

DigiComms I

- Explored the fundamentals
 - What to ensure receiver can take in information at the rate the transmitter is sending it

2

- Circuit switching
 - Easy agree channel rate
- Packet switching
 - Harder
 - Interactions in intermediate switches
- Sliding windows for flow control
 - Receiver opens and closes the window
 - Still a delay in taking effect

Flow control problem

- Consider file transfer
- Sender sends a stream of packets representing fragments of a file
- Sender should try to match rate at which receiver and network can process data
- Can't send too slow or too fast
- Too slow
 - Wastes time
- Too fast
 - Can lead to buffer overflow
- How to find the correct rate?
 - Particularly given delays on control information

Other considerations

- Simplicity
- Overhead
- Scaling
- Fairness
- Stability
- Many interesting tradeoffs
 - Overhead for stability
 - Simplicity for unfairness

Where?

- Usually at transport layer
- Also, in some cases, in datalink layer

Model

Source, sink, server, service rate, bottleneck, round trip time







Hard problems

- Choosing a descriptor at a source
- Choosing a scheduling discipline at intermediate network elements
- Admitting calls so that their performance objectives are met (call admission control).

Descriptor requirements

- Representativity
 - Adequately describes flow, so that network does not reserve too little or too much resource

11

- Verifiability
 - Verify that descriptor holds
- Preservability
 - Doesn't change inside the network
- Usability
 - Easy to describe and use for admission control

Traffic descriptors

- Usually an envelope
 - Constrains worst case behaviour
- Three uses
 - Basis for traffic contract
 - Input to regulator
 - Input to policer

Examples

- Representative, verifiable, but not useable
 - Time series of inter-arrival times
- Verifiable, preservable, and useable, but not representative
 Peak rate





- Two ways to compute it
- For networks with fixed-size packets
 - Minimum inter-packet spacing
- For networks with variable-size packets • Highest rate over all intervals of a particular duration
- Regulator for fixed-size packets
 - Timer set on packet transmission
 - If timer expires, send packet, if any
- Problem
 - Sensitive to extremes

Average rate

- Rate over some time period (window)
- Less susceptible to outliers
- Parameters: *t* and *a*
- Two types: jumping window and moving window
- Jumping window
 - Over consecutive intervals of length *t*, only *a* bits sent
 - Regulator reinitializes every interval
- Moving window
 - Over all intervals of length t, only a bits sent
 - Regulator forgets packet sent more than t seconds ago

15

Linear Bounded Arrival Process

- Source bounds # bits sent in any time interval by a linear function of time
- The number of bits transmitted in any active interval of length t is less than rt + s
- *r* is the long term rate
- s is the burst limit
- Insensitive to outliers

14











Open loop vs. closed loop

- Open loop
 - Describe traffic
 - Network admits/reserves resources
 - Regulation/policing

Closed loop

- · Can't describe traffic or network doesn't support reservation
- Monitor available bandwidth
 - + Perhaps allocated using GPS-emulation
- Adapt to it
- If not done properly either
 - Too much loss
- Unnecessary delay



Explicit vs. Implicit

- Network tells source its current rate
- Better control
- More overhead
- Implicit

23

- Endpoint figures out rate by looking at network
- Less overhead
- Ideally, want overhead of implicit with effectiveness of explicit

24



- Recall error control window
- Largest number of packet outstanding (sent but not ACKed)
- If endpoint has sent all packets in window, it must wait => slows down its rate

25

27

- Thus, window provides *both* error control and flow control
- This is called *transmission* window
- Coupling can be a problem
 - Few buffers are receiver → slow rate!

Hop-by-hop vs. end-to-end

Hop-by-hop

- First generation flow control at each link
 - Next server = sink
- Easy to implement
- End-to-end
 - Sender matches all the servers on its path
- Plusses for hop-by-hop
 - Simpler
 - Distributes overflow
 - Better control
- Plusses for end-to-end
 - Cheaper

Window vs. rate

- In adaptive rate, we directly control rate
- Needs a timer per connection
- Plusses for window
 - No need for fine-grained timer
 - Self-limiting
- Plusses for rate
 - Better control (finer grain)
 - No coupling of flow control and error control
- Rate control must be careful to avoid overhead and sending too much

26

28

On-off

- Receiver gives ON and OFF signals
- If ON, send at full speed
- If OFF, stop
- OK when RTT is small
- What if OFF is lost?
- Bursty
- Used in serial lines or LANs



What should window size be?

