



# BGP Multihoming Techniques

Philip Smith <pfs@cisco.com>

AfNOG 2011

Dar Es Salaam, Tanzania

5 June 2011

# Presentation Slides

- Available on

**<ftp://ftp-eng.cisco.com>**

**[/pfs/seminars/AfNOG2011-BGP-Multihoming.pdf](#)**

And on the AfNOG2011 website

- Feel free to ask questions any time

# Preliminaries

- Presentation has many configuration examples
  - Uses Cisco IOS CLI
- Aimed at Service Providers
  - Techniques can be used by many enterprises too

# BGP Multihoming Techniques

- **Why Multihome?**
- Definition & Options
- How to Multihome
- Preparing the Network
- Basic Multihoming
- Service Provider Multihoming
- Complex Cases & Caveats
- Using Communities
- Case Study



## Why Multihome?

**It's all about redundancy, diversity & reliability**

# Why Multihome?

- Redundancy

One connection to internet means the network is dependent on:

Local router (configuration, software, hardware)

WAN media (physical failure, carrier failure)

Upstream Service Provider (configuration, software, hardware)

# Why Multihome?

- Reliability

Business critical applications demand continuous availability

Lack of redundancy implies lack of reliability implies loss of revenue

# Why Multihome?

- Supplier Diversity

Many businesses demand supplier diversity as a matter of course

Internet connection from two or more suppliers

With two or more diverse WAN paths

With two or more exit points

With two or more international connections

**Two of everything**



# Why Multihome?

- Not really a reason, but oft quoted...
- Leverage:
  - Playing one ISP off against the other for:
    - Service Quality
    - Service Offerings
    - Availability

# Why Multihome?

- Summary:

Multihoming is easy to demand as requirement for any service provider or end-site network

But what does it really mean:

In real life?

For the network?

For the Internet?

And how do we do it?

# BGP Multihoming Techniques

- Why Multihome?
- **Definition & Options**
- How to Multihome
- Preparing the Network
- Basic Multihoming
- Service Provider Multihoming
- Complex Cases & Caveats
- Using Communities
- Case Study



# Multihoming: Definitions & Options

What does it mean, what do we need, and how do we do it?

# Multihoming Definition

- More than one link external to the local network
  - two or more links to the same ISP
  - two or more links to different ISPs
- Usually **two** external facing routers
  - one router gives link and provider redundancy only

# Autonomous System Number (ASN)

- Two ranges

0-65535	(original 16-bit range)
65536-4294967295	(32-bit range - RFC4893)

- Usage:

0 and 65535	(reserved)
1-64495	(public Internet)
64496-64511	(documentation - RFC5398)
64512-65534	(private use only)
23456	(represent 32-bit range in 16-bit world)
65536-65551	(documentation - RFC5398)
65552-4294967295	(public Internet)

- 32-bit range representation specified in RFC5396

Defines “asplain” (traditional format) as standard notation

# Autonomous System Number (ASN)

- ASNs are distributed by the Regional Internet Registries

They are also available from upstream ISPs who are members of one of the RIRs

Around 37500 are visible on the Internet

- Current 16-bit ASN allocations up to 58367 have been made to the RIRs
- Each RIR has also received a block of 32-bit ASNs
  - Out of 1400 assignments, around 1100 are visible on the Internet
- See [www.iana.org/assignments/as-numbers](http://www.iana.org/assignments/as-numbers)

# Private-AS – Application

- Applications

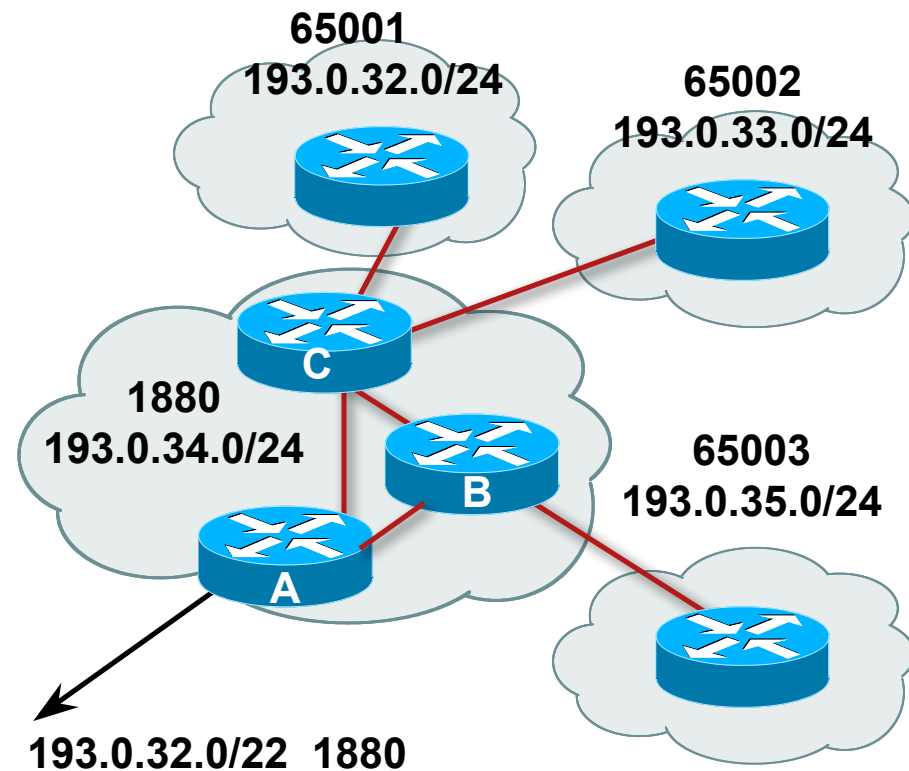
An ISP with customers multihomed on their backbone (RFC2270)

-or-

A corporate network with several regions but connections to the Internet only in the core

-or-

Within a BGP Confederation





# Private-AS – Removal

- Private ASNs MUST be removed from all prefixes announced to the public Internet
  - Include configuration to remove private ASNs in the eBGP template
- As with RFC1918 address space, private ASNs are intended for internal use
  - They should not be leaked to the public Internet
- Cisco IOS
  - `neighbor x.x.x.x remove-private-AS`

# Transit/Peering/Default

- Transit

Carrying traffic across a network

Usually **for a fee**

- Peering

Exchanging locally sourced routing information and traffic

Usually **for no fee**

Sometimes called **settlement free peering**

- Default

Where to send traffic when there is no explicit match in the routing table

# Configuring Policy

- Three BASIC Principles for IOS configuration examples throughout presentation:
  - `prefix-lists` to filter `prefixes`
  - `filter-lists` to filter `ASNs`
  - `route-maps` to apply `policy`
- Route-maps can be used for filtering, but this is more “advanced” configuration

# Policy Tools

- Local preference  
outbound traffic flows
- Metric (MED)  
inbound traffic flows (local scope)
- AS-PATH prepend  
inbound traffic flows (Internet scope)
- Communities  
specific inter-provider peering

# Originating Prefixes: Assumptions

- **MUST** announce assigned address block to Internet
- MAY also announce subprefixes – reachability is not guaranteed
- Current minimum allocation is from /20 to /24 depending on the RIR

Several ISPs filter RIR blocks on this boundary

Several ISPs filter the rest of address space according to the IANA assignments

This activity is called “Net Police” by some

# Originating Prefixes

- The RIRs publish their minimum allocation sizes per /8 address block

AfriNIC: [www.afrinic.net/docs/policies/afpol-v4200407-000.htm](http://www.afrinic.net/docs/policies/afpol-v4200407-000.htm)

