



Hamilton Institute

TCP and Network Congestion Control

Amazon Dublin, December 2005

(David Malone and Doug Leith)



Hamilton Institute

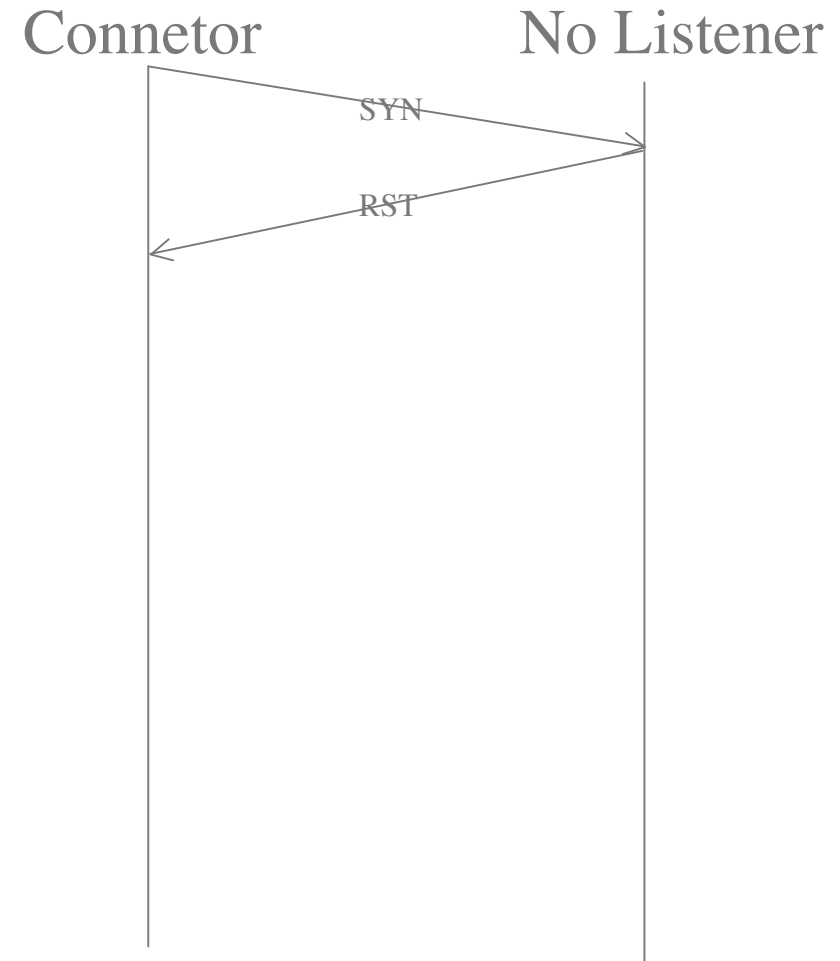
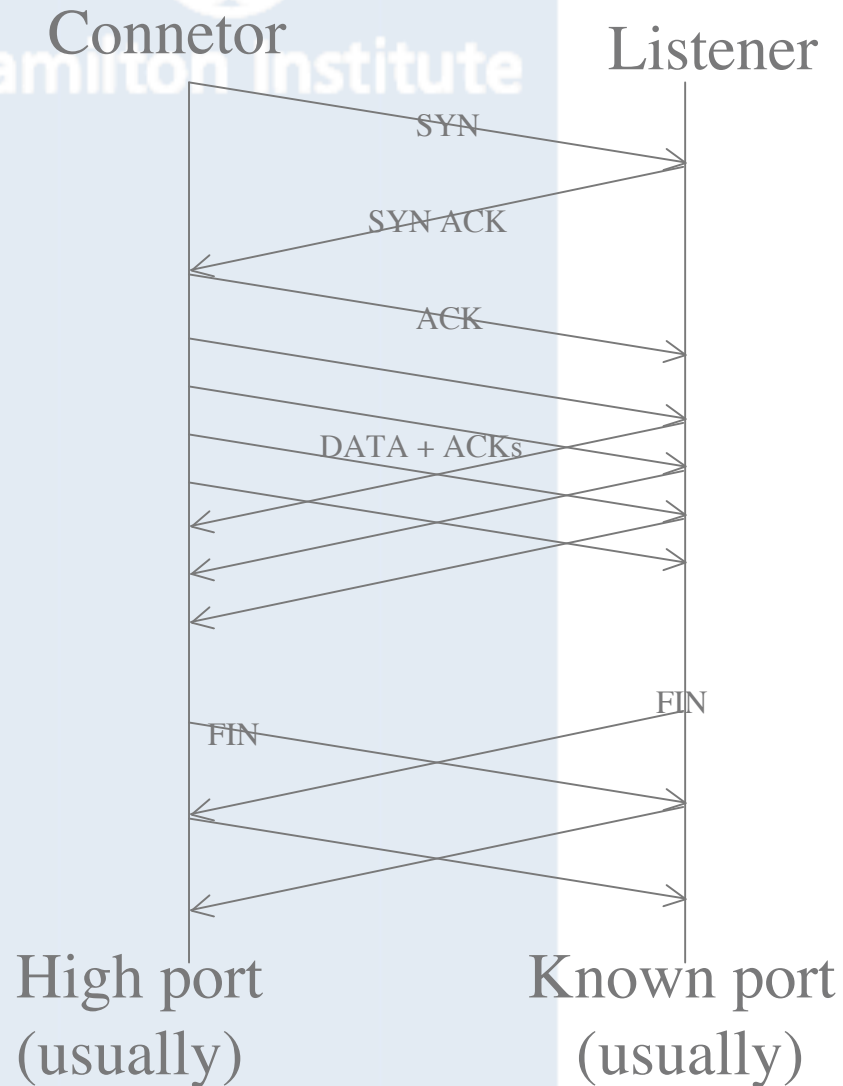
What TCP does for us

