

# Chapter 4: Network Layer

## Chapter goals:

- \* understand principles behind network layer services:
  - routing (path selection)
  - dealing with scale
  - how a router works
  - advanced topics: IPv6, multicast
- \* instantiation and implementation in the Internet

## Overview:

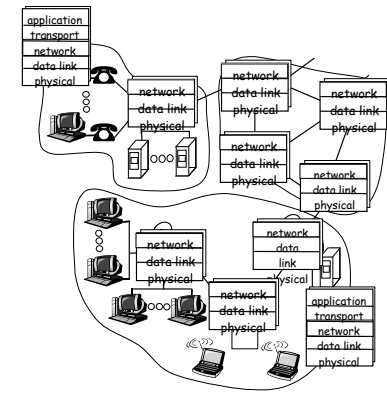
- \* network layer services
- \* routing principle: path selection
- \* hierarchical routing
- \* IP
- \* Internet routing protocols
- \* reliable transfer
  - intra-domain
  - inter-domain
- \* what's inside a router?
- \* IPv6
- \* multicast routing

## Network layer functions

- \* transport packet from sending to receiving hosts
- \* network layer protocols in every host, router

### three important functions:

- \* **path determination:** route taken by packets from source to dest. *Routing algorithms*
- \* **switching:** move packets from router's input to appropriate router output
- \* **call setup:** some network architectures require router call setup along path before data flows



## Network service model

**Q: What service model for “channel” transporting packets from sender to receiver?**

- \* guaranteed bandwidth?
- \* preservation of inter-packet timing (no jitter)?
- \* loss-free delivery?
- \* in-order delivery?
- \* congestion feedback to sender?

service abstraction

The most important abstraction provided by network layer:

virtual circuit  
or  
datagram?

## Virtual circuits

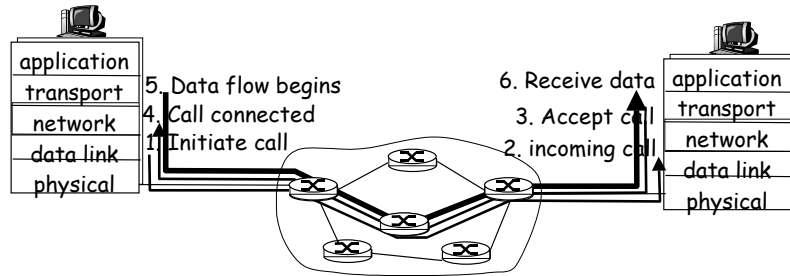
“source-to-dest path behaves much like telephone circuit”

- performance-wise
- network actions along source-to-dest path

- \* call setup, teardown for each call *before* data can flow
- \* each packet carries VC identifier (not destination host OD)
- \* every router on source-dest path s maintain “state” for each passing connection
  - transport-layer connection only involved two end systems
- \* link, router resources (bandwidth, buffers) may be *allocated* to VC
  - to get circuit-like perf.

## Virtual circuits: signaling protocols

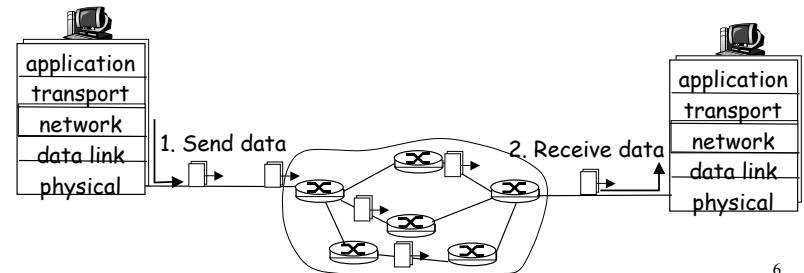
- used to setup, maintain teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet



5

## Datagram networks: the Internet model

- no call setup at network layer
- routers: no state about end-to-end connections
  - no network-level concept of "connection"
- packets typically routed using destination host ID
  - packets between same source-dest pair may take different paths



6

## Datagram or VC network: why?

### Internet

- data exchange among computers
  - "elastic" service, no strict timing req.
- "smart" end systems (computers)
  - can adapt, perform control, error recovery
  - simple inside network, complexity at "edge"
- many link types
  - different characteristics
  - uniform service difficult

### ATM

- evolved from telephony
- human conversation:
  - strict timing, reliability requirements
  - need for guaranteed service
- "dumb" end systems
  - telephones
  - complexity inside network

7

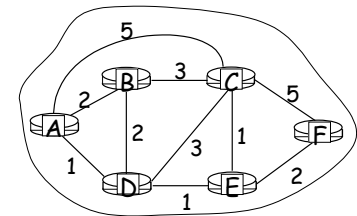
## Routing

### Routing protocol

Goal: determine "good" path (sequence of routers) thru network from source to dest.

