Welsome back!

3 Marses to 801 DISTRIBUTED SYSTEMS, NFS

Shivaram Venkataraman CS 537, Spring 2020

ADMINISTRIVIA

Project 5: Due today! ____ Nopm ___ Slip days last project Optional project - Project Extre Credit No ship days No thip days No AEFIS feedback ~ 30 7 Discussion today: Optional project (short?) 4 30 mins or 100

Final Exam -> Cumulative

AGENDA / LEARNING OUTCOMES

What are some basic building blocks for systems that span across machines?

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

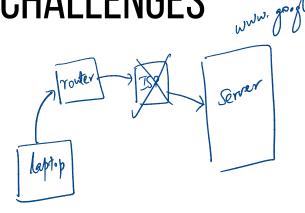
RECAP

DISTRIBUTED SYSTEMS: CHALLENGES

System failure: need to worry about partial failure

Communication failure: links unreliable

- bit errors
- packet loss
- node/link failure



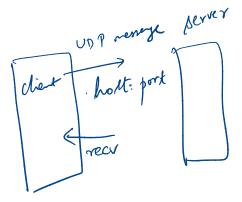
RAW MESSAGES: UDP

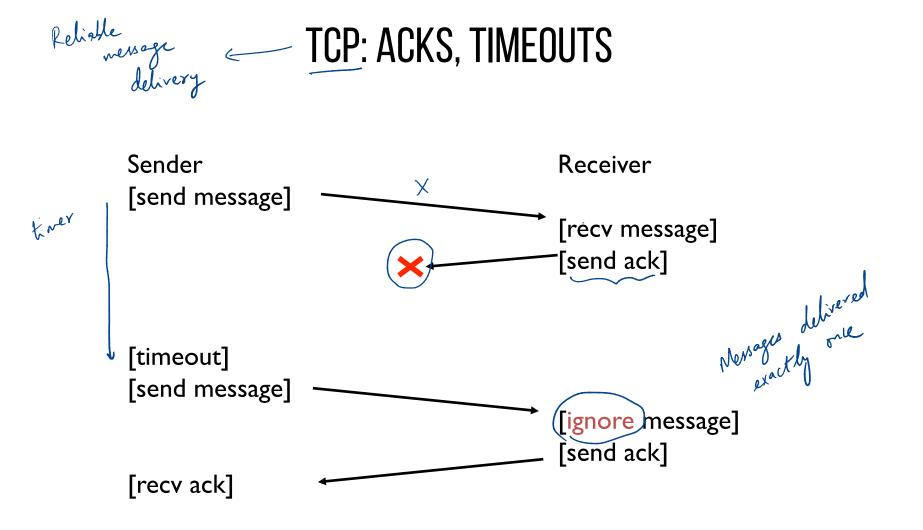
UDP : User Datagram Protocol API:

- reads and writes over socket file descriptors
- messages sent from/to ports to target a process on machine

Provide minimal reliability features:

- messages may be lost
- messages may be reordered
- messages may be duplicated
- only protection: checksums to ensure data not corrupted





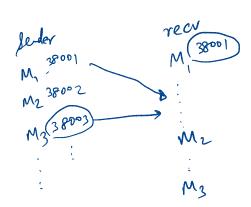
TCP: SEQUENCE NUMBERS

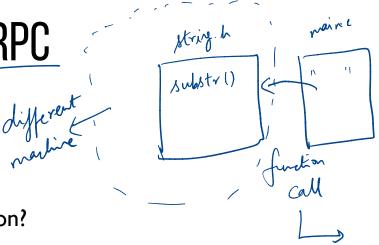
Sequence numbers

- senders gives each message an increasing unique seq number
- receiver knows it has seen all messages before N

Suppose message K is received.

