# Topic

- How do we connect nodes with a <u>switch</u> instead of multiple access
  - Uses multiple links/wires
  - Basis of modern (switched) Ethernet



Switch

# Switched Ethernet

- Hosts are wired to Ethernet switches with twisted pair
  - Switch serves to connect the hosts
  - Wires usually run to a closet

Switch

Switch ports

Twisted pair

# What's in the box?

- Remember from protocol layers:

Hub, or repeater

| Physical | Physical |
|----------|----------|

Switch

| Link | Link |
|------|------|

All look like this:

Router

| Network | Network |
|---------|---------|
| Link    | Link    |

# Inside a Hub

- All ports are wired together; more convenient and reliable than a single shared wire

# Inside a Switch

- Uses frame addresses to connect input port to the right output port; multiple frames may be switched in parallel

# Inside a Switch (2)

- Port may be used for both input and output (full-duplex)
  - Just send, no multiple access protocol

# Inside a Switch (3)

- Need buffers for multiple inputs to send to one output

Input

Output

Input Buffer

Fabric

Output Buffer

# Inside a Switch (4)

- Sustained overload will fill buffer and lead to frame loss

# Advantages of Switches

- Switches and hubs have replaced the shared cable of classic Ethernet
  - Convenient to run wires to one location
  - More reliable; wire cut is not a single point of failure that is hard to find

- Switches offer scalable performance
  - E.g., 100 Mbps per port instead of 100 Mbps for all nodes of shared cable / hub

# Switch Forwarding

- Switch needs to find the right output port for the destination address in the Ethernet frame. How?
  - Want to let hosts be moved around readily; don't look at IP

# Backward Learning

- Switch forwards frames with a port/address table as follows:
  1. To fill the table, it looks at the source address of input frames
  2. To forward, it sends to the port, or else broadcasts to all ports

# Backward Learning (2)

- 1: A sends to D



| Address | Port |
|---------|------|
| A       |      |
| B       |      |
| C       |      |
| D       |      |

# Backward Learning (3)

- 2: D sends to A



| Address | Port |
|---------|------|
| A | 1 |
| B | |
| C | |
| D | |

# Backward Learning (4)

- 3: A sends to D



| Address | Port |
|---------|------|
| A | 1 |
| B | |
| C | |
| D | 4 |

# Backward Learning (5)

- 3: A sends to D



| Address | Port |
|---------|------|
| A | 1 |
| B | |
| C | |
| D | 4 |

# Learning with Multiple Switches

- Just works with multiple switches and a mix of hubs *assuming no loo*ps, e.g., A sends to D then D sends to A

# Learning with Multiple Switches (2)

- Just works with multiple switches and a mix of hubs *assuming no loo*ps, e.g., A sends to D then D sends to A

# Learning with Multiple Switches (3)

- Just works with multiple switches and a mix of hubs *assuming no loo*ps, e.g., A sends to D then D sends to A

# Topic

- How can we connect switches in any topology so they just work
  - This is part 2 of switched Ethernet



Loops – yikes!

# Problem – Forwarding Loops

- May have a loop in the topology
  - Redundancy in case of failures
  - Or a simple mistake

- Want LAN switches to "just work"
  - Plug-and-play, no changes to hosts
  - But loops cause a problem …

Redundant Links

# Forwarding Loops (2)

- Suppose the network is started and A sends to F. What happens?



Left / Right

# Forwarding Loops (3)

- Suppose the network is started and A sends to F. What happens?
  - A → C → B, D-left, D-right
  - D-left → C-right, E, F
  - D-right → C-left, E, F
  - C-right → D-left, A, B
  - C-left → D-right, A, B
  - D-left → …
  - D-right → …



Left / Right

# Spanning Tree Solution

- Switches collectively find a <u>spanning tree</u> for the topology

  - A subset of links that is a tree (no loops) and reaches all switches

  - They switches forward as normal on the spanning tree

  - Broadcasts will go up to the root of the tree and down all the branches

# Spanning Tree (2)

Topology                    One ST                    Another ST

# Spanning Tree (3)

# Spanning Tree Algorithm

- Rules of the distributed game:
  - All switches run the same algorithm
  - They start with no information
  - Operate in parallel and send messages
  - Always search for the best solution

