

Carnegie
Mellon
University

Intro to Database
Systems (15-445/645)

Lecture #19

Multi- Version Concurrency Control

SPRING 2024 >> Prof. Jignesh Patel



ADMINISTRIVIA

Project #3 is due Sun April 7, 2024 @ 11:59pm

Final Exam

→ Thu May 2, 2024 @ 05:30pm-08:30pm

Lectures #23 and #24

→ Recorded lectures and will be posted next week

→ I'm traveling for a small group meeting on “Hardware Support for Cloud Database Systems in the Post-Moore’s Law Era”

Lecture #26: Guest Speaker from Snowflake

→ Devin Petersohn on “Beyond SQL: Dataframes in the Database”

<https://db.cs.cmu.edu/events/spring-2024-beyond-sql-dataframes-in-the-database-devin-petersohn/>

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



Oracle Rdb
*the Database for HP
OpenVMS Platform*

MULTI-VERSION CONCURRENCY CONTROL

Writers do not block readers.

Readers do not block writers.

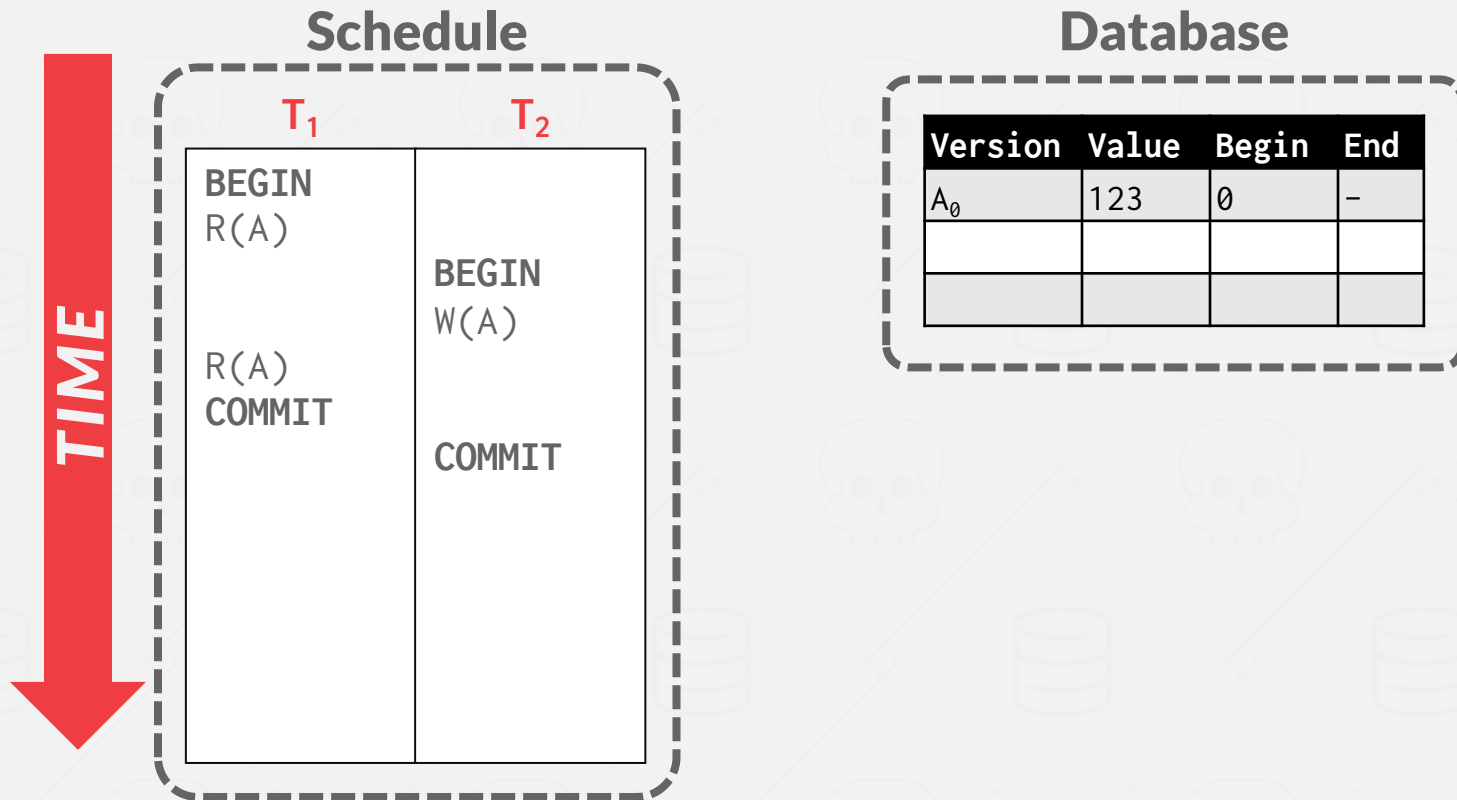
Read-only txns can read a consistent snapshot without acquiring locks.

→ Use timestamps to determine visibility.

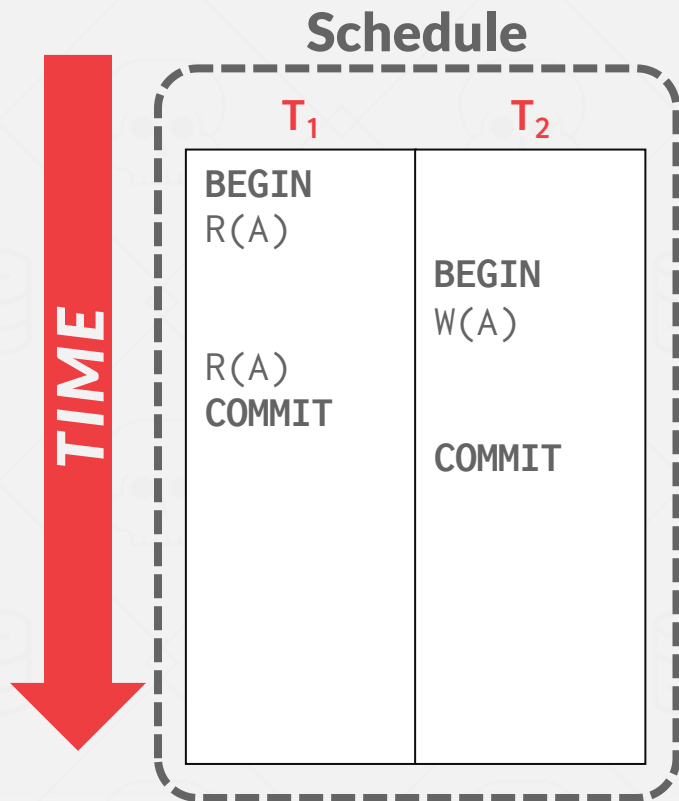
→ MVCC naturally supports Snapshot Isolation (SI).

Easily support time-travel queries.

