

BGP Multihoming Techniques



Philip Smith

<philip@nsrc.org>

SANOG 27

25th-27th January 2016

Kathmandu

Last updated 9th December 2015

Presentation Slides

- Available on
 - <http://bgp4all.com/ftp/seminars/SANOG27-BGP-Multihoming.pdf>
 - And on the SANOG27 website
- Feel free to ask questions any time

Preliminaries

- Presentation has many configuration examples
 - Uses Cisco IOS CLI
- Aimed at Network Operators
 - Techniques can be used by many enterprises and end-user networks too



BGP Multihoming Techniques

- **Why Multihome?**
- Definition & Options
- How to Multihome
- Principles & Addressing
- Basic Multihoming
- Service Provider Multihoming
- Using Communities

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?

- ❑ Changing upstream provider
- ❑ With one upstream, migration means:
 - Disconnecting existing connection
 - Moving the link to the new upstream
 - Reconnecting the link
 - Reannouncing address space
 - Break in service for end users (hours, days,...?)
- ❑ With two upstreams, migration means:
 - Bring up link with new provider (including BGP and address announcements)
 - Disconnect link with original upstream
 - No break in service for end users

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
- Principles & Addressing
- Basic Multihoming
- Service Provider Multihoming
- Using Communities

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 – RFC6793)
- 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)
 - 65552-4199999999 (public Internet)
 - 4200000000-4294967295 (private use only)
- 32-bit range representation specified in RFC5396
 - Defines “asplain” (traditional format) as standard notation

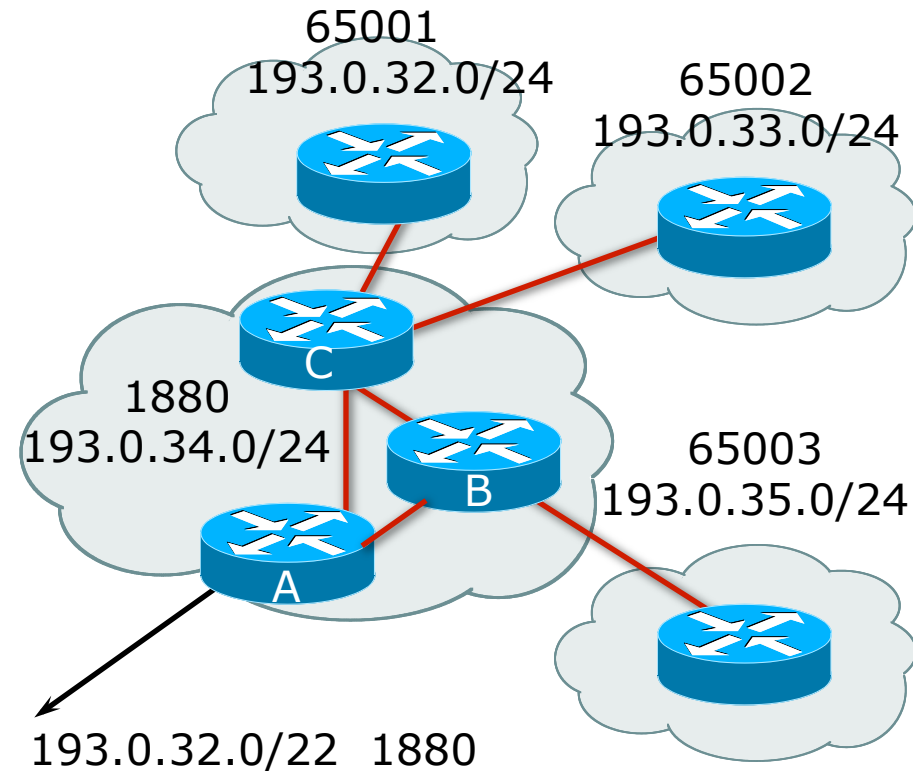
Autonomous System Number (ASN)

- ❑ ASNs are distributed by the Regional Internet Registries
 - They are also available from upstream ISPs who are members of one of the RIRs
- ❑ Current 16-bit ASN assignments up to 64297 have been made to the RIRs
 - Around 43000 16-bit ASNs are visible on the Internet
 - Around 200 left unassigned
- ❑ Each RIR has also received a block of 32-bit ASNs
 - Out of 12000 assignments, around 9200 are visible on the Internet
- ❑ See 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
```

More Definitions

- **Transit**
 - Carrying traffic across a network
 - Usually **for a fee**
- **Peering**
 - Exchanging 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)
- Subdividing Aggregates
 - Inbound traffic flows (local & 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
 - APNIC: www.apnic.net/db/min-alloc.html
 - ARIN: www.arin.net/reference/ip_blocks.html
 - LACNIC: lacnic.net/en/registro/index.html
 - RIPE NCC: 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
- 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:
 - www.team-cymru.org/Services/Bogons/routeserver.html



BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- **How to Multihome**
- Principles & Addressing
- Basic Multihoming
- Service Provider Multihoming
- Using Communities

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

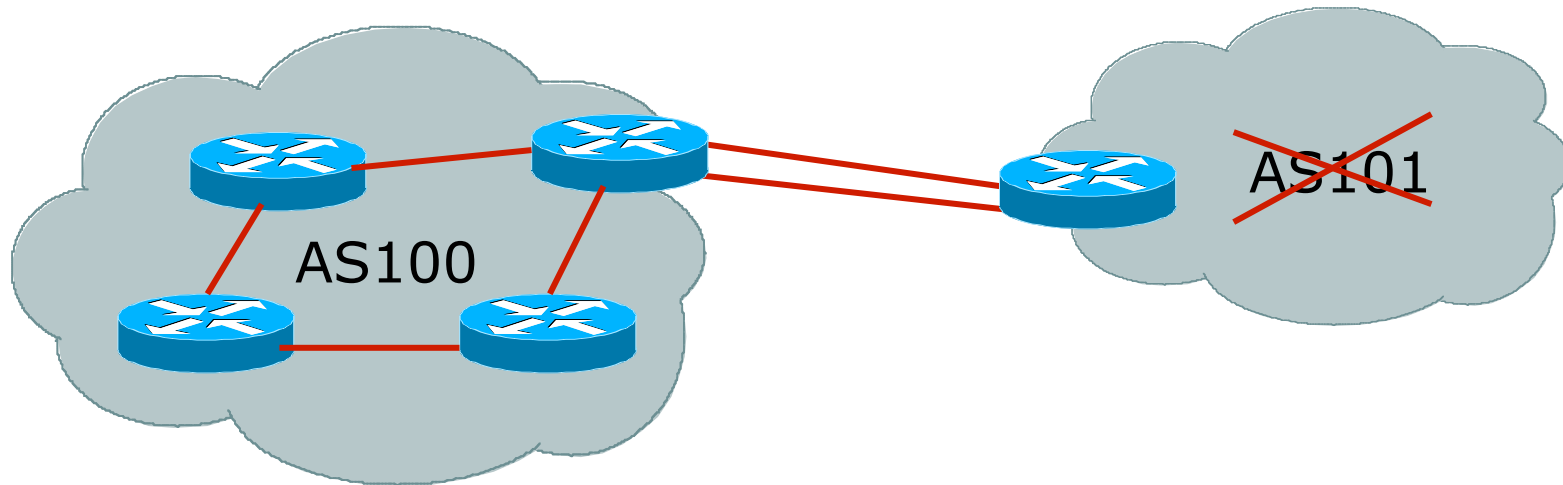
- ❑ Mistaking a transit provider's "Exchange" business for a no-cost public peering point
- ❑ Not working hard to get as much peering as possible
 - Physically near a peering point (IXP) but not present at it
 - (Transit sometimes is cheaper than peering!!)
- ❑ Ignoring/avoiding competitors because they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



Multihoming Scenarios

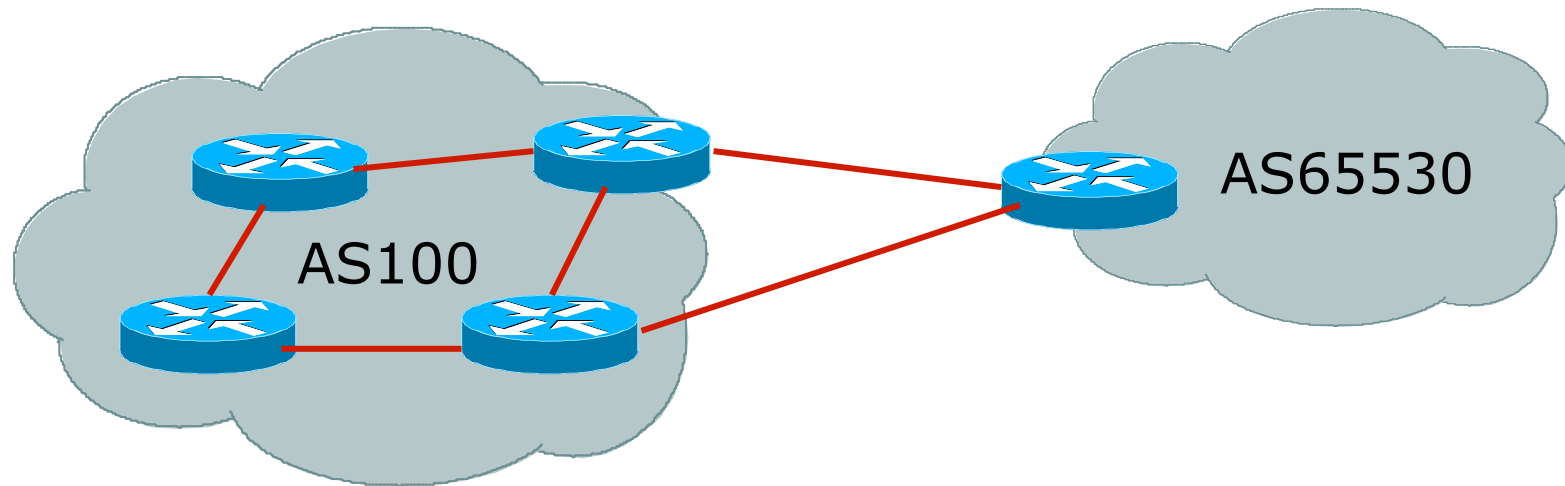
- ❑ Stub network
- ❑ Multi-homed stub network
- ❑ Multi-homed network
- ❑ Multiple sessions to another AS

Stub Network



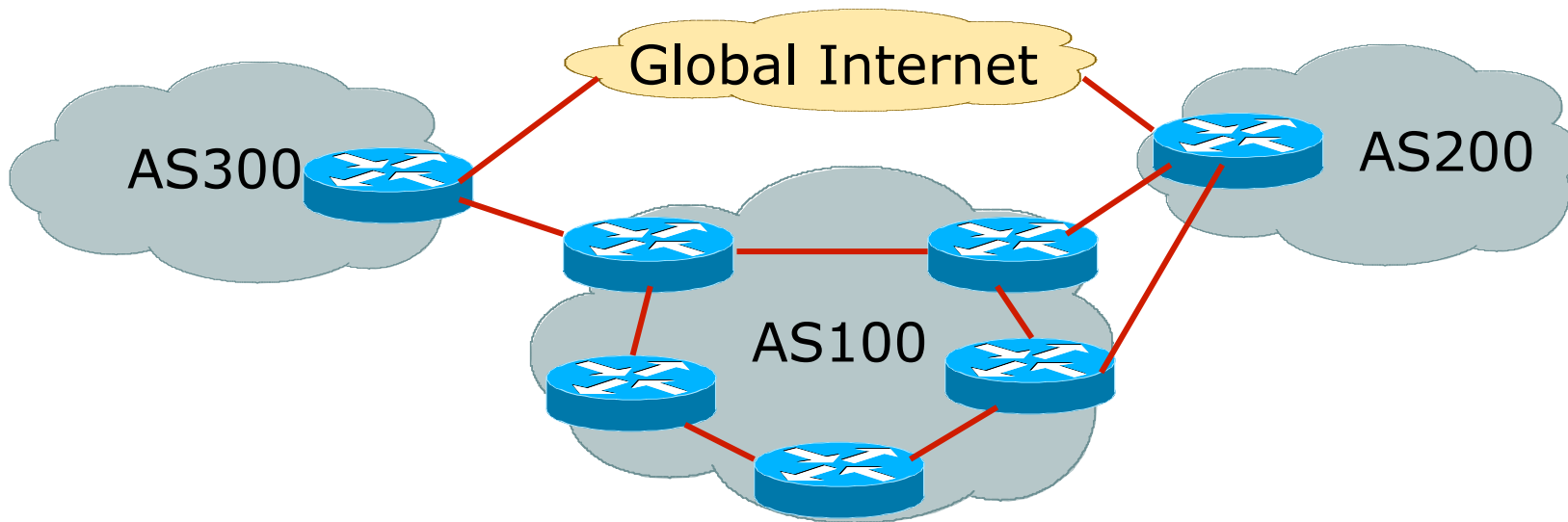
- ❑ No need for BGP
- ❑ Point static default to upstream ISP
- ❑ Upstream ISP advertises stub network
- ❑ Policy confined within upstream ISP's policy

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

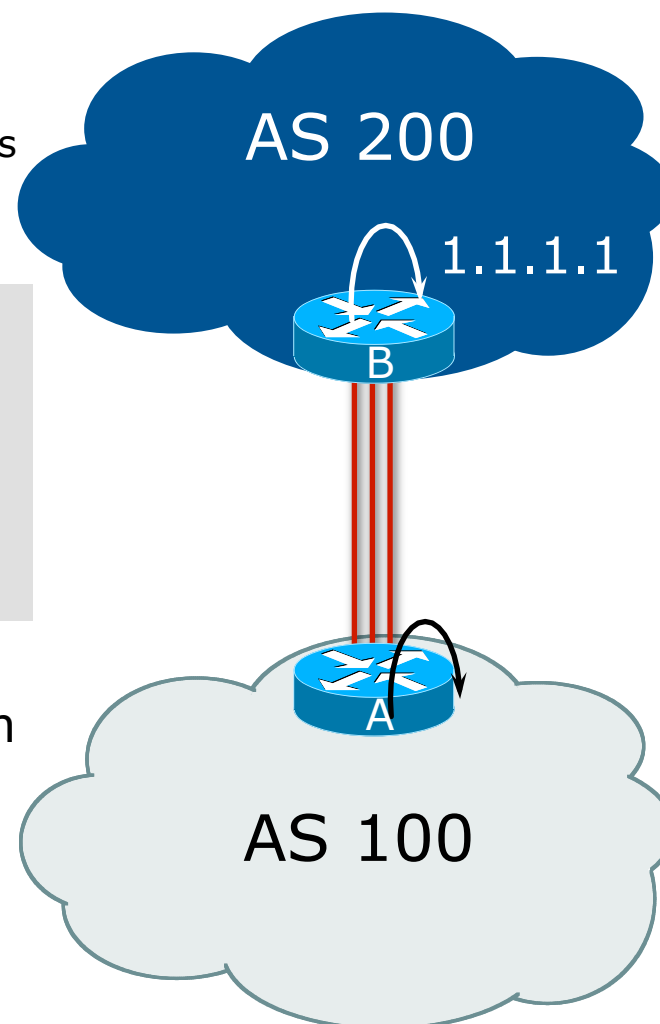
– ebgp multihop

- ❑ Use 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
```

