**School of Computing Science**
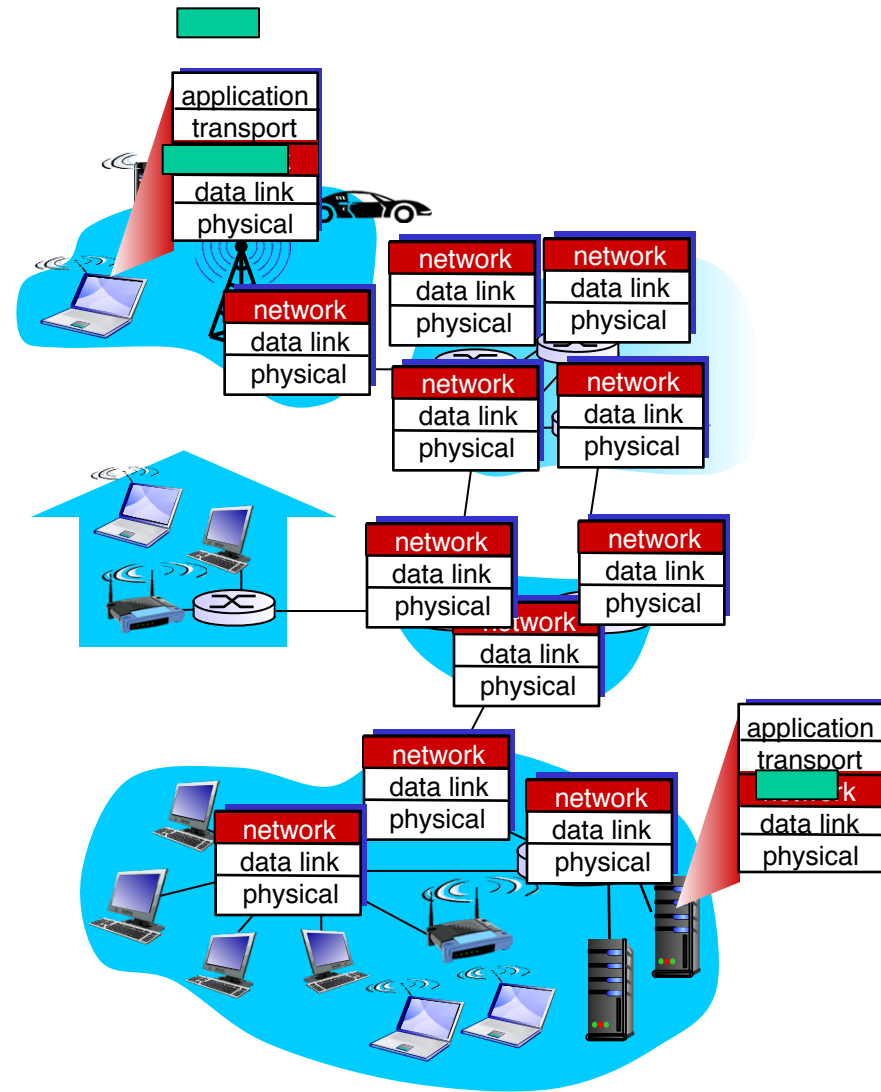**Simon Fraser University**

# CMPT 471: Networking II

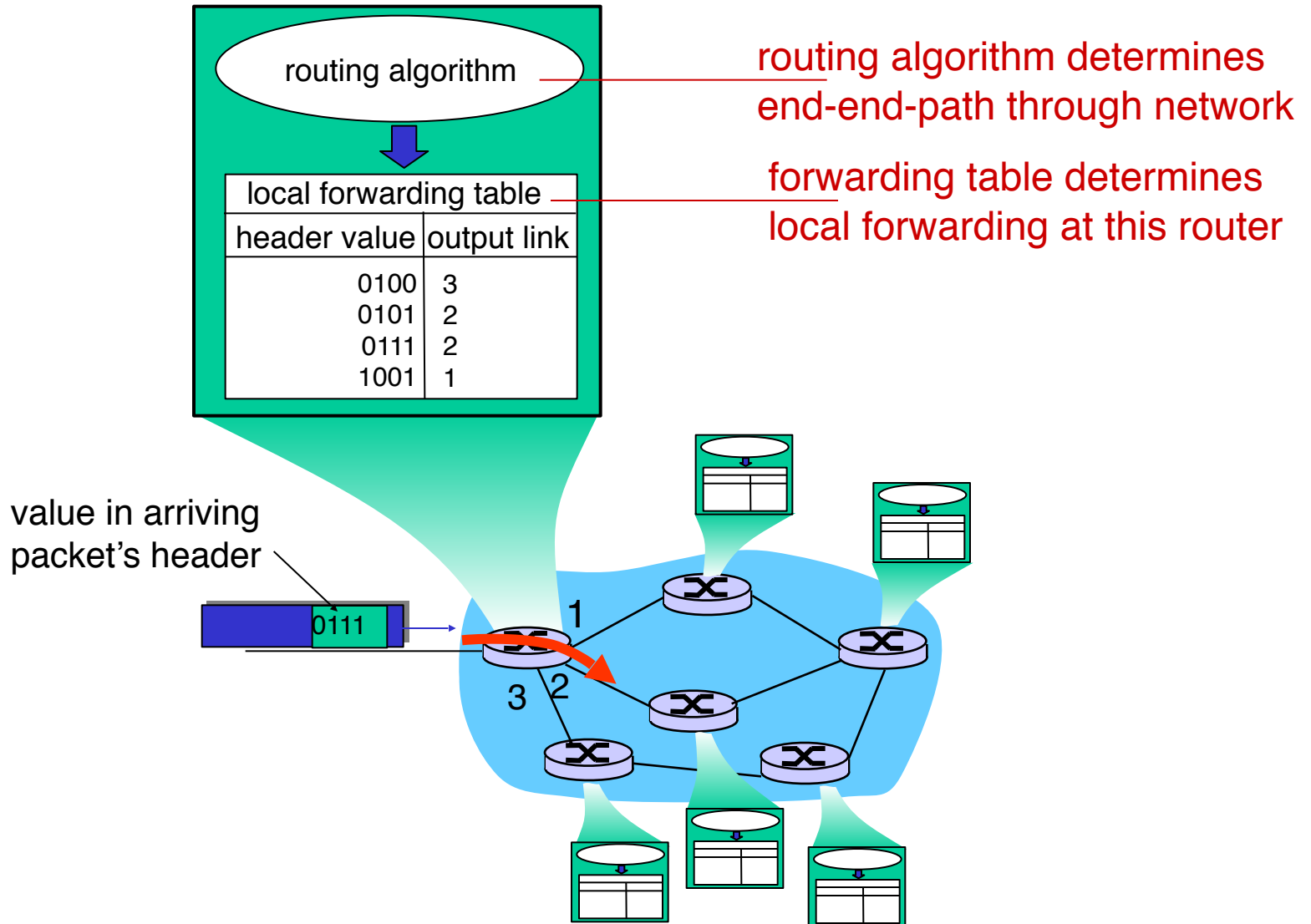# Network Layer

## Instructor: Mohamed Hefeeda

# Network layer

❖ transport segment from sending to receiving host

❖ on sending side encapsulates segments into datagrams

❖ on receiving side, delivers segments to transport layer

❖ network layer protocols in *every* host, router

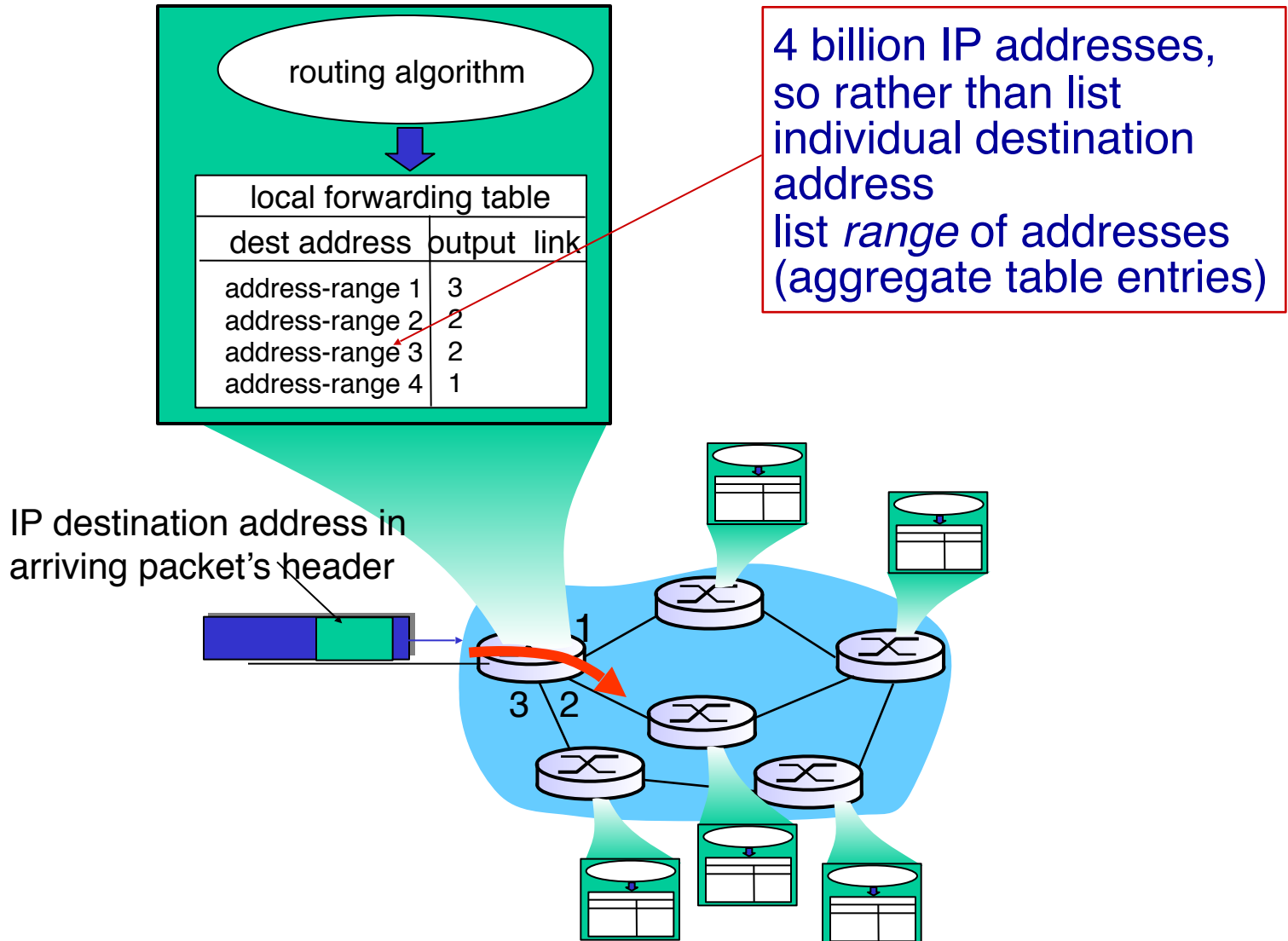❖ router examines header fields in all IP datagrams passing through it

# 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 | |
|---|---|
| header value | output link |
| 0100 | 3 |
| 0101 | 2 |
| 0111 | 2 |
| 1001 | 1 |

routing algorithm determines
end-end-path through network

forwarding table determines
local forwarding at this router

value in arriving
packet's header

0111

1

3  2

# Datagram forwarding  table

routing algorithm

local forwarding table

| dest address | output  link |
|---|---|
| address-range 1 | 3 |
| address-range 2 | 2 |
| address-range 3 | 2 |
| address-range 4 | 1 |

