# L11 : Inter-domain Routing with BGP Lecture14 Michaelmas, 2016
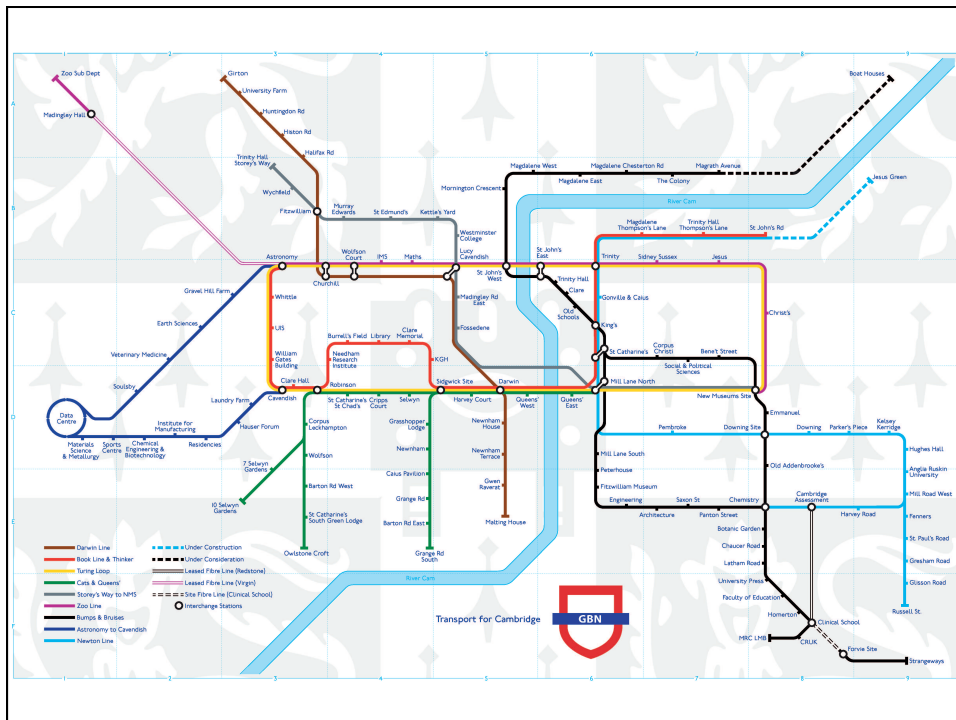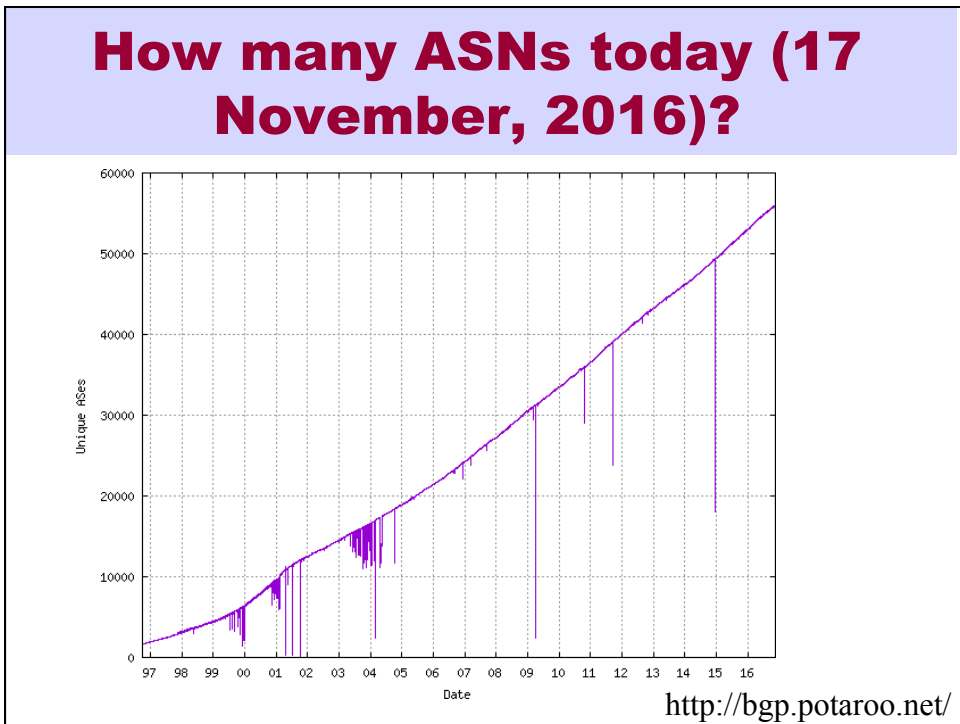