# Algorithms for Distributed Storage and Computer Forensics
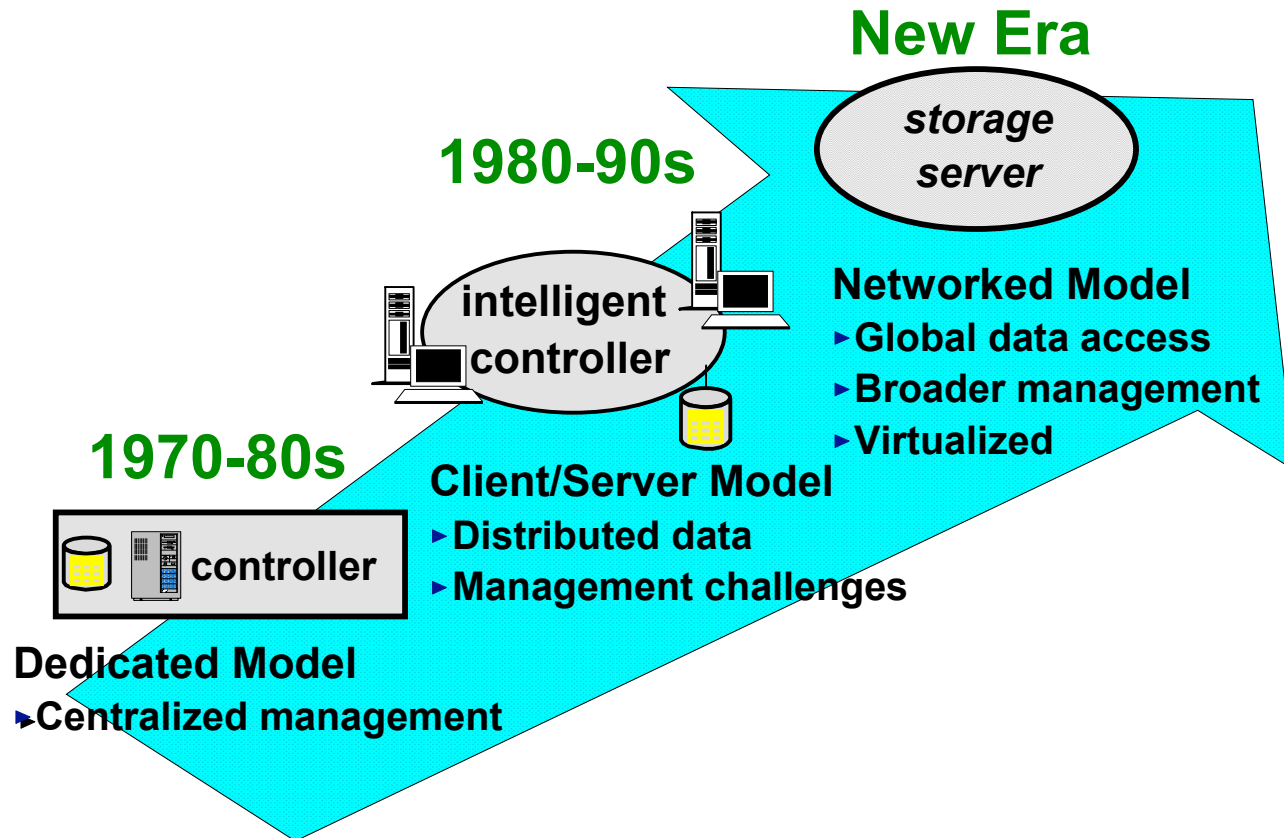
## 11 Networking
### Christian Schindelhauer

University of Freiburg
Technical Faculty
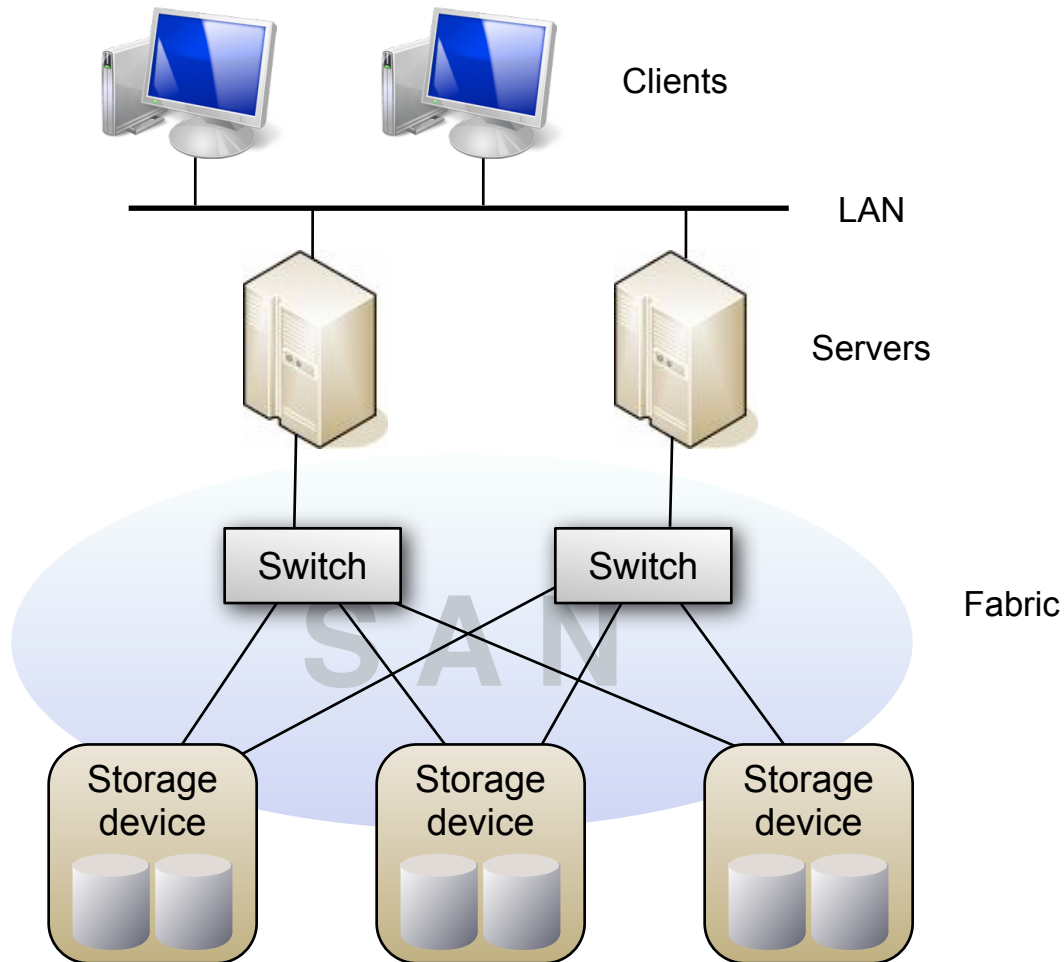Computer Networks and Telematics
Winter Semester 2011/12

**CoNe Freiburg**

**IIF**
INSTITUT FÜR
INFORMATIK
FREIBURG

# Evolution of Storage



**New Era**

*storage server*

**1980-90s**

intelligent controller

**Networked Model**
- ▶ Global data access
- ▶ Broader management
- ▶ Virtualized

**1970-80s**

controller

**Client/Server Model**
- ▶ Distributed data
- ▶ Management challenges

**Dedicated Model**
- ▶ Centralized management

[Tate, Lucchese, Moore: Introduction to Storage Area Networks, IBM 2006]

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

2

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Storage Area Network

Clients

LAN

Servers

Switch    Switch

**S A N**    Fabric

Storage device    Storage device    Storage device

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

3

Computer Networks and Telematics
University of Freiburg
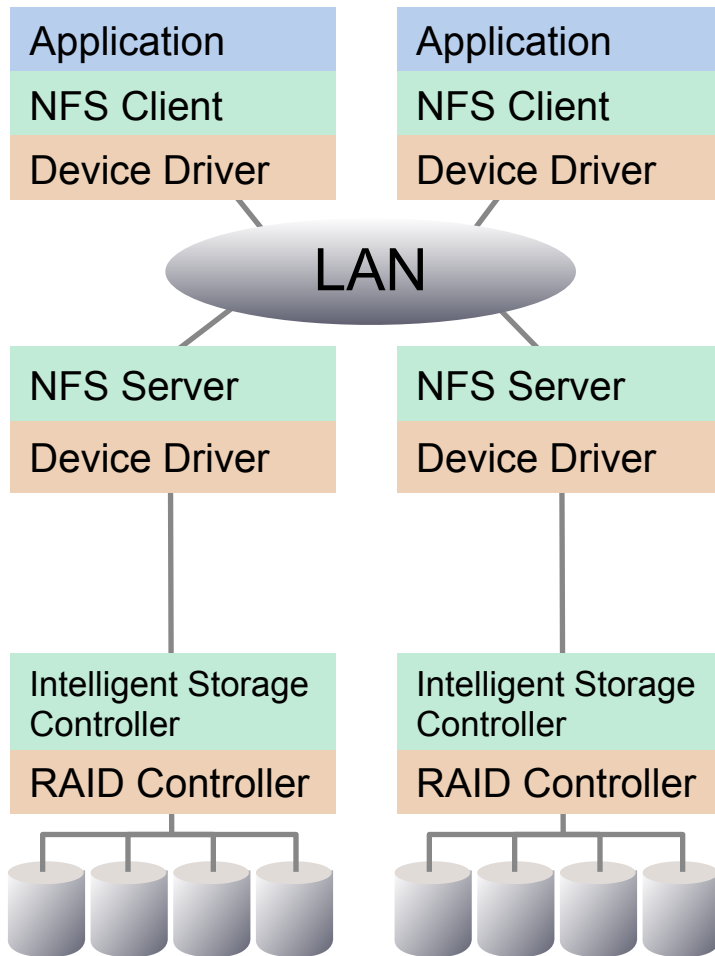Christian Schindelhauer

