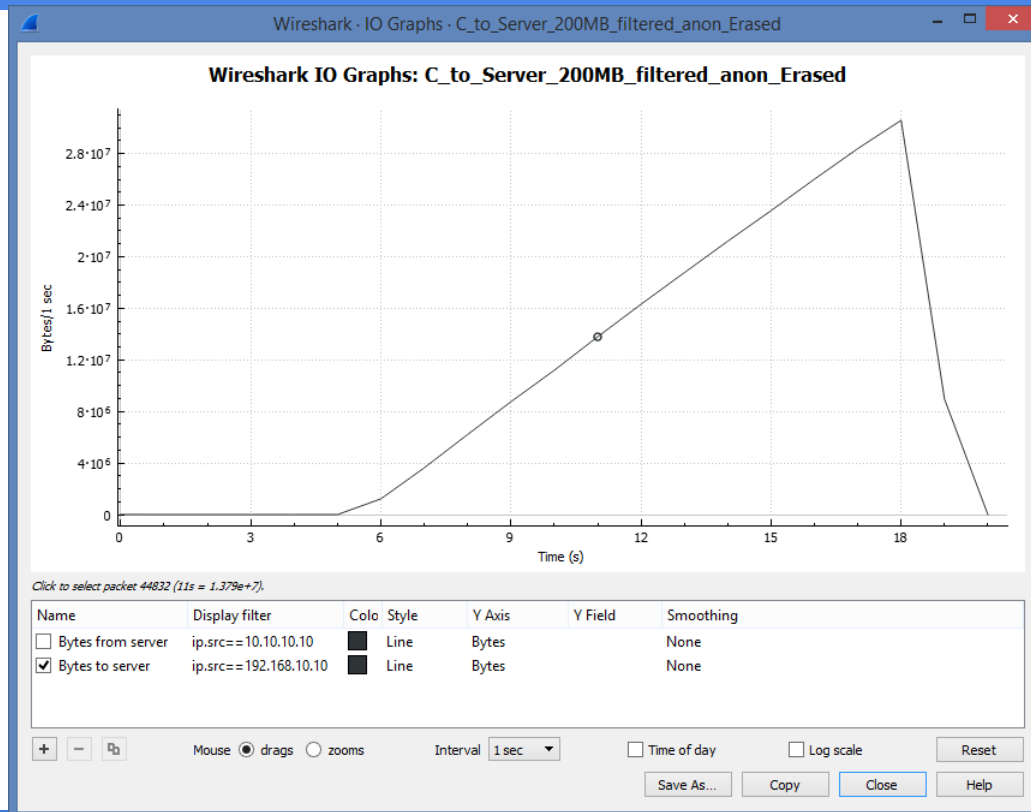# An interesting Graph?
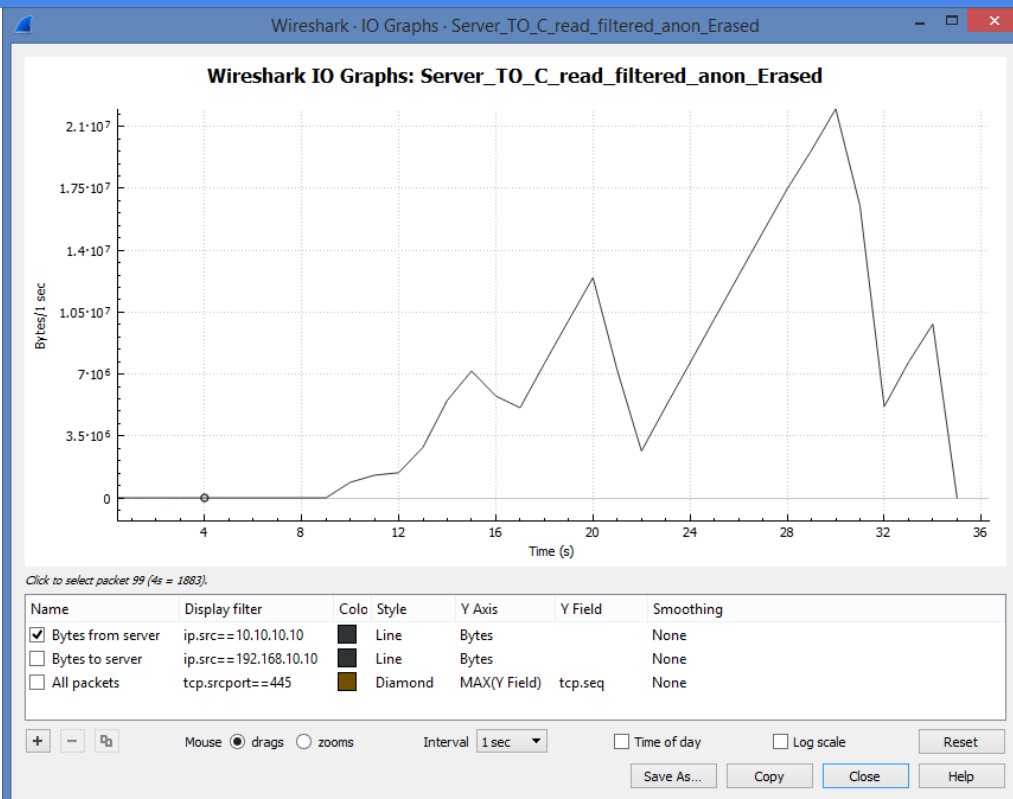
https://osqa-ask.wireshark.org/questions/55972/slow-writes-even-slower-reads-spanning-wan-to-netapp

