



Delay-based AIMD congestion Control

D. Leith¹, R. Shorten¹, G. McCullagh¹, J. Heffner², L. Dunn³, F. Baker³

¹ *Hamilton Institute, NUI Maynooth, IE*

² *Pittsburgh Supercomputing Center, US*

³ *Cisco Systems, US*



Overview

1. Environments which challenge TCP
2. Delay-based congestion control
3. Delay-based AIMD algorithm
4. Experiment Results
5. Conclusions





Environments which challenge TCP



Environments which challenge TCP

- High speed networks — large BDP
- Satellite links — very long delays
- Wireless links — non-congestive losses
- aDSL links — large latencies

Large BDPs, large queues, large transmission times, random losses.



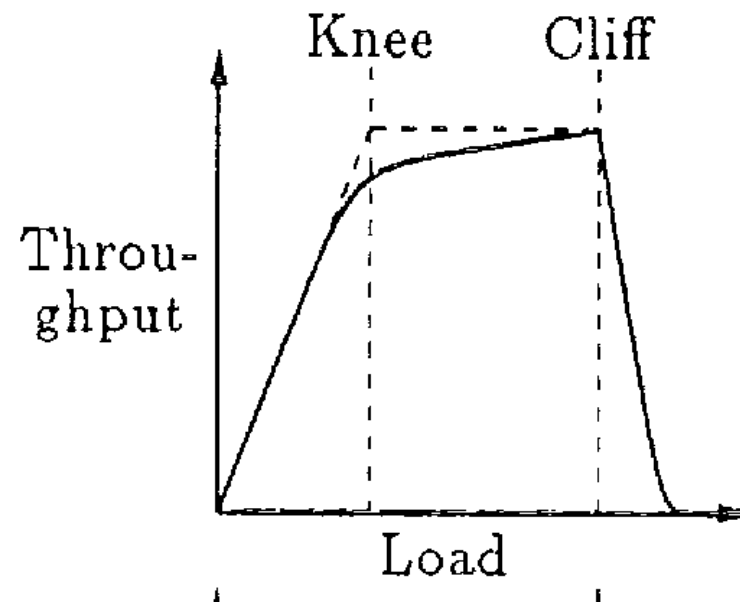
Proposed Solution

Combine ideas from high-speed and delay-based protocols

- Fast recovery to improve throughput on large BDP links
- Allow new flows to gain their share quickly, even with long transmission times
- Maintain low queueing delay



Operation “At the knee of the curve”



Jain’s “Sweet spot” around which delay is low but throughput is high





Delay-based Congestion Control



Delay-based Congestion Control

- Suggested by Brakmo et al in 1994
- Proposed the “Vegas” protocol.
- Others include FAST, Compound-TCP.
- Vegas is one of the more widely explored delay-based cc algorithms



Delay-based Congestion Control

Possible benefits:

- No congestive losses
- Doesn't fill queues, lower delays
- Lower cost per congestion event

Problems:

- Coexistence with loss-based systems
- Difficulty in accurately measuring delay, sampling rates
- Limited correlation between delay and congestion
- Delay scales with number of flows





Brief Analysis of Vegas



Hamilton Institute



Vegas

The expected and observed throughputs are:

$$r_{exp} = \frac{w}{T_{min}} \quad (1)$$

$$r_{obs} = \frac{w}{T_{min} + \tau} \quad (2)$$

$$\epsilon = r_{exp} - r_{obs} \quad (3)$$

$$= \frac{w}{T_{min}} - \frac{w}{T_{min} + \tau} \quad (4)$$



Vegas

Once per round trip-time cong. window, w is adjusted:

$$w \leftarrow \begin{cases} w + 1 & \epsilon < \alpha \\ w & \epsilon \in [\alpha, \beta] \\ w - 1 & \epsilon > \beta \end{cases} \quad (5)$$

where α and β are design parameters.



