

Transport Protocols



- UDP provides just integrity and demux
- TCP adds...
 - Connection-oriented
 - Reliable
 - Ordered
 - Byte-stream
 - Full duplex
 - · Flow and congestion controlled
- DCCP, RTP, SCTP -- not widely used.

UDP: User Datagram Protocol [RFC 768]



- "No frills," "bare bones" Internet transport protocol
- "Best effort" service, UDP segments may be:
 - Lost
 - · Delivered out of order to app
- Connectionless:
 - No handshaking between UDP sender, receiver
 - Each UDP segment handled independently of others

Why is there a UDP?

- No connection establishment (which can add delay)
- Simple: no connection state at sender, receiver
- Small header
- No congestion control: UDP can blast away as fast as desired

UDP, cont.



Dest port #

Checksum

- Often used for streaming multimedia apps
 - Loss tolerant
 - · Rate sensitive
- Other UDP uses (why?):
 - DNS
- Reliable transfer over UDP
 - Must be at application layer
 - · Application-specific error recovery

Length, in bytes of UDP segment, including header

Application data (message)

32 bits

Source port #

→Length

UDP segment format

UDP Checksum



Goal: detect "errors" (e.g., flipped bits) in transmitted segment – optional use!

Sender:

- Treat segment contents as sequence of 16-bit integers
- Checksum: addition (1's complement sum) of segment contents
- Sender puts checksum value into UDP checksum field

Receiver:

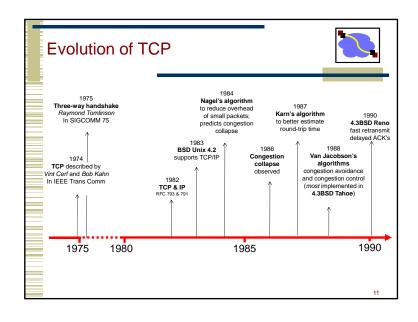
- · Compute checksum of received segment
- · Check if computed checksum equals checksum field value:
 - NO error detected
 - · YES no error detected But maybe errors nonetheless?

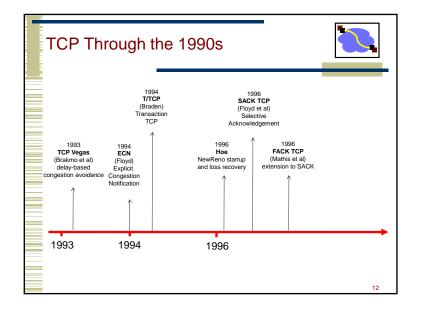
High-Level TCP Characteristics



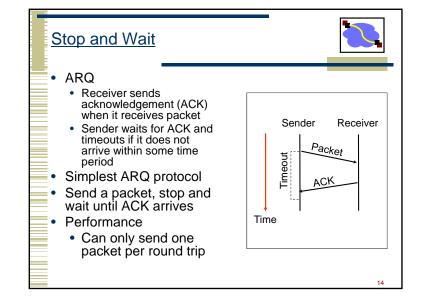
- Protocol implemented entirely at the ends
 - Fate sharing (on IP)
- Protocol has evolved over time and will continue to do so
 - Nearly impossible to change the header
 - Use options to add information to the header
 - · Change processing at endpoints
 - · Backward compatibility is what makes it TCP

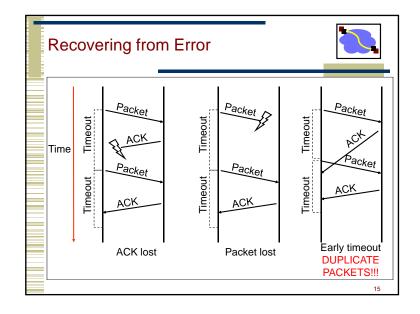
TCP Header Source port Destination port Sequence number Flags: SYN Acknowledgement FIN RESET HdrLen Flags Advertised window 0 **PUSH** Checksum Urgent pointer URG ACK Options (variable) Data

