Chapter 4 Network Layer

Chapter 4: Network Layer

- 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

Intra-AS Routing

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

RIP (Routing Information Protocol)

- included in BSD-UNIX distribution in 1982
- distance vector algorithm
 - distance metric: # hops (max = 15 hops), each link has cost 1
 - Hop: is the number of subnets traversed along the shortest path from source router to destination subnet, including the destination subnet.
 - DVs exchanged with neighbors every 30 sec in response message (called advertisement)
 - each advertisement: list of up to 25 destination subnets (in IP addressing sense)



from router A to destination subnets:

<u>subnet</u>	<u>hops</u>
u	1
V	2
W	2
Х	3
У	3
7	2

RIP: Example



routing table in router D

destination subnet	next router	# hops to dest
W	A	2
У	В	2
Z	В	7
X		1



routing table in router D

destination subnet	next router	# hops to dest
W	А	2
y y	В	2
Z	В	7
X		1



routing table in router D

destination subnet	next router	# hops to dest
W	A	2
У	В	2 5
Z	BA	75
X		1

RIP: Link Failure and Recovery

If no advertisement heard after 180 sec --> neighbor/link declared dead

- routes via neighbor invalidated
- new advertisements sent to neighbors
- neighbors in turn send out new advertisements (if tables changed)
- Ink failure info quickly (?) propagates to entire net
- poison reverse used to prevent ping-pong loops (infinite distance = 16 hops)

RIP Table processing

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



OSPF (Open Shortest Path First)

- "open": publicly available
- uses Link State algorithm
 - LS packet dissemination
 - topology map at each node
 - route computation using Dijkstra's algorithm
- OSPF advertisement carries one entry per neighbor router
- advertisements disseminated to entire AS (via flooding)
 - carried in OSPF messages directly over IP (rather than TCP or UDP

OSPF "advanced" features (not in RIP)

- security: all OSPF messages authenticated (to prevent malicious intrusion)
- multiple same-cost paths allowed (only one path in RIP)
- for each link, multiple cost metrics for different TOS (e.g., satellite link cost set "low" for best effort ToS; high for real time ToS)
- integrated uni- and multicast support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- hierarchical OSPF in large domains.



Hierarchical OSPF

* two-level hierarchy: local area, backbone.

- Ink-state advertisements only in area
- each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
- <u>area border routers:</u> "summarize" distances to nets in own area, advertise to other Area Border routers.
- * *backbone routers:* run OSPF routing limited to backbone.
- boundary routers: connect to other AS's.

Internet inter-AS routing: BGP

- BGP (Border Gateway Protocol): the de facto interdomain routing protocol
 - "glue that holds the Internet together"
- BGP provides each AS a means to:
 - eBGP: obtain subnet reachability information from neighboring ASs.
 - iBGP: propagate reachability information to all AS-internal routers.
 - determine "good" routes to other networks based on reachability information and policy.
- allows subnet to advertise its existence to rest of Internet: "I am here"

BGP basics

- BGP session: two BGP routers ("peers") exchange BGP messages:
 - advertising *paths* to different destination network prefixes ("path vector" protocol)
 - exchanged over semi-permanent TCP connections



BGP basics

- BGP session: two BGP routers ("peers") exchange BGP messages:
 - advertising *paths* to different destination network prefixes ("path vector" protocol)
 - exchanged over semi-permanent TCP connections
- when AS3 advertises a prefix to AS1:
 - AS3 promises it will forward datagrams towards that prefix
 - AS3 can aggregate prefixes in its advertisement



BGP basics: distributing path information

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



Path attributes & BGP routes

- advertised prefix includes BGP attributes
 - prefix + attributes = "route"
- two important attributes:
 - AS-PATH: contains ASs through which prefix advertisement has passed: e.g., AS 67, AS 17
 - NEXT-HOP: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- gateway router receiving route advertisement uses import policy to accept/decline
 - e.g., never route through AS x
 - policy-based routing

BGP route selection

- router may learn about more than 1 route to destination AS, selects route based on:
 - 1. local preference value attribute: policy decision
 - 2. shortest AS-PATH
 - 3. closest NEXT-HOP router: hot potato routing
 - 4. additional criteria

BGP messages

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

BGP routing policy





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

BGP routing policy (2)





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

Why different Intra- and Inter-AS routing ?

Policy:

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

Performance:

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

Chapter 4: Network Layer

- 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

Broadcast Routing

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



source duplication in-network duplication

source duplication: how does source determine recipient addresses?

