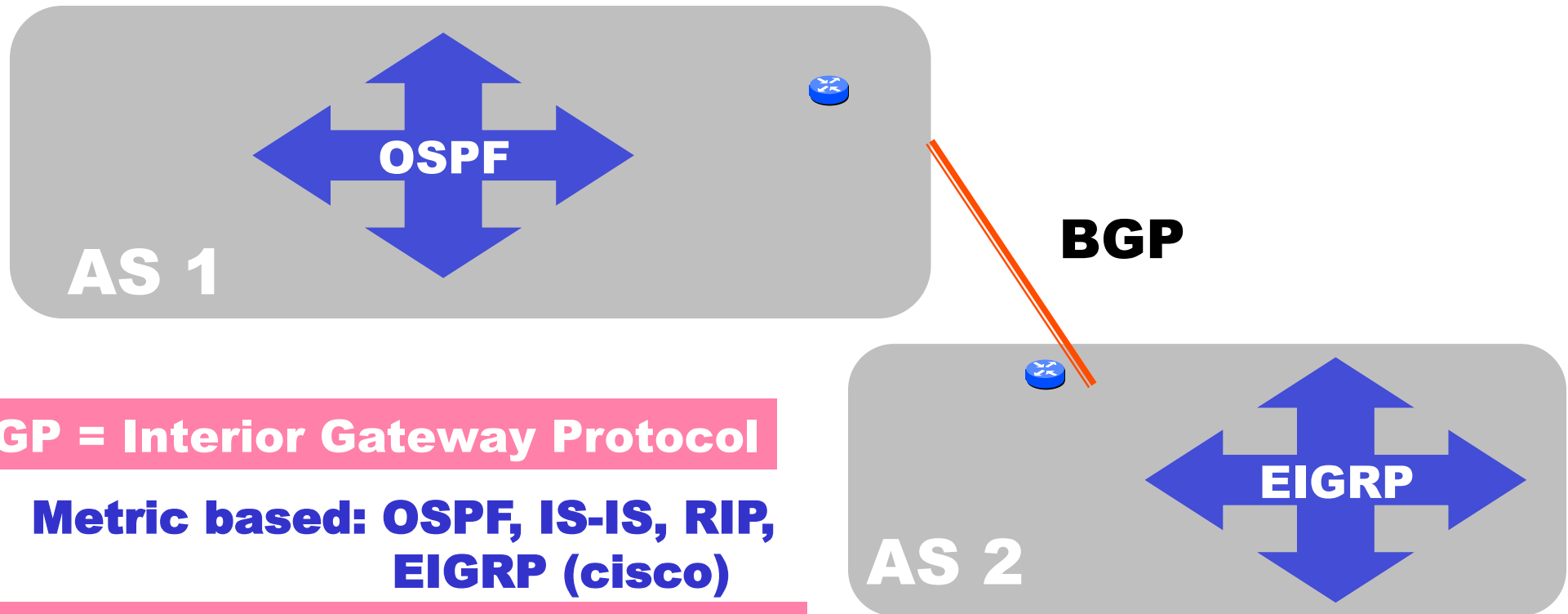# Inter-Domain Routing & BGP

# Outline

The glue that holds the Internet together : interdomain routing with The Border Gateway Protocol (BGP)

# Architecture of Dynamic Routing

**AS 1**

**OSPF**

**BGP**

**EIGRP**

**AS 2**

**IGP = Interior Gateway Protocol**

Metric based: OSPF, IS-IS, RIP, EIGRP (cisco)

**EGP = Exterior Gateway Protocol**

Policy based: BGP

The Routing Domain of BGP is the entire Internet

# Technology of Distributed Routing

## Link State

- Topology information is <u>flooded</u> 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 <u>shared</u> and <u>uniform</u>
- 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
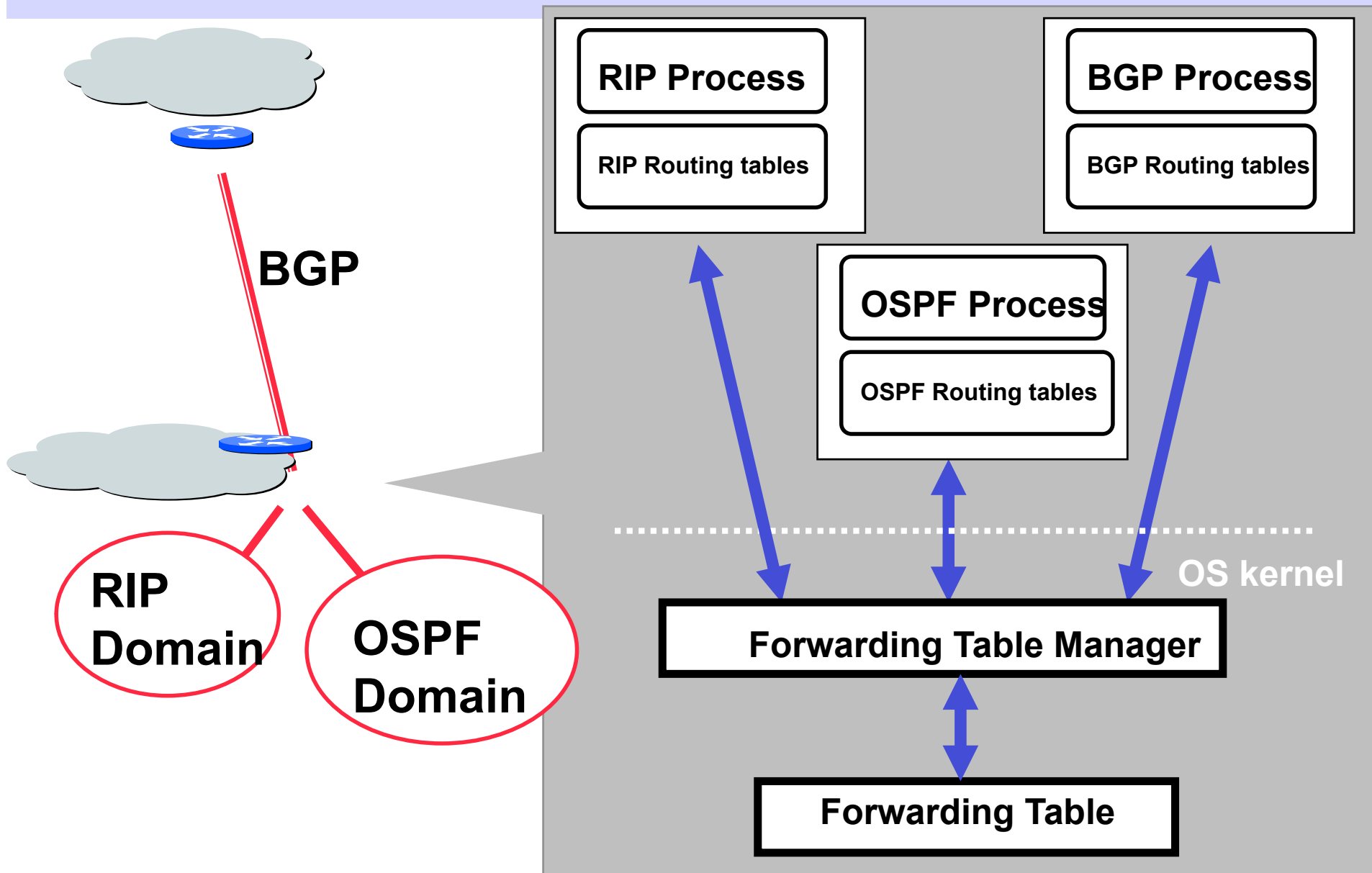
# The Gang of Four

|  | Link State | Vectoring |
|---|---|---|
| IGP | OSPF<br>IS-IS | RIP |
| EGP |  | BGP |

# Many Routing Processes Can Run on a Single Router



BGP

RIP Domain

OSPF Domain

RIP Process

RIP Routing tables

BGP Process

BGP Routing tables

OSPF Process

OSPF Routing tables

OS kernel

Forwarding Table Manager

Forwarding Table

6

# Internet Hierarchy

❖ What is an Autonomous System (AS)?

➢ A set of routers under a single technical administration, using an *intra-domain routing protocol* (IGP) and common metrics to route packets within the AS and using an *inter-domain routing protocol* (EGP) to route packets to other ASes

➢ Sometimes ASes use multiple intra-domain routing protocols and metrics, but appear as a single AS to other ASes

❖ Each AS is assigned a unique ID

# AS Numbers (ASNs)

ASNs are 16 bit values.
64512 through 65535 are "private"

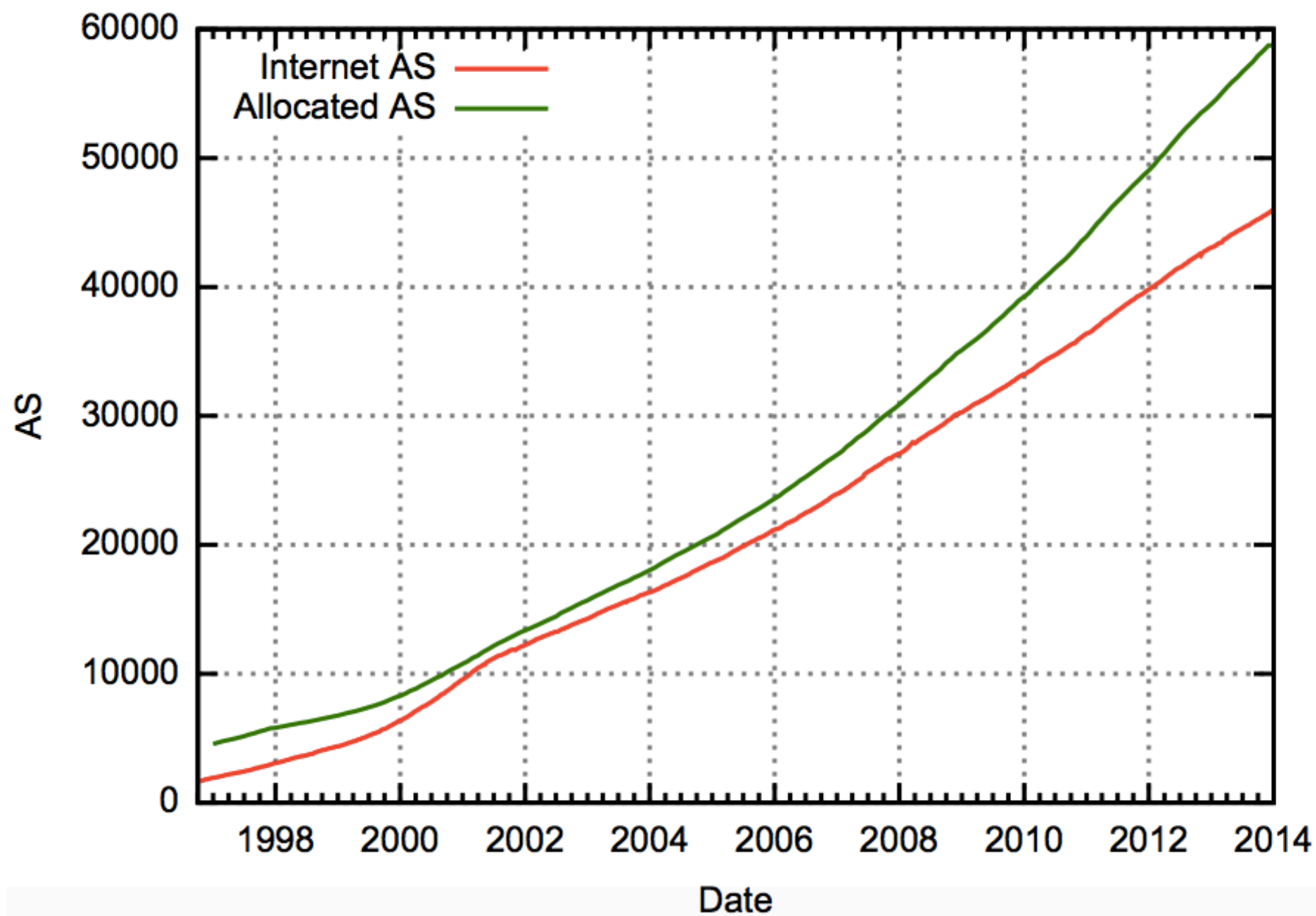Currently over 12,000 in use.