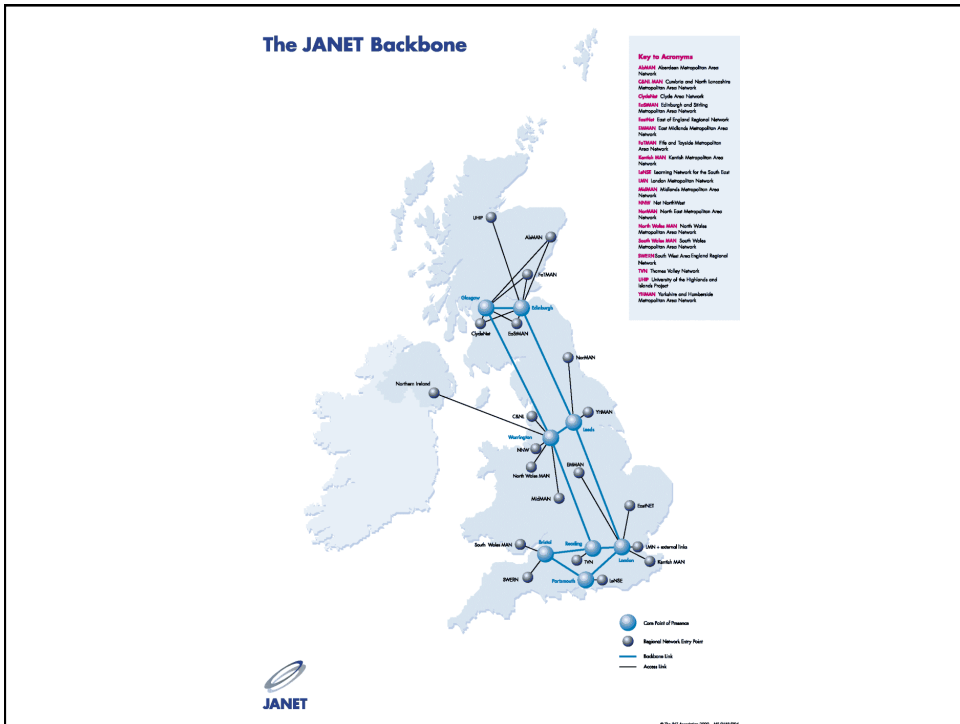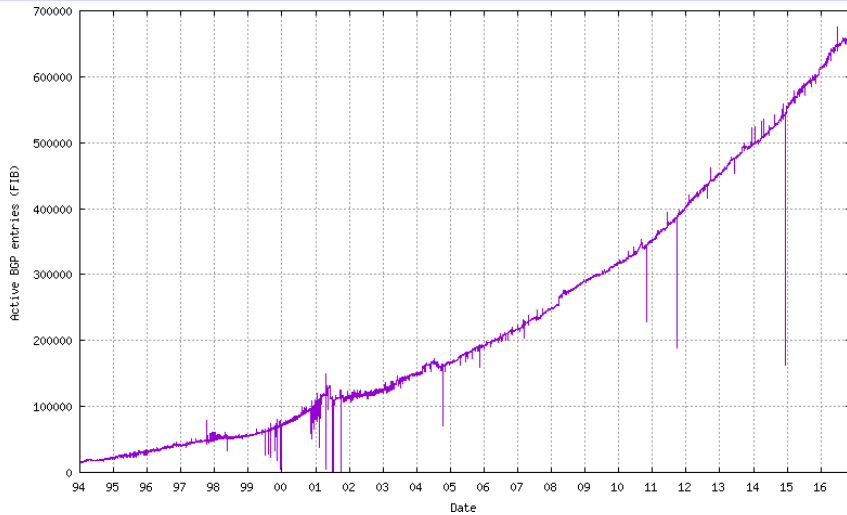
## Timothy G. Griffin
## Computer Lab
## Cambridge UK