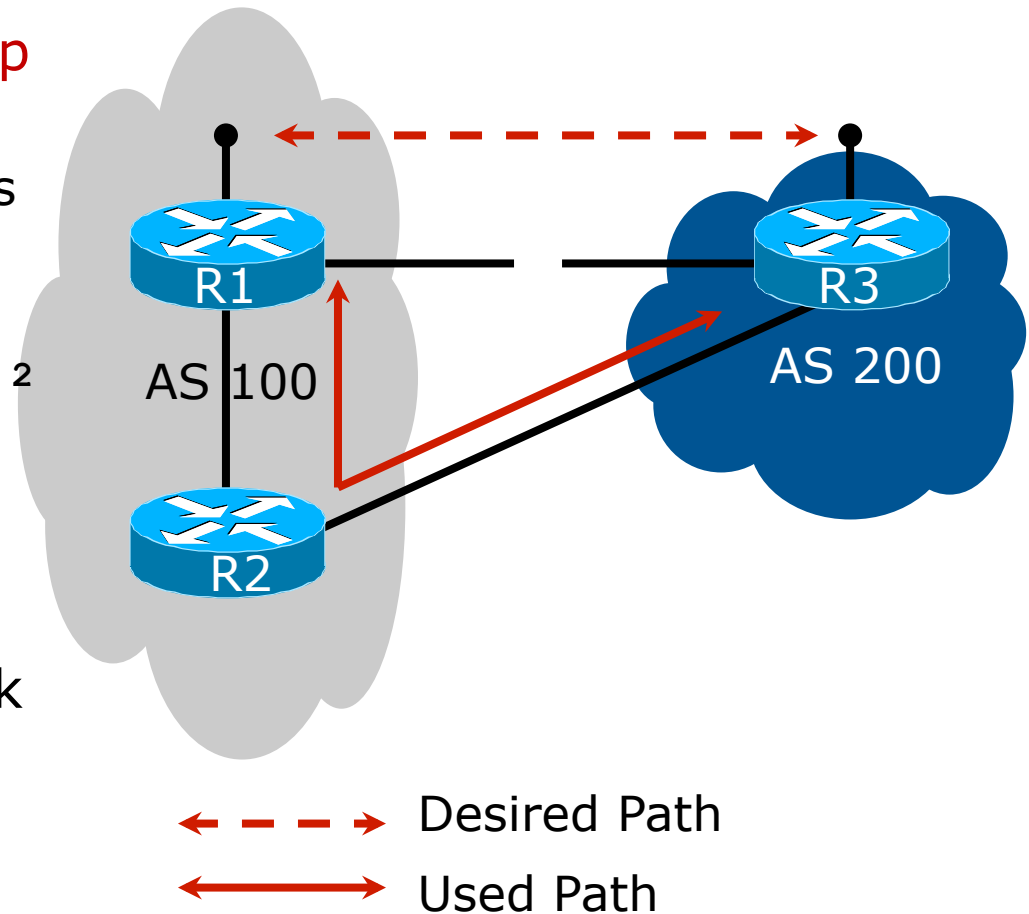
- ❑ Common error made is to point remote loopback route at IP address rather than specific link



Multiple Sessions to an AS

– ebgp multihop

- ❑ One serious eBGP-multihop caveat:
 - 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 ISP

– ebgp multihop

- Try and avoid use of ebgp-multihop unless:
 - It's absolutely necessary –or–
 - Loadsharing across multiple 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:

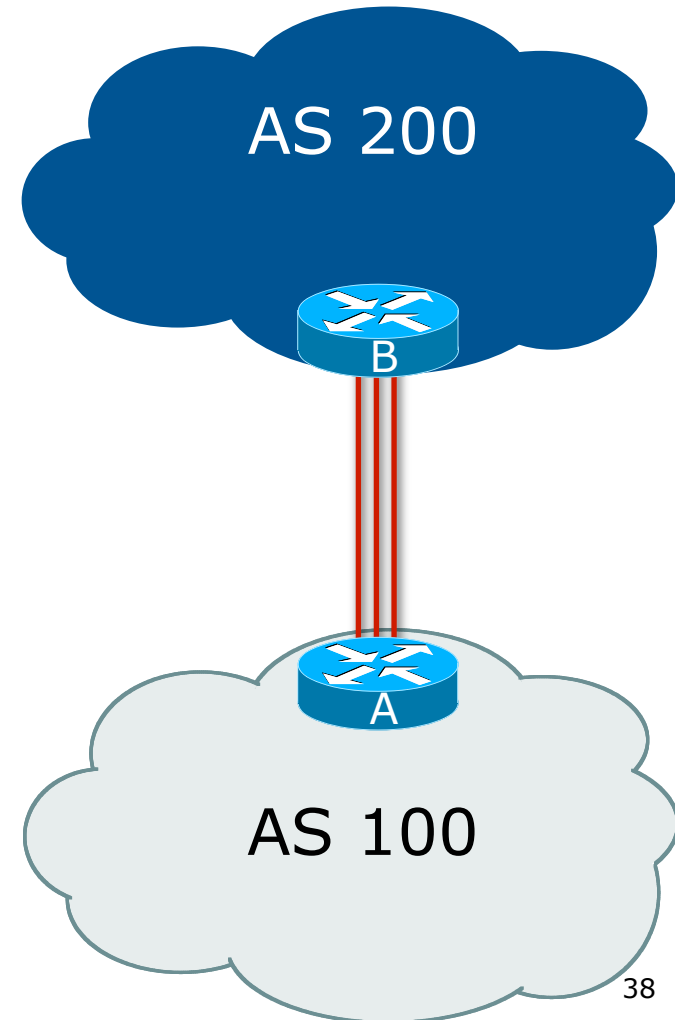
- routing loops
- failure to realise that BGP session stability problems are usually due connectivity problems between their CPE and their BGP speaker

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

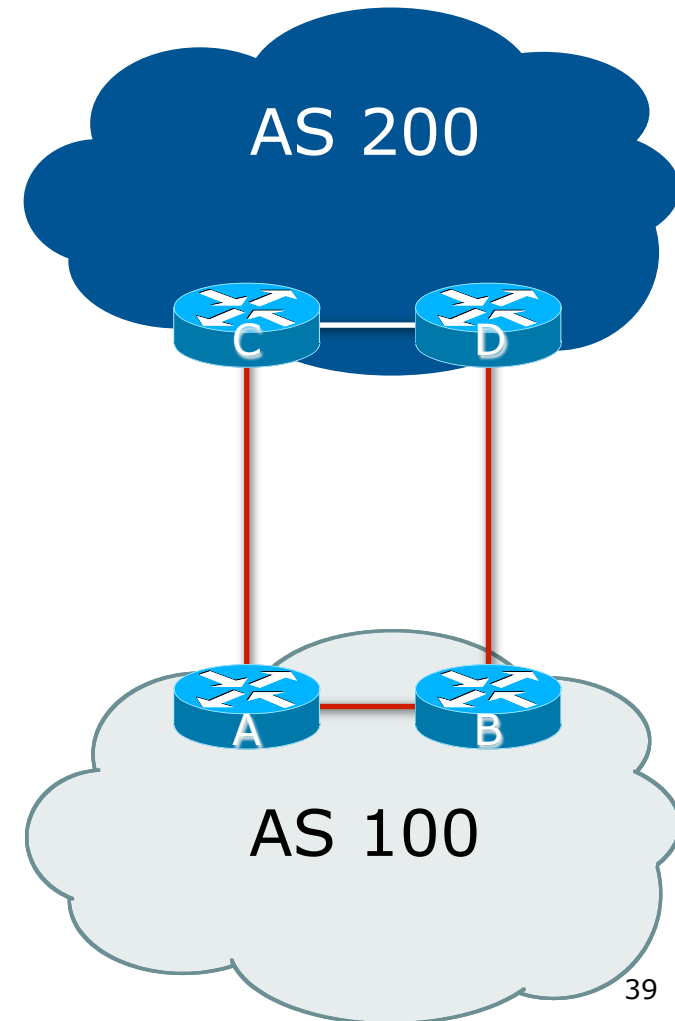
```
router bgp 100
  neighbor 1.1.2.1 remote-as 200
  neighbor 1.1.2.5 remote-as 200
  neighbor 1.1.2.9 remote-as 200
  maximum-paths 3
```



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**





BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- How to Multihome
- Principles & Addressing
- Basic Multihoming
- Service Provider Multihoming
- Using Communities

Basic Principles of Multihoming



Let's learn to walk before we try
running...

The Basic Principles

- ❑ Announcing address space attracts traffic
 - (Unless policy in upstream providers interferes)
- ❑ Announcing the ISP aggregate out a link will result in traffic for that aggregate coming in that link
- ❑ Announcing a subprefix of an aggregate out a link means that all traffic for that subprefix will come in that link, even if the aggregate is announced somewhere else
 - The most specific announcement wins!

The Basic Principles

- To split traffic between two links:
 - Announce the aggregate on both links - ensures redundancy
 - Announce one half of the address space on each link
 - (This is the first step, all things being equal)
- Results in:
 - Traffic for first half of address space comes in first link
 - Traffic for second half of address space comes in second link
 - If either link fails, the fact that the aggregate is announced ensures there is a backup path

The Basic Principles

- The keys to successful multihoming configuration:
 - Keeping traffic engineering prefix announcements independent of customer iBGP
 - Understanding how to announce aggregates
 - Understanding the purpose of announcing subprefixes of aggregates
 - Understanding how to manipulate BGP attributes
 - Too many upstreams/external paths makes multihoming harder (2 or 3 is enough!)

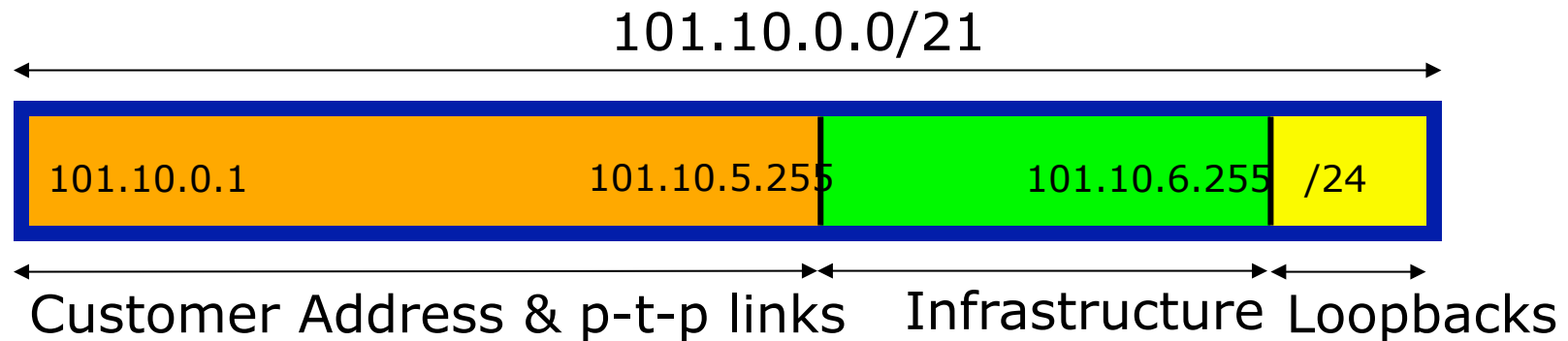
IP Addressing & Multihoming



How Good IP Address Plans
assist with Multihoming

IP Addressing & Multihoming

- ❑ IP Address planning is an important part of Multihoming
- ❑ Previously have discussed separating:
 - Customer address space
 - Customer p-t-p link address space
 - Infrastructure p-t-p link address space
 - Loopback address space

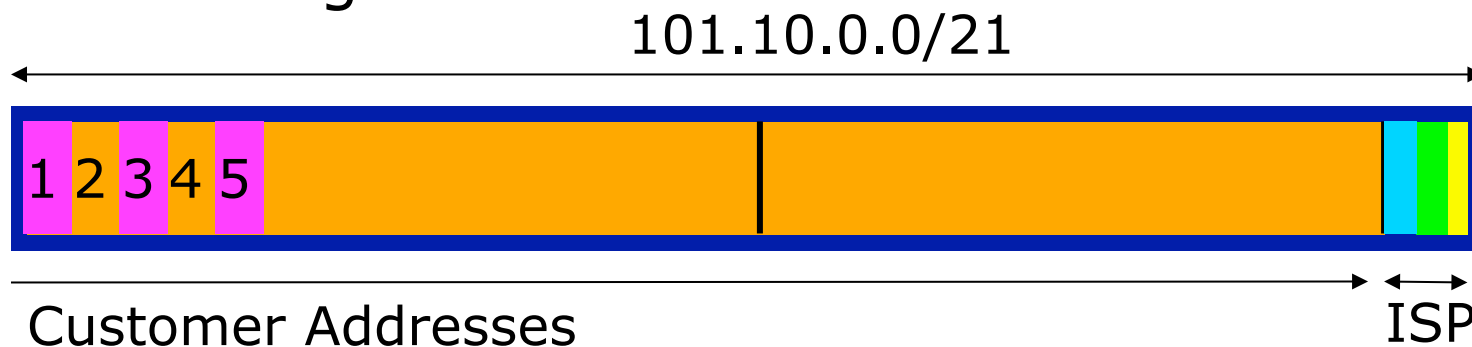


IP Addressing & Multihoming

- ❑ ISP Router loopbacks and backbone point to point links make up a small part of total address space
 - And they don't attract traffic, unlike customer address space
- ❑ Links from ISP Aggregation edge to customer router needs one /30
 - Small requirements compared with total address space
 - Some ISPs use IP unnumbered
- ❑ Planning customer assignments is a very important part of multihoming
 - Traffic engineering involves subdividing aggregate into pieces until load balancing works

