

# Transport Layer

## Felövervakning

### Övningar på Internet-checksumman

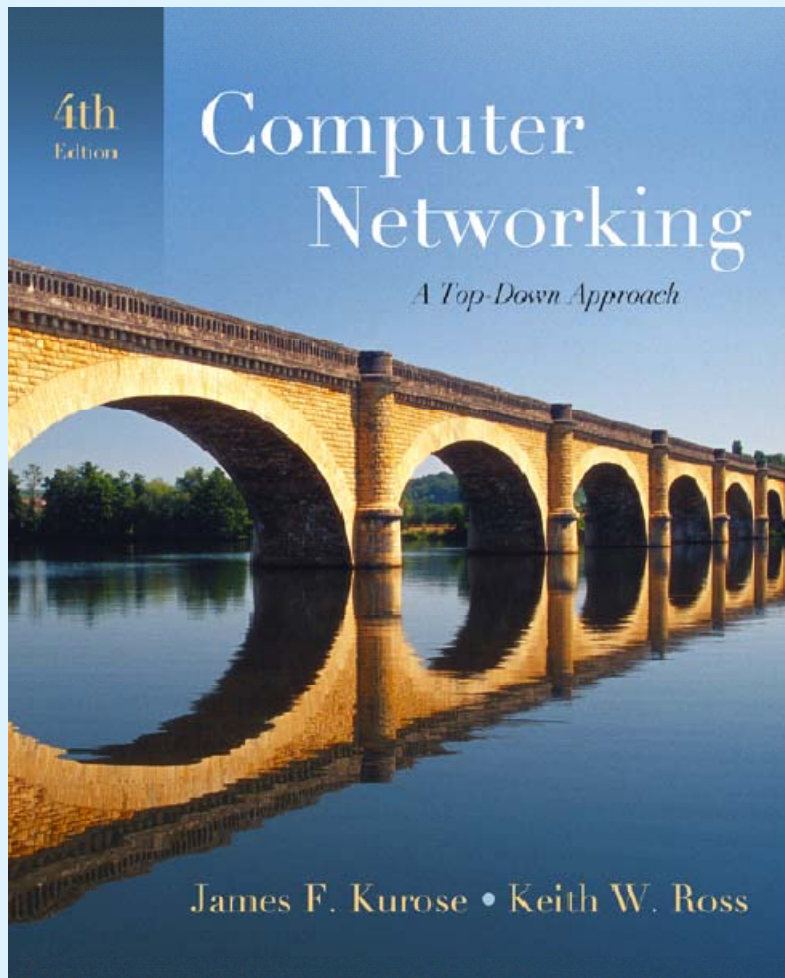
#### □ UDP

- Förbindelselöst
- "Best effort"

#### □ TCP

- Förbindelseorientat
- Tillförlitligt
- Flödesreglering
- Trafikstockningsreglering

Bildspelet omfattar till stor del bilder som hör till följande bok:



### A note on the use of these ppt slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- ☐ If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)
- ☐ If you post any slides in substantially unaltered form on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

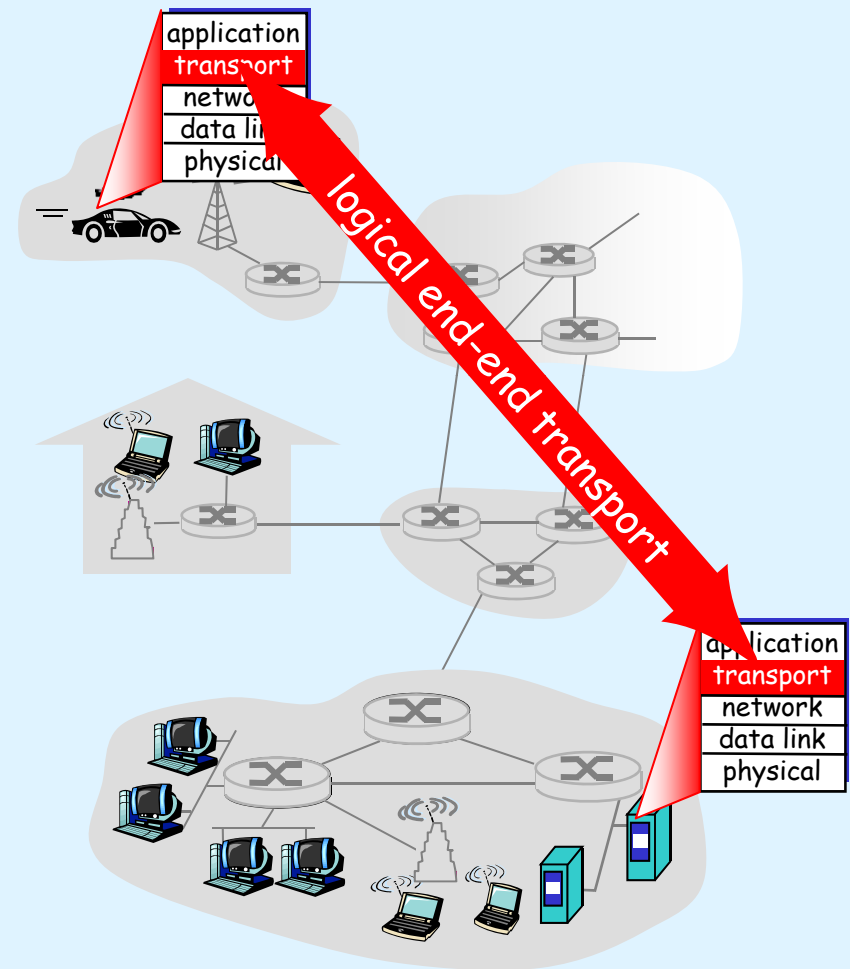
All material copyright 1996-2007  
J.F Kurose and K.W. Ross, All Rights Reserved

*Computer Networking: A Top Down Approach , 4<sup>th</sup> edition.*  
*Jim Kurose, Keith Ross, Addison-Wesley, July 2007.*

Transport Layer

# Transport services and protocols

- ❑ provide *logical communication* between app processes running on different hosts
- ❑ transport protocols run in end systems
  - send side: breaks app messages into *segments*, passes to network layer
  - receiver side: reassembles segments into messages, passes to app layer



# Transport vs. network layer

- *Transport layer:*

logical communication between processes

- relies on, enhances, network layer services

- *Network layer:*

logical communication between hosts

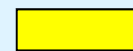
# Multiplexing/demultiplexing

## Multiplexing at send host:

gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

## Demultiplexing at rcv host:

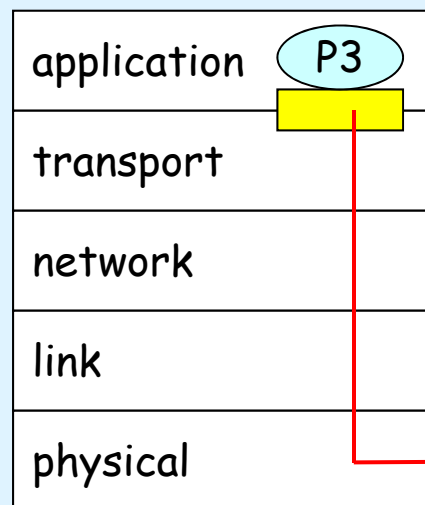
delivering received segments to correct socket