The JANET Backbone

# How many ASNs today (17 November, 2016)?



http://bgp.potaroo.net/

# How many prefixes today (17 November, 2016)?



http://bgp.potaroo.net/

# AS Graph != Internet Topology

BGP was designed to throw away information!



ICNP 2002

# BGP Route Attributes

```
Value     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]
...
255       reserved for development
```
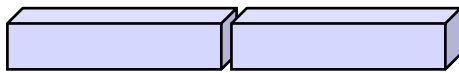
**Most important attributes**

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

Not all attributes
need to be present in
every announcement

---

# How Can Routes be Classified?
# BGP Communities

**A community value is 32 bits**

**Used for signally
within and between
ASes**

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

**community
number**

**Very powerful
BECAUSE it
has no (predefined)
meaning**

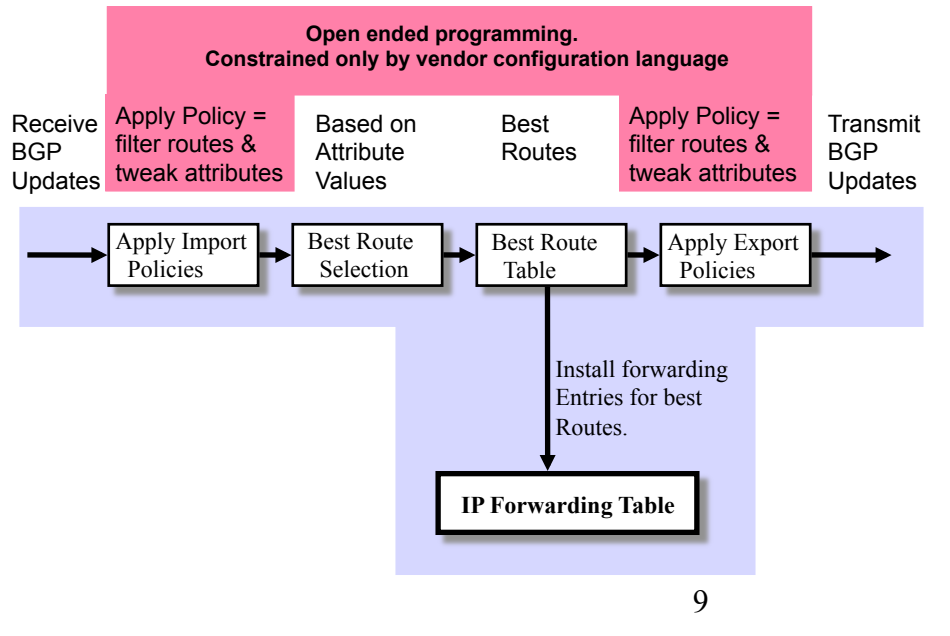**Community Attribute = a list of community values.
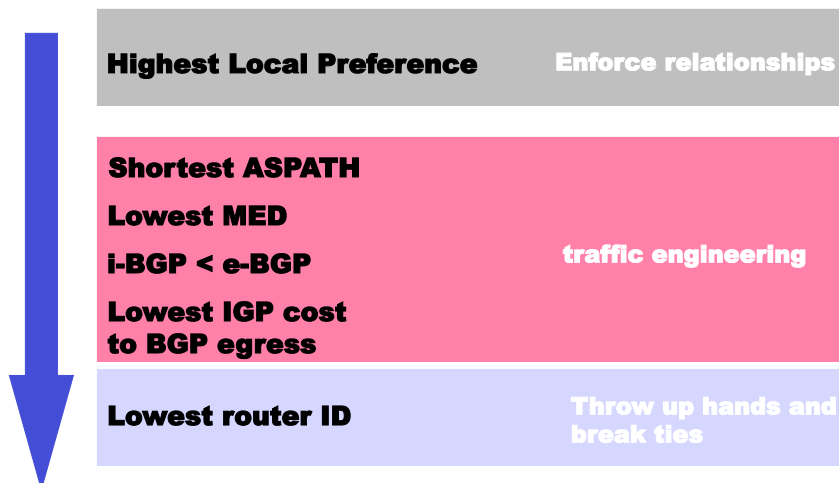(So one route can belong to multiple communities)**

