

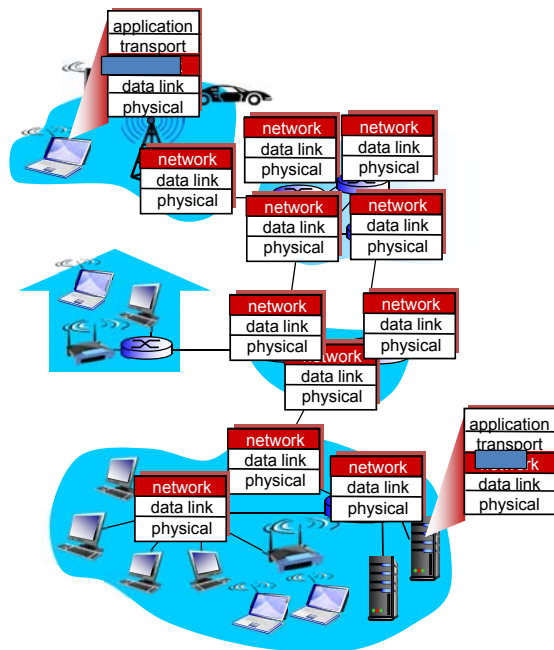
# Network Layer

- Reference
  - Chapter 04, *Computer Networking: A Top Down Approach*, 6/E, Jim Kurose, Keith Ross, Addison-Wesley
  - Adapted from part of the slides provided by the authors

CS3700 - Instructor: Dr. Zhu

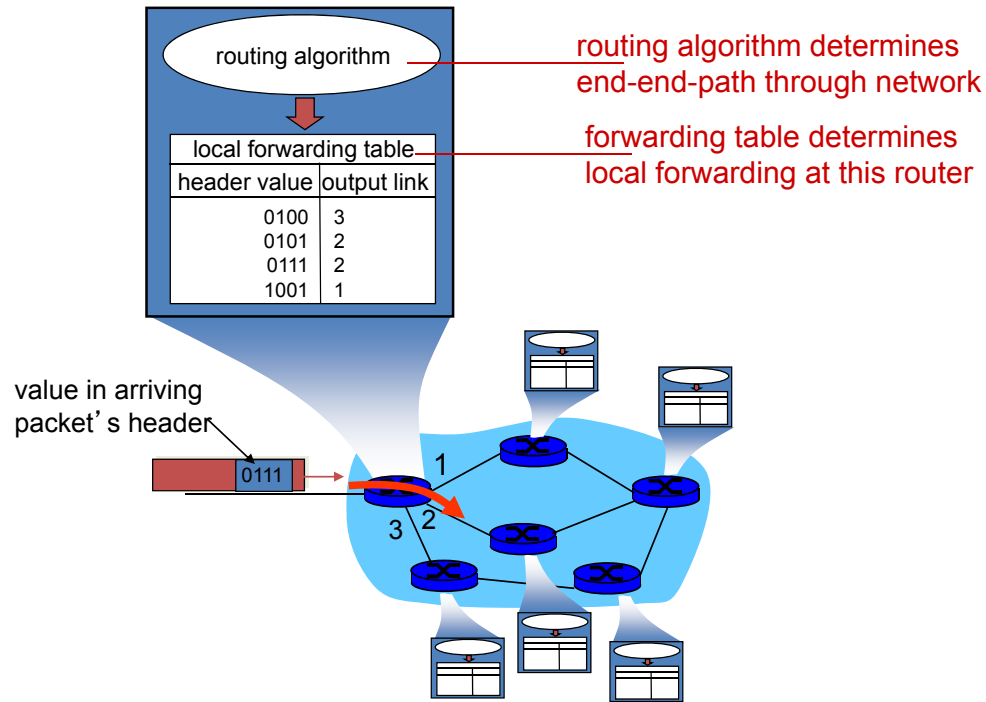
## Network layer

- 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



CS3700 - Instructor: Dr. Zhu

# Interplay between routing and forwarding



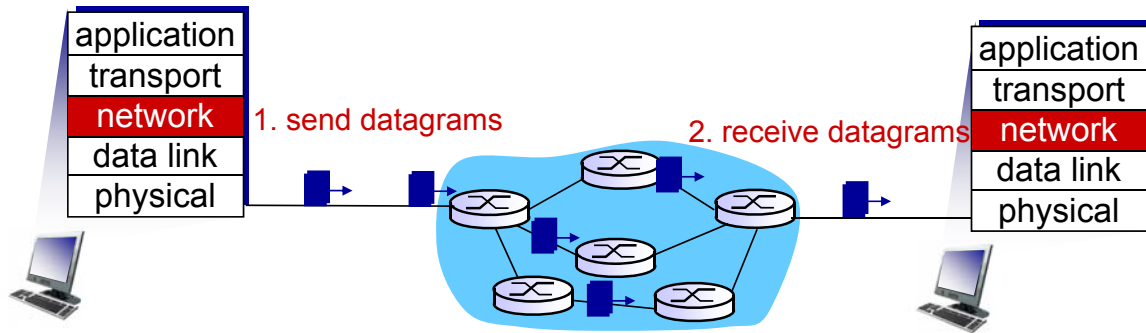
CS3700 - Instructor: Dr. Zhu

## Connection, connection-less service

- *datagram* network provides network-layer *connectionless* service
  - used in today's Internet
- *virtual-circuit* network provides network-layer *connection* service
  - used in ATM, frame-relay, X.25
  - not used in today's Internet
- analogous to TCP/UDP connection-oriented / connectionless transport-layer services, but:
  - *service*: host-to-host
  - *no choice*: network provides one or the other
  - *implementation*: in network core

CS3700 - Instructor: Dr. Zhu

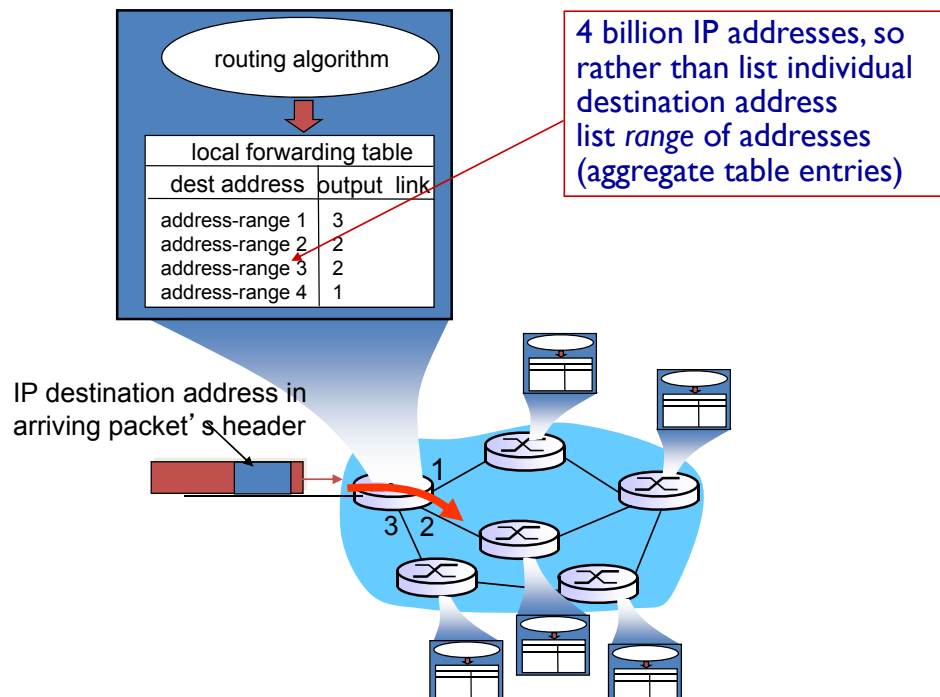
# Datagram networks



- no call setup at network layer
- routers: no state about end-to-end connections
  - no network-level concept of “connection”
- packets forwarded using destination host address

CS3700 - Instructor: Dr. Zhu

## Datagram forwarding table



CS3700 - Instructor: Dr. Zhu

# Datagram forwarding table

Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

**Q:** but what happens if ranges don't divide up so nicely?

CS3700 - Instructor: Dr. Zhu

## Longest prefix matching

### *longest prefix matching*

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

11001000 00010111 00011000 \*\*\*\*\* is excluded  
from 11001000 00010111 00011\*\*\* \*\*\*\*\*

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

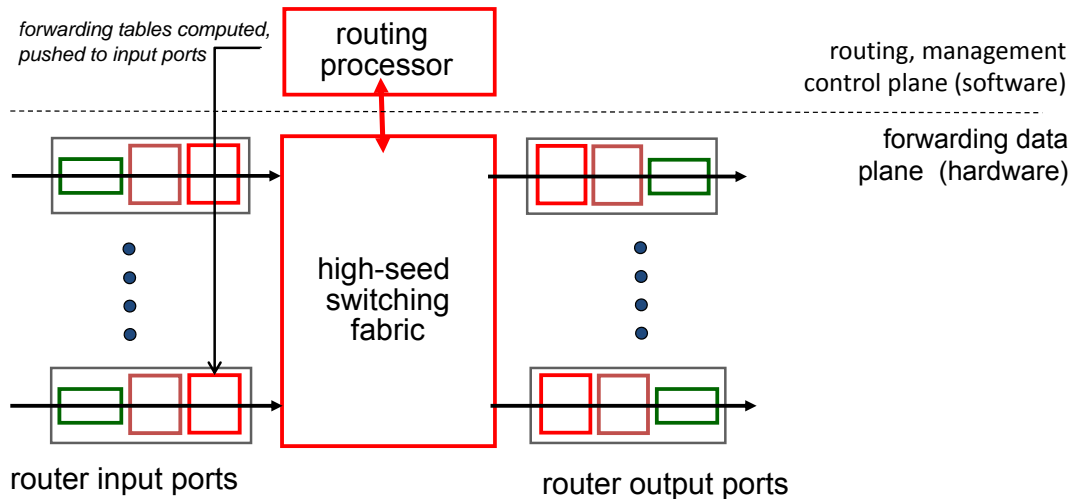
which interface?

DA: 11001000 00010111 00011000 10101010

which interface?