### Graph abstraction for routing algorithms:

- graph nodes are routers
- graph edges are physical links
  - link cost: delay, \$ cost, or congestion level



"good" path:  
typically means  
minimum cost path  
other def's possible

8

## Routing Algorithm classification

### Global or decentralized information?

#### Global:

- ✳ all routers have complete topology, link cost info
- ✳ “link state” algorithms

#### Decentralized:

- ✳ router knows physically-connected neighbors, link costs to neighbors
- ✳ iterative process of computation, exchange of info with neighbors
- ✳ “distance vector” algorithms

### Static or dynamic?

#### Static:

- ✳ route selection based on static metrics (hop count)

#### Dynamic:

- ✳ Route selection based on performance metrics (traffic load, delay, queue size)
  - periodic update
  - in response to link cost changes

9

## Two fundamental routing approaches

### ✳ Link State:

- Each node knows the whole graph
- Advertise each link to everybody
- Find path by running the same algorithm

### ✳ Distance Vector:

- Each node has a vector of distances to all others nodes
- Advertise only your vector of distances to neighbors
  - *Forward info only if your vector of distance changed*
- Determine next hop according to distance
  - *Each node only is aware of the next hop*

10

## Distance Vector Routing Algorithm

### iterative:

- ✳ continues until no nodes exchange info.
- ✳ *self-terminating*: no “signal” to stop

### asynchronous:

- ✳ nodes need *not* exchange info/iterate in lock step!

### distributed:

- ✳ each node communicates *only* with directly-attached neighbors

### Distance Table data structure

- ✳ each node has its own
- ✳ row for each possible destination
- ✳ column for each directly-attached neighbor to node
- ✳ example: in node X, for dest. Y via neighbor Z:

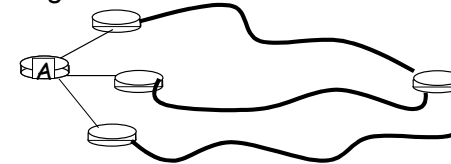
$$D^X(Y,Z) = \text{distance from X to Y, via Z as next hop}$$

$$= c(X,Z) + \min_w \{D^Z(Y,w)\}$$

11

## Main Idea

Edge + Shortest Path



### ✳ My shortest path thru each neighbor to X

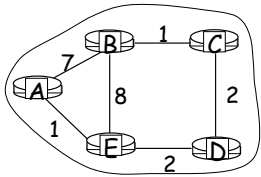
- Shortest path from a neighbor to X
- Plus my cost to reach the neighbor

### ✳ My shortest path overall

- minimum among shortest path from each neighbor

12

## Distance Table: example



$$D^E(C,D) = c(E,D) + \min_w \{D^D(C,w)\} = 2+2 = 4$$

From node E: to go to C  
What is the cost thru D?

D <sup>E</sup> ()	Via		
	A	B	D
A	①	14	5
B	7	8	⑤
C	6	9	④
D	4	11	②

13

## Distance table gives routing table

D <sup>E</sup> ()	cost to destination via			Outgoing link to use, cost	
	A	B	D	destination	
A	①	14	5	A	A,1
B	7	8	⑤	B	D,5
C	6	9	④	C	D,4
D	4	11	②	D	D,2

Distance table → Routing table

14

## Distance Vector Routing: overview

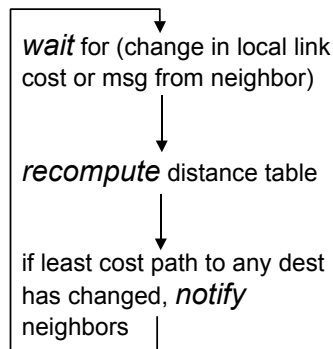
**Iterative, asynchronous:**  
each local iteration caused by:

- local link cost change
- message from neighbor: its least cost path change from neighbor

**Distributed:**

- each node notifies neighbors **only** when its least cost path to any destination changes
  - neighbors then notify their neighbors if necessary

Each node:

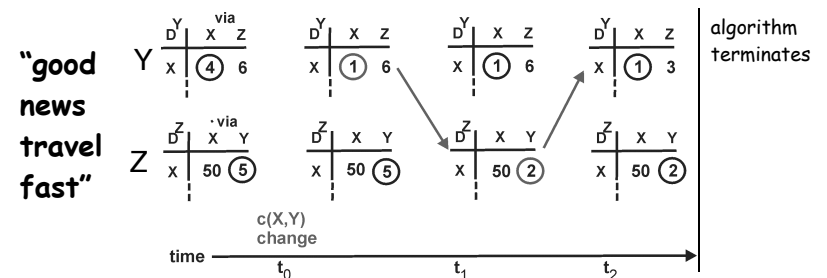
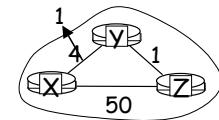


15

## Distance Vector: link cost changes

Link cost changes:

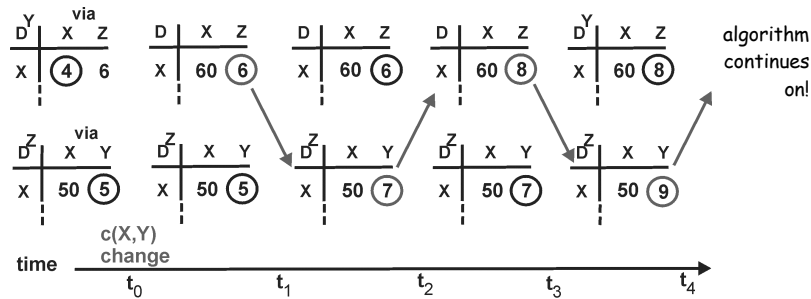
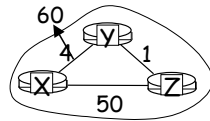
- node detects local link cost change updates distance table (line 15)
- if cost change in least cost path, notify neighbors (lines 23,24)



