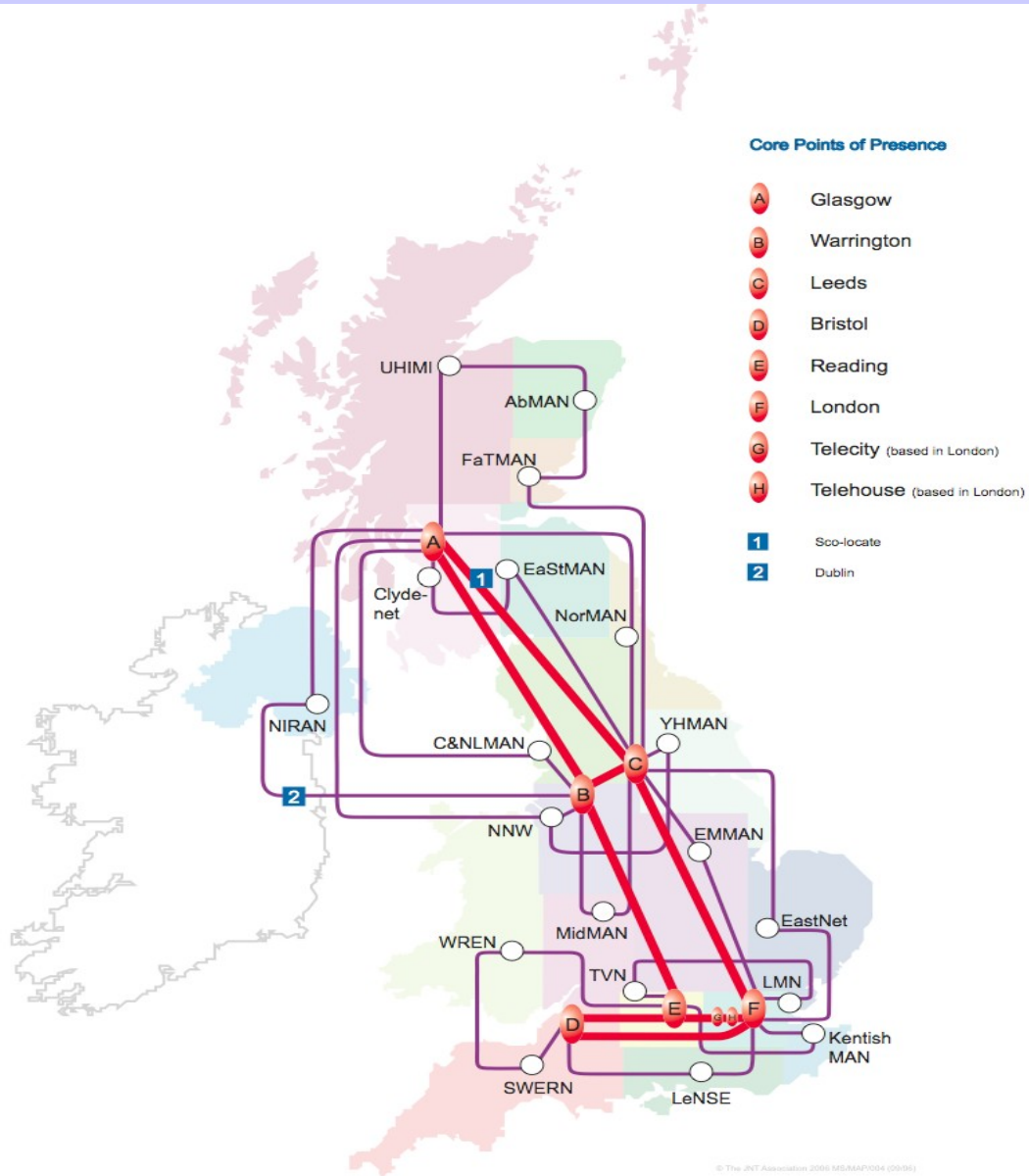


L11 : BGP

Lecture 13

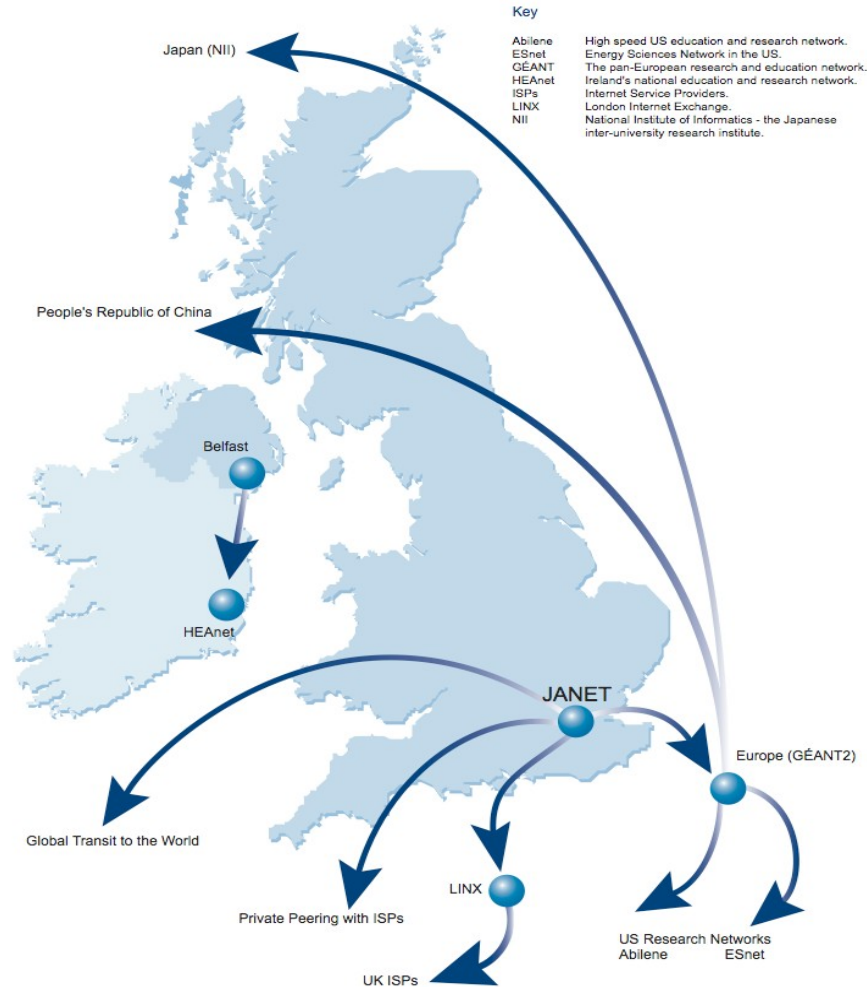
Timothy G. Griffin
Computer Lab
Cambridge UK

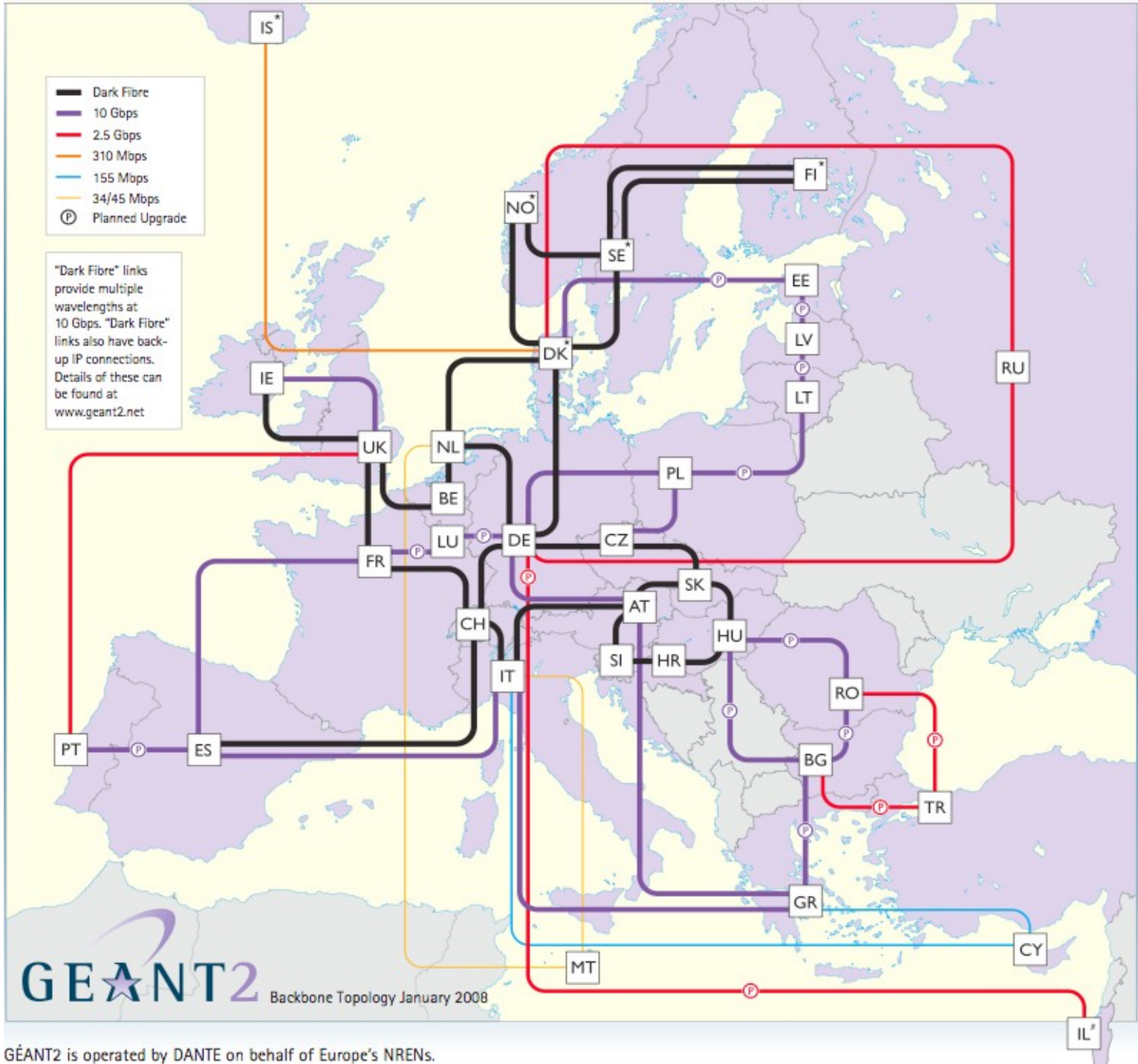
JANET



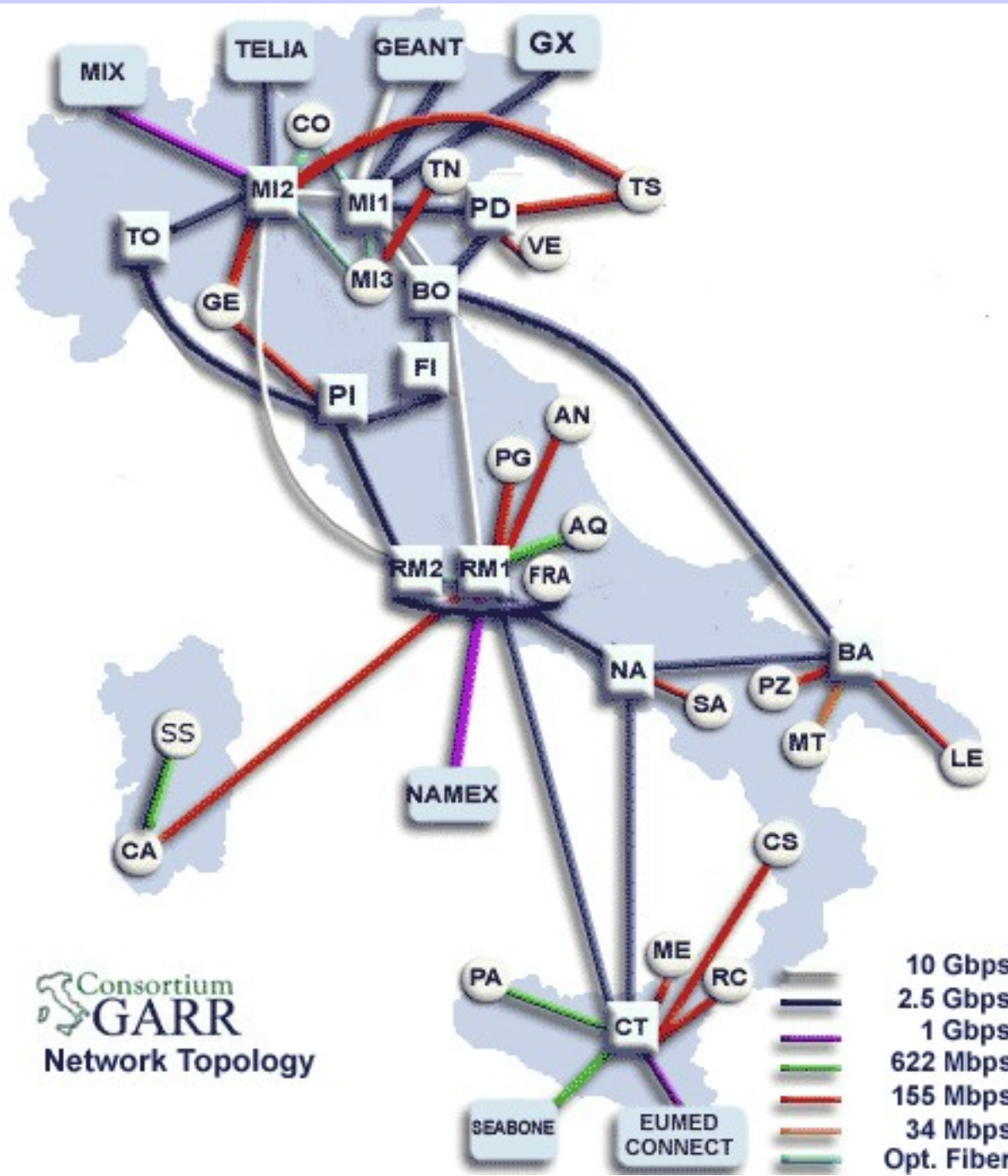
JANET and the Internet

JANET External Network Access Provision





GEANT2 is operated by DANTE on behalf of Europe's NRENs.



 Consortium
GARR
 Network Topology

RENATER-4 is deployed since september 2005



Réseau National de télécommunications
pour la technologie, l'enseignement et la Recherche



RENATER-4



Connexion à
l'Internet mondial

SFINX
Global Internet eXchange, accès aux autres
prestataires de service Internet en France

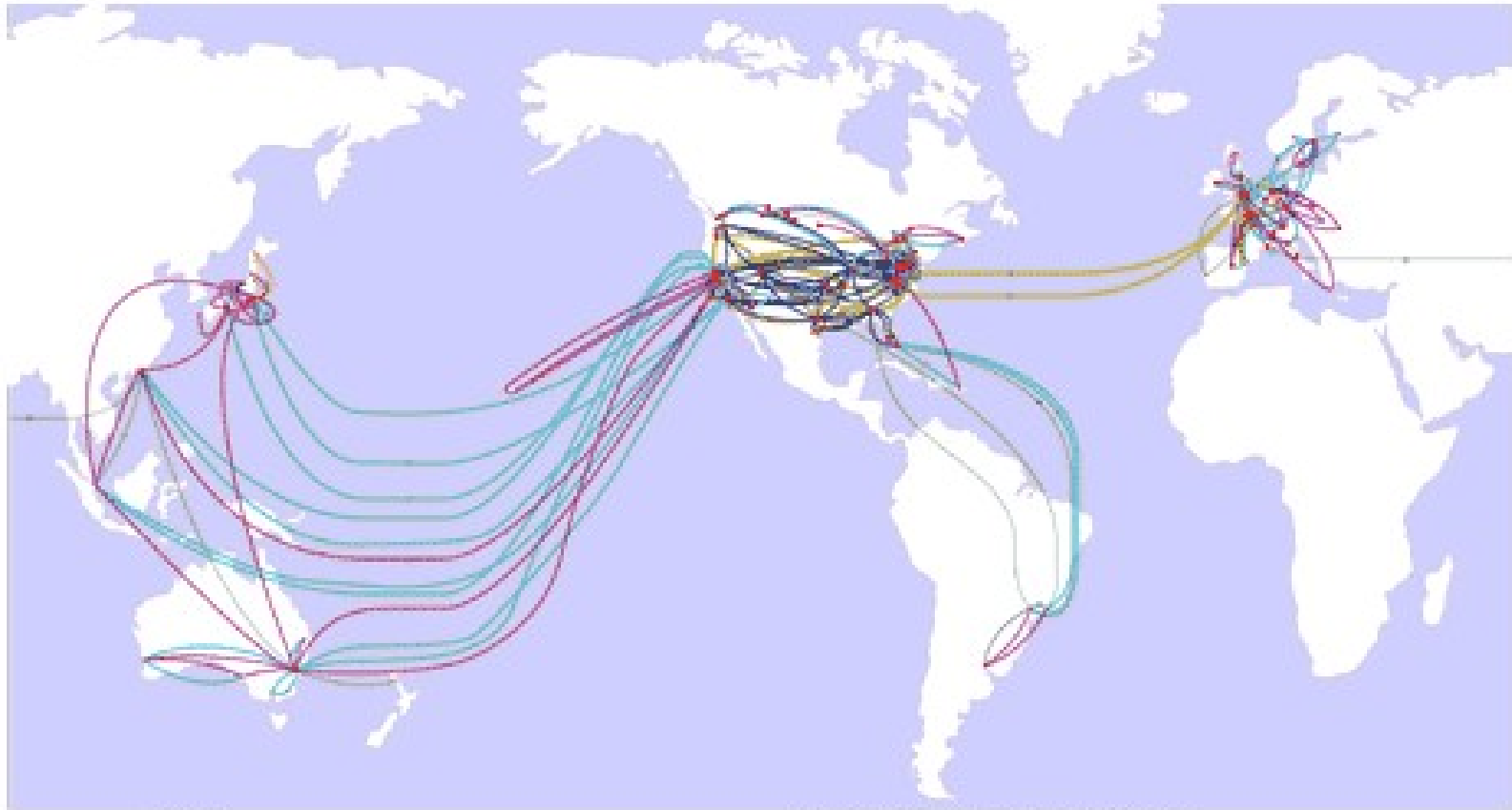
GEANT2 www.geant2.net
Connexion vers les réseaux
de la Recherche en Europe,
et les réseaux de la Recherche
des pays méditerranéens,
de la zone Asie Pacifique
de l'Amérique du sud
de l'Amérique centrale
CJRA



