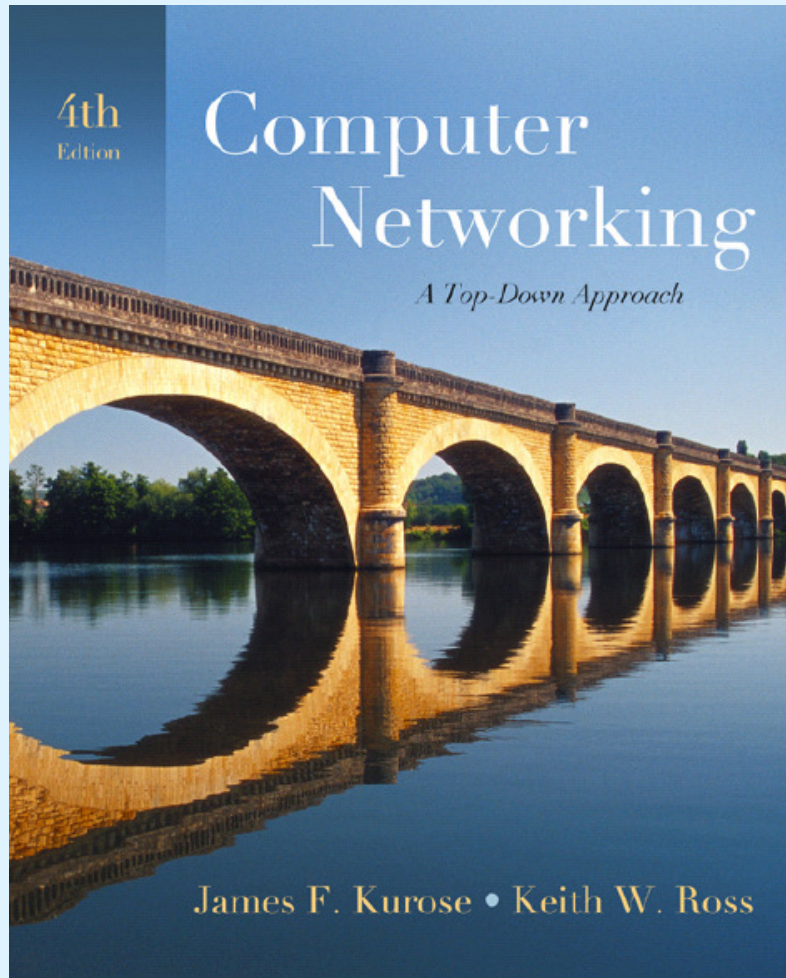# Network Layer

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

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

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

❑ If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)

❑ If you post any slides in substantially unaltered form on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.
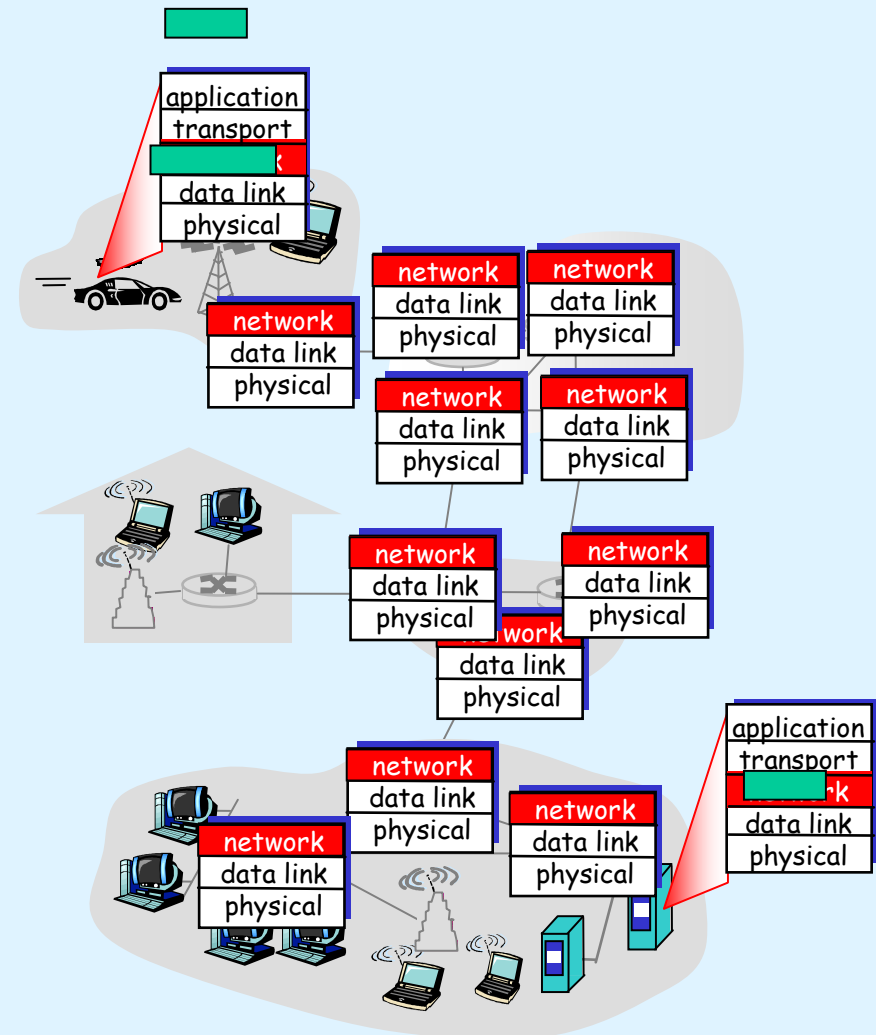
Thanks and enjoy! JFK/KWR