- Ensures a highly robust solution
  - Any topology, with no configuration
  - Adapts to link/switch failures, …

# Radia Perlman (1952–)

- Key early work on routing protocols
  - Routing in the ARPANET
  - Spanning Tree for switches (next)
  - Link-state routing (later)

- Now focused on network security

# Spanning Tree Algorithm (2)

- Outline:

  1. Elect a root node of the tree (switch with the lowest address)

  2. Grow tree as shortest distances from the root (using lowest address to break distance ties)

  3. Turn off ports for forwarding if they aren't on the spanning tree

# Spanning Tree Algorithm (3)

- Details:
  - Each switch initially believes it is the root of the tree
  - Each switch sends periodic updates to neighbors with:
    - Its address, address of the root, and distance (in hops) to root
  - Switches favors ports with shorter distances to lowest root
    - Uses lowest address as a tie for distances

Hi, I'm <u>C</u>, the root is <u>A</u>, it's <u>2</u> hops away    or (C, A, 2)

# Spanning Tree Example

- 1st round, sending:
  - A sends (A, A, 0) to say it is root
  - B, C, D, E, and F do likewise
- 1st round, receiving:
  - A still thinks is it (A, A, 0)
  - B still thinks (B, B, 0)
  - C updates to (C, A, 1)
  - D updates to (D, C, 1)
  - E updates to (E, A, 1)
  - F updates to (F, B, 1)

A,A,0    B,B,0
   C,C,0
   D,D,0
E,E,0    F,F,0

# Spanning Tree Example (2)

- 2<sup>nd</sup> round, sending
  - Nodes send their updated state
- 2<sup>nd</sup> round receiving:
  - A remains (A, A, 0)
  - B updates to (B, A, 2) via C
  - C remains (C, A, 1)
  - D updates to (D, A, 2) via C
  - E remains (E, A, 1)
  - F remains (F, B, 1)

# Spanning Tree Example (3)

- 3$^{rd}$ round, sending
  - Nodes send their updated state
- 3$^{rd}$ round receiving:
  - A remains (A, A, 0)
  - B remains (B, A, 2) via C
  - C remains (C, A, 1)
  - D remains (D, A, 2) via C-left
  - E remains (E, A, 1)
  - F updates to (F, A, 3) via B

A,A,0   B,A,2

C,A,1

D,A,2

E,A,1   F,B,1

# Spanning Tree Example (4)

- 4$^{th}$ round
  - Steady-state has been reached
  - Nodes turn off forwarding that is not on the spanning tree

- Algorithm continues to run
  - Adapts by timing out information
  - E.g., if A fails, other nodes forget it, and B will become the new root

# Spanning Tree Example (5)

- Forwarding proceeds as usual on the ST
- Initially D sends to F:


- And F sends back to D:

# Spanning Tree Example (6)

- Forwarding proceeds as usual on the ST
- Initially D sends to F:
  - D → C-left
  - C → A, B
  - A → E
  - B → F

- And F sends back to D:
  - F → B
  - B → C
  - C → D

  (hm, not such a great route)

# Where we are in the Course

- Starting the Network Layer!
  - Builds on the link layer. <u>Routers</u> send <u>packets</u> over multiple networks

| Application |
|:---:|
| Transport |
| Network |
| Link |
| Physical |

# Why do we need a Network layer?

- We can already build networks with links and switches and send frames between hosts ...

# Shortcomings of Switches

1. Don't scale to large networks
   – Blow up of routing table, broadcast



Table for all destinations in the world!

Broadcast new destinations to the whole world!

# Shortcomings of Switches (2)

2. Don't work across more than one link layer technology

   – Hosts on Ethernet + 3G + 802.11 …

# Shortcomings of Switches (3)

3. Don't give much traffic control
   – Want to plan routes / bandwidth



That was lame.