- **Yale: 29**
- **MIT: 3**
- **Harvard: 11**
- **Genuity: 1**
- **AT&T: 7018, 6341, 5074, ...**
- **UUNET: 701, 702, 284, 12199, ...**
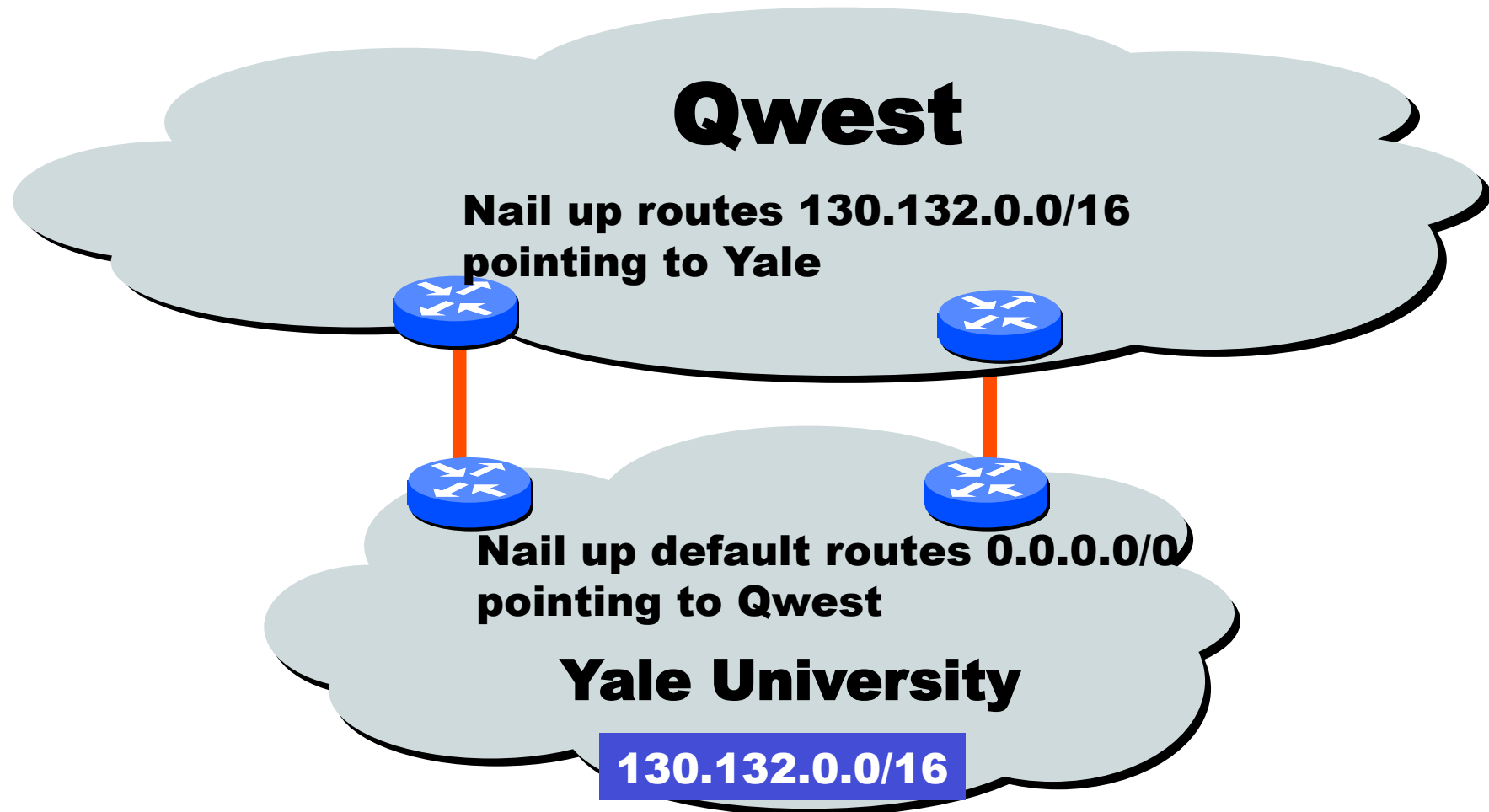- **Sprint: 1239, 1240, 6211, 6242, ...**
- **...**

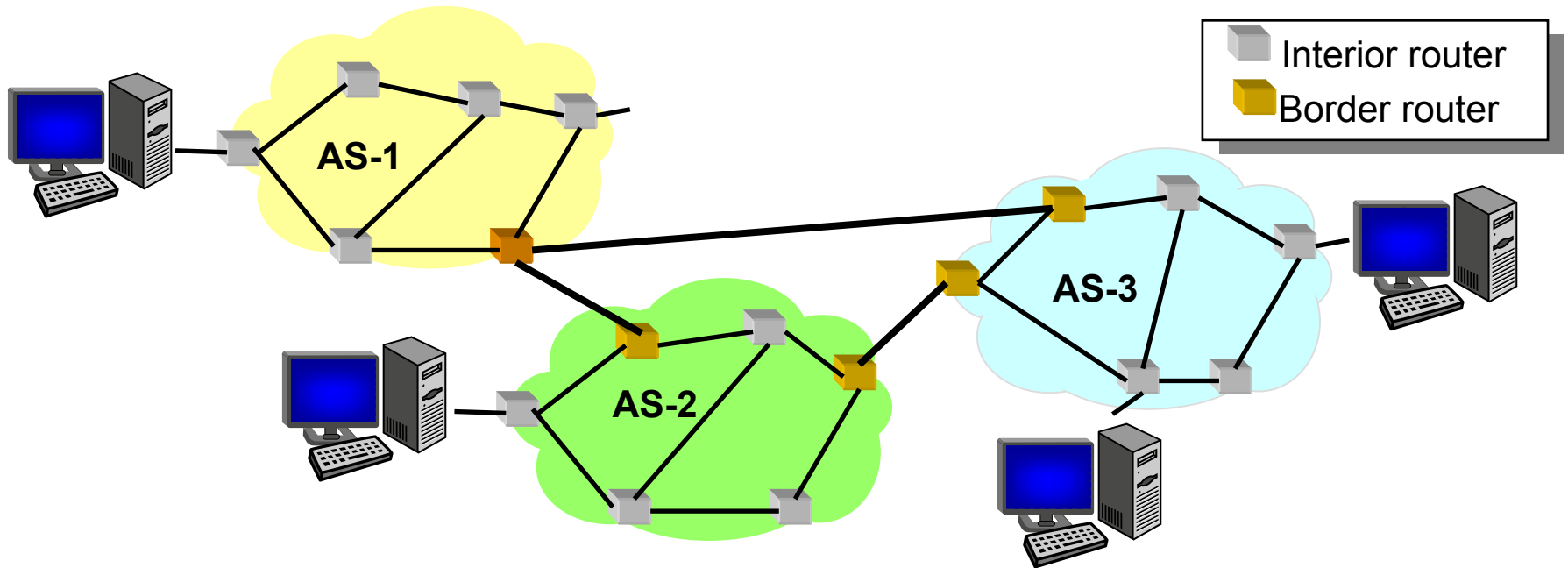ASNs represent units of routing policy

AS announced
on the Internet

# Autonomous Routing Domains Don't Always Need BGP or an ASN

## Qwest

Nail up routes 130.132.0.0/16 pointing to Yale

Nail up default routes 0.0.0.0/0 pointing to Qwest

## Yale University

130.132.0.0/16

Static routing is the most common way of connecting an autonomous routing domain to the Internet.
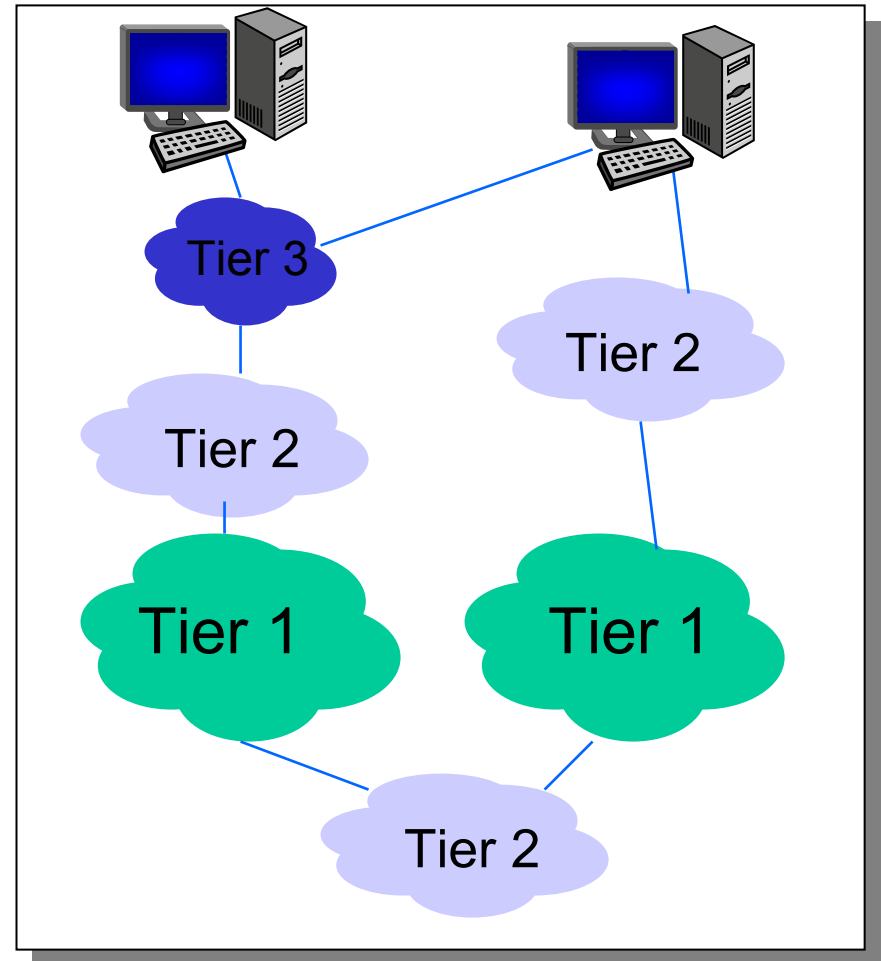This helps explain why BGP is a mystery to many ...

# Picture of the Internet



Figures by MIT OpenCourseWare.