CS3700 - Instructor: Dr. Zhu

# Router architecture overview

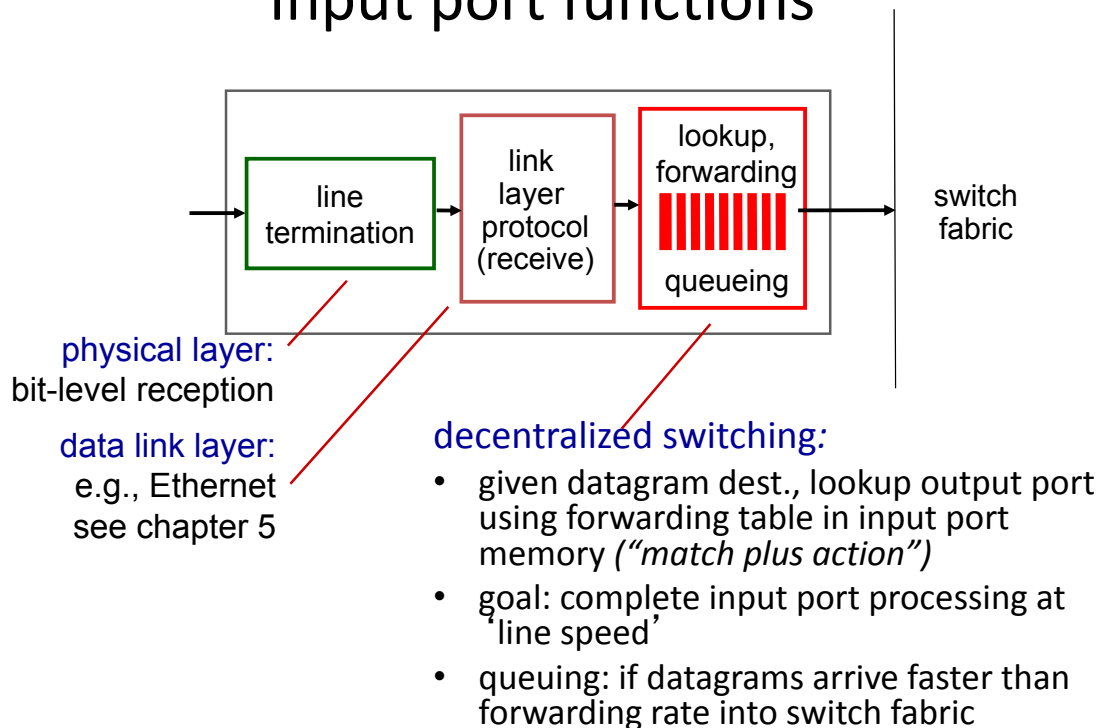


two key router functions:

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

CS3700 - Instructor: Dr. Zhu

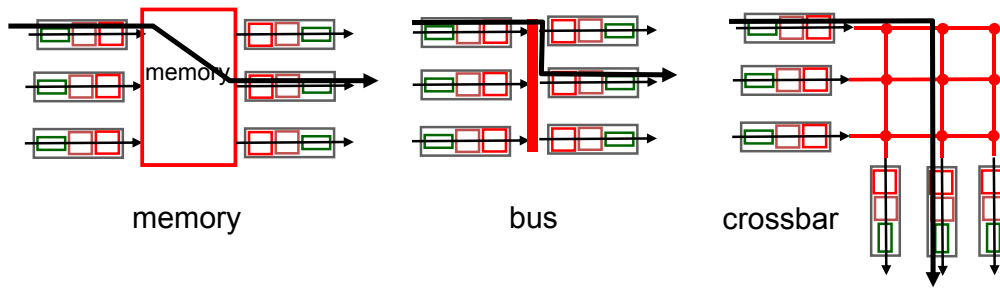
## Input port functions



CS3700 - Instructor: Dr. Zhu

# Switching fabrics

- transfer packet from input buffer to appropriate output buffer
- switching rate: rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: switching rate N times line rate desirable
- three types of switching fabrics

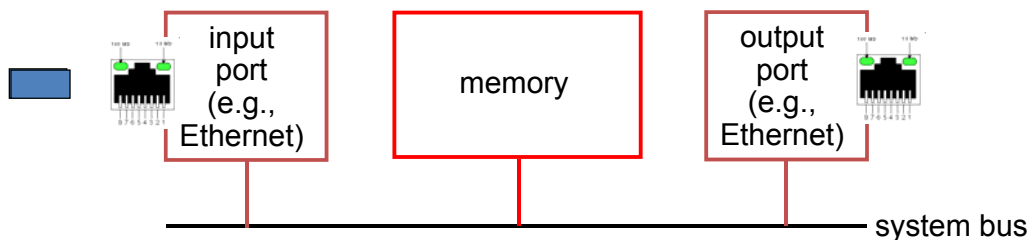


CS3700 - Instructor: Dr. Zhu

## Switching via memory

### *first generation routers:*

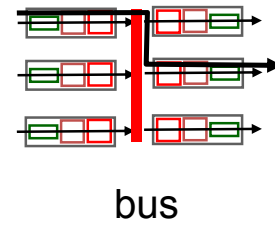
- traditional computers with switching under direct control of CPU
- packet copied to system's memory
- speed limited by memory bandwidth (2 bus crossings per datagram)



CS3700 - Instructor: Dr. Zhu

# Switching via a bus

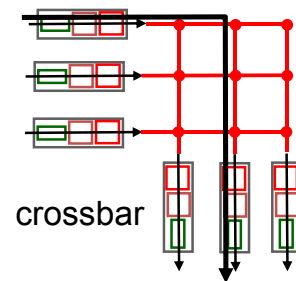
- datagram from input port memory to output port memory via a shared bus
- *bus contention*: switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers



CS3700 - Instructor: Dr. Zhu

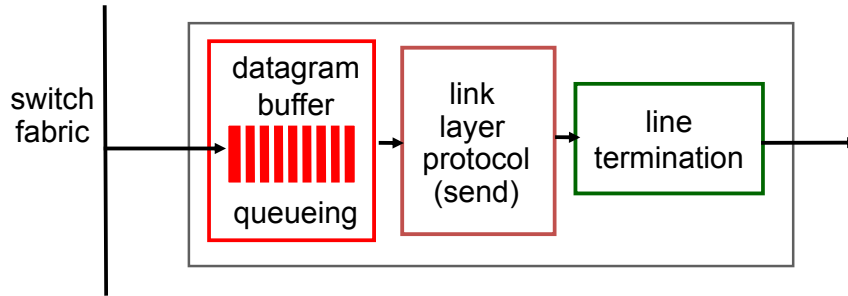
# Switching via interconnection network

- overcome bus bandwidth limitations
- banyan networks, crossbar, 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 60 Gbps through the interconnection network



CS3700 - Instructor: Dr. Zhu

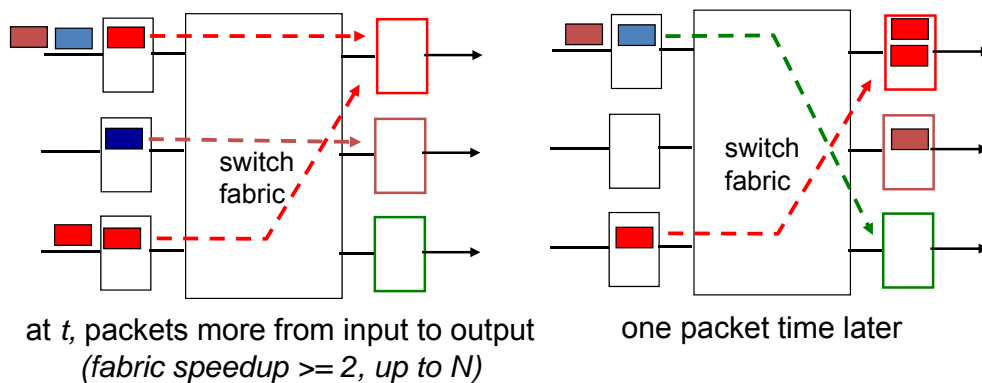
# Output ports



- *buffering* required when datagrams arrive from fabric faster than the transmission rate
- *scheduling discipline* chooses among queued datagrams for transmission

CS3700 - Instructor: Dr. Zhu

## Output port queueing

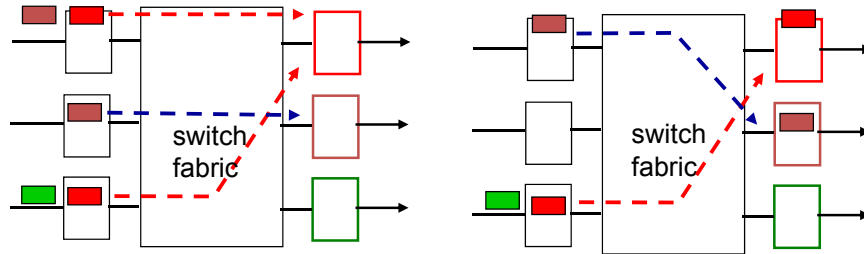


- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

CS3700 - Instructor: Dr. Zhu



# Input port queuing



output port contention: only one red datagram  
can be transferred. (fabric speedup = 1)  
*lower red packet is blocked*

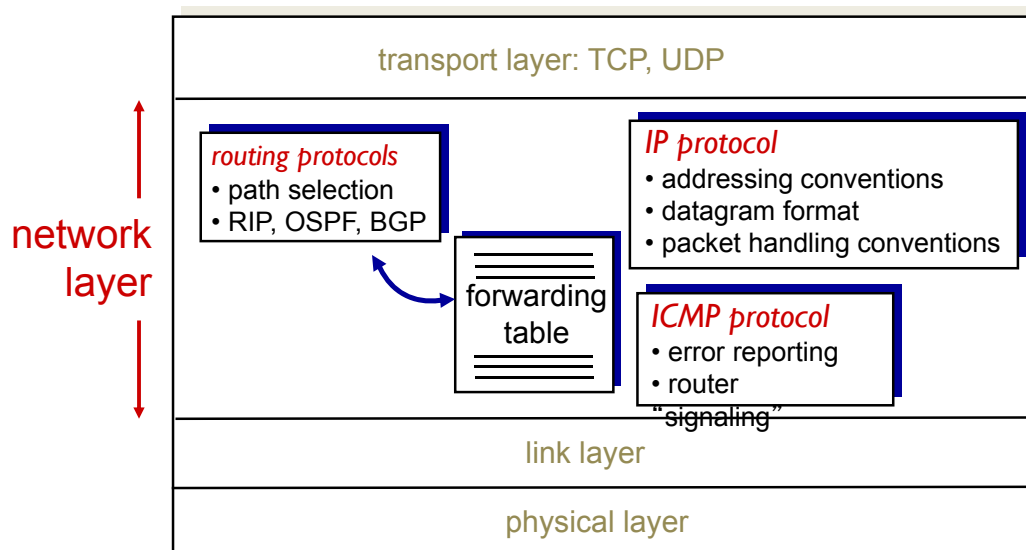
one packet time later: green packet  
experiences HOL blocking