APNIC: [www.apnic.net/db/min-alloc.html](http://www.apnic.net/db/min-alloc.html)

ARIN: [www.arin.net/reference/ip\\_blocks.html](http://www.arin.net/reference/ip_blocks.html)

LACNIC: [lacnic.net/en/registro/index.html](http://lacnic.net/en/registro/index.html)

RIPE NCC: [www.ripe.net/ripe/docs/smallest-alloc-sizes.html](http://www.ripe.net/ripe/docs/smallest-alloc-sizes.html)

Note that AfriNIC only publishes its current minimum allocation size, not the allocation size for its address blocks

- IANA publishes the address space it has assigned to end-sites and allocated to the RIRs:

[www.iana.org/assignments/ipv4-address-space](http://www.iana.org/assignments/ipv4-address-space)

- Several ISPs use this published information to filter prefixes on:

What should be routed (from IANA)

The minimum allocation size from the RIRs

## “Net Police” prefix list issues

- Meant to “punish” ISPs who pollute the routing table with specifics rather than announcing aggregates
- Impacts legitimate multihoming especially at the Internet’s edge
- Impacts regions where domestic backbone is unavailable or costs \$\$\$ compared with international bandwidth
- Hard to maintain – requires updating when RIRs start allocating from new address blocks
- Don’t do it unless consequences understood and you are prepared to keep the list current

Consider using the Team Cymru or other reputable bogon BGP feed:

<http://www.team-cymru.org/Services/Bogons/routeserver.html>

# BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- **How to Multihome**
- Preparing the Network
- Basic Multihoming
- Service Provider Multihoming
- Complex Cases & Caveats
- Using Communities
- Case Study





# How to Multihome

Choosing between transit and peer

# Transits

- Transit provider is another autonomous system which is used to provide the local network with access to other networks

Might be local or regional only

But more usually the whole Internet

- Transit providers need to be chosen wisely:

Only one                      no redundancy

Too many                      more difficult to load balance

no economy of scale (costs more per Mbps)

hard to provide service quality

- **Recommendation: at least two, no more than three**

# Common Mistakes

- ISPs sign up with too many transit providers

  - Lots of small circuits (cost more per Mbps than larger ones)

  - Transit rates per Mbps reduce with increasing transit bandwidth purchased

  - Hard to implement reliable traffic engineering that doesn't need daily fine tuning depending on customer activities

- No diversity

  - Chosen transit providers all reached over same satellite or same submarine cable

  - Chosen transit providers have poor onward transit and peering

# Peers

- A peer is another autonomous system with which the local network has agreed to exchange locally sourced routes and traffic
- Private peer
  - Private link between two providers for the purpose of interconnecting
- Public peer
  - Internet Exchange Point, where providers meet and freely decide who they will interconnect with
- **Recommendation: peer as much as possible!**

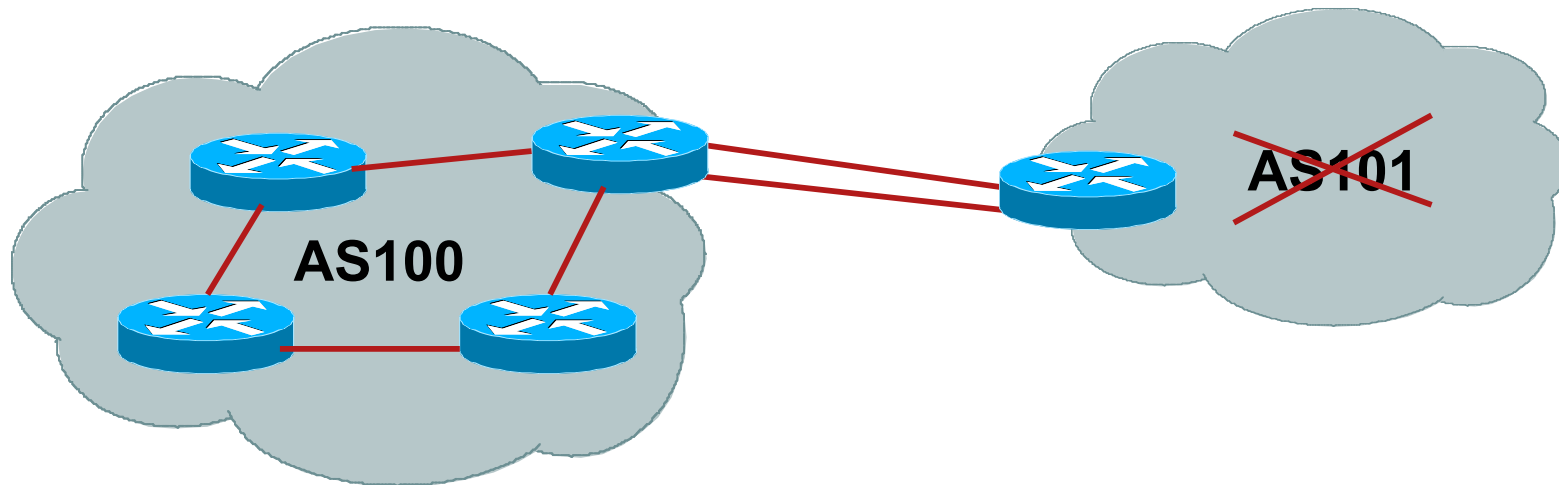
# Common Mistakes

- Mistaking a transit provider's "Exchange" business for a no-cost public peering point
- Not working hard to get as much peering as possible
  - Physically near a peering point (IXP) but not present at it  
(Transit sometimes is cheaper than peering!!)
- Ignoring/avoiding competitors because they are competition
  - Even though potentially valuable peering partner to give customers a better experience

# Multihoming Scenarios

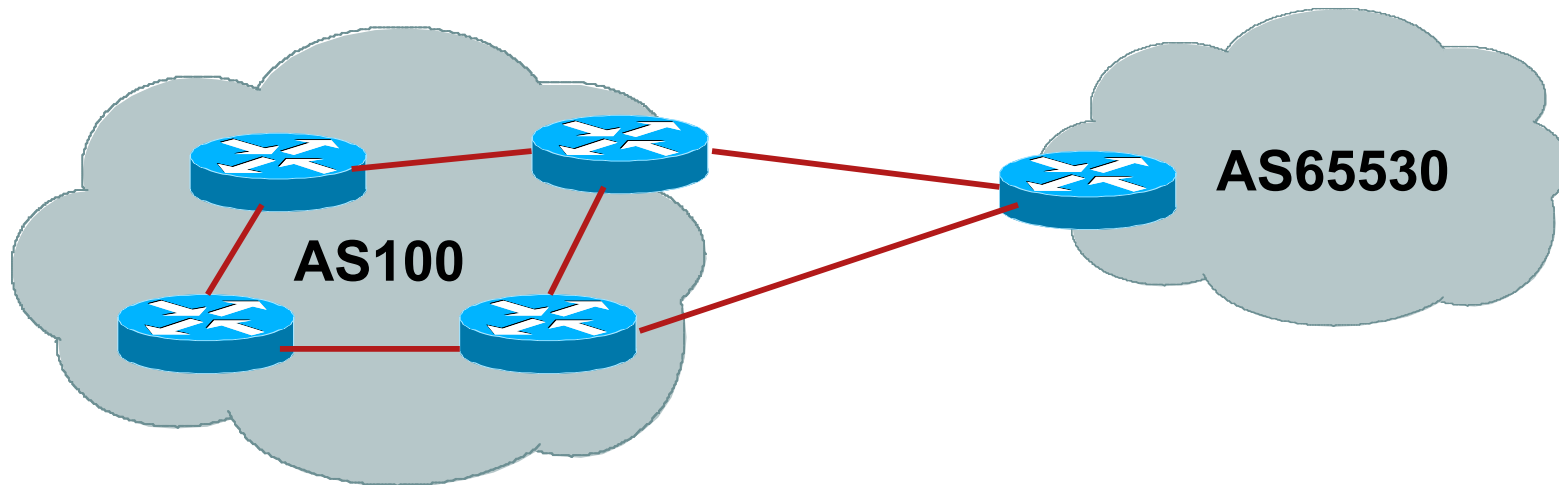
- Stub network
- Multi-homed stub network
- Multi-homed network
- Multiple sessions to another AS

