

Distributed Systems

Remote Procedure Calls

Today's Agenda

- Last time:
 - Computer networks, primarily from an application perspective
 - Protocol layering
 - Client-server architecture
 - End-to-end principle

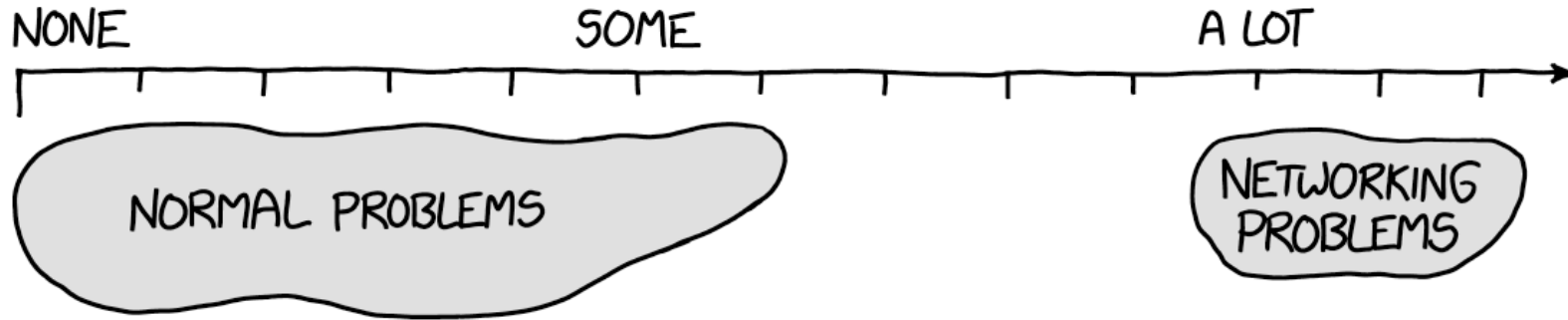
- Today:
 - TCP
 - Socket programming
 - Remote Procedure Calls

Assignment Reminder

- Start early
 - Most assignments will be simple IF you start thinking about the design early and get help in labs/office hours
 - Systems projects take $\sim 3X$ the time estimate
- Probability of email response = $O(\text{Time to deadline})$
 - Emailing day of deadline \rightarrow Can't provide helpful response
 - Canvas discussion board: increased chance of response from TA's and other students
- "Head fake": The assignments are not for teaching distributed systems.
 - They are a way to get "comfortable" with uncertainty and frustration.
 - Hence the loose specification and sometimes ambiguous descriptions

TYPES OF COMPUTER PROBLEMS

BY HOW MUCH DEBUGGING THEM MAKES YOUR BRAIN STOP WORKING



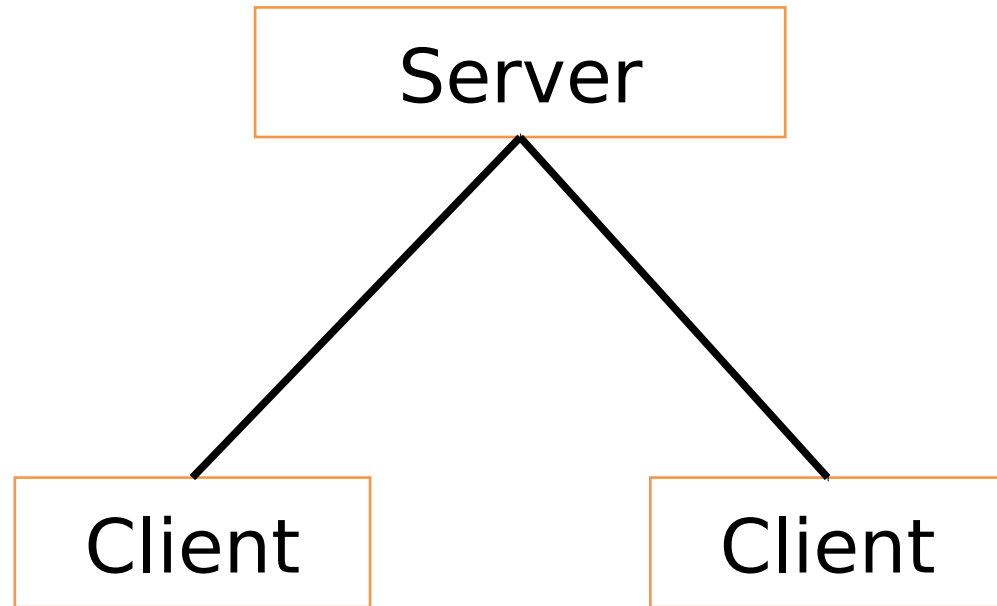
*BEFORE NOON, ODD-NUMBERED
PACKETS WERE LAGGY, BUT AFTER
NOON, EVEN-NUMBERED ONES ARE!
IT'S THE OPPOSITE OF YESTERDAY!*

ARE YOU SURE YOU'RE OKAY?

*I'M FINE AND I BELIEVE
IN GHOSTS NOW!*



Client-server architecture



Server:

- always-on host
- permanent IP address
- data centers for scaling

Clients:

- communicate with server
- may be intermittently connected
- may have dynamic IP addresses
- do not communicate directly with each other