16

## Distance Vector: link cost changes

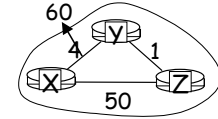
**Link cost changes:**  
bad news travels slow -  
“count to infinity” problem



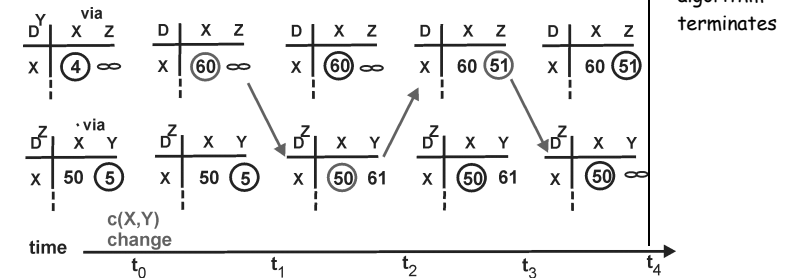
17

## Distance Vector: poisoned reverse

**Poison Reverse:** If Z routes through Y to get to X :  
Z tells Y its distance to X is infinite  
(so Y won't route to X via Z)



**Split Horizon:** don't advertise at all distance to next hop neighbor (for path)



18

## Link-State Routing Algorithm

- Net topology, link costs known to all nodes
  - accomplished via “link state broadcast”
  - all nodes have same info
- Each node computes locally: cost paths from to all other nodes
  - gives routing table for that node

19

## Dijkstra's algorithm, discussion

**Algorithm complexity:** n nodes

- each iteration: need to check all nodes, w, not in N
- $n*(n+1)/2$  comparisons:  $O(n^2)$
- more efficient implementations possible:  $O(n \log n)$

**Oscillations possible:**

- if link cost = amount of carried traffic

20

## Link State Routing

---

- ✱ **Initial state** : similar to distance vector i.e., state of link to neighbors known (up/down).
- ✱ **Goal**: To find the path of least cost to destination.
- ✱ **Basic Idea** --
  - Every node knows how to reach its neighbors.
  - Disseminate the info is disseminated to every node,
  - every node has the complete map of the network
  - All nodes execute the same algorithm for paths

21

## Two Mechanisms

---

- ✱ **Reliable dissemination of link state information**
  - The process is called reliable flooding.
- ✱ **Calculation of routes using the collected information**
  - The computation is based on Dijkstra's algorithm.

22

## Reliable Flooding

---

- ✱ **Process of making sure that all the nodes participating in the link state routing protocol get a copy of the link-state info. from all other nodes.**
- ✱ **Each node sends out link-state info. on its directly connected links.**
- ✱ **Each node that receives this, forwards it to all links**
  - except the incoming one.

23

## Link State Information

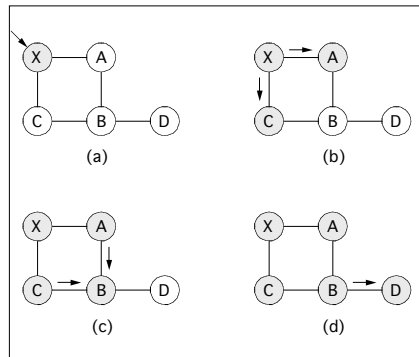
---

- ✱ **Each node creates a link-state packet (LSP) that contains:**
    - ID of the node that created LSP
    - a list of directly connected nodes and the cost to each node.
    - sequence number
    - TTL
- } → for reliability

24

## An Example

- ✳ X receives info from node Y.
- ✳ X checks to see if it already has an update for that link. If it does, it compares the sequence number in the new LSP to the one stored.
- ✳ If New seq no < Old sequence number, then, discard LSP.
- ✳ Else -- store LSP and send the LSP to all neighbors except the one that sent the LSP.



25

## Dissemination of LSPs

- ✳ LSPs are sent
  - periodically (upon the expiry of a timer order of hours)
  - Or may be triggered due to a change in topology (as in RIP).
- ✳ Topology change in a directly connected link
  - Failures detected by link layer protocol by using what are known as "HELLO" packets -- probes to determine if neighbor is alive.
- ✳ Sequence numbers help in identifying new info and TTL helps in ensuring that packets don't percolate in the network indefinitely.

26

## Graph abstraction

- ✳ Used for computation of shortest path using Dijkstra's.
- ✳ Let  $N$  denote the set of nodes in the graph.
- ✳  $I(i,j)$  denotes the non-negative cost or weight associated with the edge between nodes  $i$  and  $j$ .
- ✳  $I(i,j) = \infty$  if there is no edge between  $i$  and  $j$ .
- ✳ Remember -- each node has entire map of network.

