# The Internet Network layer

Host, router network layer functions:

**Network layer**

Transport layer: TCP, UDP

**Routing protocols**
- path selection
- RIP, OSPF, BGP
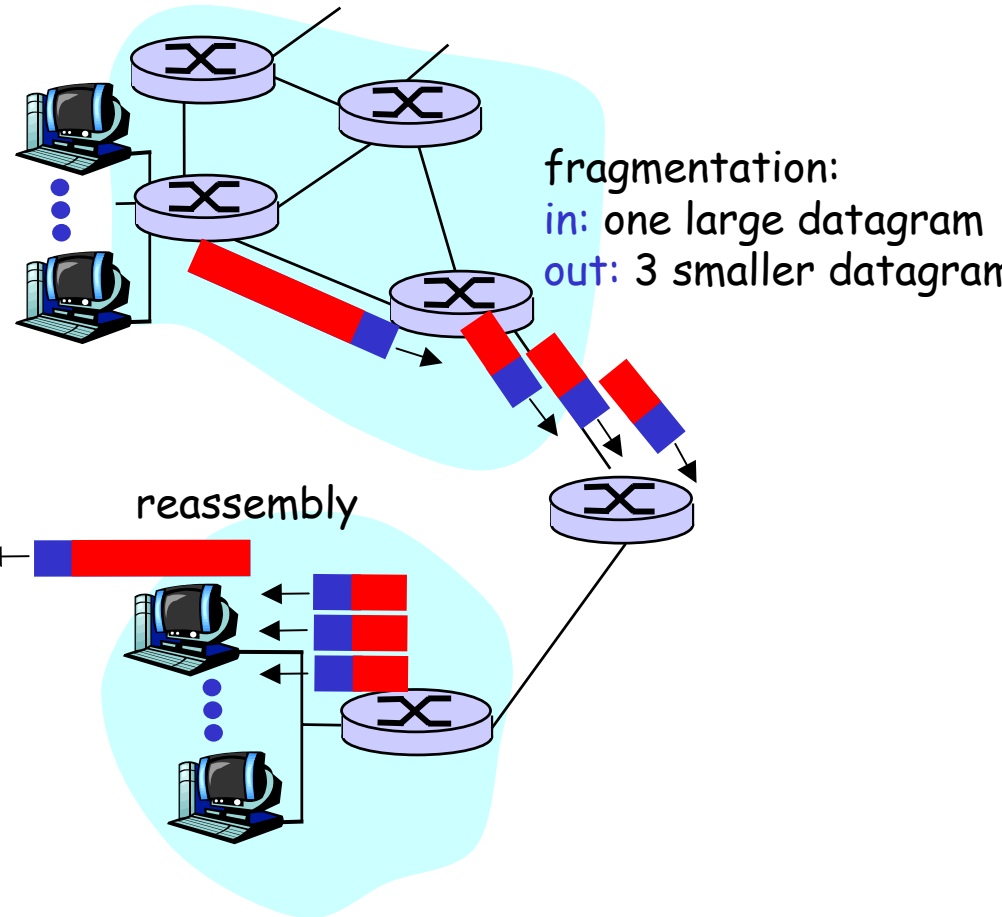
routing table

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

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

Link layer

physical layer

1

# 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

32 bits

total datagram length (bytes)

for fragmentation/ reassembly

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

| ver | head. len | type of service | length | |
| time to live | upper layer | Internet checksum | | |

16-bit identifier | flgs | fragment offset

32 bit source IP address

32 bit destination IP address

Options (if any)

data
(variable length,
typically a TCP
or UDP segment)

# 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 datagram

reassembly

# IP Fragmentation and Reassembly

| | length =4000 | ID =x | fragflag =0 | offset =0 | |
|---|---|---|---|---|---|

One large datagram becomes several smaller datagrams

| | length =1500 | ID =x | fragflag =1 | offset =0 | |
|---|---|---|---|---|---|

| | length =1500 | ID =x | fragflag =1 | offset =1480 | |
|---|---|---|---|---|---|

| | length =1040 | ID =x | fragflag =0 | offset =2960 | |
|---|---|---|---|---|---|

# ICMP: Internet Control Message Protocol

- used by hosts, routers, gateways to communication 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