# Topic

- How do routers <u>forward</u> packets?
  - We'll look at how IP does it
  - (We'll cover routing later)

# Recap

- We want the network layer to:
  - Scale to large networks
    - Using addresses with hierarchy ⎫ This lecture
  - Support diverse technologies
    - Internetworking with IP ⎫ More later
  - Use link bandwidth well
    - Lowest-cost routing ⎫ Next time

# IP Addresses

- IPv4 uses 32-bit addresses
  - Later we'll see IPv6, which uses 128-bit addresses
- Written in "dotted quad" notation
  - Four 8-bit numbers separated by dots

```
    8 bits        8 bits        8 bits        8 bits

aaaaaaaabbbbbbbbccccccccdddddddd   ↔  A.B.C.D
00010010000111110000000000000001   ↔
```

# IP Prefixes

- Addresses are allocated in blocks called <u>prefixes</u>
  - Addresses in an L-bit prefix have the same top L bits
  - There are $2^{32-L}$ addresses aligned on $2^{32-L}$ boundary

# IP Prefixes (2)

- Written in "IP address/length" notation
  - Address is lowest address in the prefix, length is prefix bits
  - E.g., 128.13.0.0/16 is 128.13.0.0 to 128.13.255.255
  - So a /24 ("slash 24") is 256 addresses, and a /32 is one address

```
00010010 00011111 00000000 xxxxxxxx ↔
```

↔ 128.13.0.0/16

# Classful IP Addressing

- Originally, IP addresses came in fixed size blocks with the class/size encoded in the high-order bits
  - They still do, but the classes are now ignored

| 0 | 8 | 16 | 24 | 32 bits |
|---|---|----|----|---------|

| 0 | | | | Class A, $2^{24}$ addresses |
| 10 | | | | Class B, $2^{16}$ addresses |
| 110 | | | | Class C, $2^{8}$ addresses |

Network portion          Host portion

# IP Forwarding

- All addresses on one network belong to the same prefix
- Node uses a table that lists the next hop for prefixes

| Prefix | Next Hop |
|---|---|
| 192.24.0.0/19 | D |
| 192.24.12.0/22 | B |

# Longest Matching Prefix

- Prefixes in the table might overlap!
  - Combines hierarchy with flexibility

- <u>Longest matching prefix</u> forwarding rule:
  - For each packet, find the longest prefix that contains the destination address, i.e., the most specific entry
  - Forward the packet to the next hop router for that prefix

# Longest Matching Prefix (2)

| Prefix | Next Hop |
|---|---|
| 192.24.0.0/19 | D |
| 192.24.12.0/22 | B |

192.24.6.0 →
192.24.14.32 →
192.24.54.0 →

192.24.63.255

/19

192.24.15.255

192.24.12.0

192.24.0.0

More specific

/22

IP address

# Host/Router Distinction

- In the Internet:
  - Routers do the routing, know which way to all destinations
  - Hosts send remote traffic (out of prefix) to nearest router

Not for my network?
Send it to the router

It's my job to know
which way to go …

# Host Forwarding Table

- Give using longest matching prefix
  - 0.0.0.0/0 is a default route that catches all IP addresses

| Prefix | Next Hop |
|---|---|
| My network prefix | Send to that IP |
| 0.0.0.0/0 | Send to my router |

# Flexibility of Longest Matching Prefix

- Can provide default behavior, with less specifics
  - To send traffic going outside an organization to a border router

- Can special case behavior, with more specifics
  - For performance, economics, security, …

# Performance of Longest Matching Prefix

- Uses hierarchy for a compact table
  - Relies on use of large prefixes

- Lookup more complex than table
  - Used to be a concern for fast routers
  - Not an issue in practice these days

# Topic

- Filling in the gaps we need to make for IP forwarding work in practice
  - Getting IP addresses (DHCP) »
  - Mapping IP to link addresses (ARP) »

What's my IP?

What link layer address do I use?