# Another interesting Graph?

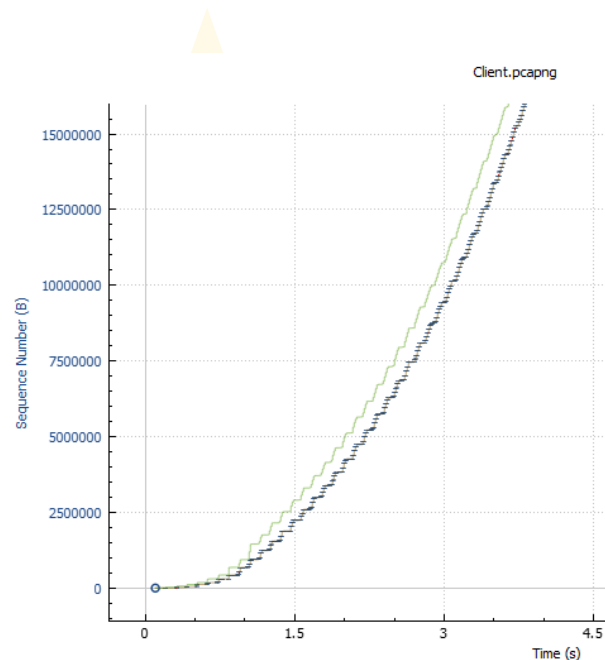https://osqa-ask.wireshark.org/questions/55972/slow-writes-even-slower-reads-spanning-wan-to-netapp

