

College of Information Studies

University of Maryland Hornbake Library Building College Park, MD 20742-4345

Routing

Session 20 INST 346 Technologies, Infrastructure and Architecture

Goals for Today

- Shortest-Path Routing
- Routers
- Border Gateway Protocol
- Analysis Group 4

aggregate routers into regions known as "autonomous systems" (AS) (a.k.a. "domains")

intra-AS routing

- routing among hosts, routers in same AS ("network")
- all routers in AS must run same intra-domain protocol
- routers in different AS can run different intra-domain routing protocol
- gateway router: at "edge" of its own AS, has link(s) to router(s) in other AS'es

inter-AS routing

- routing among AS'es
- gateways perform interdomain routing (as well as intra-domain routing)

Interconnected ASes



- forwarding table configured by both intraand inter-AS routing algorithm
 - intra-AS routing determine entries for destinations within AS
 - inter-AS & intra-AS determine entries for external destinations

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 (IS-IS protocol essentially same as OSPF)
 - IGRP: Interior Gateway Routing Protocol (Cisco proprietary for decades, until 2016)

Intra-AS Routing (OSPF)

- (Open) Shortest Path First
- A "link state" method
- First get a complete network map at each node
 - Each router floods the AS with OSPF "advertisements"
 - Advertisement: list of adjacent routers with estimated delay
- Use Dijkstra's algorithm for shortest path computation

Dijsktra'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
$$D(v) = \infty$$

node x to y;
$$= \infty$$
 if
not direct neighbors
D(v): current value
of cost of path from
source to dest. v
p(v): predecessor
node along path from

c(x,y): link cost from

N': set of nodes whose least cost path definitively known

source to v

8 **Loop**

7

- 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

5

7

H

4

3

-W

8

		D(v)	D(w)	D(X)	D(y)	D(z)
Step	> N'	p(v)	p(w)	p(x)	p(y)	p(z)
0	u	7,u	(3,u	5,u	∞	∞
1	uw	6,w		<u>(5,u</u>) 11,w	8
2	uwx	6,w			11,W	14,X
3	UWXV				10,V	14,X
4	uwxvy					(12,y)
5	uwxvyz					

 construct shortest path tree by tracing predecessor nodes D(v): current value of cost of path from source to dest. v

p(v): predecessor node along path from source to v

N': set of nodes whose least cost path definitively known

Dijkstra's algorithm: another example

St	ер	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		3,y			4,y
	4	uxyvw 🔶					4,y
	5						



D(v): current value of cost of path from source to dest. v

p(v): predecessor node along path from source to v

N': set of nodes whose least cost path definitively known

Dijkstra's algorithm: solution

resulting shortest-path tree from u:



resulting forwarding table in u:

destination	link
V	(u,v)
Х	(u,x)
У	(u,x)
W	(u,x)
Z	(u,x)
	1

Logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs) in routers to compute forwarding tables



Router architecture overview

high-level view of generic router architecture:





forwarding rate into switch fabric

Input port queuing

- fabric slower than input ports combined -> queueing may occur at input queues
 - queueing delay and loss due to input buffer overflow!
- Head-of-the-Line (HOL) blocking: queued datagram at front of queue prevents others in queue from moving forward



blocking

lower red packet is blocked

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



 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

bus

Destination-based forwarding

forwarding table				
Destination Address Range	Link Interface			
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 111111	0			
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 1111111	1			
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 1111111	2			
otherwise	3			

Q: but what happens if ranges don't divide up so nicely?

Longest prefix matching

- longest prefix matching

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** ********	0
11001000 00010111 00011000 ********	1
11001000 00010111 00011*** ********	2
otherwise	3

examples:

DA: 11001000 00010111 00010110 10100001 DA: 11001000 00010111 00011000 10101010 which interface? which interface?

Longest prefix matching

- longest prefix matching: often performed using ternary content addressable memories (TCAMs)
 - content addressable: present address to TCAM: retrieve address in one clock cycle, regardless of table size
 - Cisco Catalyst: can up ~IM routing table entries in TCAM



This slide in HUGELY important!



