Recovery

CSC443H1 Database System Technology

Last stretch for the project!



Last lecture today



Last lecture today



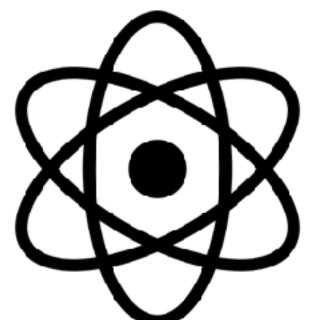
One hour lecture followed by tutorial questions

Please do the course evaluation!



ACID Transactions

Atomicity



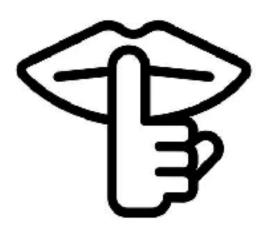
All or nothing

Consistency



Transition across consistent states

Isolation



Not corrupted by concurrency

Durability



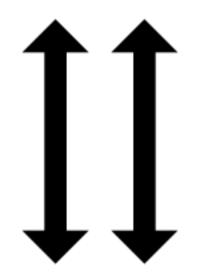
Recover from failure

What Endangers the ACID properties?

System Failure

(4)

Concurrency



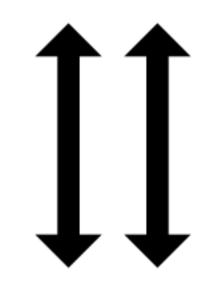
What Endangers the ACID properties?

System Failure

(4)

Today

Concurrency



Last week

Media failure

1TB

Power failure



Data center failure



Media failure

Power failure

Data center failure



RAID (covered)





Media failure

Power failure

Data center failure



RAID (covered)



Logging



Media failure

1TB

RAID (covered)

Power failure



Logging

Data center failure



Replication

Example:
Bank transfer

A = A - 100

B = B + 100



A = A - 100

Example:
Bank transfer

Power fails

B = B + 100



If the update to account A is saved to disk, money disappears from the system

A = A - 100

Example:

Bank transfer

Power fails

(4)

B = B + 100

If the update to account A is saved to disk, money disappears from the system

A = A - 100

Power fails

Example:

Bank transfer

B = B + 100

What's a high level solution?

Solution outline: (1) wrap this up as a transaction

Start transaction

$$A = A - 100$$

$$B = B + 100$$

End transaction

Solution outline: (1) wrap this up as a transaction

(2) write changes to by each transaction in a log



Start transaction

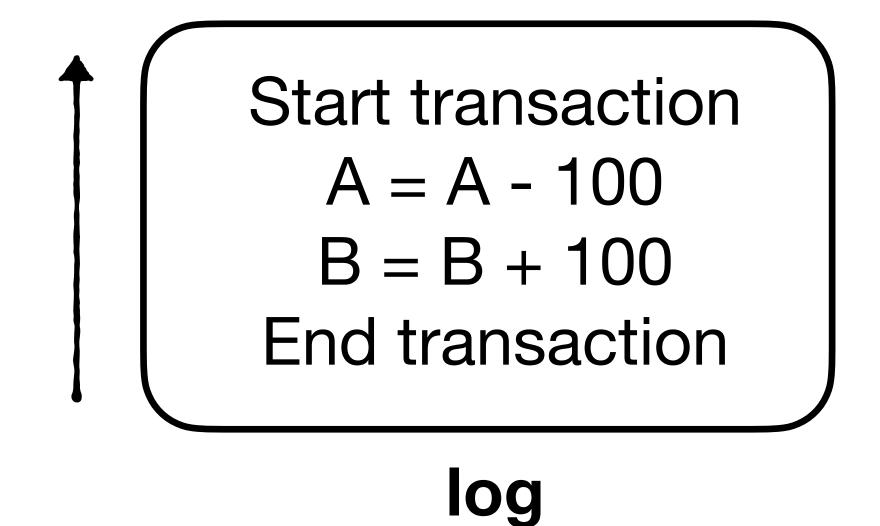
$$A = A - 100$$

$$B = B + 100$$

End transaction

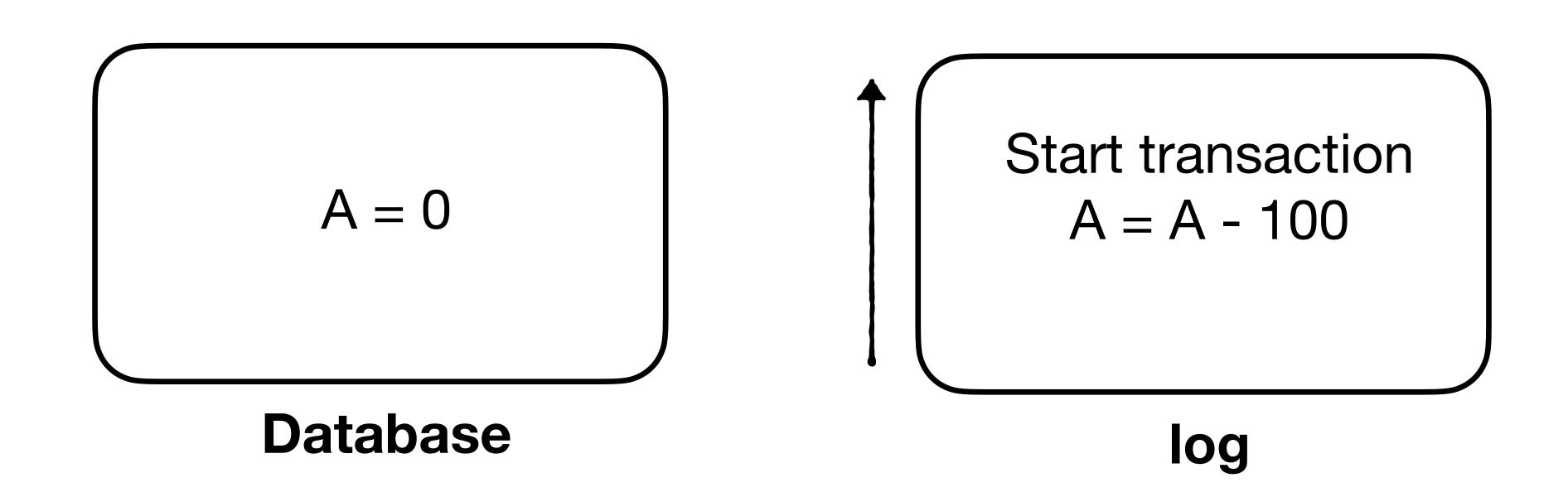
log: an append-only file in storage

- Solution outline: (1) wrap this up as a transaction
 - (2) write changes to by each transaction in a log
 - (3) after power fails, scan log and cancel effects of uncommitted transactions

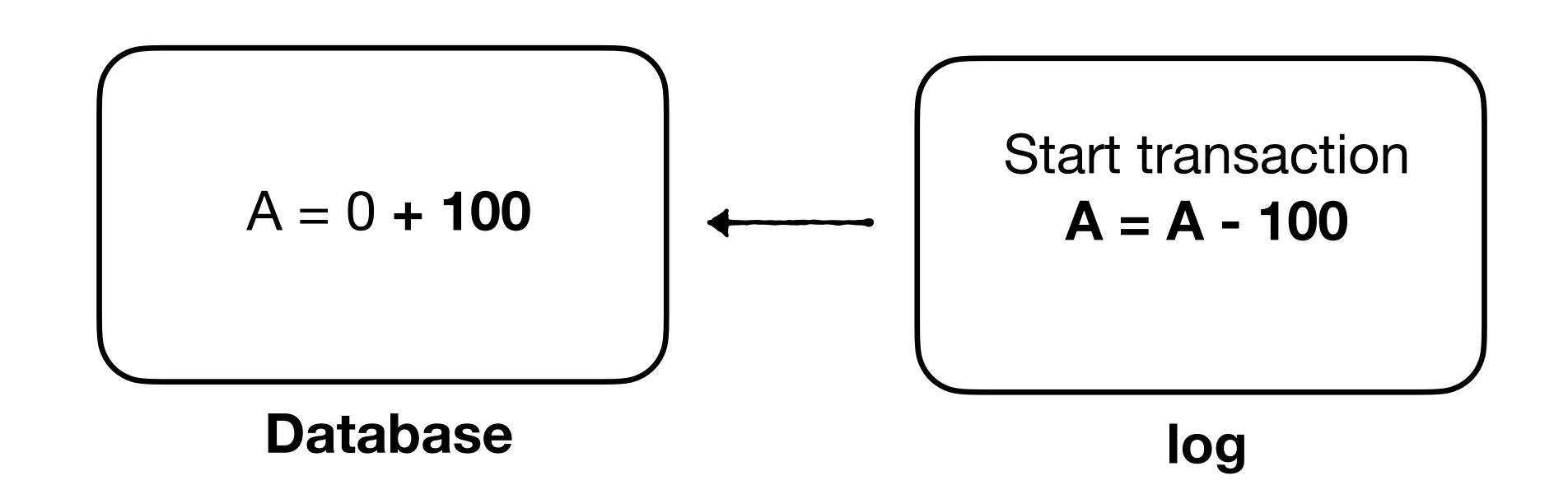


Solution outline: (1) wrap this up as a transaction

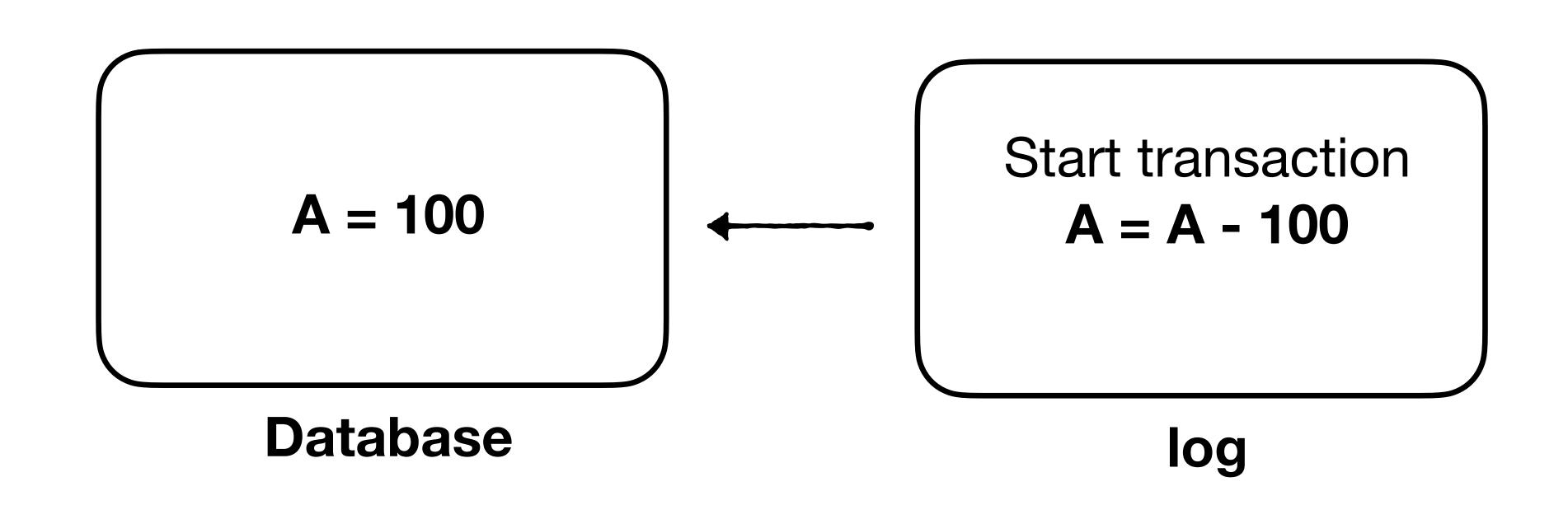
- (2) write changes to by each transaction in a log
- (3) after power fails, scan log and cancel effects of uncommitted transactions



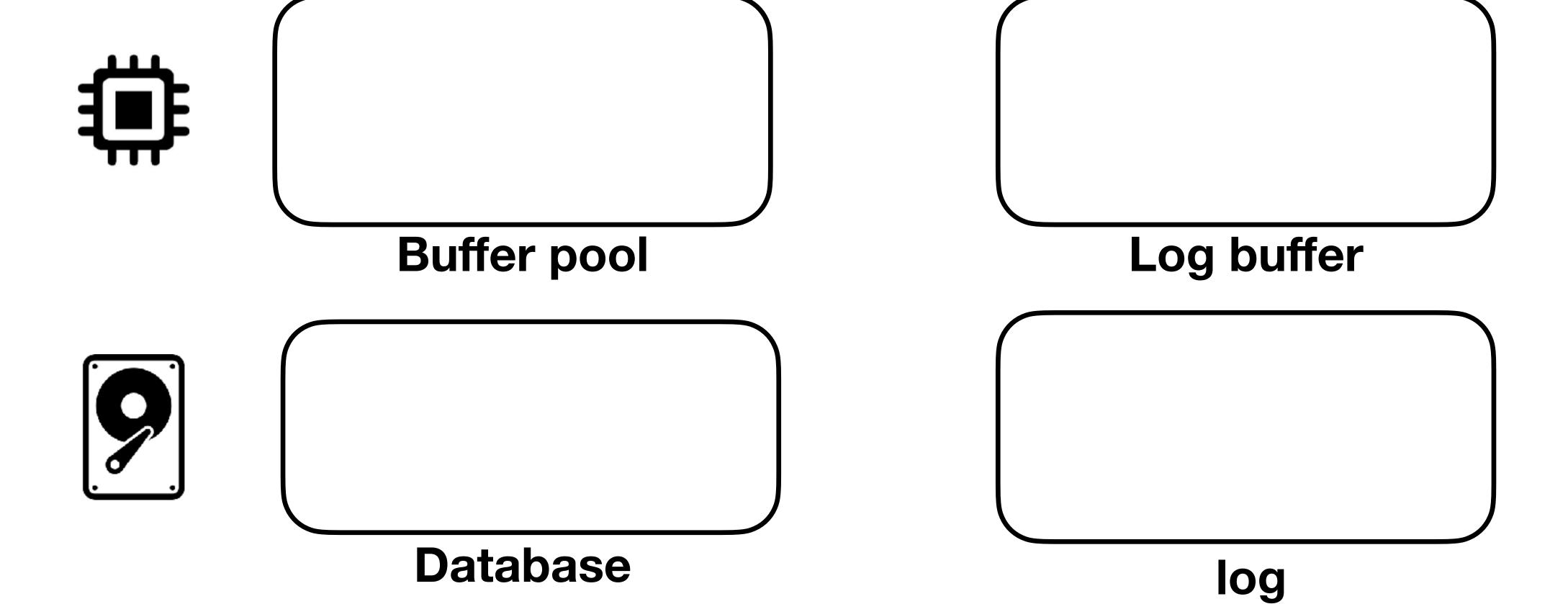
- Solution outline: (1) wrap this up as a transaction
 - (2) write changes to by each transaction in a log
 - (3) after power fails, scan log and cancel effects of uncommitted transactions



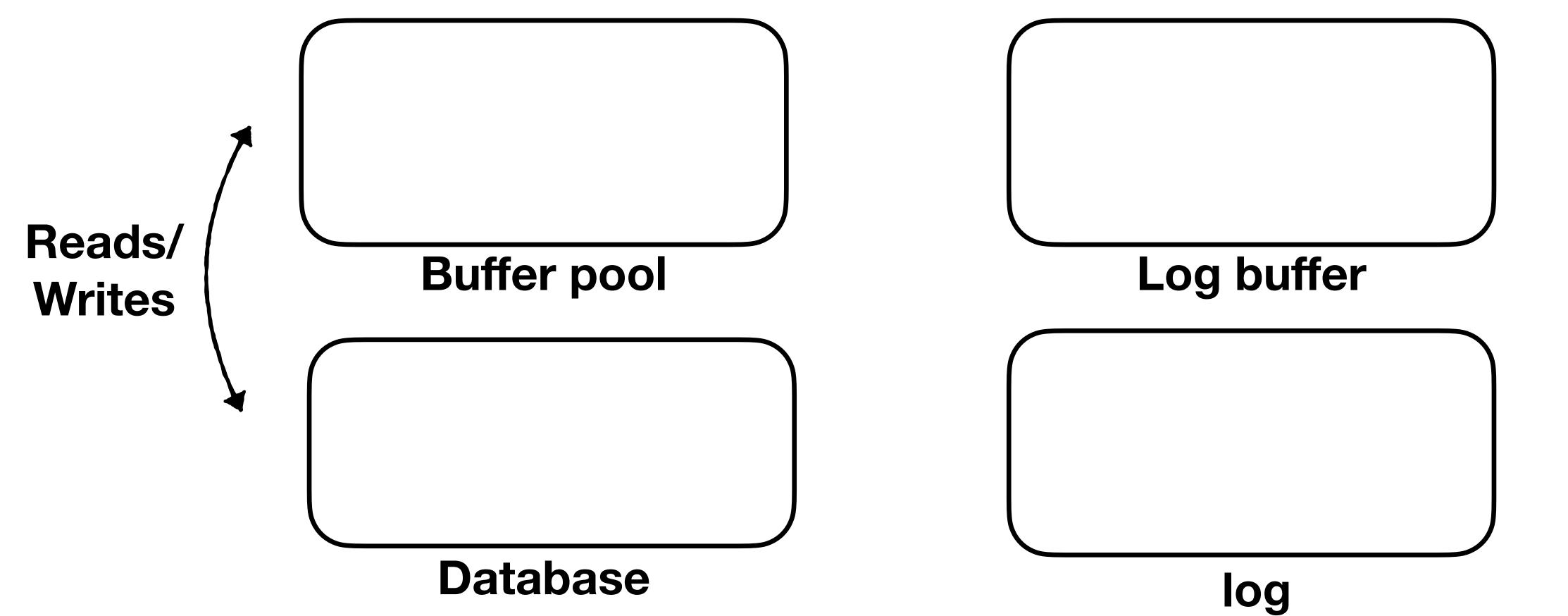
- Solution outline: (1) wrap this up as a transaction
 - (2) write changes to by each transaction in a log
 - (3) after power fails, scan log and cancel effects of uncommitted transactions