Connexion
vers les DOM-TOM

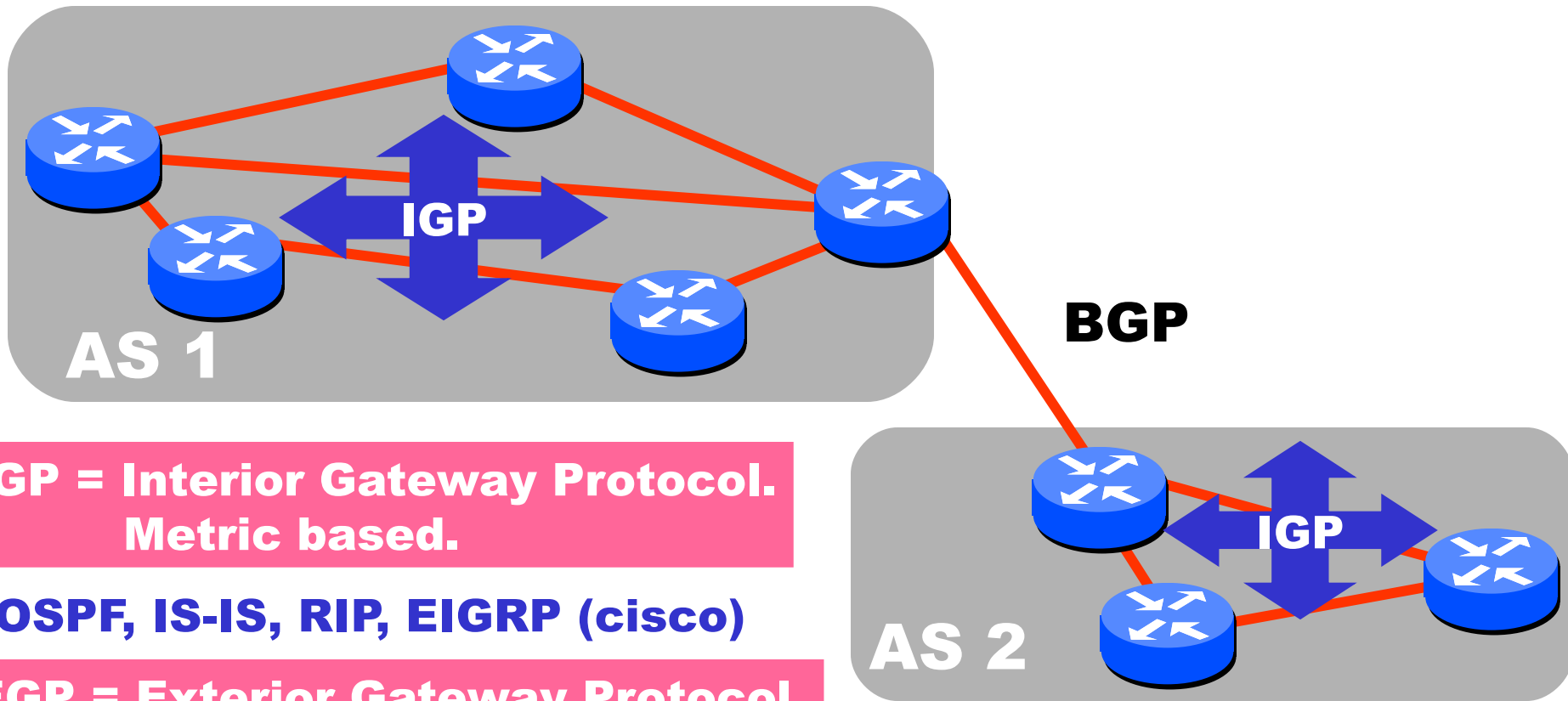
- Réseau en Ile de France
- 2.5 Gbit/s
- Liaisons projets de recherche
- Liaison projets à venir
- NR
- NRI

WorldCom (UUNet)



- | | |
|-------------------------------|--------------------------|
| — 64 Kbps | — OC12c/STM4 (622 Mbps) |
| — T1/E1 (1.5 Mbps/2 Mbps) | — OC48c/STM16 (2.5 Gbps) |
| — E3/T3/DS3 (35 Mbps/45 Mbps) | — OC192c/STM64 (10 Gbps) |
| — T2 (6 Mbps) | ● Single Hub City |
| — OC3c/STM1 (155 Mbps) | ■ Multiple Hubs City |
| | ● Data Center Hub |

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

Link State

Vectoring

IGP

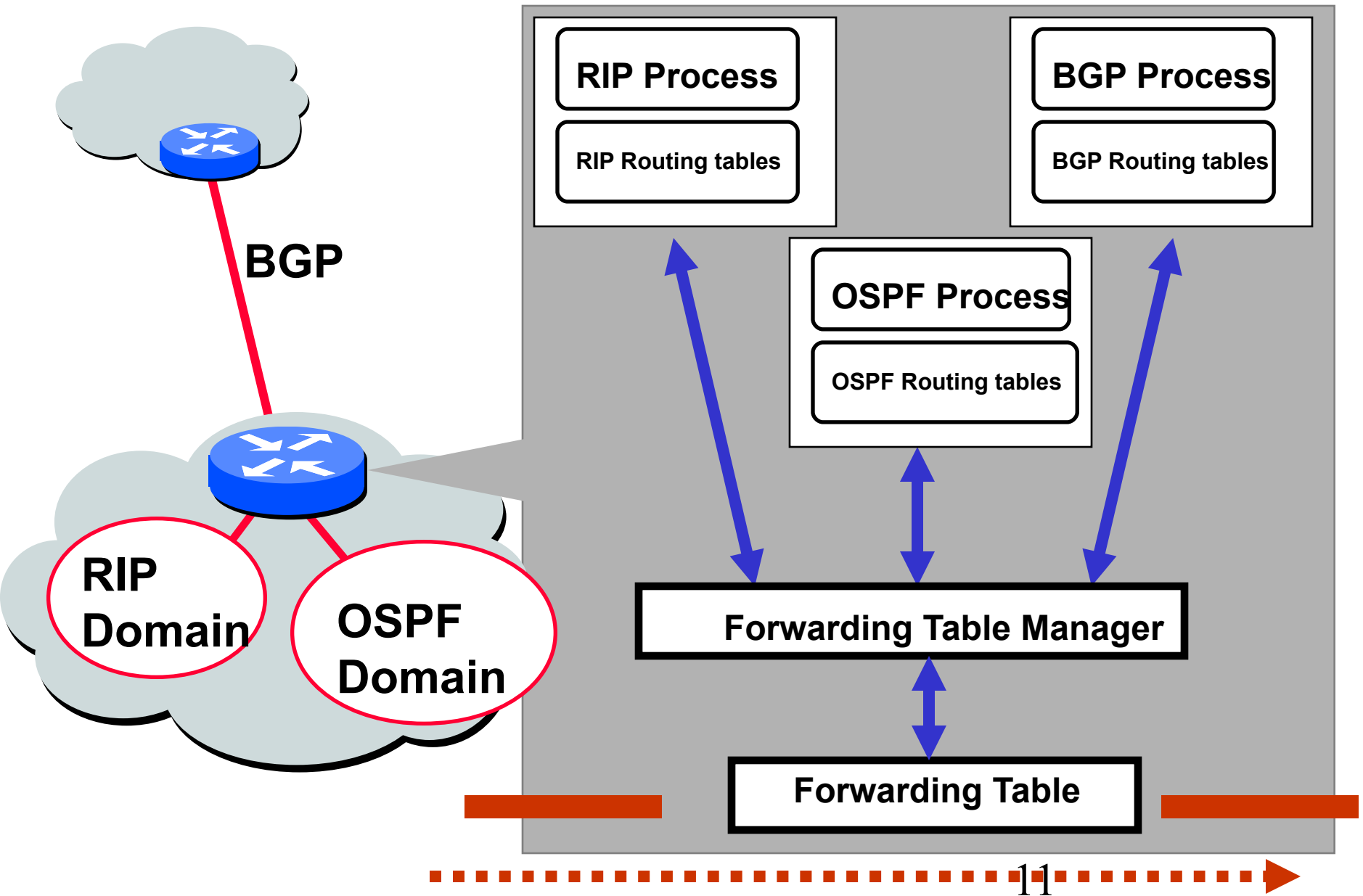
OSPF
IS-IS

RIP

EGP

BGP

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)

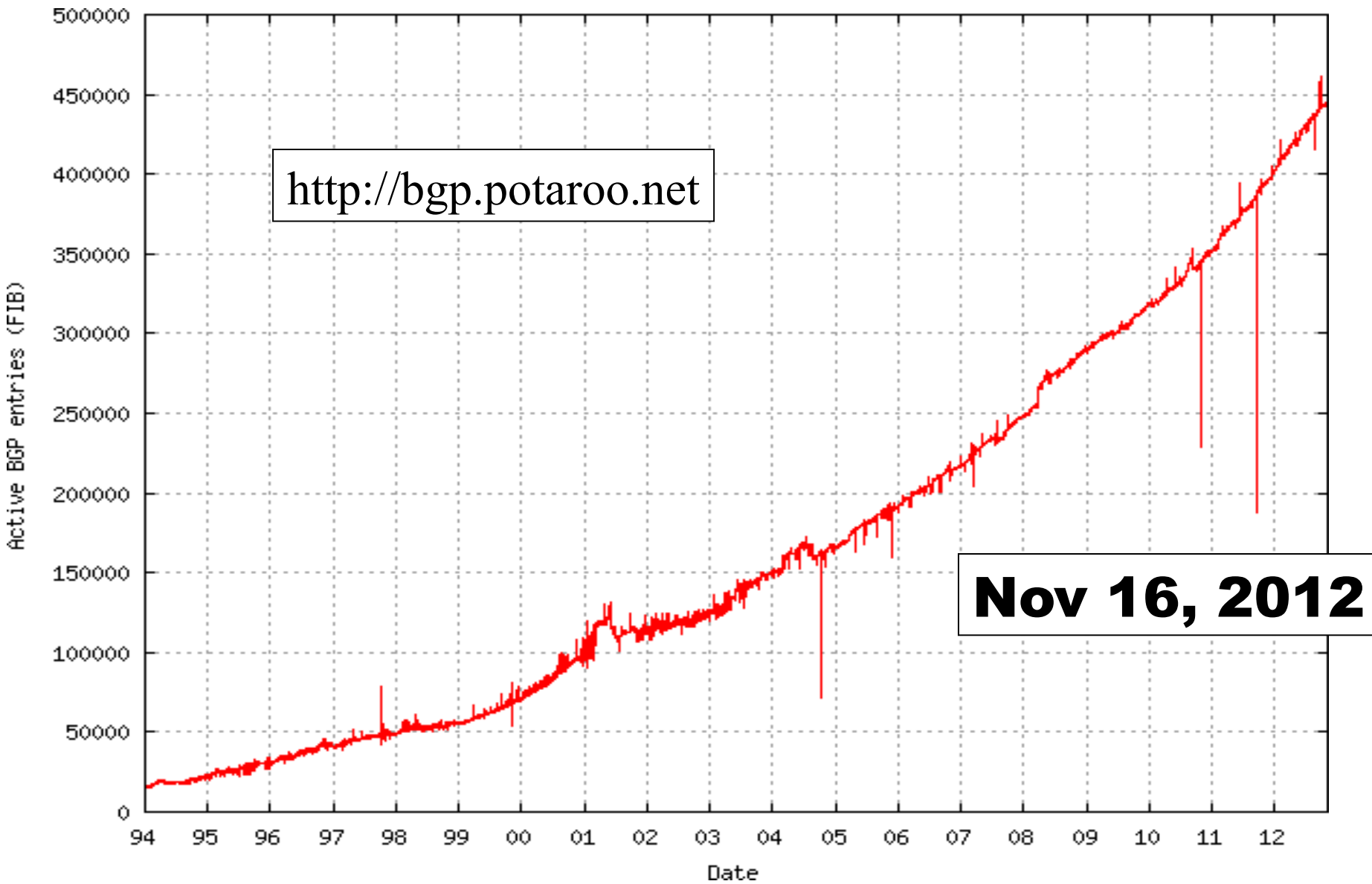
64512 through 65535 are “private”

Currently nearly 30,000 in use.

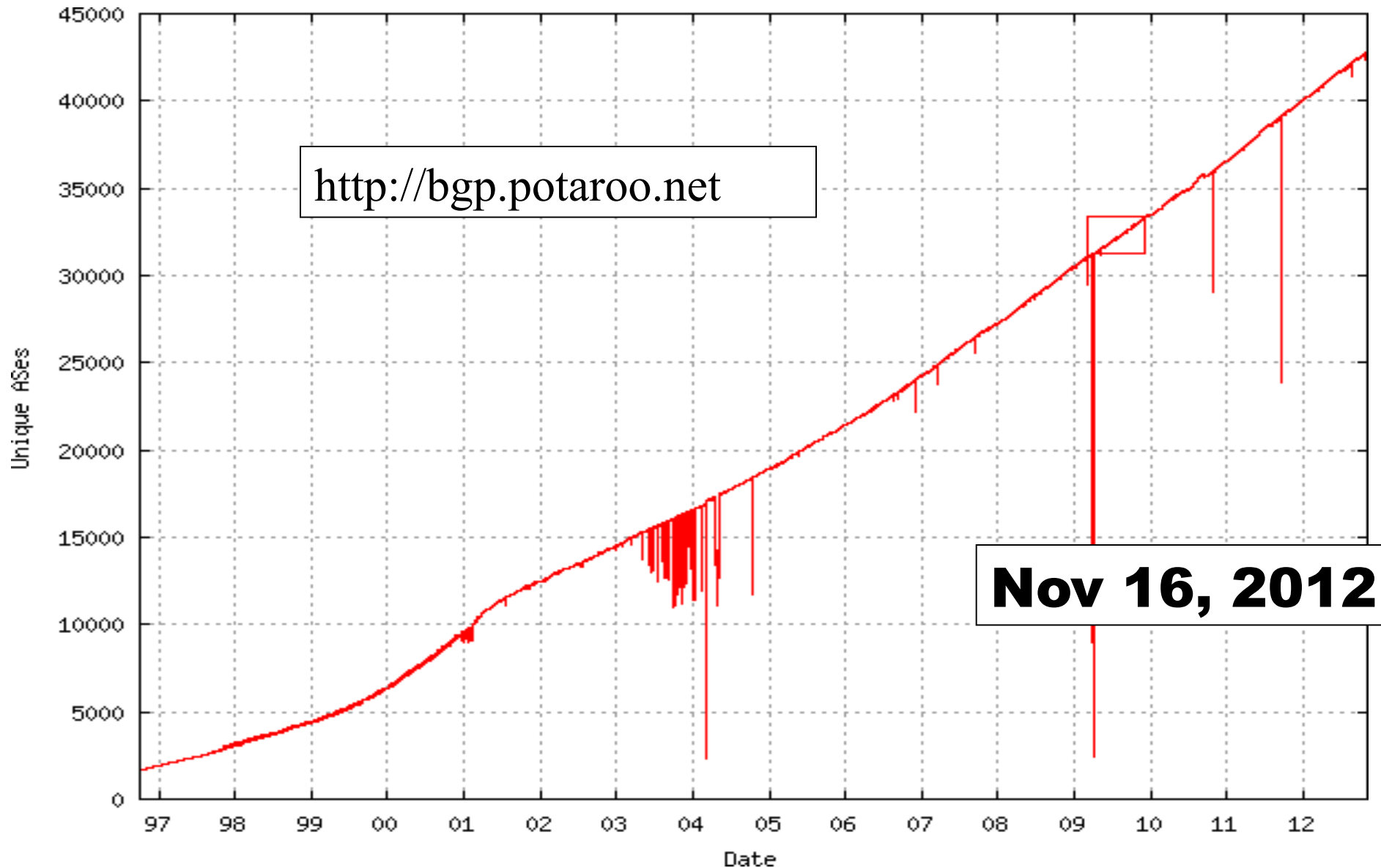
- **JANET: 786**
- **MIT: 3**
- **Harvard: 11**
- **UC San Diego: 7377**
- **AT&T: 7018, 6341, 5074, ...**
- **UUNET: 701, 702, 284, 12199, ...**
- **Sprint: 1239, 1240, 6211, 6242, ...**
- **...**

ASNs represent units of routing policy

How many prefixes are used today?

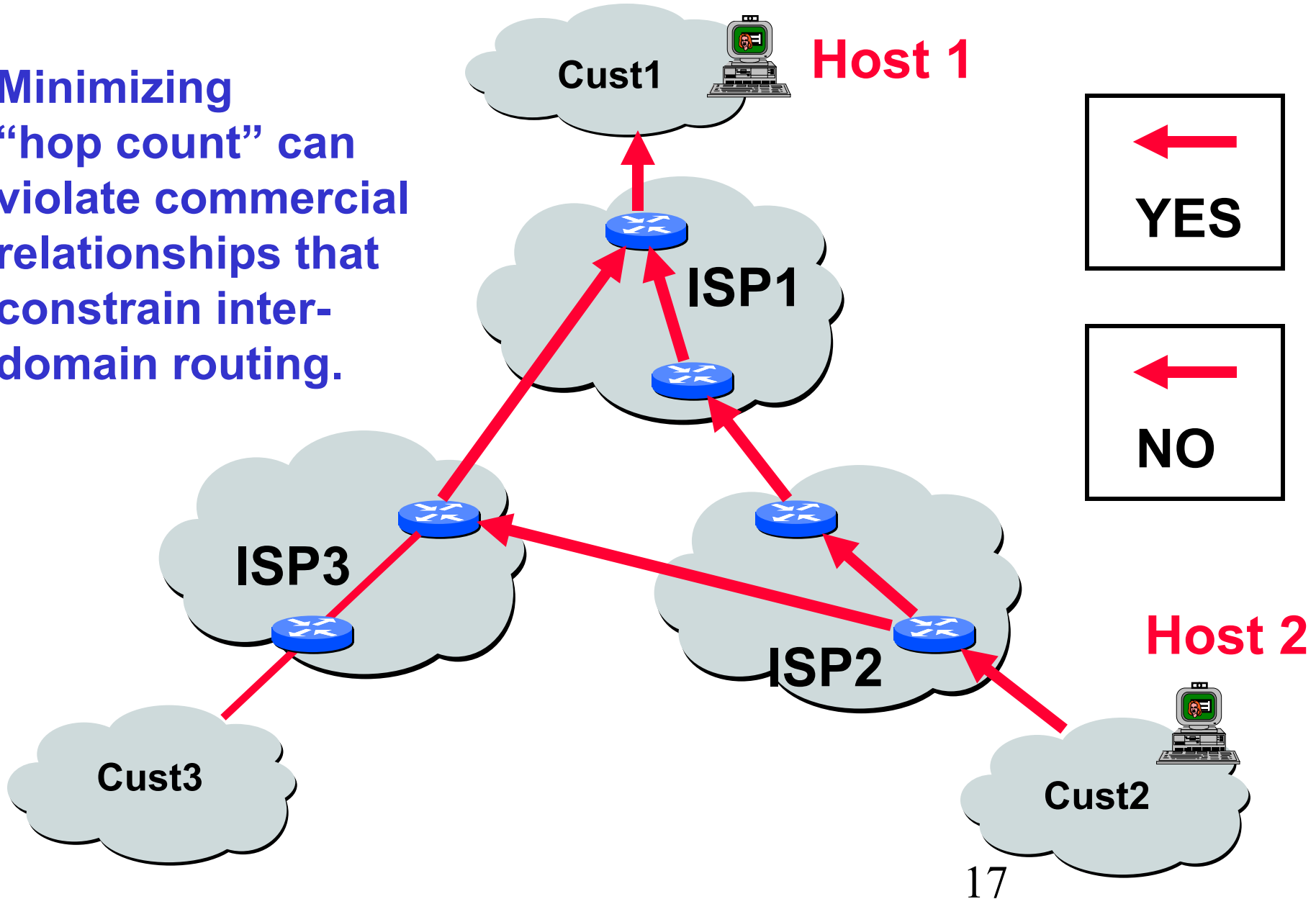


How many ASNs are used today?

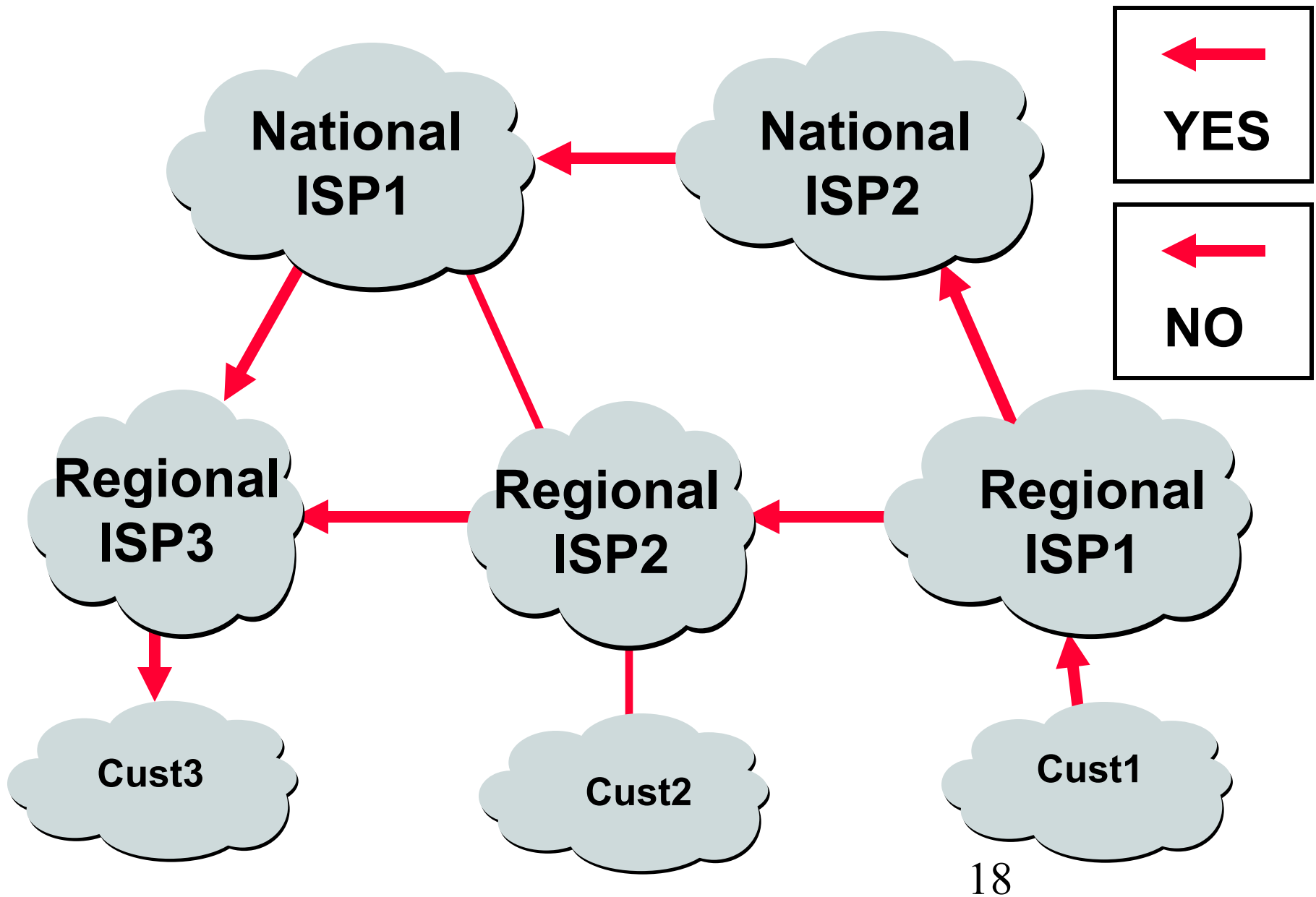


Policy-Based vs. Distance-Based Routing?