**RFC 1997 (August 1996)**

**Reserved communities**
no_export = 0xFFFFFF01: don't export out of AS
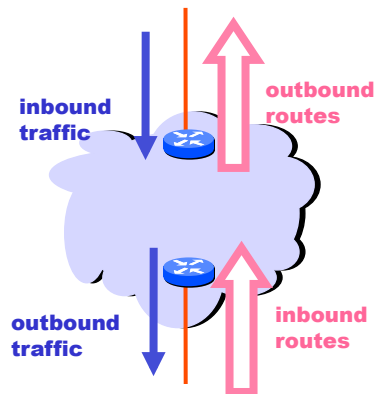no_advertise 0xFFFFFF02: don't pass to BGP neighbors

8

# 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

→ Apply Import Policies → Best Route Selection → Best Route Table → Apply Export Policies →

Install forwarding Entries for best Routes.

**IP Forwarding Table**

9

# Route Selection Summary (A lexicographic product)

**Highest Local Preference** — **Enforce relationships**

**Shortest ASPATH**

**Lowest MED**

**i-BGP < e-BGP**

**Lowest IGP cost to BGP egress**

**traffic engineering**

**Lowest router ID** — **Throw up hands and break ties**
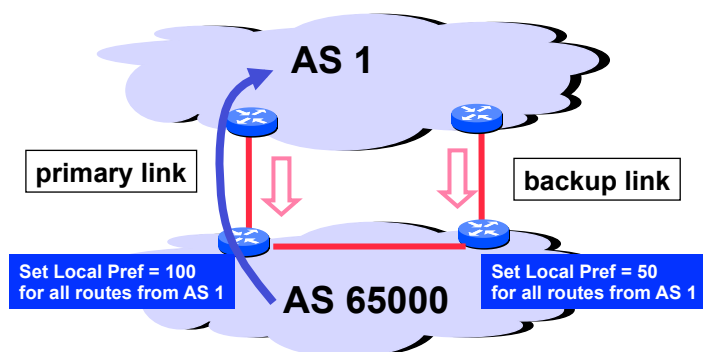
# Traffic Engineering

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

**inbound traffic**

**outbound routes**

**outbound traffic**

**inbound routes**

In general, an AS has more control over outbound traffic

---

# Implementing Backup Links with Local Preference (Outbound Traffic)

**AS 1**

**primary link**

**backup link**

Set Local Pref = 100 for all routes from AS 1

Set Local Pref = 50 for all routes from AS 1

**AS 65000**

Forces <u>outbound</u> traffic to take primary link, unless link is down.

