

# CSC358 Intro. to Computer Networks

## Lecture 8: Network Layer

Amir H. Chinaei, Winter 2016

ahchinaei@cs.toronto.edu  
http://www.cs.toronto.edu/~ahchinaei/

Many slides are (inspired/adapted) from the above source  
© all material copyright; all rights reserved for the authors



Office Hours: T 17:00–18:00 R 9:00–10:00 BA4222

TA Office Hours: W 16:00-17:00 BA3201 R 10:00-11:00 BA7172  
csc358ta@cdf.toronto.edu  
http://www.cs.toronto.edu/~ahchinaei/teaching/2016jan/csc358/

# Chapter 4: network layer

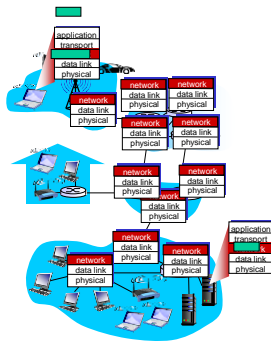
## chapter goals:

- ❖ understand principles behind network layer services:
  - network layer service models
  - forwarding versus routing
  - how a router works
  - routing (path selection)
  - broadcast, multicast
- ❖ instantiation, implementation in the Internet

Network Layer 4-2

## 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



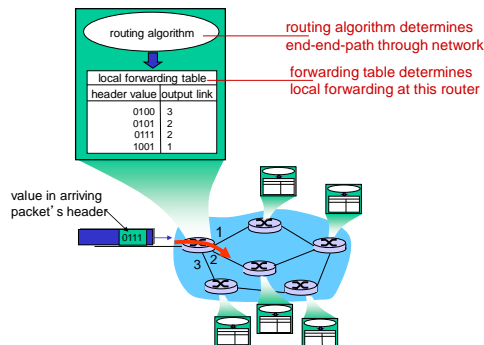
Network Layer 4-3

## Two key network-layer functions

- ❖ **forwarding**: move packets from router's input to appropriate router output
  - ❖ **routing**: determine route taken by packets from source to dest.
    - routing algorithms
- analogy:*
- ❖ **routing**: process of planning trip from source to dest
  - ❖ **forwarding**: process of getting through single interchange

Network Layer 4-4

## Interplay between routing and forwarding



Network Layer 4-5

## Connection setup

- ❖ 3<sup>rd</sup> important function in some network architectures:
  - ATM, frame relay, X.25
- ❖ before datagrams flow, two end hosts and intervening routers establish virtual connection
  - routers get involved
- ❖ network vs transport layer connection service:
  - **network**: between two hosts (may also involve intervening routers in case of VCs)
  - **transport**: between two processes

Network Layer 4-6

## Network service model

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

*example services for individual datagrams:*

- ❖ guaranteed delivery
- ❖ guaranteed delivery with less than 40 msec delay

*example services for a flow of datagrams:*

- ❖ in-order datagram delivery
- ❖ guaranteed minimum bandwidth to flow
- ❖ restrictions on changes in inter-packet spacing

Network Layer 4-7

## Network layer service models:

Network Architecture	Service Model	Guarantees ?				Congestion feedback
		Bandwidth	Loss	Order	Timing	
Internet	best effort	none	no	no	no	no (inferred via loss)
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

Network Layer 4-8

## Chapter 4: outline

4.1 introduction

4.2 *virtual circuit and datagram networks*

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Network Layer 4-9

## Connection, connection-less service

❖ *datagram* network provides network-layer *connectionless* service

❖ *virtual-circuit* network provides network-layer *connection* service

❖ 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

Network Layer 4-10

## 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 address)
- ❖ every router on source-dest path maintains “state” for each passing connection
- ❖ link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

Network Layer 4-11

## VC implementation

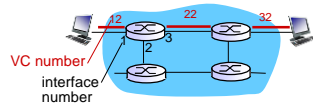
*a VC consists of:*

1. *path* from source to destination
2. *VC numbers*, one number for each link along path
3. *entries in forwarding tables* in routers along path