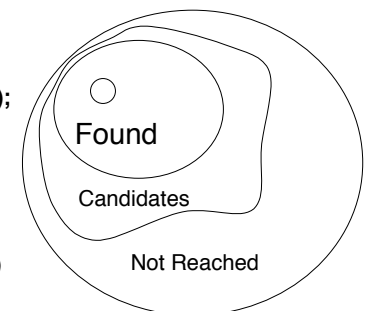
27

## The Dijkstra's Algorithm

- ✳ Let "s" the node that executes the algorithm.
- ✳ Algorithm maintains  $M$  --> set of nodes whose shorter path I found
- ✳  $C(n)$  -- cost of the path from  $s$  to  $n$ .
- ✳  $L(v,w)$  cost of link from  $v$  to  $w$
- ✳ Reminder of Dijkstra algorithm
 

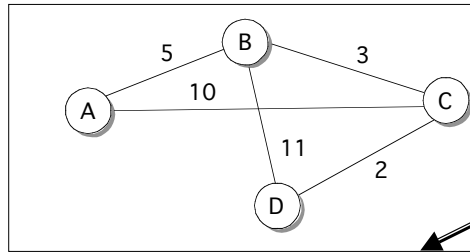
```

M = {s}
for each node n in N - {s}, C(n) = L(s,n);
while (N ≠ M),
    find node w in N - M, with cost C(w)
    M = M U {w}
    for each n in (N-M)
        C(n) = min( C(n), C(w) + L(w,n) )
            
```



28

## An Example



- Let  $s = A$ .
- Initially  $M = A$ .
- Now, minimum cost node is B (note that only two nodes can be considered). So add B to M.

• At next step, look at neighbors of A and B (we have C and D). We add C (via B).

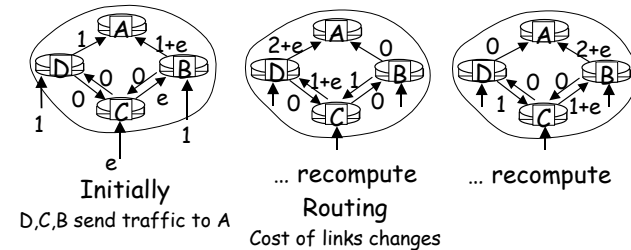
• Then, we consider D, (note that we consider the link via C as well -- finally that is added.



• Finally, the tree with links AB, BC and CD is created.

29

## Routing Oscillations with dynamic weights



✳ If link weight is dynamic (function of load)

✳ Paths can change after we route

- Example: D,C,B send traffic to A continuously
- Once we route, best path changes!

30

## Two fundamental routing approaches

### ✳ Link State:

- Each node knows the whole graph
- Advertise each link to everybody
- Find path by running the same algorithm

### ✳ Distance Vector:

- Each node has a vector of distances to all others
- Advertise only your vector of distances to neighbors
  - Forward info only if your vector of distance changed
- Determine next hop according to distance
  - Each node only is aware of the next hop

31

## Summary: Distributed Routing Techniques

### Link State

- ✳ Topology information is flooded within the routing domain
- ✳ Best end-to-end paths are computed locally at each router.
- ✳ Best end-to-end paths determine next-hops.
- ✳ Advertises: link info
- ✳ Works only if policy is shared and uniform
- ✳ Examples: OSPF, IS-IS

### Vectoring

- ✳ Each router knows little about network topology
- ✳ Only best next-hops are chosen by each router for each destination network.
- ✳ Best end-to-end paths result from composition of all next-hop choices
- ✳ Advertises: path info, distance
- ✳ Does not require uniform policies at all routers
- ✳ Examples: RIP, BGP

32

# Routing in the Global Internet

# Hierarchical Routing

Our routing study thus far - idealization  
 all routers identical  
 network “flat”  
 ... *not* true in practice

- scale: with 50 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

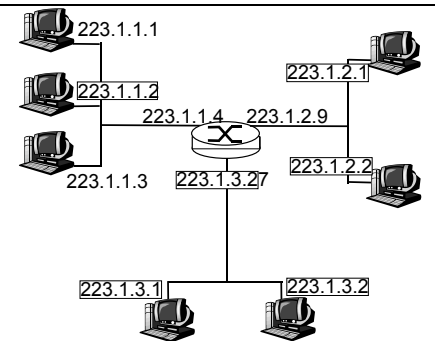
- ✱ **Aggregate routers into regions, “autonomous systems” (AS)**
- ✱ **Routers in same AS run same routing protocol**
  - “intra-AS” routing protocol
  - different AS can run different intra-AS protocol
- ✱ **Aggregatable addresses**

**gateway routers**

- ✱ special routers in AS
- ✱ run intra-AS routing protocol with all other routers in AS
- ✱ **also responsible for routing to destinations outside AS**
  - run *inter-AS routing* protocol with other gateway routers

# IP Addressing: introduction

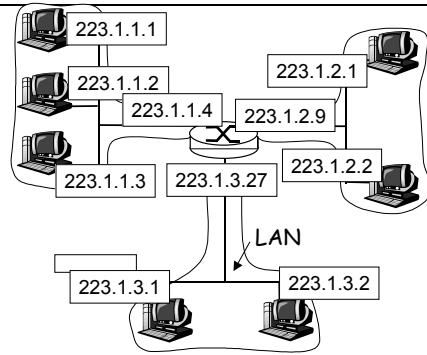
- ✱ **IP address: 32-bit identifier for host, router *interface***
- ✱ ***interface*: connection between host, router and physical link**
  - router's typically have multiple interfaces
  - host may have multiple interfaces
  - IP addresses associated with interface, not host, router