# Getting IP Addresses

- Problem:
  - A node wakes up for the first time …
  - What is its IP address? What's the IP address of its router? Etc.
  - At least Ethernet address is on NIC

Hey, where am I?

# Getting IP Addresses (2)

1. Manual configuration (old days)
   - Can't be factory set, depends on use

2. A protocol for automatically configuring addresses (DHCP) **»**
   - Shifts burden from users to IT folk

What's my IP?

Use A.B.C.D

# DHCP

- DHCP (Dynamic Host Configuration Protocol), from 1993, widely used

- It leases IP address to nodes
- Provides other parameters too
  - Network prefix
  - Address of local router
  - DNS server, time server, etc.

# DHCP Protocol Stack

- DHCP is a client-server application
  - Uses UDP ports 67, 68

| DHCP |
| :---: |
| UDP |
| IP |
| Ethernet |

# DHCP Addressing

- Bootstrap issue:
  - How does node send a message to DHCP server before it is configured?

- Answer:
  - Node sends <u>broadcast</u> messages that delivered to all nodes on the network
  - <u>Broadcast address</u> is all 1s
  - IP (32 bit): 255.255.255.255
  - Ethernet (48 bit): ff:ff:ff:ff:ff:ff

# DHCP Messages



Client     Server

One link

# DHCP Messages (2)



Client          Server

DISCOVER

Broadcast

OFFER

REQUEST

ACK

# DHCP Messages (3)

- To renew an existing lease, an abbreviated sequence is used:
  - REQUEST, followed by ACK

- Protocol also supports replicated servers for reliability

# Sending an IP Packet

- Problem:
  - A node needs Link layer addresses to send a frame over the local link
  - How does it get the destination link address from a destination IP address?

Uh oh …

My IP is 1.2.3.4

# ARP (Address Resolution Protocol)

- Node uses to map a local IP address to its Link layer addresses

Link layer

| Source Ethernet | Dest. Ethernet | Source IP | Dest. IP | Payload … |
|---|---|---|---|---|

From NIC

From ARP

From DHCP

# ARP Protocol Stack

- ARP sits right on top of link layer
  - No servers, just asks node with target IP to identify itself
  - Uses broadcast to reach all nodes

| ARP |
| --- |
| Ethernet |

# ARP Messages

Node                Target

One link

# ARP Messages (2)



Node

Target

REQUEST
Broadcast
Who has IP 1.2.3.4?

REPLY
I do at 1:2:3:4:5:6

# Discovery Protocols

- Help nodes find each other
  - There are more of them!
    - E.g., zeroconf, Bonjour

- Often involve broadcast
  - Since nodes aren't introduced
  - Very handy glue

# Topic

- IP version 6, the future of IPv4 that is now (still) being deployed

Why do I want IPv6 again?

# Internet Growth

- At least a billion Internet hosts and growing …

- And we're using 32-bit addresses!



Internet Domain Survey Host Count

Source: Internet Systems Consortium (www.isc.org)

# The End of New IPv4 Addresses

- Now running on leftover blocks held by the regional registries; much tighter allocation policies

Exhausted on 4/11 and 9/12!

| ARIN (US, Canada) |
| APNIC (Asia Pacific) |
| RIPE (Europe) |
| LACNIC (Latin Amer.) |
| AfriNIC (Africa) |

IANA (All IPs)

Exhausted on 2/11!

ISPs

Companies

End of the world ? 12/21/12?

# IP Version 6 to the Rescue

- Effort started by the IETF in 1994
  - Much larger addresses (128 bits)
  - Many sundry improvements

- Became an IETF standard in 1998
  - Nothing much happened for a decade
  - Hampered by deployment issues, and a lack of adoption incentives
  - Big push ~2011 as exhaustion looms

# IPv6 Deployment

Percentage of users accessing Google via IPv6

Time for growth!



