

Background: BGP and Inter-domain Routing (It's all about the Money)

Based on slides from P. Gill, D. Choffnes, J. Rexford, and A. Feldman
Revised 2015 & 2018 by N. Carlsson

Control plane vs. Data Plane

2

- Control:
 - Make sure that if there's a path available, data is forwarded over it
 - BGP sets up such paths at the AS-level
- Data:
 - For a destination, send packet to most-preferred next hop
 - Routers forward data along IP paths

Network Layer, Control Plane

3

- Function:
 - ▣ Set up routes between networks
- Key challenges:
 - ▣ Implementing provider policies
 - ▣ Creating stable paths

Data Plane

Application

Transport

Network

Data Link

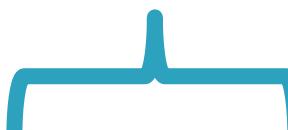
Physical

RIP

OSPF

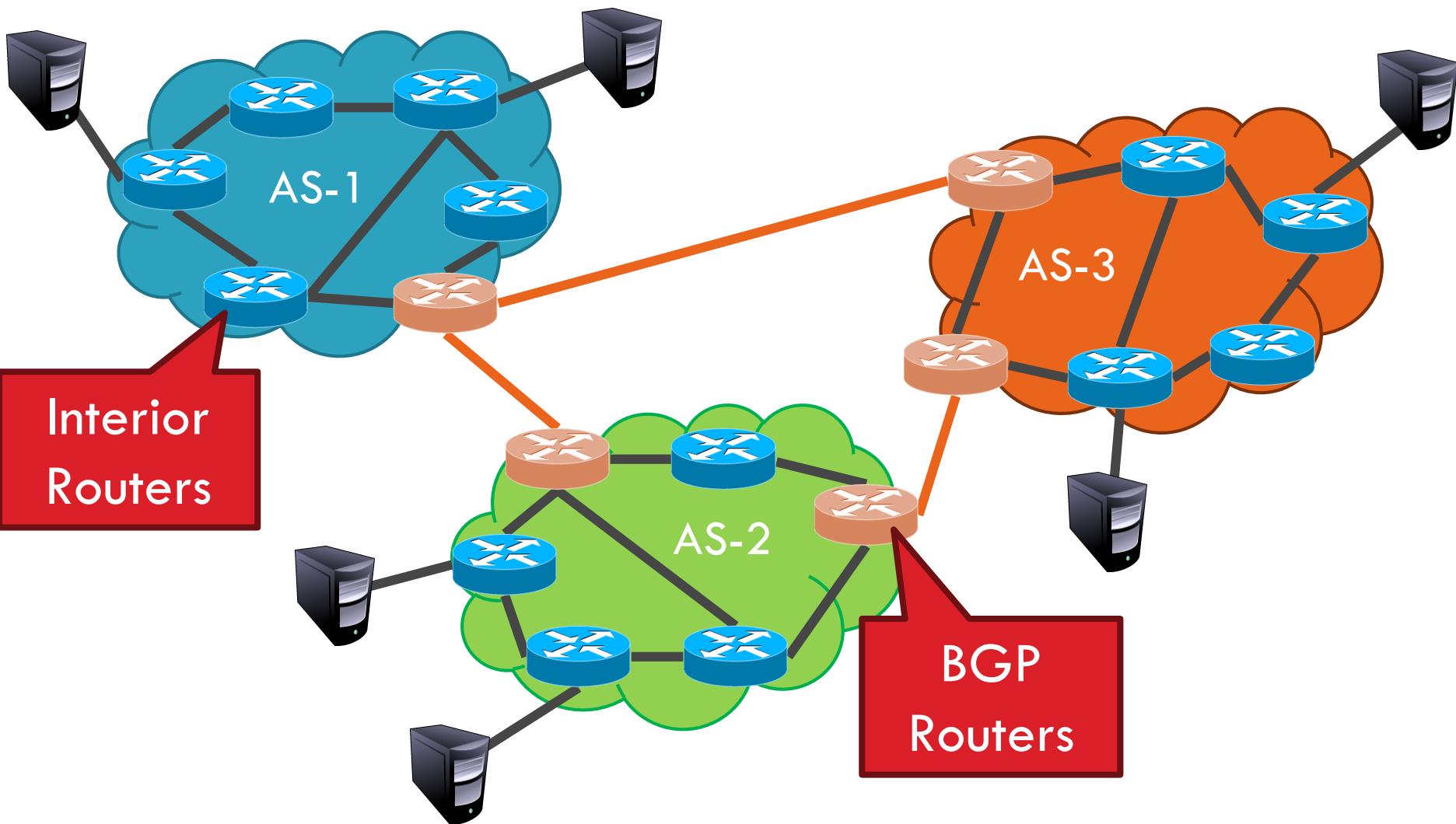
BGP

Control Plane



ASs, Revisited

5



AS Numbers

6

- Each AS identified by an ASN number
 - 16-bit values (latest protocol supports 32-bit ones)
 - 64512 – 65535 are reserved
- Currently, there are ~ 40000 ASNs
 - AT&T: 5074, 6341, 7018, ...
 - Sprint: 1239, 1240, 6211, 6242, ...
 - LIUNET: 2843 (prefix: 130.236.0.0/16)
 - Google 15169, 36561 (formerly YT), + others
 - Facebook 32934
 - North America ASs → <ftp://ftp.arin.net/info/asn.txt>

Inter-Domain Routing

7

- Global connectivity is at stake!
 - Thus, all ASs must use the same protocol
 - Contrast with intra-domain routing
- What are the requirements?
 - Scalability
 - Flexibility in choosing routes
 - Cost
 - Routing around failures
- Question: link state or distance vector?
 - Trick question: BGP is a path vector protocol

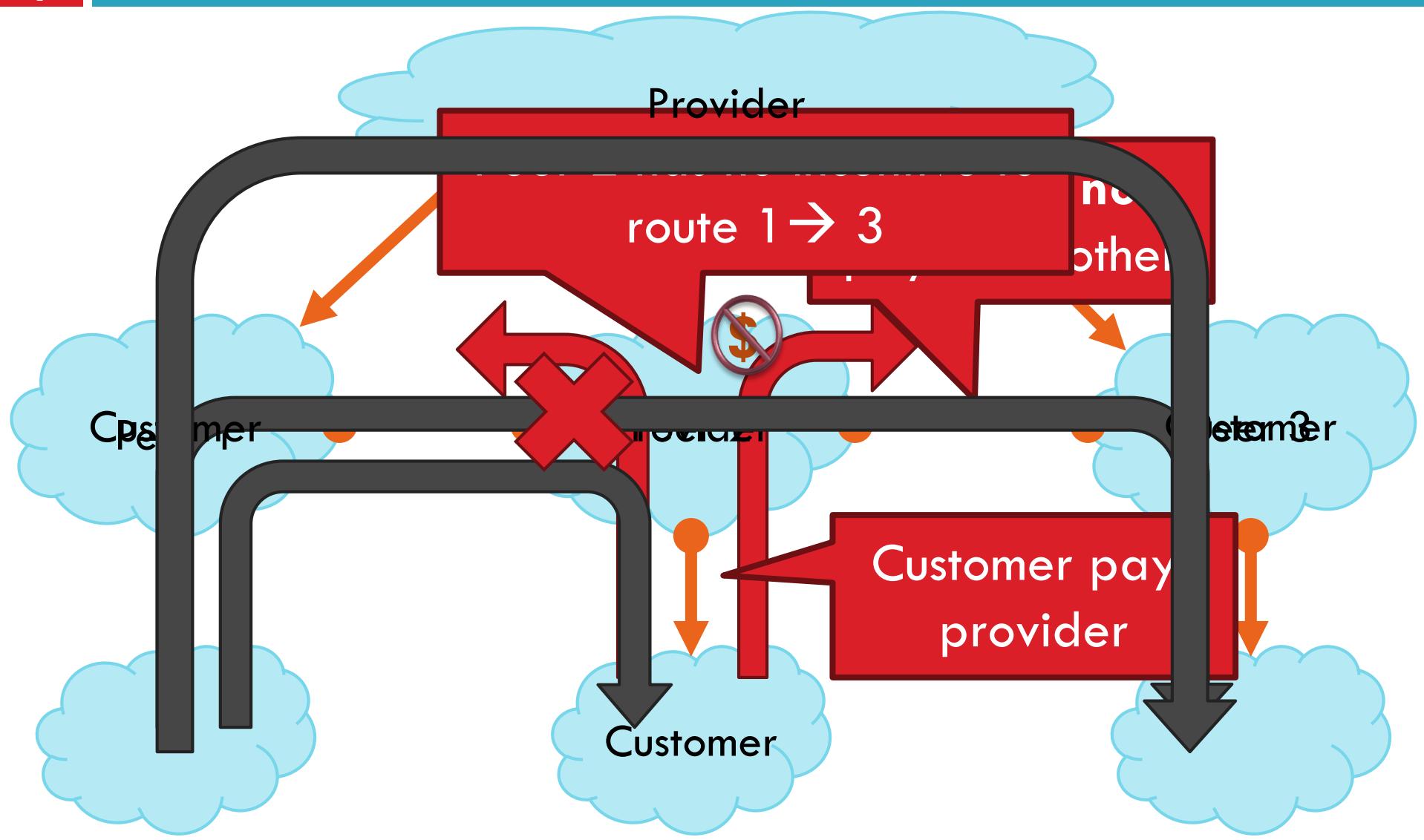
BGP

8

- Border Gateway Protocol
 - ▣ De facto inter-domain protocol of the Internet
 - ▣ Policy based routing protocol
 - ▣ Uses a Bellman-Ford path vector protocol
- Relatively simple protocol, but...
 - ▣ Complex, manual configuration
 - ▣ Entire world sees advertisements
 - Errors can screw up traffic globally
 - ▣ Policies driven by economics
 - How much \$\$\$ does it cost to route along a given path?
 - Not by performance (e.g. shortest paths)

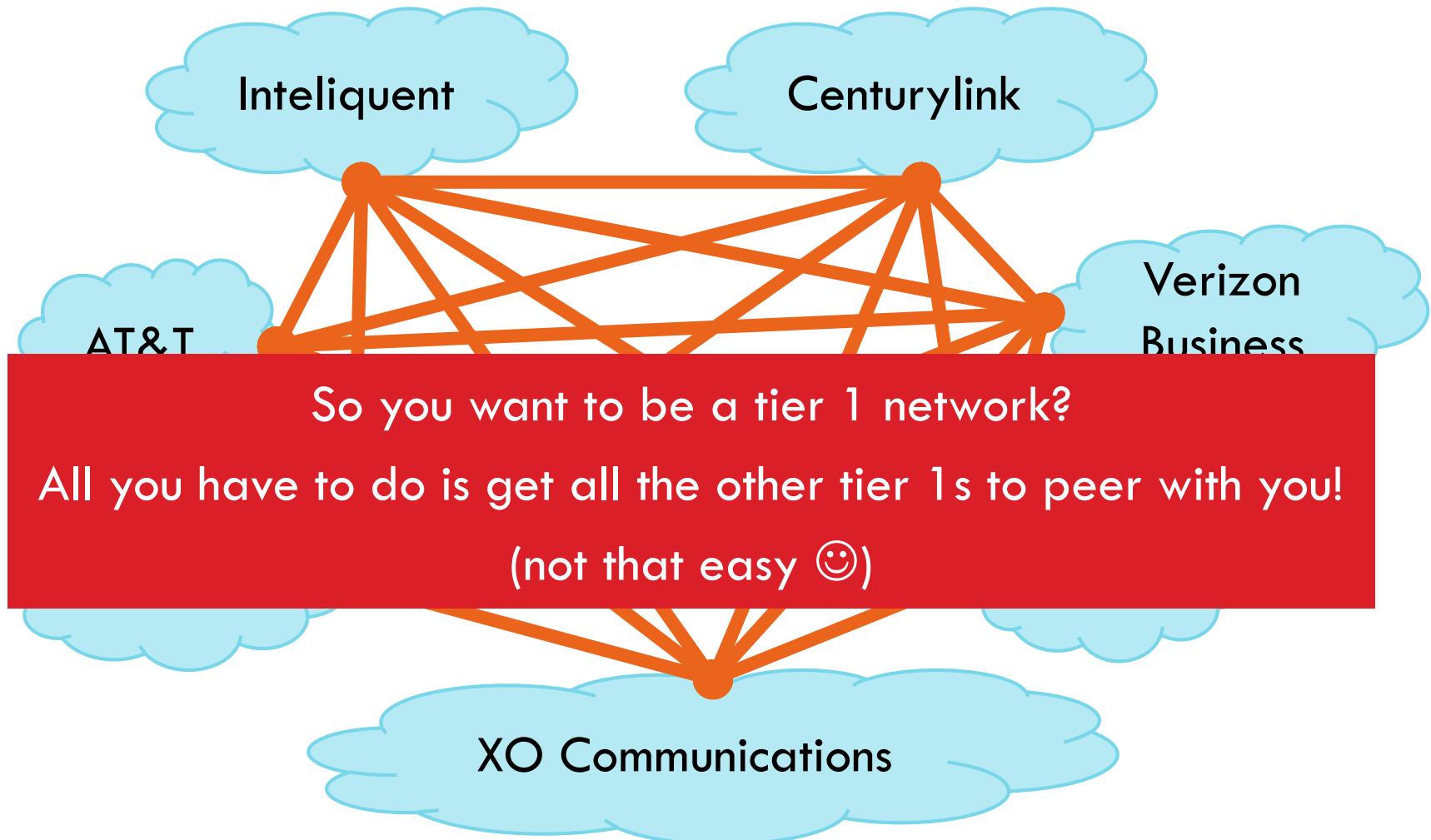
BGP Relationships

9



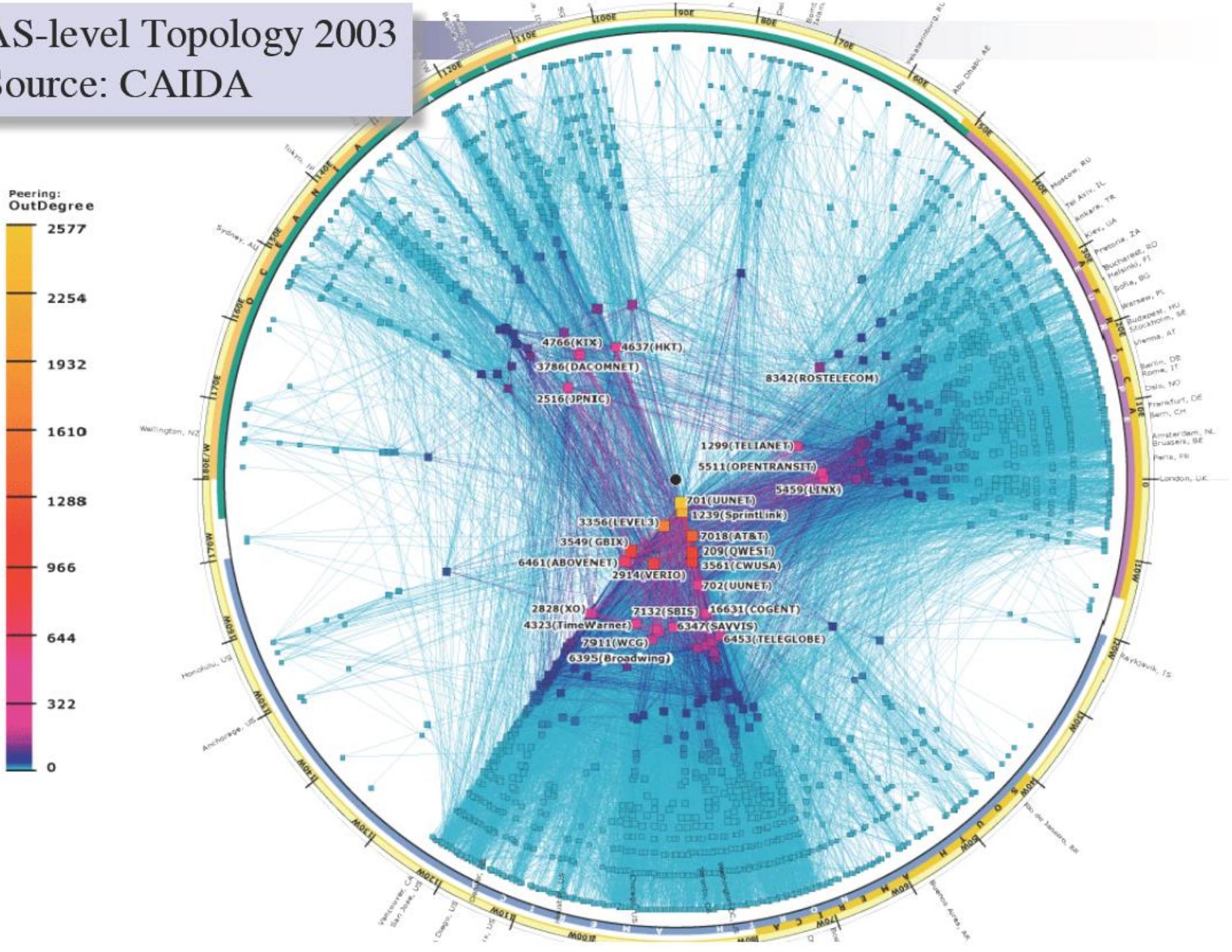
Tier-1 ISP Peering

10



AS-level Topology 2003

Source: CAIDA



Peering Wars

12

Peer

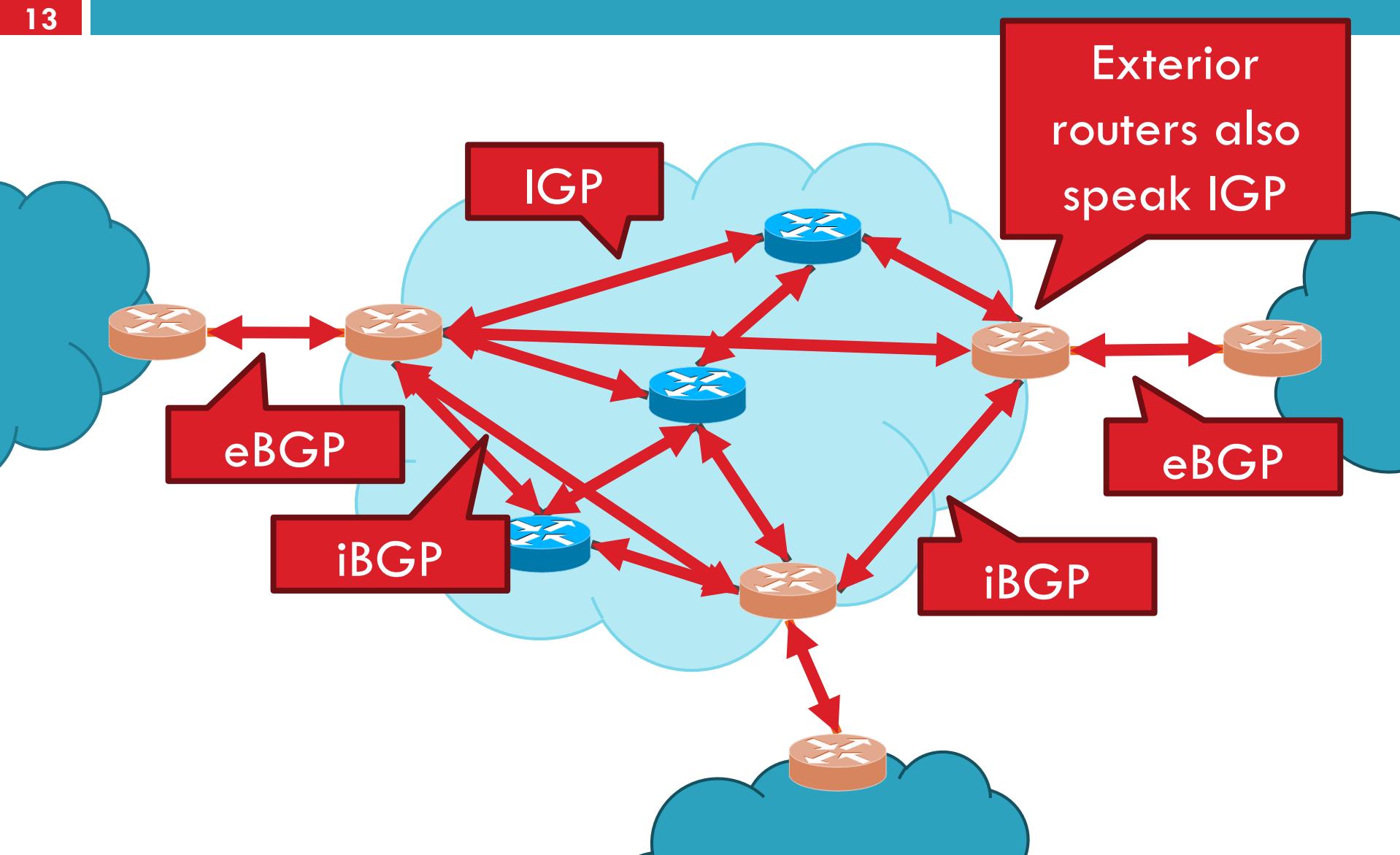
- Reduce upstream costs
- Peering struggles in the ISP world are extremely contentious
 - agreements are usually confidential
- Example: If you are a customer of my peer why should I peer with you? You should pay me too!
 - Incentive to keep relationships private!

Don't Peer

- You would rather have

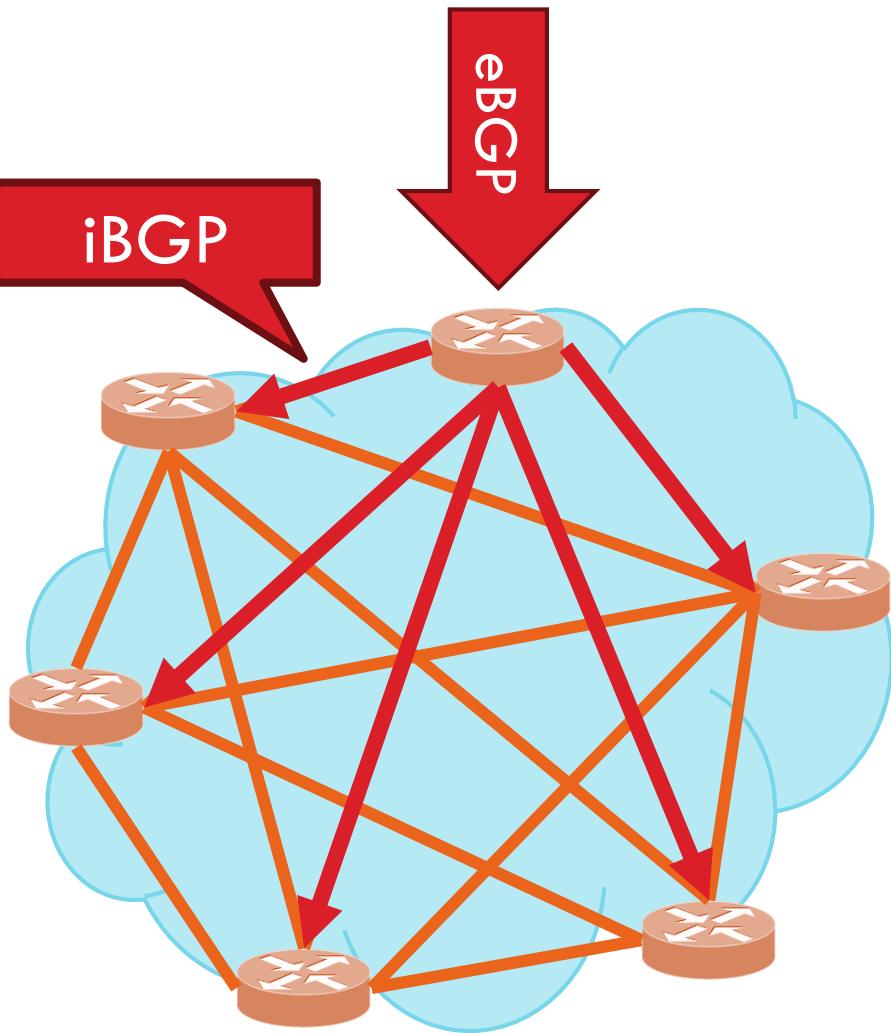
Two Types of BGP Neighbors

13



Full iBGP Meshes

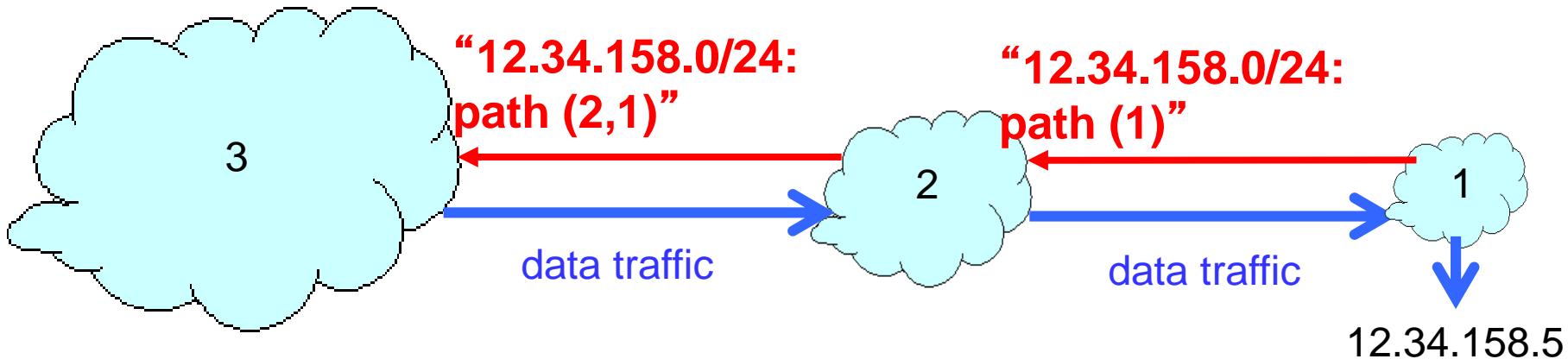
14



- Question: why do we need iBGP?
 - OSPF does not include BGP policy info
 - Prevents routing loops within the AS
- iBGP updates do not trigger announcements

Border Gateway Protocol

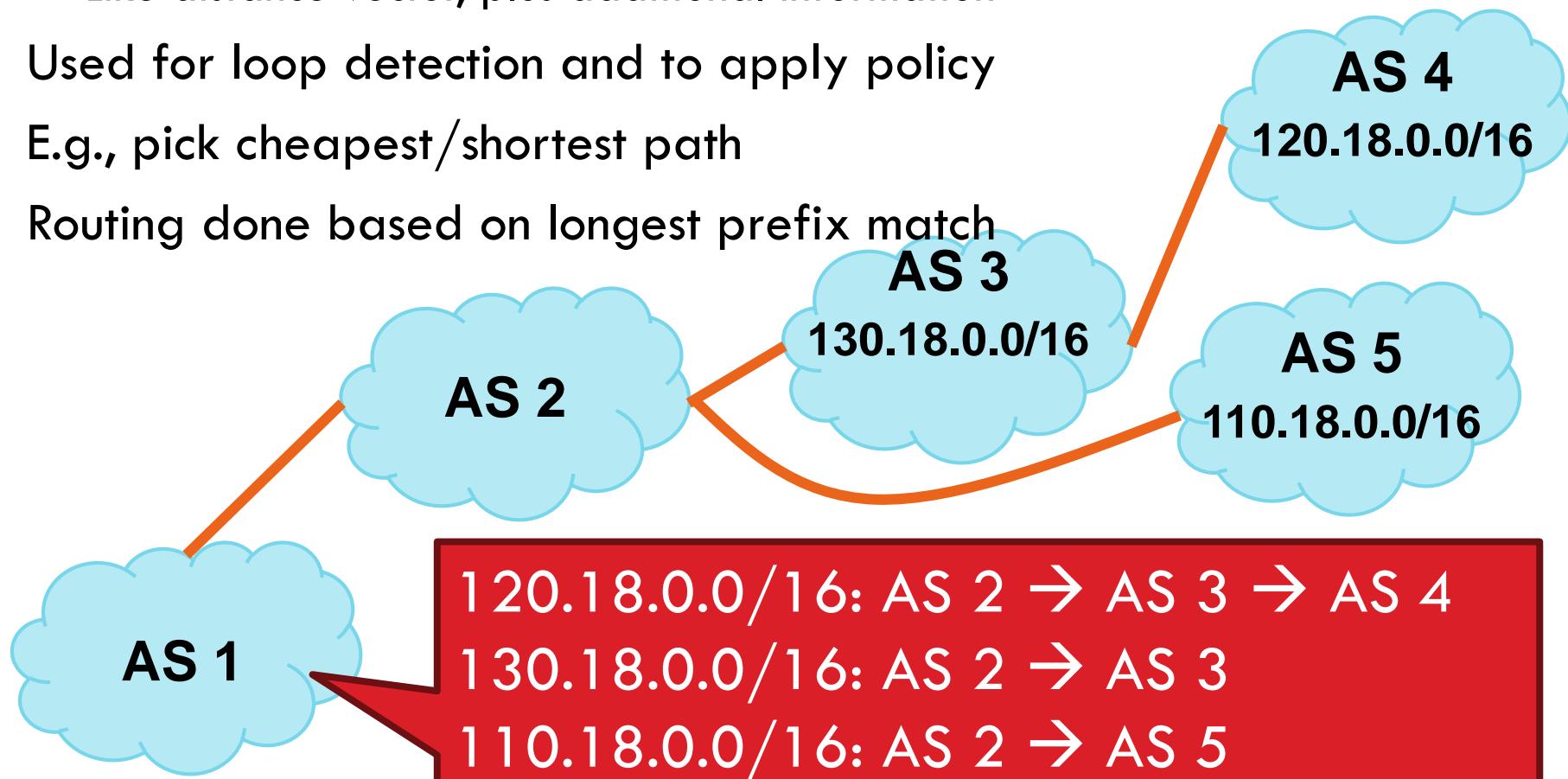
- ASes exchange info about who they can reach
 - IP prefix: block of destination IP addresses
 - AS path: sequence of ASes along the path
- Policies configured by the AS's operator
 - Path selection: which of the paths to use?
 - Path export: which neighbors to tell?