4 billion IP addresses, so rather than list individual destination address
list *range* of addresses (aggregate table entries)

IP destination address in arriving packet's header

1

3  2

# Longest prefix matching

*longest prefix matching*

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

| Destination Address Range | Link interface |
|---|---|
| 11001000 00010111 00010*** ******** | 0 |
| 11001000 00010111 00011000 ******** | 1 |
| 11001000 00010111 00011*** ******** | 2 |
| otherwise | 3 |

examples:
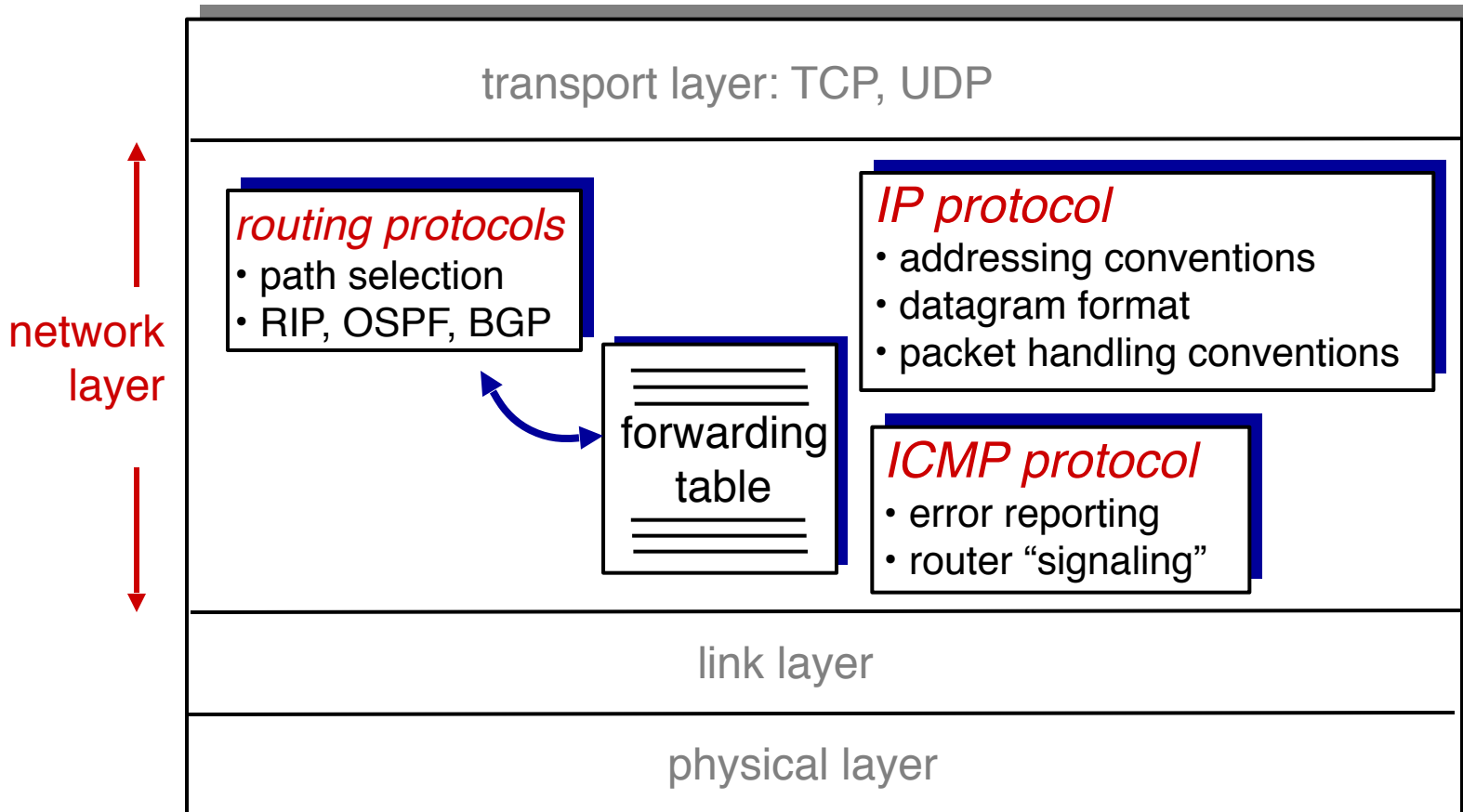
DA: 11001000  00010111  00010110  10100001      which interface?
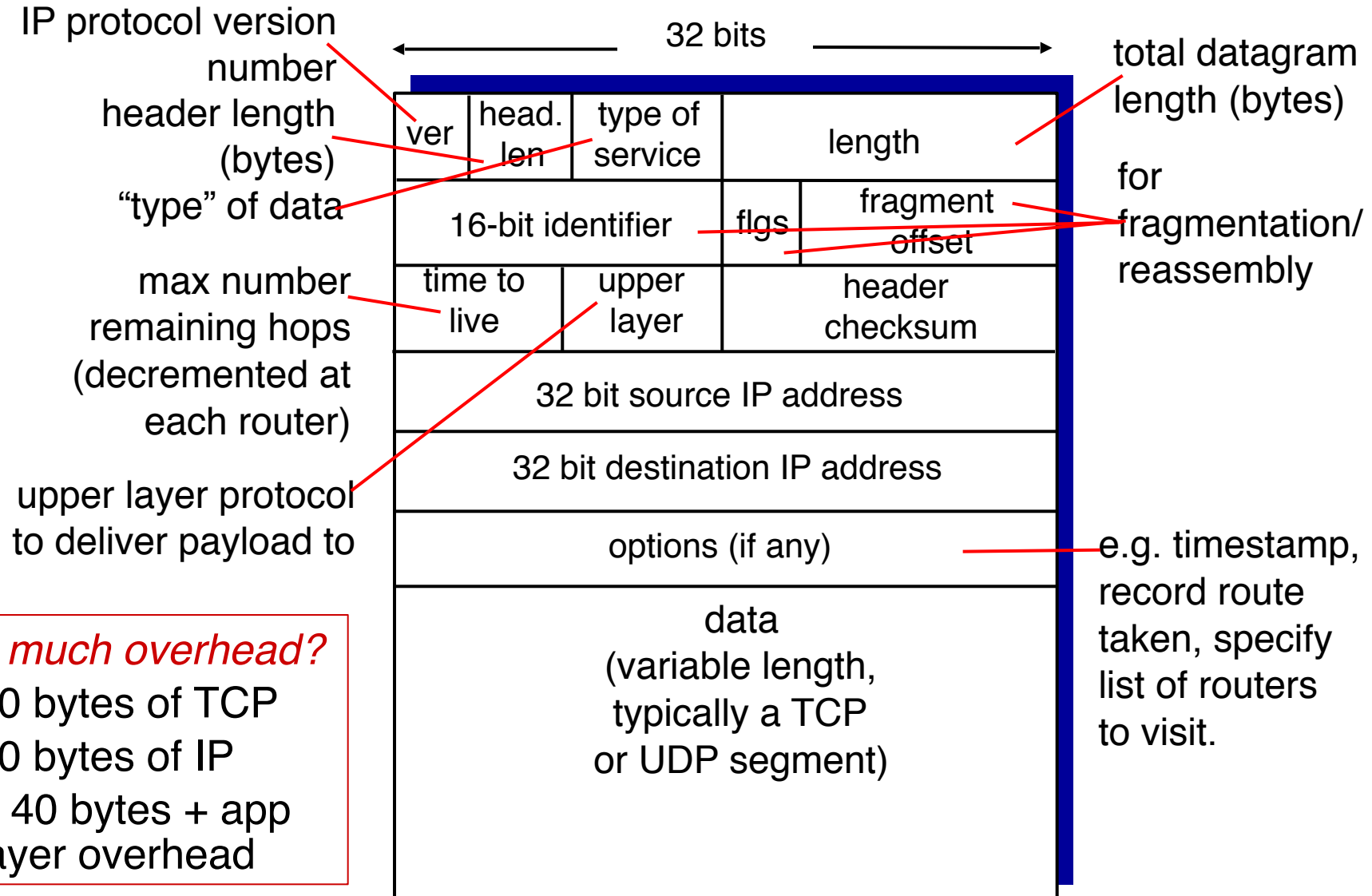
DA: 11001000  00010111  00011000  10101010      which interface?

