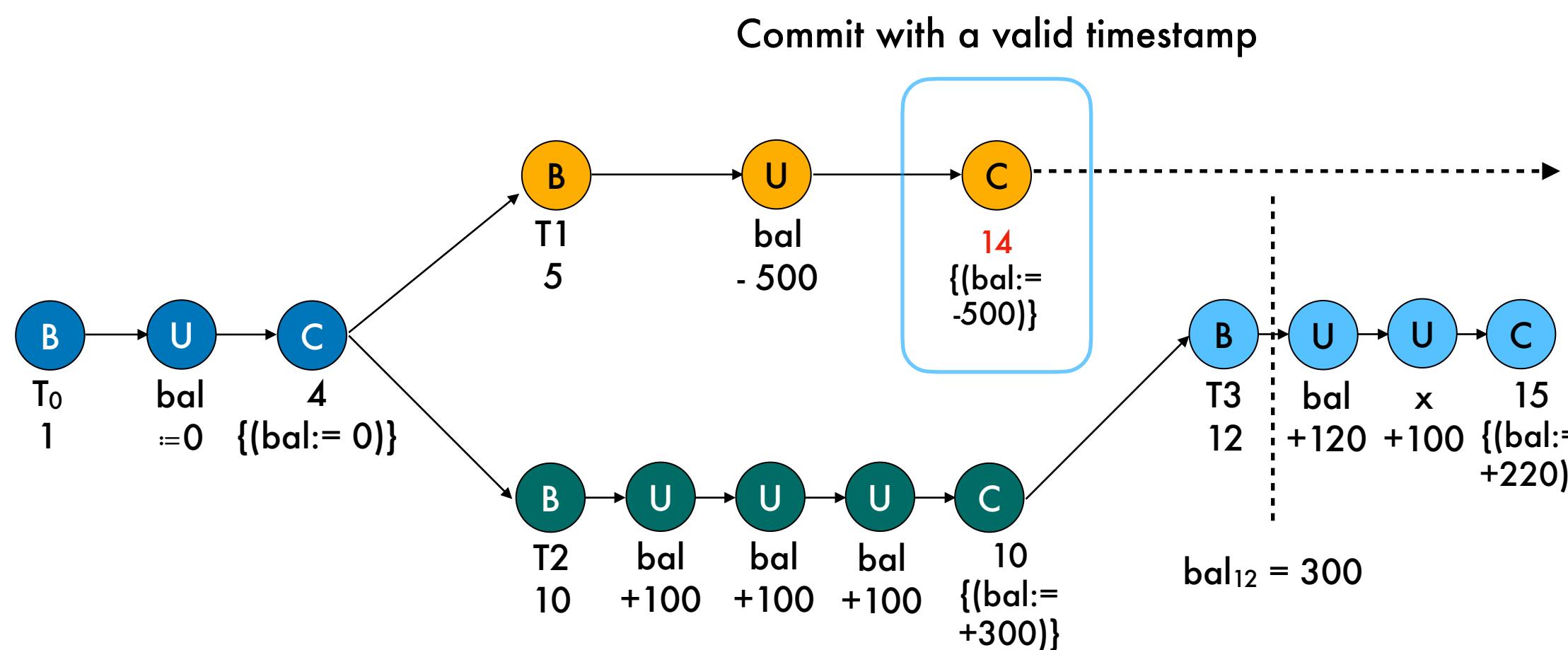
$$noInversion(n) \triangleq \forall B : n.ct \not\leq B.st$$

No Inversion Condition

Prevents transactions from choosing a commit timestamp which change causal past.

When adding a commit node, we ensure the following:

$$noInversion(n) \triangleq \forall B : n.ct \not\leq B.st$$



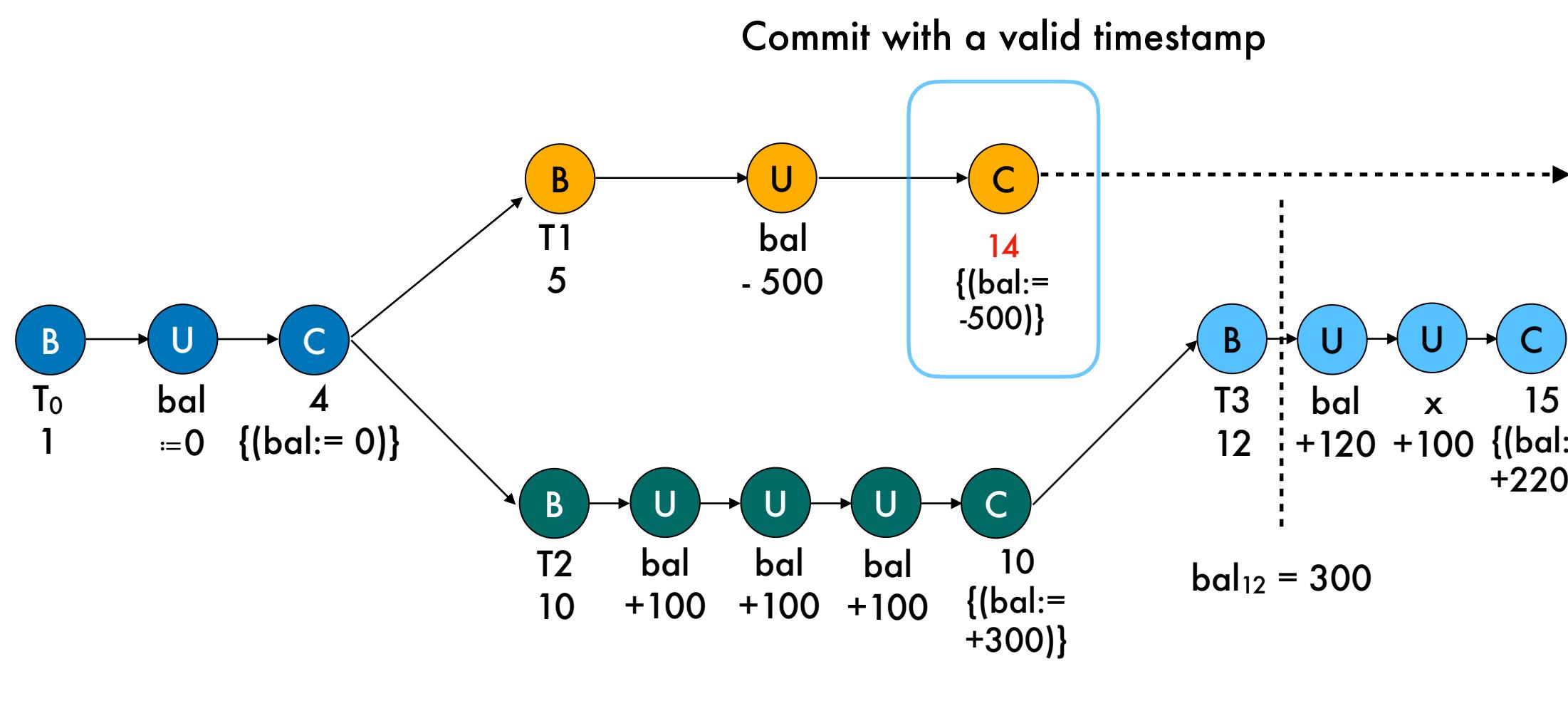
Solution 1

No Inversion Condition

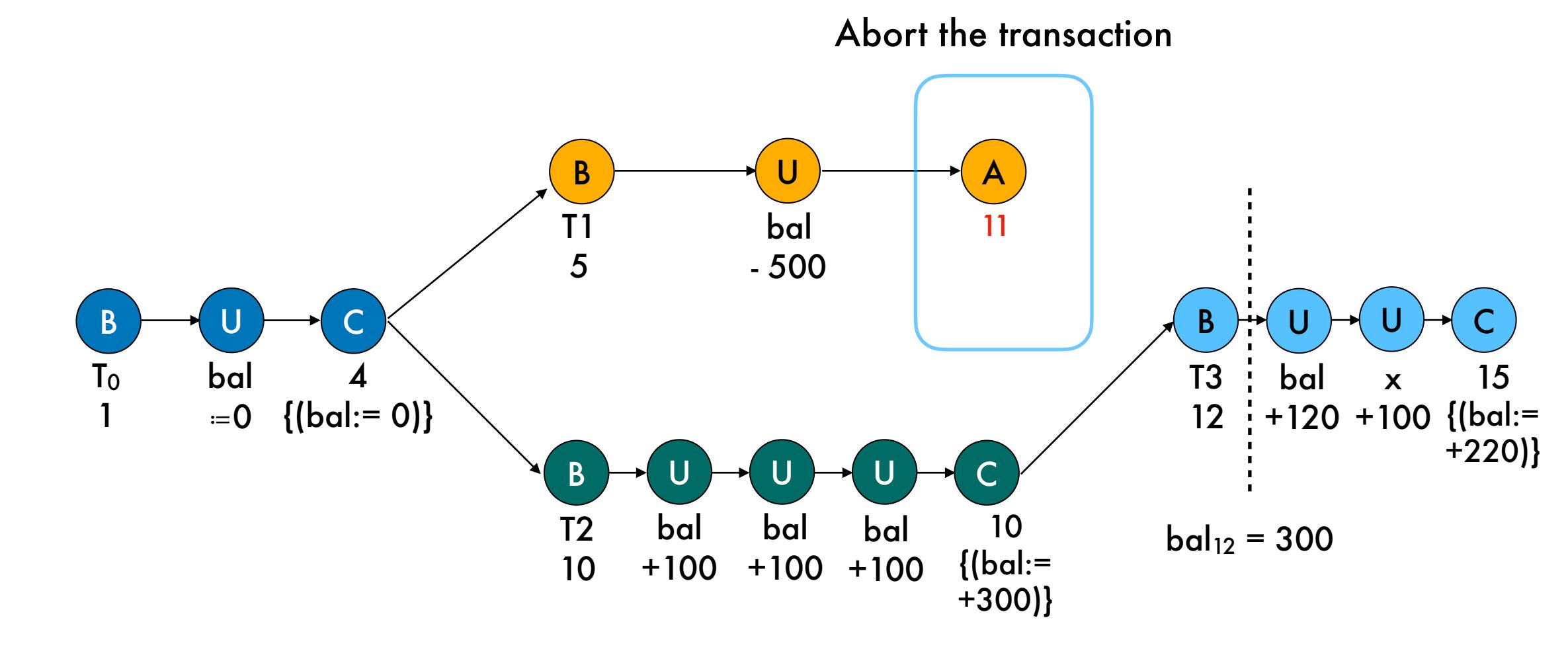
Prevents transactions from choosing a commit timestamp which change causal past.