Mittwoch, 8. Februar 12

# NAS and SAN

‣ **Network-attached Storage**

- storage device attached to a network

- access through NFS, AFS, SMB, etc. (file level)

‣ **Storage Area Network**

- storage system of interconnected storage devices
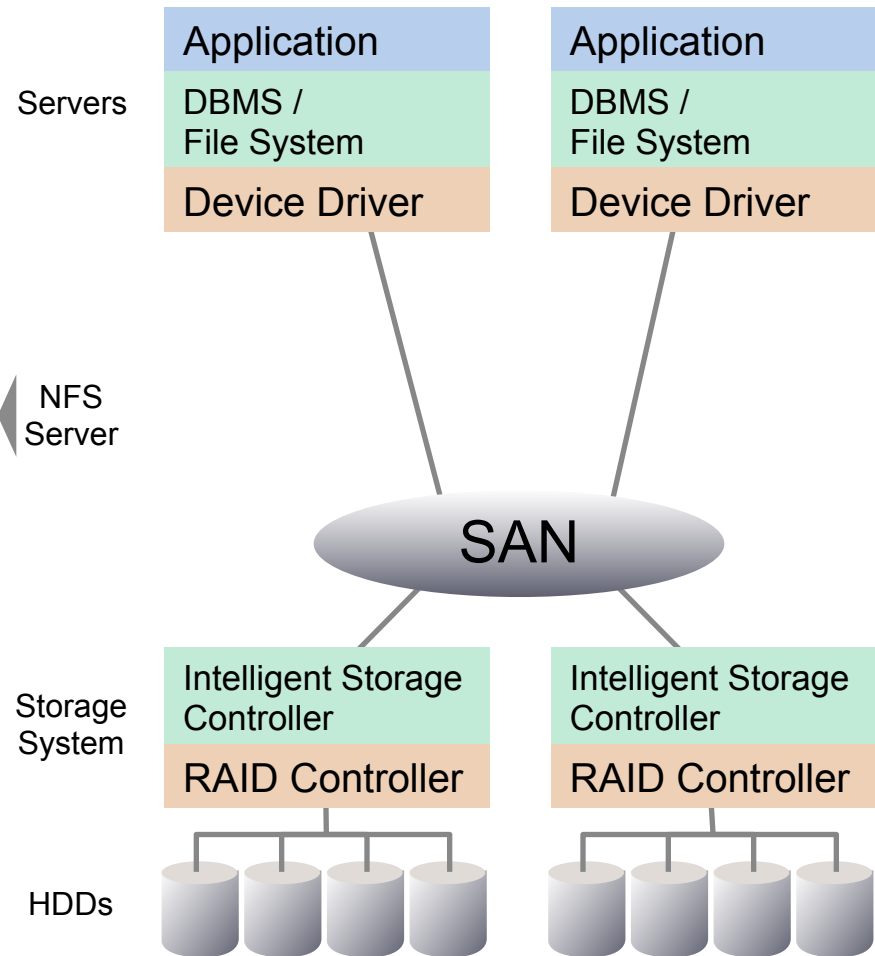
- access through FCP, iFCP, iSCSI (block level)

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

4

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# NAS and SAN

**Network Attached Storage**

| Application |
|---|
| NFS Client |
| Device Driver |

| Application |
|---|
| NFS Client |
| Device Driver |

Servers

LAN

| NFS Server |
|---|
| Device Driver |

| NFS Server |
|---|
| Device Driver |

◄ NFS Server

| Intelligent Storage Controller |
|---|
| RAID Controller |

| Intelligent Storage Controller |
|---|
| RAID Controller |

Storage System

HDDs

**Storage Area Network**

| Application |
|---|
| DBMS / File System |
| Device Driver |

| Application |
|---|
| DBMS / File System |
| Device Driver |

SAN

| Intelligent Storage Controller |
|---|
| RAID Controller |

| Intelligent Storage Controller |
|---|
| RAID Controller |

[Morris, Truskowski: The evolution of storage systems, IBM Systems Journal, 42(2), 2003]
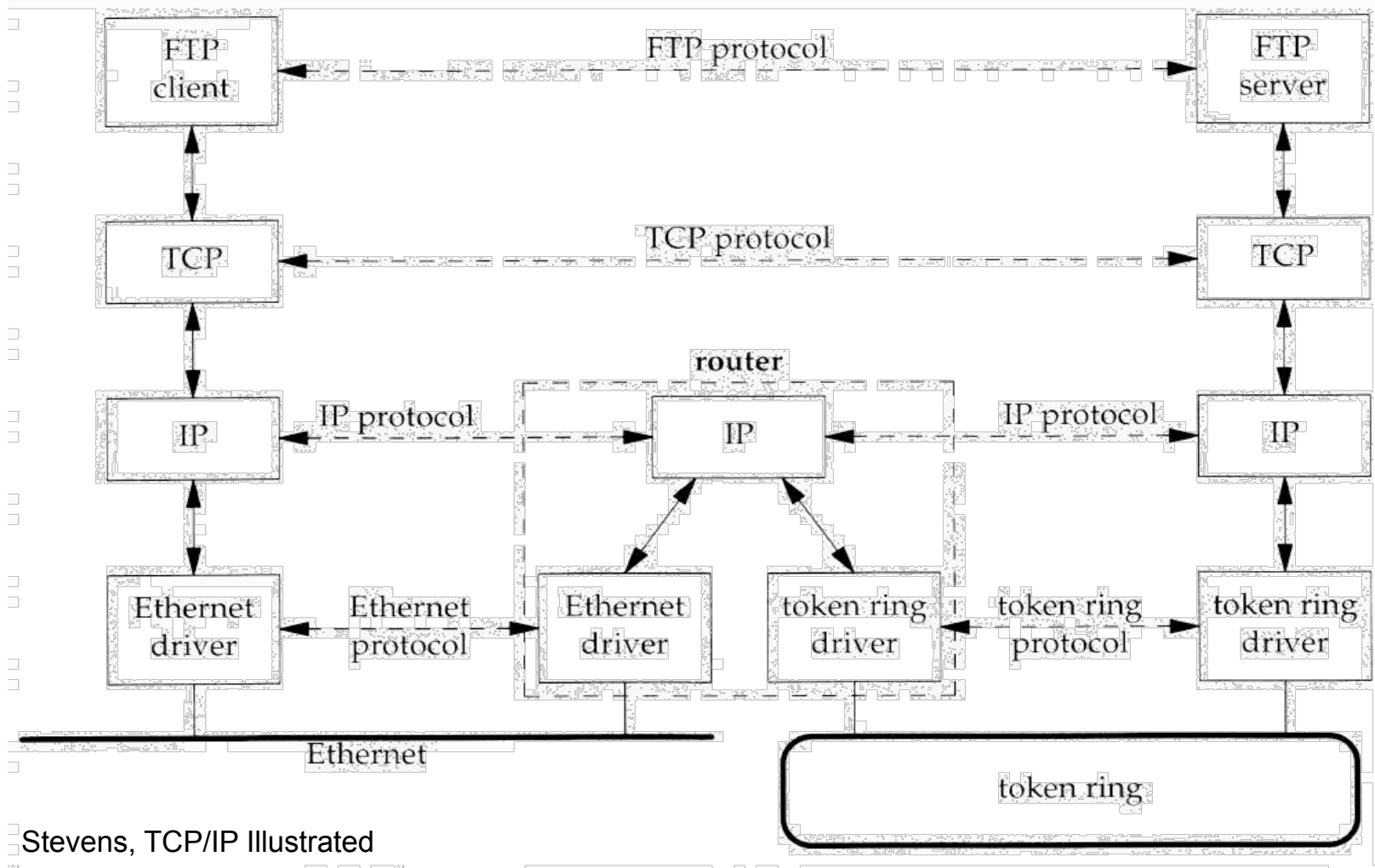
# Internet:
# An Open Network Architecture

