# Computer  Networking : Principles, Protocols and Practice

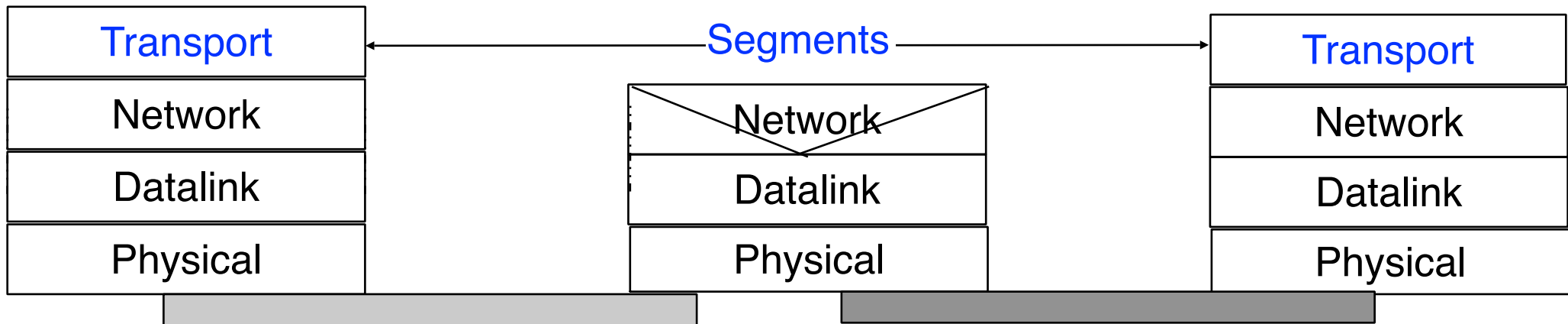## Part 3 : Transport Layer

Olivier Bonaventure
http://inl.info.ucl.ac.be/

# Module 3 : Transport Layer

⟶ **Basics**

Building a reliable transport layer
   Reliable data transmission
   Connection establishment
   Connection release

UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

# The transport layer

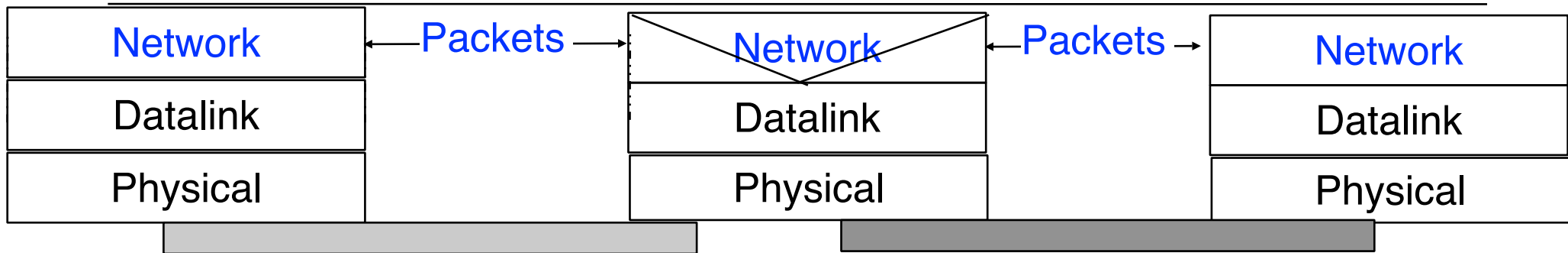| Transport | | Segments | | Transport |
|---|---|---|---|---|
| Network | | Network | | Network |
| Datalink | | Datalink | | Datalink |
| Physical | | Physical | | Physical |

# Goals
Improves the service provided by the network layer to allow it to be useable by applications
- reliability
- multiplexing

# Transport layer services
Unreliable connectionless service
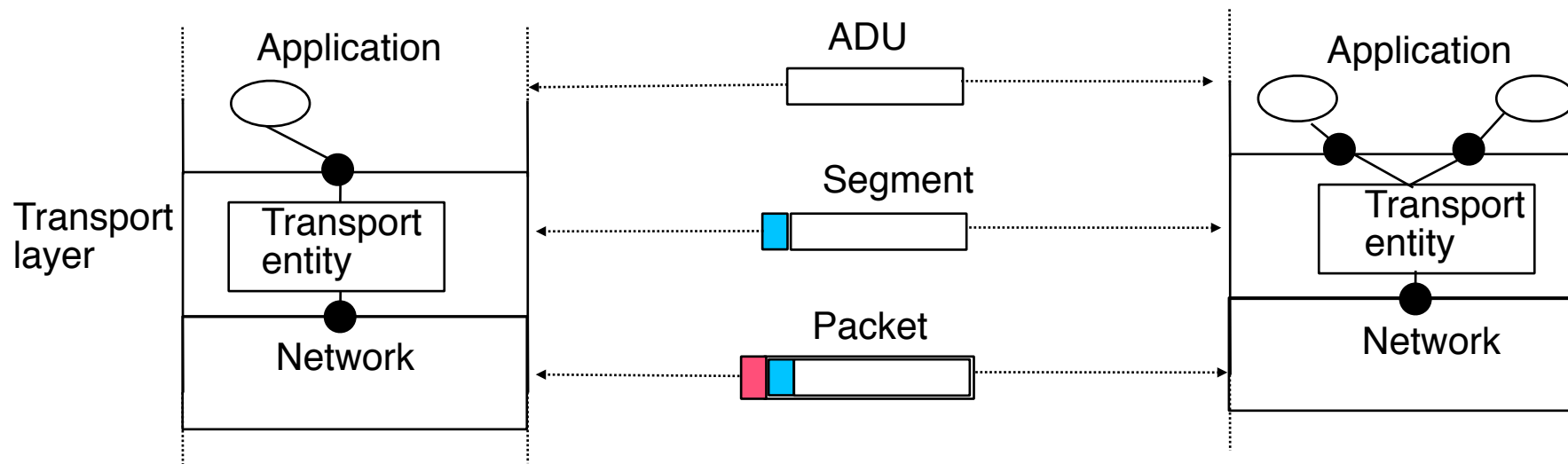Reliable connection-oriented service

# The network layer

| Network | | Network | | Network |
|---------|---|---------|---|---------|
| Datalink | ← Packets → | Datalink | ← Packets → | Datalink |
| Physical | | Physical | | Physical |

# Network layer service in Internet
## Unreliable connectionless service
Packets can be lost
Packets can suffer from transmission errors
Packet ordering is not preserved
Packet can be duplicated
Packet size is limited to about 64 KBytes

How to build a service useable by applications ?

# The transport layer

## Problems to be solved by transport layer

Transport layer must allow two <span style="color:red">applications</span> to exchange information
> This requires a method to identify the applications

The transport layer service must be useable by applications
> detection of transmission errors
> correction of transmission errors
> recovery from packet losses and packet duplications
> different types of services
>> connectionless
>> connection-oriented
>> request-response

# The transport layer (2)

## Internal organisation

The transport layer uses the service provide by the network layer

Two transport layer entities exchanges segments

# Module 3 : Transport layer

Basics

<span style="color:red">Building a reliable transport layer</span>
→     <span style="color:red">Reliable data transmission</span>
    Connection establishment
    Connection release

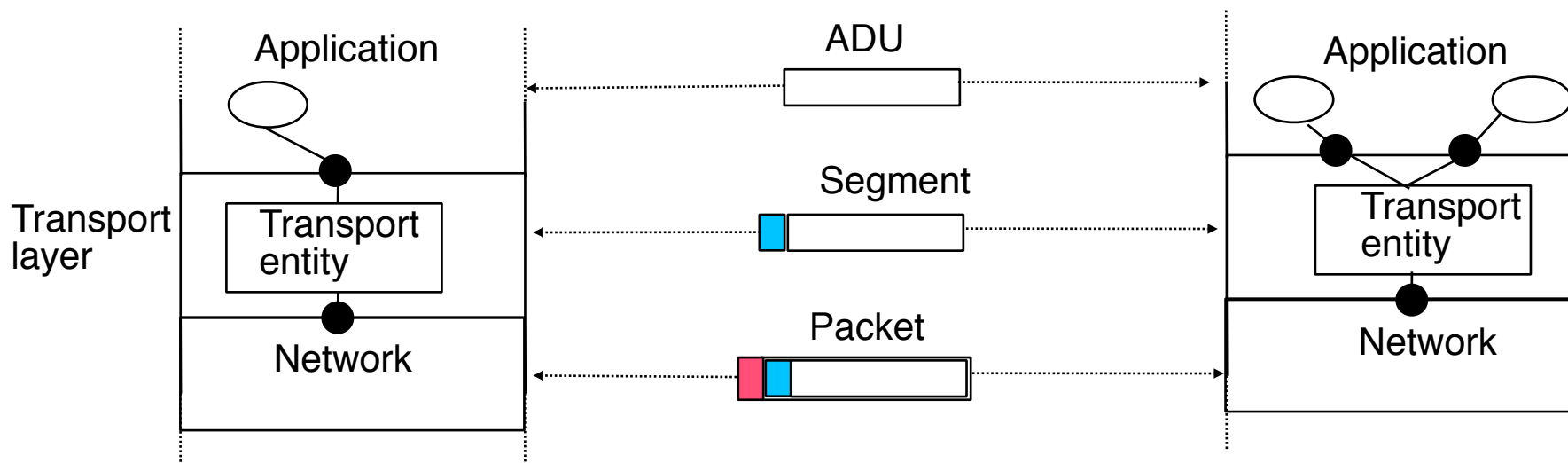UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

# Transport layer protocols

## How can we provide a reliable service in the transport layer

Hypotheses
1. The application sends small SDUs
2. The network layer provides a perfect service
    1. There are no transmission errors inside the packets
    2. No packet is lost
    3. There is no packet reordering
    4. There are no duplications of packets
3. Data transmission is unidirectional

# Transport layer protocols (2)

## Reference environment



## Notations

`data.req` and `data.ind` primitives for application/ transport interactions

`recv()` and `send()` for interactions between transport entity and network layer
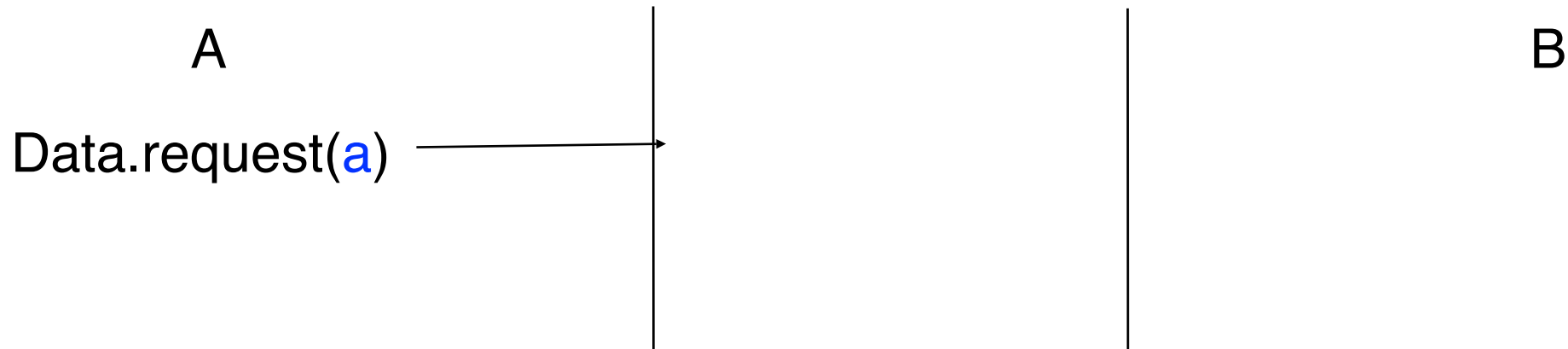
# Protocol 1 : Basics

A                                                                    B

**Principle**

Upon reception of `data.request(SDU)`, the transport entity sends a segment containing this SDU through the network layer (`send(SDU)`)

Upon reception of the contents of one packet from the network layer (`recv(SDU)`), transport entity delivers the SDU found in the packet to its user by using `data.ind(SDU)`
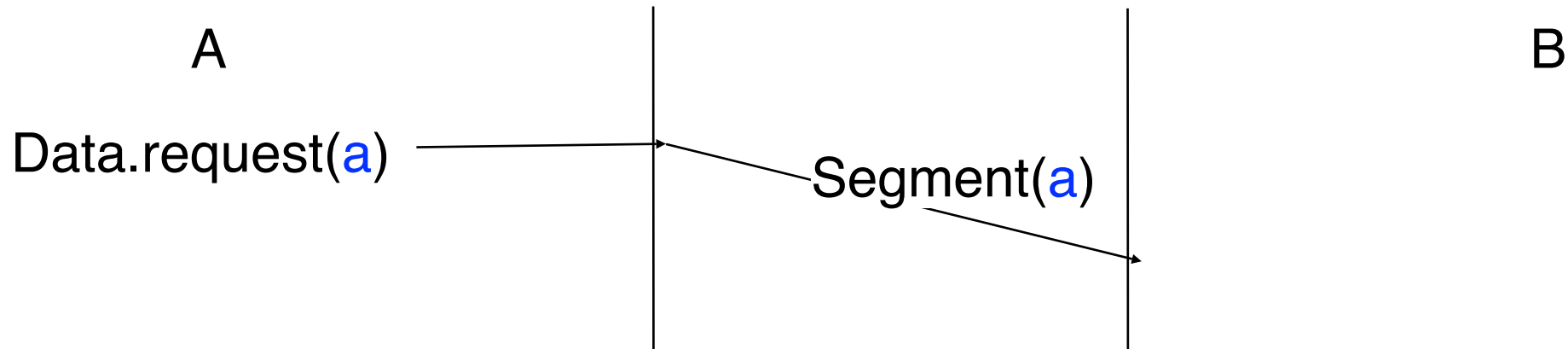
# Protocol 1 : Basics

A                                               B

Data.request(a) ⟶

Principle

    Upon reception of `data.request(SDU)`, the transport entity sends a segment containing this SDU through the network layer (`send(SDU)`)

    Upon reception of the contents of one packet from the network layer (`recv(SDU)`), transport entity delivers the SDU found in the packet to its user by using `data.ind(SDU)`
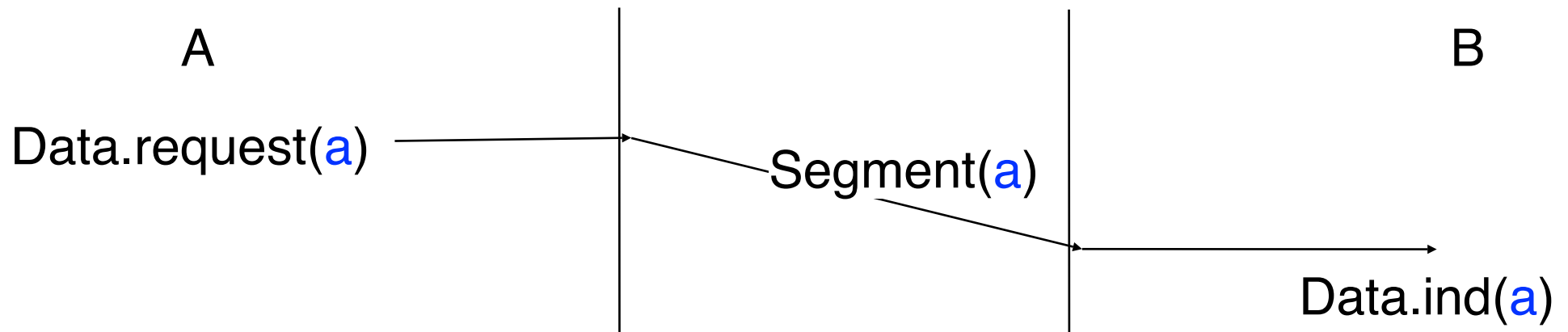
# Protocol 1 : Basics



**Principle**

Upon reception of `data.request(SDU)`, the transport entity sends a segment containing this SDU through the network layer (`send(SDU)`)

Upon reception of the contents of one packet from the network layer (`recv(SDU)`), transport entity delivers the SDU found in the packet to its user by using `data.ind(SDU)`
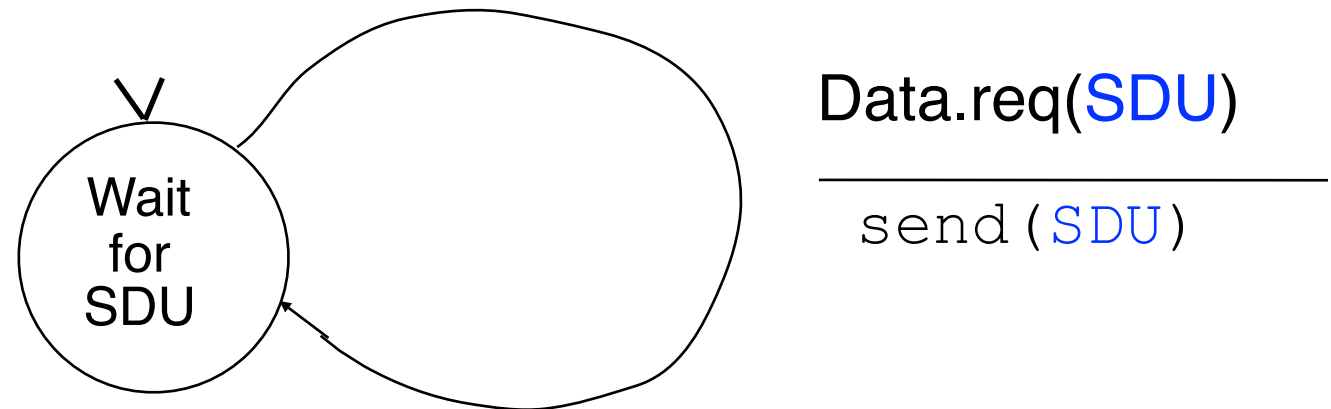
# Protocol 1 : Basics

A

B

Data.request(a) → Segment(a) →

Data.ind(a)

## Principle

Upon reception of `data.request(SDU)`, the transport entity sends a segment containing this SDU through the network layer (`send(SDU)`)

Upon reception of the contents of one packet from the network layer (`recv(SDU)`), transport entity delivers the SDU found in the packet to its user by using `data.ind(SDU)`

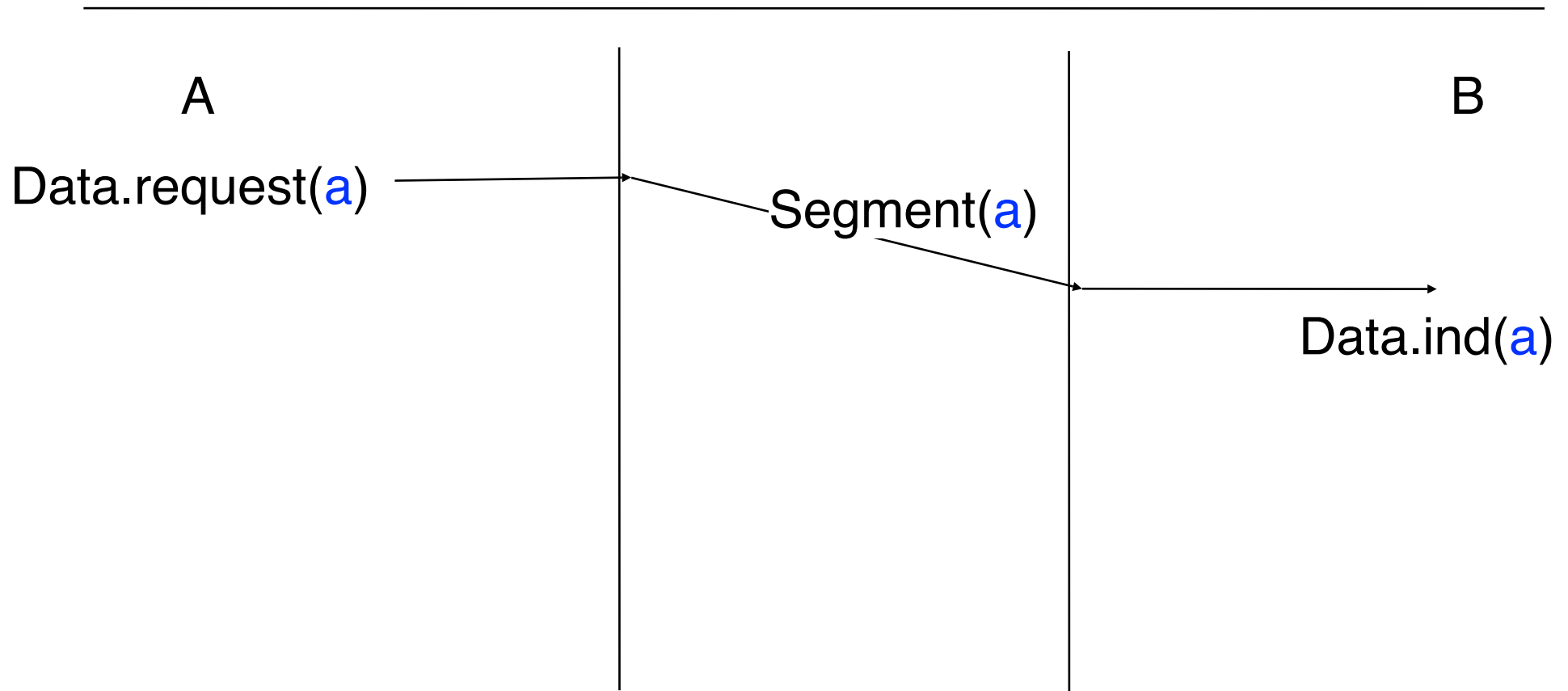# Protocol 1 as a FSM



Sender

Receiver

# Protocol 1 : Example

A

B

## Issue

What happens if the receiver is much slower than the sender ?

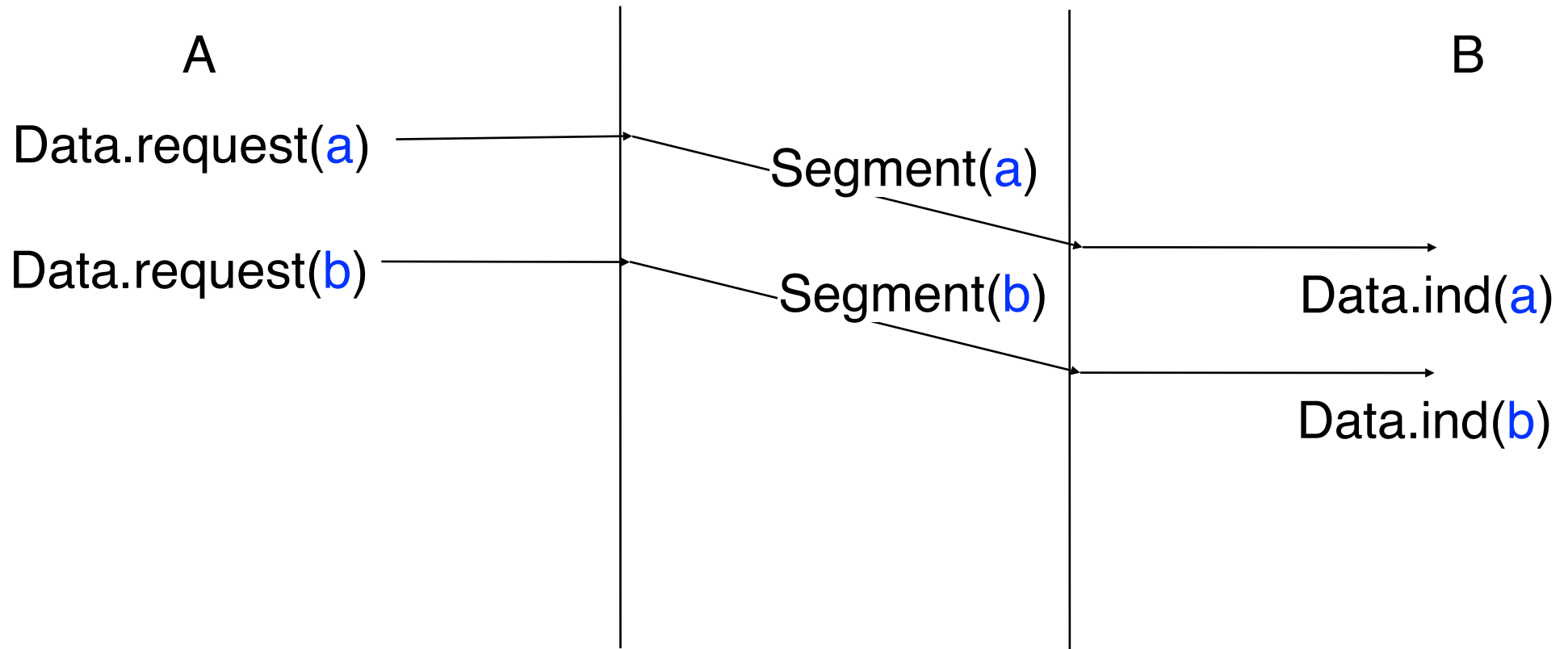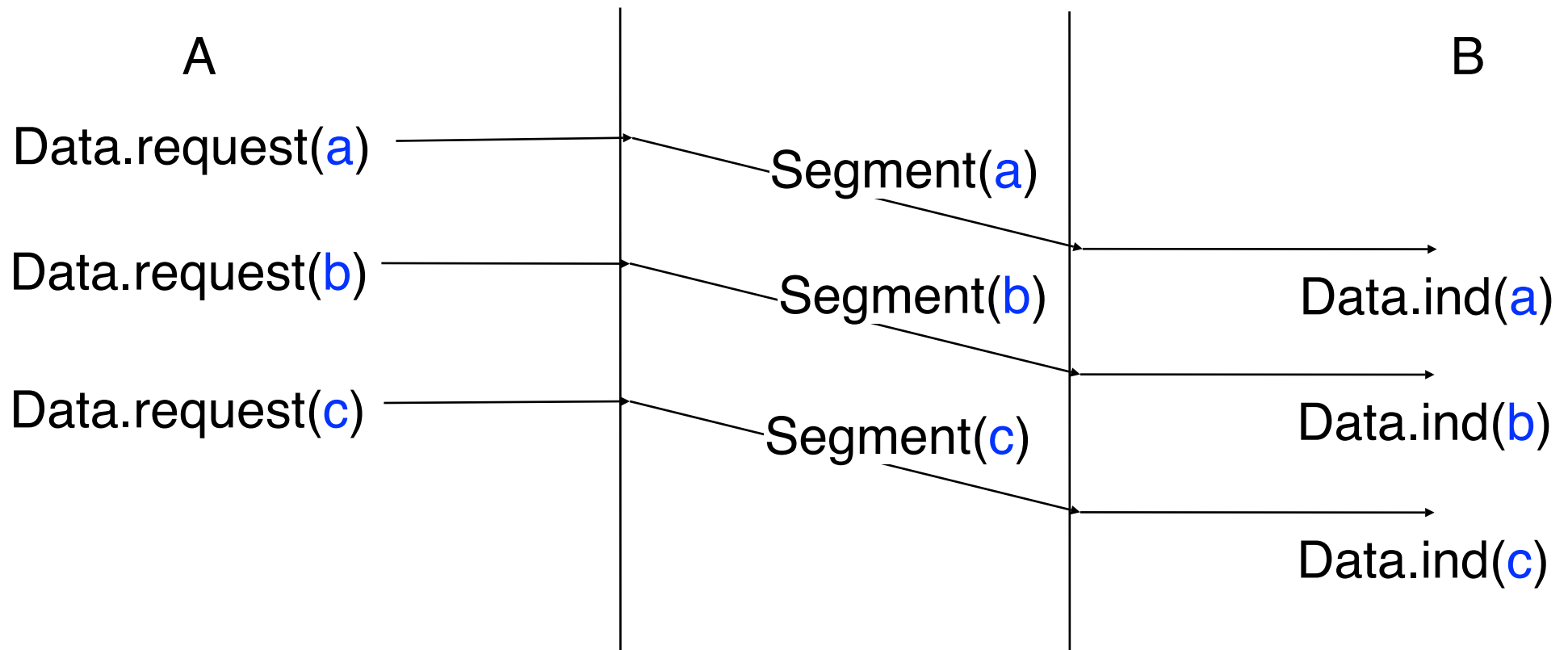e.g. receiver can process one segment per second while sender is producing 10 segments per second ?

# Protocol 1 : Example

A                                                                                        B

Data.request(a) ———————————→
                                         Segment(a)
                                                    ————————————————→
                                                                        Data.ind(a)

## Issue

What happens if the receiver is much slower than the sender ?

e.g. receiver can process one segment per second while sender is producing 10 segments per second ?

# Protocol 1 : Example

A                                               B

Data.request(a)

Segment(a)

Data.request(b)

Segment(b)

Data.ind(a)

Data.ind(b)

## Issue

What happens if the receiver is much slower than the sender ?

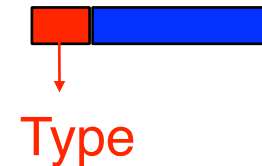e.g. receiver can process one segment per second while sender is producing 10 segments per second ?

# Protocol 1 : Example



A

Data.request(a)

Segment(a)

Data.request(b)

Segment(b)

Data.ind(a)

Data.request(c)

Segment(c)

Data.ind(b)

Data.ind(c)

B

## Issue

What happens if the receiver is much slower than the sender ?

e.g. receiver can process one segment per second while sender is producing 10 segments per second ?

# Protocol 2

## Principle
Use a control segment (OK) that is sent by the receiver after having processed the received segment

creates a feedback loop between sender and receiver

## Consequences
Two types of segments
- Data segment containing on SDU
  - Notation : D(SDU)
- Control segment
  - Notation : C(OK)

## Segment format
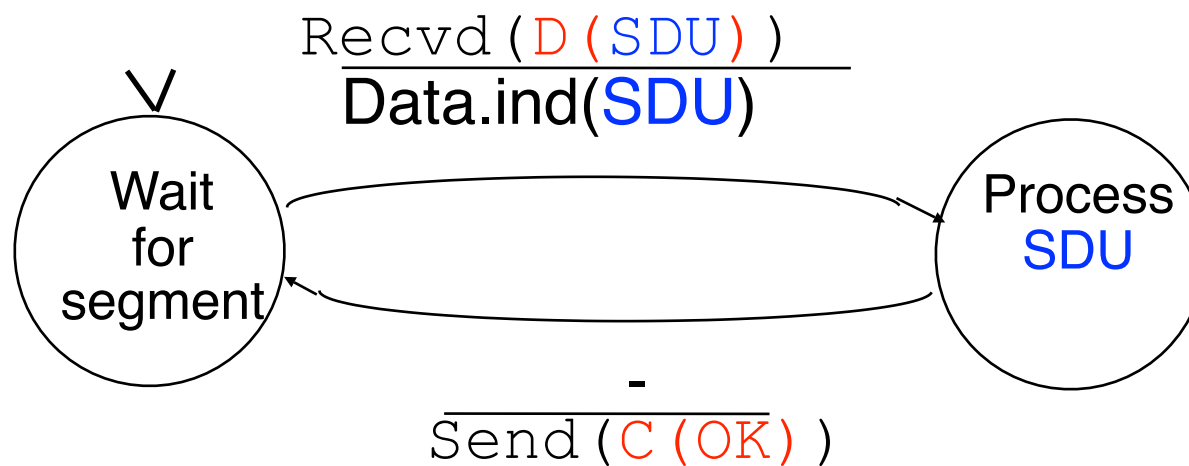At least one bit in the segment header is used to indicate the type of segment

Type

# Protocol 2 (cont.)



Sender

$$\frac{\text{Data.req}(SDU)}{\text{send}(D(SDU))}$$

Wait for SDU → Wait for OK

$$\frac{\text{Recvd}(C(OK))}{-}$$

Receiver

$$\frac{\text{Recvd}(D(SDU))}{\text{Data.ind}(SDU)}$$

Wait for segment → Process SDU
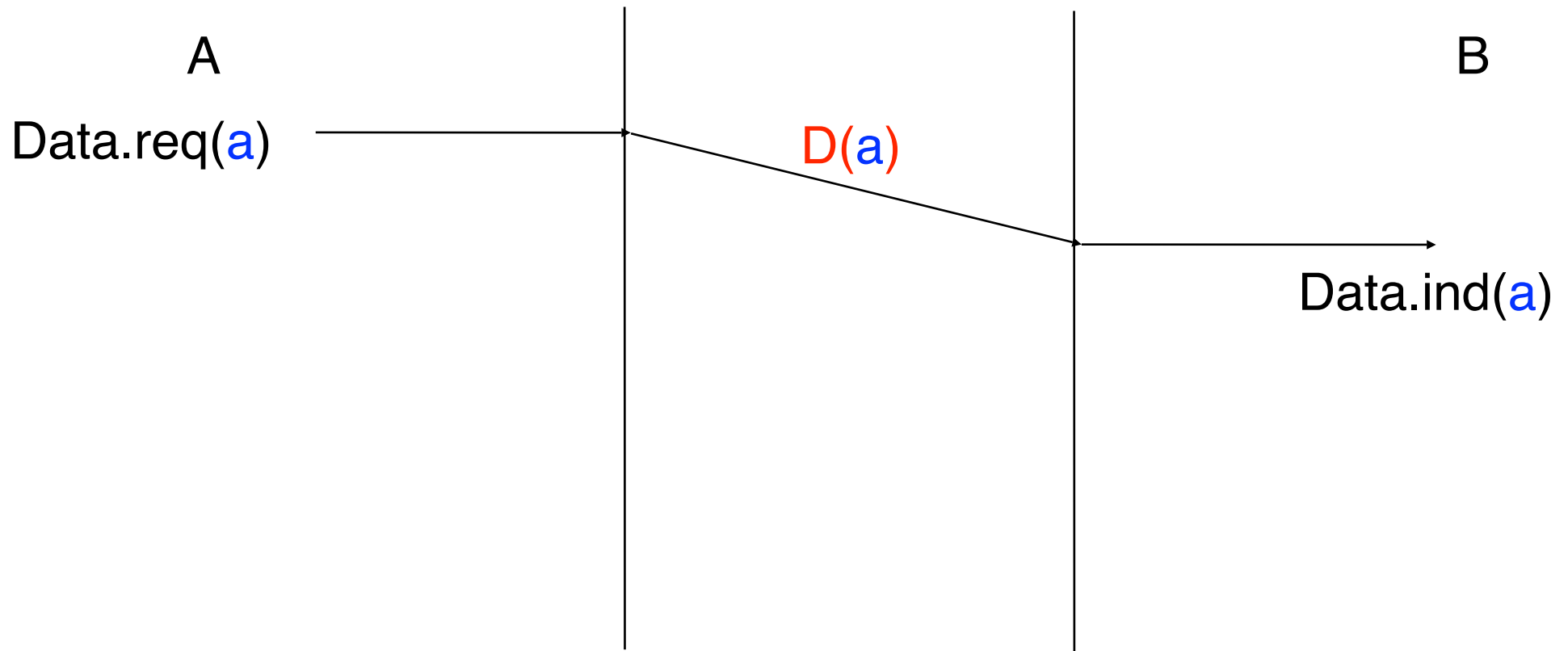
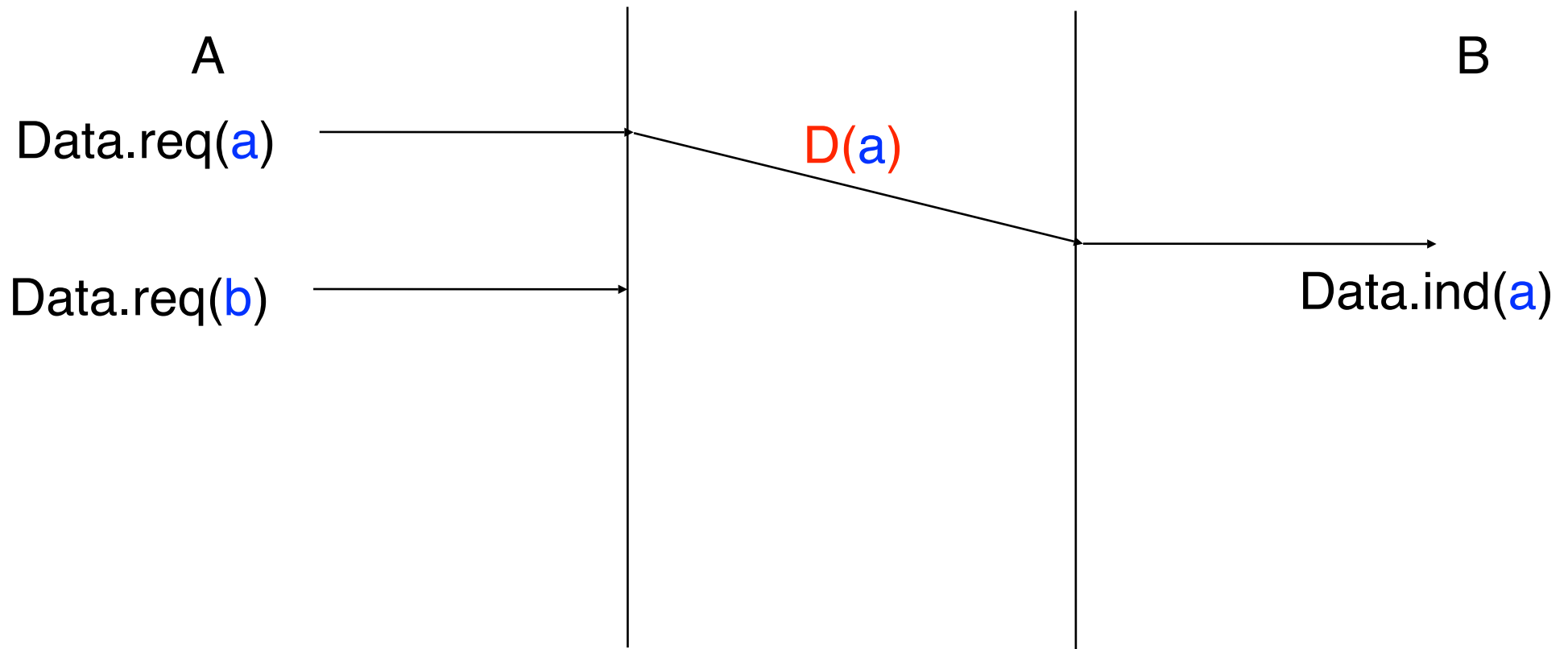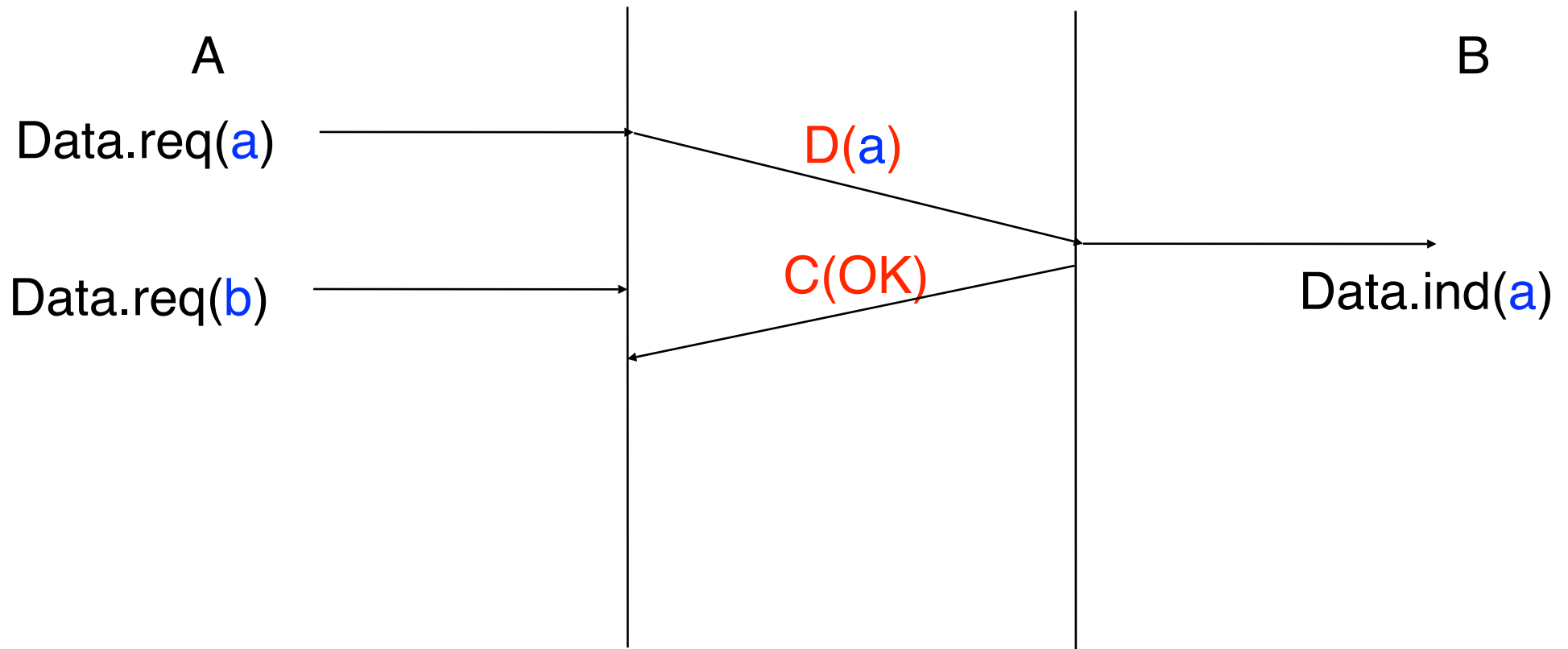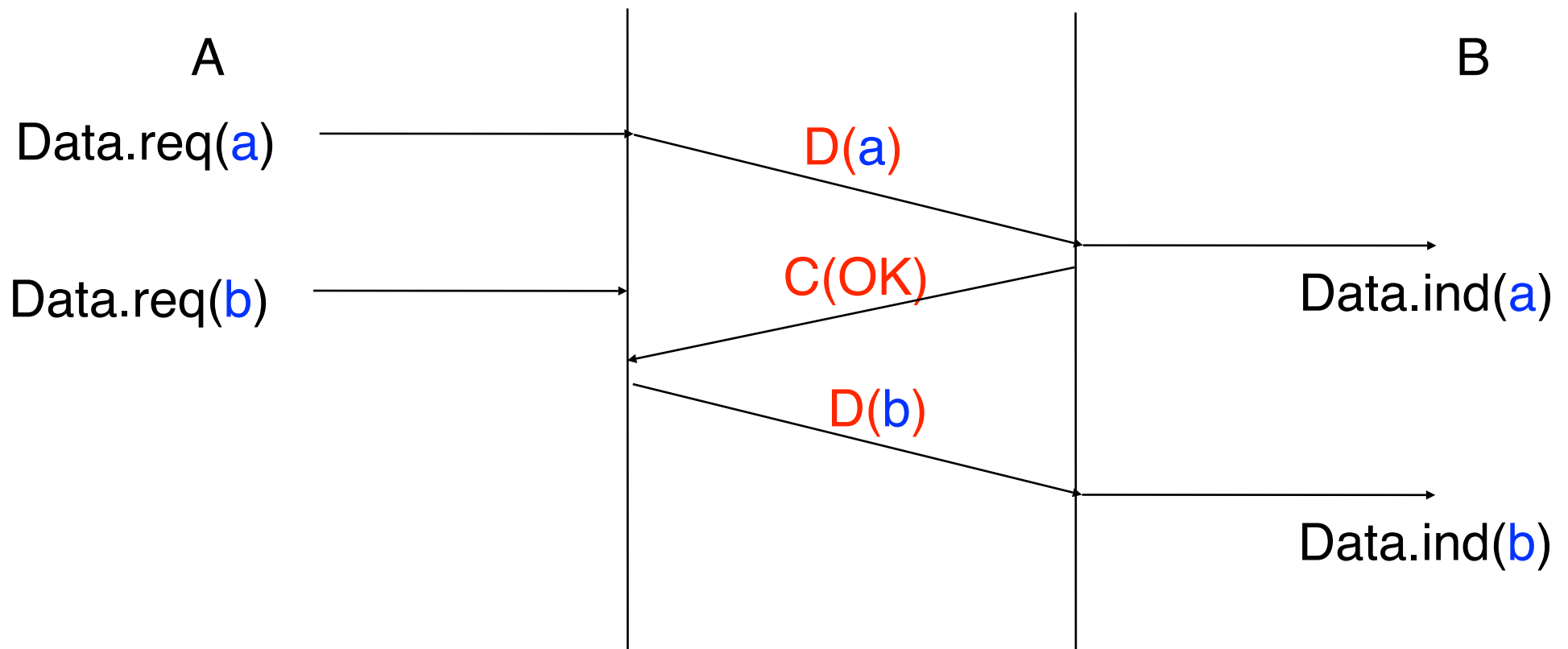$$\frac{-}{\text{Send}(C(OK))}$$

# Protocol 2 : Example

A

B

The sender only sends segments when authorised by the receiver

# Protocol 2 : Example



The sender only sends segments when authorised by the receiver

# Protocol 2 : Example

# Protocol 2 : Example

A
B

Data.req(a) ⟶ D(a)

Data.req(b) ⟶ C(OK) ⟶ Data.ind(a)
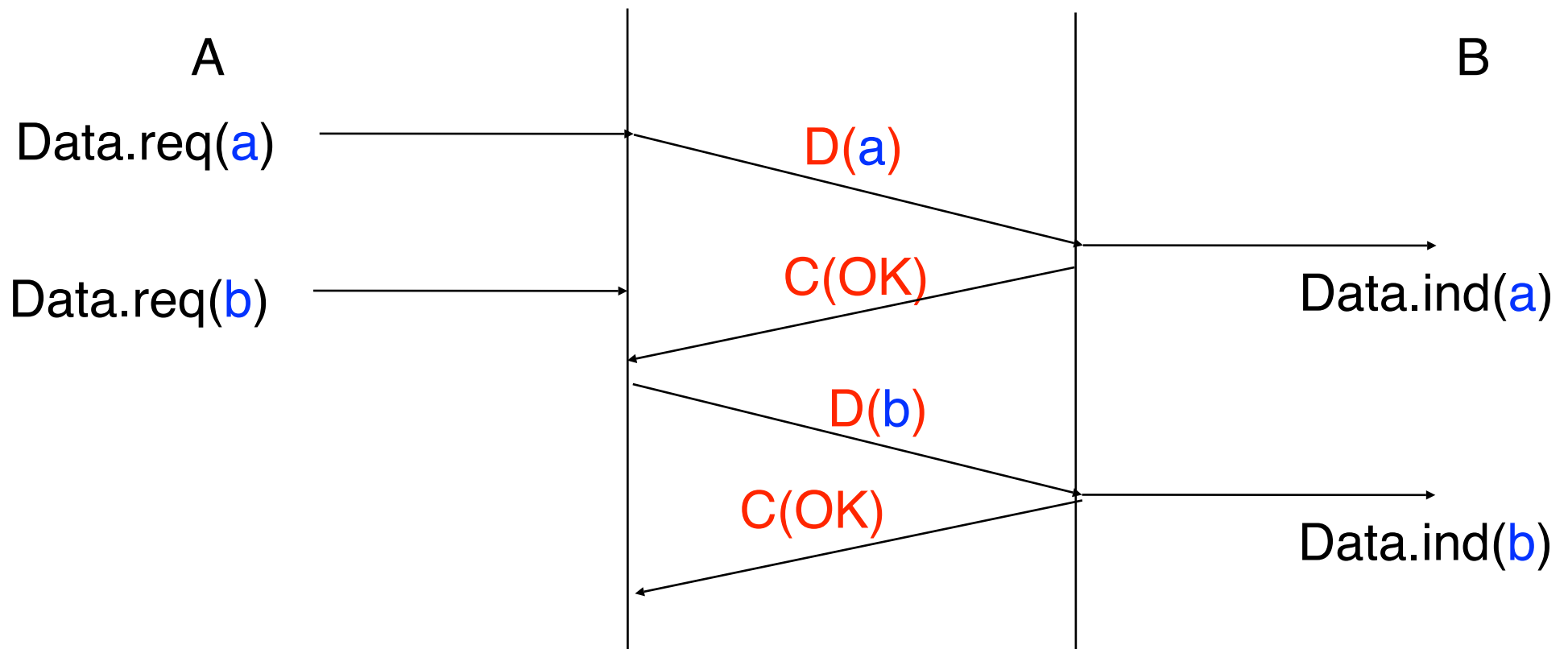
The sender only sends segments when authorised by the receiver

# Protocol 2 : Example



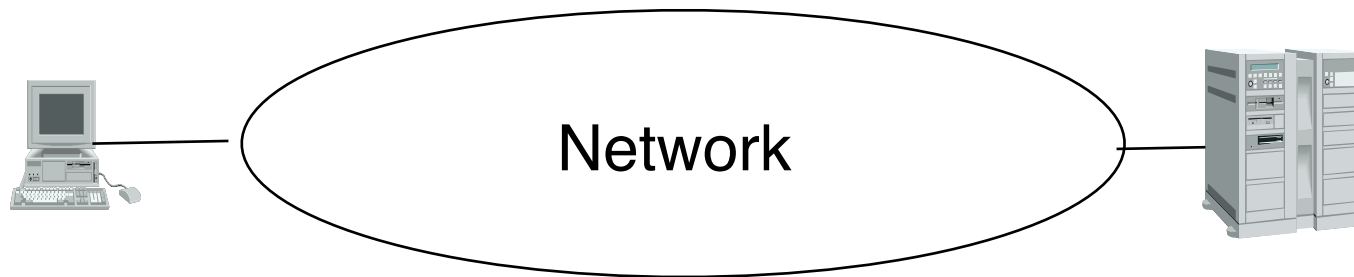The sender only sends segments when authorised by the receiver

# Protocol 2 : Example



The sender only sends segments when authorised
by the receiver

# Protocol 3

How can we provide a reliable service in the transport layer

Hypotheses
1. The application sends small SDUs
2. The network layer provides a perfect service
   1. Transmission errors are possible
   2. No packet is lost
   3. There is no packet reordering
   4. There are no duplications of packets
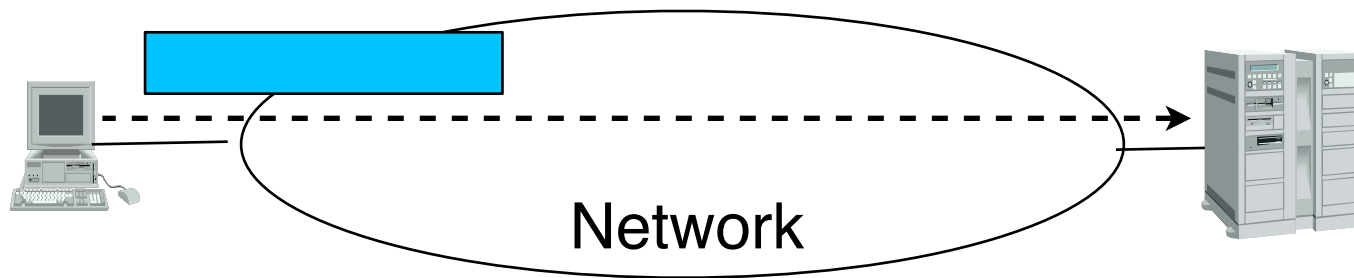3. Data transmission is unidirectional

# Transmission errors

Which types of transmission errors do we need to consider in the transport layer ?