When adding a commit node, we ensure the following:

$$\text{noInversion}(n) \triangleq \forall B : n.ct \not\leq B.st$$



Solution 1



Solution 2

Correctness

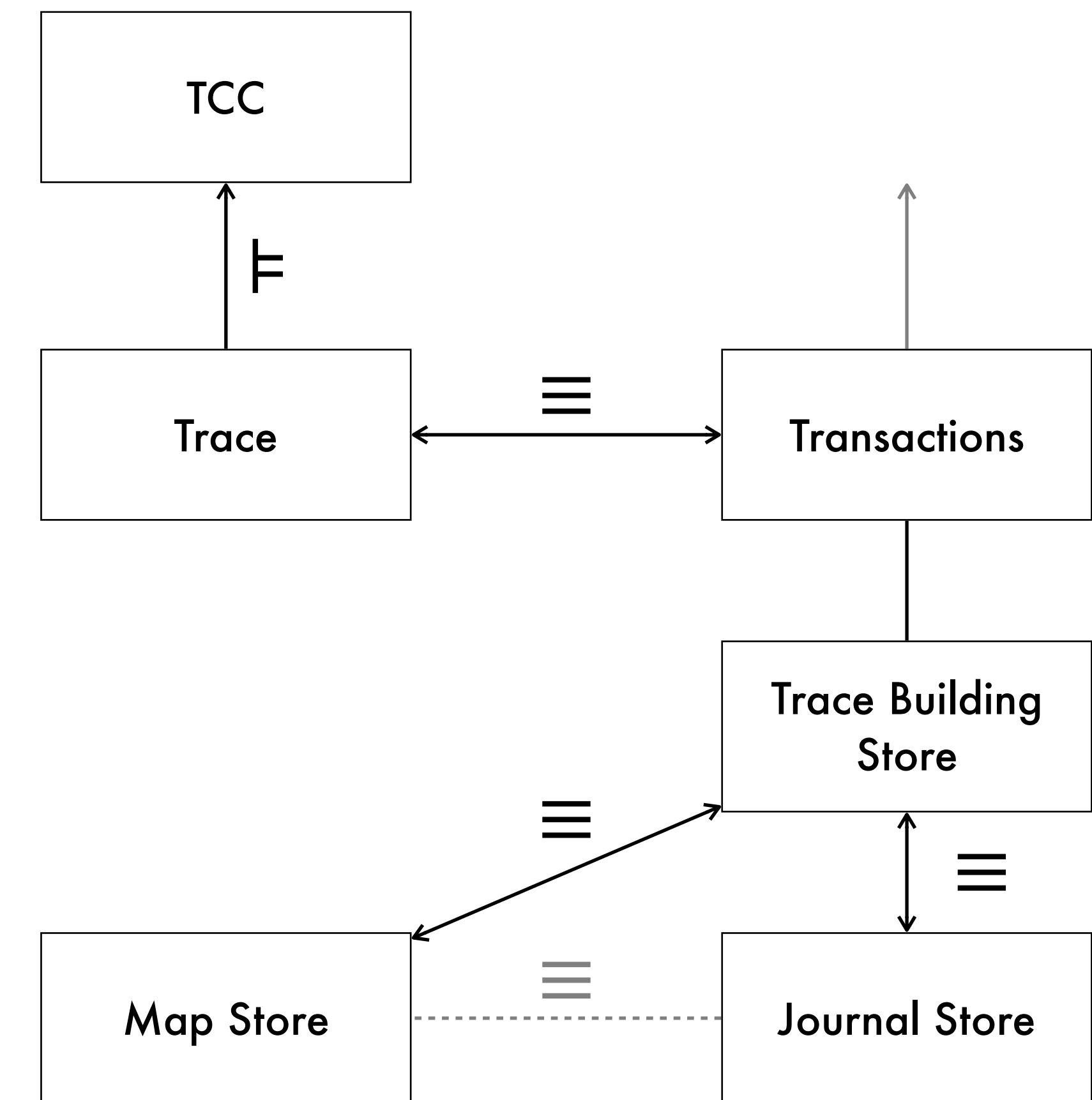
Correctness Proofs

Correct Transactions = TCC + NoInversion

1. Show that Transaction Semantics bisimulates Trace Semantics
2. Show that Trace semantics \models TCC + NoInversion
3. Using 1, Transaction semantics \models TCC + No inversion

Store Equivalence = Map and Journal are interchangeable

1. Valuation of Trace \models Lookup in Map
2. Valuation of Trace \models Lookup in Journal
3. Since valuation is the same, we get observational equivalence.



