# Network Research and Linux at the Hamilton Institute, NUIM.
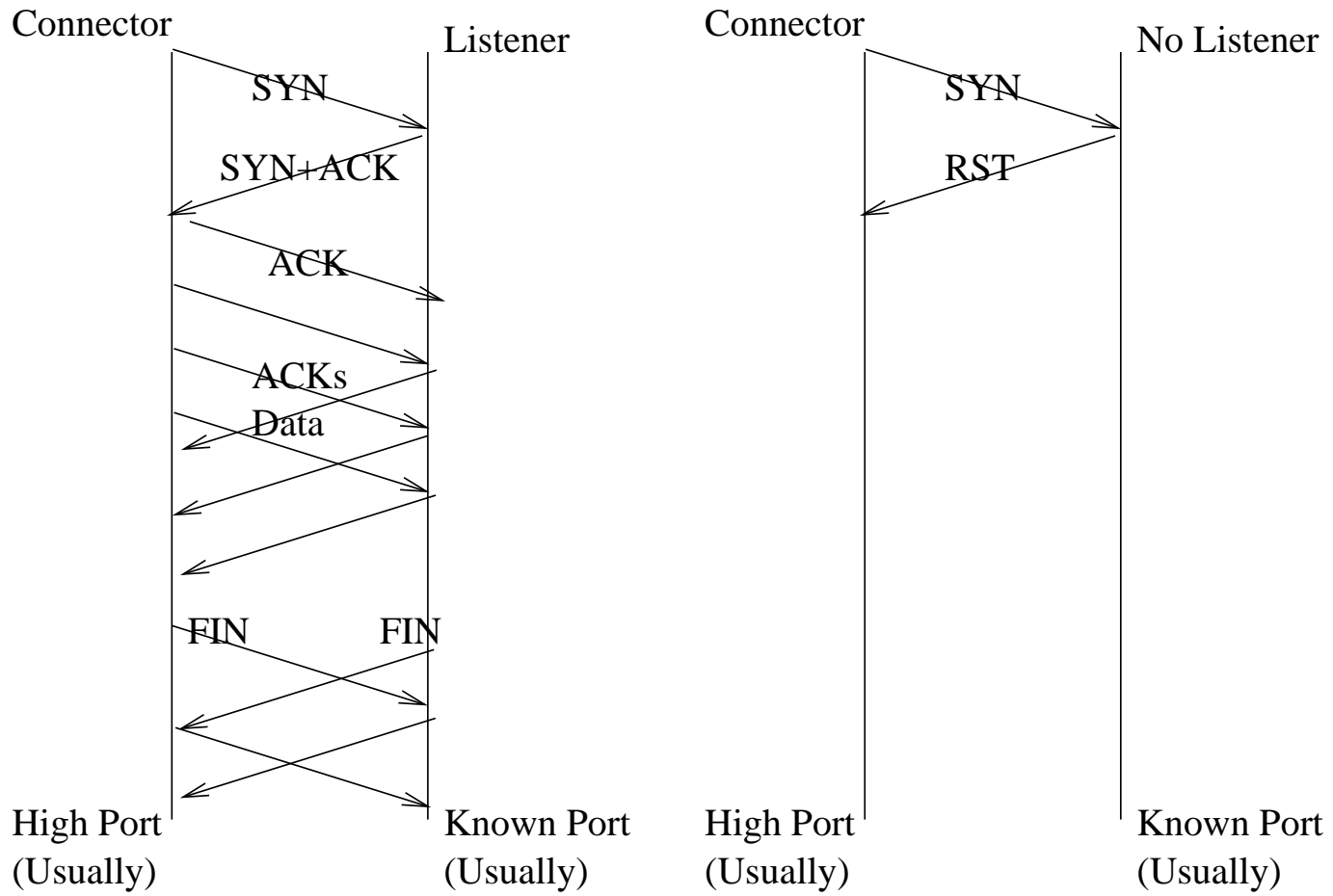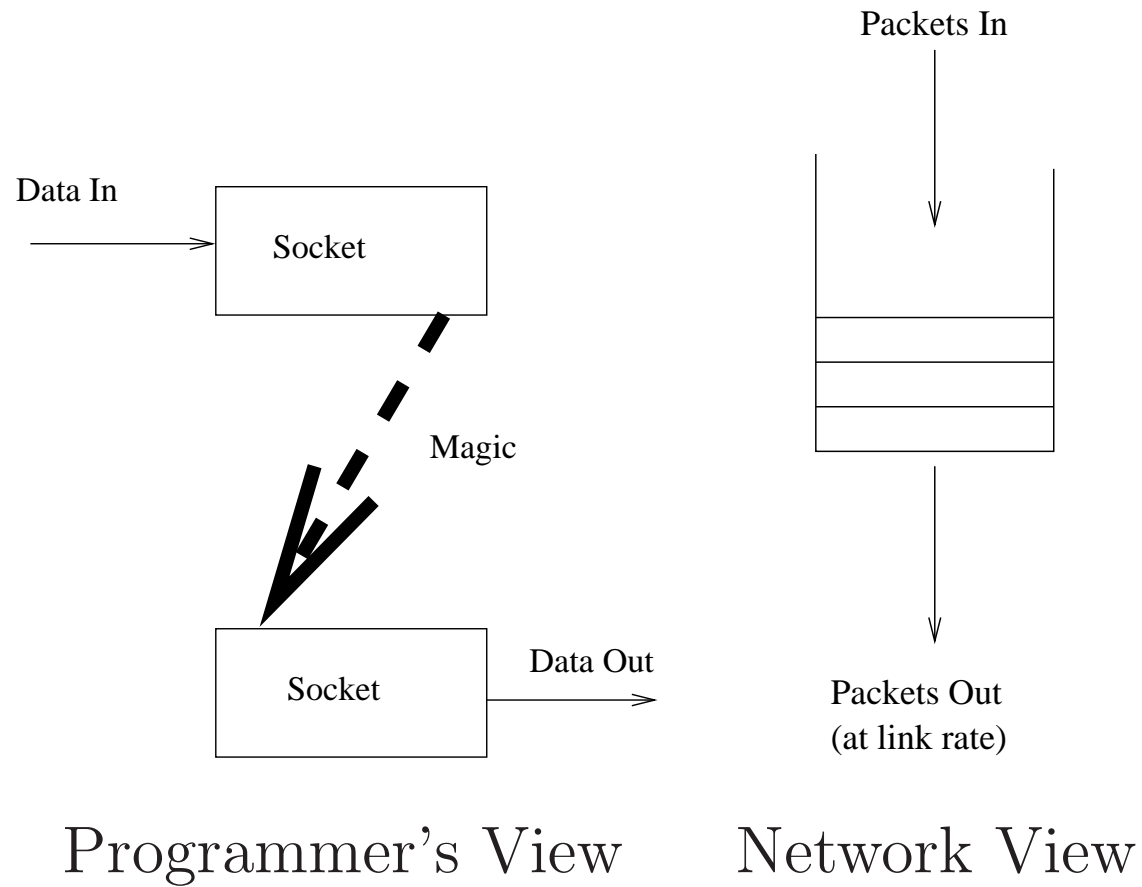
## David Malone