Network

Physical-layer transmission errors caused by nature

- Random isolated error
  - one bit is flipped in the segment
- Random burst error
  - a group of n bits inside the segment is errored
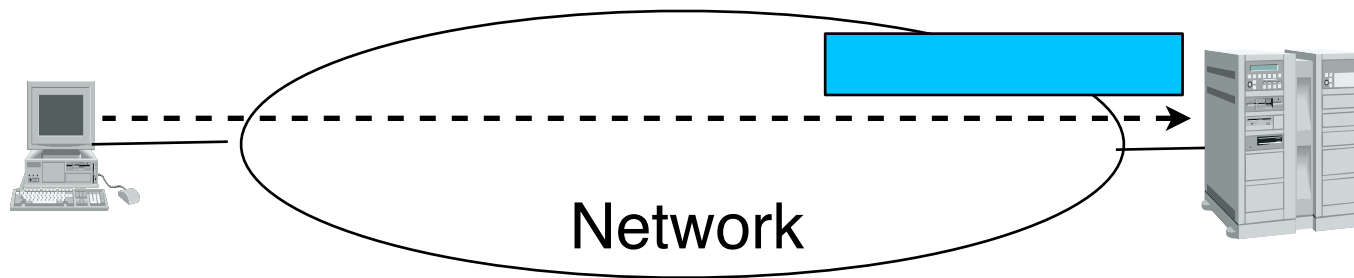    - most of the bits in the group are flipped

# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

# Security issues versus transmission errors

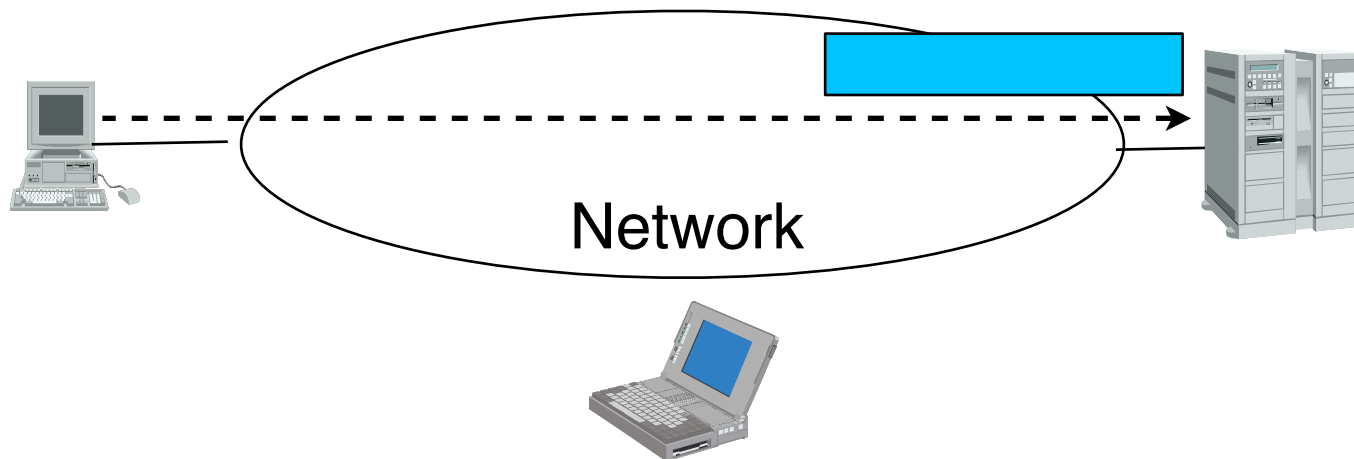Information sent over a network may become corrupted for other reasons than transmission errors



Network

# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

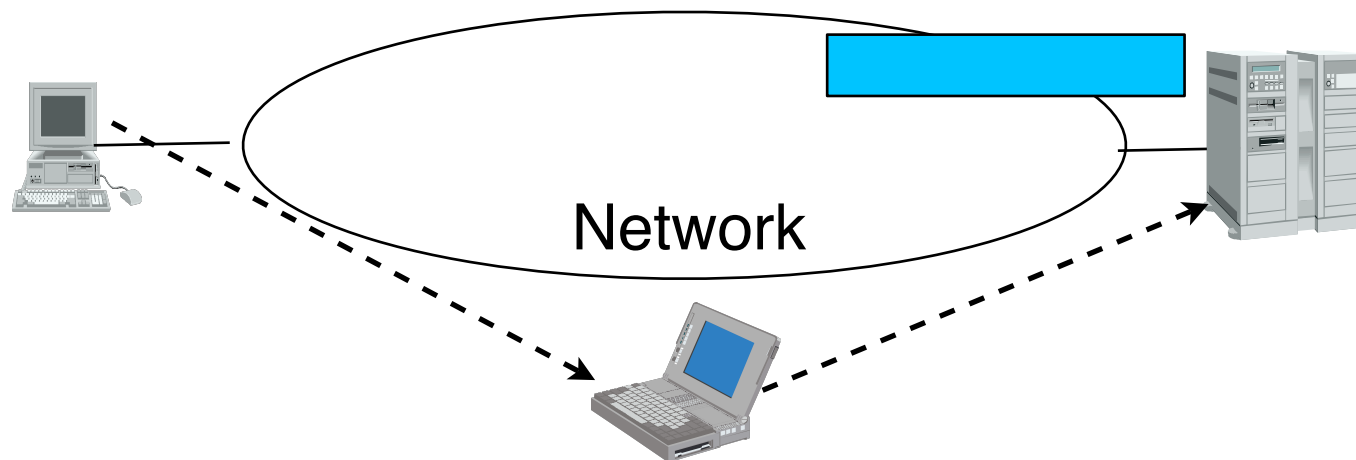# Security issues versus transmission errors

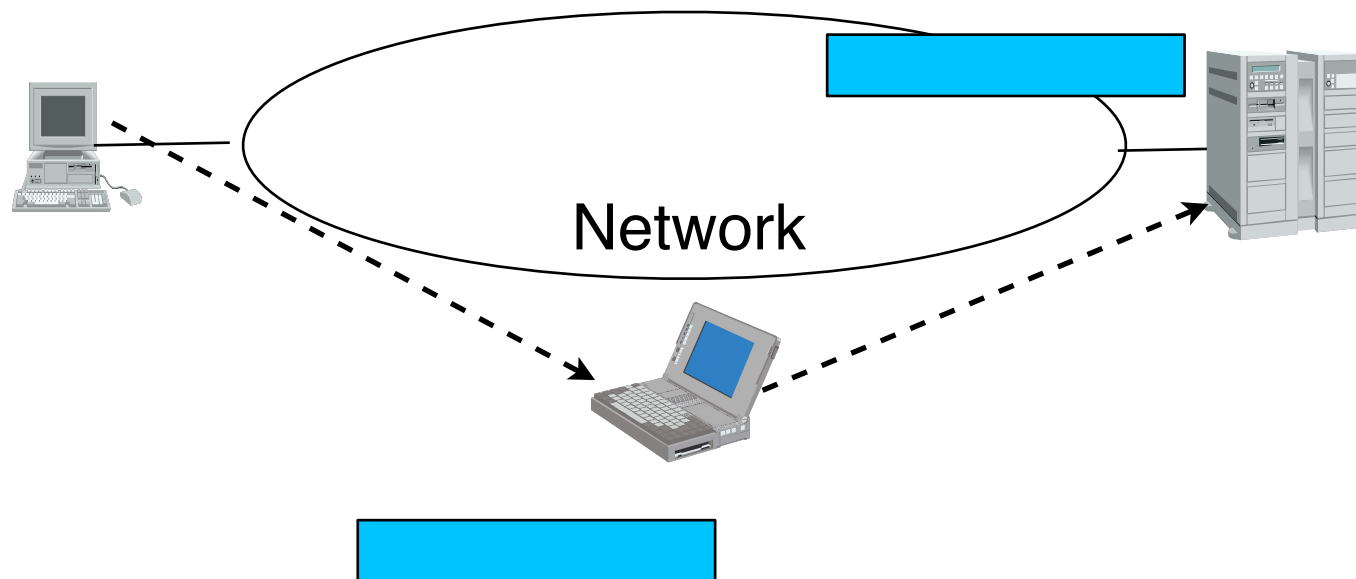Information sent over a network may become corrupted for other reasons than transmission errors



Network

# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

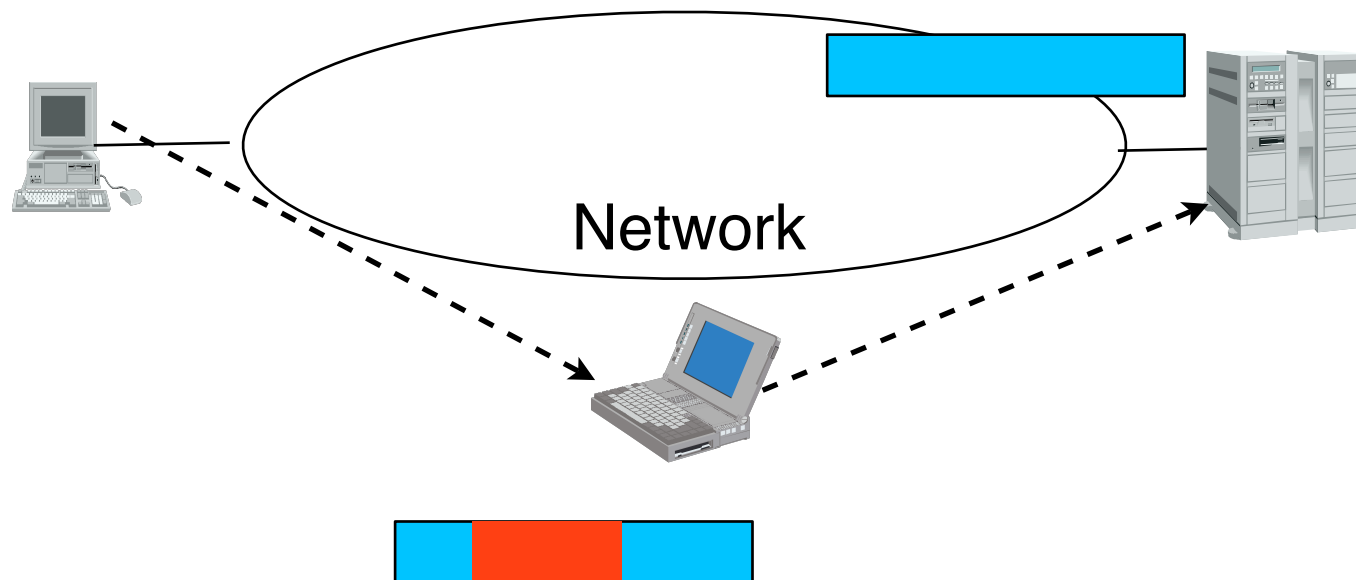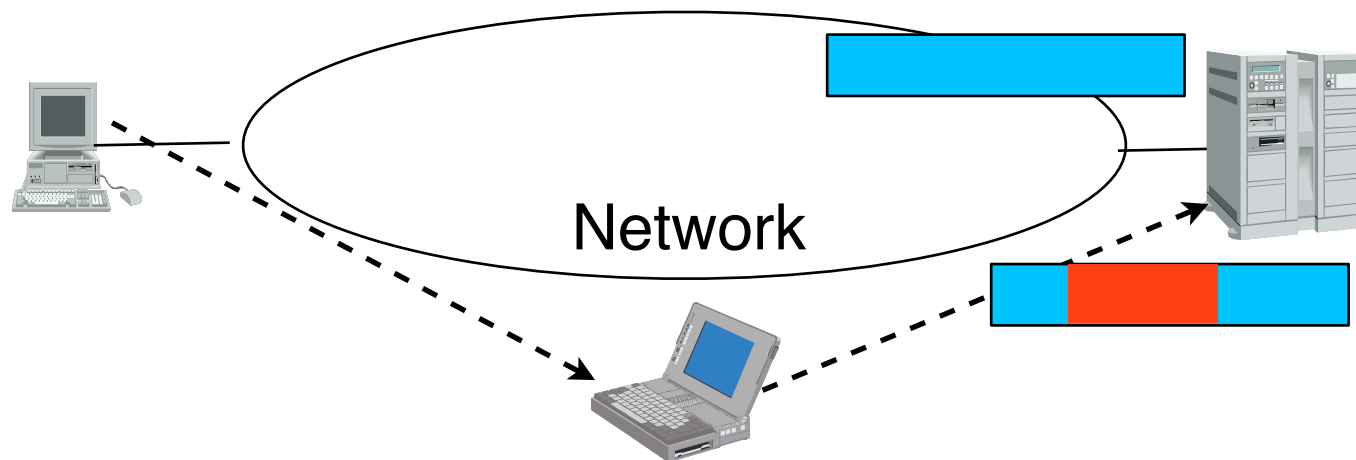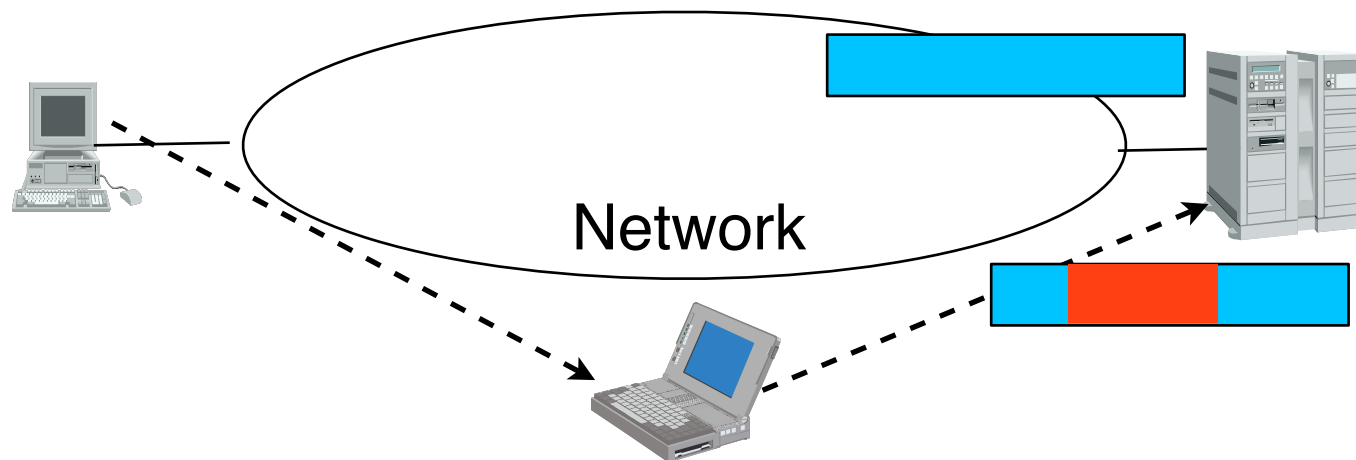# Security issues versus transmission errors

Information sent over a network may become corrupted for other reasons than transmission errors



Network

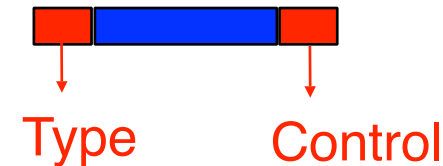These attacks are dealt by using special security protocols and mechanisms outside the transport layer

# How to detect transmission errors ?

## Principle

Sender adds some control information inside the segment

control information is computed over the entire segment and placed in the segment header or trailer

Type+ Control    Type    Control

Receiver checks that the received control information is correct by recomputing it

# Parity bits

Simple solution to detect transmission errors

Used on slow-speed serial lines
e.g. modems connected to the telephone network

Odd Parity
For each group of n bits, sender computes the n+1th bit so that the n+1 group contains an odd number of bits set to 1

Examples

```
0011010                    1101100
```

Even Parity

# Parity bits

Simple solution to detect transmission errors

Used on slow-speed serial lines
  e.g. modems connected to the telephone network

Odd Parity
  For each group of n bits, sender computes the n+1th bit so that the n+1 group contains an odd number of bits set to 1

    Examples

          0011010 0                    1101100

Even Parity

# Parity bits

Simple solution to detect transmission errors

Used on slow-speed serial lines
   e.g. modems connected to the telephone network

Odd Parity
   For each group of n bits, sender computes the
   n+1th bit so that the n+1 group contains an odd
   number of bits set to 1

      Examples

                0011010 0                    1101100 1

Even Parity

# Internet checksum

## Motivation

Internet protocols are implemented in software and we would like to have efficient algorithms to detect transmission errors that are easy to implement

## Solution

Internet checksum

Sender computes for each segment and over the entire segment the 1s complement of the sum of all the 16 bits words in the segment

Receiver recomputes the checksum over each received segment and verifies that it is correct. Otherwise, the

# Internet checksum : example

## Assume a segment composed of 48 bits

0110011001101100 0101010101010101 0000111100001111

# Internet checksum : example

## Assume a segment composed of 48 bits

0110011001101100                    0000111100001111
0101010101010101

# Internet checksum : example

## Assume a segment composed of 48 bits

```
0110011001101100                          0000111100001111
0101010101010101

1011101110111011
```

# Internet checksum : example

## Assume a segment composed of 48 bits

```
0110011001101100
0101010101010101

1011101110111011
0000111100001111
```

# Internet checksum : example

## Assume a segment composed of 48 bits

0110011001101100
0101010101010101

1011101110111011
0000111100001111

1100101011001010

# Internet checksum : example

## Assume a segment composed of 48 bits

0110011001101100
0101010101010101

1011101110111011
0000111100001111

1100101011001010                    <span style="color:red">0011010100110101</span>

# Cyclical Redundancy Check (CRC)

## Principle

Improve the performance of the Internet checksum by using polynomial codes

Sender and receiver agree on r+1 bits pattern called Generator (G)

Sender adds r bits of CRC to a d bits data segment such that the d+r bits pattern is exactly divisible by G using modulo 2 arithmetic

D * 2^r XOR R = n*G

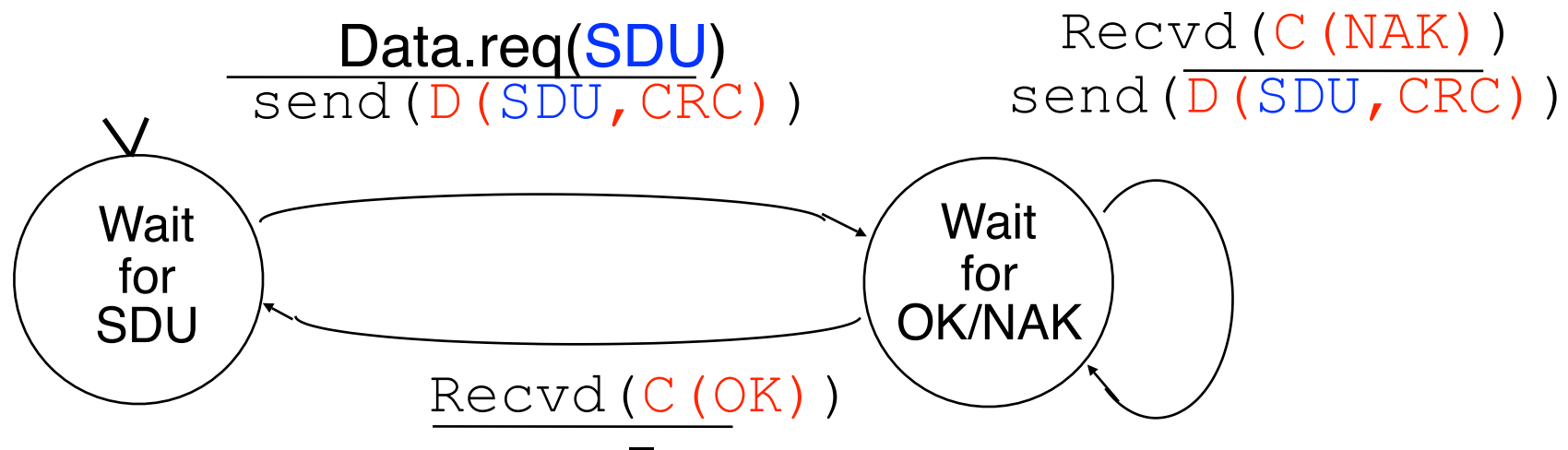| d bits | r bits |
|--------|--------|

All computations are done in modulo 2 arithmetic by using XOR

1011 + 0101 = 1110          1001 + 1101 = 0100
1011 - 0101 = 1110          1001 - 1101 = 0100

# Detection of transmission errors (2)

Behaviour of the receiver

### If the checksum is correct
Send an <span style="color:red">OK</span> control segment to the sender to
confirm the reception of the data segment
allow the sender to send the next segment

### If the checksum is incorrect
The content of the segment is corrupted and must be discarded
Send a special control segment (<span style="color:red">NAK</span>) to the sender to ask it to retransmit the corrupted data segment
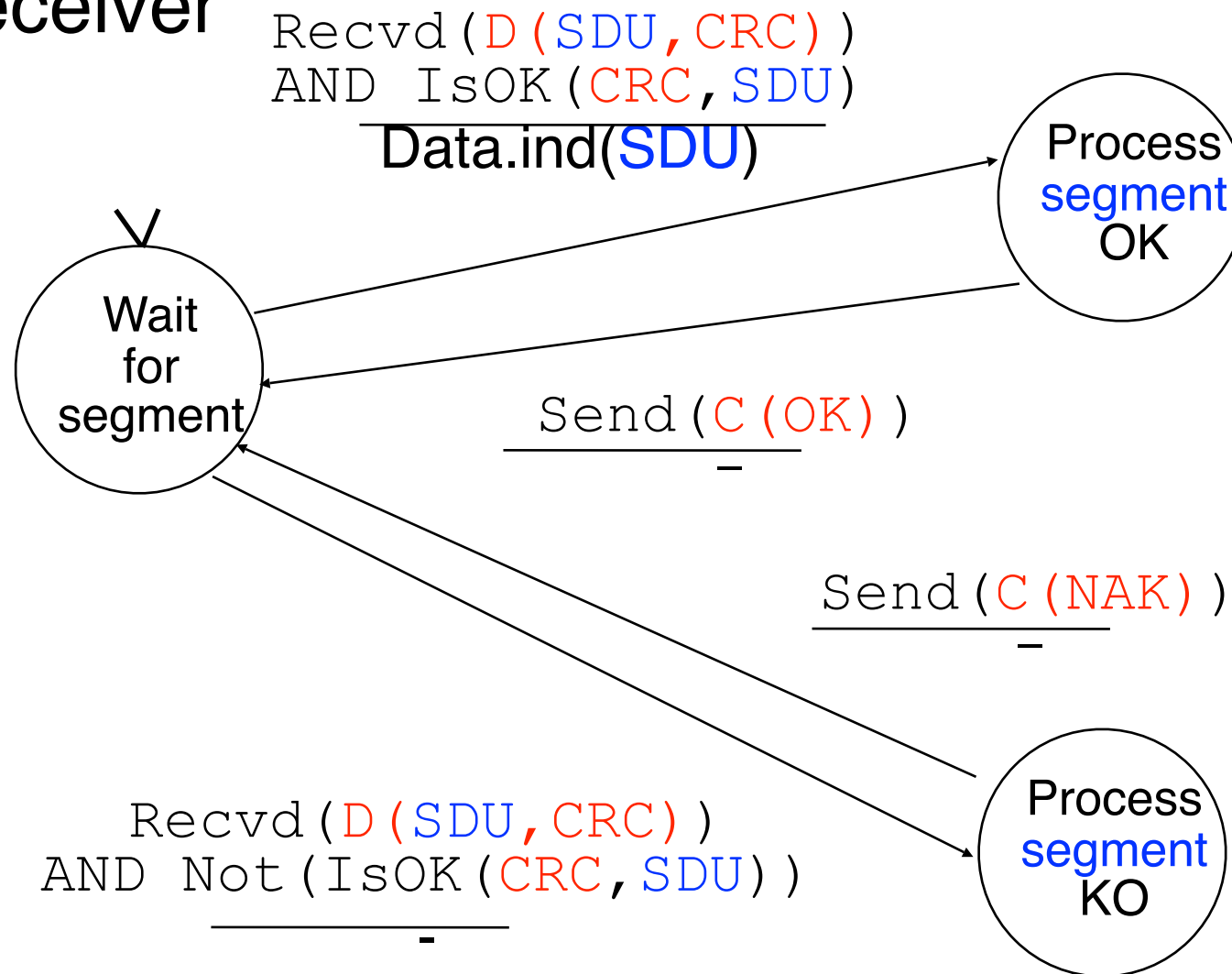
# Protocol 3a : Receiver



Receiver

Recvd(D(SDU,CRC))
AND IsOK(CRC,SDU)
―――――――――――
Data.ind(SDU)

Wait for segment

Process segment OK

Send(C(OK))
―――――――

Send(C(NAK))
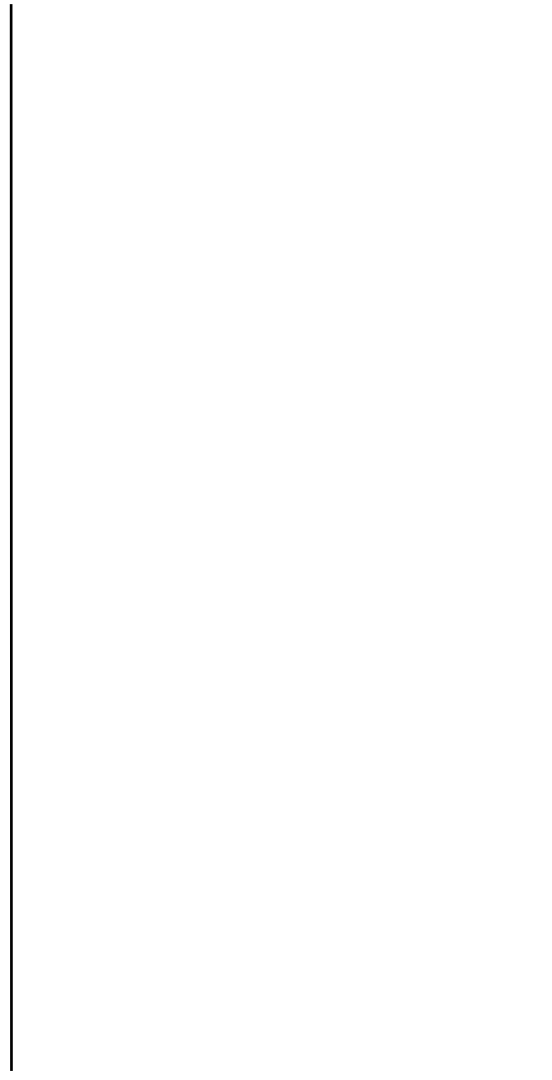――――――――

Recvd(D(SDU,CRC))
AND Not(IsOK(CRC,SDU))
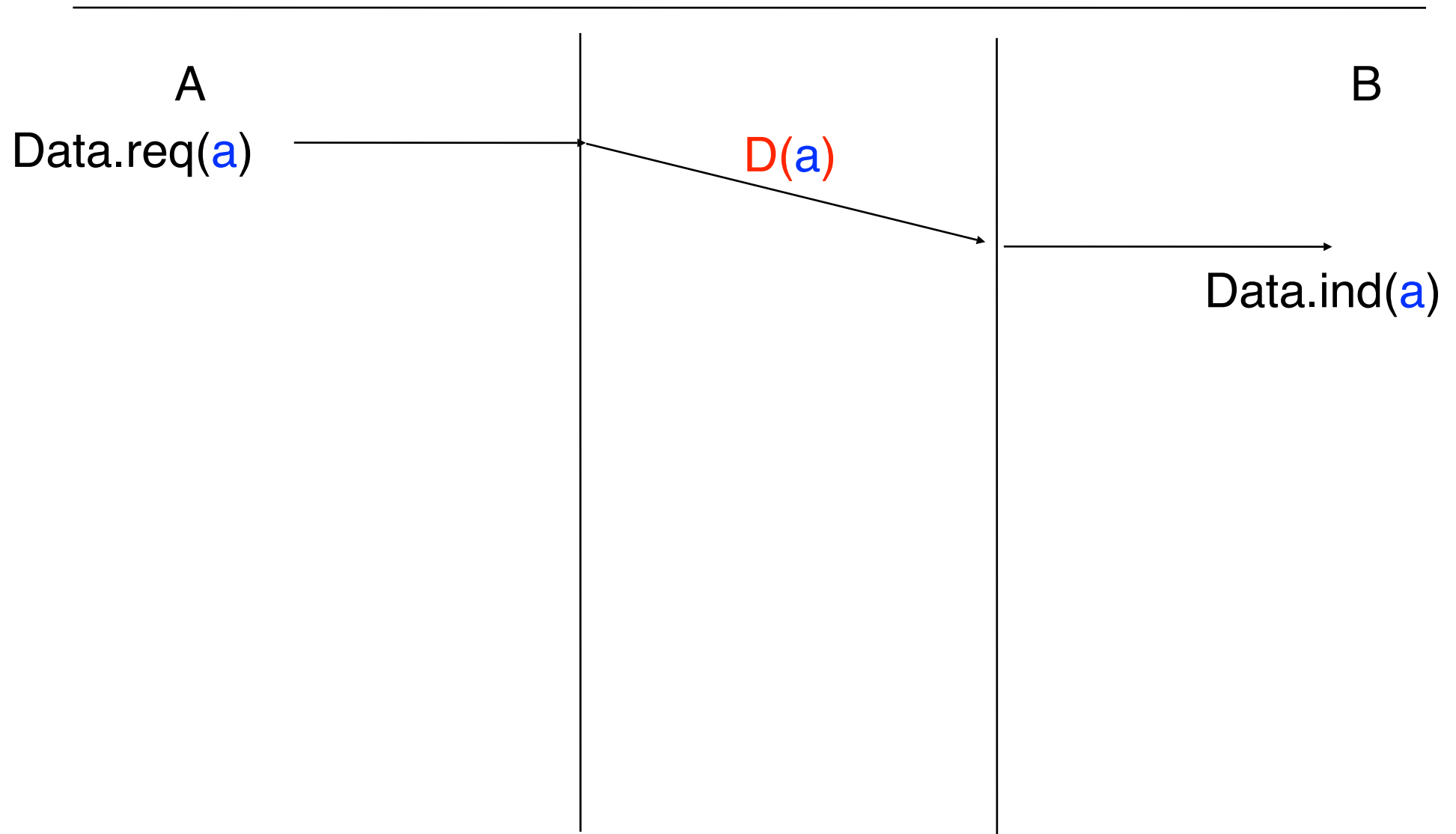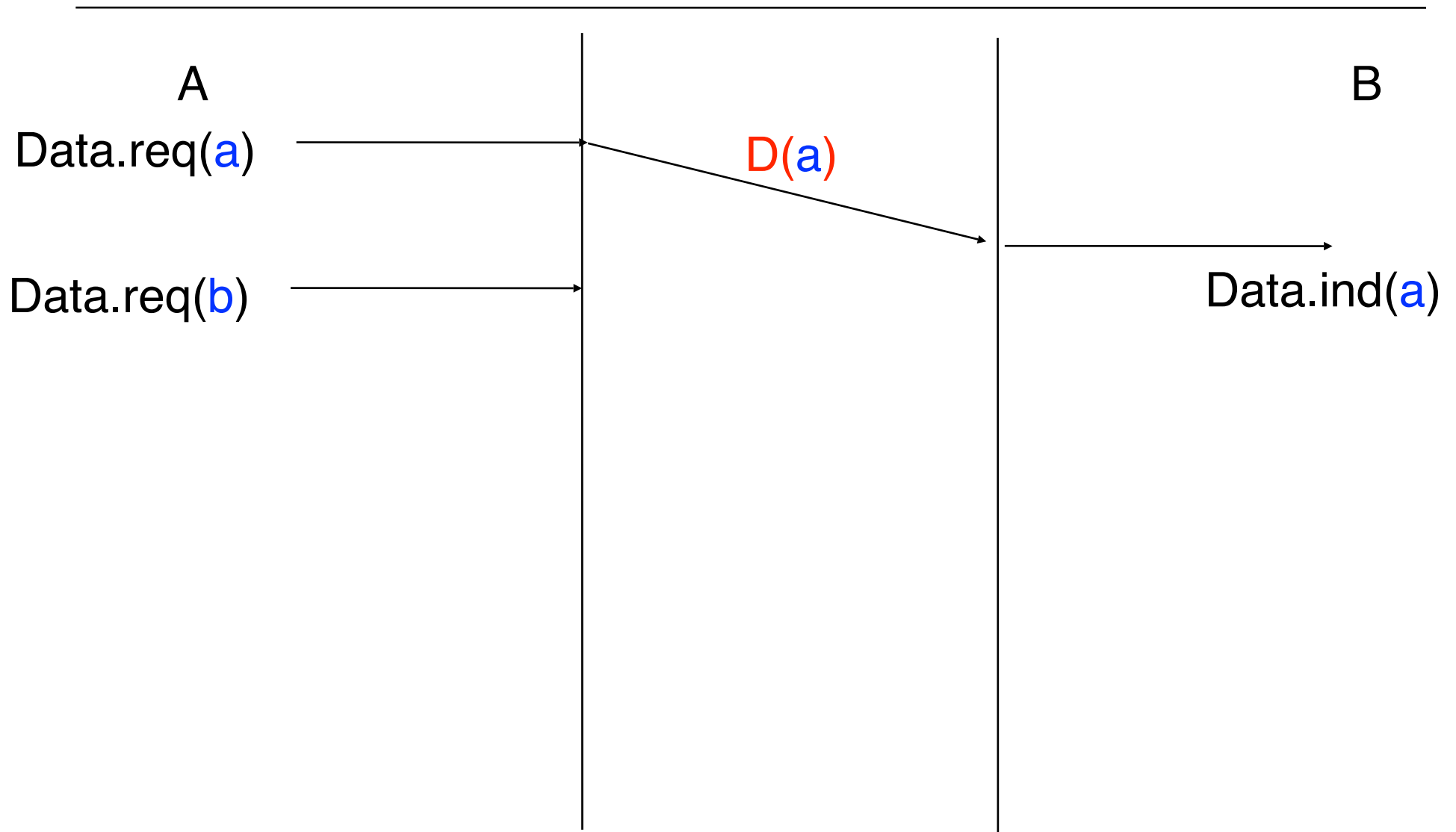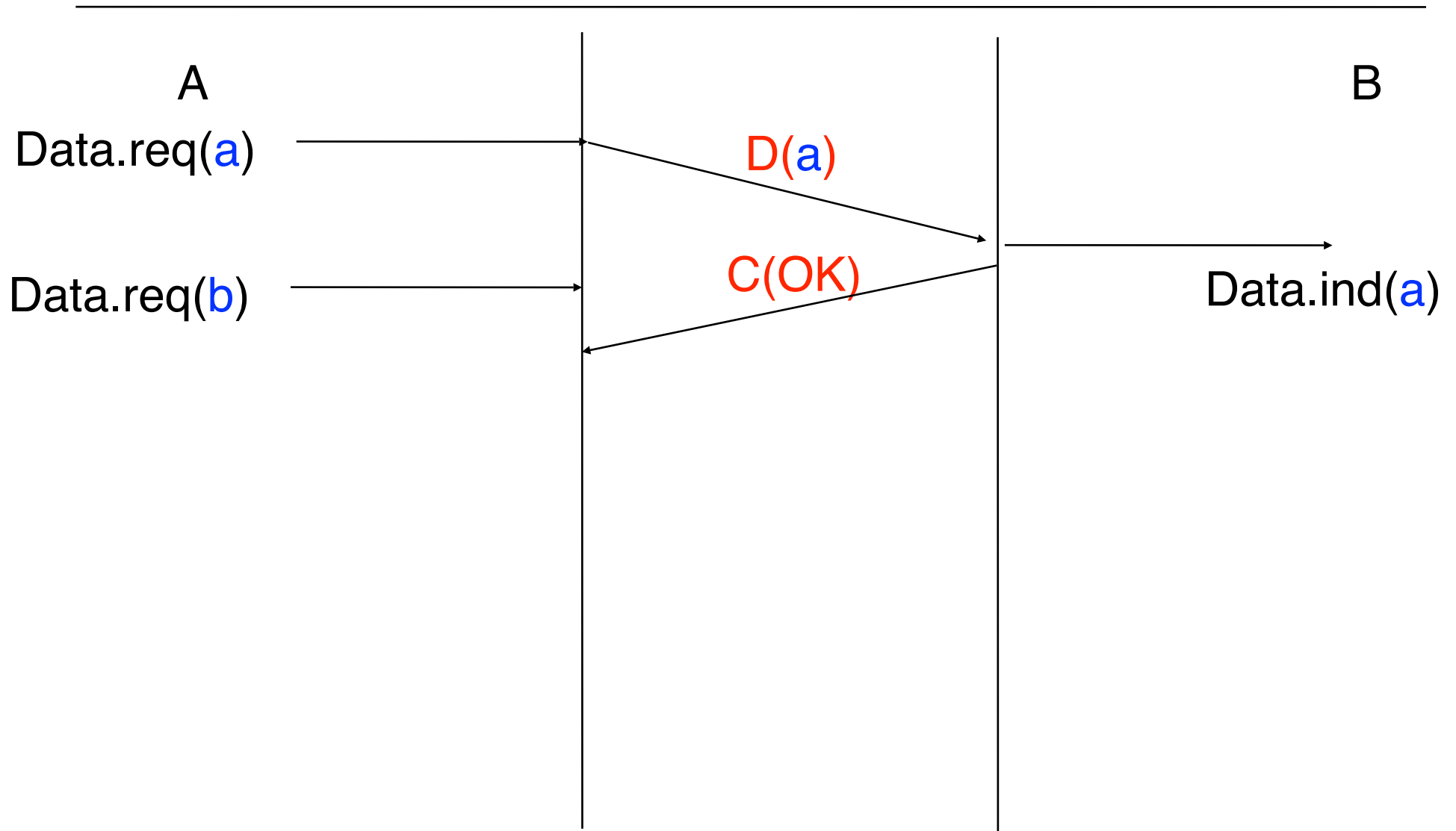―――――――――――
Process segment KO

# Protocol 3a : Example

A

B

# Protocol 3a : Example

# Protocol 3a : Example



A

Data.req(a)

Data.req(b)

D(a)

B

Data.ind(a)

# Protocol 3a : Example

A

Data.req(a)

D(a)

Data.req(b)

C(OK)

B

Data.ind(a)

# Protocol 3a : Example



A
B

Data.req(a)

D(a)

Data.ind(a)

Data.req(b)

C(OK)

Transmission error

D(b')

Invalid checksum

# Protocol 3a : Example



A

Data.req(a)

D(a)

Data.req(b)

C(OK)

Data.ind(a)

*Transmission error*

D(b')

*Invalid checksum*

C(NAK)

B

# Protocol 3a : Example

# Protocol 3a : Example



A

Data.req(a)

Data.req(b)

D(a)

C(OK)

*Transmission error*

D(b')

C(NAK)

*Invalid checksum*

*Retransmission*

D(b)

B

Data.ind(a)

Data.ind(b)

# Protocol 3a : Example



A                      B

Data.req(a)    D(a)

Data.req(b)    C(OK)          Data.ind(a)

*Transmission error*

D(b')

C(NAK)    *Invalid checksum*

*Retransmission*    D(b)

C(OK)          Data.ind(b)

# Protocol 3b

How can we provide a reliable service in the transport layer ?