MVCC - EXAMPLE #1



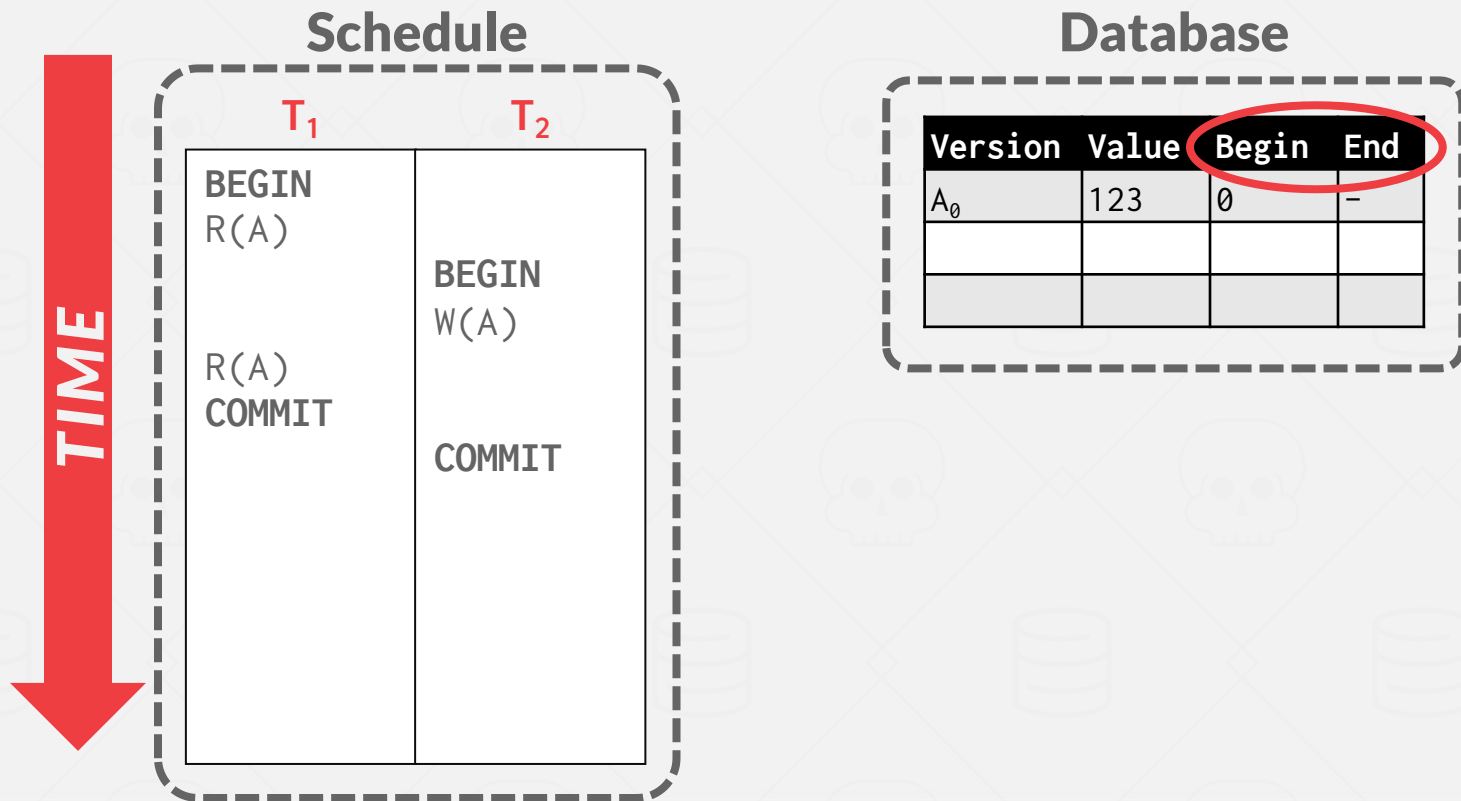
MVCC - EXAMPLE #1



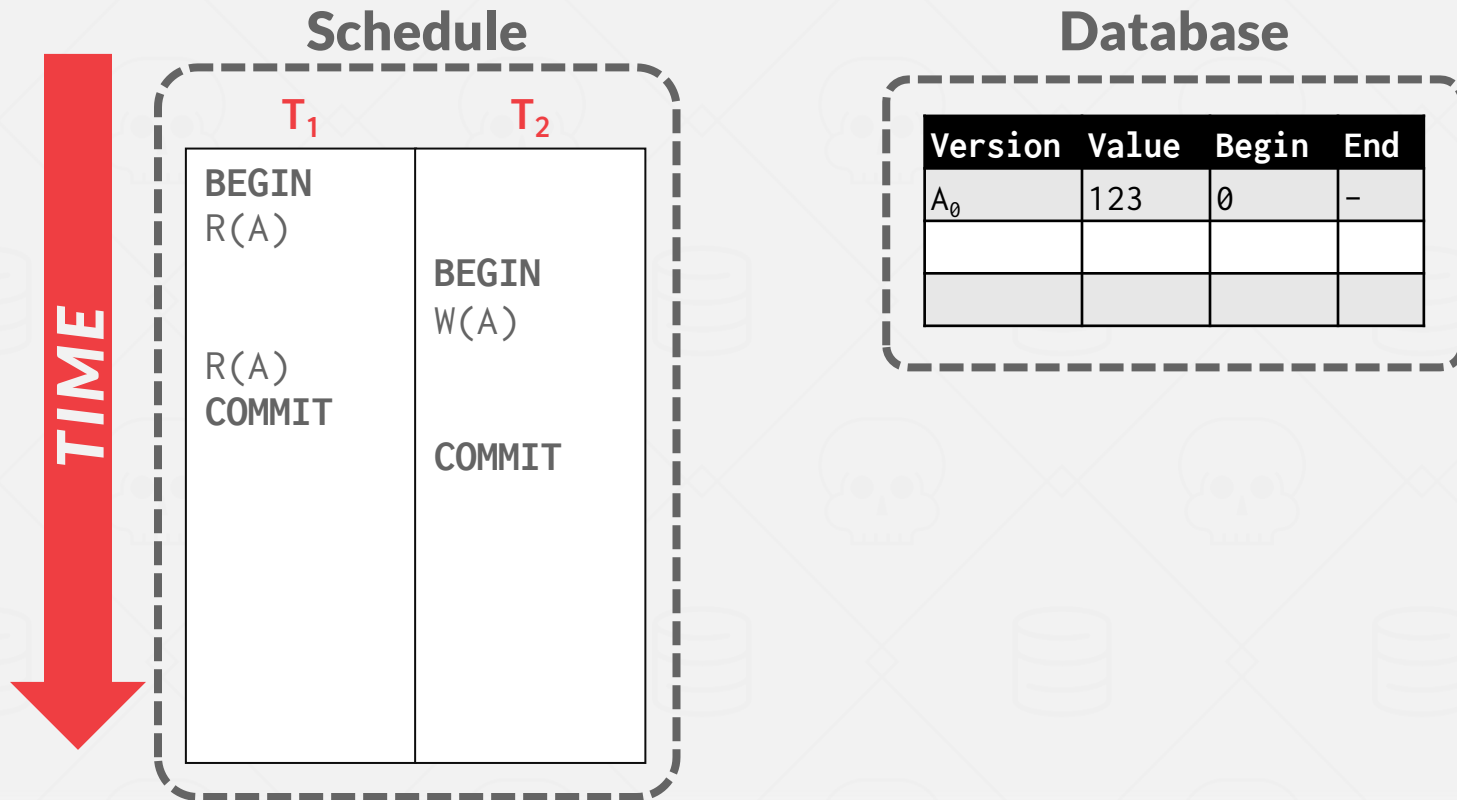
Database

Version	Value	Begin	End
A_0	123	\emptyset	-

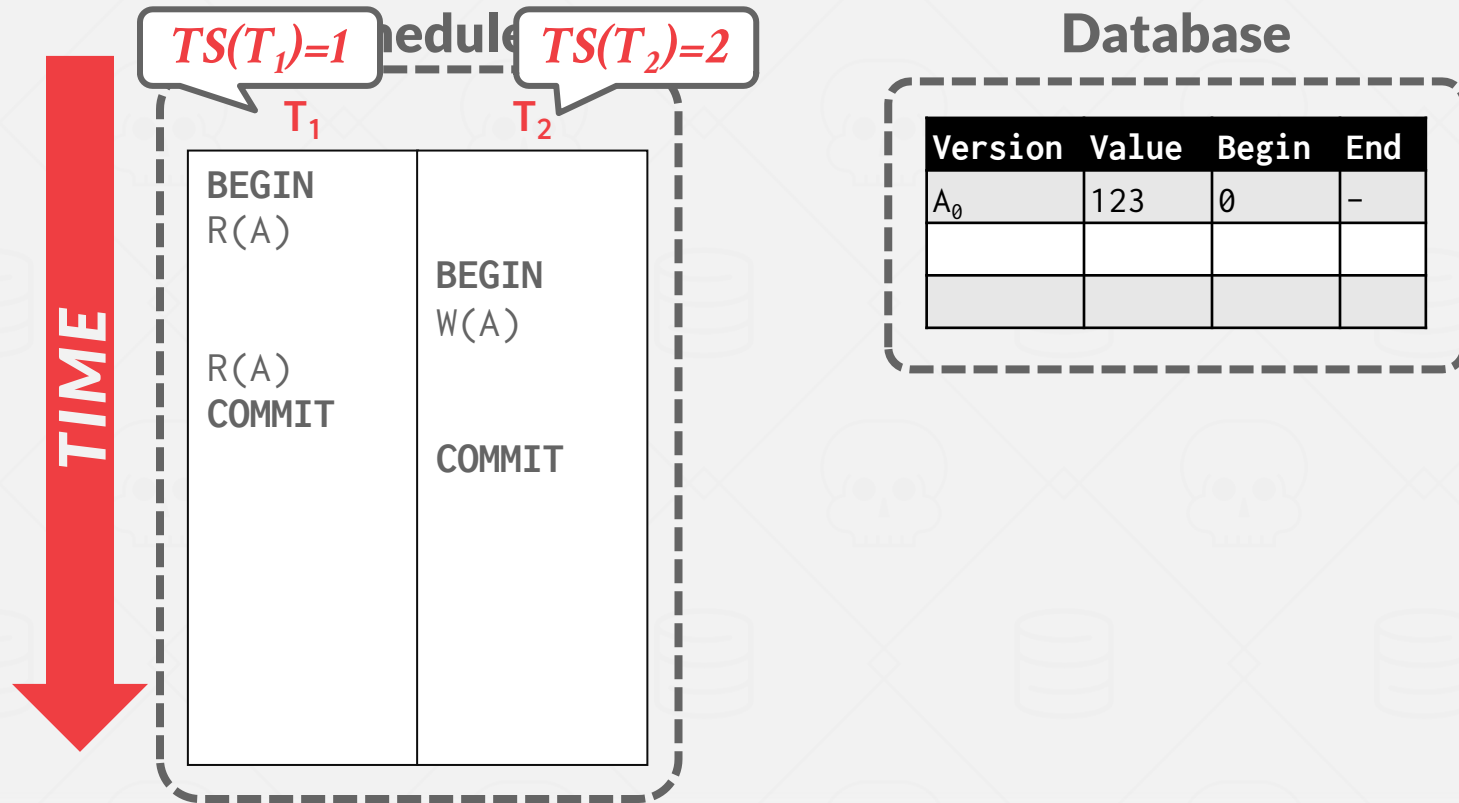
MVCC - EXAMPLE #1



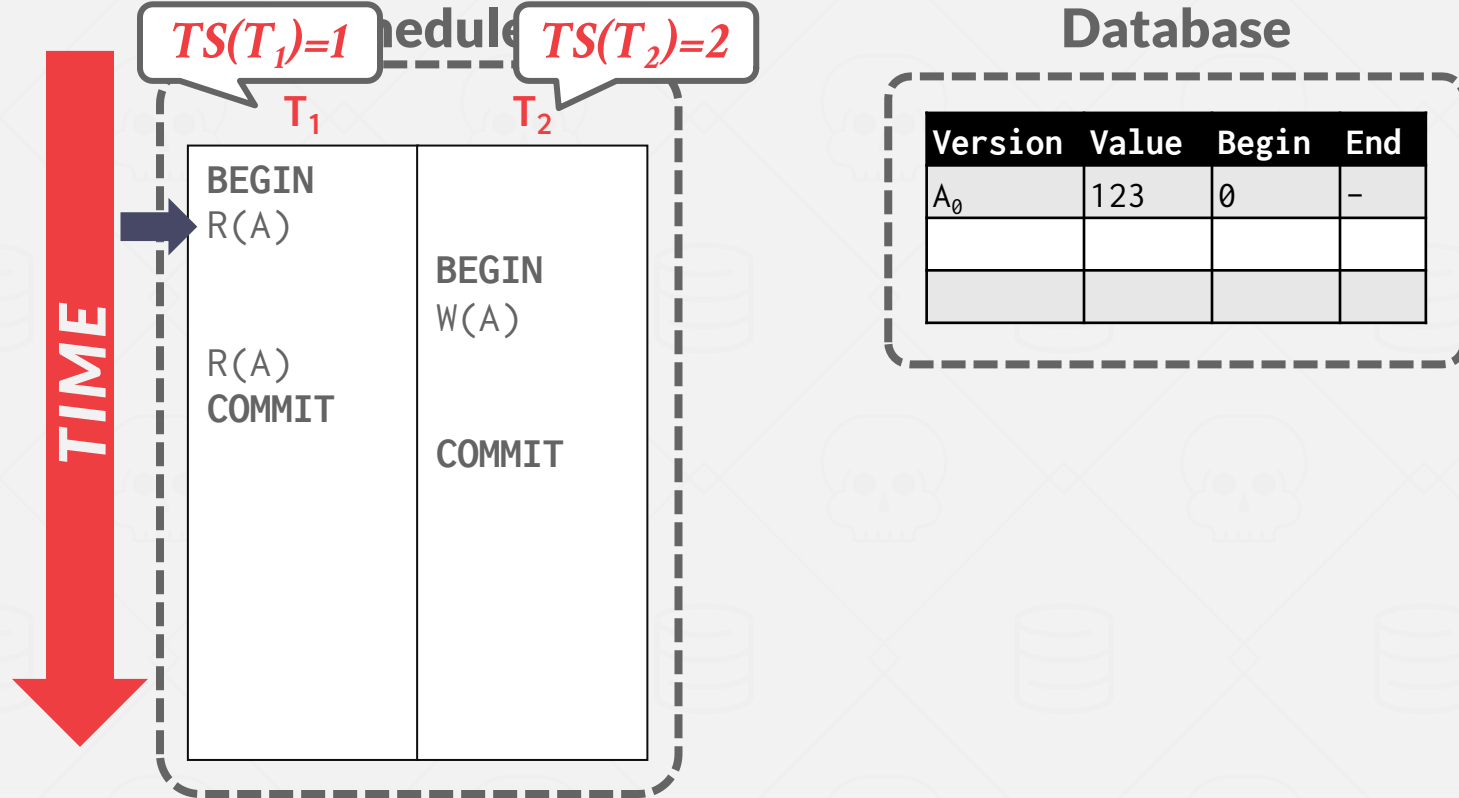
MVCC - EXAMPLE #1



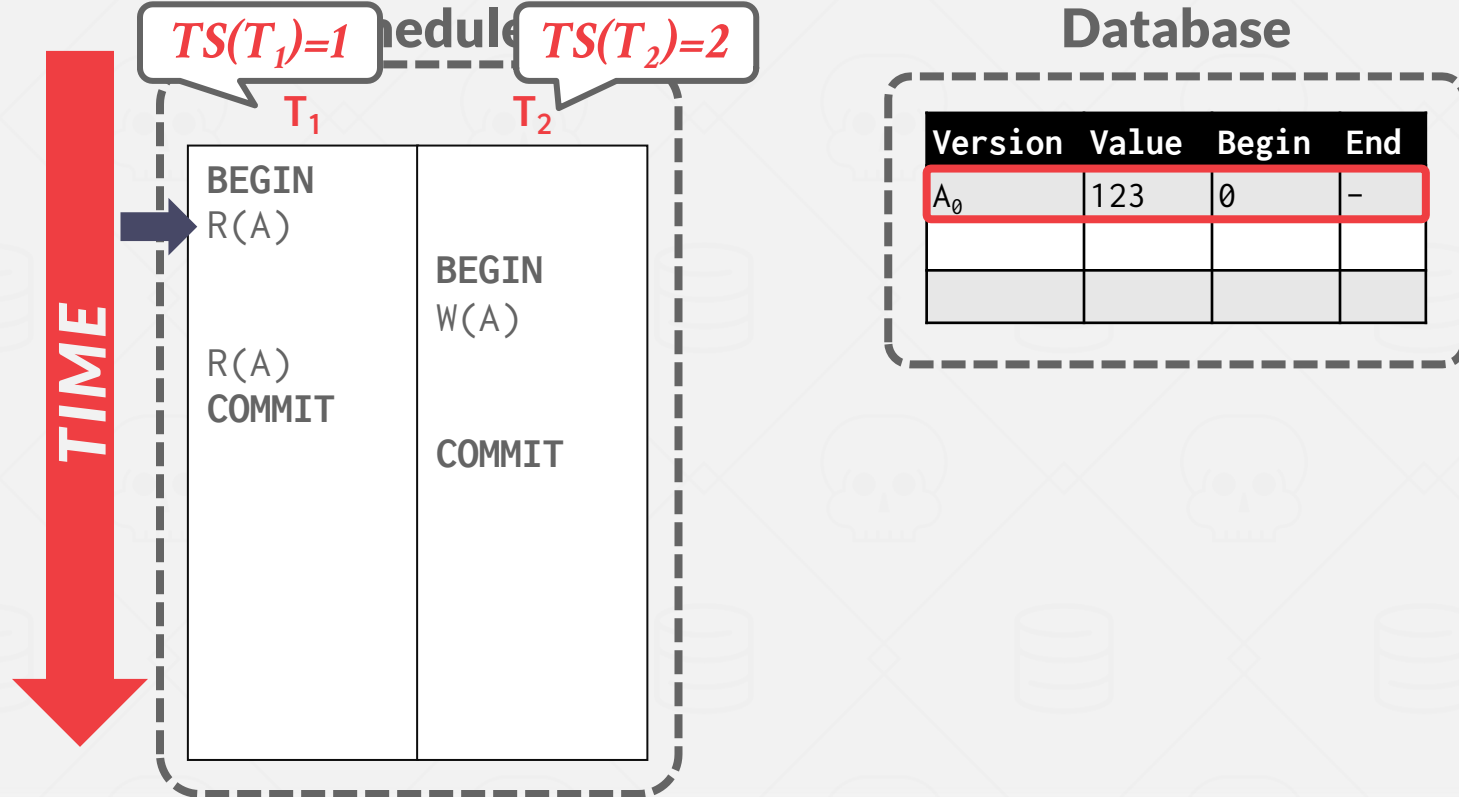
MVCC - EXAMPLE #1



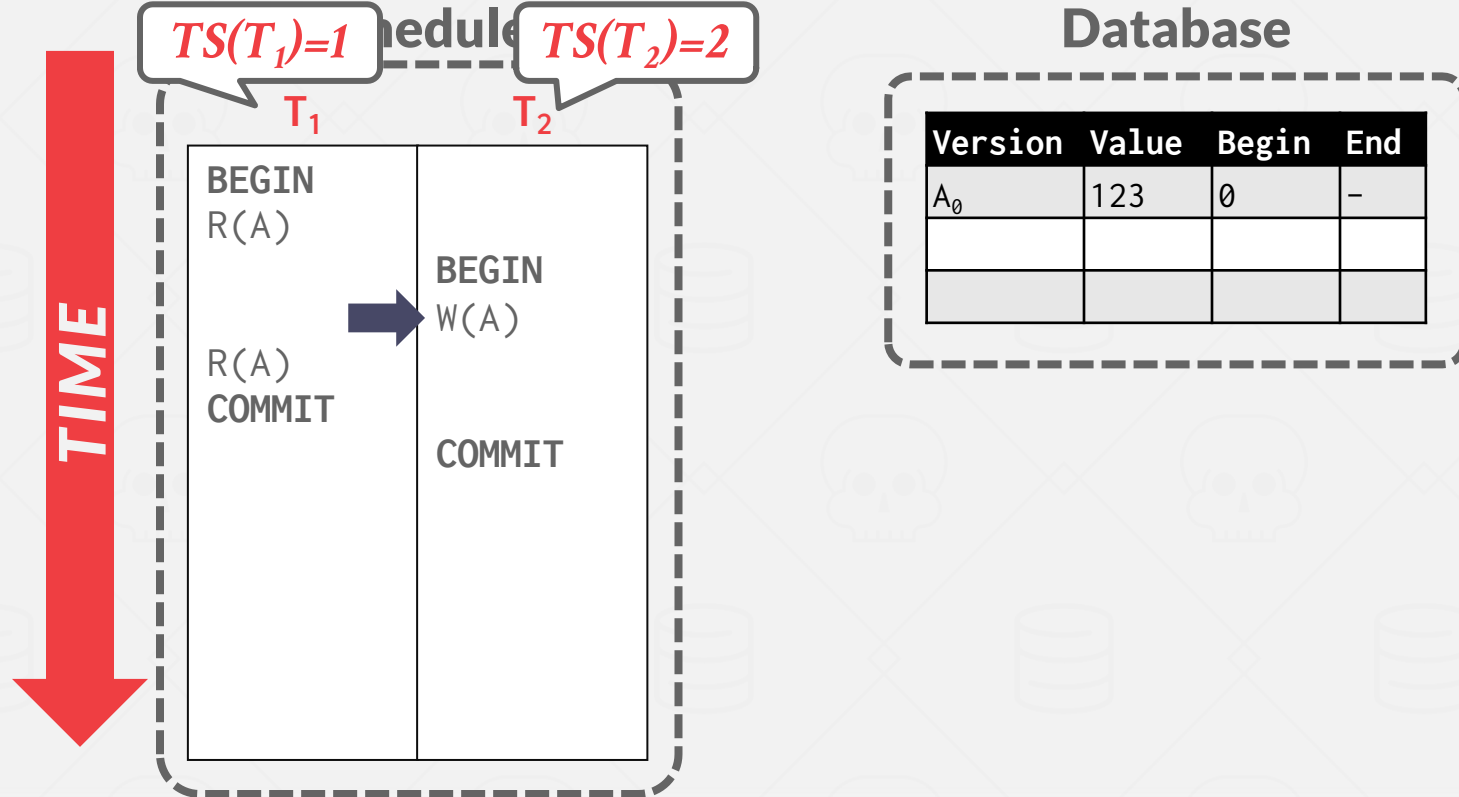
MVCC - EXAMPLE #1



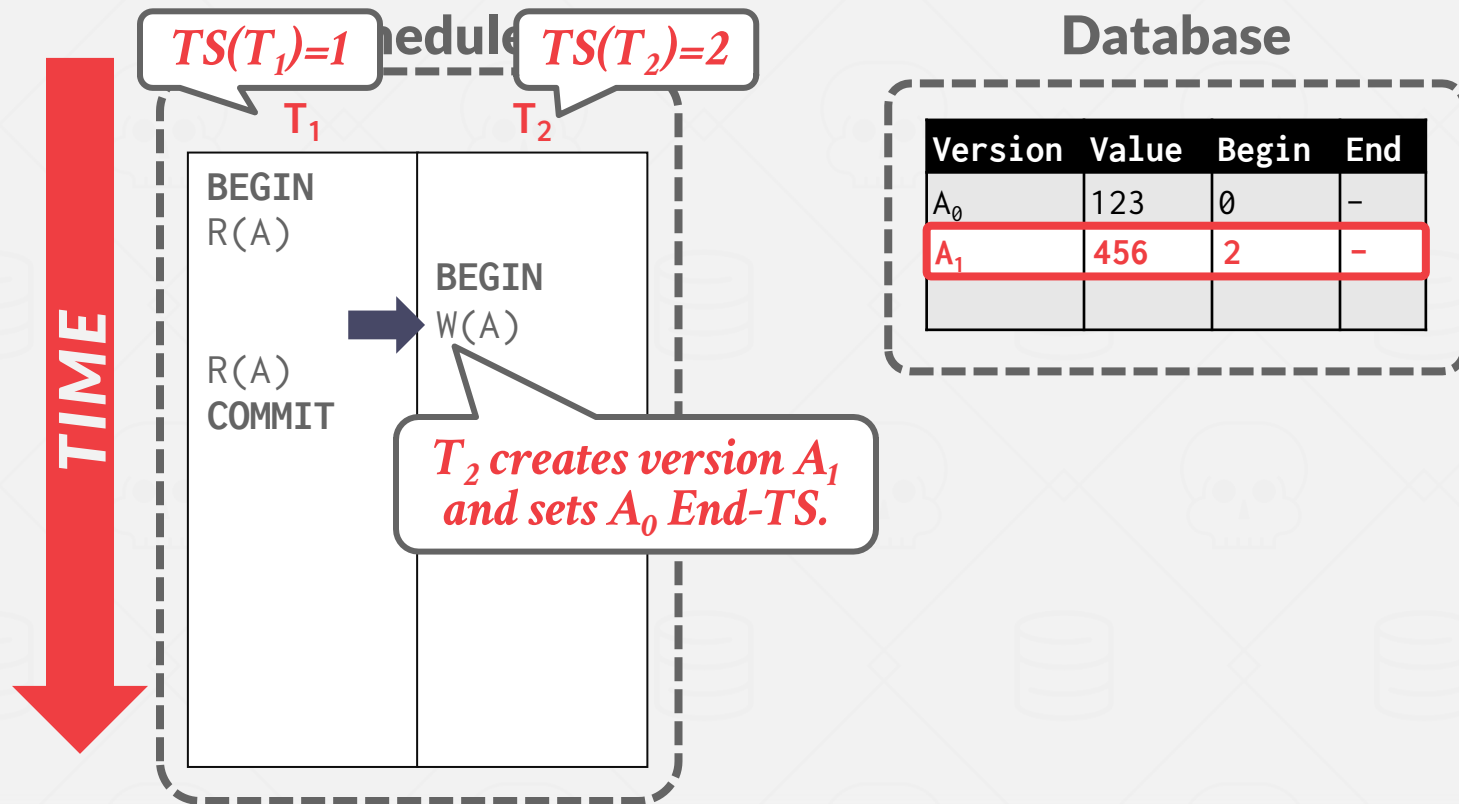
MVCC - EXAMPLE #1



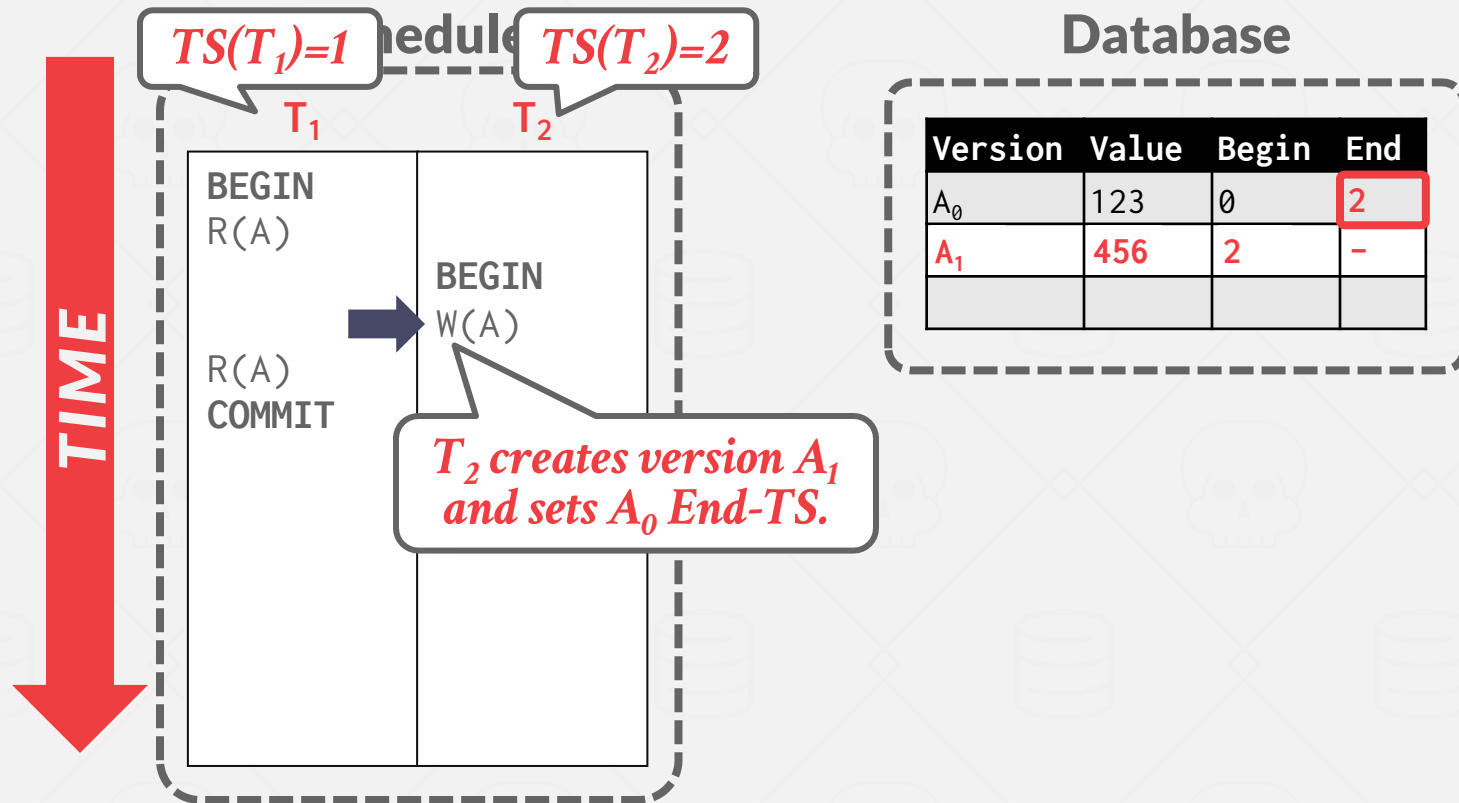
MVCC - EXAMPLE #1



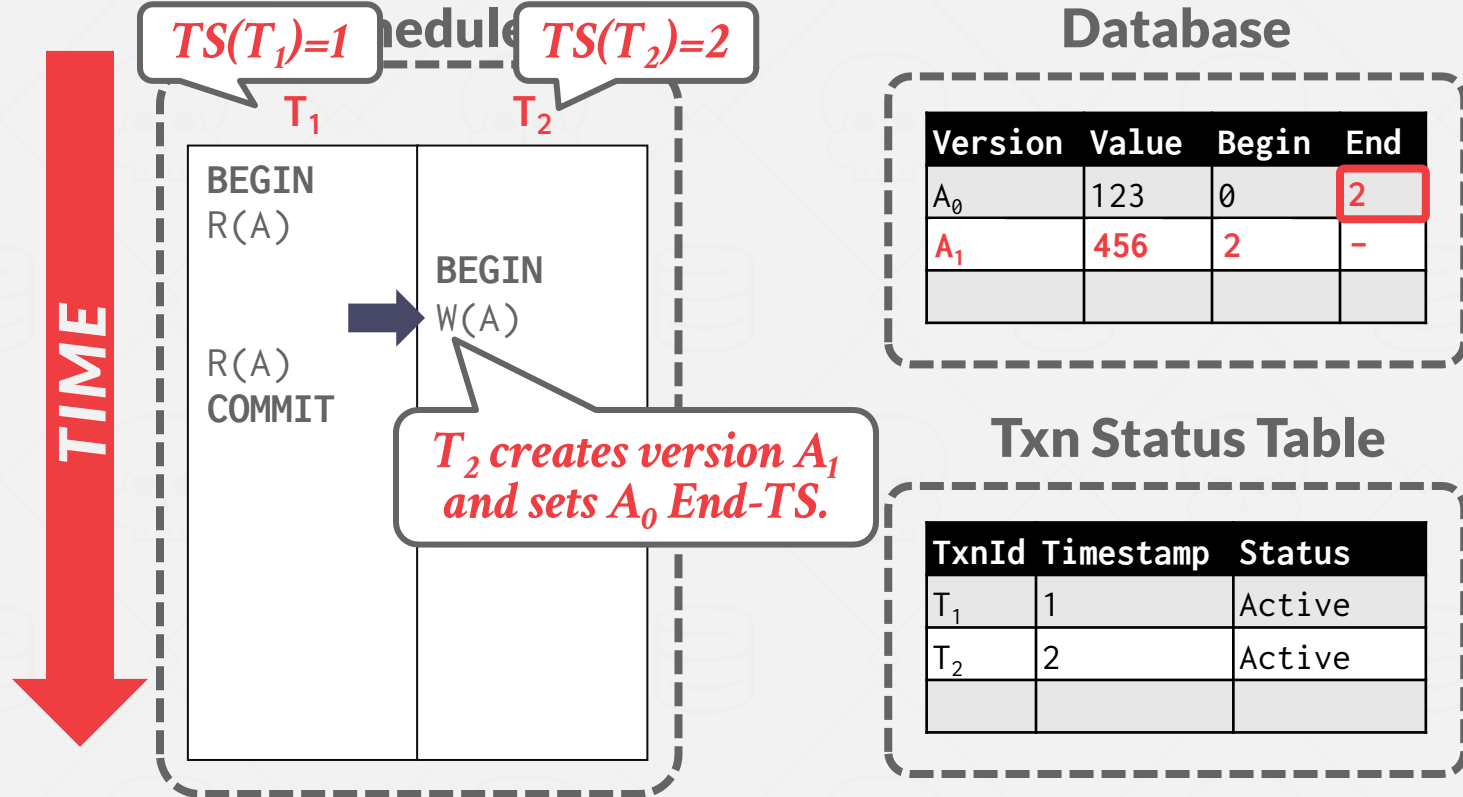
MVCC - EXAMPLE #1



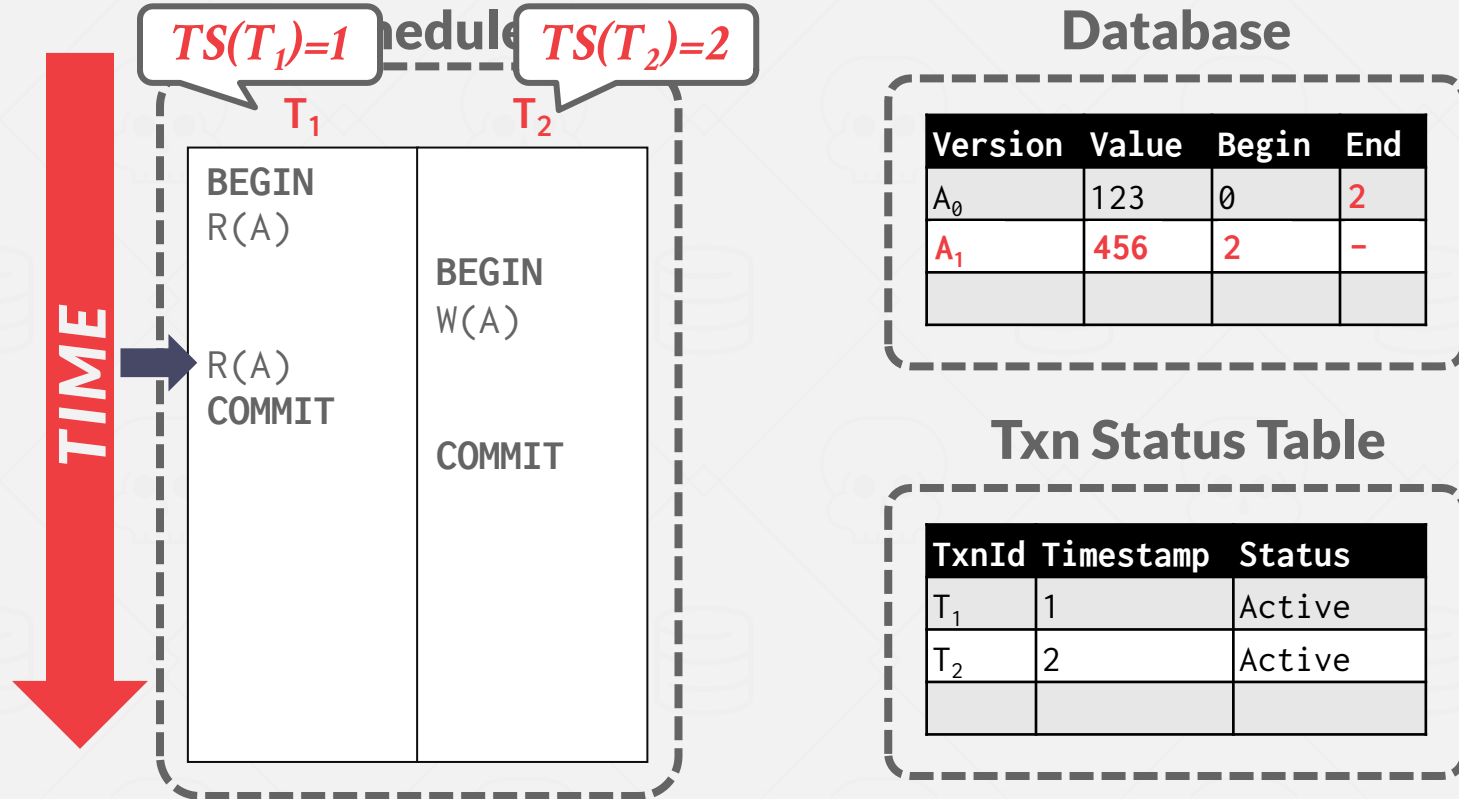
