

# **A Survey of Results on TCP/IP over ATM**

**Raj Jain**

**Raj Jain is now at  
Washington University in Saint Louis  
Jain@cse.wustl.edu**

**<http://www.cse.wustl.edu/~jain/>**

# Our Team

## q **Current:**

- q Shivkumar Kalyanaraman
- q Rohit Goyal
- q Sonia Fahmy
- q Jianping Jian

## q **Past:**

- q Fang Lu
- q Ram Viswanathan
- q Manu Vasandani
- q Arun Krishnamoorthy



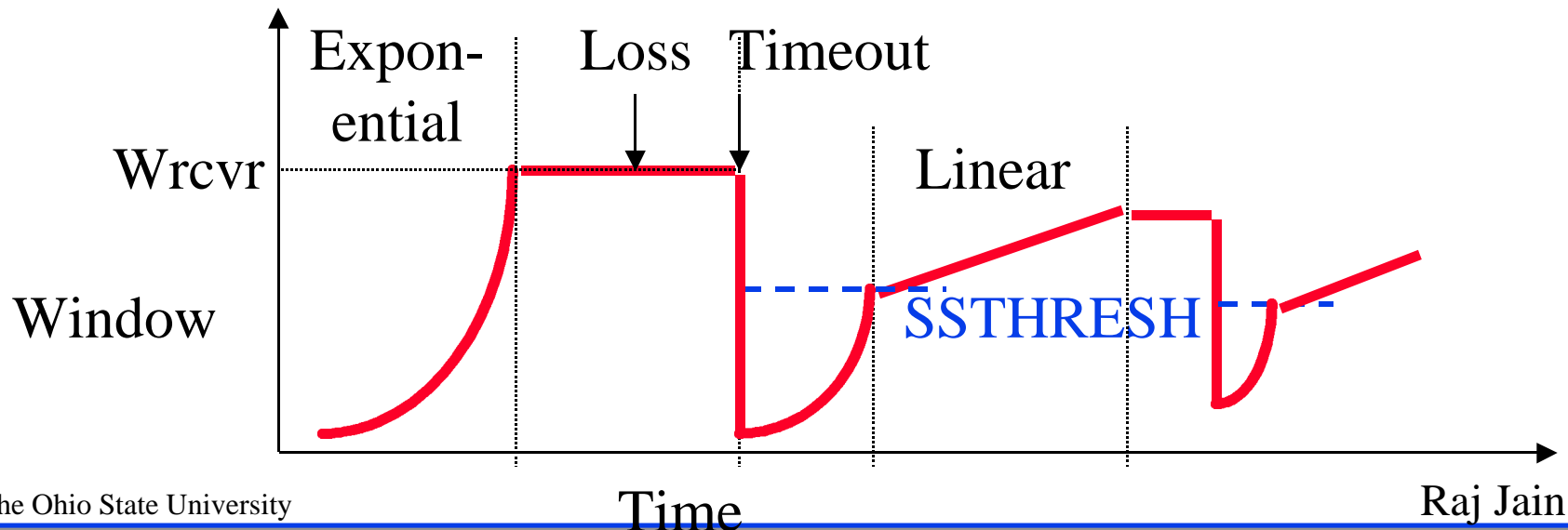
- q TCP Congestion mechanisms  
Slow start, Fast Retransmit/recovery, RED, ECN
- q TCP over ABR results
- q Tough TCP Tests: Further modifications of ERICA
- q TCP over UBR results
- q UBR+

# TCP Congestion Mechanisms

- q End-System Based:
  - q Silly Window Syndrome Avoidance
  - q Delayed Ack
  - q Slow Start Congestion Avoidance
  - q Fast Retransmit and Recovery
  - q *Selective Acknowledgment (SACK)*
- q Router Based:
  - q *Random Early Discard (RED)*
  - q *Explicit Notification (Future)*
- q Difference between routers and switches is decreasing  
It is important to understand TCP/IP mechanisms

