

CS 925

Lecture 10

TCP Congestion Control

Thursday, February 22, 2024

TCP Congestion Control

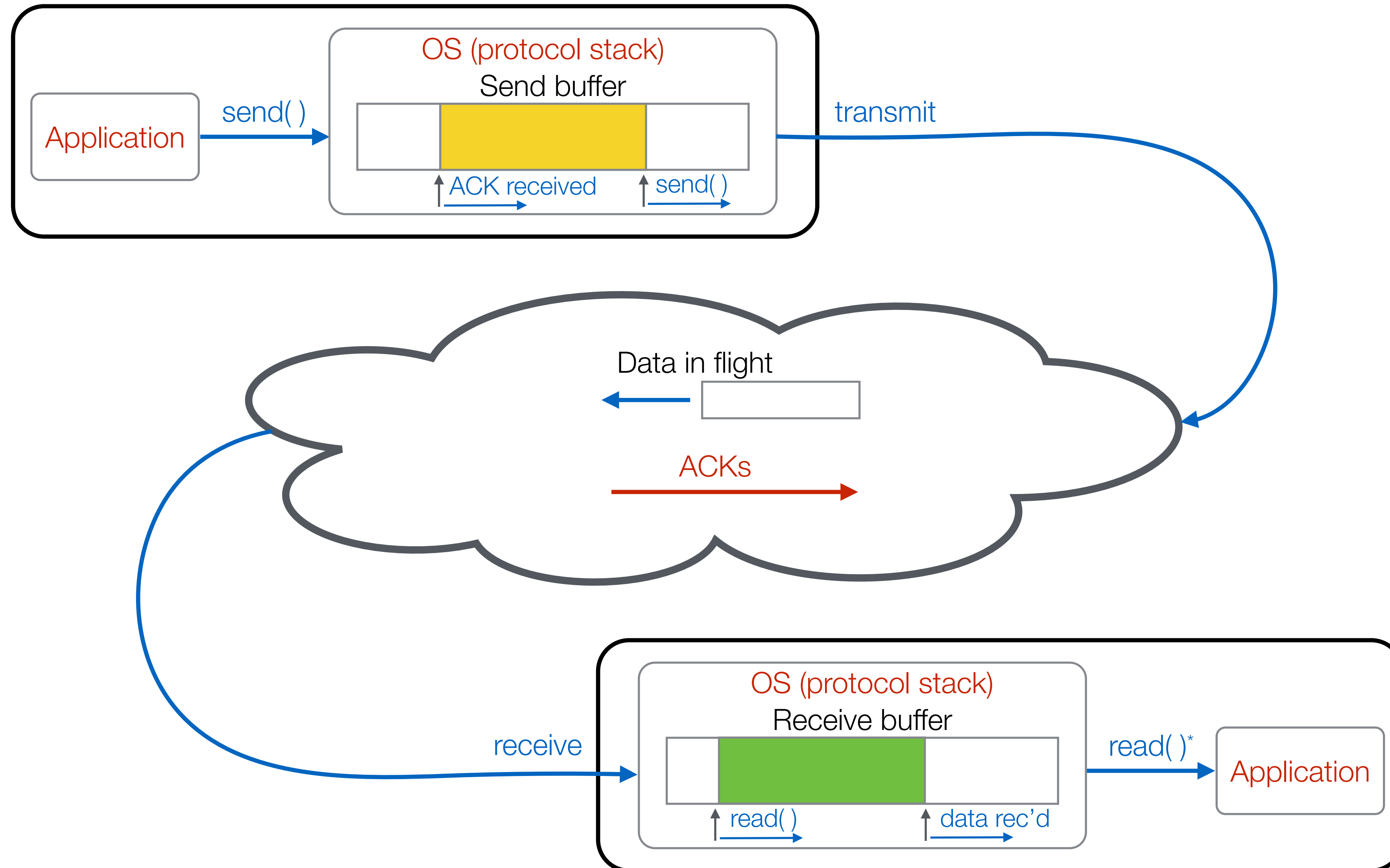
▶ Receiver congestion (flow control)

