# Networked Systems

# Distributed system

A collection of physically separate computers working together

- cheaper to build
- easier to add capabilities incrementally

Decade	Technology	\$/machine	Sales volume	Users/machine
50's	custom	\$10M	100	1000's
60's	mainframe	\$1M	10K	100's
70's	minicomputers	\$100K	1M	10's
80's	PCs	\$10K	100M	1
90's	PCs, portables, PDAs	\$1K	1B	1/10
00's	appliances	\$0.1K	10B	1/100
	cloud	\$1K	???	1/1K - 1/10K

# Diamonds and Rust

- Higher availability
  - □ if one machine fails, another can step in
- Better reliability
  - replication
- More security
  - easier to make each smaller piece secure
  - in cloud, professional system management

# Diamonds and Rust

- Higher availability
   if one machine fails, another can step in
- Better reliability
   replication
- More security
  - easier to make each smaller piece secure
  - in cloud, professional system management

may stop if any machine fails

Lower availability

# What is a distributed system?

"A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable."

Leslie Lamport



# Diamonds and Rust

Higher availability

- if one machine fails, another can step in
- Better reliability
  - redundancy through replication
- More security
  - easier to make each smaller piece secure
- 🗅 in cloud, professional system management

### S Lower availability

🗆 may stop if any machine fails

### 0

 $\square$  hard to coordinate replicas

### Less security

□ anyone in the world can break into the system

# Message transmission

- OS sees network as just another device
  - D Network Interface Controller (NIC) added to bus
  - transfer data to/from memory to NIC through DMA or memory mapped I/O



# Networks

### Ethernet

Xerox PARC, 1973

Host Host Host

Host Host

Hub 100 Mb/

□ co-invented by Bob Metcalfe, now at UT!

bandwidth: from 3Mb/s to 1 Gb/s

hub copies every bit it receives on a port to every other port multiple segments can be connected into

a larger LAN (Local Area Network)

Host Host Hub

Host Host

Host Host Host

Multiple LAN can be connected in a WAN (Wide Area Network)



# Routing

- Need to send a message to the right process on the right host
- @ Each host has a unique network ID (e.g. IP address 128.83.120.156)
  - a process can create a port on the host
    - ▶ UTCS web server: 128.83.120.39:80
- Sender finds IP address of intended receiver through Domain Name Service (DNS)
  - distributed hierarchical database
  - naps names such as nikon.cs.utexas.edu to 128.83.120.156

# Distance Vector Routing



# Address Resolution Protocol

Each adapter has a unique link layer address (MAC)

- 6 bytes, burned in ROM
- D ARP (Address Resolution Protocol) maps between link and network address (MAC and IP)



- What if Host 1 want to send a message to Host 4?
- 1) Host: contacts DNS to get IP address of Host: 2) adds TCP header specifies port 3) adds IP header 4) contacts DNS to learn router's IP address 5) uses ARP to get router's MAC address and adds header
- 1) Router strips MAC header 2) checks IP header to determine interface on which to forward > 222.222.222.220 3) uses ARP to find MAC of Host 4) adds header and forwards

# Message Loss

### Buffers overflow, interference, etc.

### Solution 1: use acks

- a) sender sends (msg, msgId) and sets timer
- b) receiver receives message and sends (ack, msgId)
- c) sender receives (ack, msgId) and clears timer
- if timer goes off, go to a) □ "at least once" semantics
- Low throughput
  - □ 1 packet/roundtrip latency ▷ 1KB per 10ms = 100 KB/s

application response as implicit ack

timeout

Optimizations

Solution2: pipeline Solution 1

multiple packets in flight

ack i acks packets 1 to i

resend unacked packets after

Naale's algorithm

delayed acks

combine small packets to reduce overheads

immediate resend on nack (or repeated acks to previously acked message)

as long as there is a sent packet for which sender has not received ack, buffe output until packet is full

for bidirectional communication, use

- ▶ good for telnet, bad for real-time applications

# Congestion: causes and costs

- Two senders, a router with infinite buffer
  - senders send at same rate
  - 🛛 linK capacity R



- Two senders and a router, with finite buffer
  - needed retransmissions due to overflow reduce effective throughput
  - unneeded retransmissions caused by large delays reduce effective throughput
- Multiple senders, multihop paths
  - dropped packets waist resources used to forward them to the place they are dropped

# **TCP Flow Control**

- Prevents sender from overflowing receiver's buffer
  - Sender maintains
  - RcvWindow estimate of buffer space at receiver
  - LastByteSent, LastByteAcked
- Receiver maintains
- 🗆 RcvBuffer
- 🗅 LastByteRead, LastByteRcvd
- To avoid overflowing
  - □ LastByteRcvd-LastByteRead ≤ RcvBuffer
- Hence, sender ensures
  - LastByteSent LastByteAcked ≤ RcvWindow = RcvBuffer (LastByteRcvd LastByteRead)



# **TCP** Congestion Control

- Both sides of a connection keep track of CongWin
  - □ LastByteRcvd-LastByteRead ≤ min{CongWin, RcvBuffer}
  - send rate: CongWin/RTT (assuming negligible retransmission and loss)
- @ If acks are received regularly, CongWin grows
- If ack loss is detected, CongWin shrinks
  - ack loss detected as
  - ⊳ timeout
    - ▹ receipt of 3 duplicate acks
- Congestion Control algorithm has five major components
  - 🗅 additive increase, multiplicative decrease
  - slow start
  - reaction to timeout events
  - round trip variance estimation
  - exponential retransmit timer backoff