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

CS3700 - Instructor: Dr. Zhu

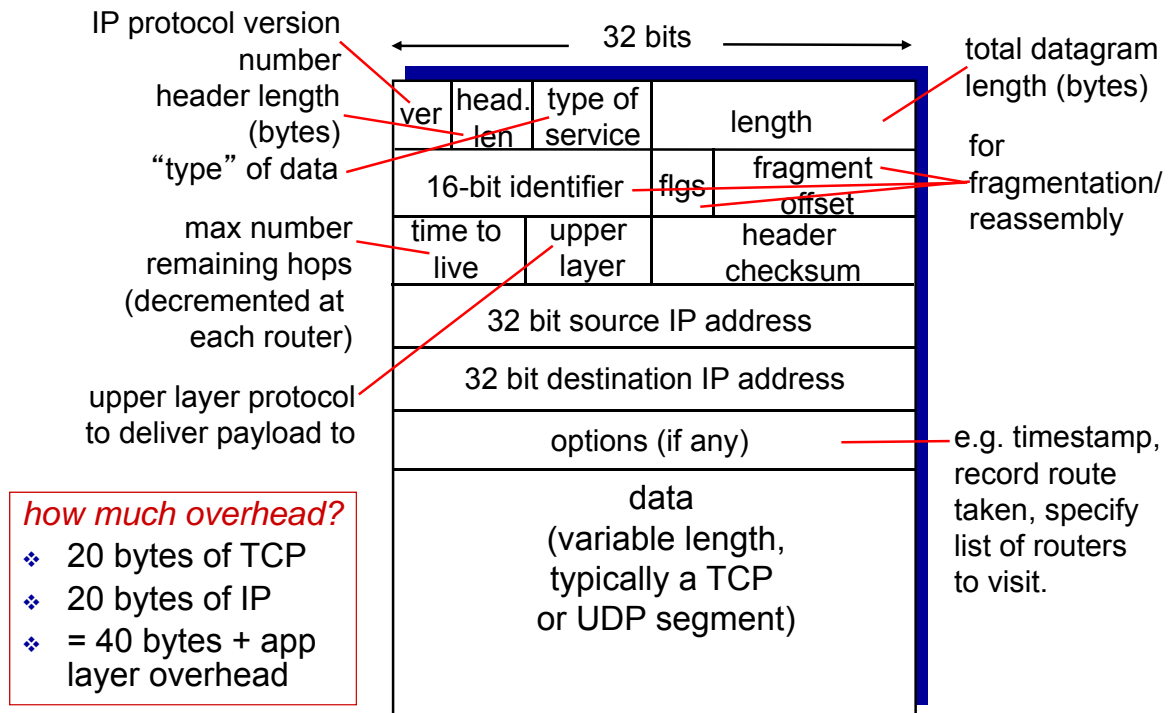
## The Internet network layer

host, router network layer functions:



CS3700 - Instructor: Dr. Zhu

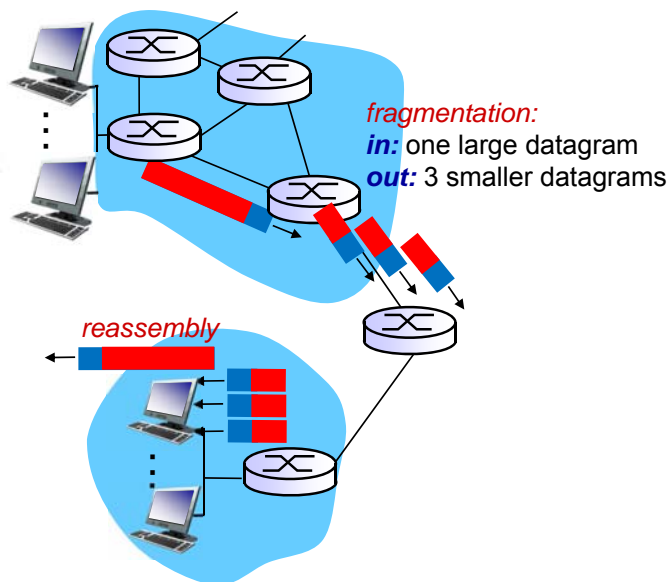
# IP datagram format



CS3700 - Instructor: Dr. Zhu

## 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



CS3700 - Instructor: Dr. Zhu

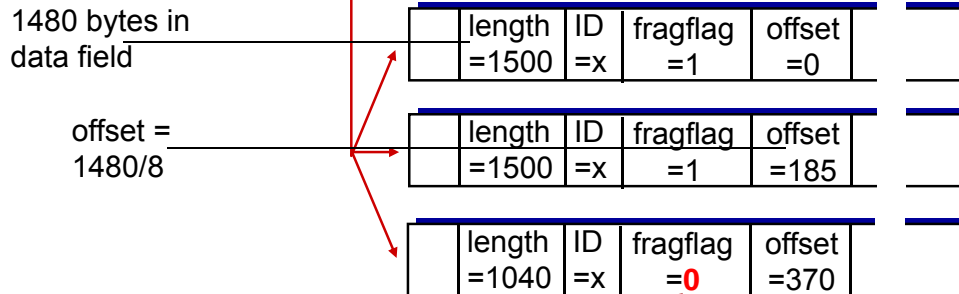
# IP fragmentation, reassembly

## example:

- 4000 byte datagram
- MTU = 1500 bytes

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

*one large datagram becomes several smaller datagrams*

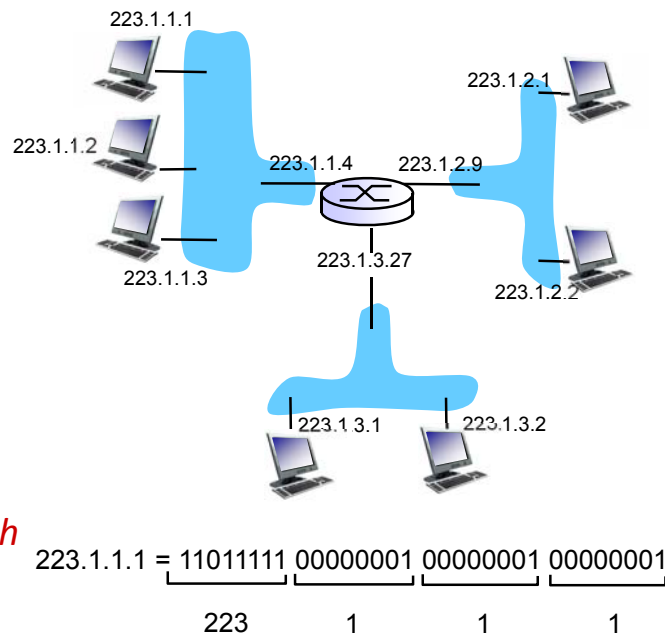


“fragflag = 0” meaning this is the last fragment

CS3700 - Instructor: Dr. Zhu

## 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**



CS3700 - Instructor: Dr. Zhu

# IP addressing: introduction

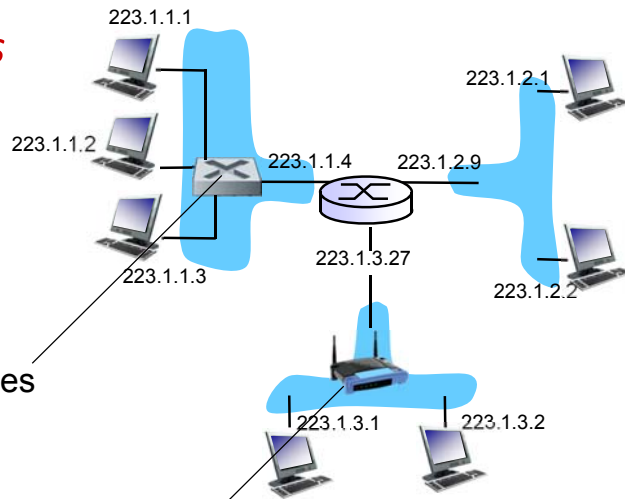
*Q: how are interfaces actually connected?*

*A: we'll learn about that in chapter 5, 6.*

*A: wired Ethernet interfaces connected by Ethernet switches*

*For now:* don't need to worry about how one interface is connected to another (with no intervening router)

*A: wireless WiFi interfaces connected by WiFi base station*



CS3700 - Instructor: Dr. Zhu

# IP addressing: introduction

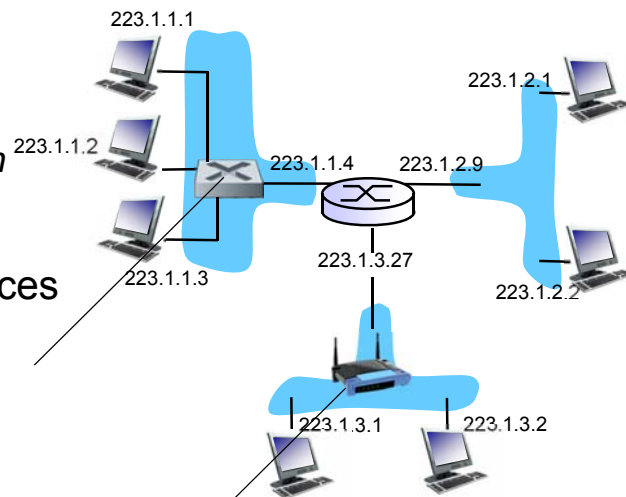
*Q: how are interfaces actually connected?*

*A: we'll learn about that in chapter 5, 6.*

*A: wired Ethernet interfaces connected by Ethernet switches*

*For now:* don't need to worry about how one interface is connected to another (with no intervening router)

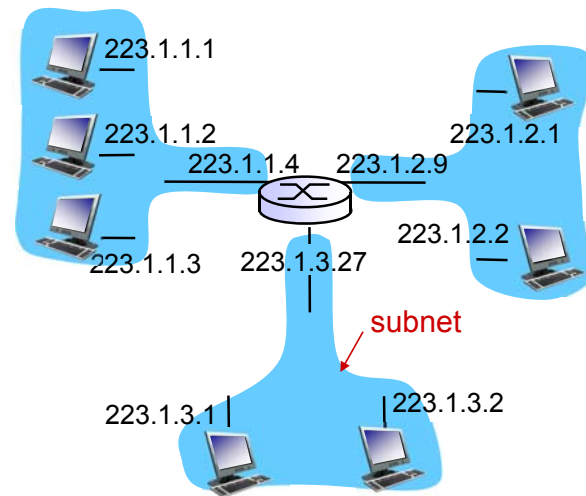
*A: wireless WiFi interfaces connected by WiFi base station*



CS3700 - Instructor: Dr. Zhu

# Subnets

- IP address:
  - subnet part - high order bits
  - host part - low order bits
- *what's a subnet?*
  - device interfaces with same subnet part of IP address
  - can physically reach each other **without intervening router**



