

BGP Best Current Practices



ISP Workshops

Configuring BGP



Where do we start?

IOS Good Practices

- ❑ ISPs should start off with the following BGP commands as a basic template:

```
router bgp 64511
```

← Replace with public ASN

```
bgp deterministic-med
```

```
distance bgp 200 200 200
```

← Make ebgp and ibgp distance the same

```
no synchronization
```

```
no auto-summary
```

- ❑ If supporting more than just IPv4 unicast neighbours

```
no bgp default ipv4-unicast
```

- is also very important and required

Cisco IOS Good Practices

- ❑ BGP in Cisco IOS is **permissive** by default
- ❑ Configuring BGP peering without using filters means:
 - All best paths on the local router are passed to the neighbour
 - All routes announced by the neighbour are received by the local router
 - Can have disastrous consequences
- ❑ **Good practice is to ensure that each eBGP neighbour has inbound and outbound filter applied:**

```
router bgp 64511
  neighbor 1.2.3.4 remote-as 64510
  neighbor 1.2.3.4 prefix-list as64510-in in
  neighbor 1.2.3.4 prefix-list as64510-out out
```

What is BGP for??



What is an IGP not for?

BGP versus OSPF/ISIS

- ❑ Internal Routing Protocols (IGPs)
 - examples are ISIS and OSPF
 - used for carrying **infrastructure** addresses
 - **NOT** used for carrying Internet prefixes or customer prefixes
 - design goal is to **minimise** number of prefixes in IGP to aid scalability and rapid convergence

BGP versus OSPF/ISIS

- ❑ BGP used internally (iBGP) and externally (eBGP)
- ❑ iBGP used to carry
 - some/all Internet prefixes across backbone
 - customer prefixes
- ❑ eBGP used to
 - exchange prefixes with other ASes
 - implement routing policy

BGP versus OSPF/ISIS

- ❑ DO NOT:
 - distribute BGP prefixes into an IGP
 - distribute IGP routes into BGP
 - use an IGP to carry customer prefixes
- ❑ **YOUR NETWORK WILL NOT SCALE**

Aggregation



Aggregation

- ❑ Aggregation means announcing the address block received from the RIR to the other ASes connected to your network
- ❑ Subprefixes of this aggregate may be:
 - Used internally in the ISP network
 - Announced to other ASes to aid with multihoming
- ❑ Unfortunately too many people are still thinking about class Cs, resulting in a proliferation of /24s in the Internet routing table
 - Note: Same is happening for /48s with IPv6

Configuring Aggregation – Cisco IOS

- ❑ ISP has 101.10.0.0/19 address block
- ❑ To put into BGP as an aggregate:

```
router bgp 64511
  network 101.10.0.0 mask 255.255.224.0
  ip route 101.10.0.0 255.255.224.0 null0
```
- ❑ The static route is a “pull up” route
 - more specific prefixes within this address block ensure connectivity to ISP’s customers
 - “longest match lookup

Aggregation

- ❑ Address block should be announced to the Internet as an aggregate
- ❑ Subprefixes of address block should **NOT** be announced to Internet unless for traffic engineering
 - See BGP Multihoming presentations
- ❑ Aggregate should be generated internally
 - Not on the network borders!

Announcing Aggregate – Cisco IOS

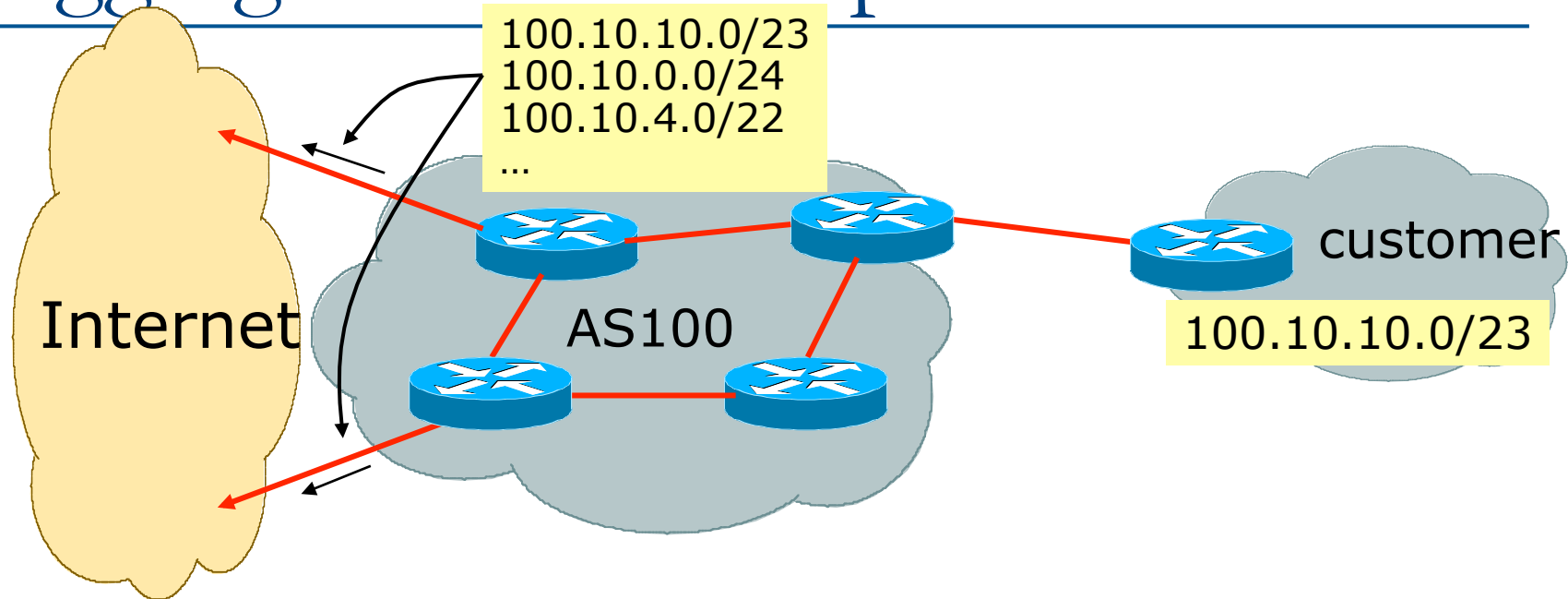
□ Configuration Example

```
router bgp 64511
  network 101.10.0.0 mask 255.255.224.0
  neighbor 102.102.10.1 remote-as 101
  neighbor 102.102.10.1 prefix-list out-filter out
!
ip route 101.10.0.0 255.255.224.0 null0
!
ip prefix-list out-filter permit 101.10.0.0/19
ip prefix-list out-filter deny 0.0.0.0/0 le 32
```

