

# Multi-Version Concurrency Control



Lecture #19



Database Systems  
15-445/15-645  
Fall 2018

AP

Andy Pavlo  
Computer Science  
Carnegie Mellon Univ.

# ADMINISTRIVIA

---

**Homework #4:** Monday Nov 12<sup>th</sup> @ 11:59pm

**Project #3:** Monday Nov 19<sup>th</sup> @ 11:59am



# 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**



# MULTI-VERSION CONCURRENCY CONTROL

---

**Writers don't block readers.**  
**Readers don't block writers.**

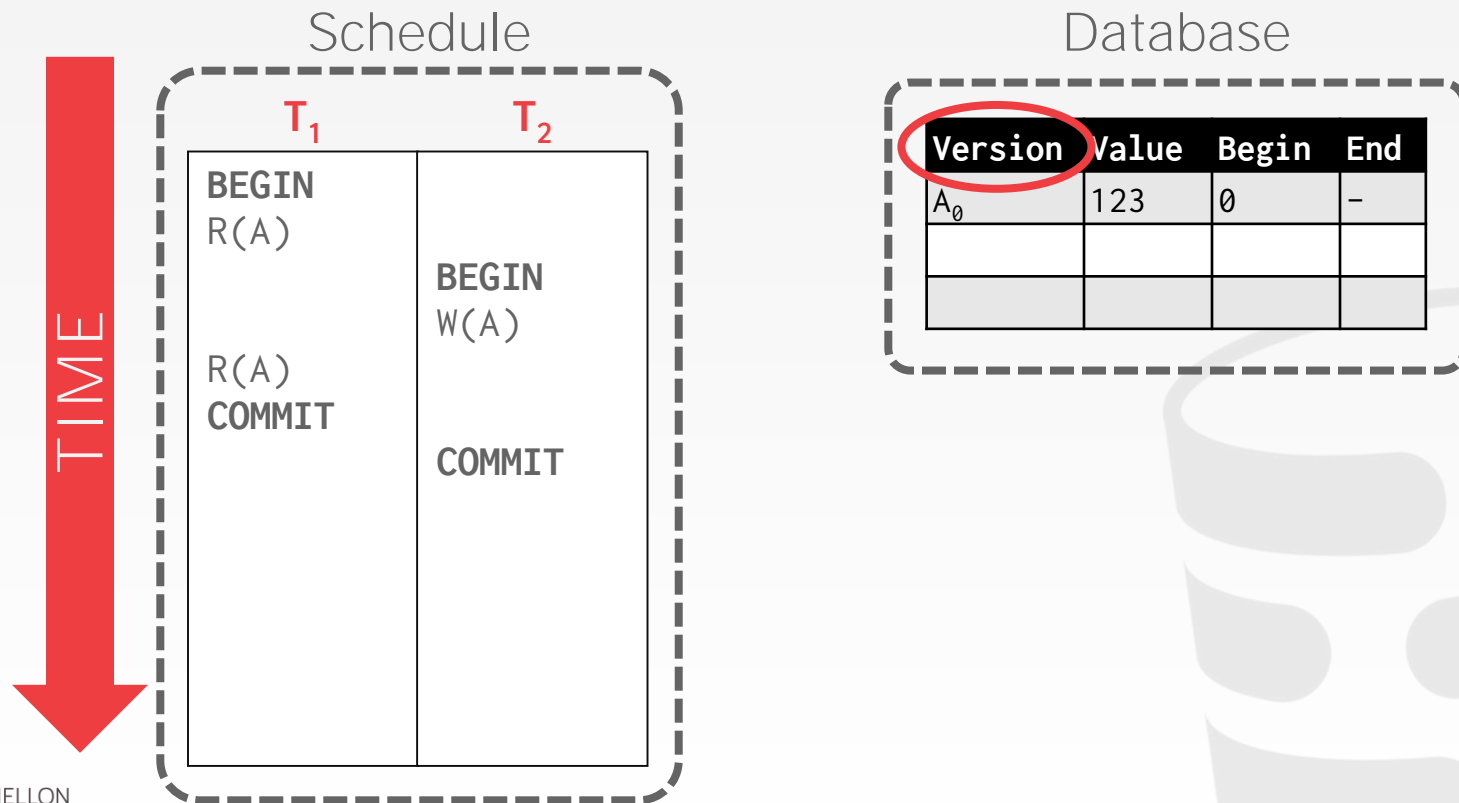
Read-only txns can read a consistent snapshot without acquiring locks.

→ Use timestamps to determine visibility.

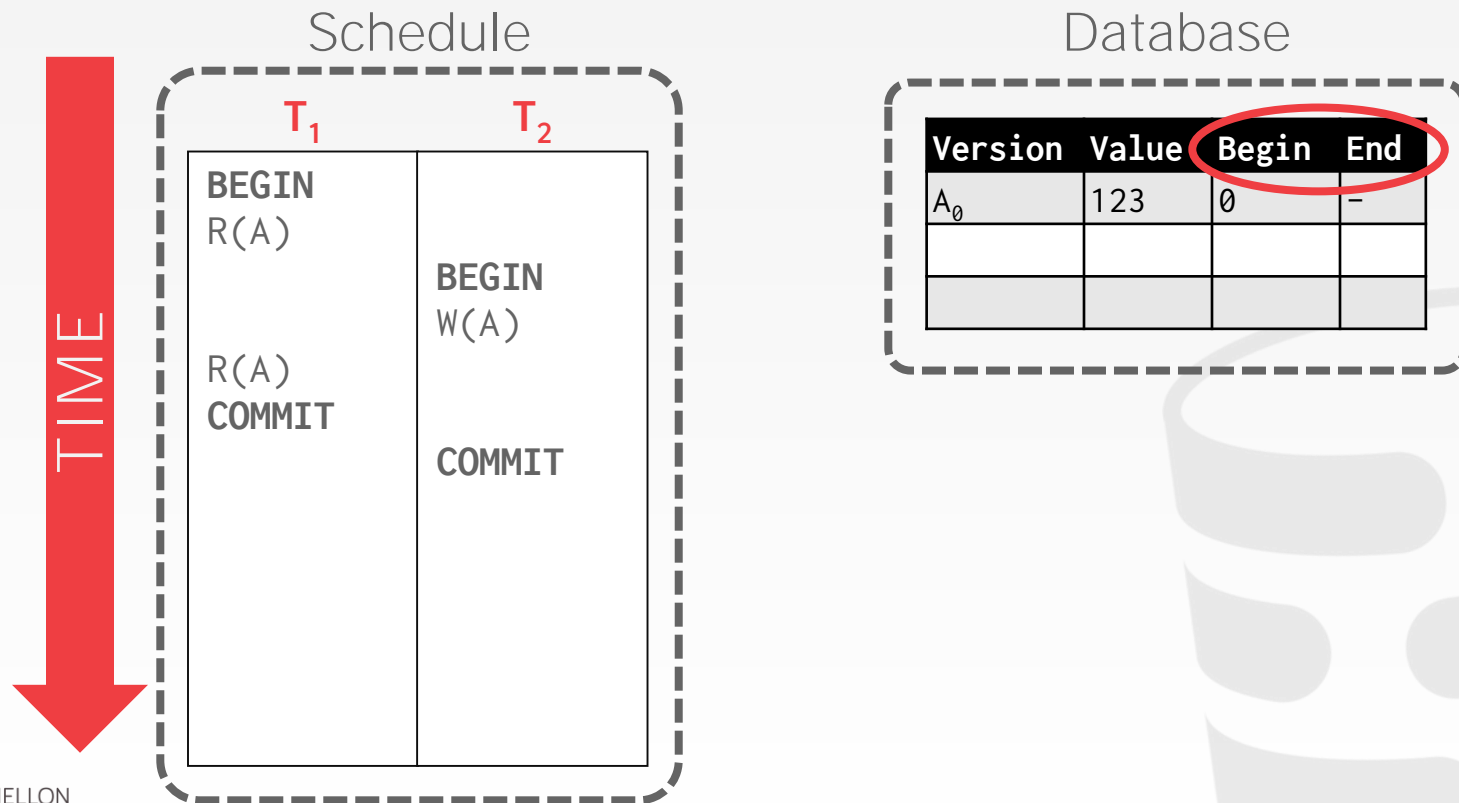
Easily support time-travel queries.



