

Database recovery

CS348

Instructor: Sujaya Maiyya

Announcements

- Assignment 3 due Nov 24th
- Final demo for projects:
 - Please reach out to your TA by Nov 21st (tomorrow) to schedule your final demos for next week

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;

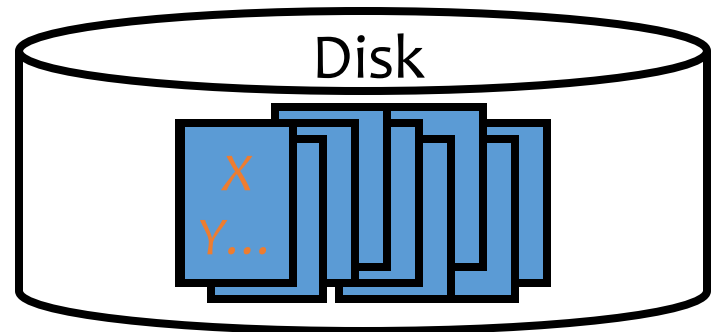
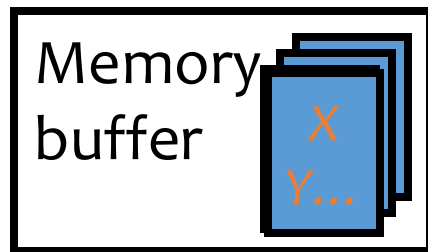
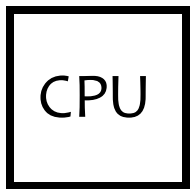
Outline

