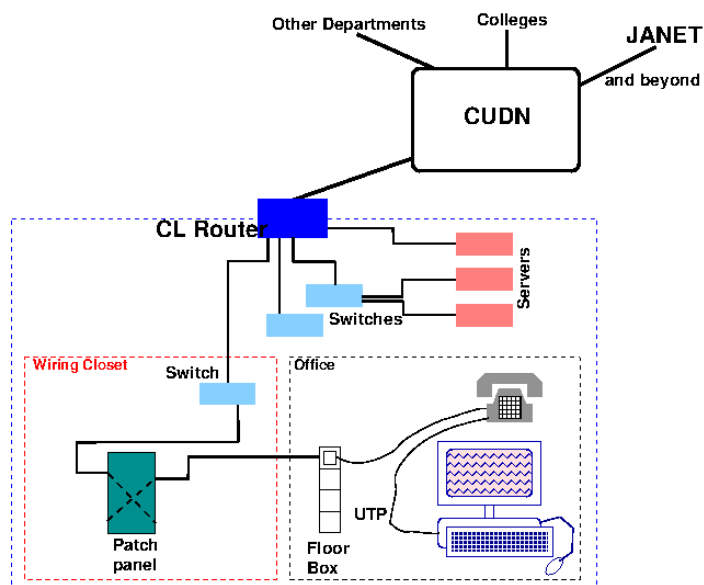


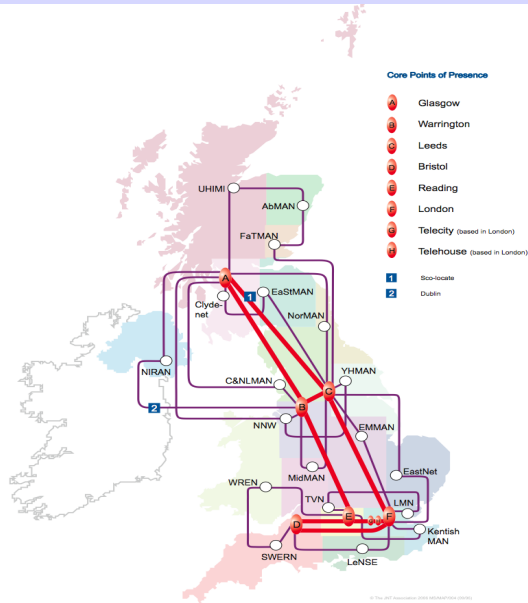
L11 : Interdomain Routing Lectures 8 & 9 November 2014

Timothy G. Griffin
Computer Lab
Cambridge UK

CL Network

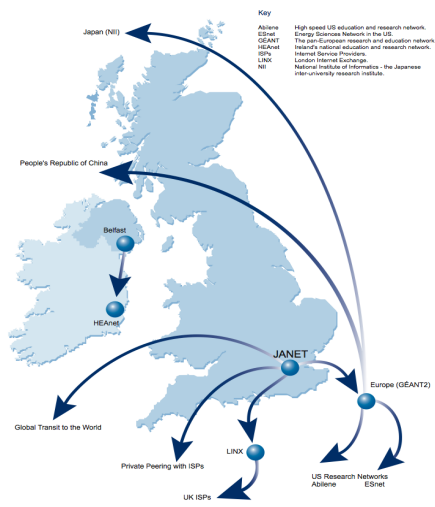


JANET

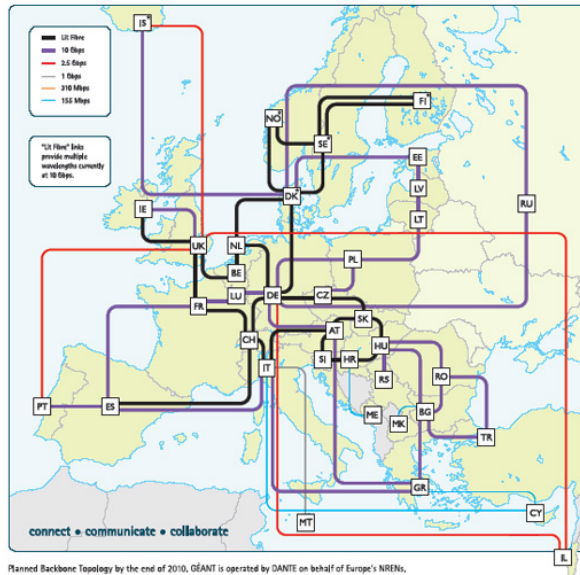


JANET in the Internet

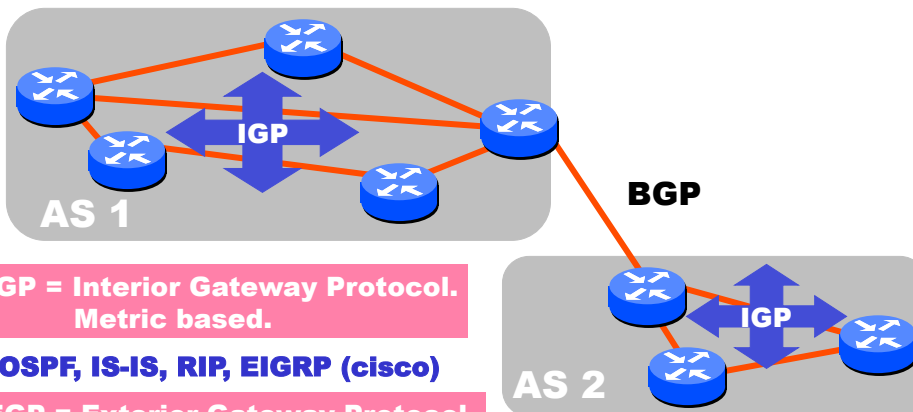
JANET External Network Access Provision



Geant = EU Research and Education Backbone



Architecture of Dynamic Routing



IGP = Interior Gateway Protocol.
Metric based.

OSPF, IS-IS, RIP, EIGRP (cisco)

EGP = Exterior Gateway Protocol.
Policy Based.

Only one: BGP

The Routing Domain of BGP is the entire Internet

Technology of Distributed Routing

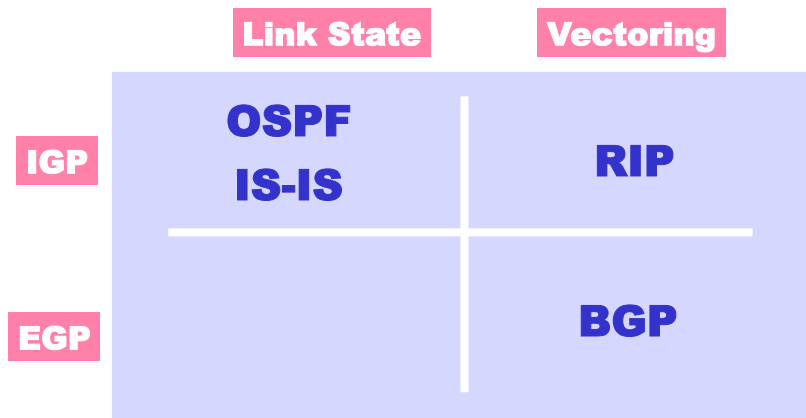
Link State

- Topology information is flooded within the routing domain
- Best end-to-end paths are computed locally at each router.
- Best end-to-end paths determine next-hops.
- Based on minimizing some notion of distance
- Works only if policy is shared and uniform
- Examples: OSPF, IS-IS