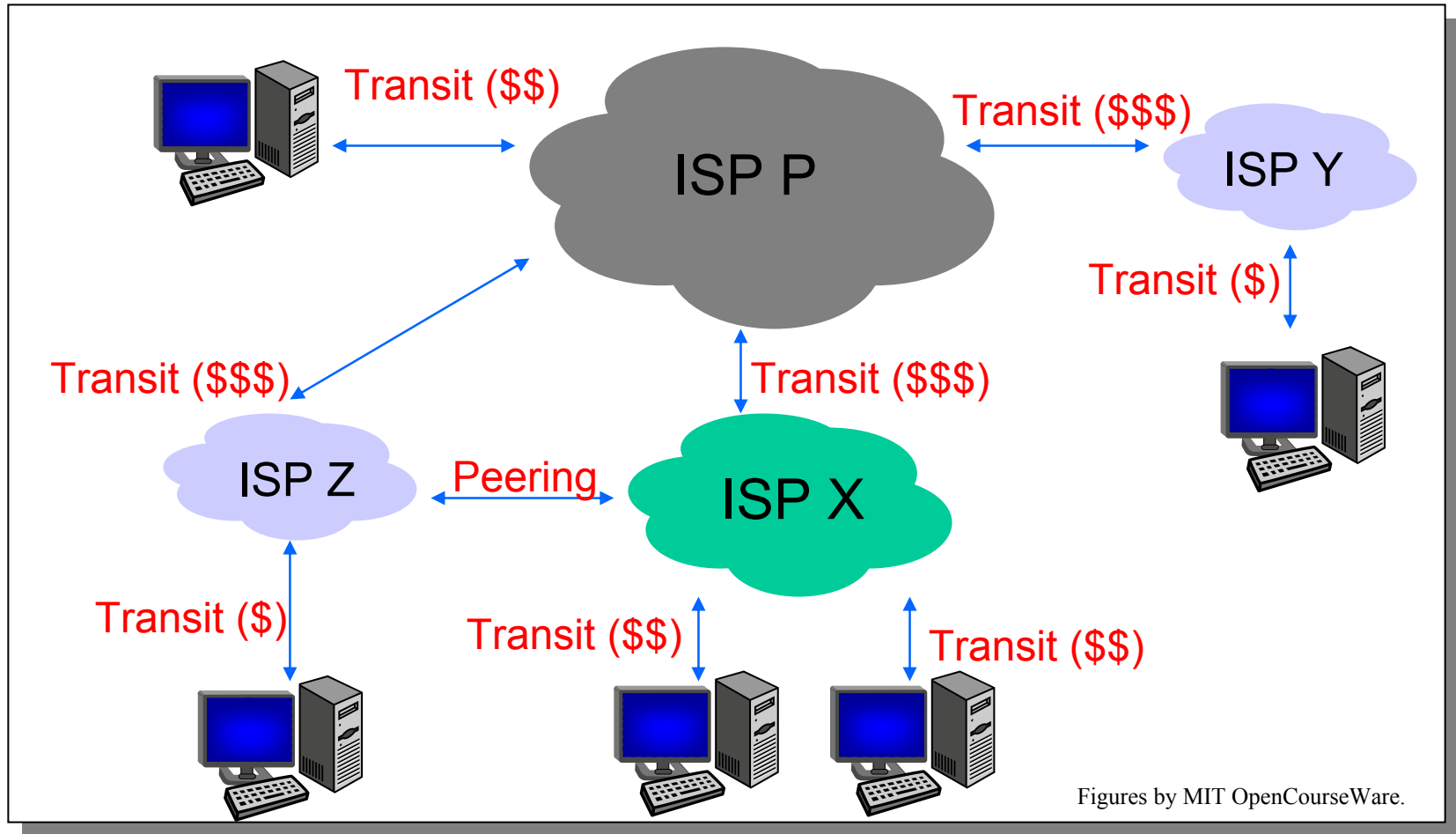- ❖ Intra-domain routing inside an AS
- ❖ Iner-domain routing between ASes

# A Logical View of the Internet

❖ Internet connectivity is provided by commercial entities called ISPs, who compete for profit yet have to cooperate to provide connectivity
  ➢ Each ISP has its own AS (sometimes multiple ASes)

❖ Not all ISPs are created equal
  ➢ Tier 1 ISP
    • "Default-free" global reachability info
  ➢ Tier 2 ISP
    • Regional or country-wide
  ➢ Tier 3 ISP
    • Local

# Inter-AS Relationship: Transit vs. Peering



Figures by MIT OpenCourseWare.

# Policy Impact on Routing

❖ **AS relationships**
  ➢ Customer-provider
  ➢ Peers

❖ **Want "Valley-free" routes**
  ➢ Number links as (+1, 0, -1) for provider, peer and customer links
  ➢ In any path, you should only see sequence of +1, followed by at most one 0, followed by sequence of -1

# Customer-Provider Hierarchy



provider ●——→ customer

◄—————— IP traffic

36

# 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



| | | |
|---|---|---|
| **peer** | ●———● | **peer** |
| **provider** | ●———▶ | **customer** |

38

# Policy-Based Routing

❖ Policies are used to force customer-provider-peer relationships, backup links, load balancing, …

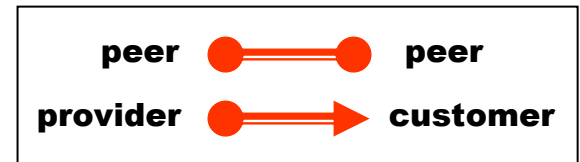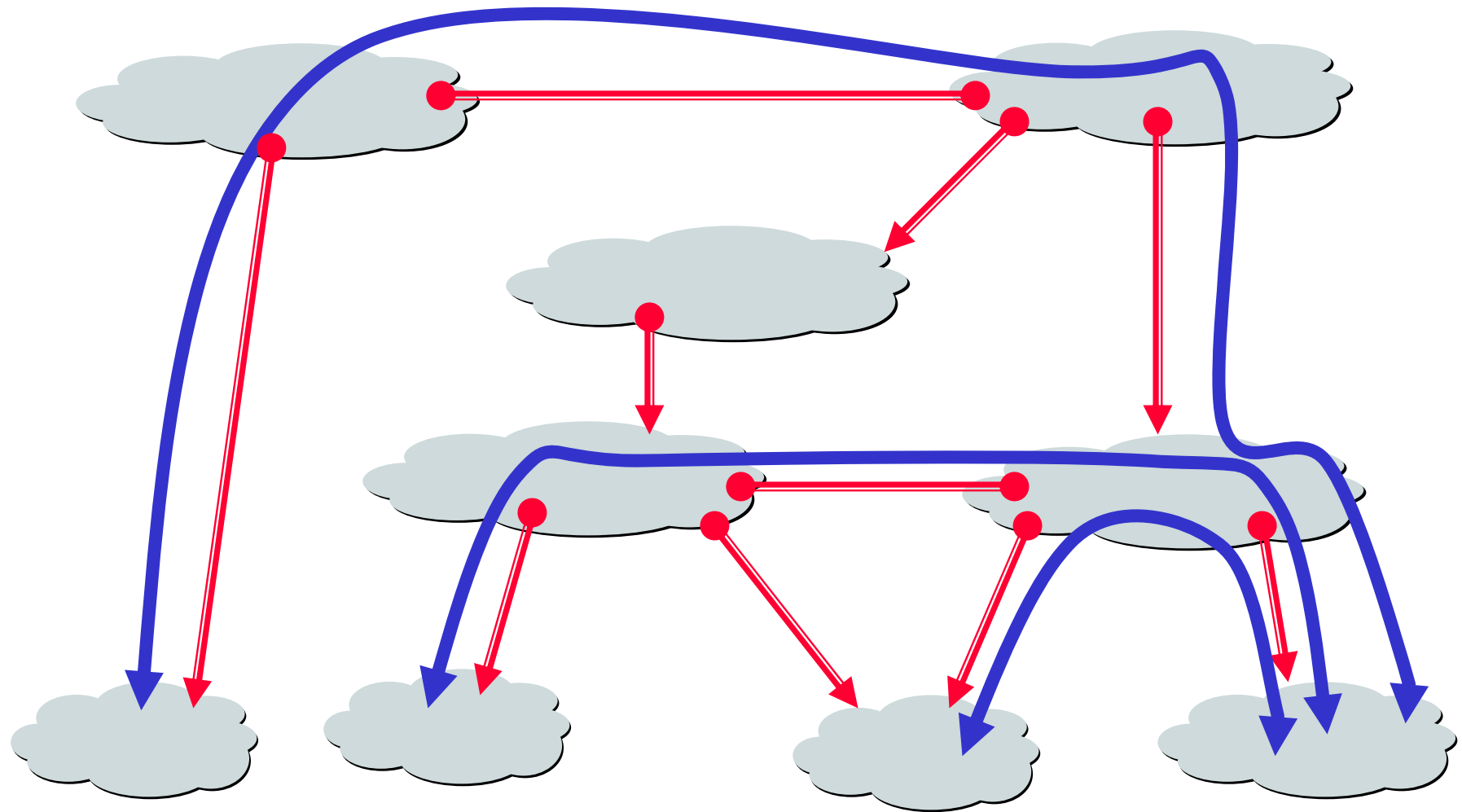❖ Can't use shortest path routing

➢ No universal metric – policy-based decisions

➢ Main characteristic of shortest path does not hold (i→x→j is shortest route, then x→ j is shortest route)

❖ Problems with distance-vector:

➢ Bellman-Ford algorithm may not converge, and may loop

❖ Problems with link state:

➢ Metric used by different routers are not the same → loops

➢ LS database too large – entire Internet

➢ May expose policies to other AS's

# BGP: Distance Vector with Path

❖ Each routing update carries the entire path
  ➢ e.g.,: destination 18.26/16 is reachable using {AS1, AS3, AS11}
❖ When AS receives a routing update
  ➢ Reject routes with loops
    • To detect loops check whether my AS is already in path
❖ AS remembers loop-free routes
❖ For each destination, the AS chooses the best route according its policies.
❖ AS advertises a neighbor routes to a subset of all the destinations, depending on its policy
  ➢ E.g., I might hide from you that I know how to get to destination X, because I don't want to deliver your messages to X
❖ AS advertises to neighbors only those routes that it uses
  ➢ Ensures that if i→x→j is the used route, then x→j is the used route
  ➢ *What happens if an AS advertises routes that it doesn't use?*
❖ Advantage:
  ➢ Metrics are local - AS chooses path, protocol ensures no loops

# Implementing Customer/Provider and Peer/Peer relationships using BGP

❖ BGP provides capability for enforcing various policies

❖ Policies are **<u>not</u>** part of BGP: they are provided to BGP as configuration information

❖ BGP enforces policies by

1. choosing paths from multiple alternatives (importing routers)

2. controlling advertisement to other AS's (exporting routes)

# Importing Routes

❖ Based on route attributes
  ➢ First, Prefer customer > peer > provider
  ➢ Then, Shortest AS PATH length
  ➢ Then, look at other route attributes

# Exporting Routes

- ❖ When an AS exports a route, others can use the AS to forward packets along that route
- ❖ Rules:
  - ➢ Export customers routes to everyone
    - why?
  - ➢ Export routes to your own addresses to everyone
    - Why?
  - ➢ Don't export routes advertised to you by your provider (may advertise them to customers)
    - Why?
  - ➢ Don't export routes advertised to you by your peer (may advertise them to customers)
    - Why?

