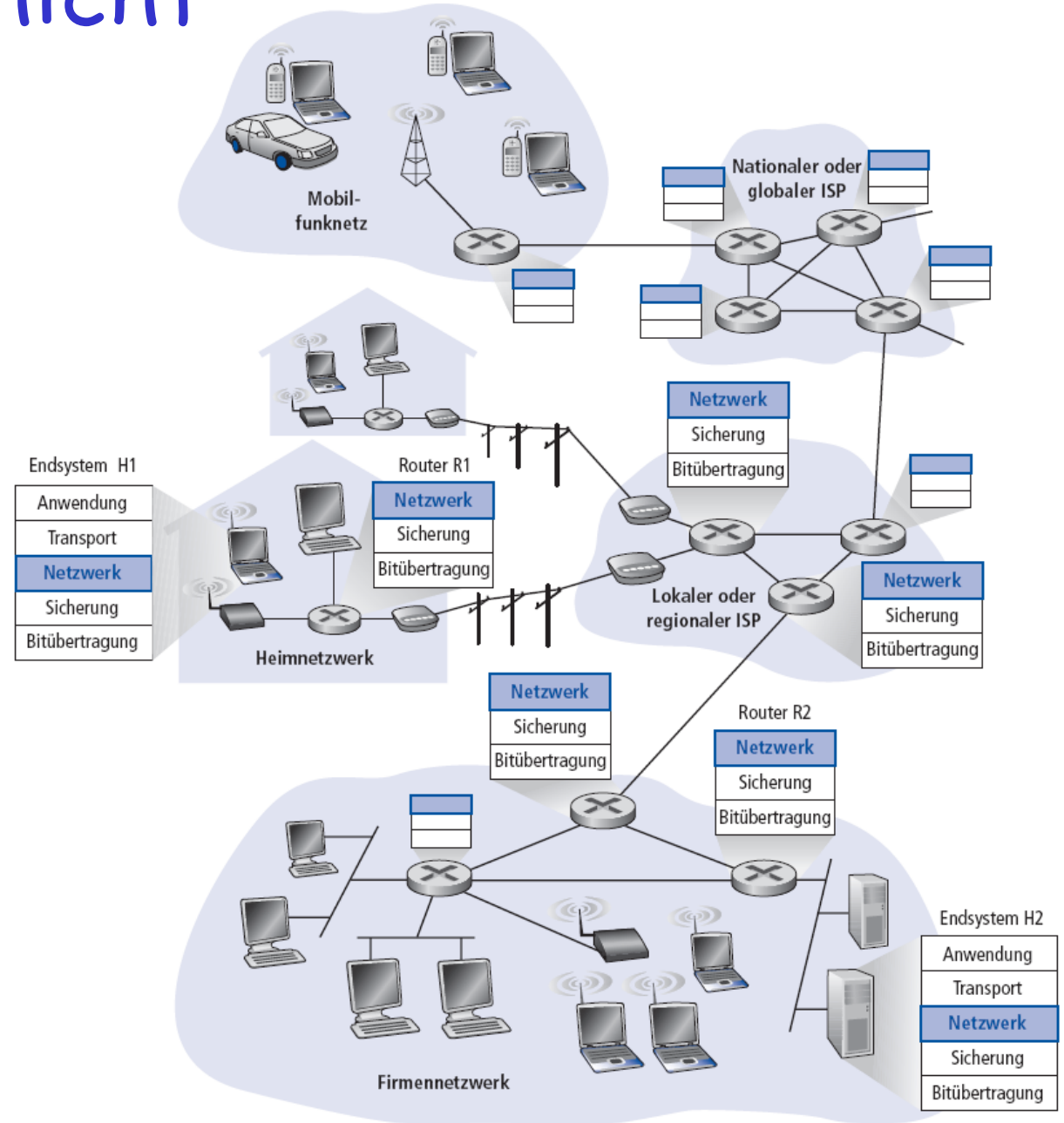


Kapitel 4:

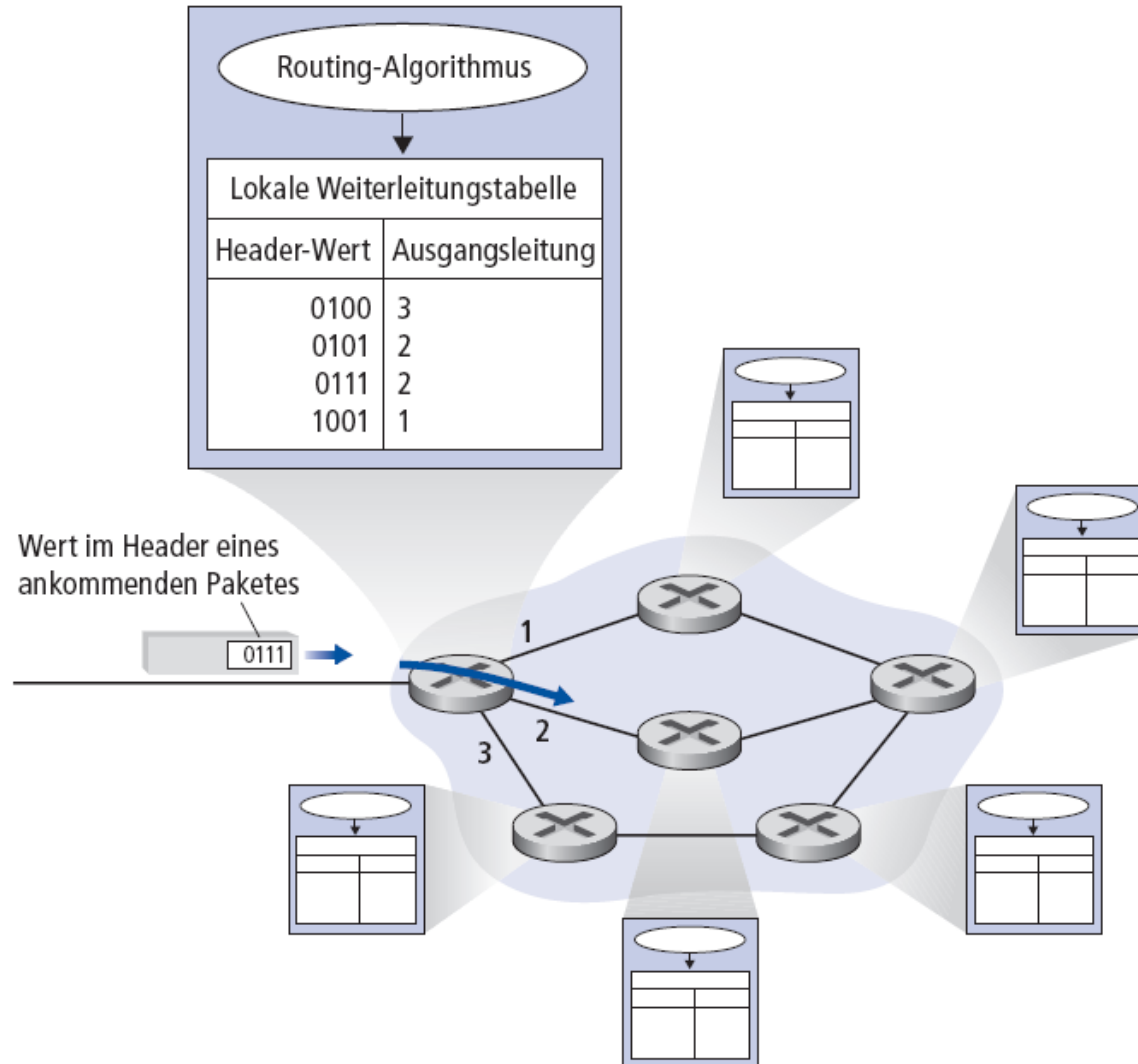
Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP - Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

Netzwerkschicht



Zusammenspiel von Routing und Forwarding



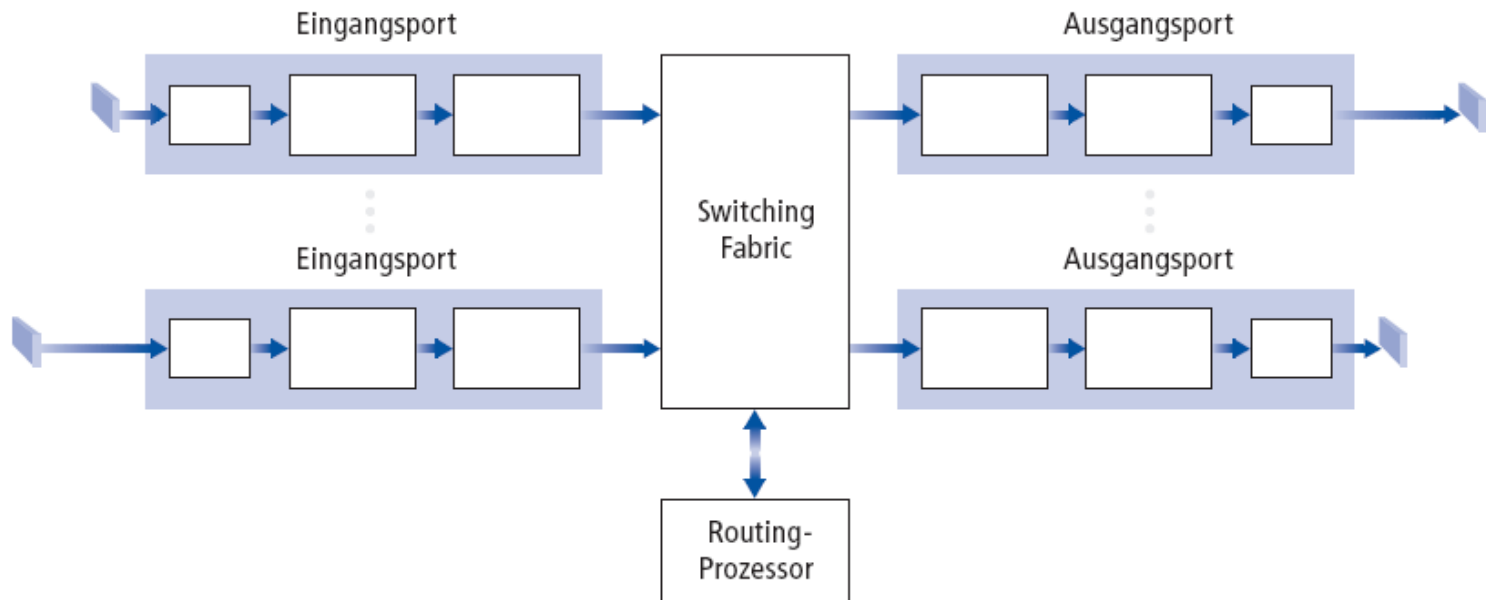
Kapitel 4:

Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

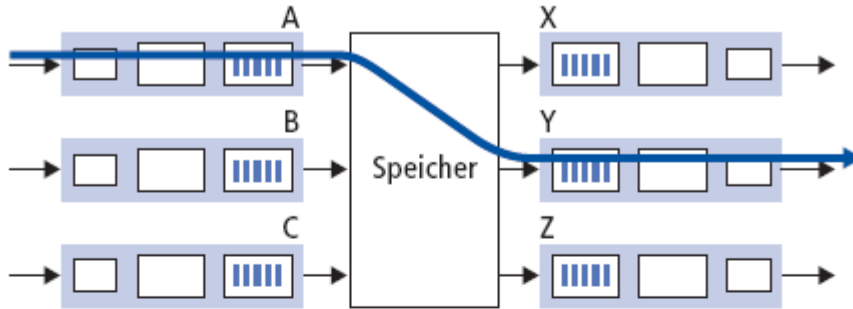
Übersicht: Routerarchitektur

- Leitungsabschluss: physikalische Schicht, Bits empfangen
- Sicherungsschicht: z.B. Ethernet (s. Kapitel 5)
- Nachschlagen, Weiterleiten, Queuing:
 - Suche nach einem geeigneten Ausgangs-port
 - Dezentral, Kopie der Routing-Tabelle (oder Teile davon) notwendig
 - Ziel: Behandlung der Pakete mit „line speed“, also mit der Geschwindigkeit der Eingangsleitung des Ports
 - Puffern von Paketen, wenn die Switching Fabric belegt ist

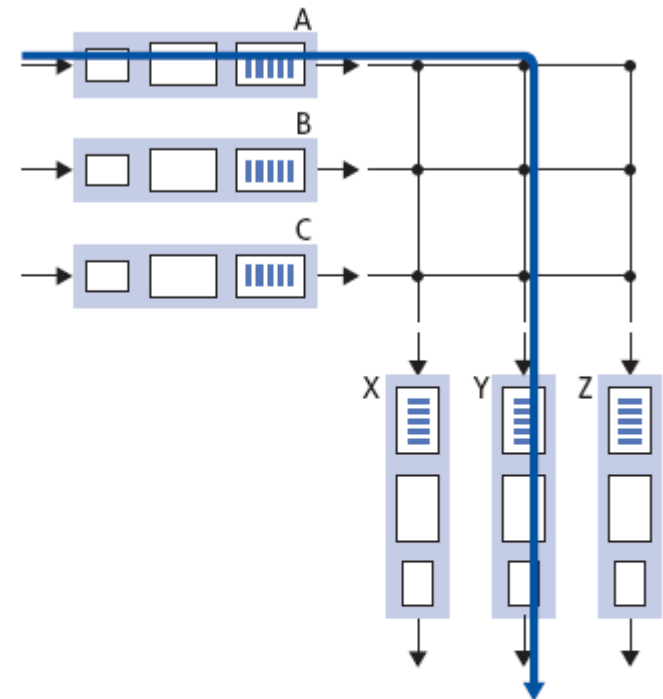


Drei verschiedene Arten von Switching Fabrics

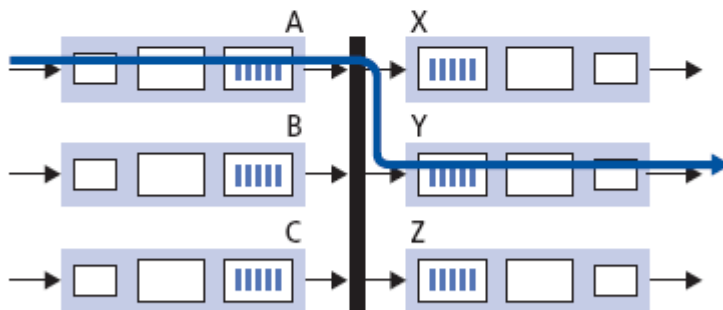
Speicher



Crossbar



Bus

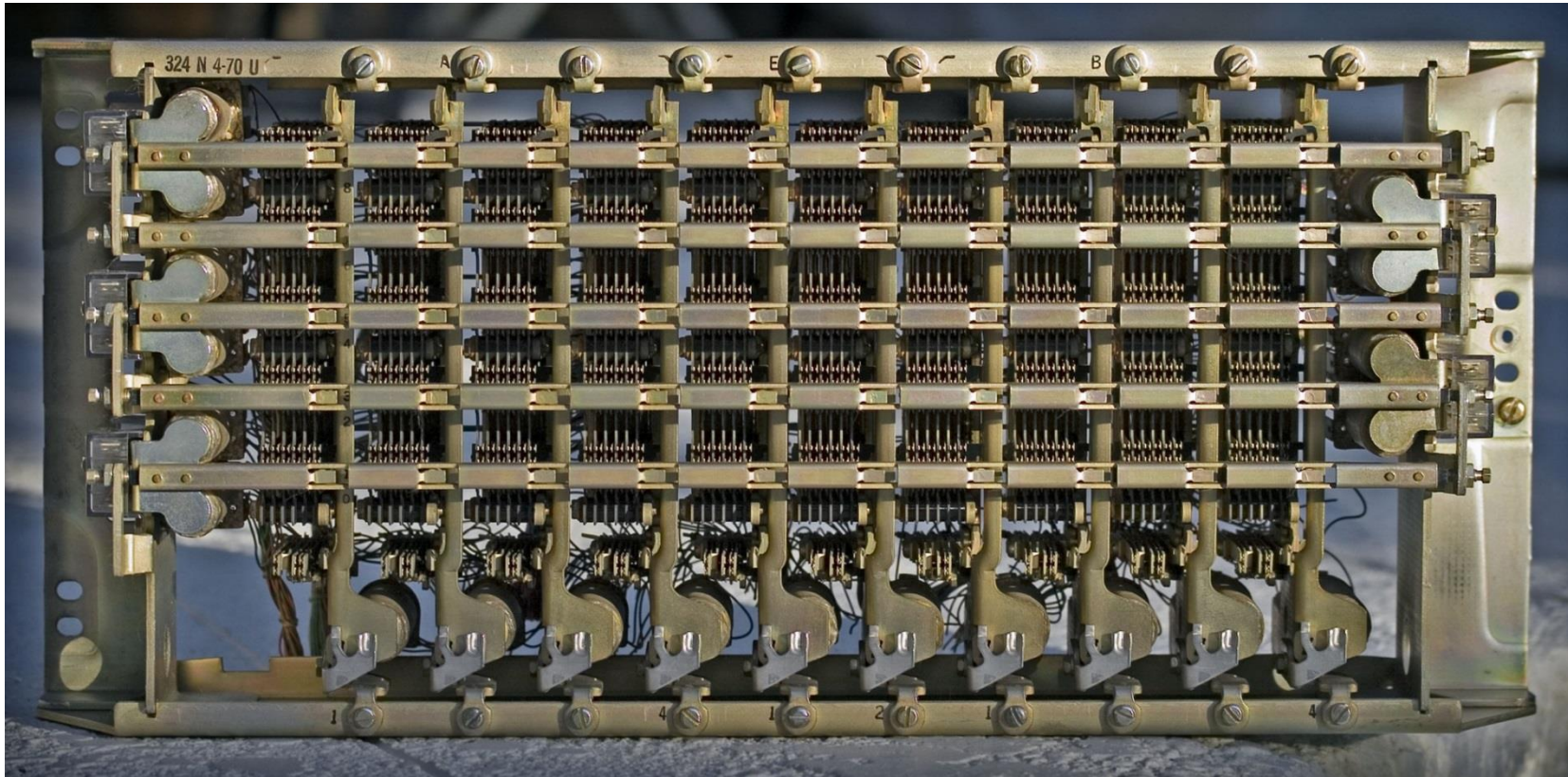


Legende:



Übersicht: Routerarchitektur

Crossbar - Switch: klassisch



Wikipedia: Western Electric 100-point six-wire Type B crossbar switch (1960)