# TCP/IP Slow Start

- q Maximum Segment Size (MSS) = 512 bytes
- q Congestion Window (CWND)
- q Window  $W = \text{Min}\{W_{rcvr}, \text{CWND}\}$
- q Slow-Start Threshold =  $\max\{2, \min\{\text{CWND}/2, W_{rcvr}\}\}$
- q Exponential until SSTHRESH:  $W = W + 1$  for every ack
- q Linear afterwards:  $W = W + 1/W$  for every ack until  $W_{rcvr}$



# ABR with Small Buffers

# srcs	TBE	Buffer Size	T1	T2	T3	T4	T5	Throughput	% of Max	CLR.
2	128	256	3.1	3.1				6.2	10.6	1.2
2	128	1024	10.5	4.1				14.6	24.9	2.0
2	512	1024	5.7	5.9				11.6	19.8	2.7
2	512	2048	8.0	8.0				16.0	27.4	1.0
5	128	640	1.5	1.4	3.0	1.6	1.6	9.1	15.6	4.8
5	128	1280	2.7	2.4	2.6	2.5	2.6	12.8	21.8	1.0
5	512	2560	4.0	4.0	4.0	3.9	4.1	19.9	34.1	0.3
5	512	5720	11.7	11.8	11.6	11.8	11.6	58.4	100.0	0.0

- q CLR has high variance
- q CLR does not reflect performance. Higher CLR does not necessarily mean lower throughput
- q CLR and throughput are one order of magnitude apart

# TCP over ABR: Observations

- q CLR in the switch is low. But, throughput is also low
- q The buffers can not be allocated based on TBE
- q Maximum queue length and TBE have little/no relationship

# TCP: Observations

- q With enough buffers in the network, TCP can automatically fill any available capacity.
- q TCP performs best when there is NO packet loss. Even a single packet loss can reduce throughput considerably.
- q Slow start limits the packet loss but loses considerable time. With TCP, you may not lose too many packets but you lose time.
- q Bursty losses cause more throughput degradation than isolated losses.
- q With low buffers, TCP does not use all the available bandwidth
- q Many duplicate packets are dropped at the destination



- q For each packet loss, much time is lost due to timer granularity  
Timer granularity is the key parameter in determining time lost

# Fast Retransmit and Recovery

- q Fast Retransmit: Three consecutive acks for the same segment  $\Rightarrow$  Loss  $\Rightarrow$  Retransmit before timeout
- q Fast Recovery: Reduce congestion window to half (instead of 1)  
 $\Rightarrow$  No new transmissions until duplicate acks arrive for the remaining half
- q Single packet loss  $\Rightarrow$  One RTT wasted  $\Rightarrow$  Not bad
- q Multiple packet loss  $\Rightarrow$  Timeout  
Timeout = Mean + 4 Stdv. of RTT or one tick  
 $\Rightarrow$  100 ms wasted  $\Rightarrow$  Really bad

# Effect of Fast Retransmit

- q Fast retransmit helps only if occasional losses  
Mild congestion or errors
- q With n packet loss, Ssthresh is reduced to half after each retransmission. Window enters the linear-increase zone even when the window is small  $\Rightarrow$  Low throughput.
- q Even with fast retransmits, there are time-outs when the losses are bursty. These time-outs are more damaging than if there is no fast retransmit since Ssthresh is low.

	Bursty Loss	Scattered Loss
With Fast-Retransmit Fast-Recovery	×	√
Without Fast-Retransmit Fast-Recovery	√	×

# Buffer Requirements for ABR: Key Factors

- q Switch Algorithm: Transient Response (settling) time
- q Round Trip Time (RTT)
- q Feedback Delay (bottleneck to source)
- q Switch Algorithm *Parameters*:
  - q Averaging Interval
  - q Target Utilization
  - q ERICA+ queue control
- q Presence and characteristics of background VBR
- q Number of VCs
- q TCP Receiver window size

# ABR Switch Buffer Requirements

- q ABR performance depends heavily upon the switch algorithm.

Following statements are based on our *modified ERICA* switch algorithm.

