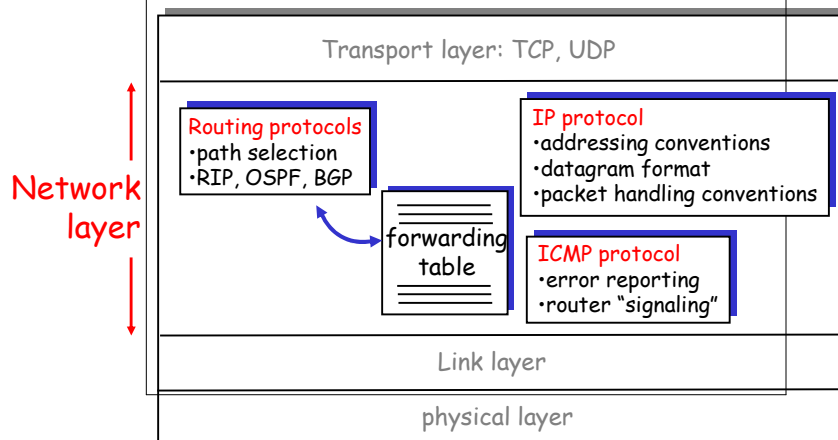


The Internet Network layer

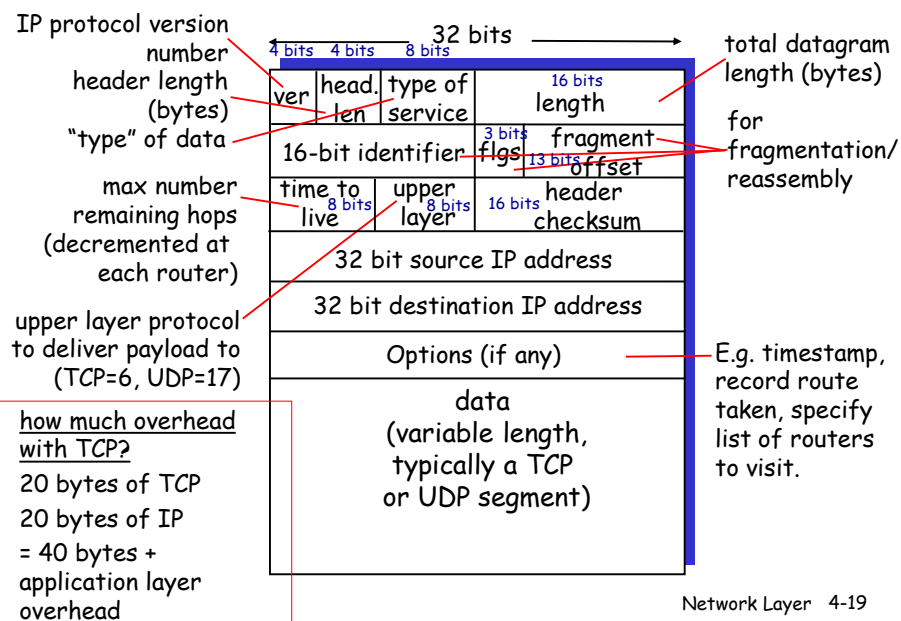
Host, router network layer functions:



Network Layer 4-18

18

IP datagram format



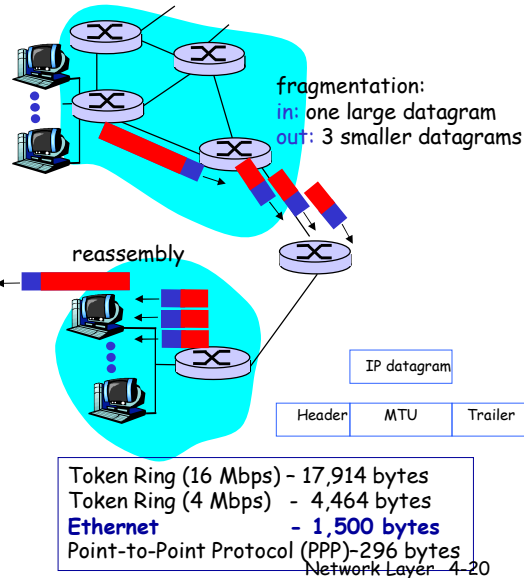
Network Layer 4-19

19

MTU = Maximum Transmission Units

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



20

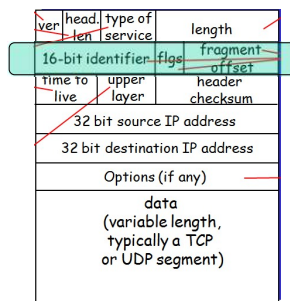
IP Fragmentation and Reassembly

Example

- 4000 byte datagram
- MTU = 1500 bytes

1480 bytes in data field

offset = 1480/8



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

One large datagram becomes several smaller datagrams

length	ID	fragflag	offset
=1500	=x	=1	=0
length	ID	fragflag	offset
=1500	=x	=1	=185
length	ID	fragflag	offset
=1040	=x	=0	=370

Network Layer 4-21

21

Network Layer : Logical Addressing

- o Communication at network layer is host-to-host
- o Computer somewhere in the world need to communicate with another computer somewhere else in the world through Internet
- o Packet transmitted by sending computer may pass through several LANs or WANs before reaching destination computer
- o We need global addressing scheme called logical addressing
- o Today, we use the term IP address to mean a logical address in network layer of TCP/IP protocol suite

Network Layer 4-22

22

IP Addresses

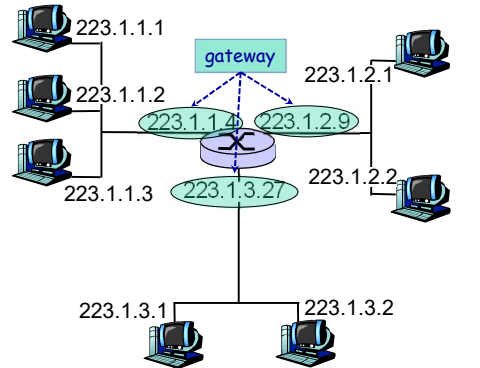
- o The Internet address are 32 bits in length
 - o Address space is 2^{32} or 4,294,967,296
 - o These addresses are referred to as IPv4 (IP version 4) addresses or simply IP address
- o The need for more addresses motivated a new design of the IP layer called new generation of IP or IPv6 (IP version 6)
 - o The Internet uses 128-bit addresses that give much greater flexibility in address location
 - o These addresses are referred to as IPv6 (IP version 6) address

