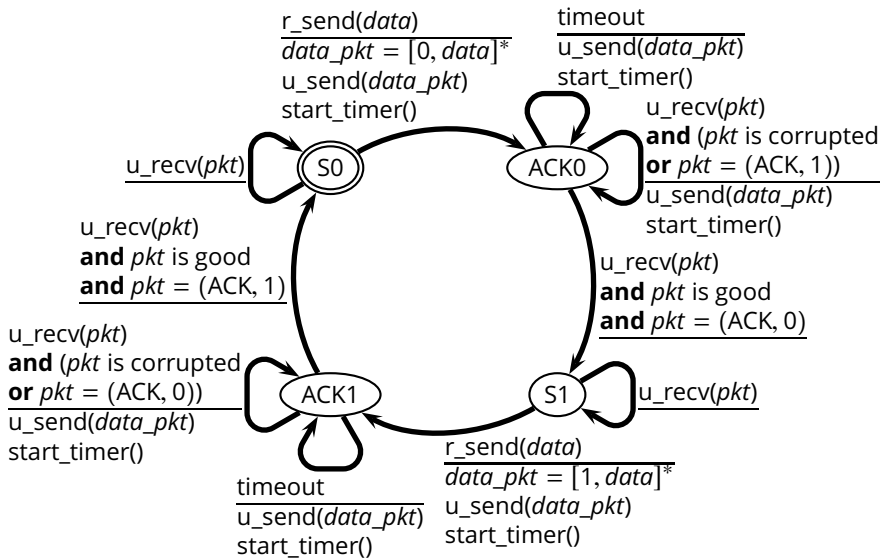# Reliable Data Transfer II

Antonio Carzaniga

Faculty of Informatics
Università della Svizzera italiana

November 10, 2016
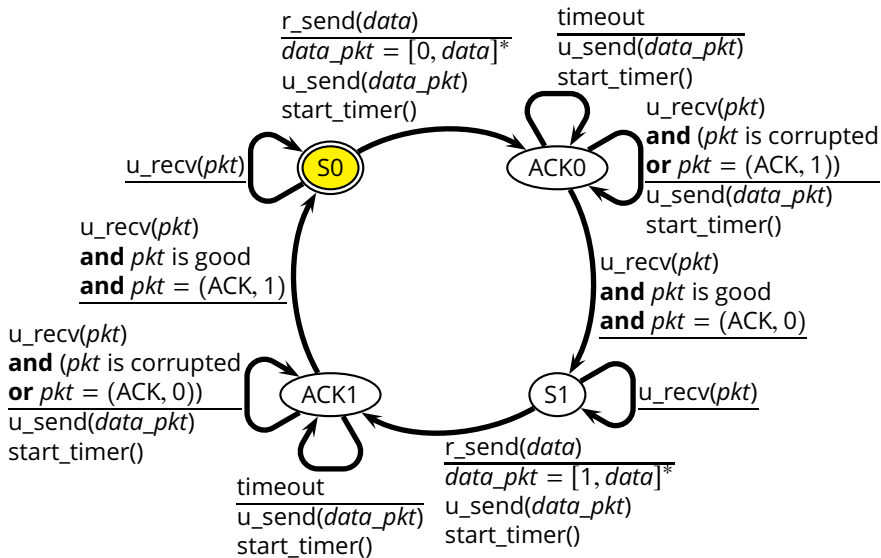
- Performance of the stop-and-wait protocol
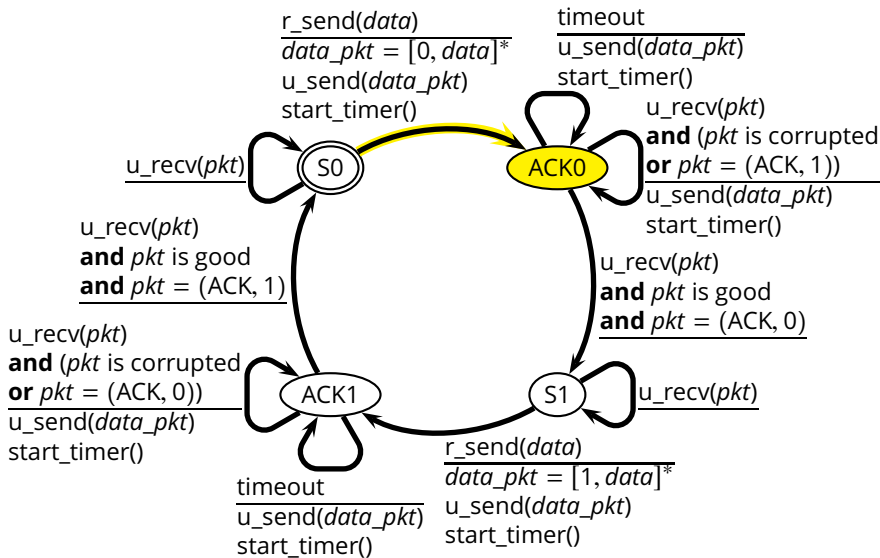
- Go-Back-N