Vegas Latency Scaling

So, at equilibrium Vegas maintains the congestion window such that

$$\alpha \leq \epsilon \leq \beta \quad (6)$$

$$\alpha \leq \frac{w}{T_{min}} - \frac{w}{T_{min} + \tau} \leq \beta \quad (7)$$

$$\alpha \leq \frac{w\tau}{T_{min}(T_{min} + \tau)} \leq \beta \quad (8)$$

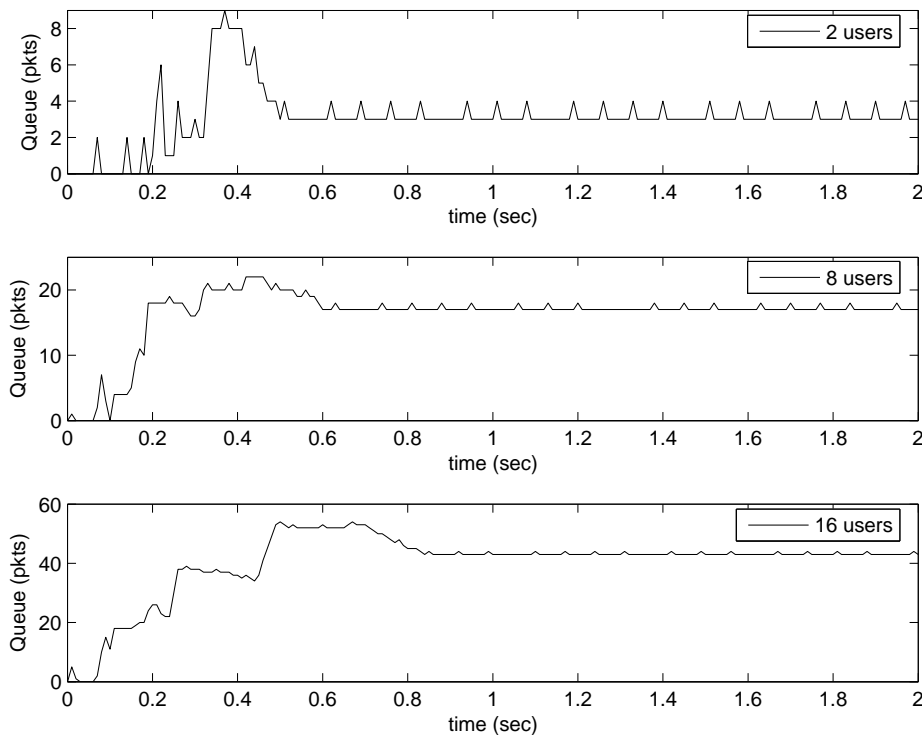
$$\alpha \leq \frac{\tau}{T_{min}} r_{obs} \leq \beta \quad (9)$$

For n flows:

$$n\alpha \leq \frac{\tau}{T_{min}} \sum r_{obs} \leq n\beta \quad (10)$$



Vegas Latency Scaling



Simulated Queue Occupancy for 2, 8 and 16 vegas flows.



Bandwidth=5Mbps, delay=30ms, queue=100 packets

Vegas

Vegas has a number of problems:

- Accurate *baseRTT* is critical to Vegas or it will underestimate the window size.
- Router queue occupancy scales with the number of flows, ie. Vegas doesn't maintain low delay, it only uses the delay as a signal.
- Vegas responds to any delay, whether or not it is the cause.

Even in ordinary network environments, with enough Vegas flows, persistent queueing occurs and queues can overflow.





Delay-Based AIMD

An Alternative Approach



Delay-based AIMD

Four main components:

1. Extra congestion event to react to queueing delay.
2. Modified β to drain queues
3. Modified α to improve congestion recovery
4. Experimental solutions for coexistence with loss-based flows.



Extra congestion event

The congestion window is updated similarly to Reno, except to add an extra congestion event where $Cwnd > w_0$, a minimum window size:

$$Cwnd \leftarrow \begin{cases} Cwnd * \beta & loss \\ Cwnd * \beta & \tau > \tau_0 \\ Cwnd + \frac{\alpha}{Cwnd} & ACK \end{cases} \quad (11)$$

τ and τ_0 are the observed and threshold queueing delay.