- Demuxes applications (using port numbers).
- Makes sure lost data is retransmitted.
- Delivers data to application in order.
- Engages in congestion control.
- Allows a little out-of-band data.
- Some weird stuff in TCP options.



Hamilton Institute

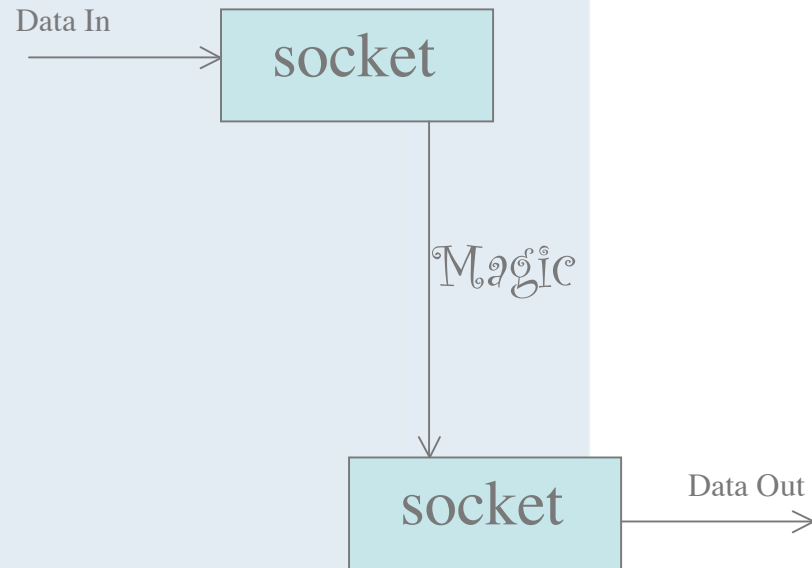
Standard Picture of TCP



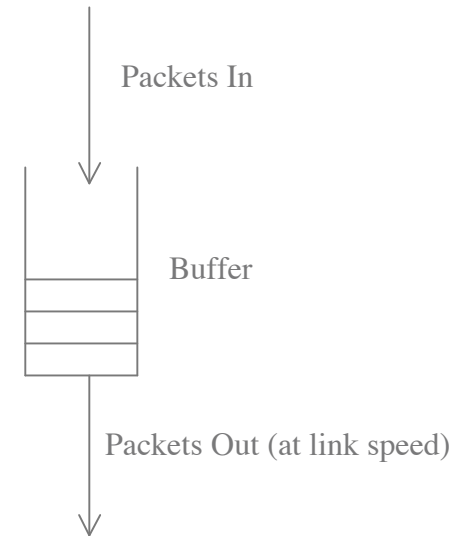


Hamilton Institute

Other views of TCP



Programmer's View



Network View



Hamilton Institute

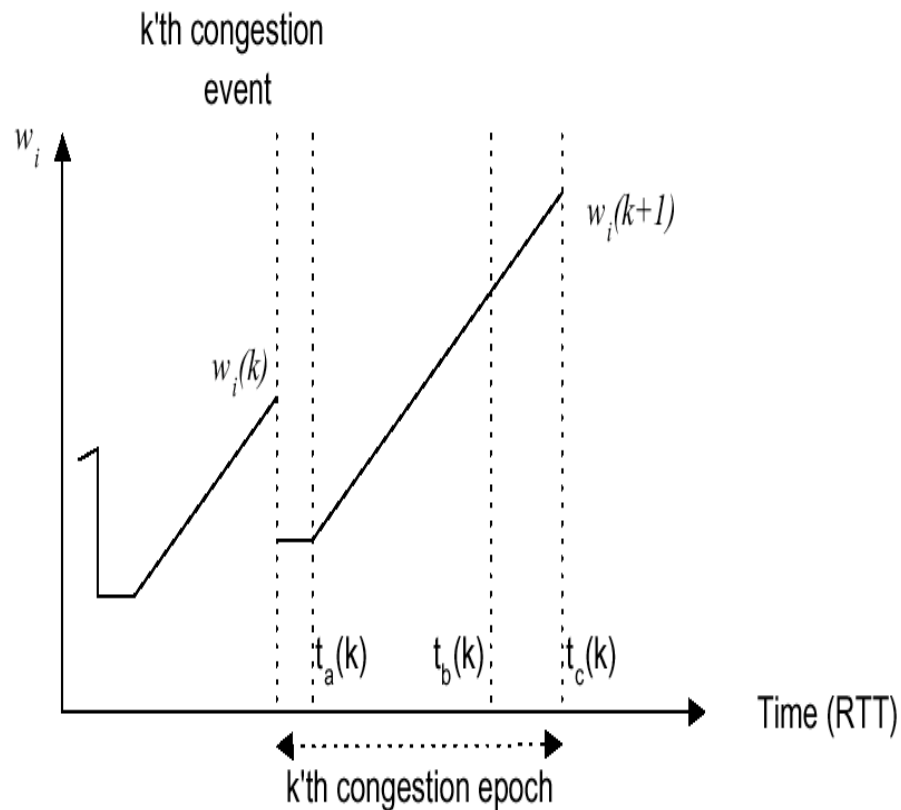
Congestion Control

- TCP controls the number of packets in the network.
- Packets are acknowledged, so flow of ACKs.
- Receiver advertises window to avoid overflow.
- Congestion window tries to adapt to network.
- Slow start mechanism to find rough link capacity.
- Congestion avoidance to gradually adapt.



Hamilton Institute

The Congestion Window



- Additive increase, multiplicative decrease (AIMD).
- To fill link need to reach $BW \times \text{Delay}$.
- Backoff by $1/2 \Rightarrow$ buffer at bottleneck link should be $BW \times \text{Delay}$.