- Selective repeat

Finite state machine diagram with states S0, ACK0, S1, ACK1.

**S0 → ACK0:** r_send(data) / $data\_pkt = [0, data]^*$ / u_send(data_pkt) / start_timer()

**ACK0 self-loop (top):** timeout / u_send(data_pkt) / start_timer()

**ACK0 self-loop (right):** u_recv(pkt) **and** (pkt is corrupted **or** pkt = (ACK, 1)) / u_send(data_pkt) / start_timer()

**S0 self-loop:** u_recv(pkt)

**ACK0 → S1:** u_recv(pkt) **and** pkt is good **and** pkt = (ACK, 0)

**S1 → ACK1:** r_send(data) / $data\_pkt = [1, data]^*$ / u_send(data_pkt) / start_timer()

**S1 self-loop:** u_recv(pkt)

**ACK1 self-loop (bottom):** timeout / u_send(data_pkt) / start_timer()

**ACK1 self-loop (left):** u_recv(pkt) **and** (pkt is corrupted **or** pkt = (ACK, 0)) / u_send(data_pkt) / start_timer()

**ACK1 → S0:** u_recv(pkt) **and** pkt is good **and** pkt = (ACK, 1)

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
u_send(*data_pkt*)
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
u_send(*data_pkt*)
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
u_send(*data_pkt*)
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

# Back to Reliable Data Tranfer

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
u_send(*data_pkt*)
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
u_send(*data_pkt*)
start_timer()

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
u_send(*data_pkt*)
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(*data*)
$\overline{data\_pkt = [0, data]^*}$
u_send(*data_pkt*)
start_timer()

timeout
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 1))
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

u_recv(*pkt*)

S0

ACK0

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 1)

u_recv(*pkt*)
**and** *pkt* is good
**and** *pkt* = (ACK, 0)

u_recv(*pkt*)
**and** (*pkt* is corrupted
**or** *pkt* = (ACK, 0))
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

ACK1

S1

u_recv(*pkt*)

timeout
$\overline{\text{u\_send(}data\_pkt)}$
start_timer()

r_send(*data*)
$\overline{data\_pkt = [1, data]^*}$
u_send(*data_pkt*)
start_timer()

r_send(data)
$\overline{data\_pkt = [0, data]^*}$
u_send(data_pkt)
start_timer()

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(pkt) **and** (pkt is corrupted **or** pkt = (ACK, 1))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

u_recv(pkt)

S0

ACK0

u_recv(pkt) **and** pkt is good **and** pkt = (ACK, 0)

u_recv(pkt) **and** pkt is good **and** pkt = (ACK, 1)

u_recv(pkt) **and** (pkt is corrupted **or** pkt = (ACK, 0))
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

ACK1

S1

u_recv(pkt)

timeout
$\overline{\text{u\_send}(data\_pkt)}$
start_timer()

r_send(data)
$\overline{data\_pkt = [1, data]^*}$
u_send(data_pkt)
start_timer()

sender          receiver

r_send($pkt_1$)

sender          receiver

r_send($pkt_1$)

u_send([$pkt_1$,0])

r_send($pkt_1$)

u_send([$pkt_1$,0])

[$pkt_1$,0]

u_send([ACK,0])

[ACK,0]

r_recv($pkt_1$)

sender

receiver

r_send($pkt_1$)

u_send([$pkt_1$,0])

[$pkt_1$,0]

u_send([ACK,0])

r_recv($pkt_1$)

[ACK,0]

r_send($pkt_2$)

u_send([$pkt_2$,1])

[$pkt_2$,1]

sender

receiver

sender                                    receiver

r_send($pkt_1$)
    u_send([$pkt_1$,0])

[$pkt_1$,0]

             u_send([ACK,0])
                   r_recv($pkt_1$)

[ACK,0]

r_send($pkt_2$)
    u_send([$pkt_2$,1])

[$pkt_2$,1]

             u_send([ACK,1])
                   r_recv($pkt_2$)

sender　　　　receiver

$[pkt_1,0]$

$[ACK,0]$

$[pkt_2,1]$

$[ACK,1]$

*utilization factor*

$$U = \frac{\ell_{pkt}/R}{2d + \ell_{pkt}/R}$$

- How do we achieve a better *utilization factor*?

