



Implementation issues for high-speed TCPs

1st October 2003

UCL

Tom Kelly

`ctk21@cam.ac.uk`

Laboratory for Communication Engineering

University of Cambridge

Motivation for new high-speed TCPs

- ⑥ Some Internet users (mainly scientific) want to perform very large bulk transfers
- ⑥ TCP congestion control performs poorly at high-speeds in wide area networks
- ⑥ Even “turning off” congestion control functions unreliably on some implementations

Problem area 1: algorithm

- ⑥ Poor performance of TCP in high bandwidth wide area networks due to TCP congestion control algorithm

Throughput	Window	Loss recovery time	Supporting loss rate
10Mbps	170pkts	17s	5.4×10^{-5}
100Mbps	1700pkts	2mins 50s	5.4×10^{-7}
1Gbps	17000pkts	28mins	5.4×10^{-9}
10Gbps	170000pkts	4hrs 43mins	5.4×10^{-11}

Characteristics of a 200ms, MTU 1500 bytes TCP connection

Problem area 2: OS implementation



- ⑥ What is causing the Linux TCP stack to become unpredictable with large windows?
 - △ Problem with PAWs implementation?
 - △ Hardware/software driver issues?
 - △ SACK implementation problems?

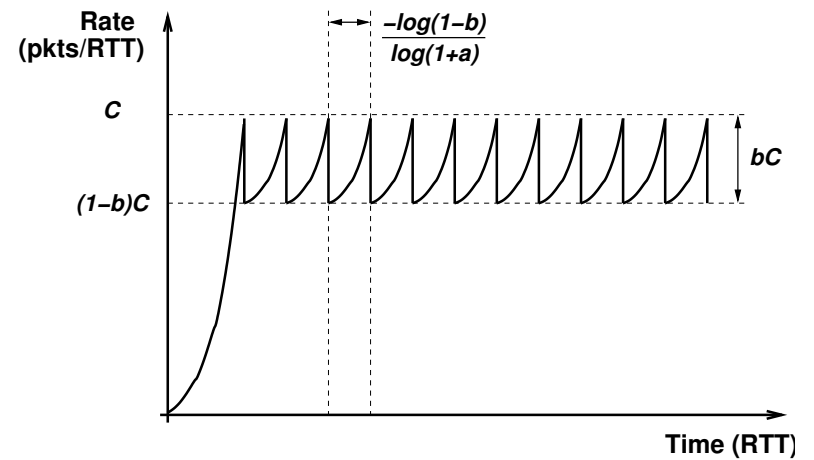
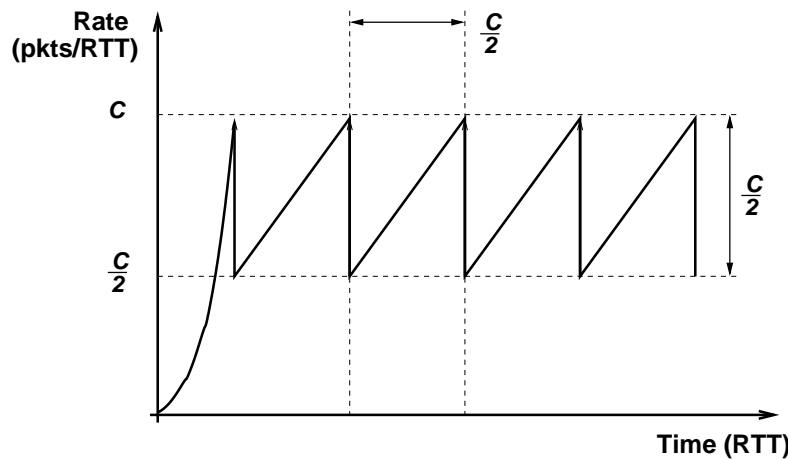
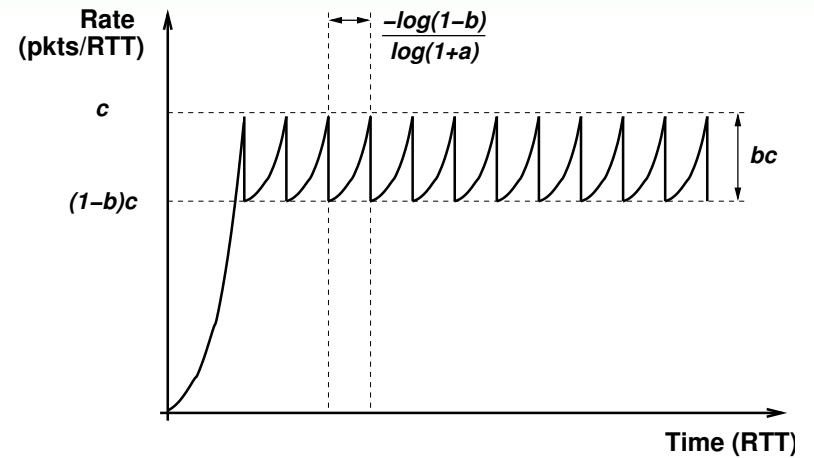
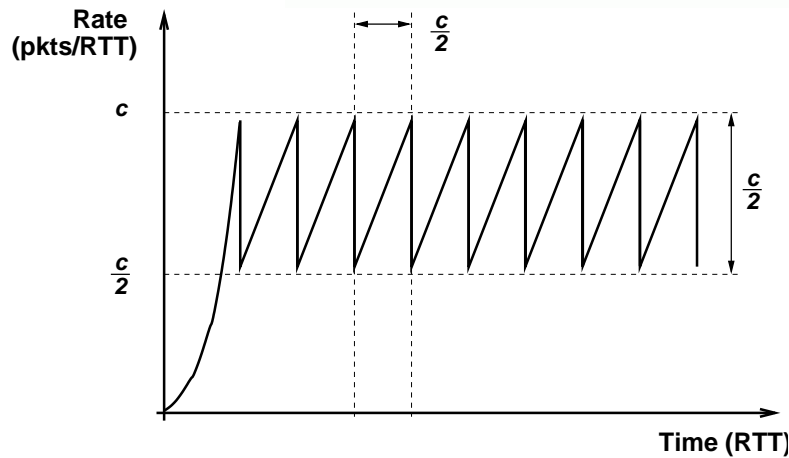
Changing the algorithm - aims and assumptions