Three Types of logging with different trade-offs

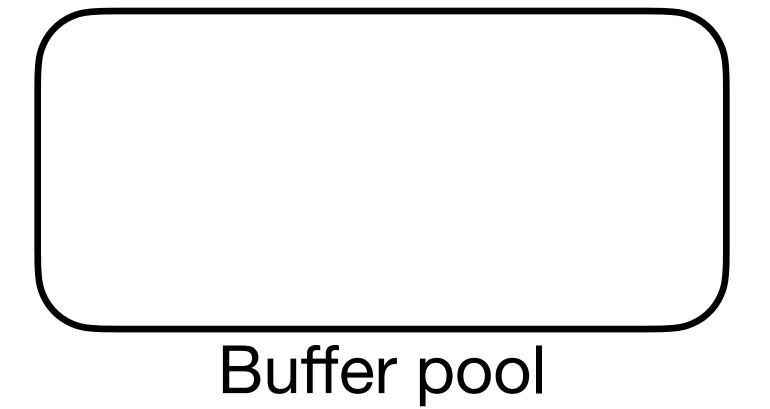


Assume relational database



$$A = A - 100$$

$$B = B + 100$$



$$A = 100$$

$$B = 100$$

Database



Log buffer

$$A = A - 100$$

$$B = B + 100$$

$$A = 100$$

$$B = 100$$

Buffer pool

$$A = 100$$

$$B = 100$$

Database

<T1, start>

Log buffer

$$A = A - 100$$

$$B = B + 100$$

Write older values to log

$$A = 0$$

$$B = 200$$

Buffer pool

$$A = 100$$

$$B = 100$$

Database

Log buffer

$$A = A - 100$$

$$B = B + 100$$

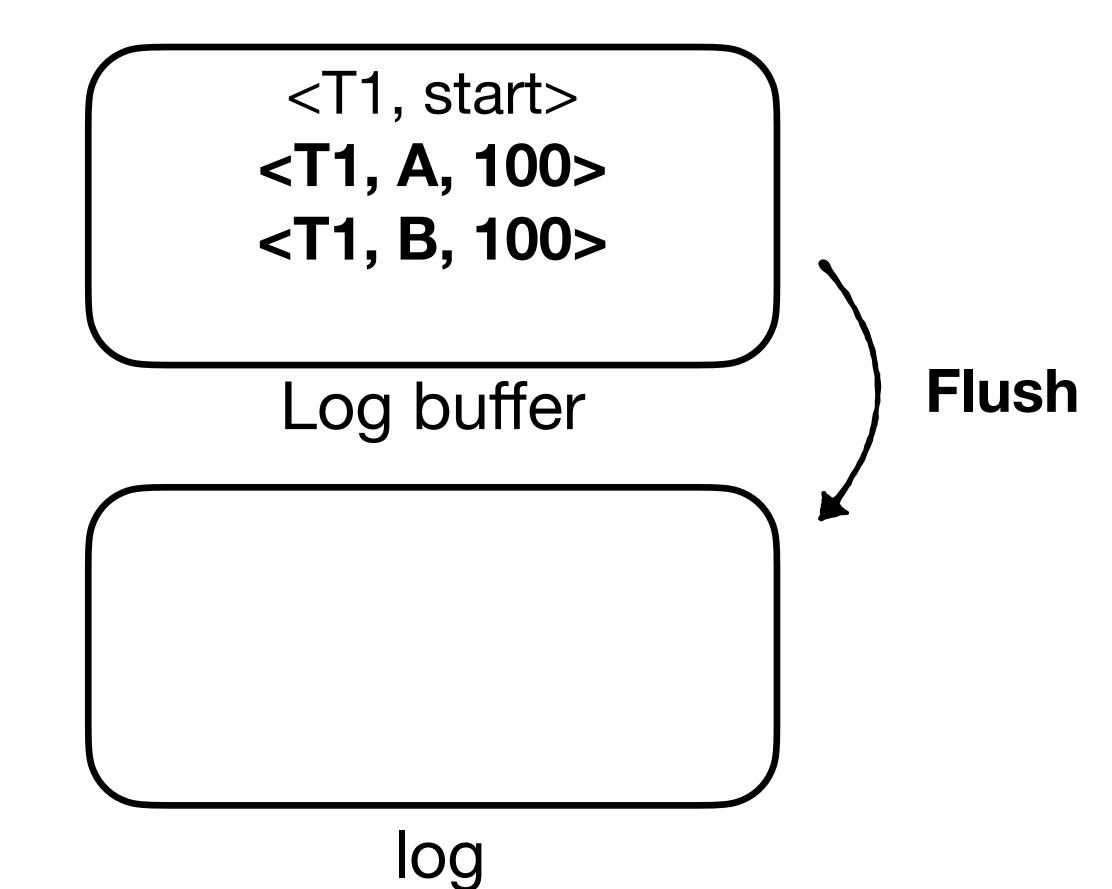
We now want to commit. First, flush the log.

$$A = 0$$

$$B = 200$$
Buffer pool

$$A = 100$$
 $B = 100$

Database



Transaction T1

$$A = A - 100$$

$$B = B + 100$$



A = 0

$$B = 200$$

Buffer pool

$$A = 100$$

$$B = 100$$

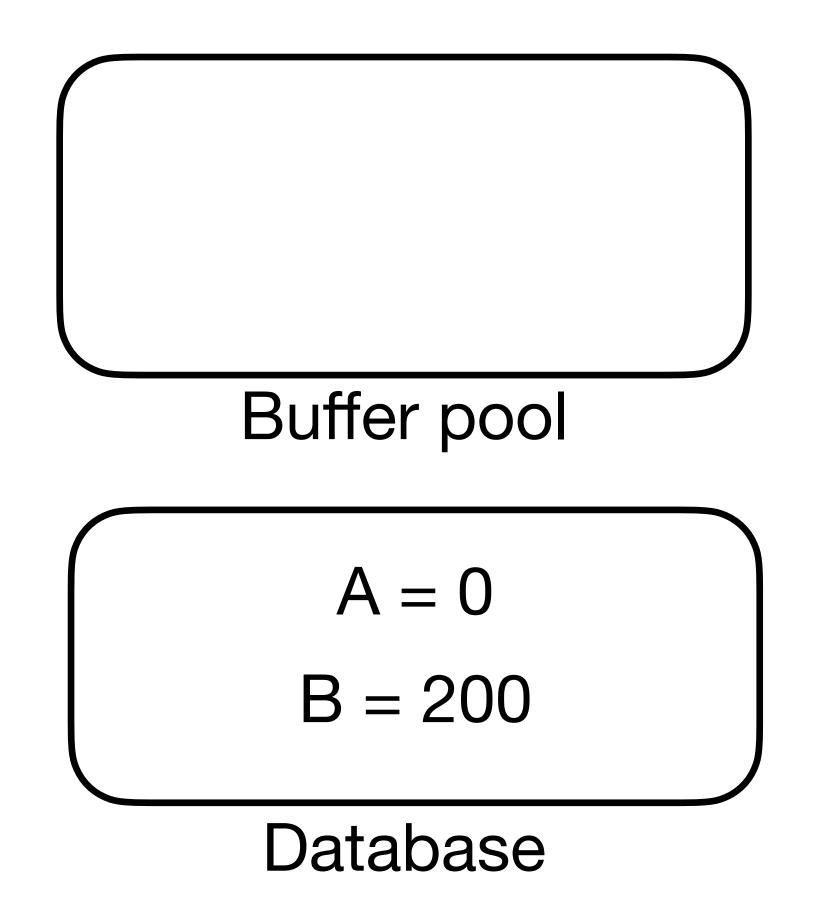
Database

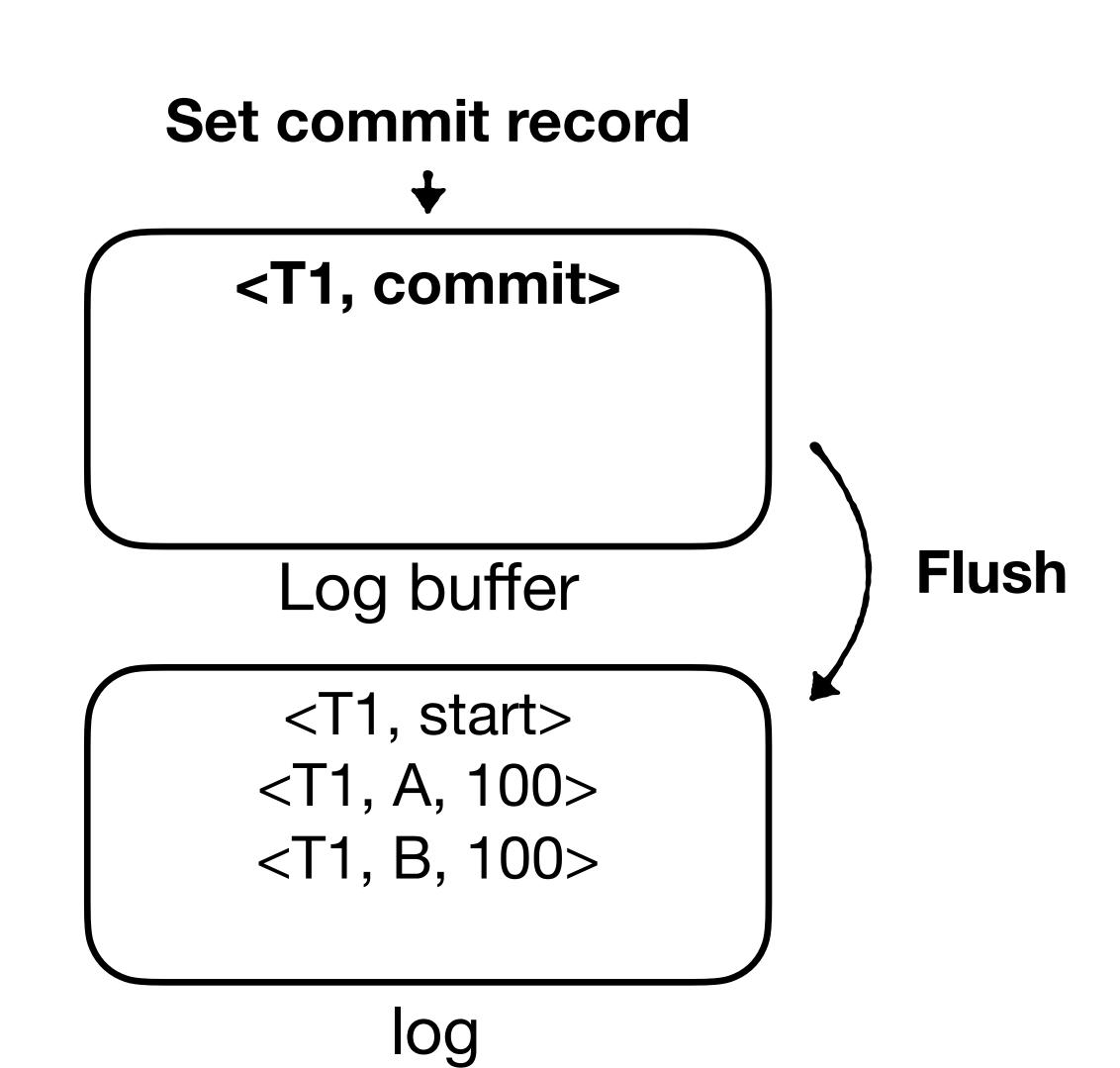
Log buffer

<T1, start>
<T1, A, 100>
<T1, B, 100>

$$A = A - 100$$

$$B = B + 100$$





$$A = A - 100$$

$$B = B + 100$$



$$A = 0$$

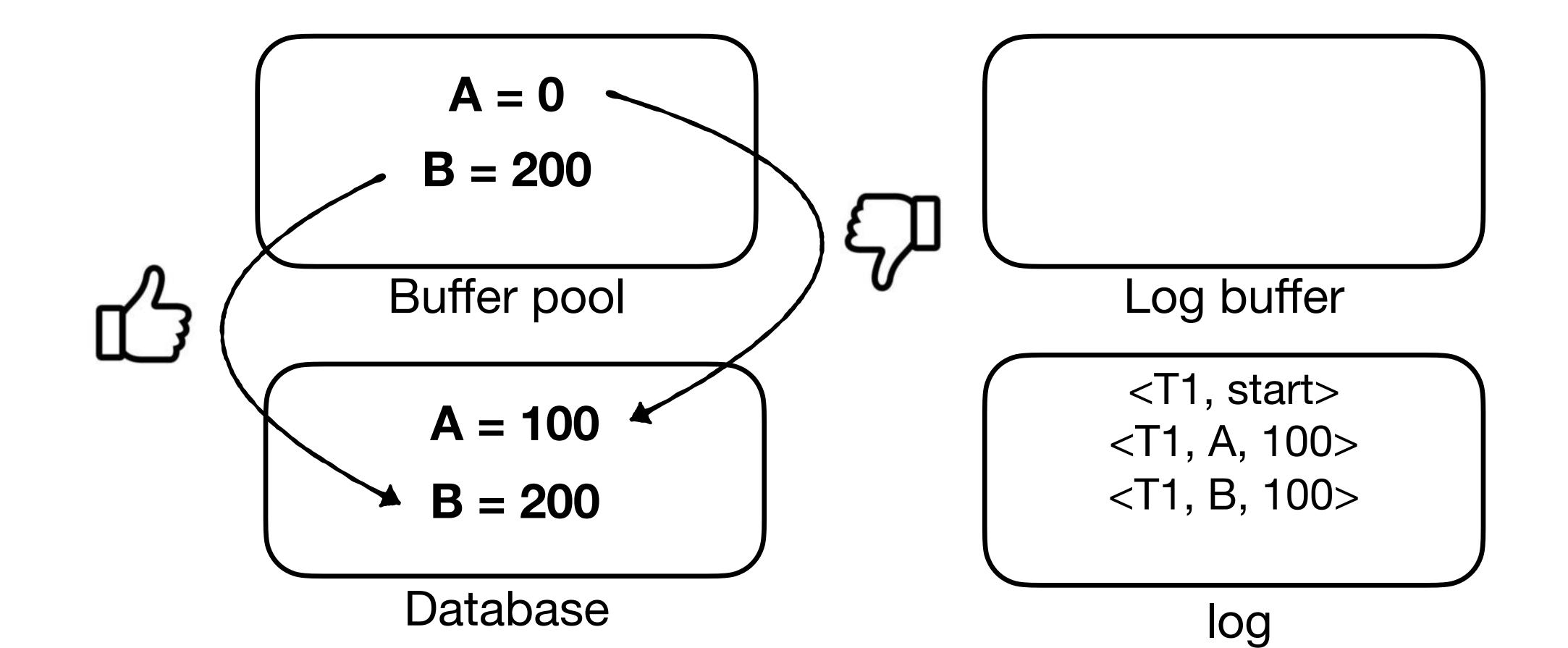
$$B = 200$$
Database

- (1) write before-image for any changed value to log buffer
- (2) flush log
- (3) force the changed data into storage
- (4) add commit record to log buffer & flush

Recovery Undo logging:

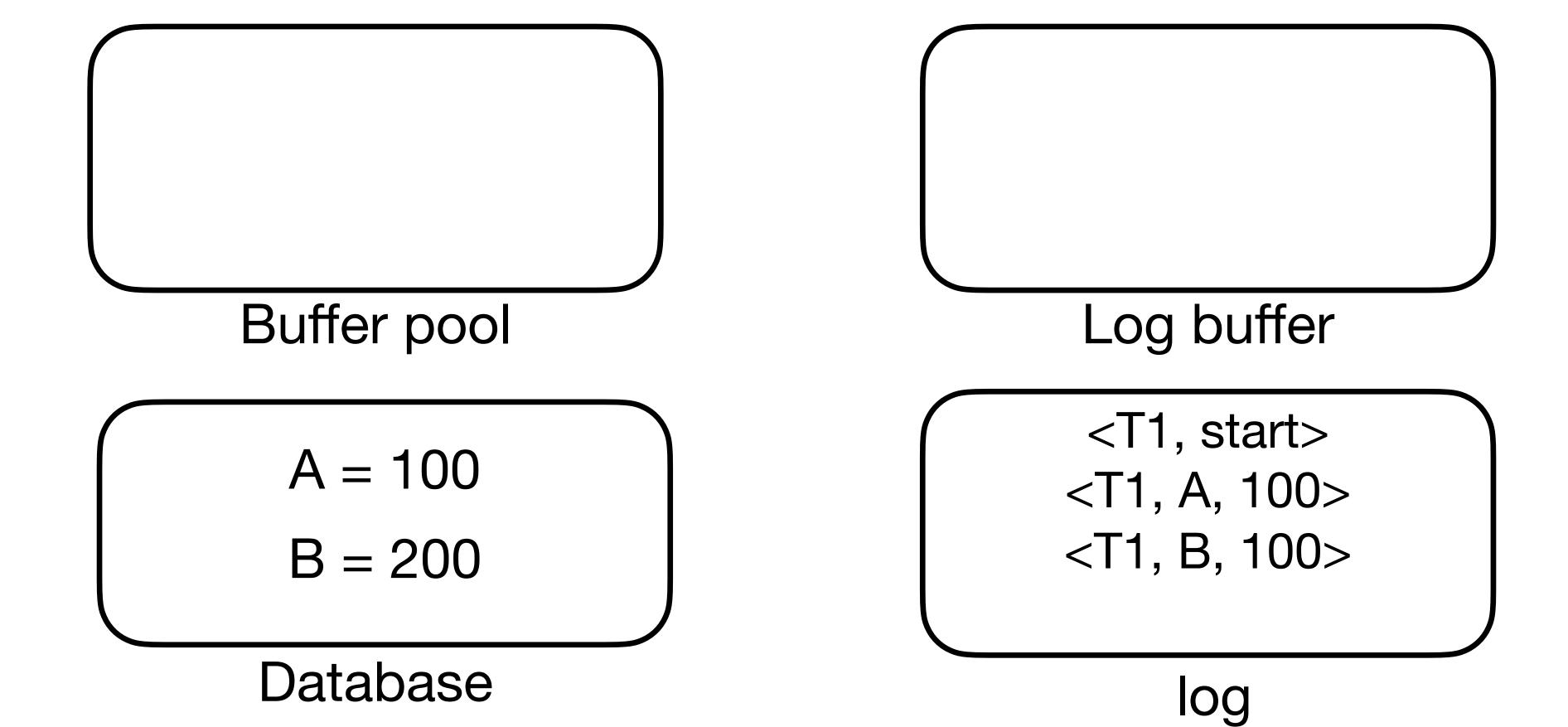


Suppose power fails after we save the value of B to storage but before we save the value of A



Recovery Undo logging:

We lose all contents in memory

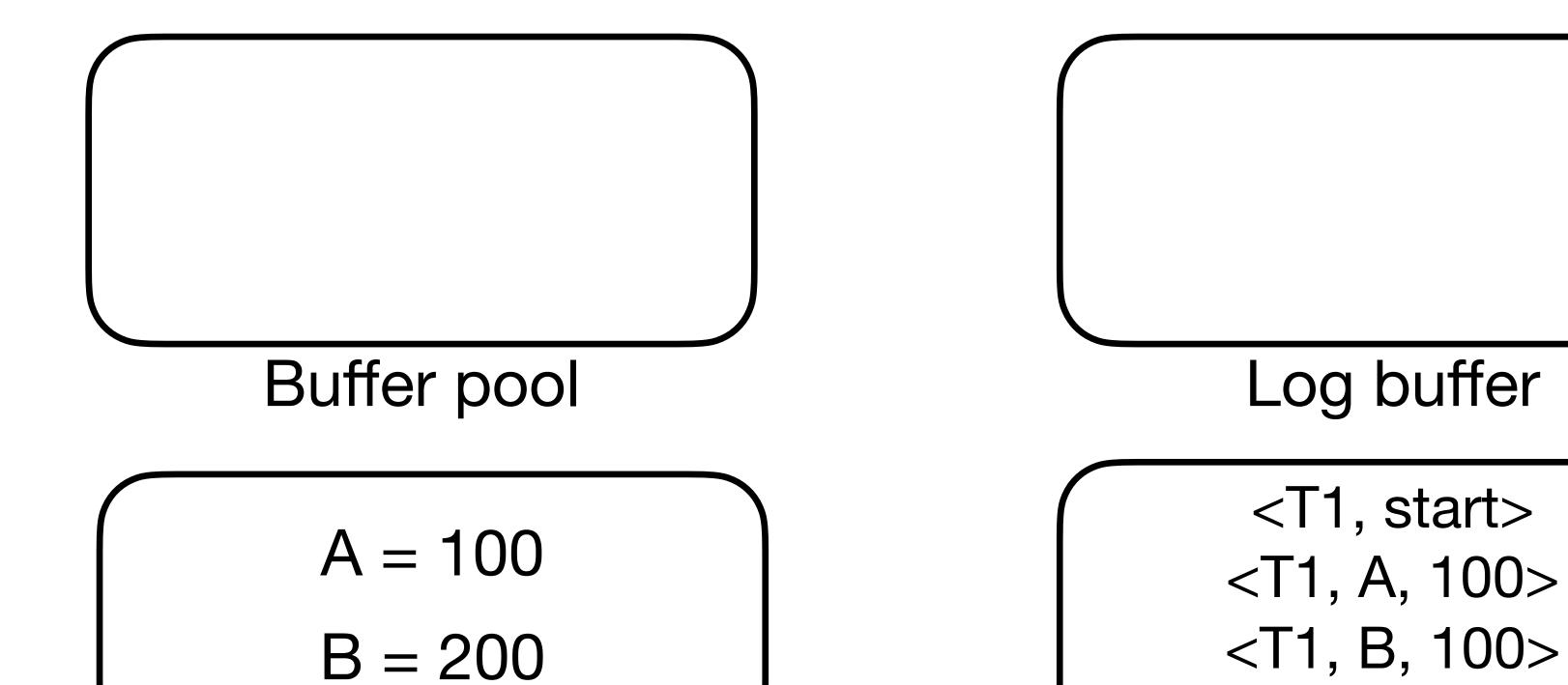


Recovery Undo logging:

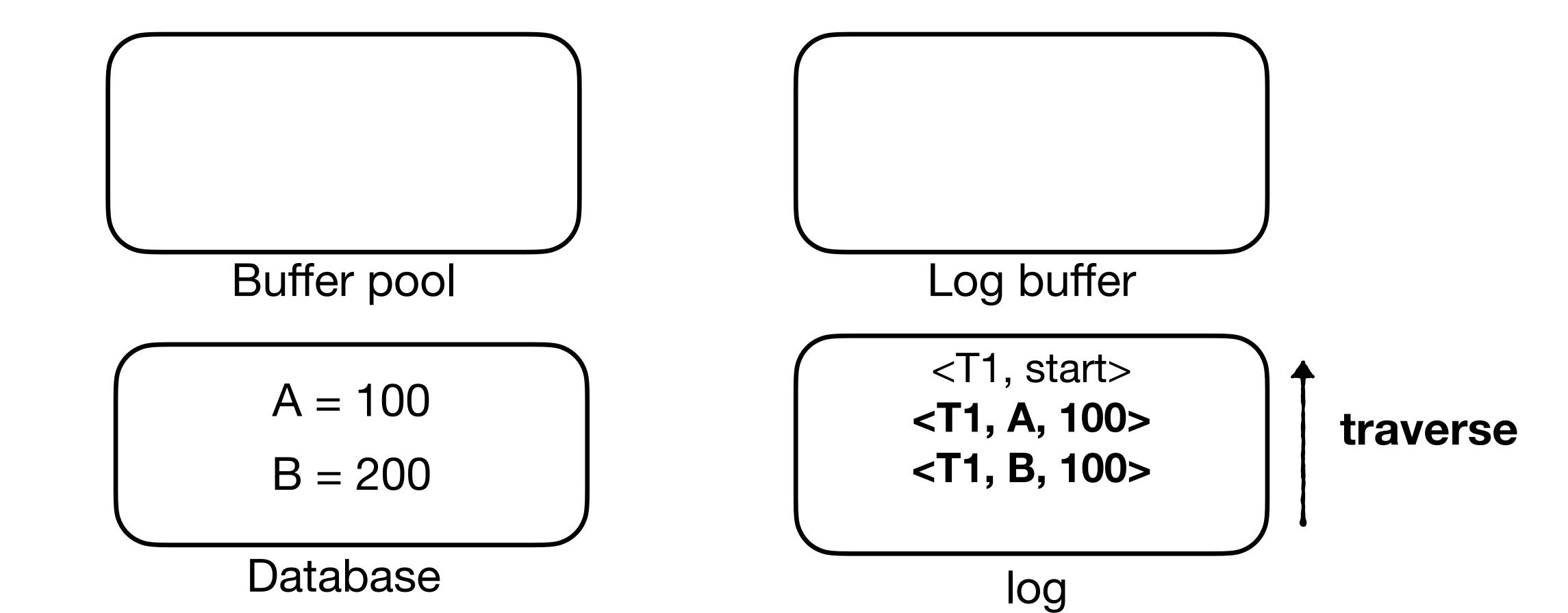
Database

We lose all contents in memory

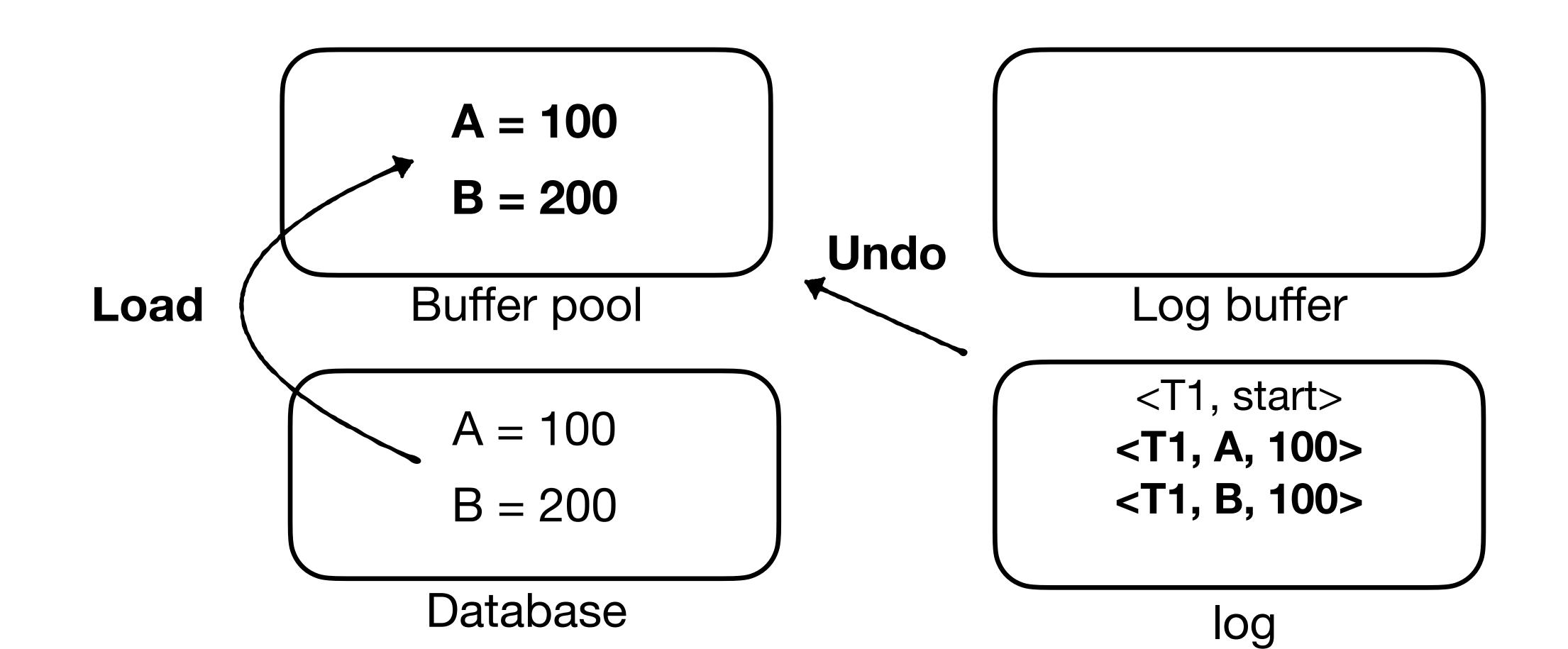
How do we recover?



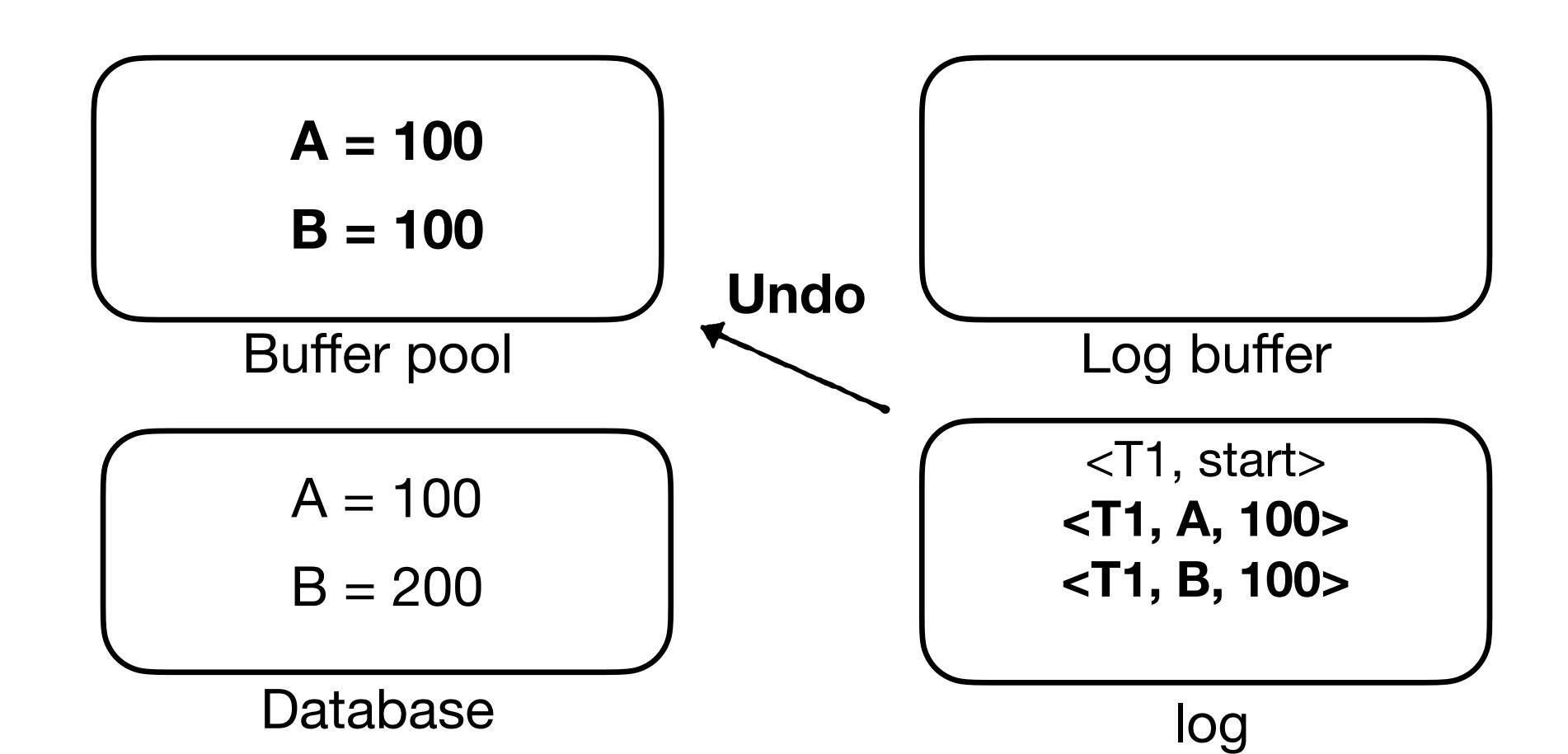
During recovery, we traverse the log backwards



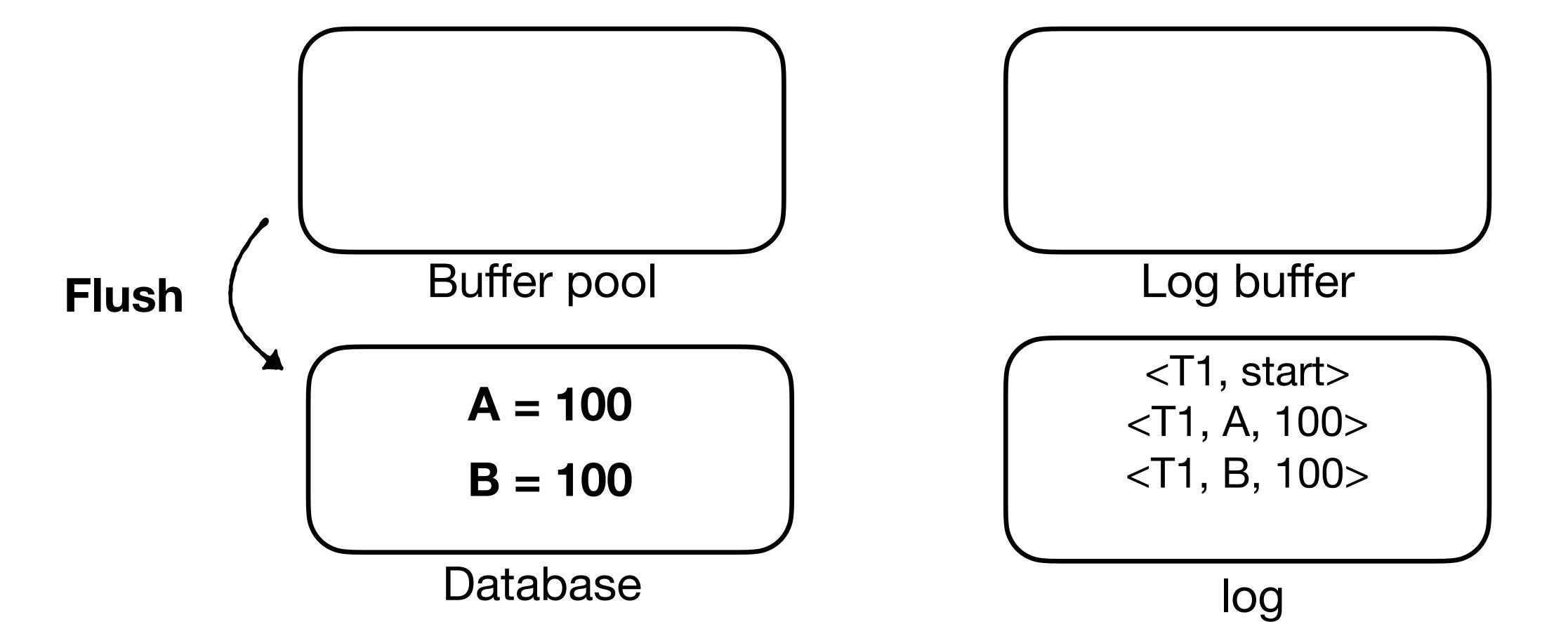
Undo any operation of an uncommitted transaction