■ How do we achieve a better *utilization factor*?

# Improving Network Usage

- How do we achieve a better *utilization factor*?

# Improving Network Usage

- How do we achieve a better *utilization factor*?

# Improving Network Usage

■ How do we achieve a better *utilization factor*?

- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement
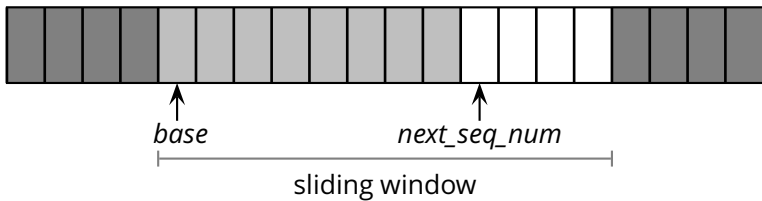
- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to $W$ unacknowledged packets in the pipeline
  - ▶ the sender's state machine gets very complex
  - ▶ we represent the sender's state with its queue of acknowledgements
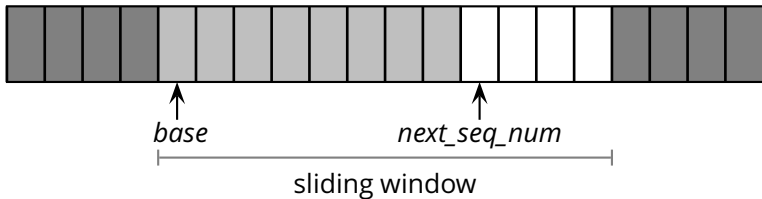
- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to $W$ unacknowledged packets in the pipeline
  - ▶ the sender's state machine gets very complex
  - ▶ we represent the sender's state with its queue of acknowledgements

- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to *W* unacknowledged packets in the pipeline
  - ▸ the sender's state machine gets very complex
  - ▸ we represent the sender's state with its queue of acknowledgements

acknowledged       pending       available   unavailable

- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to $W$ unacknowledged packets in the pipeline
  - the sender's state machine gets very complex
  - we represent the sender's state with its queue of acknowledgements

acknowledged          pending          available   unavailable

first pending
acknowledgement
*(base)*

- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to *W* unacknowledged packets in the pipeline
  - the sender's state machine gets very complex
  - we represent the sender's state with its queue of acknowledgements



acknowledged      pending      available   unavailable

first pending
acknowledgement
*(base)*

next available
sequence number
*(next_seq_num)*

- **Idea:** the sender transmits multiple packets without waiting for an acknowledgement

- Sender has up to *W* unacknowledged packets in the pipeline
  - ▸ the sender's state machine gets very complex
  - ▸ we represent the sender's state with its queue of acknowledgements



acknowledged      pending      available    unavailable

first pending
acknowledgement
*(base)*

next available
sequence number
*(next_seq_num)*

window size (*W*)

base        next_seq_num

sliding window

# Sliding Window Protocol: Sender



*base*  *next_seq_num*

sliding window

- r_send($pkt_1$)

# Sliding Window Protocol: Sender



*base*  *next_seq_num*

sliding window

- r_send($pkt_1$)
  - u_send([$pkt_1$,*next_seq_num*])

# Sliding Window Protocol: Sender



*base*  *next_seq_num*

sliding window

- r_send($pkt_1$)
    - u_send([$pkt_1$,*next_seq_num*])
    - *next_seq_num* = *next_seq_num* + 1

# Sliding Window Protocol: Sender



*A*

*base*          *next_seq_num*

sliding window

■ r_send($pkt_1$)

   ► u_send([$pkt_1$,*next_seq_num*])

   ► *next_seq_num* = *next_seq_num* + 1

■ u_recv([ACK,*A*])

