

# On the Effectiveness of Delay-based Congestion Avoidance

Ravi Prasad,  
Manish Jain,  
Constantinos Dovrolis

ravi, jain, dovrolis@cc.gatech.edu

College of Computing  
Georgia Tech

## Outline

- o Loss-based Congestion Avoidance (LCA)
  - o Inefficient in high BDP paths
- o Delay-based Congestion Avoidance (DCA)
  - o Includes TCP Vegas and FAST
- o Controllability and observability
  - o Does DCA meet these properties?
- o Four failure scenarios for DCA
  - o Small RTT variations
  - o RTT undersampling
  - o DCA in highly aggregated paths
  - o Random losses

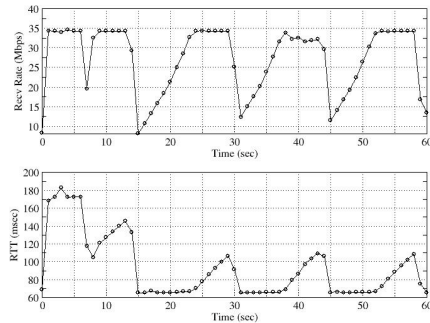
2/17/04

PfIdNet 2004

2

## Loss-based Congestion Avoidance (LCA)

- o TCP Reno follows the LCA model
- o LCA flow increases cwnd until it sees packet loss



### Consequences:

- Throughput less than available bandwidth
- Increased loss-rate
- Large delay variation
- Large buffering requirement in routers

2/17/04

PfldNet 2004

3

## TCP Reno can be inefficient in high Bandwidth-Delay Product (BDP) paths

- o Example: 1Gbps, 100msec, 1500B pkts
  - o Large window reduction upon loss
    - o Single loss  $\approx$  4000 pkt window reduction
  - o Large recovery time
    - o Recovery from 1 pkt loss  $\approx$  400 sec
  - o Need very small loss probability
    - o Loss rate must be less than  $2 \cdot 10^{-8}$
- o Side comment: previous example assumes unlimited socket buffers
- o See SOBAS for automatic socket buffer sizing to avoid losses

2/17/04

PfldNet 2004

4

## Delay-based Congestion Avoidance (DCA)

- o Early paper by R. Jain:  
ACM CCR - Oct 89
- o Follow-up protocols:
  - o CARD ('89)
  - o Tri-S ('91)
  - o DUAL ('92)
  - o TCP-Vegas ('94)
  - o TCP-BFA ('98)
  - o TCP-FAST ('03)
  - o SOBAS ('03)

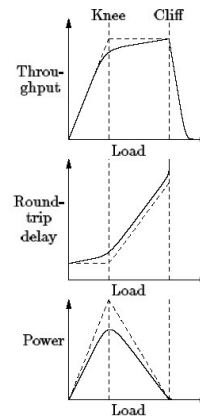


Figure 1: Network performance as a function of the load. Broken curves indicate performance with deterministic service and interarrival times.

2/17/04

PfldNet 2004

6

## TCP Vegas

- o Expected rate  $E = \text{Window} / T_{\min}$
- o Measure actual send rate  $R$
- o Adjust window based on difference  $E - R$
- o If  $(E - R)$  is larger than threshold, then:
  - o Flow has built up queue at bottleneck
  - o RTT  $T$  is larger than  $T_{\min}$
  - o Reduce window to avoid congestion

2/17/04

PfldNet 2004

7

## TCP FAST

- o Improved (stabilized) version of TCP Vegas
- o Reduce window when RTT increases
- o Window decrease factor depends on:
  - o Current window size
  - o RTT relative increase
- o As opposed to constant decrease factor of Vegas
- o Scales well in high BDP paths

2/17/04

PfldNet 2004

8

## First concerns about DCA schemes

- o Measurement studies showed little or no correlation between increased RTTs and network load
  - o Martin, et al., ToN 2003, "Delay-based congestion avoidance for TCP"
  - o Biaz and Vaidya, IMC 2003, "Is the Round-Trip Time Correlated with the Number of Packets in Flight?"
- o Previous studies showed that correlations of RTT and load are even weaker in high-bandwidth paths
  - o Interesting, because DCA is supposed to work better than LCA in high BDP paths!