## 4 November 2006

## What has TCP ever done for Us?

- Demuxes applications (using port numbers).

- Makes sure lost data is retransmitted.

- Delivers data to application in order.

- Engages in congestion control.

- Allows OOB data.

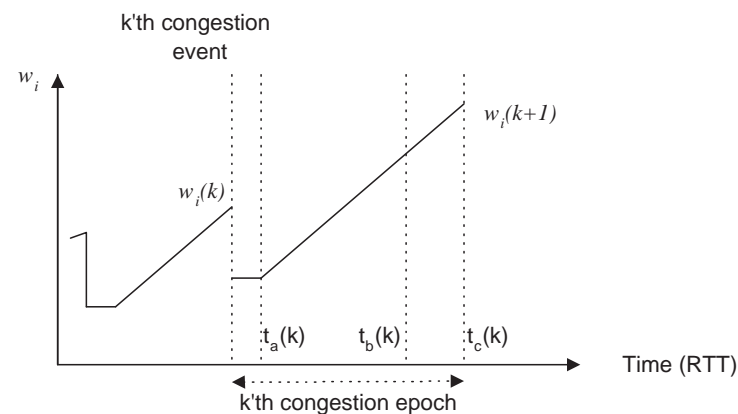- Some weird stuff with TCP options.

# Standard Picture of TCP

Connector      Listener      Connector      No Listener

SYN

SYN+ACK

ACK

ACKs

Data

FIN     FIN

SYN

RST

High Port      Known Port      High Port      Known Port
(Usually)      (Usually)      (Usually)      (Usually)

# Other Views of TCP

Packets In

Data In

Socket

Magic

Socket → Data Out

Packets Out
(at link rate)

Programmer's View          Network View

# Congestion Control

- TCP controls number of packets in network.

- Packets are acknowledged, so flow of ACKs.

- Receiver advertises window to avoid overflow.

- Congestion window tries to adapt to network.

- Slow start to roughly find capacity.

- Congestion avoidance gradually adapts.

# The Congestion Window



- Additive increase, multiplicative decrease (AIMD).

- To fill link need to reach $BW \times \tau$.

- Backoff by $1/2$, implies buffer is $BW \times \tau$.
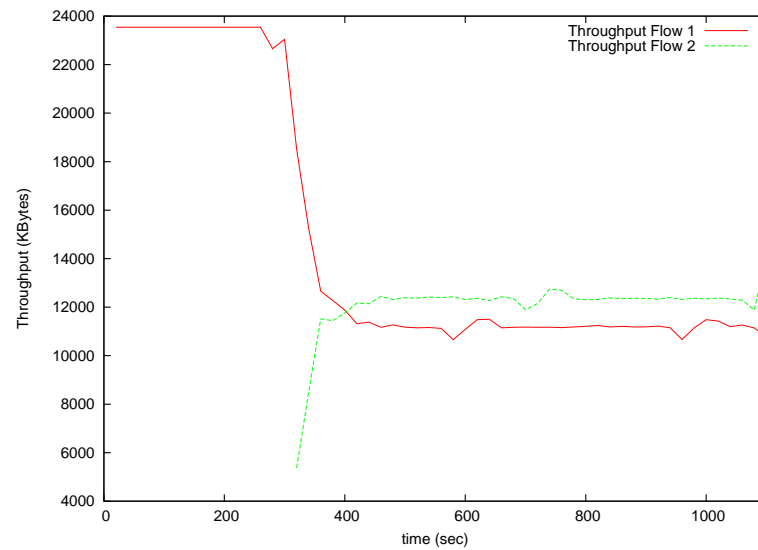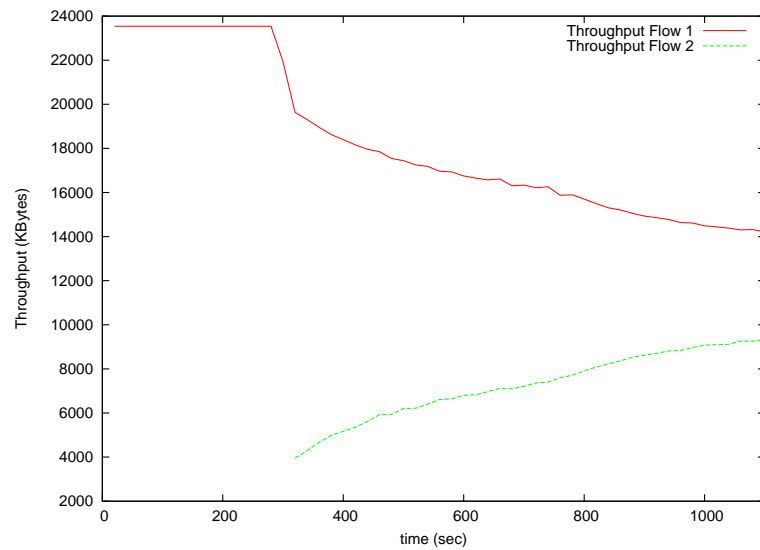
- Fairness, responsiveness, stability, . . .

## TCP On Linux

- Network stack buffers in-flight data.