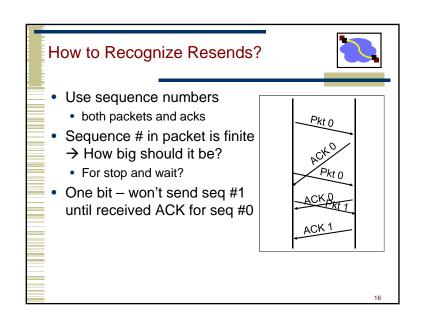


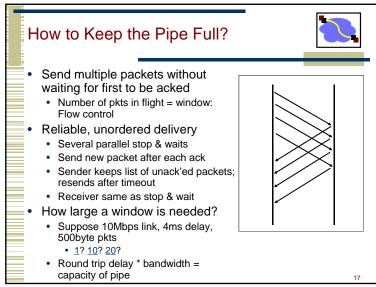


Transport introduction Error recovery & flow control TCP flow control/connection setup/data transfer TCP reliability Congestion sources and collapse Congestion control basics









Sender/Receiver State

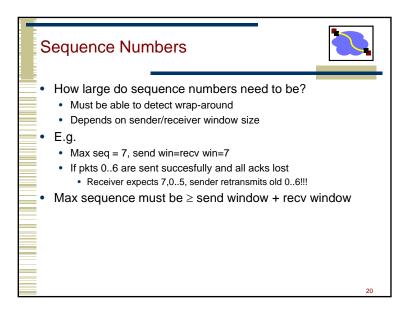
Sender

Not Usable

OK to Send

Receiver Received & Acked Acceptable Packet Not Usable

Sliding Window · Reliable, ordered delivery • Receiver has to hold onto a packet until all prior packets have arrived • Why might this be difficult for just parallel stop & wait? • Sender must prevent buffer overflow at receiver Circular buffer at sender and receiver Packets in transit ≤ buffer size Advance when sender and receiver agree packets at beginning have been received



Window Sliding - Common Case



- On reception of new ACK (i.e. ACK for something that was not acked earlier)
 - Increase sequence of max ACK received
 - · Send next packet
- On reception of new in-order data packet (next expected)
 - Hand packet to application
 - Send cumulative ACK acknowledges reception of all packets up to sequence number
 - Increase sequence of max acceptable packet

21

Loss Recovery



- On reception of out-of-order packet
 - Send nothing (wait for source to timeout)
 - Cumulative ACK (helps source identify loss)
- Timeout (Go-Back-N recovery)
 - Set timer upon transmission of packet
 - Retransmit all unacknowledged packets
- Performance during loss recovery
 - No longer have an entire window in transit
 - Can have much more clever loss recovery

22

Important Lessons



- Transport service
 - UDP → mostly just IP service
 - TCP \rightarrow congestion controlled, reliable, byte stream
- Types of ARQ protocols
 - Stop-and-wait → slow, simple
 - Go-back-n → can keep link utilized (except w/ losses)
 - Selective repeat → efficient loss recovery -- used in SACK
- Sliding window flow control
 - Addresses buffering issues and keeps link utilized

23

Good Ideas So Far...



- Flow control
 - Stop & wait
 - Parallel stop & wait
 - Sliding window
- Loss recovery
 - Timeouts
 - Acknowledgement-driven recovery (selective repeat or cumulative acknowledgement)

Outline



- · Transport introduction
- Error recovery & flow control
- TCP flow control/connection setup/data transfer
- TCP reliability
- Congestion sources and collapse
- Congestion control basics

25

More on Sequence Numbers



- 32 Bits, Unsigned → for bytes not packets!
- Why So Big?
 - · For sliding window, must have
 - |Sequence Space| > |Sending Window| + |Receiving Window|
 - No problem
 - Also, want to guard against stray packets
 - With IP, packets have maximum lifetime of 120s
 - Sequence number would wrap around in this time at 286Mbps

26

TCP Flow Control



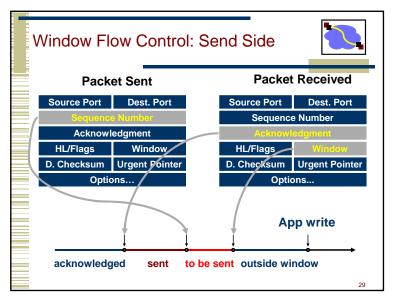
- TCP is a sliding window protocol
 - For window size *n*, can send up to *n* bytes without receiving an acknowledgement
 - When the data is acknowledged then the window slides forward
- Each packet advertises a window size
 - Indicates number of bytes the receiver has space for
- Original TCP always sent entire window
 - Congestion control now limits this