- ⑥ Make effective use of high bandwidth links
- ⑥ Changes need to be robust in a wide variety of networks and traffic conditions
 - △ L2 switches, bugs, packet corruption, reordering and jitter
- ⑥ Do not adversely damage existing network traffic
- ⑥ Do not require manual tuning to achieve reasonable performance
 - △ 80% of maximal performance for 95% of the people for the foreseeable future

The generalised Scalable TCP algorithm

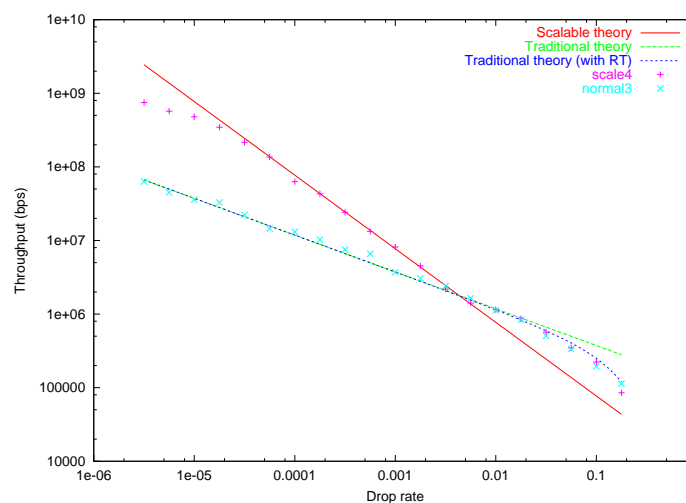
- ⑥ Let a and b be constants
 - △ for each ack in a RTT without loss:
 $wnd_r \leftarrow wnd_r + a$
 - △ for each window experiencing loss:
 $wnd_r \leftarrow wnd_r - b \times wnd_r$
- ⑥ Loss recovery times for RTT 200ms and MTU 1500 bytes
 - △ Scalable TCP: $\frac{\log(1-b)}{\log(1+a)}$ RTTs
e.g. if $a = 0.01, b = 0.125$ then it is about 2.7s
 - △ Traditional: at 50Mbps about 1min 38s, at 500Mbps about 27min 47s!