Proof Results:

1. **Transaction semantics preserve TCC + NoInversion**
2. **Journal Store and Map store are equivalent under transaction semantics**

Bisimulation example - BeginTxn

Trace semantics	Transaction semantics
Preconditions: <ul style="list-style-type: none">• τ unique• $st \neq \perp$• $p \in parents_G(n) \implies p \in \mathbb{C} \wedge p \in visited$	Preconditions: <ul style="list-style-type: none">• τ unique• $st \neq \perp$• $\{(\tau, st, _, _, _, _) \} \notin \mathcal{X}_r$ $\wedge \nexists T(\tau', st', _, _, _, ct')$: $st \geq ct' \wedge T \notin \mathcal{X}_c$
State Change: <ul style="list-style-type: none">• $\sigma' = \text{doBegin}(\sigma, \tau, st)$• $visited' = visited \cup \{\mathbb{B}(\tau, st)\}$	State Change: <ul style="list-style-type: none">• $\sigma' = \text{beginTxn}(st)$• $\mathcal{X}'_r = \mathcal{X}_r \cup \{(\tau, st, _, _, _, _) \}$
Postconditions: <ul style="list-style-type: none">• $\{\mathbb{B}(\tau, st)\} \in visited'$• $B \in \sigma'$	Postconditions: <ul style="list-style-type: none">• $\{(\tau, st, _, _, _, _) \} \in \mathcal{X}'_r$• $\text{beginTxn} \in \sigma'$
Enables: <ul style="list-style-type: none">• UPDATE• COMMIT• ABORT	Enables: <ul style="list-style-type: none">• UPDATE• COMMIT• ABORT