No cell loss for *TCP* if switch has Buffers =  $4 \times \text{RTT}$ .

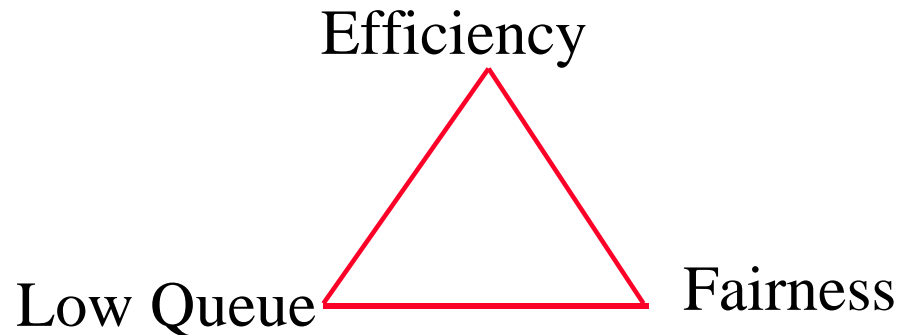
- q No loss for **any** number of TCP sources w  $4 \times \text{RTT}$  buffers.
- q No loss even with **VBR**. W/o VBR,  $3 \times \text{RTT}$  buffers will do.
- q Under many circumstances,  $1 \times \text{RTT}$  buffers may do.
- q With ABR most of the queues are at the source.  
Not much queue in the switch
- q In general:  
$$Q_{\max} = a \times \text{RTT} + b \times \text{Averaging Interval} + c \times \text{Feedback delay} + d \times \text{fn(VBR)}$$

# High Frequency VBR: Problem

- q Limit of  $1 \times \text{RTT}$  due to VBR is good for large VBR cycle times.  
TCP and ABR get enough time to adjust.
- q Faster VBR causes faster variations in available capacity.  
Neither TCP nor Switch algorithm may have time to adjust  
 $\Rightarrow$  Can lead to instability at high utilization levels.

VBR	F/b	Maximum	Total	Effici-	Fair-	
On/Off	RTT	Delay	Queue	Throughput	ency	Ness
30 ms	30	10	12359=1.12*RTT	69.60	92.65	0.9967
100 ms	30	10	13073=1.18*RTT	63.85	85.00	0.9987
10 ms	30	10	diverges			
1 ms	30	10	diverges			

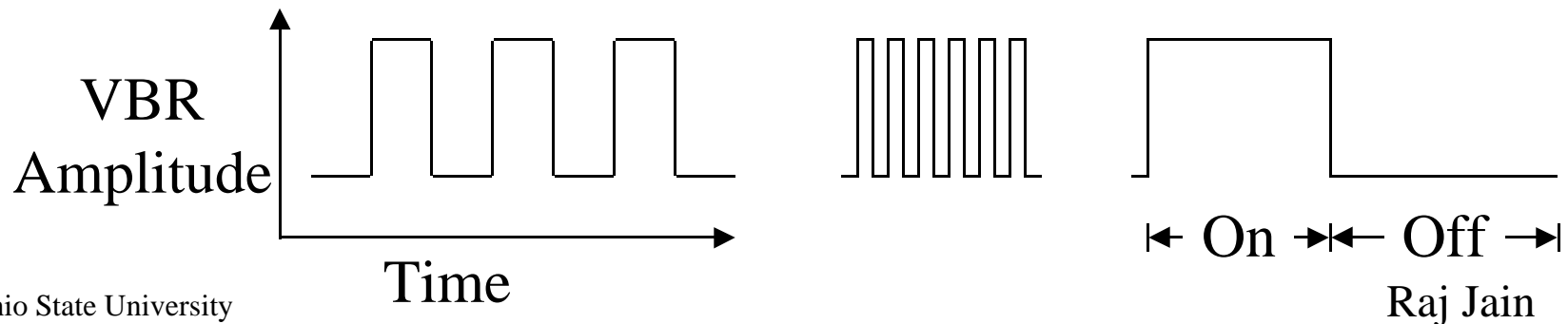
# Three Way Tradeoff



- q Buffers vs Efficiency (Utilization) vs Fairness
- q It is possible to have lower queues (lower buffer required) if the target utilization is kept low.