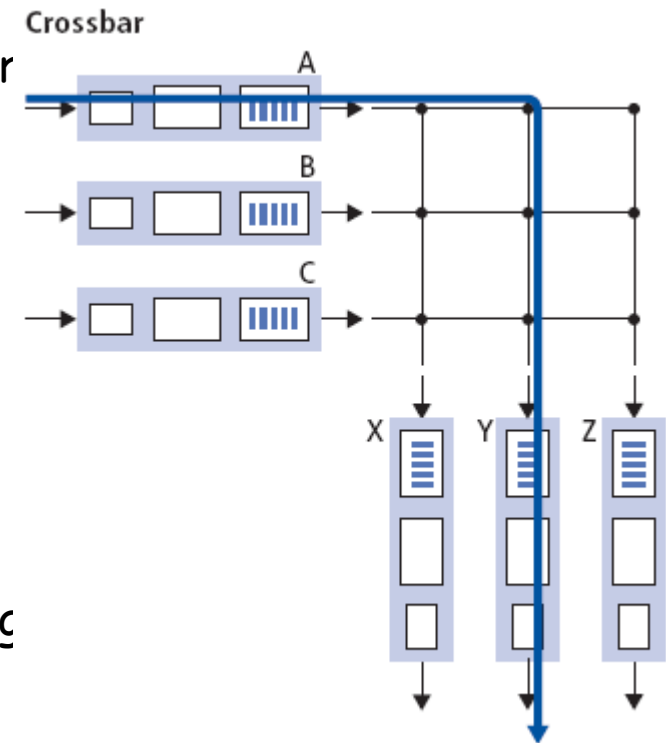
Switching über einen Bus

- Beispiel: 32-Gbps-Bus, Cisco 7600,
- ausreichend für Zugangsrouter und Router für Firmennetze (nicht geeignet im Backbone)



Switching über ein Spezialnetz

- Ports sind über ein Netzwerk miteinander verbunden
 - Beispielsweise alle Eingangsports über einen Crossbar mit allen Ausgangs-ports
 - Oder Banyan-Netzwerke
 - Technologie ursprünglich für das Verbinden mehrerer Prozessoren in einem Parallelrechner entwickelt
- Weitere Fortschritte: Zerlegen der Pakete in Zellen fester Größe, Zellen können dann schneller durch die Switching Fabric geleitet werden
- Beispiel: Cisco 12000, Switching von 60 Gbps durch das interne Netz



Switching über ein Spezialnetz

□ Beispiel: Cisco 12000,
Switching von 60 Gbps
durch das interne Netz



Switching IC



BCM88130 PRODUCT Brief

Fakultät für Feinwerk-
und Mikrotechnik,
Physikalische Technik

HOCHSCHULE
FÜR ANGEWANDTE
WISSENSCHAFTEN · FH
MÜNCHEN



HIGH-PERFORMANCE PACKET SWITCH FABRIC

FEATURES

- Scales linearly to over 10 Tbps
- Non-blocking architecture
- 100 GE ready
- 600 Gbps switched bandwidth in single device
- Central bandwidth management in single device
 - Globally managed Quality of Service (QoS)
 - Bandwidth guarantees and low latency/Jitter
 - Hierarchical VOQs with 16 COS
- High-speed 6.5 Gbps SerDes
- High-performance multicast
- Self-routing crossbar
- Reliability and availability features
 - 1+1 and load shared redundancy
 - Hardware based lossless switchover
 - Fault detection and correction
 - Graceful degradation
- Complete end-to-end application solutions

SUMMARY OF BENEFITS

- Proven fabric architecture for a range of modular platforms
- Interoperable with current QE2000 and future Queuing Engines to provide line-card future-proofing
- Central bandwidth management enforces service level agreements (SLAs) and bandwidth guarantees, including low latency and jitter services.
- Very low fabric overhead optimizes the use of backplane bandwidth.
- Line rate operation for all packet sizes with a full mesh of unicast and multicast traffic under stress
- High-performance multicast provides wirespeed non-blocking multicast while maintaining system SLAs. Streaming multicast/broadcast services such as IPTV requires these efficient multicast capabilities.
- Deep buffers provide a single control point for managing and guaranteeing QoS, and absorbing network round trip delays.
- Self-routing crossbar for chip-level autonomous operation allowing for single control point across devices and efficient multicast.

Kapitel 4:

Netzwerkschicht

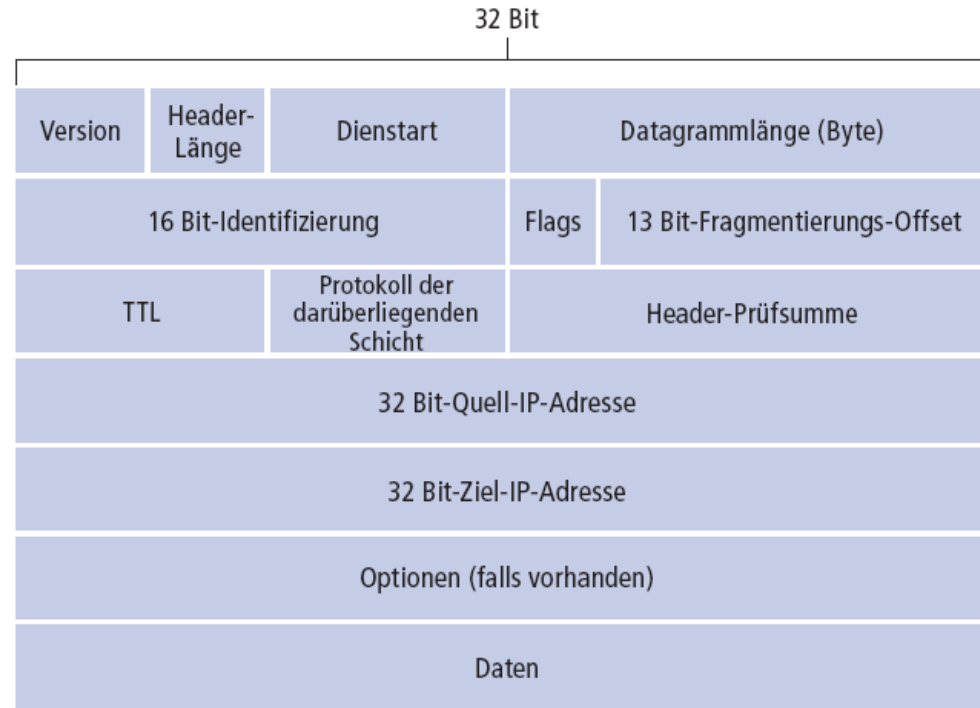
- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

Kapitel 4:

Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - **Datagrammformat**
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

IP-Datagrammformat

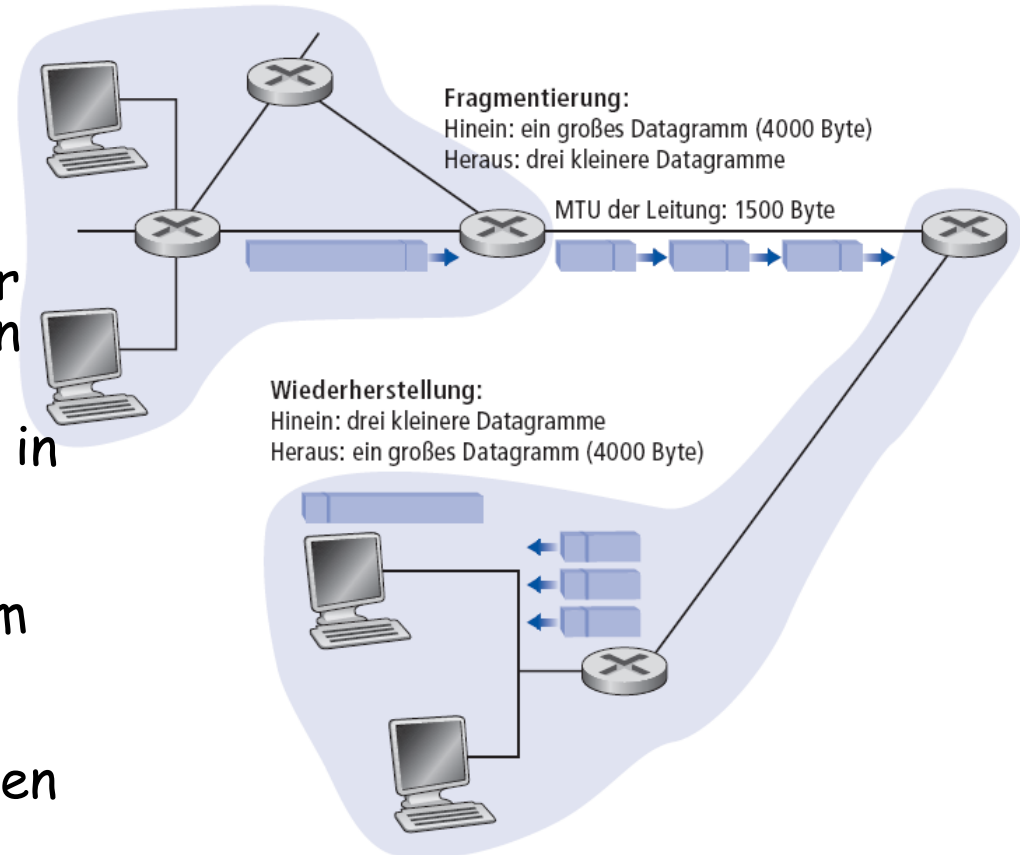


Wie viel Overhead entsteht bei Verwendung von TCP?

