<u>Chapter 4: Network Layer,</u> <u>partb</u>

The slides are adaptations of the slides available by the main textbook authors, Kurose&Ross

Interplay between routing, forwarding



Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP

Graph abstraction



Graph: G = (N,E) N = set of routers = { u, v, w, x, y, z } E = set of links ={ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) }

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

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

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

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

Routing algorithm: algorithm that finds least-sost-upgth 4-4

Routing Algorithm classification

Global or decentralized information?

Global:

- r all routers have complete topology, link cost info
- r "link state" algorithms

Decentralized:

- r router knows physicallyconnected neighbors, link costs to neighbors
- r iterative process of computation, exchange of info with neighbors
- r "distance vector" algorithms

Static or dynamic? Static:

r routes don't change (or do slowly over time)

Dynamic:

- r routes change
 - m periodic update
 - m in response to link cost changes

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing

<u>A Link-State Routing Algorithm</u>

Dijkstra's algorithm

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

Notation:

- r C(x,y): link cost from node x to y; = ∞ if not direct neighbors
- r D(v): current value of cost of path from source to dest. v
- r p(v): predecessor node along path from source to v
- r N': set of nodes whose least cost path definitively known

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
$$D(v) = \infty$$

7

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':
- 12 D(v) = min(D(v), D(w) + c(w,v))
- 13 /* new cost to v is either old cost to v or known
- 14 shortest path cost to w plus cost from w to v */
- 15 until all nodes in N'

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	2,u	5,u	1,u	∞	∞
1	ux 🔶	2,u	4,x		2,x	∞
2	UXY•	<u>2,u</u>	З,у			4,y
3	uxyv 🗸					4,y
4	uxyvw 🔶					4,y
5						



Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

destination	link
V	(u,v)
×	(u,x)
У	(u,x)
W	(u,x)
Z	(u,x)

Network Layer 4-10

Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

- r each iteration: need to check all nodes, w, not in N
- r n(n+1)/2 comparisons: O(n²)
- r more efficient implementations possible: O(nlogn)

Oscillations possible:

r e.g., link cost = amount of carried traffic



Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

r 4.5 Routing algorithms

 m Link state
 m Distance Vector
 m Hierarchical routing

 r 4.6 Routing in the

 Internet
 m RIP
 m OSPF
 m BGP

Distance Vector Algorithm

Bellman-Ford Equation

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

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

Distance Vector Algorithm

- r $D_x(y)$ = estimate of least cost from x to y
- r Node x maintains distance vector $D_x = [D_x(y): y \in N]$
- r Node x also needs to know its neighbors' distance vectors

m For each neighbor v, x knows $D_v = [D_v(y): y \in N]$

Distance vector algorithm (4)

<u>Basic idea:</u>

- r From time-to-time, each node sends its own distance vector estimate to neighbors
- r Asynchronous
- r When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation:
 D_(y) ← min_{c(x,v) + D_(y)} for each node y ∈ N
- r Under minor, natural conditions, the estimate $D_x(y)$ converges to the actual least cost $d_x(y)$

Distance Vector Algorithm (5)

Iterative, asynchronous: each local iteration caused by:

- r local link cost change
- r DV update message from neighbor

Distributed:

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

Each node:







Network Layer 4-19

Distance Vector: link cost changes

Link cost changes:

- r node detects local link cost change
- r updates routing info, recalculates distance vector



r if DV changes, notify neighbors

"good news travels fast" At time t_0 , y detects the link-cost change, updates its DV, and informs its neighbors.

At time t_1 , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.

At time t_2 , y receives z's update and updates its distance table. y's least costs do not change and hence y does *not* send any message to z.

Distance Vector: link cost changes

Link cost changes:

- r node detects local link cost change
- r updates distance table
- r if cost change in least cost path, notify neighbors





algorithm terminates

Distance Vector: link cost changes

Link cost changes:

- r good news travels fast
- r bad news travels slow (watch: loops!) - "count to infinity" problem!



60

^{4:} Network Layer 4a-22

<u>Distance Vector:count to infinity</u> problem: way out?