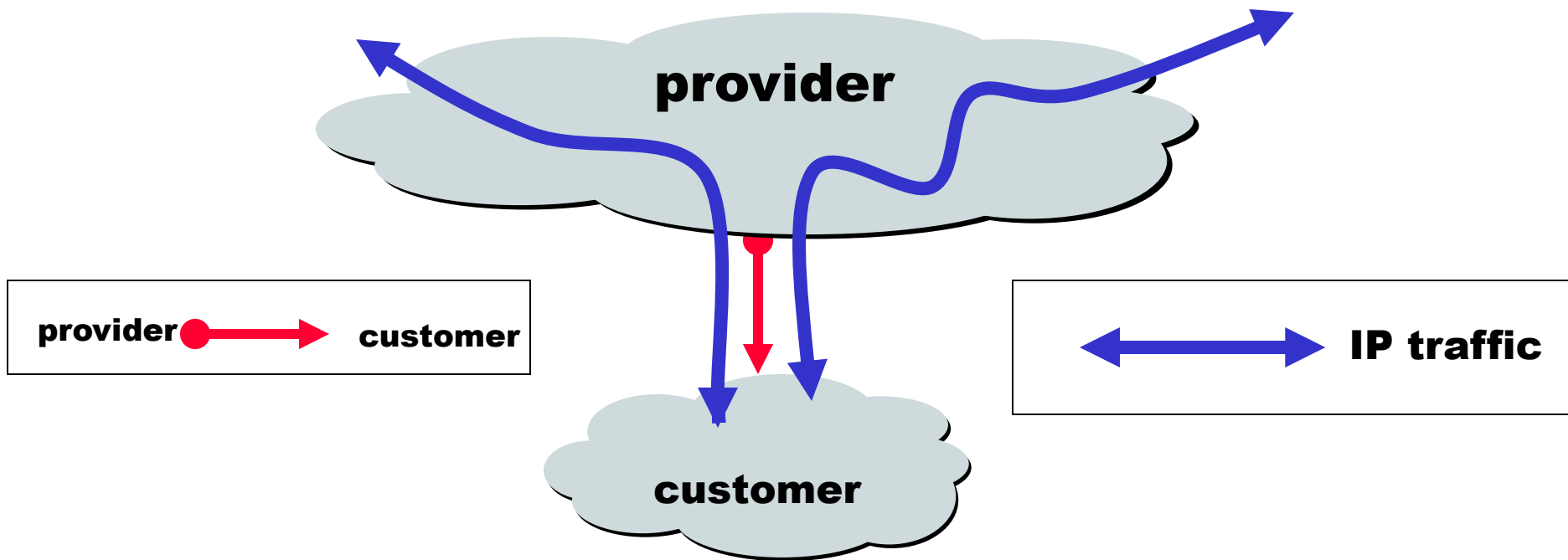
Minimizing
“hop count” can
violate commercial
relationships that
constrain inter-
domain routing.



Why not minimize “AS hop count”?

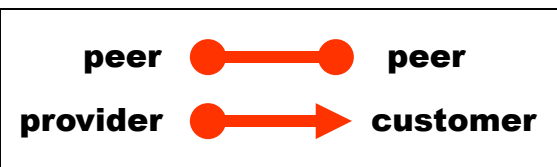
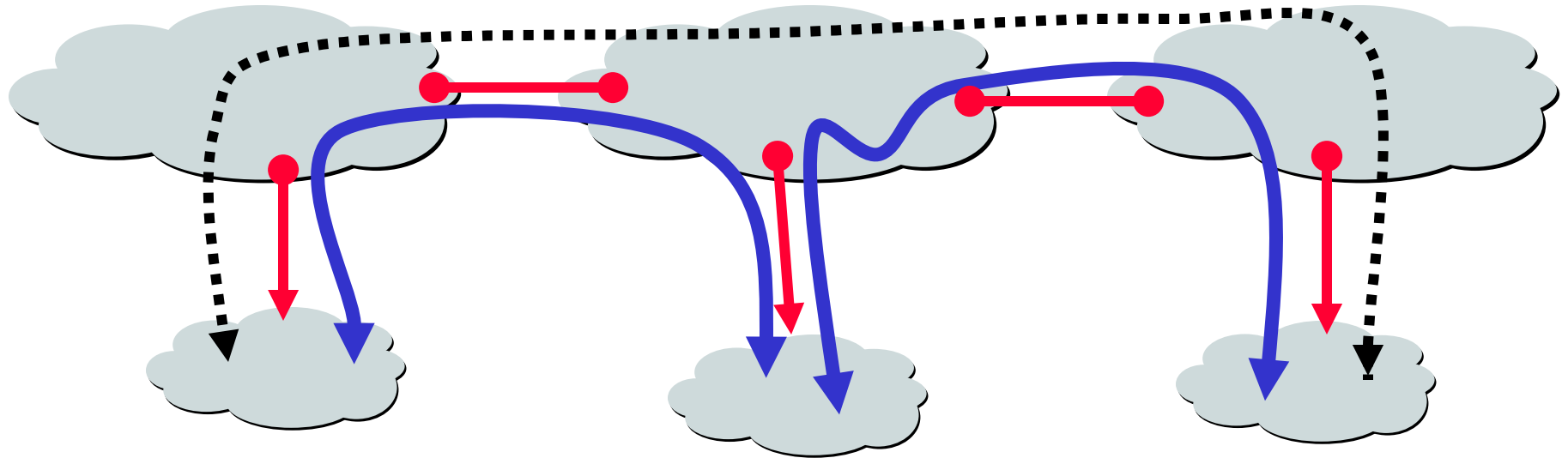


Customers and Providers



Customer pays provider for access to the Internet

The "Peering" Relationship



**traffic
allowed**



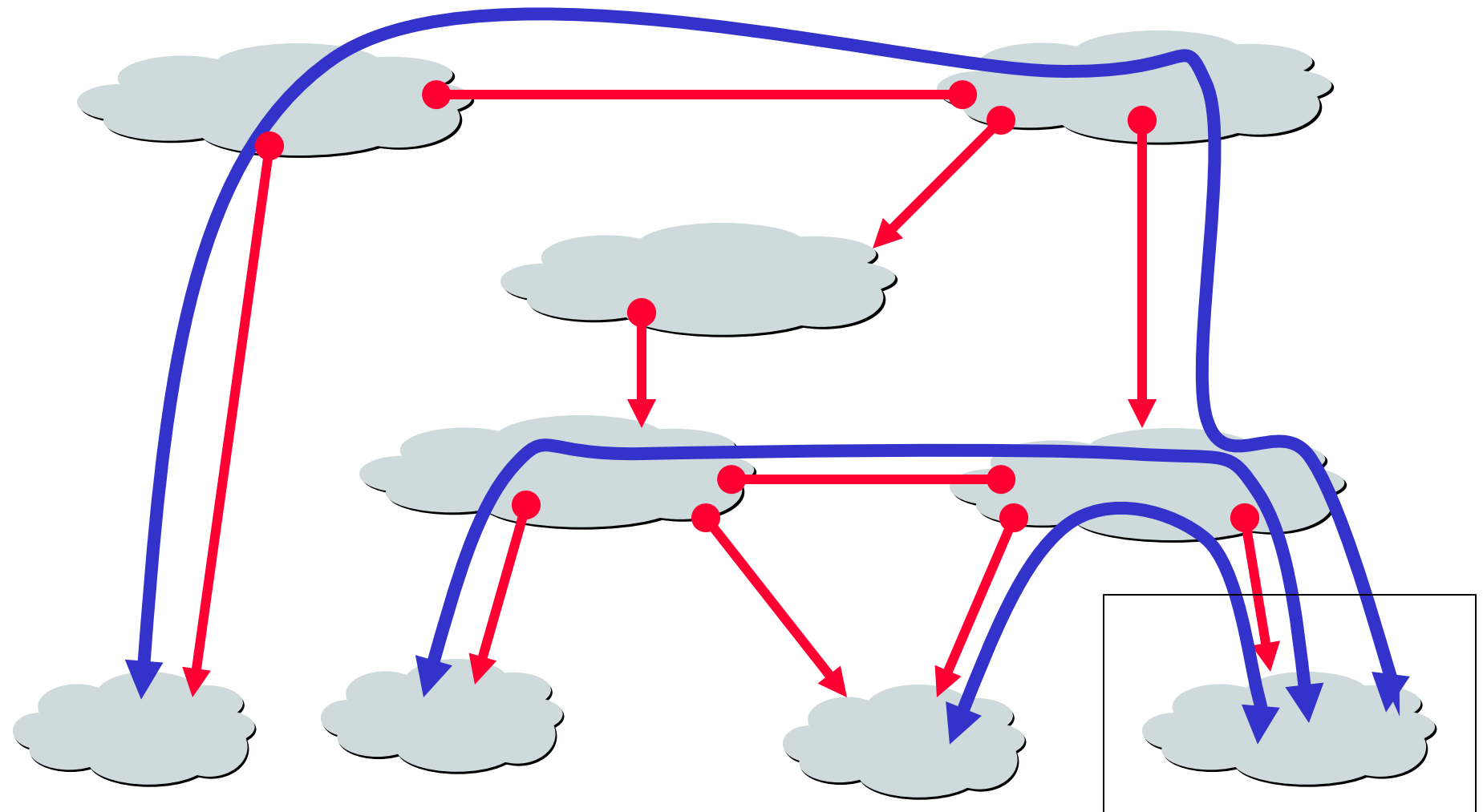
**traffic NOT
allowed**

**Peers provide transit between
their respective customers**

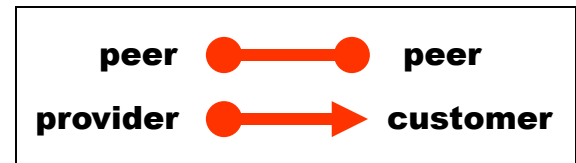
**Peers do not provide transit
between peers**

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!

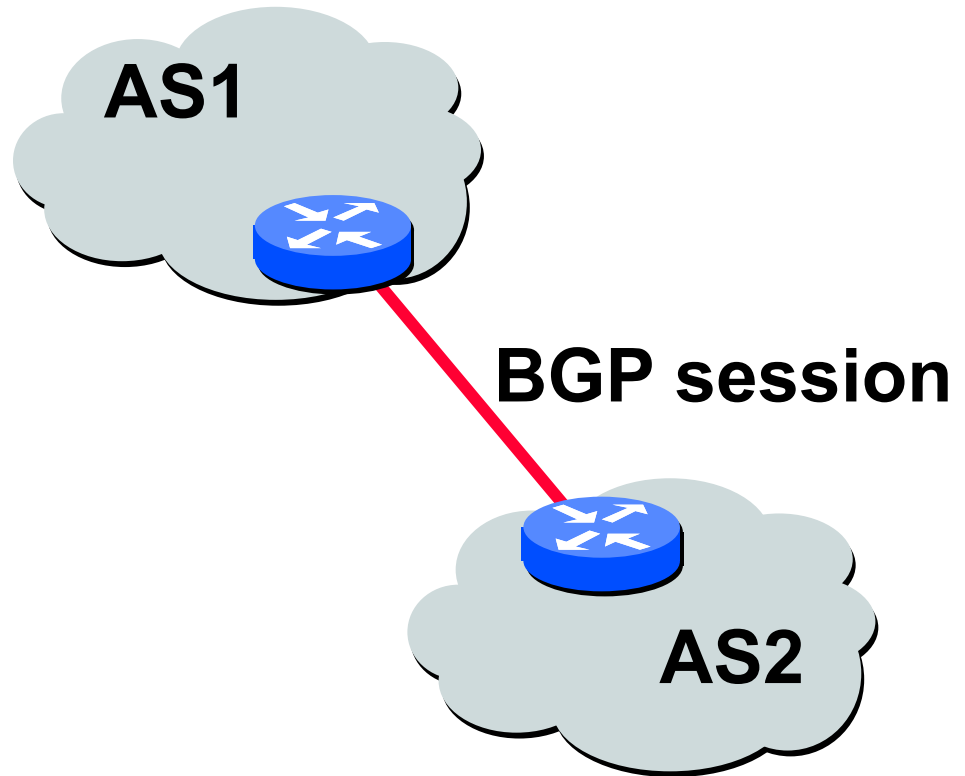
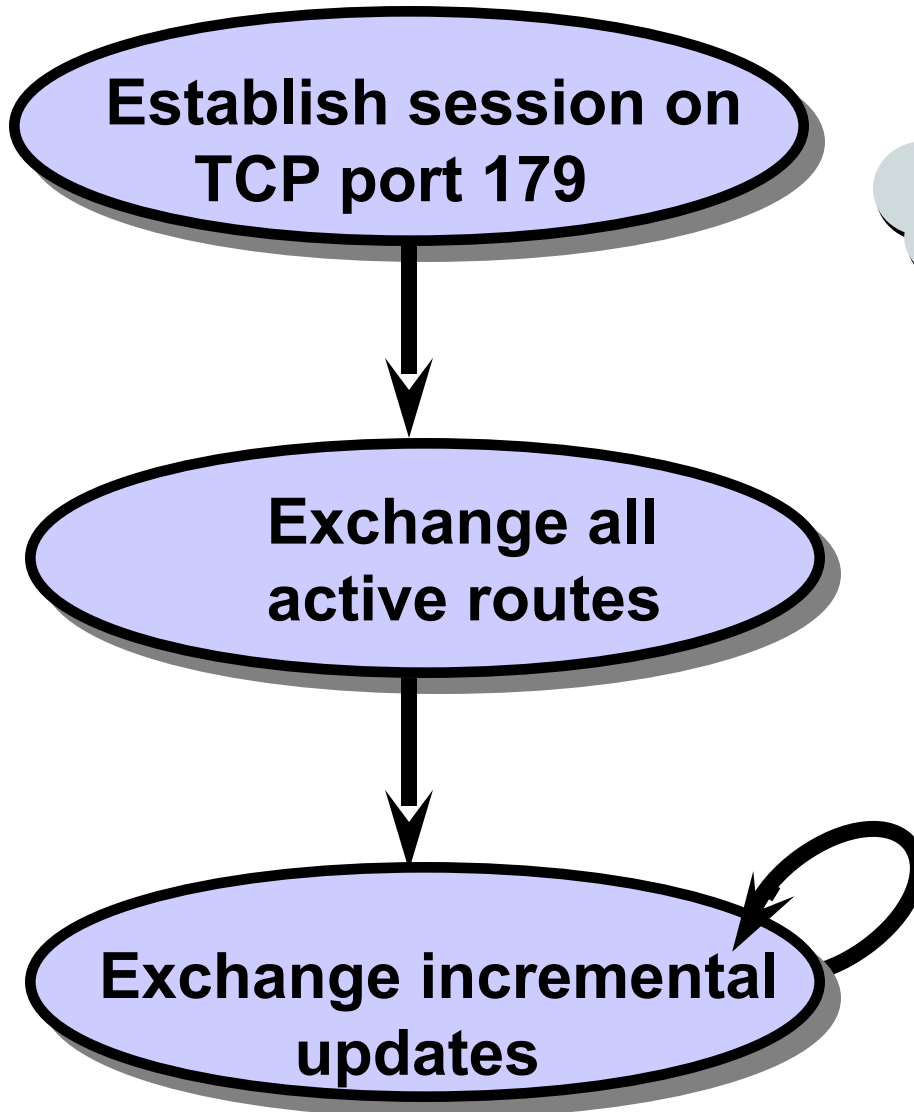
Peering agreements are often confidential.