Announcing an Aggregate

- ❑ ISPs who don't and won't aggregate are held in poor regard by community
- ❑ Registries publish their minimum allocation size
 - Now ranging from a /20 to a /24 depending on RIR
 - Different sizes for different address blocks
 - (APNIC changed its minimum allocation to /24 in October 2010)
- ❑ Until recently there was no real reason to see anything longer than a /22 prefix in the Internet
 - BUT there are currently (July 2013) >242000 /24s!
 - IPv4 run-out is starting to have an impact

Aggregation – Example

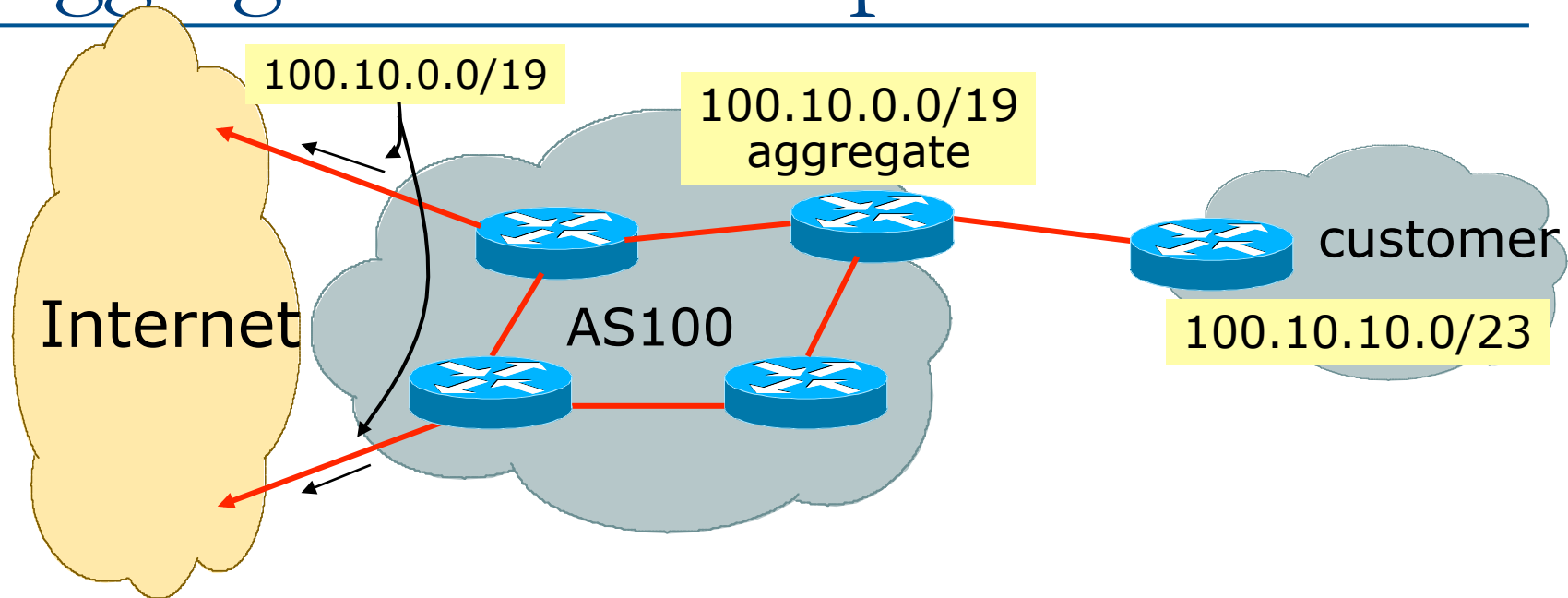


- ❑ Customer has /23 network assigned from AS100's /19 address block
- ❑ AS100 announces customers' individual networks to the Internet

Aggregation – Bad Example

- Customer link goes down
 - Their /23 network becomes unreachable
 - /23 is withdrawn from AS100's iBGP
 - Their ISP doesn't aggregate its /19 network block
 - /23 network withdrawal announced to peers
 - starts rippling through the Internet
 - added load on all Internet backbone routers as network is removed from routing table
-
- Customer link returns
 - Their /23 network is now visible to their ISP
 - Their /23 network is re-advertised to peers
 - Starts rippling through Internet
 - Load on Internet backbone routers as network is reinserted into routing table
 - Some ISP's suppress the flaps
 - Internet may take 10-20 min or longer to be visible
 - Where is the Quality of Service???

