

# CS519: Computer Networks

Lecture 4, Part 5: Mar 1, 2004  
*Internet Routing:*

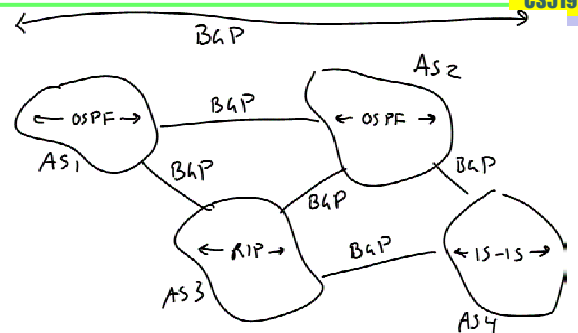
## AS's, igp, and BGP

- As we said earlier, the Internet is composed of Autonomous Systems (ASs)
  - Where each AS is a set of routers, links, and hosts
  - And is controlled by a single administration (autonomous)
  - An ISP, a large enterprise (Cornell), etc.

## AS's, igp, and BGP

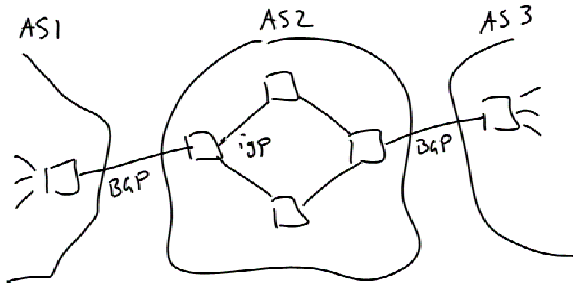
- Internally, each AS can run any routing protocol it wants
  - "interior gateway protocol", or igp
  - Examples: RIP, OSPF, IS-IS
- ASs run BGP between them
  - Border Gateway Protocol
- Though many stub ASs don't run BGP, but simply default to their ISP
  - The ISP runs BGP "on their behalf"

## AS's, igp, and BGP



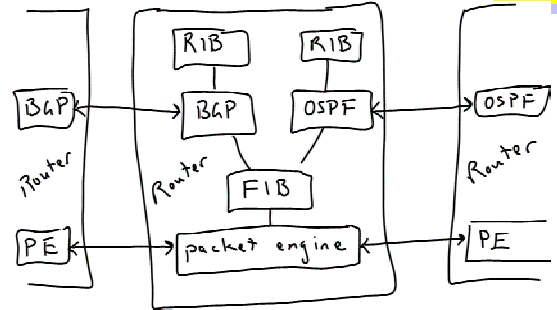
## Border routers run both BGP and the igp

CS519



## Two routing protocols, one FIB?

CS519



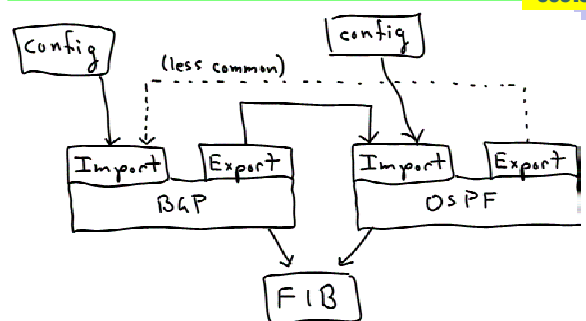
## Importing and exporting routes

CS519

- Routing algorithms have to "originate" routes
  - Which means that the routing algorithm has to "import" the route in some way other than getting it from a neighbor router
- Two ways:
  - From configuration (iface1 has prefix P...)
  - From another routing algorithm!
- Likewise, routing algorithms can "export" routes

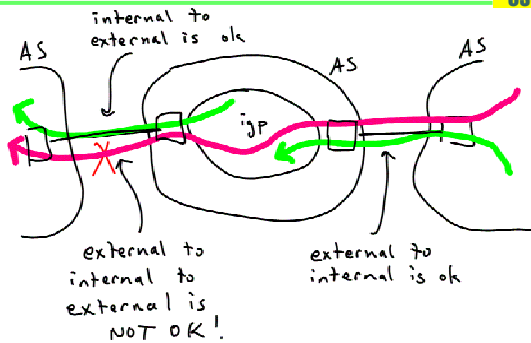
## Importing and exporting routes

CS519



## Limits to importing and exporting

CS519



## Why this limitation?

CS519

- o Semantic mismatch between BGP and igp
  - BGP is path-vector, and requires the AS-path to the destination prefix
  - igp's don't require this AS-path, and can't be expected or forced to carry it
  - Want to maintain independence between igp and BGP
- o Also, igp convergence may be slow...

