

BGP for Internet Service Providers

Philip Smith <pfs@cisco.com>

AfNOG 3, Lome, Togo

Presentation Slides

Cisco.com

- **Will be available on**
www.cisco.com/public/cons/seminars/AfNOG3
- **Feel free to ask questions any time**

BGP for Internet Service Providers

Cisco.com

- **BGP Basics (quick recap)**
- **Scaling BGP**
- **Deploying BGP in an ISP network**
- **Multihoming Examples**

BGP Basics

What is this BGP thing?

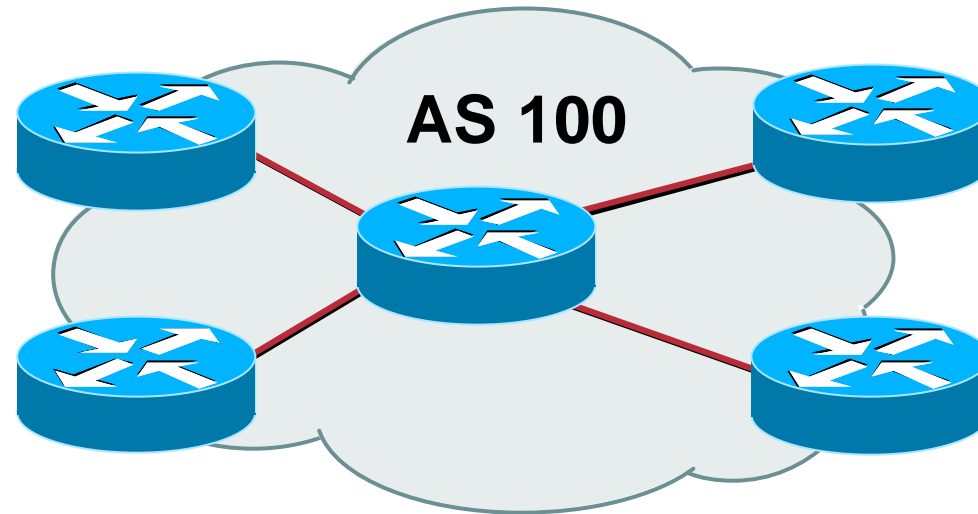
Border Gateway Protocol

Cisco.com

- **Routing Protocol used to exchange routing information between networks**
exterior gateway protocol
- **RFC1771**
work in progress to update
`draft-ietf-idr-bgp4-17.txt`

Autonomous System (AS)