- Socket buffer must be $BW \times \tau$.

- `/proc/net/core/{r,w}mem_max` $\rightarrow$ sockbuf sizes.

- `/proc/net/ipv4/tcp_{r,w,}mem` $\rightarrow$ min/def/max tcp window.

- Trade off — memory is wired, so valuable.
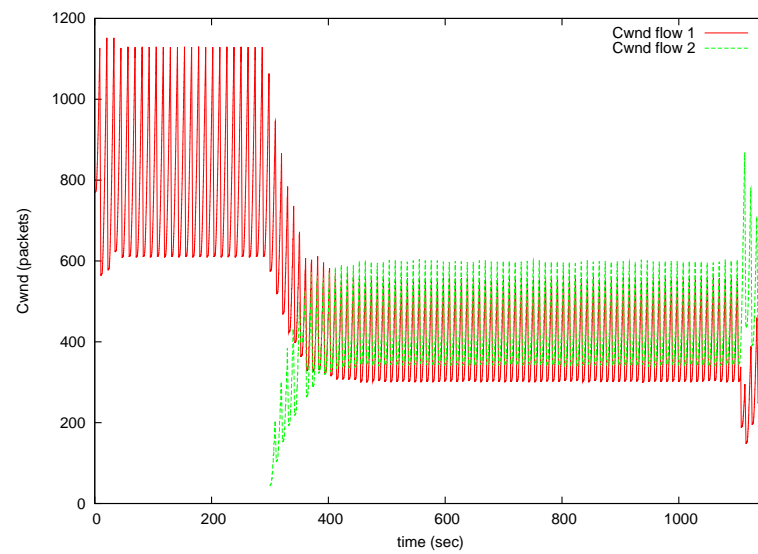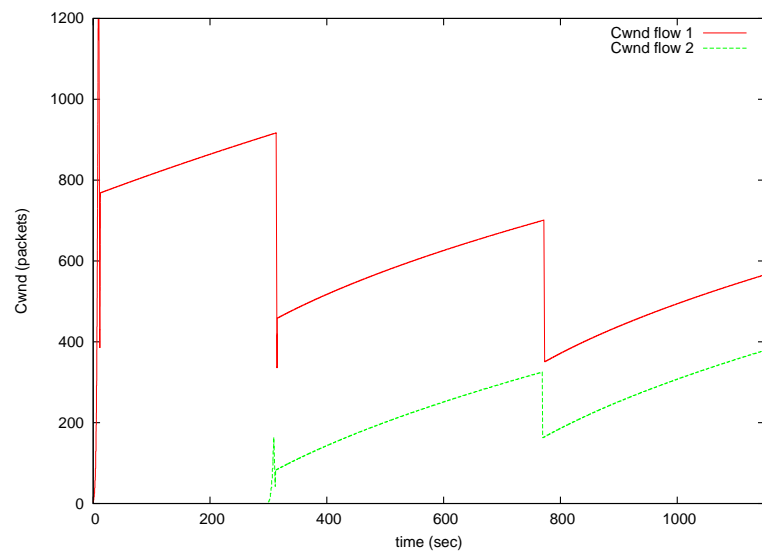
- Defaults have recently been increased.

# Research Work

- Packet loss not caused by congestion.

- Filling big $BW \times \tau$ product packet at a time.

- Bad for long-distance high-bandwidth links.

- Various solutions in pipeline (BIC, Scalable, High-Speed, FAST, H-TCP).

- Pluggable congestion control in Linux (behind `TCP_CONG_ADVANCED`).

