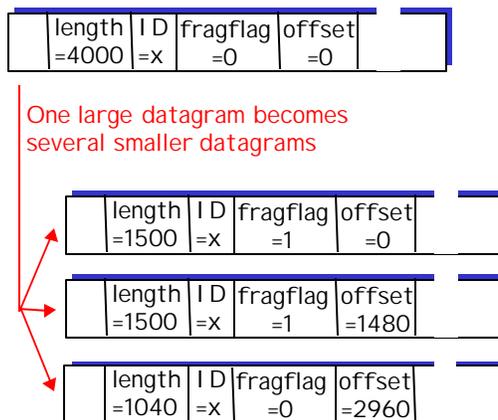


## IP Fragmentation and Reassembly



4: Network Layer 4b-3

## ICMP: Internet Control Message Protocol

<ul style="list-style-type: none"> <li>r used by hosts, routers, gateways to communication network-level information</li> <li>m error reporting: unreachable host, network, port, protocol</li> <li>m echo request/reply (used by ping)</li> <li>r network-layer "above" IP:                             <ul style="list-style-type: none"> <li>m ICMP msgs carried in IP datagrams</li> </ul> </li> <li>r <b>ICMP message:</b> type, code plus first 8 bytes of IP datagram causing error</li> </ul>	<table border="1"> <thead> <tr> <th>Type</th> <th>Code</th> <th>description</th> </tr> </thead> <tbody> <tr><td>0</td><td>0</td><td>echo reply (ping)</td></tr> <tr><td>3</td><td>0</td><td>dest. network unreachable</td></tr> <tr><td>3</td><td>1</td><td>dest host unreachable</td></tr> <tr><td>3</td><td>2</td><td>dest protocol unreachable</td></tr> <tr><td>3</td><td>3</td><td>dest port unreachable</td></tr> <tr><td>3</td><td>6</td><td>dest network unknown</td></tr> <tr><td>3</td><td>7</td><td>dest host unknown</td></tr> <tr><td>4</td><td>0</td><td>source quench (congestion control - not used)</td></tr> <tr><td>8</td><td>0</td><td>echo request (ping)</td></tr> <tr><td>9</td><td>0</td><td>route advertisement</td></tr> <tr><td>10</td><td>0</td><td>router discovery</td></tr> <tr><td>11</td><td>0</td><td>TTL expired</td></tr> <tr><td>12</td><td>0</td><td>bad IP header</td></tr> </tbody> </table>	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
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																																									