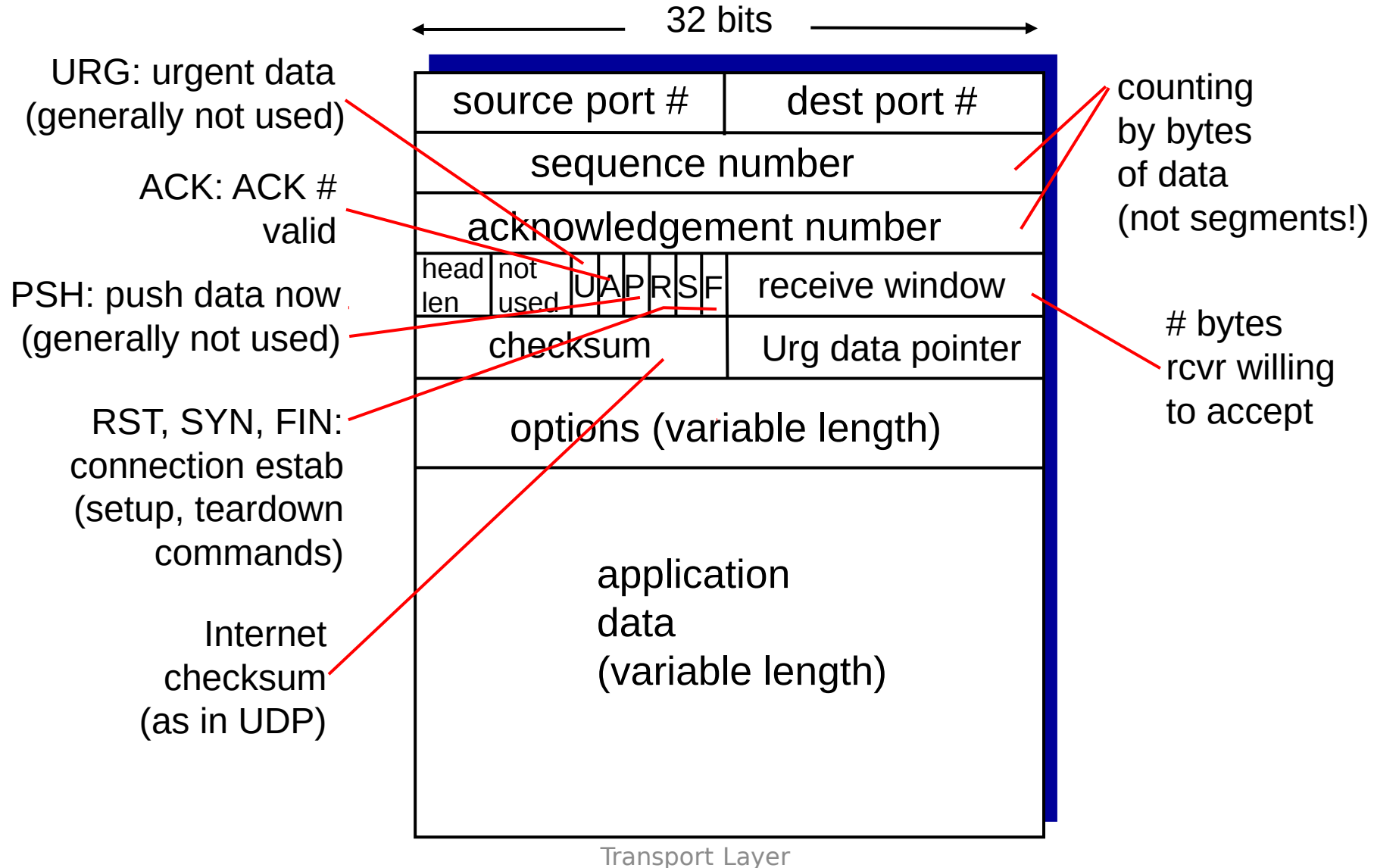
TCP: Overview

2018, 2581

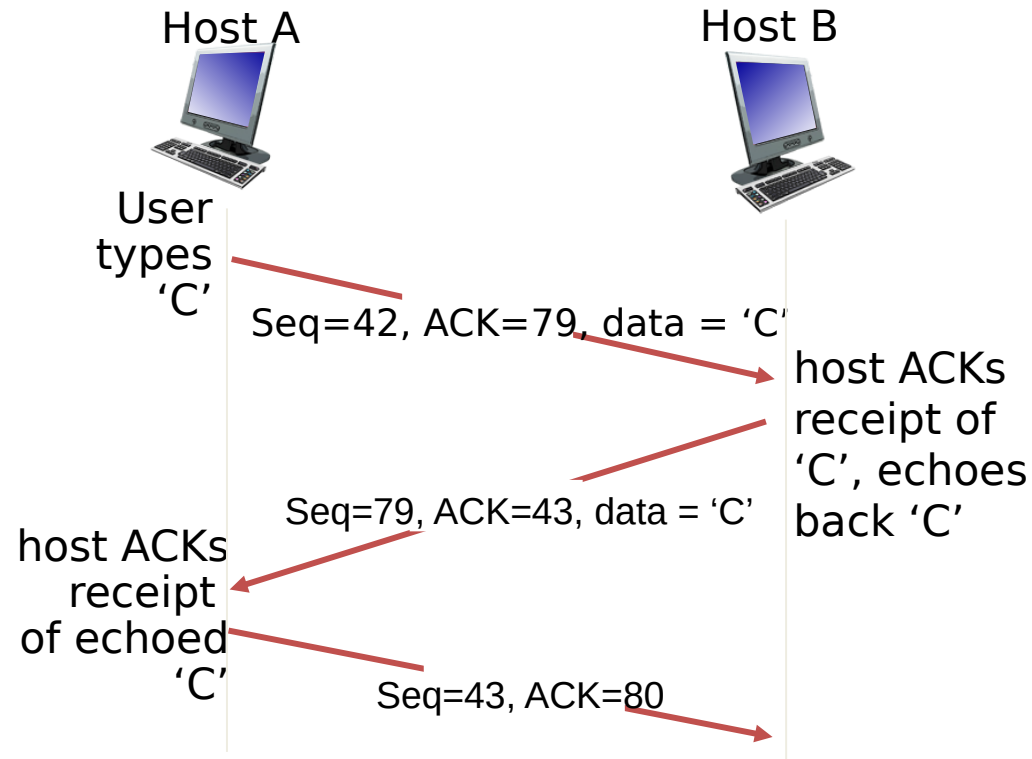
RFCs: 793,1122,1323,

- **point-to-point:**
 - one sender, one receiver
- **reliable, in-order *byte stream*:**
 - no “message boundaries”
- **pipelined:**
 - TCP congestion and flow control set window size
- ❖ **full duplex data:**
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- ❖ **connection-oriented:**
 - handshaking (exchange of control msgs) inits sender, receiver state before data exchange
- ❖ **flow controlled:**
 - sender will not overwhelm receiver

TCP segment structure



TCP seq. numbers, ACKs



simple telnet scenario

TCP seq. numbers, ACKs

sequence numbers:

- byte stream “number” of first byte in segment’s data

acknowledgements:

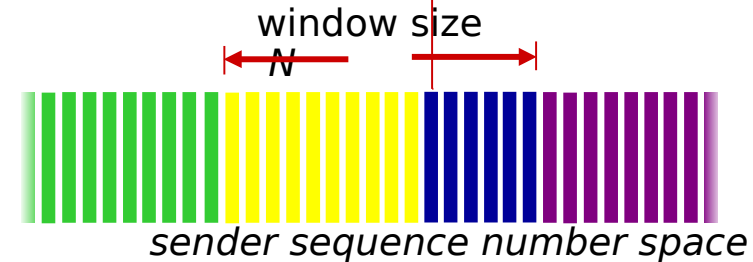
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments

- **A:** TCP spec doesn’t say, - up to implementor

outgoing segment from sender

source port #	dest port #
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer



incoming segment to sender

source port #	dest port #
sequence number	
acknowledgement number	
A	rwnd
checksum	urg pointer

TCP sender events:

data rcvd from app:

- ❖ create segment with seq #
- ❖ seq # is byte-stream number of first data byte in segment
- ❖ start timer if not already running
 - think of timer as for oldest unacked segment
 - expiration interval: `TimeoutInterval`

timeout:

- ❖ retransmit segment that caused timeout
- ❖ restart timer

ack rcvd:

- ❖ if ack acknowledges previously unacked segments
 - update what is known to be ACKed
 - start timer if there are still unacked segments

Approaches towards congestion control

two broad approaches towards congestion control:

end-end congestion control:

- ❖ no explicit feedback from network
- ❖ congestion inferred from end-system observed loss, delay
- ❖ approach taken by TCP

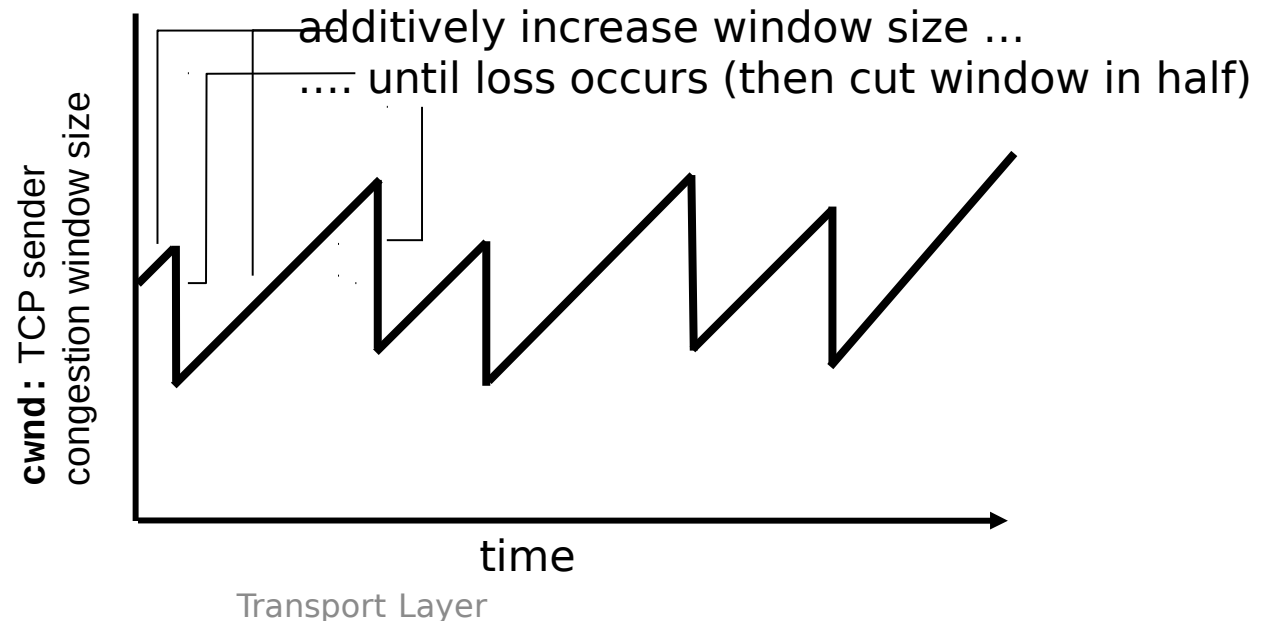
network-assisted congestion control:

- ❖ routers provide feedback to end systems
 - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - explicit rate for sender to send at

TCP congestion control: additive increase multiplicative decrease

- ❖ *approach*: sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
 - *additive increase*: increase **cwnd** by 1 MSS every RTT until loss detected
 - *multiplicative decrease*: cut **cwnd** in half after loss

AIMD saw tooth behavior: probing for bandwidth



TCP Performance

- Roughly, max throughput = Window size/RTT
- Throughput = $1/RTT * (\sqrt{2/3} * \text{packet-loss-probability})$
- TCP performance also depends on receive buffer sizes

Socket programming *with* *UDP*

UDP: no “connection” between client & server

- no handshaking before sending data
- sender explicitly attaches IP destination address and port # to each packet
- rcvr extracts sender IP address and port# from received packet

UDP: transmitted data may be lost or received out-of-order

Application viewpoint:

- UDP provides *unreliable* transfer of groups of bytes (“datagrams”) between client and server

Client/server socket interaction: UDP

server (running on serverIP)

create socket, port= x:
`serverSocket =
socket(AF_INET,SOCK_DGRAM)`

↓
read datagram from
`serverSocket`

↓
write reply to
`serverSocket`
specifying
client address,
port number

client

create socket:
`clientSocket =
socket(AF_INET,SOCK_DGRAM)`

↓
Create datagram with server IP and
port=x; send datagram via
`clientSocket`

↓
read datagram from
`clientSocket`

↓
close
`clientSocket`

Example app: UDP client

Python UDPClient

include Python's socket library

→ from socket import *

serverName = 'hostname'

serverPort = 12000

create UDP socket for server

→ clientSocket = socket(socket.AF_INET,
socket.SOCK_DGRAM)

get user keyboard input

→ message = raw_input('Input lowercase sentence:')