▸ **Concept of Robert Kahn (DARPA 1972)**

- Local networks are autonomous

  - independent

  - no WAN configuration

- **packet-based** communication

- "**best effort**" communication

  - if a packet cannot reach the destination, it will be deleted

  - the application will re-transmit

- black-box approach to connections

  - black boxes: gateways and routers

Distributed Storage Networks
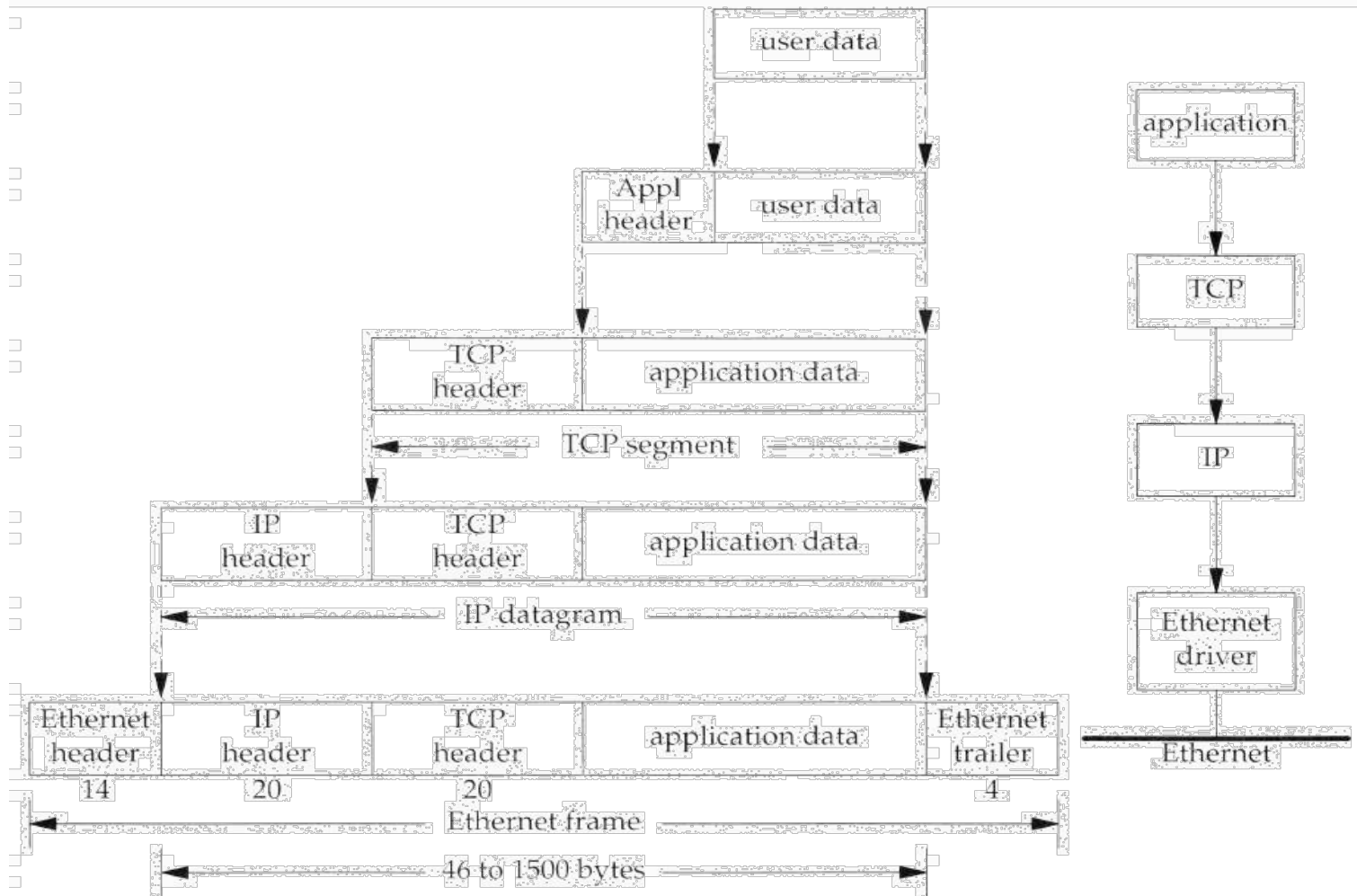and Computer Forensics
Winter 2011/12

6

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Protocols of the Internet

| | |
|---|---|
| Application | Telnet, FTP, HTTP, SMTP (E-Mail), ... |
| Transport | TCP (Transmission Control Protocol)<br><br>UDP (User Datagram Protocol) |
| Network | IP **(Internet Protocol)**<br>+ ICMP **(Internet Control Message Protocol)**<br>+ IGMP **(Internet Group Management Protoccol)** |
| Host-to-Network | LAN **(e.g. Ethernet, Token Ring etc.)** |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

7

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Example: Routing between LANs



Stevens, TCP/IP Illustrated

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

8

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Data/Packet Encapsulation



Stevens, TCP/IP Illustrated

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

9

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer
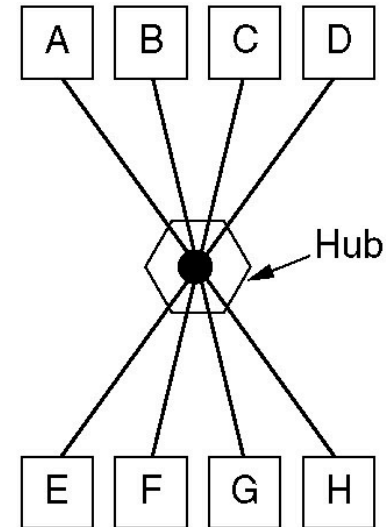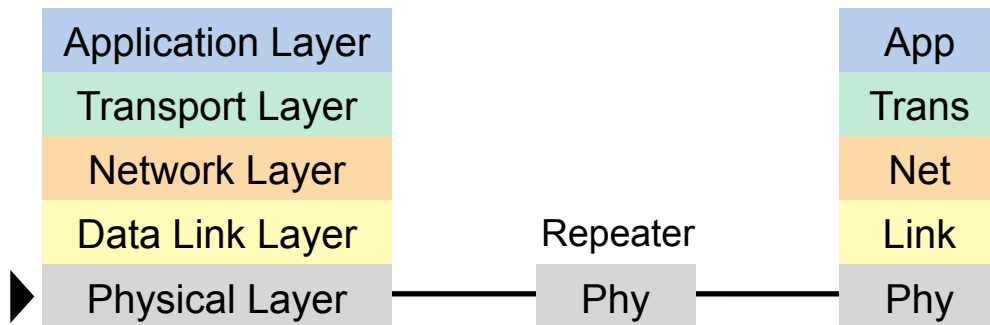
# Network Interconnections

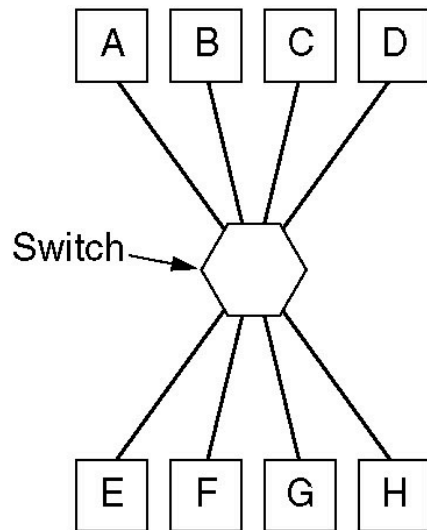| | |
|---|---|
| Application layer | Application gateway |
| Transport layer | Transport gateway |
| Network layer | Router |
| Data link layer | Bridge, switch |
| Physical layer | Repeater, hub |



[Tanenbaum, Computer Networks]

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

10

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Repeater and Hub

‣ **Receives, aplifies, re-transmits**

- only on the signal level

- Information remains untouched

| A | B | C | D |

Hub

| E | F | G | H |

| | |
|---|---|
| Application Layer | App |
| Transport Layer | Trans |
| Network Layer | Net |
| Data Link Layer | Link |
| Physical Layer | Phy |

Repeater — Phy

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

11

Computer Networks and Telematics
University of Freiburg
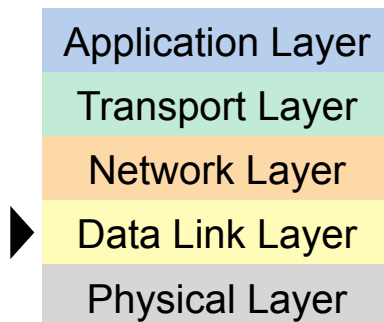Christian Schindelhauer

