



Interdomain Routing

Reading: Sections P&D 4.3,{3,4}

EE122: Intro to Communication Networks
 Fall 2006 (MW 4:00-5:30 in Donner 155)
 Vern Paxson

TAs: Dilip Antony Joseph and Sukun Kim
<http://inst.eecs.berkeley.edu/~ee122/>

Materials with thanks to Jennifer Rexford, Ion Stoica and colleagues at Princeton and UC Berkeley

1

Outline

- Why does BGP exist?
 - What is interdomain routing and why do we need it?
 - Why does BGP look the way it does?
- How does BGP work?
 - Boring details
 - Yuck

pay more attention to the “why” than the “how”

3

Routing

- Provides paths between networks
- Previous lecture presented two routing designs
 - link-state
 - distance vector
- Previous lecture assumed single domain
 - all routers have same routing metric (shortest path)
 - no privacy issues, no policy issues

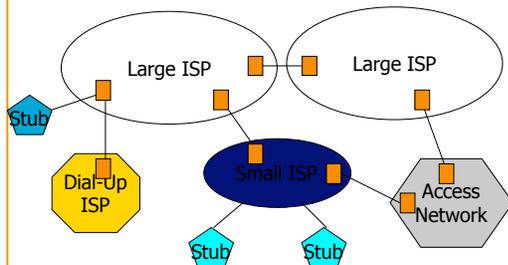
4

Internet is more complicated.....

- Internet not just unstructured collection of networks
- Internet is comprised of a set of “autonomous systems” (ASes)
 - independently run networks, some are commercial ISPs
 - currently around 20,000 ASes
- ASes are sometimes called “domains”
 - hence “interdomain routing”

5

Internet: a large number of ASes



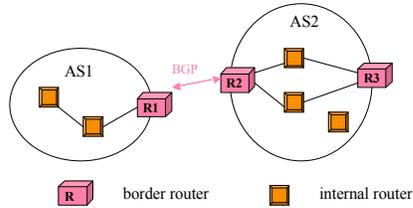
6

This adds another level in hierarchy

- Three levels in logical routing hierarchy
 - networks: reaches individual hosts
 - intradomain: routes between networks
 - interdomain: routes between ASes
- Need a protocol to route between domains
 - BGP is current standard
- Different kinds of unification
 - IP unifies network technologies
 - BGP unifies network organizations

7

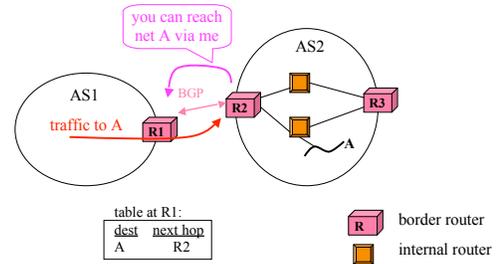
Who speaks BGP?



- Two types of routers
 - Border router (Edge), Internal router (Core)

8

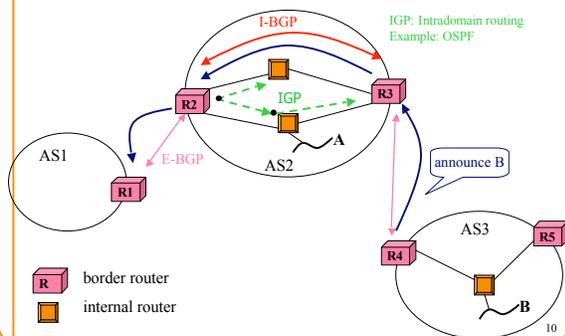
Purpose of BGP



Share connectivity information across ASes

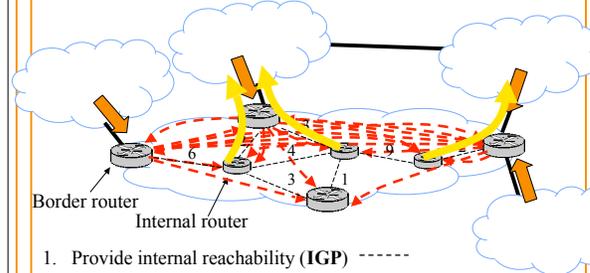
9

I-BGP and E-BGP



10

In more detail



- Provide internal reachability (IGP) -----
- Learn routes to external destinations (eBGP) →
- Distribute externally learned routes internally (iBGP) - - - - -
- Select closest egress (IGP) -----

11

Rest of lecture...

