

Welcome to the
Penultimate
lecture!

NFS

Shivaram Venkataraman
CS 537, Spring 2020

ADMINISTRIVIA

No SCLP DAYS

AEFIS feedback → May 1 ~ 35% currently

Optional project → Due Wed 10pm 70? like this to be

Final exam details → Check Piazza

↳ Canvas Quiz → Randomization

[No discussion this week!]

Two Parts

→ 2/5 ⇒ 2 extra credit points

AGENDA / LEARNING OUTCOMES

How to design a distributed file system that can survive partial failures?

What are consistency properties for such designs?

↳ Distributed systems

RECAP

TCP }
UDP } RPC

DISTRIBUTED FILE SYSTEMS

Local FS: processes on same machine access shared files

Network FS: processes on different machines access shared files in same way

Goals

Transparent access

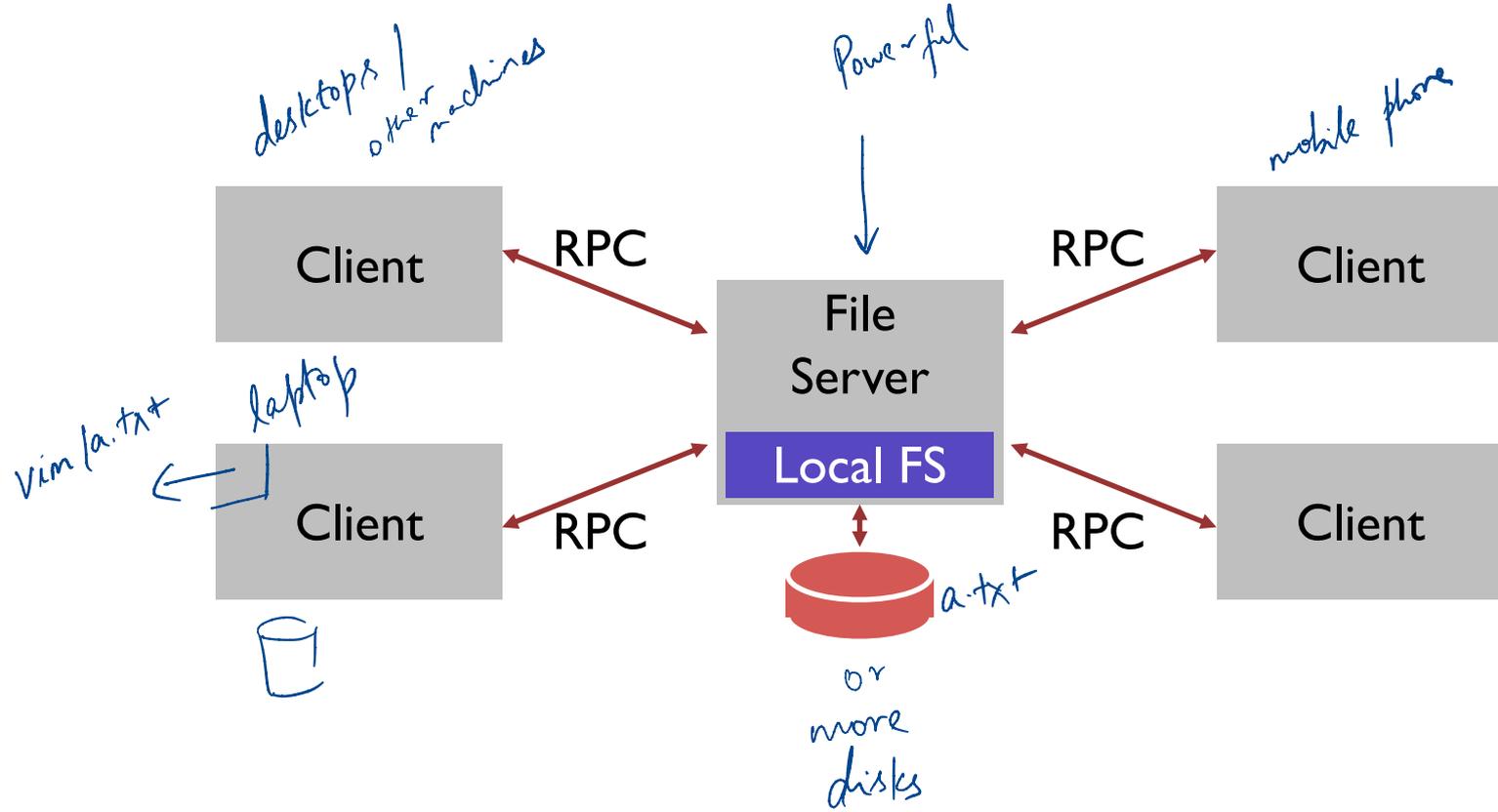
Fast + simple crash recovery

Reasonable performance?

process is
unaware that it
is using
Network FS

For server & client failure

NFS ARCHITECTURE



Naming!