# Import Routes

provider route    peer route    customer route    ISP route

From provider

From provider

From peer

From peer

(Prefer customer > peer > provider)

From customer

From customer

# Export Routes



**Legend:**
- ◆ provider route
- ✚ peer route
- ♥ customer route
- ☺ ISP route

To provider
From provider
To peer
To peer
To customer
To customer

filters block
◆ ✚

# AS Graphs Depend on Point of View

# 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)**

10

# BGP Operations (Simplified)

**Establish session on TCP port 179**

↓

**Exchange all active routes**

↓

**Exchange incremental updates**

AS1

**BGP session**

AS2

While connection is ALIVE exchange route UPDATE messages

11

# Four Types of BGP Messages

- **Open : Establish a peering session.**

- **Keep Alive : Handshake at regular intervals.**

- **Notification : Shuts down a peering session.**

- **Update : <u>Announcing</u> new routes or <u>withdrawing</u> previously announced routes.**

announcement
=
prefix + <u>attributes values</u>

# BGP 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
```
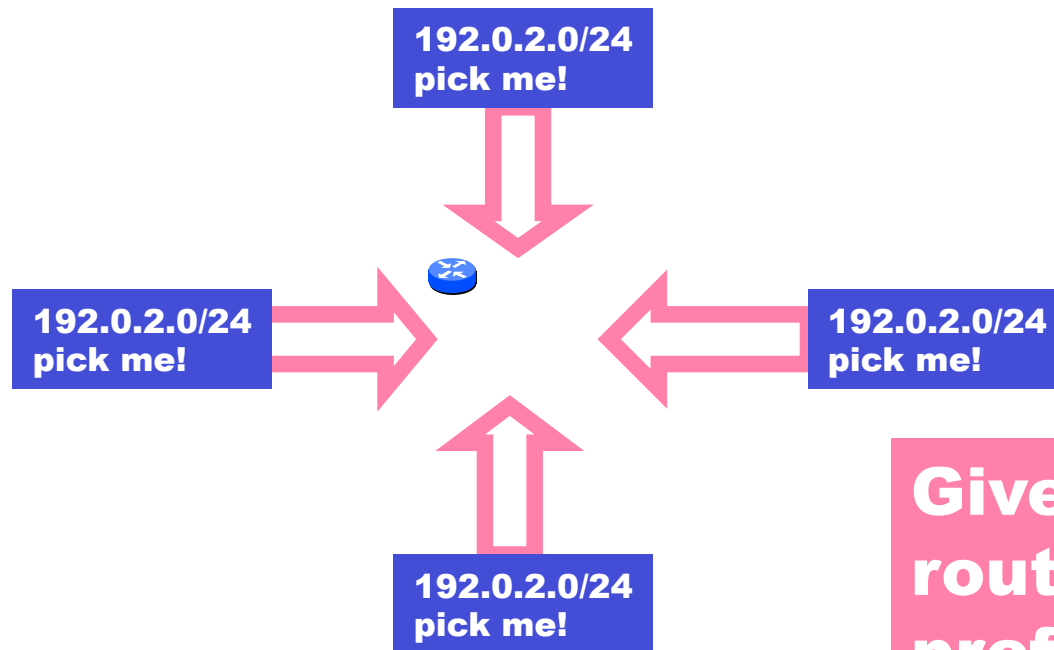
**Most important attributes**

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

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

# Attributes are Used to Select Best Routes

192.0.2.0/24 pick me!

192.0.2.0/24 pick me!

192.0.2.0/24 pick me!

192.0.2.0/24 pick me!

Given multiple routes to the same prefix, a BGP speaker must pick at most <u>one</u> best route

(Note: it could reject them all!)

# 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**

15