- **<u>This is due to Slow Start!</u>**
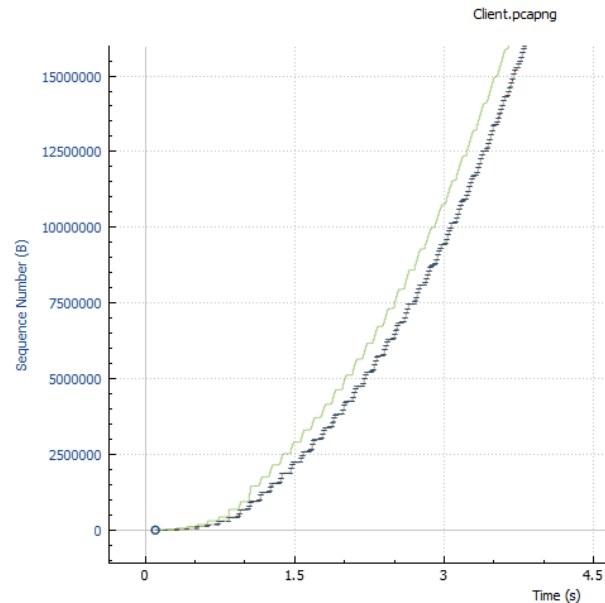

Client.pcapng

- But how does Slow Start work?
- Why is it called Slow Start?
- And why do we need Slow Start?
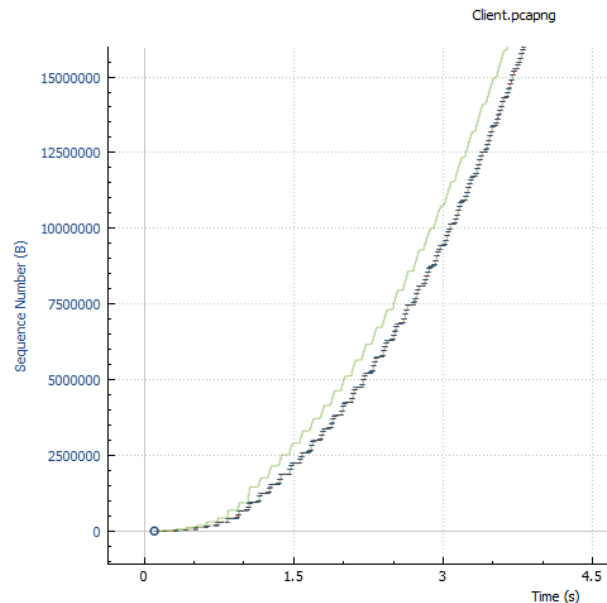
# Slow Start RFC History

- Slow Start is defined in the following RFCs
  - **RFC5681** of the year 2009 obsoletes:
    - RFC2581 of the year 1999 which obsoletes:
      - RFC2001 of the year 1997 author Richard W. Stevens



Client.pcapng

# Slow Start Variables

- We need Slow Start to determine the available bandwidth

- We do this by using the following variables:
  - Receive Window Size
    **rwnd**
  - Congestion Window:
    **cwnd**
  - Slow Start Threshhold:
    **ssthresh**
  - Initial Congestion Window Size
    **iw**
  - Sender Maximum Segment Size
    **smss**
  - Receivers Maximum Segment Size
    **rmss**



Client.pcapng
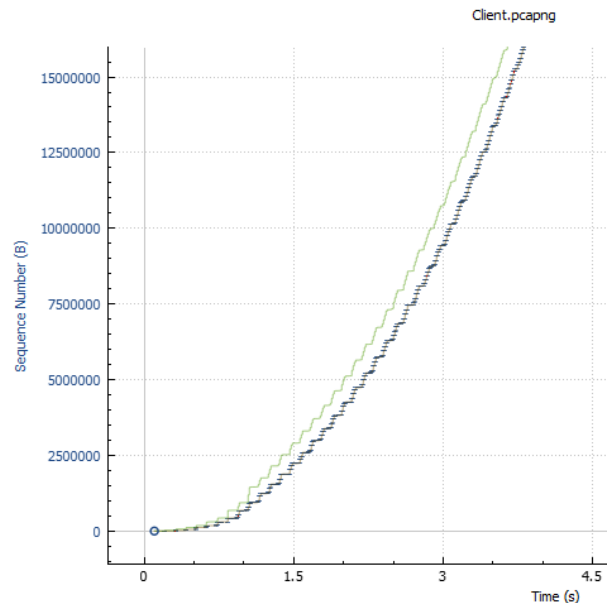
- **rwnd:**
  - The receiver's advertised window (**rwnd**) is a receiver-side limit on the amount of outstanding data.
  - The **rwnd** size can be seen in a trace.
    -> It is represented by the TCP Window Size within the TCP- header

# cwnd = Congestion Window

- **cwnd:**
  - The congestion window (**cwnd**) is a sender-side limit on the amount of data the sender can transmit into the network before receiving an acknowledgment (ACK)
  - The **cwnd** size is a stack variable which **CAN NOT** be **seen** in a trace, it just can be guessed
  - Can be made visible by the command "ss -i" on linux systems

- The **minimum of cwnd and rwnd** governs data transmission.



Client.pcapng

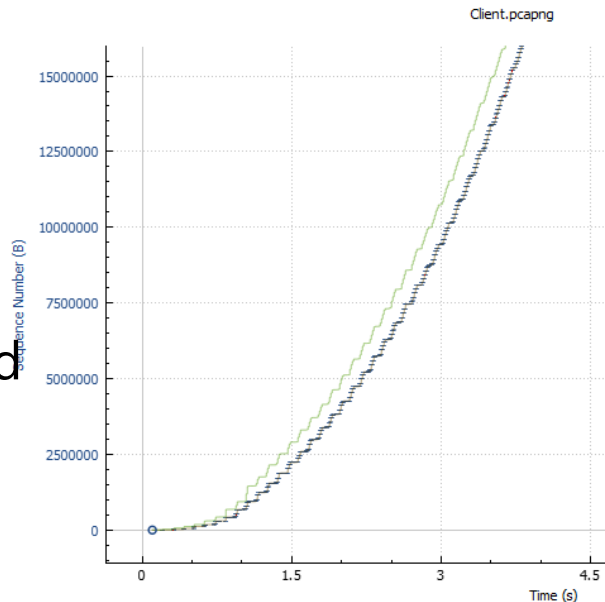- **ssthresh**
  The slow start threshold (**ssthresh**), says if slow start algorithm is used or not.

