# Computer Networks (WS23/24)

## L5: The Network Layer - Part 1

**Prof. Dr. Lin Wang**
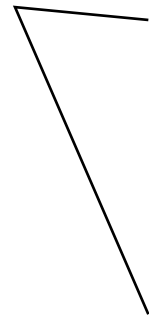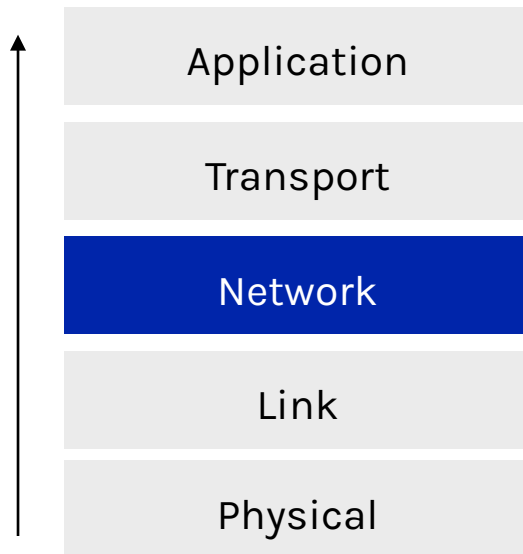
Computer Networks Group (PBNet)

Department of Computer Science
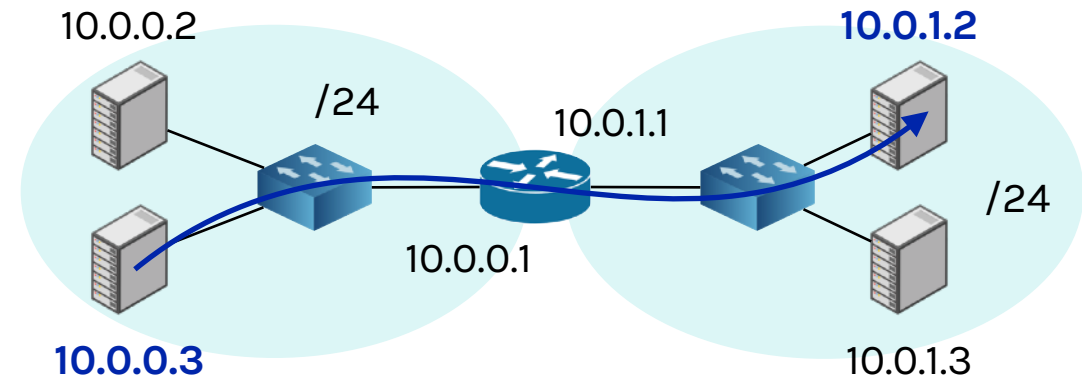
Paderborn University

Materials inspired by Shyam Gollakota and Laurent Vanbever

# Learning objectives



Application

Transport

**Network**

Link

Physical

10.0.0.2

/24

10.0.1.1

10.0.1.2

10.0.0.1

10.0.0.3

10.0.1.3

What happens when a packet travels across networks?

**Part 1**
- Inter-networking
- IP forwarding
- Network Address Translation (NAT)
- Helper protocols: ARP, DHCP, ICMP

2

# Inter-networking
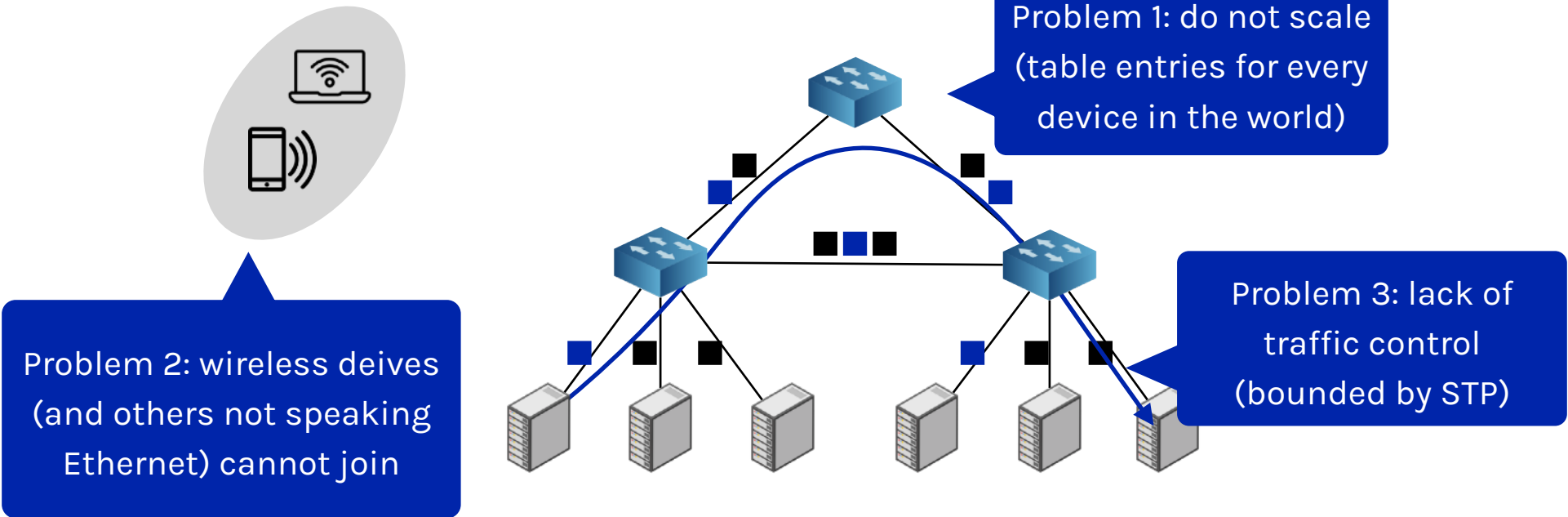
# Internet based on switched Ethernet



Problem 1: do not scale (table entries for every device in the world)

Problem 2: wireless deives (and others not speaking Ethernet) cannot join

Problem 3: lack of traffic control (bounded by STP)

4

# Inter-networking

Ethernet

WiFi

IP router

5G

Ethernet

# The network layer

# People behind it

**Pioneers: Cerf and Kahn**

- Fathers of the Interent

- In 1974, later led to TCP/IP

**Tackled the problem of interconnecting networks**

- Instead of mandating a single network technology



Vinton Cerf



Bob Kahn

ACM Turing Award 2004

# Internet reference model

**Internet Protocol (IP) is the narrow waist**

- Supports many different links below:
  Ethernet, WiFi, 4G/LTE, 5G

- Supports many application above

**IP as a lowest common denominator**

- Asks little of lower-layer networks

- Gives little as a higher-layer service

# IPv4 address

**Unique 32-bit numbers associated to a network interface (on a host or a router…)**

Usually written using dotted-quad notation

## 82 . 130 . 102 . 10

01010010    10000010    01100110    00001010

# IPv6 address

**Notation**

- 8 groups of 16 bits, each separated by colons (:)

- Leading zeros in any groups are removed

- One section of zeros is replaced by a double colon (::)

**Examples**

- 1080:0:0:0:8:800:200C:417A → 1080::8:800:200C:417A

- FF01:0:0:0:0:0:0:0101 → FF01::101

- 0:0:0:0:0:0:0:1 → ::1

# IP address assignment



1.2.3.4   5.6.7.8   2.4.6.8

1.2.3.5   5.6.7.9   2.4.6.9

LAN 1

LAN 2

WAN 1

WAN 2

IP router   IP router   IP router

1.2.3.4 ←
1.2.3.5 →
...

