

Carnegie Mellon University

Database Systems

Multi-Version
Concurrency Control

15-445/645 SPRING 2025 » PROF. JIGNESH PATEL

ADMINISTRIVIA

Project #4 is out; due Sunday April 20th @ 11:59pm
→ Recitation: Friday, April 11th in GHC 4303 from 3:00 - 4:00 PM

Final Exam is on Monday, April 28, 2025, from
05:30pm - 08:30pm
→ Early exam will not be offered. Do not make travel plans.

No OH on April 7th (I have to teach another class)

UPCOMING DATABASE TALKS

StarRocks (DB Seminar)

- Monday, March 31 @ 4:30pm
- StarRocks Query Optimizer
- Speaker: Kaisen Kang
- <https://cmu.zoom.us/j/93441451665>



Oracle (DB Seminar)

- Tuesday, April 1 @ noon
- Evolving Transactions in Oracle 23ai: New Advances in Concurrency and Consistency
- Speakers: Kishy Kumar and Akshay S. Kulkarni
- <https://cmu.zoom.us/my/jignesh>



LAST CLASS

We discussed concurrency control protocols for generating conflict serializable schedules without needing to know what queries a txn will execute.

The two-phase locking (2PL) protocol requires txns to acquire locks on database objects before they are allowed to access them.

MULTI-VERSION CONCURRENCY CONTROL

The DBMS maintains multiple **physical** versions of a single **logical** object in the database:

- When a txn writes to an object, the DBMS creates a new version of that object.
- When a txn reads an object, it reads the newest version that existed when the txn started.

MVCC HISTORY

Protocol was first proposed in 1978
MIT PhD dissertation.

First implementations was Rdb/VMS
and InterBase at DEC in early 1980s.

- Both were by Jim Starkey, co-founder of NuoDB.
- DEC Rdb/VMS is now “Oracle Rdb”.
- InterBase was open-sourced as Firebird.



Rdb/VMS

MVCC HISTORY

Protocol was first proposed in 1978
MIT PhD dissertation.

First implementations was Rdb/VMS
and InterBase at DEC in early 1980s.

- Both were by Jim Starkey, co-founder of NuoDB.
- DEC Rdb/VMS is now “Oracle Rdb”.
- InterBase was open-sourced as Firebird.



Firebird

Rdb/VMS



Oracle Rdb
the Database for HP
OpenVMS Platform

MULTI-VERSION CONCURRENCY CONTROL

Writers do not block readers.

Readers do not block writers.

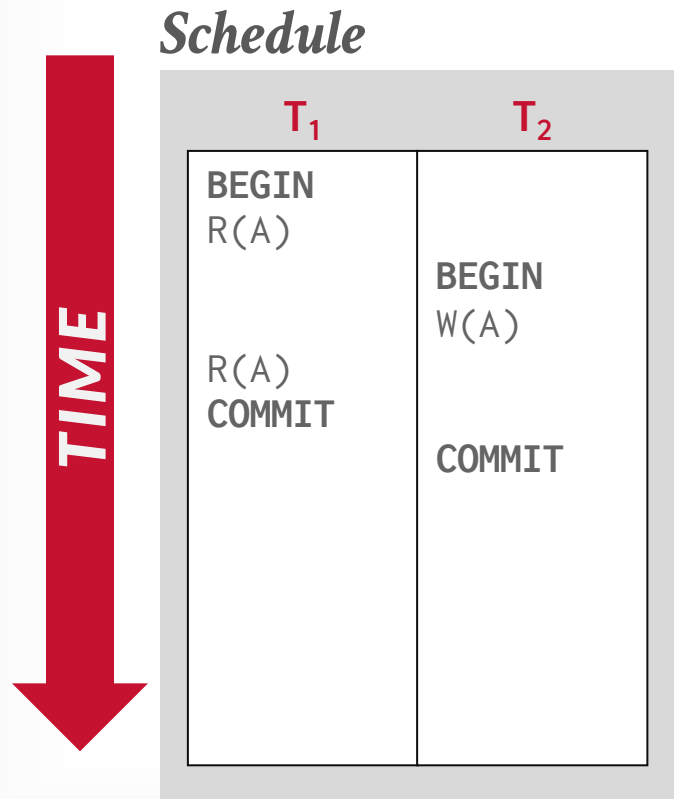
Read-only txns can read a consistent snapshot without acquiring locks.

→ Use timestamps to determine visibility.

→ MVCC naturally supports Snapshot Isolation (SI).

Multi-versioning without garbage collection allows the DBMS to support time-travel queries.

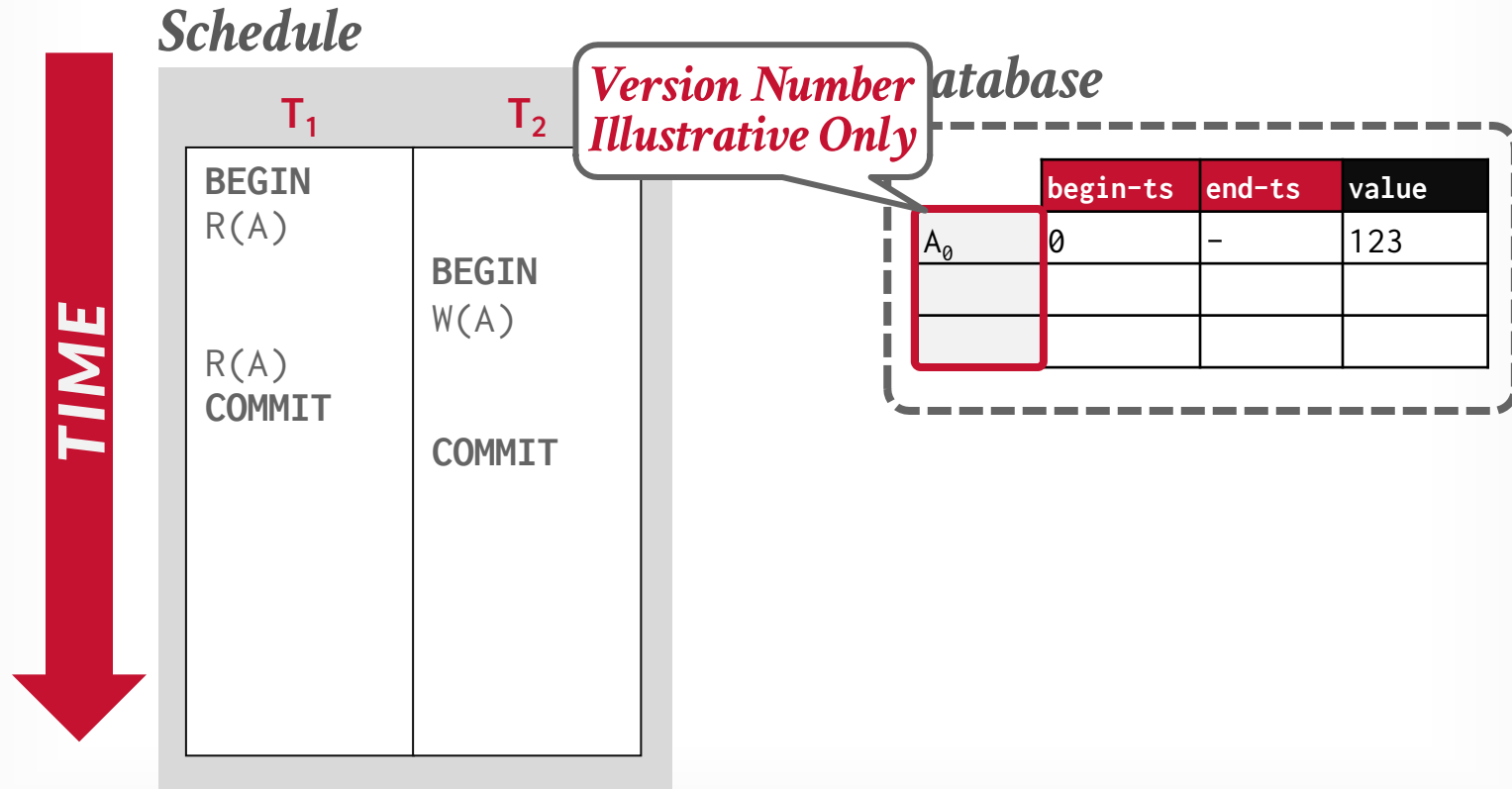
MVCC – EXAMPLE #1



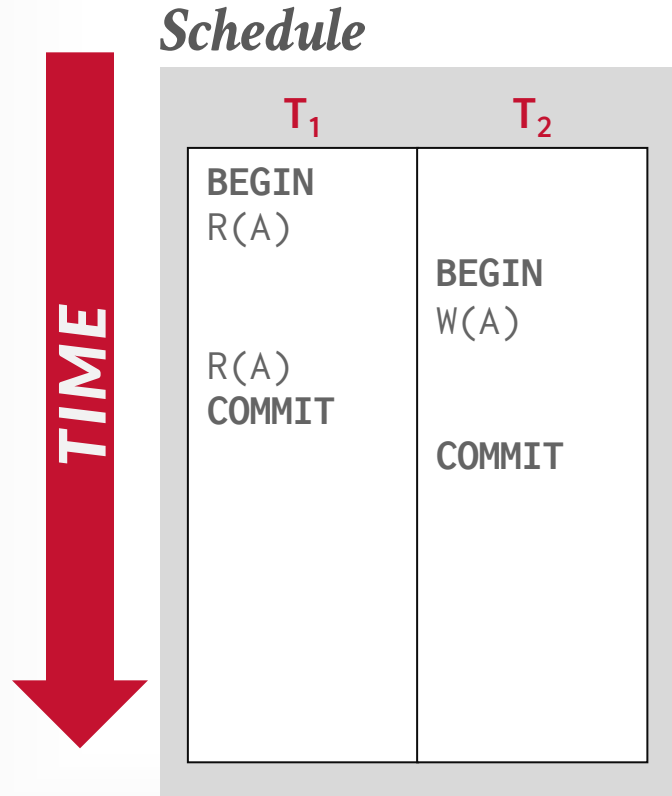
Database

	begin-ts	end-ts	value
A_0	0	-	123

MVCC – EXAMPLE #1



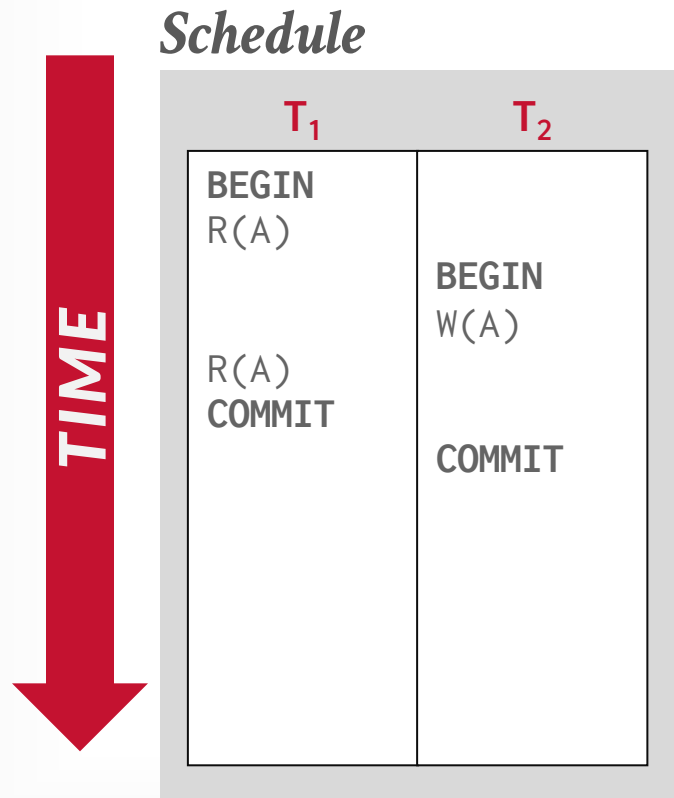
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123

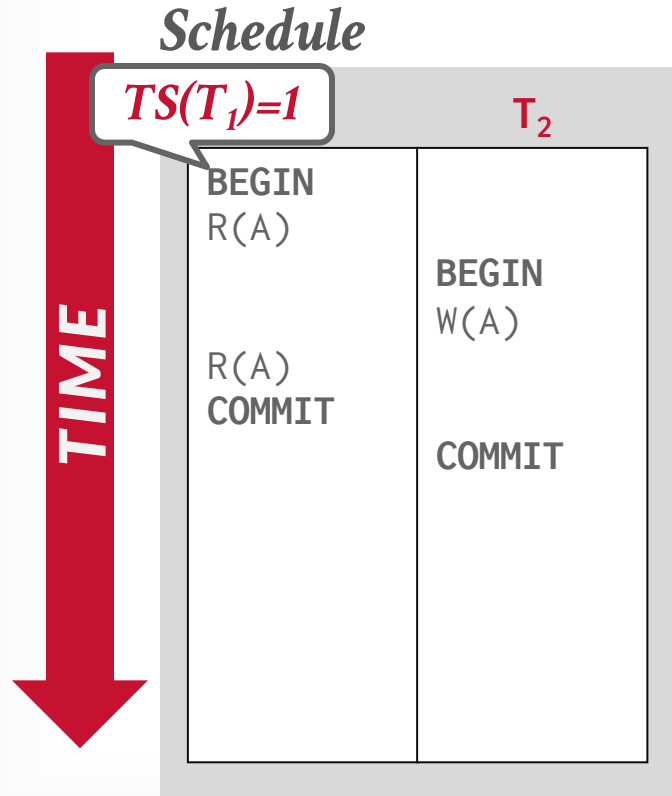
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123

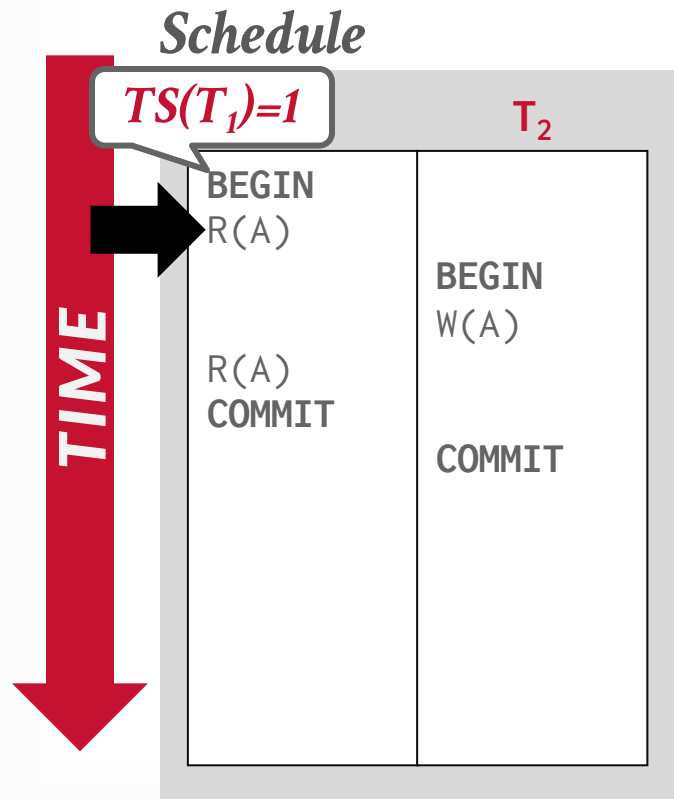
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123

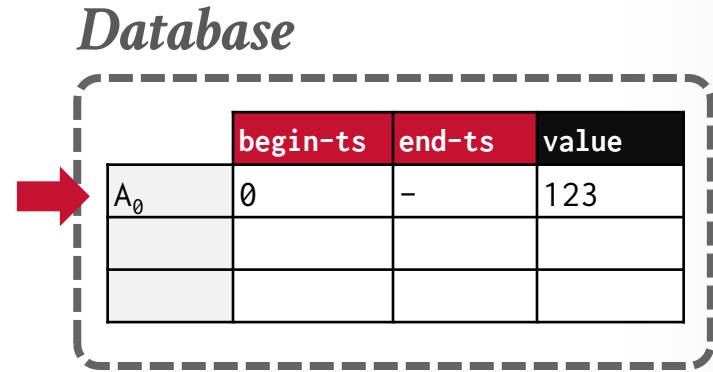
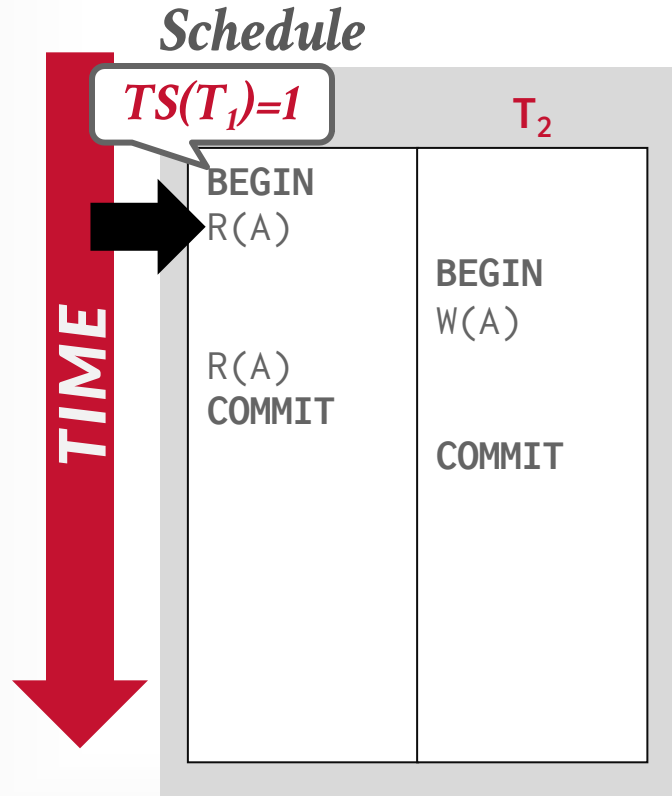
MVCC – EXAMPLE #1