- Motivate why BGP is the way it is
 - driven by two salient aspects of AS structure
- Discuss some problems with interdomain routing
- Discuss (briefly!) what a new BGP might look like
- Explain some of BGP's details
 - not fundamental, just series of specific design decisions

12

#1 ASes are autonomous

- Want to choose their own internal routing protocol
 - different algorithms and metrics
- Want freedom to route based on policy
 - “my traffic can't be carried over my competitor's network”
 - “I don't want to carry transit traffic through my network”
 - not expressible as Internet-wide “shortest path”!
- Want to keep their connections and policies private
 - would reveal business relationships, network structure

13

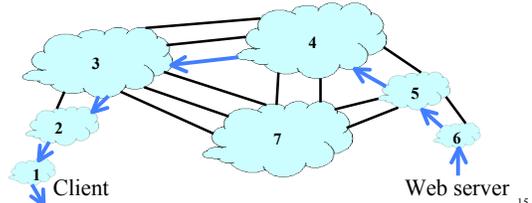
#2 ASes have business relationships

- Three kinds of relationships between ASes
 - AS A can be AS B's *customer*
 - AS A can be AS B's *provider*
 - AS A can be AS B's *peer*
- Business implications
 - customer pays provider
 - peers don't pay each other
- Policy implications
 - "When sending traffic, I prefer to route through customers over peers, and peers over providers"
 - "I don't carry traffic from one provider to another provider"

14

AS-level topology

- Destinations are IP prefixes (e.g., 12.0.0.0/8)
- Nodes are Autonomous Systems (ASes)
 - internals are hidden
- Links are connections & business relationships



15

What routing algorithm can we use?

- Key issues are *policy* and *privacy*
- Can't use shortest path
 - domains don't have any shared metric
 - *policy choices might not be shortest path*
- Can't use link state
 - would have to flood policy preferences and topology
 - *would violate privacy*

16

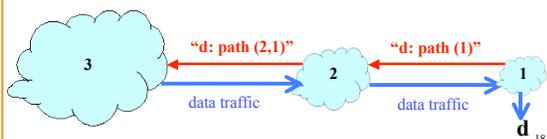
What about distance vector?

- Does not reveal any connectivity information
- But is designed to compute shortest paths
- Extend distance vector to allow policy choices?

17

Path-Vector Routing

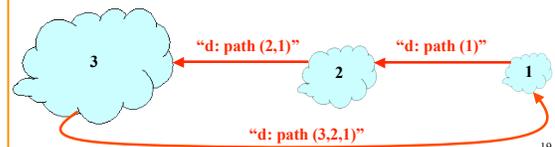
- Extension of distance-vector routing
 - Support flexible routing policies
 - Faster loop detection (no count-to-infinity)
- Key idea: advertise the entire path
 - Distance vector: send *distance metric* per dest *d*
 - Path vector: send the *entire path* for each dest *d*



18

Faster Loop Detection

- Node can easily detect a loop
 - Look for its own node identifier in the path
 - E.g., node 1 sees itself in the path "3, 2, 1"
- Node can simply discard paths with loops
 - E.g., node 1 simply discards the advertisement



19

Reachability

- In normal routing, if graph is connected then reachability is assured
- With policy routing, this does not always hold

26

Security

- An AS can claim to serve a prefix that they actually don't have a route to (blackholing traffic)
 - problem not specific to policy or path vector
 - important because of AS autonomy
- Fixable: make ASes "prove" they have a path

