# The Transmission Control Protocol

Dorgham Sisalem
GMD Fokus, Berlin
`sisalem@fokus.gmd.de`

**fo‹us**

---

# What are transmission protocols needed for?

**Addressing:** application to application addressing

**Reliable delivery:** the receiver application should receive the same data stream the source puts on the net

**Segment order maintenance:** data segments should reach the application in the same order they left the sender

**Flow control:** the data sending speed should adapt itself to the receivers speed
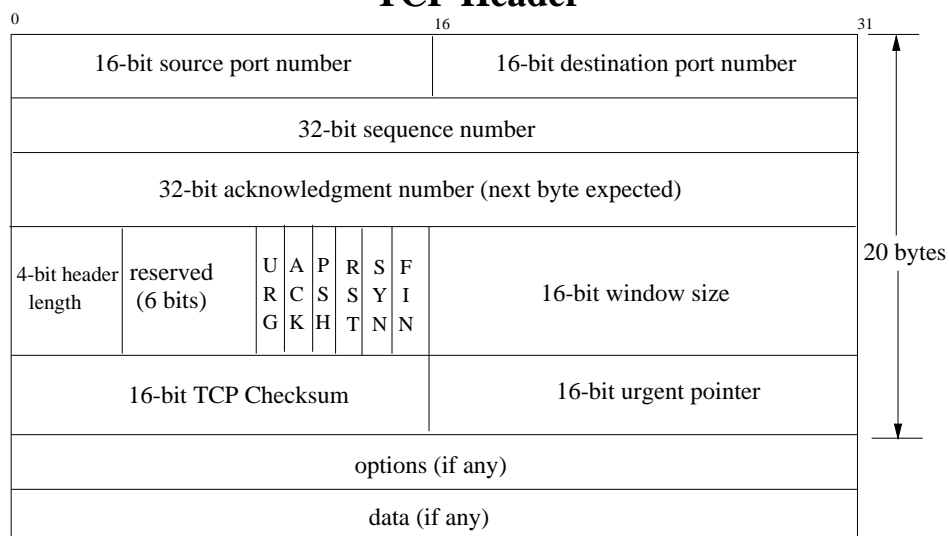
**Congestion control:** the transmission speed can not be faster than the speed of the slowest link traversed on the connections path

**Segmentation:** data is sent in segments that provide the highest throughout.

**fo‹us**

## Transmission Control Protocol

- TCP is connection oriented and full duplex.

- The maximum segment size *(MSS)*is set during connection establishment.

- Reliability is achieved using acknowledgments, round trip delay estimations and data retransmission.

- TCP uses a variable window mechanism for flow control.

- Congestion control and avoidance is reached using slow start and congestion avoidance schemes.
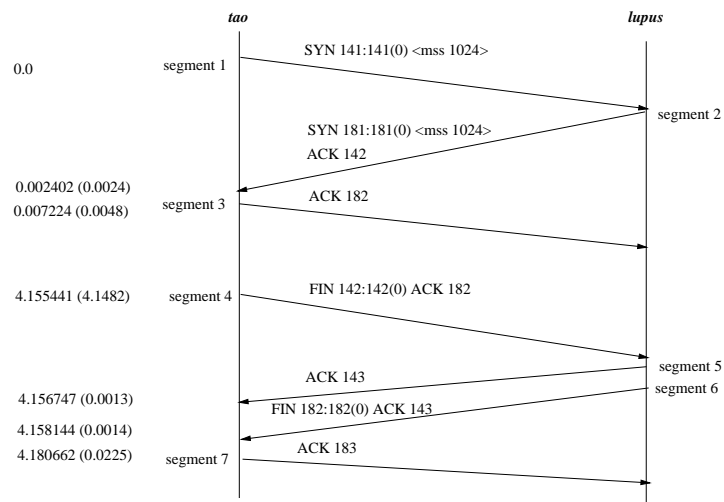
GMD fo‹us

