



Communication Networks II

Transport Layer - Protocols

www.kom.tu-darmstadt.de
www.httc.de

Prof. Dr.-Ing. **Ralf Steinmetz**

TU Darmstadt - Technische Universität Darmstadt,

Dept. of Electrical Engineering and Information Technology, Dept. of Computer Science

KOM - Multimedia Communications Lab

Merckstr. 25, D-64283 Darmstadt, Germany, Ralf.Steinmetz@KOM.tu-darmstadt.de

Tel.+49 6151 166151, Fax. +49 6151 166152

httc - Hessian Telemedia Technology Competence-Center e.V

Merckstr. 25, D-64283 Darmstadt, Ralf.Steinmetz@httc.de



Scope

KN III (Mobile Networking), Distributed Multimedia Systems (MM I and MM II), Telecooperation II,III. ...; Embedded Systems								
L5	Applications	Terminal access	File access	E-mail	Web	Peer-to-Peer	Inst.-Msg.	IP-Tel.
	Application Layer (Anwendung)							SIP & H.323
L4	Transport Layer (Transport)	Internet: UDP, TCP, SCTP			Netw. Transitions	Security	Addressing	Transport QoS - RTP
L3	Network Layer (Vermittlung)	Internet: IP						Network QoS
L2	Data Link Layer (Sicherung)	LAN, MAN High-Speed LAN						
L1	Physical Layer (Bitübertragung)	Queueing Theory & Network Calculus						
Introduction								
Legend:		KN I			KN II			



Overview

- 1. Transport Protocols and some History**
 - 1.1 ISO-OSI Transport Protocols**

- 2. Internet Transport Layer (in General & Addressing)**
 - 2.1 Port - Addressing Concept**
 - 2.2 Port - Link to Application**

- 3. UDP - User Datagram Protocol**

- 4. TCP - Transmission Control Protocol**
 - 4.1 TCP in Use & Application Areas**
 - 4.2 TCP Characteristics**
 - 4.3 Connection - Addressing**

- 5. TCP - Protocol, PDU Format, Segments**
 - 5.1 Segments, Fragmentation (and Reassembly)**
 - 5.2 Protocol with Flow Control**



Overview

6. TCP: Connections & Management

6.1 Connection Establishment

6.2 Connection Release

6.3 Connection Management Modelling

7. TCP - Foundations

7.1 Flow Control

7.2 Timer Management

7.3 Congestion Control

7.4 TCP - Further Comments

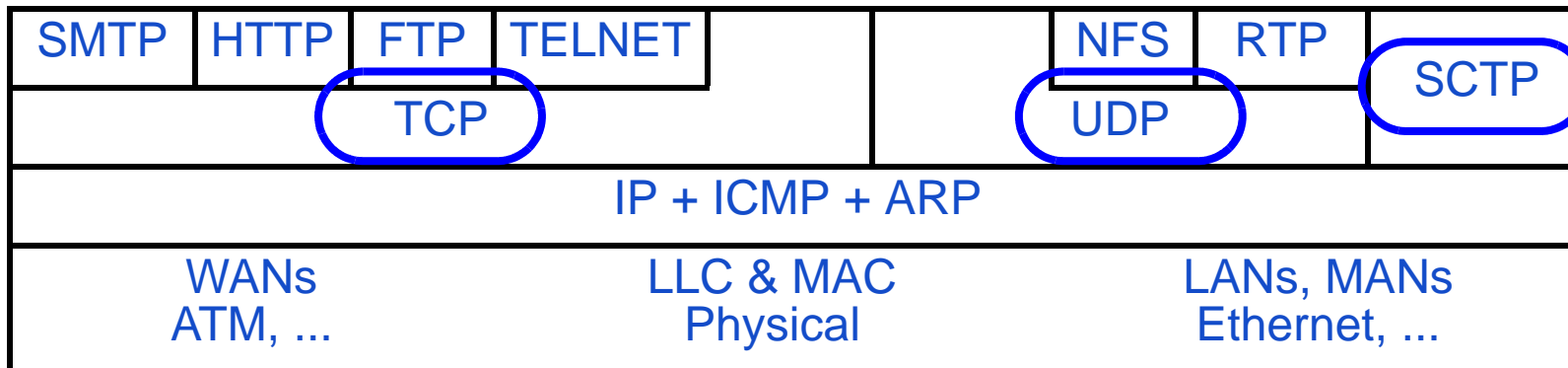
8. Stream Control Transmission Protocol (SCTP)

9. Further Development of Transport Protocols



1. Transport Protocols and some History

Internet



- **UDP** = **User Datagram Protocol**
- **TCP** = **Transmission Control Protocol**
- **SCTP** = **Stream Control Transmission Protocol**

ISO-OSI

- **practically irrelevant**
- **but show overall design space**

Other



1.1 ISO-OSI Transport Protocols

ISO (International Organization for Standardization)

<http://www.iso.ch/>

- is the world's largest developer of standards. Although ISO's principal activity is the development of technical standards, ISO standards also have important economic and social repercussions. ISO standards make a positive difference, not just to engineers and manufacturers for whom they solve basic problems in production and distribution, but to society as a whole.
- is a network of the national standards institutes of 147 countries, on the basis of one member per country

Open systems interconnection OSI

- communications standards

taking

- **a practical view:**
 - irrelevant today
 - never got sufficient support with respect to implementations etc.
- **general view:**
 - provides for an overall classification of schemes
 - shows overall design space

⇒ **we take a short look at them in a general manner**



ISO-OSI Transport Protocols & Network Services

L4 protocol depends on the quality of the L3 service (NS network service)

- **Network type A:**
typically **CONS** (connection oriented network service) on LANs
 - network is *RELIABLE*
 - network recognizes data loss as an error
 - errors are displayed to the user(N-RESET) i.e. acceptable rate of the errors
 - minor (for the user acceptable) error rate
 - network *NEVER* duplicates or manipulates packets
 - order of sent packets is *ALWAYS* maintained
- **Network type B:**
typically **CONS** on (old) WANs
 - like type A, except
 - *REMAINING ERROR RATE* (for data loss) *IS NOT ACCEPTABLE*
- **Network type C:**
typically **CLNS** (connectionless network serv.) on WANs
 - network is *UNRELIABLE*
 - errors due to losses, duplication and manipulation of packets, as well as faulty packet sequence errors possible
 - errors might remain undetected
 - transport protocol has to / should compensate for this



ISO-OSI Transport Protocols: Classes

Transport protocol classes

- 5 classes: ISO OSI TP0..TP4

Protocol Class	Network Type	Network Properties	Name
TP 0	A	Acceptable error rate Acceptable rate of displayed errors	Simple class
TP 2			Multiplexing class
TP 1	B	INacceptable error rate Acceptable rate of displayed errors	Basic error recovery
TP 3			Error recovery and multiplexing class
TP 4	C	INacceptable error rate INacceptable rate of displayed errors	Error recovery and multiplexing class



Class TP 0: simple class (A)

- mechanisms for connect and disconnect
- segmentation / reassembly
- no error, sequence or flow control
- no expedited data

Class TP 1: basic error recovery (B)

- class 0 including additional error recovery
- error recovery masks N-RESETs
 - TPDU numbering
 - TPDU storage until ACK
 - after N-RESET: resynchronization
- expedited data optional
 - important data for example have a higher priority
 - i.e. preferred processing before current data is processed



Class TP 2: multiplexing class (A)

- **class 0 including additional multiplexing capability**
- **MULTIPLEXING: several L4 connections on one L3 connection**
- **flow control optional**
- **expedited data optional**

Class TP 3: including multiplexing and error recovery (B)

- **class 1 and 2 functions combined**
- **i.e. error recovery, expedited data, multiplexing**

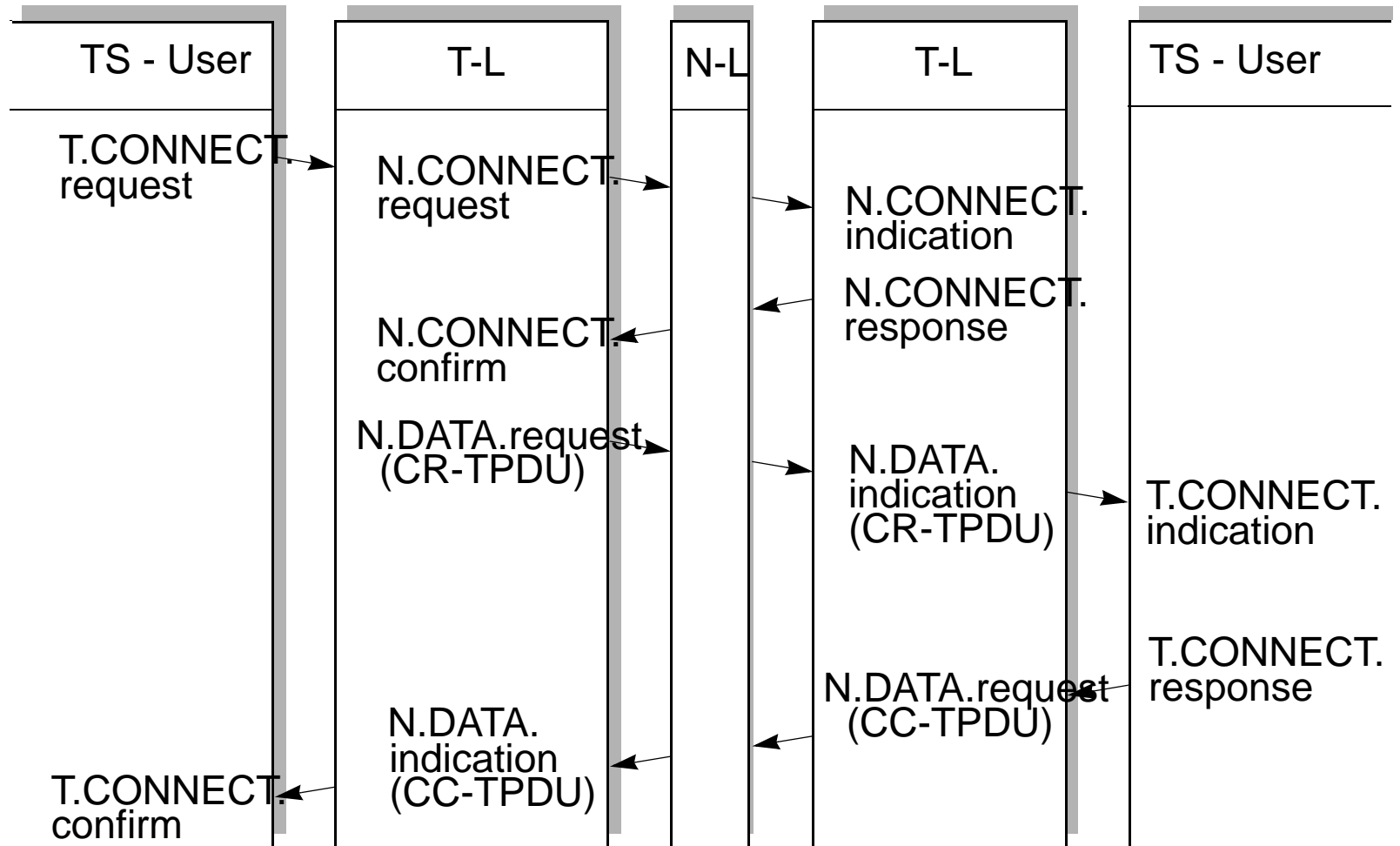
Class TP 4: error monitoring and recovery(C)

- **detects and recovers**
 - TPDU losses and TPDU duplication
 - sequence errors
- **flow control**
- **multiplexing**
- **splitting (one T connection uses several N connections)**
- **expedited data**