4: Network Layer 4b-4



## Intra-AS Routing

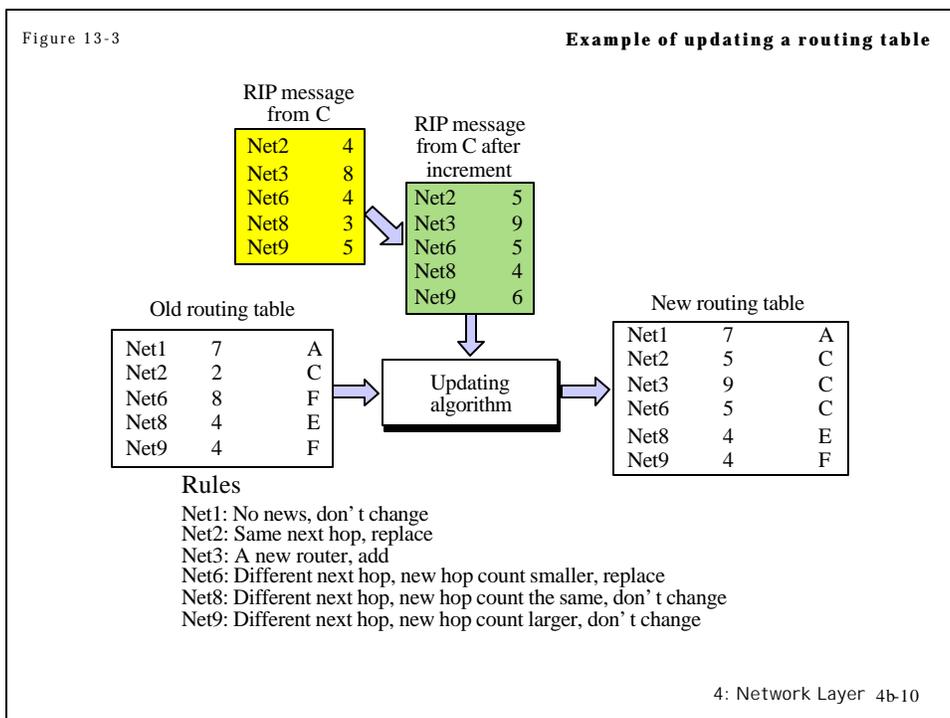
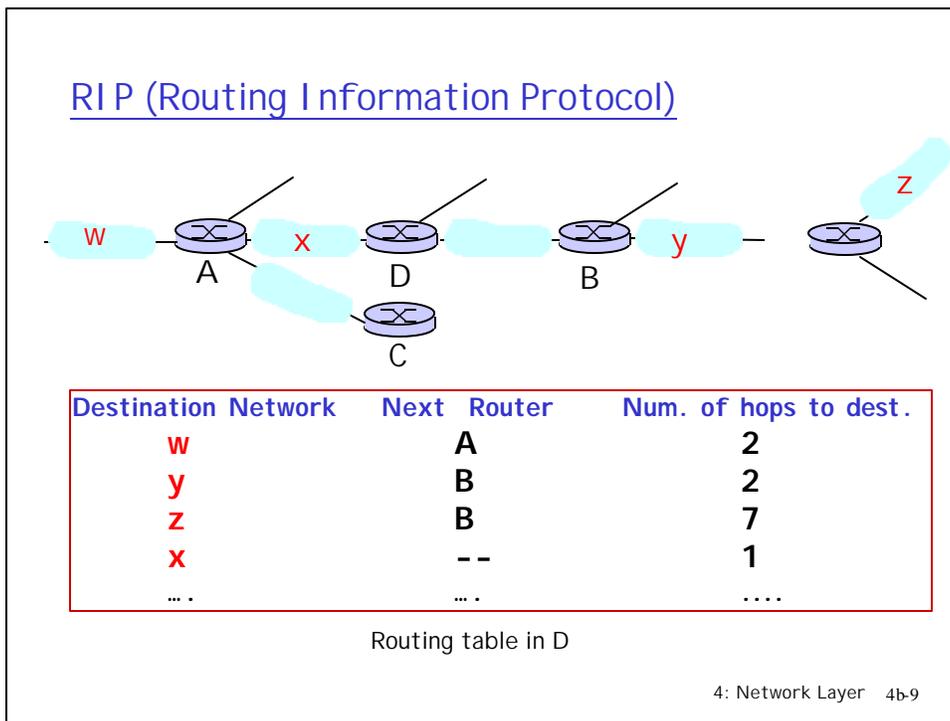
- r Also known as **Interior Gateway Protocols (IGP)**
- r Most common IGPs:
  - m RIP: Routing Information Protocol
  - m OSPF: Open Shortest Path First
  - m IGRP: Interior Gateway Routing Protocol (Cisco propr.)