- `/proc/sys/net/ipv4/tcp_congestion_control`

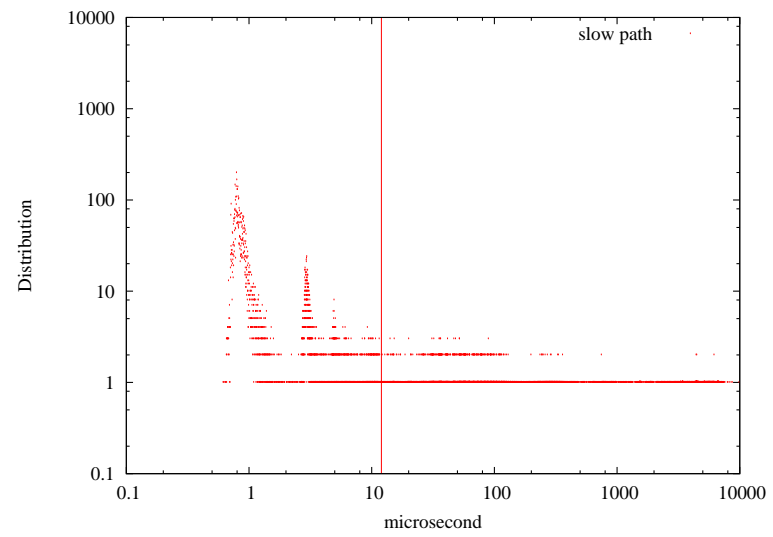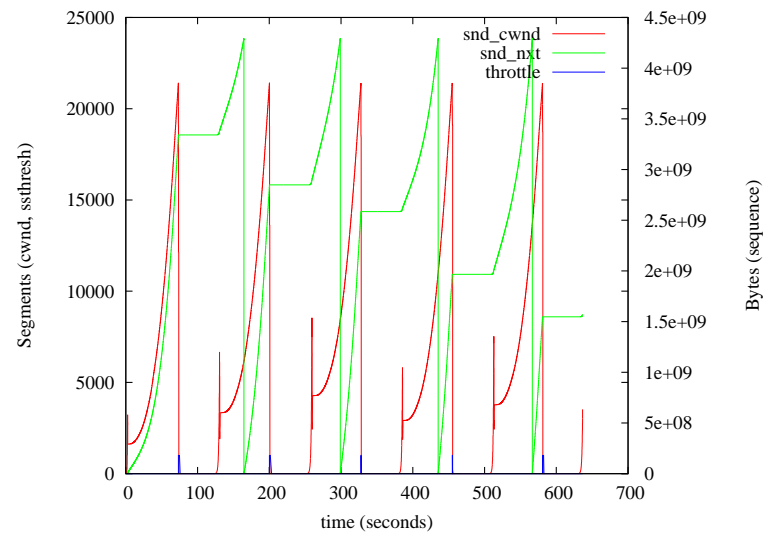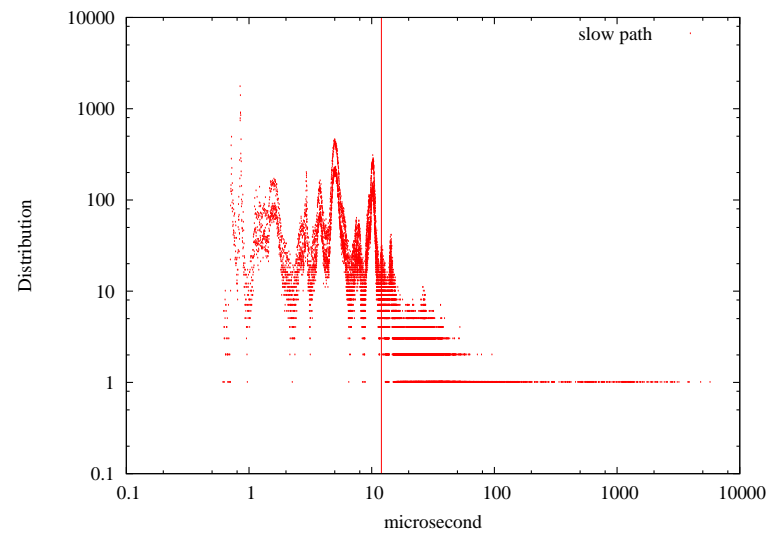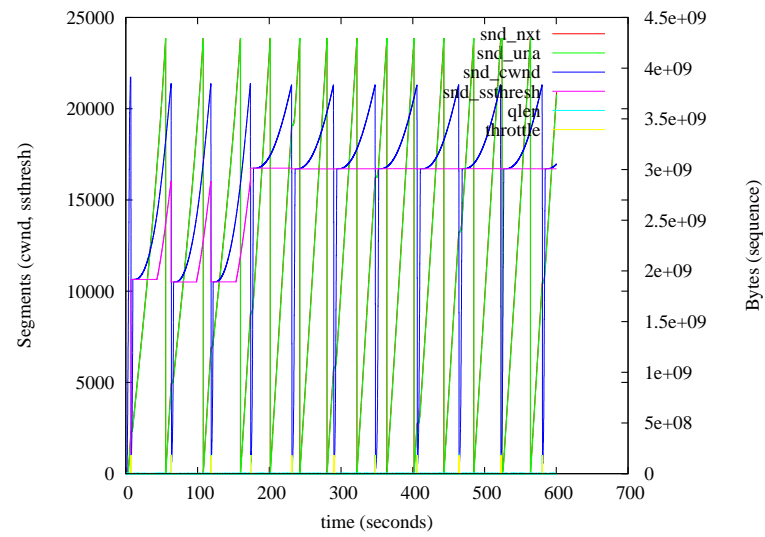- Working on other congestion detection techniques.

# Throughput

# Cwnd

## Practical Stuff

- Other issues at play, such as implementation quality.

- For example ACK processing and queueing problems.

- Testing is important: land speed records.

- Project with OSDL to build validation suite.

# Before

# After

# 802.11(e) MAC
# Summary

- After TX choose $\text{rand}(0, \text{CW} - 1)$.

- Wait until medium idle for $\text{DIFS}(50\mu s)$,

- While idle count down in slots $(20\mu s)$.

- TX when counter gets to 0, ACK after SIFS $(10\mu s)$.

- If ACK then $\text{CW} = \text{CW}_{\min}$ else $\text{CW} * = 2$.

Ideally produces even distribution of packet TX.

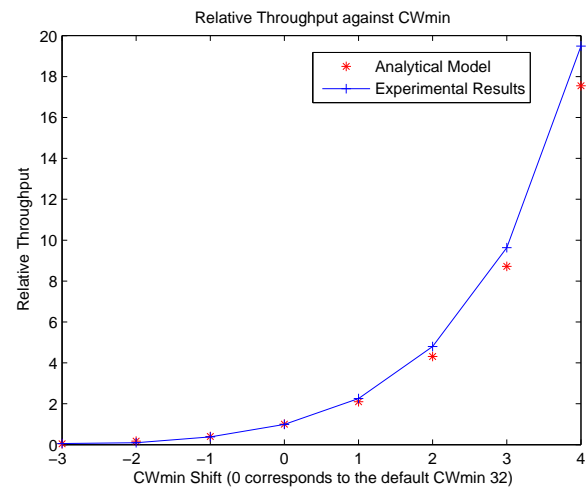In 11e have multiple queues. Each has own $\text{CW}_{\min}$, DIFS(aka AIFS) and can have TXOP.
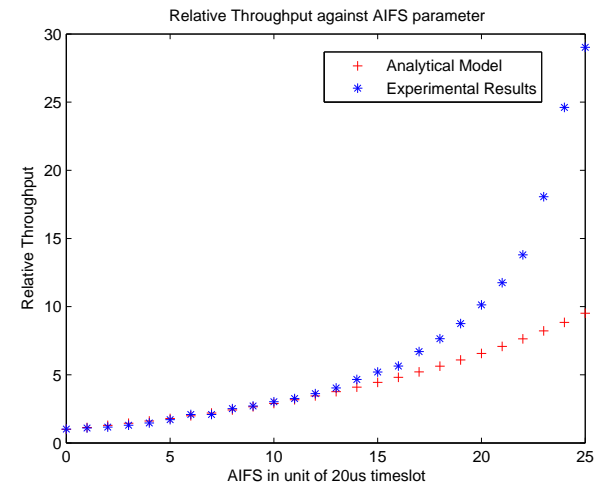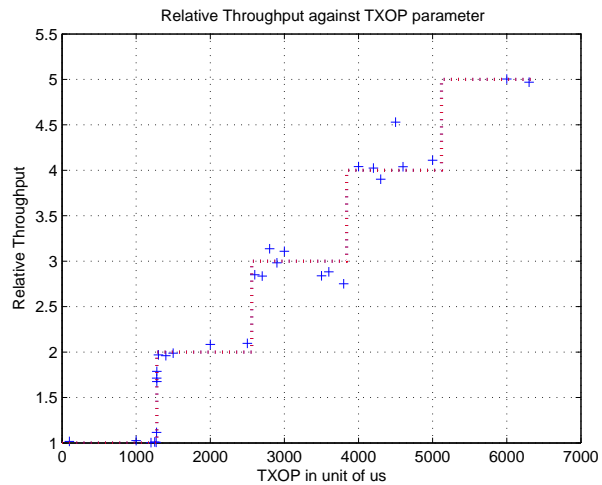
# Testbed setup

