

Lecture 17: Reliable Storage

Zhong Shao Dept. of Computer Science Yale University



File system reliability

- What can happen if disk loses power or machine software crashes?
 - Some operations in progress may complete
 - Some operations in progress may be lost
 - Overwrite of a block may only partially complete
- File system wants durability (as a minimum!)
 - Data previously stored can be retrieved (maybe after some recovery step), regardless of failure







- Moving a file between directories:
 - * Delete file from old directory
 - * Add file to new directory
- Create new file
 - * Allocate space on disk for header, data
 - * Write new header to disk
 - * Add the new file to directory

What if there is a crash in the middle? Even with write-through it can still have problems







Transaction is a group of operations

- Atomic: operations appear to happen as a group, or not at all (at logical level)
 - * At physical level, only single disk/flash write is atomic
- Durable: operations that complete stay completed
 - ${}^{\star}\,$ Future failures do not corrupt previously stored data
- Isolation: other transactions do not see results of earlier transactions until they are committed
- Consistency: sequential memory model

7

Reliability approach #1: careful ordering

- Sequence operations in a specific order
 - Careful design to allow sequence to be interrupted safely
- Post-crash recovery
 - Read data structures to see if there were any operations in progress
 - Clean up/finish as needed
- Approach taken in FAT, FFS (fsck), and many app-level recovery schemes (e.g., Word)





















Normal operation:

- Write name of each open file to app folder
- Write changes to backup file
- Rename backup file to be file (atomic operation provided by file system)
- Delete list in app folder on clean shutdown

Recovery:

- On startup, see if any files were left open
- If so, look for backup file
- If so, ask user to compare versions

17

<section-header> Careful ordering Pros Works with minimal support in the disk drive Works for most multi-step operations Cons Can require time-consuming recovery after a failure Difficult to reduce every operation to a safely interruptible sequence of writes Difficult to achieve consistency when multiple operations occur concurrently















- Instead of modifying data structures on disk directly, write changes to a journal/log
 - Intention list: set of changes we intend to make
 - Log/Journal is append-only
- Once changes are on log, safe to apply changes to data structures on disk
 - Recovery can read log to see what changes were intended
- Once changes are copied, safe to remove log





- Key idea: fix problem of how you make multiple updates to disk, by turning multiple updates into a single disk write!
- Example: money transfer from account x to account y:

Begin transaction x = x + 1 y = y - 1 Commit

- Keep "redo" log on disk of all changes in transaction.
 - A log is like a journal, never erased, record of everything you've done
 - Once both changes are on log, transactions are committed.
 - Then can "write behind" changes to disk --- if crash after commit, replay log to make sure updates get to disk















Redo logging

- ♦ Prepare
 - Write all changes (in transaction) to log
- Commit
 - Single disk write to make transaction durable
- ♦ Redo
 - Copy changes to disk
- Garbage collection
 - Reclaim space in log

Recovery

- Read log
- Redo any operations for committed transactions
- Garbage collect log





Transaction isolation	
Process A	Process B
move file from x to y mv x/file y/	grep across x and y grep x/* y/* > log
	What if grep starts after changes are logged, but before commit?

Two-phase locking

- Don't allow "unlock" before commit.
- First phase: only allowed to acquire locks (this avoids deadlock concerns).
- Second phase: all unlocks happen at commit
- Thread B can't see any of A's changes, until A commits and releases locks. This provides serializability.

Transaction isolation	
Process A	Process B
Lock x, y move file from x to y mv x/file y/ Commit and release x,y	Lock x, y, log grep across x and y grep x/* y/* > log Commit and release x, y, log Grep occurs either before or after move

Serializability

- With two phase locking and redo logging, transactions appear to occur in a sequential order (serializability)
 - Either: grep then move or move then grep
- Other implementations can also provide serializability
 - Optimistic concurrency control: abort any transaction that would conflict with serializability



Question

- Do we need the copy back?
 - What if update in place is very expensive?
 - Ex: flash storage, RAID



Storage availability

- Storage reliability: data fetched is what you stored
 Transactions, redo logging, etc.
- Storage availability: data is there when you want it
 - More disks => higher probability of some disk failing
 - Data available ~ Prob(disk working)^k
 - * If failures are independent and data is spread across k disks
 - For large k, probability system works -> 0

RAID update

- Mirroring
 - Write every mirror
- RAID-5: to write one block
 - Read old data block
 - Read old parity block
 - Write new data block
 - Write new parity block
 - * Old data xor old parity xor new data
- RAID-5: to write entire stripe
 - Write data blocks and parity

- Disk devices can lose data
 - One sector per 10^15 bits read
 - Causes:
 - * Physical wear
 - * Repeated writes to nearby tracks
- What impact does this have on RAID recovery?

51