# Stub Network



- No need for BGP
- Point static default to upstream ISP
- Upstream ISP advertises stub network
- Policy confined within upstream ISP's policy

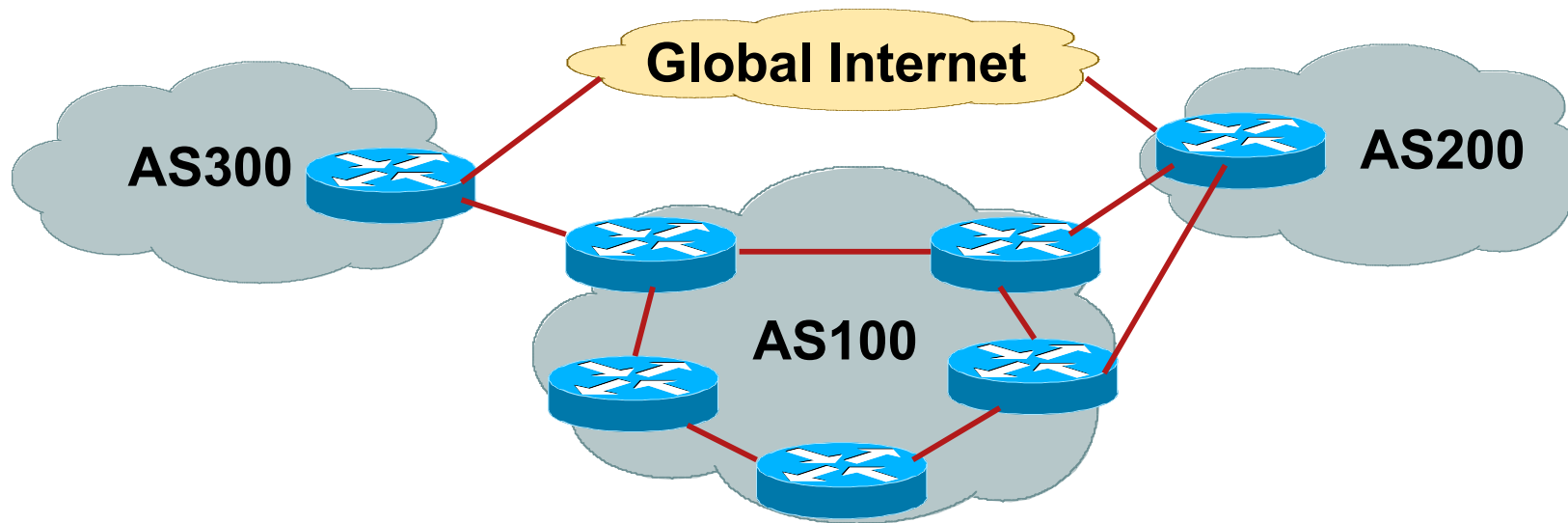
# Multi-homed Stub Network



- Use BGP (not IGP or static) to loadshare
- Use private AS (ASN > 64511)
- Upstream ISP advertises stub network
- Policy confined within upstream ISP's policy



# Multi-homed Network



- Many situations possible
  - multiple sessions to same ISP
  - secondary for backup only
  - load-share between primary and secondary
  - selectively use different ISPs

# Multiple Sessions to an AS

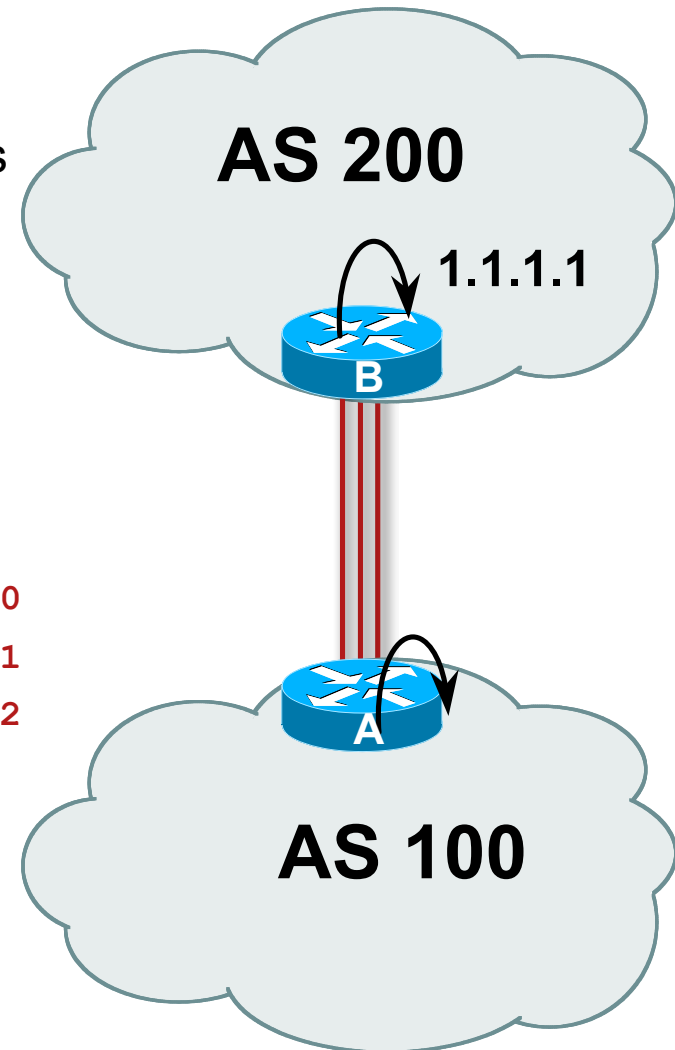
## – ebgp multihop

- Run eBGP between loopback addresses  
eBGP prefixes learned with loopback address  
as next hop

- Cisco IOS

```
router bgp 100
  neighbor 1.1.1.1 remote-as 200
  neighbor 1.1.1.1 ebgp-multihop 2
  !
  ip route 1.1.1.1 255.255.255.255 serial 1/0
  ip route 1.1.1.1 255.255.255.255 serial 1/1
  ip route 1.1.1.1 255.255.255.255 serial 1/2
```