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

Network Layer

# Network layer
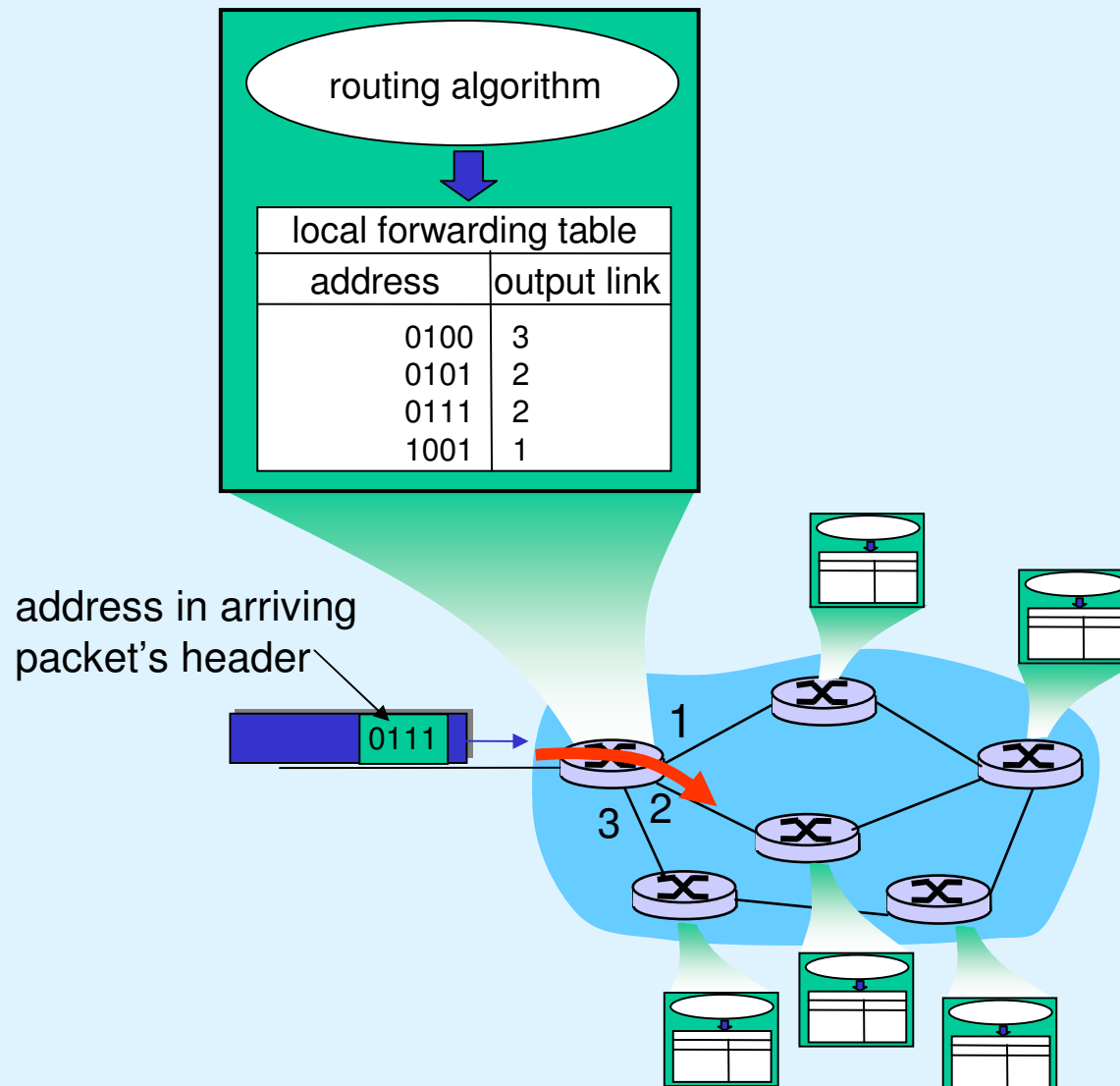
- transport segment from sending to receiving host
- on sending side encapsulates segments into datagrams
- on rcving side, delivers segments to transport layer
- network layer protocols in *every* host, router
- router examines header fields in all IP datagrams passing through it

# Two Key Network-Layer Functions

□ *forwarding:* move packets from router's input to appropriate router output

□ *routing:* determine route taken by packets from source to dest.

  ○ *routing algorithms*

# Interplay between routing and forwarding

routing algorithm

| local forwarding table | |
|---|---|
| address | output link |
| 0100 | 3 |
| 0101 | 2 |
| 0111 | 2 |
| 1001 | 1 |

address in arriving packet's header

0111

1

3  2

# Network layer connection and connection-less service

□ **datagram network provides network-layer connectionless service**

□ **VC network, e.g. ATM, provides network-layer connection service**

□ **analogous to the transport-layer services, but:**

  ○ service: host-to-host

  ○ no choice: network provides one or the other

  ○ implementation: in network core

# Forwarding table

4 billion ($2^{32}$) possible entries

Destination Address Range                                    Link Interface

11001000 00010111 00010000 00000000
through                                                                    0
11001000 00010111 00010111 11111111

11001000 00010111 00011000 00000000
through                                                                    1
11001000 00010111 00011000 11111111

11001000 00010111 00011001 00000000
through                                                                    2
11001000 00010111 00011111 11111111

otherwise                                                                  3

# Longest prefix matching

| Prefix Match | Link Interface |
|---|---|
| 11001000 00010111 00010 | 0 |
| 11001000 00010111 00011000 | 1 |
| 11001000 00010111 00011 | 2 |
| otherwise | 3 |

Examples

DA: 11001000  00010111  00010110  10100001          Which interface?

