

Hello!

# DISTRIBUTED SYSTEMS

Shivaram Venkataraman

CS 537, Spring 2023

# ADMINISTRIVIA

Project 6 grades → tomm AM

Project 7 → next week

Project 8 – Extra credit (4%) -

Midterm 3 conflicts →

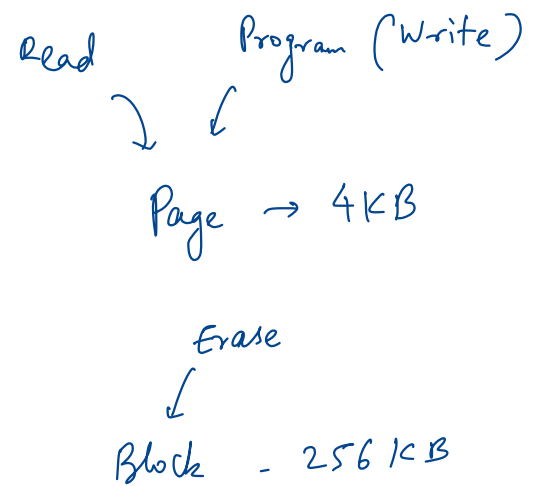
AEFIS feedback → useful

# AGENDA / LEARNING OUTCOMES

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

**RECAP**

# SSD OPERATIONS



Read a page: Retrieve contents of entire page (e.g., 4 KB)

- Cost: 25—75 microseconds
- Independent of page number, prior request offsets

Erase a block: Resets each page in the block to all 1s

- Cost: 1.5 to 4.5 milliseconds
- Much more expensive than reading!
- Allows each page to be written

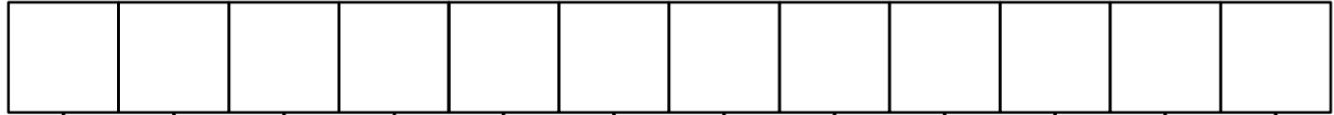
Program (i.e., write) a page: Change selected 1s to 0s

- Cost is 200 to 1400 microseconds
- Faster than erasing a block, but slower than reading a page

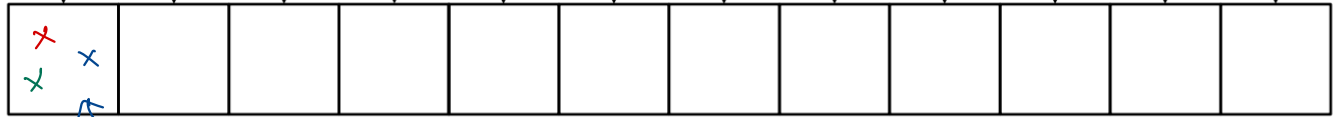
# FTL: DIRECT MAPPING

user or  
the FS sees

Logical  
pages



Physical  
pages



Flash  
Pages

Cons?

wear out !

Erase block frequently

Write amplification  
No wear-leveling

number of physical writes for 1 logical write

# FTL: LOG-BASED MAPPING

Idea: Treat the physical blocks like a log

Table: 100  $\rightarrow$  ~~0~~<sup>03</sup>    101  $\rightarrow$  1    2000  $\rightarrow$  02    Memory

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:	a1	a2	c1	a4									
State:	V	E	E	E	i	i	i	i	i	i	i	i	

1. Erase block

2. Write pages in order (log)

$\rightarrow$  Avoid read modify writes

$\rightarrow$  Avoids wear on few pages

# GARBAGE COLLECTION

Steps:

Read all pages in physical block

Write out the alive entries to the end of the log

Erase block (freeing it for later use)

Table: 100 → 4 101 → 5 2000 → 2 2001 → 3 Memory

Block:	0				1				2							
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip			
Content:	a1	a2	b1	b2	c1	c2							Chip			
State:	V	V	V	V	V	V	E	E	i	i	i	i				
	<i>block</i>															

Table: 100 → 4 101 → 5 2000 → 6 2001 → 7 Memory

Block:	0				1				2				
Page:	00	01	02	03	04	05	06	07	08	09	10	11	Flash Chip
Content:					c1	c2	b1	b2					Chip
State:	E	E	E	E	V	V	V	V	i	i	i	i	