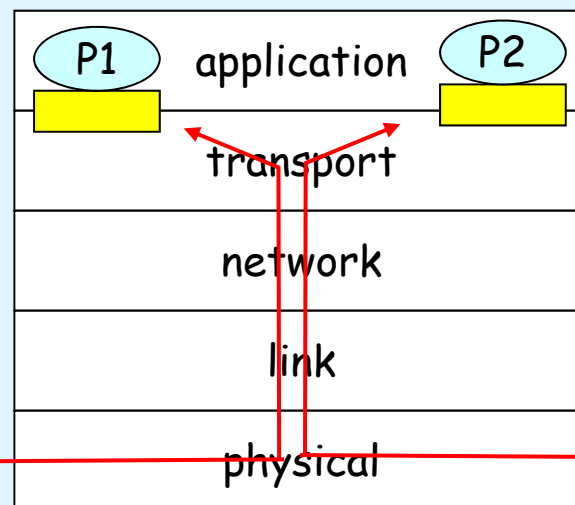
= socket



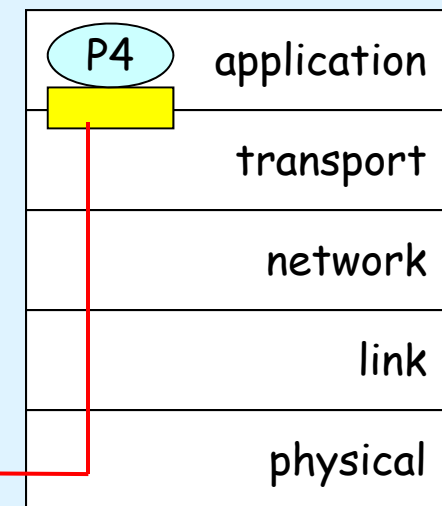
= process



host 1



host 2



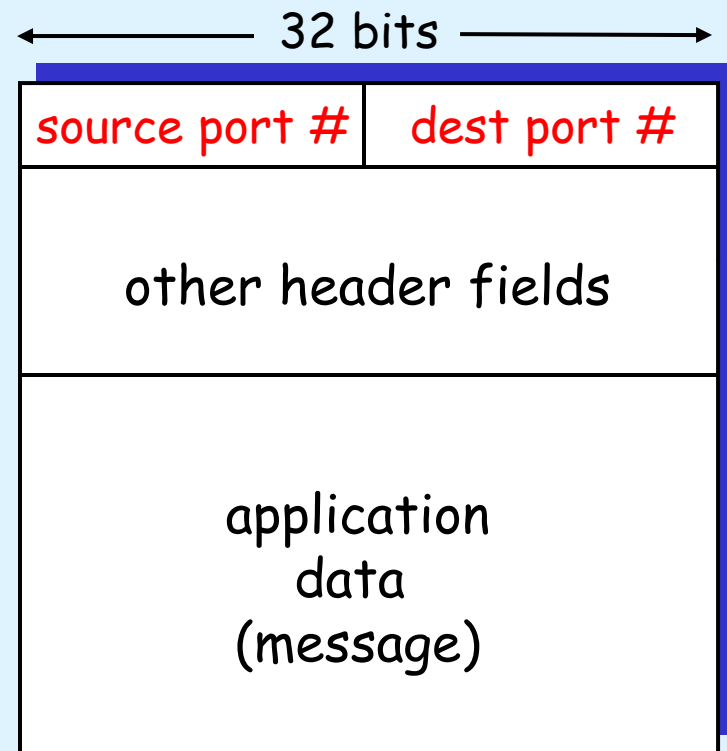
host 3

# How demultiplexing works

## □ host receives IP datagrams

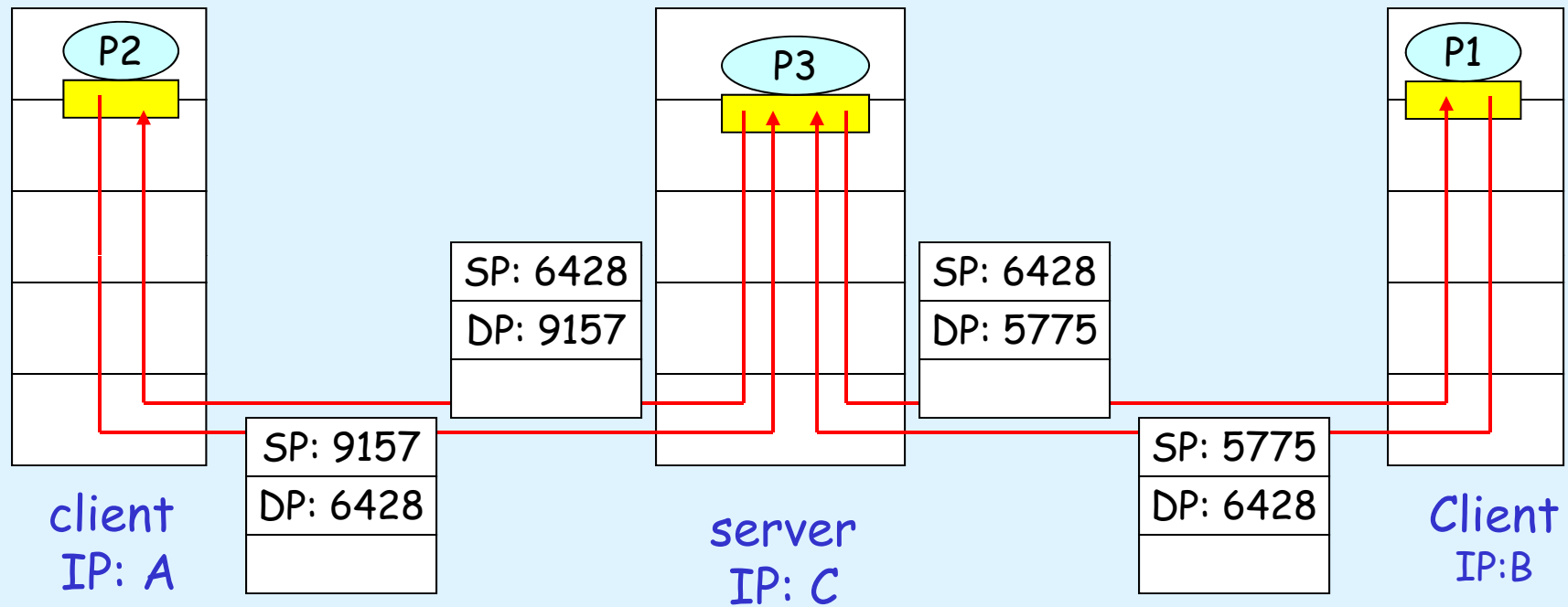
- each datagram has source IP address, destination IP address
- each datagram carries 1 transport-layer segment
- each segment has source, destination port number

## □ host uses IP addresses & port numbers to direct segment to appropriate socket



TCP/UDP segment format

# Connectionless demux (cont)



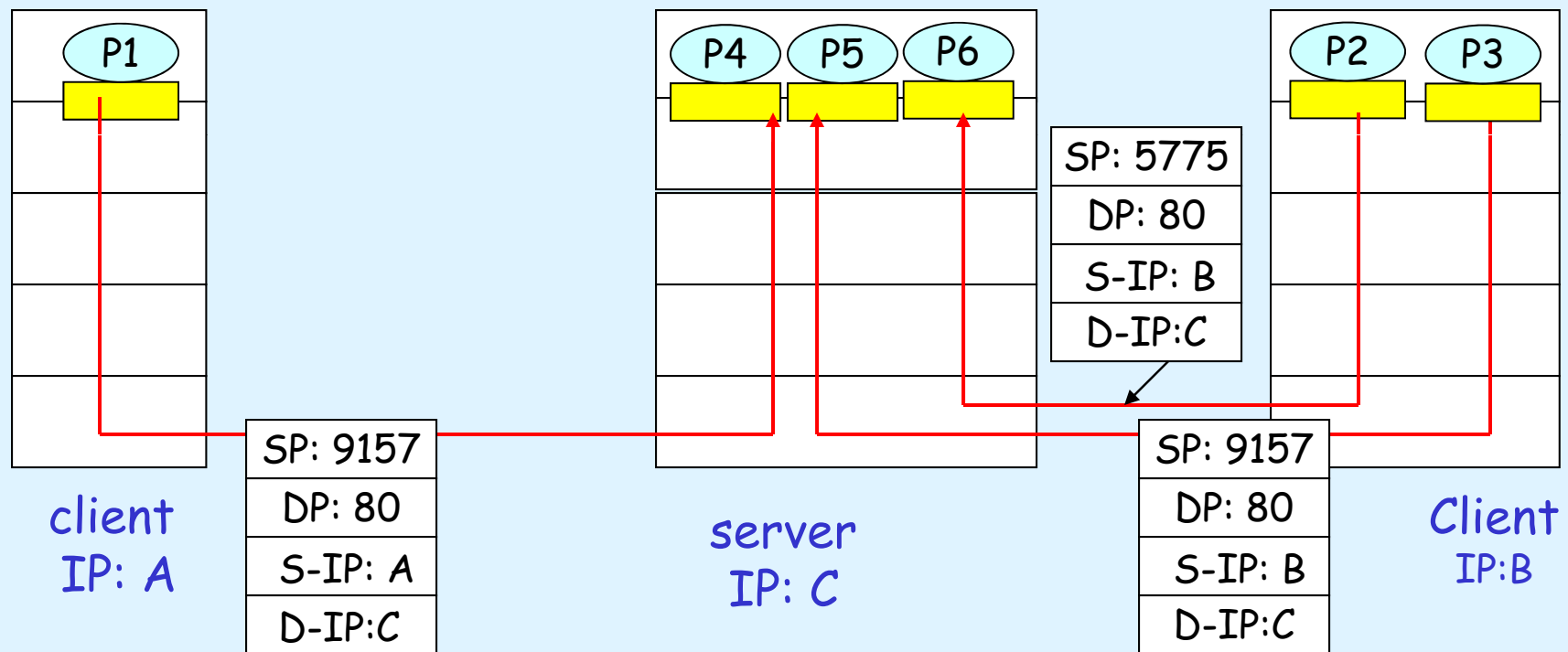
SP provides "return address"

# Connection-oriented demux

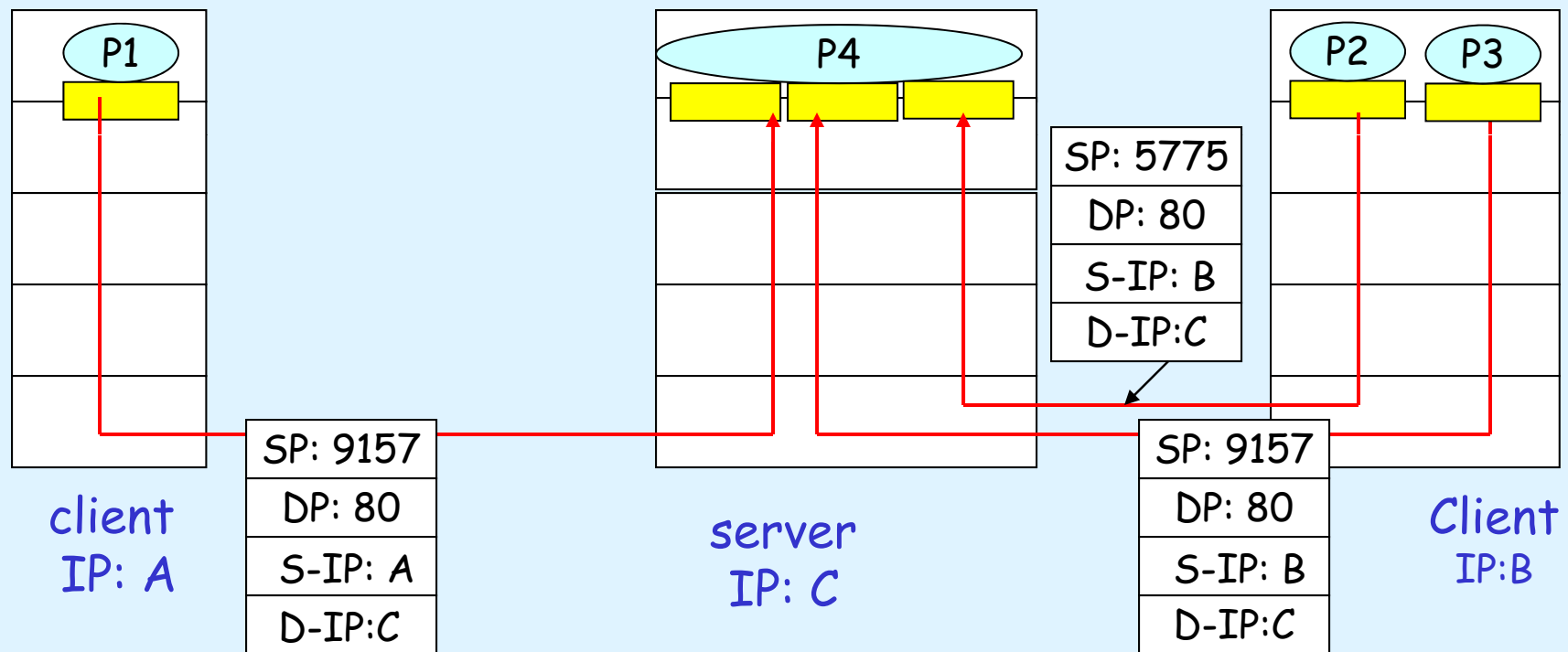
- ❑ TCP socket identified by 4-tuple:
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- ❑ recv host uses all four values to direct segment to appropriate socket
- ❑ Server host may support many simultaneous TCP sockets:
  - each socket identified by its own 4-tuple
- ❑ Web servers have different sockets for each connecting client
  - non-persistent HTTP will have different socket for each request



# Connection-oriented demux (cont)



# Connection-oriented demux: Threaded Web Server



# Felövervakning

(1)

- ❑ Sekvensnummer
- ❑ Checksumma
- ❑ Omsändningsmetod

# Felövervakning

(2)

- Manuell

Exempelvis Echo checking.

- Automatisk  
(Automatic Repeat Request, ARQ)

Upptäcka bitfel (feedback error control)

Korrigera bitfel (forward error control)

# Kvittenser för ARQ

- Positiva kvittenser

(Acknowledgement, ACK)

- Negativa kvittenser

(Negative Acknowledgement, NAK)

# Idle repeat request

- Implicit retransmission

Endast ACK.