Window Flow Control: Send Side

window

Sent and acked Sent but not acked Not yet sent

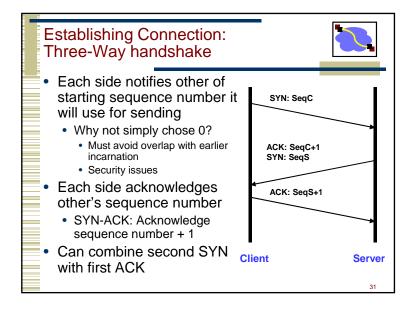
Next to be sent

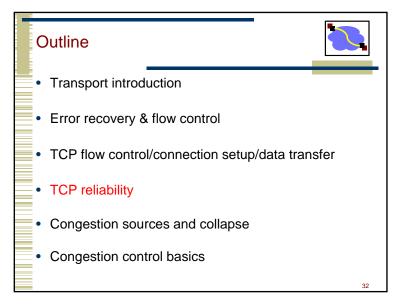


Performance Considerations



- The window size can be controlled by receiving application
 - Can change the socket buffer size from a default (e.g. 8Kbytes) to a maximum value (e.g. 64 Kbytes)
- The window size field in the TCP header limits the window that the receiver can advertise
 - 16 bits → 64 KBytes
 - 10 msec RTT → 51 Mbit/second
 - 100 msec RTT → 5 Mbit/second
 - TCP options to get around 64KB limit → scales window size





Reliability Challenges



- Congestion related losses
- · Variable packet delays
 - What should the timeout be?
- · Reordering of packets
 - How to tell the difference between a delayed packet and a lost one?

33

TCP = Go-Back-N Variant



- Sliding window with cumulative acks
 - Receiver can only return a single "ack" sequence number to the sender.
 - · Acknowledges all bytes with a lower sequence number
 - · Starting point for retransmission
 - · Duplicate acks sent when out-of-order packet received
- But: sender only retransmits a single packet.
 - Reason???
 - · Only one that it knows is lost
 - Network is congested → shouldn't overload it
- Error control is based on byte sequences, not packets.
 - Retransmitted packet can be different from the original lost packet – Why?

34

Round-trip Time Estimation



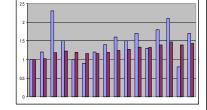
- Wait at least one RTT before retransmitting
- Importance of accurate RTT estimators:
 - · Low RTT estimate
 - · unneeded retransmissions
 - · High RTT estimate
 - poor throughput
- RTT estimator must adapt to change in RTT
 - But not too fast, or too slow!
- Spurious timeouts
 - "Conservation of packets" principle never more than a window worth of packets in flight

.

Original TCP Round-trip Estimator



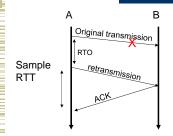
- Round trip times exponentially averaged:
 - New RTT = α (old RTT) + (1 α) (new sample)
 - Recommended value for α: 0.8 - 0.9
 - 0.875 for most TCP's

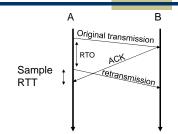


- Retransmit timer set to (b * RTT), where b = 2
 - Every time timer expires, RTO exponentially backed-off
- Not good at preventing premature timeouts

RTT Sample Ambiguity







- · Karn's RTT Estimator
 - If a segment has been retransmitted:
 - Don't count RTT sample on ACKs for this segment
 - · Keep backed off time-out for next packet
 - Reuse RTT estimate only after one successful transmission

37

Jacobson's Retransmission Timeout



- Key observation:
 - At high loads round trip variance is high
- Solution:
 - Base RTO on RTT and standard deviation
 - RTO = RTT + 4 * rttvar
 - new_rttvar = β * dev + (1- β) old_rttvar
 - Dev = linear deviation
 - Inappropriately named actually smoothed linear deviation

38

Timestamp Extension



- Used to improve timeout mechanism by more accurate measurement of RTT
- When sending a packet, insert current time into option
 - 4 bytes for time, 4 bytes for echo a received timestamp
- Receiver echoes timestamp in ACK
 - · Actually will echo whatever is in timestamp
- Removes retransmission ambiguity
 - Can get RTT sample on any packet