27

Performance

- BGP designed for policy not performance
- "Hot Potato" routing common but suboptimal
 - AS wants to hand off the packet as soon as possible
- 20% of paths inflated by at least 5 router hops
- Not clear this is a significant problem

28

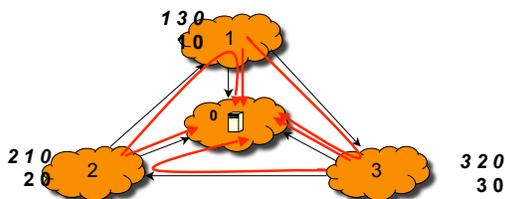
Lack of Isolation

- If there is a change in the path, the path must be re-advertised to every node upstream of the change
- Distance-vector provides more isolation

29

Persistent Oscillations due to Policies

Depends on the interactions of policies



We are back to where we started!

30

Policy Oscillations (cont'd)

- Policy autonomy vs network stability
 - focus of much recent research
- If there is no global constraint, then any degree of autonomy can lead to oscillations
 - only "shortest path" is guaranteed to be stable
- However, if policies follow normal business practices, stability is guaranteed
 - lack of cycles in business graph a global constraint

31

Redesigning BGP

- If we keep all the current constraints, not many alternative design options (at high-level)
 - Which constraints might we lift?
- Are most policies really private?
 - could use link-state for some of the routing
- Do ASes really need to see the entire path?
 - could hide some of the path, reducing updates
- Can AS structure be integrated into addressing?

32

Any Questions?

33

Rest of lecture....

- BGP details
- Stay awake as long as you can....

34

Border Gateway Protocol (BGP)

- Interdomain routing protocol for the Internet
 - Prefix-based path-vector protocol
 - Policy-based routing based on AS Paths
 - Evolved during the past 15 years

- **1989 : BGP-1 [RFC 1105]**
 - Replacement for EGP (1984, RFC 904)
- **1990 : BGP-2 [RFC 1163]**
- **1991 : BGP-3 [RFC 1267]**
- **1995 : BGP-4 [RFC 1771]**
 - Support for Classless Interdomain Routing (CIDR)

35

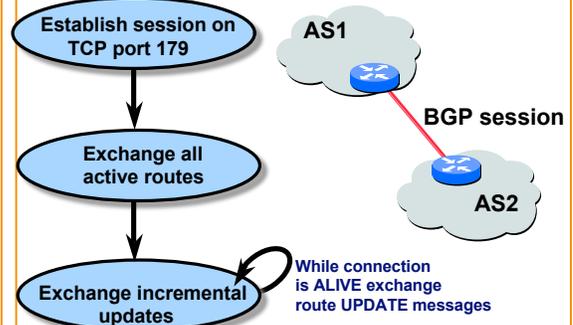
BGP Routing Table

```

ner-routes>show ip bgp
BGP table version is 6128791, local router ID is 4.2.34.165
Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
Origin codes: i - IGP, e - EGP, ? - incomplete
   Network        Next Hop        Metric LocPrf Weight Path
  * 13.0.0.0      4.0.6.142       1000   50      0 701 80 i
  * 14.0.0.0      4.24.1.35        0     100     0 i
  * 112.3.21.0/23 192.205.32.153   0      50      0 701 8 4264 6468 ?
  * e128.32.0.0/16 192.205.32.153  0      50      0 701 8 4264 6468 25 e
    
```