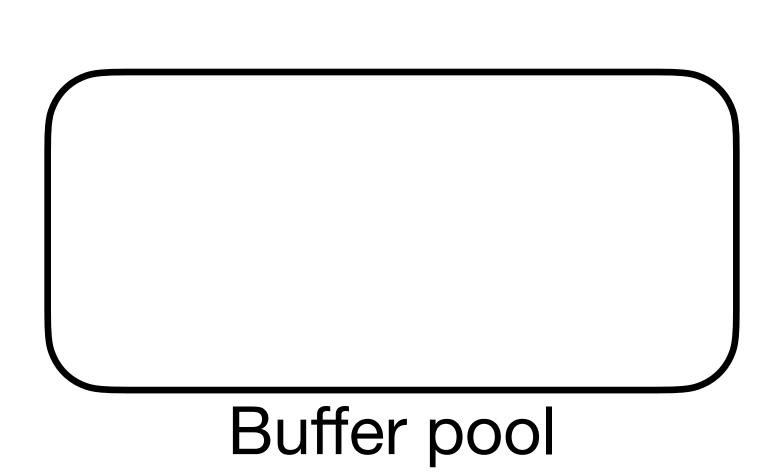


Undo any operation of an uncommitted transaction



Undo any operation of an uncommitted transaction





$$A = 100$$

$$B = 100$$
Database

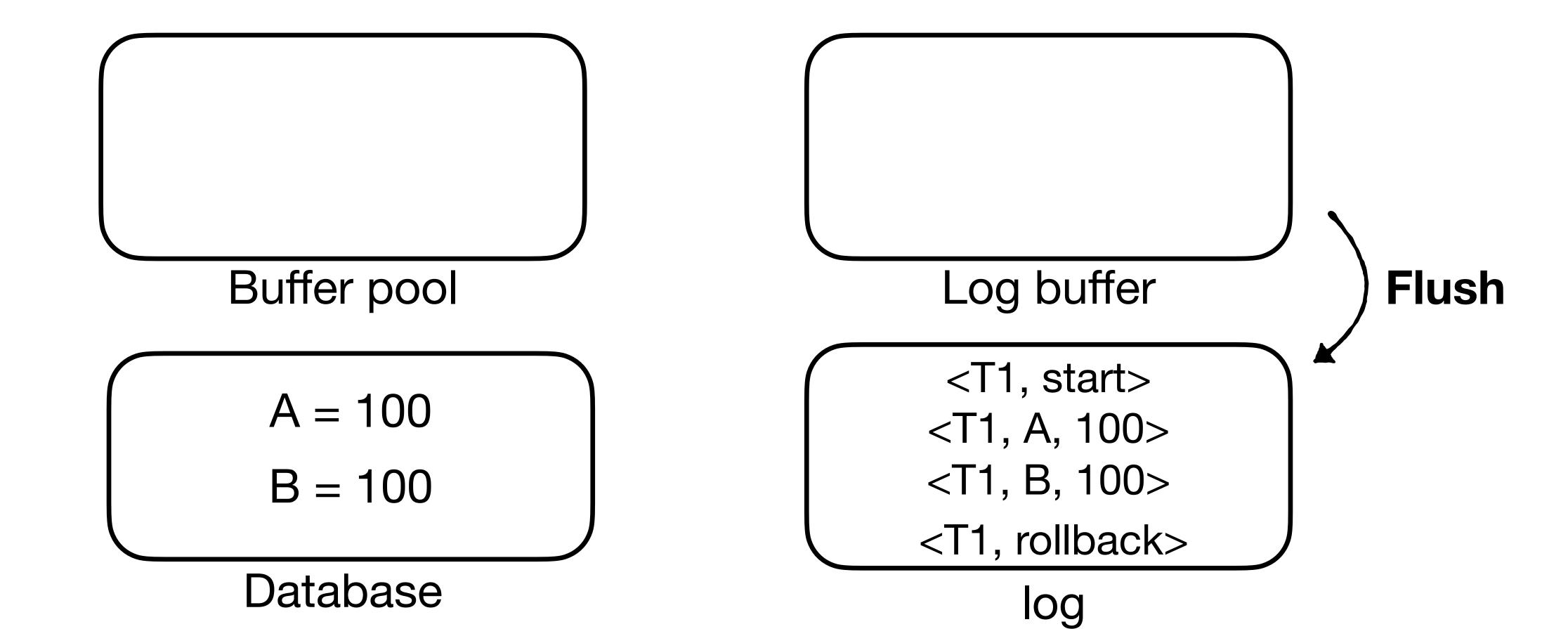
Set rollback record



Log buffer

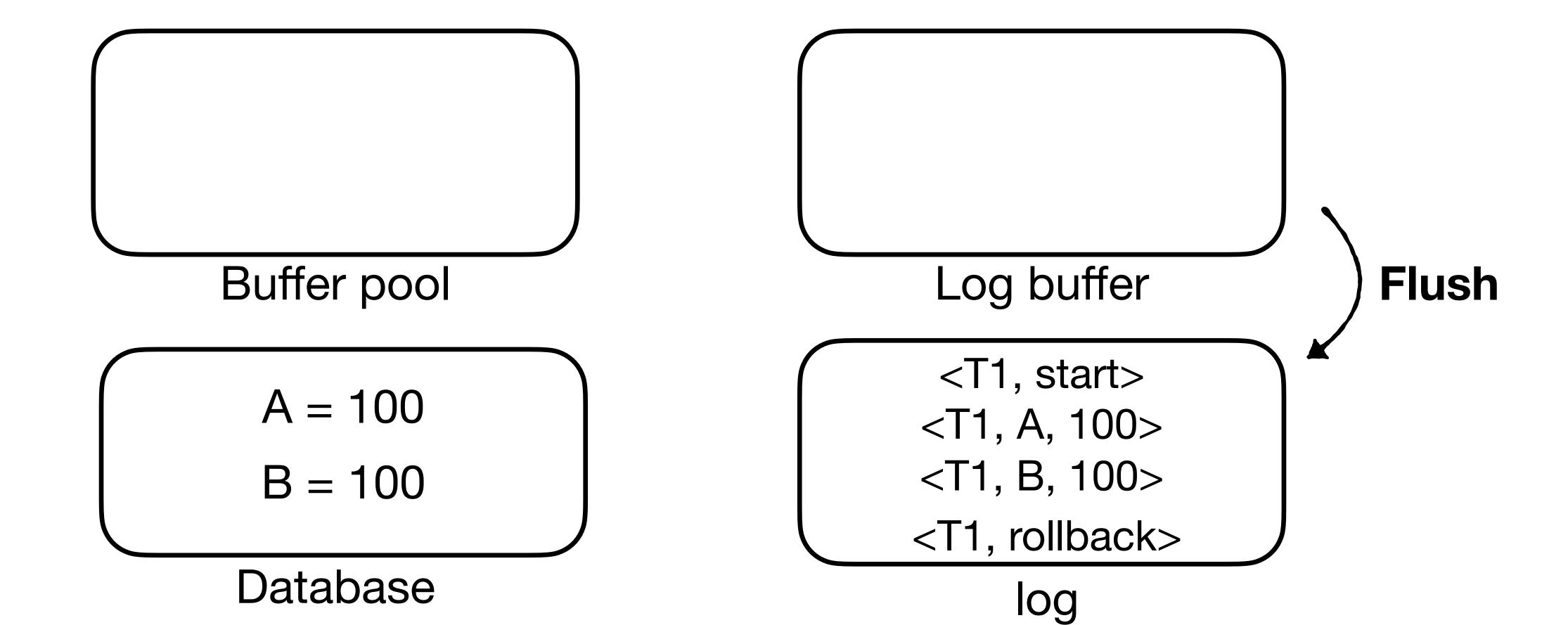
<T1, start>
<T1, A, 100>
<T1, B, 100>

Why do we need this rollback record?

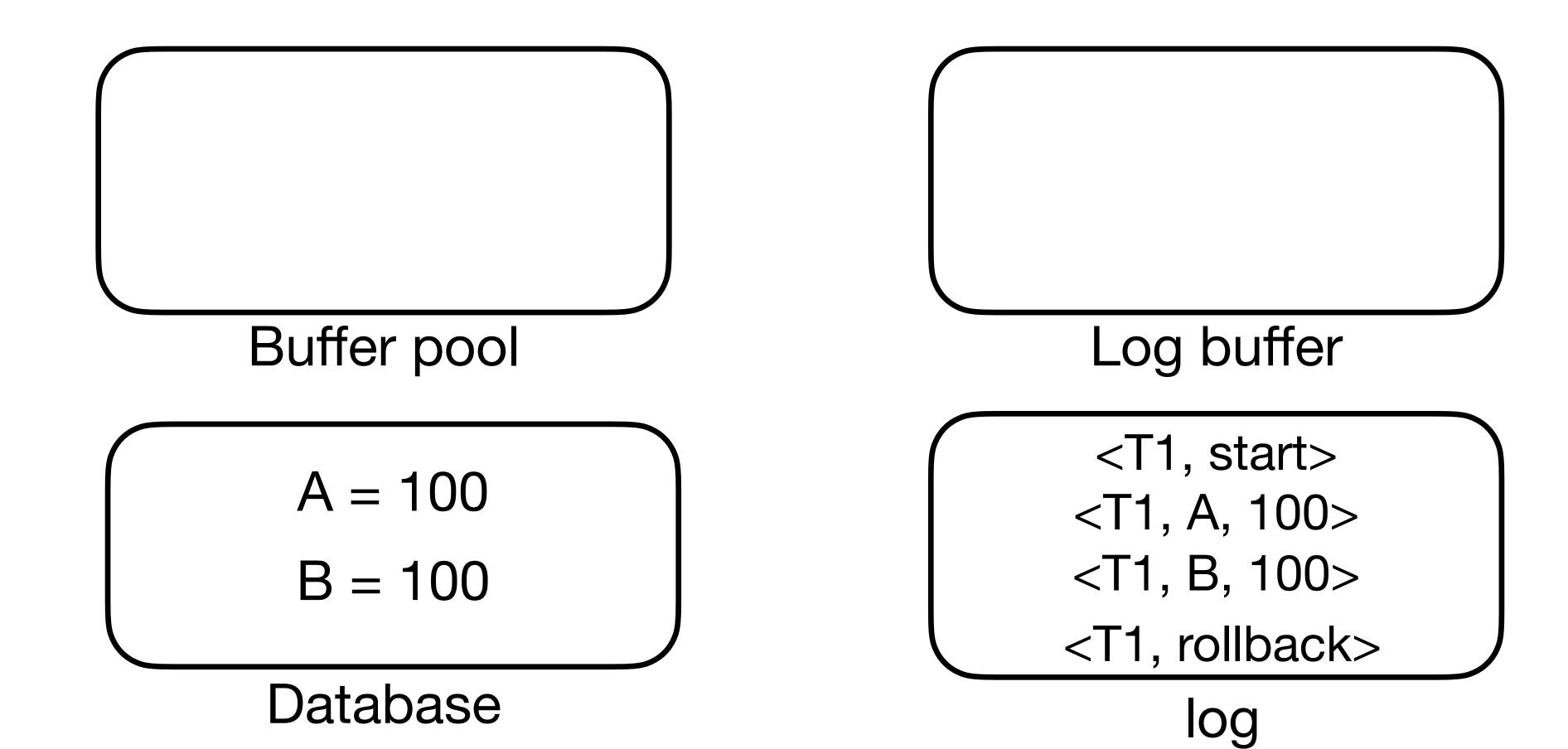


Why do we need this rollback record?

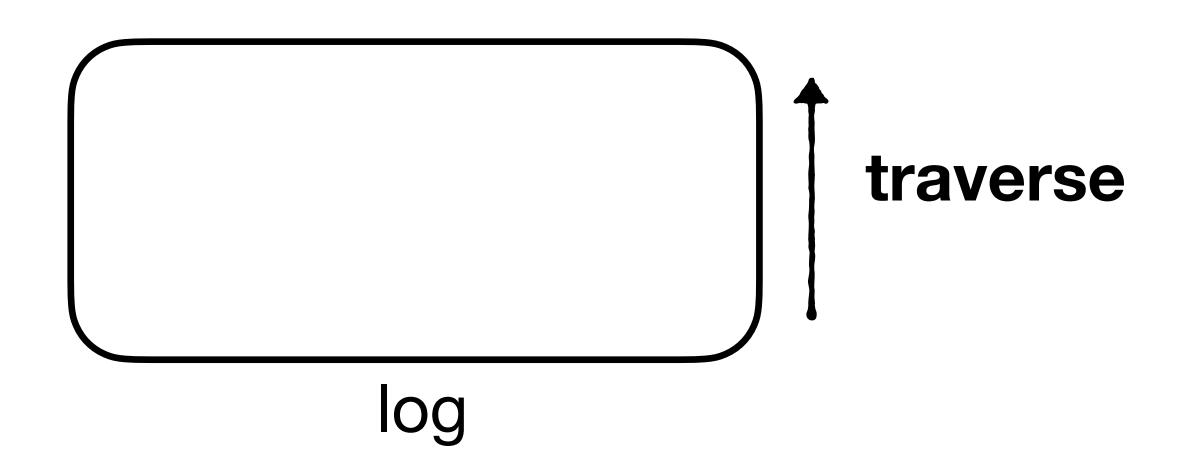
It saves us work if power fails again so we don't undo B again



We are now done. It's as if the original transaction never executed.

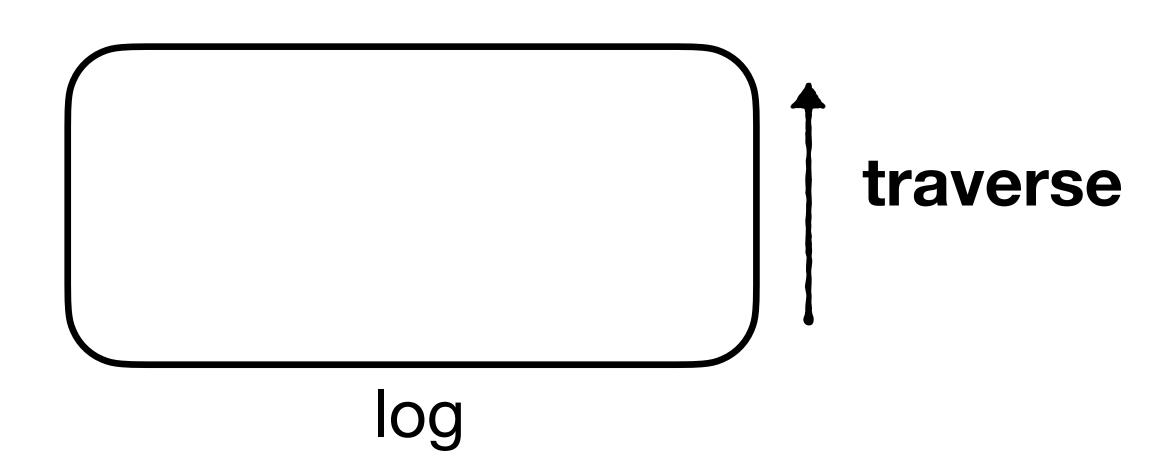


How much of the log must we traverse during recovery?



How much of the log must we traverse during recovery?

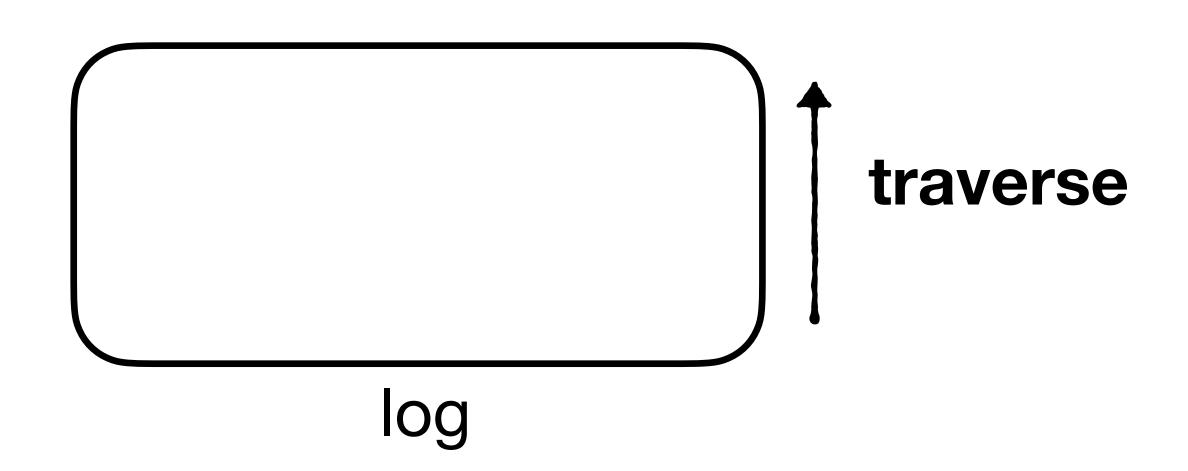
All of it, there may be a transaction from the very beginning that hasn't committed.

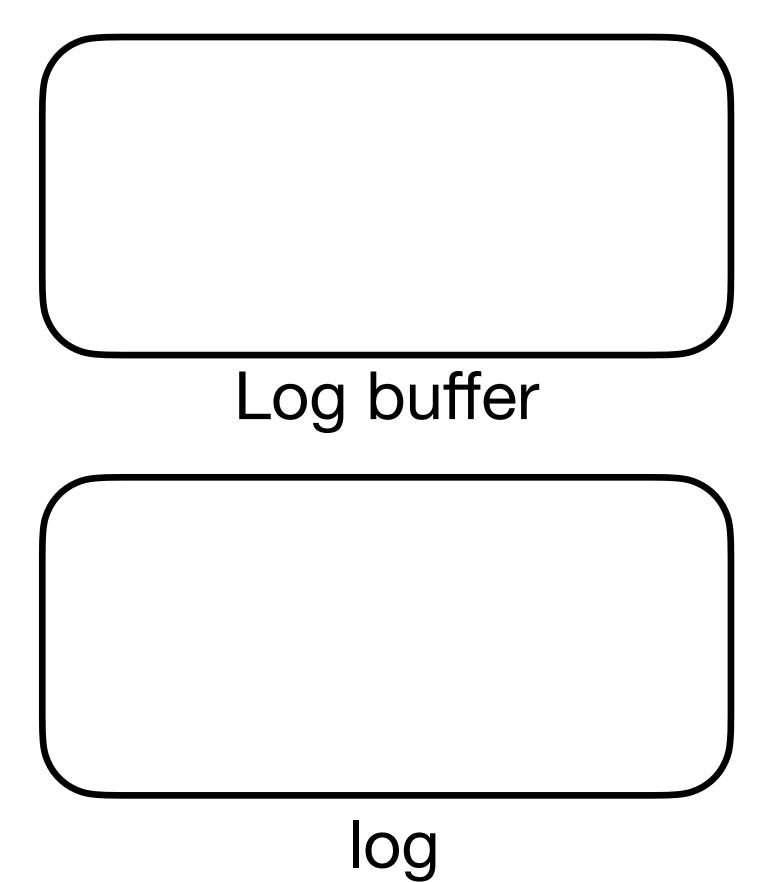


How much of the log must we traverse during recovery?

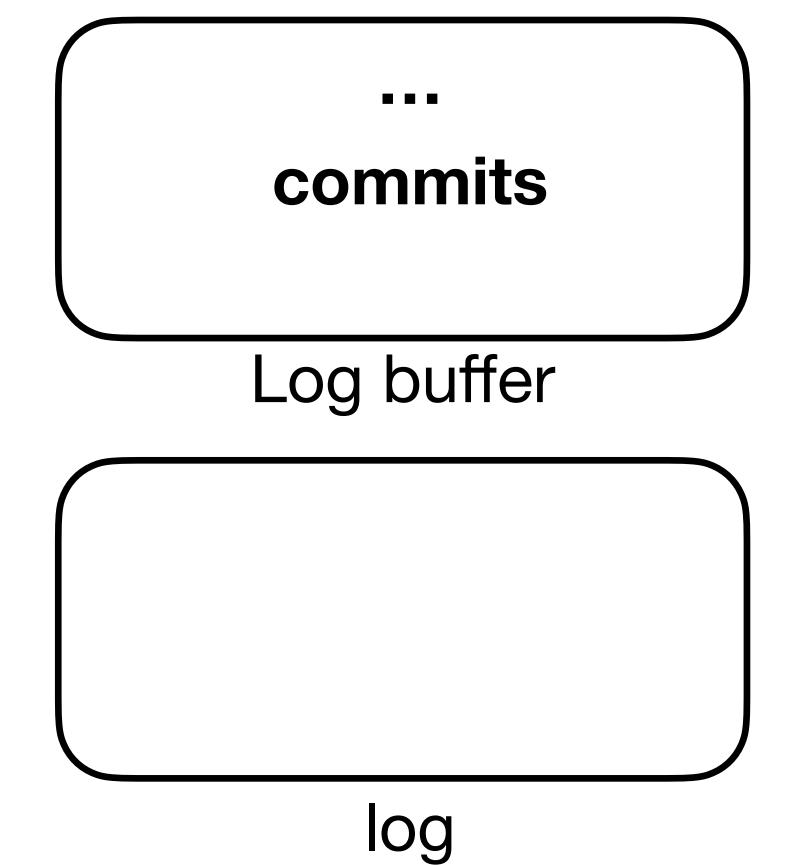
All of it, there may be a transaction from the very beginning that hasn't committed.

How can we bound recovery time without having to traverse the whole log?

