Funda : ✓

Techniques :
→ Local Nodes
→ Distributed Systems.

→ NFS.  (idempotency)

→ Leases
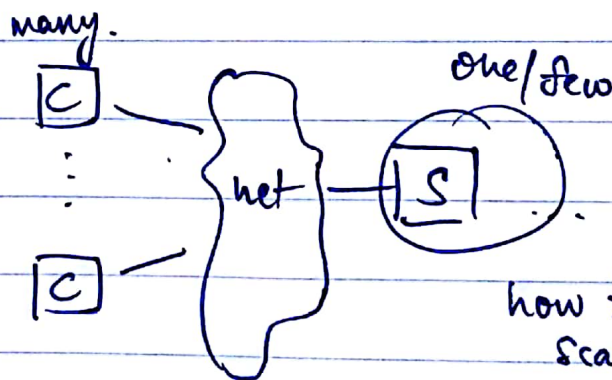concurrency
failure. ✓

* locks can't do the same
* includes time-lease.

NFS :  "Distributed"  in  1980's.
1990's

⇓

client, servers.

many.



one/few

how to scale?

1) scale out (add servers)

happened later (challenging).  ⟼  → hard → increases complexity.

2) scale up (make server more powerful).

↑CPU, ↑memory.

easy → fewer interactions.

Remzi's opinion.
Main Goal :

goals :-

=> transparency.
(performance ?
semantics ? ).

deleting open
~~and~~ files
and
writing to
them.

=> simplify failures handling (server).
fault into a
performance blip.

[ client behaviour. -> just try it
(idempotency) anything bad happens ]

$$A = \frac{MTTF}{MTTF + \overset{\frown}{MTTR}}$$

↓
reduce to ↑ availability.

[ unifies packet loss, server crash ] .

based upon :

why? open files.
handling

"statelessness".

local → no problem , enter system fails
c/s model → partial failures

**Protocol:**

             handle.

      ( file, ~~descriptor~~ )

          ↓

      ( vol #, inode #, gen # )

          ↑      ↑          ↖ useful why?

      which   which     ( across file

      file      file ?       delete

      sys ?             + reuse' inode # ).

**Ops:**

→ ~~read~~ read ( fh, offset, size ).

          =) data, error code.

→ write ( fh, [offset], size, data )

        local → os takes care of
                    offset of write.

→ lookup ( parent fh, name).

      =) fh of "name" in parent dir.

→ getattr ( fh )

      → stat() syscall about a
                    file.

## from protocol to FS :

not a one to one mapping

$\Rightarrow$ fs API $!=$ DFS protocol.
$\downarrow$
distributed

### APIs

open ( ) $\longrightarrow$ lookup

$$/x/y/z$$
$$\llcorner\uparrow\llcorner\uparrow\llcorner\uparrow$$
fh1 . fh2 .. fh3.

read ( ) $\longrightarrow$ read (s). $\left.\begin{array}{c}\\ \\\end{array}\right\}$ one or many.
write ( ) $\longrightarrow$ or write (s).

close ( ) $\longrightarrow$ ( N/A ).
only client-side

### How to handle failure ?

1) $c \overset{x}{\longrightarrow} \cdots s$   just retry

req loss

2) $c \longrightarrow s$   retry should be okay.

$x \overset{reply\ loss}{\longleftarrow}$   (idempotency)

$\overset{\downarrow (server)}{}$

server got smarter.   they cached that
they've done it
and sends
ACK.

3)  C $\longrightarrow$ S (crash)
    X

Ans : <u>retry</u> .

Be careful to retry :

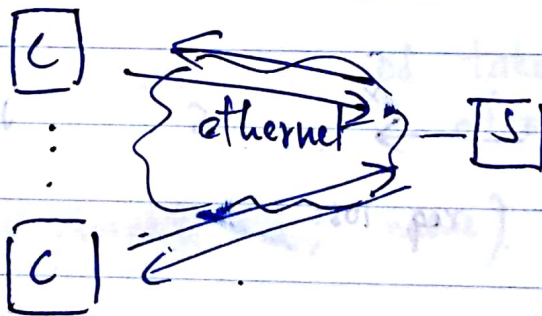$\rightarrow$  while (1) {
                 do :
                    if  success ( ) :
                        break
         }

too many retrys causing system to overload.