36

BGP Operations



37

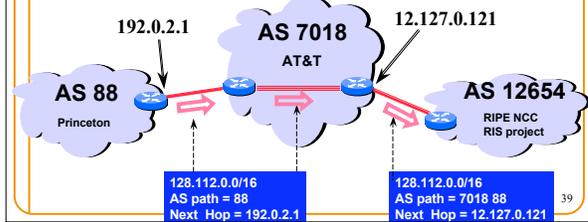
Incremental Protocol

- A node learns multiple paths to destination
 - Stores all of the routes in a routing table
 - Applies policy to select a single active route
 - ... and may advertise the route to its neighbors
- Incremental updates
 - Announcement
 - Upon selecting a new active route, add node id to path
 - ... and (optionally) advertise to each neighbor
 - Withdrawal
 - If the active route is no longer available
 - ... send a withdrawal message to the neighbors

38

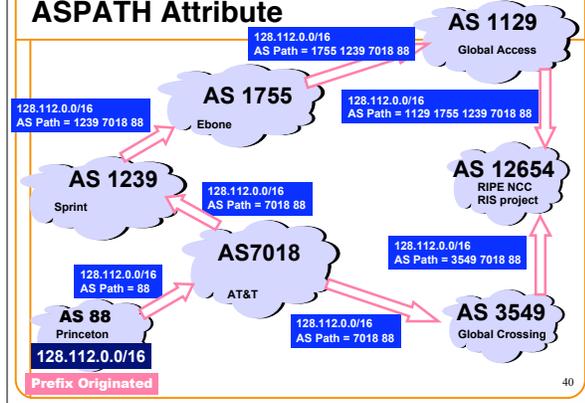
BGP Route

- Destination prefix (e.g., 128.112.0.0/16)
- Routes have attributes, including
 - AS path (e.g., "7018 88")
 - Next-hop IP address (e.g., 12.127.0.121)



39

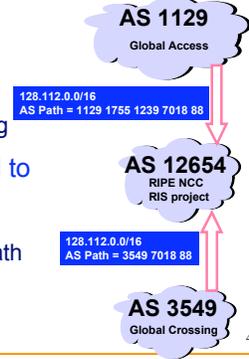
ASPATH Attribute



40

BGP Path Selection

- Simplest case
 - Shortest AS path
 - Arbitrary tie break
 - AS 12654 prefers path through Global Crossing
- But, BGP is not limited to shortest-path routing
 - Policy-based routing
 - Could choose longer path



41

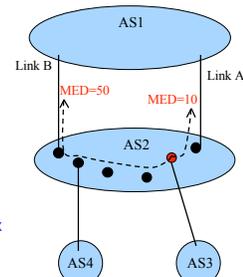
Other Routing Attributes

- Origin, MED, Local Preference,....
- Origin:
 - Who originated the announcement?
 - Where was a prefix injected into BGP?
 - IGP, BGP or Incomplete (often used for static routes)

42

Multi-Exit Discriminator (MED)

- When ASes interconnected via 2 or more links
- AS announcing prefix sets MED (AS2 in picture)
- AS receiving prefix uses MED to select link
- A way to specify how close a prefix is to the link it is announced on



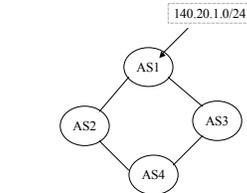
43

Local Preference

Policy choice between different AS paths

The higher the value the more preferred

Carried by IBGP, local to the AS.



BGP table at AS4:

Destination	AS Path	Local Pref
140.20.1.0/24	AS3 AS1	300
140.20.1.0/24	AS2 AS1	100

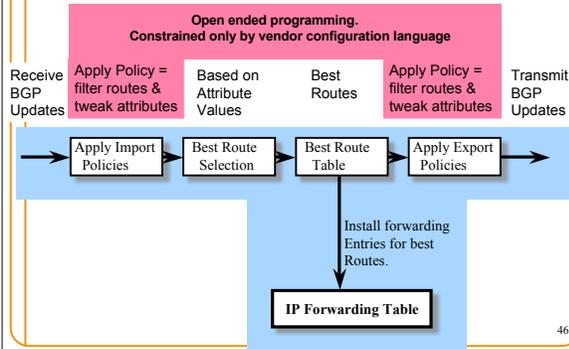
44

Choosing Best Route (simplified)

- Choose AS path with highest LOCAL_PREF
 - Preference-based routing
 - Tie: select route with shortest hop-count
- Multiple egress choices for same neighboring AS:
 - choose path with min MED value
- Among IGP paths, choose one with lowest cost
 - Finally use router ID to break the tie.