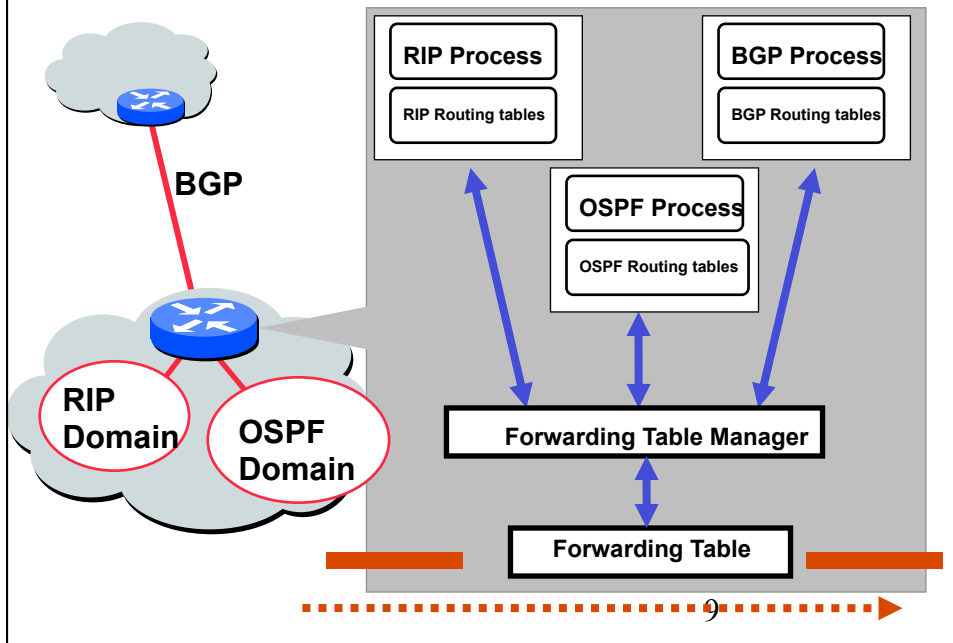
Vectoring

- Each router knows little about network topology
- Only best next-hops are chosen by each router for each destination network.
- Best end-to-end paths result from composition of all next-hop choices
- Does not require any notion of distance
- Does not require uniform policies at all routers
- Examples: RIP, BGP

The Gang of Four



Happy Packets: The Internet Does Not Exist Only to Populated Routing Tables



Autonomous Routing Domains

A collection of physical networks glued together using IP, that have a unified administrative routing policy.

- **Campus networks**
- **Corporate networks**
- **ISP Internal networks**
- ...

Autonomous Systems (ASes)

An autonomous system is an autonomous routing domain that has been assigned an Autonomous System Number (ASN).

... the administration of an AS appears to other ASes to have a single coherent interior routing plan and presents a consistent picture of what networks are reachable through it.

RFC 1930: Guidelines for creation, selection, and registration of an Autonomous System

AS Numbers (ASNs)

ASNs are 16 bit values (soon to be 32 bits)

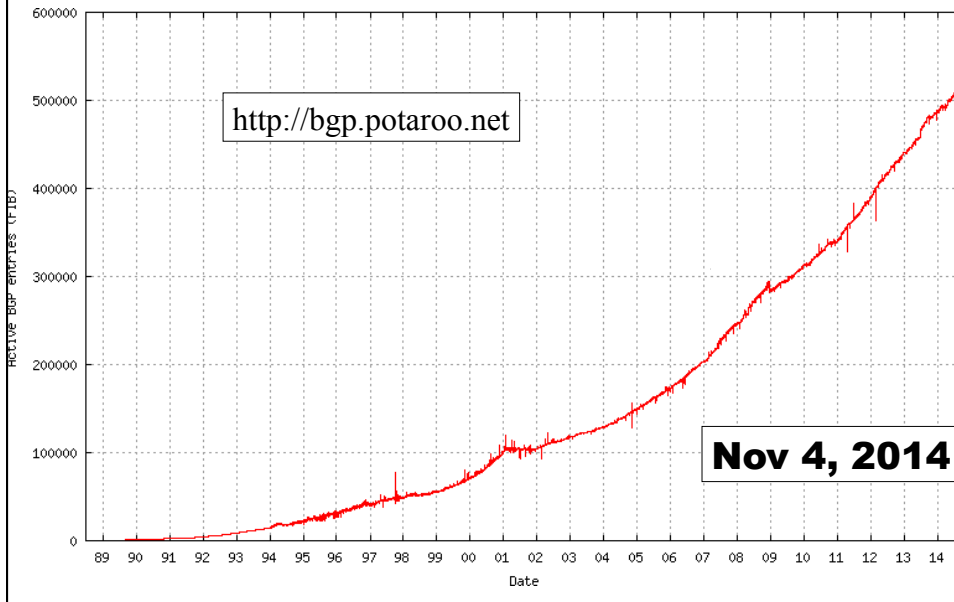
64512 through 65535 are “private”

Currently nearly 30,000 in use.

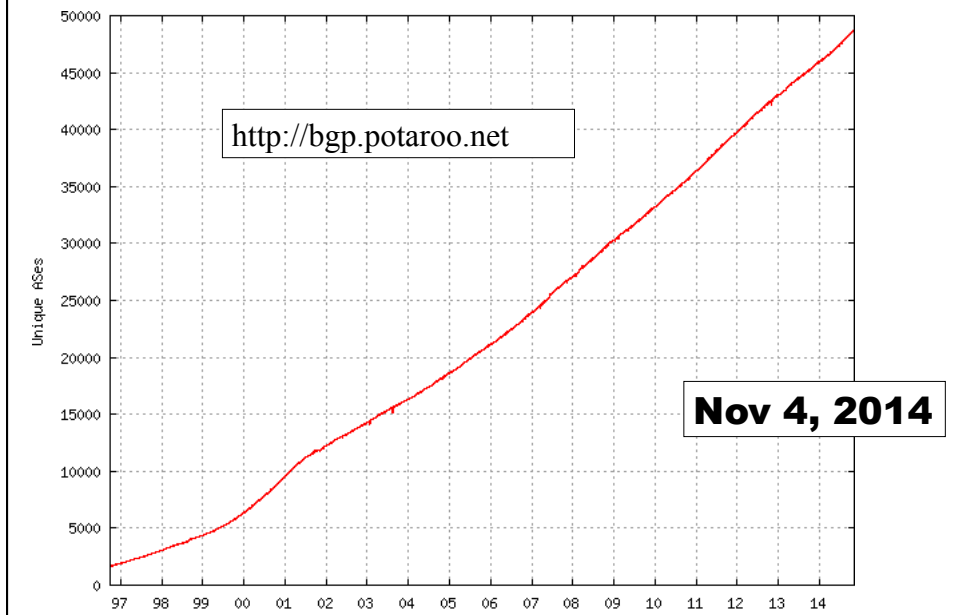
- **JANET: 786**
- **MIT: 3**
- **Harvard: 11**
- **UC San Diego: 7377**
- **AT&T: 7018, 6341, 5074, ...**
- **UUNET: 701, 702, 284, 12199, ...**
- **Sprint: 1239, 1240, 6211, 6242, ...**
- ...

ASNs represent units of routing policy

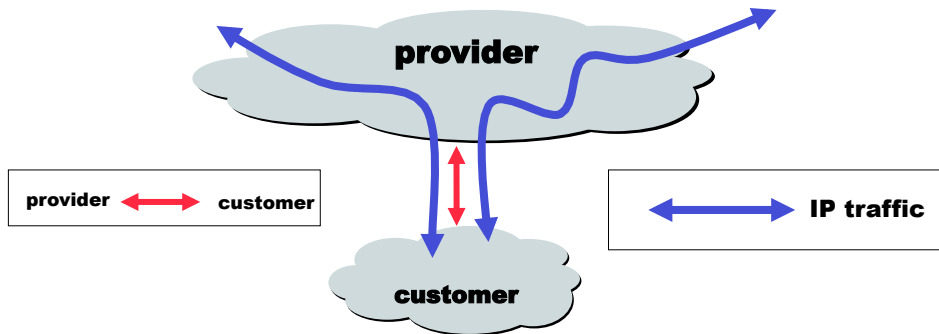
How many prefixes are used today?



How many ASNs are used today?