Multiple STA (Linux) connected to AP (Linux hostap).

| Hardware | model |
|----------|-------|
| 1× AP | Desktop PC |
| 18× STA | Soekris |
| 1× STA | Desktop PC |
| WLAN NIC | Atheros AR5212 |

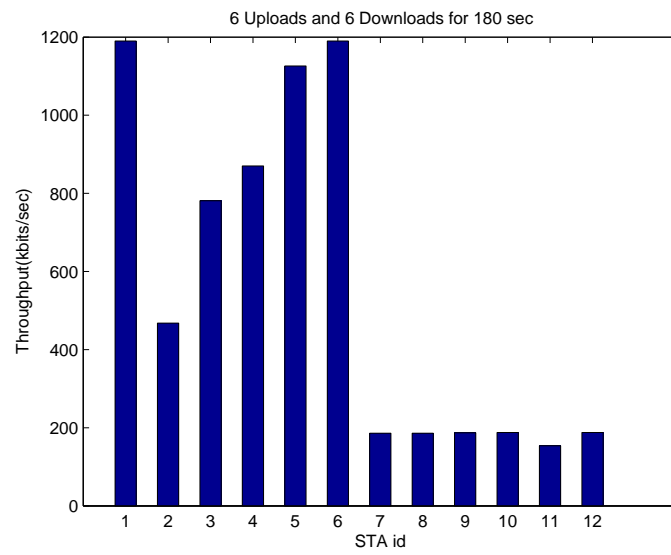External antenna, PCI interface, Madwifi driver with local patches for 11e parameter setting.

# Validation



Relative Throughput against TXOP parameter



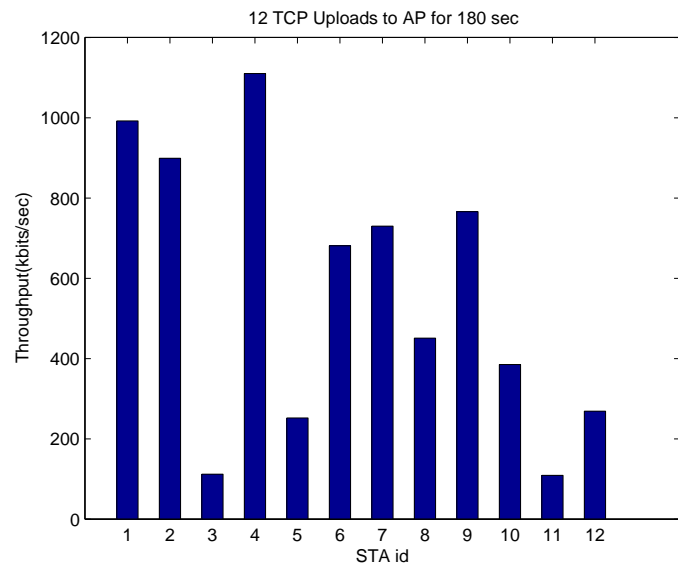Relative Throughput against AIFS parameter
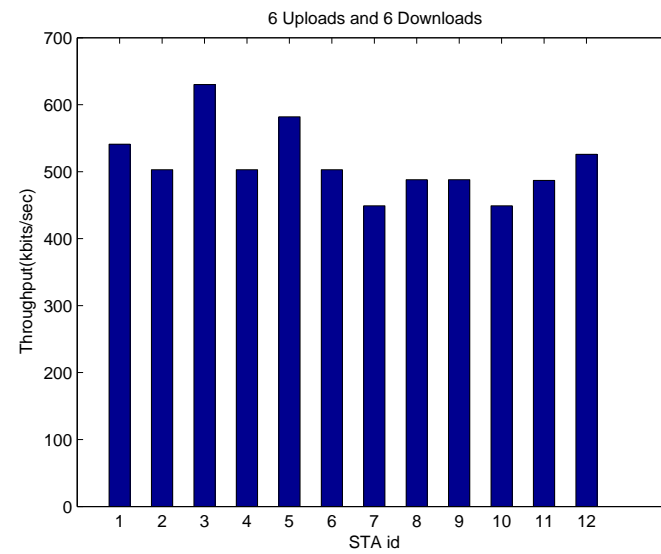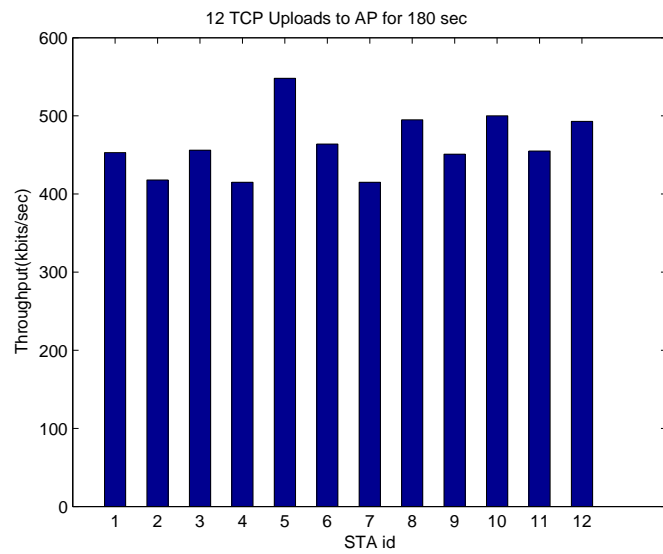


Relative Throughput against CWmin

Measure relative performance of two saturated flows while varying TXOP, AIFS and $CW_{min}$. Compare to well-known models.