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

Execution model

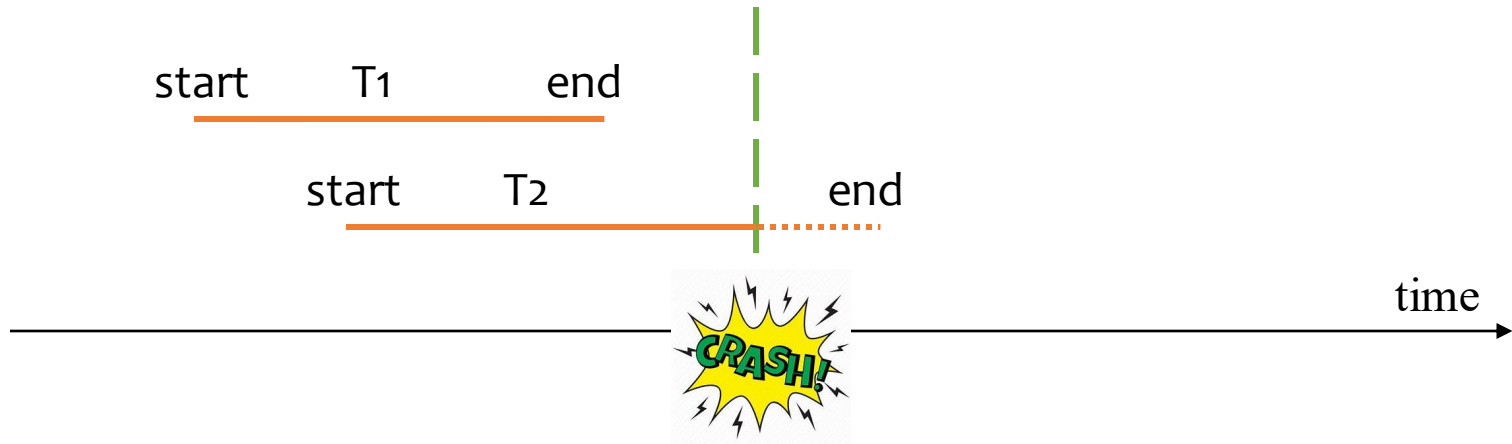
To read/write X

- 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: durability

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

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

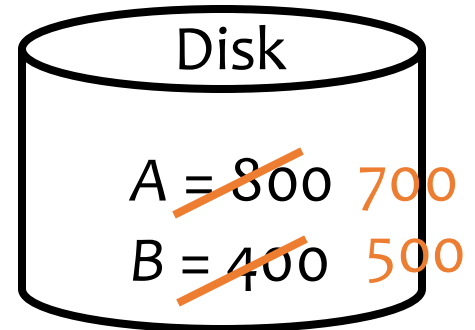
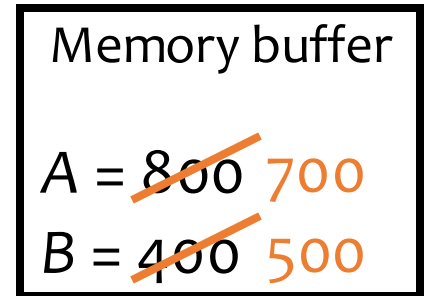
write(*A*, *a*);

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

write(*B*, *b*);

commit;

Force all writes to disk when a transaction commits



Naïve approach: durability

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

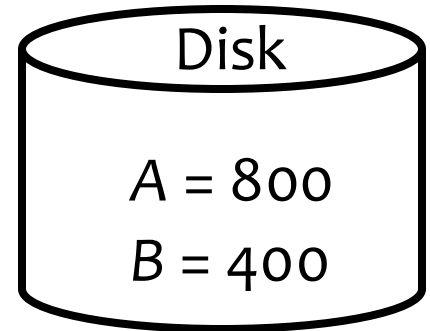
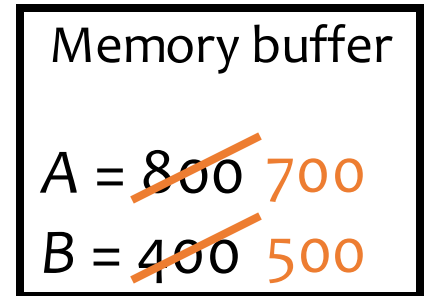
read(*A*, *a*); *a* = *a* - 100;

write(*A*, *a*);

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

write(*B*, *b*);

commit;



Force all writes to 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: 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;
```

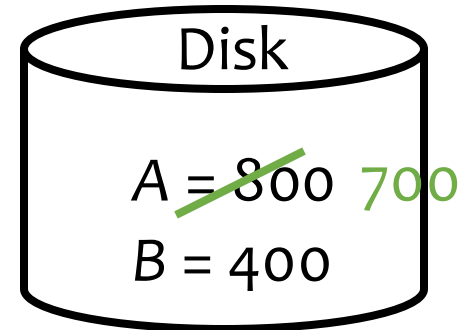
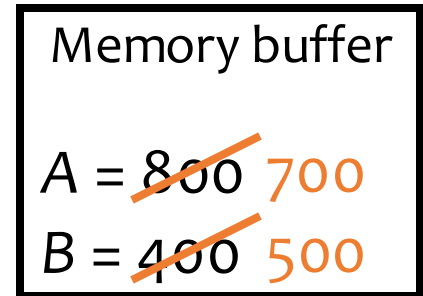
```
write(B, b);
```