- Window Size field - explicitly reported by the receiver
- TCP Window Scale Option

▶ Network congestion

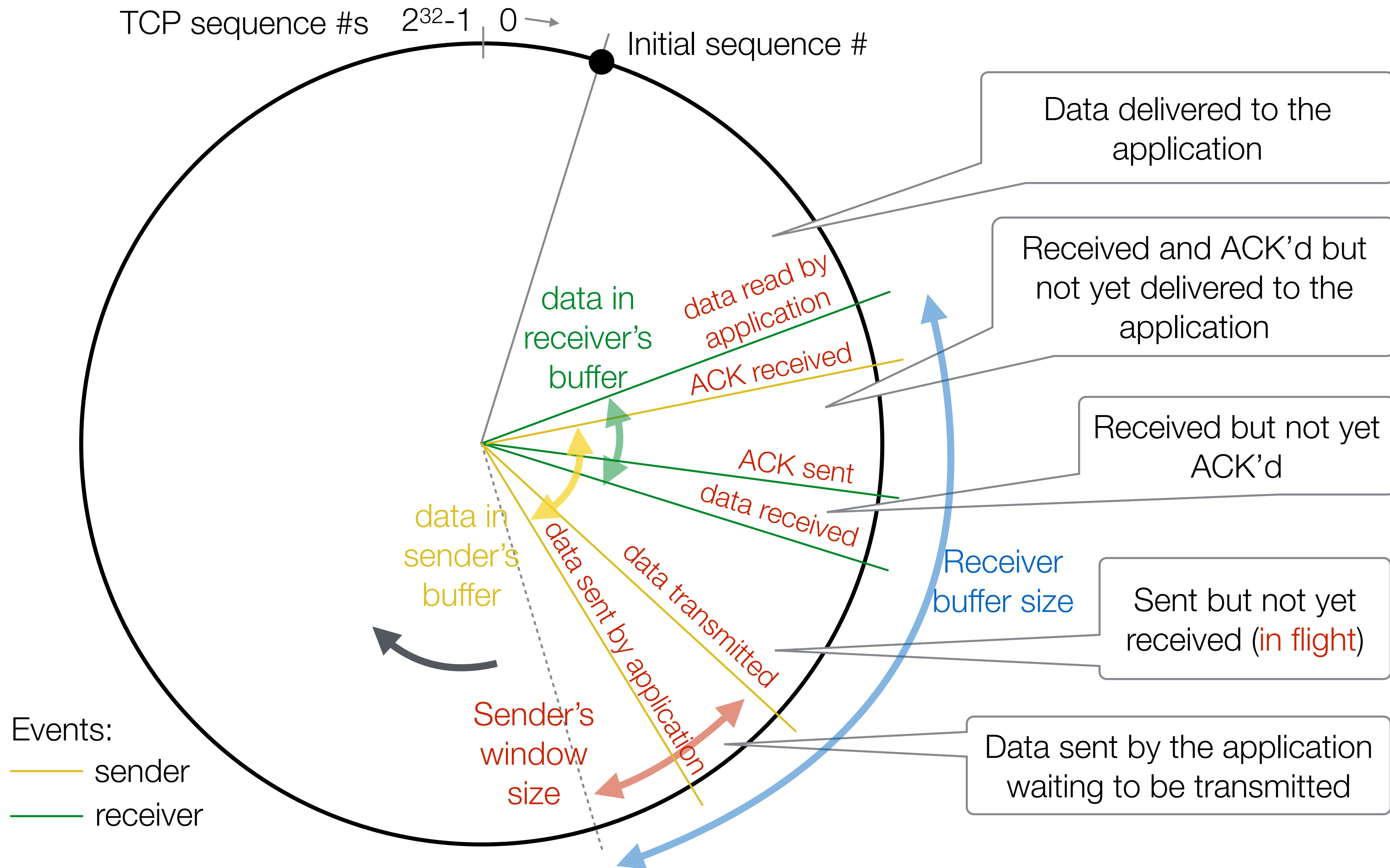
- Retransmission timeout
- Transmission window
- Set based on observed RTT and on detected packet loss