Source: Google IPv6 Statistics, 30/1/13

# IPv6

- Features large addresses
  - 128 bits, most of header
- New notation
  - 8 groups of 4 hex digits (16 bits)
  - Omit leading zeros, groups of zeros

32 bits

| Version | Diff. Serv. | Flow label | |
|---|---|---|---|
| Payload length | | Next header | Hop limit |

Source address
(16 bytes)

Destination address
(16 bytes)

Ex: 2001:0db8:0000:0000:0000:ff00:0042:8329
→

# IPv6 (2)

- Lots of other, smaller changes
  - Streamlined header processing
  - Flow label to group of packets
  - Better fit with "advanced" features (mobility, multicasting, security)

32 bits

| Version | Diff. Serv. | Flow label | | |
| Payload length | | | Next header | Hop limit |

Source address
(16 bytes)

Destination address
(16 bytes)

# IPv6 Transition

- The Big Problem:
  - How to deploy IPv6?
  - Fundamentally incompatible with IPv4

- Dozens of approaches proposed
  - Dual stack (speak IPv4 and IPv6)
  - Translators (convert packets)
  - Tunnels (carry IPv6 over IPv4) »

# Tunneling

- Native IPv6 islands connected via IPv4
  - Tunnel carries IPv6 packets across IPv4 network

# Tunneling (2)

- Tunnel acts as a single link across IPv4 network

# Tunneling (3)

- Tunnel acts as a single link across IPv4 network
  - Difficulty is to set up tunnel endpoints and routing

# Topic

- What is NAT (Network Address Translation)? How does it work?
  - NAT is widely used at the edges of the network, e.g., homes

# Layering Review

- Remember how layering is meant to work?
  - "Routers don't look beyond the IP header." Well …

# Middleboxes

- Sit "inside the network" but perform "more than IP" processing on packets to add new functionality
  - NAT box, Firewall / Intrusion Detection System

Middlebox

App

| TCP |
| IP |
| 802.11 |

App / TCP

| IP | IP |
| 802.11 | Ethernet |

| IP | IP |
| 802.11 | Ethernet |

App

| TCP |
| IP |
| 802.11 |

# Middleboxes (2)

- Advantages
  - A possible rapid deployment path when there is no other option
  - Control over many hosts (IT)

- Disadvantages
  - Breaking layering interferes with connectivity; strange side effects
  - Poor vantage point for many tasks

# NAT (Network Address Translation) Box

- NAT box connects an internal network to an external network
  - Many internal hosts are connected using few external addresses
  - Middlebox that "translates addresses"

- Motivated by IP address scarcity
  - Controversial at first, now accepted

# NAT (2)

- Common scenario:
  - Home computers use "private" IP addresses
  - NAT (in AP/firewall) connects home to ISP using a single external IP address

Unmodified computers at home

Looks like one computer outside

ISP

NAT box

# How NAT Works

- Keeps an internal/external table
  - Typically uses IP address + TCP port
  - This is address and port translation

What host thinks     What ISP thinks

| Internal IP:port | External IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |
| 192.168.1.13 : 1234 | 44.25.80.3 : 1501 |
| 192.168.2.20 : 1234 | 44.25.80.3 : 1502 |

- Need ports to make mapping 1-1 since there are fewer external IPs

# How NAT Works (2)

- Internal → External:
  - Look up and rewrite Source IP/port

Internal source

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External destination
IP=X, port=Y

Src =

Dst =

NAT box

Src =

Dst =

# How NAT Works (3)

- External → Internal
  - Look up and rewrite Destination IP/port

Internal
destination

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External
source
IP=X, port=Y

NAT box

Src =

Dst =

Src =

Dst =

# How NAT Works (4)

- Need to enter translations in the table for it to work
  - Create external name when host makes a TCP connection

Internal source

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | |

External destination
IP=X, port=Y



Src =

Dst =

NAT box

Src =

Dst =

# NAT Downsides

- Connectivity has been broken!
  - Can only send incoming packets after an outgoing connection is set up
  - Difficult to run servers or peer-to-peer apps (Skype) at home

- Doesn't work so well when there are no connections (UDP apps)

- Breaks apps that unwisely expose their IP addresses (FTP)

# NAT Upsides

- Relieves much IP address pressure
  - Many home hosts behind NATs
- Easy to deploy
  - Rapidly, and by you alone
- Useful functionality
  - Firewall, helps with privacy

- Kinks will get worked out eventually
  - "NAT Traversal" for incoming traffic

# Topic

- Defining "best" paths with link costs
  - These are <u>shortest path</u> routes



Best?