# Sliding Window Protocol: Sender



- r_send($pkt_1$)
  - u_send([$pkt_1$,$next\_seq\_num$])
  - $next\_seq\_num = next\_seq\_num + 1$

- u_recv([ACK,$A$])
  - $base = A + 1$

# Sliding Window Protocol: Sender



- r_send($pkt_1$)
    - u_send([$pkt_1$,$next\_seq\_num$])
    - $next\_seq\_num = next\_seq\_num + 1$

- u_recv([ACK,$A$])
    - $base = A + 1$
    - notice that acknewledgements are "cumulative"

# Sliding Window Protocol: Sender

- The sender remembers the first sequence number that has not yet been acknowledged
    - or the highest acknowledged sequence number

- The sender remembers the first available sequence number
    - or the highest used sequence number (i.e., sent to the receiver)

- The sender responds to three types of events

# Sliding Window Protocol: Sender

- The sender remembers the first sequence number that has not yet been acknowledged
  - or the highest acknowledged sequence number

- The sender remembers the first available sequence number
  - or the highest used sequence number (i.e., sent to the receiver)

- The sender responds to three types of events
  - *r_send():* invocation from the application layer: send more data if a sequence number is available

# Sliding Window Protocol: Sender

- The sender remembers the first sequence number that has not yet been acknowledged
  - or the highest acknowledged sequence number

- The sender remembers the first available sequence number
  - or the highest used sequence number (i.e., sent to the receiver)

- The sender responds to three types of events
  - *r_send():* invocation from the application layer: send more data if a sequence number is available
  - *ACK:* receipt of an acknowledgement: shift the window (it's a "cumulative" ACK)

# Sliding Window Protocol: Sender

- The sender remembers the first sequence number that has not yet been acknowledged
  - ▶ or the highest acknowledged sequence number

- The sender remembers the first available sequence number
  - ▶ or the highest used sequence number (i.e., sent to the receiver)

- The sender responds to three types of events
  - ▶ *r_send():* invocation from the application layer: send more data if a sequence number is available
  - ▶ *ACK:* receipt of an acknowledgement: shift the window (it's a "cumulative" ACK)
  - ▶ *timeout:* "Go-Back-N." I.e., resend all the packets that have been sent but not acknowledged

■ *init*

---

$base = 1$

$next\_seq\_num = 1$

# Sliding Window Protocol: Sender

- *init*

  $base = 1$

  $next\_seq\_num = 1$

- r_send(*data*)

  **if** $next\_seq\_num < base + W$:

      $pkt[next\_seq\_num] = [next\_seq\_num, data]^*$

      u_send($pkt[next\_seq\_num]$)

      **if** $next\_seq\_num == base$:

          start_timer()

      $next\_seq\_num = next\_seq\_num + 1$

  **else**:

      refuse_data(data)    *// block the sender*

- u_recv(*pkt*) **and** *pkt* is corrupted

- u_recv(*pkt*) **and** *pkt* is corrupted

- u_recv(ACK,*ack_num*)

  *base* = *ack_num* + 1     *// resume the sender*
  **if** *next_seq_num == base*:
      stop_timer()
  **else**:
      start_timer()

- u_recv(*pkt*) **and** *pkt* is corrupted

- u_recv(ACK,*ack_num*)

  *base* = *ack_num* + 1     *// resume the sender*
  **if** *next_seq_num* == *base*:
      stop_timer()
  **else**:
      start_timer()

- timeout

  start_timer()
  **foreach** *i* **in** *base* . . . *next_seq_num* − 1:
      u_send(*pkt*[*i*])

- Simple: as in the stop-and-wait case, the receiver maintains a counter representing the *expected sequence number*

# Sliding Window Protocol: Receiver

- Simple: as in the stop-and-wait case, the receiver maintains a counter representing the *expected sequence number*

- The receiver waits for a (good) data packet with the expected sequence number

# Sliding Window Protocol: Receiver

- Simple: as in the stop-and-wait case, the receiver maintains a counter representing the *expected sequence number*

- The receiver waits for a (good) data packet with the expected sequence number

  - acknowledges the expected sequence number

# **Sliding Window Protocol: Receiver**

- Simple: as in the stop-and-wait case, the receiver maintains a counter representing the *expected sequence number*