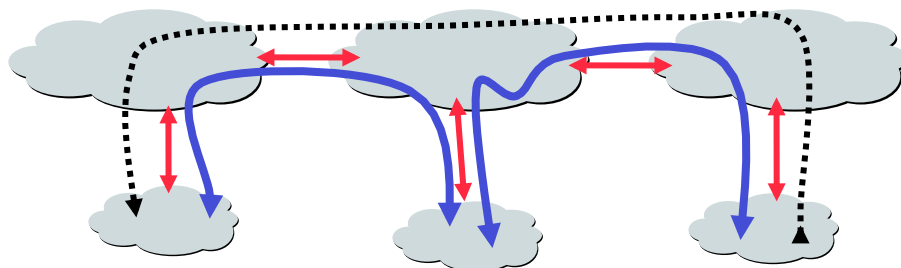


Customers and Providers



Customer pays provider for access to the Internet

The "Peering" Relationship



peer ↔ peer
provider ↔ customer

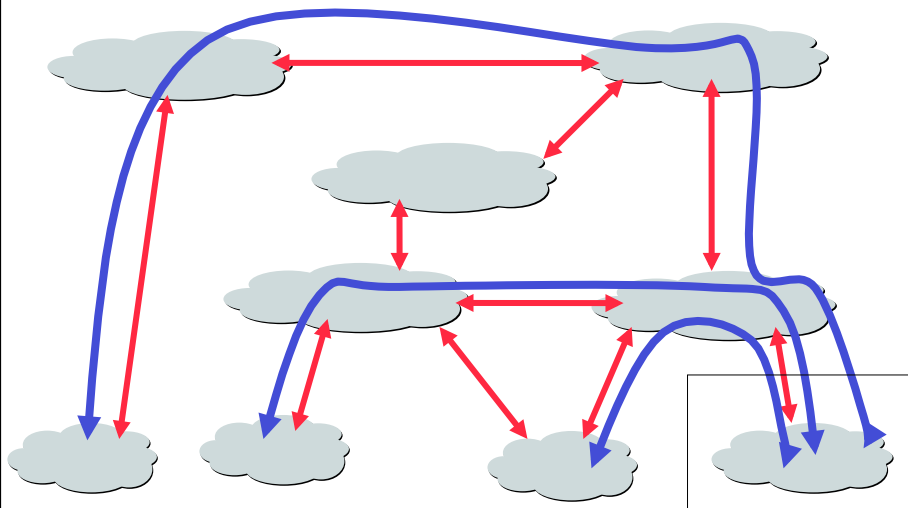
↔ traffic allowed
- - - - - traffic NOT allowed

Peers provide transit between their respective customers

Peers do not provide transit between peers

Peers (often) do not exchange \$\$\$

Peering Provides Shortcuts



Peering also allows connectivity between the customers of "Tier 1" providers.



Peering Wars

Peer

- Reduces upstream transit costs
- Can increase end-to-end performance
- May be the only way to connect your customers to some part of the Internet ("Tier 1")

Don't Peer

- You would rather have customers
- Peers are usually your competition
- Peering relationships may require periodic renegotiation

Peering struggles are by far the most contentious issues in the ISP world!

Peering agreements are often confidential.

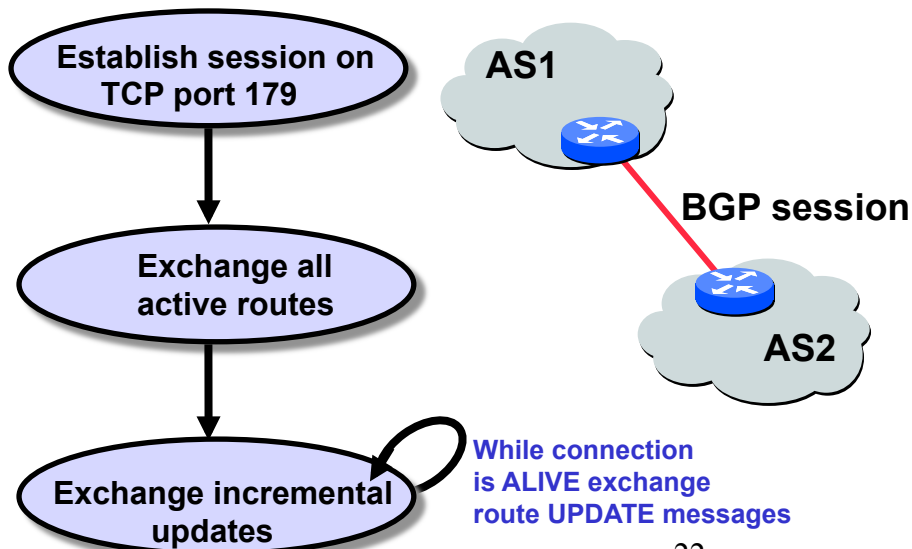
BGP-4

- **BGP = Border Gateway Protocol**
- Is a **Policy-Based** routing protocol
- Is the **de facto EGP** of today's global Internet
- Relatively simple protocol, but configuration is complex and the entire world can see, and be impacted by, your mistakes.

- **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)
- **2006 : BGP-4 [RFC 4271]**

21

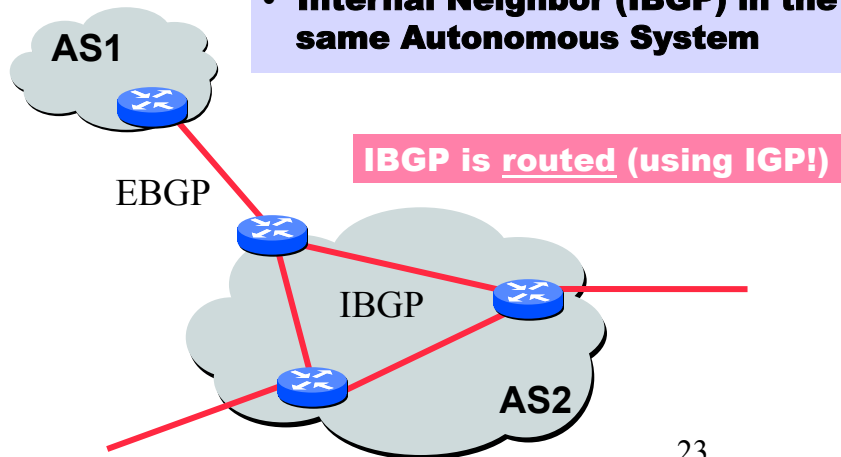
BGP Operations (Simplified)



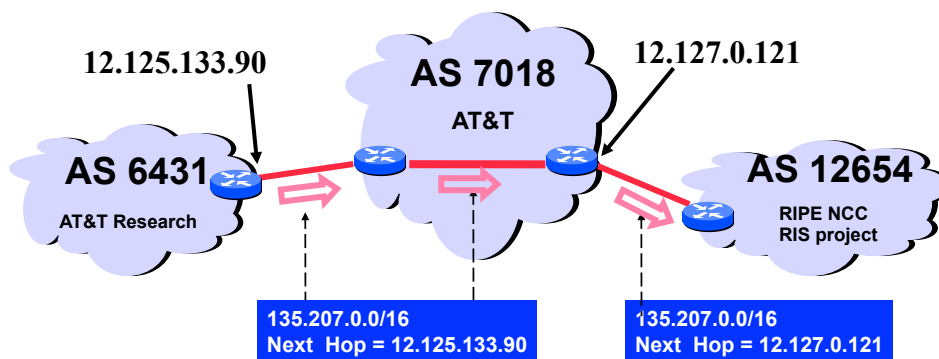
22

Two Types of BGP Sessions

- **External Neighbor (EBGP) in a different Autonomous Systems**
- **Internal Neighbor (IBGP) in the same Autonomous System**