Path Vector Protocol

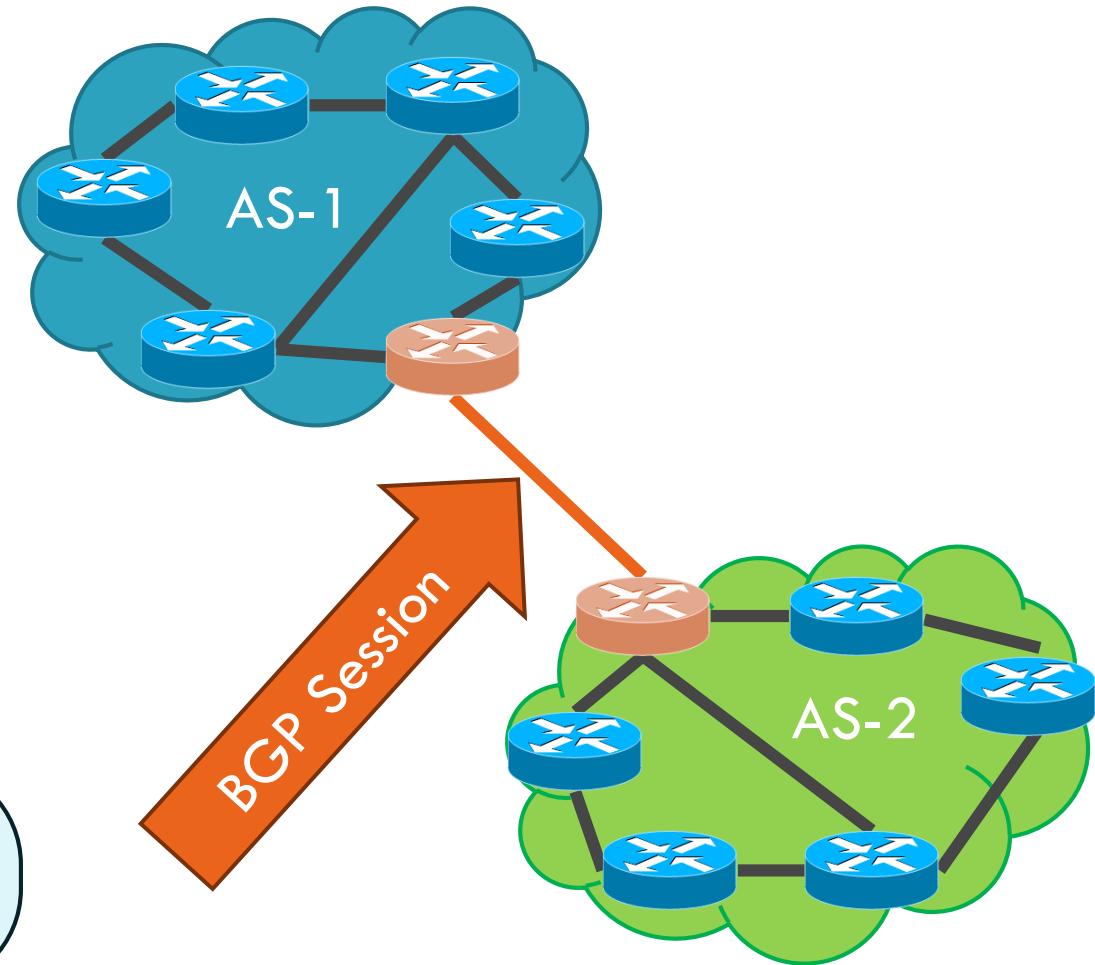
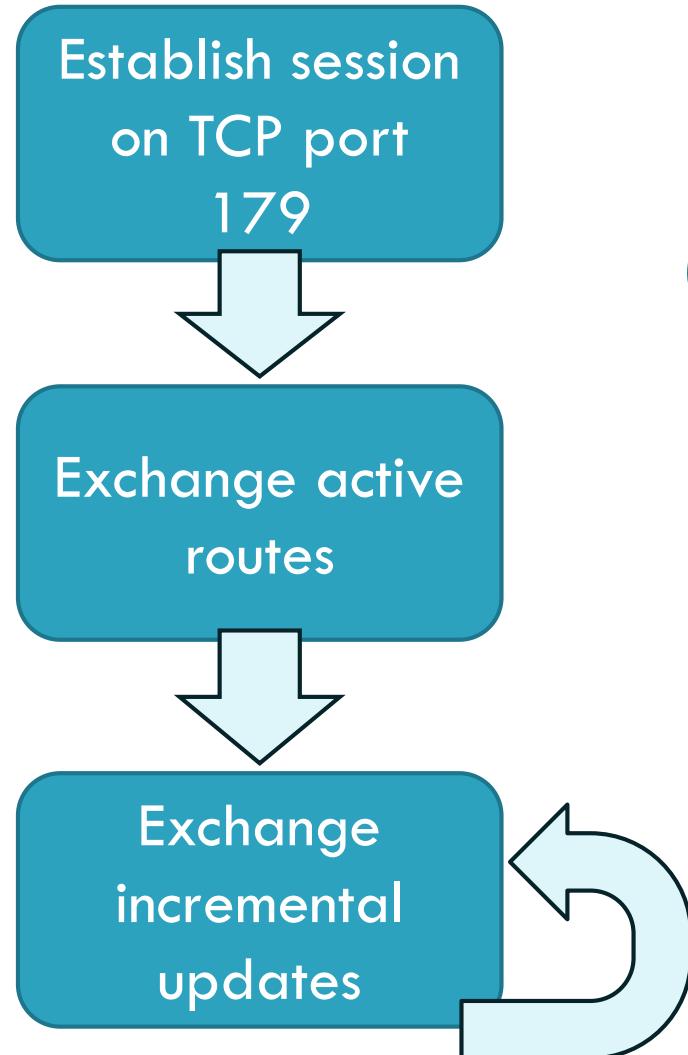
16

- AS-path: sequence of ASes a route traverses
 - Like distance vector, plus additional information
- Used for loop detection and to apply policy
- E.g., pick cheapest/shortest path
- Routing done based on longest prefix match



BGP Operations (Simplified)

17



Four Types of BGP Messages

18

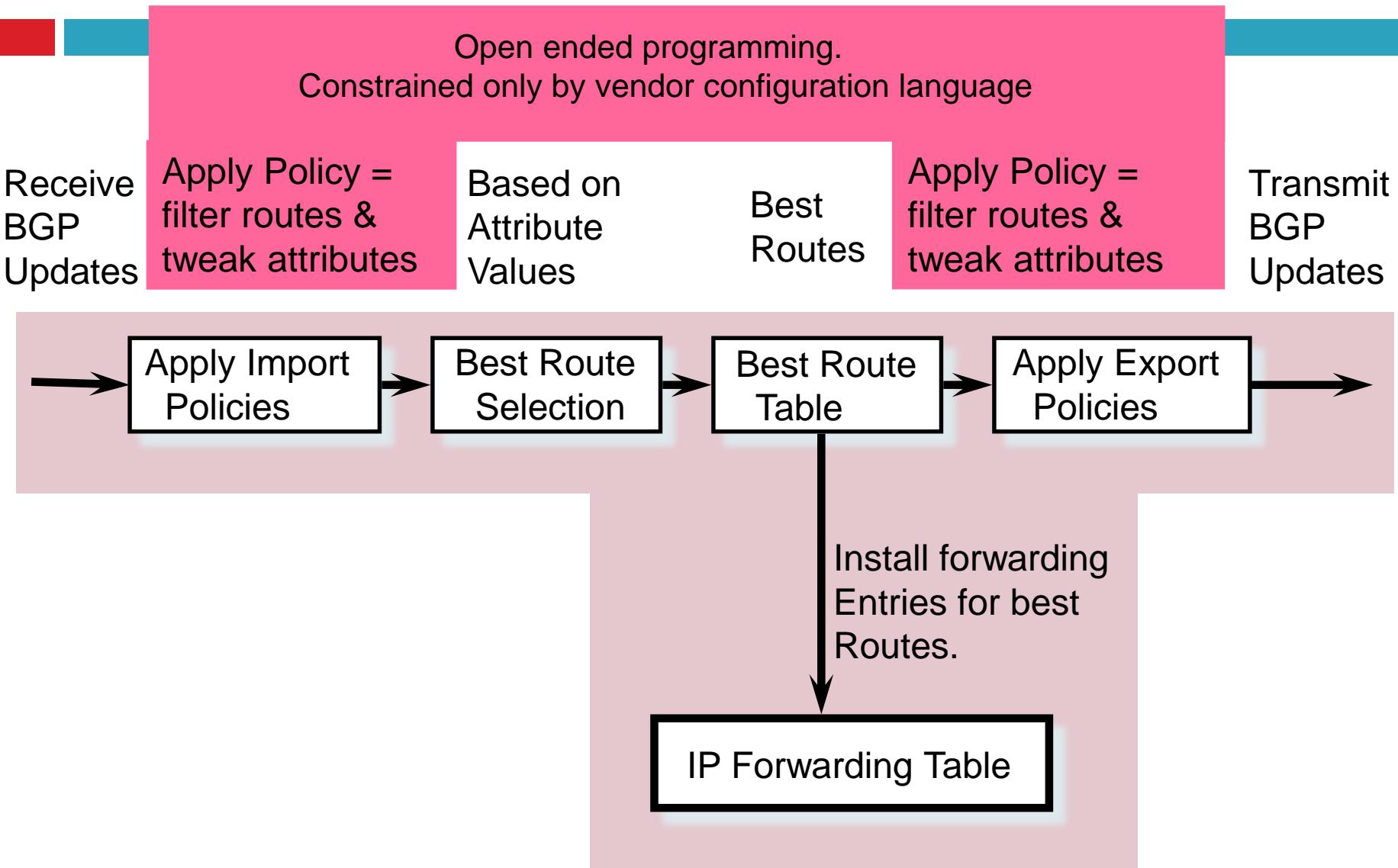
- **Open**: Establish a peering session.
- **Keep Alive**: Handshake at regular intervals.
- **Notification**: Shuts down a peering session.
- **Update**: Announce new routes or withdraw previously announced routes.

announcement = IP prefix + attributes values

Applying Policy to Routes

- Import policy
 - Q: What route advertisements do I accept?
 - Filter unwanted routes from neighbor
 - E.g. prefix that your customer doesn't own
 - Manipulate attributes to influence path selection
 - E.g., assign local preference to favored routes
- Export policy
 - Q: Which routes do I forward to whom?
 - Filter routes you don't want to tell your neighbor
 - E.g., don't tell a peer a route learned from other peer
 - Manipulate attributes to control what they see
 - E.g., make a path look artificially longer than it is

BGP Policy: Influencing Decisions



Routing Policies

- Economics
 - Enforce business relationships
 - Pick routes based on revenue and cost
 - Get traffic out of the network as early as possible
- Traffic engineering
 - Balance traffic over edge links
 - Select routes with good end-to-end performance
- Security and scalability
 - Filter routes that seem erroneous
 - Prevent the delivery of unwanted traffic
 - Limit the dissemination of small address blocks

Route Selection Summary

22

Highest Local Preference

Enforce relationships

Shortest AS Path

Lowest MED

Lowest IGP Cost to BGP Egress

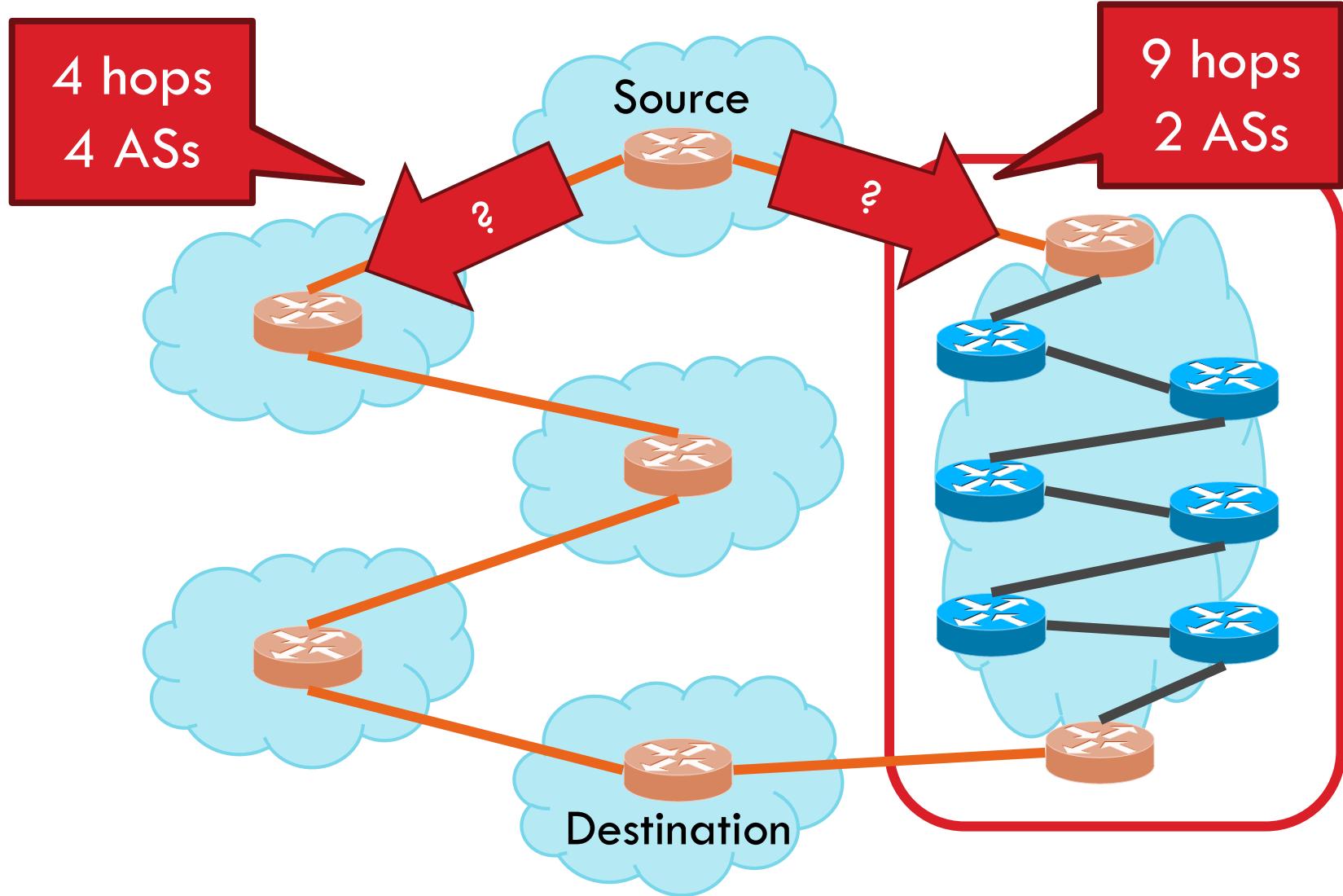
Traffic engineering

Lowest Router ID

When all else fails,
break ties

Shortest AS Path != Shortest Path

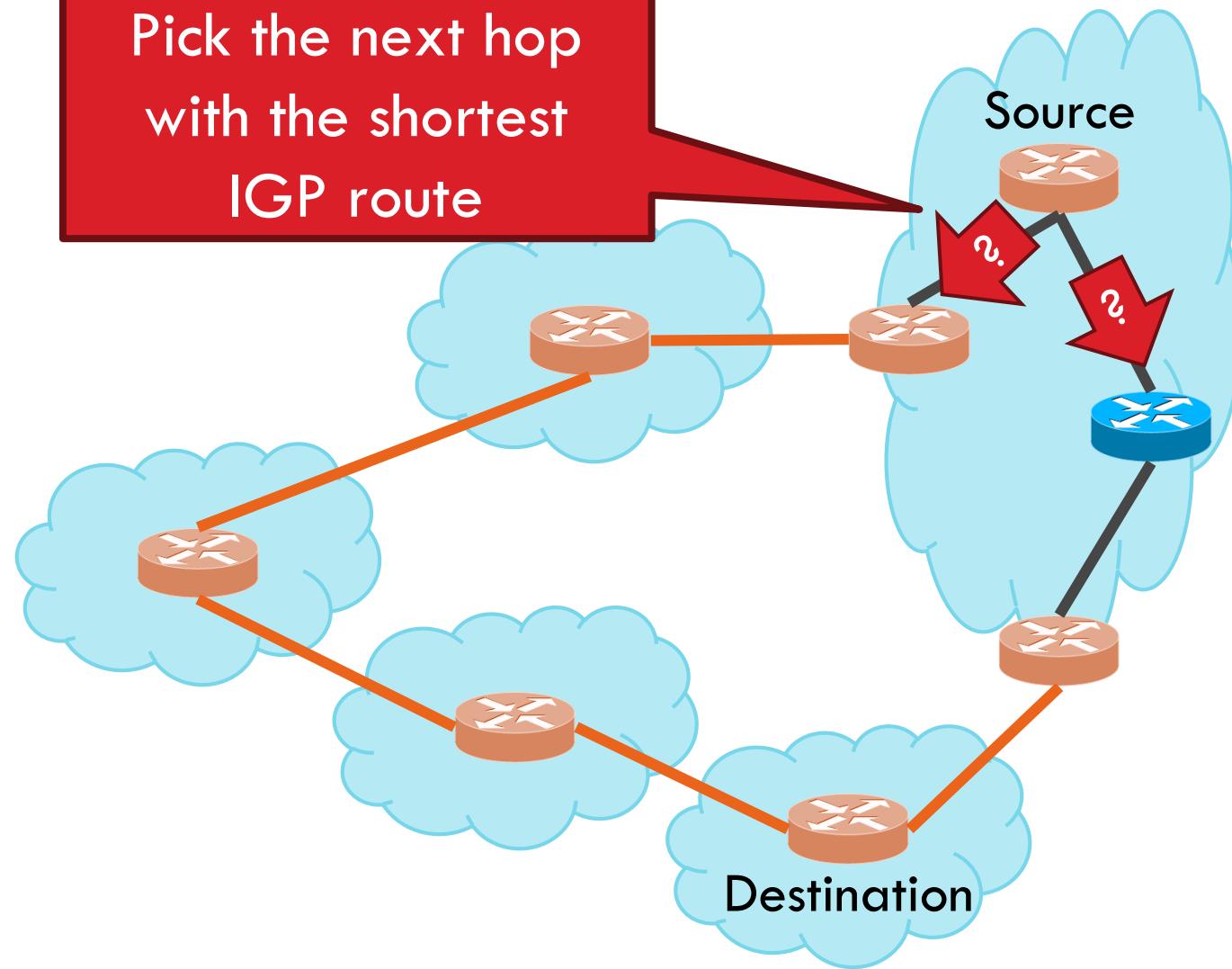
23



Hot Potato Routing

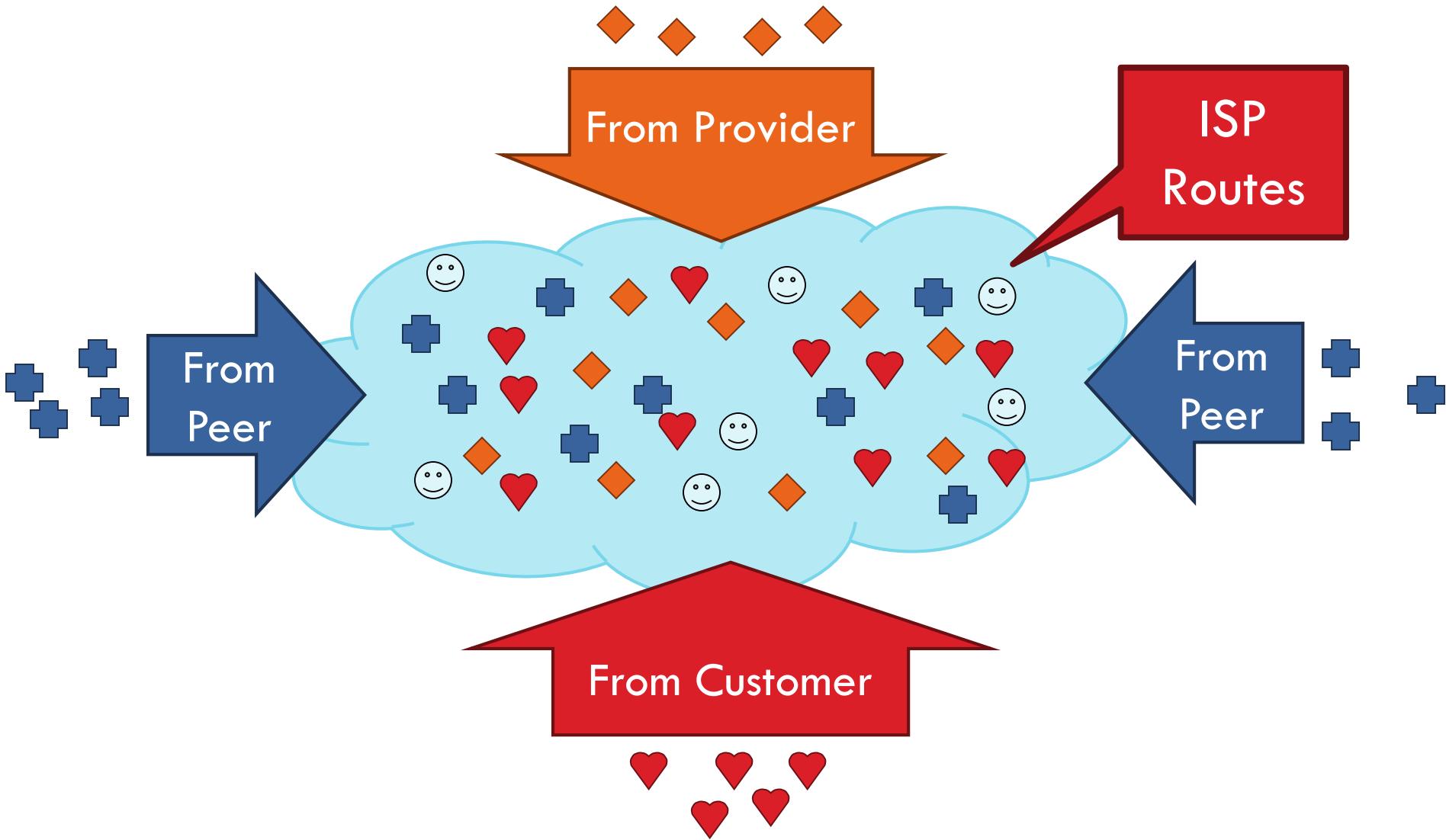
24

Pick the next hop
with the shortest
IGP route



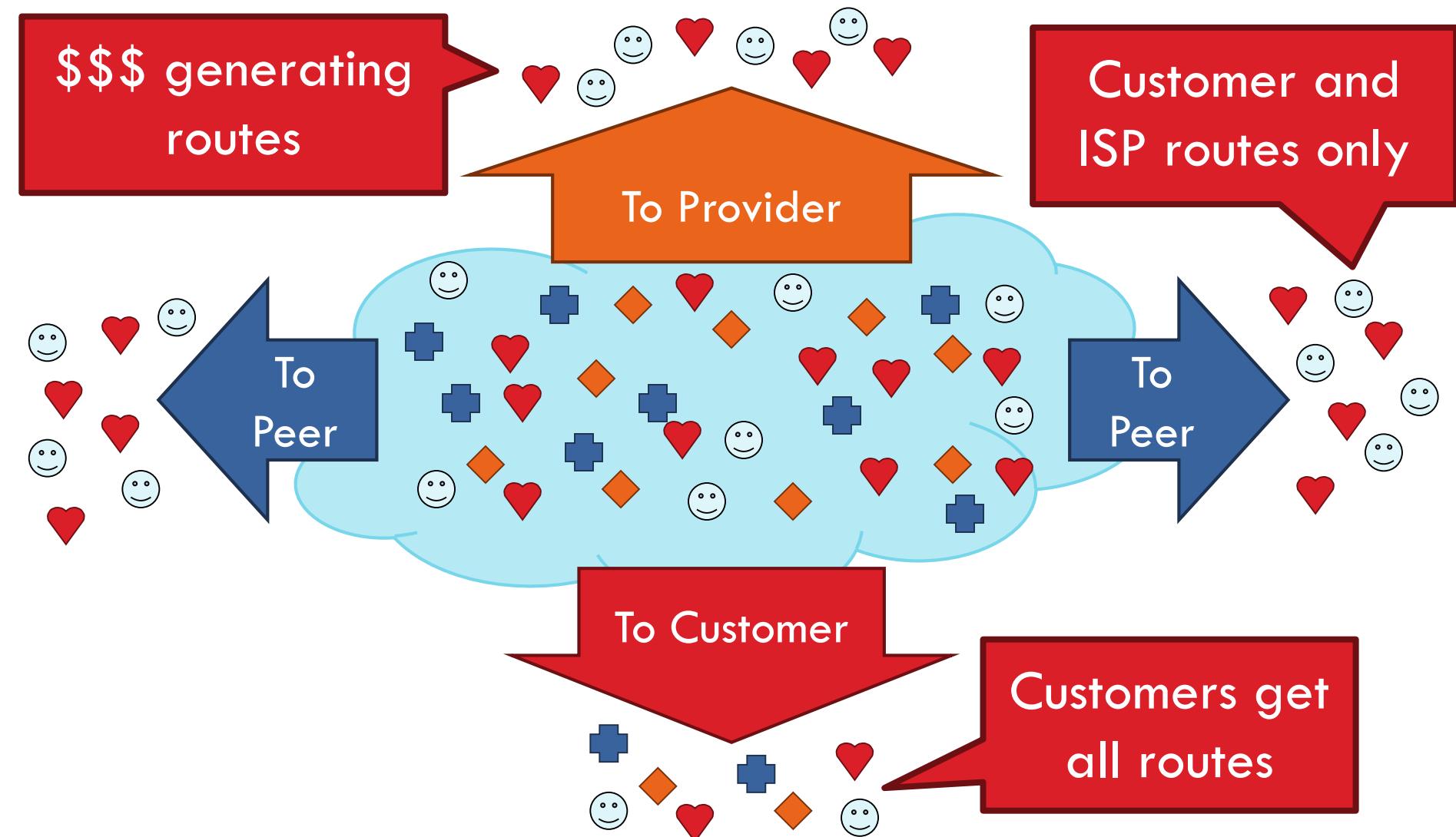
Importing Routes

25



Exporting Routes

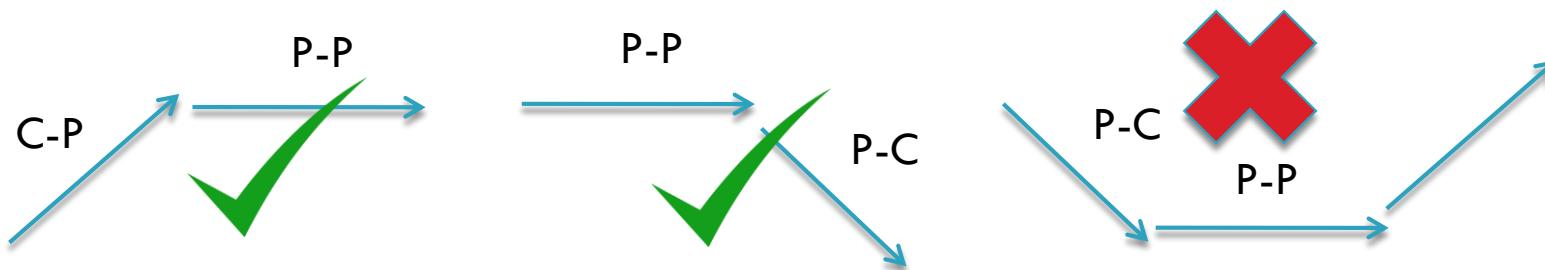
26



Modeling BGP

27

- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



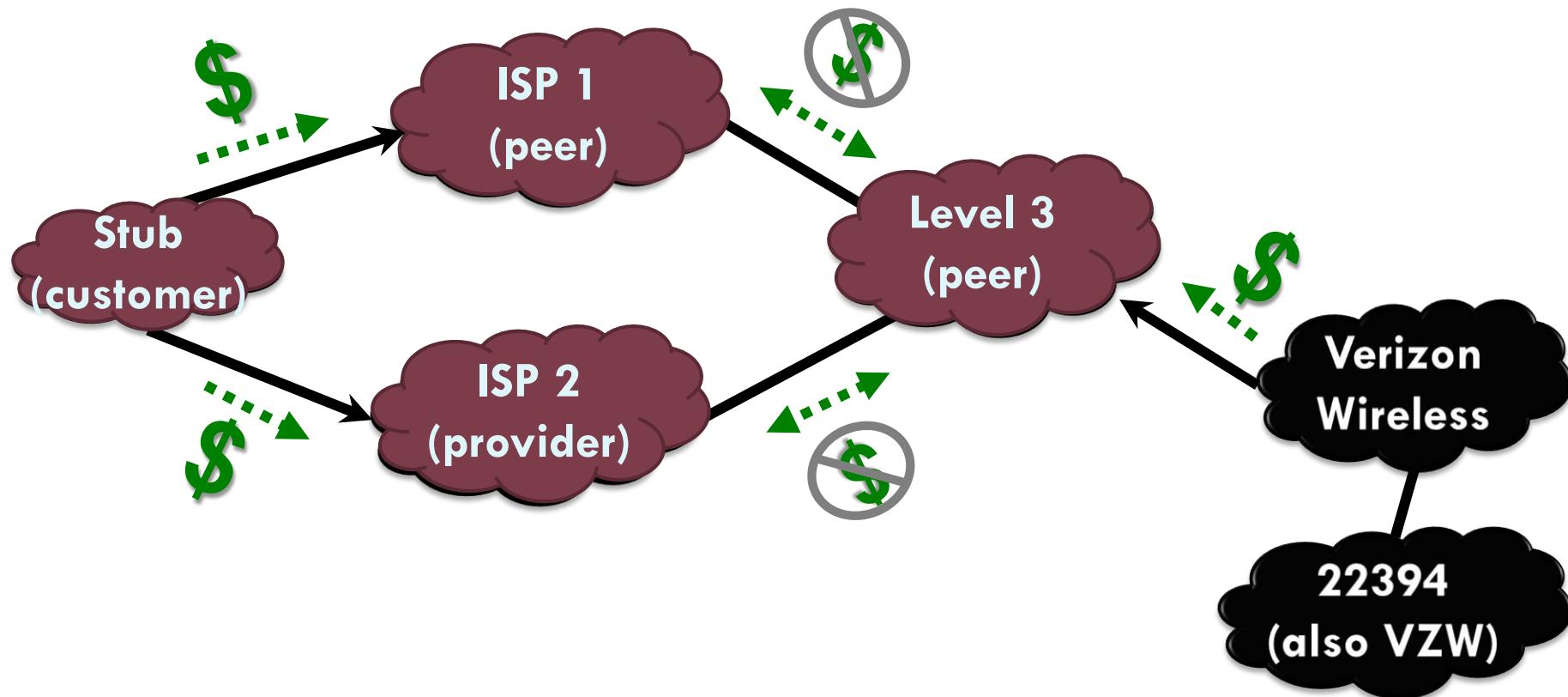
AS Relationships: It's Complicated

28

- GR Model is strictly hierarchical
 - Each AS pair has exactly one relationship
 - Each relationship is the same for all prefixes
- In practice it's much more complicated
 - Rise of widespread peering
 - Regional, per-prefix peerings
 - Tier-1's being shoved out by "hypergiants"
 - IXPs dominating traffic volume
- Modeling is very hard, very prone to error
 - Huge potential impact for understanding Internet behavior

BGP: The Internet's Routing Protocol

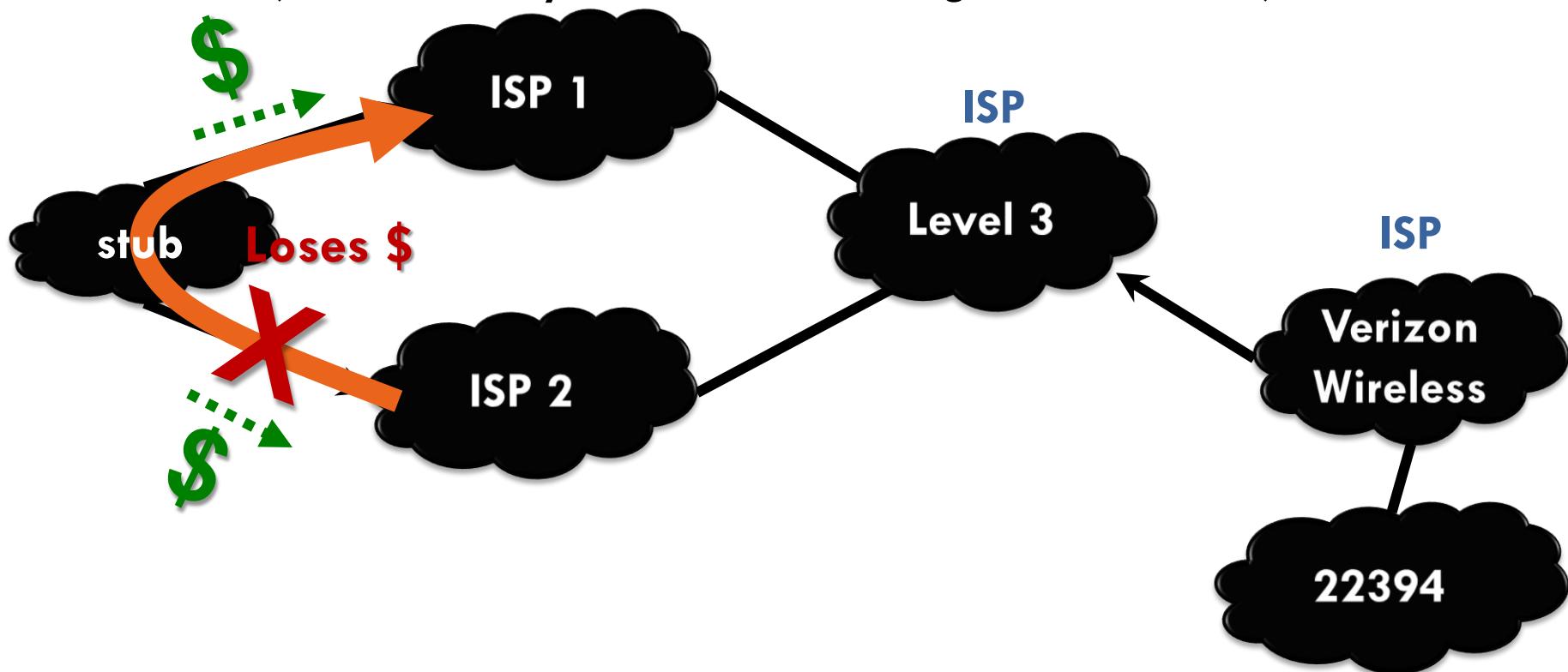
A simple model of AS-level business relationships.