- The receiver waits for a (good) data packet with the expected sequence number
  - ► acknowledges the expected sequence number
  - ► delivers the data to the application

- *init*

  $\overline{expected\_seq\_num = 1}$

  $ackpkt = [ACK, 0]^*$

## Sliding Window Protocol: Receiver

- *init*

    $\overline{expected\_seq\_num = 1}$

    $ackpkt = [ACK, 0]^*$

- u_recv([*data*, *seq_num*]) **and** good

    **and** *seq_num = expected_seq_num*

    $\overline{\text{r\_recv}(data)}$

    $ackpkt = [ACK, expected\_seq\_num]^*$

    $expected\_seq\_num = expected\_seq\_num + 1$

    u_send(*ackpkt*)

## Sliding Window Protocol: Receiver

- *init*

  $expected\_seq\_num = 1$

  $ackpkt = [ACK, 0]^*$

- u_recv([*data*, *seq_num*]) **and** good

  **and** *seq_num = expected_seq_num*

  r_recv(*data*)

  $ackpkt = [ACK, expected\_seq\_num]^*$

  $expected\_seq\_num = expected\_seq\_num + 1$

  u_send(*ackpkt*)

- u_recv([*data*, *seq_num*])

  **and** (corrupted **or** $seq\_num \neq expected\_seq\_num$)

  u_send(*ackpkt*)

- Concepts

- Concepts
  - *sequence numbers*

- Concepts
  - *sequence numbers*
  - *sliding window*

- Concepts
  - *sequence numbers*
  - *sliding window*
  - *cumulative acknowledgements*

■ Concepts

- ▸ *sequence numbers*
- ▸ *sliding window*
- ▸ *cumulative acknowledgements*
- ▸ *checksums*, *timeouts*, and *sender-initiated retransmission*

- Concepts
  - ▸ *sequence numbers*
  - ▸ *sliding window*
  - ▸ *cumulative acknowledgements*
  - ▸ *checksums*, *timeouts*, and *sender-initiated retransmission*

- Advantages: *simple*

- Concepts
  - ▸ *sequence numbers*
  - ▸ *sliding window*
  - ▸ *cumulative acknowledgements*
  - ▸ *checksums*, *timeouts*, and *sender-initiated retransmission*

- Advantages: *simple*
  - ▸ the sender maintains ***two counters*** and ***one timer***
  - ▸ the receiver maintains ***one counter***

- Concepts
  - *sequence numbers*
  - *sliding window*
  - *cumulative acknowledgements*
  - *checksums*, *timeouts*, and *sender-initiated retransmission*

- Advantages: *simple*
  - the sender maintains **two counters** and **one timer**
  - the receiver maintains **one counter**

- Disadvantages: *not optimal, not adaptive*

- Concepts
  - *sequence numbers*
  - *sliding window*
  - *cumulative acknowledgements*
  - *checksums*, *timeouts*, and *sender-initiated retransmission*

- Advantages: *simple*
  - the sender maintains ***two counters*** and ***one timer***
  - the receiver maintains ***one counter***

- Disadvantages: *not optimal*, *not adaptive*
  - the sender can fill the window without filling the pipeline

- Concepts
  - *sequence numbers*
  - *sliding window*
  - *cumulative acknowledgements*
  - *checksums*, *timeouts*, and *sender-initiated retransmission*

- Advantages: *simple*
  - the sender maintains **two counters** and **one timer**
  - the receiver maintains **one counter**

- Disadvantages: *not optimal, not adaptive*
  - the sender can fill the window without filling the pipeline
  - the receiver may buffer out-of-order packets…

- What is a good value for $W$?

- What is a good value for *W*?
  - *W* that achieves the *maximum utilization* of the connection

- What is a good value for *W*?
  - ▸ *W* that achieves the *maximum utilization* of the connection

$$\ell = stream$$
$$d = 500ms$$
$$R = 1Mb/s$$
$$W = ?$$

- What is a good value for *W*?

  ▸ *W* that achieves the *maximum utilization* of the connection

$$\ell = stream$$
$$d = 500ms$$
$$R = 1Mb/s$$
$$W = ?$$

- The problem may seem a bit underspecified. What is the (average) packet size?

$$\ell_{pkt} = 1Kb$$
$$d = 500ms$$
$$R = 1Mb/s$$
$$W = \frac{2d \times R}{\ell_{pkt}} = 1000$$

- The RTT–throughput product ($2d \times R$) is the crucial factor

- The RTT–throughput product ($2d \times R$) is the crucial factor

  - $W \times \ell_{pkt} \leq 2d \times R$

    - why $W \times \ell_{pkt} > 2d \times R$ doesn't make much sense?