Attach server name, port to message; send into socket

→ clientSocket.sendto(message,(serverName, serverPort))

read reply characters from socket into string

→ modifiedMessage, serverAddress = clientSocket.recvfrom(2048)

print modifiedMessage

print out received string and close socket

→ clientSocket.close()

Example app: UDP server

Python UDPServer

```
from socket import *
serverPort = 12000
serverSocket = socket(AF_INET, SOCK_DGRAM)
serverSocket.bind(('', serverPort))
print "The server is ready to receive"
while 1:
    message, clientAddress = serverSocket.recvfrom(2048)
    modifiedMessage = message.upper()
    serverSocket.sendto(modifiedMessage, clientAddress)
```

create UDP socket →

bind socket to local port number 12000 →

loop forever →

Read from UDP socket into message, getting client's address (client IP and port) →

send upper case string back to this client →

Socket programming *with* *TCP*

client must contact server

- server process must first be running
- server must have created socket (door) that welcomes client's contact

client contacts server by:

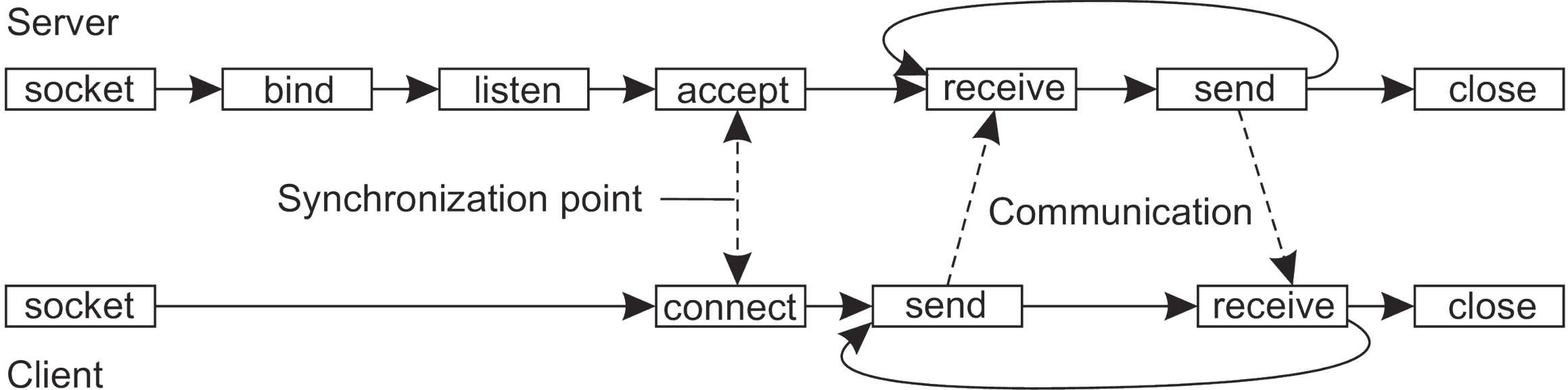
- Creating TCP socket, specifying IP address, port number of server process
- *when client creates socket:* client TCP establishes connection to server TCP

- when contacted by client, *server TCP creates new socket* for server process to communicate with that particular client
 - allows server to talk with multiple clients
 - source port numbers used to distinguish clients

application viewpoint:

TCP provides reliable, in-order byte-stream transfer ("pipe") between client and server

Socket Programming With TCP



Client/server socket interaction: TCP

server (running on `hostid`)

client

create socket,
port=`x`, for incoming
request:
`serverSocket = socket()`

wait for incoming
connection request
`connectionSocket =`
`serverSocket.accept()`

read request from
`connectionSocket`

write reply to
`connectionSocket`

close
`connectionSocket`

create socket,
connect to `hostid`, port=`x`
`clientSocket = socket()`

send request using
`clientSocket`

read reply from
`clientSocket`

close
`clientSocket`

TCP
connection setup

Example app: TCP client

Python TCPClient

create TCP socket for
server, remote port 12000

```
from socket import *
serverName = 'servername'
serverPort = 12000
clientSocket = socket(AF_INET, SOCK_STREAM)
clientSocket.connect((serverName,serverPort))
sentence = raw_input('Input lowercase sentence:')
clientSocket.send(sentence)
modifiedSentence = clientSocket.recv(1024)
print 'From Server:', modifiedSentence
clientSocket.close()
```

No need to attach server
name, port

Example app: TCP server

Python TCP Server

create TCP welcoming
socket

server begins listening for
incoming TCP requests

loop forever

server waits on accept()
for incoming requests, new
socket created on return

read bytes from socket (but
not address as in UDP)

close connection to this
client (but *not* welcoming
socket)