**Forwarding table**

For each IP address, we need an entry in the table

LAN: Local Area Network

WAN: Wide Area Network

# Two universal tricks in Computer Science

When you need…        more **flexibility**

You add…        a layer of **indirection**

When you need…        more **scalability**

You add…        a **hierarchical** structure

# Hierarchical postal addresses

| | |
|---:|:---|
| **City** | Paderborn |
| **Zip** | 33098 |
| **Street** | Pohlweg |
| **Building** | O |
| **Room (in building)** | O3-158 |
| **Name** | Lin Wang |

Nobody maintains where every single building is in the Deutsch Post system

# Hierarchical forwarding

**Step 1**     Deliver the letter to the post office responsible for the city and zip code

**Step 2**     Assign letter to the mail person covering the street

**Step 3**     Drop letter into the mailbox attached to the building

**Step 4**     Hand in the letter to the addressed person

# IP addresses are hierarchical

32 bits

01010010.10000010.01100110.00001010

**Prefix**

(Identifying the **network**)

**Suffix**

(Identifying the **hosts** in the network)

# IP prefix

82 . 130 . 102 . 0 **/24**

↑

Prefix length
(in bits)

| Prefix part | Suffix part | IP address |
|---|---|---|
| 01010010.10000010.01100110. | 00000000 | 82.130.102.0 |

**Identifies the network itself**

| | | |
|---|---|---|
| 01010010.10000010.01100110. | 00000001 | 82.130.102.1 |
| 01010010.10000010.01100110. | 00000010 | 82.130.102.2 |
| 01010010.10000010.01100110. | 00000011 | 82.130.102.3 |
| . . . . . . | | |
| 01010010.10000010.01100110. | 11111110 | 82.130.102.254 |
| 01010010.10000010.01100110. | 11111111 | 82.130.102.255 |

**Broadcast address**

Only 254 valid addresses to allocate to hosts for /24

# IP prefix with masks

Address | 82.130.102.0
01010010.10000010.01100110.00000000

Mask | 255.255.255.0
11111111.11111111.11111111.00000000

ANDing the address and the mask gives you the IP prefix

# Scalable forwarding

1.2.3.4    1.2.3.5    1.2.3.253

5.6.7.1    5.6.7.2    5.6.7.200

LAN 1

LAN 2

WAN 1    WAN 2

IP router    IP router    IP router

1.2.3.0/24 ←
5.6.7.0/24 →

Only two entries needed
with prefix matching

**Forwarding table**

# Legacy classful networking

| | Leading bits | Prefix length | #hosts | Start addr. | End addr. |
|---|---|---|---|---|---|
| **Class A** | 0 | 8 | $2^{24}$ | 0.0.0.0 | 127.255.255.255 |
| **Class B** | 10 | 16 | $2^{16}$ | 128.0.0.0 | 191.255.255.255 |
| **Class C** | 110 | 24 | $2^{8}$ | 192.0.0.0 | 223.255.255.255 |
| **Class D** (Multicast) | 1110 | | | 224.0.0.0 | 239.255.255.255 |
| **Class E** (Reserved) | 1111 | | | 240.0.0.0 | 255.255.255.255 |

Class C was too small ⇒ everyone requested class B (too big, a lot of waste)

# Classless Inter-Domain Routing (CIDR)

**Enables flexible division between network and hosts addresses**

**CIDR must specify both the address and the mask**

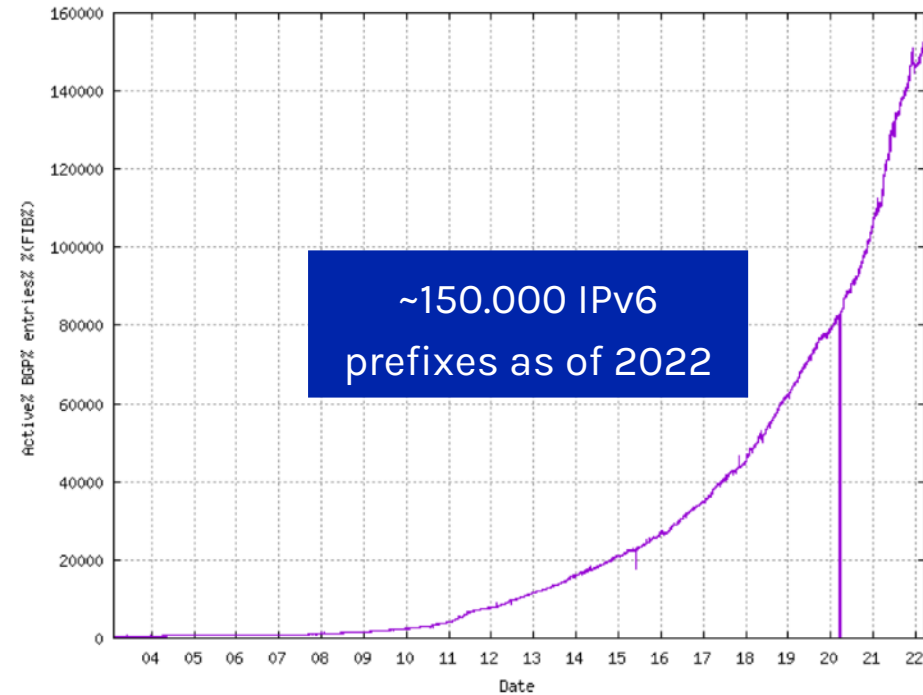- Mask in the classful address is implicit in the first address bits

- Mask in CIDR is carried by the routing algorithm, not implicitly in the address

**Example: an organization needs 500 addresses**

| Type | Allocation | Utilization |
|------|-----------|-------------|
| Classful | Class B (/16) | 1% |
| CIDR | /23 (2 Class C's) | 98% |

With CIDR, the address utilization is always higher than 50%. Why?

# IP prefixes on the Internet



> 900.000 IPv4
prefixes as of 2022



~150.000 IPv6
prefixes as of 2022

http://www.cidr-report.org/
https://www.cidr-report.org/v6/as2.0/

# Hierarchical IP address allocation



**Internet Corporation for Assigned Names and Numbers**

Reginal Internet Registries (RIRs)

America  Latin America  Europe  Asia-Pacific  Africa

ISPs / Institutions

# IP Forwarding

# What's inside an IP router?

# What's inside an IP router?

# Forwarding table on IP routers



Forwarding table

| IP prefix | Output |
|-----------|--------|
| 129.0.0.0/8 | port2 |
| 129.132.1.0/24 | port2 |
| 129.132.2.0/24 | port2 |
| 129.133.0.0/16 | port3 |

# Forwarding based on prefix matching

Packet arrival at ISP 2

`DST_ADDR = 129.0.0.1`

Forwarding table

| IP prefix | Output |
|---|---|
| **129.0.0.0/8** | **port2** |
| 129.132.1.0/24 | port2 |
| 129.132.2.0/24 | port2 |
| 129.133.0.0/16 | port3 |

**port2**

ISP 2

port3

129.0.0.0/8

ISP 1

129.133.0.0/16

129.132.1.0/24      129.132.2.0/24      129.132.4.0/24

Router at ISP 2 performs an IP lookup to find the matching prefix
→ forward the packet to port2

# Forwarding based on prefix matching

Packet arrival at ISP 2

DST_ADDR = 129.133.0.1

Forwarding table

| IP prefix | Output |
|-----------|--------|
| **129.0.0.0/8** | **port2** |
| 129.132.1.0/24 | port2 |
| 129.132.2.0/24 | port2 |
| **129.133.0.0/16** | **port3** |

port2    ISP 2    port3

129.0.0.0/8

ISP 1

129.133.0.0/16

129.132.1.0/24    129.132.2.0/24    129.132.4.0/24

Where the packet should be forwarded?

# Longest prefix matching

Packet arrival at ISP 2

`DST_ADDR = 129.133.0.1`

Forwarding table



| IP prefix | Output |
|:---:|:---:|
| 129.0.0.0/8 | port2 |
| 129.132.1.0/24 | port2 |
| 129.132.2.0/24 | port2 |
| **129.133.0.0/16** | **port3** |

Router at ISP 2 matches on the longest prefix → forward the packet to port3

port2

ISP 2

**port3**

129.0.0.0/8

ISP 1

129.133.0.0/16

129.132.1.0/24    129.132.2.0/24    129.132.4.0/24

# Hardware support for longest prefix matching

**Ternary Content Addressable Memory (TCAM)**

Address → DRAM/SRAM → Content

Content → TCAM → Address

IP address
0100110010101 →

IP prefix records:
- 010011010100
- 0100110 ????
- 0101 ????