- if K <= N, Msg K is already delivered, ignore it
- if K = N + I, first time seeing this message
- if K > N + 1? buffer this message



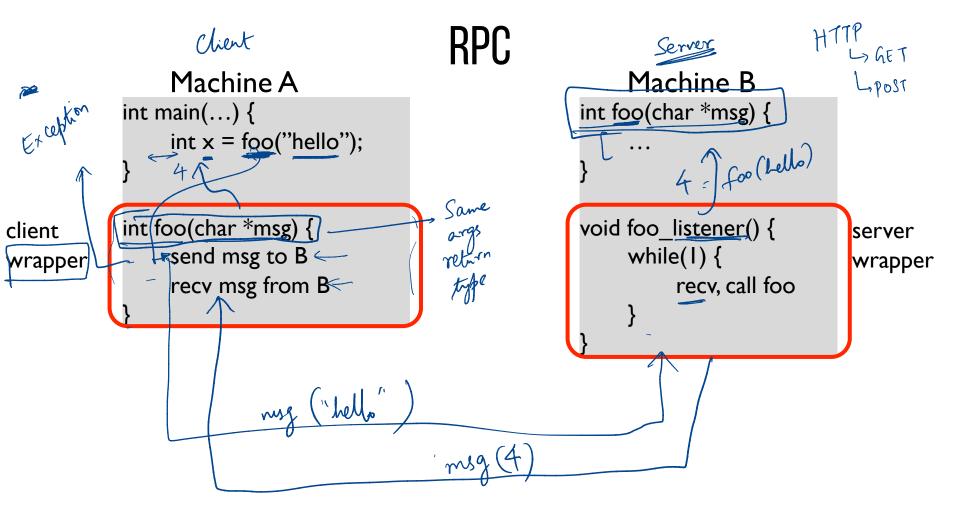


Remote Procedure Call

What could be easier than calling a function?

Approach: create wrappers so calling a function on another machine feels just like calling a local function!

Simplifies application development



int foo (string & RPC TOOLS send mag *- return volues. 3 return x;

more parallelism when sending RPCs

work to actually send / recv network deta

foo, proto file

int foo (string) bool lar (int)

for server. h

RPC packages help with two components

(I) Runtime library

- Thread pool -
- Socket listeners call functions on server

(2) Stub generation

- Create wrappers automatically
- Many tools available (rpcgen, thrift, protobufs) generate wonpik foo - client h loo client java

WRAPPER GENERATION TSON ->

fod " m Lello")

lytes

Connorly used

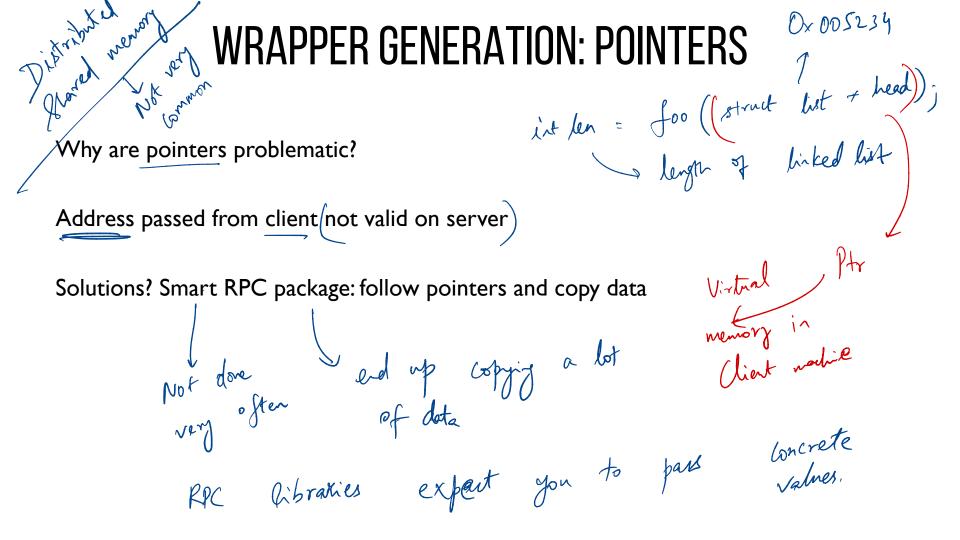
Wrappers must do conversions:

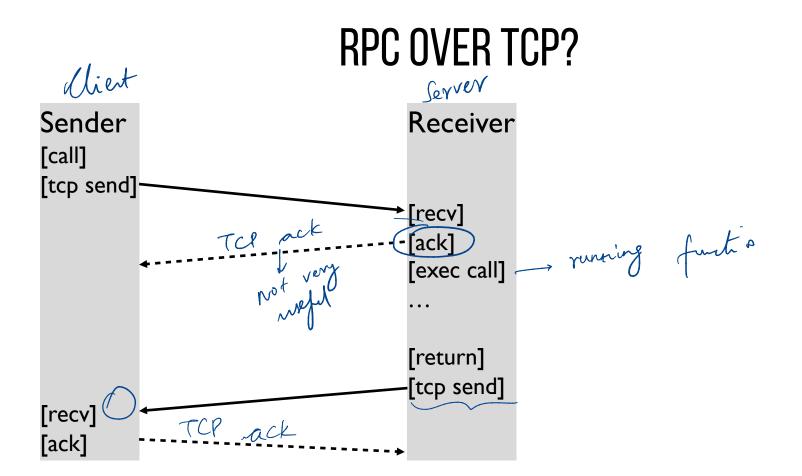
- client arguments to message
- message to server arguments
- convert server return value to message
- convert message to client return value

Need uniform endianness (wrappers do this)

Conversion is called marshaling/unmarshaling, or serializing deserializing

Serialize (object of type T) L) bytes



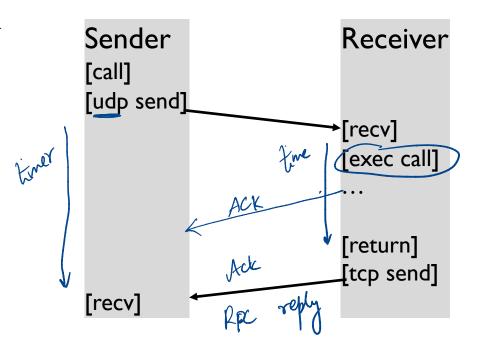


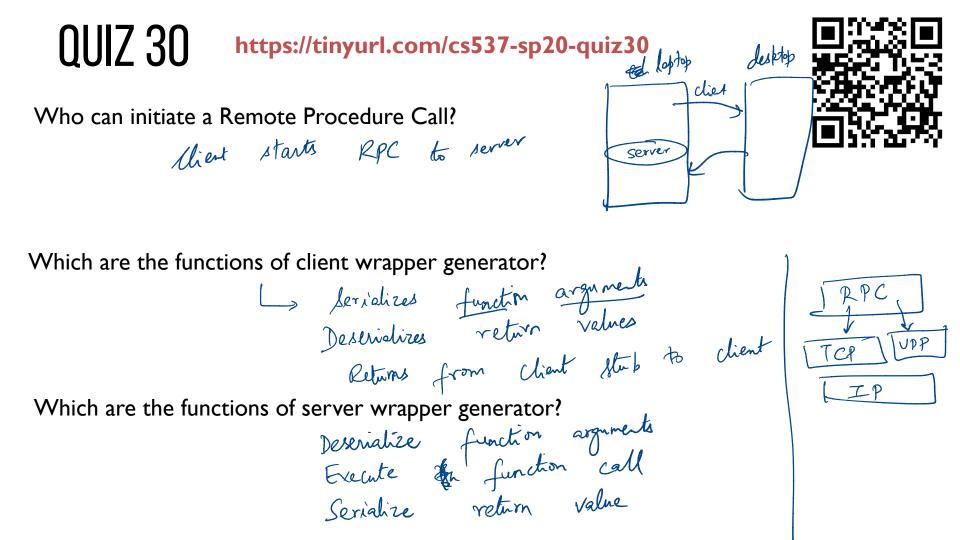
RPC OVER UDP

Strategy: use function return as implicit ACK

Piggybacking technique

What if function takes a long time? then send a separate ACK





DISTRIBUTED FILE SYSTEMS

Local FS: processes on same machine access shared files

Network FS: processes on different machines access shared files in same way laptop deskta bello^c hello^c hello^c hello^c hello^c

GOALS FOR DISTRIBUTED FILE SYSTEMS

Transparent access

- normal UNIX semantics open, read, write, stat etc.