- Explicit request

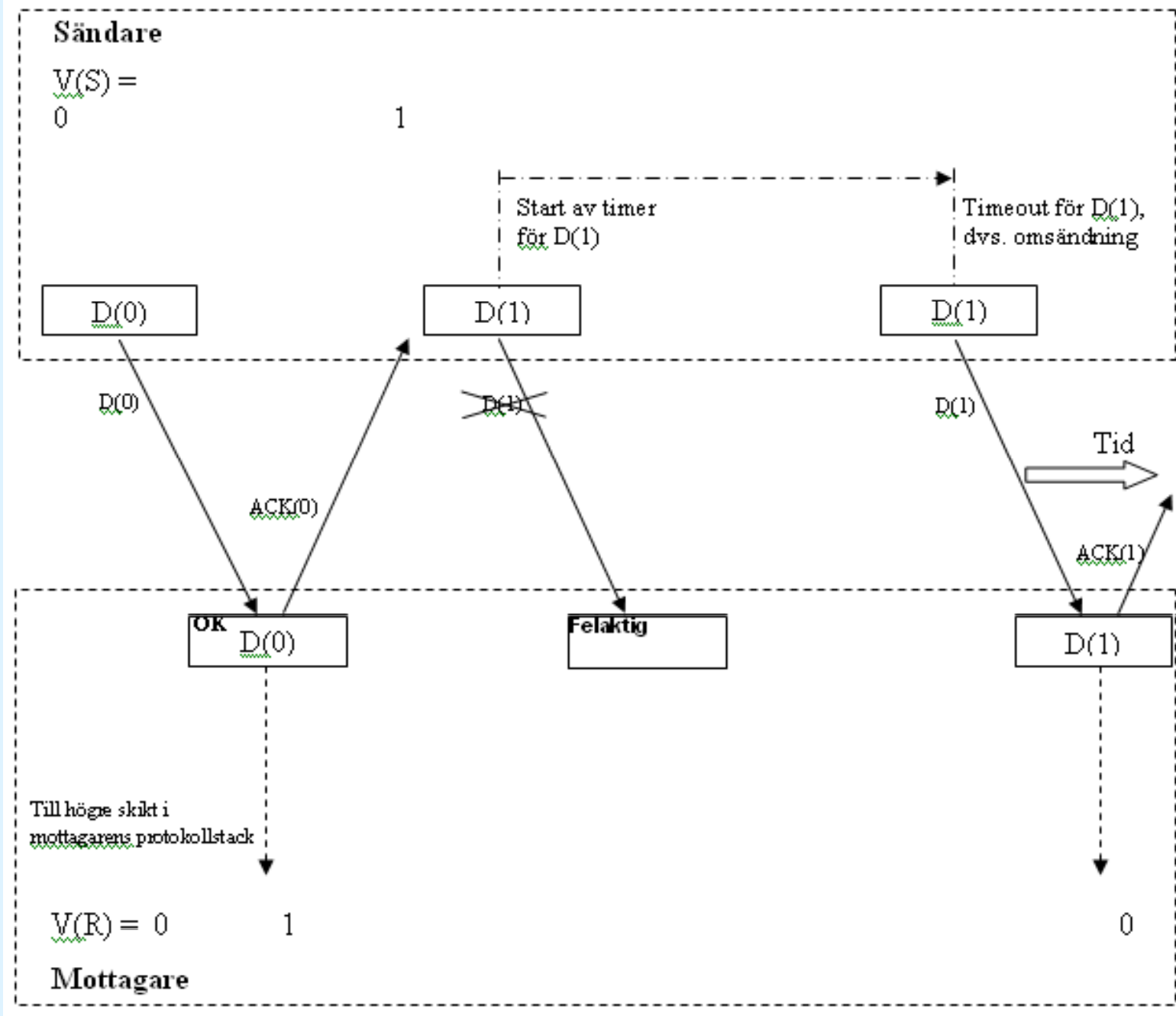
ACK och NAK.

# Implicit retransmission

- Sänd - vänta på ACK eller timeout

Om ACK, sänd nästa - ...

Om timeout, repetera föregående - ...





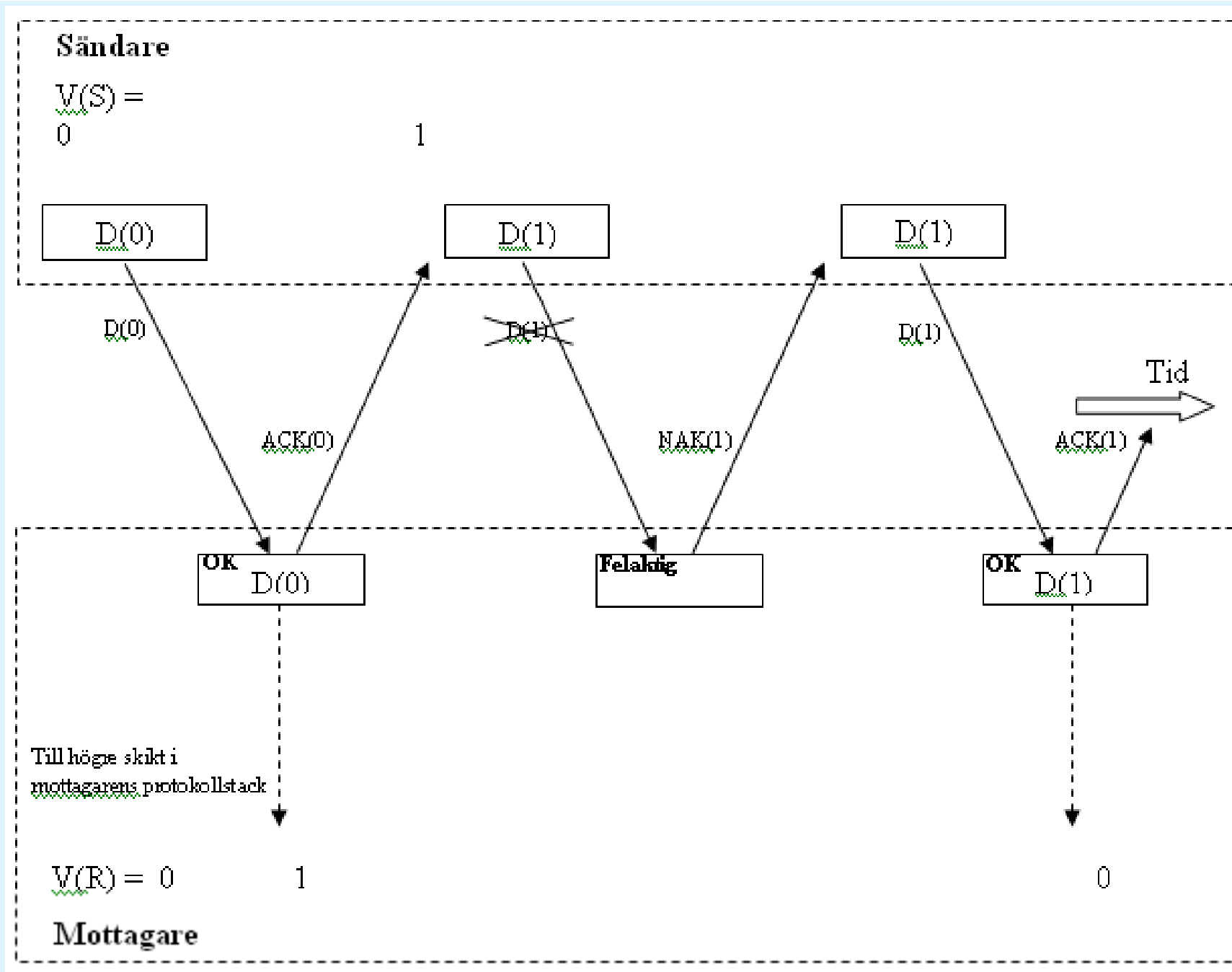
# Explicit request

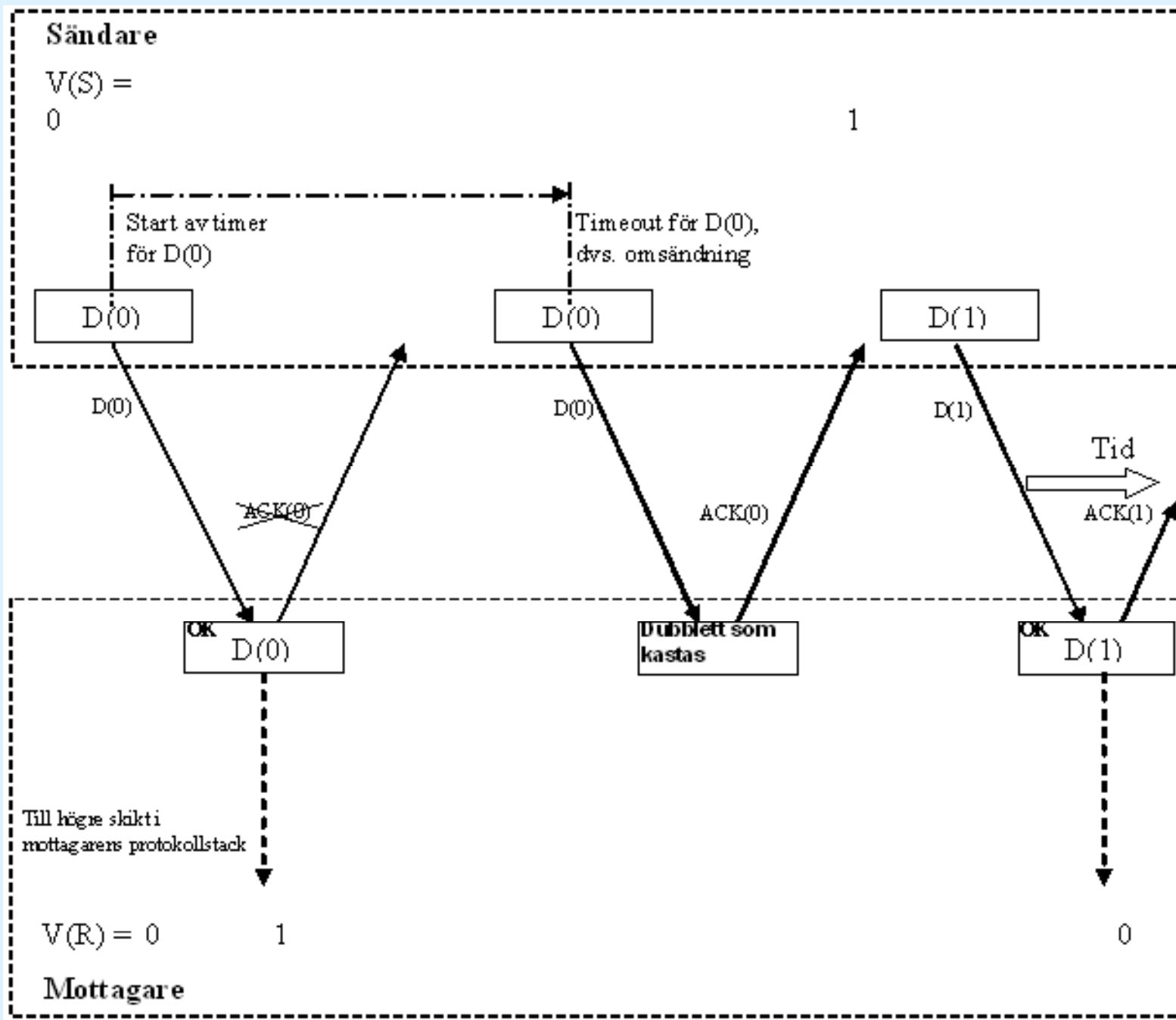
- Sänd - vänta på ACK, NAK eller timeout

Om ACK, sänd nästa - ...