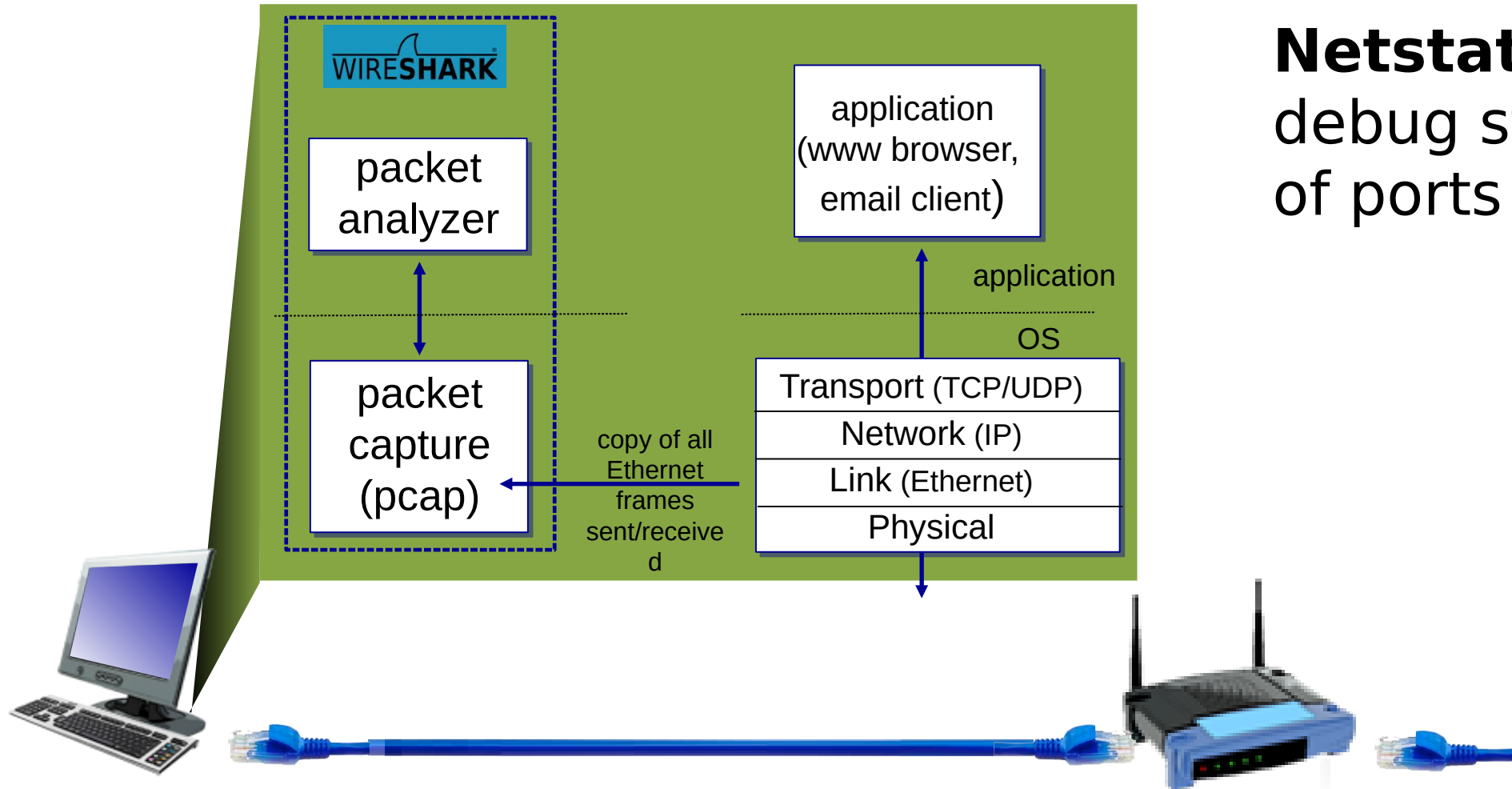
```
from socket import *
serverPort = 12000
serverSocket = socket(AF_INET,SOCK_STREAM)
serverSocket.bind(('',serverPort))
serverSocket.listen(1)
print 'The server is ready to receive'
while 1:
    connectionSocket, addr = serverSocket.accept()

    sentence = connectionSocket.recv(1024)
    capitalizedSentence = sentence.upper()
    connectionSocket.send(capitalizedSentence)
    connectionSocket.close()
```

Higher Level Networking

- Client/server code abstracted out (python's twisted framework)
- Message queues: Kafka, ZeroMQ, etc
 - Durability of messages (can persist on disk)
 - Message lifetimes (time to live)
 - Filtering, queueing policies
 - Batching policies
 - Delivery policies (at most once, at least once, etc)