Hypotheses
1. The application sends small SDUs
2. The network layer provides a perfect service
    1. Transmission errors are possible
    2. Packets can be lost
    3. There is no packet reordering
    4. There are no duplications of packets
3. Data transmission is unidirectional

2. How to deal with these problems ?

# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?

A

B

# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?



A

Data.req(a)

D(a)

B

Data.ind(a)

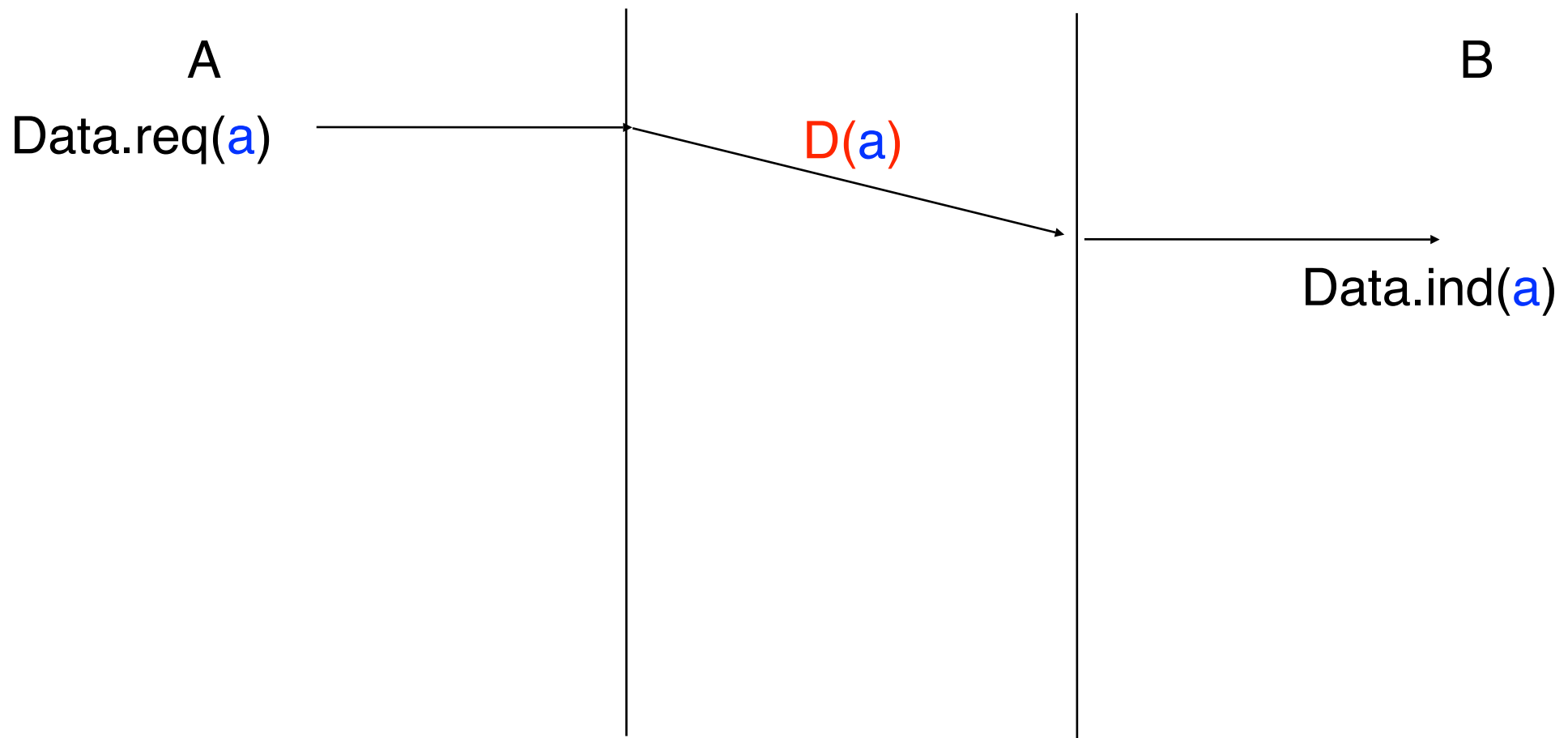# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?

# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?
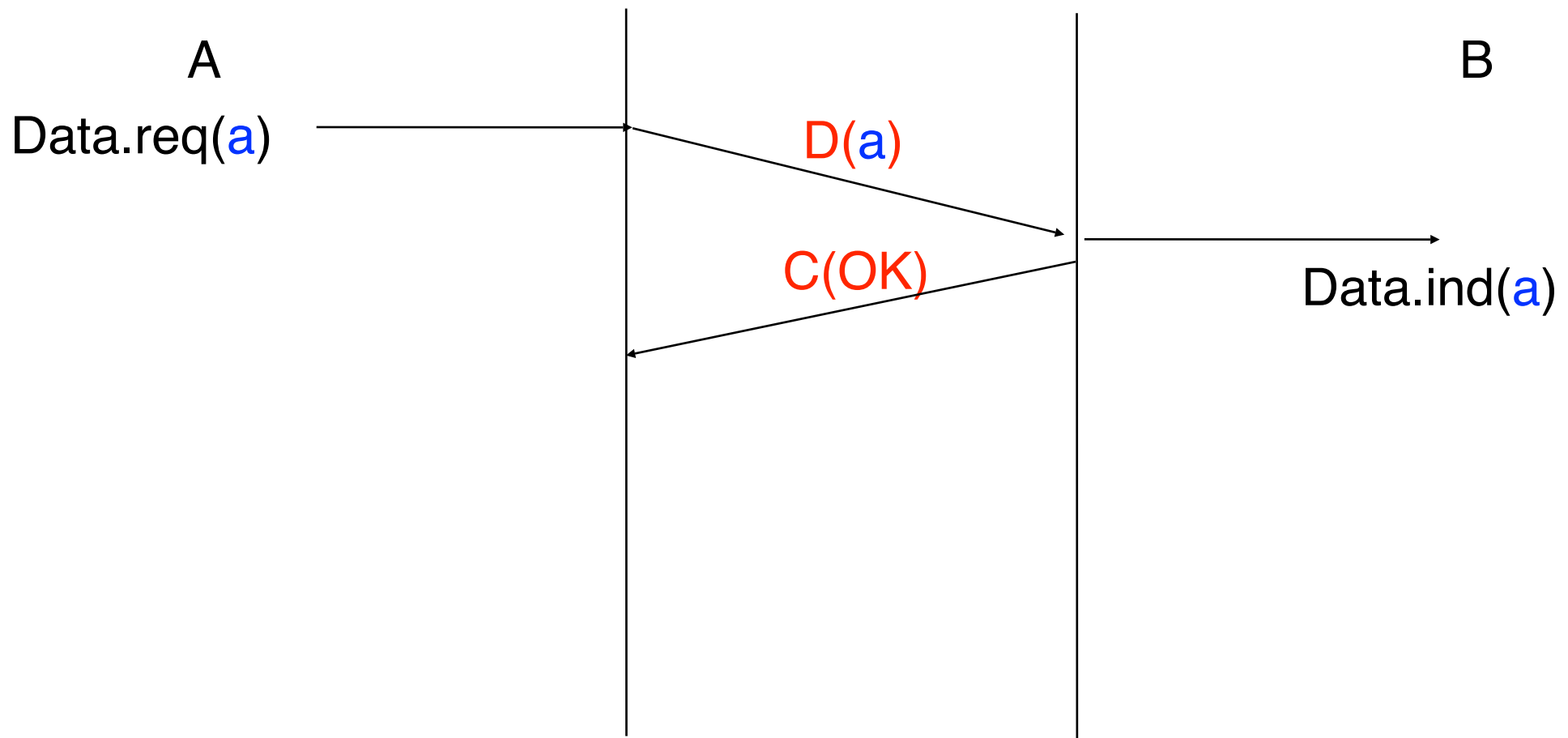
# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?

# Protocol 3a and segment losses

## How do segment losses affect protocol 3a ?
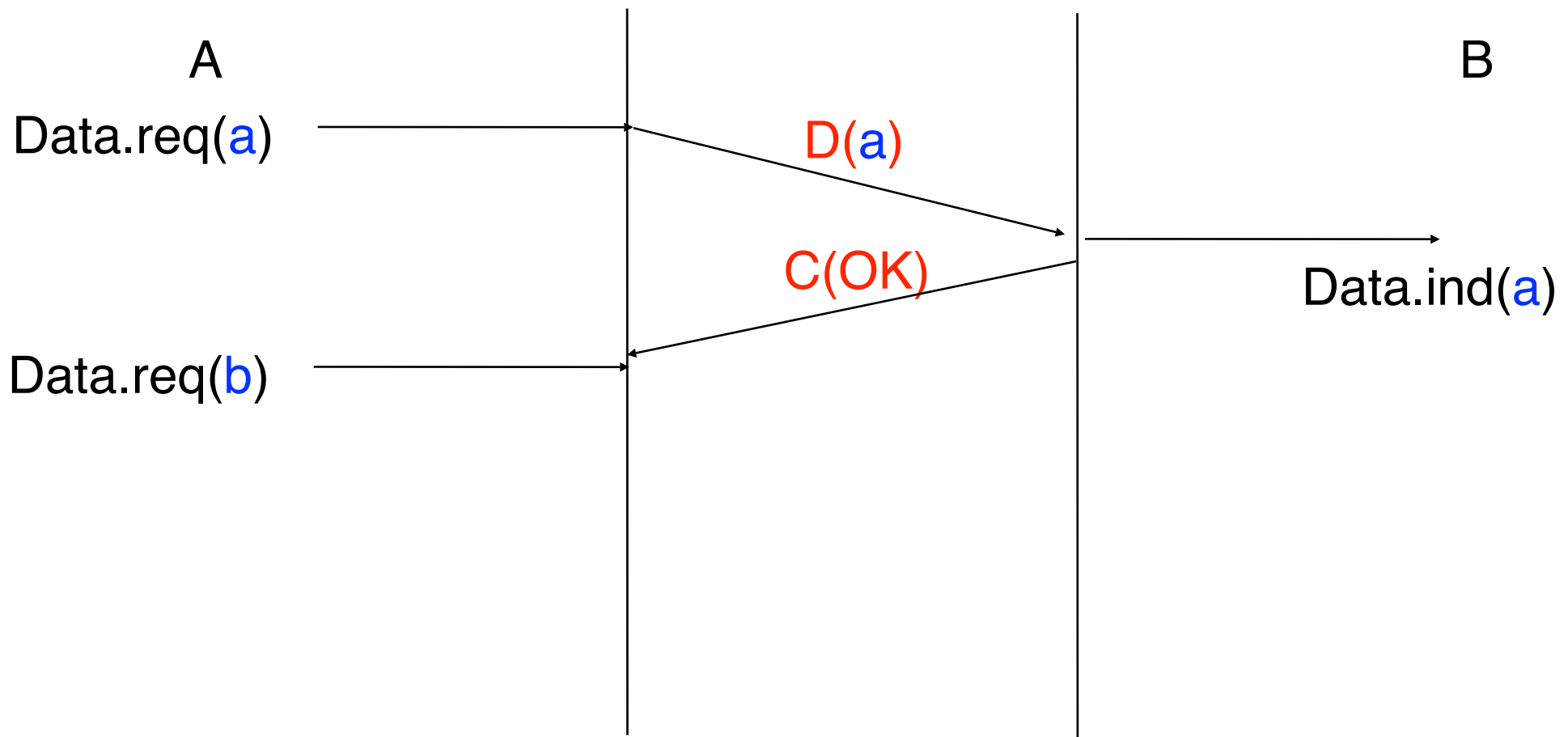


DEADLOCK

*A is waiting for a control segment*

*B is waiting for a data segment*

# Protocol 3b

## Modification to the sender
Add a retransmission timer to retransmit the lost segment after some time



No modification to the receiver

# Protocol 3b : Example

A

B

# Protocol 3b : Example

Data.req(a)

*start timer*

D(a)

B

Data.ind(a)

# Protocol 3b : Example



A

Data.req(a)

*start timer*

D(a)

C(OK)

*cancel timer*

B

Data.ind(a)

# Protocol 3b : Example



A

Data.req(a)

*start timer*

D(a)

*cancel timer*

C(OK)

Data.req(b)

*start timer*

B

Data.ind(a)

# Protocol 3b : Example



A

Data.req(a)

*start timer*

D(a)

C(OK)

*cancel timer*

Data.req(b)

*start timer*

B

Data.ind(a)

# Protocol 3b : Example

# Protocol 3b : Example

# Protocol 3b : Example

A

Data.req(a)

*start timer*

D(a)

C(OK)

*cancel timer*

Data.req(b)

*start timer*

D(b)

*timer expires*

D(b)

*Lost segment*

B

Data.ind(a)

Data.ind(b)

# Protocol 3b : Example



A

Data.req(a)

*start timer*

D(a)

C(OK)

*cancel timer*

Data.req(b)

*start timer*

D(b)

*Lost segment*

D(b)

*timer expires*

C(OK)

B

Data.ind(a)

Data.ind(b)

# Protocol 3b : Example



**A**

Data.req(a)

*start timer*

D(a)

C(OK)

*cancel timer*

Data.req(b)

*start timer*

D(b)

*Lost segment*

D(b)

*timer expires*

C(OK)

**B**

Data.ind(a)

Data.ind(b)

## Does this protocol always work ?

# Protocol 3b : Example

A

B

# Protocol 3b : Example

A

Data.req(a)

*start timer*

D(a)

*cancel timer*

C(OK)

B

Data.ind(a)

# Protocol 3b : Example



A

B

Data.req(a)

*start timer*

D(a)

Data.ind(a)

C(OK)

*cancel timer*

Data.req(b)

*start timer*

D(b)

Data.ind(b)

# Protocol 3b : Example

A                B

Data.req(a) &rarr; D(a) &rarr; Data.ind(a)

*start timer*

C(OK)

*cancel timer*

Data.req(b) &rarr; D(b) &rarr; Data.ind(b)

*start timer*

*Lost segment* &rarr; ● C(OK)

Data.ind(b)

Protocol 3b : Example

CNP3/2008.3.

© O. Bonaventure, 2008

# Protocol 3b : Example



A

Data.req(a)

start timer

D(a)

C(OK)

cancel timer

Data.req(b)

start timer

D(b)

C(OK)

Lost segment

D(b)

timer expires

C(OK)

B

Data.ind(a)

Data.ind(b)

Data.ind(b)

## How to solve this problem ?

# Alternating bit protocol

Principles of the solution
- Add sequence numbers to each data segment sent by sender
- By looking at the sequence number, the receiver can check whether it has already received this segment

Contents of each segment
- Data segments
- Control segments



Type, Seq. number    CRC

Type, Seq. number    CRC

How many bits do we need for the sequence number?
- a single bit is enough

# Alternating bit protocol Sender

# Alternating bit protocol Receiver

Recvd(D(1,SDU,CRC))
AND IsOK(CRC,SDU)
Send(C(OK1))

Recvd(D(0,SDU,CRC))
AND IsOK(CRC,SDU)
Data.ind(SDU)

Recvd(D(0,SDU,CRC))
AND IsOK(CRC,SDU)
Send(C(OK0))

-
Send(C(OK0))

Wait for D(0,...)

Process SDU0 OK

Wait for D(1,...)

Bad(D(?,SDU,CRC))
-

Bad(D(?,SDU,CRC))
-

Recvd(D(1,SDU,CRC))
AND IsOK(CRC,SDU)
Data.ind(SDU)

-
Send(C(NAK1))

-
Send(C(NAK0))

-
Send(C(OK1))

Process SDU0 KO

Process SDU1 OK

Process SDU1 KO

# Alternating bit protocol Example

A

B

# Alternating bit protocol Example

# Alternating bit protocol Example

# Alternating bit protocol Example

# Alternating bit protocol Example



A

Data.req(a)

Data.req(b)

D(0,a)

C(OK0)

D(1,b)

*Retransmission timer*

D(1,b)

*D(1,b) recvd*

C(OK1)

Data.req(c)

C(OK1)

*D(1,b) recvd*

D(0,c)

B

Data.ind(a)

Data.ind(b)

*Duplicate detected*

*Lost segment*

# Alternating bit protocol Example

# Alternating bit protocol Example

# Performance of the alternating bit protocol

What is the performance of the ABP in this case
  One-way delay : 250 msec
  Physical layer throughput : 50 kbps
  segment size : 1000 bits

# Performance of the alternating bit protocol

What is the performance of the ABP in this case
One-way delay : 250 msec
Physical layer throughput : 50 kbps
segment size : 1000 bits



1000/50000 = 20 msec

Data

Ack

250 msec

520 msec

# Performance of the alternating bit protocol

What is the performance of the ABP in this case
- One-way delay : 250 msec
- Physical layer throughput : 50 kbps
- segment size : 1000 bits

$1000/50000 = 20$ msec

Data

Ack

250 msec

520 msec

-> Performance is function of
*bandwidth * round-trip-time*

# How to improve the alternating bit protocol ?

Use a pipeline

Principle

The sender should be allowed to send more than one segment while waiting for an acknowledgement from the receiver

# How to improve
# the alternating bit protocol ?

Use a pipeline
Principle
　　The sender should be allowed to send more than
　　one segment while waiting for an
　　acknowledgement from the receiver

# How to improve the alternating bit protocol ? (2)

## Modifications to alternating bit protocol

### Sequence numbers inside each segment
Each data segment contains its own sequence number
Each control segment indicates the sequence number of the data segment being acknowledged (OK/NAK)

### Sender
Needs enough buffers to store the data segments that have not yet been acknowledged to be able to retransmit them if required

### Receiver
Needs enough buffers to store the out-of-sequence segments

# How to improve the alternating bit protocol ? (2)

## Modifications to alternating bit protocol

### Sequence numbers inside each segment
Each data segment contains its own sequence number
Each control segment indicates the sequence number of the data segment being acknowledged (OK/NAK)

### Sender
Needs enough buffers to store the data segments that have not yet been acknowledged to be able to retransmit them if required

### Receiver
Needs enough buffers to store the out-of-sequence segments

*How to avoid an overflow of the receiver's buffers ?*

# Sliding window

## Principle

Sender keeps a list of all the segments that it is allowed to send

`sending_window`



... 0 1 2 3 4 5 *6 7* 8 9 10 11 12 13 14 15 ....

Acked segments

Unacknowledged segments

Available seq. nums

Forbidden seq. num.

Receiver also maintains a receiving window with the list of acceptable sequence number

receiving_window

Sender and receiver must use compatible windows

sending_window ≤ receiving window

For example, window size is a constant for a given protocol or negotiated during connection establishment phase

# Sliding windows : example

## Sending and receiving window : 3 segments

A

Sending window

<span style="color:red">[</span> <span style="color:green">0 1 2</span> <span style="color:red">]</span> 3 4 5 6 7 8

B

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

© O. Bonaventure, 2008

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

© O. Bonaventure, 2008

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

© O. Bonaventure, 2008

# Sliding windows : example

## Sending and receiving window : 3 segments

# Sliding windows : example

## Sending and receiving window : 3 segments

# Encoding sequence numbers

## Problem

How many bits do we have in the segment header to encode the sequence number

N bits means $2^N$ different sequence numbers

## Solution

place inside each transmitted segment its sequence number modulo $2^N$

The same sequence number will be used for several different segments

be careful, this could cause problems...

## Sliding window

List of consecutive sequence numbers (modulo $2^N$) that the sender is allowed to transmit

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

A

Sending window

<span style="color:red">[</span><span style="color:green">0 1 2</span><span style="color:red">]</span>3

B

OK(0)

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field



A

Sending window

```
0 1 2 3
0 1 2 3
0 1 2 3
0 1 2 3


0 1 2 3

0 1 2 3
```

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

C(OK0)

C(OK1)

B

Data.ind(a)

Data.ind(b)

Data.ind(c)

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Sliding window : second example

## 3 segments sending and receiving window
### Sequence number encoded as 2 bits field

# Reliable transfer with a sliding window

How to provide a reliable data transfer with a sliding window
- How to react upon reception of a control segment ?
- Sender's and receiver's behaviours

Basic solutions

Go-Back-N
- simple implementation, in particular on receiving side
- throughput will be limited when losses occur

Selective Repeat
- more difficult from an implementation viewpoint
- throughput can remain high when limited losses occur

# GO-BACK-N

## Principle
### Receiver must be as simple as possible

### Receiver
- Only accepts consecutive in-sequence data segments
- Meaning of control segments
  - Upon reception of data segment
    - OKX means that all data segments, up to and including X have been received correctly
    - NAKX means that the data segment whose sequence number is X contained an error or was lost

### Sender
- Relies on a retransmission timer to detect segment losses
- Upon expiration of retransmission time or arrival of a NAK segment : retransmit all the unacknowledged data segments
  - the sender may thus retransmit a segment that was already received correctly but out-of-sequence at destination

# Go-Back-N : Receiver

State variable

`next` : sequence number of expected data segment



```
Recvd(D(next,SDU,CRC))
AND NOT(IsOK(CRC,SDU))
        discard(SDU);
send(C(NAK,next,CRC));
```

```
Recvd(D(next,SDU,CRC))
   AND IsOK(CRC,SDU)
     Data.ind(SDU)
```

Wait

Process
SDU
OK

```
        -
send(C(OK,next,CRC));
   next=(next+1);
```

```
   Recvd(D(t<>next,SDU,CRC))
      AND IsOK(CRC,SDU)
        discard(SDU);
send(C(OK,(next-1) ,CRC));
```

# Go-Back-N : Sender

State variables

base : sequence number of oldest data segment

seq : first available sequence number

W : size of sending window



```
Recvd(C(?,?,CRC))
and NOT( CRCOK(C(?,?,CRC)))
```

```
Data.req(SDU)
AND ( seq < (base+w) )
if (seq==base) { start_timer ; }
insert_in_buffer(SDU);
send(D(seq,SDU,CRC));
seq=seq+1 ;
```

Wait

```
[ Recvd(C(NAK,?,CRC))
 and CRCOK(C(NAK,?,CRC))]
       or timer expires
for (i=base;i<seq; i=i+1)
{ send(D(i,SDU,CRC)); }
restart_timer();
```

```
Recvd(C(OK,t,CRC))
and CRCOK(C(OK,t,CRC))
base=(t+1);
if (base==seq)
{ cancel_timer();}
else
{ restart_timer(); }
```

# Go-Back-N : Example

A

Sending window

<span style="color:green">0 1 2</span>3

B

# Go-Back-N : Example

A

Sending window

0 1 2 3

0 1 2 3

Data.req(a)

D(0,a)

B

Data.ind(a)

# Go-Back-N : Example

# Go-Back-N : Example



A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

B

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

Transmission error

Data.ind(a)

Invalid CRC, discarded

Not expected seq num, discarded

# Go-Back-N : Example

A

Sending window

0 1 2 3
0 1 2 3
0 1 2 3
0 1 2 3

0 1 2 3

B

Data.req(a)

Data.req(b)

D(0,a)

Data.req(c)

D(1,b)

*Transmission error*

Data.ind(a)

D(2,c)

*Invalid CRC, discarded*

C(OK,0)

*Not expected seq num, discarded*

CNP3/2008.3.

© O. Bonaventure, 2008

# Go-Back-N : Example

A                                                                 B

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

Transmission error

Data.ind(a)

Invalid CRC, discarded

C(OK,0)

C(NAK,1)

Not expected seq num, discarded

# Go-Back-N : Example

# Go-Back-N : Example

# Go-Back-N : Example

# Go-Back-N : Example

# Go-Back-N : Example



A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

B

Data.req(a)

Data.req(b)

D(0,a)

Data.req(c)

D(1,b)

*Transmission error*

Data.ind(a)

D(2,c)

*Invalid CRC, discarded*

C(OK,0)

C(NAK,1)

*Not expected seq num, discarded*

*Retransmission*
Data.req(d)

C(OK,0)

D(1,b)

Data.req(e)

D(2,c)

Data.ind(b)

D(3,d)

Data.ind(c)

Data.ind(d)

*Sending window is full*
*Application blocked*
*e will be accepted*
*and sent later*

# Go-Back-N : Example (2)

A

Sending window

0 1 2 3

B

# Go-Back-N : Example (2)

A

Sending window

0 1 2 3

*0* 1 2 3

*0 1* 2 3

Data.req(a)

D(0,a)

Data.ind(a)

B

# Go-Back-N : Example (2)

A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

Data.req(a)

Data.req(b)

D(0,a)

D(1,b)

Segment lost

B

Data.ind(a)

# Go-Back-N : Example (2)

**A**

Sending window

0 1 2 3

*0* 1 2 3

*0 1* 2 3

*0 1 2* 3

**B**

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

*Segment lost*

Data.ind(a)

*Not expected seq num, discarded*

# Go-Back-N : Example (2)



A

Sending window

0 1 2 3
0 1 2 3
0 1 2 3
0 1 2 3

0 1 2 3

B

Data.req(a)

Data.req(b)

D(0,a)

Data.req(c)

D(1,b)

Segment lost

D(2,c)

Data.ind(a)

C(OK,0)

Not expected seq num, discarded

# Go-Back-N : Example (2)

A

Sending window

0 1 2 3
0 1 2 3
0 1 2 3
0 1 2 3

0 1 2 3

B

Data.req(a)

Data.req(b)

D(0,a)

Data.req(c)

D(1,b)

*Segment lost*

D(2,c)

Data.ind(a)

C(OK,0)

*Not expected seq num, discarded*

C(OK,0)

# Go-Back-N : Example (2)



A

Sending window

`0 1 2` 3
`0 1 2` 3
`0 1 2` 3
`0 1 2` 3

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

*Segment lost*

D(2,c)

Data.ind(a)

0 `1 2 3`

*Retransmission timer expires*

C(OK,0)

*Not expected seq num, discarded*

C(OK,0)

0 `1 2 3`

0 `1 2 3`

D(1,b)

D(2,c)

Data.ind(b)

Data.ind(c)

B

# Go-Back-N : Example (2)



A

Sending window

`0 1 2` 3

`0 1 2` 3

`0 1 2` 3

`0 1 2` 3

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

*Segment lost*

Data.ind(a)

0 `1 2 3`

*Retransmission timer expires*

C(OK,0)

*Not expected seq num, discarded*

Data.req(d)

C(OK,0)

D(1,b)

0 `1 2 3`

0 `1 2 3`

D(2,c)

Data.ind(b)

D(3,d)

Data.ind(c)

Data.ind(d)

B

# Go-Back-N : Example (2)

# Go-Back-N : Example (2)

# Selective Repeat

---

## Receiver

Uses a buffer to store the segments received out of sequence and reorder their content

Receiving window

... 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 ....

Acknowledged segments

Out of sequence segments

Acceptable segments

Segments outside receiving window

# Selective Repeat

## Receiver

Uses a buffer to store the segments received out of sequence and reorder their content

Receiving window



... 0 1 2 3 4 5 **6 7 8 9 10** 11 12 13 14 15 ....

Acknowledged segments

Out of sequence segments

Acceptable segments

Segments outside receiving window

### Semantics of the control segments

OKX

The segments up to and including sequence number X have been received

NAKX

The segment with sequence number X was errored

## Sender

Upon detection of an errored or lost segment, sender retransmits only this segment

may require one retransmission timer per segment

# Selective-Repeat : Receiver

State variable

`next` : sequence number of expected data segment

`Last` : last received in-sequence segment

```
Recvd(D(t,SDU,CRC))
AND NOT(IsOK(CRC,SDU))
        discard(SDU);
send(C(NAK,t,CRC));
```

```
Recvd(D(t,SDU,CRC))
  AND IsOK(CRC,SDU)
insert_in_buffer(SDU);
```

Wait

Process
SDU
OK

For all in sequence segments inside buffer
Data.ind(SDU);
```
slide the sliding window;
update next and last
send(C(OK,(next-1)));
```

# Selective Repeat : Sender

State variables
    `base` : sequence number of oldest unacknowledged segment
    `seq` : first free sequence number
    `W` : size of sending window

```
         Recvd(C(?,?,CRC))
and NOT( CRCOK(C(?,?,CRC)))
```

```
            Data.req(SDU)
     AND ( window not full )
start_timer(seq) ;
insert_in_buffer(SDU);
send(D(seq,SDU,CRC));
seq=(seq+1) ;
```

**Wait**

```
    [ Recvd(C(NAK,t,CRC))
 and CRCOK(C(NAK,t,CRC))]
        or  timer (t) expires
send(D(t,SDU,CRC)); }
restart_timer(t);
```

```
     Recvd(C(OK,t,CRC))
  and CRCOK(C(OK,t,CRC))
For all segments i<=t
    cancel_timer(t);
slide sliding window to
the right;
```

# Selective Repeat : Example

A

Sending window

<span style="color:green">0 1 2</span>3

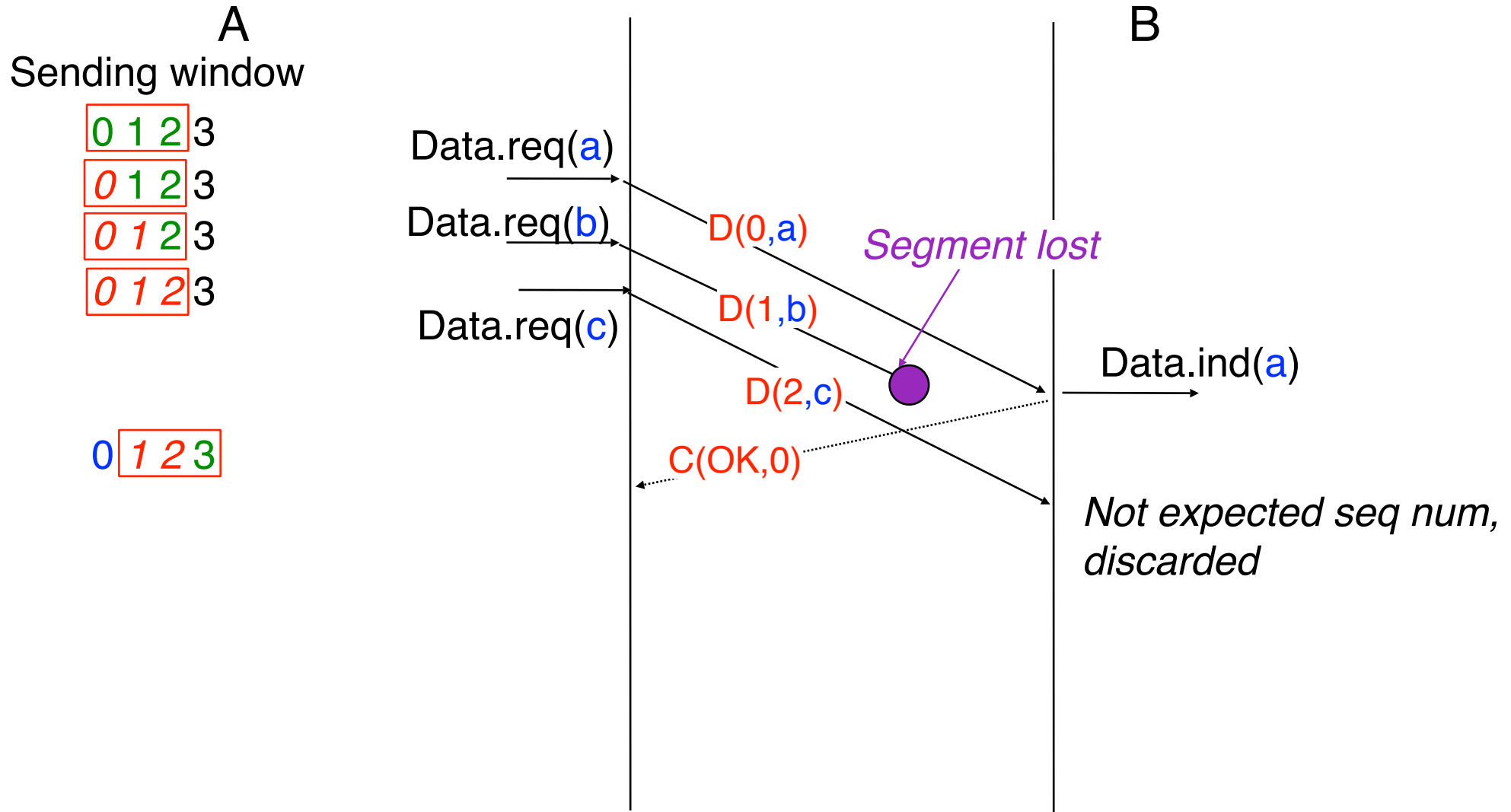<span style="color:blue">0</span> <span style="color:red">*1 2*</span> <span style="color:green">3</span>

B

Rec. window

<span style="color:green">0 1 2</span>3

# Selective Repeat : Example

A

Sending window

0 1 2 3

0 1 2 3

Data.req(a)

D(0,a)

Data.ind(a)

0 1 2 3

B

Rec. window

0 1 2 3

0 1 2 3

# Selective Repeat : Example

A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

Data.req(a)

Data.req(b)

D(0,a)

D(1,b)

*Transmission error*

Data.ind(a)

*Discard segment*

0 1 2 3

B

Rec. window

0 1 2 3

0 1 2 3

CNP3/2008.3.

© O. Bonaventure, 2008

# Selective Repeat : Example

A
Sending window
Rec. window

Data.req(a)
Data.req(b)
Data.req(c)

D(0,a)
D(1,b)
D(2,c)

Transmission error

Data.ind(a)
Discard segment

CNP3/2008.3.

© O. Bonaventure, 2008

# Selective Repeat : Example

# Selective Repeat : Example

# Selective Repeat : Example

# Selective Repeat : Example

# Selective Repeat : Example

# Selective Repeat : Example



A

Sending window

| 0 1 2 | 3
| 0 1 2 | 3
| 0 1 2 | 3
| 0 1 2 | 3

0 | 1 2 3 |

0 | 1 2 3 |

0 | 1 2 3 |

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

*Transmission error*

C(OK,0)

*Retransmission*    C(NAK,1)

C(OK,0)

D(1,b)

C(OK,**2**)

B

Rec. window

| 0 1 2 | 3

Data.ind(a)    0 | 1 2 3 |

*Discard segment*

*Segment stored*   0 | 1 2 3 |

Data.ind(b)   0 | 1 2 3 |

Data.ind(c)   | 0 1 2 3 |

# Selective Repeat : Example



A

Sending window

0 1 2 3
0 1 2 3
0 1 2 3
0 1 2 3

0 1 2 3
0 1 2 3
0 1 2 3

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

Transmission error

Retransmission

Data.req(d)

D(1,b)

D(3,d)

C(OK,0)

C(NAK,1)

C(OK,0)

C(OK,**2**)

B

Rec. window

0 1 2 3

Data.ind(a)

0 1 2 3

Discard segment

Segment stored    0 1 2 3

Data.ind(b)    0 1 2 3
Data.ind(c)    0 1 2 3

Data.ind(d)

0 1 2 3

# Selective Repeat : Example

# Selective Repeat : Example

© O. Bonaventure, 2008

# Selective Repeat : Example (2)

A
Sending window

0 1 2 3

0 1 2 3

B
Rec. window

0 1 2 3

# Selective Repeat : Example (2)

A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

Data.req(a)

D(0,a)

Data.ind(a)

B

Rec. window

0 1 2 3

0 1 2 3

0 1 2 3

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)



A

Sending window

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

0 1 2 3

B

Rec. window

0 1 2 3

Data.req(a)

Data.req(b)

Data.req(c)

D(0,a)

D(1,b)

D(2,c)

Lost segment

Data.ind(a)

0 1 2 3

C(OK,0)

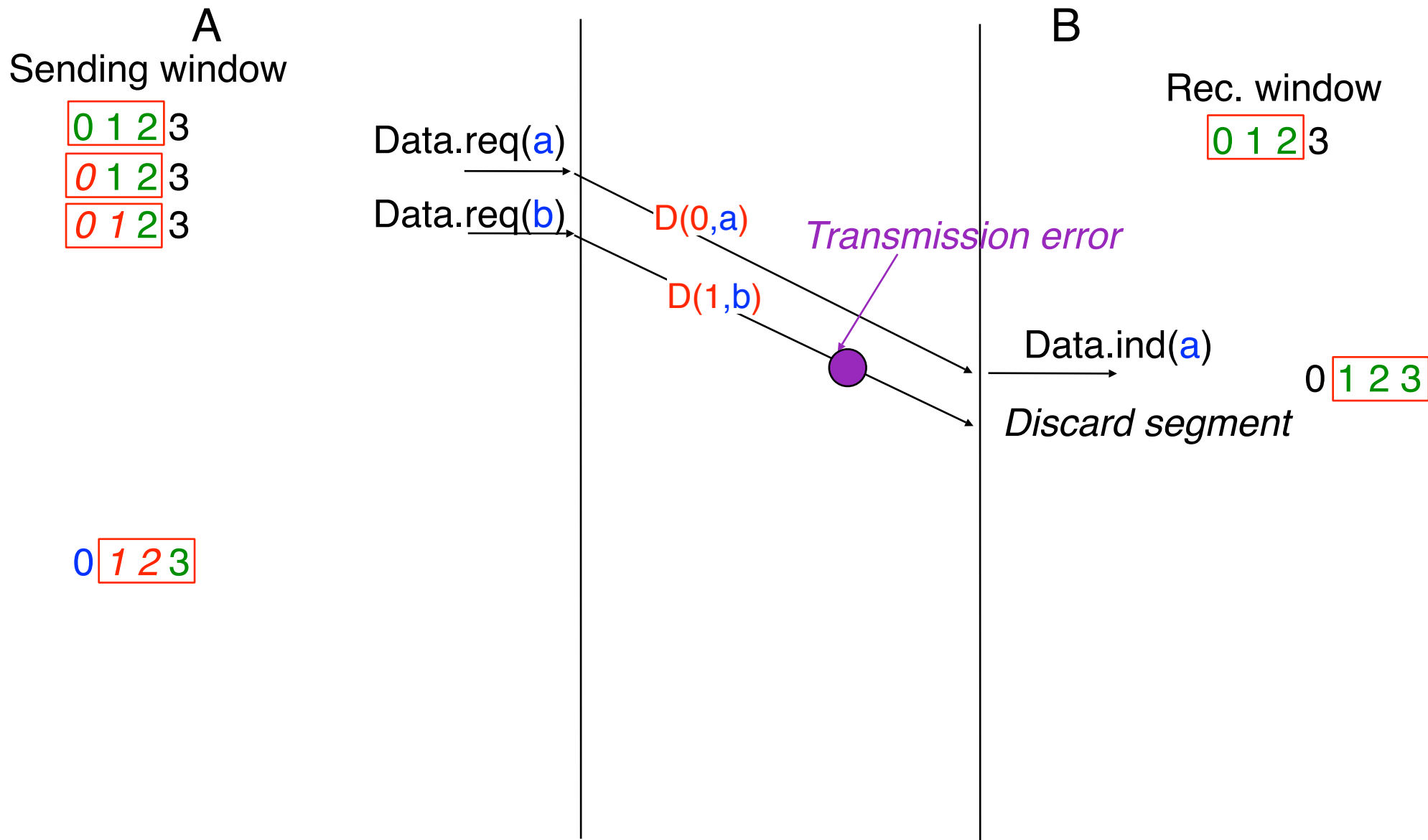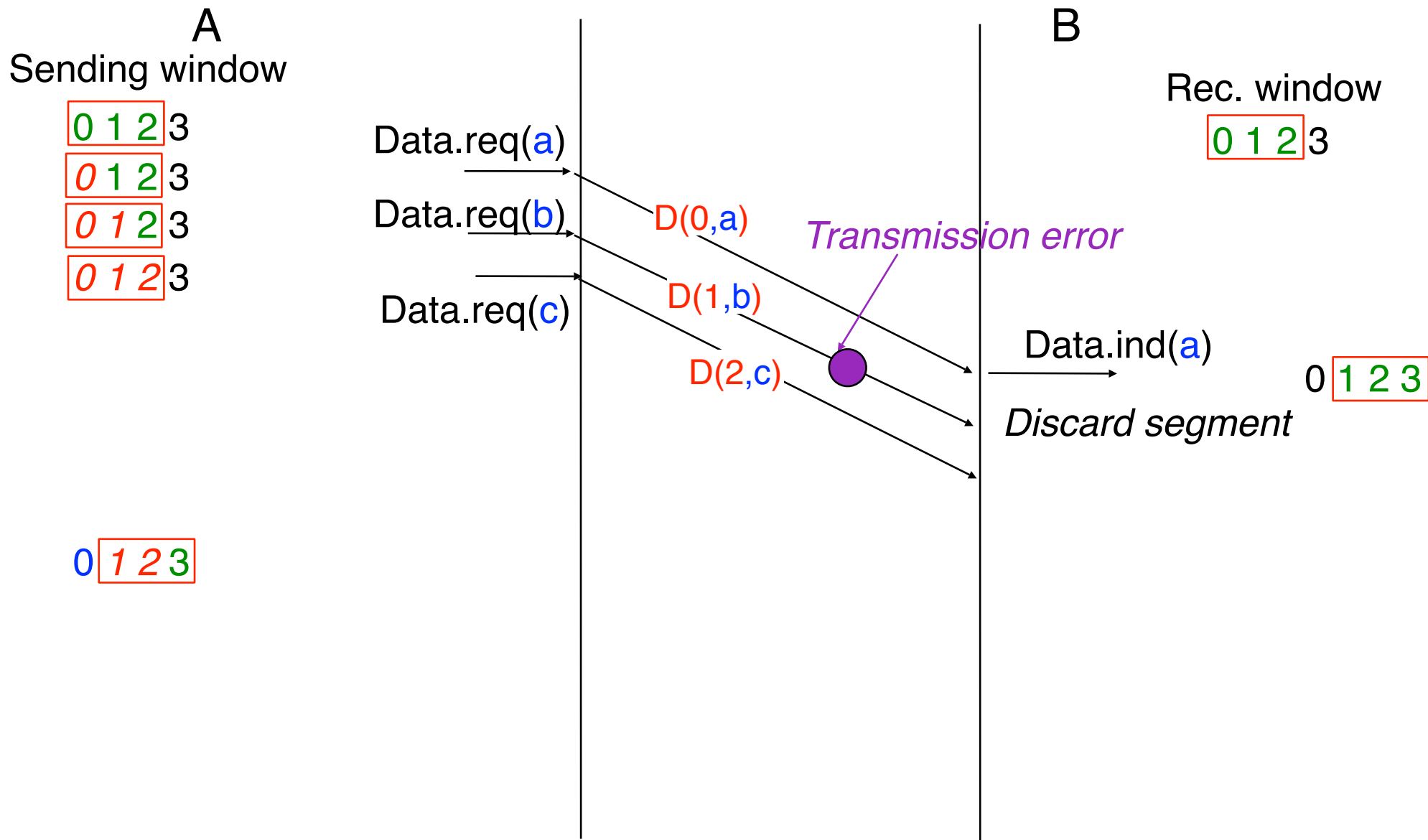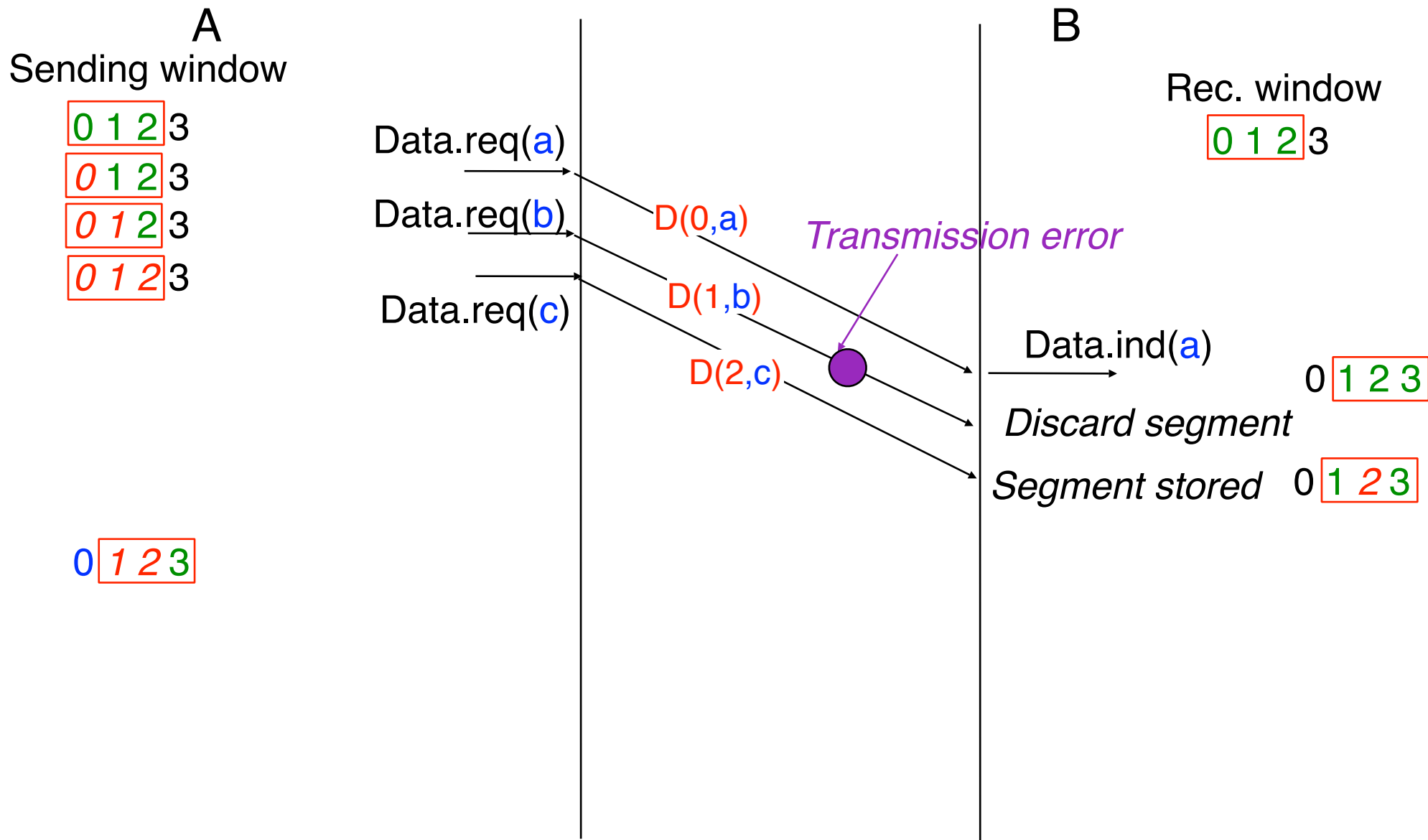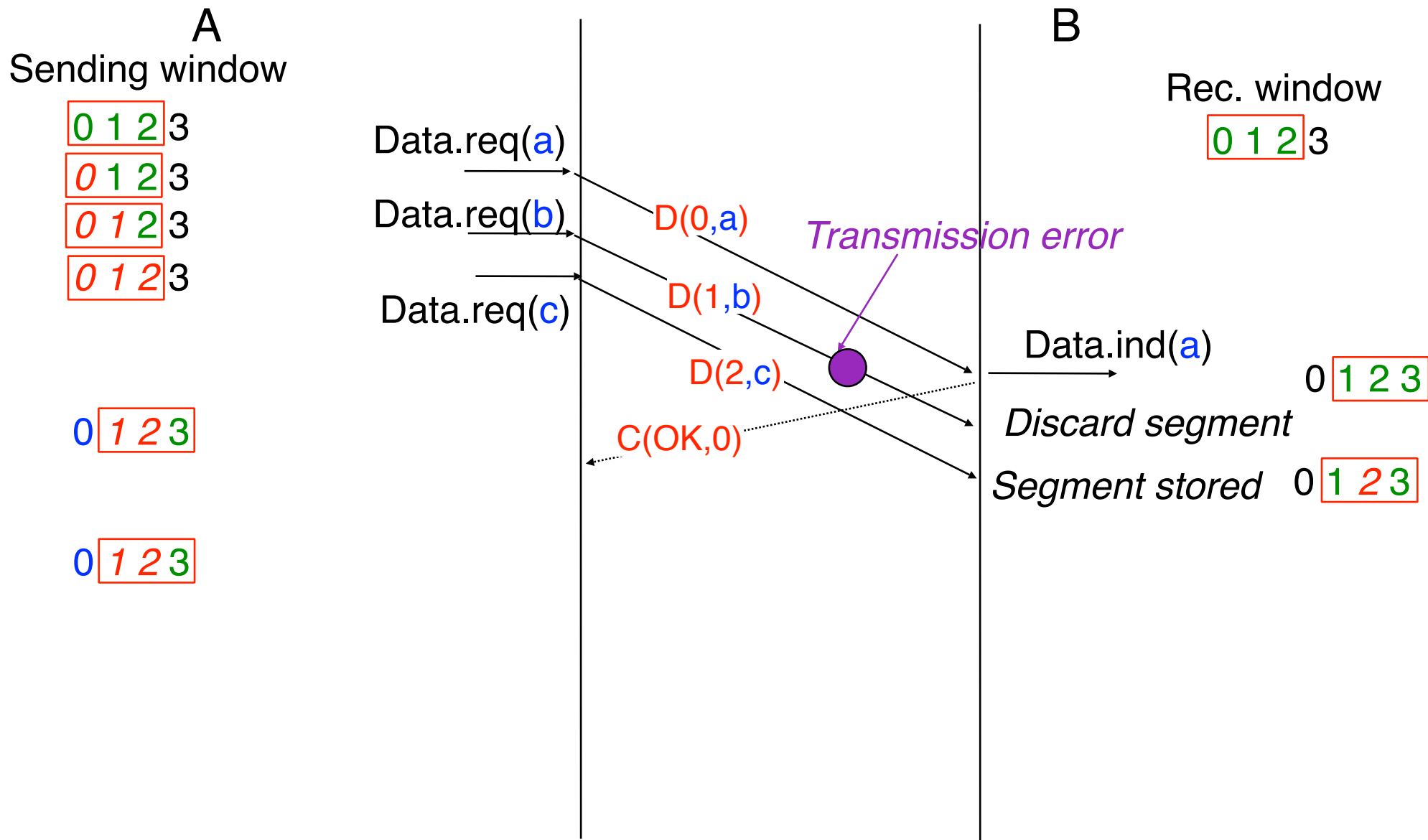Segment stored    0 1 2 3

C(OK,0)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)

# Selective Repeat : Example (2)



CNP3/2008.3.