Network Layer 4-23

23

IPv4 Addressing: introduction

- o IPv4 address: 32-bit identifier for *host*, *router interface*
- o *Interface*: connection between *host/router* and physical link
 - o router's typically have multiple interfaces
 - o host typically has one interface
 - o IP addresses associated with *each* interface



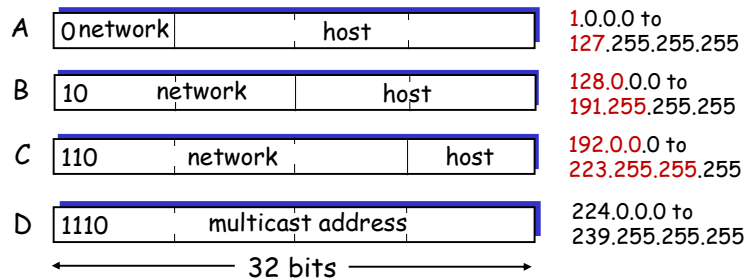
223.1.1.1 = $\underbrace{11011111}_{223} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00000001}_1$

Network Layer 4-24

24

IP Addresses "class-full" addressing:

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



	First byte	Second byte	Third byte	Fourth byte
Class A	0			
Class B	10			
Class C	110			
Class D	1110			
Class E	1111			

a. Binary notation

	First byte	Second byte	Third byte	Fourth byte
Class A	0-127			
Class B	128-191			
Class C	192-223			
Class D	224-239			
Class E	240-255			

b. Dotted-decimal notation

Network Layer 4-25

25

Class	Number of Blocks	Block Size	Application
A	128	16,777,216	Unicast
B	16,384	65,536	Unicast
C	2,097,152	256	Unicast
D	1	268,435,456	Multicast
E	1	268,435,456	Reserved

Class	Network Octets (blanks in the IP address are used for octets identifying hosts)	Total Number of Possible Networks or Licenses	Host Octets (blanks in IP address are used for octets identifying networks)	Total Number of Possible IP Addresses in Each Networks
A	0.____.____.____ to 127.____.____.____	128	____.0.0.1 to ____.255.255.254	16,777,214
B	128.0.____.____ to 191.255.____.____	64x256 16,384	____.____.0.1 to ____.____.255.254	65,534
C	192.0.0.____ to 223.255.255.____	32x256x256 2,097,152	____.____.____.1 to ____.____.____.254	254

Network Layer 4-26

26

Address for Private Networks

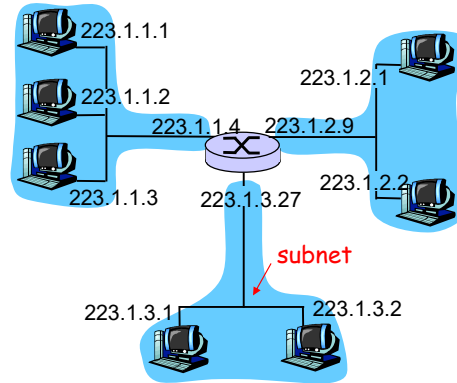
	Range			Total
Class A	10.0.0.0	to	10.255.255.255	2^{24}
Class B	172.16.0.0	to	172.31.255.255	2^{20}
Class C	192.168.0.0	to	192.168.255.255	2^{16}

Network Layer 4-27

27

Subnets

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



network consisting of 3 subnets

Class C

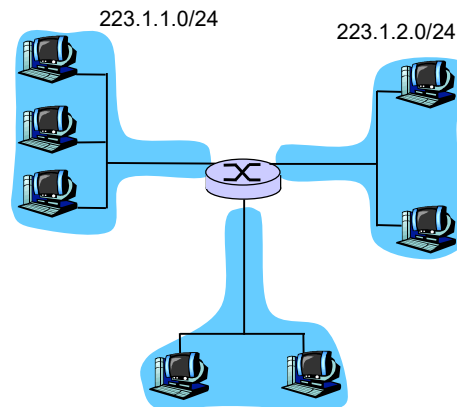
Network Layer 4-28

28

Subnets

Recipe

- o To determine the subnets, detach each interface from its host or router, creating islands of isolated networks.
- o Each isolated network is called a subnet.



Subnet mask: /24

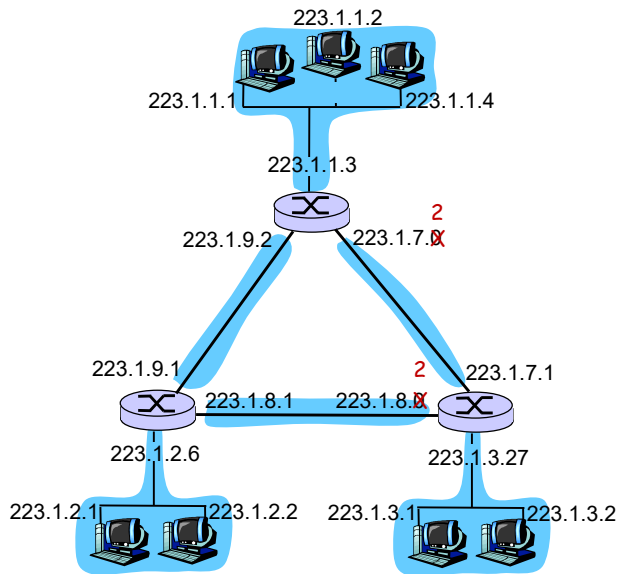
Class C

Network Layer 4-29

29

Subnets

How many?