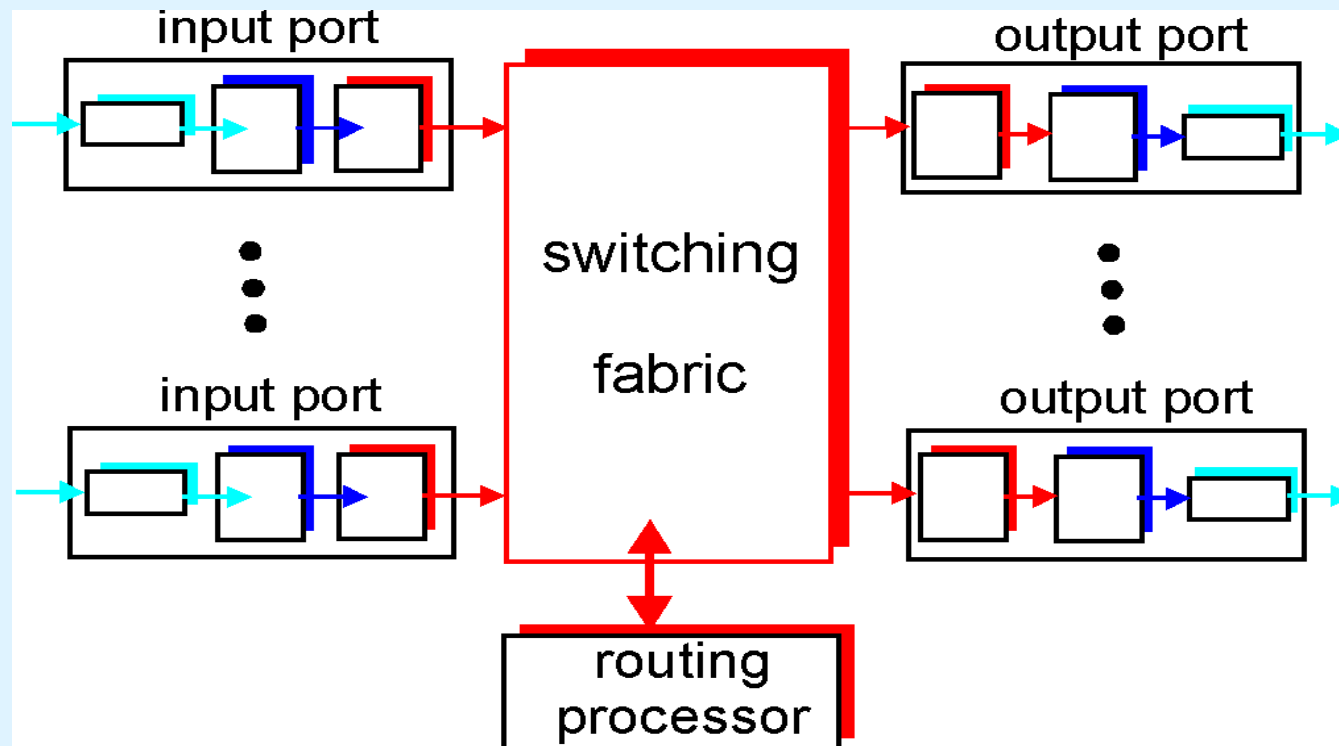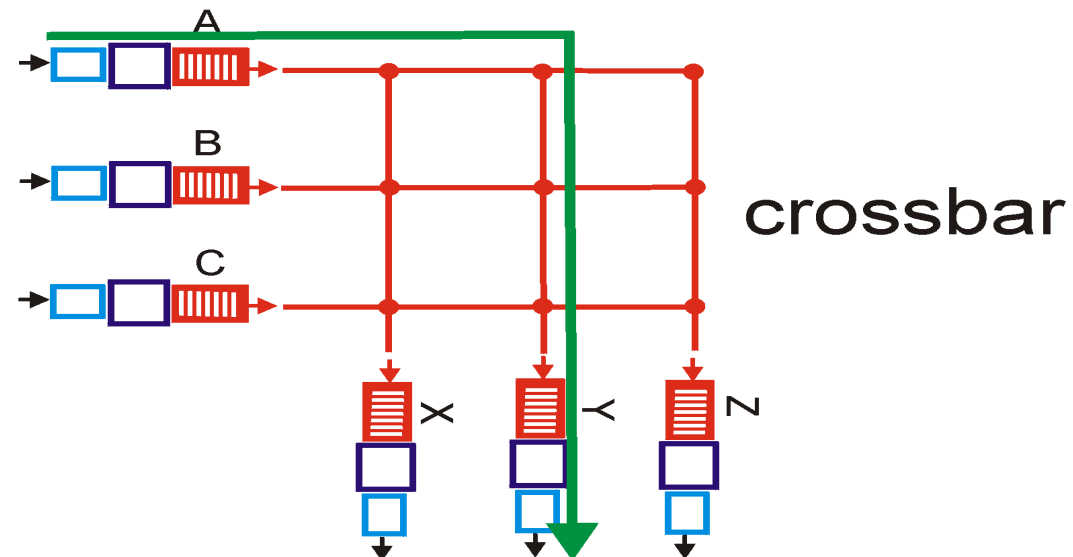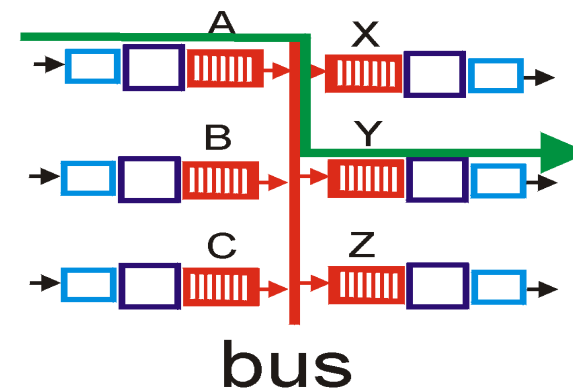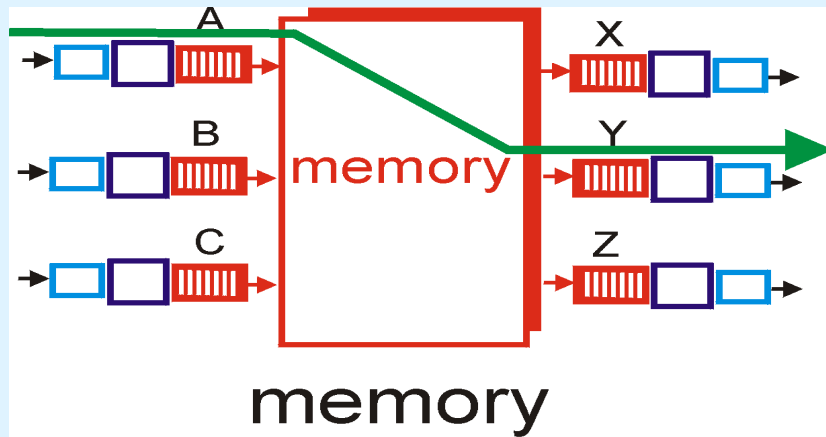DA: 11001000  00010111  00011000  10101010          Which interface?