From Model to Implementation

Specification as Pseudocode

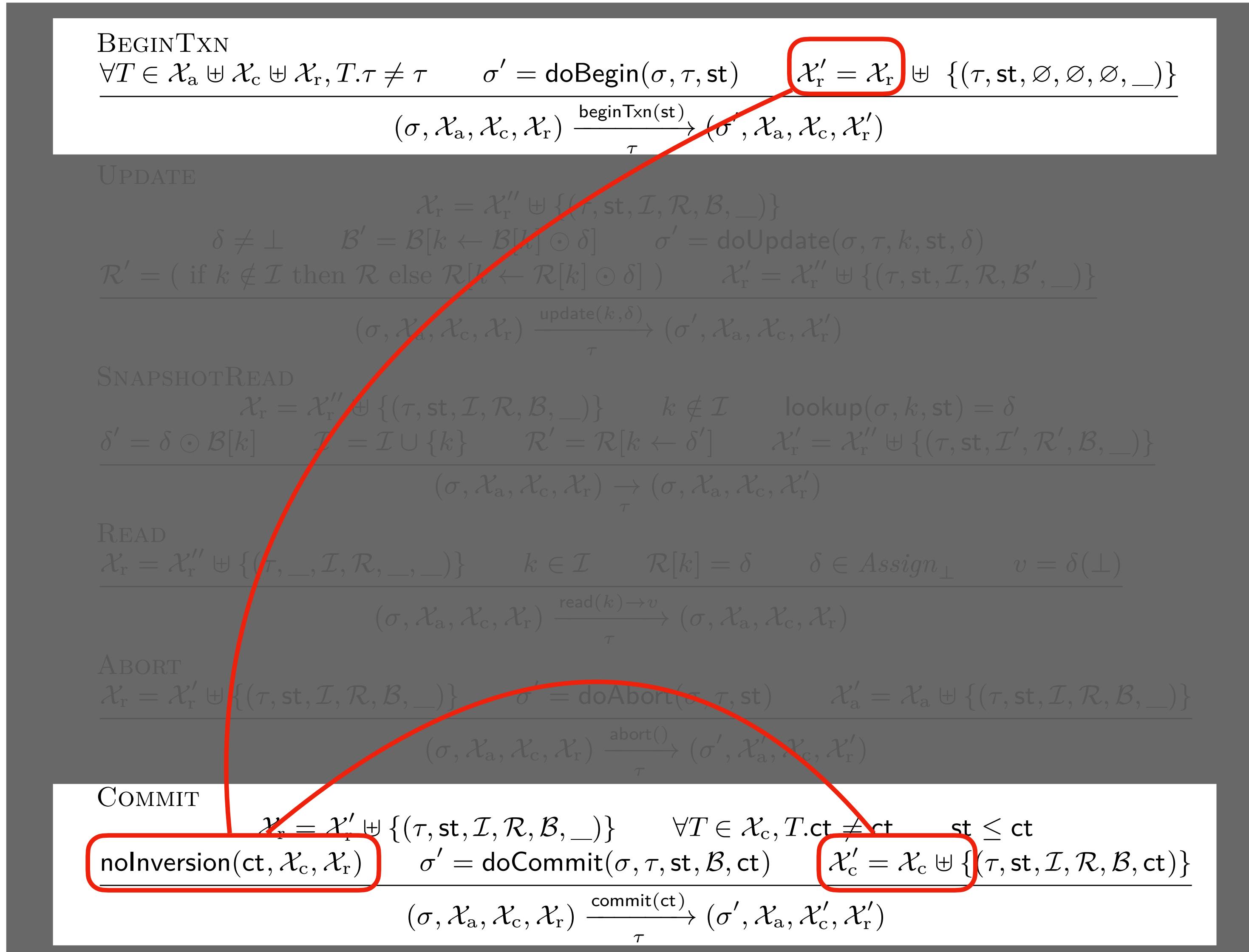
$$\begin{array}{c}
 \text{BEGINTXN} \\
 \frac{\forall T \in \mathcal{X}_a \uplus \mathcal{X}_c \uplus \mathcal{X}_r, T.\tau \neq \tau \quad \sigma' = \text{doBegin}(\sigma, \tau, \text{st}) \quad \mathcal{X}'_r = \mathcal{X}_r \uplus \{(\tau, \text{st}, \emptyset, \emptyset, \emptyset, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{beginTxn(st)}} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{UPDATE} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \delta \neq \perp \quad \mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta] \quad \sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta) \quad \mathcal{R}' = (\text{if } k \notin \mathcal{I} \text{ then } \mathcal{R} \text{ else } \mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]) \quad \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}', _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{update}(k, \delta)} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{SNAPSHOTREAD} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad k \notin \mathcal{I} \quad \text{lookup}(\sigma, k, \text{st}) = \delta \quad \delta' = \delta \odot \mathcal{B}[k] \quad \mathcal{I}' = \mathcal{I} \cup \{k\} \quad \mathcal{R}' = \mathcal{R}[k \leftarrow \delta'] \quad \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}', \mathcal{R}', \mathcal{B}, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{READ} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, _, \mathcal{I}, \mathcal{R}, _, _) \} \quad k \in \mathcal{I} \quad \mathcal{R}[k] = \delta \quad \delta \in \text{Assign}_{\perp} \quad v = \delta(\perp)}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{read}(k) \rightarrow v} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r)}
 \\[10pt]
 \text{ABORT} \\
 \frac{\mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \sigma' = \text{doAbort}(\sigma, \tau, \text{st}) \quad \mathcal{X}'_a = \mathcal{X}_a \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{abort}()} (\sigma', \mathcal{X}'_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{COMMIT} \\
 \frac{\mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \forall T \in \mathcal{X}_c, T.\text{ct} \neq \text{ct} \quad \text{st} \leq \text{ct} \quad \text{noInversion}(\text{ct}, \mathcal{X}_c, \mathcal{X}_r) \quad \sigma' = \text{doCommit}(\sigma, \tau, \text{st}, \mathcal{B}, \text{ct}) \quad \mathcal{X}'_c = \mathcal{X}_c \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, \text{ct}) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{commit(ct)}} (\sigma', \mathcal{X}_a, \mathcal{X}'_c, \mathcal{X}'_r)}
 \end{array}$$