BGP-4

- **BGP** = **B**order **G**ateway **P**rotocol
- 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]**

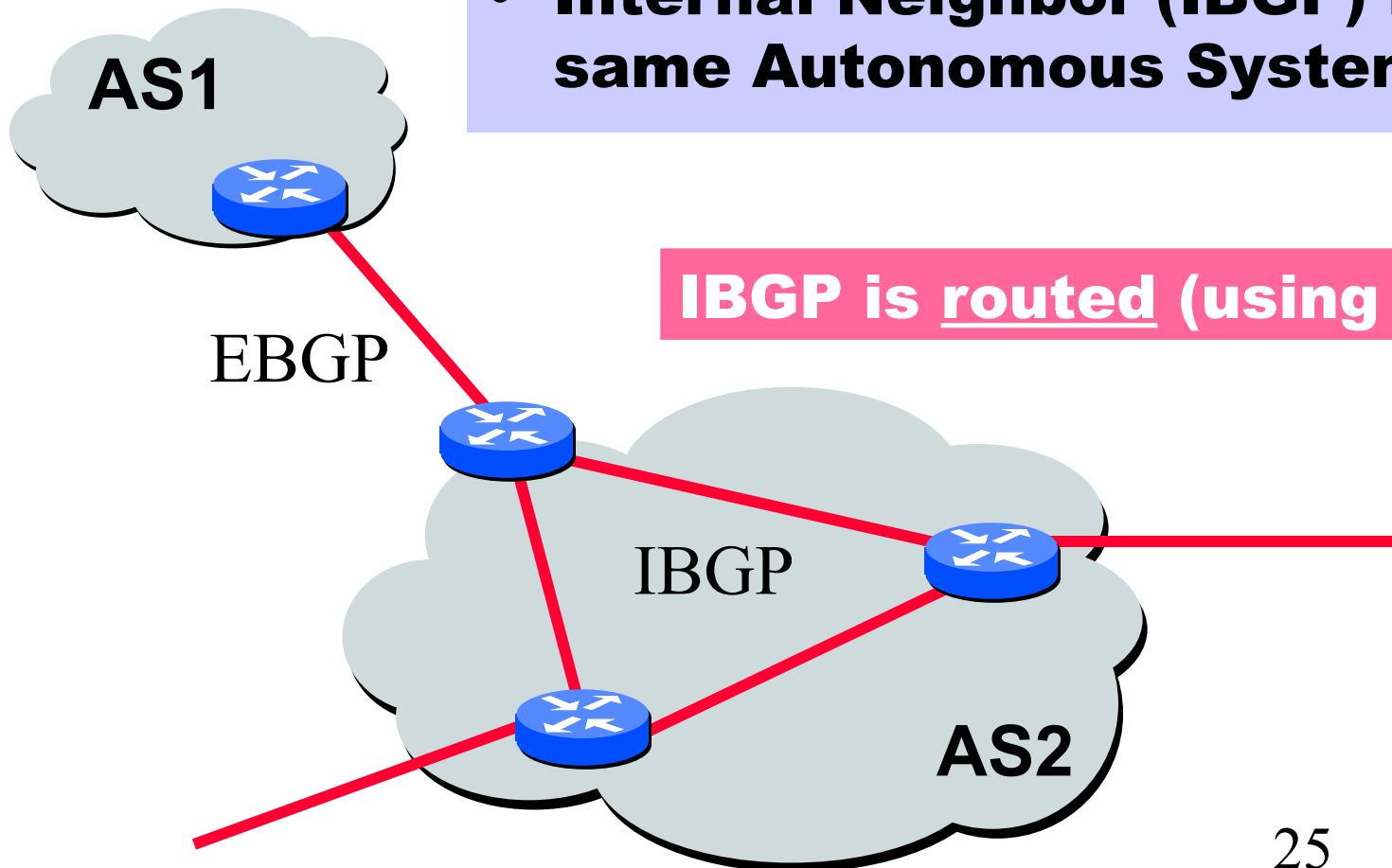
BGP Operations (Simplified)



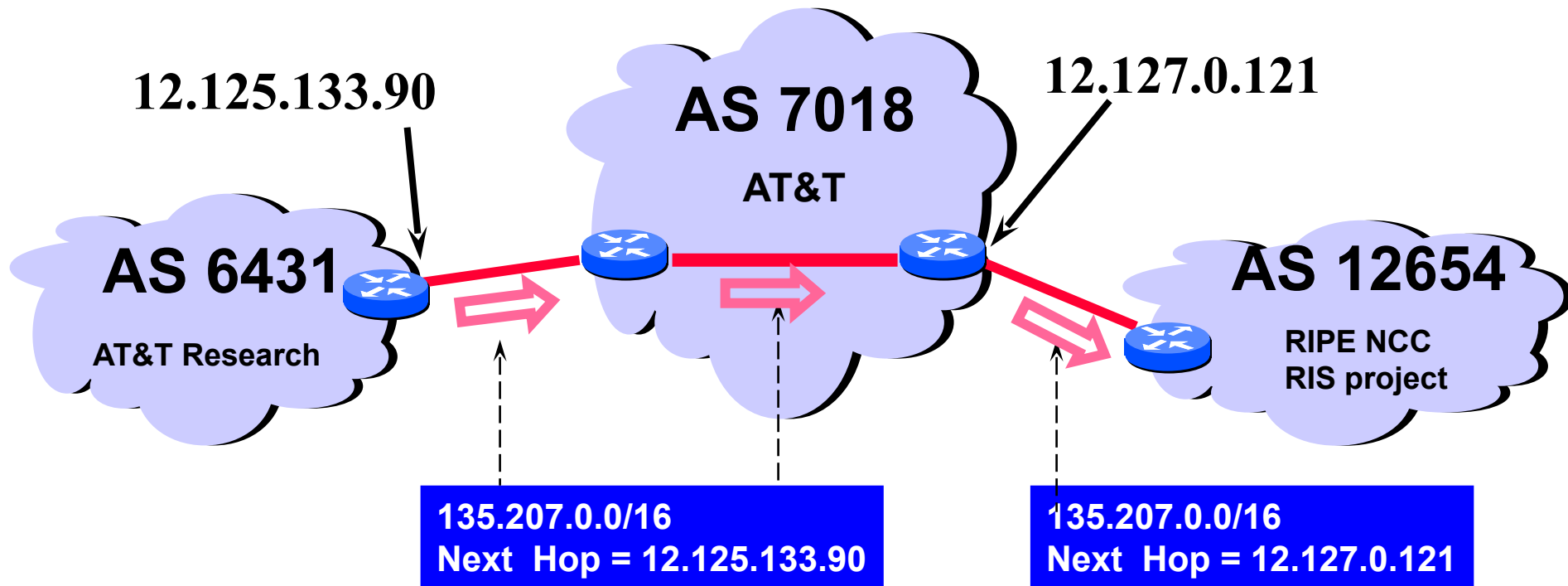
While connection is ALIVE exchange route UPDATE messages

Two Types of BGP Sessions

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

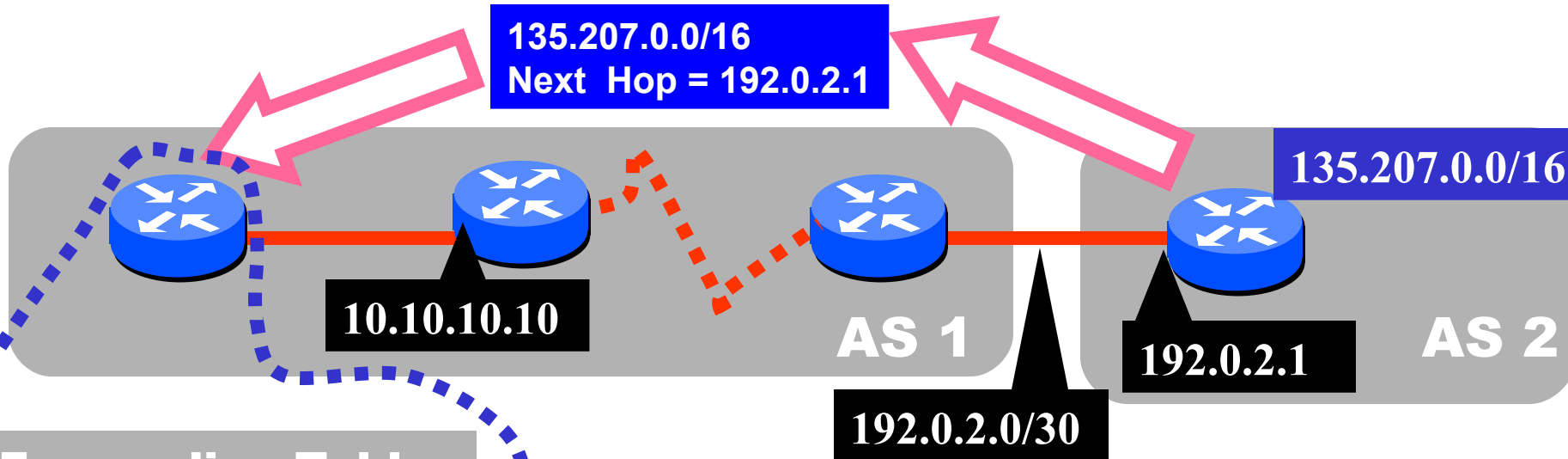


BGP Next Hop Attribute



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



Forwarding Table

destination	next hop
192.0.2.0/30	10.10.10.10

+

EGP

destination	next hop
135.207.0.0/16	192.0.2.1

Forwarding Table

destination	next hop
135.207.0.0/16	10.10.10.10
192.0.2.0/30	10.10.10.10

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]
10	CLUSTER_LIST [RFC2796]
11	DPA [Chen]
12	ADVERTISER [RFC1863]
13	RCID_PATH / CLUSTER_ID [RFC1863]
14	MP_REACH_NLRI [RFC2283]
15	MP_UNREACH_NLRI [RFC2283]
16	EXTENDED COMMUNITIES [Rosen]
17-255	reserved for development

**Most
important
attributes**

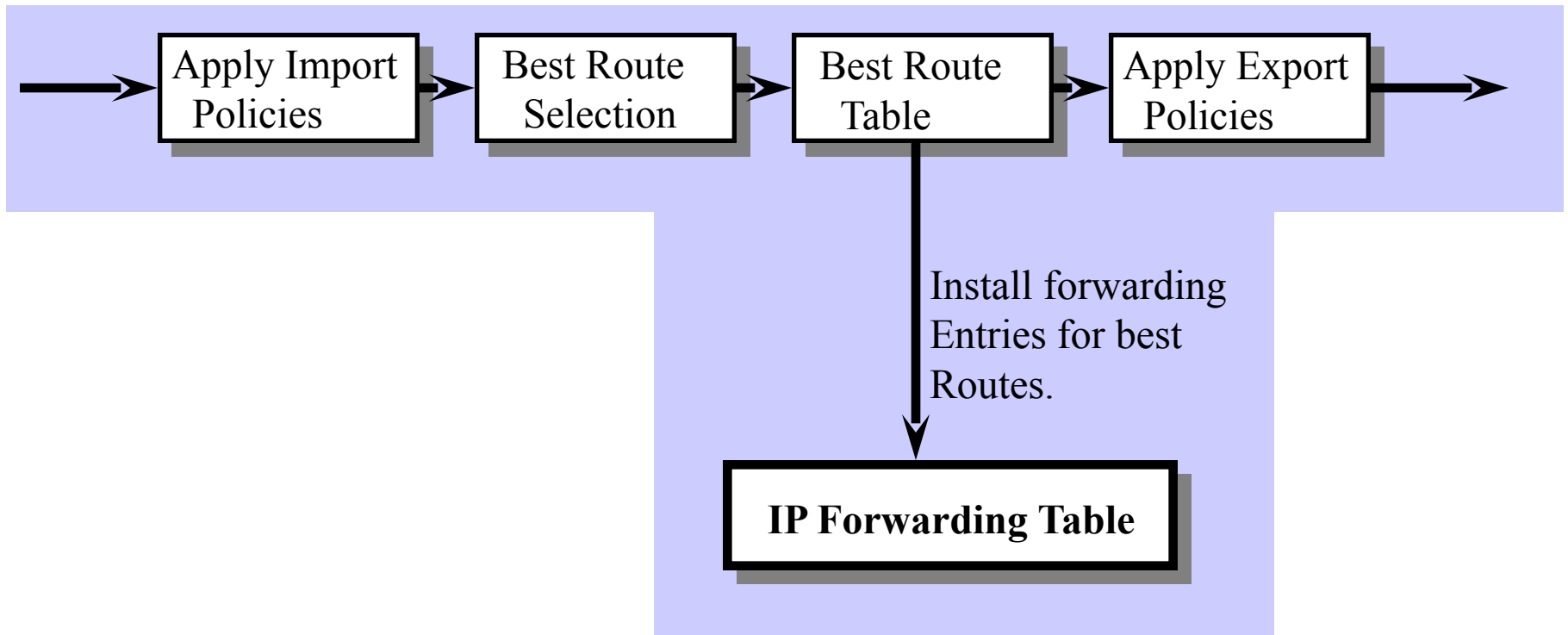
From IANA: <http://www.iana.org/assignments/bgp-parameters>

**Not all attributes
need to be present in
every announcement**

BGP Route Processing

Open ended programming.
Constrained only by vendor configuration language

Receive BGP Updates Apply Policy = filter routes & tweak attributes Based on Attribute Values Best Routes Apply Policy = filter routes & tweak attributes Transmit BGP Updates



