

# Database recovery

CS348 Spring 2024

Instructor: Sujaya Maiyya

Sections: **002 & 003 only**

# Announcements

- Assignment 3 due July 19<sup>th</sup>
- Final demo for projects:
  - Option 1: Online live demo with the TA
  - Option 2: Send a recording to the TA
- Send your choice to your TA by **July 22<sup>nd</sup>** (sooner is better)
  - Lose 2 points otherwise

# Review

- ACID

- **Atomicity**: TX's are either completely done or not done at all
- **Consistency**: TX's should leave the database in a consistent state
- **Isolation**: TX's must behave as if they are executed in isolation
- **Durability**: Effects of committed TX's are resilient against failures

- SQL transactions

**BEGIN TRANSACTION**

**SELECT ...;**

**UPDATE ...;**

**ROLLBACK | COMMIT;**

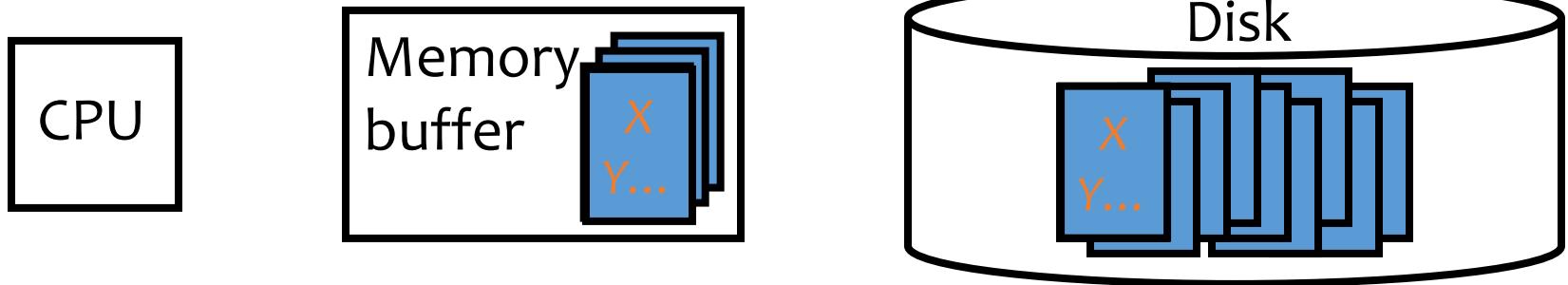
# Outline

- Recovery – atomicity and durability
  - Naïve approaches
  - Logging for undo and redo

# Execution model

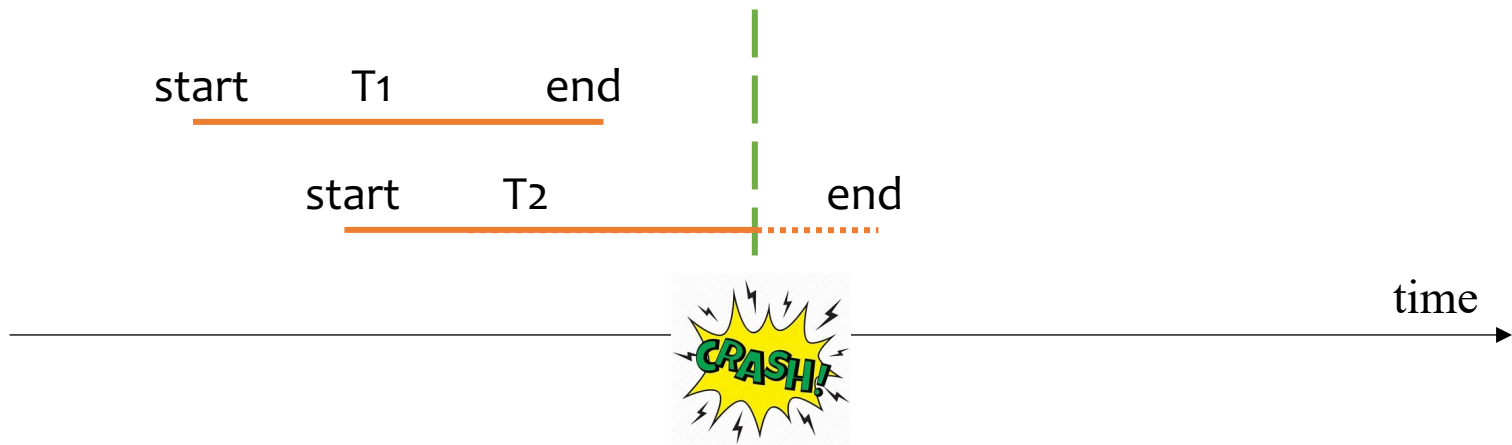
To read/write  $X$

- The disk block containing  $X$  must be first brought into memory
- $X$  is read/written in memory
- The memory block containing  $X$ , if modified, must be written back (flushed) to disk eventually



# Failures

- System crashes right after a transaction  $T_1$  commits; **but not all effects of  $T_1$  were written to disk**
  - How do we complete/redo  $T_1$  (**durability**)?
- System crashes in the middle of a transaction  $T_2$ ; **partial effects of  $T_2$  were written to disk**
  - How do we undo  $T_2$  (**atomicity**)?



