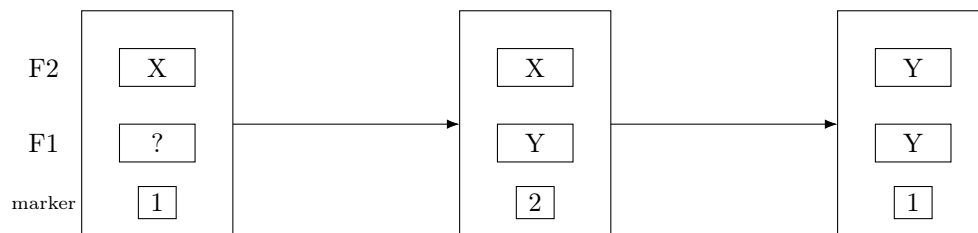


To discuss risks like this, we need to define a few key terms;

- error := a mistake by the designer or user
- fault := a latent problem in design
- failure := the problem caused when the fault trap springs
- durability := the system's ability for data to survive limited failure
- atomicity
 1. *All-or-Nothing Atomicity*
 - From the point of view of a function's invoker, the sequence either:
 - (a) Completes
 - (b) Aborts such that it appears the action was never started (back out)
 - How can we give the READ function AoN atomicity?
 - (a) *Blocking Read*: wait and sets return address before READ
 - (b) *Non-Blocking Read*: kernel returns if stream is empty
 - (c) Non-Atomic: READ waits and blocks until char is delivered
 2. *Before-or-After Atomicity*
 - From the view of a function's invokers, the result is the same as if the actions occurred completely before or completely after one another
 - If two actions have before or after atomicity, they are *serializable* := there exists some serial order of those concurrent actions that would, if followed, lead to the same ending state
 3. *Sequential/external time consistency*
 - If it appears to the outside that the events occurred in a certain order, the correct result is as if they were executed in this orders.
 - This is only required in some cases.

So take a test editor, say EMACS, and assume it is writing to blocks; what happens if the power goes out while we are overwriting our data? This leads us to our Golden Rule of Atomicity: never delete your only copy.



Our first attempt at atomicity: we create a new file and write the new data there; on completion we switch which is the active file.

- But what if we lost power mid-swap? We wouldn't know which is the active data?
We can solve this in a probabilistic sense with checksums post fail
- But what if it a block write isn't atomic? Write from A to B is 3 stages: $A \rightarrow ? \rightarrow B$
We can make an atomic write using 3 blocks!

	File Data Contents						
F1	A	?	B	B	B	B	B
F2	A	A	A	?	B	B	B
F3	A	A	A	A	A	?	B
	time →						

Algorithm:

1. write 3 in series
2. use best 2/3 to choose
3. choose block 0 if they all differ

Since we have atomic block write, we can discuss file system robustness on a larger scale:

We use a Lampion-Sturgis failure model:

1. storage writes may fail
2. storage writes may corrupt other blocks
3. storage blocks may decay spontaneously

4. a read can detect a bad block (using checksums)
5. errors are rare
6. reports can be done in time

We also need to establish a few file system robustness invariants:

1. every block serves at most one purpose — failure allows a program to overwrite another program’s data
2. all referenced blocks are properly initialized — failure allows an uninitialized block to look like a pointer to a random block
3. all referenced blocks are marked as “used” — failure allows data to be overwritten
4. all non-referenced blocks are marked as “free” — failure leaves an unused block marked as used → data leak

The failure of most of these would mean the loss of data—except number 4, so our failure hinges on number 4.

Take a risky operation. . . say `rename(“d/a”, “e/b”)`

- `fsck()` will catch our errors, but it is insanely slow!
- It prioritizes the inode data and moves unreferenced inodes to a lost & found
- File permissions are set to kernel only

This could be done in one of a few ways:

BAD Algorithm:

1. block a → RAM
2. block b → RAM
3. update blocks
4. block a → flash
5. block b → flash

GOOD Algorithm:

1. read block a into RAM
2. read inode into RAM
3. read block b into RAM
4. update blocks and increment link count
5. write inode to flash
6. write b to flash
7. write a to flash
8. link count -= 1 and write inode to flash

BAD:

- Failure between 4 & 5 would lose data!