In-network duplication

- flooding: when node receives broadcast packet, sends copy to all neighbors
 - problems: cycles & broadcast storm
- controlled flooding: node only broadcasts pkt if it hasn't broadcast same packet before
 - node keeps track of packet ids already broadcasted
 - or reverse path forwarding (RPF): only forward packet if it arrived on shortest path between node and source
- spanning tree
 - No redundant packets received by any node

Spanning Tree

- First construct a spanning tree
- Nodes forward copies only along spanning tree



(a) Broadcast initiated at A



(b) Broadcast initiated at D

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



(b) Constructed spanning tree

Multicast Routing: Problem Statement

- Goal: find a tree (or trees) connecting routers having local mcast group members
 - tree: not all paths between routers used
 - source-based: different tree from each sender to receivers
 - shared-tree: same tree used by all group members



Shared tree

Source-based trees

Approaches for building mcast trees

Approaches:

- source-based tree: one tree per source
 - shortest path trees
 - reverse path forwarding
- group-shared tree: group uses one tree
 - minimal spanning (Steiner)
 - center-based trees

...we first look at basic approaches, then specific protocols adopting these approaches

Shortest Path Tree

- mcast forwarding tree: tree of shortest path routes from source to all receivers
 - Dijkstra's algorithm



LEGEND

- router with attached group member
- router with no attached group member
- i link used for forwarding, i indicates order link added by algorithm

Reverse Path Forwarding

- rely on router's knowledge of unicast shortest path from it to sender
- each router has simple forwarding behavior:

if (mcast datagram received on incoming link on shortest path back to center)
 then flood datagram onto all outgoing links
 else ignore datagram

Reverse Path Forwarding: example



LEGEND



- router with no attached group member
 - datagram will be forwarded
- datagram will not be forwarded
- result is a source-specific reverse SPT
 - may be a bad choice with asymmetric links

Reverse Path Forwarding: pruning

- forwarding tree contains subtrees with no multicast group members
 - no need to forward datagrams down subtree
 - "prune" msgs sent upstream by router with no downstream group members



LEGEND

- router with attached group member
- router with no attached group member
 - prune message
 - links with multicast forwarding

Shared-Tree: Steiner Tree

- Steiner Tree: minimum cost tree connecting all routers with attached group members
- problem is NP-complete
- excellent heuristics exists
- not used in practice:
 - computational complexity
 - information about entire network needed
 - monolithic: rerun whenever a router needs to join/leave

Center-based trees

- single delivery tree shared by all
- one router identified as "center" of tree
- * to join:
 - edge router sends unicast join-msg addressed to center router
 - join-msg "processed" by intermediate routers and forwarded towards center
 - *join-msg* either hits existing tree branch for this center, or arrives at center
 - path taken by *join-msg* becomes new branch of tree for this router

Center-based trees: an example

Suppose R6 chosen as center:



LEGEND

- e ra
 - router with attached group member
 - Fouter with no attached group member
 - path order in which join messages generated

Internet Multicasting Routing: DVMRP

- DVMRP: distance vector multicast routing protocol, RFC1075
- flood and prune: reverse path forwarding, source-based tree
 - RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers
 - no assumptions about underlying unicast
 - initial datagram to mukticast group flooded everywhere via RPF
 - routers not wanting group: send upstream prune messages

DVMRP: continued...

- Soft state: DVMRP router periodically (1 min.) "forgets" branches are pruned:
 - mcast data again flows down unpruned branch
 - downstream router: reprune or else continue to receive data
- routers can quickly regraft to tree
 - following IGMP join at leaf
- odds and ends
 - commonly implemented in commercial routers
 - Mbone routing done using DVMRP

Tunneling

Q: How to connect "islands" of multicast routers in a "sea" of unicast routers?



physical topology

logical topology

- mcast datagram encapsulated inside "normal" (non-multicastaddressed) datagram
- normal IP datagram sent thru "tunnel" via regular IP unicast to receiving mcast router
- receiving mcast router unencapsulates to get mcast datagram

PIM: Protocol Independent Multicast

- not dependent on any specific underlying unicast routing algorithm (works with all)
- two different multicast distribution scenarios :

<u>Dense:</u>

- group members densely packed, in "close" proximity.
- bandwidth more plentiful

<u>Sparse:</u>

- # networks with group members small wrt # interconnected networks
- group members "widely dispersed"
- bandwidth not plentiful

Consequences of Sparse-Dense Dichotomy:

<u>Dense</u>

- group membership by routers assumed until routers explicitly prune
- *data-driven* construction
 on mcast tree (e.g., RPF)
- bandwidth and nongroup-router processing profligate

<u>Sparse</u>:

- no membership until routers explicitly join
- receiver- driven
 construction of mcast tree
 (e.g., center-based)
- bandwidth and non-grouprouter processing conservative

PIM- Dense Mode

flood-and-prune RPF, similar to DVMRP but

- underlying unicast protocol provides RPF info for incoming datagram
- less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm
- has protocol mechanism for router to detect it is a leaf-node router

PIM - Sparse Mode

- center-based approach
- router sends join msg to rendezvous point (RP)
 - intermediate routers update state and forward join
- after joining via RP, router can switch to source-specific tree
 - increased performance: less concentration, shorter paths



PIM - Sparse Mode

sender(s):

- unicast data to RP, which distributes down RP-rooted tree
- RP can extend mcast tree upstream to source
- RP can send stop msg if no attached receivers
 - "no one is listening!"