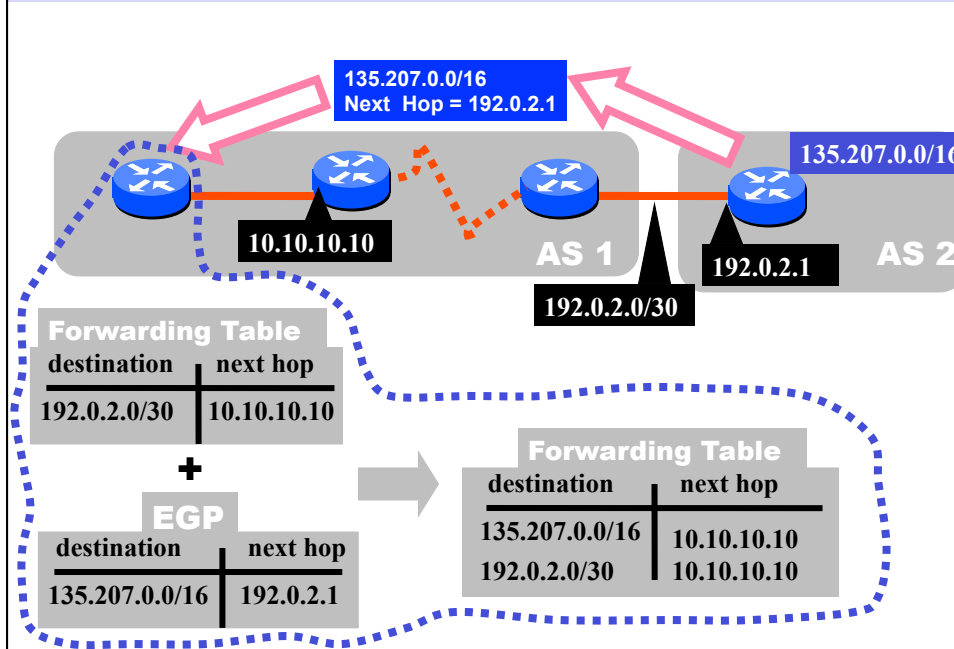


BGP Next Hop Attribute



Every time a route announcement crosses an AS boundary, the Next Hop attribute is changed to the IP address of the border router that announced the route.

Join EGP with IGP For Connectivity



Four Types of BGP Messages

- **Open** : Establish a peering session.
- **Keep Alive** : Handshake at regular intervals.
- **Notification** : Shuts down a peering session.
- **Update** : Announcing new routes or withdrawing previously announced routes.

announcement
=
prefix + attributes values

BGP Attributes

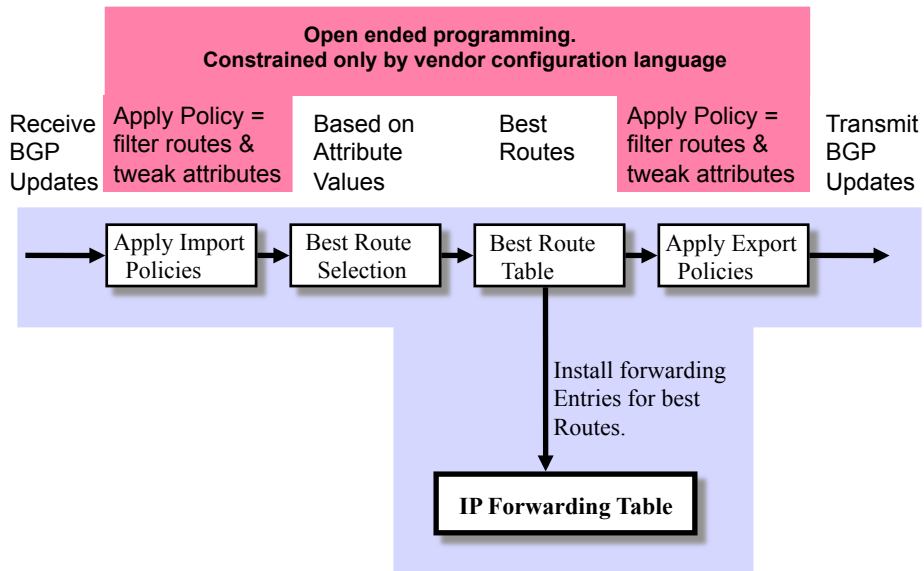
Code	Reference
1 ORIGIN	[RFC1771]
2 AS_PATH	[RFC1771]
3 NEXT_HOP	[RFC1771]
4 MULTI_EXIT_DISC	[RFC1771]
5 LOCAL_PREF	[RFC1771]
6 ATOMIC_AGGREGATE	[RFC1771]
7 AGGREGATOR	[RFC1771]
8 COMMUNITY	[RFC1997]
9 ORIGINATOR_ID	[RFC2796]
0 CLUSTER_LIST	[RFC2796]
1 DPA	[Chen]
2 ADVERTISER	[RFC1863]
3 RCID_PATH / CLUSTER_ID	[RFC1863]
4 MP_REACH_NLRI	[RFC2283]
5 MP_UNREACH_NLRI	[RFC2283]
6 EXTENDED COMMUNITIES	[Rosen]
7-255 reserved for development	

Most important attributes

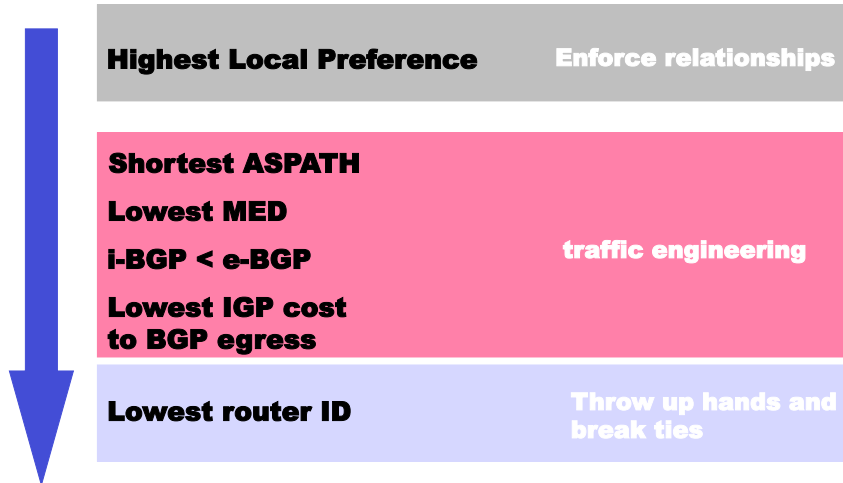
From IANA: <http://www.iana.org/assignments/bgp-parameters>

Not all attributes need to be present in every announcement

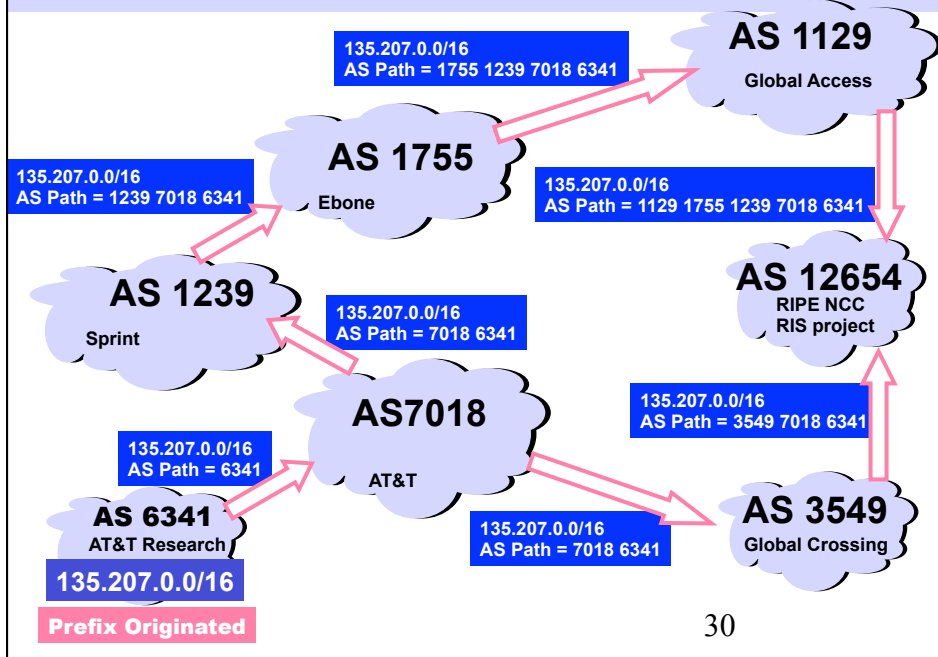
BGP Route Processing



Route Selection Summary

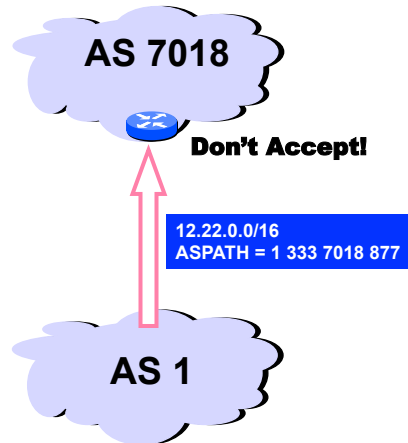


AS PATH Attribute



Interdomain Loop Prevention

BGP at AS YYY will never accept a route with ASPATH containing YYY.