- ❖ packet belonging to VC carries VC number (rather than dest address)
- ❖ VC number can be changed on each link.
  - new VC number comes from forwarding table

Network Layer 4-12

## VC forwarding table



forwarding table in northwest router:

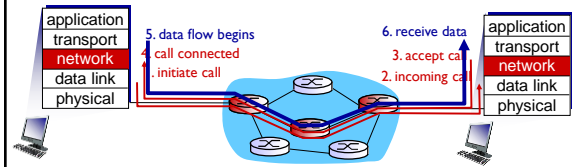
Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...	...	...	...

VC routers maintain connection state information!

Network Layer 4-13

## Virtual circuits: signaling protocols

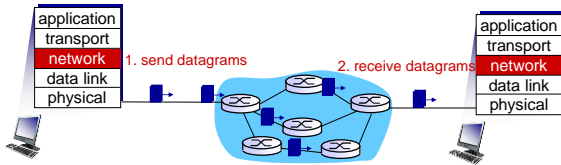
- used to setup, maintain, teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet



Network Layer 4-14

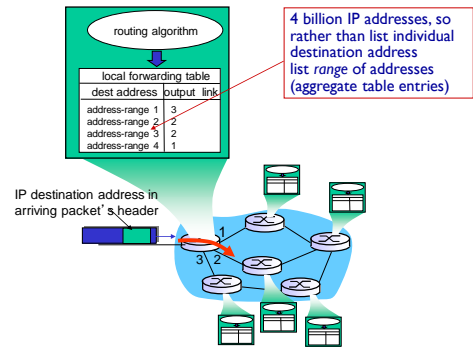
## Datagram networks

- 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



Network Layer 4-15

## Datagram forwarding table



Network Layer 4-16

## 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?

Network Layer 4-17

## Longest prefix matching

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

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

DA: 11001000 00010111 0001**0110 10100001** which interface?  
 DA: 11001000 00010111 0001**1000 10101010** which interface?

Network Layer 4-18

## Datagram or VC network: why?

### Internet (datagram)

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

### ATM (VC)

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

Network Layer 4-19

## Chapter 4: outline

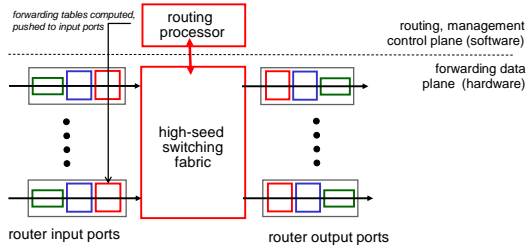
- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what’s inside a router
- 4.4 IP: Internet Protocol
  - datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 routing algorithms
  - link state
  - distance vector
  - hierarchical routing
- 4.6 routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 broadcast and multicast routing

Network Layer 4-20

## Router architecture overview

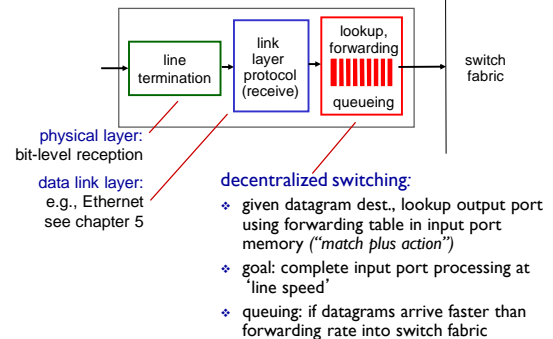
two key router functions:

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



Network Layer 4-21

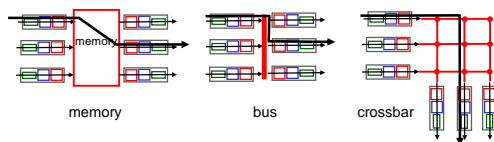
## Input port functions



Network Layer 4-22

## 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

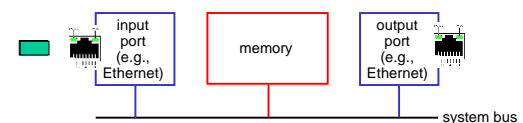


Network Layer 4-23

## Switching via memory

### first generation routers:

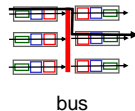
- ❖ traditional computers with switching under direct control of CPU
- ❖ packet copied to system’s memory
- ❖ speed limited by memory bandwidth (2 bus crossings per datagram)



Network Layer 4-24

## Switching via a bus

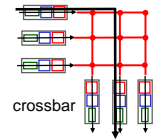
- ❖ datagram from input port memory to output port memory via a shared bus
- ❖ **bus contention**: switching speed limited by bus bandwidth
- ❖ 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers



Network Layer 4-25

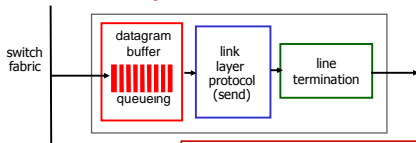
## Switching via interconnection network

- ❖ 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



Network Layer 4-26

## Output ports



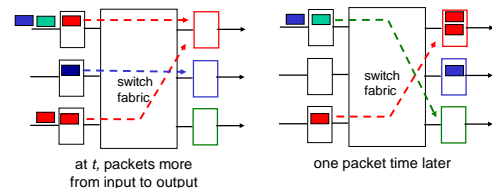
Datagram (packets) can be lost due to congestion, lack of buffers

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

Priority scheduling – who gets best performance, network neutrality

Network Layer 4-27

## Output port queueing



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

Network Layer 4-28

## How much buffering?

- ❖ RFC 3439 rule of thumb: average buffering equal to “typical” RTT (say 200 msec) times link capacity C
  - e.g., C = 10 Gbps link: 2 Gbit buffer
- ❖ recent recommendation: with N flows, buffering equal to

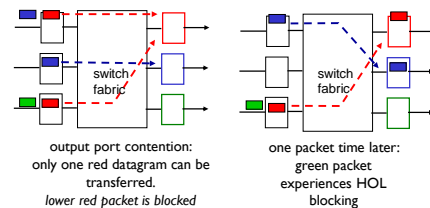
$$\frac{RTT \cdot C}{\sqrt{N}}$$

- ❖ e.g., if 4 flows in example above: 1 Gbit buffer

Network Layer 4-29

## Input port queueing

- ❖ 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



output port contention: only one red datagram can be transferred. lower red packet is blocked

one packet time later: green packet experiences HOL blocking

Network Layer 4-30

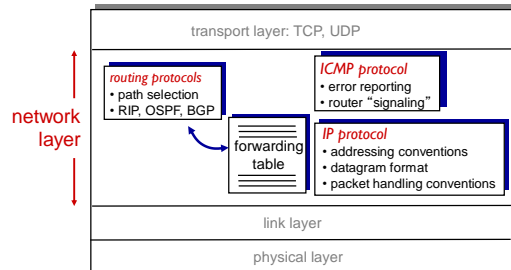
## Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol**
  - datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 routing algorithms
  - link state
  - distance vector
  - hierarchical routing
- 4.6 routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 broadcast and multicast routing

Network Layer 4-31

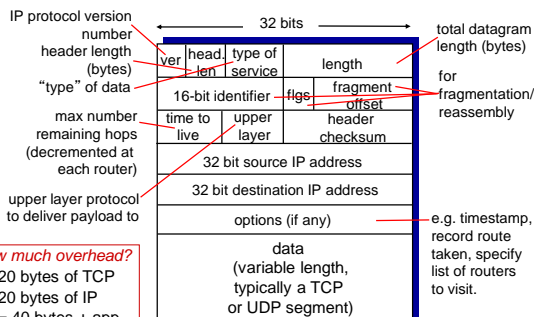
## The Internet network layer

host, router network layer functions:



Network Layer 4-32

## IP datagram format



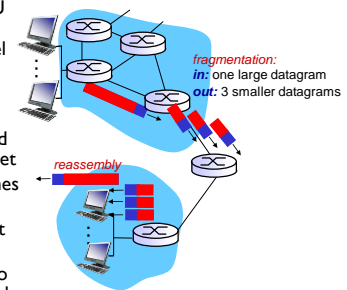
**how much overhead?**

- ❖ 20 bytes of TCP
- ❖ 20 bytes of IP
- ❖ = 40 bytes + app layer overhead

Network Layer 4-33

## 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

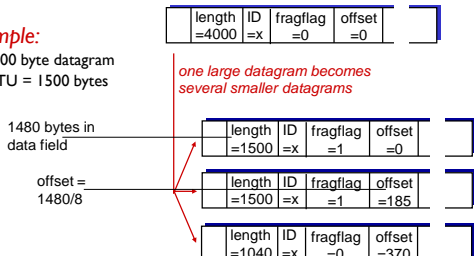


Network Layer 4-34

## IP fragmentation, reassembly

**example:**

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes



Network Layer 4-35

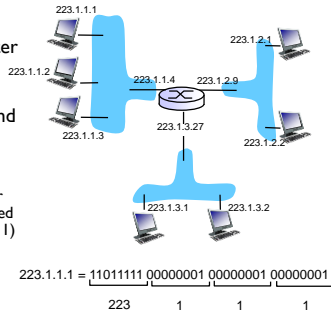
## Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol**
  - datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 routing algorithms
  - link state
  - distance vector
  - hierarchical routing
- 4.6 routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 broadcast and multicast routing

Network Layer 4-36

## 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**



Network Layer 4-37

## IP addressing: introduction

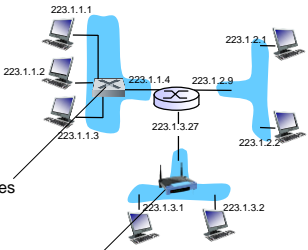
**Q:** how are interfaces actually connected?

**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)

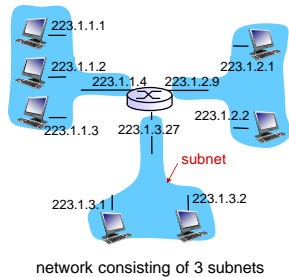
**A:** wireless WiFi interfaces connected by WiFi base station



Network Layer 4-38

## 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*



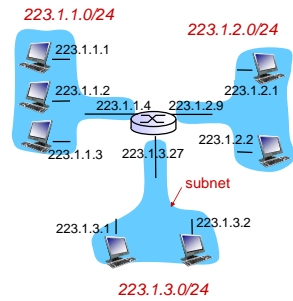
network consisting of 3 subnets

Network Layer 4-39

## 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**

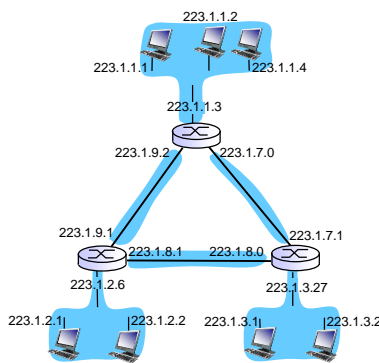


subnet mask: /24

Network Layer 4-40

## Subnets

how many?



Network Layer 4-41

## 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



Network Layer 4-42

## IP addresses: how to get one?

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"

Network Layer 4-43

## DHCP: Dynamic Host Configuration Protocol

**goal:** allow host to *dynamically* obtain its IP address from network server when it joins network

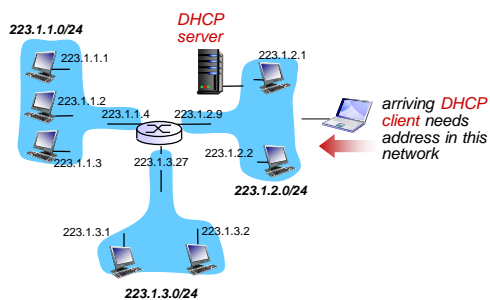
- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/"on")
- support for mobile users who want to join network (more shortly)