- Fairness (responsiveness, stability, ...)



Hamilton Institute

Basic TCP tuning

- Network stack has to buffer in-flight data.
- Need $BW \times \text{delay}$ sized sockbuf!
- `/proc/net/core/{r,w}mem_max` store something like sockbuf sizes.
- `/proc/net/ipv4/tcp_{r,w,}mem` store min/def/max for tcp windows.
- (or sockopt `SO_SNDBUF/SO_RCVBUF`).
- For large transfers, crank up to few MB.
- Memory will be wired, so need to use balance.



Hamilton Institute

Existing TCP mods

- Window scaling.
- TCP traditionally ACKs last contiguous byte. SACK - transmits information about gaps.
- TCP usually uses drops as feedback signal. ECN allows use of few bits in IP header.
- Timestamps: more accurate RTT estimates.
- MD5 checksums.



Problems for Congestion Control

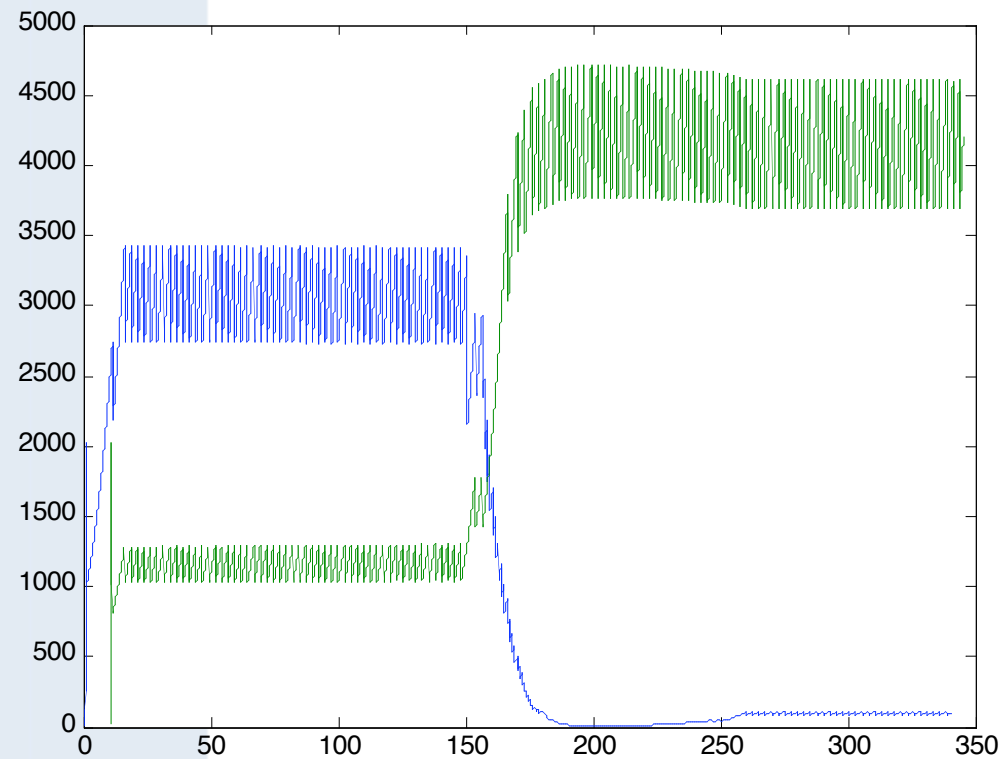
Hamilton Institute

- Packet loss caused by other factors.
- Filling a big link at one-packet-per-round-trip.
- The combination is bad for high speed long distance links.
- Problem was flagged up: various solutions being studied (BIC, Scalable TCP, High-Speed TCP, FAST TCP, H-TCP, ...)