→ Match encoder → Address → Decoder → Output port

# Simplifying the forwarding table

**A child prefix can be filtered from the table if it shares the same output as its parent**

129.132.1.0/24

Child

129.133.0.0/16

Child

Child

129.132.2.0/24

Parent: 129.0.0.0/8

| IP prefix | Output | |
|---|---|---|
| 129.0.0.0/8 | port2 | |
| 129.132.1.0/24 | port2 | ✖ |
| 129.132.2.0/24 | port2 | ✖ |
| 129.133.0.0/16 | port3 | |

| IP prefix | Output |
|---|---|
| 129.0.0.0/8 | port2 |
| 129.133.0.0/16 | port3 |

# IPv4 packet format

32 bits

| Version | HLEN | Type of service | Total length | | |
|---|---|---|---|---|---|
| Identification | | | Flags | Fragment offset | |
| Time to live | | Protocol | Header checksum | | |
| Source IP address | | | | | |
| Destination IP address | | | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

# IPv4 packet format

| Version | HLEN | Type of service | Total length | | |
|---------|------|-----------------|--------------|--|--|
| IPv4: 4, IPv6: 6 ation | | | Flags | Fragment offset | |
| Time to live | | Protocol | Header checksum | | |
| Source IP address | | | | | |
| Destination IP address | | | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

# IPv4 packet format

| Version | HLEN | Type of service | | Total length | |
|---------|------|-----------------|--|--------------|--|
| | | | | | Fragment offset |
| | Time | | | eader checksum | |
| | | Source IP address | | | |
| | | Destination IP address | | | |
| | | Options (if any) | | | |
| | | Payload (data) | | | |

The number of 32-bit words in the header, typically set to 5 (20 bytes header)

# IPv4 packet format

| Version | HLEN | Type of service | Total length | | |
|---------|------|-----------------|--------------|---|---|
| Identification | | | Flags | Fragment offset | |
| Time to live | | | | checksum | |
| Destination IP address | | | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

Allows different packets to be treated differently, e.g., low latency for voice, high bandwidth for video

# IPv4 packet format

| Version | HLEN | Service type | Total length | | |
|---------|------|--------------|--------------|--|--|
| Identification | | | Flags | Fragment offset | |
| Time to live | | Prot... | ...um | | |
| Source IP address | | | | | |
| Destination IP address | | | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

Number of bytes in the entire packet,
with a maximum of 65.535 bytes

# IPv4 packet format

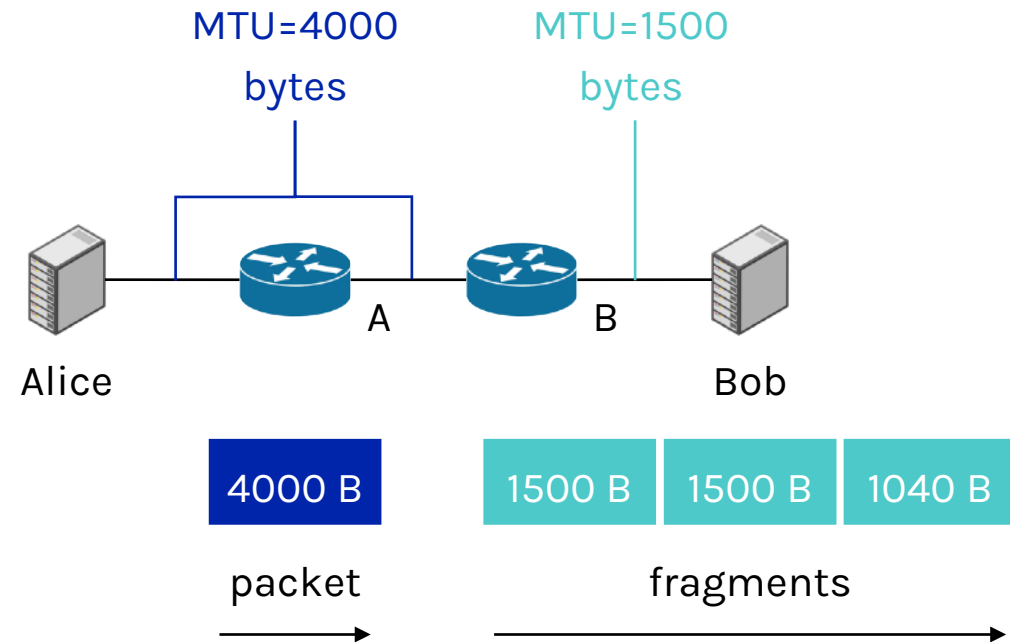| Version | HLEN | Service type | Total length | | |
|---------|------|--------------|--------------|--------|------------------|
| Identification | | | Flags | | Fragment offset |
| Time to live | | Protocol | Header checksum | | |
| | | | Used when packets get fragmented | | |
| | | Destination IP address | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

# Maximum transmission unit (MTU)

**MTU is the maximum number of bytes a link can carry as one unit**

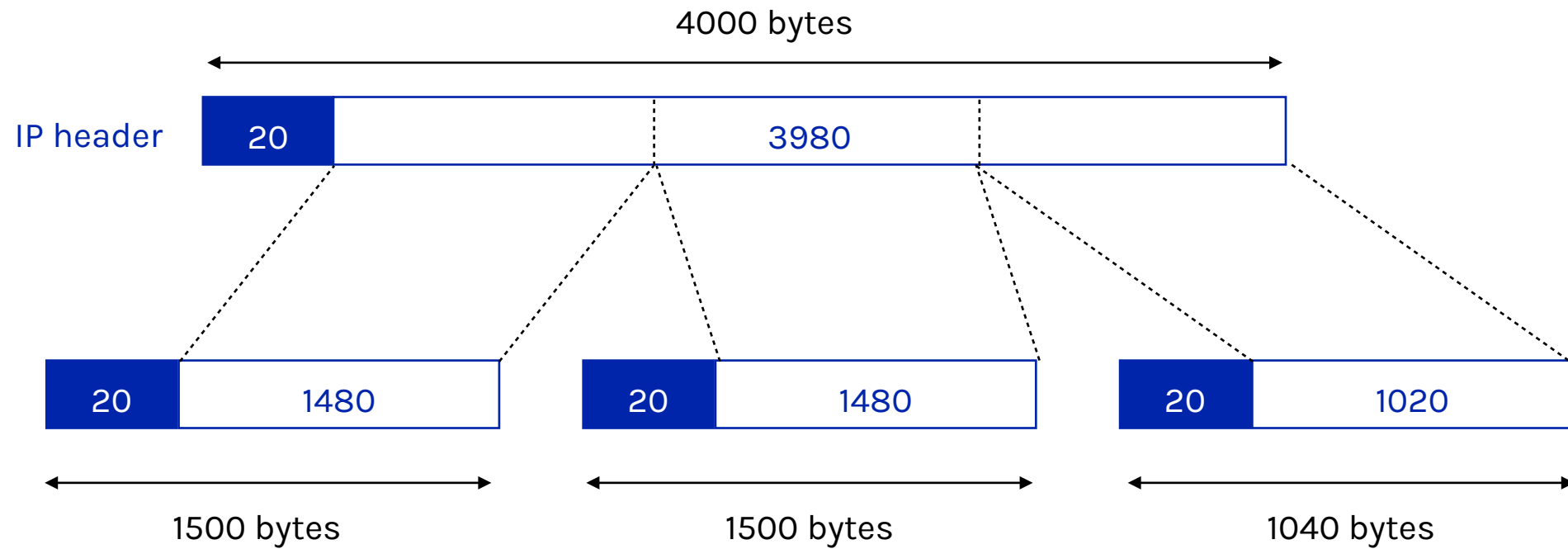- 1500 bytes for normal Ethernet, 9000 for Jumbo frames

**A router fragments a packet if outgoing link MTU < packet size**

**Fragmented packets are recomposed at the destination**

- Why not directly in the network?

MTU=4000 bytes          MTU=1500 bytes

Alice                    A        B        Bob

| 4000 B |
packet

| 1500 B | 1500 B | 1040 B |
fragments

# IP fragmentation

4000 bytes

| IP header | 20 | 3980 |

| 20 | 1480 |

1500 bytes

| 20 | 1480 |

1500 bytes

| 20 | 1020 |

1040 bytes

# IPv4 packet format

| Version | HLEN | Service type | Total length | |
|---------|------|--------------|--------------|--|
| Identification | | | Flags | Fragment offset |
| Time to live | | Protocol | Header checksum | |

Uniquely identify the fragments of a particular packet

Destination IP address

Options (if any)

Payload (data)

# IPv4 packet format

| Version | HLEN | Service type | Total length | | |
|---------|------|--------------|--------------|---|---|
| Identification | | | Flags | Fragment offset | |
| Time to live | | Protocol | Header checksum | | |
| | | | | | |
| | | | | | |
| Options (if any) | | | | | |
| Payload (data) | | | | | |

Used for putting back the fragments in the right order in case of reordering

# IPv4 packet format

| Version | HLEN | Service type | Total length |
|---------|------|--------------|--------------|
| Identification | | Flags | Fragment offset |
| Time to live | Protocol | | Header checksum |
| | | | |
| Destination IP address | | | |
| Options (if any) | | | |
| Payload (data) | | | |

Whether or not there are more fragments coming

# IPv4 packet format

| Version | HLEN | Service type | Total length | | |
|---------|------|--------------|--------------|---|---|
| Identification | | | Flags | Fragment offset | |
| Time to live | | Protocol | Header checksum | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

