# CS162
# Operating Systems and
# Systems Programming
# Lecture 22

# Networking II

April 23, 2008
Prof. Anthony D. Joseph
http://inst.eecs.berkeley.edu/~cs162

---

## Review: Point-to-point networks



- **Point-to-point network:** a network in which every physical wire is connected to only two computers
- **Switch:** a bridge that transforms a shared-bus (broadcast) configuration into a point-to-point network.
- **Hub:** a multiport device that acts like a repeater broadcasting from each input to every output
- **Router:** a device that acts as a junction between two networks to transfer data packets among them.

---

## Sequence Numbers

- **Ordered Messages**
  - Several network services are best constructed by ordered messaging
    » Ask remote machine to first do x, then do y, etc.
  - Unfortunately, underlying network is packet based:
    » Packets are routed one at a time through the network
    » Can take different paths or be delayed individually
  - IP can reorder packets! $P_0, P_1$ might arrive as $P_1, P_0$
- **Solution: Queue out of order packets at destination**
  - Need to hold onto packets to undo misordering
  - Total degree of reordering impacts queue size
- Ordered messages on top of unordered ones:
  - Assign sequence numbers to packets
    » 0,1,2,3,4…..
    » If packets arrive out of order, reorder before delivering to user application
    » For instance, hold onto #3 until #2 arrives, etc.
  - Sequence numbers are specific to particular connection
    » Reordering among connections normally doesn't matter
  - If restart connection, need to make sure use different range of sequence numbers than previously…

---

## Goals for Today

- Networking
  - Protocols
  - Reliable Messaging
    » TCP windowing and congestion avoidance
  - Two-phase commit

Note: Some slides and/or pictures in the following are adapted from slides ©2005 Silberschatz, Galvin, and Gagne. Many slides generated from my lecture notes by Kubiatowicz.

---

## Reliable Message Delivery: the Problem
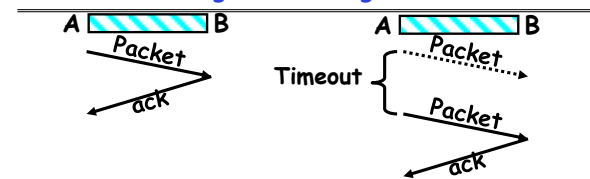
- **All physical networks can garble and/or drop packets**
  - Physical media: packet not transmitted/received
    - » If transmit close to maximum rate, get more throughput – even if some packets get lost
    - » If transmit at lowest voltage such that error correction just starts correcting errors, get best power/bit
  - Congestion: no place to put incoming packet
    - » Point-to-point network: insufficient queue at switch/router
    - » Broadcast link: two host try to use same link
    - » In any network: insufficient buffer space at destination
    - » Rate mismatch: what if sender sends faster than receiver can process?
- **Reliable Message Delivery**
  - Reliable messages on top of unreliable packets
  - Need some way to make sure that packets actually make it to receiver
    - » Every packet received at least once
    - » Every packet received only once
  - Can combine with ordering: every packet received by process at destination exactly once and in order

## Using Acknowledgements



- **How to ensure transmission of packets?**
  - Detect garbling at receiver via checksum, discard if bad
  - Receiver acknowledges (by sending "ack") when packet received properly at destination
  - Timeout at sender: if no ack, retransmit
- **Some questions:**
  - If the sender doesn't get an ack, does that mean the receiver didn't get the original message?
    - » No
  - What it ack gets dropped? Or if message gets delayed?
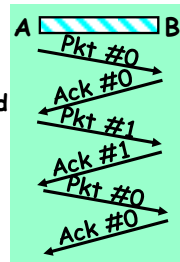    - » Sender doesn't get ack, retransmits. Receiver gets message twice, acks each.

## How to deal with message duplication

- **Solution: put sequence number in message to identify re-transmitted packets**
  - Receiver checks for duplicate #'s; Discard if detected
- **Requirements:**
  - Sender keeps copy of unack'ed messages
    - » Easy: only need to buffer messages
  - Receiver tracks possible duplicate messages
    - » Hard: when ok to forget about received message?
- **Alternating-bit protocol:**
  - Send one message at a time; don't send next message until ack received
  - Sender keeps last message; receiver tracks sequence # of last message received
- **Pros: simple, small overhead**
- **Con: Poor performance**
  - Wire can hold multiple messages; want to fill up at (wire latency × throughput)
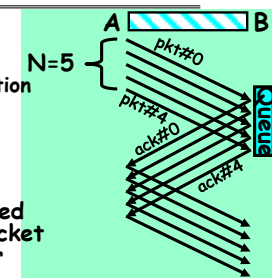- **Con: doesn't work if network can delay or duplicate messages arbitrarily**