Network Layer 4-30

30

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



200.23.16.0/23

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 4-31

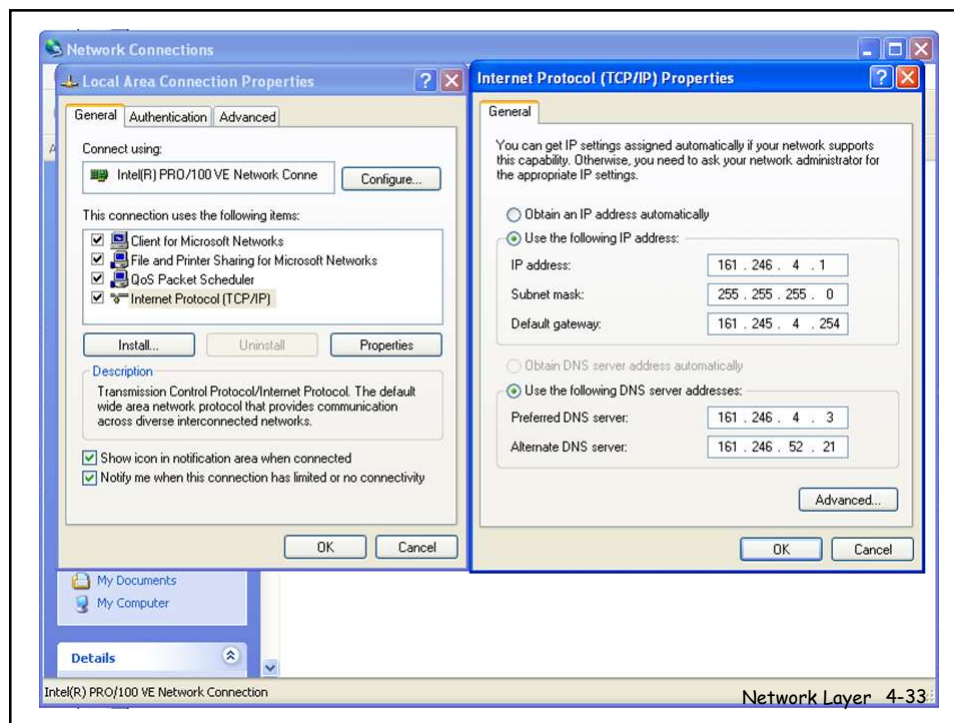
31

IP addresses: how to get one?

- o **Q:** How does *host* get IP address?
- o hard-coded by system admin in a file
 - o **Windows:**
 - o control-panel->network connections->properties
->Internet Protocol (TCP/IP)
 - o **UNIX:** /etc/rc.config
- o **DHCP:** Dynamic Host Configuration Protocol:
dynamically get address from as server
 - o "plug-and-play"
 - o allow host to *dynamically* obtain its IP address from
network server when it joins network

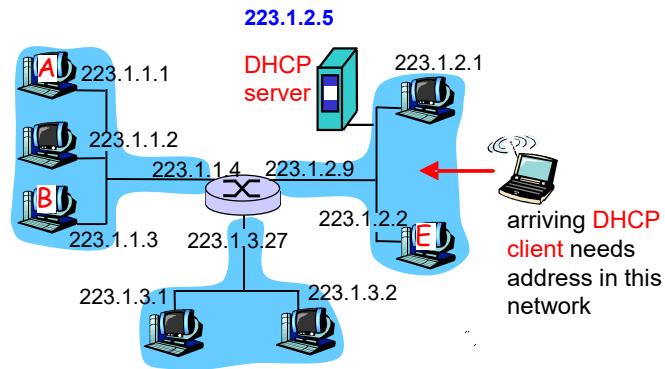
Network Layer 4-32

32



33

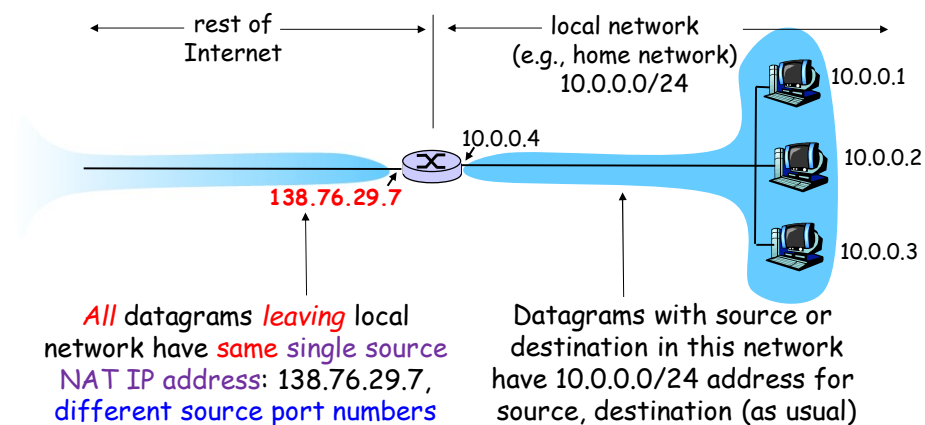
DHCP client-server scenario



Network Layer 4-35

35

NAT: Network Address Translation



	Range		Total
Class A	10.0.0.0	to 10.255.255.255	2^{24}
Class B	172.16.0.0	to 172.31.255.255	2^{20}
Class C	192.168.0.0	to 192.168.255.255	2^{16}

Network Layer 4-41

41

NAT: Network Address Translation

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

Network Layer 4-42

42

NAT: Network Address Translation

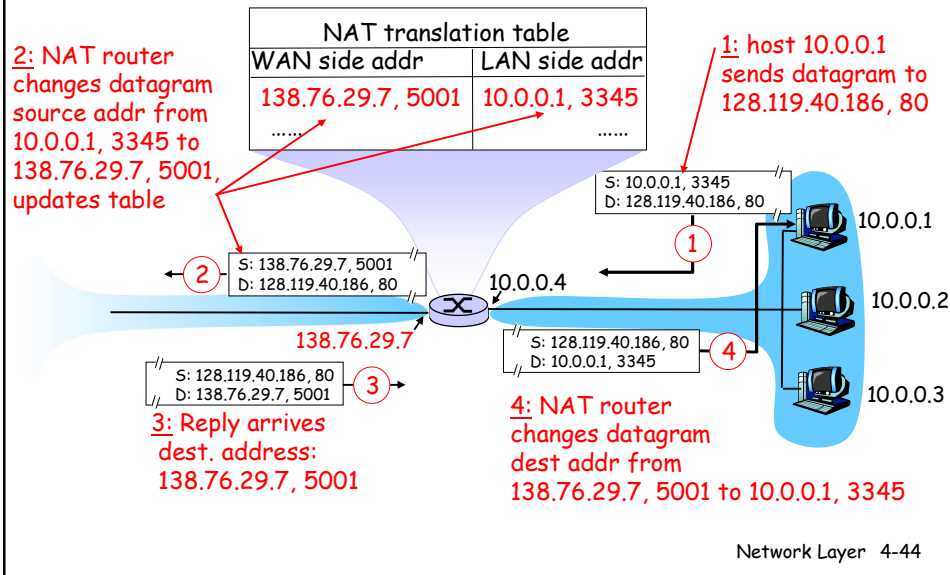
Implementation: NAT router must:

- o 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 address.
- o remember (in NAT translation table) every (source IP address, port #) to (NAT IP address, new port #) translation pair
- o 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 4-43