TCP buffering and data flow

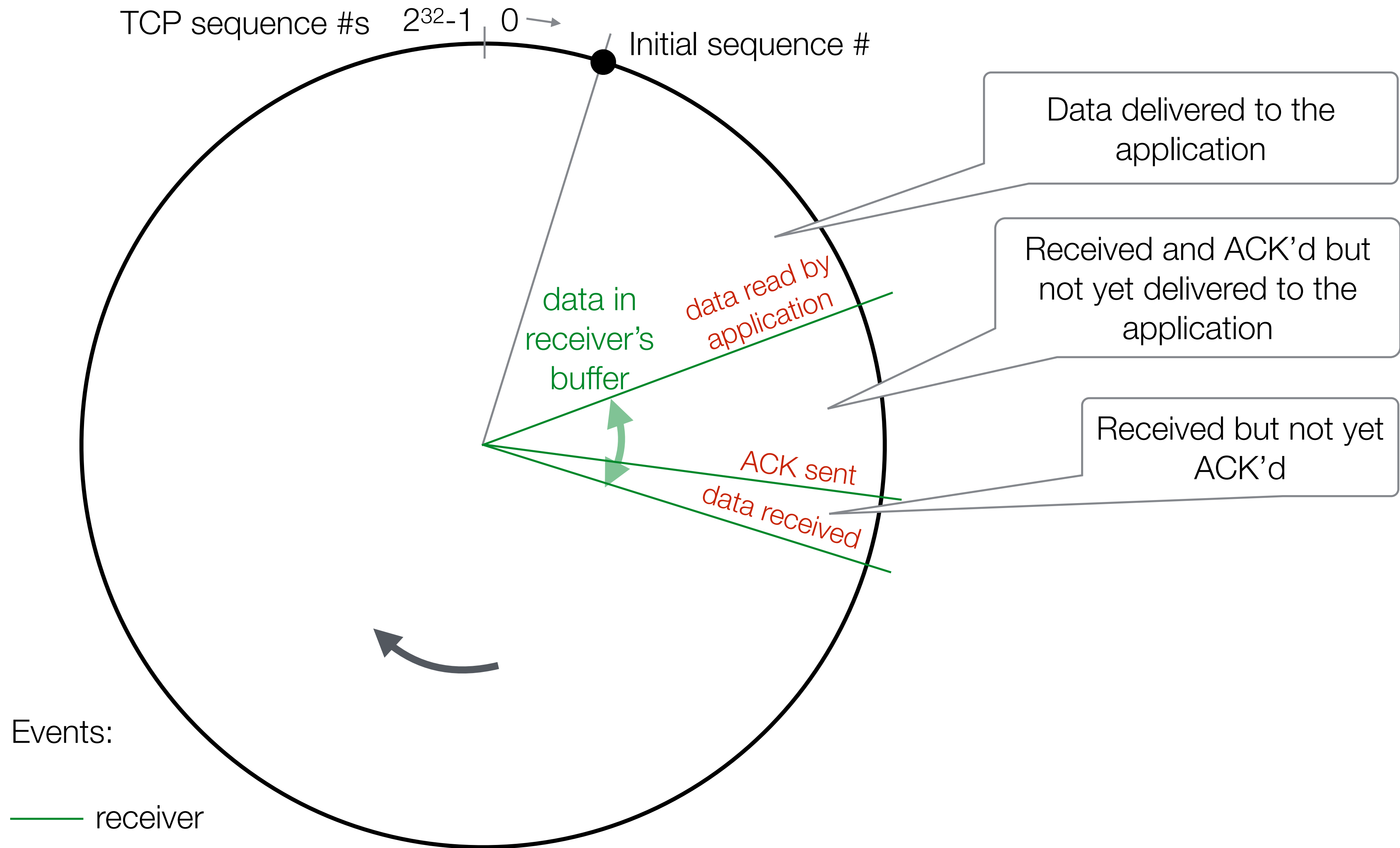


(*) many APIs call the `read()` operation "receive" (eg: `recv()`), `read` is used here to avoid confusion with receiving data on an interface