# ABR Test Cases

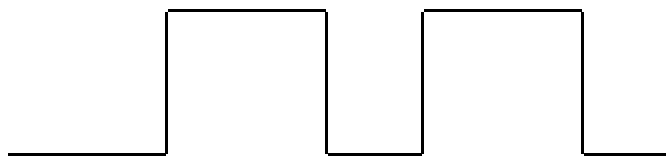
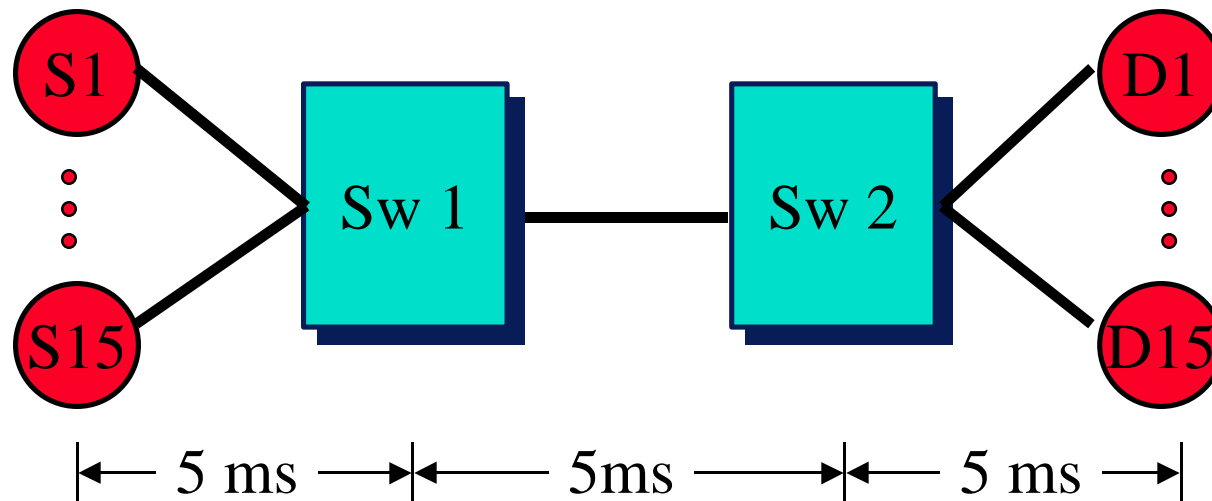
- q Configurations: n sources, parking lot, ...
- q ABR Traffic: Infinite, **bursty**
- q Background traffic: without and **with VBR**
- q High Layer: non-TCP, **TCP**
- q RTT mix: Similar RTTs, **varying RTTs**
- q VBR Period: Large, **medium**, small  
(compared to feedback delay)
- q VBR Duty cycle: 0.9, **0.8**, **0.7**, ...





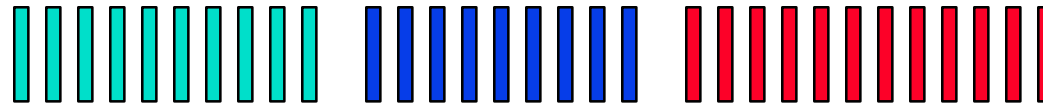
# A Tough Test Case

- q 15 TCP sources, with 10 ms VBR period with a duty cycle of 0.7 (7 ms on, 3 ms off),  
10 ms feedback delay,  
30 ms RTT



# Flocking Effect

- q All cells of a VC are often seen together.
- q There is clustering of sources.
- q Not all sources are seen all the time.



# ERICA Modifications

- q Boundary Cases:
  - q No ABR cells received
    - $\Rightarrow$  No active sources ( $N=0$ )
    - Fairshare =  $\infty$ ?