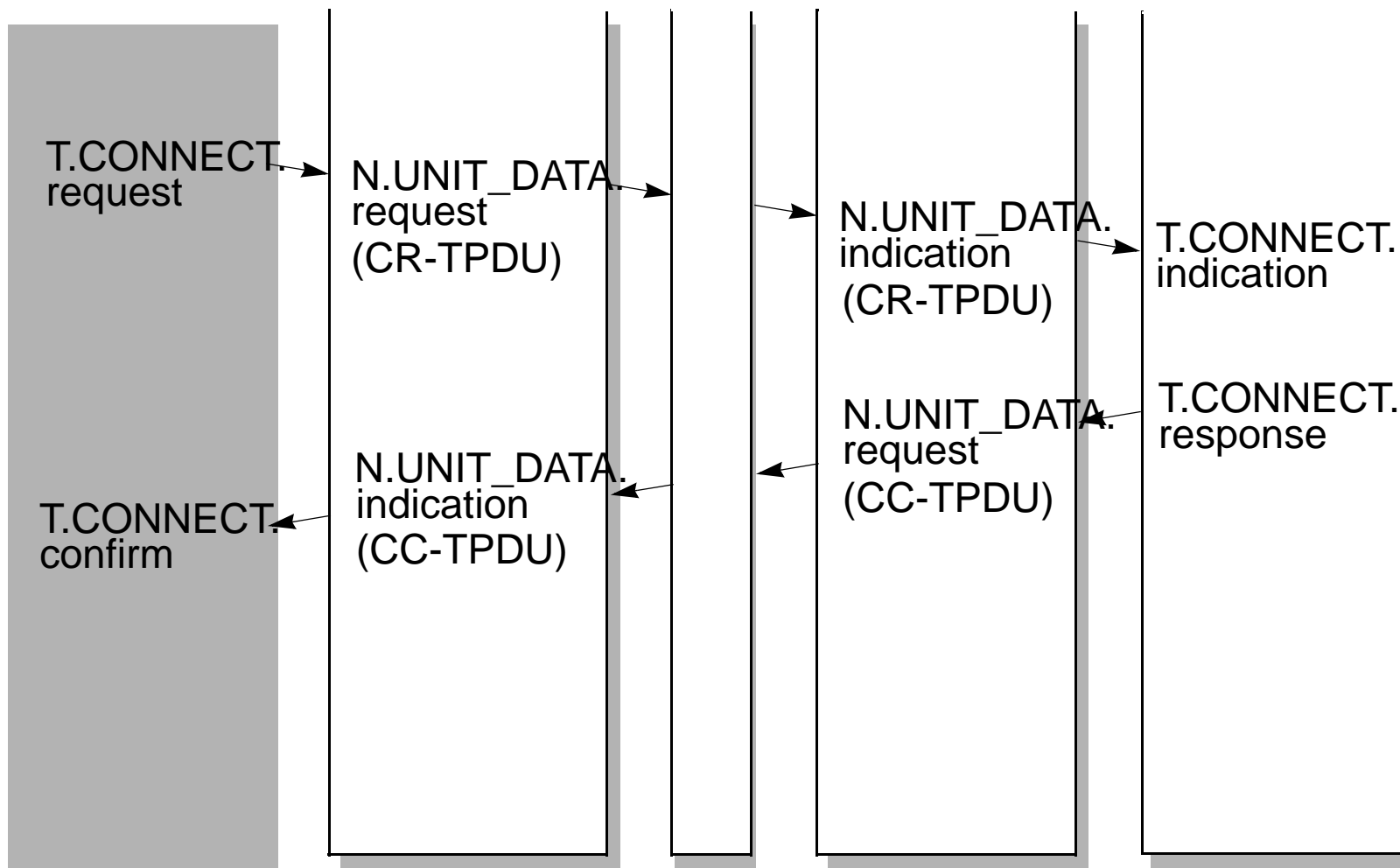
ISO OSI Transport Protocols: Model

Based on connection oriented network service





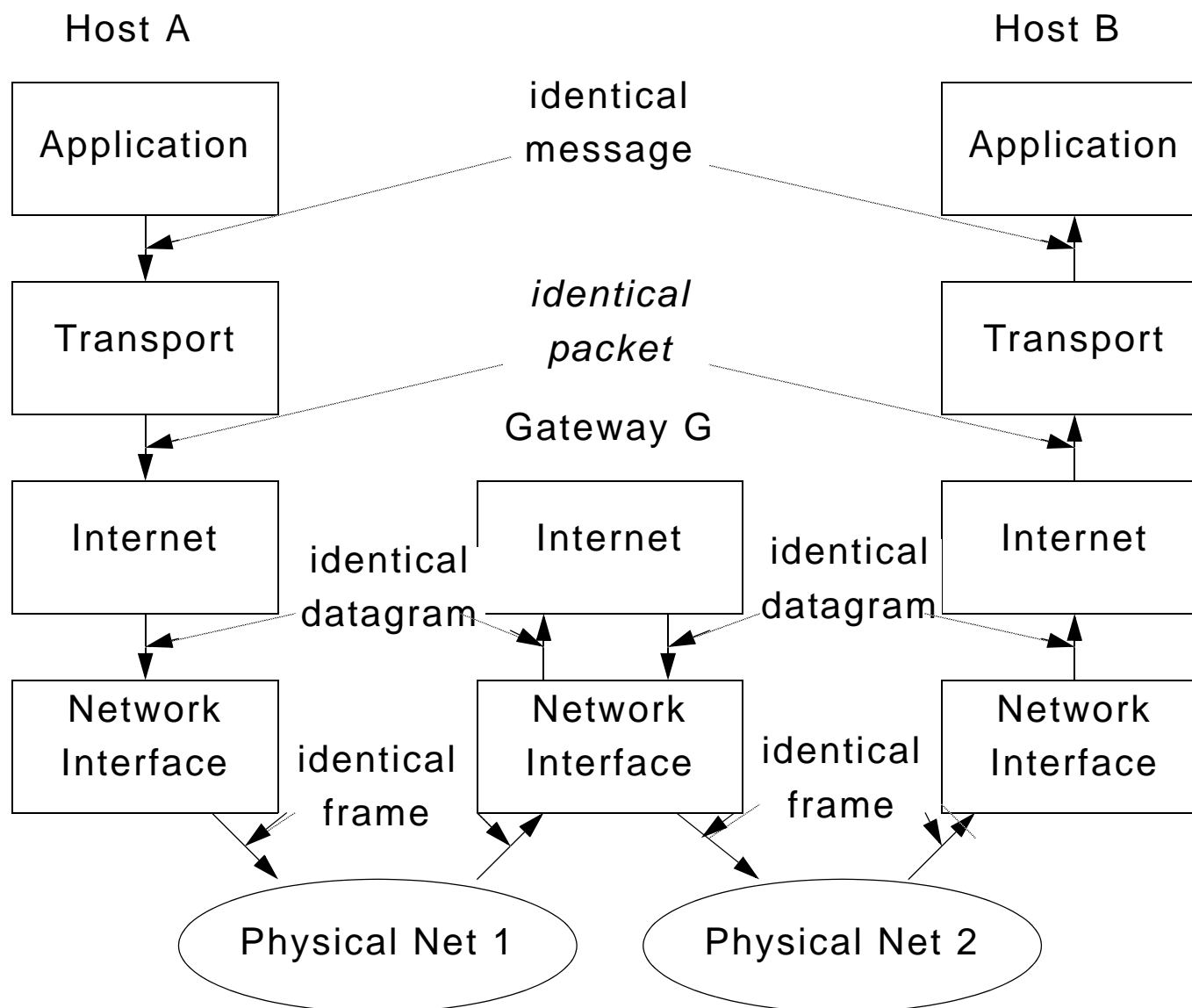
Based on connectionless network service





2. Internet Transport Layer (in General & Addressing)

Internet architecture





Well-Known Internet Protocols

SMTTP	HTTP	FTP	TELNET			NFS	RTP	SCTP
TCP					UDP			
IP + ICMP + ARP								
WANs ATM, ...			LLC & MAC Physical			LANs, MANs Ethernet, ...		

- ARP = Address Resolution Protocol
- FTP = File Transfer Protocol
- HTTP = Hypertext Transfer Protocol
- IP = Internet Protocol
- ICMP = Internet Control Message Protocol
- LLC = Logical Link Control
- MAC = Media Access Control
- NFS = Network File System
- SMTTP = Simple Mail Transfer Protocol
- TELNET = Remote Login Protocol
- SCTP = STREAM CONTROL TRANSMISSION PROTOCOL
- TCP = TRANSMISSION CONTROL PROTOCOL
- UDP = USER DATAGRAM PROTOCOL

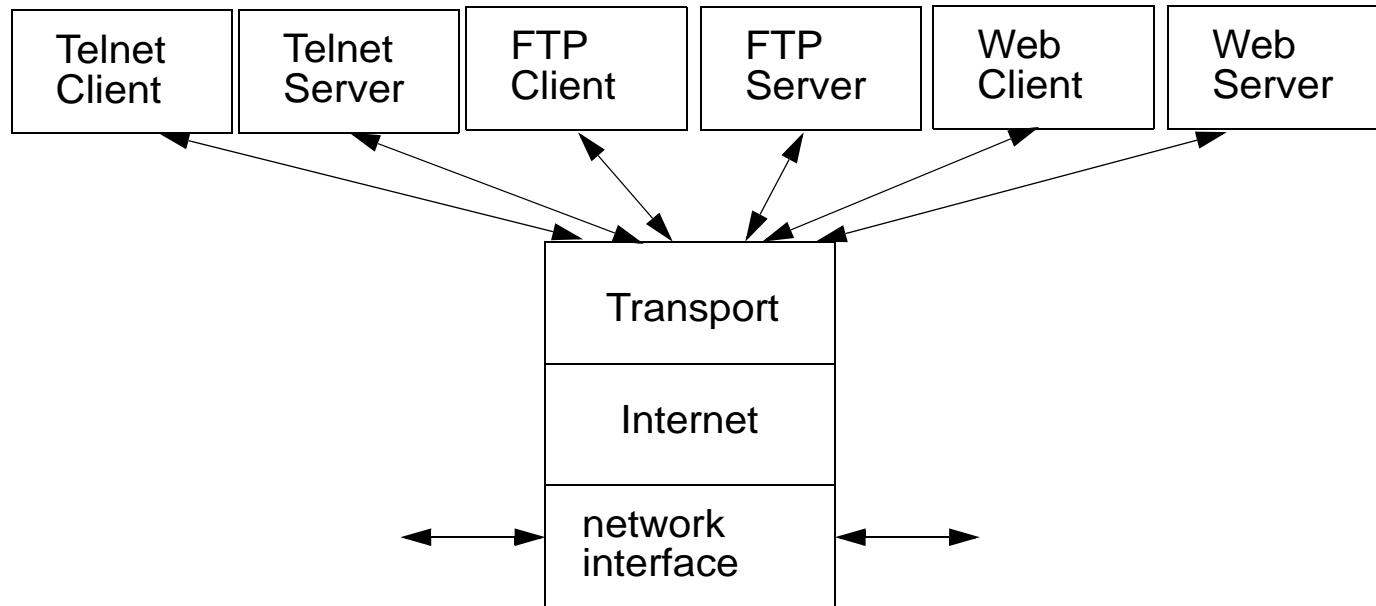
And further protocols including

IMAP= Interactive Mail Access Protocol, NNTP=Network News Transfer Protocol
POP3= Post Office Protocol



Transport Layer

Application



- **communication between applications required**
- **application communicate**
 - locally by interprocess communication
 - between system via **TRANSPORT SERVICES**

Transport layer

- **interprocess communication via communication networks**

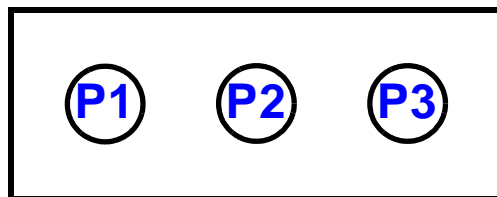
Internet Protocol IP

- **enables endsystem-to-endsystem communication**

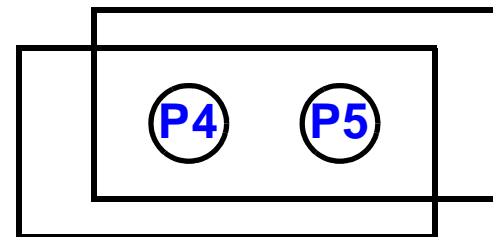


2.1 Port - Addressing Concept

Service A



Service B



Service C

3 types of identifiers: names, addresses and routes

[Shoch 78]

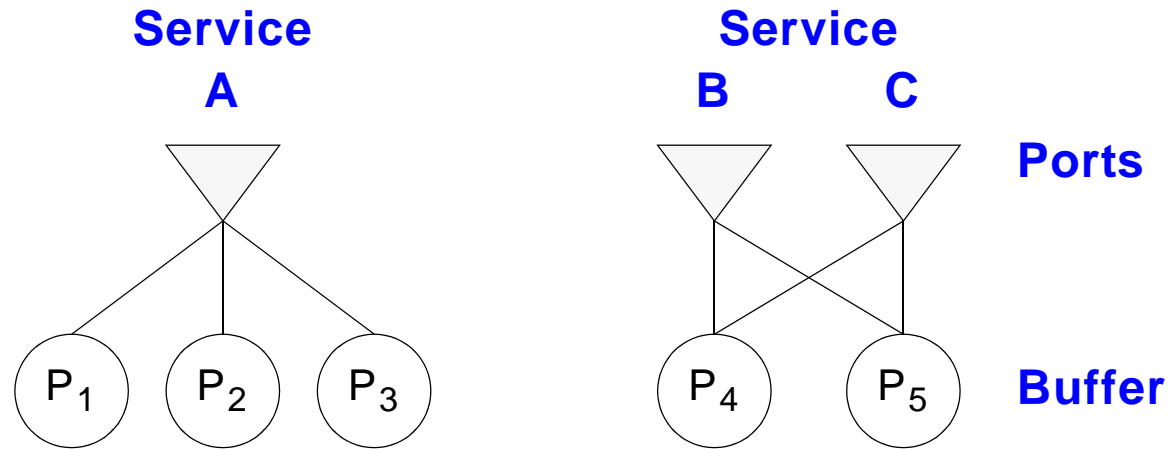
“The **NAME** of a resource indicates WHAT we seek, an **ADDRESS** indicates WHERE it is, and a **ROUTE** tells HOW TO GET THERE

- **address identifies**
 - type of service or application
- **addressing by process number is unsuitable**
 - processes are generated/terminated dynamically, i.e. the process number rarely known
 - relationship “service - process” not fix
 - 1 process can supply multiple services
 - various processes can provide same service

⇒ **Concept of an abstract communication endpoint: Port**



Communication Ports



Service

- related to exactly one single port

Port access

- asynchronous or
- synchronous

Port

- associated with buffer



Reserved Port Numbers

www.kom.tu-darmstadt.de
www.httc.de

Decimal	Keyword	UNIX Keyword	Description
0			Reserved
1	TCPMUX		TCP Multiplex
5	RJE		Remote Job Entry
7	ECHO	echo	Echo
9	DISCARD	discard	Discard
11	USERS	systat	Active Users
13	DAYTIME	daytime	Daytime
15		netstat	Network status program
17	QUOTE	qotd	Quote of the Day
19	CHARGEN	chragen	Character Generator
20	FTP-DATA	FTP-DATA	FILE TRANSFER PROTOCOL (DATA)
21	FTP	FTP	FILE TRANSFER PROTOCOL
23	TELNET	TELNET	TERMINAL CONNECTIONS
25	SMTP	SMTP	SIMPLE MAIL TRANSFER PROTOCOL
37	TIME	time	Time
42	NAMESERVER	name	Host Name Server

- **TCP and UDP have their own assignments**
 - this table shows some examples for TCP



Reserved Port Numbers

(2)

Decimal	Keyword	UNIX Keyword	Description
43	NICNAME	whois	Who is
53	DOMAIN	nameserver	Domain Name Server
77		rje	any private rje service
79	FINGER	finger	Finger
80	HTTP	HTTP	WORLD WIDE WEB
101	HOSTNAME	hostname	NIC Host Name Server
102	ISO-TSAP	iso-tsap	ISO TSAP
103	X400	x400	X.400 Mail Service
104	X400-SND	x400-snd	X.400 Mail Sending
110	POP3	POP3	REMOTE EMAIL ACCESS
111	SUN RPC	sunrpc	SUN Remote Procedure Call
113	AUTH	auth	Authentication Service
117	UUCP-PATH	uucp-path	UUCP Path Services
119	NNTP	nntp	USENET News Transfer Protocol
129	PWDGEN		Password Generator Protocol
139	NETBIOS-SSN		NETBIOS Session Protocol
160-1023	Reserved		



2.2 Port - Link to Application

Application

- **example**
 - decompression of video data
 - read process from database or file system
- **implementation of application**
 - process, thread
- **interface to communication systems**
 - buffers with predefined access mechanisms

Sender and receiver create

- **stream or**
- **socket**
 - several connections share a socket
 - address: IP address of the endsystem
 - address: 16-bit port number
 - 0..1024: predefined ports, “well known”
 - additional ones managed dynamically

Example:

192.169.100.17:80 socket with
IP address 192.169.100.17 and port no. 80



3. UDP - User Datagram Protocol

Specification:

- **RFC 768**

UDP is a simple transport protocol

- **unreliable**
- **connectionless**
- **message-oriented**

⇒ **UDP is mostly IP with short transport header**

- source and destination port
- ports allow for dispatching of messages to receiver process

Characteristics

- **no flow control**
 - application may transmit as fast as it can / want and the network permits
- **no error control or retransmission**
 - no guarantee about packet sequencing
 - packet delivery to receiver not ensured
 - possibility of duplicated packets
- **may be used with broadcast / multicasting and streaming**



UDP: Message Format

Sender port

- 16 bit sender identification
- optional
- response may be sent there

Receiver port

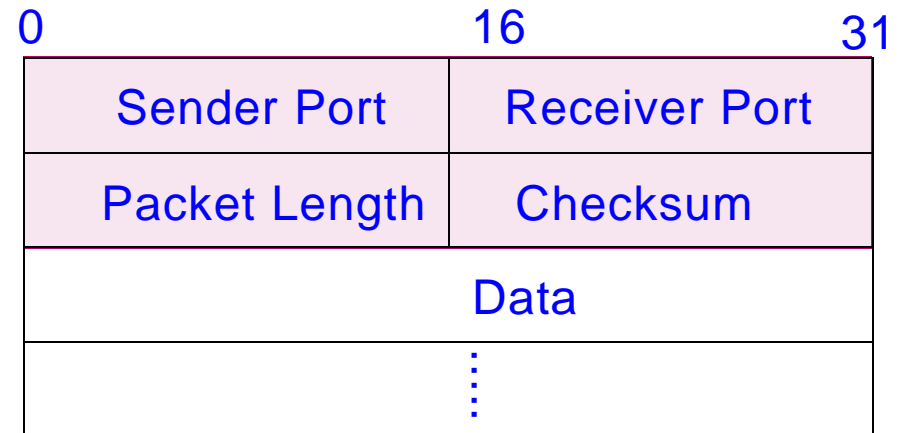
- receiver identification

Packet length

- in byte (including UDP header)
- minimum: 8 (byte)
 - i.e. header without data

Checksum

- of header and data for error detection
- use of checksum optional



 UDP Header



UDP: Message Format - Checksum

Purpose

- error detection (header and data)