Hamilton Institute

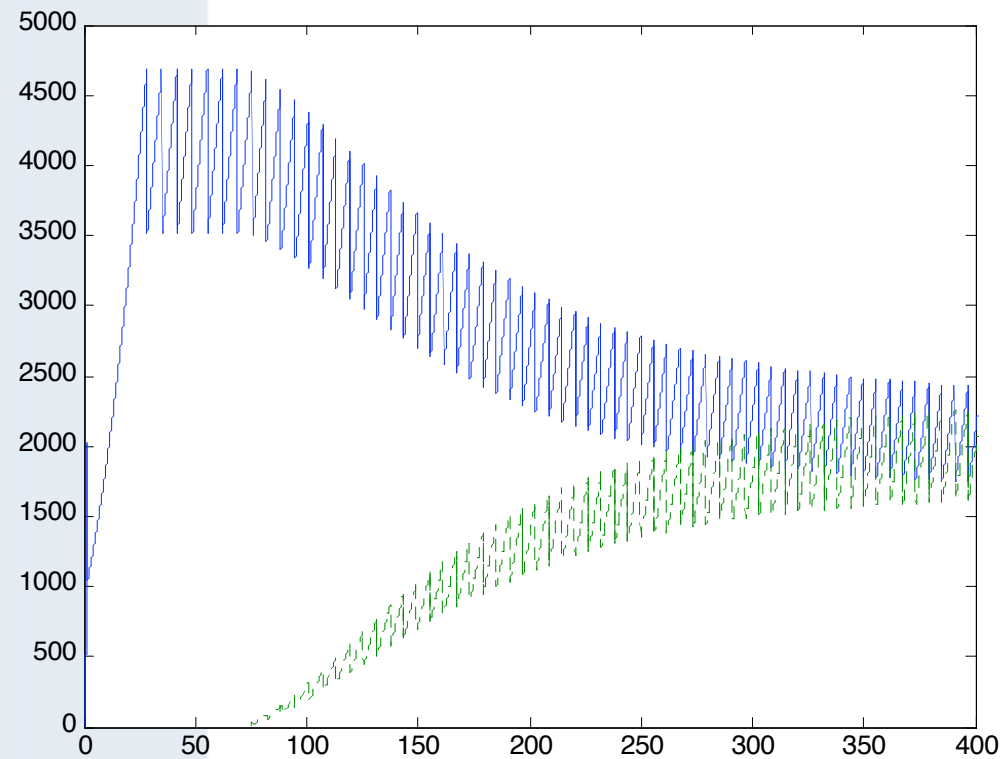
Stability issues





Hamilton Institute

Convergence Issues

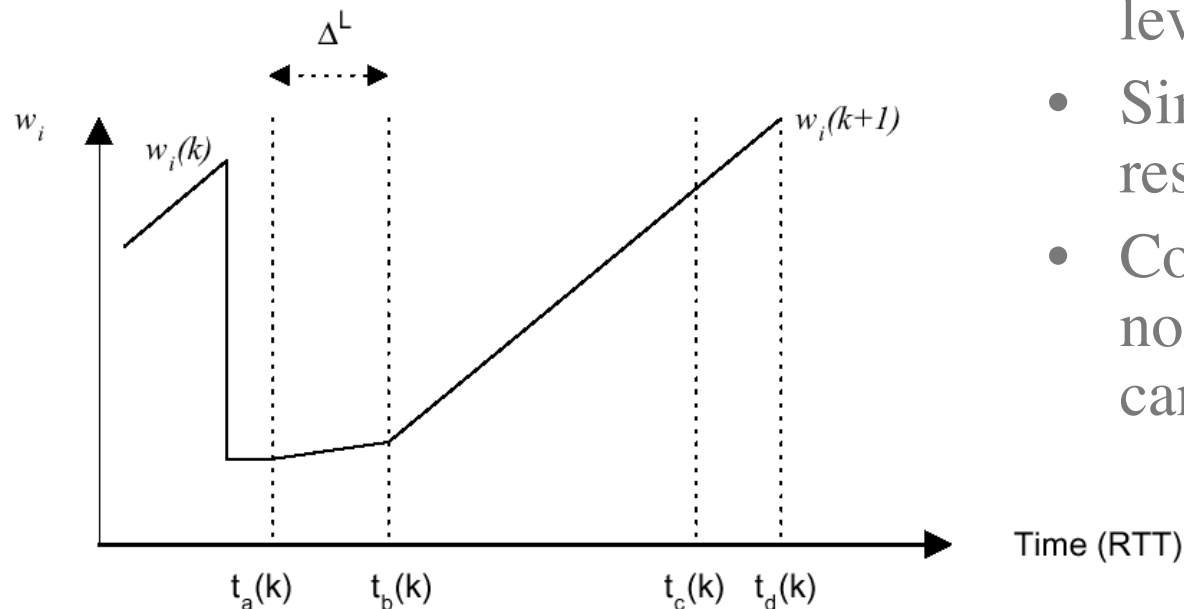




Hamilton Institute

H-TCP

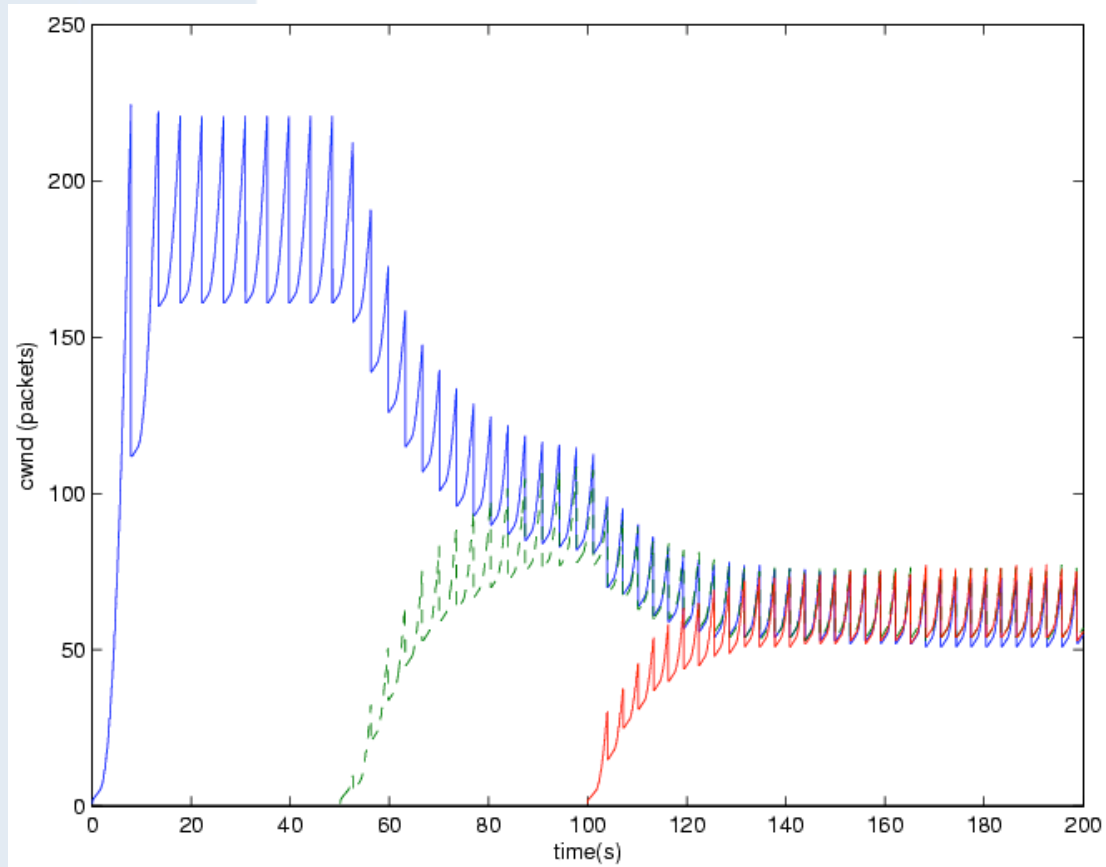
- Aim to make small changes we can analyse.
- Rate of increase depends on how long since last backoff.
- New flows compete on level playing field.
- Similar fairness & responsiveness.
- Competes fairly with normal TCP where it can compete.





Hamilton Institute

Quicker Convergence





Hamilton Institute

More Linux tuning

- As of 2.6.13 Linux allows you to choose the congestion control technique.
- Hidden behind TCP_CONG_ADVANCED.
- Can use `/proc/sys/net/ipv4/tcp_congestion_control`
- Older versions can disable/enable with `/proc/sys/net/ipv4/tcp_{bic,vegas_cong_avoid,westwood}`
- Some bugs fixed recently, so new kernels useful.



