

TCP Forensics

Identifying and debugging issues with
TCP in FreeBSD

Lawrence Stewart

lastewart@swin.edu.au

Centre for Advanced Internet Architectures (CAIA)
Swinburne University of Technology





- 1 Who is this guy?
- 2 TCP Congestion Control
- 3 Tools
- 4 Case studies
- 5 Wrapping Up

Detailed outline (section 1 of 5)



1 Who is this guy?

1 Who is this guy?

2 TCP Congestion Control

3 Tools

4 Case studies

5 Wrapping Up

Who is this guy (and who let him past security)?



- BEng (Telecomms and Internet Technologies) 1st class honours / BSci (Comp Sci and Software Eng) (2001-2006)
- Centre for Advanced Internet Architectures, Swinburne University (2003-2007)
 - Research assistant/engineer during/after studies
 - <http://caia.swin.edu.au/>
- Currently a PhD candidate in telecomms eng at CAIA
 - Strong focus on transport protocols
 - <http://caia.swin.edu.au/cv/lstewart/>
- Addicted to FreeBSD since 2003
 - Experimental research, software development, home networking, servers and personal desktops

Detailed outline (section 2 of 5)



1 Who is this guy?

2 TCP Congestion Control

3 Tools

4 Case studies

5 Wrapping Up

2 TCP Congestion Control

- TCP jargon
- Brief recap
- Where are we today
- Open issues



cwnd congestion window

MSS maximum segment size

ssthresh slow start threshold

ACK TCP acknowledgment

RTT round trip time

BDP bandwidth-delay product

RFC request for comment

CC congestion control

tcpcb TCP control block

RTO Retransmit timeout



- “Congestion Avoidance and Control” (ACM SIGCOMM-88)
 - Seminal paper on TCP CC
- RFC 1122 (1989)
 - Required a TCP to use Van Jacobsen’s CC algorithms
- RFC 2582/3782 (1999/2004)
 - NewReno TCP
- RFC 4614 (2006)
 - A good summary of TCP related RFCs

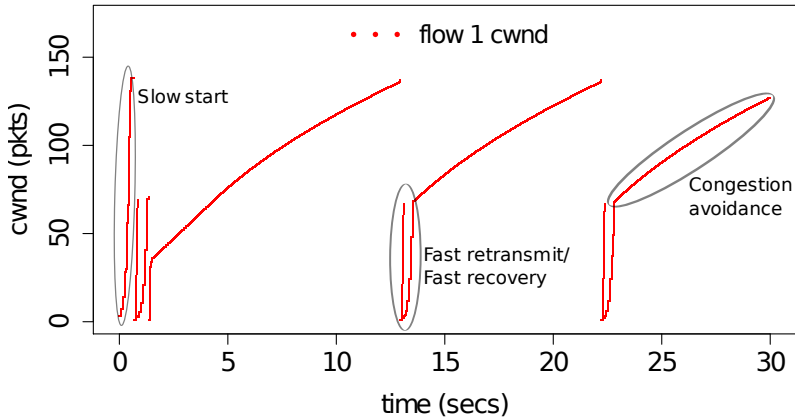


- Core CC algorithms ¹
 - Slow start
 - Congestion avoidance
 - Fast retransmit
 - Fast recovery
- Tweaks and additions along the way
 - SACK, ABC, ECN, window scaling, timestamps, etc.

¹See RFC2001



Vanilla FreeBSD 7.0 – 80 RTT, 10Mbps



Where are we today



- NewReno is defacto standard with warts e.g.
 - Fat pipes
 - Wireless
- State of the CC union
 - BSD uses NewReno
 - Linux uses CUBIC
 - Windows Vista uses CTCP
- TCP/IP stack enhancements
 - CSO/TSO/LRO/TOE
 - Socket buffer autotuning
- Are we going to hell in a handbasket?