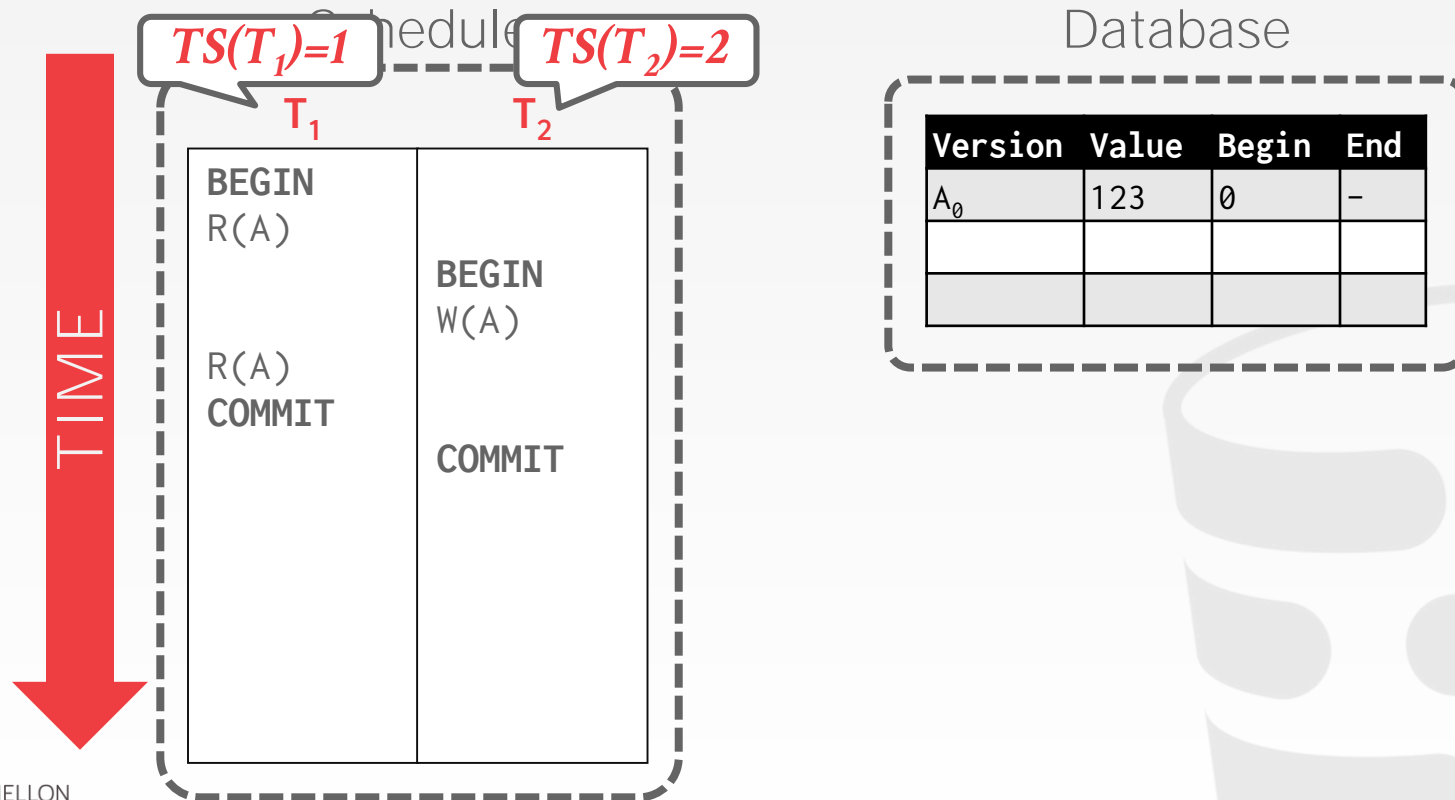
# MVCC – EXAMPLE #1



# MVCC – EXAMPLE #1

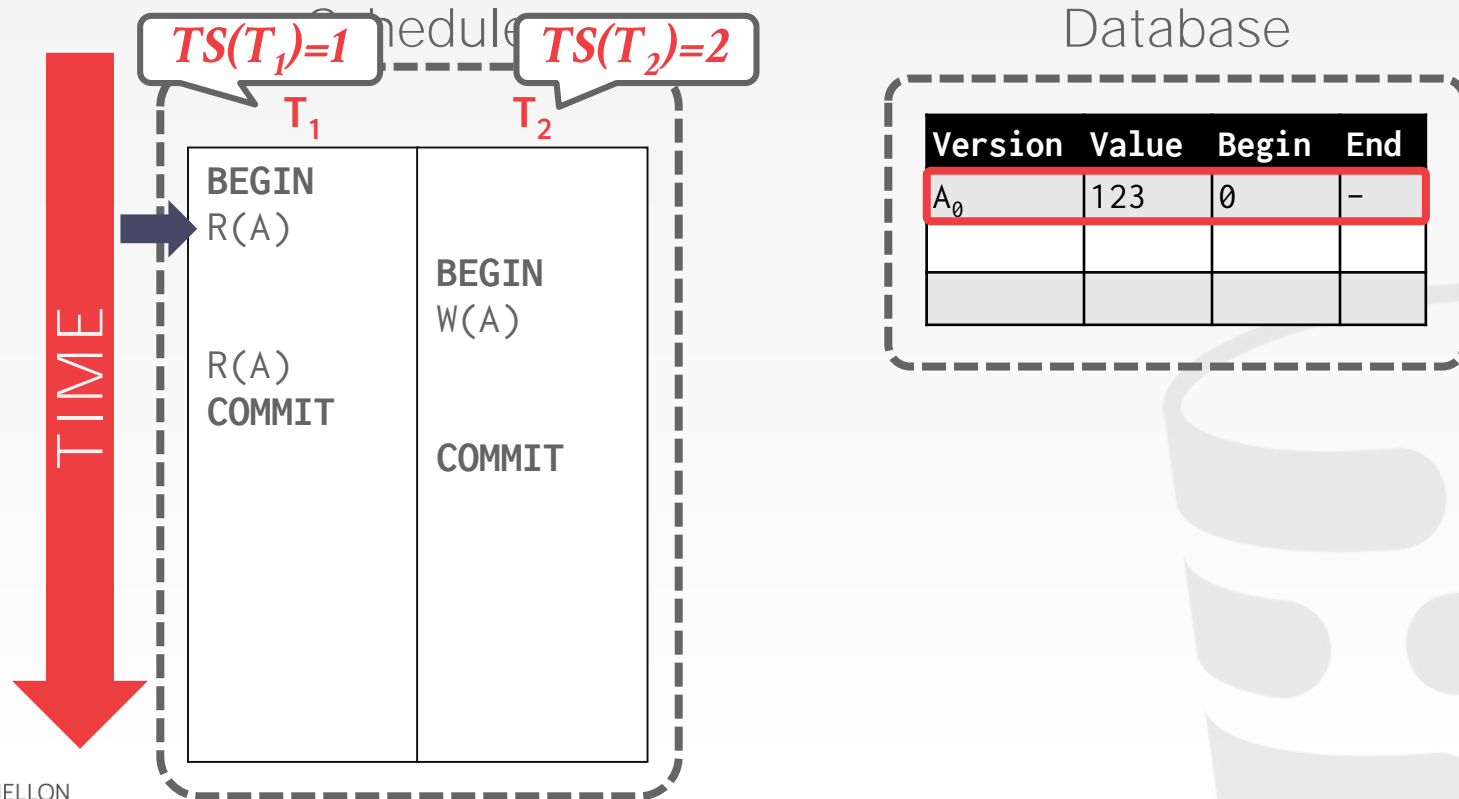


# MVCC – EXAMPLE #1

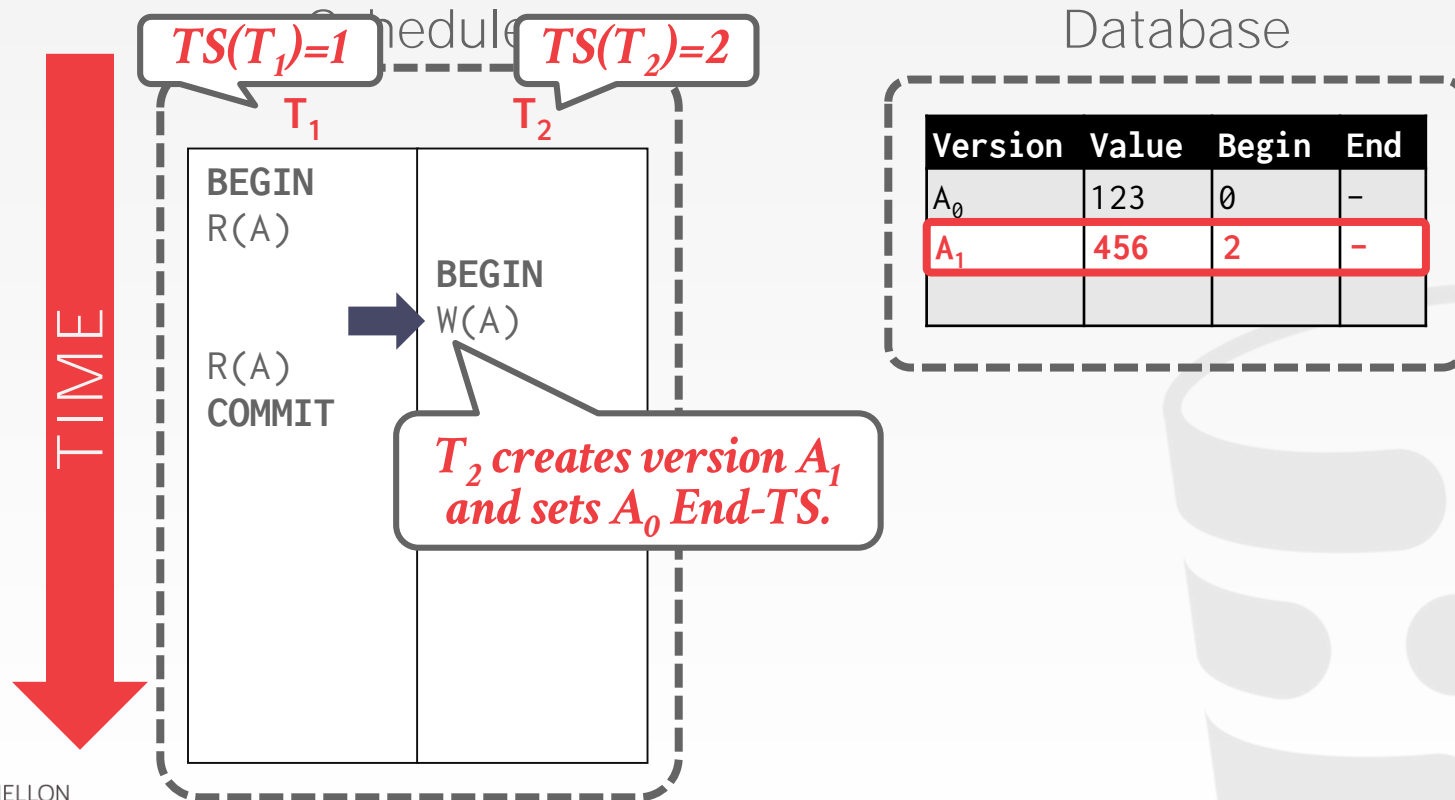




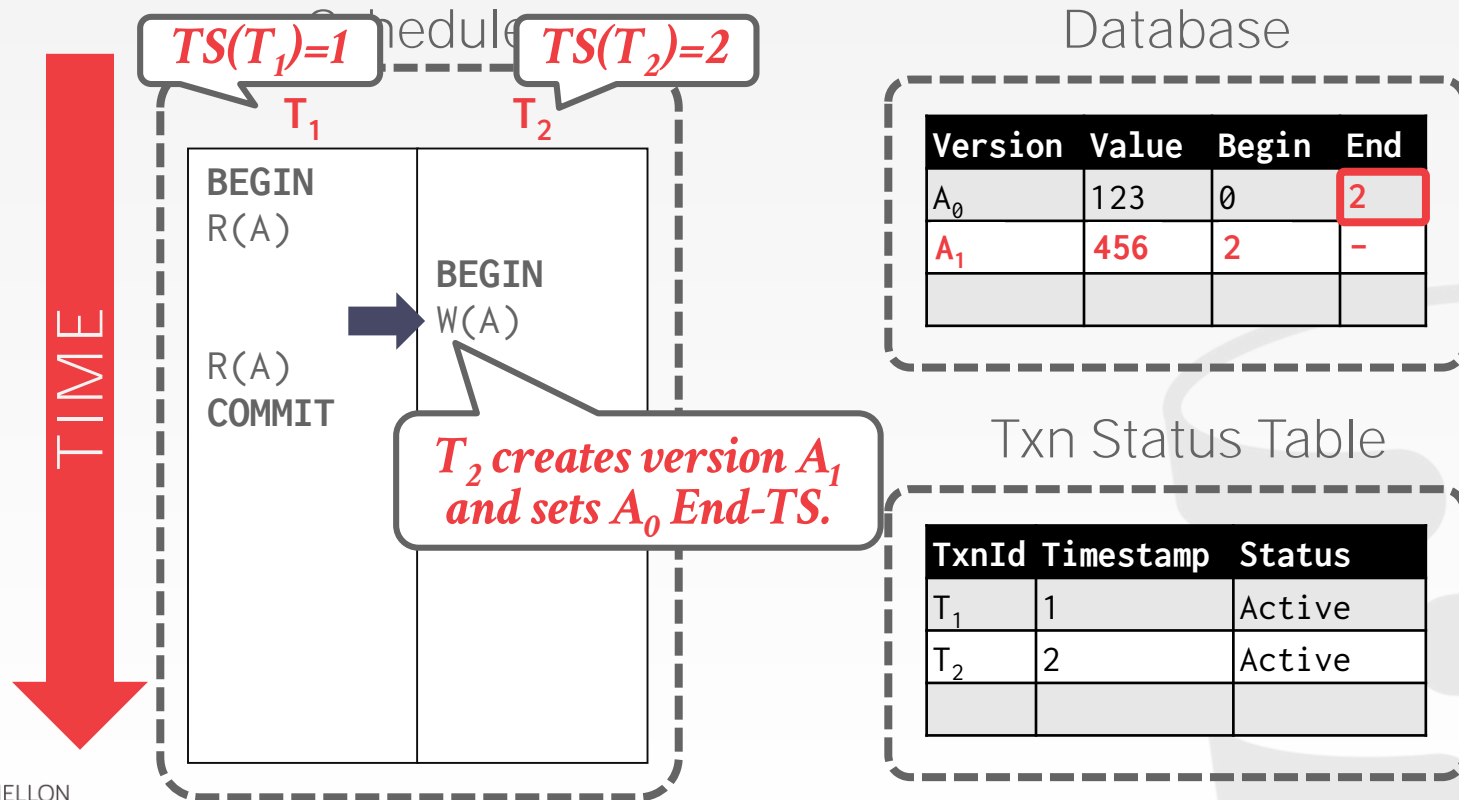
# MVCC – EXAMPLE #1



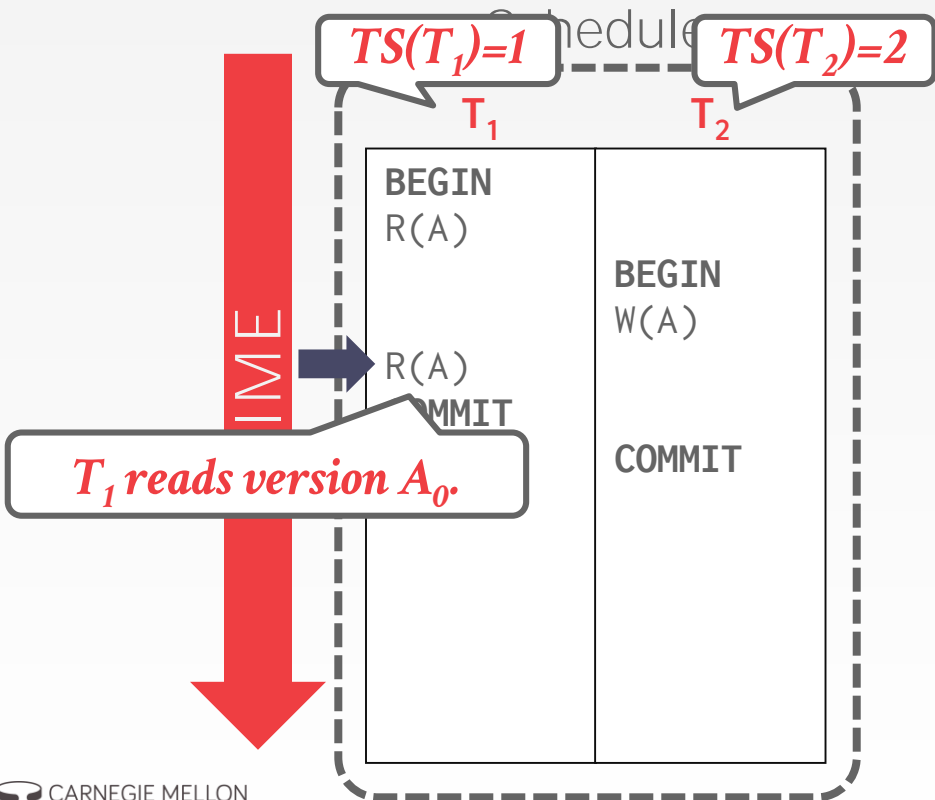