Om NAK eller timeout, repetera föregående - ...

- NAK används av mottagaren för att snabbt få ett meddelande (ram, paket, segment etc.) repeterat.



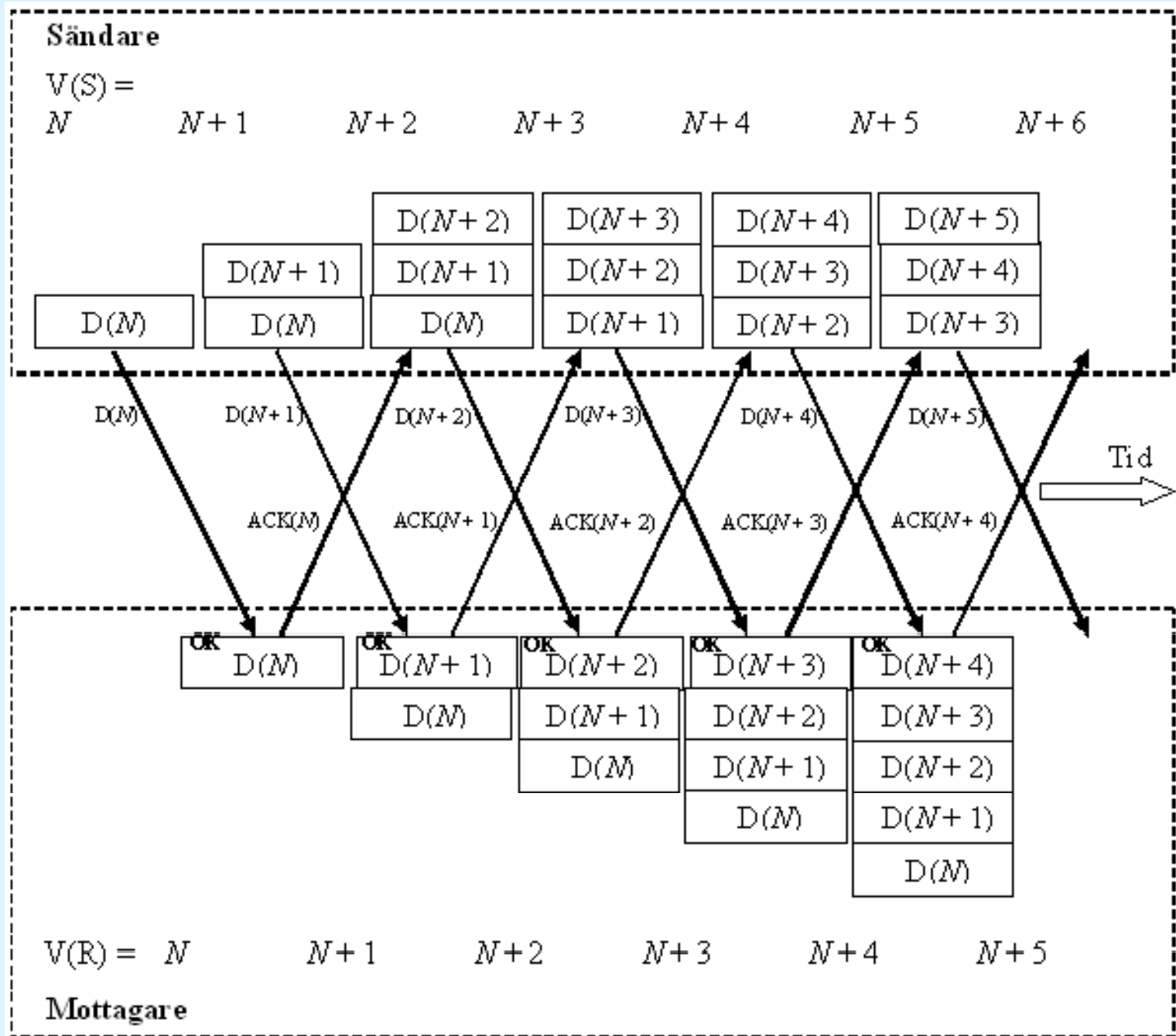


# Continuous repeat request (1)

- ❑ Omsändningsmetoderna med stora sändnings- och mottagningsbuffertar
- ❑ Skicka meddelanden (ramar, paket, segment etc.) som i strömmar utan att behöva invänta varje ACK, NAK eller timeout.

# Continuous repeat request (2)

- ❑ Selective repeat with implicit retransmission
- ❑ Selective repeat with explicit request
- ❑ Go-back-N



## Selective repeat with implicit retransmission (1)

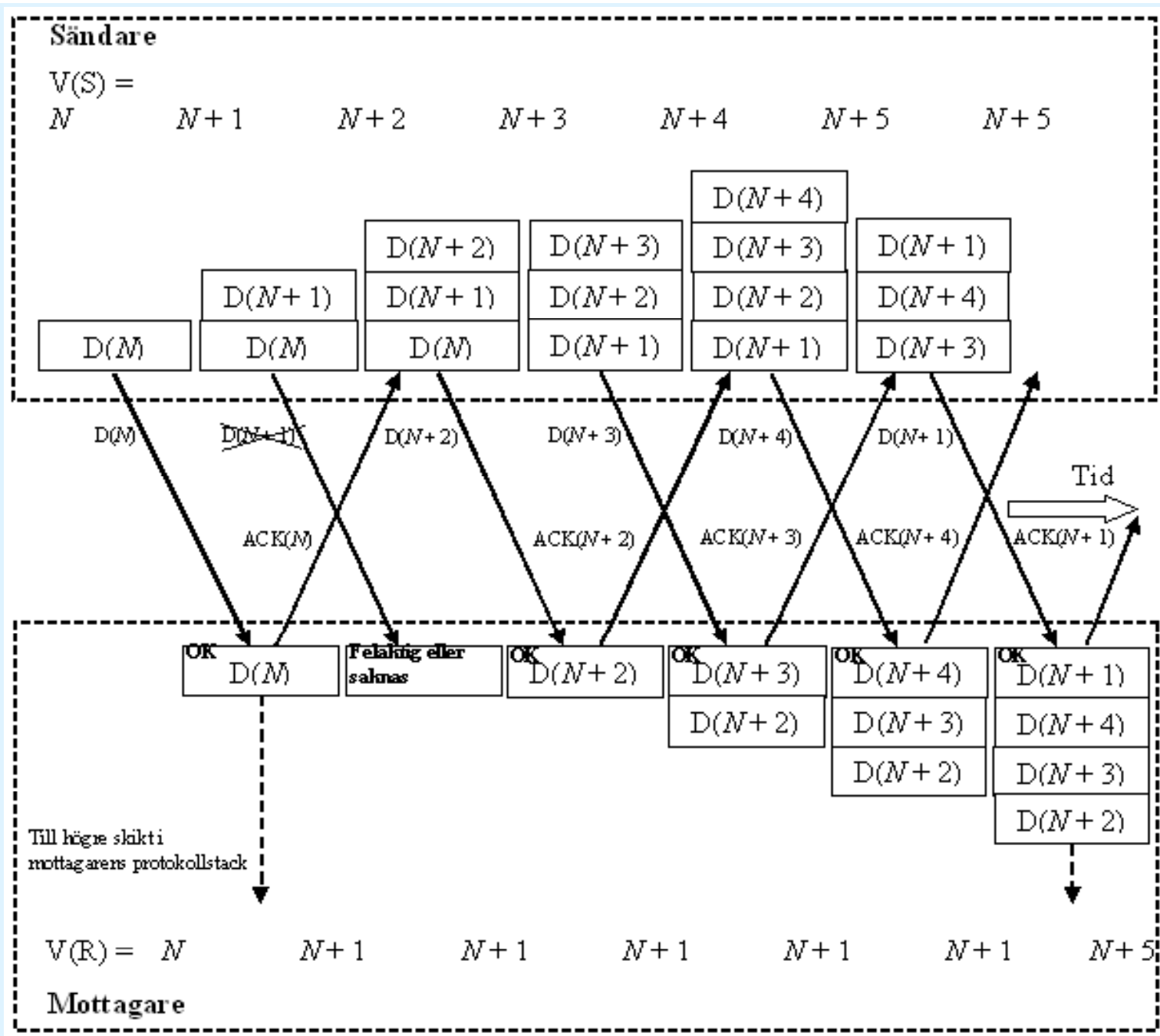
- ❑ Enbart felaktiga meddelanden (ramar, paket, segment etc.) sänds om vid begäran.
- ❑ Implicit retransmission betyder att NAK inte används.