```
commit;
```



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 this: some writes are on disk before *T* commits

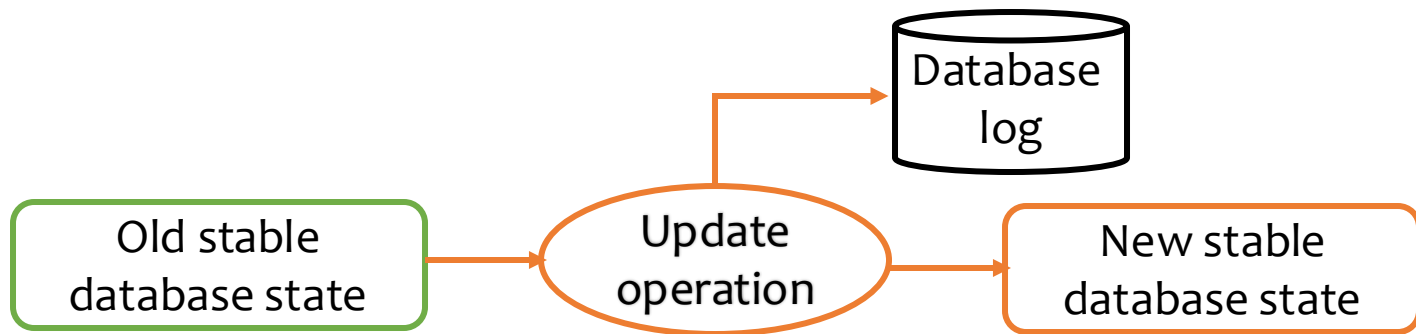
If system crashes before *T1* commits, there is no way to undo the changes

Naïve approach

- **Approach 1:** When a transaction commits, all writes of this transaction must be reflected on disk
 - Ensures durability
 - ☞ Problem: Lots of **random writes** hurt performance
- **Approach 2:** Writes of a transaction can only be flushed to disk at commit time
 - Ensures atomicity
 - ☞ Problem: Holding on to all dirty blocks **requires lots of memory**

Logging

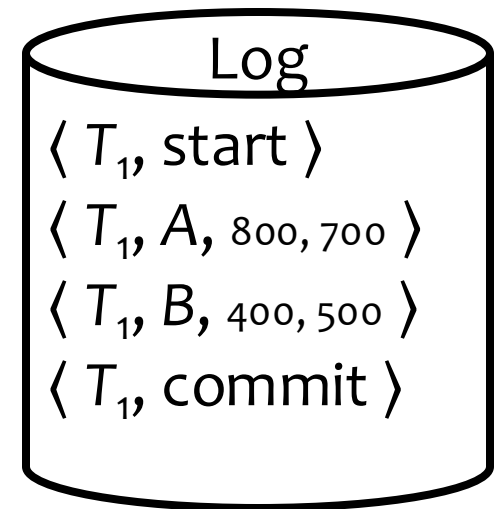
- **Database log**: sequence of **log records**, recording all changes made to the database, written to stable storage (e.g., disk) during normal operation