# MVCC – EXAMPLE #1



# MVCC – EXAMPLE #1



# MVCC – EXAMPLE #1



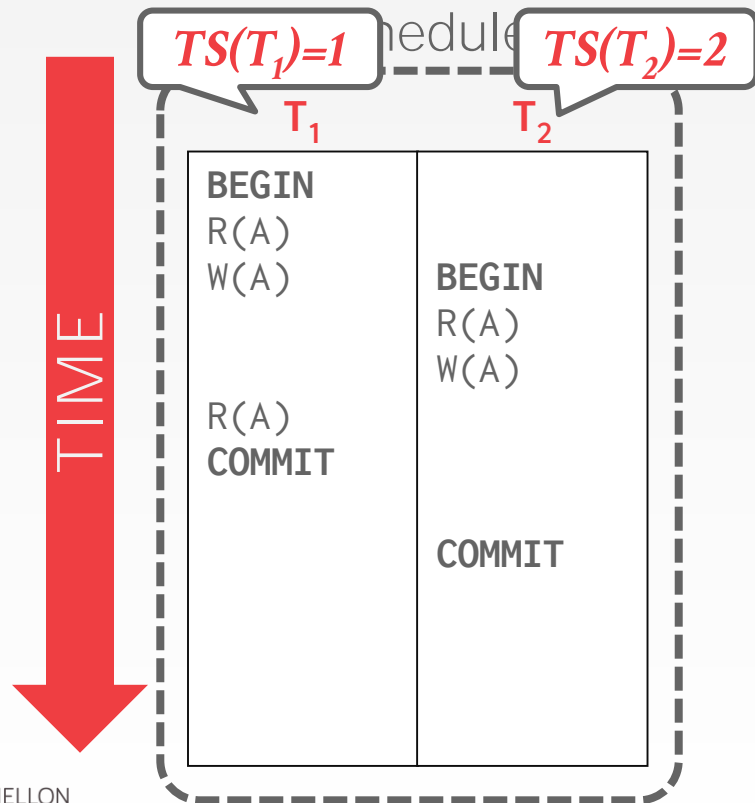
Database

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 #2



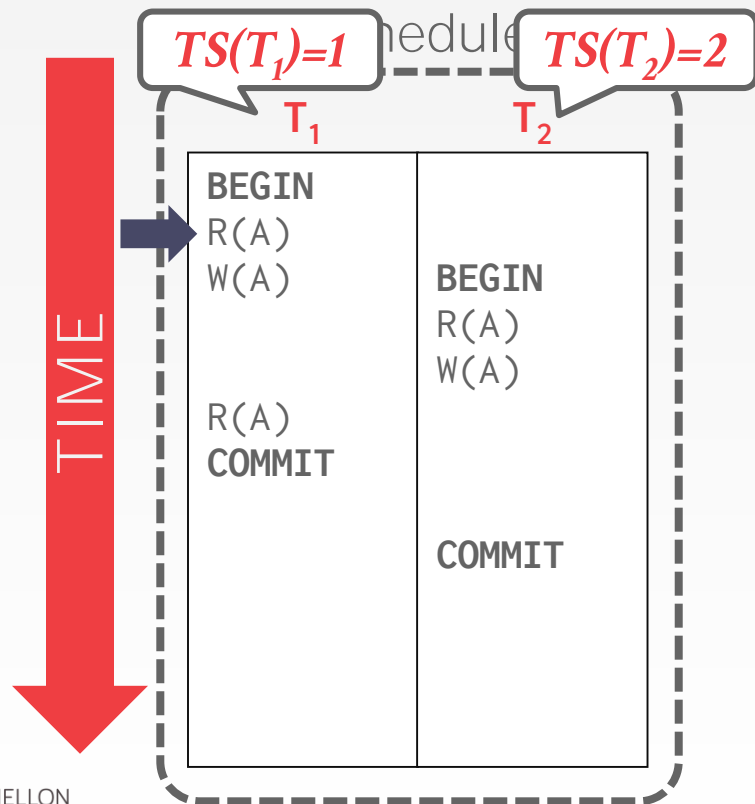
Database

Version	Value	Begin	End
$A_0$	123	0	

Txn Status Table

TxnId	Timestamp	Status
$T_1$	1	Active

# MVCC – EXAMPLE #2



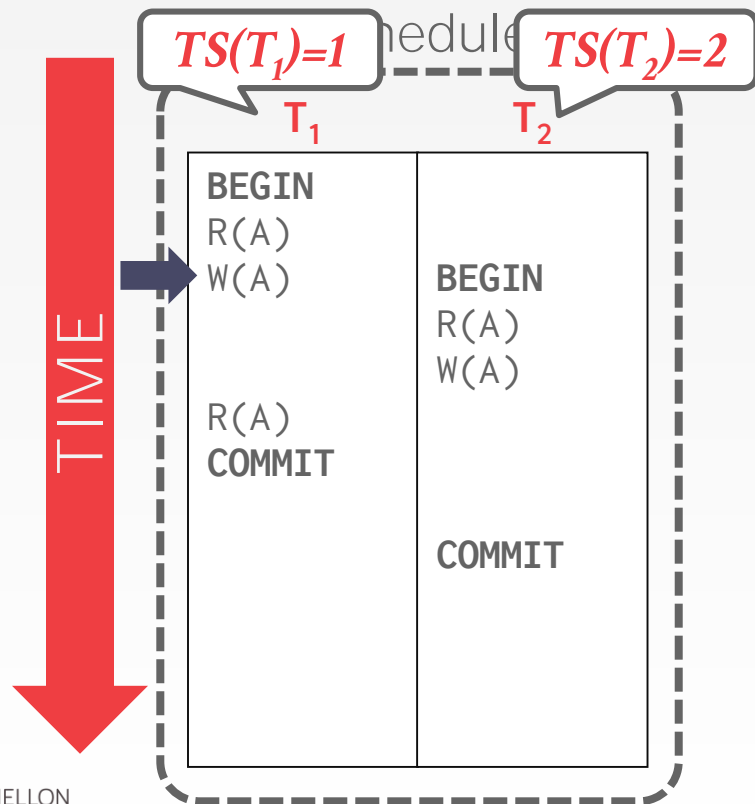