MVCC - EXAMPLE #1



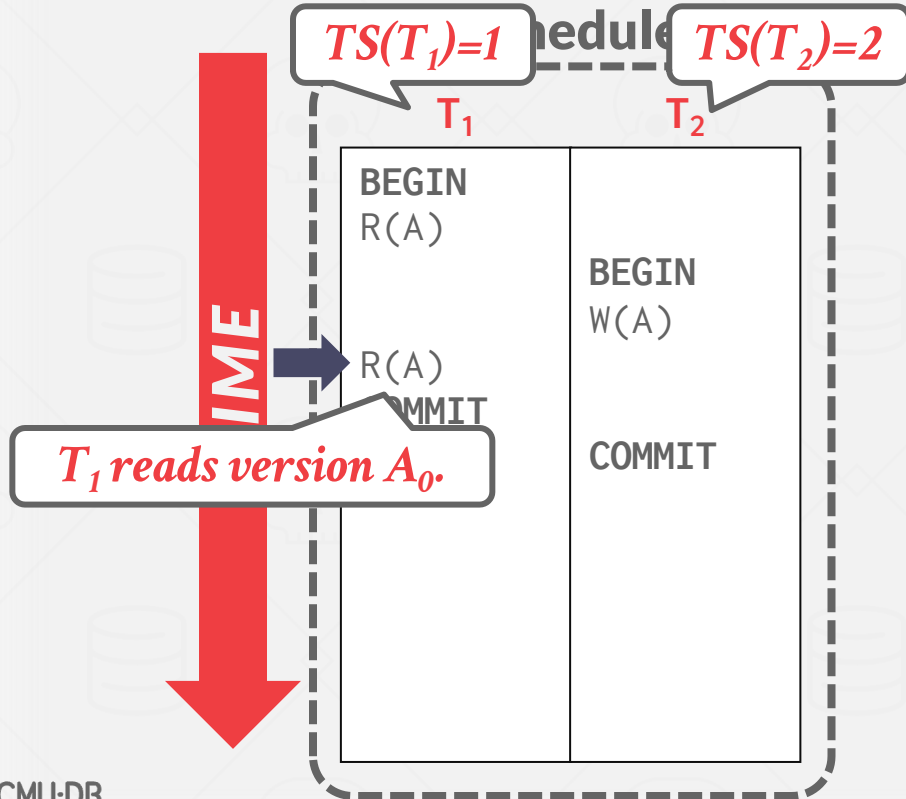
MVCC - EXAMPLE #1



MVCC - EXAMPLE #1



MVCC - EXAMPLE #1



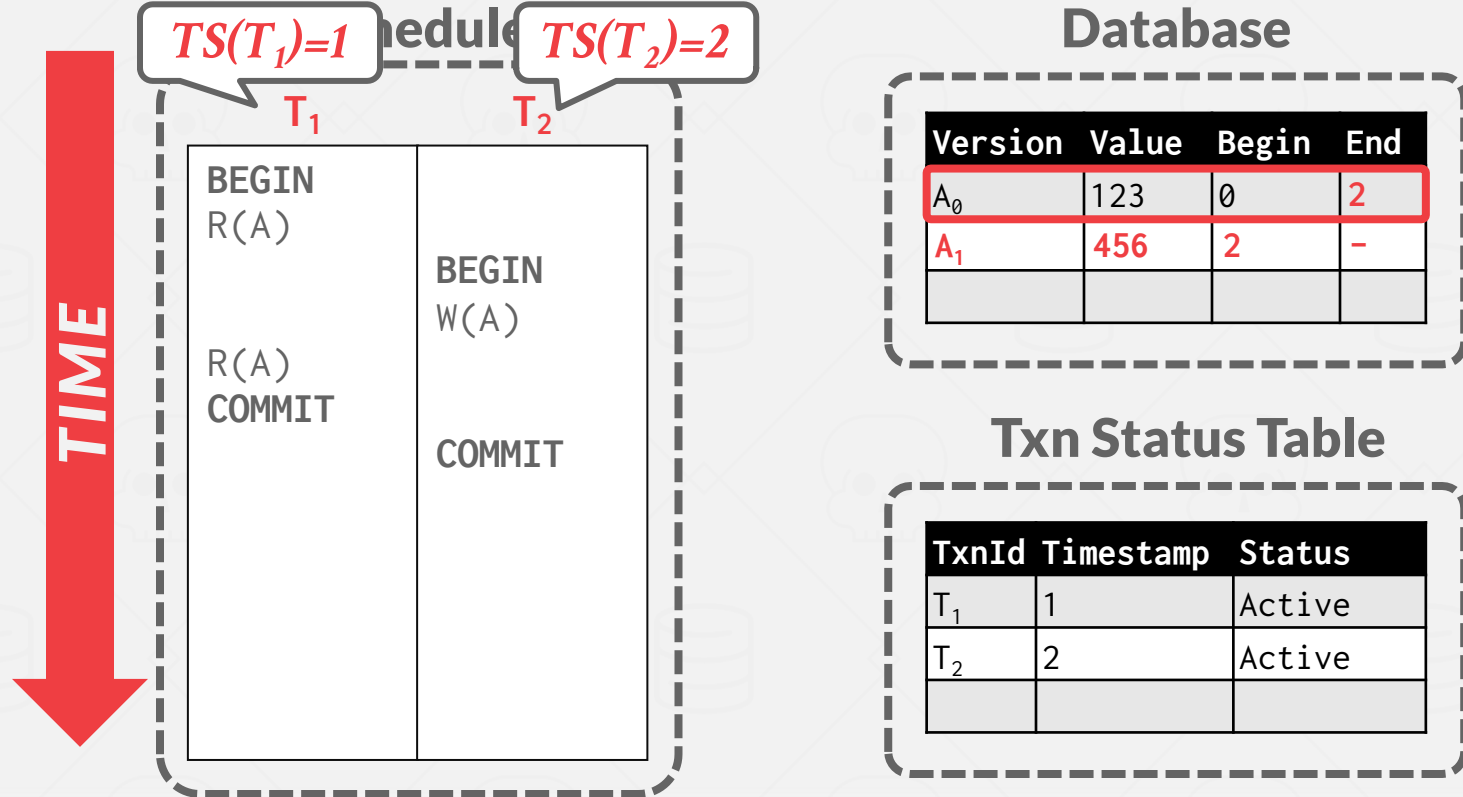
Database

Version	Value	Begin	End
A ₀	123	0	2
A ₁	456	2	-

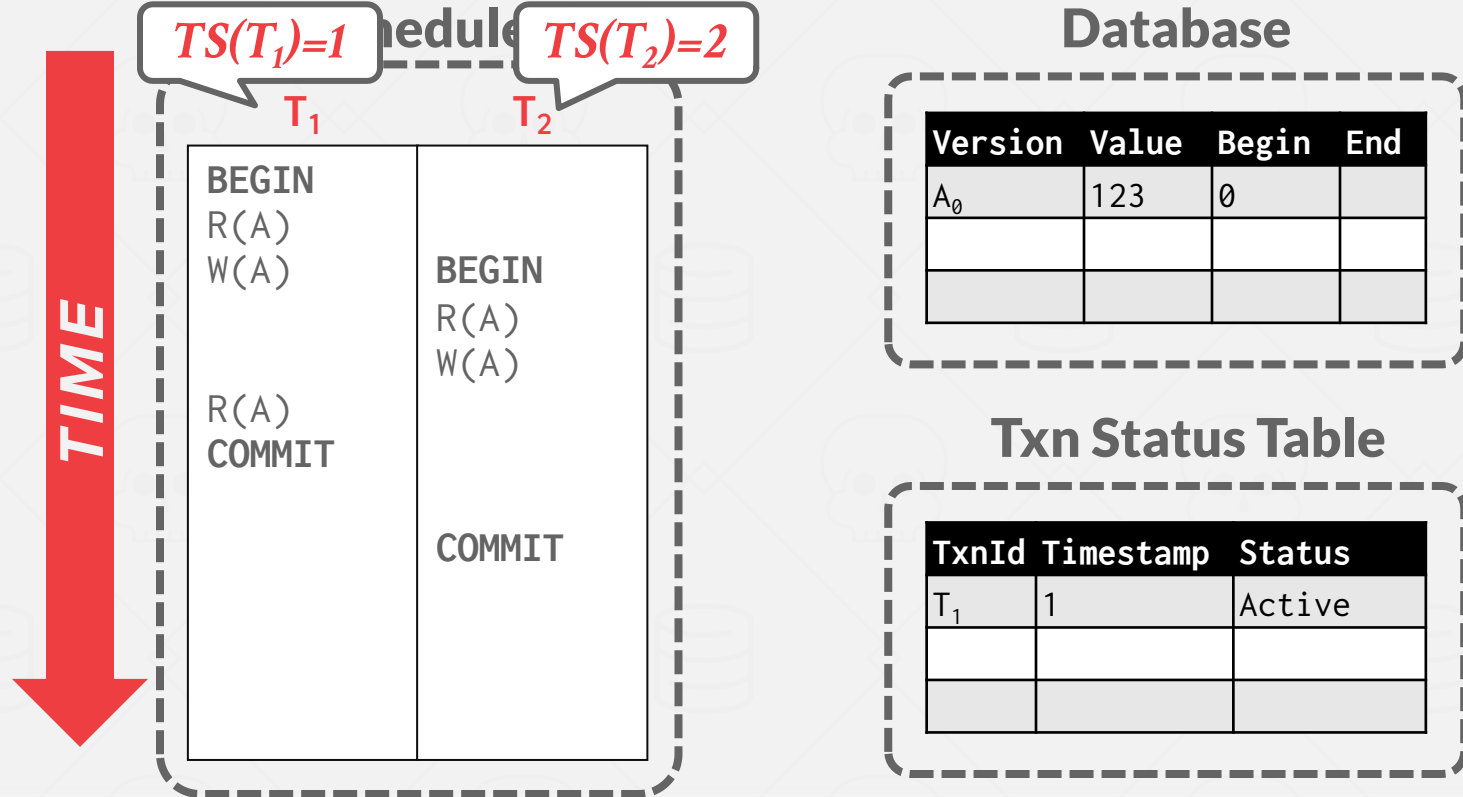
Txn Status Table

TxnId	Timestamp	Status
T ₁	1	Active
T ₂	2	Active

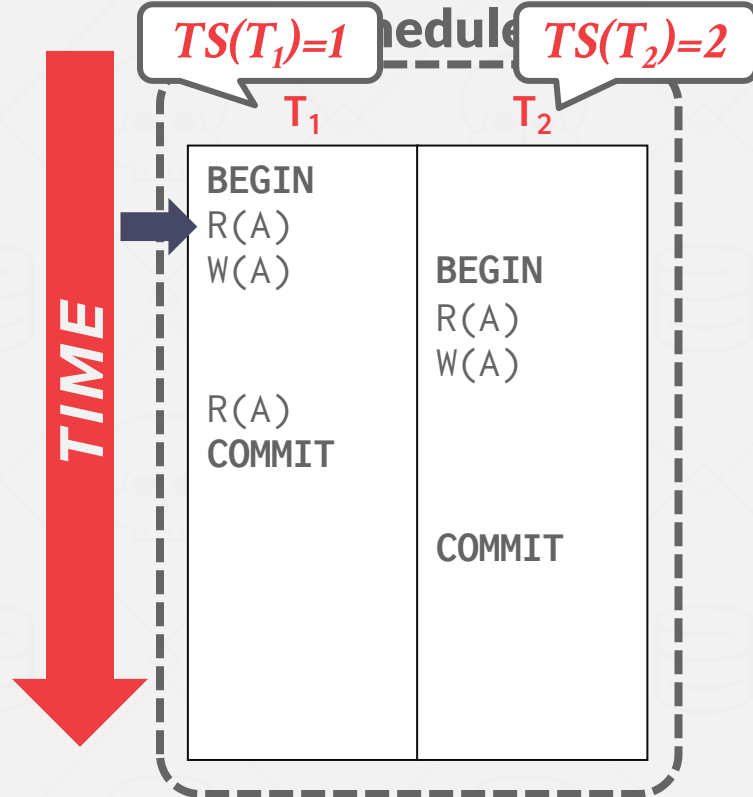
MVCC - EXAMPLE #1



MVCC - EXAMPLE #2



MVCC - EXAMPLE #2



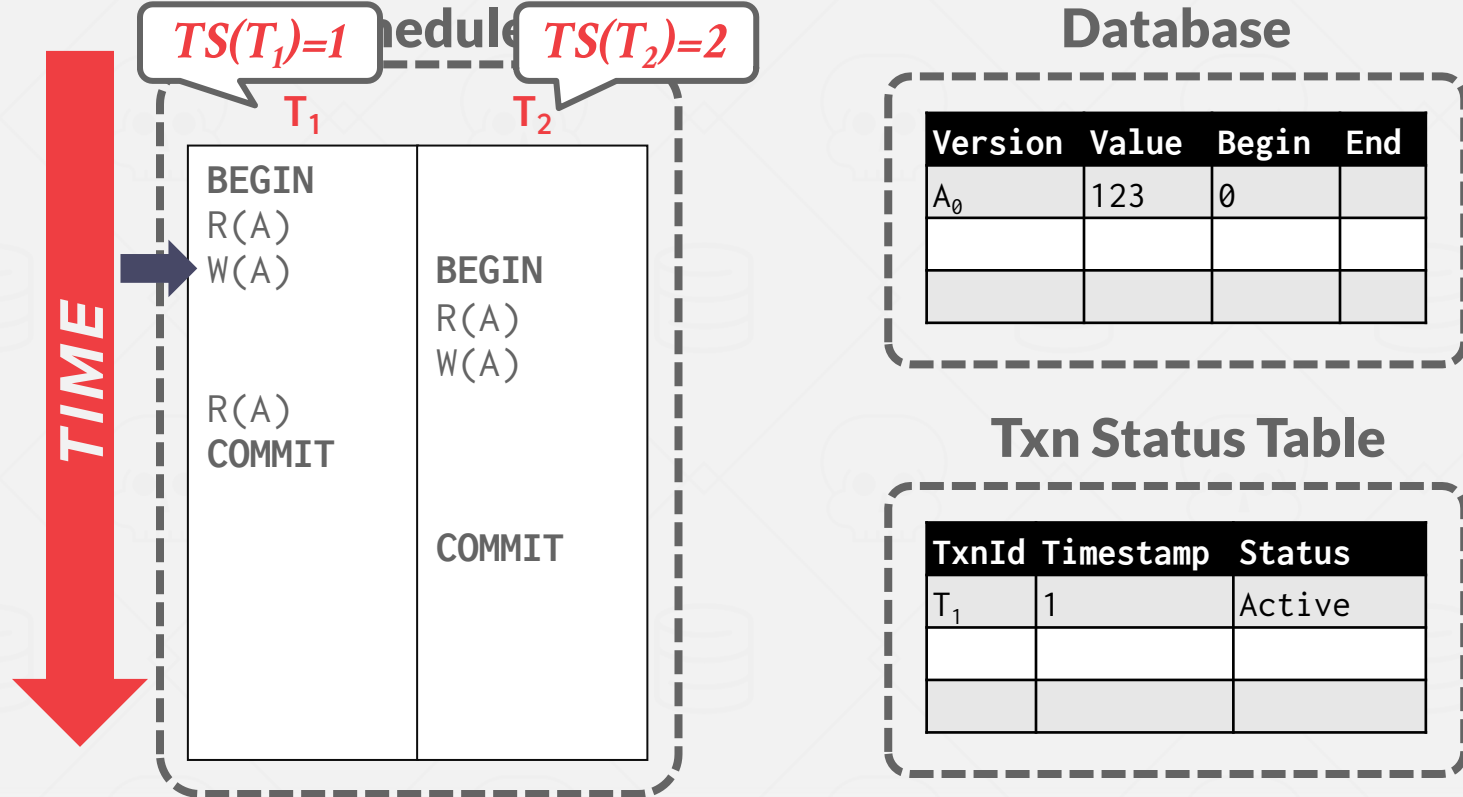
Database

Version	Value	Begin	End
A_0	123	0	

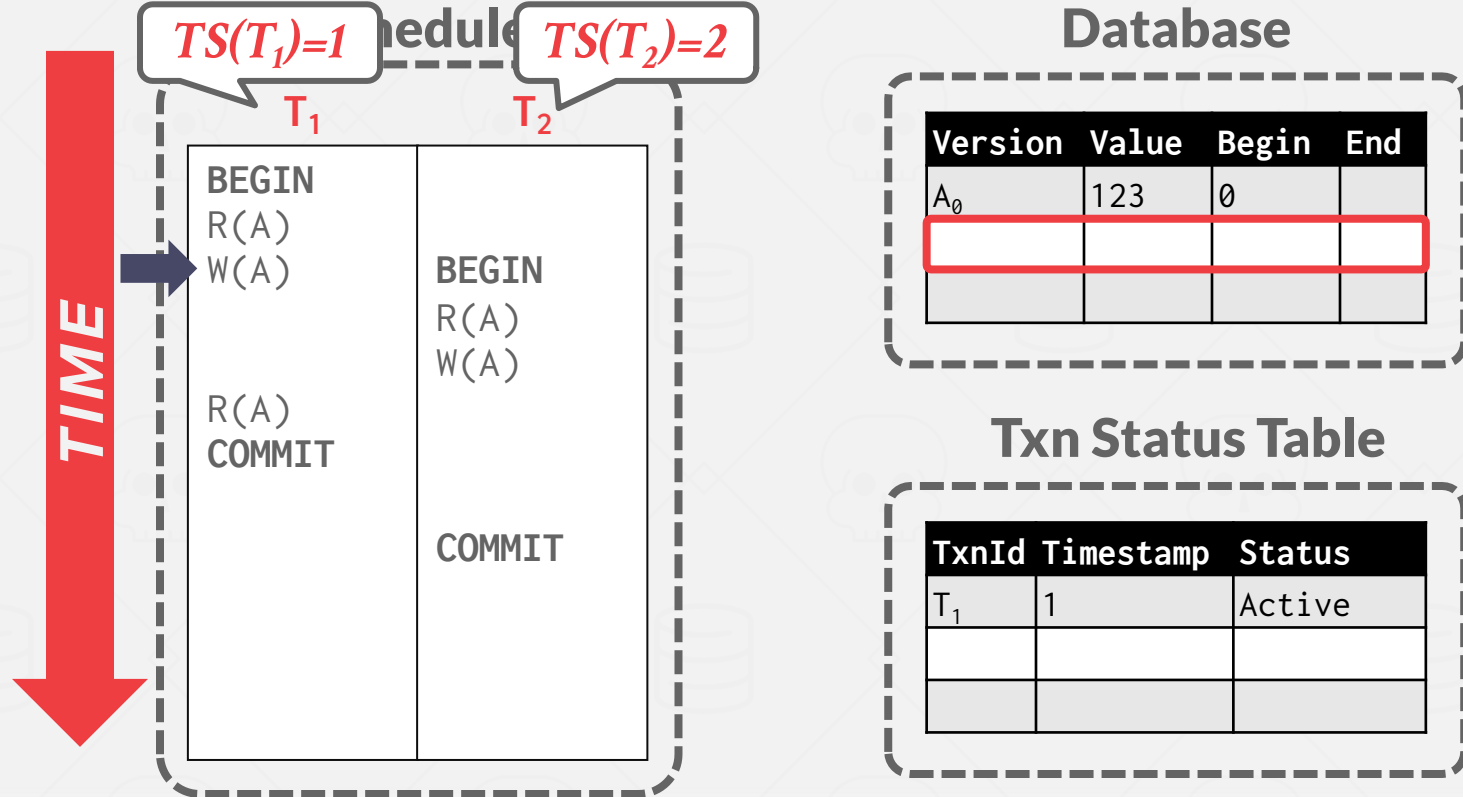
Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active