Specification as Pseudocode

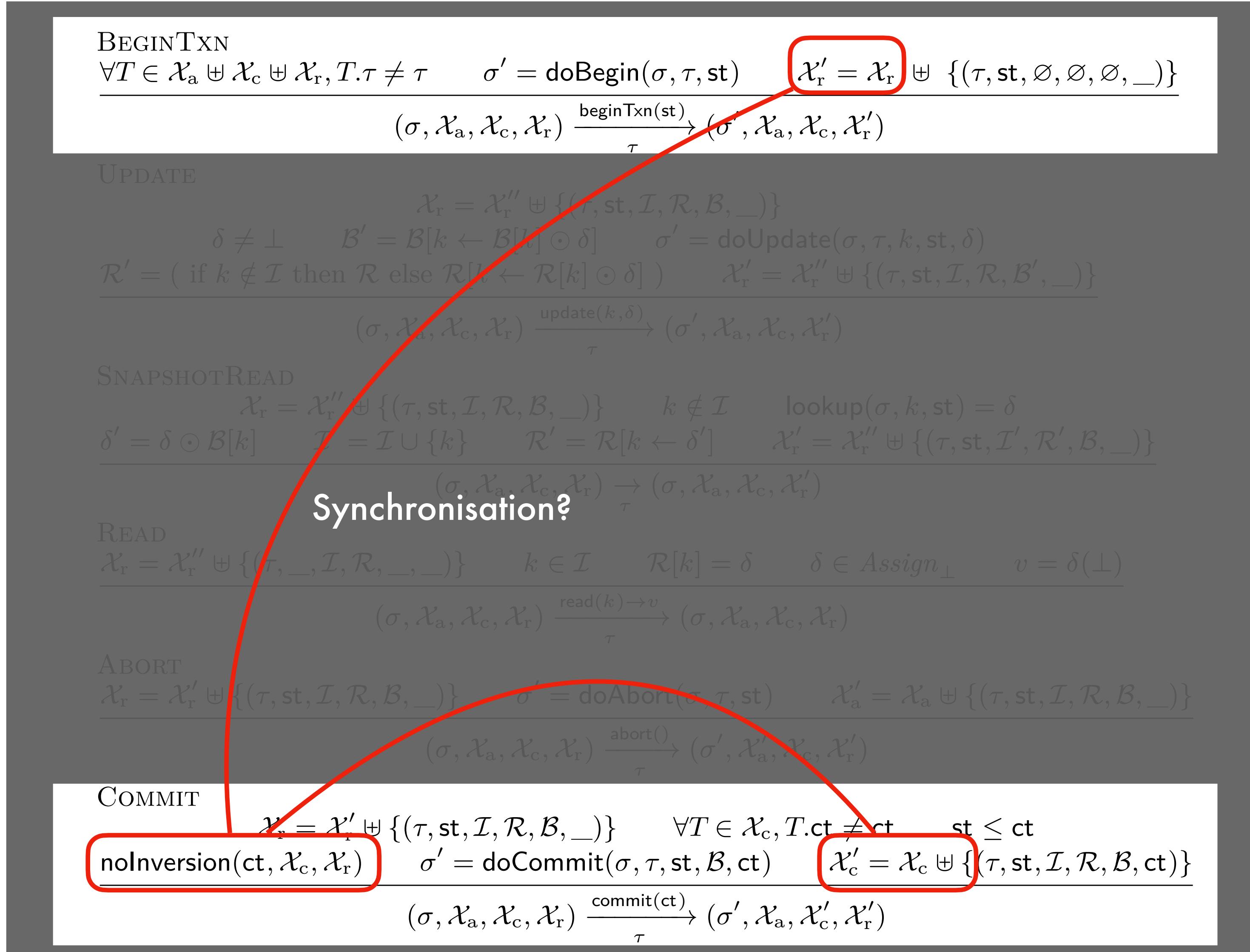
$$\begin{array}{c}
 \text{BEGINTXN} \\
 \frac{\forall T \in \mathcal{X}_a \uplus \mathcal{X}_c \uplus \mathcal{X}_r, T.\tau \neq \tau \quad \sigma' = \text{doBegin}(\sigma, \tau, \text{st}) \quad \mathcal{X}'_r = \mathcal{X}_r \uplus \{(\tau, \text{st}, \emptyset, \emptyset, \emptyset, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{beginTxn(st)}} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{UPDATE} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \delta \neq \perp \quad \mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta] \quad \sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta) \quad \mathcal{R}' = (\text{if } k \notin \mathcal{I} \text{ then } \mathcal{R} \text{ else } \mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]) \quad \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}', _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{update}(k, \delta)} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{SNAPSHOTREAD} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad k \notin \mathcal{I} \quad \text{lookup}(\sigma, k, \text{st}) = \delta \quad \delta' = \delta \odot \mathcal{B}[k] \quad \mathcal{I}' = \mathcal{I} \cup \{k\} \quad \mathcal{R}' = \mathcal{R}[k \leftarrow \delta'] \quad \mathcal{X}'_r = \mathcal{X}''_r \uplus \{(\tau, \text{st}, \mathcal{I}', \mathcal{R}', \mathcal{B}, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{READ} \\
 \frac{\mathcal{X}_r = \mathcal{X}''_r \uplus \{(\tau, _, \mathcal{I}, \mathcal{R}, _, _) \} \quad k \in \mathcal{I} \quad \mathcal{R}[k] = \delta \quad \delta \in \text{Assign}_{\perp} \quad v = \delta(\perp)}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{read}(k) \rightarrow v} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r)}
 \\[10pt]
 \text{ABORT} \\
 \frac{\mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \sigma' = \text{doAbort}(\sigma, \tau, \text{st}) \quad \mathcal{X}'_a = \mathcal{X}_a \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{abort}()} (\sigma', \mathcal{X}'_a, \mathcal{X}_c, \mathcal{X}'_r)}
 \\[10pt]
 \text{COMMIT} \\
 \boxed{\frac{\mathcal{X}_r = \mathcal{X}'_r \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad \forall T \in \mathcal{X}_c, T.\text{ct} \neq \text{ct} \quad \text{st} \leq \text{ct} \quad \text{noInversion}(\text{ct}, \mathcal{X}_c, \mathcal{X}_r) \quad \sigma' = \text{doCommit}(\sigma, \tau, \text{st}, \mathcal{B}, \text{ct}) \quad \mathcal{X}'_c = \mathcal{X}_c \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, \text{ct}) \}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{commit(ct)}} (\sigma', \mathcal{X}_a, \mathcal{X}'_c, \mathcal{X}'_r)}}
 \end{array}$$

Atomic / all-or-nothing

Specification as Pseudocode



Specification as Pseudocode



Implementing the Specification

Academic Importance:

Faithful translation of specification into code => Map and Journal Stores + Composition logic

- A lot of asserts and tests with validation and fault injection

Real World Use:

LSMTree: Emulating RocksDB => CobbleDB

- Compose multiple Map and Journal stores to produce caches, checkpoints, memtables etc.
- 4K LoC total of Java << RocksDB ≈ 300 KLoC
- Confidence in Correctness
- Sub-par Performance

Lessons Learned

High and steep mountain!