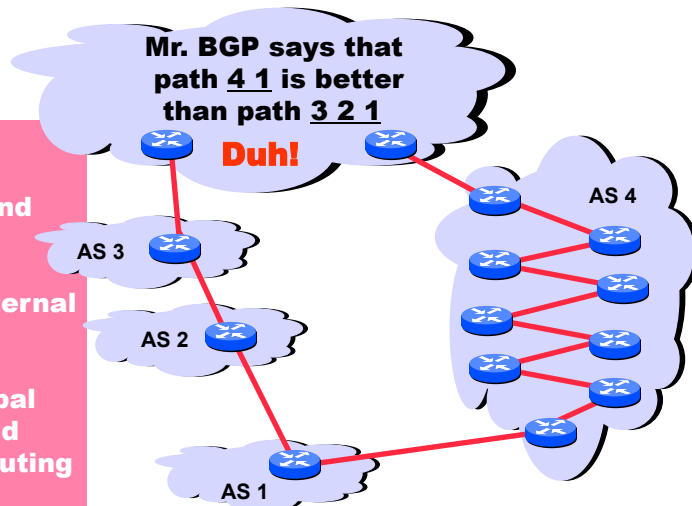


31

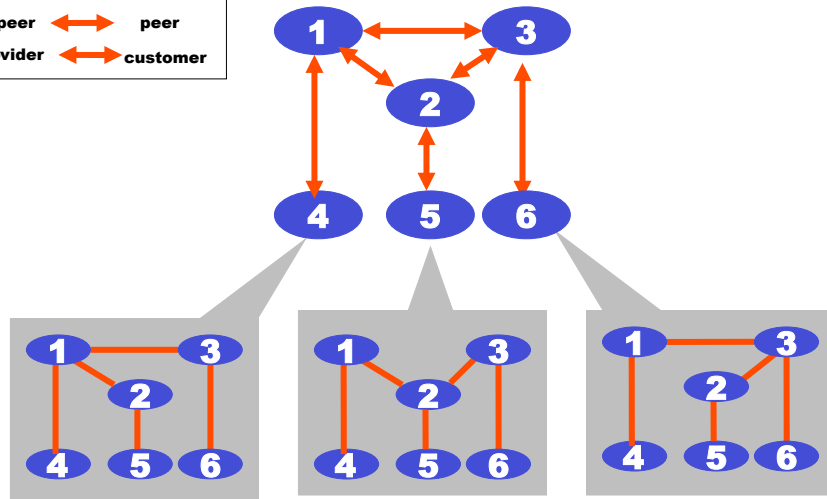
Shorter Doesn't Always Mean Shorter

In fairness: could you do this "right" and still scale?

Exporting internal state would dramatically increase global instability and amount of routing state

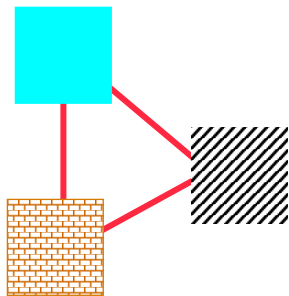


AS Graphs Depend on Point of View

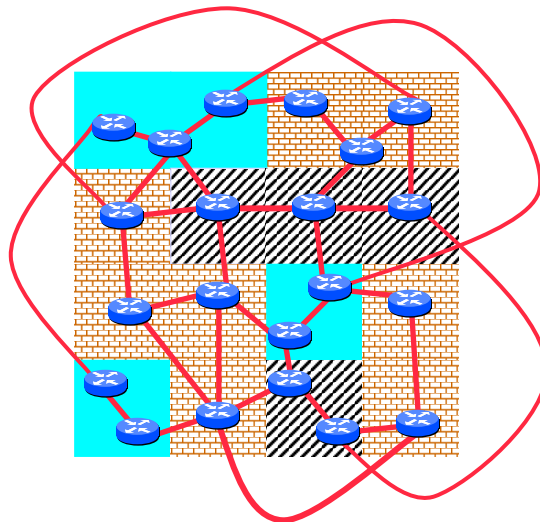


AS Graphs Do Not Show “Topology”!

BGP was designed to throw away information!



The AS graph may look like this.



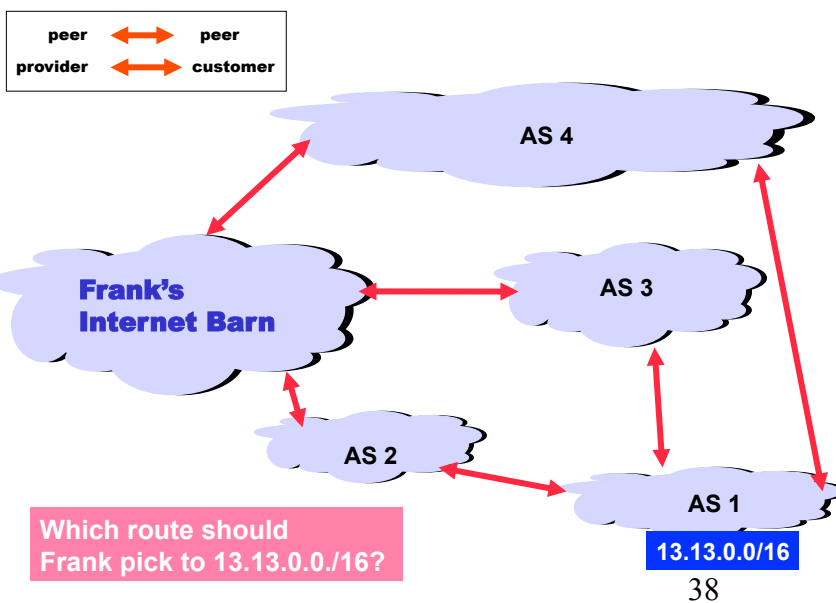
Reality may be closer to this...

Implementing Customer/Provider and Peer/Peer relationships

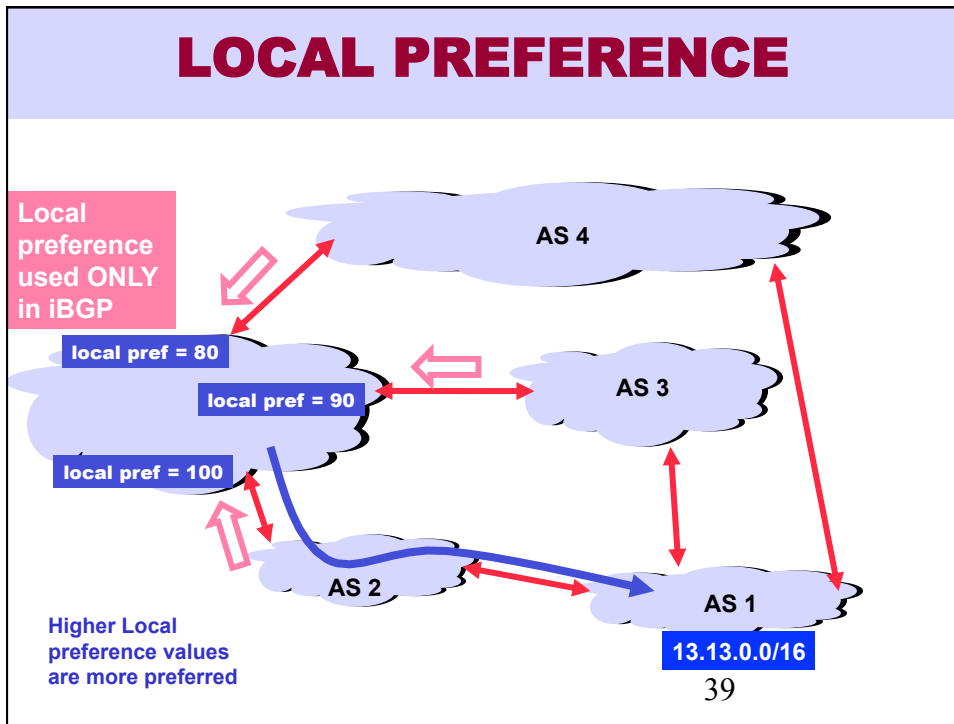
Two parts:

- Enforce transit relationships
 - Export all (best) routes to customers
 - Send only own and customer routes to all others
- Enforce order of route preference
 - provider < peer < customer

So Many Choices

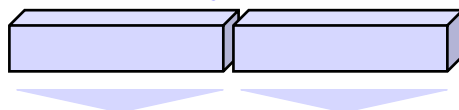


LOCAL PREFERENCE



How Can Routes be Classified? BGP Communities!