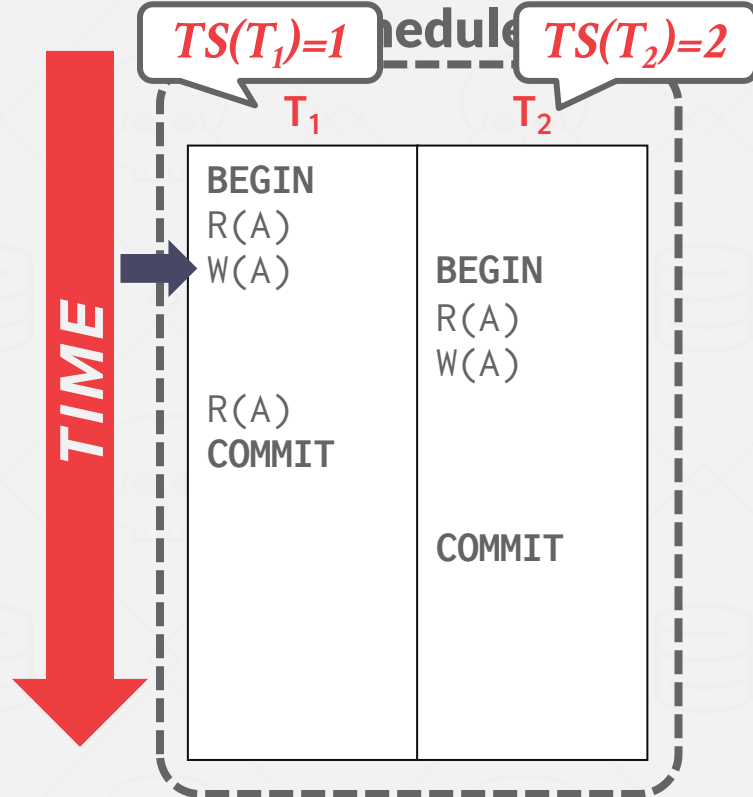
MVCC - EXAMPLE #2



MVCC - EXAMPLE #2



MVCC - EXAMPLE #2



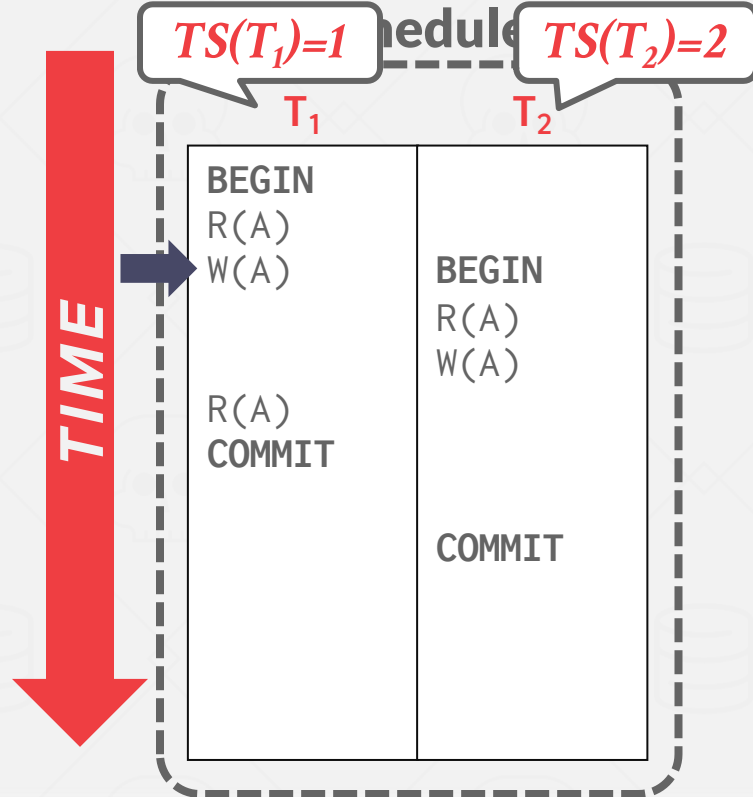
Database

Version	Value	Begin	End
A_0	123	0	
A_1	456	1	-

Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active

MVCC - EXAMPLE #2



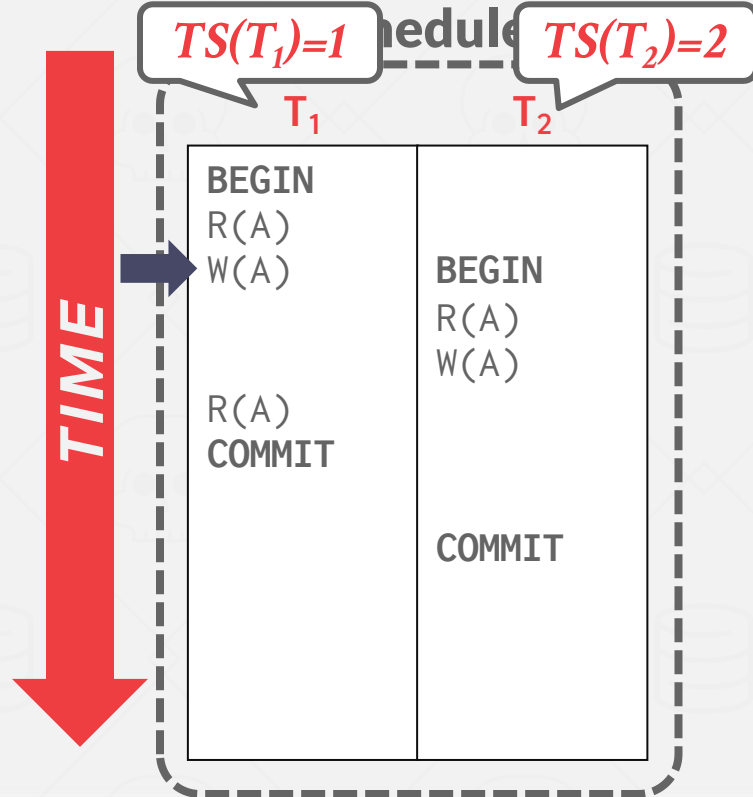
Database

Version	Value	Begin	End
A_0	123	0	
A_1	456	1	-

Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active

MVCC - EXAMPLE #2



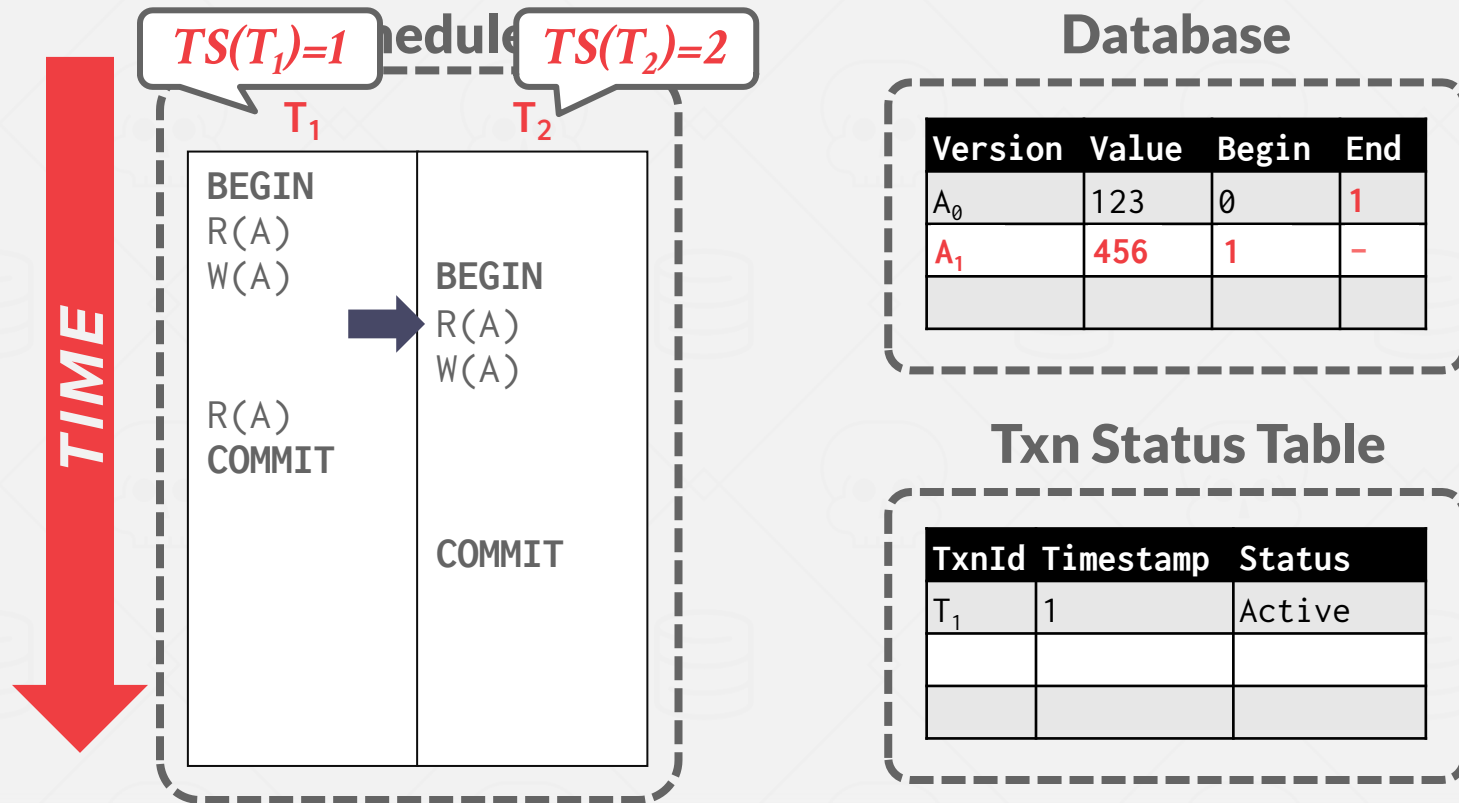
Database

Version	Value	Begin	End
A_0	123	0	1
A_1	456	1	-

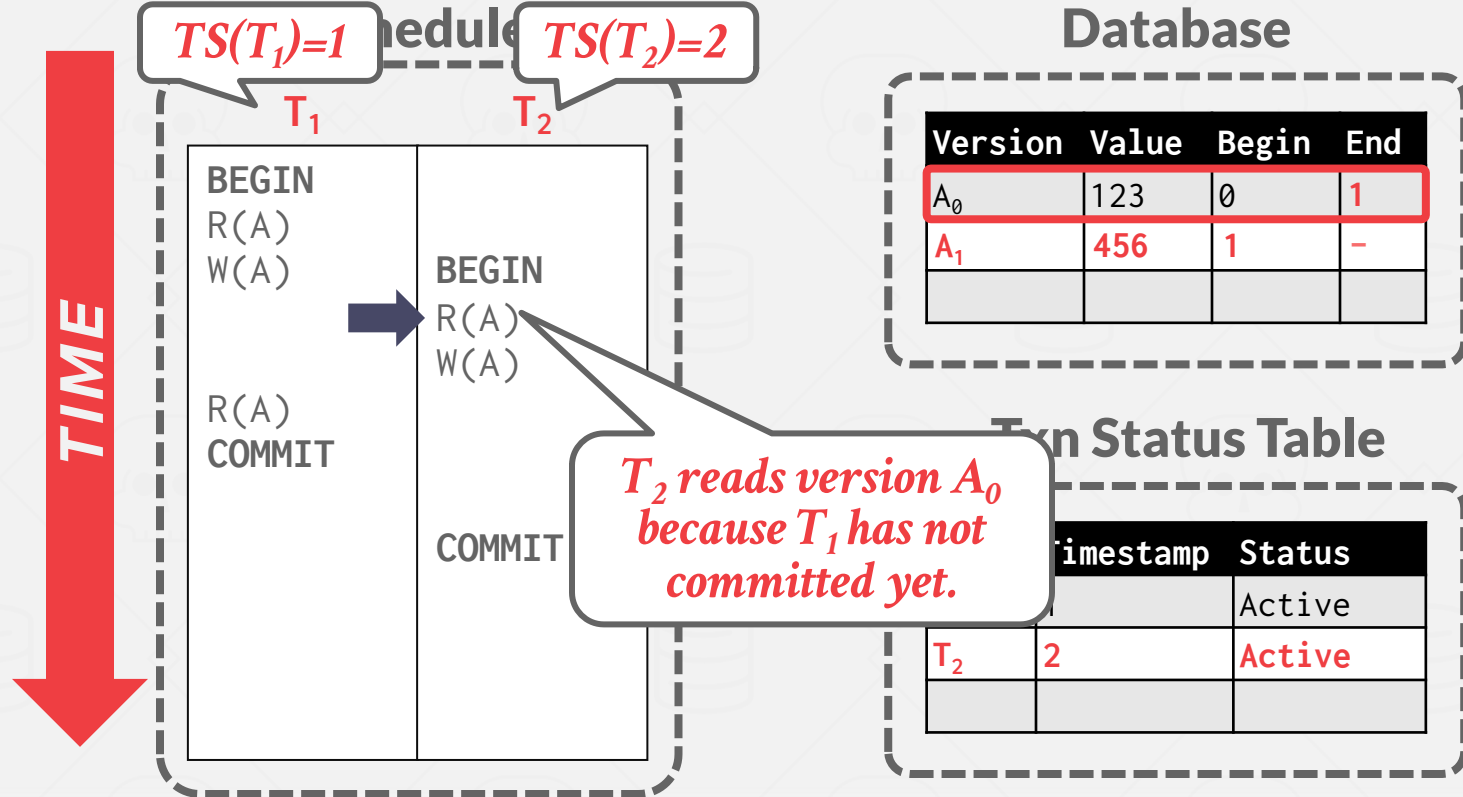
Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active