Mittwoch, 8. Februar 12

# Switch



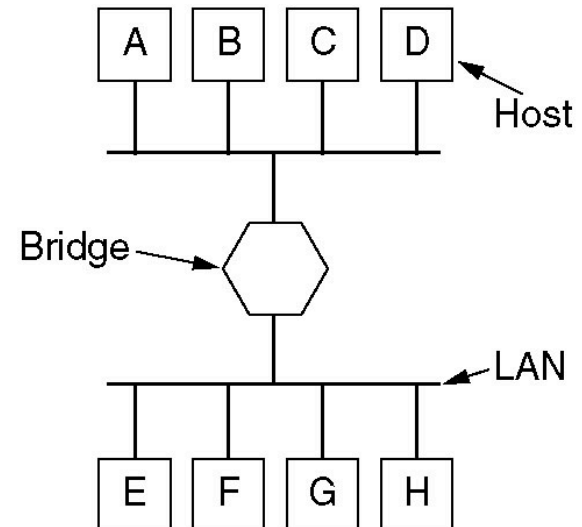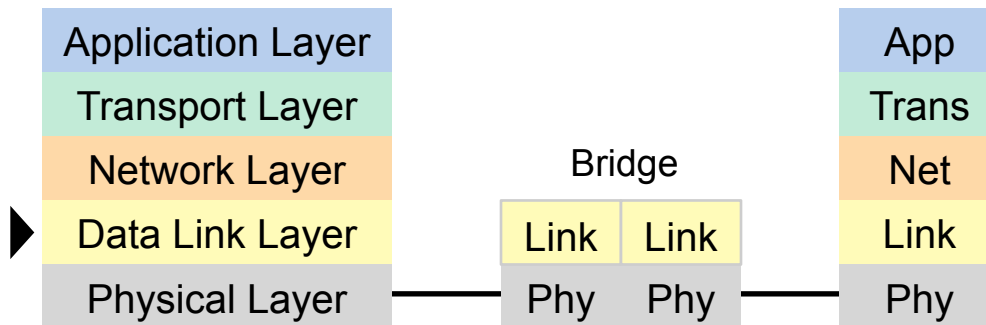‣ **Connection of multiple network segments**

- frames are forwarded only to the target segment

- collisions are not repeated

- store & forward (w. error correction)

- cut through switching: forwarding starts after the header is read

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

12

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Bridge

‣ **Connection of two network segments**

- different access methods
- multiport bridge similar to switch

| Application Layer | | | | App |
|---|---|---|---|---|
| Transport Layer | | | | Trans |
| Network Layer | | Bridge | | Net |
| Data Link Layer | Link | Link | | Link |
| Physical Layer | Phy | Phy | | Phy |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

13

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Routing

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

14

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Why do we need a network layer?

‣ **Local Networks can be connected by hubs, switches, bridges**
  - Problems:
    - Hubs propagate collisions
    - Switching: Inefficient collection of routing information
    - Problem of broadcasting
    - Internet connects >> 10 Mio. local networks

‣ **In large networks, routing information becomes necessary**
  - How is it collected?
  - How are packets forwarded?

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

15

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Routing Tables and Packet Forwarding

‣ **IP Routing Table**

- contains for each destination the address of the next gateway
- destination: host computer or sub-network
- default gateway

‣ **Packet Forwarding**

- IP packet (datagram) contains start IP address and destination IP address
    - if destination = my address then hand over to higher layer
    - if destination in routing table then forward packet to corresponding gateway
    - if destination IP subnet in routing table then forward packet to corresponding gateway
    - otherwise, use the default gateway

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

16

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Routing Table (Distance Vector)



Router

A    B    C    D

F    G

E         H

I    J    K    L

(a)

[Tanenbaum, Computer Networks]

New estimated delay from J

| To | A | I | H | K | | Line |
|----|----|----|----|----|----|----|
| A | 0 | 24 | 20 | 21 | 8 | A |
| B | 12 | 36 | 31 | 28 | 20 | A |
| C | 25 | 18 | 19 | 36 | 28 | I |
| D | 40 | 27 | 8 | 24 | 20 | H |
| E | 14 | 7 | 30 | 22 | 17 | I |
| F | 23 | 20 | 19 | 40 | 30 | I |
| G | 18 | 31 | 6 | 31 | 18 | H |
| H | 17 | 20 | 0 | 19 | 12 | H |
| I | 21 | 0 | 14 | 22 | 10 | I |
| J | 9 | 11 | 7 | 10 | 0 | – |
| K | 24 | 22 | 22 | 0 | 6 | K |
| L | 29 | 33 | 9 | 9 | 15 | K |
| | JA delay is 8 | JI delay is 10 | JH delay is 12 | JK delay is 6 | New routing table for J | |

Vectors received from J's four neighbors

(b)

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

17

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# IPv4 Packet Header



|←——————————————— 32 Bits ———————————————→|

| Version | IHL | Type of service | | Total length | | |
|---------|-----|-----------------|--|--------------|--|--|
| Identification | | | | DF | MF | Fragment offset |
| Time to live | | Protocol | | Header checksum | | |
| Source address | | | | | | |
| Destination address | | | | | | |
| Options (0 or more words) | | | | | | |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

18

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# IP Packet Forwarding

‣ **IP -Paket (datagram) contains...**

• TTL (Time-to-Live): Hop count limit

• Start IP Address

• Destination IP Address

‣ **Packet Handling**

• Reduce TTL (Time to Live) by 1

• If TTL ≠ 0  then forward packet according to routing table

• If TTL = 0 or forwarding error (buffer full etc.):

    - delete packet

    - if packet is not an ICMP Packet then

        ∗ sende ICMP Packet with

            · start = current IP Address

            · destination = original start IP Address

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

19

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Static and Dynamic Routing
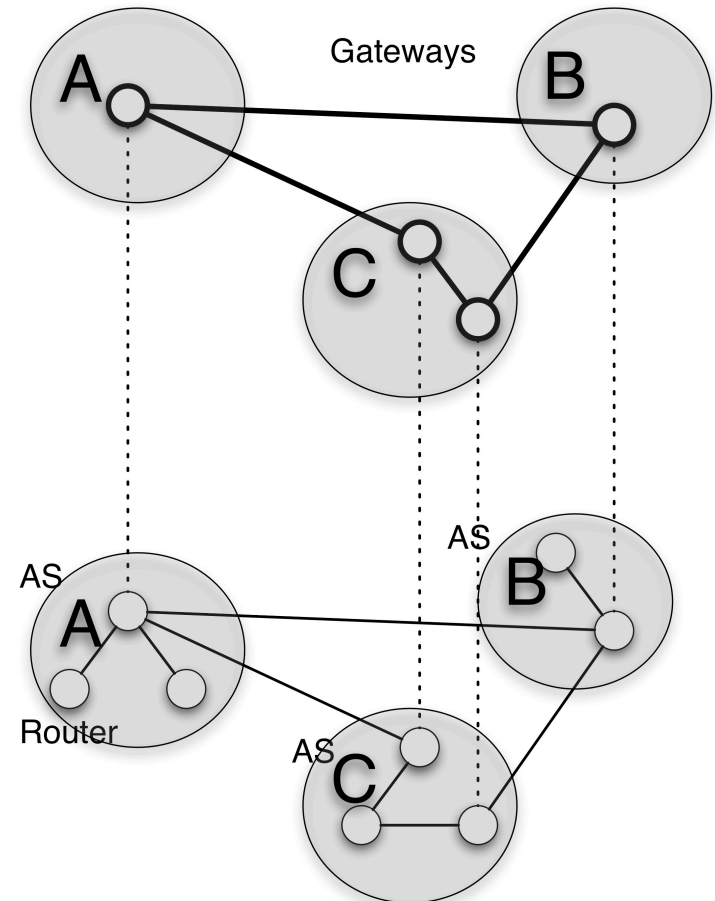
‣ **Static Routing**
  - Routing table created manually
  - used in small LANs

‣ **Dynamic Routing**
  - Routing table created by Routing Algorithm
  - **Centralized**, e.g. Link State
    - Router knows the complete network topology
  - **Decentralized**, e.g. Distance Vector
    - Router knows gateways in its local neighborhood

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

20

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Hierarchical Routing

‣ **Internet consists of Autonomous Systems (AS)**

- example: uni-freiburg.de

‣ **Intra-AS-Routing (Interior Gateway Protocol)**

- z.B. RIP, OSPF, IGRP, ...

