

Transaction Processing: Recovery

CPS 216
Advanced Database Systems

Failures

- System crashes in the middle of a transaction T ; partial effects of T were written to disk
 - How do we undo T (atomicity)?
- System crashes right after a transaction T commits; not all effects of T were written to disk
 - How do we complete T (durability)?
- Media fails; data on disk corrupted
 - How do we reconstruct the database (durability)?

4

Review

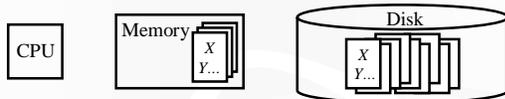
- ACID
 - Atomicity
 - Consistency
 - Isolation → Concurrency control
 - Durability → Recovery

2

Logging

- Log
 - Sequence of log records, recording all changes made to the database
 - Written to stable storage (e.g., disk) during normal operation
 - Used in recovery
- Hey, one change turns into two!
 - Isn't it bad for performance?
 - But writes are sequential (append to the end of log)
 - Can use dedicated disk(s) to improve performance,

Execution model



- $\text{input}(X)$: copy the disk block containing object X to memory
 - Issued by transactions
- $\text{read}(X, v)$: read the value of X into a local variable v (execute $\text{input}(X)$ first if necessary)
- $\text{write}(X, v)$: write value v to X in memory (execute $\text{input}(X)$ first if necessary)
 - Issued by DBMS
- $\text{output}(X)$: write the memory block containing X to disk

3

Undo logging

- Basic idea
 - Every time you modify something on disk, record its old value in the log
 - If system crashes, undo the writes of partially executed transactions by restoring the old values

6

Undo logging example

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

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

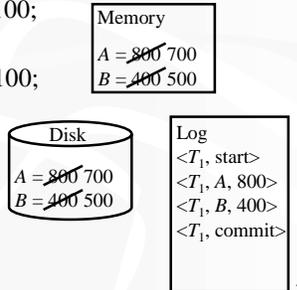
write(A, a);

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

write(B, b);

output(A);

output(B);



7

Another technicality

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

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

write(A, a);

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

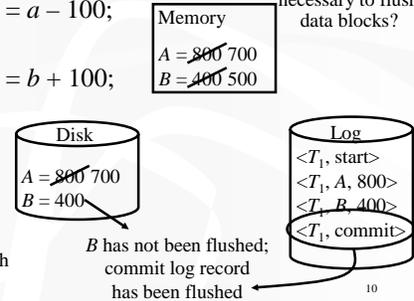
write(B, b);

output(A);

output(B);

System crash

When is it necessary to flush data blocks?



10

One technicality

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

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

write(A, a);

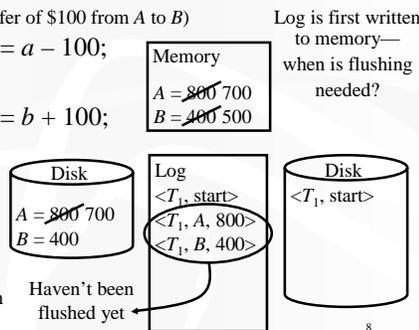
read(B, b); $b = b + 100$;

write(B, b);

output(A);

output(B);

System crash



8

Force

- Recap of the situation to be avoided
 - T_1 has committed (the log says so)
 - Not all effects of T_1 have been flushed disk
 - Because there is no redo information in the log, we cannot redo the rest of T_1
 - So perhaps we should try redo logging?
- Solution: force
 - Before the commit record of a transaction is flushed to log, all writes of this transaction must be reflected on disk

11

WAL

- Recap of the situation to be avoided
 - T_1 has not completed yet
 - A is modified on disk already
 - But there is no log record for A
 - Cannot undo the modification of A!
- Solution: WAL (Write-Ahead Logging)
 - Before any database object X is modified on disk, the log record pertaining to X must be flushed

9

Undo logging rules

- For every write, generate undo log record containing the old value being overwritten
 - $\langle T_i, X, \text{old_value_of_X} \rangle$
 - Typically (assuming physical logging)
 - T_i : transaction id
 - X: physical address of X (block id, offset)
 - old_value_of_X : bits
- WAL
- Force