Unplanned IP addressing

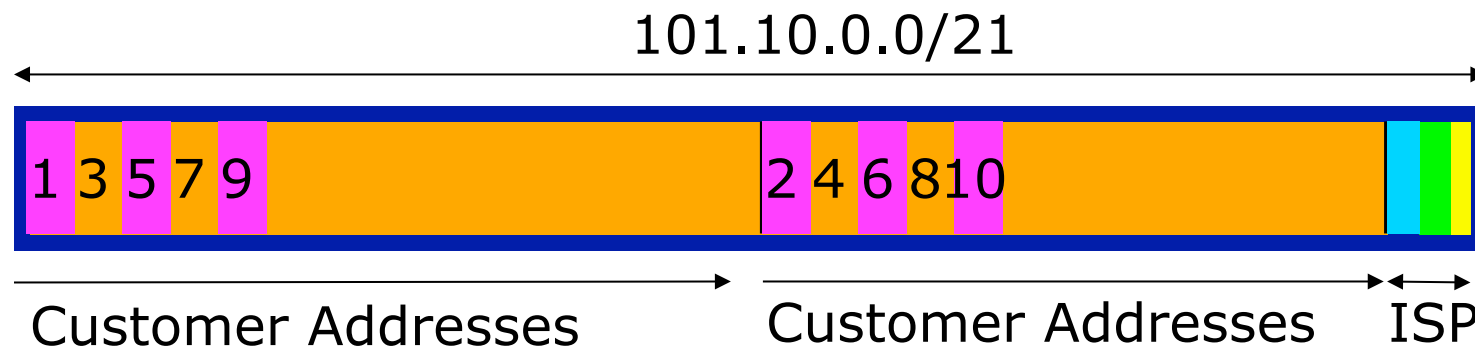
- ISP fills up customer IP addressing from one end of the range:



- Customers generate traffic
 - Dividing the range into two pieces will result in one /22 with all the customers, and one /22 with just the ISP infrastructure the addresses
 - No loadbalancing as all traffic will come in the first /22
 - Means further subdivision of the first /22 = harder work

Planned IP addressing

- If ISP fills up customer addressing from both ends of the range:



- Scheme then is:
 - First customer from first /22, second customer from second /22, third from first /22, etc
- This works also for residential versus commercial customers:
 - Residential from first /22
 - Commercial from second /22

Planned IP Addressing

- ❑ This works fine for multihoming between two upstream links (same or different providers)
- ❑ Can also subdivide address space to suit more than two upstreams
 - Follow a similar scheme for populating each portion of the address space
- ❑ Don't forget to always announce an aggregate out of each link



BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- How to Multihome
- Principles & Addressing
- **Basic Multihoming**
- Service Provider Multihoming
- Using Communities

Basic Multihoming



Let's try some simple worked examples...

Basic Multihoming

- No frills multihoming
- Will look at two cases:
 - Multihoming with the same ISP
 - Multihoming to different ISPs
- Will keep the examples easy
 - Understanding easy concepts will make the more complex scenarios easier to comprehend
 - All assume that the site multihoming has a /19 address block

Basic Multihoming

- This type is most commonplace at the edge of the Internet
 - Networks here are usually concerned with inbound traffic flows
 - Outbound traffic flows being “nearest exit” is usually sufficient
- Can apply to the leaf ISP as well as Enterprise networks

Basic Multihoming



Multihoming to the Same ISP

Basic Multihoming:

Multihoming to the same ISP

- Use BGP for this type of multihoming
 - Use a private AS (ASN > 64511)
 - There is no need or justification for a public ASN
 - Making the nets of the end-site visible gives no useful information to the Internet
- Upstream ISP proxy aggregates
 - In other words, announces only your address block to the Internet from their AS (as would be done if you had one statically routed connection)

Two links to the same ISP

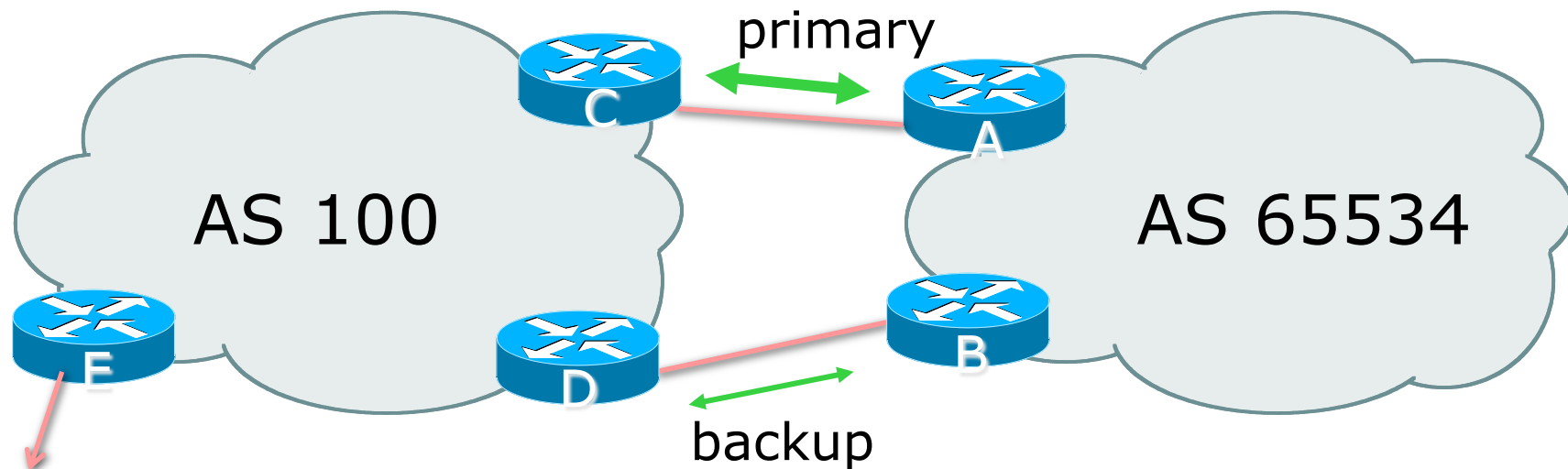


One link primary, the other link
backup only

Two links to the same ISP (one as backup only)

- Applies when end-site has bought a large primary WAN link to their upstream a small secondary WAN link as the backup
 - For example, primary path might be an E1, backup might be 64kbps

Two links to the same ISP (one as backup only)



- AS100 removes private AS and any customer subprefixes from Internet announcement

Two links to the same ISP (one as backup only)

- Announce /19 aggregate on each link
 - primary link:
 - Outbound – announce /19 unaltered
 - Inbound – receive default route
 - backup link:
 - Outbound – announce /19 with increased metric
 - Inbound – received default, and reduce local preference
- When one link fails, the announcement of the /19 aggregate via the other link ensures continued connectivity

Two links to the same ISP (one as backup only)

□ Router A Configuration

```
router bgp 65534
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.2 remote-as 100
  neighbor 122.102.10.2 description RouterC
  neighbor 122.102.10.2 prefix-list aggregate out
  neighbor 122.102.10.2 prefix-list default in
  !
  ip prefix-list aggregate permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

Two links to the same ISP (one as backup only)

□ Router B Configuration

```
router bgp 65534
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.6 remote-as 100
  neighbor 122.102.10.6 description RouterD
  neighbor 122.102.10.6 prefix-list aggregate out
  neighbor 122.102.10.6 route-map med10-out out
  neighbor 122.102.10.6 prefix-list default in
  neighbor 122.102.10.6 route-map lp-low-in in
```

!

..next slide

Two links to the same ISP (one as backup only)

```
ip prefix-list aggregate permit 121.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 121.10.0.0 255.255.224.0 null0
!
route-map med10-out permit 10
  set metric 10
!
route-map lp-low-in permit 10
  set local-preference 90
!
```

Two links to the same ISP (one as backup only)

□ Router C Configuration (main link)

```
router bgp 100
  neighbor 122.102.10.1 remote-as 65534
  neighbor 122.102.10.1 default-originate
  neighbor 122.102.10.1 prefix-list Customer in
  neighbor 122.102.10.1 prefix-list default out
!
ip prefix-list Customer permit 121.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
```


Two links to the same ISP (one as backup only)

□ Router D Configuration (backup link)

```
router bgp 100
  neighbor 122.102.10.5 remote-as 65534
  neighbor 122.102.10.5 default-originate
  neighbor 122.102.10.5 prefix-list Customer in
  neighbor 122.102.10.5 prefix-list default out
!
ip prefix-list Customer permit 121.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
```

Two links to the same ISP (one as backup only)

❑ Router E Configuration

```
router bgp 100
  neighbor 122.102.10.17 remote-as 110
  neighbor 122.102.10.17 remove-private-AS
  neighbor 122.102.10.17 prefix-list Customer out
```

!

```
ip prefix-list Customer permit 121.10.0.0/19
```

- ❑ Router E removes the private AS and customer's subprefixes from external announcements
- ❑ Private AS still visible inside AS100

Two links to the same ISP

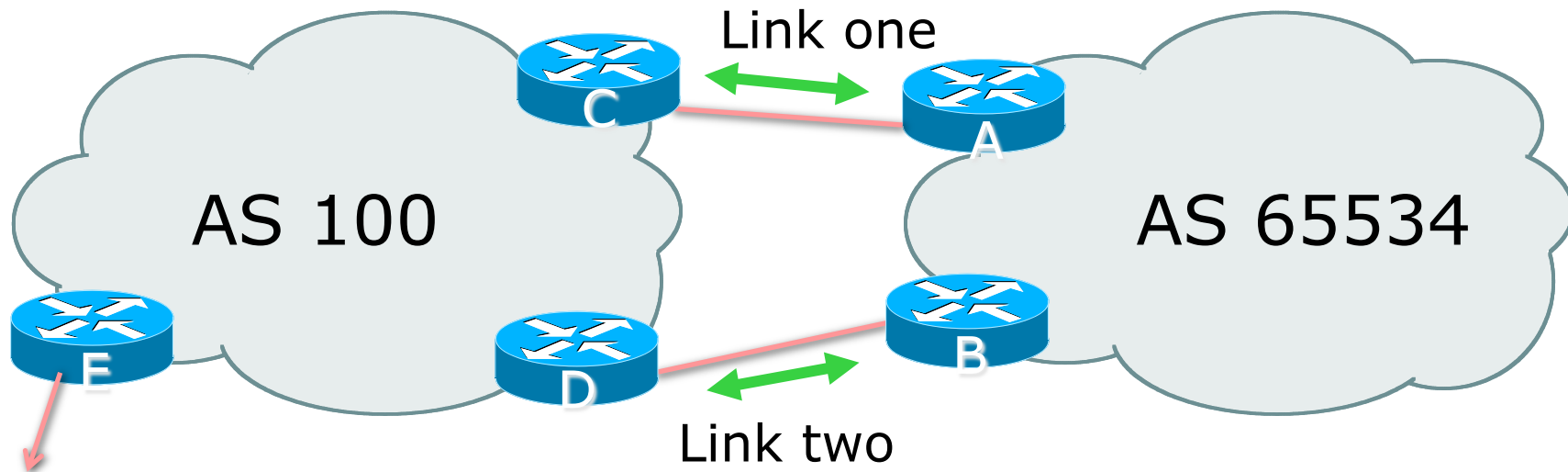


With Loadsharing

Loadsharing to the same ISP

- More common case
- End sites tend not to buy circuits and leave them idle, only used for backup as in previous example
- This example assumes equal capacity circuits
 - Unequal capacity circuits requires more refinement – see later

Loadsharing to the same ISP



- ❑ Border router E in AS100 removes private AS and any customer subprefixes from Internet announcement

Loadsharing to the same ISP

- ❑ Announce /19 aggregate on each link
- ❑ Split /19 and announce as two /20s, one on each link
 - Basic inbound loadsharing
 - Assumes equal circuit capacity and even spread of traffic across address block
- ❑ Vary the split until “perfect” loadsharing achieved
- ❑ Accept the default from upstream
 - Basic outbound loadsharing by nearest exit
 - Okay in first approximation as most ISP and end-site traffic is inbound

Loadsharing to the same ISP (with redundancy)

□ Router A Configuration

```
router bgp 65534
  network 121.10.0.0 mask 255.255.224.0
  network 121.10.0.0 mask 255.255.240.0
  neighbor 122.102.10.2 remote-as 100
  neighbor 122.102.10.2 prefix-list as100-a out
  neighbor 122.102.10.2 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list as100-a permit 121.10.0.0/20
ip prefix-list as100-a permit 121.10.0.0/19
!
ip route 121.10.0.0 255.255.240.0 null0
ip route 121.10.0.0 255.255.224.0 null0
```

Loadsharing to the same ISP (with redundancy)

- ❑ Router C Configuration

```
router bgp 100
  neighbor 122.102.10.1 remote-as 65534
  neighbor 122.102.10.1 default-originate
  neighbor 122.102.10.1 prefix-list Customer in
  neighbor 122.102.10.1 prefix-list default out
  !
  ip prefix-list Customer permit 121.10.0.0/19 le 20
  ip prefix-list default permit 0.0.0.0/0
```

- ❑ Router C only allows in /19 and /20 prefixes from customer block
- ❑ Router D configuration is identical

Loadsharing to the same ISP (with redundancy)

❑ Router E Configuration

```
router bgp 100
  neighbor 122.102.10.17 remote-as 110
  neighbor 122.102.10.17 remove-private-AS
  neighbor 122.102.10.17 prefix-list Customer out
!
ip prefix-list Customer permit 121.10.0.0/19
```

❑ Private AS still visible inside AS100

Loadsharing to the same ISP (with redundancy)

- Default route for outbound traffic?
 - Originate the default route in the IGP on the Border routers
 - Rely on IGP metrics for nearest exit
 - IGP originates default route as long as BGP puts default route in RIB
 - e.g. on router A using OSPF:

```
router ospf 65534
  default-information originate
```

- e.g. on router A using ISIS:

```
router isis as65534
  default-information originate
```

Loadsharing to the same ISP

- Loadsharing configuration is only on customer router
- Upstream ISP has to
 - Remove customer subprefixes from external announcements
 - Remove private AS from external announcements
- Could also use BGP communities

Two links to the same ISP

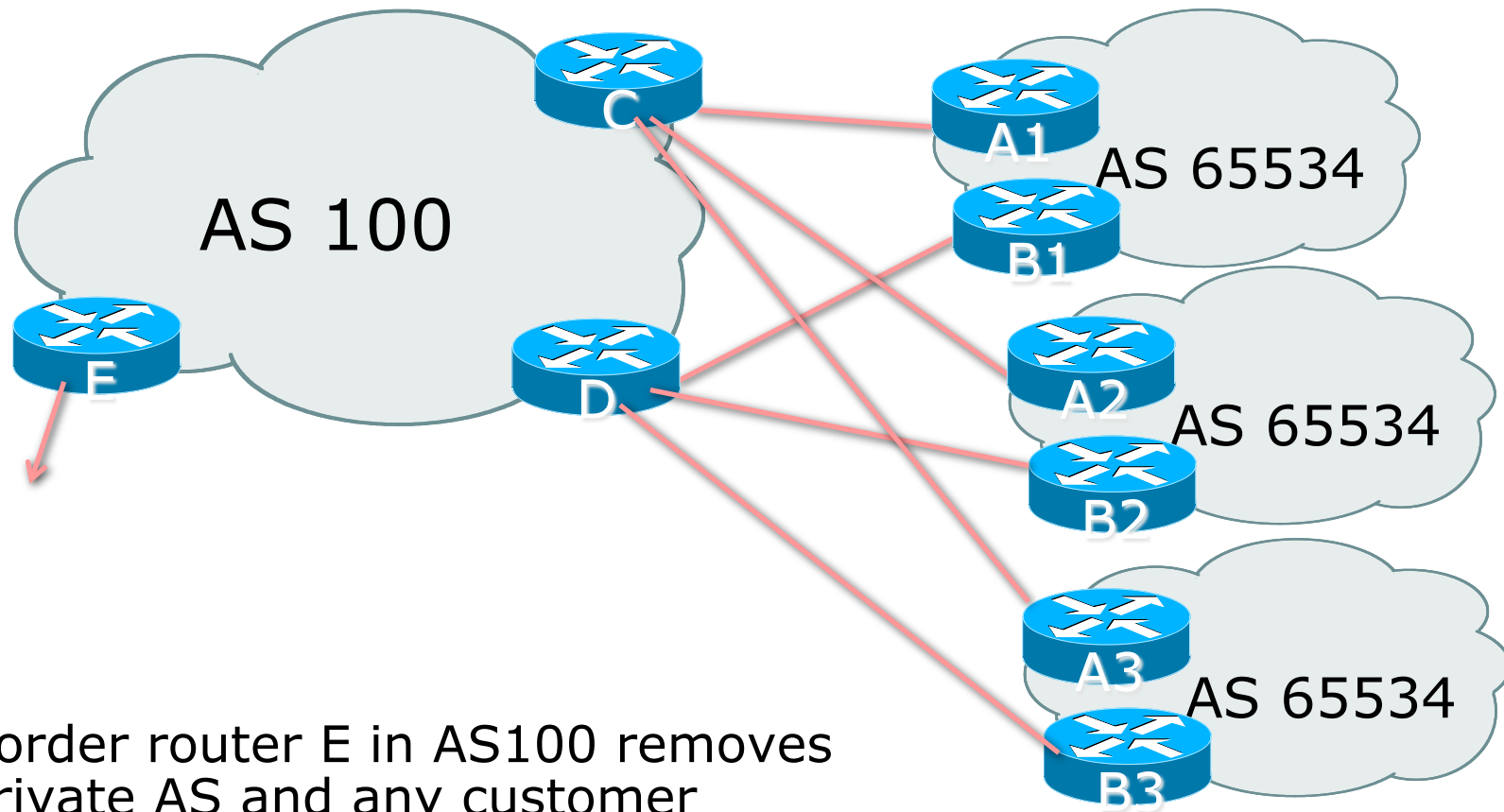


Multiple Dualhomed Customers
(RFC2270)

Multiple Dualhomed Customers (RFC2270)

- Unusual for an ISP just to have one dualhomed customer
 - Valid/valuable service offering for an ISP with multiple PoPs
 - Better for ISP than having customer multihome with another provider!
- Look at scaling the configuration
 - ⇒ Simplifying the configuration
 - Using templates, peer-groups, etc
 - Every customer has the same configuration (basically)

Multiple Dualhomed Customers (RFC2270)



- ❑ Border router E in AS100 removes private AS and any customer subprefixes from Internet announcement

Multiple Dualhomed Customers (RFC2270)

- ❑ Customer announcements as per previous example
- ❑ Use the same private AS for each customer
 - Documented in RFC2270
 - Address space is not overlapping
 - Each customer hears default only
- ❑ Router A_n and B_n configuration same as Router A and B previously

Multiple Dualhomed Customers (RFC2270)

□ Router A1 Configuration

```
router bgp 65534
  network 121.10.0.0 mask 255.255.224.0
  network 121.10.0.0 mask 255.255.240.0
  neighbor 122.102.10.2 remote-as 100
  neighbor 122.102.10.2 prefix-list as100-a out
  neighbor 122.102.10.2 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list as100-a permit 121.10.0.0/20
ip prefix-list as100-a permit 121.10.0.0/19
!
ip route 121.10.0.0 255.255.240.0 null0
ip route 121.10.0.0 255.255.224.0 null0
```


Multiple Dualhomed Customers (RFC2270)

□ Router C Configuration

```
router bgp 100
  neighbor bgp-customers peer-group
  neighbor bgp-customers remote-as 65534
  neighbor bgp-customers default-originate
  neighbor bgp-customers prefix-list default out
neighbor 122.102.10.1 peer-group bgp-customers
neighbor 122.102.10.1 description Customer One
neighbor 122.102.10.1 prefix-list Customer1 in
neighbor 122.102.10.9 peer-group bgp-customers
neighbor 122.102.10.9 description Customer Two
neighbor 122.102.10.9 prefix-list Customer2 in
```

Multiple Dualhomed Customers (RFC2270)

```
neighbor 122.102.10.17 peer-group bgp-customers
neighbor 122.102.10.17 description Customer Three
neighbor 122.102.10.17 prefix-list Customer3 in
!
ip prefix-list Customer1 permit 121.10.0.0/19 le 20
ip prefix-list Customer2 permit 121.16.64.0/19 le 20
ip prefix-list Customer3 permit 121.14.192.0/19 le 20
ip prefix-list default permit 0.0.0.0/0
```

- ❑ Router C only allows in /19 and /20 prefixes from customer block

Multiple Dualhomed Customers (RFC2270)

□ Router E Configuration

- assumes customer address space is not part of upstream's address block

```
router bgp 100
  neighbor 122.102.10.17 remote-as 110
  neighbor 122.102.10.17 remove-private-AS
  neighbor 122.102.10.17 prefix-list Customers out
!
ip prefix-list Customers permit 121.10.0.0/19
ip prefix-list Customers permit 121.16.64.0/19
ip prefix-list Customers permit 121.14.192.0/19
```

□ Private AS still visible inside AS100

Multiple Dualhomed Customers (RFC2270)

- ❑ If customers' prefixes come from ISP's address block
 - Do **NOT** announce them to the Internet
 - Announce ISP aggregate only
- ❑ Router E configuration:

```
router bgp 100
  neighbor 122.102.10.17 remote-as 110
  neighbor 122.102.10.17 prefix-list my-aggregate out
!
ip prefix-list my-aggregate permit 121.8.0.0/13
```

Multihoming Summary

- ❑ Use private AS for multihoming to the same upstream
- ❑ Leak subprefixes to upstream only to aid loadsharing
- ❑ Upstream router E configuration is identical across all situations

Basic Multihoming



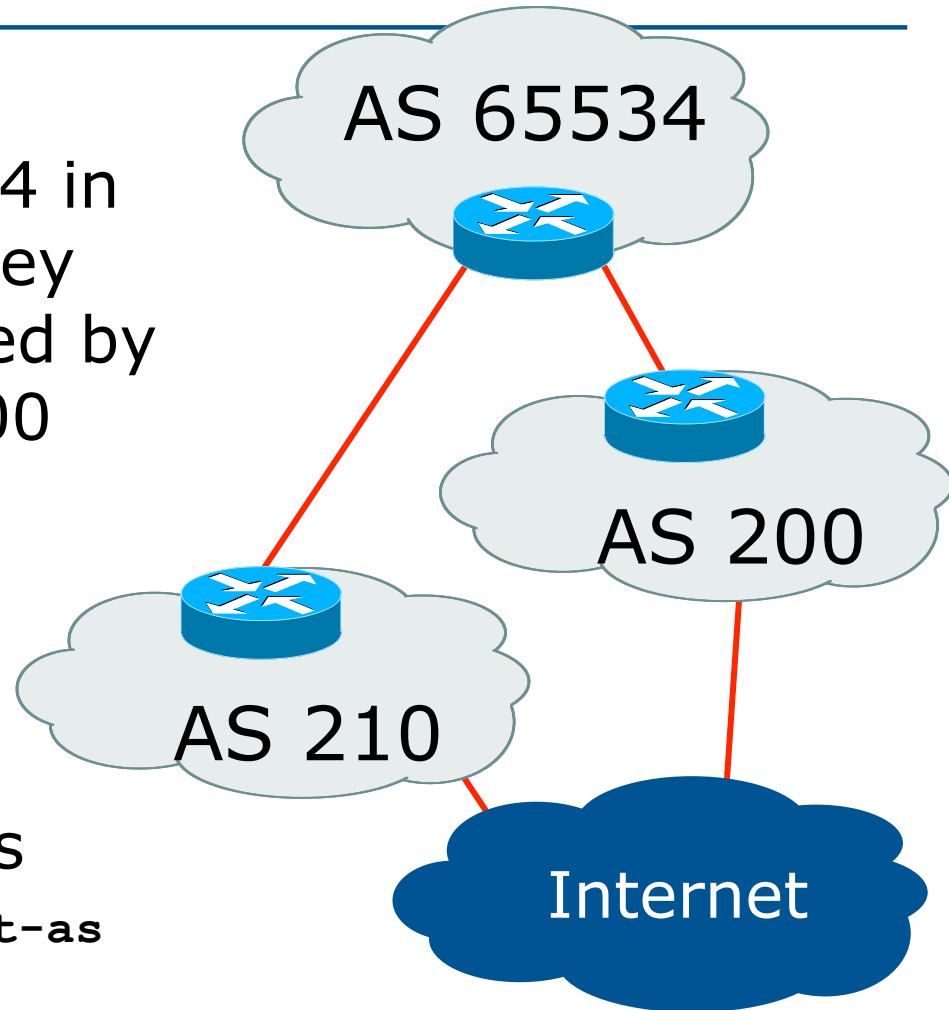
Multihoming to Different ISPs

Two links to different ISPs

- Use a Public AS
 - Or use private AS if agreed with the other ISP
 - But some people don't like the "inconsistent-AS" which results from use of a private-AS
- Address space comes from
 - Both upstreams *or*
 - Regional Internet Registry
 - NB. Very hard to multihome with address space from both upstreams due to typical operational policy in force to day
- Configuration concepts very similar to those used for two links to the same AS

Inconsistent-AS?

- ❑ Viewing the prefixes originated by AS65534 in the Internet shows they appear to be originated by both AS210 and AS200
 - This is NOT bad
 - Nor is it illegal
- ❑ Cisco IOS command is
`show ip bgp inconsistent-as`

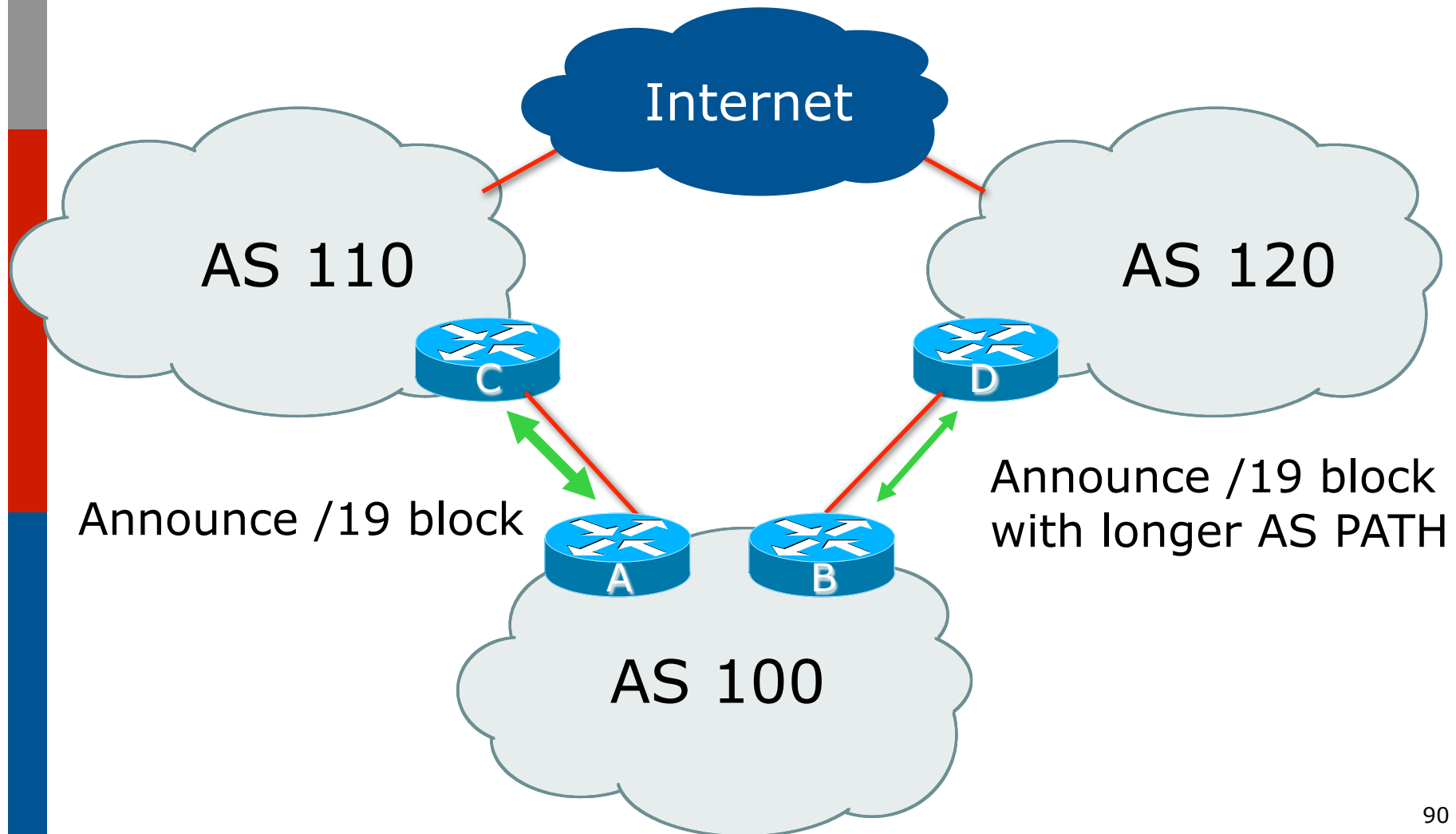


Two links to different ISPs



One link primary, the other link
backup only

Two links to different ISPs (one as backup only)



Two links to different ISPs (one as backup only)

- Announce /19 aggregate on each link
 - primary link makes standard announcement
 - backup link lengthens the AS PATH by using AS PATH prepend
- When one link fails, the announcement of the /19 aggregate via the other link ensures continued connectivity

Two links to different ISPs (one as backup only)

□ Router A Configuration