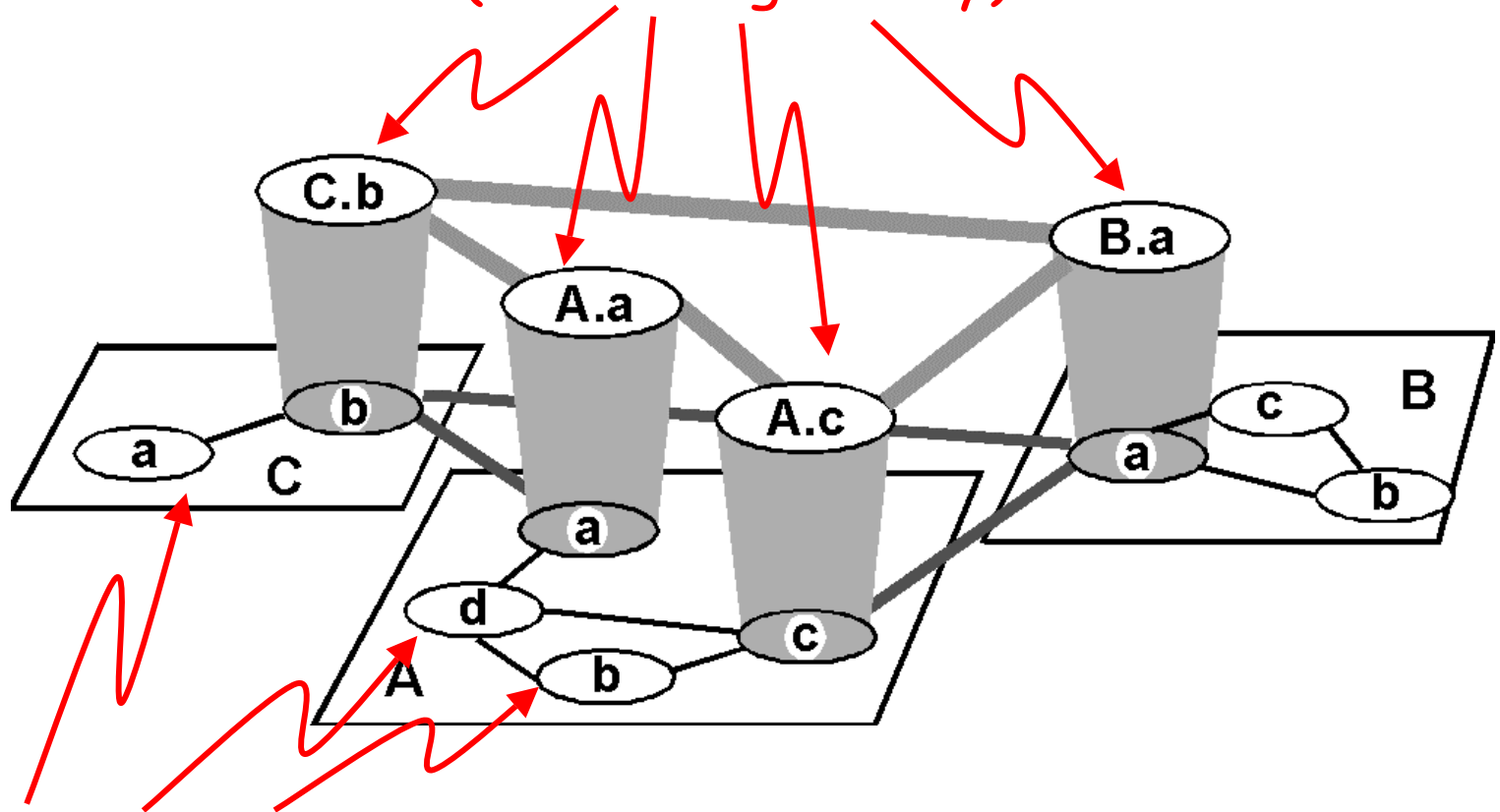| Type | Code | description |
|------|------|-------------|
| 0 | 0 | echo reply (ping) |
| 3 | 0 | dest. network unreachable |
| 3 | 1 | dest host unreachable |
| 3 | 2 | dest protocol unreachable |
| 3 | 3 | dest port unreachable |
| 3 | 6 | dest network unknown |
| 3 | 7 | dest host unknown |
| 4 | 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 |

# Routing in the Internet

- The Global Internet consists of Autonomous Systems (AS) interconnected with each other:
  - **Stub AS**: small corporation
  - **Multihomed AS**: large corporation (no transit)
  - **Transit AS**: provider

- Two-level routing:
  - **Intra-AS:** administrator is responsible for choice
  - **Inter-AS:** unique standard

# Internet AS Hierarchy

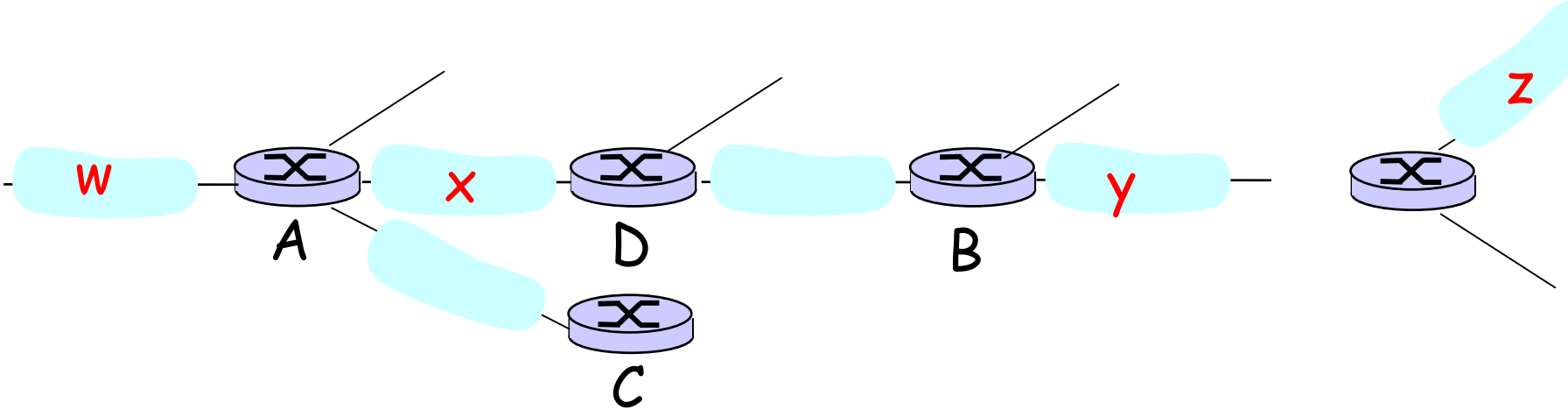Intra-AS border (exterior gateway) routers



Inter-AS interior (gateway) routers

7

# Intra-AS Routing

- Also known as Interior Gateway Protocols (IGP)
- Most common IGPs:

  - RIP: Routing Information Protocol

  - OSPF: Open Shortest Path First

  - IGRP: Interior Gateway Routing Protocol (Cisco propr.)