- 20 Byte für den TCP-Header, 20 Byte für den IP-Header
- = 40 Byte + Overhead auf der Anwendungsschicht

IP-Fragmentierung

- Links haben eine Maximalgröße für Rahmen
- Diese nennt man Maximum Transmission Unit (MTU)
- Verschiedene Links haben unterschiedliche MTUs
- IP-Datagramme müssen unter Umständen aufgeteilt werden
 - Aufteilung (Fragmentierung) erfolgt in den Routern
 - Zusammensetzen (Reassembly) erfolgt beim Empfänger
 - IP-Header enthält die notwendigen Informationen hierzu



IP-Fragmentierung

Beispiel: IP-Datagramm mit 4000 Byte (inklusive 20 Byte IP-Header),
MTU des nächsten Links = 1500 Byte

Fragment	Bytes	ID	Offset	Flag
1. Fragment	1.480 Byte im Datenfeld des IP-Datagramms	Identifizierung = 777	Offset = 0 (d.h., die Daten sollten beginnend bei Byte 0 eingefügt werden)	Flag = 1 (d.h., da kommt noch mehr)
2. Fragment	1.480 Datenbytes	Identifizierung = 777	Offset = 185 (d.h., die Daten sollten bei Byte 1.480 beginnend eingefügt werden; beachten Sie, dass $185 \cdot 8 = 1.480$)	Flag = 1 (d.h., da kommt noch mehr)
3. Fragment	1.020 Datenbytes (= $3.980 - 1.480 - 1.480$)	Identifizierung = 777	Offset = 370 (d.h., die Daten sollten beginnend bei Byte 2.960 eingefügt werden; beachten Sie, dass $370 \cdot 8 = 2.960$)	Flag = 0 (d.h., es ist das letzte Fragment)

IP-Fragmentierung

Beispiel: ping -l 1600 www.web.de

351	100.354998000	192.168.2.118	212.227.222.8	IPv4	1514	Fragmented IP protocol (proto=ICMP 1, off=0, ID=77d3) [Reassembled in #352]
352	100.355008000	192.168.2.118	212.227.222.8	ICMP	162	Echo (ping) request id=0x0001, seq=479/57089, ttl=128
353	100.402886000	212.227.222.8	192.168.2.118	IPv4	1506	Fragmented IP protocol (proto=ICMP 1, off=0, ID=b7ef) [Reassembled in #355]
354	100.403052000	212.227.222.8	192.168.2.118	IPv4	162	Fragmented IP protocol (proto=ICMP 1, off=1480, ID=b7ef) [Reassembled in #355]
355	100.403333000	212.227.222.8	192.168.2.118	ICMP	42	Echo (ping) reply id=0x0001, seq=479/57089, ttl=58
356	101.356895000	192.168.2.118	212.227.222.8	IPv4	1514	Fragmented IP protocol (proto=ICMP 1, off=0, ID=77d4) [Reassembled in #357]
357	101.356907000	192.168.2.118	212.227.222.8	ICMP	162	Echo (ping) request id=0x0001, seq=480/57345, ttl=128

Frame 352: 162 bytes on wire (1296 bits), 162 bytes captured (1296 bits) on interface 0

Ethernet II, Src: IntelCor_a3:85:1c (58:94:6b:a3:85:1c), Dst: AvM_e2:5b:8b (00:04:0e:e2:5b:8b)

Internet Protocol Version 4, Src: 192.168.2.118 (192.168.2.118), Dst: 212.227.222.8 (212.227.222.8)

- Version: 4
- Header length: 20 bytes
- Differentiated Services Field: 0x00 (DSCP 0x00: Default; ECN: 0x00: Not-ECT (Not ECN-Capable Transport))
- Total Length: 148
- Identification: 0x77d3 (30675)
- Flags: 0x00
- Fragment offset: 1480
- Time to live: 128
- Protocol: ICMP (1)
- Header checksum: 0x4bd2 [correct]
- Source: 192.168.2.118 (192.168.2.118)
- Destination: 212.227.222.8 (212.227.222.8)
- [Source GeoIP: Unknown]
- [Destination GeoIP: Unknown]
- [2 IPv4 Fragments (1608 bytes): #351(1480), #352(128)]

Internet Control Message Protocol

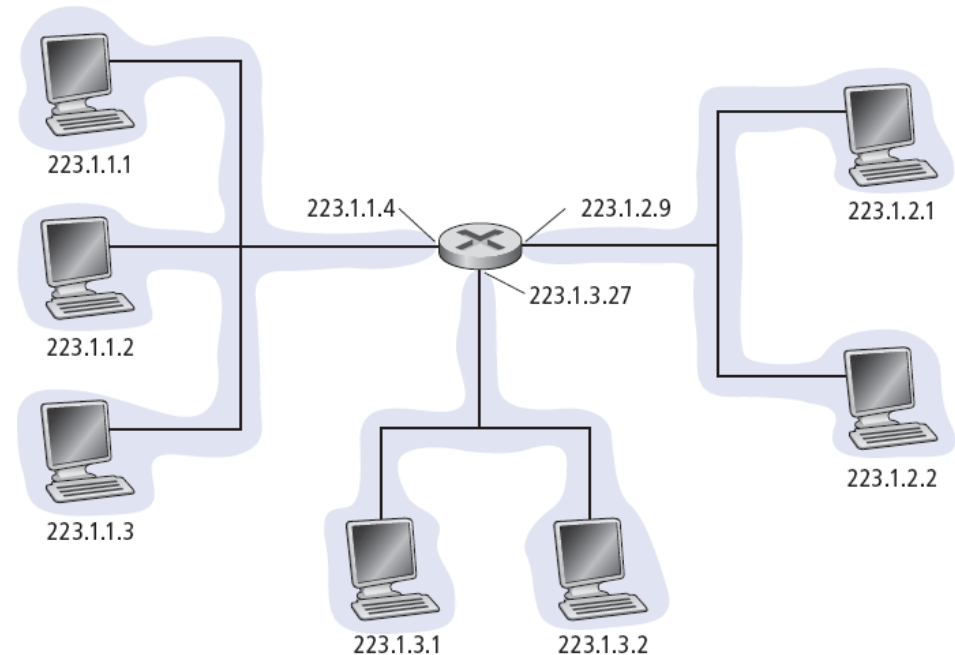
Kapitel 4:

Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

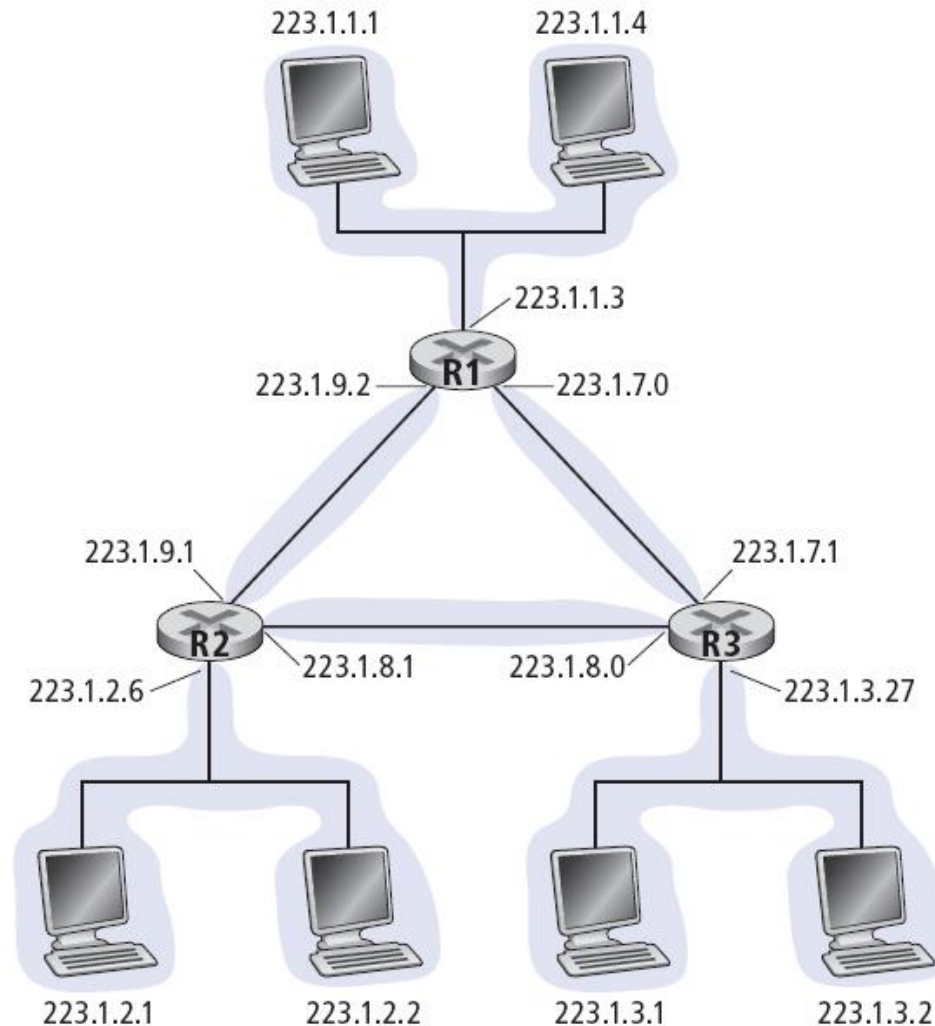
IP-Adressierung - Grundlagen

- IP-Adresse: 32-Bit-Kennung für das Interface (Schnittstelle) eines Endsystems oder eines Routers
- Interface: Verbindung zwischen dem System und dem Link
 - Wird normalerweise durch eine Netzwerkkarte bereitgestellt
 - Router haben typischerweise mehrere Interfaces
 - Endsysteme können ebenfalls mehrere Interfaces haben
 - Jedes Interface besitzt eine IP-Adresse