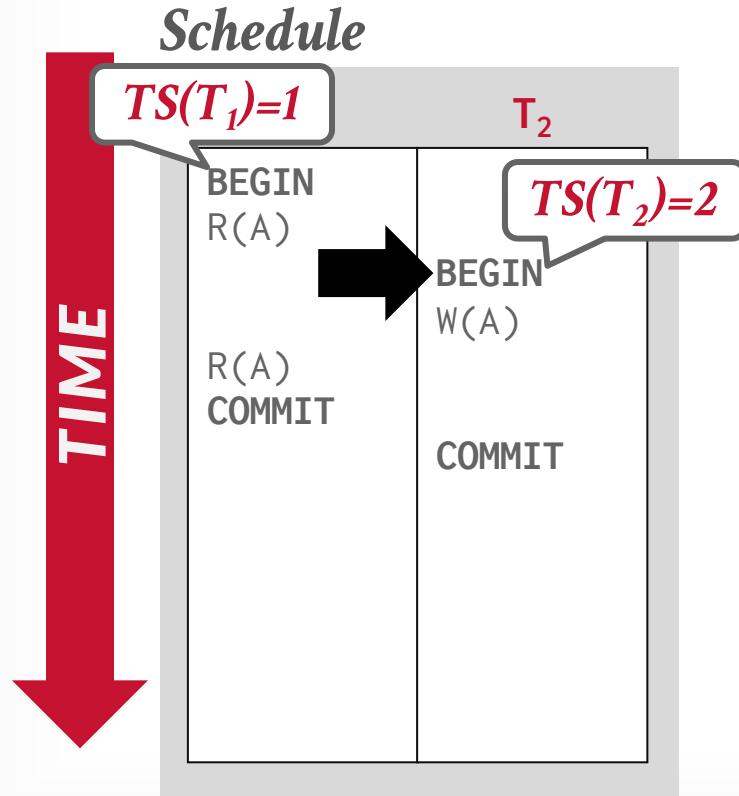
Database

	begin-ts	end-ts	value
A_0	0	-	123

MVCC – EXAMPLE #1



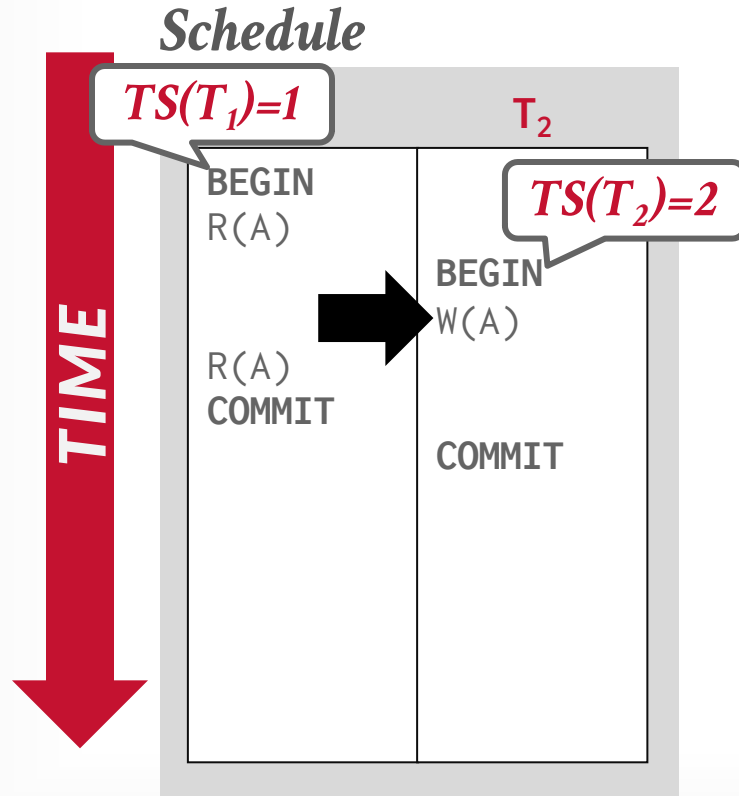
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123

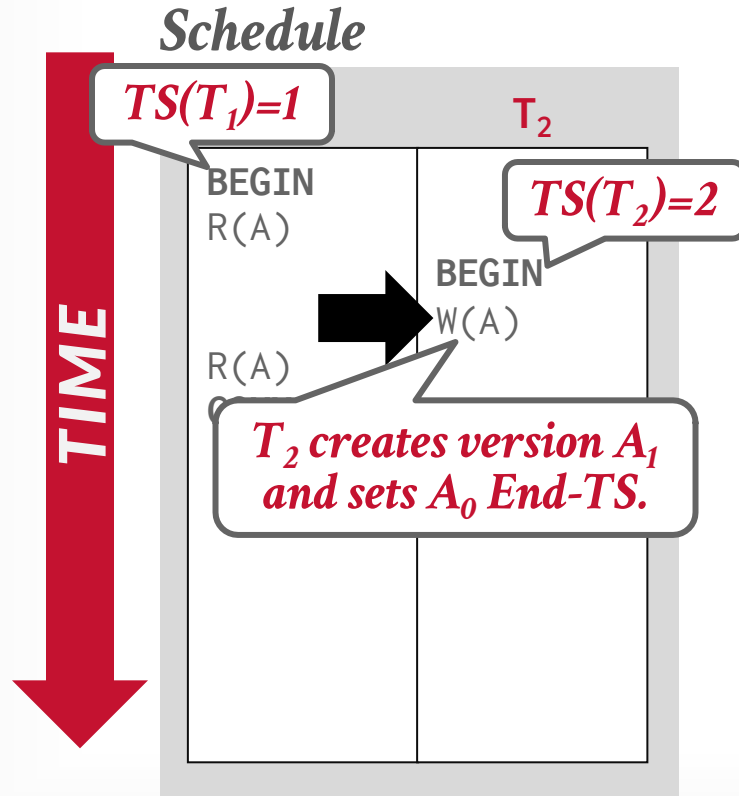
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123

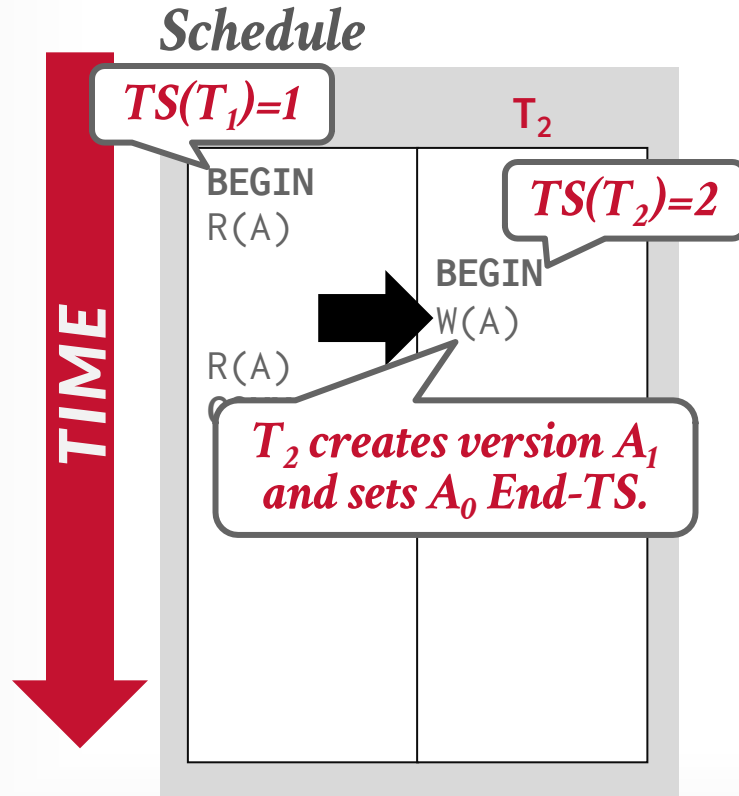
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	-	123
A_1	2	-	456

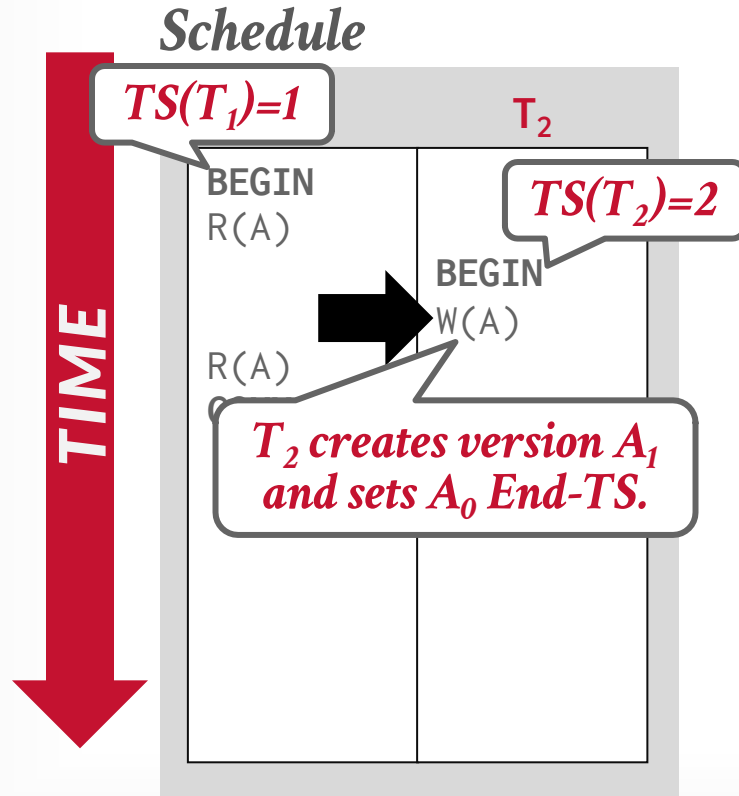
MVCC – EXAMPLE #1



Database

	begin-ts	end-ts	value
A_0	0	2	123
A_1	2	-	456

MVCC – EXAMPLE #1



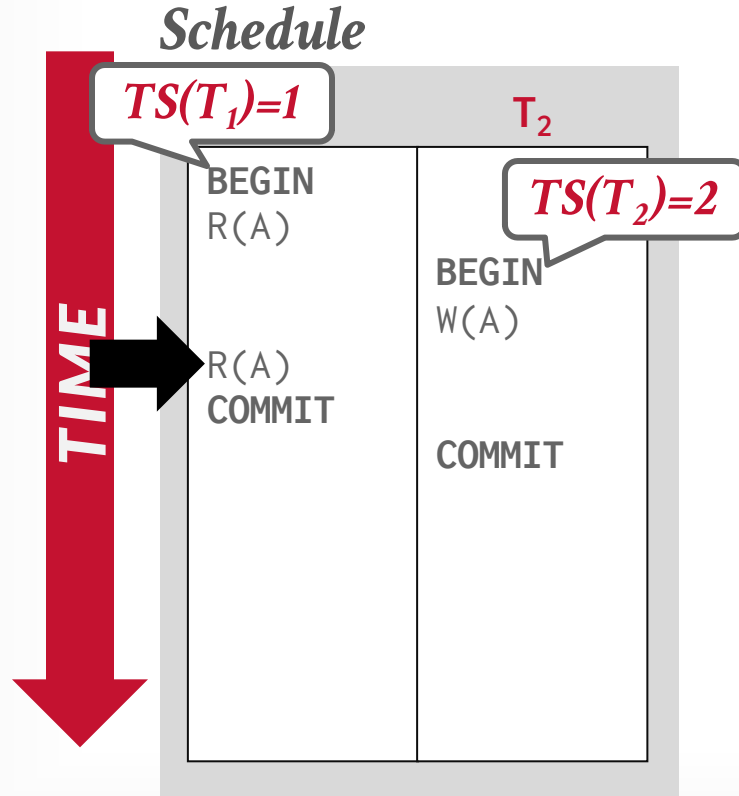
Database

	begin-ts	end-ts	value
A_0	0	2	123
A_1	2	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active
T_2	2	Active

MVCC – EXAMPLE #1



Database

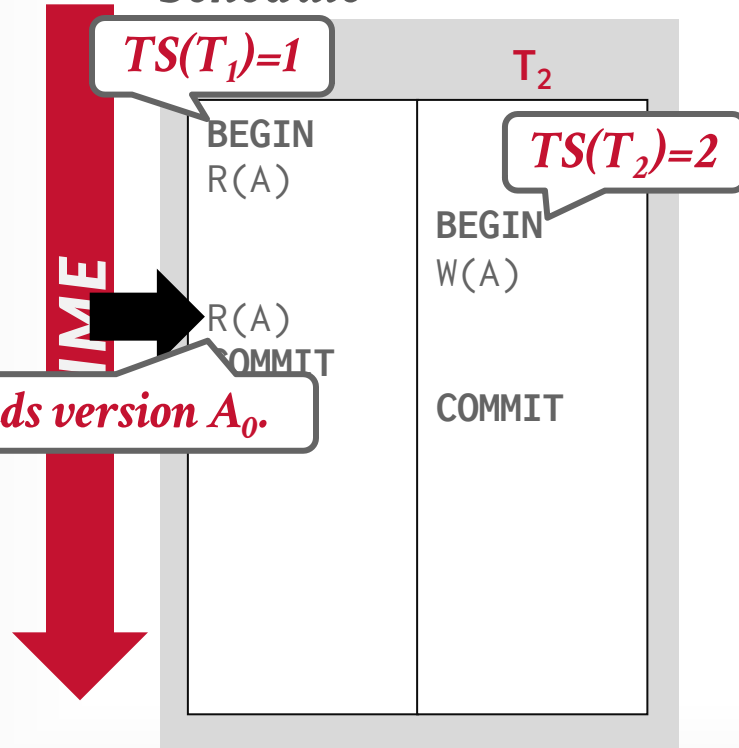
	begin-ts	end-ts	value
A_0	0	2	123
A_1	2	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active
T_2	2	Active

MVCC – EXAMPLE #1

Schedule



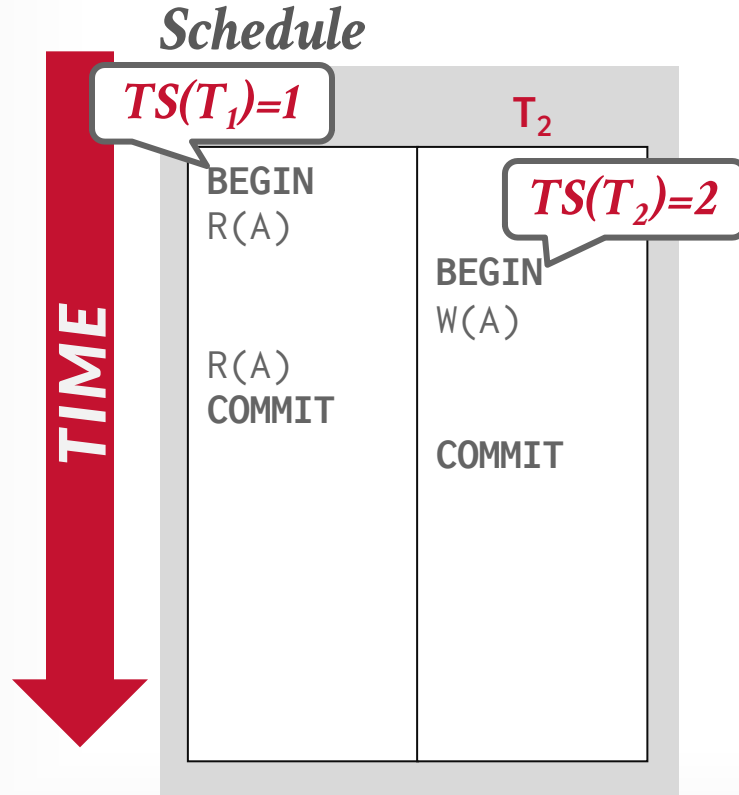
Database

	begin-ts	end-ts	value
A_0	0	2	123
A_1	2	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active
T_2	2	Active

MVCC – EXAMPLE #1



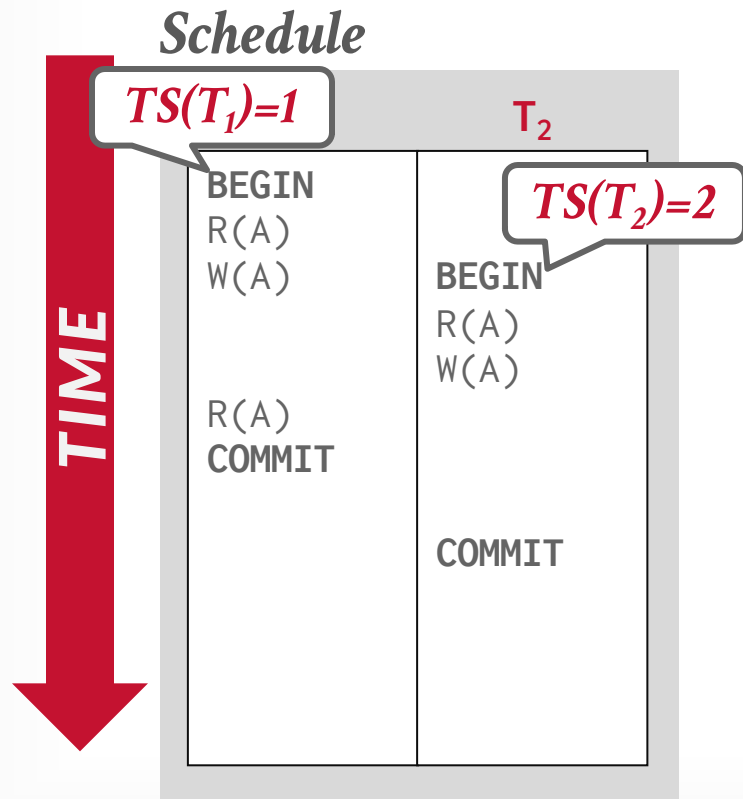
Database

	begin-ts	end-ts	value
A_0	0	2	123
A_1	2	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active
T_2	2	Active