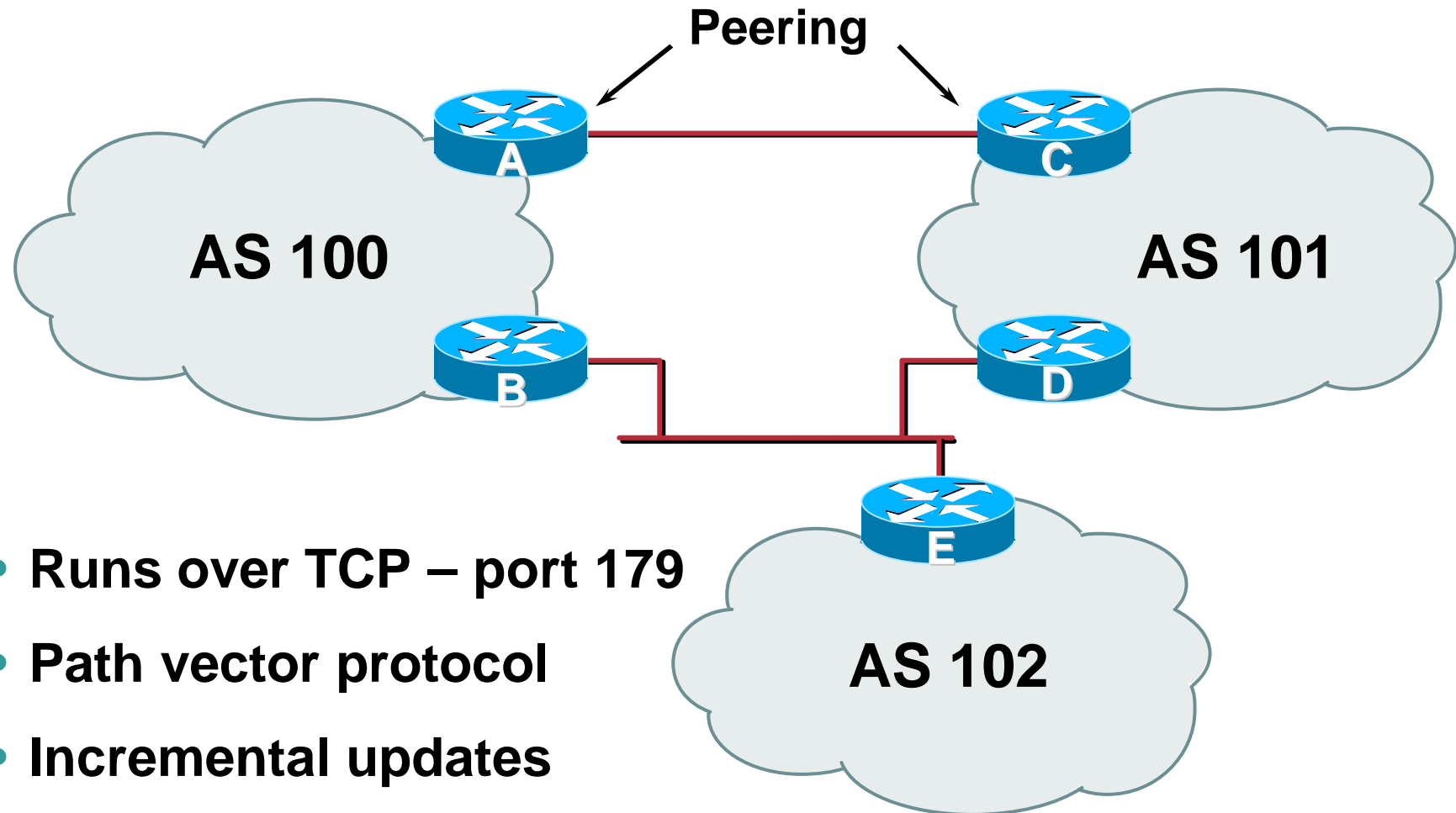
Cisco.com



- **Collection of networks with same routing policy**
- **Single routing protocol**
- **Usually under single ownership, trust and administrative control**

BGP Basics

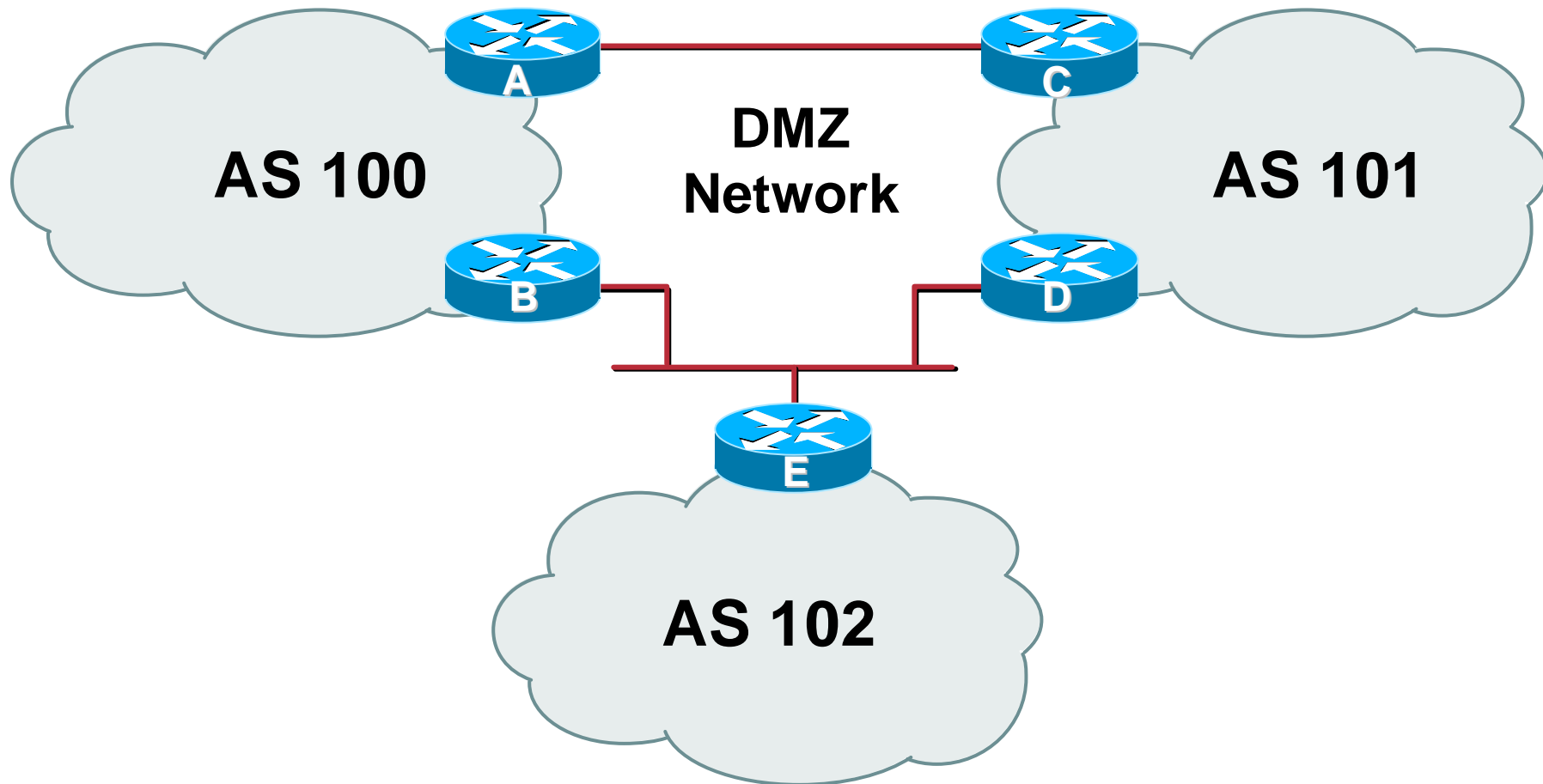
Cisco.com



- Runs over TCP – port 179
- Path vector protocol
- Incremental updates
- “Internal” & “External” BGP

Demarcation Zone (DMZ)