We'll talk about <u>inbound</u> traffic soon …
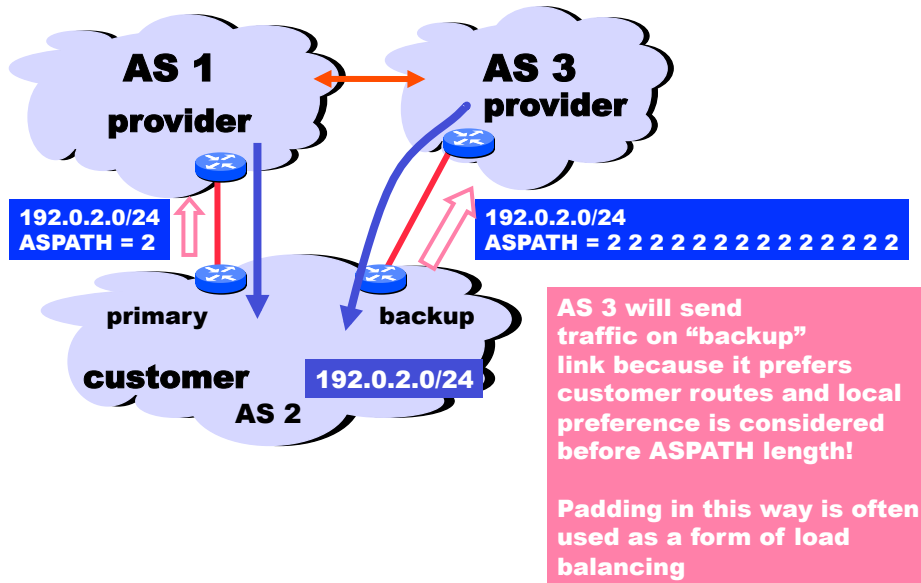
12

# Multihomed Backups
## (Outbound Traffic)

**AS 1**
provider

**AS 3**
provider

| primary link |

| backup link |

Set Local Pref = 100
for all routes from AS 1

Set Local Pref = 50
for all routes from AS 3

**AS 2**

Forces <u>outbound</u> traffic to take primary link, unless link is down.

13

---

# Shedding Inbound Traffic with ASPATH Padding.  Yes, this is an ugly hack ...

**AS 1**    provider

192.0.2.0/24
ASPATH = 2

192.0.2.0/24
ASPATH = 2  2  2

primary

backup

Padding will (usually)
force inbound
traffic from AS 1
to take primary link

**customer**    192.0.2.0/24

AS 2

14

# ... But Padding Does Not Always Work

**AS 1**
**provider**

**AS 3**
**provider**

| 192.0.2.0/24 |
| ASPATH = 2 |

| 192.0.2.0/24 |
| ASPATH = 2 2 2 2 2 2 2 2 2 2 2 2 2 2 |

**primary**

**backup**

**customer**

**AS 2**

| 192.0.2.0/24 |

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!

**AS 1**
**provider**

**AS 3**
**provider**

AS 3: normal customer local pref is 100, peer local pref is 90

| 192.0.2.0/24 |
| ASPATH = 2 |

| 192.0.2.0/24 |
| ASPATH = 2 |
| COMMUNITY = 3:70 |

**primary**

**backup**

**customer**

**AS 2**

| 192.0.2.0/24 |

Customer import policy at AS 3:
If 3:90 in COMMUNITY then
   set local preference to 90
If 3:80 in COMMUNITY then
   set local preference to 80
If 3:70 in COMMUNITY then
   set local preference to 70

16

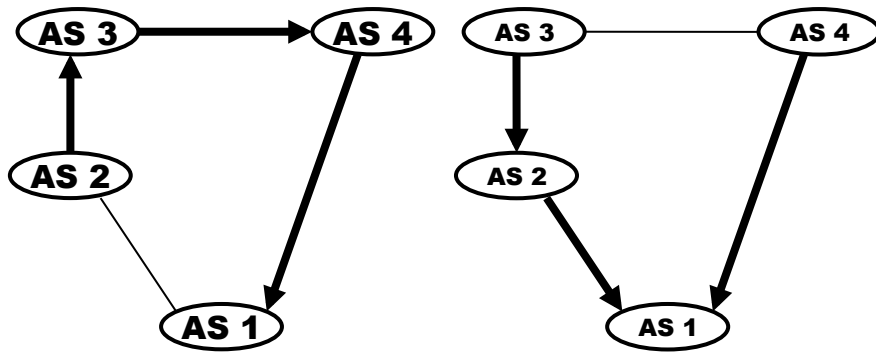## Problem scenario with inter-domain routing in the Internet (BGP)

