# TCP Enhancements

## Kimmo Raatikainen

**kimmo.raatikainen@cs.helsinki.fi**

---

# Lession Outline

- TCP options recommended for 2.5g3g
- Linux TCP implementation

# Recommendations

- Appropriate Window Size (Sender & Receiver)
  - Bandwidth Delay Product (BDP) of the end-to-end path
  - the window scale option can be used to overcome the 64 kB limitation.
- Increased Initial Window (Sender)
  - the initial CWND (congestion window):
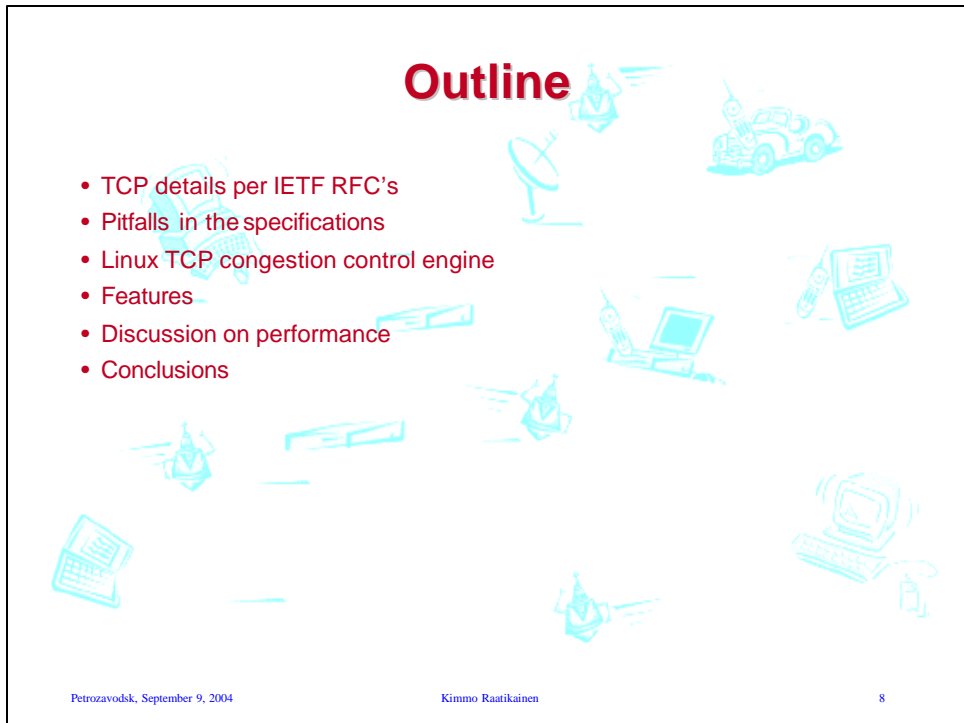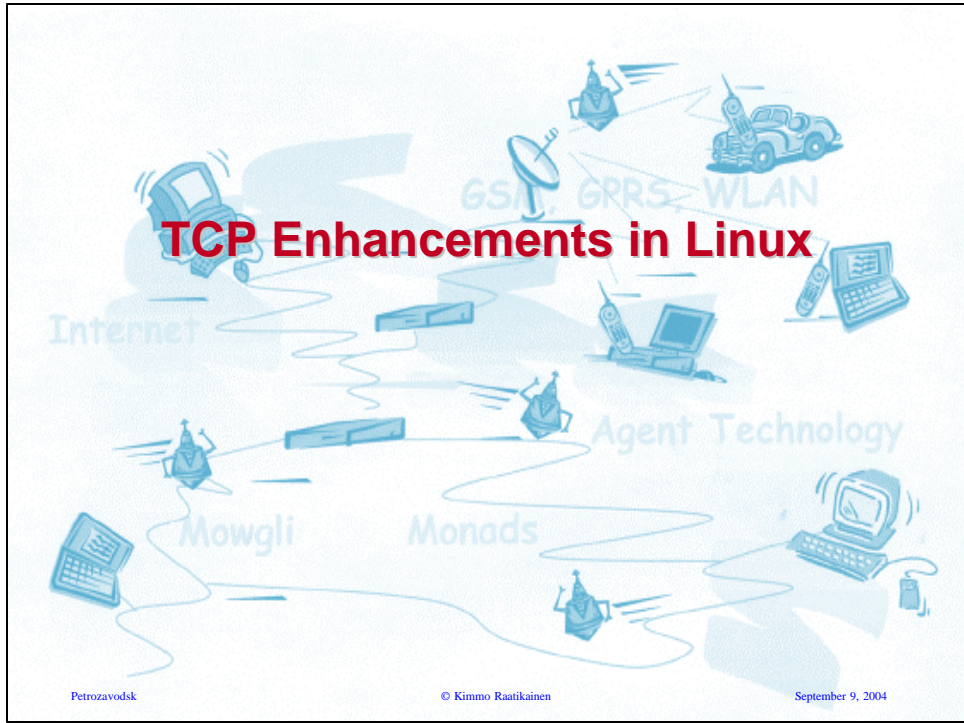  - min (4*MSS, max (2*MSS, 4380 bytes))