```
router bgp 130
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote-as 100
  neighbor 122.102.10.1 prefix-list aggregate out
  neighbor 122.102.10.1 prefix-list default in
  !
  ip prefix-list aggregate permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

Two links to different ISPs (one as backup only)

□ Router B Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 120.1.5.1 remote-as 120
  neighbor 120.1.5.1 prefix-list aggregate out
  neighbor 120.1.5.1 route-map as120-prepend out
  neighbor 120.1.5.1 prefix-list default in
  neighbor 120.1.5.1 route-map lp-low in
!
```

...next slide...

Two links to different ISPs (one as backup only)

```
ip route 121.10.0.0 255.255.224.0 null0
!  
ip prefix-list aggregate permit 121.10.0.0/19  
ip prefix-list default permit 0.0.0.0/0  
!  
route-map as120-prepend permit 10  
  set as-path prepend 100 100 100  
!  
route-map lp-low permit 10  
  set local-preference 80  
!
```

Two links to different ISPs (one as backup only)

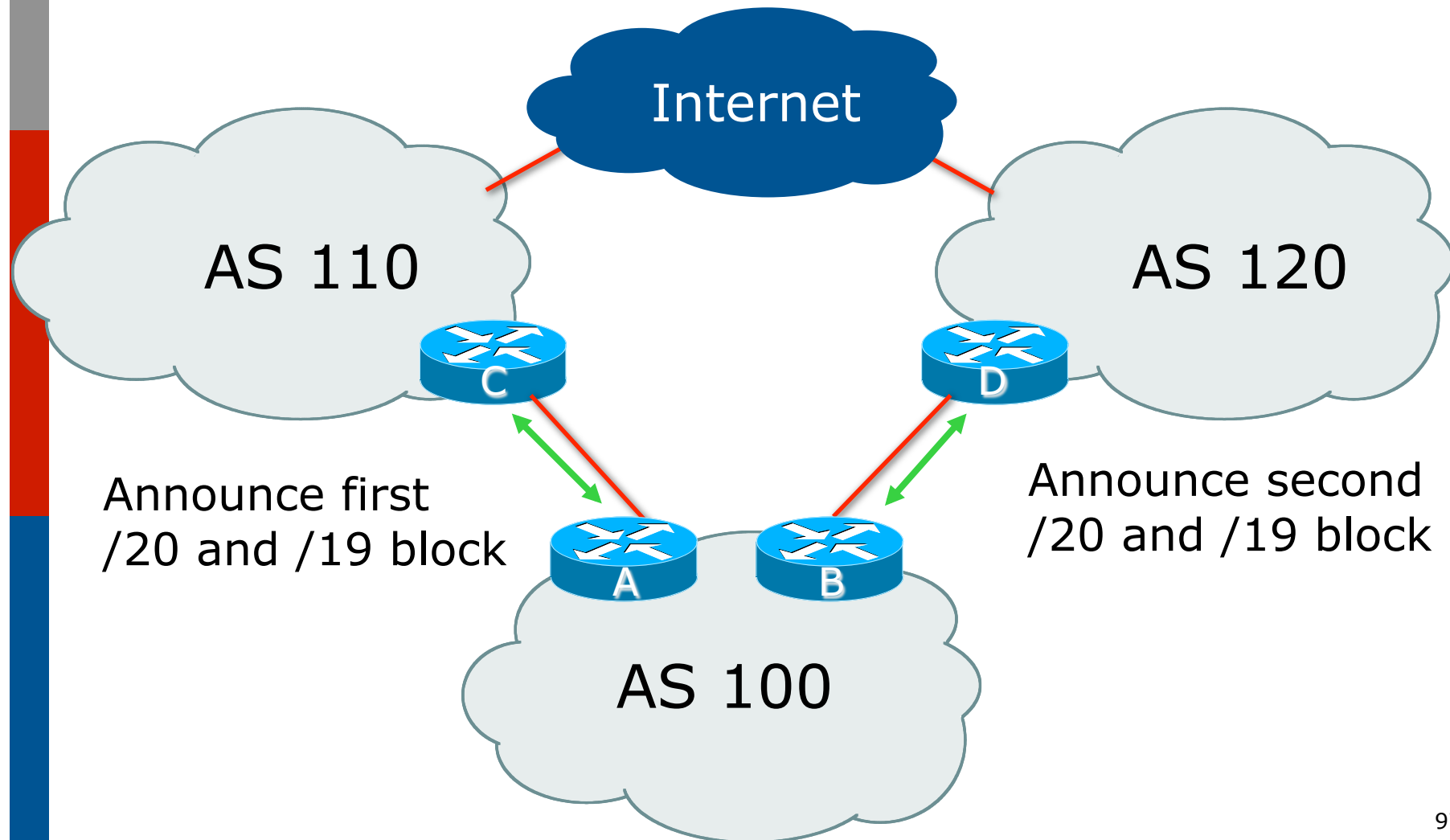
- ❑ Not a common situation as most sites tend to prefer using whatever capacity they have
 - (Useful when two competing ISPs agree to provide mutual backup to each other)
- ❑ But it shows the basic concepts of using local-prefs and AS-path prepends for engineering traffic in the chosen direction

Two links to different ISPs



With Loadsharing

Two links to different ISPs (with loadsharing)



Two links to different ISPs (with loadsharing)

- Announce /19 aggregate on each link
- Split /19 and announce as two /20s, one on each link
 - basic inbound loadsharing
- When one link fails, the announcement of the /19 aggregate via the other ISP ensures continued connectivity

Two links to different ISPs (with loadsharing)

□ Router A Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  network 121.10.0.0 mask 255.255.240.0
  neighbor 122.102.10.1 remote-as 110
  neighbor 122.102.10.1 prefix-list as110-out out
  neighbor 122.102.10.1 prefix-list default in
!
ip route 121.10.0.0 255.255.224.0 null0
ip route 121.10.0.0 255.255.240.0 null0
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list as110-out permit 121.10.0.0/20
ip prefix-list as110-out permit 121.10.0.0/19
```

Two links to different ISPs (with loadsharing)

□ Router B Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  network 121.10.16.0 mask 255.255.240.0
  neighbor 120.1.5.1 remote-as 120
  neighbor 120.1.5.1 prefix-list as120-out out
  neighbor 120.1.5.1 prefix-list default in
!
ip route 121.10.0.0 255.255.224.0 null0
ip route 121.10.16.0 255.255.240.0 null0
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list as120-out permit 121.10.0.0/19
ip prefix-list as120-out permit 121.10.16.0/20
```

Two links to different ISPs (with loadsharing)

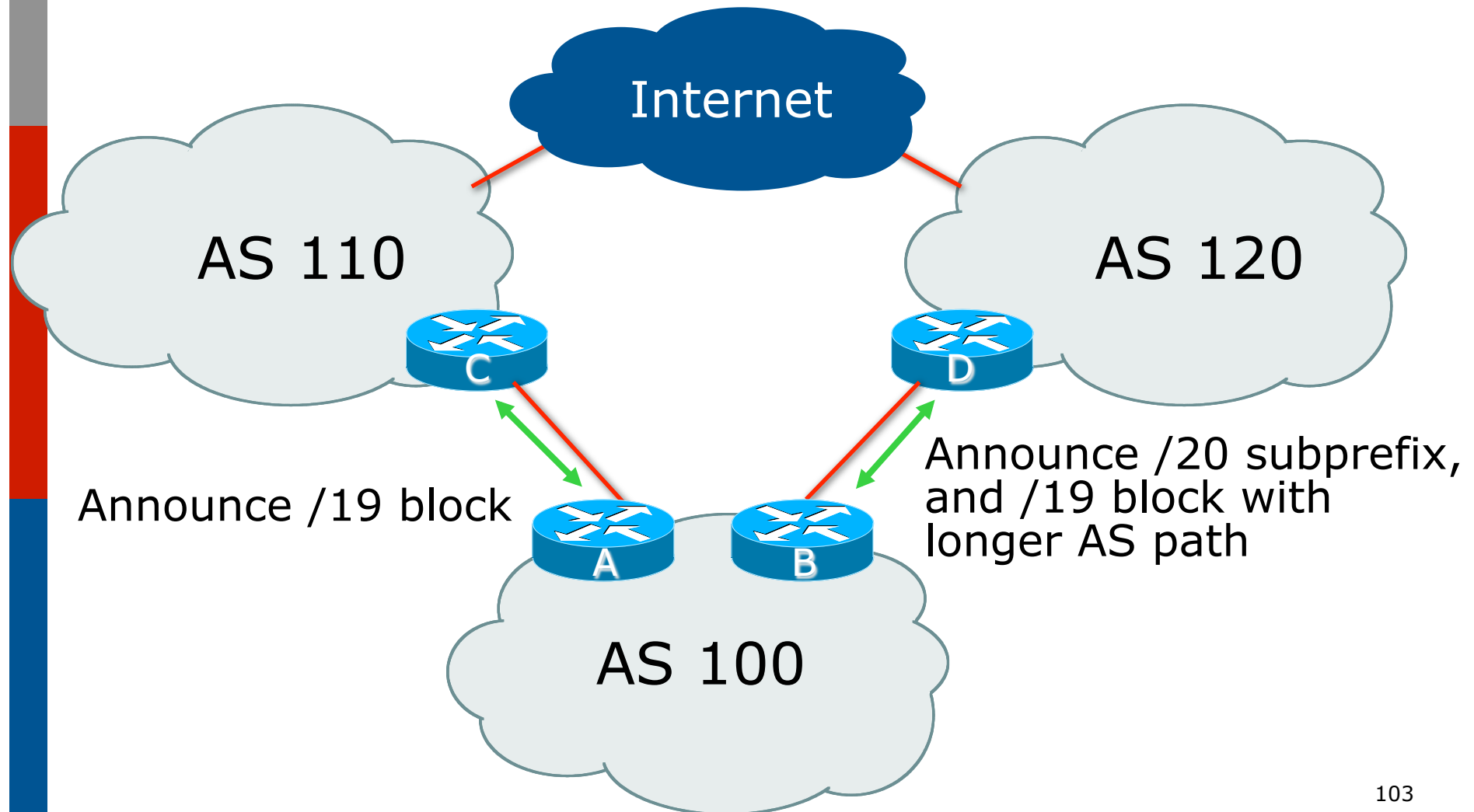
- Loadsharing in this case is very basic
- But shows the first steps in designing a load sharing solution
 - Start with a simple concept
 - And build on it...!

Two links to different ISPs



More Controlled Loadsharing

Loadsharing with different ISPs



Loadsharing with different ISPs

- Announce /19 aggregate on each link
 - On first link, announce /19 as normal
 - On second link, announce /19 with longer AS PATH, and announce one /20 subprefix
 - controls loadsharing between upstreams and the Internet
- Vary the subprefix size and AS PATH length until “perfect” loadsharing achieved
- Still require redundancy!

Loadsharing with different ISPs

□ Router A Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote-as 110
  neighbor 122.102.10.1 prefix-list default in
  neighbor 122.102.10.1 prefix-list as110-out out
!
ip route 121.10.0.0 255.255.224.0 null0
!
ip prefix-list as110-out permit 121.10.0.0/19
!
ip prefix-list default permit 0.0.0.0/0
```

Loadsharing with different ISPs

□ Router B Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  network 121.10.16.0 mask 255.255.240.0
  neighbor 120.1.5.1 remote-as 120
  neighbor 120.1.5.1 prefix-list default in
  neighbor 120.1.5.1 prefix-list as120-out out
  neighbor 120.1.5.1 route-map agg-prepend out
!
ip route 121.10.0.0 255.255.224.0 null0
ip route 121.10.16.0 255.255.240.0 null0
!
...next slide...
```

Loadsharing with different ISPs