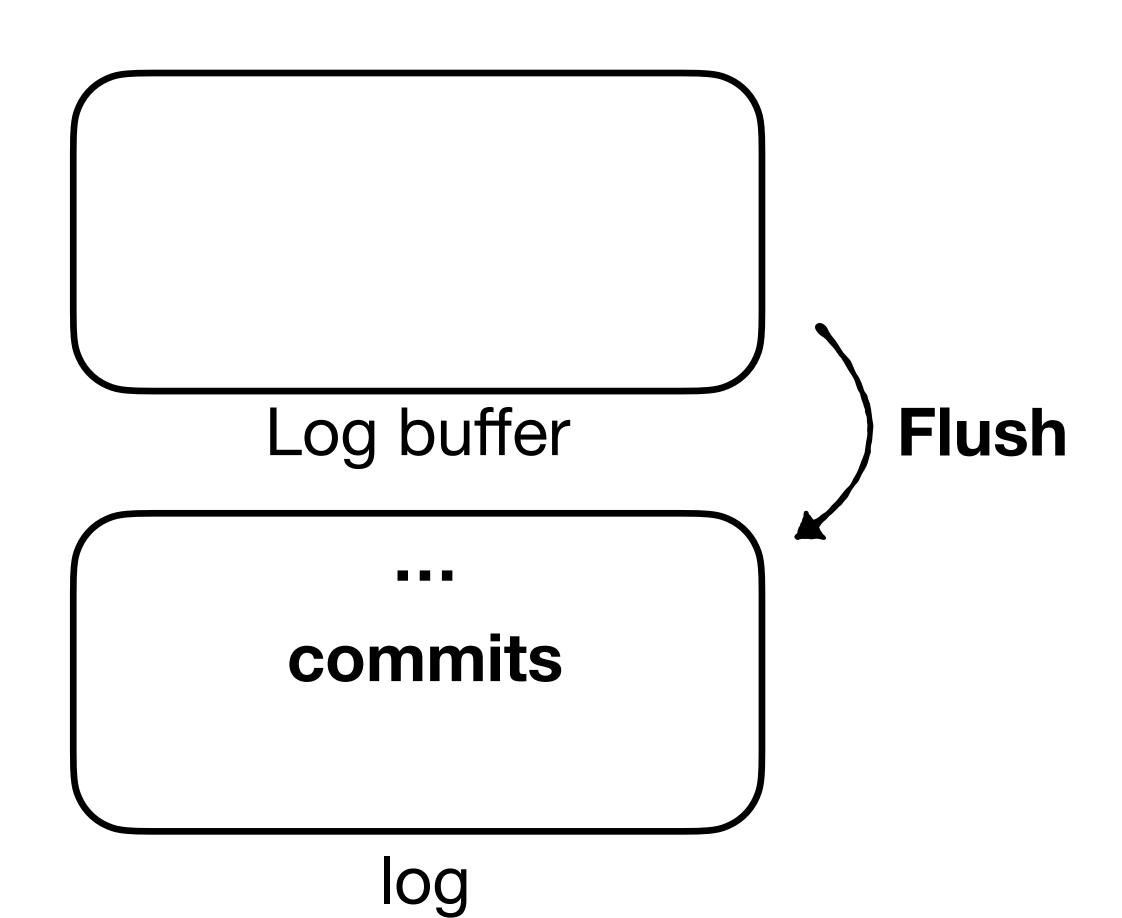




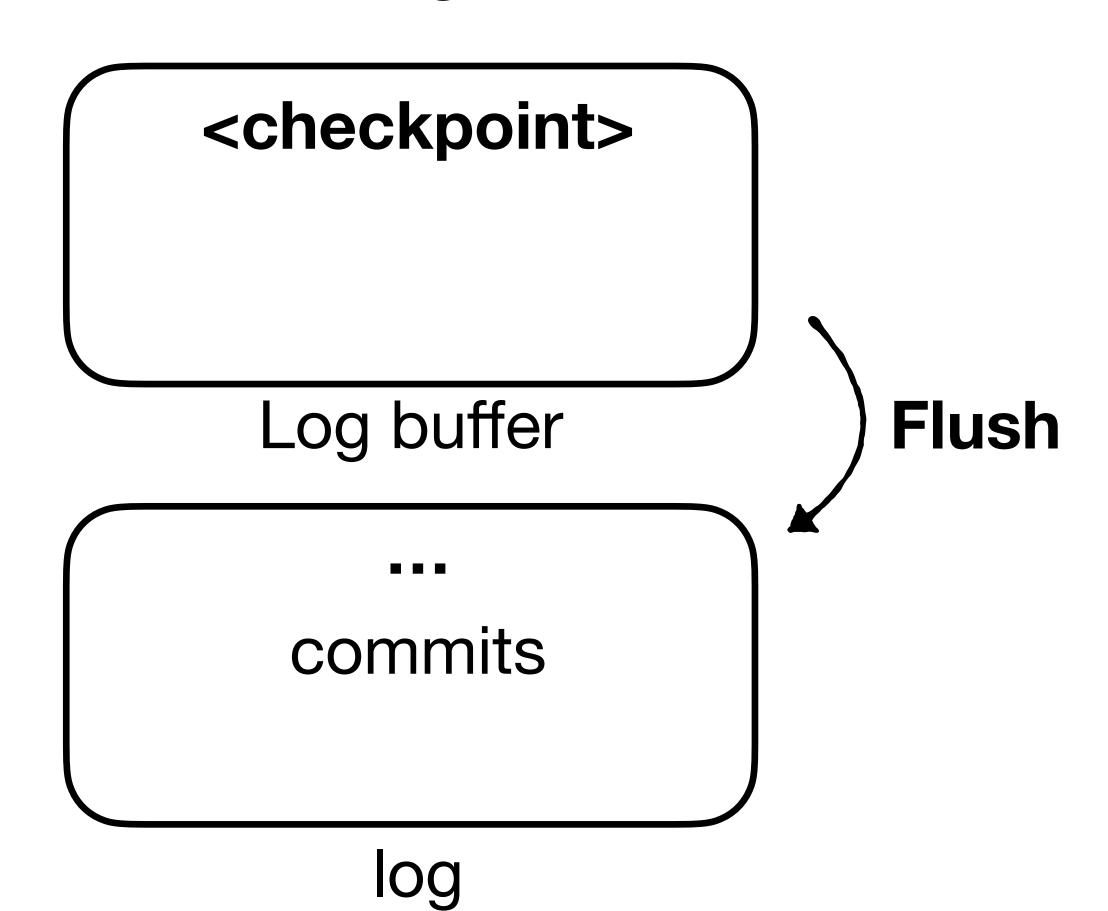
(1) Stop accepting new transactions & wait until all existing ones commit



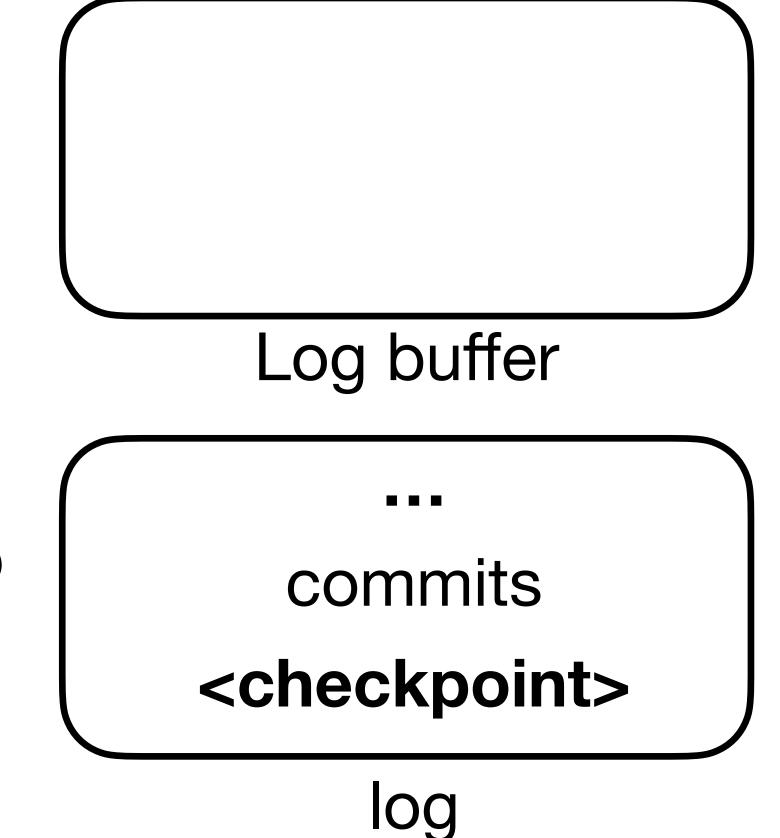
- (1) Stop accepting new transactions & wait until all existing ones commit
- (2) flush log to storage



- (1) Stop accepting new transactions & wait until all existing ones commit
- (2) flush log to storage
- (3) add checkpoint record to log buffer and flush again



- (1) Stop accepting new transactions & wait until all existing ones commit
- (2) flush log to storage
- (3) add checkpoint record to log buffer and flush again



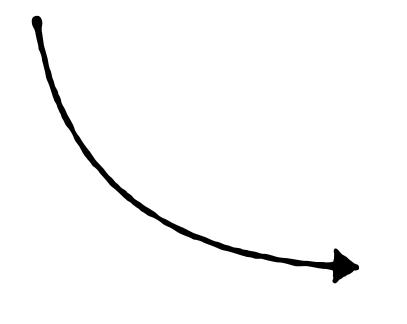
During next recovery, only traverse log up to first checkpoint record we encounter

Problems with UNDO Logging?

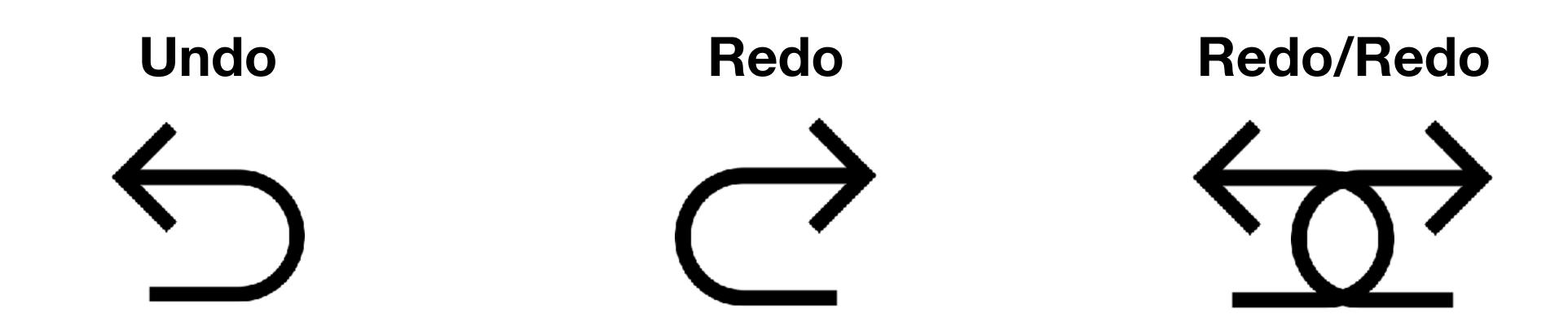
- (1) write before-image for any changed value to log buffer
- (2) flush log
- (3) force the changed data into storage
- (4) add commit record to log buffer & flush

Problems with UNDO Logging?

Many random I/Os



- (1) write before-image for any changed value to log buffer
- (2) flush log
- (3) force the changed data into storage
- (4) add commit record to log buffer & flush



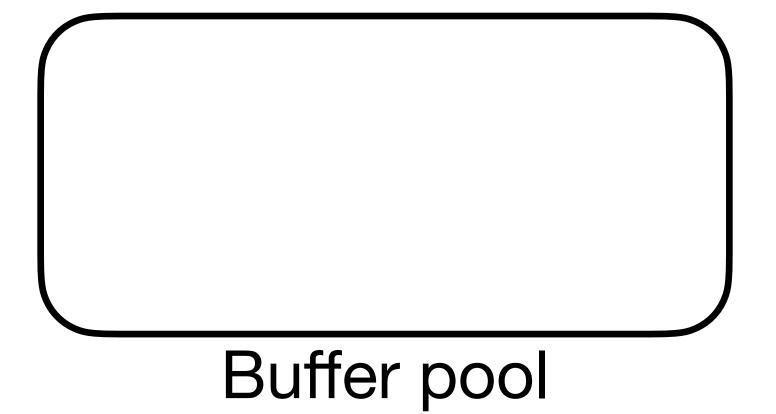
random I/Os

- (1) write after-image for any changed value to log buffer
- (2) keep changed data in buffer pool
- (3) add commit record to log buffer & flush
- (4) flush changed data to storage at leisure

Transaction

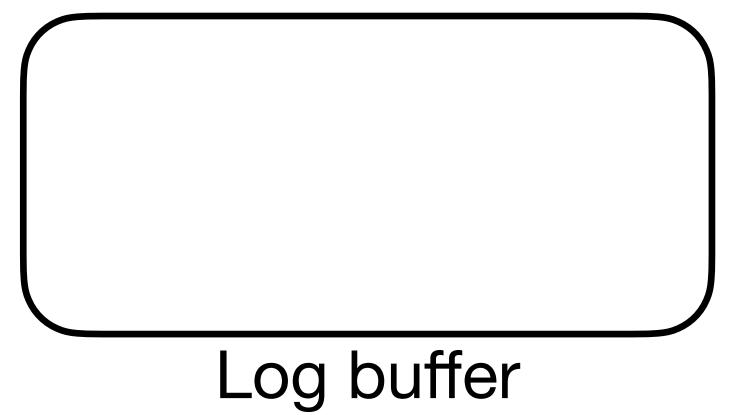
$$A = A - 100$$

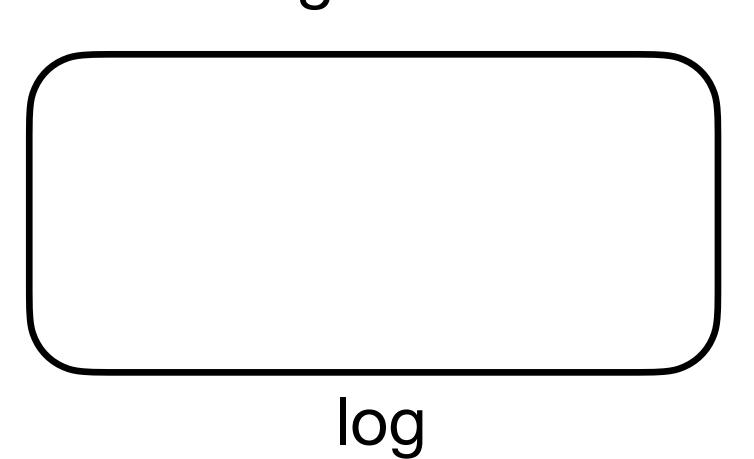
$$B = B + 100$$



$$B = 100$$

Database





Transaction

$$A = A - 100$$

$$B = B + 100$$

$$A = 100$$

B = 100

Buffer pool

$$A = 100$$

$$B = 100$$

Database

<T1, start>

Log buffer

$$A = A - 100$$

$$B = B + 100$$

Write new values to log

$$A = 0$$

$$B = 200$$

<T1, start>

<T1, A, 0>

<T1, B, 200>

Buffer pool

Log buffer

$$A = 100$$

$$B = 100$$

Database

$$A = A - 100$$

$$B = B + 100$$

Write commit record

(diff from undo logging)