STRATEGY 1

Attempt: Wrap regular UNIX system calls using RPC

`open()` on client calls `open()` on server

`open()` on server returns `fd` back to client

`read(fd)` on client calls `read(fd)` on server

`read(fd)` on server returns data back to client

```
int fd = open("foo", O_RDONLY);
```

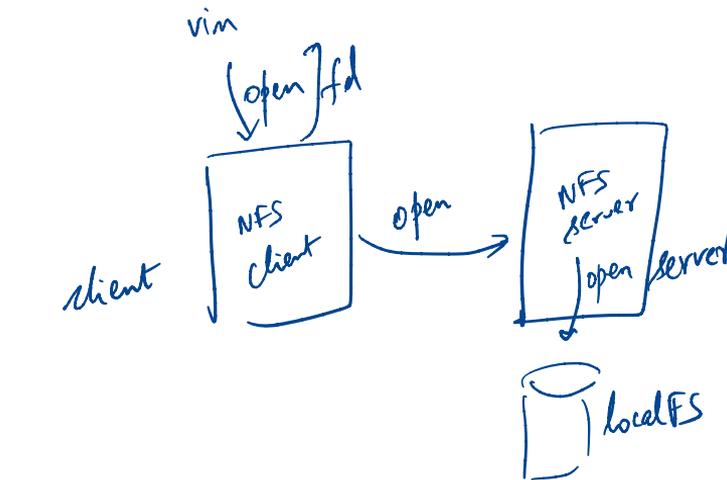
```
read(fd, buf, MAX);
```

...

```
read(fd, buf, MAX);
```

← Server crash!

server comes back up



*fd store inode num
and offset
⇒ FD table is lost
offset is lost*

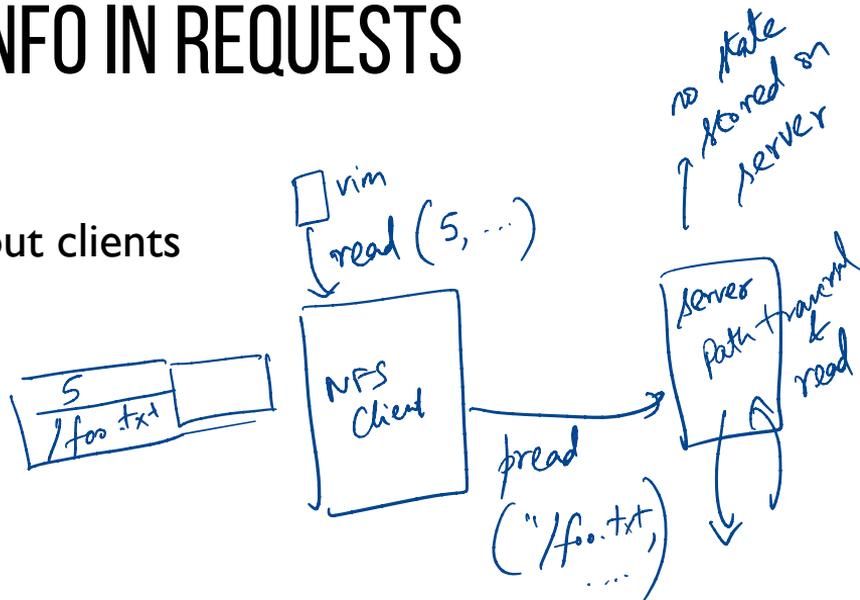
STRATEGY 2: PUT ALL INFO IN REQUESTS

“Stateless” protocol: server maintains no state about clients

Need API change. One possibility:

```
pread(char *path, buf, size, offset);
```

```
pwrite(char *path, buf, size, offset);
```



Specify path and offset each time. Server need not remember anything from clients.

Pros? Server can crash and reboot transparently to clients

Cons? Too many path lookups.

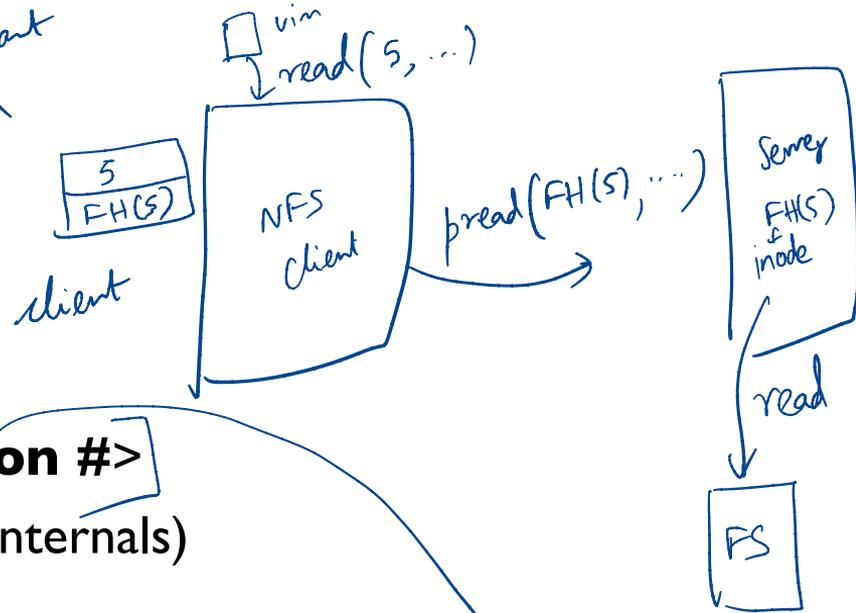
STRATEGY 3: FILE HANDLES

```
fh = open(char *path);  
pread(fh, buf, size, offset);  
pwrite(fh, buf, size, offset);
```

File Handle = \langle volume ID, inode #, generation # \rangle
Opaque to client (client should not interpret internals)

Stateless. We don't have path traversal on every pread, pwrite

make sure FH is not stale



Client

Server

fd = open("/foo", ...);
Send LOOKUP (rootdir FH, "foo")

RPC

Receive LOOKUP reply
allocate file desc in open file table
store foo's FH in table
store current file position (0)
return file descriptor to application

read(fd, buffer, MAX);
Index into open file table with fd
get NFS file handle (FH)
use current file position as offset
Send READ (FH, offset=0, count=MAX)

pread
RPC

Receive READ reply
update file position (+bytes read)
set current file position = MAX
return data/error code to app

message

Receive LOOKUP request
look for "foo" in root dir
return foo's FH + attributes

reply

cache FH



Receive READ request
use FH to get volume/inode num
read inode from disk (or cache)
compute block location (using offset)
read data from disk (or cache)
return data to client

block contents in reply

file handle for / is well known (mount)

home / shivaram / foo.txt

FH

Lookup (FH (/), "home")
Lookup (FH ("/home"), "shivaram")

similar to path traversal but is each message the lookup over network.