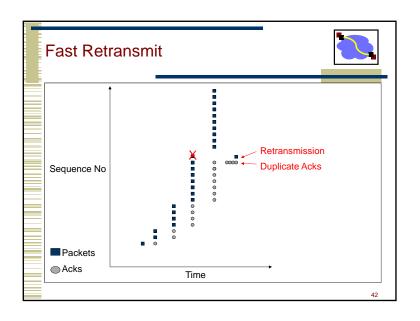
39

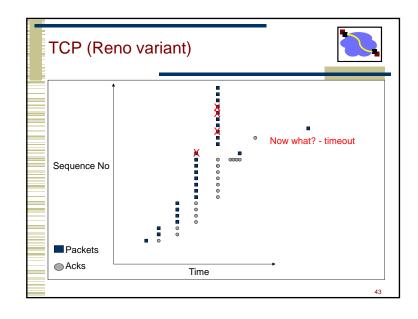
Timer Granularity

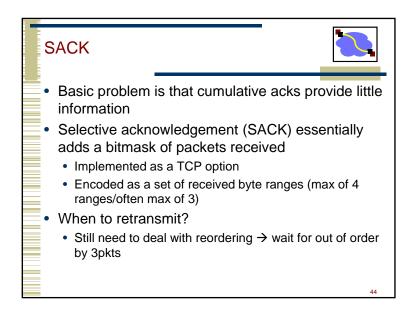


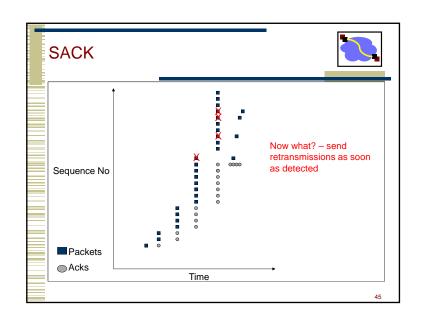
- Many TCP implementations set RTO in multiples of 200,500,1000ms
- Why?
 - Avoid spurious timeouts RTTs can vary quickly due to cross traffic
 - Reduce timer expensive timer interrupts on hosts
- What happens for the first couple of packets?
 - Pick a very conservative value (seconds)

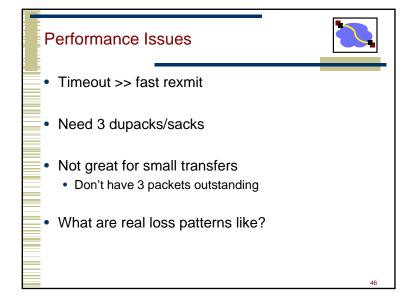
Fast Retransmit -- Avoiding Timeouts • What are duplicate acks (dupacks)? • Repeated acks for the same sequence • When can duplicate acks occur? • Loss • Packet re-ordering • Window update – advertisement of new flow control window • Assume re-ordering is infrequent and not of large magnitude • Use receipt of 3 or more duplicate acks as indication of loss • Don't wait for timeout to retransmit packet