# Naïve approach: Force -- durability

*T1* (balance transfer of \$100 from *A* to *B*)

```
read(A, a); a = a - 100;
```

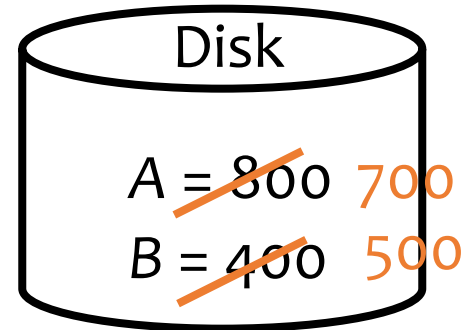
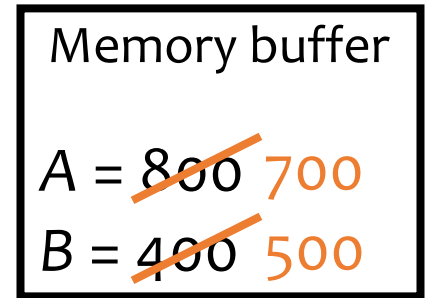
```
write(A, a);
```

```
read(B, b); b = b + 100;
```

```
write(B, b);
```

```
commit;
```

**Force:** all writes must be reflected on disk when a transaction commits



# Naïve approach: Force -- durability

*T1* (balance transfer of \$100 from *A* to *B*)

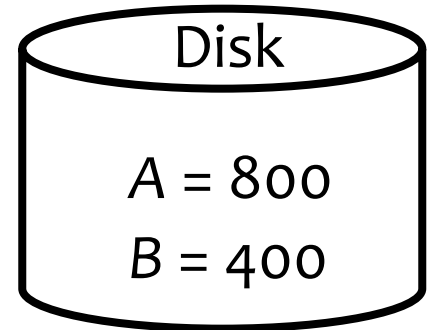
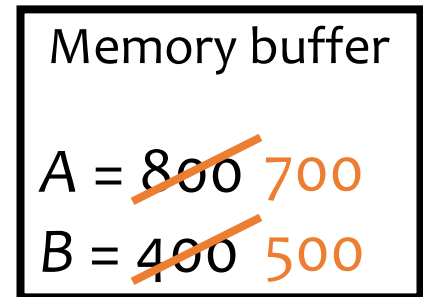
```
read(A, a); a = a - 100;
```

```
write(A, a);
```

```
read(B, b); b = b + 100;
```

```
write(B, b);
```

```
commit;
```



**Force:** all writes must be reflected on disk when a transaction commits

Without force: not all writes are on disk when *T1* commits **Bad!**

If system crashes right after *T1* commits, effects of *T1* will be lost



# Naïve approach: No steal -- atomicity

*T1* (balance transfer of \$100 from *A* to *B*)

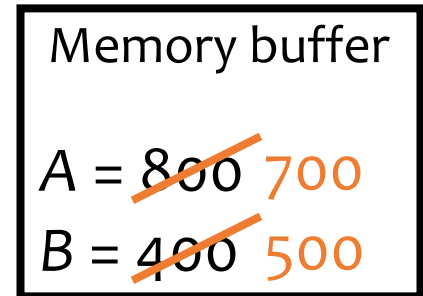
```
read(A, a); a = a - 100;
```

```
write(A, a);
```

```
read(B, b); b = b + 100;
```

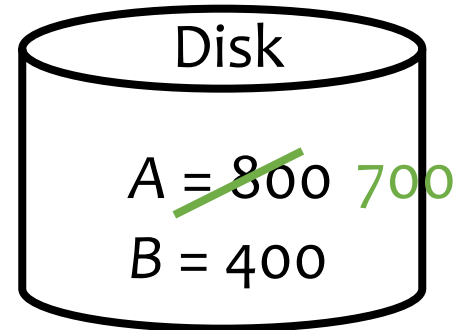
```
write(B, b);
```

```
commit;
```



**No steal:** Writes of a transaction can only be flushed to disk at commit time:

- e.g.  $A=700$  cannot be flushed to disk before commit.



Bad!

With steal: some writes are on disk before *T* commits

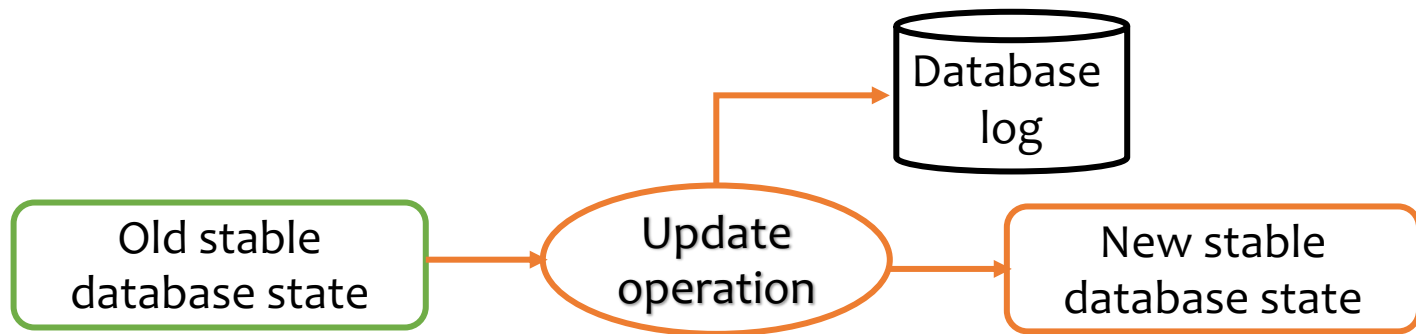
If system crashes before *T1* commits, there is no way to undo the changes