Aggregation – Example



- ❑ Customer has /23 network assigned from AS100's /19 address block
- ❑ AS100 announced /19 aggregate to the Internet

Aggregation – Good Example

- ❑ Customer link goes down
 - their /23 network becomes unreachable
 - /23 is withdrawn from AS100's iBGP
 - ❑ /19 aggregate is still being announced
 - no BGP hold down problems
 - no BGP propagation delays
 - no damping by other ISPs
-
- ❑ Customer link returns
 - ❑ Their /23 network is visible again
 - The /23 is re-injected into AS100's iBGP
 - ❑ The whole Internet becomes visible immediately
 - ❑ Customer has Quality of Service perception

Aggregation – Summary

- Good example is what everyone should do!
 - Adds to Internet stability
 - Reduces size of routing table
 - Reduces routing churn
 - Improves Internet QoS for **everyone**
- Bad example is what too many still do!
 - Why? Lack of knowledge?
 - Laziness?

Separation of iBGP and eBGP

- ❑ Many ISPs do not understand the importance of separating iBGP and eBGP
 - iBGP is where all customer prefixes are carried
 - eBGP is used for announcing aggregate to Internet and for Traffic Engineering
- ❑ Do **NOT** do traffic engineering with customer originated iBGP prefixes
 - Leads to instability similar to that mentioned in the earlier bad example
 - Even though aggregate is announced, a flapping subprefix will lead to instability for the customer concerned
- ❑ **Generate traffic engineering prefixes on the Border Router**

The Internet Today (July 2013)

□ Current Internet Routing Table Statistics

■ BGP Routing Table Entries	460126
■ Prefixes after maximum aggregation	186940
■ Unique prefixes in Internet	228811
■ Prefixes smaller than registry alloc	161061
■ /24s announced	242050
■ ASes in use	44587

Efforts to improve aggregation

□ The CIDR Report

- Initiated and operated for many years by Tony Bates
- Now combined with Geoff Huston's routing analysis
 - www.cidr-report.org
 - (covers both IPv4 and IPv6 BGP tables)
- Results e-mailed on a weekly basis to most operations lists around the world
- Lists the top 30 service providers who could do better at aggregating

□ RIPE Routing WG aggregation recommendations

- IPv4: RIPE-399 — www.ripe.net/ripe/docs/ripe-399.html
- IPv6: RIPE-532 — www.ripe.net/ripe/docs/ripe-532.html

Efforts to Improve Aggregation

The CIDR Report

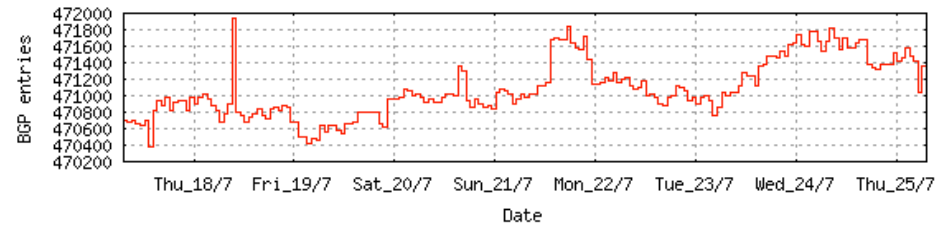
- ❑ Also computes the size of the routing table assuming ISPs performed optimal aggregation
- ❑ Website allows searches and computations of aggregation to be made on a per AS basis
 - Flexible and powerful tool to aid ISPs
 - Intended to show how greater efficiency in terms of BGP table size can be obtained without loss of routing and policy information
 - Shows what forms of origin AS aggregation could be performed and the potential benefit of such actions to the total table size
 - Very effectively challenges the traffic engineering excuse

Status Summary

Table History

Date	Prefixes	CIDR Aggregated
18-07-13	470980	267519
19-07-13	470686	267677
20-07-13	470966	267992
21-07-13	470849	265755
22-07-13	471142	266003
23-07-13	470970	266883
24-07-13	471645	267187
25-07-13	471515	266540

Plot: [BGP Table Size](#)



AS Summary

44720	Number of ASes in routing system
18453	Number of ASes announcing only one prefix
4219	Largest number of prefixes announced by an AS
	AS7029 : WINDSTREAM - Windstream Communications Inc
117330144	Largest address span announced by an AS (/32s)
	AS4134 : CHINANET-BACKBONE No.31,Jin-rong Street

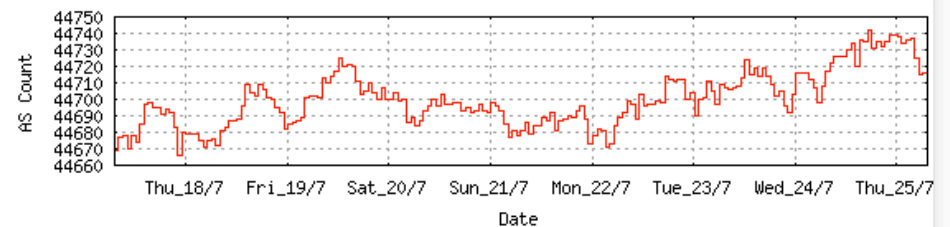
Plot: [AS count](#)

Plot: [Average announcements per origin AS](#)

Report: [ASes ordered by originating address span](#)

Report: [ASes ordered by transit address span](#)

Report: [Autonomous System number-to-name mapping \(from Registry WHOIS data\)](#)



Announced Prefixes

Rank	AS	Type	Originate	Addr Space (pfx)	Transit	Addr space (pfx)	Description
15	AS6389		ORG+TRN Originate:	29761280 /7.17	Transit:	936704 /12.16	BELLSOUTH-NET-BLK - BellSouth.net Inc.

Aggregation Suggestions

This report does not take into account conditions local to each origin AS in terms of policy or traffic engineering requirements, so this is an approximate guideline as to aggregation possibilities.