# Route Selection Summary

**Highest Local Preference**          Enforce relationships

**Shortest ASPATH**

**Lowest MED**

**i-BGP < e-BGP**                     traffic engineering

**Lowest IGP cost
to BGP egress**

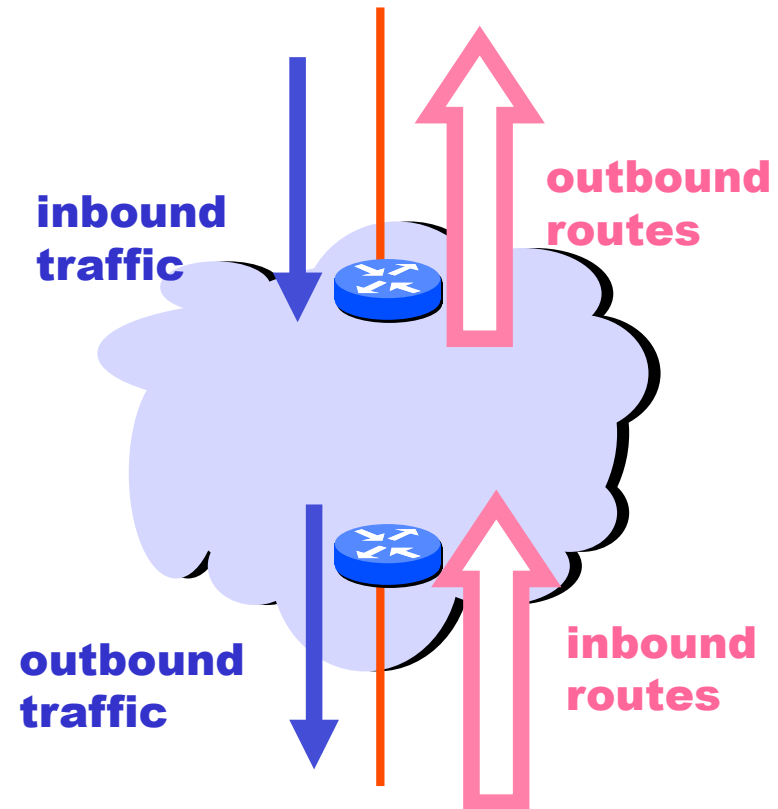**Lowest router ID**                  Throw up hands and
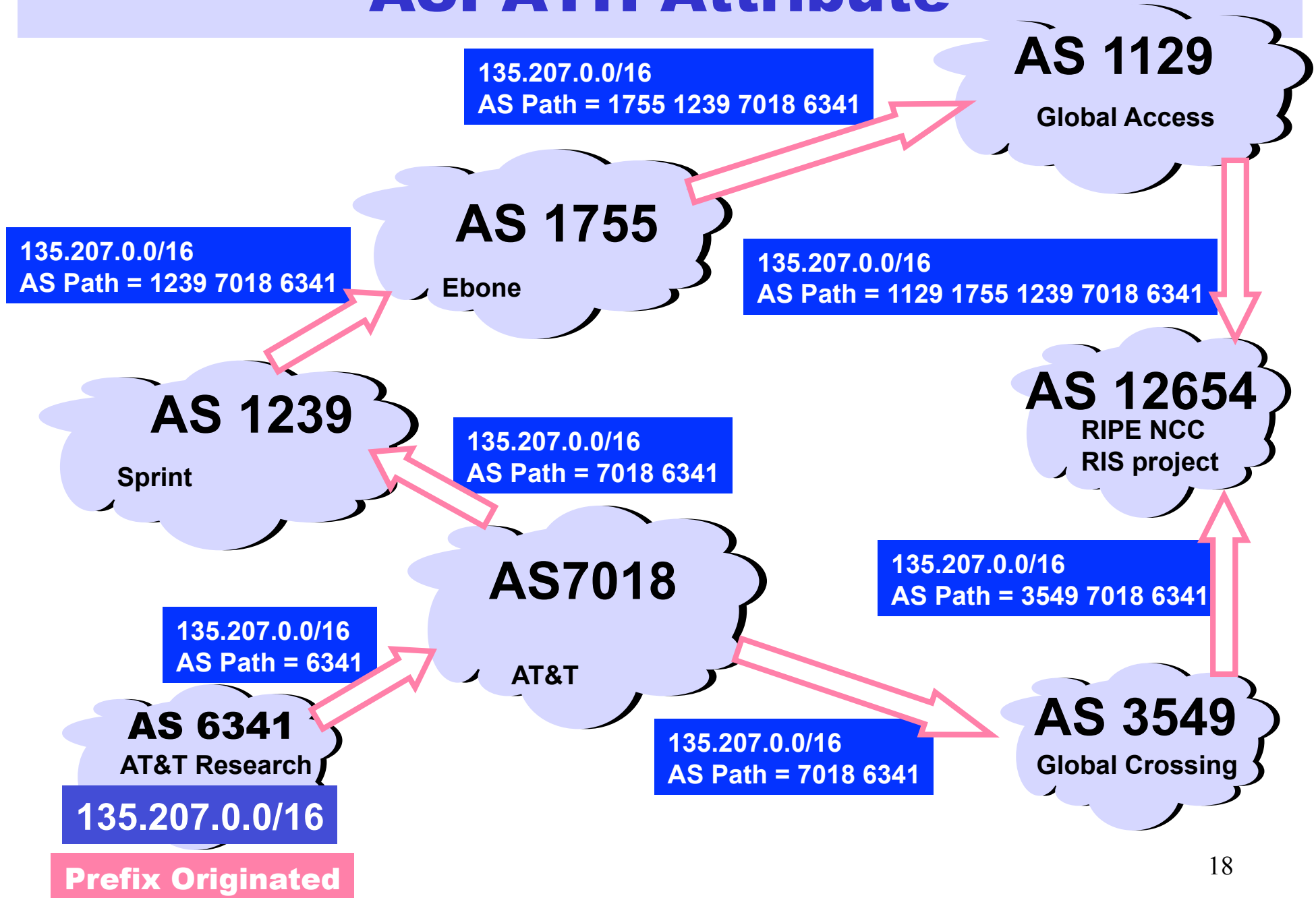                                      break ties

# Tweak Tweak Tweak

- **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
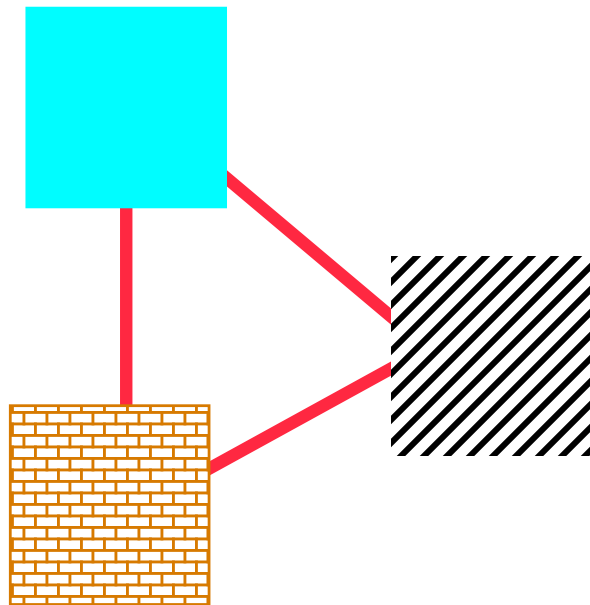  - Tweak attributes on <u>inbound</u> routes to influence best route selection

In general, an AS has more control over outbound traffic

inbound traffic

outbound routes

outbound traffic

inbound routes

# ASPATH Attribute

**AS 1129**
Global Access

135.207.0.0/16
AS Path = 1755 1239 7018 6341

**AS 1755**
Ebone

135.207.0.0/16
AS Path = 1239 7018 6341

135.207.0.0/16
AS Path = 1129 1755 1239 7018 6341

**AS 1239**
Sprint

**AS 12654**
RIPE NCC
RIS project

135.207.0.0/16
AS Path = 7018 6341

**AS7018**
AT&T

135.207.0.0/16
AS Path = 6341

135.207.0.0/16
AS Path = 3549 7018 6341

**AS 6341**
AT&T Research

**135.207.0.0/16**

**Prefix Originated**

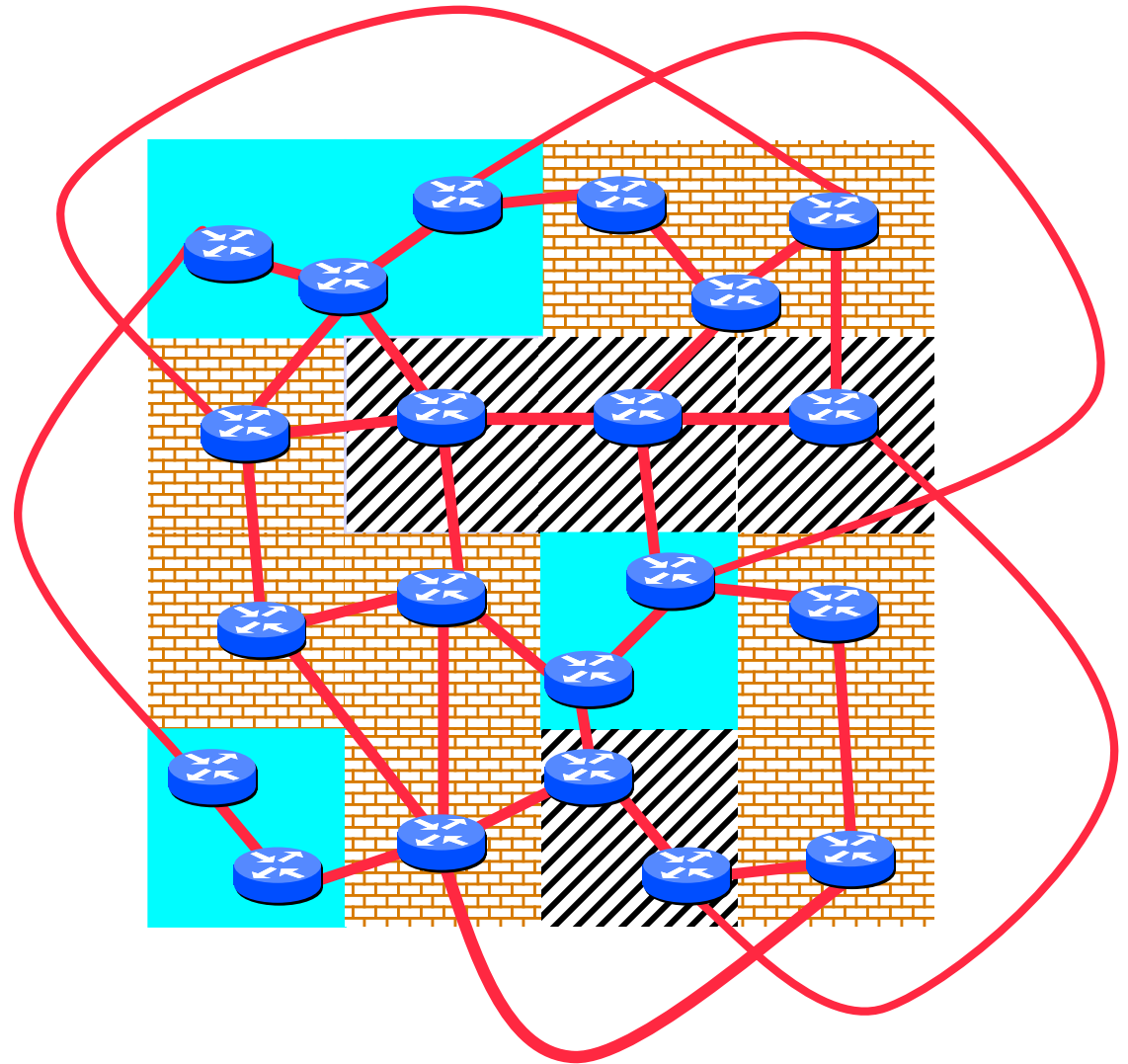135.207.0.0/16
AS Path = 7018 6341

**AS 3549**
Global Crossing

18

# 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...
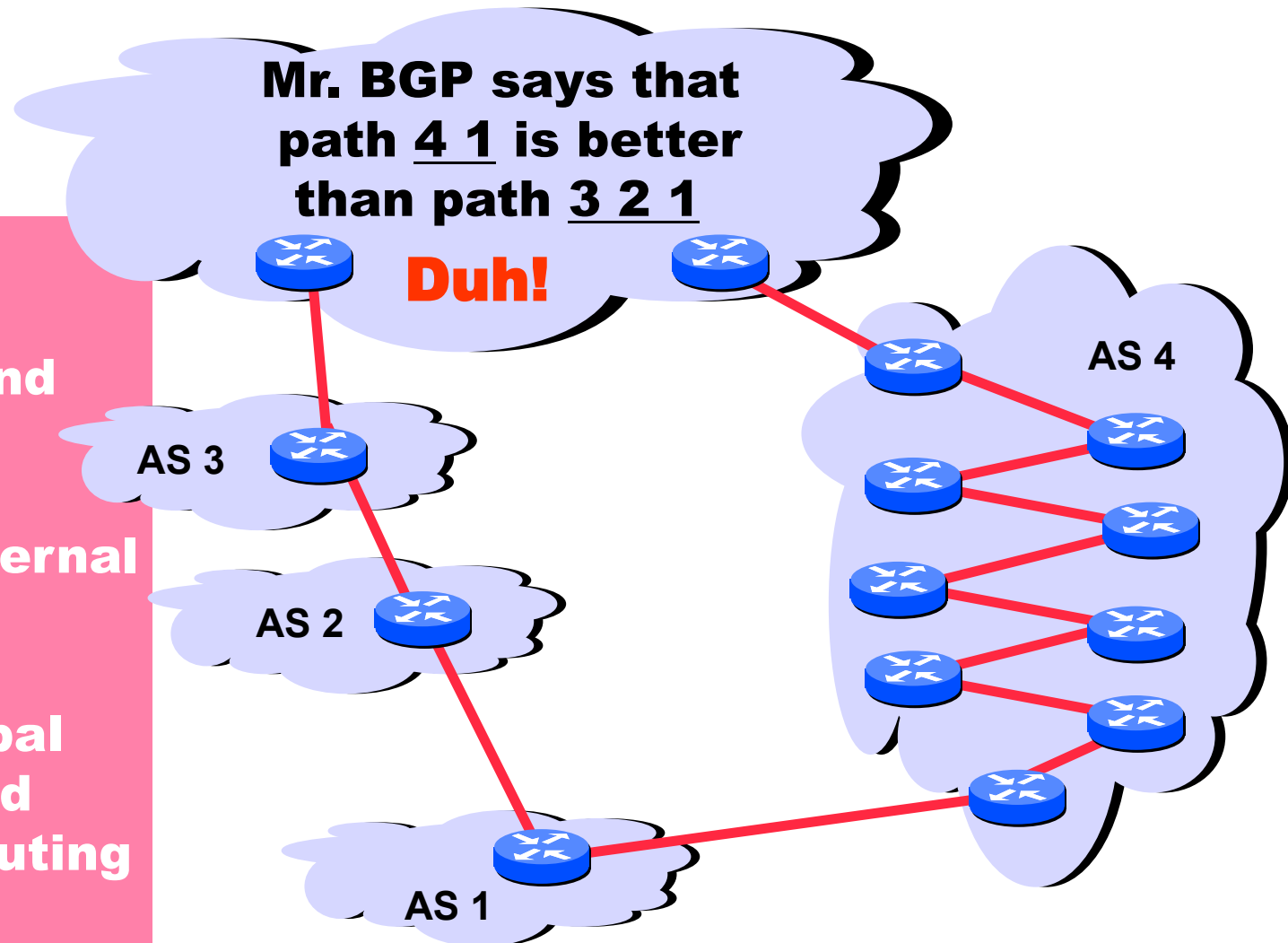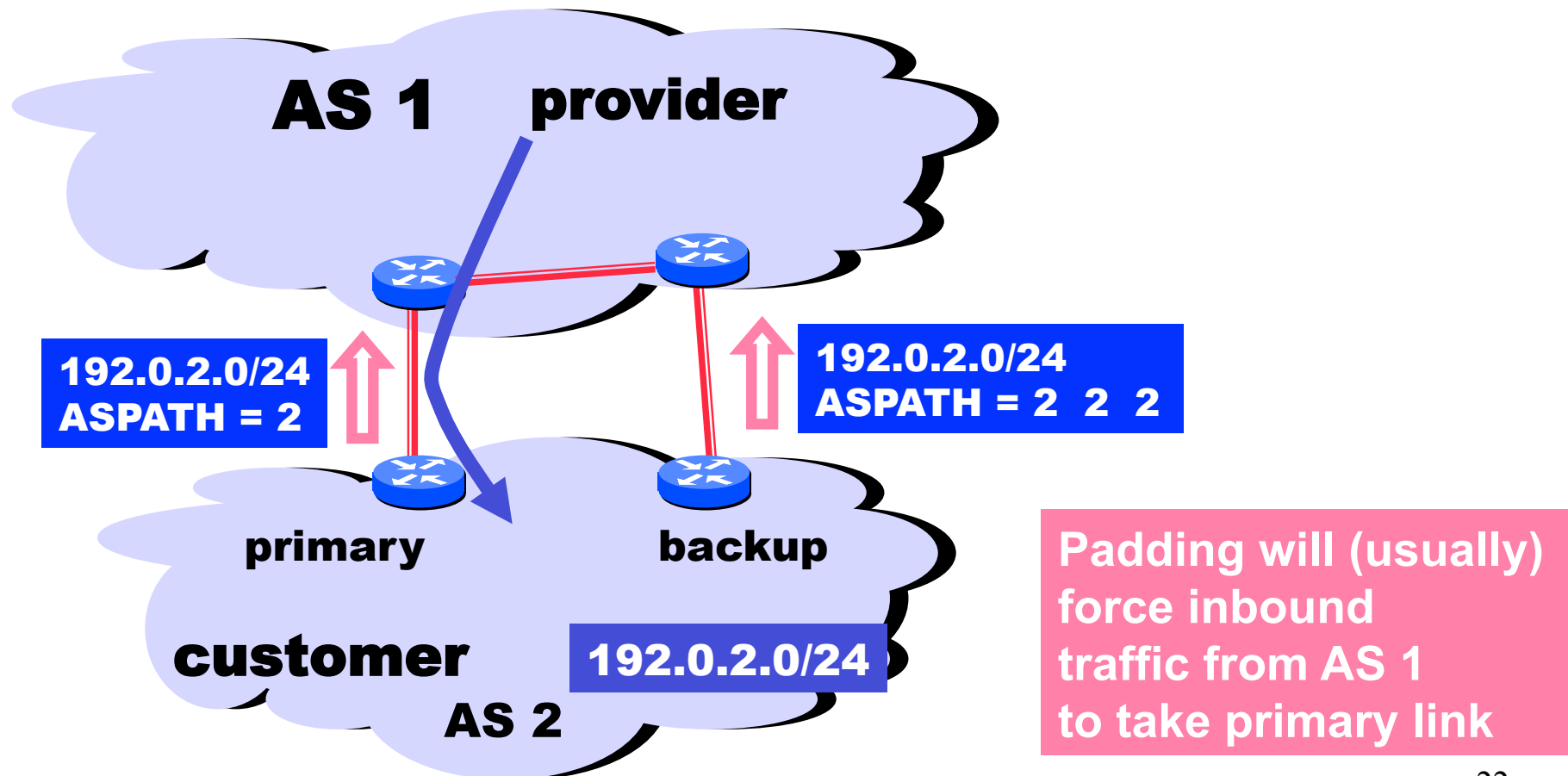
# Shorter Doesn't Always Mean Shorter

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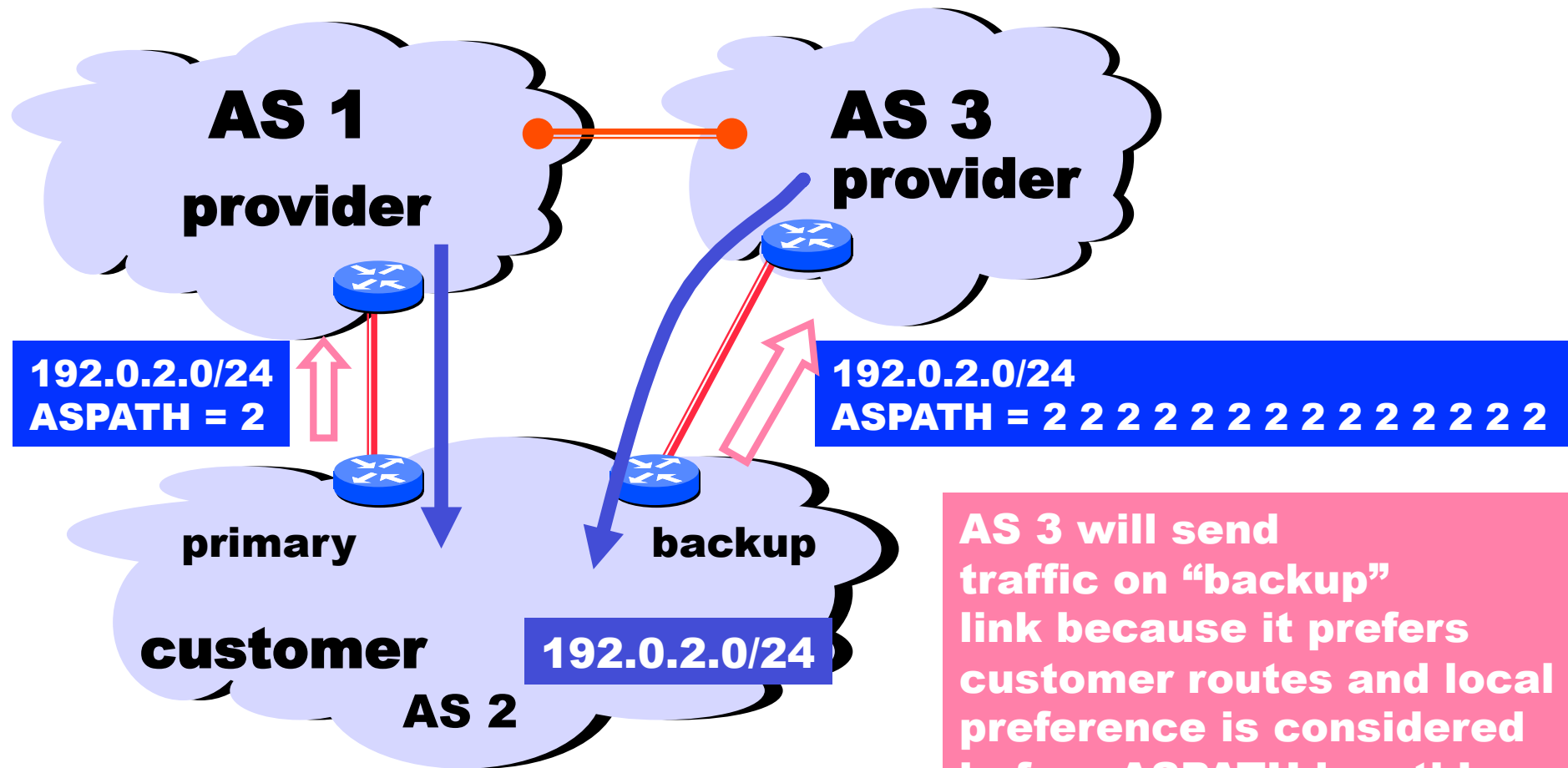