Database

Version	Value	Begin	End
$A_0$	123	0	

Txn Status Table

TxnId	Timestamp	Status
$T_1$	1	Active

# MVCC – EXAMPLE #2



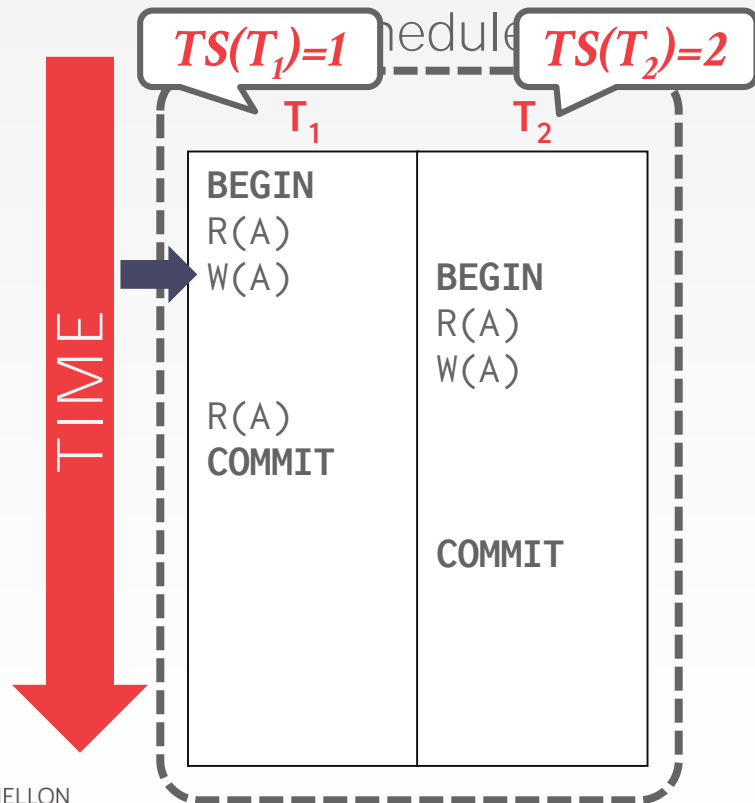
Database

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

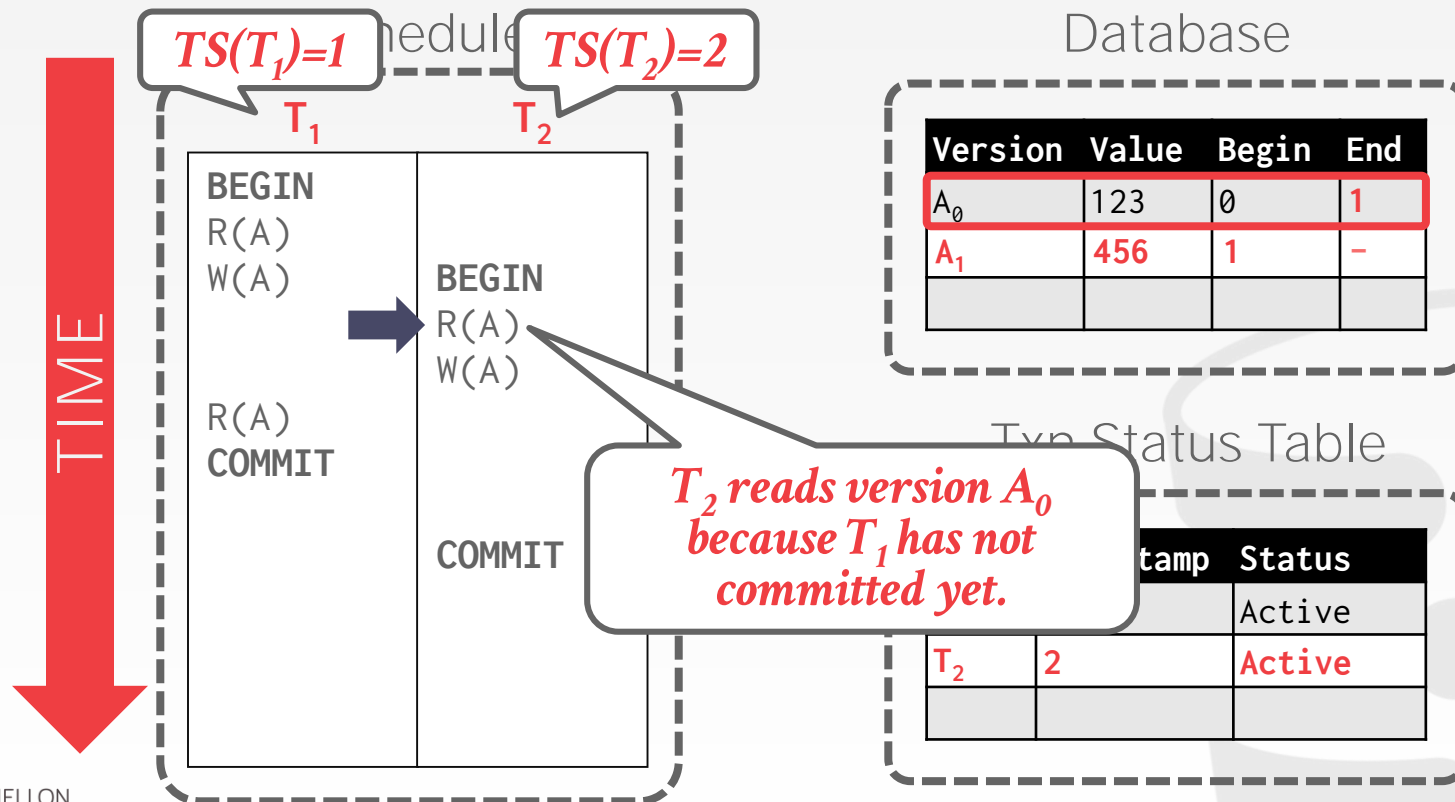
Version	Value	Begin	End
$A_0$	123	0	1
$A_1$	456	1	-