BGP: The Internet's Routing Protocol (2)

A stub is an AS with no customers that never transits traffic.

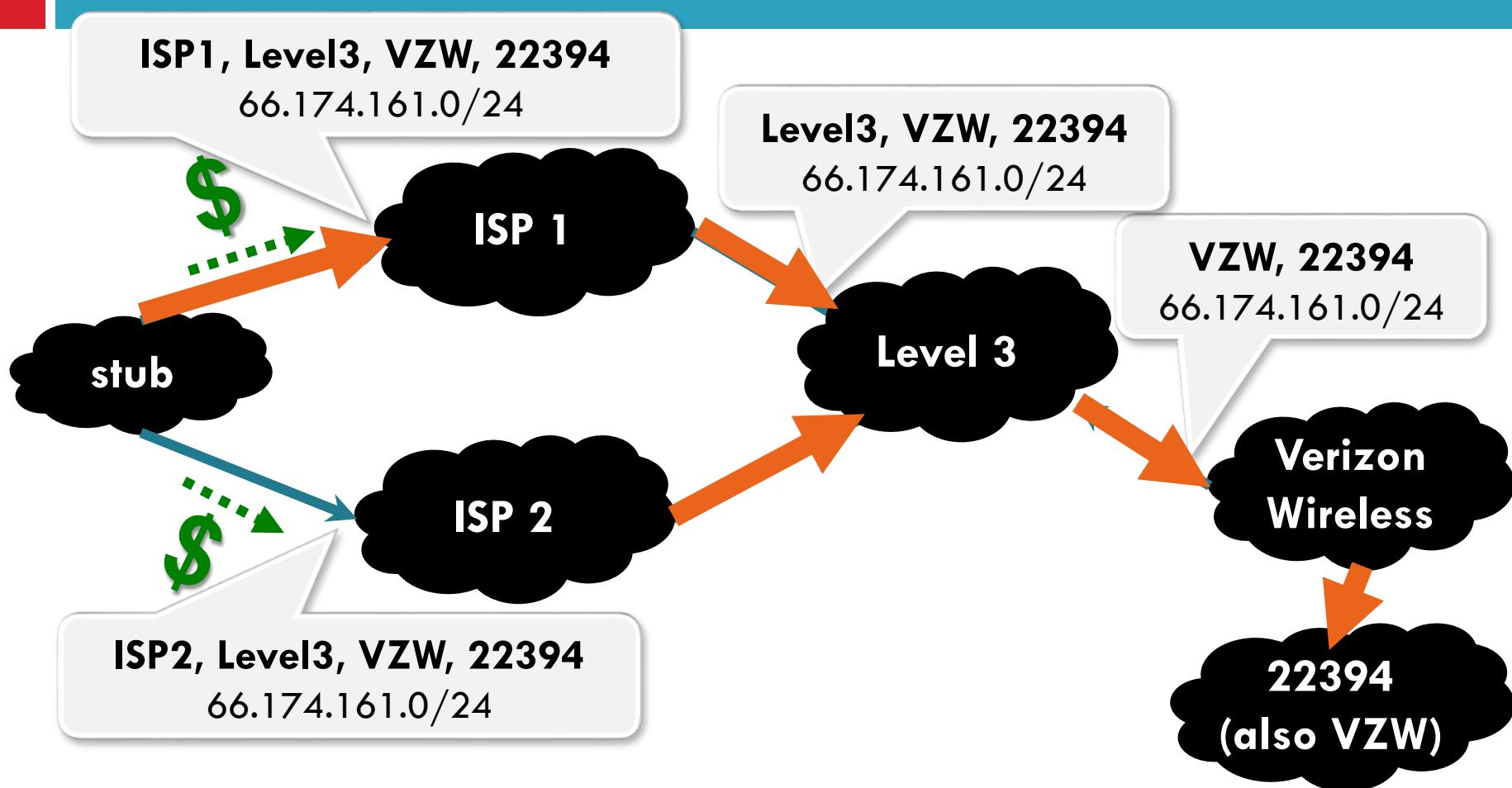
(Transit = carry traffic from one neighbor to another)



85% of ASes are stubs!
We call the rest (15%) ISPs.

BGP: The Internet's Routing Protocol (3)

BGP sets up paths from ASes to destination IP prefixes.



A model of BGP routing policies:

Prefer cheaper paths. Then, prefer shorter paths.

Standard model of Internet routing

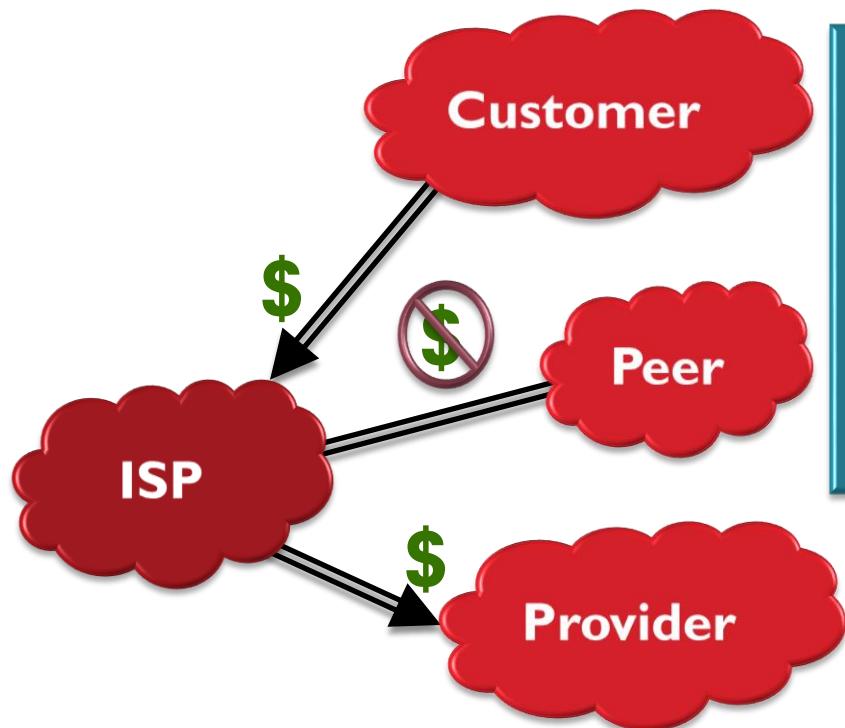
33

- Proposed by Gao & Rexford 12 years ago
- Based on practices employed by a large ISP
- Provide an intuitive model of path selection and export policy

Standard model of Internet routing

34

- Proposed by Gao & Rexford 12 years ago



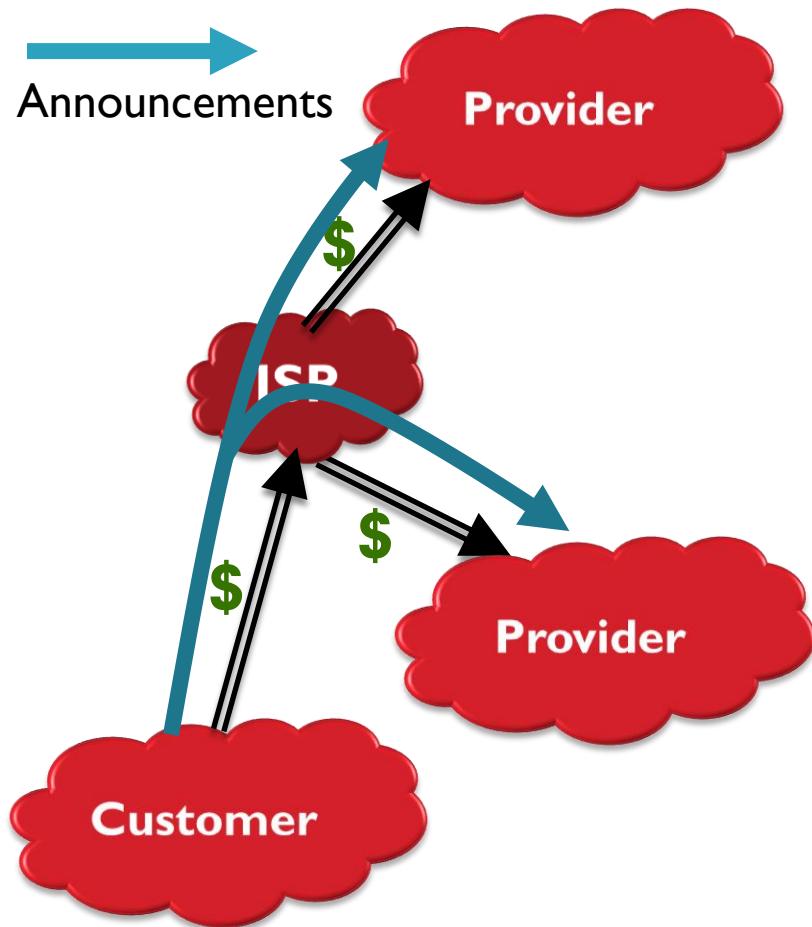
Path Selection:

- I. LocalPref: Prefer customer paths over peer paths over provider paths
2. Prefer shorter paths
3. Arbitrary tiebreak

Standard model of Internet routing

35

- Proposed by Gao & Rexford 12 years ago



Path Selection:

- I. LocalPref: Prefer customer paths over peer paths over provider paths
2. Prefer shorter paths
3. Arbitrary tiebreak

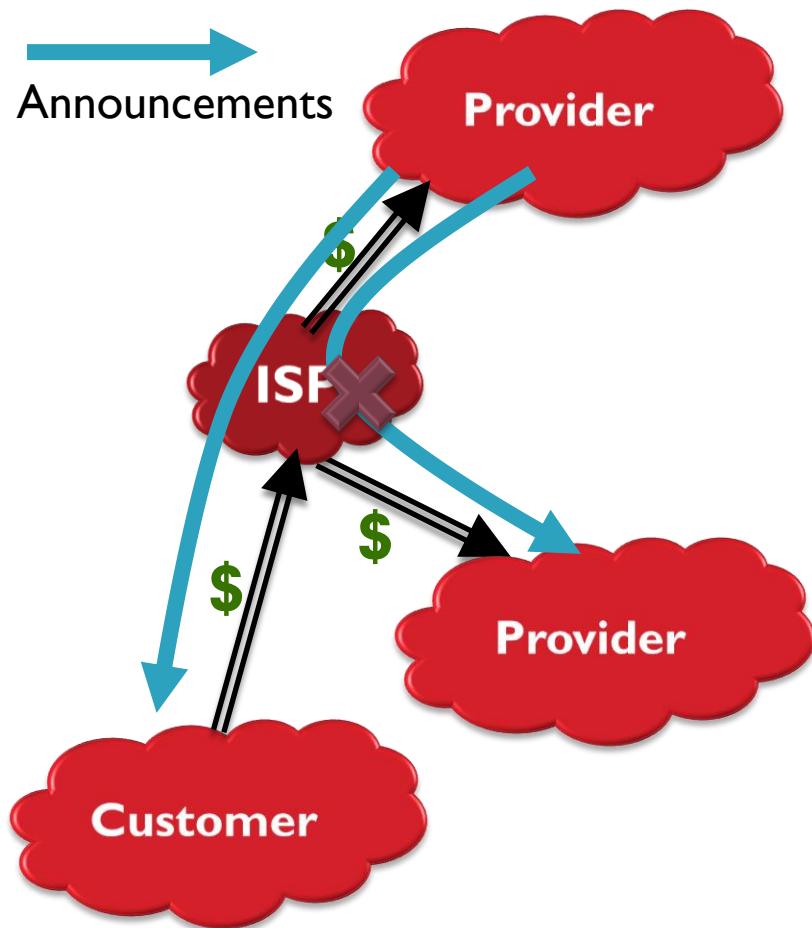
Export Policy:

1. Export customer path to all neighbors.
2. Export peer/provider path to all customers.

Standard model of Internet routing

36

- Proposed by Gao & Rexford 12 years ago



Path Selection:

- I. LocalPref: Prefer customer paths over peer paths over provider paths
2. Prefer shorter paths
3. Arbitrary tiebreak

Export Policy:

1. Export customer path to all neighbors.
2. Export peer/provider path to all customers.

BGP-related Hijacks

Normal operation

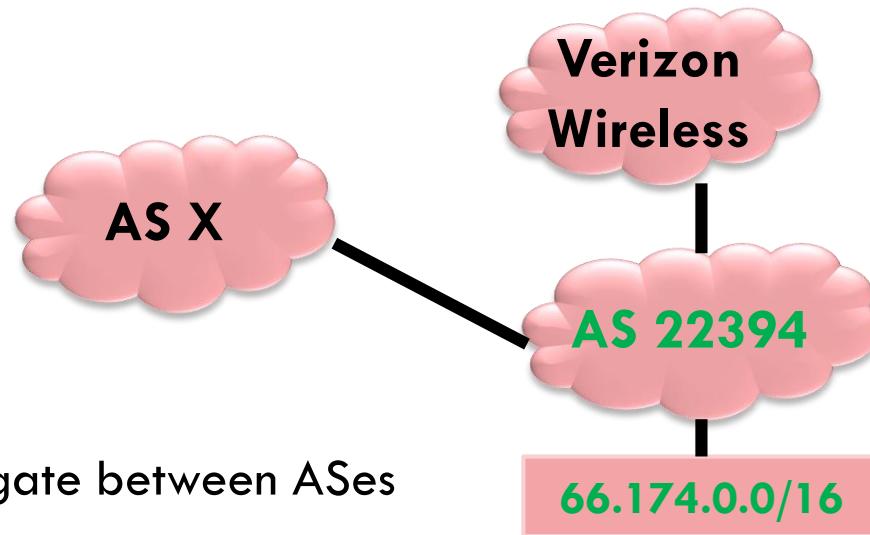
- Origin AS announces prefix
- Route announcements propagate between ASes
- Helps ASes learn about “good” paths to reach prefix



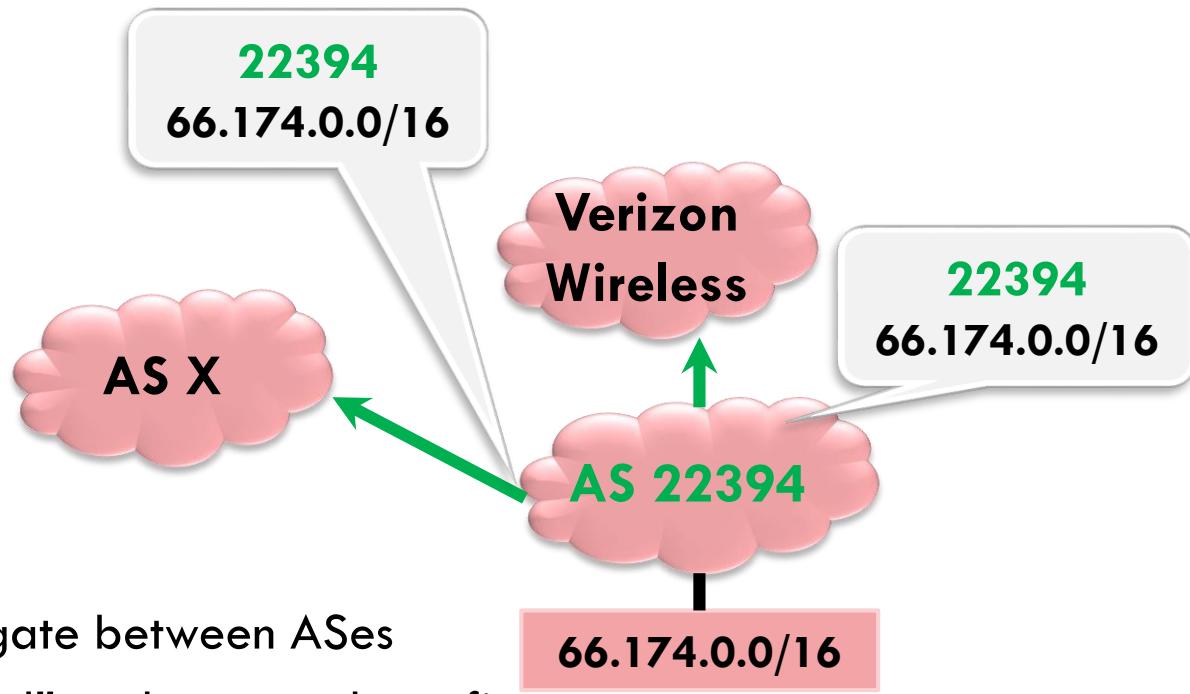
BGP-related Hijacks

Normal operation

- Origin AS announces prefix
- Route announcements propagate between ASes
- Helps ASes learn about “good” paths to reach prefix



BGP-related Hijacks



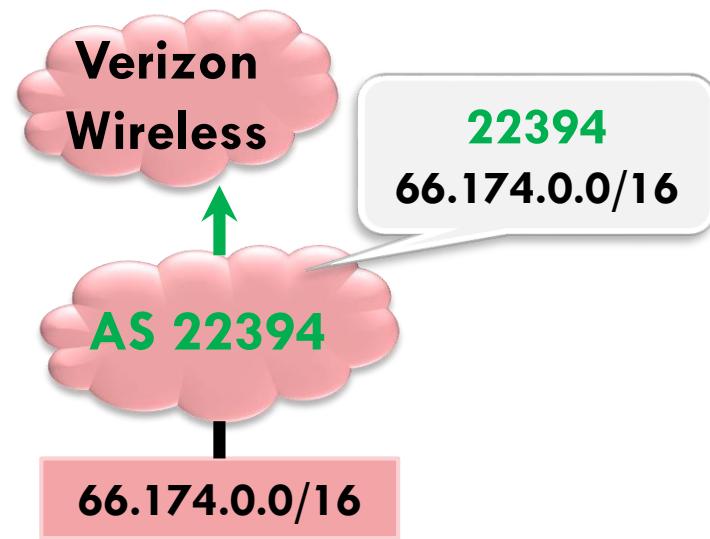
Normal operation

- Origin AS announces prefix
- Route announcements propagate between ASes
- Helps ASes learn about “good” paths to reach prefix

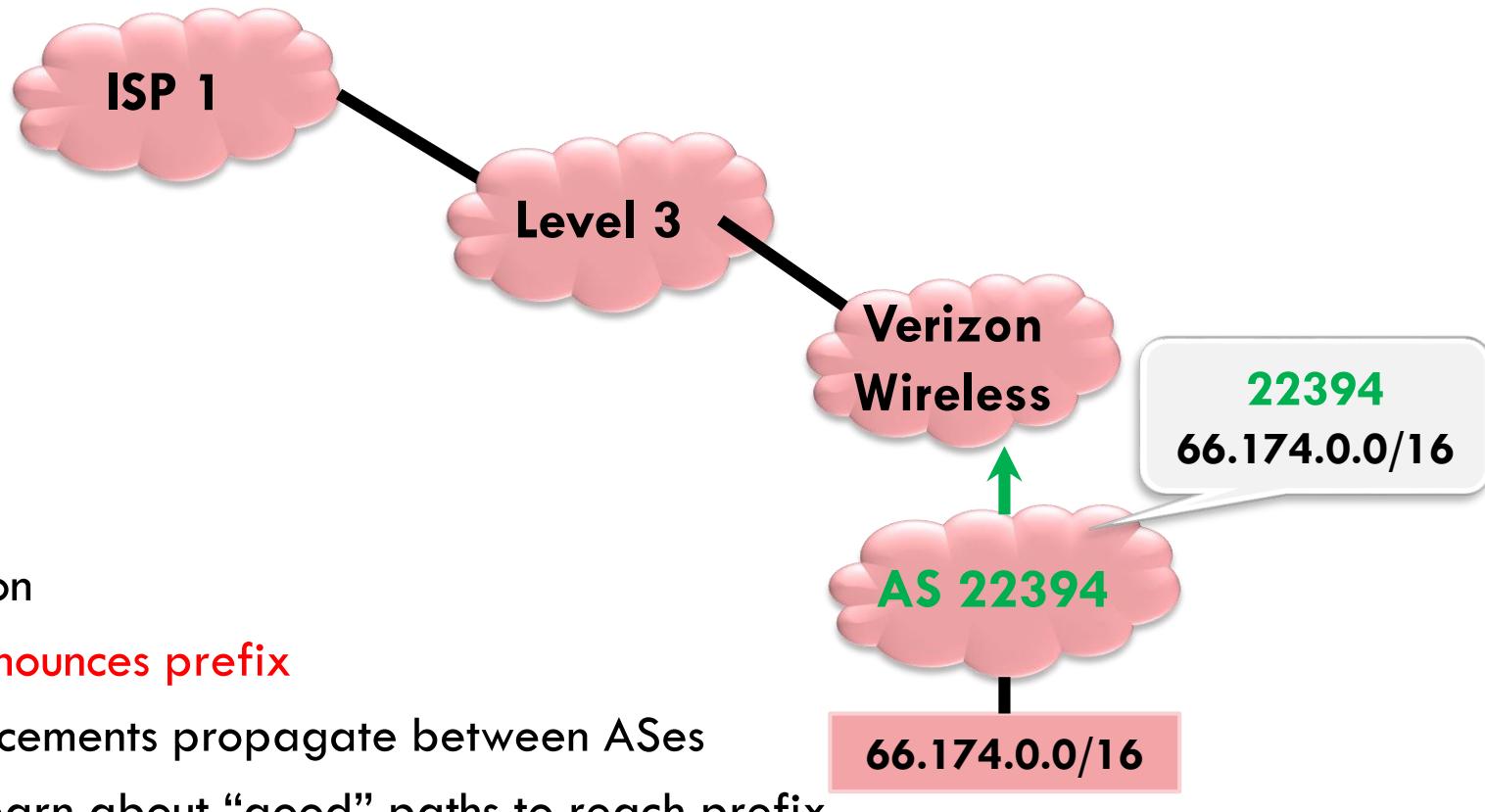
BGP-related Hijacks

Normal operation

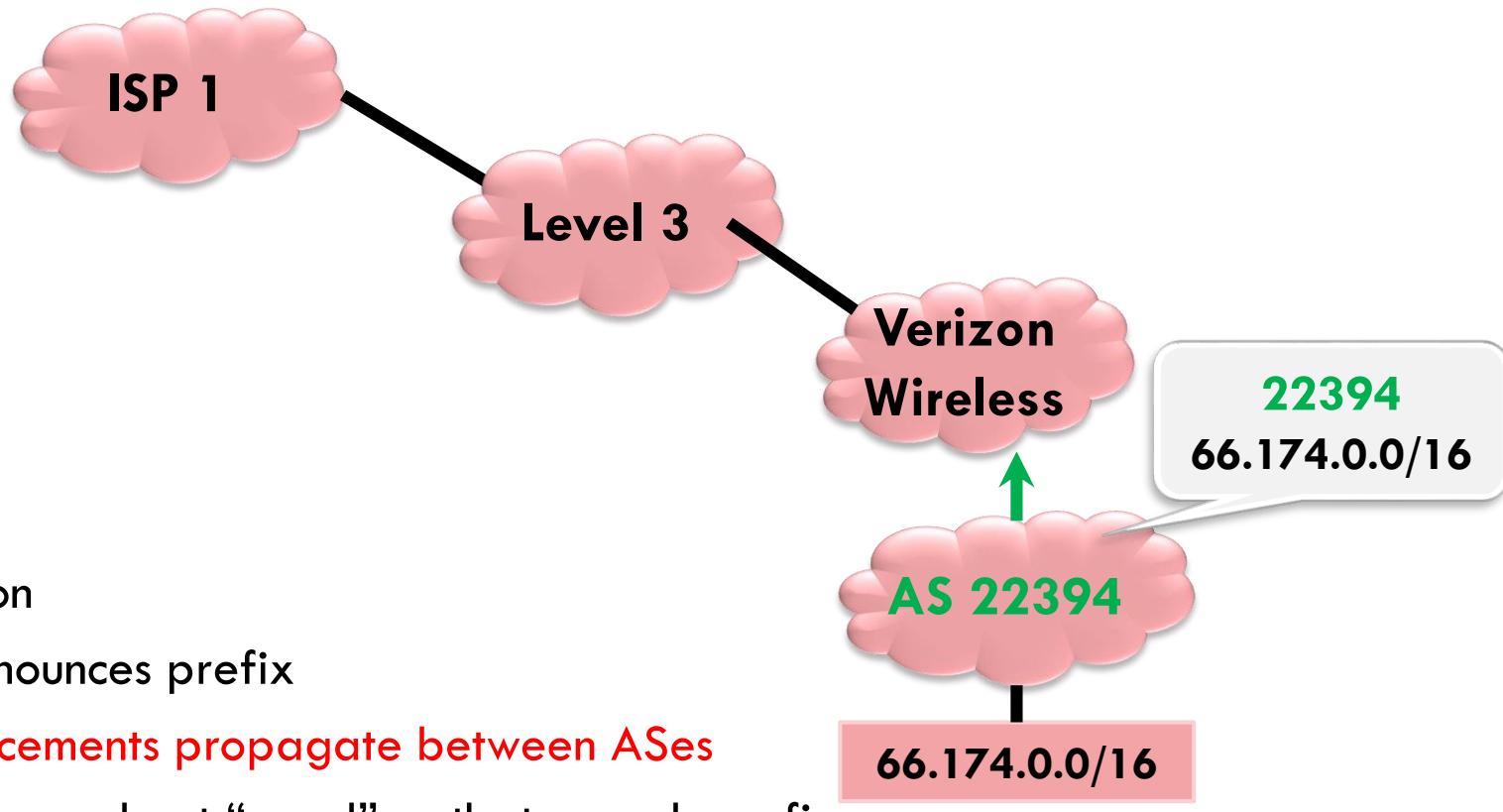
- Origin AS announces prefix
- Route announcements propagate between ASes
- Helps ASes learn about “good” paths to reach prefix