The Scalable TCP algorithm



Fairness

- ⑥ Choose a legacy window size, $lwnd$
- ⑥ When $cwnd > lwnd$ use the Scalable TCP algorithm
- ⑥ When $cwnd \leq lwnd$ use traditional TCP algorithm



- ⑥ Same argument used in the HighSpeed TCP proposal
- ⑥ Fixing $lwnd$, fixes the ratio $\frac{a}{b}$

Variance and Convergence

- ⑥ increasing b , more variable flows but faster backoff
- ⑥ increasing a , instability but more aggressive ramp up
- ⑥ $lwnd = 16$, $a = 0.01$, and $b = 0.125$ represents a good trade off of concerns

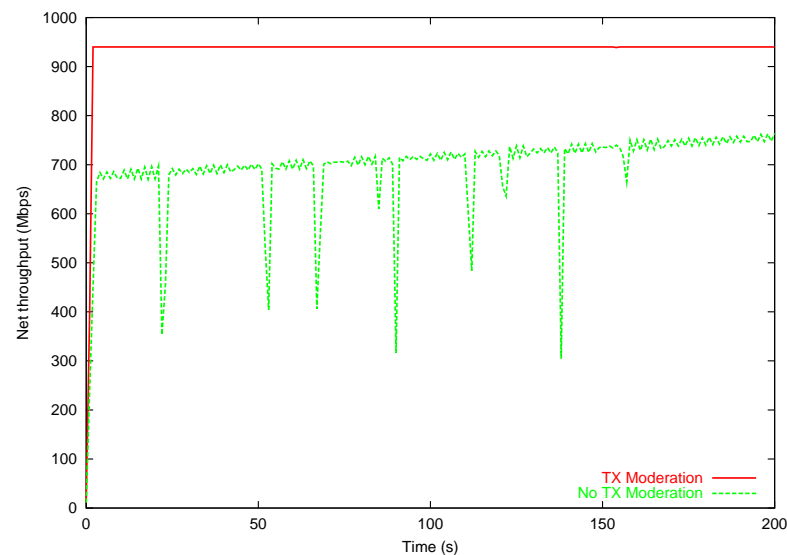
b	a	Rate CoV	Loss recovery time	Rate halving time	Rate doubling time
$\frac{1}{2}$	$\frac{1}{25}$	0.50	$17.7T_r$ (3.54s)	T_r (0.20s)	$17.7T_r$ (3.54s)
$\frac{1}{4}$	$\frac{1}{50}$	0.35	$14.5T_r$ (2.91s)	$2.41T_r$ (0.48s)	$35T_r$ (7.00s)
$\frac{1}{8}$	$\frac{1}{100}$	0.25	$13.4T_r$ (2.68s)	$5.19T_r$ (1.04s)	$69.7T_r$ (13.9s)
$\frac{1}{16}$	$\frac{1}{200}$	0.18	$12.9T_r$ (2.59s)	$10.7T_r$ (2.15s)	$139T_r$ (27.8s)

Linux TCP implementation

- ⑥ Includes all standard high-speed TCP extensions; PAWs, timestamps, SACK
- ⑥ Includes some experimental non-standard features:
 - △ Forward acknowledgement (FACK) to capture flight size during recovery
 - △ Rate-halving; send packet every other acknowledgement during recovery
 - △ Aggressive RTO checking on sent segments when receiving duplicate acknowledgements
 - △ Mechanisms for undoing congestion window decreases if thought to be due to bogus loss detection

Impact of driver TX interrupts

- ⑥ Default Linux SysKonnnect does no transmit interrupt moderation
- ⑥ By altering the driver TX interrupts can be moderated



Linux NAPI driver model

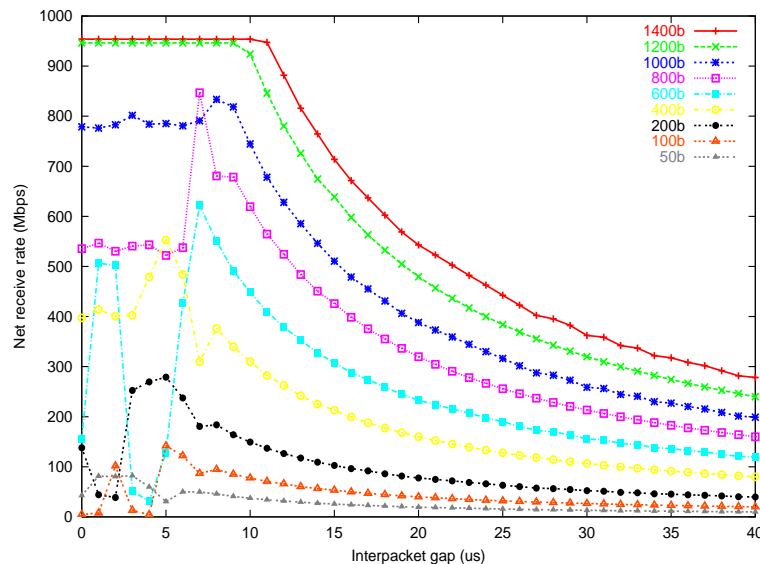
- ⑥ Around for some time in 2.5.x and incorporated in 2.4.20
- ⑥ On receiving a packet, NIC raises interrupt
- ⑥ Driver switches off RX interrupts and schedules RX DMA ring poll
- ⑥ Frames are pulled off DMA ring and is processed up to application
- ⑥ When all frames are processed RX interrupts are re-enabled
- ⑥ Dramatic reduction in RX interrupts under load