CAN NFS PROTOCOL INCLUDE APPEND?

```
fh = open(char *path);
```

```
pread(fh, buf, size, offset); ✓
```

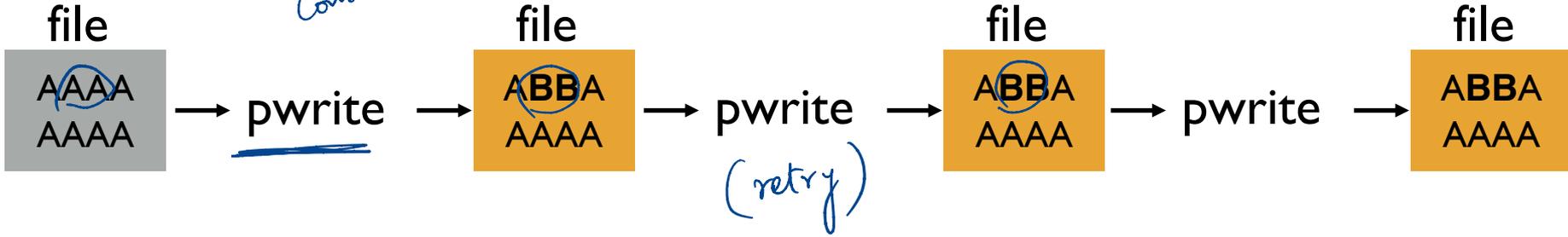
```
pwrite(fh, buf, size, offset); ✓
```

```
[ append(fh, buf, size); ]
```

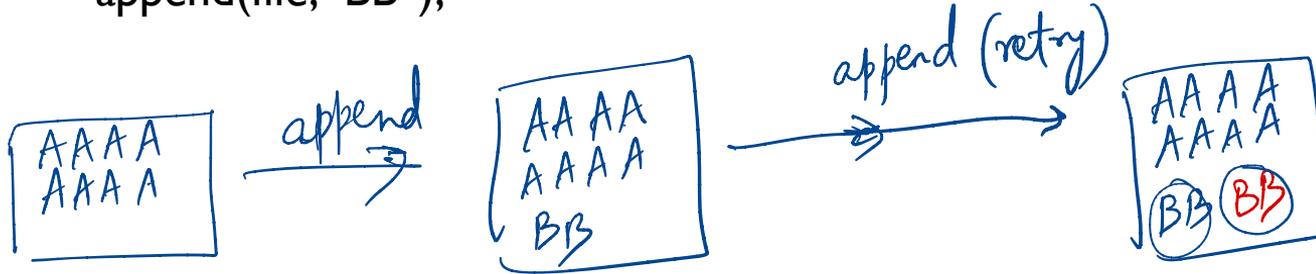
PWRITE VS APPEND

`pwrite(file, "BB", 2, 2);`

Contents
offset



`append(file, "BB");`



Client doesn't know if server crashed before/after doing the operation

IDEMPOTENT OPERATIONS

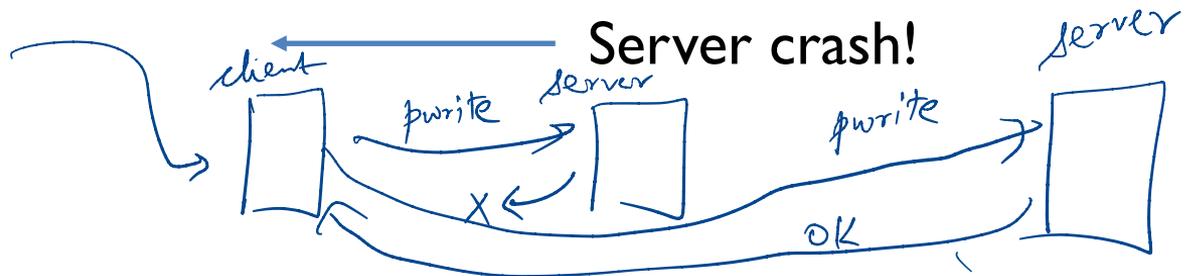
Solution: Design API so no harm to executing function more than once

If f() is idempotent, then:

→ f() has the same effect as f(); f(); ... f(); f()

Retrying functions is only safe if f is idempotent

```
int fd = open("foo", O_RDONLY);  
read(fd, buf, MAX);  
write(fd, buf, MAX);  
...
```



WHAT OPERATIONS ARE IDEMPOTENT?

Idempotent

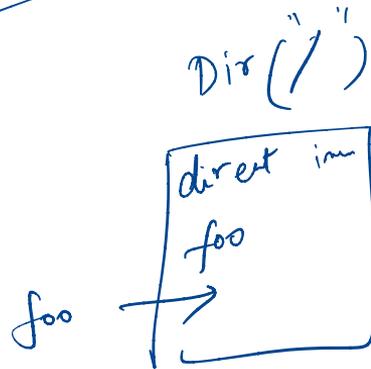
- any sort of read that doesn't change anything
- pwrite → specified offset

Not idempotent

- append

What about these?