Time-to-live (TTL) is decremented by 1 at each router and the packet is discarded if TTL reaches 0

Default TTL values:
Windows: 128
Linux/Mac: 64
(Can be used for OS fingerprinting)

# IPv4 packet format

| Version | HLEN | Service type | | Total length | |
|---------|------|--------------|--|--------------|--|
| Identification | | | Flags | Fragment offset | |
| Time to live | | Protocol | | Header checksum | |
| | | | | | |
| | | | | | |
| | | Options (if any) | | | |
| Payload (data) | | | | | |

Identifying the higher-level protocol carried in the packet: "6" for TCP, "17" for UDP

# IPv4 packet format

| Version | HLEN | Service type | Total length | |
|---|---|---|---|---|
| Identification | | | Flags | Fragment offset |
| Time to live | | Protocol | **Header checksum** | |
| Source IP address | | | | |
| Des... | | | | |
| | | | | |
| Payload (data) | | | | |

Internet checksum calculated in 16 bits (does not protect the payload)

# IPv4 packet format

| Version | HLEN | Service type | Total length | |
|---|---|---|---|---|
| Identi... | | IP options include: record route, strict source route, loose source route, timestamp, traceroute, route alert. For security reasons, there are often disactivated. | ...ment offset | |
| Time to live | | | ...ecksum | |
| | | | | |
| Destination IP address | | | | |
| **Options (if any)** | | | | |
| Payload (data) | | | | |

IPv4 addresses have been exhausted, but they still account for most of the Internet traffic

# IPv6 is simpler than IPv4

**Removed**

- Fragmentation
- Checksum                    ⊢──────────→   Leave problems to the end-host
- Header length   ⊢──────────────→   Simplify handling

**Added**

- New options mechanism   ⊢──────→   Simplify handling
- Expanded addresses
- Flow label   ⊢────────────→   Flexibility

# IPv4 vs. IPv6

# IPv6 options

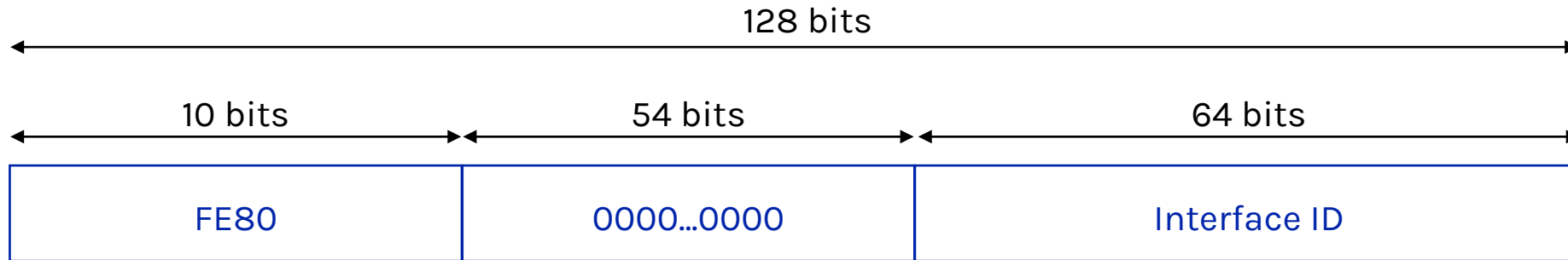**Enables to insert arbitrary options in the header (see RFC 2460)**

# IPv6 unicast address

**Hierarchically allocated, similar to global IPv4 addresses**

128 bits

| N bits | M bits | 128-N-M bits |
|---|---|---|
| Global routing prefix | Subnet ID | Interface ID |
| Identifies the ISP responsible for the address | A subnet or a customer of this ISP | Usually 64 bits, based on the MAC address (EUI-64, deprecated), obtained from a server |

Currently, only 2000::/3 is used for global unicast; all addresses are in the range of 2000 to 3FFFF

# IPv6 link-local address

**Same as private IPv4 addresses**

128 bits

| 10 bits | 54 bits | 64 bits |
|---|---|---|
| FE80 | 0000…0000 | Interface ID |

```
en0: flags=8863<UP,BROADCAST,SMART,RUNNING,SIMPLEX,MULTICAST> mtu 1500
        options=400<CHANNEL_IO>
        ether 6c:7e:67:d7:f5:71
        inet6 fe80::c8c:a9ae:f3b1:8b9d%en0 prefixlen 64 secured scopeid 0xf
        inet 192.168.2.104 netmask 0xffffff00 broadcast 192.168.2.255
        inet6 2003:d0:271b:2f58:859:fcea:5d83:f3bb prefixlen 64 autoconf secured
        inet6 2003:d0:271b:2f58:f082:f7a3:7b05:eb99 prefixlen 64 autoconf temporary
        inet6 2003:d0:271b:2fbf:82b:fb86:b0ba:9a79 prefixlen 64 autoconf secured
        inet6 2003:d0:271b:2fbf:8588:7d45:e7cc:6c51 prefixlen 64 autoconf temporary
        inet6 2003:d0:271b:2fb8:40d:e3fb:2690:5fb9 prefixlen 64 autoconf secured
        inet6 2003:d0:271b:2fb8:38d2:1256:1047:aa04 prefixlen 64 autoconf temporary
        nd6 options=201<PERFORMNUD,DAD>
        media: autoselect
        status: active
```
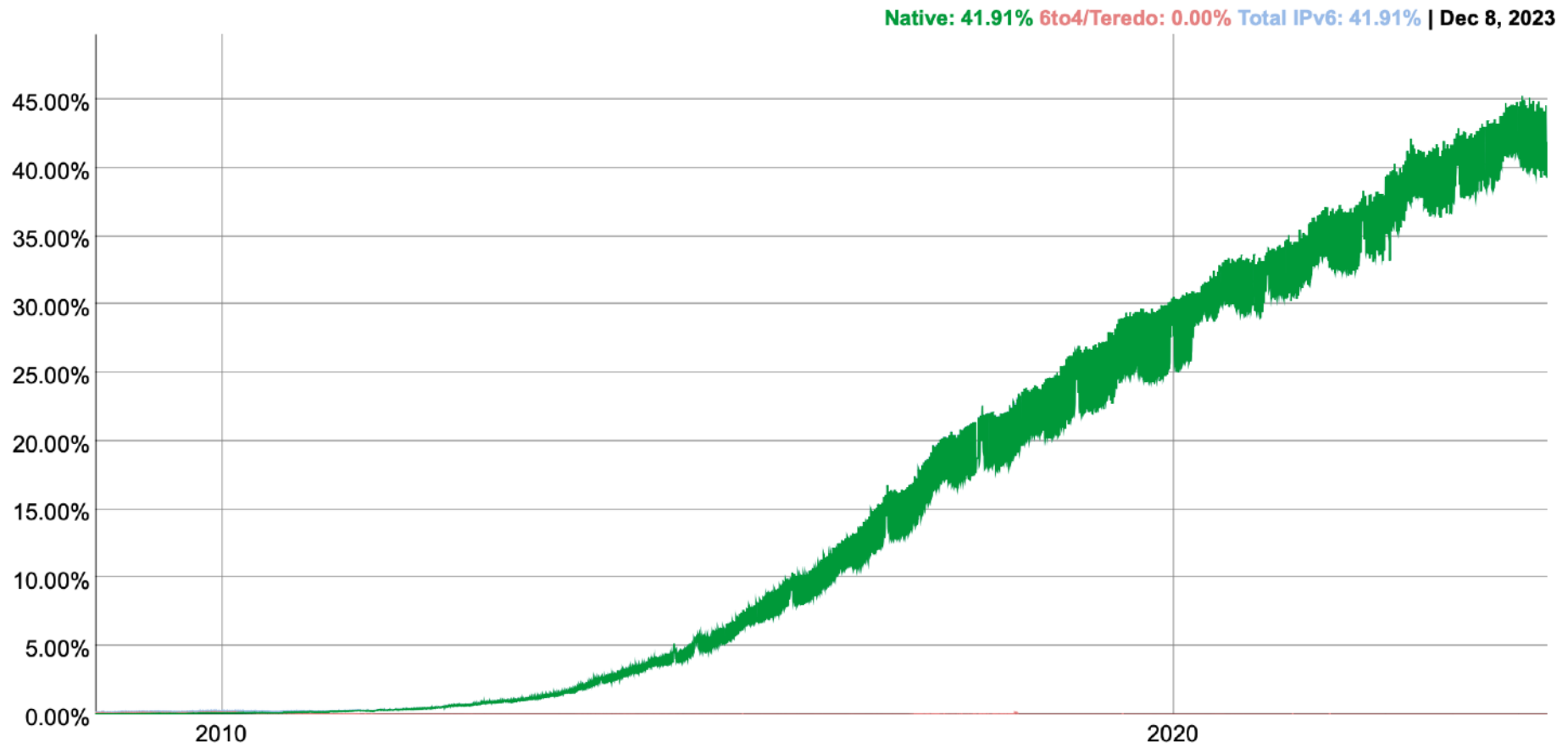
Each host/router must generate a link-local address for each of its interfaces