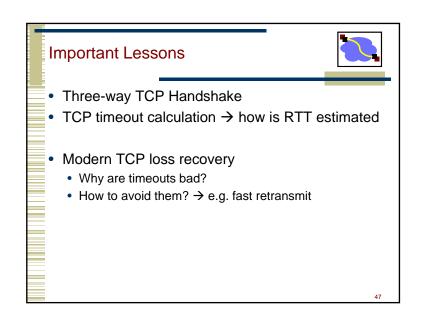


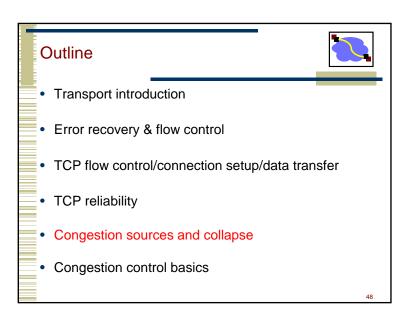


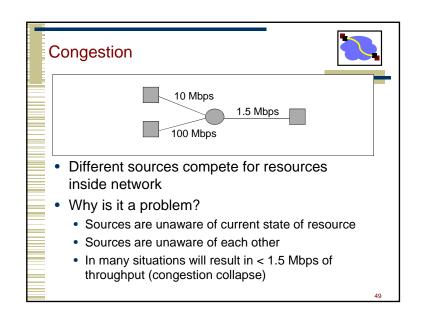


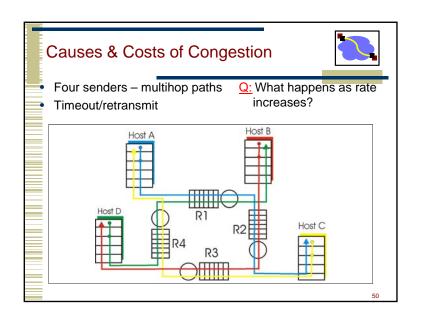


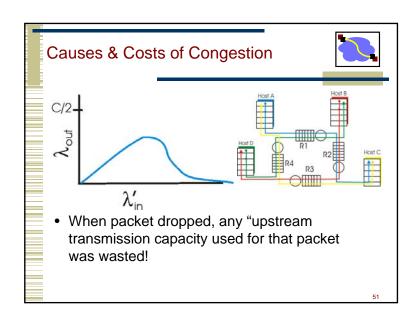


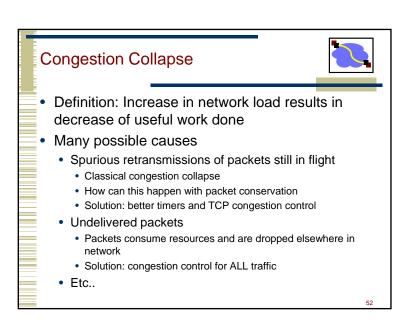












Other Congestion Collapse Causes



- Fragments
 - · Mismatch of transmission and retransmission units
 - Solutions
 - Make network drop all fragments of a packet (early packet discard in ATM)
 - Do path MTU discovery
- Control traffic
 - Large percentage of traffic is for control
 - · Headers, routing messages, DNS, etc.
- Stale or unwanted packets
 - · Packets that are delayed on long queues
 - · "Push" data that is never used

53

Where to Prevent Collapse?



- Can end hosts prevent problem?
 - · Yes, but must trust end hosts to do right thing
 - E.g., sending host must adjust amount of data it puts in the network based on detected congestion
- Can routers prevent collapse?
 - · No, not all forms of collapse
 - · Doesn't mean they can't help
 - · Sending accurate congestion signals
 - Isolating well-behaved from ill-behaved sources

54

Congestion Control and Avoidance



- A mechanism which:
 - · Uses network resources efficiently
 - Preserves fair network resource allocation
 - · Prevents or avoids collapse
- · Congestion collapse is not just a theory
 - Has been frequently observed in many networks

_

Approaches For Congestion Control



Two broad approaches towards congestion control:

End-to-end

Network-assisted

- No explicit feedback from network
- Congestion inferred from end-sys tem observed loss, delay
- Approach taken by TCP
- Routers provide feedback to end systems
 - Explicit rate sender should send at
 - Single bit indicating congestion (SNA, DEC bit, TCP/IP ECN, ATM)
- Problem: makes routers complicated

Example: TCP Congestion Control



- · Very simple mechanisms in network
 - · FIFO scheduling with shared buffer pool
 - · Feedback through packet drops
- TCP interprets packet drops as signs of congestion and slows down
 - This is an assumption: packet drops are not a sign of congestion in all networks
 - . E.g. wireless networks
- Periodically probes the network to check whether more bandwidth has become available.

Outline



- Transport introduction
- Error recovery & flow control
- TCP flow control/connection setup/data transfer
- TCP reliability
- Congestion sources and collapse
- Congestion control basics

Objectives



- Simple router behavior
- Distributedness
- Efficiency: $X_{knee} = \Sigma x_i(t)$
- Fairness: (Σx_i)²/n(Σx_i²)
- Power: (throughput^α/delay)
- Convergence: control system must be stable





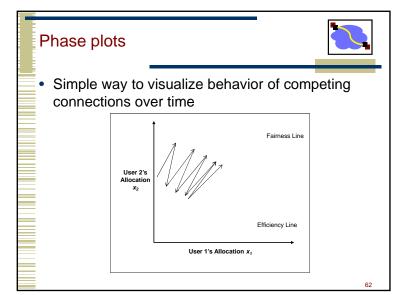
- Let's assume window-based control
- Reduce window when congestion is perceived
 - How is congestion signaled?
 - · Either mark or drop packets
 - When is a router congested?
 - Drop tail queues when queue is full
 - Average queue length at some threshold
- Increase window otherwise
 - Probe for available bandwidth how?

Linear Control



- Many different possibilities for reaction to congestion and probing
 - Examine simple linear controls
 - Window(t + 1) = a + b Window(t)
 - Different a_i/b_i for increase and a_d/b_d for decrease
- Supports various reaction to signals
 - Increase/decrease additively
 - Increased/decrease multiplicatively
 - Which of the four combinations is optimal?

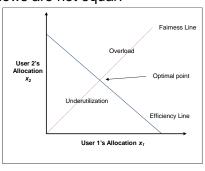
61



Phase plots



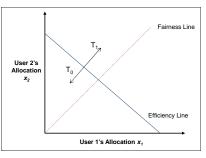
- What are desirable properties?
- What if flows are not equal?