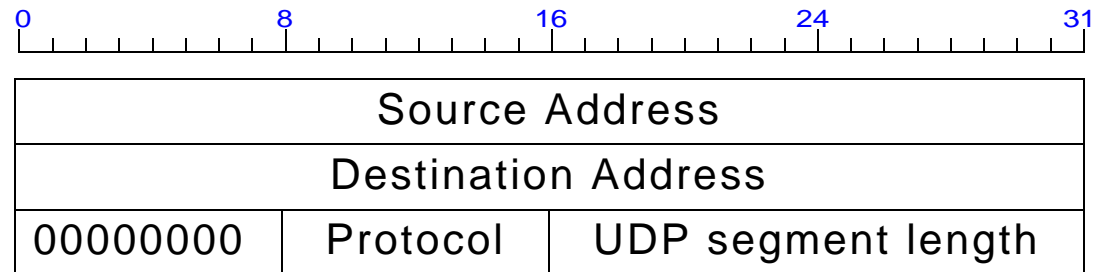
Same algorithm as IP

- one's complement of sum of 16-bit halfwords in one's complement arithmetic

UDP checksum includes

- UDP header (checksum field initially set to 0)
- data
- pseudoheader

- part of IP header
 - source IP address
 - destination IP address
 - protocol
 - length of (UDP) data
- allows to detect misdelivered UDP messages



Use of checksum optional

- i.e., if checksum contains only "0"s, it is not used
 - transmit 0xFFFF if calculated checksum is 0



UDP: Ranges of Application

Benefits of UDP

- **needs very few resources (in endsystems, in each data unit)**
- **no connection establishment (hence, no overhead, faster transmission)**
- **simple implementation possible (e.g. suitable for boot PROM)**
- **applications can precisely control**
 - packet flow
 - error handling / reliability
 - timing

Suitable for simple client-server interactions, i.e. typically

- **1 request packet from client to server**
- **1 response packet from server to client**

Used by e.g.

- **DNS: Domain Name Service**
- **SNMP: Simple Network Management Protocol**
- **system status**
- **bootstrap protocol**
- **TFTP: Trivial File Transfer Protocol (but not ftp!)**
- **NFS: Network File System**
- **RTP: Real-time Transport Protocol**



4. TCP - Transmission Control Protocol

Overview

TCP: is the transport protocol (the major Internet transport protocol)

Motivation: network with connectionless service

- **packets and messages may be**
 - duplicated, in wrong order, faulty
 - i.e., with such service only, each application would have to provide recovery
 - error detection and correction
 - **network or service can**
 - impose packet length
 - define additional requirements to optimize data transmission
 - i.e., application would have to be adapted
- ⇒ **TCP is the Internet transport protocol providing**
- reliable end-to-end byte stream over an unreliable internetwork

History

- **RFC 793: originally**
- **RFC 1122 and RFC 1323: errors corrected, enhancements implemented**



What is TCP?

TCP is

- **a communication protocol**
 - not a piece of software
 - (compare: programming languages, compilers)

TCP specifies

- **data and control information formats**
- **procedures for**
 - flow control
 - error detection and correction
 - connect and disconnect
- **as a primary abstraction**
 - a connection
 - not just the relationships of ports (as a queue, like UDP)

TCP does not specify

- **the interface to the application (sockets, streams)**



4.1 TCP in Use & Application Areas

Each machine supporting TCP has a TCP transport entity

- **library procedure**
- **user process**
- **part of kernel**

TCP transport entity manages

- **TCP streams**
- **interfaces to IP layer**

TCP transport entity on sending side

- **acceptes user data streams for local processes**
- **splits them into pieces ≤ 64 KB**
 - typically 1460 bytes
 - (to fit into single Ethernet frame with IP and TCP headers)
 - sends each piece as separate IP datagram

TCP transport entity on receiving side

- **gets TCP data from datagrams received at host**
- **reconstructs original byte streams**

TCP must ensure reliability

- **IP layer doesn't guarantee that datagrams will be delivered properly / in order (TCP must handle this, e.g. timeout and retransmit / reorder)**



Benefits of TCP

- **reliable data transmission**
- **efficient data transmission despite complexity**
 - (up to 8Mbps on 10Mbps ethernet)
- **can be used with LAN and WAN for**
 - low data rates (e.g. interactive terminal) and
 - high data rates (e.g. file transfer)

Disadvantages when compared with UDP

- **higher resource requirements**
 - buffering, status information, timer usage
- **connection set-up and disconnect necessary**
 - even with short data transmissions

Applications

- **file transfer (FTP)**
- **interactive terminal (Telnet)**
- **e-mail (SMTP)**
- **X-Windows**



4.2 TCP Characteristics

Data stream oriented

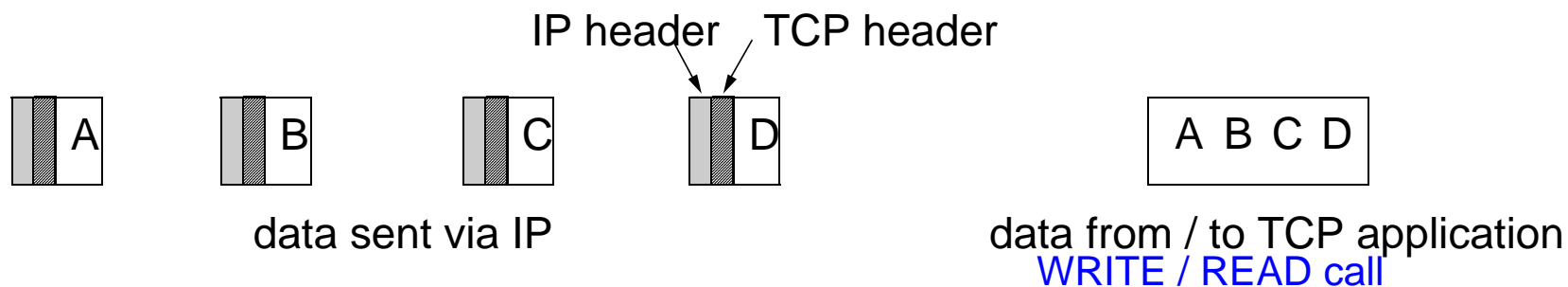
- TCP transfers serial byte stream
- maintains sequential order

Unstructured byte stream

- application often has to transmit more structured data
- TCP does not support such groupings into (higher) structures within byte stream

Buffered data transmission

- byte stream not message stream: message boundaries are not preserved
 - no way for receiver to detect the unit(s) in which data were written



- for transmission the sequential data stream is

- divided into segments
- delayed if necessary (to collect data)
 - for more efficient transmission (e.g. network utilization)



Virtual connection

- **connection established between communication parties before data transmission**

Two-way communications (fully duplex)

- **data may be transmitted simultaneously in both directions over a TCP connection**

Point-to-point

- **each connection has exactly two endpoints**

Reliable

- **fully ordered, fully reliable**
 - **sequence maintained**
 - **no data loss, no duplicates, no modified data**



TCP Characteristics: Some Protocol Elements & Features

(3)

Error detection

- through checksum

Piggybacking:

- control information and data can be transmitted within the same segment

Out-of-band data: expedited data

- important information sent to receiver
- i.e. should get to receiver's application before data that was sent earlier

PUSH operation

- data is not stored in a buffer
- sent immediately and immediately made available to application on receiver's end

example <CR> end of line at terminal emulation

Urgent flag

- send and transfer data to application immediately

example <Ctrl C>
arrival interrupts receiver's application



TCP Characteristics (missing)

(missing) Characteristics

- **no broadcast**
 - no possibility to address all applications
 - with connect, however, not necessarily sensible
- **no multicasting**
 - group addressing not possible
- **no QoS parameters**
 - not suited for different media characteristics
- **no real-time support**
 - no correct treatment/communications of audio or video possible
 - e.g. no forward error correction



4.3 Connection - Addressing

TCP service obtained via service endpoints on sender and receiver

- typically socket
- socket number consists of:
 - IP address of host and
 - 16-bit local number (port)

Port

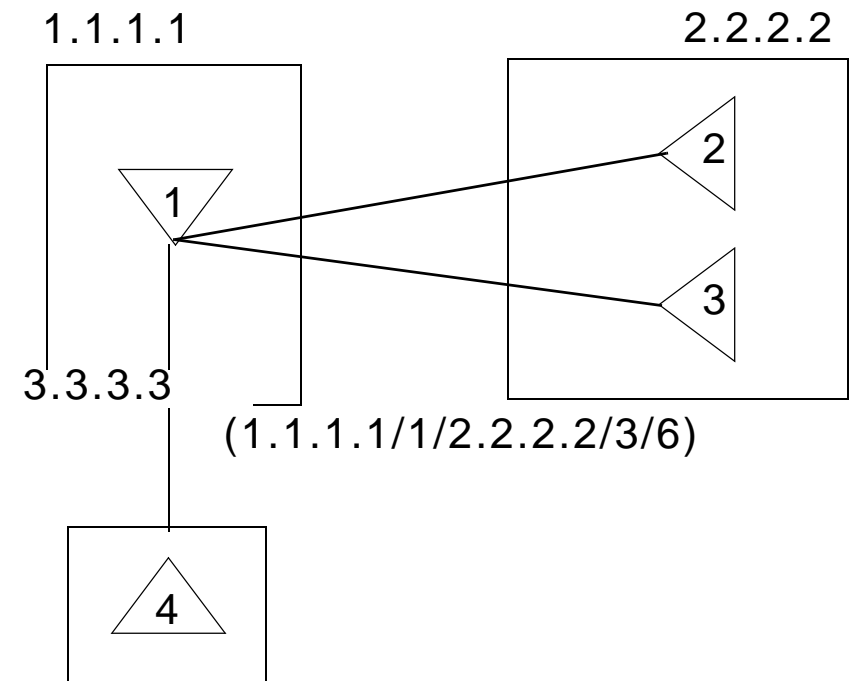
- TCP's name for TSAP

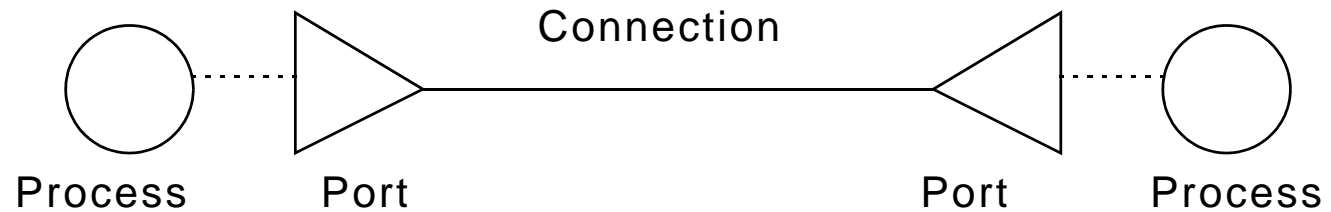
TCP connection is clearly defined by a quintuple consisting of

- IP address of *SENDER* and *RECEIVER*
- port address of *SENDER* and *RECEIVER*
- TCP protocol identifier

⇒ Applications can use the same local ports for several connections

IP addr.sender/port sender/ ..
(1.1.1.1/1/2.2.2.2/2/6)
../IP addr.rec/port rec/TCPid





Passive open:

- **process indicates that it would accept connect request**

Active open:

- **process requests a connection**

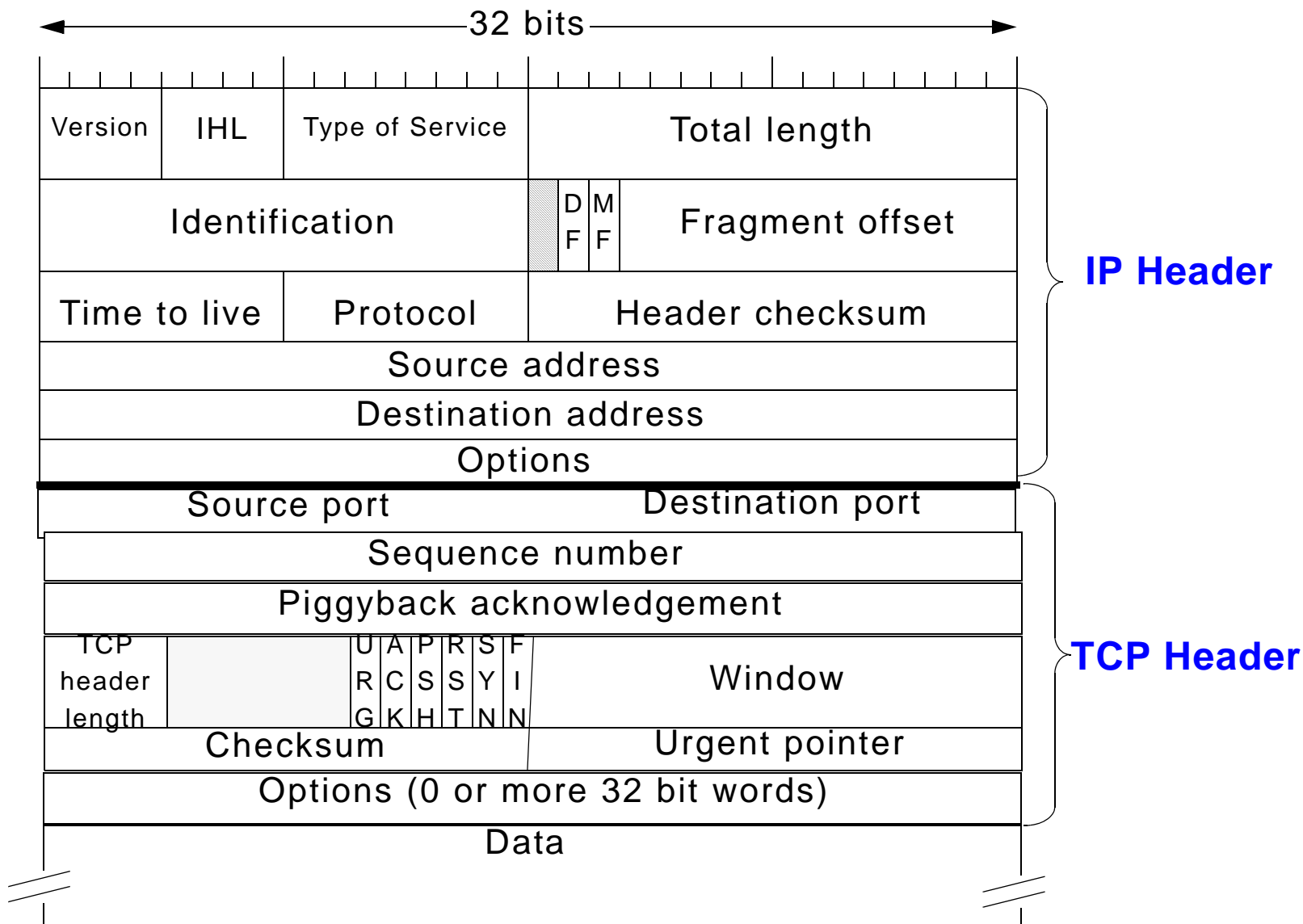
Addressing:

- **port number + protocol identification**
 - clearly identifies entity in the ES
- **IP address**
 - clearly identifies ES



5. TCP - Protocol, PDU Format, Segments

TCP/IP Header Format





Segments

TCP entities exchange data in form of *SEGMENTS*

TCP segment consists of

- **fixed 20 byte header (plus optional part)**
- **zero or more data bytes**

TCP software (entity) decides

- **segment size to be used**
 - data from several writes can be accumulated into one segment
 - or data from one write can be split into several segments
- **limits**
 - each segment (including TCP header) must fit into 65515 byte IP payload
 - segment must fit into maximum transfer unit (MTU) of visited networks
 - each network may have MTU, depending on L2 technology used
 - often 1500 byte (Ethernet payload size), typical upper bound on segm. size
 - further on (if really needed) IP will "fragment" packets if they are too large for visited networks



Identification of TCP-Packets i.e. Segments

A key feature:

- every byte on TCP connection has its own 32-bit sequence number

Separate sequence numbers used for

- data
- acknowledgements
- window mechanism

Remember:

- 32-bit sequence number space was big in early days of the Internet
- nowadays, it can be consumed very fast



Protocol with Flow Control

TCP uses: *SLIDING WINDOW* protocol

- **sender starts timer when it transmits segment**
- **receiving TCP entity sends back a segment**
 - with data if any, otherwise without data
 - acknowledgment number equal to next sequence number it **expects** to receive
- **if sender's timer goes off before acknowledgement arrives, segment is retransmitted**



5.1 Segments, Fragmentation (and Reassembly)

Segments

- **TCP *DATA STREAM* split into segments**
 - ***SEGMENTS* sent as IP packets**

Fragments

- **IP packets are split (if necessary) into *FRAGMENTS* in order to adapt them to underlying networks**

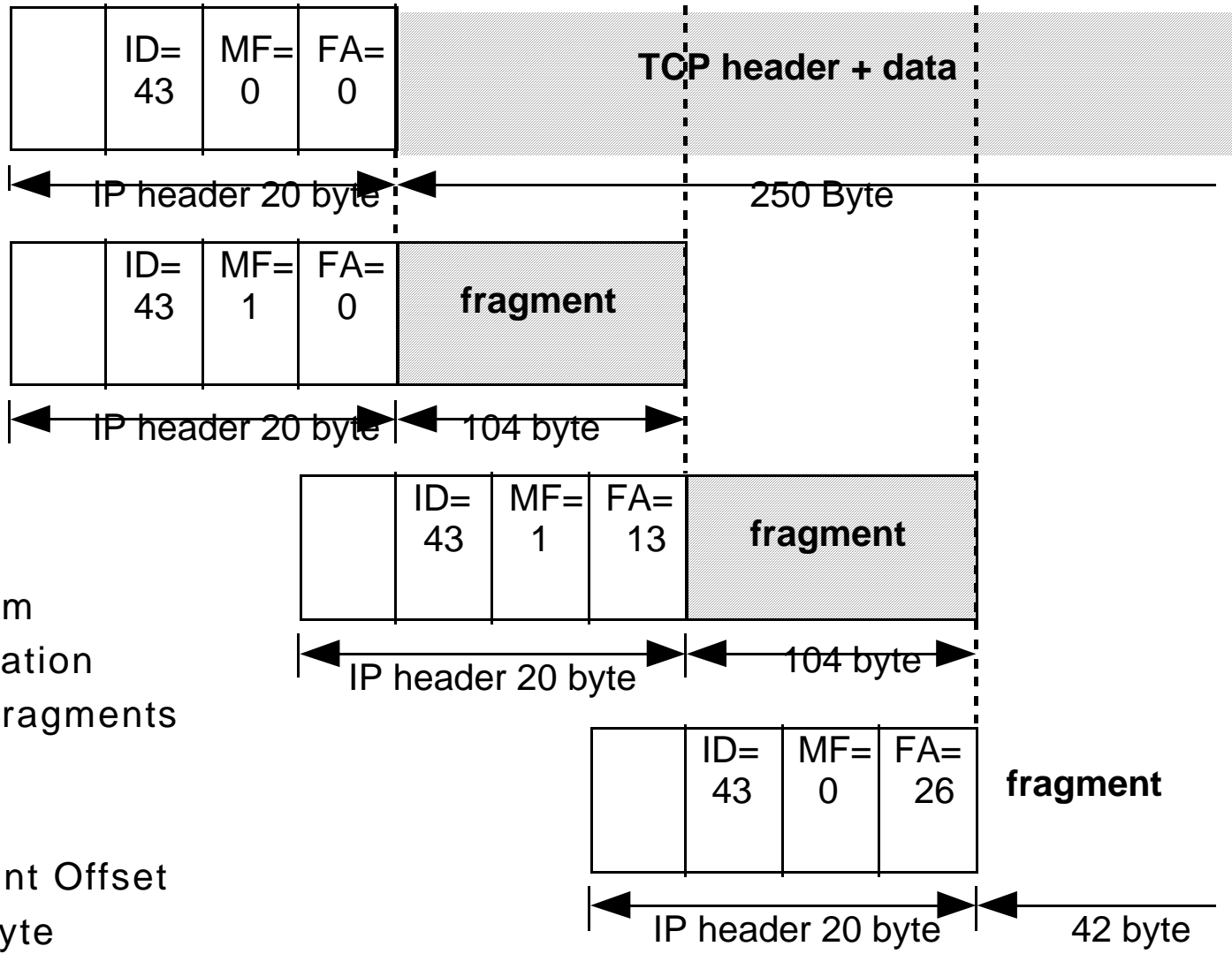
Transport layer

- **reassembles segments and fragments**



Segments, Fragmentation (and Reassembly) (2)

www.kom.tu-darmstadt.de
www.httc.de



ID: Datagram Identification
 MF: More Fragments
 0: no
 1: yes
 FA: Fragment Offset
 n: n*8 byte



5.2 Protocol with Flow Control

SOURCE PORT	
DESTINATION PORT	local endpoints of connection
SEQUENCE NUMBER	Number of transmitted bytes (each byte of the “message” is numbered)
ACKNOWLEDGEMENT	Byte number referring to the acknowledgement, specifies next byte expected
HLEN	Length of the header in 32 bit words. Needed since Options field is of variable length indicates start of data within segment (in 32-bit words)
RESERVED	not used
FLAGS	Bits from left to right
1. URG	Urgent Pointer is being used
2. ACK	AckNo is valid (if ACK=0 then no acknowledgment in segment)
3. PSH	data transferred with PUSH. Receiver should give received data to application immediately, not waiting & buffering until more / full buffer has been got
4. RST	Reset: connection is being reset
5. SYN	used to establish connections (synchronize seq. numbers) SYN=1 & ACK=0: connection request SYN=1 & ACK=1: connection accept
6. FIN	release connection



WINDOW SIZE

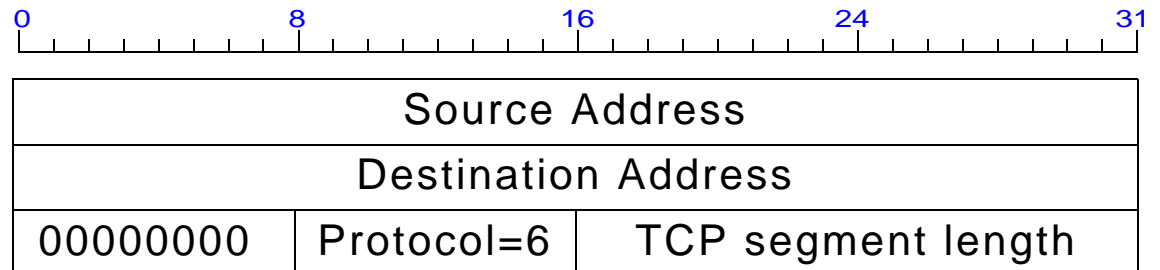
Buffer size in bytes
 Used for variable-sized sliding window.
 Window size field indicates #bytes sender may transmit beginning with byte acknowledged.
 Window size = 0 is valid: bytes up to ACK-1 received but receiver does not want more data at the moment.
 Later permission for more data by sending segment with same ACK and non-zero window size

CHECKSUM

Checksum header, data, and pseudoheader for calculation: TCP checksum field is set to 0 and data field padded with additional zero byte. If length is odd one's complement of sum of 16-bit halfwords in one's complement arithmetic.
 Receiver's calculation on entire segment including checksum field should result in 0

- pseudoheader

- part of IP header
 - source IP addr.
 - dest. IP addr.
 - protocol
 - length of TCP segment (including header)
- allows to detect misdelivered packets
- but violates protocol hierarchy





URGENT POINTER Byte offset to the current sequence number at which important data starts

OPTIONS for optional facilities, not covered by regular header, e.g., allows to specify max. TCP payload host may accept.

Without option, it defaults to 536 byte payload.

Window scale option (to deal with 64 KB window limitation).

Shift window size field up to 14 bits to left maximum window then 2^{30} byte.

Selective repeat instead of go-back-n

NAKs to allow receiver to request a specific segment



6. TCP: Connections & Management

Notion of connections

How connections work

How connections are "remembered"



6.1 Connection Establishment

One passive & one active side

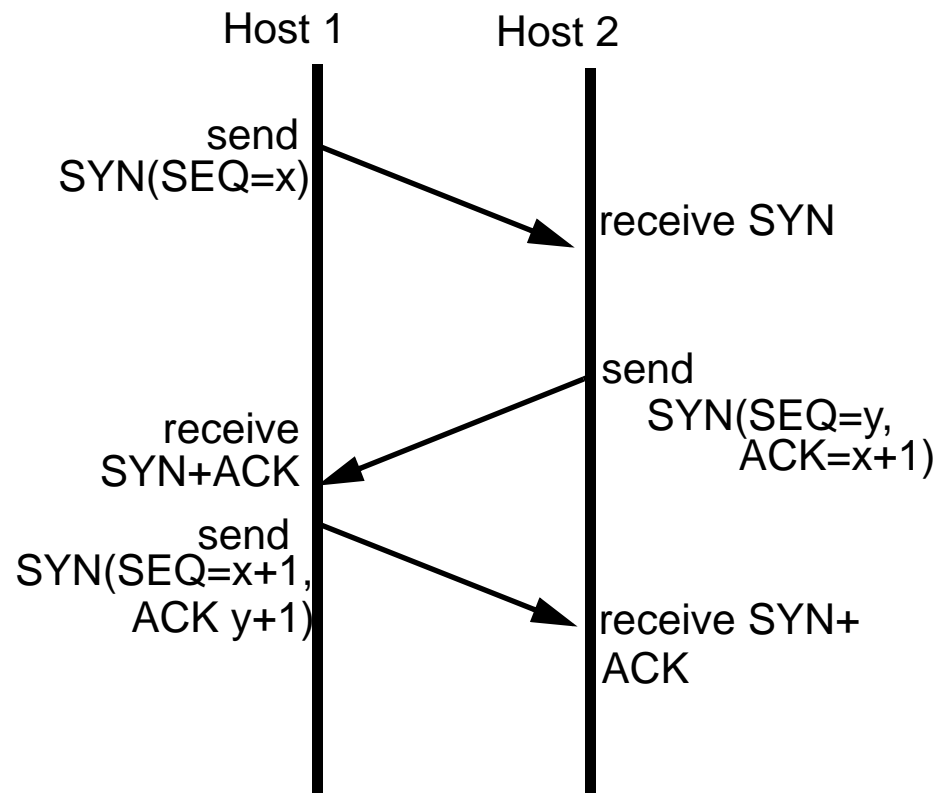
- **server:** wait for incoming connection using **LISTEN** and **ACCEPT**
- **client:** **CONNECT** (specifying IP addr. and port, max. TCP segment size)

Three-Way-Handshake

- **CONNECTING** through 3 packets

Comments

- **when establishing a connection**
 - initial sequence numbers of both partners are also exchanged and acknowledged
 - initial seq.# is not 0
- **SYN segment consumes one byte of sequence space**
 - in order to be acknowledged unambiguously





Connection Establishment

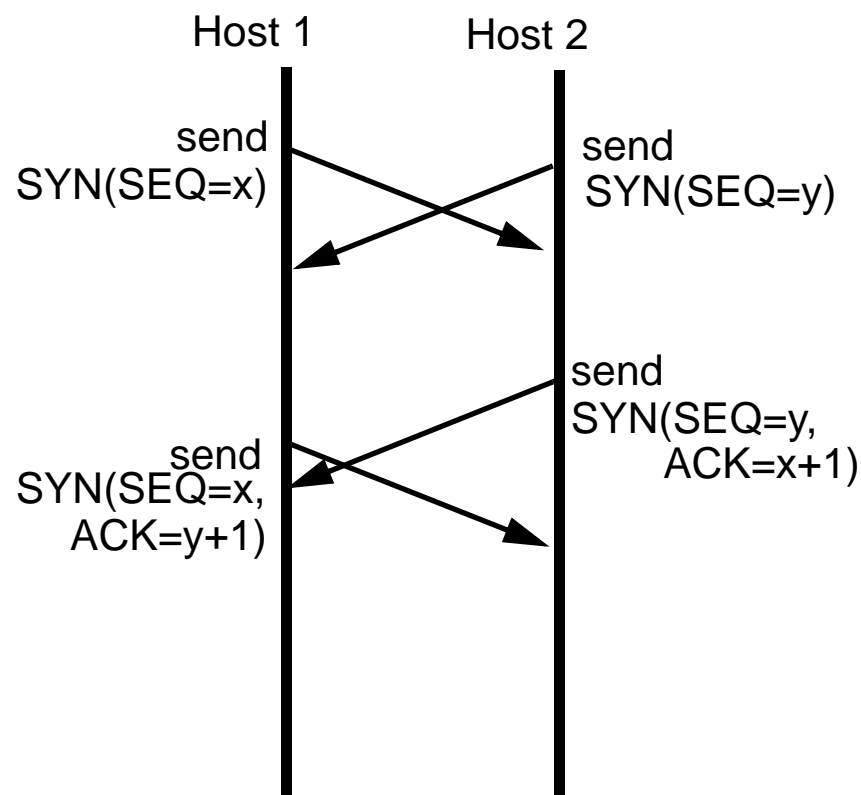
(2)

Comments

- **if on server side no process is waiting on port (no process did LISTEN)**
 - reply segment with RST bit set is send to reject connection attempt
- **process listening on port may accept or reject**

Call collision

- **still only one single connection will be established even when**
 - both partners actively try to establish a connection simultaneously





6.2 Connection Release

Connection release for pairs of simplex connections

- **each direction is released independently of other**

Connection release by either side sending a segment with FIN bit set

- **no more data to be transmitted**
- **when FIN is acknowledged, this direction is shut down for new data**

Directions are released independently:

- **other direction may still be open**
- **full release of connection if both directions have been shut down**



Connection Release

(2)