- The RTT–throughput product ($2d \times R$) is the crucial factor

  - $W \times \ell_{pkt} \leq 2d \times R$

    - why $W \times \ell_{pkt} > 2d \times R$ doesn't make much sense?

  - maximum channel utilization when $W \times \ell_{pkt} = 2d \times R$

  - $2d \times R$ can be thought of as the *capacity* of a connection

- Let's consider a fully utilized connection

■ Let's consider a fully utilized connection

$$
\begin{aligned}
\ell_{pkt} &= 1Kb \\
d &= 500ms \\
R &= 1Mb/s \\
W &= \frac{R \times d}{\ell_{pkt}} = 1000
\end{aligned}
$$

- Let's consider a fully utilized connection

$$
\begin{aligned}
\ell_{pkt} &= 1\,Kb \\
d &= 500\,ms \\
R &= 1\,Mb/s \\
W &= \frac{R \times d}{\ell_{pkt}} = 1000
\end{aligned}
$$

- What happens if the first packet (or acknowledgement) is lost?

- Let's consider a fully utilized connection

$$\ell_{pkt} = 1Kb$$
$$d = 500ms$$
$$R = 1Mb/s$$
$$W = \frac{R \times d}{\ell_{pkt}} = 1000$$

- What happens if the first packet (or acknowledgement) is lost?

- Sender retransmits the entire content of its buffers

- Let's consider a fully utilized connection

$$\ell_{pkt} = 1Kb$$
$$d = 500ms$$
$$R = 1Mb/s$$
$$W = \frac{R \times d}{\ell_{pkt}} = 1000$$

- What happens if the first packet (or acknowledgement) is lost?

- Sender retransmits the entire content of its buffers

  ▸ $W \times \ell_{pkt} = 2d \times R = 1Mb$
  ▸ retransmitting 1Mb to recover 1Kb worth of data isn't exactly the best solution. Not to mention conjestions...

- Let's consider a fully utilized connection

$$
\begin{aligned}
\ell_{pkt} &= 1Kb \\
d &= 500ms \\
R &= 1Mb/s \\
W &= \frac{R \times d}{\ell_{pkt}} = 1000
\end{aligned}
$$

- What happens if the first packet (or acknowledgement) is lost?

- Sender retransmits the entire content of its buffers

  ▸ $W \times \ell_{pkt} = 2d \times R = 1Mb$
  ▸ retransmitting 1Mb to recover 1Kb worth of data isn't exactly the best solution. Not to mention conjestions...

- Is there a better way to deal with retransmissions?

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

# Selective Repeat

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

  - sender maintains a vector of acknowledgement flags

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

  - ▸ sender maintains a vector of acknowledgement flags
  - ▸ receiver maintains a vector of acknowledged falgs

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

  - ▶ sender maintains a vector of acknowledgement flags
  - ▶ receiver maintains a vector of acknowledged falgs
  - ▶ in fact, receiver maintains a buffer of out-of-order packets

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

  - ▶ sender maintains a vector of acknowledgement flags
  - ▶ receiver maintains a vector of acknowledged falgs
  - ▶ in fact, receiver maintains a buffer of out-of-order packets
  - ▶ sender maintains a timer for each pending packet

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

  - ‣ sender maintains a vector of acknowledgement flags
  - ‣ receiver maintains a vector of acknowledged falgs
  - ‣ in fact, receiver maintains a buffer of out-of-order packets
  - ‣ sender maintains a timer for each pending packet
  - ‣ sender resends a packet when its timer expires

- **Idea:** have the sender retransmit only those packets that it suspects were lost or corrupted

    - sender maintains a vector of acknowledgement flags
    - receiver maintains a vector of acknowledged falgs
    - in fact, receiver maintains a buffer of out-of-order packets
    - sender maintains a timer for each pending packet
    - sender resends a packet when its timer expires
    - sender slides the window when the lowest pending sequence number is acknowledged

*base*       *next_seq_num*

sliding window

*base*  *next_seq_num*

sliding window

- r_send($pkt_1$)