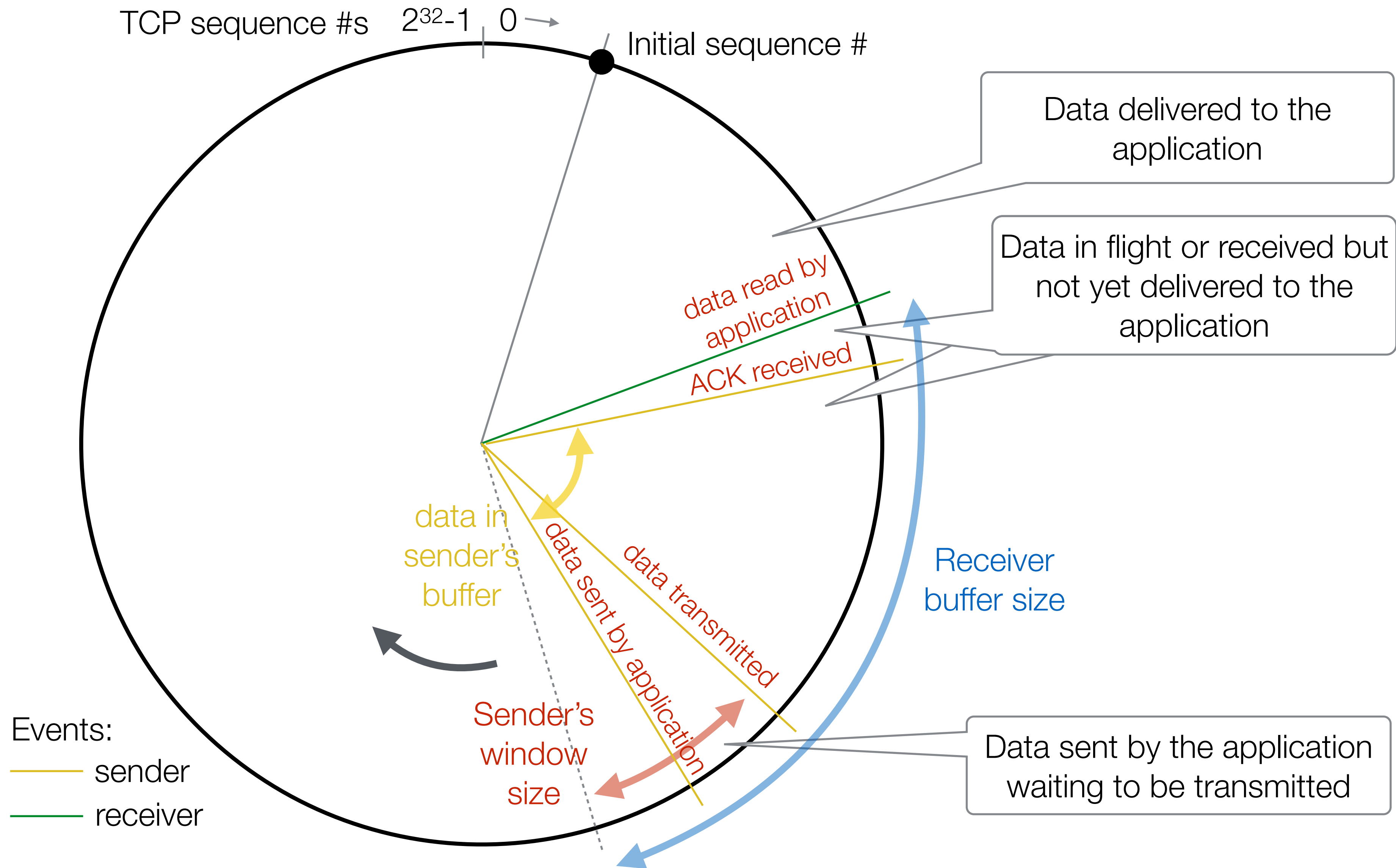
TCP Sliding Window



TCP Sliding Window



TCP Sliding Window



TCP SACK

- ▶ TCP Selective Acknowledgment (SACK)
 - negotiated during connection open
 - additional info about **received data** on top of TCP's cumulative acknowledgment (ACK# field)
 - a TCP option negotiated when a connection is opened

```
Options: (24 bytes), Maximum segment size, No-Operation
  > TCP Option - Maximum segment size: 1440 bytes
  > TCP Option - No-Operation (NOP)