GOOD:

- Instruction (5-6): only old link exists but `lc = 2`
- Instruction (6-7): the old and new copy exist and `lc = 2`
- Instruction (7-8): only new link exists but `lc = 2`

The good one is better because we would prefer to lose space over data!

We can now look at abstracting a single block write into multiple blocks:

0.0.1 COMMIT RECORDS

- document writes such that writes are atomic
- put commit records and all writes into a well recognized location

0.0.2 JOURNALING

Journal

...	A'	B'	CR	DR
-----	----	----	----	----

Cell Storage

			B
	A		

Algorithm:

1. Write A' to journal
2. Write B' to journal
3. Write CR to journal (BEGIN)
4. Copy A' to cell storage (CHANGEA)
5. Copy B' to cell storage (CHANGEB)
6. Write DR to journal (OUTCOME)

We then reorganize, and we're done.

Consequences of Failure:

Pre BEGIN: No effect.

Post BEGIN: Never initiated.

Post OUTCOME: Write Complete.

We keep cells in RAM and only copy to disk on write

- Wastes storage and will eventually run off disks
- + Solves many inconsistency problems
- + If mostly writing, avoids seeks

Since we can fail during reboot, our recovery strategy must also be *idempotent*, := execution one time is the same as executing any number of times.

We have to write each entry twice... what if we do a large write? Failure seems likely.

⇒ this is true, but most apps think between writes, so long writes are rare.

What do we do if a directory is replaced with a file & there is a crash?

- Linux uses the idea of “revoke records” to record a negative change in extv4
- Linux also writes

Is this viable for flash?

No — the performance benefit of this is that we ignore seek costs for write, but we know where we are going to write here.

We have two major journaling options: (the example uses write-ahead)

Write Ahead (ext4 option)

1. log changes in data
2. write changes in data
3. write the commit record
4. write the done record

RECOVERY:

- replay commit records

DOWNSIDE:

- the start of our replay can be hard to find

BENEFITS:

- more likely to save last action
- do not need to keep multiple versions

Write Behind

1. log old data values
2. write the commit record
3. write changes in data
4. write the done record

RECOVERY:

- undo the changes

DOWNSIDE:

- we must copy all of our data on each write

BENEFITS:

- old data cached ⇒ small pre-write benefit
- more conservative ⇒ more recovery
- no data searching ⇒ faster recovery

We need to be able to handle the failure of low level operations. We describe two interactions:

1. cascading aborts
 - if a low level operation fails, the high level one fails
 - very automate-able

- Example: Write-Ahead Protocol Cascading Abort:
 - (a) While logging planned writes, error occurs
 - (b) Abort record
 - (c) Send cascading aborts to higher level functions.
- 2. compensating actions
 - if a low level operation fails, the higher one makes up for it
 - very flexible
 - Example: Write-Ahead Protocol Compensating Action:
 - (a) Log planned writes
 - (b) Commit record
 - (c) While writing to disk, error occurs. After reboot compensating actions continue writing to cell memory using data written in the log.

There are many different types of corruption which can occur in a file system:

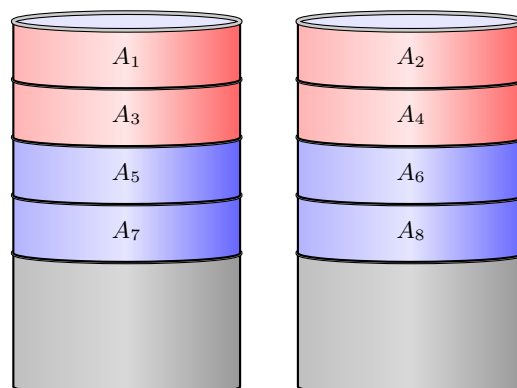
1. gamma rays can flip a bit
 - this happens roughly once a week
 - ECC memory with a parity bit can fix single flips and catch doubles
2. Drive Failure
 - this includes physical damage to the drive
 - companies give an annualized fail rate to give % fail chance for a year of operation
 - we can use SMART to check condition metrics; if we find a bad sector, we have replacement ones
3. User error
 - ex) `rm *` (remove all files) instead of `rm *.o` (remove all object files)
4. OER Errors
 - include configuration errors, which are the majority of errors
 - application and OS errors also qualify

We can catch and fix power failure with our old method. We can catch drive failure with a log structured file system, but to fix it we need to make a copy.