MVCC – EXAMPLE #2



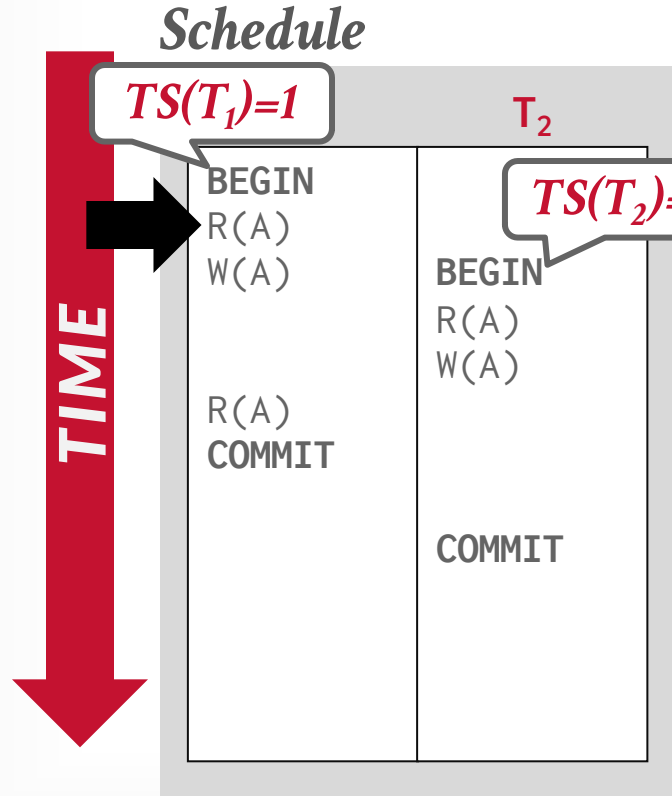
Database

	begin-ts	end-ts	value
A ₀	0	-	123

Txn Status Table

txnid	timestamp	status
T ₁	1	Active

MVCC – EXAMPLE #2



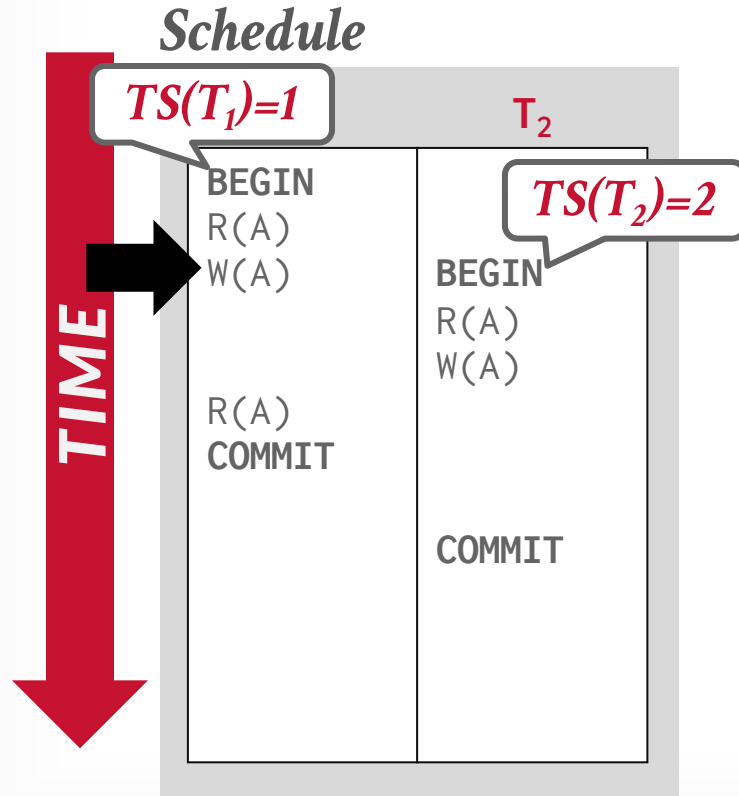
Database

	begin-ts	end-ts	value
A_0	0	-	123

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



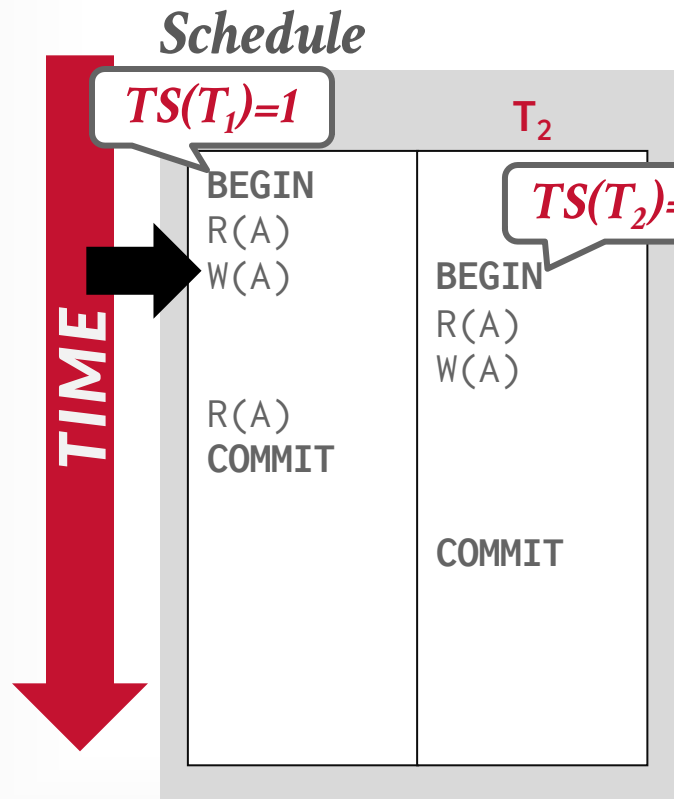
Database

	begin-ts	end-ts	value
A_0	0	-	123

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



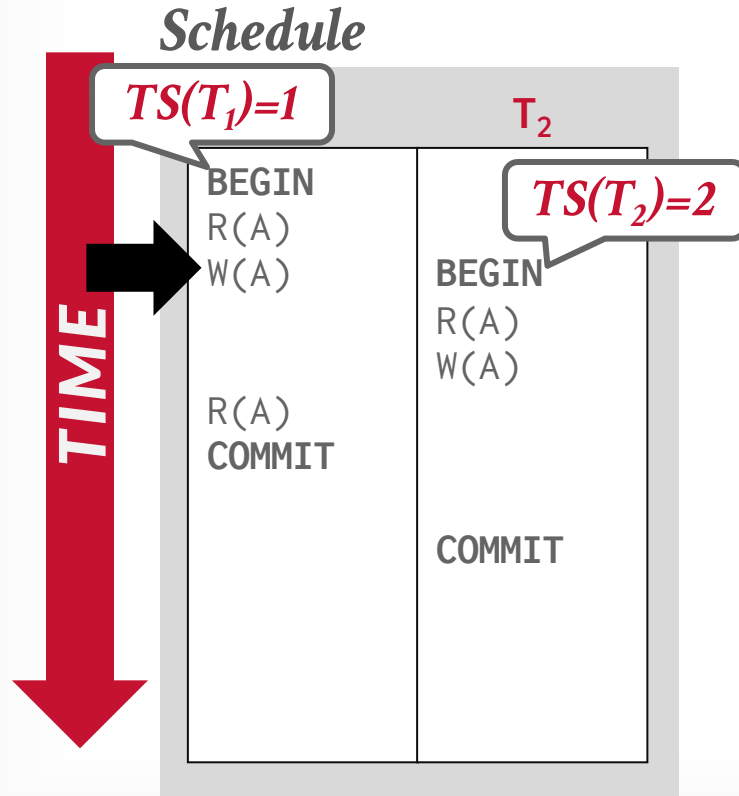
Database

	begin-ts	end-ts	value
A_0	0	-	123

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



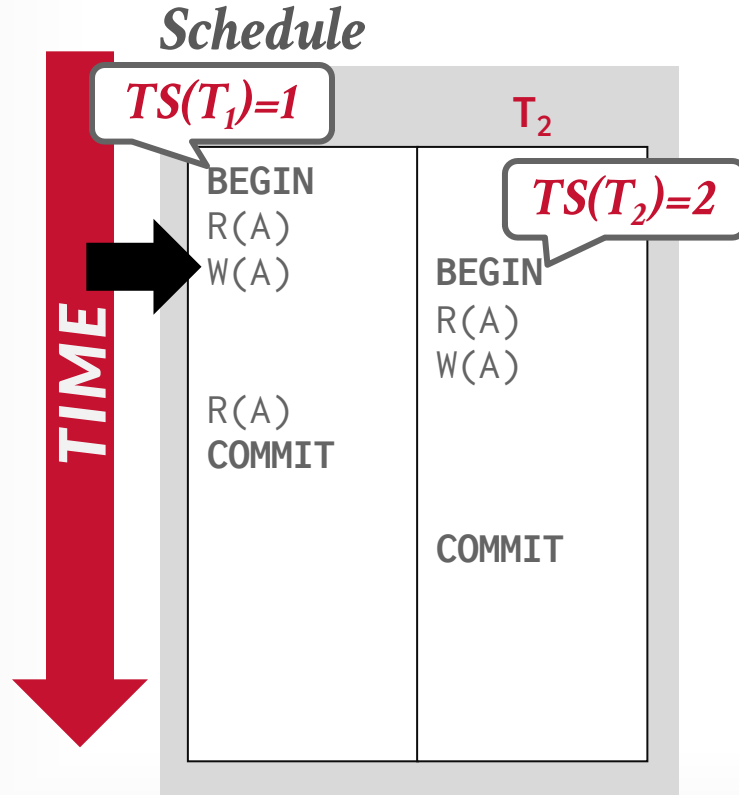
Database

	begin-ts	end-ts	value
A_0	0	-	123
A_1	1	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



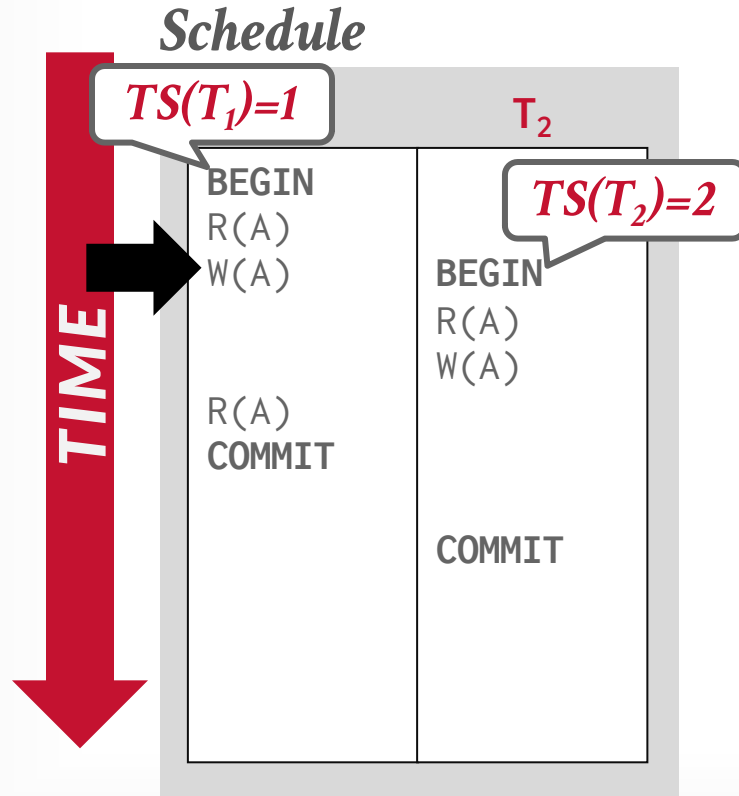
Database

	begin-ts	end-ts	value
A_0	0	-	123
A_1	1	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



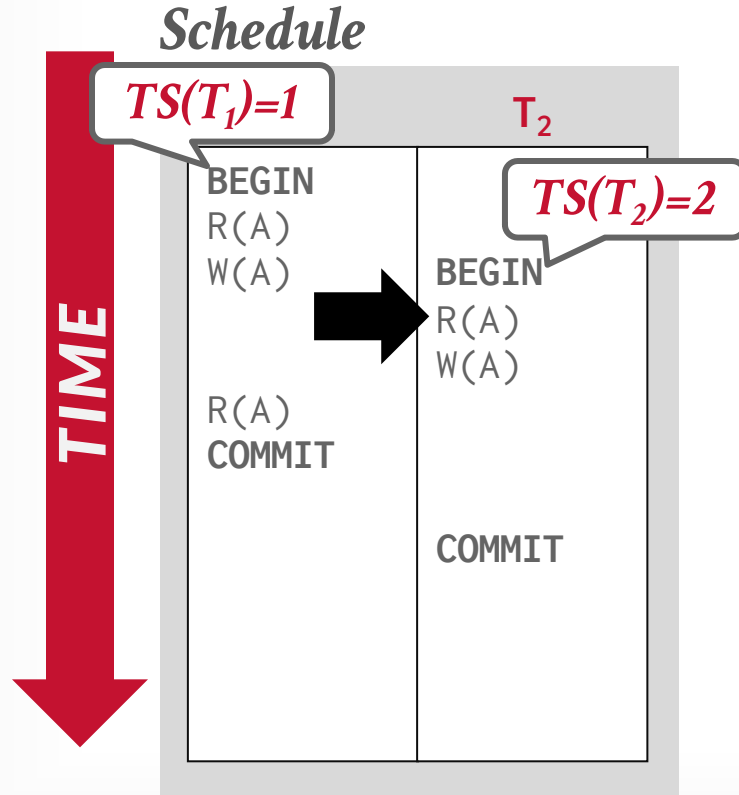
Database

	begin-ts	end-ts	value
A_0	0	1	123
A_1	1	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active