# Before



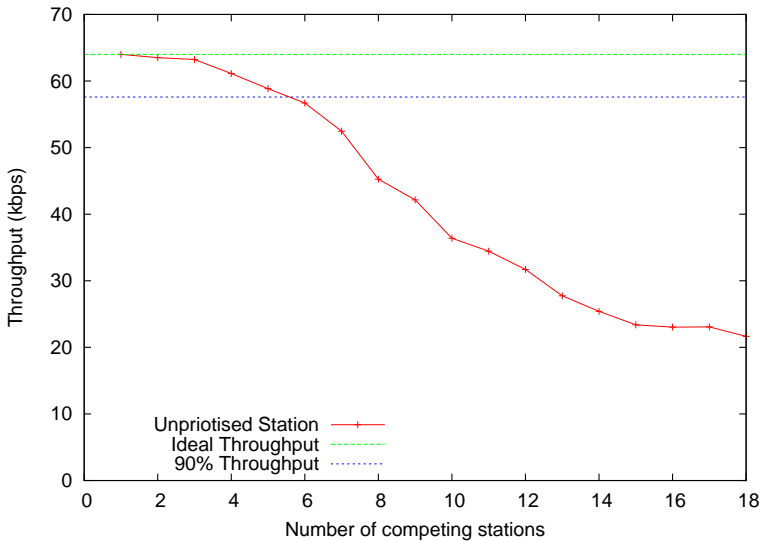12 TCP Uploads to AP for 180 sec

6 Uploads and 6 Downloads for 180 sec

# After



12 TCP Uploads to AP for 180 sec
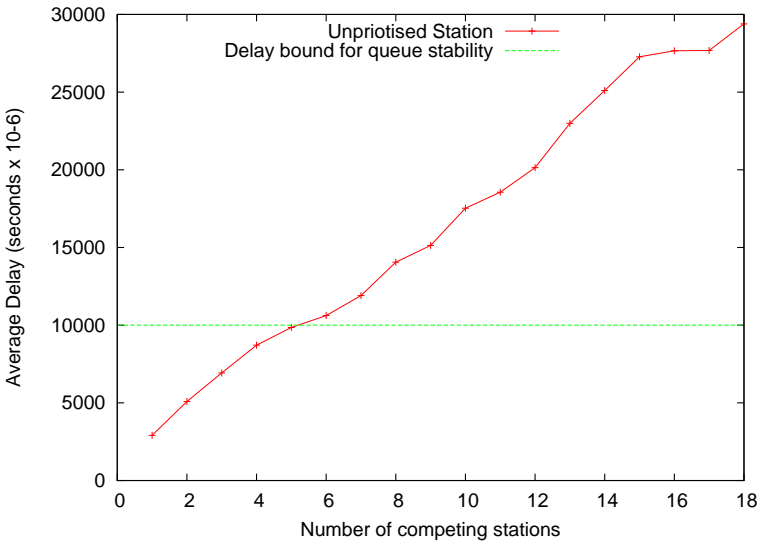
6 Uploads and 6 Downloads

# Unprioritised Voice
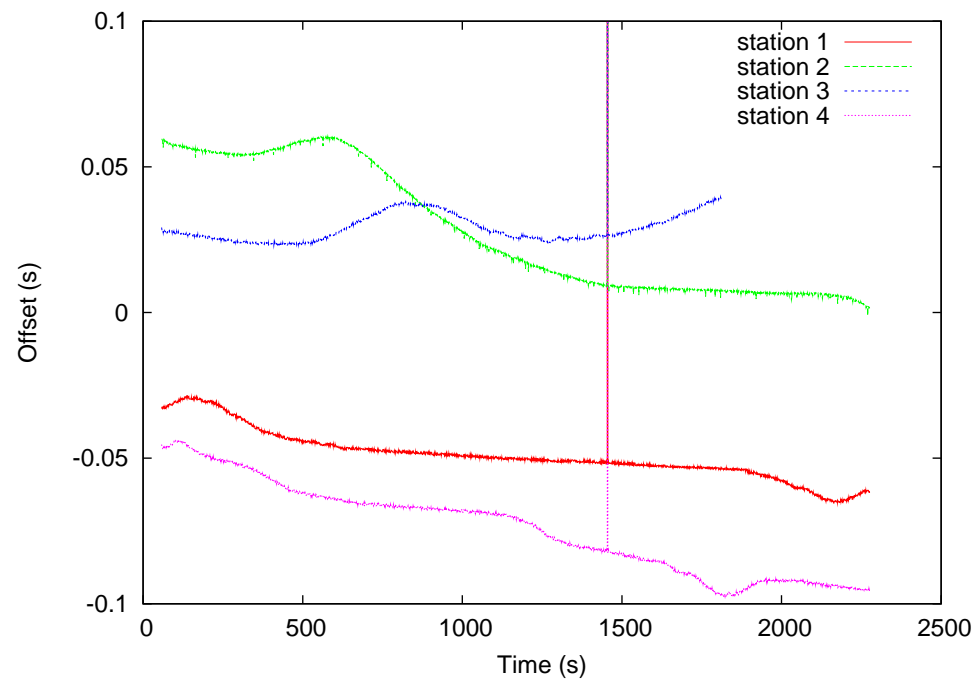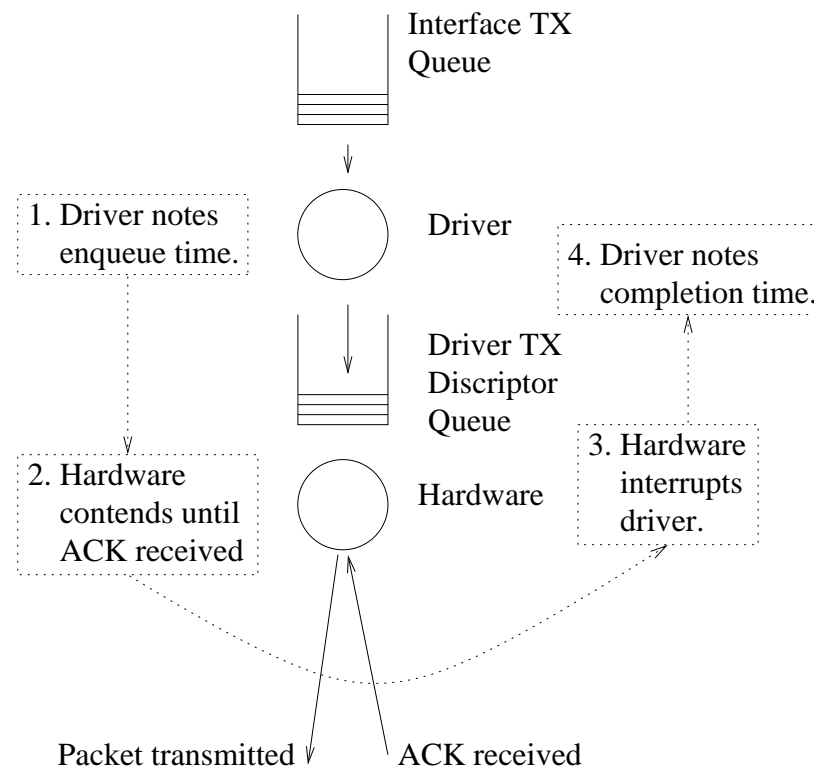