# The Internet network layer
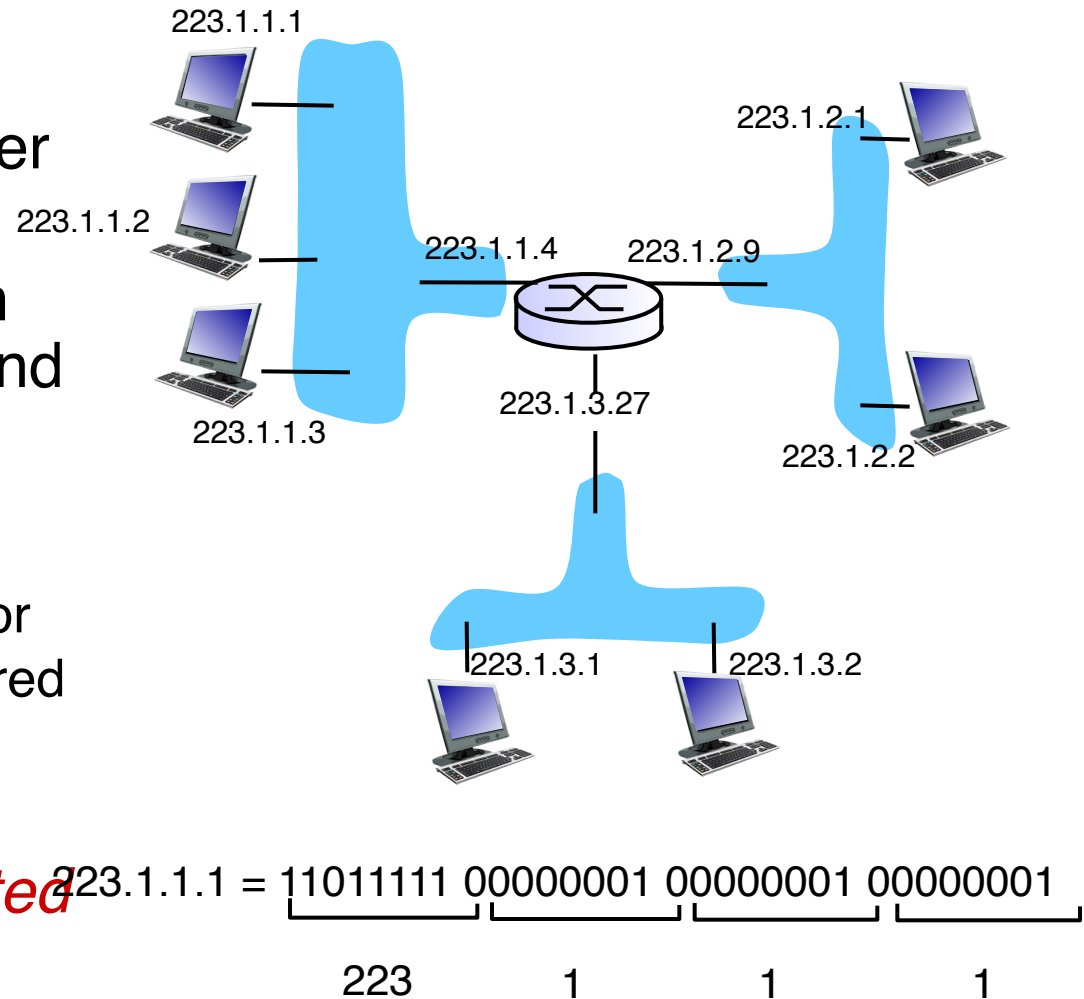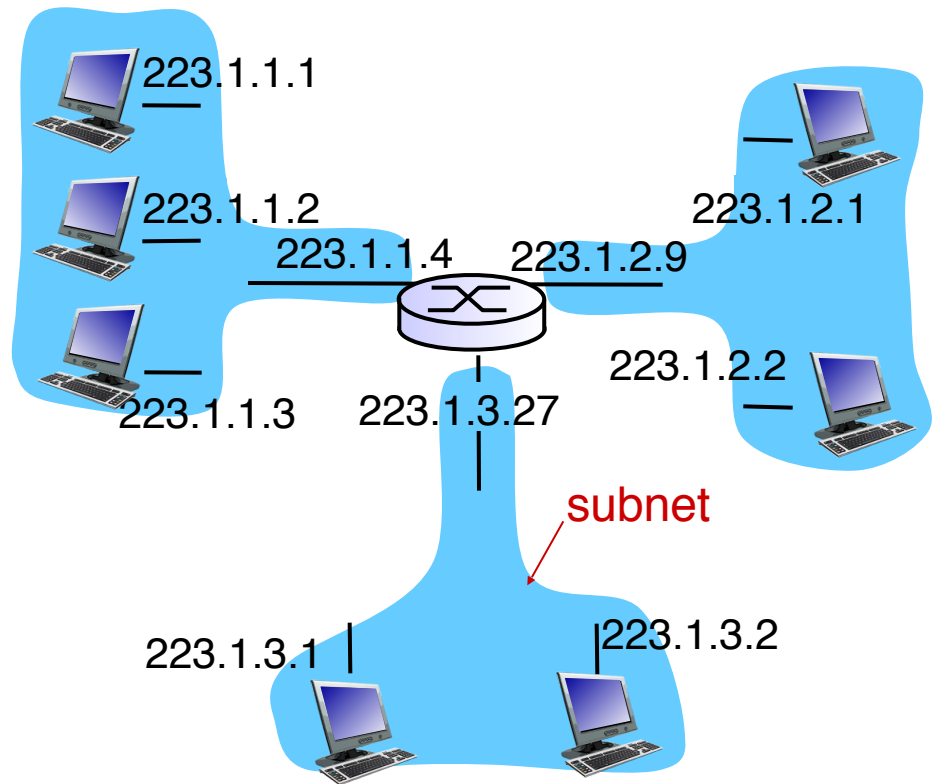
host, router network layer functions:



- **network layer** (vertical span indicator)
- transport layer: TCP, UDP
  - *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

# IP 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

**how much overhead?**
- ❖ 20 bytes of TCP
- ❖ 20 bytes of IP
- ❖ = 40 bytes + app layer overhead

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) | | | | |

total datagram length (bytes)

for fragmentation/ reassembly

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

# IP addressing: introduction

❖ *IP address:* 32-bit identifier for host, router *interface*

❖ *interface:* connection between host/router and physical link
  ▪ router's typically have multiple interfaces
  ▪ host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)

❖ *IP addresses associated with each interface*

223.1.1.1

223.1.1.2

223.1.1.4      223.1.2.9

223.1.2.1

223.1.3.27

223.1.1.3

223.1.2.2

223.1.3.1      223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001

        223        1        1        1

# Subnets

❖ IP address:
  ▪ subnet part - high order bits
  ▪ host part - low order bits

❖ *what's a subnet ?*
  ▪ device <u>interfaces</u> with same subnet part of IP address
  ▪ can physically reach each other *without intervening router*



network consisting of 3 subnets

# Subnets

## how many?
❖ 6

*recipe*

❖ to determine the subnets, detach each interface from its host or router, creating islands of isolated networks

❖ each isolated network is called a *subnet*



223.1.1.2

223.1.1.1

223.1.1.4

223.1.1.3

223.1.9.2    223.1.7.0

223.1.9.1

223.1.8.1    223.1.8.0

223.1.7.1

223.1.2.6

223.1.3.27

223.1.2.1    223.1.2.2    223.1.3.1    223.1.3.2

# IP addressing: CIDR

CIDR: Classless InterDomain Routing
  - subnet portion of address of arbitary length
  - address format: a.b.c.d/x, where x is # bits in subnet portion of address (called mask)



subnet part ⟵——————————————⟶  host part ⟵——⟶

11001000  00010111  00010000  00000000

200.23.16.0/23

# ICMP: internet control message protocol

❖ used by hosts & routers to communicate network-level information
  ▪ error reporting: unreachable host, network, port, protocol
  ▪ echo request/reply (used by ping)

❖ network-layer "above" IP:
  ▪ ICMP msgs carried in IP datagrams

❖ ICMP message: type, code plus first 8 bytes of IP datagram causing error

| Type | Code | description |
|------|------|-------------|
| 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 |

# MAC addresses and ARP

❖ 32-bit IP address:
  ▪ *network-layer* address for interface
  ▪ used for layer 3 (network layer) forwarding

❖ MAC (or LAN or physical or Ethernet) address:
  ▪ *used 'locally" to get frame from one interface to another physically-connected interface (same network, in IP-addressing sense)*
  ▪ 48 bit MAC address (for most LANs) burned in NIC ROM, also sometimes software settable
  ▪ e.g.: 1A-2F-BB-76-09-AD

  hexadecimal (base 16) notation
  (each "number" represents 4 bits)

# ARP: address resolution protocol

**ARP:** Maps IP address to MAC address

**ARP table:** each IP node (host, router) on LAN has table

- IP/MAC address mappings for some LAN nodes:

< IP address; MAC address; TTL>

- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

137.196.7.78

1A-2F-BB-76-09-AD

137.196.7.23

137.196.7.14

LAN

71-65-F7-2B-08-53

58-23-D7-FA-20-B0

0C-C4-11-6F-E3-98

137.196.7.88

# ARP protocol: same LAN

❖ A wants to send datagram to B
   ▪ B's MAC address not in A's ARP table.