Txn Status Table

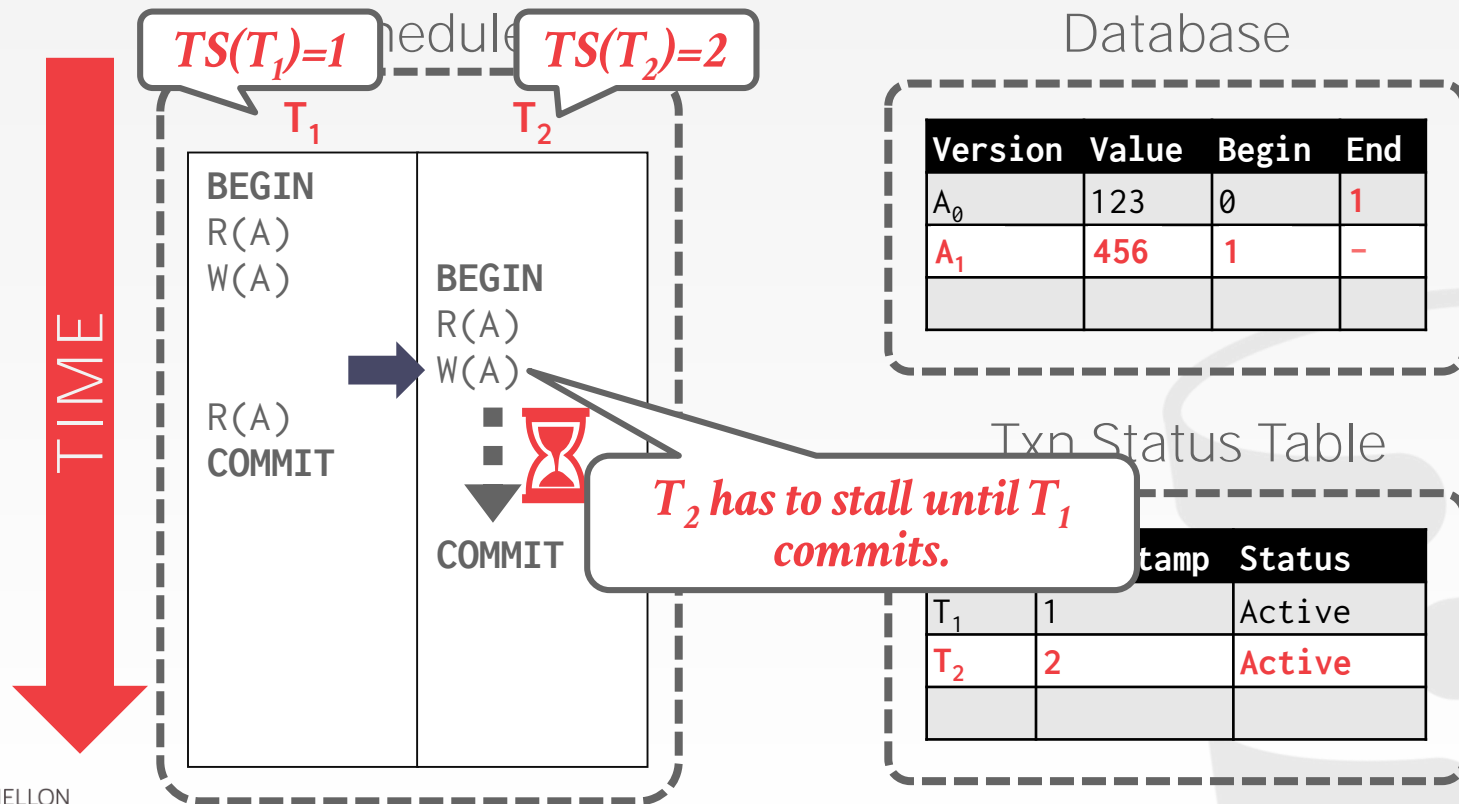
TxnId	Timestamp	Status
$T_1$	1	Active



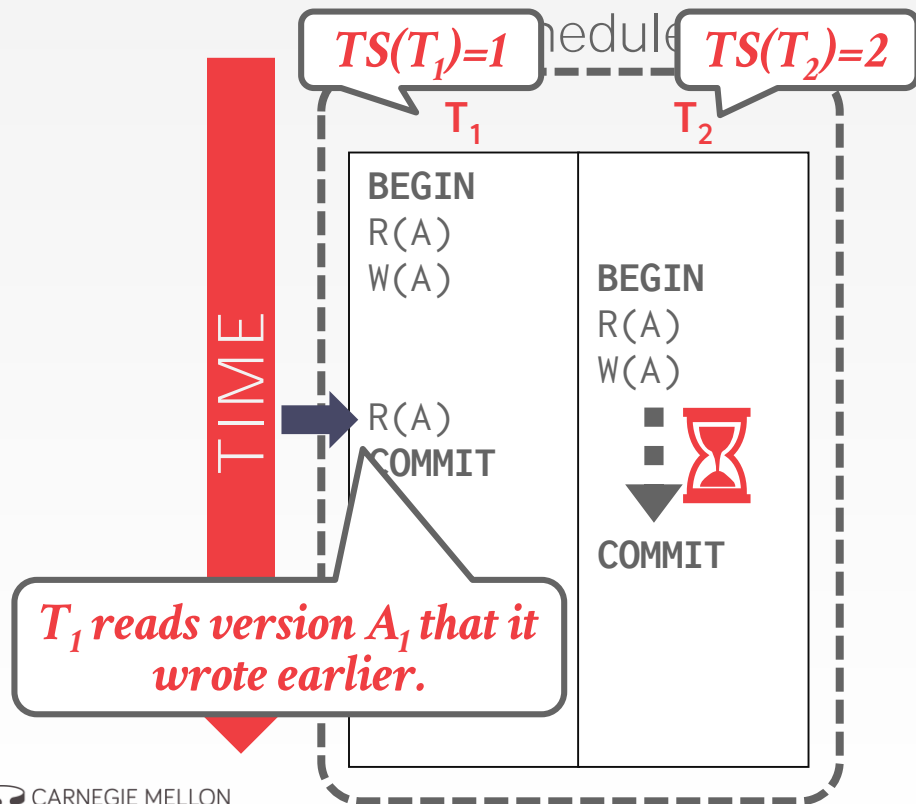
# MVCC – EXAMPLE #2



# MVCC – EXAMPLE #2



# MVCC – EXAMPLE #2



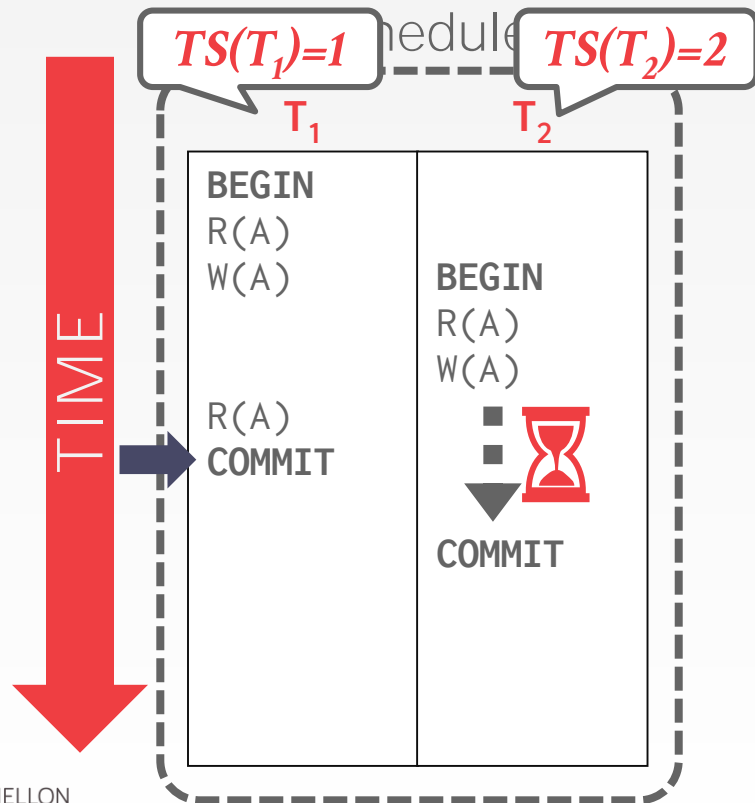
Database

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



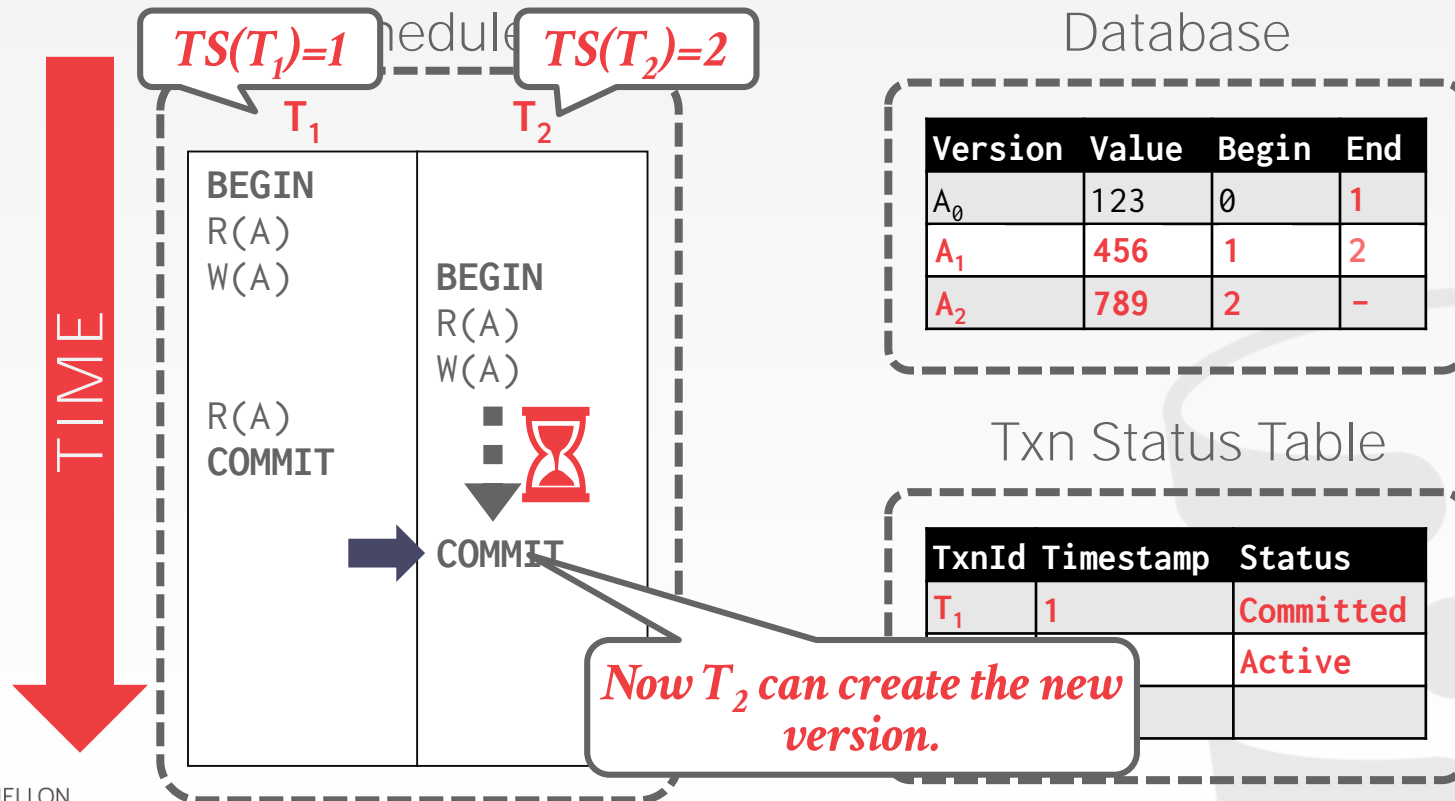
## Database

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



# 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



# 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.





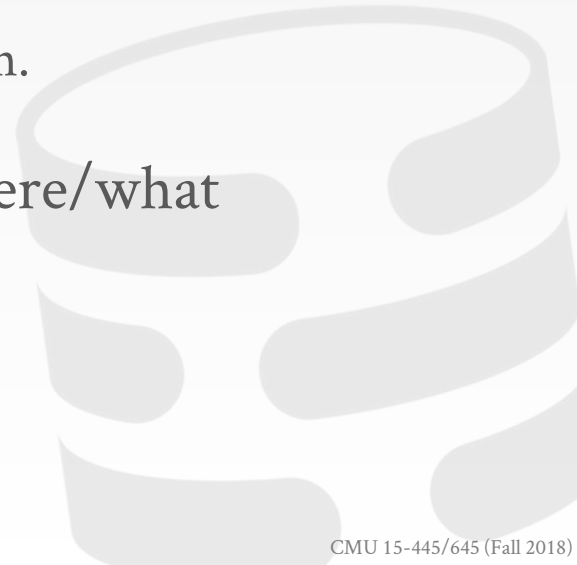
# VERSION STORAGE

---

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**

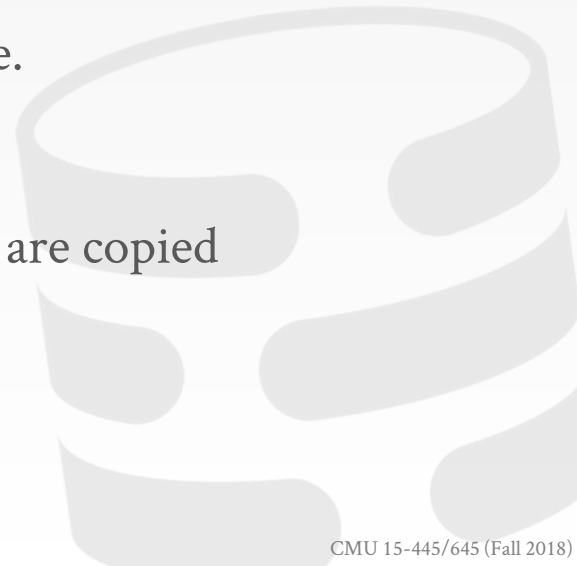
→ New versions are appended to the same table space.

## **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.

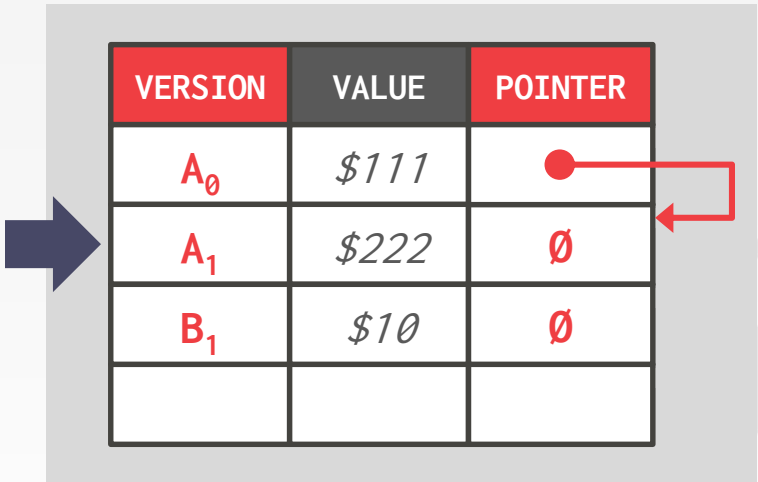



# APPEND-ONLY STORAGE

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

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

*Main Table*



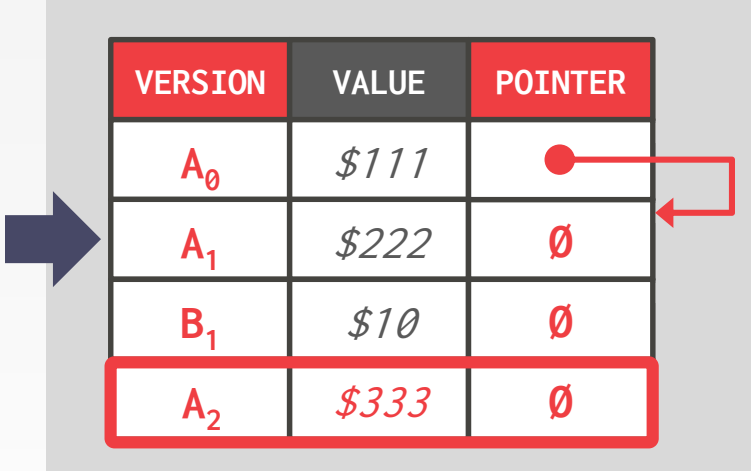
VERSION	VALUE	POINTER
$A_0$	\$111	
$A_1$	\$222	$\emptyset$
$B_1$	\$10	$\emptyset$

# APPEND-ONLY STORAGE

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

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

## *Main Table*



VERSION	VALUE	POINTER
$A_0$	\$111	●
$A_1$	\$222	$\emptyset$
$B_1$	\$10	$\emptyset$
$A_2$	\$333	$\emptyset$

# APPEND-ONLY STORAGE

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

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

*Main Table*

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

# APPEND-ONLY STORAGE

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

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

## *Main Table*

VERSION	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)**

- Just append new version to end of the chain.
- Have to traverse chain on look-ups.


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

- Have to update index pointers for every new version.
- Don't have to traverse chain on look ups.