MVCC – EXAMPLE #2



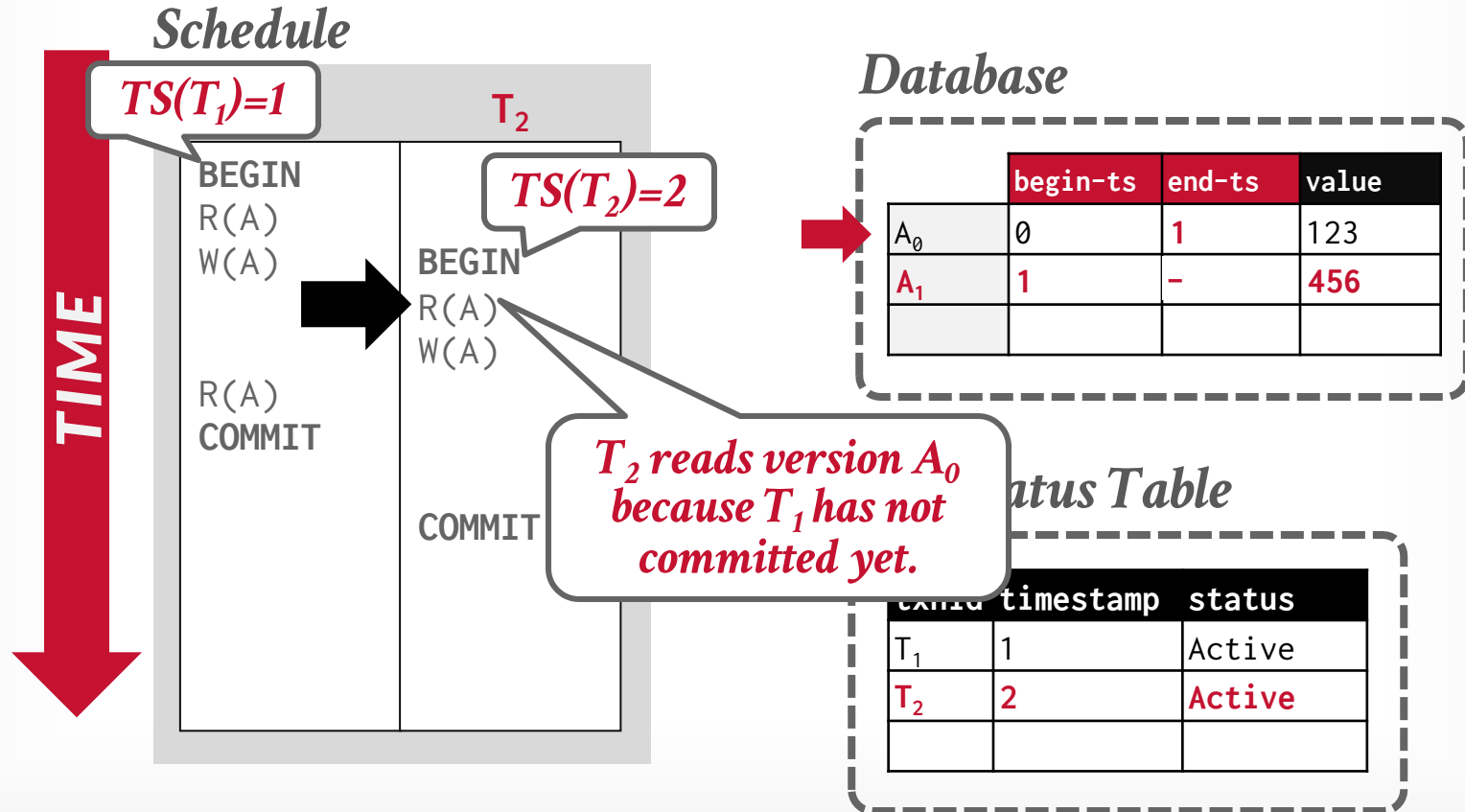
Database

	begin-ts	end-ts	value
A ₀	0	1	123
A ₁	1	-	456

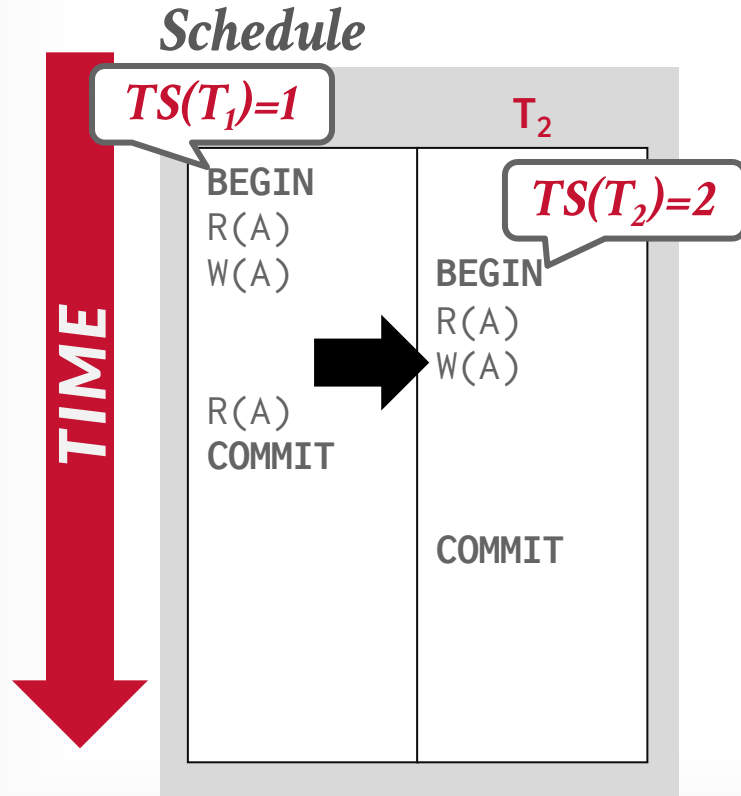
Txn Status Table

txnid	timestamp	status
T ₁	1	Active

MVCC – EXAMPLE #2



MVCC – EXAMPLE #2



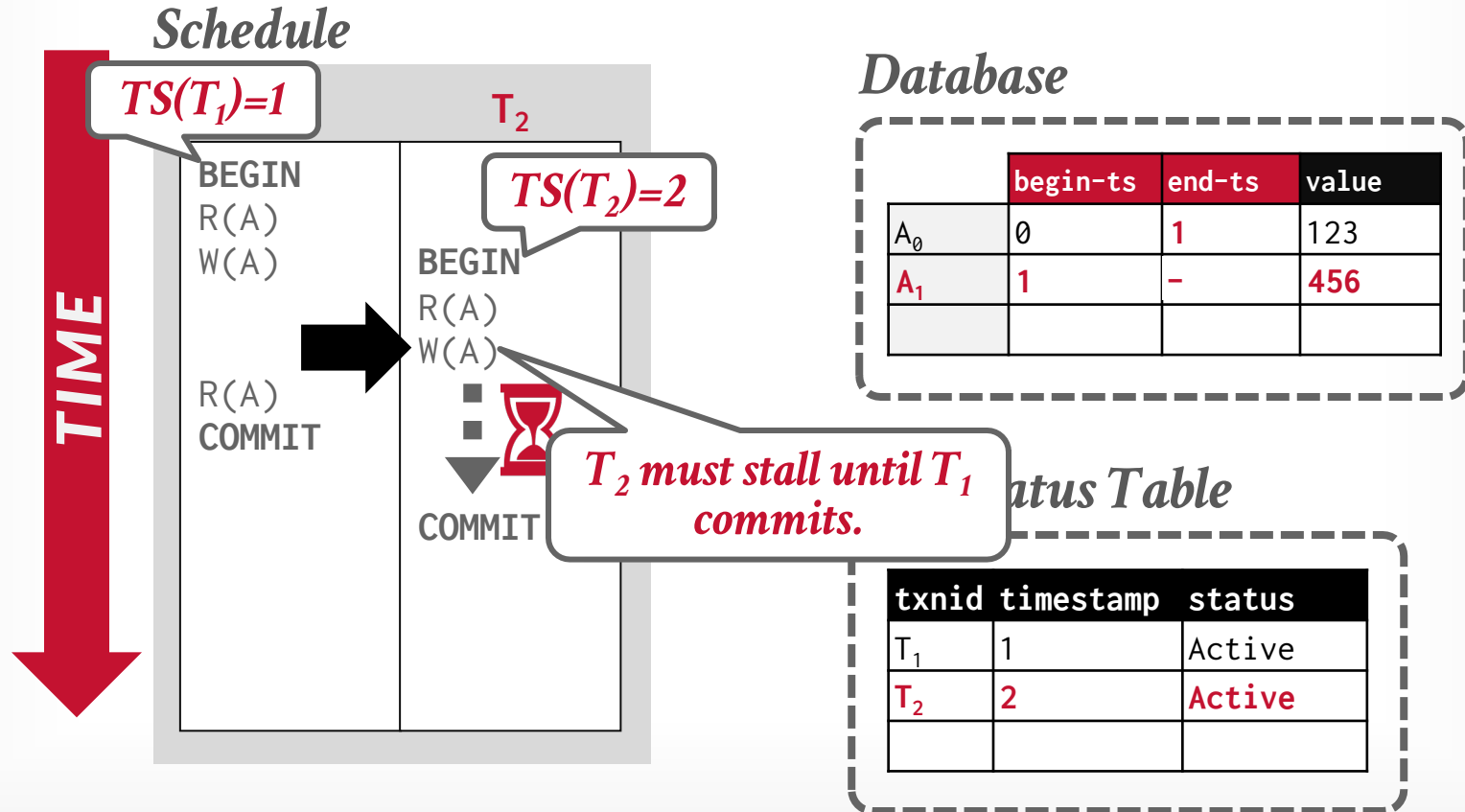
Database

	begin-ts	end-ts	value
A ₀	0	1	123
A ₁	1	-	456

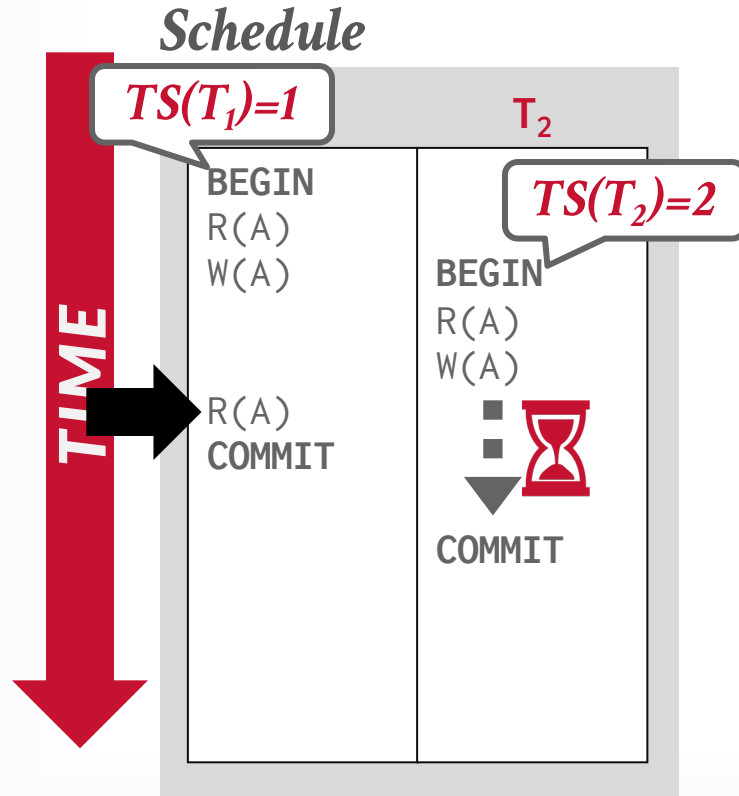
Txn Status Table

txnid	timestamp	status
T ₁	1	Active
T ₂	2	Active

MVCC – EXAMPLE #2



MVCC – EXAMPLE #2



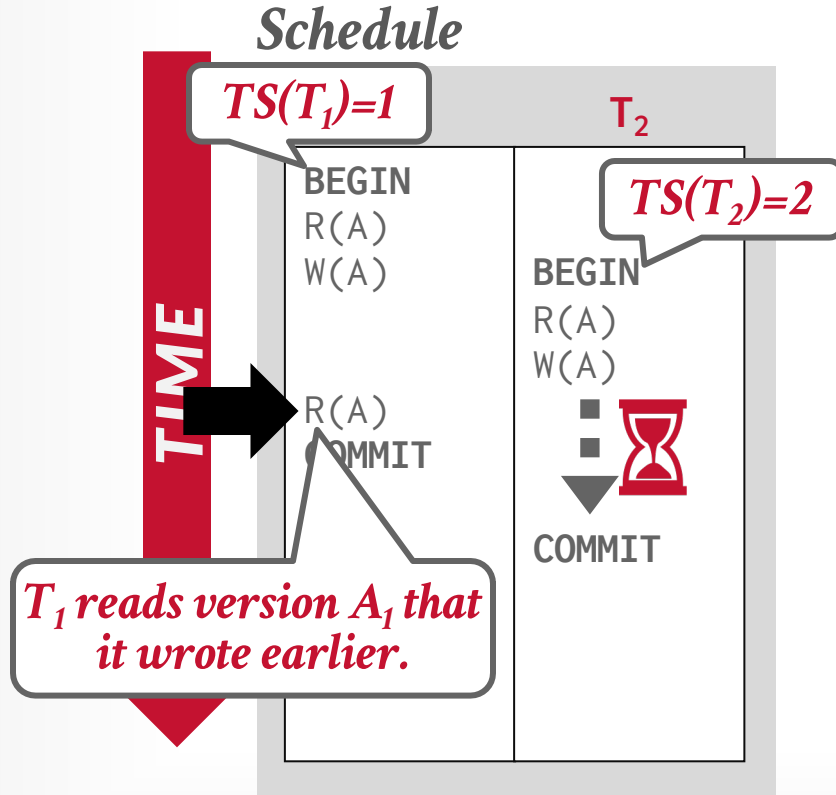
Database

	begin-ts	end-ts	value
A ₀	0	1	123
A ₁	1	-	456

Txn Status Table

txnid	timestamp	status
T ₁	1	Active
T ₂	2	Active

MVCC – EXAMPLE #2



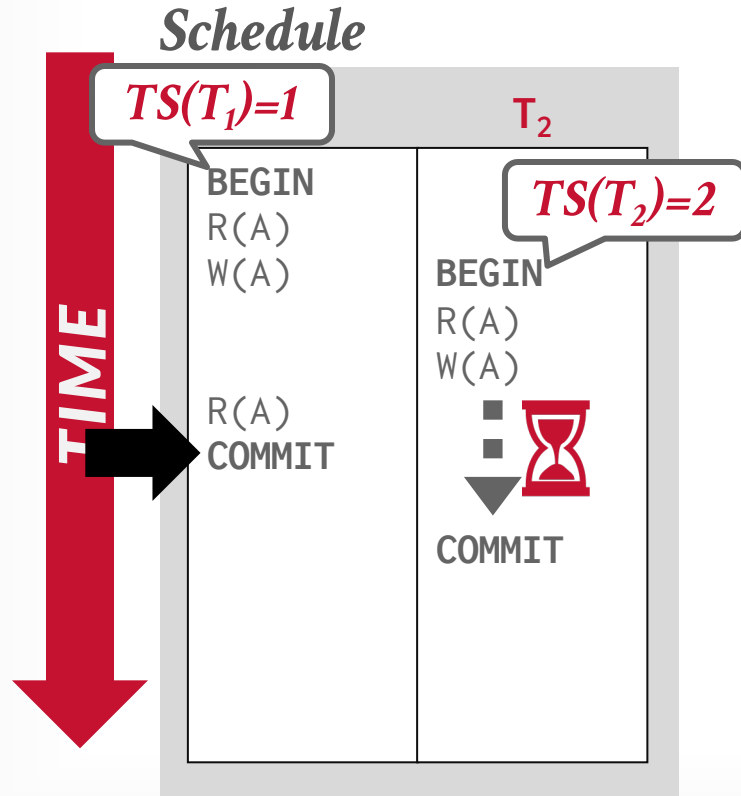
Database

	begin-ts	end-ts	value
A_0	0	1	123
A_1	1	-	456

Txn Status Table

txnid	timestamp	status
T_1	1	Active
T_2	2	Active

MVCC – EXAMPLE #2



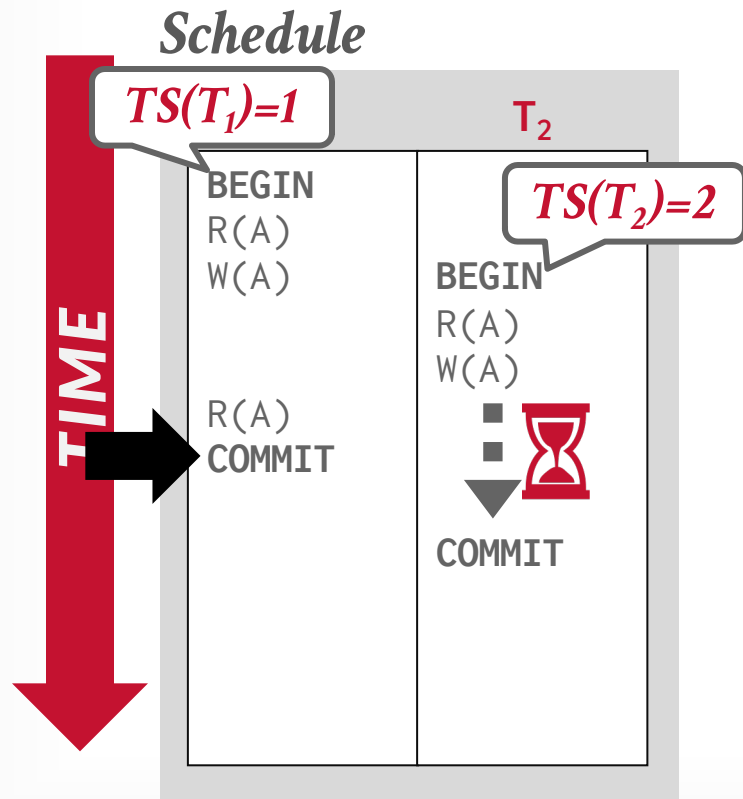
Database

	begin-ts	end-ts	value
A ₀	0	1	123
A ₁	1	-	456

Txn Status Table

txnid	timestamp	status
T ₁	1	Active
T ₂	2	Active

MVCC – EXAMPLE #2



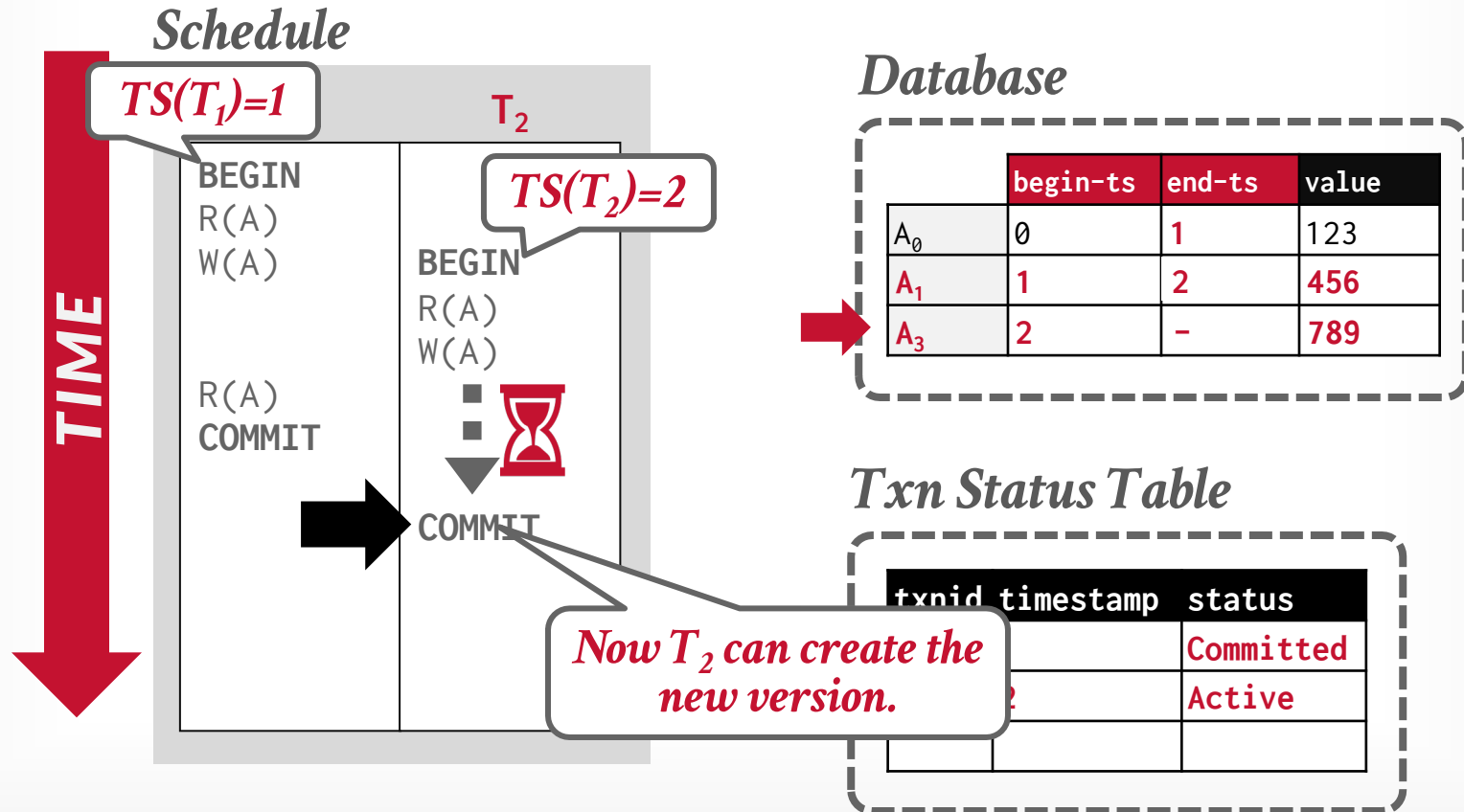
Database

	begin-ts	end-ts	value
A ₀	0	1	123
A ₁	1	-	456

Txn Status Table

txnid	timestamp	status
T ₁	1	Committed
T ₂	2	Active

MVCC – EXAMPLE #2



SNAPSHOT ISOLATION (SI)

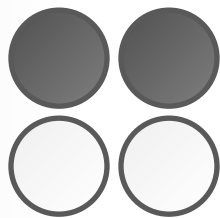
When a txn starts, it sees a consistent snapshot of the database that existed when that the txn started.

→ No torn writes from active txns.

→ If two txns update the same object, then first writer wins.

SI is susceptible to the Write Skew Anomaly.