$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_{1} \underbrace{00000001}_{1} \underbrace{00000001}_{1}$$

# IP Addressing

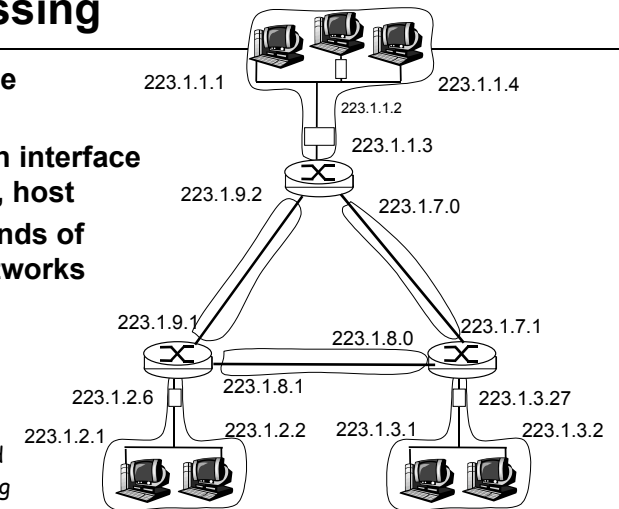
- IP address:**
  - network part (high order bits)
  - host part (low order bits)
- What's a network ?**  
 (from IP address perspective)
  - device interfaces with same network part of IP address
  - can physically reach each other without intervening router



network consisting of 3 IP networks  
 (for IP addresses starting with 223,  
 first 24 bits are network address)

# IP Addressing

- How to find the networks?**
- Detach each interface from router, host
  - create "islands of isolated networks"

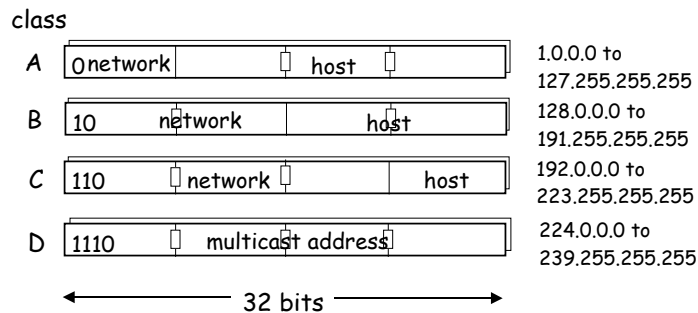


Interconnected system consisting of six networks

# IP Addresses

given notion of "network", let's re-examine IP addresses:

"class-full" addressing:

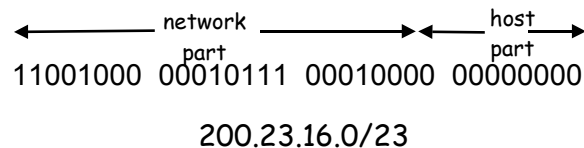


# IP addressing: need for change

- Problems with classful addressing:**
  - inefficient use of address space, address space exhaustion
  - e.g., class B net allocated enough addresses for 65K hosts, even if only 2K hosts in that network

## IP addressing: CIDR

- ☛ CIDR: Classless InterDomain Routing
- ☛ network portion of address of arbitrary length
- ☛ address format: a.b.c.d/x, where x is # bits in network portion of address



41

## IP Addresses and Prefixes

- ☛ IP addresses have 32 bits: 4 octets of bits (IPv4)
- ☛ A prefix is a group of IP addresses
- ☛ 128.32.101.5 is an IP address (32 bits)
- ☛ 128.32.0.0/16 is a prefix of the 16 first bits:
  - 128.32.0.0 – 128.32.255.255 (2<sup>16</sup> addresses)
- ☛ 128.32.4.0/24 is a prefix of the 24 first bits - longer

42

## IP addresses: how to get one?

### Hosts (host portion):

- ☛ hard-coded by system admin in a file
- ☛ DHCP: Dynamic Host Configuration Protocol: dynamically get address: “plug-and-play”
  - host broadcasts “DHCP discover” msg
  - DHCP server responds with “DHCP offer” msg
  - host requests IP address: “DHCP request” msg
  - DHCP server sends address: “DHCP ack” msg

43

## IP addresses: how to get one?

### Network (network portion):

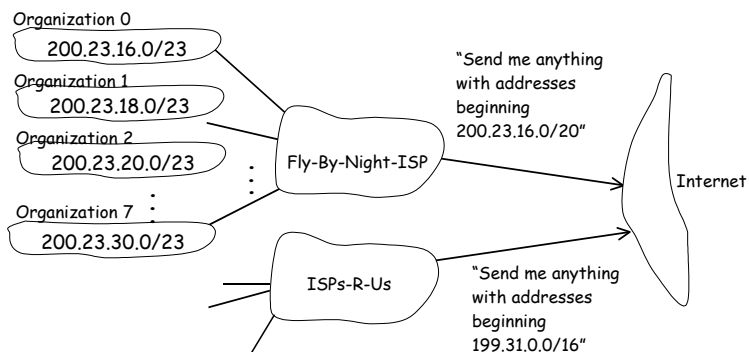
- ☛ get allocated portion of ISP’s address space:

ISP's block	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/20
Organization 0	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/23
Organization 1	<u>11001000</u>	<u>00010111</u>	<u>00010010</u>	00000000	200.23.18.0/23
Organization 2	<u>11001000</u>	<u>00010111</u>	<u>00010100</u>	00000000	200.23.20.0/23
...	.....	.....	.....	.....	.....
Organization 7	<u>11001000</u>	<u>00010111</u>	<u>00011110</u>	00000000	200.23.30.0/23

44

## Hierarchical addressing: route aggregation

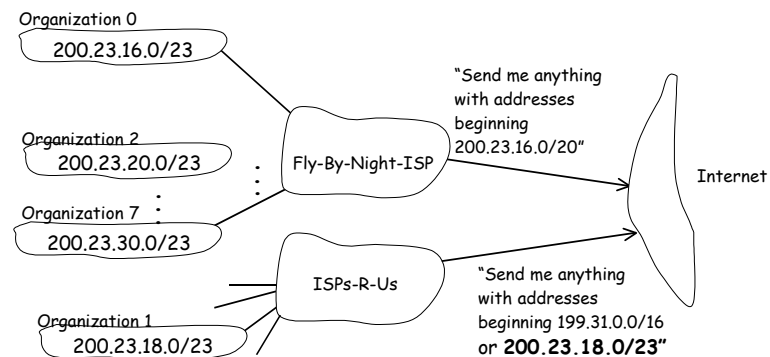
Hierarchical addressing allows efficient advertisement of routing information:



45

## Hierarchical addressing: more specific routes

ISPs-R-Us has a more specific route to Organization 1



46

## Routing is Based on Prefixes

- ✱ A BGP Routing table has prefixes for entries
- ✱ For a IP address of a packet, find longest match
- ✱ Example: packet IP 128.32.101.1
- ✱ Matching:
- ✱ 128.1.1.4 matches the first 8 bits – no match!
- ✱ 128.32.0.0/16 match for 16 bits
- ✱ 128.32.101.0/24 is a longer match

47

## Prefix Matching in More Detail

- ✱ For a IP address of a packet, find longest match
- ✱ Example: Compare
  - packet IP 128.32.101.1
  - With 128.32.0.0/16
  - IP : 01000000. 00100000. 01100101 .00000001
  - Mask : 11111111. 11111111. 00000000 .00000000
  - AND : 01000000. 00100000. 00000000 .00000000
  - Prefix : 01000000. 00100000. 00000000. 00000000
  - Equal? Yes
  - We have a match of length 16

48

## IP addressing: the last word...

---

**Q: How does an ISP get block of addresses?**

**A: ICANN: Internet Corporation for Assigned**

### Names and Numbers

- allocates addresses
- manages DNS
- assigns domain names, resolves disputes

49

## ICMP: Reporting errors

---

- ✱ **What happens if things cannot proceed as expected?**
- ✱ **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)

50

## ICMP: Internet Control Message Protocol

---

- ✱ **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

51

## Let's see how things work in practice

---

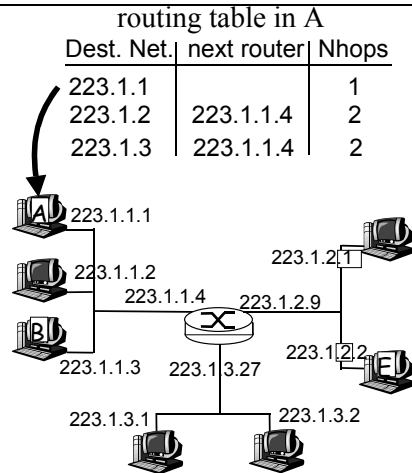
52

## Getting a datagram from source to dest.

### IP datagram:

misc	source	dest	data
fields	IP addr	IP addr	

datagram remains unchanged, as it travels source to destination  
addr fields of interest here



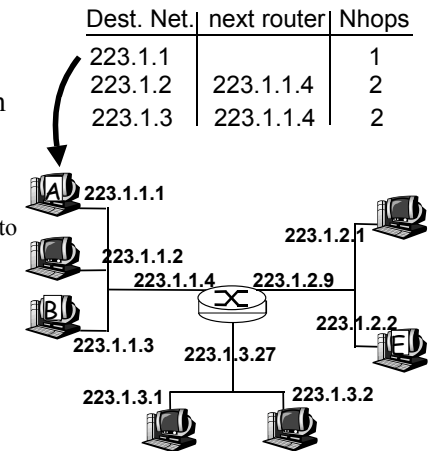
53

## Getting a datagram from source to dest.

misc	223.1.1.1	223.1.1.3	data
fields			

Starting at A, given IP datagram addressed to B:

look up net. address of B  
find B is on same net. as A  
link layer will send datagram directly to B inside link-layer frame  
B and A are directly connected