# TIME-TRAVEL STORAGE

## Main Table



VERSION	VALUE	POINTER
$A_2$	\$222	● →
$B_1$	\$10	

## Time-Travel Table

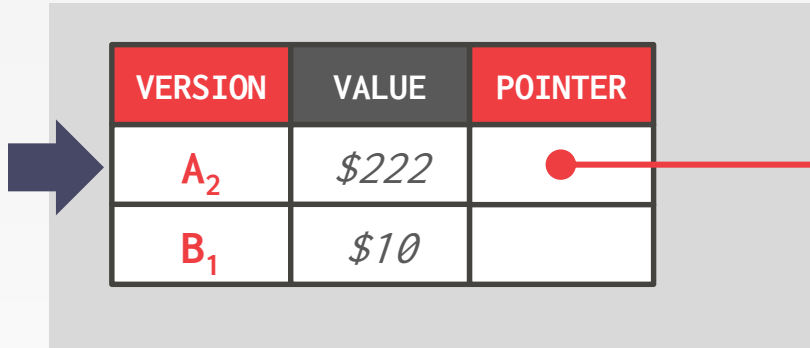
VERSION	VALUE	POINTER
$A_1$	\$111	∅

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



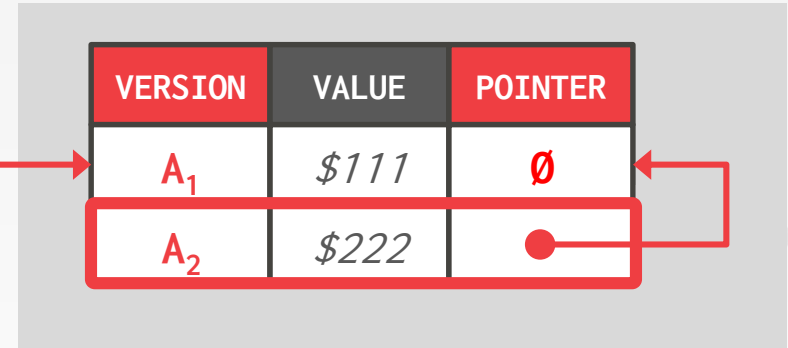
# TIME-TRAVEL STORAGE

## Main Table



VERSION	VALUE	POINTER
A <sub>2</sub>	\$222	●
B <sub>1</sub>	\$10	

## Time-Travel Table

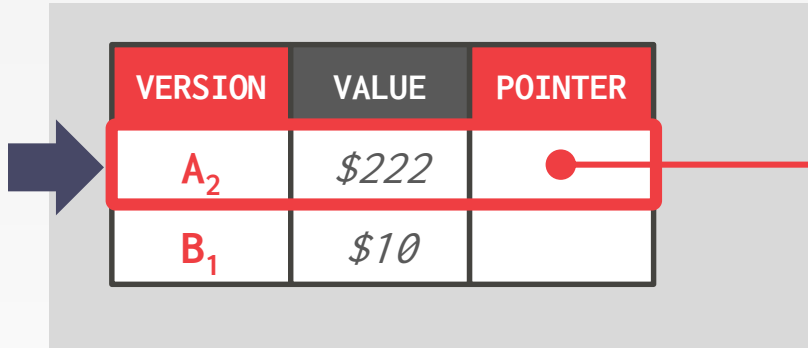


VERSION	VALUE	POINTER
A <sub>1</sub>	\$111	∅
A <sub>2</sub>	\$222	●

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

# TIME-TRAVEL STORAGE

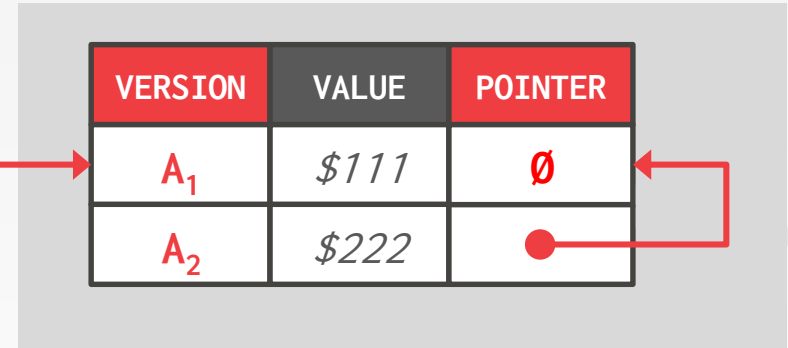
## Main Table



VERSION	VALUE	POINTER
A <sub>2</sub>	\$222	●
B <sub>1</sub>	\$10	

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

## Time-Travel Table

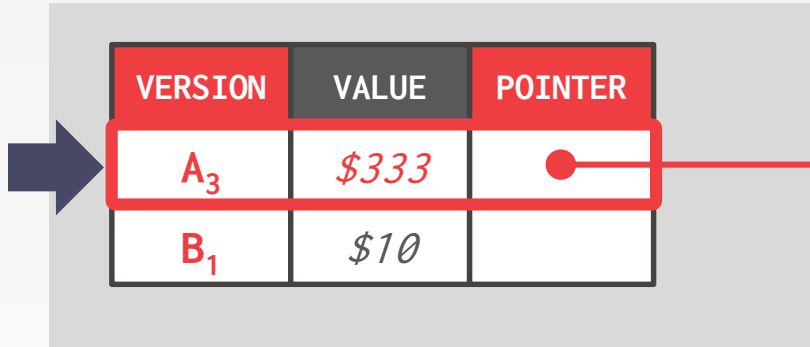


VERSION	VALUE	POINTER
A <sub>1</sub>	\$111	∅
A <sub>2</sub>	\$222	●

Overwrite master version in the main table.  
Update pointers.

# TIME-TRAVEL STORAGE

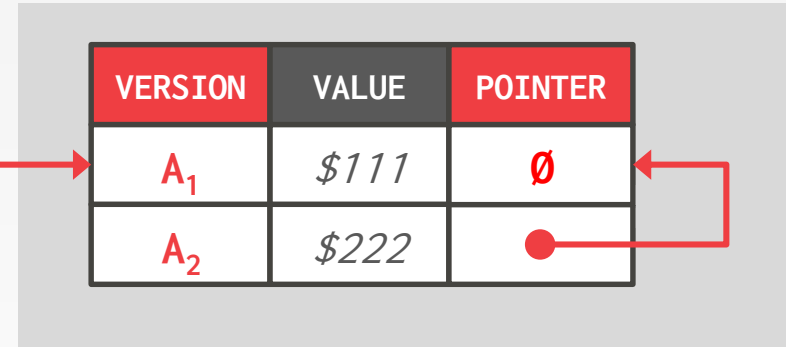
## Main Table



VERSION	VALUE	POINTER
A <sub>3</sub>	\$333	● →
B <sub>1</sub>	\$10	

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

## Time-Travel Table



VERSION	VALUE	POINTER
A <sub>1</sub>	\$111	∅
A <sub>2</sub>	\$222	● →

Overwrite master version in the main table.  
Update pointers.

# TIME-TRAVEL STORAGE

## Main Table

VERSION	VALUE	POINTER
A <sub>3</sub>	\$333	● →
B <sub>1</sub>	\$10	

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


## Time-Travel Table

VERSION	VALUE	POINTER
A <sub>1</sub>	\$111	∅
A <sub>2</sub>	\$222	● →

Overwrite master version in the main table.  
Update pointers.

# DELTA STORAGE

## Main Table



VERSION	VALUE	POINTER
A <sub>1</sub>	\$111	
B <sub>1</sub>	\$10	

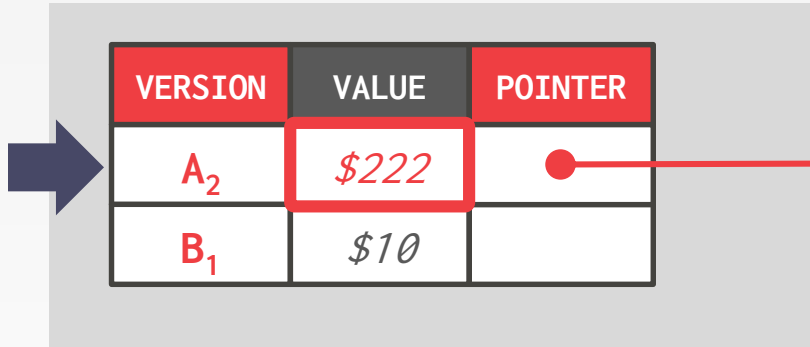
## Delta Storage Segment



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

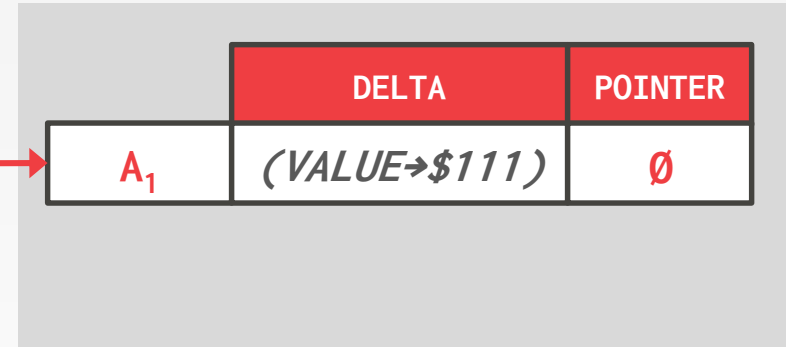
# DELTA STORAGE

## Main Table



VERSION	VALUE	POINTER
A <sub>2</sub>	\$222	● →
B <sub>1</sub>	\$10	

## Delta Storage Segment

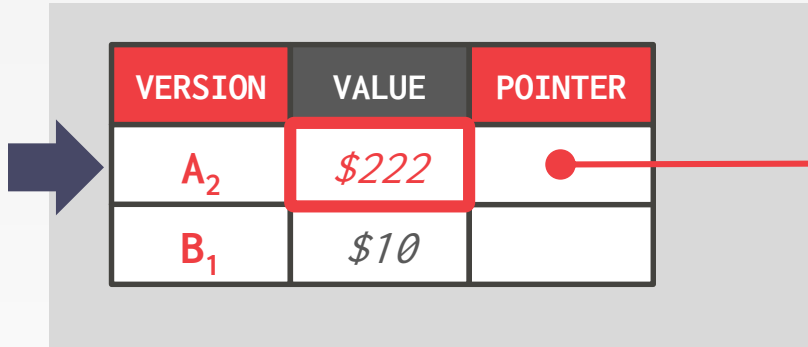


	DELTA	POINTER
A <sub>1</sub>	(VALUE→\$111)	∅

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

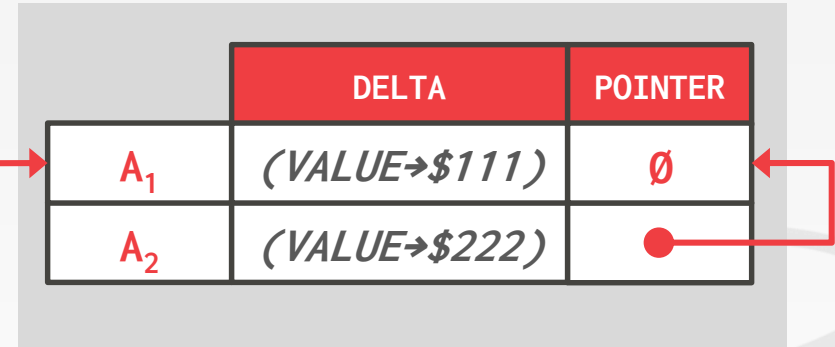
# DELTA STORAGE

## Main Table



VERSION	VALUE	POINTER
A <sub>2</sub>	\$222	● →
B <sub>1</sub>	\$10	

## Delta Storage Segment



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

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

# DELTA STORAGE

## Main Table

VERSION	VALUE	POINTER
A <sub>3</sub>	\$333	● →
B <sub>1</sub>	\$10	

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

## Delta Storage Segment

	DELTA	POINTER
A <sub>1</sub>	(VALUE->\$111)	∅
A <sub>2</sub>	(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

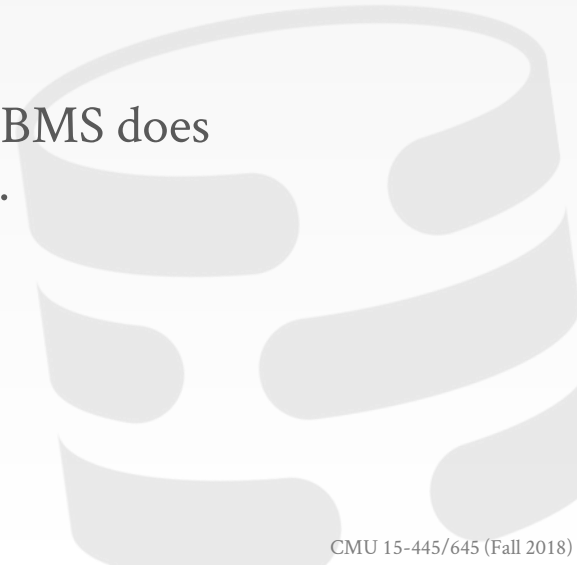
---

## 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

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



VERSION	BEGIN	END
$A_{100}$	1	9
$B_{100}$	1	9
$B_{101}$	10	20

## Background Vacuuming:

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

# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



VERSION	BEGIN	END
$A_{100}$	1	9
$B_{100}$	1	9
$B_{101}$	10	20

## Background Vacuuming:

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

# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



VERSION	BEGIN	END
$A_{100}$	1	9
$B_{100}$	1	9
$B_{101}$	10	20

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

# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



VERSION	BEGIN	END
$B_{101}$	10	20

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

# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



*Dirty Page BitMap*

VERSION	BEGIN	END
$B_{101}$	10	20

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

# TUPLE-LEVEL GC

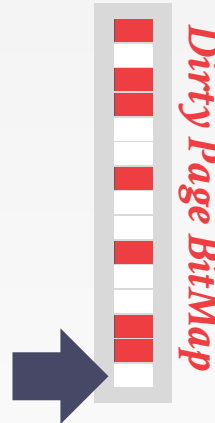
*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$

*Vacuum*



VERSION	BEGIN	END
$B_{101}$	10	20

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

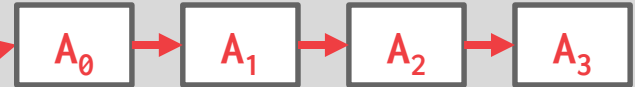


# TUPLE-LEVEL GC

*Thread #1*

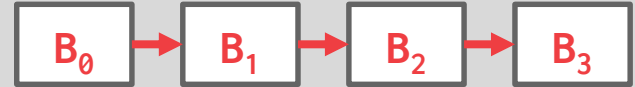
$TS(T_1)=12$

GET(A)