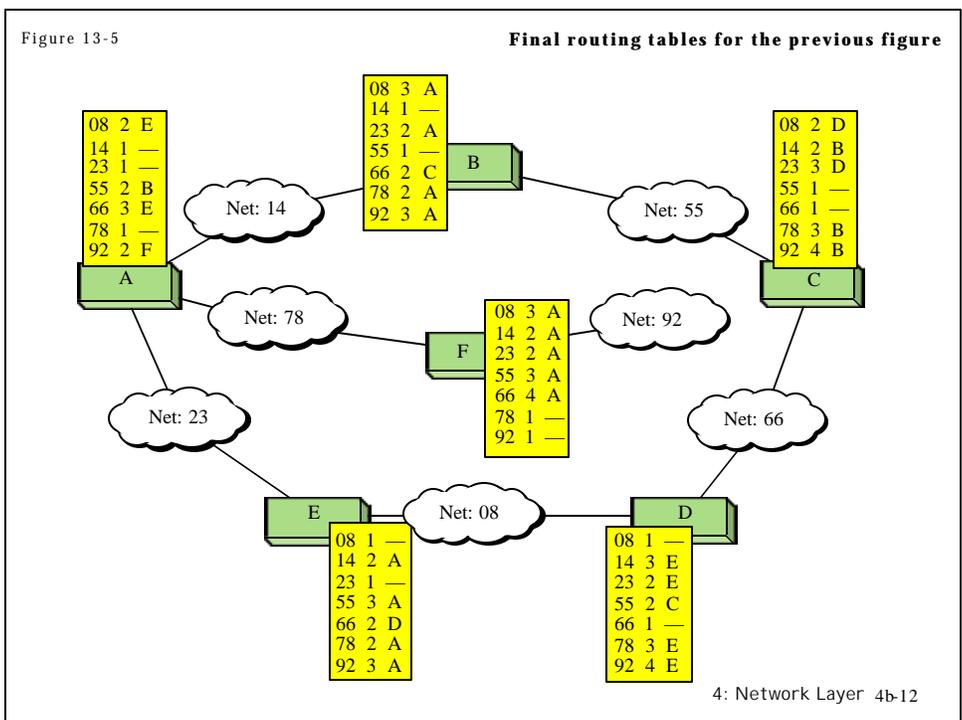
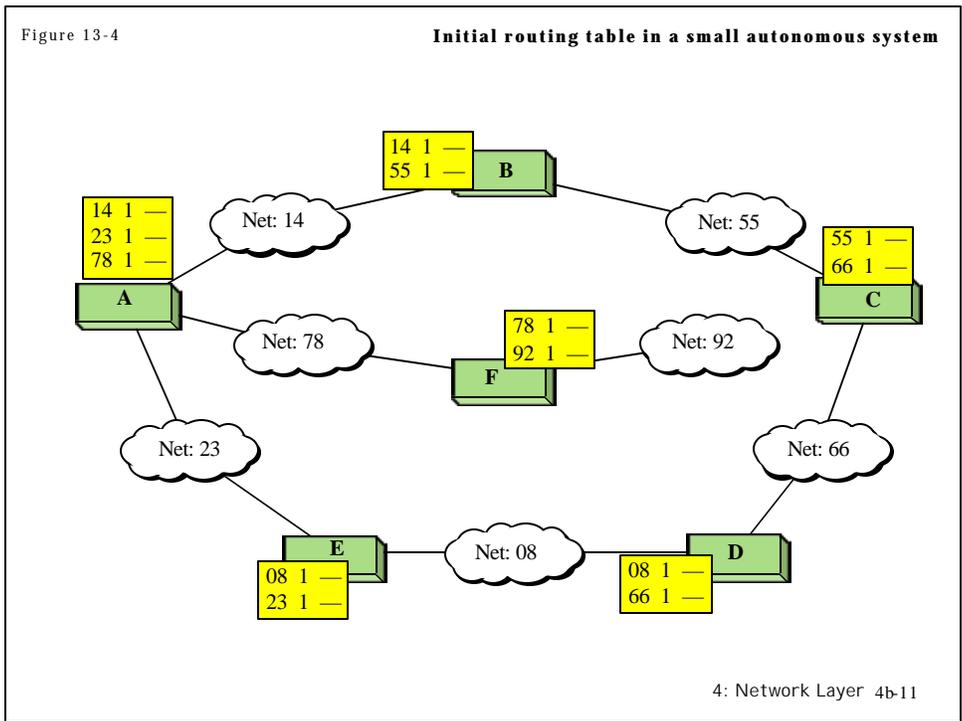
4: Network Layer 4b-7

## RIP ( Routing Information Protocol)

- r Distance vector algorithm
- r Included in BSD-UNIX Distribution in 1982
- r Distance metric: # of hops (max = 15 hops)
  - m *Can you guess why?*
- r Distance vectors: exchanged every 30 sec via Response Message (also called **advertisement**)
- r Each advertisement: route to up to 25 destination nets