# Router Architecture Overview

## Two key router functions:

☐ run routing algorithms/protocol (RIP, OSPF, BGP)

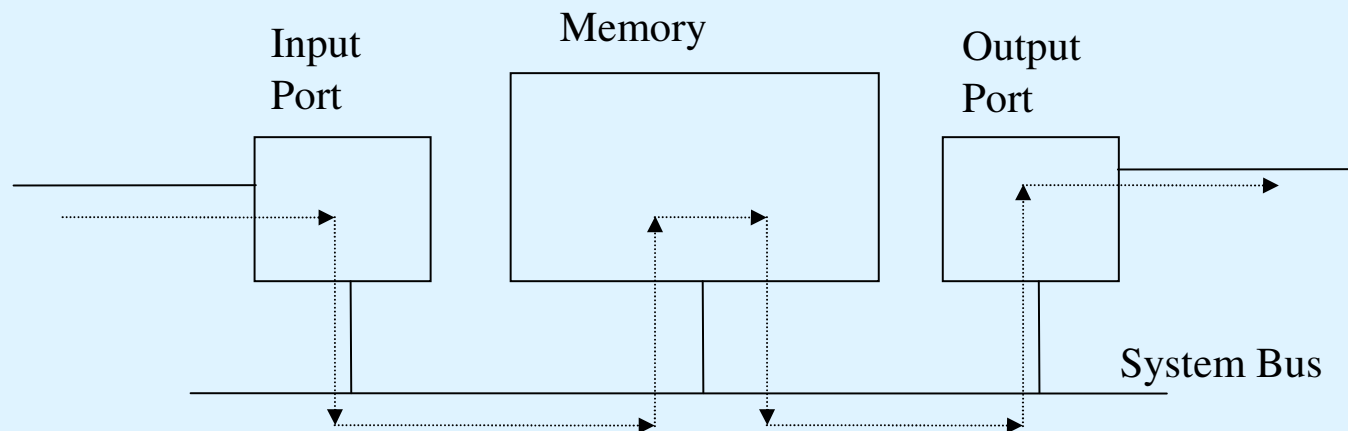☐ *forwarding* datagrams from incoming to outgoing link

# Three types of switching fabrics



memory

bus

crossbar

# Switching Via Memory

First generation routers:

□ traditional computers with switching under direct control of CPU

□ packet copied to system's memory

□ speed limited by memory bandwidth (2 bus crossings per datagram)

Input Port

Memory

Output Port

System Bus

# Switching Via a Bus


bus

□ datagram from input port memory to output port memory via a shared bus

□ bus contention: switching speed limited by bus bandwidth

□ 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

# Switching Via An Interconnection Network

❒ overcome  bus bandwidth limitations

❒ Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor

❒ advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.

❒ Cisco 12000: switches 60 Gbps through the interconnection network

# Output Ports



□ *Buffering* required when datagrams arrive from fabric faster than the transmission rate

□ *Scheduling discipline* chooses among queued datagrams for transmission

# Output port queueing



Output Port Contention at Time *t*

One Packet Time Later

❐ **buffering when arrival rate via switch exceeds output line speed**

❐ *queueing (delay) and loss due to output port buffer overflow!*

# How much buffering?

□ Rule of thumb (RFC 3439): average buffering equal to B = RTT·C

  ○ e.g., $RTT_{typical}$ = 250 msec, C = 10 Gbps link:

$$B = RTT \cdot C = 2.5 \text{ Gbit}$$

□ Recent recommendation: with $N$ flows, buffering equal to

$$B = \frac{RTT \cdot C}{\sqrt{N}}$$

# Input Port Queuing

□ Fabric slower than input ports combined -> queueing may occur at input queues

□ Head-of-the-Line (HOL) blocking: queued datagram at front of queue prevents others in queue from moving forward

□ *queueing delay and loss due to input buffer overflow!*



output port contention
at time t – only one red
packet can be transferred

green packet
experiences HOL blocking

# The Internet Network layer

Host, router network layer functions:



Transport layer: TCP, UDP

Network layer

**Routing protocols**
• path selection
• RIP, OSPF, BGP

forwarding table

**IP protocol**
• addressing conventions
• datagram format
• packet handling conventions

**ICMP protocol**
• error reporting
• router "signaling"

Link layer

physical layer

# IPv4 datagram format

IP protocol version number

header length (bytes)

"type" of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

total datagram length (bytes)

for fragmentation/ reassembly

E.g. timestamp, record route taken, specify list of routers to visit.