*Thread #2*

$TS(T_2)=25$



**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.

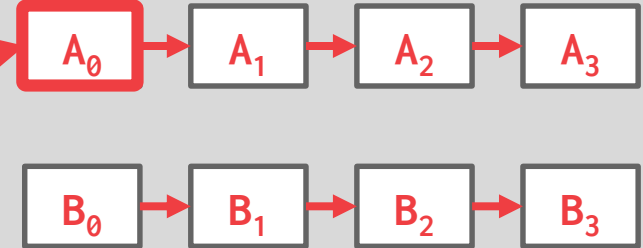
# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$



**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.

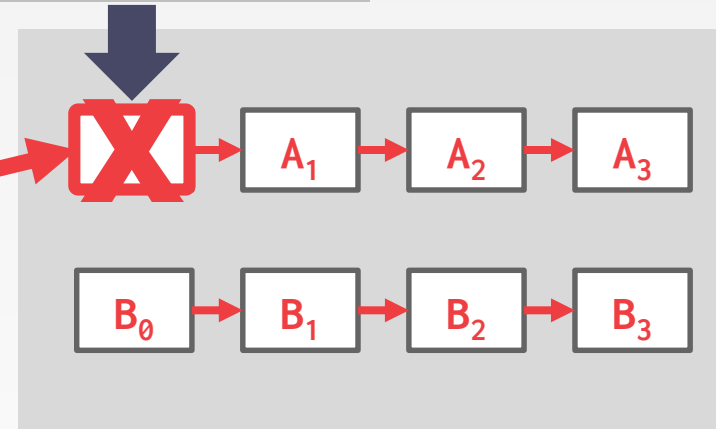
# TUPLE-LEVEL GC

*Thread #1*

$TS(T_1)=12$

*Thread #2*

$TS(T_2)=25$



**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

*Thread #1*

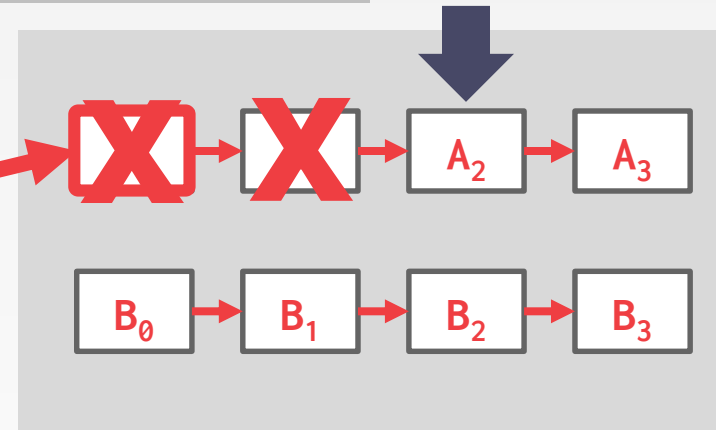
$TS(T_1)=12$

GET(A)



*Thread #2*

$TS(T_2)=25$



## 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

*Thread #1*

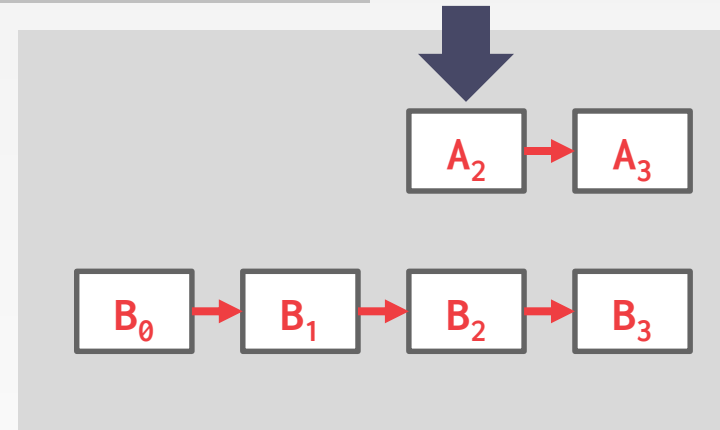
$TS(T_1)=12$

GET(A)



*Thread #2*

$TS(T_2)=25$



**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

*Thread #1*

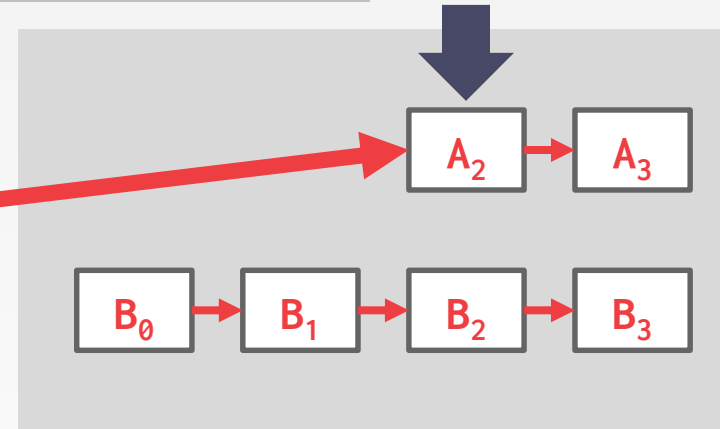
$TS(T_1)=12$

GET(A)



*Thread #2*

$TS(T_2)=25$



**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.

# TRANSACTION-LEVEL GC

---

Each txn keeps track of its read/write set.

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



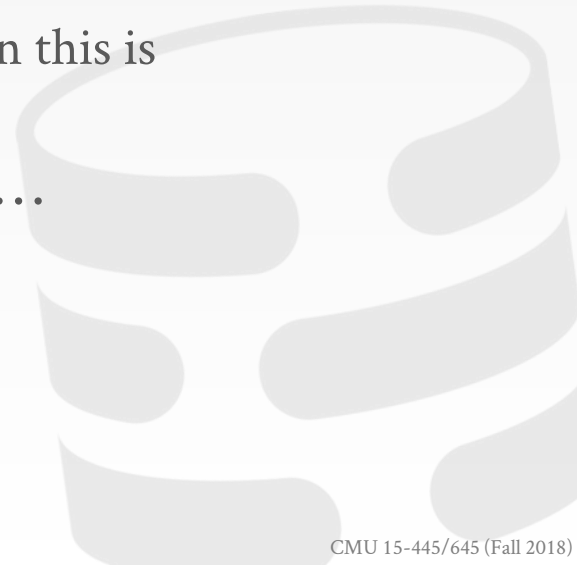
# INDEX MANAGEMENT

---

Primary key indexes point to version chain head.

- How often the DBMS has to 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 an **DELETE** followed by an **INSERT**.

Secondary indexes are more complicated...





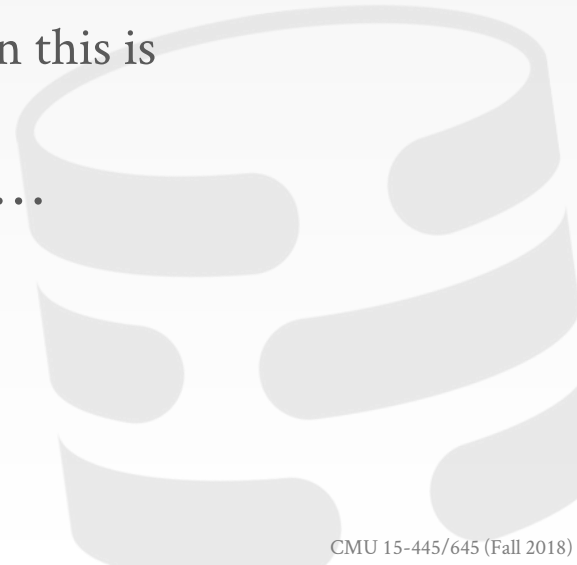
# INDEX MANAGEMENT

---

Primary key indexes point to version chain head.

- How often the DBMS has to 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 an **DELETE** followed by an **INSERT**.

Secondary indexes are more complicated...