↳ live pages end of the log



# SSD VS HDD PERFORMANCE

small requests

SLC vs. MLC  
FTL layer

Actual being realized as sequential operations

SSDs are good at random access

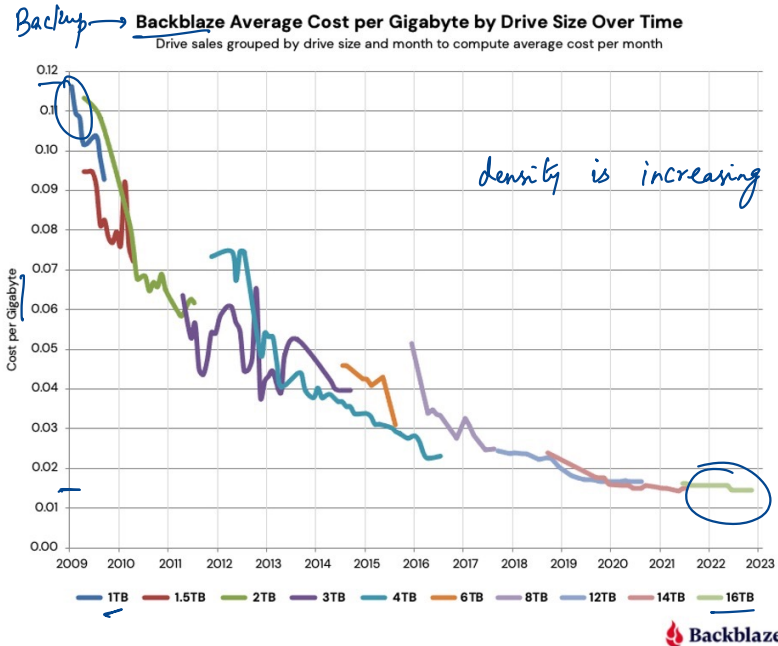
Device	Random		Sequential	
	Reads (MB/s)	Writes (MB/s)	Reads (MB/s)	Writes (MB/s)
Samsung 840 Pro SSD	103	287	421	384
Seagate 600 SSD	84	252	424	374
Intel SSD 335 SSD	39	222	344	354
Seagate Savvio 15K.3 HDD	2	2	223	223

much lower

random write

large granularity

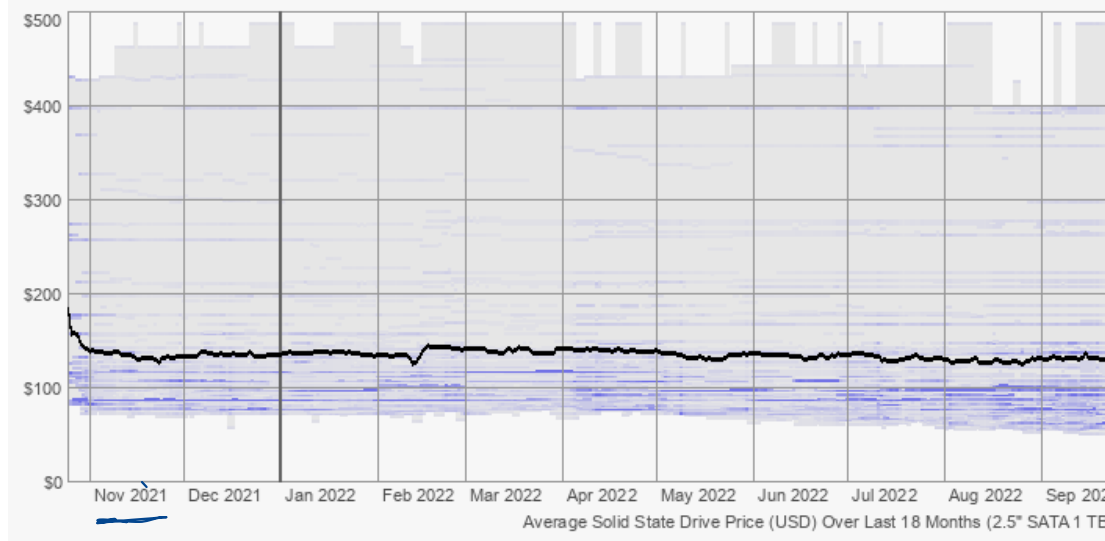
# SSD VS HDD COST



Backblaze

= ~1.5 cents / GB

*10x more expensive*



1TB ~ \$150 on average  
~15 cents / GB

# PERSISTENCE SUMMARY

FS API ← applications  
open,



Managing I/O devices is a significant part of OS!