4: Network Layer 4b-8





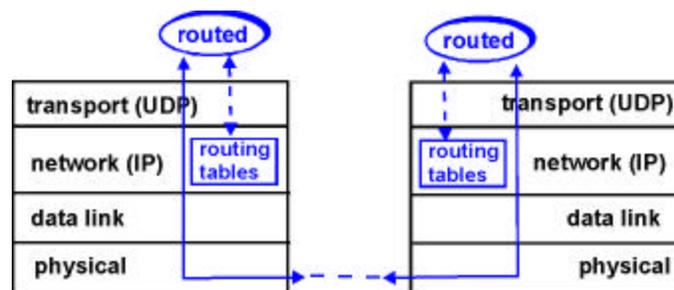
### RIP: Link Failure and Recovery

- If no advertisement heard after 180 sec --> neighbor/link declared dead
  - m routes via neighbor invalidated
  - m new advertisements sent to neighbors
  - m neighbors in turn send out new advertisements (if tables changed)
  - m link failure info quickly propagates to entire net
  - m poison reverse used to prevent ping-pong loops (infinite distance = 16 hops)

4: Network Layer 4b-13

### RIP Table processing

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



4: Network Layer 4b-14

RIP Table example (continued)

Router: *giroflée.eurocom.fr*

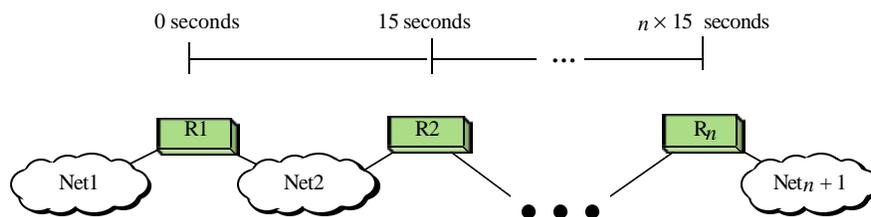
Destination	Gateway	Flags	Ref	Use	Interface
127.0.0.1	127.0.0.1	UH	0	26492	lo0
192.168.2.	192.168.2.5	U	2	13	fa0
193.55.114.	193.55.114.6	U	3	58503	le0
192.168.3.	192.168.3.5	U	2	25	qaa0
224.0.0.0	193.55.114.6	U	3	0	le0
default	193.55.114.129	UG	0	143454	

