

15-441 Computer Networking Lecture 10 – Routers and Routing

Peter Steenkiste
Departments of Computer Science and
Electrical and Computer Engineering

15-441 Networking, Spring 2008
<http://www.cs.cmu.edu/~dga/15-441/S08>

1

Outline

- ICMP
- How do Routers Works?
- Routing
- Distance vector

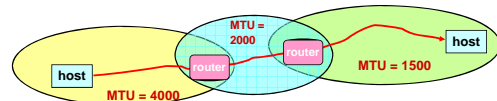
2

Internet Control Message Protocol (ICMP)

- Short messages used to send error & other control information
- Examples
 - » Ping request / response
 - Can use to check whether remote host reachable
 - » Destination unreachable
 - Indicates how packet got & why couldn't go further
 - » Flow control
 - Slow down packet delivery rate
 - » Redirect
 - Suggest alternate routing path for future messages
 - » Router solicitation / advertisement
 - Helps newly connected host discover local router
 - » Timeout
 - Packet exceeded maximum hop limit

3

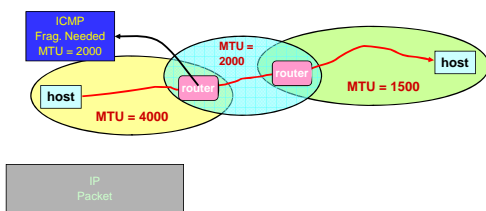
IP MTU Discovery with ICMP



- Typically send series of packets from one host to another
- Typically, all will follow same route
 - » Routes remain stable for minutes at a time
- Makes sense to determine path MTU before sending real packets
- Operation
 - » Send max-sized packet with "do not fragment" flag set
 - » If encounters problem, ICMP message will be returned
 - "Destination unreachable: Fragmentation needed"
 - Usually indicates MTU encountered

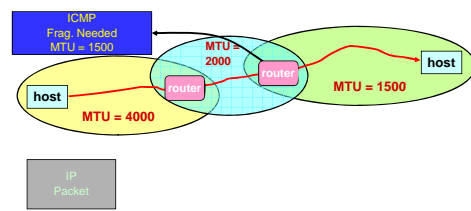
4

IP MTU Discovery with ICMP



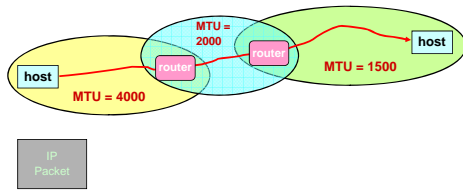
5

IP MTU Discovery with ICMP



6

IP MTU Discovery with ICMP



- » When successful, no reply at IP level
 - “No news is good news”
- » Higher level protocol might have some form of acknowledgement

7

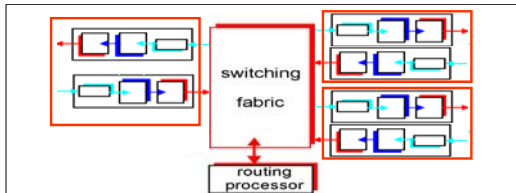
Outline

- ICMP
- How do Routers Works?
- Routing
- Distance vector

8

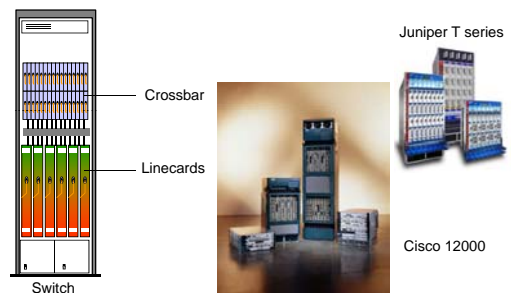
Router Architecture: Key Functions

- Run routing algorithms/protocol (RIP, OSPF, BGP)
 - » Done by routing processor
- **Switching** datagrams from incoming to outgoing link
 - » Common case handled by line cards



9

Router Physical Layout



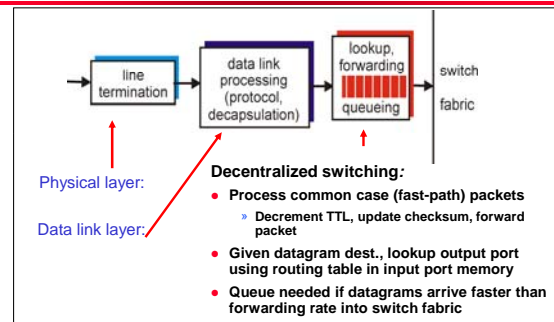
10

Line Cards

- Often uses special purpose hardware (e.g. ASICs)
- Network interface cards
- **Fast path (common-case) processing**
 - » Decrement TTL
 - » Recompute checksum
 - » Forward to next hop line card
 - Forwarding engine