‣ **Inter-AS-Routing (Exterior Gateway Protocol)**

- between Gateways
- decentralized
- everybody can define a metric
- z.B. BGP

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

21

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Intra-AS Routing

‣ **Inter-AS**
- Routing Information Protocol (RIP)
  - Distance Vector Algorithmus
  - Metric = hop count
  - exchange of distance vectors (by UDP)
- Interior Gateway Routing Protocol (IGRP)
  - successor of RIP
  - different routing metrics (delay, bandwidth)
- Open Shortest Path First (OSPF)
  - Link State Routing (every router knows the topology)
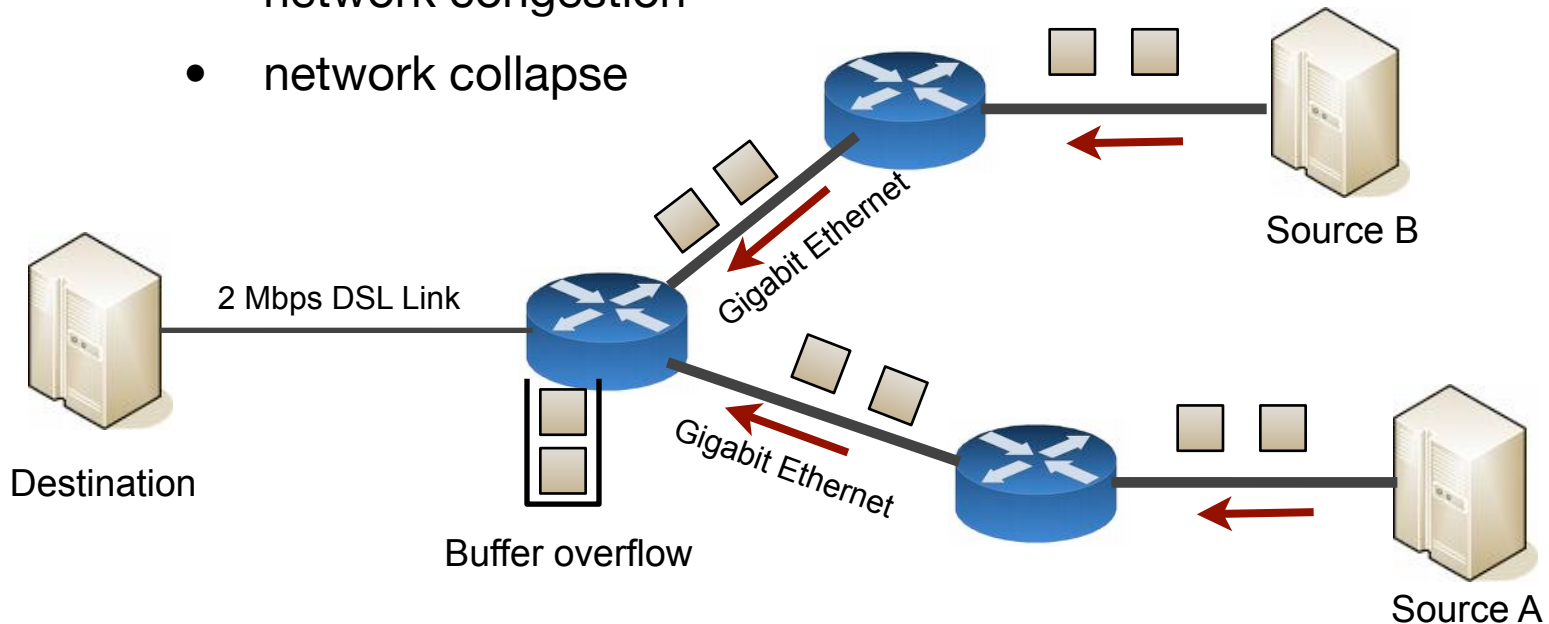  - Route calculation by Dijkstra's shortest path algorithm

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

22

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Inter-AS Routing

‣ **Problems of Inter-AS Routing**

- AS may reject packets

- Political consideration: Routing through other contries?

- Routing metrics of different AS are not compatible

  - path optimization impossible

  - Inter-AS Routing tries to achieve reachability

- Currently, Inter-Domain Router know more than 140.000 Networks
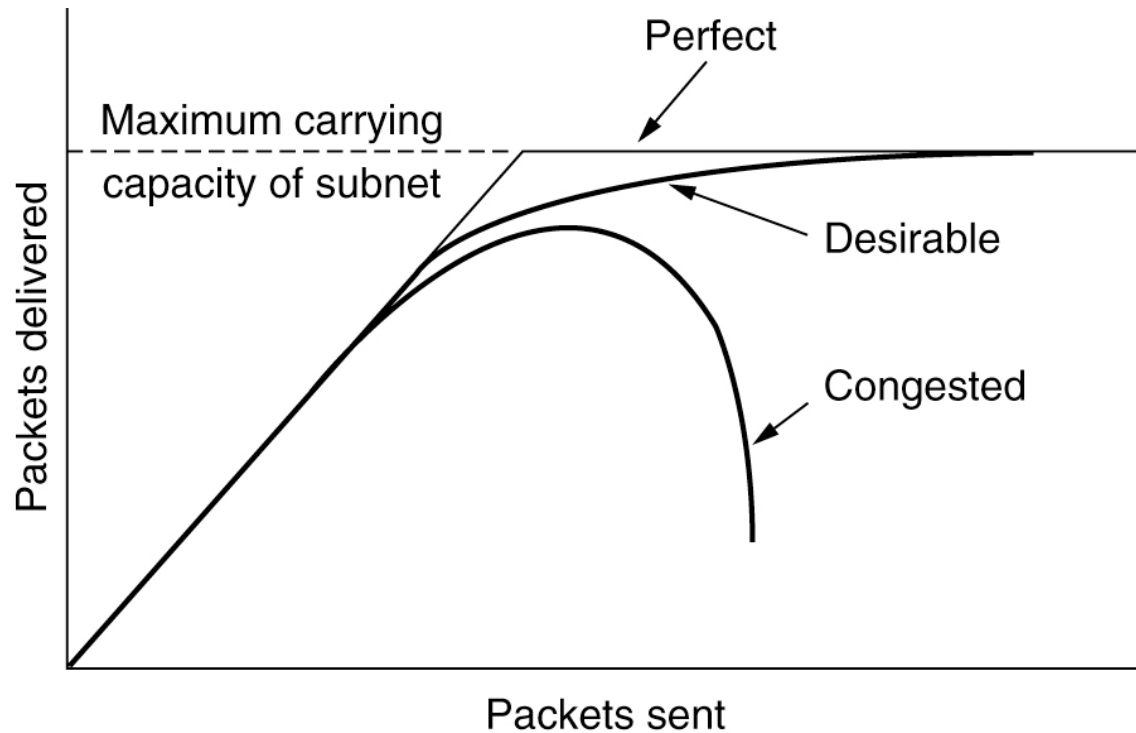

‣ **Border Gateway Protocol (BGP)**

- Path-Vector Protocol

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

23

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Network Congestion

‣ **(Sub-)Networks have limited bandwidth**

‣ **Injecting too many packets leads to**

- network congestion

- network collapse

2 Mbps DSL Link

Gigabit Ethernet

Gigabit Ethernet

Source B

Source A

Destination

Buffer overflow

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

24

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Congestion and capacity

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

25

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer
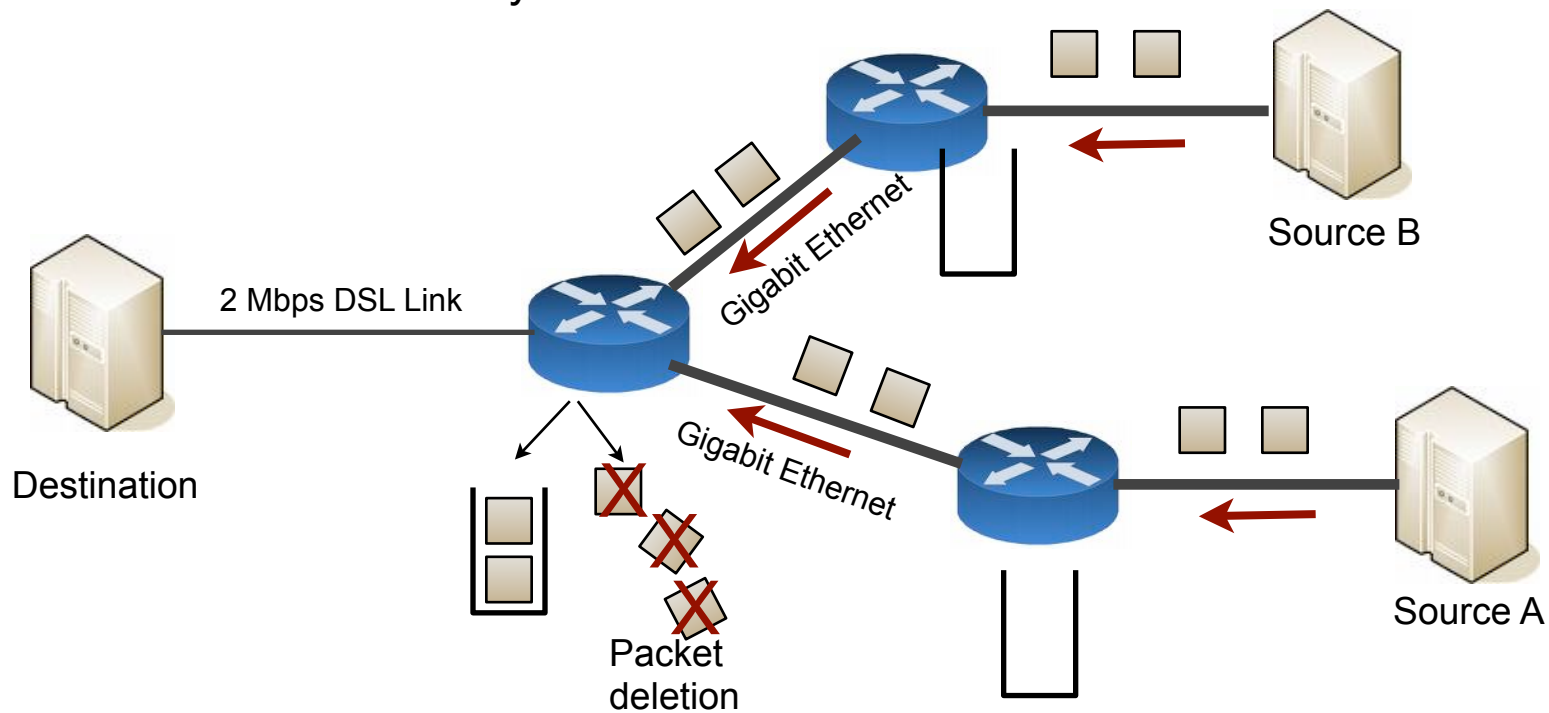
# Congestion Prevention

| Layer | Policies |
|---|---|
| Transport | • Retransmission policy<br>• Out-of-order caching policy<br>• Acknowledgement policy<br>• Flow control policy<br>• Timeout determination |
| Network | • Virtual circuits versus datagram inside the subnet<br>• Packet queueing and service policy<br>• Packet discard policy<br>• Routing algorithm<br>• Packet lifetime management |
| Data link | • Retransmission policy<br>• Out-of-order caching policy<br>• Acknowledgement policy<br>• Flow control policy |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