8

# RIP ( Routing Information Protocol)

- Distance vector algorithm
- Included in BSD-UNIX Distribution in 1982
- Distance metric: # of hops (max = 15 hops)
  - *Can you guess why?*

- Distance vectors: exchanged every 30 sec via Response Message (also called **advertisement**)
- Each advertisement: route to up to 25 destination nets

# RIP (Routing Information Protocol)



| Destination Network | Next Router | Num. of hops to dest. |
|---|---|---|
| w | A | 2 |
| y | B | 2 |
| z | B | 7 |
| x | -- | 1 |
| …. | …. | …. |

Routing table in D

# RIP: Link Failure and Recovery

If no advertisement heard after 180 sec --> neighbor/link declared dead

- routes via neighbor invalidated
- new advertisements sent to neighbors
- neighbors in turn send out new advertisements (if tables changed)
- link failure info quickly propagates to entire net
- poison reverse used to prevent ping-pong loops (infinite distance = 16 hops)

# RIP Table processing

- RIP routing tables managed by a**pplication-level** process called route-d (daemon)
- advertisements sent in UDP packets, periodically repeated

# RIP Table example (continued)

Router: *giroflee.eurocom.fr*

```
Destination          Gateway              Flags  Ref    Use    Interface
-------------------- -------------------- -----  ----- ------  ----------
127.0.0.1            127.0.0.1              UH       0  26492   lo0
192.168.2.           192.168.2.5            U        2     13   fa0
193.55.114.          193.55.114.6           U        3  58503   le0
192.168.3.           192.168.3.5            U        2     25   qaa0
224.0.0.0            193.55.114.6           U        3      0   le0
default              193.55.114.129         UG       0 143454
```

- Three attached class C networks (LANs)
- Router only knows routes to attached LANs
- Default router used to "go up"
- Route multicast address: 224.0.0.0
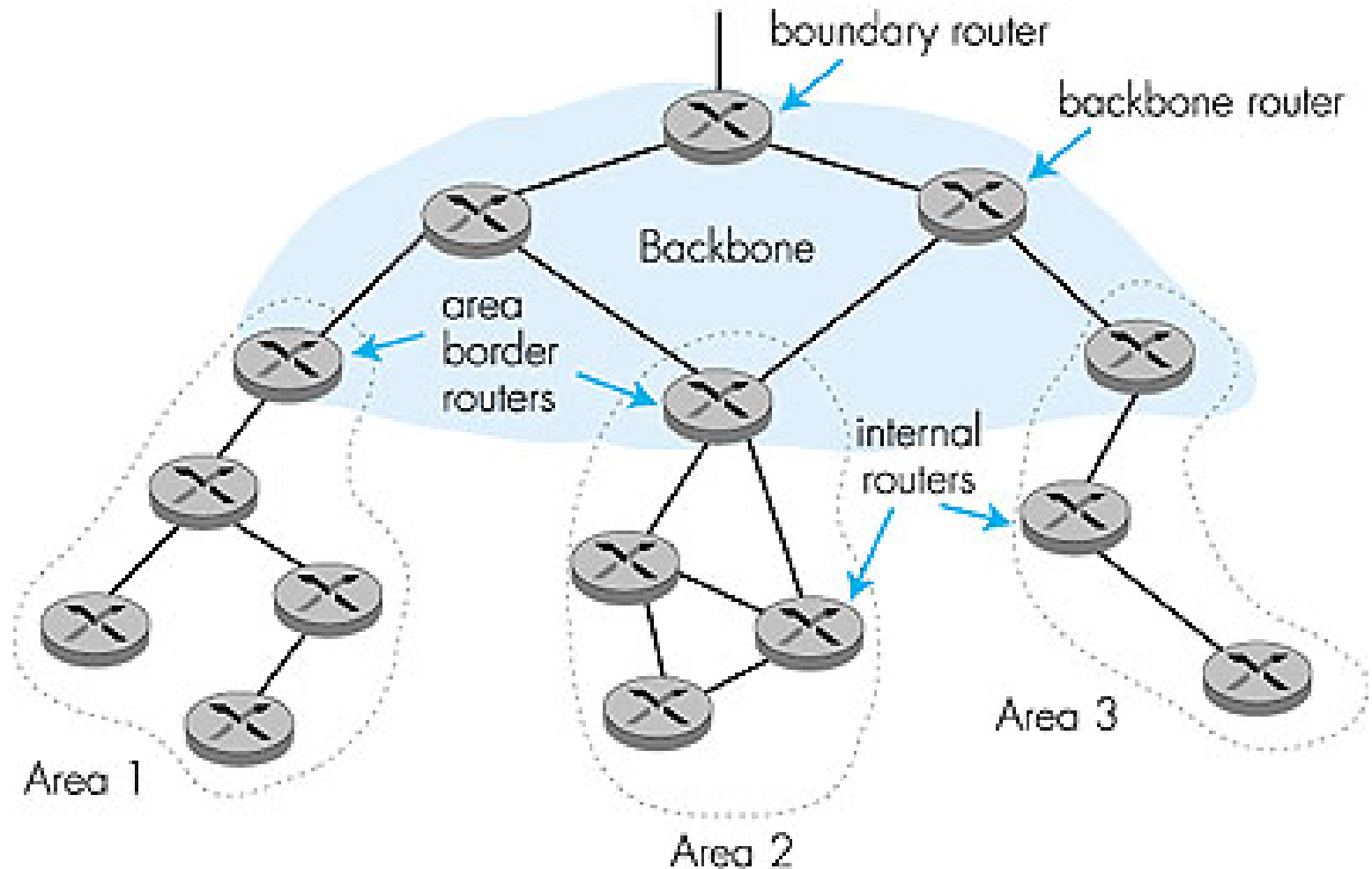- Loopback interface (for debugging)

# OSPF (Open Shortest Path First)

- "open": publicly available
- Uses Link State algorithm
  - LS packet dissemination
  - Topology map at each node
  - Route computation using Dijkstra's algorithm

- OSPF advertisement carries one entry per neighbor router
- Advertisements disseminated to entire AS (via flooding)