45

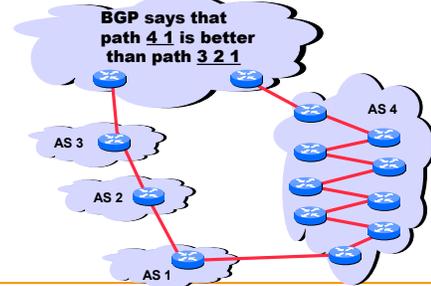
BGP Route Processing



46

AS is Not a Single Node

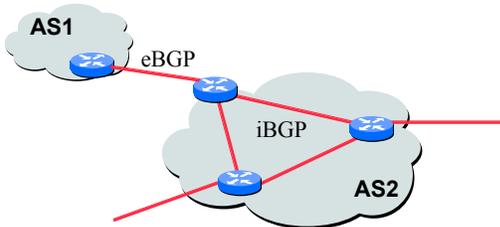
- AS path length can be misleading
 - An AS may have many router-level hops



47

An AS is Not a Single Node

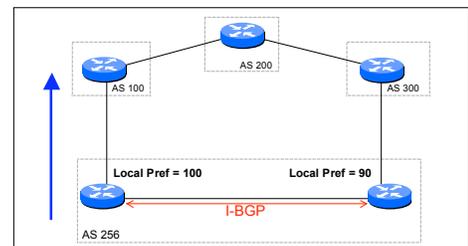
- Multiple routers in an AS
 - Need to distribute BGP information within the AS
 - Internal BGP (iBGP) sessions between routers



48

Internal BGP and Local Preference

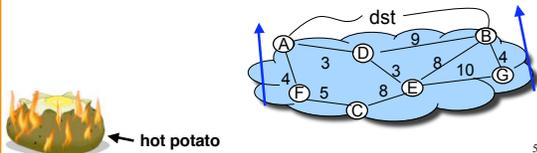
- Example
 - Both routers prefer the path through AS 100 on the left
 - ... even though the right router learns an external path



49

Hot-Potato (Early-Exit) Routing

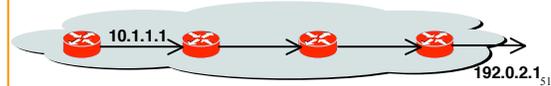
- Hot-potato routing
 - Each router selects the closest egress point
 - ... based on the path cost in intradomain protocol
- Somewhat in conflict with MED



50

Joining BGP and IGP Information

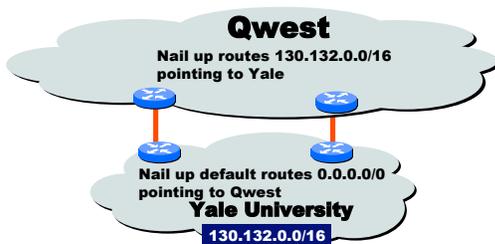
- Border Gateway Protocol (BGP)
 - Announces reachability to external destinations
 - Maps a destination prefix to an egress point
 - 128.112.0.0/16 reached via 192.0.2.1
- Interior Gateway Protocol (IGP)
 - Used to compute paths within the AS
 - Maps an egress point to an outgoing link
 - 192.0.2.1 reached via 10.1.1.1



51

Some Routers Don't Need BGP

- Customer that connects to a single upstream ISP
 - The ISP can introduce the prefixes into BGP
 - ... and the customer can simply default-route to the ISP



52

Summary

- BGP is essential to the Internet
 - ties different organizations together
- Poses fundamental challenges....
 - leads to use of path vector approach
- ...and myriad details

53