- 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*);

Memory buffer

A = ~~800~~ 700

B = ~~400~~ 500

Disk

A = 800

B = 400

Log

⟨ *T*₁, start ⟩

⟨ *T*₁, *A*, 800, 700 ⟩

⟨ *T*₁, *B*, 400, 500 ⟩

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

Undo/redo logging example cont.

*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*);

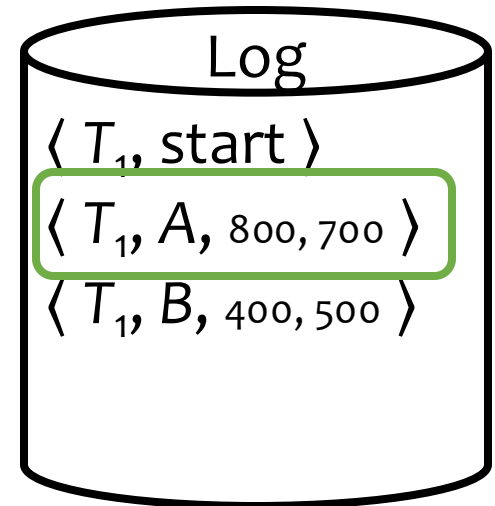
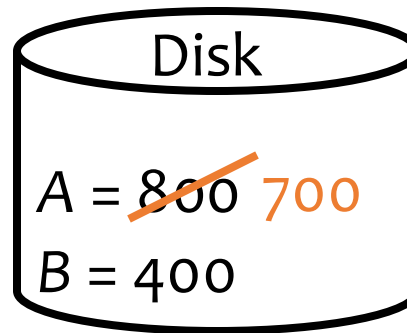


Memory buffer

~~*A* = 800~~ 700

~~*B* = 400~~ 500

Can flush
before commit



If system crashes before *T*₁ commits, we have the old value of *A* stored on the log to **undo** *T*₁

Undo/redo logging example cont.

*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);

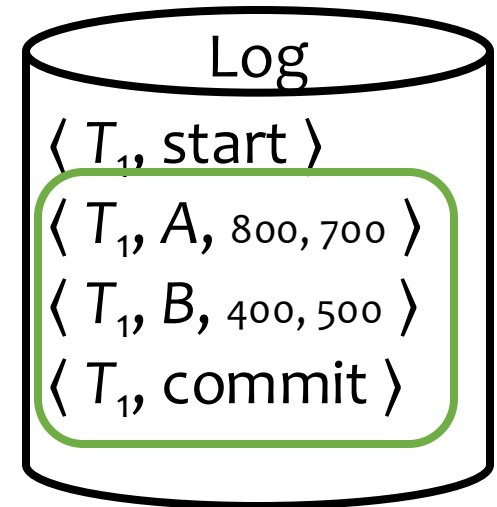
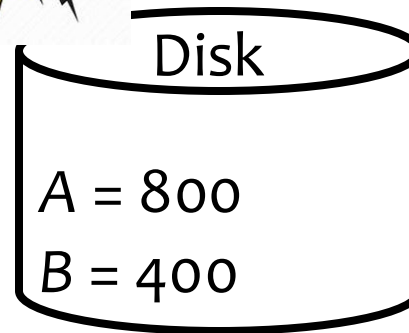
commit;

Memory buffer

A = ~~800~~ 700