# Naïve approach

- **Force**: When a transaction commits, all writes of this transaction must be reflected on disk
  - Ensures durability
  - ☞ Problem of force: Lots of **random writes** hurt performance
- **No steal**: Writes of a transaction can only be flushed to disk at commit time
  - Ensures atomicity
  - ☞ Problem of no steal: Holding on to all dirty blocks **requires lots of memory**

# Logging

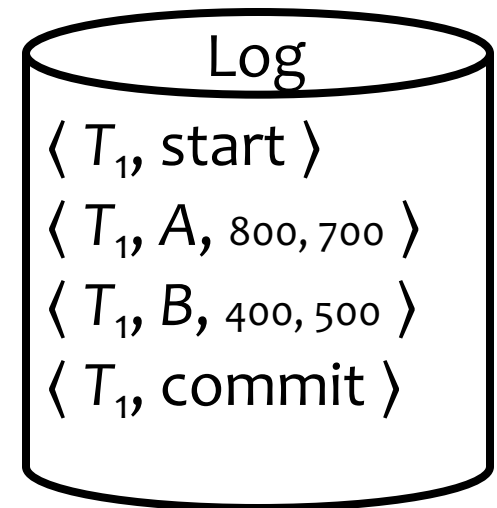
- **Database log**: sequence of **log records**, recording all changes made to the database, written to stable storage (e.g., disk) during normal operation



- Hey, one change turns into two—bad for performance?
  - But writes to log are **sequential** (append to the end of log)

# Log format

- When a transaction  $T_i$  starts
  - $\langle T_i, \text{start} \rangle$
- Record values before and after each modification:
  - $\langle T_i, X, \text{old\_value\_of\_}X, \text{new\_value\_of\_}X \rangle$
  - $T_i$  is transaction id
  - $X$  identifies the data item
- A transaction  $T_i$  is committed when its commit log record is written to disk
  - $\langle T_i, \text{commit} \rangle$



# When to write log records into stable store?

- Before  $X$  is modified or after?
- **Write-ahead logging (WAL)**: Before  $X$  is modified on disk, the log record pertaining to  $X$  must be flushed
- Without WAL, system might crash after  $X$  is modified on disk but before its log record is written to disk—no way to undo

# Undo/redo logging example

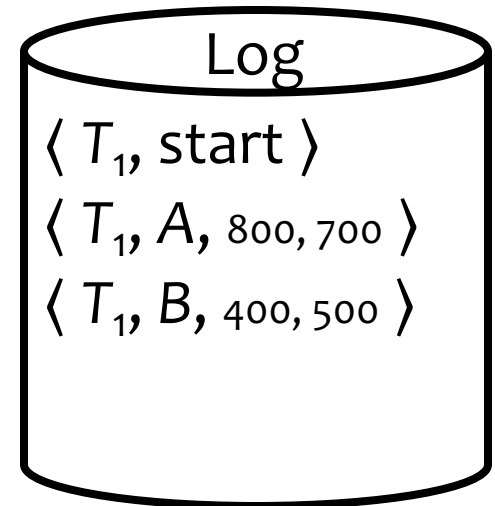
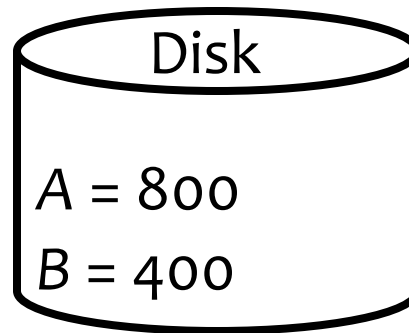
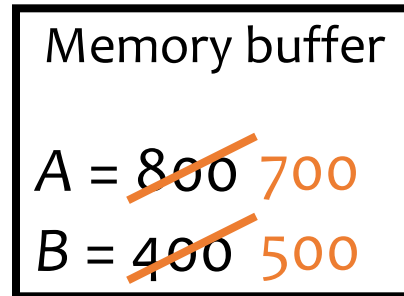
$T_1$  (balance transfer of \$100 from  $A$  to  $B$ )

read( $A, a$ );  $a = a - 100$ ;

write( $A, a$ );

read( $B, b$ );  $b = b + 100$ ;

write( $B, b$ );



WAL: Before  $A, B$  are modified on disk, their log info must be flushed

# Undo/redo logging example cont.

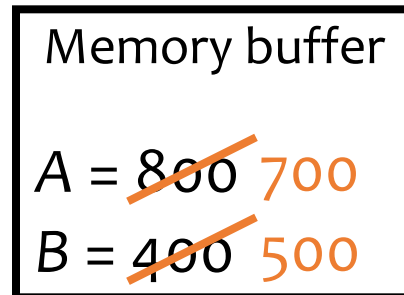
*T1* (balance transfer of \$100 from *A* to *B*)

read(*A*, *a*); *a* = *a* - 100;