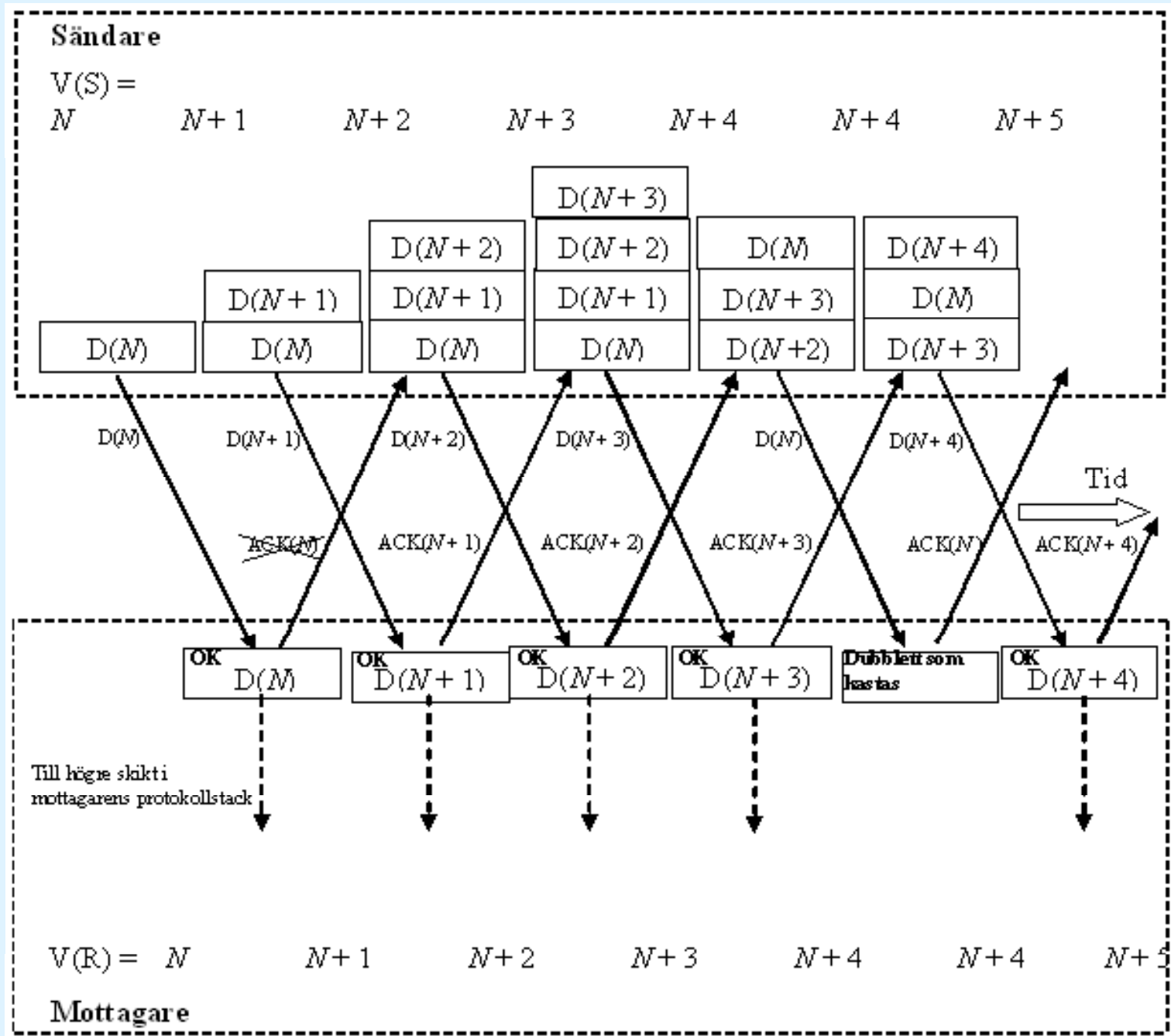
## Selective repeat with implicit retransmission (2)

Speciella egenskaper:

- ❑ Om meddelande  $N+1$  kvitteras innan meddelande  $N$ , sänds meddelande  $N$  på nytt.
- ❑ Fel på ACK uppfattas av sändaren som ingen ACK.
- ❑ Dubbla försändelser kan hanteras av mottagaren.
- ❑ Det behövs inga timers för att kontrollera ACK:ar.  
(Däremot kan timers användas för kunna att upptäcka avbrutna sändningar.)







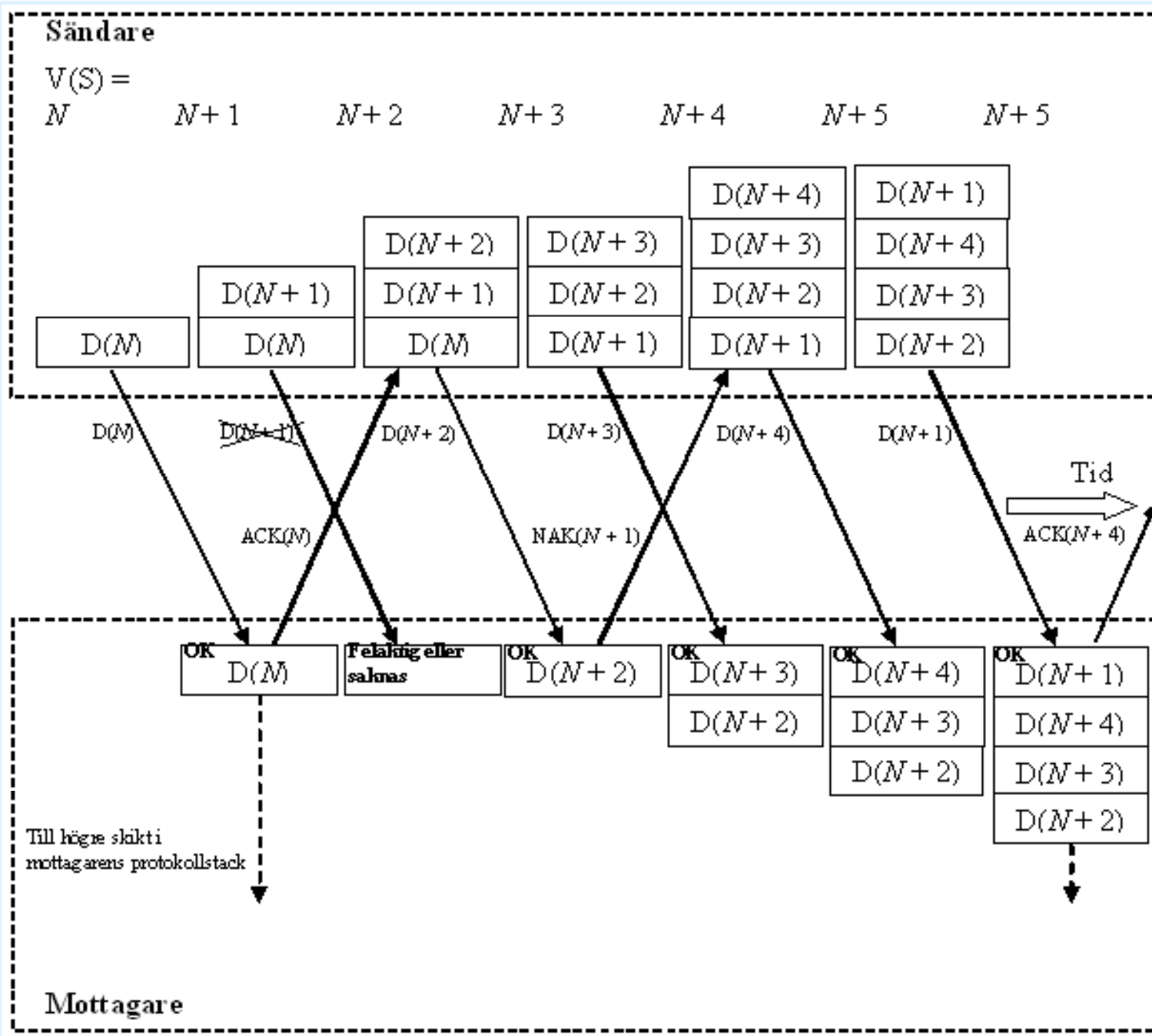
## Selective repeat with explicit request (1)

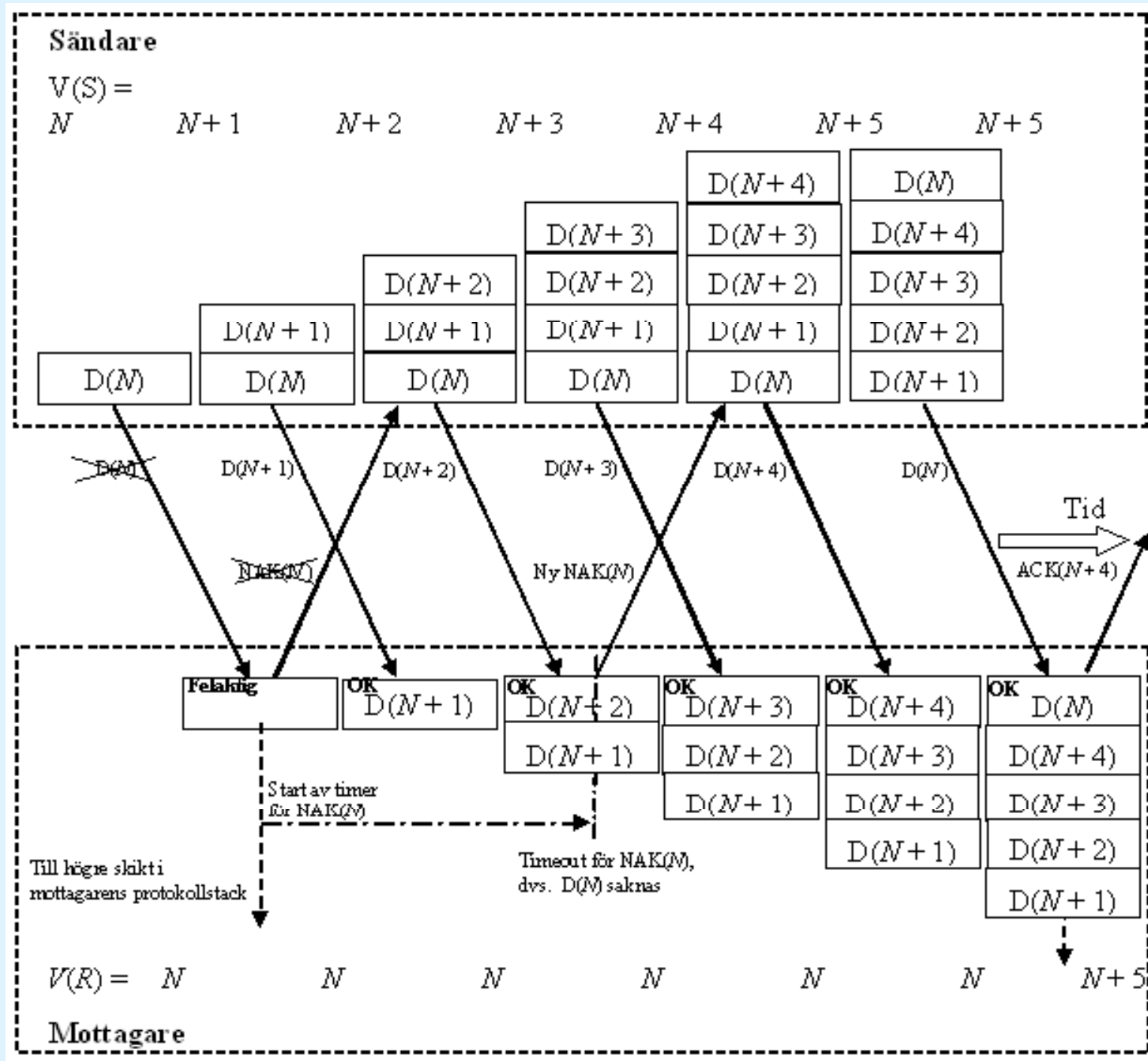
- ❑ Enbart felaktiga meddelanden (ramar, paket, segment etc.) sänds om vid begäran.
- ❑ Explicit request betyder att NAK används.