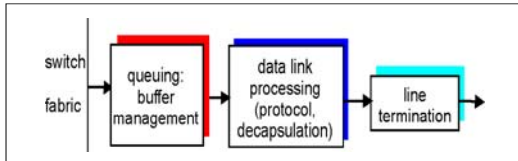
11

Line Card: Input Port



12

Line Card: Output Port



- Queuing required when datagrams arrive from fabric faster than the line transmission rate

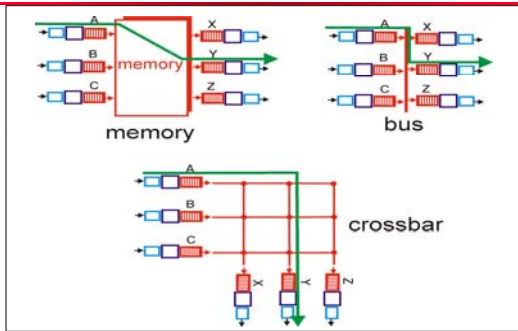
13

Router Processor

- Runs routing protocol and downloads forwarding table to forwarding engines
- Performs “slow” path processing
 - » ICMP error messages
 - » IP option processing
 - » Fragmentation
 - » Packets destined to router

14

Three Types of Switching Fabrics

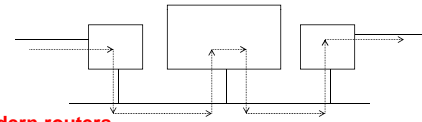


15

Switching Via a Memory

First generation routers → looked like PCs

- Packet copied by system's (single) CPU
- Speed limited by memory bandwidth (2 bus crossings per datagram)



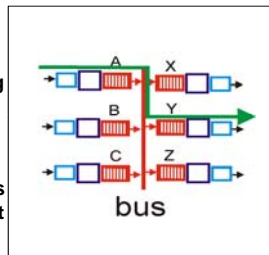
Modern routers

- Input port processor performs lookup, copy into memory
- Cisco Catalyst 8500

16

Switching Via a Bus

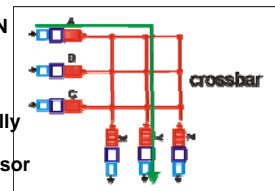
- Datagram from input port memory to output port memory via a shared bus
- **Bus contention:** switching speed limited by bus bandwidth
- 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)



17

Switching Via an Interconnection Network

- Overcome bus bandwidth limitations
- Crossbar provides full NxN interconnect
 - » Expensive
- Banyan networks & other interconnection nets initially developed to connect processors in multiprocessor
 - » Typically less capable than complete crossbar
- Cisco 12000: switches Gbps through the interconnection network



18

Buffering

- Suppose we have N inputs and M outputs
 - » Multiple packets for same output → output contention
 - » Switching fabric may force different inputs to wait → Switch contention
- Solution – buffer packets when/where needed: input, switch, or output
- What happens when these buffers fill up?
 - » Packets are **THROWN AWAY!!** This is where packet loss comes from

19

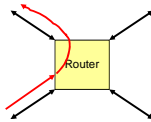
Outline

- ICMP
- How do Routers Works?
- **Routing**
- Distance vector

20

IP Forwarding versus Routing

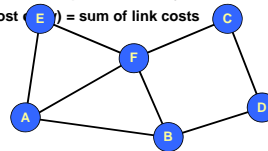
- The Story So Far...
 - » IP addresses are structure to reflect Internet structure
 - » IP packet headers carry these addresses
 - » When Packet Arrives at Router
 - Examine header to determine intended destination
 - Look up in table to determine next hop in path
 - Send packet out appropriate port
- How do we generate the forwarding table?