32 bits

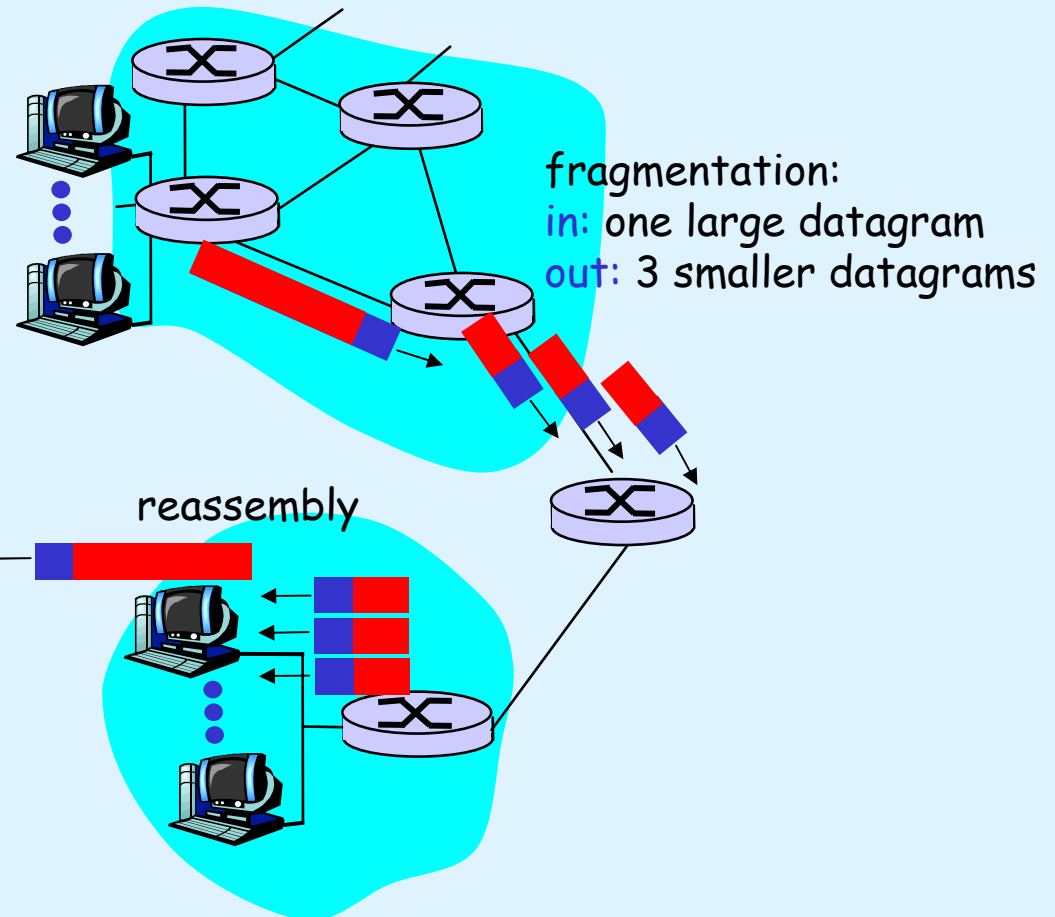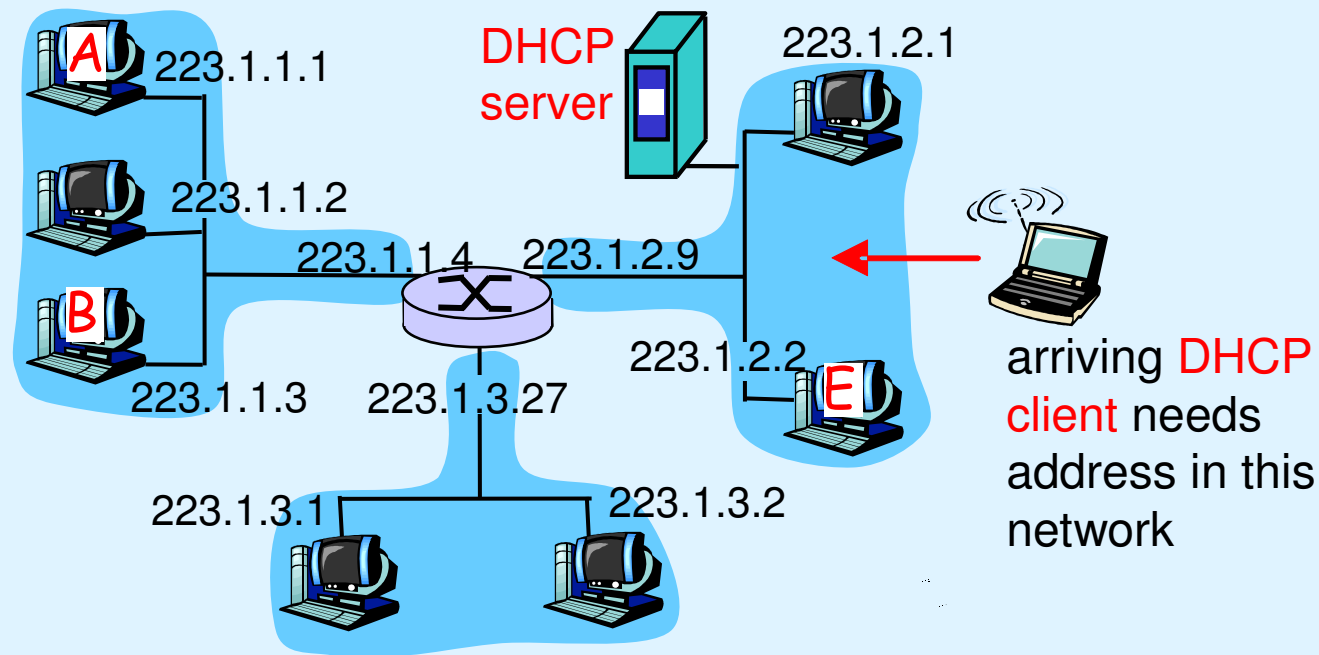| ver | head. len | type of service | length | | |
|---|---|---|---|---|---|
| 16-bit identifier | | | flgs | fragment offset | |
| time to live | upper layer | | header checksum | | |
| 32 bit source IP address | | | | | |
| 32 bit destination IP address | | | | | |
| Options (if any) | | | | | |
| data (variable length, typically a TCP or UDP segment) | | | | | |