- The **initial** value of **ssthresh** SHOULD be set **HIGH ENOUGH**
  (e.g., to the size of the largest possible advertised window)



Client.pcapng

# IW = Initial Congestion Window

- **IW** MUST be set using the following guidelines as an upper bound.
    - If **SMSS <= 1095 bytes**:
      **IW** = **4 * SMSS** bytes and **MUST NOT** be more than **4 segments**

    - If (**SMSS > 1095 bytes**) and (**SMSS <= 2190 bytes**):
      **IW** = **3 * SMSS** bytes and **MUST NOT** be more than **3 segments**

    - If **SMSS > 2190 bytes**:
      **IW** = **2 * SMSS** bytes and **MUST NOT** be **more** than **2 segments**

# When do we use Slow Start

- If **cwnd** < **sstresh** then
  **Slow Start** is used

- If **cwnd** > **sstresh** then
  **congestion avoidance** is used



Client.pcapng

**SMSS**      = 1     * 1000 Bytes
**IW**          = 4     * 1000 Bytes
**SSTRESH**   = 64   * 1000 Bytes

CWND = 4 * 1000 Bytes

Sender

Receiver

**SMSS** = 1 * 1000 Bytes
**IW** = 4 * 1000 Bytes
**SSTRESH** = 64 * 1000 Bytes



CWND = 4 * 1000 Bytes

Sender

Data1
Data2
Data3
Data4

Receiver

**SMSS** = 1 * 1000 Bytes
**IW** = 4 * 1000 Bytes
**SSTRESH** = 64 * 1000 Bytes



Sender

CWND = 4 * 1000 Bytes

CWND = 8 * 1000 Bytes

Data1 Data2 Data3 Data4

ACK4+1

Receiver

**SMSS** = 1 * **1000 Bytes**
**IW** = 4 * **1000 Bytes**
**SSTRESH** = 64 * **1000 Bytes**

SMSS = 1 * 1000 Bytes
IW = 4 * 1000 Bytes
SSTRESH = 64 * 1000 Bytes

# Slow Start cgwnd Growth

| Time [cycles] | sstresh [1000 Bytes] | cwnd [1000 Bytes] | SMSS [1000 Bytes] | rwnd [1000 bytes] |
|---|---|---|---|---|
| 1 | 128 | 4 | 1 | 256 |
| 2 | 128 | 8 | 1 | 256 |
| 3 | 128 | 16 | 1 | 256 |
| 4 | 128 | 32 | 1 | 256 |
| 5 | 128 | 64 | 1 | 256 |
| 6 | 128 | 128 | 1 | 256 |
| 7 | 128 | 128 | 1 | 256 |
| 8 | 128 | 128 | 1 | 256 |
| 9 | 128 | 128 | 1 | 256 |



Exponetial Slow Start behaviour

**SMSS** = 1 * 1000 Bytes
**IW** = 4 * 1000 Bytes
**SSTRESH** = 64 * 1000 Bytes

CWND = 4 *
1000 Bytes

Sender

Data1
Data2
Data3
Data4

Receiver

**SMSS** = 1 * 1000 Bytes
**IW** = 4 * 1000 Bytes
**SSTRESH** = 64 * 1000 Bytes

CWND = 4 * 1000 Bytes

Sender

Data1
Data2
Data3
ACK1+1
ACK4+1

Receiver

SMSS = 1 * 1000 Bytes
IW = 4 * 1000 Bytes
SSTRESH = 64 * 1000 Bytes



CWND = 4 * 1000 Bytes
CWND = 5 * 1000 Bytes
CWND = 8 * 1000 Bytes

Sender

Data1
Data2
Data3
ACK1+1

ACK4+1

Data5
Data6
Data7
Data8
Data9
Data10
ACK8+1
Data11
Data12

Receiver

SMSS        = 1     * 1000 Bytes
IW          = 4     * 1000 Bytes
SSTRESH     = 64    * 1000 Bytes