Rank	AS	AS Name	Current	Wthdw	Aggte	Annce	Redctn	%
2	AS6389	BELLSOUTH-NET-BLK - BellSouth.net Inc.	2985	2921	7	71	2914	97.62%

Prefix	AS Path	Aggregation Suggestion
12.81.90.0/23	4777 2516 3356 7018 6389	
12.81.120.0/24	4777 2516 3356 7018 6389	
12.83.3.0/24	4777 2516 3356 7018 6389	
12.83.5.0/24	4777 2516 3356 7018 6389	
12.83.7.0/24	4777 2516 3356 7018 6389	
65.0.0.0/12	4777 2516 3356 7018 6389	
65.0.0.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.0.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.40.0/22	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.50.0/23	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.64.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.128.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.192.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.0.224.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.0.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.32.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.64.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.128.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.224.0/20	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.1.240.0/20	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.2.0.0/16	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.2.0.0/17	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.2.128.0/17	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.3.224.0/19	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.4.64.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.4.192.0/18	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.1.0/24	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.12.0/22	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.16.0/22	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.20.0/23	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.21.0/24	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.22.0/23	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.24.0/22	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389
65.5.28.0/22	4777 2516 3356 7018 6389 - Withdrawn	- matching aggregate 65.0.0.0/12 4777 2516 3356 7018 6389

Announced Prefixes

Rank	AS	Type	Originate	Addr Space (pfx)	Transit	Addr space (pfx)	Description
184	AS18566	ORG+TRN	Originate:	2788864 /10.59	Transit:	50944 /16.36	COVAD - Covad Communications Co.

Aggregation Suggestions

This report does not take into account conditions local to each origin AS in terms of policy or traffic engineering requirements, so this is an approximate guideline as to aggregation possibilities.

Rank	AS	AS Name	Current	Withdw	Aggte	Annce	Redctn	%
8	AS18566	COVAD - Covad Communications Co.	2067	1822	223	468	1599	77.36%

Prefix	AS Path	Aggregation Suggestion
64.81.16.0/22	4777 2516 3356 18566	
64.81.20.0/22	4777 2516 4565 18566	
64.81.22.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.20.0/22 4777 2516 4565 18566
64.81.24.0/21	4777 2516 3356 18566	+ Announce - aggregate of 64.81.24.0/22 (4777 2516 3356 18566) and 64.81.28.0/22 (4777 2516 3356 18566)
64.81.24.0/22	4777 2516 3356 18566	- Withdrawn - aggregated with 64.81.28.0/22 (4777 2516 3356 18566)
64.81.28.0/22	4777 2516 3356 18566	- Withdrawn - aggregated with 64.81.24.0/22 (4777 2516 3356 18566)
64.81.32.0/20	4777 2516 4565 18566	
64.81.32.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.33.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.34.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.35.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.36.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.37.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.38.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.39.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.40.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.44.0/24	4777 2516 4565 18566	- Withdrawn - matching aggregate 64.81.32.0/20 4777 2516 4565 18566
64.81.48.0/20	4777 2516 3356 18566	
64.81.48.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.49.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.50.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.51.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.52.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.53.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.54.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.55.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.56.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.57.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.58.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.59.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.60.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.61.0/24	4777 2516 3356 18566	- Withdrawn - matching aggregate 64.81.48.0/20 4777 2516 3356 18566
64.81.64.0/20	4777 2516 3356 18566	

Importance of Aggregation

- ❑ Size of routing table
 - Router Memory is not so much of a problem as it was in the 1990s
 - Routers routinely carry over 1 million prefixes
- ❑ Convergence of the Routing System
 - This is a problem
 - Bigger table takes longer for CPU to process
 - BGP updates take longer to deal with
 - BGP Instability Report tracks routing system update activity
 - bgpupdates.potaroo.net/instability/bgpupd.html

The BGP Instability Report

The BGP Instability Report is updated daily. This report was generated on 25 July 2013 06:29 (UTC+1000)

50 Most active ASes for the past 7 days

RANK	ASN	UPDs	%	Prefixes	UPDs/Prefix	AS NAME
1	18403	51249	1.66%	599	85.56	FPT-AS-AP The Corporation for Financing & Promoting Technology
2	9829	39127	1.26%	1537	25.46	BSNL-NIB National Internet Backbone
3	10620	36190	1.17%	2701	13.40	Telmex Colombia S.A.
4	8402	33575	1.08%	1822	18.43	CORBINA-AS OJSC "Vimpelcom"
5	28573	29960	0.97%	3016	9.93	NET Serviços de Comunicação S.A.
6	27738	28301	0.91%	576	49.13	Ecuadortelecom S.A.
7	4538	25037	0.81%	536	46.71	ERX-CERNET-BKB China Education and Research Network Center
8	15003	24391	0.79%	854	28.56	NOBIS-TECH - Nobis Technology Group, LLC
9	10428	22688	0.73%	7	3241.14	CWV-NETWORKS - The College of West Virginia
10	50710	20310	0.66%	239	84.98	EARTHLINK-AS EarthLink Ltd. Communications&Internet Services
11	17974	19721	0.64%	2626	7.51	TELKOMNET-AS2-AP PT Telekomunikasi Indonesia
12	33770	18472	0.60%	76	243.05	KDN
13	9416	17757	0.57%	65	273.18	MULTIMEDIA-AS-AP Hoshin Multimedia Center Inc.
14	4775	17051	0.55%	127	134.26	GLOBE-TELECOM-AS Globe Telecoms
15	3356	13728	0.44%	1105	12.42	LEVEL3 Level 3 Communications
16	36998	13025	0.42%	1819	7.16	SDN-MOBITEL
17	8151	12695	0.41%	1284	9.89	Uninet S.A. de C.V.
18	14287	12590	0.41%	63	199.84	TRIAD-TELECOM - Triad Telecom, Inc.
19	7552	12207	0.39%	1191	10.25	VIETEL-AS-AP Vietel Corporation
20	45899	11481	0.37%	374	30.70	VNPT-AS-VN VNPT Corp
21	13188	11383	0.37%	838	13.58	BANKINFORM-AS TOV "Bank-Inform"
22	52280	11306	0.37%	6	1884.33	INTERNEXA Chile S.A.
23	34969	10904	0.35%	8	1363.00	PASJONET-AS Pasjo.Net Sp, z o.o.
24	7044	10440	0.34%	1470	8.88	FRONTIER AND CITIZENS Frontier Communications of America, Inc.