Debugging Networks: Packet Capture

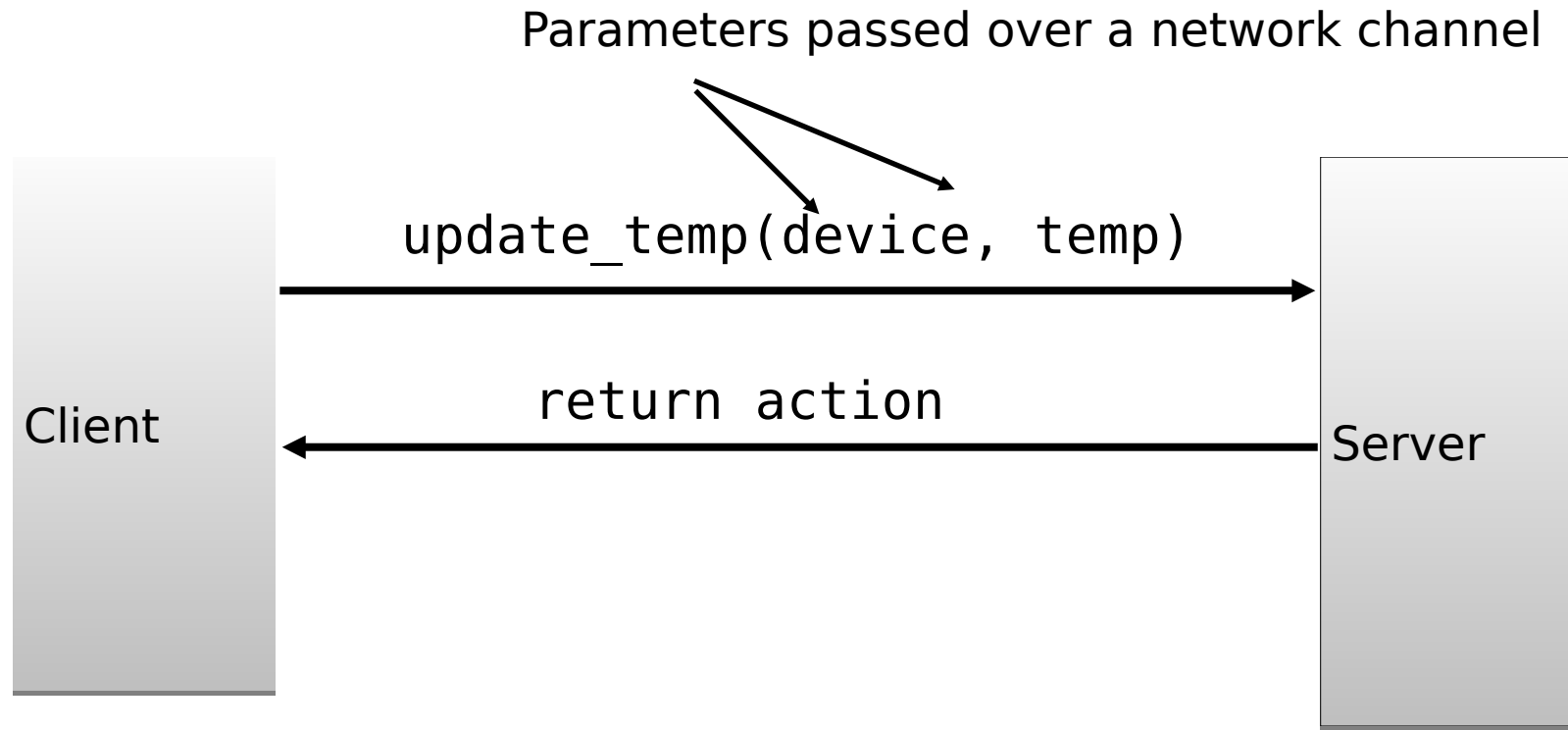


Remote Procedure Calls

Remote Procedure Calls

- Procedure (function) calls a well known and understood mechanism for transfer of data and control within a program/process
- Remote Procedure Calls : extend conventional local calls to work *across* processes.
 - Processes may be running on different machines
 - Allows communication of data via function parameters and return values
 - RPC invocations also serve as notifications (transfer of control)

RPC Example



RPC Advantages

- Clean and simple to understand semantics similar to local procedure calls
- Generality: all languages have local procedure calls
 - RPC libraries augment the procedure call interface to make RPCs appear similar to local calls
- Abstraction for a common client/server communication pattern

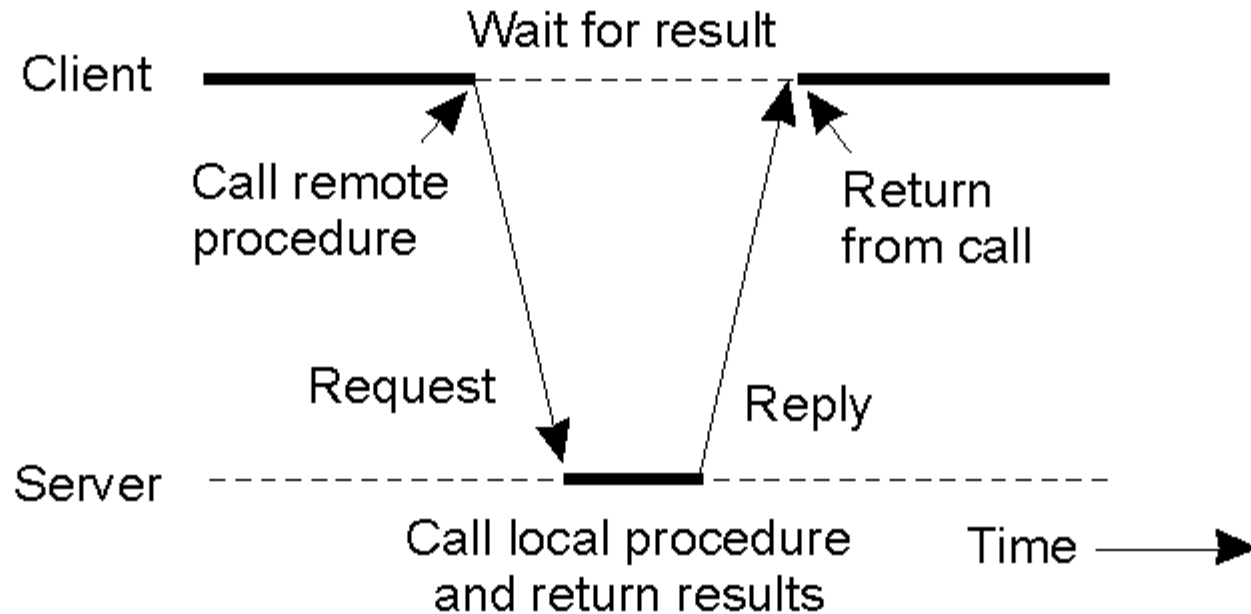
```
push_temp(name) {  
    t = get_current_temp();  
    return update_temp (name, t); //RPC  
}
```

Challenges

- RPCs impose new challenges not faced in local calls
- How to pass parameters?
 - Passing data over a network raises issues like endian-ness
 - Pointers: machines may not share an address space
- How to deal with machine failures?
 - Local procedures are assumed to always run
 - A remote machine running an RPC may face crashes, network issues
 - Need to consider failure semantics in RPC implementations
- How to integrate RPCs with existing language runtimes?
 - Seamless local and remote calls
 - Integrate RPCs with language caller/callee interface