© O. Bonaventure, 2008

# Buffer management

## Problem



A transport entity may support many transport connections at the same time

How can we share the available buffer among these connections ?

The number of connections changes with time

Some connections require large buffers while others can easily use smaller ones

`ftp` versus `telnet`

# Buffer management (2)

## Principle

Adjust the size of the receiving window according to the amount of buffering available on the receiver

Allow the receiver to advertise its current receiving window size to the sender

## New information carried in control segments

`win` indicates the current receiving window's size

## Changes to sender

Sending window : `swin` (function of available memory)

Keep in a state variable the receiving window advertised by the receiver : `rwin`

At any time, the sender is only allowed to send data segments whose sequence number fits inside

min(`rwin`, `swin`)

# Buffer management (3)

A

Swin=3, rwin=1

`0` 1 2 3

B

Rwin=1

`0` 1 2 3

# Buffer management (3)

A                  B

Swin=3, rwin=1            Rwin=1

0 1 2 3     Data.req(a)

D(0,a)

0 1 2 3

Data.ind(a)

# Buffer management (3)

A                   B

Swin=3, rwin=1             Rwin=1

0 1 2 3      Data.req(a)

D(0,a)

0 1 2 3

Data.ind(a)     0 1 2 3



CNP3/2008.3.

© O. Bonaventure, 2008

# Buffer management (3)



A

Swin=3, rwin=1

[0] 1 2 3

Data.req(a)

D(0,a)

B

Rwin=1

[0] 1 2 3

Data.ind(a)      0 [1] 2 3

C(OK,0, w=1)

# Buffer management (3)



A

Swin=3, rwin=1

0 1 2 3

Data.req(a)

D(0,a)

Swin=3, rwin=1

0 1 2 3

C(OK,0, w=1)

B

Rwin=1

0 1 2 3

Data.ind(a)    0 1 2 3

# Buffer management (3)

# Buffer management (3)

# Buffer management (3)

# Buffer management (3)



A — Swin=3, rwin=1

B — Rwin=1

Data.req(a)
D(0,a)
Data.ind(a)

Swin=3, rwin=1

C(OK,0, w=1)
Data.req(b)
C(OK,0,w=3)
D(1,b)

*2 new buffers become available*

Swin=3, rwin=3

Data.req(c)
D(2,c)
Data.req(d)
C(OK,1,w=3)
Data.ind(b)
D(3,d)

# Buffer management (4)

A

Swin=3, rwin=1

<span style="color:green">0</span> 1 2 3

B

Rwin=1

<span style="color:green">0</span> 1 2 3

*Receiver cannot handle segment immediately*

Data.ind(<span style="color:blue">a</span>)

<span style="color:red">OK(0)</span>

# Buffer management (4)

# Buffer management (4)



A

Swin=3, rwin=1

0 1 2 3

Data.req(a)

D(0,a)

C(OK,0, w=0)

B

Rwin=1

0 1 2 3

*Receiver cannot handle segment immediately*

Data.ind(a)

# Buffer management (4)



A     B

Swin=3, rwin=1     Rwin=1

0 1 2 3

Data.req(a)

D(0,a)

0 1 2 3

*Receiver cannot handle segment immediately*

C(OK,0, w=0)

Data.ind(a)

Swin=3, rwin=0

0 1 2 3

Data.req(b)

*Window blocked*
*No transmission possible*

# Buffer management (4)

A
　　Swin=3, rwin=1

0 1 2 3
　　　　　　　Data.req(a)

D(0,a)

B
Rwin=1

0 1 2 3

*Receiver cannot handle segment immediately*

C(OK,0, w=0)
　　Data.ind(a)

Swin=3, rwin=0
　　　　　　Data.req(b)

0 1 2 3

*Window blocked
No transmission possible*

*Lost segment*

C(OK,0,w=3)

CNP3/2008.3.

© O. Bonaventure, 2008

# Buffer management (4)



A

Swin=3, rwin=1

0 1 2 3

Data.req(a)

D(0,a)

B

Rwin=1

0 1 2 3

*Receiver cannot handle segment immediately*

Data.ind(a)

C(OK,0, w=0)

Swin=3, rwin=0

0 1 2 3

Data.req(b)

*Window blocked*
*No transmission possible*

*Lost segment*

*2 new buffers are available*

C(OK,0,w=3)

0 1 2 3

# Buffer management (4)

# Buffer management (4)

A | B

Swin=3, rwin=1 | Rwin=1

⬚0 1 2 3    Data.req(a)

D(0,a)

⬚0 1 2 3

*Receiver cannot handle segment immediately*

Data.ind(a)

C(OK,0, w=0)

Swin=3, rwin=0    Data.req(b)

0 1 2 3

*Lost segment*    *2 new buffers are available*

*Window blocked*
*No transmission possible*

0 1 ⬚2 3

C(OK,0,w=3)

*Waits for control segment*    *Waits for data segment*

How to recover from deadlock ?
Persitence timer on receiver, resend control segment after
timer expiration

# Duplication and reordering

How can we provide a reliable service in the transport layer ?

Hypotheses

1. The application sends small SDUs
2. The network layer provides a perfect service
    1. Transmission errors are possible
    2. Packets can be lost
    3. Packet reordering is possible
    4. Packets can be duplicated
3. Data transmission is unidirectional

2. How to deal with these problems ?

# Duplication and reordering (2)

## Problem

A late segment could be confused with a valid segment

A B

# Duplication and reordering (2)

Problem

A late segment could be confused with a valid segment

# Duplication and reordering (2)

## Problem
A late segment could be confused with a valid segment

# Duplication and reordering (2)

## Problem
A late segment could be confused with a valid segment

# Duplication and reordering (2)

## Problem

A late segment could be confused with a valid segment

# Duplication and reordering (2)

## Problem

A late segment could be confused with a valid segment

# Duplication and reordering (3)

How to deal with duplication and reordering ?

Possible provided that segments do not remain forever inside the network

Constraint on network layer
- A packet cannot remain inside the network for more than MSL seconds

## Principle of the solution
Only one segment carrying sequence number x can be transmitted during MSL seconds
- upper bound on maximum throughput

# Bidirectional flow

## How can we allow both hosts to transmit data ?

## Principle
Each host sends both control and data segments

### Piggybacking
Place control fields inside the data segments as well (e.g. window, ack number) so that data segments also carry control information
Reduces the transmission overhead



SDU

Type : D or C
CRC
Seq : segment's sequence number
Ack : sequence number of the last received in-order segment

# Bidirectional flow
# Example

A                                                    B

OK(0)

# Bidirectional flow
# Example

# Bidirectional flow
# Example

# Bidirectional flow Example



A

B

Data.req(a)

Data.req(b)

D(0,0,a)

D(1,0,b)

*Error*

Data.req(w)

Data.ind(a)

*Discarded*

Data.ind(w)

D(5,0,w)

*D(5,0,w) acks D(0,0,a)*

# Bidirectional flow Example



A

Data.req(a)

Data.req(b)

Data.req(c)

Data.ind(w)

D(5,0,w) acks D(0,0,a)

D(0,0,a)

D(1,0,b)

D(2,0,c)

D(5,0,w)

*Error*

B

Data.req(w)

Data.ind(a)

*Discarded*

*Segment -> buffer*

# Bidirectional flow Example



Data.req(a)

Data.req(b)

D(0,0,a)

Data.req(c)

D(1,0,b)

*Error*

Data.req(w)

D(2,0,c)

Data.ind(a)

Data.ind(w)

D(5,0,w)

*Discarded*

*D(5,0,w) acks D(0,0,a)*

C(NAK,1)

*Retransmission*

*Segment -> buffer*

D(1,5,b)

Data.ind(b)

Data.ind(c)

C(OK,2)

C(OK,3)

# Bidirectional flow Example

CNP3/2008.3.

© O. Bonaventure, 2008

# Bidirectional flow
# Example

# Byte stream service

## How to provide a byte stream service ?

Principle
- Sender splits the byte stream in segments
- Receiver delivers the payload of the received in-sequence segments to its user
- Usually each octet of the byte stream has its own sequence number and the segment header contains the sequence number of the first byte of the payload
  - In this case, window sizes are often also expressed in bytes

# Byte stream service (2)

A                                                                    B

OK(0)

# Byte stream service (2)



A

B

Data.req(abcdef)

D(0,ab)

Data.ind(ab)

C(OK,1)

# Byte stream service (2)

A                                                      B

Data.req(abcdef)

D(0,ab)

*Lost segment*

D(2,cd)

Data.ind(ab)

C(OK,1)

# Byte stream service (2)

# Byte stream service (2)

# Byte stream service (2)

# Byte stream service (2)

# Byte stream service (2)

# Module 3 : Transport Layer

Basics

<span style="color:red">Building a reliable transport layer</span>
    Reliable data transmission
→   <span style="color:red">Connection establishment</span>
    Connection release

UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

# Transport connection establishment

How to open a transport connection between two transport entities ?

The transport layer uses the imperfect network layer service
- Transmission errors are possible
- Segments can get lost
- Segments can get reordered
- Segments can be duplicated

Hypothesis
- We will first assume that a single transport connection needs to be established between the two transport entities

# Simple solution

**Principle**

   2 control segments

      CR is used to request a connection establishment

      CA is used to acknowledge a connection establishment

   Is this sufficient with an imperfect network layer service ?

# Simple solution

Connect.req       CR       Connect.ind

## Principle
### 2 control segments
CR is used to request a connection establishment
CA is used to acknowledge a connection establishment

### Is this sufficient with an imperfect network layer service ?

# Simple solution

Connect.req     →     CR     →     Connect.ind
Connect.resp

## Principle
### 2 control segments
CR is used to request a connection establishment
CA is used to acknowledge a connection establishment

Is this sufficient with an imperfect network layer service ?

# Simple solution



`Connect.req` — CR → `Connect.ind`

`Connect.conf` ← CA — `Connect.resp`

*Connection established*

*Connection established*

## Principle

2 control segments
   CR is used to request a connection establishment
   CA is used to acknowledge a connection establishment

Is this sufficient with an imperfect network layer service ?

# Simple solution (2)

How to deal with losses and transmission errors ?

Control segments must be protected by CRC or checksum

Retransmission timer is used to protect against segment losses segments

# Simple solution (2)

How to deal with losses and transmission errors ?

- Control segments must be protected by CRC or checksum
- Retransmission timer is used to protect against segment losses segments

CR

Connect.req() ⟶ 🔴

# Simple solution (2)

How to deal with losses and transmission errors ?

Control segments must be protected by CRC or checksum

Retransmission timer is used to protect against segment losses segments



```
                          CR
Connect.req() ───────►             ●

Retransmission    CR
timer expires ───────────────────────────►  Connect.ind()
```

# Simple solution (2)

How to deal with losses and transmission errors ?

- Control segments must be protected by CRC or checksum
- Retransmission timer is used to protect against segment losses segments

```
Connect.req()          CR
                              ●

Retransmission          CR
timer expires
                                      Connect.ind()
                                      Connect.resp

                              Connection established
```

# Simple solution (2)

How to deal with losses and transmission errors ?

Control segments must be protected by CRC or checksum

Retransmission timer is used to protect against segment losses segments

```
Connect.req()              CR

                  Retransmission
                  timer expires    CR
                                              Connect.ind()
                                              Connect.resp

Connect.conf()                            Connection established
                           CA
```

# Simple solution (2)

How to deal with losses and transmission errors ?

Control segments must be protected by CRC or checksum

Retransmission timer is used to protect against segment losses segments

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

Connect.req() → CR → Connect.ind()

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one



```
Connect.req()                          CR
                                                      Connect.ind()
                                                      Connect.resp
Connect.conf()
                                       CA
First connection established              First connection established
```

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one

# Connection establishment

How to deal with duplicated or delayed packets ?

A duplicated CR should not lead to the establishment of two transport connections instead of a single one



`Connect.req()` — CR → `Connect.ind()`

`Connect.resp`

`Connect.conf()` ← CA

*First connection established*

*First connection established*

*First connection stopped*

*First connection stopped*

CR

Old previous CR

CA

How to detect duplicates ?

D

# Connection establishment (2)

How to detect duplicates ?

Principles

The network layer guarantees by its protocols and internal organisation that a packet and its duplicates will not live forever inside the network

No packet will survive more than MSL seconds inside the network

Transport entities rely on a local clock to detect duplicated connection establishment requests

# Connection establishment (3)

## Transport clock

Maintained by each transport entity
usually implemented as a k-bits counter
$2^k$ * clock cycle >> MSL

Must continue to count even if the transport entity stops or reboots

Transport clocks are not synchronised
neither with other transport clocks nor with realtime

# Connection establishment (3)

## Transport clock

### Maintained by each transport entity
usually implemented as a k-bits counter

$2^k$ * clock cycle >> MSL

Must continue to count even if the transport entity stops or reboots

Transport clocks are not synchronised

neither with other transport clocks nor with realtime



Transport clock

$2^k$-1

MSL

Time

# Three way handshake

Host A                                                      Host B

# Three way handshake

Host A                                                    Host B

Sequence number x read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

CR (seq=x)

# Three way handshake

Host A                                                                 Host B

Sequence number x read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

CR (seq=x)

Sequence number y read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(y)

CA  (seq=y, ack=x)

# Three way handshake

Host A                                                    Host B

Sequence number x read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

CR (seq=x)

Sequence number y read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(y)

CA  (seq=y, ack=x)

Received CA acknowledges CR
Send CA to ack received CA
**Local state :**
Connection to B :
- established
- current_seq = x

CA (seq=x, ack=y)

*Connection established*

# Three way handshake

Host A                                                            Host B

Sequence number x read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

Received CA acknowledges CR
Send CA to ack received CA
**Local state :**
Connection to B :
- established
- current_seq = x

*Connection established*

CR (seq=x)

Sequence number y read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(y)

CA  (seq=y, ack=x)

CA (seq=x, ack=y)

**Local state :**
Connection to A :
- established
- current_seq=y

*Connection established*

# Three way handshake

**Host A**

**Host B**

Sequence number x read from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

CR (seq=x)

Sequence number y read from local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(y)

Received CA acknowledges CR
Send CA to ack received CA
**Local state :**
Connection to B :
- established
- current_seq = x

CA (seq=y, ack=x)

*Connection established*

CA (seq=x, ack=y)

**Local state :**
Connection to A :
- established
- current_seq=y

*Connection established*

D(x)

The sequence numbers used for the data segments will start from x

# Three way handshake

**Host A**                                    **Host B**

Sequence number x read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (x)
- Start retransmission timer

CR (seq=x)

Sequence number y read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(y)

Received CA acknowledges CR
Send CA to ack received CA
**Local state :**
Connection to B :
- established
- current_seq = x

CA (seq=y, ack=x)

CA (seq=x, ack=y)

**Local state :**
Connection to A :
- established
- current_seq=y

*Connection established*

***Connection established***

The sequence numbers used
for the data segments will start
from x

D(x)

The sequence numbers
used for the data segments
will start from y

D(y)

## What happens with duplicates
### Duplicate CR

Host A

Host B

# Three way handshake (2)

What happens with duplicates
Duplicate CR

Host A

CR (seq=z)

Host B

# Three way handshake (2)

## What happens with duplicates
### Duplicate CR

Host A | CR (seq=z) | Host B

Sequence number y read from local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(y)

## What happens with duplicates
### Duplicate CR

Host A                                                          Host B

CR (seq=z)

Sequence number y read from
local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(y)

CA  (seq=y, ack=z)

# Three way handshake (2)

## What happens with duplicates
### Duplicate CR

Host A

CR (seq=z)

Host B

CA (seq=y, ack=z)

Sequence number y read from local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(y)

**Local state :**
No connection to B
Send REJECT to cancel connection establishment

# Three way handshake (2)

## What happens with duplicates
### Duplicate CR

Host A    CR (seq=z)    Host B

Sequence number y read from
local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(y)

CA (seq=y, ack=z)

**Local state :**
No connection to B
Send REJECT to cancel
connection establishment

REJECT (ack=y)

Connection cancelled

*No connection is established*

# Three way handshake (3)

Host A                                                    Host B

Current state does not contain
a CR with seq=x

# Three way handshake (3)

Sequence number $z$ read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR ($z$)
- Start retransmission timer

Current state does not contain
a CR with seq=$x$

CR (seq=$z$)

# Three way handshake (3)

Sequence number $z$ read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR ($z$)
- Start retransmission timer

Current state does not contain
a CR with seq=$x$

CR (seq=$z$)

CA  (seq=$y$, ack=$x$)

# Three way handshake (3)

Sequence number $z$ read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR ($z$)
- Start retransmission timer

CR (seq=$z$)

Current state does not contain
a CR with seq=$x$

CA  (seq=$y$, ack=$x$)

REJECT (ack=$y$)

Current state does not contain
a segment with seq=$y$
REJECT ignored

# Three way handshake (3)

Host A                                                                                    Host B

Sequence number z read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (z)
- Start retransmission timer

CR (seq=z)

Current state does not contain
a CR with seq=x

CA  (seq=y, ack=x)

REJECT (ack=y)

Current state does not contain
a segment with seq=y
REJECT ignored

Retransmission timer
expires

CR (seq=z)

Sequence number w read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(w)

CA  (seq=w, ack=z)

# Three way handshake (3)

Sequence number z read
from local transport clock
**Local state :**
Connection to B :
- Wait for ack for CR (z)
- Start retransmission timer

Current state does not contain
a CR with seq=x

CR (seq=z)

CA (seq=y, ack=x)

REJECT (ack=y)

Current state does not contain
a segment with seq=y
REJECT ignored

Retransmission timer
expires

CR (seq=z)

Received CA acknowledges CR
Send CA to ack received CA
**Local state :**
Connection to B :
- established
- current_seq = z

CA (seq=w, ack=z)

Sequence number w read from
local transport clock
CA sent to ack CR
**Local state :**
Connection to A :
- Wait for ack for CA(w)

CA (seq=z, ack=w)

*Connection established*

    

# Three way handshake (4)

Another scenario

Host A                                                        Host B

# Three way handshake (4)

## Another scenario

Host A                                                           Host B

CR (seq=z)

# Three way handshake (4)

## Another scenario

Host A                                                                Host B

CR (seq=z)

Sequence number w read from
local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(w)

CA (seq=w, ack=z)

# Three way handshake (4)

## Another scenario

Host A                                                          Host B

CR (seq=z)

Sequence number w read from
local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(w)

Current state does not contain
a CR with seq=z

CA (seq=w, ack=z)

REJECT (ack=w)

# Three way handshake (4)

## Another scenario



Host A

Host B

CR (seq=z)

Sequence number w read from local transport clock
Acknowledges CR segment
**Local state :**
Connection to A :
- Wait for ack for CA(w)

CA (seq=w, ack=z)

Current state does not contain a CR with seq=z

REJECT (ack=w)

CA (seq=z, ack=y)

Invalid CA received from A
Send REJECT

REJECT (ack=z)

*No connection is established*

# Module 3 : Transport Layer

Basics

Building a reliable transport layer
  Reliable data transmission
  Connection establishment
→ Connection release

UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

# Connection release

A transport connection can be used in both directions

Types of connection release

Abrupt connection release
One of the transport entities closes both directions of data transfer
can lead to losses of data

Graceful release
Each transport entity closes its own direction of data transfer
connection will be closed once all data has been correctly delivered

# Abrupt release

**Connection closed**

**Connection closed**

# Abrupt release

# Abrupt release

# Abrupt release

# Abrupt release



CR (seq=z)

CA (seq=w, ack=z)

CA (seq=z, ack=w)

Data.req()   D

Data.req()   D    Data.ind()

Disc.req()

Disc.req()   Connection closed

DR

**Connection closed**

**These segments will not be delivered !**

# Abrupt release (2)

A transport layer entity may itself be forced to release a transport connection

    the same data segment has been transmitted multiple times without receiving an acknowledgement
the network layer reports that the destination host is not reachable anymore
the transport layer entity does not have enough resources available to support this connection (e.g. not enough memory)

In this case, the transport layer entity will perform an abrupt disconnection

# Graceful shutdown

Principle
Each entity closes its own direction of data transfer once all its data have been sent

# Graceful shutdown

Principle
Each entity closes its own direction of data transfer once all its data have been sent

DISCONNECT.req (A-B)

DR(A-B)

DISCONNECT.ind(A-B)

# Graceful shutdown

## Principle
Each entity closes its own direction of data transfer once all its data have been sent

DISCONNECT.req (A-B)

DR(A-B)

DISCONNECT.ind(A-B)

*Incoming connection (A->B) closed*

# Graceful shutdown

## Principle
Each entity closes its own direction of data transfer once all its data have been sent

DISCONNECT.req (A-B)

DR(A-B)

DISCONNECT.ind(A-B)

*Incoming connection (A->B) closed*

ACK

DISCONNECT.conf(A-B)

*Outgoing connection (A->B) closed*

# Graceful shutdown

## Principle
Each entity closes its own direction of data transfer once all its data have been sent

DISCONNECT.req (A-B)

DR(A-B)

DISCONNECT.ind(A-B)

*Incoming connection (A->B) closed*

ACK

DISCONNECT.conf(A-B)

*Outgoing connection (A->B) closed*

DISCONNECT.req(B-A)

DR(B-A)

DISCONNECT.ind(B-A)

*Incoming connection (B->A) closed*

DISCONNECT.conf(A-B)

ACK

*Outgoing connection (B->A) closed*

# Reliability of the transport layer

## Limitations

Transport layer provides a reliable data transfer
<span style="color:red">during the lifetime of the transport connection</span>
<span style="color:red">If a connection is gracefully shutdown, then all the data sent of this connection have been received correctly</span>
data transfer may be unreliable (e.g. loss of segments) if the connection is abruptly released

## Transport layer does not recover itself from abrupt connection releases

Possible solutions
Application reopens the connection and restarts the data transfer
Session Layer
Transaction processing

# Module 3 : Transport layer

Basics

Building a reliable transport layer

→ UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

# A simple transport protocol

## User Datagram Protocol (UDP)
The simplest transport protocol

## Goal
Allow applications to exchange small SDUs by relying on the IP service
- on most operating systems, sending raw IP packets requires special privileges while any application can use directly the transport service

## Constraint
The implementation of the UDP transport entity should remain as simple as possible

# UDP : design choices

Which mechanisms inside UDP ?

Application identification
Several applications running on the same host must be able to use the UDP service

Solution
Source port to identify sending application
Destination port to identify receiving application
Each UDP segment contains both the source and the destination ports

Detection of transmission errors

# UDP protocol

2 UDP entities exchange UDP segments

UDP segment format



32 bits

| Source Port | Destination port |
|---|---|
| UDP length | UDP Checksum |
| Payload | |

8 bytes

# UDP protocol

2 UDP entities exchange UDP segments

UDP segment format

Used to identify the application that sent this segment on sending host

| 32 bits | |
|---|---|
| *Source Port* | *Destination port* |
| UDP length | **UDP Checksum** |
| Payload | |

8 bytes

# UDP protocol

## 2 UDP entities exchange UDP segments

### UDP segment format

Used to identify the application that sent this segment on sending host

Used to identify the application that will receive this segment on destination host

32 bits

8 bytes

| Source Port | Destination port |
|---|---|
| UDP length | UDP Checksum |
| Payload | |

# UDP protocol

## 2 UDP entities exchange UDP segments

### UDP segment format

Used to identify the application that will receive this segment on destination host

Used to identify the application that sent this segment on sending host

32 bits

| Source Port | Destination port |
|---|---|
| UDP length | **UDP Checksum** |
| Payload | |

8 bytes

Checksum computed over the entire UDP segment and part of the IP header to detect transmission errors. 0 means that the sender did not compute a checksum.

# UDP protocol

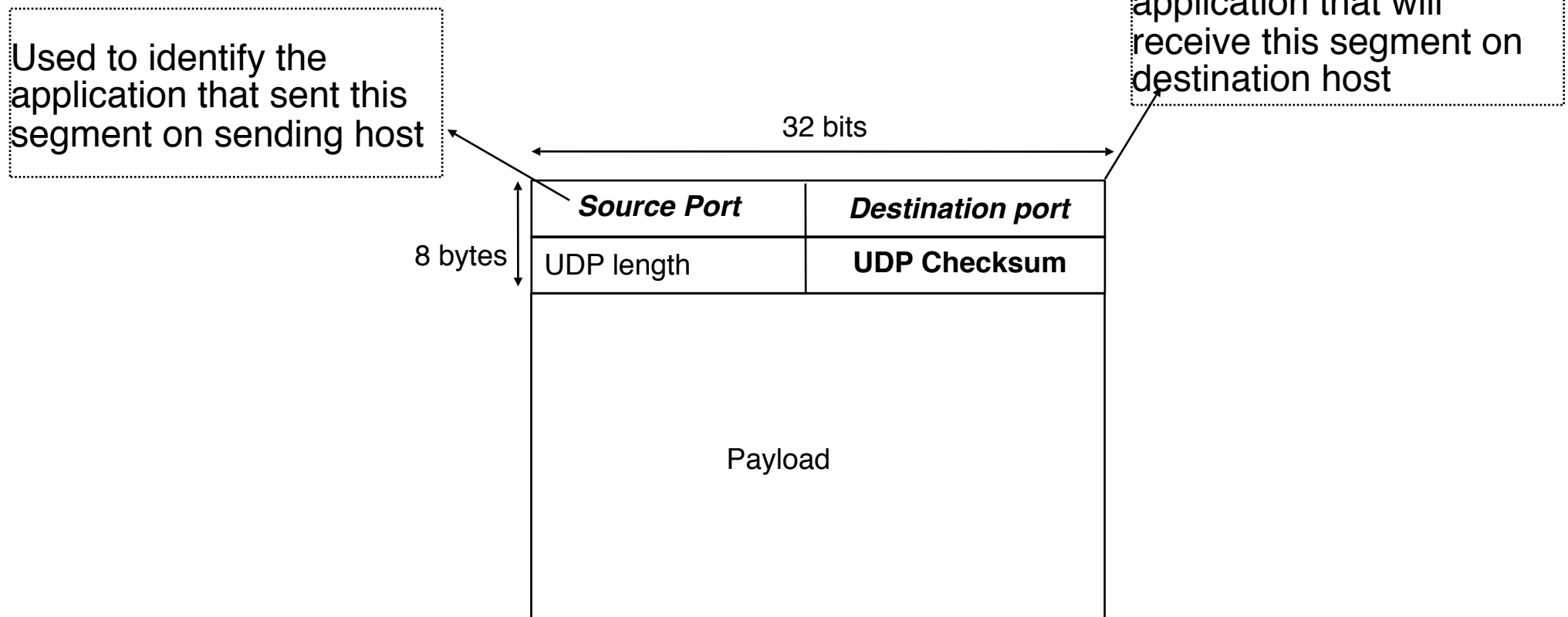## 2 UDP entities exchange UDP segments

## UDP segment format

Used to identify the application that sent this segment on sending host

Used to identify the application that will receive this segment on destination host

32 bits

| Source Port | Destination port |
|---|---|
| UDP length | UDP Checksum |
| Payload | |

8 bytes

Checksum computed over the entire UDP segment and part of the IP header to detect transmission errors. 0 means that the sender did not compute a checksum.

Constraint
Each UDP segment must fit inside a single IP packet

# UDP Protocol (2)

## Utilisation of the UDP ports

Request →

Client

Source port        : 1234
Destination port: 5678

Server

# UDP Protocol (2)

## Utilisation of the UDP ports

Request →

Client

```
 Source port        : 1234
Destination port: 5678
```

Server

```
Source port        : 5678
Destination port: 1234
```

← Response

# Limitations of the UDP service

## Limitations

Maximum length of UDP SDUs depends on maximum size of IP packets

Unreliable connectionless service
SDUs can get lost but transmission errors will be detected

UDP does not preserve ordering

UDP does not detect nor prevent duplication

# Usage of UDP

Request-response applications where requests and responses are short and short delay is required or used in LAN environments
- DNS
- Remote Procedure Call
- NFS
- Games

Multimedia transfer were reliable delivery is not necessary and retransmissions would cause too long delays
- Voice over IP
- Video over IP

# Module 3 : Transport Layer

Basics

Building a reliable transport layer

UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol

→        TCP connection establishment
       TCP connection release
       Reliable data transfer
       Congestion control

# TCP

Transmission Control Protocol

Provides a reliable byte stream service

Characteristics of the TCP service

- TCP connections
- Data transfer is reliable
  - no loss
  - no errors
  - no duplications
- Data transfer is bidirectional
- TCP relies on the IP service
- TCP only supports unicast

# TCP connection

## How to identify a TCP connection

Address of the source application
- IP Address of the source host
- TCP port number of the application on source host

Address of the destination application
- IP Address of the destination host
- TCP port number of the application on destination host

## Each TCP segment contains the identification of the connection it belongs to
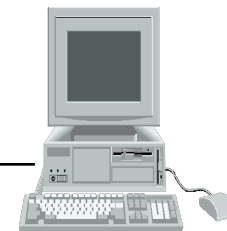
# TCP connection (2)

## Usage of TCP port numbers

Request →

**Client : C**

Source Port        : 1234
Destination Port: 5678

**Server : S**

Source Port        : 5678
Destination Port: 1234

← Response

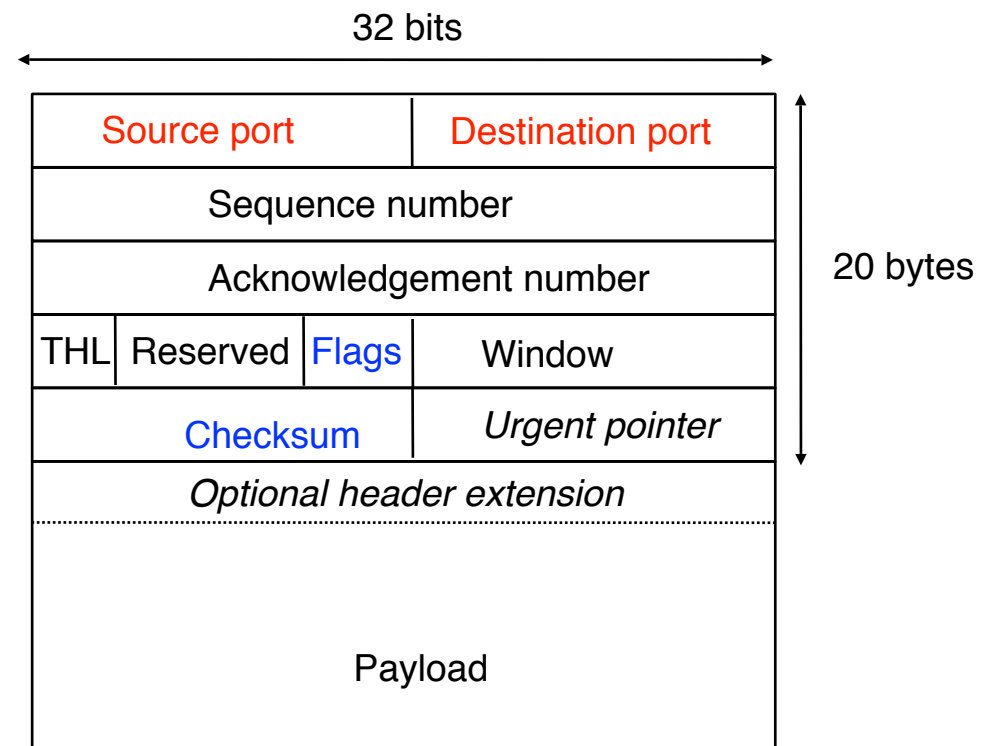| Established TCP connections on client | | | |
|---|---|---|---|
| Local IP | Remote IP | Local Port | Remote Port |
| C | S | 1234 | 5678 |

| Established TCP connections on server | | | |
|---|---|---|---|
| Local IP | Remote IP | Local Port | Remote Port |
| S | C | 5678 | 1234 |

# TCP protocol

## Single segment format

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgement number | |

20 bytes

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

*Optional header extension*

Payload
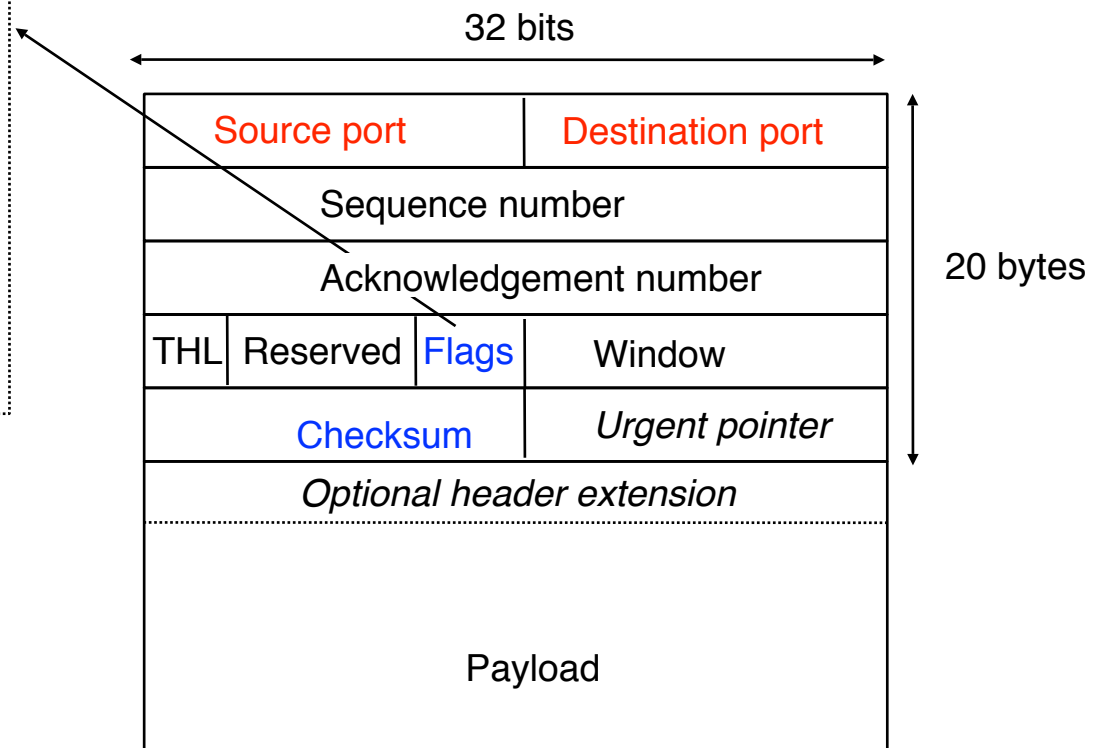
# TCP protocol

## Single segment format

Flags :
used to indicate the function of a segment
**SYN :** used during establishment
**FIN :** used during connection release
**RST :** used in case of problems
**ACK :** if true, means that the Acknowledgement number inside the segment is valid

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgement number | |

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|
| Optional header extension | |

Payload
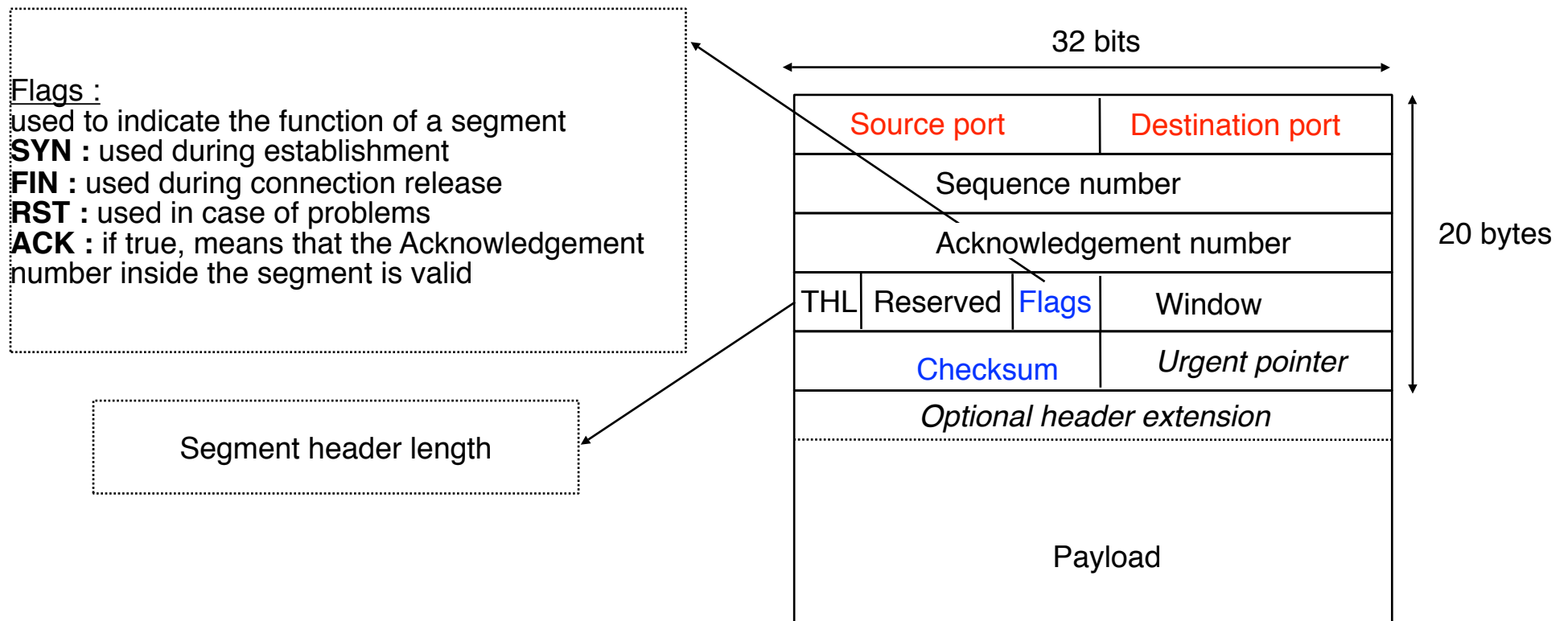
20 bytes

# TCP protocol

## Single segment format

Flags :
used to indicate the function of a segment
**SYN :** used during establishment
**FIN :** used during connection release
**RST :** used in case of problems
**ACK :** if true, means that the Acknowledgement
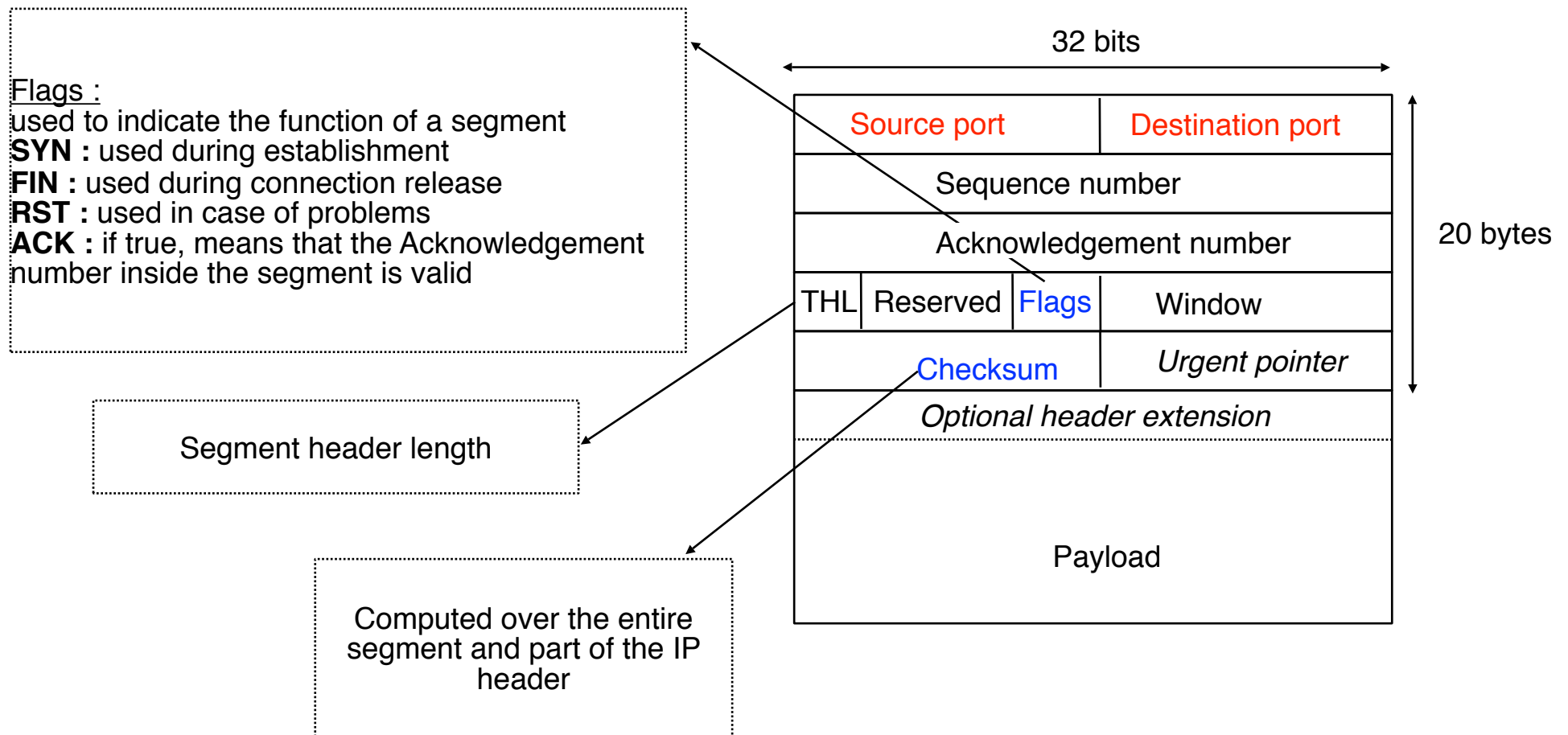number inside the segment is valid