50 Most active Prefixes for the past 7 days

RANK	PREFIX	UPDs	%	Origin AS – AS NAME
1	190.211.175.0/24	11701	0.33%	28032 – INTERNEXA S.A. 52280 – INTERNEXA Chile S.A.
2	92.246.207.0/24	10031	0.28%	48612 – RTC-ORENBURG-AS CJSC "Comstar-Regions"
3	203.118.232.0/21	8889	0.25%	9416 – MULTIMEDIA-AS-AP Hoshin Multimedia Center Inc.
4	203.118.224.0/21	8704	0.24%	9416 – MULTIMEDIA-AS-AP Hoshin Multimedia Center Inc.
5	192.58.232.0/24	8587	0.24%	6629 – NOAA-AS - NOAA
6	222.127.0.0/24	8241	0.23%	4775 – GLOBE-TELECOM-AS Globe Telecoms
7	120.28.62.0/24	8167	0.23%	4775 – GLOBE-TELECOM-AS Globe Telecoms
8	12.43.218.0/24	7536	0.21%	10428 – CWV-NETWORKS - The College of West Virginia
9	199.248.240.0/24	7536	0.21%	10428 – CWV-NETWORKS - The College of West Virginia
10	205.166.165.0/24	7536	0.21%	10428 – CWV-NETWORKS - The College of West Virginia
11	65.90.49.0/24	7304	0.21%	3356 – LEVEL3 Level 3 Communications
12	62.84.76.0/24	6502	0.18%	42334 – BBP-AS Broadband Plus s.a.l.
13	69.38.178.0/24	4642	0.13%	19406 – TWRS-MA - Towerstream I, Inc.
14	64.187.64.0/23	4143	0.12%	16608 – KENTEC - Kentec Communications, Inc.
15	115.170.128.0/17	4068	0.11%	4847 – CNIX-AP China Networks Inter-Exchange
16	211.214.206.0/24	3996	0.11%	9854 – KTO-AS-KR KTO
17	206.105.75.0/24	3560	0.10%	6174 – SPRINTLINK8 - Sprint
18	208.16.110.0/24	3560	0.10%	6174 – SPRINTLINK8 - Sprint
19	64.187.64.0/24	3248	0.09%	16608 – KENTEC - Kentec Communications, Inc.
20	213.133.192.0/24	2928	0.08%	13208 – NEWTELSOLUTIONS-AS Newtel Ltd
21	213.133.193.0/24	2928	0.08%	13208 – NEWTELSOLUTIONS-AS Newtel Ltd
22	178.61.252.0/23	2892	0.08%	21050 – FAST-TELCO Fast Telecommunications Company W.L.L.
23	2.93.235.0/24	2851	0.08%	8402 – CORBINA-AS OJSC "Vimpelcom"
25	84.205.66.0/24	2505	0.07%	12654 – RIPE-NCC-RIS-AS Reseaux IP Europeens Network Coordination Centre (RIPE NCC)
26	208.73.244.0/22	2460	0.07%	14287 – TRIAD-TELECOM - Triad Telecom, Inc.
27	208.88.232.0/21	2460	0.07%	14287 – TRIAD-TELECOM - Triad Telecom, Inc.
28	216.162.0.0/20	2460	0.07%	14287 – TRIAD-TELECOM - Triad Telecom, Inc.
29	208.78.116.0/22	2458	0.07%	14287 – TRIAD-TELECOM - Triad Telecom, Inc.

Receiving Prefixes



Receiving Prefixes

- ❑ There are three scenarios for receiving prefixes from other ASNs
 - Customer talking BGP
 - Peer talking BGP
 - Upstream/Transit talking BGP
- ❑ Each has different filtering requirements and need to be considered separately

Receiving Prefixes: From Customers

- ❑ ISPs should only accept prefixes which have been assigned or allocated to their downstream customer
- ❑ If ISP has assigned address space to its customer, then the customer IS entitled to announce it back to his ISP
- ❑ If the ISP has NOT assigned address space to its customer, then:
 - Check in the five RIR databases to see if this address space really has been assigned to the customer
 - The tool: `whois -h jwhois.apnic.net x.x.x.0/24`
 - ❑ (jwhois queries all RIR databases)

Receiving Prefixes: From Customers

- Example use of whois to check if customer is entitled to announce address space:

```
$ whois -h whois.apnic.net 202.12.29.0
inetnum:          202.12.28.0 - 202.12.29.255
netname:          APNIC-AP
descr:            Asia Pacific Network Information Centre
descr:            Regional Internet Registry for the Asia-Pacific
descr:            6 Cordelia Street
descr:            South Brisbane, QLD 4101
descr:            Australia
country:          AU
admin-c:          AIC1-AP
tech-c:           NO4-AP
mnt-by:           APNIC-HM
mnt-irt:           IRT-APNIC-AP
changed:          hm-changed@apnic.net
status:           ASSIGNED PORTABLE
changed:          hm-changed@apnic.net 20110309
source:           APNIC
```