## Better messaging: Window-based acknowledgements

- **Window based protocol (TCP):**
  - Send up to N packets without ack
    - » Allows pipelining of packets
    - » Window size (N) < queue at destination
  - Each packet has sequence number
    - » Receiver acknowledges each packet
    - » Ack says "received all packets up to sequence number X"/send more
- **Acks serve dual purpose:**
  - Reliability: Confirming packet received
  - Flow Control: Receiver ready for packet
    - » Remaining space in queue at receiver can be returned with ACK
- **What if packet gets garbled/dropped?**
  - Sender will timeout waiting for ack packet
    - » Resend missing packets⇒ Receiver gets packets out of order!
  - Should receiver discard packets that arrive out of order?
    - » Simple, but poor performance
  - Alternative: Keep copy until sender fills in missing pieces?
    - » Reduces # of retransmits, but more complex
- **What if ack gets garbled/dropped?**
  - Timeout and resend just the un-acknowledged packets

## Transmission Control Protocol (TCP)

**Stream in:**  ...zyxwvuts → [Router] — [Router] → **Stream out:** gfedcba
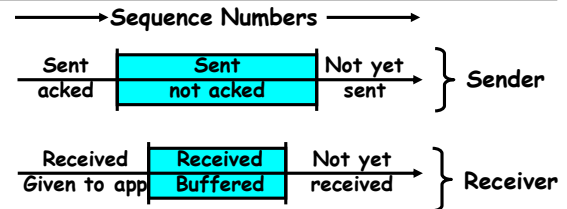
- Transmission Control Protocol (TCP)
  - TCP (IP Protocol 6) layered on top of IP
  - Reliable byte stream between two processes on different machines over Internet (read, write, flush)
- TCP Details
  - Fragments byte stream into packets, hands packets to IP
    » IP may also fragment by itself
  - Uses window-based acknowledgement protocol (to minimize state at sender and receiver)
    » "Window" reflects storage at receiver – sender shouldn't overrun receiver's buffer space
    » Also, window should reflect speed/capacity of network – sender shouldn't overload network
  - Automatically retransmits lost packets
  - Adjusts rate of transmission to avoid congestion
    » A "good citizen"

---

## TCP Windows and Sequence Numbers

→ Sequence Numbers →

| Sent acked | Sent / not acked | Not yet sent | } Sender |

| Received Given to app | Received / Buffered | Not yet received | } Receiver |

- Sender has three regions:
  - Sequence regions
    » sent and ack'ed
    » Sent and not ack'ed
    » not yet sent
  - Window (colored region) adjusted by sender
- Receiver has three regions:
  - Sequence regions
    » received and ack'ed (given to application)
    » received and buffered
    » not yet received (or discarded because out of order)

---

## Window-Based Acknowledgements (TCP)

100    140    190    230    260    300    340    380  400

Seq:100 Size:40 | Seq:140 Size:50 | Seq:190 Size:40 | Seq:230 Size:30 | Seq:260 Size:40 | Seq:300 Size:40 | Seq:340 Size:40 | Seq:380 Size:20

Seq:100 → A:100/300
Seq:140 → A:140/260
Seq:190 → A:190/210
Seq:230 → A:190/140
Seq:260 → A:190/100
Seq:300 → A:190/60
Seq:190 Retransmit! → A:340/60
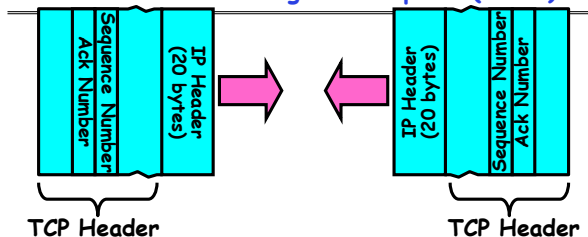Seq:340 → A:380/20
Seq:380 → A:400/0

---

## Administrivia

- Project #4 design deadline is Thu 5/1 at 11:59pm
  - You need to create an account (ASAP!)
    » Let us know if you have problems creating an account
  - Code deadline is Wed 5/14

- Final Exam – May 21st, 12:30-3:30pm
  - Email conflicts to **cs162@cory** by Wed 4/23 at 5pm

- Final Topics: Any suggestions?
  - Please send them to me…

- Thank you for the anonymous web comments!