# ERICA Modifications (Cont)

- q Average number of sources
- q Average load factor = ABR Input rate/ABR capacity
  - q Average ABR Input Rate
    - = Number of cells/averaging interval
    - : Average Number of ABR cells
    - : Average Averaging interval
  - q Average VBR usage
    - : Average Number of VBR cells
    - : Average Averaging interval
- q Averaging  $\Rightarrow$  Decisions based on longer timer
  - $\Rightarrow$  Slower response
  - $\Rightarrow$  Buffer requirements are over  $4 \times \text{RTT}$

# TCP Over Plain UBR

- q Low throughput
- q Unfair
- q Anomalies: More receiver buffer  $\Rightarrow$  Lower throughput  
Due to Silly window avoidance + Delayed Ack
- q Solution: Min sender buffer size should be  $3 \times \text{MSS}$

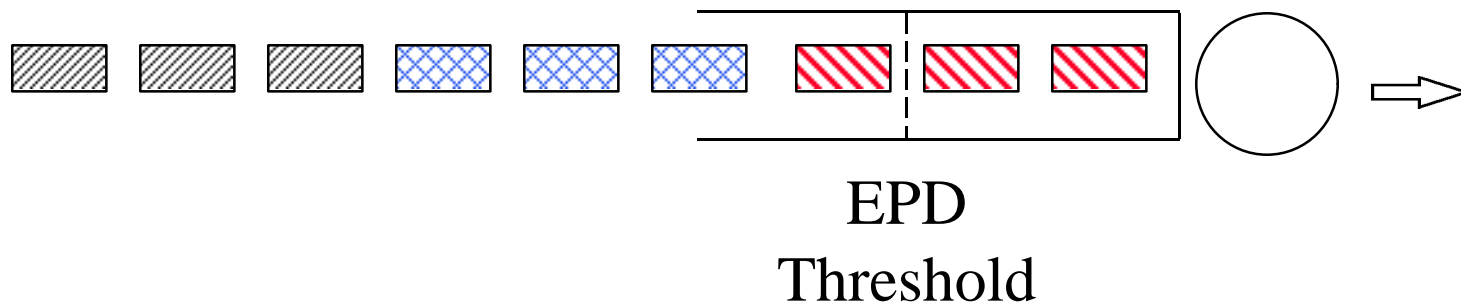
Ref: Comer

# TCP/IP over UBR: Improvements

- q Switch Based Mechanisms:
  - q PPD
  - q EPD
  - q EPD + per-VC Accounting
  - q EPD + per-VC queueing
- q Source Based Mechanisms:
  - q CLP Probe
  - q Cell Pacing
  - q Smaller Segments

# PPD and EPD

- q Plain ATM: Discard all cells if  $Q > \text{threshold}$
- q Partial Packet Discard:  
Discard all cells of a packet if one cell dropped  
 $Q > \text{threshold}$
- q Early Packet Discard:  
Discard all cells of new packets if  $Q > \text{threshold}$



# PPD vs EPD

- q Plain ATM  $\Rightarrow$  Many packets dropped
- q Dropping all cells of a packet is better than dropping randomly  
 $\Rightarrow$  PPD is better than plain UBR
- q Never drop the EOM cell of a packet.  
It results in two packet losses.
- q EPD  $\Rightarrow$  Even fewer packets dropped  $\Rightarrow$  better throughput
- q Plain ATM  $\ll$  PPD  $\ll$  EPD
- q EPD improves efficiency but not fairness