$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00000001}_1$$

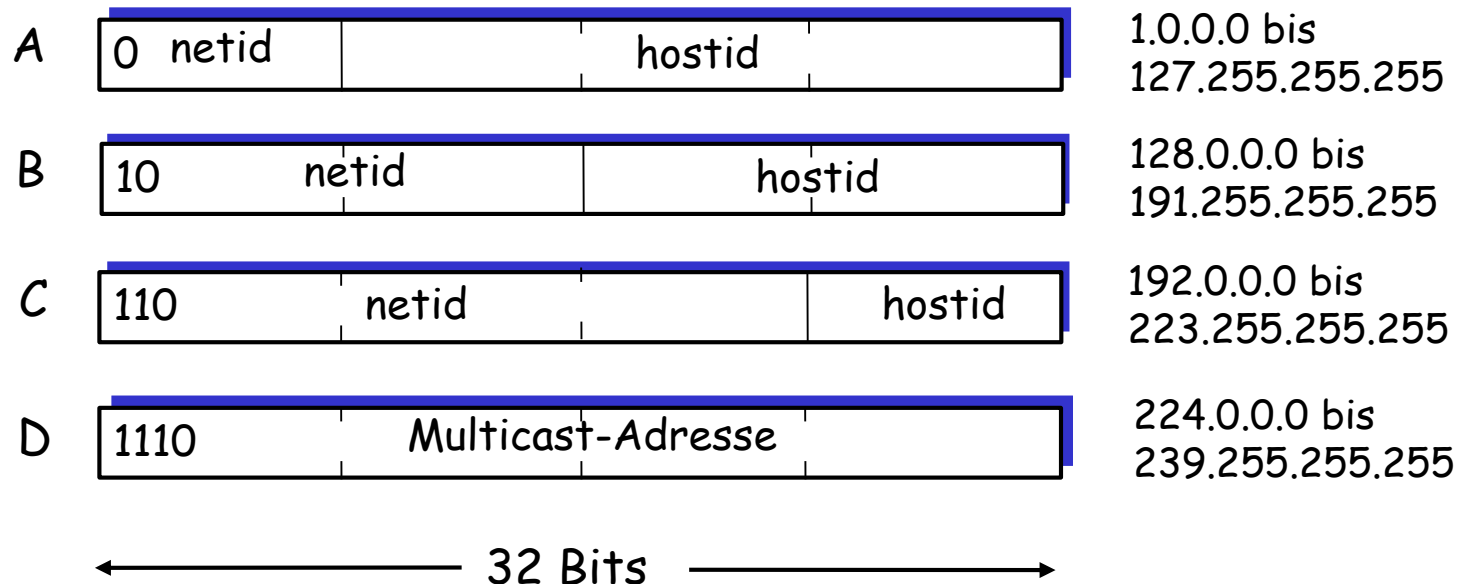
Wie viele Subnetzwerke sehen Sie?



IP-Adressierung - Adressklassen

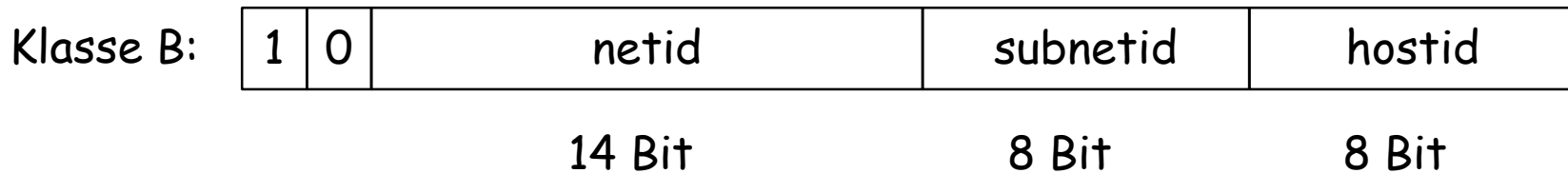
- Früher wurden IP-Adressen in Adressklassen aufgeteilt
- Die Klasse bestimmte das Verhältnis der Längen netid/hostid
- Dies nennt man „classfull“ addressing oder auch klassenbasierte Adressierung

Klasse



Adressierung von Subnetzen I

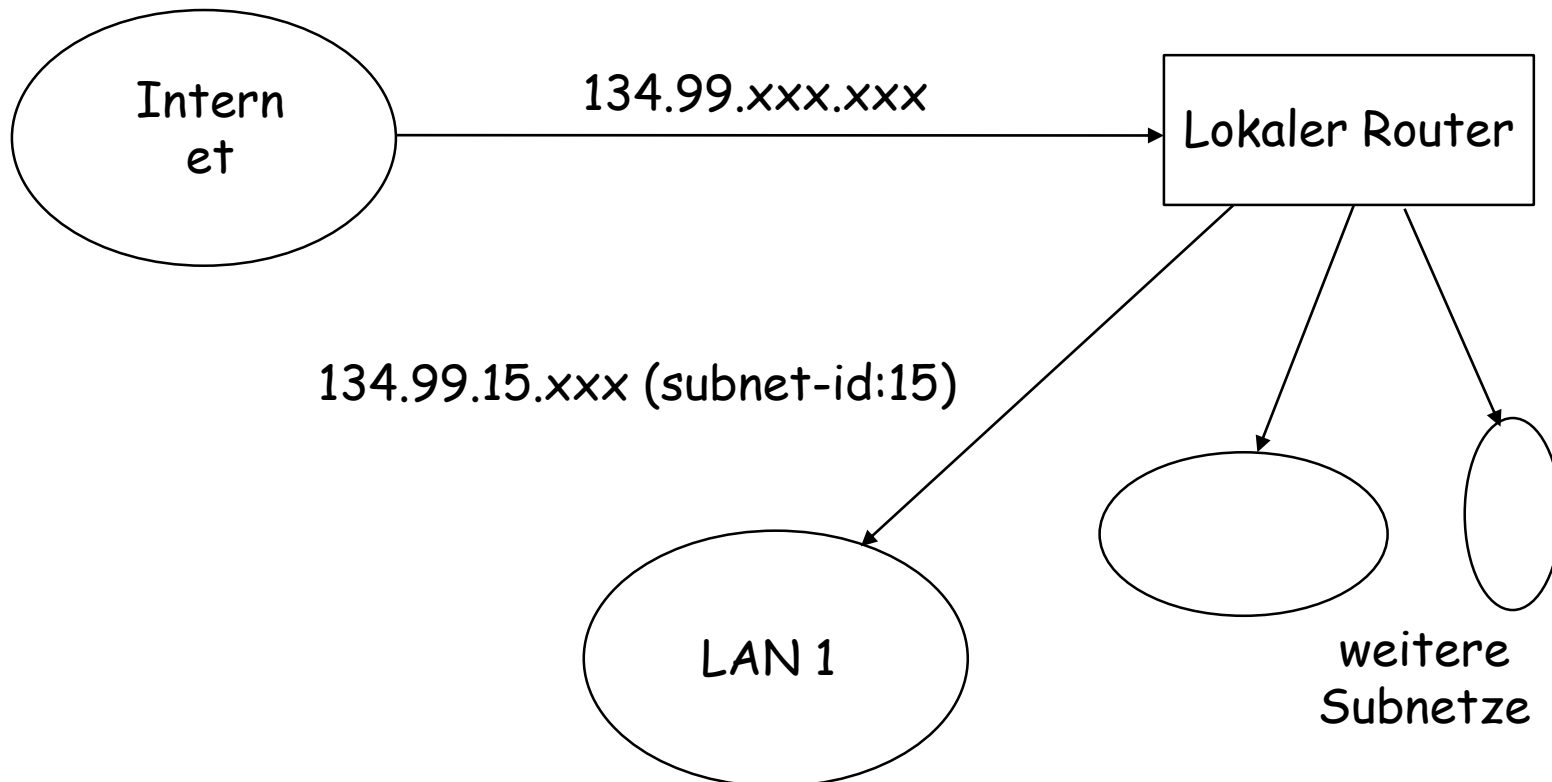
- Klasse-A- und -B-Adressen haben Platz für mehr Endsysteme, als man in einem Netzwerk sinnvoll unterbringen kann
- Daher teilt man die hostid weiter auf, z.B. so:



- Die Unterteilung (subnetid, hostid) ist eine lokale Entscheidung und wird von der Organisation vorgenommen, der die netid zugeordnet wurde

Adressierung von Subnetzen II

- Die subnetid ist außerhalb des Netzwerkes, für das sie verwendet wird, nicht sichtbar:



Adressierung von Subnetzen III

- Subnetzmaske (subnet mask)
 - Wird für jede IP-Adresse eines Systems im System gespeichert
 - Sie identifiziert, welcher Teil der Adresse zur subnetid und welcher zur hostid gehört

	16 Bit	8 Bit	8 Bit
Beispiel (Class B)	1111111111111111	11111111	00000000

subnet mask:
0xffffffff00=255.255.255.0 oder
auch /24

- Die eigene IP-Adresse in Verbindung mit der Subnetzmaske erlaubt Rückschlüsse darüber, wo sich eine andere IP-Adresse befindet:
 - im selben Subnetz (also direkt erreichbar)
 - im selben Netzwerk, aber in einem anderen Subnetz
 - in einem anderen Netzwerk