Route Selection Summary



Highest Local Preference

Enforce relationships

Shortest AS PATH

Lowest MED

i-BGP < e-BGP

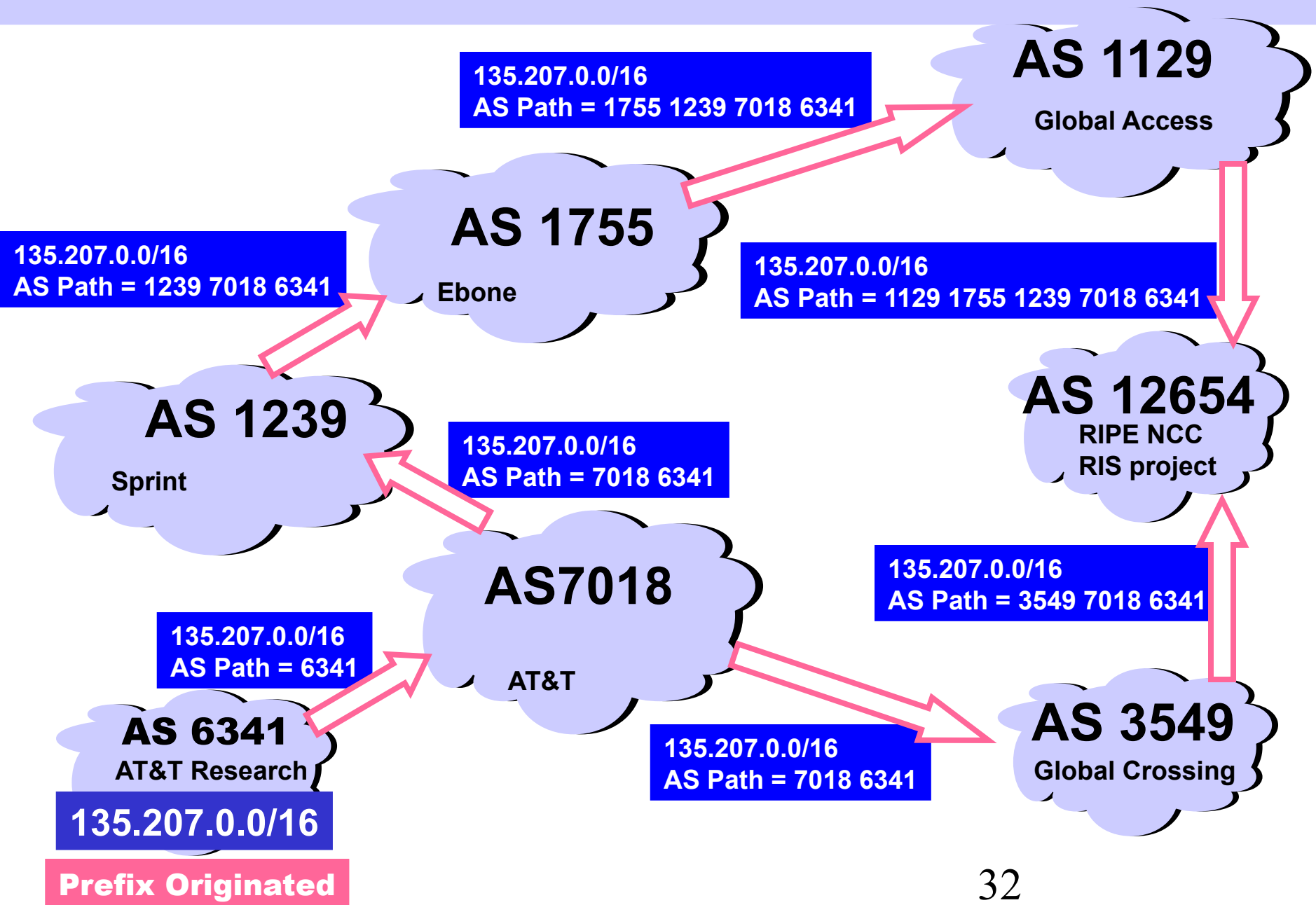
**Lowest IGP cost
to BGP egress**

traffic engineering

Lowest router ID

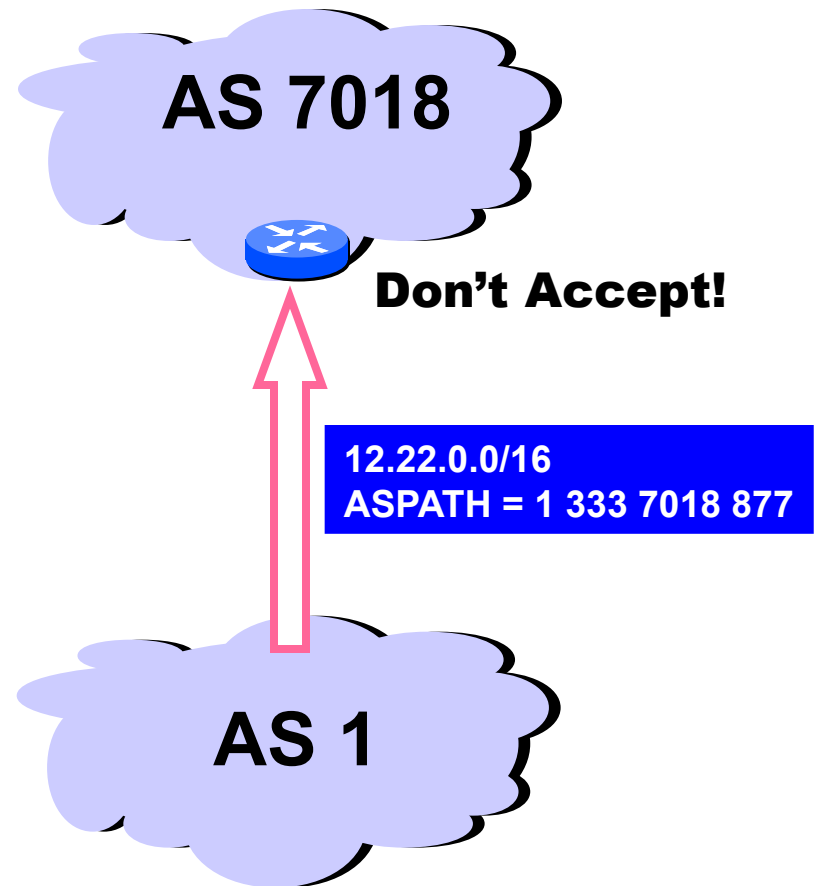
**Throw up hands and
break ties**

ASPATH Attribute



Interdomain Loop Prevention

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



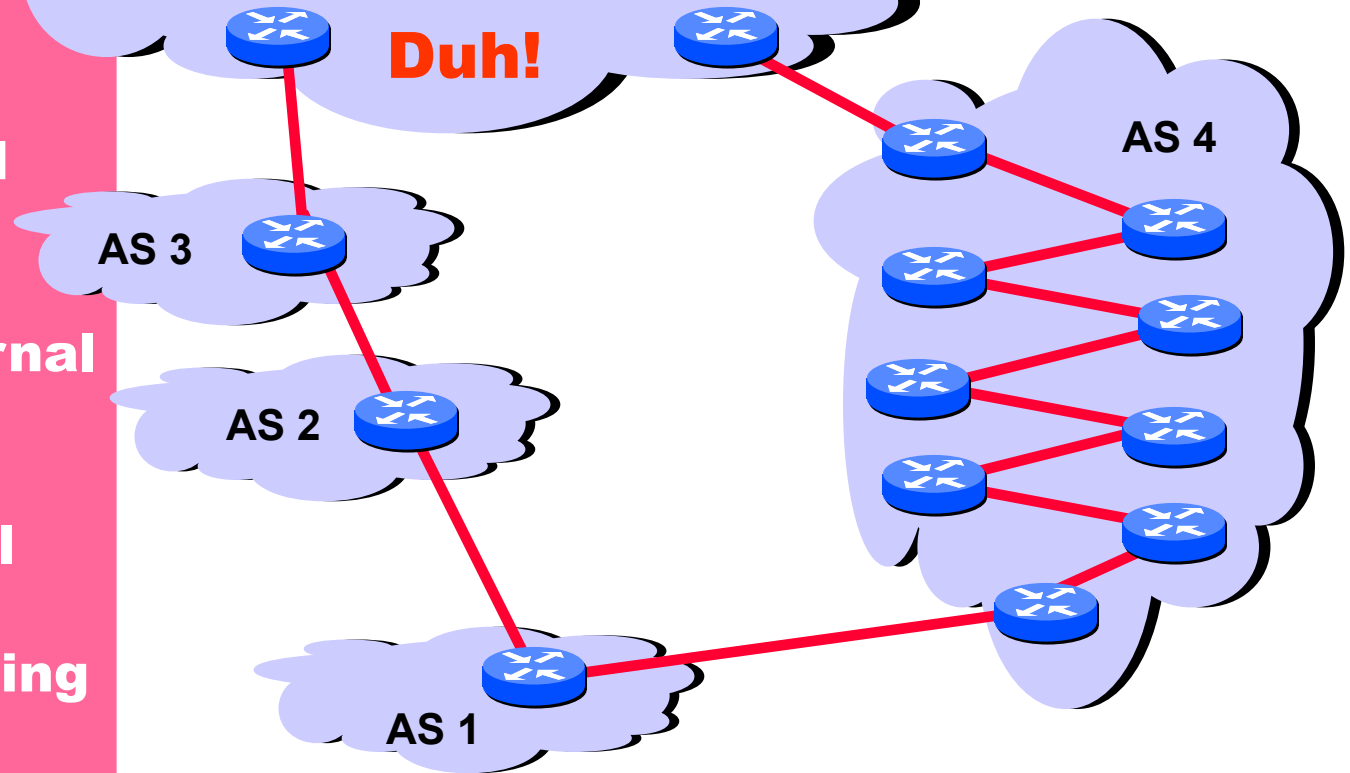
Shorter Doesn't Always Mean Shorter

Mr. BGP says that path 4 1 is better than path 3 2 1

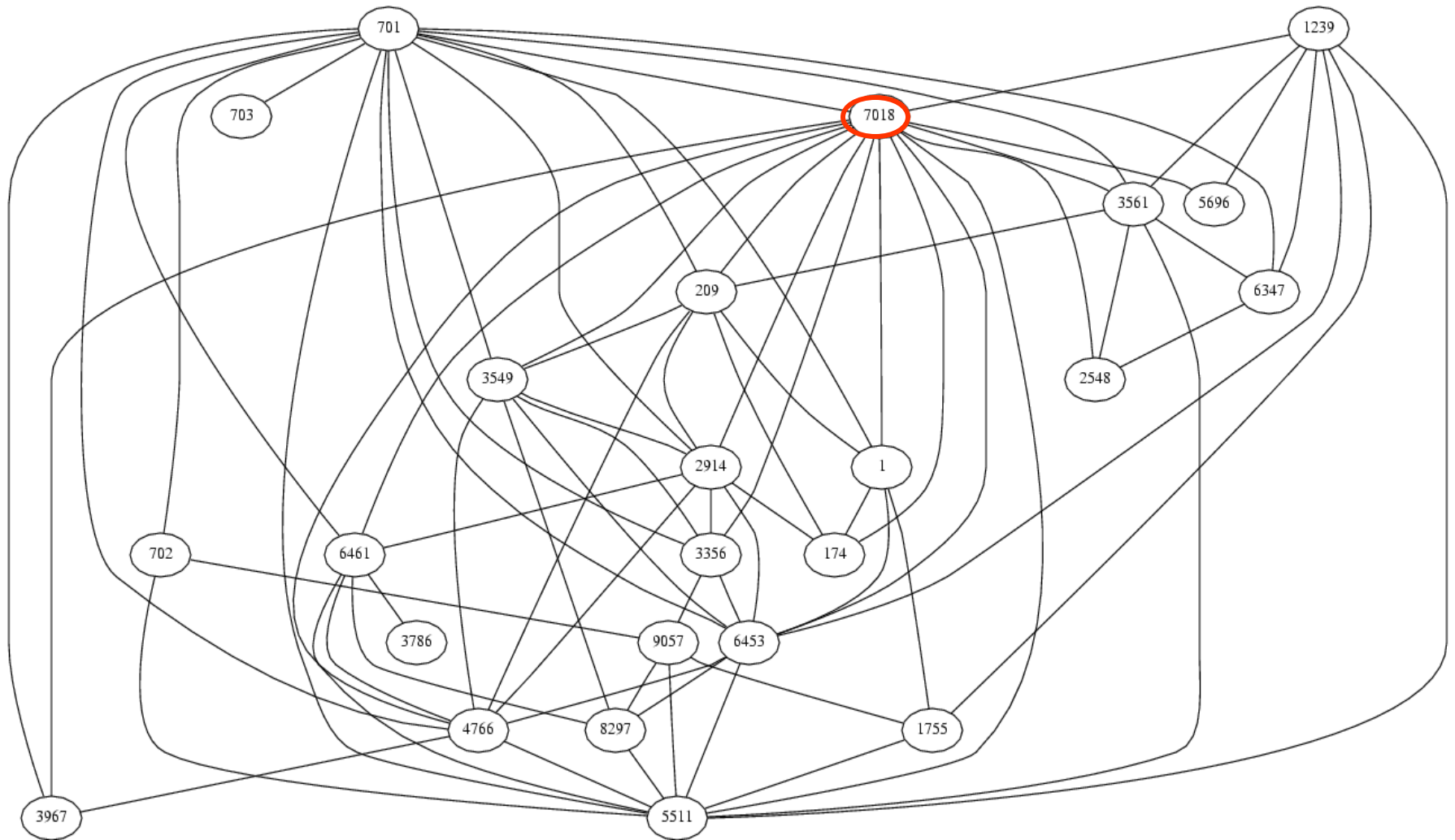
Duh!

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

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

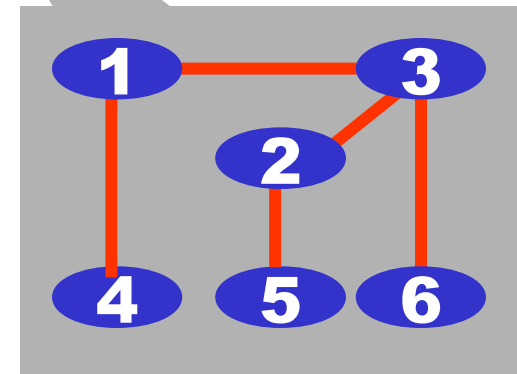
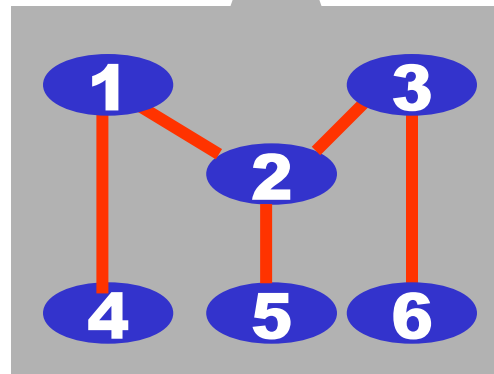
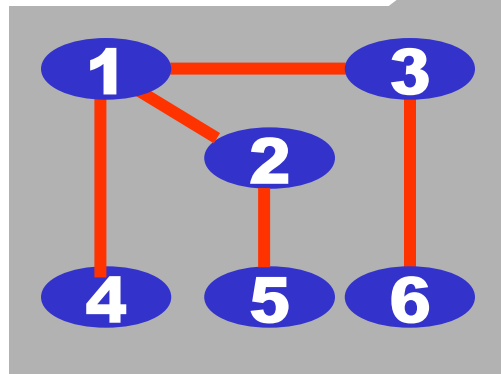
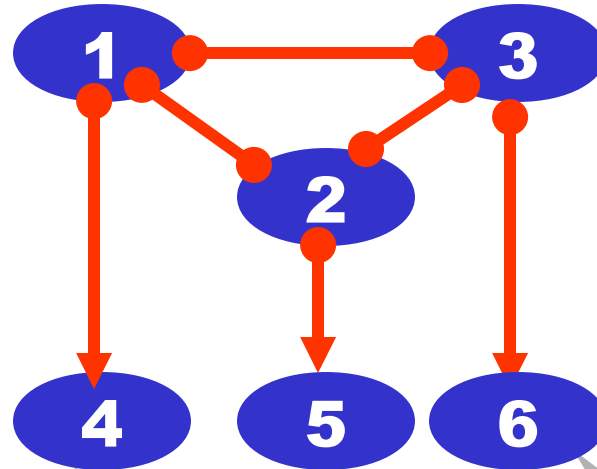
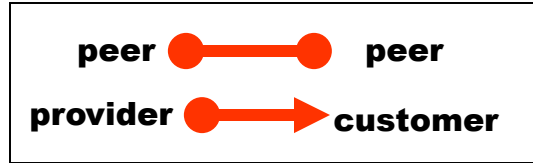


AS Graphs Can Be Fun



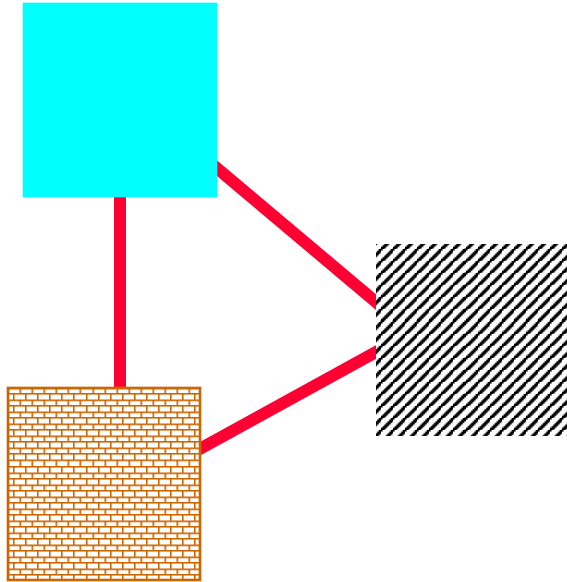
The subgraph showing all ASes that have more than 100 neighbors in full graph of 11,158 nodes. July 6, 2001. **Point of view: AT&T route-server**

AS Graphs Depend on Point of View

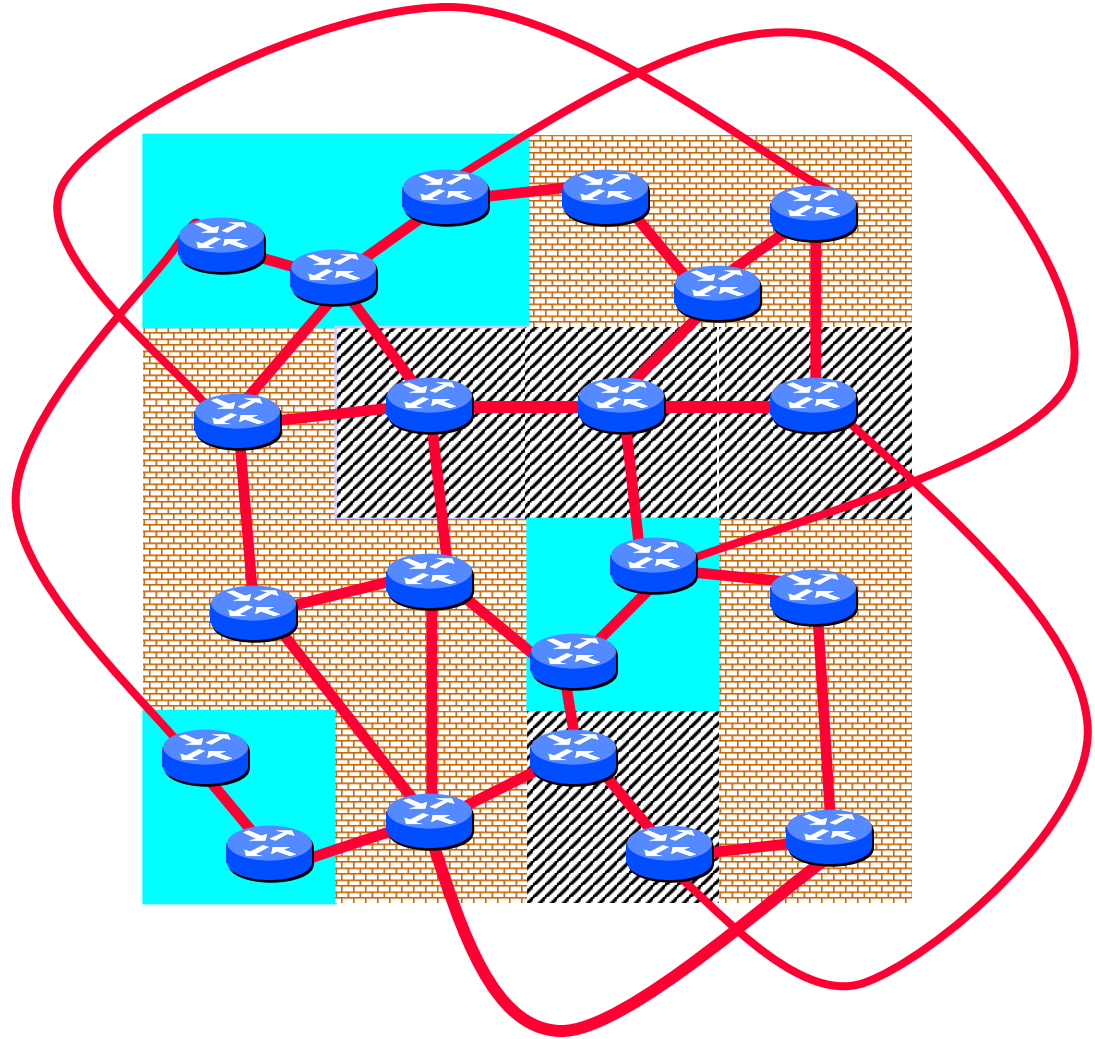


AS Graphs Do Not Show “Topology”!

BGP was designed to throw away information!



The AS graph may look like this.



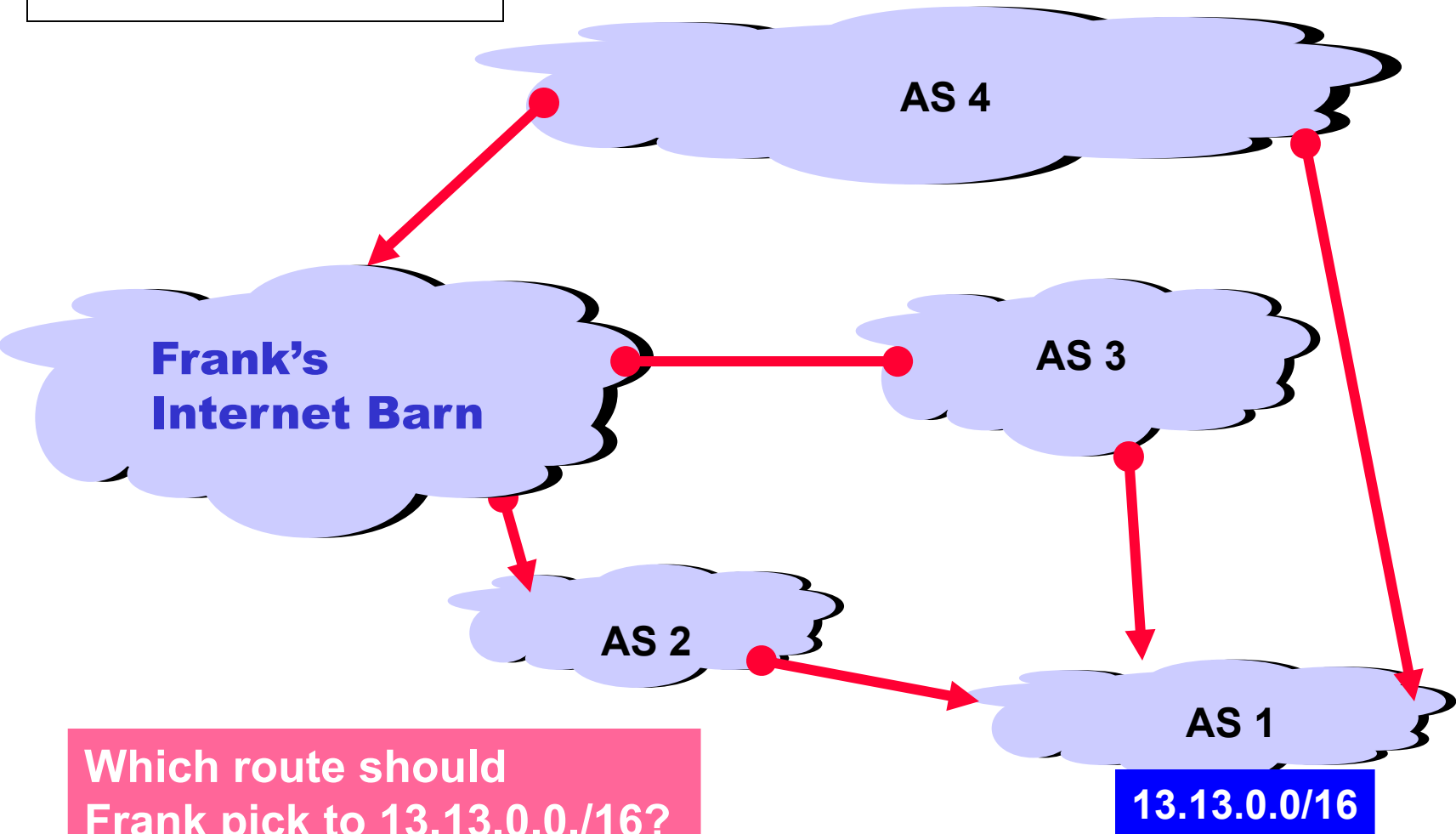
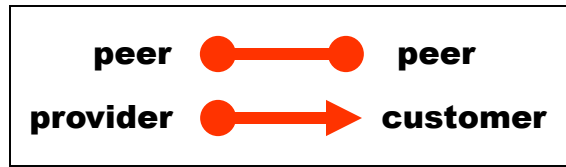
Reality may be closer to this...

Implementing Customer/Provider and Peer/Peer relationships

Two parts:

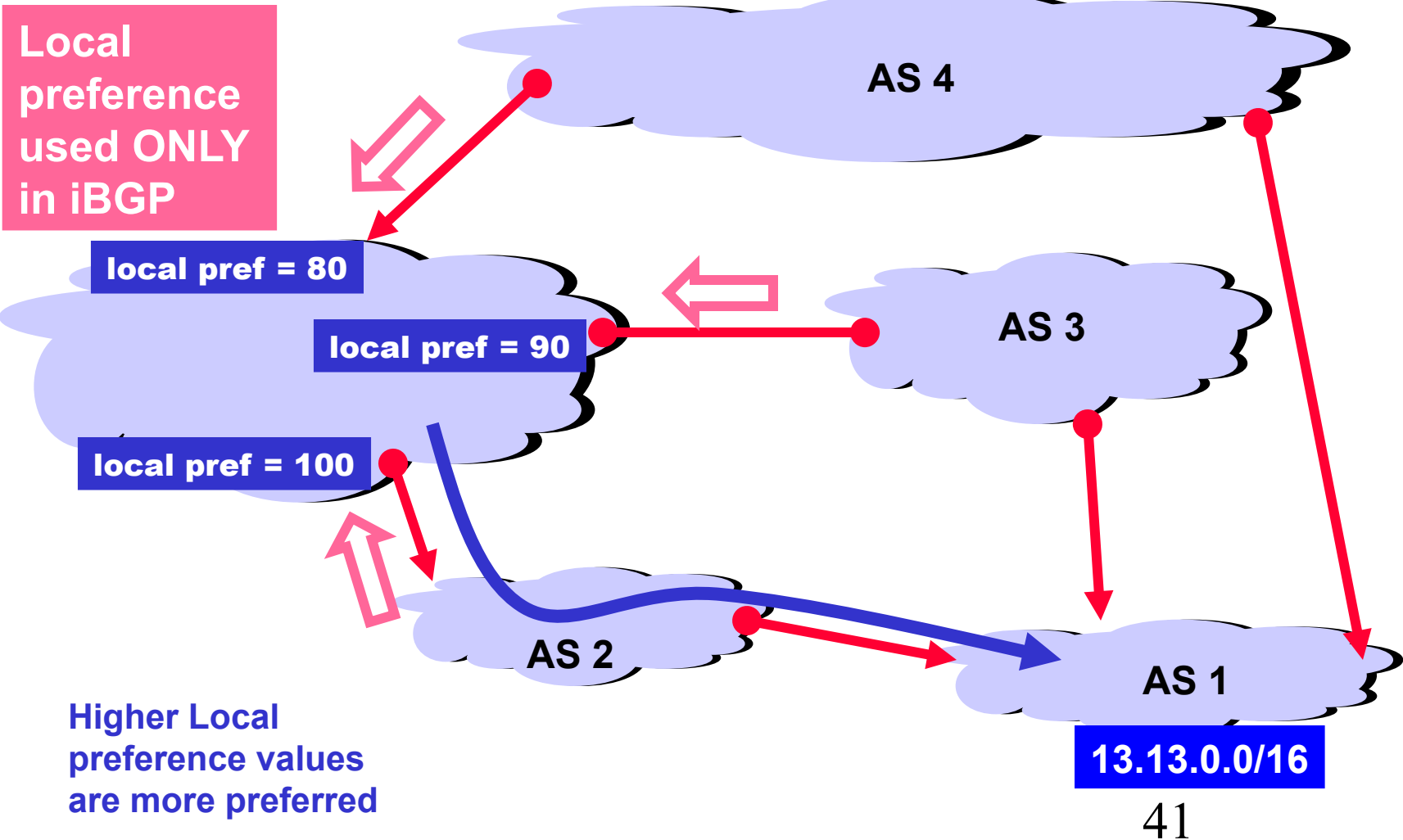
- Enforce transit relationships
 - Export all (best) routes to customers
 - Send only own and customer routes to all others
- Enforce order of route preference
 - provider < peer < customer

So Many Choices



Which route should Frank pick to 13.13.0.0/16?