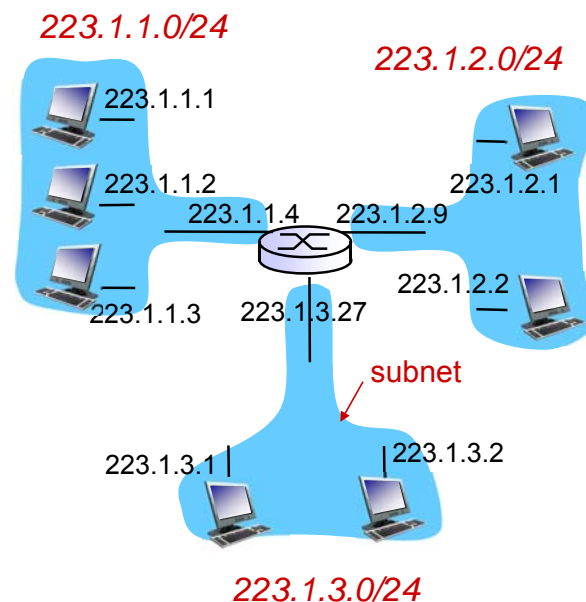
network consisting of 3 subnets

CS3700 - Instructor: Dr. Zhu

# Subnets

## *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*



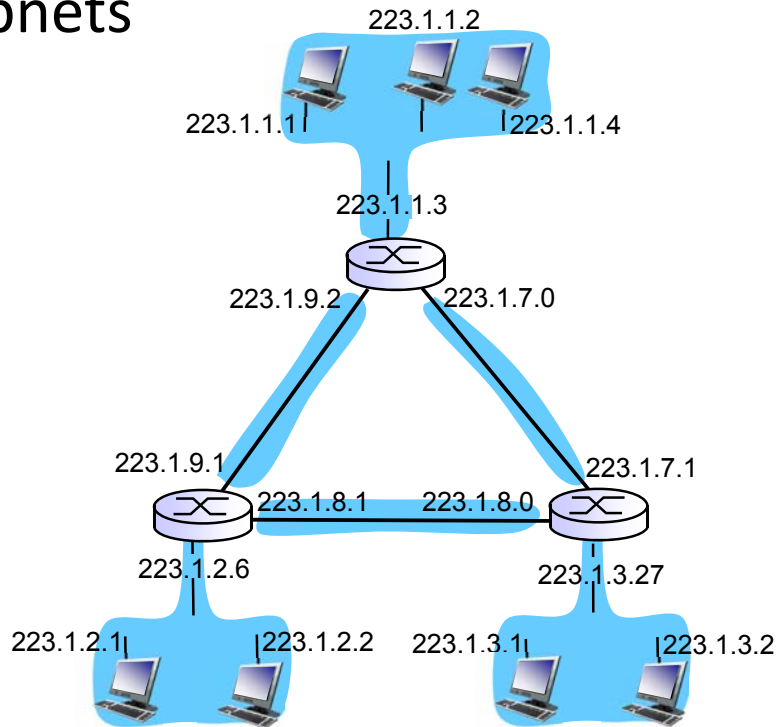
Network Prefix: 223.1.1.0/24

*<first address of a network>/<bit-length of the prefix>*

CS3700 - Instructor: Dr. Zhu

# Subnets

how many?



CS3700 - Instructor: Dr. Zhu

## IP addressing: CIDR

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

← subnet part → host part →  
11001000 00010111 00010000 00000000  
**200.23.16.0/23**

**Q:** How does a *host* get IP address?

- hard-coded by system admin in a file
  - Windows: control-panel->network->configuration->tcp/ip->properties
  - UNIX: /etc/rc.config
- **DHCP: Dynamic Host Configuration Protocol:** dynamically get address from as server
  - “plug-and-play”

CS3700 - Instructor: Dr. Zhu

# IP addresses: how to get one?

**Q:** how does *network* get subnet part of IP addr?

**A:** gets allocated portion of its provider ISP's address space

ISP's block      11001000 00010111 00010000 00000000    200.23.16.0/20

Organization 0    11001000 00010111 00010000 00000000    200.23.16.0/23

Organization 1    11001000 00010111 00010010 00000000    200.23.18.0/23

Organization 2    11001000 00010111 00010100 00000000    200.23.20.0/23

...

....

....

....

Organization 7    11001000 00010111 00011110 00000000    200.23.30.0/23

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

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

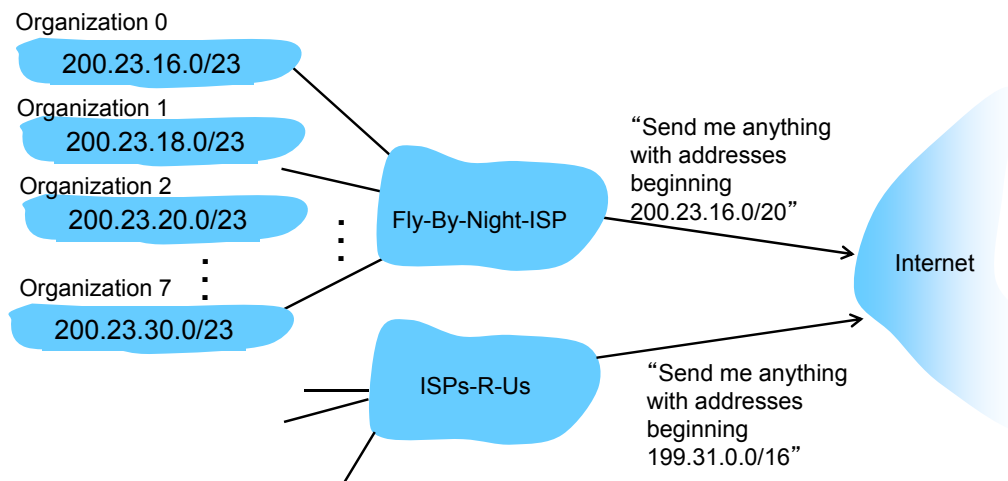
Names and Numbers <http://www.icann.org/>

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

CS3700 - Instructor: Dr. Zhu

## Hierarchical addressing: route aggregation

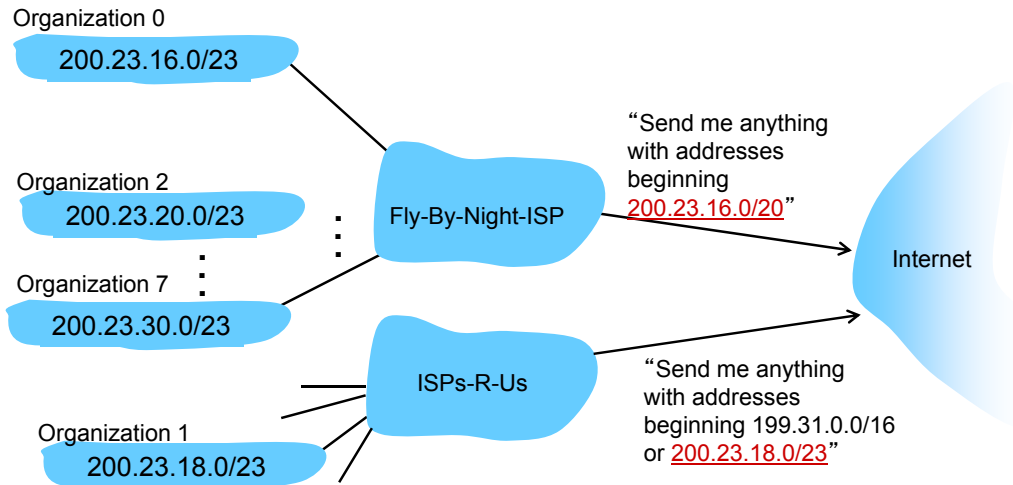
hierarchical addressing allows efficient advertisement of routing information:



CS3700 - Instructor: Dr. Zhu

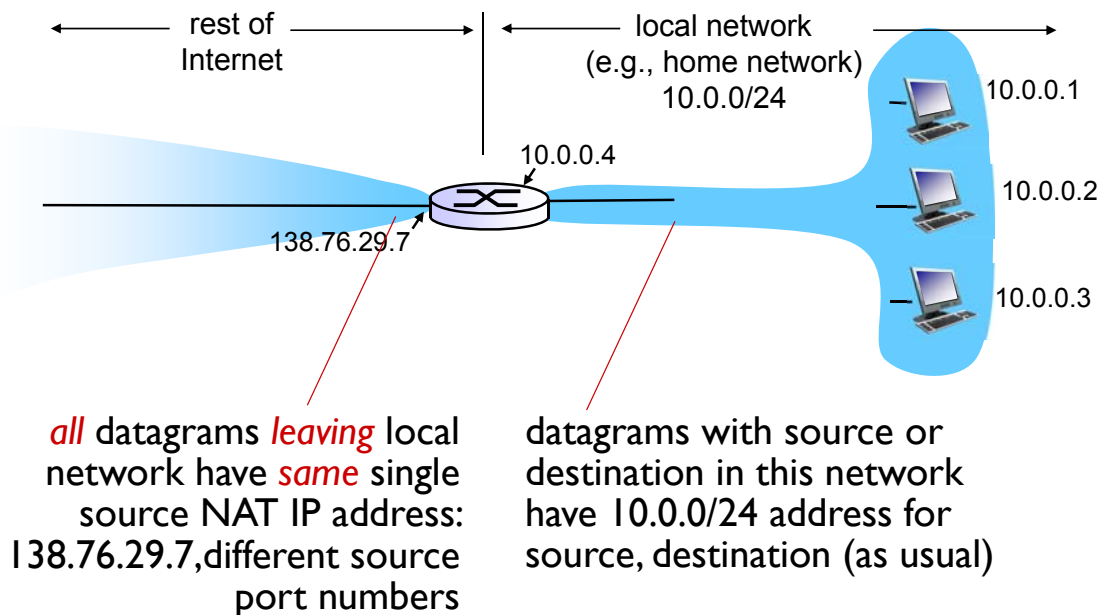
# Hierarchical addressing: more specific routes

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



CS3700 - Instructor: Dr. Zhu

## NAT: network address translation



CS3700 - Instructor: Dr. Zhu



# NAT: network address translation

- *motivation*: local network uses just one IP address as far as outside world is concerned:
  - range of addresses not needed from ISP: just one IP address for all devices
  - can change addresses of devices in local network without notifying outside world
  - can change ISP without changing addresses of devices in local network
  - devices inside local net not explicitly addressable, visible by outside world (a security plus)

CS3700 - Instructor: Dr. Zhu

# NAT: network address translation

- *implementation*: NAT router must:
  - *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)  
... remote clients/servers will respond using (NAT IP address, new port #) as destination addr
  - *remember (in NAT translation table)* every (source IP address, port #) to (NAT IP address, new port #) translation pair
  - *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

CS3700 - Instructor: Dr. Zhu

# NAT: Private Network

- RFC 1918, RFC 4193
- A private network uses private IP Addresses:
  - not globally delegated
  - not allocated to any specific organization
  - IP packets addressed by them cannot be transmitted onto the public Internet
  - anyone may use these addresses without approval from a Regional Internet Registry (RIR).
  - A NAT (Network Address Translator) router or a proxy server must be used to connect a private network to the Internet