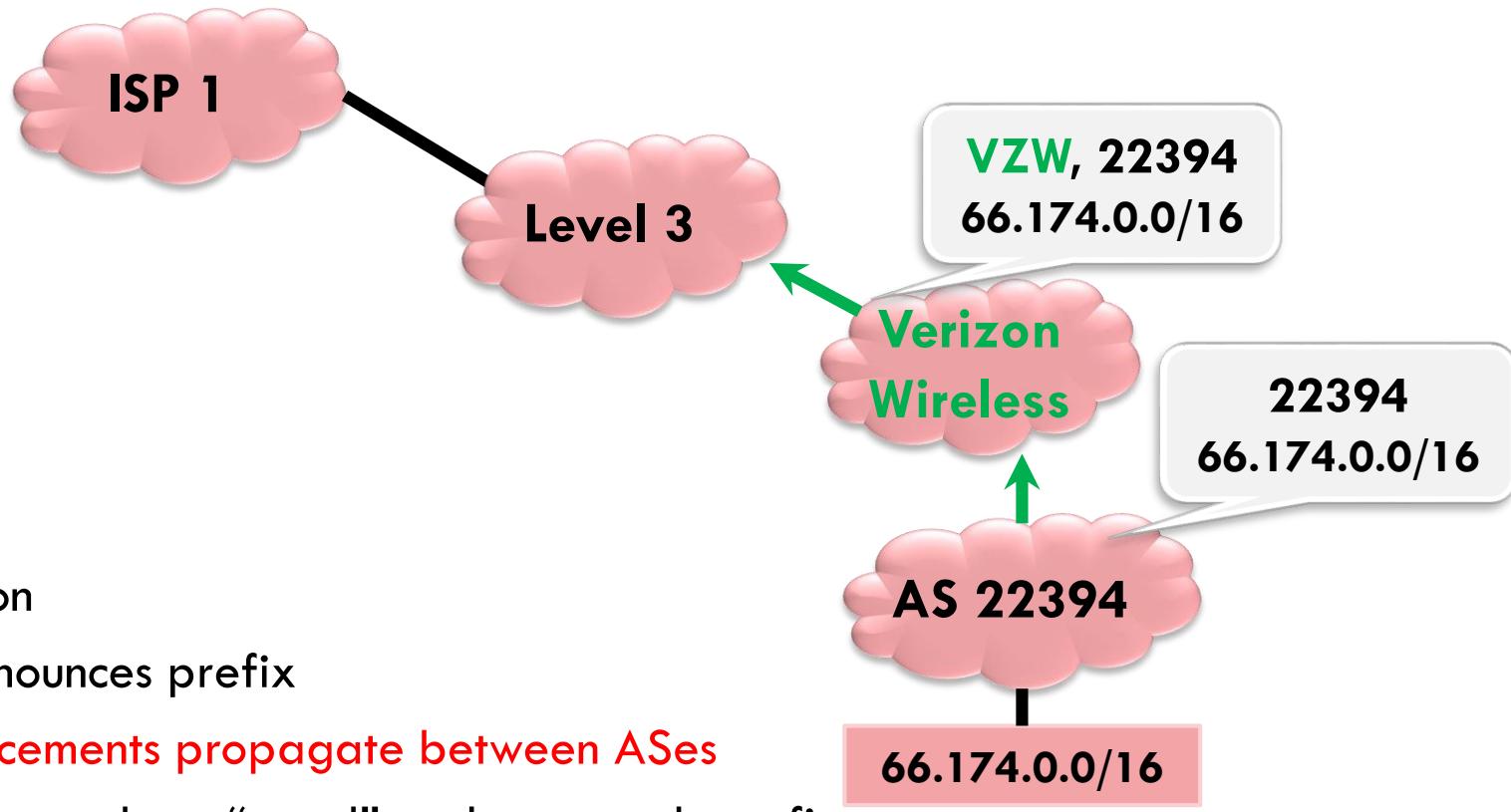
BGP-related Hijacks



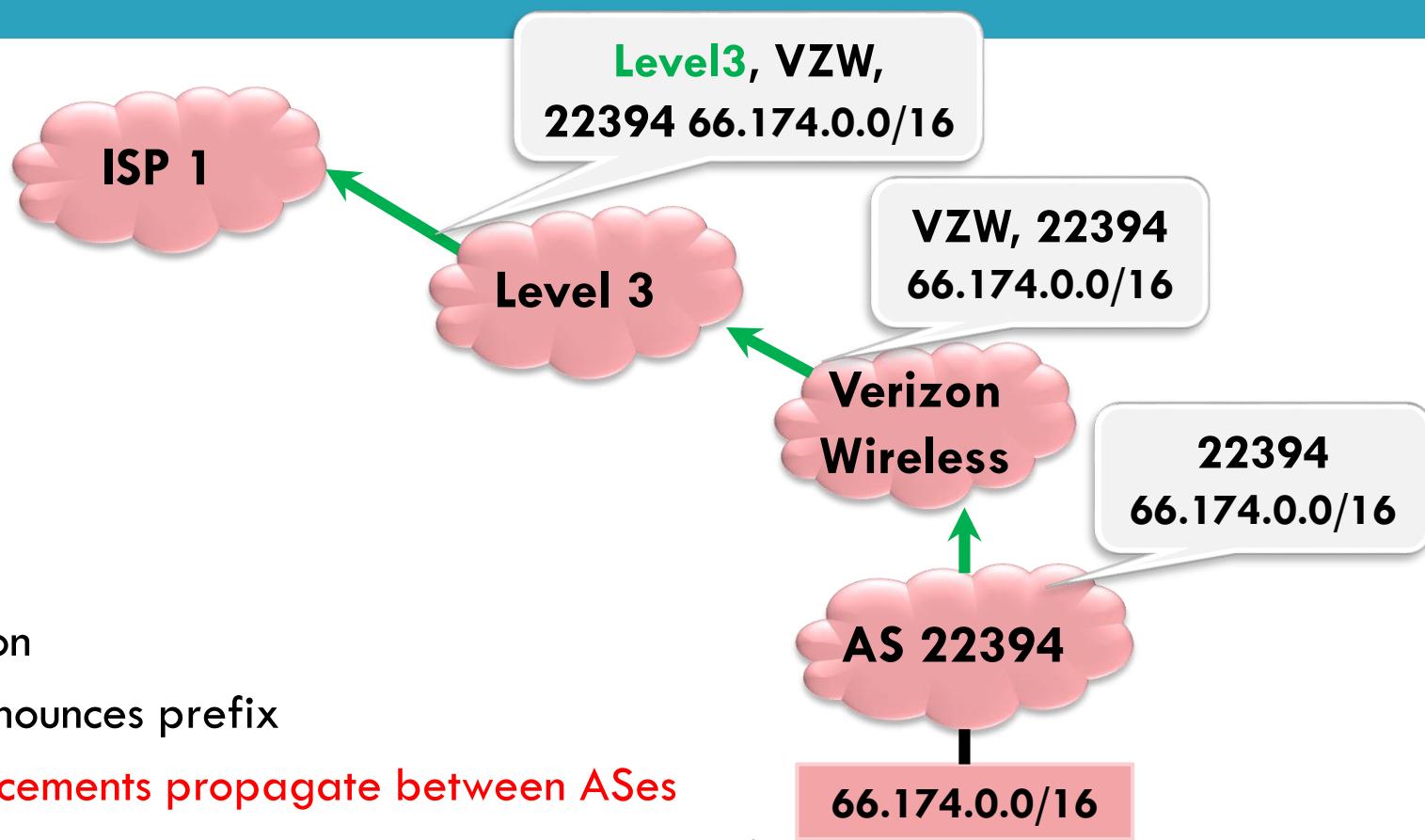
BGP-related Hijacks



BGP-related Hijacks



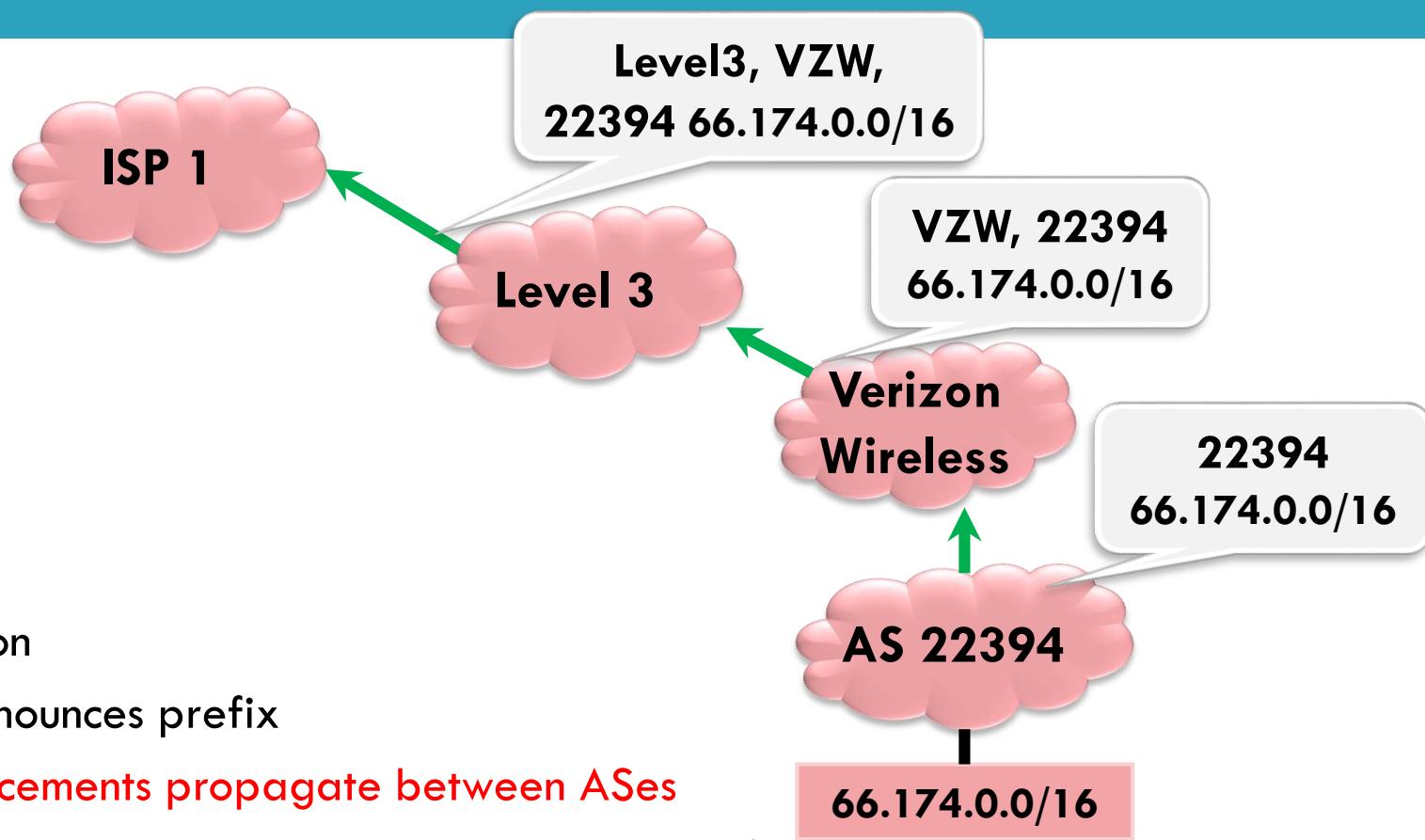
BGP-related Hijacks



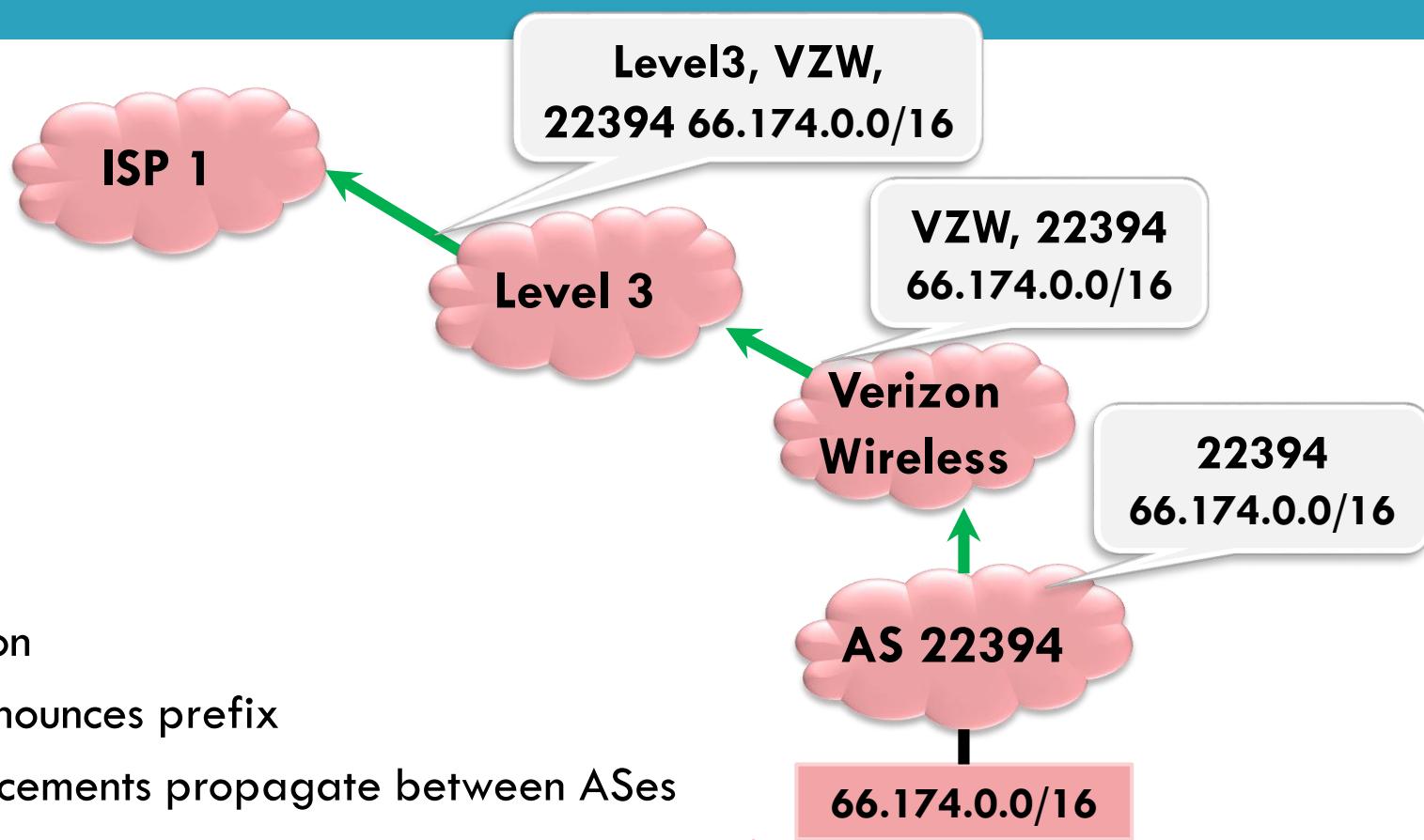
Normal operation

- Origin AS announces prefix
- **Route announcements propagate between ASes**
- Helps ASes learn about “good” paths to reach prefix

BGP-related Hijacks



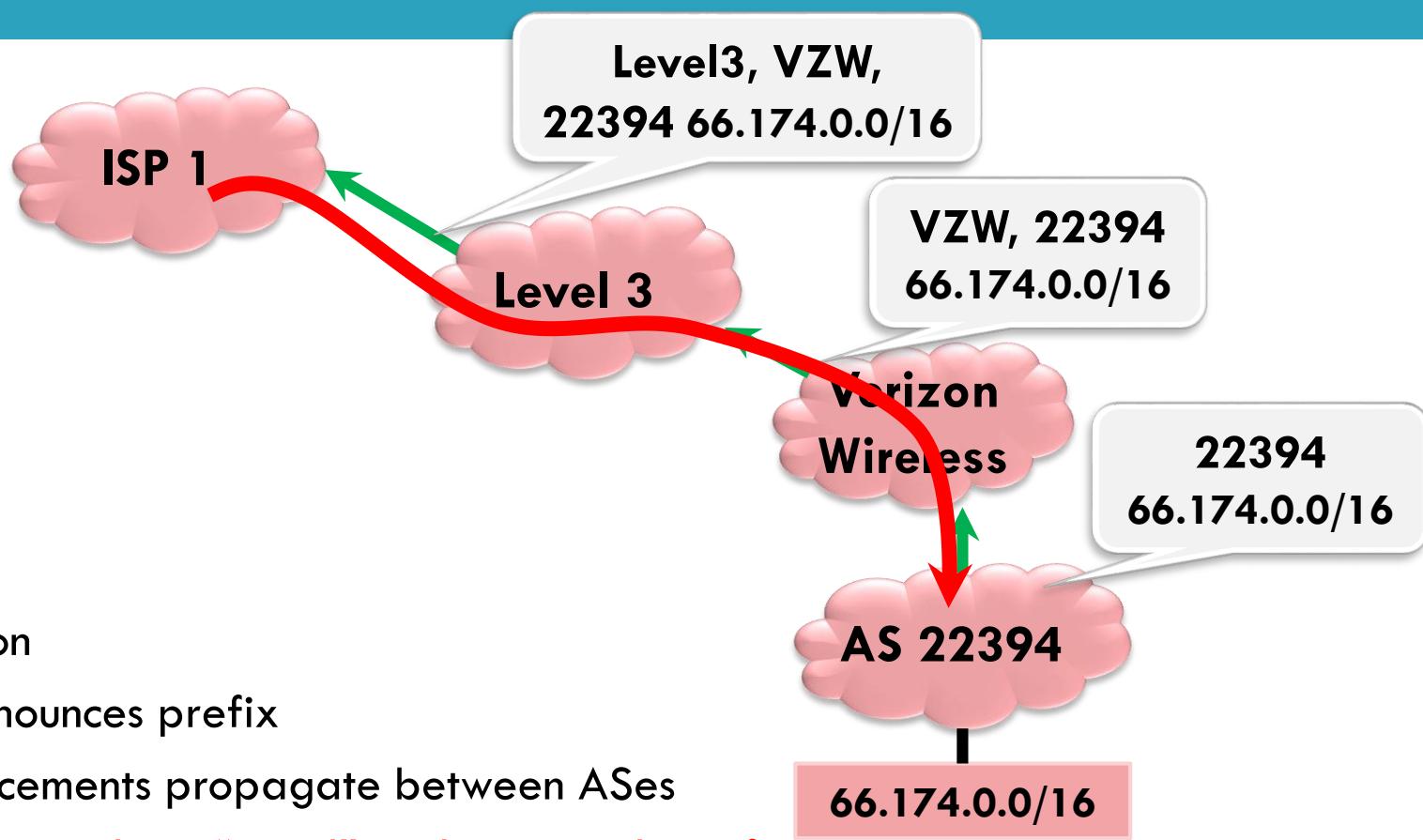
BGP-related Hijacks



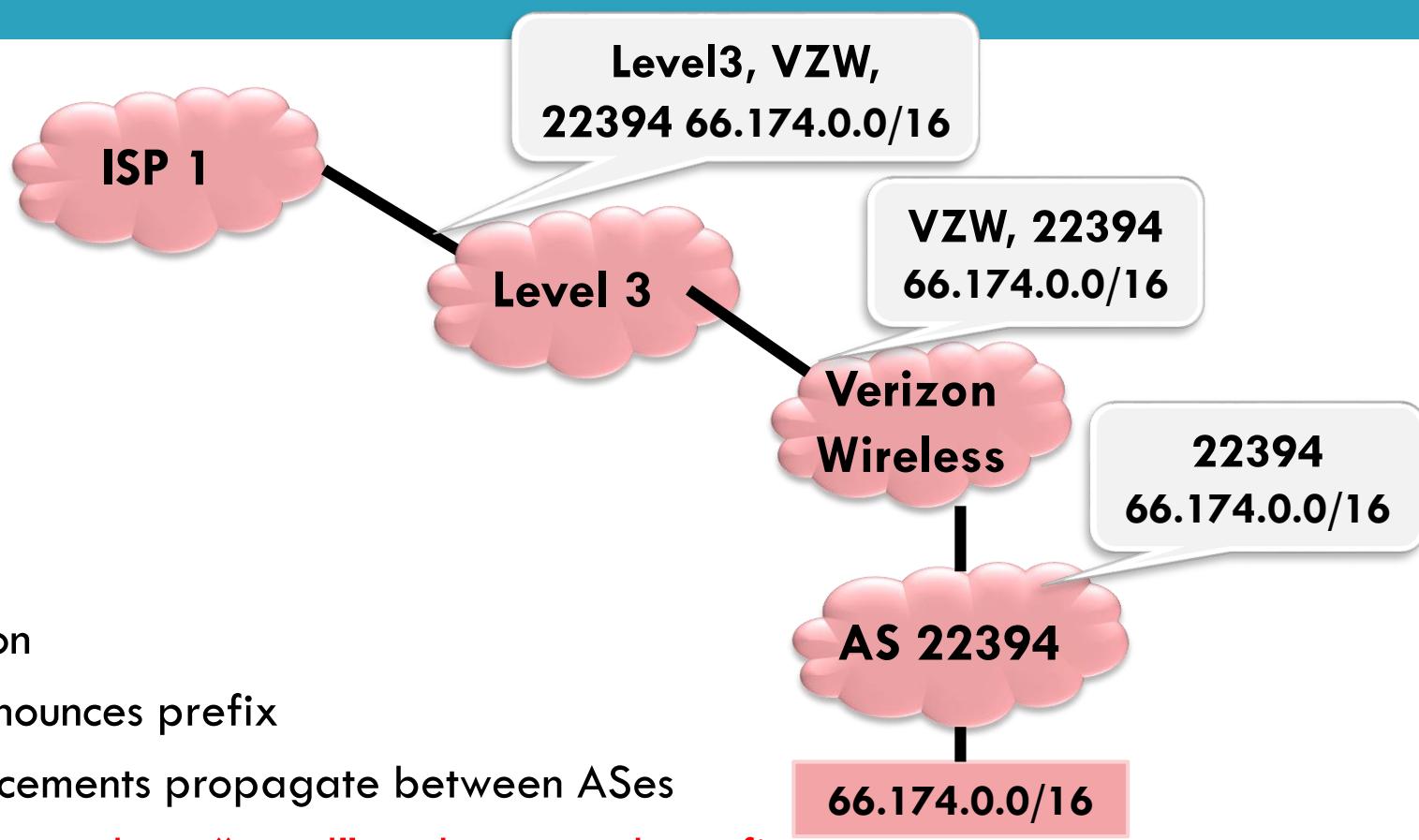
Normal operation

- Origin AS announces prefix
- Route announcements propagate between ASes
- Helps ASes learn about “good” paths to reach prefix