43

NAT: Network Address Translation



44

Address Mapping

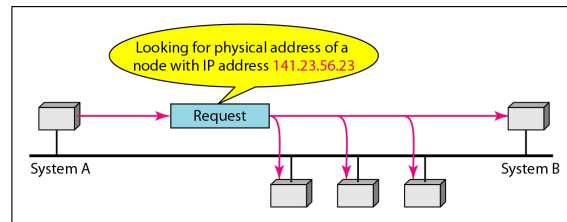
- o Delivery of packet to host or router requires two levels of addressing: **logical address** and **physical address**
- o We need to be able to map a logical address to its corresponding physical address and vice versa.
- o Mapping Logical Address to Physical Address can be done by **Address Resolution Protocol (ARP)** RFC 826



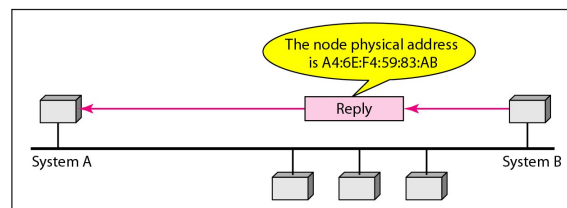
Network Layer 4-49

49

Address Resolution Protocol (ARP)



a. ARP request is broadcast



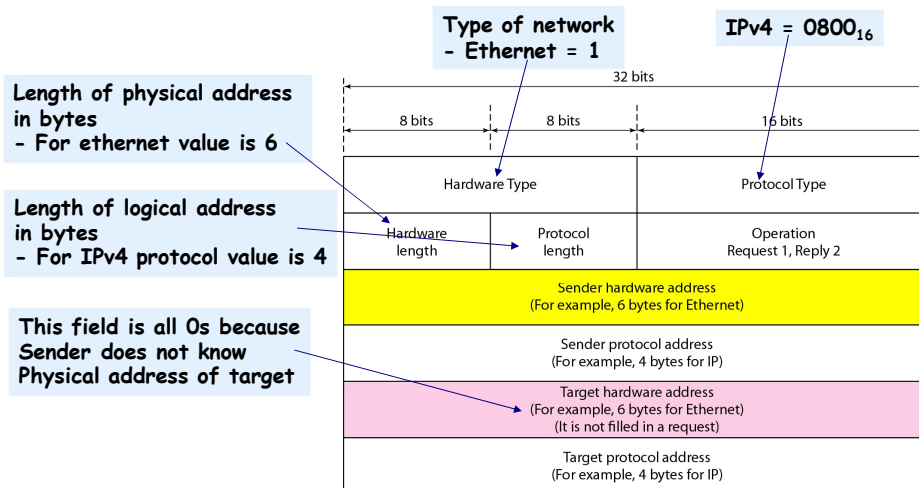
b. ARP reply is unicast

ARP Operation

Network Layer 4-50

50

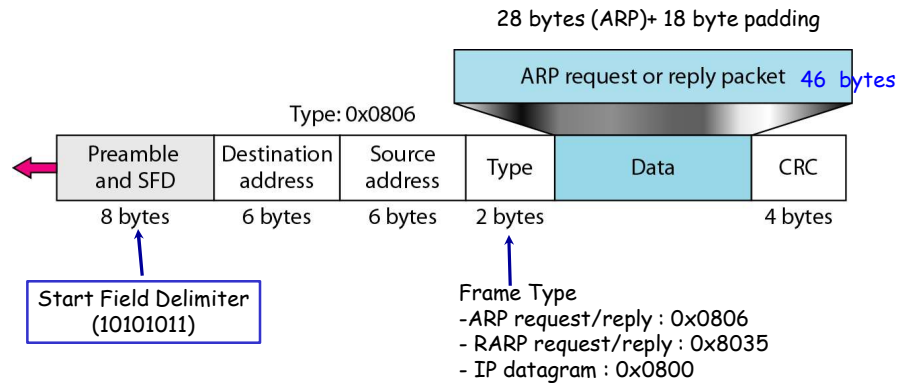
ARP Packet



Network Layer 4-51

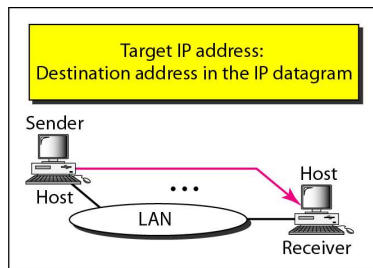
51

Encapsulating of ARP packet

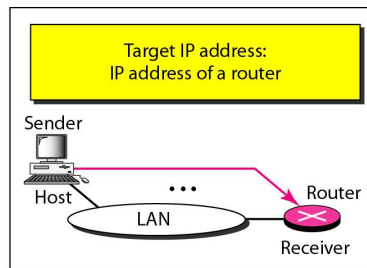


Network Layer 4-52

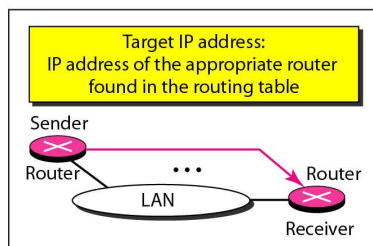
52



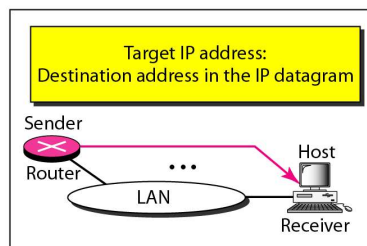
Case 1. A host has a packet to send to another host on the same network.



Case 2. A host wants to send a packet to another host on another network. It must first be delivered to a router.



Case 3. A router receives a packet to be sent to a host on another network. It must first be delivered to the appropriate router.