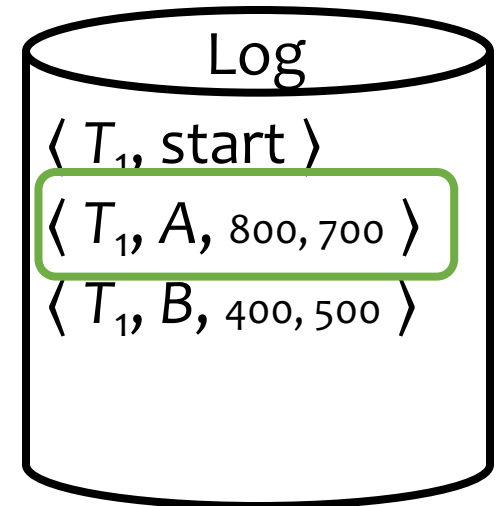
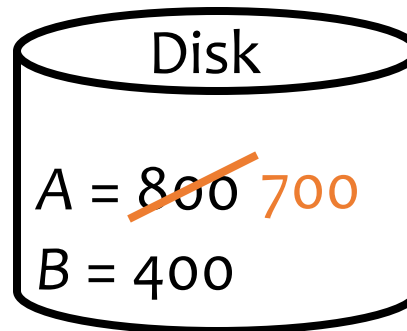
write(*A*, *a*);

read(*B*, *b*); *b* = *b* + 100;

write(*B*, *b*);



Steal: can flush  
before commit



If system crashes before *T1* commits, we have the old value of *A* stored on the log to **undo** *T1*

# Undo/redo logging example cont.

*T1* (balance transfer of \$100 from *A* to *B*)

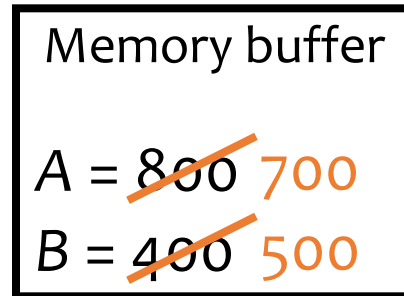
read(*A*, *a*); *a* = *a* - 100;

write(*A*, *a*);

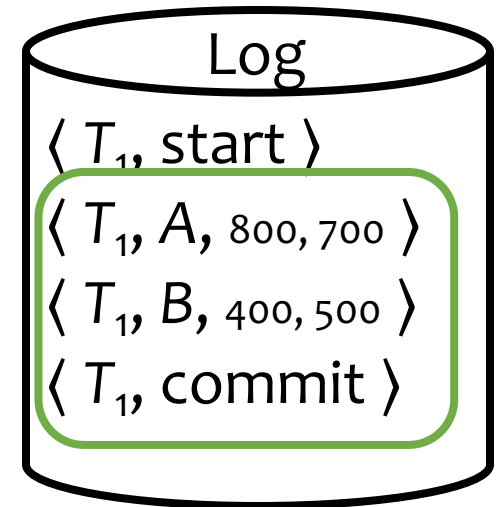
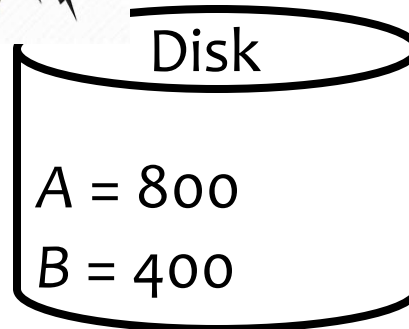
read(*B*, *b*); *b* = *b* + 100;

write(*B*, *b*);

commit;



No force: can flush  
after commit

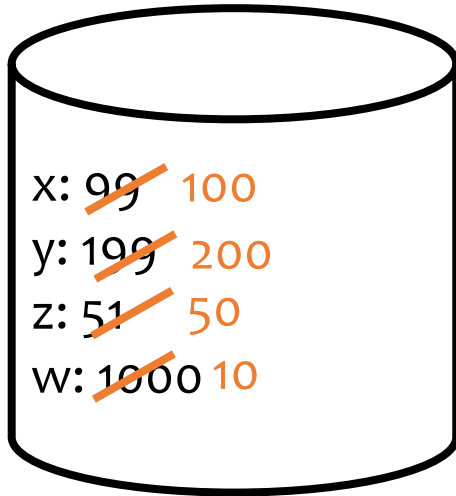


If system crashes before we flush the changes of *A*, *B* to the disk, we have their new committed values on the log to **redo** *T1*



# Log example - redo

- Redo phase:



List of active transactions at crash:

T1 T2 T3

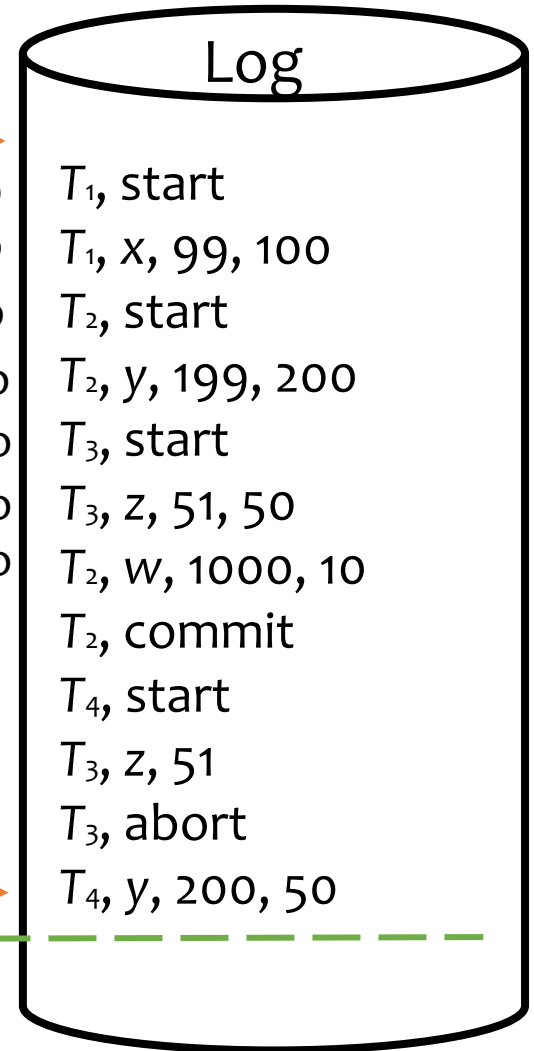


Start of log



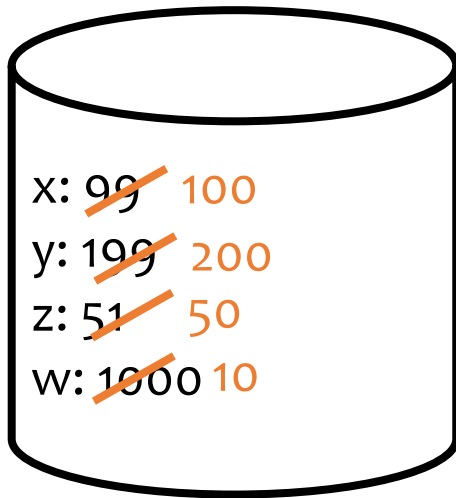
redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
T<sub>2</sub>, commit  
T<sub>4</sub>, start  
T<sub>3</sub>, z, 51  
T<sub>3</sub>, abort  
T<sub>4</sub>, y, 200, 50

End of log



# Log example

- Redo phase:



List of active transactions at crash:

~~T1~~ ~~T2~~ T3

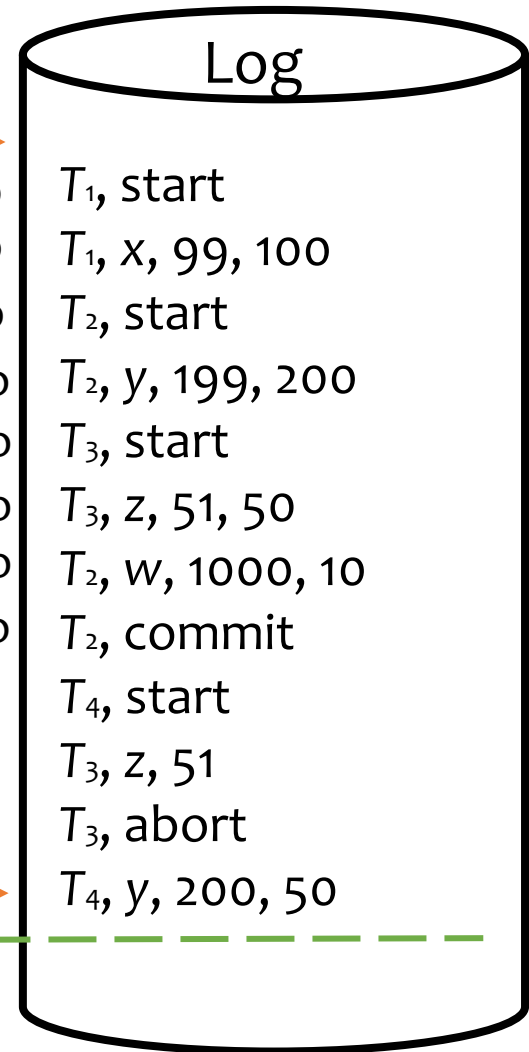


Start of log



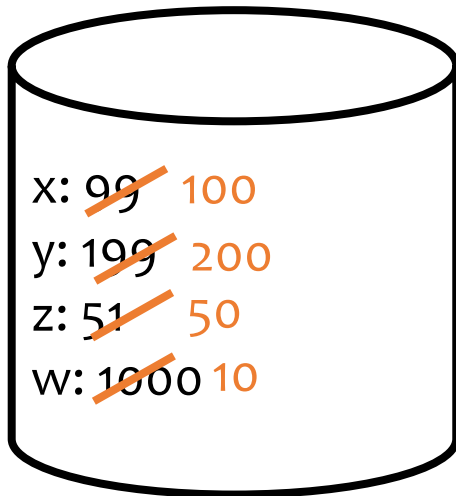
redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
redo T<sub>2</sub>, commit  
T<sub>4</sub>, start  
T<sub>3</sub>, z, 51  
T<sub>3</sub>, abort  
T<sub>4</sub>, y, 200, 50

End of log



# Log example

- Redo phase:



List of active transactions at crash:

T1 ~~T2~~ T3 T4

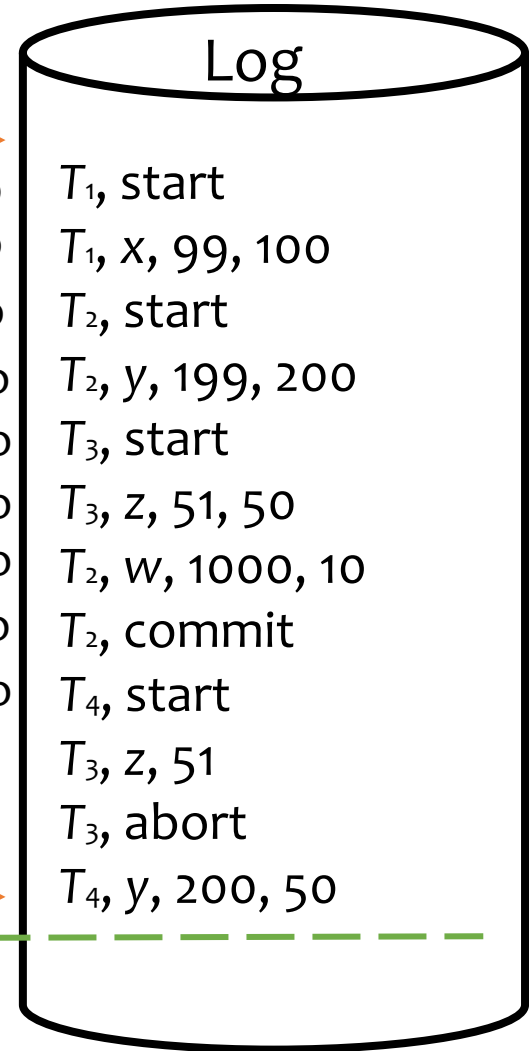


Start of log



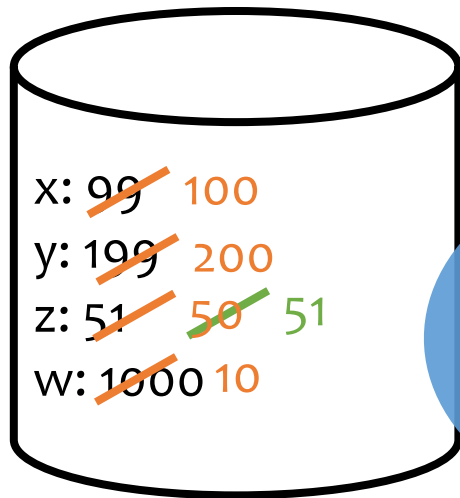
redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
redo T<sub>2</sub>, commit  
redo T<sub>4</sub>, start  
T<sub>3</sub>, z, 51  
T<sub>3</sub>, abort  
T<sub>4</sub>, y, 200, 50

End of log



# Log example

- Redo phase:



When txn manager receives abort, it logs reverse operations before abort

List of active transactions at crash:

T1 ~~T2~~ T3 T4



Start of log



redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
redo T<sub>2</sub>, commit  
redo T<sub>4</sub>, start  
redo T<sub>3</sub>, z, 51

End of log

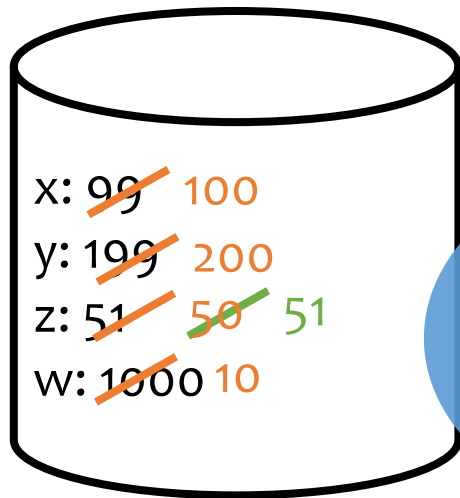


Log

T<sub>3</sub>, abort  
T<sub>4</sub>, y, 200, 50

# Log example

- Redo phase:



List of active transactions at crash:

T1 ~~T2~~ ~~T3~~ T4

When txn manager receives abort, it logs reverse operations before abort

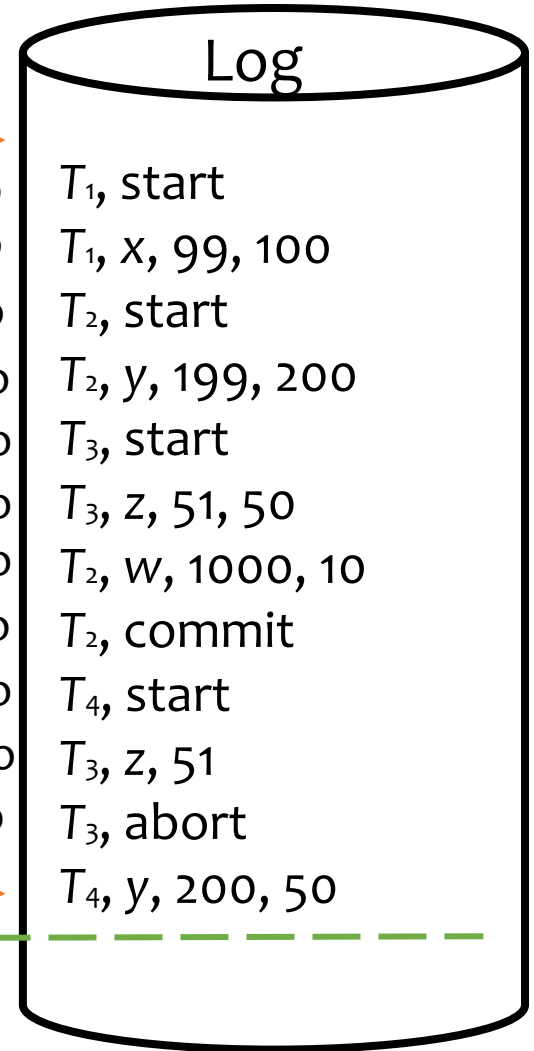


Start of log



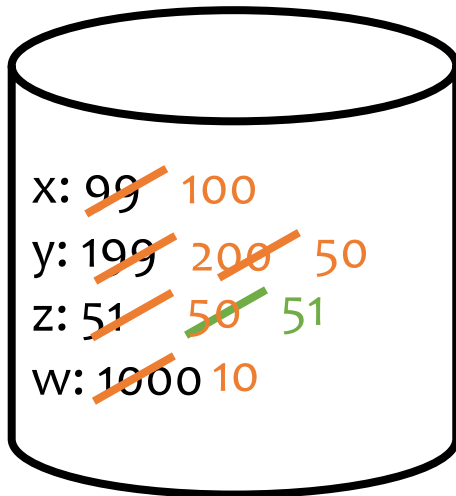
redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
redo T<sub>2</sub>, commit  
redo T<sub>4</sub>, start  
redo T<sub>3</sub>, z, 51  
redo T<sub>3</sub>, abort  
redo T<sub>4</sub>, y, 200, 50

End of log



# Log example

- Redo phase:



List of active transactions at crash:

T1 ~~T2~~ ~~T3~~ T4

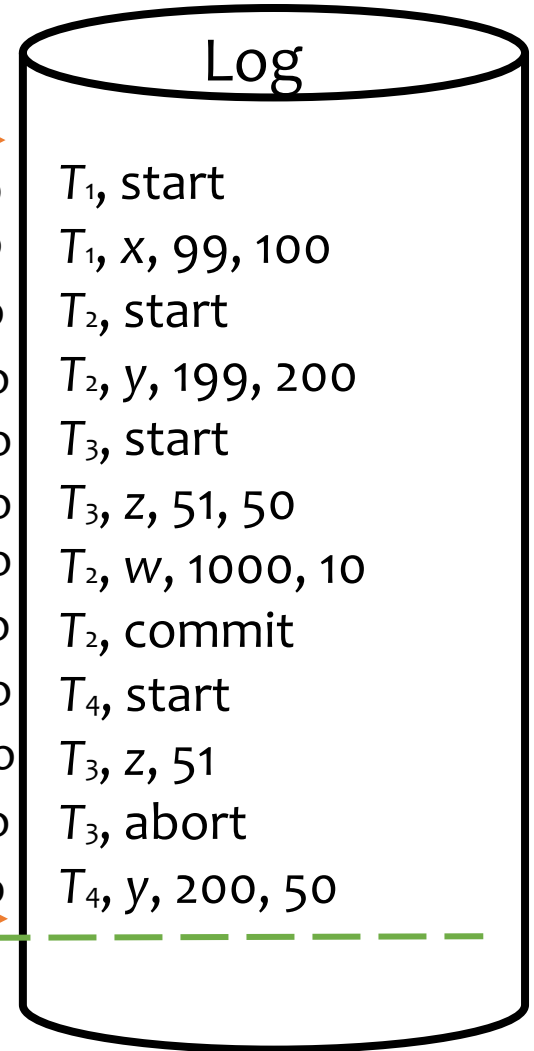


Start of log



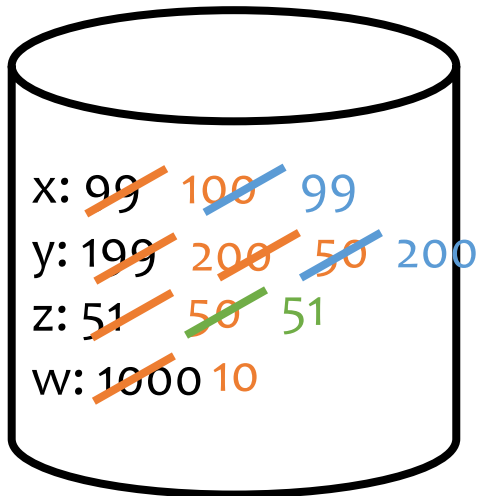
redo T<sub>1</sub>, start  
redo T<sub>1</sub>, x, 99, 100  
redo T<sub>2</sub>, start  
redo T<sub>2</sub>, y, 199, 200  
redo T<sub>3</sub>, start  
redo T<sub>3</sub>, z, 51, 50  
redo T<sub>2</sub>, w, 1000, 10  
redo T<sub>2</sub>, commit  
redo T<sub>4</sub>, start  
redo T<sub>3</sub>, z, 51  
redo T<sub>3</sub>, abort  
redo T<sub>4</sub>, y, 200, 50