- mkdir →
- creat →

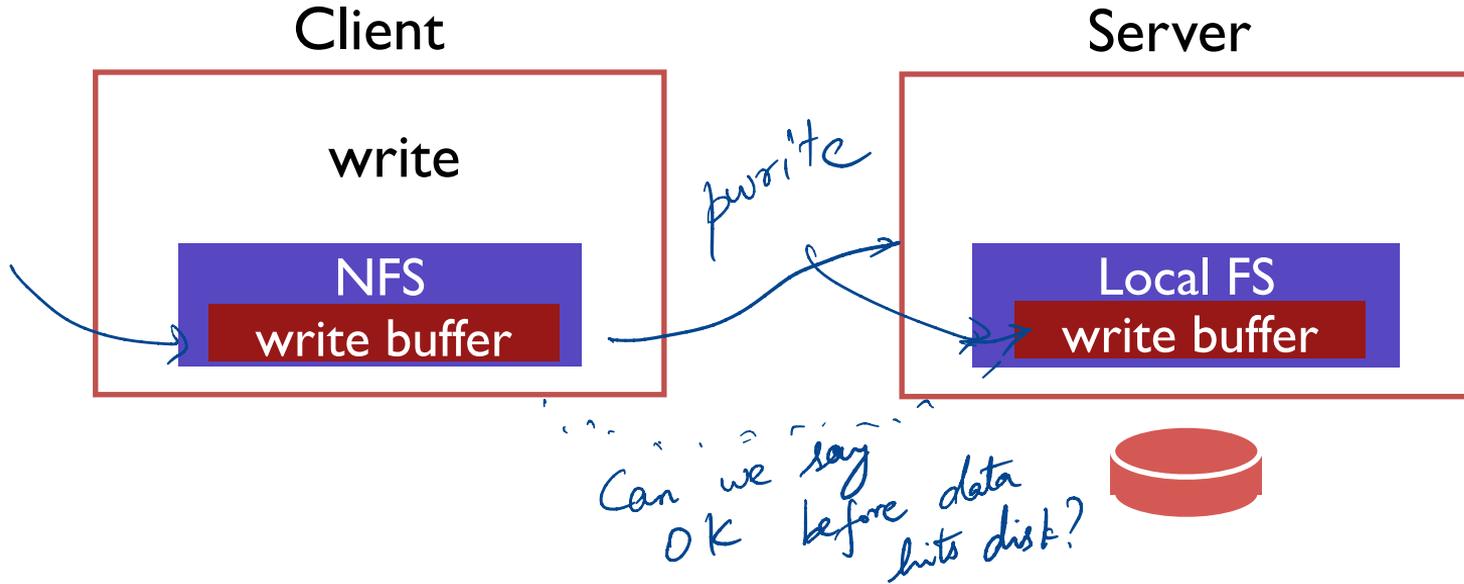


error: /foo already exists

```
rc
0 ← mkdir ("/foo")
1 ← mkdir ("/foo")
```

not part of NFS API, while pwrite is

WRITE BUFFERS



Server acknowledges write before write is pushed to disk;
What happens if server crashes?

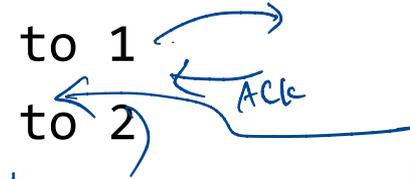
SERVER WRITE BUFFER LOST

client:

write A to 0

write B to 1

(write C to 2)



wait / retry
kill server
is back

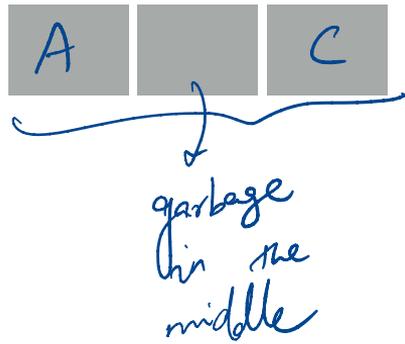
Crashes?



server mem:



server disk:



Block
0 1 2

server acknowledges write before write is pushed to disk

SERVER WRITE BUFFER LOST

Client:

write A to 0

write B to 1

write C to 2

← ABC

write X to 0

← XBC

write Y to 1

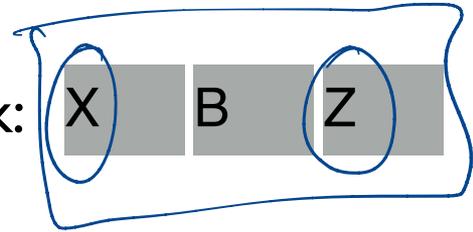
← XYC

write Z to 2

server mem:



server disk:

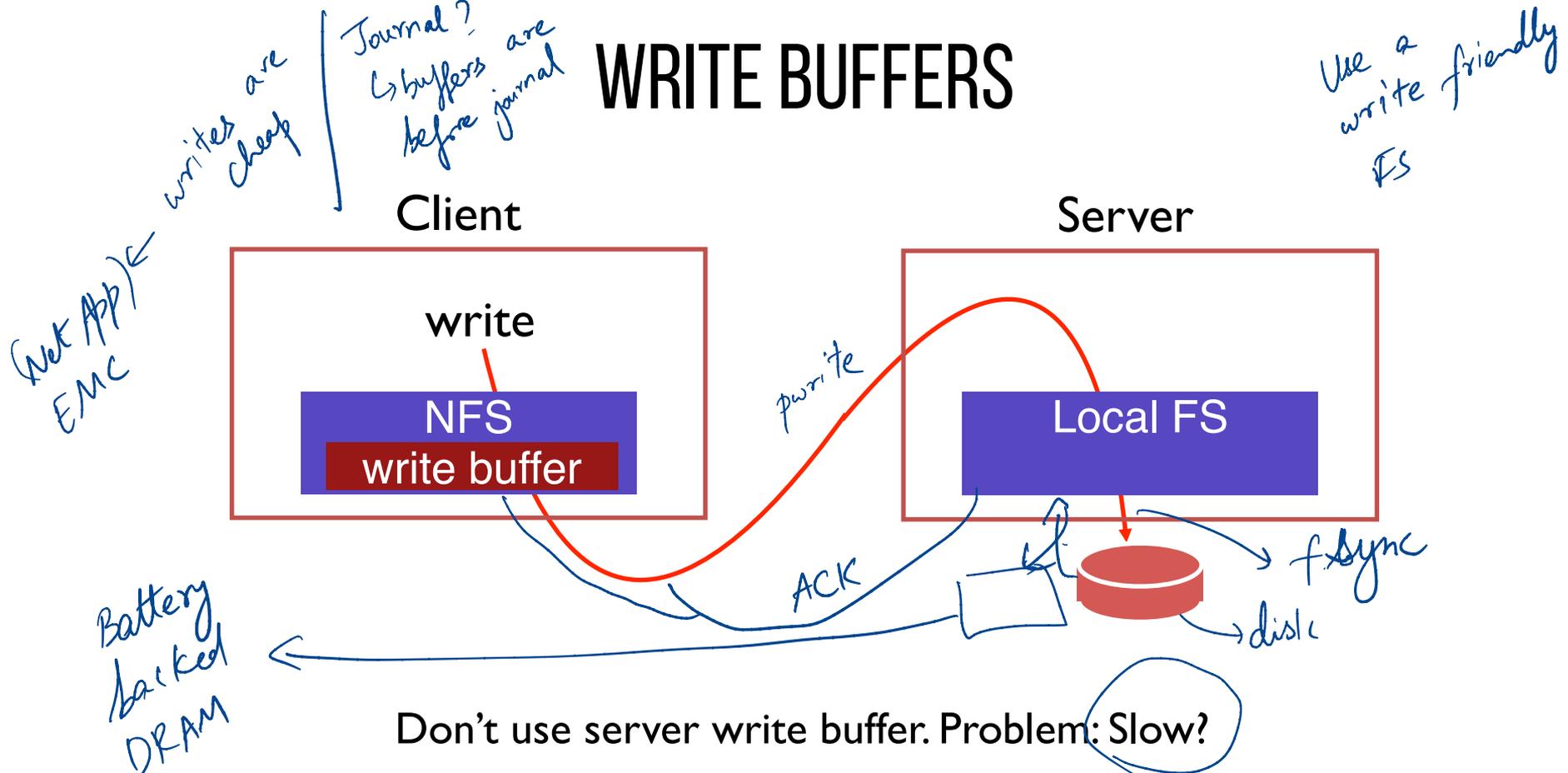


Problem:

No write failed, but disk state doesn't match any point in time

Solutions?

WRITE BUFFERS



Don't use server write buffer. Problem: Slow?

Use persistent write buffer (more expensive)

QUIZ 31

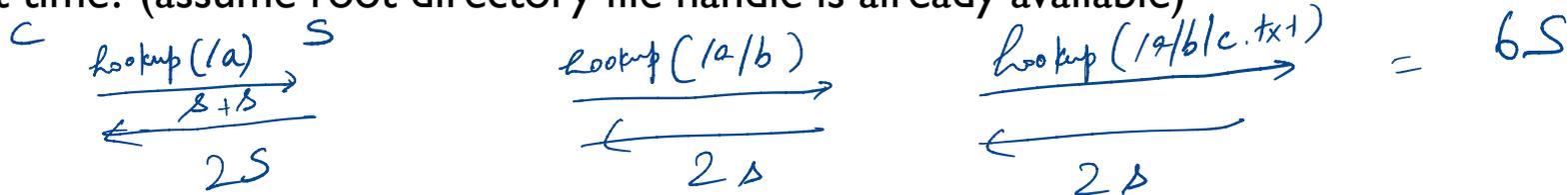
<https://tinyurl.com/cs537-sp20-quiz31>

CSL \rightarrow AFS



The only costs to worry about are network costs. Assume ("small") messages takes S units of time, whereas a "bigger" message (e.g., size of a block=4KB) takes B units. If a message is larger than 4KB, it takes longer (2B for 8KB).

1. How long does it take to open a 100-block (400 KB) file called /a/b/c.txt for the first time? (assume root directory file handle is already available)



2. How long does it take to read the whole file?



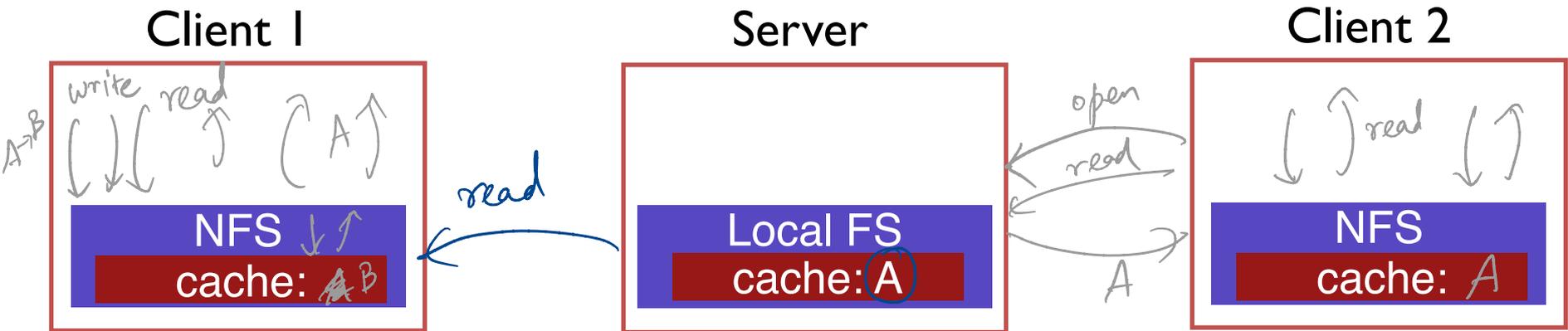
CACHE CONSISTENCY

NFS can cache data in three places:

- server memory
- client disk
- client memory

How to make sure all versions are in sync?

DISTRIBUTED CACHE



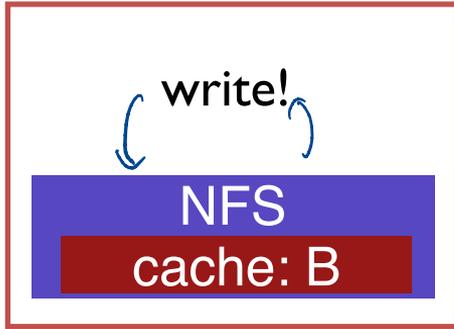
Using client cache
⇒ no network latency

① If ~~A~~ we have a write when do we flush that to server?

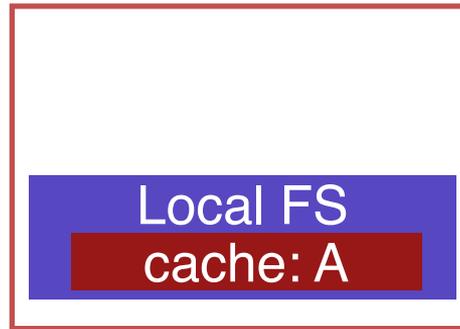
② How do we detect that cache Client 2 is STALE?

CACHE

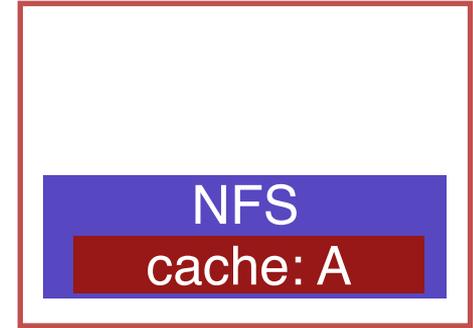
Client 1



Server



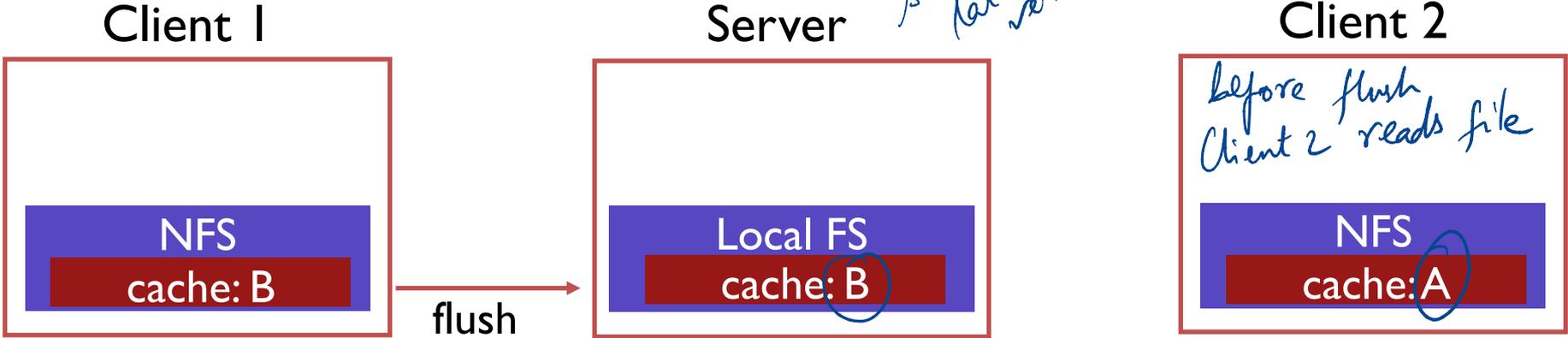
Client 2



“Update Visibility” problem: server doesn’t have latest version

What happens if Client 2 (or any other client) reads data?