2/17/04

PfldNet 2004

9

## DCA can fail when:

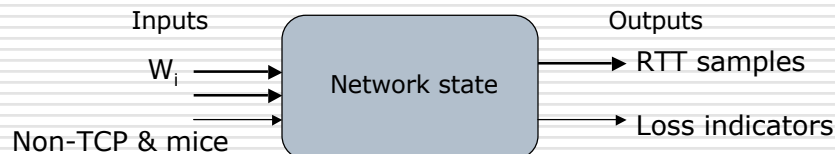
- o RTT variations cannot be reliably measured
  - o Cause: small network buffers or noisy RTTs
- o DCA flow undersamples RTT variations
  - o Cause: flow throughput is small relative to cross traffic throughput
- o DCA flow cannot affect RTT
  - o Cause: DCA flow competes with many other flows in highly aggregated path
- o DCA flow cannot avoid losses
  - o Cause: random losses or congestive losses due to cross traffic bursts

2/17/04

PfldNet 2004

10

## Another look at the problem..



- o Network state:
  - o Tight link queue size  $Q(t)$  and available bandwidth  $A(t)$
- o Input variables:
  - o Window  $W_i(t)$  for flow  $i$  (DCA or LCA)
  - o Instantaneous rate of non-TCP flows
- o Output variables:
  - o Sampled RTT  $T_i(t)$  and loss indicator  $L_i(t)$  for flow  $i$

2/17/04

PfldNet 2004

11

## Controllability & Observability

- o Major concepts in theory of control systems
- o Controllability:
  - o A state variable (or output) is controllable if it can be driven to a certain level by an input variable
- o Observability:
  - o A state variable is observable if it directly affects one or more outputs
- o A DCA flow should be able to:
  - o **Observe** the queue size (state variable) through sampled RTT signal (output)
  - o **Control** the queue size through send window (input)
- o Under which conditions will a DCA flow meet (or not meet) these properties?

2/17/04

PfIdNet 2004

12

## RTT Signal-to-Noise ratio

- o Minimum RTT:  $T_{\min}$
- o Maximum RTT:  $T_{\min} + B_{\dagger}/C_{\dagger}$ 
  - o  $B_{\dagger}$ : Tight link buffer
  - o  $C_{\dagger}$ : Tight link capacity
- o We need to consider RTT noise  $n(t)$ 
  - o Random queueing at:
    - o Non-tight links in forward path
    - o Reverse path
  - o End-host timestamping resolution & OS noise
- o Measured RTT:  $T(t) = T_{\min} + Q_{\dagger}(t)/C_{\dagger} + n(t)$ 
  - o What if noise > signal?

2/17/04

PfIdNet 2004

13

## RTT Signal-to-Noise ratio (cont')

- o Signal:  $Q_+(t)/C_+$
- o We have a problem if  $B_+/C_+ = O(\text{noise})$ 
  - o RTT variation will not be measured accurately
  - o DCA flow will not avoid congestive losses
  - o Queue size is **not observable**, because the output signal (RTT) is too weak
- o Lesson:
  - o DCA effectiveness depends on link buffer sizes
  - o Small buffers and RTT measurement noise can **break observability**

2/17/04

PfIdNet 2004

14

## RTT undersampling

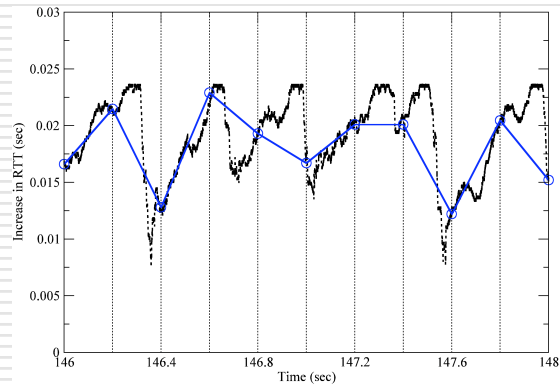
- o Suppose that at time  $t$ :
  - o Available bandwidth =  $A$
  - o New cross-traffic flow starts with rate  $R_x > A$
  - o Tight link buffer will fill up after  $B_+/(R_x - A)$
- o DCA flow with throughput  $R$  samples RTT at rate:
  - o  $L/R$ ,  $L$ : packet size
- o DCA undersamples RTT, and may not detect queue buildup, if
  - o  $L/R = O(B_+/(R_x - A))$
  - o Or,  $R \ll (R_x - A)$
  - o Queue size **not observable** because input is too **sparse**