## TCP Header

| 0 | | | | | | | | 16 | 31 |
|---|---|---|---|---|---|---|---|---|---|
| 16-bit source port number | | | | | | | | 16-bit destination port number | |
| 32-bit sequence number | | | | | | | | | |
| 32-bit acknowledgment number (next byte expected) | | | | | | | | | |
| 4-bit header length | reserved (6 bits) | U R G | A C K | P S H | R S T | S Y N | F I N | 16-bit window size | |
| 16-bit TCP Checksum | | | | | | | | 16-bit urgent pointer | |
| options (if any) | | | | | | | | | |
| data (if any) | | | | | | | | | |

20 bytes

- The most common option is the maximum segment size option.

GMD fo‹us

## Connection Establishment and Termination

- Connection establishment is done with a three way handshake

*tao*　　　　　　　　　　　　*lupus*

SYN 141:141(0) <mss 1024>

0.0　　　　　　segment 1

　　　　　　　　　　　　　　　　　　　　segment 2

SYN 181:181(0) <mss 1024>
ACK 142

0.002402 (0.0024)
0.007224 (0.0048)　　segment 3　　ACK 182

4.155441 (4.1482)　　segment 4　　FIN 142:142(0) ACK 182

　　　　　　　　　　　　　　　　　　　　segment 5
ACK 143　　　　　　　　　　　　　segment 6
4.156747 (0.0013)
FIN 182:182(0) ACK 143
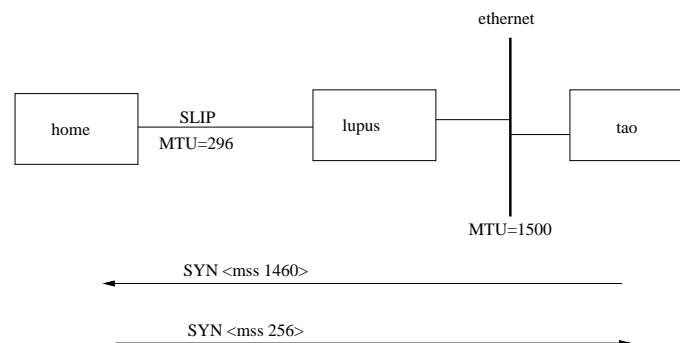4.158144 (0.0014)
4.180662 (0.0225)　　segment 7　　ACK 183

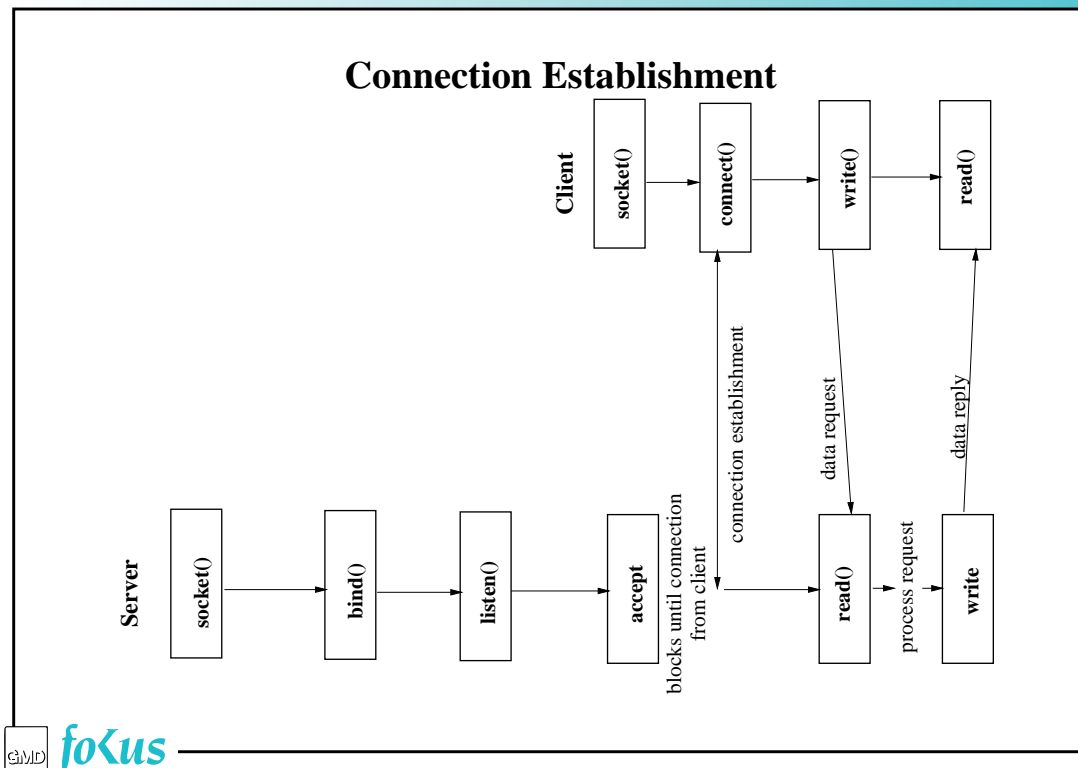- Each side can just close its transmission side resulting in a half close.