Cisco.com



- **Shared network between ASes**

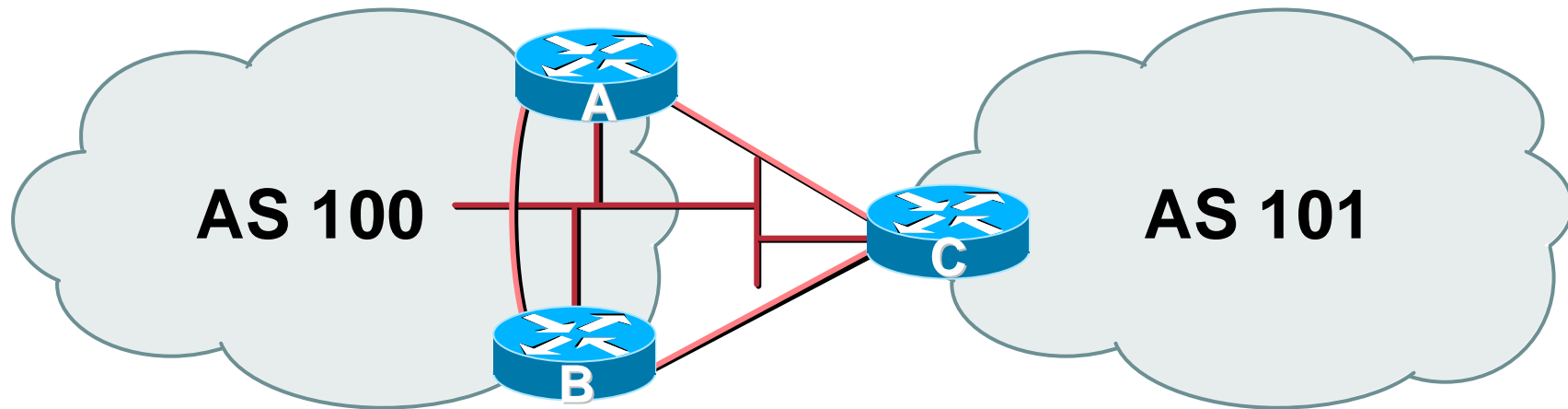
BGP General Operation

Cisco.com

- **Learns multiple paths via internal and external BGP speakers**
- **Picks the best path and installs in the forwarding table**
- **Best path is sent to external BGP neighbours**
- **Policies applied by influencing the best path selection**

External BGP Peering (eBGP)

Cisco.com



- **Between BGP speakers in different AS**
- **Should be directly connected**
- **Never** run an IGP between eBGP peers

Configuring External BGP

Cisco.com

Router A in AS100

```
interface ethernet 5/0
ip address 222.222.10.2 255.255.255.240
router bgp 100
  network 220.220.8.0 mask 255.255.252.0
  neighbor 222.222.10.1 remote-as 101
  neighbor 222.222.10.1 prefix-list RouterC in
  neighbor 222.222.10.1 prefix-list RouterC out
```

Router C in AS101

```
interface ethernet 1/0/0
ip address 222.222.10.1 255.255.255.240
router bgp 101
  network 220.220.16.0 mask 255.255.240.0
  neighbor 222.222.10.2 remote-as 100
  neighbor 222.222.10.2 prefix-list RouterA in
  neighbor 222.222.10.2 prefix-list RouterA out
```

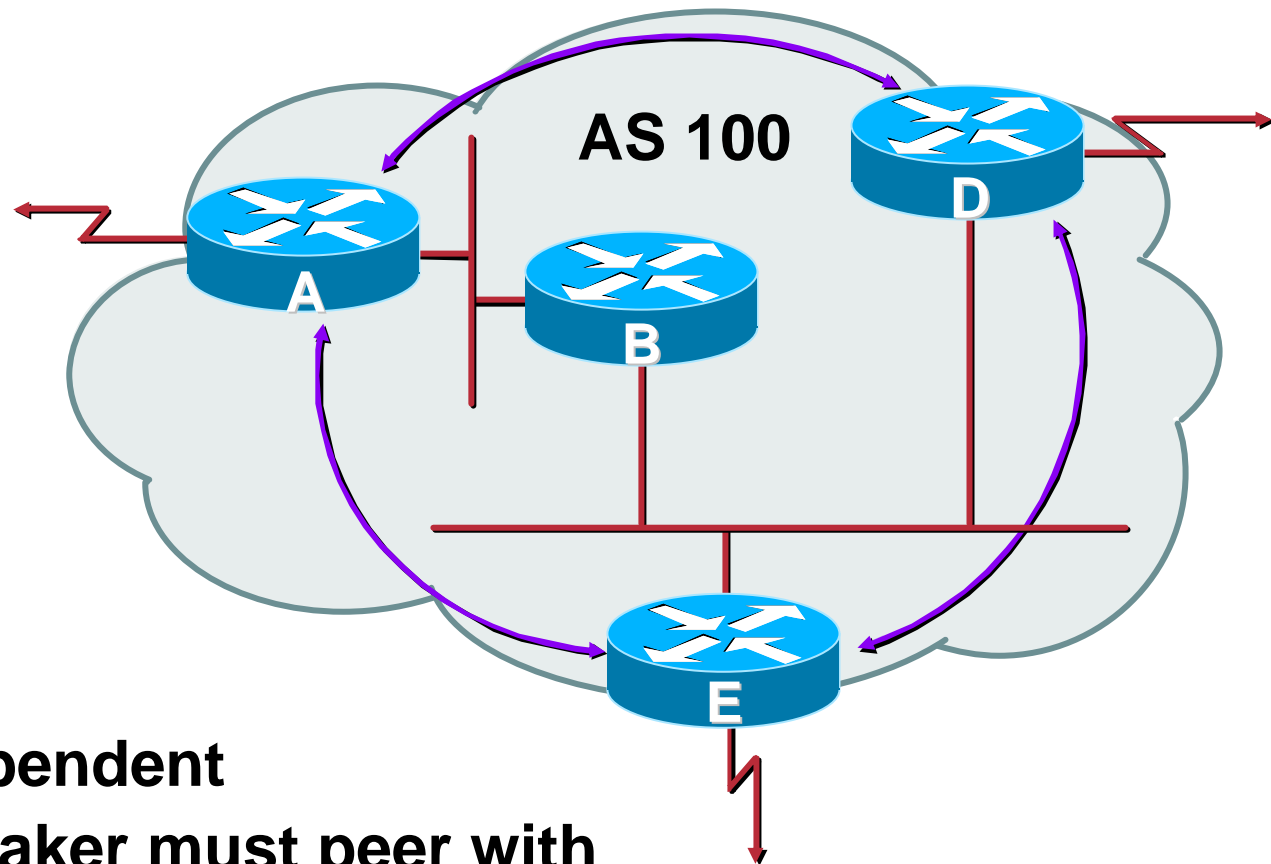
Internal BGP (iBGP)