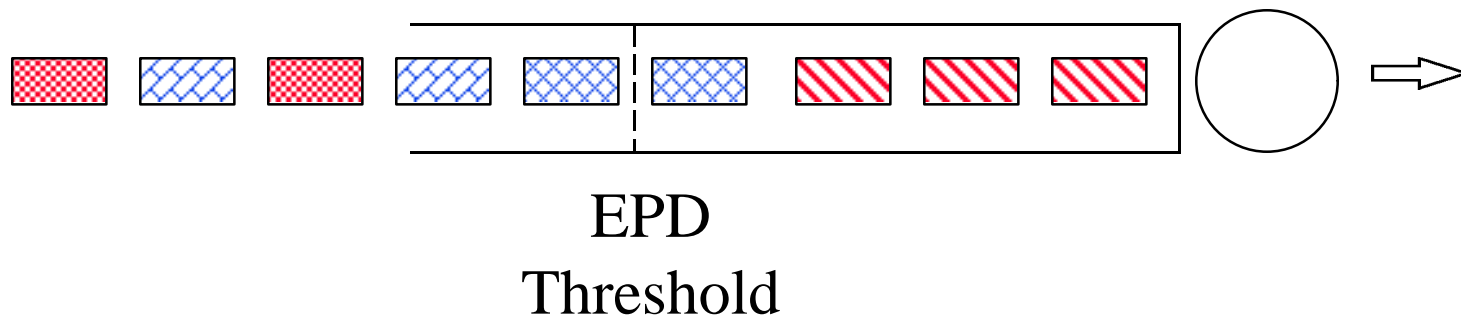


# UBR Switch Buffer Requirements

- q Switch queues may be as high as the sum of TCP windows  
No cell loss for TCP if Buffers =  $\Sigma$  TCP receiver window
- q Required buffering depends upon the number of sources.
- q TCP receiver window  $\geq$  RTT for full throughput with 1 source.
- q Unfairness in many cases.
- q Fairness can be improved by proper buffer allocation, drop policies, and scheduling.
- q Drop policies are more critical (than ABR) for good throughput
- q No starvation  $\Rightarrow$  Lower throughput shows up as increased file transfer times = Lower capacity

# EPD + Per-VC Accounting

- q Selective EPD: Select only high rate VCs  
Fast Buffer Allocation Scheme
- q EPD: Drop all packets if queue  $X > \text{threshold } R \Rightarrow \text{Unfair}$
- q No per-VC queueing  $\Rightarrow$  All VCs share a single FIFO queue
- q per-VC accounting  $\Rightarrow$  track  $X_i$  and  $N$
- q  $N = \#$  of non-zero  $X_i$ 's

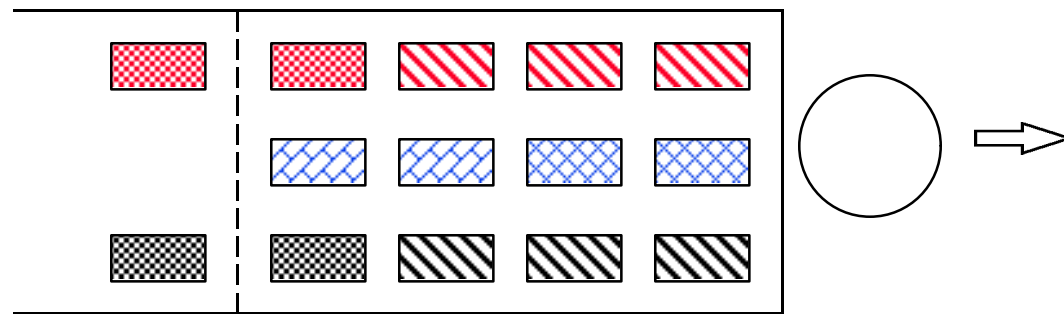


- q If  $X > \text{threshold}$ , drop next arriving packet if  $X_i \geq \text{fn}(X, N, K, R)$
- q Where  $K = \text{Total number of buffers}$
- q Drop if  $X_i/(X/N) \geq Z(1 + (K-X)/(X-R))$
- q Here  $Z = \text{parameter between } 0.5 \text{ and } 1$
- q Note that packets from more and more flows are dropped as queue  $X$  increases
- q Siu has analyzed a EPD + Simpler per-VC accounting
- q If  $X > \text{threshold}$ , drop next arriving packet if  $X_i/(X/N) \geq Z$
- q **Conclusion:** Per-VC accounting improves fairness
- q Other Ideas:
  - q Do not drop successive packets
  - q Drop from queues not tails  $\Rightarrow$  earlier effect