An interface can have multiple IPv6 addresses

# IPv6 adoption



Native: **41.91%** 6to4/Teredo: **0.00%** Total IPv6: **41.91%** | **Dec 8, 2023**

# IPv6 deployment challenges

Requires every device to support it (all routers, end-hosts, middleboxes, applications…)

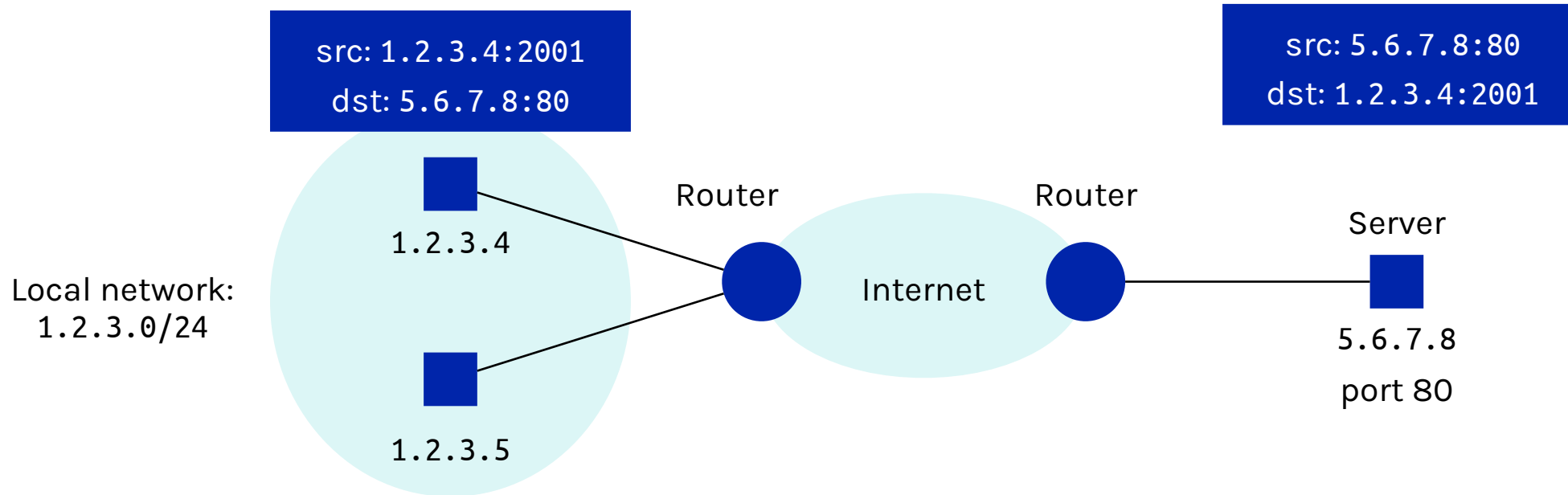Most of IPv6 features were back-ported to IPv4 (no obvious advantages in using IPv6)

Network Address Translation (NAT) is working well
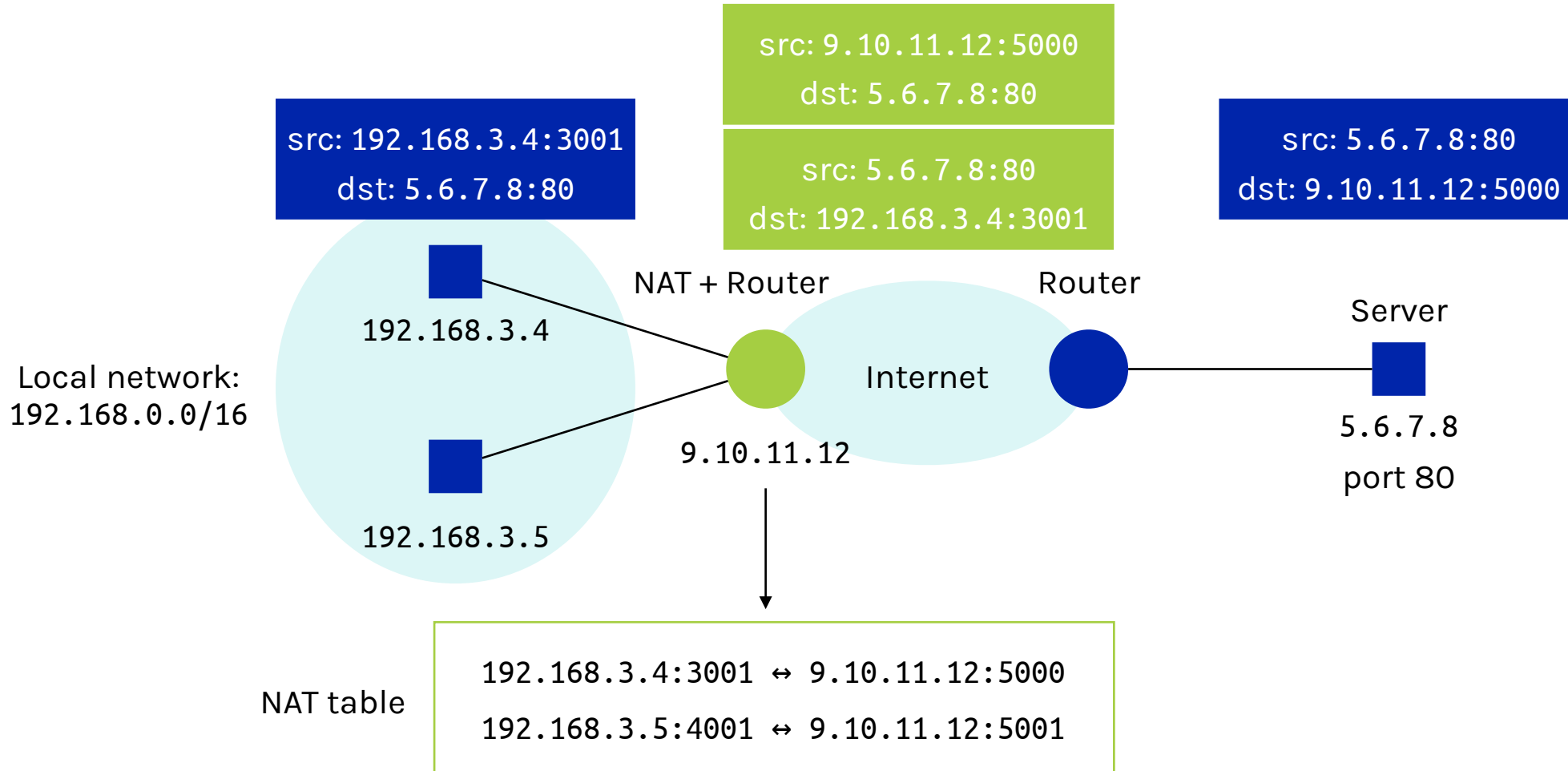
# Network Address Translation