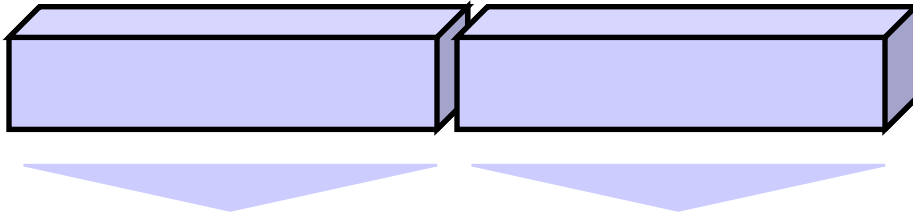
LOCAL PREFERENCE



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

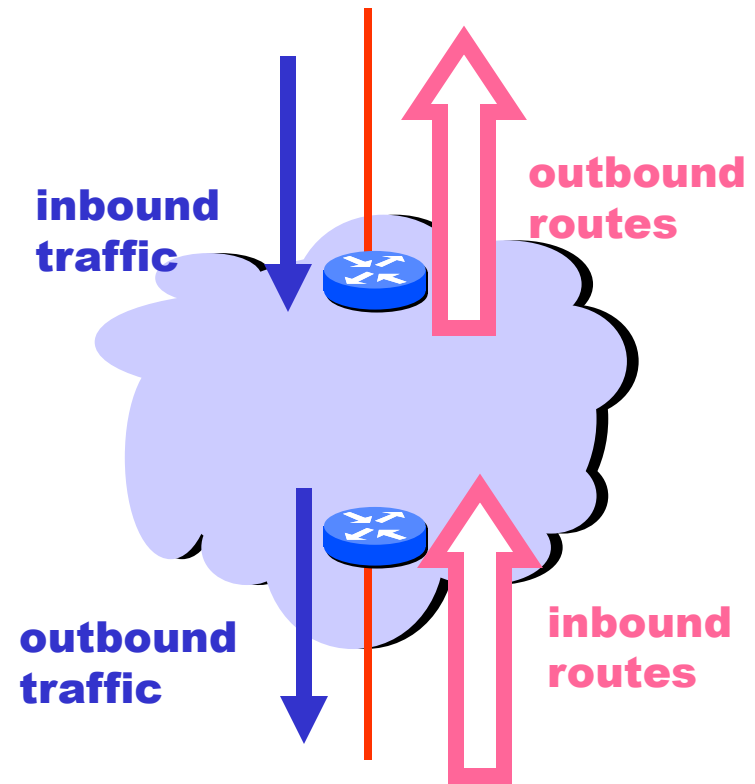
no_export = 0xFFFFFFFF01: don't export out of AS

no_advertise 0xFFFFFFFF02: don't pass to BGP neighbors

RFC 1997 (August 1996)

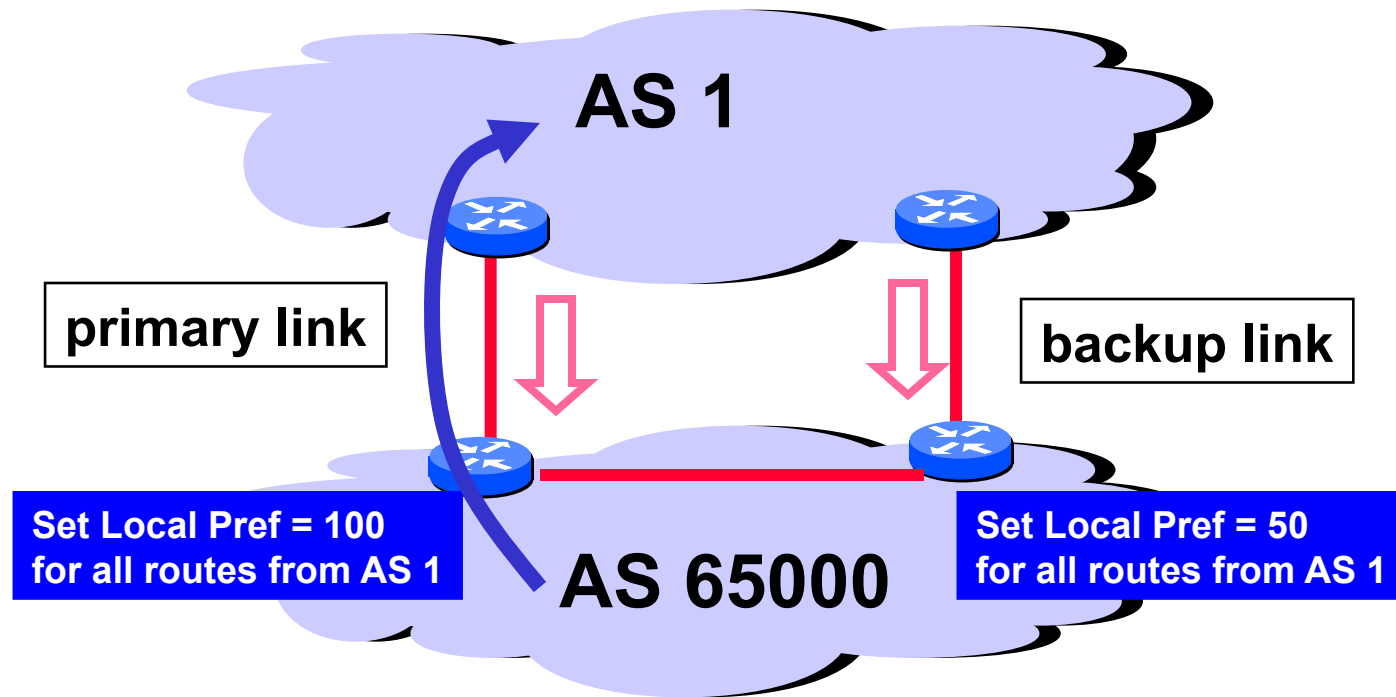
Tweak Tweak Tweak (TE)

- 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

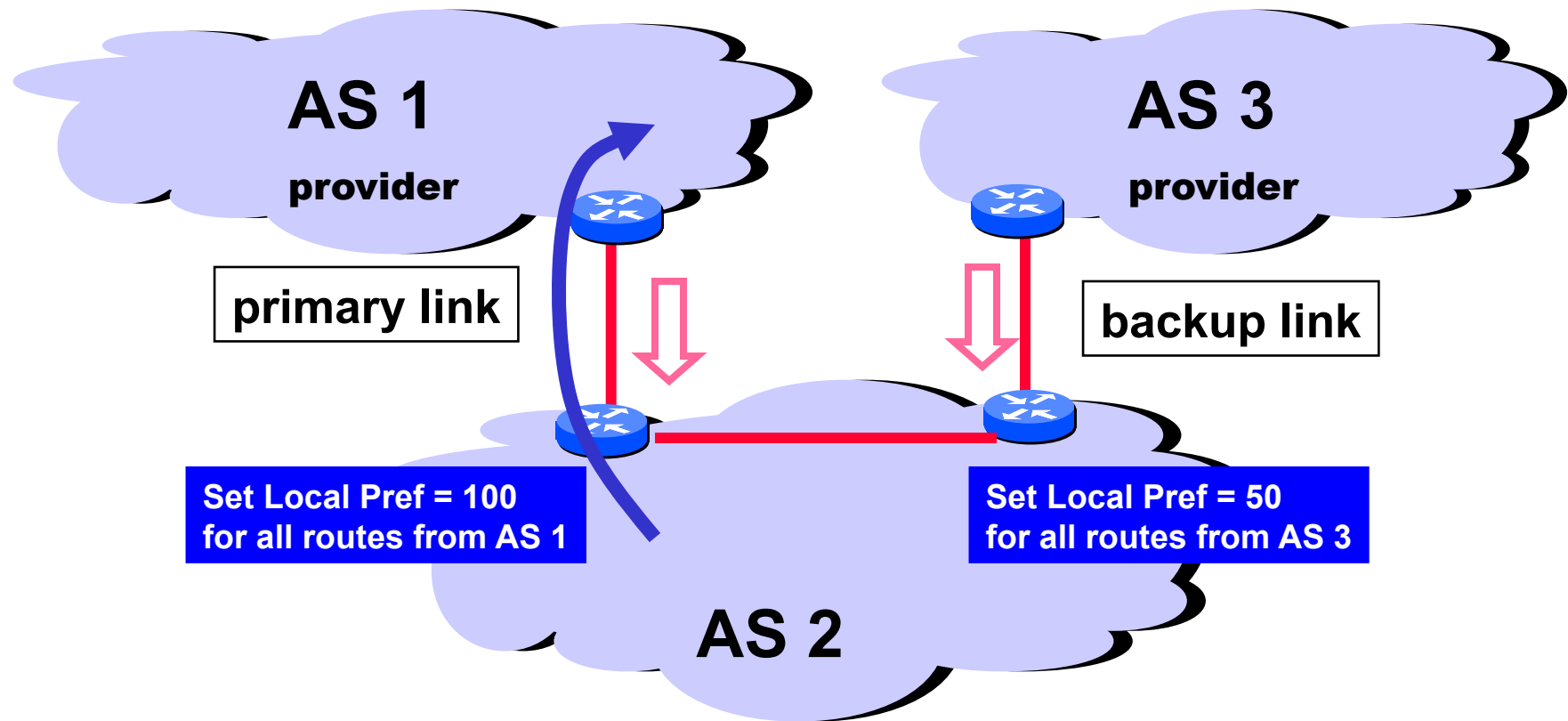
Implementing Backup Links with Local Preference (Outbound Traffic)



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

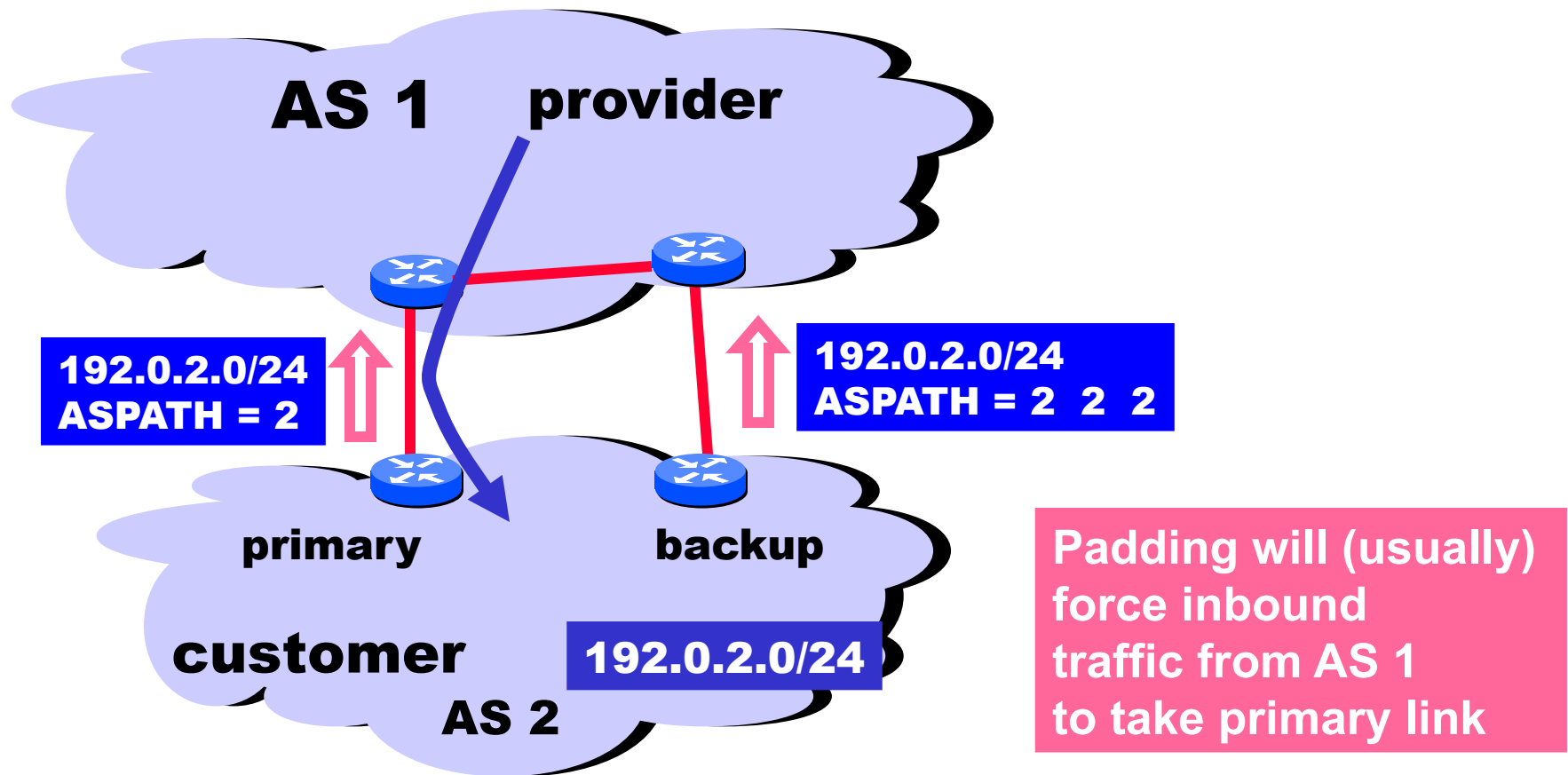
We'll talk about inbound traffic soon ...

Multihomed Backups (Outbound Traffic)

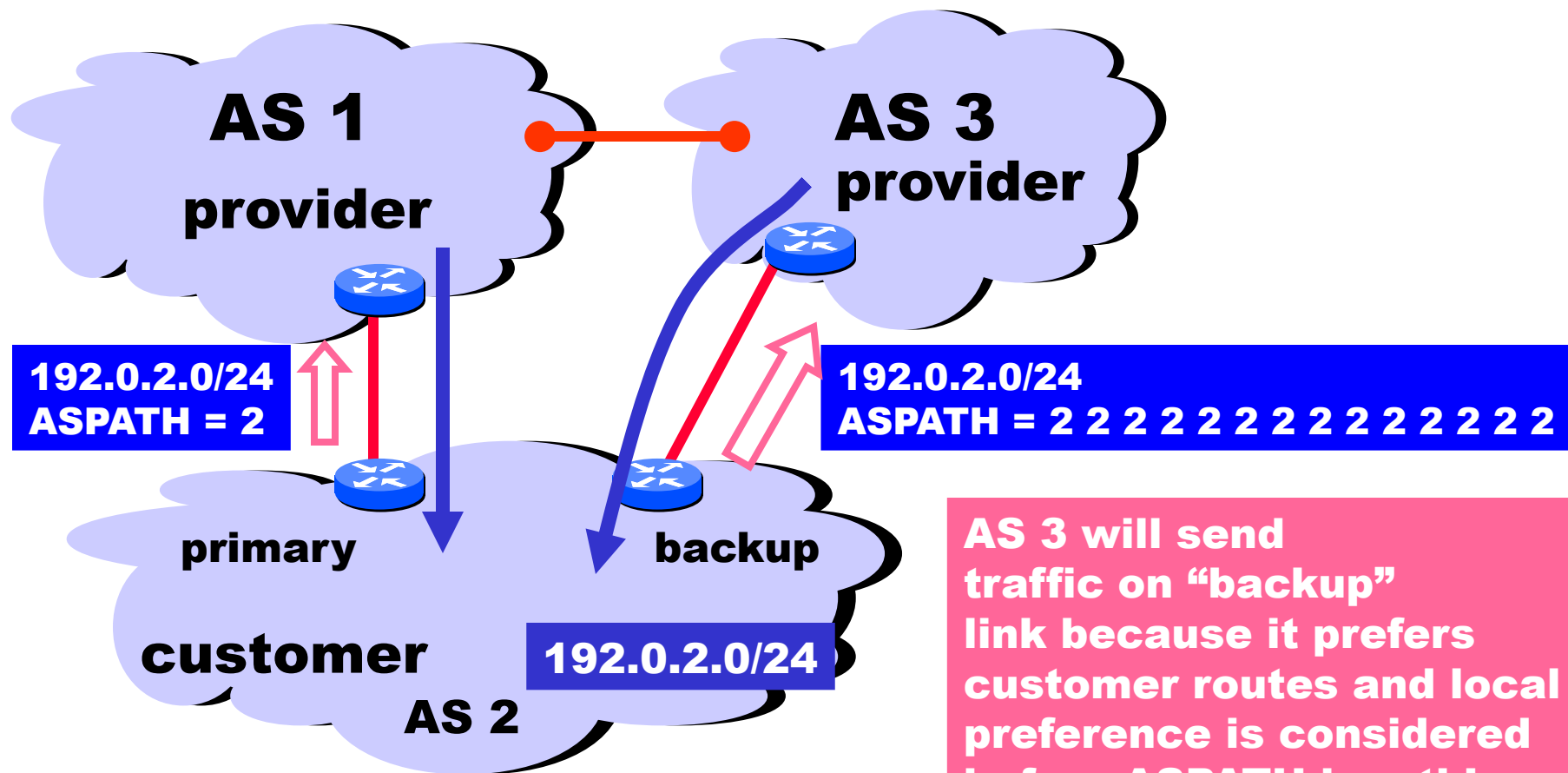


Forces outbound traffic to take primary link, unless link is down.

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



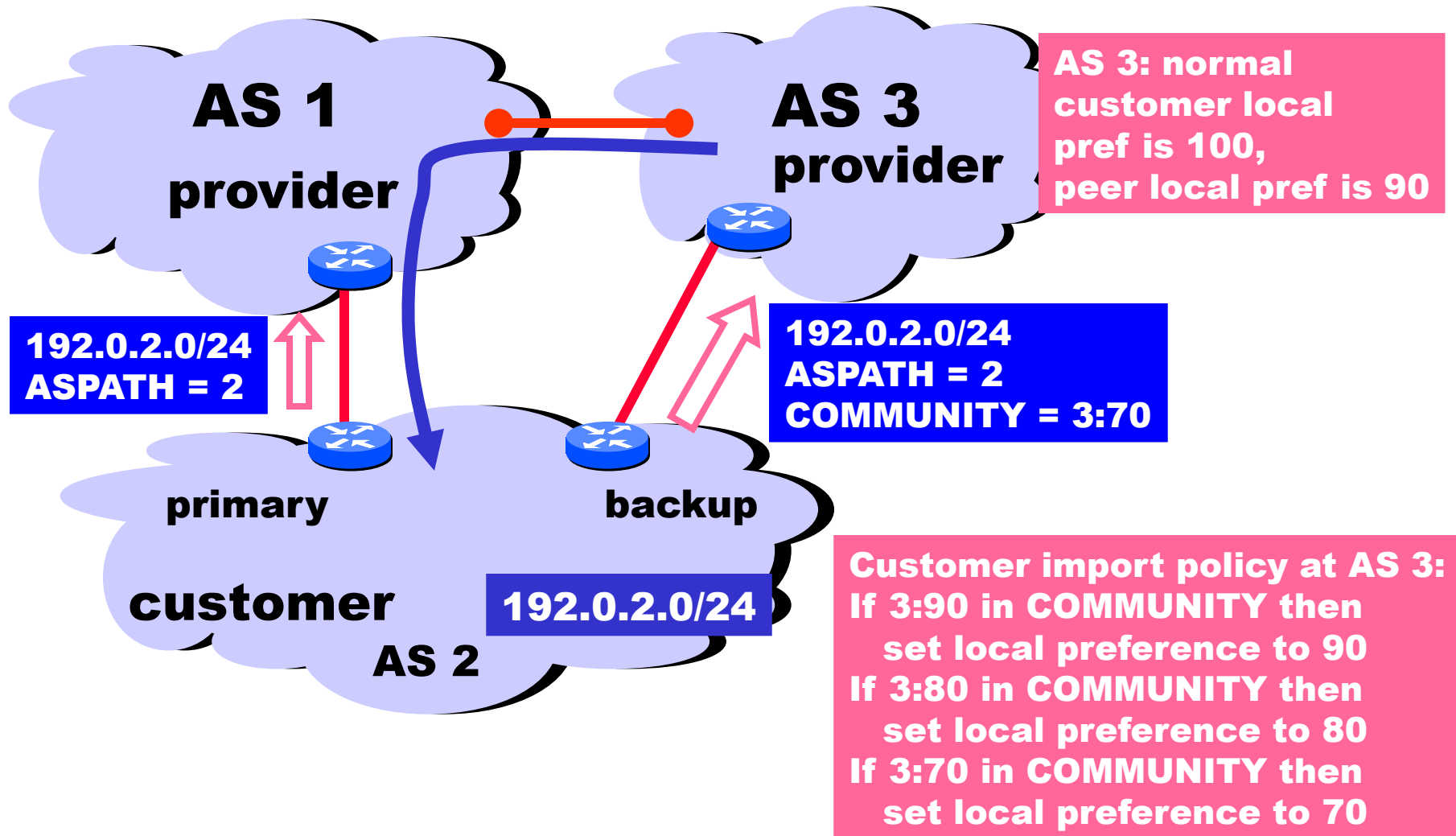
... But Padding Does Not Always Work



AS 3 will send traffic on “backup” link because it prefers customer routes and local preference is considered before ASPATH length!