Disk drives: storage media with specific geometry

SSDs: Pages, Blocks

Interface

Disk Drives / SSD

Filesystems: OS provided API to access disk

Simple FS: FS layout with SB, Bitmaps, Inodes, Datablocks

FFS: Split simple FS into groups. Key idea: put inode, data close to each other

LFS: Puts data where it's fastest to write, hope future reads cached in memory

<https://www.eecs.harvard.edu/~margo/papers/usenix95-lfs/supplement/>

FSCK, Journaling

# DISTRIBUTED SYSTEMS

# WHAT IS A DISTRIBUTED SYSTEM?

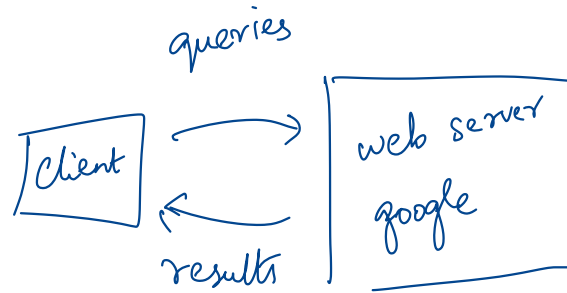
*A distributed system is one where a machine I've never heard of can cause my program to fail.*

— Leslie Lamport → Turing award winner

Definition: More than one machine working together to solve a problem

Examples:

- client/server: web server and web client
- cluster: page rank computation



# WHY GO DISTRIBUTED?

- Increase throughput



- Connectivity → share or send data

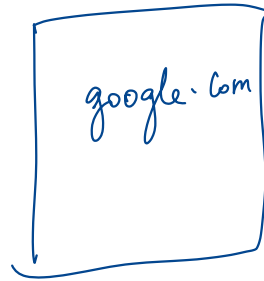


more performance

- Specialized machines → better performance

- Increase fault tolerance

↳ Tolerate failures



x if 1 machine fails

# WHY GO DISTRIBUTED?

More computing power

More storage capacity

Fault tolerance

Data sharing

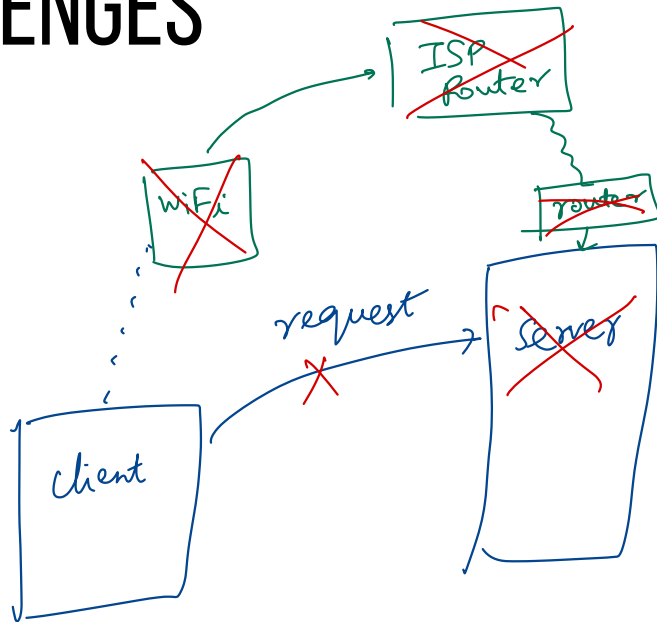
# NEW CHALLENGES

System failure: need to worry about **partial** failure

Communication failure: links unreliable

- bit errors
- packet loss
- node/link failure

just this  
request  
was  
lost





# COMMUNICATION OVERVIEW

Raw messages: UDP

Reliable messages: TCP

Remote procedure call: RPC

→ Acronyms

640 Networking

# 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

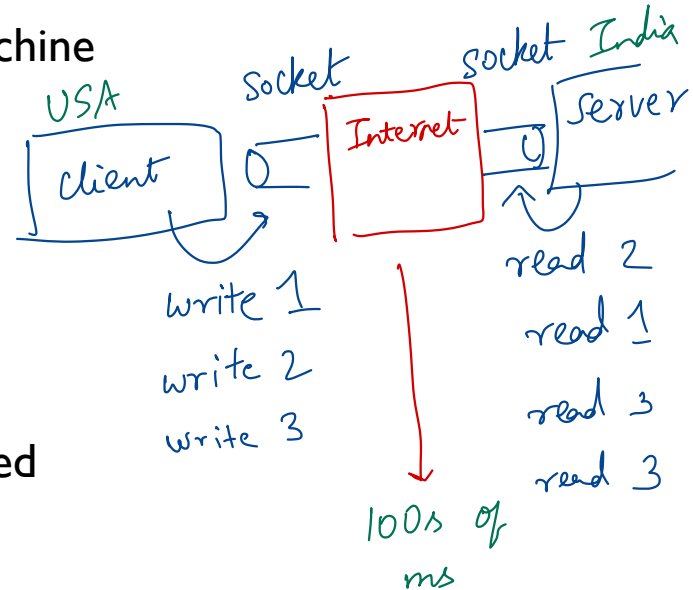
Provide minimal reliability features:

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



very lightweight

file ↔ socket



# RAW MESSAGES: UDP

## Advantages

- Lightweight
- Some applications make better reliability decisions themselves (e.g., video conferencing programs)

## Disadvantages

- More difficult to write applications correctly

Server → video bytes → client

reliably → retry same frame

skip frames → responsiveness

# NOT A QUIZ?

Course feedback: <https://aefis.wisc.edu>

# RELIABLE MESSAGES: LAYERING STRATEGY

TCP: Transmission Control Protocol



*mostly commonly used  
protocol*

Using software to build

reliable logical connections over unreliable physical connections

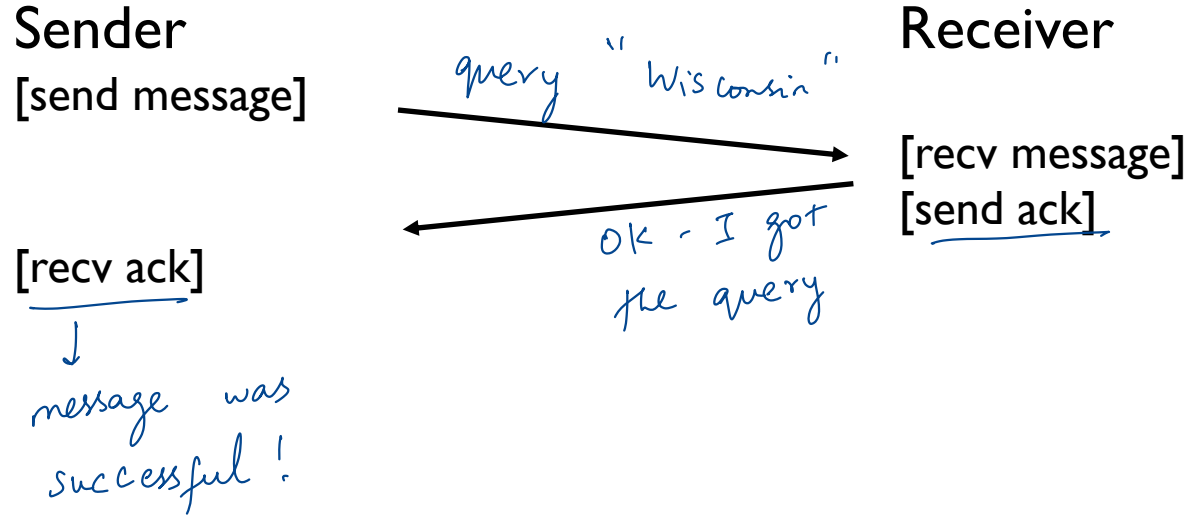


*interface*



*links, machines unreliable*

# TECHNIQUE #1: ACK → Acknowledgement



Ack: Sender knows message was received  
What to do about message loss?

# TECHNIQUE #2: TIMEOUT

↓  
send("Wisconsin")

TCP

Sender

[send message]  
[start timer]

... waiting for ack ...