**Exporting internal state would dramatically increase global instability and amount of routing state**

AS 4

AS 3

AS 2

AS 1

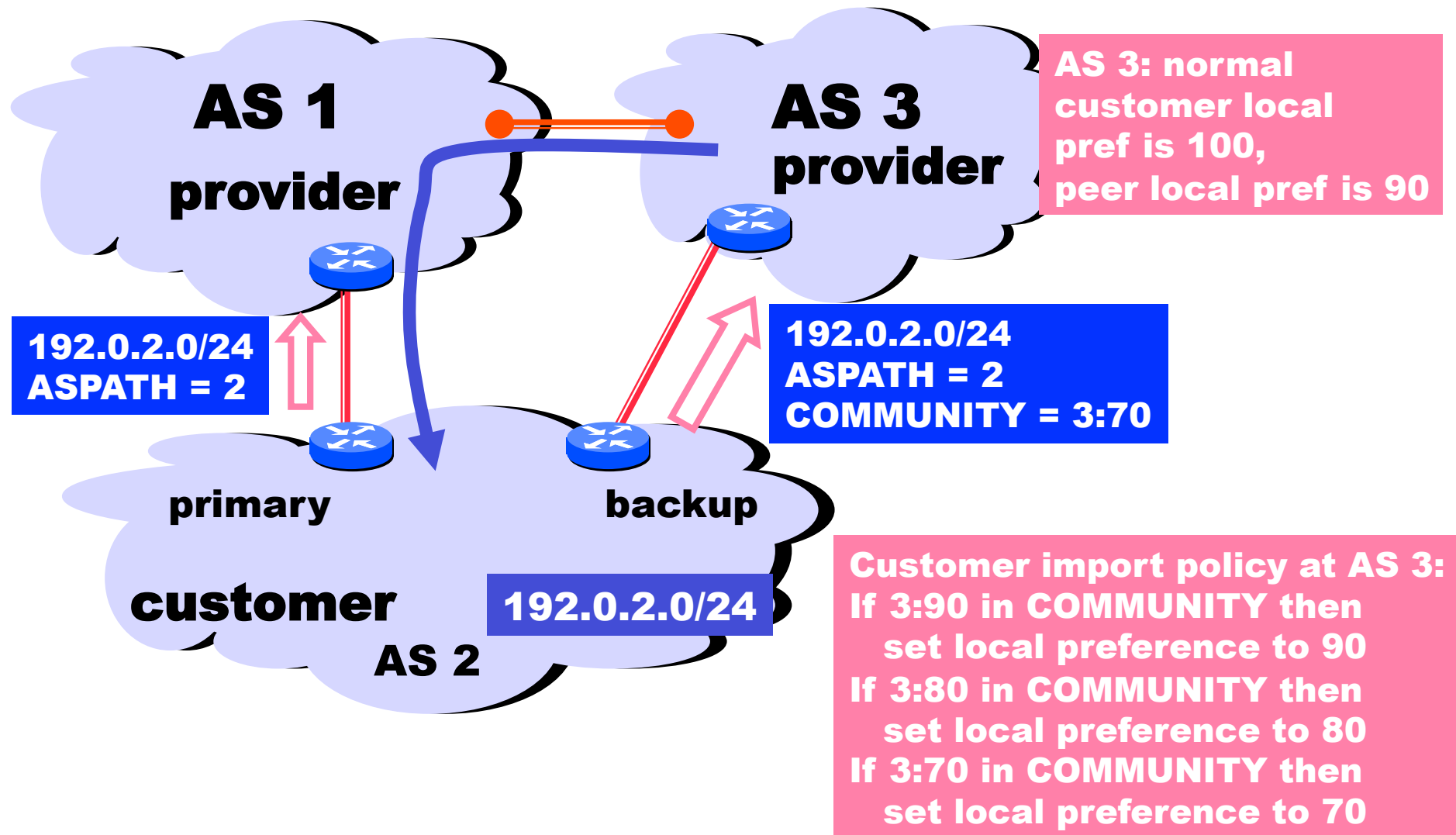# Shedding Inbound Traffic with ASPATH Padding Hack



AS 1    provider

192.0.2.0/24
ASPATH = 2

192.0.2.0/24
ASPATH = 2  2  2

primary    backup

customer

192.0.2.0/24

AS 2

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

# Padding May Not Shut Off All Traffic

**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
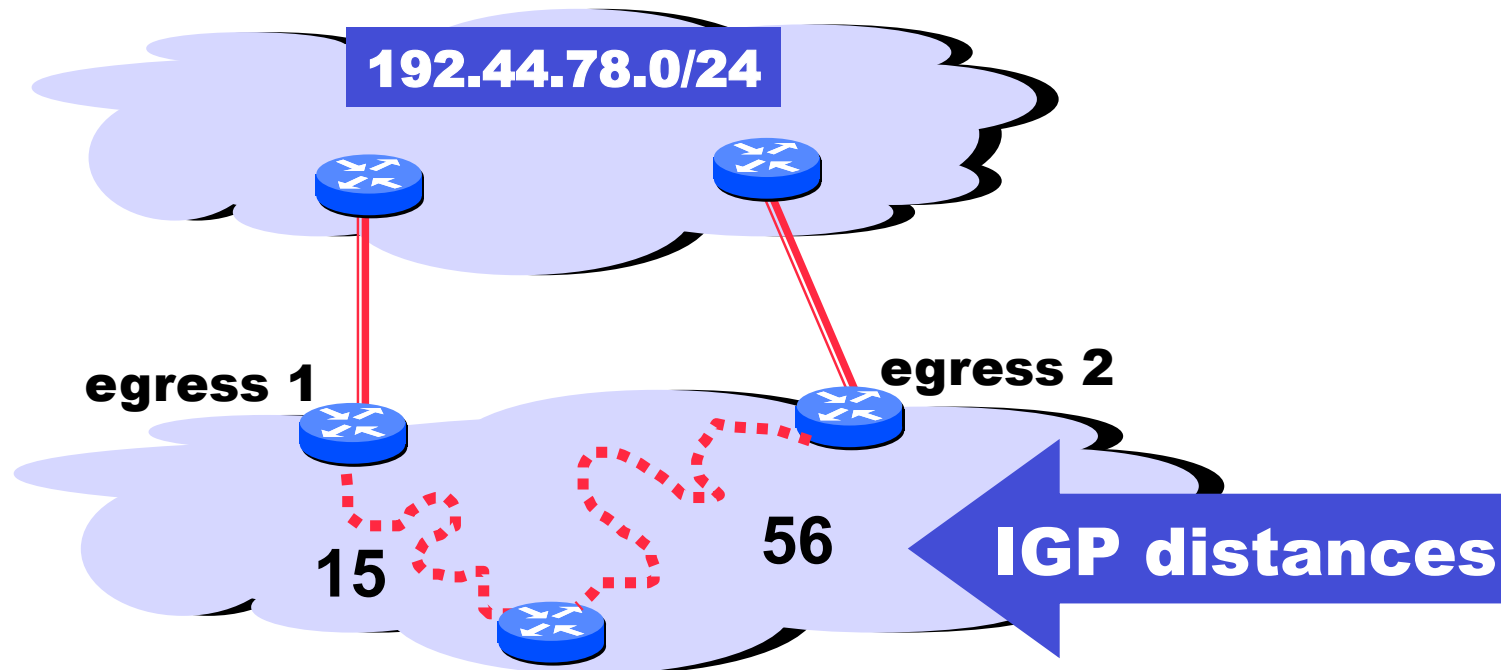
**primary**

**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

192.0.2.0/24

AS 2

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

24

# Hot Potato Routing: Go for the Closest Egress Point
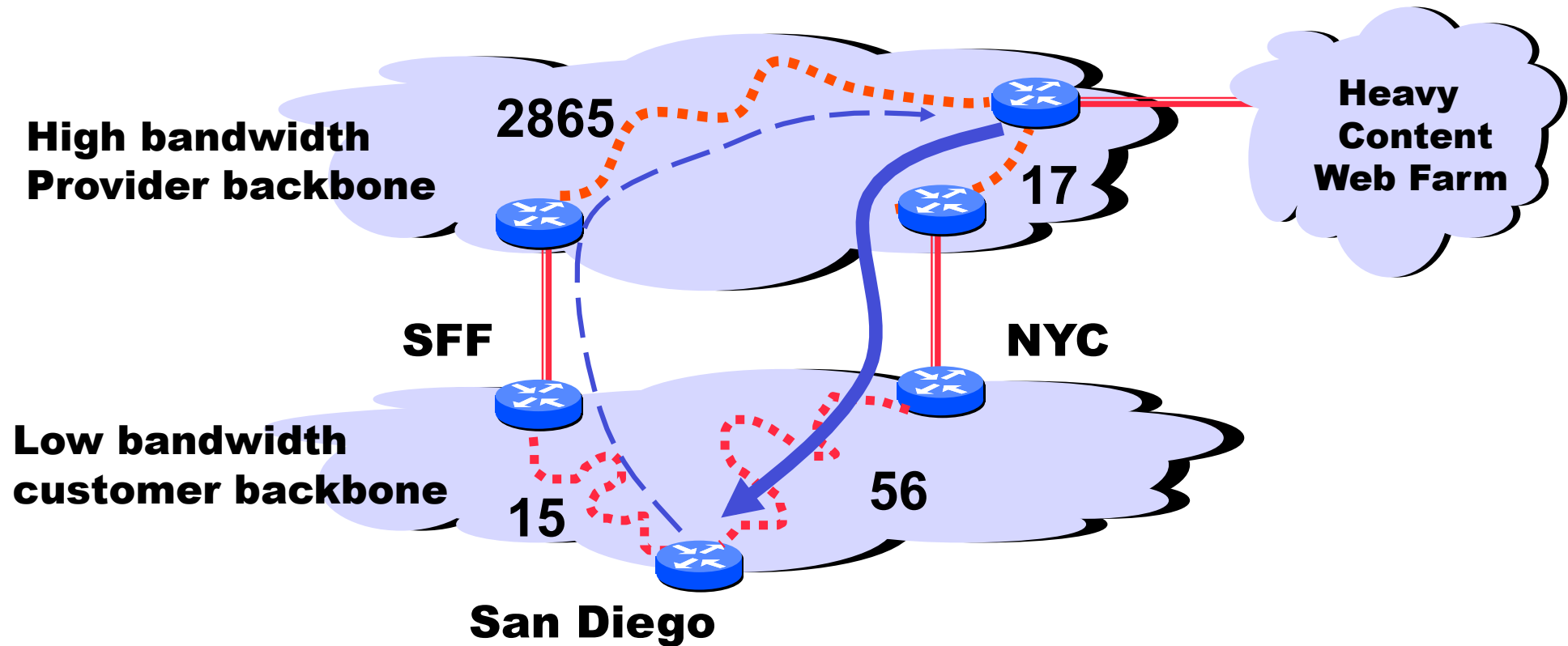
192.44.78.0/24

egress 1

egress 2

15

56

IGP distances

This Router has two BGP routes to 192.44.78.0/24.

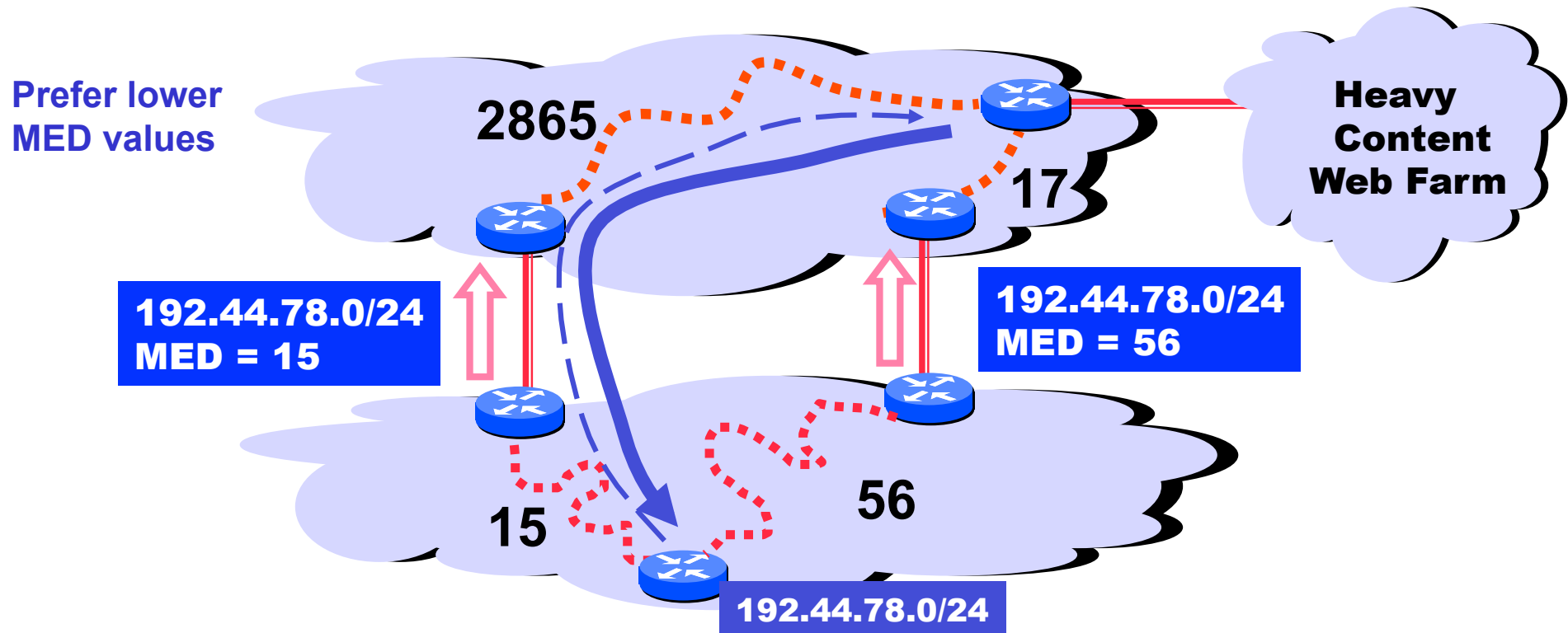Hot potato: get traffic off of your network as Soon as possible.  Go for egress 1!

# Getting Burned by the Hot Potato



High bandwidth
Provider backbone

2865

Heavy
Content
Web Farm

17

SFF

NYC

Low bandwidth
customer backbone

15

56

San Diego

**Many customers want their provider to carry the bits!**

- - - → tiny http request
──→ huge http reply

26

# Cold Potato Routing with MEDs (Multi-Exit Discriminator Attribute)



Prefer lower MED values

2865

17

Heavy Content Web Farm

192.44.78.0/24
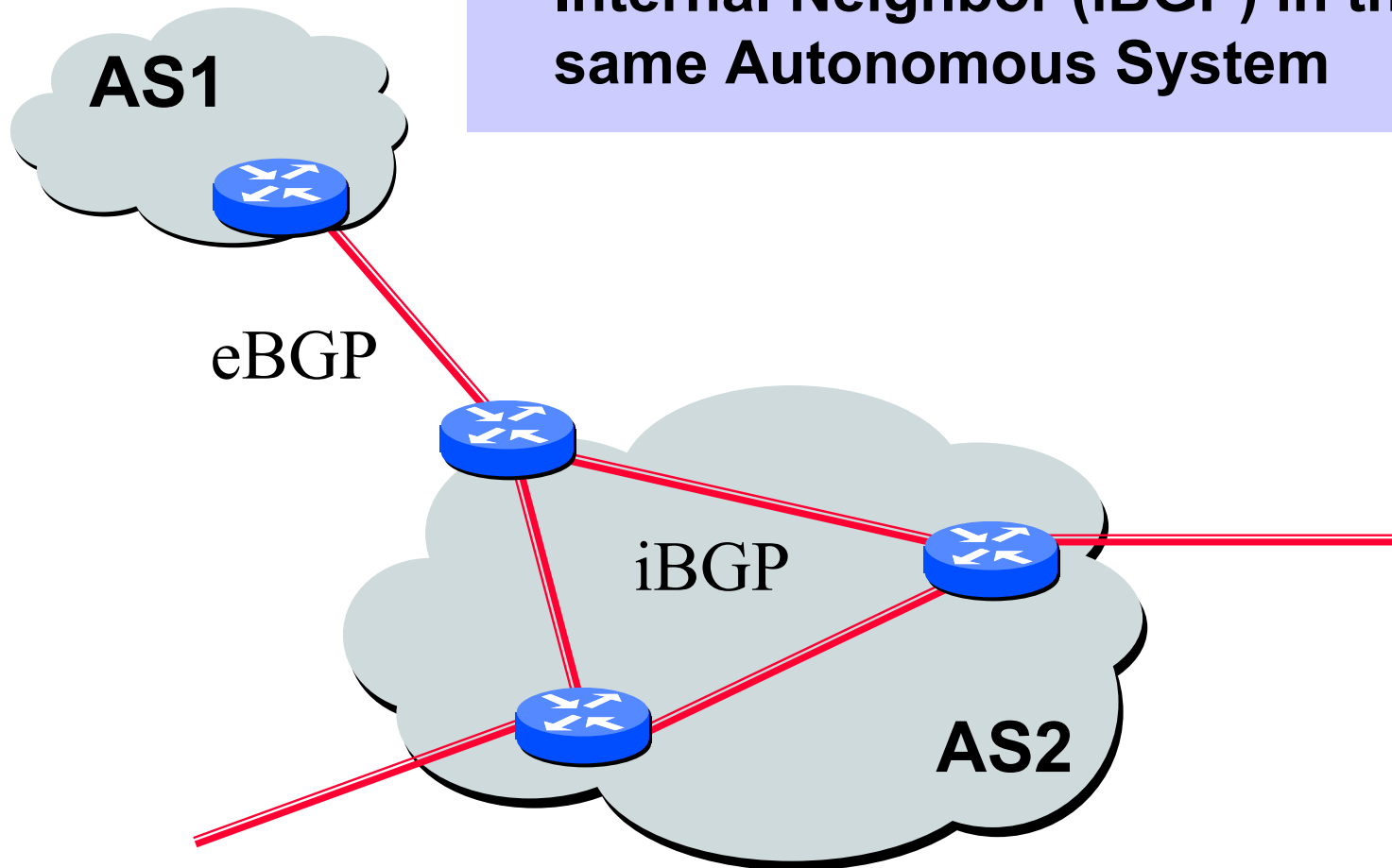MED = 15

192.44.78.0/24
MED = 56

15

56

192.44.78.0/24