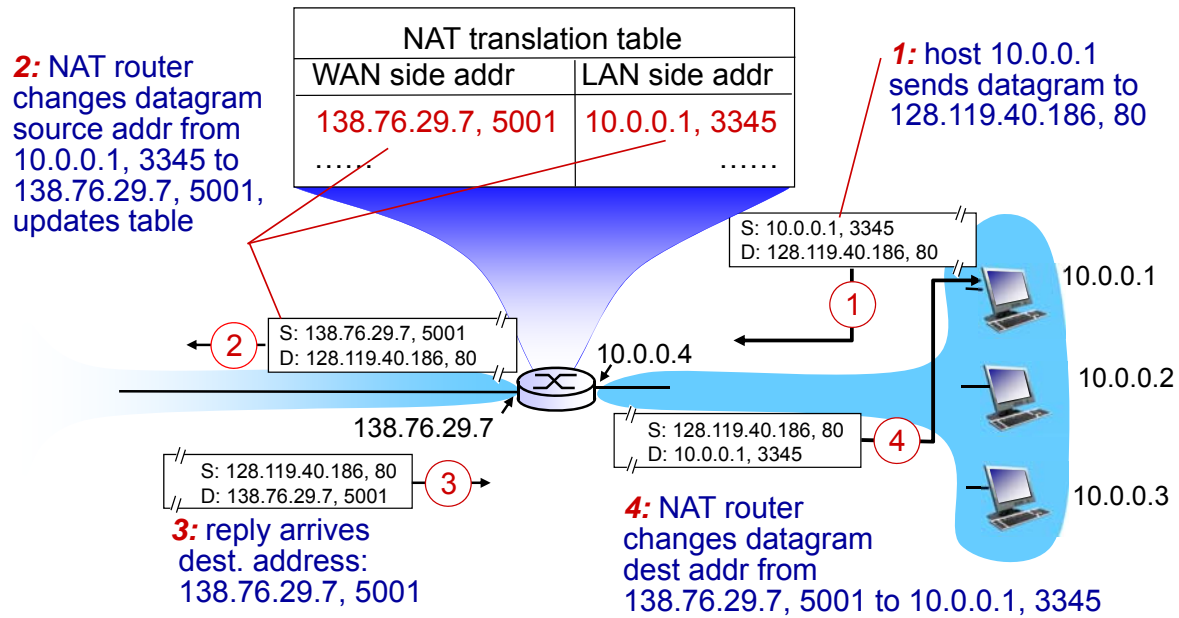
CS3700 - Instructor: Dr. Zhu

## NAT: Private IPv4 Address Spaces

<b>RFC 1918 name</b>	<b>IP address range</b>	<b>number of addresses</b>	<b>Classful description</b>	<b>largest CIDR block (subnet mask)</b>	<b>host id size</b>	<b>mask bits</b>
24-bit block	10.0.0.0 - 10.255.255.255	16,777,216	Single class A network	10.0.0.0/8 (255.0.0.0)	24 bits	8 bits
20-bit block	172.16.0.0 - 172.31.255.255	1,048,576	16 contiguous class B network	172.16.0.0/12 (255.240.0.0)	20 bits	12 bits
16-bit block	192.168.0.0 - 192.168.255.255	65,536	256 contiguous class C network	192.168.0.0/16 (255.255.0.0)	16 bits	16 bits

CS3700 - Instructor: Dr. Zhu

# NAT: network address translation



CS3700 - Instructor: Dr. Zhu

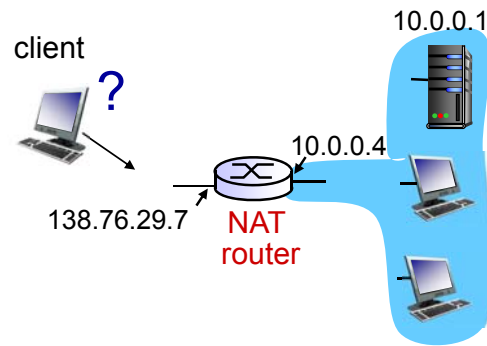
# NAT: network address translation

- 16-bit port-number field:
  - 60,000 simultaneous connections with a single LAN-side address!
- NAT is controversial:
  - NAT traversal problem
  - routers should only process up to layer 3
  - violates end-to-end argument
    - NAT possibility must be taken into account by app designers, e.g., P2P applications
  - address shortage should instead be solved by IPv6

CS3700 - Instructor: Dr. Zhu

# NAT traversal problem

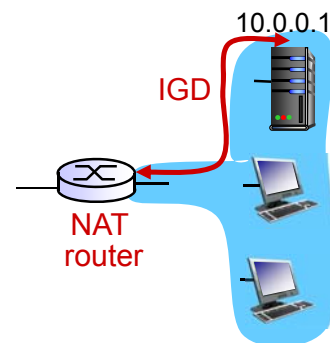
- client wants to connect to server with address 10.0.0.1
  - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - only one externally visible NATed address: 138.76.29.7
- **solution1:** statically configure NAT to forward incoming connection requests at given port to server
  - e.g., (123.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000



CS3700 - Instructor: Dr. Zhu

# NAT traversal problem

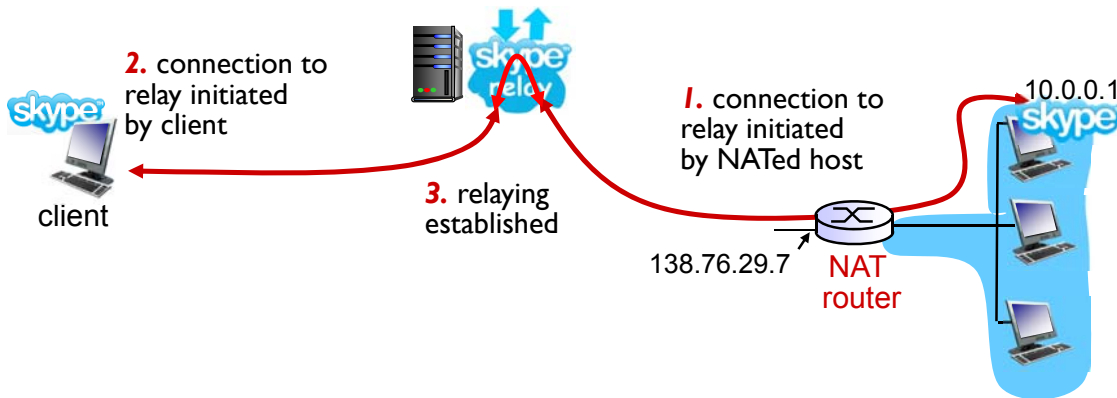
- **solution 2:** Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATed host to:
    - learn public IP address (138.76.29.7)
    - add/remove port mappings (with lease times)
- i.e., automate static NAT port map configuration



CS3700 - Instructor: Dr. Zhu

# NAT traversal problem

- **solution 3:** relaying (used in Skype)
  - NATed client establishes connection to relay
  - external client connects to relay
  - relay bridges packets between to connections



CS3700 - Instructor: Dr. Zhu

## ICMP: internet control message protocol

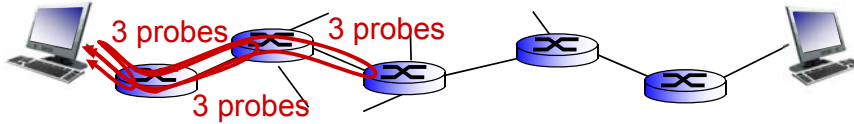
- used by hosts & routers to communicate network-level information
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- network-layer “above” IP:
  - ICMP msgs carried in IP datagrams
- **ICMP message:** type, code plus first 8 bytes of IP datagram causing 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

CS3700 - Instructor: Dr. Zhu

# Traceroute and ICMP

- source sends series of UDP segments to dest
    - first set has TTL=1
    - second set has TTL=2, etc.
    - unlikely port number
  - when  $n$ th set of datagrams arrives to  $n$ th router:
    - router discards datagrams
    - and sends source ICMP messages (type 11, code 0)
    - ICMP messages includes name of router & IP address
  - when ICMP messages arrives, source records RTTs
- stopping criteria:*
- UDP segment eventually arrives at destination host
  - destination returns ICMP “port unreachable” message (type 3, code 3)
  - source stops



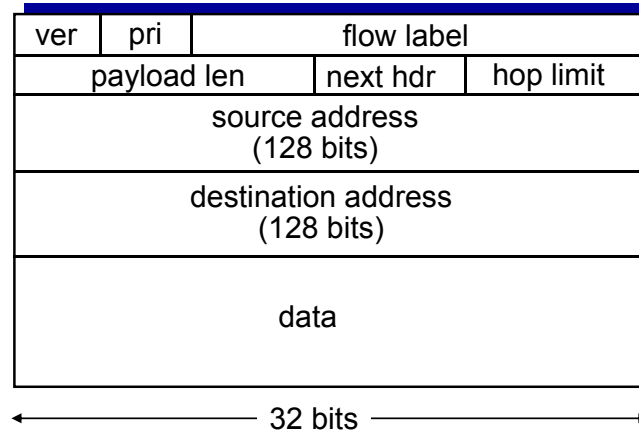
CS3700 - Instructor: Dr. Zhu

# 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*:
  - *options*: allowed, but outside of header, indicated by “Next Header” field
  - E.g., identify upper layer protocol for data