Segment header length

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgement number | |

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

*Optional header extension*

Payload

20 bytes

# TCP protocol

## Single segment format

Flags :
used to indicate the function of a segment
**SYN :** used during establishment
**FIN :** used during connection release
**RST :** used in case of problems
**ACK :** if true, means that the Acknowledgement
number inside the segment is valid

Segment header length

Computed over the entire segment and part of the IP header

32 bits

| Source port | Destination port |
|---|---|
| Sequence number | |
| Acknowledgement number | |

| THL | Reserved | Flags | Window |
|---|---|---|---|

| Checksum | Urgent pointer |
|---|---|

Optional header extension

Payload

20 bytes

# TCP connection establishment

## Three-way handshake

# TCP connection establishment

## Three-way handshake

CONNECT.req

Initial sequence number (x)
read from TCP's clock

SYN(seq=x)

CONNECT.ind

# TCP connection establishment

## Three-way handshake



CONNECT.req

Initial sequence number ($x$)
read from TCP's clock

SYN(seq=$x$)

CONNECT.ind
CONNECT.resp

CONNECT.conf

SYN+ACK(ack=$x+1$,seq=$y$)

*Connection established*

Initial sequence number ($y$)
read from TCP's clock

# TCP connection establishment

## Three-way handshake



CONNECT.req

Initial sequence number (x)
read from TCP's clock

SYN(seq=x)

CONNECT.ind
CONNECT.resp

CONNECT.conf

SYN+ACK(ack=x+1,seq=y)

Initial sequence number (y)
read from TCP's clock

*Connection established*

The sequence numbers of all
segments A->B will start at x+1

ACK(seq=x+1, ack=y+1)

*Connection established*

The sequence numbers of all
segments B->A will start at y+1

Remarks
   Setting the SYN flag in a segment consumes one sequence number
   The ACK flag is set only when the acknowledgement field contains a
   valid value
   The default recommendation for the TCP clock is to be incremented by
   1 at least after 4 microseconds and after each TCP connection
   establishment

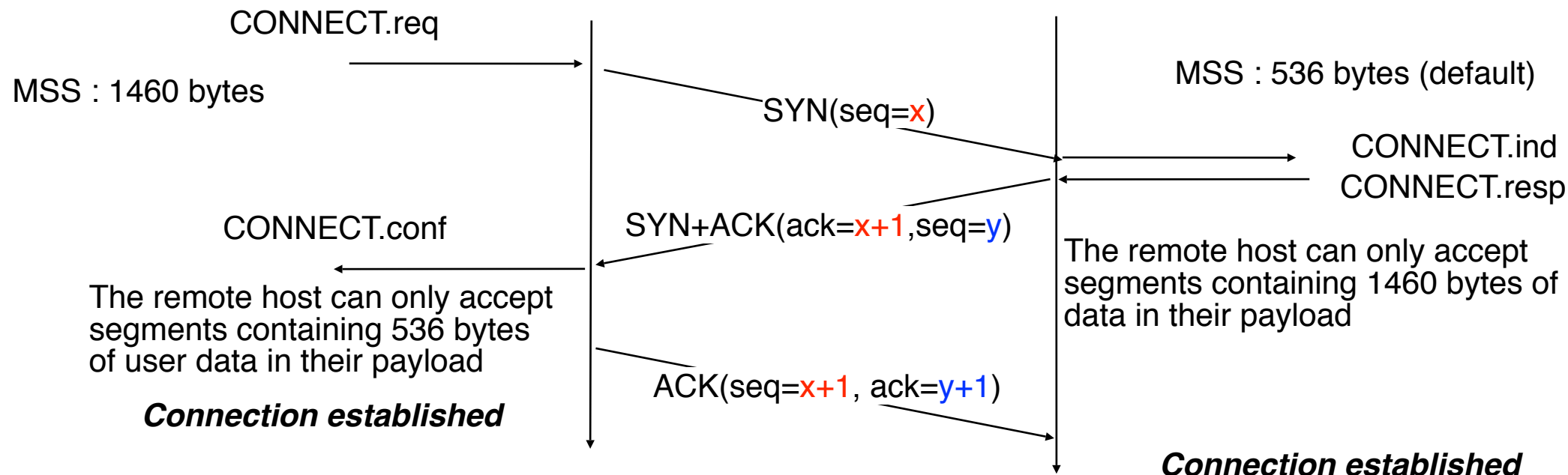# TCP connection establishment (2)

## Option negotiation

### During the opening of a connection, it is possible to negotiate the utilisation of TCP extensions

#### Option encoded inside the optional part of TCP header

Maximum segment size (MSS)
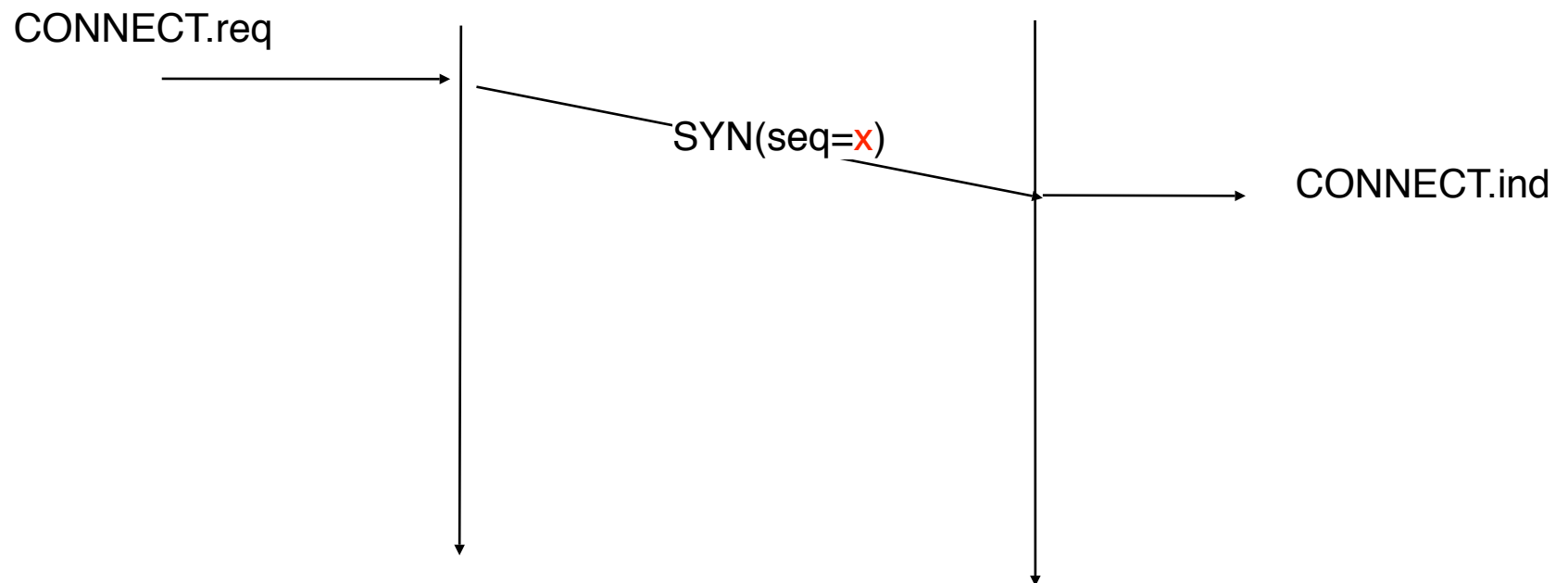RFC1323 timestamp extensions
Selective Acknowledgments

MSS : 1460 bytes

MSS : 536 bytes (default)

# TCP connection establishment (2)

## Option negotiation
### During the opening of a connection, it is possible to negotiate the utilisation of TCP extensions
#### Option encoded inside the optional part of TCP header
##### Maximum segment size (MSS)
##### RFC1323 timestamp extensions
##### Selective Acknowledgments

CONNECT.req

MSS : 1460 bytes

MSS : 536 bytes (default)

SYN(seq=x)

CONNECT.ind

The remote host can only accept segments containing 1460 bytes of data in their payload

# TCP connection establishment (2)

Option negotiation
During the opening of a connection, it is possible to negotiate the utilisation of TCP extensions
Option encoded inside the optional part of TCP header
Maximum segment size (MSS)
RFC1323 timestamp extensions
Selective Acknowledgments

CONNECT.req

MSS : 1460 bytes

MSS : 536 bytes (default)

SYN(seq=x)

CONNECT.ind
CONNECT.resp

CONNECT.conf

SYN+ACK(ack=x+1,seq=y)

The remote host can only accept segments containing 1460 bytes of data in their payload

The remote host can only accept segments containing 536 bytes of user data in their payload

*Connection established*

# TCP connection establishment (2)

Option negotiation

During the opening of a connection, it is possible to negotiate the utilisation of TCP extensions

Option encoded inside the optional part of TCP header

Maximum segment size (MSS)

RFC1323 timestamp extensions

Selective Acknowledgments

CONNECT.req

MSS : 1460 bytes

MSS : 536 bytes (default)

SYN(seq=x)

CONNECT.ind
CONNECT.resp

CONNECT.conf

SYN+ACK(ack=x+1,seq=y)

The remote host can only accept segments containing 1460 bytes of data in their payload

The remote host can only accept segments containing 536 bytes of user data in their payload

*Connection established*

ACK(seq=x+1, ack=y+1)

*Connection established*

# TCP connection establishment (3)

## Rejection of connection establishment

© O. Bonaventure, 2008

# TCP connection establishment (3)

## Rejection of connection establishment



CONNECT.req

SYN(seq=x)

CONNECT.ind

# TCP connection establishment (3)

## Rejection of connection establishment

# TCP connection establishment (3)

## Rejection of connection establishment



A TCP entity should never send a RST segment
upon reception of another RST segment

# TCP connection establishment (4)

## Simultaneous establishment

# TCP connection establishment (4)

## Simultaneous establishment

CONNECT.req

SYN(seq=x)

# TCP connection establishment (4)

## Simultaneous establishment

CONNECT.req

SYN(seq=x)

SYN(seq=y)

CONNECT.req

# TCP connection establishment (4)

## Simultaneous establishment



CONNECT.req

SYN(seq=x)

SYN(seq=y)

CONNECT.req

SYN+ACK(seq=x, ack=y+1)

CONNECT.conf

*Connection established*

# TCP connection establishment (4)

## Simultaneous establishment

CONNECT.req

CONNECT.req

SYN(seq=x)

SYN(seq=y)

SYN+ACK(seq=x, ack=y+1)

SYN+ACK(seq=y, ack=x+1)

CONNECT.conf

CONNECT.conf

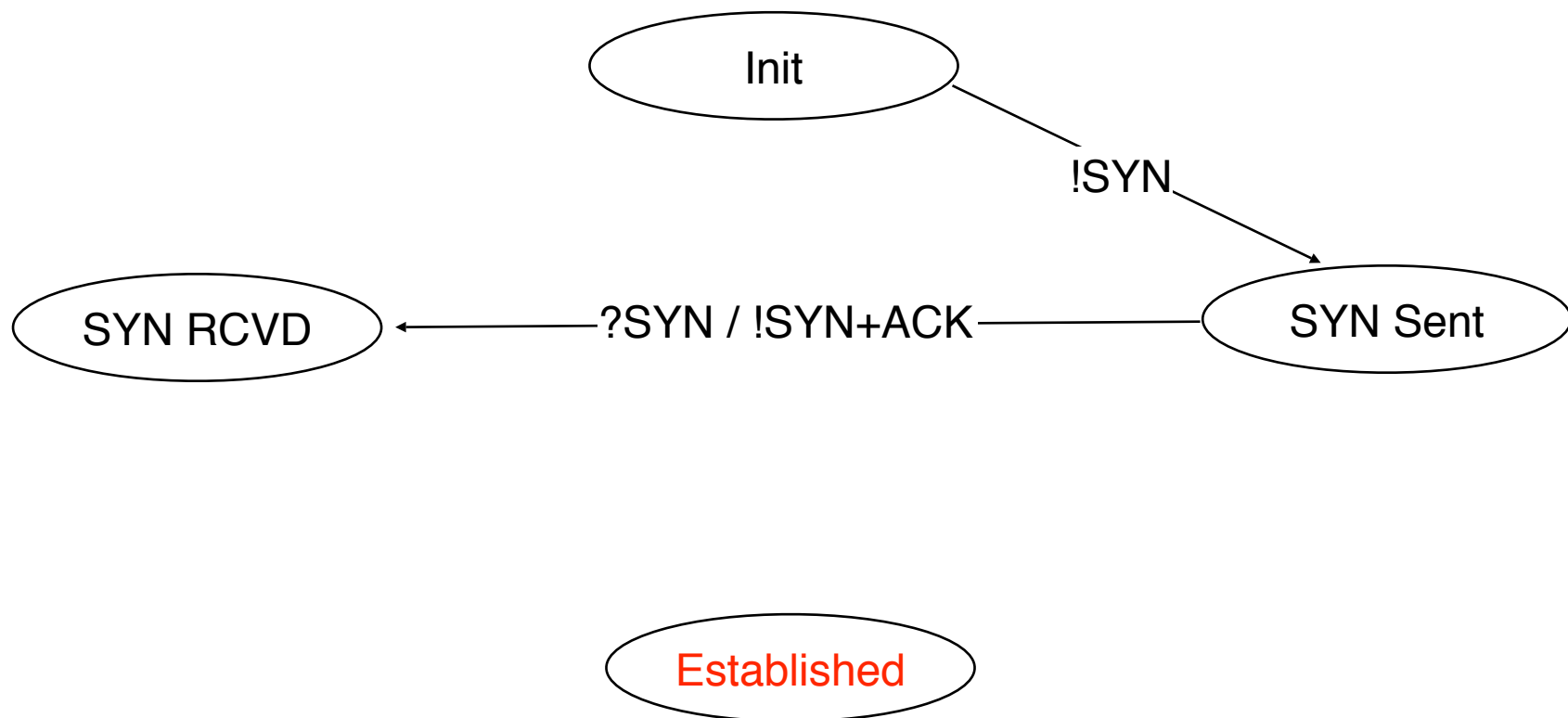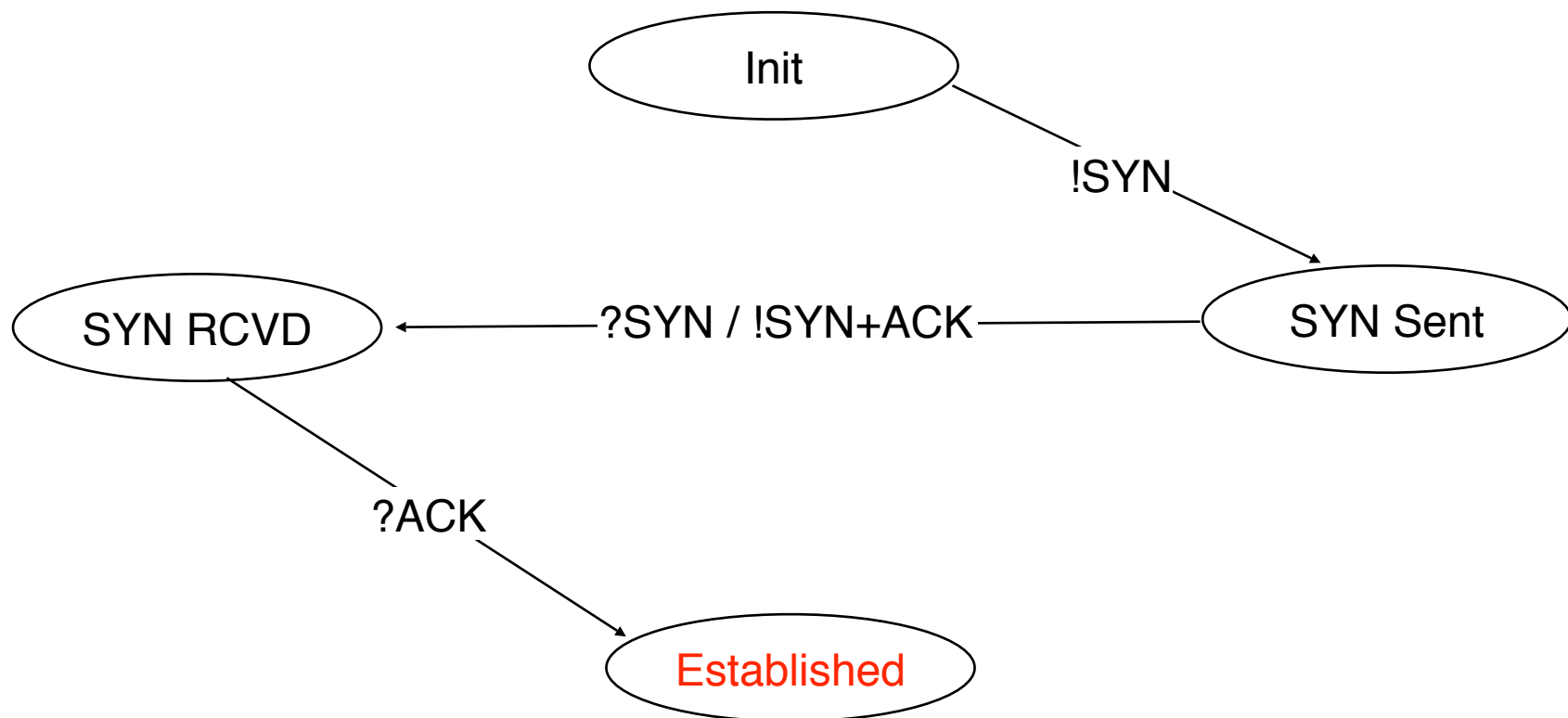*Connection established*

*Connection established*

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (5)

## Representation as a finite state machine

## Representation as a finite state machine

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (5)
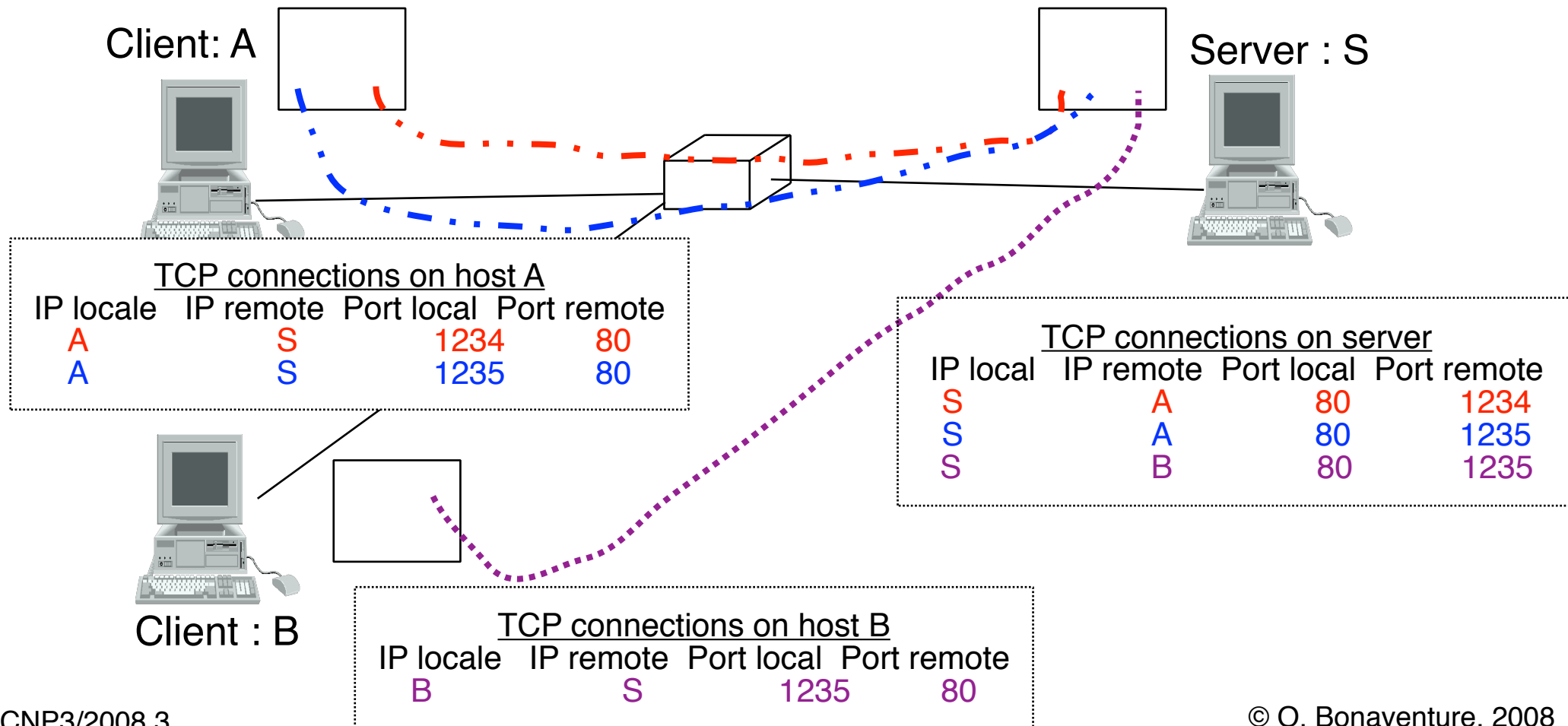
## Representation as a finite state machine

Init

SYN RCVD

SYN Sent

Established

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (5)

## Representation as a finite state machine

# TCP connection establishment (6)

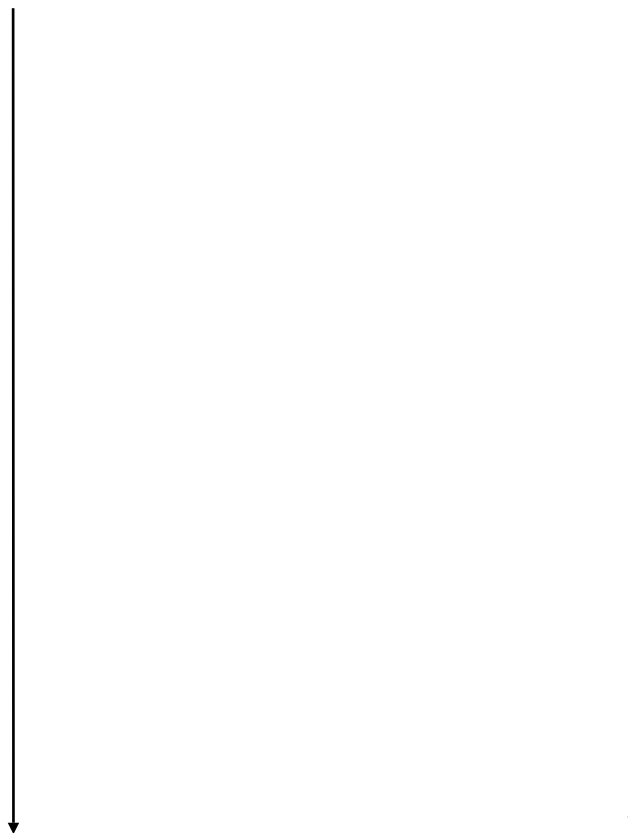How to open several TCP connections at the same time ?



**Client: A**

**Server : S**

### TCP connections on host A

| IP locale | IP remote | Port local | Port remote |
|-----------|-----------|------------|-------------|
| A | S | 1234 | 80 |
| A | S | 1235 | 80 |

### TCP connections on server

| IP local | IP remote | Port local | Port remote |
|----------|-----------|------------|-------------|
| S | A | 80 | 1234 |
| S | A | 80 | 1235 |
| S | B | 80 | 1235 |

**Client : B**

### TCP connections on host B

| IP locale | IP remote | Port local | Port remote |
|-----------|-----------|------------|-------------|
| B | S | 1235 | 80 |

# Module 3 : Transport Layer

Basics

Building a reliable transport layer

UDP : a simple connectionless transport protocol

TCP : a reliable connection oriented transport protocol
>    TCP connection establishment
→    TCP connection release
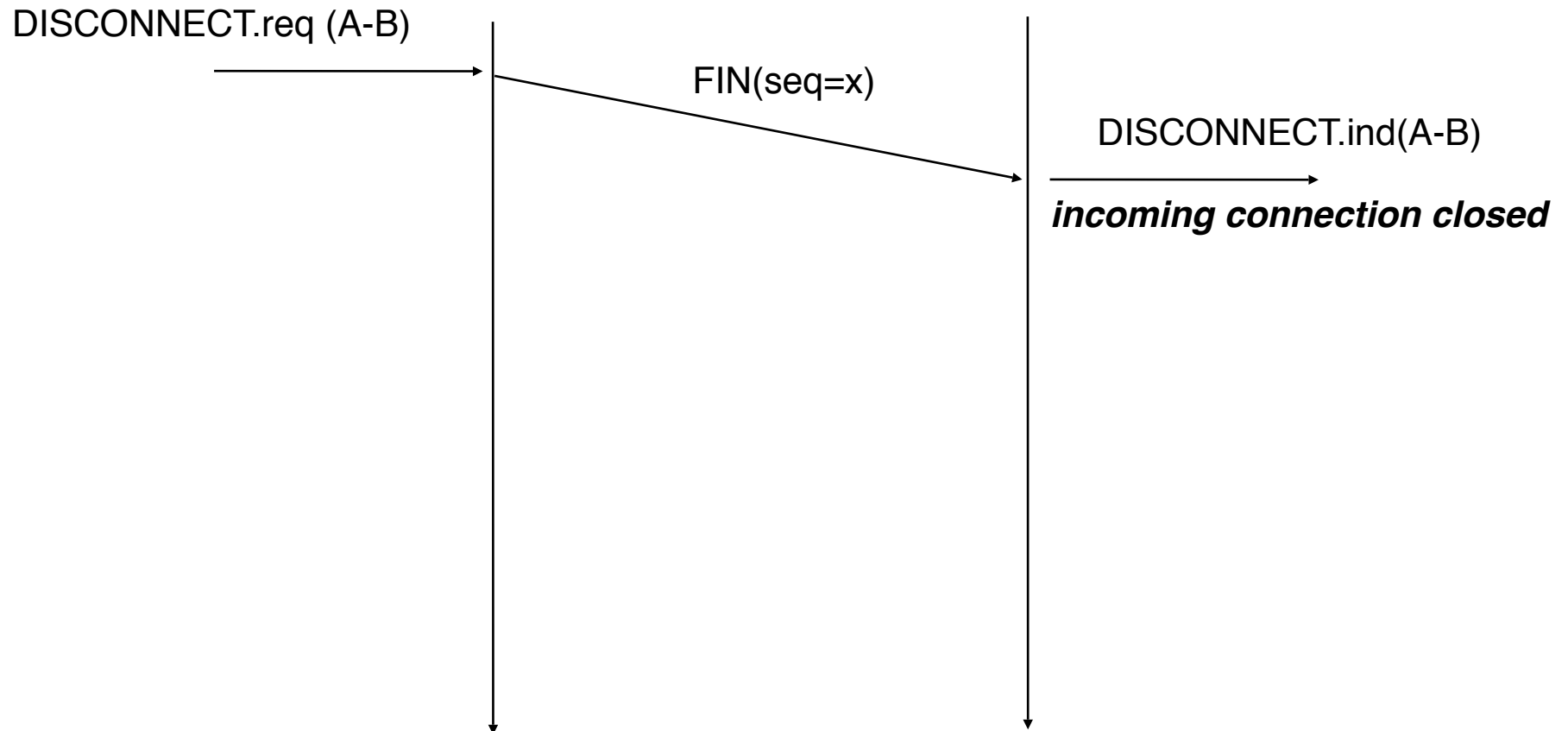    Reliable data transfer
    Congestion control

# TCP connection release

## Graceful shutdown of a TCP connection
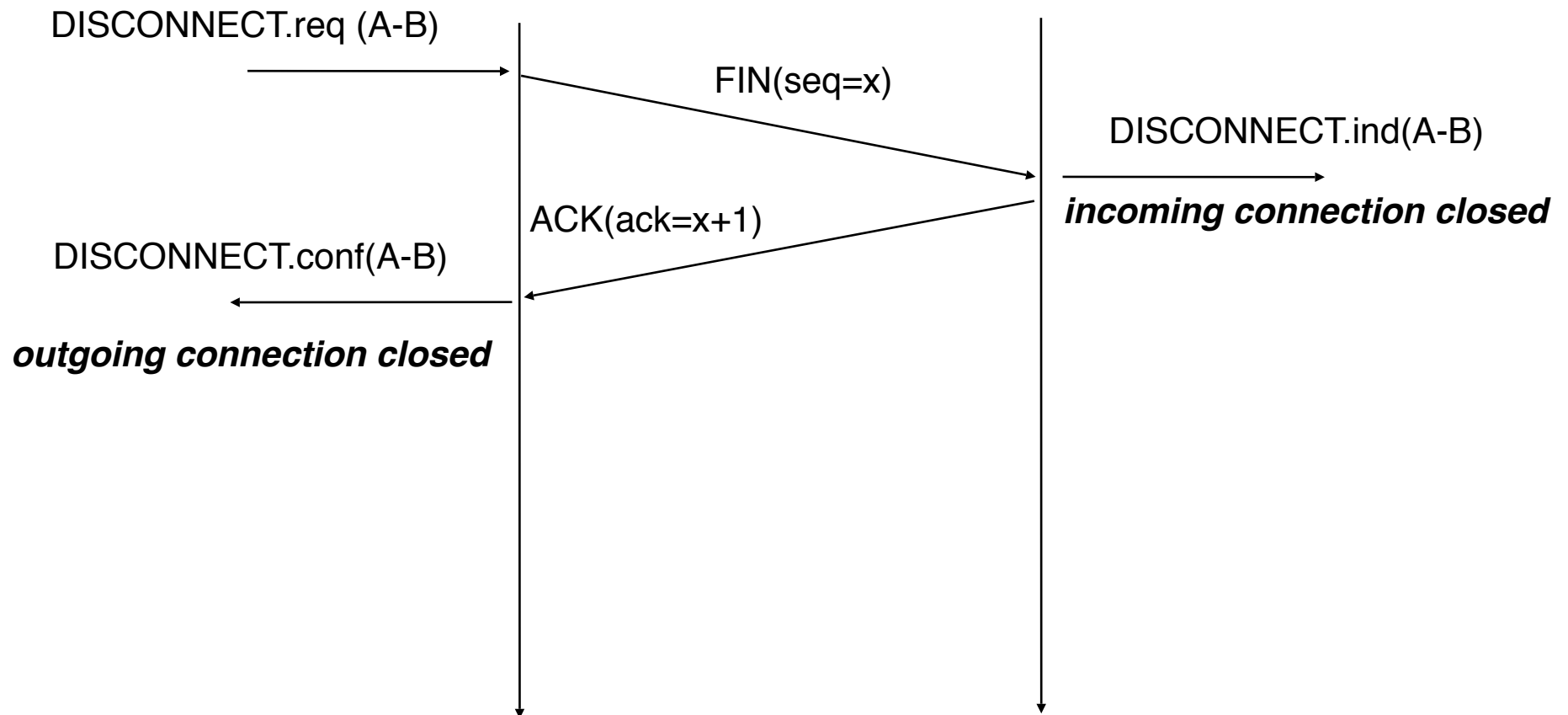
# TCP connection release
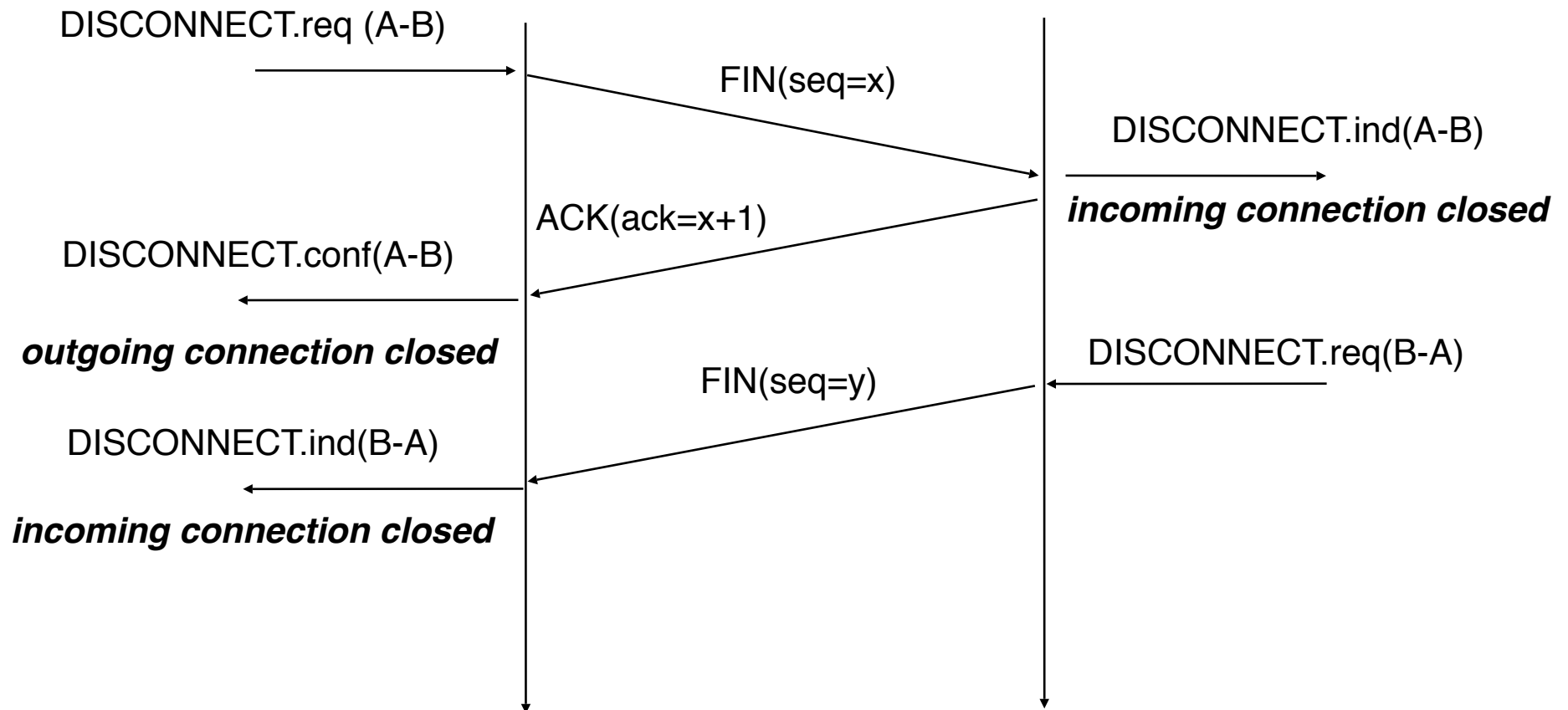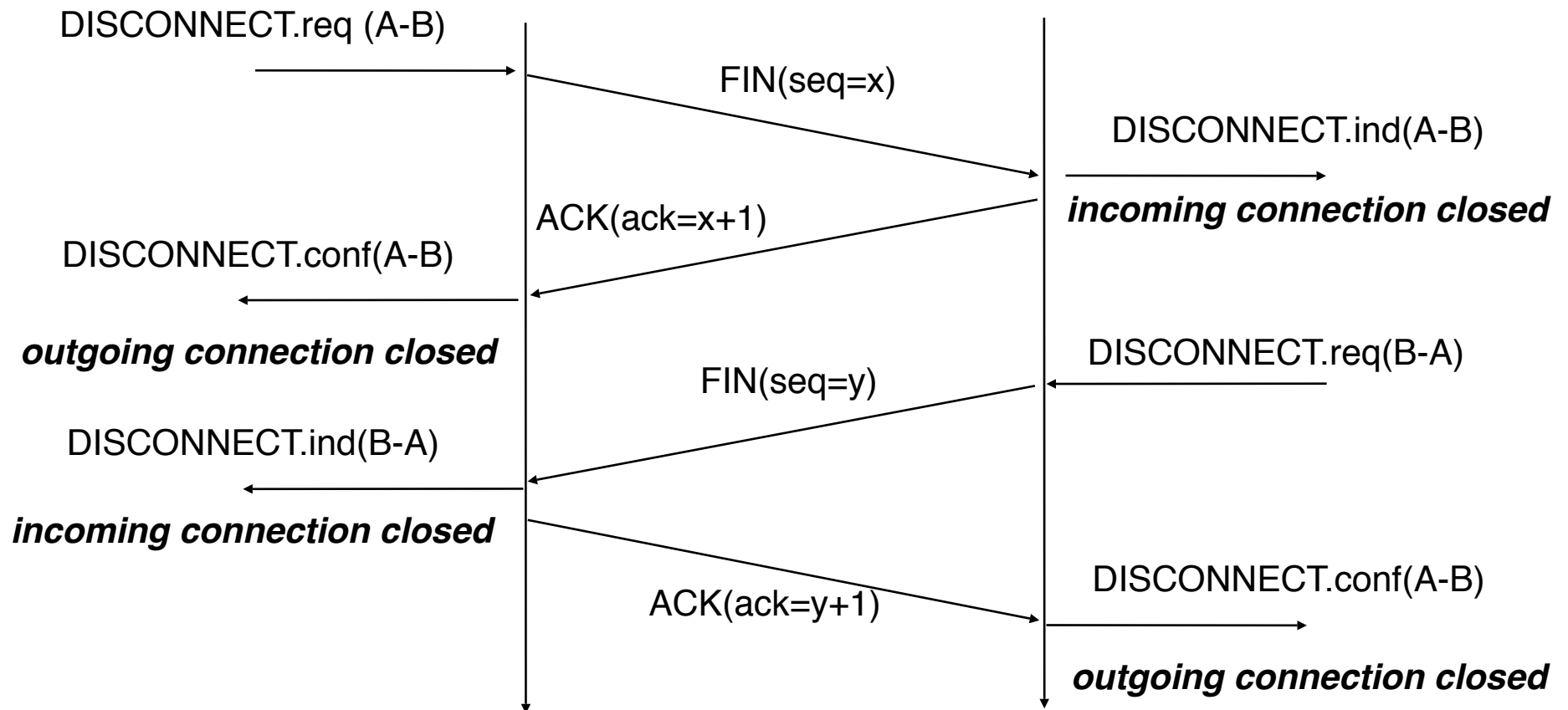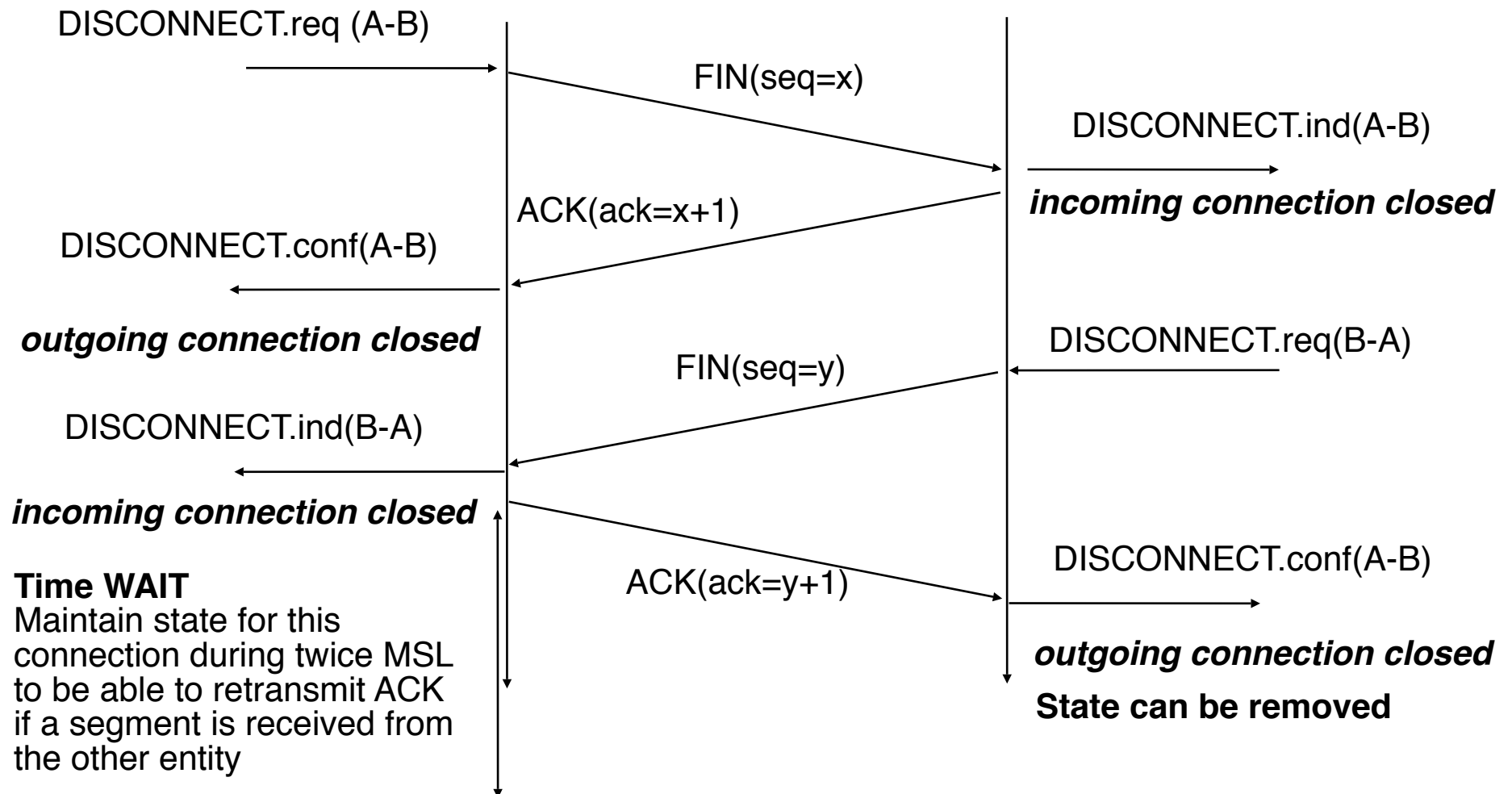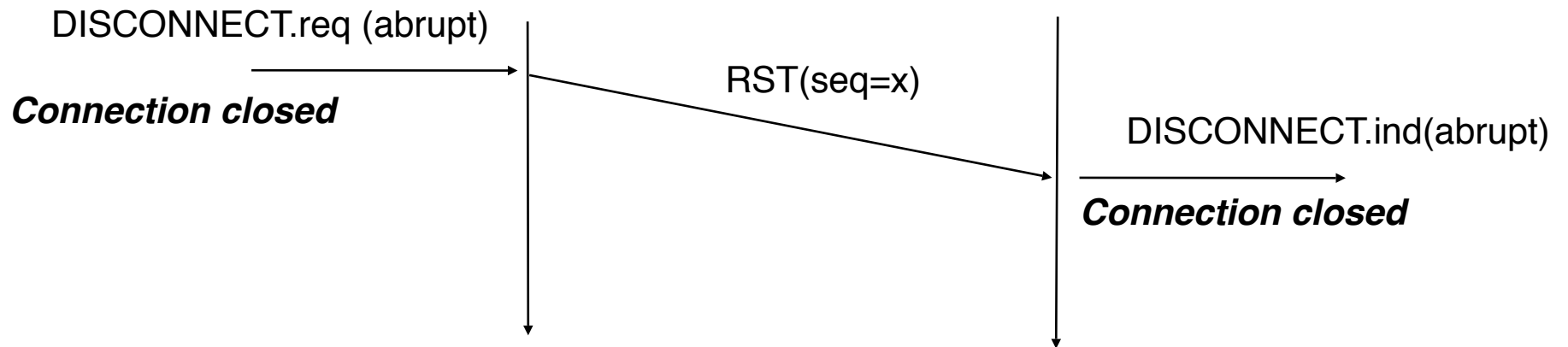
## Graceful shutdown of a TCP connection