Specification is useful for system implementors:

- Compact, aids reasoning; no forgetting edge cases
- Read as pseudocode
- Identify concurrency, crash-tolerance issues

New: rigorous store composition

- Explanatory value
- Ease of implementation



Contributions

Formalised:

- Rich Datatypes: Incremental, Convergent
- Transactional trace, valuation
- Store behaviour: map, journal
- Safety: $\text{transactions} \cap \text{store} \rightarrow \text{valuation}$
 - Satisfied for map, journal
- Client-driven transaction semantics
 - Bisimulates traces
 - Map, Journal safe

Ongoing:

- Specification in Coq and Dafny to extract correct by design code into OCaml (functional) and Java (procedural)
- Correct-by-construction implementation via asserts and test suites
- Better performance, within model improving data structures used to implement spec.

Appendices and References

RocksDB by Composition

Live level

Jrnl₁ (immutable): [t₃, t₆]

(k ₄ , δ __)	(k ₃ , δ __)	(k ₁ , δ __)	(k ₂ , δ __)
t ₃	t ₄	t ₅	t ₆

Jrnl₂ (immutable): [t₇, t₁₀]

(k ₇ , δ __)	(k ₂ , δ __)	(k ₄ , δ __)	(k ₂ , δ __)
t ₇	t ₈	t ₉	t ₁₀

Jrnl₃ (active): [t₁₁, now]

(k ₁ , δ __)	(k ₁ , δ __)		
t ₁₁	t ₁₂	t ₁₃	t ₁₄

Memtable₁ (immutable): [t₃, t₆]

(k ₁ , δ __)	(k ₂ , δ __)	(k ₃ , δ __)	(k ₄ , δ __)
------------------------------------	------------------------------------	------------------------------------	------------------------------------

Memtable₂ (immutable): [t₇, t₁₀]

(k ₂ , δ __)	(k ₄ , δ __)	(k ₇ , δ __)
------------------------------------	------------------------------------	------------------------------------

Memtable₃ (active): [t₁₁, now]

(k ₁ , δ __)		
------------------------------------	--	--

Level 0

SST₁: [t₁, t₂])

(k ₁ , δ __)	(k ₃ , δ __)	(k ₅ , δ __)	(k ₆ , δ __)	(k ₇ , δ __)	(k ₈ , δ __)
------------------------------------	------------------------------------	------------------------------------	------------------------------------	------------------------------------	------------------------------------

SST₂: [t₂, t₃])

(k ₃ , δ __)	(k ₄ , δ __)	(k ₅ , δ __)	(k ₈ , δ __)
------------------------------------	------------------------------------	------------------------------------	------------------------------------

Level 1: [t₀ = 0, t₁)

SST₁: keys 0–10

(k ₀ , := __)	(k ₁ , := __)	(k ₂ , := __)	(k ₃ , := __)	(k ₈ , := __)	(k ₉ , := __)
-------------------------------------	-------------------------------------	-------------------------------------	-------------------------------------	-------------------------------------	-------------------------------------

Indices indicate order: t₁ older than t₂, k₁ less than k₂.
Journals and SSTs are persistent; Memtables are volatile.

SST₁: keys 11+

(k ₁₇ , := __)	(k ₁₉ , := __)	(k ₂₀ , := __)	(k ₂₁ , := __)	(k ₂₂ , := __)	(k ₇₇ , := __)
--------------------------------------	--------------------------------------	--------------------------------------	--------------------------------------	--------------------------------------	--------------------------------------

Transaction Semantics - Read

SNAPSHOTREAD

$$\begin{array}{c}
 \mathcal{X}_r = \mathcal{X}_r'' \uplus \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \quad k \notin \mathcal{I} \quad \text{lookup}(\sigma, k, \text{st}) = \delta \\
 \delta' = \delta \odot \mathcal{B}[k] \quad \mathcal{I}' = \mathcal{I} \cup \{k\} \quad \mathcal{R}' = \mathcal{R}[k \leftarrow \delta'] \quad \mathcal{X}'_r = \mathcal{X}_r'' \uplus \{(\tau, \text{st}, \mathcal{I}', \mathcal{R}', \mathcal{B}, _) \} \\
 \hline
 (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)
 \end{array}$$

READ

$$\begin{array}{c}
 \mathcal{X}_r = \mathcal{X}_r'' \uplus \{(\tau, _, \mathcal{I}, \mathcal{R}, _, _) \} \quad k \in \mathcal{I} \quad \mathcal{R}[k] = \delta \quad \delta \in \text{Assign}_\perp \quad v = \delta(\perp) \\
 \hline
 (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{read}(k) \rightarrow v} (\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r)
 \end{array}$$