Beispiel für die Verwendung von Subnetzmasken

- Gegeben:
 - Eigene IP-Adresse: 134.155.48.10
 - Subnetzmaske: 255.255.255.0
 - Adresse A: 134.155.48.96, Adresse B: 134.155.55.96
- Überprüfen der beiden Adressen:
 - $134.155.48.10 \ \& \ 255.255.255.0 = 134.155.48.0$
 - $134.155.48.96 \ \& \ 255.255.255.0 = 134.155.48.0$ identisch, gleiches Subnetz
 - $134.155.55.96 \ \& \ 255.255.255.0 = 134.155.55.0$ verschieden, anderes Subnetz

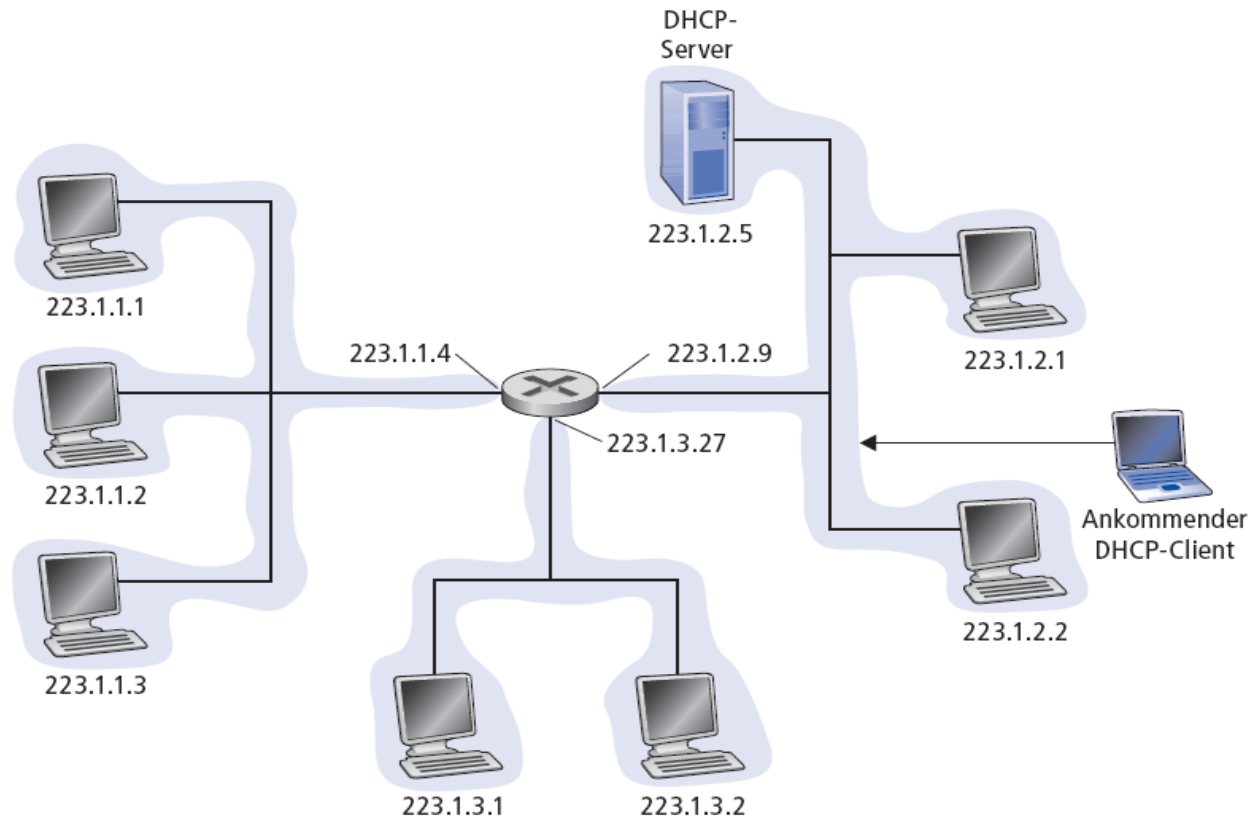
Subnetzmasken variabler Länge

- Problem: Gegeben sei ein Klasse-C-Netzwerk, welches in zwei Subnetze mit 50 Endsystemen und ein Subnetz mit 100 Endsystemen unterteilt werden soll.
- Das funktioniert nicht mit einer einzelnen Subnetzmaske!
 - 255.255.255.128: zwei Netze mit je 128 hostids
 - 255.255.255.192: vier Netze mit je 64 hostids
- Lösung: Subnetzmasken variabler Länge
 - Unterteile den Adressraum zunächst mit der kürzeren Subnetzmaske (1 Bit im Beispiel)
 - Unterteile eine Hälfte davon weiter mit der längeren Subnetzmaske (2 Bit im Beispiel)
 - Resultat: Subnetze verschiedener Größe

Frage: Wie bekommt ein Host seine IP-Adresse?

- Durch manuelle Konfiguration:
 - IP-Adresse
 - Subnetzmaske
 - Weitere Parameter
- **DHCP:** Dynamic Host Configuration Protocol: dynamisches Beziehen der Adresse von einem Server
 - "Plug-and-Play"

DHCP-Szenario



DHCP-Szenario

DHCP verwendet UDP.

DHCP-Nachrichten werden an die MAC-Broadcast-Adresse geschickt.

Es gibt ein Feld, in dem eine eindeutige Kennung des Clients verpackt ist. Dies ist meist die MAC-Adresse.

DHCP-Server:
223.1.2.5



Ankommender
Client



DHCP-Discover

src: 0.0.0.0, 68
dest: 255.255.255.255, 67
DHCPDISCOVER
yiaddr: 0.0.0.0
transaction ID: 654

DHCP-Offer

src: 223.1.2.5, 67
dest: 255.255.255.255, 68
DHCPOFFER
yiaddr: 223.1.2.4
transaction ID: 654
DHCP server ID: 223.1.2.5
Lifetime: 3600 secs

DHCP-Request

src: 0.0.0.0, 68
dest: 255.255.255.255, 67
DHCPREQUEST
yiaddr: 223.1.2.4
transaction ID: 655
DHCP server ID: 223.1.2.5
Lifetime: 3600 secs

DHCP-ACK

src: 223.1.2.5, 67
dest: 255.255.255.255, 68
DHCPACK
yiaddr: 223.1.2.4
transaction ID: 655
DHCP server ID: 223.1.2.5
Lifetime: 3600 secs

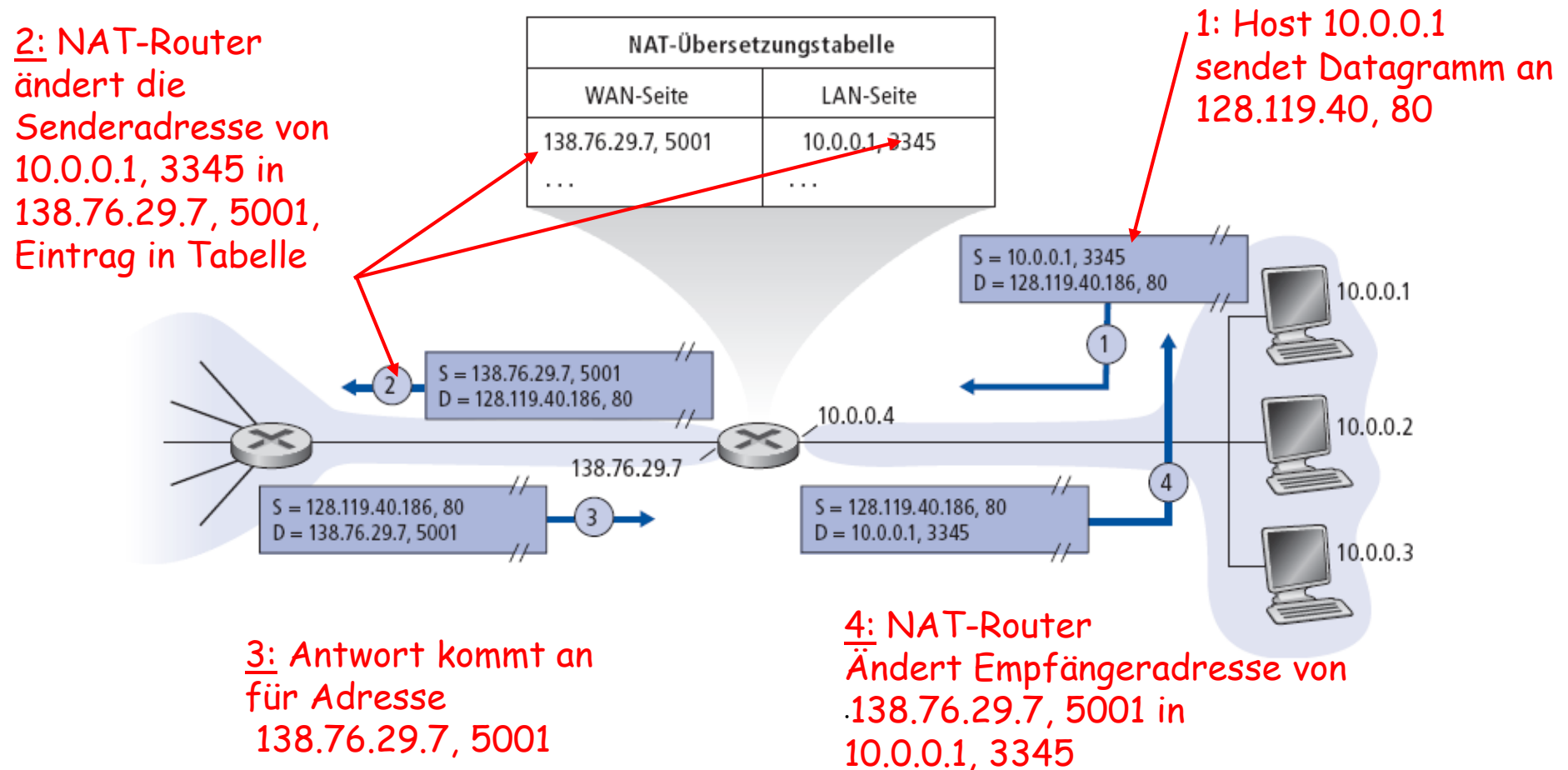
Zeit

Zeit

Besondere Adressen (RFC 3330)