# OSPF "advanced" features (not in RIP)

- Security: all OSPF messages authenticated (to prevent malicious intrusion); TCP connections used
- Multiple same-cost paths allowed (only one path in RIP)
- For each link, multiple cost metrics for different TOS (eg, satellite link cost set "low" for best effort; high for real time)
- Integrated uni- and multicast support:
  - Multicast OSPF (MOSPF) uses same topology data base as OSPF
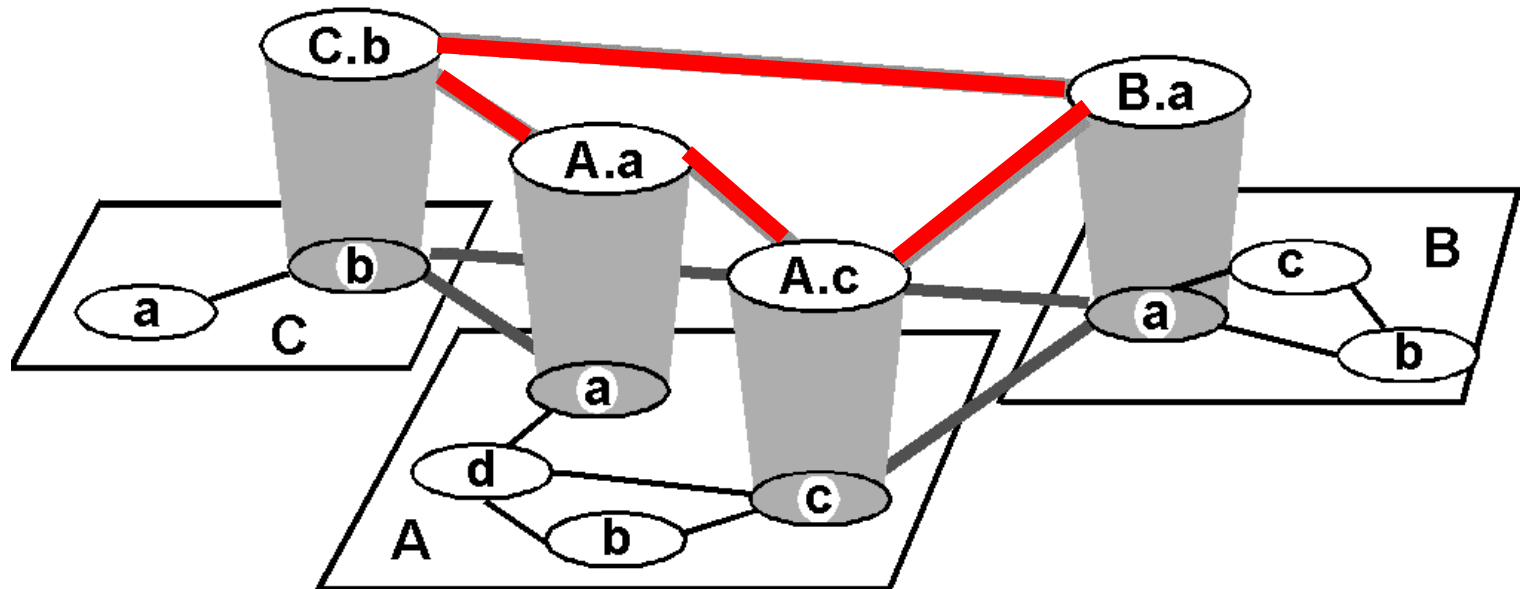- Hierarchical OSPF in large domains.

# Hierarchical OSPF

# Hierarchical OSPF

- Two-level hierarchy: local area, backbone.
    - Link-state advertisements only in area
    - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
- **Area border routers:** "summarize" distances to nets in own area, advertise to other Area Border routers.
- **Backbone routers:** run OSPF routing limited to backbone.
- **Boundary routers:** connect to other ASs.

# IGRP (Interior Gateway Routing Protocol)

- CISCO proprietary; successor of RIP (mid 80s)
- Distance Vector, like RIP
- several cost metrics (delay, bandwidth, reliability, load etc)
- uses TCP to exchange routing updates
- Loop-free routing via Distributed Updating Alg. (DUAL) based on *diffused computation*

# Inter-AS routing

# Internet inter-AS routing: BGP

- BGP (Border Gateway Protocol): *the* de facto standard

- **Path Vector** protocol:
  - similar to Distance Vector protocol
  - each Border Gateway broadcast to neighbors (peers) *entire path* (I.e, sequence of ASs) to destination
  - E.g., Gateway X may send its path to dest. Z:

  Path (X,Z) = X,Y1,Y2,Y3,…,Z

# Internet inter-AS routing: BGP

*Suppose:* gateway X send its path to peer gateway W

- W may or may not select path offered by X
    - cost, policy (don't route via competitors AS), loop prevention reasons.
- If W selects path advertised by X, then:

$$\text{Path (W,Z) = w, Path (X,Z)}$$

- Note: X can control incoming traffic by controling it route advertisements to peers:
    - e.g., don't want to route traffic to Z -> don't advertise any routes to Z

# Internet inter-AS routing: BGP

- BGP messages exchanged using TCP.
- BGP messages:
  - **OPEN:** opens TCP connection to peer and authenticates sender
  - **UPDATE:** advertises new path (or withdraws old)
  - **KEEPALIVE** keeps connection alive in absence of UPDATES; also ACKs OPEN request
  - **NOTIFICATION:** reports errors in previous msg; also used to close connection