Systematic disconnect by 4 segments

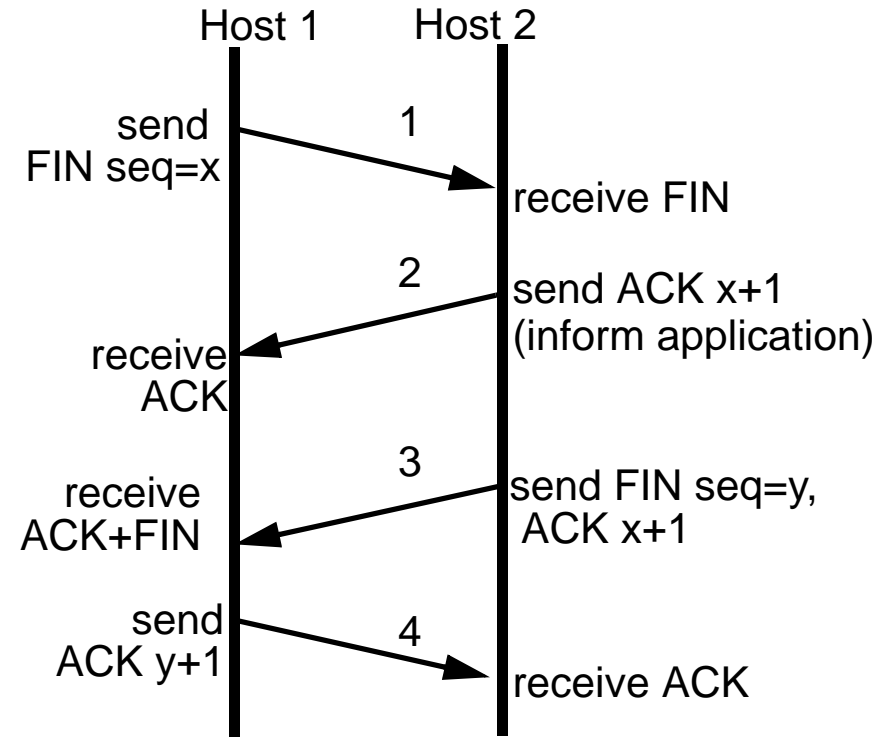
- **between 2nd and 3rd**
 - host 2 can still send data to host 1

Reduction of packet number possible:

- **first ACK and second FIN may be contained in same segment (3 segments instead of 4)**

CONNECTION INTERRUPT: Opposite side cannot transmit data anymore

- **immediate acknowledgement, release of all resources**
- **data in transit may be lost**





6.3 Connection Management Modelling

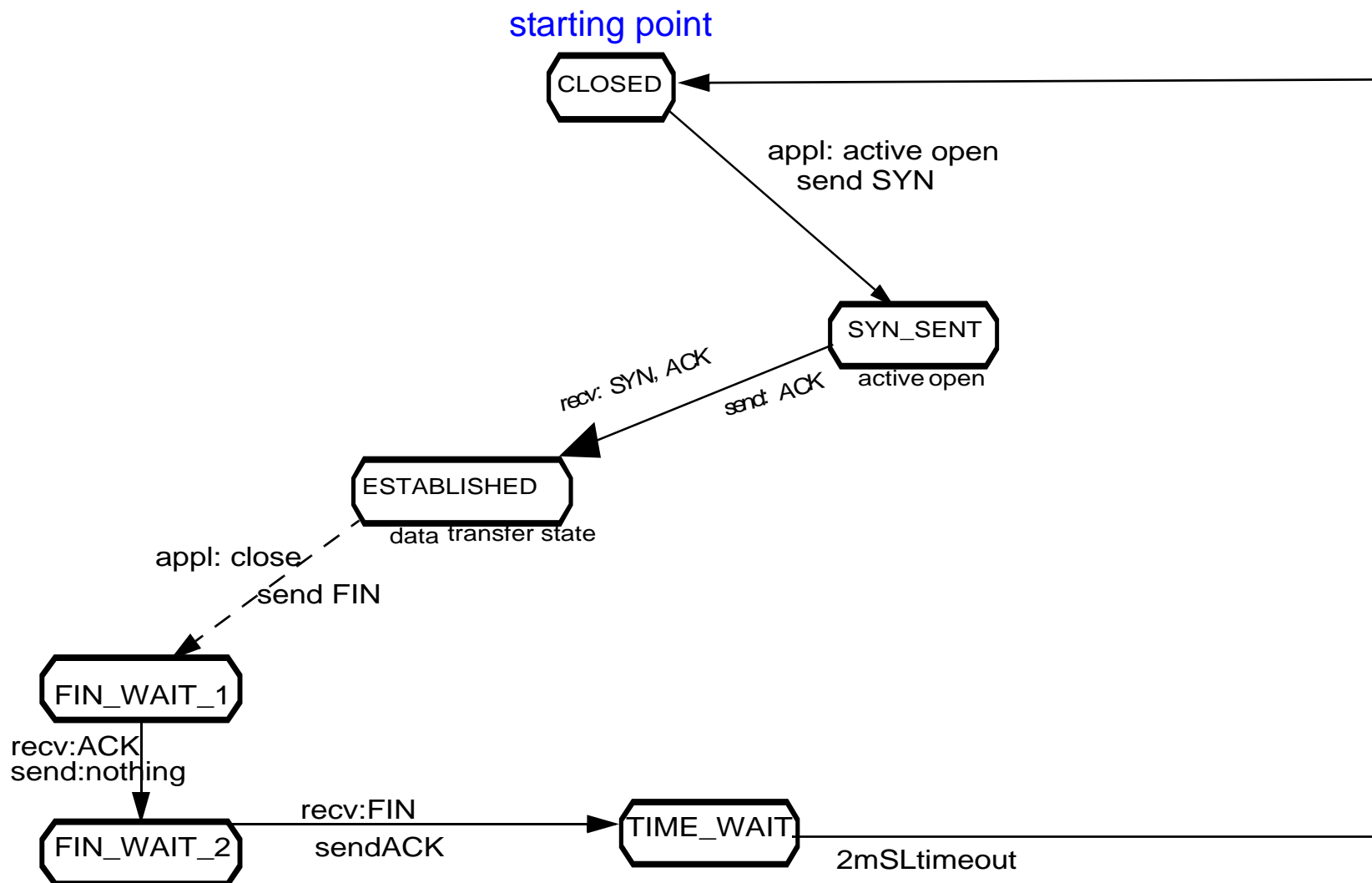
States

State	Description
CLOSED	No connection is active or pending
LISTEN	Server is waiting for an incoming call
SYN RCVD	Connection request has arrived; wait for ACK
SYN SENT	Application has started to open a connection
ESTABLISHED	Normal data transfer state
FIN WAIT 1	Application has said it is finished
FIN WAIT 2	The other side has agreed to release
TIMED WAIT	Wait for all packets to die off
CLOSING	Both sides have tried to close simultaneously
CLOSE WAIT	The other side has initiated a release
LAST ACK	Wait for all packets to die off



Connection Management Modelling: Finite State Machine

Typical sequence of TCP states visited by client TCP

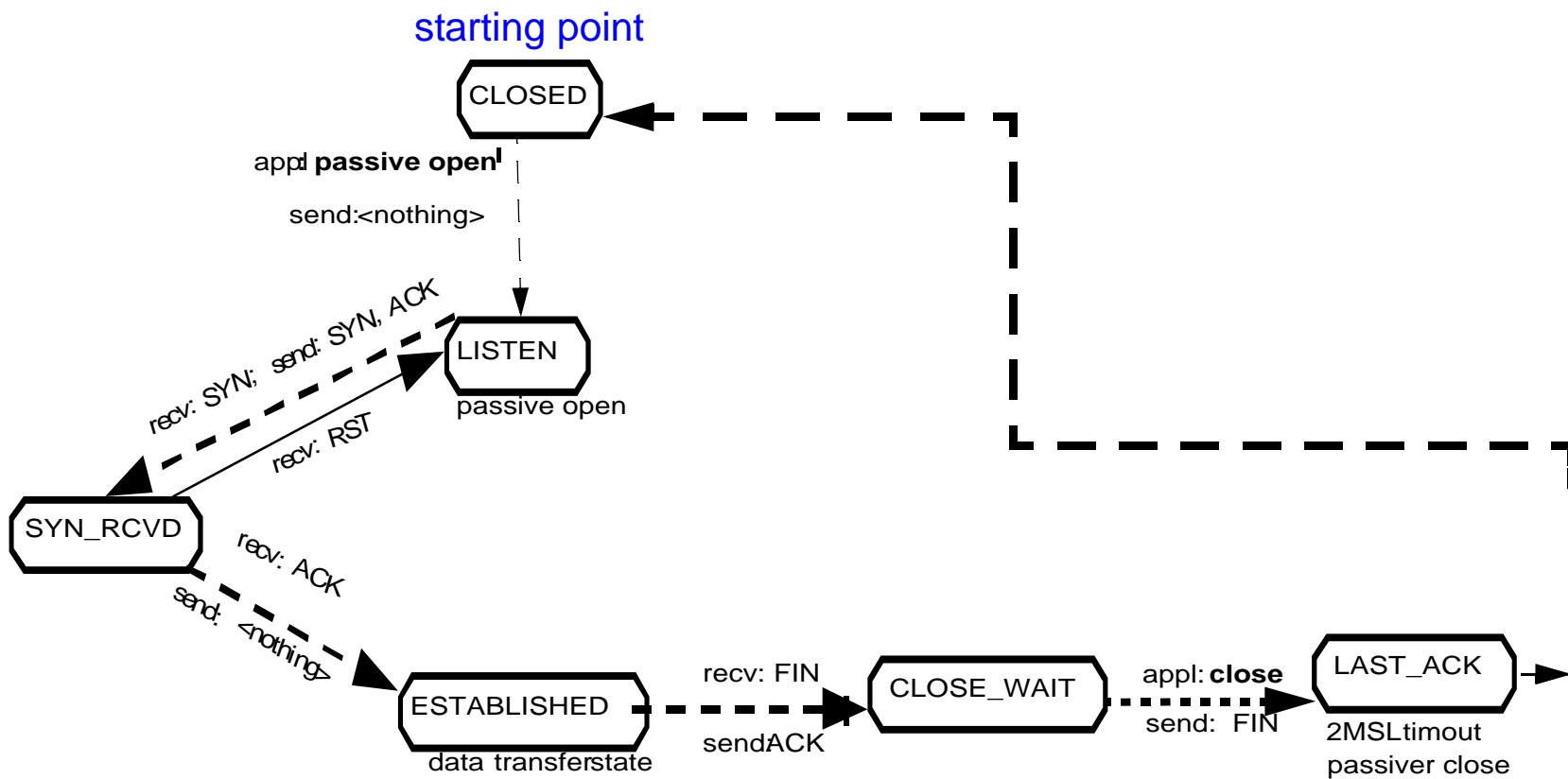




Connection Management Modelling: Finite State Machine

(2)

Typical sequence of TCP states visited by server-side TCP





7. TCP - Foundations

Some principles related to the TCP basics

- **flow control**
- **round trip time RTT**
- **congestion control**
 - Additive Increase: Example Slow Start
 - Multiplicative Decrease

And 2 problem areas: (to be considered separately)

- **receiver capacity** (“actual window”)
- **network capacity** (“congestion window”)



7.1 Flow Control

Flow control: “sliding window” mechanism

- **acknowledgement and sequence number**
 - acknowledgments refer to byte positions
 - not to whole segment
 - sequence numbers refer to the byte position of a TCP connection
- **positive acknowledgement**
- **cumulative acknowledgements**
 - byte position in the data stream up to which all data has been received correctly
 - reduction of overhead

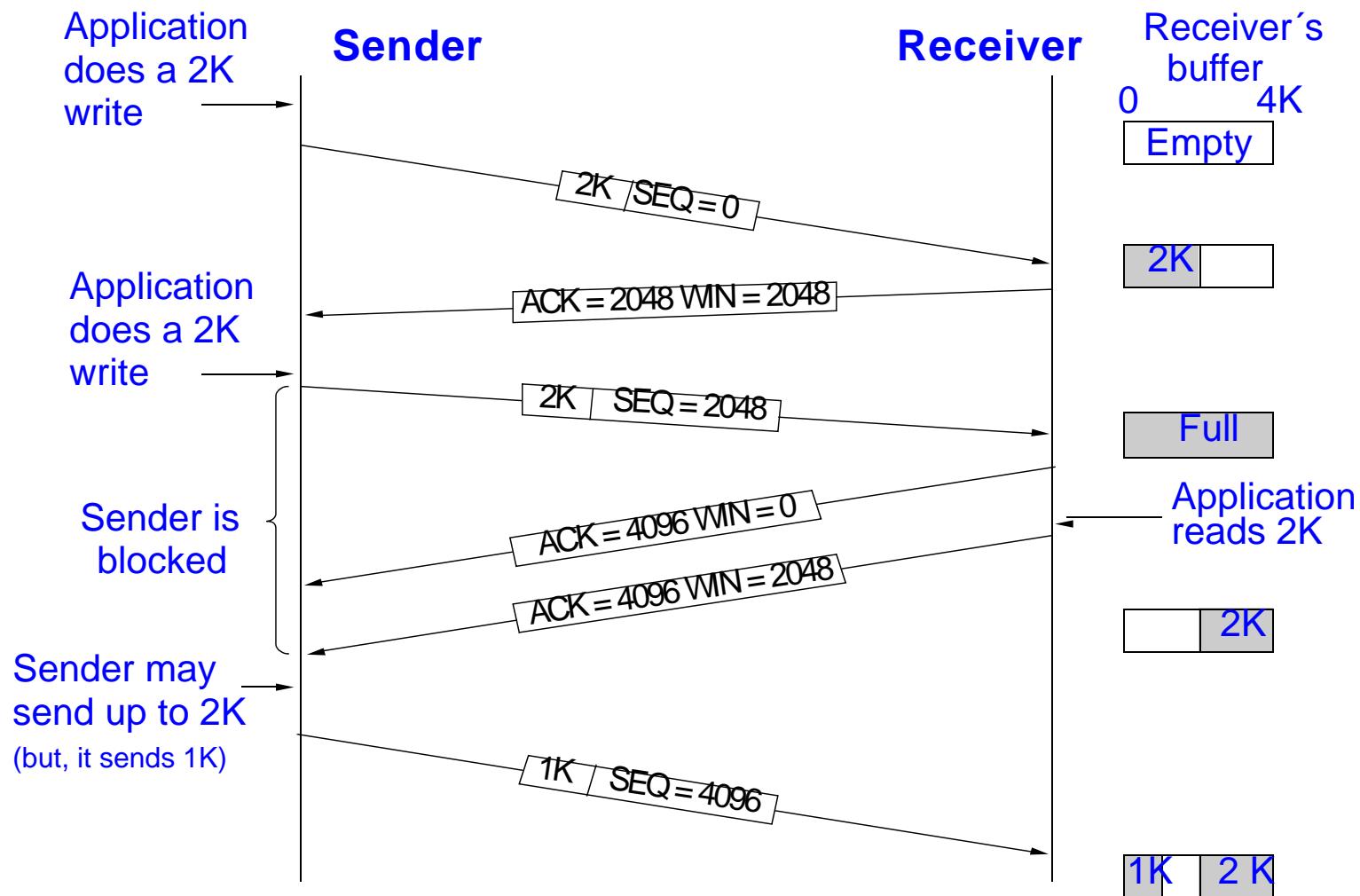
Variable window sizes (credit mechanism)

- **implementation**
 - the actual window size is also transmitted with each acknowledgement
 - permits dynamic adjustment to existing buffer