# EPD + Per-VC Queueing

- q Accept the next packet if  $X_i/(X/N) < Z$
- q Round-robin scheduling  $\Rightarrow$  Fairness further improved
- q However, more VC's have packets dropped  $\Rightarrow$  Lower total throughput

Ref: Siu



EPD  
Threshold

# CLP Probe

- q Idea:
  - q Use probe packets with CLP bit set to sense network congestion
  - q If probe makes it then increase window.
  - q otherwise slow-start
- q Whenever window is increased, the next packet is sent with CLP set
- q Throughput improved from 53% to 85%

Ref: Perloff and Reiss

# Cell Pacing

- q Use lower than link rate
- q Using the right rate changed throughput from 0.9% to 68%
- q Even with multiple sources throughput changed from 1.6% to 52%
- q How to select the right rate?

Ref: Ewy, et al

# Effect of Segment Size

- q Large segments  $\Rightarrow$  Large retransmissions
- q If buffering is the bottleneck, smaller segments  $\Rightarrow$  better throughput
- q If processing is the bottleneck, smaller segments  $\Rightarrow$  More overhead  $\Rightarrow$  Less throughput
- q Buffering was small in initial switches
- q Need buffering equal to several round-trips

Ref: Ewy, et al

# Selective Acknowledgment

- q Allows receiver to indicate multiple blocks of received segments.
- q Receiver indicates lower-edge and upper-edge of all received segments
- q Senders can retransmit only the missing segments.
- q There is no need for multiple timeouts or duplicates





# Random Early Discard

- q Exponential averaging
  - q Bursty traffic
    - ⇒ Instantaneous queue can be high or low
    - ⇒ time averaging
  - q  $Q_{avg} = (1-\alpha)Q_{avg} + \alpha q$  if  $q > 0$
  - q  $Q_{avg} = (1-\alpha)^\beta Q_{avg}$  otherwise,  $\beta = f(\text{idle time})$
- q Two thresholds:
  - q  $\text{Min} < Q_{avg} < \text{Max} \Rightarrow$  mark (drop) arriving packet with probability  $p$
  - q  $pb = pb_{max} \times (Q_{avg} - \text{Min}) / (\text{Max} - \text{min})$
  - q  $p = pb / (1 - \text{count} \times pb)$

- q Count = Number of packets since the last mark  
⇒ Marking probability increases as more packets arrive
- q Count reset after marking a packet
- q  $Q_{avg} \geq \text{Max\_threshold} \Rightarrow \text{Drop all}$

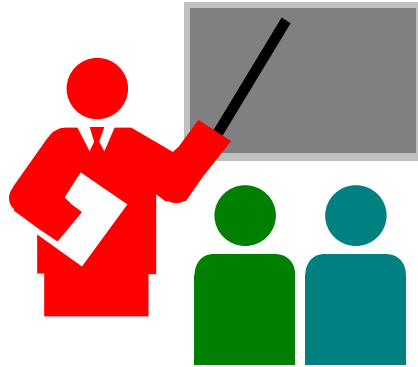
Ref: Floyd and Jacobson

# Explicit Notification

- q Routers send ICMP messages (or set bit) when  $Q_{avg} > \text{threshold}$
- q Sources respond to ECN once per round trip  $\Rightarrow$  Ignore others
- q Halve  $cwnd$  and  $ssthresh$  on first ECN
- q Do not respond to succeeding ECNs until all outstanding packets have been acked (i.e., one round trip)
- q Do not reduce  $cwnd$  or  $ssthresh$  after timeout or 3 duplicate acks, if ECN action taken in this round trip.

Ref: S. Floyd, "TCP and Explicit Congestion Notification," LBL Tech Report.