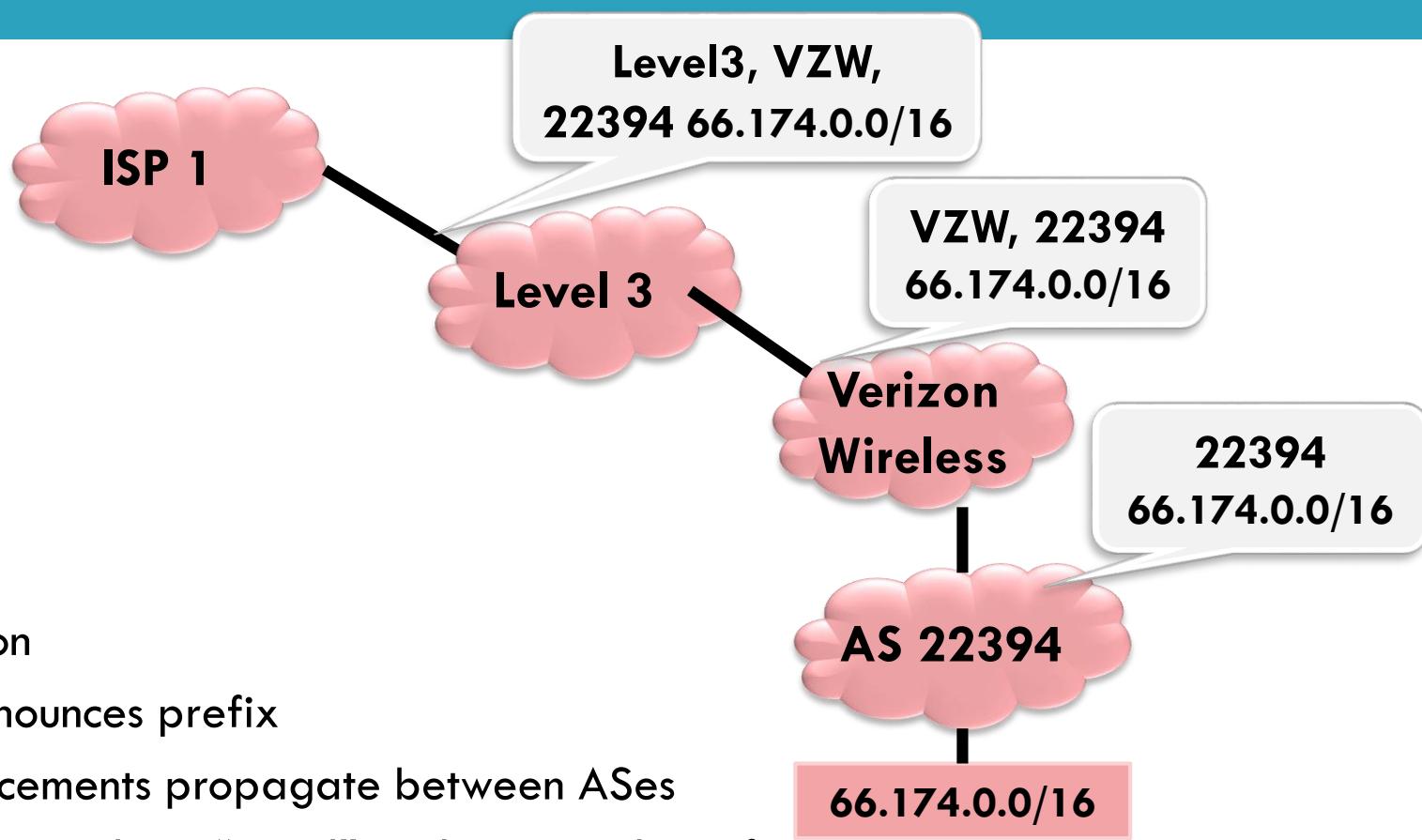
BGP-related Hijacks



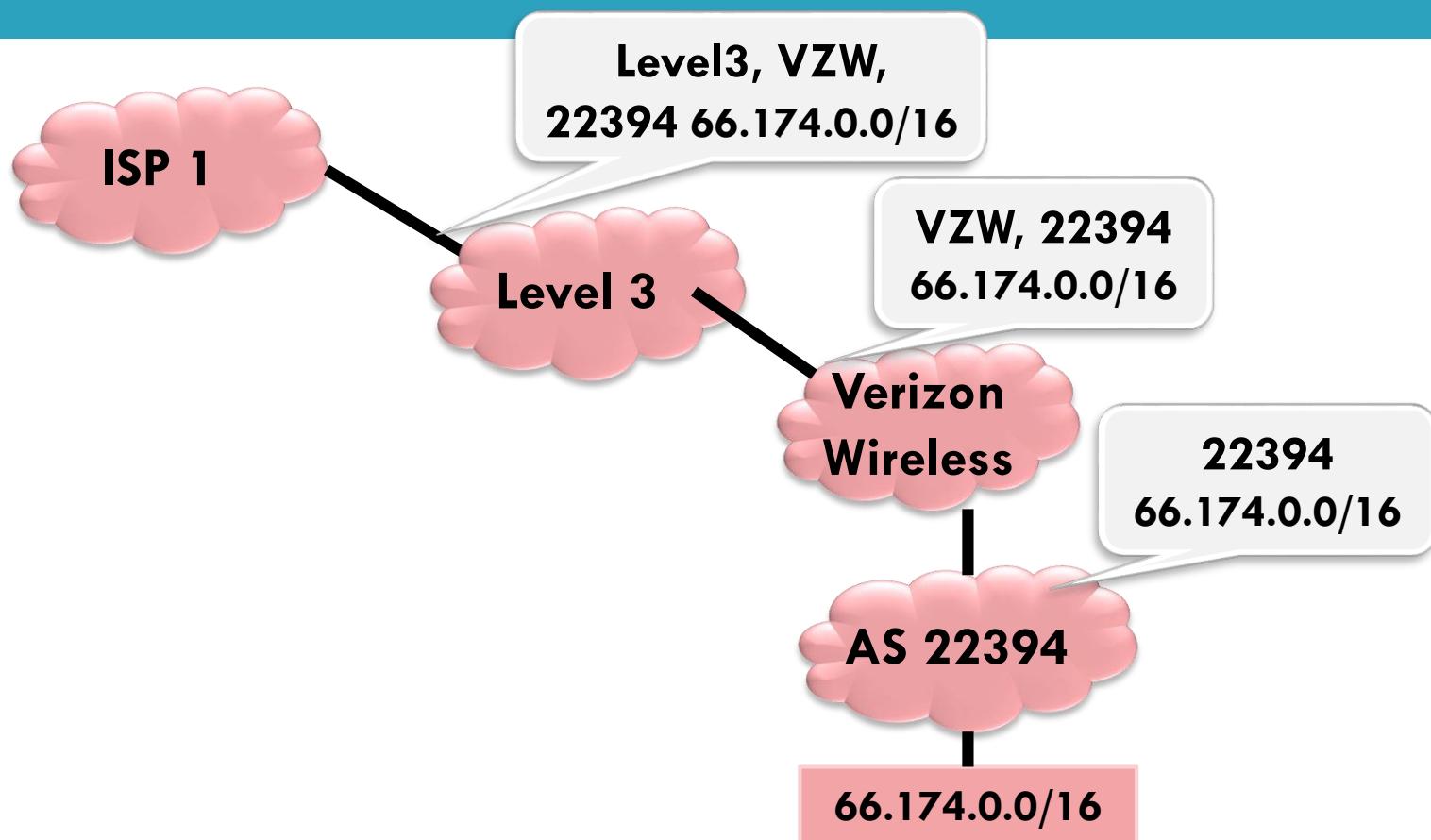
BGP-related Hijacks



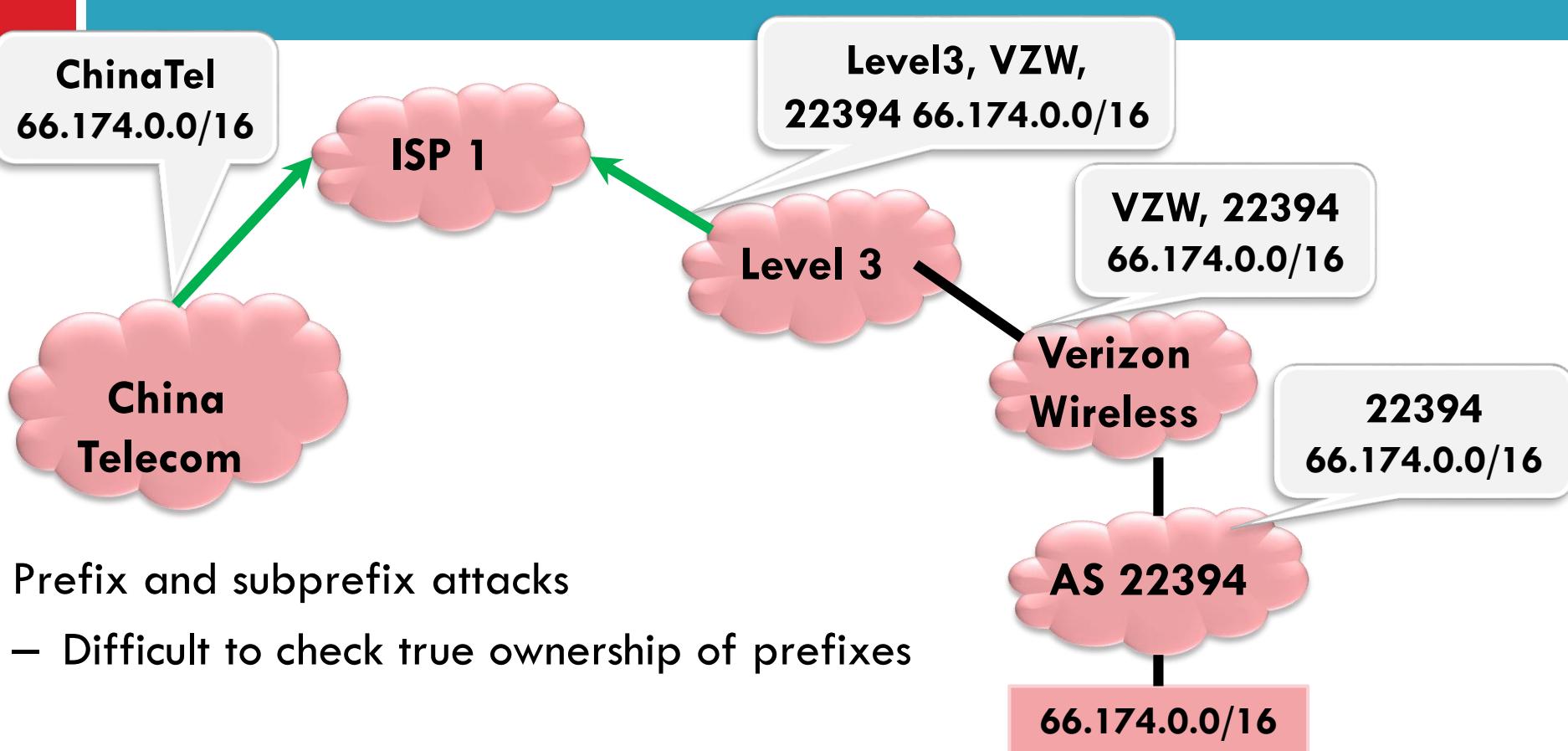
BGP-related Hijacks



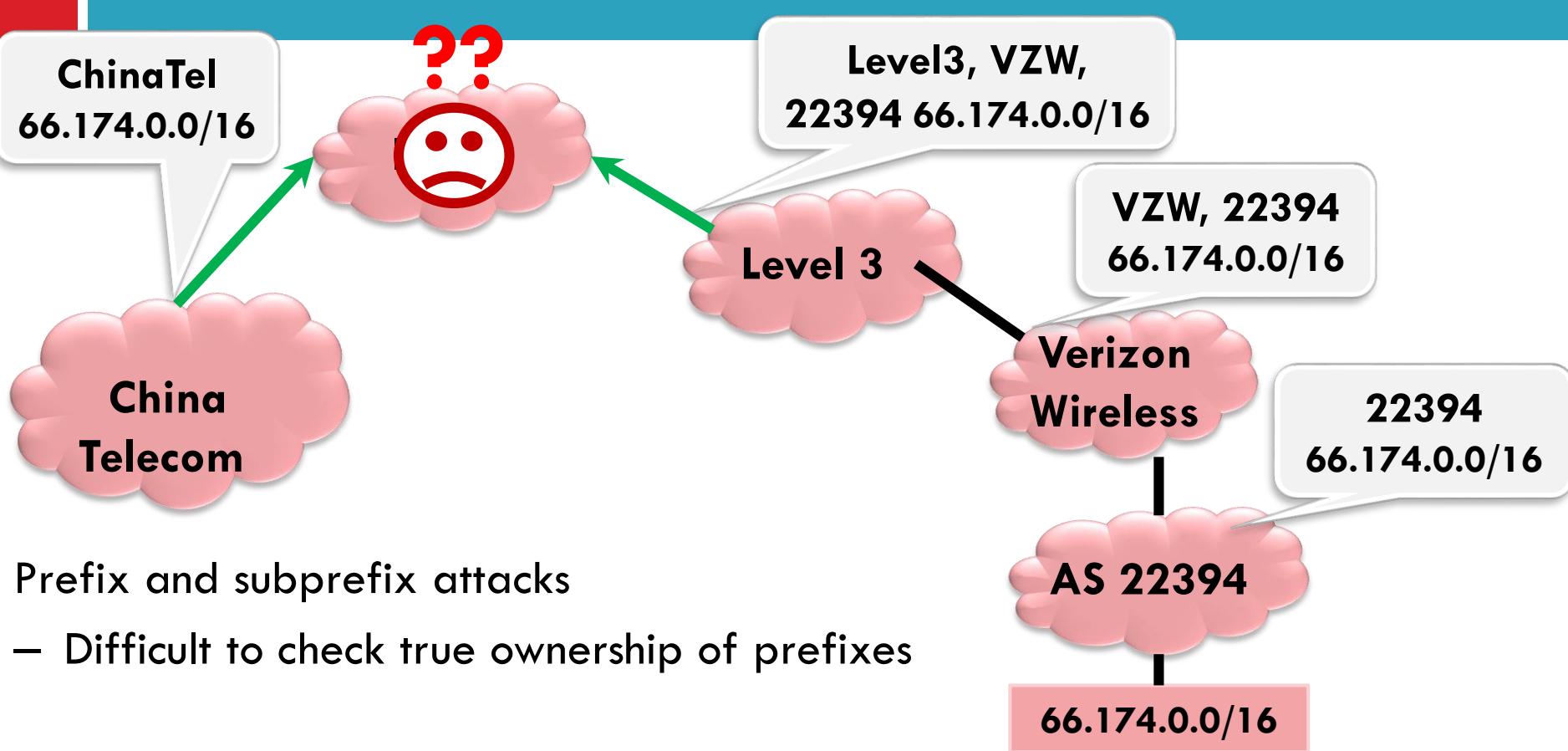
BGP-related Hijacks



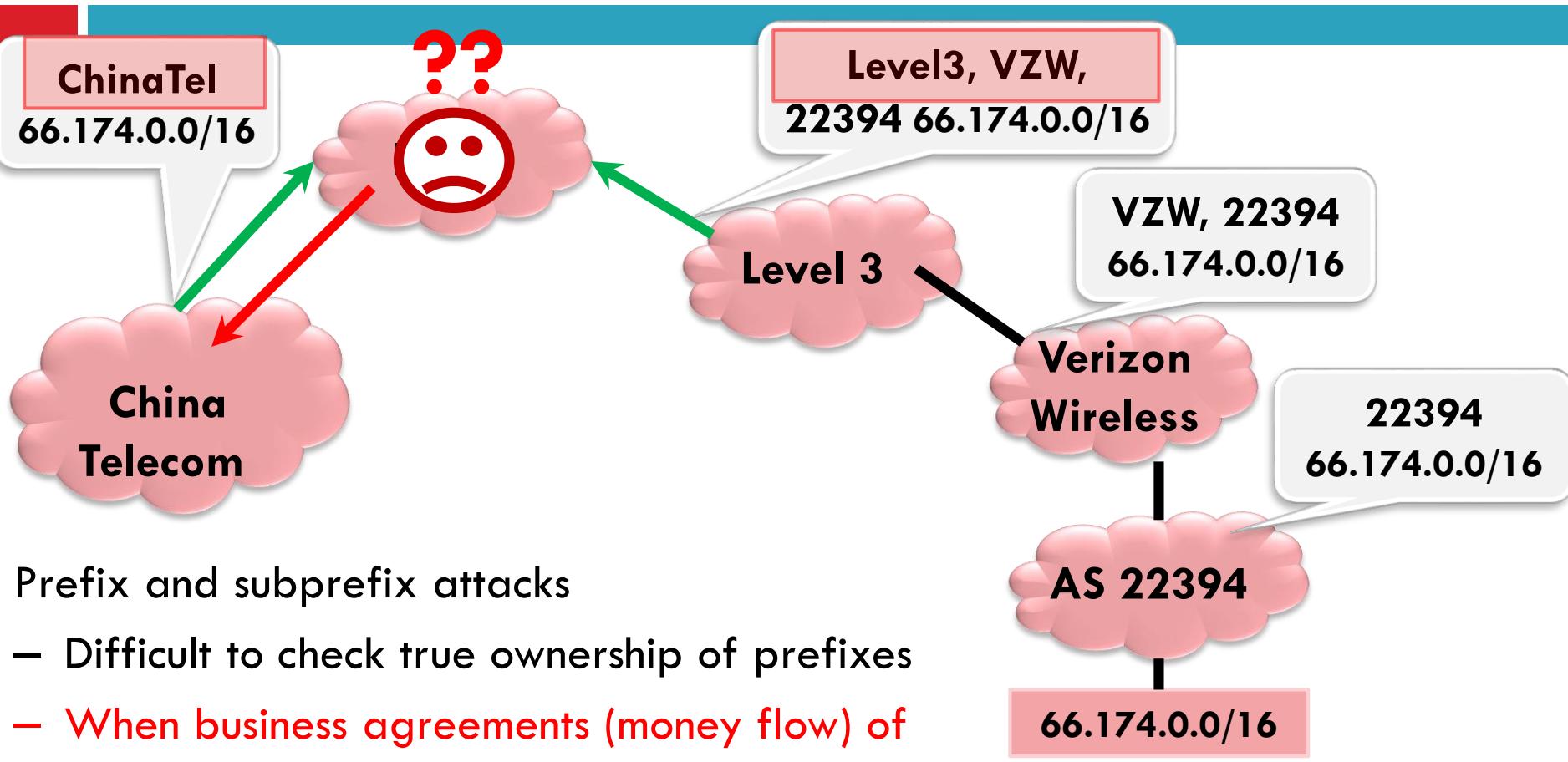
BGP-related Hijacks



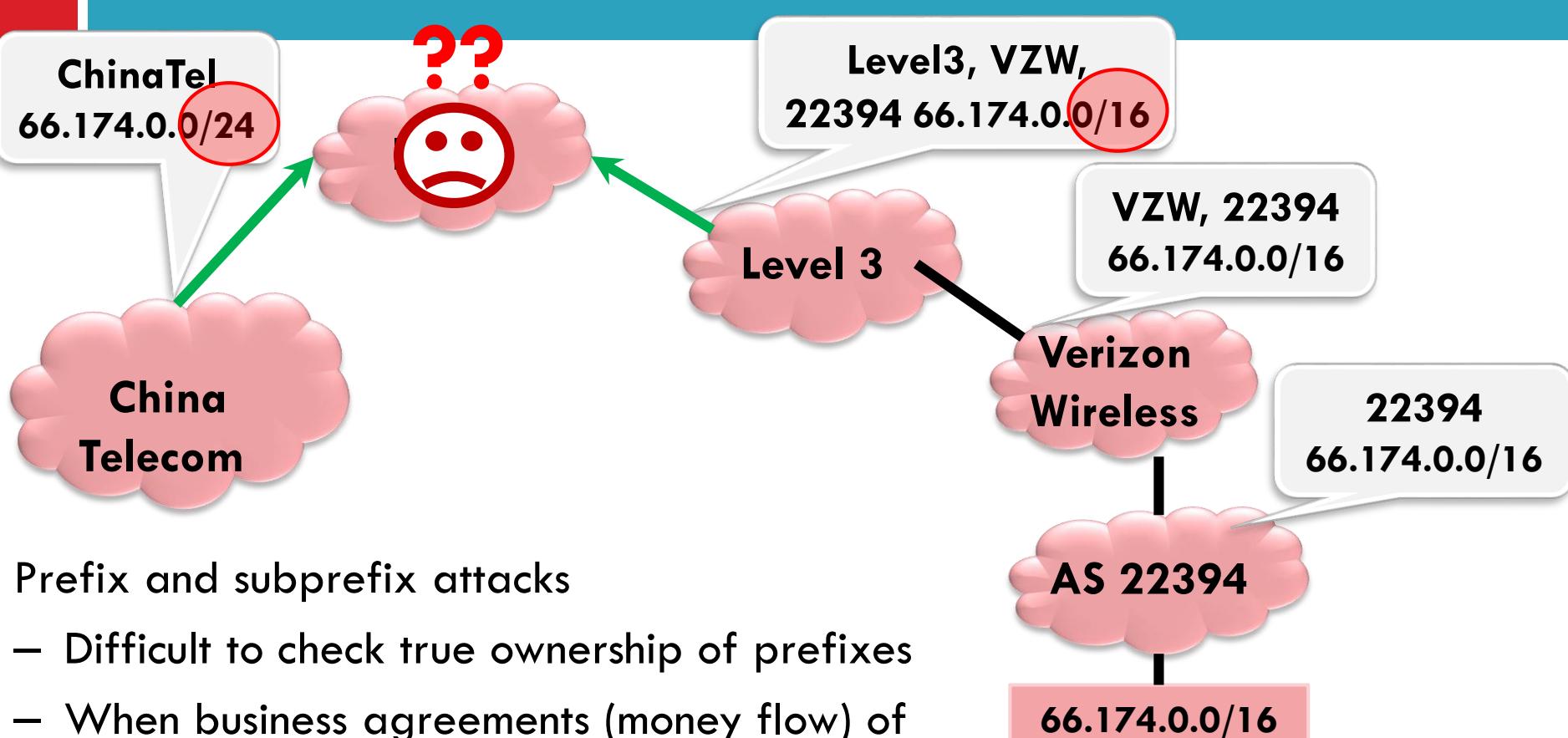
BGP-related Hijacks



BGP-related Hijacks



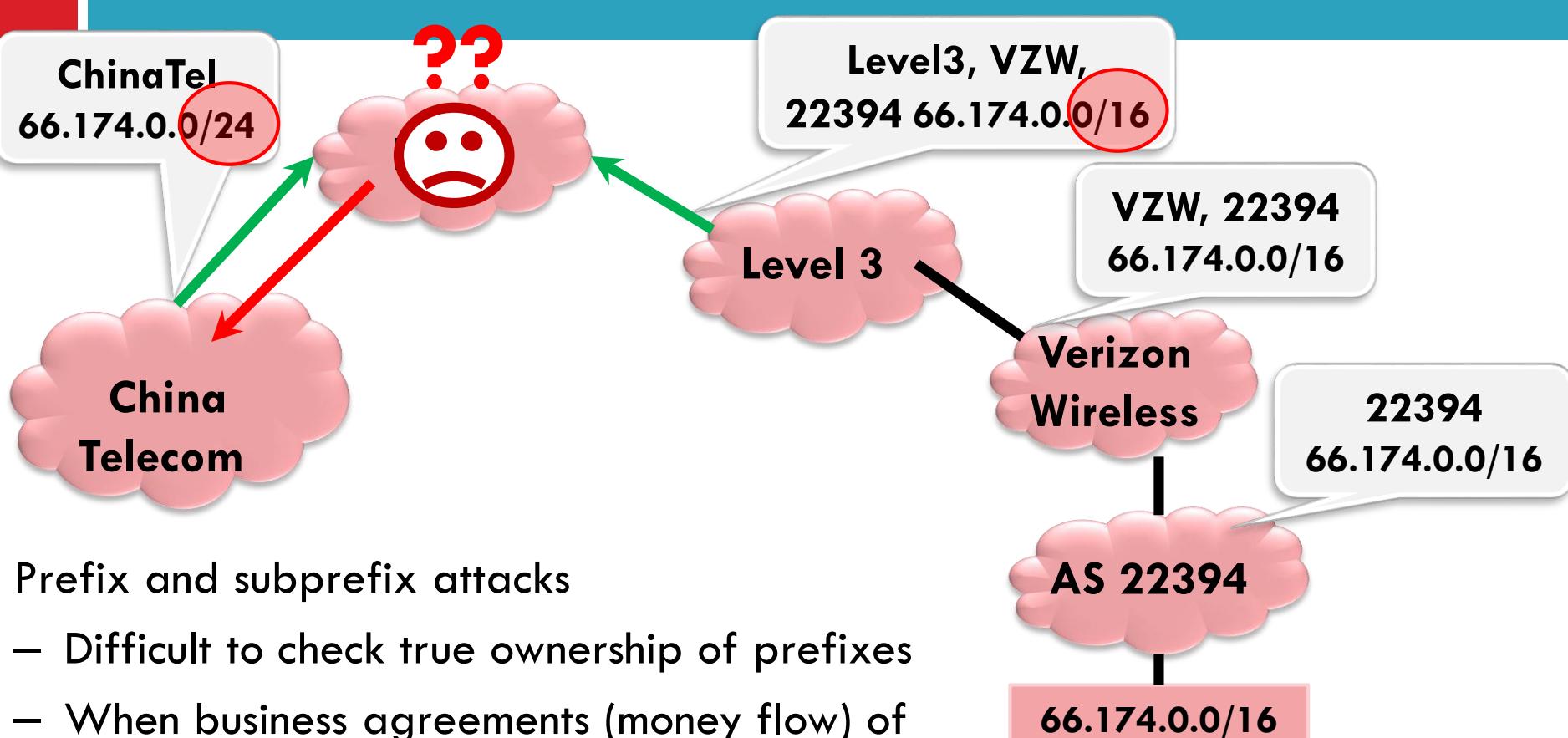
BGP-related Hijacks



Prefix and subprefix attacks

- Difficult to check true ownership of prefixes
- When business agreements (money flow) of same type, typically pick “shorter” path
- Or more specific prefix (subprefix attack)

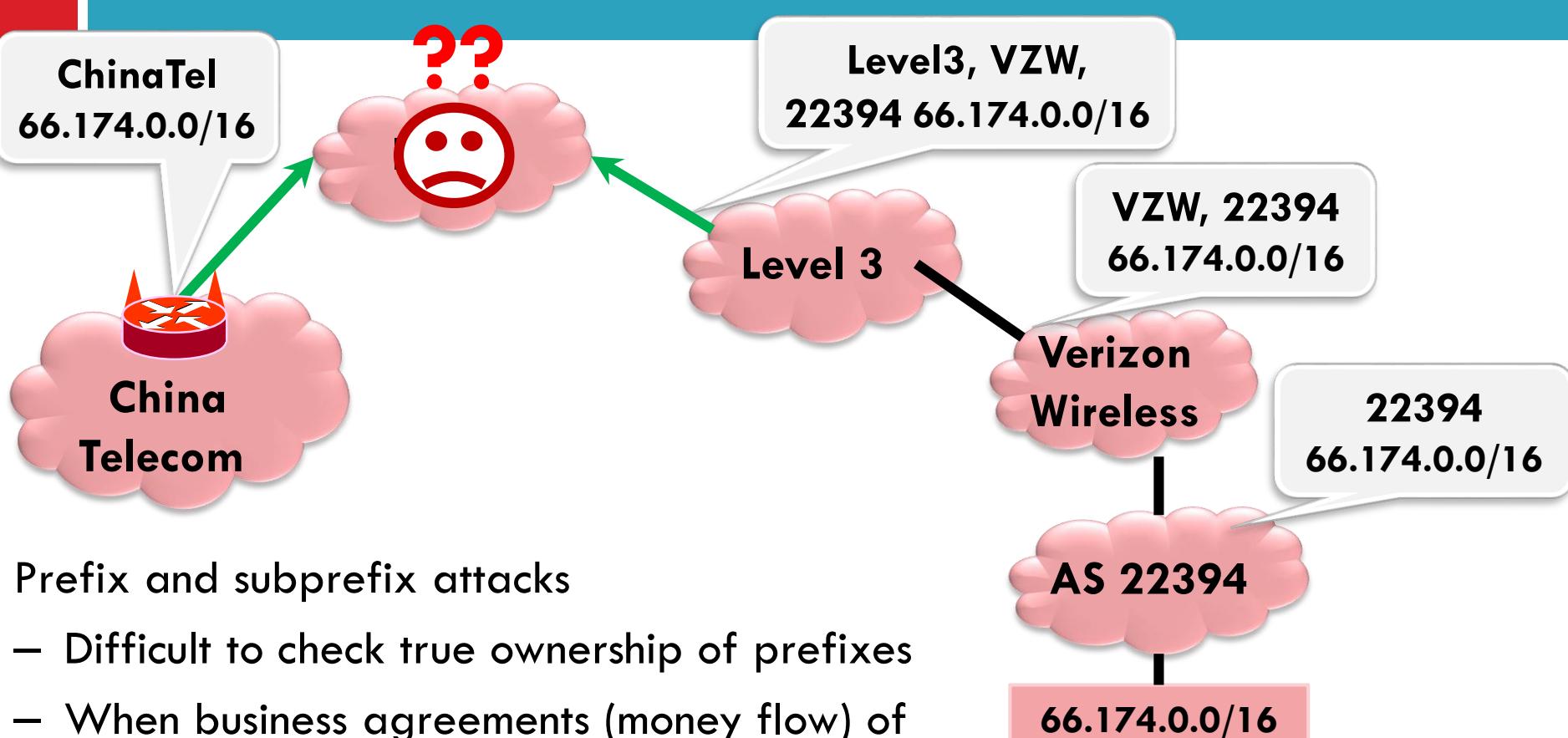
BGP-related Hijacks



Prefix and subprefix attacks

- Difficult to check true ownership of prefixes
- When business agreements (money flow) of same type, typically pick “shorter” path
- Or more specific prefix (subprefix attack)

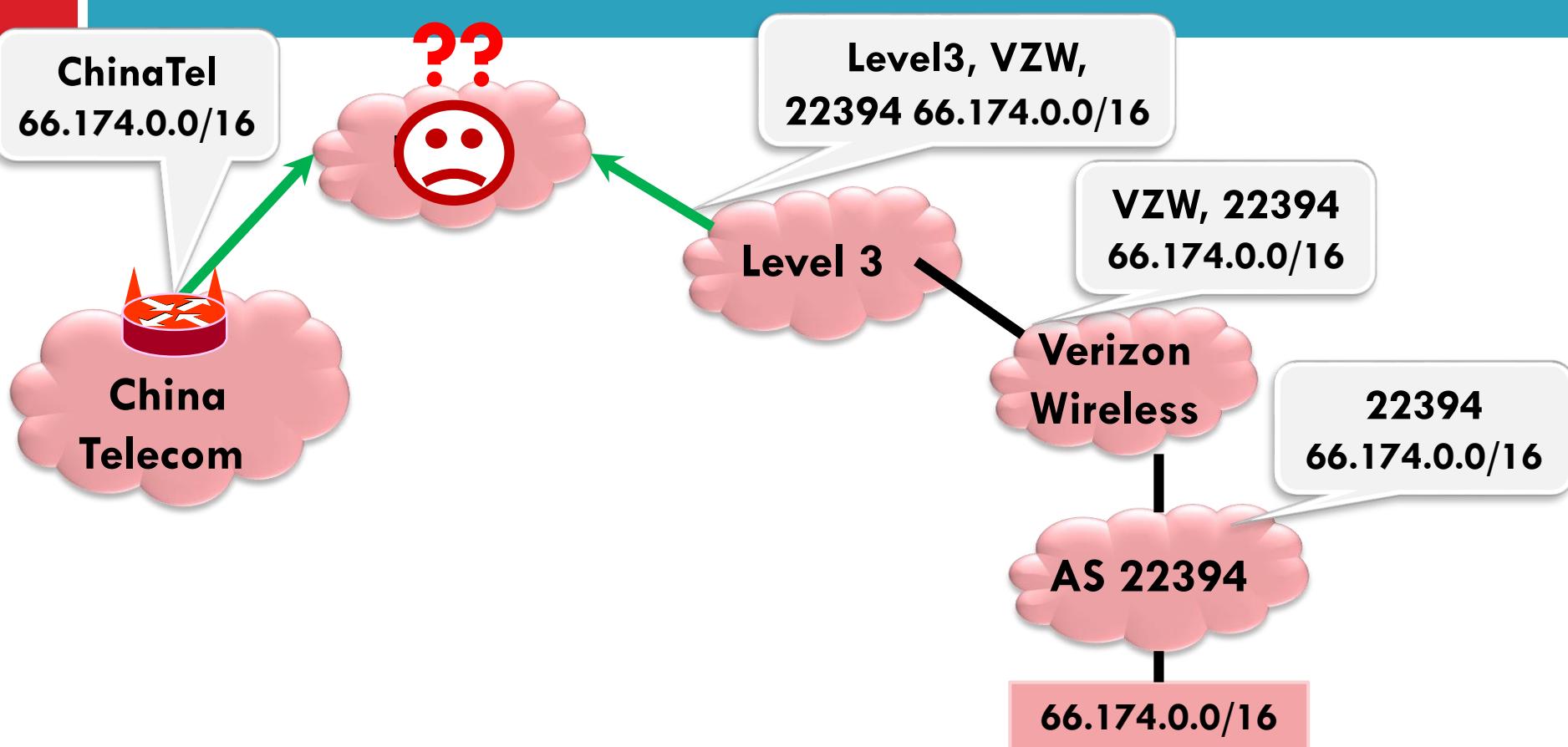
BGP-related Hijacks



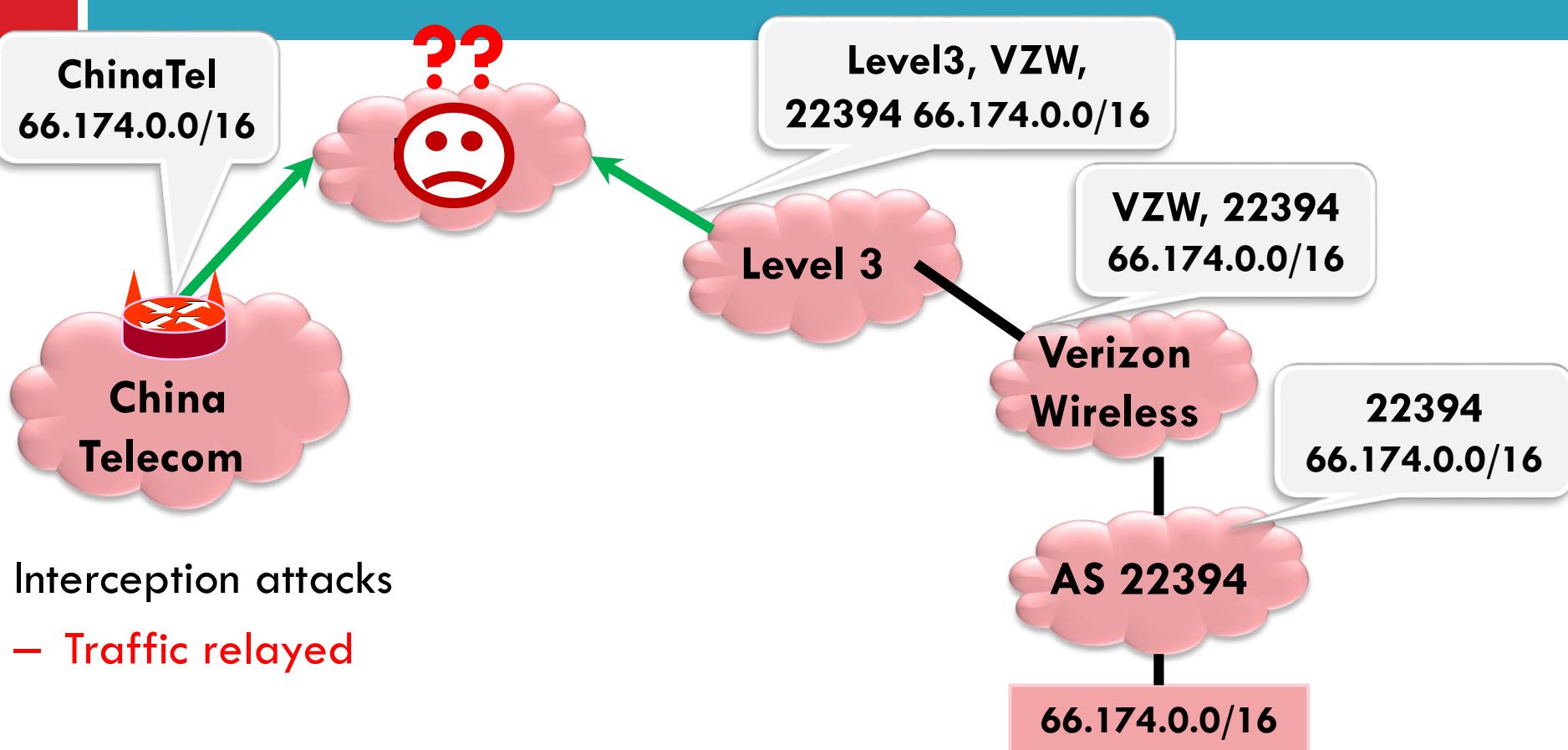
Prefix and subprefix attacks

- Difficult to check true ownership of prefixes
- When business agreements (money flow) of same type, typically pick “shorter” path
- Or more specific prefix (subprefix attack)
- Apr. 2010: ChinaTel announces 50K prefixes