Portable – means its an assignment to the customer, the customer can announce it to you

Receiving Prefixes: From Customers

- Example use of whois to check if customer is entitled to announce address space:

```
$ whois -h whois.ripe.net 193.128.0.0
inetnum:          193.128.0.0 - 193.133.255.255
netname:          UK-PIPEX-193-128-133
descr:           Verizon UK Limited
country:         GB
org:             ORG-UA24-RIPE
admin-c:         WERT1-RIPE
tech-c:          UPHM1-RIPE
status:          ALLOCATED UNSPECIFIED
remarks:         Please send abuse notification to abuse@uk.uu.net
mnt-by:          RIPE-NCC-HM-MNT
mnt-lower:       AS1849-MNT
mnt-routes:      AS1849-MNT
mnt-routes:      WCOM-EMEA-RICE-MNT
mnt-irt:         IRT-MCI-GB
source:          RIPE # Filtered
```

ALLOCATED – means that this is Provider Aggregatable address space and can only be announced by the ISP holding the allocation (in this case Verizon UK)

Receiving Prefixes from customer: Cisco IOS

- ❑ For Example:
 - downstream has 100.50.0.0/20 block
 - should only announce this to upstreams
 - upstreams should only accept this from them
- ❑ Configuration on upstream

```
router bgp 100
  neighbor 102.102.10.1 remote-as 101
  neighbor 102.102.10.1 prefix-list customer in
!
ip prefix-list customer permit 100.50.0.0/20
```

Receiving Prefixes: From Peers

- A peer is an ISP with whom you agree to exchange prefixes you originate into the Internet routing table
 - Prefixes you accept from a peer are only those they have indicated they will announce
 - Prefixes you announce to your peer are only those you have indicated you will announce

Receiving Prefixes: From Peers

- Agreeing what each will announce to the other:

- Exchange of e-mail documentation as part of the peering agreement, and then ongoing updates

OR

- Use of the Internet Routing Registry and configuration tools such as the IRRToolSet

www.isc.org/sw/IRRToolSet/

Receiving Prefixes from peer: Cisco IOS

- ❑ For Example:
 - Peer has 220.50.0.0/16, 61.237.64.0/18 and 81.250.128.0/17 address blocks
- ❑ Configuration on local router

```
router bgp 100
  neighbor 102.102.10.1 remote-as 101
  neighbor 102.102.10.1 prefix-list my-peer in
!
ip prefix-list my-peer permit 220.50.0.0/16
ip prefix-list my-peer permit 61.237.64.0/18
ip prefix-list my-peer permit 81.250.128.0/17
ip prefix-list my-peer deny 0.0.0.0/0 le 32
```

Receiving Prefixes:

From Upstream/Transit Provider

- ❑ Upstream/Transit Provider is an ISP who you pay to give you transit to the **WHOLE** Internet
- ❑ Receiving prefixes from them is not desirable unless really necessary
 - Traffic Engineering – see BGP Multihoming presentations
- ❑ Ask upstream/transit provider to either:
 - originate a default-route
 - OR
 - announce one prefix you can use as default

Receiving Prefixes: From Upstream/Transit Provider

❑ Downstream Router Configuration

```
router bgp 100
  network 101.10.0.0 mask 255.255.224.0
  neighbor 101.5.7.1 remote-as 101
  neighbor 101.5.7.1 prefix-list infilter in
  neighbor 101.5.7.1 prefix-list outfilter out
!
ip prefix-list infilter permit 0.0.0.0/0
!
ip prefix-list outfilter permit 101.10.0.0/19
```


Receiving Prefixes: From Upstream/Transit Provider

□ Upstream Router Configuration

```
router bgp 101
  neighbor 101.5.7.2 remote-as 100
  neighbor 101.5.7.2 default-originate
  neighbor 101.5.7.2 prefix-list cust-in in
  neighbor 101.5.7.2 prefix-list cust-out out
!
ip prefix-list cust-in permit 101.10.0.0/19
!
ip prefix-list cust-out permit 0.0.0.0/0
```

Receiving Prefixes:

From Upstream/Transit Provider

- ❑ If necessary to receive prefixes from any provider, care is required.
 - Don't accept default (unless you need it)
 - Don't accept your own prefixes
- ❑ For IPv4:
 - Don't accept private (RFC1918) and certain special use prefixes:
<http://www.rfc-editor.org/rfc/rfc5735.txt>
 - Don't accept prefixes longer than /24 (?)
- ❑ For IPv6:
 - Don't accept certain special use prefixes:
<http://www.rfc-editor.org/rfc/rfc5156.txt>
 - Don't accept prefixes longer than /48 (?)

Receiving Prefixes: From Upstream/Transit Provider

- ❑ Check Team Cymru's list of "bogons"
www.team-cymru.org/Services/Bogons/http.html
- ❑ For IPv4 also consult:
www.rfc-editor.org/rfc/rfc6441.txt
- ❑ For IPv6 also consult:
www.space.net/~gert/RIPE/ipv6-filters.html
- ❑ Bogon Route Server:
www.team-cymru.org/Services/Bogons/routeserver.html
 - Supplies a BGP feed (IPv4 and/or IPv6) of address blocks which should not appear in the BGP table