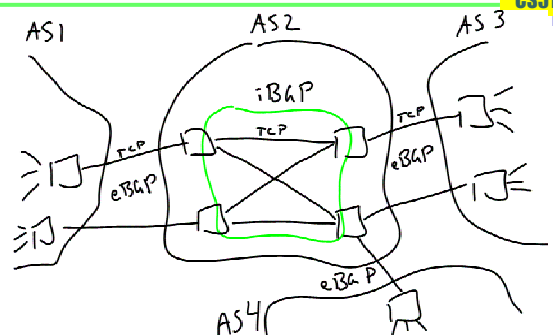
## iBGP and eBGP (interior and exterior)

CS519

- o BGP avoids dependence on igp by running both between ASs (exterior, eBGP) and within ASs (interior, iBGP)
- o iBGP runs over TCP

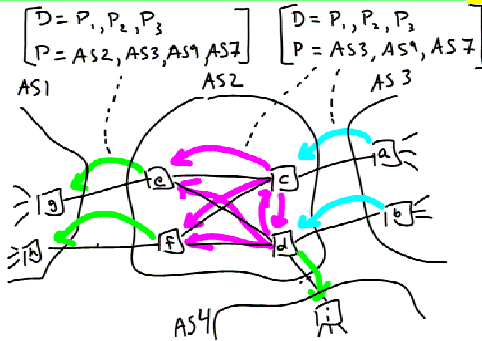
## iBGP and eBGP (interior and exterior)

CS519



## iBGP and eBGP (interior and exterior)

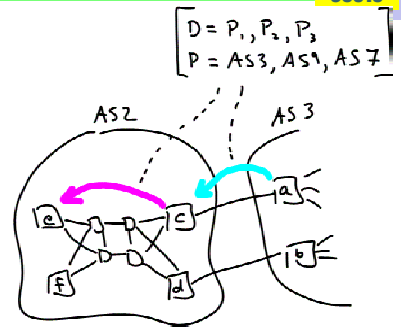
CS519



## Next hops are weird in iBGP

CS519

What does it mean for e to have c as its next hop to P1, when there are multiple routers between e and c???



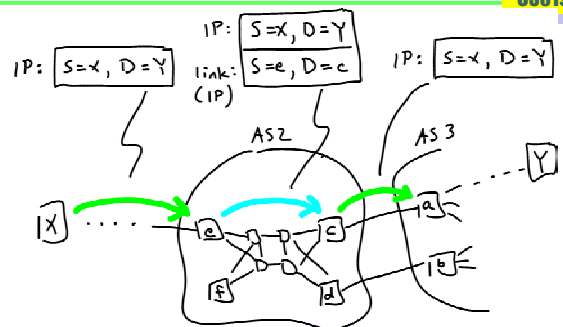
## One option: tunnel across AS

CS519

- iBGP speakers form IP tunnels across the AS
  - IP over IP
    - (perhaps with GRE between them, but lets not get into this now)
  - This creates a "link" between the two iBGP speakers
  - Remember, IP doesn't care what subnet technology it runs over, even if that subnet is IP!!!

## iBGP next hop using IP-in-IP tunneling

CS519



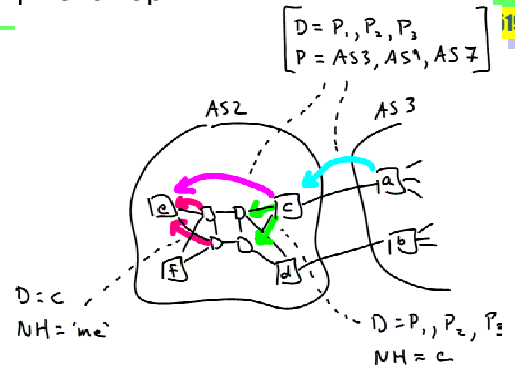
## Another option, use both BGP and igp RIBs