DISCONNECT.req (A-B)

FIN(seq=x)

DISCONNECT.ind(A-B)

*incoming connection closed*

# TCP connection release

## Graceful shutdown of a TCP connection

DISCONNECT.req (A-B)

FIN(seq=x)

DISCONNECT.ind(A-B)

*incoming connection closed*

ACK(ack=x+1)

DISCONNECT.conf(A-B)

*outgoing connection closed*

# TCP connection release

## Graceful shutdown of a TCP connection

# TCP connection release

## Graceful shutdown of a TCP connection

# TCP connection release
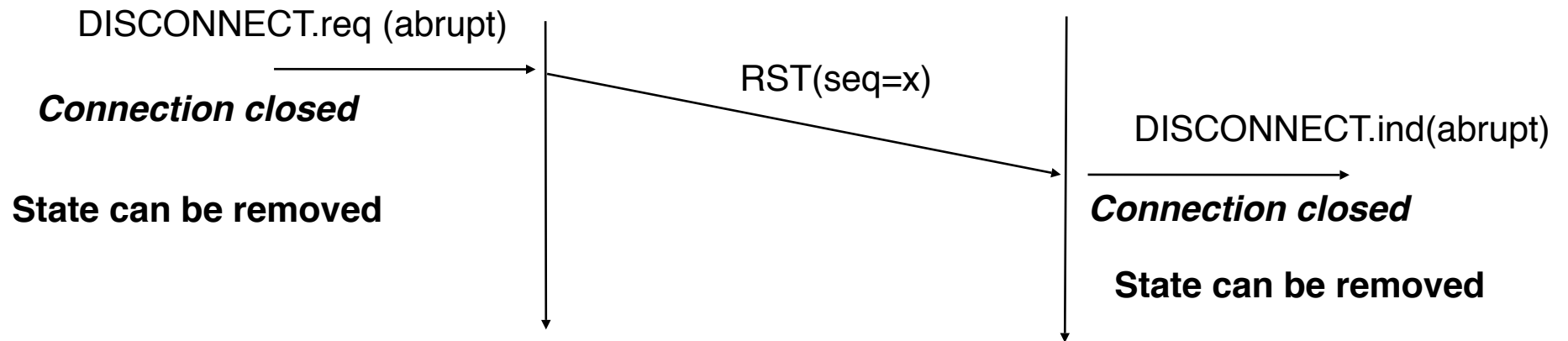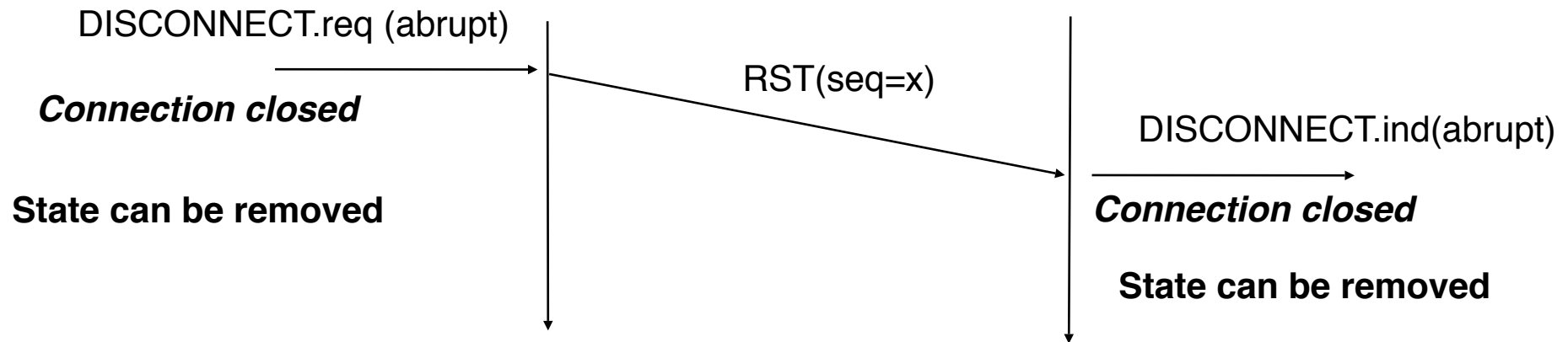
## Graceful shutdown of a TCP connection



DISCONNECT.req (A-B)

FIN(seq=x)

DISCONNECT.ind(A-B)

*incoming connection closed*

ACK(ack=x+1)

DISCONNECT.conf(A-B)

*outgoing connection closed*

DISCONNECT.req(B-A)

FIN(seq=y)

DISCONNECT.ind(B-A)

*incoming connection closed*

DISCONNECT.conf(A-B)

**Time WAIT**
Maintain state for this
connection during twice MSL
to be able to retransmit ACK
if a segment is received from
the other entity

ACK(ack=y+1)

*outgoing connection closed*

**State can be removed**

# Abrupt TCP connection release

# Abrupt TCP connection release

DISCONNECT.req (abrupt)

**Connection closed**

RST(seq=x)

DISCONNECT.ind(abrupt)

**Connection closed**

# Abrupt TCP connection release

# Abrupt TCP connection release

DISCONNECT.req (abrupt)

RST(seq=x)

***Connection closed***

DISCONNECT.ind(abrupt)

**State can be removed**

***Connection closed***

**State can be removed**

Data segments can be lost during such an abrupt release

No entity needs to wait in TIME_WAIT state after such a release

anyway, any segment received when there is no state causes the transmission of a RST segment
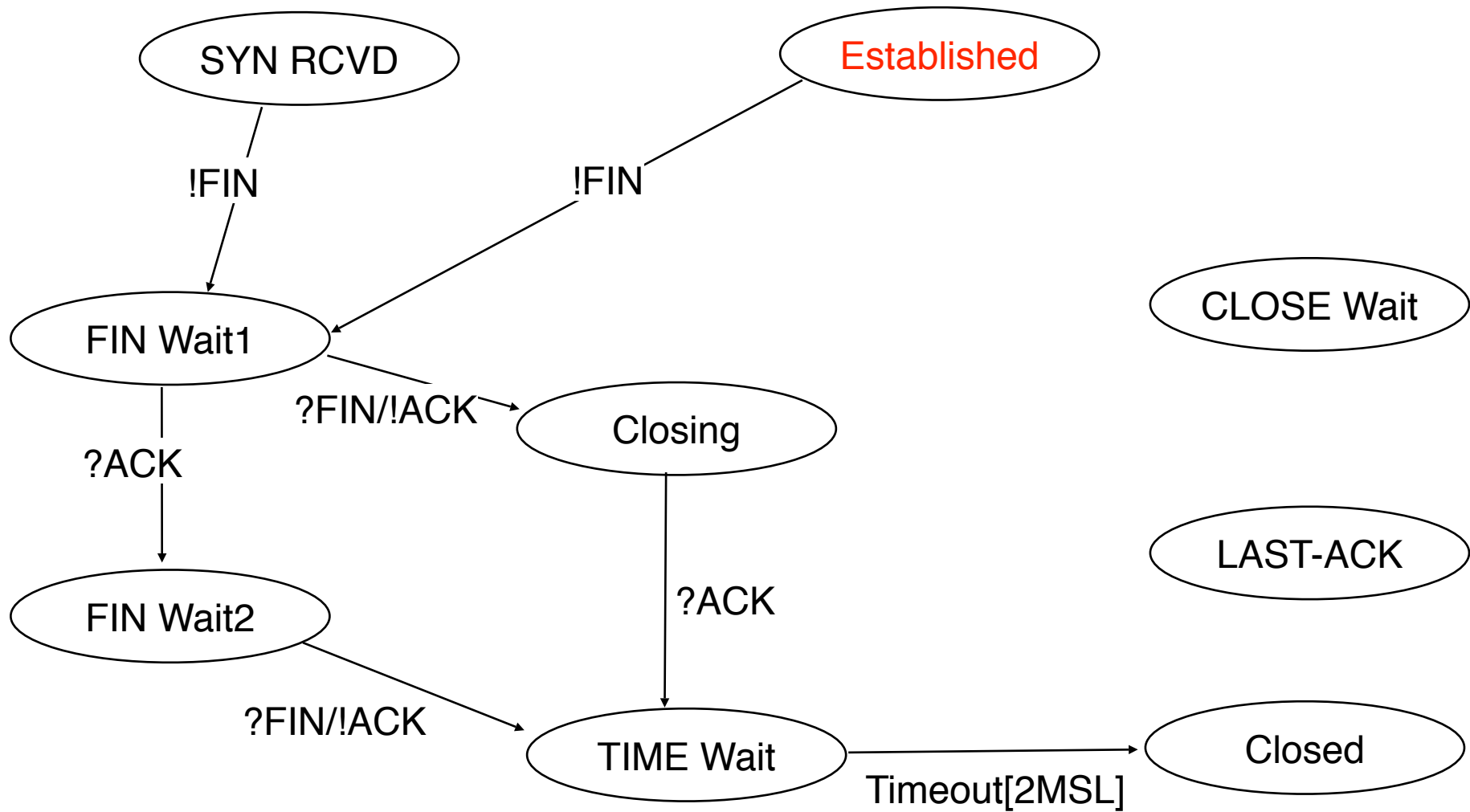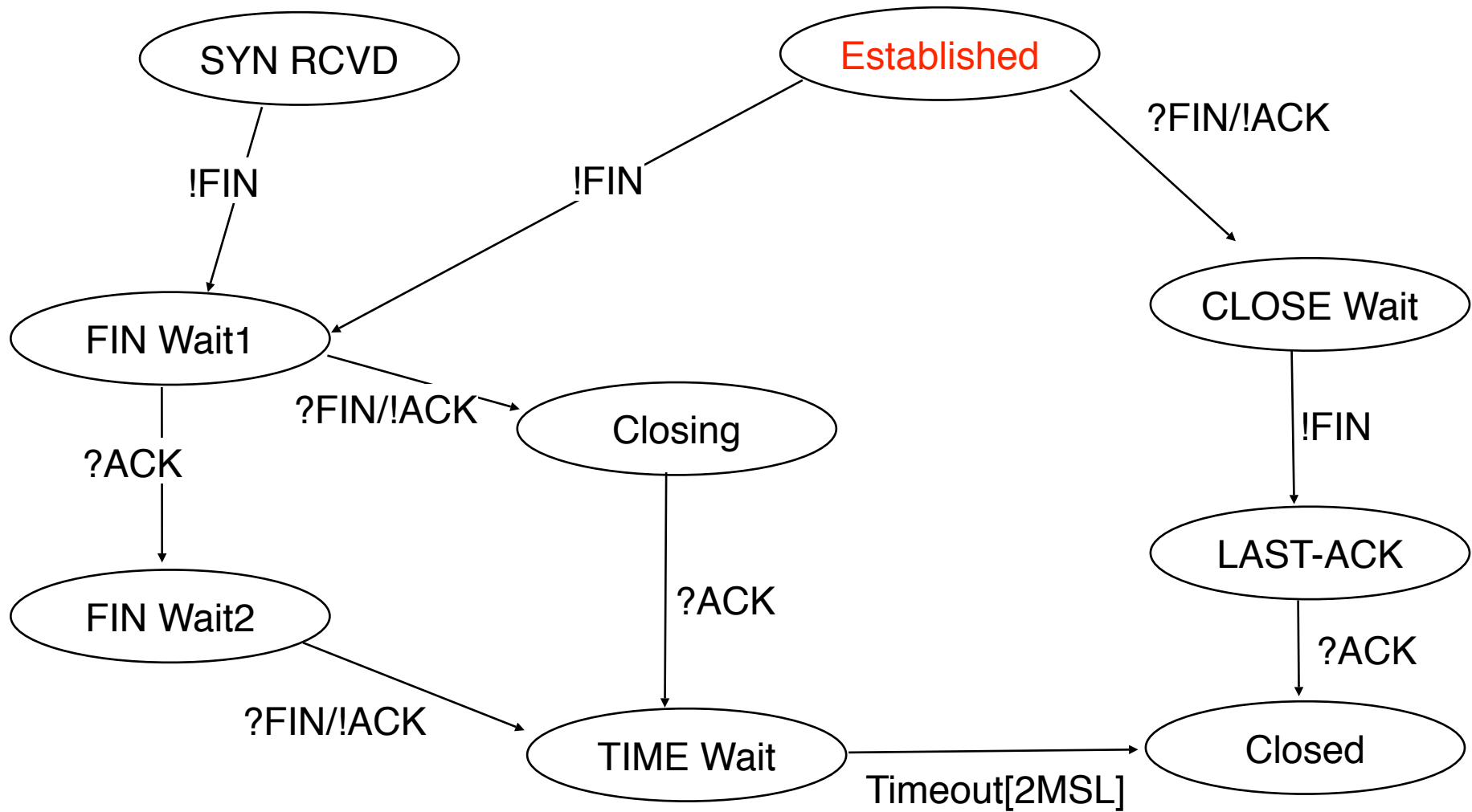
# TCP connection release

# TCP connection release

# TCP connection release

# TCP connection release
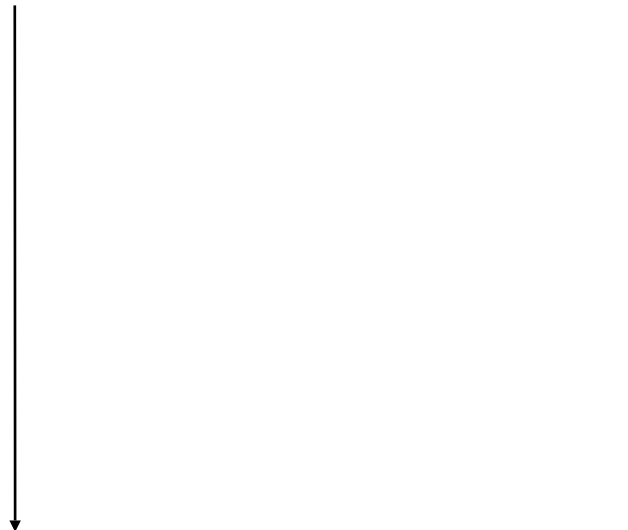
# TCP connection release

# Module 3 : Transport Layer

Basics

Building a reliable transport layer

UDP : a simple connectionless transport protocol

<span style="color:red">TCP : a reliable connection oriented transport protocol</span>

    TCP connection establishment

    TCP connection release

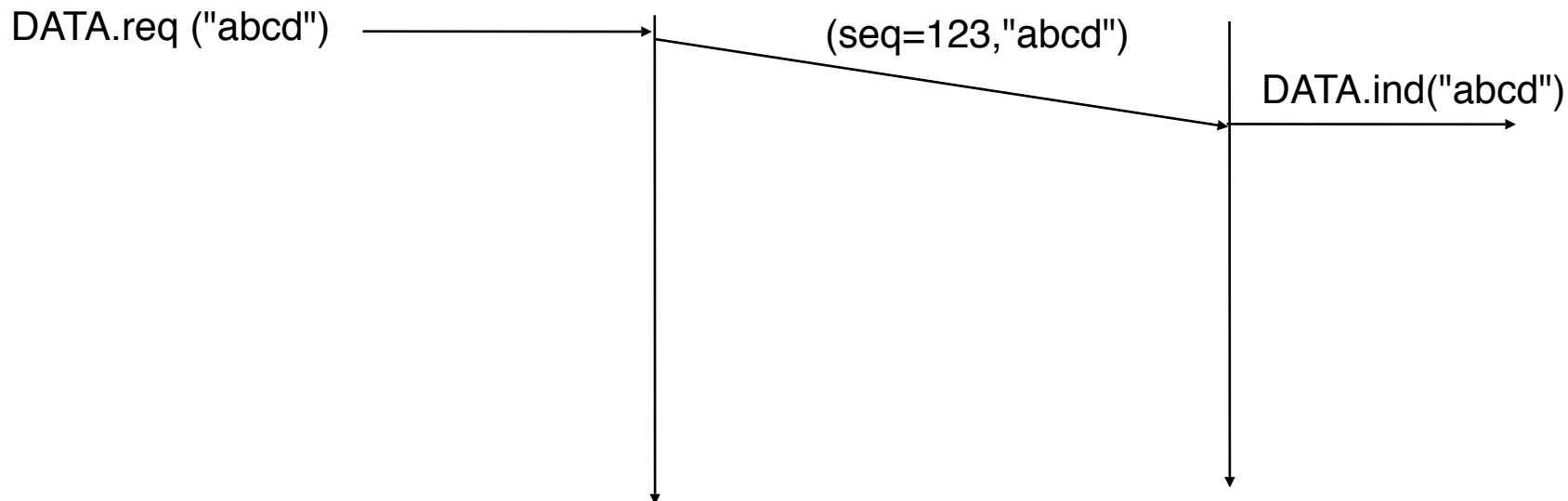⟶   <span style="color:red">Reliable data transfer</span>

    Congestion control

# Reliable data transfer

Each TCP segment contains

16 bits checksum
used to detect transmission errors affecting paylaod

32 bits sequence number (one byte=one seq. number)
used by sender to delimitate sent segments
used by receiver to reorder received segments

32 bits acknowledgement number
used (when ACK flag is 1) by receiver to advertise the sequence
number of the next expected byte (last byte received in sequence+1)

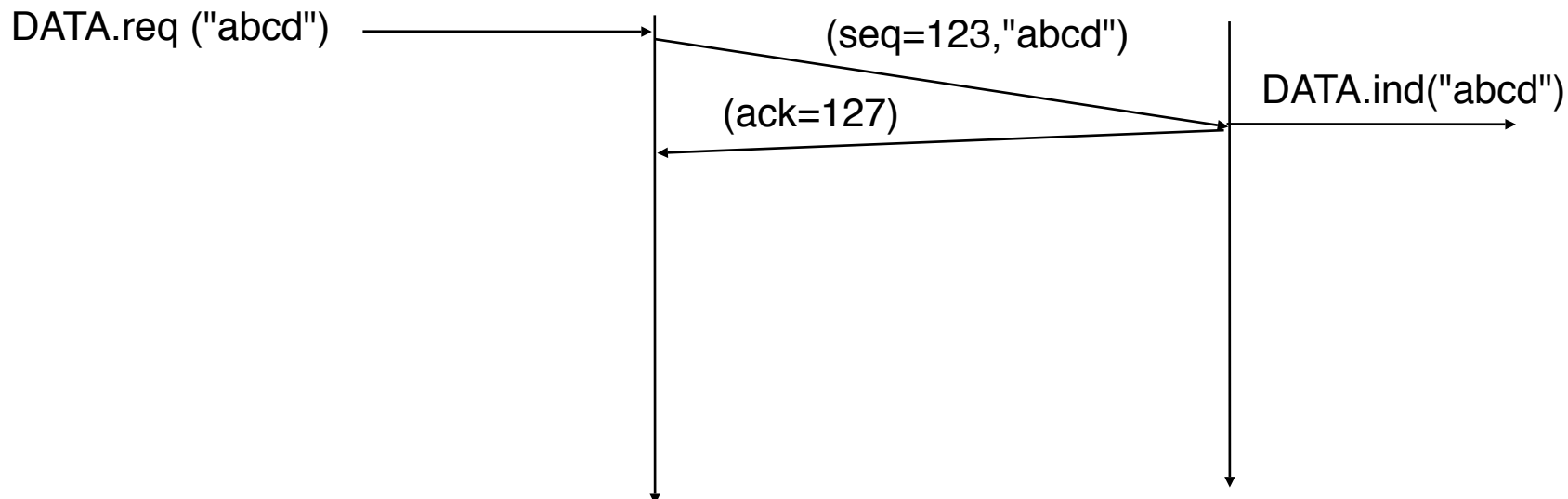# Reliable data transfer

Each TCP segment contains

16 bits checksum
used to detect transmission errors affecting paylaod

32 bits sequence number (one byte=one seq. number)
used by sender to delimitate sent segments
used by receiver to reorder received segments

32 bits acknowledgement number
used (when ACK flag is 1) by receiver to advertise the sequence
number of the next expected byte (last byte received in sequence+1)

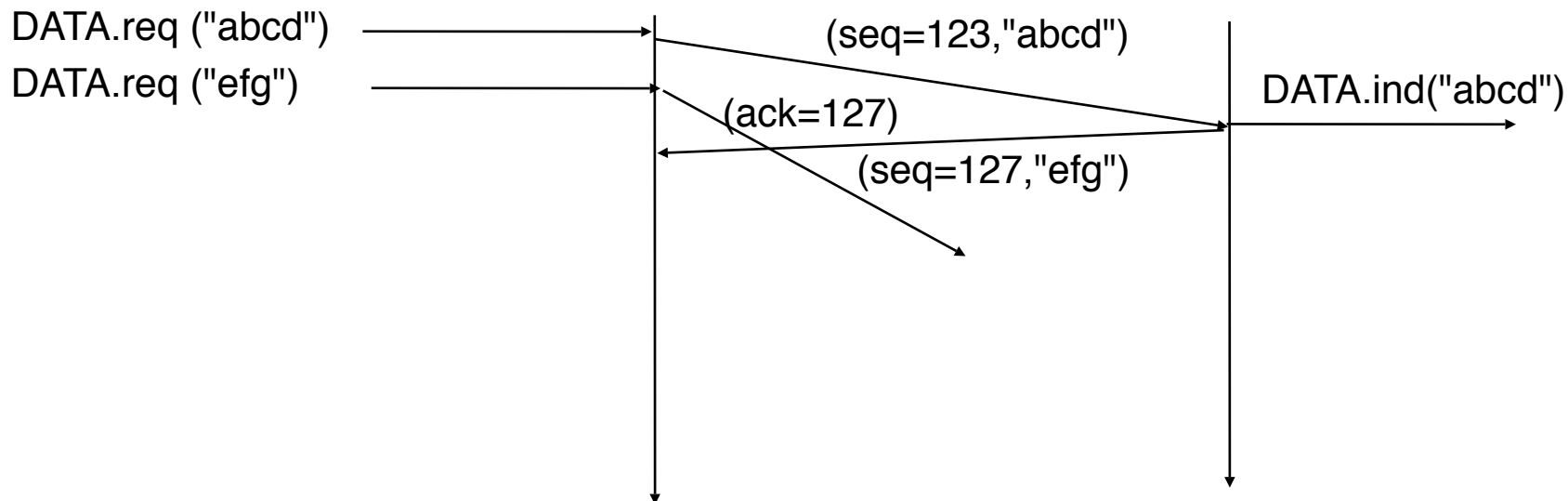DATA.req ("abcd")

(seq=123,"abcd")

DATA.ind("abcd")

# Reliable data transfer

Each TCP segment contains

16 bits checksum
used to detect transmission errors affecting paylaod

32 bits sequence number (one byte=one seq. number)
used by sender to delimitate sent segments
used by receiver to reorder received segments

32 bits acknowledgement number
used (when ACK flag is 1) by receiver to advertise the sequence
number of the next expected byte (last byte received in sequence+1)

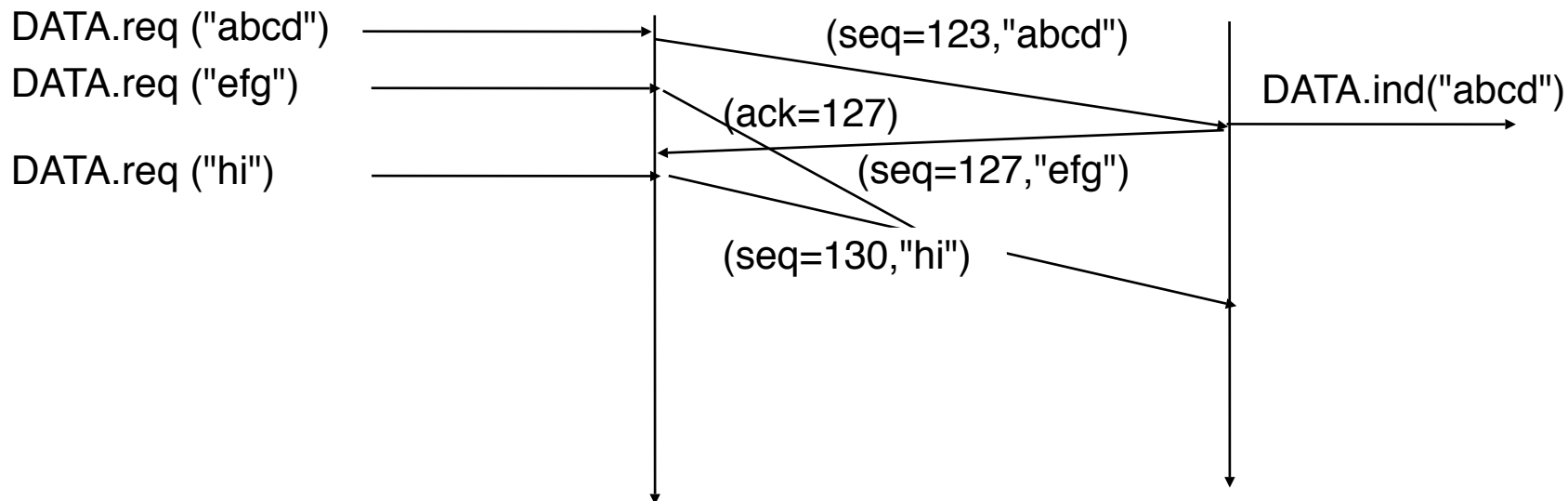DATA.req ("abcd")

(seq=123,"abcd")

DATA.ind("abcd")

(ack=127)

# Reliable data transfer

Each TCP segment contains
- 16 bits checksum
  - used to detect transmission errors affecting paylaod
- 32 bits sequence number (one byte=one seq. number)
  - used by sender to delimitate sent segments
  - used by receiver to reorder received segments
- 32 bits acknowledgement number
  - used (when ACK flag is 1) by receiver to advertise the sequence
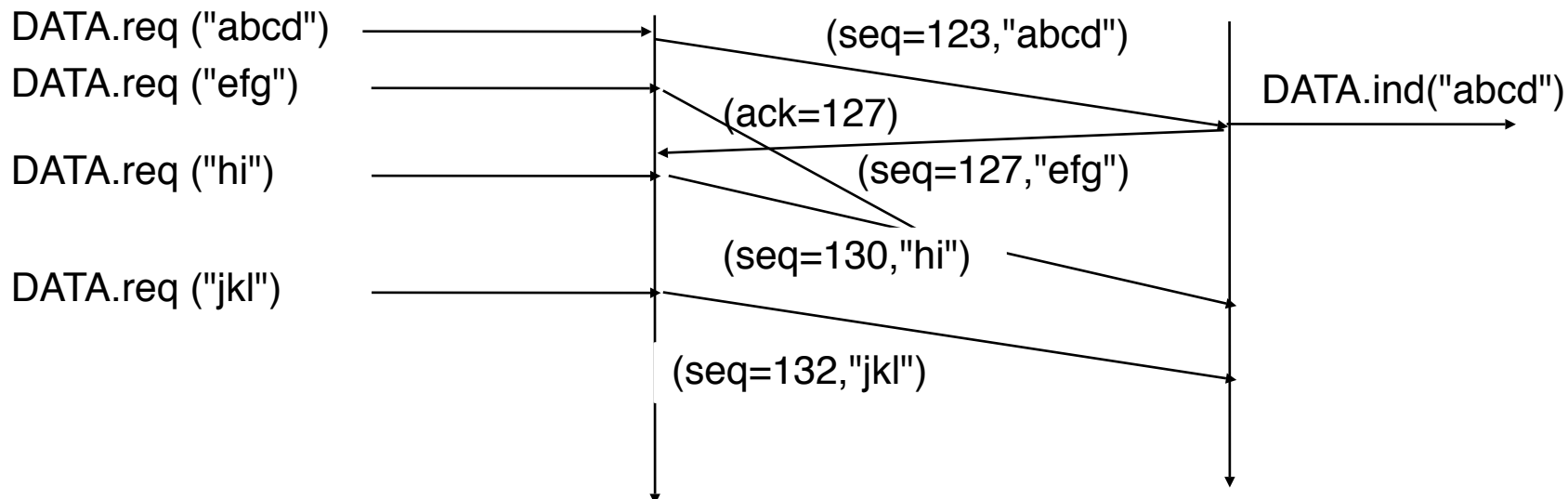  - number of the next expected byte (last byte received in sequence+1)

DATA.req ("abcd")             (seq=123,"abcd")

DATA.req ("efg")                                DATA.ind("abcd")

(ack=127)

(seq=127,"efg")

# Reliable data transfer

Each TCP segment contains

16 bits checksum
used to detect transmission errors affecting paylaod

32 bits sequence number (one byte=one seq. number)
used by sender to delimitate sent segments
used by receiver to reorder received segments

32 bits acknowledgement number
used (when ACK flag is 1) by receiver to advertise the sequence
number of the next expected byte (last byte received in sequence+1)

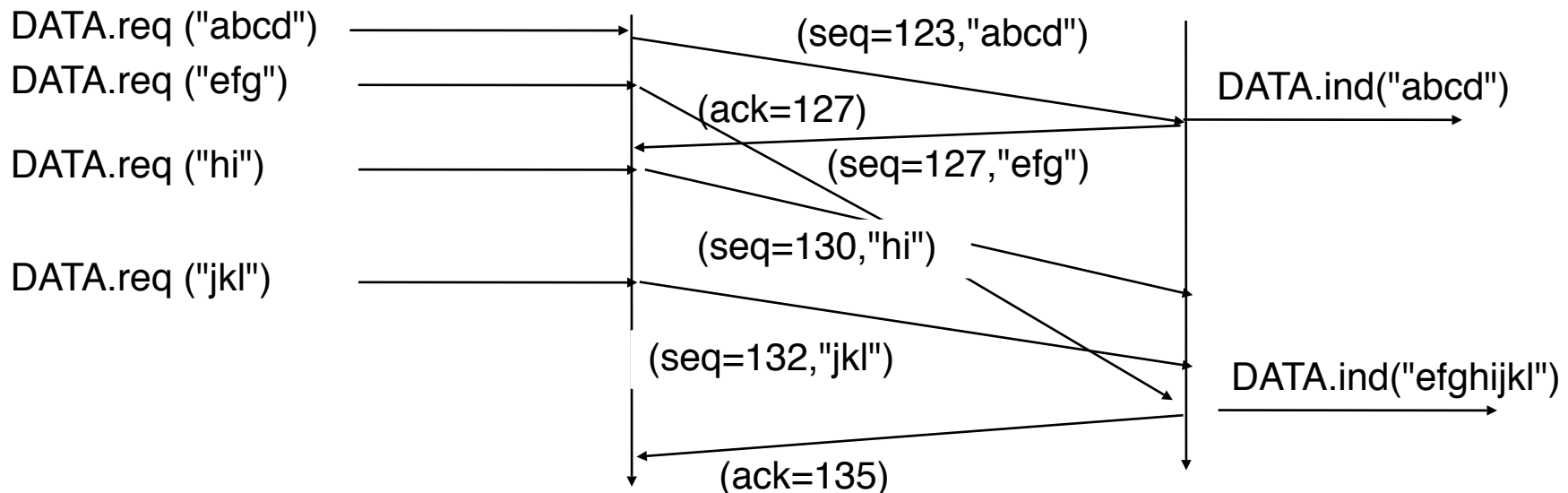DATA.req ("abcd")          (seq=123,"abcd")

DATA.req ("efg")                              DATA.ind("abcd")
                    (ack=127)

DATA.req ("hi")           (seq=127,"efg")

              (seq=130,"hi")

# Reliable data transfer

Each TCP segment contains
- 16 bits checksum
    used to detect transmission errors affecting paylaod
- 32 bits sequence number (one byte=one seq. number)
    used by sender to delimitate sent segments
    used by receiver to reorder received segments
- 32 bits acknowledgement number
    used (when ACK flag is 1) by receiver to advertise the sequence
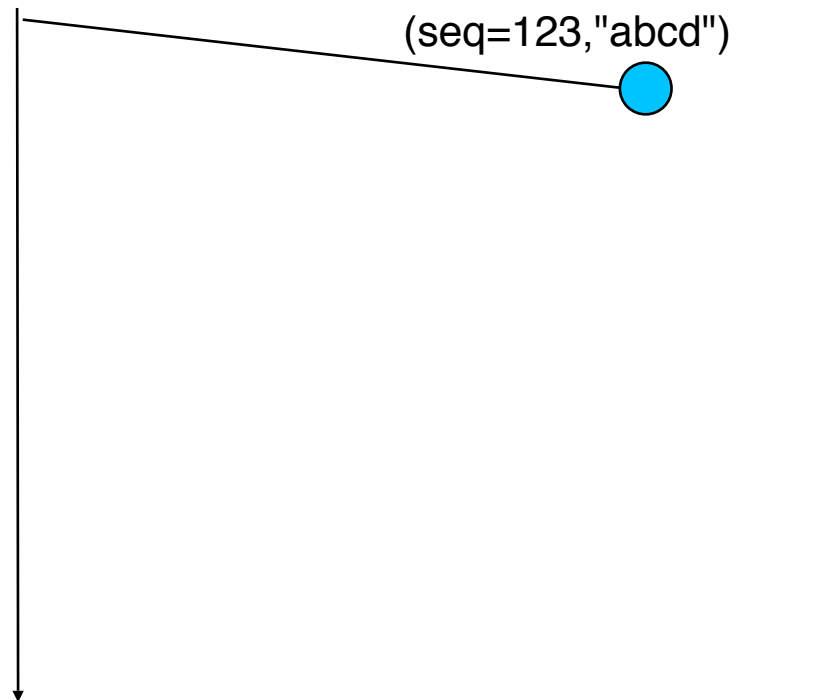    number of the next expected byte (last byte received in sequence+1)

# Reliable data transfer

Each TCP segment contains

16 bits checksum
used to detect transmission errors affecting paylaod

32 bits sequence number (one byte=one seq. number)
used by sender to delimitate sent segments
used by receiver to reorder received segments

32 bits acknowledgement number
used (when ACK flag is 1) by receiver to advertise the sequence
number of the next expected byte (last byte received in sequence+1)

DATA.req ("abcd") → (seq=123,"abcd")

DATA.req ("efg") →

DATA.ind("abcd")

(ack=127)

DATA.req ("hi") → (seq=127,"efg")

(seq=130,"hi")

DATA.req ("jkl") →

(seq=132,"jkl")

DATA.ind("efghijkl")

(ack=135)

# Reliable data transfer

## How to deal with segment losses ?

### TCP uses a retransmission timer

If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments

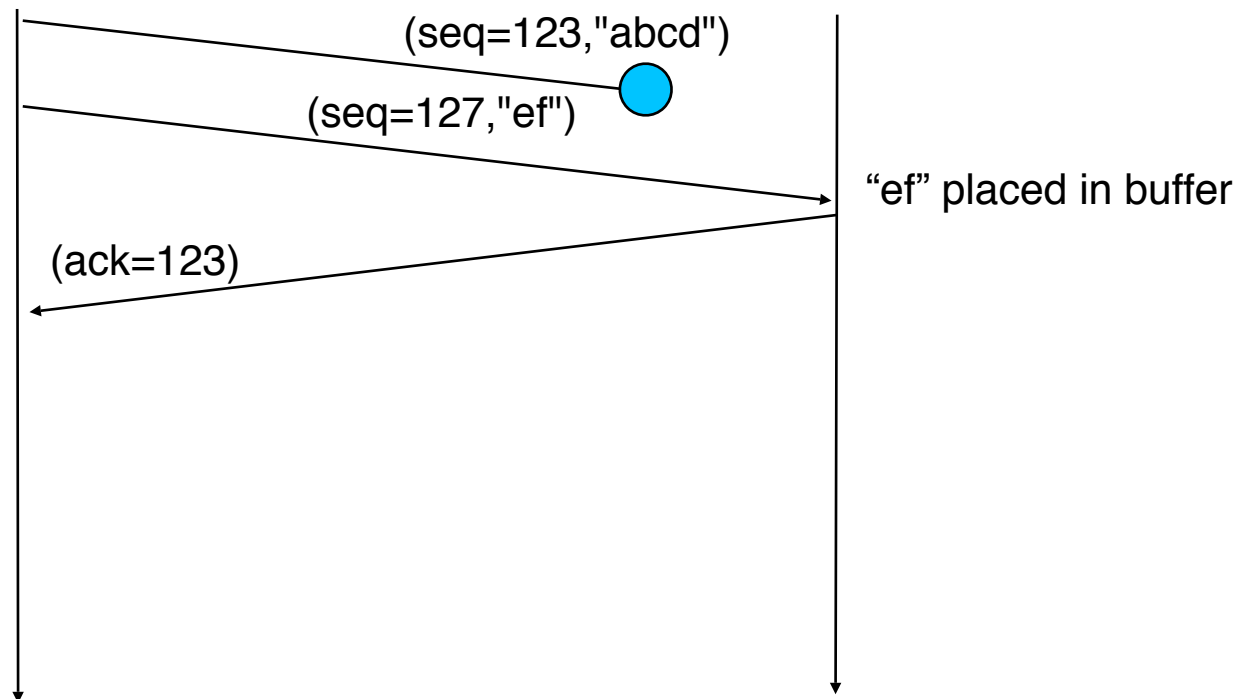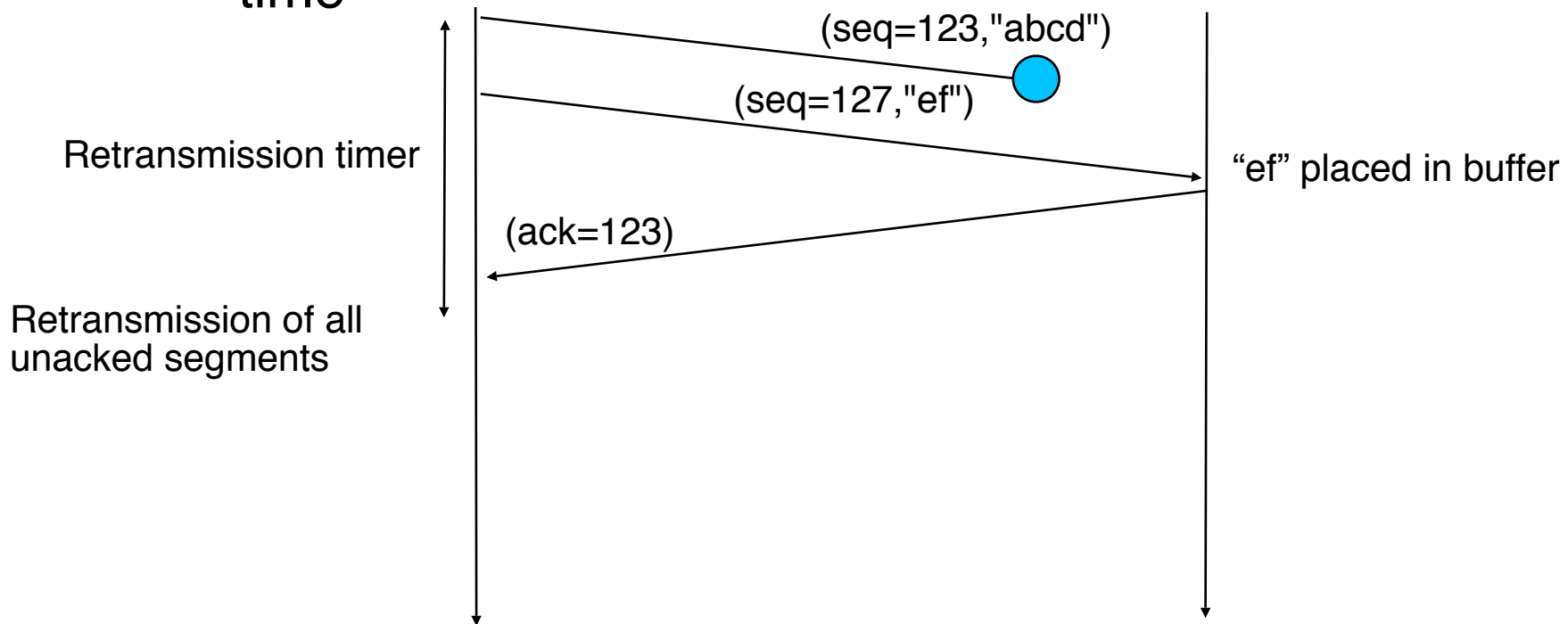usually a single retransmission timer is running at a given time

# Reliable data transfer

## How to deal with segment losses ?

### TCP uses a retransmission timer

If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments
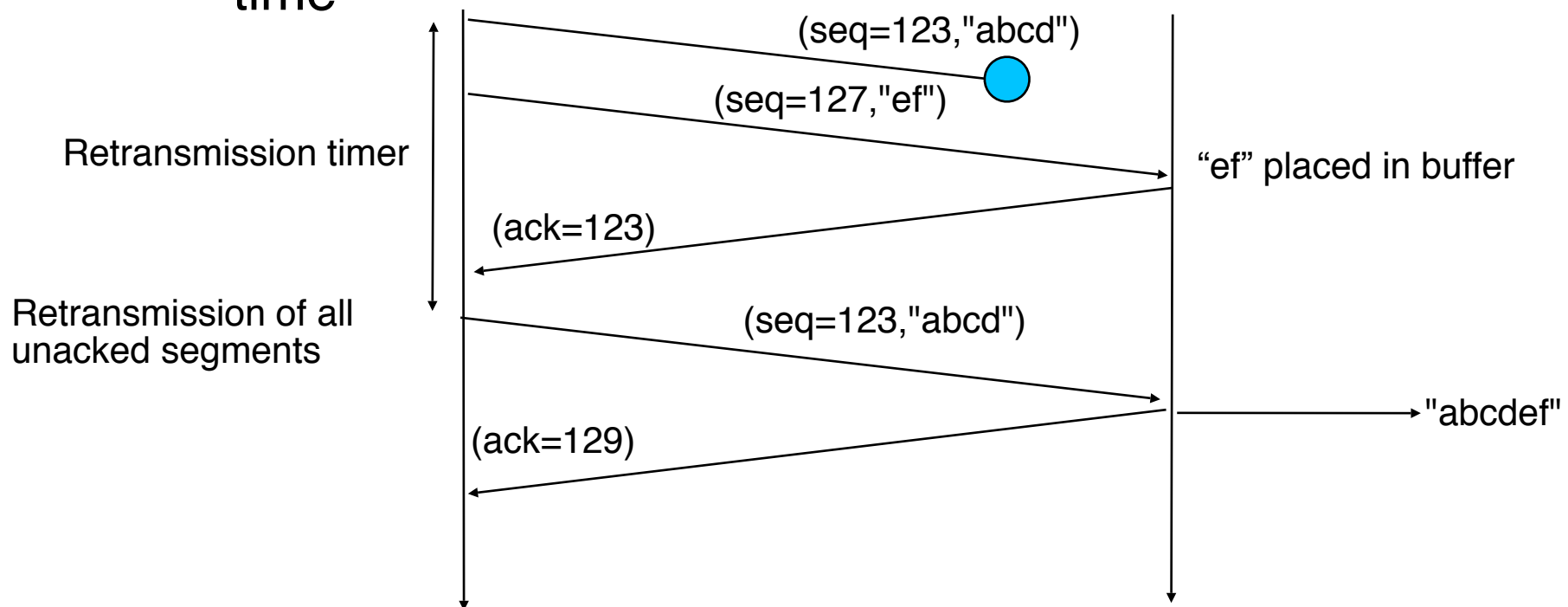usually a single retransmission timer is running at a given time

(seq=123,"abcd")