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

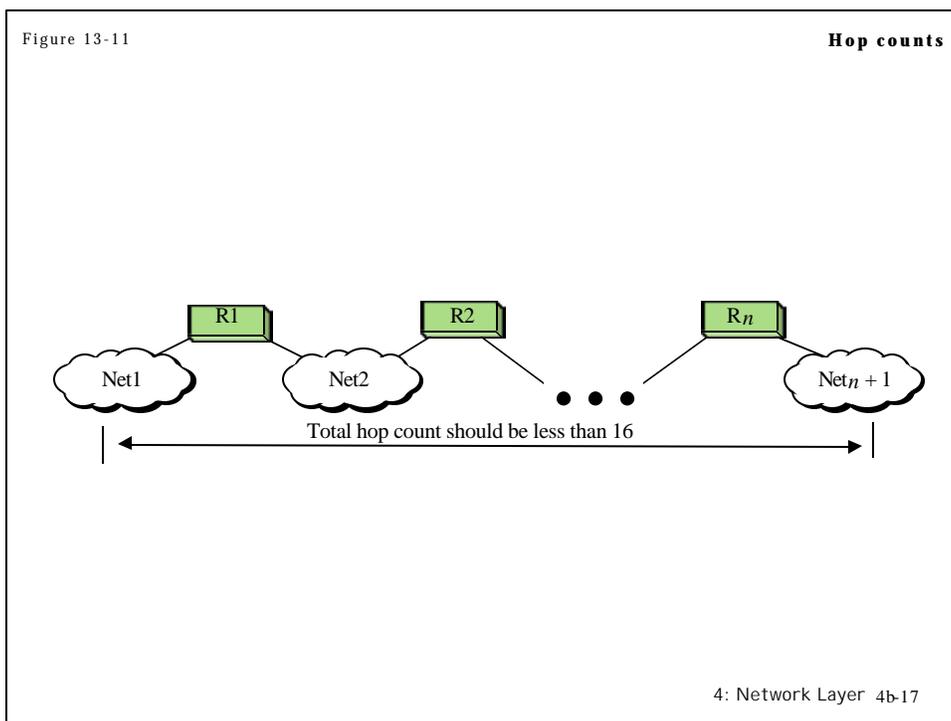
4: Network Layer 4b-15

Figure 13-10

**Slow convergence**



4: Network Layer 4b-16



## OSPF (Open Shortest Path First)

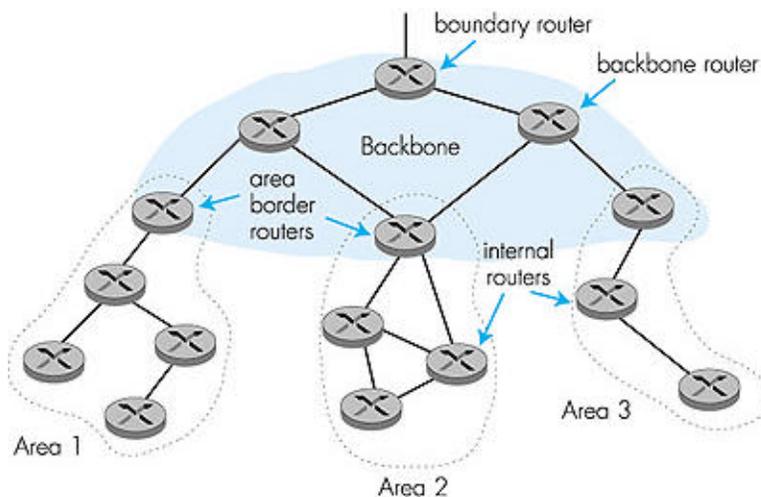
- r "open": publicly available
- r Uses Link State algorithm
  - m LS packet dissemination
  - m Topology map at each node
  - m Route computation using Dijkstra's algorithm
- r OSPF advertisement carries one entry per neighbor router
- r Advertisements disseminated to **entire** AS (via flooding)

## OSPF "advanced" features (not in RIP)

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

4: Network Layer 4b-19

## Hierarchical OSPF



4: Network Layer 4b-20

## Hierarchical OSPF

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

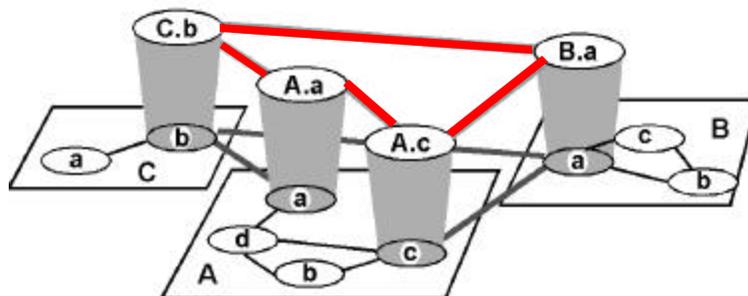
4: Network Layer 4b-21

## IGRP (Interior Gateway Routing Protocol)

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

4: Network Layer 4b-22

## Inter-AS routing



4: Network Layer 4b-23

## Internet inter-AS routing: BGP

- r **BGP (Border Gateway Protocol)**: *the de facto standard*
- r **Path Vector** protocol:
  - m similar to Distance Vector protocol
  - m each Border Gateway broadcast to neighbors (peers) *entire path* (i.e., sequence of ASs) to destination
  - m E.g., Gateway X may send its path to dest. Z:

$$\text{Path}(X,Z) = X, Y_1, Y_2, Y_3, \dots, Z$$

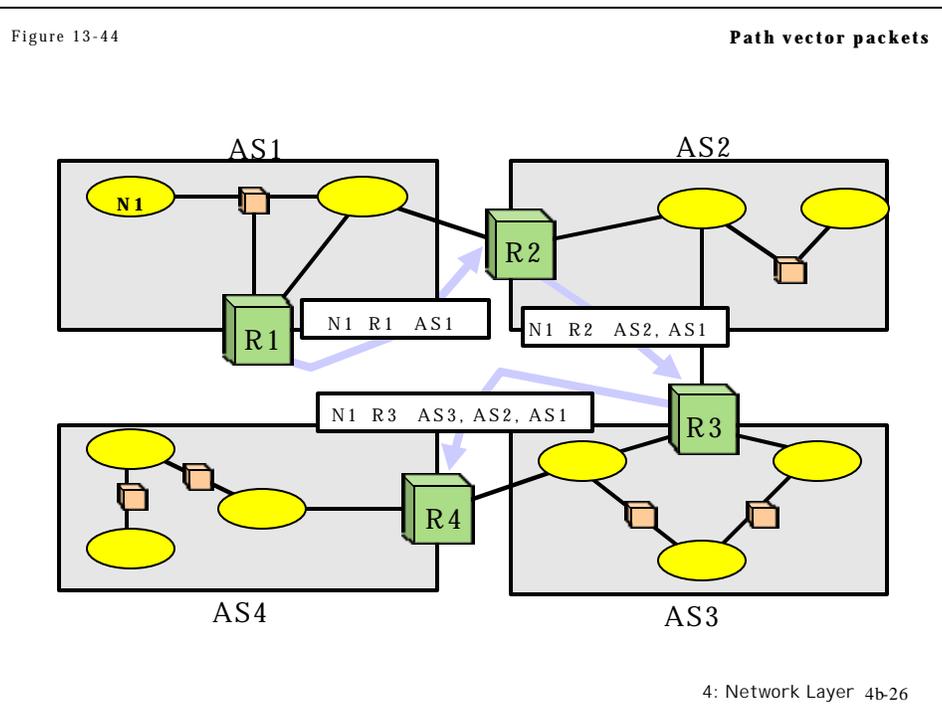
4: Network Layer 4b-24

## Internet inter-AS routing: BGP

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

- r W may or may not select path offered by X
  - m cost, policy (don't route via competitors AS), loop prevention reasons.
- r If W selects path advertised by X, then:  
     Path (W,Z) = w, Path (X,Z)
- r Note: X can control incoming traffic by controlling its route advertisements to peers:
  - m e.g., don't want to route traffic to Z -> don't advertise any routes to Z

4: Network Layer 4b-25



## Internet inter-AS routing: BGP

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

4: Network Layer 4b-27

## Why different Intra- and Inter-AS routing ?

### **Policy:**

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

### **Scale:**

- r hierarchical routing saves table size, reduced update traffic

### **Performance:**

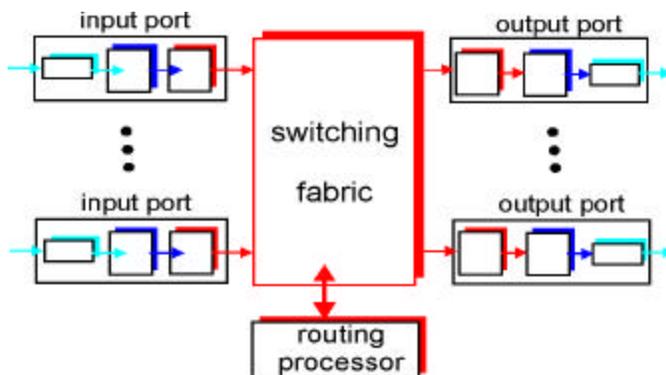
- r Intra-AS: can focus on performance
- r Inter-AS: policy may dominate over performance

4: Network Layer 4b-28

## Router Architecture Overview

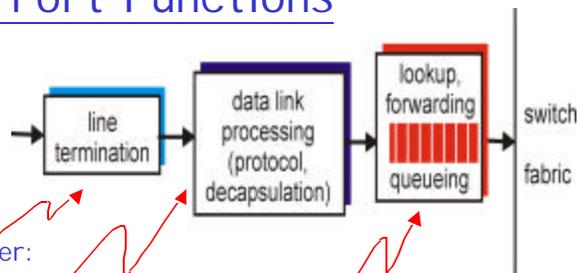
Two key router functions:

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



4: Network Layer 4b-29

## Input Port Functions



Physical layer:  
bit-level reception

Data link layer:  
e.g., Ethernet  
see chapter 5

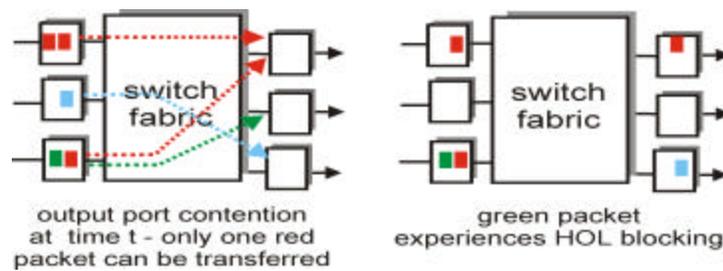
### Decentralized switching:

- r given datagram dest., lookup output port using routing table in input port memory
- r goal: complete input port processing at 'line speed'
- r queuing: if datagrams arrive faster than forwarding rate into switch fabric

4: Network Layer 4b-30

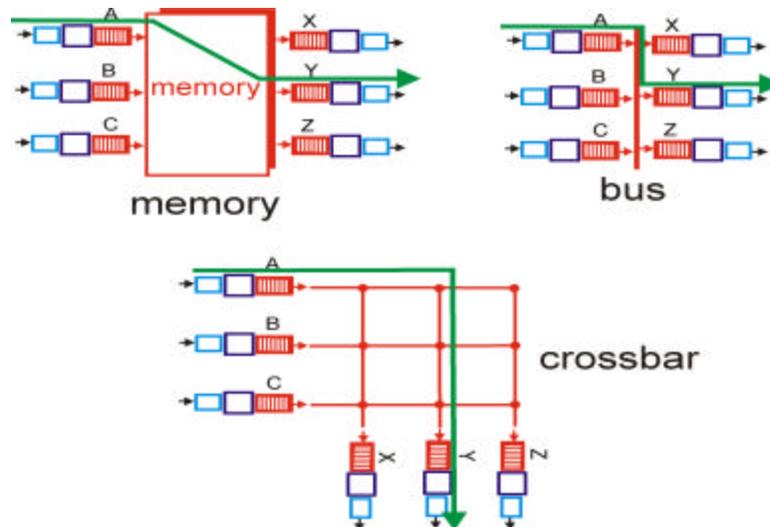
## Input Port Queuing

- r Fabric slower than input ports combined -> queuing may occur at input queues
- r **Head-of-the-Line (HOL) blocking**: queued datagram at front of queue prevents others in queue from moving forward
- r *queuing delay and loss due to input buffer overflow!*



4: Network Layer 4b-31

## Three types of switching fabrics

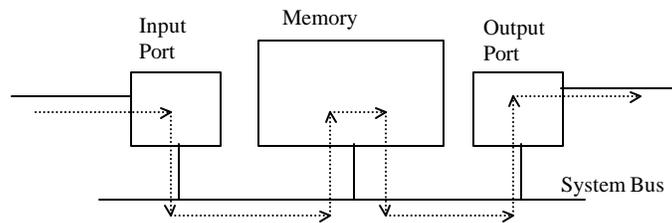


4: Network Layer 4b-32

## Switching Via Memory

### First generation routers:

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

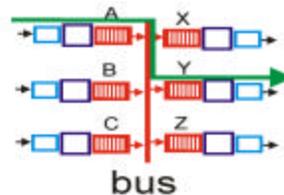


### Modern routers:

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

4: Network Layer 4b-33

## Switching Via Bus



- r datagram from input port memory to output port memory via a shared bus
- r **bus contention:** switching speed limited by bus bandwidth
- r 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)