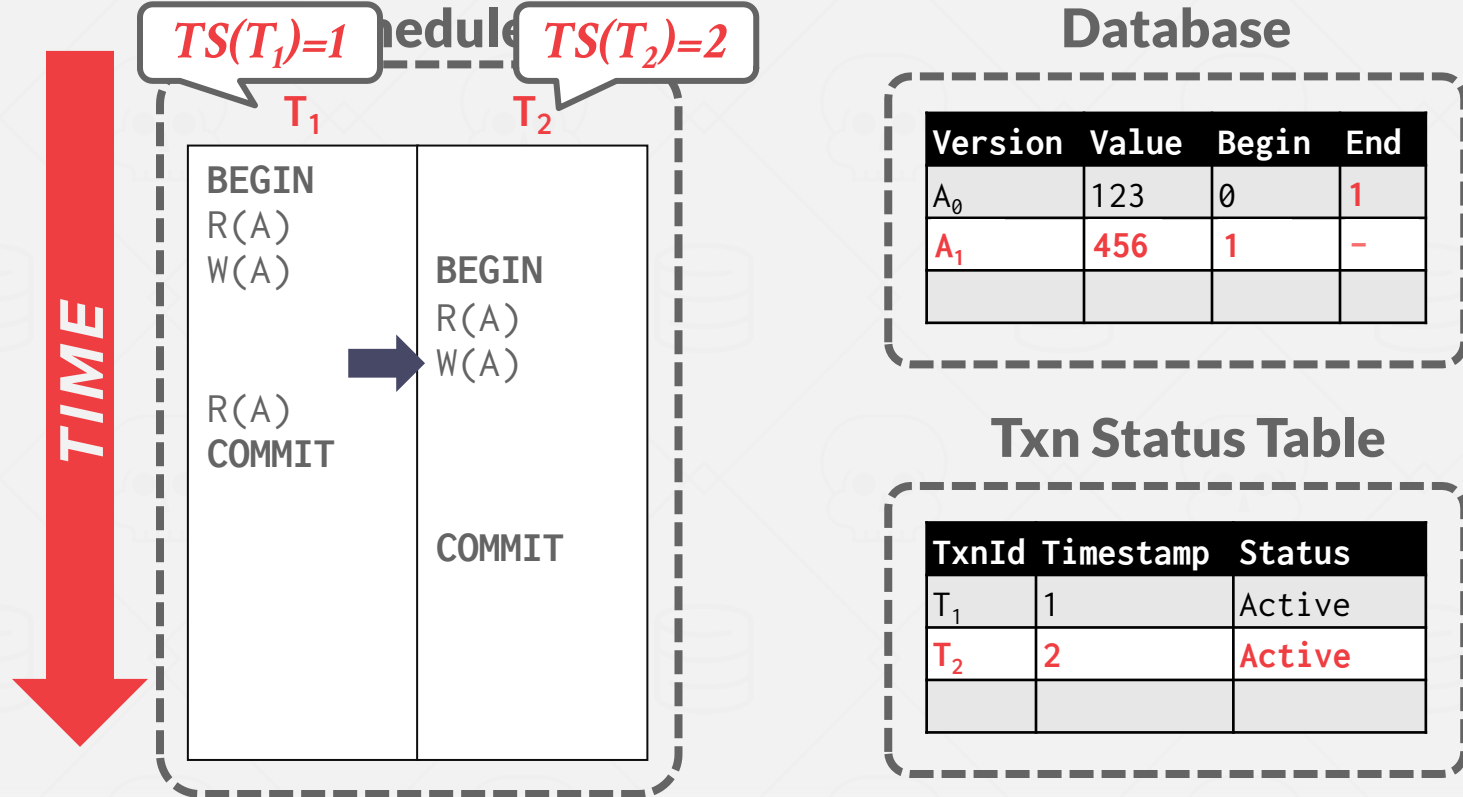
MVCC - EXAMPLE #2



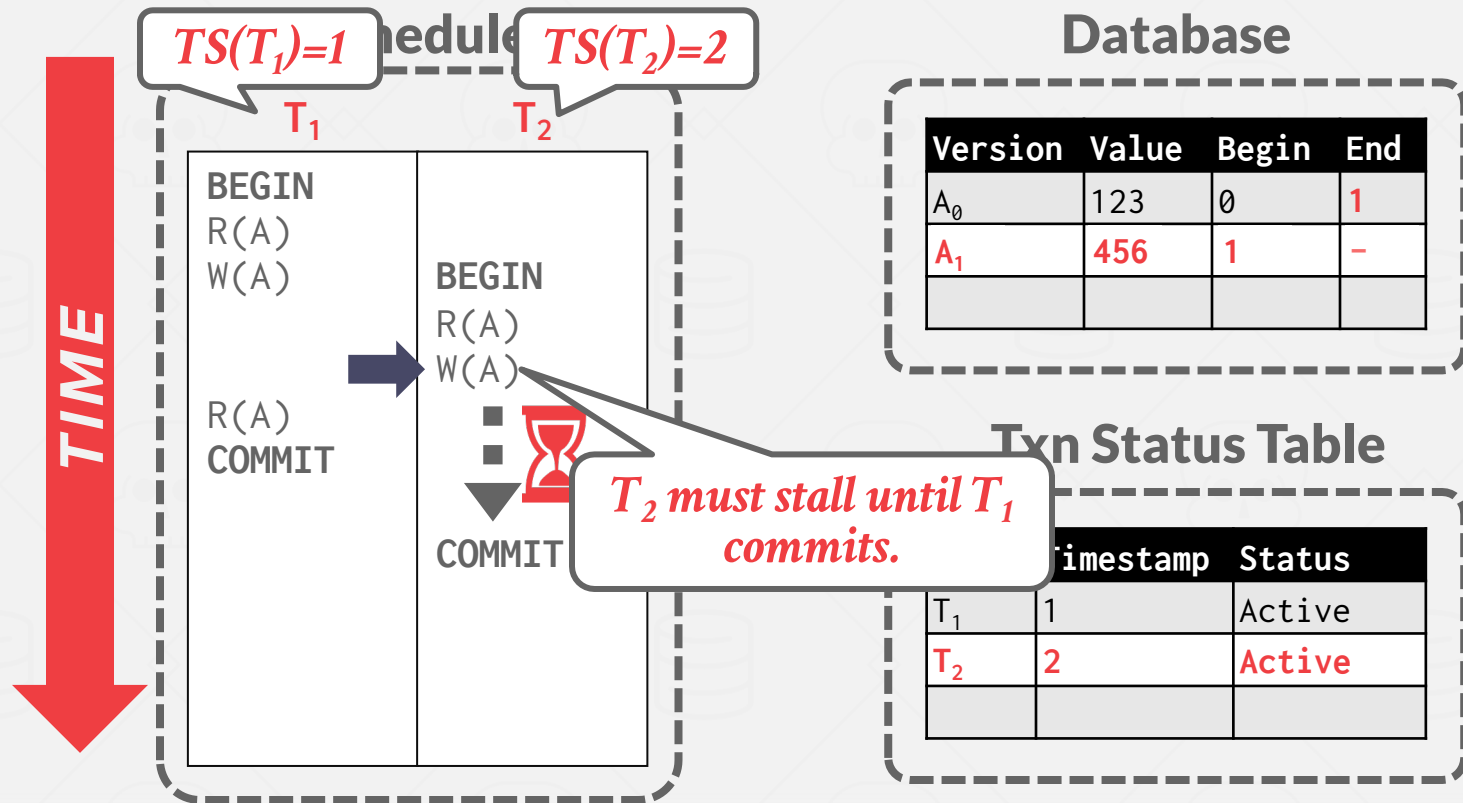
MVCC - EXAMPLE #2



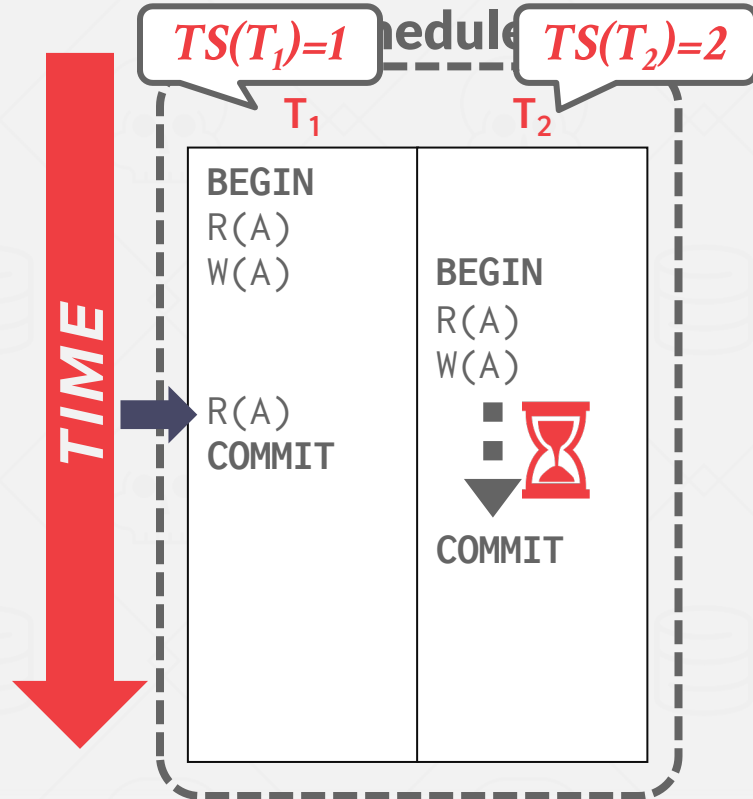
MVCC - EXAMPLE #2



MVCC - EXAMPLE #2



MVCC - EXAMPLE #2



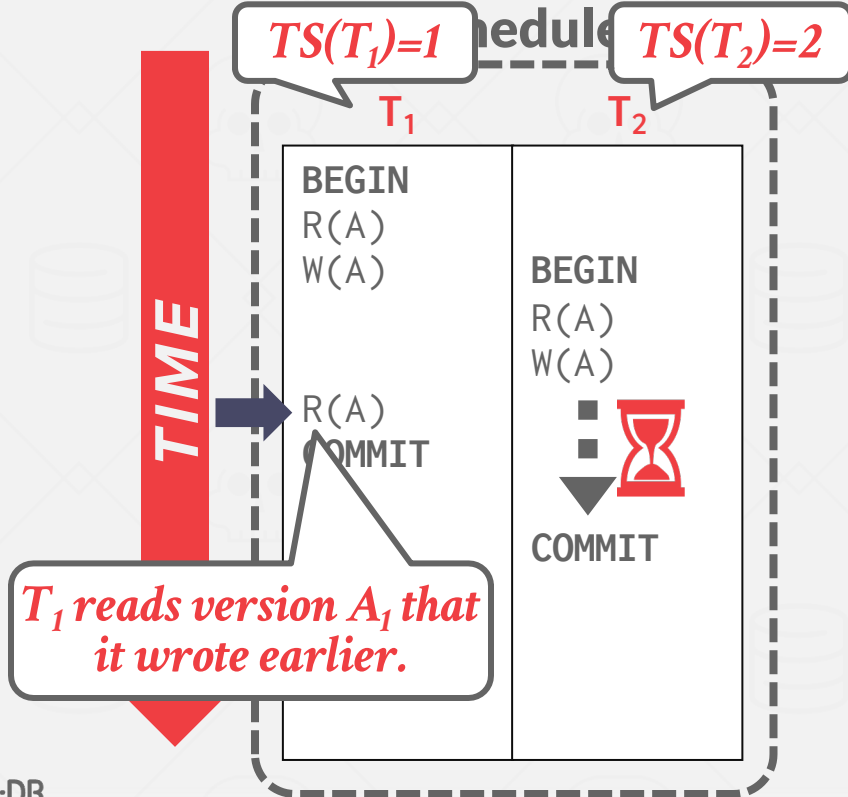
Database

Version	Value	Begin	End
A_0	123	0	1
A_1	456	1	-

Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active
T_2	2	Active

MVCC - EXAMPLE #2



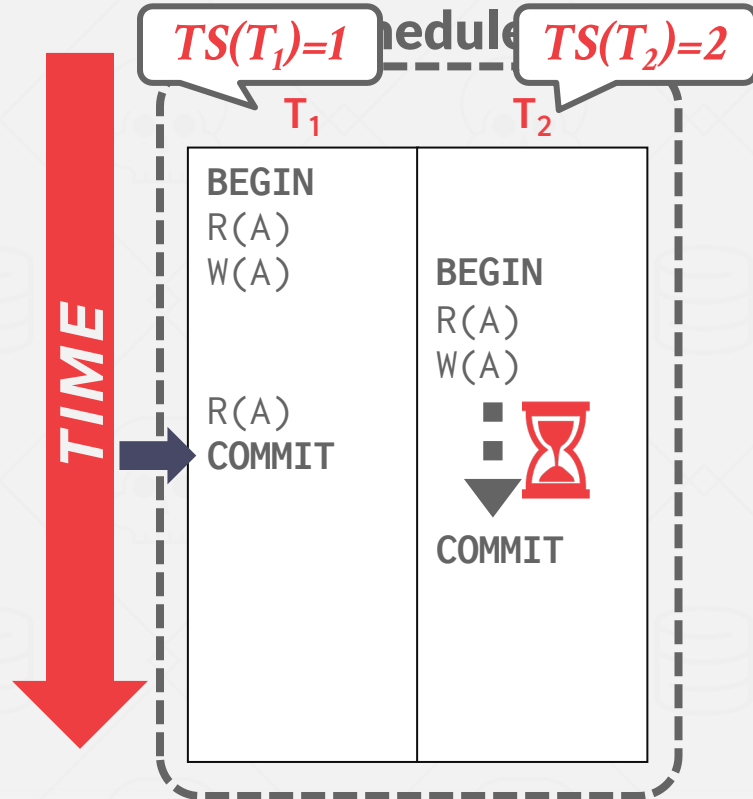
Database

Version	Value	Begin	End
A ₀	123	0	1
A ₁	456	1	-

Txn Status Table

TxnId	Timestamp	Status
T ₁	1	Active
T ₂	2	Active

MVCC - EXAMPLE #2



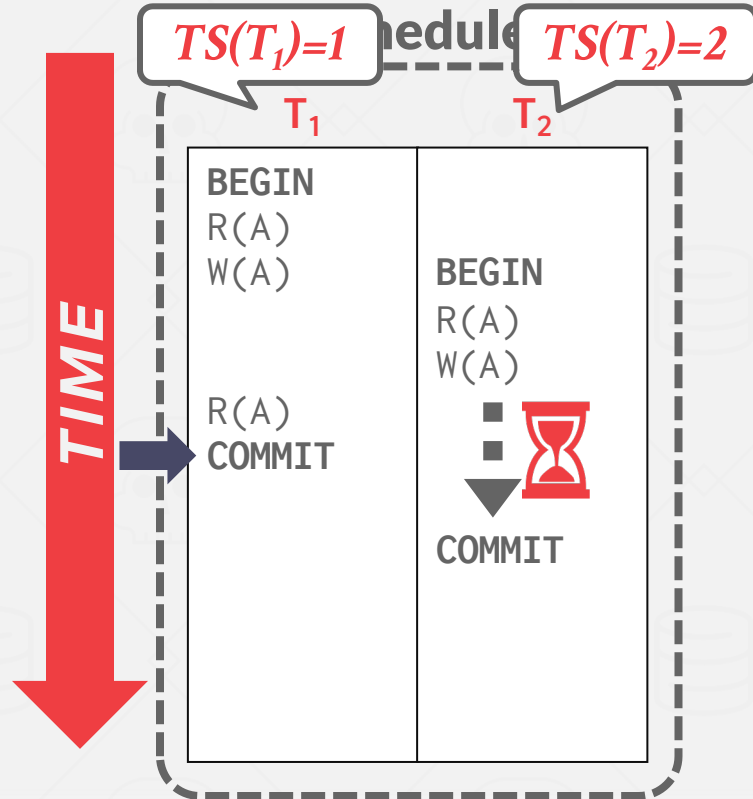
Database

Version	Value	Begin	End
A_0	123	0	1
A_1	456	1	-

Txn Status Table

TxnId	Timestamp	Status
T_1	1	Active
T_2	2	Active

MVCC - EXAMPLE #2



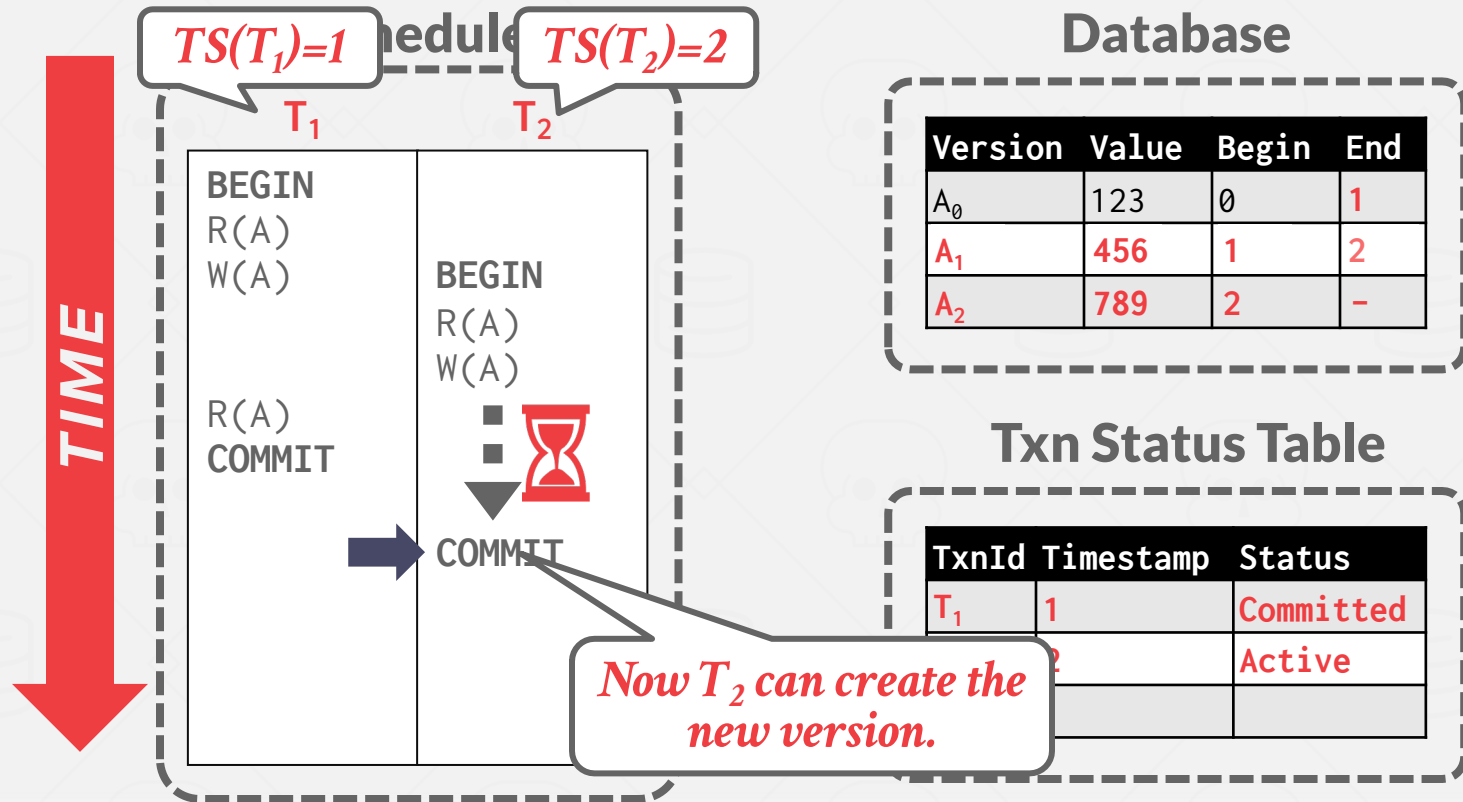
Database

Version	Value	Begin	End
A_0	123	0	1
A_1	456	1	-

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.

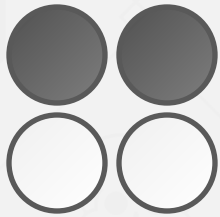
→ If two txns update the same object, then first writer wins.

SI is susceptible to the Write Skew Anomaly.

WRITE SKEW ANOMALY

Txn #1

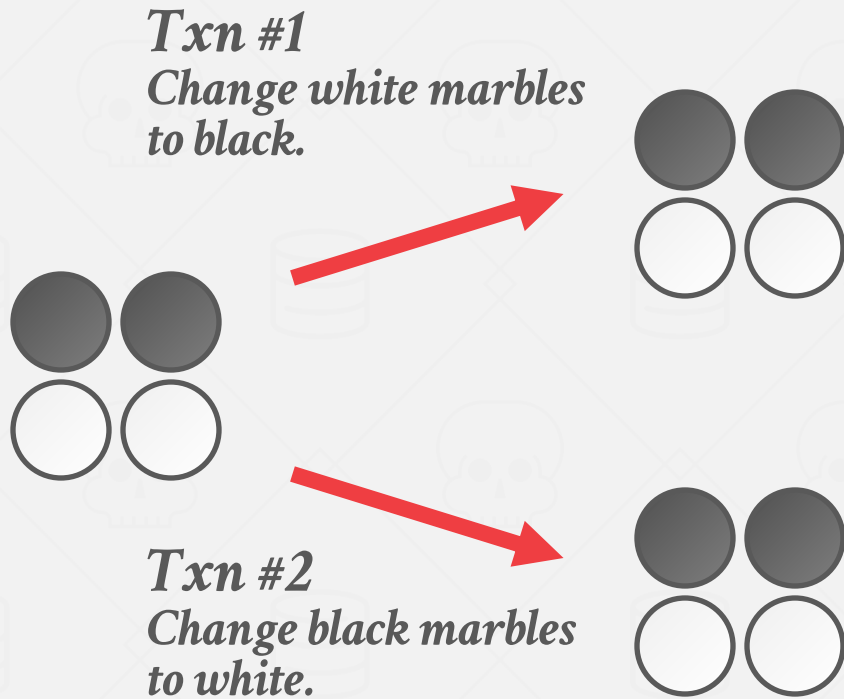
*Change white marbles
to black.*



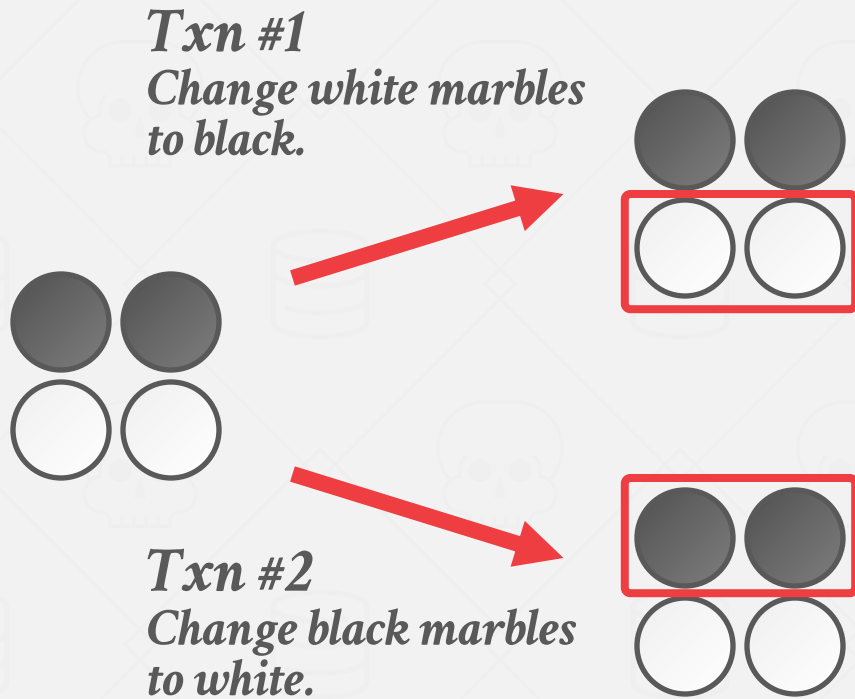
Txn #2

*Change black marbles
to white.*

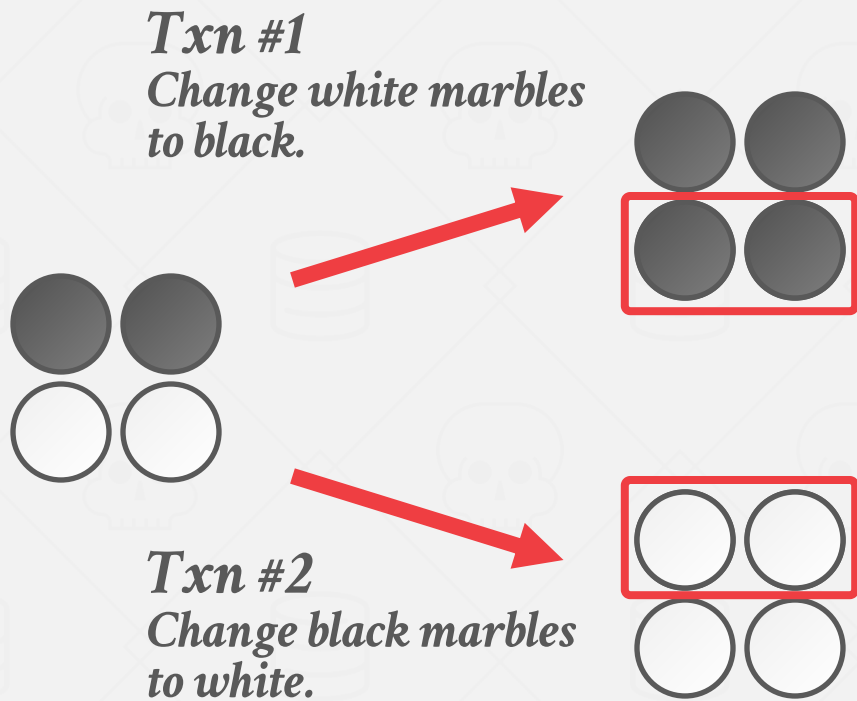
WRITE SKEW ANOMALY



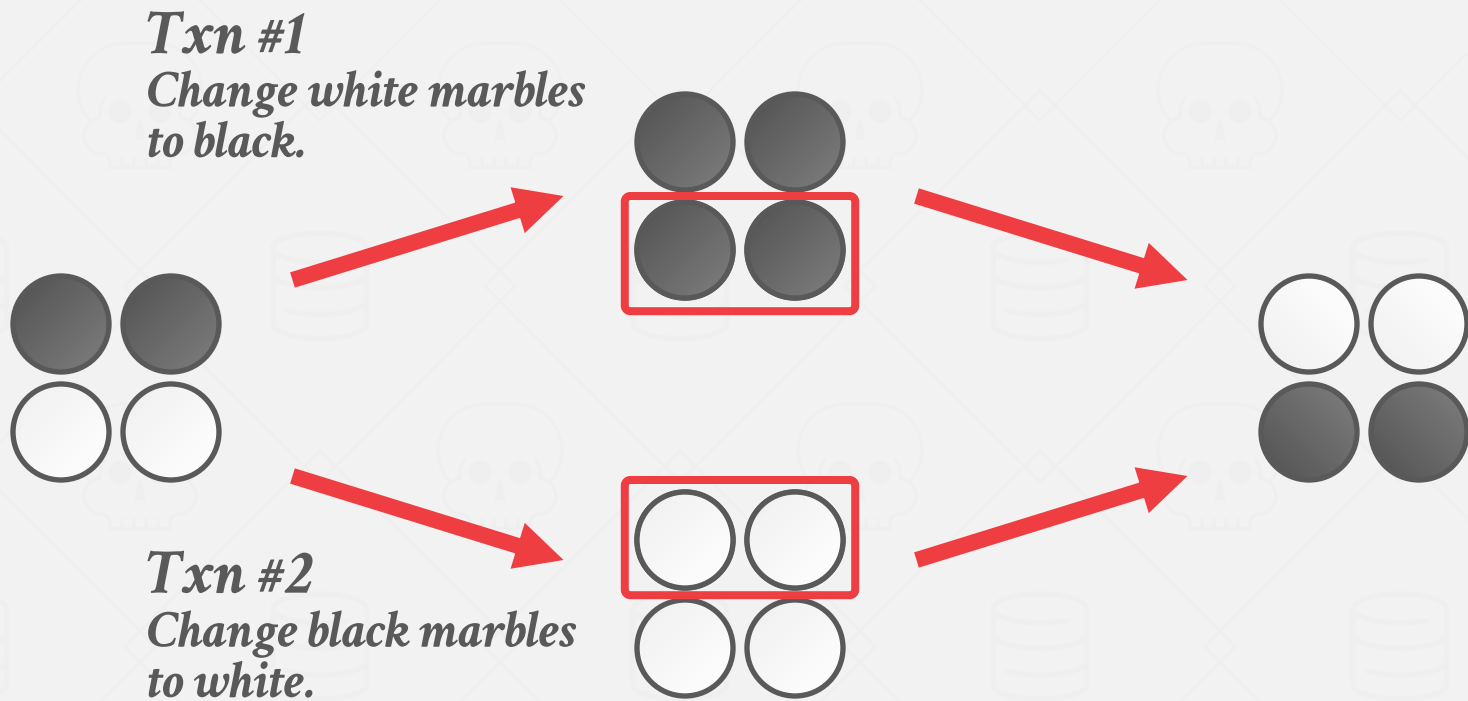
WRITE SKEW ANOMALY



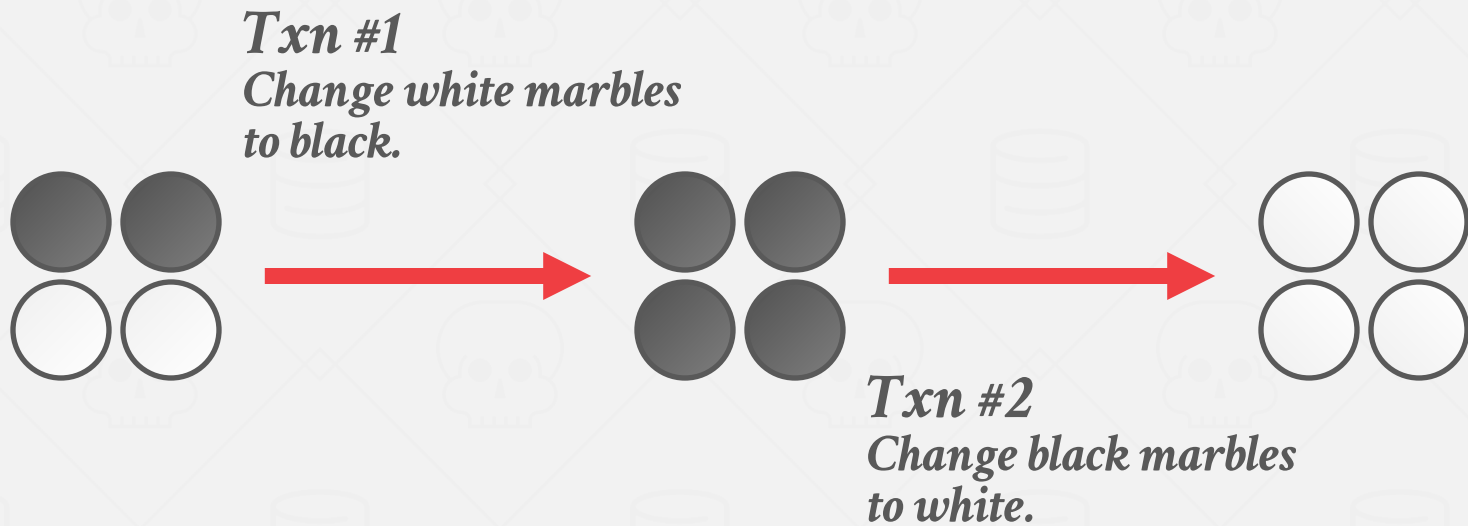
WRITE SKEW ANOMALY



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.