Draining the Queue

β , the backoff factor, is designed so as to empty the queue at every congestion event.

$$\beta = \delta \frac{RTT_{min}}{RTT(t)} \quad (12)$$

In practice, an extra factor $\delta \sim 0.9$ is added to ensure RTT_{min} is regularly seen.



Recovering Quickly

The increase function, α is quadratic in time since the last back off, as in H-TCP and is balanced against β to maintain fairness.

$$q = \min[1, 1 + 10(\delta - 1) + 0.5(\delta - 1)^2] \quad (13)$$

$$\alpha = 2(1 - \beta)q \quad (14)$$



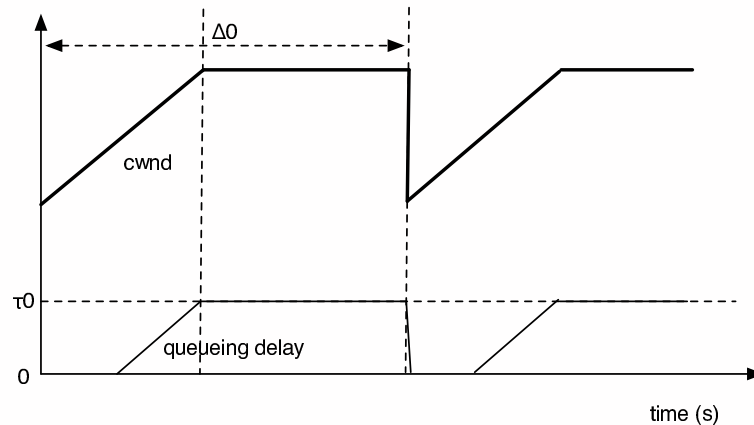
Coexistence with loss-based flows

Two experimental approaches thus far:

- Sliding delay threshold, τ_0
- Probabilistic losses at network endpoints



Coexistence with loss-based flows



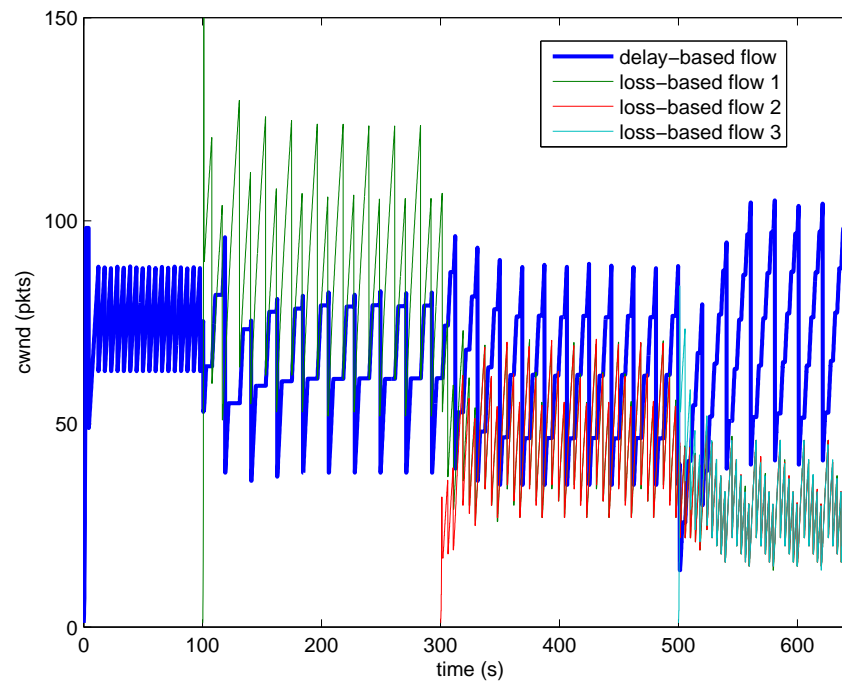
- Sliding delay threshold ($0 < \gamma < 1$)

$$\tau_0 = (1 - \gamma)\tau_0 + \gamma(RTT_{max} - RTT_{min}) \quad (15)$$

- When $\tau > \tau_0$, stop increasing but delay back off until time since last back off $> \Delta_0$.