B = ~~400~~ 500



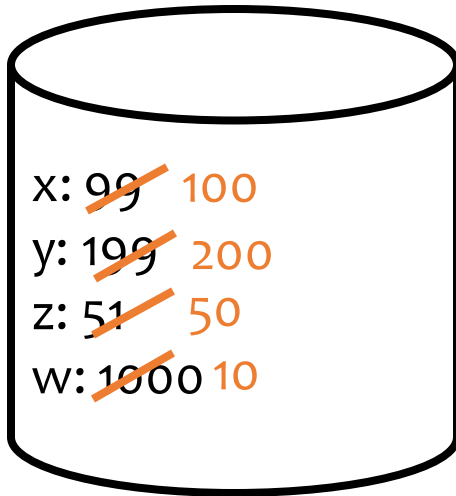
Can flush
after commit



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:

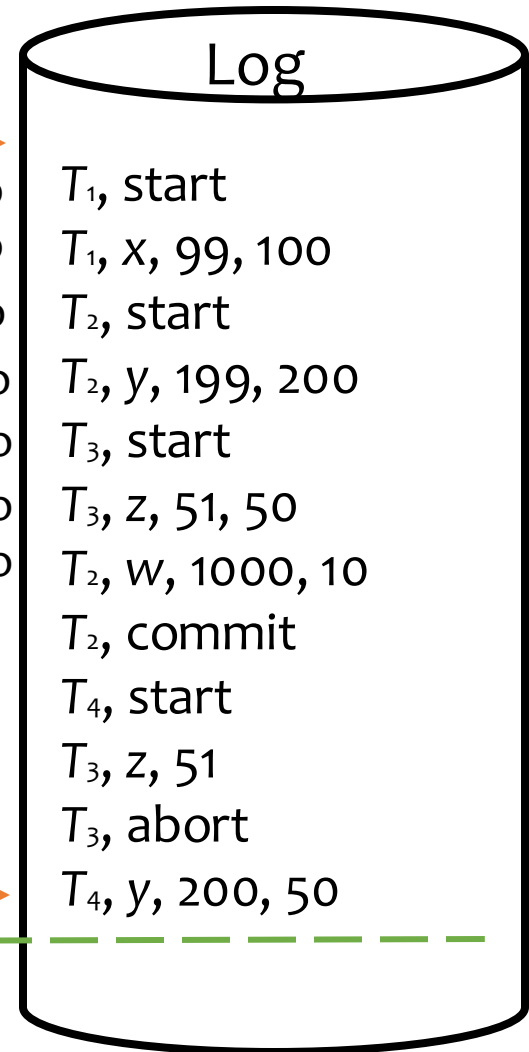
T1 T2 T3



Start of log

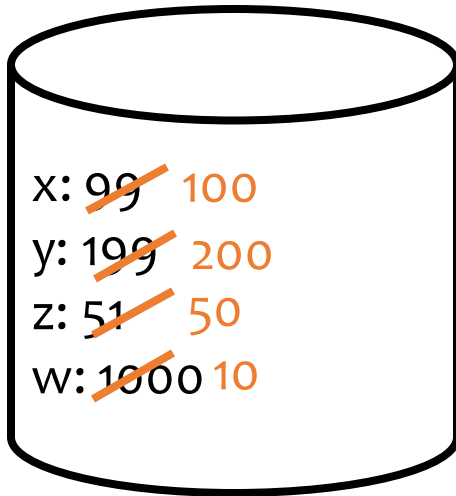
redo
redo
redo
redo
redo
redo
redo

End of log



Log example

- Redo phase:



List of active transactions at crash:

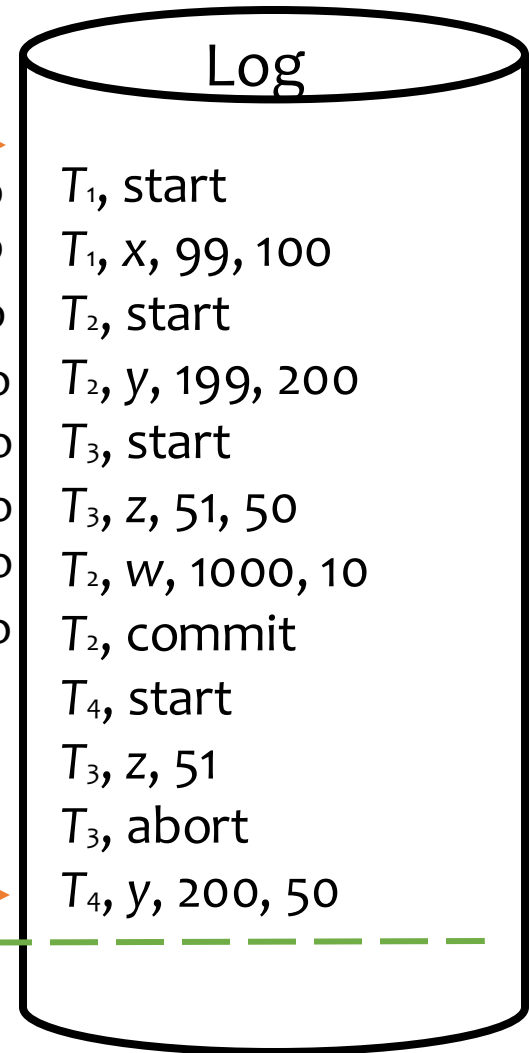
T1 ~~T2~~ T3



Start of log

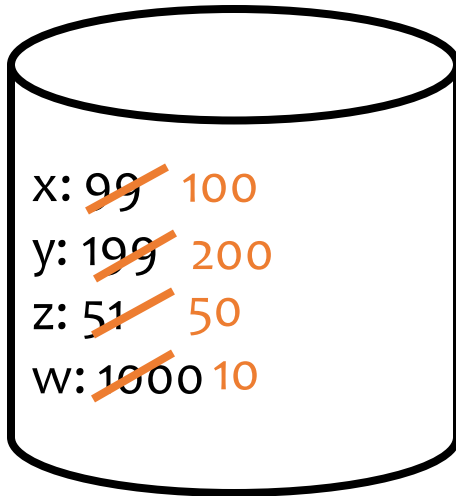
redo
redo
redo
redo
redo
redo
redo
redo

End of log



Log example

- Redo phase:



List of active transactions at crash:

T1 ~~T2~~ T3 T4

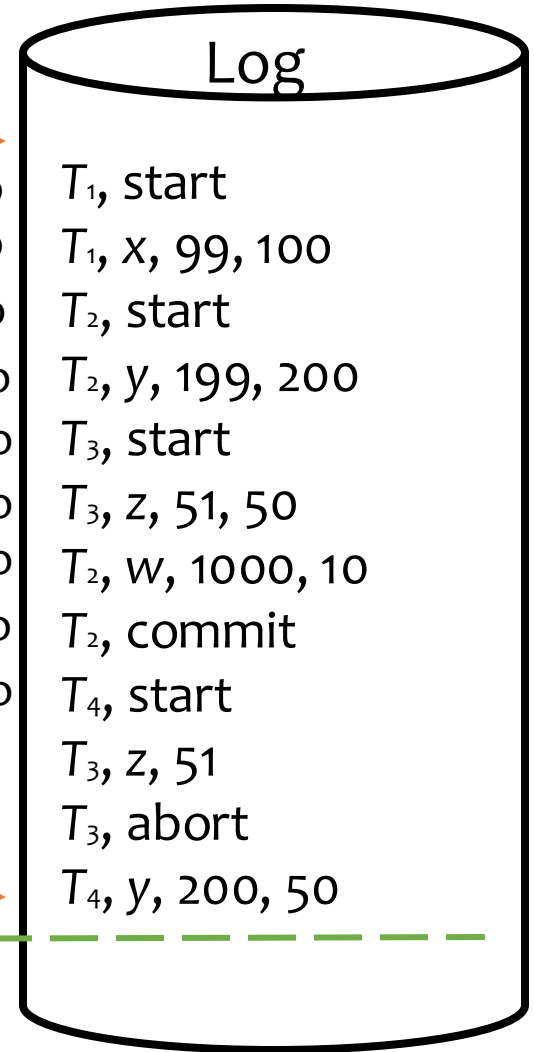


Start of log

redo T₁, start
redo T₁, x, 99, 100
redo T₂, start
redo T₂, y, 199, 200
redo T₃, start
redo T₃, z, 51, 50
redo T₂, w, 1000, 10
redo T₂, commit
redo T₄, start