2/17/04

PfIdNet 2004

15

## RTT undersampling (cont')



- o Lesson: DCA may not work for low-throughput flows (relative to the cross traffic flows)

2/17/04

PfIdNet 2004

16

## Highly aggregated traffic

- o Cross traffic rate:  $R_c = C_T - R$ 
  - o Aggregate of many flows
  - o Each flow is small share of aggregate
- o DCA flow rate:  $R \ll R_c$
- o Tight link queue size, and RTT variations, are not controlled by DCA flow
  - o RTT variations appear as erratic ("random")
  - o Queue size is not controllable because input signal due to any single flow is too weak

2/17/04

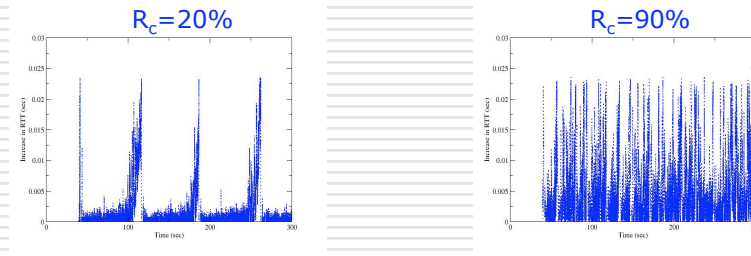
PfIdNet 2004

17



## Highly aggregated traffic (cont')

- o Cross traffic
  - o TCP flows of size 15-20pkts
  - o Constant arrival rate of flows



- o Lesson: DCA may not work in paths that carry many relatively small flows

2/17/04

PfIdNet 2004

18

## Random losses

- o A DCA flow can avoid self-induced congestive losses
  - o Losses caused and experienced by that flow
- o A DCA flow cannot avoid congestive losses caused by high-rate bursts from other flows
  - o Queue size controlled by other flows
  - o Losses appear as "random" to DCA flow
- o But, a DCA flow should still react to such losses to avoid congestion collapse
  - o How should DCA flows react to losses?

2/17/04

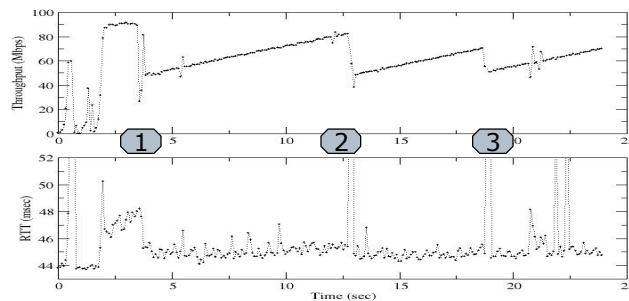
PfIdNet 2004

19

o They often do not play fair with TCP Reno in that aspect

## Random losses (cont')

- o Losses are not always preceded by increasing RTTs
- o RTTs may have increased, but not observed by DCA flow



- DCA flow could avoid losses at 1, but not at 2 or 3
- Lesson: DCA schemes should still expect and react to congestive losses

2/17/04

PfIdNet 2004

20

## Concluding remarks

- o DCA is probably more efficient than LCA (Reno) under certain conditions:
  - o A few high-throughput flows
  - o Well-buffered links
  - o No significant RTT noise
  - o No "random" losses
- o Most simulation studies meet previous conditions

o The real Internet does not...

2/17/04

PfIdNet 2004

21

## Concluding remarks (cont')

- o More research is needed to establish the robustness of DCA schemes
  - o Robustness vs efficiency trade-off?
- o Robustness in terms of:
  - o Different buffer sizes, link capacities, flow RTTs
  - o Heterogeneous traffic models and applications
- o Realistic noise sources in RTT measurements

2/17/04

PfdNet 2004

22