- Let bottleneck service rate along path = b packets/sec
- Let round trip time = R sec
- Let flow control window = w packets
- Sending rate is w packets in R seconds = w/R
- To use bottleneck $w/R > b \rightarrow w > bR$
- This is the bandwidth delay product or optimal window size



















- Works with FIFO
 - but requires per-connection state (demand)
- Simple implementation
- But:
 - Assumes cooperation!
 - Has conservative window increase policy

TCP details

- Window starts at 1
- Increases exponentially for a while, then linearly
- Exponentially → doubles every RTT
- Linearly → increases by 1 every RTT
- During exponential phase, every ACK results in window increase by 1
- During linear phase, window increases by 1 when # ACKs = window size

41

43

- Exponential phase is called *slow start*
- Linear phase is called *congestion avoidance*

More TCP details

- On a loss, current window size is stored in a variable called slow start threshold or ssthresh
- Switch from exponential to linear (slow start to congestion avoidance) when window size reaches threshold
- Loss detected either with timeout or fast retransmit (duplicate cumulative ACKs)
- Many versions of TCP
 - Tahoe: in both cases, drop window to 1
 - Reno: on timeout, drop window to 1, and on fast retransmit drop window to half previous size (also, increase window on subsequent ACKs)

42

44

- Vegas, New Reno
- (CU)BIC: Used in Linux 2.6.8 (2.6.19)

TCP vs. DECbit

- Both use dynamic window flow control and additive-increase multiplicative decrease
- TCP uses implicit measurement of congestion
 Probe a black box
- Operates at the *cliff*
- Source does not filter information

Evaluation

- Effective over a wide range of bandwidths
- A lot of operational experience
- Weaknesses
 - Loss → overload? (not always: e.g. wireless)
 - Overload \rightarrow self-blame, problem with FCFS
 - Overload detected only on a loss
 - In steady state, source induces loss
 - Needs at least bR/3 buffers per connection



NETwork Block Transfer (NETBLT)

- First rate-based flow control scheme
- Separates error control (window) and flow control (no *coupling*)
- So, losses and retransmissions do not affect the flow rate
- Application data sent as a series of buffers, each at a given rate
- Rate = (burst size at burst rate) so granularity of control = burst
- Initially, no adjustment of rates
- Later, if received rate < sending rate, multiplicatively decrease rate

47

• Change rate only once per buffer \rightarrow slow













ATM Forum End-to-End Rate-based Flow Control (EERC)

- Similar to DECbit, but send a whole cell's worth of info instead of one bit
- Sources periodically send a Resource Management (RM) cell with a rate request

53

55

- Typically once every 32 cells
- Each server fills in RM cell with current share, if less
- Source sends at this rate

Comparison with DECbit

- Sources know exact rate
- Non-zero initial cell-rate → conservative increase can be avoided
- Interoperation between ER/CI switches

ATM Forum EERC details

- Source sends Explicit Rate (ER) in RM cell
- Switches compute source share in an unspecified manner (allows competition)
- Current rate = Allowed Cell Rate = ACR
- If ER > ACR then ACR = ACR + RIF × PCR else ACR = ER
 (PCR is Peak Cell Rate, RIF is Rate Increase Factor)
- If switch does not change ER, then use DECbit idea
 - If CI bit set, ACR = ACR (1 RDF)
- If ER < AR, AR = ER
- Allows interoperability of a sort
- If idle 500 ms, reset rate to Initial cell rate
- If no RM cells return for a while, ACR = ACR × (1-RDF)

Problems

- RM cells sitting in the data path is a mess
- Updating sending rate based on RM cell can be hard
- Interoperability comes at the cost of reduced efficiency (as bad as DECbit)
- Computing ER is difficult

54



Hybrid flow control

- Source gets a minimum rate, but can use more
- All problems of both open loop and closed loop flow control
- Resource partitioning problem
 - What fraction can be reserved?
 - How?