Padding in this way is often used as a form of load balancing

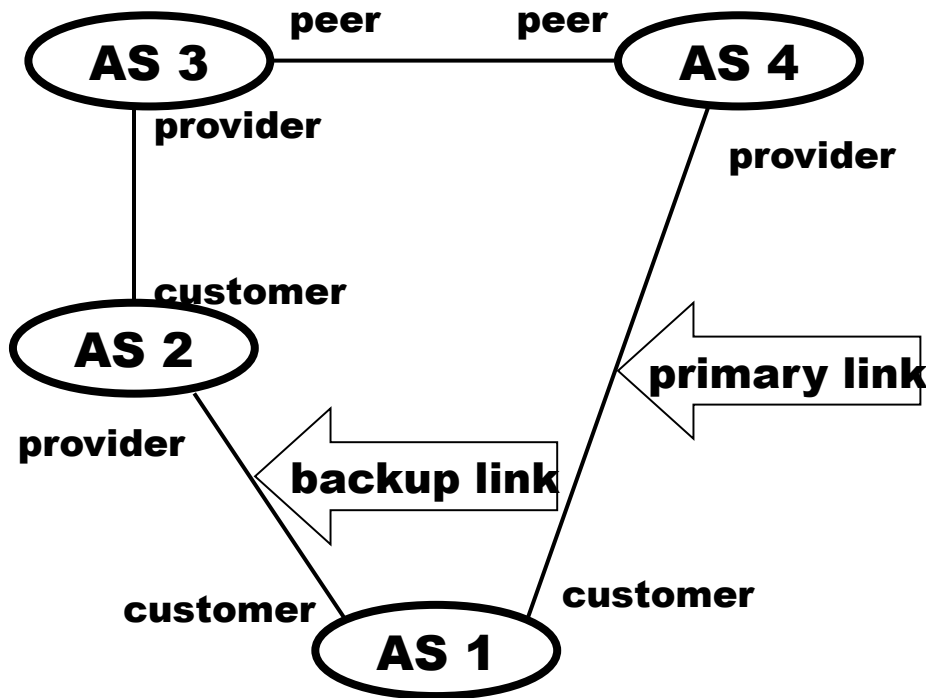
COMMUNITY Attribute to the Rescue!



What is a BGP Wedgie?

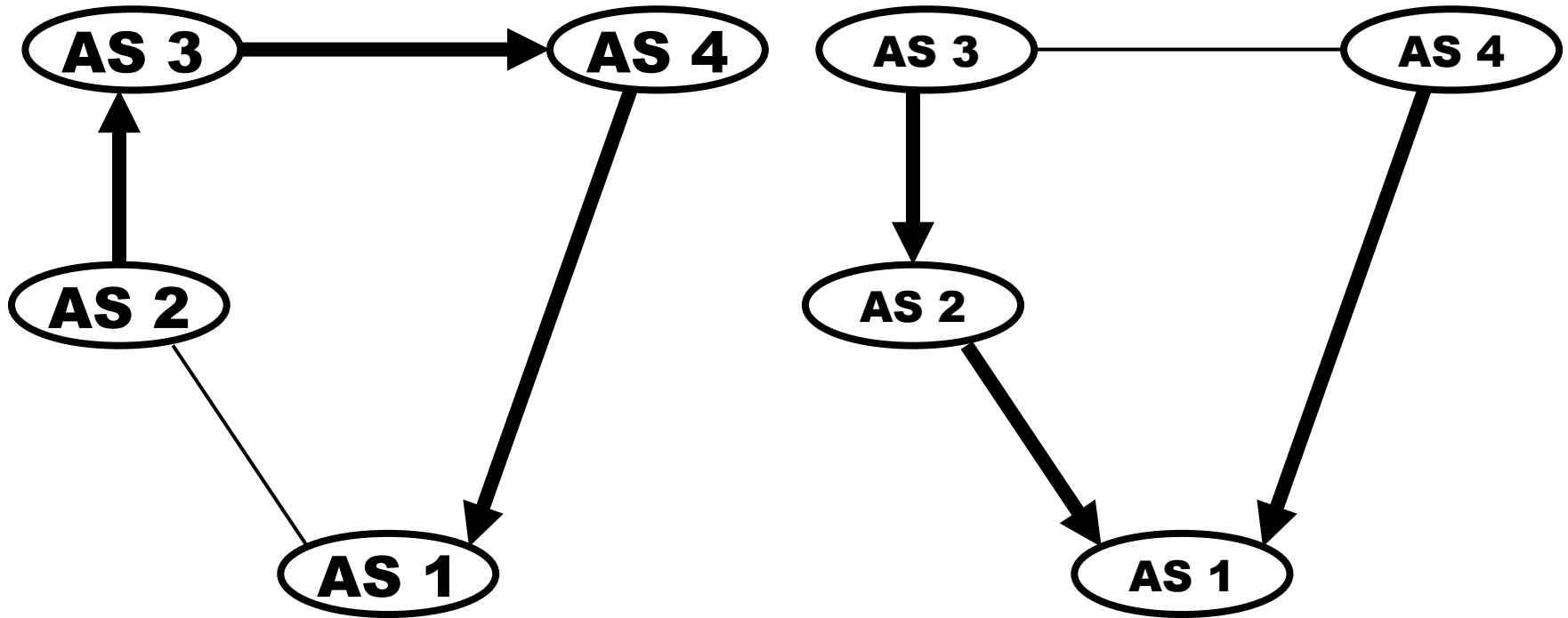
-
- The diagram consists of two large curly braces on the left side. The outer brace is labeled 'full wedgie' and spans the entire vertical range of the list. The inner brace is labeled '3/4 wedgie' and spans the top three items of the list, indicating that these three items are the primary components of a '3/4 wedgie'.
- BGP policies make sense locally
 - Interaction of local policies allows multiple stable routings
 - Some routings are consistent with intended policies, and some are not
 - If an unintended routing is installed (BGP is “wedged”), then manual intervention is needed to change to an intended routing
 - When an unintended routing is installed, no single group of network operators has enough knowledge to debug the problem

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...



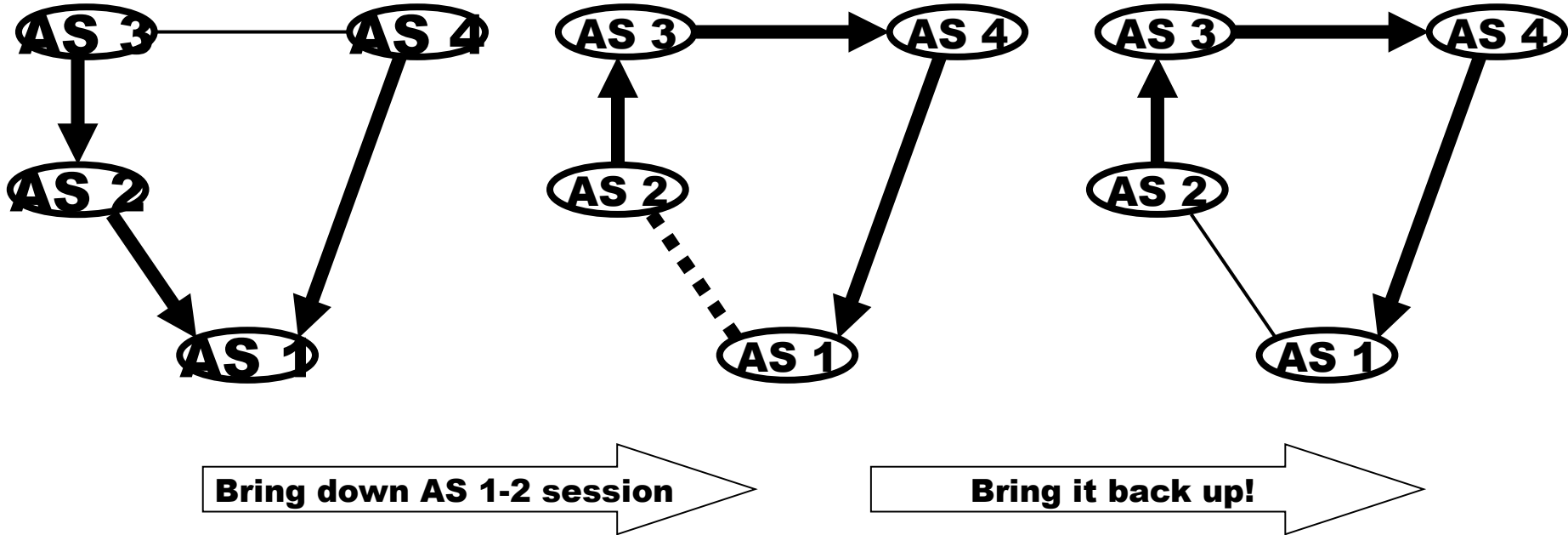
Intended Routing

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

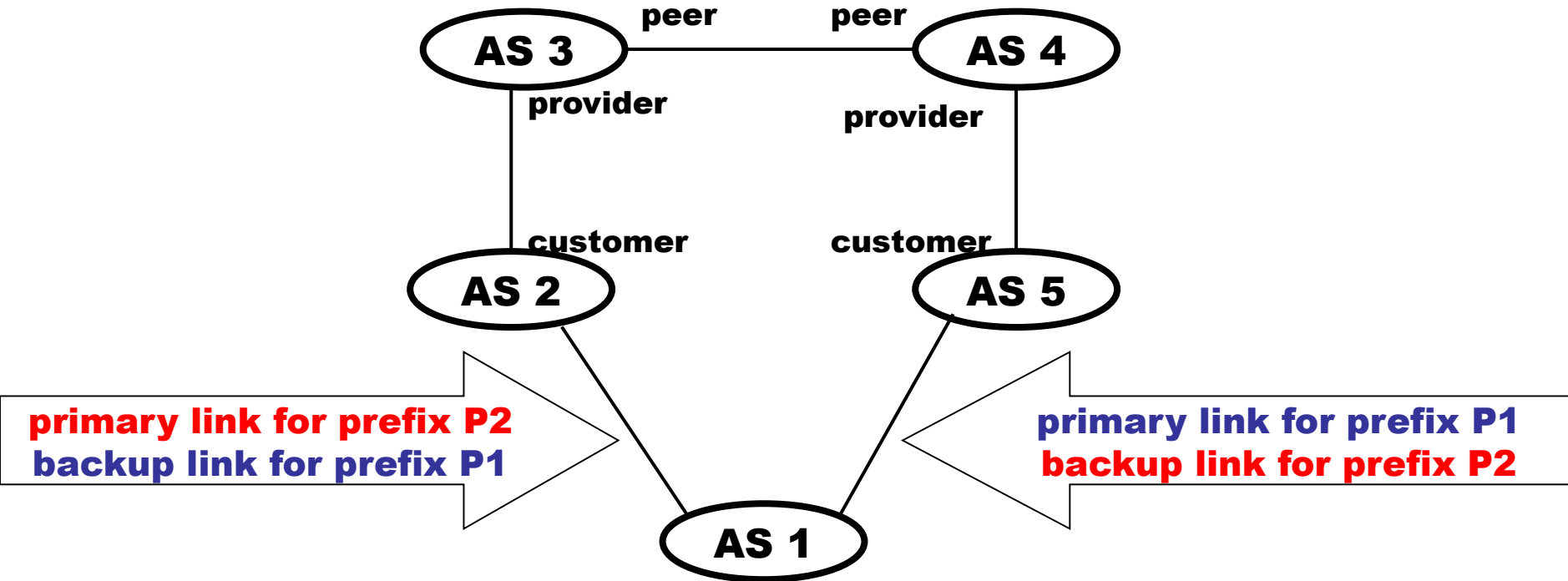


- Requires manual intervention
- Can be done in AS 1 or AS 2

What the heck is going on?

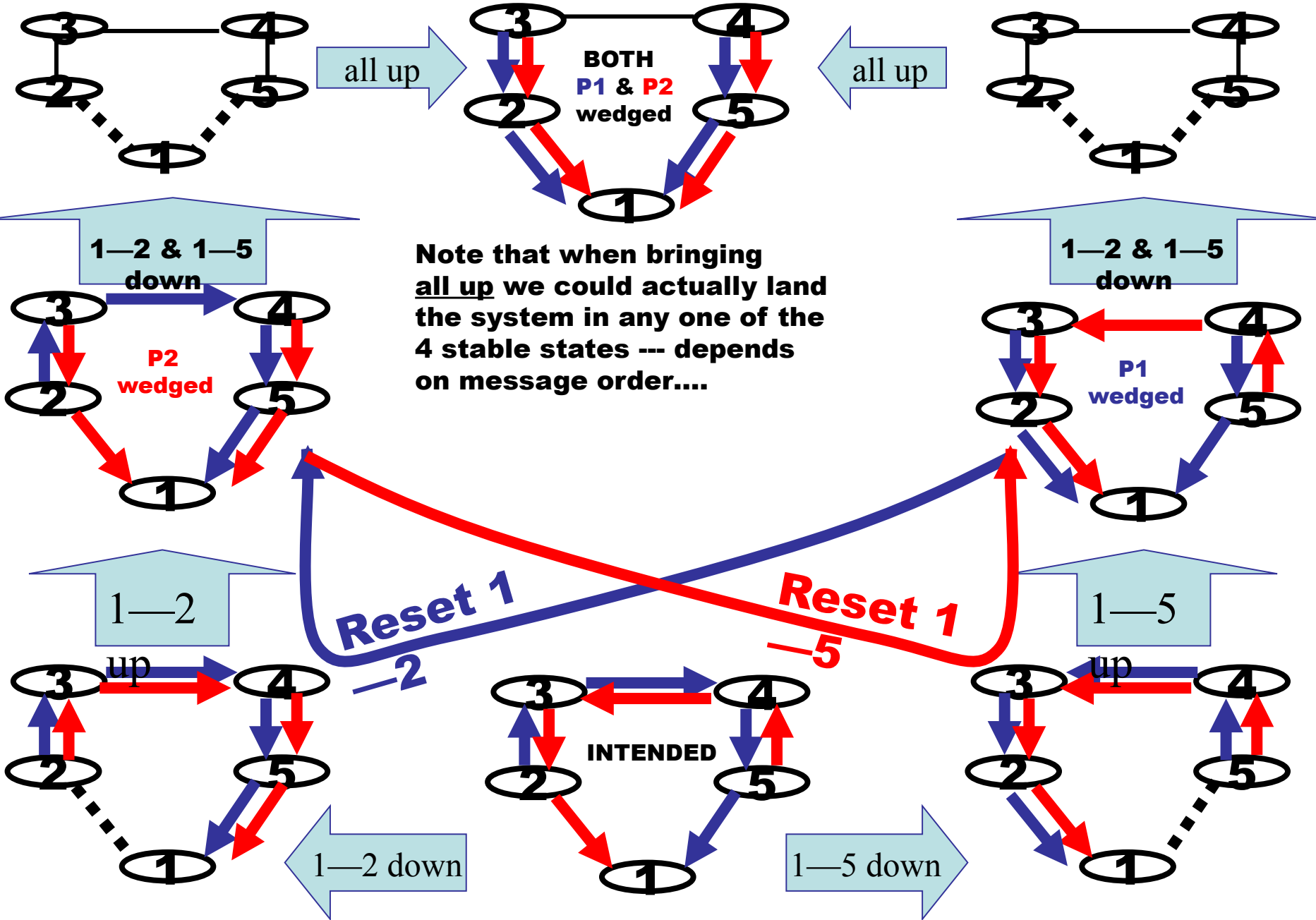
- There is no guarantee that a BGP configuration has a unique routing solution.
 - When multiple solutions exist, the (unpredictable) order of updates will determine which one is wins.
- There is no guarantee that a BGP configuration has any solution!
 - And checking configurations NP-Complete
 - Lab 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



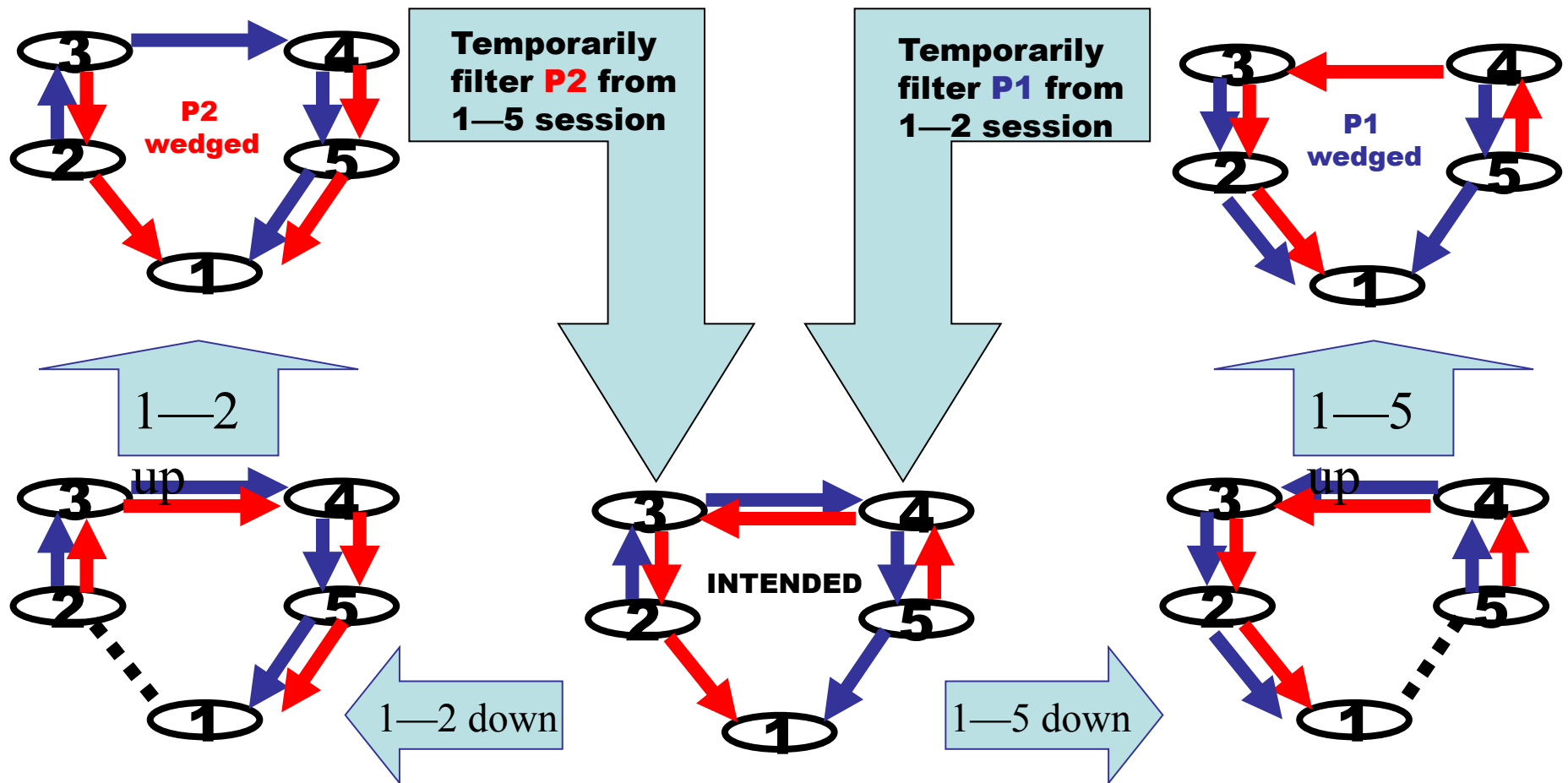
Simple session reset may not work!!

Can't un-wedge with session resets!