# Why different Intra- and Inter-AS routing ?

## Policy:

- Inter-AS: admin wants control over how its traffic routed, who routes through its net.
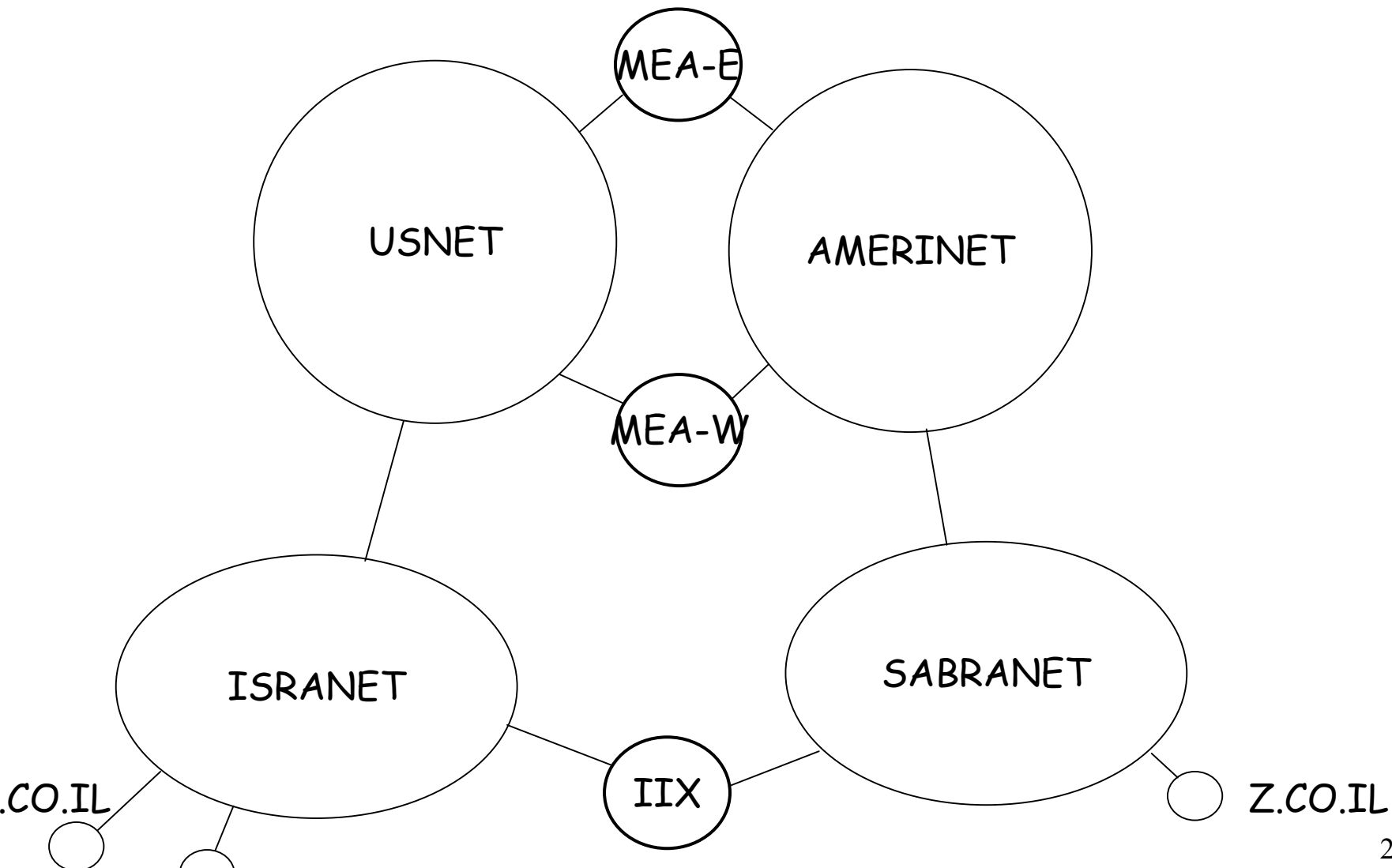- Intra-AS: single admin, so no policy decisions needed