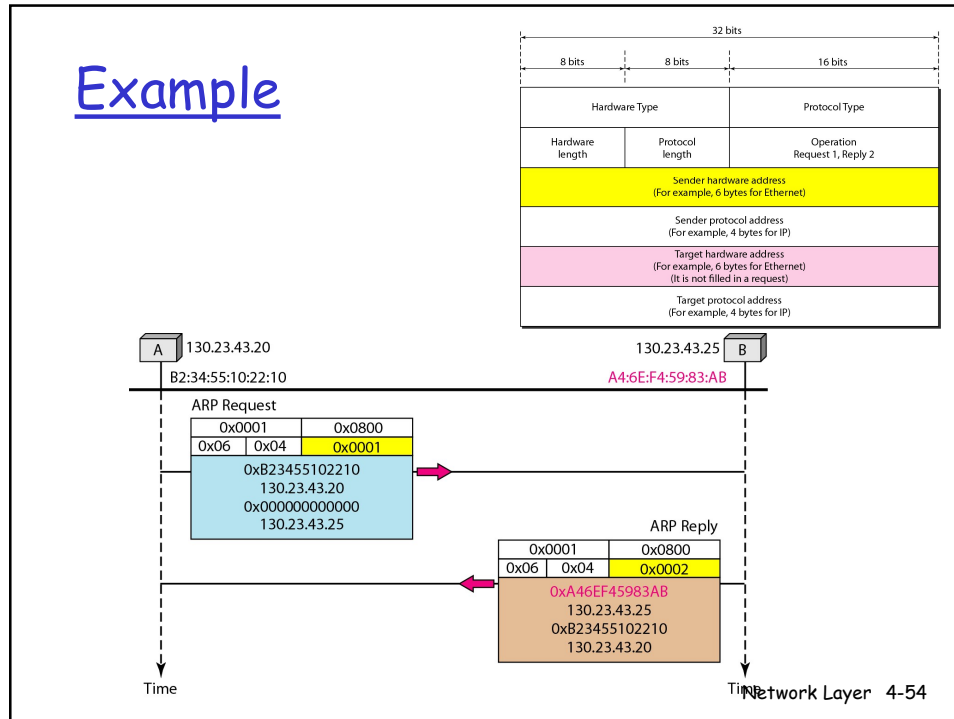


Case 4. A router receives a packet to be sent to a host on the same network.

Network Layer 4-53

53

Example



54

Internet Control Message Protocol (ICMP)

- o IP provides **unreliable** and **connectionless** datagram delivery
- o IP protocol is a **best-effort delivery service** that **delivers datagram** from original **source** to final **destination**
- o Two deficiencies
 - o **Lack of error control**
 - o No error-reporting or error-correcting mechanism
 - o **Lack of assistance mechanism for host and management queries**
 - o Host sometimes **needs** to **determine** if **router** or **another host** is **alive**
 - o Sometimes a network administrator **needs** **information** from **another host** or **router**

Network Layer 4-57

57

ICMP : Type of Messages

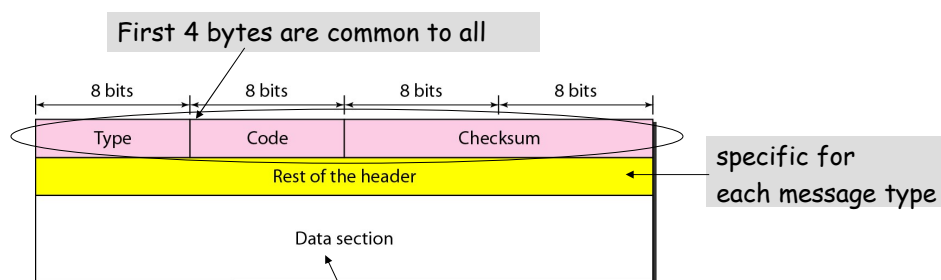
- o ICMP message are divided into two broad categories
 - o Error-reporting message
 - o Report problems that router or host (destination) may encounter when it processes IP packet
 - o Query message
 - o Help host or network manager get specific information from router or another host
 - o Ex. Nodes can discover their neighbors
 - o Hosts can discover and learn about routers on their network

Network Layer 4-58

58

ICMP : Message Format

- o ICMP message has 8-byte header and variable-size data section



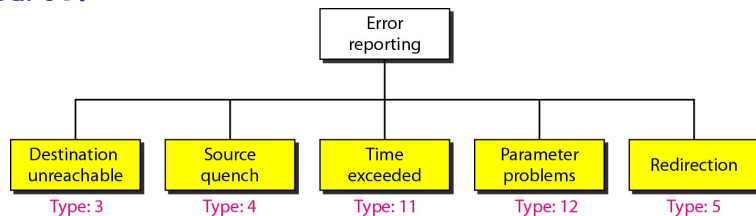
- o In error messages carries information for finding original packet that had error
- o In query messages data section carries extra information based on type of query

Network Layer 4-59

59

ICMP : Error Reporting

- **ICMP** always reports error messages to original source.

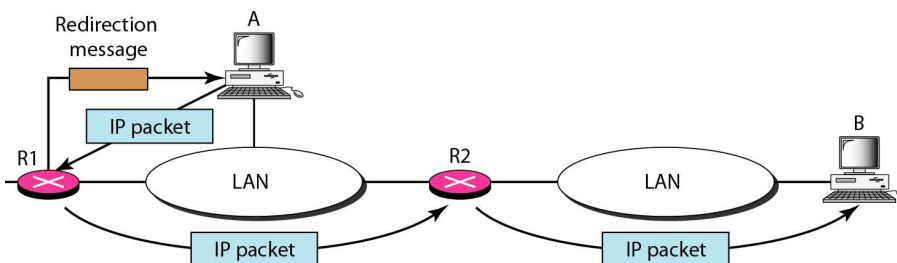


Type	Code	description
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
11	0	TTL expired
12	0	bad IP header

Network Layer 4-60

60

Redirection Concept



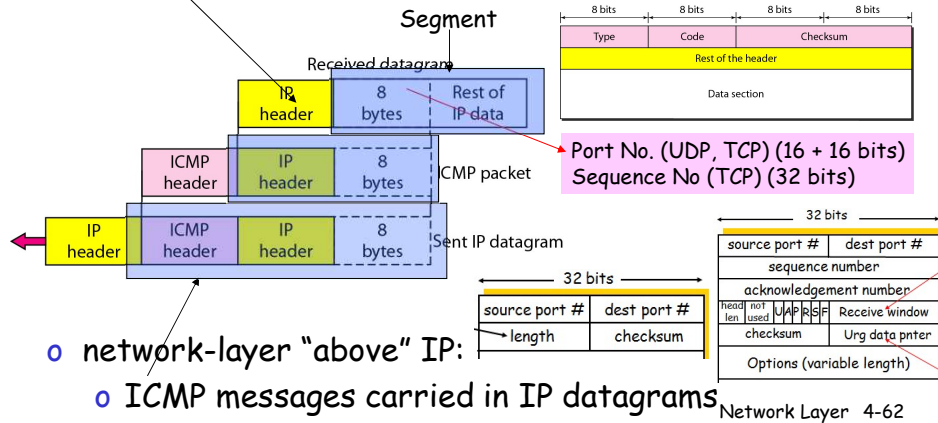
- Host A wants to send a datagram to host B
- Router R2 is obviously **most efficient routing** choice, but host A did not choose router R2
- Datagram goes to R1 instead
- Router R1, after consulting its table, **finds that packet should have gone to R2**
- R1 sends packet to R2 and, at the same time, **sends a redirection message** to host A
- Host A's routing table can now be updated

Network Layer 4-61

61

ICMP : Error Reporting (continued)