- High-speed CC algorithms ²
 - FAST, HS-TCP, H-TCP, CTCP, CUBIC, etc.
- Delay based CC algorithms
- How do we compare and evaluate CC algorithms?
- CSO/TSO/LRO/TOE obscure behaviours
- Testing/verification of TCP/IP stack behaviour

²Nice summary:

<http://kb.pert.geant2.net/PERTKB/TcpHighSpeedVariants>

Detailed outline (section 3 of 5)



1 Who is this guy?

2 TCP Congestion Control

3 Tools

4 Case studies

5 Wrapping Up

3 Tools

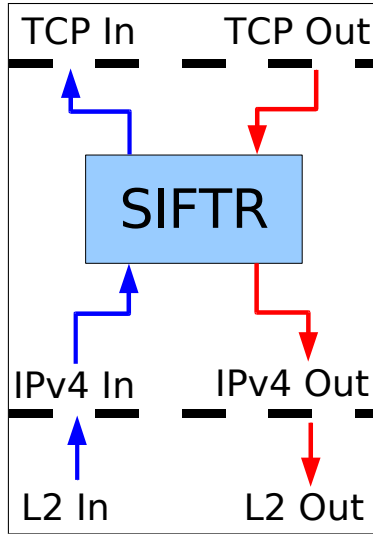
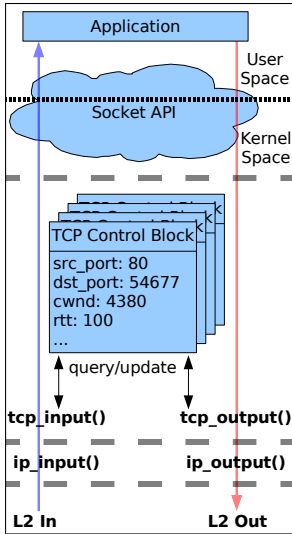
- SIFTR
- SIFTR demo
- Dummynet
- Others

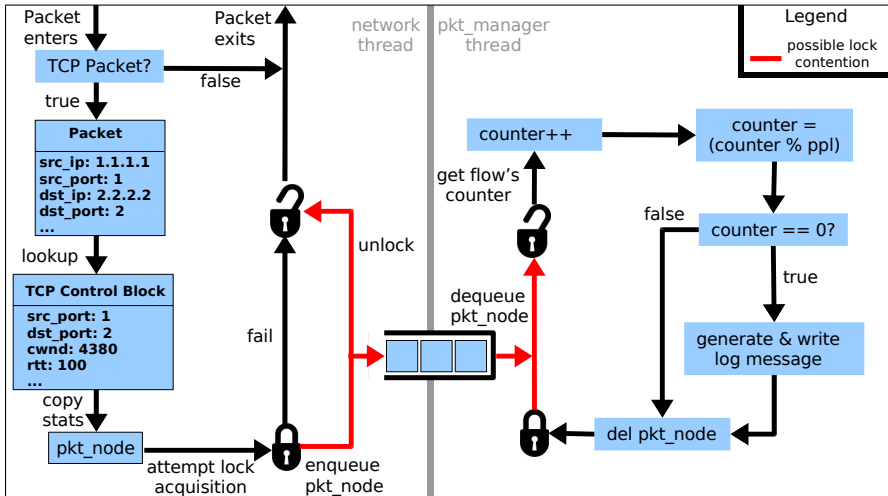


- Statistical Information For TCP Research
- FreeBSD [5,6,7,8] kernel module
- BSD licenced source ³
- Similar base concept to Web100
- Event triggered (not poll based)
- Currently logs 25 different variables to file as CSV data ⁴

³Available from: <http://caia.swin.edu.au/urp/newtcp/tools.html>

⁴See README in SIFTR distribution for specific details







Let's see what we can see!