[timer goes off]  
[send message]

[recv ack]

Receiver



[recv message]  
[send ack]

30s

retry →  
same  
message

send returns

# TIMEOUT

Adaptively configure timeout

↳ 640

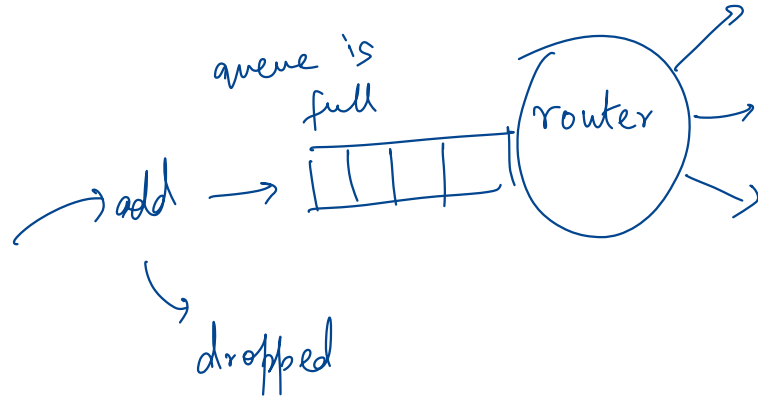
How long to wait?

Too long? →

- System feels unresponsive

Too short?

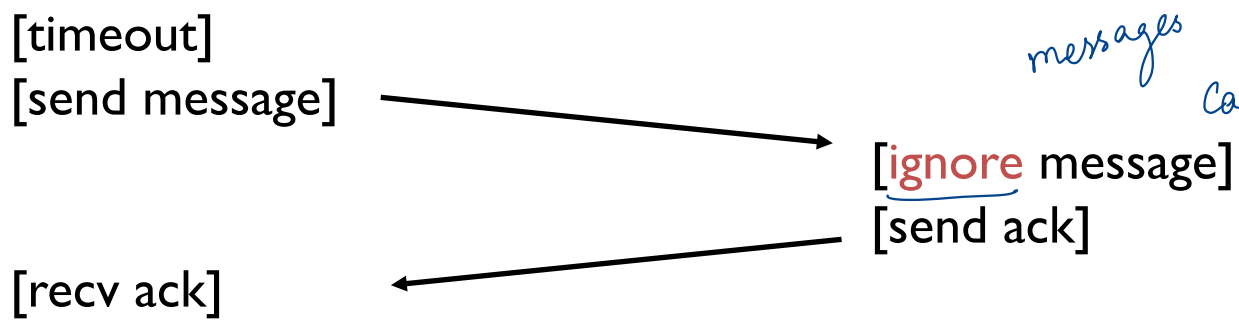
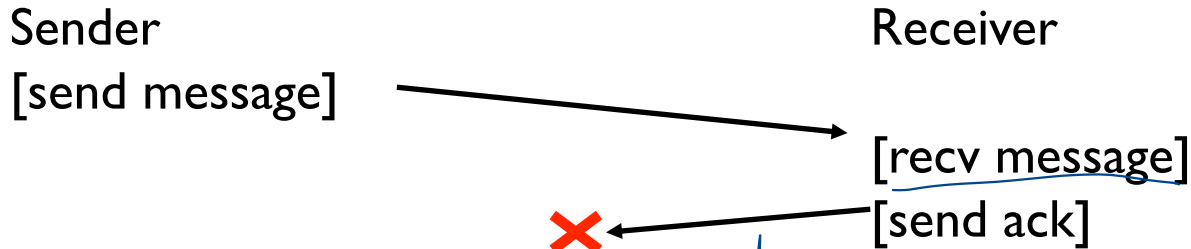
- Messages needlessly re-sent →
- Messages may have been dropped due to overloaded server. Resending makes overload worse!



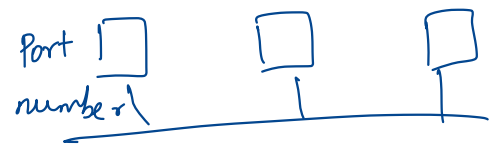


Exactly once semantics → <sup>dropped</sup> x duplicated → <sup>retry</sup>

# LOST ACK PROBLEM



messages need an ID  
can be used to ignore messages already delivered



# SEQUENCE NUMBERS

→ avoid duplication  
prevent reordering

Packets = 64 bytes



## Sequence numbers

- sender gives each message an increasing unique seq number  $Seq.$
- receiver knows it has seen all messages before  $N$

$N \rightarrow$  last message delivered  
=  $Recv$

msg 1  
seq-no: 2052 → 2052 =

Suppose message  $K$  is received.

- if  $K \leq N$ , Msg  $K$  is already delivered, ignore it
- if  $K = N + 1$ , first time seeing this message
- if  $K > N + 1$ ?