**DHCP overview:**

- host broadcasts "DHCP discover" msg [optional]
- DHCP server responds with "DHCP offer" msg [optional]
- host requests IP address: "DHCP request" msg
- DHCP server sends address: "DHCP ack" msg

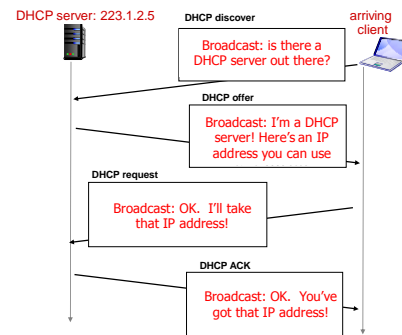
Network Layer 4-44

## DHCP client-server scenario



Network Layer 4-45

## DHCP client-server scenario



Network Layer 4-46

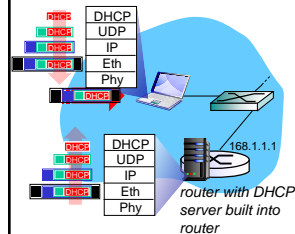
## DHCP: more than IP addresses

DHCP can return more than just allocated IP address on subnet:

- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

Network Layer 4-47

## DHCP: example

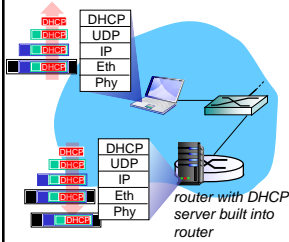


- ❖ connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP
- ❖ DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.1 Ethernet
- ❖ Ethernet frame broadcast (dest: FFFFFFFF) on LAN, received at router running DHCP server
- ❖ Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

Network Layer 4-48



## DHCP: example



- ❖ DHCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- ❖ encapsulation of DHCP server, frame forwarded to client, demuxing up to DHCP at client
- ❖ client now knows its IP address, name and IP address of DNS server, IP address of its first-hop router

Network Layer 4-49

## IP addresses: how to get one?

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

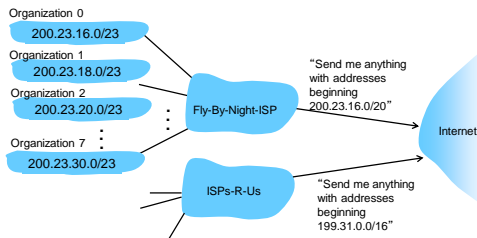
**A:** gets allocated portion of its provider ISP's address space

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

Network Layer 4-50

## Hierarchical addressing: route aggregation

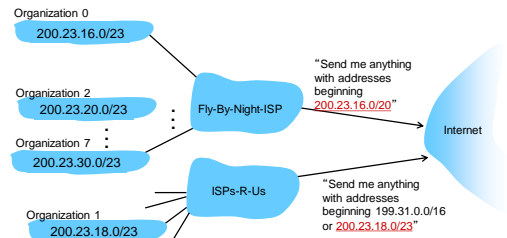
hierarchical addressing allows efficient advertisement of routing information:



Network Layer 4-51

## Hierarchical addressing: more specific routes

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



Network Layer 4-52

## IP addressing: the last word...

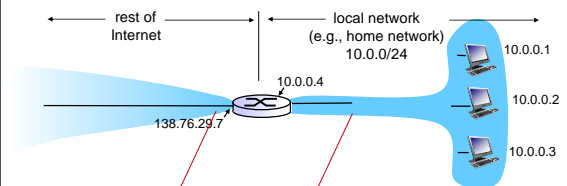
**Q:** how does an ISP get block of addresses?