$$A = 0$$

$$B = 200$$

Buffer pool

$$A = 100$$

$$B = 100$$

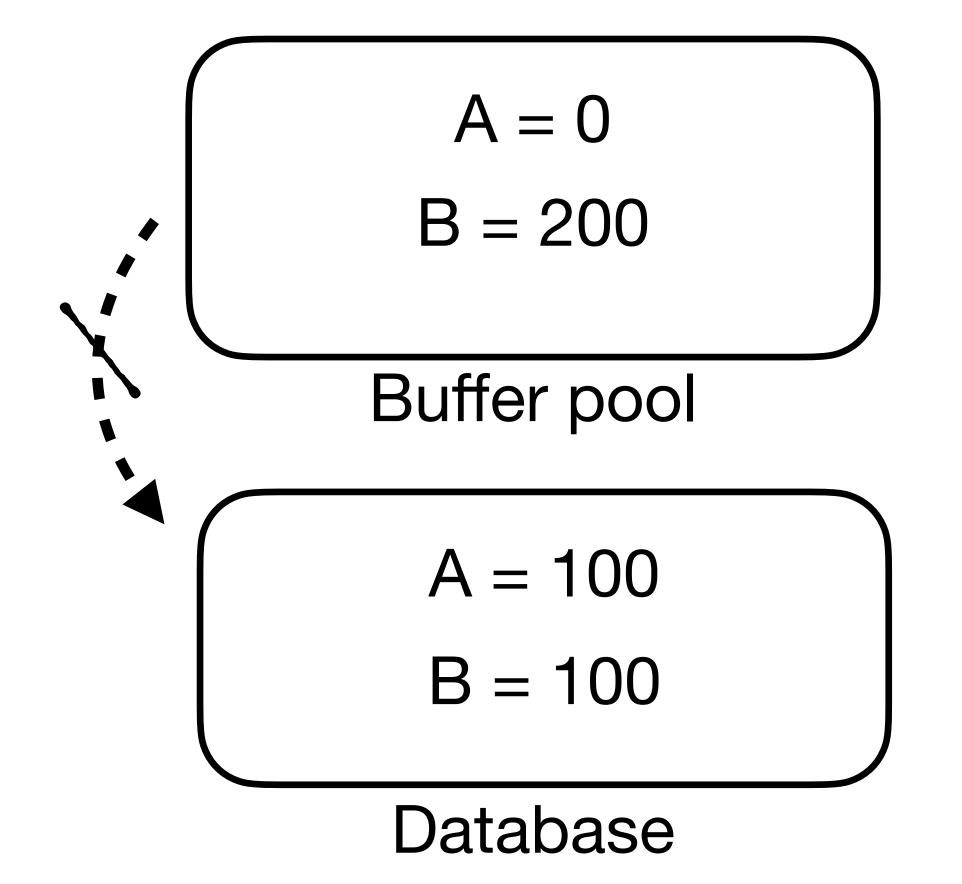
Database

Log buffer

$$A = A - 100$$

$$B = B + 100$$

Not allowed to evict until log buffer flushes

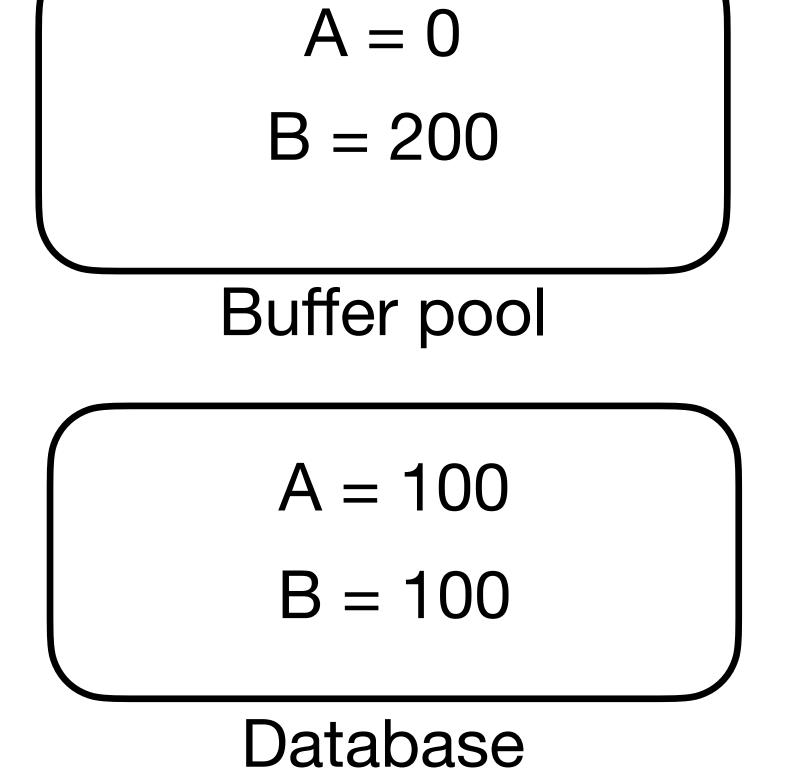


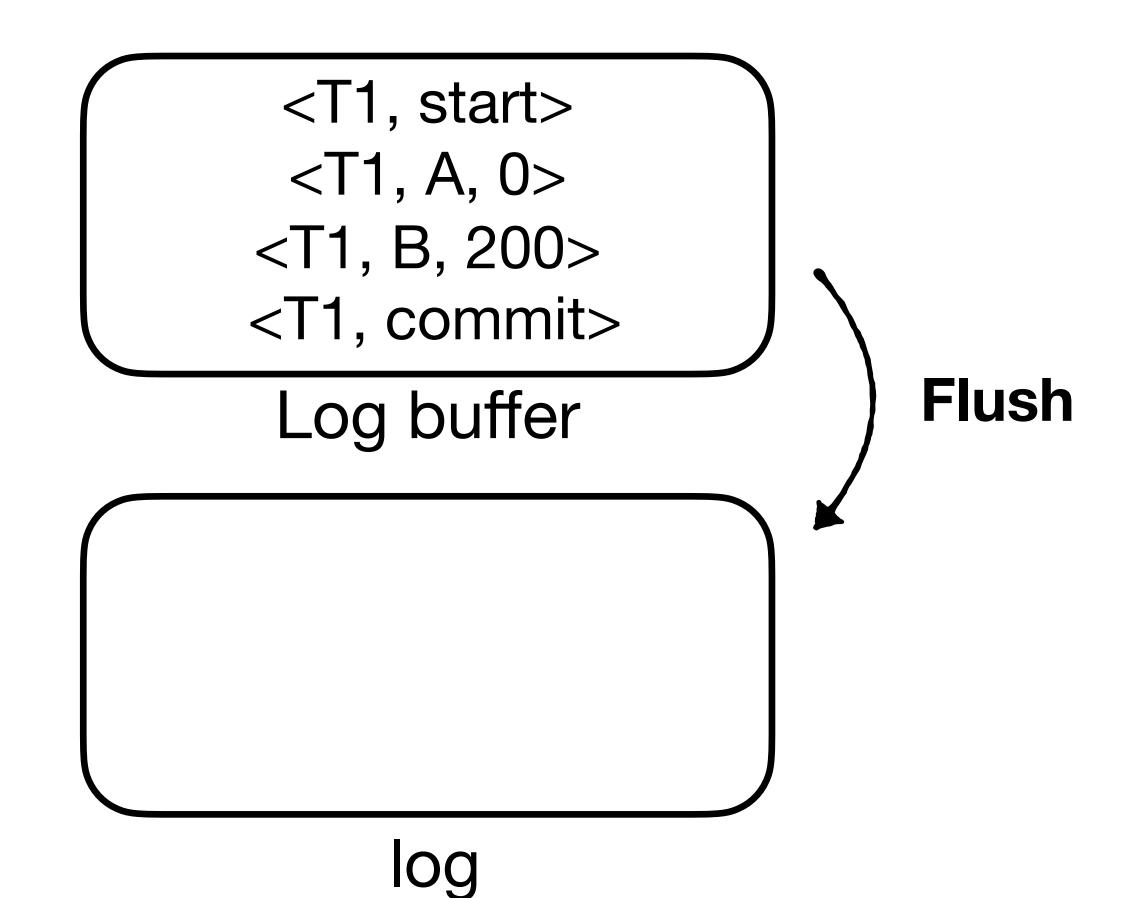
Log buffer

Transaction

$$A = A - 100$$

$$B = B + 100$$

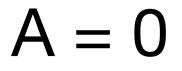




Transaction

$$A = A - 100$$

$$B = B + 100$$



B = 200

Buffer pool

$$A = 100$$

B = 100

Database

Log buffer

<T1, start>

<T1, A, 0>

<T1, B, 200>

<T1, commit>

log

Flush

$$A = A - 100$$

$$B = B + 100$$

Transaction is now committed

$$A = 0$$

$$B = 200$$

Buffer pool

$$A = 100$$

$$B = 100$$

Database

Log buffer

<T1, start>

<T1, A, 0>

<T1, B, 200>

<T1, commit>

Transaction

$$A = A - 100$$

$$B = B + 100$$

Now allowed to evict at leisure

A = 0

B = 200

evict (

Buffer pool

A = 100

B = 100

Database

Log buffer

<T1, start>

<T1, A, 0>

<T1, B, 200>

<T1, commit>

$$A = A - 100$$

$$B = B + 100$$

Note that we don't have to force all changes to storage at once

(A stays in memory in this example)

$$A = 100$$

$$B = 200$$

Database

<T1, start>
<T1, A, 0>
<T1, B, 200>
<T1, commit>

log

Transaction

$$A = A - 100$$

$$B = B + 100$$

The buffer pool can evict autonomously based on clock or LRU:)



Buffer pool

$$A = 100$$

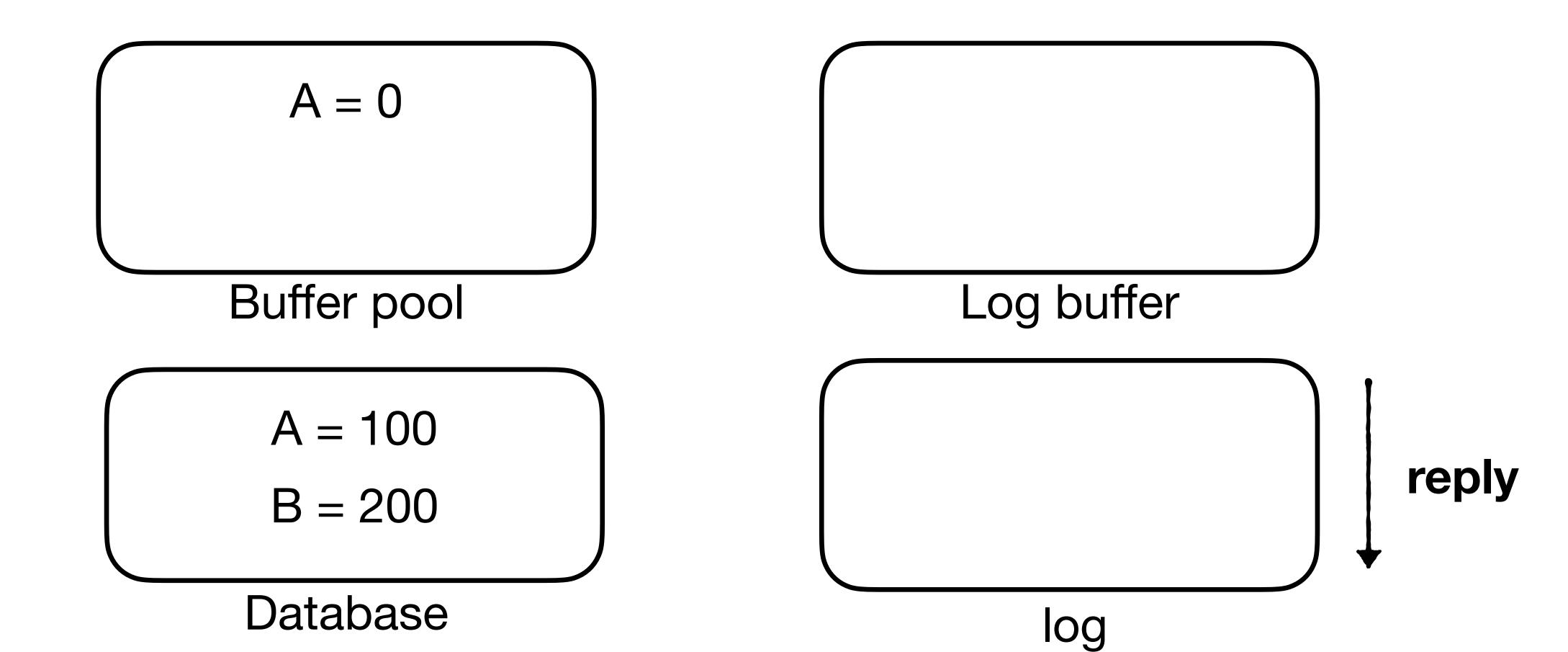
$$B = 200$$

Database



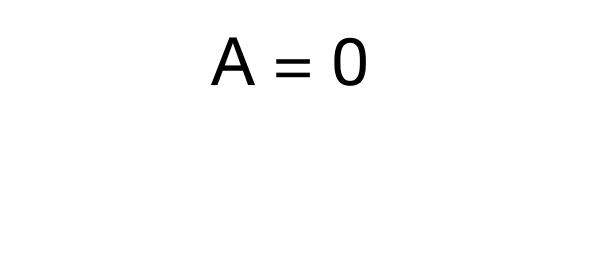
Log buffer

To recover, traverse log forward, replying effects of all committed transactions





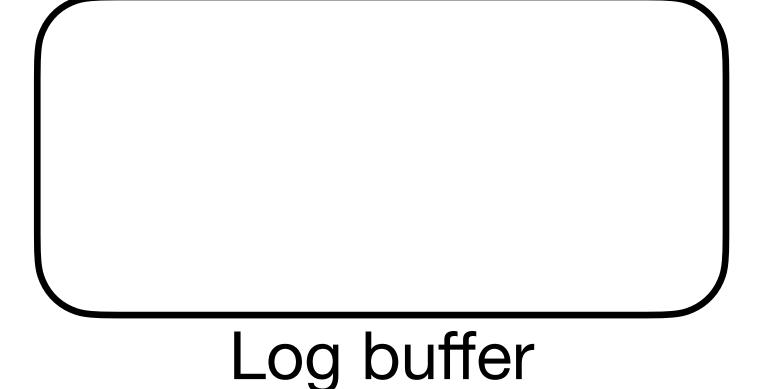
Suppose the transaction was not committed and power failed

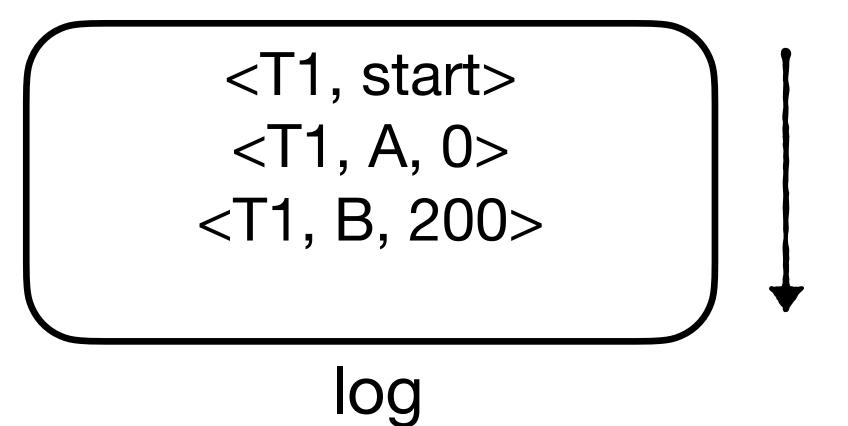


Buffer pool

$$A = 100$$

$$B = 100$$
Database

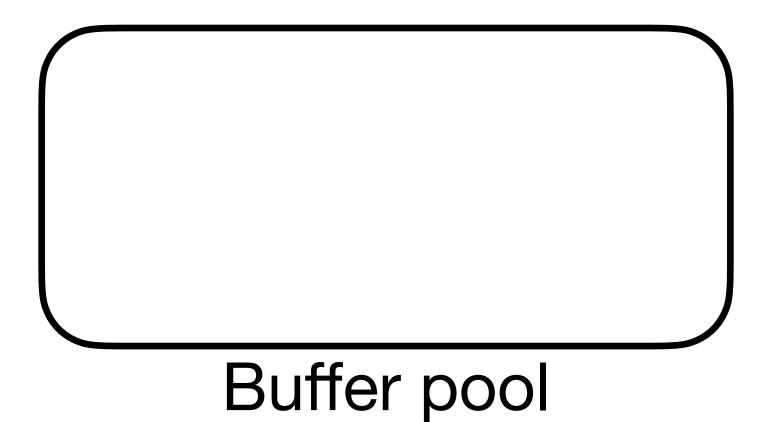






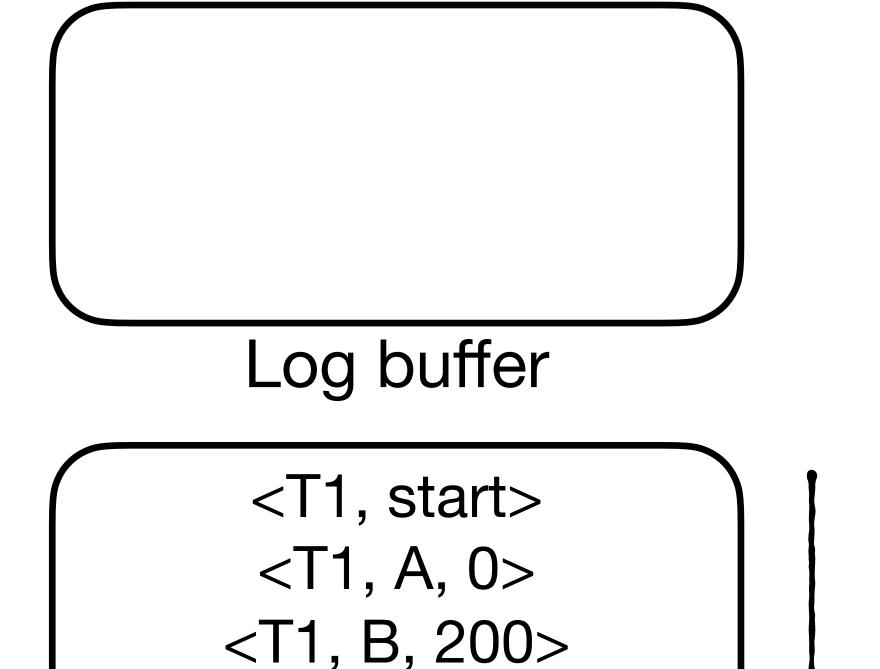
Suppose the transaction was not committed and power failed

If there is no commit, nothing to do as we know modified versions didn't reach storage

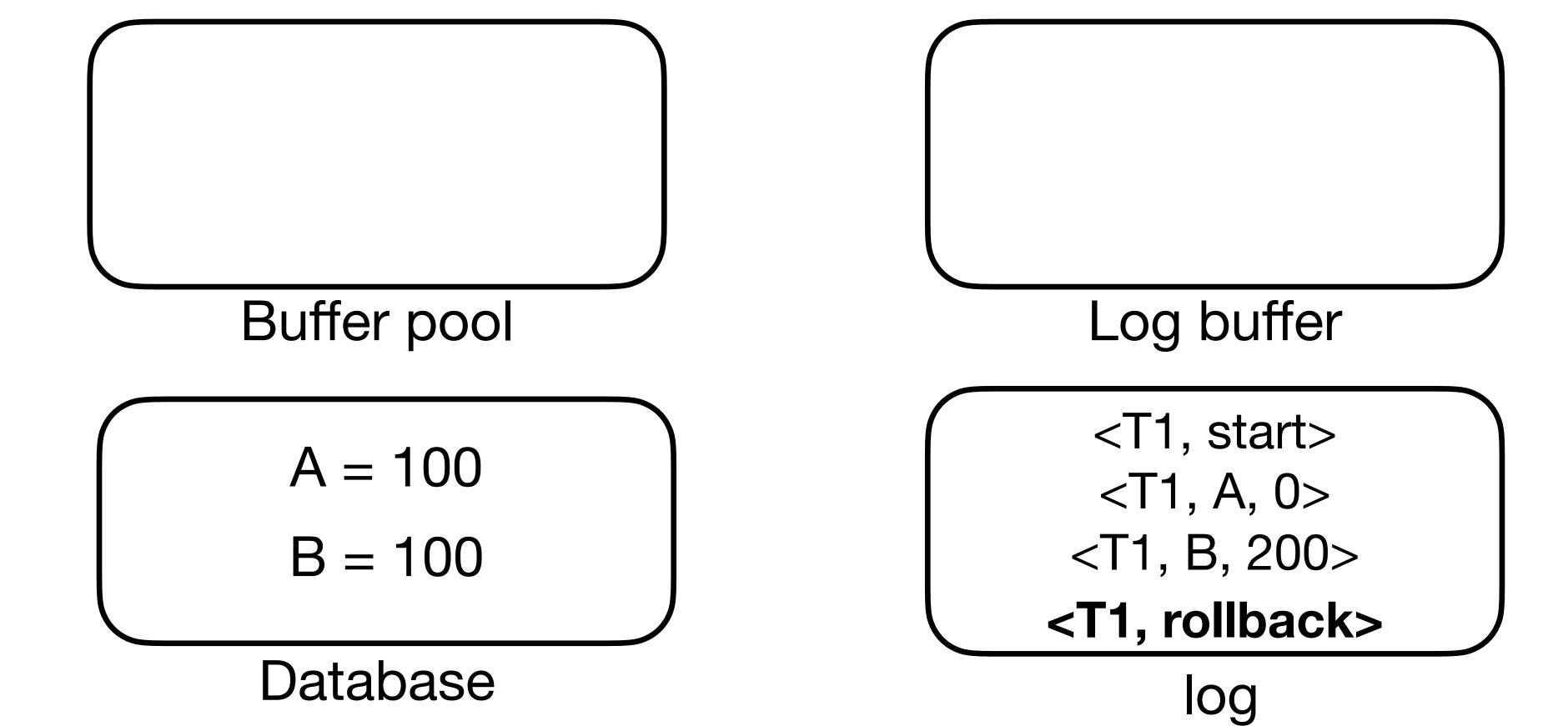


$$A = 100$$

$$B = 100$$
Database

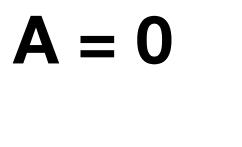


Just add rollback record





Now, suppose the transaction committed but power failed before all changes reached storage

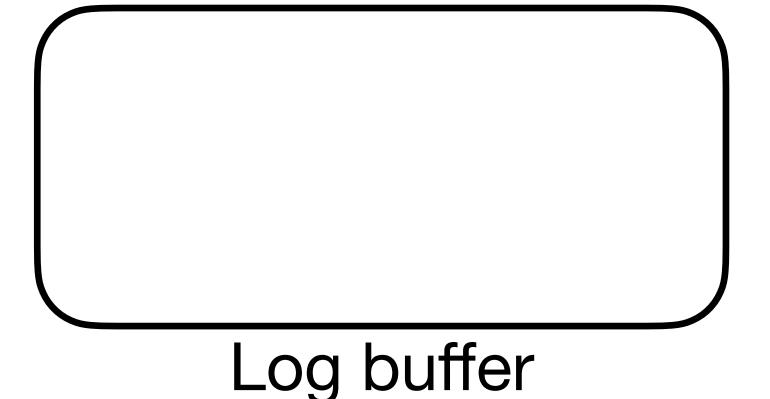


Buffer pool

$$A = 100$$

$$B = 200$$

Database





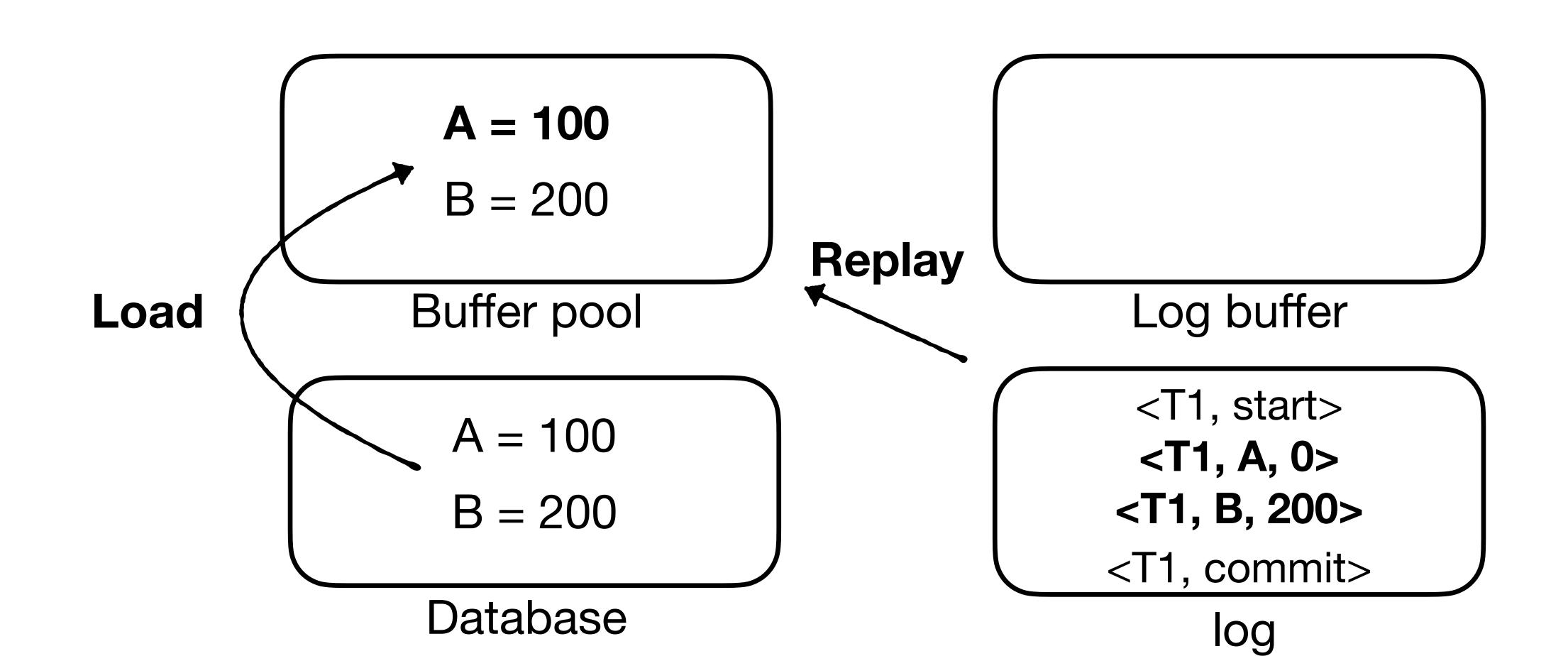
Now, suppose the transaction committed but power failed before all changes reached storage

How to recover?