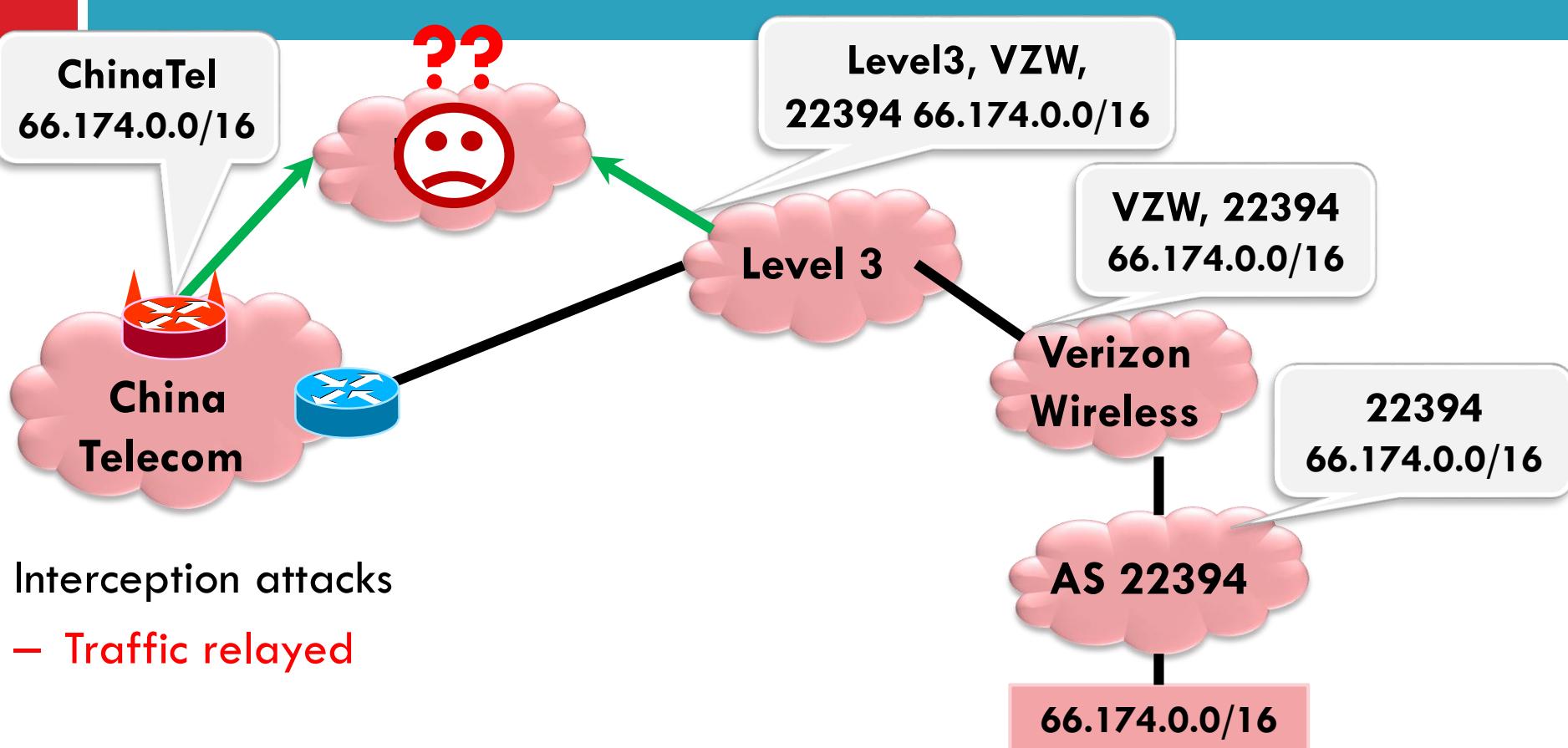
BGP-related Hijacks



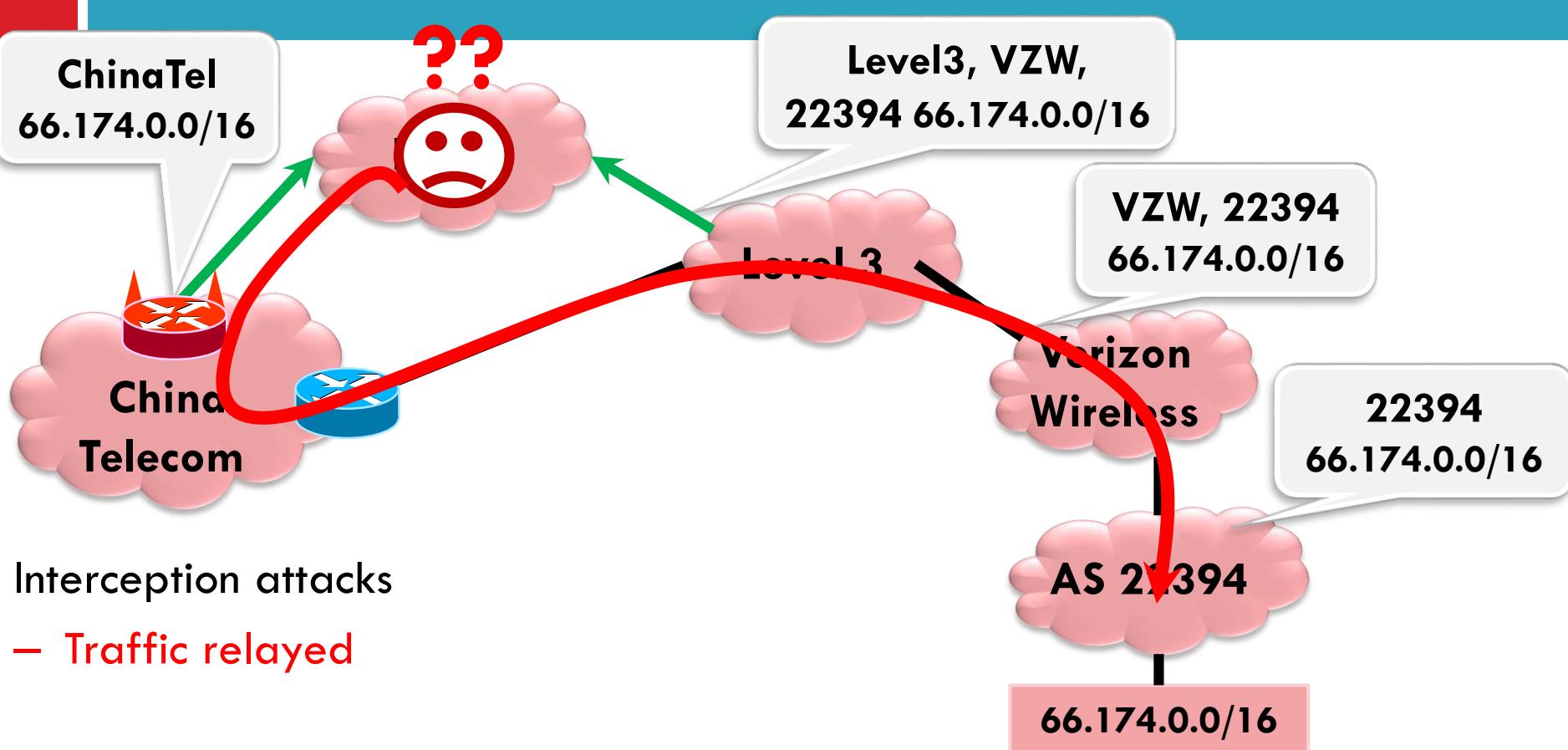
BGP-related Hijacks



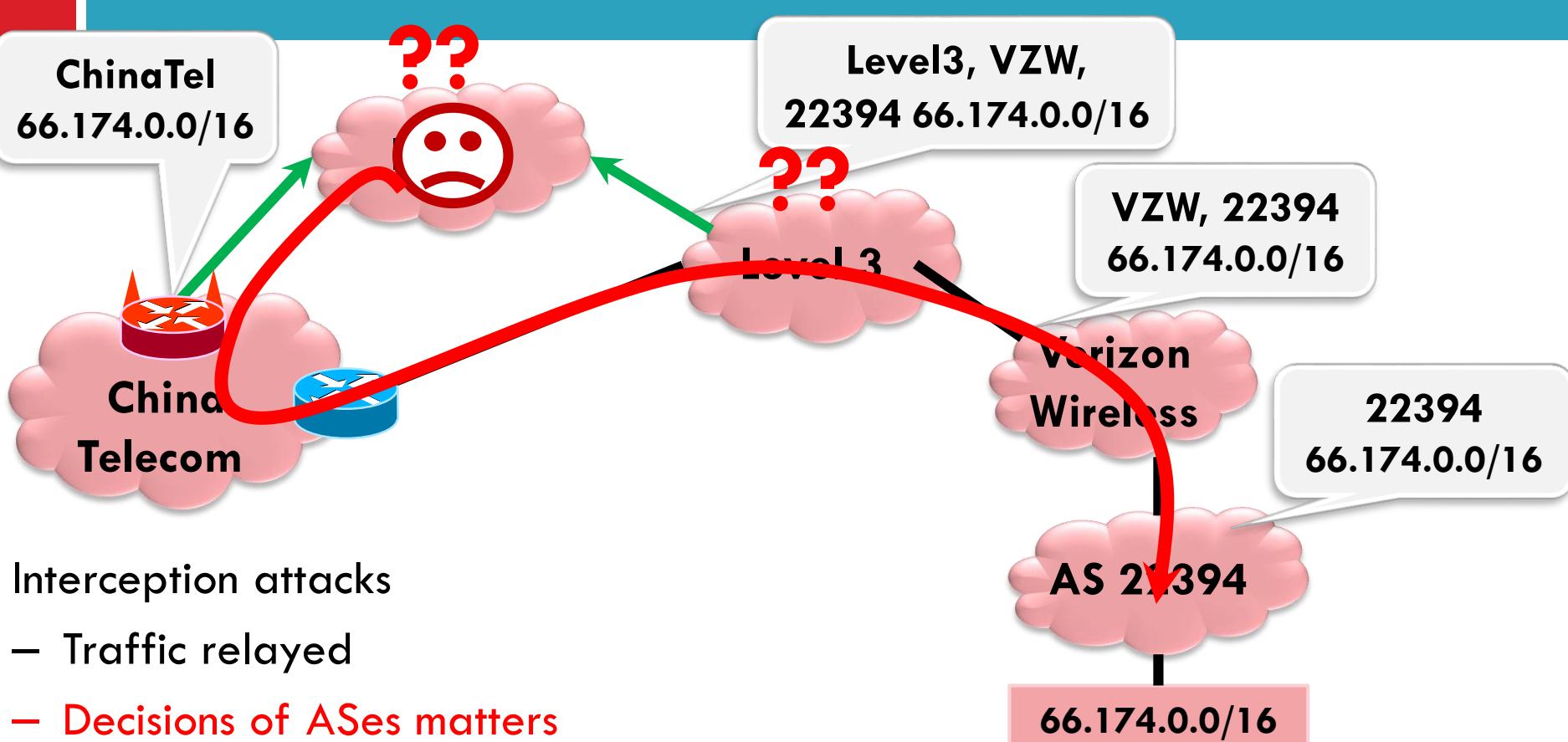
BGP-related Hijacks



BGP-related Hijacks



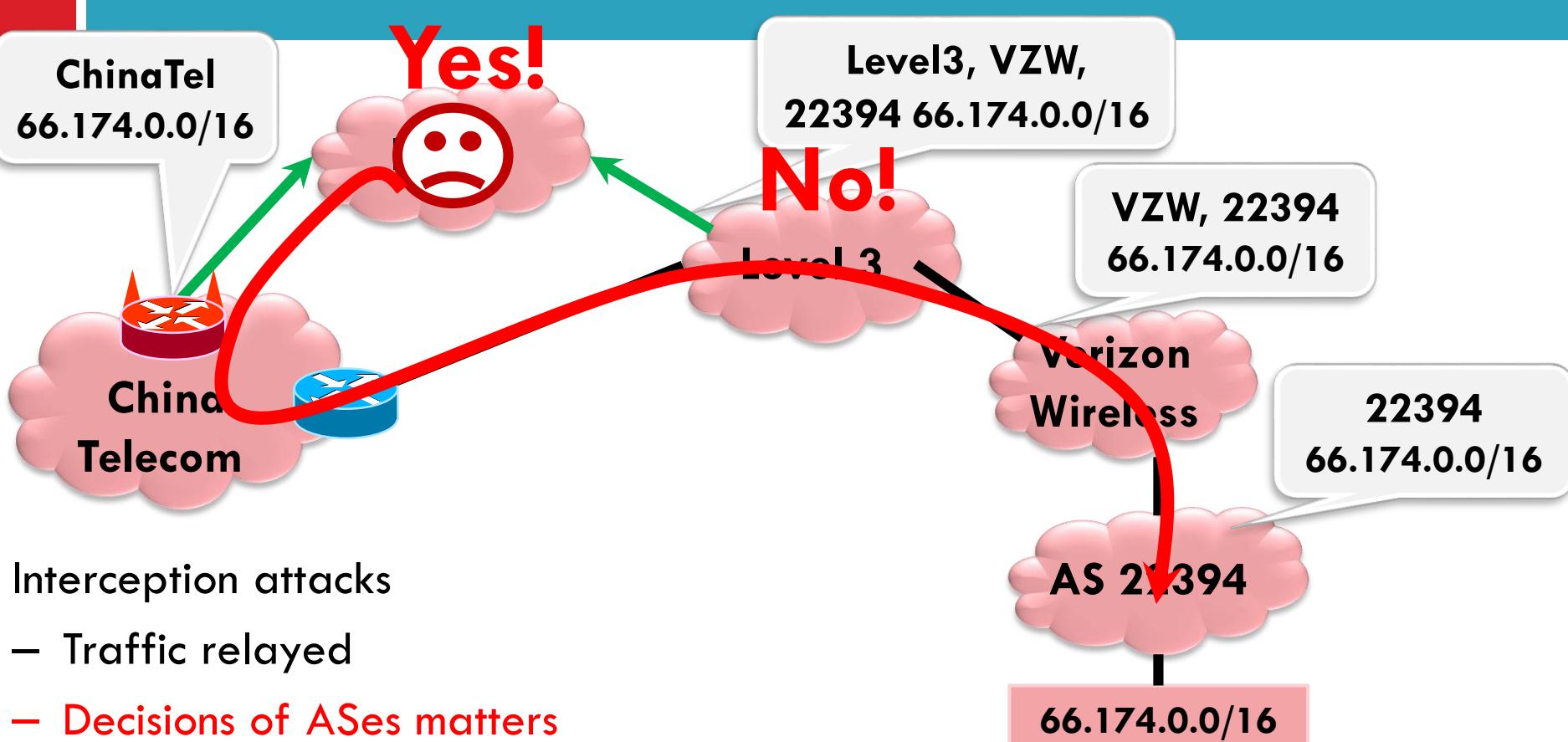
BGP-related Hijacks



Interception attacks

- Traffic relayed
- Decisions of ASes matters
 - E.g., selection of ChinaTel path
- Collaboration important

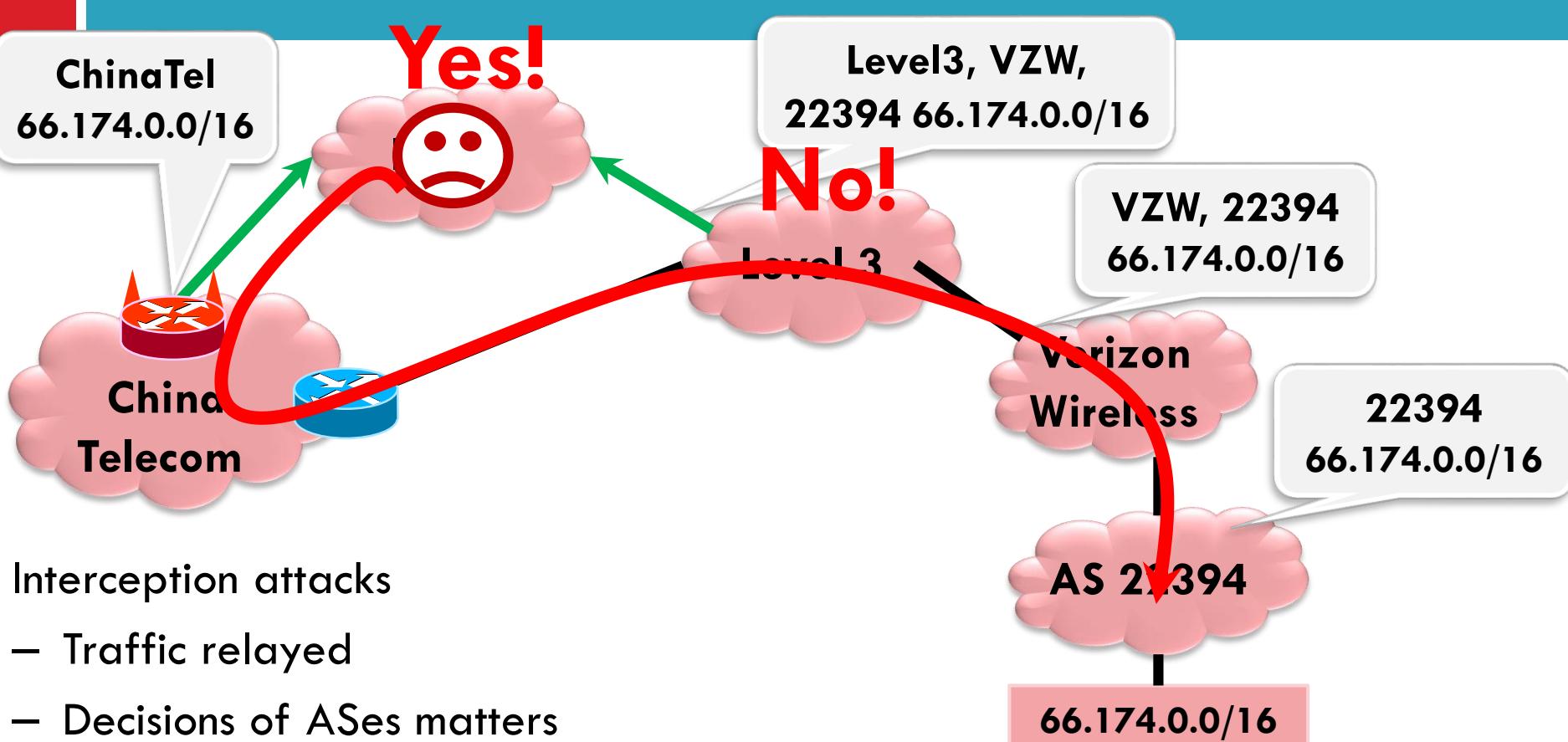
BGP-related Hijacks



Interception attacks

- Traffic relayed
- Decisions of ASes matters
 - E.g., selection of ChinaTel path
- Collaboration important

BGP-related Hijacks



Example attacks

64



Traceroute Path 1: from Guadalajara, Mexico to Washington, D.C. via Belarus

LEGEND: ● → NORMAL • → HIJACKED

START 1. Guadalajara, Mexico

1. Guadalajara, Mexico

2. Monterrey, Mexico

3. Laredo, TX

4. •

5. •

6. London, •

7. Moscow, Russia

8. Minsk, Belarus

• renesys

Internet Traffic from U.S. Government Websites Was Redirected Via Chinese Networks

By Joshua Rhett Miller / Published November 16, 2010 / FoxNews.com



- “Characterizing Large-scale Routing Anomalies: A Case Study of the China Telecom Incident”, Hiran et al., Proc. PAM 2013

Conventional Wisdom (i.e., lies)

66

- Internet is a global scale end-to-end network
 - Packets transit (mostly) unmodified
 - Value of network is global addressability /reachability
- Broad distribution of traffic sources / sinks
- An Internet “core” exists
 - Dominated by a dozen global transit providers (tier 1)
 - Interconnecting content, consumer and regional providers

Does this still hold?

67

- Emergence of ‘hyper giant’ services



- How much traffic do these services contribute?
- Hard to answer!
 - Reading on Web page: Labovitz 2010 tries to look at this.

Change in Carrier Traffic Demands

68

- In 2007 top ten match “Tier 1” ISPs
- In 2009 global transit carry significant volumes
 - But Google and Comcast join the list
 - Significant fraction of ISP A traffic is Google transit

Rank	2007 Top Ten	%
1	ISP A	5.77
2	ISP B	4.55
3	ISP C	3.35
4	ISP D	3.2
5	ISP E	2.77
6	ISP F	2.6
7	ISP G	2.24
8	ISP H	1.82
9	ISP I	1.35
10	ISP J	1.23

Rank	2009 Top Ten	%
1	ISP A	9.41
2	ISP B	5.7
3	Google	5.2
4	-	
5	-	
6	Comcast	3.12
7	-	
8	-	
9	-	
10	-	

Based on analysis of anonymous ASN (origin/transit) data (as a weighted average % of all Internet Traffic). Top ten has NO direct relationship to study participation.

Market intuition

69

- Commoditization of IP and hosting/CDN
 - ▣ Drop in price of transit
 - ▣ Drop in price of video/CDN
 - ▣ Economics of scale → Cloud computing
- Consolidation
 - ▣ Big get bigger (economics of scale)
 - ▣ Acquisitions (e.g., Google + YT)
- New economic models
 - ▣ Paid peering, paid content
- Disintermediation
 - ▣ Direct connections between content + consumer
 - ▣ Cost + performance considerations

New applications + ways to access them

70

Rank	Upstream		Downstream		Aggregate	
	Application	Share	Application	Share	Application	Share
1	BitTorr					
2	HTTP					
3	SSL					
4	Netflix					
5	YouTub					
6	Skype					
7	Facebo					
8	FaceTi					
9	Dropbox					
10	iTunes					

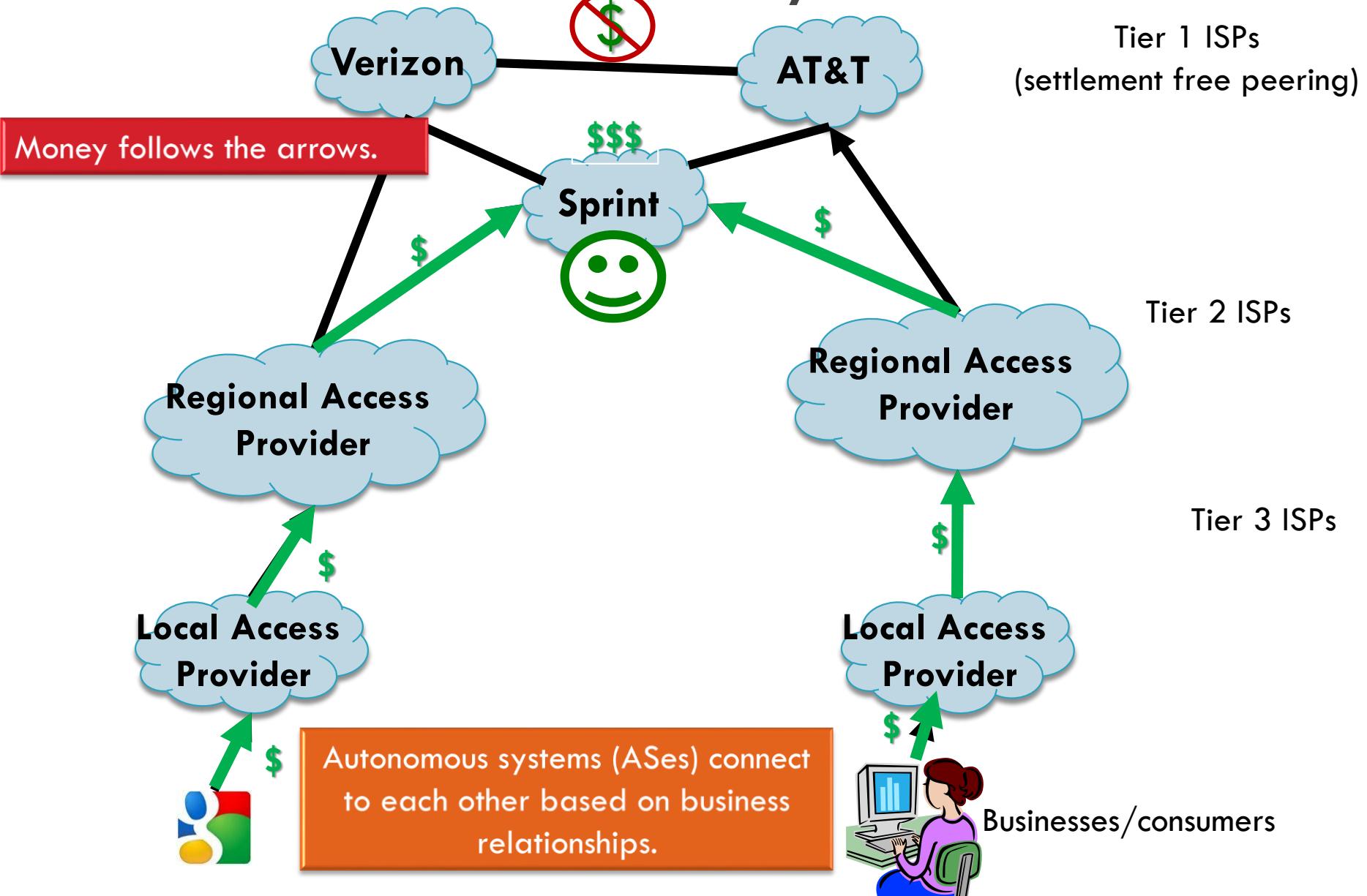
Rank	Upstream		Downstream		Aggregate	
	Application	Share	Application	Share	Application	Share
1	Facebook	26.95%	YouTube	17.61%	YouTube	17.26%
2	SSL	12.49%	Facebook	14.03%	Facebook	14.76%
3	HTTP	11.80%	HTTP	12.70%	HTTP	12.59%
4	YouTube	3.77%	MPEG	8.64%	MPEG	7.77%
5	Instagram	3.47%	SSL	6.52%	SSL	7.25%
6	BitTorrent	2.09%	Google Market	5.27%	Google Market	4.78%
7	MPEG	1.70%	Pandora Radio	5.15%	Pandora Radio	4.72%
8	Pandora Radio	1.61%	Netflix	5.05%	Netflix	4.55%
9	Gmail	1.61%	Instagram	3.49%	Instagram	3.49%
10	iCloud	1.56%	iTunes	3.10%	iTunes	2.84%
		65.50%		78.46%		77.17%



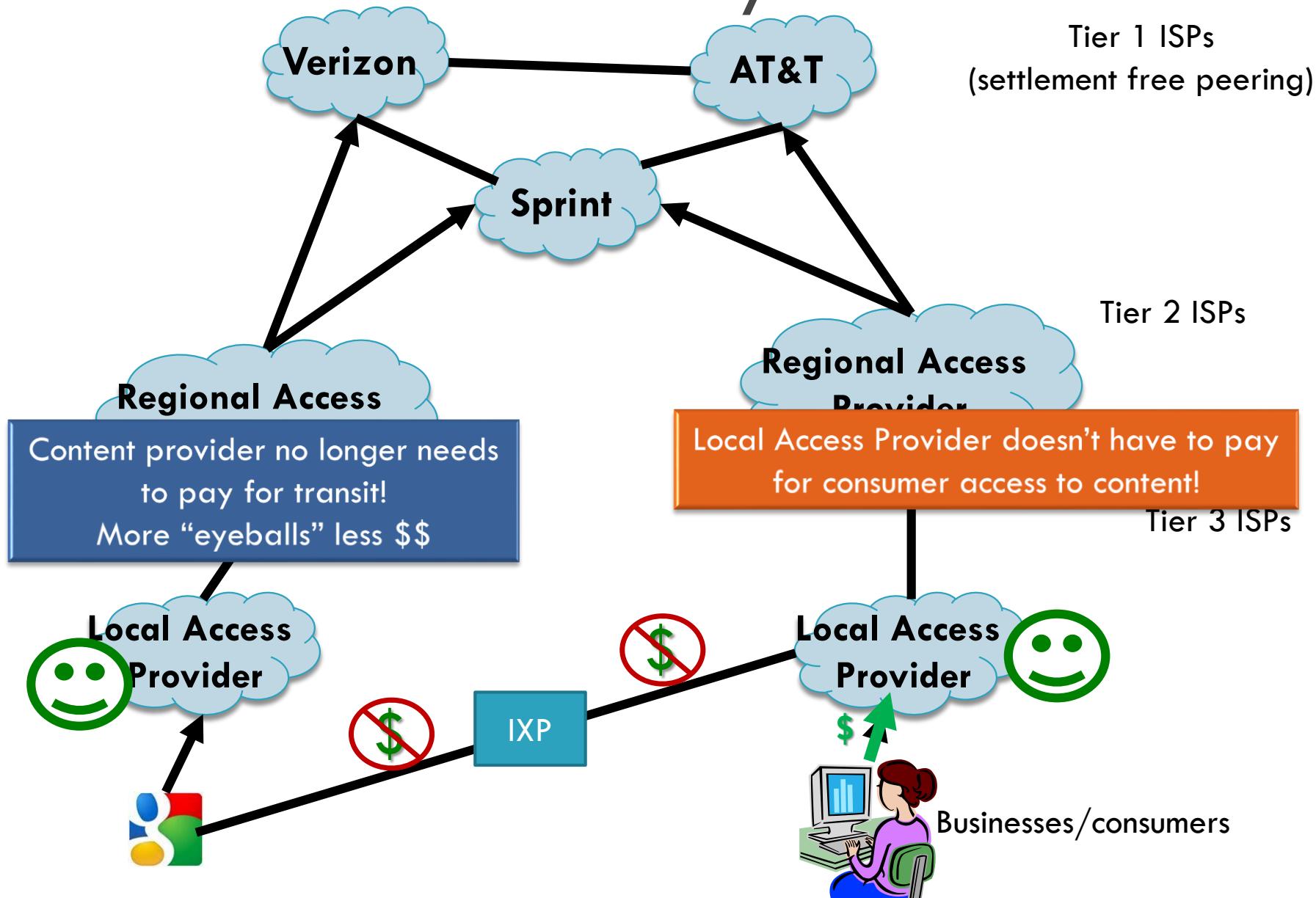
Table 4 - Top 10 Peak Period Applications - North America, Mobile Access

Fixed vs. Mobile Usage

The shift from hierarchy to flat

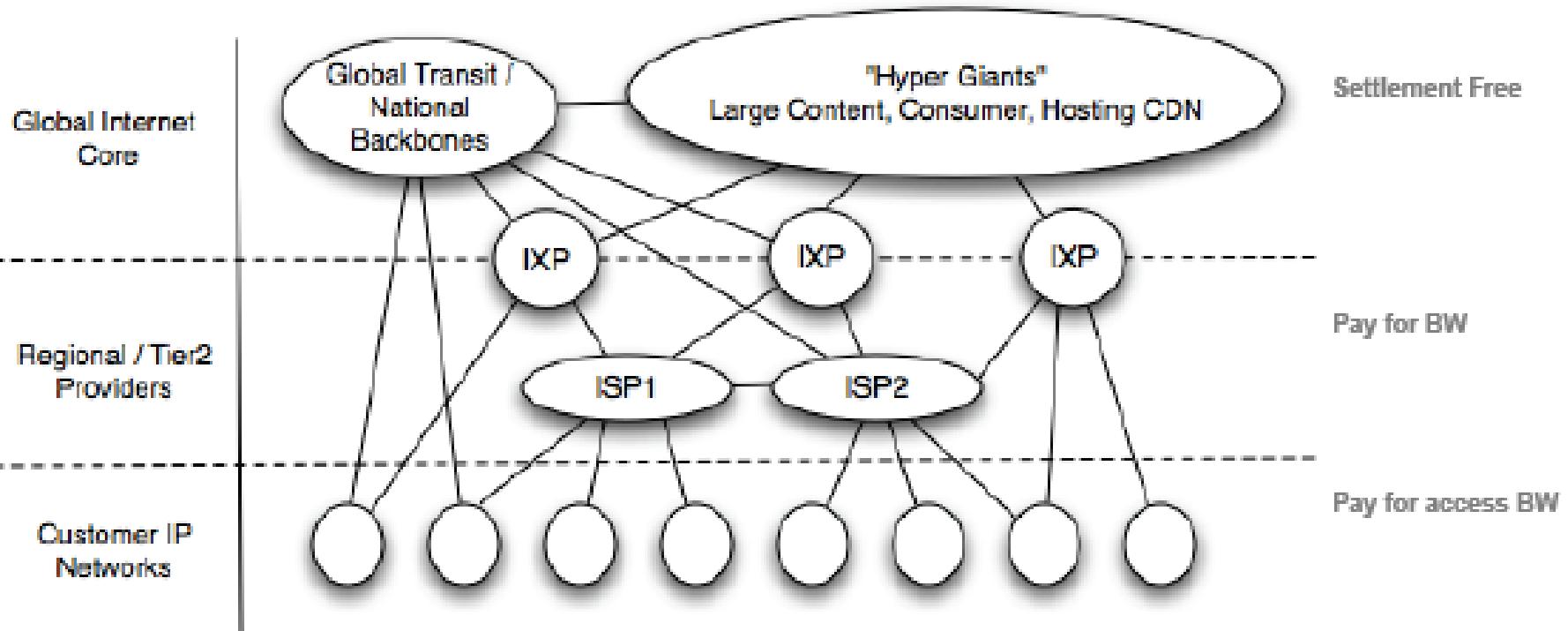


The shift from hierarchy to flat



A new Internet model

73

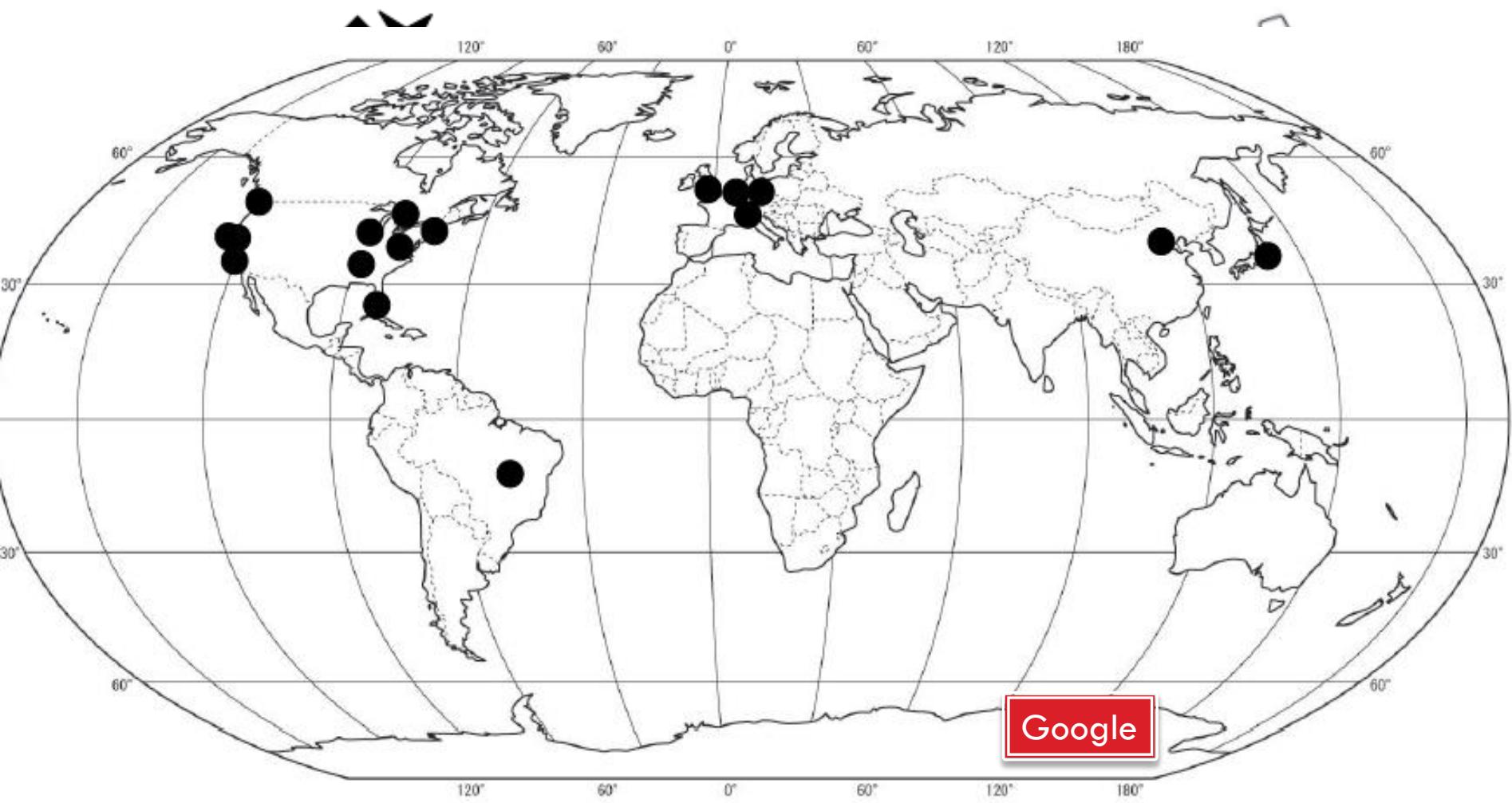


- Flatter and much more densely interconnected Internet
- Disintermediation between content and "eyeball" networks
- New commercial models between content, consumer and transit

An initial map of connectivity

The Flattening Internet Topology: Natural Evolution, Unsightly Barnacles or Contrived Collapse?, Proc. PAM 2008

74



How do ASes connect?