# The Internet before NAT

**Every machine connected to the Internet had a unique IP**



src: 1.2.3.4:2001
dst: 5.6.7.8:80

src: 5.6.7.8:80
dst: 1.2.3.4:2001

1.2.3.4

Router

Router

Server

Local network:
1.2.3.0/24

Internet

5.6.7.8

port 80

1.2.3.5

# The Internet before NAT

**Every machine connected to the Internet had a unique IP**

src: 9.10.11.12:5000
dst: 5.6.7.8:80

src: 5.6.7.8:80
dst: 192.168.3.4:3001

src: 192.168.3.4:3001
dst: 5.6.7.8:80

src: 5.6.7.8:80
dst: 9.10.11.12:5000

NAT + Router

Router

Server

192.168.3.4

Internet

Local network:
192.168.0.0/16

9.10.11.12

5.6.7.8

port 80

192.168.3.5

NAT table

192.168.3.4:3001 ↔ 9.10.11.12:5000

192.168.3.5:4001 ↔ 9.10.11.12:5001

**57**

# NAT (dis)advantages

**Better privacy/anonymization**

- All hosts in one network get the same public IP

- But hosts may still be identified by cookies, browser version,…
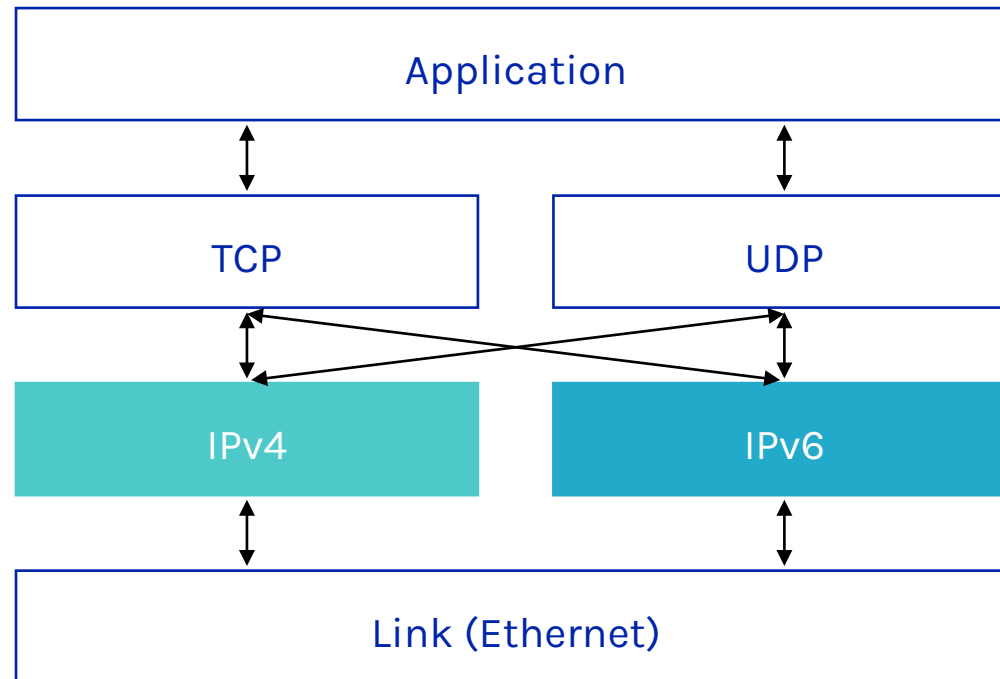
**Better security**

- From the outside you cannot directly reach the hosts

- Problematic for applications like online gaming

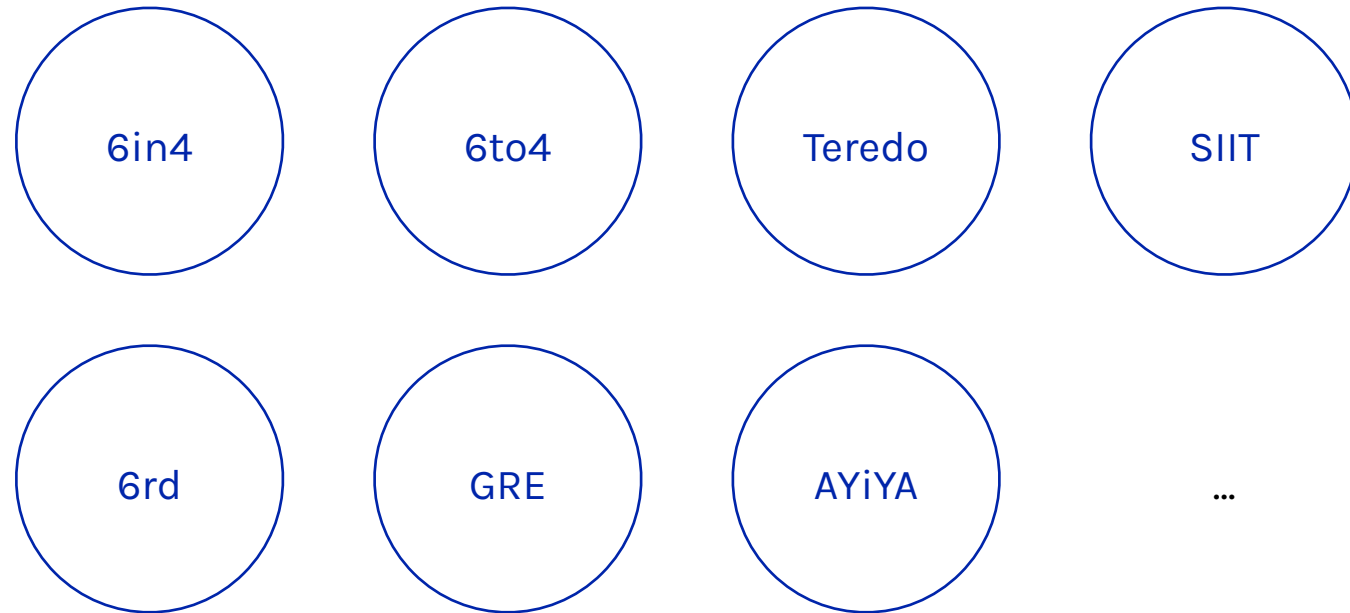**Limited scalability (limited mapping table)**

- Example: WiFi access problems in public places often due to full NAT table

# IPv4 to IPv6 transition

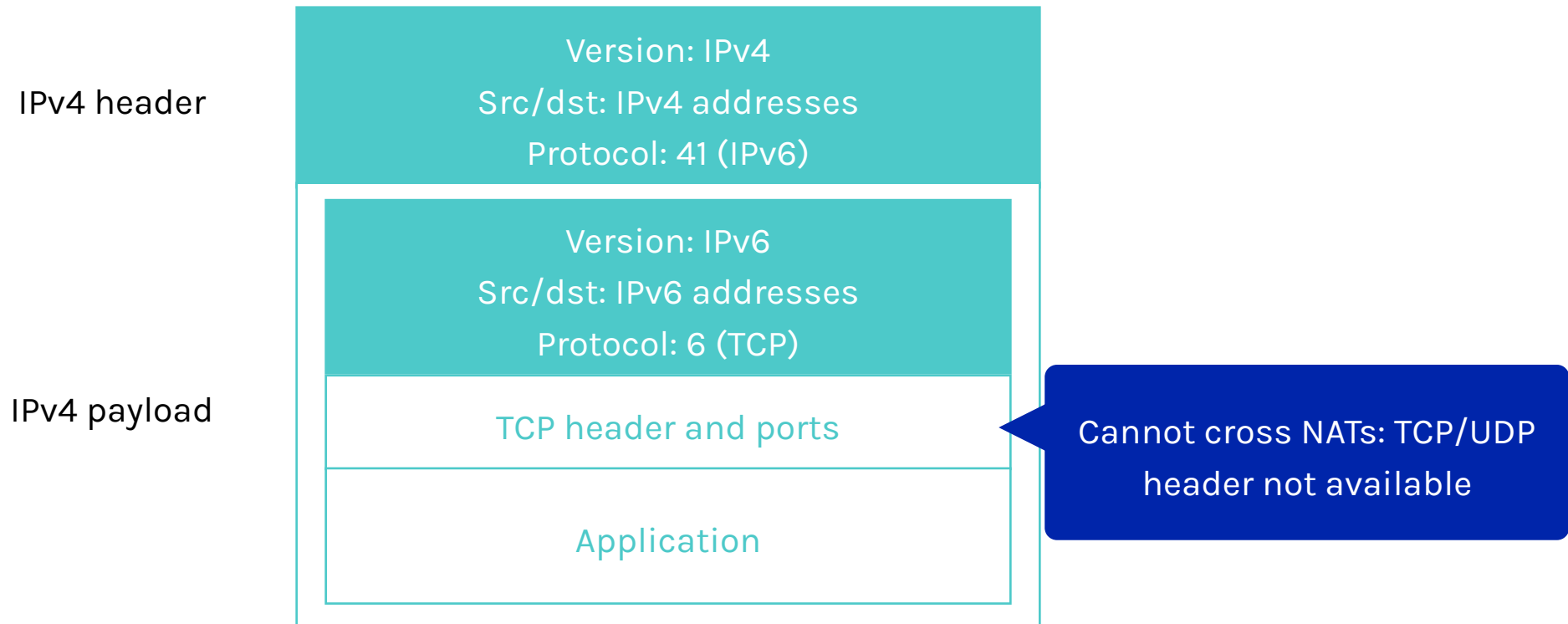**Duo stack approach used in many OSes and applications**
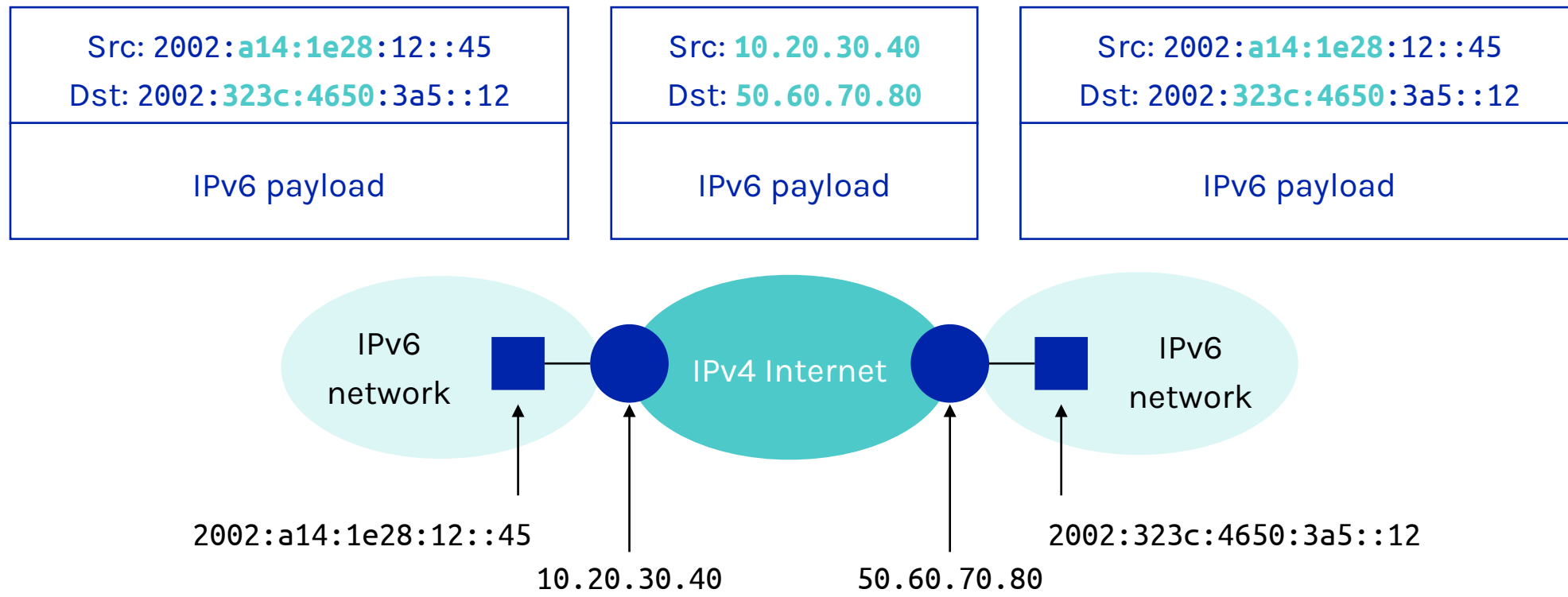
# IPv4 to IPv6 transition mechanisms