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.


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

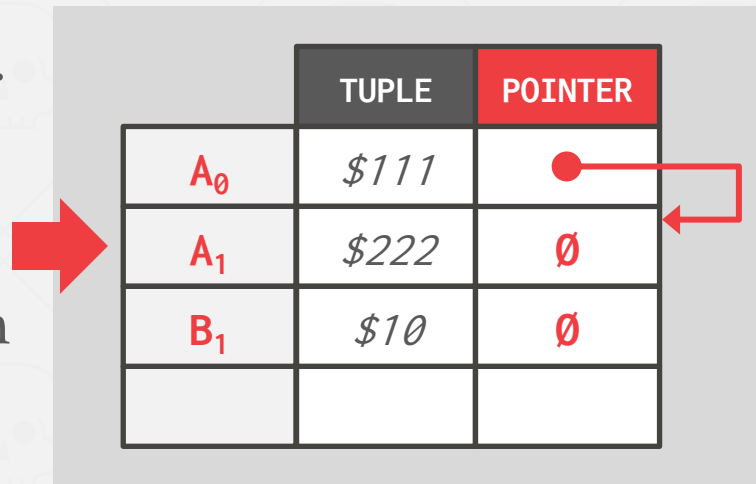
	TUPLE	POINTER
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset

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



The diagram shows a table with three columns: TUPLE, POINTER, and an unlabeled column. The rows contain the following data:

	TUPLE	POINTER	
A_0	\$111	●	→
A_1	\$222	∅	←
B_1	\$10	∅	

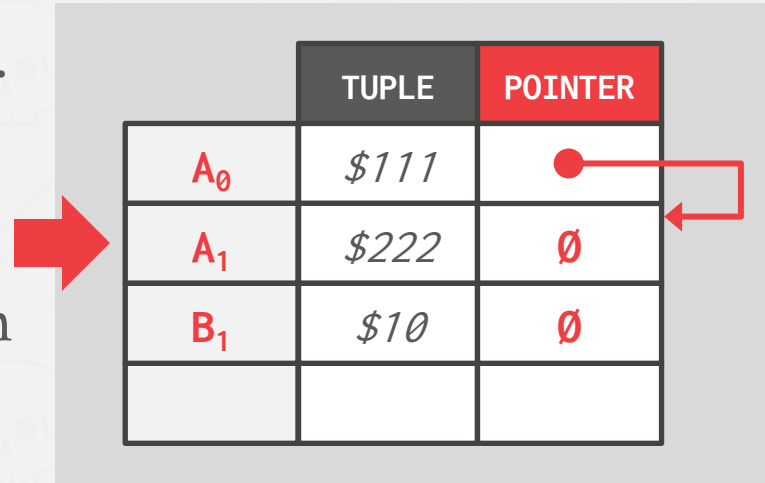
A red arrow points from the left towards the table. A red dot in the pointer cell of the first row is connected by a red line to a red arrow pointing to the right. Another red arrow points from the right towards the pointer cell of the second row.

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



The diagram shows a table with three columns: TUPLE, POINTER, and an unlabeled column. The rows contain the following data:

	TUPLE	POINTER	
A_0	\$111	●	→
A_1	\$222	∅	←
B_1	\$10	∅	

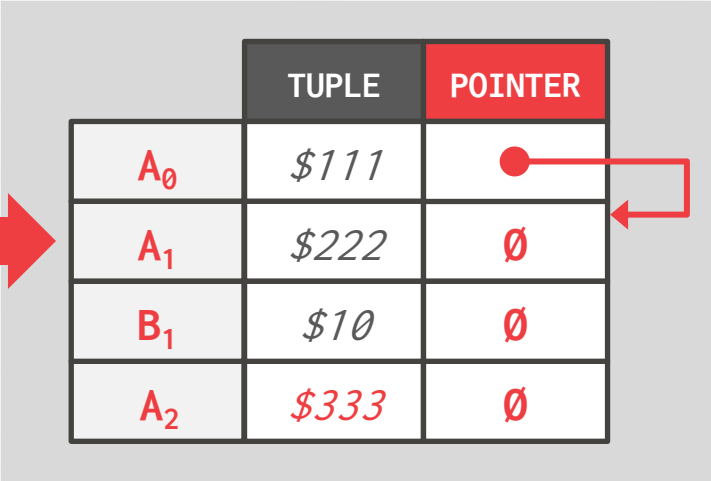
A red arrow points from the left towards the table. A red dot in the pointer cell of the first row has a line extending to the right, which then turns down and left to point at the first row of the second column.

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



The diagram shows a table with four rows and three columns. The columns are labeled 'TUPLE' and 'POINTER'. The rows contain tuple identifiers (A₀, A₁, B₁, A₂), their values (\$111, \$222, \$10, \$333), and their pointers (a red dot, and three empty sets). A red arrow points from the pointer cell of A₀ to the pointer cell of A₁. A large red arrow points from the text 'On every update...' to the table.

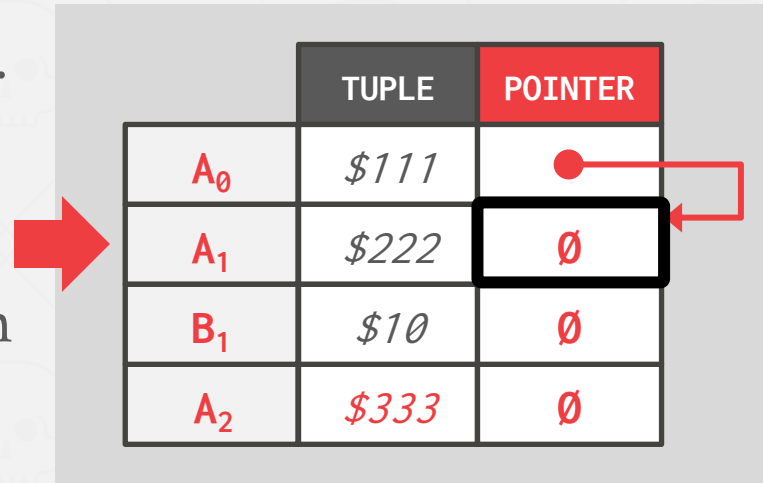
	TUPLE	POINTER
A ₀	\$111	●
A ₁	\$222	∅
B ₁	\$10	∅
A ₂	\$333	∅


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



	TUPLE	POINTER
A_0	\$111	
A_1	\$222	\emptyset
B_1	\$10	\emptyset
A_2	\$333	\emptyset

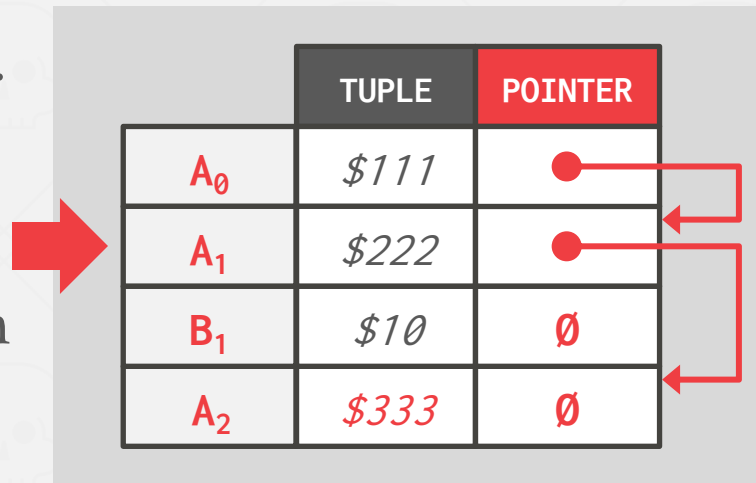
APPEND-ONLY STORAGE

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



A diagram illustrating the Main Table structure. A large red arrow points from the text on the left to the table. The table has two columns: 'TUPLE' and 'POINTER'. The rows represent different versions of tuples. Red dots in the 'POINTER' column are connected by red lines to the 'TUPLE' column of the previous row, showing a chain of updates. The first row is A₀ with value \$111 and a pointer to the second row. The second row is A₁ with value \$222 and a pointer to the third row. The third row is B₁ with value \$10 and an empty pointer. The fourth row is A₂ with value \$333 and an empty pointer.

	TUPLE	POINTER
A ₀	\$111	●
A ₁	\$222	●
B ₁	\$10	∅
A ₂	\$333	∅

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.

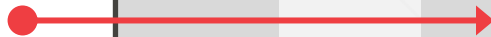
TIME-TRAVEL STORAGE

Main Table

	TUPLE	POINTER
A_2	\$222	●
B_1	\$10	


Time-Travel Table

	TUPLE	POINTER
A_1	\$111	∅



TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A_2	\$222	●
B_1	\$10	

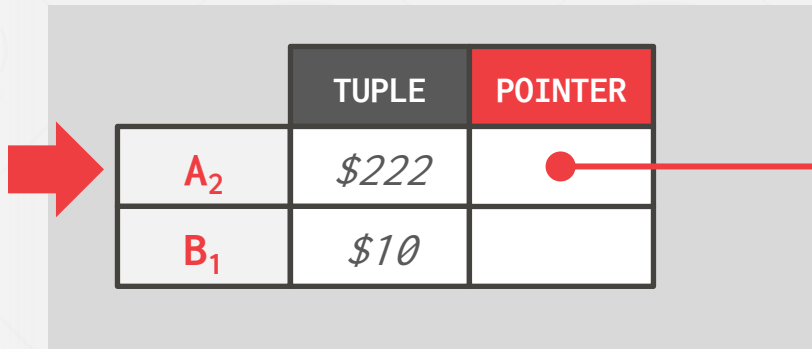
Time-Travel Table

	TUPLE	POINTER
A_1	\$111	∅

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

TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A₂	\$222	● →
B₁	\$10	

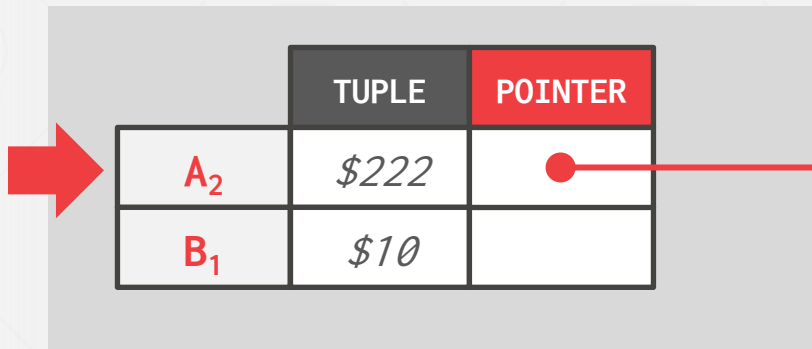
Time-Travel Table

	TUPLE	POINTER
A₁	\$111	∅

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

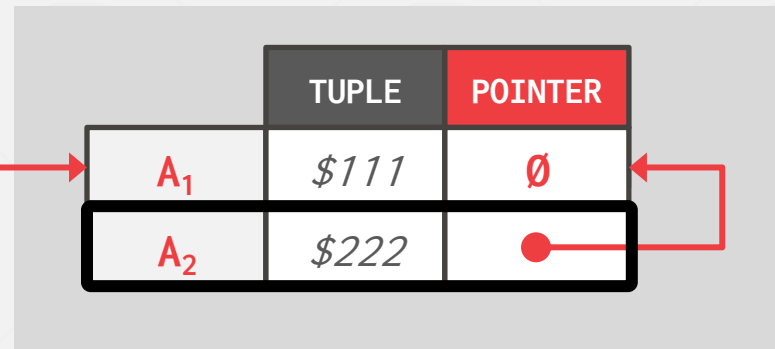
TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A_2	\$222	●
B_1	\$10	

Time-Travel Table

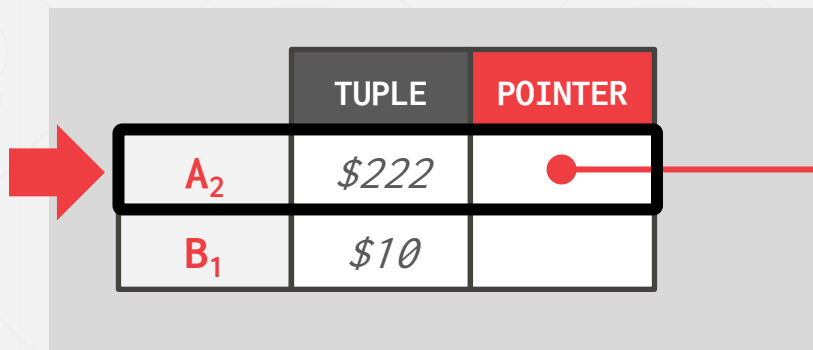


	TUPLE	POINTER
A_1	\$111	\emptyset
A_2	\$222	●

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

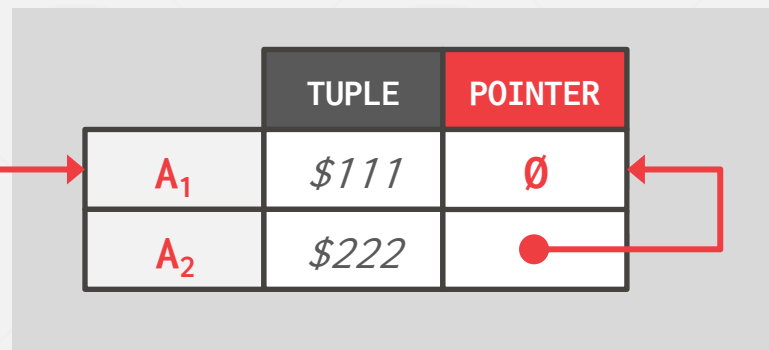
TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A₂	\$222	● →
B₁	\$10	

Time-Travel Table



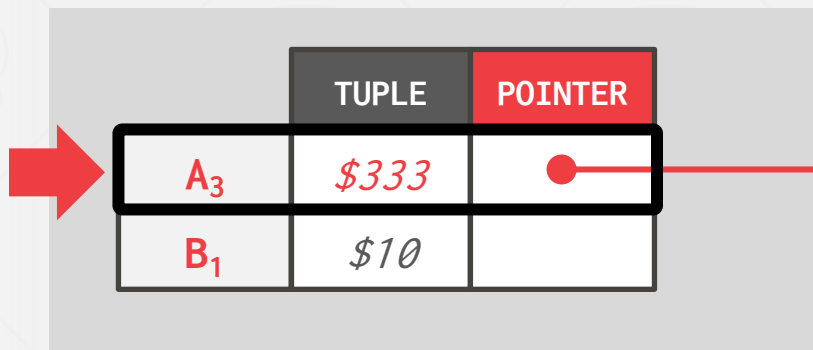
	TUPLE	POINTER
A₁	\$111	∅
A₂	\$222	● →

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

Overwrite master version in the main table and update pointers.

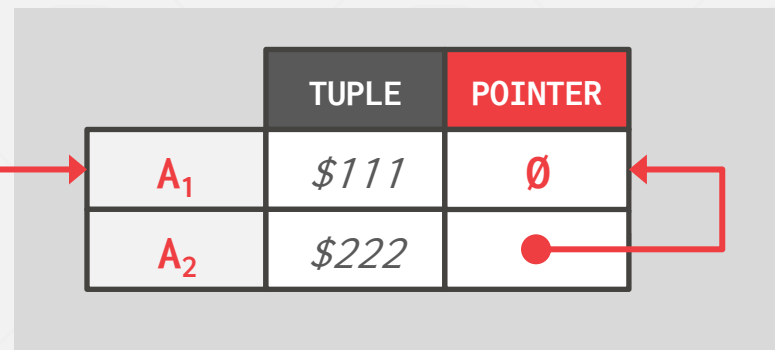
TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A₃	<i>\$333</i>	● →
B₁	<i>\$10</i>	

Time-Travel Table




	TUPLE	POINTER
A₁	<i>\$111</i>	∅ ←
A₂	<i>\$222</i>	● ←

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

Overwrite master version in the main table and update pointers.


TIME-TRAVEL STORAGE

Main Table



	TUPLE	POINTER
A_3	$\$333$	●
B_1	$\$10$	

Time-Travel Table



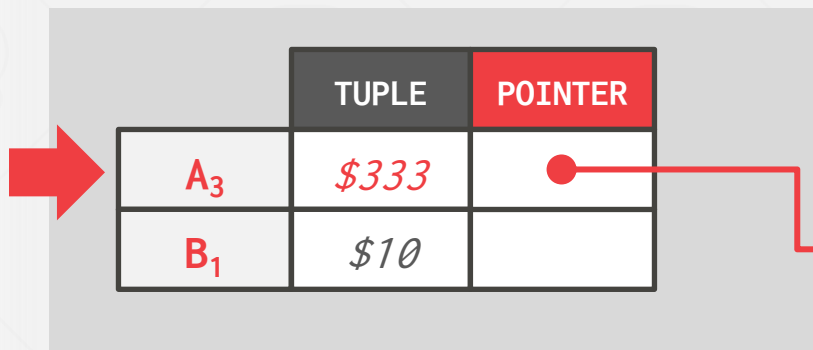
	TUPLE	POINTER
A_1	$\$111$	\emptyset
A_2	$\$222$	●

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

Overwrite master version in the main table and update pointers.

TIME-TRAVEL STORAGE

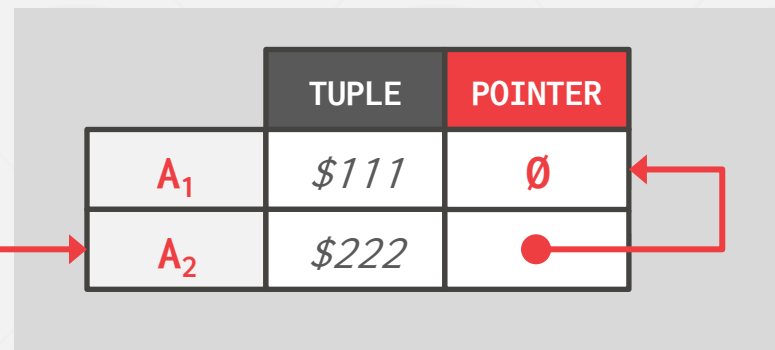
Main Table



The Main Table is a 2x2 grid with columns 'TUPLE' and 'POINTER'. The first row contains 'A₃' and '\$333'. The second row contains 'B₁' and '\$10'. A red arrow points to the first row. A red dot in the 'POINTER' cell of the first row has a line that goes right, then down, then left to point at the first row of the Time-Travel Table.

	TUPLE	POINTER
→	A ₃	\$333
	B ₁	\$10

Time-Travel Table



The Time-Travel Table is a 2x2 grid with columns 'TUPLE' and 'POINTER'. The first row contains 'A₁' and '\$111'. The second row contains 'A₂' and '\$222'. A red arrow points to the second row. A red dot in the 'POINTER' cell of the second row has a line that goes right, then up, then left to point at the second row of the Main Table.

	TUPLE	POINTER
	A ₁	\$111
←	A ₂	\$222

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

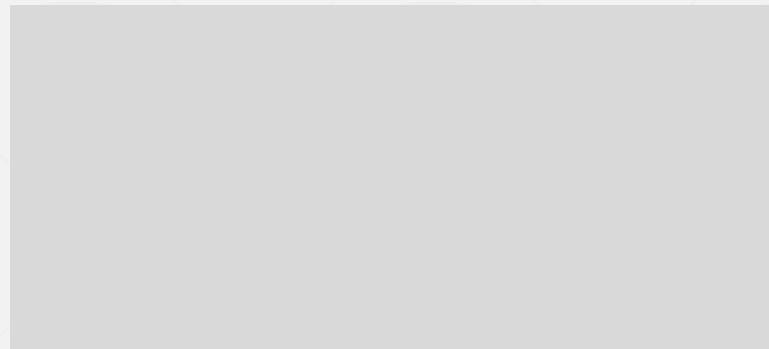
Overwrite master version in the main table and update pointers.