Coexistence with loss-based flows

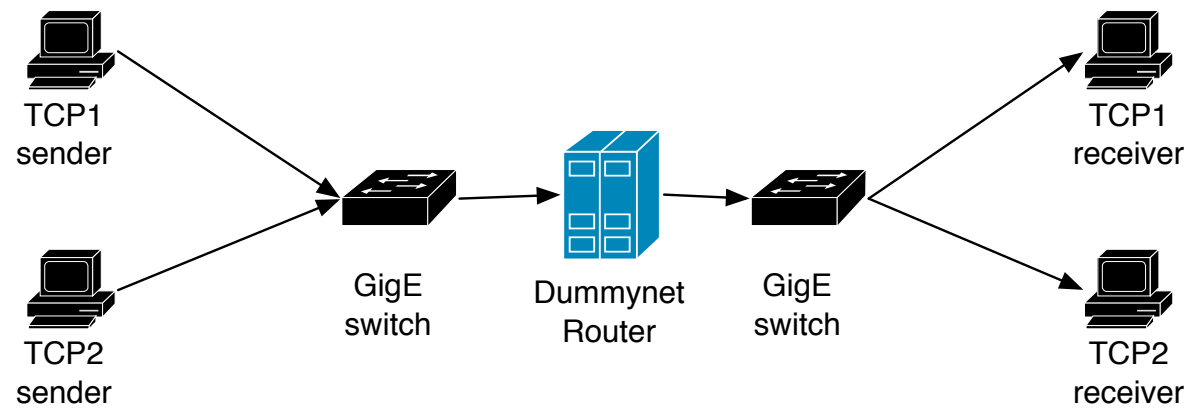


Coexistence of 1 delay and 3 loss-based flows (ns simulation, propagation delay 100ms, bandwidth 10Mbps, 100 packet queue).