how much overhead with TCP?

- ❑ 20 bytes of TCP
- ❑ 20 bytes of IP
- ❑ = 40 bytes + app layer overhead

# IP Fragmentation & Reassembly

- network links have MTU (max.transfer size) - largest possible link-level frame.
  - different link types, different MTUs
- large IP datagram divided ("fragmented") within net
  - one datagram becomes several datagrams
  - "reassembled" only at final destination
  - IP header bits used to identify, order related fragments

fragmentation:
in: one large datagram
out: 3 smaller datagrams

reassembly

# DHCP client-server scenario



A  223.1.1.1

223.1.1.2

B

223.1.1.3    223.1.3.27

223.1.1.4    223.1.2.9

DHCP server    223.1.2.1

223.1.2.2    E

223.1.3.1    223.1.3.2
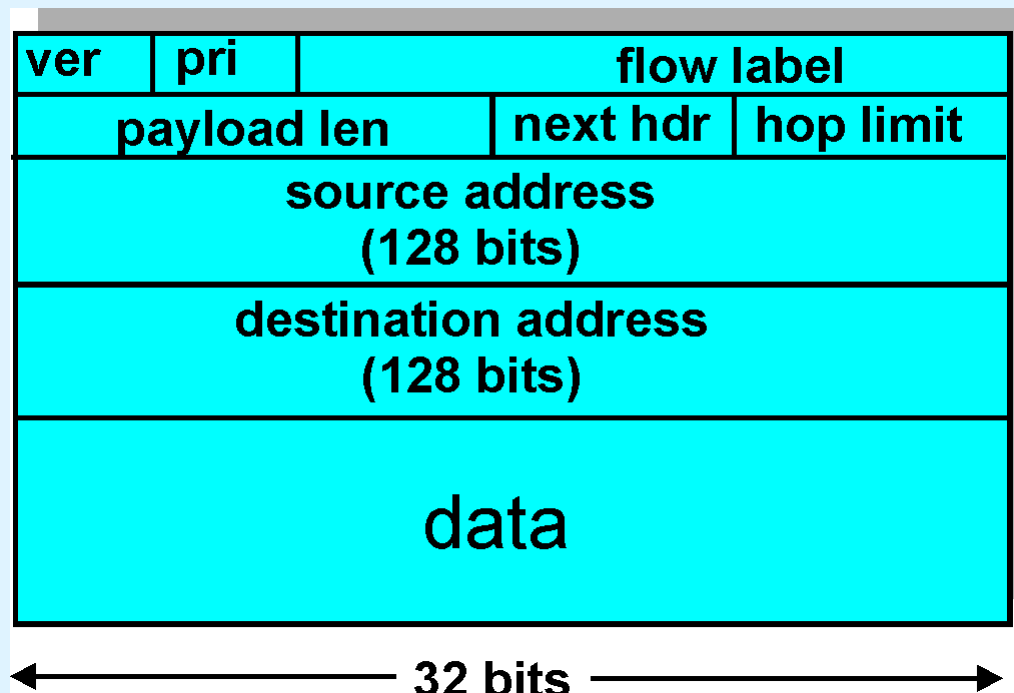
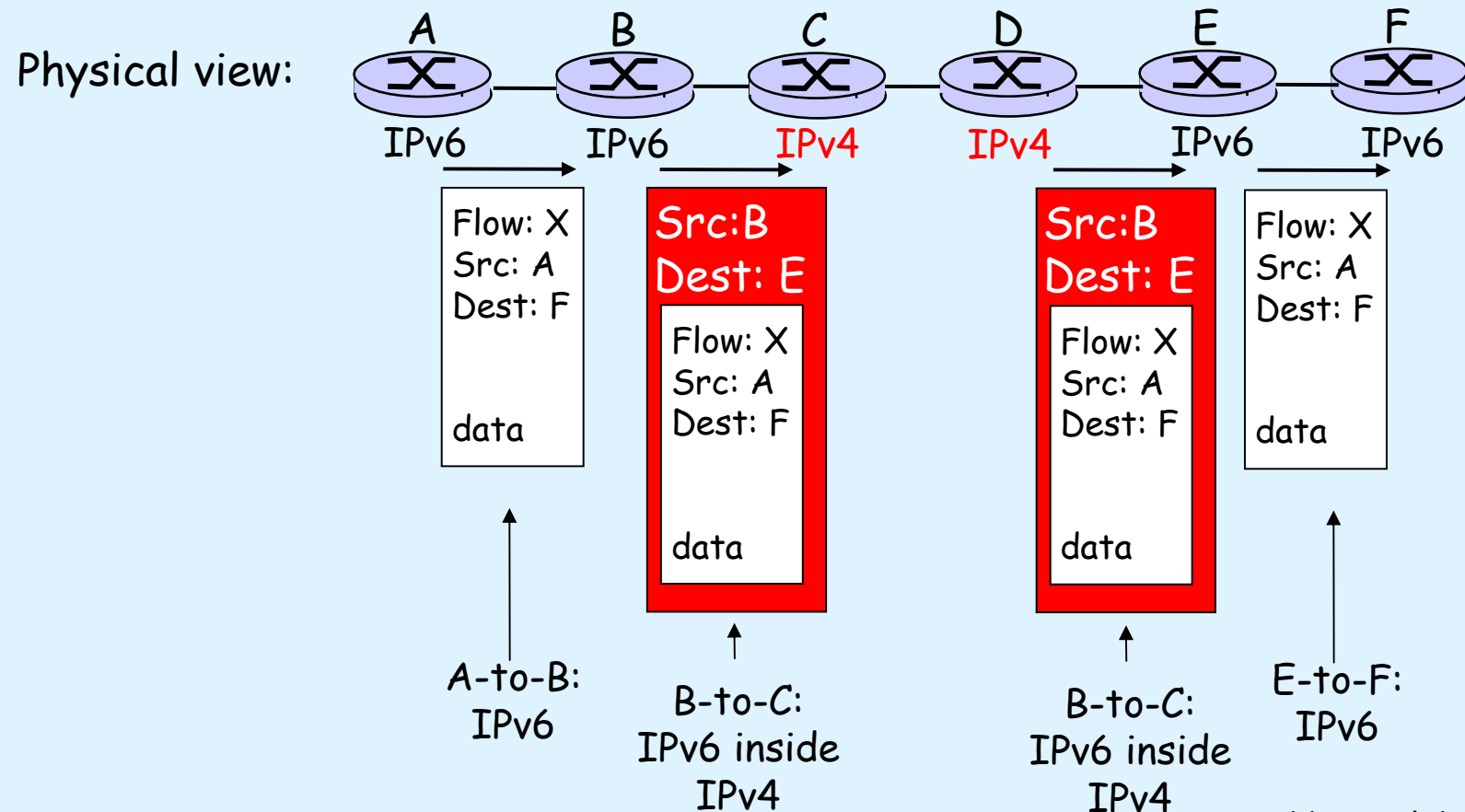arriving DHCP client needs address in this network