=) [<u>back</u>off  mechanisum  needed]
                     required.
         $\downarrow$
    wait  longer and longer after each
                 retry.

<u>Performance problems :</u>

[C]
 ⋮        ⊰ ethernet ⊱ —[S] .
[C]

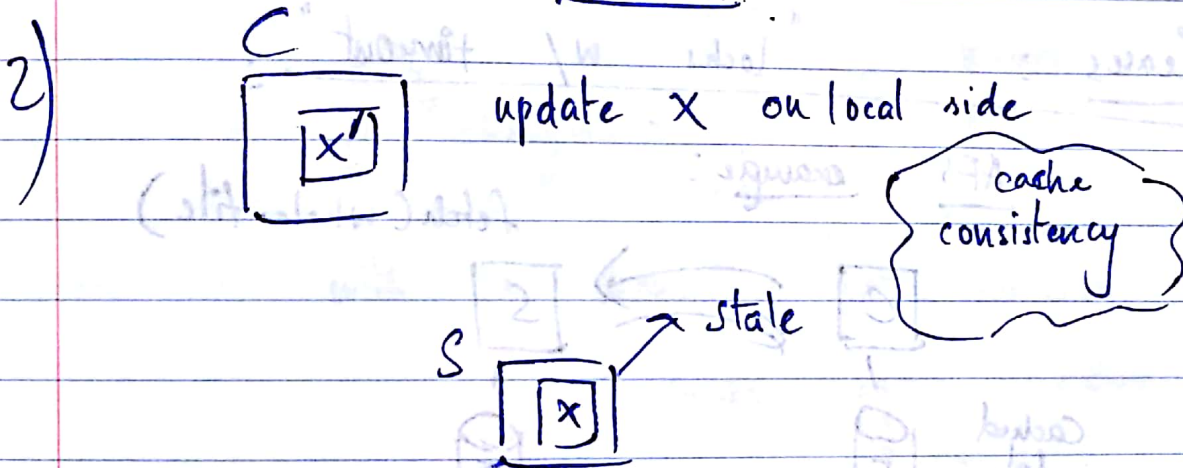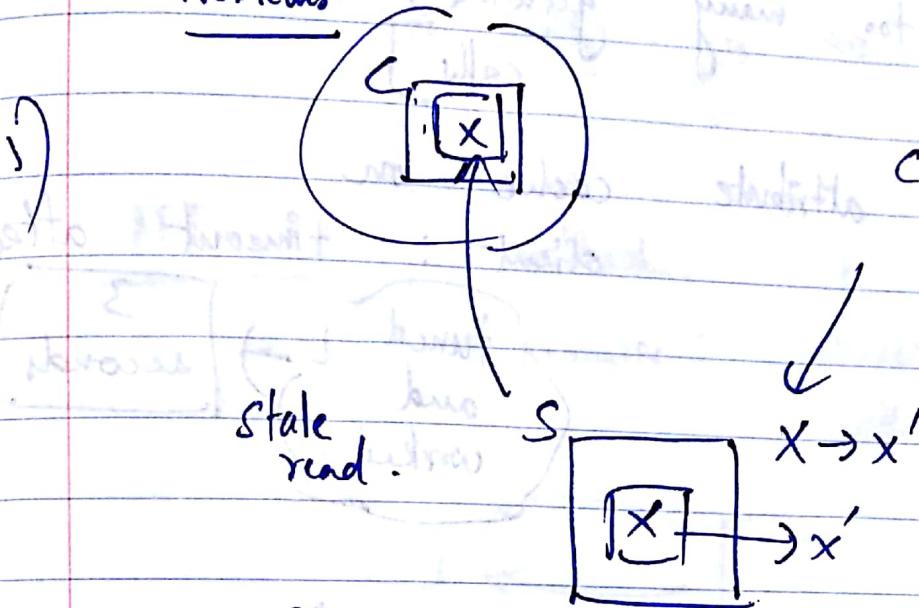NFS        read caching  to prevent server
                                  access everything.
       [client-side
              caching ] .