### Throughput



### Delay

# Measuring Delay

- Want to measure one-way MAC delay.

- NTP slow and insufficiently accurate.
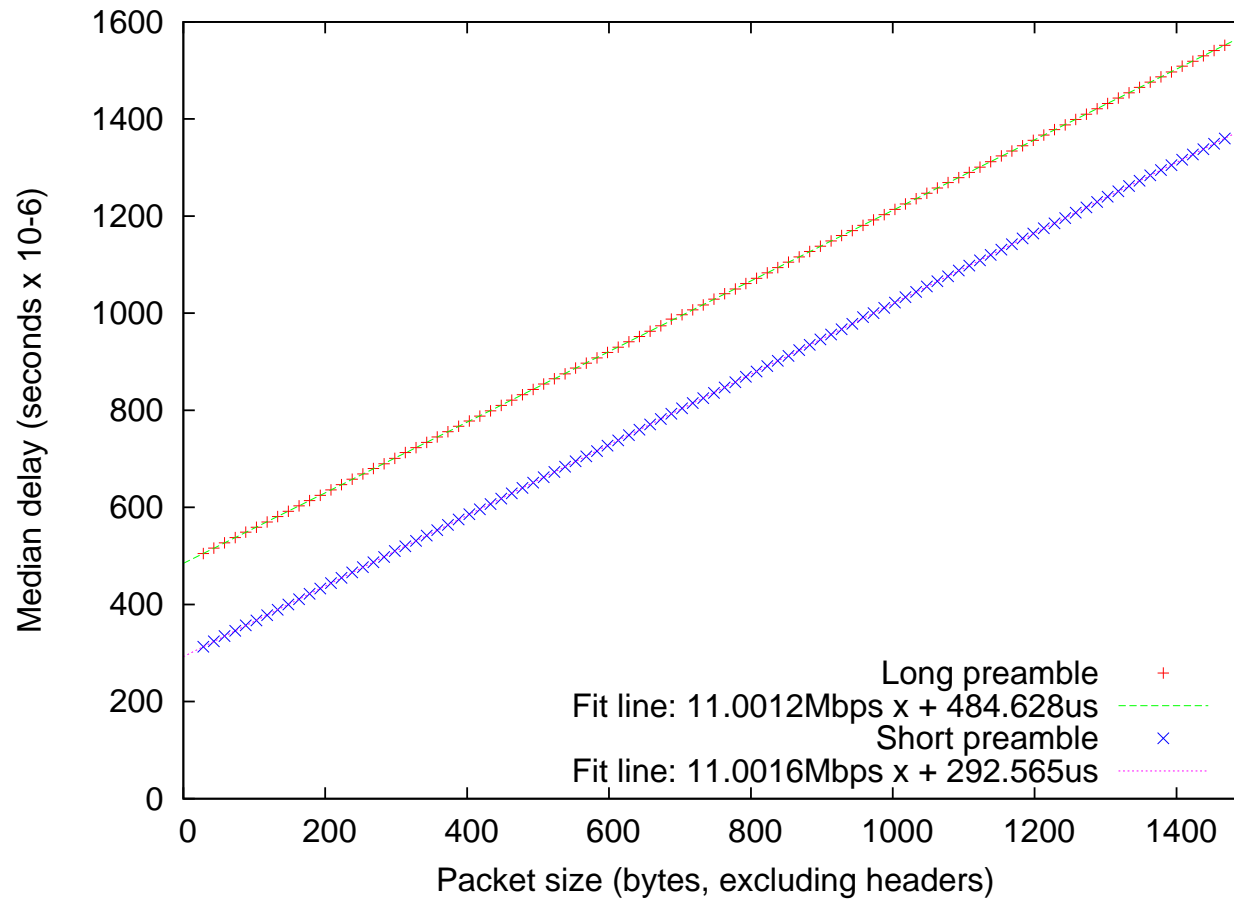
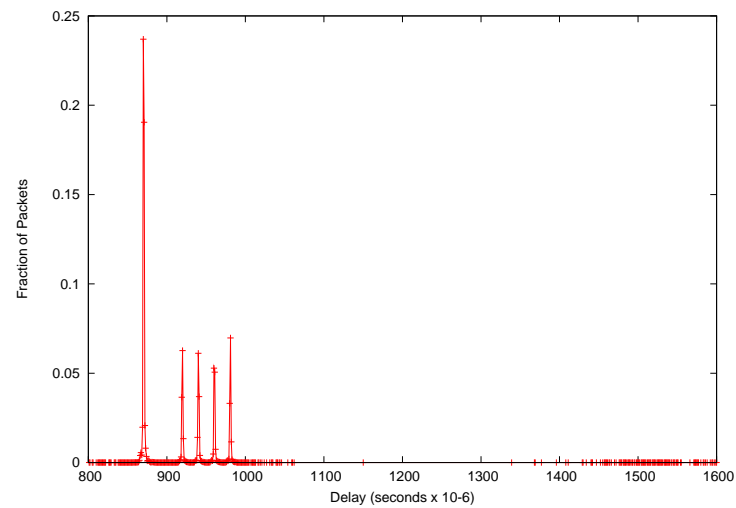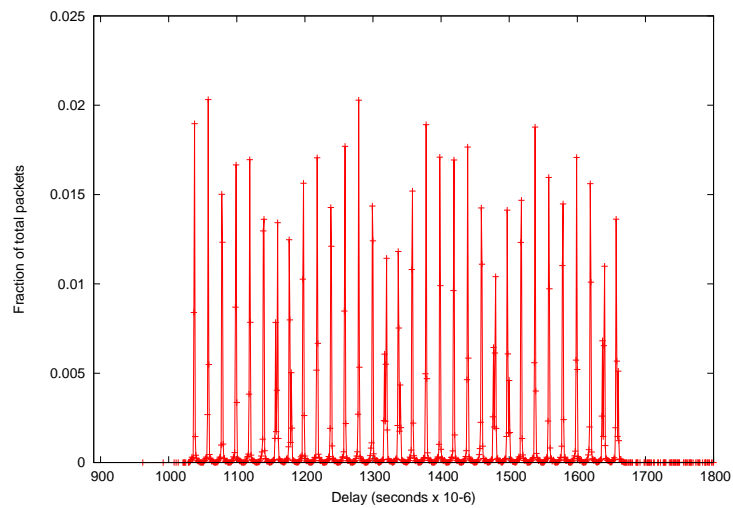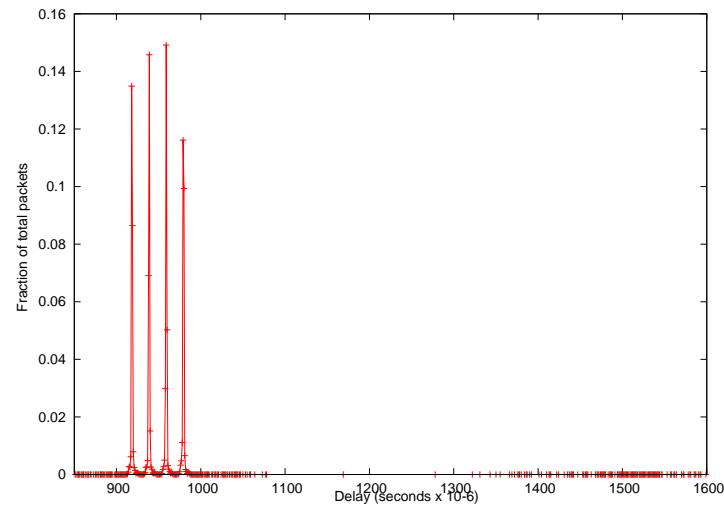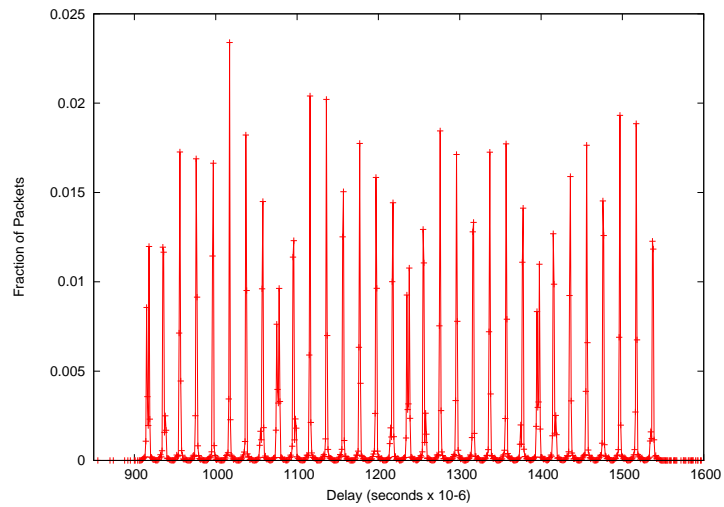- Simultaneously observable TX better, largish noise.