# Summary



- q Packet loss results in a significant degradation in TCP throughput. For best throughput, TCP needs no loss.
- q With enough buffers, ABR may guarantee zero loss for any number of TCP sources.
- q Performance of ABR depends on the switch algorithm
- q For zero loss, UBR need buffers =  $\Sigma$  receiver windows
- q PPD  $\ll$  EPD  $\ll$  Selective EPD
- q ABR vs UBR issue is that of ATM end-to-end vs backbone

# TCP Congestion: References

- q W. Richard Stevens, "TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms," Internet Draft draft-stevens-tcpca-spec-00.txt, February 1996
- q V. Jacobson and M. Karels, "Congestion Avoidance and Control," Proc. SIGCOMM'88, August 1988, pp. 314-329.
- q S. Floyd and V. Jacobson, "Random Early Detection Gateways for Congestion Avoidance," IEEE/ACM Transactions on Networking, Vol. 1, August 1993, pp. 397-413.
- q J. Nagle, "Congestion Control in IP/TCP Internetworks," ACM Computer Communications Review, Vol. 14, No. 4, October 1984, pp. 11-17.

- q V. Jacobson, R. Braden, and D. Borman, "TCP Extensions for High Performance," Internet RFC 1323, May 1992.
- q S. Floyd, "TCP and Explicit Congestion Notification,"
- q S. Floyd, "TCP and Successive Fast Retransmits," LBL Technical Report, May 1995, available <ftp://ftp.ee.lbl.gov/papers/fastretrans.ps>
- q C. Villamizar and C. Song, "High Performance TCP in ANSNET," Computer Communications Review, October 1994, pp. 45-60.
- q J. Hoe, "Improving the start-up behavior of a congestion control scheme for TCP," ACM SIGCOMM 96, to appear.
- q L. Zhang, S. Shenker, and D.D. Clark, "Observations on the dynamics of a congestion control algorithm: the effects of two-way traffic," SIGCOMM'91, September 1991, pp. 133-147.

# TCP over ATM: Our Papers/Contributions

All our past ATM forum contributions, papers and presentations can be obtained on-line at <http://www.cis.ohio-state.edu/~jain/>

- q S. Kalyanaraman, R. Jain, S. Fahmy, R. Goyal, F. Lu and S. Srinidhi, ``Performance of TCP/IP over ABR," To appear Globecom'96.
- q R. Jain, et al, "Buffer Requirements for TCP over ABR," ATM Forum/96-0517, April 1996.
- q R. Jain, et al, "Performance of TCP over UBR and buffer requirements," ATM Forum/96-0518, April 1996.
- q R. Jain, et al, "TBE and TCP/IP traffic," ATM Forum/96-0177, February 1996.

# TCP over ATM: References

- q A. Romanow and S. Floyd, "Dynamics of TCP Traffic over ATM Networks," IEEE Journal on Selected Areas in Communications, Vol. 13, No. 4, May 1995, pp. 633-641, [ftp://ftp.ee.lbl.gov/papers/tcp\\_atm.ps.Z](ftp://ftp.ee.lbl.gov/papers/tcp_atm.ps.Z)
- q J. Heinanen and K. Kilkki, "A Fair Buffer Allocation Scheme," Telecom Finland Draft 17 March 1995.
- q H. Li, K-Y Siu, and H-Y Tzeng, "TCP Performance over ABR and UBR Services in ATM," Proc. IPCCC'96, March 1996.
- q D. E. Comer and J. C. Lin, "TCP Buffering and Performance over an ATM Network," Internetworking: Research and Experience, Vol. 6, 1995, pp. 1-13.



- q M. Perloff and K. Reiss, "Improvements to TCP Performance in High-Speed ATM Networks," Communications of ACM, February 1995, pp. 90-100.
- q B.j. Ewy, et al, "TCP/ATM Experiences in the MAGIC Testbed,"
- q L. Kalampoukas and A. Varma, "Performance of TCP over Multi-Hop ATM Networks: A Comparative Study of ATM-Layer Congestion Control Schemes," Technical Report, UCSC-CRL-95-13, in <ftp://ftp.cse.ucsc.edu/>