```
route-map agg-prepend permit 10
  match ip address prefix-list aggregate
  set as-path prepend 100 100
!
route-map agg-prepend permit 20
!
ip prefix-list default permit 0.0.0.0/0
!
ip prefix-list as120-out permit 121.10.0.0/19
ip prefix-list as120-out permit 121.10.16.0/20
!
ip prefix-list aggregate permit 121.10.0.0/19
!
```

Loadsharing with different ISPs

- ❑ This example is more commonplace
- ❑ Shows how ISPs and end-sites subdivide address space frugally, as well as use the AS-PATH prepend concept to optimise the load sharing between different ISPs
- ❑ Notice that the /19 aggregate block is **ALWAYS** announced



BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- How to Multihome
- Principles & Addressing
- Basic Multihoming
- “BGP Traffic Engineering”
- Using Communities

Service Provider Multihoming



BGP Traffic Engineering

Service Provider Multihoming

- Previous examples dealt with loadsharing inbound traffic
 - Of primary concern at Internet edge
 - What about outbound traffic?
- Transit ISPs strive to balance traffic flows in both directions
 - Balance link utilisation
 - Try and keep most traffic flows symmetric
 - Some edge ISPs try and do this too
- The original “Traffic Engineering”

Service Provider Multihoming

- Balancing outbound traffic requires inbound routing information
 - Common solution is “full routing table”
 - Rarely necessary
 - Why use the “routing mallet” to try solve loadsharing problems?
 - “Keep It Simple” is often easier (and \$\$\$ cheaper) than carrying N-copies of the full routing table

Service Provider Multihoming

MYTHS!!

Common MYTHS

1. **You need the full routing table to multihome**
 - People who sell router memory would like you to believe this
 - Only true if you are a transit provider
 - Full routing table can be a significant hindrance to multihoming
2. **You need a BIG router to multihome**
 - Router size is related to data rates, not running BGP
 - In reality, to multihome, your router needs to:
 - Have two interfaces,
 - Be able to talk BGP to at least two peers,
 - Be able to handle BGP attributes,
 - Handle at least one prefix
3. **BGP is complex**
 - In the wrong hands, yes it can be! Keep it Simple!

Service Provider Multihoming: Some Strategies

- Take the prefixes you need to aid traffic engineering
 - Look at NetFlow data for popular sites
- Prefixes originated by your immediate neighbours and their neighbours will do more to aid load balancing than prefixes from ASNs many hops away
 - Concentrate on local destinations
- Use default routing as much as possible
 - Or use the full routing table with care

Service Provider Multihoming

- Examples
 - One upstream, one local peer
 - One upstream, local exchange point
 - Two upstreams, one local peer
 - Three upstreams, unequal link bandwidths
- Require BGP and a public ASN
- Examples assume that the local network has their own /19 address block

Service Provider Multihoming

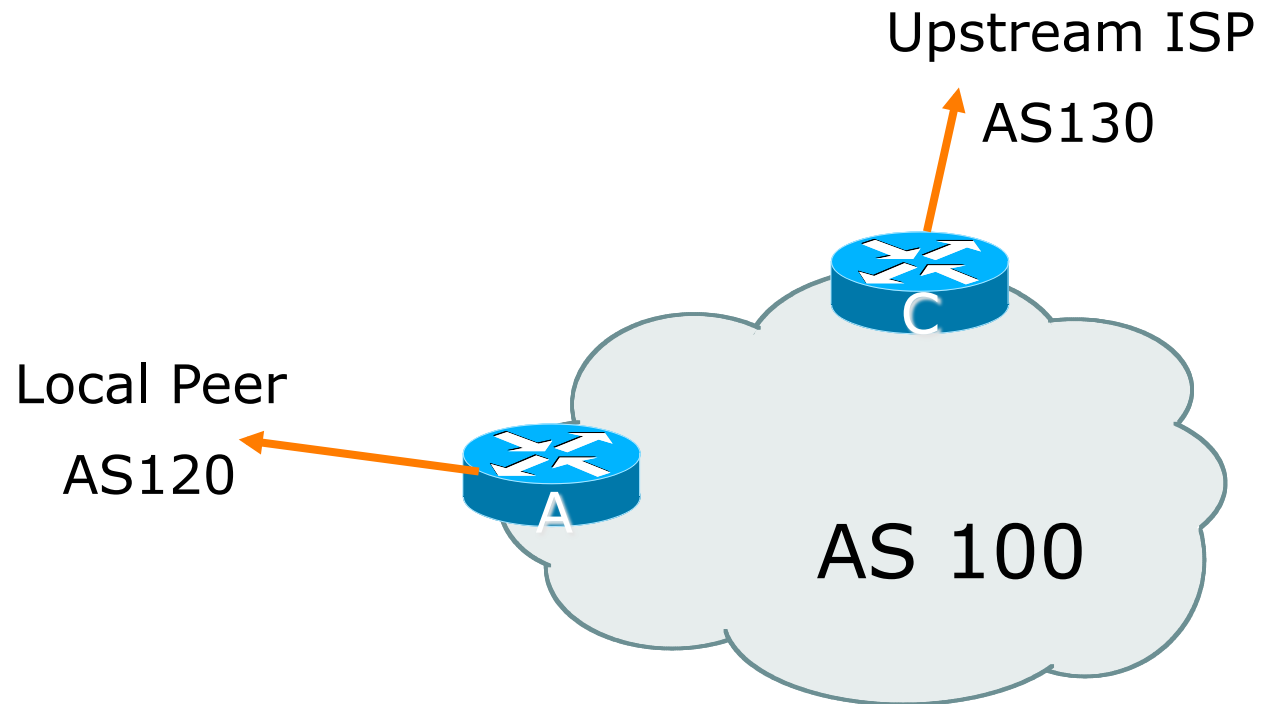


One upstream, one local peer

One Upstream, One Local Peer

- ❑ Very common situation in many regions of the Internet
- ❑ Connect to upstream transit provider to see the “Internet”
- ❑ Connect to the local competition so that local traffic stays local
 - Saves spending valuable \$ on upstream transit costs for local traffic

One Upstream, One Local Peer





One Upstream, One Local Peer

- Announce /19 aggregate on each link
- Accept default route only from upstream
 - Either 0.0.0.0/0 or a network which can be used as default
- Accept all routes from local peer

One Upstream, One Local Peer

□ Router A Configuration

Prefix filters
inbound



```
router bgp 100
 network 121.10.0.0 mask 255.255.224.0
 neighbor 122.102.10.2 remote-as 120
 neighbor 122.102.10.2 prefix-list my-block out
 neighbor 122.102.10.2 prefix-list AS120-peer in
!
ip prefix-list AS120-peer permit 122.5.16.0/19
ip prefix-list AS120-peer permit 121.240.0.0/20
!
ip prefix-list my-block permit 121.10.0.0/19
!
ip route 121.10.0.0 255.255.224.0 null0 250
```


One Upstream, One Local Peer

□ Router A – Alternative Configuration

```
router bgp 100
 network 121.10.0.0 mask 255.255.224.0
 neighbor 122.102.10.2 remote-as 120
 neighbor 122.102.10.2 prefix-list my-block out
 neighbor 122.102.10.2 filter-list 10 in
!
ip as-path access-list 10 permit ^(120_)+$
!
ip prefix-list my-block permit 121.10.0.0/19
!
ip route 121.10.0.0 255.255.224.0 null0
```

AS Path filters –
more “trusting”

One Upstream, One Local Peer

□ Router C Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote-as 130
  neighbor 122.102.10.1 prefix-list default in
  neighbor 122.102.10.1 prefix-list my-block out
  !
  ip prefix-list my-block permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

One Upstream, One Local Peer

- Two configurations possible for Router A
 - Filter-lists assume peer knows what they are doing
 - Prefix-list higher maintenance, but safer
 - Some ISPs use **both**
- Local traffic goes to and from local peer, everything else goes to upstream



Aside:

Configuration Recommendations

□ Private Peers

- The peering ISPs exchange prefixes they originate
 - Sometimes they exchange prefixes from neighbouring ASNs too
- ### □ Be aware that the private peer eBGP router should carry only the prefixes you want the private peer to receive
- Otherwise they could point a default route to you and unintentionally transit your backbone

Service Provider Multihoming

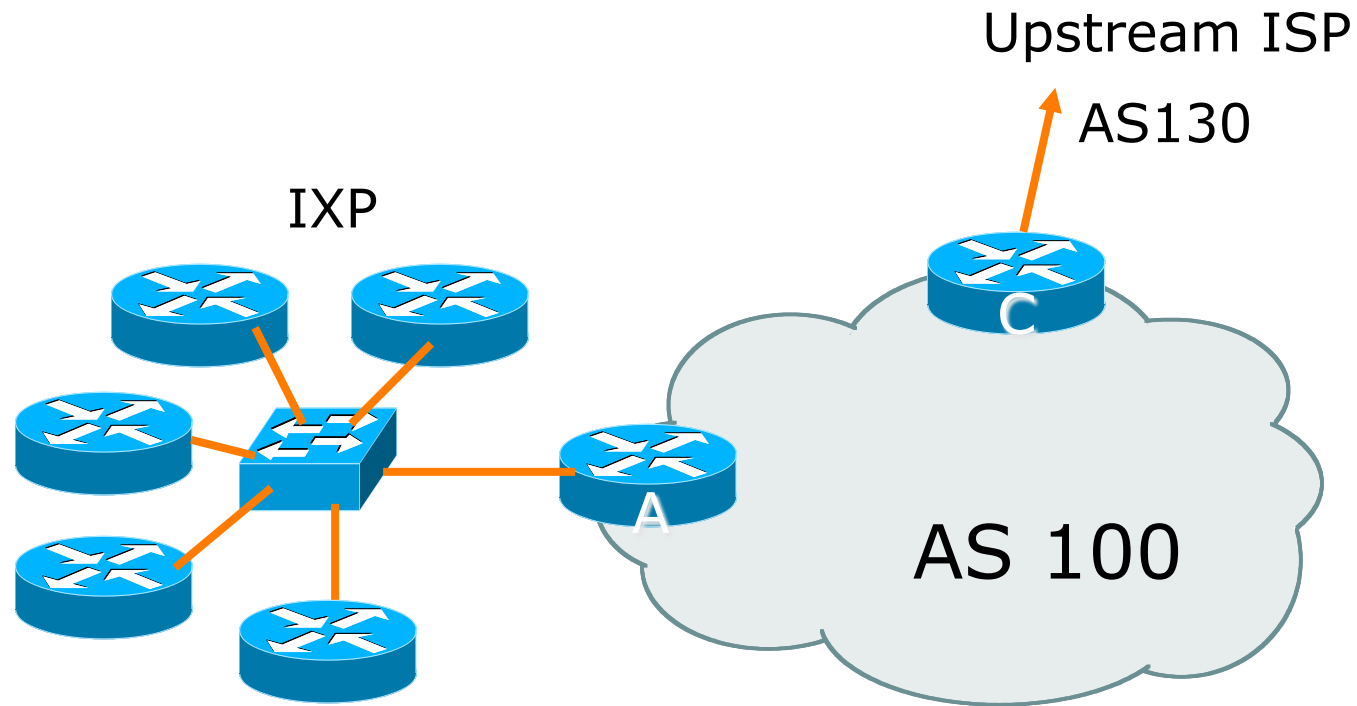


One upstream, Local Exchange
Point

One Upstream, Local Exchange Point

- ❑ Very common situation in many regions of the Internet
- ❑ Connect to upstream transit provider to see the “Internet”
- ❑ Connect to the local Internet Exchange Point so that local traffic stays local
 - Saves spending valuable \$ on upstream transit costs for local traffic

One Upstream, Local Exchange Point



One Upstream, Local Exchange Point

- ❑ Announce /19 aggregate to every neighbouring AS
- ❑ Accept default route only from upstream
 - Either 0.0.0.0/0 or a network which can be used as default
- ❑ Accept all routes originated by IXP peers

One Upstream, Local Exchange Point

□ Router A Configuration

```
interface fastethernet 0/0
  description Exchange Point LAN
  ip address 120.5.10.1 mask 255.255.255.224
!
router bgp 100
  neighbor ixp-peers peer-group
  neighbor ixp-peers prefix-list my-block out
  neighbor ixp-peers remove-private-AS
  neighbor ixp-peers send-community
  neighbor ixp-peers route-map set-local-pref in
!
```

...next slide

One Upstream, Local Exchange Point

```
neighbor 120.5.10.2 remote-as 200
neighbor 120.5.10.2 peer-group ixp-peers
neighbor 120.5.10.2 prefix-list peer200 in
neighbor 120.5.10.3 remote-as 201
neighbor 120.5.10.3 peer-group ixp-peers
neighbor 120.5.10.3 prefix-list peer201 in
neighbor 120.5.10.4 remote-as 202
neighbor 120.5.10.4 peer-group ixp-peers
neighbor 120.5.10.4 prefix-list peer202 in
neighbor 120.5.10.5 remote-as 203
neighbor 120.5.10.5 peer-group ixp-peers
neighbor 120.5.10.5 prefix-list peer203 in
```

...next slide

One Upstream, Local Exchange Point