Poisoned reverse:

- r 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)



<u>Comparison of LS and DV algorithms</u>

Message complexity

- r <u>LS:</u> with n nodes, E links, O(nE) msgs sent
- r <u>DV:</u> exchange between neighbors only

Speed of Convergence

- r LS: O(n²) algorithm
 m may have oscillations
- r <u>DV</u>: convergence time varies
 - m may be routing loopsm count-to-infinity problem

Robustness: what happens if router malfunctions?

<u>LS:</u>

- m node can advertise incorrect *link* cost
- m each node computes only its own table

DV:

- m DV node can advertise incorrect *path* cost
- m each node's table used by others
 - error propagates thru network

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing

Hierarchical Routing

Recall:

- r all routers identical
- r network "flat"
- ... not true in practice

scale: with 200 million destinations:

- r can't store all dest's in routing tables!
- r routing table exchange would swamp links!

administrative autonomy

- r internet = network of networks
- r each network admin may want to control routing in its own network

<u>Hierarchical Routing:</u> Interconnected ASes

Intra-AS

Routing

algorithm

Forwarding

table

<u>Gateway router</u>

20

AS2

20

AS1

Inter-AS

Routing

algorithm

Direct link to router in another AS

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

Network Layer 4-27

regions, "autonomous systems" (AS) routers in same AS

AS3

aggregate routers into

run same routing protocol

3b

- m "intra-AS" routing protocol
- routers in different
 AS can run different
 intra-AS routing
 protocol

Inter-AS tasks

- r suppose router in AS1 receives datagram destined outside of AS1:
 - m router should forward packet to gateway router, but which one?

<u>AS1 must:</u>

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

Job of inter-AS routing!



Example 1: Setting forwarding table in router 1d

- suppose AS1 learns (via inter-AS protocol) that subnet
 reachable via AS3 (gateway 1c) but not via AS2.
- r inter-AS protocol propagates reachability info to all internal routers.
- r router 1d determines from intra-AS routing info that its interface I is on the least cost path to 1c.

m installs forwarding table entry (x, I)



Example 2: Choosing among multiple ASes

- r now suppose AS1 learns from inter-AS protocol that subnet *x* is reachable from AS3 *and* from AS2.
- r to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest X.

m this is also job of inter-AS routing protocol!



Example 2: Choosing among multiple ASes

- r now suppose AS1 learns from inter-AS protocol that subnet *x* is reachable from AS3 *and* from AS2.
- r to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest ×.

m this is also job of inter-AS routing protocol!

r hot potato routing: send packet towards closest of the two routers.



Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP

Intra-AS Routing

- r also known as Interior Gateway Protocols (IGP)
- r most common Intra-AS routing protocols:
 - m RIP: Routing Information Protocol
 - m OSPF: Open Shortest Path First
 - m IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing



Routing table in D

Distance vectors: advertised every 30 sec (no advertisement heard after 180 sec --> neighbor/link declared dead)
 4: Network Layer 4b-

RIP Table processing

- r RIP routing tables managed by application-level process called route-d (daemon)
- r advertisements sent in UDP packets, periodically repeated



Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP

OSPF (Open Shortest Path First)

- r "open": publicly available
- r Uses Link State algorithm (configurable edge-costs)
 - Advertisements disseminated to entire AS (via flooding), via IP packets (unlike RIP)
- r OSPF "advanced" features (Note: features of the standardized protocol, not the algorithm) -not in RIP
 - m Security: all OSPF messages authenticated (to prevent malicious intrusion)
 - m Multiple same-cost paths allowed (only one path in RIP)
 - m multiple cost metrics for different TypeOfService (eg, satellite link cost "low" for best effort; high for real time)
 - m Integrated uni- and multicast support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF

m Hierarchical OSPF in large domains.

Hierarchical OSPF boundary router backbone router Backbone area border routers internal routers Area 3 Area 1 Area 2

<u>Hierarchical OSPF</u>

- r Two-level hierarchy: local area, backbone.
 m Link-state advertisements only in area
 m each node has detailed area topology; only know direction (shortest path) to nets in other areas.
- r Area border routers: "summarize" distances to nets in own area, advertise to other Area Border routers.
- r **Backbone routers:** run OSPF routing limited to backbone.
- r Boundary routers: connect to other ASs.