- 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

## Simple 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 it's upstream provider (AS 3 routes)
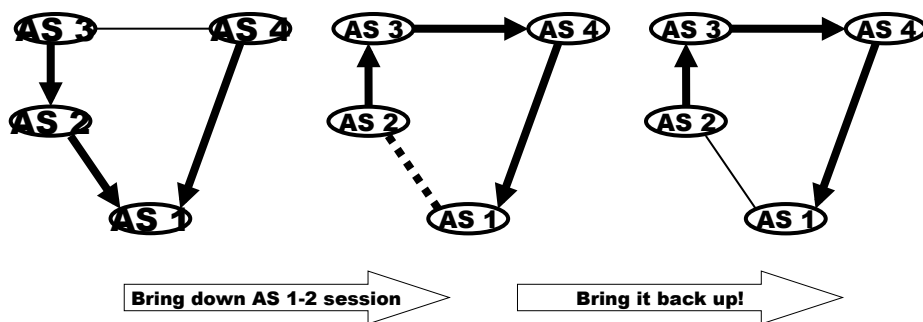
# And the Routings are...



## 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

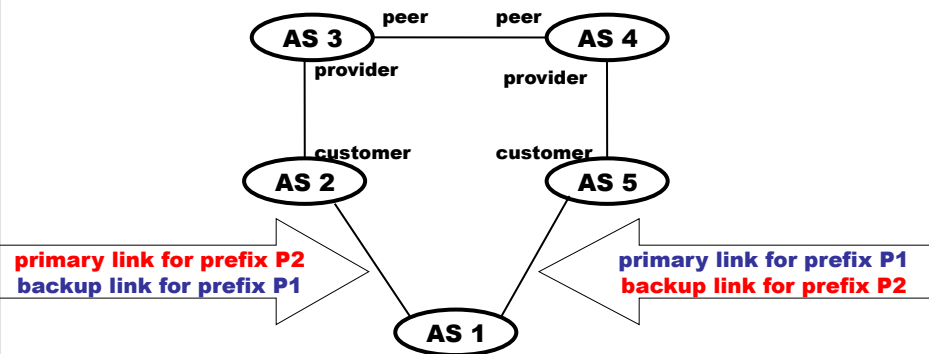Note: This is easy to reach from the intended routing just by "bouncing" the BGP session on the primary link.

---

# Recovery



Bring down AS 1-2 session →

Bring it back up! →

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

# What 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
- Complex policies (weights, communities setting preferences, and so on) increase chances of routing anomalies.
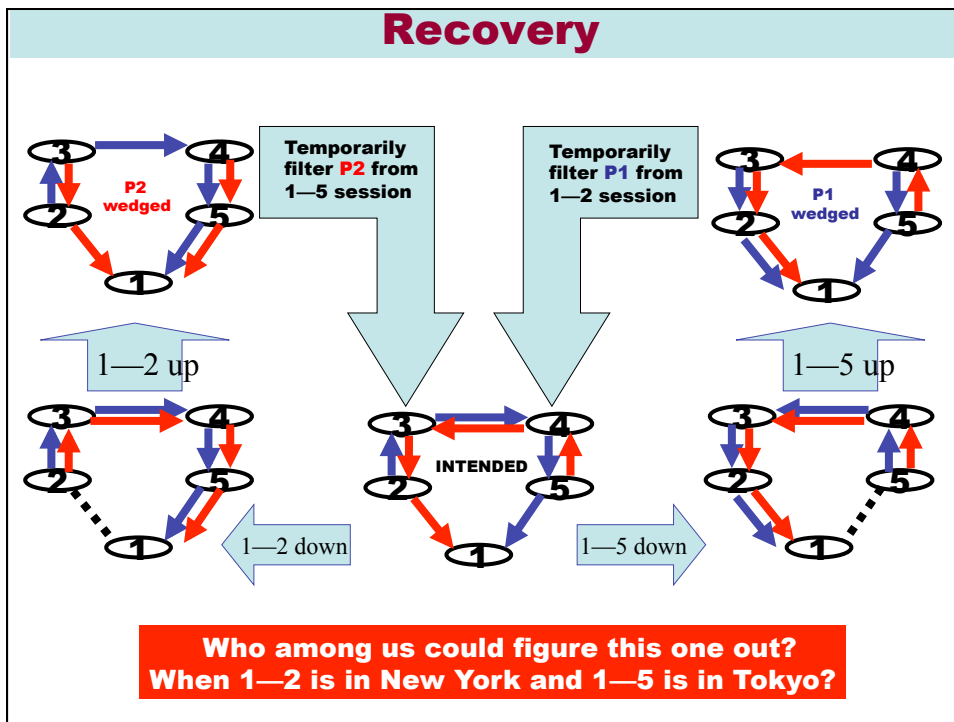  - … yet this is the current trend!