# Reliable data transfer

## How to deal with segment losses ?
### TCP uses a retransmission timer

If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments
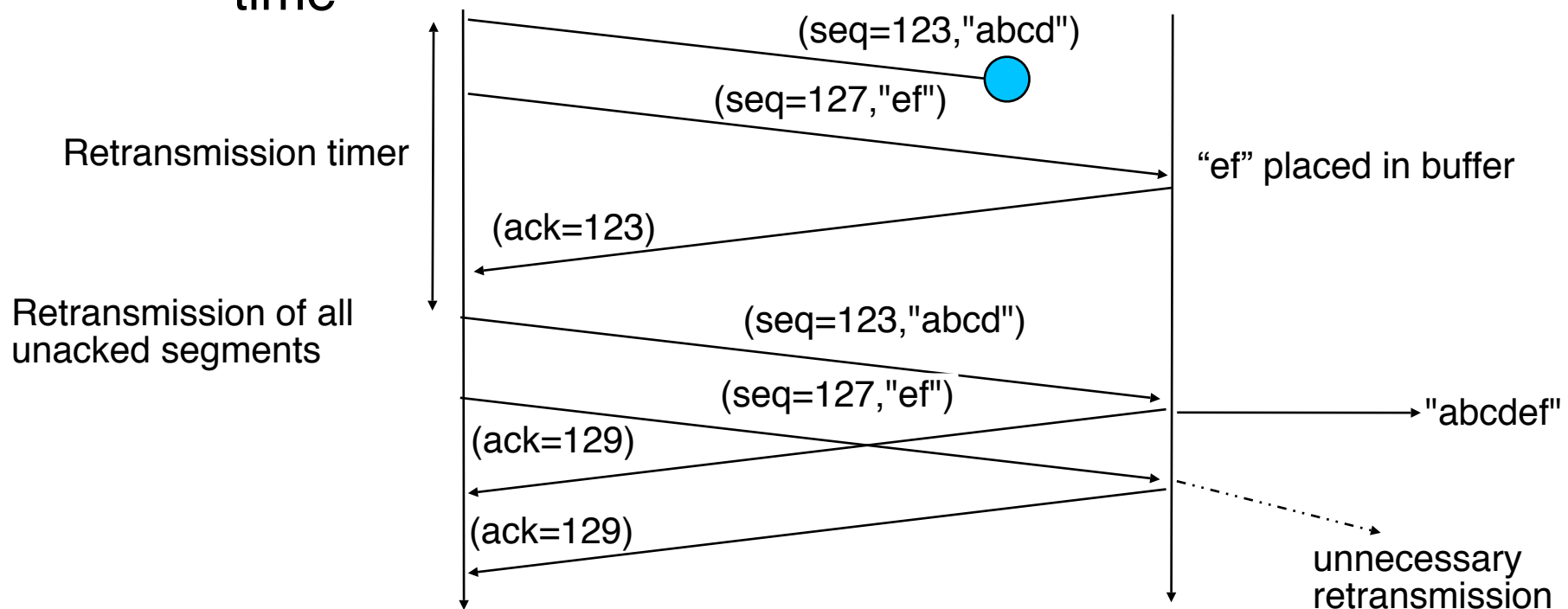usually a single retransmission timer is running at a given time



(seq=123,"abcd")

(seq=127,"ef")

"ef" placed in buffer

(ack=123)

# Reliable data transfer

## How to deal with segment losses ?
### TCP uses a retransmission timer
If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments
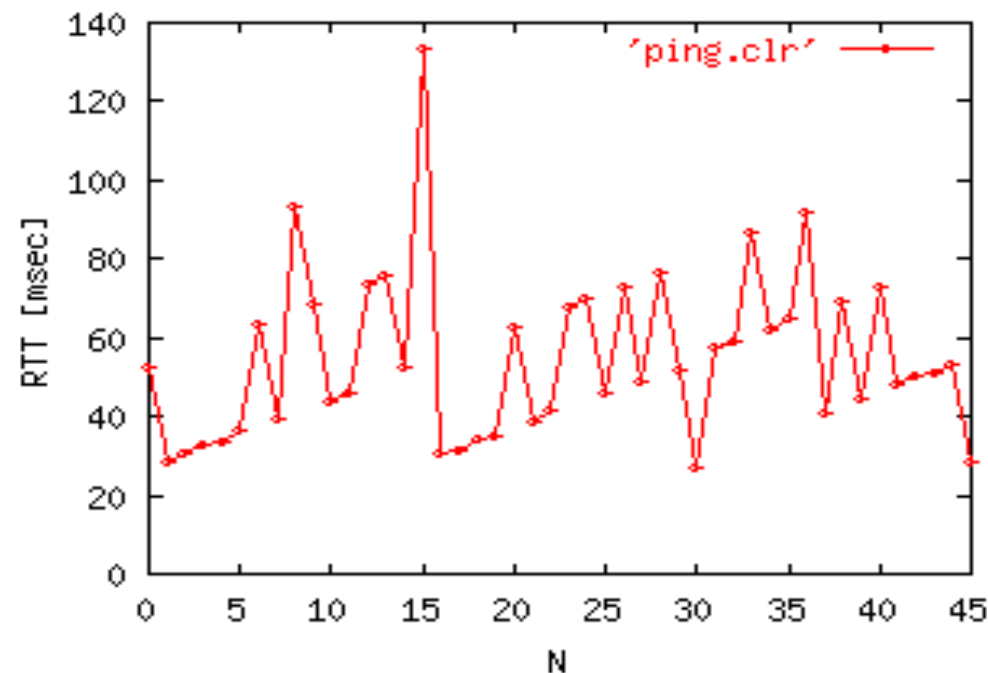usually a single retransmission timer is running at a given time

(seq=123,"abcd")

(seq=127,"ef")

Retransmission timer

"ef" placed in buffer

(ack=123)

Retransmission of all
unacked segments

# Reliable data transfer

## How to deal with segment losses ?
### TCP uses a retransmission timer
If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments
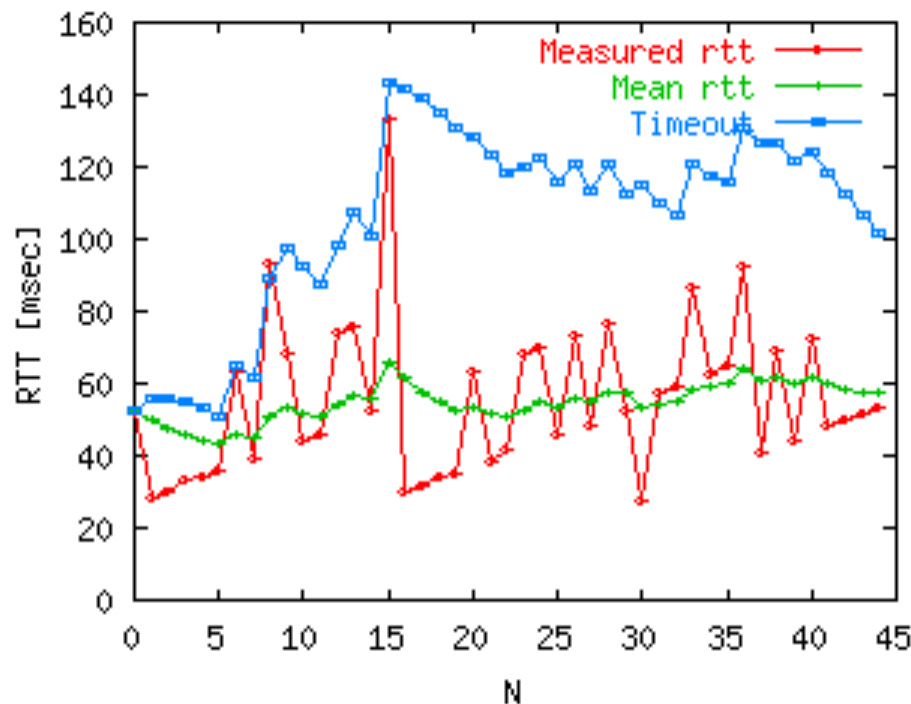usually a single retransmission timer is running at a given time



Retransmission timer

(seq=123,"abcd")

(seq=127,"ef")

"ef" placed in buffer

(ack=123)

Retransmission of all unacked segments

(seq=123,"abcd")

"abcdef"

(ack=129)

# Reliable data transfer

## How to deal with segment losses ?
### TCP uses a retransmission timer
If the retransmission timer expires, TCP performs go-back-n and retransmits all the unacknowledged segments
usually a single retransmission timer is running at a given time



Retransmission timer

(seq=123,"abcd")

(seq=127,"ef")

"ef" placed in buffer

(ack=123)

Retransmission of all unacked segments

(seq=123,"abcd")

(seq=127,"ef")

"abcdef"

(ack=129)

(ack=129)

unnecessary retransmission

# Retransmission timer

## How to compute it ?
### Issue
round-trip-time may change frequently during the lifetime of a TCP connection

# Retransmission timer

TCP's retransmission timer
  One timer per connection
    timer = mean(rtt) + 4*std_dev(rtt)
    Estimation of the mean
      $est\_mean(rtt) = (1-\alpha)*est\_mean(rtt) + \alpha*rtt\_measured$
    Estimation of the standard deviation of the rtt
      $est\_std\_dev = (1-\beta)*est\_std\_dev + \beta*|rtt\_measured - est\_mean(rtt)|$

# Round-trip-time estimation

## Problem

### How to measure rtt after retransmissions ?



measured rtt

(seq=120,"xyz")

(ack=123)

which is the good one ?   Timer

(seq=123,"abcd")

(seq=123,"abcd")

(ack=128)

## Solution (Karn/Partridge)

1. Do not measure rtt of retransmitted segments

# Round-trip-time estimation (2)

## Improvement to Karn/Partridge
### Add a timestamp in each segment sent
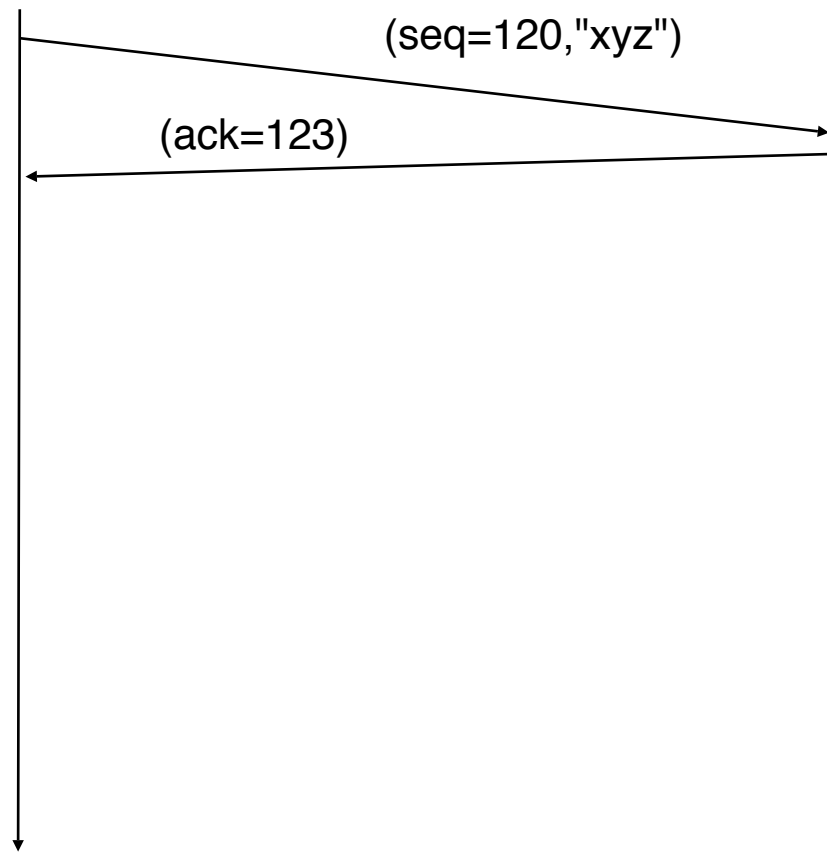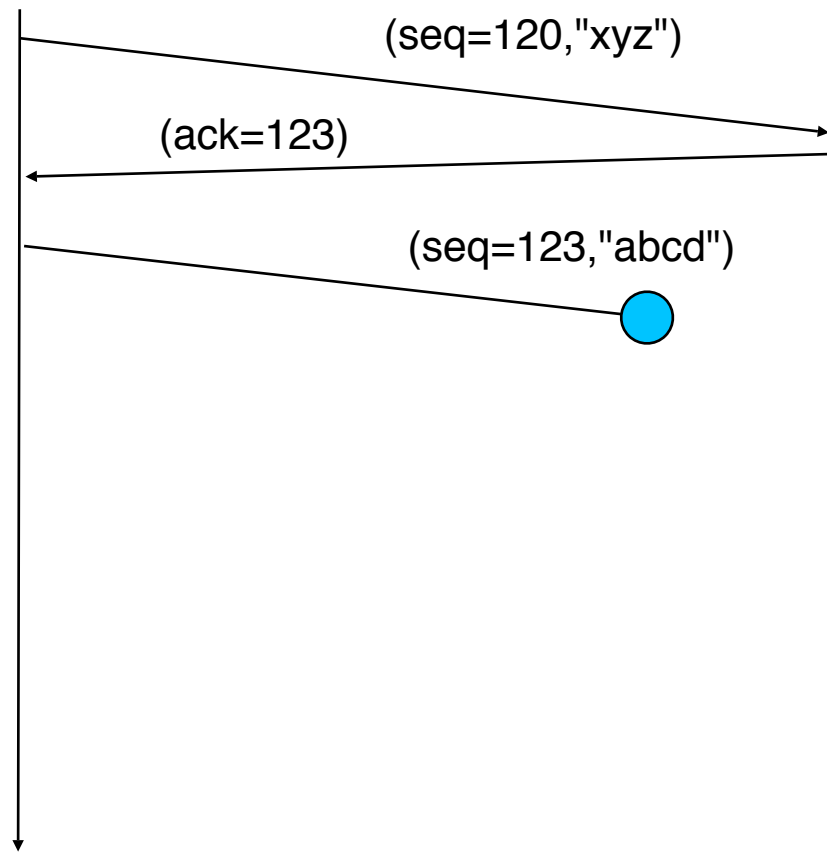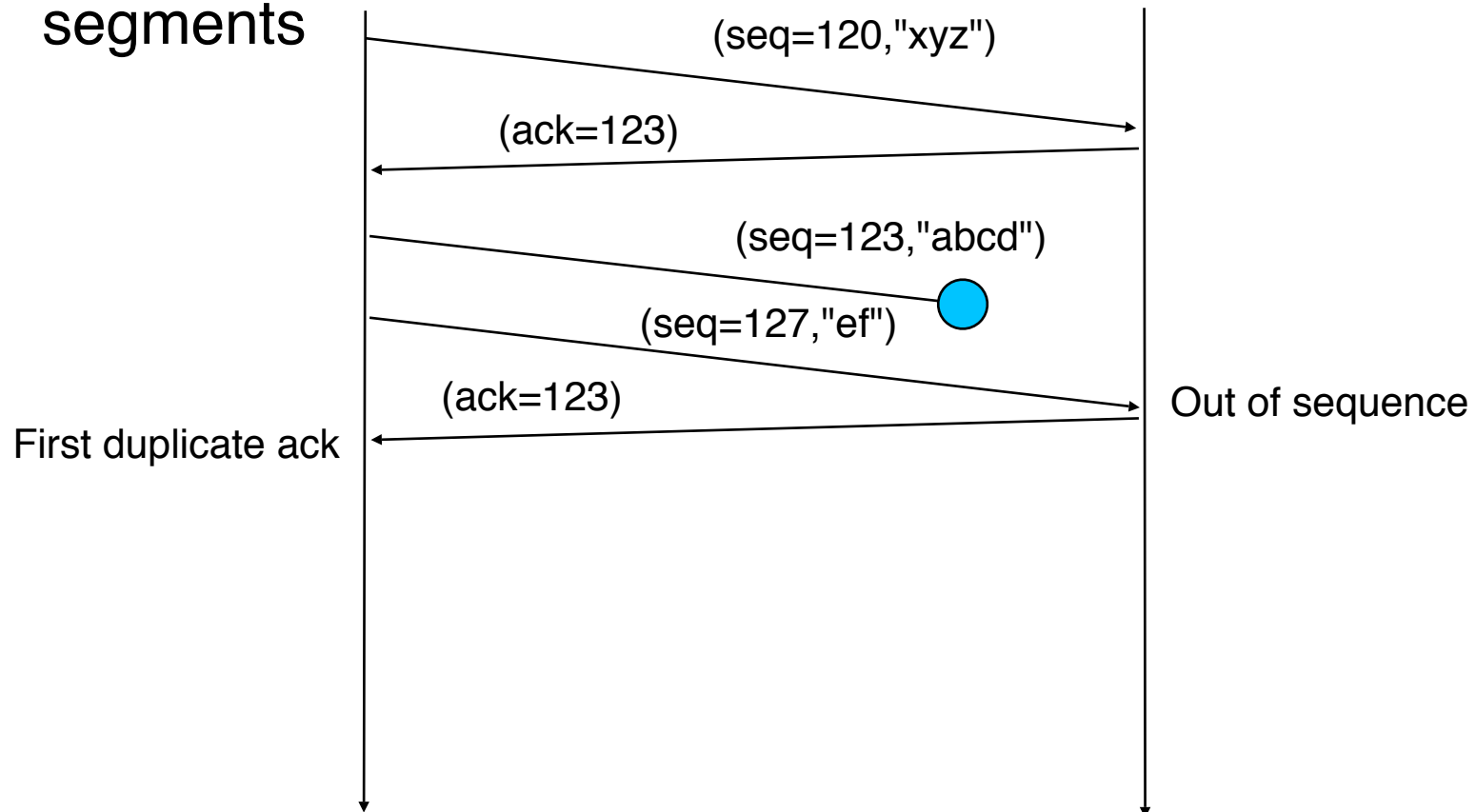#### TS and TSEcho (RFC1323)



© O. Bonaventure, 2008

CNP3/2008.3.

# Improving the reliable data transfer

How to improve the reaction to segment losses ?

TCP receiver should send an ack everytime an out-of-sequence segment is received

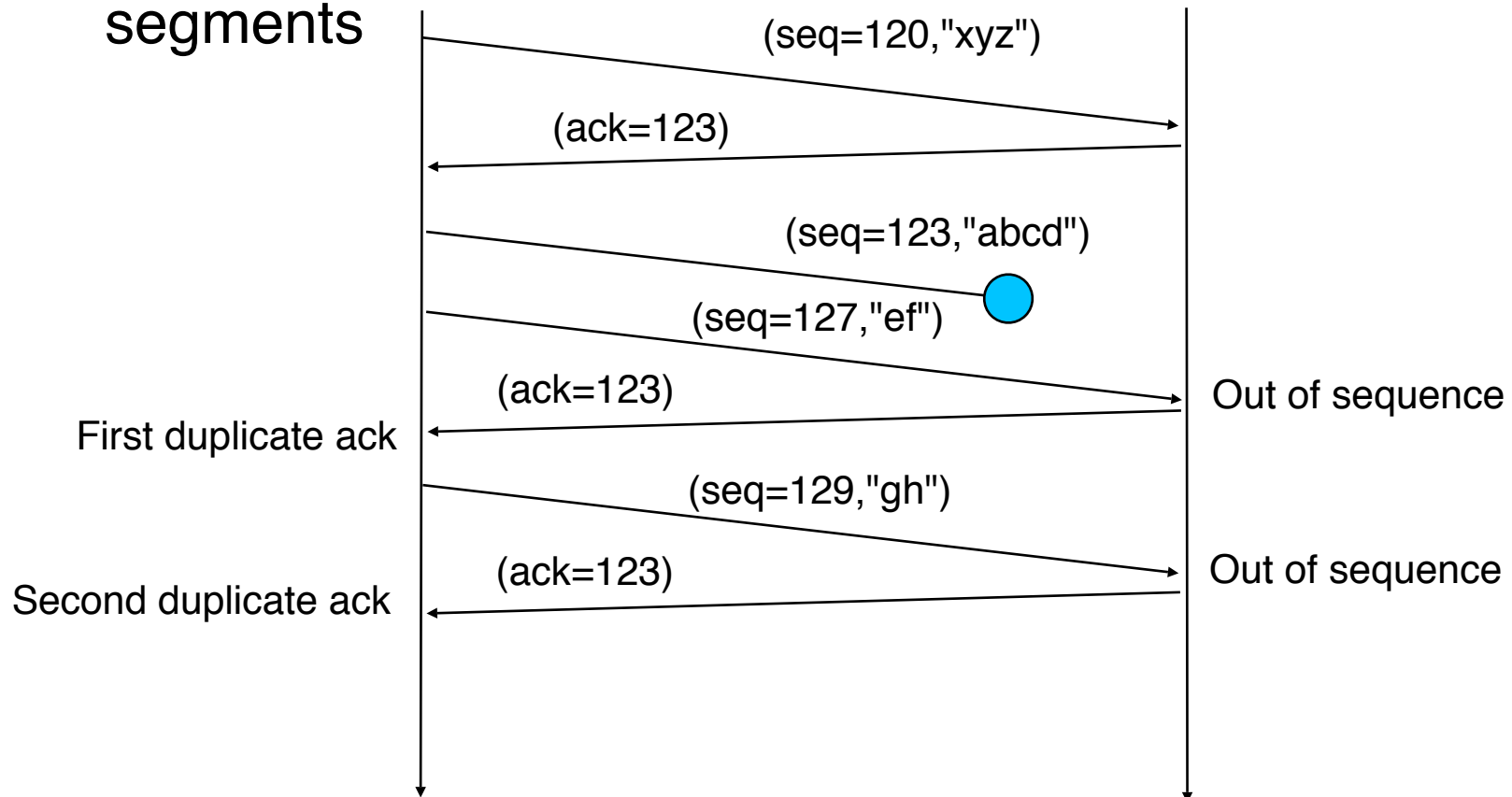Heuristic : a segment is considered lost after three duplicate segments

# Improving the reliable data transfer

How to improve the reaction to segment losses ?
TCP receiver should send an ack everytime an out-of-sequence segment is received

Heuristic : a segment is considered lost after three duplicate segments

(seq=120,"xyz")

(ack=123)

# Improving the reliable data transfer

How to improve the reaction to segment losses ?
TCP receiver should send an ack everytime an out-of-sequence segment is received

Heuristic : a segment is considered lost after three duplicate segments
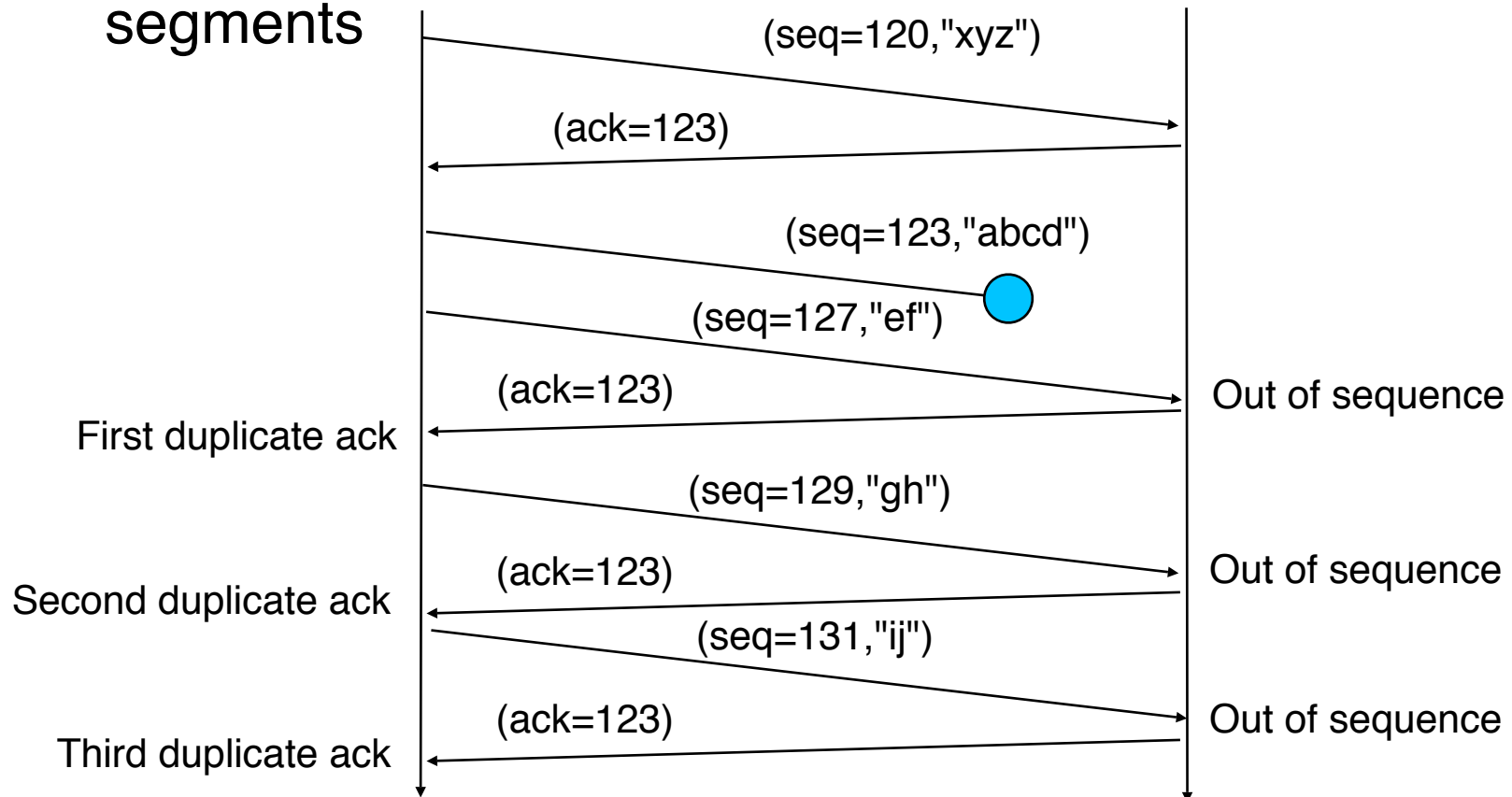


(seq=120,"xyz")

(ack=123)

(seq=123,"abcd")

# Improving the reliable data transfer

## How to improve the reaction to segment losses ?

TCP receiver should send an ack everytime an out-of-sequence segment is received

Heuristic : a segment is considered lost after three duplicate segments

# Improving the reliable data transfer

How to improve the reaction to segment losses ?
TCP receiver should send an ack everytime an out-of-sequence segment is received

Heuristic : a segment is considered lost after three duplicate segments

# Improving the reliable data transfer

How to improve the reaction to segment losses ?
TCP receiver should send an ack everytime an out-of-sequence segment is received

Heuristic : a segment is considered lost after three duplicate segments
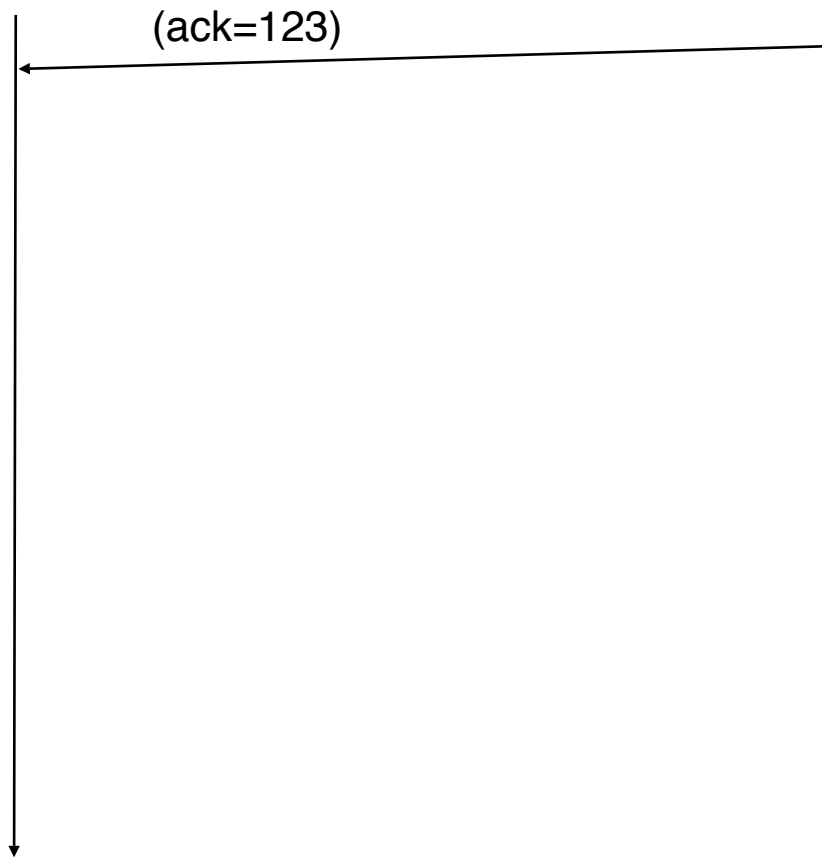
# Improving the reliable data transfer

## How to retransmit the lost segments

Upon reception of three duplicate acks, retransmit <span style="color:red">the first unacked segment</span>
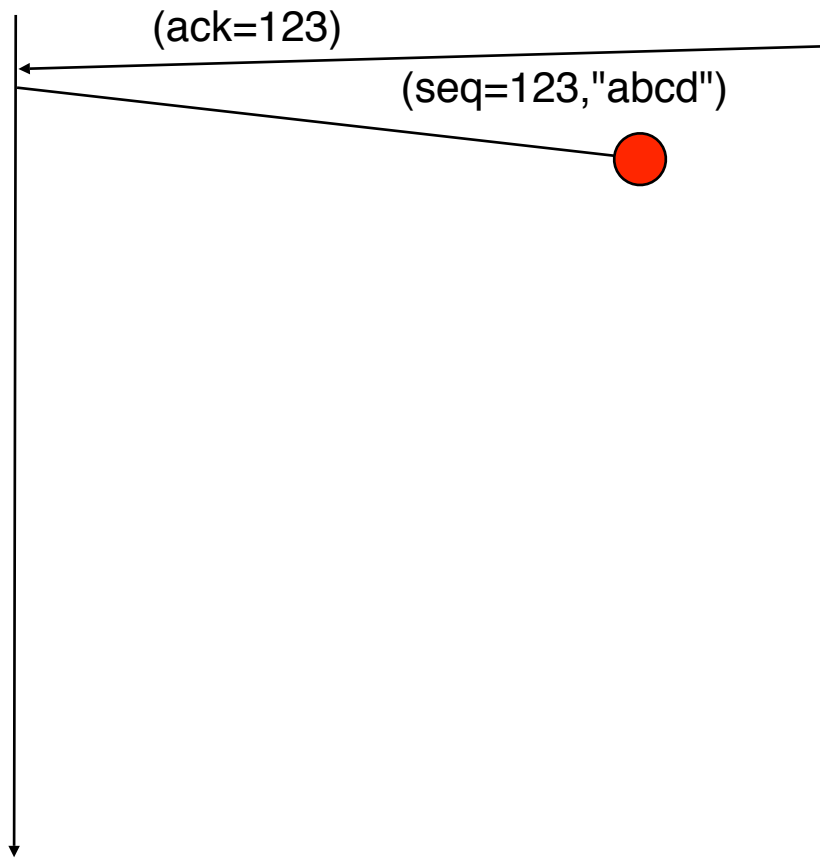
Fast retransmit, used by most TCP implementations

# Improving the reliable data transfer

## How to retransmit the lost segments

Upon reception of three duplicate acks, retransmit the first unacked segment
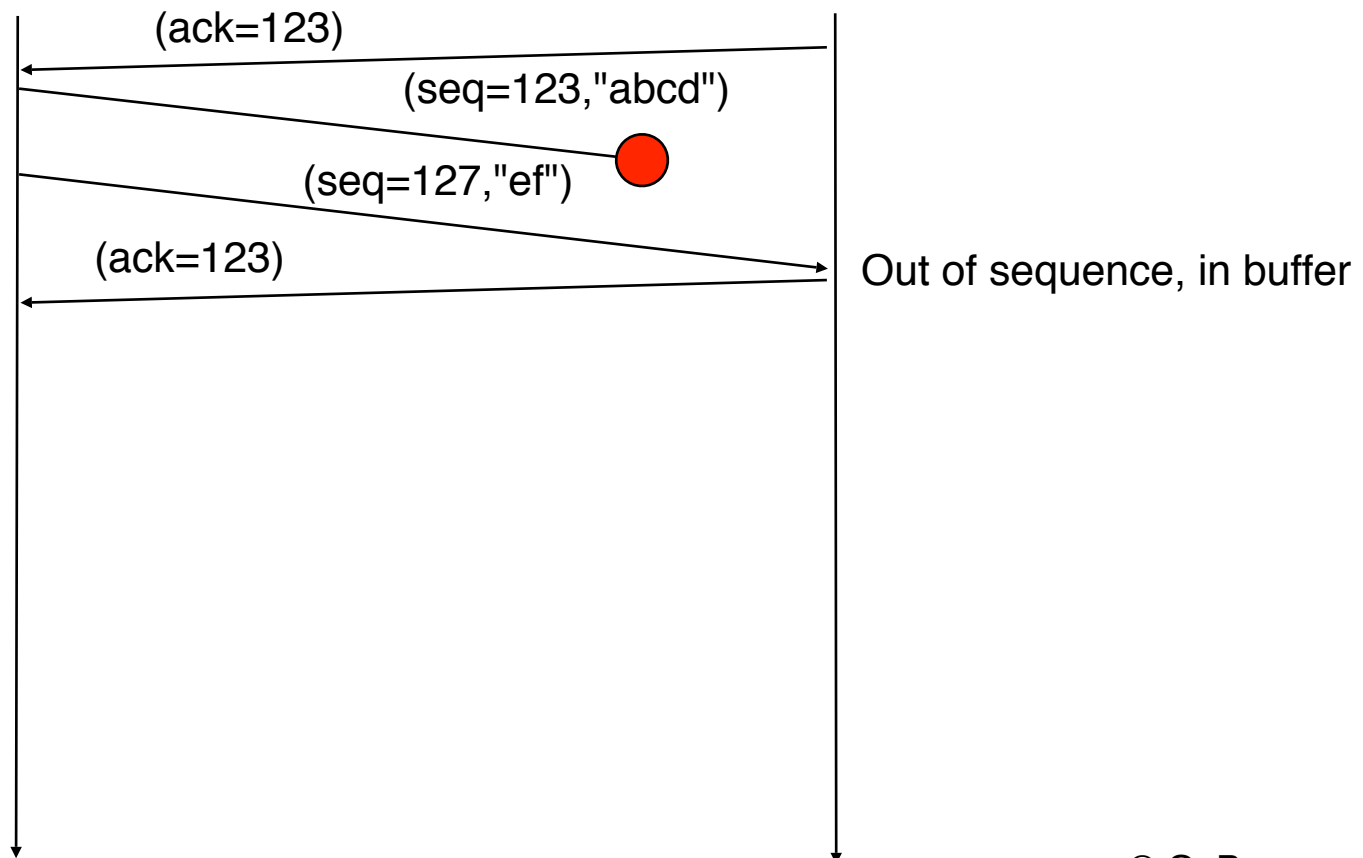
Fast retransmit, used by most TCP implementations

(ack=123)

# Improving the reliable data transfer

## How to retransmit the lost segments
Upon reception of three duplicate acks, retransmit
<u>the first unacked segment</u>
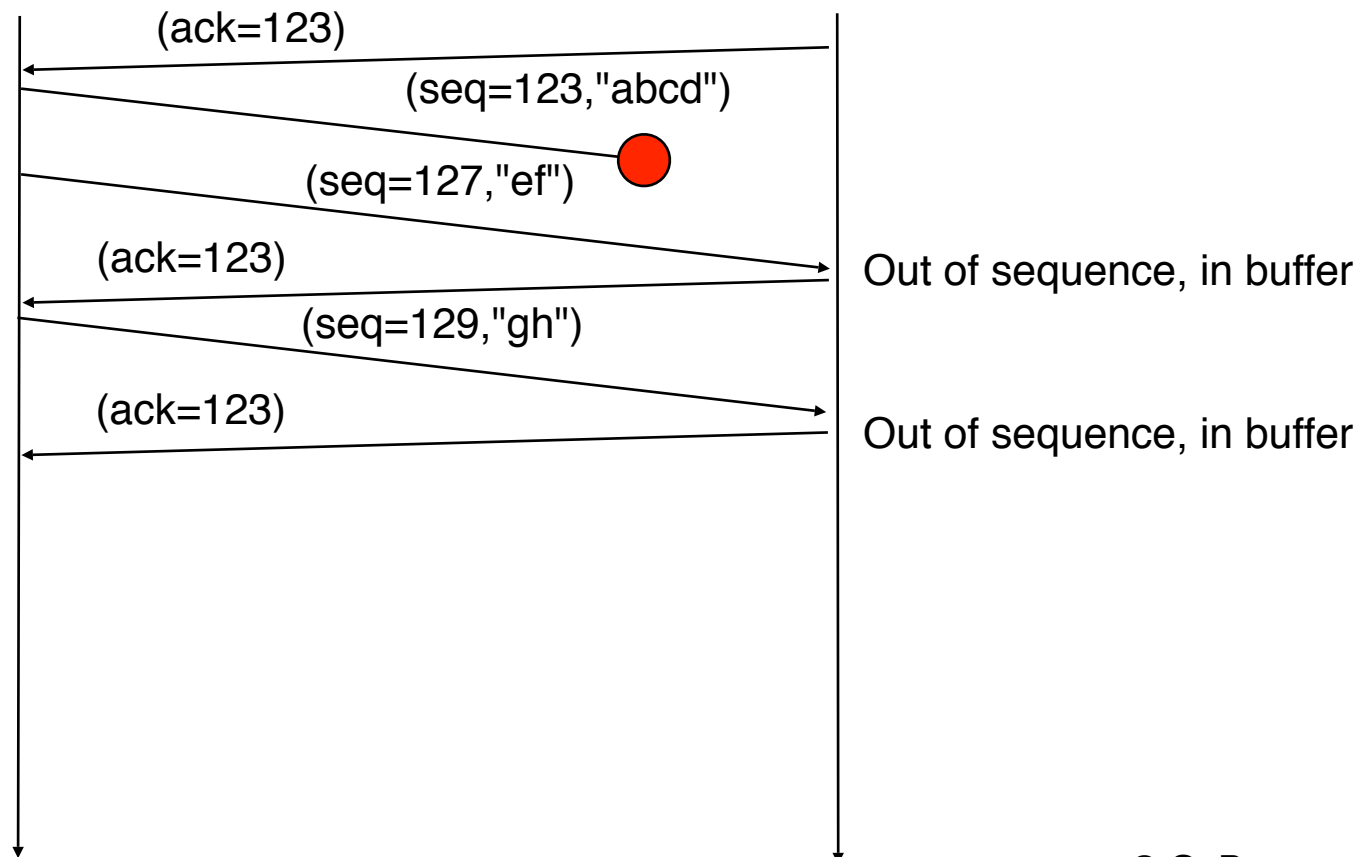Fast retransmit, used by most TCP implementations

# Improving the reliable data transfer

## How to retransmit the lost segments
Upon reception of three duplicate acks, retransmit <u>the first unacked segment</u>
Fast retransmit, used by most TCP implementations



(ack=123)

(seq=123,"abcd")

(seq=127,"ef")

(ack=123)

Out of sequence, in buffer

# Improving the reliable data transfer

## How to retransmit the lost segments

Upon reception of three duplicate acks, retransmit <u>the first unacked segment</u>
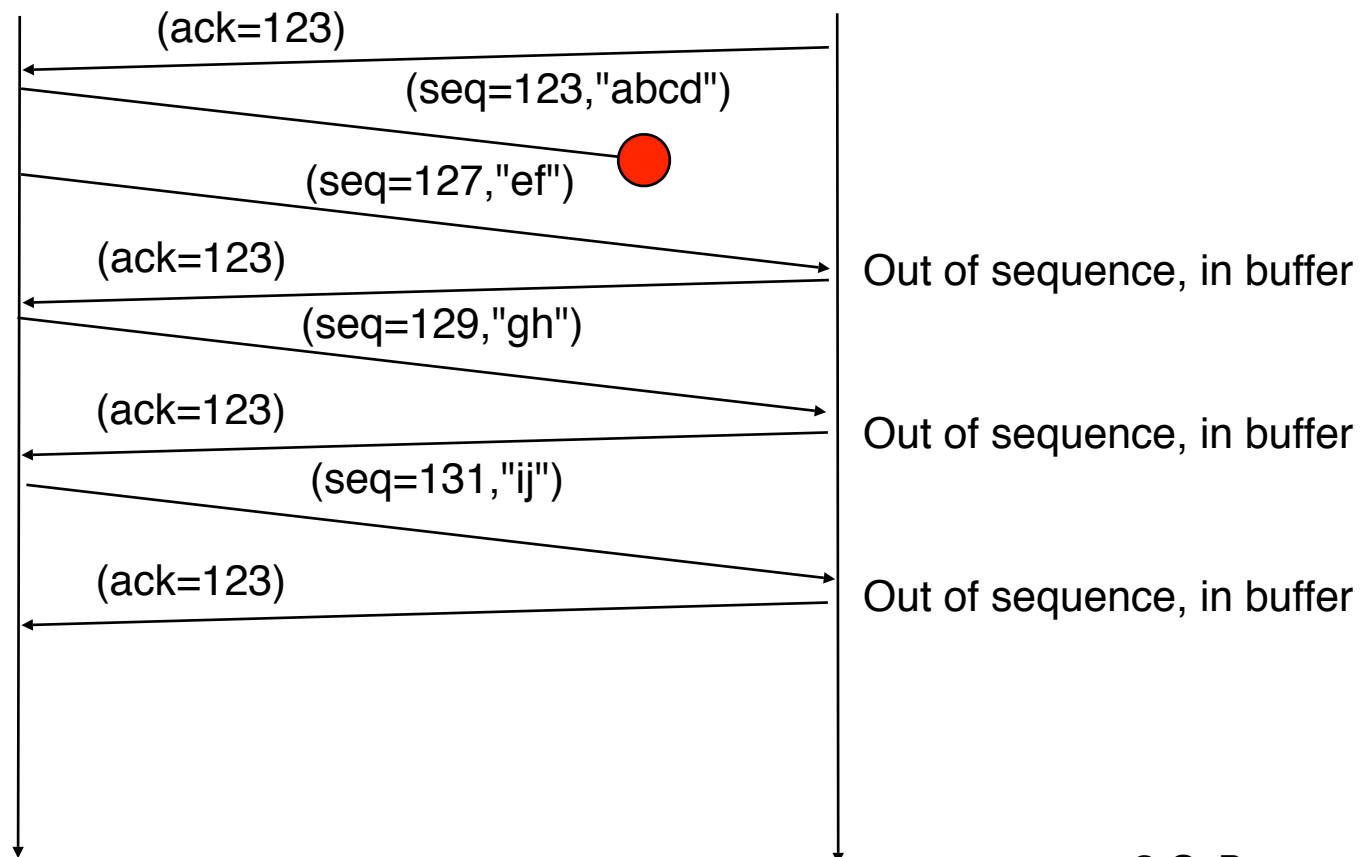
Fast retransmit, used by most TCP implementations