CS519

- o iBGP “resolves” the iBGP next hop to its igp next hop
  - iBGP computes its next hop
  - iBGP looks into igp RIB to determine igp next hop to iBGP next hop
  - This becomes the actual next hop
- o iBGP must advertise external prefixes into the igp

## iBGP using igp RIB to resolve next hop

19



## BGP security model

CS519

- o Authentication is hop-by-hop, like OSPF
- o But threat is much worse, because no single organization controls all of BGP
- o So, BGP uses policy to help prevent bogus routes
  - BGP routers have an expectation of what they should hear from where

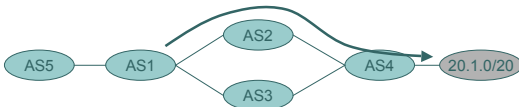
## BGP policies

CS519

- o Who to peer with (which ASs)
- o What routes to originate
- o What routes to import (prevent bogus advertisements)
- o What routes to export (and how to aggregate them)
- o What paths to prefer
  - Shorter AS paths
  - Some ASs preferred over others
    - The big ASs (UUnet, AT&T, etc.)
    - Primary versus backup transit AS

## BGP policy limitation (hop by hop policy decisions)

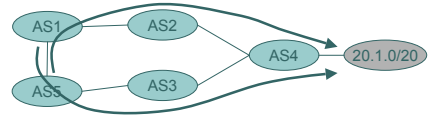
CS519



AS1 chooses AS2 as the path to 20.1.0/20.  
 AS5 is forced to accept the choice of AS1  
 (If AS5 really doesn't like it, it should find a new peer)

## BGP policy conflict

CS519

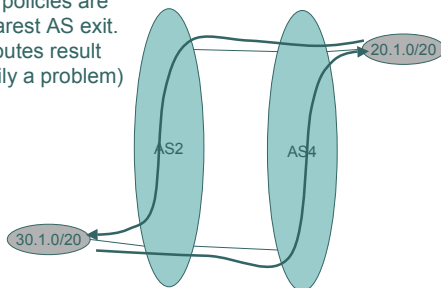


AS5 policy is to prefer route to AS4 via AS2  
 AS1 policy is to prefer route to AS4 via AS3  
Both policies cannot be satisfied

## Hot potato routing

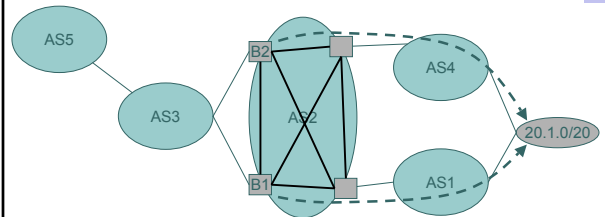
CS519

AS2 and AS4 policies are to route to nearest AS exit.  
 Asymmetric routes result  
 (not necessarily a problem)



## Misconfigured policies may lead to oscillation

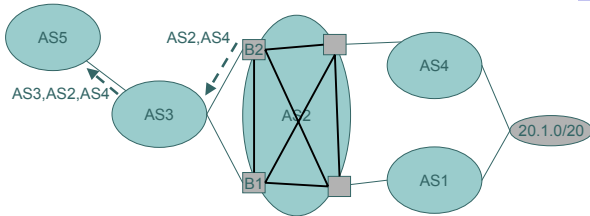
CS519



B2 configured to prefer AS4  
 B1 configured to prefer AS1

## Misconfigured policies may lead to oscillation

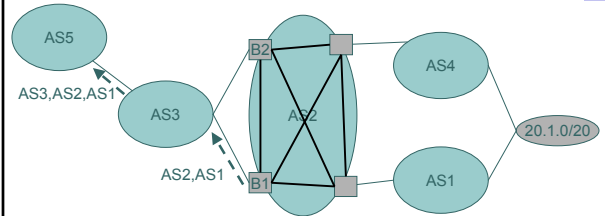
CS519



B2 (periodically) updates AS3 with path AS2,AS4

## Misconfigured policies may lead to oscillation

CS519



B1 (periodically) updates AS3 with path AS2,AS1  
With each period AS3 advertises a different route

## Other route flapping

CS519

- A link continuously goes up and down
  - The update for this is propagated throughout the internet
- Mid-90's these kinds of problems were severe
  - 1996: 45,000 prefixes, 1,500 unique AS paths, 1,300 ASs, 3-6 million BGP update messages/day
    - 6 updates per prefix per hour!
    - (Labovitz et. al.)