         + <u>write buffering</u> .

           send bulk requests from client
helpful?  1) one ~~pade~~ chunk of requests.
          2) redundant write or actions

## Problems:

1)

stale
read.

$S$  $x \to x'$

$\boxed{x} \to x'$

$C$

2)

$\boxed{x'}$   update x on local side

cache
consistency

$S$  $\boxed{x}$ → stale

## Solutions :

1) C before use of → S
   cached:
   ask if file changed ?      check
   get dttr ( )              before
                             use.

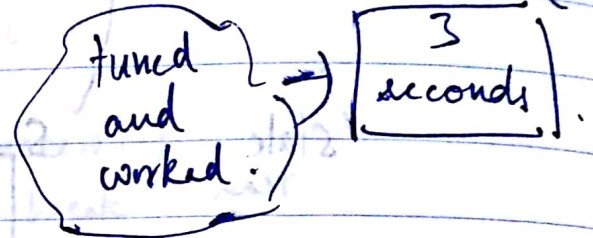2) flush-on-close :
   open
   [ write —
     ⋮    ]
   close        → force update
                   on
                   close . to server
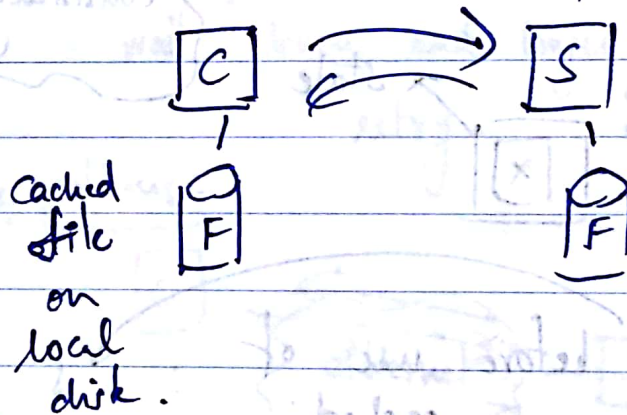
last problems :

[ too many getattr() calls ].

=) attribute cache on client : timeout after

tuned and worked := [ 3 seconds ].

Leases : "locks w/ timeout".

AFS example :

fetch ( whole file )



C ⇄ S

Cached file F on local disk.

Reads & Write on local.

on close :-, if file has changed flush to server
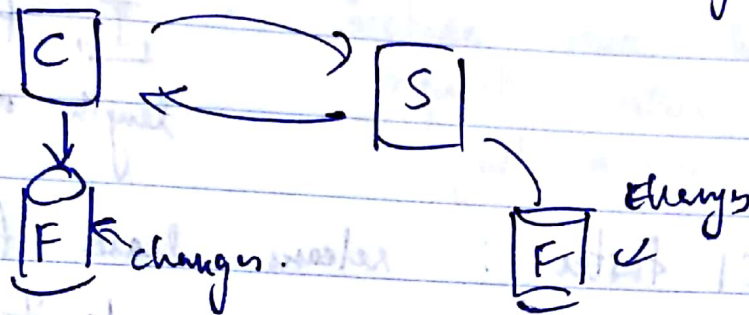
cache on mem or disk ?

↓
better
+ persistant (on crash)
+ bigger.

**AFS v1:** cache validation ?

Test Auth : is this file ~~cache~~ valid ?

**AFS v2:** callback.

promise : server notifies client with changes to file.



with a cache : not stateless anymore.

server crash handling is different.

=> complex.

client crash => how-to notify.