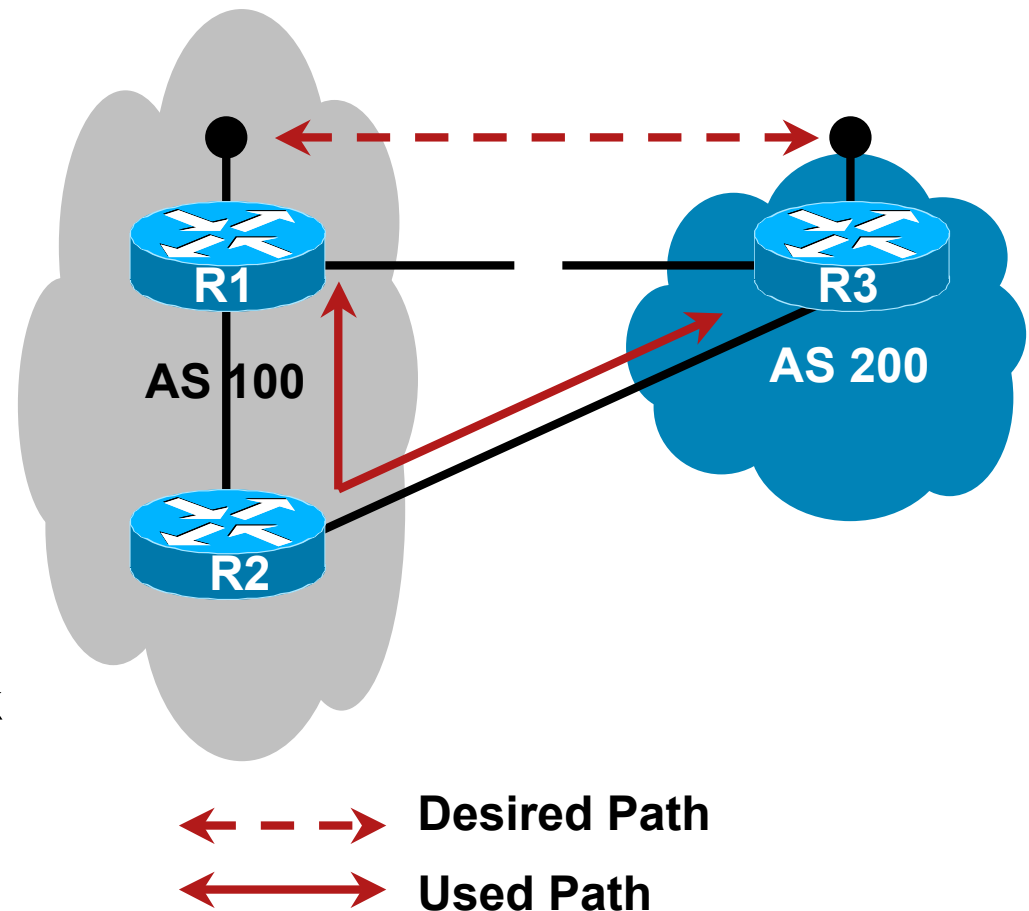
- Common error made is to point remote  
loopback route at IP address rather than  
specific link



# Multiple Sessions to an AS

## – ebgp multihop

- One eBGP-multihop gotcha:  
R1 and R3 are eBGP peers that are loopback peering  
Configured with:  
`neighbor x.x.x.x ebgp-multihop 2`  
If the R1 to R3 link goes down the session could establish via R2
- Usually happens when routing to remote loopback is dynamic, rather than static pointing at a link



# Multiple Sessions to an AS

## – ebgp multihop

- Avoid the use of ebgp-multihop unless:
  - There is simply no alternative
  - or–
  - Loadsharing across multiple parallel links
- Many ISPs discourage its use, for example:

**We will run eBGP multihop, but do not support it as a standard offering because customers generally have a hard time managing it due to:**

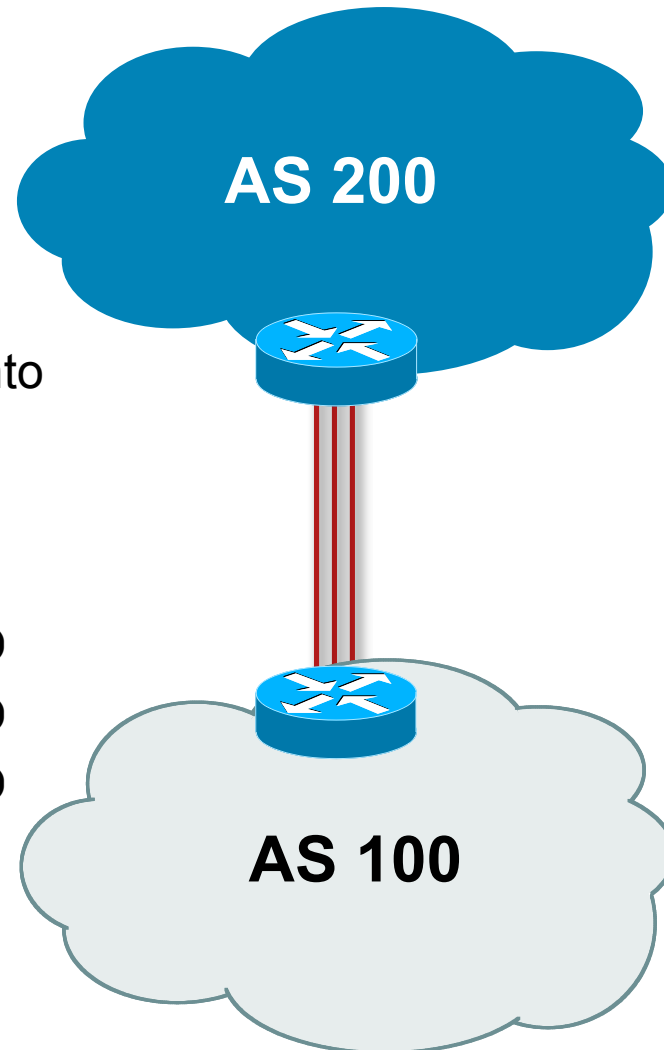
- routing loops
- failure to realise that BGP session stability problems are usually due connectivity problems between their CPE and their BGP speaker

# Multiple Sessions to an AS

## – bgp multi path

- Three BGP sessions required
- Platform limit on number of paths (could be as little as 6)
- Full BGP feed makes this unwieldy  
3 copies of Internet Routing Table goes into the FIB

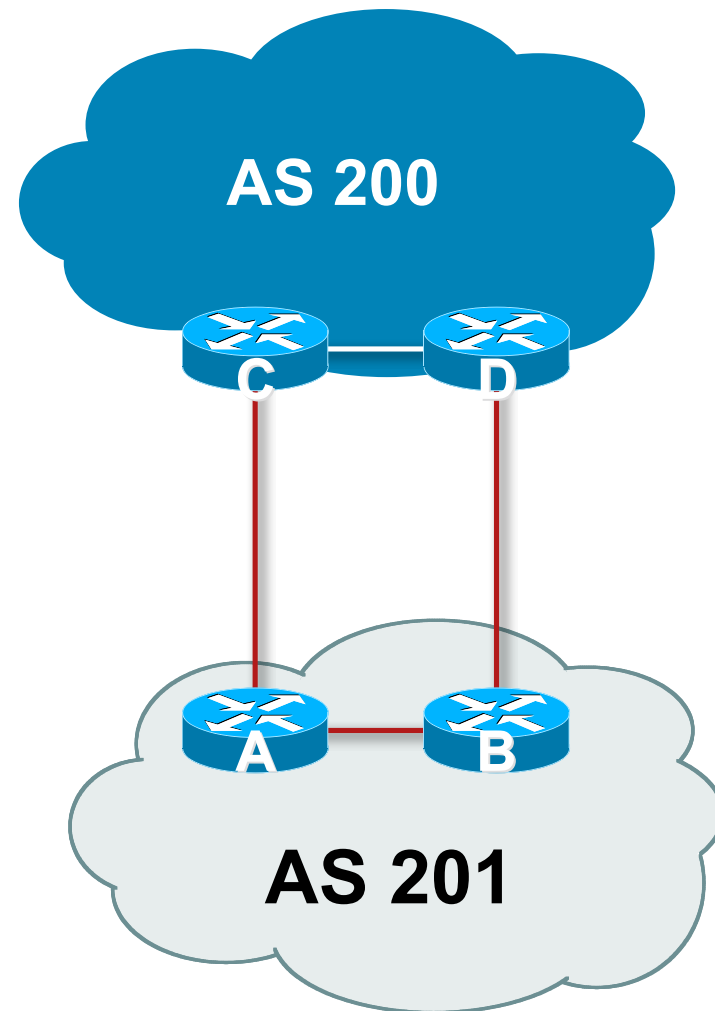
```
router bgp 100
  neighbor 1.1.2.1 remote-as 200
  neighbor 1.1.2.5 remote-as 200
  neighbor 1.1.2.9 remote-as 200
  maximum-paths 3
```