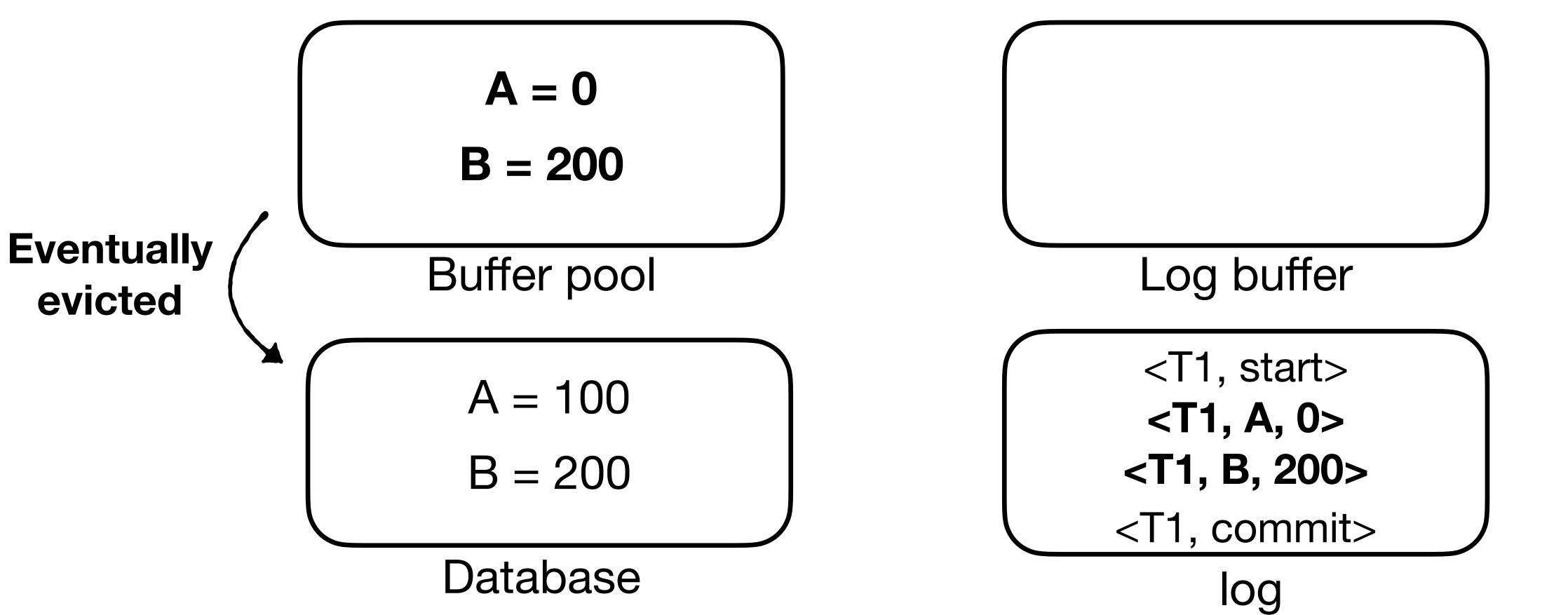


A = 100 B = 200Database

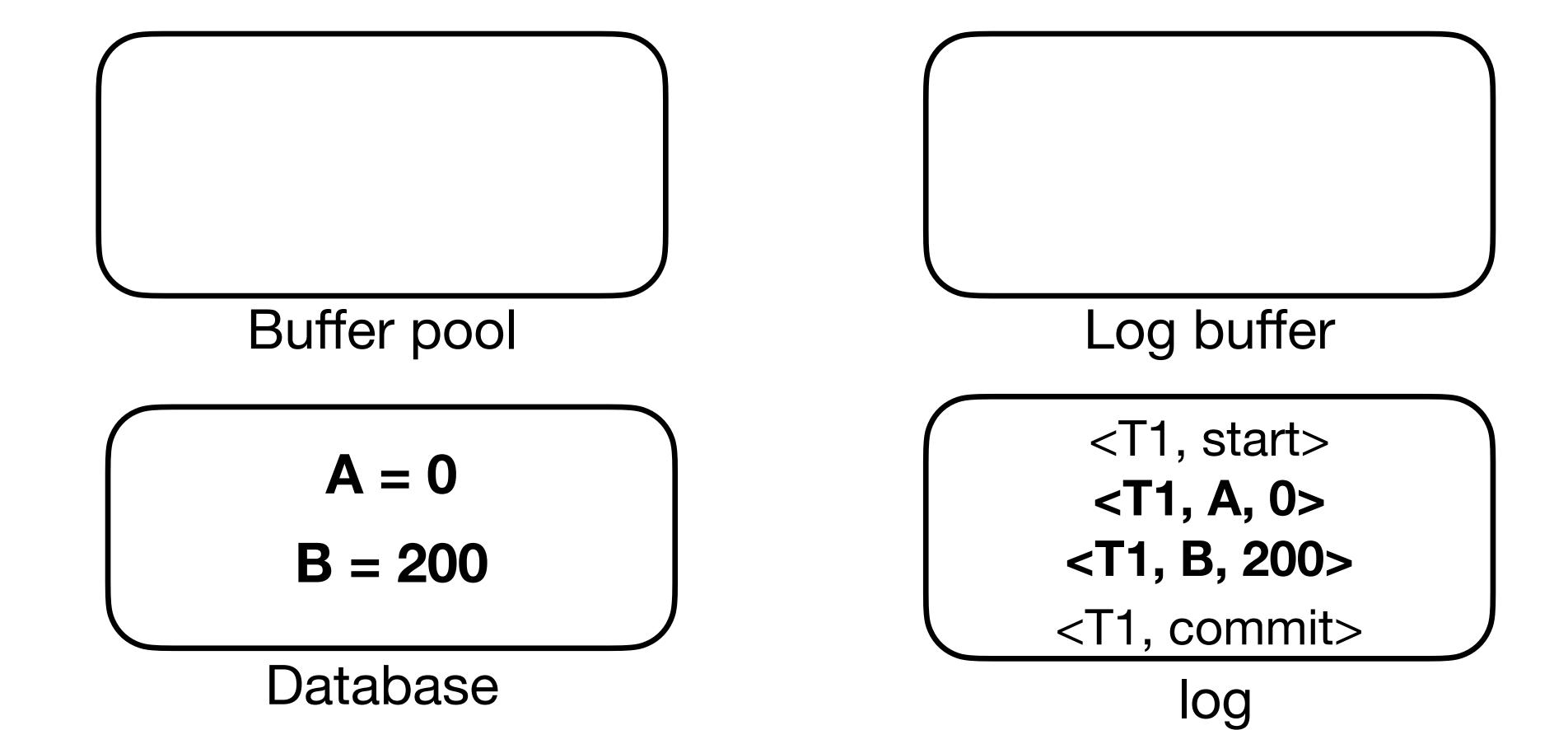
Now, suppose the transaction committed but power failed before all changes reached storage



Now, suppose the transaction committed but power failed before all changes reached storage



Now, suppose the transaction committed but power failed before all changes reached storage

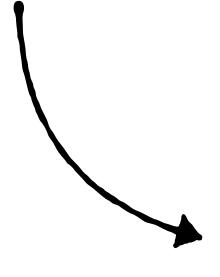


Any problems with Redo logging?

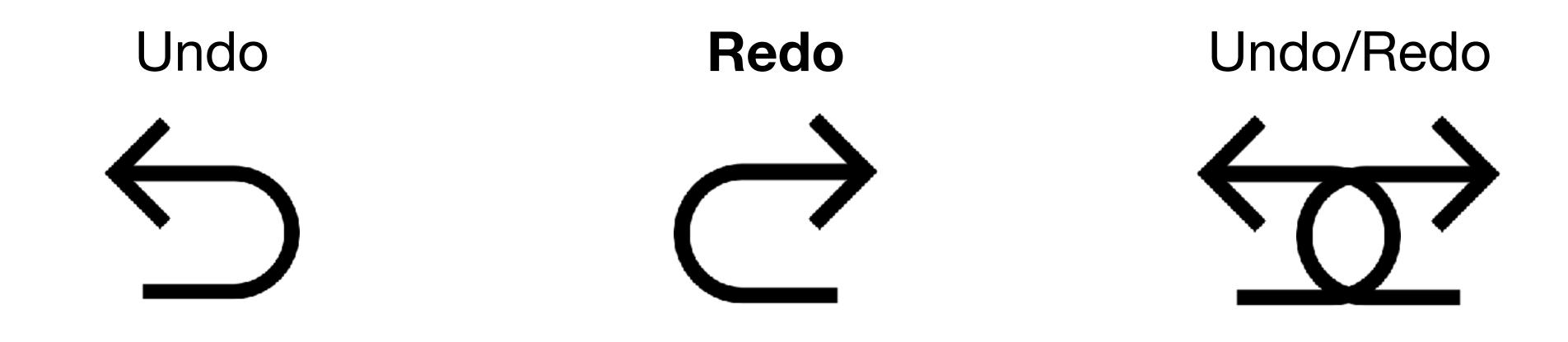
- (1) write after-image for any changed value to log buffer
- (2) keep changed data in buffer pool
- (3) add commit record to log buffer & flush
- (4) flush changed data to storage at leisure

Any problems with Redo logging?

Holds up memory



- (1) write after-image for any changed value to log buffer
- (2) keep changed data in buffer pool
- (3) add commit record to log buffer & flush
- (4) flush changed data to storage at leisure



random I/Os

holds memory

Undo

Redo

Redo/Redo

random I/Os

holds memory

Addresses both problems!

Undo/Redo logging rules

- (1) write before and after-image for any changed value to log buffer
- (2) before modifying item on disk, must flush log record
- (3) when transaction is finished, flush commit record

Undo/Redo logging rules

- (1) write before and after-image for any changed value to log buffer
- (2) before modifying item on disk, must flush log record
- (3) when transaction is finished, flush commit record

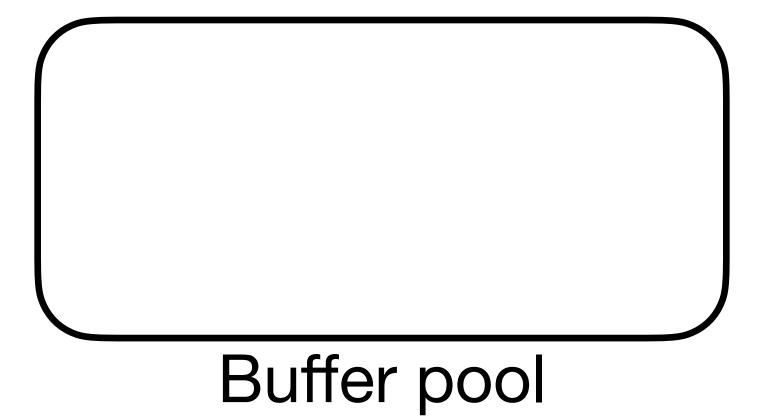
Undo/Redo logging rules

- (1) write before and after-image for any changed value to log buffer
- (2) before modifying item on disk, must flush log record
- (3) when transaction is finished, flush commit record