26

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Congestion Prevention by Routers

‣ **IP Routers drop packets**

- Tail dropping

- Random Early Detection

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

27

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# The Transport Layer
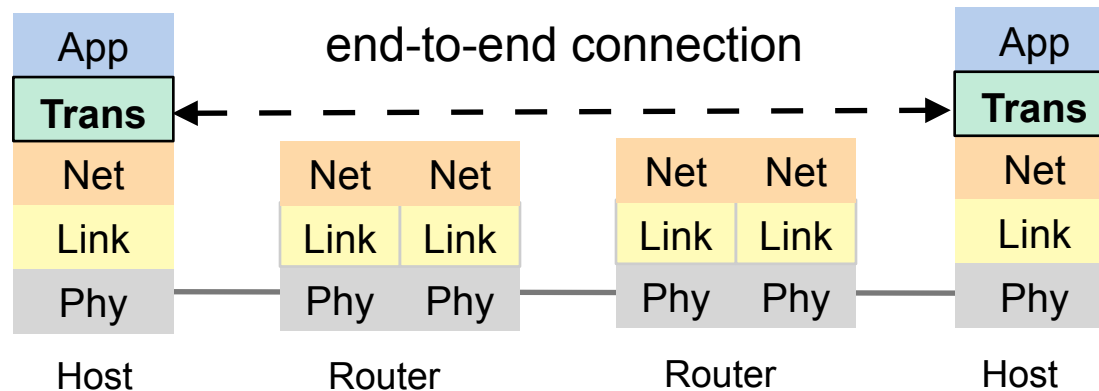
‣ **TCP (Transmission Control Protocol**

- connection-oriented
- delivers a stream of bytes
- reliable and ordered

‣ **UDP (User Datagram Protocol)**

- delivery of datagrams
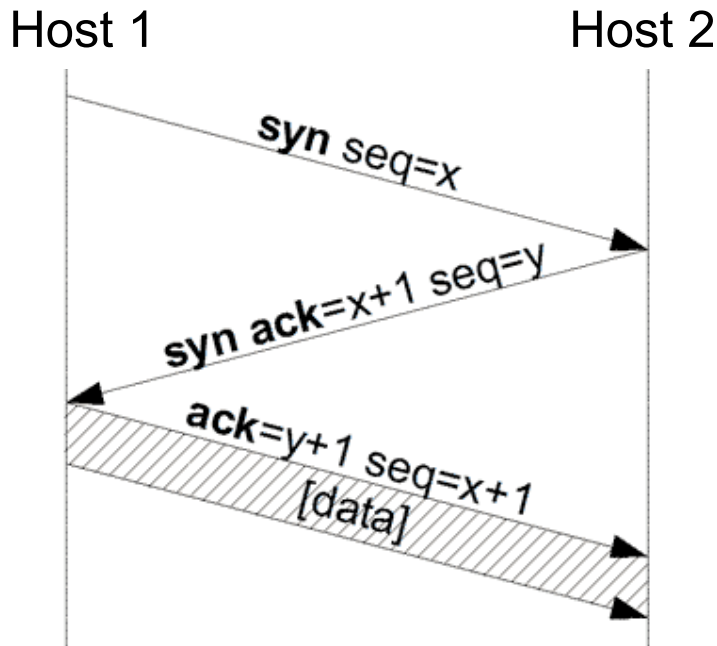- connectionless, unreliable, unordered

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

28

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# The Transmission Control Protocol (TCP)

‣ Connection-oriented

‣ Reliable delivery of a byte stream

  • fragmentation and reassembly (*TCP segments*)

  • acknowledgements and retransmission

‣ In-order delivery, duplicate detection

  • sequence numbers

‣ Flow control and congestion control

  • window-based (receiver window, congestion window)

‣ **challenge**: IP (network layer) packets can be dropped, delayed, delivered out-of-order ...

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

29

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# TCP Connections

**Connection establishment**                **Connection termination**

Host 1                  Host 2        Host 1                  Host 2

**syn** seq=x

**syn ack**=x+1 seq=y

**ack**=y+1 seq=x+1
[data]

**fin**, seq=x

**ack**=x+1

**fin**, seq=y

**ack**=y+1

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

30

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Flow control and congestion control



[Tanenbaum, Computer Networks]

(a)

(b)

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

31

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Flow Control

## acknowledgements and window management

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

32

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

# Retransmissions

‣ Retransmissions are triggered, if acknowledgements do not arrive

... but how to decide that?

‣ Measurement of the **round trip time (RTT)**

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

33

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Retransmissions and RTT

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

34

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Congestion revisited

‣ **IP Routers drop packets**

‣ **TCP has to react, e.g. lower the packet injection rate**



TCP

2 Mbps DSL Link

Gigabit Ethernet

Gigabit Ethernet

Destination

Source B

Source A

TCP

Packet
deletion

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

35

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Congestion revisited

| App | | Net | | Net | Net | | App |
|-----|---|-----|---|-----|-----|---|-----|
| **Trans** | | | | | | | **Trans** |
| Net | | **Congestion!** | | | | | Net |
| Link | | Link | | Link | Link | | Link |
| Phy | | Phy | | Phy | Phy | | Phy |
| Host | | Router | | Router | | | Host |

**from a transport layer perspective:**

App   ?   ?   ?   App

Trans ← – – – – – – – – – – – – – Trans    no ACKs
                                            received

Net   Net   Net   Net   Net   Net

Link  Link  Link  Link  Link  Link

Phy   Phy   Phy   Phy   Phy   Phy

Host   Router   Router   Host

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

36

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

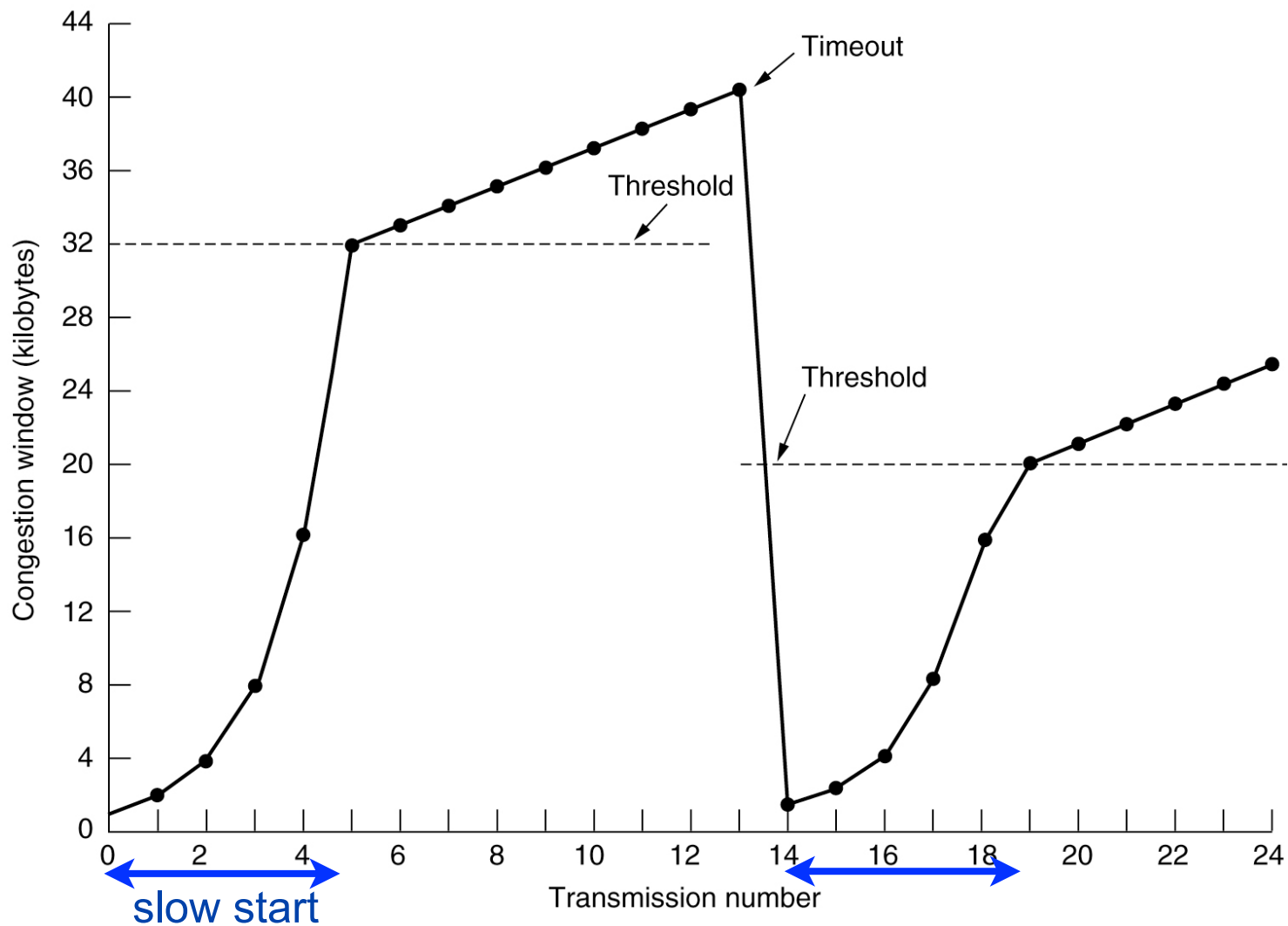# Data rate adaption and the congestion window