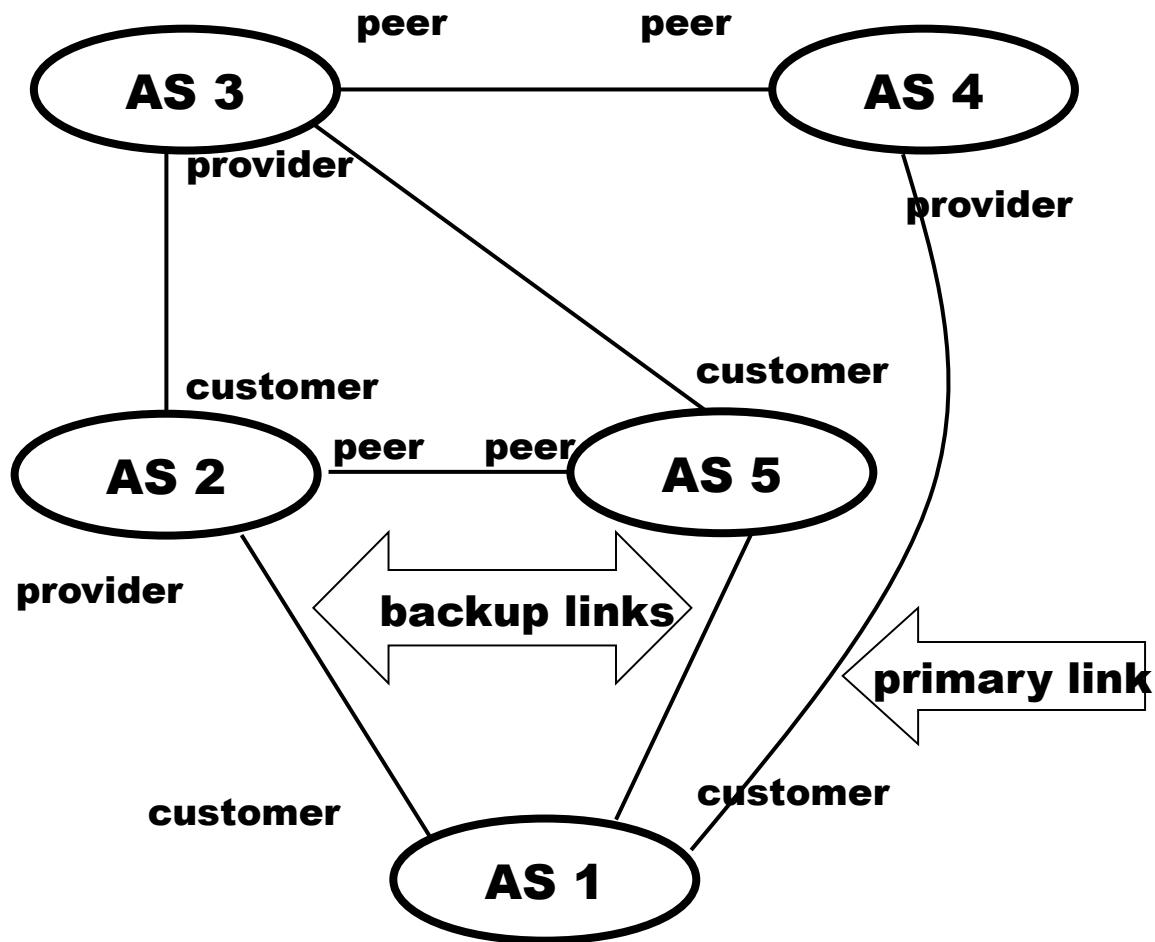
Note that when bringing all up we could actually land the system in any one of the 4 stable states --- depends on message order....

Recovery



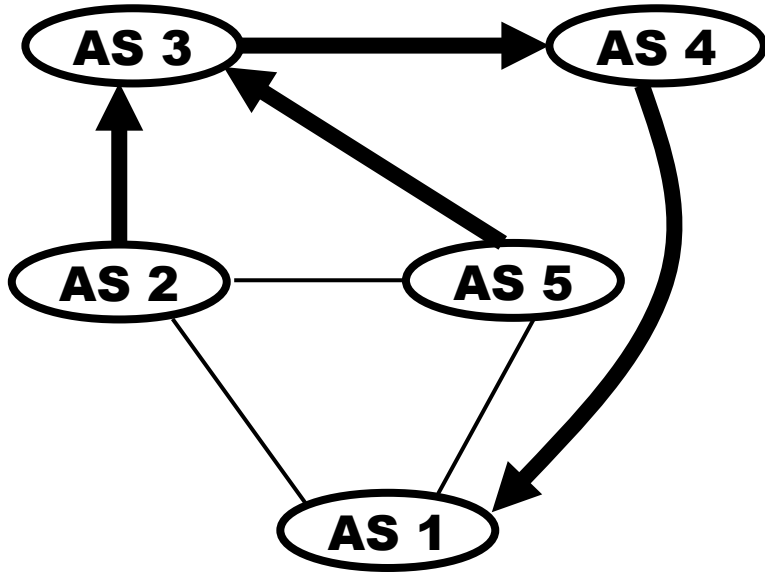
**Who among us could figure this one out?
When 1-2 is in New York and 1-5 is in Tokyo?**

Full Wedgie Example

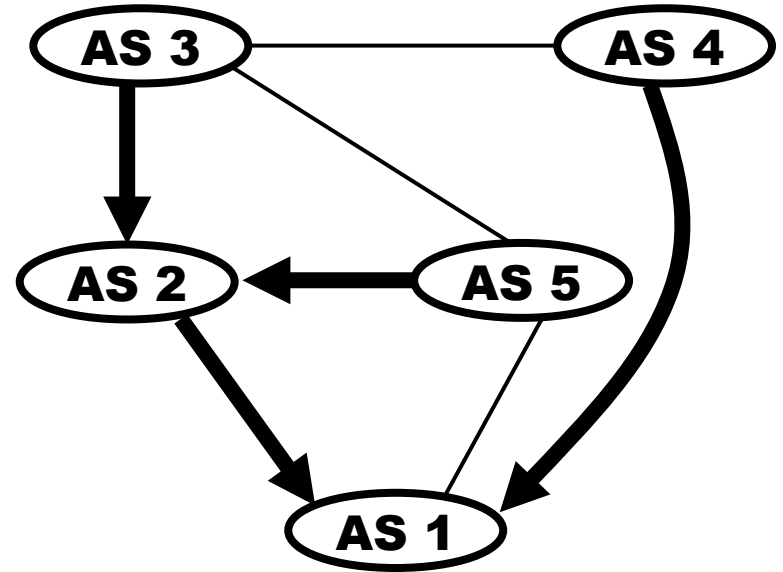


- AS 1 implements backup links by sending AS 2 and AS 3 a “depref me” communities.
- AS 2 implements its community so that the resulting local pref is below that of its upstream providers and its peers (AS 3 and AS 5 routes)
- AS 5 implements its community so that the resulting local pref is below its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

And the Routings are...

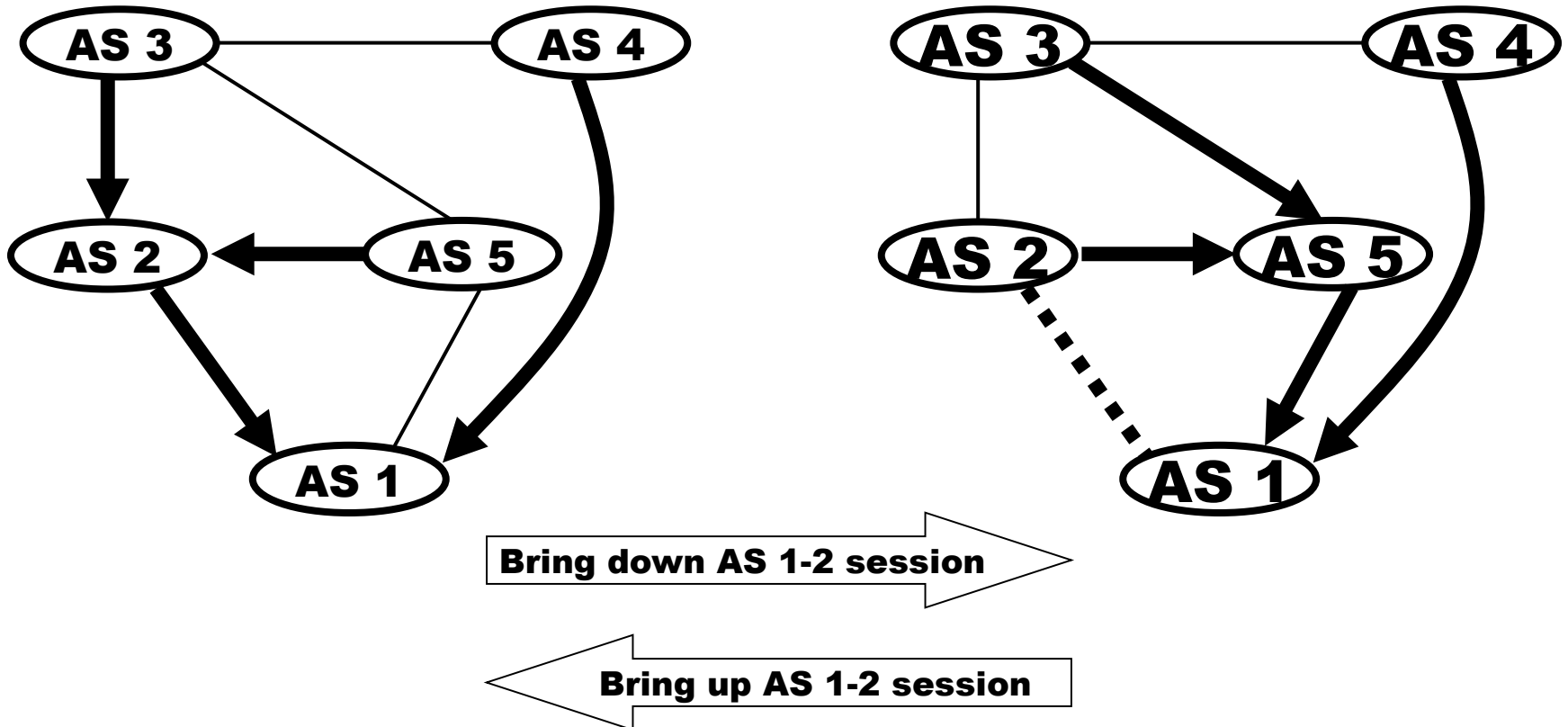


Intended Routing

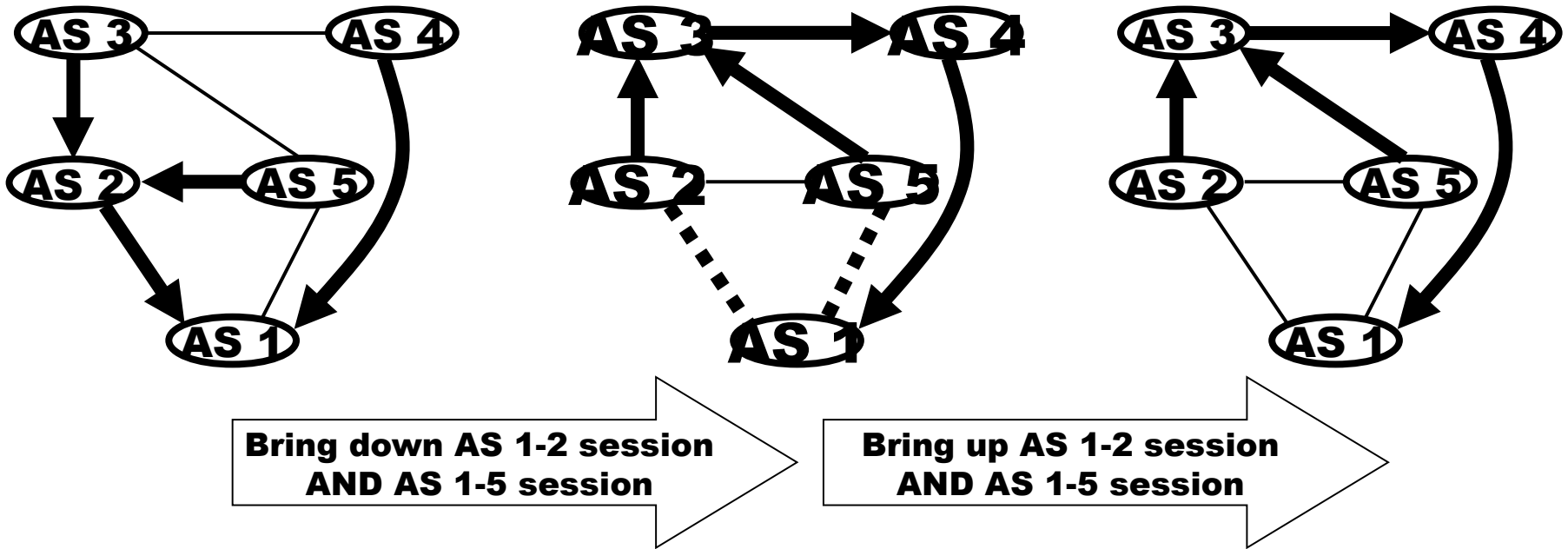


Unintended Routing

Resetting 1–2 does not help!!



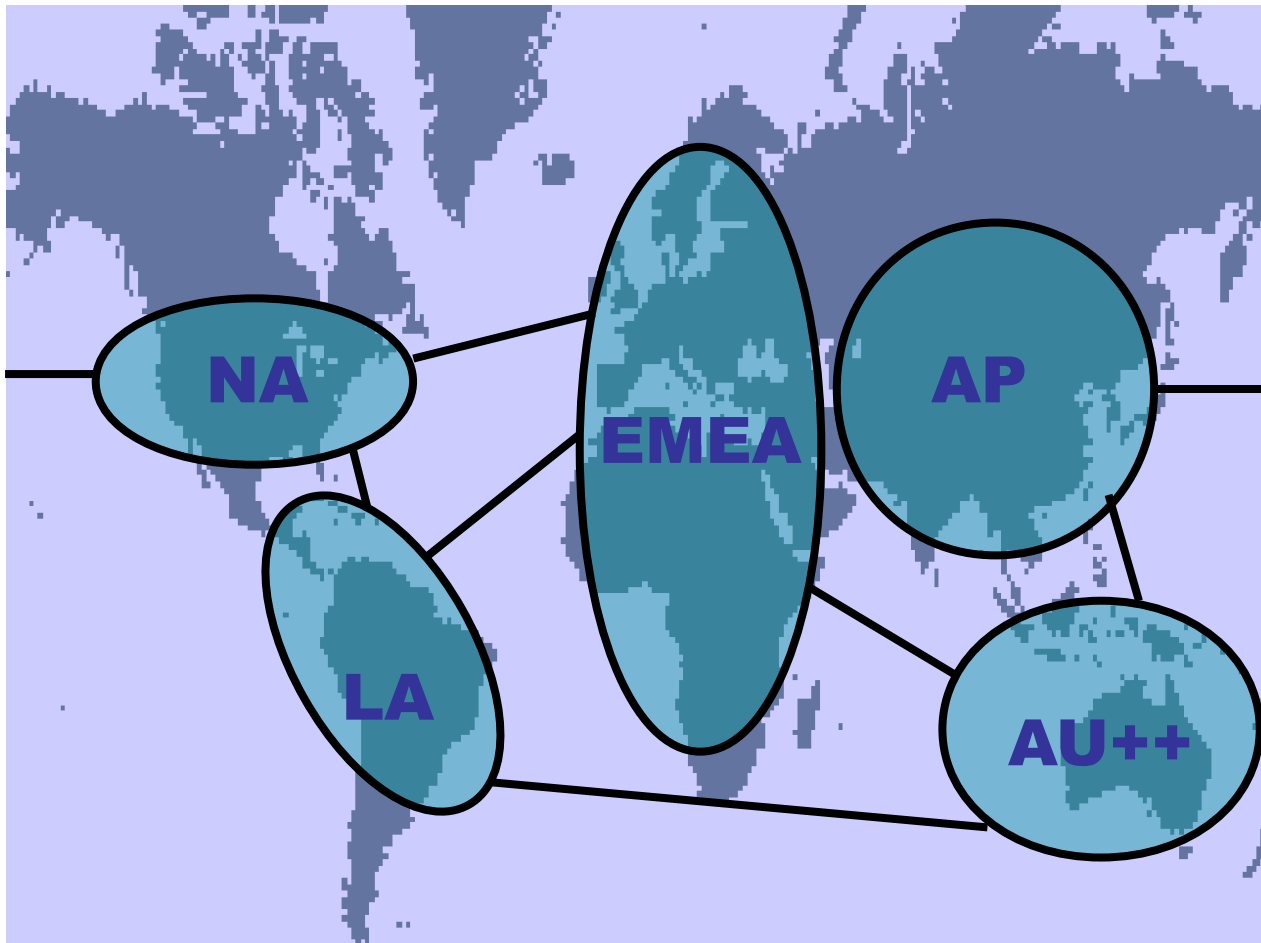
Recovery



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

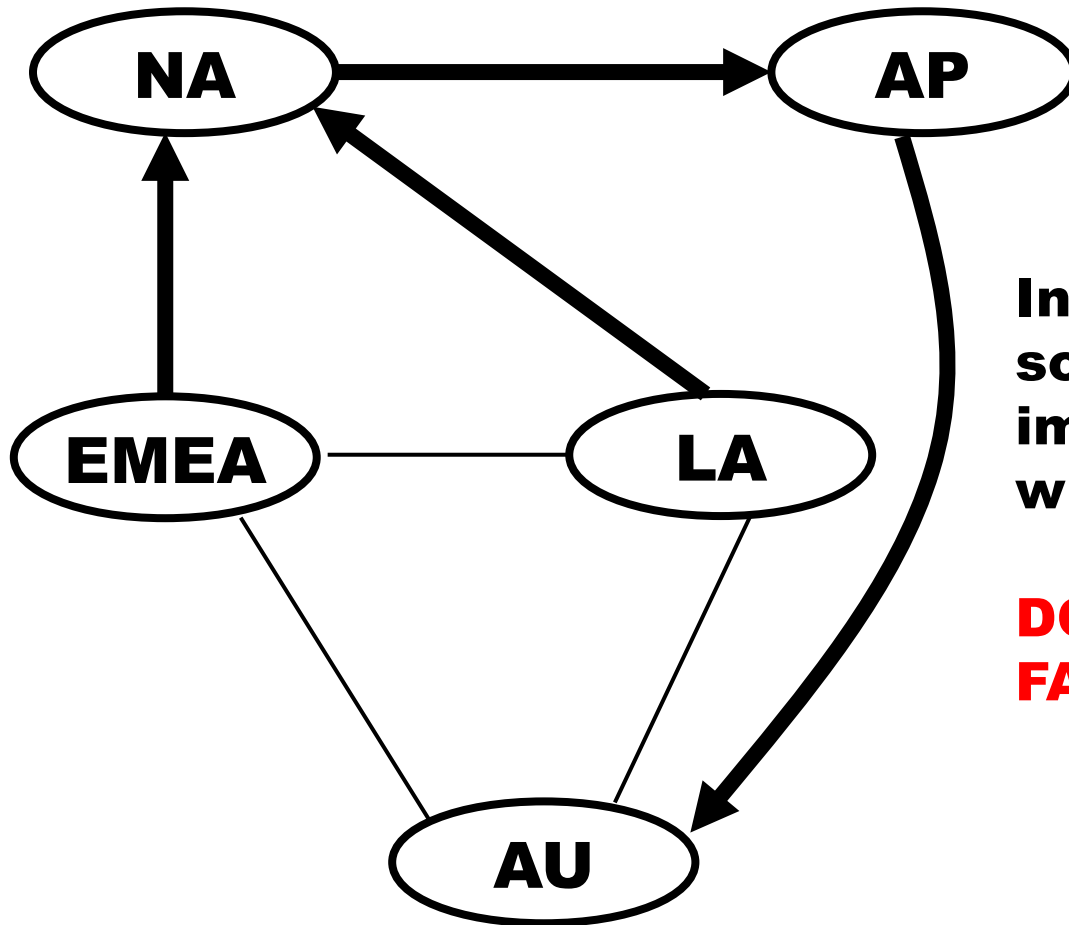
Try to convince AS 5 and AS 1 that their session has been reset (or filtered) even though it is not associated with an active route!

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.