- o All error messages contain data section that includes
 - o IP header of the original datagram plus
 - o the first 8 bytes of data in that datagram



62

Type: 3	Code: 0 to 15	Checksum
Unused (All 0s)		
Part of the received IP datagram including IP header plus the first 8 bytes of datagram data		

Destination Unreachable

Type: 4	Code: 0	Checksum
Unused (All 0s)		
Part of the received IP datagram including IP header plus the first 8 bytes of datagram data		

Source Quench

0 : Time to Live exceeded in Transit
1 : Fragment Reassembly Time Exceeded

Type: 11	Code: 0 or 1	Checksum
Unused (All 0s)		
Part of the received IP datagram including IP header plus the first 8 bytes of datagram data		

TTL Expired

Network Layer 4-63

63

Code 0 : Error in one of header field, value in pointer field points to byte with problem
 Code 1 : required part of option is missing

Type: 12	Code: 0 or 1	Checksum
Pointer	Unused (All 0s)	
Part of the received IP datagram including IP header plus the first 8 bytes of datagram data		

Parameter Problem

Type: 5	Code: 0 to 3	Checksum
IP address of the target router		
Part of the received IP datagram including IP header plus the first 8 bytes of datagram data		

Redirection message format

Code 0 - Redirection for network (or Subnet) -specific route
 Code 1 - Redirection for host-specific route
 Code 2 - same as code 0 , based on specified type of service
 Code 3 - same as code 1, based on specified type of service

Network Layer 4-64

64

ICMP : Query

Type	Code	description
0	0	echo reply (ping)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery

- Query message is encapsulated in IP packet, which in turn is encapsulated in data link layer frame
- In this case, no bytes of original IP are included in message



Network Layer 4-65

65

Echo-request and echo-reply messages

8: Echo request
0: Echo reply

Type: 8 or 0	Code: 0	Checksum
Identifier		Sequence number
Optional data Sent by the request message; repeated by the reply message		

Network Layer 4-66

66

Route discovery and Route Advertisement

Route discovery

Type: 10	Code: 0	Checksum
Identifier		Sequence number

Route advertisement

number of router advertisements in message

Type: 9	Code: 0	Checksum
Number of addresses	Address entry size	Lifetime
IPv4 address router		
Router address 1		
Address preference 1		
Router address 2		
Address preference 2		
⋮		

information for each router address entry in the list. The value is normally set to 2 (router address + preference level).

Ranking
0 - default
80000000₁₆

Ranking of router and used to select router as default router

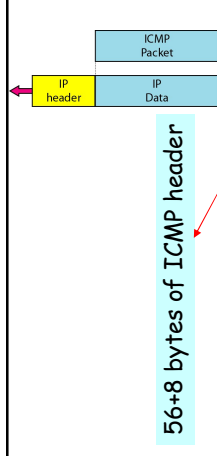
Network Layer 4-67

67

Ping

- o Ping program is used to find whether a host is alive or not

56+8 bytes of ICMP header + 20 bytes of IP header



```

$ ping fhda.edu
PING fhda.edu (153.18.8.1) 56 (84) bytes of data:
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=0    ttl=62    time=1.91 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=1    ttl=62    time=2.04 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=2    ttl=62    time=1.90 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=3    ttl=62    time=1.97 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=4    ttl=62    time=1.93 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=5    ttl=62    time=2.00 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=6    ttl=62    time=1.94 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=7    ttl=62    time=1.94 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=8    ttl=62    time=1.97 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=9    ttl=62    time=1.89 ms
64 bytes from tiptoe.fhda.edu (153.18.8.1): icmp_seq=10   ttl=62    time=1.98 ms

--- fhda.edu ping statistics ---
11 packets transmitted, 11 received, 0% packet loss, time 10103ms
rtt min/avg/max = 1.899/1.955/2.041 ms
Network Layer 4-68
  
```

68

Traceroute

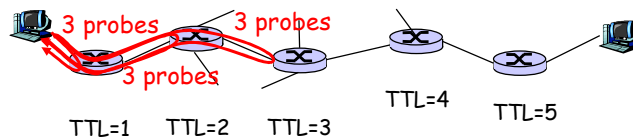
- o Source sends series of UDP segments to dest
 - o First has TTL =1
 - o Second has TTL=2, etc.
 - o Unlikely port number
- o When nth datagram arrives to nth router:
 - o Router discards datagram
 - o And sends to source an ICMP message (type 11, code 0) : TTL expired
 - o Message includes name of router & IP address
- o When ICMP message arrives, source calculates RTT
- o Traceroute does this 3 times
- o Stopping criterion
- o UDP segment eventually arrives at destination host
- o Destination returns ICMP "port unreachable" packet (type 3, code 3 : Destination Port Unreachable)
- o When source gets this ICMP, stops.

Network Layer 4-69

69

"Real" Internet delays and routes

- o What do "real" Internet delay & loss look like?
- o Traceroute program (tracert for windows) :
provides delay measurement from source to router along end-end Internet path towards destination.
For all i :
 - o sends three packets that will reach router i on path towards destination
 - o router i will return packets to sender
 - o sender times interval between transmission and reply.



Network Layer 4-70

70

"Real" Internet delays and routes

traceroute: gaia.cs.umass.edu to www.eurecom.fr

Three delay measurements from gaia.cs.umass.edu to cs-gw.cs.umass.edu

```

1 cs-gw (128.119.240.254) 1 ms 1 ms 2 ms
2 border1-rt-fa5-1-0.gw.umass.edu (128.119.3.145) 1 ms 1 ms 2 ms
3 cht-vbns.gw.umass.edu (128.119.3.130) 6 ms 5 ms 5 ms
4 jn1-at1-0-0-19.wor.vbns.net (204.147.132.129) 16 ms 11 ms 13 ms
5 jn1-so7-0-0-0.wae.vbns.net (204.147.136.136) 21 ms 18 ms 18 ms
6 abilene-vbns.abilene.ucaid.edu (198.32.11.9) 22 ms 18 ms 22 ms
7 nycm-wash.abilene.ucaid.edu (198.32.8.46) 22 ms 22 ms 22 ms
8 62.40.103.253 (62.40.103.253) 104 ms 109 ms 106 ms
9 de2-1.de1.de.geant.net (62.40.96.129) 109 ms 102 ms 104 ms
10 de.fr1.fr.geant.net (62.40.96.50) 113 ms 121 ms 114 ms
11 renater-gw.fr1.fr.geant.net (62.40.103.54) 112 ms 114 ms 112 ms
12 nio-n2.cssi.renater.fr (193.51.206.13) 111 ms 114 ms 116 ms
13 nice.cssi.renater.fr (195.220.98.102) 123 ms 125 ms 124 ms
14 r3t2-nice.cssi.renater.fr (195.220.98.110) 126 ms 126 ms 124 ms
15 eurecom-valbonne.r3t2.ft.net (193.48.50.54) 135 ms 128 ms 133 ms
16 194.214.211.25 (194.214.211.25) 126 ms 128 ms 126 ms
17 ***
18 ***
19 fantasia.eurecom.fr (193.55.113.142) 132 ms 128 ms 136 ms
    
```

trans-oceanic link

* means no response (probe lost, router not replying)