Address Block	Present Use	Reference
0.0.0.0/8	"This" Network	[RFC1700, page 4]
10.0.0.0/8	Private-Use Networks	[RFC1918]
14.0.0.0/8	Public-Data Networks	[RFC1700, page 181]
24.0.0.0/8	Cable Television Networks	--
39.0.0.0/8	Reserved but subject to allocation	[RFC1797]
127.0.0.0/8	Loopback	[RFC1700, page 5]
128.0.0.0/16	Reserved but subject to allocation	--
169.254.0.0/16	Link Local	--
172.16.0.0/12	Private-Use Networks	[RFC1918]
191.255.0.0/16	Reserved but subject to allocation	--
192.0.0.0/24	Reserved but subject to allocation	--
192.0.2.0/24	Test-Net	
192.88.99.0/24	6to4 Relay Anycast	[RFC3068]
192.168.0.0/16	Private-Use Networks	[RFC1918]
198.18.0.0/15	Network Interconnect Device Benchmark Testing	[RFC2544]
223.255.255.0/24	Reserved but subject to allocation	--
224.0.0.0/4	Multicast	[RFC3171]
240.0.0.0/4	Reserved for Future Use	[RFC1700, page 4]

NAT: Network Address Translation

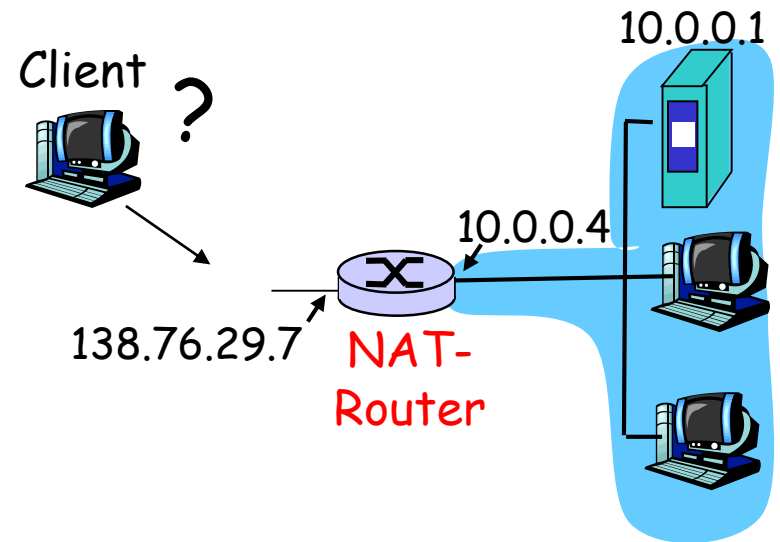


NAT: Network Address Translation

- 16-Bit-Port Number-Feld:
 - Mehr als 60.000 gleichzeitige Verbindungen mit einer IP-Adresse
- NAT ist nicht unumstritten:
 - Router sollten nur Informationen der Schicht 3 verwenden
 - Verletzung des sogenannten Ende-zu-Ende-Prinzips (end-to-end principle):
 - Transparente Kommunikation von Endsystem zu Endsystem, im Inneren des Netzes wird nicht an den Daten „herumgepfuscht“
 - Bei NAT: Der Anwendungsentwickler muss die Präsenz von NAT-Routern berücksichtigen. Beispiele:
 - Verwenden der IP-Adresse als weltweit eindeutige Nummer
 - Verwenden von UDP
 - NAT dient hauptsächlich der Bekämpfung der Adressknappheit im Internet. Dies sollte besser über IPv6 (s. später) erfolgen

Durchqueren von NAT

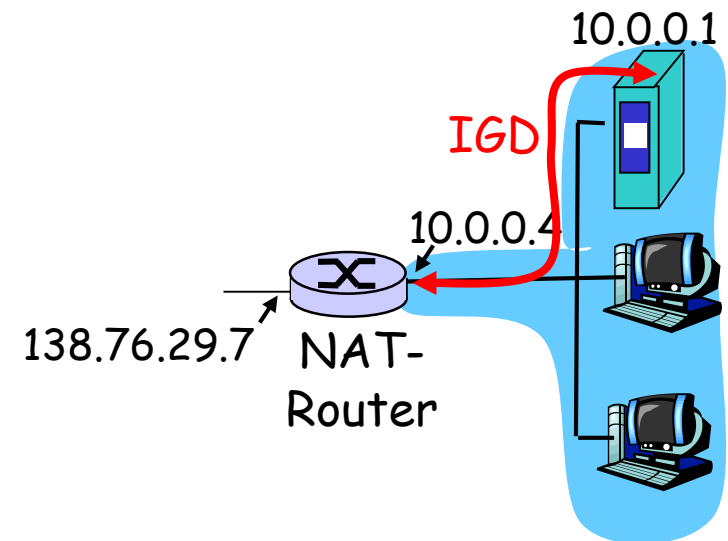
- Engl. NAT traversal
- Der Client möchte den Server mit der Adresse 10.0.0.1 kontaktieren
 - Die Adresse 10.0.0.1 ist eine lokale Adresse und kann nicht als Adresse im globalen Internet verwendet werden
 - Die einzige nach außen sichtbare Adresse ist: 138.76.29.7
- Lösung 1: Statische Konfiguration von NAT, so dass eingehende Anfragen geeignet weitergeleitet werden
 - Beispiel: (123.76.29.7, Port 2500) wird immer an 10.0.0.1, Port 25000 weitergeleitet



Durchqueren von NAT

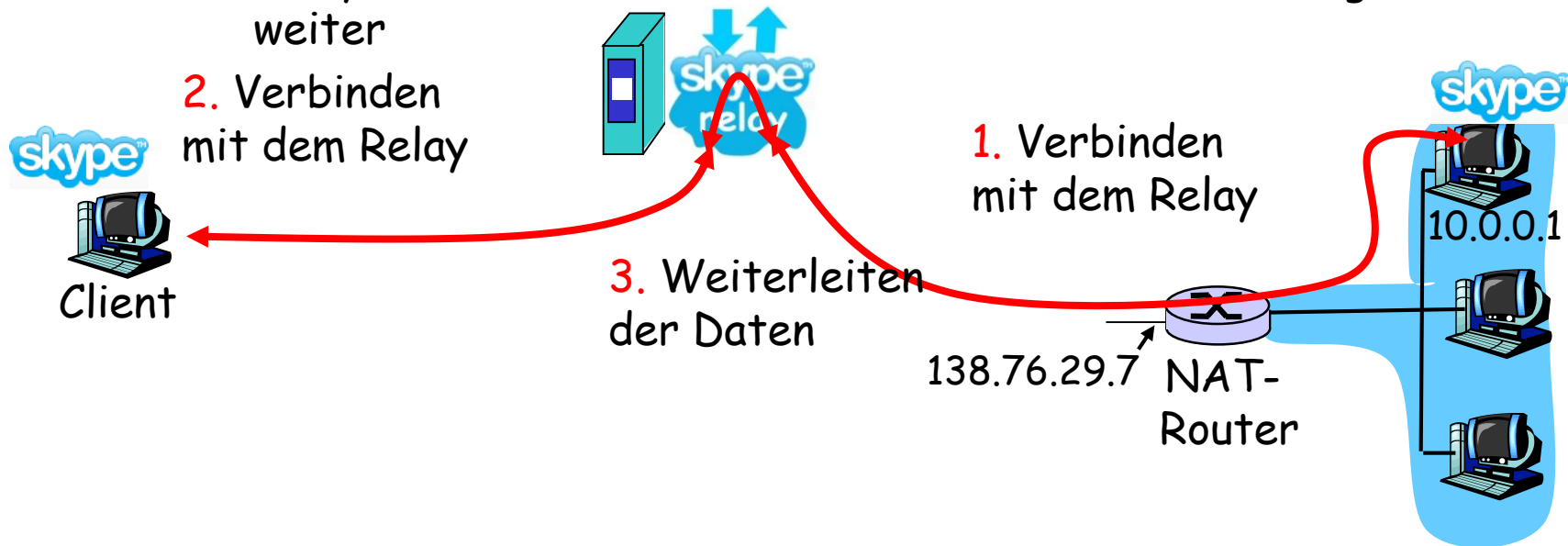
- Lösung 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Dies ermöglicht dem Host hinter dem NAT Folgendes:
 - v Herausfinden der öffentlichen IP-Adresse des NAT-Routers (138.76.29.7)
 - v Kennenlernen existierender Abbildungen in der NAT-Tabelle
 - v Einträge in die NAT-Tabelle einfügen oder aus ihr löschen

Das heißt automatische
Konfiguration von statischen
NAT-Einträgen



Durchqueren von NAT

- Lösung 3: Relaying (von Skype verwendet)
 - Server hinter einem NAT-Router baut eine Verbindung zu einem Relay auf (welches nicht hinter einem NAT-Router liegt)
 - Client baut eine Verbindung zum Relay auf
 - Relay leitet die Pakete vom Client zum Server und umgekehrt weiter