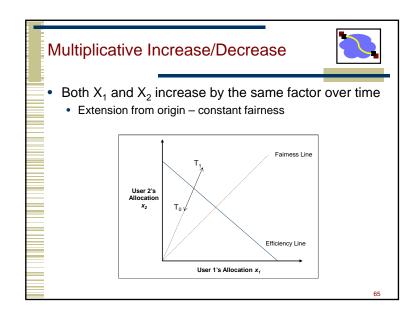


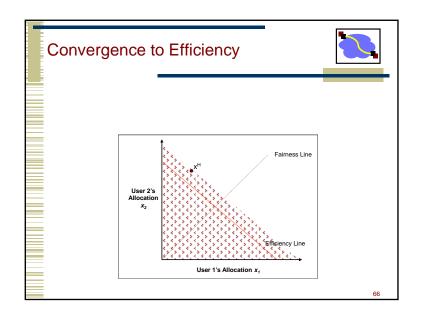
Additive Increase/Decrease

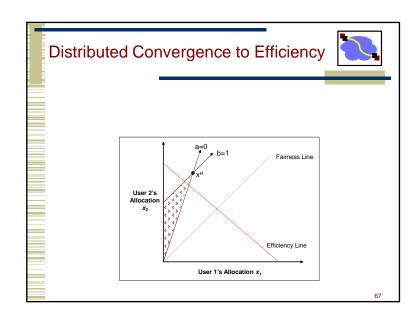


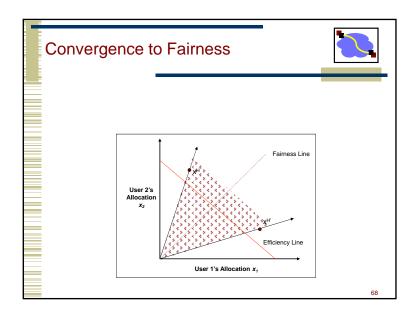
- Both X₁ and X₂ increase/decrease by the same amount over time
 - Additive increase improves fairness and additive decrease reduces fairness

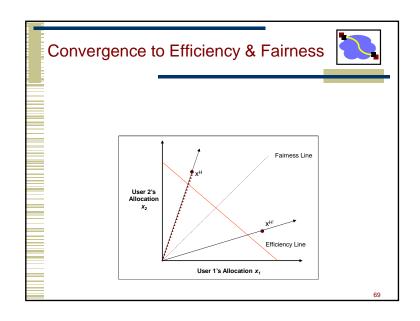


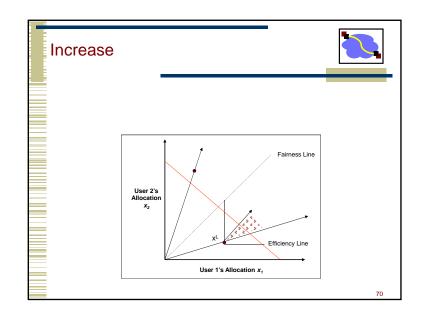


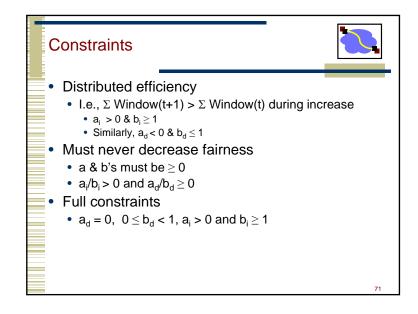


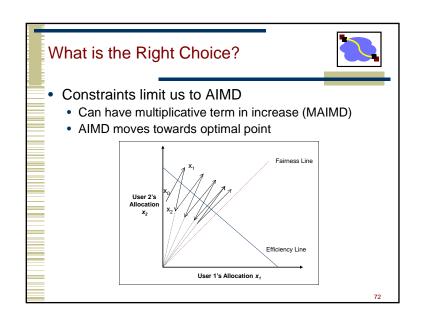












Questions



- Fairness why not support skew → AIMD/GAIMD analysis
- More bits of feedback → DECbit, XCP, Vegas
- Guess # of users → hard in async system, look at loss rate?
- Stateless vs. stateful design
- Wired vs. wireless
- Non-linear controls → Bionomial

73

TCP Congestion Control



- Congestion Control
- RED
- Assigned Reading
 - [FJ93] Random Early Detection Gateways for Congestion Avoidance
 - [TFRC] Equation-Based Congestion Control for Unicast Applications (two sections)