↳ buffer it until  $N+1$  arrives

deliver & increment  $K$

msg 2  
seq-no: 2053 →  
↓  
⋮  
 $K$   
incoming message

# TCP

TCP: Transmission Control Protocol

Most popular protocol based on seq nums

Buffers messages so arrive in order

Timeouts are adaptive

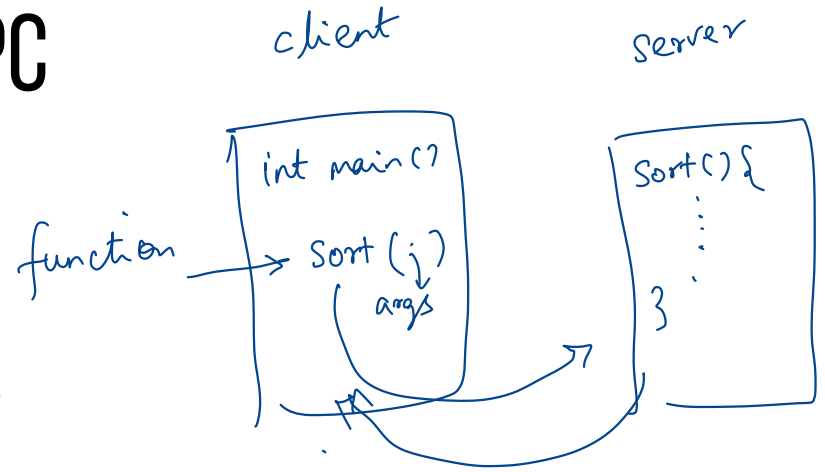
# COMMUNICATIONS OVERVIEW

Raw messages: UDP

Reliable messages: TCP

Remote procedure call: RPC → *Software abstraction*

# RPC



## 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!

RPC library → provides abstraction

# RPC

## Machine A

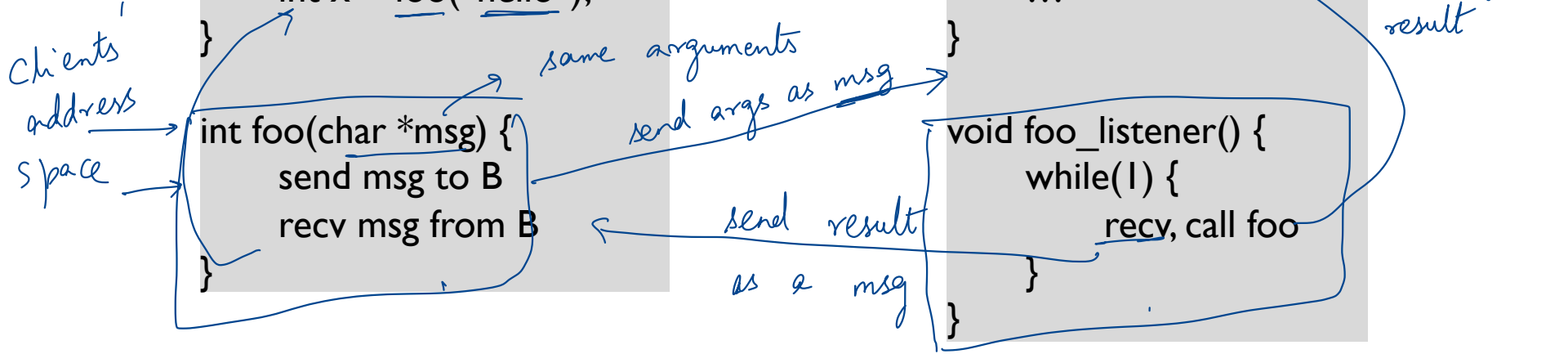
```
int main(...) {  
    int x = foo("hello");  
}
```

```
int foo(char *msg) {  
    send msg to B  
    rcv msg from B  
}
```

## Machine B

```
int foo(char *msg) {  
    ...  
}
```

```
void foo_listener() {  
    while(1) {  
        rcv, call foo  
    }  
}
```



# RPC

## Machine A

```
int main(...) {  
    int x = foo("hello");  
}
```

```
int foo(char *msg) {  
    send msg to B  
    rcv msg from B  
}
```

## Machine B

```
int foo(char *msg) {  
    ...  
}
```

```
void foo_listener() {  
    while(1) {  
        rcv, call foo  
    }  
}
```