## Connection Establishment

- *tao* sends a SYN segment with an initial sequence number (ISN) and and the maximum segment size (MSS) it is willing to receive.

- *lupus* replies with a SYN segment acknowledging ISN and announcing its MSS.

- MSS can be at the most as large as the interface segment size minus 40 bytes.

ethernet

home　　SLIP
MTU=296　　lupus　　　　tao

MTU=1500

SYN <mss 1460>

SYN <mss 256>

# Connection Establishment

**Client**

socket() → connect() → write() → read()

connection establishment

data request

data reply

**Server**

socket() → bind() → listen() → accept

blocks until connection from client

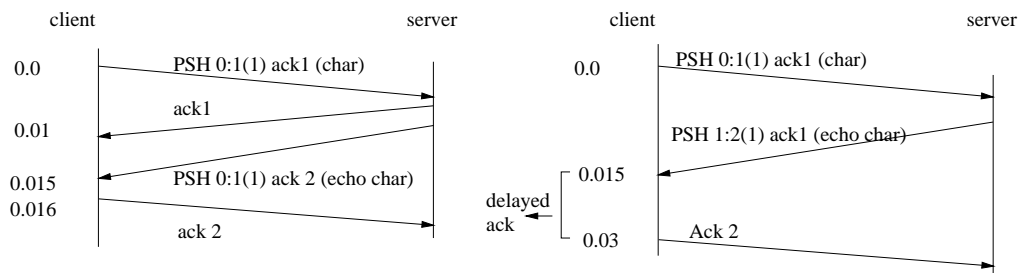read() → process request → write

GMD fo‹us

# Connection Termination

- A sender terminates its part of the connection by sending a FIN segment.

- After acknowledging the FIN the receiver can still send data on its part of the connection *(half close)*.

- A connection can be aborted with a RST segment.

GMD fo‹us

# Interactive Data Transfer

- Data received from the application is usually sent in segments of MSS.

- In the case of interactive applications -rlogin, telnet- the sender can force the delivery of small packets using the PSH (push) flag.

- With delayed acknowledgments the receiver delays sending the acknowledgments until it has some data to send or a 200 ms timer expires.

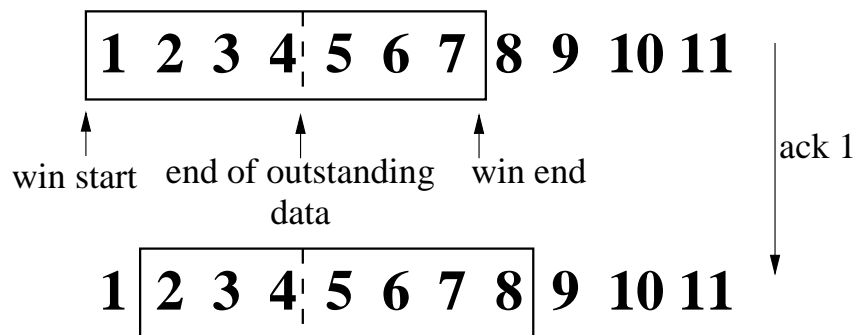| client | | server | client | | server |
|---|---|---|---|---|---|
| 0.0 | PSH 0:1(1) ack1 (char) → | | 0.0 | PSH 0:1(1) ack1 (char) → | |
| 0.01 | ← ack1 | | | PSH 1:2(1) ack1 (echo char) → | |
| 0.015 | ← PSH 0:1(1) ack 2 (echo char) | | 0.015 | | |
| 0.016 | ack 2 → | | delayed ack ← 0.03 | Ack 2 → | |