Hamilton Institute

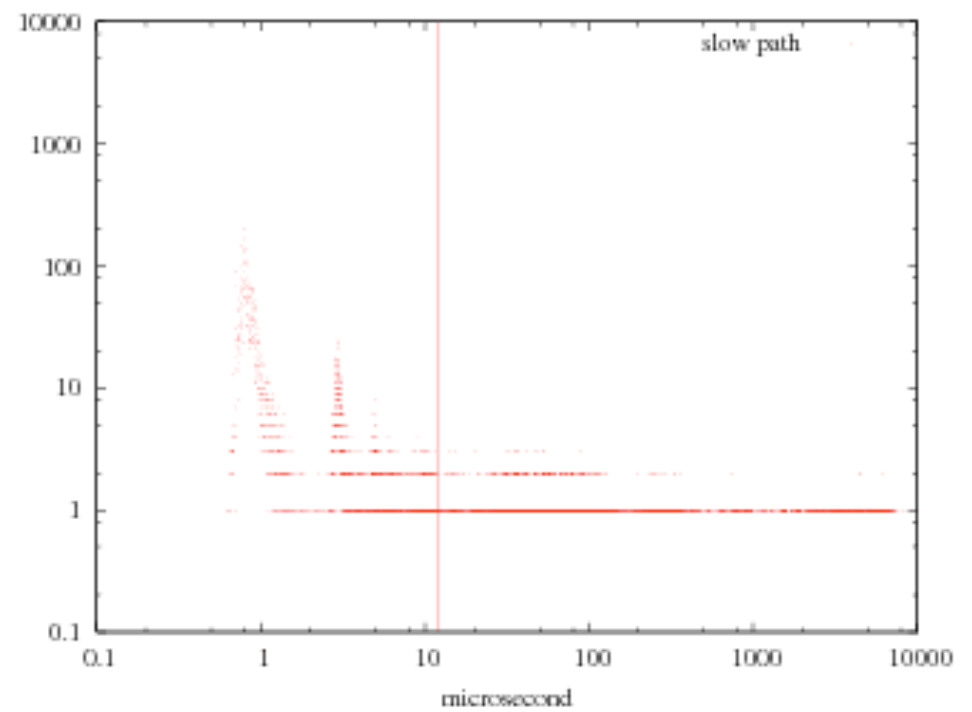
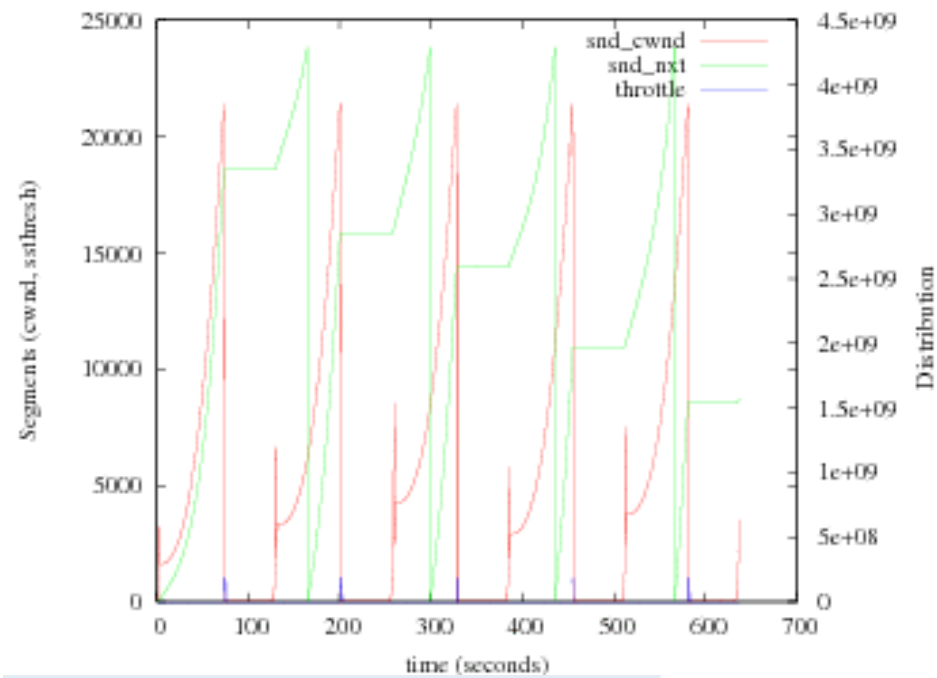
Practical Issues

- Congestion control isn't the only issue.
- Implementation is important.
- Testing is important: land speed records.
- Web 100 project to instrument Linux.
- Important stack tuning to be done.
- <http://www.psc.edu/networking/projects/pathdiag/>
- <http://www.csm.ornl.gov/~dunigan/netperf/web100.html>
- <http://www.psc.edu/networking/projects/hpn-ssh/>