# Recommendations

- Limited Transmit (Sender)
  - RFC3042, Limited Transmit, extends Fast Retransmit/Fast Recovery for TCP connections with small congestion windows that are not likely to generate the three duplicate acknowledgements required to trigger Fast Retransmit.
  - TCP over 2.5G/3G implementations SHOULD implement Limited Transmit
- IP MTU Larger than Default
- Path MTU Discovery (Sender & Intermediate Routers)

# Recommendations

- Selective Acknowledgments (Sender & Receiver)
    - TCP over 2.5G/3G SHOULD support SACK.
    - In the absence of SACK feature, the TCP should use NewReno RFC2582
- Explicit Congestion Notification (Sender, Receiver & Intermediate Routers)
    - TCP over 2.5G/3G SHOULD support ECN.

# Recommendations

- TCP Timestamps Option (Sender & Receiver)
    - TCP SHOULD use the TCP Timestamps option
- Disabling RFC1144 TCP/IP Header Compression (Wireless Host)

# TCP Enhancements in Linux

# Outline

- TCP details per IETF RFC's
- Pitfalls in the specifications
- Linux TCP congestion control engine
- Features
- Discussion on performance
- Conclusions

# TCP Basics

- Slow start, congestion avoidance
- Receiver generates duplicate ACKs when data is missing
- Fast retransmit at third duplicate ACK
- Fast recovery to keep the "ACK clock" in pace
  - Standard Reno (RFC 2581) or NewReno (RFC 2582)
- Without SACK at most one retransmission in RTT
- Retransmission Timer adjusted smoothly based on measured round-trip times
  - SRTT + 4 * RTTVAR

# Some TCP Enhancements

- SACK: allow several retransmissions in RTT
  - acknowledge separate blocks of received data
  - conservative: "holes" are still outstanding
  - Forward ACKs (FACK): "holes" are considered lost
- D-SACK: report duplicate segments using SACK
- Timestamps: measure RTT for retransmissions
- Eifel: report unnecesary retransmissions using timestamps
- ECN: Explicit Congestion Notification
- Limited transmit: Avoid timeouts with small window

# Discussion on Specifications

- RFC 2581 & RFC 2582: Congestion Control
  - Cwnd is artificially increased on duplicate ACKs. It does not correspond to real number of segments allowed to be in flight

- SACK congestion control draft

  > in flight = SND.NXT – SND.UNA

  - Separate document that assumes SACK is in use
  - Cwnd is not artificially increased
  - We need to implement both? Nah...

- RFC 2988 does not work well with high-

  > in flight = SND.NXT – SND.UNA – SACKed

  - No one sees this, because RTTs are generally below 1000ms

---

# RFC 2988: RTO Calculation

> $RTTVAR \leftarrow \frac{3}{4} * RTTVAR + \frac{1}{4} * | SRTT - MRTT |$
>
> $SRTT \leftarrow 7/8 * SRTT + 1/8 * MRTT$
>
> $RTO \leftarrow \max(1000ms, SRTT + 4 * RTTVAR)$