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing

Internet inter-AS routing: BGP

- r BGP (Border Gateway Protocol): *the* de facto standard
- r BGP provides each AS a means to:
 - 1. Obtain subnet reachability information from neighboring ASs.
 - 2. Propagate reachability information to all ASinternal routers.
 - 3. Determine "good" routes to subnets based on reachability information and policy.
- r allows subnet to advertise its existence to rest of Internet: "*I am here*"

BGP basics

- r pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: BGP sessions
 m External, internal: eBGP, iBGP
 m BCP consistence peed not composed to physical
 - m BGP sessions need not correspond to physical links.
- when AS2 advertises a prefix (e.g. subnet) to AS1:
 m AS2 promises it will forward datagrams towards that prefix.
 - m AS2 can aggregate prefixes in its advertisement



Distributing reachability info

- r using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
 - m 1c can then use iBGP do distribute new prefix info to all routers in AS1
 - m 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session
- r when router learns of new prefix, it creates entry for prefix in its forwarding table.



BGP: routing

 Path Vector protocol (similar to Distance Vector): each Border Gateway advertises *entire path* (I.e, sequence of ASs) to destination

Suppose: gateway X send its path to peer gateway W

- m W may or may not select path offered by X
 - cost, policy (don't route via competitor's AS), loop prevention reasons.
- m If W selects path advertised by X, then:

Path (W,Z) = w, Path (X,Z)

- Mote: X can control incoming traffic by controling its route advertisements to peers:
 - e.g., don't want to route traffic to Z -> don't advertise any routes to Z

Path attributes & BGP routes

- r advertised prefix includes BGP attributes. m prefix + attributes = "route"
- r two important attributes:
 - m AS-PATH: contains ASs through which prefix advertisement has passed: e.g, AS 67, AS 17
 - m NEXT-HOP: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- r when gateway router receives route advertisement, uses import policy to accept/decline.

BGP route selection

- r router may learn about more than 1 route to some prefix. Router must select route.
- r elimination rules:
 - 1. local preference value attribute: policy decision
 - 2. shortest AS-PATH
 - 3. closest NEXT-HOP router: hot potato routing
 - 4. additional criteria

BGP messages

- r BGP messages exchanged using TCP.
- r BGP messages:
 - m OPEN: opens TCP connection to peer and authenticates sender
 - m UPDATE: advertises new path (or withdraws old)
 - MEEPALIVE keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - m NOTIFICATION: reports errors in previous msg; also used to close connection

BGP routing policy: example



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

<u>BGP routing policy: example (cont)</u>



legend: provider network customer network:

- r A advertises path AW to B
- r Badvertises path BAW to X
- r Should B advertise path BAW to C?
 - m No way! B gets no "revenue" for routing CBAW since neither W nor C are B's customers
 - m B wants to force C to route to w via A
 - m B wants to route only to/from its customers! Network Layer 4-50

Why different Intra- and Inter-AS routing?

Policy:

- r Inter-AS: admin wants control over how its traffic routed, who routes through its net.
- r Intra-AS: single admin, so no policy decisions needed Scale:
- r hierarchical routing saves table size, reduced update traffic

Performance:

- r Intra-AS: can focus on performance
- r Inter-AS: policy may dominate over performance

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6

- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP

<u>Review questions for this part</u>

- r Most commonly used routing protocols in the Internet?
 - m What algorithms they use? Why?
 - m What else besides algorithms choices is important? (hint: think about policies) Why?

Broadcast routing

deliver packets from source to all other nodes
source duplication is inefficient:



source duplication: how does source determine recipient addresses?

In-network duplication

- *flooding:* when node receives broadcast packet, sends copy to all neighbors
 m problems: cycles & broadcast storm
- r controlled flooding: node only broadcasts pkt if it hasn't broadcast same packet before
 - m node keeps track of packet ids already broadacsted
 - m or reverse path forwarding (RPF): only forward packet if it arrived on shortest path between node and source
- r *spanning tree*:

m no redundant packets received by any node



first construct a spanning tree

nodes then forward/make copies only along spanning tree



(a) broadcast initiated at A



Spanning tree: creation

center node

- each node sends unicast join message to center node
 - message forwarded until it arrives at a node already belonging to spanning tree



(a) stepwise construction of spanning tree (center: E)



(b) constructed spanning tree