- Repeatable tests makes things easier
- Patched dummynet to:
 - Include a new pipe option “pls” for predicatable packet loss
 - Provide detailed per-packet logging functionality
- New variables added to the `net.inet.ip.dummynet sysctl` tree:
 - `log_enable`: Enable/disable detailed logging
 - `dn_logfile`: Set path to logfile
 - `max_queue_size_pkts`: Set the maximum queue size for packet-based pipes
 - `max_queue_size_bytes`: Set the maximum queue size for byte-based pipes



- I've personally found these invaluable:
 - Iperf
 - tcpstat
 - R
 - ns-2

Detailed outline (section 4 of 5)



1 Who is this guy?

2 TCP Congestion Control

3 Tools

4 Case studies

5 Wrapping Up

4 Case studies

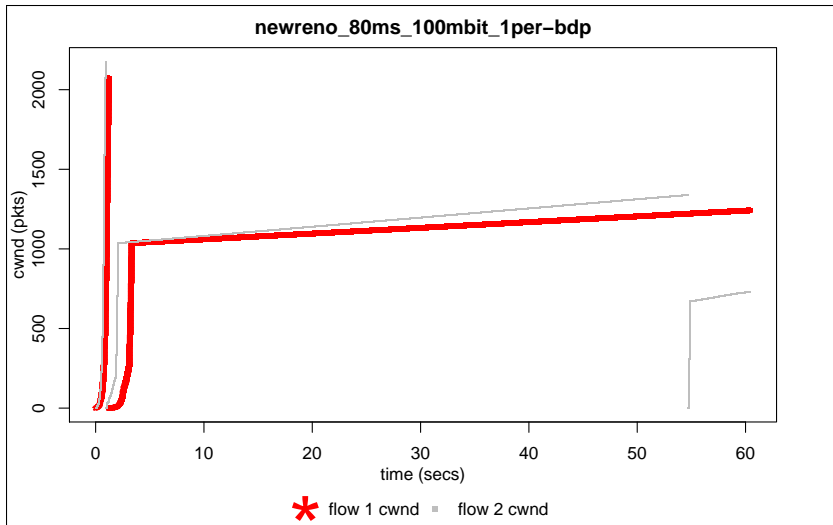
- The case for ABC
- Skeletons in the (New?)Reno closet
- Holey SACK

The case for ABC

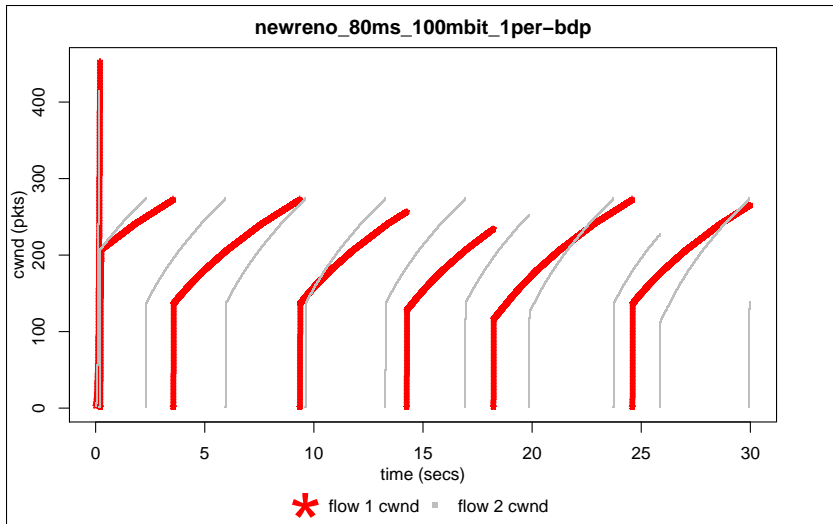


- Appropriate Byte Counting (RFC 3465)
- Change cwnd increase to be based on # ACK'ed bytes
- Removes error introduced by current cwnd increase approximation:
 - For each ack, $cwnd_+ = \alpha / cwnd$, where $\alpha = 1$ for (New)Reno