TCP Flow Control: Example

www.kom.tu-darmstadt.de
www.httc.de





TCP Flow Control: Special Cases

Optimization for low throughput rate

- **problem: Telnet (and ssh) connection to interactive editor reacting on every keystroke**
 - 1 character typed requires up to 162 byte
 - data: 20 bytes TCP header, 20 bytes IP header, 1 byte payload
 - ACK: 20 bytes TCP header, 20 bytes IP header
 - editor echoes character:
20 bytes TCP header, 20 bytes IP header, 1 byte payload
 - ACK: 20 bytes TCP header, 20 bytes IP header
- **approach often used**
 - delay acknowledgment and window update by 500 ms (hoping for more data)

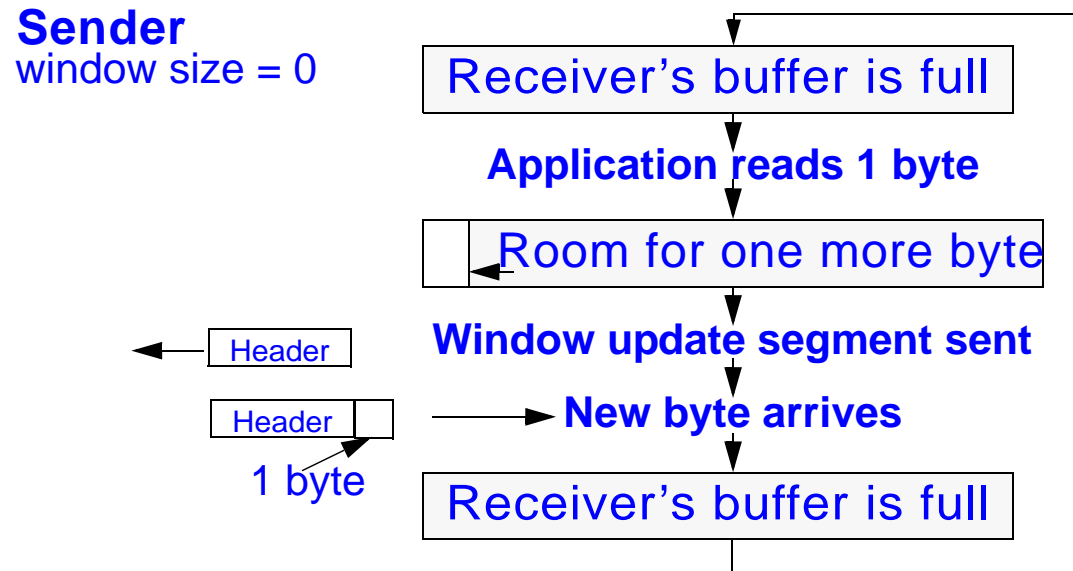
Nagle's algorithm, 1984

- **algorithm**
 - send first byte immediately
 - keep on buffering bytes until first byte has been acknowledged
 - then send all buffered characters in one TCP segment and start buffering again
- **comment**
 - effect at e.g. X-windows: jerky pointer movements



Silly window syndrome (Clark, 1982)

- **Problem:**
 - data on sending side arrives in large blocks
 - but receiving side reads data at one byte at a time only



- **Clark's solution:**
 - prevent receiver from sending window update for 1 byte
 - certain amount of space must be available in order to send window update
 - $\min(X, Y)$
X = maximum segment size announced during connection establishment
Y = buffer / 2



7.2 Timer Management

TCP uses several timers for different purposes

Retransmission timer as most important one

- **timer is set when segment is sent**
- **if ACK arrives before timer expires: timer is stopped**
- **if timer expires before ACK arrives: segment is retransmitted**
 - and new timer is set

Question: How long should the timeout interval be?



Timer Management

(2)

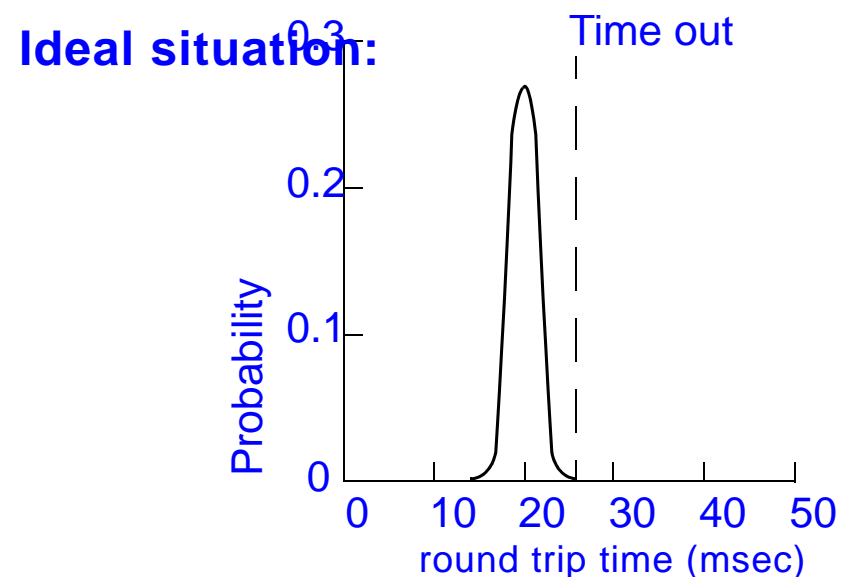
Timeout after a certain time period

Situation in the Internet:
packet communication time may vary immensely

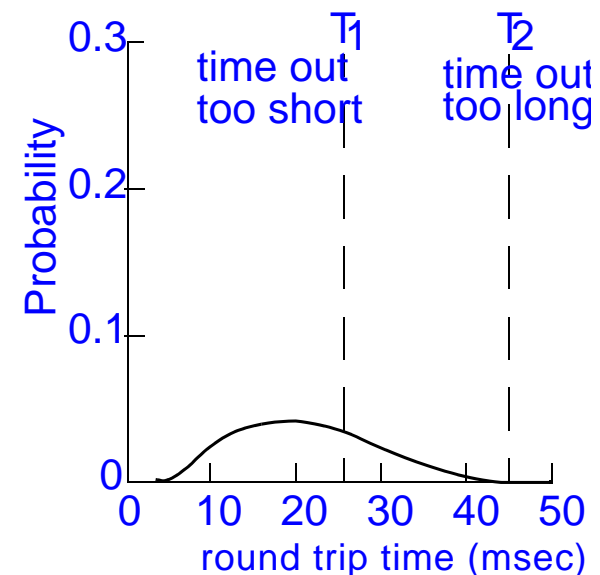
- **timeout too short**
 - fast reaction due to any errors
 - but correct packets retransmitted
- **timeout too long**
 - too slow reaction due to errors
 - less additional traffic

Better: adaptive timeout intervals

- **timeout period (until retransmit)**
 - continuous adaptation to pre-calculated, expected waiting period
- **waiting period increases for each retransmission of the same segment**



Example of real situation





Relevant parameters

- **RTT: Round Trip Time (or Round Trip Sample)**
 - time between
 - sending of octet and the receiving of respective acknowledgment
- **timeout: initiates retransmission**

Algorithm (Jacobson, 1988)

$$RTT = (\alpha \times Old_RTT) + ((1 - \alpha) \times New_Round_Trip_Sample)$$

$$Timeout = \beta \times RTT$$

α Smoothing factor-

- 0..1 determines importance of the history
- 0: only current/last value is relevant
- **7/8 TYPICAL VALUE**
- 1: only old values are relevant

β (later/today proportional to the average deviation)

- 1: faster error correction
- 2: original value in implementations
- **COMMENT:** usually hard to define (as fixed value)



RTT Calculation in Case of Error Recovery

How to handle dynamic RTT estimation when segment must be resent?

Situation: ambiguous measurement

- t_1 : retransmit initiated because of timeout
- t_2 : acknowledgement arrives at the receiver's

Question:

- **does confirmation refer to**
 - 1. data which has previously been thought as being lost?
or
 - 2. data which has been sent last?
- **problem: no differentiation possible**

This influences the definition/parameter setting of the timeout

- **by the RTT measured**



RTT Calculation in Case of Error Recovery

(2)

Situation: ambiguous measuring

1ST assumption: reference to the data sent first (originally)

• **but situation is following**

- short-term increase in the delay
 - for any unknown reason
- but not really any loss of data
 - acknowledgement arrives shortly after timeout

• **result**

- New-Round_Trip_Sample: has very high values
- RTT increases

• **sender**

- now, if sender detects (at later time) packets that have been lost
 - retransmits at a (very) later point in time
 - consequently acknowledgements are also returned later

• **global effect:**

- RTT may increase even more and more



Situation: ambiguous measuring

2ND assumption: reference to last retransmitted data

- **but situation occurred due to**
 - Internet in fact losses data
- **result**
 - New_Round_Trip_Sample: has smaller value
 - RTT is reduced
- **sender**
 - may set time interval which may be too short
 - the acknowledgement of the first packages will again further reduce RTT
- **global effect:**
 - RTT gets too small
 - usually at 1/2 of the actual value
 - i.e. on the average each segment will be sent twice



RTT Calculation in Case of Error Recovery

(4)

Fix: (Phil) Karn's algorithm

- **ignore the measurements of lost packets**
 - but this would not reflect actual system changes in a correct manner
- **previous algorithm**

$$RTT = (\alpha \times Old_RTT) + ((1 - \alpha) \times New_Round_Trip_Sample)$$

$$Timeout = \beta \times RTT$$

- **timer backoff strategy**

$$New_Timeout = \gamma \times TimeOut$$

- increase timeout for each lost packet (actually each segment) until a segment has arrived
- γ -value: typically 2
- **comment:**
 - originated from Phil Karn
 - Karn used TCP based on (very unreliable) radio transmission



Further improvements

- **Karn's algorithm**

- bad adjustment whenever
 - variation of consecutive delays is large
 - e.g. short-long-short-long ...

- **high load means**

- often large fluctuations within the delays
- larger deviation (as well as average runtime deviation)
"höhere Standardabweichung (und mittlere Abweichung)"

- **therefore deviation to be used (instead of β)**

- previously

$$RTT = (\alpha \times \text{Old_RTT}) + ((1 - \alpha) \times \text{New_Round_Trip_Sample})$$

$$\text{Timeout} = \beta \times RTT$$

- new algorithm

$$\text{Timeout} = \text{Function}(\text{AverageRuntimeDeviation})$$



7.3 Congestion Control

(understandably)
diverse objectives

- **endsystem**

- optimize its own throughput
- possibly at the expense of other endsystems

- **network**

- optimizes overall throughput

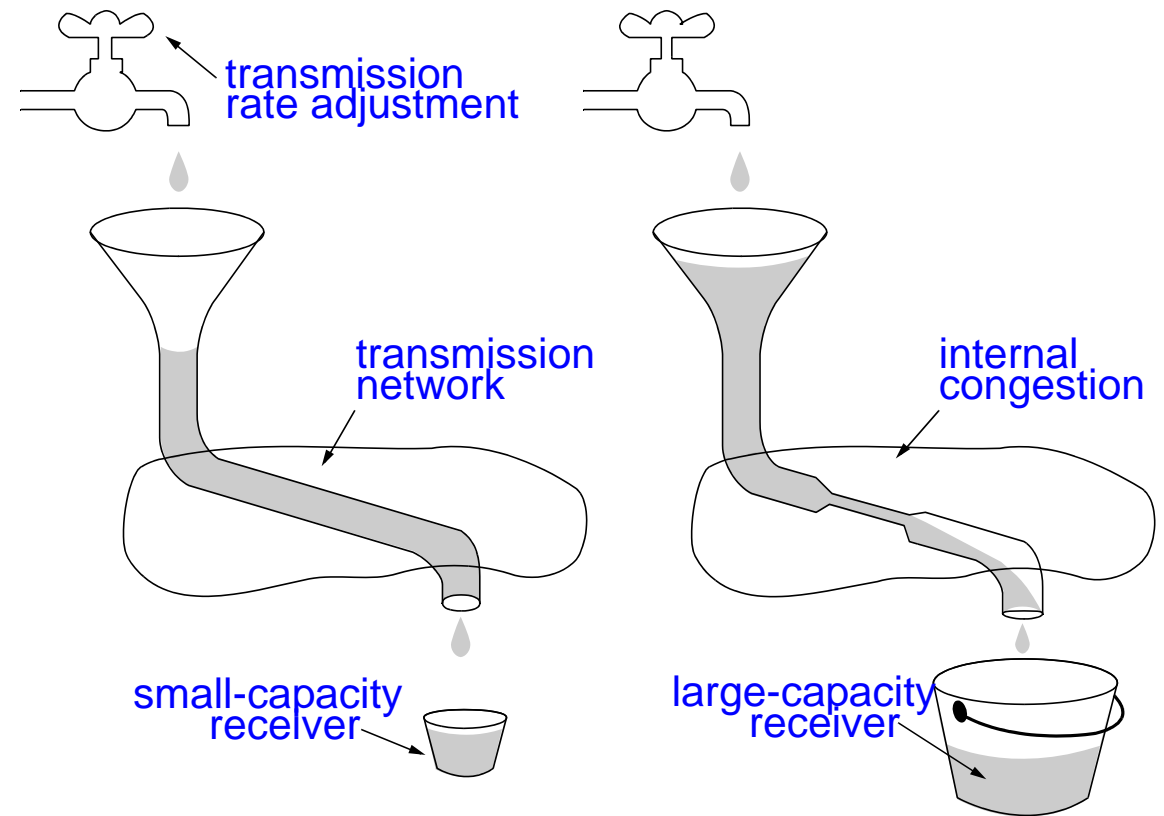
in example:

- **two different problems**

- receiver capacity
- network capacity

- **cannot be distinguished easily at all places**

- **should be differentiated**



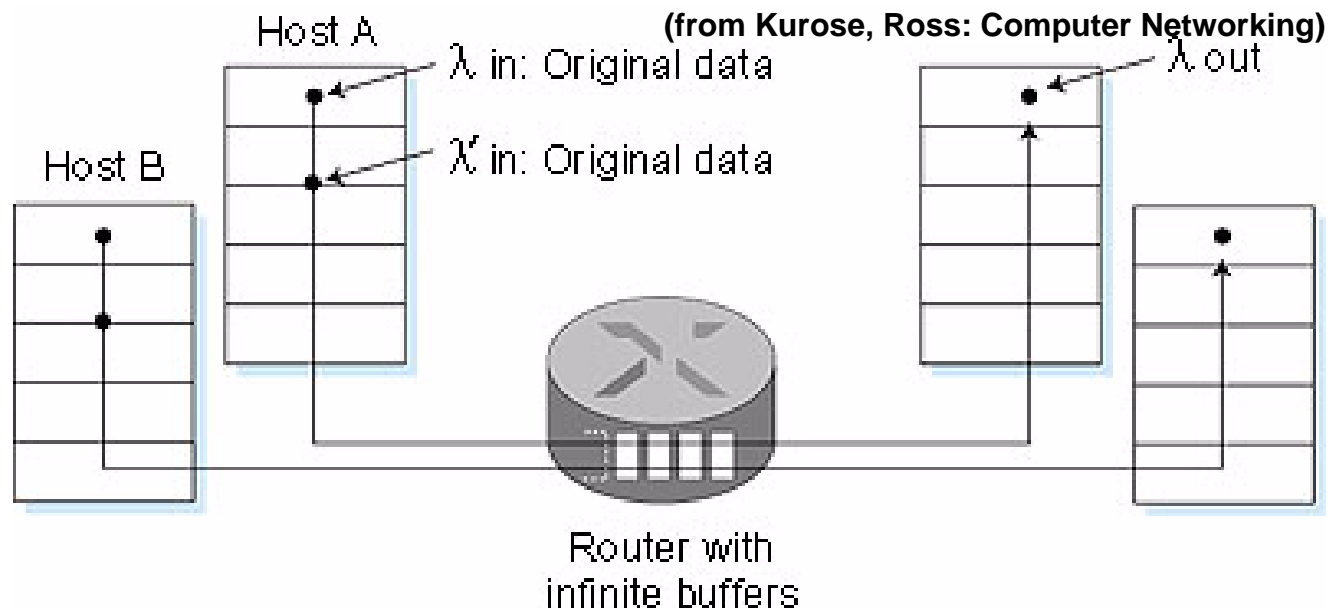


Congestion Considerations

Problem statement: throughput (easier) and delay (more critical)

e.g.