## Selective repeat with explicit request (2)

Speciella egenskaper:

- ❑  $ACK(N)$  kvitterar samtliga meddelanden till och med meddelande  $N$ .
- ❑ Mottagaren har timer för att förhindra uteblivet meddelande vid ev. fel på NAK.





# Go-back-N

(1)

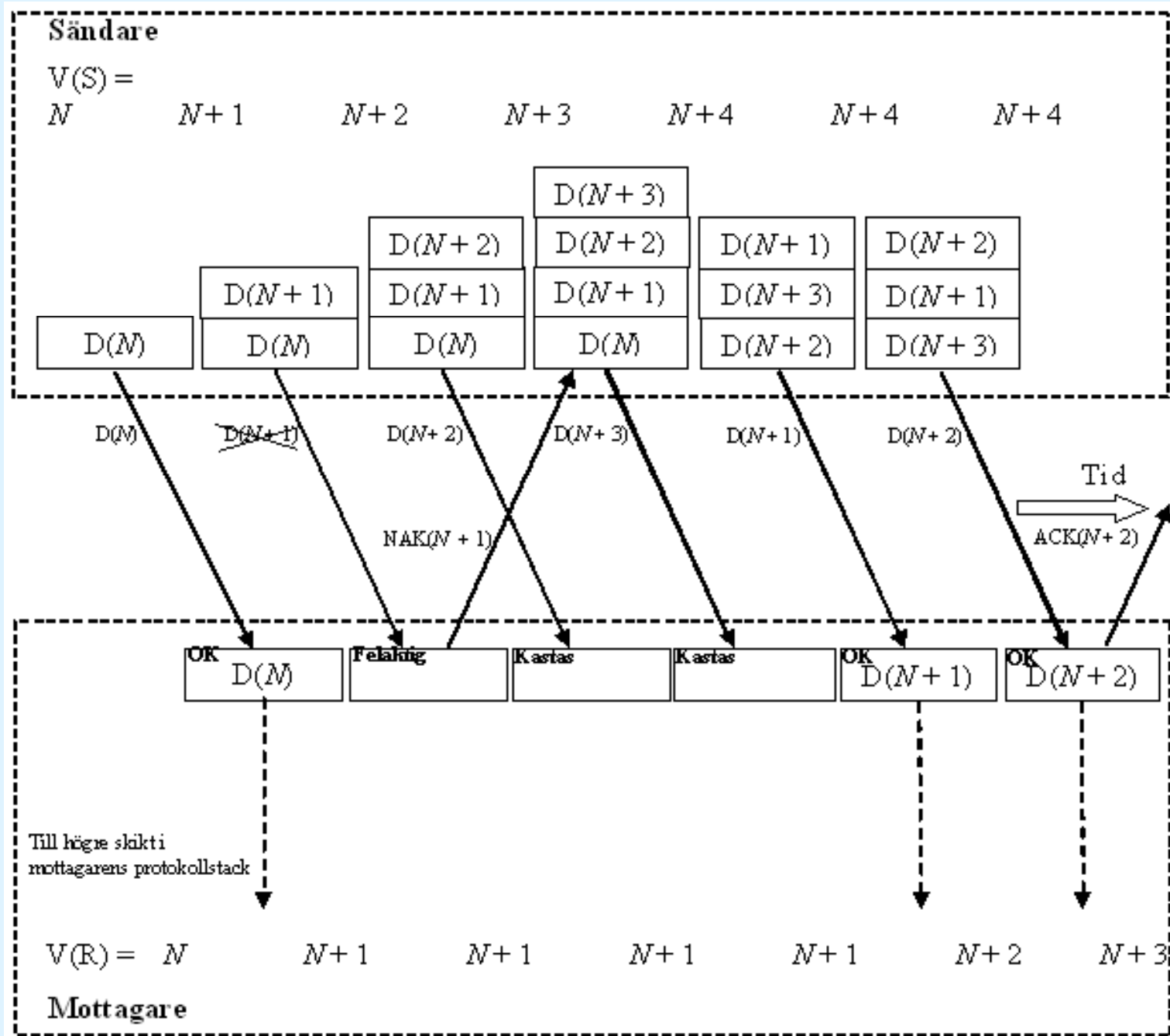
- ❑ Repeterar samtliga meddelanden från och med ett felaktigt.
- ❑ Fördel: mottagningsbufferten behöver plats för ett enda meddelande.

# Go-back-N

(2)

- ❑ Både ACK och NAK används.
- ❑  $\text{NAK}(N)$  medför omsändning från och med meddelande  $N$ .
- ❑  $\text{ACK}(N)$  och  $\text{NAK}(N + 1)$  ger implicit kvittering för alla meddelanden till och med meddelande  $N$ .
- ❑ Varje meddelande som sänds har en timer hos sändaren så att utebliven ACK/NAK ger omsändning.





# Piggyback acknowledgement

- ❑ Kvittenser inbäddade i nyttomeddelanden
- ❑ Onödig trafik reduceras

<u>V(S)</u>	<u>V(R)</u>
3	1

<u>V(S)</u>	<u>V(R)</u>
0	3

$D(N(S) = 3, N(R) = 1)$

1	4
---	---

$D(N(S) = 1, N(R) = 4)$

4	2
---	---

$D(N(S) = 4, N(R) = 2)$

2	5
---	---

Om  $N(R) = 1$ , skulle detta vara en begäran om repetition av ram nr. 1.

# Hammingavstånd

□  $HD = \min(\text{alla avstånd i en mängd av kodord})$

□ Om  $HD = u + 1$ ,

kan  $u$  st. bitfel upptäckas i ett felaktigt kodord.

□ Om  $HD = 2r + 1$ ,

kan  $r$  st. bitfel korrigeras i ett felaktigt kodord.

## Exempel på blockparitet (blocksumma)

00000010

= STX

10101000

1:a tecknet

00100000

2:a tecknet

10101101

3:e tecknet

11100011

4:e tecknet

$\oplus$  10000011

= ETX

BCC = 11000111

# Cyclic Redundancy Check (CRC)

- Används i många typer av LAN

Exempelvis i Ethernet.

- Felsökning

- Felrättning

CRC integrerad med felrättande kod

# Generatorpolynom

□ Bestämmer checksummans tillförlitlighet

□ CRC-12:  $G = x^{12} + x^{11} + x^3 + x^2 + 1$

□ CRC-16:  $G = x^{16} + x^{15} + x^2 + 1$

□ CRC-CCITT:  $G = x^{16} + x^{12} + x^5 + 1$

□ CRC-32:  $G = x^{32} + x^{26} + x^{23} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$

# CRC-beräkningar

$$\frac{D \cdot 2^n}{G} = Q \oplus \frac{R}{G} \qquad \frac{13}{4} = 3 + \frac{1}{4}$$

Sändaren skickar

$$T = D \cdot 2^n \oplus R$$

Mottagaren testar

$$\frac{T'}{G} = Q \oplus \frac{E}{G}$$

Om  $E = 0$ , felfri övering.



# Feltyper med CRC

Med ett bra val av generatorpolynom kan mottagaren upptäcka följande:

- ❑ alla en-bits fel
- ❑ alla två-bitars fel
- ❑ alla udda antal bitar fel
- ❑ alla felskurar  $< n + 1$  bitar
- ❑ de flesta felskurar  $\geq n + 1$  bitar

# Reed-Solomon (RS) (1)

- Används för datakommunikation och lagringsmedia

Det gör CRC också!

- Felsökning

- Felrättning

# Reed-Solomon (RS) (2)

□  $RS(n, k)$

Antal data ( $s$ -bitars symboler):  $k$

Antal data inklusive checksummorna:  $n$

Antal checksummor ( $s$ -bitars symboler):  
 $n - k$

# Reed-Solomon (RS) (3)

- Maximalt antal data som kan rättas ( $t$ )

$$2t = n - k$$

dvs.

$$t = (n - k)/2$$

# Flödesreglering

- Programstyrd
- Hårdvarustyrd
- Fönsterteknik

# Fönsterteknik

Omsändningsmetod	Sändningsfönster	Mottagningsfönster
<u>Idle repeat request</u>	1	1
<u>Selective repeat</u>	$k/2$	$k/2$
<u>Go-back-N</u>	$k - 1$	1

Antal sekvensnummer:  $k = 2^n$

# TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

## □ point-to-point:

- one sender, one receiver

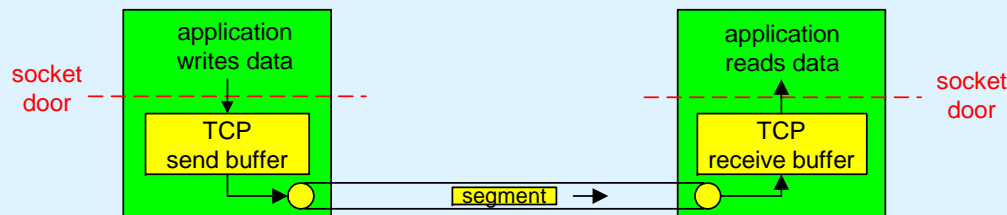
## □ reliable, in-order *byte stream*:

- no "message boundaries"

