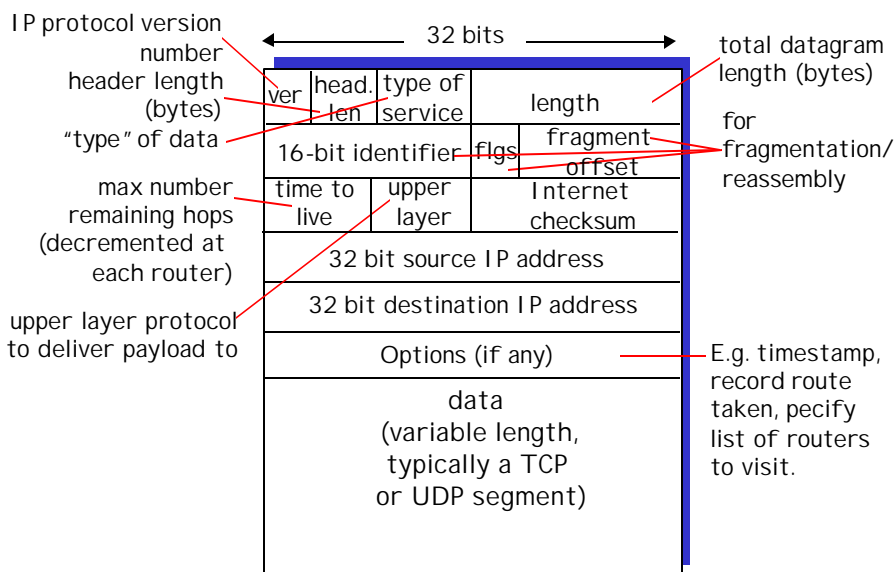## Routing in the Internet

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

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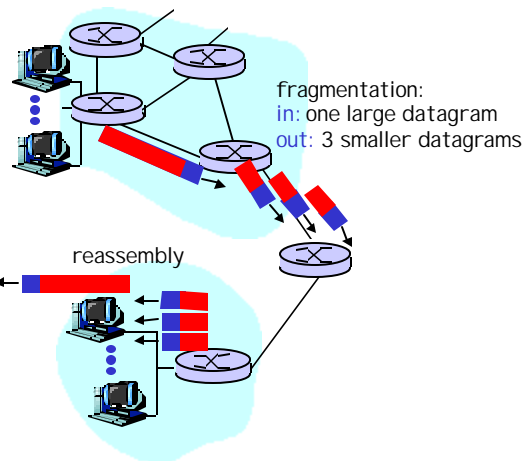  m **Inter-AS:** unique standard

4: Network Layer    4a-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

| ver | head. len | type of service | length |
| 16-bit identifier | flgs | fragment offset |
| time to live | upper layer | Internet 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, pecify list of routers to visit.

4: Network Layer    4a-2

# IP Fragmentation & Reassembly

r  network links have MTU
   (max.transfer size) - largest
   possible link-level frame.

   m  different link types,
      different MTUs

r  large IP datagram divided
   ("fragmented") within net

   m  one datagram becomes
      several datagrams

   m  "reassembled" only at final
      destination

   m  IP header bits used to
      identify, order related
      fragments

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

reassembly

4: Network Layer     4a-3

# 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 | |
|---|---|---|---|---|

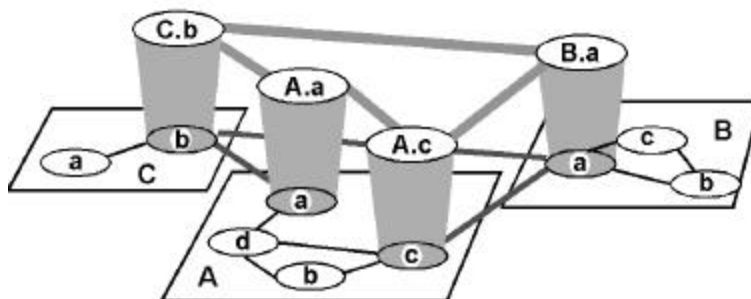| length =1040 | ID =x | fragflag =0 | offset =2960 | |
|---|---|---|---|---|

4: Network Layer     4a-4

# ICMP: Internet Control Message Protocol

r   used by hosts, routers, gateways to communication network-level information

  m   error reporting: unreachable host, network, port, protocol

  m   echo request/reply (used by ping)

r   network-layer "above" IP:

  m   ICMP msgs carried in IP datagrams

r   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 |

4: Network Layer    4a-5

# Internet AS Hierarchy



4: Network Layer    4a-6

## 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.)
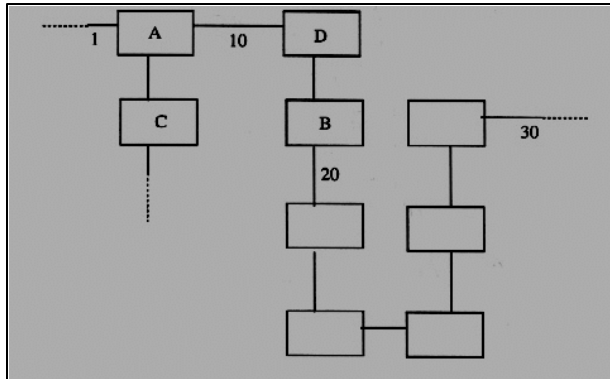
4: Network Layer   4a-7

## RIP ( Routing Information Protocol)

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

r  Distance vector: exchanged every 30 sec via a Response Message (also called **Advertisement**)
r  Each Advertisement contains up to 25 destination nets