6in4

6to4

Teredo

SIIT

6rd

GRE

AYiYA

…

# 6in4

**Transmits IPv6 packets over statically configured IPv4 tunnels**

IPv4 header
> Version: IPv4
> Src/dst: IPv4 addresses
> Protocol: 41 (IPv6)

IPv4 payload
> Version: IPv6
> Src/dst: IPv6 addresses
> Protocol: 6 (TCP)
>
> TCP header and ports
>
> Application

Cannot cross NATs: TCP/UDP header not available

# 6to4

**Transmits IPv6 packets over IPv4 networks without explicit tunnels**

| |
|---|
| Src: 2002:**a14:1e28**:12::45 |
| Dst: 2002:**323c:4650**:3a5::12 |
| IPv6 payload |

| |
|---|
| Src: **10.20.30.40** |
| Dst: **50.60.70.80** |
| IPv6 payload |

| |
|---|
| Src: 2002:**a14:1e28**:12::45 |
| Dst: 2002:**323c:4650**:3a5::12 |
| IPv6 payload |

IPv6 network    IPv4 Internet    IPv6 network

2002:a14:1e28:12::45

10.20.30.40

50.60.70.80

2002:323c:4650:3a5::12

# Special IPv6 addresses in 6to4

| 16 bits | 32 bits | 16 bits | 64 bits |
|---------|---------|---------|---------|
| 2002 | **IPv4 address** | Subnet | Host ID |

IPv4       192.15.3.73

              **c0.0f.03.49**

IPv6      2002:**c00f:0349**::/49

# Virtual Private Network (VPN)

R1                R2

Network 1
(home)

Internet

Network 2
(UPB)

131.65.12.3                            131.234.5.1

| Dst = 2.x |
| --- |
| IP payload |

| Dst = 131.234.5.1 |
| --- |
| Dst = 2.x |
| IP payload 🔒 |

| Dst = 2.x |
| --- |
| IP payload |

# Helper Protocols

# Address Resolution Protocol (ARP)

**ARP query**  Whoever has the IP of `10.0.0.4`, please respond to me

**ARP response**  That is me, and here is my MAC address `88:b2:2f:54:1a:0f`



`1a:23:f9:cd:06:9b`

`10.0.0.1`

`5c:66:ab:90:75:b1`

`10.0.0.2`

`49:bd:d2:c7:56:2a`

`10.0.0.3`

`88:b2:2f:54:1a:0f`

`10.0.0.4`

ARP table on end-devices

# Internet Control Message Protocol (ICMP)

**ICMP is a companion protocol to IP**

- Sits on top of IP (IP protocol = 1)

**Provides error report and testing**

- Error is at router while forwarding
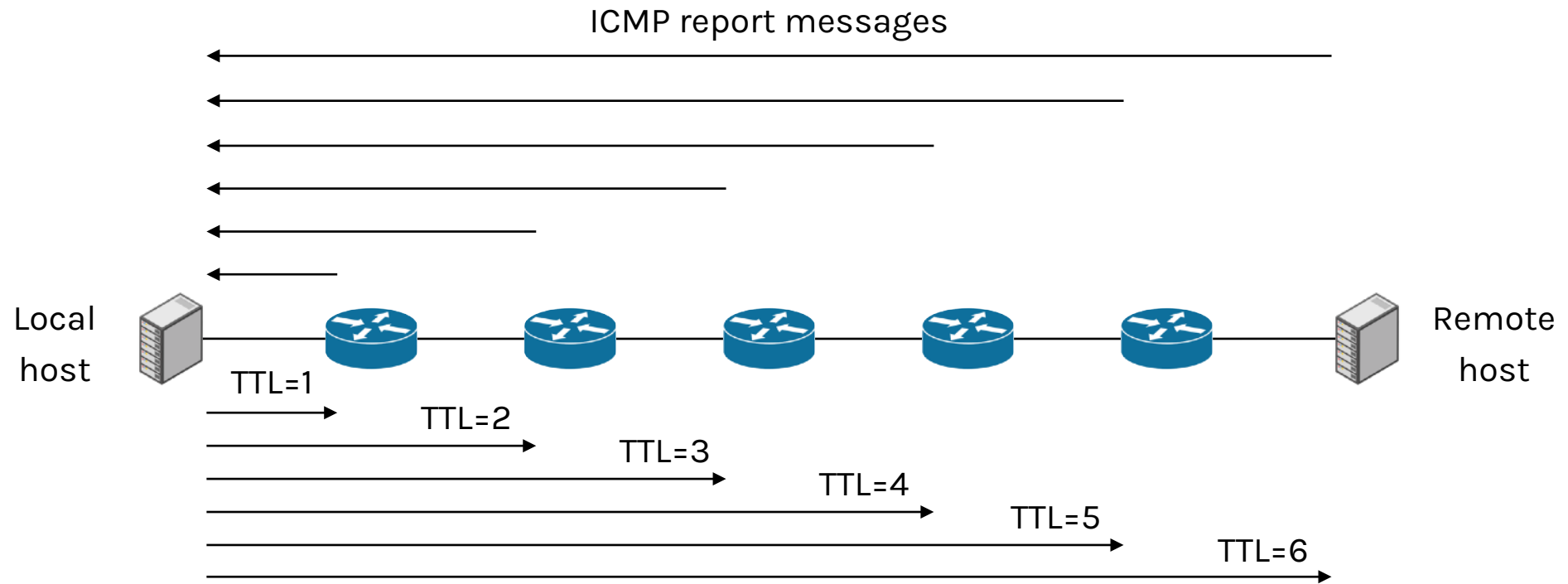
- Usual tools: ping, traceroute…