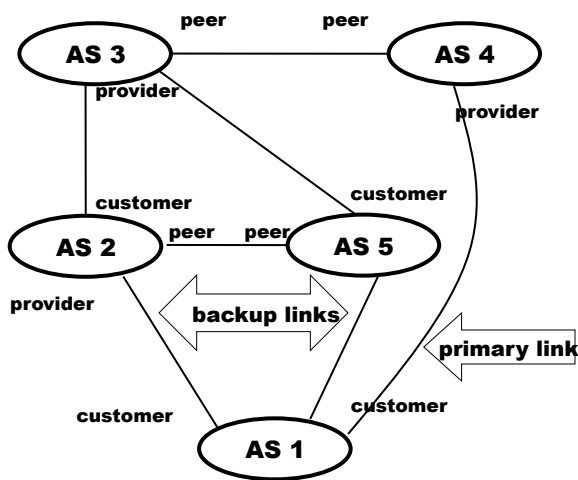# Load Balancing Example



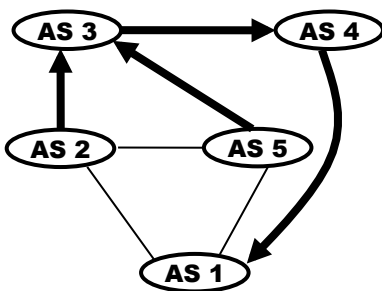Simple session reset my not work!!

# Can't un-wedge with session resets!

all up

BOTH
P1 & P2
wedged

all up

1—2 & 1—5
down

P2
wedged

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

1—2 & 1—5
down

P1 wedged

1—2 up

Reset 1
—2

Reset 1
—5

1—5 up

1—2 down

INTENDED

1—5 down



# Recovery

P2
wedged

Temporarily
filter P2 from
1—5 session

Temporarily
filter P1 from
1—2 session

P1
wedged

1—2 up

1—5 up

1—2 down

INTENDED

1—5 down

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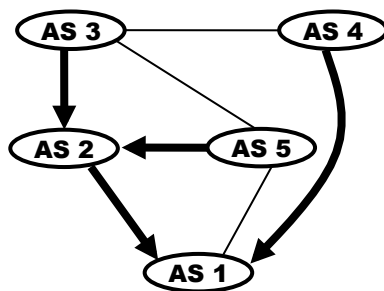When 1—2 is in New York and 1—5 is in Tokyo?**

