

RACK: a new fast recovery

Yuchung Cheng <ycheng@google.com>

Fast recovery should use time not counters

RFC3517: after `dupthresh` DUPACKs

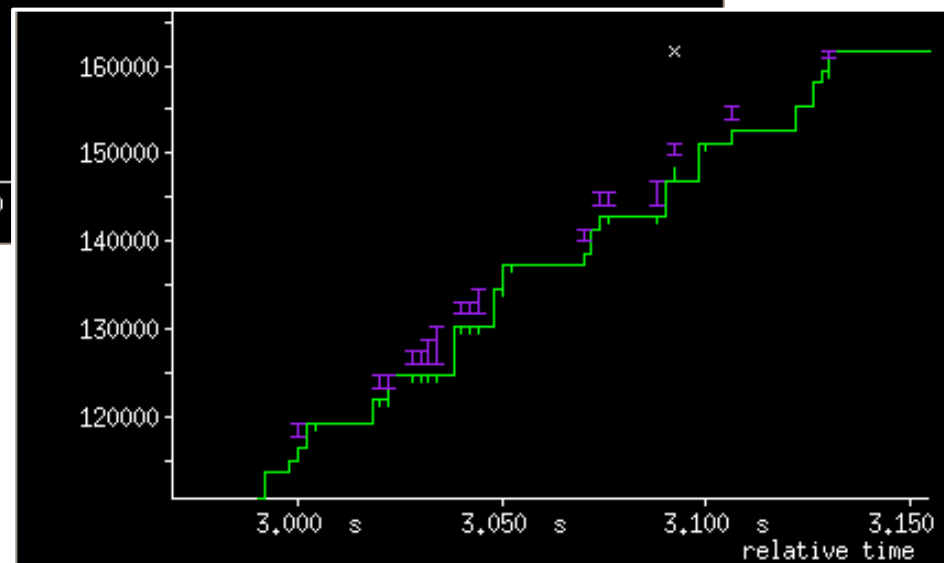
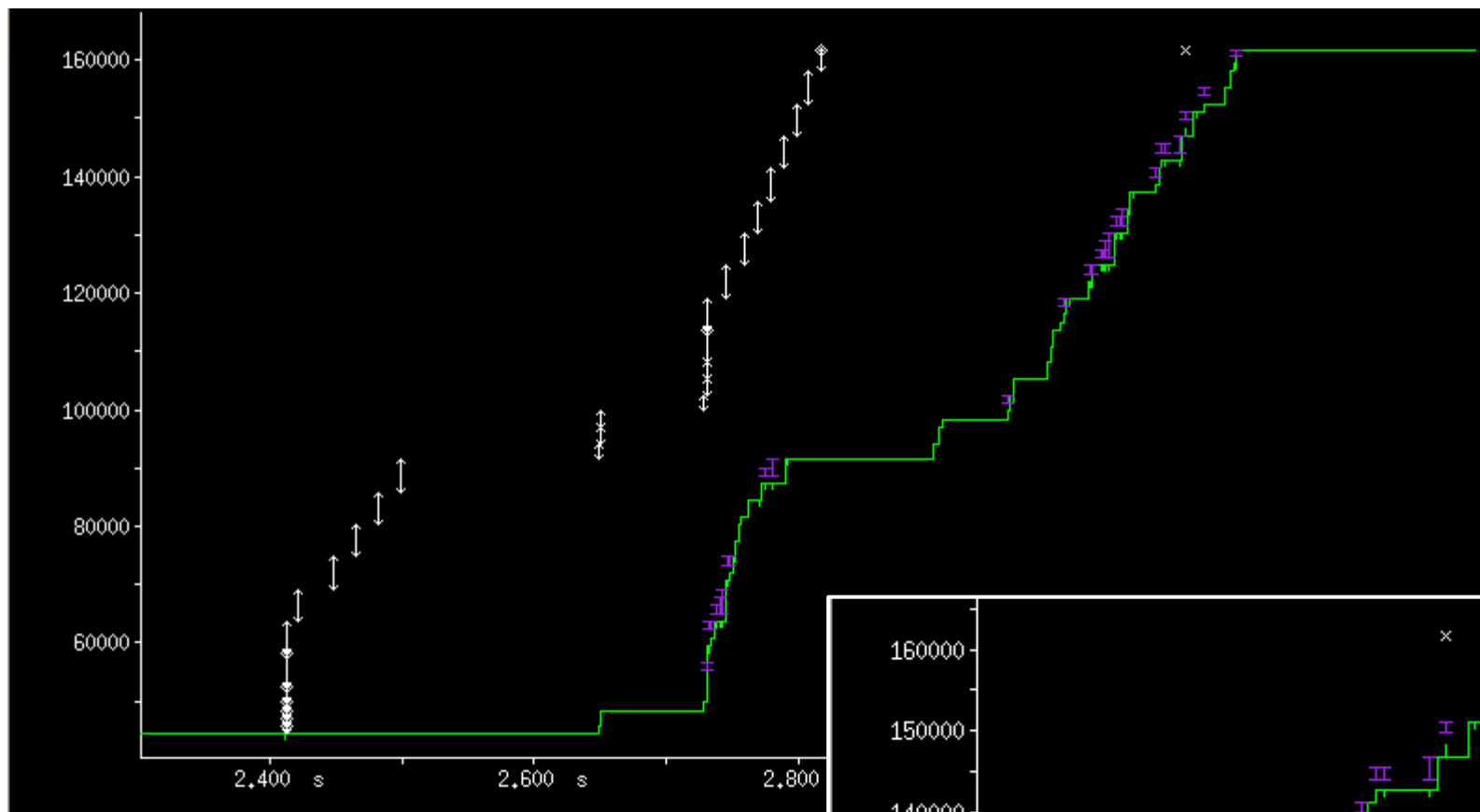
FAACK: after `dupthresh`*MSS highest seq is sacked