fo<us

---

# Interactive Data Transfer

- Sending a lot of small segments can add congestion to a wide area network.

- *Nagle Algorithm*: a sender can at most have one outstanding small segment, that has not yet been acknowledged.

- All data arriving at TCP from the application are queued until the currently outstanding segment is acknowledged.
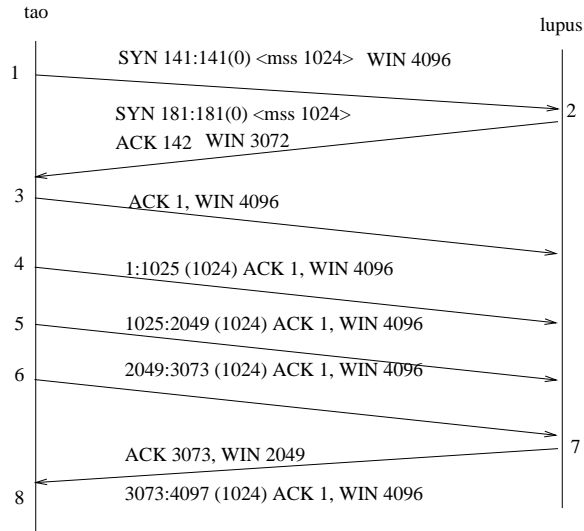
fo<us

## Flow Control in TCP

- TCP uses a sliding window mechanism to adjust the senders transmission speed to that of the receiver.

- The sliding window permits the sending of multiple segments before waiting for an acknowledgment.

- Ack segments indicate the last correctly received byte and the number of bytes the receiver is still willing to accept.

**1 2 3 4 5 6 7 8 9 10 11**

win start   end of outstanding   win end
data

ack 1

**1 2 3 4 5 6 7 8 9 10 11**

fo‹us

---

## Flow Control in TCP

- Ack segments indicate the last correctly received byte and the number of bytes the receiver is still willing to accept.

tao             lupus

1    SYN 141:141(0) <mss 1024>  WIN 4096

    SYN 181:181(0) <mss 1024>    2
    ACK 142  WIN 3072

3    ACK 1, WIN 4096

4    1:1025 (1024) ACK 1, WIN 4096

5    1025:2049 (1024) ACK 1, WIN 4096

6    2049:3073 (1024) ACK 1, WIN 4096

                 7
    ACK 3073, WIN 2049

8    3073:4097 (1024) ACK 1, WIN 4096

fo‹us

## Acknowledgments and Retransmission

- A TCP receiver always acknowledges the last correctly received byte.

- After sending a segment the sender starts a timer.

- If the timer expires before receiving an acknowledgment for the sent segment the segment is considered lost and must be retransmitted.