# What are "Best" paths anyhow?

- Many possibilities:
  - Latency, avoid circuitous paths
  - Bandwidth, avoid slow links
  - Money, avoid expensive links
  - Hops, to reduce switching

- But only consider topology
  - Ignore workload, e.g., hotspots

# Shortest Paths

We'll approximate "best" by a cost function that captures the factors

– Often call lowest "shortest"

1. Assign each link a cost (distance)

2. Define best path between each pair of nodes as the path that has the lowest total cost (or is shortest)

3. Pick randomly to any break ties

# Shortest Paths (2)

- Find the shortest path A → E

- All links are bidirectional, with equal costs in each direction
  - Can extend model to unequal costs if needed

# Shortest Paths (3)

- ABCE is a shortest path
- dist(ABCE) = 4 + 2 + 1 = 7

- This is less than:
  - dist(ABE) = 8
  - dist(ABFE) = 9
  - dist(AE) = 10
  - dist(ABCDE) = 10

# Shortest Paths (4)

- Optimality property:
  - Subpaths of shortest paths are also shortest paths

- ABCE is a shortest path
  - →So are ABC, AB, BCE, BC, CE

# Sink Trees

- Sink tree for a destination is the union of all shortest paths towards the destination
  - Similarly source tree

- Find the sink tree for E

# Sink Trees (2)

- Implications:
  - Only need to use destination to follow shortest paths
  - Each node only need to send to the next hop

- Forwarding table at a node
  - Lists next hop for each destination
  - Routing table may know more

# Topic

- ## How to compute shortest paths given the network topology

  - ### With Dijkstra's algorithm



Source tree for E

# Edsger W. Dijkstra (1930-2002)

- Famous computer scientist
  - Programming languages
  - Distributed algorithms
  - Program verification

- Dijkstra's algorithm, 1969
  - Single-source shortest paths, given network with non-negative link costs

By Hamilton Richards, CC-BY-SA-3.0, via Wikimedia Commons

# Dijkstra's Algorithm

Algorithm:

- Mark all nodes tentative, set distances from source to 0 (zero) for source, and ∞ (infinity) for all other nodes

- While tentative nodes remain:
  - Extract N, a node with lowest distance
  - Add link to N to the shortest path tree
  - Relax the distances of neighbors of N by lowering any better distance estimates

# Dijkstra's Algorithm (2)

- Initialization



We'll compute
shortest paths
from A

# Dijkstra's Algorithm (3)

- Relax around A

# Dijkstra's Algorithm (4)

- Relax around B



Distance fell!

F **7**

4

2

**7**
G

3

E **8**

3

10

2

0

4

∞ D

A    4    B **4**

1

2

H

∞    3    C **6**

2

# Dijkstra's Algorithm (5)

- Relax around C

# Dijkstra's Algorithm (6)

- Relax around G (say)



Didn't fall …

# Dijkstra's Algorithm (7)

- Relax around F (say)



Relax has no effect

# Dijkstra's Algorithm (8)

- Relax around E

# Dijkstra's Algorithm (9)

- Relax around D

# Dijkstra's Algorithm (10)

- Finally, H … done

# Dijkstra Comments

- Finds shortest paths in order of increasing distance from source
  - Leverages optimality property

- Runtime depends on efficiency of extracting min-cost node
  - Superlinear in network size (grows fast)

- Gives complete source/sink tree
  - More than needed for forwarding!
  - But requires complete topology

# Topic

- How to compute shortest paths in a distributed network
  - The Distance Vector (DV) approach