DELTA STORAGE

Main Table


	VALUE	POINTER
A_1	\$111	
B_1	\$10	

Delta Storage Segment



DELTA STORAGE

Main Table




	VALUE	POINTER
A₁	\$111	
B₁	\$10	

Delta Storage Segment

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

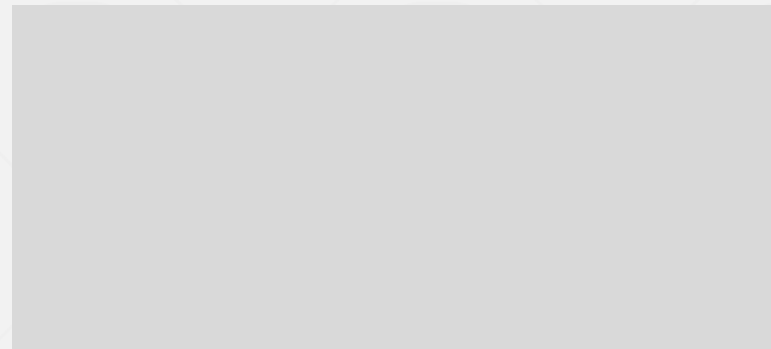
DELTA STORAGE

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	

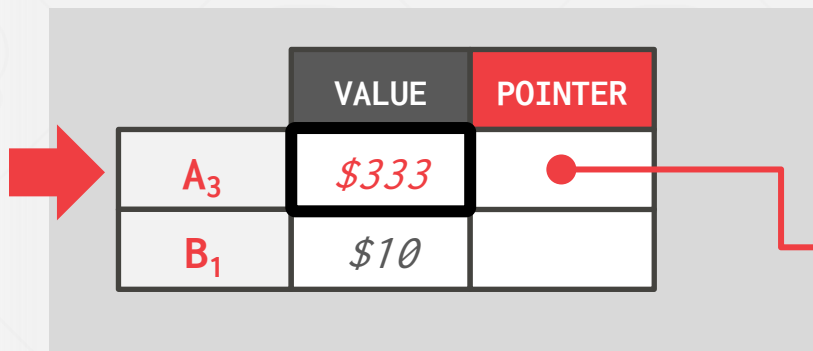
Delta Storage Segment

	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.

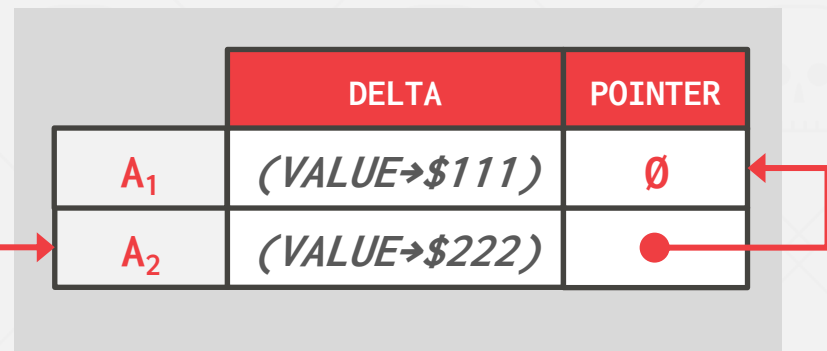
DELTA STORAGE

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	

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

Delta Storage Segment

	DELTA	POINTER
A ₁	(VALUE->\$111)	∅
A ₂	(VALUE->\$222)	●

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

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>A_{100}</i>	<i>1</i>	<i>9</i>
<i>B_{100}</i>	<i>1</i>	<i>9</i>
<i>B_{101}</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

Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

Vacuum



	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

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

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.

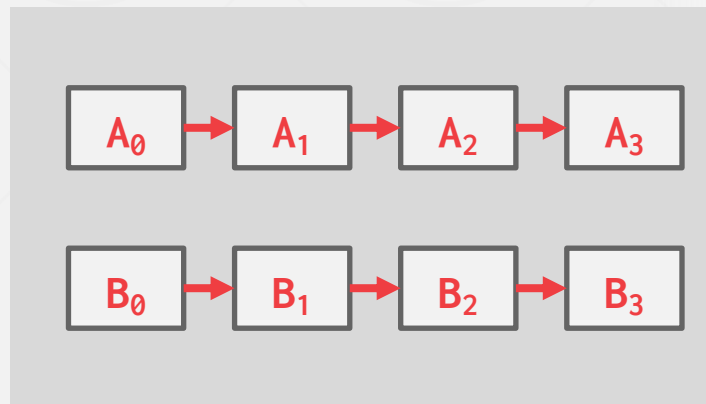
TUPLE-LEVEL GC

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

$T_{id}=12$

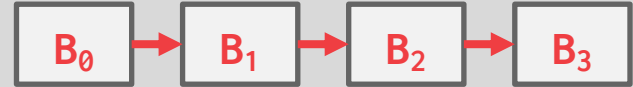
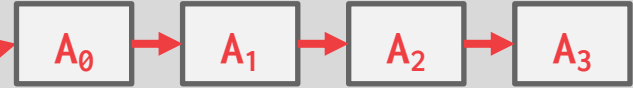
GET(A)



INDEX

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

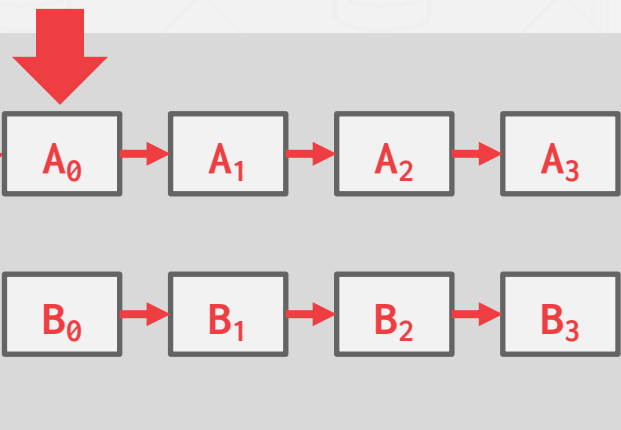
Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

GET(A)



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.

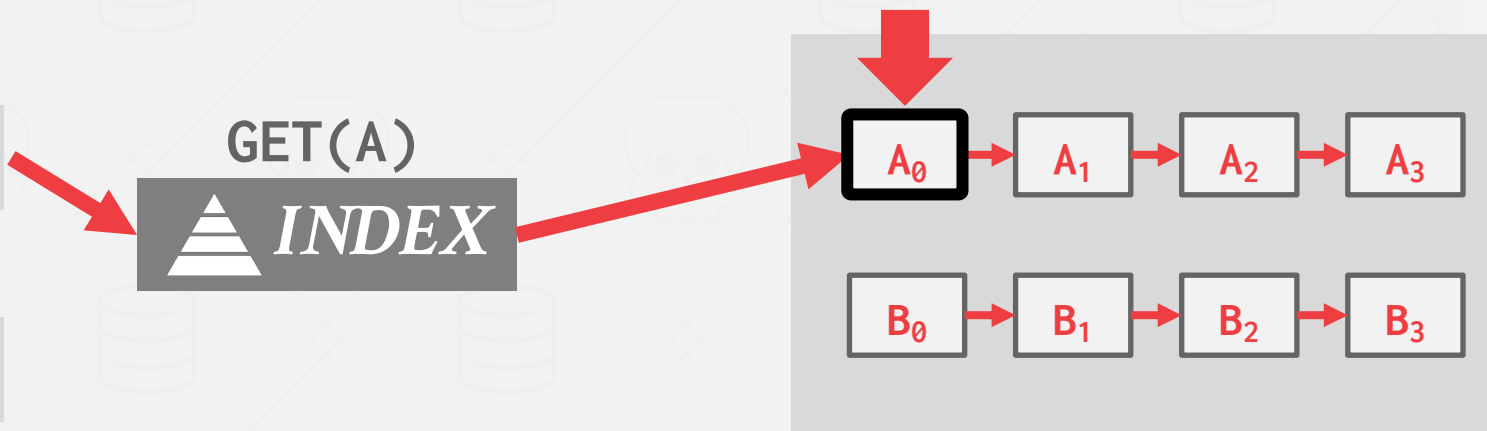
TUPLE-LEVEL GC

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.

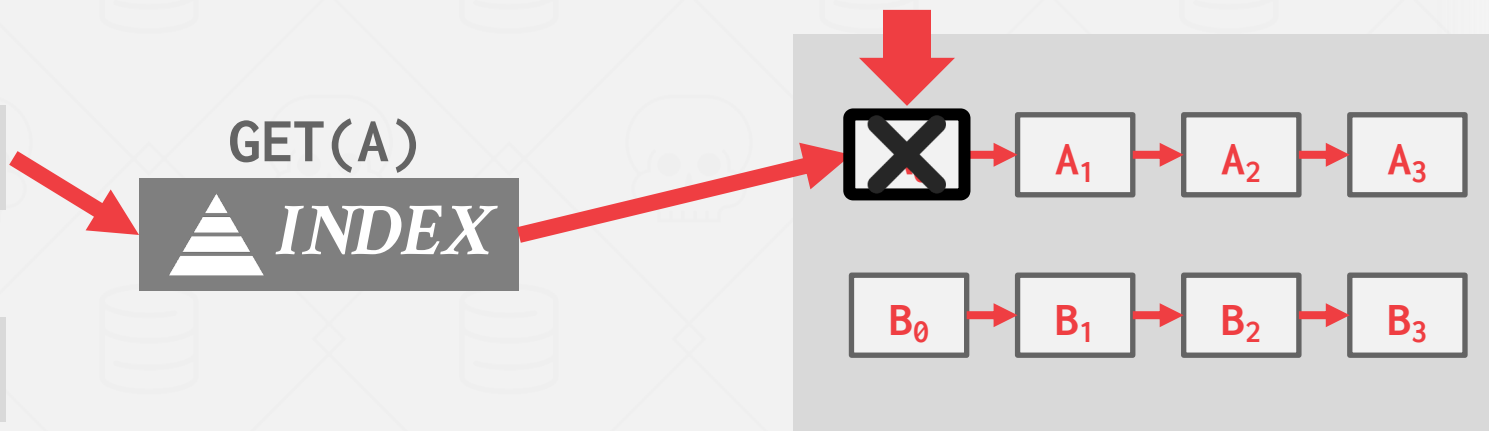
TUPLE-LEVEL GC

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.

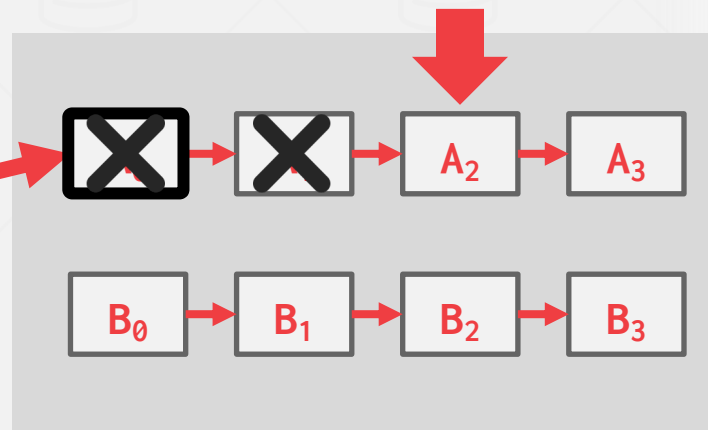
TUPLE-LEVEL GC

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

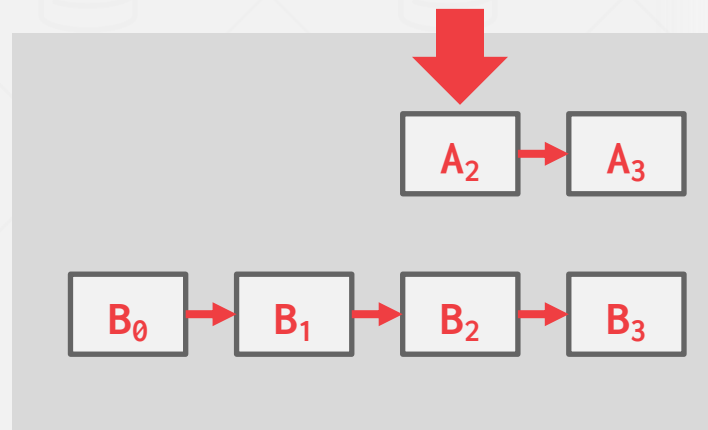
$T_{id}=12$

GET(A)



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

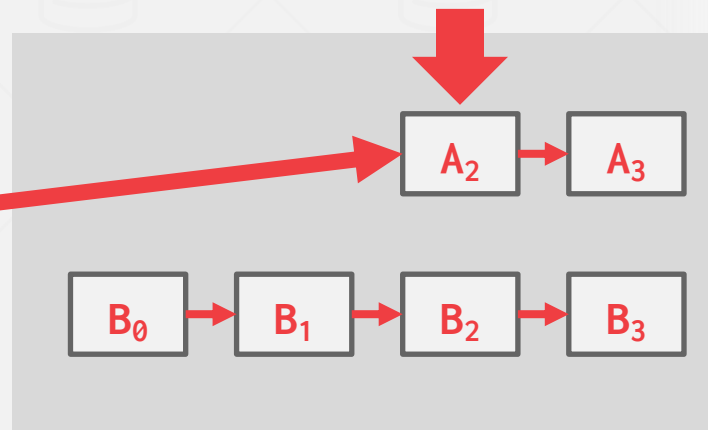
Txn #1

$T_{id}=12$

Txn #2

$T_{id}=25$

GET(A)



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	DATA
A_2	1	∞	-
B_6	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)

	BEGIN-TS	END-TS	DATA
A_2	1	∞	-
B_6	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



	BEGIN-TS	END-TS	DATA
A_2	1	∞	-
B_6	8	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)



	BEGIN-TS	END-TS	DATA
A_2	1	10	-
B_6	8	∞	-
A_3	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



Old Versions

A₂

	BEGIN-TS	END-TS	DATA
A₂	1	10	-
B₆	8	∞	-
A₃	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)

Old Versions

A₂

	BEGIN-TS	END-TS	DATA
A₂	1	10	-
B₆	8	∞	-
A₃	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)



UPDATE(B)

Old Versions

A₂



	BEGIN-TS	END-TS	DATA
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	DATA
A ₂	1	10	-
B ₆	8	10	-
A ₃	10	∞	-
B ₇	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10



UPDATE(A)



UPDATE(B)

Old Versions

A_2

B_6

	BEGIN-TS	END-TS	DATA
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)

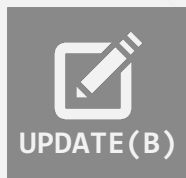
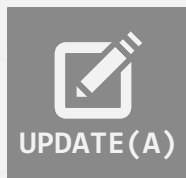
	BEGIN-TS	END-TS	DATA
A₂	1	10	-
B₆	8	10	-
A₃	10	∞	-
B₇	10	∞	-

TRANSACTION-LEVEL GC

Txn #1

BEGIN @ 10
COMMIT @ 15

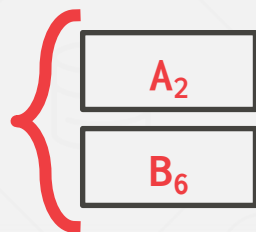
Old Versions



	BEGIN-TS	END-TS	DATA
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**.

The image shows a presentation slide from Uber Engineering. The slide title is "WHY UBER ENGINEERING SWITCHED FROM POSTGRES TO MYSQL" by Evan Klitzke, dated July 20, 2016. It features a diagram illustrating index management. The diagram shows a "Secondary Index" with columns A, B, C, and D. Below it is a "Primary Index" with columns 1, 2, 3, and 4. Arrows indicate that secondary index A points to primary index 1, B to 2, C to 3, and D to 4. Below the primary index is a "Disk" with four green bars representing data blocks at addresses 76, 103, 107, and 211. Arrows show that primary index 1 points to disk address 76, 2 to 103, 3 to 107, and 4 to 211.

Secondary indexes are more complicated...

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



PRIMARY INDEX



SECONDARY INDEX



INDEX POINTERS



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS

GET(A) ↓



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS

GET(A) ↓



PRIMARY INDEX



SECONDARY INDEX

Record Id

A₄

A₃

A₂

A₁

*Append-Only
Newest-to-Oldest*

INDEX POINTERS



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS

↓ GET(A)