❖ A broadcasts ARP query packet, containing B's IP address
   ▪ dest MAC address = FF-FF-FF-FF-FF-FF
   ▪ all nodes on LAN receive ARP query

❖ B receives ARP packet, replies to A with its (B's) MAC address
   ▪ frame sent to A's MAC address (unicast)

❖ A caches (saves) IP-to-MAC address pair in its ARP table until information becomes old (times out)
   ▪ soft state: information that times out (goes away) unless refreshed

❖ ARP is "plug-and-play":
   ▪ nodes create their ARP tables *without* *intervention from net administrator*

# IPv6: motivation

❖ *initial motivation:* 32-bit address space soon to be completely allocated.

❖ additional motivation:
- header format helps speed processing/forwarding
- header changes to facilitate QoS

*IPv6 datagram format:*
- fixed-length 40 byte header
- no fragmentation allowed

# IPv6 datagram format

*priority:* identify priority among datagrams in flow
*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 →

# Other changes from IPv4

* *checksum*: removed entirely to reduce processing time at each hop
* *options:* allowed, but outside of header, indicated by "Next Header" field
* *ICMPv6:* new version of ICMP
  * additional message types, e.g. "Packet Too Big"
  * multicast group management functions

# Transition from IPv4 to IPv6

❖ not all routers can be upgraded simultaneously
  ▪ no "flag days"
  ▪ how will network operate with mixed IPv4 and IPv6 routers?
❖ *tunneling:* IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers

IPv4 header fields

IPv4 source, dest addr

IPv6 header fields

IPv6 source dest addr

UDP/TCP payload

IPv4 payload

IPv6 datagram

IPv4 datagram

# Tunneling

logical view:

A  
IPv6

B  
IPv6

*IPv4 tunnel*  
*connecting IPv6 routers*

E  
IPv6

F  
IPv6

physical view:

A  
IPv6

B  
IPv6

C  
IPv4

D  
IPv4

E  
IPv6

F  
IPv6

# Tunneling

logical view:

A      B         *IPv4 tunnel*         E      F

*connecting IPv6 routers*

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

# Routing Algorithms

❖ **Needed to populate forwarding tables**

❖ **They run in "control plane"**
  - Typically invoked in the order of 10s of seconds or whenever a change in network topology happens

  - They are much slower than forwarding algorithms that run in "control plane" at "wire speed" (micro/nano seconds)

# Routing Algorithms

❖ **Problem solved by routing algorithms:**

**Find optimal path between any two points in the network (graph)**

- ➔ use graph algorithms (shortest path)

❖ **If Network = Internet ➔ huge graph**
  - And sub graphs (sub nets) controlled by different entities

❖ **How do we solve this problem?**

# Hierarchical Routing

❖ **Solve routing problem in two levels:**

❖ **Intra AS (Autonomous System)**
  - **Use any algorithm, based on admin**
  - **graph algorithms**

❖ **Inter-ASes**
  - **Use global, standard, routing (BGP)**

❖ **Forwarding tables are set by both:**
  - intra-AS ➔ sets entries for internal destinations
  - inter-AS & intra-AS sets entries for external destinations

3c · 3a · 3b · AS3

2c · 2a · 2b · AS2

1c · 1a · 1b · 1d · AS1

Intra-AS Routing algorithm

Inter-AS Routing algorithm

Forwarding table

# Intra-AS Routing

❖ also known as *interior gateway protocols (IGP)*

❖ most common intra-AS routing protocols:
- **RIP:** Routing Information Protocol
  - Distance vector, Bellman-Ford algorithm, distributed
  - Old and small networks
- **OSPF:** Open Shortest Path First
  - Link state, Dijkstra's algorithm, centralized
  - Most current networks
- **IGRP:** Interior Gateway Routing Protocol
  - Cisco proprietary

# Internet inter-AS routing: BGP

❖ BGP (Border Gateway Protocol): *the* de facto inter-domain routing protocol
   ▪ "glue that holds the Internet together"

❖ BGP provides each AS a means to:
   ▪ eBGP: obtain subnet reachability information from neighboring ASs.
   ▪ iBGP: propagate reachability information to all AS-internal routers.
   ▪ determine "good" routes to other networks based on reachability information and policy.

❖ allows subnet to advertise its existence to rest of Internet: *"I am here"*

# BGP basics

❖ **BGP session:** two BGP routers ("peers") exchange BGP messages:
  - advertising *paths* to different destination network prefixes ("path vector" protocol)
  - exchanged over semi-permanent TCP connections

❖ when AS3 advertises a prefix to AS1:
  - AS3 *promises* it will forward datagrams towards that prefix
  - AS3 can aggregate prefixes in its advertisement

# BGP basics: distributing path information

- ❖ using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
  - ▪ 1c can then use iBGP do distribute new prefix info to all routers in AS1
  - ▪ 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session
- ❖ when router learns of new prefix, it creates entry for prefix in its forwarding table.