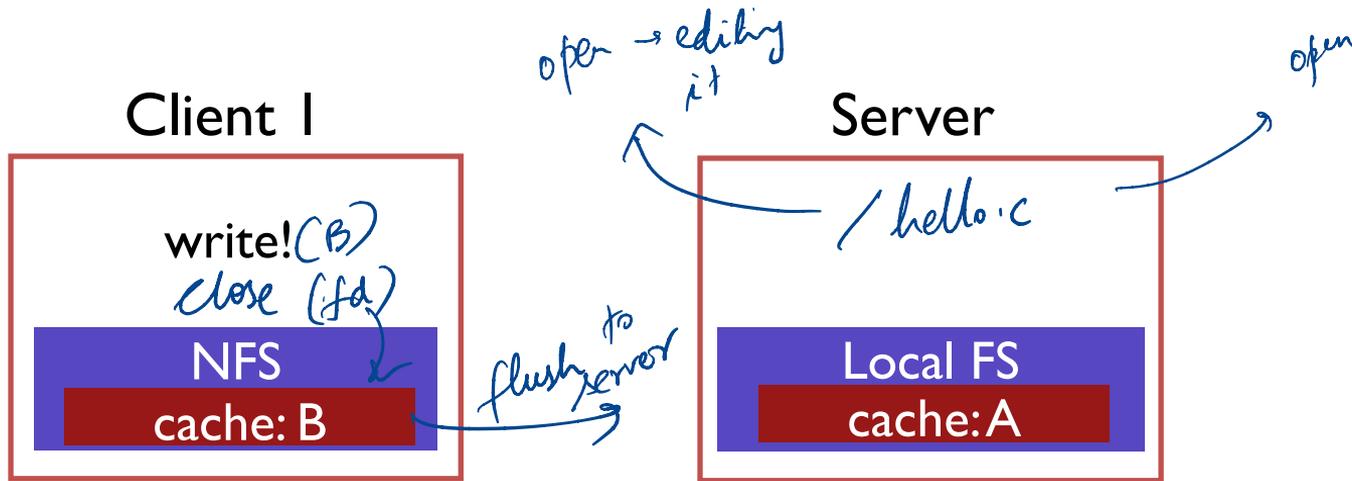
CACHE



“Stale Cache” problem: client 2 doesn't have latest version

What happens if Client 2 reads data?

PROBLEM 1: UPDATE VISIBILITY



When client buffers a write, how can server (and other clients) see update?

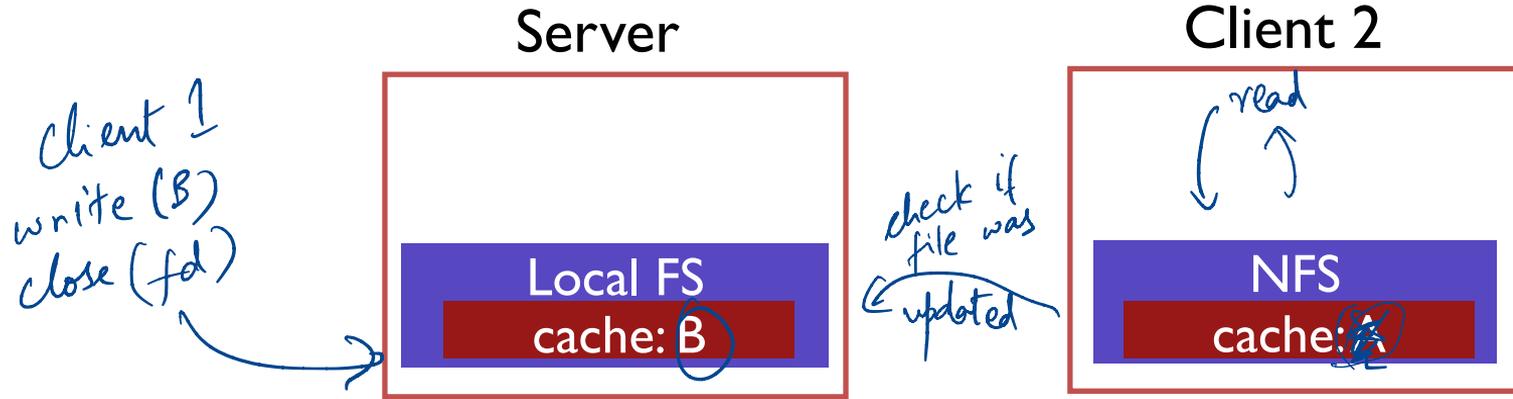
Client flushes cache entry to server

When should client perform flush?

NFS solution: flush on fd close

worst case guarantee

PROBLEM 2: STALE CACHE

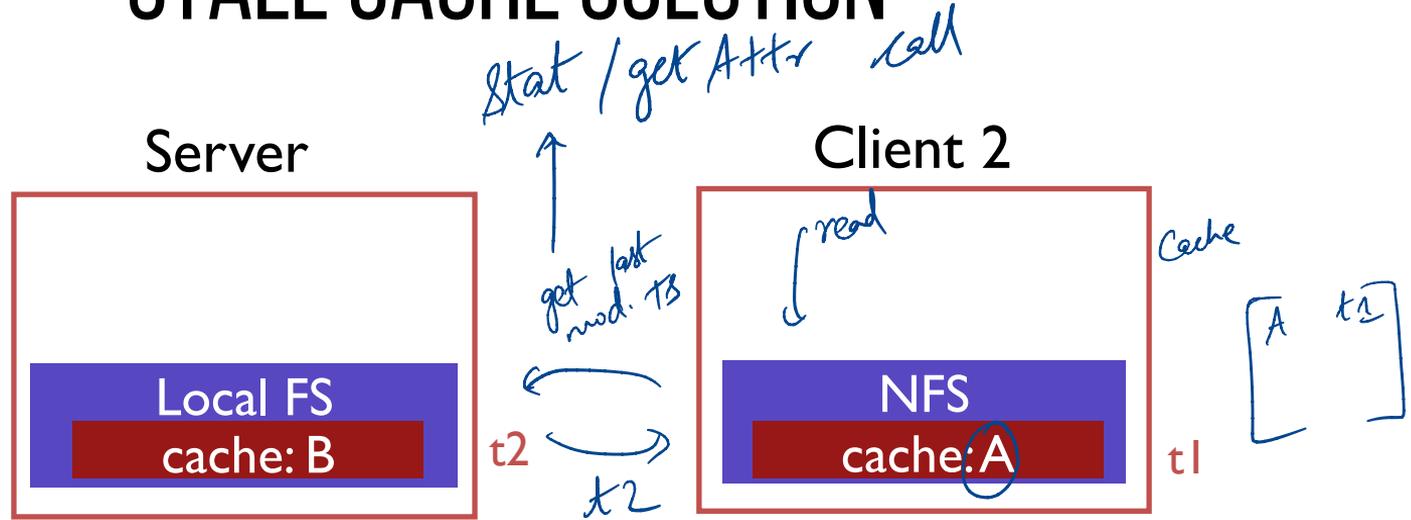


Problem: Client 2 has stale copy of data; how can it get the latest?