# Multiple Sessions to an AS

## – bgp attributes & filters

- Simplest scheme is to use defaults
- Learn/advertise prefixes for better control
- Planning and some work required to achieve loadsharing
  - Point default towards one ISP
  - Learn selected prefixes from second ISP
  - Modify the number of prefixes learnt to achieve acceptable load sharing
- **No magic solution**





# Multihoming

# Inbound Load Balancing

- Traffic comes into a network because of the address spaced announced by that network

Loadbalancing is achieved by manipulating outbound announcements

- How?

Announcing aggregate (an all links, always)

Carefully leaking a subprefix or two of that aggregate

Varying the size of leaked subprefixes

Using AS-PATH prepend carefully



# Outbound Load Balancing

- Traffic goes out of a network based on the addresses announced into it

Load balancing is achieved by manipulating these inbound routing announcements

- Achieved by:

Default route from one upstream, full table plus default from the other upstream

And then throwing most of the latter away until traffic is balanced

“Throwing” away achieved by selective AS-PATH filtering

Also the use of local-preference on selective paths

# Multihoming

- Inbound and Outbound Load balancing is known as “Traffic Engineering”
- Configuration examples covered in the AR-E Workshop  
Consult the AR-E Workshop materials



# Using Communities for Multihoming

# Multihoming and Communities

- The BGP community attribute is a very powerful tool for assisting and scaling BGP Multihoming
- Most major ISPs make extensive use of BGP communities:
  - Internal policies
  - Inter-provider relationships (MED replacement)
  - Customer traffic engineering

# Using BGP Communities

- Four scenarios are covered:
  - Use of RFC1998 traffic engineering
  - Extending RFC 1998 ideas for even greater customer policy options
  - Community use in ISP backbones
  - Customer Policy Control (aka traffic engineering)



# RFC1998

An example of how ISPs use communities...

# RFC1998

- Informational RFC
- Describes how to implement loadsharing and backup on multiple inter-AS links
  - BGP communities used to determine local preference in upstream's network
- Gives control to the customer
  - Means the customer does not have to phone upstream's technical support to adjust traffic engineering needs
- Simplifies upstream's configuration
  - simplifies network operation!

# RFC1998

- RFC1998 Community values are defined to have particular meanings
- ASx:100                    set local preference 100  
    Make this the preferred path
- ASx :90                    set local preference 90  
    Make this the backup if dualhomed on ASx
- ASx :80                    set local preference 80  
    The main link is to another ISP with same AS path length
- ASx :70                    set local preference 70  
    The main link is to another ISP



# RFC1998

- Upstream ISP defines the communities mentioned
- Their customers then attach the communities they want to use to the prefix announcements they are making
- For example:
  - If upstream is AS 100
  - To declare a particular path as a backup path, their customer would announce the prefix with community 100:70 to AS100
  - AS100 would receive the prefix with the community 100:70 tag, and then set local preference to be 70

# RFC1998

- Sample Customer Router Configuration

```
router bgp 130
  neighbor x.x.x.x remote-as 100
  neighbor x.x.x.x description Backup ISP
  neighbor x.x.x.x route-map as100-out out
  neighbor x.x.x.x send-community
!
ip as-path access-list 20 permit ^$
!
route-map as100-out permit 10
  match as-path 20
  set community 100:70
!
```

# RFC1998

- Sample ISP Router Configuration

```
router bgp 100
  neighbor y.y.y.y remote-as 130
  neighbor y.y.y.y route-map customer-policy-in in
!
! Homed to another ISP
ip community-list 7 permit 100:70
! Homed to another ISP with equal ASPATH length
ip community-list 8 permit 100:80
! Customer backup routes
ip community-list 9 permit 100:90
!
```

# RFC1998

```
route-map customer-policy-in permit 10
  match community 7
  set local-preference 70
!
route-map customer-policy-in permit 20
  match community 8
  set local-preference 80
!
route-map customer-policy-in permit 30
  match community 9
  set local-preference 90
!
route-map customer-policy-in permit 40
  set local-preference 100
!
```

# RFC1998

- RFC1998 was the inspiration for a large variety of differing community policies implemented by ISPs worldwide
- There are no “standard communities” for what ISPs do
- But best practices today consider that ISPs should use BGP communities extensively for multihoming support of traffic engineering
- Look in the ISP AS Object in the IRR for documented community support



## Service Provider use of Communities

RFC1998 was so inspiring...

# Background

- RFC1998 is okay for “simple” multihoming situations
- ISPs create backbone support for many other communities to handle more complex situations
  - Simplify ISP BGP configuration
  - Give customer more policy control

# ISP BGP Communities

- There are no recommended ISP BGP communities apart from RFC1998  
The five standard communities  
[www.iana.org/assignments/bgp-well-known-communities](http://www.iana.org/assignments/bgp-well-known-communities)
- Efforts have been made to document from time to time  
[totem.info.ucl.ac.be/publications/papers-elec-versions/draft-quoitin-bgp-comm-survey-00.pdf](http://totem.info.ucl.ac.be/publications/papers-elec-versions/draft-quoitin-bgp-comm-survey-00.pdf)  
But so far... nothing more... ☹️  
Collection of ISP communities at [www.onesc.net/communities](http://www.onesc.net/communities)  
[www.nanog.org/meetings/nanog40/presentations/BGPcommunities.pdf](http://www.nanog.org/meetings/nanog40/presentations/BGPcommunities.pdf)
- ISP policy is usually published  
On the ISP's website  
Referenced in the AS Object in the IRR



# Typical ISP BGP Communities

- X:80            set local preference 80  
Backup path
- X:120          set local preference 120  
Primary path (over ride BGP path selection default)
- X:1            set as-path prepend X  
Single prepend when announced to X's upstreams
- X:2            set as-path prepend X X  
Double prepend when announced to X's upstreams
- X:3            set as-path prepend X X X  
Triple prepend when announced to X's upstreams
- X:666          set ip next-hop 192.0.2.1  
Blackhole route - very useful for DoS attack mitigation

# Sample Router Configuration (1)

```
router bgp 100
  neighbor y.y.y.y remote-as 130
  neighbor y.y.y.y route-map customer-policy-in in
  neighbor z.z.z.z remote-as 200
  neighbor z.z.z.z route-map upstream-out out
!
ip community-list 1 permit 100:1
ip community-list 2 permit 100:2
ip community-list 3 permit 100:3
ip community-list 4 permit 100:80
ip community-list 5 permit 100:120
ip community-list 6 permit 100:666
!
ip route 192.0.2.1 255.255.255.255 null0
```

Customer BGP



Upstream BGP



Black hole route  
(on all routers)



## Sample Router Configuration (2)

```
route-map customer-policy-in permit 10
  match community 4
  set local-preference 80
!
route-map customer-policy-in permit 20
  match community 5
  set local-preference 120
!
route-map customer-policy-in permit 30
  match community 6
  set ip next-hop 192.0.2.1
!
route-map customer-policy-in permit 40
...etc...
```

## Sample Router Configuration (3)

```
route-map upstream-out permit 10
  match community 1
  set as-path prepend 100
!
route-map upstream-out permit 20
  match community 2
  set as-path prepend 100 100
!
route-map upstream-out permit 30
  match community 3
  set as-path prepend 100 100 100
!
route-map upstream-out permit 40
...etc...
```

within 3 business days of receipt of the request.