Kapitel 4:

Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - **ICMP**
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

ICMP: Internet Control Message Protocol

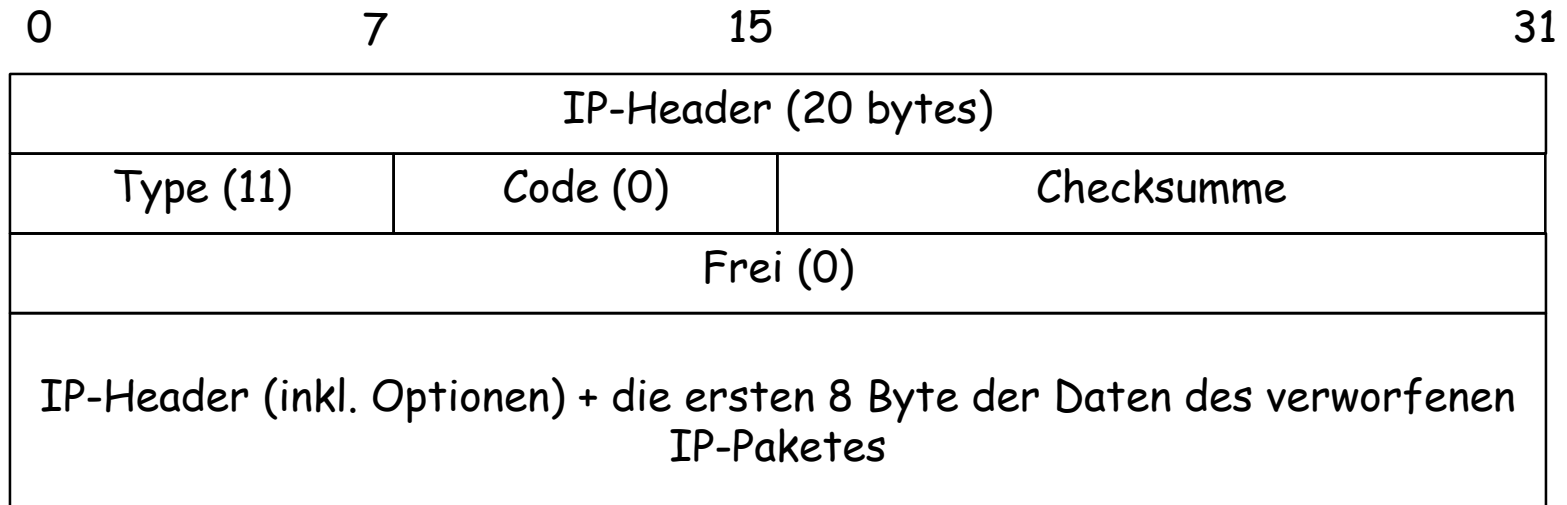
- Wird von Hosts und Routern verwendet, um Informationen über das Netzwerk selbst zu verbreiten
 - Fehlermeldungen: Host, Netzwerk, Port, Protokoll nicht erreichbar
 - Echo-Anforderung und Antwort (von ping genutzt)
- Gehört zur Netzwerkschicht, wird aber in IP-Datagrammen transportiert
- ICMP-Nachricht:** Type, Code und die ersten 8 Byte des IP-Datagramms, welches die Nachricht ausgelöst hat

<u>Type</u>	<u>Code</u>	<u>Beschreibung</u>
0	0	echo reply (ping)
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery
11	0	TTL expired
12	0	bad IP header

Traceroute und ICMP

- Aufgabe:
 - Traceroute bestimmt Informationen über alle Router, die auf dem Weg zu einer IP-Adresse liegen
 - Dabei wird auch die Round-Trip-Zeit zu jedem Router bestimmt
- Funktionsweise:
 - Traceroute schickt ein UDP-Paket an die Adresse, für die der Weg untersucht werden soll; TTL im IP-Header wird auf 1 gesetzt
 - Der erste Router verwirft das IP-Paket (TTL = 1!) und schickt eine ICMP-Time-Exceeded-Fehlermeldung an den Absender
 - Traceroute wiederholt dies mit TTL = 2 etc.

ICMP-Time-Exceeded-Nachricht



Kapitel 4:

Netzwerkschicht

- 4.1 Einleitung
- 4.2 Aufbau eines Routers
- 4.3 IP: Internet Protocol
 - Datagrammformat
 - IPv4-Adressierung
 - ICMP
 - IPv6
- 4.4 statisches Routen
- 4.5 Routing-Algorithmen
 - Link State
 - Distance Vector
 - Hierarchisches Routing
- 4.6 Routing im Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast- und Multicast-Routing

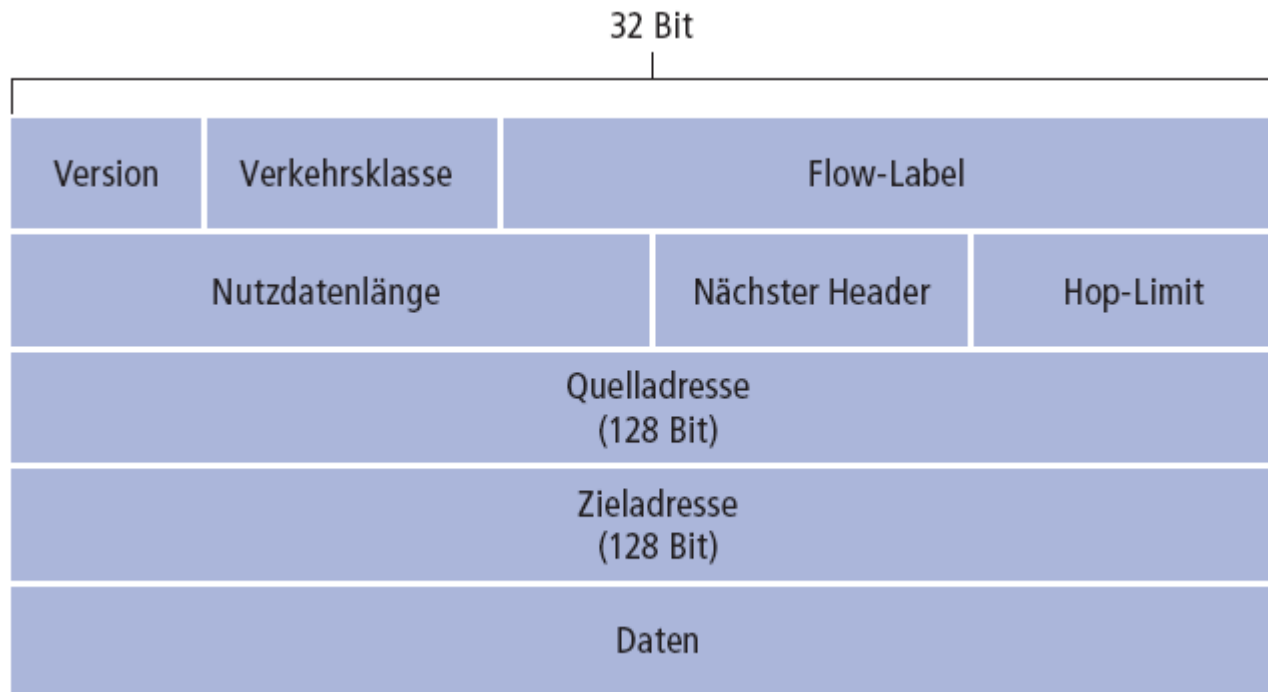
- **Ursprüngliche Motivation:** 32-Bit-Adressen werden in naher Zukunft komplett zugeteilt sein
 - NAT hat dies ein wenig verzögert, aber das grundlegende Problem nicht gelöst
 - Beispiel: Was passiert, wenn jedes Handy eine feste IP-Adresse bekommen soll?
- **Weitere Motivation:**
 - Vereinfachtes Header-Format für eine schnellere Verarbeitung in den Routern
 - Header soll Dienstgütemechanismen (Quality of Service, QoS) unterstützen
- **IPv6-Datagrammformat:**
 - Header fester Länge (40 Byte)
 - keine Fragmentierung in den Routern

IPv6-Header

Verkehrsklasse: Priorisierung von Datagrammen

Flow Label: Identifikation von zusammengehörigen Flüssen von Datagrammen (z.B. ein Voice-over-IP-Telefonat)

Nächster Header: An welches Protokoll sollen die Daten im Datenteil übergeben werden? Beispiel: TCP!



Weitere Veränderungen gegenüber IPv4

- **Checksumme:** entfernt, um die Verarbeitung in den Routern zu erleichtern
- **Optionen:** als separate Header, die auf den IP-Header folgen
 - Werden durch das "Nächster Header"-Feld angezeigt
 - Einfachere Behandlung in Hosts und Routern
- **ICMPv6:** neue Version von ICMP
 - Zusätzliche Pakettypen, z.B. "Packet Too Big"
 - Funktionen zur Verwaltung von Multicast-Gruppen (später mehr)

Übergang von IPv4 zu IPv6

- Es können nicht alle Router gleichzeitig umgestellt werden
 - Wie kann ein Netzwerk funktionieren, in dem sowohl IPv4- als auch IPv6-Router vorhanden sind?
- *Tunneling*: IPv6 wird im Datenteil von IPv4-Datagrammen durch das klassische IPv4-Netzwerk transportiert

Tunneling

Logische Sicht



Reale Situation