0.1 Redundant Array of Independent Disks (RAID)

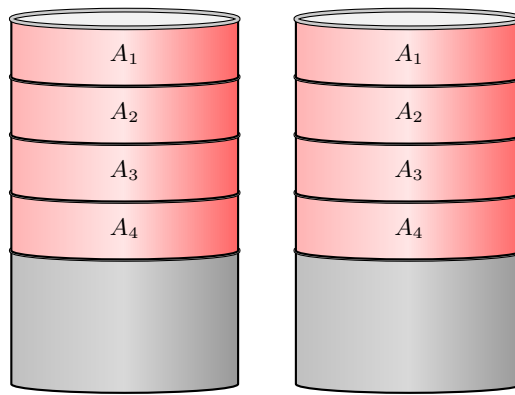
RAID was originally developed to save money by aggregating many small drives into a large one. It is now one of our most useful memory tools!

RAID-0



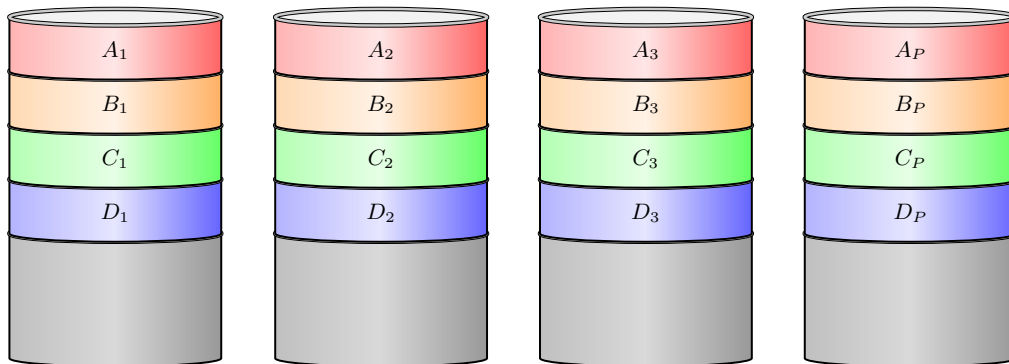
- We concatenate physical disks to form a large virtual one.
- We stripe the data to increase concurrency for read and write.

RAID-1



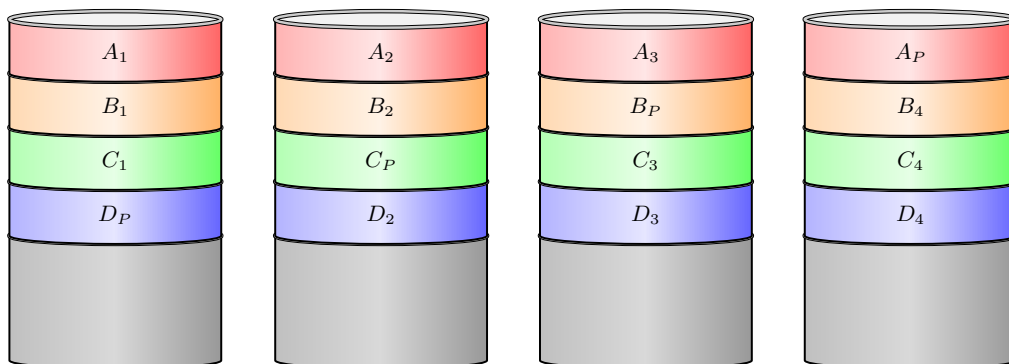
- We mirror half of our disks.
- We double write speed and half storage.
- We can now survive with half our data lost.

RAID-4



- We reserve one of our disks for *parity bits*. Parity bits are calculated with exclusive-or (\oplus)
- When a drive fails, we go into degraded mode;
 - We especially need to guarantee robustness here.
 - We can use a hot spare to replace a corrupt drive.
 - We can now add drives easily.
 - We can now recover from a drive loss.

RAID-5



- We stripe within a RAID 4 system.
- This is faster but more difficult to add drives.
- The book says these have almost entirely replaced RAID 4, but Eggert now seems to disagree.

All RAID uses a full-stop model; on detection of an error the operation stops. It can use checksums that include the data location