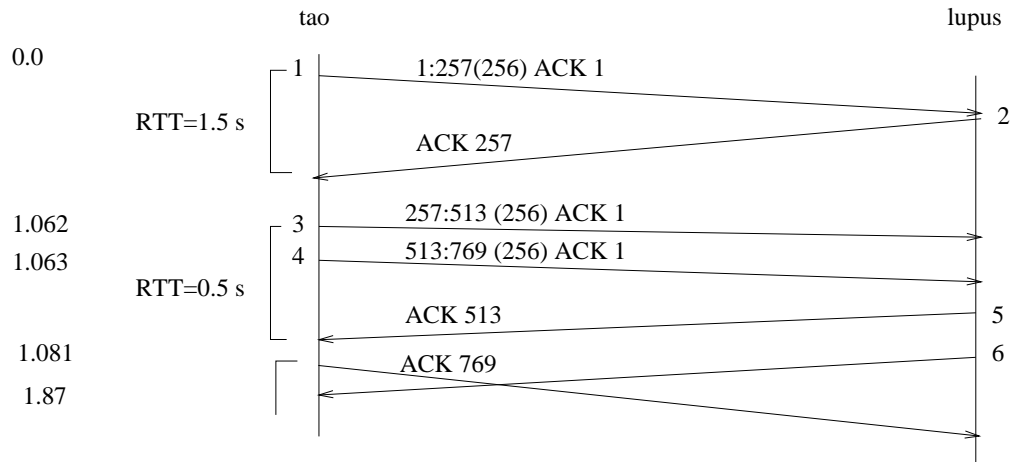- The timeout value is calculated dynamically according to the measured round trip times *(RTT)*.

$$
\begin{aligned}
\text{Err} &= \text{RTT} - A & A &= \text{smoothed RTT}\\
A &= A + g\,\text{Err} & \text{gain } g &= 1/8\\
D &= D + h\left(|\text{Err}| - D\right) & D &= \text{smoothed mean deviation}\\
\text{RTO} &= A + 4D
\end{aligned}
$$

**fo‹us**

## Round-Trip Time Measurement

- TCP implementations use a 500-ms clock for time measurements and timeout determination.

- Only one measurement is done at a time.

- At the start of a measurement a counter is set to 0 and is then incremented every time the 500-ms TCP timer is invoked and the number of the sent segment is remembered.

- Only after acknowledging the sent segment can a new measurement start.

- After a retransmission the timeout value is not updated until an acknowledgment for a segment arrives that was not retransmitted *(Karns's algorithm)*.
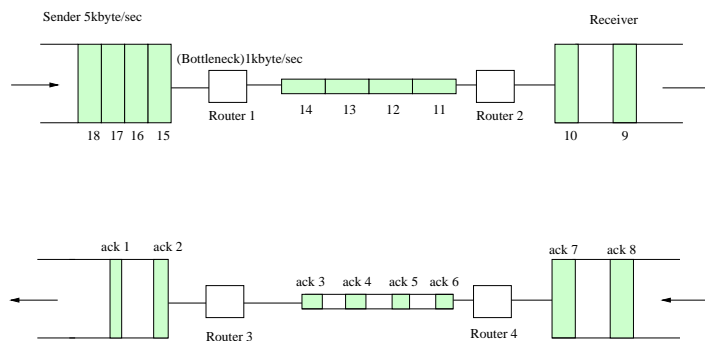
**fo‹us**

# Round-Trip Time Measurement

- As the 500-ms timer is used for determining the RTT the values used for updating the timeout value might differ up to $\pm 500$ ms from the actual value.



# Congestion Control in TCP

- A connection's rate is determined as transmission window/round trip time.

- When the sum of the connection rates over a link is higher than the link's rate segments can be dropped.

- TCP uses packet drops and timeouts as congestion indication.

## Slow Start and Congestion Avoidance

- To avoid congestion in advance, the sender must adapt its transmission window to the available link bandwidth.

- On connection establishment TCP uses a window of the size of 1 MSS *Congestion Window*.

- The congestion window is increased by 1 MSS for each acknowledged segment.

- At any time the sender has has a transmission window of

  transmission window = min (advertised window, congestion window)

fo‹us

## Slow Start and Congestion Avoidance

- With the slow start scheme the congestion window is exponentially increased.

- This can quickly congest the network and cause packet drops.

- After a timeout the congestion window is set again to 1 MSS.

- Slow start is reused but only until the congestion window reaches half of its value before the timeout.

- Afterwards the congestion window is increased only by 1/congestion window for each acknowledged segment *(congestion avoidance)*.

fo‹us

# Fast Retransmission and Fast Recovery

- Using only timeouts as loss indication leads to long idle periods.