CWND = 4 *
1000 Bytes

CWND = 5 *
1000 Bytes

CWND = 8 *
1000 Bytes

CWND = 12
* 1000 Bytes

Sender

Data1
Data2
Data3
ACk1+1

ACk4+1

Data5
Data6
Data7
Data8
Data9
Data10
Data11
ACk8+1
Data12

Receiver

SMSS = 1 * 1000 Bytes
IW = 4 * 1000 Bytes
SSTRESH = 64 * 1000 Bytes

CWND = 4 * 1000 Bytes
CWND = 5 * 1000 Bytes
CWND = 8 * 1000 Bytes
CWND = 12 * 1000 Bytes
CWND = 16 * 1000 Bytes

Sender

Receiver

# Slow Start Variation 7 / 7

SMSS     = 1    * 1000 Bytes
IW       = 4    * 1000 Bytes
SSTRESH  = 64   * 1000 Bytes

# Slow Start Formulas

Growthrate of the cwnd:
- per window basis
**cwnd = cwnd * 2**
- per ACK basis
**cwnd = cwnd +1**

# DEMO
## SlowStart.pcapng

# TCP Reno History

- Evolution of TCP Reno
  - TCP Tahoe        (1988)
  - TCP Reno          (1990) RFC2001
  - TCP New Reno (1999)
    - Frirstly described in RFC 2582
    - Now **RFC 6582**

- TCP Tahoe and Reno where originally the codenames of their initial Berkeley Unix Distributions



Wireshark IO Graphs: Server_TO_C_read_filtered_anon_Erased

# TCP Reno

- TCP Tahoe
  - Slow Start
  - Fast Retransmission
  - Additive Increase Multiplicative Decrease (AIMD)
- TCP Reno
  - Fast Recovery
- TCP New Reno
  - Partial ACK



Wireshark IO Graphs: Server_TO_C_read_filtered_anon_Erased

- Selective Acknnowledgements (**SACK**)
  - Is no direct Part of TCP Reno History
  - It is a „side improvement" which is defined in **RFC7323** (RFC1323 <- old one)

**SMSS**      **= 1**     **\* 1000 Bytes**
**IW**          **= 2**     **\* 1000 Bytes**

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

Sender

Data64 Data65 Data66 Data67 Data68 Data69 Data70 Data71

Receiver

SMSS = 1 * 1000 Bytes
IW = 2 * 1000 Bytes

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

**3 received Duplicate ACKs will cause a „Fast Retransmission"!**

Sender

Receiver

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

ACK64 + 1
ACK65 + 1
ACK66 + 1
ACK66 + 1 DUP#1
ACK66 + 1 DUP#2
ACK66 + 1 DUP#3
ACK66 + 1 DUP#4

SMSS      = 1    * 1000 Bytes
IW        = 2    * 1000 Bytes

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 2 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 3 * 1000 Bytes

cwnd = cwnd +1

Sender

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

ACK64 + 1
ACK65 + 1
ACK66 + 1

ACK66 + 1 DUP#1
ACK66 + 1 DUP#2
ACK66 + 1 DUP#3
ACK66 + 1 DUP#4

Data67 FastRetrans

Data68 Retransn...

ACK67 + 1
ACK68 + 1

Receiver

SMSS     = 1     * 1000 Bytes
IW       = 2     * 1000 Bytes

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 2 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 3 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 4 * 1000 Bytes

cwnd = cwnd +1

Sender

Receiver

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

ACK64 + 1
ACK65 + 1
ACK66 + 1

ACK66 + 1 DUP#1
ACK66 + 1 DUP#2
ACK66 + 1 DUP#3
ACK66 + 1 DUP#4

Data67 FastRetrans

Data68 Retransmit

ACK67 + 1
ACK68 + 1

Data71 Retransmit
Data70 Retransmit
Data69 Retransmit

ACK69 + 1
ACK70 + 1
ACK71 + 1

Data72
Data73
Data74
Data75

SMSS          = 1      * 1000 Bytes
IW            = 2      * 1000 Bytes

SSTRESH = 8
* 1000 Bytes

CWND = 8 *
1000 Bytes

Sender

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

Receiver

SMSS   = 1   * 1000 Bytes
IW     = 2   * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 4 * 1000 Bytes

CWND = 5

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

| SSTRESH = ½ CWND |
| --- |
| **CWND = CWND + 1** by receiving on each DUP ACK after the DUP ACK #3 |

Sender

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

ACK64 + 1
ACK65 + 1
ACK66 + 1
ACK66 + 1 DUP#1
ACK66 + 1 DUP#2
ACK66 + 1 DUP#3
ACK66 + 1 DUP#4

Receiver

SMSS = 1 * 1000 Bytes
IW = 2 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes

CWND = 4 * 1000 Bytes

CWND = 5

SSTRESH = 8 * 1000 Bytes

CWND = 8 * 1000 Bytes

CWND = 4 * 1000 Bytes

**CWND = SSTRESH**
by receiving ACK for retransmitted Packet 67
-> Fast Recovery has ended

Sender

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

ACK64 + 1
ACK65 + 1
ACK66 + 1

ACK66 + 1 DUP#1
ACK66 + 1 DUP#2
ACK66 + 1 DUP#3
ACK66 + 1 DUP#4

Data67 Fast Retransmit
Data68 Retransmit

ACK67 + 1
ACK68 + 1

Receiver

# TCP Reno: Fast Recovery 5 / 5

SMSS = 1 * 1000 Bytes
IW = 2 * 1000 Bytes

SSTRESH = 8 * 1000 Bytes
CWND = 8 * 1000 Bytes

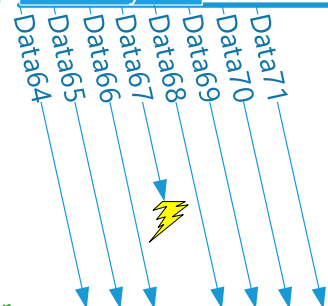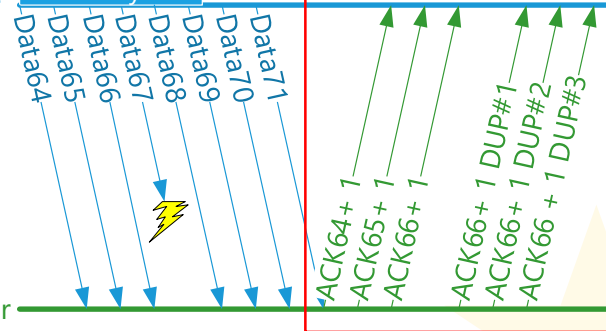SSTRESH = 4 * 1000 Bytes
CWND = 4 * 1000 Bytes
CWND = 5
CWND = 6

SSTRESH = 4 * 1000 Bytes
CWND = 4 * 1000 Bytes

SSTRESH = 4 * 1000 Bytes
CWND = 5 * 1000 Bytes

CWND = CWND + 1
Congestion Avoidance Mode

Sender

Data64
Data65
Data66
Data67
Data68
Data69
Data70
Data71

Data67 Fast Retransmit
Data68 Retransmit

Data69 Retransmit
Data70 Retransmit
Data71 Retransmit
Data72
Data73
Data74
Data75

ACK64 + 1
ACK65 + 1
ACK66 + 1
ACK66 + 1
ACK66 + 1
ACK66 + 1
ACK66 + 1
ACK67 + 1
ACK68 + 1

ACK69 + 1
ACK70 + 1
ACK71 + 1
ACK72 + 1

Receiver

- Selective Acknowledgements (SACK)
  - Allows us to send DUP ACK to a specific Paket
    - E.g. we send 31, 32, 33, 34
    - Packet 32 is lost
    - With SACK we can say:
      - We are awaiting 35 but we are still missing 32