Transaction T1

$$A = A - 100$$

$$B = B + 100$$



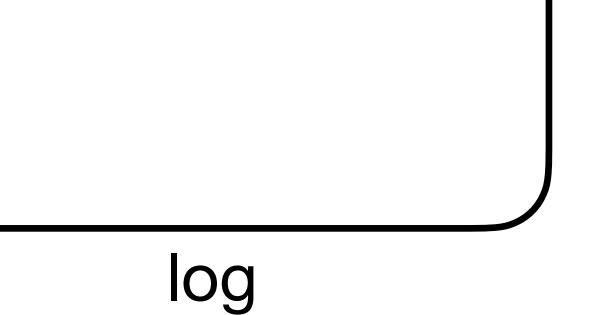
$$A = 100$$

$$B = 100$$

Database



Log buffer



$$A = A - 100$$

$$B = B + 100$$

$$A = 100$$

$$B = 100$$

Buffer pool

$$A = 100$$

$$B = 100$$

Database

<T1, start>

Log buffer

$$A = A - 100$$

$$B = B + 100$$

Write before & after image to log

A = 0

B = 100

Buffer pool

A = 100

B = 100

Database

<T1, start>
<T1, A, 100, 0>

Log buffer

$$A = A - 100$$

$$B = B + 100$$

Suppose log flushes due to other transaction



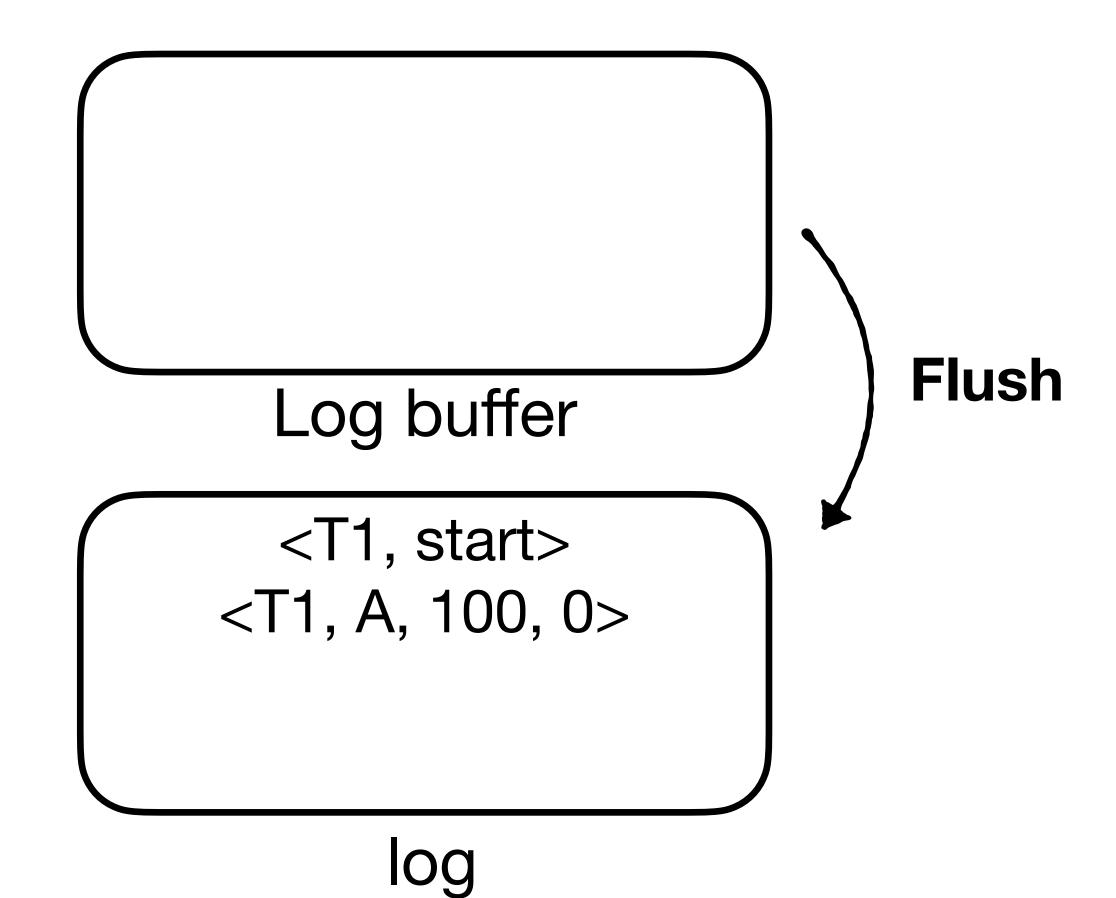
B = 100

Buffer pool

$$A = 100$$

$$B = 100$$

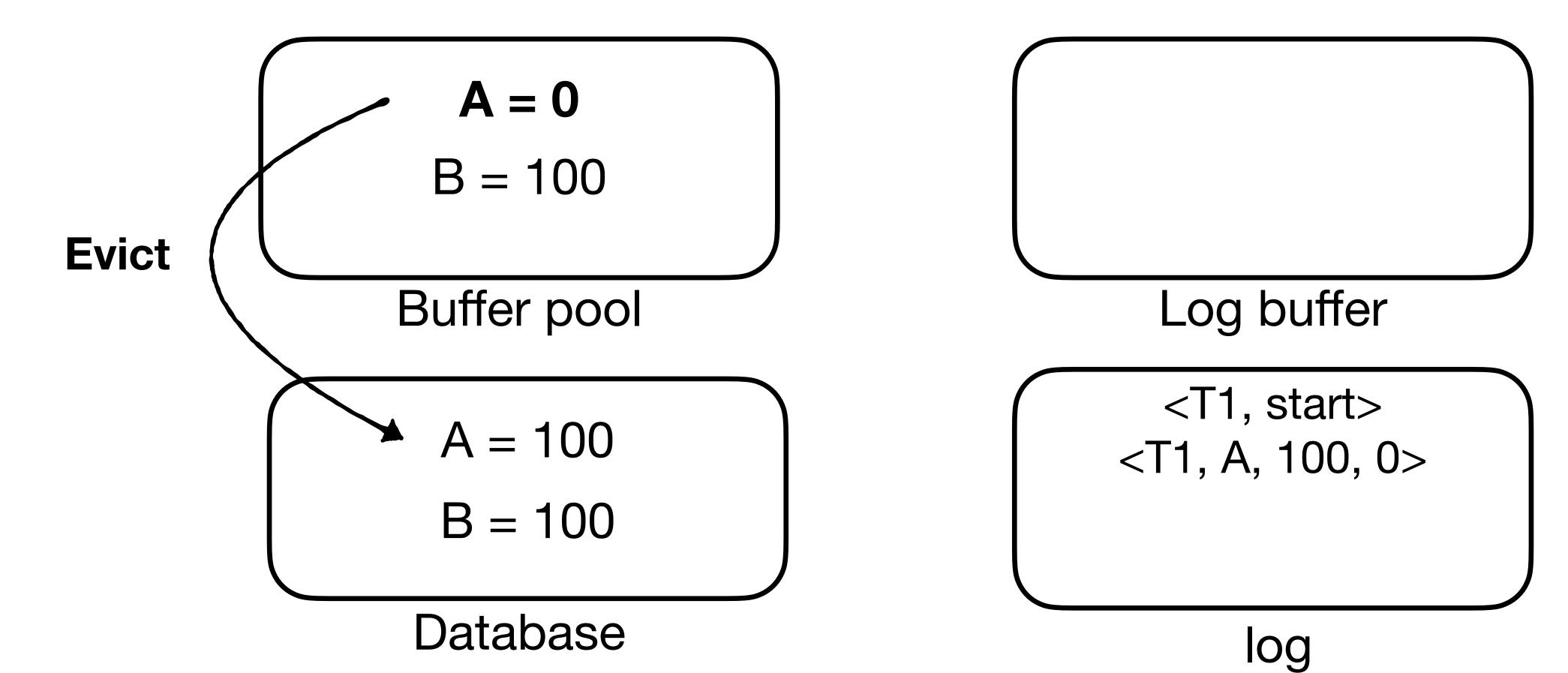
Database



$$A = A - 100$$

$$B = B + 100$$

We are now allowed to evict A if buffer pool needs to



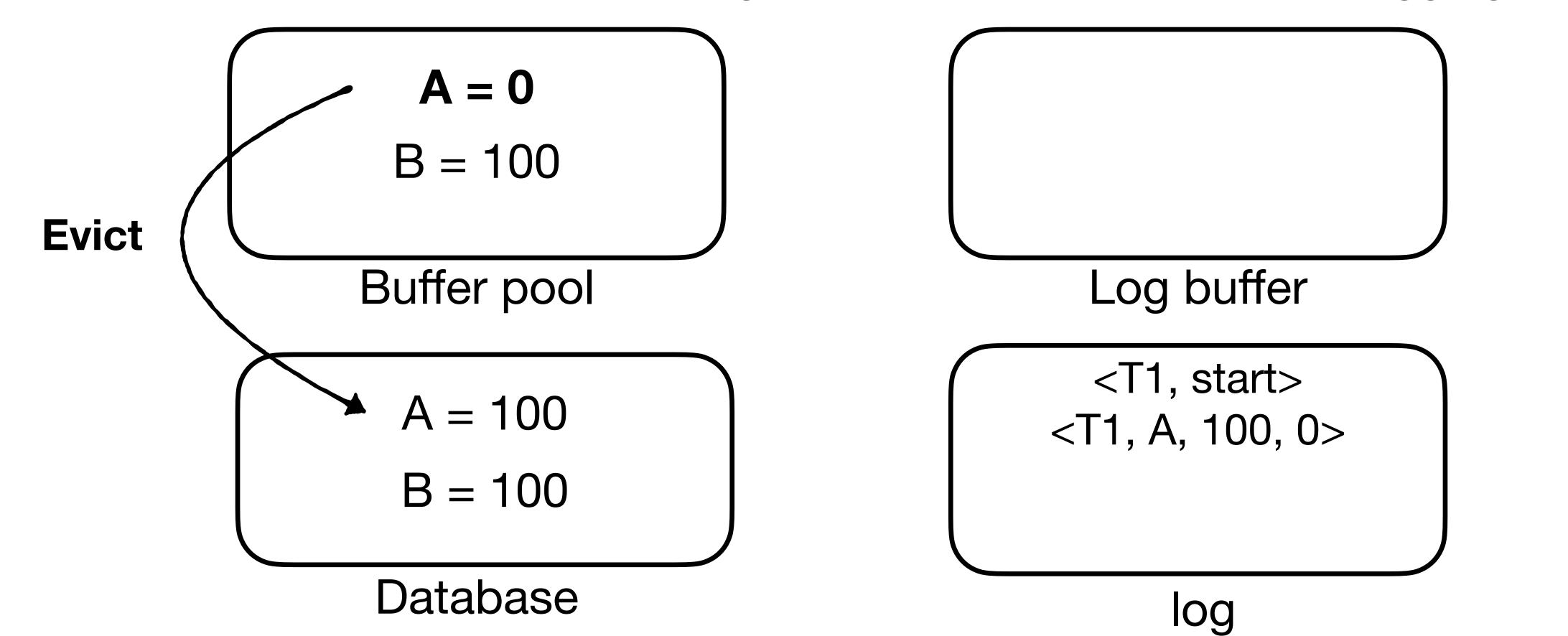
Transaction T1

$$A = A - 100$$

$$B = B + 100$$

We are now allowed to evict A if buffer pool needs to

we do not need to keep changed data in memory like undo logging:)



$$A = A - 100$$

$$B = B + 100$$



If power fails at this point, we could undo change to A using log

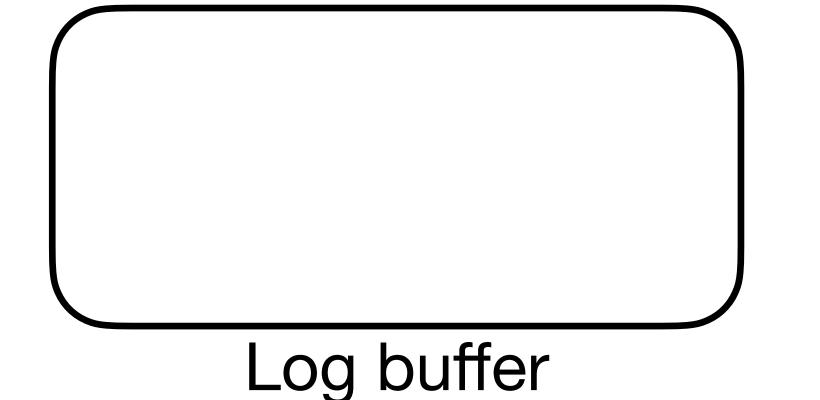
B = 100

Buffer pool

A = 0

B = 100

Database



$$A = A - 100$$

$$B = B + 100$$

Otherwise, continue with the transaction

B = 200

Buffer pool

A = 0

B = 100

Database

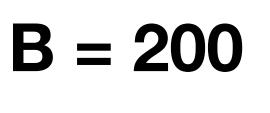
<T1, B, 100, 200> <T1, commit>

Log buffer

<T1, start>
<T1, A, 100, 0>

$$A = A - 100$$

$$B = B + 100$$

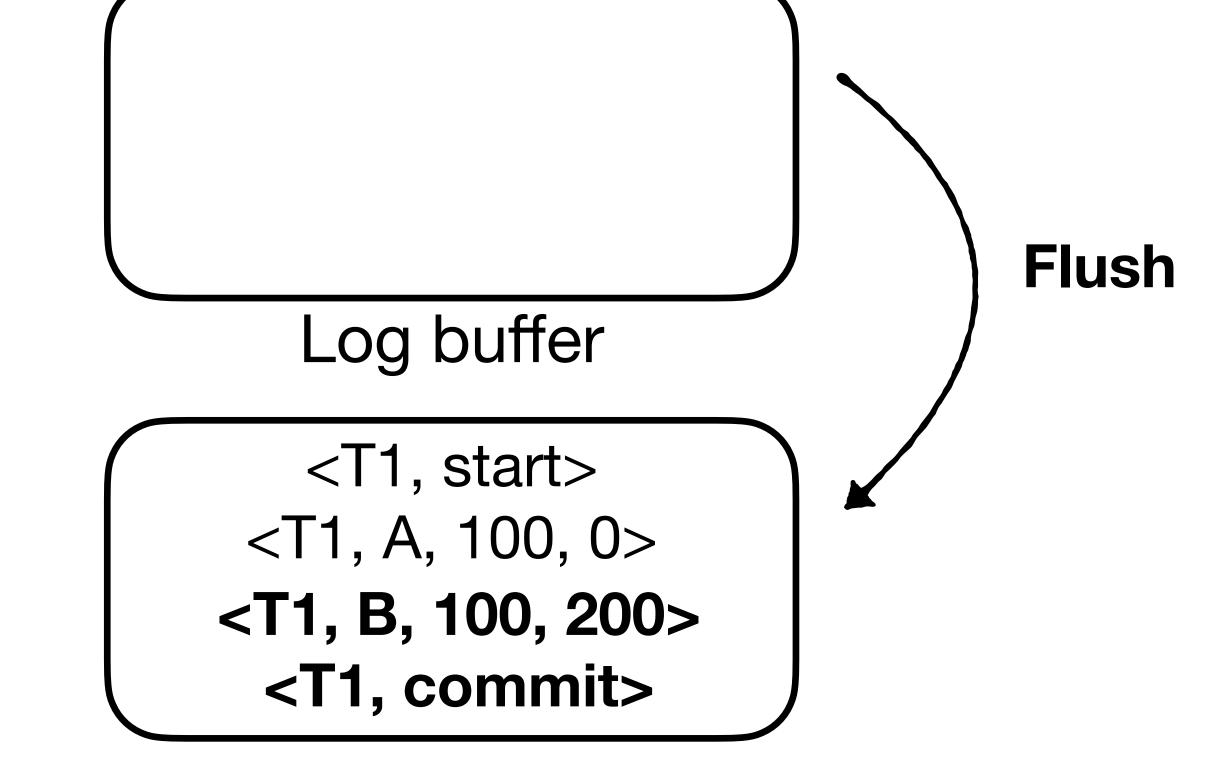


Buffer pool

$$A = 0$$

$$B = 100$$

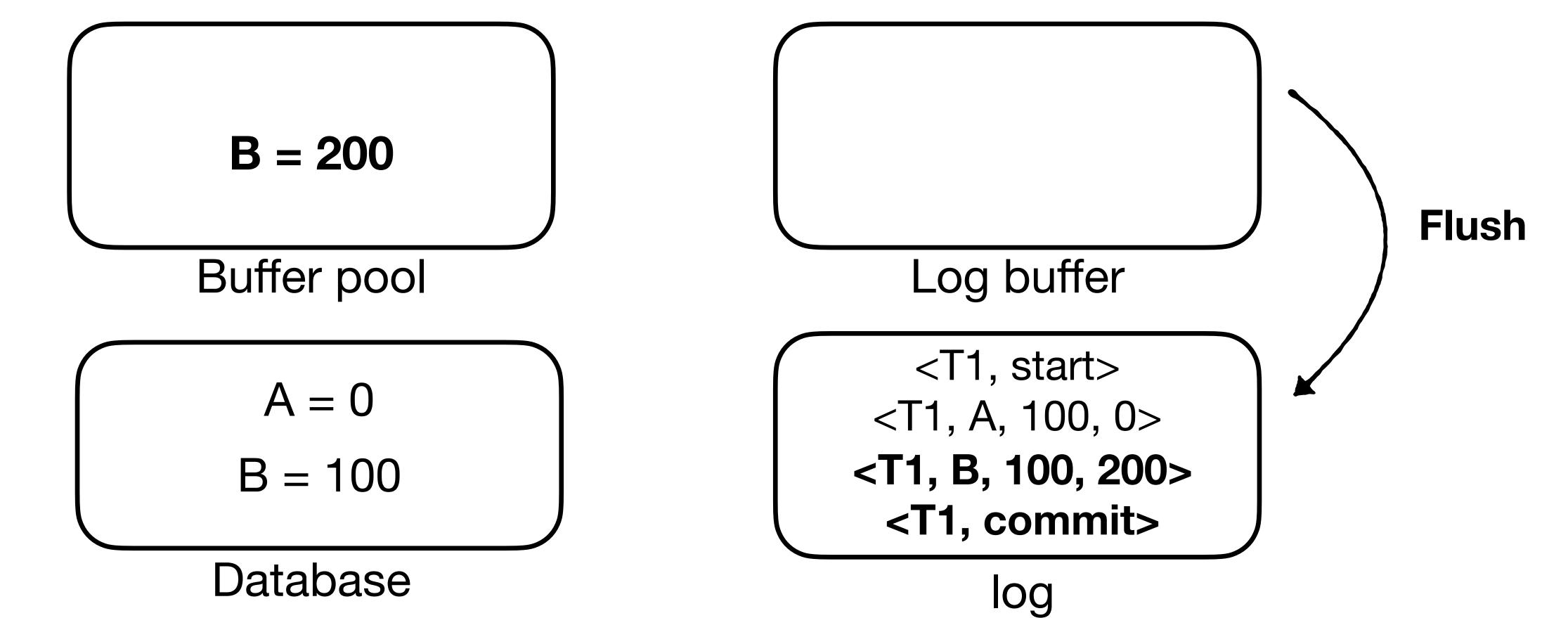
Database



$$A = A - 100$$

$$B = B + 100$$

Do do not now need to force changed data into storage like undo logging:)



$$A = A - 100$$

$$B = B + 100$$



Suppose power now fails before we save B to storage

B = 200

Buffer pool

0 = A

B = 100

Database

Log buffer

<T1, start>
<T1, A, 100, 0>
<T1, B, 100, 200>
<T1, commit>

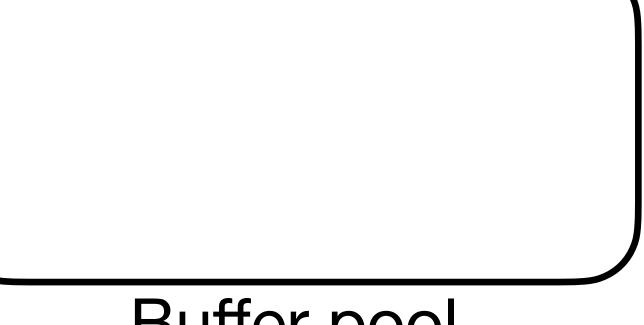
$$A = A - 100$$

$$B = B + 100$$



Suppose power now fails before we save B to storage

How to recover B?

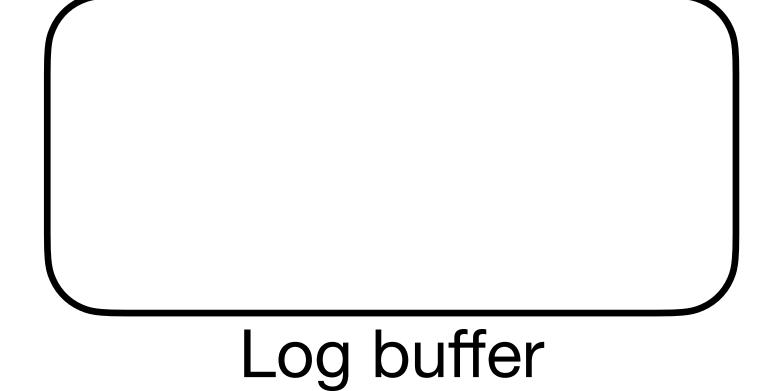


Buffer pool

$$A = 0$$

$$B = 100$$

Database



Transaction T1

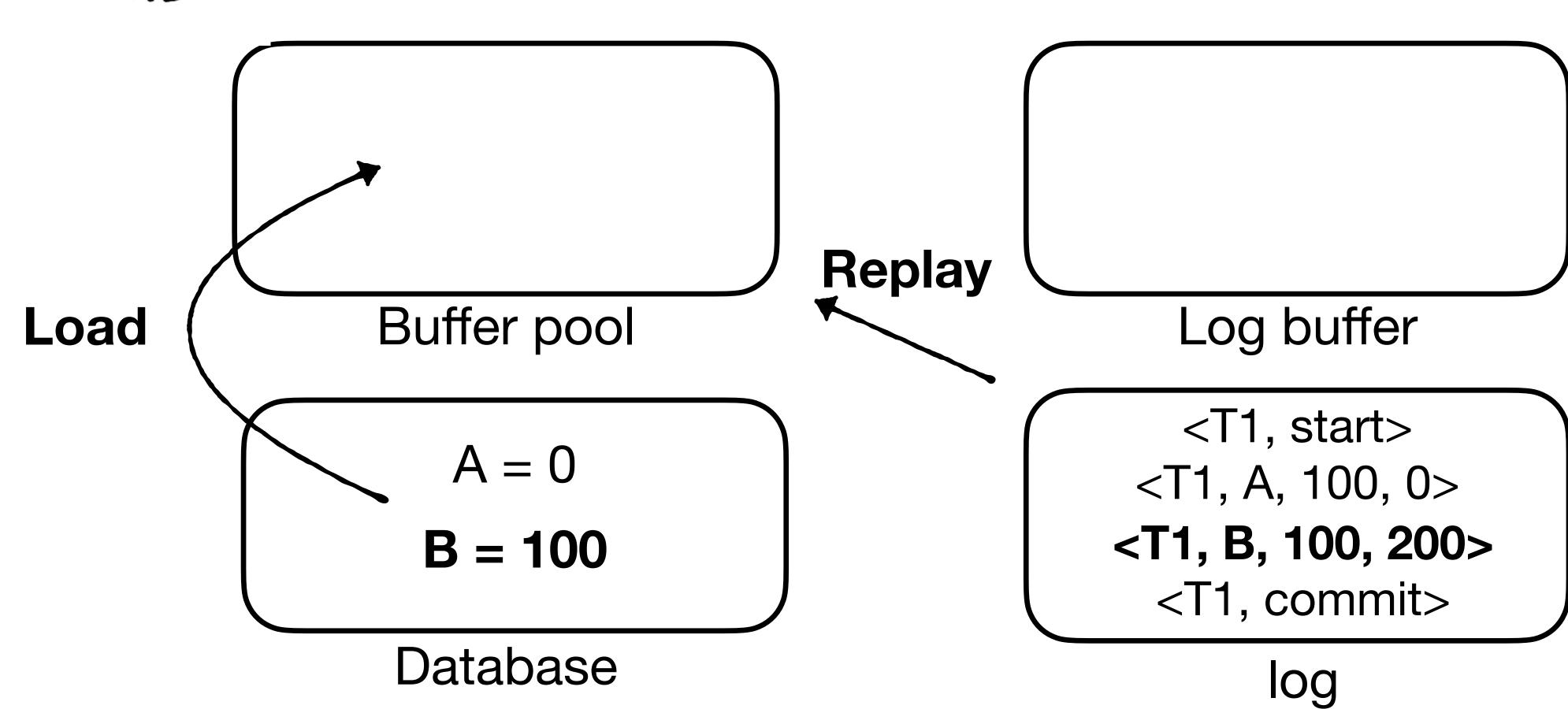
$$A = A - 100$$

$$B = B + 100$$



Suppose power now fails before we save B to storage

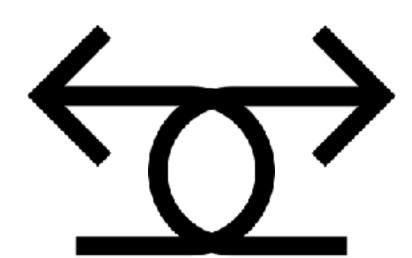
How to recover B?



Undo

Redo

Redo/Redo



random I/Os

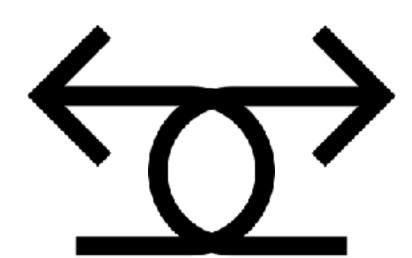
holds memory

Address both problems!

Undo

Redo

Redo/Redo



random I/Os

holds memory

Address both problems!

Tutorial on Recovery beings now:)