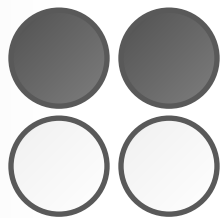
WRITE SKEW ANOMALY



WRITE SKEW ANOMALY

Txn #1

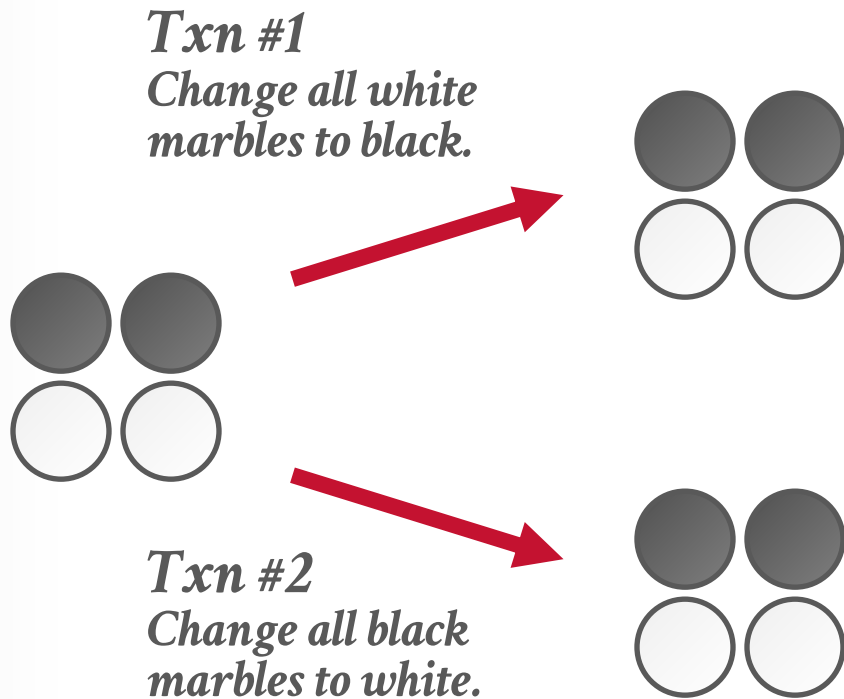
*Change all white
marbles to black.*



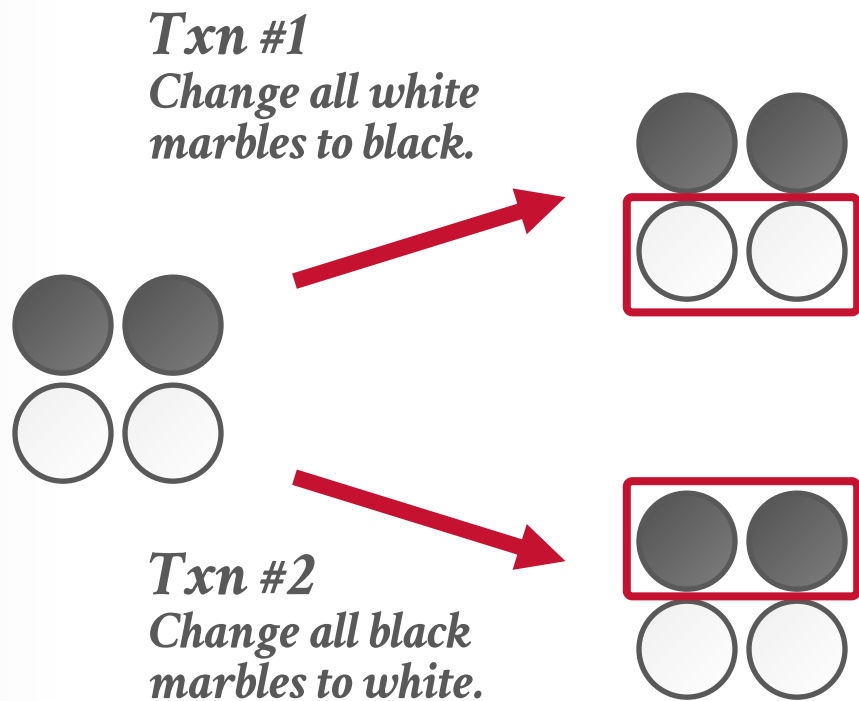
Txn #2

*Change all black
marbles to white.*

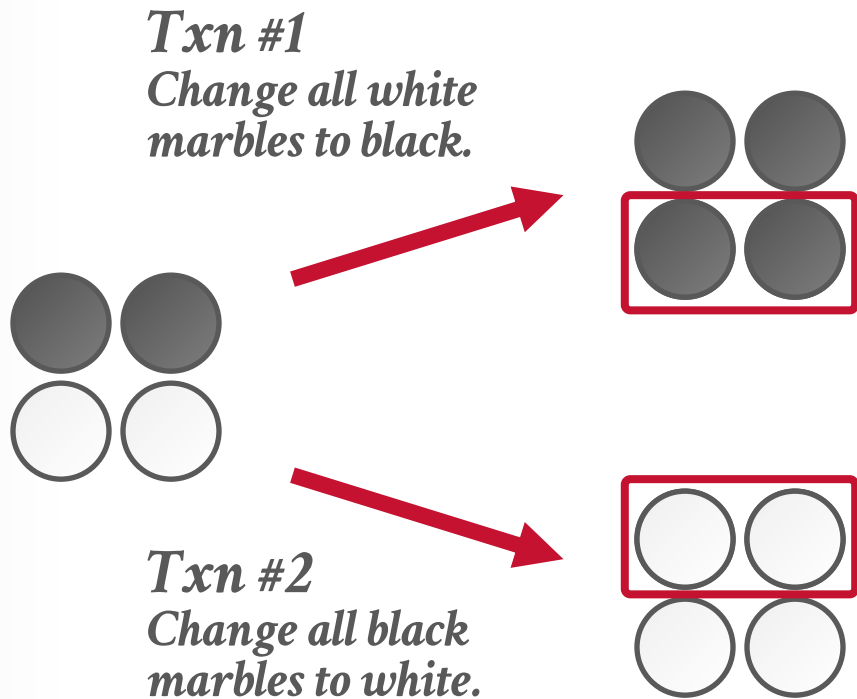
WRITE SKEW ANOMALY



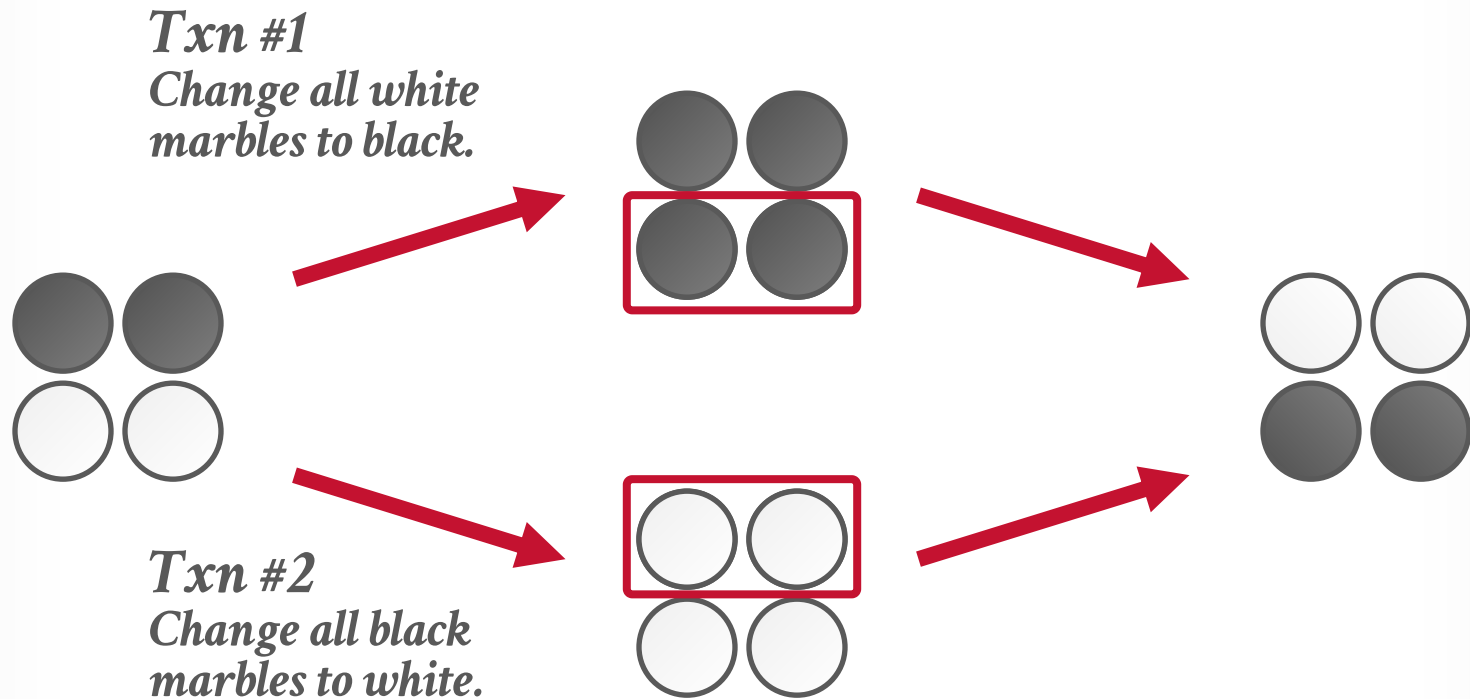
WRITE SKEW ANOMALY



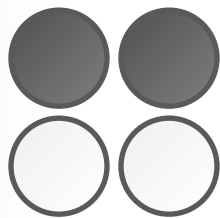
WRITE SKEW ANOMALY



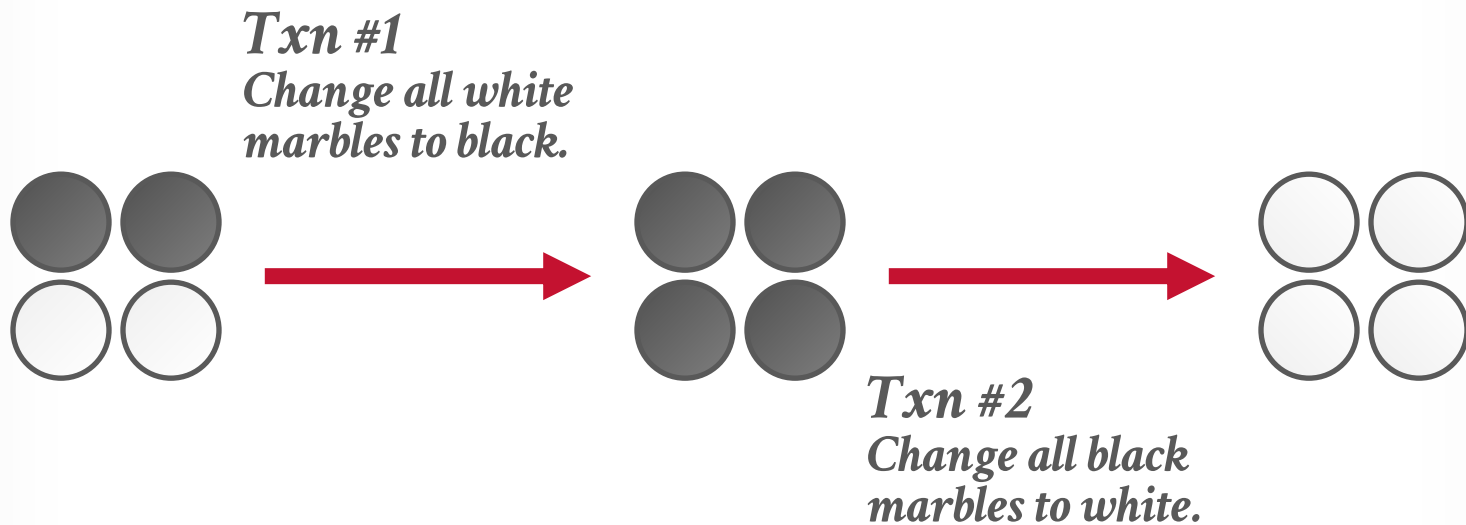
WRITE SKEW ANOMALY



WRITE SKEW ANOMALY



WRITE SKEW ANOMALY



MULTI-VERSION CONCURRENCY CONTROL

MVCC is more than just a concurrency control protocol. It completely affects how the DBMS manages transactions and the database.

MULTI-VERSION CONCURRENCY CONTROL

MVCC is more than just a concurrency control protocol. It completely affects how the DBMS manages transactions and the database.



MVCC DESIGN DECISIONS

Concurrency Control Protocol

Version Storage

Garbage Collection

Index Management

Deletes

CONCURRENCY CONTROL PROTOCOL

Approach #1: Timestamp Ordering

→ Assign txns timestamps that determine serial order.

Approach #2: Optimistic Concurrency Control

→ Three-phase protocol from last class.

→ Use private workspace for new versions.

Approach #3: Two-Phase Locking

→ Txns acquire appropriate lock on physical version before they can read/write a logical tuple.

VERSION STORAGE

The DBMS uses the tuples' pointer field to create a **version chain** per logical tuple.

- This allows the DBMS to find the version that is visible to a particular txn at runtime.
- Indexes always point to the “head” of the chain.

Different storage schemes determine where/what to store for each version.

VERSION STORAGE

Approach #1: Append-Only Storage

→ New versions are appended to the same table space.

Approach #2: Time-Travel Storage

→ Old versions are copied to separate table space.

Approach #3: Delta Storage

→ The original values of the modified attributes are copied into a separate delta record space.

VERSION STORAGE



*Don't
Do This!*

Approach #1: Append-Only Storage

→ New versions are appended to the same table space.

Approach #2: Time-Travel Storage

→ Old versions are copied to separate table space.

Approach #3: Delta Storage


→ The original values of the modified attributes are copied into a separate delta record space.

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table

	value	pointer
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table

	value	pointer
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table


	value	pointer
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table

	value	pointer
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset
A_2	\$333	\emptyset

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table

	value	pointer
A_0	\$111	●
A_1	\$222	∅
B_1	\$10	∅
A_2	\$333	∅

APPEND-ONLY STORAGE

All the physical versions of a logical tuple are stored in the same table space. The versions are inter-mixed.

On every update, append a new version of the tuple into an empty space in the table.

Main Table

	value	pointer
A_0	\$111	●
A_1	\$222	●
B_1	\$10	\emptyset
A_2	\$333	\emptyset

VERSION CHAIN ORDERING

Approach #1: Oldest-to-Newest (O2N)

- Append new version to end of the chain.
- Must traverse chain on look-ups.

Approach #2: Newest-to-Oldest (N2O)

- Must update index pointers for every new version.
- Do not have to traverse chain on look-ups.
- This is typically the better approach because most txns only care about the newest version.

TIME-TRAVEL STORAGE

Main Table

	value	pointer
A₂	\$222	●
B₁	\$10	

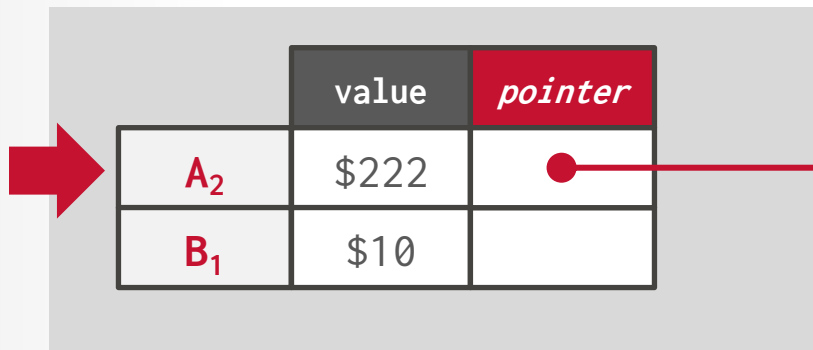
Time-Travel Table

	value	pointer
A₁	\$111	∅



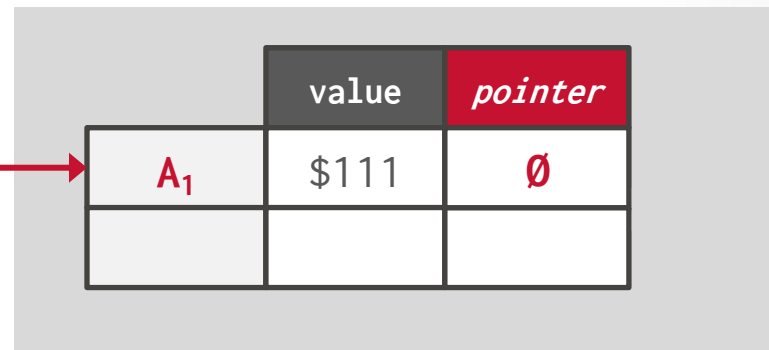
TIME-TRAVEL STORAGE

Main Table



	value	pointer
A₂	\$222	●
B₁	\$10	

Time-Travel Table

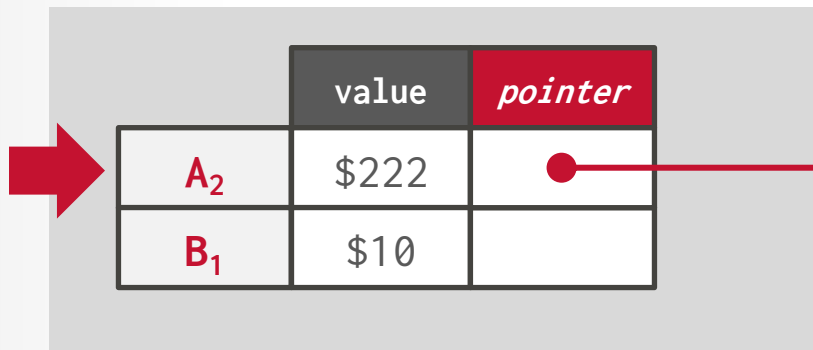


	value	pointer
A₁	\$111	∅

On every update, copy the current version to the time-travel table. Update pointers.

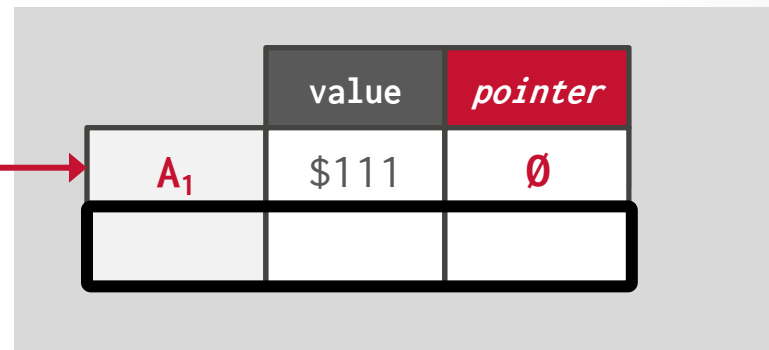
TIME-TRAVEL STORAGE

Main Table



	value	pointer
A₂	\$222	●
B₁	\$10	

Time-Travel Table

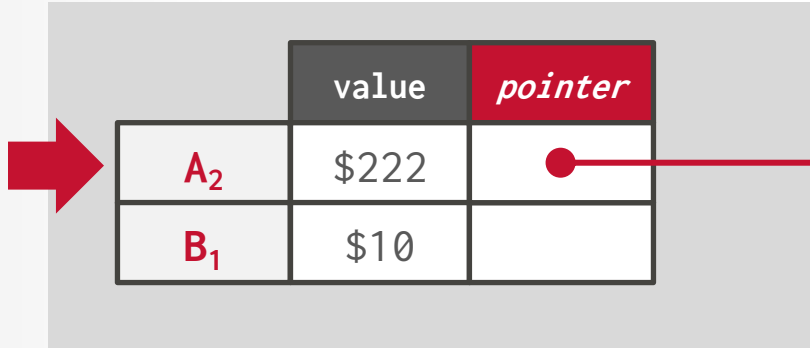


	value	pointer
A₁	\$111	∅

On every update, copy the current version to the time-travel table. Update pointers.

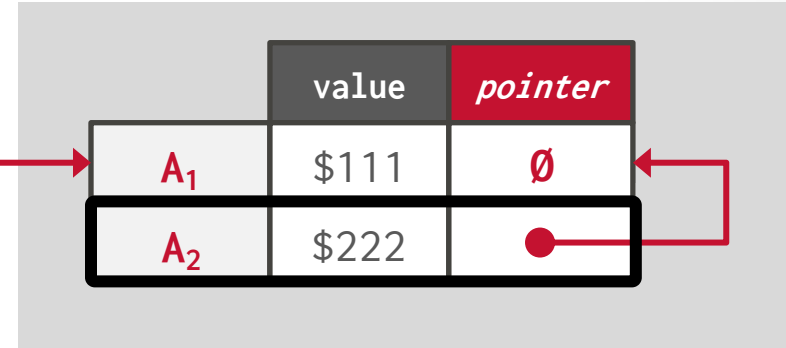
TIME-TRAVEL STORAGE

Main Table



	value	pointer
A ₂	\$222	●
B ₁	\$10	

Time-Travel Table

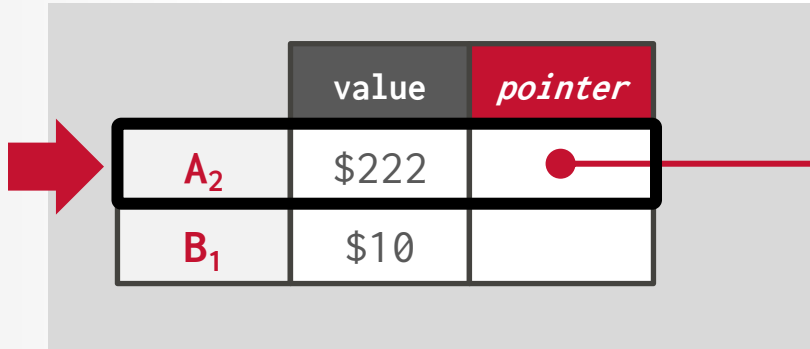


	value	pointer
A ₁	\$111	∅
A ₂	\$222	●

On every update, copy the current version to the time-travel table. Update pointers.