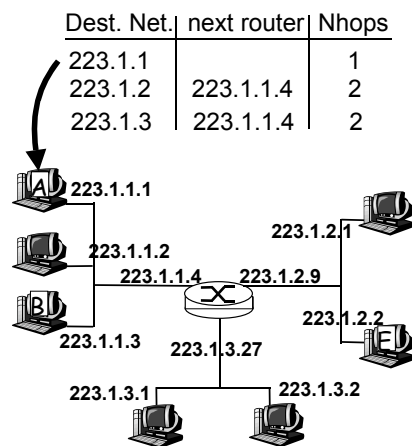
54

## Getting a datagram from source to dest.

misc	223.1.1.1	223.1.2.2	data
fields			

Starting at A, dest. E:  
look up network address of E  
E on *different* network

A, E not directly attached  
routing table: next hop router to E is 223.1.1.4  
link layer sends datagram to router 223.1.1.4 inside link-layer frame  
datagram arrives at 223.1.1.4  
continued.....



55

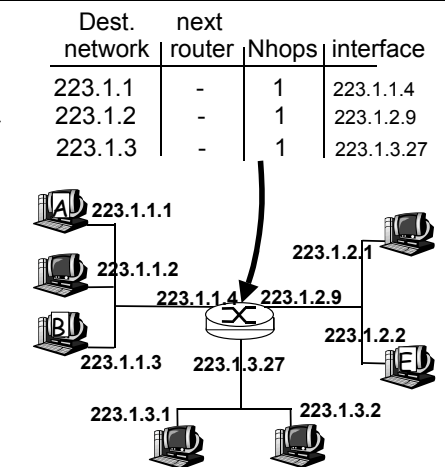
## Getting a datagram from source to dest.

misc	223.1.1.1	223.1.2.2	data
fields			

Arriving at 223.1.4, destined for 223.1.2.2

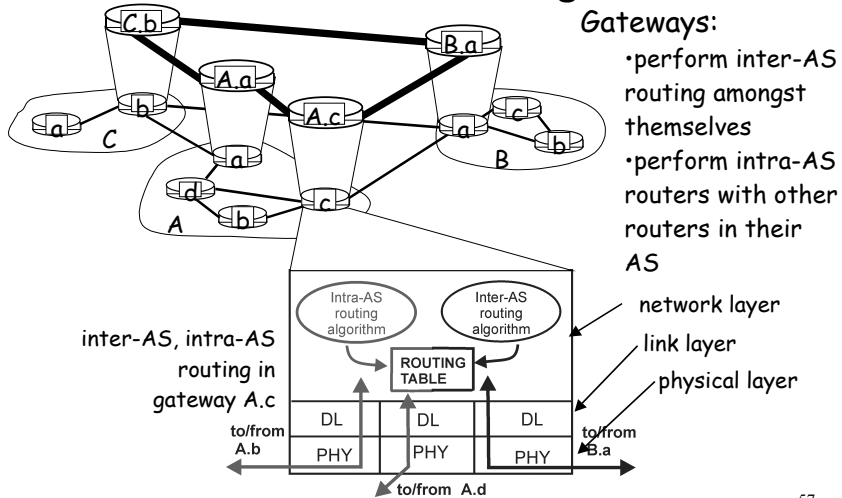
look up network address of E  
E on *same* network as router's interface 223.1.2.9

router, E directly attached  
link layer sends datagram to 223.1.2.2 inside link-layer frame via interface 223.1.2.9  
datagram arrives at 223.1.2.2!!!  
(hooray!)



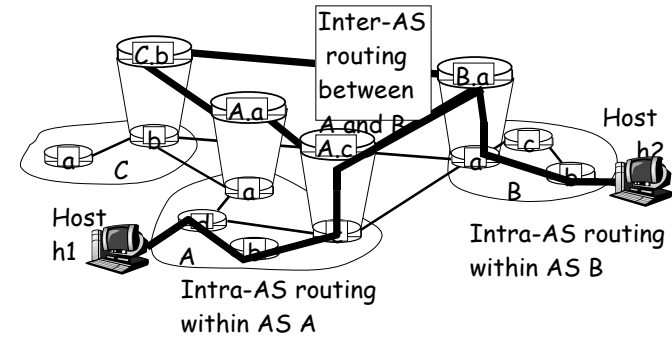
56

## Intra-AS and Inter-AS routing



57

## Intra-AS and Inter-AS routing

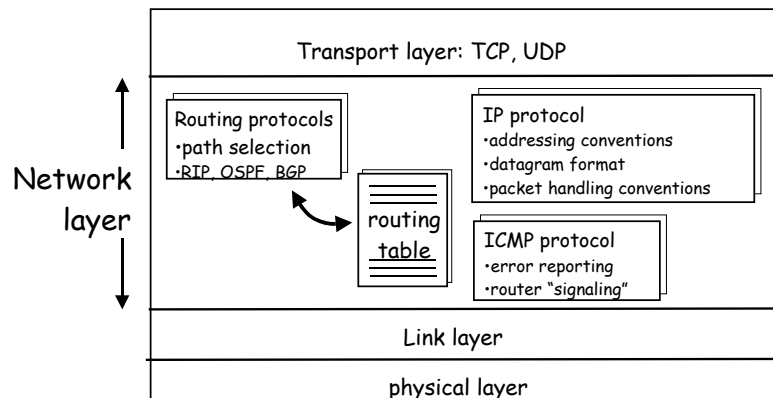


We'll examine specific inter-AS and intra-AS Internet routing protocols shortly

58

## The Internet Network layer

Host, router network layer functions:



59

## Measurements in the Internet

- ✱ Difficulties in measuring
- ✱ Measuring tools (traceroute)
- ✱ Misc issues

60

## Measuring and Modeling Is not Easy

---

- ✱ **Constantly changing environment**
- ✱ **How much data is enough**
  - Recently: we need to measure more than 24h!
- ✱ **How frequently should I be measuring?**
- ✱ **Are the measurements representative?**

61

## Operation versus Measurements

---

- ✱ **Operators do not care about**
  - Measurements
  - Academic Research
- ✱ **Why?**
  - Takes away resources
  - Can create problems
  - Complicates their lives
- ✱ **Luckily, there are measurement centers**
  - CAIDA, NLANR, routeviews, RIPE

62

## Types of Measurement Tools

---

- ✱ **Application level:**
  - Install application agents at two measuring entries
  - REALITI tool from UCR
  - More control over process
- ✱ **Network level:**
  - Use the Internet control functionality (ICMP)
  - Trick the network to provide information

63

## Ping: the tool

---

- ✱ **Uses ICMP ECHO\_REQUEST datagram to elicit an ICMP ECHO\_RESPONSE from a host or gateway**
- ✱ **Reports**
  - Round trip time
  - Packets loss
- ✱ **Many available options: packet type, size etc**
- ✱ **Limitation: >1sec measurement frequency**
- ✱ **Read manual: man ping**

64

## Traceroute: the tool

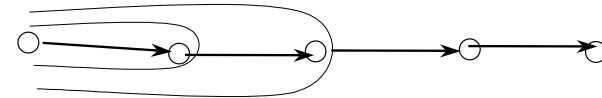
---

- ✦ **Traceroute measures**
  - the path and the round trip time
- ✦ **Traceroute: ingenious (ab)use of the network layer by Van Jacobson**
- ✦ **Main ideas:**
  - send “bad” packets to receive ICMP: “packet died”
  - Recursive probing to identify the path
  - Send three packets at a time
- ✦ **Read manual: man traceroute**

65

## The ingenuity of traceroute

---

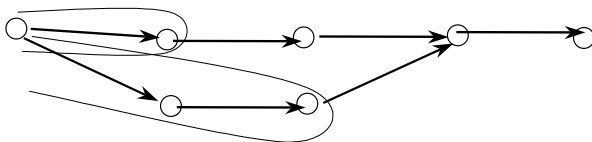


- ✦ **Send a packet for every hop of the path**
- ✦ **Set TTL = 1, packet expires, ICMP returns**
- ✦ **Increase TTL by one, and repeat**
- ✦ **At the destination, port number is wrong: return an ICMP packet, port not found**

66

## Traceroute: Some Limitations

---



- ✦ **In traceroute, you may be exploring multiple paths without knowing it**
- ✦ **Delays for each part of the path correspond to different measurements: ie they don't sum up**

67

## Identifying The Router Topology

---

- ✦ **Several efforts rely on multiple traceroute**
  - Govindan et al INFOCOM 2000
  - Cheswick and Burch Internet Mapping Project
- ✦ **Main idea:**
  - Do thousands of traceroutes
  - Collect all adjacent nodes
  - Generate a graph

68

## Router Graphs: A Complication

---

- ✱ **Routers have multiple IP addresses**
  - One for each interface
- ✱ **How do we resolve this?**
- ✱ **Only heuristics exist [Govindan]**
- ✱ **Heuristic: Send packets to one interface and hope that they will respond with the other interface**
  - Typically, router responds with IP of interface the packet came on

69

## End

---

70

## Proposals

---

- ✱ **Overall: decent**
- ✱ **Not enough motivation/background**
  - All related papers
- ✱ **Not enough thought on what you will do**
  - Spend an evening thinking what you will do and how
- ✱ **Try to clarify goals early and talk to me**
- ✱ **Photocopy them, and give me the original**
  - It's the contract

71

## Practical Tips

---

- ✱ **The earlier, the better**
- ✱ **Talk to me early**
- ✱ **Look for tools, data, previous work**
  - It can save you a lot of time in the long run
- ✱ **Try to focus on a topic**

72

## Side Note

---

### ✿ Important things in research:

- Asking the right questions
- Identifying the right topic
- Context of research
  - *Motivation*
  - *Previous and related work*
  - *Importance of problem*
- Thoroughness of work

73

## Quiz

---

- ✿ **What is the main role of routing?**
- ✿ **Describe a centralized routing protocol**
  - Main operation functions
  - Comment on complexity, robustness
- ✿ **Describe a table driven routing protocol**
  - Main functions
  - Comment on complexity
- ✿ **How can the Internet routing scale?**
  - Describe the elements that make it work
  - What does a router need to keep in its memory
  - Describe a routing table
- ✿ **What happens to a packet when it arrives at a router?**
  - What kind of “hardware” does it go through?
  - Where does the delay and loss come into play?

74