**This means that MEDs must be considered BEFORE IGP distance!**

Note1 : some providers will not listen to MEDs
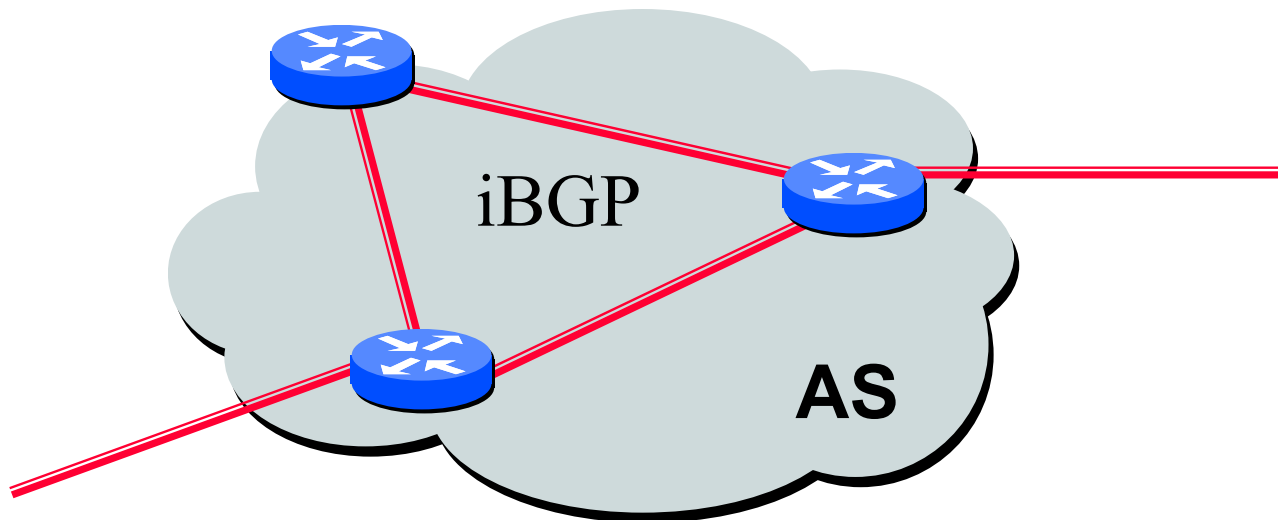
Note2 : MEDs need not be tied to IGP distance

# Two Types of BGP Neighbor Relationships

- **External Neighbor (eBGP) in a different Autonomous Systems**
- **Internal Neighbor (iBGP) in the same Autonomous System**

**AS1**

eBGP

iBGP

**AS2**

# iBGP

❖ AS has more than one router participating in eBGP

❖ iBGP is run between BGP routers in the same AS to allow all of them to obtain a complete and consistent view of external routes

iBGP

AS

# Internal BGP (iBGP)

❖ Same messages as eBGP

❖ Different rules about re-advertising prefixes:

  ➢ Prefix learned from eBGP can be advertised to iBGP neighbor and vice-versa, but

  ➢ Prefix learned from one iBGP neighbor cannot be advertised to another iBGP neighbor

    • Reason: no AS PATH within the same AS and thus danger of looping.

# We learned

❖ Inter-domain routing uses policy

❖ As a result, routing is not a simple optimization of a single number which can be done using shortest path algorithms

❖ BGP is designed to route based on policies

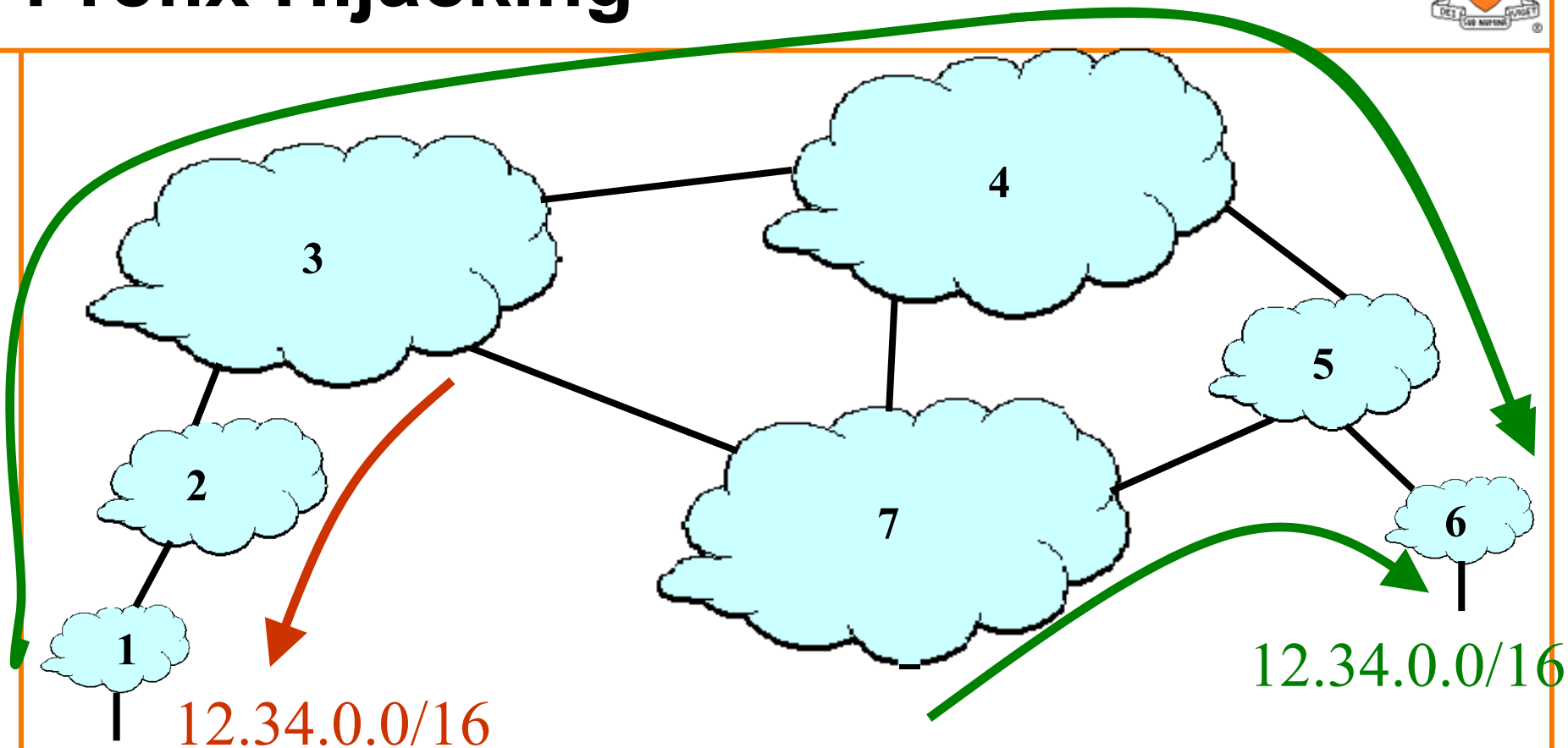# Validity of the routing information: Origin authentication

# IP Address Ownership and Hijacking

- IP address block assignment
  - Regional Internet Registries (ARIN, RIPE, APNIC)
  - Internet Service Providers

- Proper origination of a prefix into BGP
  - By the AS who owns the prefix
  - … or, by its upstream provider(s) in its behalf

- However, what's to stop someone else?
  - Prefix hijacking: another AS originates the prefix