TIME-TRAVEL STORAGE

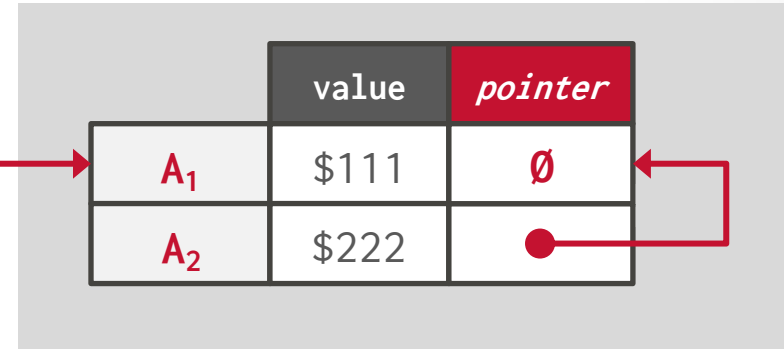
Main Table



	value	pointer
A₂	\$222	● →
B₁	\$10	

On every update, copy the current version to the time-travel table. Update pointers.

Time-Travel Table



	value	pointer
A₁	\$111	∅
A₂	\$222	● →

Overwrite master version in the main table and update pointers.

TIME-TRAVEL STORAGE

Main Table

	value	pointer
A ₃	\$333	● →
B ₁	\$10	

On every update, copy the current version to the time-travel table. Update pointers.

Time-Travel Table

	value	pointer
A ₁	\$111	∅
A ₂	\$222	● →

Overwrite master version in the main table and update pointers.

TIME-TRAVEL STORAGE

Main Table

	value	pointer
A ₃	\$333	●
B ₁	\$10	

On every update, copy the current version to the time-travel table. Update pointers.

Time-Travel Table

	value	pointer
A ₁	\$111	∅
A ₂	\$222	●

Overwrite master version in the main table and update pointers.

TIME-TRAVEL STORAGE

Main Table

	value	pointer
A ₃	\$333	●
B ₁	\$10	

On every update, copy the current version to the time-travel table. Update pointers.

Time-Travel Table

	value	pointer
A ₁	\$111	∅
A ₂	\$222	●

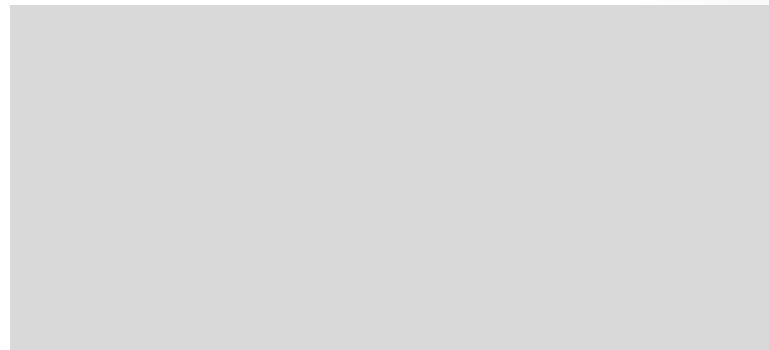
Overwrite master version in the main table and update pointers.

DELTA STORAGE

Main Table


	value	<i>pointer</i>
A₁	\$111	
B₁	\$10	

Delta Storage Segment



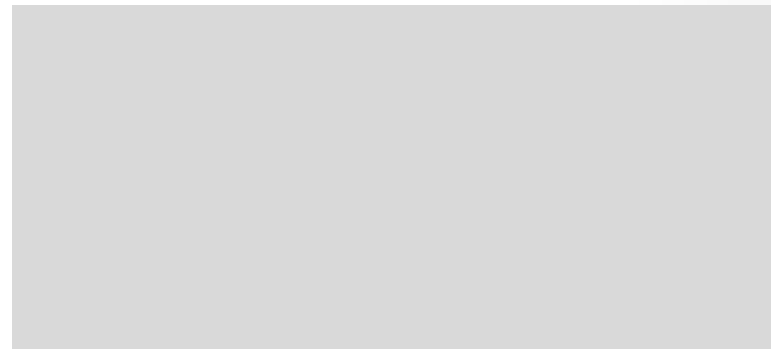
DELTA STORAGE

Main Table



	value	pointer
A ₁	\$111	
B ₁	\$10	

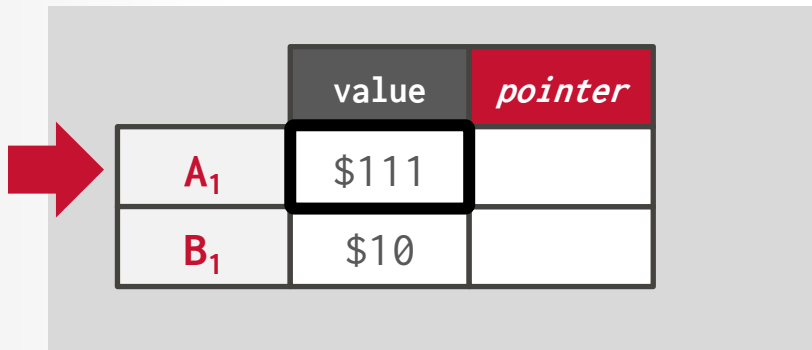
Delta Storage Segment



On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

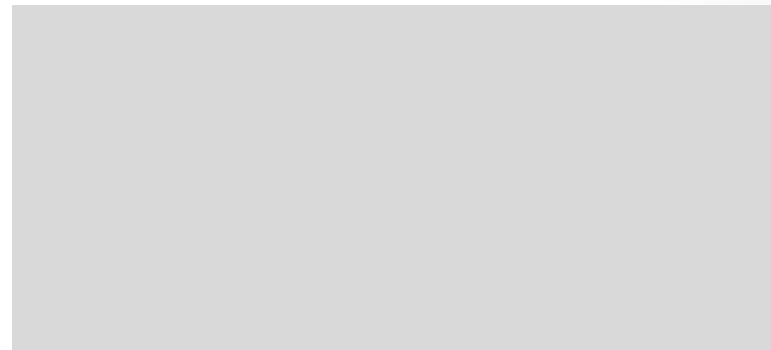
DELTA STORAGE

Main Table



	value	pointer
A ₁	\$111	
B ₁	\$10	

Delta Storage Segment



On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₁	\$111	
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE->\$111)	∅

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₂	\$222	● →
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE→\$111)	∅

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₂	\$222	●
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE->\$111)	∅
A ₂	(VALUE->\$222)	●

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₂	\$222	●
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE->\$111)	∅
A ₂	(VALUE->\$222)	●

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₃	\$333	●
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE->\$111)	∅
A ₂	(VALUE->\$222)	●

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

DELTA STORAGE

Main Table

	value	pointer
A ₃	\$333	●
B ₁	\$10	

Delta Storage Segment

	delta	pointer
A ₁	(VALUE->\$111)	∅
A ₂	(VALUE->\$222)	●

On every update, copy only the column values that were modified to the delta storage and overwrite the master version.

Txns can recreate old versions by applying the delta in reverse order.

GARBAGE COLLECTION

The DBMS needs to remove reclaimable physical versions from the database over time.

- No active txn in the DBMS can “see” that version (SI).
- The version was created by an aborted txn.

Two additional design decisions:

- How to look for expired versions?
- How to decide when it is safe to reclaim memory?

GARBAGE COLLECTION

The DBMS needs to remove reclaimable physical versions from the database over time.

- No active txn in the DBMS can “see” that version (SI).
- The version was created by an aborted txn.

Two additional design decisions:

- How to look for expired versions?
- How to decide when it is safe to reclaim memory?

GARBAGE COLLECTION

Approach #1: Tuple-level

- Find old versions by examining tuples directly.
- Background Vacuuming vs. Cooperative Cleaning

Approach #2: Transaction-level

- Txns keep track of their old versions so the DBMS does not have to scan tuples to determine visibility.

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

Vacuum



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

TUPLE-LEVEL GC

Txn #1

T_{id}=12

Txn #2

T_{id}=25

Vacuum



	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC

Txn #1

T_{id}=12

Txn #2

T_{id}=25

Vacuum



	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

Vacuum



	begin-ts	end-ts
<i>A₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₀</i>	<i>1</i>	<i>9</i>
<i>B₁₀₁</i>	<i>10</i>	<i>20</i>

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC

Txn #1

T_{id}=12

Txn #2

T_{id}=25

Vacuum

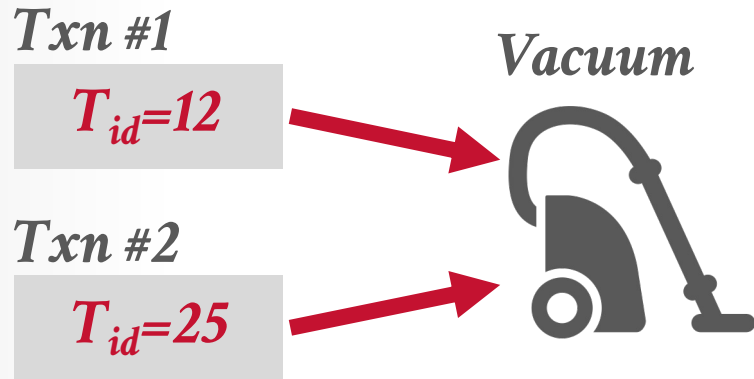


	begin-ts	end-ts
B₁₀₁	10	20

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC



	begin-ts	end-ts
B_{101}	10	20

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

Vacuum



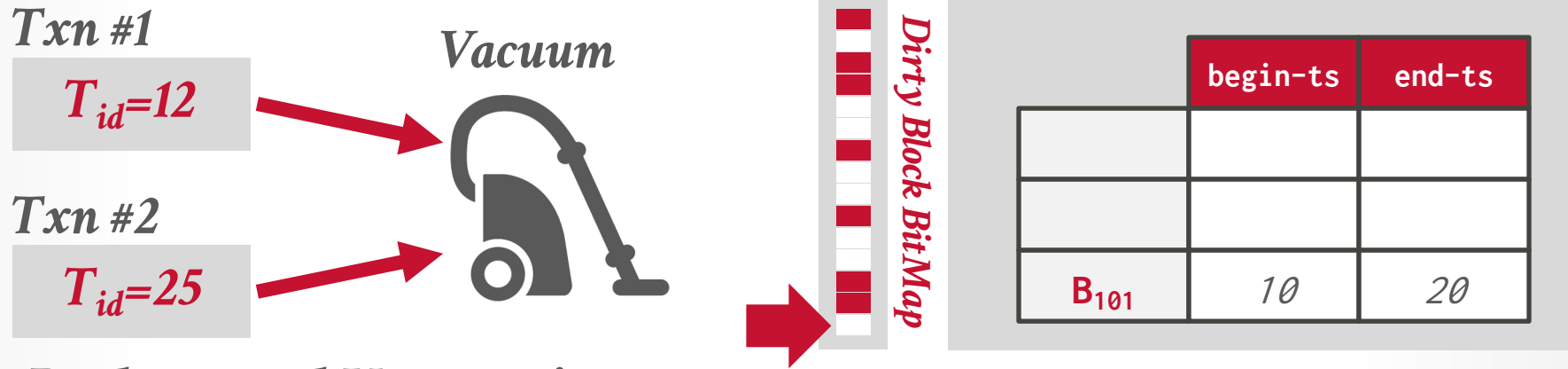
Dirty Block BitMap

	begin-ts	end-ts
B_{101}	10	20

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

TUPLE-LEVEL GC

Txn #1

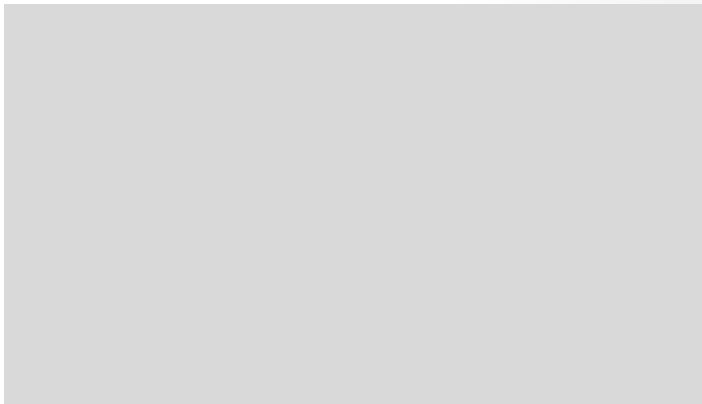
$T_{id}=12$

Txn #2

$T_{id}=25$

Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.



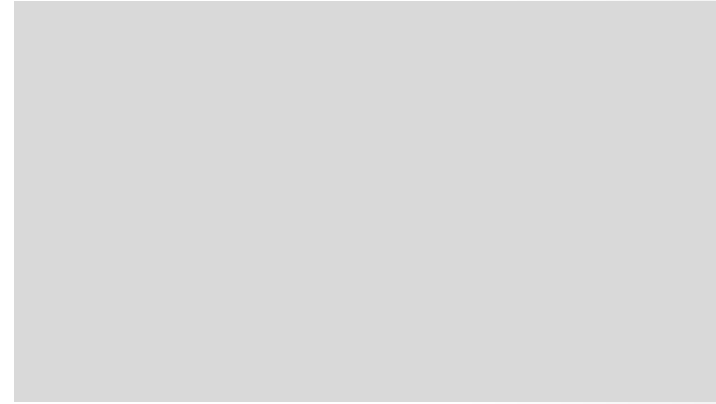
TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

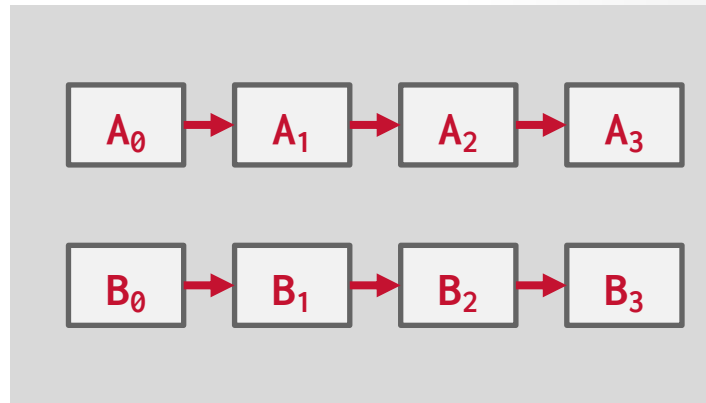
TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

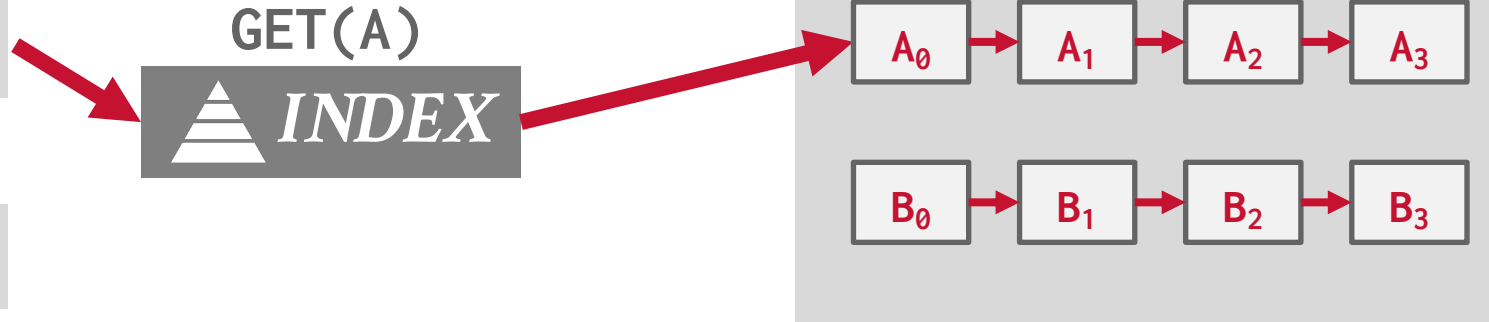
TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$



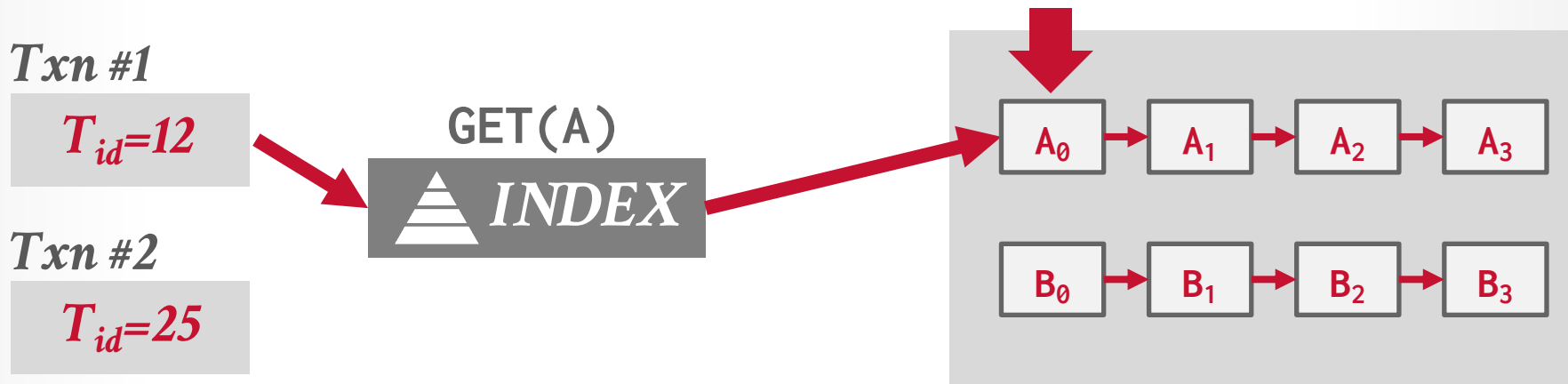
Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC



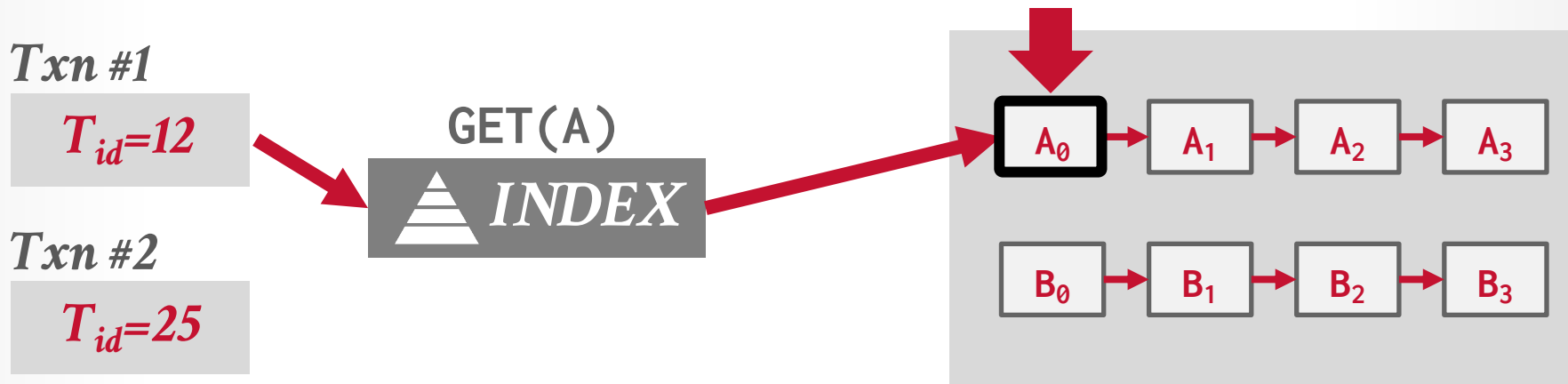
Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC



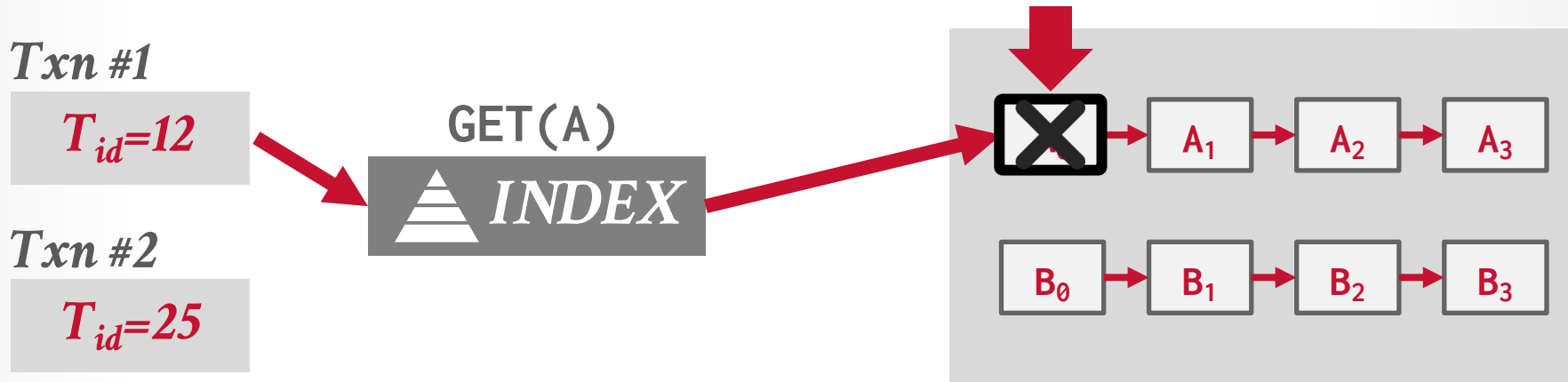
Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC

Txn #1

$T_{id}=12$

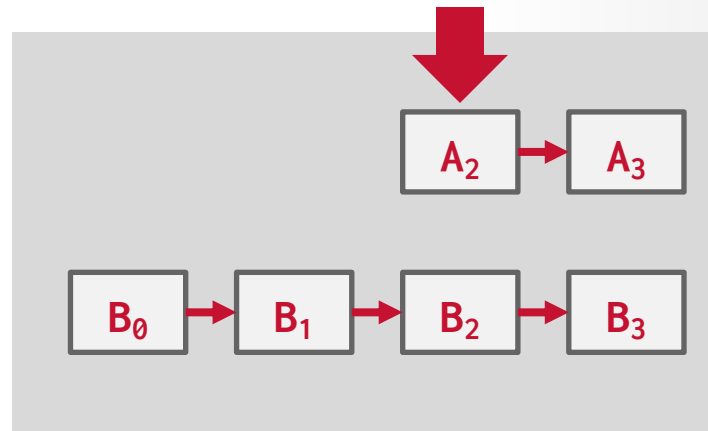
Txn #2

$T_{id}=25$



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.



Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TUPLE-LEVEL GC



Background Vacuuming:

Separate thread(s) periodically scan the table and look for reclaimable versions. Works with any storage.

Cooperative Cleaning:

Worker threads identify reclaimable versions as they traverse version chain. Only works with O2N.

TRANSACTION-LEVEL GC

Each txn keeps track of its read/write set.

On commit/abort, the txn provides this information to a centralized vacuum worker.

The DBMS periodically determines when all versions created by a finished txn are no longer visible.

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10

	begin-ts	end-ts	value
A₂	1	∞	-
B₆	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



	begin-ts	end-ts	value
A_2	1	∞	-
B_6	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



	begin-ts	end-ts	value
A_2	1	∞	-
B_6	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



	begin-ts	end-ts	value
A_2	1	10	-
B_6	8	∞	-
A_3	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



Old Versions

A₂

	begin-ts	end-ts	value
A₂	1	10	-
B₆	8	∞	-
A₃	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



Old Versions

A_2

	begin-ts	end-ts	value
A_2	1	10	-
B_6	8	∞	-
A_3	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)



UPDATE(B)

Old Versions

A₂



	begin-ts	end-ts	value
A₂	1	10	-
B₆	8	∞	-
A₃	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10

Old Versions

A₂



UPDATE(A)



UPDATE(B)



	begin-ts	end-ts	value
A₂	1	10	-
B₆	8	10	-
A₃	10	∞	-
B₇	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10

Old Versions

A_2

B_6



UPDATE(A)



UPDATE(B)

	begin-ts	end-ts	value
A_2	1	10	-
B_6	8	10	-
A_3	10	∞	-
B_7	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10
COMMIT @ 15

Old Versions

A₂

B₆



UPDATE(A)



UPDATE(B)

	begin-ts	end-ts	value
A₂	1	10	-
B₆	8	10	-
A₃	10	∞	-
B₇	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10
COMMIT @ 15

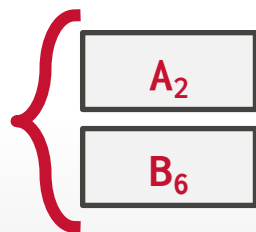
Old Versions



	begin-ts	end-ts	value
A_2	1	10	-
B_6	8	10	-
A_3	10	∞	-
B_7	10	∞	-

Vacuum

$TS < 10$



INDEX MANAGEMENT

Primary key indexes point to version chain head.

- How often the DBMS must update the pkey index depends on whether the system creates new versions when a tuple is updated.
- If a txn updates a tuple's pkey attribute(s), then this is treated as a **DELETE** followed by an **INSERT**.

Secondary indexes are more complicated...

IND

UBER Engineering

JOIN THE TEAM

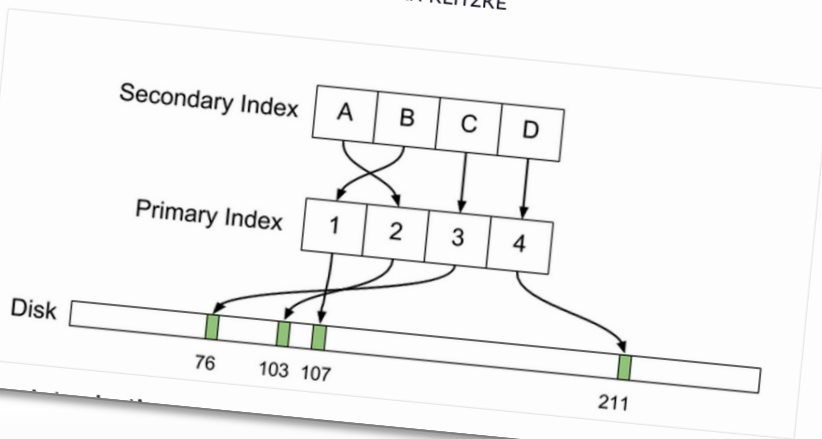
MEET THE PEOPLE

ARCHITECTURE

WHY UBER ENGINEERING SWITCHED FROM POSTGRES TO MYSQL

JULY 26, 2016

BY EVAN KLITZKE



Primary key in

→ How often the
on whether the
is updated.

→ If a txn updates
treated as a D

Secondary in

SECONDARY INDEXES

Approach #1: Logical Pointers

- Use a fixed identifier per tuple that does not change.
- Requires an extra indirection layer.
- Primary Key vs. Tuple Id

Approach #2: Physical Pointers

- Use the physical address to the version chain head.

INDEX POINTERS: APPEND-ONLY



PRIMARY INDEX



SECONDARY INDEX



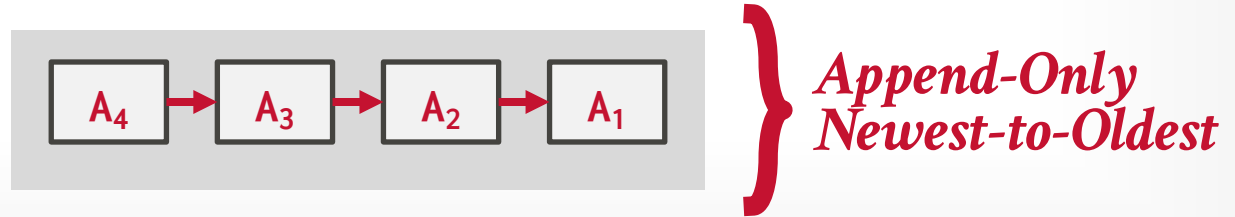
INDEX POINTERS: APPEND-ONLY



PRIMARY INDEX



SECONDARY INDEX



INDEX POINTERS: APPEND-ONLY

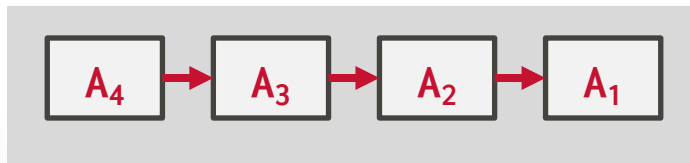
GET(A) ↓



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS: APPEND-ONLY

GET(A) ↓

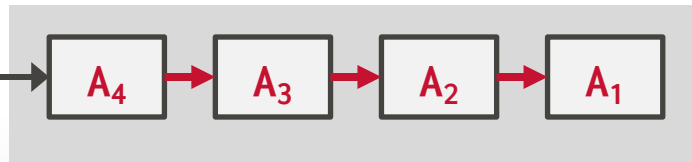


PRIMARY INDEX



SECONDARY INDEX

Record Id



*Append-Only
Newest-to-Oldest*

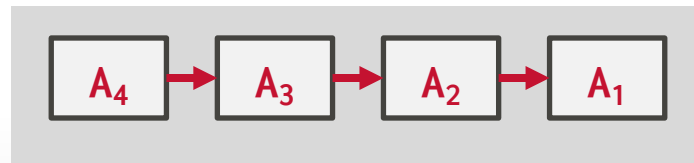
INDEX POINTERS: APPEND-ONLY



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS: APPEND-ONLY

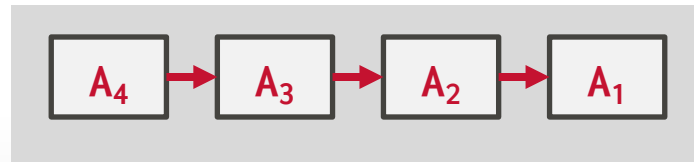
↓ GET(A)