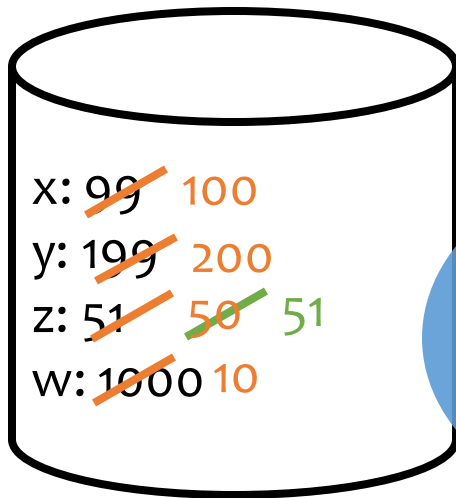
End of log

T₃, z, 51
T₃, abort
T₄, y, 200, 50



Log example

- Redo phase:



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

List of active transactions at crash:

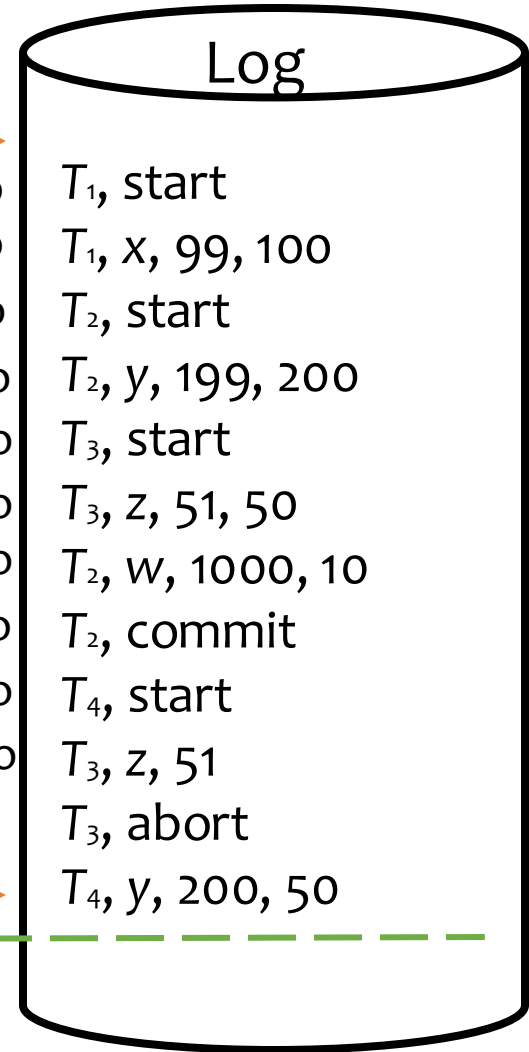
T1 ~~T2~~ T3 T4



Start of log

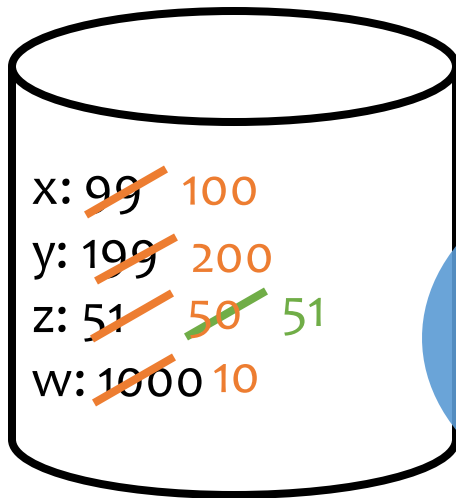
redo
redo
redo
redo
redo
redo
redo
redo
redo

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

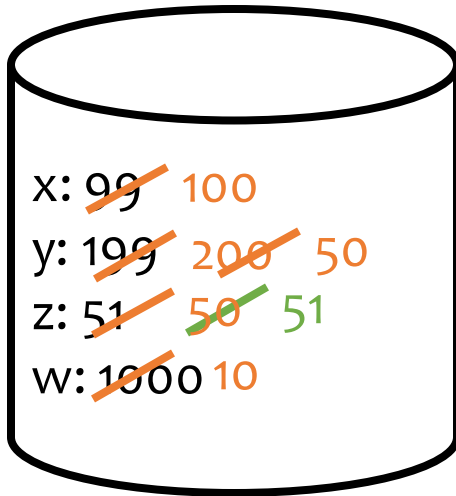
redo T₁, start
redo T₁, x, 99, 100
redo T₂, start
redo T₂, y, 199, 200
redo T₃, start
redo T₃, z, 51, 50
redo T₂, w, 1000, 10
redo T₂, commit
redo T₄, start
redo T₃, z, 51
redo T₃, abort
redo T₄, y, 200, 50

End of log

Log

Log example

- Redo phase:



List of active transactions at crash:

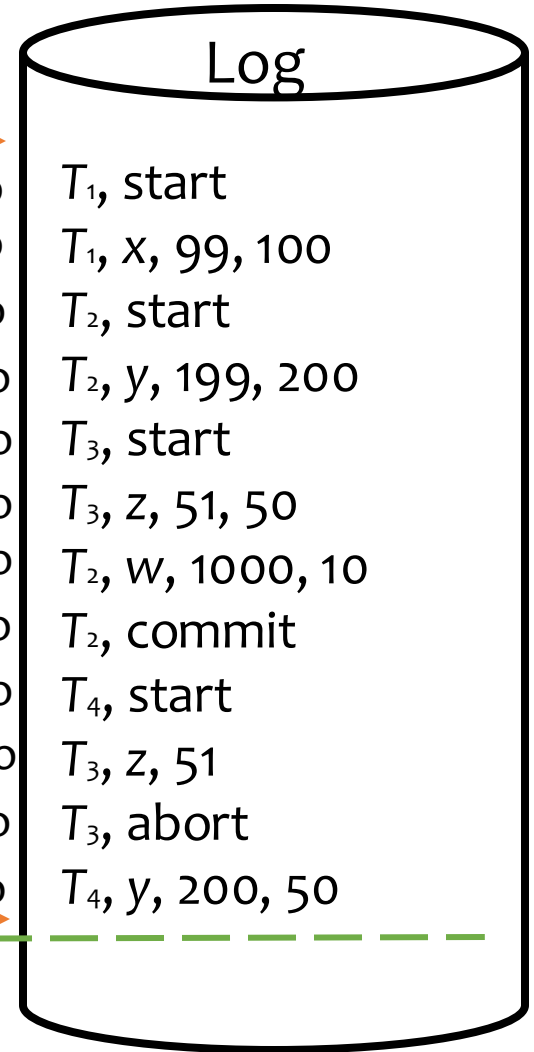
T1 ~~T2~~ ~~T3~~ T4