CS3700 - Instructor: Dr. Zhu

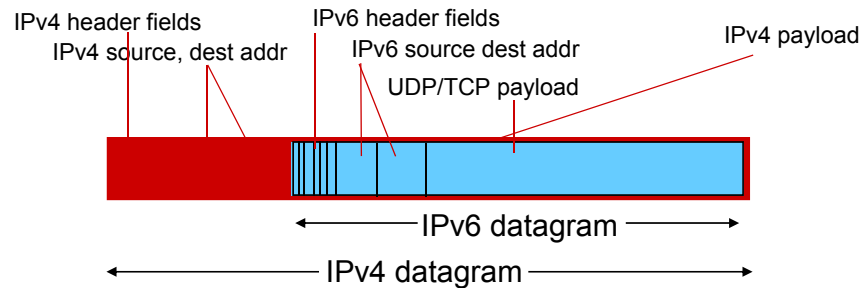
## Other changes from IPv4

- *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

CS3700 - Instructor: Dr. Zhu

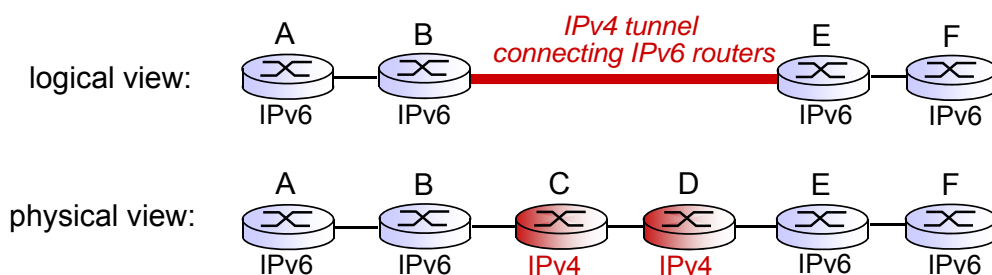
# Transition from IPv4 to IPv6

- 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



CS3700 - Instructor: Dr. Zhu

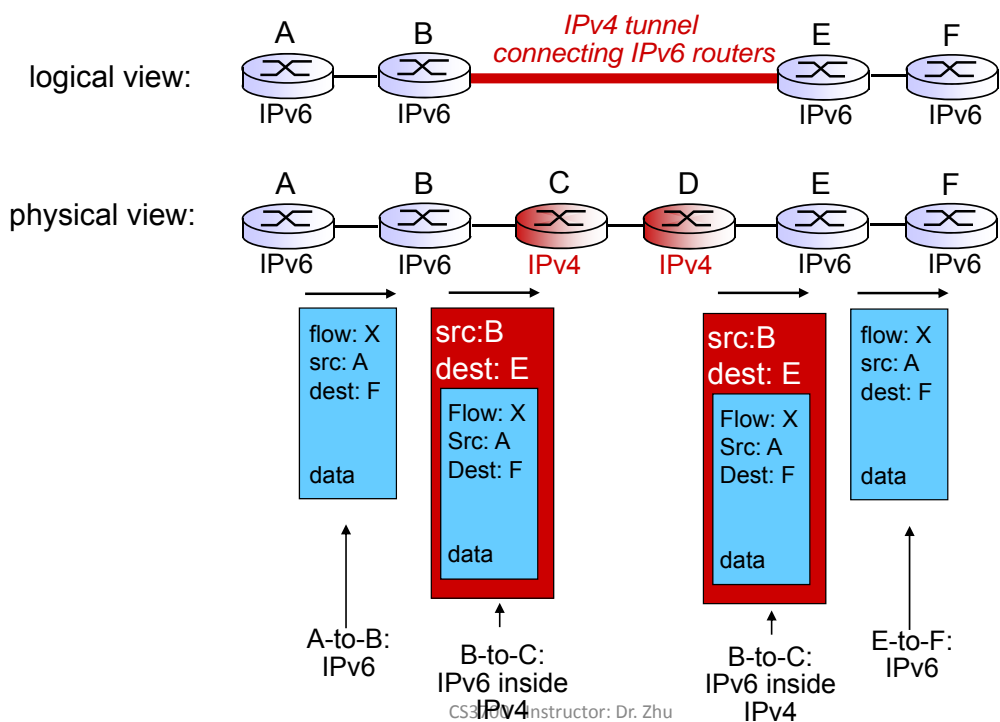
## Tunneling



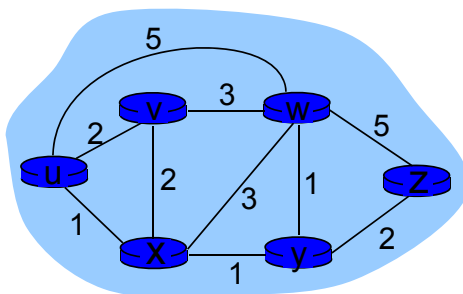
CS3700 - Instructor: Dr. Zhu



# Tunneling



## Graph abstraction



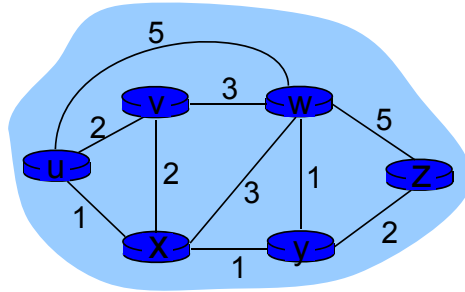
graph:  $G = (N, E)$

$N$  = set of routers =  $\{ u, v, w, x, y, z \}$

$E$  = set of links =  $\{ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) \}$

aside: graph abstraction is useful in other network contexts, e.g., P2P, where  $N$  is set of peers and  $E$  is set of TCP connections

# Graph abstraction: costs



$c(x, x')$  = cost of link  $(x, x')$   
e.g.,  $c(w, z) = 5$

cost could always be 1, or  
inversely related to bandwidth,  
or inversely related to  
congestion

cost of path  $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

*key question:* what is the least-cost path between u and z ?  
*routing algorithm:* algorithm that finds that least cost path

CS3700 - Instructor: Dr. Zhu

## Routing algorithm classification

*Q: 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

*Q: static or dynamic?*

*static:*

- routes change slowly over time

*dynamic:*

- routes change more quickly
  - periodic update
  - in response to link cost changes

CS3700 - Instructor: Dr. Zhu

# A Link-State Routing Algorithm

## *Dijkstra's algorithm*

- net topology, link costs known to all nodes
  - accomplished via “link state broadcast”
  - all nodes have same info
- computes least cost paths from one node (“source”) to all other nodes
  - i.e., shortest-path tree
  - gives *forwarding table* for that node
- iterative: after  $k$  iterations, know least cost path to  $k$  dest.'s

## *notation:*

- $c(x,y)$ : link cost from node  $x$  to  $y$ ;  $= \infty$  if not direct neighbors
- $D(i)$ : current value of cost of path from source to dest.  $i$
- $p(i)$ : predecessor node along path from source to  $i$
- $(x,y)$ : edge (i.e., link) from node  $x$  to  $y$
- $N'$ : set of nodes whose least cost path definitively known
- $Y'$ : set of edges currently known to be in shortest-path tree rooted at “source” node

## *Iterative update:*

- add a node from  $N - N'$ , that has a shortest path from source node, using only nodes in  $N'$  as intermediates

CS3700 - Instructor: Dr. Zhu

## Dijkstra's Algorithm

### 1 *Initialization:*

```
2  $N' = \{u\}; \quad Y' = \emptyset$  (empty set)
3 for all nodes  $i$ 
4   if  $i$  adjacent to  $u$ 
5     then  $D(i) = c(u,i), p(i) = u$ 
6   else  $D(i) = \infty$ 
```

7

### 8 *Loop*

```
9 find  $k$  not in  $N'$  such that  $D(k)$  is a minimum
10 add node  $k$  to  $N'$ 
11 add edge  $(p(k), k)$  to  $Y'$ 
12 update  $D(i)$  and  $p(i)$  for all  $i$  adjacent to  $k$  and not in  $N'$  :
13   if  $D(k) + c(k,i) < D(i)$ 
14     then  $D(i) = D(k) + c(k,i)$  and  $p(i) = k$ 
15 /* new cost to  $i$  is either old cost to  $i$  or known
16    shortest path cost to  $k$  plus cost from  $k$  to  $i$  */
```

```
17 until all nodes in  $N'$ 
```

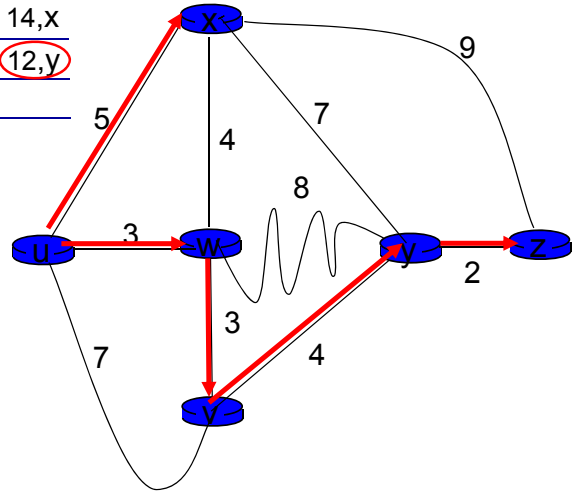
CS3700 - Instructor: Dr. Zhu

# Dijkstra's algorithm: example

Step	N'	D(v) p(v)	D(w) p(w)	D(x) p(x)	D(y) p(y)	D(z) p(z)
0	u	7,u	3,u	5,u	$\infty$	$\infty$
1	uw	6,w		5,u	11,w	$\infty$
2	uwx	6,w			11,w	14,x
3	uwxv				10,v	14,x
4	uwxvy				12,y	
5	uwxvyz					

## notes:

- Q: How Y' is updated step-by-step?
- construct shortest path tree by tracing predecessor nodes
- ties can exist (can be broken arbitrarily)

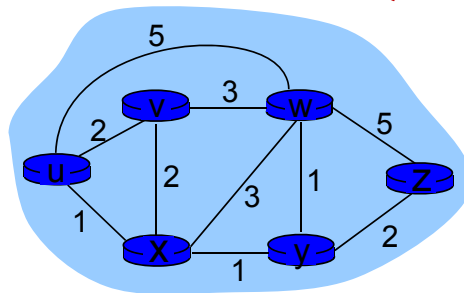


CS3700 - Instructor: Dr. Zhu

# Dijkstra's algorithm: another example

Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	$\infty$	$\infty$
1	ux	2,u	4,x		2,x	$\infty$
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					

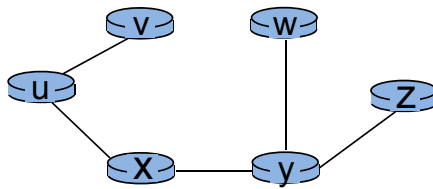
Q: How Y' is updated step-by-step?



CS3700 - Instructor: Dr. Zhu

# Dijkstra's algorithm: another example (cont.)

resulting shortest-path tree from u:



Q: Which result of Dijkstra's algorithm represents such shortest-path tree?

Y'

Q: What is the algorithm of creating such forwarding table in the source node u using the results of Dijkstra's algorithm?

resulting forwarding table in u:

destination	link
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

for each node i in N' but i != u

j = i

while (p(j) != u)

j = p(j)

end while

add entry dest = i, link = (u, j) into fwd table

end for

CS3700 - Instructor: Dr. Zhu

## Distance vector algorithm

*Bellman-Ford equation (dynamic programming)*

let

$d_x(y) := \text{cost of least-cost path from } x \text{ to } y$

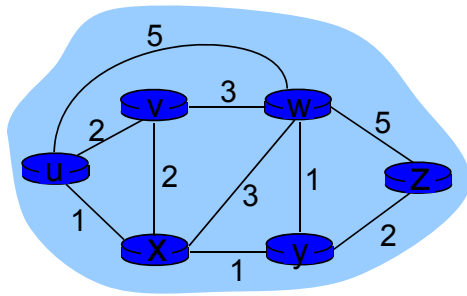
then

$$d_x(y) = \min \{ c(x,i) + d_i(y) \}$$

$\min$  taken over all neighbors i of x  
 $c(x,i)$  cost to neighbor i  
 $d_i(y)$  cost from neighbor i to destination y

CS3700 - Instructor: Dr. Zhu

# Bellman-Ford example



clearly,  $d_v(z) = 5$ ,  $d_x(z) = 3$ ,  $d_w(z) = 3$

B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z), \\ c(u,x) + d_x(z), \\ c(u,w) + d_w(z) \} \\ = \min \{ 2 + 5, \\ 1 + 3, \\ 5 + 3 \} = 4$$

node achieving minimum, **node x**, is next hop in shortest path, used in the forwarding table at **node u** to the destination **node z**. i.e., at **node u**,

destination	link
...	
z	(u, x)
...	