## Today much improved

CS519

- Better policy tools
- Better software
- Lots of damping
- But still, advances in BGP lead to new policy bugs
  - Route reflectors published in 2000 (RFC2796)
  - Inconsistent route reflectors problem published in 2002 (RFC3345)

## Policy Tools

CS519

- Routing Policy Specification Language (RPSL) (RFC 2280)
  - Earlier policy languages exist
- Language to define BGP policies
  - Peers, import, export, route preference, aggregation
- Posted at Routing Registries (RIPE, RADB, etc.)
- Tools created to look for policy inconsistencies (within AS and across ASs)
- Tools created to match measured reality (BGP tables, traceroute) with policy expectations
  - RAToolSet, USC/ISI

## Lots of Damping

CS519

- Stop advertising certain prefixes if they go up and down a lot
  - Improve stability
  - Lower overhead
- RIPE guidelines:
  - Don't dampen until after 4<sup>th</sup> flap in a row (in 50 minutes)
  - /24: dampen 60 minutes
  - /22,/23, dampen 30-45 minutes
  - </22, dampen 10-30 minutes

## Lots of Damping

CS519

- Helps the internet, but means that you can go away for a long time
  - Because of some problem in the middle!
- Most damping is done on routes that you don't care about
  - Poorly managed small ISPs
- Routes through major ISPs tend to be very stable
  - Your favorite web sites

## Effect of BGP policies on path quality

CS519

- Ramesh Govindan study (USC)
- Methodology:
  - Learn real physical topology with traceroutes, deduce actual AS connectivity
    - Imperfect, but not bad
  - Examine used "policy topology" from BGP tables, RADB (routing registry) database
  - Compare the two



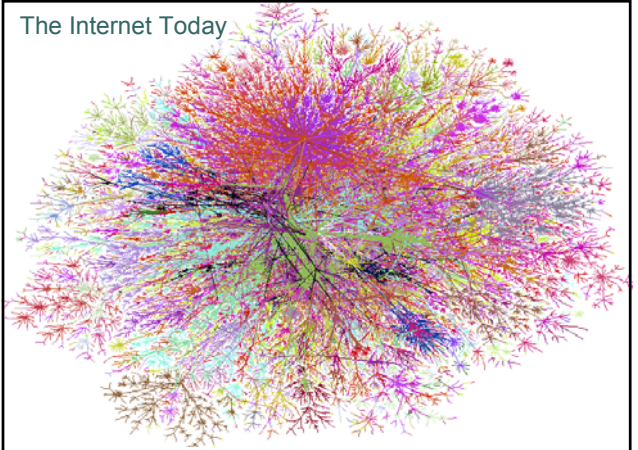
## Effect of BGP policies on path quality

CS519

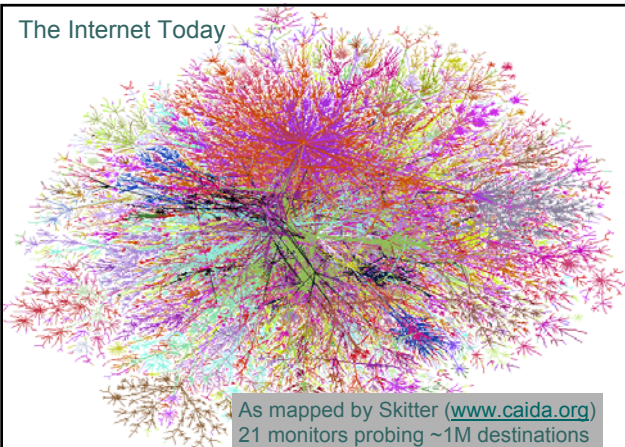
### Results:

- About  $\frac{1}{2}$  of the paths are longer than shortest path
- 20% of policy paths are 50% or more longer
- 20% of policy paths are 5 hops or more longer
- Policy tends to push paths through major backbones rather than possibly shorter routes
  - (But shorter routes may not be better routes!)

## The Internet Today



## The Internet Today



## BGP Routing Table Growth

519