 \rightarrow Fast + simple crash recovery: both clients and file server may crash

Reasonable performance?

NETWORK FILE SYSTEM: NFS < Sun microsystems

NFS: more of a protocol than a particular file system

Many companies have implemented NFS: Oracle/Sun, NetApp, EMC, IBM

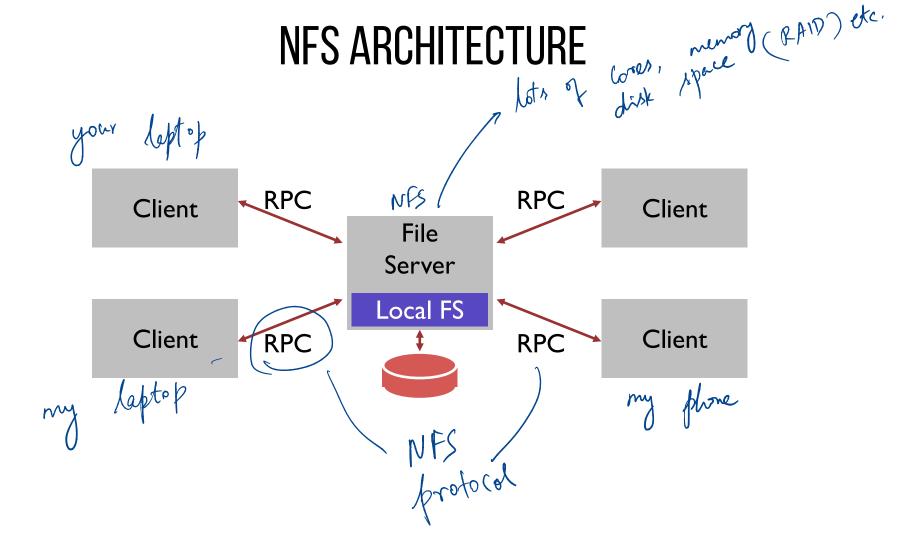
We're looking at NFSv2 NFSv4 has many changes Why look at an older protocol? Simpler, focused goals