Experimental SysKonnnect NAPI driver implemented

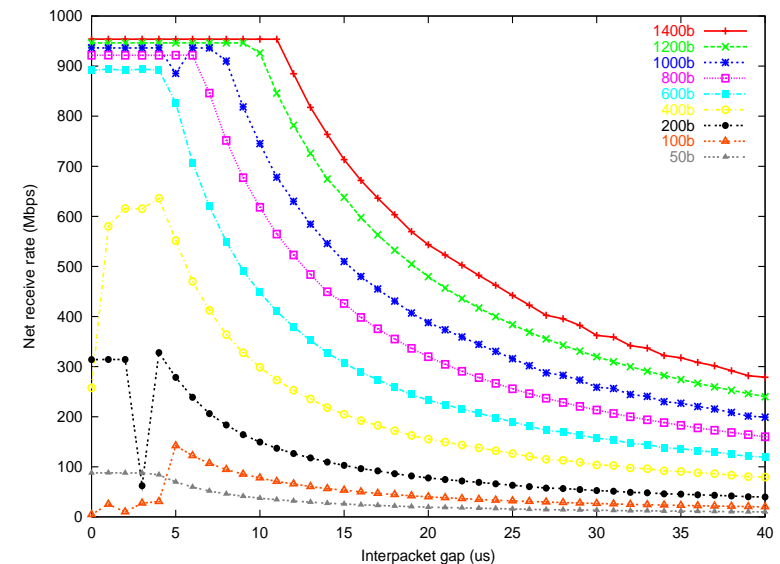
- ⑥ No spec sheet for PCI card ASIC since SysKonnnect was bought by Marvell
- ⑥ Still some RX flagged interrupts appearing; appears benign but makes me suspect there is a bug somewhere
- ⑥ Bottom line is improved performance under heavy load

NAPI receiver results

- ⑥ 2.4Ghz machines connected through router with 2.4.20 sender using TX inter rupt moderation



2.4.19 non-NAPI receiver



2.4.20 NAPI receiver

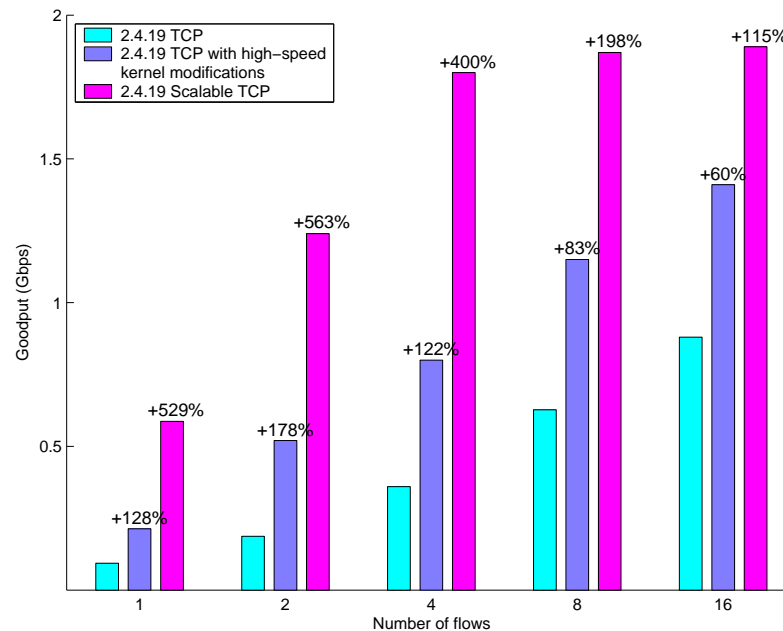
- ⑥ Better throughput for NAPI receiver under load
- ⑥ Some strange behavior with 100b and 50b packets...