4: Network Layer   4a-8

## RIP (Routing Information Protocol)



| Destination Network | Next Router | Num. of hops to dest. |
|---|---|---|
| 1 | A | 2 |
| 20 | B | 2 |
| 30 | B | 7 |
| 10 | -- | 1 |
| …. | …. | |

4: Network Layer 4a-9

## RIP: Link Failure and Recovery

r If no advertisement heard after 180 sec, neighbor/link dead

r Routes via the neighbor are invalidated; new advertisements sent to neighbors

r Neighbors in turn send out new advertisements if their tables changed

r Link failure info quickly propagates to entire net

r Poison reverse used to prevent ping-pong loops (infinite distance = 16 hops)

4: Network Layer 4a-10

## RIP Table processing

r  RIP routing tables managed by an **application** process called route-d (daemon)

r  Advertisements encapsulated in UDP packets (no reliable delivery required; advertisements are periodically repeated)

4: Network Layer  4a-11

## RIP Table processing



4: Network Layer  4a-12

## RIP Table example (continued)

RIP Table example
  (at router *giroflee.eurocom.fr*):

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  4a-13

## RIP Table example

```
  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
```

4: Network Layer  4a-14

## OSPF (Open Shortest Path First)

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

r OSPF advertisement carries one entry per neighbor router
r Advertisements disseminated to entire Autonomous System (via flooding)
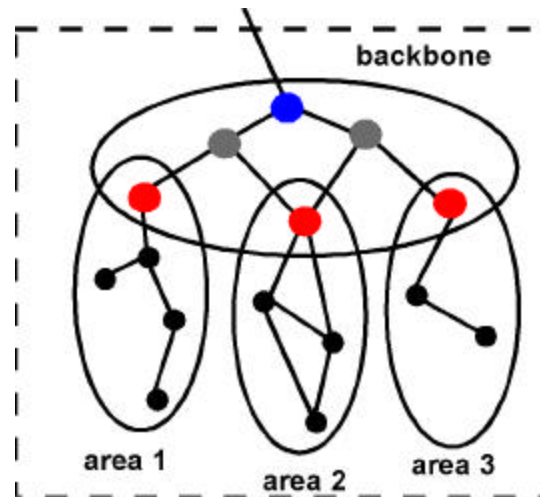
4: Network Layer  4a-15

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

r Security: all OSPF messages are 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  4a-16

## Hierarchical OSPF



backbone

area 1
area 2
area 3

4: Network Layer  4a-17

## Hierarchical OSPF

r  Two-level hierarchy: local area and backbone.
r  Link-state advertisements do not leave respective areas.
r  Nodes in each area have detailed area topology; they only know direction (shortest path) to networks in other areas.
r  **Area Border routers** "summarize" distances  to networks in the area and advertise them to other Area Border routers.
r  **Backbone routers** run an OSPF routing alg limited to the backbone.
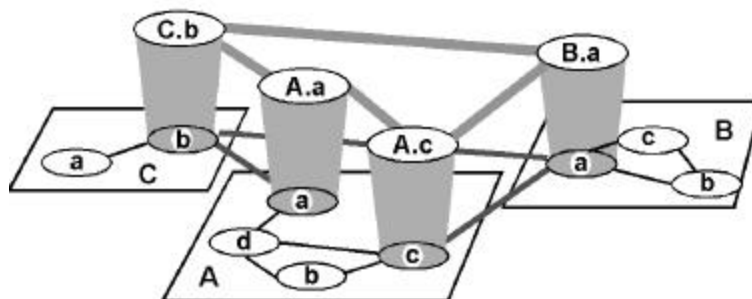r  **Boundary routers** connect to other ASs.

4: Network Layer  4a-18

## 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  routing tables exchanged only when costs change
r  Loop-free routing achieved by using a Distributed Updating Alg. (DUAL) based on *diffused computation*
r   In DUAL, after a distance increase, the routing table is *frozen* until all affected nodes have learned of the change.

4: Network Layer  4a-19

## Inter-AS routing



4: Network Layer  4a-20

## Inter-AS routing (cont)

r  BGP (Border Gateway Protocol): the de facto standard
r  **Path Vector** protocol: and extension of Distance Vector
r  Each Border Gateway broadcast to neighbors (peers) the entire path (ie, sequence of ASs) to destination
r  For example, Gateway X may store the following path to destination Z:

Path (X,Z) = X,Y1,Y2,Y3,…,Z

4: Network Layer  4a-21

## Inter-AS routing (cont)

r  Now,  suppose Gwy X send its path to peer Gwy W
r  Gwy W may or may not select the path offered by Gwy X, because of cost, policy ($$$$) or loop prevention reasons.
r  If Gwy W selects the path advertised by Gwy X, then:

Path (W,Z) = w, Path (X,Z)

Note: path selection based not so much on cost (eg,# of AS hops), but mostly on administrative and policy issues (e.g., do not route packets through competitor's AS)

4: Network Layer  4a-22

## Inter-AS routing (cont)

r Peers exchange BGP messages using TCP.

r OPEN msg opens TCP connection to peer and authenticates sender

r UPDATE msg advertises new path (or withdraws old)

r KEEPALIVE msg keeps connection alive in absence of UPDATES; it also serves as ACK to an OPEN request

r NOTIFICATION msg reports errors in previous msg; also used to close a connection

4: Network Layer  4a-23

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

r **Policy**: Inter is concerned with policies (which provider we must select/avoid, etc). Intra is contained in a single organization, so, no policy decisions necessary

r **Scale**: Inter provides an extra level of routing table size and routing update traffic reduction above the Intra layer

r **Performance**: Intra is focused on performance metrics; needs to keep costs low. In Inter it is difficult to propagate performance metrics efficiently (latency, privacy etc). Besides, policy related information is more meaningful.

We need **BOTH!**                    4: Network Layer  4a-24