## WHAT YOU CAN CONTROL

### AS-PATH PREPENDS

Sprint allows customers to use AS-path prepending to adjust route preference on the network. Such prepending will be received and passed on properly without notifying Sprint of your change in announcements.

Additionally, Sprint will prepend AS1239 to eBGP sessions with certain autonomous systems depending on a received community. Currently, the following ASes are supported: 1668, 209, 2914, 3300, 3356, 3549, 3561, 4635, 701, 7018, 702 and 8220.

String	Resulting AS Path to ASXXX
--------	----------------------------

65000:XXX	Do not advertise to ASXXX
-----------	---------------------------

65001:XXX	1239 (default) ...
-----------	--------------------

65002:XXX	1239 1239 ...
-----------	---------------

65003:XXX	1239 1239 1239 ...
-----------	--------------------

65004:XXX	1239 1239 1239 1239 ...
-----------	-------------------------

String	Resulting AS Path to ASXXX in Asia
--------	------------------------------------

65070:XXX	Do not advertise to ASXXX
-----------	---------------------------

65071:XXX	1239 (default) ...
-----------	--------------------

65072:XXX	1239 1239 ...
-----------	---------------

65073:XXX	1239 1239 1239 ...
-----------	--------------------

65074:XXX	1239 1239 1239 1239 ...
-----------	-------------------------

String	Resulting AS Path to ASXXX in Europe
--------	--------------------------------------

65050:XXX	Do not advertise to ASXXX
-----------	---------------------------

65051:XXX	1239 (default) ...
-----------	--------------------

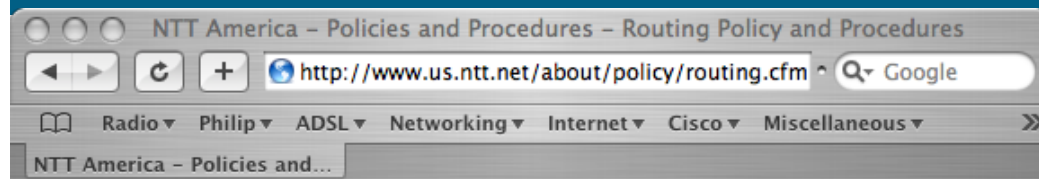
65052:XXX	1239 1239 ...
-----------	---------------

65053:XXX	1239 1239 1239 ...
-----------	--------------------

65054:XXX	1239 1239 1239 1239 ...
-----------	-------------------------

# ISP Examples: Sprint

More info at  
[https://www.sprint.net/index.php?p=policy\\_bgp](https://www.sprint.net/index.php?p=policy_bgp)



## Some ISP Examples: NTT

### BGP customer communities

#### Customers wanting to alter local preference on their routes.

NTT Communications BGP customers may choose to affect our local preference on their routes by marking their routes with the following communities:

Community	Local-pref	Description
(default)	120	customer
2914:450	96	customer fallback
2914:460	98	peer backup
2914:470	100	peer
2914:480	110	customer backup
2914:490	120	customer default

#### Customers wanting to alter their route announcements to other customers.

NTT Communications BGP customers may choose to prepend to all other NTT Communications BGP customers with the following communities:

Community	Description
2914:411	prepends o/b to customer 1x
2914:412	prepends o/b to customer 2x
2914:413	prepends o/b to customer 3x

#### Customers wanting to alter their route announcements to peers.

NTT Communications BGP customers may choose to prepend to all NTT Communications peers with the following communities:

Community	Description
2914:421	prepends o/b to peer 1x
2914:422	prepends o/b to peer 2x

More info at  
[www.us.ntt.net/about/policy/routing.cfm](http://www.us.ntt.net/about/policy/routing.cfm)

# ISP Examples:

## Verizon Business Europe

```
aut-num: AS702
descr: Verizon Business EMEA - Commercial IP service provider in Eur
remarks: VzBi uses the following communities with its customers:
        702:80      Set Local Pref 80 within AS702
        702:120     Set Local Pref 120 within AS702
        702:20      Announce only to VzBi AS'es and VzBi customers
        702:30      Keep within Europe, don't announce to other VzBi AS
        702:1       Prepend AS702 once at edges of VzBi to Peers
        702:2       Prepend AS702 twice at edges of VzBi to Peers
        702:3       Prepend AS702 thrice at edges of VzBi to Peers
Advanced communities for customers
        702:7020     Do not announce to AS702 peers with a scope of
                    National but advertise to Global Peers, European
                    Peers and VzBi customers.
        702:7001     Prepend AS702 once at edges of VzBi to AS702
                    peers with a scope of National.
        702:7002     Prepend AS702 twice at edges of VzBi to AS702
                    peers with a scope of National.
(more)
```

# ISP Examples:

## Verizon Business Europe

(more)

```
702:7003 Prepend AS702 thrice at edges of VzBi to AS702
        peers with a scope of National.
702:8020 Do not announce to AS702 peers with a scope of
        European but advertise to Global Peers, National
        Peers and VzBi customers.
702:8001 Prepend AS702 once at edges of VzBi to AS702
        peers with a scope of European.
702:8002 Prepend AS702 twice at edges of VzBi to AS702
        peers with a scope of European.
702:8003 Prepend AS702 thrice at edges of VzBi to AS702
        peers with a scope of European.
```

-----  
Additional details of the VzBi communities are located at:  
<http://www.verizonbusiness.com/uk/customer/bgp/>  
-----


```
mnt-by: WCOM-EMEA-RICE-MNT
source: RIPE
```



# Some ISP Examples


## BT Ignite

```
aut-num:      AS5400
descr:        BT Ignite European Backbone
remarks:
remarks:      Community to
remarks:      Not announce      To peer:      Community to
remarks:                                             AS prepend 5400
remarks:      5400:1000 All peers & Transits      5400:2000
remarks:
remarks:      5400:1500 All Transits      5400:2500
remarks:      5400:1501 Sprint Transit (AS1239)      5400:2501
remarks:      5400:1502 SAVVIS Transit (AS3561)      5400:2502
remarks:      5400:1503 Level 3 Transit (AS3356)      5400:2503
remarks:      5400:1504 AT&T Transit (AS7018)      5400:2504
remarks:      5400:1506 GlobalCrossing Trans (AS3549) 5400:2506
remarks:
remarks:      5400:1001 Nexica (AS24592)      5400:2001
remarks:      5400:1002 Fujitsu (AS3324)      5400:2002
remarks:      5400:1004 C&W EU (1273)      5400:2004
<snip>
notify:       notify@eu.bt.net
mnt-by:       CIP-MNT
source:       RIPE
```



## Some ISP Examples Level 3

```
aut-num:      AS3356
descr:        Level 3 Communications
<snip>
remarks:      -----
remarks:      customer traffic engineering communities - Suppression
remarks:      -----
remarks:      64960:XXX - announce to AS XXX if 65000:0
remarks:      65000:0   - announce to customers but not to peers
remarks:      65000:XXX - do not announce at peerings to AS XXX
remarks:      -----
remarks:      customer traffic engineering communities - Prepending
remarks:      -----
remarks:      65001:0   - prepend once  to all peers
remarks:      65001:XXX - prepend once  at peerings to AS XXX
<snip>
remarks:      3356:70   - set local preference to 70
remarks:      3356:80   - set local preference to 80
remarks:      3356:90   - set local preference to 90
remarks:      3356:9999 - blackhole (discard) traffic
<snip>
mnt-by:        LEVEL3-MNT
source:        RIPE
```



And many  
many more!