Cisco.com

- **BGP peer within the same AS**
- **Not required to be directly connected**
- **iBGP speakers need to be fully meshed**
 - they originate connected networks**
 - they do not pass on prefixes learned from other iBGP speakers**

Internal BGP Peering (iBGP)

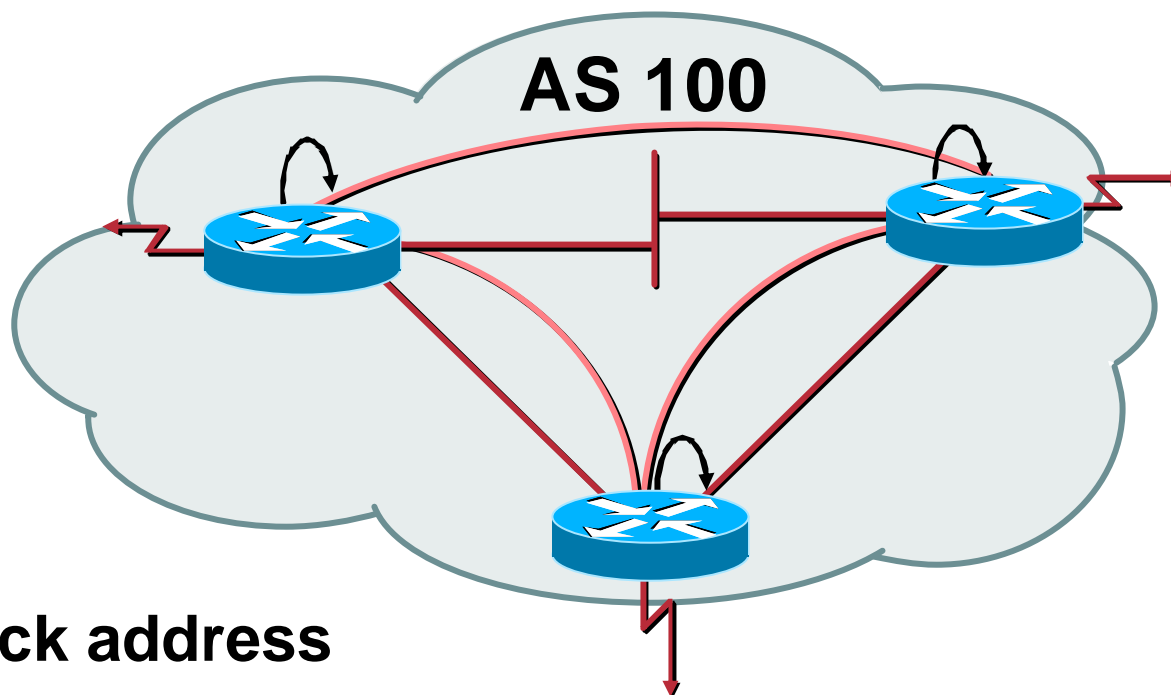
Cisco.com



- Topology independent
- Each iBGP speaker must peer with every other iBGP speaker in the AS

Peering to Loop-back Address

Cisco.com



- **Peer with loop-back address**
Loop-back interface does not go down – ever!
- **iBGP session is not dependent on state of a single interface**
- **iBGP session is not dependent on physical topology**

Configuring Internal BGP

Router A

```
interface loopback 0
ip address 215.10.7.1 255.255.255.255
router bgp 100
  network 220.220.1.0
  neighbor 215.10.7.2 remote-as 100
  neighbor 215.10.7.2 update-source loopback0
  neighbor 215.10.7.3 remote-as 100
  neighbor 215.10.7.3 update-source loopback0
```