*client*

*Developers!  
users*

client wrapper

server wrapper

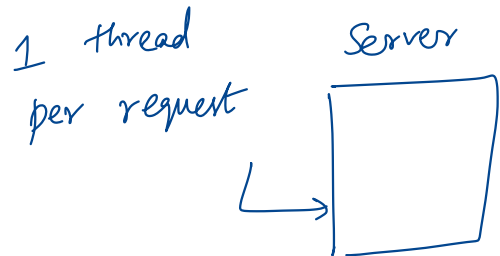
*value :: int → 4 bytes*

*ptr :: int\**

*↳ (\*int) :: int → serialize*

*part of  
RPC library*

# RPC TOOLS



RPC packages help with two components

## (1) Runtime library

- Thread pool
- Socket listeners call functions on server

so that you can handle multiple requests at same time

## (2) Stub generation

- Create wrappers automatically
- Many tools available (rpcgen, thrift, protobufs)

serializing / de-serializing arguments / results



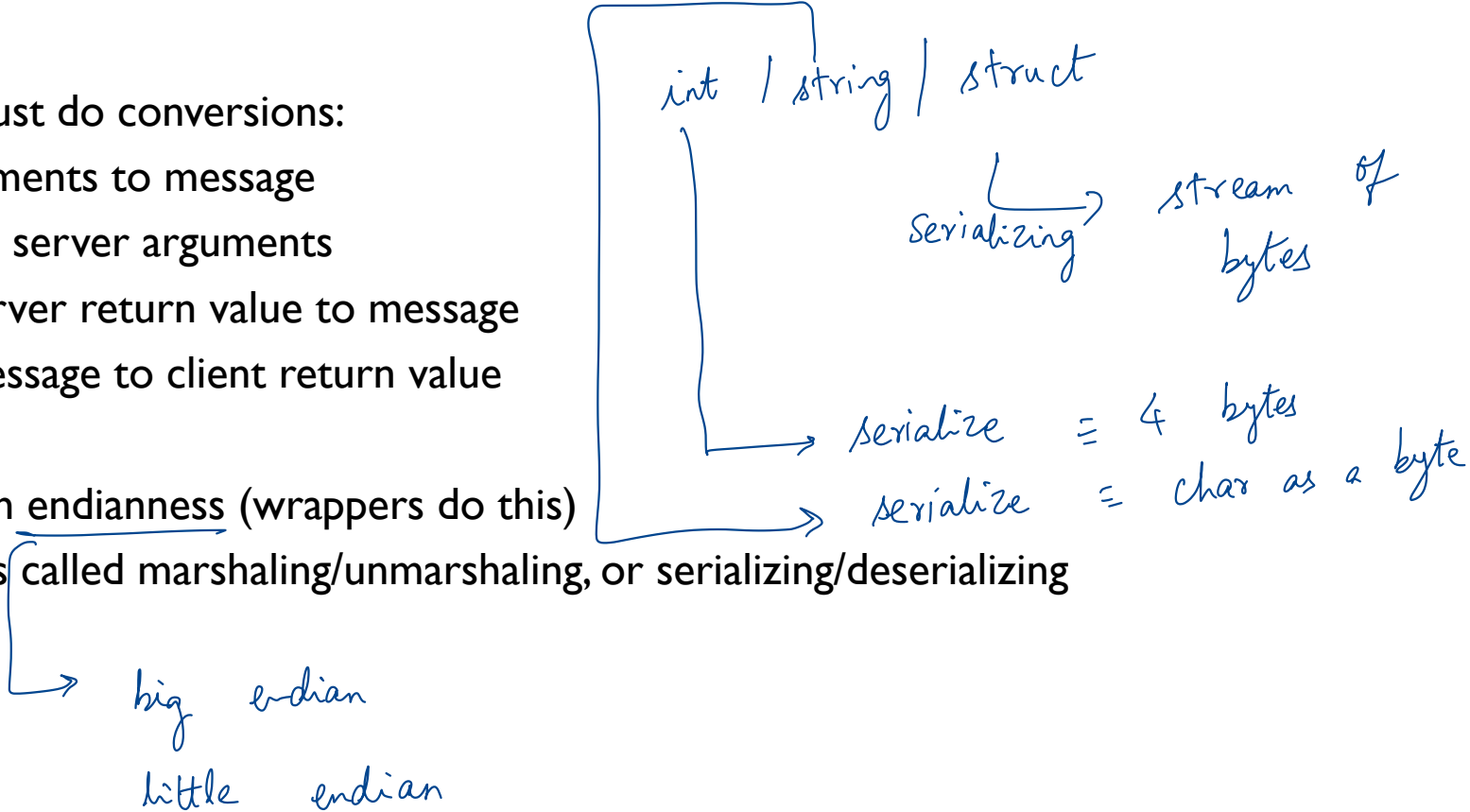