# SECONDARY INDEXES

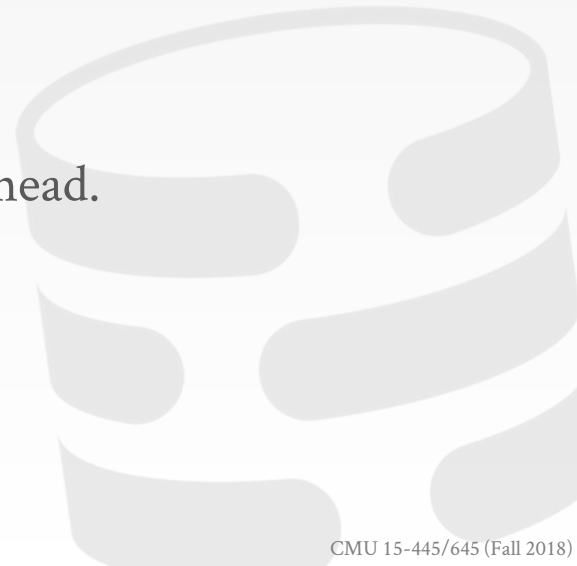
---

## **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

GET(A) ↓

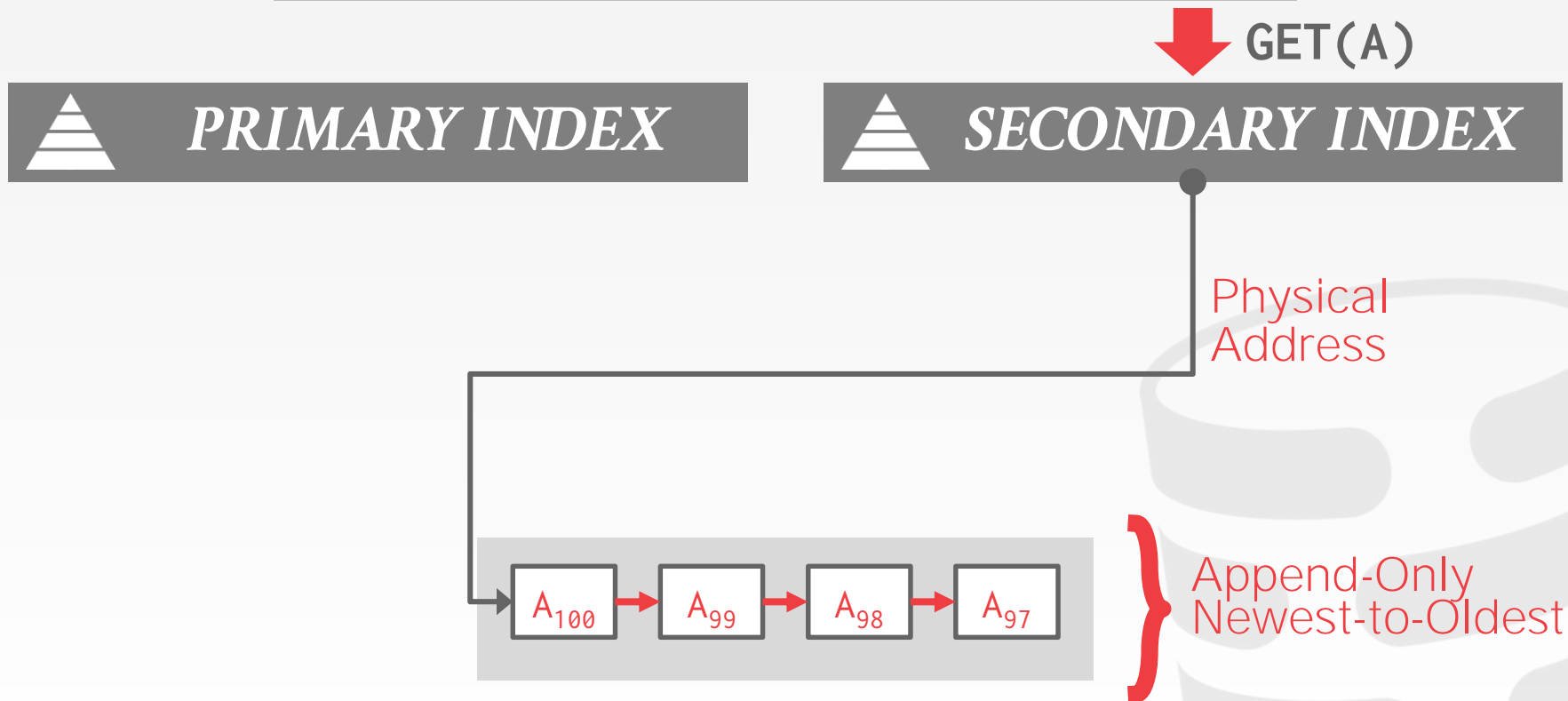


Physical  
Address

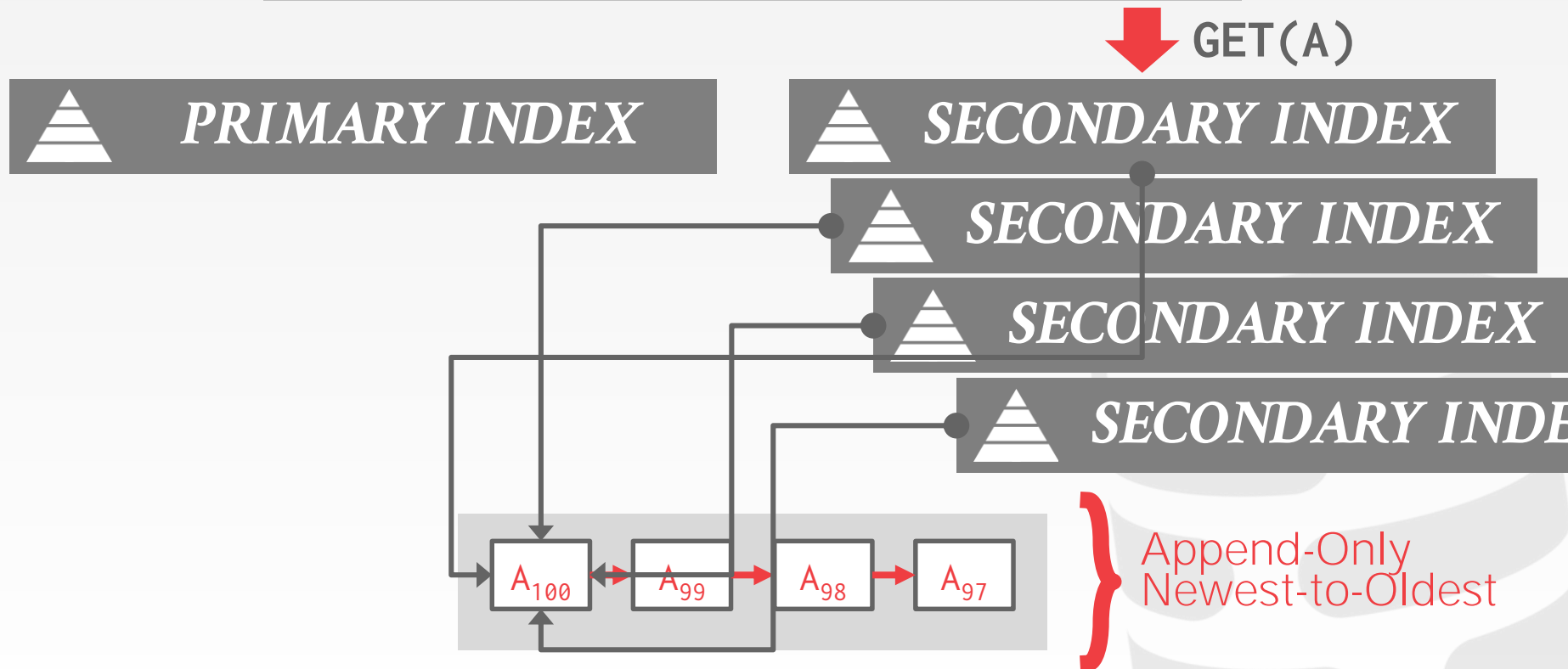


Append-Only  
Newest-to-Oldest

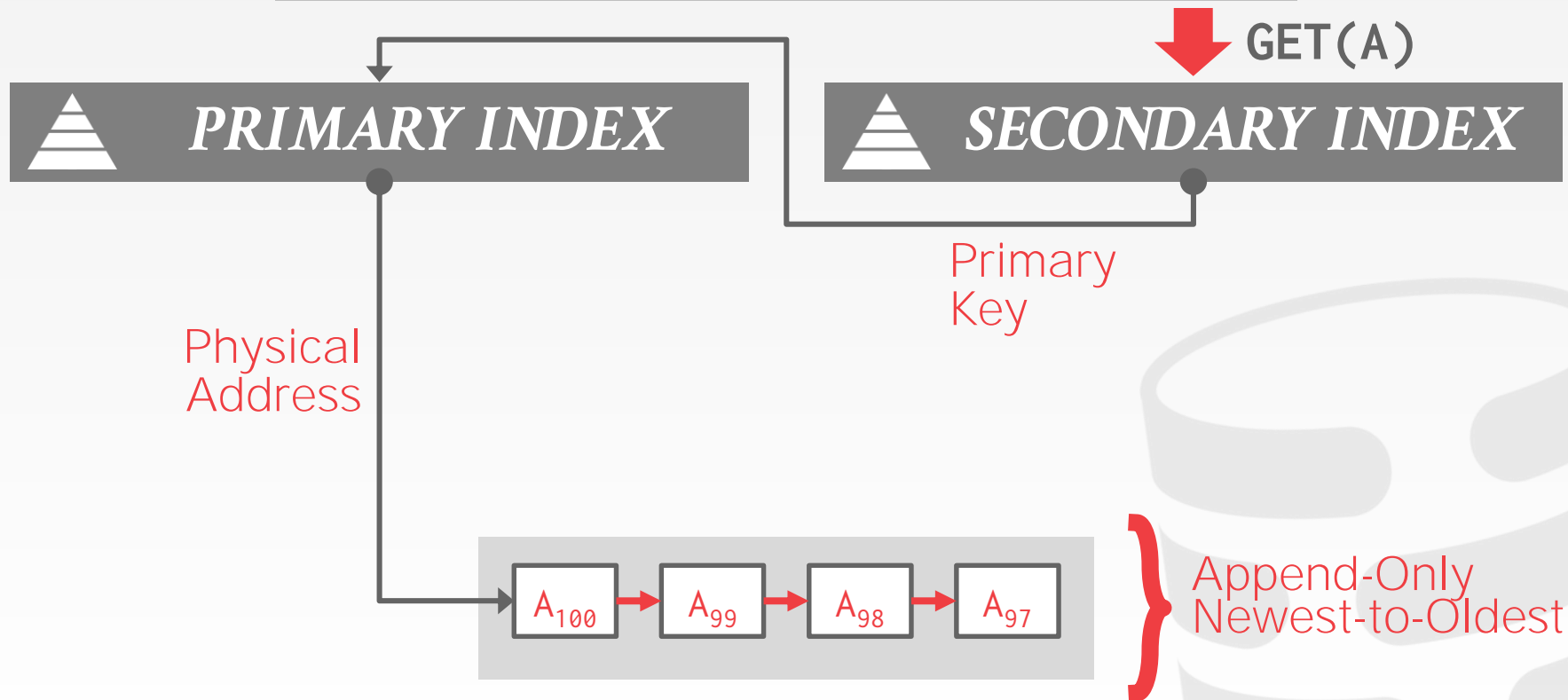
# INDEX POINTERS



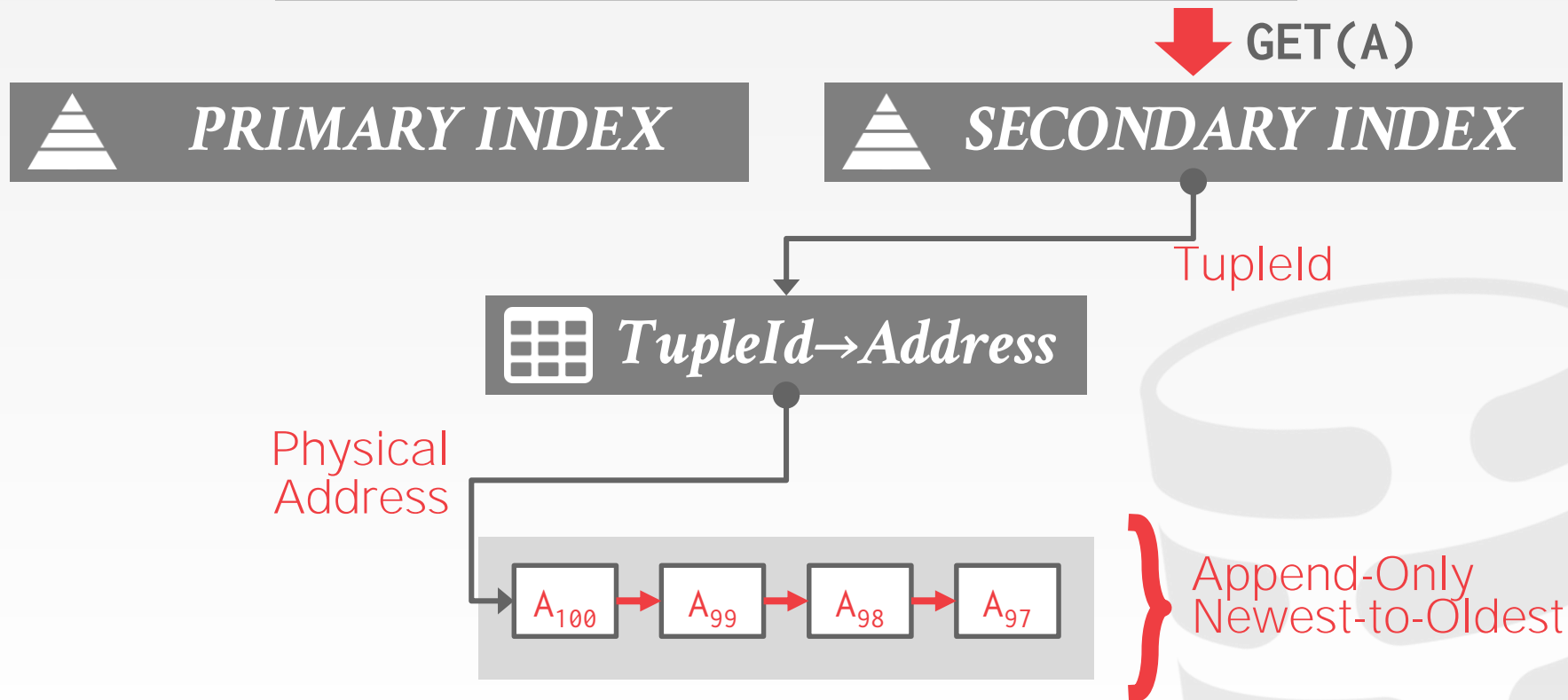
# INDEX POINTERS



# INDEX POINTERS



# INDEX POINTERS



# MVCC IMPLEMENTATIONS

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



# 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 & Recovery