Scaling loss detection and recovery algorithms

- ⑥ SACK block processing and segment retransmission both involve trawling the send queue
- ⑥ Trawling the send queue can be $O(cwnd)$ for each acknowledgement
 - △ The queues are there to avoid copying packets
- ⑥ A fix (hack) is to exploit likely fastpath
 - △ Packets delivery in order
 - △ SACK blocks in acks only change in first block
 - △ Cache pointers and assume incremental changes each ack

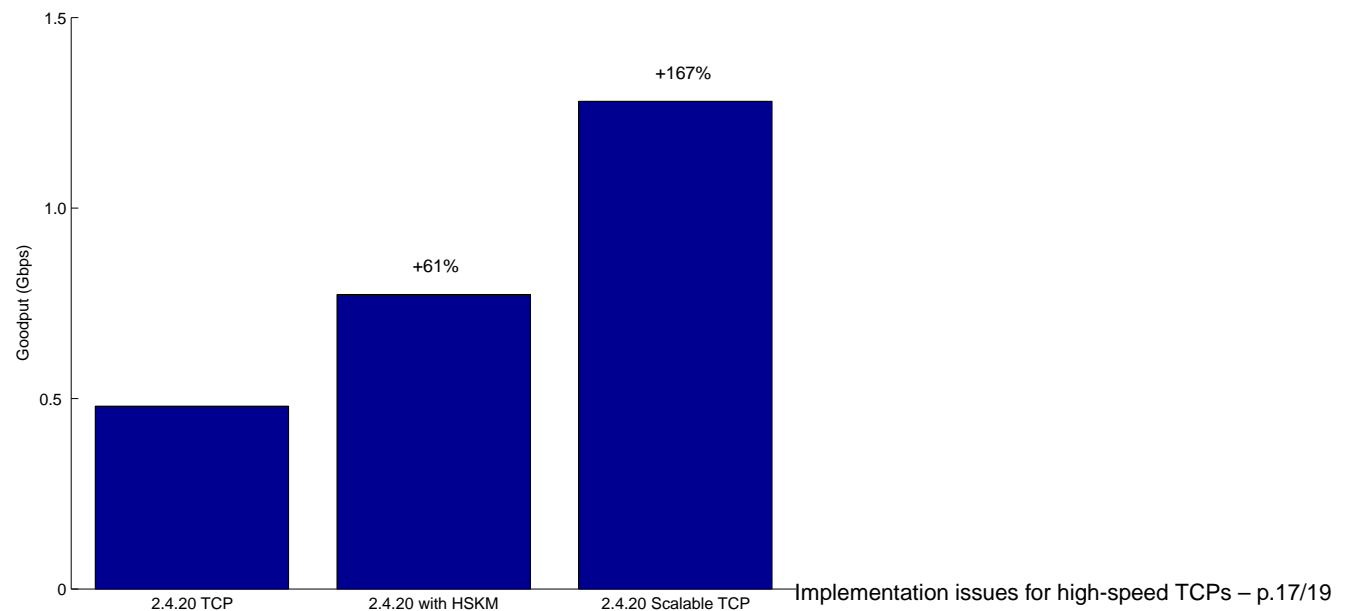
Bulk throughput

- ⑥ DataTAG 2.4Gbps link and minimal buffers (2048/40)
- ⑥ Flows transfer 2 gigabytes and start again for 1200s



Web traffic results

- ⑥ DataTAG 2.4Gbps link and minimal buffers (2048/40)
 - △ 4 bulk concurrent flows across 2 machines for 1200s
 - △ 4200 concurrent web users across 3 machines
- ⑥ No change in web traffic with and without bulk transfers in all scenarios



Problems

⑥ Result set is small!

- △ Difficult to conduct controlled implementation experiments

⑥ Linux TCP implementation a mess for high-speed

- △ Should split data segments from packet headers and protocol state (e.g. OpenBSD)
- △ Need scatter-gather I/O to do this with minimal copies

⑥ Scalable TCP

- △ Synchronisation problems due to design, worse than TCP but simulations don't match reality
- △ Which workloads and topologies m?

Conclusion

- ⑥ Linux implementation can be greatly improved for high-speed operation
- ⑥ Scalable TCP an easy evolution from the traditional TCP AMID scheme which can improve performance
- ⑥ Much more to be done deciding between schemes; HSTCP, Vegas/FAST, Westwood, etc.
- ⑥ Freely available working code
<http://www-lce.eng.cam.ac.uk/~ctk21/scalable>