Router B

```
interface loopback 0
ip address 215.10.7.2 255.255.255.255
router bgp 100
  network 220.220.5.0
  neighbor 215.10.7.1 remote-as 100
  neighbor 215.10.7.1 update-source loopback0
  neighbor 215.10.7.3 remote-as 100
  neighbor 215.10.7.3 update-source loopback0
```

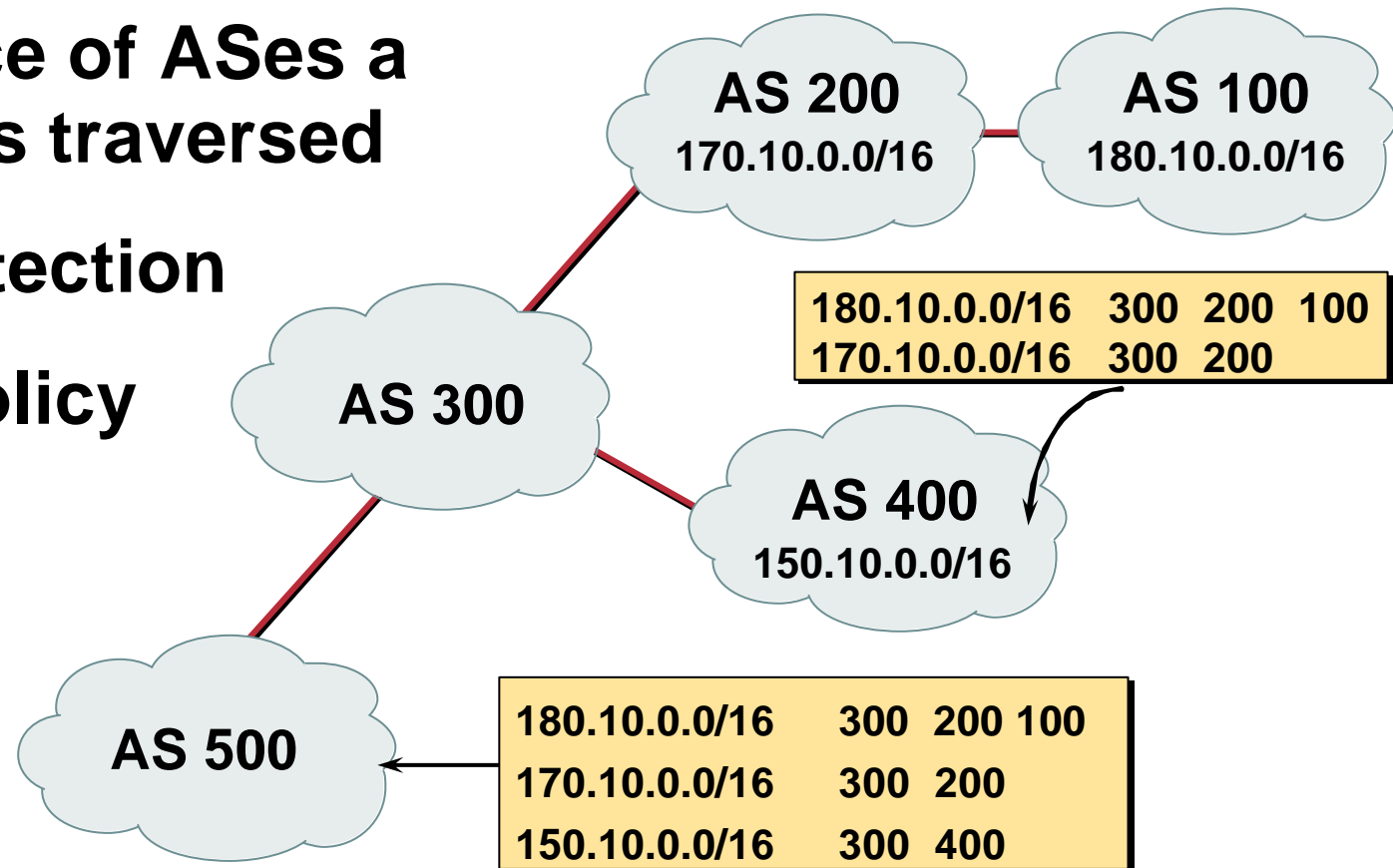
BGP Attributes

Recap

AS-Path

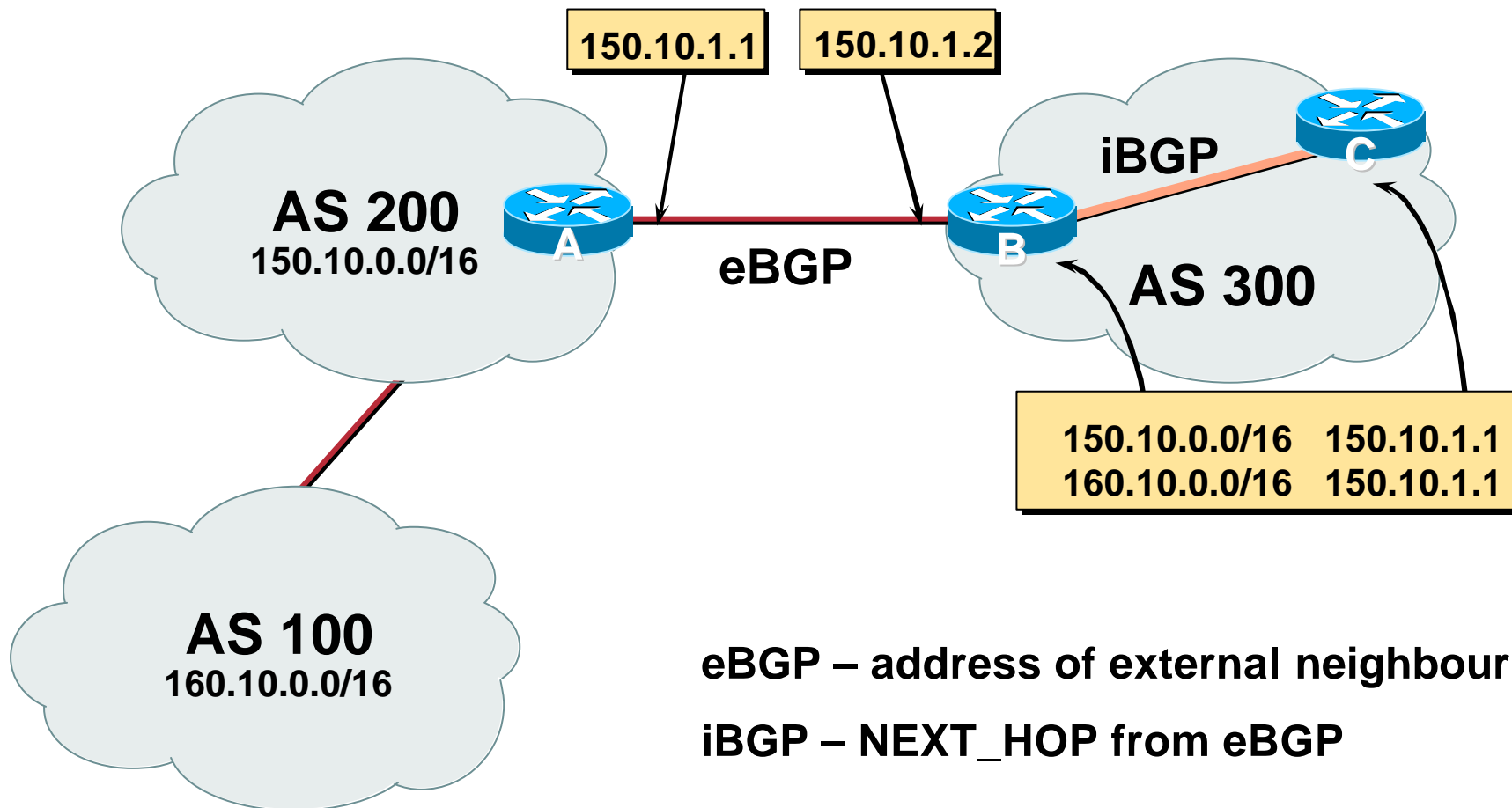
Cisco.com

- Sequence of ASes a route has traversed
- Loop detection
- Apply policy



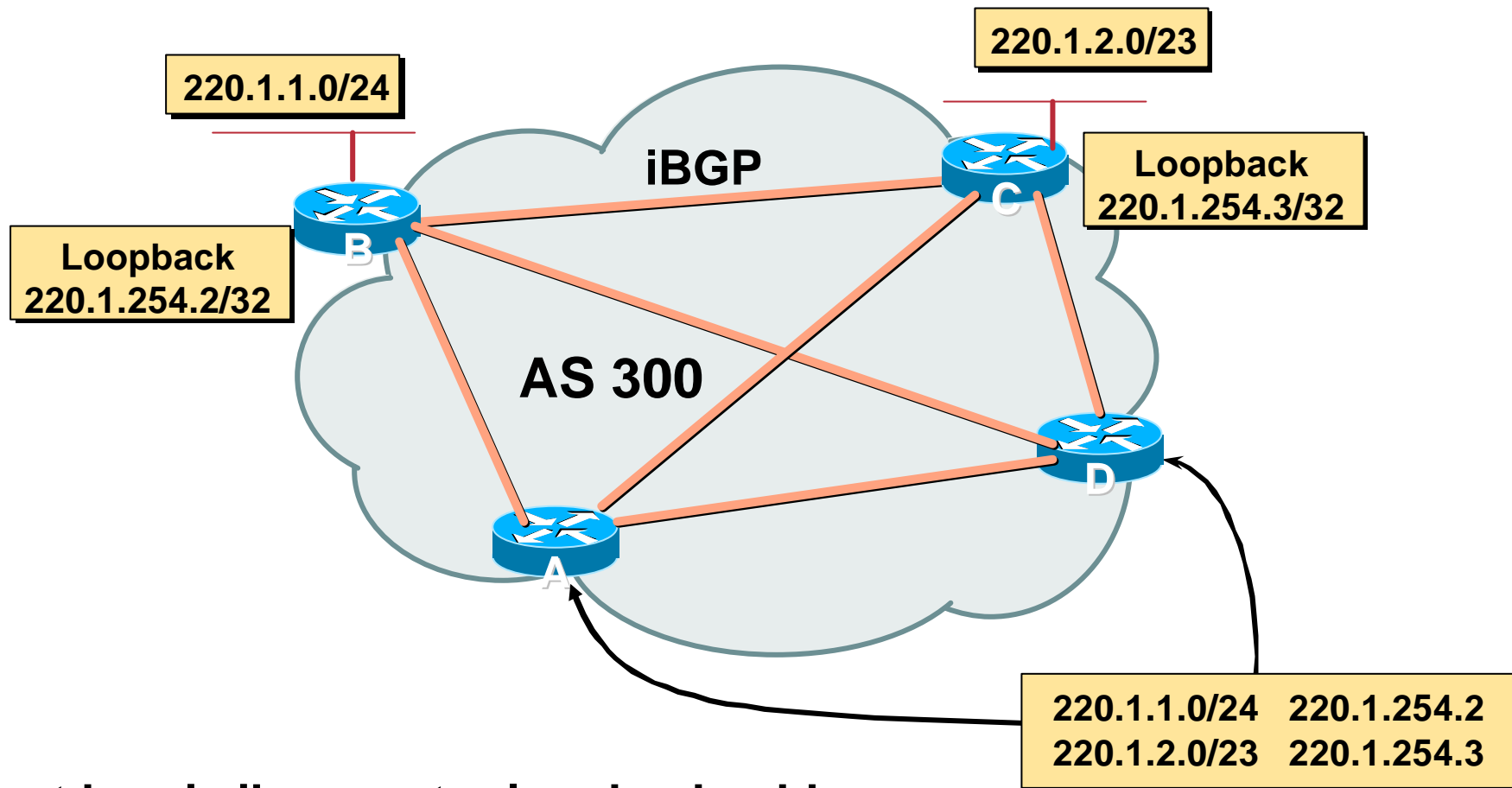
Next Hop

Cisco.com



iBGP Next Hop

Cisco.com

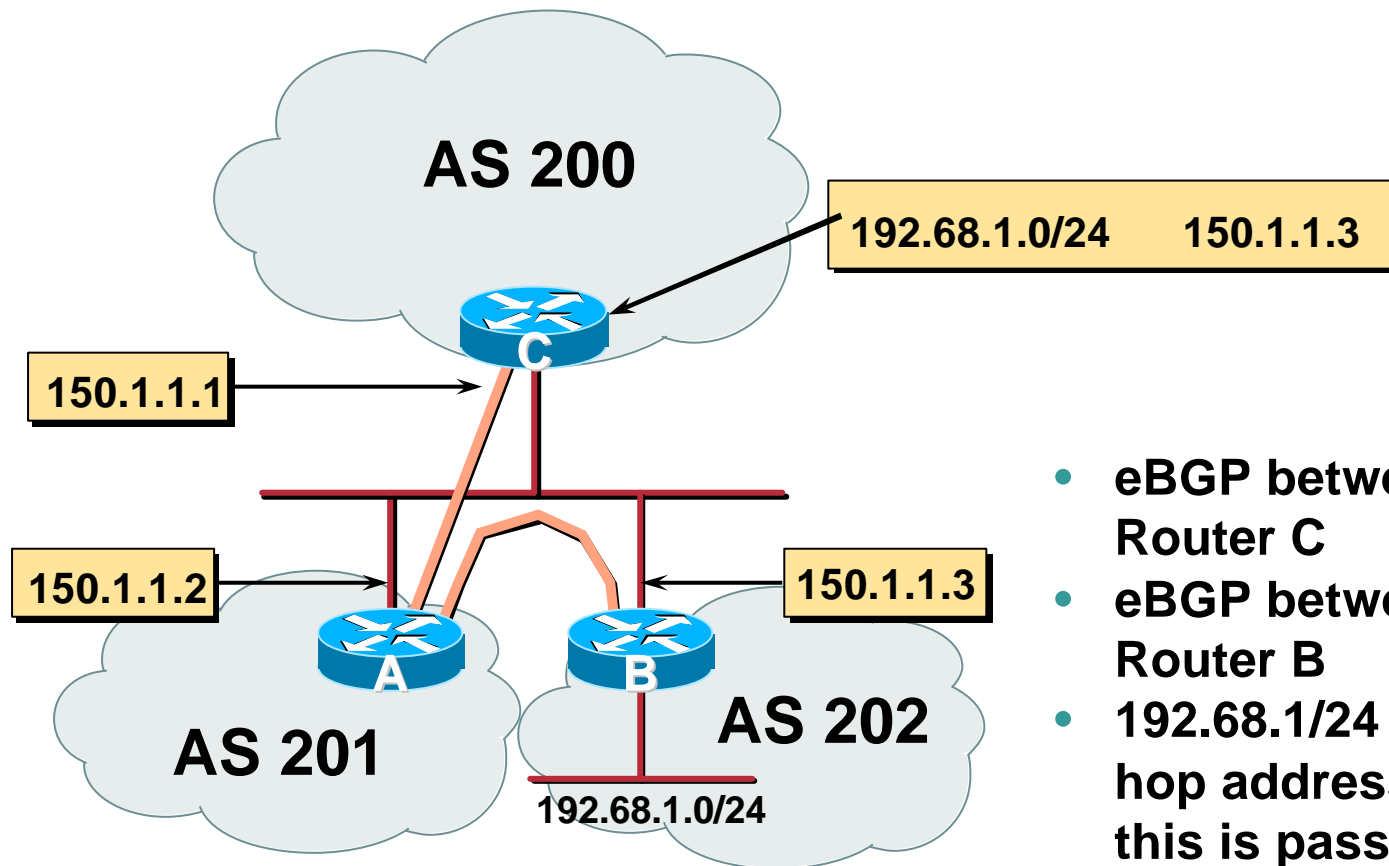


Next hop is ibgp router loopback address

Recursive route look-up

Third Party Next Hop

Cisco.com



- eBGP between Router A and Router C
- eBGP between Router A and Router B
- 192.68.1/24 prefix has next hop address of 150.1.1.3 – this is passed on to Router C instead of 150.1.1.2