12

Recovery with an undo log

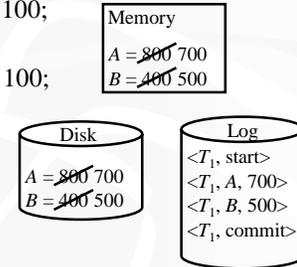
- Identify U , the set of active transactions at time of crash
 - Log contains $\langle T, \text{start} \rangle$, but neither $\langle T, \text{commit} \rangle$ nor $\langle T, \text{abort} \rangle$
- Process log backward Why?
 - For each $\langle T, X, \text{old_value} \rangle$ where T is in U , issue $\text{write}(X, \text{old_value})$ Why?
- For each T in U , append $\langle T, \text{abort} \rangle$ to the end of the log

13

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);
output(A);
output(B);
```



16

Additional issues with undo logging

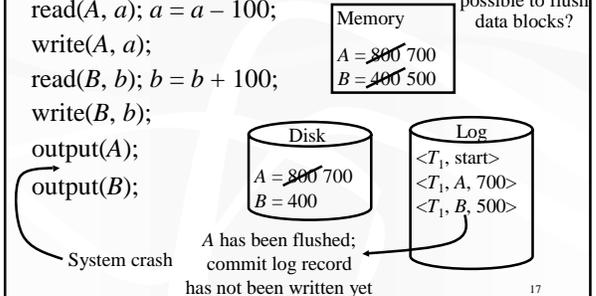
- Failure during recovery?
 - No problem, run recovery procedure again
 - Undo is idempotent!
- Can you truncate log?
 - Yes, after a successful recovery
 - Or, truncate any prefix that contain no log records for active transactions

14

One technicality

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);
output(A);
output(B);
```



17

Redo logging

- Basic idea
 - Every time you modify something on disk, record its new value (which you are writing)
 - If system crashes, redo the writes of committed transactions and ignore those that did not commit

15

No steal

- Recap of the situation to be avoided
 - T_1 has not completed yet
 - A is modified on disk already
 - There is a log record for A (i.e., WAL is followed)
 - Because there is no undo information in that log record, we cannot undo the modification of A !
 - Maybe undo/redo combined?
- Solution: no steal
 - Writes can be flushed only at commit time
 - Requires keeping all dirty blocks in memory—other transactions cannot steal any memory blocks

18

Redo logging rules

- For every write, generate redo log record containing the new value being written
 $\langle T, X, new_value_of_X \rangle$
- Do not modify any database objects on disk before you have flushed all log records for this transaction (including the commit record)
 - That is, WAL and no steal

19

Checkpointing

- Naïve approach:
 - Stop accepting new transactions (lame!)
 - Finish all active transactions
 - Take a database dump
 - Now safe to truncate the redo log
- Fuzzy checkpointing
 - Example later

22

Recovery with a redo log

- Identify C , the set of all committed transactions (those with commit log record)
- Process log **forward** Why?
 - For each $\langle T, X, new_value \rangle$ where T is in C , issue $write(X, new_value)$ Why is $output(X)$ unnecessary here?
- For each incomplete transaction T (with neither commit nor abort log record), append $\langle T, abort \rangle$ to the end of the log

20

Summary of redo and undo logging

- Undo logging—immediate write
 - Force
 - Excessive disk I/Os
 - Imagine many small transactions updating the same block!
- Redo logging—deferred write
 - No steal
 - High memory requirement
 - Imagine a big transaction updating many blocks

23

Additional issues with redo logging

- Failure during recovery?
 - No problem—redo is idempotent!
- Extremely slow recovery process!
 - I transferred the balance last year...
- Can you truncate log?
 - No, unless...

21

Logging taxonomy

Assuming each transaction modifies just one block and locking is at the block level

| | | |
|----------|--------------|-------------------|
| | no steal | steal |
| force | no logging! | undo logging |
| no force | redo logging | undo/redo logging |

Next!