# Delay Technique

- Transmission not complete until MAC ACK.

- Hardware supports interrupt after ACK.

Interface TX
Queue

1. Driver notes
enqueue time.

Driver

4. Driver notes
completion time.

Driver TX
Discriptor
Queue

2. Hardware
contends until
ACK received

Hardware

3. Hardware
interrupts
driver.

Packet transmitted          ACK received

22

# Validation



*Figure: Median delay (seconds x 10-6) vs Packet size (bytes, excluding headers)*

Legend:
- Long preamble (+)
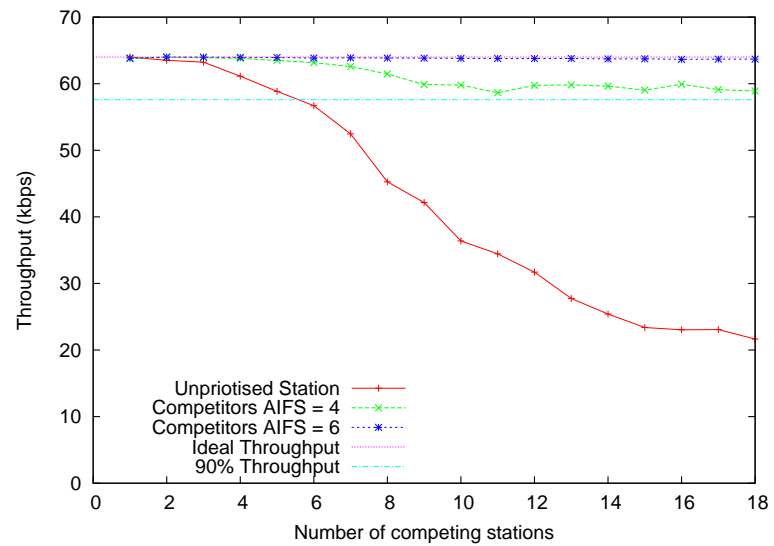- Fit line: 11.0012Mbps x + 484.628us
- Short preamble (x)
- Fit line: 11.0016Mbps x + 292.565us

# AIFS Impact



Throughput



Delay