RPC Semantics

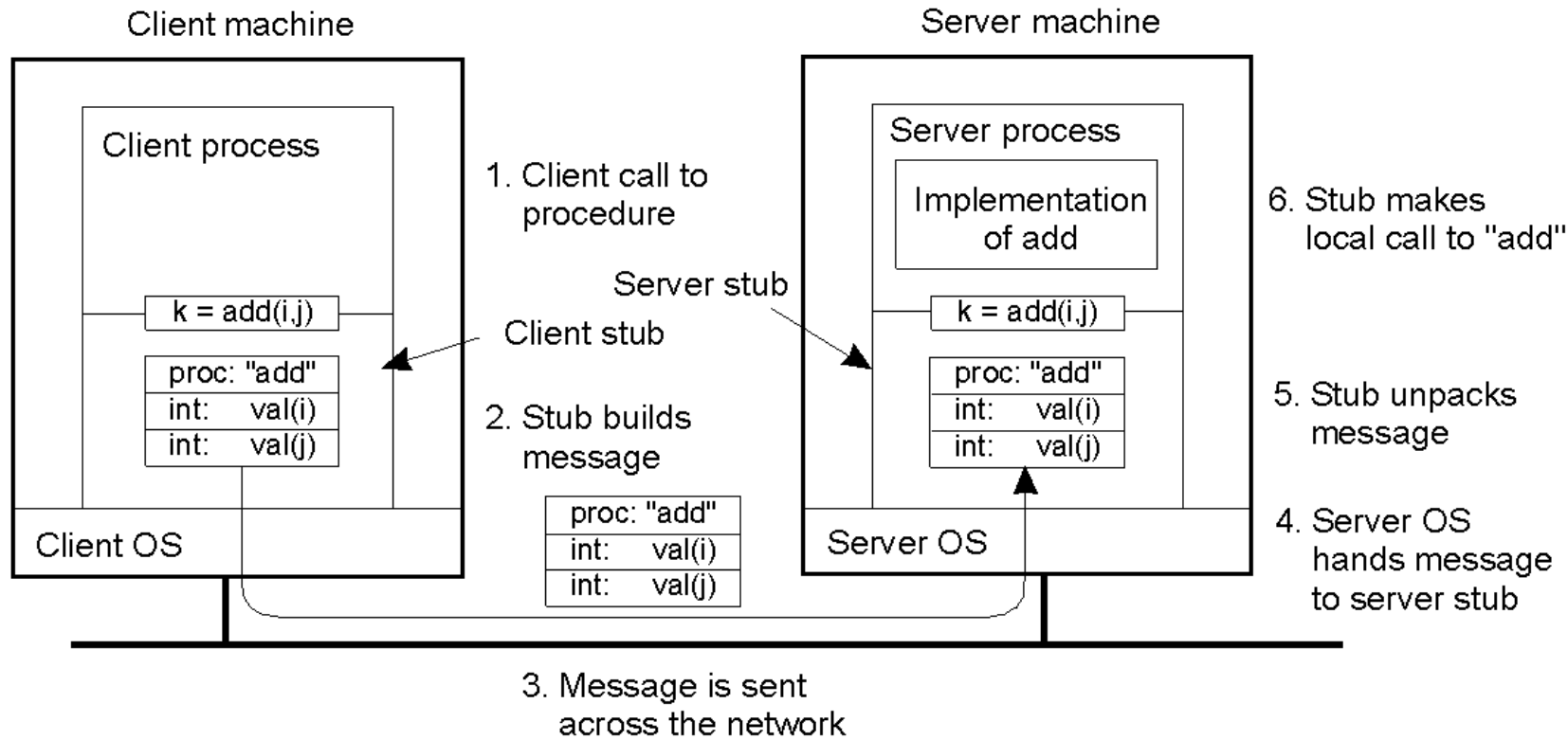
- Usually, RPCs are blocking
 - Thus, also useful for synchronization



How RPCs Work

- Each process has 2 additional components:
 - Code stubs
 - RPC runtime
- Code stubs “translate” local calls remote calls
 - Pack/unpack parameters
- RPC runtime transmits these translated calls over the network
 - Wait for result

How RPCs Work



Parameter Passing

- Local procedure parameter passing
 - Call-by-value
 - Call-by-reference: arrays, complex data structures
- Remote procedure calls simulate this through:
 - Stubs – proxies
 - Flattening – marshalling
 - Serializing local, in-memory representation
- Related issue: global variables are not allowed in RPCs

Client And Server Stubs

- Client makes procedure call (just like a local procedure call) to the client stub
- Server is written as a standard procedure
- Stubs take care of packaging arguments and sending messages
- Packaging parameters is called *marshalling*
- Stub compiler generates stub automatically from specs in an Interface Definition Language (IDL)
 - Simplifies programmer task

Steps of RPC

1. Client procedure calls client stub in normal way
2. Client stub builds message, calls local OS
3. Client's OS sends message to remote OS
4. Remote OS gives message to server stub
5. Server stub unpacks parameters, calls server
6. Server does work, returns result to the stub
7. Server stub packs it in message, calls local OS
8. Server's OS sends message to client's OS
9. Client's OS gives message to client stub
10. Stub unpacks result, returns to client

Marshalling

- Problem: different machines have different data formats
 - Intel: little endian, SPARC: big endian
- Solution: use a cross-platform, general, standard representation
 - Convert in-memory object representation to a standardized “wire” format
 - Example: external data representation (XDR)
- Problem: how do we pass pointers?
 - If it points to a well-defined data structure, pass a copy and the server stub passes a pointer to the local copy
- What about data structures containing pointers?
 - Prohibit
 - Dereference and send (used by most RPC implementations)
 - Chase pointers over network
- Marshalling: transform parameters/results into a byte stream (serialization of parameters)

Binding

- Problem: how does a client locate a server?
 - How does caller code locate and call the callee
 - Use bindings (similar to how symbols are bound to variables during run-time in local programs)
- Server
 - Export server interface during initialization
 - Send name, version no, unique identifier, handle (address) to binder
- Client
 - First RPC: send message to binder to import server interface
 - Binder: check to see if server has exported interface
 - Return handle and unique identifier to client

Binding Comments

- Binding can be at **run-time**
 - Better handling of partial failures (clients can try other advertised end-points, protocols, etc.)
 - Increased dynamism
- Exporting and importing incurs overheads
- Binder can be a bottleneck
 - Use multiple binders
- Binder can do load balancing

Failure Semantics

- *Client unable to locate server*: return error
- *Lost request messages*: simple timeout mechanisms
- *Lost replies*: timeout mechanisms
 - Make operation idempotent
 - Use sequence numbers, mark retransmissions
- *Server failures*: did failure occur before or after operation?
 - At least once semantics / Idempotent (SUNRPC)
 - At most once
 - No guarantee
 - Exactly once: desirable but difficult to achieve

More Failure Semantics

- *Client failure*: what happens to the server computation?
 - Referred to as an *orphan*
 - *Extermination*: log at client stub and explicitly kill orphans
 - Overhead of maintaining disk logs
 - *Reincarnation*: Divide time into epochs between failures and delete computations from old epochs
 - *Gentle reincarnation*: upon a new epoch broadcast, try to locate owner first (delete only if no owner)
 - *Expiration*: give each RPC a fixed quantum T ; explicitly request extensions
 - Periodic checks with client during long computations