A community value is 32 bits



By convention, first 16 bits is ASN indicating who is giving it an interpretation

community number

Used for signaling within and between ASes

Very powerful BECAUSE it has no (predefined) meaning

**Community Attribute = a list of community values.
(So one route can belong to multiple communities)**

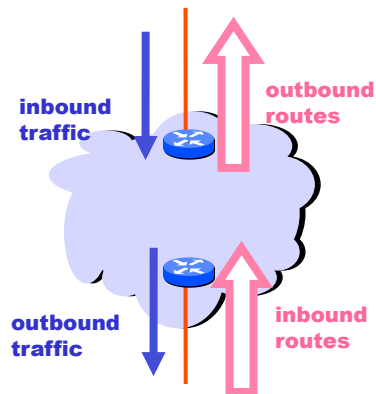
Reserved communities

no_export = 0xFFFFF01: don't export out of AS
no_advertise 0xFFFFF02: don't pass to BGP neighbors

RFC 1997 (August 1996)

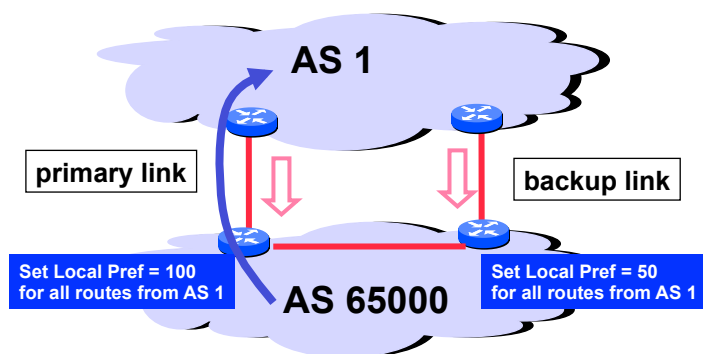
Tweak Tweak Tweak (TE)

- For inbound traffic
 - Filter outbound routes
 - Tweak attributes on outbound routes in the hope of influencing your neighbor's best route selection
- For outbound traffic
 - Filter inbound routes
 - Tweak attributes on inbound routes to influence best route selection



In general, an AS has more control over outbound traffic

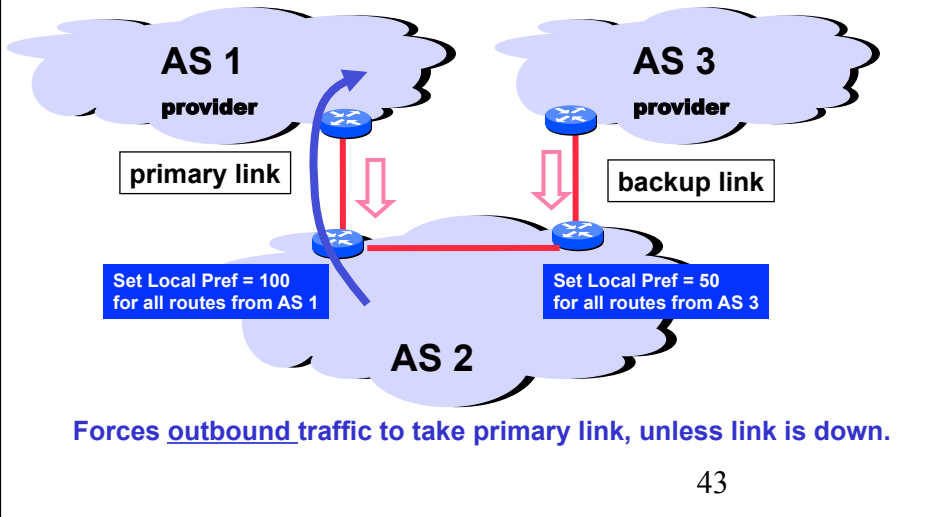
Implementing Backup Links with Local Preference (Outbound Traffic)



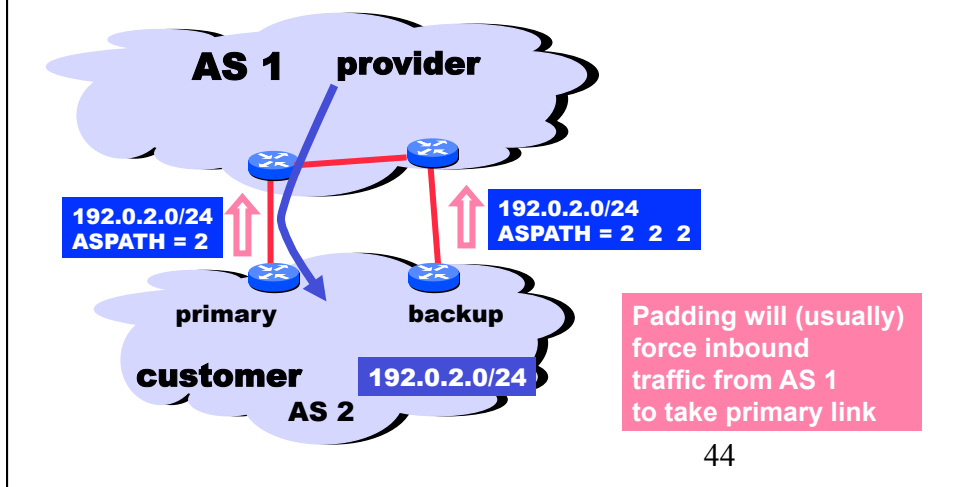
Forces outbound traffic to take primary link, unless link is down.

We'll talk about inbound traffic soon ...

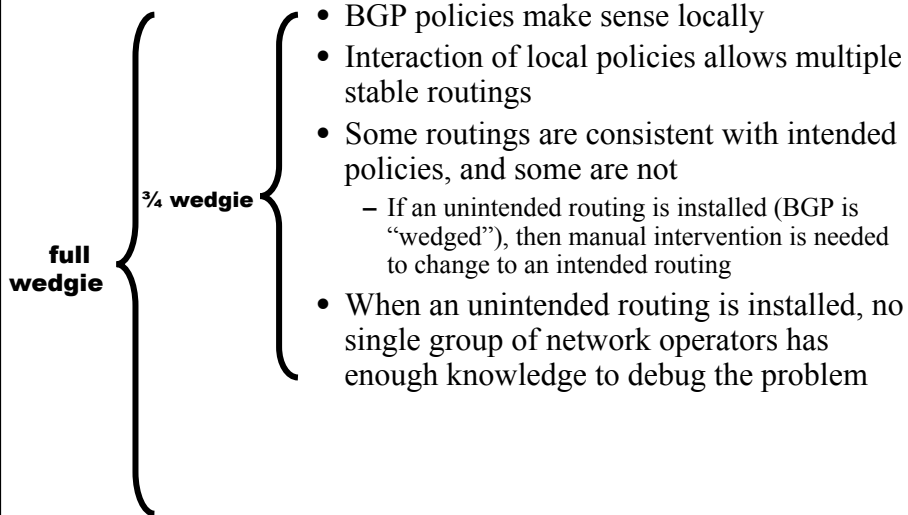
Multihomed Backups (Outbound Traffic)



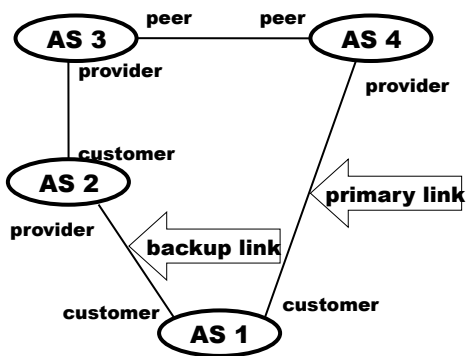
Shedding Inbound Traffic with ASPATH Padding. Yes, this is a Glorious Hack ...