CS3700 - Instructor: Dr. Zhu

## Distance vector algorithm

- Node x maintains distance vector  $\mathbf{D}_x = [D_x(j): j \in N]$  and link vector  $\mathbf{L}_x = [L_x(j): j \in N]$  (forwarding table)
  - $D_x(j)$  = estimate of least cost from x to j
  - $L_x(j) = (x, n_x(j))$ , node  $n_x(j)$  is the neighboring node achieving least cost, i.e., next hop in the path with least cost from node x to node j
  - $N$ : set of nodes in the network (maybe including node x itself)
- Meanwhile, node x also:
  - knows cost to each neighbor i:  $c(x, i)$
  - maintains its neighbors' distance vectors. For each neighbor i, node x maintains  $\mathbf{D}_i = [D_i(j): j \in N]$

CS3700 - Instructor: Dr. Zhu

# Distance vector algorithm

## *key idea:*

- from time-to-time, each node sends its own distance vector estimate to neighbors
- when  $x$  receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$D_x(j) \leftarrow \min_i \{c(x, i) + D_i(j)\}$  for each node  $j \in N$ ,  
where  $\min$  is taken over all neighbors  $i$  of node  $x$

$L_x(j) \leftarrow (x, k)$  for each node  $j \in N$ ,  
where node  $k$  is the neighbor achieving the  $\min$

- under minor, natural conditions, the estimate  $D_x(y)$  converge to the actual least cost  $d_x(y)$

CS3700 - Instructor: Dr. Zhu

# Distance vector algorithm

## *iterative, asynchronous:*

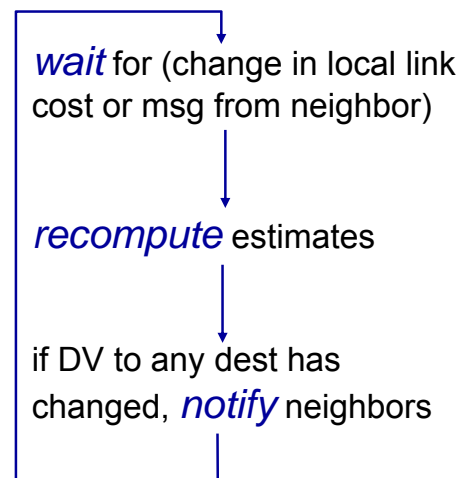
each local iteration caused by:

- local link cost change
- DV update message from neighbor

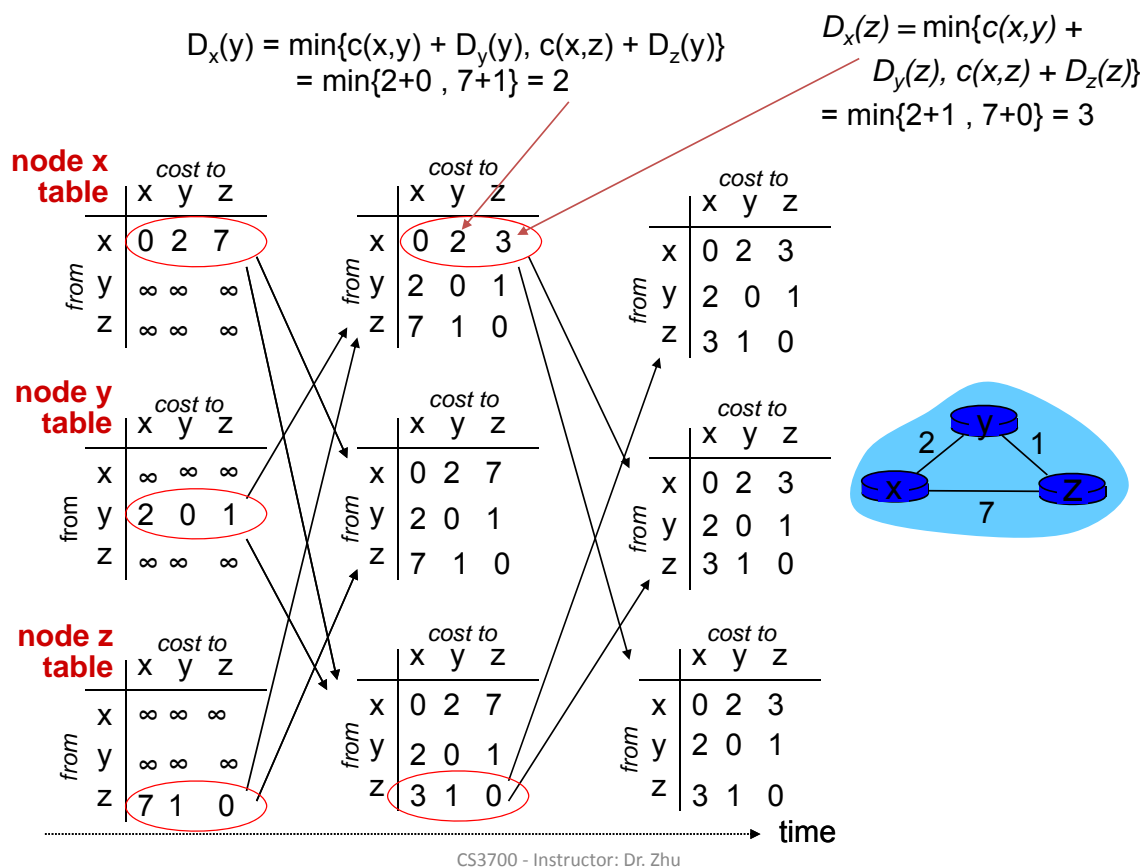
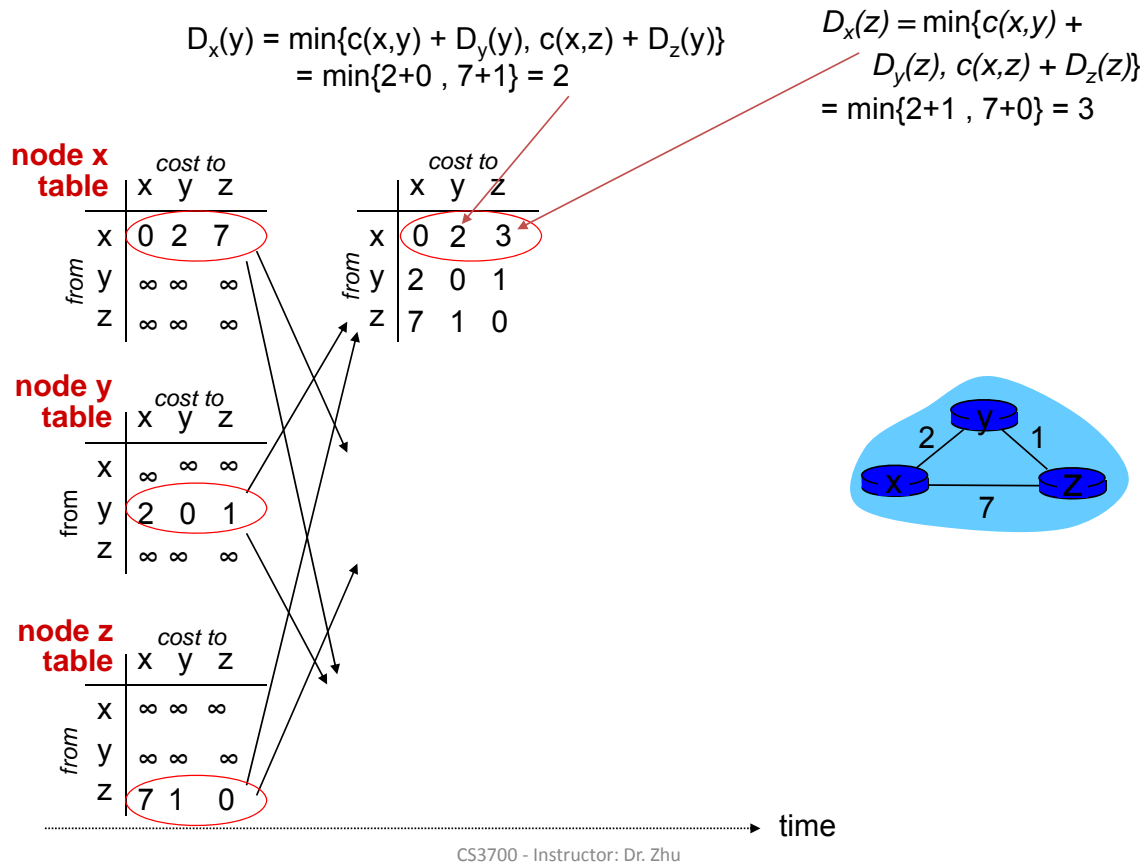
## *distributed:*

- each node notifies neighbors *only* when its DV changes
  - neighbors then notify their neighbors if necessary

## *each node:*



CS3700 - Instructor: Dr. Zhu

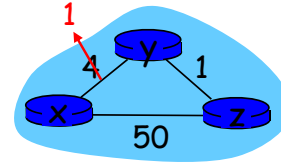




# Distance vector: link cost changes

## link cost changes:

- node detects local link cost change
- updates routing info, recalculates distance vector
- if DV changes, notify neighbors



“good news travels fast”

$t_0$ : y detects link-cost change, updates its DV, informs its neighbors.

$t_1$ : z receives update from y, updates its table, computes new least cost to x, sends its neighbors its DV.

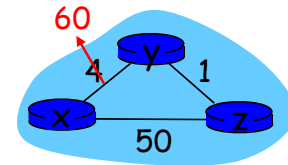
$t_2$ : y receives z's update, updates its distance table. y's least costs do not change, so y does not send a message to z.