PRIMARY INDEX

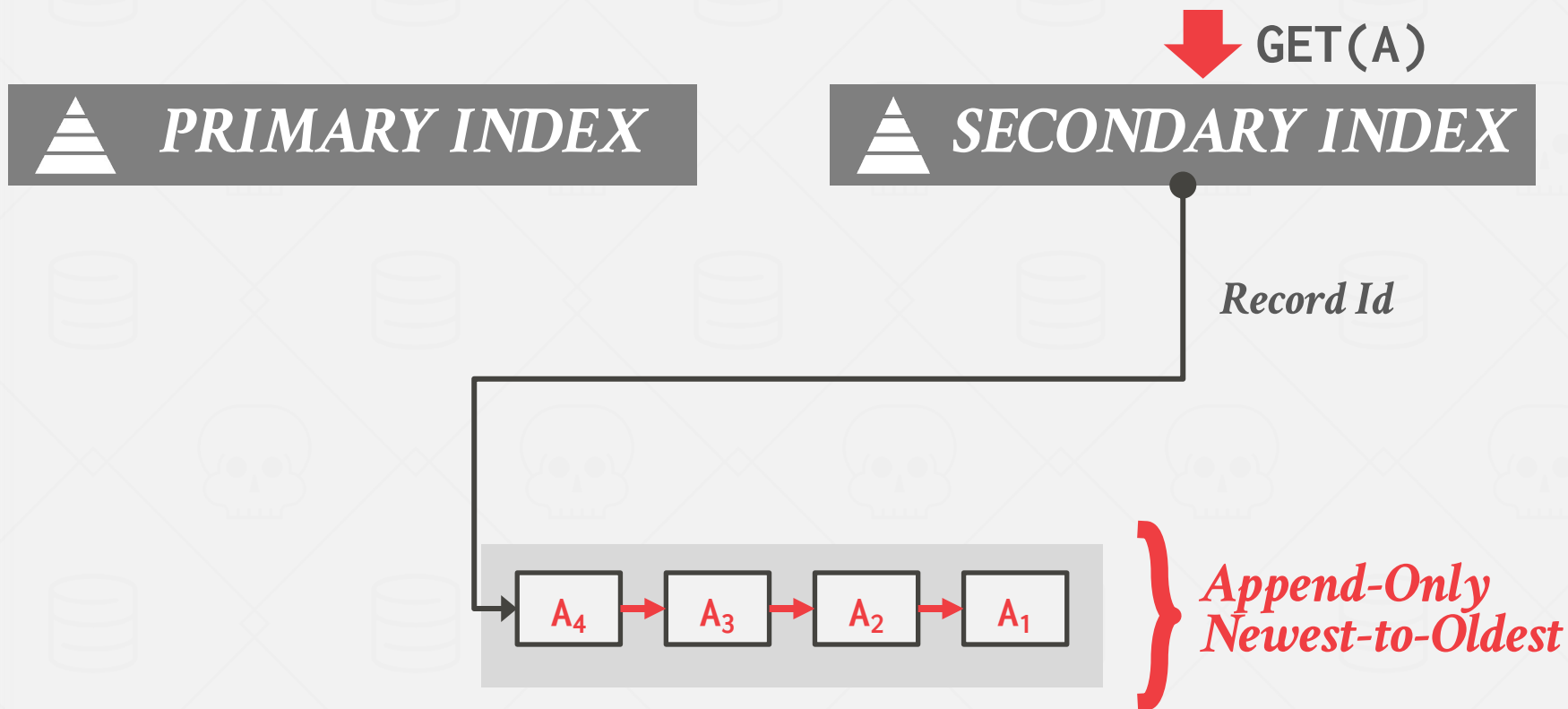


SECONDARY INDEX

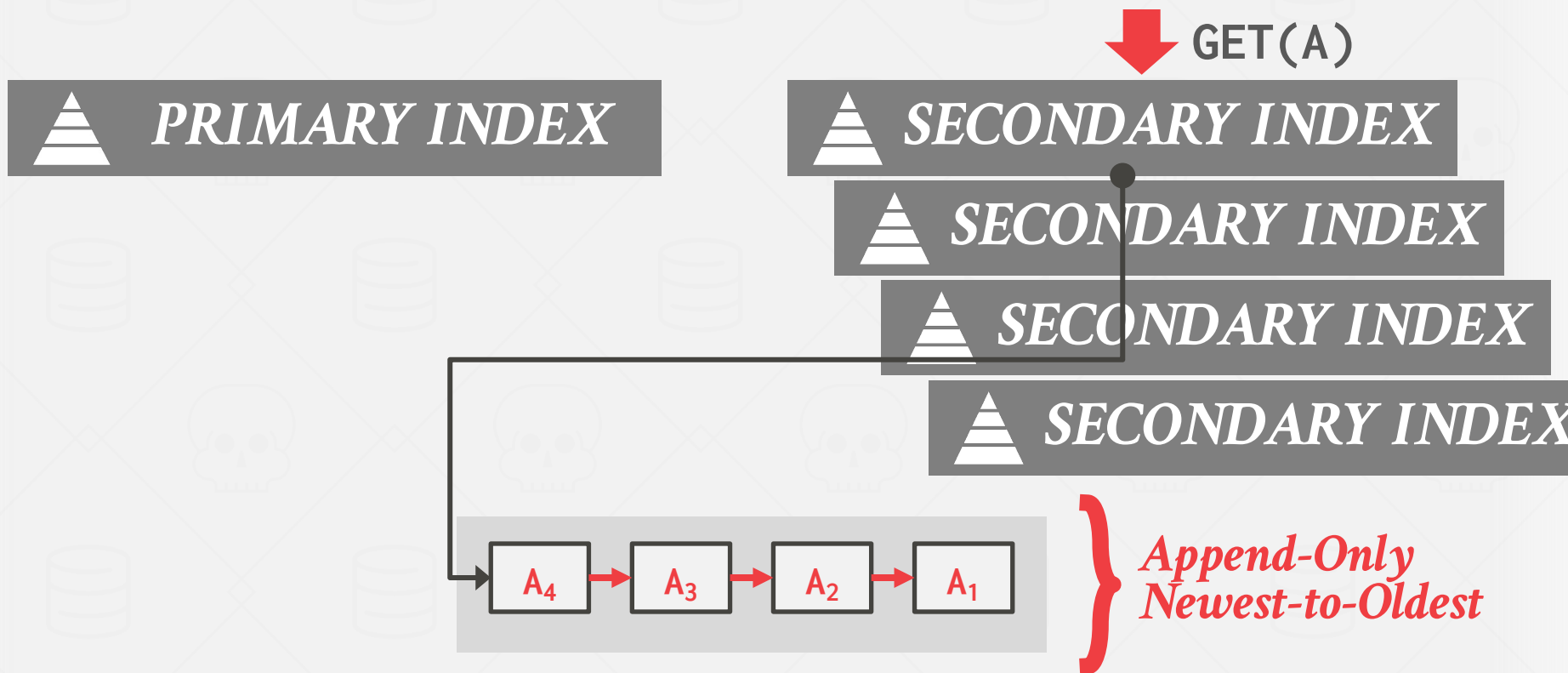


*Append-Only
Newest-to-Oldest*

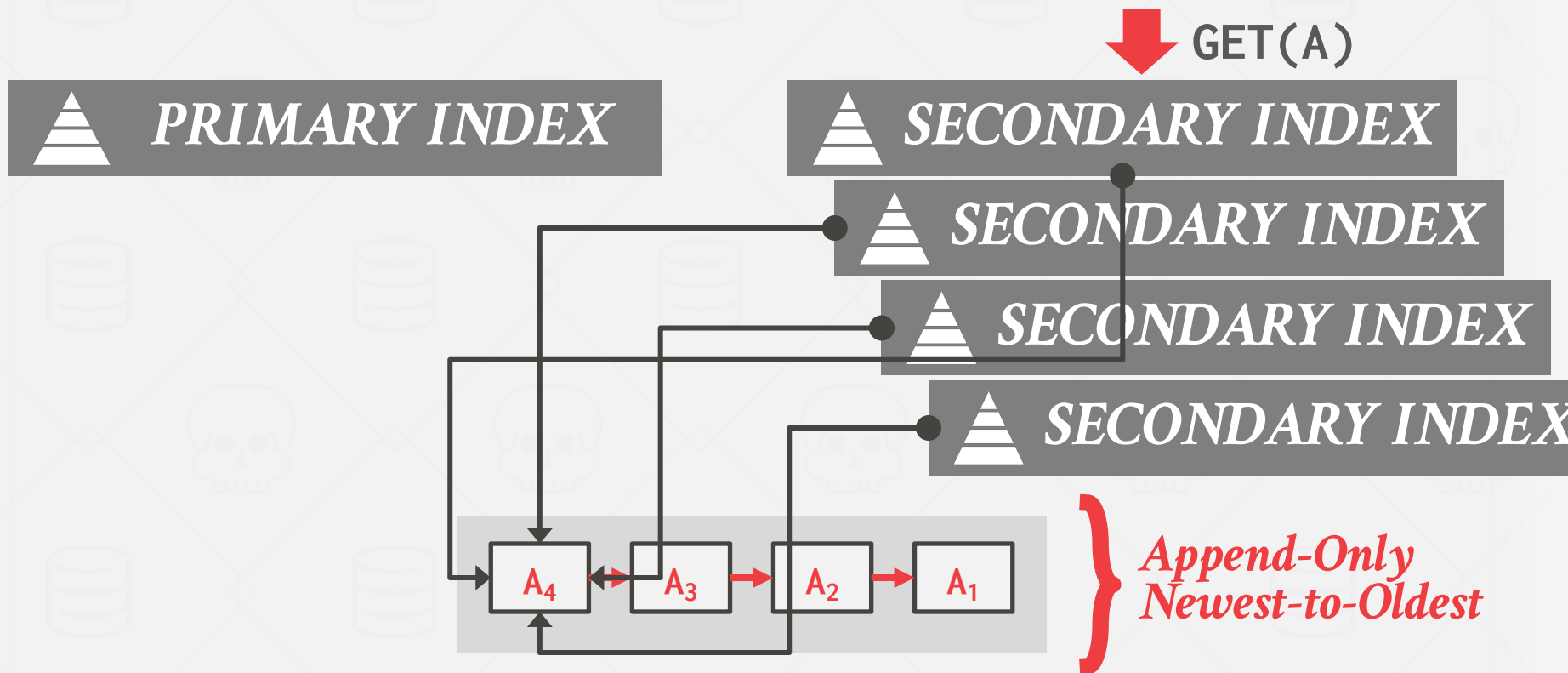
INDEX POINTERS



INDEX POINTERS



INDEX POINTERS



INDEX POINTERS

↓ GET(A)



PRIMARY INDEX

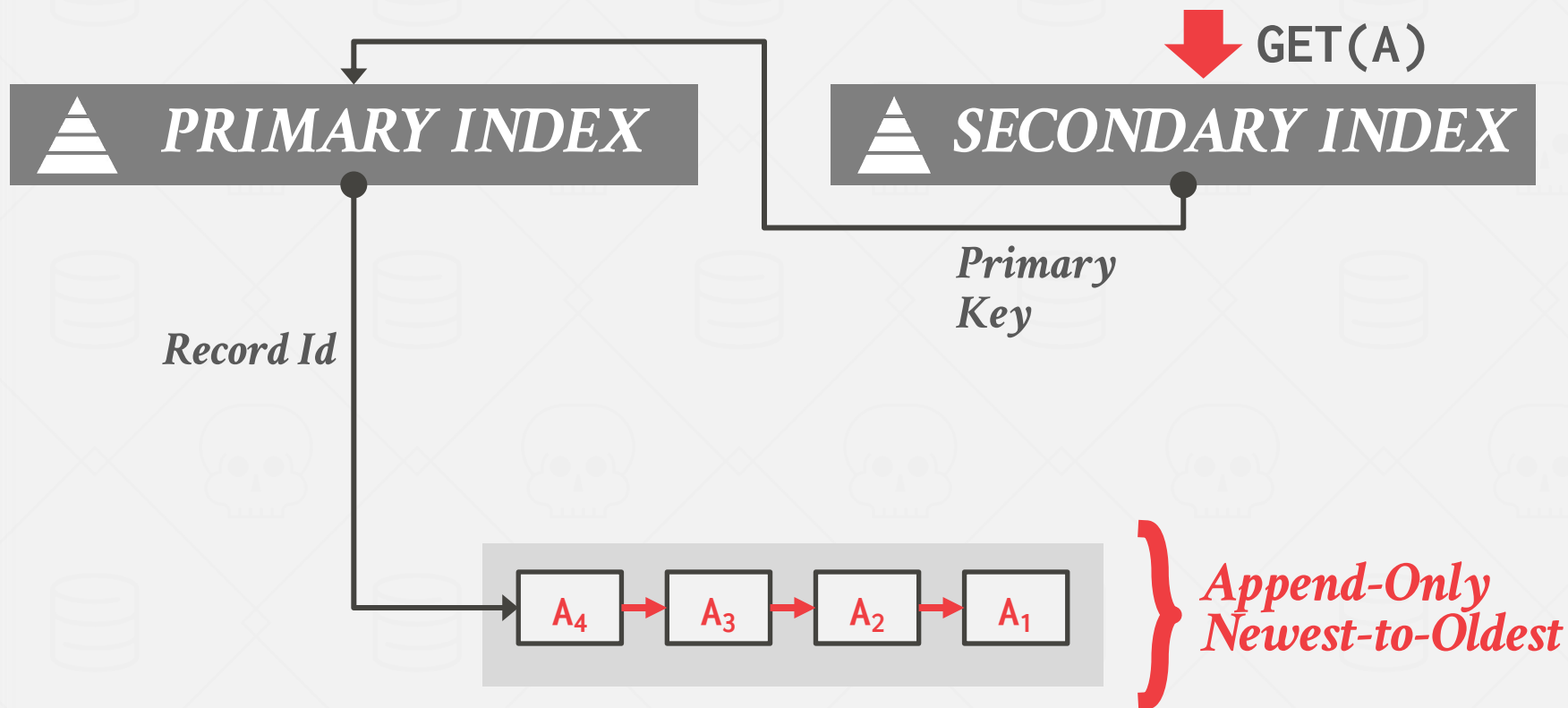


SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

INDEX POINTERS



INDEX POINTERS

↓ GET(A)



PRIMARY INDEX



SECONDARY INDEX



*Append-Only
Newest-to-Oldest*

MVCC INDEXES

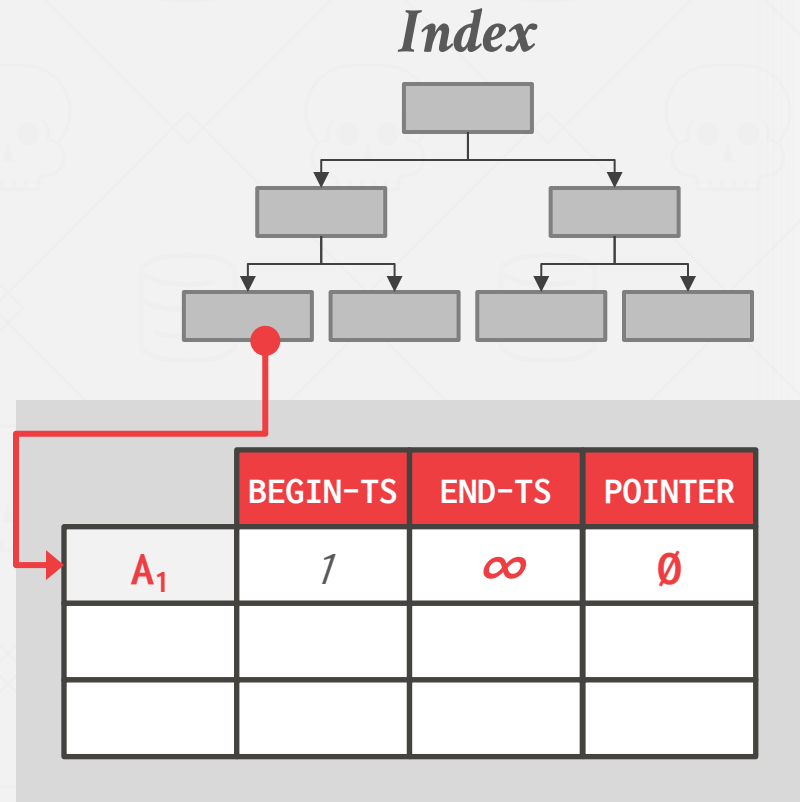
MVCC DBMS indexes (usually) do not store version information about tuples with their keys.

→ 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

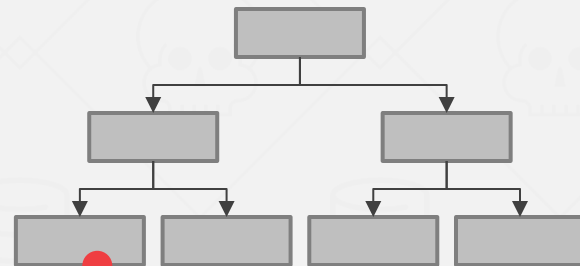
Txn #1

BEGIN @ 10



READ(A)

Index



	BEGIN-TS	END-TS	POINTER
<i>A₁</i>	<i>1</i>	<i>∞</i>	<i>∅</i>

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10

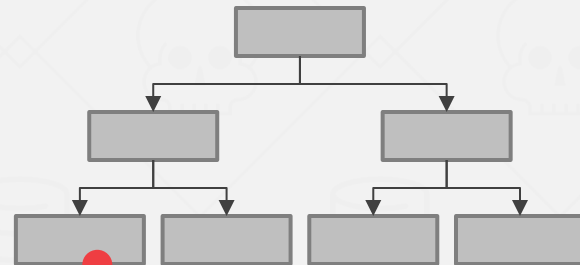


Txn #2

BEGIN @ 20



Index



	BEGIN-TS	END-TS	POINTER
A_1	1	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)

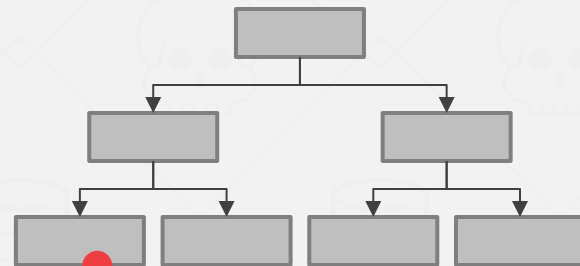
Txn #2

BEGIN @ 20



UPDATE(A)

Index



	BEGIN-TS	END-TS	POINTER
A ₁	1	20	
A ₂	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10

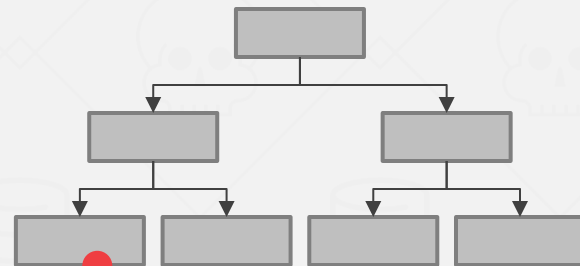


Txn #2

BEGIN @ 20



Index



	BEGIN-TS	END-TS	POINTER
A ₁	1	20	
X	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



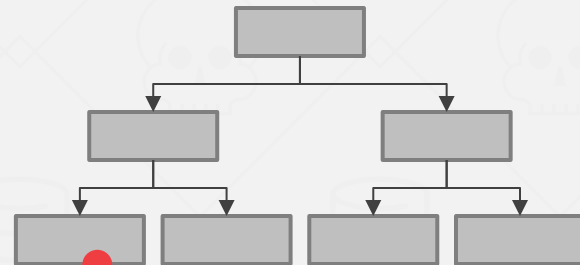
Txn #2

BEGIN @ 20

COMMIT @ 25



Index



	BEGIN-TS	END-TS	POINTER
A₁	1	20	
X	20	∞	\emptyset

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



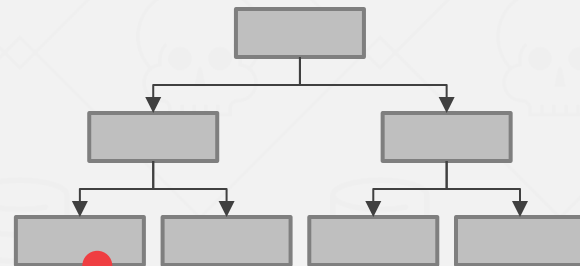
Txn #2

BEGIN @ 20

COMMIT @ 25



Index



	BEGIN-TS	END-TS	POINTER
<i>A₁</i>	1	20	
	20	20	<i>∅</i>

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



Txn #2

BEGIN @ 20

COMMIT @ 25

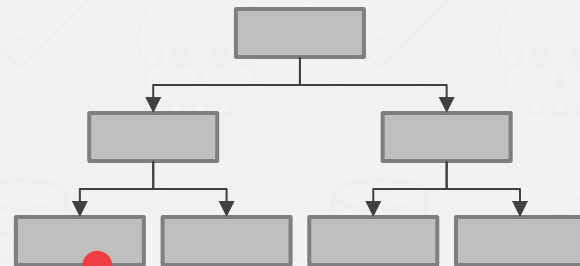


Txn #3

BEGIN @ 30



Index



	BEGIN-TS	END-TS	POINTER
A ₁	1	20	
X	20	20	∅

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



Txn #2

BEGIN @ 20

COMMIT @ 25

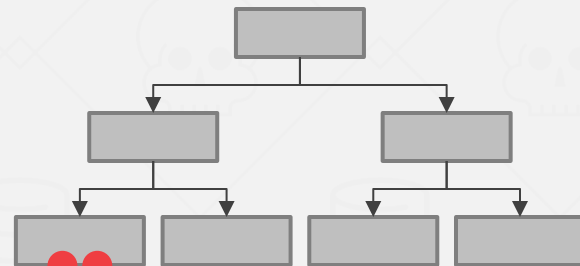


Txn #3

BEGIN @ 30



Index



	BEGIN-TS	END-TS	POINTER
A ₁	1	20	●
A₁	20	20	∅
A ₁	30	∞	∅

MVCC DUPLICATE KEY PROBLEM

Txn #1

BEGIN @ 10



READ(A)



READ(A)

Txn #2

BEGIN @ 20

COMMIT @ 25



UPDATE(A)



DELETE(A)

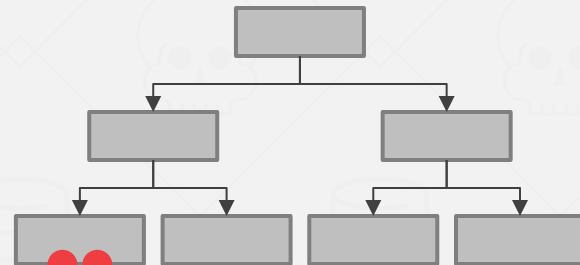
Txn #3

BEGIN @ 30



INSERT(A)

Index



	BEGIN-TS	END-TS	POINTER
A ₁	1	20	
A₁	20	20	∅
A ₁	30	∞	∅

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.

GLOBAL CLOCKS

Spanner: TrueTime and external consistency

TrueTime is a highly available, distributed clock that is provided to applications on all Google servers¹ (#1). TrueTime enables applications to generate monotonically increasing timestamps: an application can compute a timestamp T that is guaranteed to be greater than any timestamp T' if T finished being generated before T' started being generated. This guarantee holds across all servers and all timestamps.

This feature of TrueTime is used by Spanner to assign timestamps to transactions. Specifically, every transaction is assigned a timestamp that reflects the instant at which Spanner considers it to have occurred. Because Spanner uses multi-version concurrency control, the ordering guarantee on timestamps enables clients of Spanner to perform consistent reads across an entire database (even across multiple Cloud [regions](#) (/about/locations)) without blocking writes.

External consistency

Spanner provides clients with the strictest concurrency-control guarantees for transactions, which is called external consistency² (#2). Under external consistency, the system behaves as if all transactions were executed sequentially, even though Spanner actually runs them across multiple servers (and possibly in multiple datacenters) for higher performance and availability. In addition if one transaction completes before another transaction starts to commit, the system guarantees that clients can never see a state that includes the effect of the second transaction but not the first. Intuitively, Spanner is semantically indistinguishable from a single-machine database. Even though it provides such strong guarantees, Spanner enables applications to achieve performance comparable to databases that provide weaker guarantees (in return for higher performance). For example, like databases that support snapshot isolation, Spanner allows writes to proceed without being blocked by read-only transactions, but without exhibiting the anomalies that snapshot isolation allows.

External consistency greatly simplifies application development. For example, suppose that you have created a banking application on Spanner and one of your customers starts with \$50 in their checking account and \$50 in their savings account. Your application then begins a workflow in which it first commits a transaction T_1 to deposit \$200 into the savings account, and then issues a second transaction T_2 to debit \$150 from the checking account. Further, assume that at the end of the day, negative balances in one account are covered automatically from other accounts, and a customer

<https://cloud.google.com/spanner/docs/true-time-external-consistency>

[AWS Compute Blog](#)

It's About Time: Microsecond-Accurate Clocks on Amazon EC2 Instances

by [Chris Munns](#) | on 16 NOV 2023 | in [Amazon EC2](#), [Compute](#) | [Permalink](#) | [Share](#)

This post is written by [Josh Levinson](#), AWS Principal Product Manager and [Julien Ridoux](#), AWS Principal Software Engineer

Today, we announced that [we improved the Amazon Time Sync Service](#) to microsecond-level clock accuracy on supported Amazon EC2 instances. This new capability adds a local reference clock to your EC2 instance and is designed to deliver clock accuracy in the low double-digit microsecond range within your instance's guest OS software. This post shows you how to connect to the improved clocks on your EC2 instances. This post also demonstrates how you can measure your clock accuracy and easily generate and compare timestamps from your EC2 instances with [ClockBound](#), an open source daemon and library.

In general, it's hard to achieve high-fidelity clock synchronization due to hardware limitations and network variability. While customers have depended on the Amazon Time Sync Service to provide one millisecond clock accuracy, workloads that need microsecond-range accuracy, such as financial trading and broadcasting, require customers to maintain their own time infrastructure, which is a significant operational burden, and expensive. Other clock-sensitive applications that run on the cloud, including distributed databases and storage, have to incorporate message exchange delays with wait periods, data locks, or transaction journaling to maintain consistency at scale.

With global and reliable microsecond-range clock accuracy, you can now migrate and modernize your most time-sensitive applications in the cloud and retire your burdensome on-premises time infrastructure. You can also simplify your applications and increase their throughput by leveraging the high-accuracy timestamps to determine the ordering of events and transactions on workloads across instances, [Availability Zones](#), and [Regions](#). Additionally, you can audit the improved Amazon Time Sync Service to measure and monitor the expected microsecond-range accuracy.

New improvements to Amazon Time Sync Service

The new local clock source can be accessed over the existing Amazon Time Sync Service's Network Time Protocol (NTP) IPv4 and IPv6 endpoints, or by configuring a new Precision Time Protocol (PTP) reference clock device, to get the best accuracy possible. It's important to note that both NTP and the new PTP Hardware Clock (PHC) device share the same highly accurate source of time. The new PHC device is part of the [AWS Nitro System](#), so it is directly accessible on supported bare metal and virtualized Amazon EC2 instances without using any customer resources.

<https://aws.amazon.com/blogs/compute/its-about-time-microsecond-accurate-clocks-on-amazon-ec2-instances/>

NEXT CLASS

Logging and recovery!