  - BGP does not verify that the AS is authorized
  - Registries of prefix ownership are inaccurate

12

# Prefix Hijacking



12.34.0.0/16

12.34.0.0/16

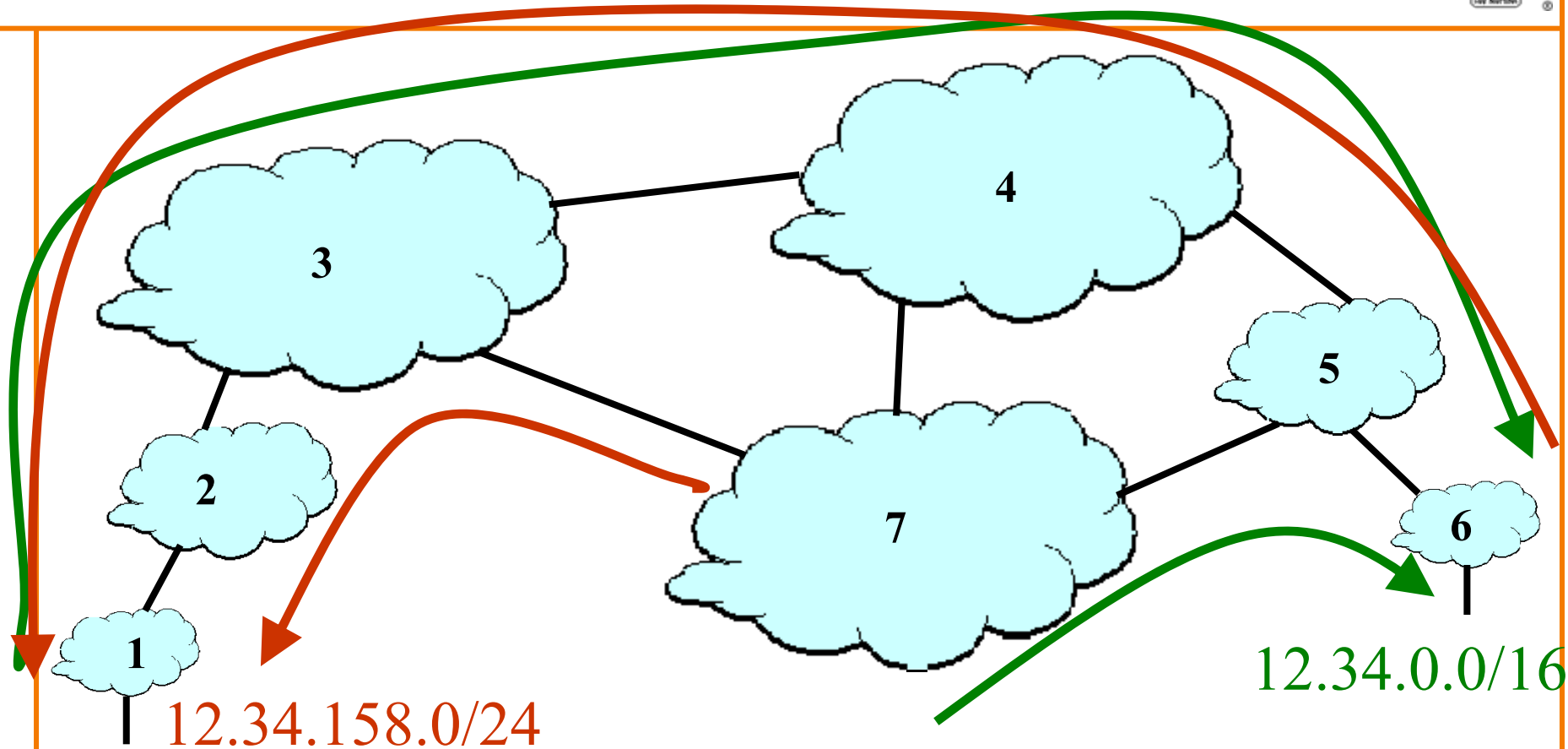- Consequences for the affected ASes
  - Blackhole: data traffic is discarded
  - Snooping: data traffic is inspected, and then redirected
  - Impersonation: data traffic is sent to bogus destinations

13

# Hijacking is Hard to Debug

- Real origin AS doesn't see the problem
  - Picks its own route
  - Might not even learn the bogus route

- May not cause loss of connectivity
  - E.g., if the bogus AS snoops and redirects
  - … may only cause performance degradation

- Or, loss of connectivity is isolated
  - E.g., only for sources in parts of the Internet

- Diagnosing prefix hijacking
  - Analyzing updates from many vantage points
  - Launching traceroute from many vantage points

14

# Sub-Prefix Hijacking



12.34.158.0/24

12.34.0.0/16

- Originating a more-specific prefix
  - Every AS picks the bogus route for that prefix
  - Traffic follows the longest matching prefix

15

# How to Hijack a Prefix

- The hijacking AS has
  - Router with eBGP session(s)
  - Configured to originate the prefix

- Getting access to the router
  - Network operator makes configuration mistake
  - Disgruntled operator launches an attack
  - Outsider breaks in to the router and reconfigures

- Getting other ASes to believe bogus route
  - Neighbor ASes not filtering the routes
  - … e.g., by allowing only expected prefixes
  - But, specifying filters on *peering* links is hard