- New Reno introduces Partial ACK
  - Workaround when SACK is not supported

# Tahoe and Reno Graphs

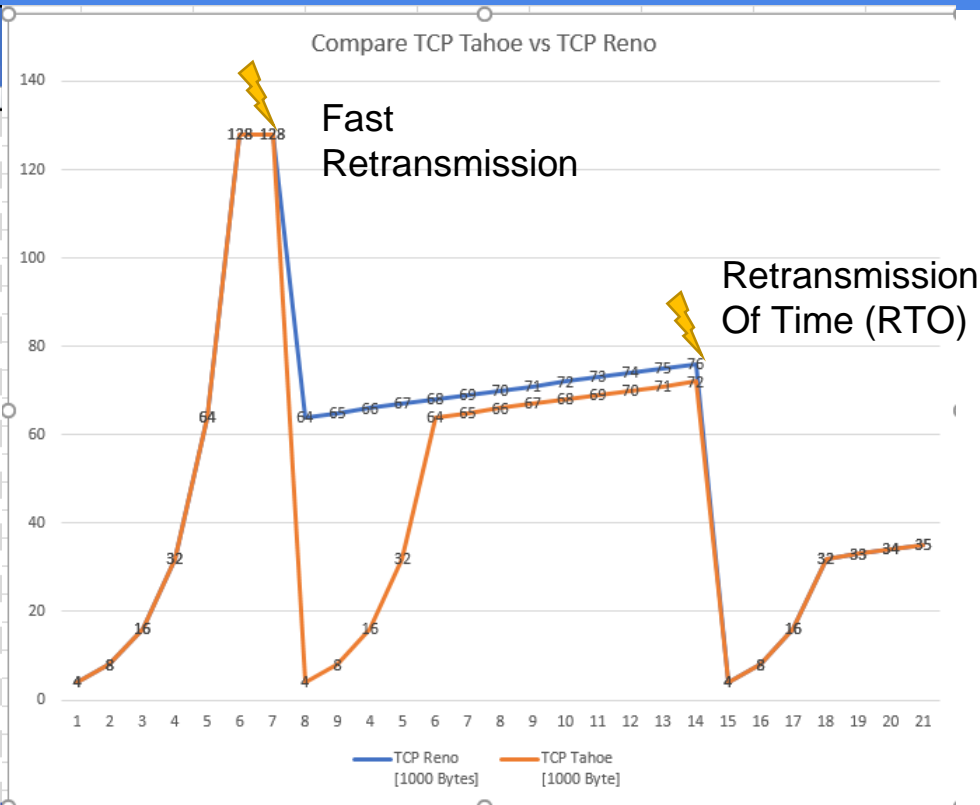| Time [cycles] | TCP Reno [1000 Bytes] | TCP Reno cwnd [1000 Bytes] | TCP Tahoe [1000 Byte] | TCP Tahoe cwnd [1000 Byte] | sstresh [1000 Bytes] | SMSS [1000 Bytes] | rwnd [1000 bytes] |
|---|---|---|---|---|---|---|---|
| 1 | 4 | 4 | 4 | 4 | 128 | 1 | 256 |
| 2 | 8 | 8 | 8 | 8 | 128 | 1 | 256 |
| 3 | 16 | 16 | 16 | 16 | 128 | 1 | 256 |
| 4 | 32 | 32 | 32 | 32 | 128 | 1 | 256 |
| 5 | 64 | 64 | 64 | 64 | 128 | 1 | 256 |
| 6 | 128 | 128 | 128 | 128 | 128 | 1 | 256 |
| 7 | 128 | 128 | 128 | 128 | 128 | 1 | 256 |
| 8 | 64 | 64 | 4 | 4 | 64 | 1 | 256 |
| 9 | 65 | 65 | 8 | 8 | 64 | 1 | 256 |
| 4 | 66 | 66 | 16 | 16 | 64 | 1 | 256 |
| 5 | 67 | 67 | 32 | 32 | 64 | 1 | 256 |
| 6 | 68 | 68 | 64 | 64 | 64 | 1 | 256 |
| 7 | 69 | 69 | 65 | 65 | 64 | 1 | 256 |
| 8 | 70 | 70 | 66 | 66 | 64 | 1 | 256 |
| 9 | 71 | 71 | 67 | 67 | 64 | 1 | 256 |
| 10 | 72 | 72 | 68 | 68 | 64 | 1 | 256 |
| 11 | 73 | 73 | 69 | 69 | 64 | 1 | 256 |
| 12 | 74 | 74 | 70 | 70 | 64 | 1 | 256 |
| 13 | 75 | 75 | 71 | 71 | 64 | 1 | 256 |
| 14 | 76 | 76 | 72 | 72 | 64 | 1 | 256 |
| 15 | 4 | 4 | 4 | 4 | 32 | 1 | 256 |
| 16 | 8 | 8 | 8 | 8 | 32 | 1 | 256 |
| 17 | 16 | 16 | 16 | 16 | 32 | 1 | 256 |
| 18 | 32 | 32 | 32 | 32 | 32 | 1 | 256 |
| 19 | 33 | 33 | 33 | 33 | 32 | 1 | 256 |
| 20 | 34 | 34 | 34 | 34 | 32 | 1 | 256 |
| 21 | 35 | 35 | 35 | 35 | 32 | 1 | 256 |



Compare TCP Tahoe vs TCP Reno

Fast Retransmission

Retransmission Of Time (RTO)

TCP Reno [1000 Bytes] — TCP Tahoe [1000 Byte]