RACK: after `reorder_delay` past when some pkt sent later is delivered

Reordering is about delay

1. TE shifts entire flow to a shorter path
2. Network forwards packets on different paths or out-of-order

Low correlation to flight size



Algorithm

`rack_snt`: last xmit time of the most recently (s)acked pkt

On receiving an ACK, for each pkt `p`:

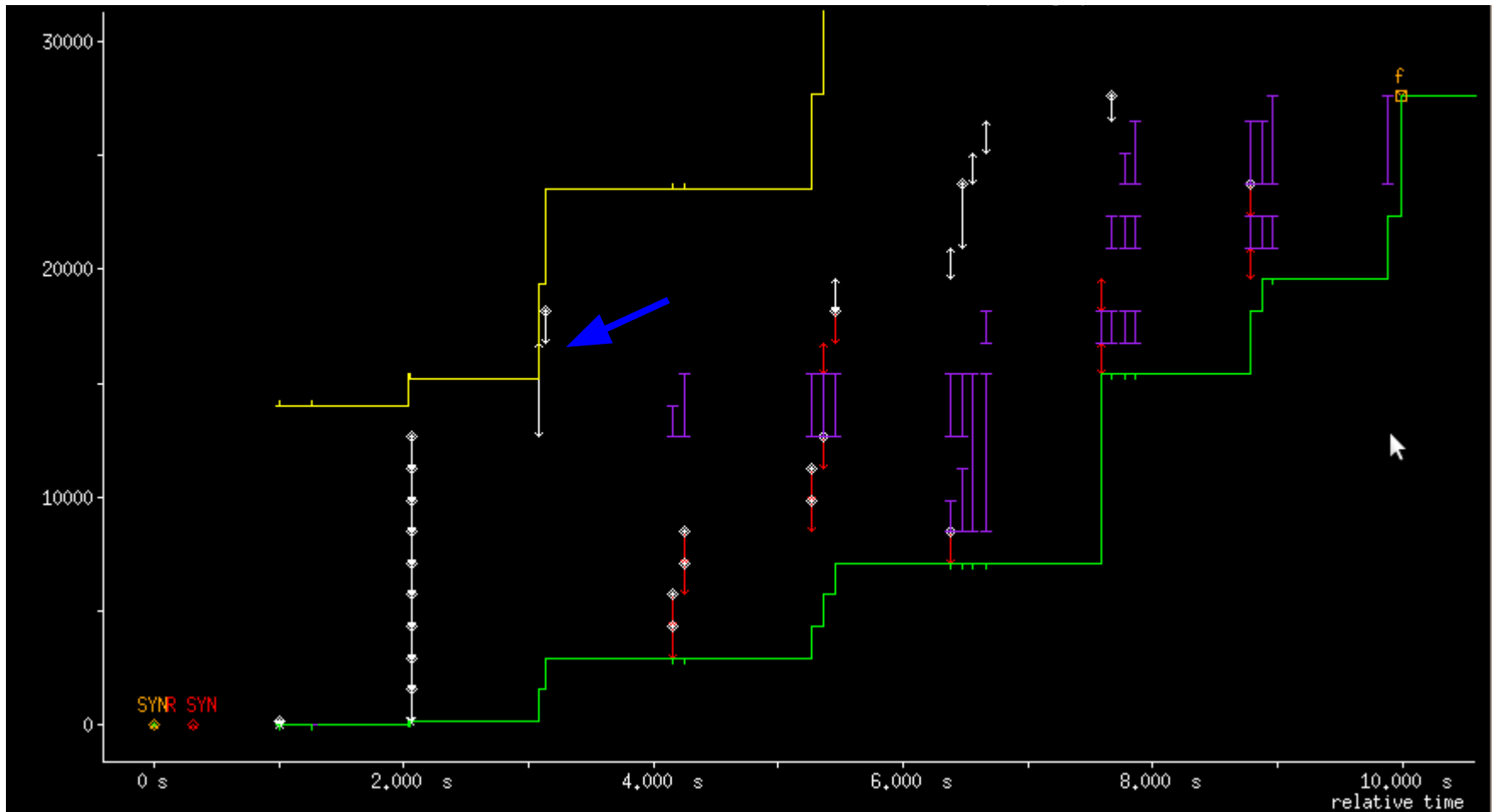
`p_snt`: last xmit time of `p`

Mark `P` lost if $\text{rack_snt} - \text{p_snt} > \text{reo_delay}$

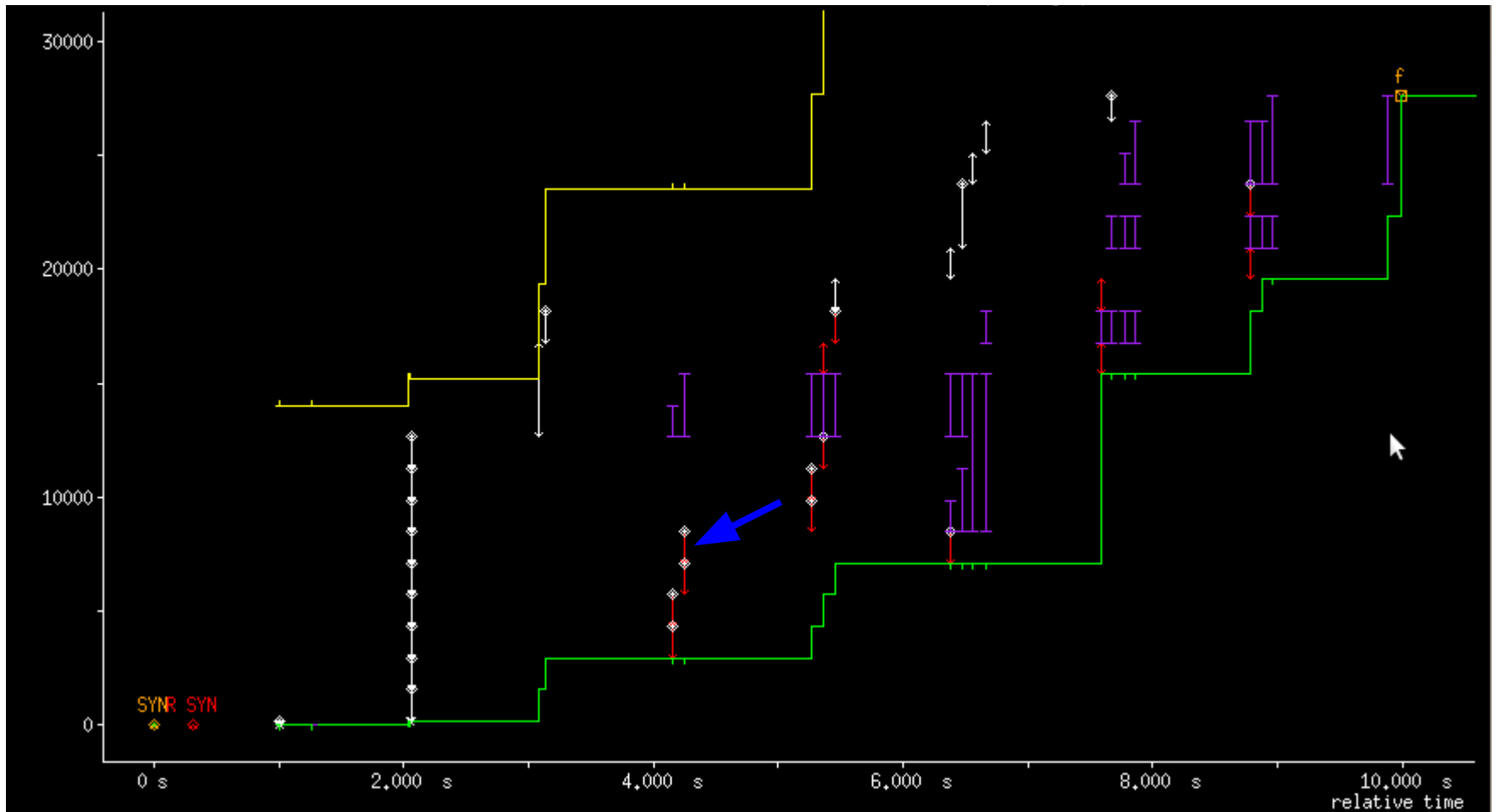
`reo_delay` starts with 0 and dynamically increase based on reordering interval

RACK works for original *and* retransmitted data

Recover tail drops (blue arrow)



Recover lost retransmits (blue arrow)



Replace all counter-based heuristics

Classic/dynamic dupthresh

RFC 3517 - SACK recovery

RFC 5827 - Early retransmit

FAK

Lost retransmit

Thin-stream dupack

Work in progress

Dynamic reo_delay

Integrate TLP draft-dukkipati-tcpm-tcp-loss-probe-01

Linux patches and IETF draft under way