Hamilton Institute

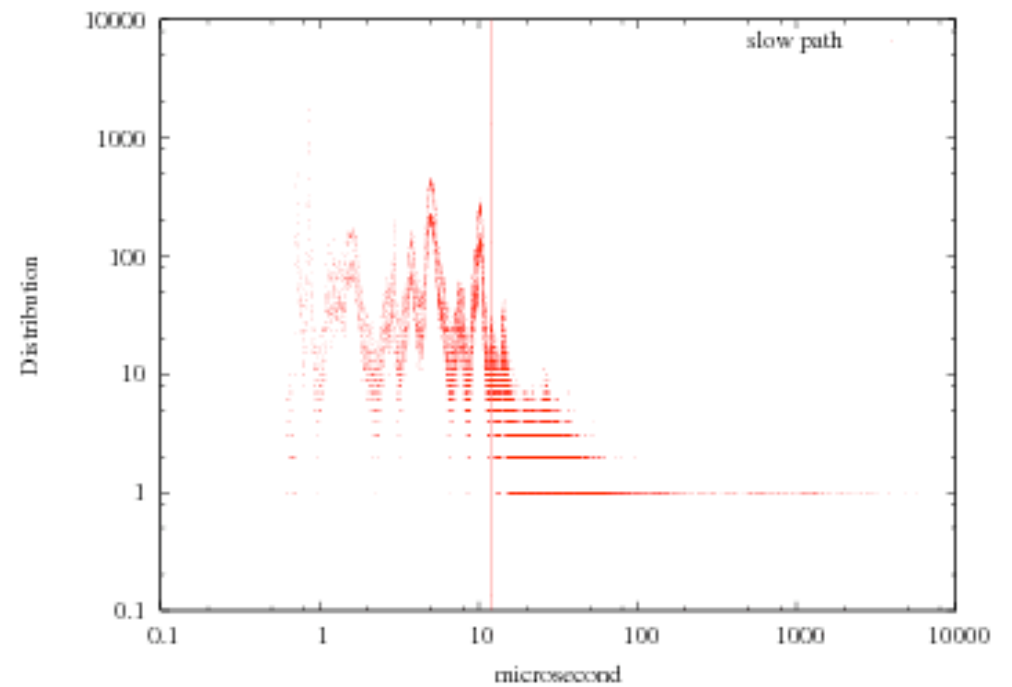
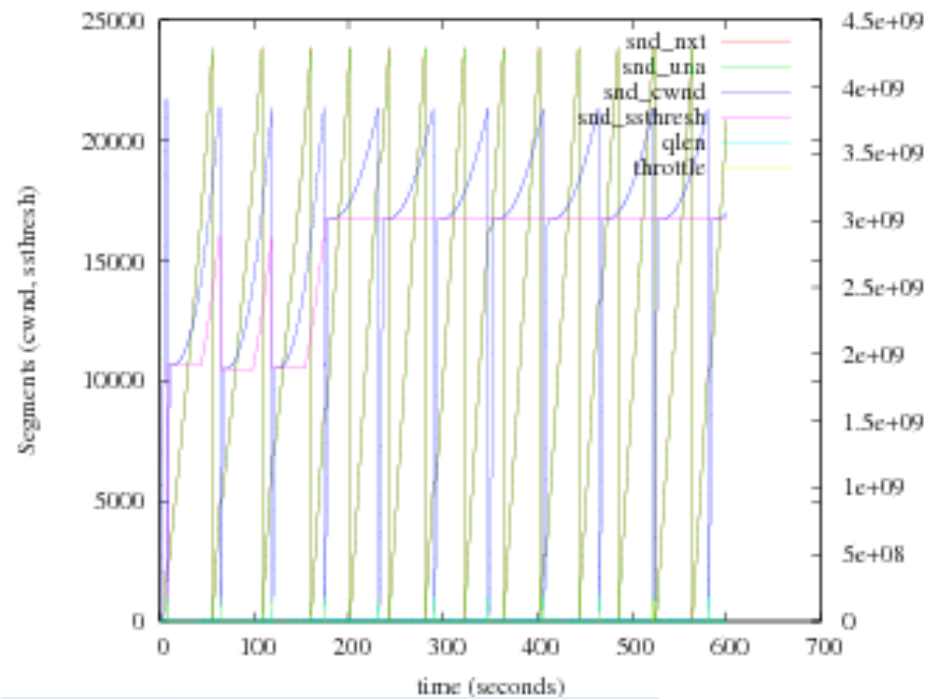
Before





Hamilton Institute

After





Hamilton Institute

Other issues

- High speed is important (packet switched vs. circuit switched), but not for everyone.
- Sizing router buffers important to everyone (cost, QoS, optics).
- Wireless interesting - random losses.
- Other interesting wireless issues too.
- Many flows don't leave slow start => Quick start.
- For really small flows handshake is too much: T/TCP2.



Hamilton Institute

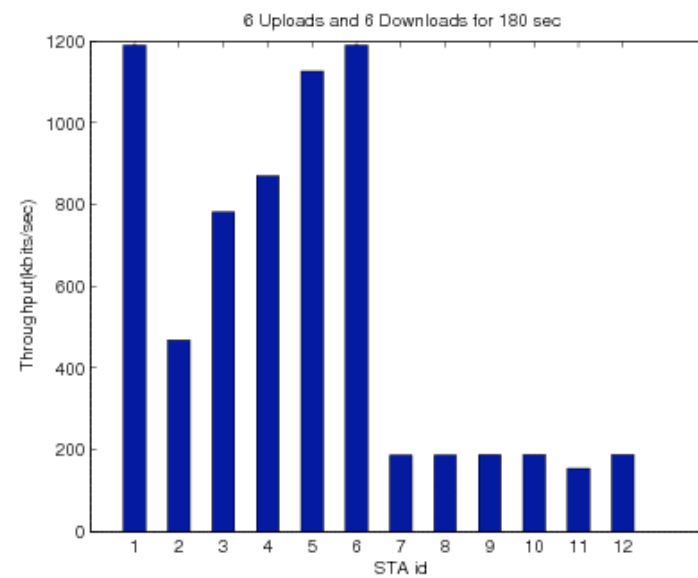
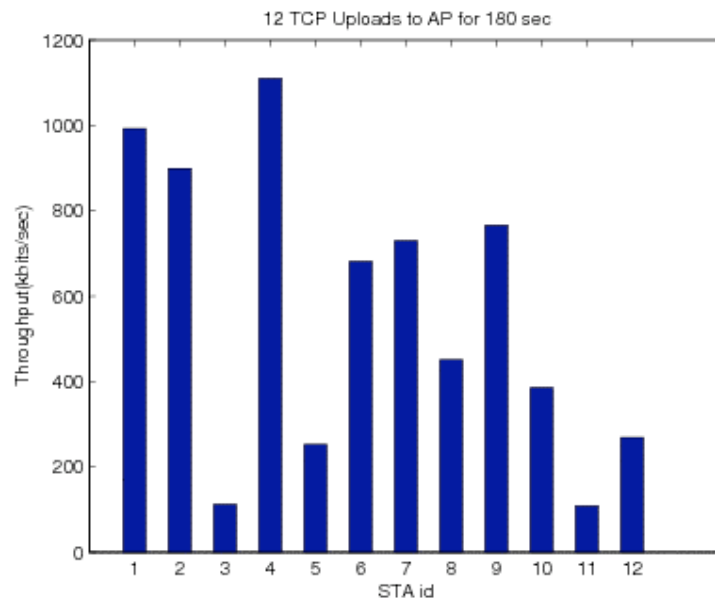
Thanks!