The case for ABC



The case for ABC

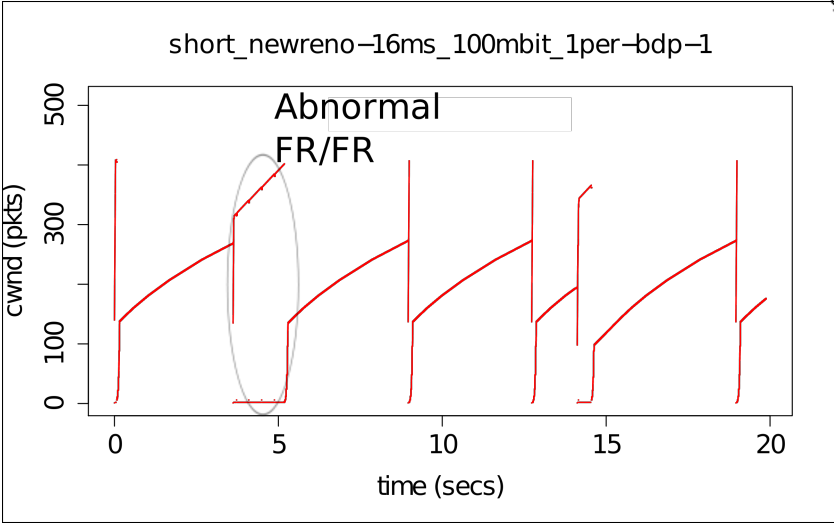


Skeletons in the (New?)Reno closet



- Bug suspected to be hiding in `tcp_output()`
- Characterised by sender going quiet for an RTT during FR/FR
- Causes very laboured FR/FR which can significantly impact performance

Skeletons in the (New?)Reno closet

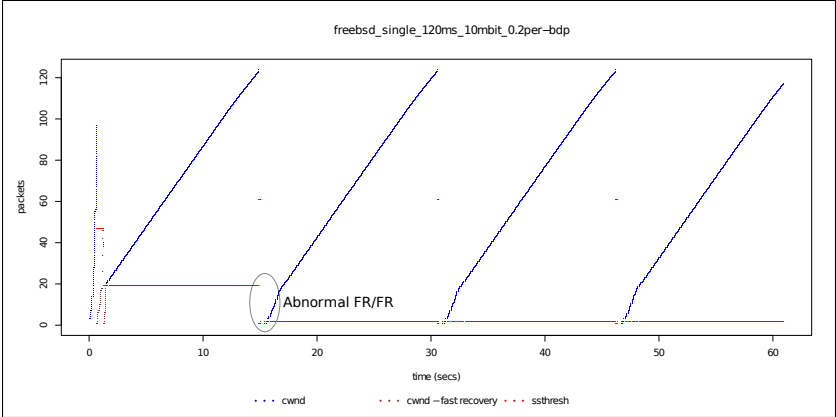




- Bug in SACK implementation ⁵
- Characterised by emission of 0 byte data packet by sender
- Causes an unnecessary RTO
- Triggered by packet loss to a receiver with small reassembly queue when cwnd is large
- May possibly be resolved by fixing the more general (New?)Reno bug

⁵Package of data demonstrating problem: [http:](http://caia.swin.edu.au/urp/newtcp/freebsd_sack_issue-1.0.tar.gz)

[//caia.swin.edu.au/urp/newtcp/freebsd_sack_issue-1.0.tar.gz](http://caia.swin.edu.au/urp/newtcp/freebsd_sack_issue-1.0.tar.gz)



Detailed outline (section 5 of 5)



1 Who is this guy?

2 TCP Congestion Control

3 Tools

4 Case studies

5 Wrapping Up

5 Wrapping Up

- Ideas for future work
- Acknowledgements
- Questions



- TCP specific:
 - RTT estimator
 - Reassembly queue autotuning (Andre)
 - Share CC between TCP/SCTP (Randall et. al.)
 - Implement more CC algorithms
 - Comprehensive RFC compliance check
- TCP/IP stack in general:
 - Framework for dealing with CSO/TSO/LRO/TOE
 - Testing framework



- The FreeBSD Foundation



- Dan Langille et. al.
- FreeBSD community

- Cisco Systems





Questions?