What is a BGP Wedgie?

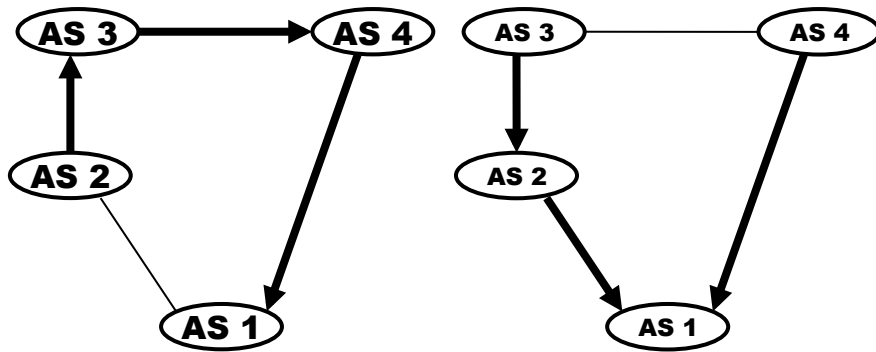


$\frac{3}{4}$ Wedgie Example



- AS 1 implements backup link by sending AS 2 a “depref me” community.
- AS 2 implements this community so that the resulting local pref is below that of routes from its upstream provider (AS 3 routes)

And the Routings are...



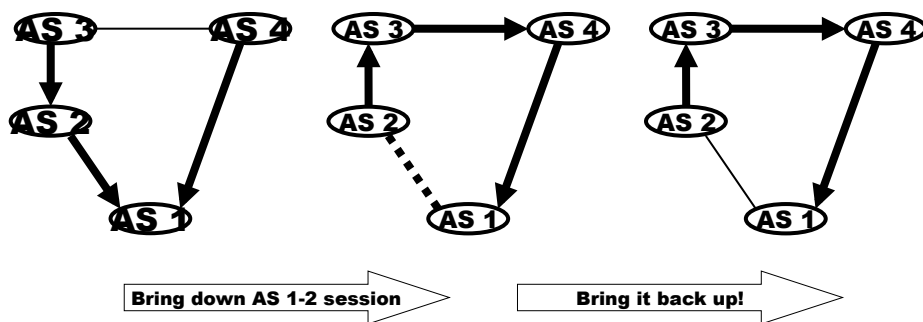
Intended Routing

Note: this would be the **ONLY** routing if AS2 translated its "depref me" community to a "depref me" community of AS 3

Unintended Routing

Note: This is easy to reach from the intended routing just by "bouncing" the BGP session on the primary link.

Recovery

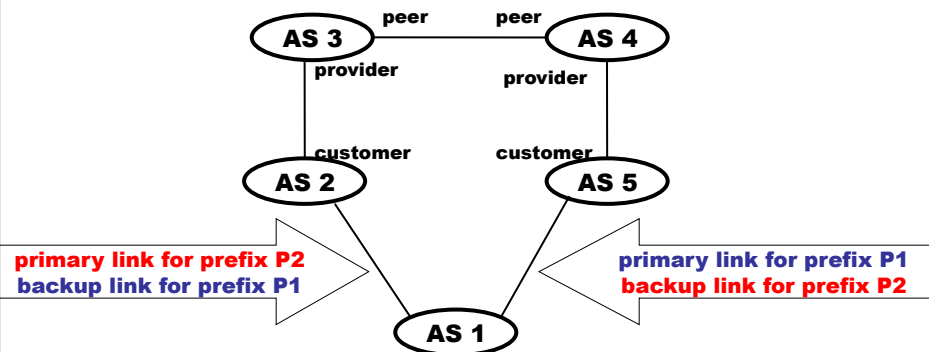


- Requires manual intervention
- Can be done in AS 1 or AS 2

What the heck is going on?

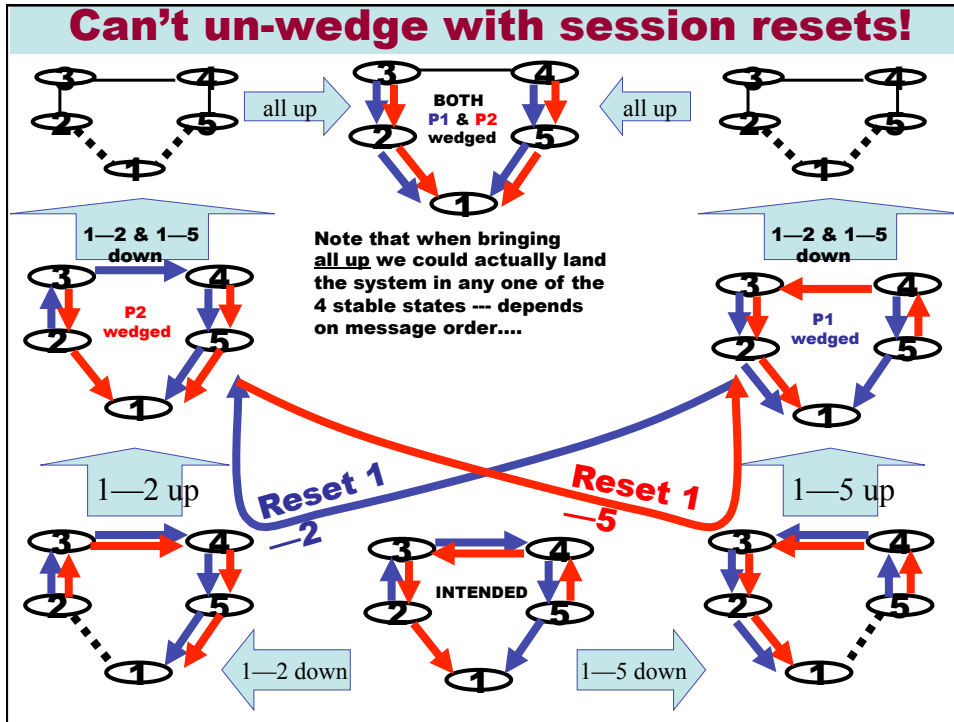
- There is no guarantee that a BGP configuration has a unique routing solution.
 - When multiple solutions exist, the (unpredictable) order of updates will determine which one is wins.
- There is no guarantee that a BGP configuration has any solution!
 - And checking configurations NP-Complete
 - Lab demonstrations of BGP configs never converging
- Complex policies (weights, communities setting preferences, and so on) increase chances of routing anomalies.
 - ... yet this is the current trend!