‣ **Sender does not use the maximum segment size in the beginning**

‣ **Congestion window (cwnd)**
  - used on the sender size
  - sending window: min {wnd,cwnd} (wnd = receiver window)
  - S: segment size
  - Initialization:
    - cwnd ← S
  - For each received acknowledgement:
    - cwnd ← cwnd + S
  - ...until a packet remains unacknowledged

**Sender** → **Receiver**

Segment 1
ACK: Segment 1

Segment 2
Segment 3
ACK: Segment 3

Segment 4
Segment 5
Segment 6
Segment 7
ACK: Segment 5
ACK: Segment 7

Segment 8
Segment 9
Segment 10
⋮

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

37

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Slow Start of TCP Tahoe

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

38

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# The AIMD principle

‣ **TCP uses basically the following mechanism
to adapt the data rate x (#packets sent per RTT):**
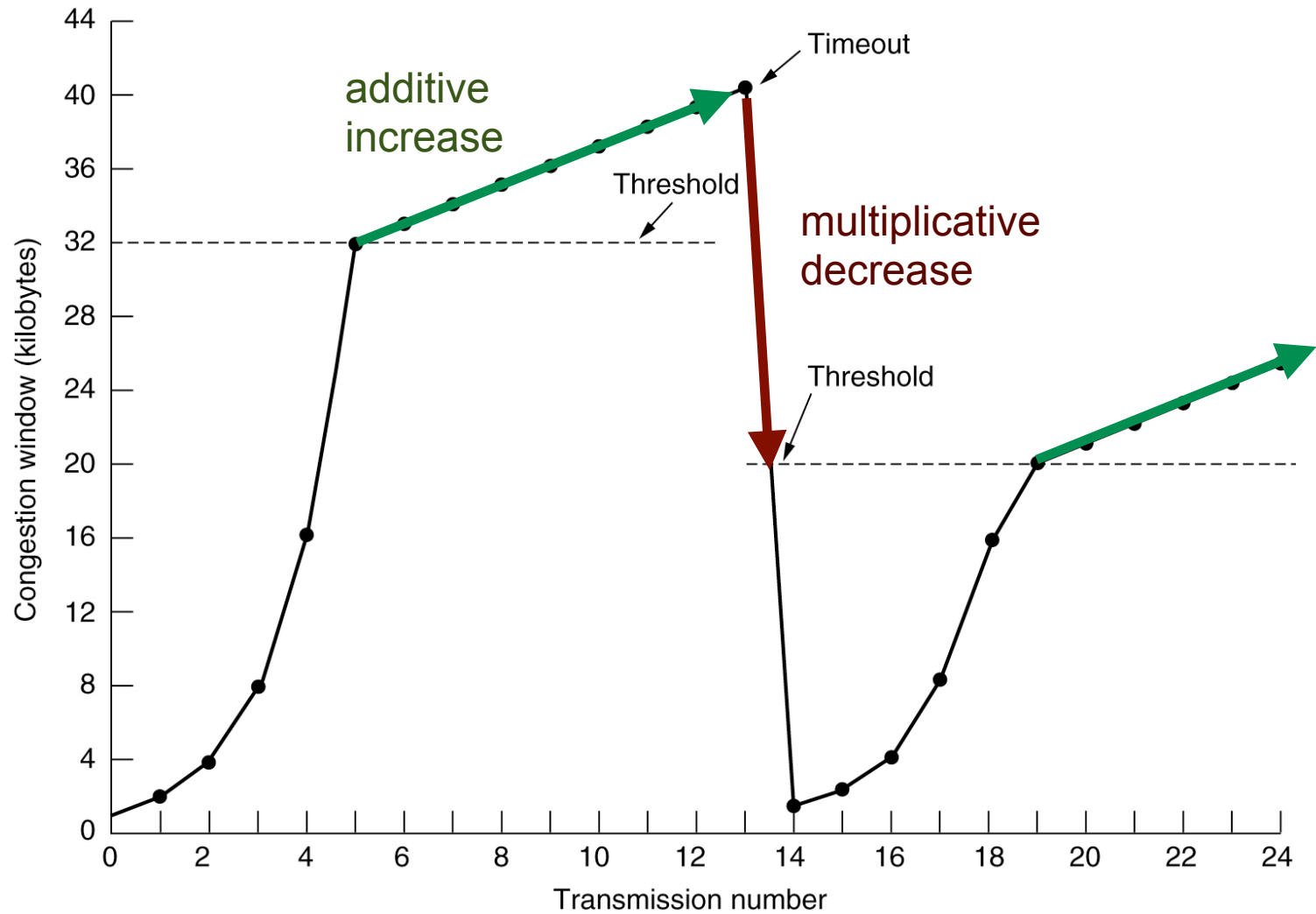
- Initialization:

$$x \leftarrow 1$$

- on packet loss: multiplicative decrease (MD)

$$x \leftarrow x/2$$

- if the acknowledgement for a segment arrives, perform additive increase (AI)

$$x \leftarrow x + 1$$

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

39

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# AIMD

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

40

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Example of TCP Reno



Fast Retransmit

Fast Recovery

Slow Start

Additive Increase

Multiplicative Decrease

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

41

Computer Networks and Telematics
University of Freiburg
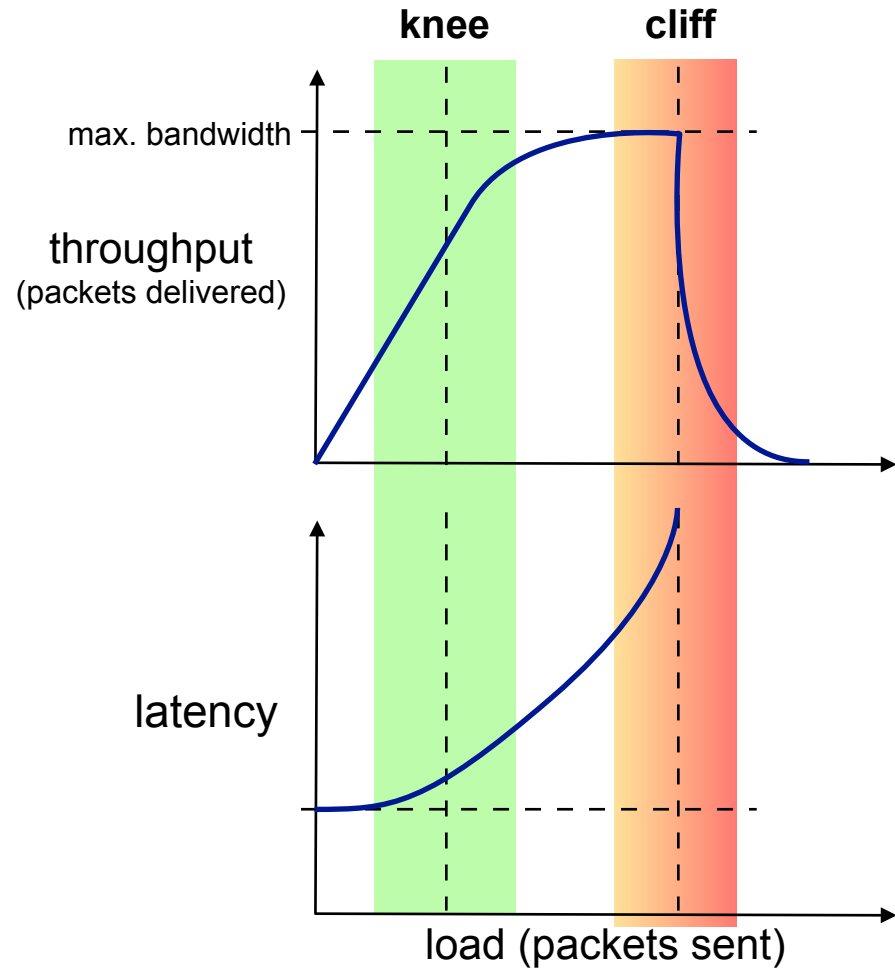Christian Schindelhauer