## Scale:

- hierarchical routing saves table size, reduced update traffic

## Performance:

- Intra-AS: can focus on performance
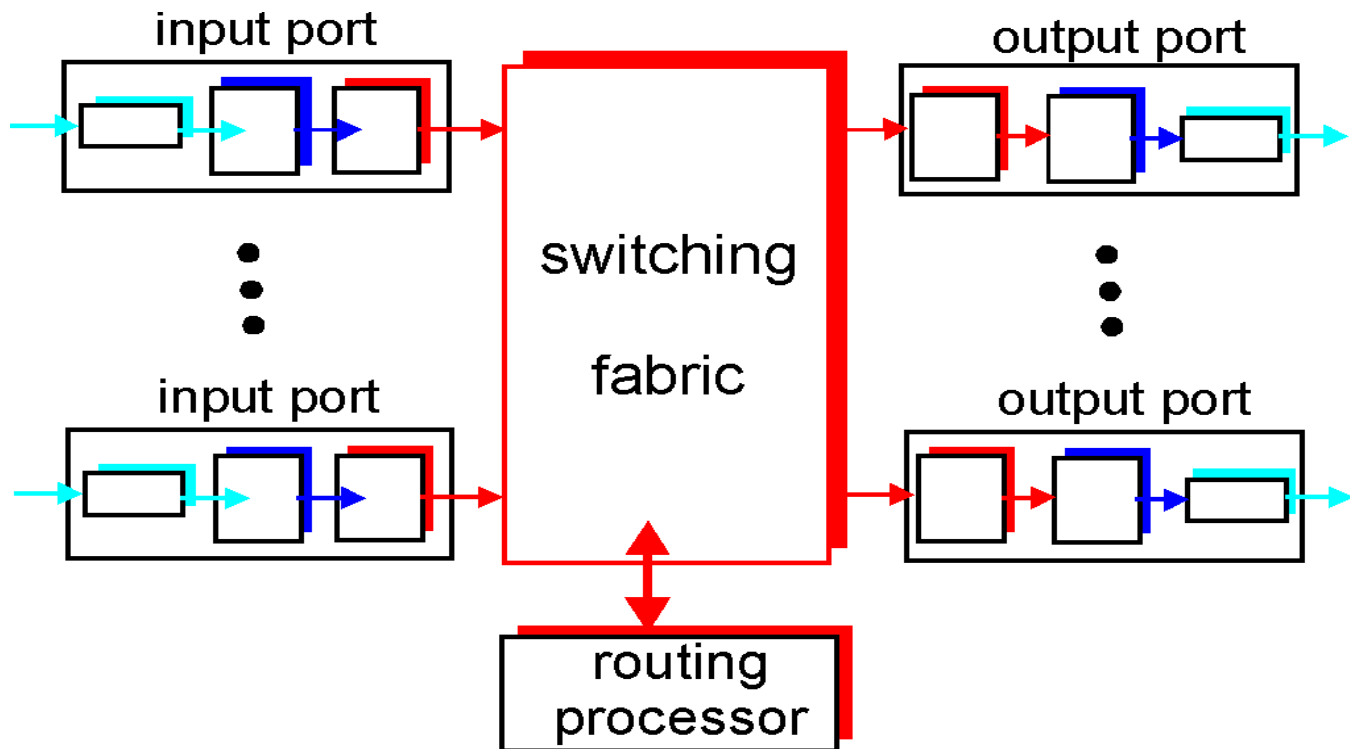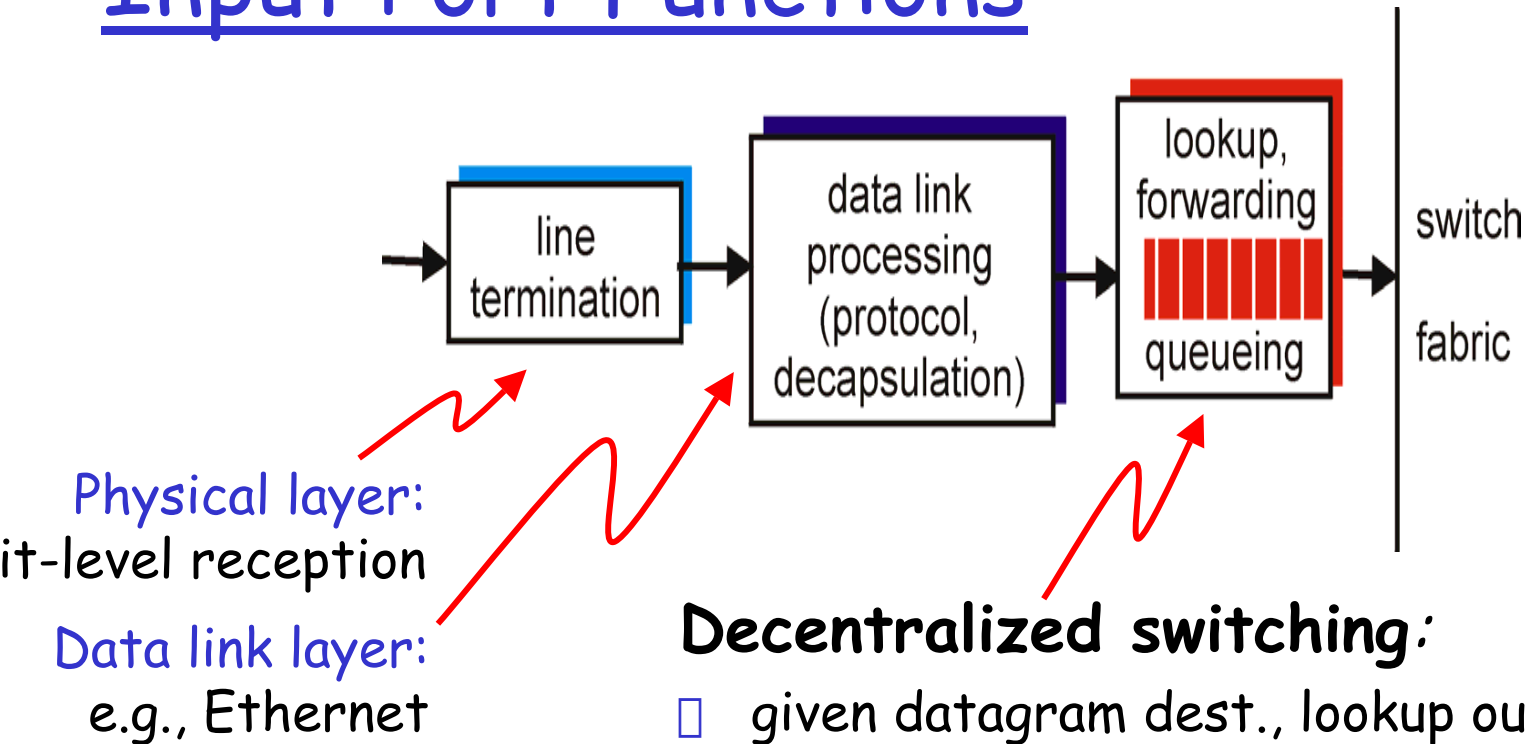- Inter-AS: policy may dominate over performance

# BGP Policy Example



Diagram showing interconnected networks: MEA-E connects USNET and AMERINET at top. USNET and AMERINET are also connected via MEA-W. USNET connects to ISRANET, and AMERINET connects to SABRANET. ISRANET connects to IIX which connects to SABRANET. .CO.IL nodes connect to ISRANET, and Z.CO.IL connects to SABRANET.

# Router Architecture Overview

Two key router functions:

- run routing algorithms/protocol (RIP, OSPF, BGP)
- *switching* datagrams from incoming to outgoing link

# Input Port Functions



**Physical layer:**
bit-level reception
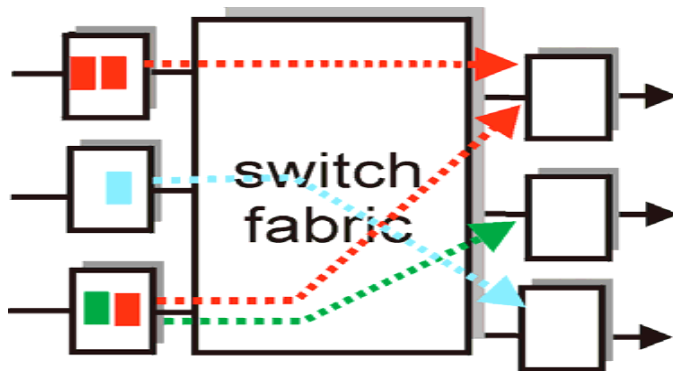
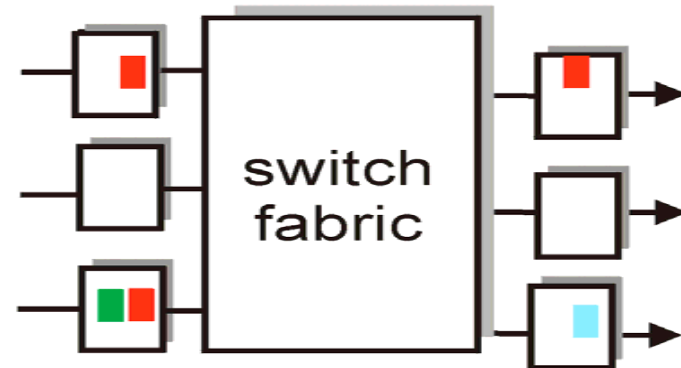**Data link layer:**
e.g., Ethernet