**A:** ICANN: Internet Corporation for Assigned Names and Numbers <http://www.icann.org/>

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

Network Layer 4-53

## NAT: network address translation



all datagrams leaving local network have same single source NAT IP address: 138.76.29.7, different source port numbers

datagrams with source or destination in this network have 10.0.0/24 address for source, destination (as usual)

Network Layer 4-54

## 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)

Network Layer 4-55

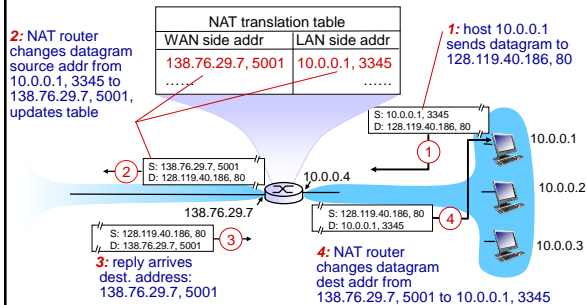
## NAT: network address translation

**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

Network Layer 4-56

## NAT: network address translation



Network Layer 4-57

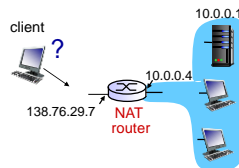
## NAT: network address translation

- ❖ 16-bit port-number field:
  - 60,000 simultaneous connections with a single LAN-side address!
- ❖ NAT is controversial:
  - 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

Network Layer 4-58

## NAT traversal problem

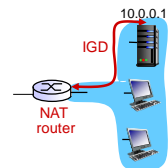
- ❖ client wants to connect to server with address 10.0.0.1
  - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - only one externally visible NATed address: 138.76.29.7
- ❖ **solution 1:** statically configure NAT to forward incoming connection requests at given port to server
  - e.g., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000



Network Layer 4-59

## NAT traversal problem

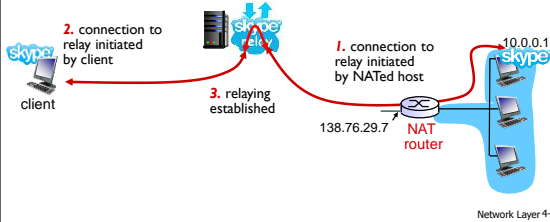
- ❖ **solution 2:** Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATed host to:
    - ❖ learn public IP address (138.76.29.7)
    - ❖ add/remove port mappings (with lease times)
- i.e., automate static NAT port map configuration



Network Layer 4-60

## 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



Network Layer 4-61

## Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
  - datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 routing algorithms
  - link state
  - distance vector
  - hierarchical routing
- 4.6 routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 broadcast and multicast routing

Network Layer 4-62

## 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

Network Layer 4-63

## 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



Network Layer 4-64

## 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

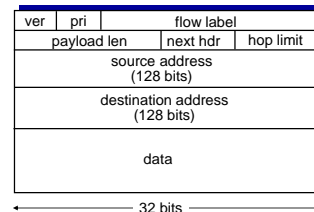
### IPv6 datagram format:

- fixed-length 40 byte header
- no fragmentation allowed

Network Layer 4-65

## 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:** identify upper layer protocol for data



Network Layer 4-66

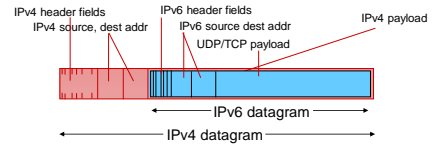
## Other changes from IPv4

- ❖ **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

Network Layer 4-67

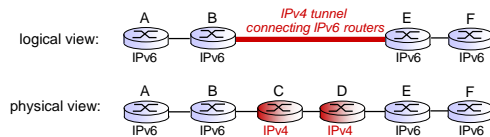
## Transition from IPv4 to IPv6

- ❖ 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



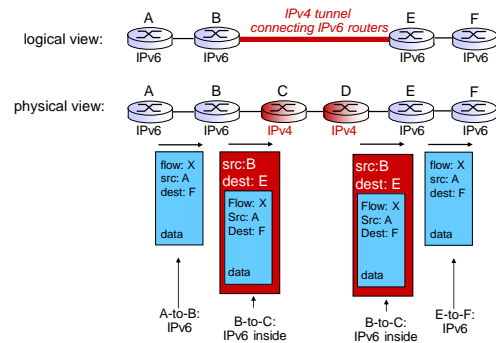
Network Layer 4-68

## Tunneling



Network Layer 4-69

## Tunneling



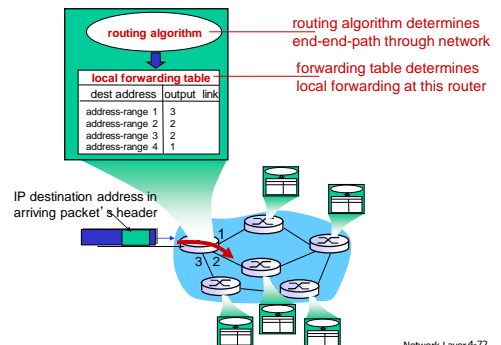
Network Layer 4-70

## Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
  - datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 routing algorithms
  - link state
  - distance vector
  - hierarchical routing
- 4.6 routing in the Internet
  - RIP
  - OSPF
  - BGP
- 4.7 broadcast and multicast routing

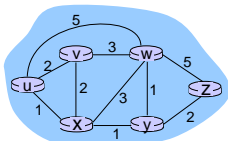
Network Layer 4-71

## Interplay between routing, forwarding



Network Layer 4-72

## Graph abstraction



graph:  $G = (N, E)$

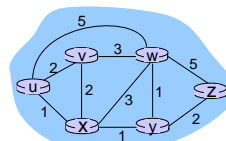
$N$  = set of routers = { u, v, w, x, y, z }

$E$  = set of links = { (u,v), (u,x), (u,w), (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

Network Layer 4-73

## Graph abstraction: costs



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

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

Network Layer 4-74

## 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

**Q: Load sensitive?**

- To reflect current level of congestion

Network Layer 4-75

## 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
  - 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
- $D(v)$ : current value of cost of path from source to dest. v
- $p(v)$ : predecessor node along path from source to v
- $N'$ : set of nodes whose least cost path definitively known

Network Layer 4-76

## Dijkstra's Algorithm

1 **Initialization:**

- $N' = \{u\}$
- for all nodes v
- if v adjacent to u
- then  $D(v) = c(u, v)$
- else  $D(v) = \infty$

8 **Loop**

- find w not in  $N'$  such that  $D(w)$  is a minimum
- add w to  $N'$
- update  $D(v)$  for all v adjacent to w and not in  $N'$ :  
 $D(v) = \min(D(v), D(w) + c(w, v))$
- /\* new cost to v is either old cost to v or known shortest path cost to w plus cost from w to v \*/
- until all nodes in  $N'$**

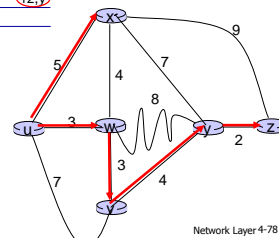
Network Layer 4-77

## 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	6, w	11, w	$\infty$	
2	uwx	6, w	6, w	11, w	14, x	
3	uwxy	6, w	6, w	10, y	14, x	
4	uwxyv	6, w	6, w	10, y	12, y	
5	uwxyvz	6, w	6, w	10, y	12, y	12, y

**notes:**

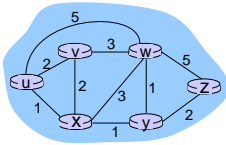
- construct shortest path tree by tracing predecessor nodes
- ties can exist (can be broken arbitrarily)



Network Layer 4-78

## Dijkstra's algorithm: another example

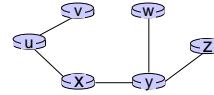
Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	$\infty$	$\infty$
1	ux	2,u	4,x	$\infty$	2,x	$\infty$
2	uxy	2,u	3,y	$\infty$	4,y	$\infty$
3	uxyv	2,u	3,y	$\infty$	4,y	4,y
4	uxyvw	2,u	3,y	$\infty$	4,y	4,y
5	uxyvwz	2,u	3,y	$\infty$	4,y	4,y