21

Graph Model

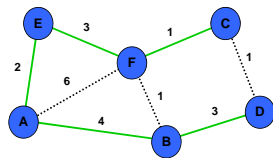
- Represent each router as node
- Direct link between routers represented by edge
 - » Symmetric links ⇒ undirected graph
- Edge “cost” $c(x,y)$ denotes measure of difficulty of using link
 - » delay, \$ cost, or congestion level
- Task
 - » Determine least cost path from every node to every other node
 - Path cost $C(x,y)$ = sum of link costs



22

Routes from Node A

Dest	Cost	Next Hop
A	0	A
B	4	B
C	6	E
D	7	B
E	2	E
F	5	E



- Properties
 - » Some set of shortest paths forms tree
 - Shortest path spanning tree
 - » Solution not unique
 - E.g., A-E-F-C-D also has cost 7

23

Ways to Compute Shortest Paths

- Centralized
 - » Collect graph structure in one place
 - » Use standard graph algorithm
 - » Disseminate routing tables
- Link-state
 - » Every node collects complete graph structure
 - » Each computes shortest paths from it
 - » Each generates own routing table
- Distance-vector
 - » No one has copy of graph
 - » Nodes construct their own tables iteratively
 - » Each sends information about its table to neighbors

24

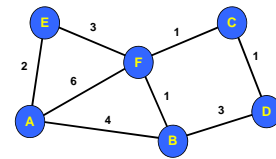
Outline

- ICMP
- How do Routers Works?
- Routing
- **Distance vector**

25

Distance-Vector Method

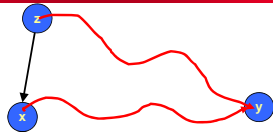
Initial Table for A		
Dest	Cost	Next Hop
A	0	A
B	4	B
C	∞	-
D	∞	-
E	2	E
F	6	F



- **Idea:** At any time, have cost/next hop of best known path to destination
 - » Use cost ∞ when no path is known
- **Initially**
 - » Only have entries for directly connected nodes

26

Distance-Vector Update



- **Update(x,y,z)**
 - $d \leftarrow c(x,z) + d(z,y)$ # Cost of path from x to y with first hop z
 - if $d < d(x,y)$
 - # Found better path
 - return d,z # Updated cost / next hop
 - else
 - return $d(x,y), \text{nexthop}(x,y)$ # Existing cost / next hop

27

Algorithm

- Bellman-Ford algorithm
- Repeat
 - For every node x
 - For every neighbor z
 - For every destination y
 - $d(x,y) \leftarrow \text{Update}(x,y,z)$
- Until converge

28

Start

Table for A					Table for B				
Dest	Cost	Next Hop	Dest	Cost	Next Hop				
A	0	A	A	4	A				
B	4	B	B	0	B				
C	∞	-	C	∞	-				
D	∞	-	D	3	D				
E	2	E	E	∞	-				
F	6	F	F	1	F				

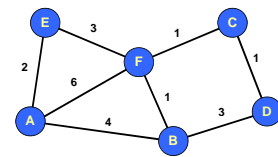


Table for C			Table for D			Table for E			Table for F		
Dest	Cost	Next Hop	Dest	Cost	Next Hop	Dest	Cost	Next Hop	Dest	Cost	Next Hop
A	∞	-	A	∞	-	A	2	A	A	6	A
B	∞	-	B	3	B	B	∞	-	B	1	B
C	0	C	C	1	C	C	∞	-	C	1	C
D	1	D	D	0	D	D	∞	-	D	∞	-
E	∞	-	E	∞	-	E	0	E	E	3	E
F	1	F	F	∞	-	F	3	F	F	0	F

29

Iteration #1

Table for A					Table for B				
Dest	Cost	Next Hop	Dest	Cost	Next Hop				
A	0	A	A	4	A				
B	4	B	B	0	B				
C	7	F	C	2	F				
D	7	B	D	3	D				
E	2	E	E	4	F				
F	5	E	F	1	F				

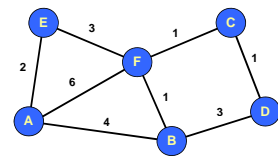


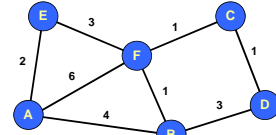
Table for C			Table for D			Table for E			Table for F		
Dest	Cost	Next Hop	Dest	Cost	Next Hop	Dest	Cost	Next Hop	Dest	Cost	Next Hop
A	7	F	A	7	B	A	2	A	A	5	B
B	2	F	B	3	B	B	4	F	B	1	B
C	0	C	C	1	C	C	4	F	C	1	C
D	1	D	D	0	D	D	∞	-	D	2	C
E	4	F	E	∞	-	E	0	E	E	3	E
F	1	F	F	2	C	F	3	F	F	0	F

30

Iteration #2

Table for A		Table for B	
Dst	Cost	Src	Cost
A	0	A	4
B	4	B	0
C	6	E	2
D	7	B	3
E	2	E	4
F	5	E	1