**Decentralized switching:**

- given datagram dest., lookup output port using routing table in input port memory

- goal: complete input port processing at 'line speed'

- queuing: if datagrams arrive faster than forwarding rate into switch fabric

# 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

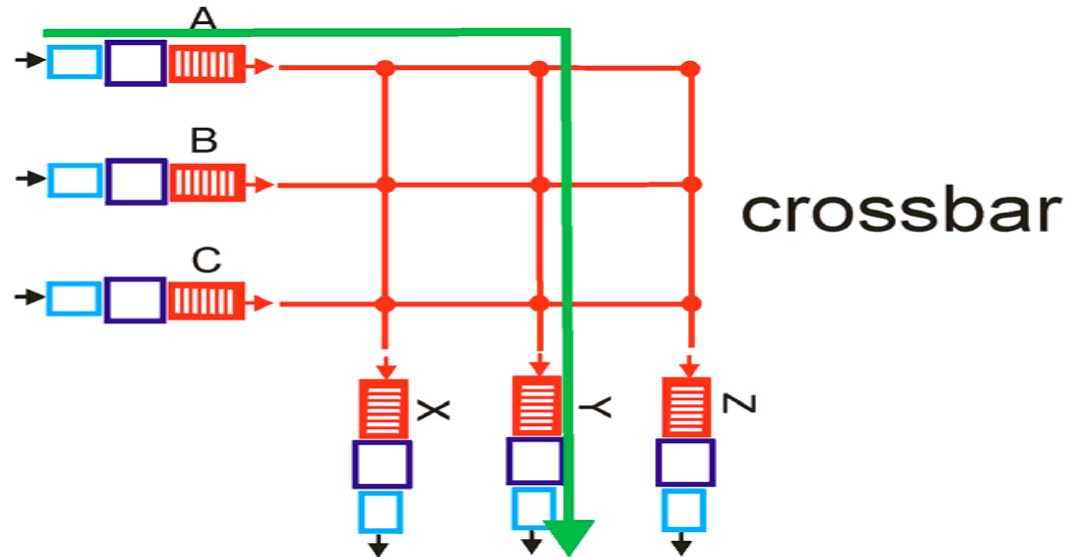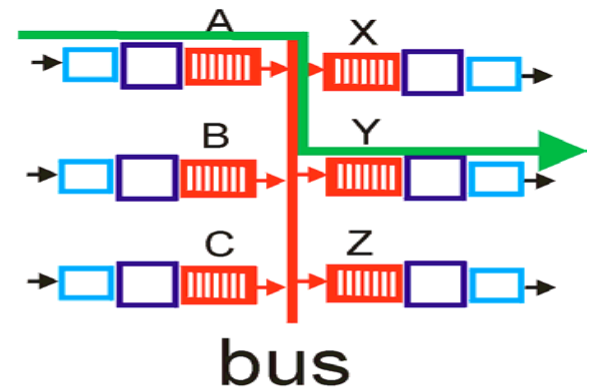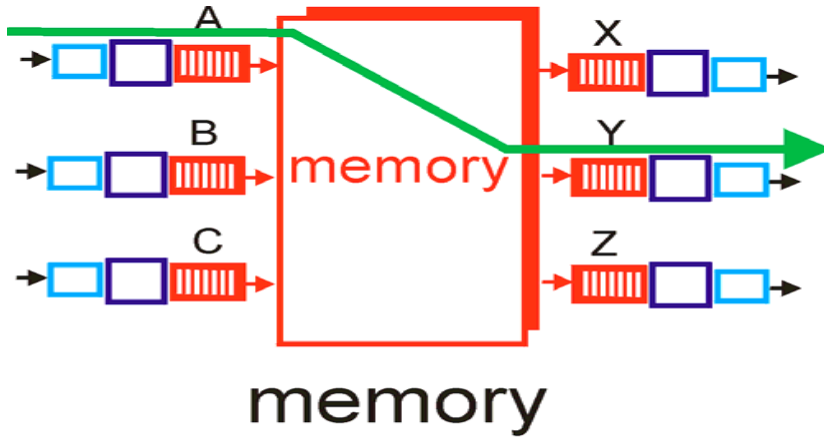- *queueing delay and loss due to input buffer overflow!*