## Selective Acknowledgement Option (SACK)



TCP Header       TCP Header

- Vanilla TCP Acknowledgement
  - Every message encodes Sequence number and Ack
  - Can include data for forward stream and/or ack for reverse stream
- Selective Acknowledgement
  - Acknowledgement information includes not just one number, but rather ranges of received packets
  - Must be specially negotiated at beginning of TCP setup
    » Not widely in use (although in Windows since Windows 98)

4/23/08     Joseph CS162 ©UCB Spring 2008     Lec 22.13

---

## Congestion Avoidance

- Congestion
  - How long should timeout be for re-sending messages?
    » Too long→wastes time if message lost
    » Too short→retransmit even though ack will arrive shortly
  - Stability problem: more congestion ⇒ ack is delayed ⇒ unnecessary timeout ⇒ more traffic ⇒ more congestion
    » Closely related to window size at sender: too big means putting too much data into network
- How does the sender's window size get chosen?
  - Must be less than receiver's advertised buffer size
  - Try to match the rate of sending packets with the rate that the slowest link can accommodate
  - Sender uses an adaptive algorithm to decide size of N
    » Goal: fill network between sender and receiver
    » Basic technique: slowly increase size of window until acknowledgements start being delayed/lost
- TCP solution: "slow start" (start sending slowly)
  - If no timeout, slowly increase window size (throughput) by 1 for each ack received
  - Timeout ⇒ congestion, so cut window size in half
  - "*Additive Increase, Multiplicative Decrease*"

4/23/08     Joseph CS162 ©UCB Spring 2008     Lec 22.14

---

## Sequence-Number Initialization

- How do you choose an initial sequence number?
  - When machine boots, ok to start with sequence #0?
    » No: could send two messages with same sequence #!
    » Receiver might end up discarding valid packets, or duplicate ack from original transmission might hide lost packet
  - Also, if it is possible to predict sequence numbers, might be possible for attacker to hijack TCP connection
- Some ways of choosing an initial sequence number:
  - Time to live: each packet has a deadline.
    » If not delivered in X seconds, then is dropped
    » Thus, can re-use sequence numbers if wait for all packets in flight to be delivered or to expire
  - Epoch #: uniquely identifies *which* set of sequence numbers are currently being used
    » Epoch # stored on disk, Put in every message
    » Epoch # incremented on crash and/or when run out of sequence #
  - Pseudo-random increment to previous sequence number
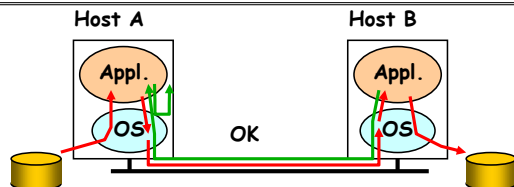    » Used by several protocol implementations

4/23/08     Joseph CS162 ©UCB Spring 2008     Lec 22.15

---

## "End-to-End Arguments in System Design" (Saltzer, Reed, and Clark)

- Most influential paper about placing functionality
  - "Sacred Text" of the Internet
    » Endless disputes about what it means
    » Everyone cites it as supporting their position
- Some applications have end-to-end performance requirements:
  - Reliability, security, …

- Implementing these in the network is very hard:
  - Every step along the way must be fail-proof

- Hosts:
  - Can satisfy the requirement without the network
  - Can't depend on the network

4/23/08     Joseph CS162 ©UCB Spring 2008     Lec 22.16

## Example: Reliable File Transfer

Host A          Host B



- **Solution 1: make each step reliable, and then concatenate them**
  - Solution 1 not complete (e.g., misbehaving net element)
    » What happens if any network element misbehaves?
    » Receiver has to do the check anyway!
- **Solution 2: end-to-end check and retry**
  - Solution 2 is complete
    » Full functionality can be entirely implemented at application layer with **no** need for reliability from lower layers

## E2E Summary

- **Implementing this functionality in the network:**
  - Doesn't reduce host implementation complexity
  - Does increase network complexity
  - Probably imposes delay and overhead on all applications, even if they don't need functionality

- **However, implementing in network can enhance performance in some cases**
  - Such as a very lossy link

- **Layering is a good way to organize networks**
  - Unified Internet layer decouples apps from networks
  - E2E argument encourages us to keep IP simple
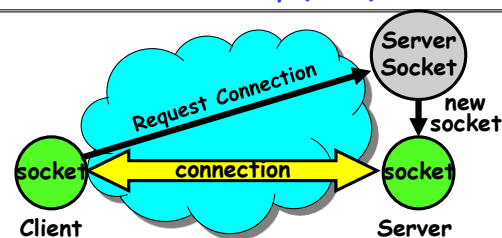  - Commercial realities may undo all of this...