- RTO estimator decays rapidly
- When measured RTT drops, RTO goes up
- No one cares, because
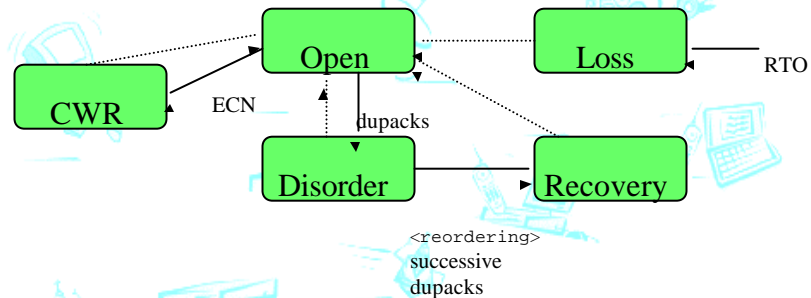  - Min limit of 1000ms
  - Coarse-grain timers

# Linux Approach

$$\text{in flight} = \text{packets\_out} - \text{sacked\_out} - \text{lost\_out} + \text{retrans\_out}$$

- Common congestion control with Reno, SACK, FACK
- *sacked_out:* # of segments surely left network
  - SACK: number of SACKed segments
  - Reno: number of duplicate ACKs
- *lost_out*: # of segments suspected lost
  - SACK & Reno: first unacknowledged is considered lost
  - FACK: holes between SACKs are considered lost
- scoreboard markings are updated accordingly
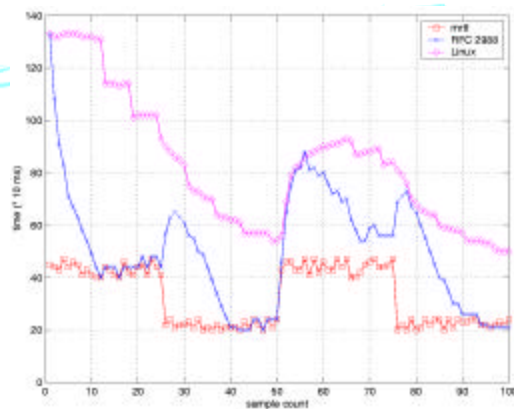
---

# CA States



- **<reordering> is adjusted when unnecessary retransmission is detected**
  - **by default 3**
- **Window is increased in *Open* and *Loss* states**
- **Window is decreased in *CWR* and *Recovery* states**

# Features

- Implements Explicit Congestion Notification (ECN)
- Congestion window is decreased steadily every second ACK in *CWR* and *Recovery* states
  - as in "rate-halving"
- Disorder state implements "Limited transmit" in practice
- Congestion window validation: If congestion window is not fully used for a while, it is reduced
- Congestion control state is cached for future connections

# Linux Retransmission Timer

- Based on RFC 2988
- min. RTO = 200 ms
- min. RTTVAR = 50 ms
- RTTVAR reduced once per round-trip time
  - but increased instantly
- if RTT drops significantly, RTTVAR weight is reduced to 1/32
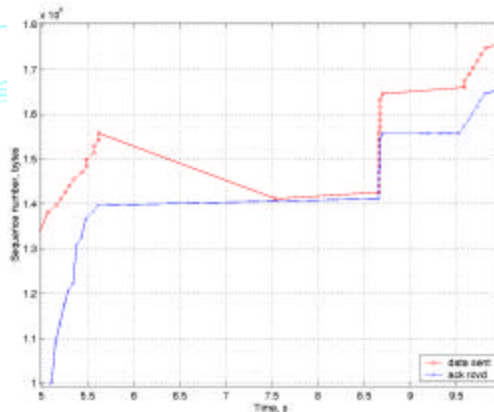
8

# Congestion Window Undoing

- TCP sender can make false retransmits, e.g. due to
  - false RTOs caused by unexpected delay
  - dupacks caused by reordering in network
- False retransmits can be detected by using
  - TCP timestamps: receiver echoes timestamp of original segment after retransmission
  - D-SACKs: a retransmitted segment is acknowledged in cumulative ACK and in D-SACK
- After detecting false retransmission the sender sets
  - cwnd <- max(cwnd, ssthresh * 2)
  - ssthresh <- prior_ssthresh