24

Undo/redo logging

- Log both old and new values
 $\langle T_i, X, \text{old_value_of_}X, \text{new_value_of_}X \rangle$
- WAL
- Steal: If chosen for replacement, modified memory blocks can be flushed to disk anytime
- No-force: When a transaction commits, modified memory blocks are not forced to disk
- Buffer manager has complete freedom!

25

Recovery: analysis and redo phase

- Need to determine U , the set of active transactions at time of crash
- Scan log backward to find the last end-checkpoint record and follow the pointer to find the corresponding $\langle \text{start-checkpoint } S \rangle$
- Initially, let U be S
- Scan forward from that start-checkpoint to end of the log
 - For a log record $\langle T, \text{start} \rangle$, add T to U
 - For a log record $\langle T, \text{commit} \mid \text{abort} \rangle$, remove T from U
 - For a log record $\langle T, X, \text{old}, \text{new} \rangle$, issue $\text{write}(X, \text{new})$
 - Repeats history!

28

Undo/redo logging example

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

$\text{read}(A, a); a = a - 100;$

$\text{write}(A, a);$

$\text{read}(B, b); b = b + 100;$

$\text{write}(B, b);$

No output operations here—they are up to the buffer manager!

Anytime after corresponding log records are flushed

| Memory |
|--------------------------------------|
| $A = 800, 700$ |
| $B = 400, 500$ |

| Disk |
|--------------------------------------|
| $A = 800, 700$ |
| $B = 400, 500$ |

| Log |
|--------------------------------------|
| $\langle T_1, \text{start} \rangle$ |
| $\langle T_1, A, 800, 700 \rangle$ |
| $\langle T_1, B, 400, 500 \rangle$ |
| $\langle T_1, \text{commit} \rangle$ |

- So when is T_1 really committed?
 - When its commit log record is flushed to disk

26

Recovery: 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 (part of the repeating-history paradigm)
 - Log $\langle T, \text{abort} \rangle$ when all effects of T have been undone
- An optimization
 - Each log record stores a pointer to the previous log record for the same transaction; follow the pointer chain during undo
- Is it possible that undo overwrites the effect of a committed transaction?
 - Not if strict 2PL!

29

Fuzzy checkpointing

- Determine S , the set of currently active transactions, and log $\langle \text{begin-checkpoint } S \rangle$
- Flush all modified memory blocks at your leisure
 - Regardless whether they are written by committed or uncommitted transactions (but do follow WAL)
- Log $\langle \text{end-checkpoint } \text{begin-checkpoint_location} \rangle$
- Between begin and end, continue processing old and new transactions

27

Physical versus logical logging

- Physical logging (what we have assumed so far)
 - Log before and after images of data
- Logical logging
 - Log operations (e.g., insert a row into a table)
 - Smaller log records
 - An insertion could cause rearrangement of things on disk
 - Or trigger hundreds of other events
 - Sometimes necessary
 - Assume row-level rather than page(block)-level locking
 - Data might have moved to another block at time of undo!
 - Much harder to make redo/undo idempotent

30

Selective redo?

- Possible optimization for our recovery procedure:
 - Selectively redo only committed transactions
 - Lots of algorithms do it (some even undo before redo)
- What is the catch?
 - $T_1.op_1, T_2.op_1, T_1.op_2 (T_1.commit)$
 - Repeating history: $T_1.op_1, T_2.op_1, T_1.op_2, \text{undo}(T_2.op_1)$
 - Exactly the same as normal transaction abort
 - Selective redo: $T_1.op_1, T_1.op_2, \text{undo}(T_2.op_1)$
 - What if $T_2.op_1$ produced some side effects that $T_1.op_2$ relies on?
 - Not possible with page-level locking and physical logging
 - In general hard to guarantee

31

ARIES

- Same basic ideas: steal, no force, WAL
- Three phases: analysis, redo, undo
 - Repeats history
- CLR (Compensation Log Record) for transaction aborts
- More efficient than our simple algorithm
 - Redo/undo on an object is only performed when necessary
 - Each disk block records the last writer
 - Can take advantage of a partial checkpoint
 - Recovery can start from any start-checkpoint, not necessarily one that corresponds to an end-checkpoint

32