Service Delivery Architecture

Lecture 4a Srinivas Narayana http://www.cs.rutgers.edu/~sn624/553-S25



Improving Performance on the Internet



Internet Routing



Distributed routing

Control plane





The Internet is a large federated network

Several autonomously run organizations (AS'es): No one "boss"

Organizations cooperate, but also compete



The Internet is a large federated network

Several autonomously run organizations: No one "boss"

Organizations cooperate, but also compete

AT&T

Verizon

Comcast

Message exchanges must not reveal internal network details.

Algorithm must work with "incomplete" information about its neighbors' internal topology.

((p))

RUTGERS

The Internet is a large federated network

AT&T

Verizon

Comcast

Internet today: > 70,000 unique autonomous networks

Internet routers: > 800,000 forwarding table entries

Keep messages & tables as small as possible. Don't flood

Algorithm must be incremental: don't recompute the whole table on every message exchanged.

((_))

RUTGERS

Inter-domain Routing

- The Internet uses Border Gateway Protocol (BGP)
- All AS'es speak BGP. It is the glue that holds the Internet together
- BGP is a path vector protocol



Q1. BGP Messages



"I can reach X"

Dst: 128.1.2.0/24

AS path: AS2, X

2a

Exchange paths: path vector

AS 2

2b

2d

"I am here."

AS path: X

2c

Dst: 128.1.2.0/24

- Routing Announcements or Advertisements No internal link or topology
 - "I am here" or "I can reach here"
 - Occur over a TCP connection (BGP session) between routers
- Route announcement = destination + attributes

1b

1d

- Destination: IP prefix
- Route Attributes:
 - AS-level path
 - Next hop
 - Several others: origin, MED, community, etc.

1a

An AS promises to use advertised path to reach destination

1c

• Only route changes are advertised after BGP session established

Q1. Next Hop



- Next hop conceptually denotes the first router interface that begins the AS-level path
 - The meaning of this attribute is context-dependent
- In an announcement arriving from a different AS (eBGP), next hop is the router in the next AS which sent the announcement
 - Example: Next Hop of the eBGP announcement reaching 1c is 2a



Q1. Next Hop



- Suppose router 1c receives a path advertisement
- Router 1c will propagate the announcement inside the AS using iBGP
- The next hop of this (iBGP) announcement is set to 1c
 - In particular, the next hop is an AS1 internal address



Q2. The algorithm



- A BGP router does *not* consider every routing advertisement it receives by default to make routing decisions
 - An import policy determines whether a route is even considered a candidate
- Once imported, the router performs route selection

Programmed by network

- A BGP router does not propagate its chosen path to a operator destination to all other AS'es by default
 - An export policy determines whether a (chosen) path can be advertised to other AS'es and routers

Policy considerations make BGP paths very different from "the most efficient" paths



Policy arises from business relationships

- Customer-provider relationships:
 - E.g., Rutgers is a customer of AT&T
- Peer-peer relationships:
 - E.g., Verizon is a peer of AT&T
- Business relationships depend on where connectivity occurs
 - "Where", also called a "point of presence" (PoP)
 - e.g., customers at one PoP but peers at another
 - Internet-eXchange Points (IXPs) are large PoPs where ISPs come together to connect with each other (often for free)



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 customers (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 announce to B a route to C



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

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

Impact of BGP export policies

- Based on your location on the Internet,
- Some paths aren't visible or usable, even if they physically exist
- Many efficient paths could be eliminated from actual use
- Focus on financial incentives, not efficient end-to-end paths

BGP Import Policy

- Remove common misconfigurations or problematic routes
- Loops
- Too-specific prefixes (e.g., anything longer than /24)



Suppose an ISP wants to minimize costs by avoiding routing through its providers when possible.

- Suppose C announces path Cy to x
- Further, y announces a direct path ("y") to x
- Then x may choose not to import the path Cy to y since it has a peer path ("y") towards y

Q2. BGP Route Selection



- When a router imports more than one route to a destination IP prefix, it selects route based on:
 - 1. local preference value attribute (import policy decision -- set by network admin)
 - 2. shortest AS-PATH
 - 3. closest NEXT-HOP router
 - 4. Several additional criteria: You can read up on the full, complex, list of criteria, e.g., at https://www.cisco.com/c/en/us/support/docs/ip/border-gateway-protocol-bgp/13753-25.html

Example of route selection

- Suppose AS A and B are connected to each other both in North America (NA) and in Europe (EU)
- A source in NA wants to reach a destination in EU
- There are two paths available
 - Assume same local preference
 - Same AS path length
- Closest next hop-router: choose path via B1 rather than B2



Example of route selection

- Choosing closest next-hop results in early exit routing
 - Try to exit the local AS as early as possible
 - Also called hot potato routing
- Reduce resource use within local AS
 - potentially at the expense of another AS
- A potential reason for inefficiency



BGP Path Selection

Approaches to bring flexibility: Flexible control logic for path selection (Google, Facebook) Detour/overlay routing (Akamai)

- Local preference, shortest AS path, closest NEXT HOP, etc.
- Not capacity aware
- Not performance aware
- Not aware of the length of the path (in # routers)
 - The protocol does not even incorporate precise perf/capacity info
- Financial incentive, not end-to-end performance, heavily determines peering and capacity Traceroute Path 1: from Guadalajara, Mexico to Washington, D.C. via Belarus
- Only a single path per destination
- Can be slow to converge
- Vulnerable to bugs and malice



Application Architecture

Lecture 4b Srinivas Narayana http://www.cs.rutgers.edu/~sn624/553-S25



Components of an Internet Service

Routers

Storage

App compute and communication patterns

Modularized applications

Endpoints

Interconnect: Routers

Data Center

Servers

Web Servers



Often the first app point where a user request lands







How does one design a web server?

• Process other requests while waiting for one to finish



process socket

 $IP_B + port_B$

 $bind(IPaddr_{B}, port_{B})$



listen()

accept()

recv()/send()/..

Parallelism

- Process requests in parallel with other requests
- One design: multiprocessing/multithreading (MP/MT)