**Solution:**

**leases** rights with timeout.
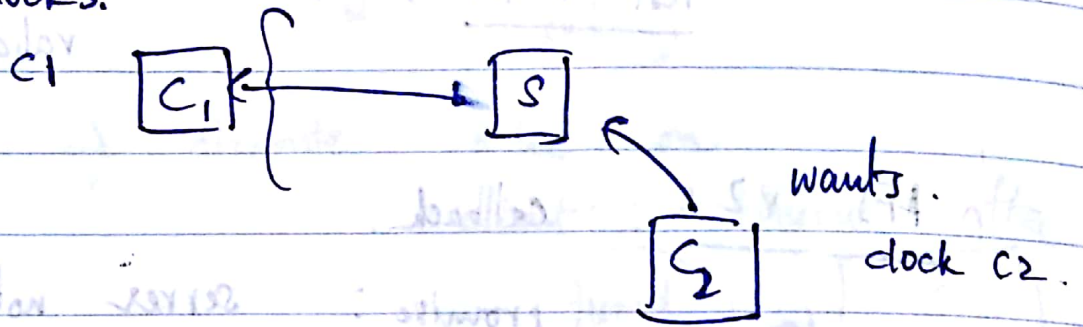
⇓

key to handling failure.

common :

acquire → use it.

(maybe renew periodically).

release ( or just let timeout).
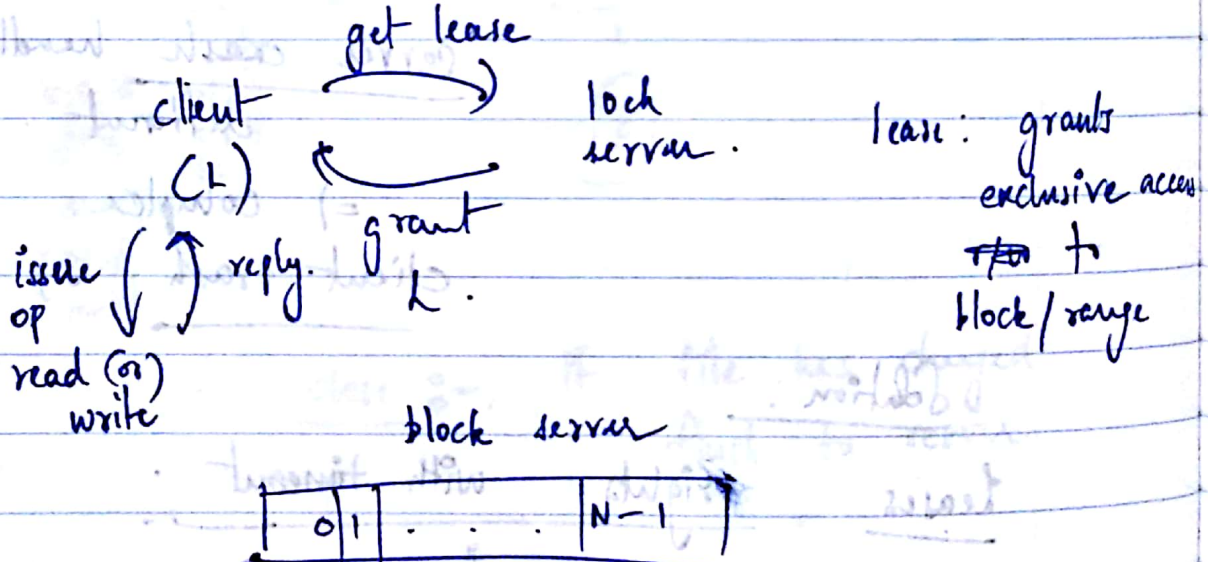
=) when lease safely grant to
   another client ?

clocks.

C1 [C₁] ← ─────── [S]

                              [C₂]  wants.
                                    clock C2.

1) server waits, wait how long?

$$I, + \Delta \quad \text{clock skew}$$
length of lease.  that
                  we
                  tolerate.

(Clock {
  C1 faster : releases lease faster
  C1 slow :  server grants it to C2 and
             C1 thinks it has lease.

              get lease
  client ────⌒────→  lock
  (L)    ←──⌒──      server.        lease: grants
  issue  )reply. grant                    exclusive access
  op   ↑         L.                       to
  read (or)                               block/range
  write         block server
         [ 0 1 │ . . . │ N-1 ]

         leases to make sure ops to
block server is serial.

**Problem:**

<u>Race condition</u>

{ client requests even after lease
  expires on server. and $s$ grants
  lease to another client.

<u>not serial</u>. anymore.

**how to handle?**

=) <u>fence</u> : include lease info in
            request. rather than
               just relying on timing.