## □ pipelined:

- TCP congestion and flow control set window size

## □ *send & receive buffers*



## □ full duplex data:

- bi-directional data flow in same connection
- MSS: maximum segment size

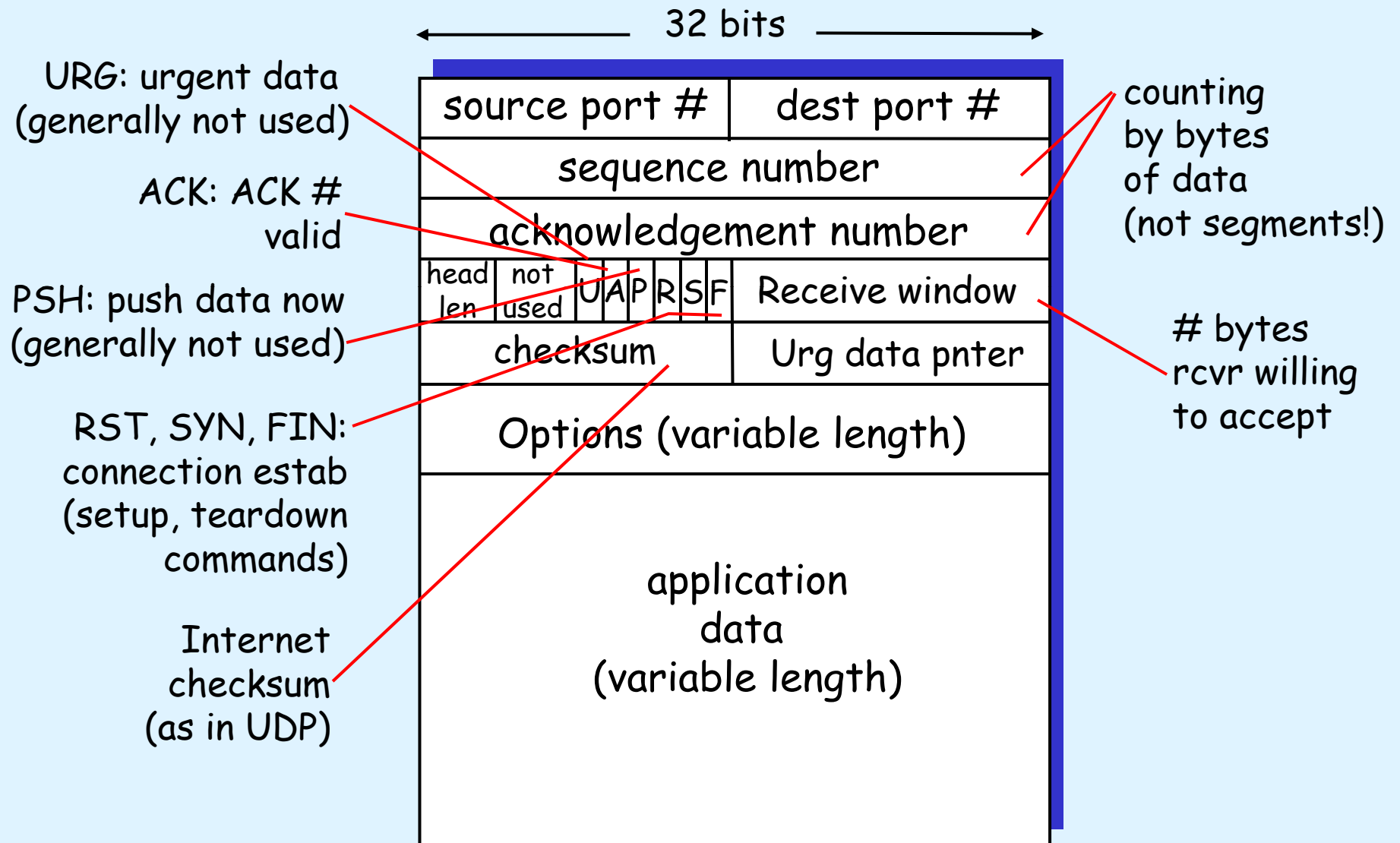
## □ connection-oriented:

- handshaking (exchange of control msgs) init's sender, receiver state before data exchange

## □ flow controlled:

- sender will not overwhelm receiver

# TCP segment structure





# TCP Connection Management

Recall: TCP sender, receiver establish "connection" before exchanging data segments

- initialize TCP variables:
  - seq. #s
  - buffers, flow control info (e.g. RcvWindow)

## Three way handshake:

Step 1: client host sends TCP SYN segment to server

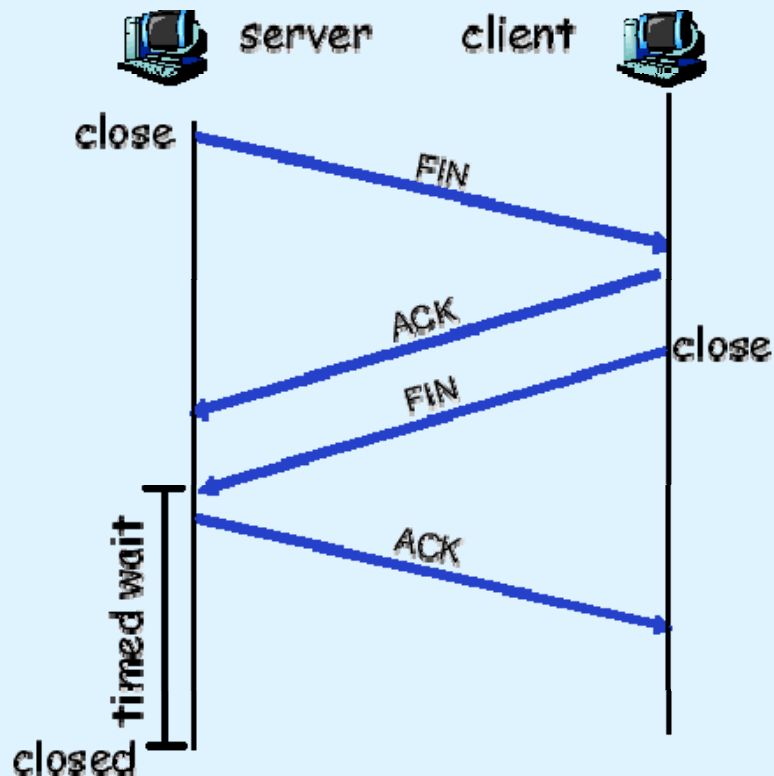
- specifies initial seq #
- no data

Step 2: server host receives SYN, replies with SYNACK segment

- server allocates buffers
- specifies server initial seq. #

Step 3: client receives SYNACK, replies with ACK segment, which may contain data

# TCP Connection Management (cont.)



Step 1: server end system sends TCP FIN control segment to client

Step 2: client receives FIN, replies with ACK. Closes connection, sends FIN.

Step 3: server receives FIN, replies with ACK.

- Enters "timed wait" - will respond with ACK to received FINs

Step 4: client, receives ACK. Connection closed.

# TCP Round Trip Time and Timeout

Q: how to set TCP timeout value?

- ❑ longer than RTT
  - but RTT varies
- ❑ too short: premature timeout
  - unnecessary retransmissions
- ❑ too long: slow reaction to segment loss

Q: how to estimate RTT?

- ❑ **SampleRTT**: measured time from segment transmission until ACK receipt
- ❑ **SampleRTT** will vary, want estimated RTT "smoother"
  - average several recent measurements, not just current **SampleRTT**

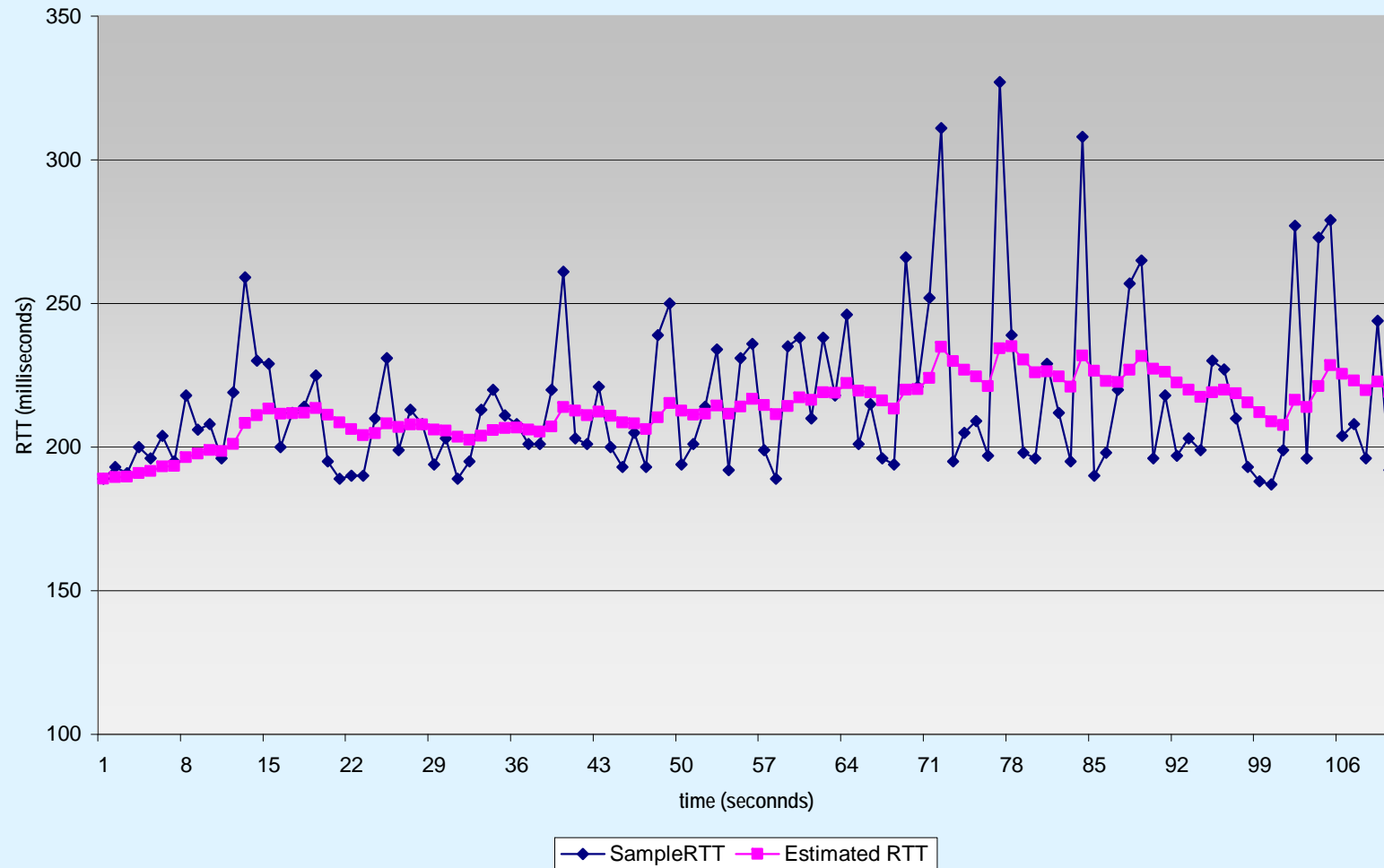
# TCP Round Trip Time and Timeout

$$\text{EstimatedRTT}_{\text{new}} = (1 - \alpha) \cdot \text{EstimatedRTT}_{\text{old}} + \alpha \cdot \text{SampleRTT}$$

- ❑ Exponential weighted moving average
- ❑ influence of past sample decreases exponentially fast
- ❑ typical value:  $\alpha = 0.125$

# Example RTT estimation:

RTT: gaia.cs.umass.edu to fantasia.eurecom.fr



# TCP Round Trip Time and Timeout

## Setting the timeout

- ❑ EstimatedRTT plus “safety margin”
  - large variation in EstimatedRTT -> larger safety margin
- ❑ first estimate of how much SampleRTT deviates from EstimatedRTT:

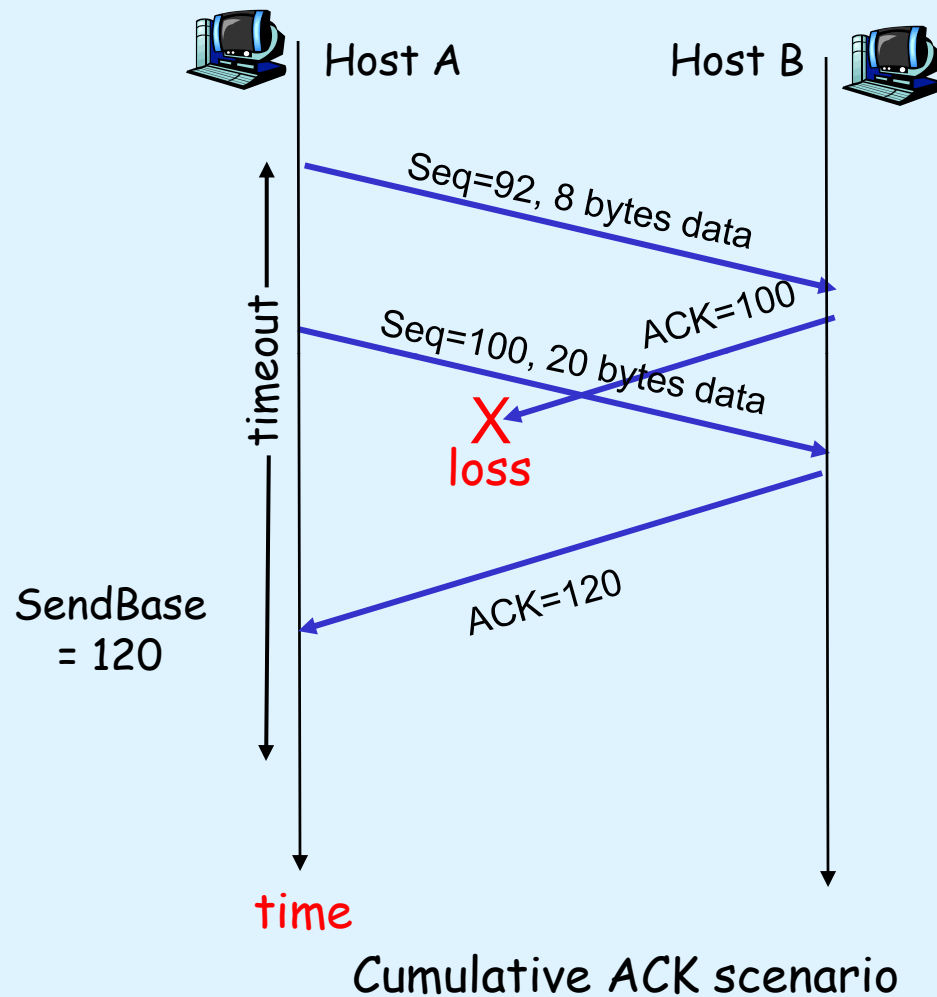
$$\text{DevRTT}_{\text{new}} = (1 - \beta) \cdot \text{DevRTT}_{\text{old}} + \beta \cdot |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically,  $\beta = 0.25$ )

Then set timeout interval:

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 \cdot \text{DevRTT}$$

# TCP retransmission scenarios (more)

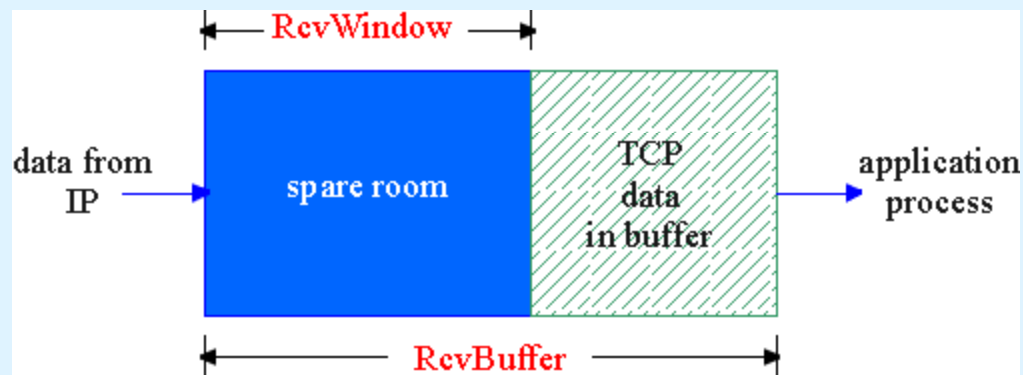


# Fast Retransmit

- ❑ Time-out period often relatively long:
  - long delay before resending lost packet
- ❑ Detect lost segments via duplicate ACKs.
  - Sender often sends many segments back-to-back
  - If segment is lost, there will likely be many duplicate ACKs.
- ❑ If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
  - fast retransmit: resend segment before timer expires



# TCP Flow control: how it works



- spare room in buffer

= RcvWindow

=  $\text{RcvBuffer} - [\text{LastByteRcvd} - \text{LastByteRead}]$

- Rcvr advertises spare room by including value of RcvWindow in segments
- Sender limits unACKed data to RcvWindow
  - guarantees receive buffer doesn't overflow

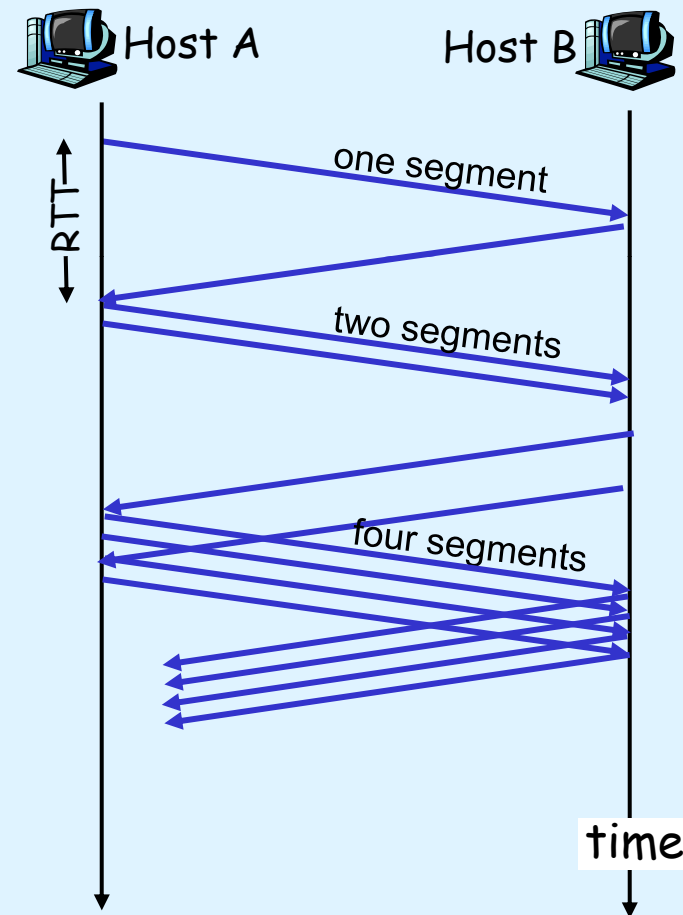
# Principles of Congestion Control

## Congestion:

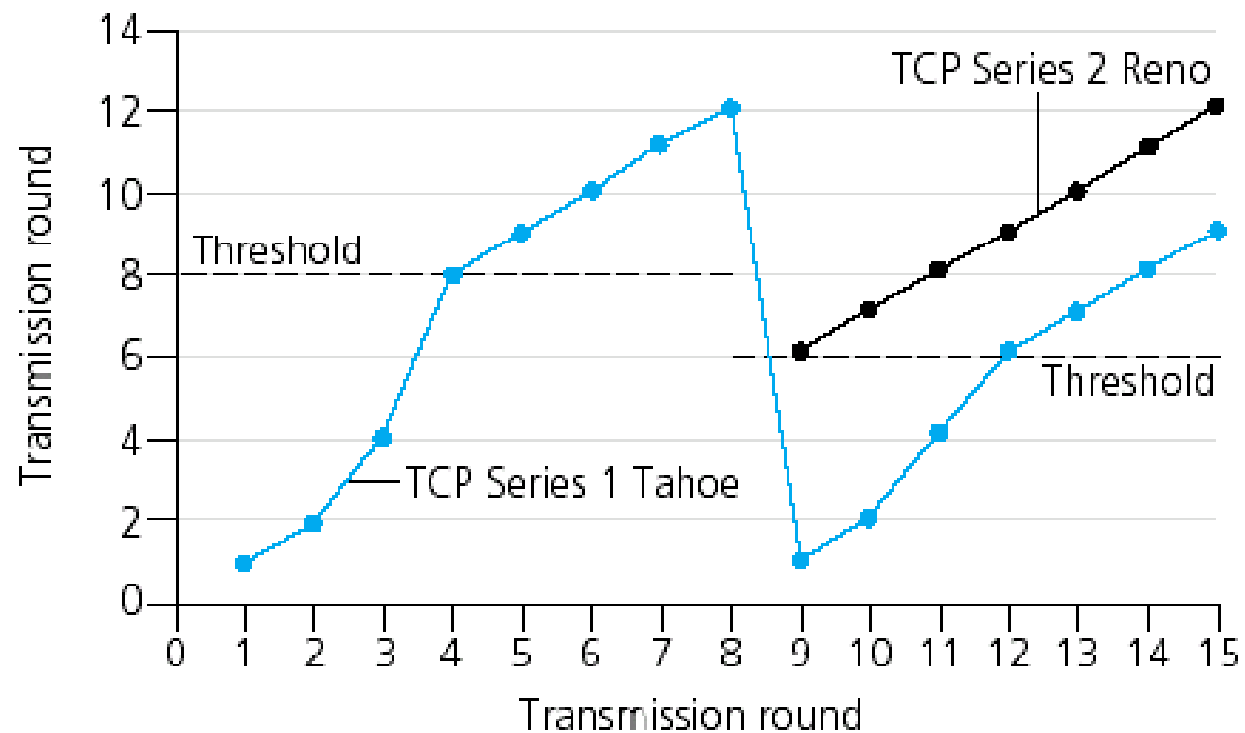
- ❑ informally: "too many sources sending too much data too fast for *network* to handle"
- ❑ different from flow control!
- ❑ manifestations:
  - lost packets (buffer overflow at routers)
  - long delays (queueing in router buffers)
- ❑ a top-10 problem!

# TCP Slow Start (more)

- When connection begins, increase rate exponentially until first loss event:
  - double CongWin every RTT
  - done by incrementing CongWin for every ACK received
- Summary: initial rate is slow but ramps up exponentially fast



# Refinement



## Summary: TCP Congestion Control

- ❑ When CongWin is below Threshold, sender in **slow-start** phase, window grows exponentially.
- ❑ When CongWin is above Threshold, sender is in **congestion-avoidance** phase, window grows linearly.
- ❑ When a **triple duplicate ACK** occurs, Threshold set to  $\text{CongWin}/2$  and CongWin set to Threshold.
- ❑ When **timeout** occurs, Threshold set to  $\text{CongWin}/2$  and CongWin is set to 1 MSS.