# Additive increase, multiplicative decrease

- Multiplicative decrease
  - half CongWin after loss event\* (until one reaches 1 MSS)
- Additive increase
  - increase CongWin by 1 MSS every roundtrip
    - ▷ on each received ack, increase by MSS x (MSS/CongWin)



# Slow Start

- At start, CongWin set to 1 MSS
- Too slow to grow linearly
  - D during slow start phase, exponential growth...
  - ▶ CongWin grows by 1 MSS for each received ack
  - □ until first loss event occurs

# Reaction to time out events

TCP actually treats loss events differently

- $\hfill\square$  on receipt of three duplicate acks
  - multiplicative decrease, additive increase
- on a timeout
  - ▷ drop to 1 MSS
  - slow start mode until CongWin reaches a threshold (typically half of CongWin before loss)
  - » additive increase after threshold reached

### Why the difference?

- $\hfill\square$  old versions of TCP (TCP Tahoe) resorted to slow start also on receipt of duplicate acks, but...
- 🛛 …receipt of acks, even if duplicate, shows that <u>some</u> packets get to destination
  - ▶ fast recovery eliminates slow start phase (TCP Reno)

### Recent proposals (TCP Vegas)

- detect congestion in the router <u>before</u> packet loss occurs
- lower rate linearly when imminent packet loss is detected
  - Ionger RTT indicates greater congestion

# RTT variance and retransmit backoff

- Original TCP set timeout to twice the estimated RTT
   but with high load (above 75%) RTT can vary by 16X
   lots of unnecessary sends under load (great...)
- Current versions set timeout to
   a estimated RTT + 4 x MeanDev(RTT)
- Retransmit backoff is exponential
  - provably needed for stability
  - 1 this is why web browser stalls for 5 sec, then for 10 then...
    - ▶ hint: hit reload if page no there after 5 sec

# TCP friendly rate control

Measure loss probability L and RTT

Solution Use them as parameters to model TCP throughput

 TCP friendly protocols use a congestion control mechanism that consumes no more bandwidth than

> 1.22 x MSS RTT x √L

Average throughput of a connection =  $\frac{1.22 \times MSS}{RTT \times \sqrt{L}}$ 

# Process coordination: two fundamental approaches

Process 1 Communication and ٤. 5 synchronization based on... Stack Stack shared memory assume processes/threads can read & write a set of shared memory Globals locations Heap implicit inter-process communication explicit synchronization Code b difficult to provide across machine boundaries message passing explicit inter process communication process process ▷ implicit synchronization send(message) receive(message)

# Send and Receive Primitives



# Semantics of message passing

### send(receiver,message)

### Synchronization

		Blocking	Non-blocking
Naming	Explicit (single)	Send message to receiver Wait until message is accepted	Send message to receiver
	Implicit (group)	Broadcast message to all receivers. Wait until message is accepted by all	Broadcast message to all receivers

# Semantics of message passing

receive(sender,message)

### Synchronization

	Blocking		Non-blocking	
Naming	Explicit (single)	Wait for a message from sender	If there is a message from sender then receive it, else continue	
	Implicit (group)	Wait for a message from any sender	If there is a message from any sender then receive it, else continue	

# Buffering messages

### No buffering

- sender must wait until receiver receives message
- rendezvous on each message
- Bounded buffer
  - 🛛 finite size
  - sender blocks on buffer full
- Onbounded buffer
  - "infinite" size
  - sender never blocks

# Direct Communication

### A single buffer at the receiver

- multiple processes may send messages to the receiver
- receiving from a specific sender requires searching entire buffer
- a A buffer at each sender
  - sender may send messages to multiple receivers
  - receiving a message requires searching through whole buffer



# Indirect Communication

- Mailbox abstraction
  - many-to-many communication
  - requires open/close of mailbox
- Buffering
  - buffer, mutex and condition variables at the mailbox



# Limitations of message passing

- Easy for OS, hard for programmer
  - programmer must code synchronization
  - programmer may have to code format conversions, flow control, error control
  - no dynamic resource discovery

## Remote Procedure Call (RPC) Birrell & Nelson, 1984

- RPC mechanism
  - □ hides message passing I/O from programmer
  - □ (almost) a procedure call-but on the server
- RPC invocation
  - calling process (client) is suspended
  - procedure parameters passed across network to called process (server)
  - return parameters send back across network
  - calling process resumes



# More about RPC

- Similarities between procedure call and RPC
  - $\square$  paramenters  $\leftrightarrow$  request message
  - $\square$  result  $\leftrightarrow$  reply message
  - $\square$  name of procedure  $\leftrightarrow$  passed in request message
  - $\square$  return address  $\leftrightarrow$  mailbox of client
- Implementation issues
  - stub generation
    - types of arguments and return value ▷ can be automated
    - requires signature of the procedure
  - binding (how client locates a server)
    - static (fixed at compile time)
    - » dynamic (at run time, with the help of a name service)
      - if server fails, automatic fail over

# Problems with RPC

- Achieves location transparency, except for
  - performance
    - ▶ cost of procedure call << same machine RPC << network RPC
  - n failures (message loss, machine crash)

### three execution semantics

### if succeeds: at least one execution if fails: none, partial, multiple if succeeds: exactly once if fails: none, partial, one "exactly once" if succeeds: exactly once

if fails: none,

### implemented by retransmission on timeout works only if operation is idempotent "at least once" (SUN RPC) -\_"at most once" (Java RMI) implemented by server filtering duplicate requests

consistency/replication

can't automatically update an object replicated at multiple machines

impossible

what would I need to do so?

□ security