76

- Point of Presence (PoP)
 - Usually a room or a building (windowless)
 - One router from one AS is physically connected to the other
 - Often in big cities
 - Establishing a new connection at PoPs can be expensive
- Internet eXchange Points
 - Facilities dedicated to providing presence and connectivity for large numbers of ASes
 - Many fewer IXPs than PoPs
 - Economies of scale

IIXPs Definition

77

Industry definition (according to Euro-IX)

A physical network infrastructure operated by a single entity with the purpose to **facilitate** the **exchange** of Internet traffic between **Autonomous Systems**

The number of Autonomous Systems connected should be at least three and there **must** be a **clear** and **open policy** for others to **join**.

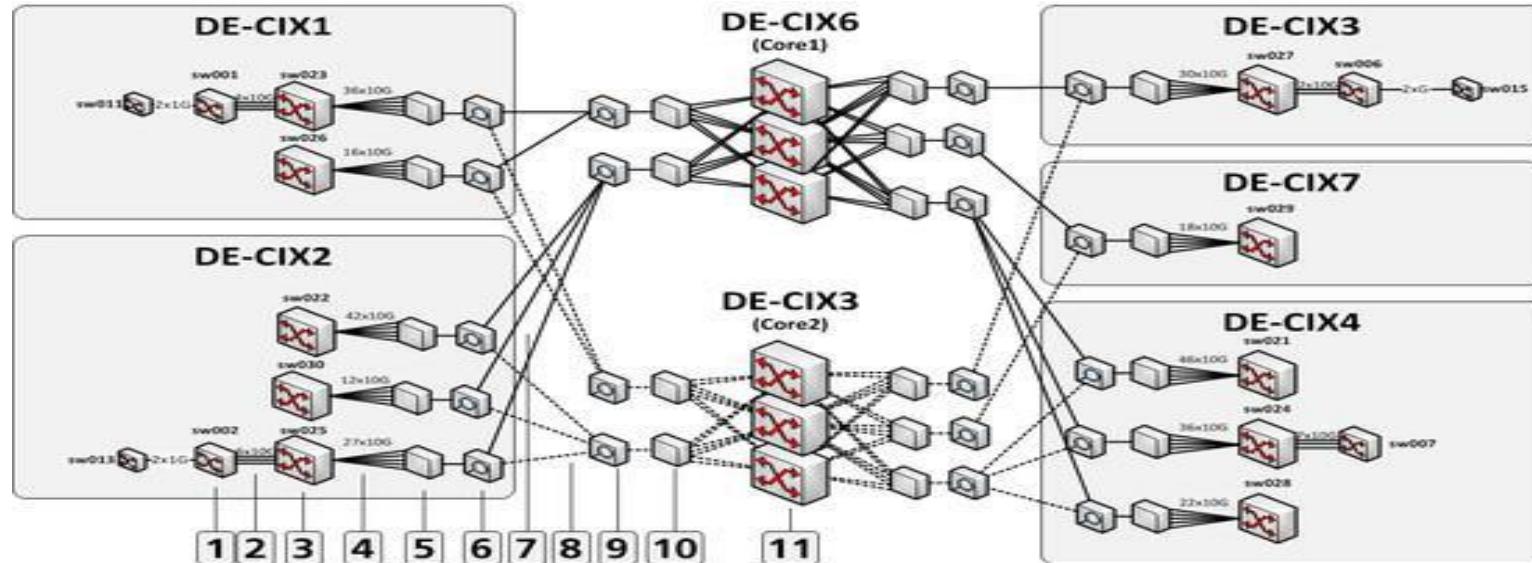
Internet eXchange Points

78



Inside an IXP

79



Robust infrastructure
with redundancy

IXPs worldwide

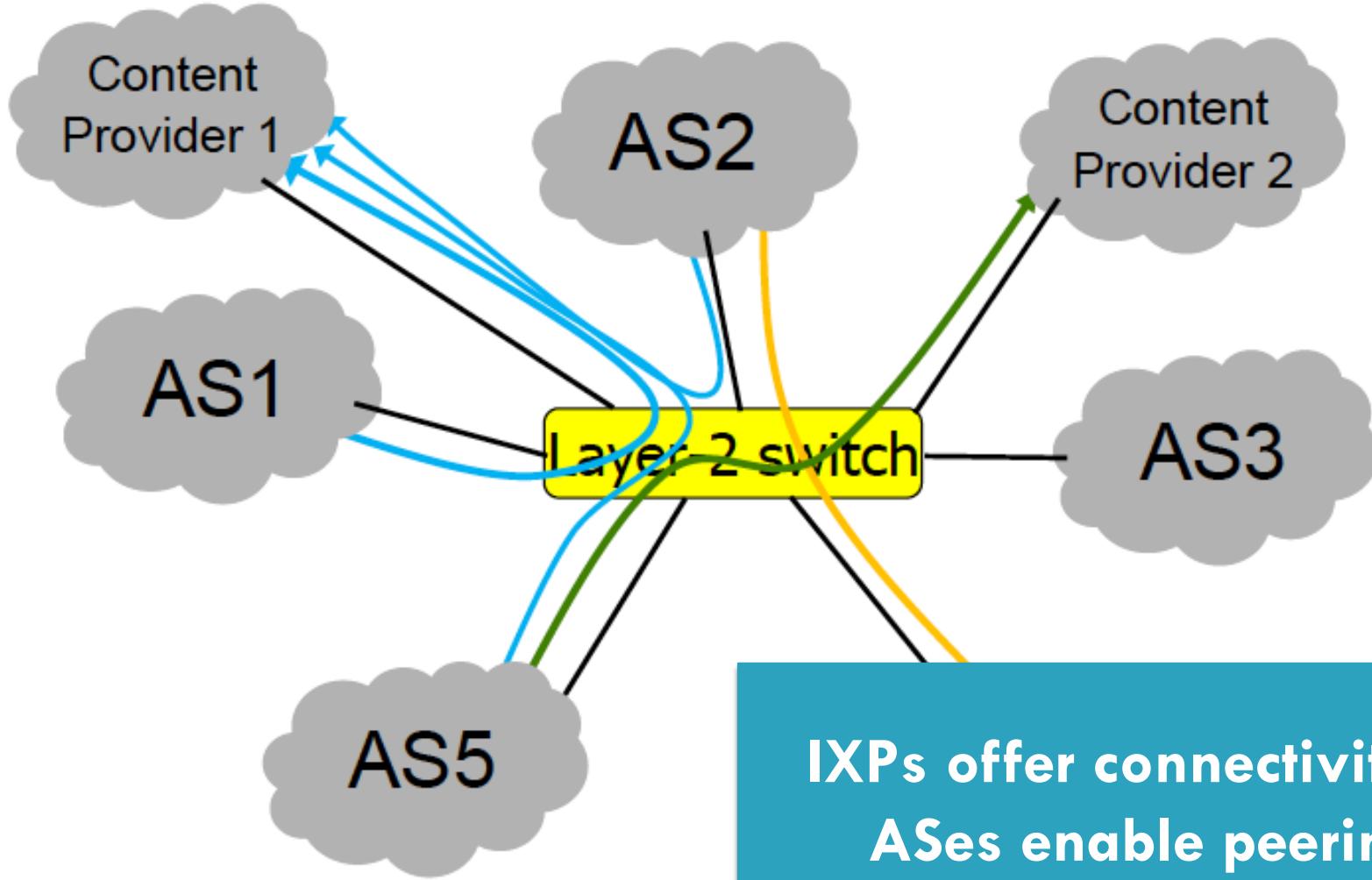
80

<https://prefix.pch.net/applications/ixpdir/>



Structure

81



IXPs -- Peering

82

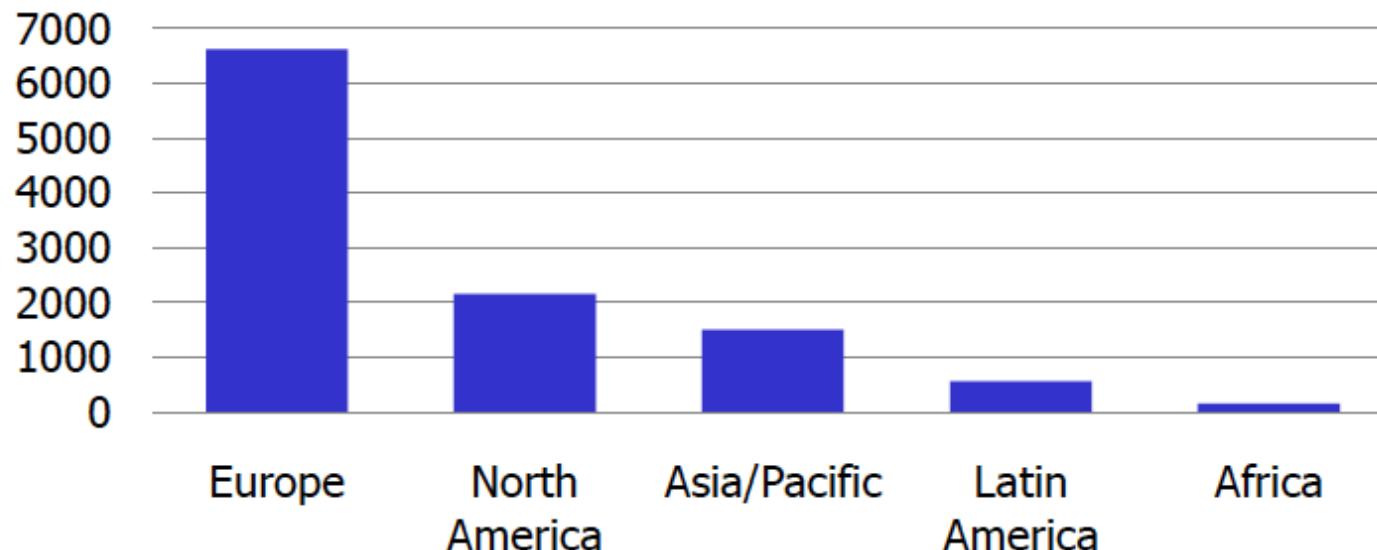
- Peering – Why? E.g., Giganews:
 - “Establishing open peering arrangements at **neutral Internet Exchange Points** is a **highly desirable** practice because the Internet Exchange members are able to **significantly improve latency, bandwidth, fault-tolerance, and the routing of traffic** between themselves at **no additional costs.**”
- IXPs – Four types of peering policies
 - Open Peering – inclination to peer with anyone, anywhere
 - Most common!
 - Selective Peering – Inclination to peer, with some conditions
 - Restrictive Peering – Inclination not to peer with any more entities
 - No Peering – No, prefer to sell transit
 - <http://drpeering.net/white-papers/Peering-Policies/Peering-Policy.html>

IXPs – Publicly available information

83

- Generally known: # IXPs ~ 350 worldwide
- Somewhat known: # ASes per IXP up to 500
- Less known: # ASes ~ 11,000 worldwide

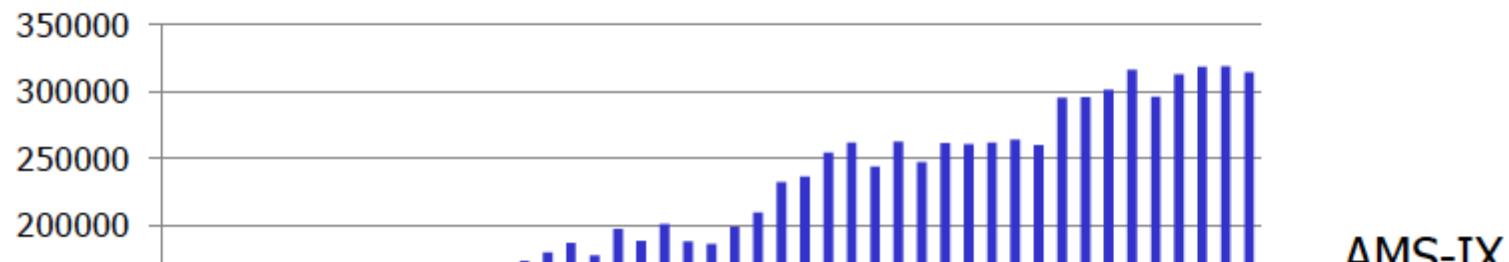
IXP Member ASes by region



IXPs- Publicly available information

84

- Generally known: # IXPs ~ 350 worldwide
- Somewhat known: # ASes per IXP up to 500
- Less known: # ASes ~ 11,000 worldwide
- Even less known: IXPs =~ Tier-1 ISP traffic



Unknown: # of peerings at IXPs

Interesting observations

85

- Myth 1: **Tier-1's don't public peer at IXPs**
 - Fact: All Tier-1's are members at IXP and do public peering
 - Tier-1's typically use a "restrictive" peering policy
 - Most IXP members use an "open" peering policy
- Myth 2: **Establishing peerings at IXPs is cumbersome**
 - Fact: Many IXPs make it very easy for its members to establish public peerings with other members
 - „Handshake agreements“
 - Use of IXP's route server is offered as free value-added service
 - Use of multi-lateral peering agreements
- Myth 3: **IXP peering links are for backup**
 - Fact: Most peering links at our IXP see traffic
 - Most of the public peering links see traffic
 - Does not include traffic on the private peering links at IXP

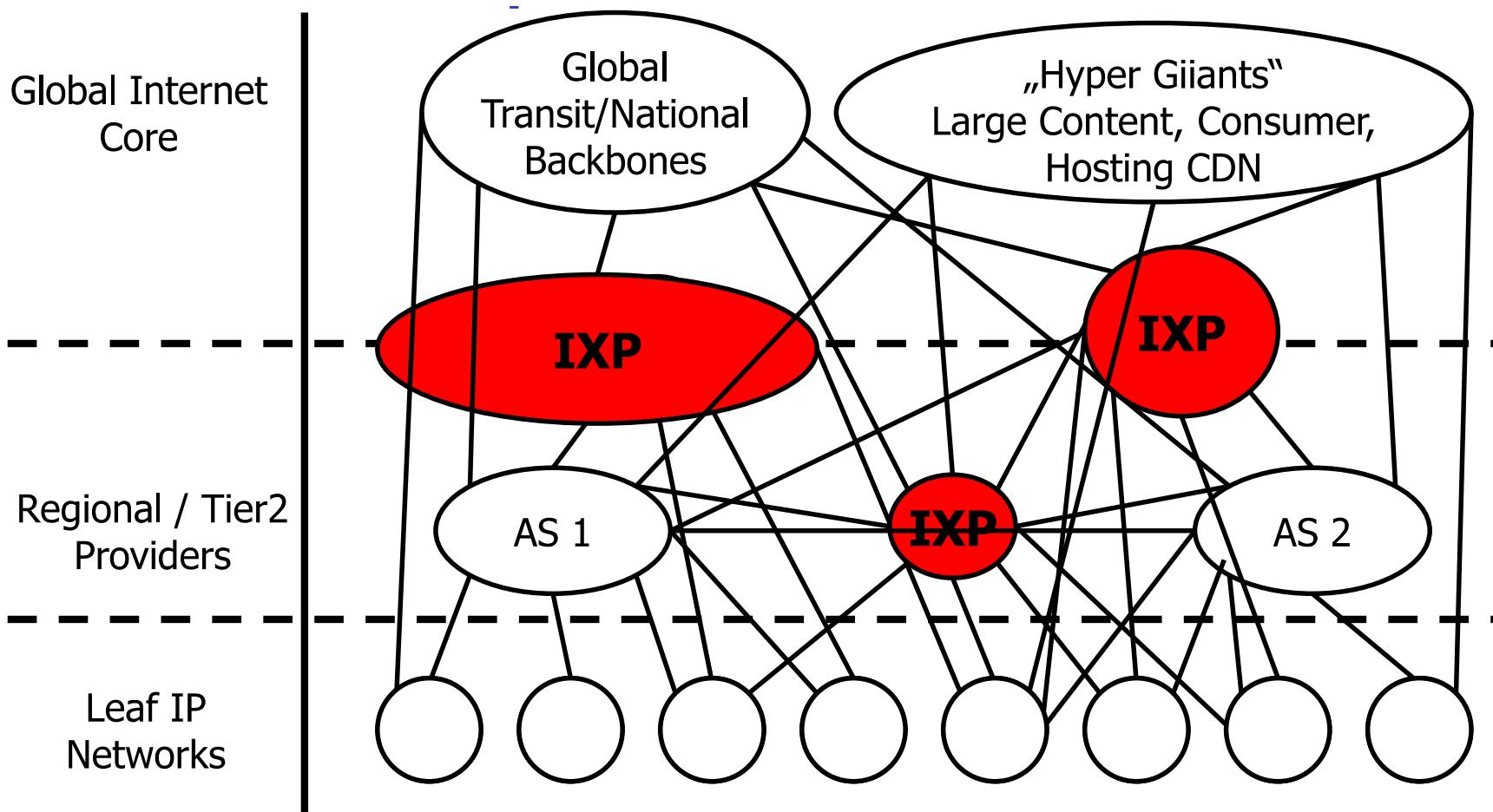
Interesting observations (2)

86

- Myth 4: IXP^s are not interesting
 - Fact: As interesting as large ASes and big content
- Myth 5: IXP^s are very different from ASes
 - Fact: Large IXPs start to look more and more like ASes
 - Offering SLAs (DE-CIX in 2008, AMS-IX in 2011)
 - Support for IXP resellers (e.g., AS43531 – IX Reach)
 - Going overseas (AMS-IX starting a site in Hong Kong)
 - Extensive monitoring capabilities
 - IXP-specific traffic matrix vs. AS-specific traffic matrix

Revised model 2012+

87



Inter-Domain Routing Summary

89

- BGP4 is the only inter-domain routing protocol currently in use world-wide
- Issues?
 - ▣ Lack of security
 - ▣ Ease of misconfiguration
 - ▣ Poorly understood interaction between local policies
 - ▣ Poor convergence
 - ▣ Lack of appropriate information hiding
 - ▣ Non-determinism
 - ▣ Poor overload behavior

Why are these still issues?

90

- Backward compatibility
- Buy-in / incentives for operators
- Stubbornness

Very similar issues to IPv6 deployment

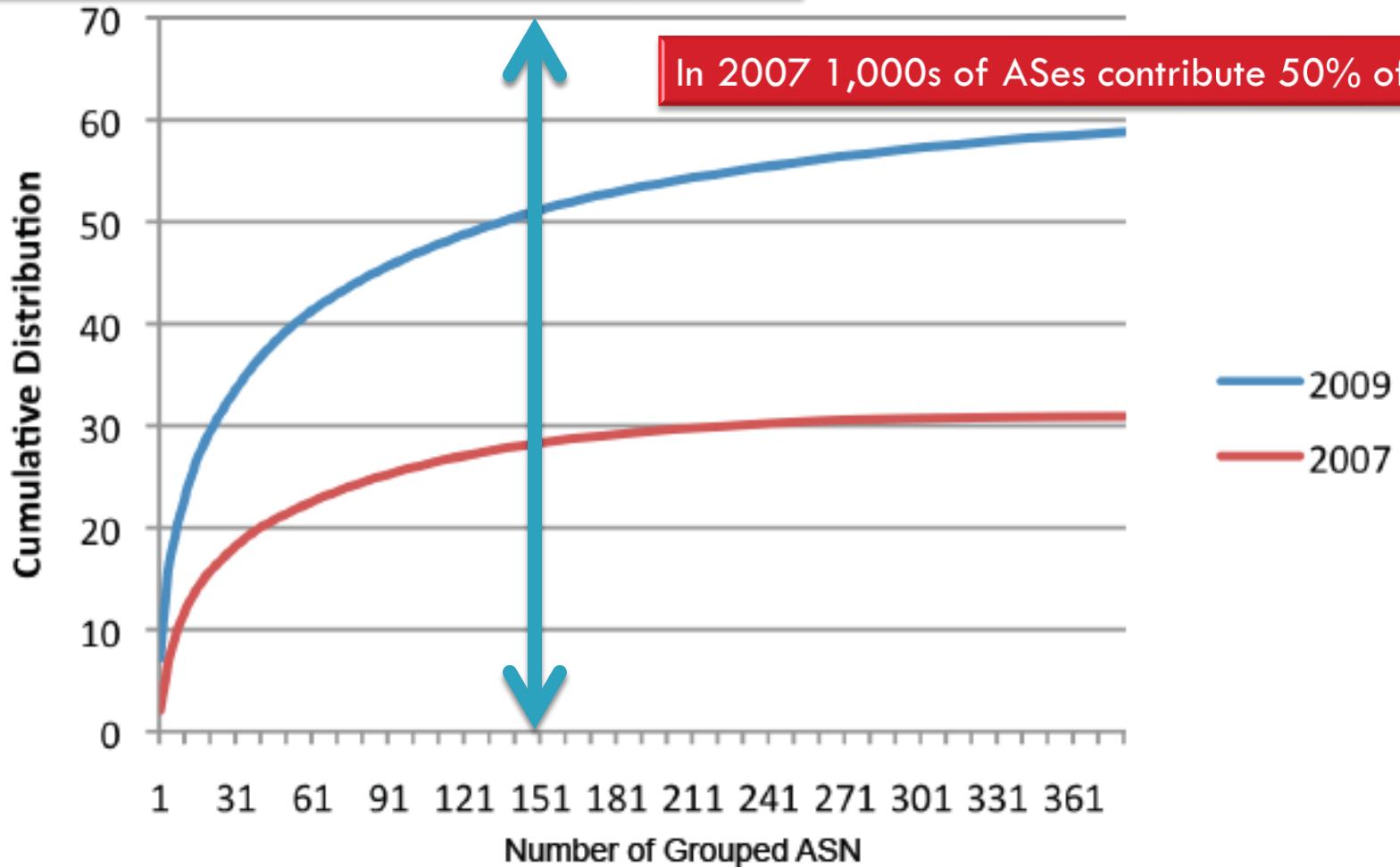
More slides ...

Consolidation of Content

93

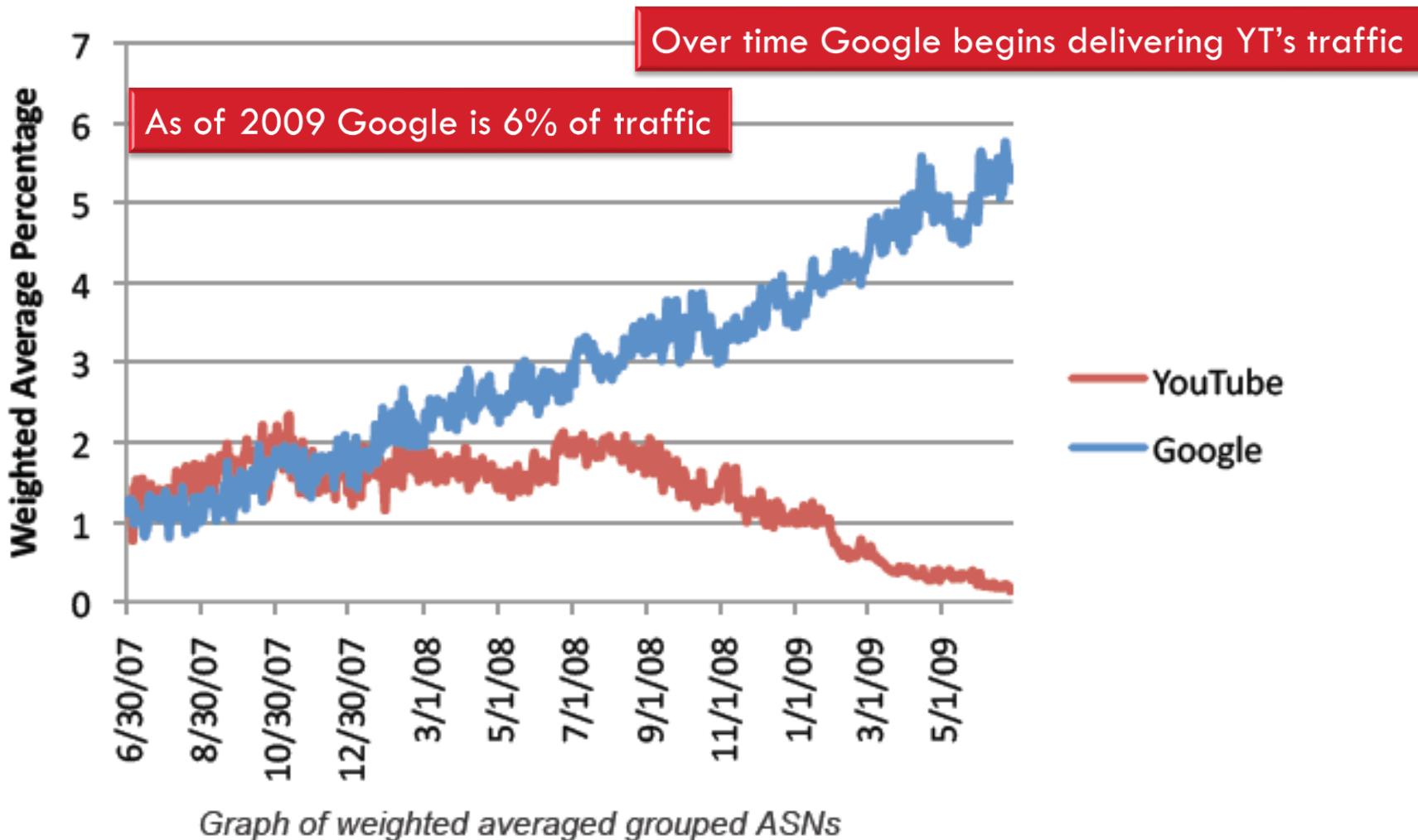
In 2009, 150 ASes contribute 50% of traffic!