  > TCP Option - Window scale: 5 (multiply by 32)
  > TCP Option - No-Operation (NOP)
  > TCP Option - No-Operation (NOP)
  > TCP Option - Timestamps: TSval 819433519, TSecr 0
    Kind: Time Stamp Option (8)
    Length: 10
    Timestamp value: 819433519
    Timestamp echo reply: 0
  > TCP Option - SACK permitted
  > TCP Option - End of Option List (EOL)
```

Wireshark decode of a TCP SYN packet

TCP SACK

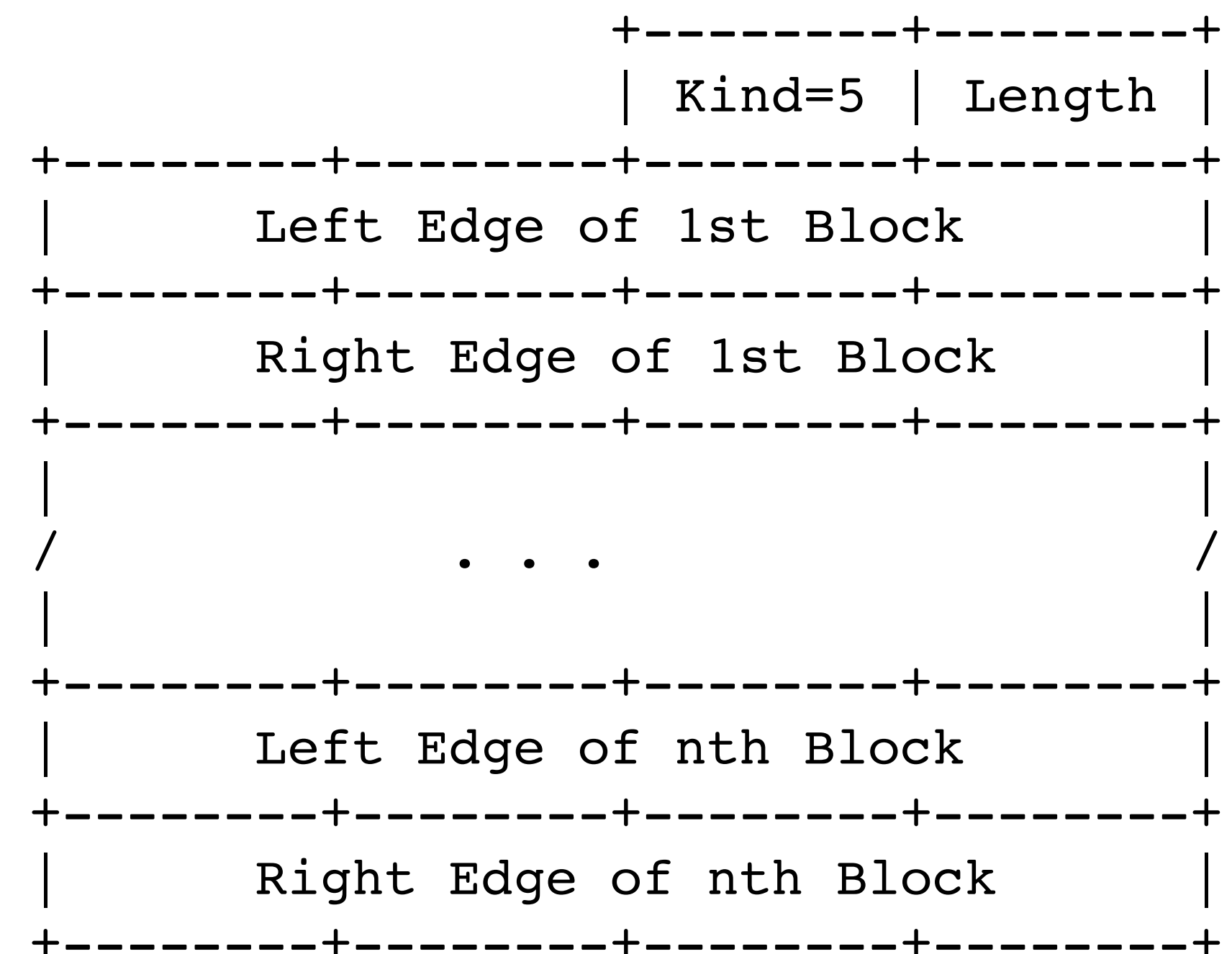
▶ TCP Selective Acknowledgment (SACK)