# Creating your own community policy

- Consider creating communities to give policy control to customers
  - Reduces technical support burden
  - Reduces the amount of router reconfiguration, and the chance of mistakes
  - Use previous ISP and configuration examples as a guideline



# Using Communities for Backbone Scaling

Scaling BGP in the ISP backbone...

# Communities for iBGP

- ISPs tag prefixes learned from their BGP and static customers with communities
  - To identify services the customer may have purchased
  - To identify prefixes which are part of the ISP's PA space
  - To identify PI customer addresses
  - To control prefix distribution in iBGP
  - To control prefix announcements to customers and upstreams (amongst several other reasons)

# Service Identification

- ISP provides:
  - Transit via upstreams
  - Connectivity via major IXP
  - Connectivity to private peers/customers
- Customers can buy all or any of the above access options
  - Each option is identified with a unique community
- ISP identifies whether address space comes from their PA block or is their customers' own PI space
  - One community for each

# Community Definitions

100:1000	AS100 aggregates
100:1001	AS100 aggregate subprefixes
100:1005	Static Customer PI space
100:2000	Customers who get Transit
100:2100	Customers who get IXP access
100:2200	Customers who get BGP Customer access
100:3000	Routes learned from the IXP

```
ip community-list 10 permit 100:1000
ip community-list 11 permit 100:1001
ip community-list 12 permit 100:1005
ip community-list 13 permit 100:2000
ip community-list 14 permit 100:2100
ip community-list 15 permit 100:2200
ip community-list 16 permit 100:3000
```

# Aggregates and Static Customers into BGP

```
router bgp 100
  network 100.10.0.0 mask 255.255.224.0 route-map as100-prefixes
  redistribute static route-map static-to-bgp
!
ip prefix-list as100-block permit 100.10.0.0/19 le 32
!
route-map as100-prefixes permit 10
  set community 100:1000
!
route-map static-to-bgp permit 10
  match ip address prefix-list as100-block
  set community 100:1001
route-map static-to-bgp permit 20
  set community 100:1005
```

Aggregate community set

Aggregate subprefixes community set

PI community is set



# Service Identification

- AS100 has four classes of BGP customers
  - Full transit (upstream, IXP and BGP customers)
  - Upstream only
  - IXP only
  - BGP Customers only
- For BGP support, easiest IOS configuration is to create a peer-group for each class (can also use peer-templates to simplify further)
  - Customer is assigned the peer-group of the service they have purchased
  - Simple for AS100 customer installation engineer to provision

# BGP Customers - creating peer-groups

```
router bgp 100
  neighbor full-transit peer-group
  neighbor full-transit route-map customers-out out
  neighbor full-transit route-map full-transit-in in
  neighbor full-transit default-originate
  neighbor transit-up peer-group
  neighbor transit-up route-map customers-out out
  neighbor transit-up route-map transit-up-in in
  neighbor transit-up default-originate
  neighbor ixp-only peer-group
  neighbor ixp-only route-map ixp-routes out
  neighbor ixp-only route-map ixp-only-in in
  neighbor bgpcust-only peer-group
  neighbor bgpcust-only route-map bgp-cust-out out
  neighbor bgpcust-only route-map bgp-cust-in in
```

# BGP Customers - creating route-maps

```
route-map customers-out permit 10  
  match ip community 10
```

Customers only get AS100  
aggregates and default route

```
route-map full-transit-in permit 10  
  set community 100:2000 100:2100 100:2200
```

```
route-map transit-up-in permit 10  
  set community 100:2000
```

Full transit go everywhere

```
route-map ixp-routes permit 10  
  match ip community 10 12 13 14 16
```

Customers buying IXP access  
only get aggregates, static & full  
transit customers and IXP routes

```
route-map ixp-only-in permit 10  
  set community 100:2100
```

```
route-map bgp-cust-out permit 10  
  match ip community 10 12 13 15
```

```
route-map bgp-cust-in permit 10  
  set community 100:2200
```

Customers buying BGP customer  
access only get aggregates,  
static & full transit customers  
and other BGP customers

# BGP Customers - configuring customers

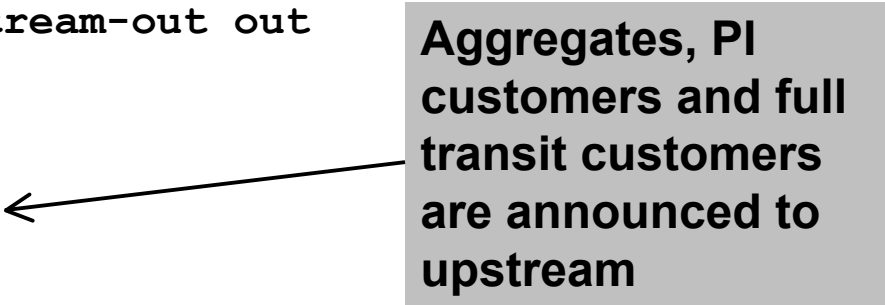
```
router bgp 100
  neighbor a.a.a.a remote-as 200
  neighbor a.a.a.a peer-group full-transit
  neighbor a.a.a.a prefix-list as200cust-in
  neighbor b.b.b.b remote-as 300
  neighbor b.b.b.b peer-group transit-up
  neighbor b.b.b.b prefix-list as300cust-in
  neighbor c.c.c.c remote-as 400
  neighbor c.c.c.c peer-group ixp-only
  neighbor c.c.c.c prefix-list as400cust-in
  neighbor d.d.d.d remote-as 500
  neighbor d.d.d.d peer-group bgpcust-only
  neighbor d.d.d.d prefix-list as500cust-in
```

Customers are simply dropped into the appropriate peer-group depending on the service they paid for

Note the specific per-customer inbound filters

# BGP Customers - configuring upstream

```
router bgp 100
  neighbor x.x.x.x remote-as 130
  neighbor x.x.x.x prefix-list full-routes in
  neighbor x.x.x.x route-map upstream-out out
!
route-map upstream-out permit 10
  match ip community 10 12 13
!
! IP prefix-list full-routes is the standard bogon
! prefix filter - or use a reputable bogon route-service such
! as that offered by Team Cymru
```



# BGP Customers - configuring IXP peers

```
router bgp 100
  neighbor y.y.y.1 remote-as 901
  neighbor y.y.y.1 route-map ixp-peers-out out
  neighbor y.y.y.1 route-map ixp-peers-in in
  neighbor y.y.y.1 prefix-list AS901-peer in
  neighbor y.y.y.2 remote-as 902
  neighbor y.y.y.2 route-map ixp-peers-out out
  neighbor y.y.y.2 route-map ixp-peers-in in
  neighbor y.y.y.2 prefix-list AS902-peer in
!
route-map ixp-peers-out permit 10
  match ip community 10 12 13 14
!
route-map ixp-peers-in permit 10
  set community 100:3000
```

Aggregates, PI  
customers full  
transit and IXP  
customers are  
announced to  
the IXP