——————— eBGP session

– – – – – iBGP session

AS3

AS1

AS2

other networks

other networks

# BGP route selection

❖ router may learn about more than 1 route to destination AS, selects route based on:
1. local preference value attribute: policy decision
2. shortest AS-PATH
3. closest NEXT-HOP router: hot potato routing
4. additional criteria

# Why different Intra-, Inter-AS routing ?

*policy:*
- ❖ inter-AS: admin wants control over how its traffic routed, who routes through its net.
- ❖ intra-AS: single admin, so no policy decisions needed

*scale:*
- ❖ hierarchical routing saves table size, reduced update traffic

*performance:*
- ❖ intra-AS: can focus on performance
- ❖ inter-AS: policy may dominate over performance

# Multicast/Broadcast routing

❖ **Broadcast: deliver packets from source to all other nodes**
❖ **Multicast: deliver packets to subset of nodes**

❖ source duplication is inefficient:



source
duplication

in-network
duplication

# Multicast routing: problem statement

*goal:* find a tree (or trees) connecting routers having local mcast group members

❖ **Two approaches:**

❖ *shared-tree:* same tree used by all group members

❖ *source-based:* different tree from each sender to rcvrs



shared tree



source-based trees

*legend*

group member

not group member

router with a group member

router without group member

# Approaches for building mcast trees

approaches:

❖ *source-based tree:* one tree per source
  - shortest path trees
  - reverse path forwarding

❖ *group-shared tree:* group uses one tree
  - minimal spanning (Steiner)
  - center-based trees

…we first look at basic approaches, then specific protocols adopting these approaches

# Shortest path tree

❖ mcast forwarding tree: tree of shortest path routes from source to all receivers
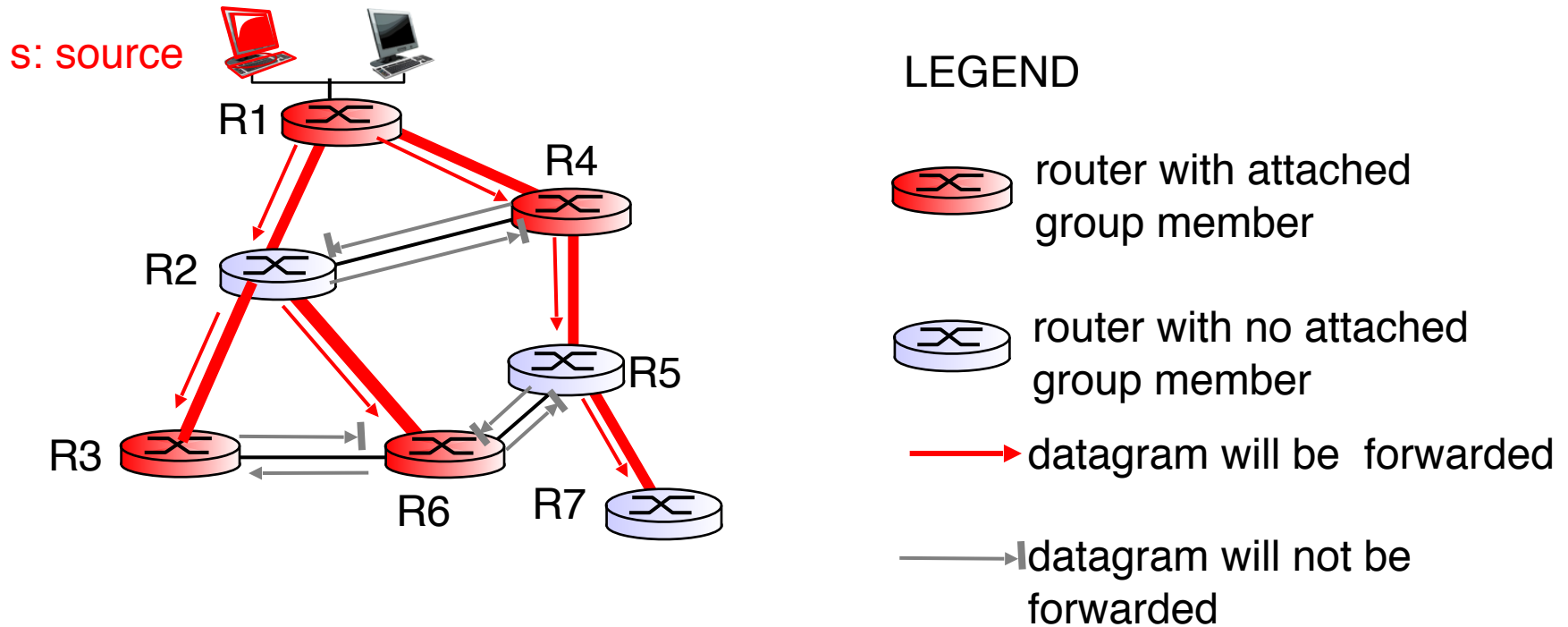  ▪ Dijkstra's algorithm

s: source

LEGEND

router with attached group member

router with no attached group member

(i) link used for forwarding, i indicates order link added by algorithm

R1 1 2
R4
R2 5
3 4
R5
R3 6
R6 R7

# Reverse path forwarding

❖ rely on router's knowledge of unicast shortest path from it  to sender

❖ each router has simple forwarding behavior:

*if* (mcast datagram received on incoming link on
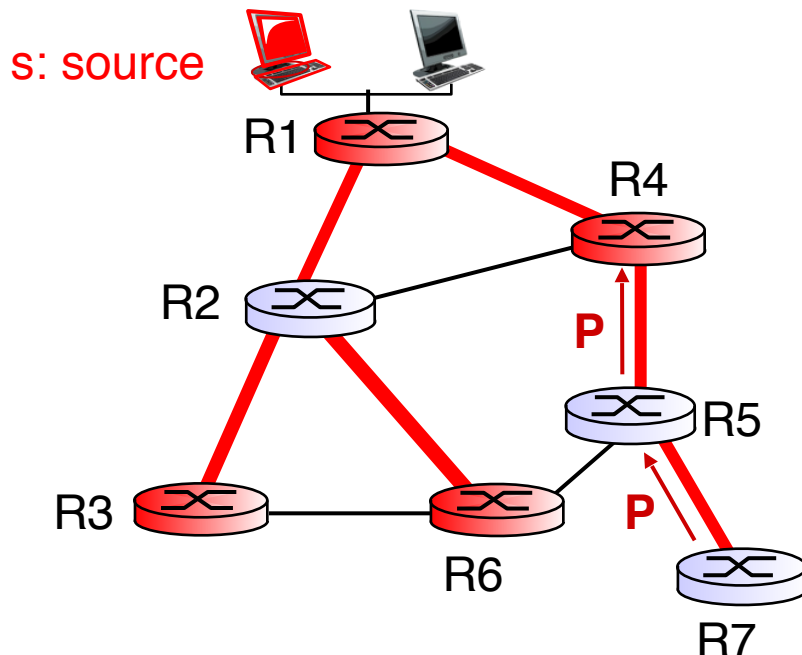  shortest path back to source)
  *then* flood datagram onto all outgoing links
  *else* ignore datagram