End of log



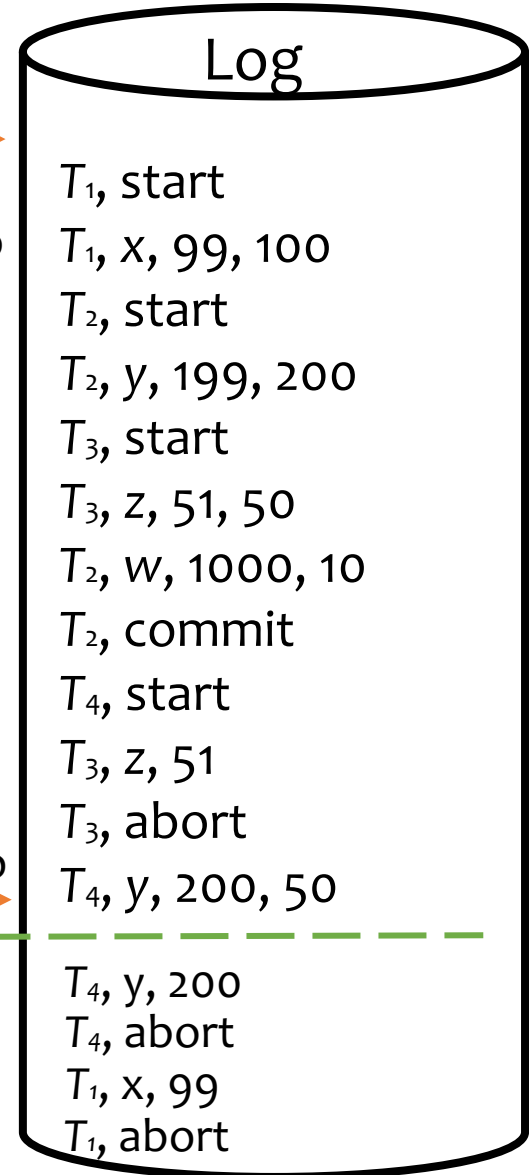
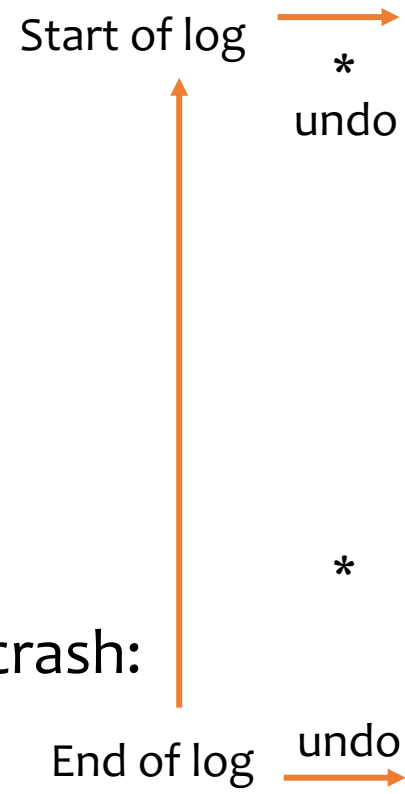
# Log example - Undo

- Undo phase: T1, T4



List of active transactions at crash:

T1 ~~T2~~ ~~T3~~ T4



# Undo/redo logging

- U: used to track the set of active transactions at crash
- Redo phase: scan **forward** to end of the log
  - For a log record  $\langle T, \text{start} \rangle$ , add  $T$  to  $U$
  - For a log record  $\langle T, X, \text{old}, \text{new} \rangle$ , issue  $\text{write}(X, \text{new})$
  - For a log record  $\langle T, \text{commit} \mid \text{abort} \rangle$ , remove  $T$  from  $U$ 
    - If *abort*, undo changes of  $T$  i.e., add  $\langle T, X, \text{old} \rangle$  before logging *abort*

👉 Basically repeats history!
- Undo phase: scan log **backward**
  - Undo the effects of transactions in  $U$
  - That is, for each log record  $\langle T, X, \text{old}, \text{new} \rangle$  where  $T$  is in  $U$ , issue  $\text{write}(X, \text{old})$ , and log this operation too, i.e., add  $\langle T, X, \text{old} \rangle$
  - Log  $\langle T, \text{abort} \rangle$  when all effects of  $T$  have been undone



# Checkpointing

- Shortens the amount of log that needs to be undone or redone when a failure occurs
- Assumption: Txns cannot perform any update actions, such as writing to a buffer block or writing a log record, while a checkpoint is in progress
- Steps:
  - Output to the disk all modified buffer blocks
  - Add to log: **<checkpoint L>**, where L is a list of txns active at the time of the checkpoint
- After a system crash has occurred, the system examines the log to find the last **<checkpoint L>** record
  - The redo operations will start from the checkpoint record
  - The undo operations will start from the end of the log until the list of active transactions is empty

# Summary

- Recovery: undo/redo logging
  - Normal operation: write-ahead logging, no force, steal
  - Recovery: first redo (forward), and then undo (backward)
  
- Next lecture:
  - Other forms of durability: data replication
  - Atomicity when data is stored on different machines
  - Data privacy