Network Layer 4-79

## Dijkstra's algorithm: example (2)

resulting shortest-path tree from u:



resulting forwarding table in u:

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

Network Layer 4-80

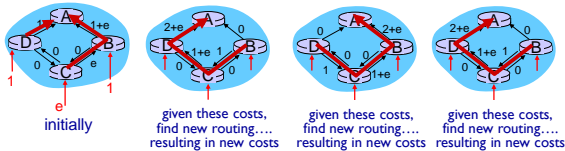
## Dijkstra's algorithm, discussion

**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:**

- e.g., support link cost equals amount of carried traffic:



Network Layer 4-81

## Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Network Layer 4-82

## Distance vector algorithm

**Bellman-Ford equation (dynamic programming)**

let

$d_x(y)$  := cost of least-cost path from x to y

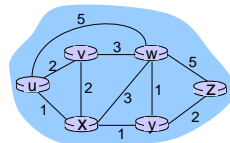
then

$$d_x(y) = \min_v \{ c(x,v) + d_v(y) \}$$

$c(x,v)$ : cost to neighbor v  
 $d_v(y)$ : cost from neighbor v to destination y  
 $\min$  taken over all neighbors v of x

Network Layer 4-83

## Bellman-Ford example



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

B-F equation says:

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

node achieving minimum is next hop in shortest path, used in forwarding table

Network Layer 4-84

## Distance vector algorithm

- ❖  $D_x(y)$  = estimate of least cost from x to y
  - x maintains distance vector  $D_x = [D_x(y): y \in N]$
- ❖ node x:
  - knows cost to each neighbor v:  $c(x,v)$
  - maintains its neighbors' distance vectors. For each neighbor v, x maintains  $D_v = [D_v(y): y \in N]$

Network Layer 4-85

## 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:

$$D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\} \text{ for each node } y \in N$$

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

Network Layer 4-86

## Distance vector algorithm

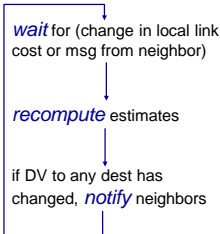
### iterative, asynchronous:

- each local iteration caused by:
  - ❖ local link cost change
  - ❖ DV update message from neighbor

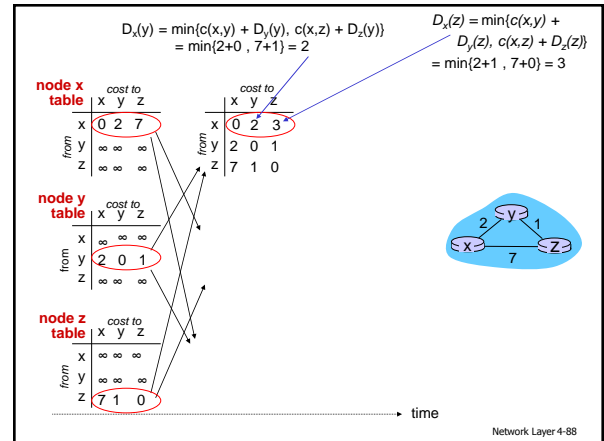
### distributed:

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

### each node:

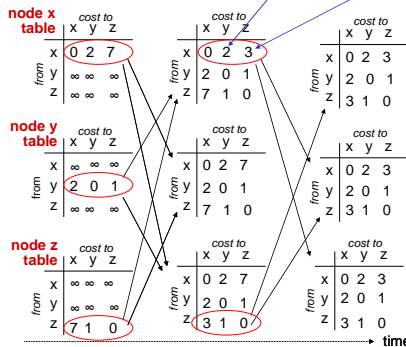


Network Layer 4-87



$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\} = \min\{2+0, 7+1\} = 2$$

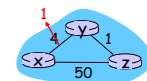
$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\} = \min\{2+1, 7+0\} = 3$$



## 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.

Network Layer 4-90