Receiving IPv4 Prefixes

```
router bgp 100
  network 101.10.0.0 mask 255.255.224.0
  neighbor 101.5.7.1 remote-as 101
  neighbor 101.5.7.1 prefix-list in-filter in
  !
  ip prefix-list in-filter deny 0.0.0.0/0                ! Default
  ip prefix-list in-filter deny 0.0.0.0/8 le 32          ! Network Zero
  ip prefix-list in-filter deny 10.0.0.0/8 le 32         ! RFC1918
  ip prefix-list in-filter deny 100.64.0.0/10 le 32      ! RFC6598 shared addr
  ip prefix-list in-filter deny 101.10.0.0/19 le 32      ! Local prefix
  ip prefix-list in-filter deny 127.0.0.0/8 le 32        ! Loopback
  ip prefix-list in-filter deny 169.254.0.0/16 le 32     ! Auto-config
  ip prefix-list in-filter deny 172.16.0.0/12 le 32      ! RFC1918
  ip prefix-list in-filter deny 192.0.2.0/24 le 32       ! TEST1
  ip prefix-list in-filter deny 192.168.0.0/16 le 32     ! RFC1918
  ip prefix-list in-filter deny 198.18.0.0/15 le 32      ! Benchmarking
  ip prefix-list in-filter deny 198.51.100.0/24 le 32    ! TEST2
  ip prefix-list in-filter deny 203.0.113.0/24 le 32     ! TEST3
  ip prefix-list in-filter deny 224.0.0.0/3 le 32        ! Multicast
  ip prefix-list in-filter deny 0.0.0.0/0 ge 25          ! Prefixes >/24
  ip prefix-list in-filter permit 0.0.0.0/0 le 32
```

Receiving IPv6 Prefixes

```
router bgp 100
  network 2020:3030::/32
  neighbor 2020:3030::1 remote-as 101
  neighbor 2020:3030::1 prefix-list v6in-filter in
  !
  ipv6 prefix-list v6in-filter permit 2001::/32           ! Teredo
  ipv6 prefix-list v6in-filter deny 2001::/32 le 128      ! Teredo subnets
  ipv6 prefix-list v6in-filter deny 2001:db8::/32 le 128  ! Documentation
  ipv6 prefix-list v6in-filter permit 2002::/16           ! 6to4
  ipv6 prefix-list v6in-filter deny 2002::/16 le 128      ! 6to4 subnets
  ipv6 prefix-list v6in-filter deny 2020:3030::/32 le 128 ! Local Prefix
  ipv6 prefix-list v6in-filter deny 3ffe::/16 le 128      ! Old 6bone
  ipv6 prefix-list v6in-filter permit 2000::/3 le 48      ! Global Unicast
  ipv6 prefix-list v6in-filter deny ::/0 le 128
```

Receiving Prefixes

- ❑ Paying attention to prefixes received from customers, peers and transit providers assists with:
 - The integrity of the local network
 - The integrity of the Internet
- ❑ Responsibility of all ISPs to be good Internet citizens

Prefixes into iBGP



Injecting prefixes into iBGP

- ❑ Use iBGP to carry customer prefixes
 - don't use IGP
- ❑ Point static route to customer interface
- ❑ Use BGP network statement
- ❑ As long as static route exists (interface active), prefix will be in BGP

Router Configuration: network statement

□ Example:

```
interface loopback 0
  ip address 215.17.3.1 255.255.255.255
!
interface Serial 5/0
  ip unnumbered loopback 0
  ip verify unicast reverse-path
!
ip route 215.34.10.0 255.255.252.0 Serial 5/0
!
router bgp 100
  network 215.34.10.0 mask 255.255.252.0
```

Injecting prefixes into iBGP

- ❑ Interface flap will result in prefix withdraw and reannounce
 - use `"ip route . . . permanent"`
- ❑ Many ISPs redistribute static routes into BGP rather than using the network statement
 - Only do this if you understand why

Router Configuration:

redistribute static

□ Example:

```
ip route 215.34.10.0 255.255.252.0 Serial 5/0
!
router bgp 100
 redistribute static route-map static-to-bgp
<snip>
!
route-map static-to-bgp permit 10
 match ip address prefix-list ISP-block
 set origin igp
<snip>
!
ip prefix-list ISP-block permit 215.34.10.0/22 le 30
```

Injecting prefixes into iBGP

- ❑ Route-map ISP-block can be used for many things:
 - setting communities and other attributes
 - setting origin code to IGP, etc
- ❑ Be careful with prefix-lists and route-maps
 - absence of either/both means all statically routed prefixes go into iBGP

Summary

□ Best Practices Covered:

- When to use BGP
- When to use ISIS/OSPF
- Aggregation
- Receiving Prefixes
- Prefixes into BGP

Configuration Tips



Of passwords, tricks and
templates

iBGP and IGP

Reminder!

- ❑ Make sure loopback is configured on router
 - iBGP between loopbacks, NOT real interfaces
- ❑ Make sure IGP carries loopback /32 address
- ❑ Consider the DMZ nets:
 - Use unnumbered interfaces?
 - Use next-hop-self on iBGP neighbours
 - Or carry the DMZ /30s in the iBGP
 - Basically keep the DMZ nets out of the IGP!

iBGP: Next-hop-self

- ❑ BGP speaker announces external network to iBGP peers using router's local address (loopback) as next-hop
- ❑ Used by many ISPs on edge routers
 - Preferable to carrying DMZ /30 addresses in the IGP
 - Reduces size of IGP to just core infrastructure
 - Alternative to using unnumbered interfaces
 - Helps scale network
 - Many ISPs consider this "best practice"

Limiting AS Path Length

- ❑ Some BGP implementations have problems with long AS_PATHS
 - Memory corruption
 - Memory fragmentation
- ❑ Even using AS_PATH prepends, it is not normal to see more than 20 ASes in a typical AS_PATH in the Internet today
 - The Internet is around 5 ASes deep on average
 - Largest AS_PATH is usually 16-20 ASNs