NFS solution:

- Clients recheck if cached copy is current before using data

STALE CACHE SOLUTION



Client cache records time when data block was fetched (t_1)

Before using data block, client does a STAT request to server

- get's last modified timestamp for this file (t_2) (not block...)
- compare to cache timestamp
- refetch data block if changed since timestamp ($t_2 > t_1$)

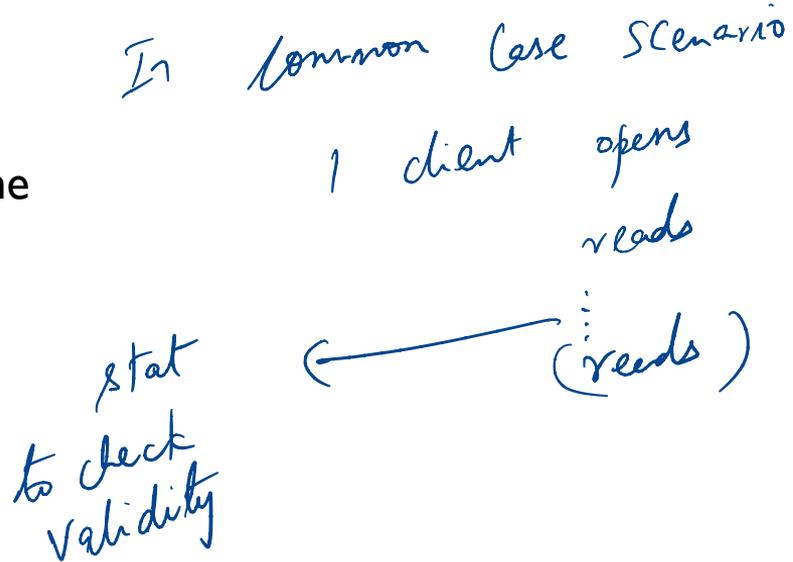
timestamp is at file granularity

MEASURE THEN BUILD

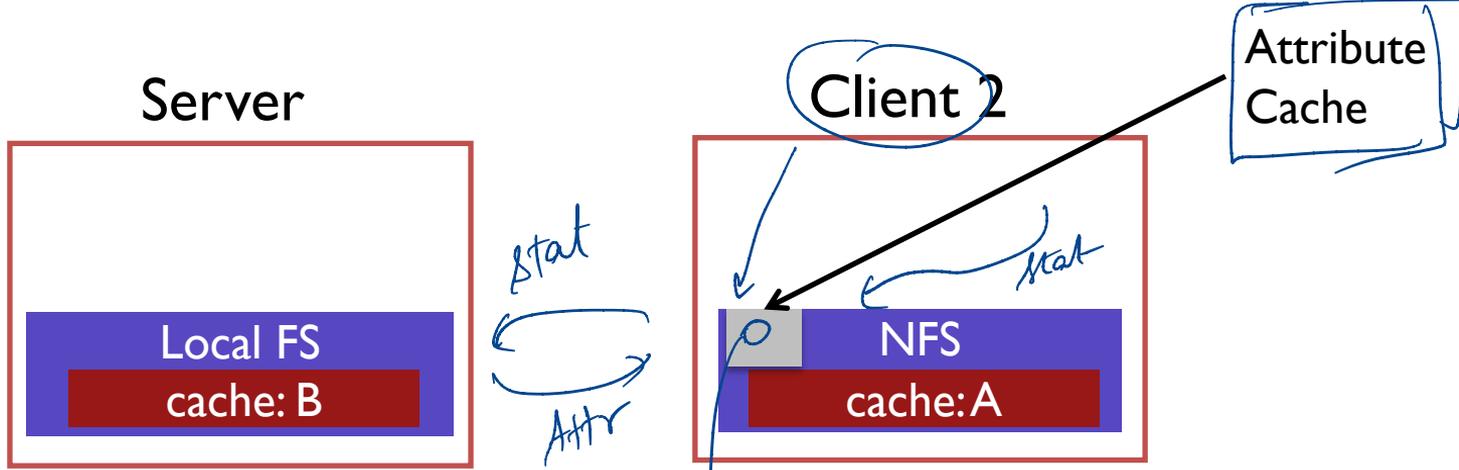
NFS developers found stat accounted for 90% of server requests

Why?

Because clients frequently recheck cache



REDUCING STAT CALLS



Solution: cache results of stat calls

Partial Solution:

Make stat cache entries expire after a given time (e.g., 3 seconds) (discard t2 at client 2)

What is the consequence?

cache reads could be stale for 3s!

if attr cache entry is older than 3s then call stat again

NFS SUMMARY

*Clients pulling
stat/ timestamps
from server*

NFS handles client and server crashes very well; robust APIs that are:

- stateless: servers don't remember clients
- idempotent: doing things twice never hurts

Caching and write buffering is harder, especially with crashes

Problems:

- Consistency model is odd (client may not see updates until 3s after file closed)
- Scalability limitations as more clients call stat() on server

NEXT STEPS

Next class: Review, Looking forward

Optional project due Wed

AEFIS feedback