- 2 senders
- router with infinite buffer
- no error recovery
- no flow control

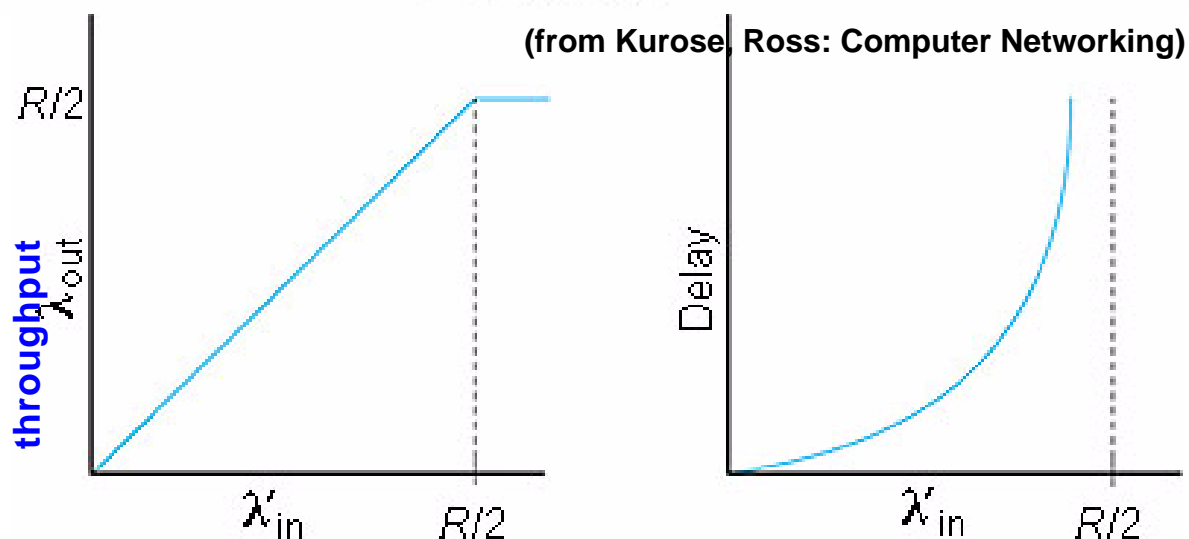


throughput per connection

- scales linear

delay

- scales exponentially
- large queueing delays as consequence of packet-arrival rate near link capacity !





Congestion Considerations

(2)

General problem

- congestion \Rightarrow more delays at ES
- higher delays \Rightarrow retransmissions (because of timeouts)
- i.e. \Rightarrow additional load increases

i.e. 2 problem areas: (to be considered separately)

- **receiver capacity** (“actual window”)
- **network capacity** (“congestion window”)

i.e. to be applied to the send window

$$(\text{valid send window}) = \text{Min} (\text{actual window}, \text{congestion window})$$

TCP strategy

- **TCP reduces transfer rate at high network load**
- **methods**
 - “slow start” generalization: additive increase
 - “congestion avoidance” generalization: multiplicative decrease



Congestion - Additive Increase: Example Slow Start

Objective

- **to avoid immediate network overload**
 - after just recent overload has finished
 - after timeout
 - at the start of a data transmission

Method: congestion window

- defined/initiated at connection set-up
 - initial max. window size = 1 segment
 - threshold (based on previous max. window size)
- max. size of segment to be transferred
 - to be doubled as long as
 1. segment size below a certain threshold and
 2. all ACKs arrive in time (i.e. correct segment transfer)(each burst acknowledged, i.e. successfully transmitted, doubles congestion window)
- to be linearly increased
 1. segment size above a certain threshold and
 2. all ACKs arrive in time (i.e. correct segment transfer)



Delay - Multiplicative Decrease: Congestion Avoidance

Threshold

- adaptive
- parameter in addition to the actual and the congestion window

Assumption

- threshold, i.e. adaptation to the network: “sensible window size”

Use: On timeout

- threshold is set to half of current congestion window
- congestion window is reset to one maximum segment
- use slow start to determine what the network can handle
 - exponential growth stops when threshold is hit
 - from there congestion window grows linearly (1 segment) on successful transmission

Congestion window and threshold

- **if**
 - congestion window $<$ threshold
 - exponential growth: congestion window doubled
 - congestion window \geq threshold
 - linear growth: increase congestion window by 1 segment each



Congestion Control: Example

Parameters:

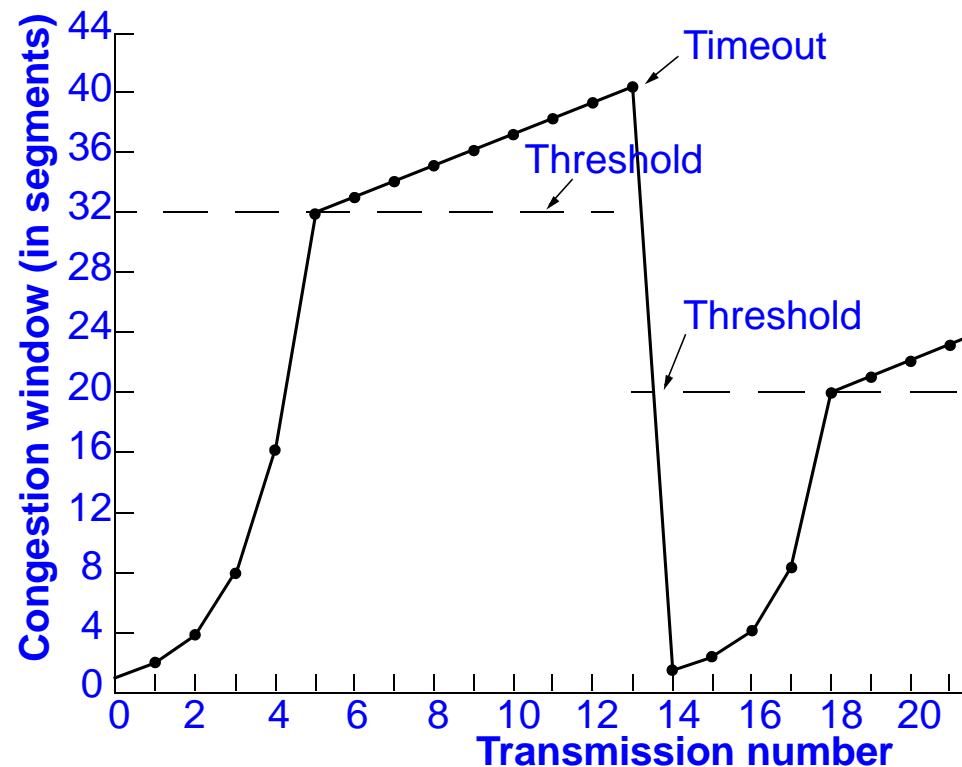
- maximum segment size = 1024 bytes
- initial congestion window = 64 KB

Areas

- **previously (transm. nr. = 0)**
 - congestion window = 64 KB
 - timeout: i.e. threshold 32 KB
 - congestion window = 1 KB
- **exponential area:**
 - “slow start”
- **linear area**
 - “congestion avoidance”

Note: but always

- **observe receiver’s window size, i.e. use minimum of**
 - congestion window
 - actual send window (receiver status)

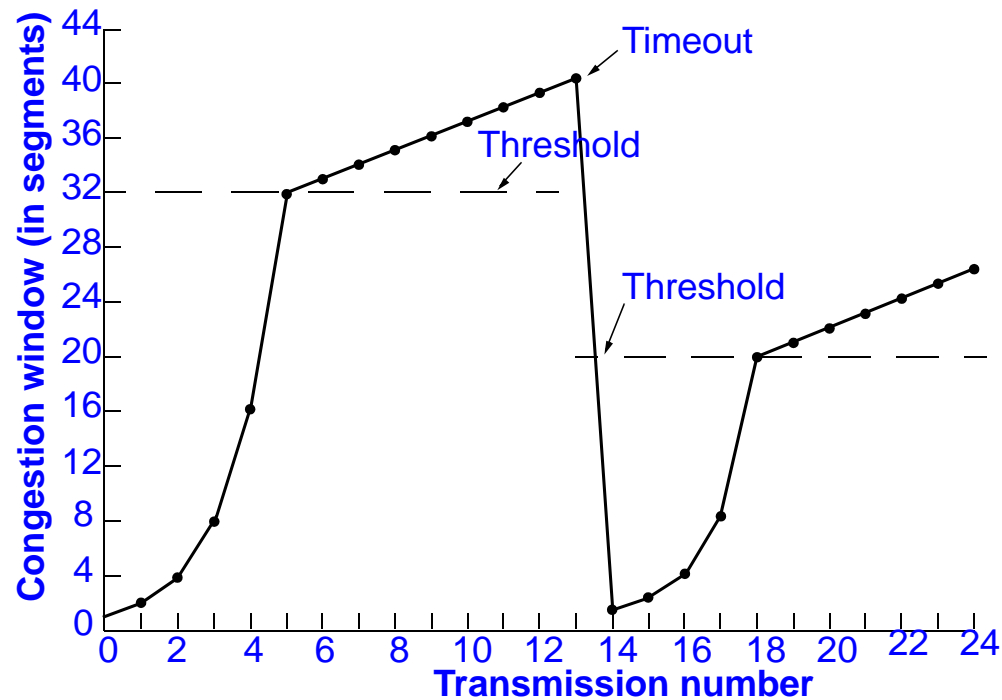
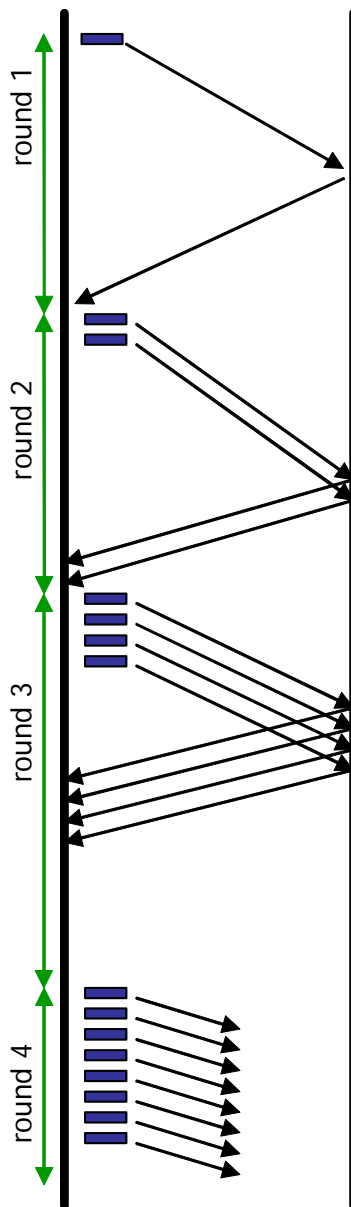




Congestion Control: Example

(2)

time diagram
of
exponential
area





7.4 TCP - Further Comments

Some parameters

- **65.536 byte max. per segment**
- **IP recommended value TTL interval 2 min**

Optimization for low throughput rate

- **problem**
 - 1 byte data requires 162 byte incl. ACK
(if, at any given time, it shows up just by itself)
- **algorithm**
 - acknowledgment delayed by 500 msec because of window adaptation
- **comment**
 - often part of TCP implementation



TCP Friendliness - TCP Compatible

A TCP connection's throughput is bounded

- **wmax** - maximum retransmission window size
- **RTT** - round-trip time

Congestion windows size changes

- **AIMD** (additive increase, multiple decrease) algorithm

TCP is said to be fair

- **Streams that share a path will reach an equal share**

A protocol is TCP-friendly if

- **Colloquial:**
 - It long-term average throughput is not bigger than TCP's
- **Formal:**
 - Its arrival rate is at most some constant over the square root of the packet loss rate



Application Area - Wireless

TCP characteristics

- if data is late, TCP assumes that network is congested
- therefore “slow start” algorithm

Mobile communications channel

- is unreliable and may contain losses

TCP's behavior on radio channels

- TCP slows down transmission, inefficient
- here it would be better to increase the rate



Bandwidth-Delay Product

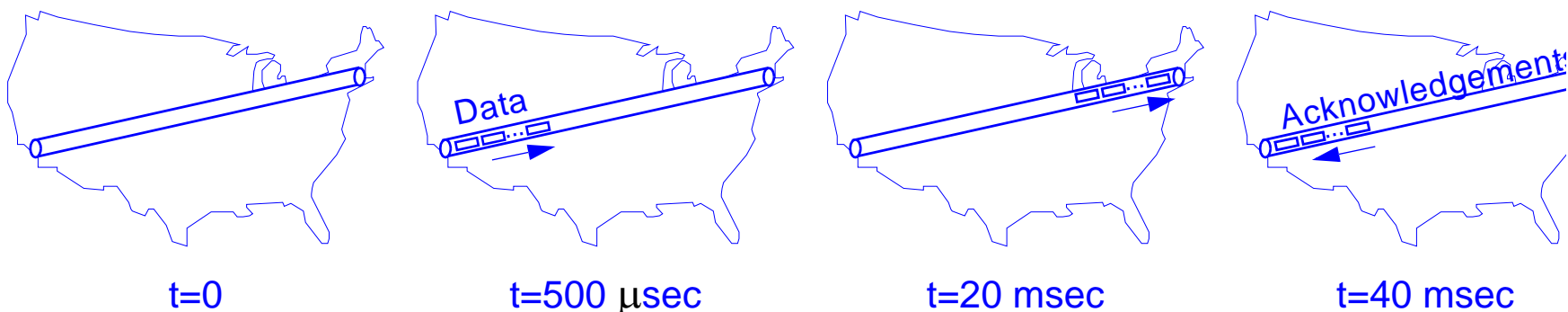
$$\text{BANDWIDTH [BITS/SEC]} * \text{ROUND-TRIP DELAY [SEC]}$$

Useful parameter for network performance analysis

- Capacity of pipe from sender to receiver and back (in bits)

Example:

- Transmission from San Diego to Boston
 - sending 64 KB burst (receiver buffer 64 KB), link: 1 Gbps
 - one-way propagation delay (speed-of-light in fiber): 20 msec



- Bandwidth-delay product: **40 MILLION BIT**
- i.e.: sender would have to transmit burst of 40 million bits to keep pipe busy till ACK