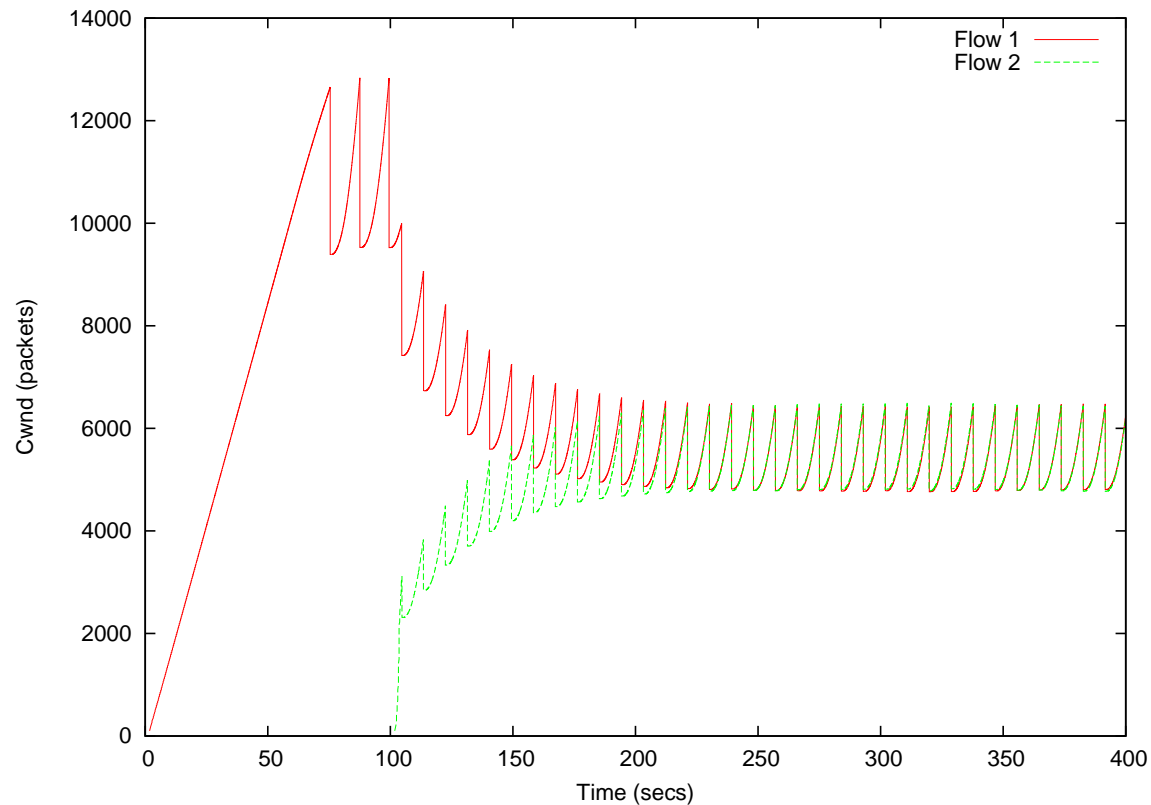


Results

- Delay-based in isolation
- Results from tests on dumbbell topology testbed, linux 2.6.18 source and destination hosts, freebsd dummysnet router.



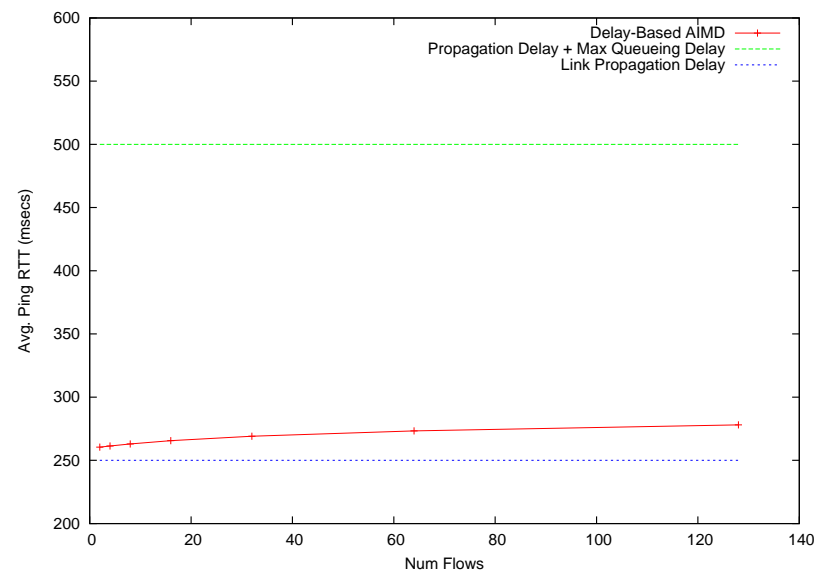
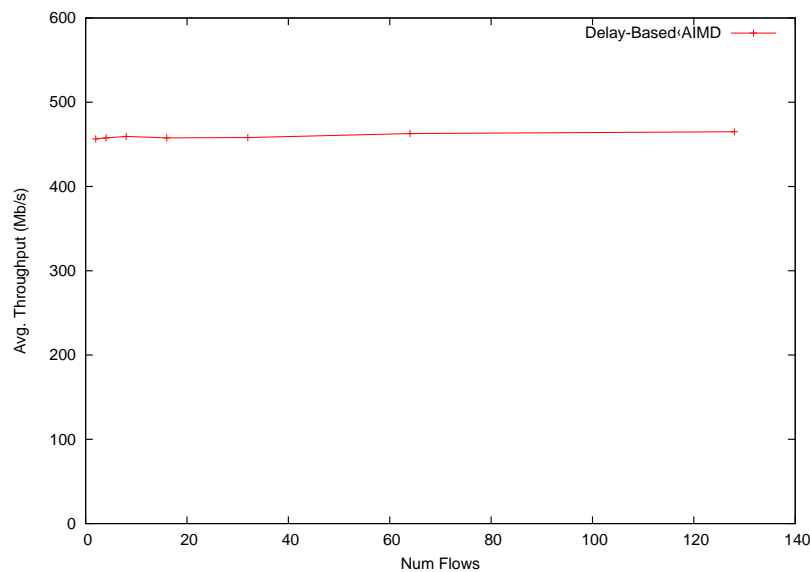
High Speed Link



Convergence of DB-AIMD following startup of a second flow.
500Mbps link rate, 250ms RTT, 250ms of buffering.