# Service Identification

- While the community set up takes a bit of thought and planning, once it is implemented:
  - eBGP configuration with customers is simply a case of applying the appropriate peer-group
  - eBGP configuration with IXP peers is simply a case of announcing the appropriate community members to the peers
  - eBGP configuration with upstreams is simply a case of announcing the appropriate community members to the upstreams
- All BGP policy internally is now controlled by communities
  - No prefix-lists, as-path filters, route-maps or other BGP gymnastics are required

# What about iBGP itself?

- We've made good use of communities to handle customer requirements

But what about iBGP

- Most ISPs deploy Route Reflectors as a means of scaling iBGP
- In transit networks:
  - Core routers (the Route Reflectors) carry the full BGP table
  - Edge/Aggregation routers carry domestic prefixes & customers



# iBGP core router/route reflector

```
router bgp 100
  neighbor rrc peer-group
  neighbor rrc descr Route Reflector Clients
  neighbor rrc remote-as 100
  neighbor rrc route-map ibgp-filter out
  neighbor rrc send-community
  neighbor ibgp-peer peer-group
  neighbor ibgp-peer Standard iBGP peers
  neighbor ibgp-peer remote-as 100
  neighbor ibgp-peer send-community
  neighbor n.n.n.a peer-group ibgp-peer
  neighbor n.n.n.b peer-group rrc
!
route-map ibgp-filter permit 10
  match community 10 11 12 13 14 15 16
!
```

The filter to restrict client iBGP to just domestic prefixes

Must NOT forget to send community to iBGP peers

Allow all prefixes coming from the domestic network & IXP

## iBGP in the core

- Notice that the filtering of iBGP from the core to the edge is again achieved by a simple route-map applying a community match

No prefix-lists, as-path filters or any other complicated policy

Once the prefix belongs to a certain community, it has the access across the backbone determined by the community policy in force



# Using Communities for Customers Policy

Giving policy control to customers...

# Customer Policy Control

- ISPs have a choice on how to handle policy control for customers
- No delegation of policy options:
  - Customer has no choices
  - If customer wants changes, ISP Technical Support handles it
- Limited delegation of policy options:
  - Customer has choices
  - ISP Technical Support does not need to be involved
- BGP Communities are the only viable way of offering policy control to customers

# Policy Definitions

- Typical definitions:

Nil	No community set, just announce everywhere
X:1	1x prepend to all BGP neighbours
X:2	2x prepend to all BGP neighbours
X:3	3x prepend to all BGP neighbours
X:80	Local pref 80 on customer prefixes
X:120	Local pref 120 on customer prefixes
X:666	Black hole this route please!
X:5000	Don't announce to any BGP neighbour
X:5AA0	Don't announce to BGP neighbour AA
X:5AAB	Prepend B times to BGP neighbour AA

# Policy Implementation

- The BGP configuration for the initial communities was discussed at the start of this slide set
- But the new communities, X:5MMN, are worth covering in more detail

The ISP in AS X documents the BGP transits and peers that they have (MM can be 01 to 99)

The ISP in AS X indicates how many prepends they will support (N can be 1 to 9, but realistically 4 prepends is usually enough on today's Internet)

Customers then construct communities to do the prepending or announcement blocking they desire

- If a customer tags a prefix announcement with:  
100:5030 don't send prefix to BGP neighbour 03  
100:5102 2x prepend prefix announcement to peer 10

# Community Definitions

- Example: ISP in AS 100 has two upstreams. They create policy based on previously slide to allow no announce and up to 3 prepends for their customers

```
ip community-list 100 permit 100:5000
ip community-list 101 permit 100:5001
ip community-list 102 permit 100:5002
ip community-list 103 permit 100:5003
ip community-list 110 permit 100:5010
ip community-list 111 permit 100:5011
ip community-list 112 permit 100:5012
ip community-list 113 permit 100:5013
ip community-list 120 permit 100:5020
ip community-list 121 permit 100:5021
ip community-list 122 permit 100:5022
ip community-list 123 permit 100:5023
```

← Don't announce anywhere

← Single prepend to all

← Don't announce to peer 1

← Single prepend to peer 2

# Creating route-maps - neighbour 1

```
route-map bgp-neigh-01 deny 10  
  match ip community 100 110
```

Don't announce these  
prefixes to neighbour 01

!

```
route-map bgp-neigh-01 permit 20  
  match ip community 101 111  
  set as-path prepend 100
```

Single prepend of these  
prefixes to neighbour 01

!

```
route-map bgp-neigh-01 permit 30  
  match ip community 102 112  
  set as-path prepend 100 100
```

Double prepend of these  
prefixes to neighbour 01

!

```
route-map bgp-neigh-01 permit 40  
  match ip community 103 113  
  set as-path prepend 100 100 100
```

Triple prepend of these  
prefixes to neighbour 01

!

```
route-map bgp-neigh-01 permit 50
```

All other prefixes  
remain untouched



## Creating route-maps - neighbour 2

```
route-map bgp-neigh-02 deny 10
  match ip community 100 120
```

Don't announce these  
prefixes to neighbour 02

!

```
route-map bgp-neigh-02 permit 20
  match ip community 101 121
  set as-path prepend 100
```

Single prepend of these  
prefixes to neighbour 02

!

```
route-map bgp-neigh-02 permit 30
  match ip community 102 122
  set as-path prepend 100 100
```

Double prepend of these  
prefixes to neighbour 02

!

```
route-map bgp-neigh-02 permit 40
  match ip community 103 123
  set as-path prepend 100 100 100
```

Triple prepend of these  
prefixes to neighbour 02

!

```
route-map bgp-neigh-02 permit 50
```

All other prefixes  
remain untouched

# ISP's BGP configuration

```
router bgp 100
  neighbor a.a.a.a remote-as 200
  neighbor a.a.a.a route-map bgp-neigh-01 out
  neighbor a.a.a.a route-map policy-01 in
  neighbor b.b.b.b remote-as 300
  neighbor b.b.b.b route-map bgp-neigh-02 out
  neighbor b.b.b.b route-map policy-02 in
```

- The route-maps are then applied to the appropriate neighbour
- As long as the customer sets the appropriate communities, the policy will be applied to their prefixes

# Customer BGP configuration

```
router bgp 600
  neighbor c.c.c.c remote-as 100
  neighbor a.a.a.a route-map upstream out
  neighbor a.a.a.a prefix-list default in
!
route-map upstream permit 10
  match ip address prefix-list blockA
  set community 100:5010 100:5023
route-map upstream permit 20
  match ip address aggregate
```

- This will:

- 3x prepend of blockA towards their upstream's 2nd BGP neighbour
  - Not announce blockA towards their upstream's 1st BGP neighbour
  - Let the aggregate through with no specific policy

# Customer Policy Control

- Notice how much flexibility a BGP customer could have with this type of policy implementation
- Advantages:
  - Customer has flexibility
  - ISP Technical Support does not need to be involved
- Disadvantages
  - Customer could upset ISP loadbalancing tuning
- Advice
  - This kind of policy control is very useful, but should only be considered if appropriate for the circumstances



# Conclusion

# Communities

- Communities are fun! 😊
- And they are extremely powerful tools
- Think about community policies, e.g. like the additions described here
- Supporting extensive community usage makes customer configuration easy
- Watch out for routing loops!



# Summary

# Summary

- Multihoming is not hard, really...

**Keep It Simple & Stupid!**

- Full routing table is rarely required

Defaults and careful filtering are just as effective and are not a resource hog

- Splitting your address space into /24s (or /48s for IPv6) will not improve your traffic engineering





# BGP Multihoming Techniques

End of Tutorial 😊