Load Balancing Example

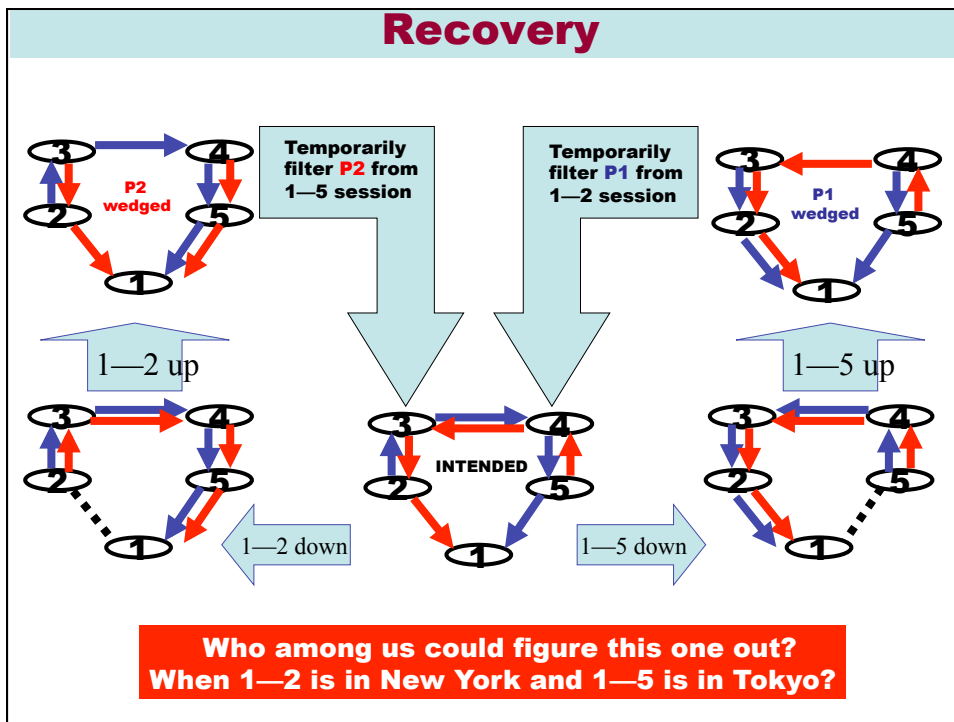


Simple session reset my not work!!

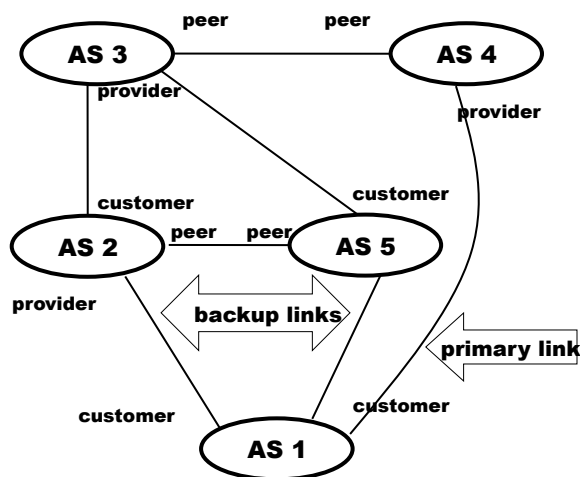
Can't un-wedge with session resets!



Recovery

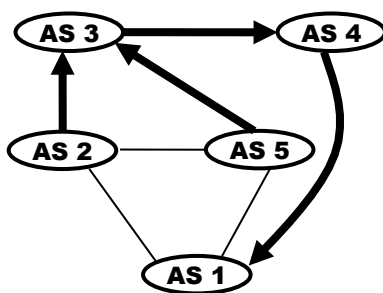


Full Wedgie Example

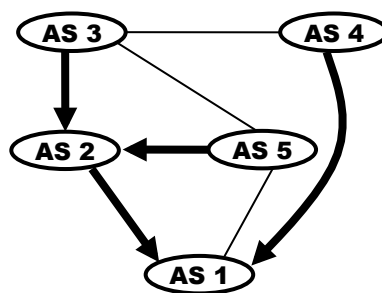


- AS 1 implements backup links by sending AS 2 and AS 3 a “depref me” communities.
- AS 2 implements its community so that the resulting local pref is below that of its upstream providers and its peers (AS 3 and AS 5 routes)
- AS 5 implements its community so that the resulting local pref is below its peers (AS 2) but above that of its providers (AS 3)

And the Routings are...

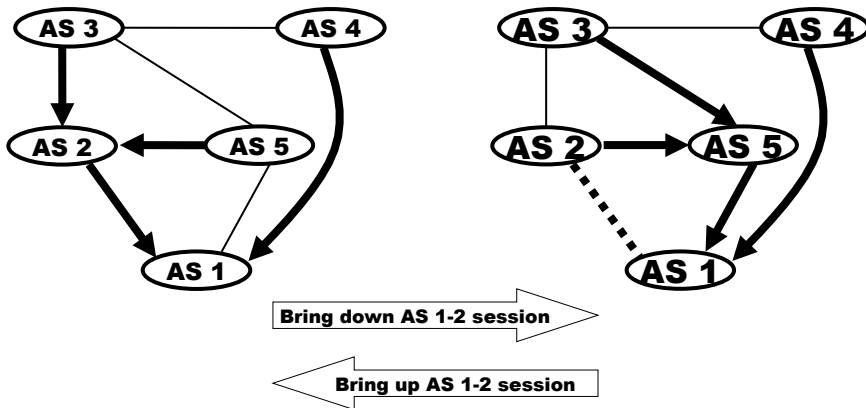


Intended Routing

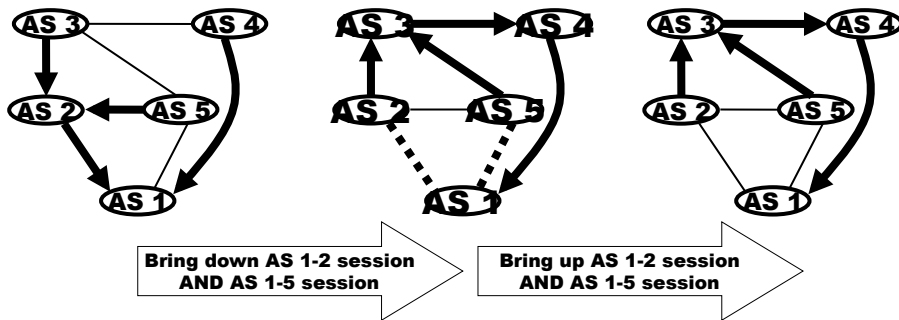


Unintended Routing

Resetting 1—2 does not help!!



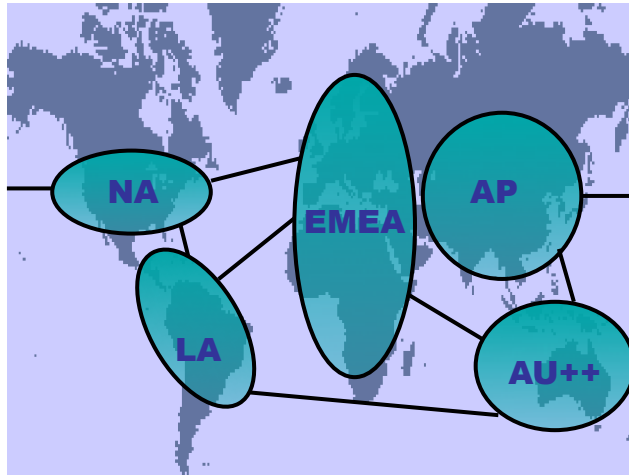
Recovery



A lot of "non-local" knowledge is required to arrive at this recovery strategy!

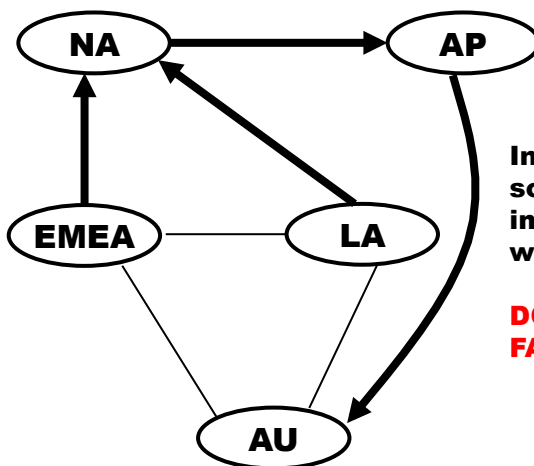
Try to convince AS 5 and AS 1 that their session has been reset (or filtered) even though it is not associated with an active route!

That Can't happen in MY network!!



An "normal" global global backbone (ISP or Corporate Intranet) implemented with 5 regional ASes

The Full Wedgie Example, in a new Guise



Intended Routing for some prefixes in AU, implemented with communities.

DOES THIS LOOK FAMILIAR??

Message: Same problems can arise with "traffic engineering" across regional networks.