PRIMARY INDEX

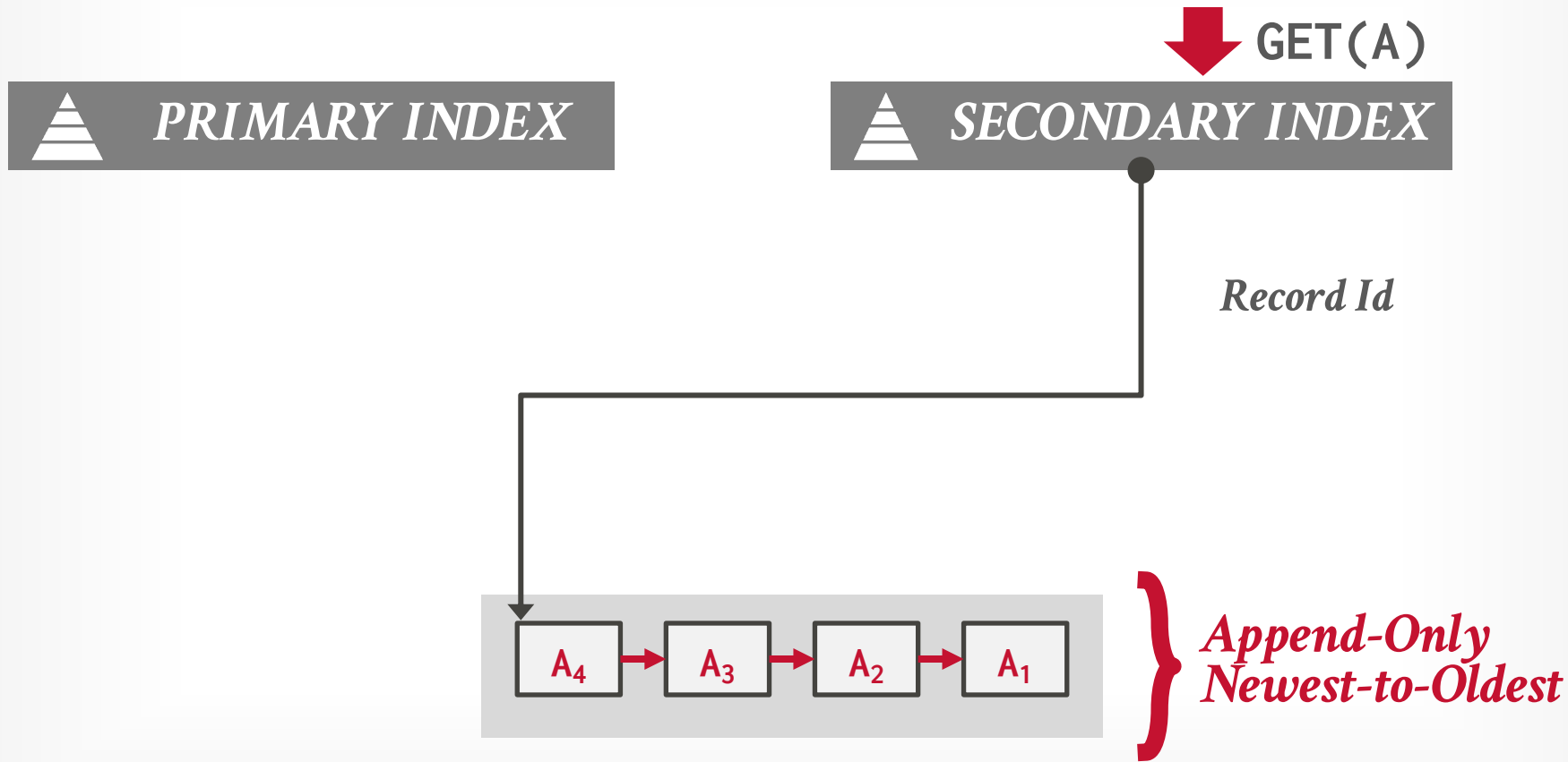


SECONDARY INDEX

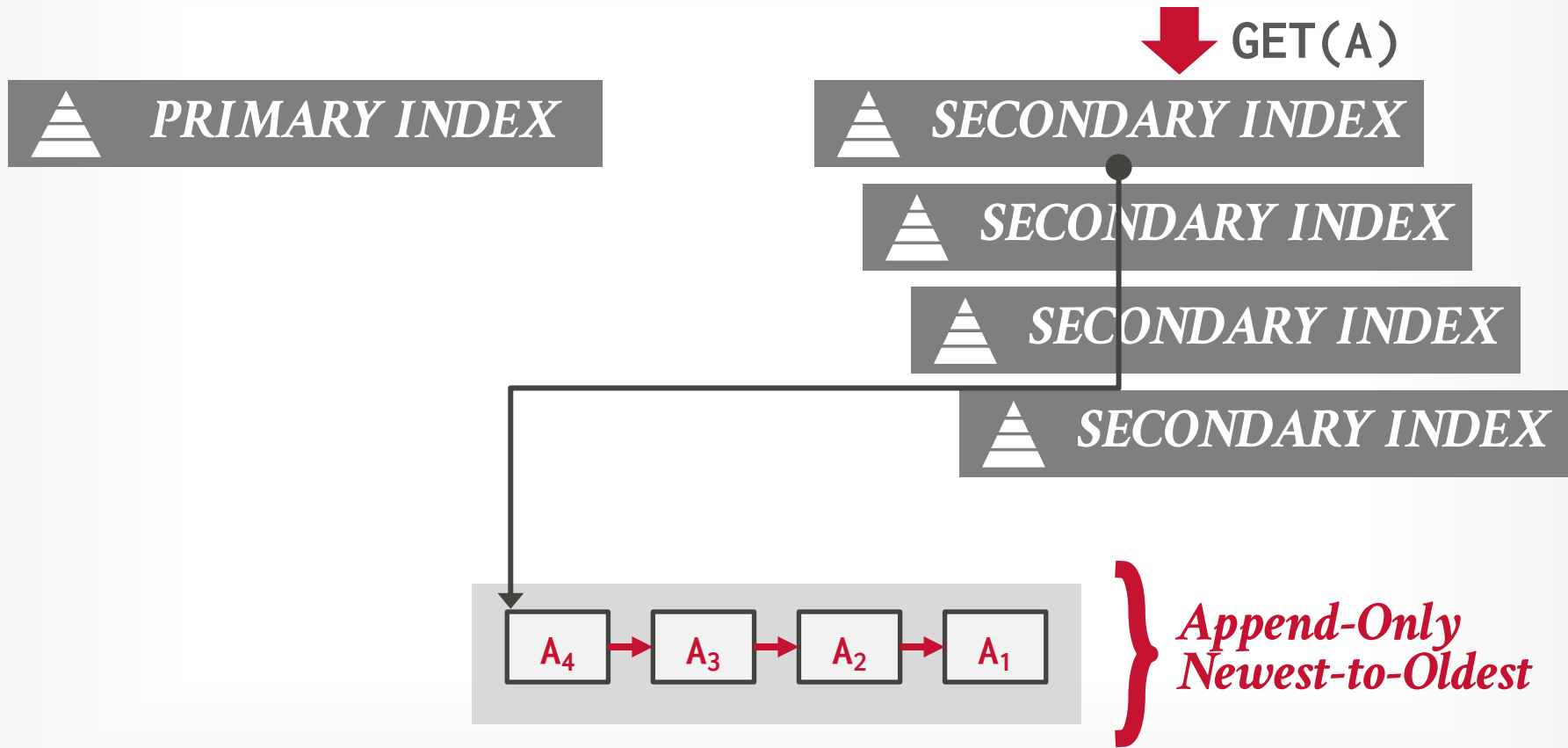


*Append-Only
Newest-to-Oldest*

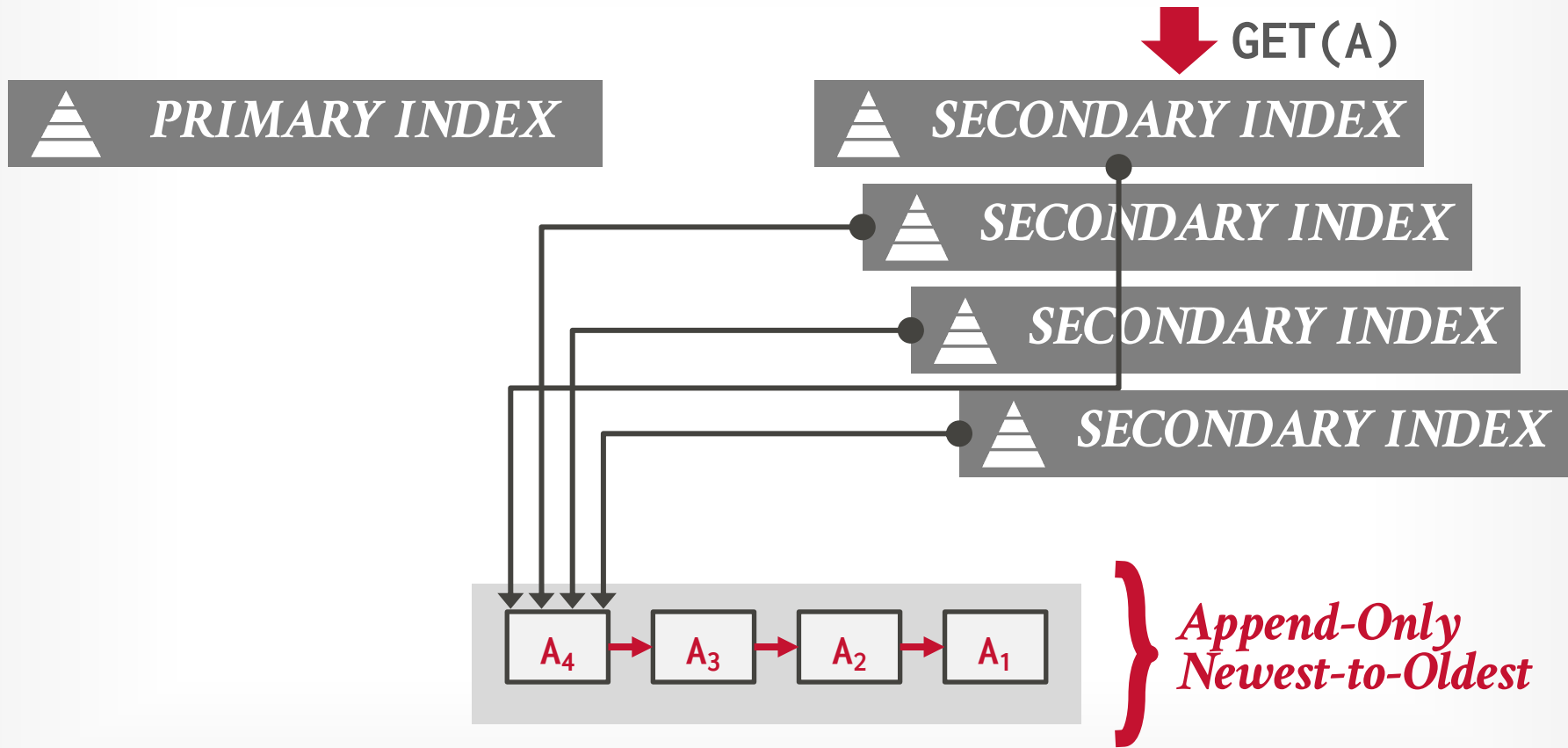
INDEX POINTERS: APPEND-ONLY



INDEX POINTERS: APPEND-ONLY



INDEX POINTERS: APPEND-ONLY



INDEX POINTERS: APPEND-ONLY

↓ GET(A)



PRIMARY INDEX

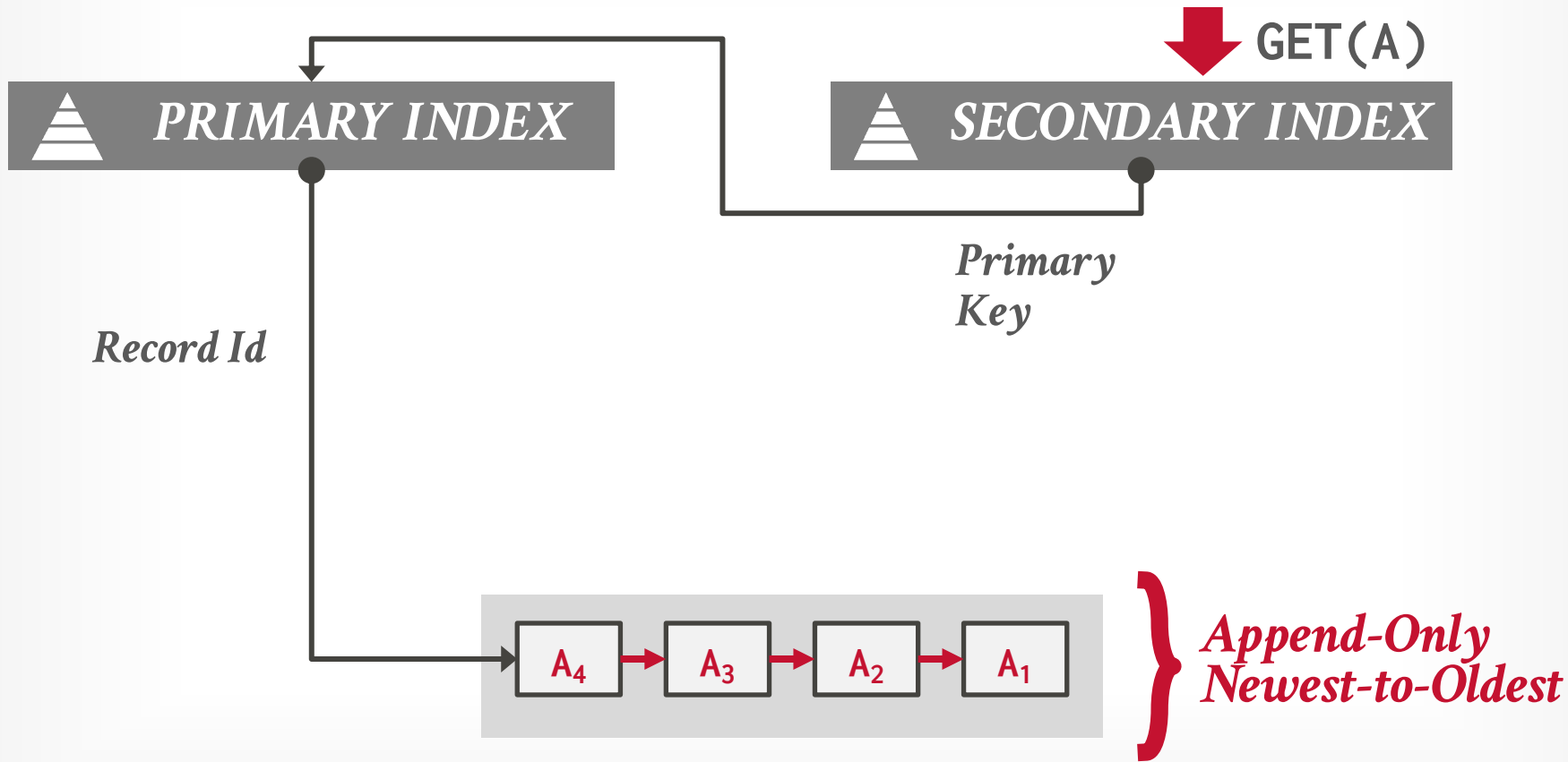


SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS: APPEND-ONLY



INDEX POINTERS: APPEND-ONLY

↓ GET(A)



PRIMARY INDEX



SECONDARY INDEX



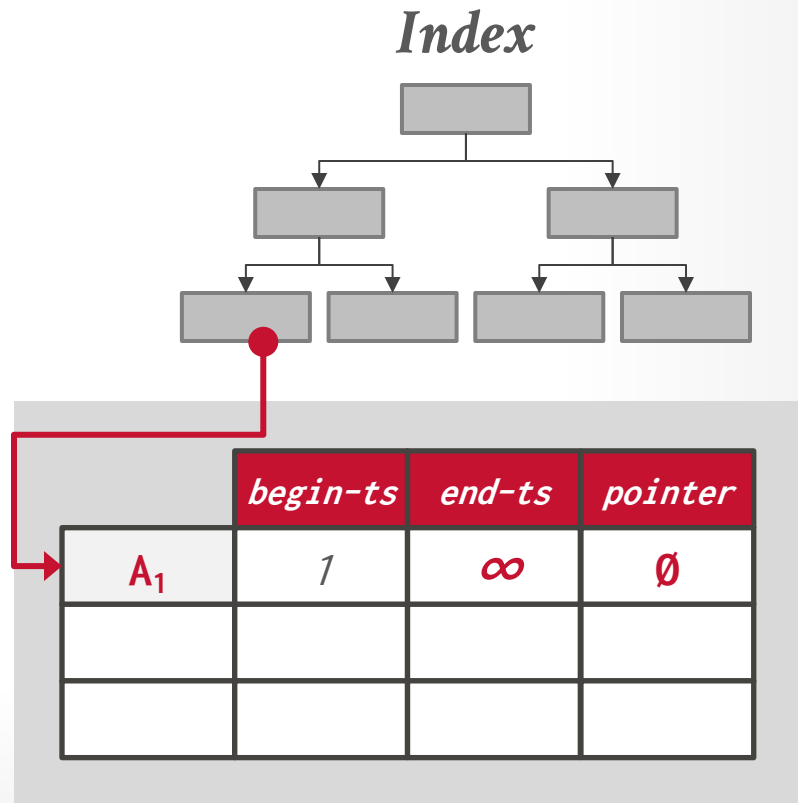
*Append-Only
Newest-to-Oldest*

MVCC INDEXES

MVCC DBMS indexes (usually) do not store version information about tuples with their keys.
→ Exception: Index-organized tables (e.g., MySQL)

Every index must support duplicate keys from different snapshots:
→ The same key may point to different logical tuples in different snapshots.

MVCC DUPLICATE KEY PROBLEM



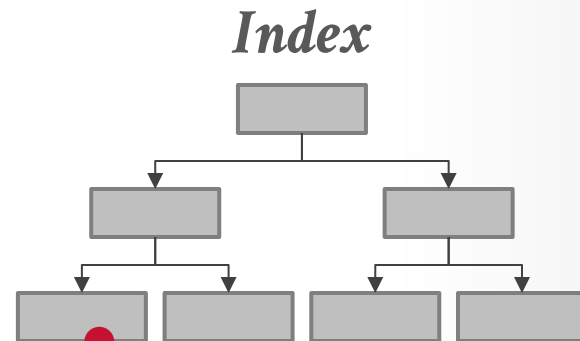
MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
<i>A₁</i>	<i>1</i>	<i>∞</i>	<i>∅</i>

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

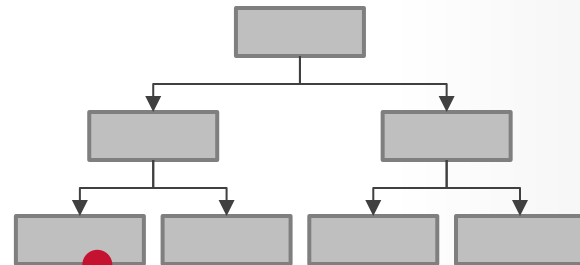
Txn #2

BEGIN @ 20



UPDATE(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
<i>A₁</i>	<i>1</i>	<i>∞</i>	<i>∅</i>

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

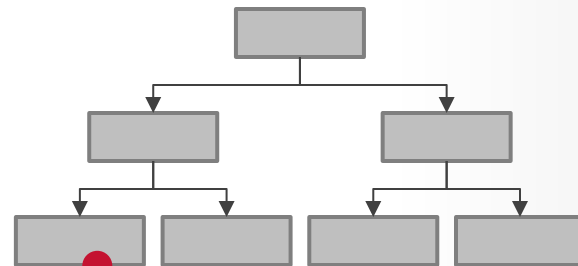
Txn #2

BEGIN @ 20



UPDATE(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
A_1	1	∞	●
A_2	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

Txn #2

BEGIN @ 20

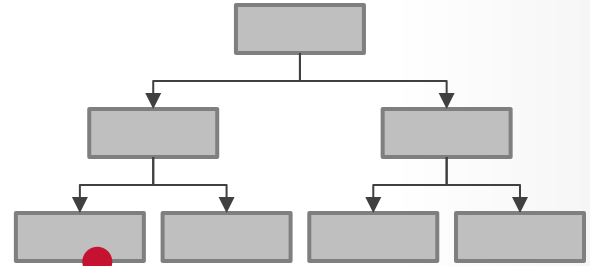


UPDATE(A)



DELETE(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
<i>A</i> ₁	1	∞	
	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

Txn #2

BEGIN @ 20

COMMIT @ 25

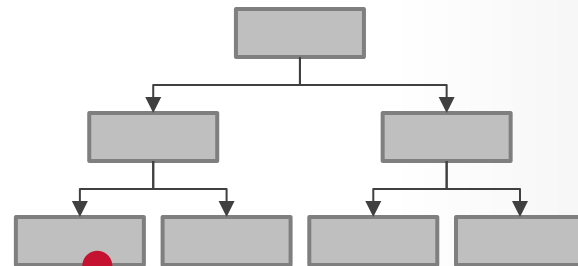


UPDATE(A)



DELETE(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
<i>A</i> ₁	1	∞	
	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

Txn #2

BEGIN @ 20

COMMIT @ 25

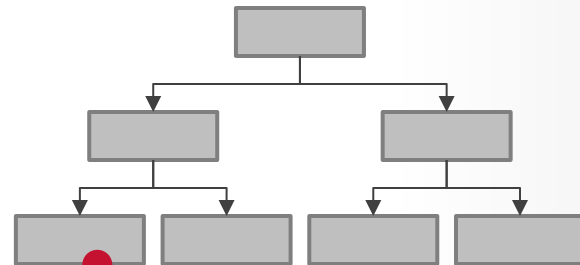


UPDATE(A)



DELETE(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
A₁	1	20	
	20	20	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



Txn #2

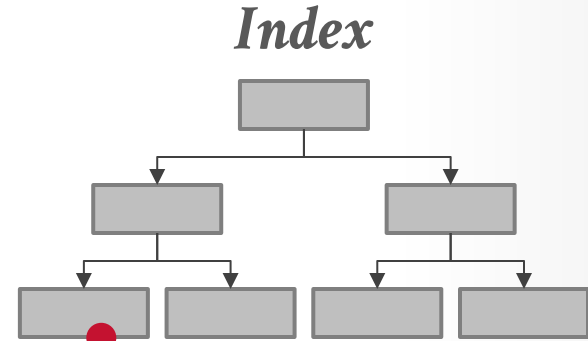
BEGIN @ 20

COMMIT @ 25



Txn #3

BEGIN @ 30



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
<i>A</i> ₁	1	20	
	20	20	∅

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

Txn #2

BEGIN @ 20

COMMIT @ 25



UPDATE(A)



DELETE(A)

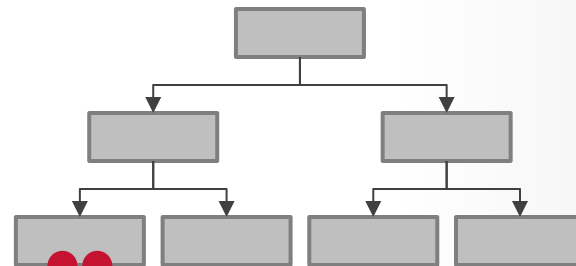
Txn #3

BEGIN @ 30



INSERT(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
A_1	1	20	
X	20	20	\emptyset
A_3	30	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)



READ(A)

Txn #2

BEGIN @ 20

COMMIT @ 25



UPDATE(A)



DELETE(A)

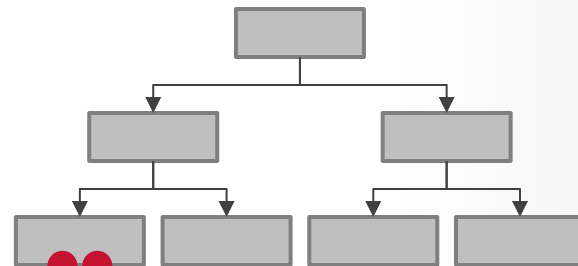
Txn #3

BEGIN @ 30



INSERT(A)

Index



	<i>begin-ts</i>	<i>end-ts</i>	<i>pointer</i>
A_1	1	20	
X	20	20	\emptyset
A_3	30	∞	\emptyset

MVCC INDEXES

Each index's underlying data structure must support the storage of non-unique keys.

Use additional execution logic to perform conditional inserts for pkey / unique indexes.

→ Atomically check whether the key exists and then insert.

Workers may get back multiple entries for a single fetch. They then must follow the pointers to find the proper physical version.

MVCC DELETES

The DBMS physically deletes a tuple from the database only when all versions of a logically deleted tuple are not visible.

- If a tuple is deleted, then there cannot be a new version of that tuple after the newest version.
- No write-write conflicts / first-writer wins

We need a way to denote that tuple has been logically delete at some point in time.

MVCC DELETES

Approach #1: Deleted Flag

- Maintain a flag to indicate that the logical tuple has been deleted after the newest physical version.
- Can either be in tuple header or a separate column.

Approach #2: Tombstone Tuple

- Create an empty physical version to indicate that a logical tuple is deleted.
- Use a separate pool for tombstone tuples with only a special bit pattern in version chain pointer to reduce the storage overhead.

MVCC IMPLEMENTATIONS

	<i>Protocol</i>	<i>Version Storage</i>	<i>Garbage Collection</i>	<i>Indexes</i>
Oracle	MV2PL	Delta	Vacuum	Logical
Postgres	MV-2PL/MV-TO	Append-Only	Vacuum	Physical
MySQL-InnoDB	MV-2PL	Delta	Vacuum	Logical
HYRISE	MV-OCC	Append-Only	-	Physical
Hekaton	MV-OCC	Append-Only	Cooperative	Physical
MemSQL (2015)	MV-OCC	Append-Only	Vacuum	Physical
SAP HANA	MV-2PL	Time-travel	Hybrid	Logical
NuoDB	MV-2PL	Append-Only	Vacuum	Logical
HyPer	MV-OCC	Delta	Txn-level	Logical
CockroachDB	MV-2PL	Delta (LSM)	Compaction	Logical

CONCLUSION

MVCC is the widely used scheme in DBMSs.
Even systems that do not support multi-statement txns (e.g., NoSQL) use it.

NEXT CLASS

Logging and recovery!