Start of log

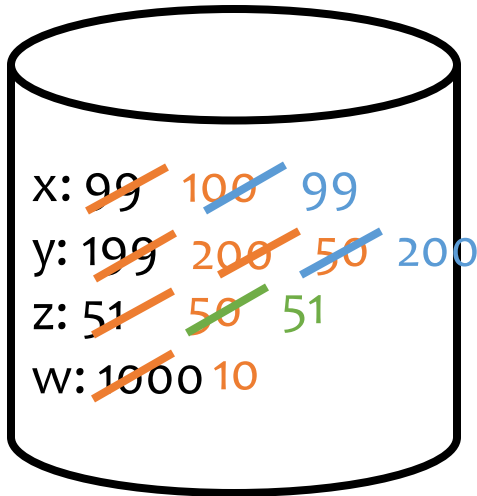


End of log



Log example - Undo

- Undo phase: T1, T4



List of active transactions at crash:

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



Start of log

undo

End of log

undo

Log

*
T₁, start
T₁, x, 99, 100
T₂, start
T₂, y, 199, 200
T₃, start
T₃, z, 51, 50
T₂, w, 1000, 10
T₂, commit
*
T₄, start
T₃, z, 51
T₃, abort
T₄, y, 200, 50

T₄, y, 200
T₄, abort
T₁, x, 99
T₁, abort

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 `write(X, new)`
 - For a log record $\langle T, \text{commit} \mid \text{abort} \rangle$, remove T from U
 - ☞ Basically repeats (some) 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 `write(X, 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
 - ☞ Basically reverts (some) history!

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
 - Recovery: first redo (forward), and then undo (backward)
- No lecture next week – schedule your project presentations
- Finals review Dec 2nd
- Final exam: Dec 7th 7:30PM-10PM 😞
 - Cheat sheet with I/O costs and cardinality estimates will be provided

Student Course Perceptions Surveys

Your chance to share your learning experience.
Your feedback is important!

Who has access to SCP results?

- Written comments: only the course instructor
- Numerical ratings: course instructor and academic leaders

How are SCP results used?

- Help instructors improve teaching and courses
- Inform pay and tenure decisions
- Contribute to decisions about program improvement and future teaching assignments

<https://perceptions.uwaterloo.ca/>

Difficulties?
Contact kabecker@uwaterloo.ca



Quick



Anonymous



Impactful