Network Layer 4-71

71

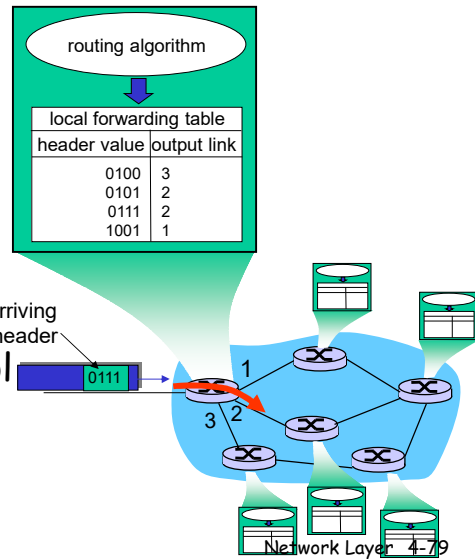
Routing Algorithms

o Link state

- o Ex. Open Shortest Path First (OSPF)

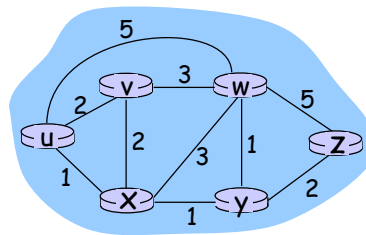
o Distance Vector

- o Ex. Routing Information Protocol (RIP)



79

Graph abstraction



Graph: $G = (N, E)$

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

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

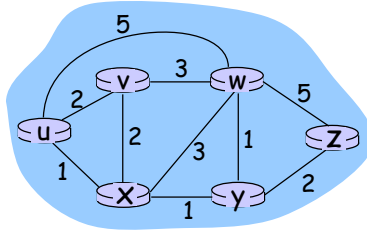
Remark: Graph abstraction is useful in other network contexts

Example: P2P, where N is set of peers and E is set of TCP connections

Network Layer 4-82

82

Graph abstraction: costs



• $c(x, x') = \text{cost of link } (x, x')$

- e.g., $c(w, z) = 5$

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

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

Question: What's the least-cost path between u and z ?

Routing algorithm: algorithm that finds least-cost path

Network Layer 4-83

83

Routing Algorithm classification

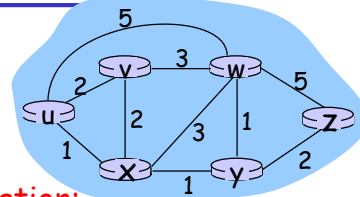
- Global or decentralized information?
- Global:
 - all routers have complete topology, link cost information
 - "link state" algorithms
- Decentralized:
 - router knows physically-connected neighbors, link costs to neighbors
 - iterative process of computation, exchange of information with neighbors
 - "distance vector" algorithms
- Static or dynamic?
- Static:
 - routes change slowly over time
- Dynamic:
 - routes change more quickly
 - periodic update
 - in response to link cost changes

Network Layer 4-84

84

A Link-State Routing Algorithm

- o **Dijkstra's algorithm**
- o **network topology, link costs** known to all nodes
 - o accomplished via "link state **broadcast**"
 - o **all nodes** have **same information**
- o computes least cost paths from one node ('source') to all other nodes
 - o gives **forwarding table** for that node
- o **iterative**: after k iterations, know least cost path to k destination's



- o **Notation:**
- o $c(x,y)$: link cost from node x to y; $= \infty$ if not direct neighbors
- o $D(v)$: current value of **cost** of path from **source** to destination v
- o $p(v)$: **predecessor node** along path from **source** to v
- o N' : set of nodes whose least cost path definitively known

Network Layer 4-86

86

Dijkstra's Algorithm

1 Initialization:

- 2 $N' = \{u\}$
- 3 for all nodes v
- 4 if v adjacent to u
- 5 then $D(v) = c(u,v)$
- 6 else $D(v) = \infty$

หา Cost จาก U ไปยัง Node ข้างเคียงก่อน

8 Loop

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

เลือก node ที่ให้ cost น้อยที่สุด เป็นสมาชิกของ N'

Network Layer 4-87

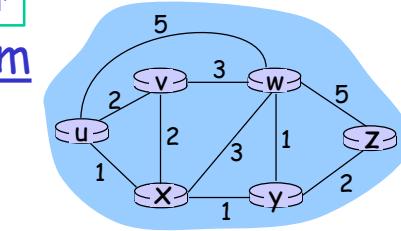
87

P : โหนดใดๆ ในเครือข่ายที่ไม่เป็นสมาชิกของ N'

Dijkstra's Algorithm

1 Initialization:

- 2 $N' = \{u\}$
- 3 for all nodes P
- 4 if P adjacent to u
- 5 then $D(P) = c(u, P)$
- 6 else $D(P) = \infty$



หา Cost จาก U ไปยัง Node ข้างเคียงก่อน

8 Loop

- 9 find Q not in N' such that $D(Q)$ is a minimum
- 10 add Q to N'
- 11 update $D(P)$ for all P adjacent to Q and not in N' :
- 12 $D(P) = \min(D(P), D(Q) + c(Q, P))$
- 13 /* new cost to P is either old cost to P or known
- 14 shortest path cost to Q plus cost from Q to P */
- 15 until all nodes in N'

เลือก node ที่ให้ cost น้อยที่สุด เป็นสมาชิกของ N'

Q : โหนดใดๆ ในเครือข่ายที่ไม่เป็นสมาชิกของ N' และมีค่า Cost น้อยที่สุด

Network Layer 4-88

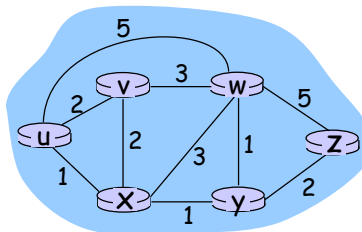
88

Dijkstra's algorithm: example

$D(P)$: current value of cost of path from source to destination P

$p(P)$: predecessor node along path from source to P

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	∞	∞
1	ux	2, u	4, x		2, x	∞
2	uxy	2, u	3, y			4, y
3	uxyv		3, y			4, y
4	uxyvw					4, y
5	uxyvwz					