## Use of TCP: Sockets

- **Socket: an abstraction of a network I/O queue**
  - Embodies one side of a communication channel
    » Same interface regardless of location of other end
    » Could be local machine (called "UNIX socket") or remote machine (called "network socket")
  - First introduced in 4.2 BSD UNIX: big innovation at time
    » Now most operating systems provide some notion of socket
- **Using Sockets for Client-Server (C/C++ interface):**
  - On server: set up "server-socket"
    » Create socket, Bind to protocol (TCP), local address, port
    » Call listen(): tells server socket to accept incoming requests
    » Perform multiple accept() calls on socket to accept incoming connection request
    » Each successful accept() returns a new socket for a new connection; can pass this off to handler thread
  - On client:
    » Create socket, Bind to protocol (TCP), remote address, port
    » Perform connect() on socket to make connection
    » If connect() successful, have socket connected to server

## Socket Setup (Con't)



- **Things to remember:**
  - Connection requires 5 values:
    [ Src Addr, Src Port, Dst Addr, Dst Port, Protocol ]
  - Often, Src Port "randomly" assigned
    » Done by OS during client socket setup
  - Dst Port often "well known"
    » 80 (web), 443 (secure web), 25 (sendmail), etc
    » Well-known ports from 0—1023

# BREAK

---

## Socket Example (Java)

```
server:
     //Makes socket, binds addr/port, calls listen()
     ServerSocket sock = new ServerSocket(6013);
     while(true) {
        Socket client = sock.accept();
        PrintWriter pout = new
           PrintWriter(client.getOutputStream(),true);

        pout.println("Here is data sent to client!");
        …
        client.close();
     }

client:
     // Makes socket, binds addr/port, calls connect()
     Socket sock = new Socket("169.229.60.38",6013);
     BufferedReader bin =
        new BufferedReader(
           new InputStreamReader(sock.getInputStream));
     String line;
     while ((line = bin.readLine())!=null)
        System.out.println(line);
     sock.close();
```

---

## Distributed Applications

- **How do you actually program a distributed application?**
  - Need to synchronize multiple threads, running on different machines
    » No shared memory, so cannot use test&set



  - One Abstraction: send/receive messages
    » Already atomic: no receiver gets portion of a message and two receivers cannot get same message
- **Interface:**
  - Mailbox (`mbox`): temporary holding area for messages
    » Includes both destination location and queue
  - `Send(message,mbox)`
    » Send message to remote mailbox identified by `mbox`
  - `Receive(buffer,mbox)`
    » Wait until `mbox` has message, copy into buffer, and return
    » If threads sleeping on this `mbox`, wake up one of them

---

## Using Messages: Send/Receive behavior

- **When should `send(message,mbox)` return?**
  - When receiver gets message? (i.e. ack received)
  - When message is safely buffered on destination?
  - Right away, if message is buffered on source node?
- **Actually two questions here:**
  - When can the sender be sure that the receiver actually received the message?
  - When can sender reuse the memory containing message?
- **Mailbox provides 1-way communication from T1→T2**
  - T1→buffer→T2
  - Very similar to producer/consumer
    » Send = V, Receive = P
    » However, can't tell if sender/receiver is local or not!

---

## Messaging for Producer-Consumer Style

- **Using send/receive for producer-consumer style:**

```
Producer:
  int msg1[1000];
  while(1) {
    prepare message;
    send(msg1,mbox);
  }
Consumer:
  int buffer[1000];
  while(1) {
    receive(buffer,mbox);
    process message;
  }
```

**Send Message**

**Receive Message**

- **No need for producer/consumer to keep track of space in mailbox: handled by send/receive**
  - **One of the roles of the window in TCP: window is size of buffer on far end**
  - **Restricts sender to forward only what will fit in buffer**

## Messaging for Request/Response communication

- **What about two-way communication?**
  - **Request/Response**
    - » **Read a file stored on a remote machine**
    - » **Request a web page from a remote web server**
  - **Also called: client-server**
    - » **Client ≡ requester, Server ≡ responder**
    - » **Server provides "service" (file storage) to the client**
- **Example: File service**

```
Client: (requesting the file)
  char response[1000];

  send("read rutabaga", server_mbox);
  receive(response, client_mbox);

Server: (responding with the file)
  char command[1000], answer[1000];

  receive(command, server_mbox);
  decode command;
  read file into answer;
  send(answer, client_mbox);
```

**Request File**

**Get Response**

**Receive Request**

**Send Response**

## Conclusion

- **Layering:** building complex services from simpler ones

- **Ordered messages:**
  - **Use sequence numbers and reorder at destination**

- **Reliable messages:**
  - **Use Acknowledgements**
  - **Want a window larger than 1 in order to increase throughput**

- **TCP:** Reliable byte stream between two processes on different machines over Internet (read, write, flush)