*base*      *next_seq_num*

sliding window

- r_send($pkt_1$)
  - u_send([$pkt_1$,*next_seq_num*])
  - start_timer(*next_seq_num*)

*base*     *next_seq_num*

sliding window

- ■ r_send($pkt_1$)
  - ▸ u_send([$pkt_1$,*next_seq_num*])
  - ▸ start_timer(*next_seq_num*)
  - ▸ *next_seq_num* = *next_seq_num* + 1

# Selective Repeat: Sender



- r_send($pkt_1$)
  - u_send([$pkt_1$,*next_seq_num*])
  - start_timer(*next_seq_num*)
  - *next_seq_num* = *next_seq_num* + 1

- u_recv([ACK,*A*])

# Selective Repeat: Sender



*A*

*base*       *next_seq_num*

sliding window

- r_send(*pkt*$_1$)
  - u_send([*pkt*$_1$,*next_seq_num*])
  - start_timer(*next_seq_num*)
  - *next_seq_num* = *next_seq_num* + 1

- u_recv([ACK,*A*])
  - *acks*[*A*] = 1     *// remember that A was ACK'd*

# Selective Repeat: Sender



*A*

*base*          *next_seq_num*

sliding window

- r_send($pkt_1$)

  - u_send([$pkt_1$,*next_seq_num*])
  - start_timer(*next_seq_num*)
  - *next_seq_num* = *next_seq_num* + 1

- u_recv([ACK,*A*])

  - *acks*[*A*] = 1          *// remember that A was ACK'd*
  - acknewledgements are no longer "cumulative"

received          acceptable          not usable

*rcv_base*

sliding window

- u_recv([$pkt_1, X_1$]) **and** $rcv\_base \leq X_1 < rcv\_base + W$

received | $X_1$ | acceptable | not usable

rcv_base

sliding window

- u_recv([$pkt_1$,$X_1$]) **and** $rcv\_base \leq X_1 < rcv\_base + W$
  - $buffer[X_1] = pkt_1$
  - u_send([$ACK, X_1$]*)    *// no longer a "cumulative" ACK*

- u_recv($[pkt_2, X_2]$) **and** $rcv\_base \leq X_2 < rcv\_base + W$
  - $buffer[X_2] = pkt_2$
  - u_send($[ACK, X_2]^*$)

received      $X_2$      acceptable      not usable

rcv_base

sliding window

- u_recv($[pkt_2, X_2]$) **and** $rcv\_base \leq X_2 < rcv\_base + W$
  - $buffer[X_2] = pkt_2$
  - u_send($[ACK, X_2]^*$)
  - **if** $X_2$ == $rcv\_base$:

## Selective Repeat: Receiver



$X_2$  $B$

received   acceptable   not usable

rcv_base

sliding window

- u_recv($[pkt_2, X_2]$) **and** $rcv\_base \leq X_2 < rcv\_base + W$
  - $buffer[X_2] = pkt_2$
  - u_send($[ACK, X_2]^*$)
  - **if** $X_2$ == $rcv\_base$:
        $B = first\_missing\_seq\_num()$
        **foreach** $i$ **in** $rcv\_base \ldots B - 1$:
            r_recv($buffer[i]$)

## Selective Repeat: Receiver



- u_recv([$pkt_2$,$X_2$]) **and** $rcv\_base \leq X_2 < rcv\_base + W$
    - $buffer[X_2] = pkt_2$
    - u_send([$ACK$,$X_2$]*)
    - **if** $X_2$ == $rcv\_base$:
        $B = first\_missing\_seq\_num()$
        **foreach** $i$ **in** $rcv\_base \ldots B - 1$:
            r_recv($buffer[i]$)
        $rcv\_base = B$

base

next_seq_num

sliding window

- Timeout for sequence number *T*

- Timeout for sequence number $T$
  - u_send($[pkt[T], T]^*$)

*base*      *next_seq_num*

sliding window

- u_recv([ACK,*A*])

- u_recv([ACK,$A$])
  - $acks[A] = 1$

- u_recv([ACK,*A*])
  - *acks*[*A*] = 1
  - **if** *A* == *base*:

*base*        *next_seq_num*

sliding window

■ u_recv([ACK,*A*])
- ▸ *acks*[*A*] = 1
- ▸ **if** *A* == *base*:
    *base* = *first_missing_ack_num*()