```
!  
ip prefix-list my-block permit 121.10.0.0/19  
ip prefix-list peer200 permit 122.0.0.0/19  
ip prefix-list peer201 permit 122.30.0.0/19  
ip prefix-list peer202 permit 122.12.0.0/19  
ip prefix-list peer203 permit 122.18.128.0/19  
!  
route-map set-local-pref permit 10  
  set local-preference 150  
!
```

One Upstream, Local Exchange

- ❑ Note that Router A does not generate the aggregate for AS100
 - If Router A becomes disconnected from backbone, then the aggregate is no longer announced to the IX
 - BGP failover works as expected
- ❑ Note the inbound route-map which sets the local preference higher than the default
 - This ensures that BGP Best Path for local traffic will be across the IXP
 - (And allows for future case where operator may need to take BGP routes from their upstream(s))

One Upstream, Local Exchange Point

□ Router C Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote-as 130
  neighbor 122.102.10.1 prefix-list default in
  neighbor 122.102.10.1 prefix-list my-block out
  !
  ip prefix-list my-block permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

One Upstream, Local Exchange Point

- Note Router A configuration
 - Prefix-list higher maintenance, but safer
 - No generation of AS100 aggregate
- IXP traffic goes to and from local IXP, everything else goes to upstream



Aside:

IXP Configuration Recommendations

- IXP peers
 - The peering ISPs at the IXP exchange prefixes they originate
 - Sometimes they exchange prefixes from neighbouring ASNs too
- Be aware that the IXP border router should carry only the prefixes you want the IXP peers to receive and the destinations you want them to be able to reach
 - Otherwise they could point a default route to you and unintentionally transit your backbone
- If IXP router is at IX, and distant from your backbone
 - Don't originate your address block at your IXP router

Service Provider Multihoming



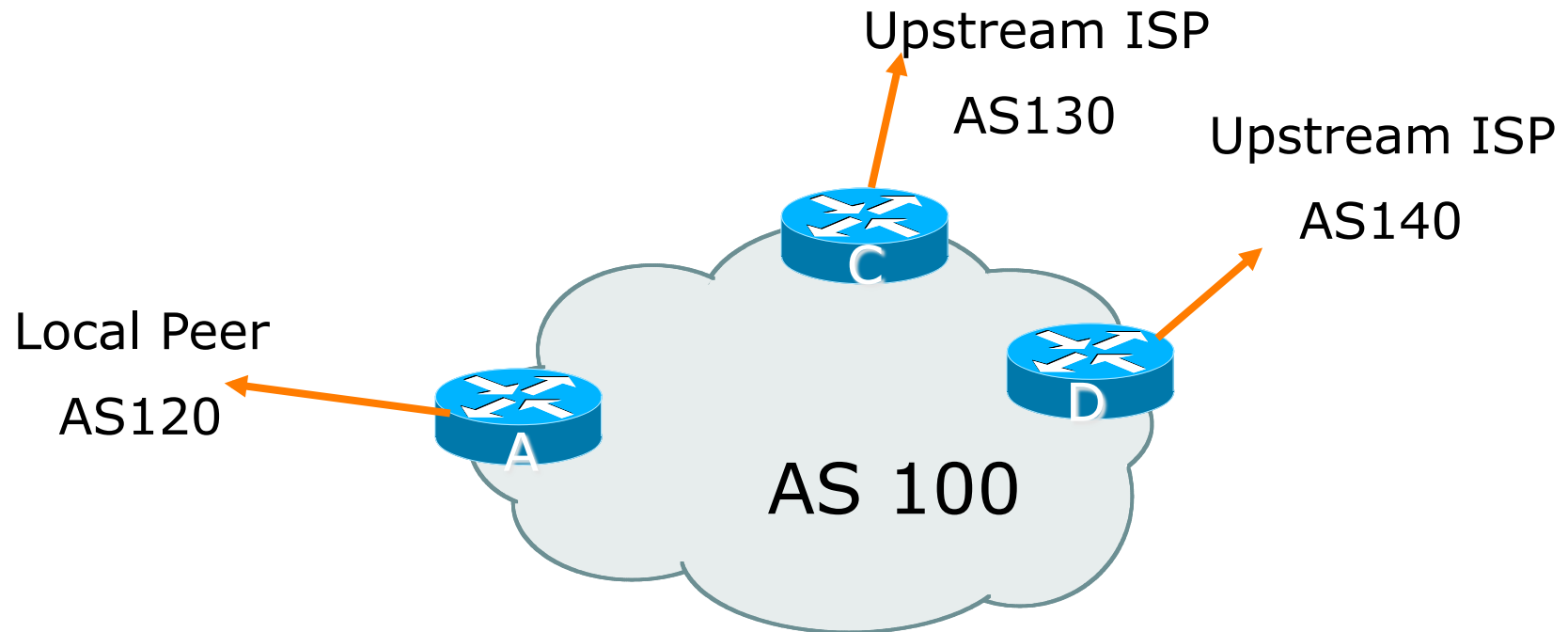
Two Upstreams, One local peer



Two Upstreams, One Local Peer

- Connect to both upstream transit providers to see the “Internet”
 - Provides external redundancy and diversity – the reason to multihome
- Connect to the local peer so that local traffic stays local
 - Saves spending valuable \$ on upstream transit costs for local traffic

Two Upstreams, One Local Peer



Two Upstreams, One Local Peer

- Announce /19 aggregate on each link
- Accept default route only from upstreams
 - Either 0.0.0.0/0 or a network which can be used as default
- Accept all routes from local peer
- Note separation of Router C and D
 - Single edge router means no redundancy
- Router A
 - Same routing configuration as in example with one upstream and one local peer

Two Upstreams, One Local Peer

□ Router C Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote-as 130
  neighbor 122.102.10.1 prefix-list default in
  neighbor 122.102.10.1 prefix-list my-block out
  !
  ip prefix-list my-block permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

Two Upstreams, One Local Peer

□ Router D Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.5 remote-as 140
  neighbor 122.102.10.5 prefix-list default in
  neighbor 122.102.10.5 prefix-list my-block out
  !
  ip prefix-list my-block permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

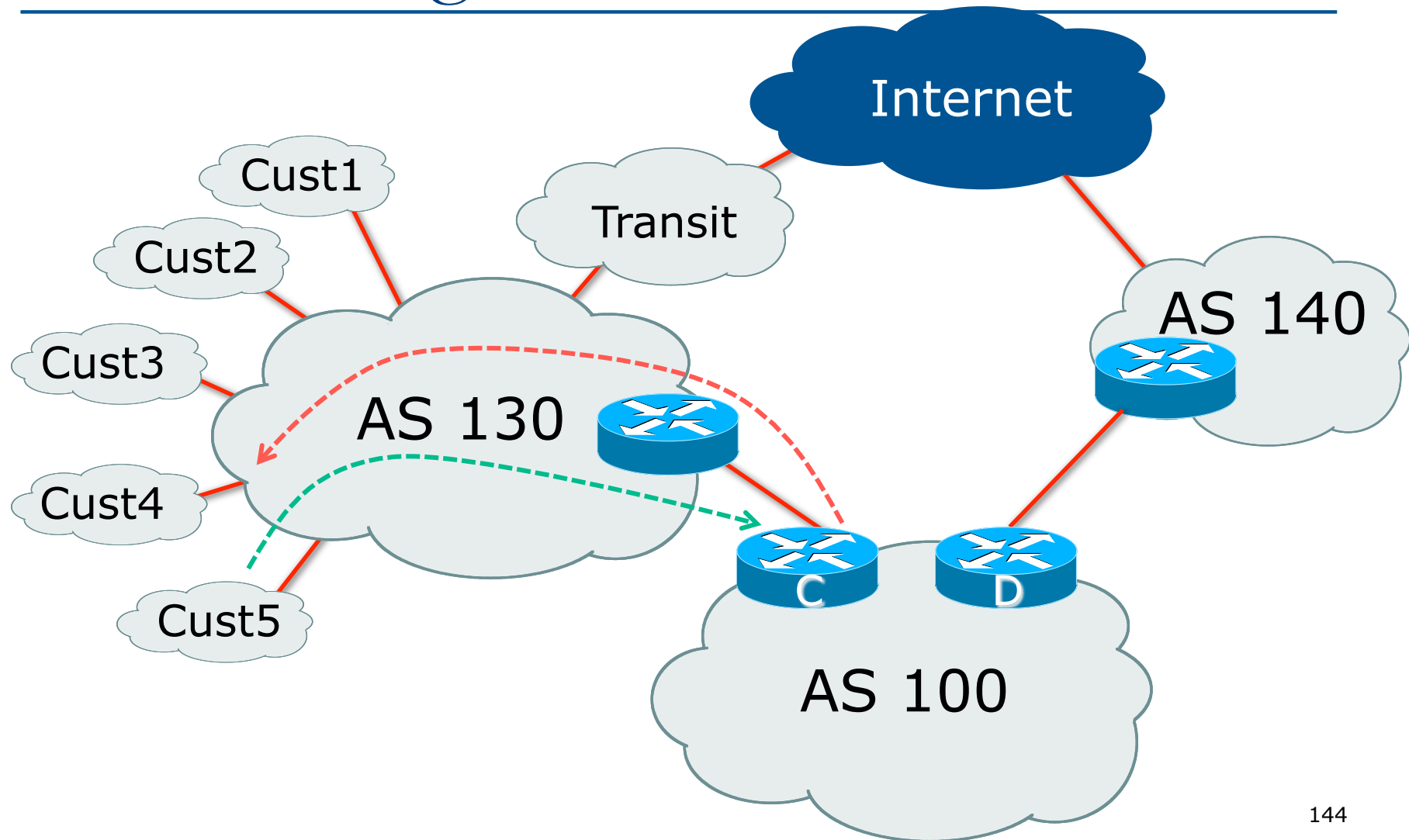
Two Upstreams, One Local Peer

- ❑ This is the simple configuration for Router C and D
- ❑ Traffic out to the two upstreams will take nearest exit
 - Inexpensive routers required
 - This is not useful in practice especially for international links
 - Loadsharing needs to be better

Two Upstreams, One Local Peer

- Better configuration options:
 - Accept full routing from both upstreams
 - **Expensive & unnecessary!**
 - Accept default from one upstream and some routes from the other upstream
 - **The way to go!**

Loadsharing with different ISPs

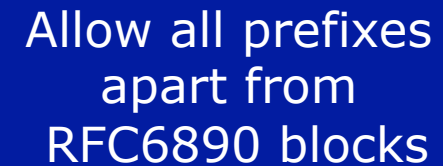


Two Upstreams, One Local Peer

Full Routes

□ Router C Configuration

Allow all prefixes
apart from
RFC6890 blocks



```
router bgp 100
 network 121.10.0.0 mask 255.255.224.0
 neighbor 122.102.10.1 remote-as 130
 neighbor 122.102.10.1 prefix-list rfc6890-deny in
 neighbor 122.102.10.1 prefix-list my-block out
 neighbor 122.102.10.1 route-map AS130-loadshare in
!
ip prefix-list my-block permit 121.10.0.0/19
!
! See http://tools.ietf.org/html/rfc6890
...next slide
```

Two Upstreams, One Local Peer

Full Routes


```
ip route 121.10.0.0 255.255.224.0 null0
!  
ip as-path access-list 10 permit ^(130_)+$  
ip as-path access-list 10 permit ^(130_)+_[0-9]+$  
!  
route-map AS130-loadshare permit 10  
  match ip as-path 10  
  set local-preference 120  
!  
route-map AS130-loadshare permit 20  
  set local-preference 80  
!
```

Two Upstreams, One Local Peer

Full Routes

□ Router D Configuration

Allow all prefixes
apart from
RFC6890 blocks



```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.5 remote-as 140
  neighbor 122.102.10.5 prefix-list rfc6890-deny in
  neighbor 122.102.10.5 prefix-list my-block out
!
ip prefix-list my-block permit 121.10.0.0/19
!
! See http://tools.ietf.org/html/rfc6890
```

Two Upstreams, One Local Peer

Full Routes

- Router C configuration:
 - Accept full routes from AS130
 - Tag prefixes originated by AS130 and AS130's neighbouring ASes with local preference 120
 - Traffic to those ASes will go over AS130 link
 - Remaining prefixes tagged with local preference of 80
 - Traffic to other all other ASes will go over the link to AS140
- Router D configuration same as Router C without the route-map

Two Upstreams, One Local Peer

Full Routes

- Full routes from upstreams
 - Summary of routes received:

ASN	Full Routes		Partial Routes
AS140	570000	@ lp 100	
AS130	30000	@ lp 120	
	540000	@ lp 80	
Total	1140000		

Two Upstreams, One Local Peer

Full Routes

- Full routes from upstreams
 - Expensive – needs lots of memory and CPU
 - Need to play preference games
 - Previous example is only an example – real life will need improved fine-tuning!
 - Previous example doesn't consider inbound traffic – see earlier in presentation for examples

Two Upstreams, One Local Peer

Partial Routes: Strategy

- Ask one upstream for a default route
 - Easy to originate default towards a BGP neighbour
- Ask other upstream for a full routing table
 - Then filter this routing table based on neighbouring ASN
 - E.g. want traffic to their neighbours to go over the link to that ASN
 - Most of what upstream sends is thrown away
 - Easier than asking the upstream to set up custom BGP filters for you

Two Upstreams, One Local Peer

Partial Routes

Router C Configuration

```
router bgp 100
 network 121.10.0.0 mask 255.255.224.0
 neighbor 122.102.10.1 remote-as 130
 neighbor 122.102.10.1 prefix-list rfc6890-deny in
 neighbor 122.102.10.1 prefix-list my-block out
 neighbor 122.102.10.1 filter-list 10 in
 neighbor 122.102.10.1 route-map tag-default-low in
!
```

Allow all prefixes
apart from
RFC6890 blocks

AS filter list filters
prefixes based on
origin ASN

Two Upstreams, One Local Peer

Partial Routes

```
ip prefix-list my-block permit 121.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 121.10.0.0 255.255.224.0 null0
!
ip as-path access-list 10 permit ^(130_)+$
ip as-path access-list 10 permit ^(130_)+_[0-9]+$
!
route-map tag-default-low permit 10
  match ip address prefix-list default
  set local-preference 80
!
route-map tag-default-low permit 20
!
```

Two Upstreams, One Local Peer

Partial Routes

□ Router D Configuration

```
router bgp 100
  network 121.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.5 remote-as 140
  neighbor 122.102.10.5 prefix-list default in
  neighbor 122.102.10.5 prefix-list my-block out
  !
  ip prefix-list my-block permit 121.10.0.0/19
  ip prefix-list default permit 0.0.0.0/0
  !
  ip route 121.10.0.0 255.255.224.0 null0
```

Two Upstreams, One Local Peer

Partial Routes

- Router C configuration:
 - Accept full routes from AS130
 - (or get them to send less)
 - Filter ASNs so only AS130 and AS130's neighbouring ASes are accepted
 - Allow default, and set it to local preference 80
 - Traffic to those ASes will go over AS130 link
 - Traffic to other all other ASes will go over the link to AS140
 - If AS140 link fails, backup via AS130 – and vice-versa

Two Upstreams, One Local Peer

Partial Routes

- Partial routes from upstreams
 - Summary of routes received:

ASN	Full Routes		Partial Routes	
AS140	570000	@ lp 100	1	@lp 100
AS130	30000	@ lp 120	30000	@lp 100
	540000	@ lp 80	1	@lp 80
Total	1140000		30002	

Distributing Default route with IGP

❑ Router C IGP Configuration

```
router ospf 100
default-information originate metric 30
!
```

❑ Router D IGP Configuration

```
router ospf 100
default-information originate metric 10
!
```

- ❑ Primary path is via Router D, with backup via Router C
 - Preferred over carrying default route in iBGP

Two Upstreams, One Local Peer

Partial Routes

- Partial routes from upstreams
 - Not expensive – only carry the routes necessary for loadsharing
 - Need to filter on AS paths
 - Previous example is only an example – real life will need improved fine-tuning!
 - Previous example doesn't consider inbound traffic – see earlier in presentation for examples

Aside:

Configuration Recommendation

- When distributing internal default by iBGP or OSPF/ISIS
 - Make sure that routers connecting to private peers or to IXPs do **NOT** carry the default route
 - Otherwise they could point a default route to you and unintentionally transit your backbone
 - Simple fix for Private Peer/IXP routers:

```
ip route 0.0.0.0 0.0.0.0 null0
```

Service Provider Multihoming



Three upstreams, unequal
bandwidths

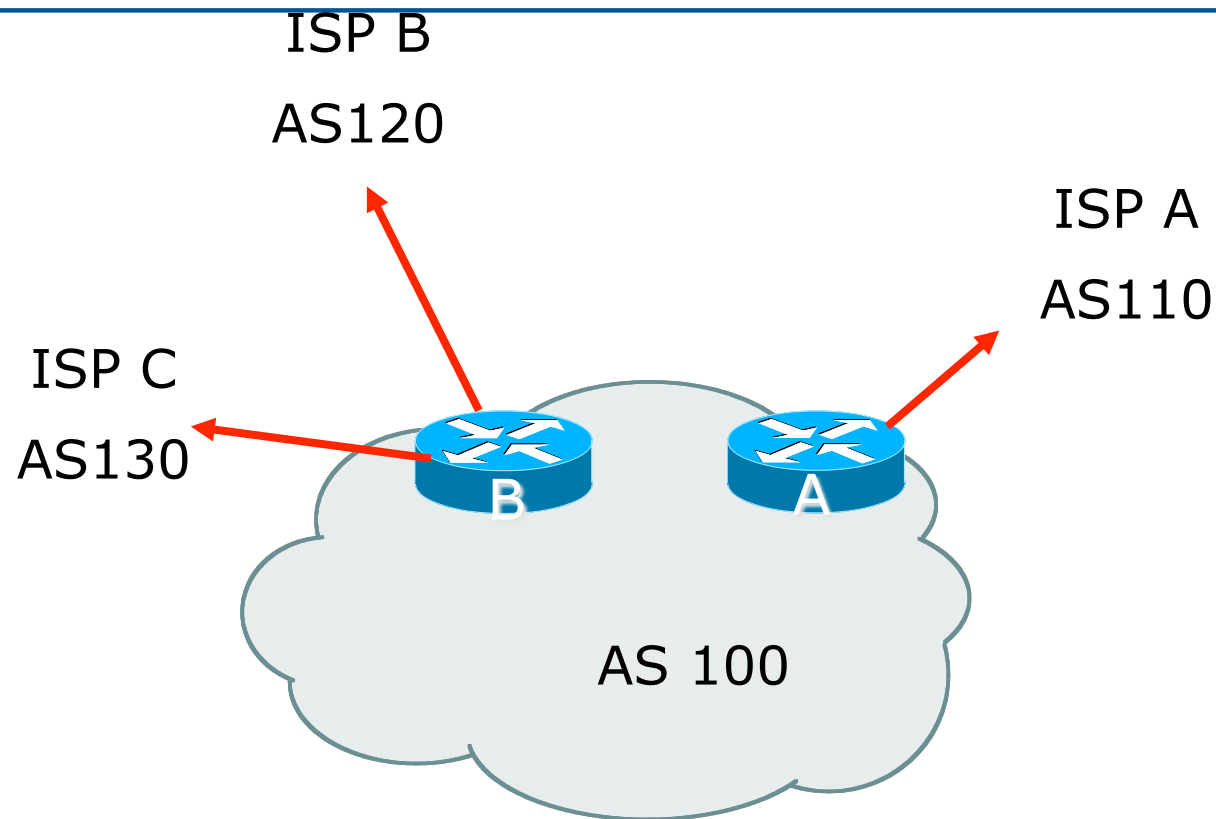
Three upstreams, unequal bandwidths

- Autonomous System has three upstreams
 - 16Mbps to ISP A
 - 8Mbps to ISP B
 - 4Mbps to ISP C
- What is the strategy here?
 - One option is full table from each
 - 3x 570k prefixes \Rightarrow 1710k paths
 - Other option is partial table and defaults from each
 - How??

Strategy

- Two external routers (gives router redundancy)
 - Do **NOT** need three routers for this
- Connect biggest bandwidth to one router
 - Most of inbound and outbound traffic will go here
- Connect the other two links to the second router
 - Provides maximum backup capacity if primary link fails
- Use the biggest link as default
 - Most of the inbound and outbound traffic will go here
- Do the traffic engineering on the two smaller links
 - Focus on regional traffic needs

Diagram



- ❑ Router A has 16Mbps circuit to ISP A
- ❑ Router B has 8Mbps and 4Mbps circuits to ISPs B&C

Outbound load-balancing strategy

- Available BGP feeds from Transit providers:
 - Full table
 - Customer prefixes and default
 - Default Route
- These are the common options
 - Very rare for any provider to offer anything different
 - Very rare for any provider to customise BGP feed for a customer

Outbound load-balancing strategy

- Accept only a default route from the provider with the **largest** connectivity, ISP A
 - Because most of the traffic is going to use this link
- If ISP A won't provide a default:
 - Still run BGP with them, but discard all prefixes
 - Point static default route to the upstream link
 - Distribute the default in the IGP
- Request the full table from ISP B & C
 - Most of this will be thrown away
 - (“Default plus customers” is not enough)

Outbound load-balancing strategy

- How to decide what to keep and what to discard from ISPs B & C?
 - Most traffic will use ISP A link — so we need to find a good/useful subset
- Discard prefixes transiting the global transit ISPs
 - Global transit ISPs generally appear in most non-local or regional AS-PATHs
- Discard prefixes with ISP A's ASN in the path
 - Makes more sense for traffic to those destinations to go via the link to ISP A

Outbound load-balancing strategy

- Global Transit (Tier-1) ISPs at the time of this exercise included:

209	CenturyLink	(Qwest)
701	VerizonBusiness	(UUNET)
1229	TeliaSonera	(Telia)
1239	Softbank	(Sprint)
1668	AOL TDN	
2914	NTT America	(NTT/Verio)
3549	Level 3	(GlobalCrossing)
3356	Level 3	
3561	CenturyLink	(Savvis, ex C&W)
7018	AT&T	

ISP B peering Inbound AS-PATH filter