- With the *fast retransmission* scheme the receiver acknowledges each out of order segment with an ack of the last correctly received segment.

- Receiving 3 duplicate acks triggers at the source the retransmission of the last acked segment.

- In the older TCP versions the same actions taken after a timeout are used in this case as well.

- in TCP versions using *fast recovery* the congestion window is only reduced by half after each loss.

**fo‹us**

# Congestion Example

- Both source and receiver have buffers up to 8192 bytes.

- The router has a buffer of 2128 bytes.

- The link has a bandwidth of 2128 bytes/sec.

- MSS=1024.

- The configuration has a round trip delay of 1 sec.



**fo‹us**

# Congestion Example: Slow Start

*lupus* | *tao*

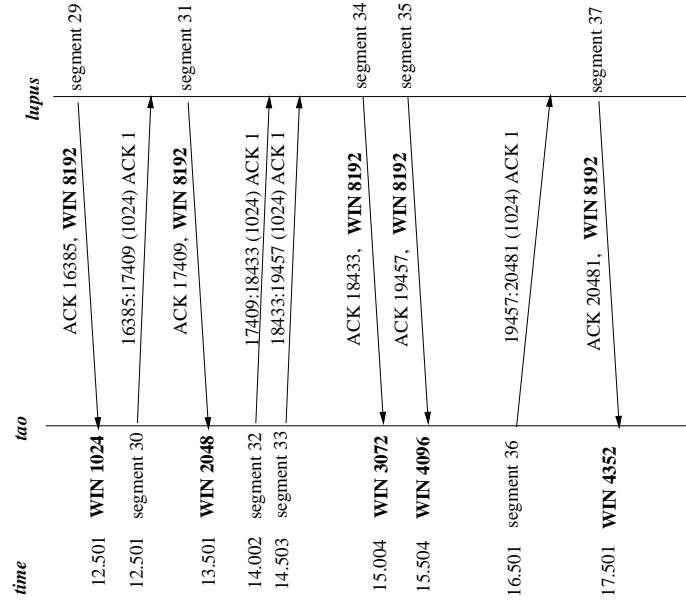| time | | |
|---|---|---|
| | 0:1025 (1024) ACK 1 | segment 3, **win=1024** |
| 1 | | |
| | ACK 1025, **WIN 8192** | **win=2048** |
| 2 | segment 4 | |
| 2.0001 | 1025:2049 (1024) ACK 1 | segment 5 |
| 2.0002 | 2049:3073 (1024) ACK 1 | segment 6 |
| | ACK 2049, **WIN 8192** | **win=3072** |
| 3.001 | segment 7 | |
| | ACK 3073, **WIN 8192** | **win=4096** |
| 4.0003 | segment 8 | |
| 4.0003 | 3073:4097 (1024) ACK 1 | segment 9 |
| 4.5004 | 4097:5121(1024) ACK 1 | segment 10 |
| 5.5005 | 5121:6145 (1024) ACK 1 | segment 11 |
| 6.5006 | 6145:7169 (1024) ACK 1 | segment 12 |
| | ACK 4097, **WIN 8192** | **win=5121** |
| 5.0004 | ACK 5121, **WIN 8192** | **win=6144** |
| 5.5006 | ACK 6145, **WIN 8192** | **win=7168** |
| 6.5004 | ACK 7169, **WIN 8192** | **win=8192** |
| 7.5006 | segment 13 / segment 14 / segment 13 / segment 14 | |
| 8.0007 | 7169:8193 (1024) ACK 1 | segment 15 |
| 8.5008 | 9217:10241 (1024) ACK 1 | segment 16 |
| 8.5009 | 10241:11265 (1024) ACK 1 | segment 18 |