# Improving the reliable data transfer

## How to retransmit the lost segments
Upon reception of three duplicate acks, retransmit
<u>the first unacked segment</u>
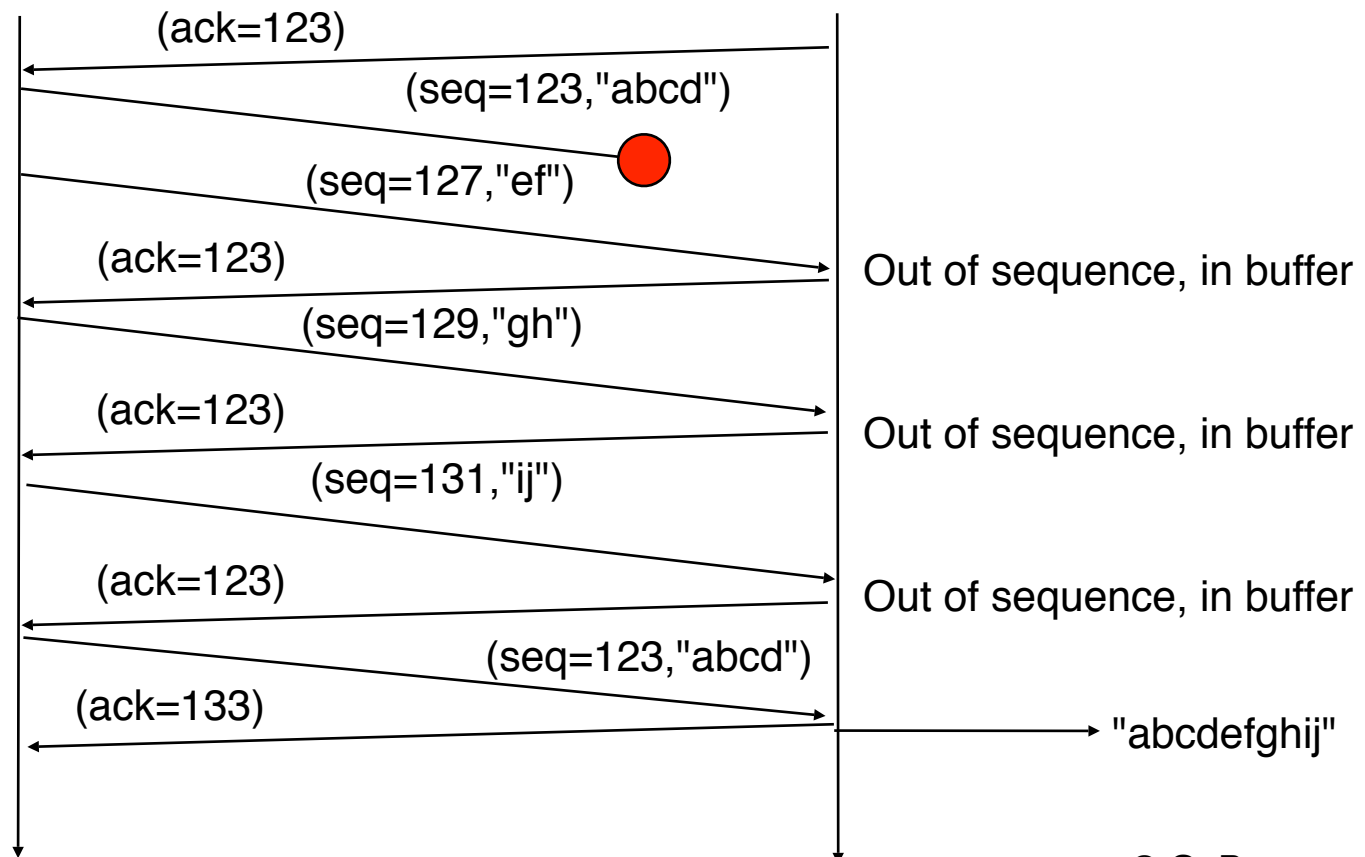Fast retransmit, used by most TCP implementations
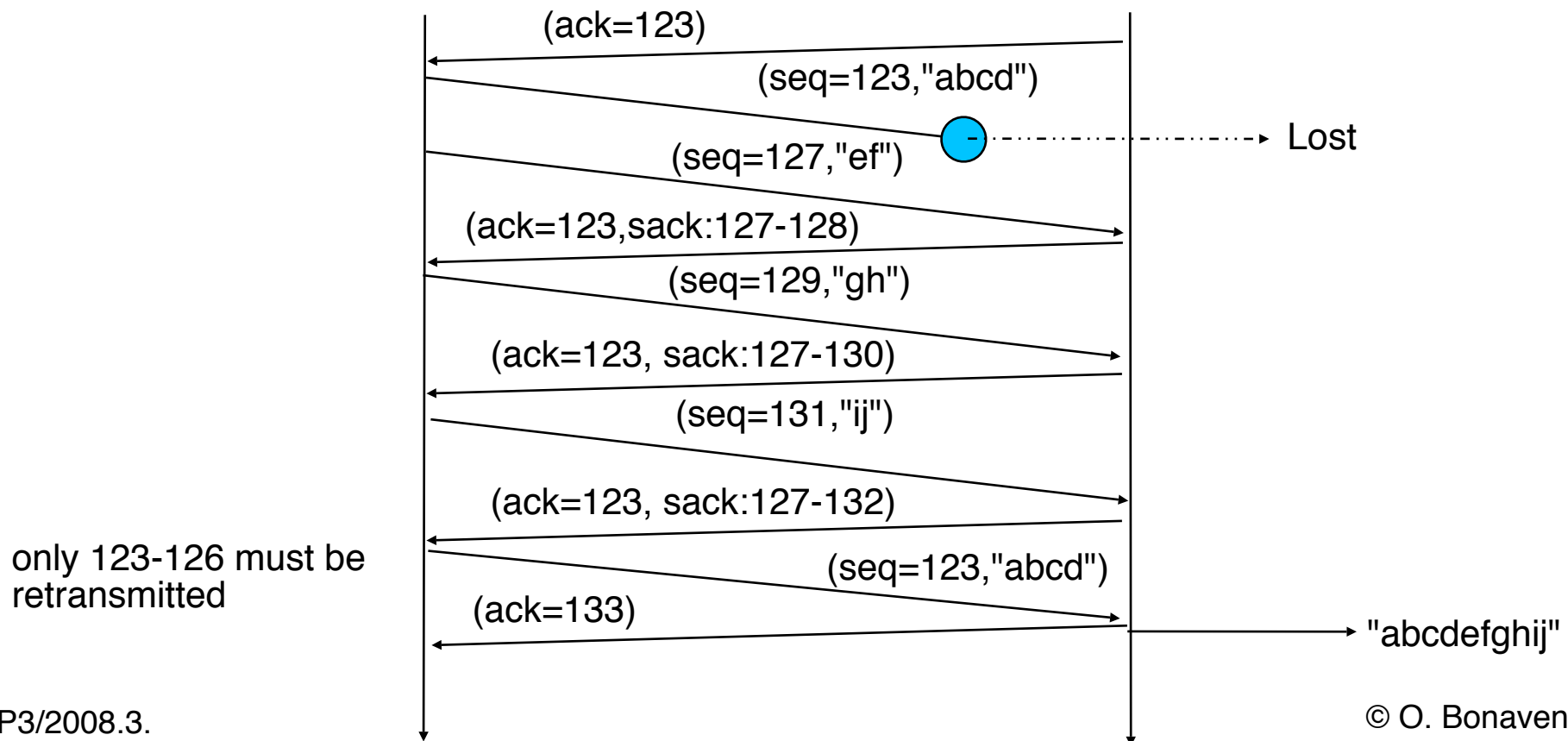
# Improving the reliable data transfer

How to retransmit the lost segments
   Upon reception of three duplicate acks, retransmit
   <u>the first unacked segment</u>
      Fast retransmit, used by most TCP implementations

© O. Bonaventure, 2008

# Improving the reliable data transfer

## Selective acknowledgement

sack:[seq1-seq2];[seq3-seq4]

# Sending acknowledgements

When to send a pure ack ?
   upon reception of a data segment
   inside data segments in the other direction
   (piggyback)

TCP tradeoff
   Insequence segment arrival
      If there is no ack waiting to be sent, start ack timer (50 msec) and wait wait until
         transmission of a data segment (piggyback)
         expiration of acktimer
      If there is already an ack waiting
         send pure ack immediately
   Out-of-sequence segment
      send ack immediately

# Flow control

Goal : protect the receiver's buffers

Principle

    Advertise receiving window in all segments

    State variables maintained by each TCP entity

       last_ack, swin, rwin

Last_ack=122, swin=100, rwin=4
To transmit : abcdefghijklm

# Flow control

Goal : protect the receiver's buffers
Principle
  Advertise receiving window in all segments
  State variables maintained by each TCP entity
    last_ack, swin, rwin
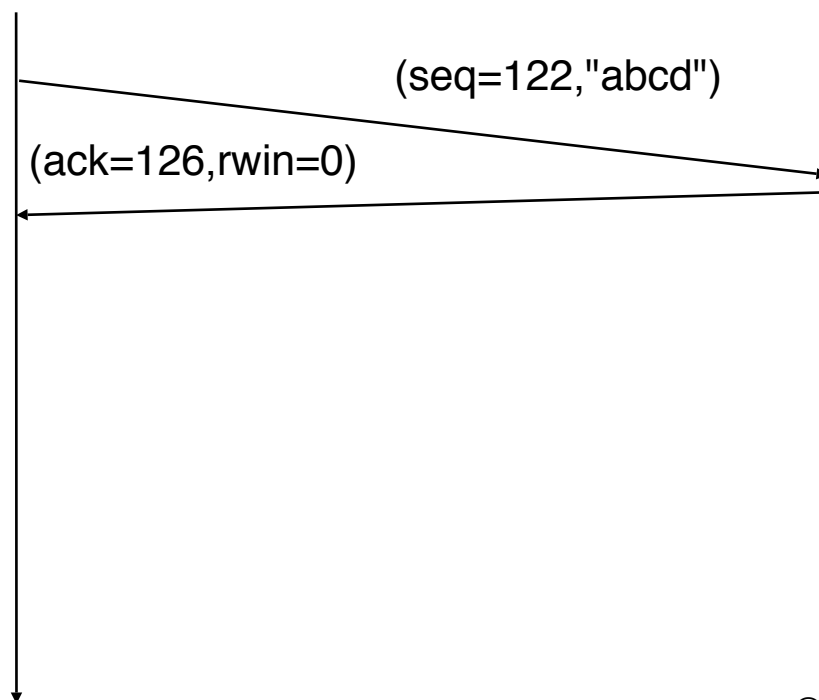
Last_ack=122, swin=100, rwin=4
To transmit : abcdefghijklm

(seq=122,"abcd")

Last_ack=122, swin=96, rwin=0

(ack=126,rwin=0)

Last_ack=126, swin=100, rwin=0

# Flow control

Goal : protect the receiver's buffers

Principle

Advertise receiving window in all segments

State variables maintained by each TCP entity

last_ack, swin, rwin
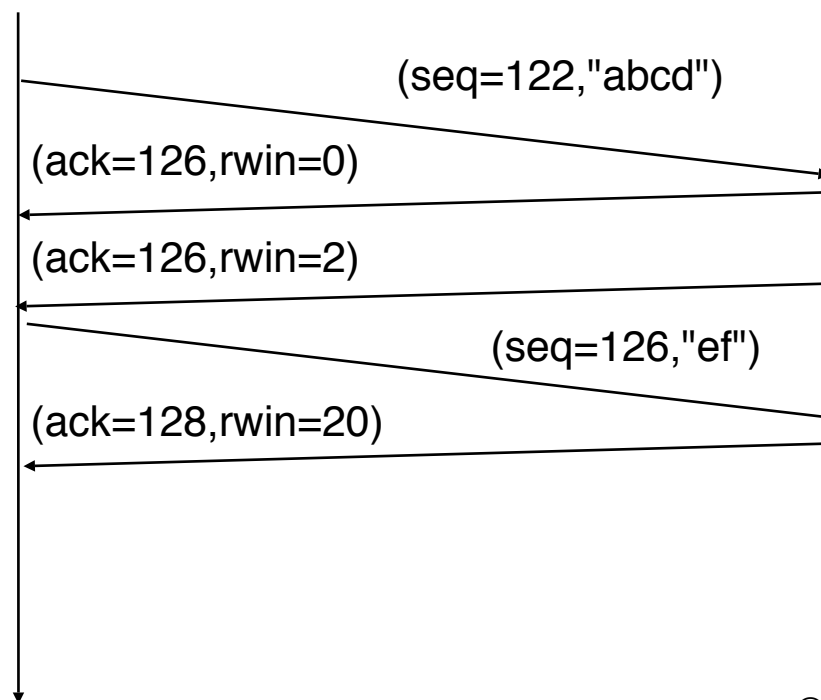
Last_ack=122, swin=100, rwin=4
To transmit : abcdefghijklm

Last_ack=122, swin=96, rwin=0

(seq=122,"abcd")

(ack=126,rwin=0)

Last_ack=126, swin=100, rwin=0

(ack=126,rwin=2)

Last_ack=126, swin=100, rwin=2
Last_ack=126, swin=98, rwin=0

(seq=126,"ef")

(ack=128,rwin=20)

Last_ack=128, swin=100, rwin=20

# Flow control

Goal : protect the receiver's buffers

Principle

Advertise receiving window in all segments

State variables maintained by each TCP entity

last_ack, swin, rwin

Last_ack=122, swin=100, rwin=4
To transmit : abcdefghijklm

Last_ack=122, swin=96, rwin=0
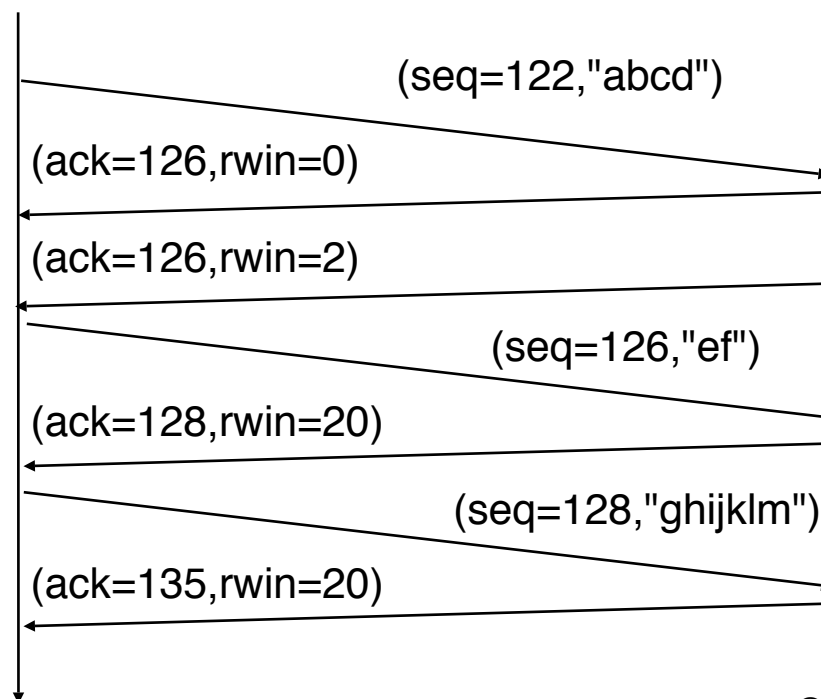
(seq=122,"abcd")

(ack=126,rwin=0)

Last_ack=126, swin=100, rwin=0

(ack=126,rwin=2)

Last_ack=126, swin=100, rwin=2
Last_ack=126, swin=98, rwin=0

(seq=126,"ef")

(ack=128,rwin=20)

Last_ack=128, swin=100, rwin=20

Last_ack=128, swin=93, rwin=13

(seq=128,"ghijklm")

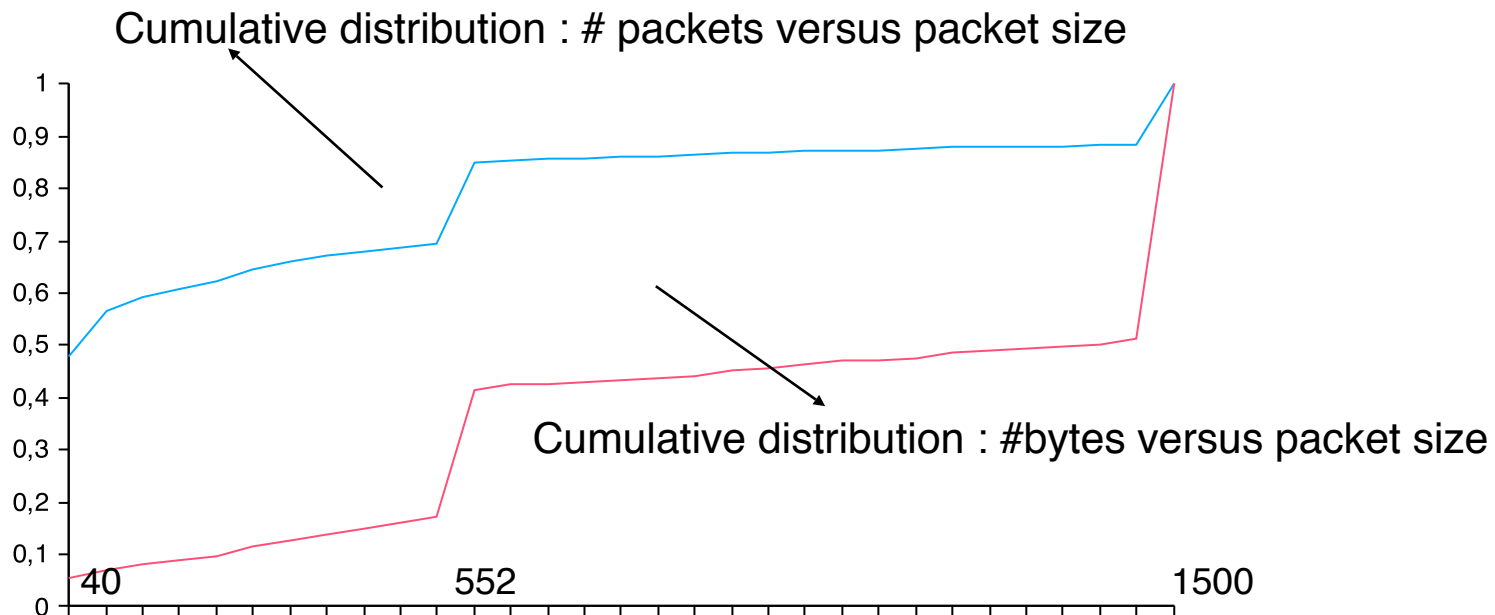(ack=135,rwin=20)

Last_ack=135, swin=100, rwin=20

# Flow control (2)

Limitations

TCP uses a 16 bits window field in the segment header

Maximum window size for normal TCP : 65535 bytes

Extension RFC1323 for larger windows

After having transmitted a window full of data, TCP sender must remain idle waiting for ack

Maximum throughput on TPC connection

~ window / round-trip-time

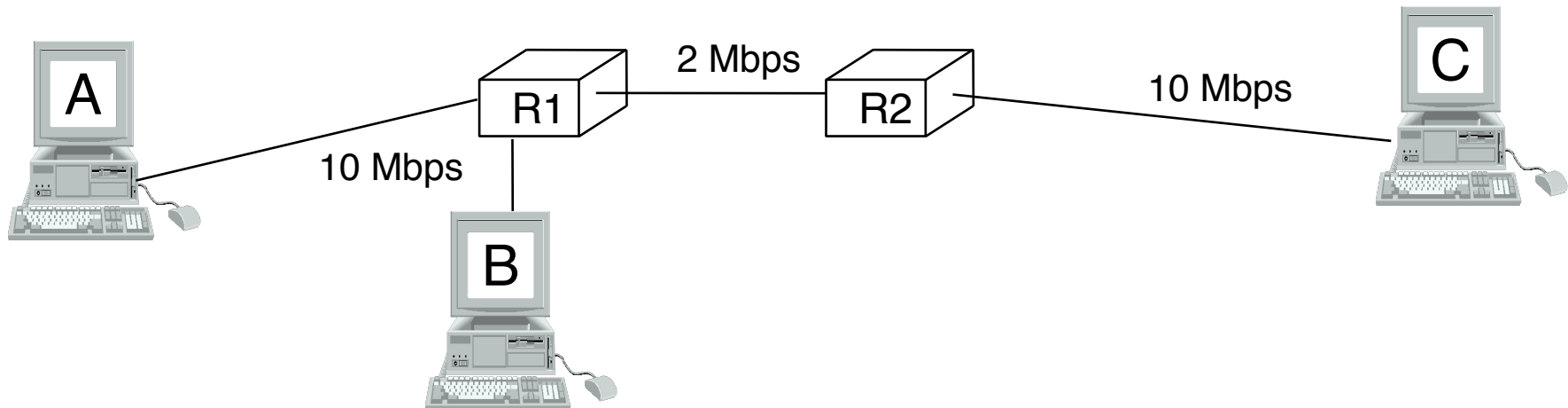| rtt Window | 1 msec | 10 msec | 100 msec |
|---|---|---|---|
| 8 Kbytes | 65.6 Mbps | 6.5 Mbps | 0.66 Mbps |
| 64 Kbytes | 524.3 Mbps | 52.4 Mbps | 5.2 Mbps |

# Transmission of data and control segments

## Nagle algorithm
### A new data segment can be sent provided that
- This is a maximum sized segment (MSS bytes)
- There are currently no unacknowledged bytes

## Consequence
### Most TCP/IP packets are small or MSS-sized

Cumulative distribution : # packets versus packet size

Cumulative distribution : #bytes versus packet size

# Module 3 : Transport Layer

Basics

Building a reliable transport layer

UDP : a simple connectionless transport protocol

<span style="color:red">TCP : a reliable connection oriented transport protocol</span>
    TCP connection establishment
    TCP connection release
    Reliable data transfer
→    <span style="color:red">Congestion control</span>
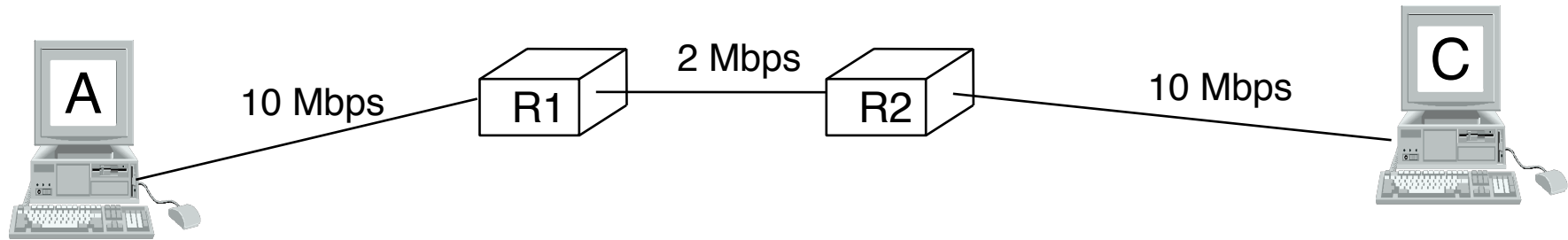
# Congestion in TCP/IP networks



**TCP/IP networks are heterogeneous**
A can send at 10 Mbps to B
B can send at 2 Mbps to C

**How to share the network among multiple hosts ?**
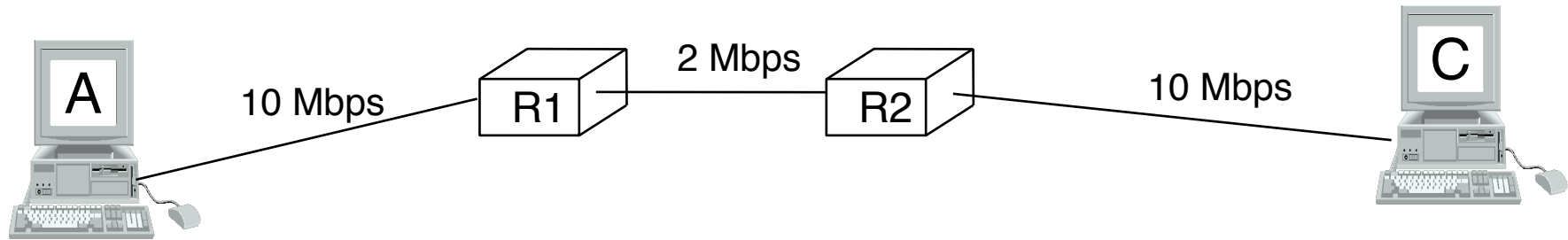A and B send data to C at the same time

# Congestion in TCP/IP networks

Possible solutions

The network indicates explicitly the bandwidth allocated to each host

network sends regularly control information to hosts

Available Bit Rate in ATM networks

Endhosts measure the state of the network and adapt their bandwidth to the network state

Endhosts must be able to measure the amount of congestion inside the network
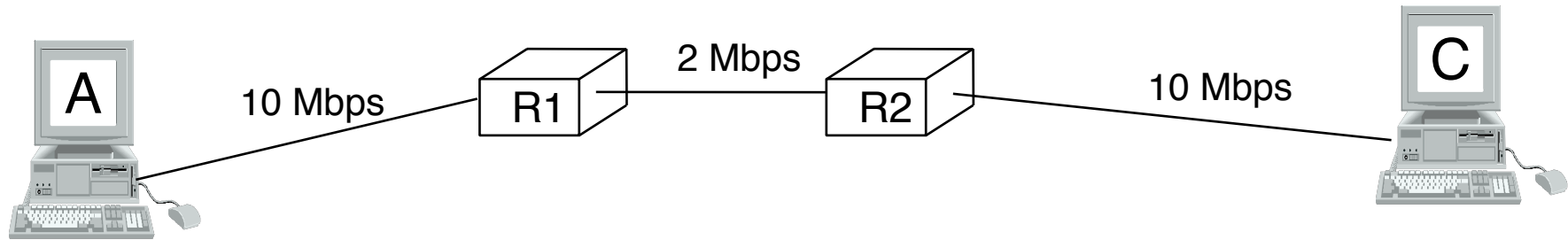
Solution used by TCP in the Internet

# Simple congestion



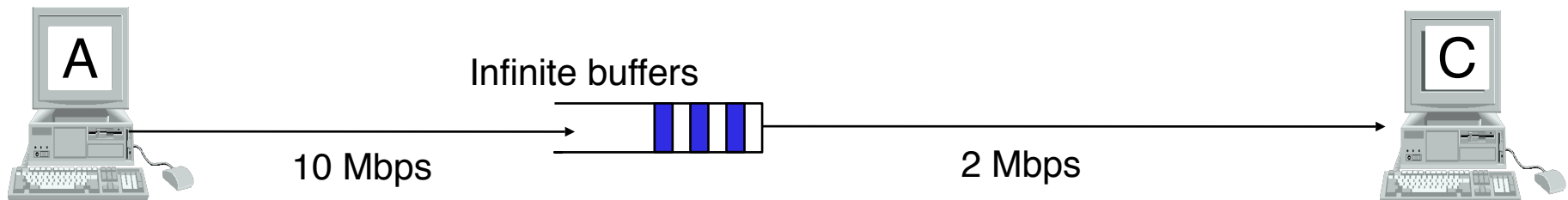A —— 10 Mbps —— R1 —— 2 Mbps —— R2 —— 10 Mbps —— C

# Simple congestion

A —— 10 Mbps —— R1 —— 2 Mbps —— R2 —— 10 Mbps —— C

## Simplified model

# Simple congestion



10 Mbps    R1    2 Mbps    R2    10 Mbps    C

A

## Simplified model



A    Infinite buffers    C

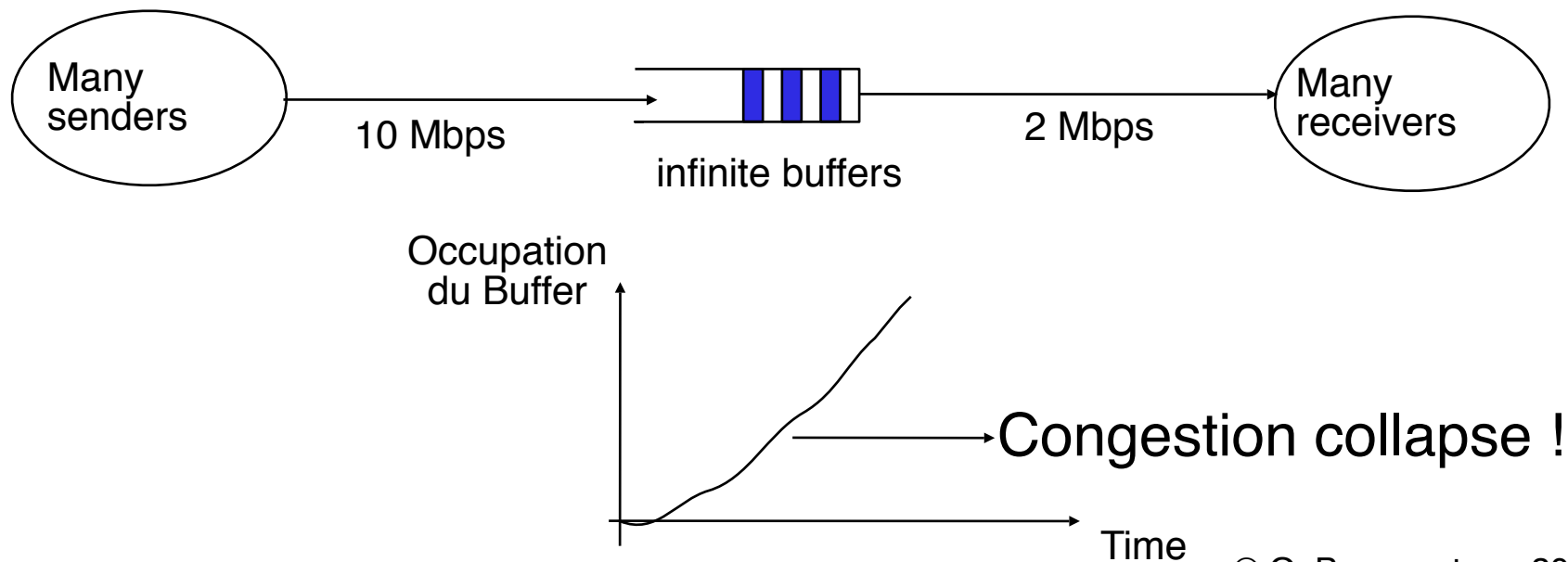10 Mbps    2 Mbps

# Simple congestion



Buffers de R1

10 Mbps

2 Mbps

## TCP self-clocking

# Simple congestion

## TCP self-clocking

Can be sufficient when a single TCP connection uses a low bandwidth link if the intermediate buffer can store a window full of segments
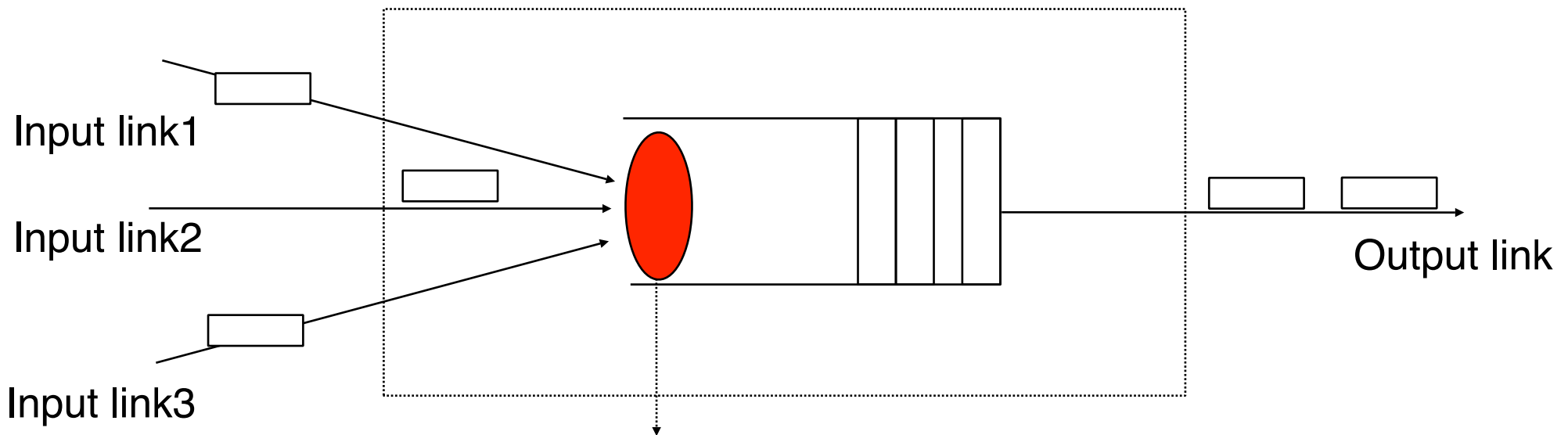What happens if several TCP connections need to share one link

# TCP congestion control

How to adapt a TCP connection to the network state ?

How to measure the current congestion state ?

TCP uses segment losses in routers as an implicit indication of congestion

This is valid in most environments besides some wireless networks where transmission errors can cause segment losses

Adapt the bandwidth of the TCP connection

TCP adapts its transmission rate by using a new congestion window (cwnd) which is controlled by the sender based on the current congestion status
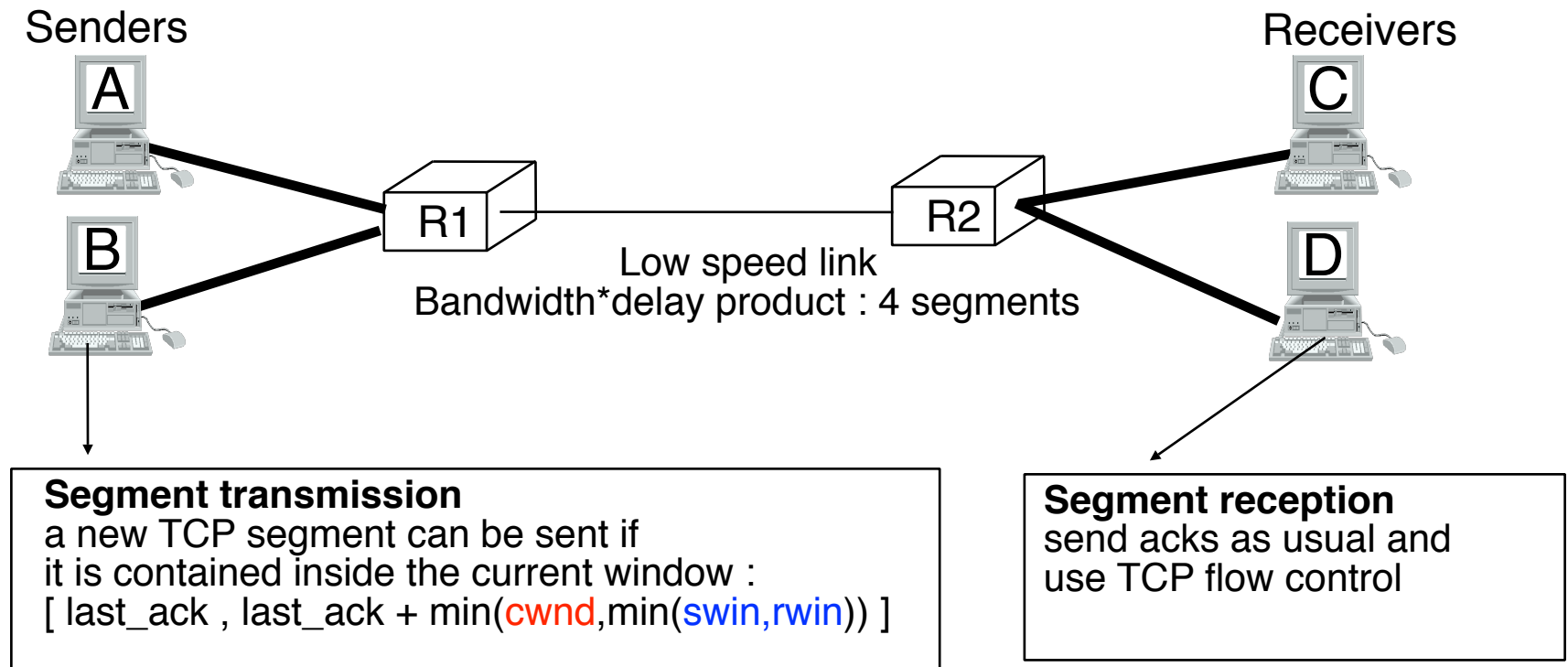
# A simple router

## Output buffer



**Buffer acceptance algorithm**
When a packet arrives in the output buffer, decides whether the packet is accepted or discarded

## taildrop
the arriving packet is discarded if the buffer is full
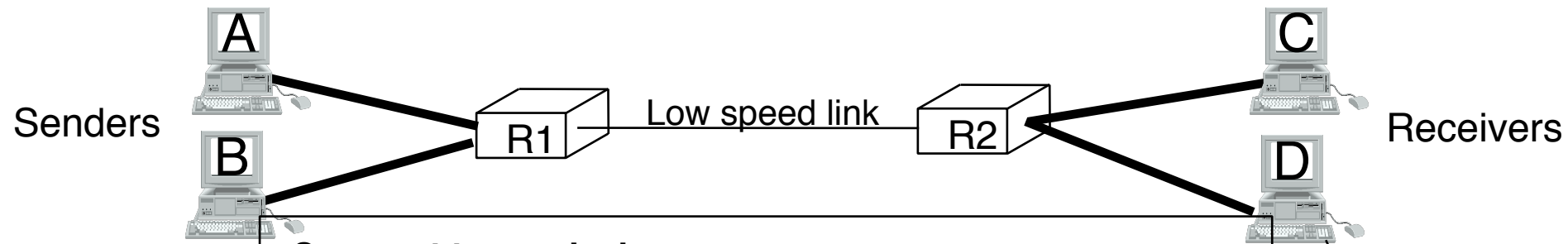
# TCP congestion control

Senders

A

B

R1

R2

Low speed link
Bandwidth*delay product : 4 segments

Receivers

C

D

**Segment transmission**
a new TCP segment can be sent if
it is contained inside the current window :
[ last_ack , last_ack + min(cwnd,min(swin,rwin)) ]

**Segment reception**
send acks as usual and
use TCP flow control

| Cwnd[B] cwnd[A] | 1 | 2 | 3 |
|---|---|---|---|
| 1 | underused | underused | ok, but unfair |
| 2 | underused | ok, fair | congestion |
| 3 | ok, but unfair | congestion | congestion |

## How to dynamically update cwnd ?

# TCP congestion control

---

## Additive Increase / Multiplicative Decrease



Senders

A

B

R1    Low speed link    R2

Receivers

C

D

**Segment transmission**
a new TCP segment can be sent if
it is contained inside the current window :
[ last_ack , last_ack + min(cwnd,min(swin,rwin)) ]

**Ack reception : Additive Increase**
// <u>congestion avoidance</u>
if (network not congested ) // no segment losses
{  // increase slowly cwnd
   // increase cwnd by one mss every rtt
   cwnd = cwnd+ mss*(mss/cwnd);
}
**Congestion : Segment loss : Multiplicative decrease**
// retransmsit lost segments
ssthresh=max(cwnd/2,2*MSS);
cwnd=MSS;

**Segment reception**
send acks as usual and
use TCP flow control

# TCP congestion control

## How to select cwnd when connection starts ?
### Congestion avoidance increases cwnd slowly

**Initialisation :**
cwnd = MSS;
ssthresh= swin;

**Ack reception :**
if (network not congested ) // no segment losses
{
  if (cwnd < ssthresh)
   { // increase quickly cwnd
     // double cwnd  every rtt
     cwnd = cwnd+ MSS;
  }
  else
  { // increase slowly cwnd
     // increase cwnd by one mss every rtt
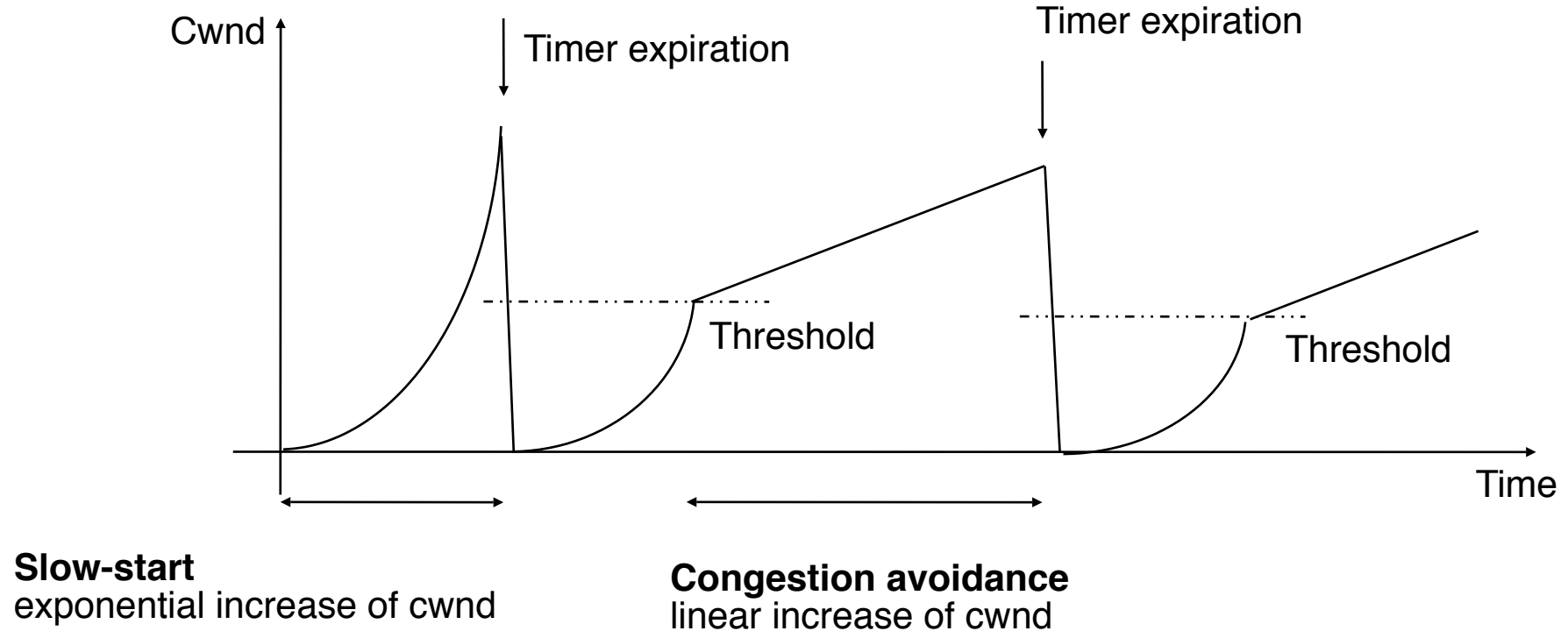     cwnd = cwnd+ mss*(mss/cwnd);
  }
}

Slowstart
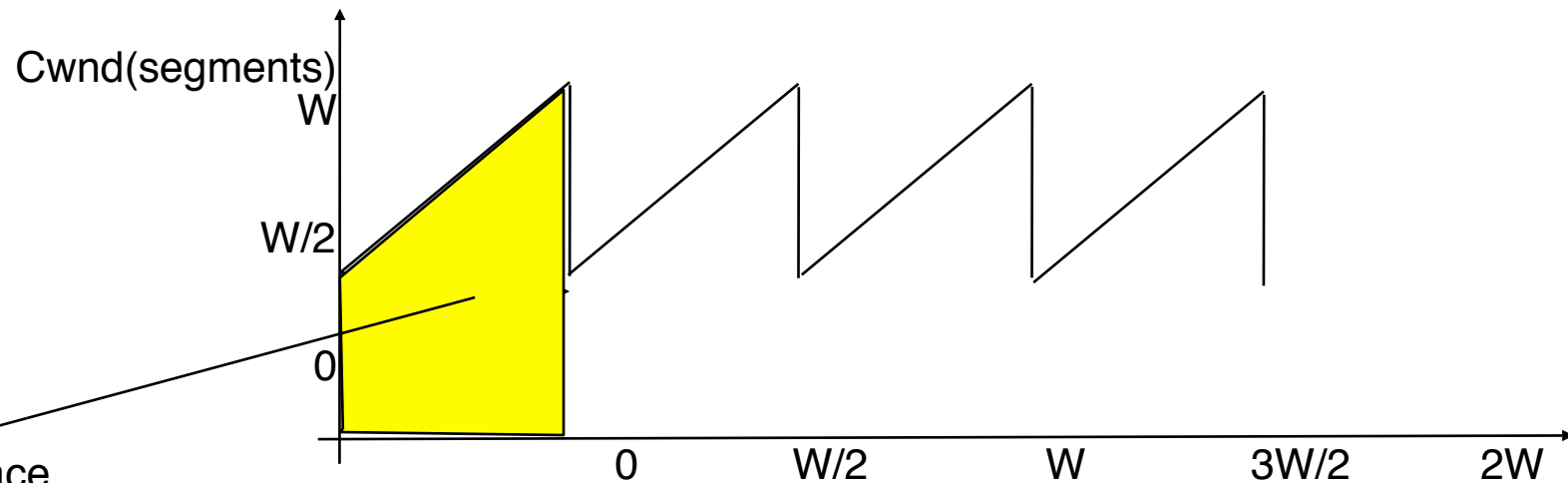
Congestion avoidance

# TCP congestion control

## Example

# TCP congestion control

How to react to segment losses ?

Two different types of multiplicative decrease

Severe loss [several lost segments]

wait until expiration of retransmission timer

sstresh=ssthresh/2

retransmit lost segments

slow-start (until cwnd=ssthresh)

congestion avoidance

Isolated loss [a single lost segment]

fast retransmit can recover from lost segment

If a single segment was lost : fast recovery

retransmit lost segment

sstresh=cwnd / 2

cwnd=ssthresh ; congestion avoidance

# TCP congestion control

## Simplified model
Assume that all segment losses are periodic and the every 1/p segment is lost



Cwnd(segments)
W

W/2

0

Surface

0    W/2    W    3W/2    2W    time(rtt)

It can be shown that the throughput of a TCP connection can be approximated by :

$$\left(\frac{W}{2}\right)^2 + \frac{1}{2}\left(\frac{W}{2}\right)^2 = \frac{1}{p}$$

$$BW < Min\left[\frac{Window}{RTT}, \left(\frac{MSS}{RTT}\right)\frac{k}{\sqrt{p}}\right]$$

Maximum throughput without losses

Throughput with losses/congestion