# Distance Vector Routing

- Simple, early routing approach
  - Used in ARPANET, and RIP

- One of two main approaches to routing
  - Distributed version of Bellman-Ford
  - Works, but very slow convergence after some failures

- Link-state algorithms are now typically used in practice
  - More involved, better behavior

# Distance Vector Setting

Each node computes its forwarding table
in a distributed setting:

1. Nodes know only the cost to their
   neighbors; not the topology

2. Nodes can talk only to their neighbors
   using messages

3. All nodes run the same algorithm
   concurrently

4. Nodes and links may fail, messages
   may be lost

# Distance Vector Algorithm

Each node maintains a vector of distances (and next hops) to all destinations

1. Initialize vector with 0 (zero) cost to self, ∞ (infinity) to other destinations
2. Periodically send vector to neighbors
3. Update vector for each destination by selecting the shortest distance heard, after adding cost of neighbor link
   - Use the best neighbor for forwarding

# Distance Vector (2)

- Consider from the point of view of node A
  - Can only talk to nodes B and E

| To | Cost |
|----|------|
| A | 0 |
| B | ∞ |
| C | ∞ |
| D | ∞ |
| E | ∞ |
| F | ∞ |
| G | ∞ |
| H | ∞ |

Initial vector →

# Distance Vector (3)

- First exchange with B, E; learn best 1-hop routes

| To | B says | E says |
|---|---|---|
| A | ∞ | ∞ |
| B | 0 | ∞ |
| C | ∞ | ∞ |
| D | ∞ | ∞ |
| E | ∞ | 0 |
| F | ∞ | ∞ |
| G | ∞ | ∞ |
| H | ∞ | ∞ |

→

| B +4 | E +10 |
|---|---|
| ∞ | ∞ |
| 4 | ∞ |
| ∞ | ∞ |
| ∞ | ∞ |
| ∞ | 10 |
| ∞ | ∞ |
| ∞ | ∞ |
| ∞ | ∞ |

→

| A's Cost | A's Next |
|---|---|
| 0 | -- |
| 4 | B |
| ∞ | -- |
| ∞ | -- |
| 10 | E |
| ∞ | -- |
| ∞ | -- |
| ∞ | -- |

Learned better route

# Distance Vector (4)

- Second exchange; learn best 2-hop routes

| To | B says | E says |
|----|--------|--------|
| A | 4 | 10 |
| B | 0 | 4 |
| C | 2 | 1 |
| D | ∞ | 2 |
| E | 4 | 0 |
| F | 3 | 2 |
| G | 3 | ∞ |
| H | ∞ | ∞ |

→

| B +4 | E +10 |
|------|-------|
| 8 | 20 |
| 4 | 14 |
| 6 | 11 |
| ∞ | 12 |
| 8 | 10 |
| 7 | 12 |
| 7 | ∞ |
| ∞ | ∞ |

→

| A's Cost | A's Next |
|----------|----------|
| 0 | -- |
| 4 | B |
| 6 | B |
| 12 | E |
| 8 | B |
| 7 | B |
| 7 | B |
| ∞ | -- |

# Distance Vector (4)

- Third exchange; learn best 3-hop routes

| To | B says | E says |
|---|---|---|
| A | 4 | 8 |
| B | 0 | 3 |
| C | 2 | 1 |
| D | 4 | 2 |
| E | 3 | 0 |
| F | 3 | 2 |
| G | 3 | 6 |
| H | 5 | 4 |

→

| B +4 | E +10 |
|---|---|
| 8 | 18 |
| 4 | 13 |
| 6 | 11 |
| 8 | 12 |
| 7 | 10 |
| 7 | 12 |
| 7 | 16 |
| 9 | 14 |

→

| A's Cost | A's Next |
|---|---|
| 0 | -- |
| 4 | B |
| 6 | B |
| 8 | B |
| 7 | B |
| 7 | B |
| 7 | B |
| 9 | B |