Limiting AS Path Length

- Some announcements have ridiculous lengths of AS-paths:

```
*> 3FFE:1600::/24          22 11537 145 12199 10318
    10566 13193 1930 2200 3425 293 5609 5430 13285 6939
    14277 1849 33 15589 25336 6830 8002 2042 7610 i
```

This example is an error in one IPv6 implementation

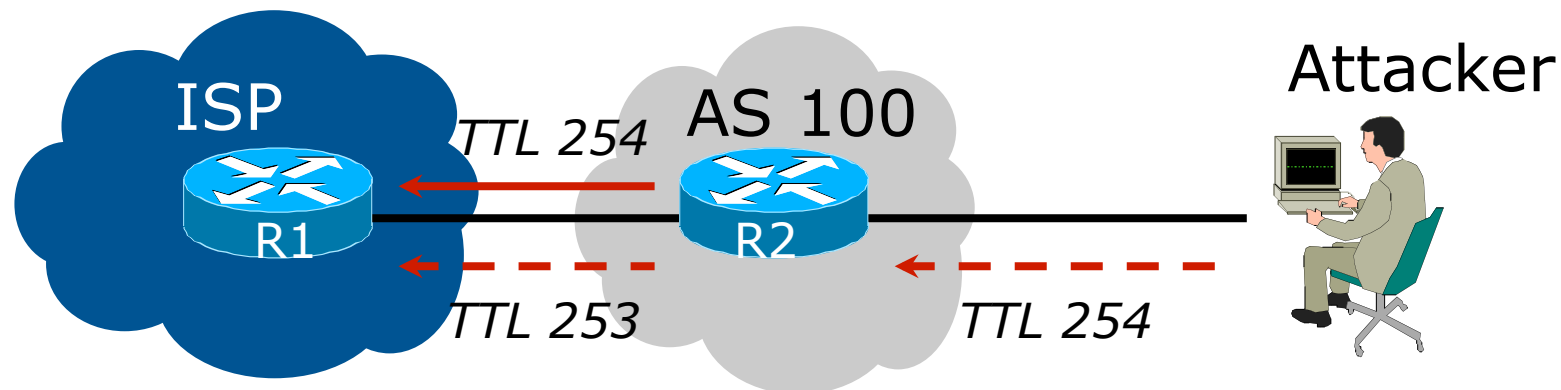
```
*> 96.27.246.0/24          2497 1239 12026 12026 12026
    12026 12026 12026 12026 12026 12026 12026 12026 12026
    12026 12026 12026 12026 12026 12026 12026 12026 12026
    12026 i
```

This example shows 21 prepends (for no obvious reason)

- If your implementation supports it, consider limiting the maximum AS-path length you will accept

BGP TTL “hack”

- ❑ Implement RFC5082 on BGP peerings
 - (Generalised TTL Security Mechanism)
 - Neighbour sets TTL to 255
 - Local router expects TTL of incoming BGP packets to be 254
 - No one apart from directly attached devices can send BGP packets which arrive with TTL of 254, so any possible attack by a remote miscreant is dropped due to TTL mismatch



BGP TTL “hack”

- TTL Hack:
 - Both neighbours must agree to use the feature
 - TTL check is much easier to perform than MD5
 - (Called BTSH – BGP TTL Security Hack)
- Provides “security” for BGP sessions
 - In addition to packet filters of course
 - MD5 should still be used for messages which slip through the TTL hack
 - See www.nanog.org/mtg-0302/hack.html for more details

Templates

- ❑ Good practice to configure templates for everything
 - Vendor defaults tend not to be optimal or even very useful for ISPs
 - ISPs create their own defaults by using configuration templates
- ❑ eBGP and iBGP examples follow
 - Also see Team Cymru's BGP templates
 - ❑ <http://www.team-cymru.org/ReadingRoom/Documents/>

iBGP Template

Example

- ❑ iBGP between loopbacks!
- ❑ Next-hop-self
 - Keep DMZ and external point-to-point out of IGP
- ❑ Always send communities in iBGP
 - Otherwise accidents will happen
- ❑ Hardwire BGP to version 4
 - Yes, this is being paranoid!

iBGP Template

Example continued

- ❑ Use passwords on iBGP session
 - Not being paranoid, **VERY** necessary
 - It's a secret shared between you and your peer
 - If arriving packets don't have the correct MD5 hash, they are ignored
 - Helps defeat miscreants who wish to attack BGP sessions
- ❑ Powerful preventative tool, especially when combined with filters and the TTL "hack"

eBGP Template

Example

- ❑ BGP damping
 - Do **NOT** use it unless you understand the impact
 - Do **NOT** use the vendor defaults without thinking
- ❑ Remove private ASes from announcements
 - Common omission today
- ❑ Use extensive filters, with “backup”
 - Use as-path filters to backup prefix filters
 - Keep policy language for implementing policy, rather than basic filtering
- ❑ Use password agreed between you and peer on eBGP session

eBGP Template

Example continued

- ❑ Use maximum-prefix tracking
 - Router will warn you if there are sudden increases in BGP table size, bringing down eBGP if desired
- ❑ Limit maximum as-path length inbound
- ❑ Log changes of neighbour state
 - ...and monitor those logs!
- ❑ Make BGP admin distance higher than that of any IGP
 - Otherwise prefixes heard from outside your network could override your IGP!!

Summary

- ❑ Use configuration templates
- ❑ Standardise the configuration
- ❑ Be aware of standard “tricks” to avoid compromise of the BGP session
- ❑ Anything to make your life easier, network less prone to errors, network more likely to scale
- ❑ It's all about scaling – if your network won't scale, then it won't be successful

BGP Best Current Practices



ISP Workshops