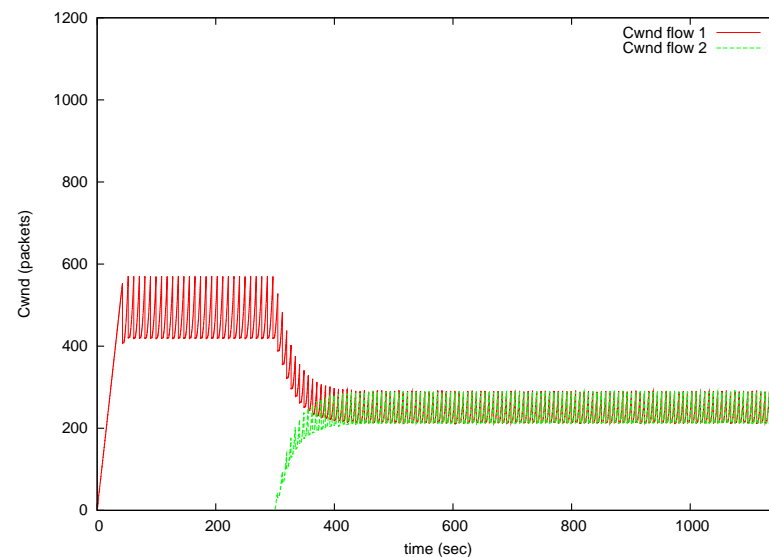
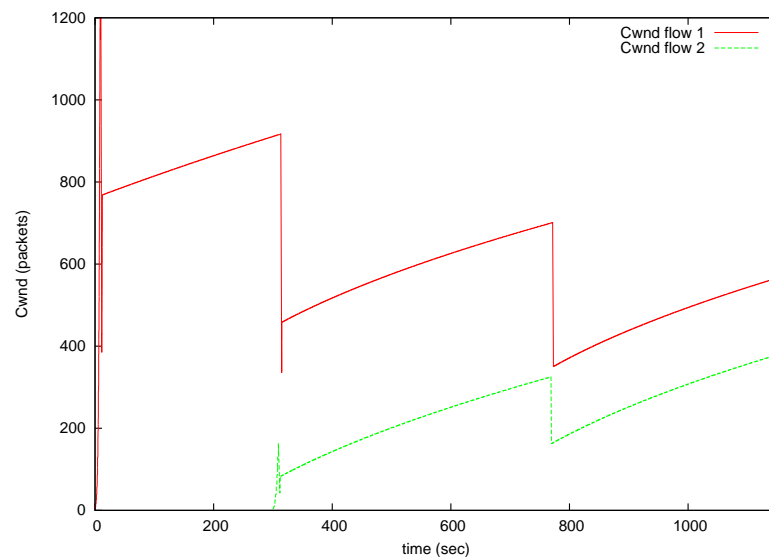
High Speed Link



Link utilisation and Delay vs number of flows with DB-AIMD
(using $\tau_0=50\text{ms}$). 500Mbps link rate, 250ms RTT,
bandwidth-delay product of buffering.