# Reverse path forwarding: example



s: source

LEGEND

router with attached group member

router with no attached group member

→ datagram will be forwarded

→ datagram will not be forwarded

❖ result is a source-specific *reverse* SPT

   ▪ may be a bad choice with asymmetric links
   (assumes shortest path R1→R2 is the same as R2→R1)

# Reverse path forwarding: pruning

❖ forwarding tree contains subtrees with no mcast group members
  ■ no need to forward datagrams down subtree
  ■ "prune" msgs sent upstream by router with no downstream group members

s: source

R1

R4

R2

P

R5

R3

R6

P

R7

LEGEND

router with attached group member

router with no attached group member

P → prune message
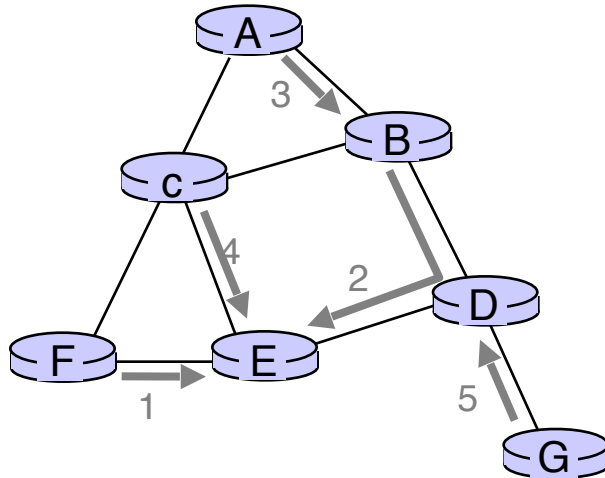
links with multicast forwarding

# Shared-tree: Steiner tree

❖ *Steiner tree:* minimum cost tree connecting all routers with attached group members

❖ problem is NP-complete

❖ excellent heuristics exist

❖ not used in practice:
  ▪ computational complexity
  ▪ information about entire network needed
  ▪ monolithic: rerun whenever a router needs to join/leave
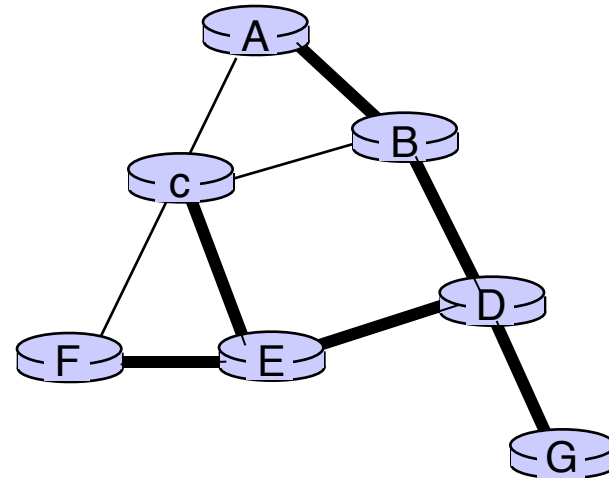
# Center-based tree (Heuristic)

❖ single delivery tree shared by all
❖ one router identified as *"center"* of tree
❖ to join:
  ▪ edge router sends unicast *join-msg* addressed to center router
  ▪ *join-msg* "processed" by intermediate routers and forwarded towards center
  ▪ *join-msg* either hits existing tree branch for this center, or arrives at center
  ▪ path taken by *join-msg* becomes new branch of tree for this router

# Center-based tree: example

❖ Chose a center node

❖ each node sends unicast join message to center node
  - message forwarded until it arrives at a node already belonging to spanning tree



(a) stepwise construction of spanning tree (center: E)

(b) constructed spanning tree

# Internet Multicasting Routing: DVMRP

❖ DVMRP: distance vector multicast routing protocol, RFC1075

❖ *flood and prune:* reverse path forwarding, source-based tree
  - RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers
  - no assumptions about underlying unicast
  - initial datagram to mcast group flooded everywhere via RPF
  - routers not wanting group: send upstream prune msgs

# DVMRP: continued…

❖ *soft state:* DVMRP router periodically (1 min.) "forgets" branches:
  ▪ mcast data again flows down unpruned branch
  ▪ downstream router: reprune or else continue to receive data

❖ routers can quickly regraft to tree
  ▪ following IGMP join at leaf

❖ DVMRP: commonly implemented in commercial router

# PIM: Protocol Independent Multicast

❖ not dependent on any specific underlying unicast routing algorithm (works with all)

❖ two different multicast distribution scenarios :

### *dense:*

❖ group members densely packed, in "close" proximity.

❖ bandwidth more plentiful

### *sparse:*

❖ # networks with group members small wrt # interconnected networks

❖ group members "widely dispersed"

❖ bandwidth not plentiful

# Consequences of sparse-dense dichotomy:

## *dense*

❖ group membership by routers *assumed* until routers explicitly prune

❖ *data-driven* construction on mcast tree

❖ bandwidth and non-group-router processing *wasted*

❖ *Uses flood and prune RPF*

## *sparse*:

❖ no membership until routers explicitly join

❖ *receiver- driven* construction of mcast tree

❖ bandwidth and non-group-router processing *conservative*

❖ Uses center-based tree

# Summary

❖ **Network layer: forwarding and routing**

❖ **Routing: hierarchical**
  ▪ **intra-AS: local, optimal**
  ▪ **and Inter-AS (BGP): global, policy based**

❖ **Routing and protocols for Broadcast and Multicast**