# Undoing on TCP Timestamps
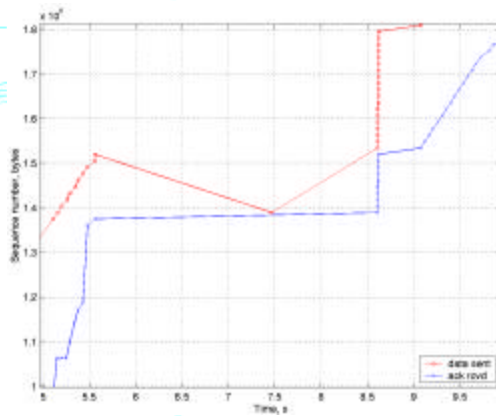
**Without timestamps**

- A 3-second excessive delay occurs on 256Kbps link
- Triggers RTO, but ACKs for original segments arrive after RTO
- congestion window is halved
- 65 KB acknowledged between 5 and 10 s.
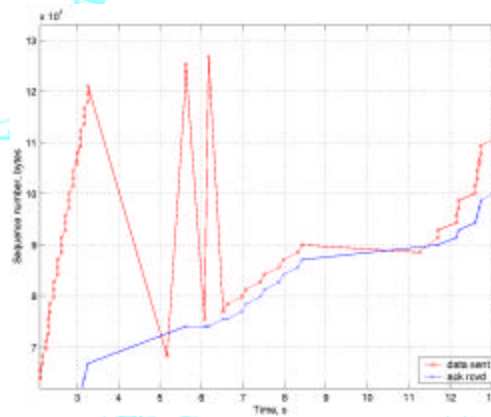
9

# Undoing on TCP Timestamps

**With timestamps**

- Next ACK after RTO echoes timestamp of original segment
- Spurious timeout is detected
  - continue by transmitting new data
  - revert recent changes on congestion control parameters
- 75 KB acknowledged between 5 and 10 s.

# Undoing Can Fail

- Link outage: One window of data segmenents and ACKs are dropped
- ACKs echo latest timestamp that updated window
- Because ACKs are lost, sender thinks new ACK acknowledged earlier data
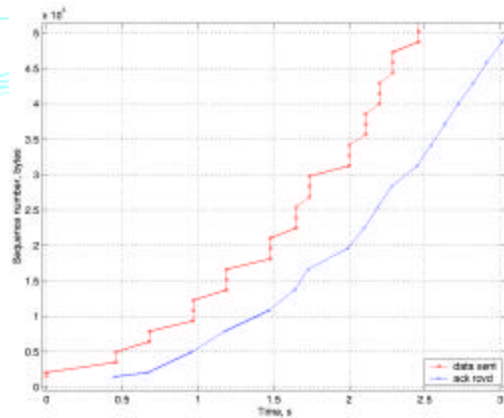  - Declares RTO spurious

# Delayed Acknowledgements

- Delayed acknowledgements should be used by TCP receiver
- Linux receiver measures interarrival times and adjusts delay timer accordingly
  - goal is to get an ACK out for every second segment
- Quick acknowledgements can be used at the beginning of the connection
  - causes the sender to increase the window faster
  - No more than (advwin / 2) quick acknowledgements are allowed to avoid silly windows

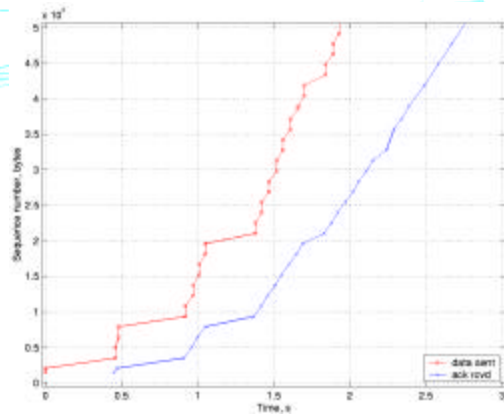# Effect of Quick Acks

**Without quickacks**

- 256 Kbps , 200 ms delay => BW*delay more than 12 KB
- 4-5 round-trips until the link is fully utilized
- every second segment is acknowledged
- 50 KB transmitted in 2.5 seconds

# Effect of Quick Acks

### With quickacks

- For the first 32 KB every segment is acknowledged
- 50 KB transmitted in 2 seconds

---

# Concluding Remarks

- Implementation follows packet conservation in practice
  - congestion window always holds a valid value
  - counters try to estimate how many packets really are outstanding
- If the data structures tracking outstanding packets and suspected losses are detected incorrect, undoing takes place
- Retransmission timer tries to avoid the pitfalls of the original algorithm