# Advanced 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 it's 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)
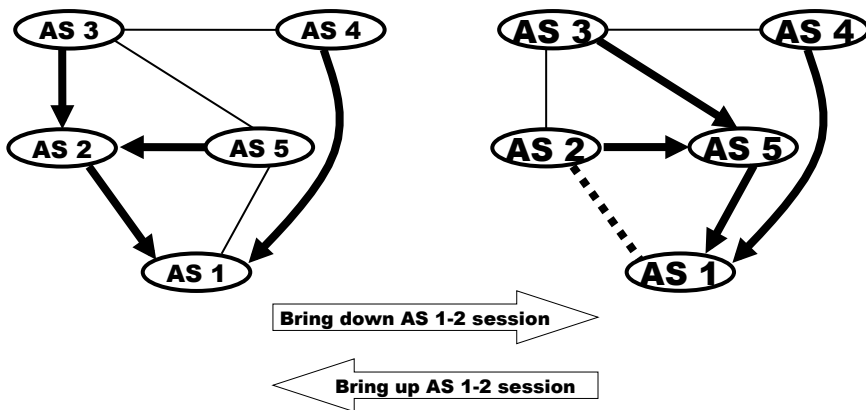
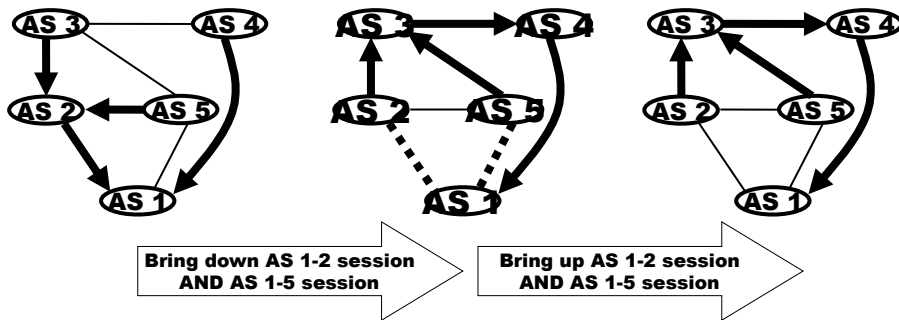# And the Routings are...



**Intended Routing**          **Unintended Routing**

# Resetting 1—2 does not help!!



Bring down AS 1-2 session

Bring up AS 1-2 session

# Recovery



Bring down AS 1-2 session
AND AS 1-5 session

Bring up AS 1-2 session
AND AS 1-5 session
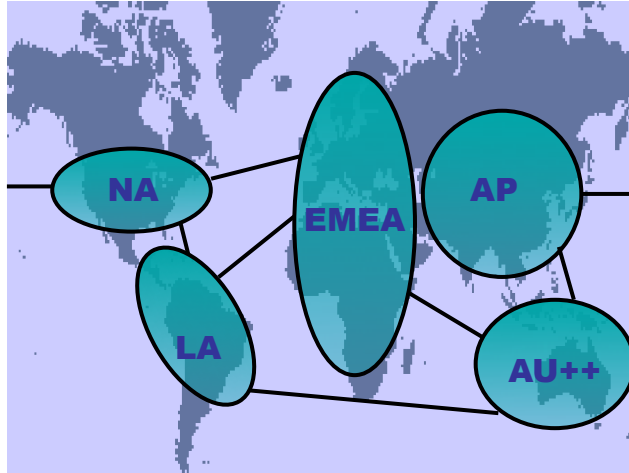
A lot of "non-local" knowledge is required to arrive at
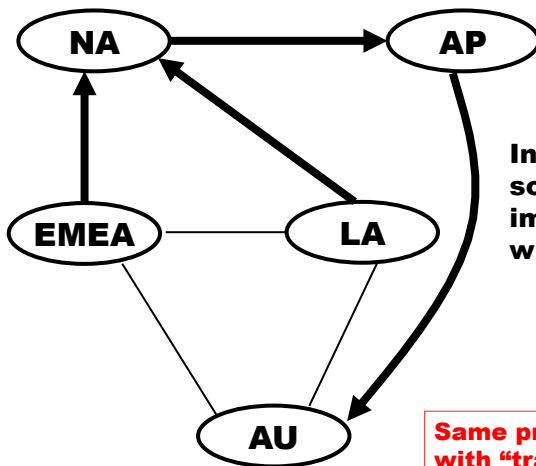this recovery strategy!

Try to convince AS 5 and AS 1 that their session has be
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



**Does this look familiar?**

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

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