Receiver window must be \geq bandwidth-delay product

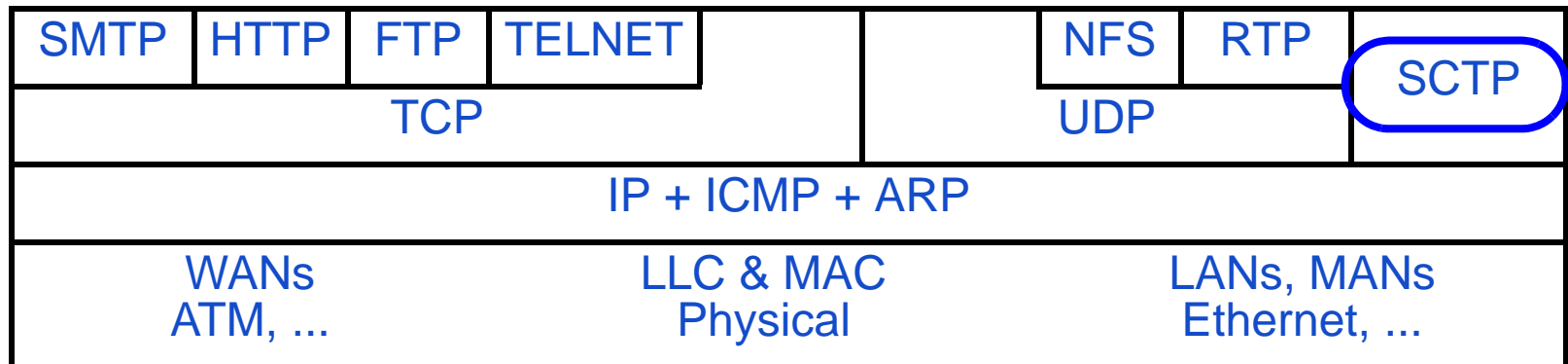
- for good performance



8. Stream Control Transmission Protocol (SCTP)

For signalling with high reliability but, low overhead

⇒ Additional transport protocol from IETF (add. to TCP+UDP)



Specification in

- **RFC 2960**
 - **Stream Control Transmission Protocol**
- **RFC 2719**
 - **Architectual Framework for Signaling Transport**
- **RFC 3057**
 - **ISDN Q.921-User Adaptation Layer**
- **see also**

<http://www.sctp.de/>



SCTP: Motivation

TCP too limited for some applications:

- **e.g., transport signaling from PSTN networks (SS7) over IP-based networks**

Goals:

- **initial goal:**
 - replace SS7 signaling in PSTN with SCTP
- **now:**
 - SCTP as a universal transport protocol (e.g., for HTTP)
- **future: replacing TCP??**

Examples:

- **Strict order-of-transmission delivery of data with multiple streams**
 - ⇒ **partial order within a stream of multiplexed streams sufficient**
- **Stream-orientation of TCP inconvenient**
 - application must set record markings
 - ⇒ **better: message-orientation**
- **TCP cannot deal with multi-homing**
 - i.e., one server with several IP addresses
- **TCP is vulnerable to DoS attacks**
 - e.g., SYN flooding



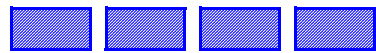
SCTP Concepts: Association and Streams

Connection-oriented:

- **concept of 'association'**

- bi-directional
- generalization of TCP-connections:
 - each association endpoint can have several IP addresses (multi-homing)
 - each association can contain several streams (multi-streaming)

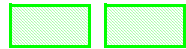
Stream A:



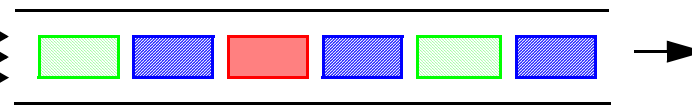
Stream B:



Stream C:



Association



- Stream: sequence of user messages to be delivered in order (up to 216 per direction)

⇒ in contrast to the notion of 'stream' of TCP

- **Reliable data transfer**

- confirmed, no duplicates, error-free



SCTP Concepts: Strict vs. Partial Order

Strictly ordered delivery optional

- packets of a stream within an association are delivered in order
 - partial order
- optional: retain order between packets of all streams
 - strict order
- **effects**
 - strict order: data transmission stalled if one stream is stalled
 - partial order: transmission for non-stalled streams can continue

Stream A:

1 2 3

Stream B:

1

Stream C:

1 2

Association sent
to the network

2 3 1 2 1 1

2 3 1 2 1 1

STRICT ORDER

Association arriving
at the receiver

PARTIAL ORDER

1 3 2 1 2 1

Example: HTTP with multiple embedded files (images)

- **Order of arrival of image data not relevant**
- **Retrieving the text can continue even if loading the image is blocked**
 - e.g., if the image is located on a different server which is highly loaded



SCTP: Further Concepts

Message segmentation according to path-MTU

- **Path-MTU:**
 - maximum transfer unit supported on the path between the endpoints
- **Path-MTU discovery mechanism as specified in RFC 1191**

Test whether the communication partner is alive: Heartbeats

Flow control and congestion control similar to TCP (Selective Ack,...)

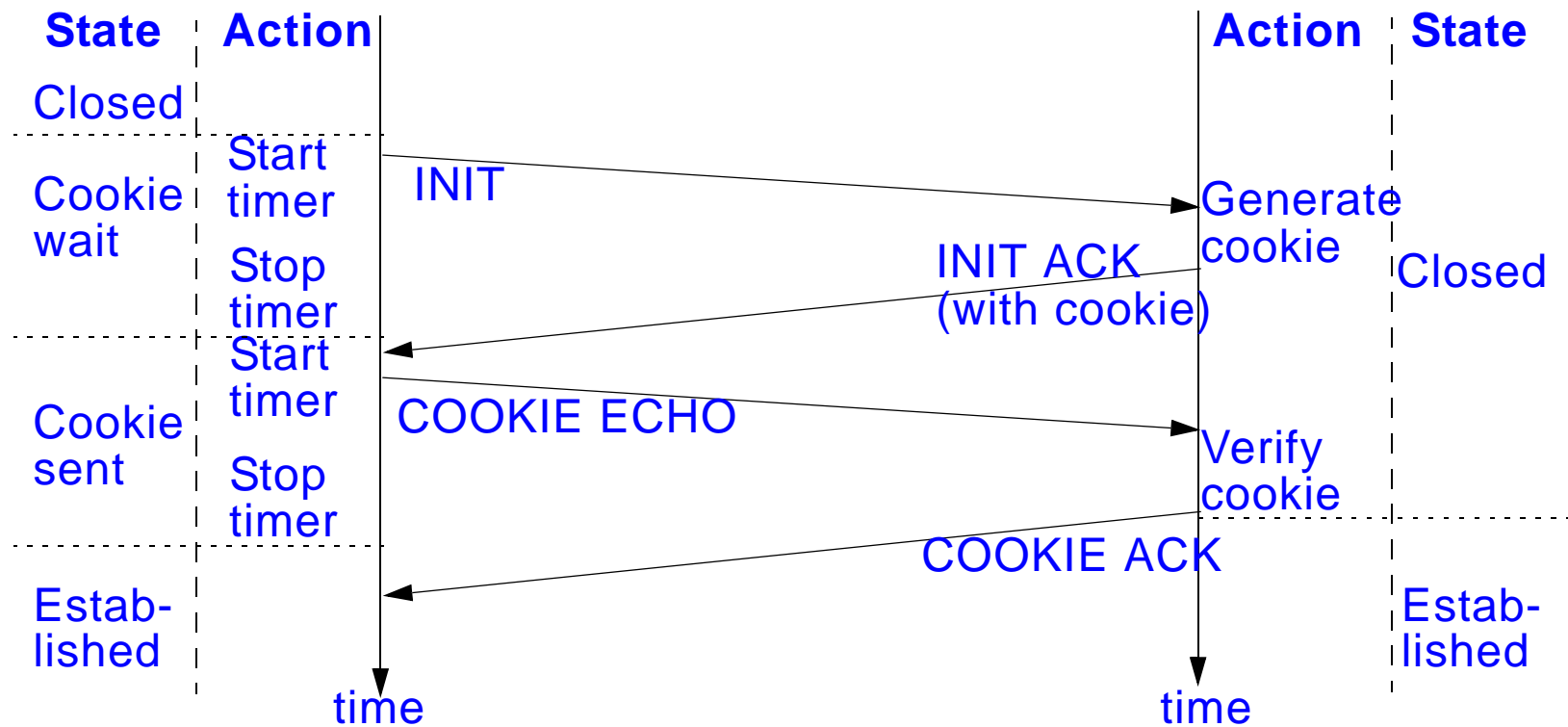
- **Coexistence with TCP**

Security means:

1. 32-bit checksum (Adler-32, CRC-32 under discussion)
2. 4-way handshake using cookies against DoS attacks



SCTP: 4-way Handshake



No half-open states as in TCP

No state information kept at the station receiving the 'INIT' message

- **no vulnerability for SYN flooding**
- **state information established only after the third step, the 'COOKIE' message**

To increase efficiency

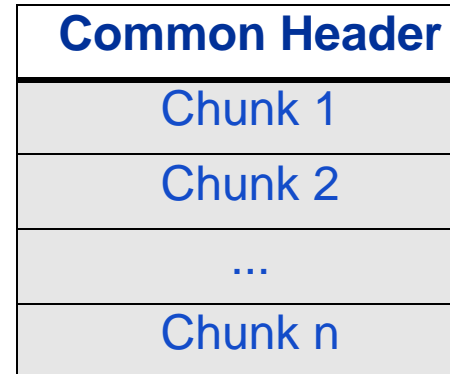
- **user data can be sent already with the 'COOKIE' and 'COOKIE ACK' messages**



SCTP: Message format

Multiplexing of several user messages: Chunk Bundling

- Chunk: part of an SCTP packet belonging to a single stream



Common Header:

Source Port	Destination Port
Verification Tag	
Checksum	

- **source port / destination port (2 Byte each): As in TCP or UDP**
- **verification tag: for validation of the sender of the SCTP message**
 - protection against blind attacks (unauthorized shutdown of association)
- **checksum:**
 - Adler-32: currently proposed in the RFC
 - CRC-32: proposed recently, better error detection properties for small packets



SCTP: Chunk Format

General format of a chunk:

Chunk Type	Chunk Flags	Chunk Length
Chunk Value		

- **Chunk Type:**
 - type of information in the Chunk Value field
 - e.g., type=0: payload data; type=1: INIT message,...
- **Chunk Flags:**
 - depend on the chunk type
- **Chunk Length:**
 - length of the chunk in bytes
 - including type, flags, length



SCTP: INIT Chunk Format

Type=1	Flags: None	Chunk Length
Initiate Tag		
Advertised Receiver Credit Window		
Number of outbound streams	Number of inbound streams	
Initial Transmission Sequence Number		
Optional/Variable-Length Parameters		

- **Initiate Tag:** random number used in all subsequent messages
 - protect against blind attacks
- **Advertised Receiver Credit Window:**
 - dedicated buffer space reserved for the association
- **Number of outbound streams:**
 - the sender of this INIT chunk wants to open
- **Number of inbound streams:**
 - maximum the sender of this INIT message can support
- **Variable-Length parameter**
 - a.o.: list of IP addresses (**MULTI-HOMING!**) being part of the association



SCTP: DATA Chunk Format

Example: Format of a DATA chunk (i.e., non SCTP control message)

Type=0	Flags: U B E	Chunk Length
Transmission Sequence Number (TSN)		
Stream Identifier S	Stream Sequence Number n	
Payload Protocol Identifier		
User Data (seq n of Stream S)		

- **Flags (one bit each):**
 - **U**nordered bit: '1' indicates unordered DATA chunk, '0' an ordered DATA chunk
 - **B**eginning bit: '1' first fragment of a user message
 - **E**nd bit: '1' last fragment of a user message
- **Transmission Sequence Number (32 bits) => 4,300 Mio numbers**
 - unique per stream
 - used for acknowledgments and duplicate recognition
 - acknowledgements are given per message (selective ACK)
 -
- ...



SCTP: DATA Chunk Format

(2)

Type=0	Flags: U B E	Chunk Length
Transmission Sequence Number (TSN)		
Stream Identifier S	Stream Sequence Number n	
Payload Protocol Identifier		
User Data (seq n of Stream S)		

- ...
- **Stream Identifier (16 bits)**
- **Stream Sequence Number (16 bits)**
 - unique per stream
 - used to assure sequenced delivery within a stream
 - separate acknowledgement mechanism from sequenced delivery
 - If fragmentation is necessary, stream sequence number is the same for all fragments
- **Payload Protocol Identifier (e.g., RTP)**
- **User Data: the actual user data to send via SCTP**
-



9. Further Development of Transport Protocols

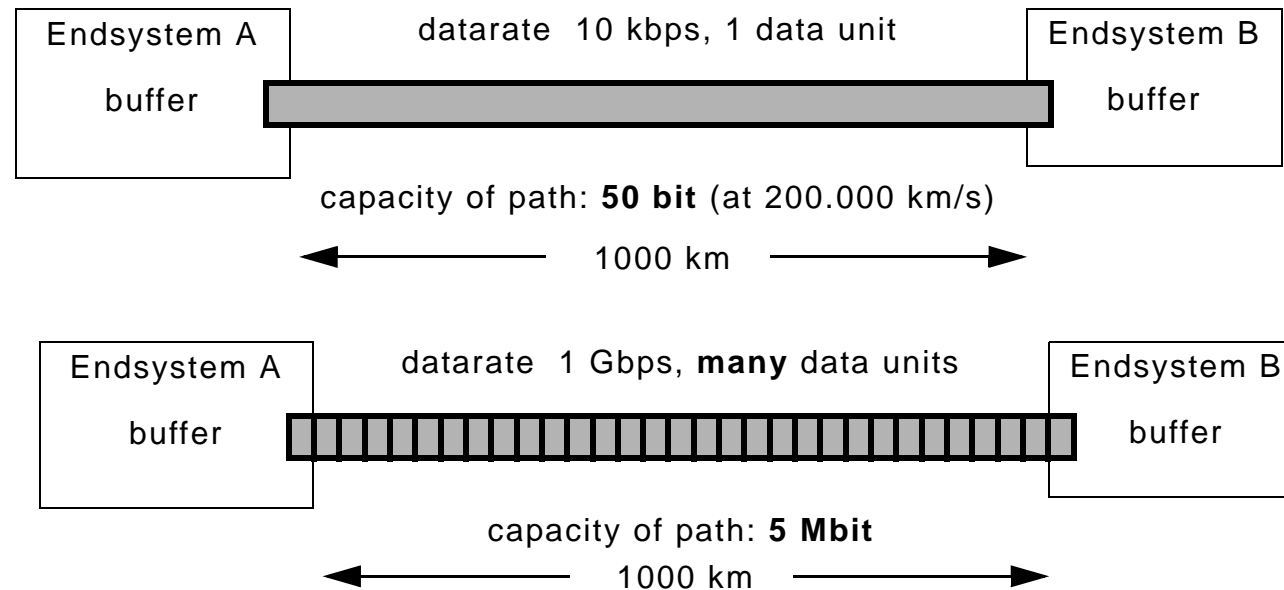
Motivation

- networks and applications have changed

Networks

- higher data rates
- also farther distances (e.g. also via satellite)
- networks for data storage

$$\text{Data amount} = \frac{\text{Data rate} \times \text{Distance}}{\text{Velocity of Propagation}}$$



Bandwidth-Delay Product increases

$$\text{BANDWIDTH [BITS/SEC]} * \text{ROUND-TRIP DELAY [SEC]}$$



Applications

- **broadcast and group communication**
- **services and service transitions**
- **varieties of media**
 - text, graphics, ..
 - animation, audio, video, ..