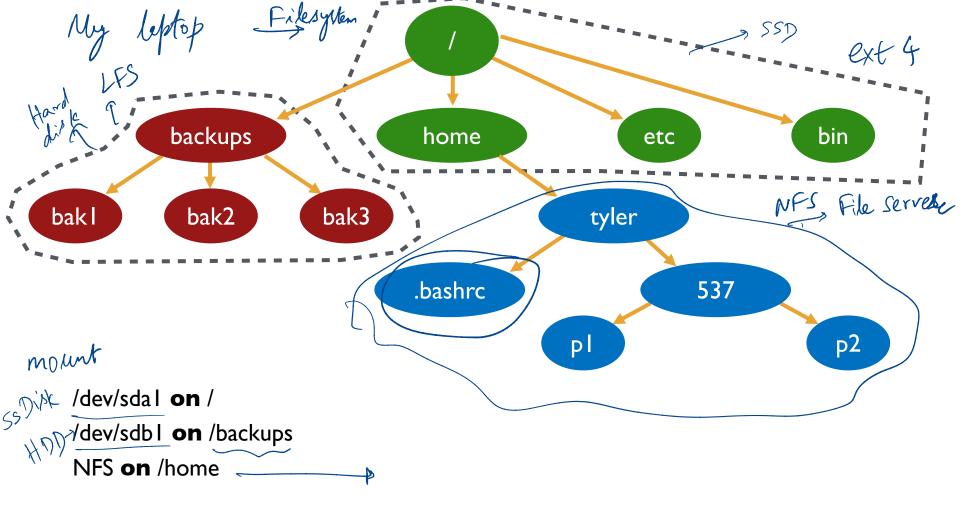
OVERVIEW

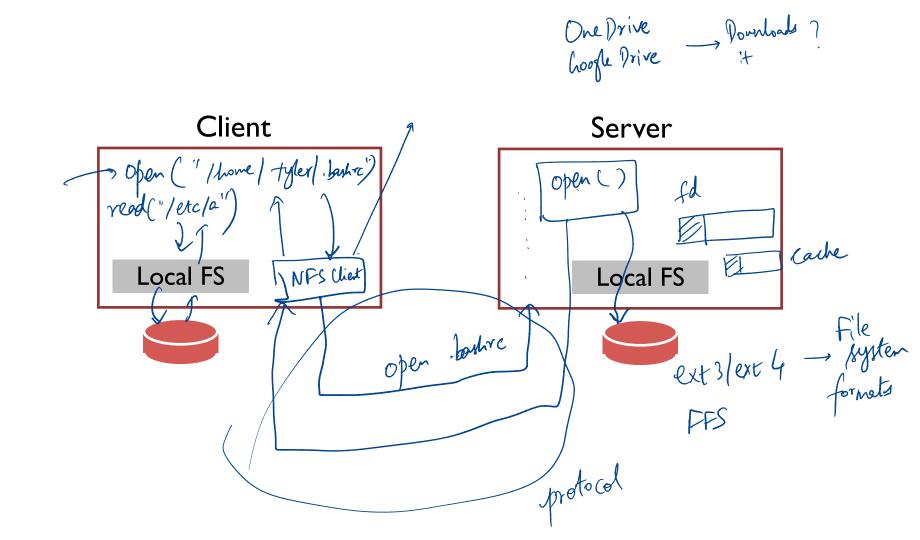
---- Architecture

Write Buffering

Cache







OVERVIEW

Architecture

Network API

Write Buffering

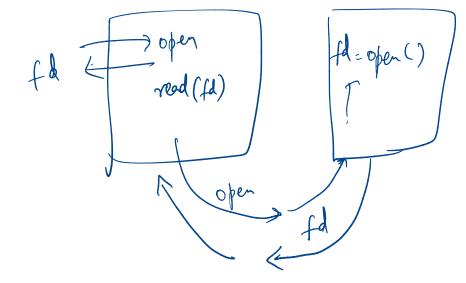
Cache

STRATEGY 1

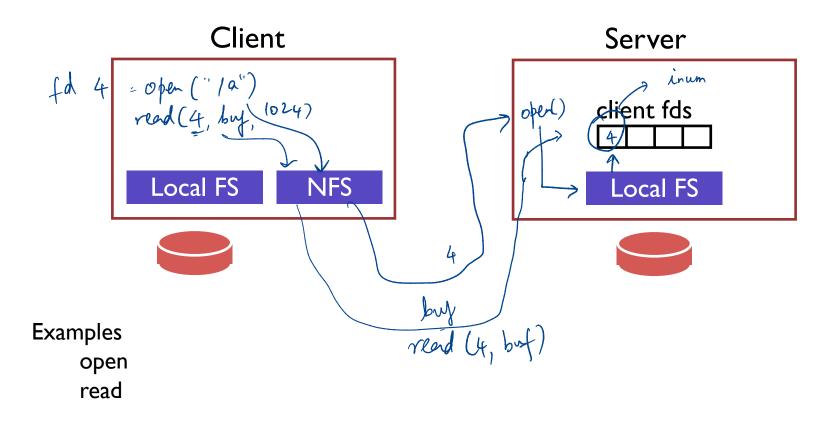
Attempt: Wrap regular UNIX system calls using RPC