Mittwoch, 8. Februar 12

# Throughput and Latency

‣ **Congested situation (cliff):**

- high load

- low throughput

- all data packets are lost

‣ **Desired situation (knee):**

- high load

- high throughput

- few data packets get lost

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

42

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# TCP vs. UDP

‣ **TCP reduces data rate**

‣ **UDP does not!**

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

43

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

# TCP - Conclusion

‣ Connection-oriented, reliable,
  in-order delivery of a byte stream

‣ Flow control and congestion control

  • Fairness among TCP streams

  • Unfair behavior of other protocols, e.g. UDP

  • Impact on latency

  • Tweaking the congestion avoidance mechanism has an
    impact on other applications

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

44

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Storage networking

‣ **Fibre Channel**

- standard connection for SANs

- Medium: fibre-optic but also twisted pair

- Protocol: channel-like transport of SCSI commands

- Topologies: From point-to-point to networks

- Advantages: flexible connectivity, networking capabilities

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

45

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Fibre Channel Protocol (FCP)

‣ **Transport protocol for SCSI commands**

‣ **Layered architecture**

| | | |
|---|---|---|
| | **Application** | |
| **FC-4** | **Protocol Mapping Layer** (Multimedia / Channels / Networks) | Upper Layers |
| **FC-3** | **Common Services** | |
| **FC-2** | **Signaling Protocol** (Framing / Flow Control) | Physical Layers |
| **FC-1** | **Transmission Protocol** (Encode/Decode) | |
| **FC-0** | **Physical Interface / Media** | |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

46

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# FCP Layers

| FC4 | **Protocol Mapping Layer** | encapsulation of other protocols |
|-----|---------------------------|----------------------------------|
| FC3 | **Common Services** | encryption, striping, RAID, etc. |
| FC2 | **Framing and Signalling** | data transport, routing |
| FC1 | **Transmission Protocol** | 8b/10b encoding and decoding |
| FC0 | **Physical Layer** | medium |

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

47

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Fibre Channel Topologies
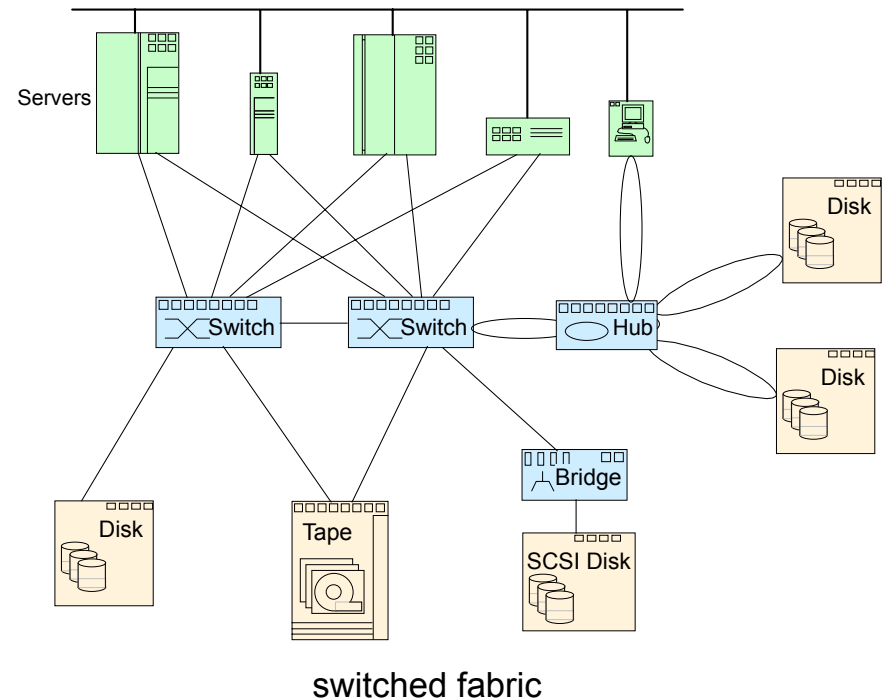
‣ **Point-to-Point**

  • connection of 2 nodes

‣ **Arbitrated Loop (FC-AL)**

  • shared bus of up to 126 nodes

‣ **Switched Fabric (FC-SW)**

  • interconnection network

  • routing and transport protocols



switched fabric

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

48

Computer Networks and Telematics
University of Freiburg
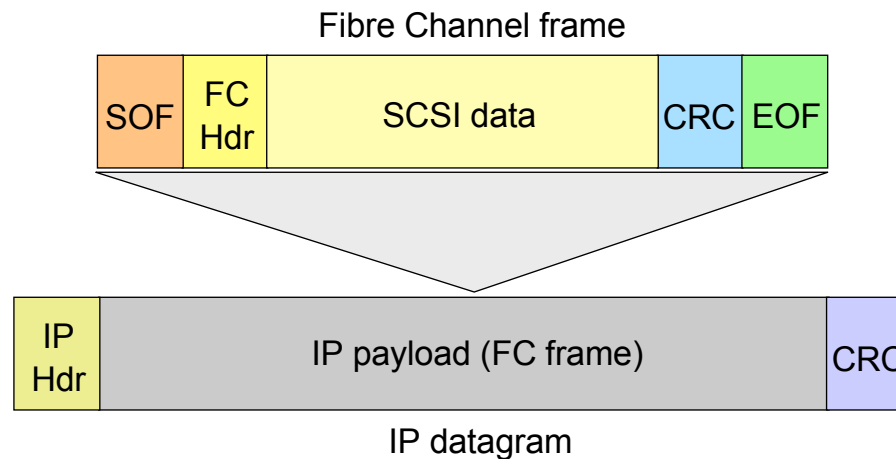Christian Schindelhauer

Mittwoch, 8. Februar 12

# Network Storage Types

‣ **Direct attached storage (DAS)**

  • traditional storage

‣ **Network attached storage (NAS)**

  • storage attached to another computer accessible at file level over LAN or WAN

‣ **Storage area network (SAN)**

  • specialized network providing other computers with storage capacity with access on block-addressing level

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

49

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# IP storage networking protocols
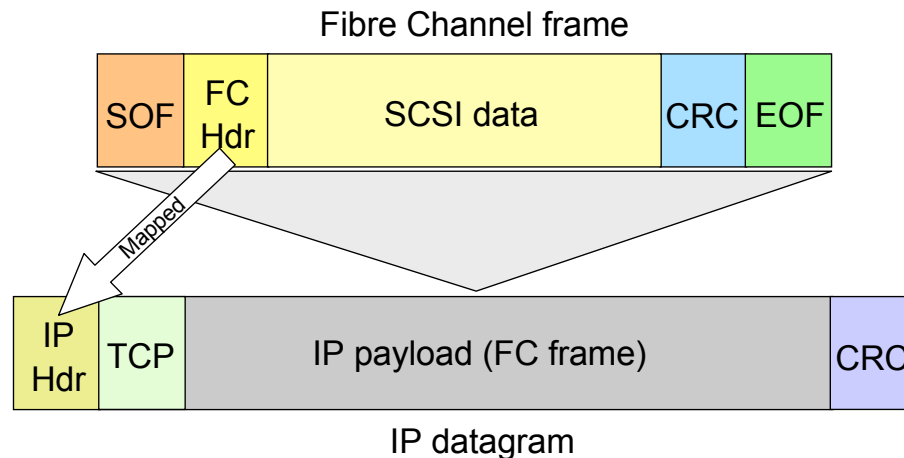
‣ **Fibre Channel over IP (FCIP)**

- Tunneling data between SAN devices through IP networks

- based on TCP connections

- links SAN devices and switch fabrics over IP networks

- Merging switch fabrics over IP links problematic
  (frequent switch reconfigurations because of link unreliability)

Fibre Channel frame

| SOF | FC Hdr | SCSI data | CRC | EOF |
|-----|--------|-----------|-----|-----|

| IP Hdr | IP payload (FC frame) | CRC |
|--------|----------------------|-----|

IP datagram

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

50

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# IP storage networking protocols

‣ **Internet Fibre Channel Protocol (iFCP)**

- Fibre Channel switch fabric services over IP networks

- based on TCP connections

- uses IP routing and switching

- can replace the Fibre Channel switch fabric

Fibre Channel frame

| SOF | FC Hdr | SCSI data | CRC | EOF |

Mapped

| IP Hdr | TCP | IP payload (FC frame) | CRC |

IP datagram

Distributed Storage Networks
and Computer Forensics
Winter 2011/12

51

Computer Networks and Telematics
University of Freiburg
Christian Schindelhauer

Mittwoch, 8. Februar 12

# Algorithms and Methods for Distributed Storage

## 6 Networking

Albert-Ludwigs-Universität Freiburg
Institut für Informatik
Rechnernetze und Telematik
Wintersemester 2008/09