```
ip as-path access-list 1 deny _209_
ip as-path access-list 1 deny _701_
ip as-path access-list 1 deny _1239_
ip as-path access-list 1 deny _3356_
ip as-path access-list 1 deny _3549_
ip as-path access-list 1 deny _3561_
ip as-path access-list 1 deny _2914_
ip as-path access-list 1 deny _7018_
!
ip as-path access-list 1 deny _ISPA_
ip as-path access-list 1 deny _ISPC_
!
ip as-path access-list 1 permit _ISPB$
ip as-path access-list 1 permit _ISPB_[0-9]+$
ip as-path access-list 1 permit _ISPB_[0-9]+_[0-9]+$
ip as-path access-list 1 permit _ISPB_[0-9]+_[0-9]+_[0-9]+$
ip as-path access-list 1 deny .*
```

Don't need ISPA and ISPC prefixes via ISPB

Outbound load-balancing strategy: ISP B peering configuration

- ❑ Part 1: Dropping Global Transit ISP prefixes
 - This can be fine-tuned if traffic volume is not sufficient
 - (More prefixes in = more traffic out)
- ❑ Part 2: Dropping prefixes transiting ISP A & C network
- ❑ Part 3: Permitting prefixes from ISP B, their BGP neighbours, and their neighbours, and their neighbours
 - More AS_PATH permit clauses, the more prefixes allowed in, the more egress traffic
 - Too many prefixes in will mean more outbound traffic than the link to ISP B can handle

Outbound load-balancing strategy

- ❑ Similar AS-PATH filter can be built for the ISP C BGP peering
- ❑ If the same prefixes are heard from both ISP B and C, then establish proximity of their origin ASN to ISP B or C
 - e.g. ISP B might be in Japan, with the neighbouring ASN in Europe, yet ISP C might be in Europe
 - Transit to the ASN via ISP C makes more sense in this case

Inbound load-balancing strategy

- ❑ The largest outbound link should announce just the aggregate
- ❑ The other links should announce:
 - a) The aggregate with AS-PATH prepend
 - b) Subprefixes of the aggregate, chosen according to traffic volumes to those subprefixes, and according to the services on those subprefixes
- ❑ Example:
 - Link to ISP B could be used just for Broadband/Dial customers — so number all such customers out of one contiguous subprefix
 - Link to ISP C could be used just for commercial leased line customers — so number all such customers out of one contiguous subprefix

Router A: eBGP Configuration

Example

```
router bgp 100
  network 100.10.0.0 mask 255.255.224.0
  neighbor 122.102.10.1 remote 110
  neighbor 122.102.10.1 prefix-list default in
  neighbor 122.102.10.1 prefix-list aggregate out
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list aggregate permit 100.10.0.0/19
!
```

Router B: eBGP Configuration

Example

```
router bgp 100
  network 100.10.0.0 mask 255.255.224.0
  neighbor 120.103.1.1 remote 120
  neighbor 120.103.1.1 filter-list 1 in
  neighbor 120.103.1.1 prefix-list ISP-B out
  neighbor 120.103.1.1 route-map to-ISP-B out
  neighbor 121.105.2.1 remote 130
  neighbor 121.105.2.1 filter-list 2 in
  neighbor 121.105.2.1 prefix-list ISP-C out
  neighbor 121.105.2.1 route-map to-ISP-C out
!
ip prefix-list aggregate permit 100.10.0.0/19
!
...next slide
```

Router B: eBGP Configuration

Example

```
ip prefix-list ISP-B permit 100.10.0.0/19
ip prefix-list ISP-B permit 100.10.0.0/21
!
ip prefix-list ISP-C permit 100.10.0.0/19
ip prefix-list ISP-C permit 100.10.28.0/22
!
route-map to-ISP-B permit 10
  match ip address prefix-list aggregate
  set as-path prepend 100
!
route-map to-ISP-B permit 20
!
route-map to-ISP-C permit 10
  match ip address prefix-list aggregate
  set as-path prepend 100 100
!
route-map to-ISP-C permit 20
```

/21 to ISP B
"dial customers"

/22 to ISP C
"biz customers"

e.g. Single
prepend on ISP B
link

e.g. Dual prepend
on ISP C link

What about outbound backup?

- We have:
 - Default route from ISP A by eBGP
 - Mostly discarded full table from ISPs B&C
- Strategy:
 - Originate default route by OSPF on Router A (with metric 10) — link to ISP A
 - Originate default route by OSPF on Router B (with metric 30) — links to ISPs B & C
 - Plus on Router B:
 - Static default route to ISP B with distance 240
 - Static default route to ISP C with distance 245
 - When link goes down, static route is withdrawn

Outbound backup: steady state

- Steady state (all links up and active):
 - Default route is to Router A — OSPF metric 10
 - (Because default learned by eBGP \Rightarrow default is in RIB \Rightarrow OSPF will originate default)
 - Backup default is to Router B — OSPF metric 20
 - eBGP prefixes learned from upstreams distributed by iBGP throughout backbone
 - (Default can be filtered in iBGP to avoid “RIB failure error”)

Outbound backup: failure examples

- Link to ISP A down, to ISPs B&C up:
 - Default route is to Router B — OSPF metric 20
 - (eBGP default gone from RIB, so OSPF on Router A withdraws the default)
- Above is true if link to B or C is down as well
- Link to ISPs B & C down, link to ISP A is up:
 - Default route is to Router A — OSPF metric 10
 - (static defaults on Router B removed from RIB, so OSPF on Router B withdraws the default)

Other considerations

- ❑ Default route should not be propagated to devices terminating non-transit peers and customers
- ❑ Rarely any need to carry default in iBGP
 - Best to filter out default in iBGP mesh peerings
- ❑ Still carry other eBGP prefixes across iBGP mesh
 - Otherwise routers will follow default route rules resulting in suboptimal traffic flow
 - Not a big issue because not carrying full table

Router A: iBGP Configuration

Example

```
router bgp 100
  network 100.10.0.0 mask 255.255.224.0
  neighbor ibgp-peers peer-group
  neighbor ibgp-peers remote-as 100
  neighbor ibgp-peers prefix-list ibgp-filter out
  neighbor 100.10.0.2 peer-group ibgp-peers
  neighbor 100.10.0.3 peer-group ibgp-peers
!
ip prefix-list ibgp-filter deny 0.0.0.0/0
ip prefix-list ibgp-filter permit 0.0.0.0/0 le 32
!
```



bandwidths:

Summary

- Example based on many deployed working multihoming/loadbalancing topologies
- Many variations possible — this one is:
 - Easy to tune
 - Light on border router resources
 - Light on backbone router infrastructure
 - Sparse BGP table \Rightarrow faster convergence



BGP Multihoming Techniques

- Why Multihome?
- Definition & Options
- How to Multihome
- Principles & Addressing
- Basic Multihoming
- “BGP Traffic Engineering”
- **Using Communities**

Using Communities for BGP Traffic Engineering



How they are used in practice
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

- 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
- 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
 - But so far... nothing more... ☹
 - Collection of ISP communities at www.onesc.net/communities
 - NANOG Tutorial: 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...
```

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 ...

ISP Example: Sprint

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 ...

More info at https://www.sprint.net/index.php?p=policy_bgp

String Resulting AS Path to ASXXX in North America

65010:XXX	Do not advertise to ASXXX
-----------	---------------------------

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
65520:nnnn	50	only within country <nnnn> (see country list below)
65530:nnnn	50	only within region <nnnn> (see region list below)
2914:435	50	only beyond the connected country
2914:436	50	only beyond the connected region
2914:450	96	customer fallback
2914:460	98	peer backup
2914:470	100	peer
2914:480	110	customer backup
2914:490	120	customer default
2914:666		blackhole

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
2914:423	prepends o/b to peer 3x
2914:429	do not advertise to any peer
2914:439	do not advertise to any peer outside region

Note: 2914 is the ASN prepend in all cases. If used, 654xx:nnn overrides 655xx:nnn and 2914:429, 655xx:nnn overrides the 2914:42x communities.

Customers wanting to alter their route announcements to selected peers.

NTT Communications BGP customers may choose to prepend to selected peers with the following communities, where *nnn* is the peer's ASN:

ISP Example: NTT

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

ISP Example: Verizon Europe

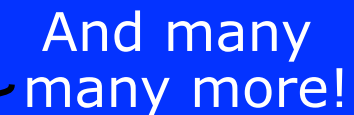
```
aut-num:          AS702
descr:           Verizon Business EMEA - Commercial IP service provider in Europe
<snip>
remarks:         -----
                 Verizon Business filters out inbound prefixes longer than /24.
                 We also filter any networks within AS702:RS-INBOUND-FILTER.
                 -----
                 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's
                 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.
<snip>
```

← And many more!

ISP Example: Telia

```
aut-num:          AS1299
descr:           TeliaSonera International Carrier
<snip>
remarks:         -----
remarks:         BGP COMMUNITY SUPPORT FOR AS1299 TRANSIT CUSTOMERS:
remarks:         Community Action (default local pref 200)
remarks:         -----
remarks:         1299:50 Set local pref 50 within AS1299 (lowest possible)
remarks:         1299:150 Set local pref 150 within AS1299 (equal to peer, backup)
remarks:         European peers
remarks:         Community Action
remarks:         -----
remarks:         1299:200x All peers Europe incl:
remarks:         1299:250x Sprint/1239
remarks:         1299:251x Savvis/3561
remarks:         1299:252x NTT/2914
remarks:         1299:253x Zayo/Abovenet/6461
remarks:         1299:254x FT/5511
remarks:         1299:255x GBLX/3549
remarks:         1299:256x Level3/3356
<snip>
remarks:         Where x is number of prepends (x=0,1,2,3) or do NOT announce (x=9)
```

And many
many more!



ISP Example: BT Ignite

```
aut-num:      AS5400
descr:       BT Ignite European Backbone
<snip>
remarks:     The following BGP communities can be set by BT
remarks:     BGP customers to affect announcements to major peers.
remarks:
remarks:     5400:NXXX
remarks:     N=1          not announce
remarks:     N=2          prepend an extra "5400 5400" on announcement
remarks:     Valid values for XXX:
remarks:     000          All peers and transits
remarks:     500          All transits
remarks:     503          Level3 AS3356
remarks:     509          Telia AS1299
remarks:     510          NTT Verio AS2914
remarks:     002          Sprint AS1239
remarks:     003          Savvis AS3561
remarks:     004          C&W AS1273
remarks:     005          Verizon EMEA AS702
remarks:     014          DTAG AS3320
remarks:     016          Opentransit AS5511
remarks:     018          GlobeInternet Tata AS6453
remarks:     023          Tinet AS3257
remarks:     027          Telia AS1299
remarks:     045          Telecom Italia AS6762
remarks:     073          Eurorings AS286
remarks:     169          Cogent AS174
<snip>
```

And many
more!



ISP Example: Level3

```
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
remarks:         65002:0   - prepend twice  to all peers
remarks:         65002:XXX - prepend twice  at peerings to AS XXX
<snip>
remarks:         -----
remarks:         customer traffic engineering communities - LocalPref
remarks:         -----
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:         -----
remarks:         customer traffic engineering communities - Blackhole
remarks:         -----
remarks:         3356:9999 - blackhole (discard) traffic
<snip>
```

And 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 the previous ISP and configuration examples as a guideline

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 previously
- ❑ 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 ← Don't announce anywhere
ip community-list 101 permit 100:5001 ← Single prepend to all
ip community-list 102 permit 100:5002
ip community-list 103 permit 100:5003
ip community-list 110 permit 100:5010 ← Don't announce to peer 1
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 ← Single prepend to peer 2
ip community-list 122 permit 100:5022
ip community-list 123 permit 100:5023
```

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 c.c.c.c route-map upstream out
  neighbor c.c.c.c 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 prefix-list 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
 - A default is often just as good
 - If customers want 570k IPv4 prefixes, charge them money for it 😊

BGP Multihoming Techniques



Philip Smith

<philip@nsrc.org>

SANOG 27

25th-27th January 2016

Kathmandu