Table for C		Table for D		Table for E		Table for F	
Dst	Cost	Src	Cost	Dst	Cost	Src	Cost
A	6	F	7	A	2	A	5
B	2	F	3	B	4	F	1
C	0	C	1	C	4	F	1
D	1	D	0	D	5	F	2
E	4	F	5	E	0	E	3
F	1	F	2	C	3	F	0

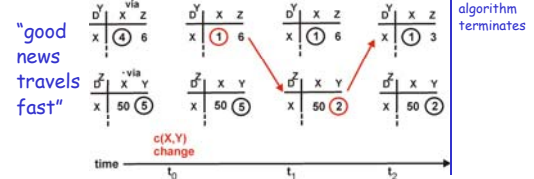
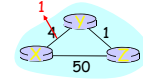


31

Distance Vector: Link Cost Changes

Link cost changes:

- Node detects local link cost change
- Updates distance table
- If cost change in least cost path, notify neighbors

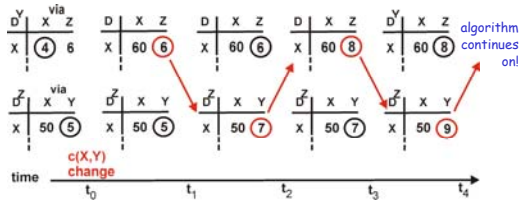
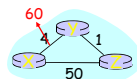


32

Distance Vector: Link Cost Changes

Link cost changes:

- Good news travels fast
- Bad news travels slow - "count to infinity" problem!

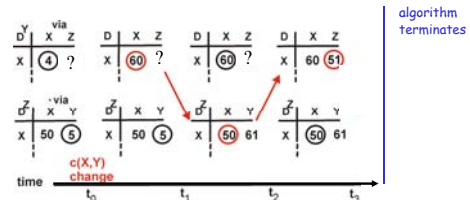
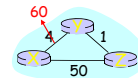


33

Distance Vector: Split Horizon

If Z routes through Y to get to X :

- No need for Z to tell Y about its route to X
- Z does not advertise its route to X back to Y

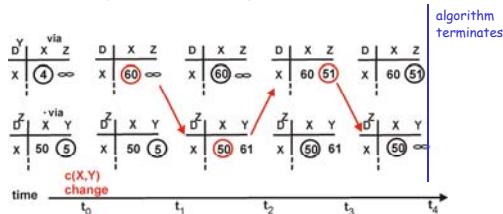
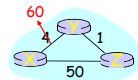


34

Distance Vector: Poison Reverse

If Z routes through Y to get to X :

- Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- Eliminates some possible timeouts with split horizon
- Will this completely solve count to infinity problem?



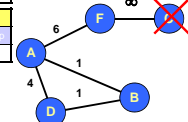
35

Poison Reverse Failures

Table for A		Table for B		Table for D		Table for F		
Dst	Cost	Src	Cost	Dst	Cost	Src	Cost	
C	7	F	8	A	C	9	B	C

Table for A	
Dst	Cost
C	*

Table for F	
Dst	Cost
C	*



- Iterations don't converge
- "Count to infinity"
- Solution
 - » Make "infinity" smaller
 - » What is upper bound on maximum path length?

Table for A	
Dst	Cost
C	19

Table for D	
Dst	Cost
B	

36

Routing Information Protocol (RIP)

- **Earliest IP routing protocol (1982 BSD)**
 - » Current standard is version 2 (RFC 1723)
- **Features**
 - » Every link has cost 1
 - » "Infinity" = 16
 - Limits to networks where everything reachable within 15 hops
- **Sending Updates**
 - » Every router listens for updates on UDP port 520
 - » RIP message can contain entries for up to 25 table entries

37

RIP Updates

- **Initial**
 - » When router first starts, asks for copy of table for every neighbor
 - » Uses it to iteratively generate own table
- **Periodic**
 - » Every 30 seconds, router sends copy of its table to each neighbor
 - » Neighbors use to iteratively update their tables
- **Triggered**
 - » When every entry changes, send copy of entry to neighbors
 - Except for one causing update (split horizon rule)
 - » Neighbors use to update their tables

38

RIP Staleness / Oscillation Control

- **Small Infinity**
 - » Count to infinity doesn't take very long
- **Route Timer**
 - » Every route has timeout limit of 180 seconds
 - Reached when haven't received update from next hop for 6 periods
 - » If not updated, set to infinity
 - » Soft-state refresh → important concept!!!
- **Behavior**
 - » When router or link fails, can take minutes to stabilize

39