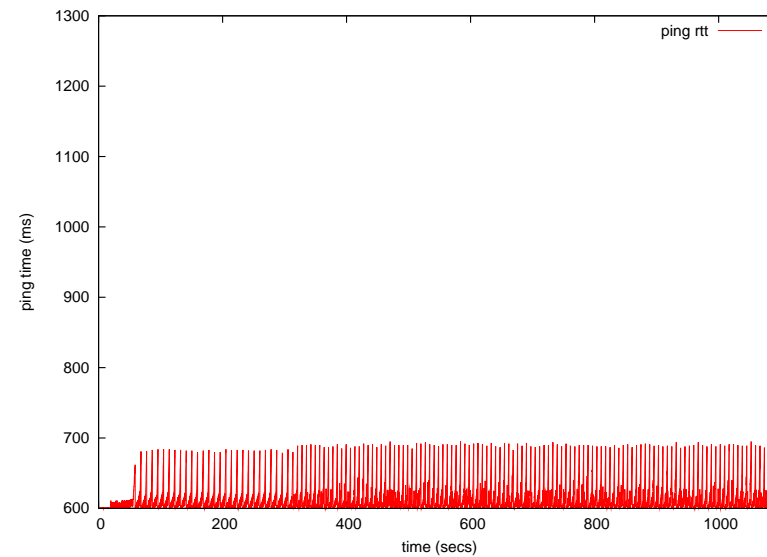
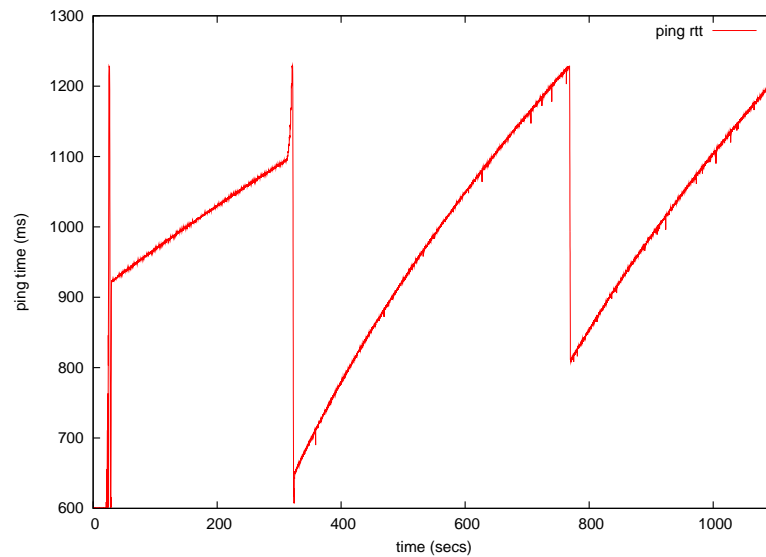
GeoSat Link



Congestion windows following startup of a second flow. Reno (left) and DB-AIMD (right). Note Reno's 15 min. congestion epoch duration. 10Mbps bandwidth, 600ms propagation delay, 600ms queue