CS3700 - Instructor: Dr. Zhu

# Distance vector: link cost changes

## link cost changes:

- node detects local link cost change
  - $t_0$ : y detects the change
  - $t_1$ : y updates  $D_y(x)$  as 6 via (y, z) because y now has  $D_z(x) = 5$  and  $c(y, z) = 1$  now
  - $t_2$ : y notifies z and z updates  $D_z(x)$  as 7 via (z, y) because z now has  $D_y(x) = 6$  and  $c(z, y) = 1$
  - $t_3$ : z notifies y and y updates  $D_y(x)$  as 8 via (y, z) because y now has  $D_z(x) = 7$  and  $c(y, z) = 1$
  - $t_4$ : y notifies z and z updates  $D_z(x)$  as 9 via (z, y) because z now has  $D_y(x) = 8$  and  $c(z, y) = 1$
- **bad news travels slow** - “count to infinity” problem!
- 44 iterations before algorithm stabilizes
- A routing loop exists before algorithm stabilizes
  - At node y, a datagram to x is bounced between y and z



CS3700 - Instructor: Dr. Zhu

# Hierarchical routing

*scale:* with 600 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:*

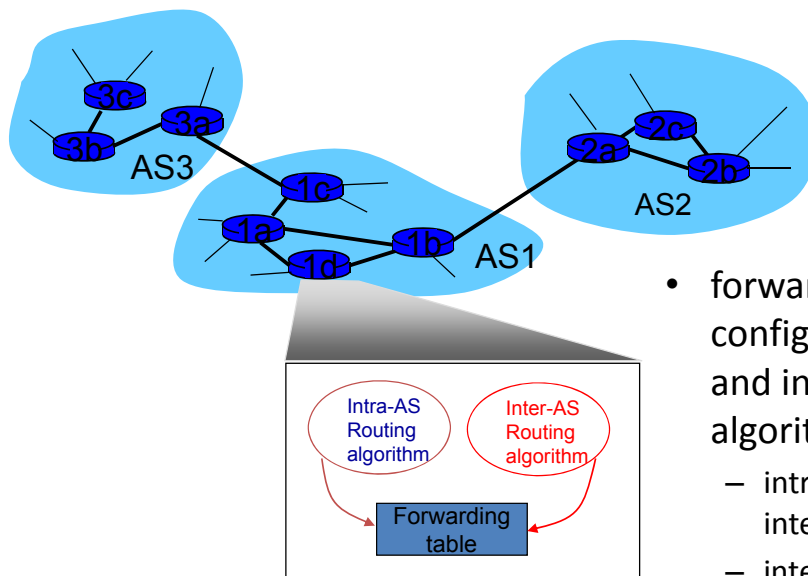
- aggregate routers into regions, “autonomous systems” (AS)
- routers in same AS run same routing protocol
  - “intra-AS” routing protocol
  - routers in different AS can run different intra-AS routing protocol

*gateway router:*

- at “edge” of its own AS
- has link to router in another AS

CS3700 - Instructor: Dr. Zhu

## Interconnected ASes



- forwarding table configured by both intra- and inter-AS routing algorithm
  - intra-AS sets entries for internal dests
  - inter-AS & intra-AS sets entries for external dests

CS3700 - Instructor: Dr. Zhu

# Inter-AS tasks

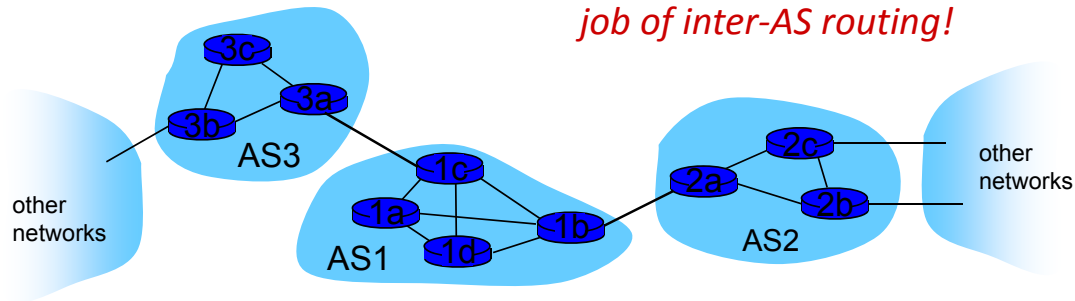
- suppose router in AS1 receives datagram destined outside of AS1:

- router should forward packet to gateway router, but which one?

*AS1 must:*

- learn which dests are reachable through AS2, which through AS3
- propagate this reachability info to all routers in AS1

*job of inter-AS routing!*



CS3700 - Instructor: Dr. Zhu

## Intra-AS Routing

- also known as *interior gateway protocols (IGP)*
- most common intra-AS routing protocols:
  - RIP: Routing Information Protocol
  - **OSPF: Open Shortest Path First**
  - IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

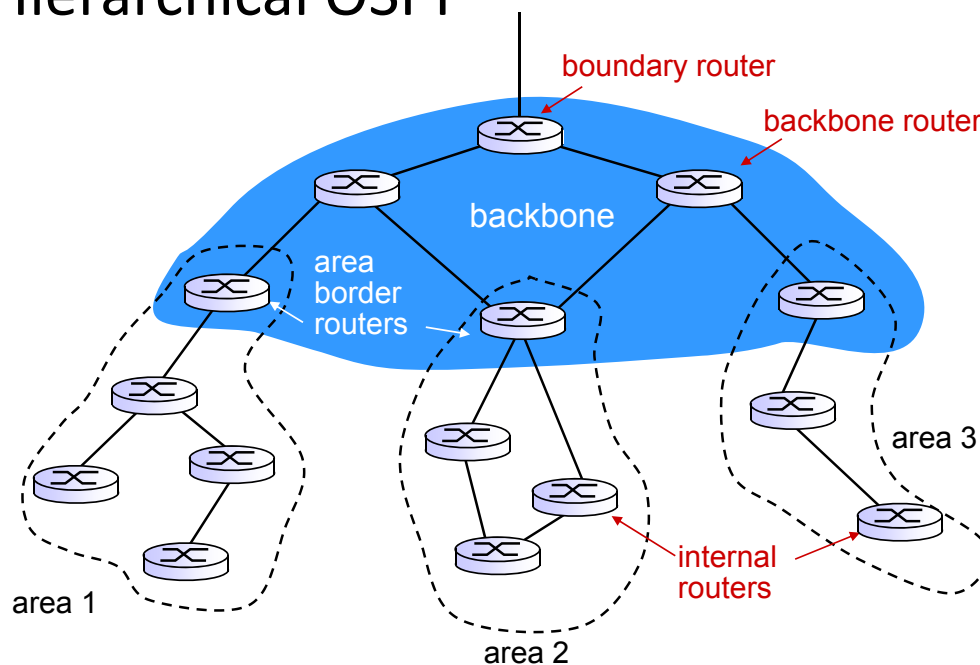
CS3700 - Instructor: Dr. Zhu

# 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
- advertisements flooded to *entire* AS
  - carried in OSPF messages directly over IP (rather than TCP or UDP)
- *IS-IS routing* (Intermediate System to Intermediate System) protocol: nearly identical to OSPF

CS3700 - Instructor: Dr. Zhu

## Hierarchical OSPF



CS3700 - Instructor: Dr. Zhu

# 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 AS’ s.

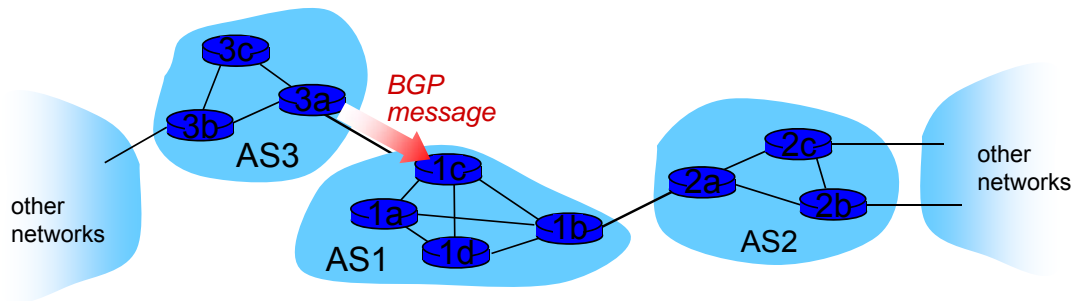
CS3700 - Instructor: Dr. Zhu

## 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”*

CS3700 - Instructor: Dr. Zhu

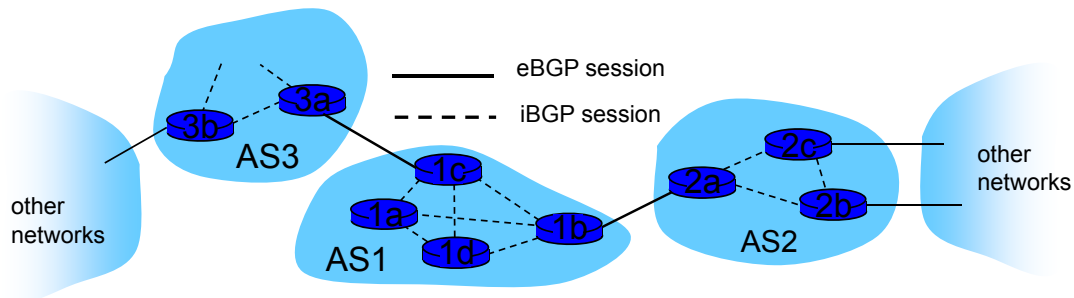
# BGP basics



- **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

CS3700 - Instructor: Dr. Zhu

## BGP basics: distributing path information



- using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
  - 1c can then use iBGP to 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.

CS3700 - Instructor: Dr. Zhu

# Path attributes and BGP routes

- advertised prefix includes BGP attributes
  - prefix + attributes = “route”
- two important attributes:
  - **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g., AS 67, AS 17
  - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- gateway router receiving route advertisement uses **import policy** to accept/decline
  - e.g., never route through AS x
  - **policy-based** routing

CS3700 - Instructor: Dr. Zhu

## 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

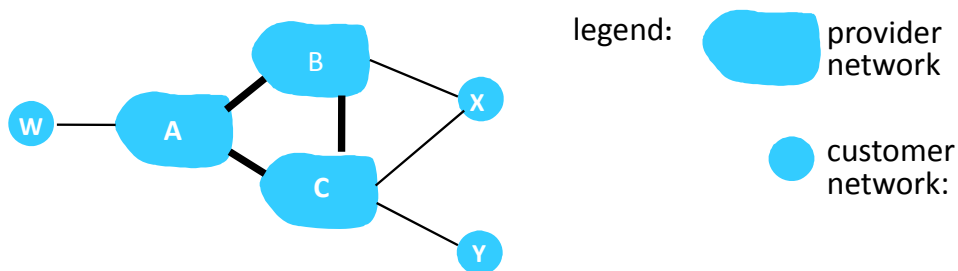
CS3700 - Instructor: Dr. Zhu

# BGP messages

- BGP messages exchanged between peers over TCP connection
- 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

CS3700 - Instructor: Dr. Zhu

## BGP routing policy

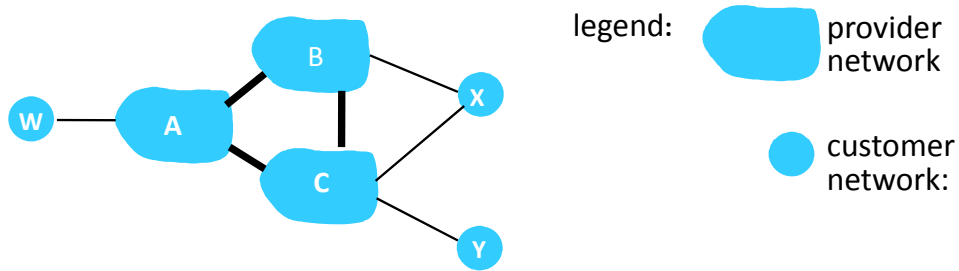


- A,B,C are *provider networks*
- X,W,Y are customer (of provider networks)
- X is *dual-homed*: attached to two networks
  - X does not want to route from B via X to C
  - .. so X will not advertise to B a route to C

CS3700 - Instructor: Dr. Zhu



# BGP routing policy



- A advertises path **AW** to B
- B advertises path **BAW** to X
- Should B advertise path **BAW** to C?
  - No way! B gets no “revenue” for routing CBAW since neither W nor C are B’s customers
  - B wants to force C to route to w via A
  - B wants to route only to/from its customers!