# IPv6 Header

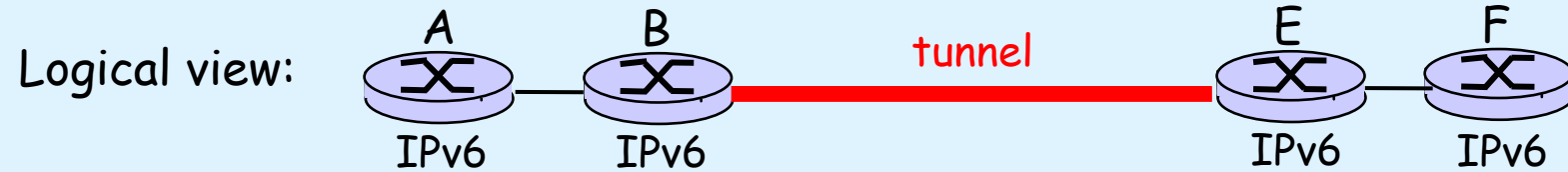*Traffic class:* 8-bit traffic class ("pri" in image).
*Flow label:* identify datagrams in same "flow"
(concept of"flow" not well defined).
*Next header:* identify upper layer protocol for data.

| ver | pri | flow label | | |
|---|---|---|---|---|
| payload len | | | next hdr | hop limit |
| source address (128 bits) | | | | |
| destination address (128 bits) | | | | |
| data | | | | |

← **32 bits** →

# Tunneling

Logical view:

A    B      tunnel      E    F

IPv6   IPv6            IPv6   IPv6

Physical view:

A    B    C    D    E    F

IPv6   IPv6   IPv4   IPv4   IPv6   IPv6

Flow: X
Src: A
Dest: F


data

Src:B
Dest: E

Flow: X
Src: A
Dest: F


data

Src:B
Dest: E

Flow: X
Src: A
Dest: F


data

Flow: X
Src: A
Dest: F


data

A-to-B:
IPv6

B-to-C:
IPv6 inside
IPv4

B-to-C:
IPv6 inside
IPv4

E-to-F:
IPv6

# Link-State Routing Algorithm (LS)

## Dijkstras algoritm

☐ Global

OBS! Algoritmen körs lokalt i varje router och är på så sätt även **decentraliserad** till skillnad mot central routing i en speciell maskin.

☐ Iterativ

- o "Billigaste väg" väljs ut till varje router.

- o Nätet spänns upp router för router (steg för steg).

- o Nästa router som läggs till det logiska nätet måste kunna nås från det nuvarande logiska nätet.

- o Routrar som inte kan nås från det nuvarande "logiska nätet" sägs ha oändlig kostnad.

- o Då alla rourar ingår i det logiska nätet är algoritmen klar.

☐ Algoritmen beskrivs med...

- o matematisk notation

- o tabell

- o grafer

☐ Används av protokollet OSPF

# Distance Vector Algorithm (DV)

**Bellman-Fords ekvation (dynamisk prog.)**

- Decentraliserad

- Iterativ

- Asynkron

  Uppdateringen görs nödvändigtvis inte i takt med andra routrar.