---

# Congestion Example: Fast Retransmission

*lupus* | *tao*

| time | | |
|---|---|---|
| 8.501 | 11265:12289 (1024) ACK 1 | segment 19 |
| 8.5011 | 12289:13313 (1024) ACK 1 | segment 20 |
| 9.0001 | ACK 12289, **WIN 8192** | |
| | segment 21 | |
| 9.0001 | 13313:14337 (1024) ACK 1 | segment 22 |
| 10.002 | 14337:15361 (1024) ACK 1 | segment 23 |
| 10.503 | 15361:16385 (1024) ACK 1 | segment 24 |
| | ACK 12289, **WIN 7168** | segment 25 |
| 10.0001 | ACK 12289, **WIN 6144** | segment 26 |
| 11.002 | ACK 12289, **WIN 5120** | segment 27 |
| 11.501 | | |
| 11.501 | 12289:13313 (1024) ACK 1 | segment 28 |

# Congestion Example: Fast Retransmission

lupus

segment 29
segment 31
segment 34
segment 35
segment 37

ACK 16385, **WIN 8192**
16385:17409 (1024) ACK 1
ACK 17409, **WIN 8192**
17409:18433 (1024) ACK 1
18433:19457 (1024) ACK 1
ACK 18433, **WIN 8192**
ACK 19457, **WIN 8192**
19457:20481 (1024) ACK 1
ACK 20481, **WIN 8192**

tao

**WIN 1024**
segment 30
**WIN 2048**
segment 32
segment 33
**WIN 3072**
**WIN 4096**
segment 36
**WIN 4352**

time

12.501
12.501
13.501
14.002
14.503
15.004
15.504
16.501
17.501

# Congestion Example: Fast Recovery

lupus

segment 29
segment 33
segment 34
segment 35
segment 36

ACK 16385, **WIN 8192**
16385:17409 (1024) ACK 1
17409:18433 (1024) ACK 1
18433:19457 (1024) ACK 1
19457:20481 (1024) ACK 1
ACK 17409, **WIN 8192**
ACK 18433, **WIN 8192**
ACK 19457, **WIN 8192**
ACK 20481, **WIN 8192**

tao

**WIN 4096**
segment 30
segment 31
segment 32
segment 32
**WIN 4352**
**WIN 4593**
**WIN 4821**
**WIN 5038**

time

12.501
12.505
13.002
13.503
14.503
13.505
14.004
14.504
15.504

## Silly Window Syndrome and Probe Packets

- It is possible for the advertised window size to go to 0.

- After freeing some buffer, the receiver sends an update message with the size of the available buffer.

- After receiving an ack with WIN=0 the sender starts a persist timer.

- On the expiration of the persist timer a small packet of 1 byte payload is sent to see if a window update message got lost -such a packet is called *probe packet*.

- To avoid sending small packets the receiver must not advertise small segments, i.e., segments smaller than MSS *(silly window syndrome)*.

fo‹us

## TCP Future and Performance

- The capacity of a link is measured as

$$capacity \; = \; bandwidth*RTT$$

- The throughput of TCP is limited to

$$throughput \; = \; \frac{max\ window\ size}{RTT}$$

- Using a window scale option improves performance on long fat pipes.

- Updating the RTO value every RTT leads to aliasing effects.

- More accurate timeout calculations can be reached using a time stamp option

fo‹us

## T/TCP

- Lots of TCP transactions consist simply of a request to a server and a reply to the client.

- This simple transaction require the sending of 10 segments.

- Due to the connection establishment and termination a simple transaction requires at least two RTT times plus the processing time required at the server.

**fo‹us**

## T/TCP

- To distinguish between consecutive transactions a connection count option is added to the header.

- To avoid unnecessary overhead a host might maintain a per-host cache of the last seen timeout value, MSS, window size and the CC value.

- A client can combine the SYN, FIN, data request and the current CC value in one segment.

- If the received CC value is larger than the cached CC the server can combine the SYN, FIN, ACK of the sender's SYN and the reply in one segment.

- The client acks the server's SYN and FIN in one segment.

- This minimal transactions reduces the time needed for the transaction to a minimum of RTT plus the processing time at the server.

**fo‹us**