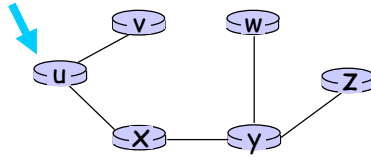
ต้องการหา Shortest Path จาก Node u ไปยังทุก Node

Network Layer 4-89

89

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 4-90

90

Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)

Define

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

Then

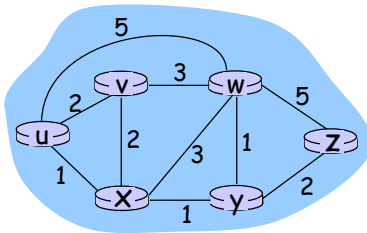
$$d_x(y) = \min_v \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors v of x

Network Layer 4-93

93

Bellman-Ford example



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

B-F equation says:

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

$$= \min \{ 2 + 5, 1 + 3, 5 + 3 \} = 4$$

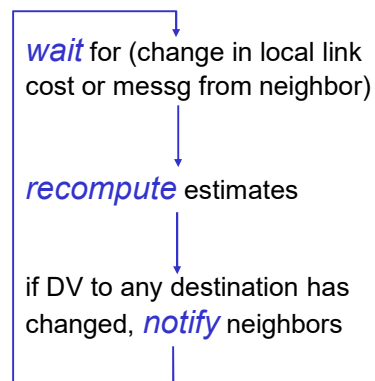
Node that achieves minimum is **next**
hop in shortest path → forwarding table

Network Layer 4-94

94

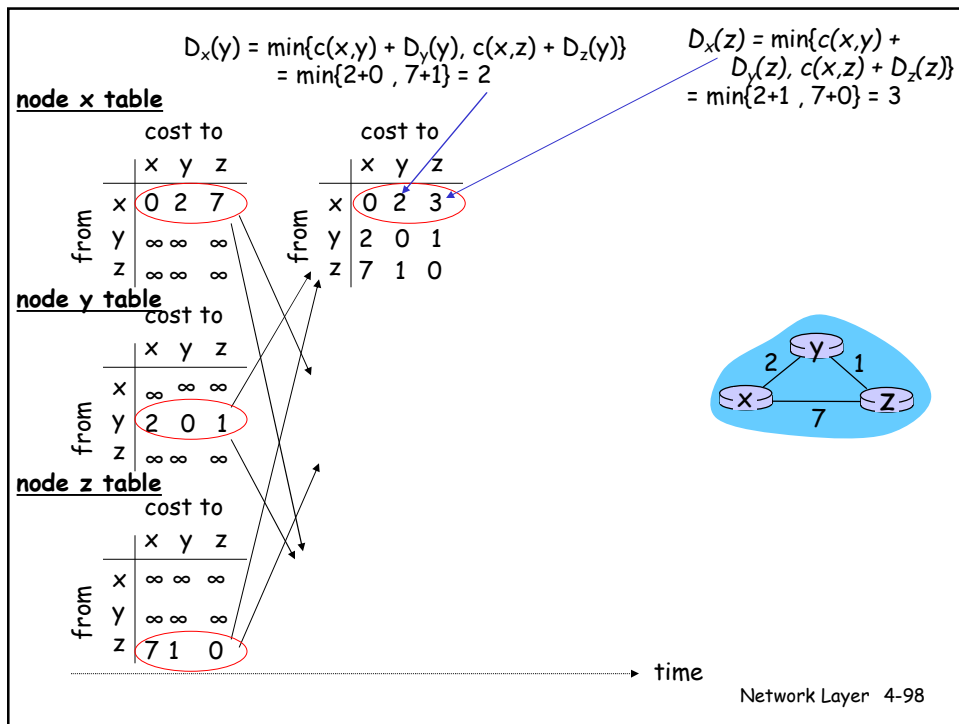
Distance Vector Algorithm

- o **Iterative, asynchronous:** Each node:
each local iteration caused by:
 - o local link cost change
 - o DV update message from neighbor
- o **Distributed:**
 - o each node notifies neighbors only when its DV changes
 - o neighbors then notify their neighbors if necessary

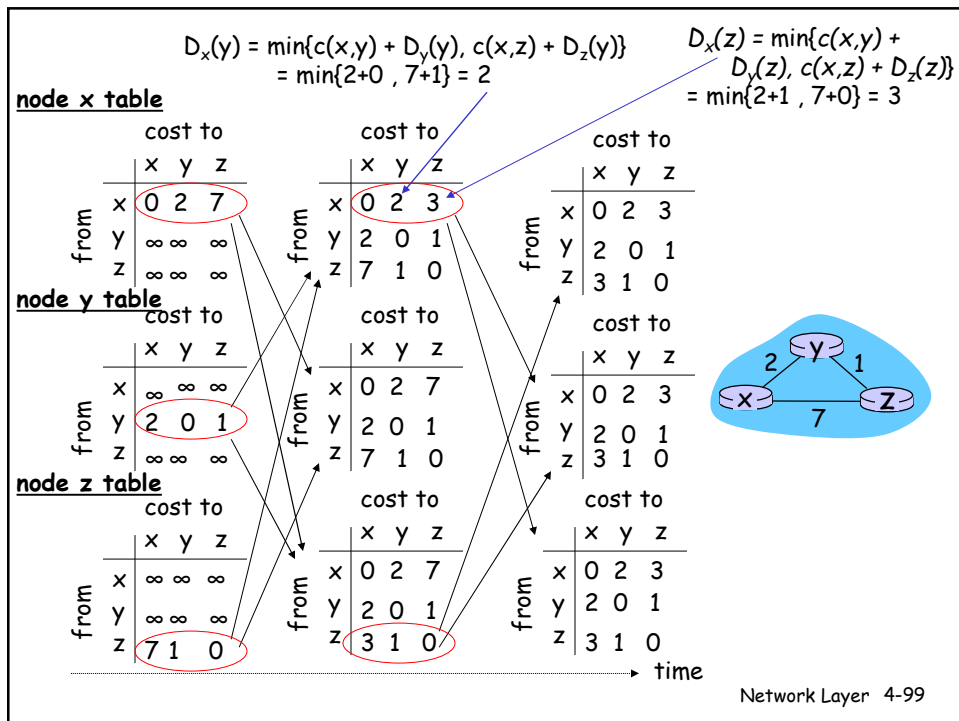


Network Layer 4-97

97



98



99