- buffering required from fabric faster rate
 Datagram (packets) can be lost due to congestion, lack of buffers
- scheduling datagrams

Priority scheduling – who gets best performance, network neutrality

Output port queueing



- buffering when arrival rate via switch exceeds output line speed
- queueing (delay) and loss due to output port buffer overflow!

How much buffering?

- RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
 - e.g., C = 10 Gpbs link: 2.5 Gbit buffer
- recent recommendation: with N flows, buffering equal to

Scheduling policies

- scheduling: choose next packet to send on link
- FIFO (first in first out) scheduling: send in order of arrival to queue
 - real-world example?
 - discard policy: if packet arrives to full queue: who to discard?
 - *tail drop*: drop arriving packet
 - *priority*: drop/remove on priority basis
 - *random*: drop/remove randomly



Scheduling policies

Weighted Fair Queuing (WFQ):

- generalized Round Robin
- each class gets weighted amount of service in each cycle





Hierarchical OSPF

- *two-level hierarchy:* local area, backbone.
 - link-state advertisements only in area
 - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
- area border routers: "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' es.

Inter-AS routing is different

policy:

- intra-AS: single admin, so single consistent policy
- inter-AS: each admin wants control over how its traffic routed and who routes through its AS

performance:

- intra-AS: can focus on performance
- inter-AS: policy may dominate over performance

Inter-AS tasks

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

AS1 must:

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



Internet inter-AS routing: BGP

- BGP (Border Gateway Protocol): the de facto inter-domain routing protocol
 - "glue that holds the Internet together"
- BGP provides each AS a means to:
 - eBGP: obtain subnet reachability information from neighboring ASes
 - **iBGP:** propagate reachability information to all ASinternal routers.
 - determine "good" routes to other networks based on reachability information and policy
- allows subnet to advertise its existence to rest of Internet: "1 am here"

eBGP, iBGP connections





gateway routers run both eBGP and iBGP protools

BGP basics

- BGP session: two BGP routers ("peers") exchange BGP messages over semi-permanent TCP connection:
 - advertising paths to different destination network prefixes (BGP is a "path vector" protocol)
- when AS3 gateway router 3a advertises path AS3,X to AS2 gateway router 2c:
 - AS3 promises to AS2 it will forward datagrams towards X



Path attributes and BGP routes

- advertised prefix includes BGP attributes
 - prefix + attributes = "route"
- two important attributes:
 - AS-PATH: list of ASes through which prefix advertisement has passed
 - NEXT-HOP: indicates specific internal-AS router to nexthop AS
- Policy-based routing:
 - gateway receiving route advertisement uses import policy to accept/decline path (e.g., never route through AS Y).
 - AS policy also determines whether to *advertise* path to other other neighboring ASes

BGP path advertisement



- AS2 router 2c receives path advertisement AS3,X (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path AS3,X, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path AS2, AS3, X to AS1 router 1c

BGP path advertisement



gateway router may learn about multiple paths to destination:

- AS1 gateway router 1C learns path AS2,AS3,X from 2a
- AS1 gateway router 1C learns path AS3,X from 3a
- Based on policy, AS1 gateway router 1C chooses path AS3, X, and advertises path within AS1 via iBGP

BGP: achieving policy via advertisements



Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A advertises path Aw to B and to C
- B chooses not to advertise BAw to C:
 - B gets no "revenue" for routing CBAw, since none of C, A, w are B's customers
 - C does not learn about CBAw path
- C will route CAw (not using B) to get to w

BGP: achieving policy via advertisements



Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

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

BGP route selection

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

Hot Potato Routing



- 2d learns (via iBGP) it can route to X via 2a or 2c
- hot potato routing: choose local gateway that has least intradomain cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!

Network Layer Summary

- IPv4 addresses
 - Hierarchical structure (subnet mask)
- Routing
 - Hierarchical structure (Autonomous Systems)
- Routers
 - Structure (input queue, switch, output queue)
 - Routing tables (hierarchical structure)
- Network layer packets
 - IPv4, IPv6