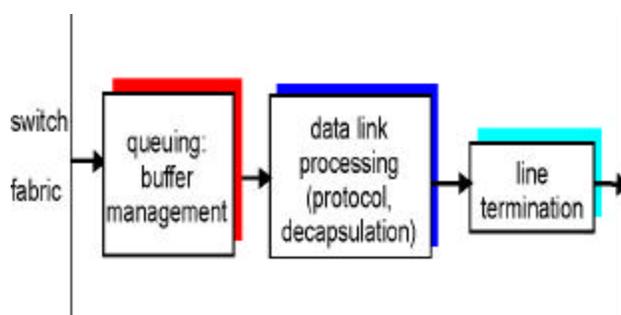
4: Network Layer 4b-34

## Switching Via An Interconnection Network

- r overcome bus bandwidth limitations
- r Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor
- r Advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- r Cisco 12000: switches Gbps through the interconnection network

4: Network Layer 4b-35

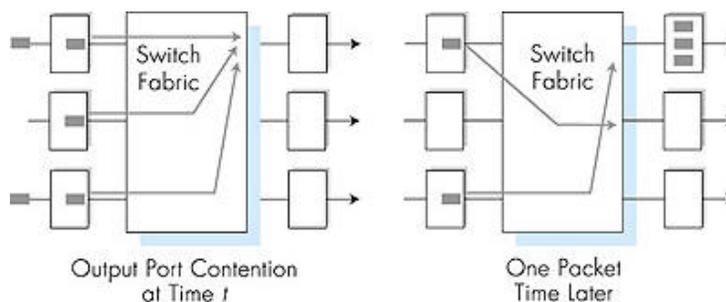
## Output Ports



- r *Buffering* required when datagrams arrive from fabric faster than the transmission rate
- r *Scheduling discipline* chooses among queued datagrams for transmission

4: Network Layer 4b-36

## Output port queueing



- r buffering when arrival rate via switch exceeds output line speed
- r *queueing (delay) and loss due to output port buffer overflow!*

4: Network Layer 4b-37

## IPv6

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

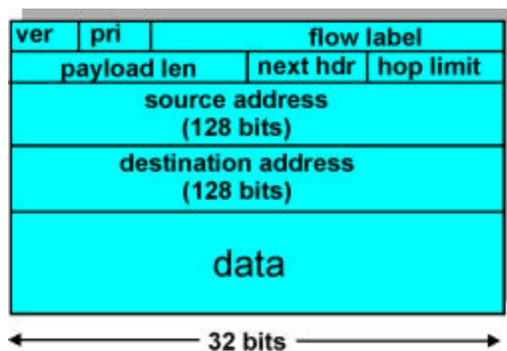
4: Network Layer 4b-38

## 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



4: Network Layer 4b-39

## Other Changes from IPv4

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

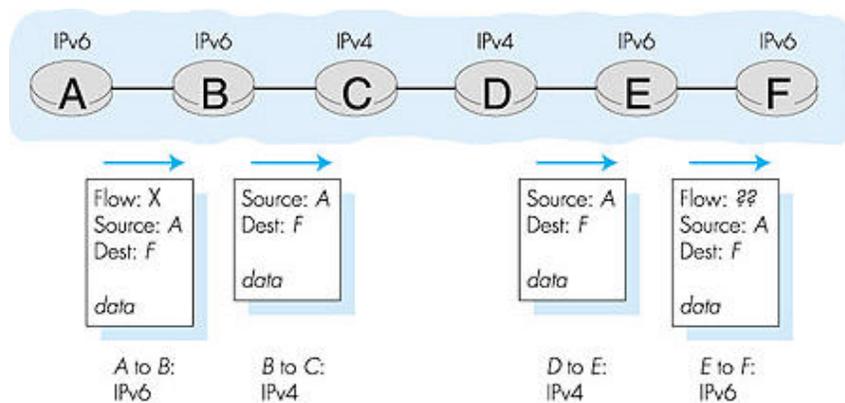
4: Network Layer 4b-40

## Transition From IPv4 To IPv6

- r Not all routers can be upgraded simultaneously
  - m no "flag days"
  - m How will the network operate with mixed IPv4 and IPv6 routers?
- r Two proposed approaches:
  - m *Dual Stack*: some routers with dual stack (v6, v4) can "translate" between formats
  - m *Tunneling*: IPv6 carried as payload in IPv4 datagram among IPv4 routers

4: Network Layer 4b-41

## Dual Stack Approach



4: Network Layer 4b-42