In 2007 1,000s of ASes contribute 50% of traffic



Case Study: Google

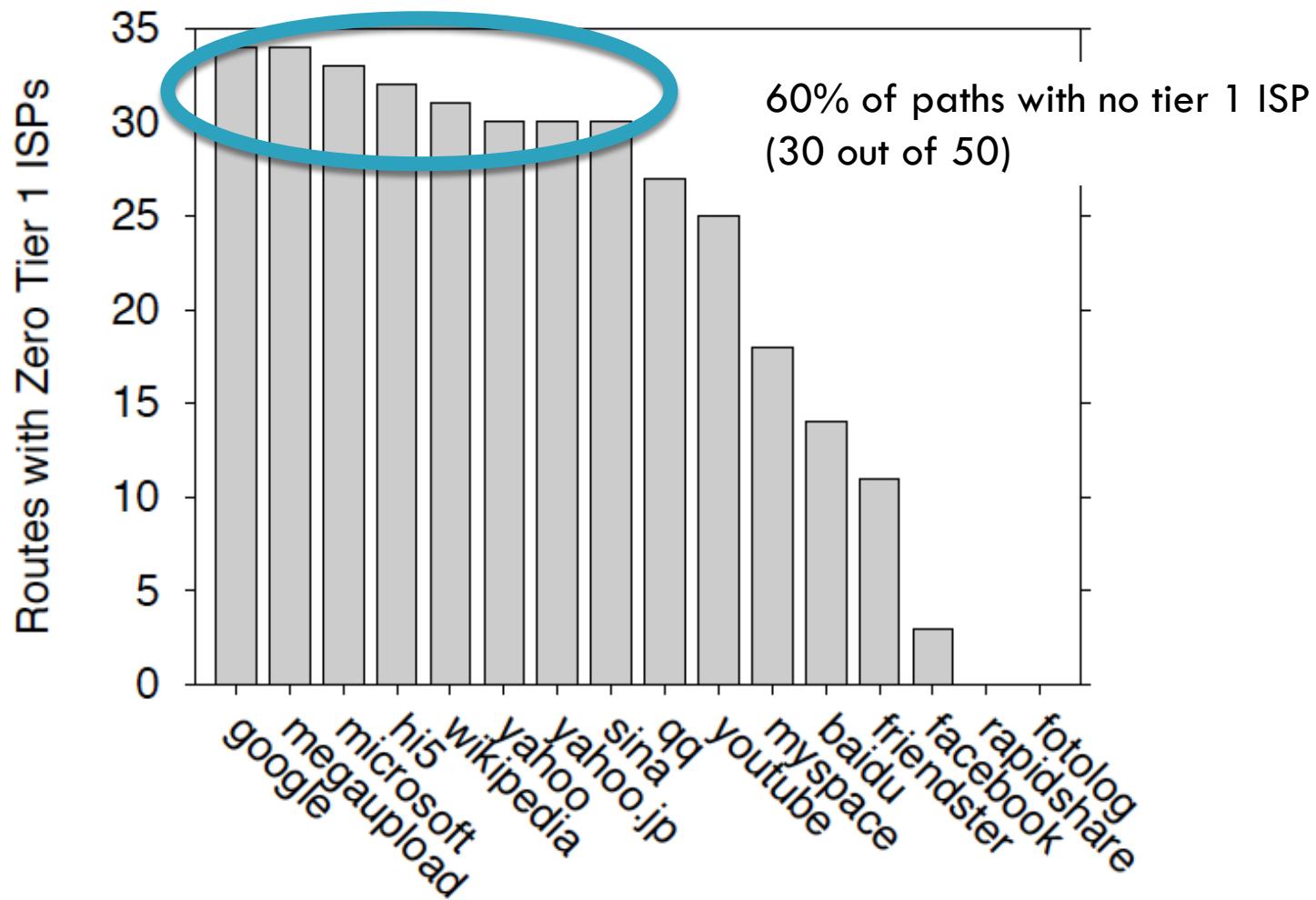
94



Flattening: Paths with no Tier 1s

The Flattening Internet Topology: Natural Evolution, Unsightly Barnacles or Contrived Collapse?, Proc. PAM 2008

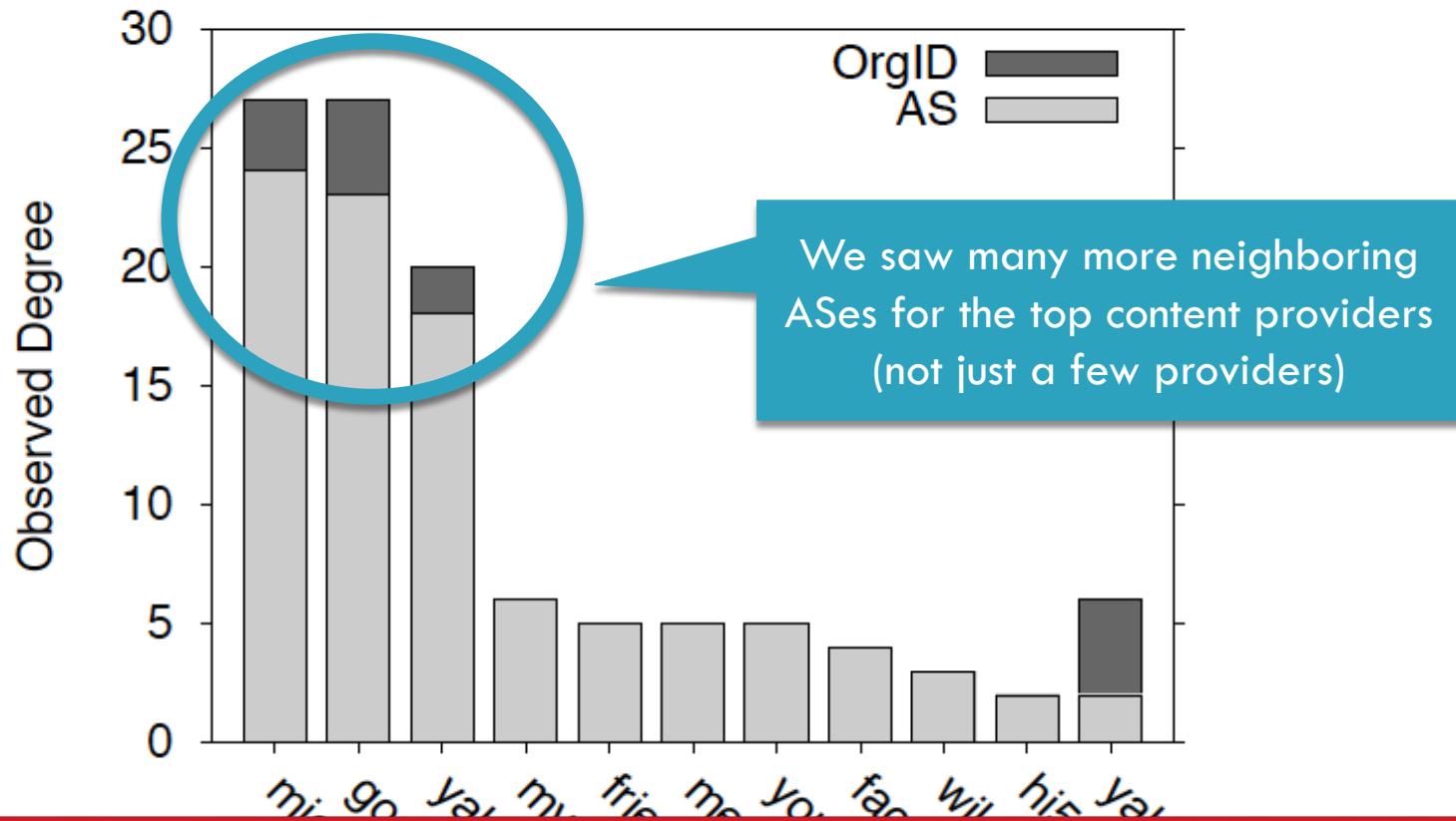
95



Relative degree of top content providers

The Flattening Internet Topology: Natural Evolution, Unsightly Barnacles or Contrived Collapse?, Proc. PAM 2008

96



These numbers are actually way lower than the true degree of these ASes

What Problem is BGP Solving?

98

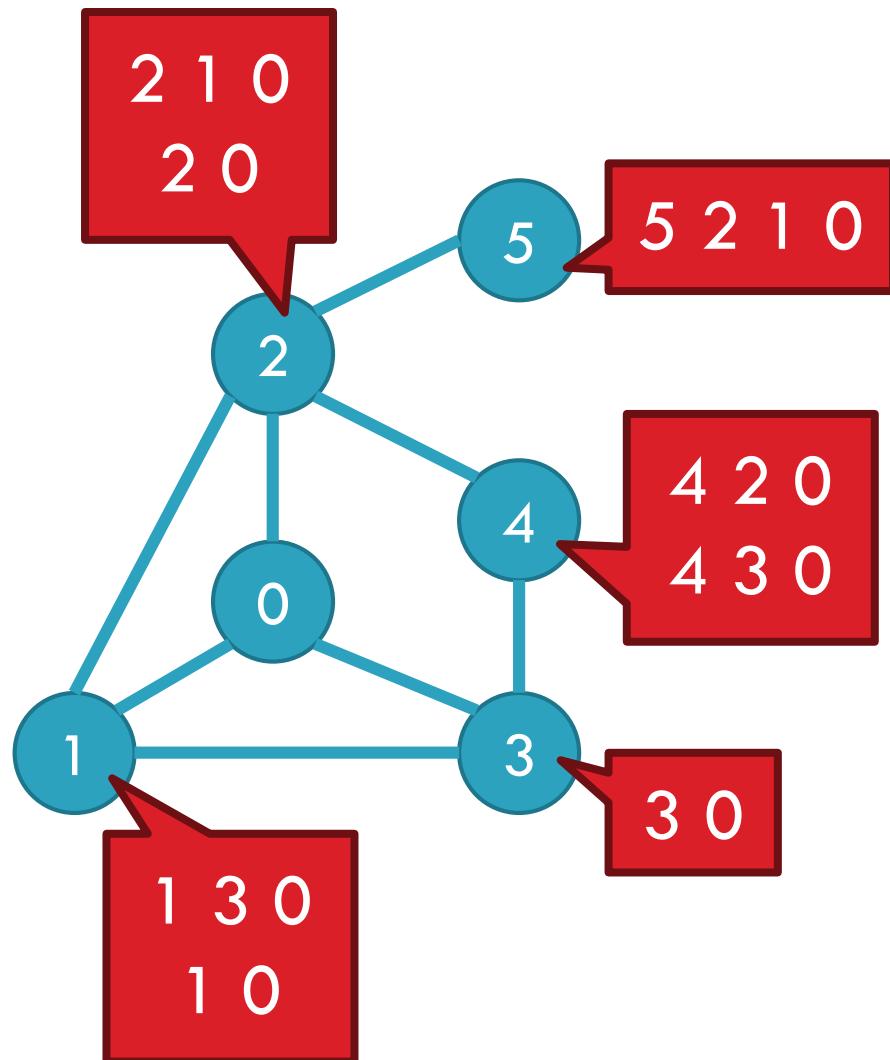
Underlying Problem	Distributed Solution
Shortest Paths	RIP, OSPF, IS-IS, etc.
???	BGP

- Knowing ??? can:
 - Aid in the analysis of BGP policy
 - Aid in the design of BGP extensions
 - Help explain BGP routing anomalies
 - Give us a deeper understanding of the protocol

The Stable Paths Problem

99

- An instance of the SPP:
 - Graph of nodes and edges
 - Node 0, called the origin
 - A set of permitted paths from each node to the origin
 - Each set of paths is ranked



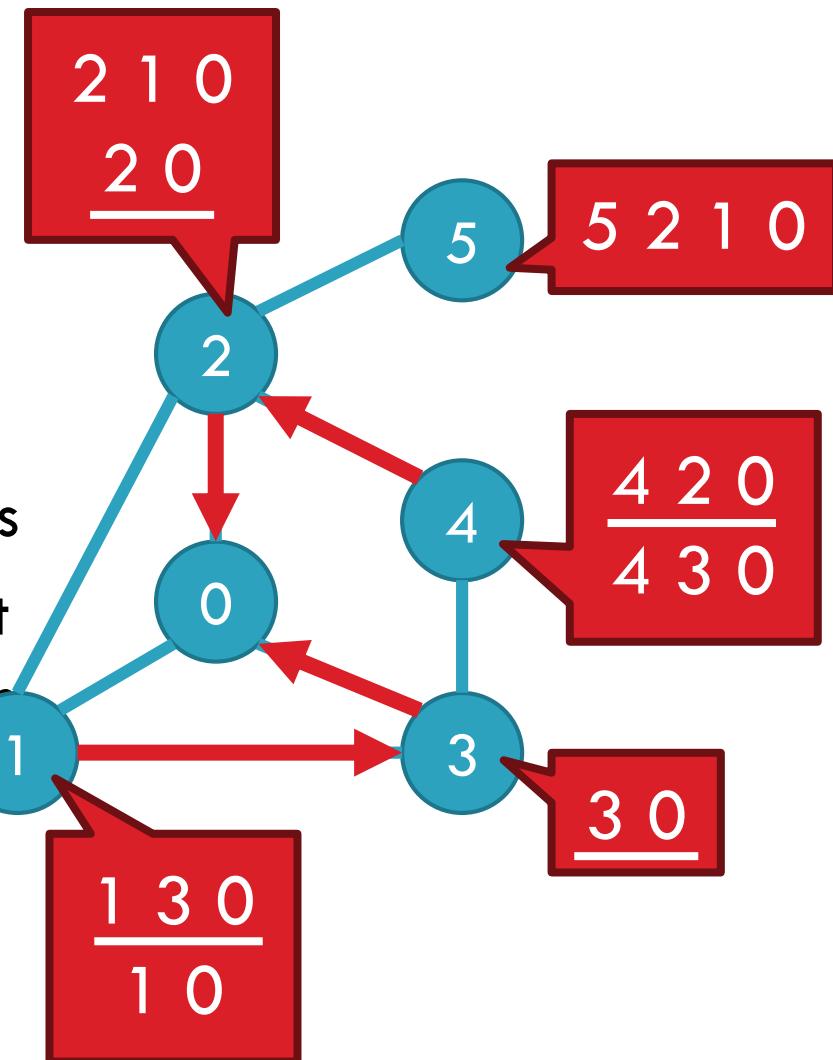
A Solution to the SPP

100

- A solution is an assignment of permitted paths to each node

Solutions need not use
the shortest paths, or
form a spanning tree

their neighbors



Simple SPP Example

101

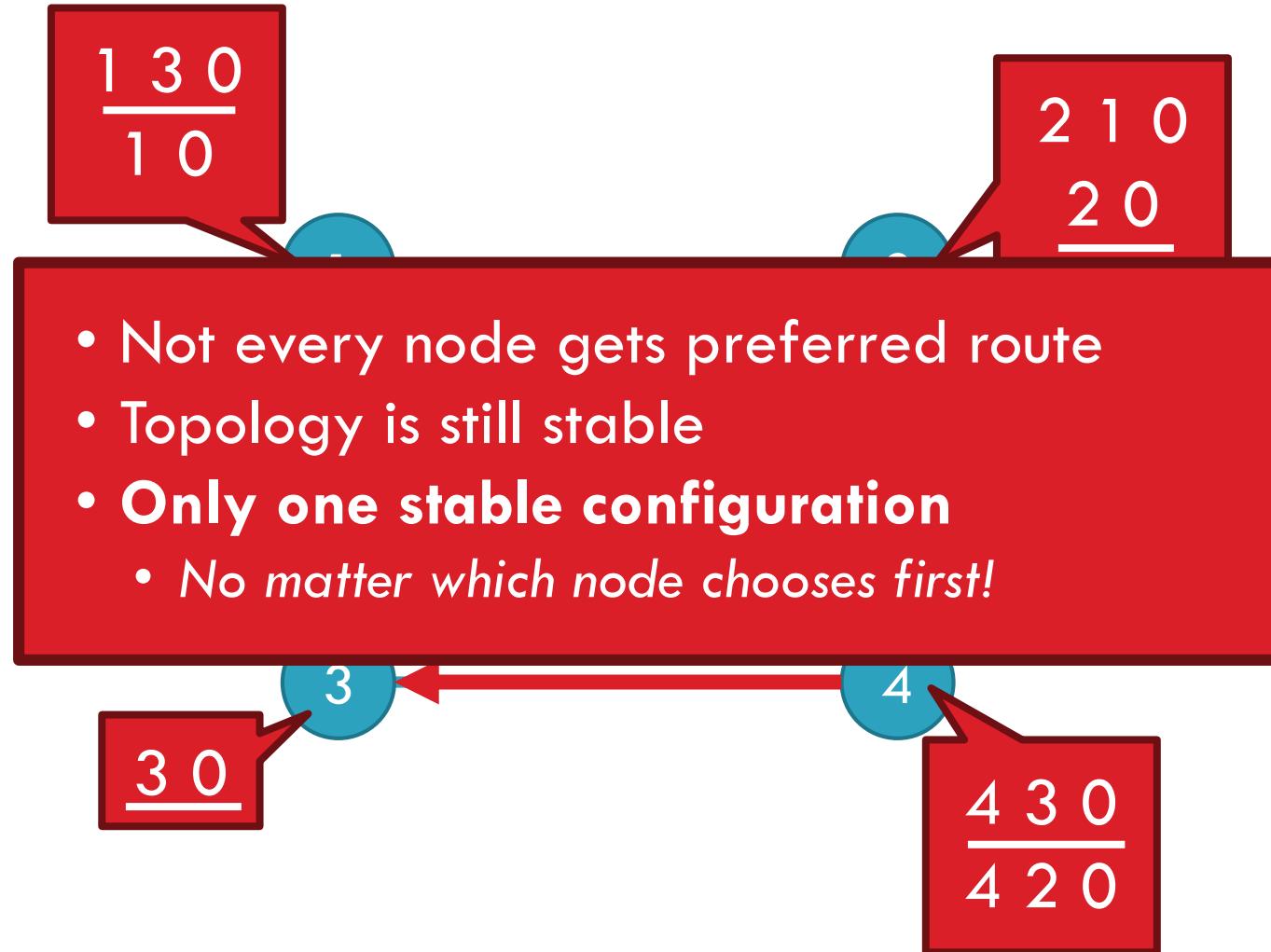


- Each node gets its preferred route
- Totally stable topology



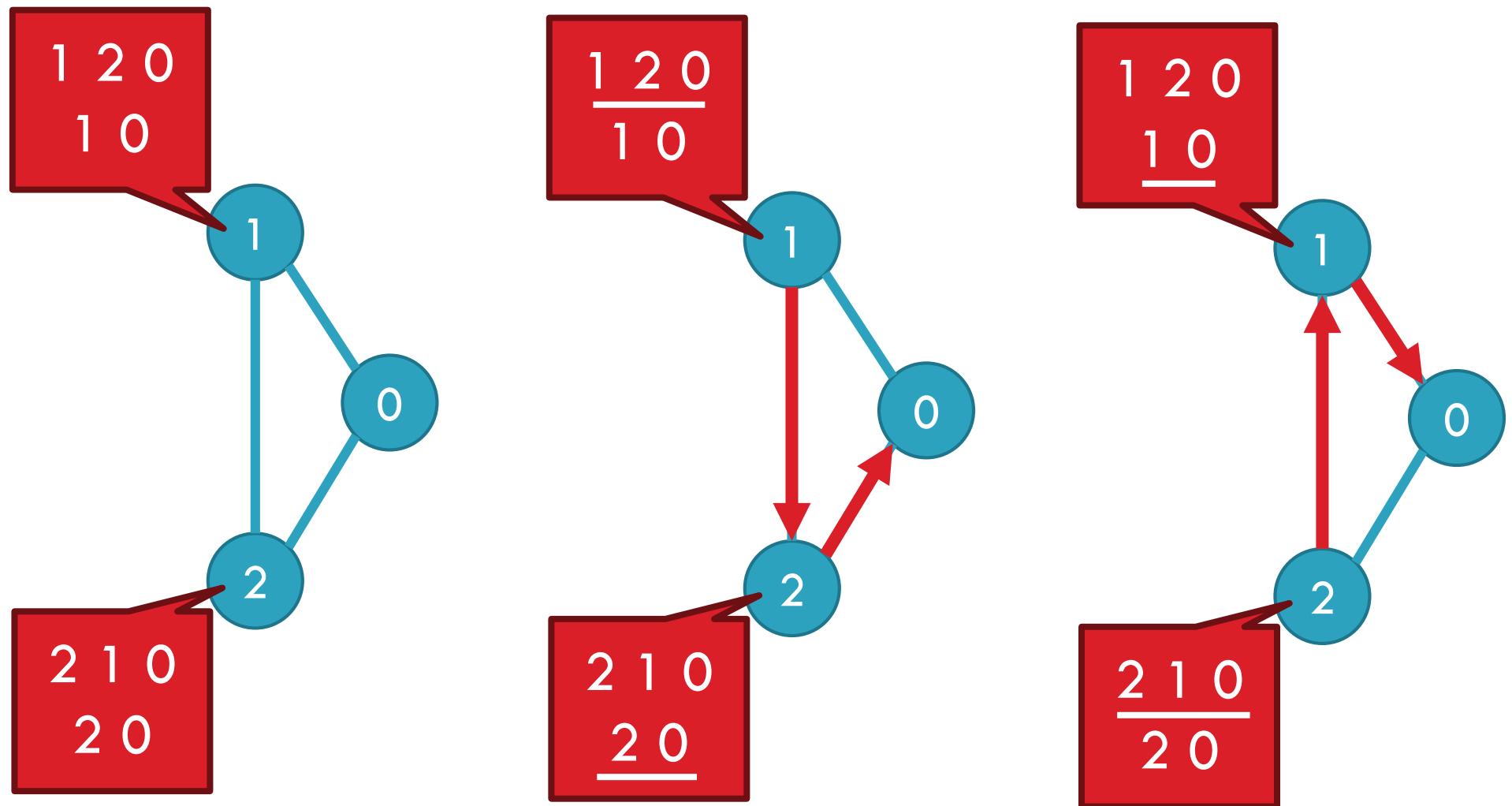
Good Gadget

102



SPP May Have Multiple Solutions

103



Bad Gadget

104

1 3 0

- That was only one round of oscillation!
- This keeps going, infinitely
- Problem stems from:
 - Local (not global) decisions
 - Ability of one node to improve its path selection

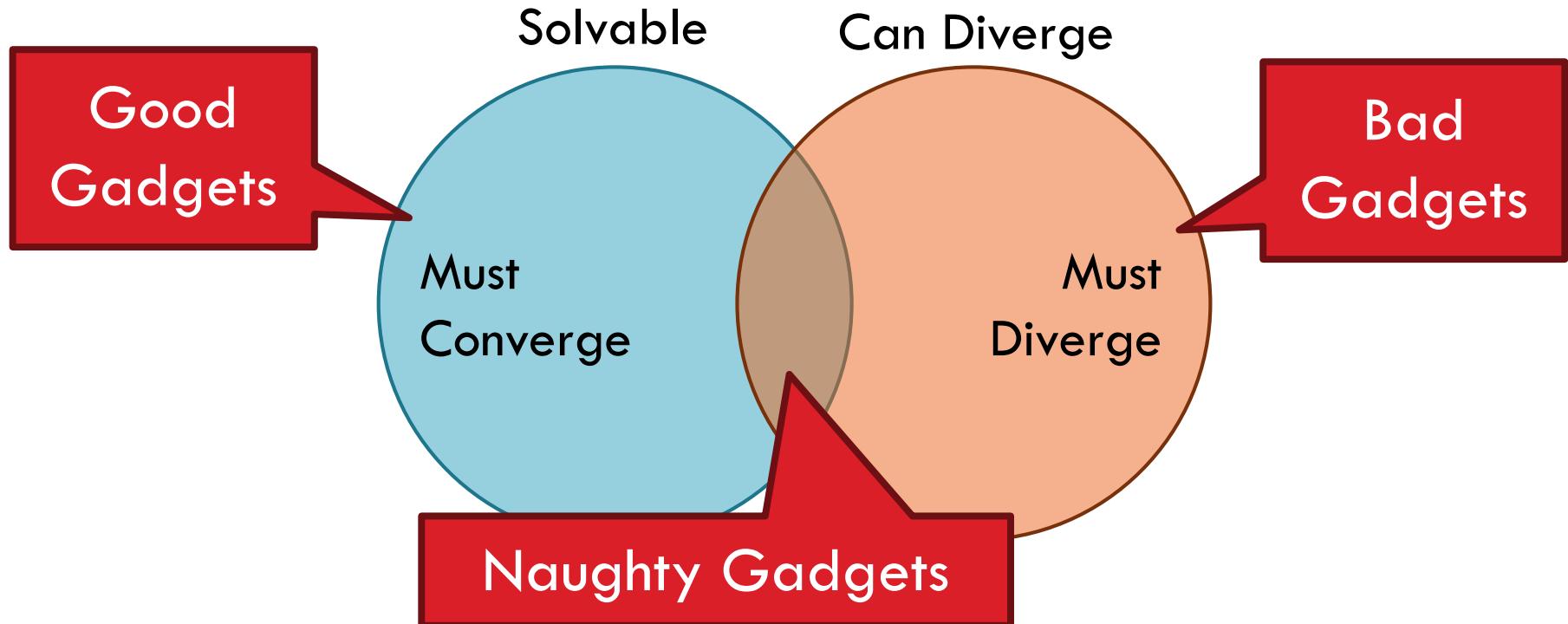
—

4 3 0

SPP Explains BGP Divergence

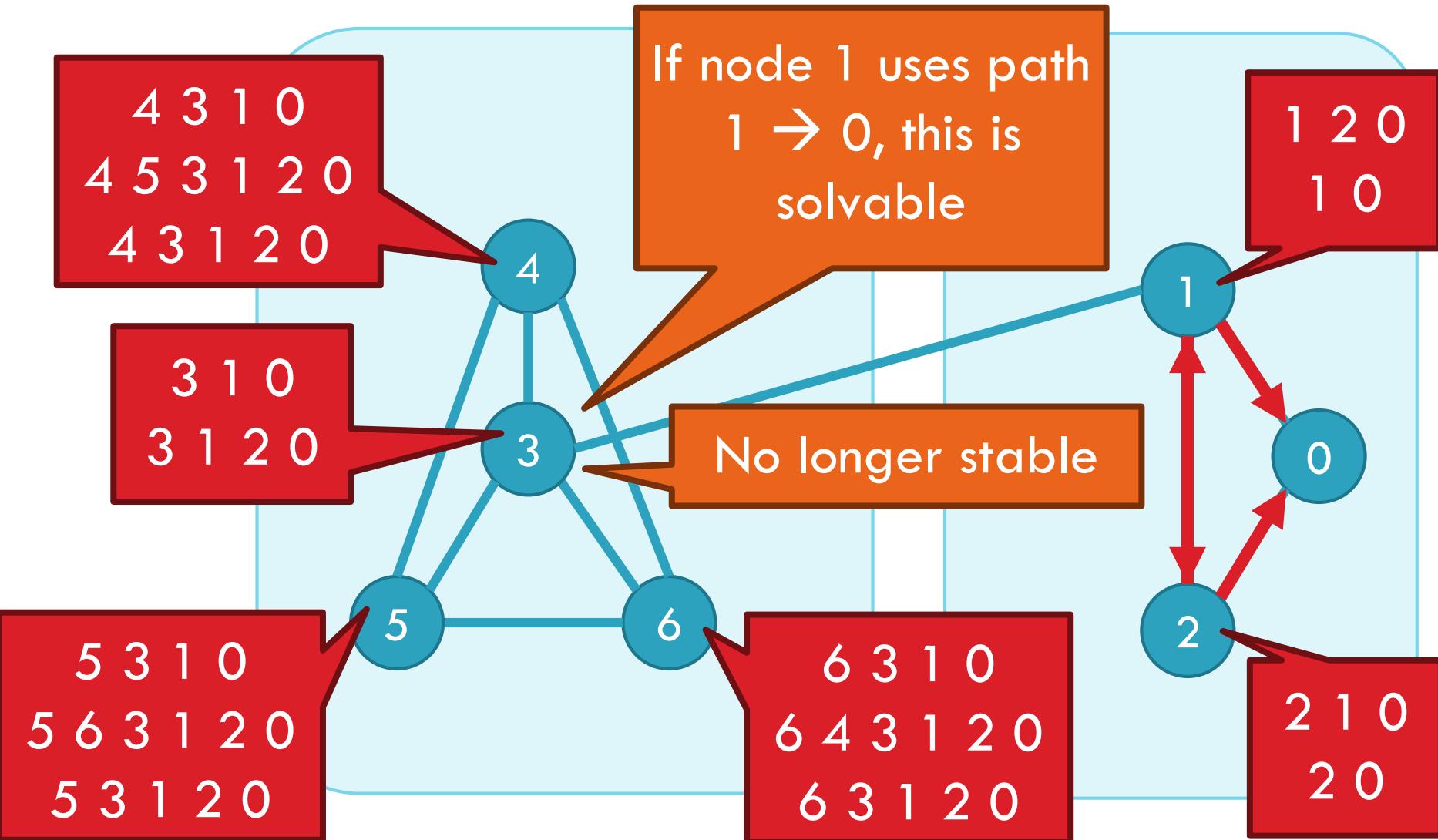
105

- BGP is **not** guaranteed to converge to stable routing
 - Policy inconsistencies may lead to “livelock”
 - Protocol oscillation



BGP is Precarious

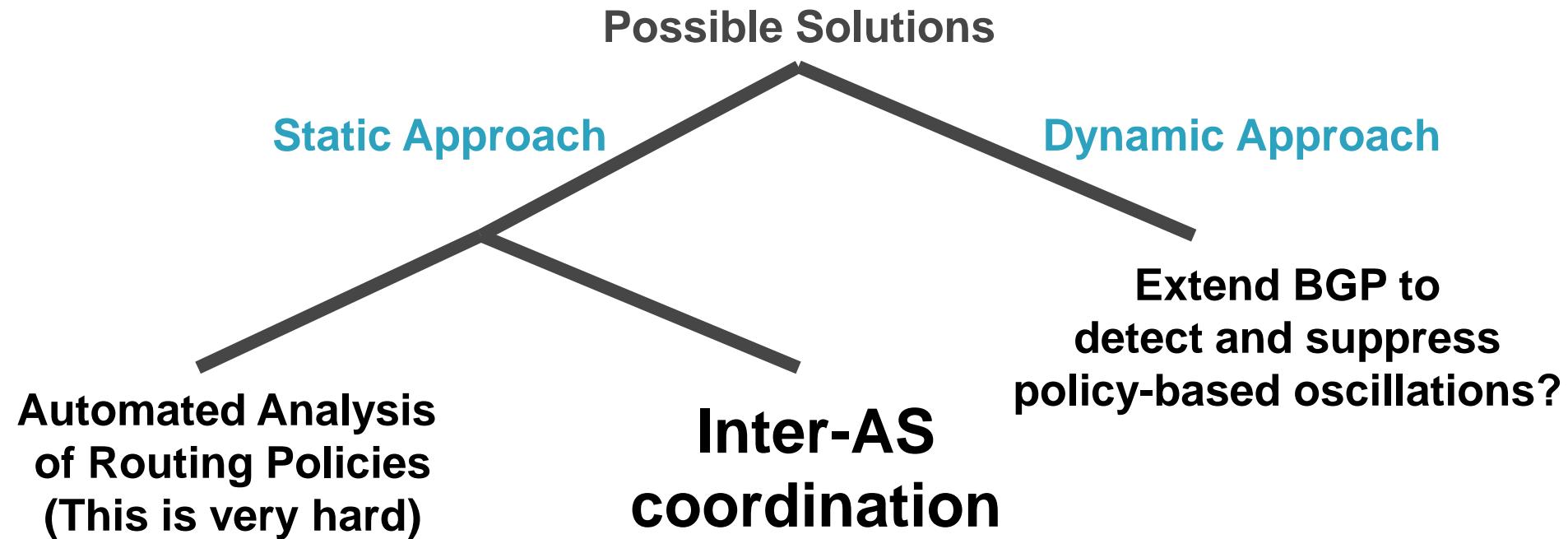
107



Can BGP Be Fixed?

108

- Unfortunately, SPP is NP-complete



These approaches are complementary