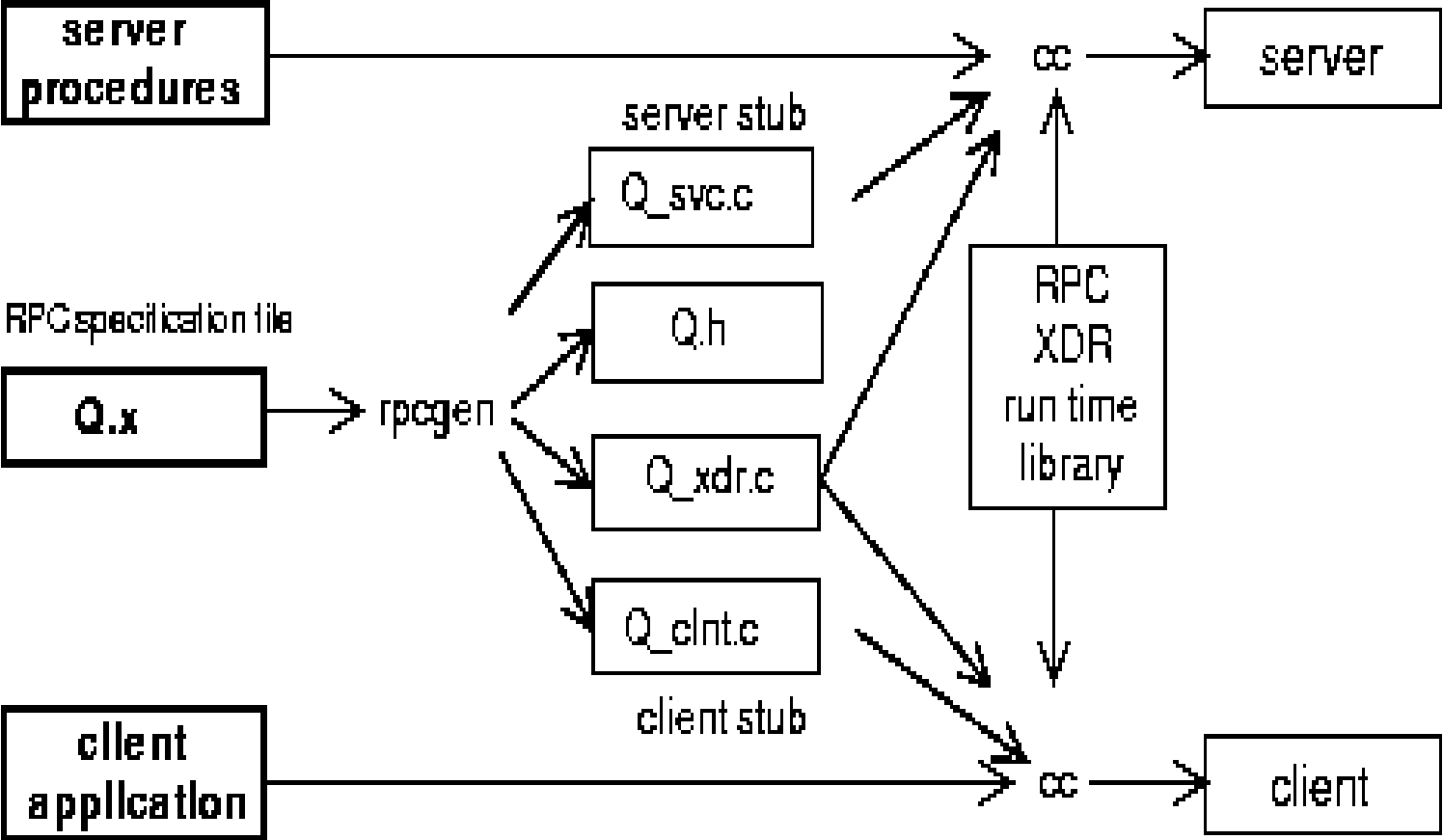
Implementation Issues

- Choice of protocol [affects communication costs]
 - Use existing protocol (UDP) or design from scratch
 - Packet size restrictions
 - Reliability in case of multiple packet messages
 - Flow control
- Copying costs are dominant overheads
 - Need at least 2 copies per message
 - From client to NIC and from server NIC to server
 - As many as 7 copies
 - Stack in stub – message buffer in stub – kernel – NIC – medium – NIC – kernel – stub – server

Sun RPC

- One of the most widely used RPC systems
- Developed for use with NFS (Network File System)
- Built on top of UDP or TCP
 - TCP: stream is divided into records
 - UDP: max packet size < 8912 bytes
 - UDP: timeout plus limited number of retransmissions
 - TCP: return error if connection is terminated by server
- Multiple arguments marshaled into a single structure
- At-least-once semantics if reply received, at-least-zero semantics if no reply. With UDP tries at-most-once
- Use SUN's eXternal Data Representation (XDR)
 - Big endian order for 32 bit integers, handle arbitrarily large data structures

Sun RPC program structure



Modern RPCs and Protocol Buffers

- Many distributed systems use RPCs today (like Mesos)
- Common paradigm: serialize function calls in some serialization format (XML, JSON,...) and send over HTTP (xmlrpc, etc.)
- HTTP servers unpacks and executes the remote call
 - POST <http://foo.com/api/function-name> {arg1:x, arg2:y}

Protocol Buffers

- Relatively new (2008) serialization format from Google
- Binary format. Faster than JSON/XML
 - Con: Not self documenting

```
message Point {  
    int32 x = 1 ; //Field "tags", since names are not included in the message  
    int32 y = 2 ;  
    String name = 3 ;  
}  
Repeated Points point = 4 ; //List/array
```

- Getters and setter methods created for each message during compilation (protoc)
- Access via `msg.fieldname()` (for example, `point.x()`)
- Multiple languages supported

gRPC: A Modern RPC Framework

- “Service” : Function declaration
 - Unary: Single response for a request
 - Streaming: Multiple streaming requests result in single response
- Uses HTTP/2 as transport
 - Messages are just POST requests. Request name is URI, params is content
 - Can multiplex multiple requests onto single TCP connection
- At-most-once failure semantics, but other schemes using retries possible
- Can use load balancers
- GRPC Python: <https://www.semantics3.com/blog/a-simplified-guide-to-grpc-in-python-6c4e25f0c506/>

Summary

- RPCs make distributed computations look like local computations
- Issues:
 - Parameter passing
 - Binding
 - Failure handling
- Case Study: SUN RPC