- Distribuerad

  - o decentraliserad

  - o annonsering

- Algoritmen beskrivs med…

  - o matematisk notation

  - o nodtabeller

  - o grafer

- Används av protokollet RIP
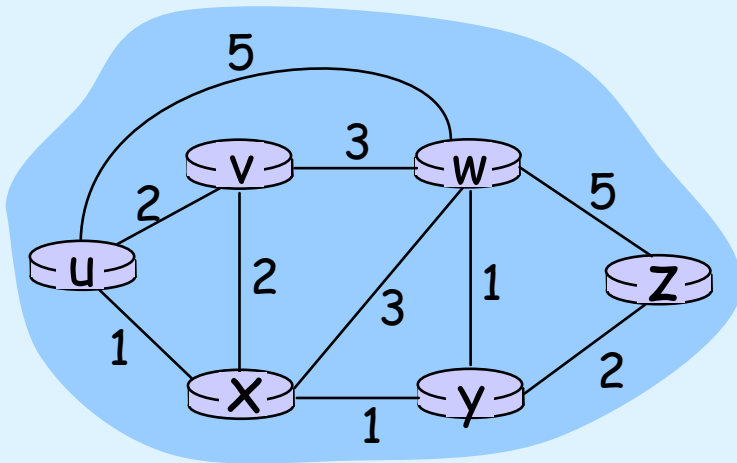
# Bellman-Ford Equation  (dynamic programming)

Define

$d_x(y)$ := cost of least-cost path from x to y

Then

$$d_x(y) = \min_v \{c(x,v) + d_v(y) \}$$

where min is taken over all neighbors v of x

# Bellman-Ford example



Clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z),$$
$$c(u,x) + d_x(z),$$
$$c(u,w) + d_w(z) \}$$
$$= \min \{2 + 5,$$
$$1 + 3,$$
$$5 + 3\} = 4$$

Node that achieves minimum is next
hop in shortest path ➜ forwarding table

# DV                                    (1)

<u>Basic idea:</u>

□ Each node periodically sends its own distance vector estimate to neighbors

□ When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$$D_x(y) \leftarrow \min_v\{c(x,v) + D_v(y)\} \quad \text{for each node } y \in N$$

□ Under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$

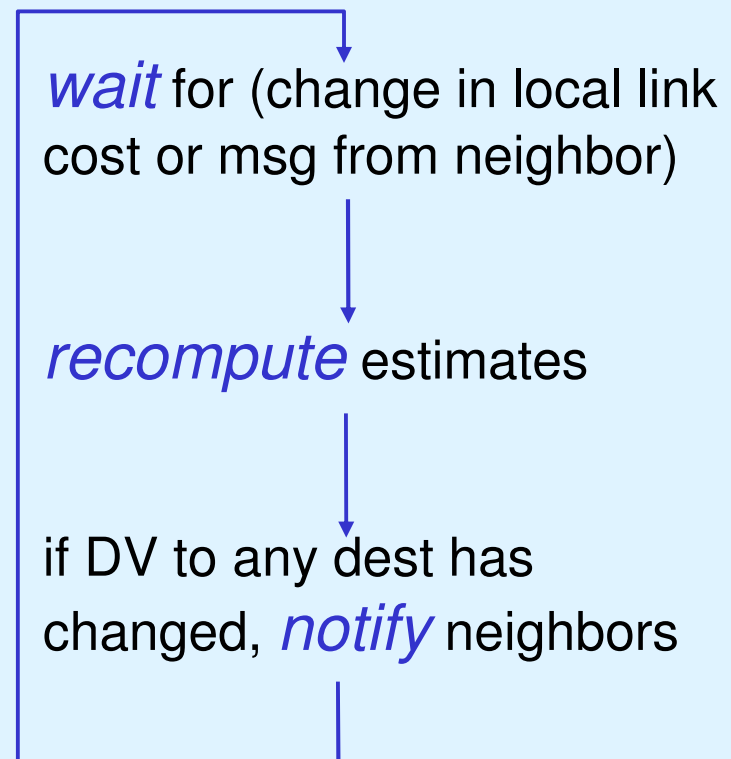# DV                                      (2)

**Iterative, asynchronous:**
each local iteration caused by:

☐  local link cost change

☐  DV update message from neighbor

**Distributed:**

☐  each node notifies neighbors *only* when its DV changes
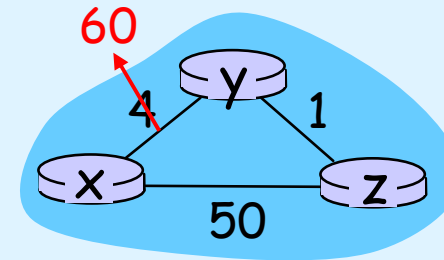
  ○  neighbors then notify their neighbors if necessary

**Each node:**

*wait* for (change in local link cost or msg from neighbor)

↓

*recompute* estimates

↓

if DV to any dest has changed, *notify* neighbors

# DV: link cost changes

## Link cost changes:

❒ good news travels fast

❒ bad news travels slow - "count to infinity" problem!
(In figure: 44 iterations before algorithm stabilizes: see text)

## Poisoned reverse:

❒ If Z routes through Y to get to X :

○ Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)

❒ will this completely solve count to infinity problem?

# Hierarchical Routing          (1)

Our routing study thus far - idealization
- all routers identical
- network "flat"

... *not* true in practice

**scale:** with 200 million destinations:
- can't store all dest's in routing tables!
- routing table exchange would swamp links!

**administrative autonomy**
- internet = network of networks
- each network admin may want to control routing in its own network

# Hierarchical Routing          (2)

- aggregate routers into regions, "autonomous systems" (AS)
- routers in same AS run same routing protocol
  - "intra-AS" routing protocol
  - routers in different AS can run different intra-AS routing protocol

Gateway router

- Direct link to router in another AS