Questions?



Hamilton Institute

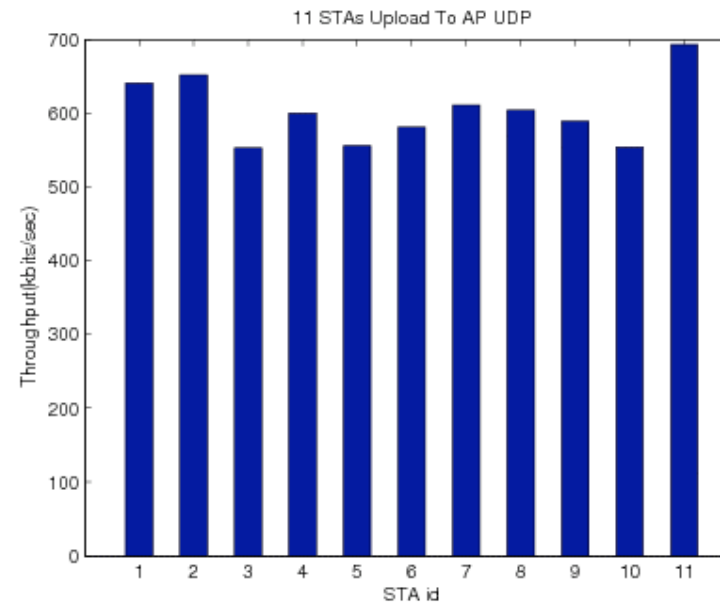
802.11 AP Before





Hamilton Institute

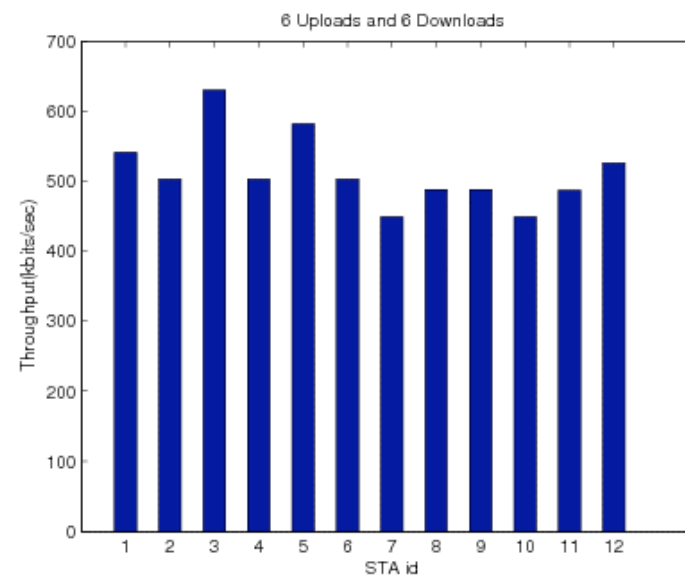
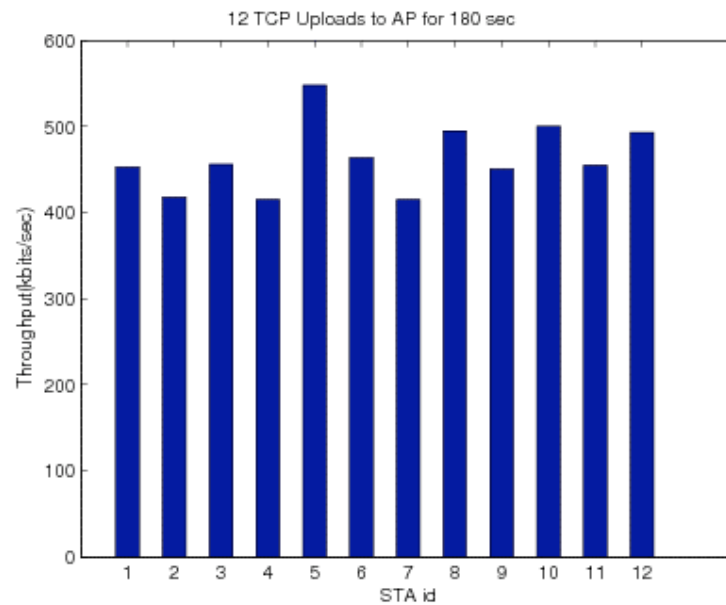
Baseline





Hamilton Institute

Adjusting 802.11e Parameters





Hamilton Institute

More Questions?