# ICMP message format

Portion of offending packet,
starting with its IP header

| Src=router, Dst=A<br>Protocol=1 | Type=X, Code=Y | Src=A, Dst=B<br>XXXXXXXX |
|---|---|---|
| IP header | ICMP header | ICMP data |

| Type / code | Name | Usage |
|---|---|---|
| 3 / 0 or 1 | Dest. unreachable (net or host) | Lack of connectivity |
| 3 / 4 | Dest. unreachable (fragment) | Path MTU discovery |
| 11 / 0 | Time exceeded (transit) | Traceroute |
| 8 or 0 / 0 | Echo request or reply | Ping (testing, not error) |

# Traceroute



ICMP report messages

Local host
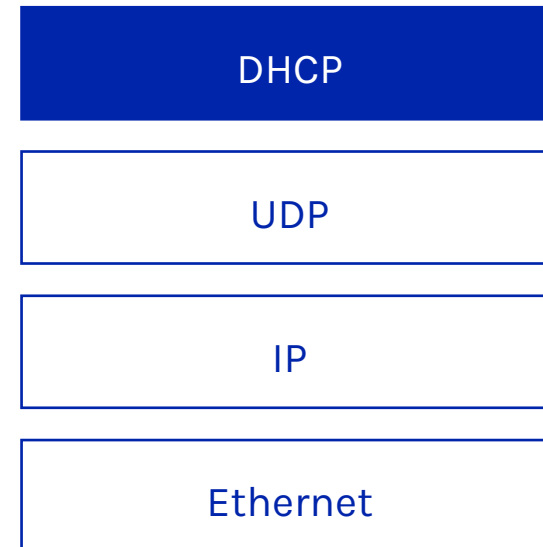
Remote host

TTL=1
TTL=2
TTL=3
TTL=4
TTL=5
TTL=6

# Dynamic Host Configuration Protocol (DHCP)

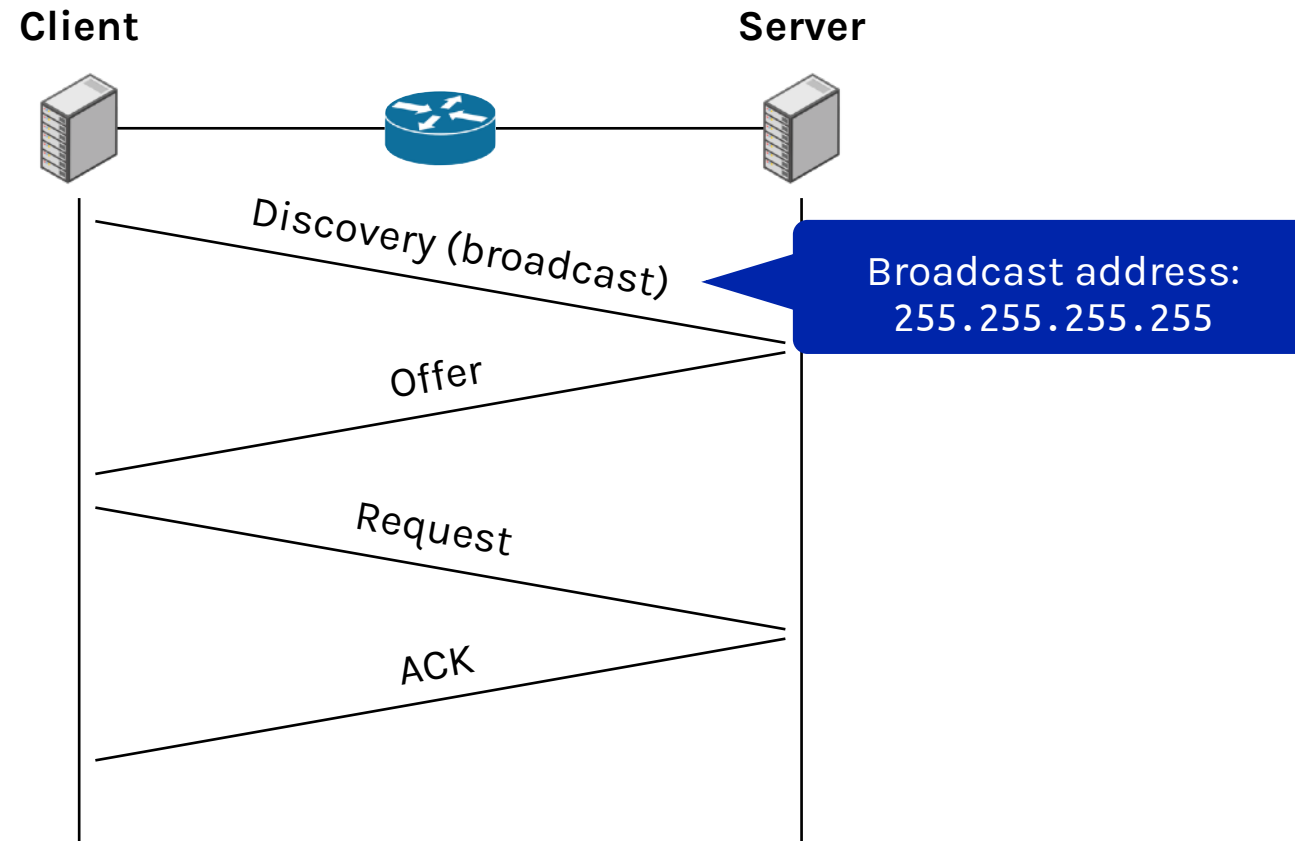**Leases IP addresses to nodes and provides other parameters**

- Network prefix

- Address of local router (gateway)

- DNS server, time server,…

**DHCP is a client-server application**

- Uses UDP ports 67, 68

| DHCP |
| --- |
| UDP |
| IP |
| Ethernet |

# DHCP messages

# Summary

## Inter-networking

- Internet narrow waist

- IP address and prefix

## IP forwarding

- Router architecture
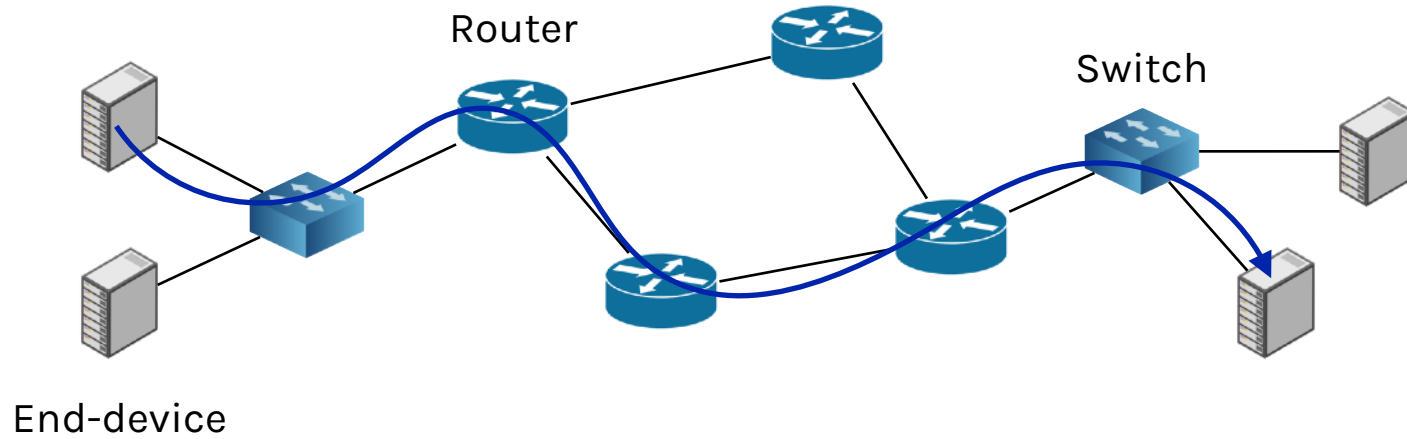
- Prefix matching

- IPv4 packet format

- IPv6

## Network Address Translation

- NAT ideas

- IPv4 to IPv6 transition

## Helper protocols

- ARP

- ICMP and traceroute

- DHCP

# Next time: network layer

Router

Switch

End-device

How to construct the routing path and
navigate through the Internet?

# Further reading material

**Andrew S. Tanenbaum, David J. Wetherall. Computer Networks (5th edition).**

- Section 5.5: Internetworking

**Larry Peterson, Bruce Davie. Computer Networks: A Systems Approach.**

- Section 3.3 Internet (IP)