# DEMO

- **FastRetransmission yes/no (SlowStart)**

- **Linear Growth**

- **ZigSawGraph**

# TCP Reno Summary

- **Fast Retransmit causes Fast Recovery:**
  - **cwnd = ½ cwnd + 1 per every DUP ACK > #3**
  - **ssthresh = max (BytesInFlight / 2,
              2*SMSS)**
  - Fast Recovery ends with receiving ACK for all lost packets with
    - cwnd = sstresh
  - **CWND growth after that**
    - **CWND = CWND + 1** per Window cycle
- **RTO causes:**
  - **ssthresh = max (BytesInFlight / 2,
              2*SMSS)**
  - **cwnd = Initial Window**

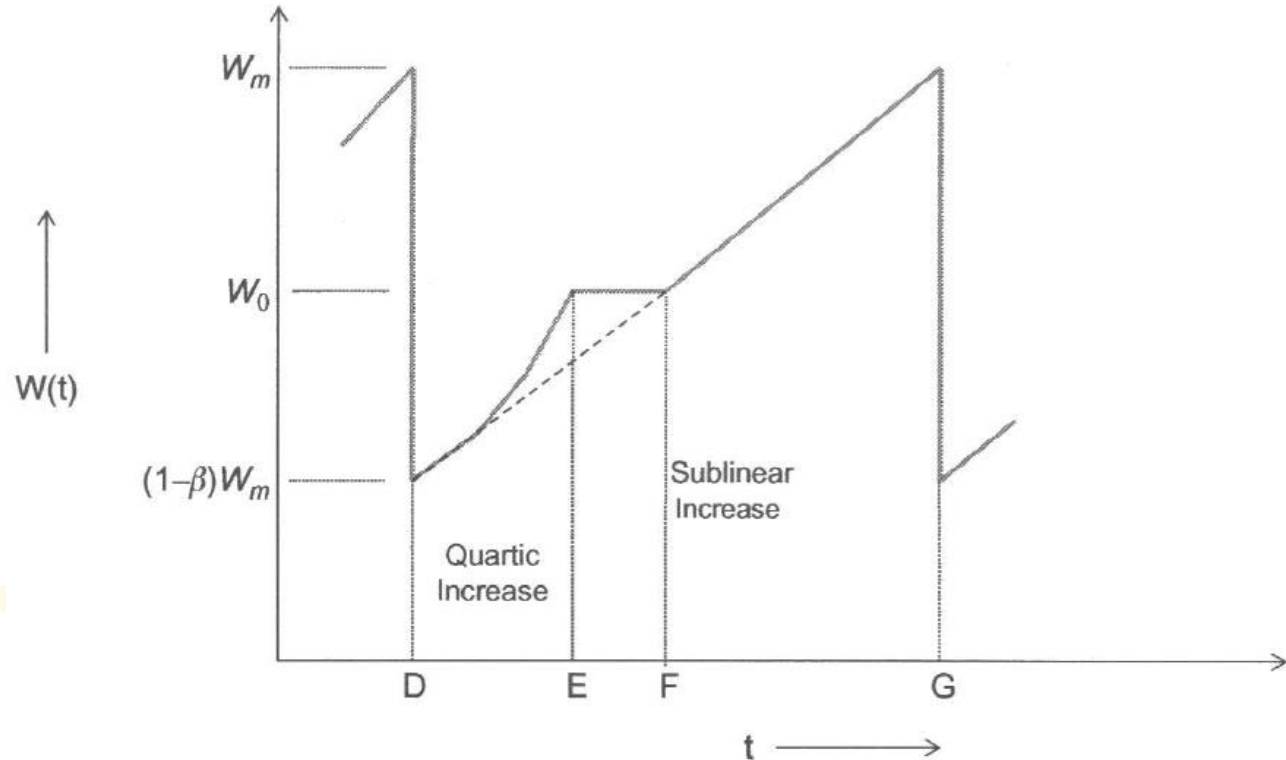- If **cwnd** < **sstresh** then
  **Slow Start** is used

- If **cwnd** > **sstresh** then
  **congestion avoidance** is used

- After Packet Loss **sstresh** can never grow larger then **½ of cwnd**

- **1 Retransmitted can significant Slow Down whole Session**

- **CTCP**

- BIC



Additive Increase | Binary Search | Max Probing

# Thank You!

- Thank you!

- I am happy about **FEEDBACK** at Sharkfest Europe Guidebook

# About me?

- Christian Reusch
  - Analyzing Networks since 1999
  - Twitter: @crnetpackets
  - Web:    crnetpackets.com

  - If you like you can send my Traces and I try to answer
    - creusch@crnetworks.com