# Distance Vector (5)

- Subsequent exchanges; converged

| To | B says | E says |
|----|--------|--------|
| A  | 4      | 7      |
| B  | 0      | 3      |
| C  | 2      | 1      |
| D  | 4      | 2      |
| E  | 3      | 0      |
| F  | 3      | 2      |
| G  | 3      | 6      |
| H  | 5      | 4      |

→

| B +4 | E +10 |
|------|-------|
| 8    | 17    |
| 4    | 13    |
| 6    | 11    |
| 8    | 12    |
| 7    | 10    |
| 7    | 12    |
| 7    | 16    |
| 9    | 14    |

→

| A's Cost | A's Next |
|----------|----------|
| 0        | --       |
| 4        | B        |
| 6        | B        |
| 8        | B        |
| 8        | B        |
| 7        | B        |
| 7        | B        |
| 9        | B        |

# Topic

- How to compute shortest paths in a distributed network
  - The Link-State (LS) approach

# Link-State Routing

- One of two approaches to routing
  - Trades more computation than distance vector for better dynamics

- Widely used in practice
  - Used in Internet/ARPANET from 1979
  - Modern networks use OSPF and IS-IS

# Link-State Setting

Nodes compute their forwarding table in the same distributed setting as for distance vector:

1. Nodes know only the cost to their neighbors; not the topology
2. Nodes can talk only to their neighbors using messages
3. All nodes run the same algorithm concurrently
4. Nodes/links may fail, messages may be lost

# Link-State Algorithm

Proceeds in two phases:

1. Nodes <u>flood</u> topology in the form of link state packets

   – Each node learns full topology

2. Each node computes its own forwarding table

   – By running Dijkstra (or equivalent)

# Phase 1: Topology Dissemination

- Each node floods <u>link state packet</u> (LSP) that describes their portion of the topology

Node E's LSP flooded to A, B, C, D, and F

| Seq. # | |
|---|---|
| A | 10 |
| B | 4 |
| C | 1 |
| D | 2 |
| F | 2 |

# Phase 2: Route Computation

- Each node has full topology
  - By combining all LSPs

- Each node simply runs Dijkstra
  - Some replicated computation, but finds required routes directly
  - Compile forwarding table from sink/source tree
  - That's it folks!

# Forwarding Table

## Source Tree for E (from Dijkstra)



## E's Forwarding Table

| To | Next |
|----|------|
| A  | C    |
| B  | C    |
| C  | C    |
| D  | D    |
| E  | --   |
| F  | F    |
| G  | F    |
| H  | C    |

# Handling Changes

- On change, flood updated LSPs, and re-compute routes
  - E.g., nodes adjacent to failed link or node initiate

B's LSP

| Seq. # | |
|---|---|
| A | 4 |
| C | 2 |
| E | 4 |
| F | 3 |
| G | ∞ |

F's LSP

| Seq. # | |
|---|---|
| B | 3 |
| E | 2 |
| G | ∞ |

Failure!

G XXXX

# Handling Changes (2)

- Link failure
  - Both nodes notice, send updated LSPs
  - Link is removed from topology

- Node failure
  - All neighbors notice a link has failed
  - Failed node can't update its own LSP
  - But it is OK: all links to node removed

# Handling Changes (3)

- Addition of a link or node
  - Add LSP of new node to topology
  - Old LSPs are updated with new link

- Additions are the easy case ...

# Link-State Complications

- Things that can go wrong:
  - Seq. number reaches max, or is corrupted
  - Node crashes and loses seq. number
  - Network partitions then heals
- Strategy:
  - Include age on LSPs and forget old information that is not refreshed

- Much of the complexity is due to handling corner cases (as usual!)

# DV/LS Comparison

| Goal | Distance Vector | Link-State |
|---|---|---|
| Correctness | Distributed Bellman-Ford | Replicated Dijkstra |
| Efficient paths | Approx. with shortest paths | Approx. with shortest paths |
| Fair paths | Approx. with shortest paths | Approx. with shortest paths |
| Fast convergence | Slow – many exchanges | Fast – flood and compute |
| Scalability | Excellent – storage/compute | Moderate – storage/compute |