# WRAPPER GENERATION

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



big endian

little endian

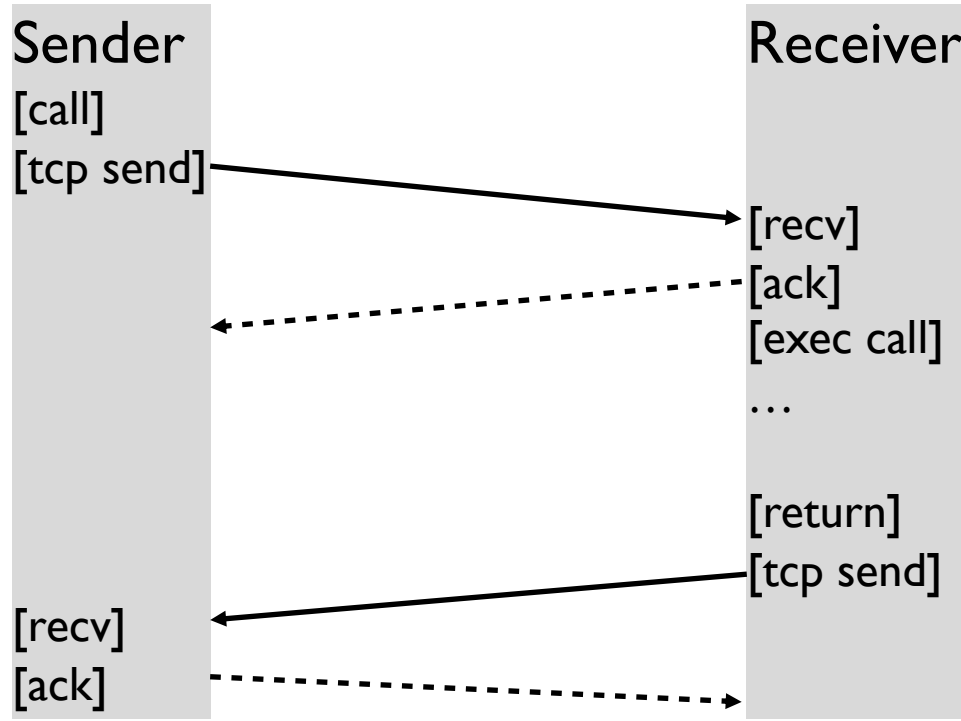
# WRAPPER GENERATION: POINTERS

Why are pointers problematic?

Address passed from client not valid on server

Solutions? Smart RPC package: follow pointers and copy data

# RPC OVER TCP?

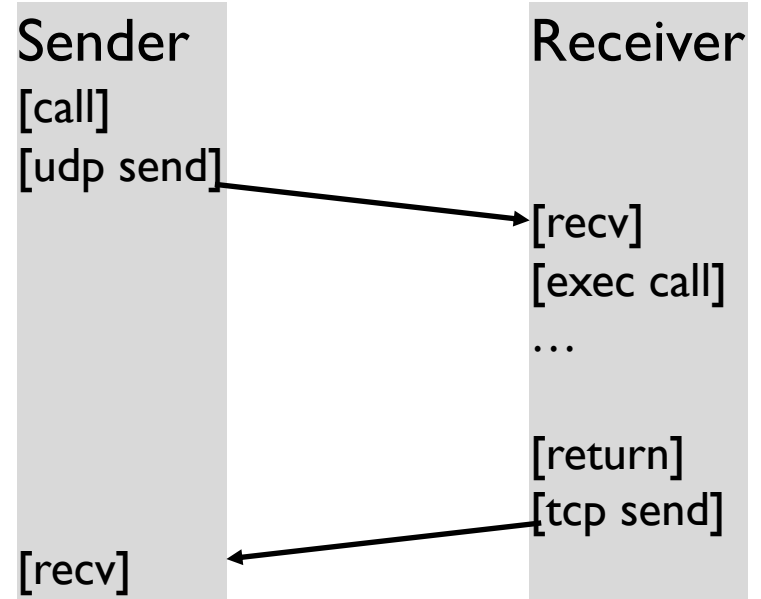


# RPC OVER UDP

Strategy: use function return as implicit ACK

Piggybacking technique

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



# NEXT STEPS

Distributed Filesystems