- negotiated during connection open
- additional info about **received data** on top of TCP's cumulative acknowledgment (ACK# field)
- a TCP option negotiated when a connection is opened

TCP SACK Option:

Kind: 5

Length: Variable



Source: RFC2018

Network Congestion Control

- ▶ ACK self-clocking
- ▶ Retransmission timer management
- ▶ Additive Increase Multiplicative Decrease (AIMD)
- ▶ Slow start mechanism
- ▶ ...

Retransmission Timeout

Initialization:

RFC 6298

$$\text{RTO} \leftarrow 1 \text{ sec}$$

After the first measurement:

$$\text{SRTT} \leftarrow R$$

$$\text{RTTVAR} \leftarrow R/2$$

$$\text{RTO} \leftarrow \text{SRTT} + \max(G, K * \text{RTTVAR})$$

After subsequent measurements:

$$\text{RTTVAR} \leftarrow (1 - \text{beta}) * \text{RTTVAR} + \text{beta} * |\text{SRTT} - R'|$$

$$\text{SRTT} \leftarrow (1 - \text{alpha}) * \text{SRTT} + \text{alpha} * R'$$

$$\text{RTO} \leftarrow \text{SRTT} + \max(G, K * \text{RTTVAR})$$

Where:

R - first RTT measurement

R' - subsequent RTT measurement

RTTVAR - RTT variance

SRTT - smoothed RTT estimate

RTO - retransmission timeout

G - clock granularity

Recommended values:

alpha=1/8, beta=1/4, K=4

Exponential Back-off

RTO after a timeout:

Recommended value: $q = 2$

$$\text{RTO} \leftarrow q * \text{RTO}$$

This a **congestion control mechanism** since retransmissions are delayed after packet loss detected. The delay is increasing **exponentially** with more packet losses.

TCP Timestamp

▶ Question:

- ACK for what?

▶ RTTM - RTT Management

- TCP option, two 4-byte values
- *TS value (TSval)* - current “timer” value
- *TS echo reply value (TSecr)* - most recently received TSval (only if it acknowledges new data)

```
▼ Options: (24 bytes), Maximum segment size, No-Operation
  > TCP Option - Maximum segment size: 1440 bytes
  > TCP Option - No-Operation (NOP)
  > TCP Option - Window scale: 5 (multiply by 32)
  > TCP Option - No-Operation (NOP)
  > TCP Option - No-Operation (NOP)
  ▼ TCP Option - Timestamps: TSval 819433519, TSecr 0
    Kind: Time Stamp Option (8)
    Length: 10
    Timestamp value: 819433519
    Timestamp echo reply: 0
  > TCP Option - SACK permitted
  > TCP Option - End of Option List (EOL)
```

Wireshark decode of a TCP SYN packet