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

green packet
experiences HOL blocking

# Three types of switching fabrics



memory

bus

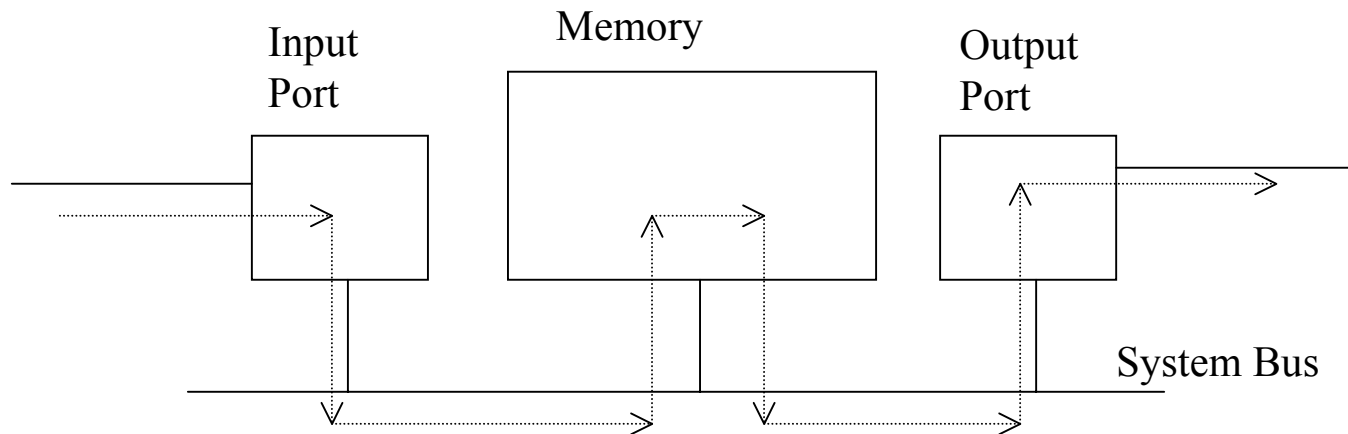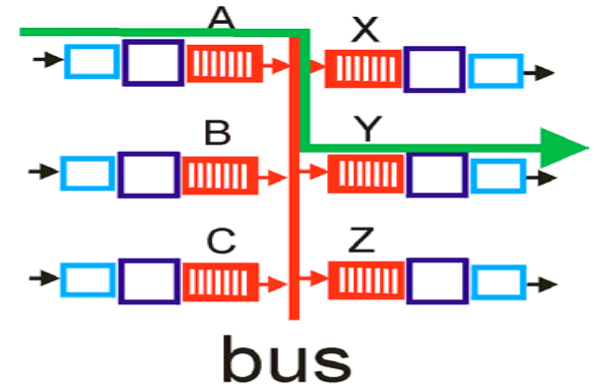crossbar

# Switching Via Memory

First generation routers:

- packet copied by system's (single) CPU
- speed limited by memory bandwidth (2 bus crossings per datagram)

Input Port     Memory     Output Port

System Bus

Modern routers:

- input port processor performs lookup, copy into memory
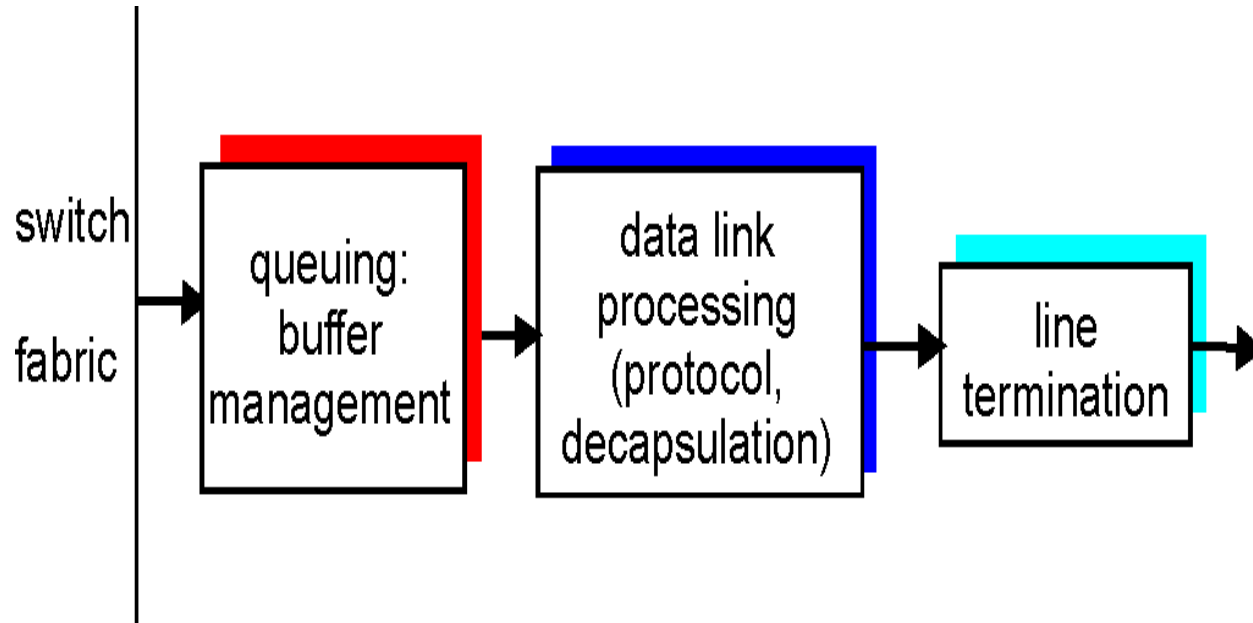- Cisco Catalyst 8500

# Switching Via Bus



bus

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

- bus contention:  switching speed limited by bus bandwidth

- 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)

# 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 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 exceeeds ouput line speed
- *queueing (delay) and loss due to output port buffer overflow!*

# IPv6

- Initial motivation: 32-bit address space completely allocated by 2008.
- Additional motivation:
  - header format helps speed processing/forwarding
  - header changes to facilitate QoS
  - new "anycast" address: route to "best" of several replicated servers
- IPv6 datagram format:
  - fixed-length 40 byte header
  - no fragmentation allowed

# IPv6 Header (Cont)

*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 simultaneous
  - no "flag days"
  - How will the network operate with mixed IPv4 and IPv6 routers?
- Two proposed approaches:
  - *Dual Stack*: some routers with dual stack (v6, v4) can "translate" between formats
  - *Tunneling:* IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Dual Stack Approach

# Tunneling

Logical view

IPv6  IPv6        tunnel           IPv6  IPv6

A    B ═══════════════════════ E    F

Physical view

IPv6 inside IPv4 where needed

IPv6   IPv6   IPv4   IPv4   IPv6   IPv6

A    B    C    D    E    F

```
Flow: X          Source: B          Source: B          Flow: X
Source: A        Dest: E            Dest: E            Source: A
Dest: F                                                Dest: F
                   Flow: X            Flow: X
data               Source: A          Source: A        data
                   Dest: F            Dest: F
```

A to B:          data               data               E to F:
IPv6                                                    IPv6

B to C:                             B to C:
IPv4                                IPv4
(encapsulating                      (encapsulating
IPv6)                               IPv6)