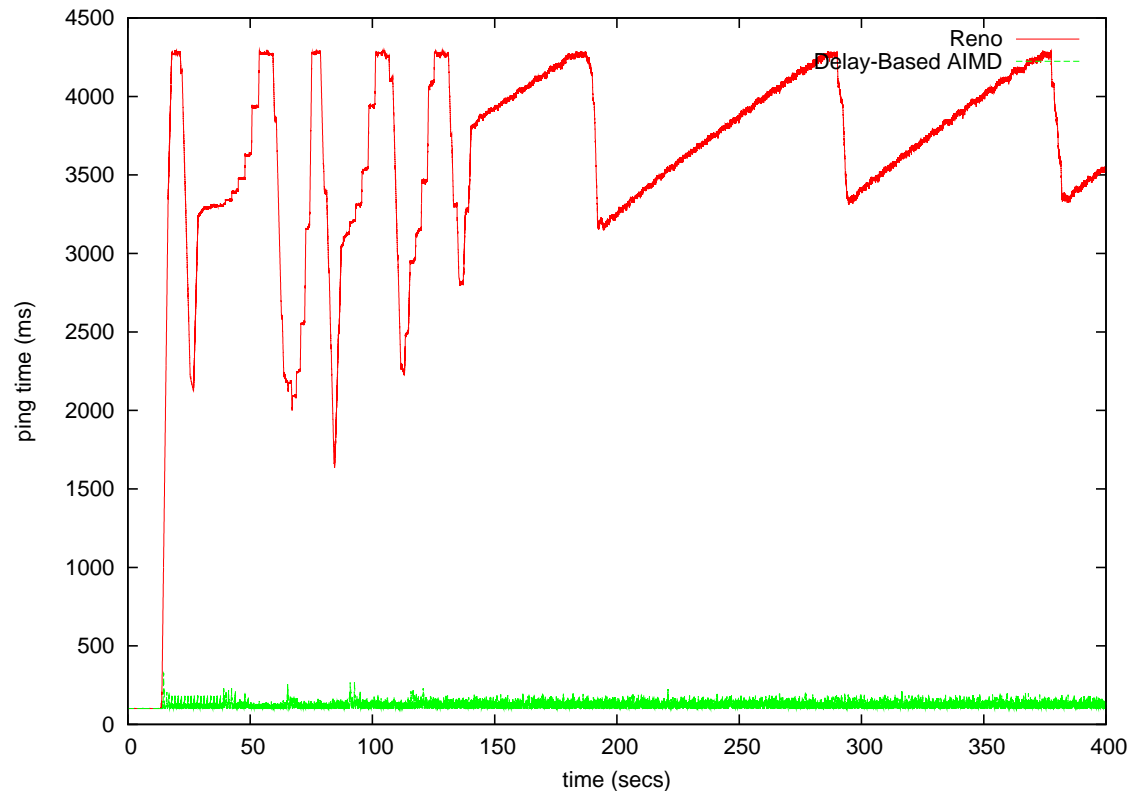
GeoSat Link



Ping times following startup of a second flow with Reno (left) and DB-AIMD (right). 10Mbps bandwidth, 600ms propagation delay, 600ms queue.



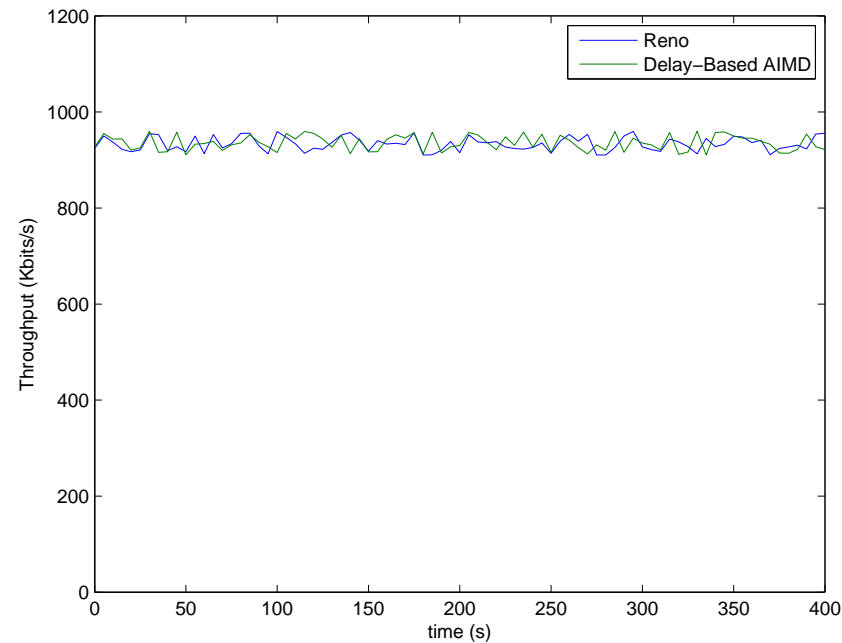
DSL Link



Ping times for a link with 1Mbps bandwidth, 100ms propagation delay, 512KB queue. Staggered flows start every 30 secs



DSL Link



Total throughput of multiple staggered flows with Reno and DB-AIMD. 1Mbps bandwidth, 100ms propagation delay, 512KB queue.



Review of Design Objectives

- At each congestion event, drain the queue completely
... yes
- Maintain low queueing delay throughout operation
... yes
- Quick startup of new flows
... yes
- Full utilisation of large BDP links
... almost
- Coexistence
... feasible?



Outstanding Issues

- Reverse-path queueing
- Several queues with competing traffic — multiple bottlenecks
- Importance of correlation between delay and congestion
- Spurious delay signals
- Route changes



Conclusions

- Unlike previous delay-based schemes, DB-AIMD actively drains router queues, lowering delay
- Coexistence may be possible
- Vegas may have some issues but delay-based schemes should not be ignored
- Some overlap in work by: MC Weigle, K Jeffay, F Donelson Smith - Computer Communications, 2005, also Westwood, FAST