Pre Conditions for Read:

- 1. Transaction should be running
- 2. The key should be initialised, otherwise perform snapshot read

State Changes: None

Pre Conditions for Snapshot read:

- 1. Transaction should be running
- 2. The key should not be initialised

State Changes:

- 1. Add key to the set of initialised keys
- 2. Add value of the key to the read buffer
- 3. Update buffers in Txn Descriptor

Transaction Semantics - Update

UPDATE

$$\frac{\begin{array}{c} \mathcal{X}_r = \mathcal{X}_r'' \uplus \{(\tau, st, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \\ \delta \neq \perp \quad \mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta] \quad \sigma' = \text{doUpdate}(\sigma, \tau, k, st, \delta) \\ \mathcal{R}' = (\text{ if } k \notin \mathcal{I} \text{ then } \mathcal{R} \text{ else } \mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]) \quad \mathcal{X}'_r = \mathcal{X}_r'' \uplus \{(\tau, st, \mathcal{I}, \mathcal{R}, \mathcal{B}', _) \} \end{array}}{(\sigma, \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}_r) \xrightarrow[\tau]{\text{update}(k, \delta)} (\sigma', \mathcal{X}_a, \mathcal{X}_c, \mathcal{X}'_r)}$$

Pre Conditions for Update:

- 1. Transaction should be running
- 2. Effect should not be null

State Changes:

- 1. Add the effect of update to write buffer
- 2. Add the effect of update to read buffer
- 3. Add update record to the store
- 4. Update buffers in the Txn descriptor

Bisimulation example - Update

UPDATE

Trace semantics	Transaction semantics
Preconditions: <ul style="list-style-type: none"> • $p \in \text{parents}_G(n)$ $\wedge p.\tau = n.\tau \wedge p.\text{st} = n.\text{st}$ $\implies (p \in \mathbb{U} \cup \mathbb{B}) \wedge p \in \text{visited}$ 	Preconditions: <ul style="list-style-type: none"> • $\{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \in \mathcal{X}_r$ • $\delta \neq \perp$
State Change: <ul style="list-style-type: none"> • $\sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta)$ • $\text{visited}' = \text{visited} \cup \{\mathbb{U}(\tau, \text{st})\}$ 	State Change: <ul style="list-style-type: none"> • $\sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta)$ • $\mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta]$ • $\mathcal{R}' = (\text{if } k \notin \mathcal{I}$ then \mathcal{R} else $\mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]$) • $\mathcal{X}'_r = \mathcal{X}_r \cup \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}', \mathcal{B}', _) \}$
Postconditions: <ul style="list-style-type: none"> • $\{\mathbb{U}(\tau, \text{st})\} \in \text{visited}'$ • $U \in \sigma'$ 	Postconditions: <ul style="list-style-type: none"> • $\{(\tau, \text{st}, \mathcal{I}, \mathcal{R}', \mathcal{B}', _) \} \in \mathcal{X}'_r$ • $\text{update} \in \sigma'$

Bisimulation example - Abort

UPDATE

Trace semantics	Transaction semantics
Preconditions: <ul style="list-style-type: none"> • $p \in \text{parents}_G(n)$ $\wedge p.\tau = n.\tau \wedge p.\text{st} = n.\text{st}$ $\implies (p \in \mathbb{U} \cup \mathbb{B}) \wedge p \in \text{visited}$ 	Preconditions: <ul style="list-style-type: none"> • $\{(\tau, \text{st}, \mathcal{I}, \mathcal{R}, \mathcal{B}, _) \} \in \mathcal{X}_r$ • $\delta \neq \perp$
State Change: <ul style="list-style-type: none"> • $\sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta)$ • $\text{visited}' = \text{visited} \cup \{\mathbb{U}(\tau, \text{st})\}$ 	State Change: <ul style="list-style-type: none"> • $\sigma' = \text{doUpdate}(\sigma, \tau, k, \text{st}, \delta)$ • $\mathcal{B}' = \mathcal{B}[k \leftarrow \mathcal{B}[k] \odot \delta]$ • $\mathcal{R}' = (\text{if } k \notin \mathcal{I}$ then \mathcal{R} else $\mathcal{R}[k \leftarrow \mathcal{R}[k] \odot \delta]$) • $\mathcal{X}'_r = \mathcal{X}_r \cup \{(\tau, \text{st}, \mathcal{I}, \mathcal{R}', \mathcal{B}', _) \}$
Postconditions: <ul style="list-style-type: none"> • $\{\mathbb{U}(\tau, \text{st})\} \in \text{visited}'$ • $U \in \sigma'$ 	Postconditions: <ul style="list-style-type: none"> • $\{(\tau, \text{st}, \mathcal{I}, \mathcal{R}', \mathcal{B}', _) \} \in \mathcal{X}'_r$ • $\text{update} \in \sigma'$