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

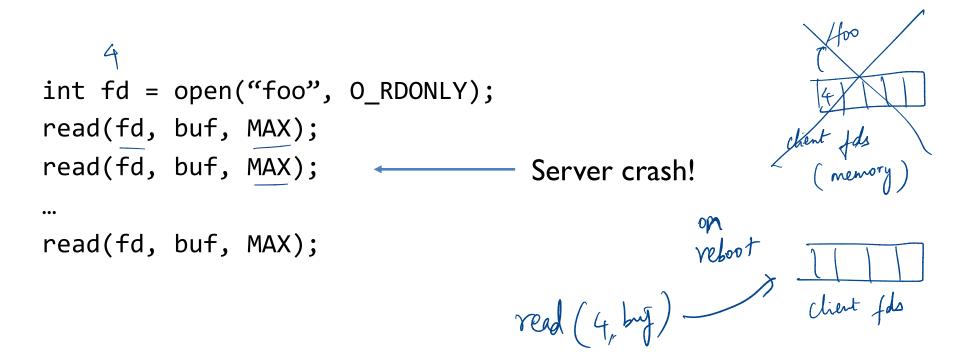
read(fd) on client calls read(fd) on server read(fd) on server returns data back to client



FILE DESCRIPTORS



STRATEGY 1: WHAT ABOUT CRASHES



POTENTIAL SOLUTIONS

I. Run some crash recovery protocol upon reboot

- Complex

- 2. Persist fds on server disk.
 - Slow
 - What if client crashes? When can fds be garbage collected?

STRATEGY 2: PUT ALL INFO IN REQUESTS

Use "stateless" protocol!

- server maintains no state about clients
- server still keeps other state, of course

Local FS

Server doent have to remember which FD is open by which client

STRATEGY 2: PUT ALL INFO IN REQUESTS

"Stateless" protocol: server maintains no state about clients

 INEED API change. One possibility:
 pread (char *path, buf, size, offset);
 pread (/ home / tyle/ / bachrc , buf, 100, 0)

 pwrite(char *path, buf, size, offset);
 pread (/ home / tyle/ / bachrc , buf, 100, 0)

 FD
 retry pread (

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

 Pros? Server Can crash and related transparently Cons? Too many path traversals

*path). **STRATEGY 3: FILE HANDLES** *will bort of the API*

, used to track inde

on server side

fh = open(char *path);
pread(<u>fh</u>, <u>buf</u>, <u>size</u>, <u>offset</u>);
pwrite(<u>fh</u>, <u>buf</u>, <u>size</u>, <u>offset</u>);

File Handle = <volume ID, inode #, generation #>

Opaque to client (client should not interpret internals) Which we have a cost NFS volume fixed traversal cost NFS volume fixed traversal cost the insole of fixed traversal cost is in fixed to fix the fixed the f

```
NFSPROC GETATTR
  expects: file handle
  returns: attributes
NFSPROC SETATTR
  expects: file handle, attributes
  returns: nothing
NFSPROC LOOKUP
  expects: directory file handle, name of file/directory to look up
  returns: file handle
NFSPROC READ
  expects: file handle, offset, count
  returns: data, attributes
NFSPROC WRITE
  expects: file handle, offset, count, data
  returns: attributes
NFSPROC CREATE
  expects: directory file handle, name of file, attributes
  returns: nothing
NFSPROC REMOVE
  expects: directory file handle, name of file to be removed
  returns: nothing
NFSPROC MKDIR
  expects: directory file handle, name of directory, attributes
  returns: file handle
NFSPROC_RMDIR
  expects: directory file handle, name of directory to be removed
  returns: nothing
NFSPROC READDIR
  expects: directory handle, count of bytes to read, cookie
  returns: directory entries, cookie (to get more entries)
```

Client

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

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

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

D, E, F partitions (pread (fh, ") NEXT STEPS Mount VFS Layer File handle int fh: Open ("/foo") Next class: More NFS server P5 is due TODAY! **AEFIS** feedback scope of the NFS mount file handle Optional project discussion 👍 Ash -> secure shell Commando on remate files -/backent / home Commends 1 server NFS Server 2 pochi NFS serve 2