Next Hop (summary)

Cisco.com

- **IGP should carry route to next hops**
- **Recursive route look-up**
- **Unlinks BGP from actual physical topology**
- **Allows IGP to make intelligent forwarding decision**

Origin

Cisco.com

- **Conveys the origin of the prefix**
- **“Historical” attribute**
- **Influences best path selection**
- **Three values: IGP, EGP, incomplete**
 - IGP – generated by BGP network statement**
 - EGP – generated by EGP**
 - incomplete – redistributed from another routing protocol**

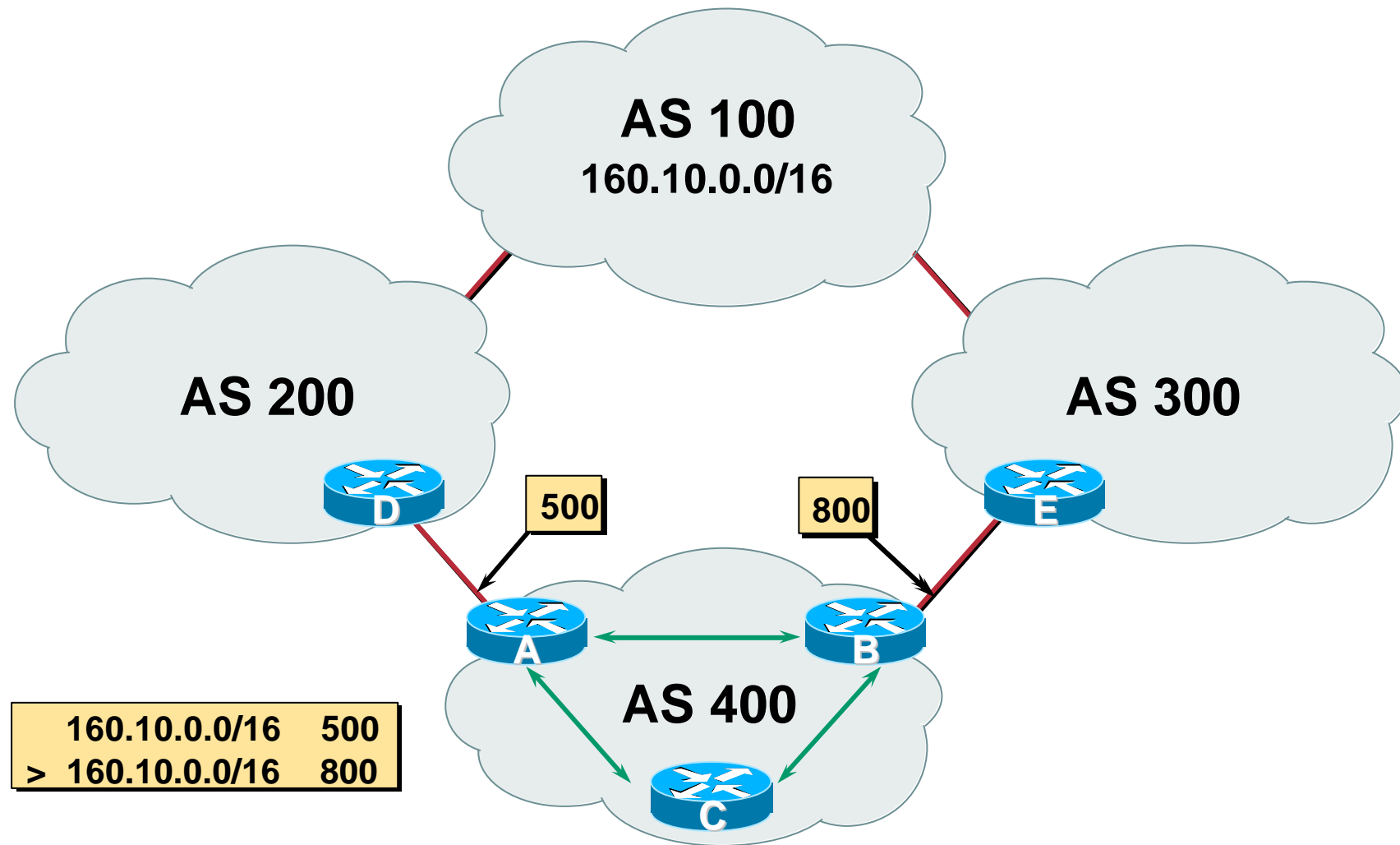
Aggregator

Cisco.com

- **Conveys the IP address of the router/BGP speaker generating the aggregate route**
- **Useful for debugging purposes**
- **Does not influence best path selection**

Local Preference

Cisco.com



Local Preference

Cisco.com

- **Local to an AS – non-transitive**
Default local preference is 100
- **Used to influence BGP path selection**
determines best path for *outbound* traffic
- **Path with highest local preference wins**

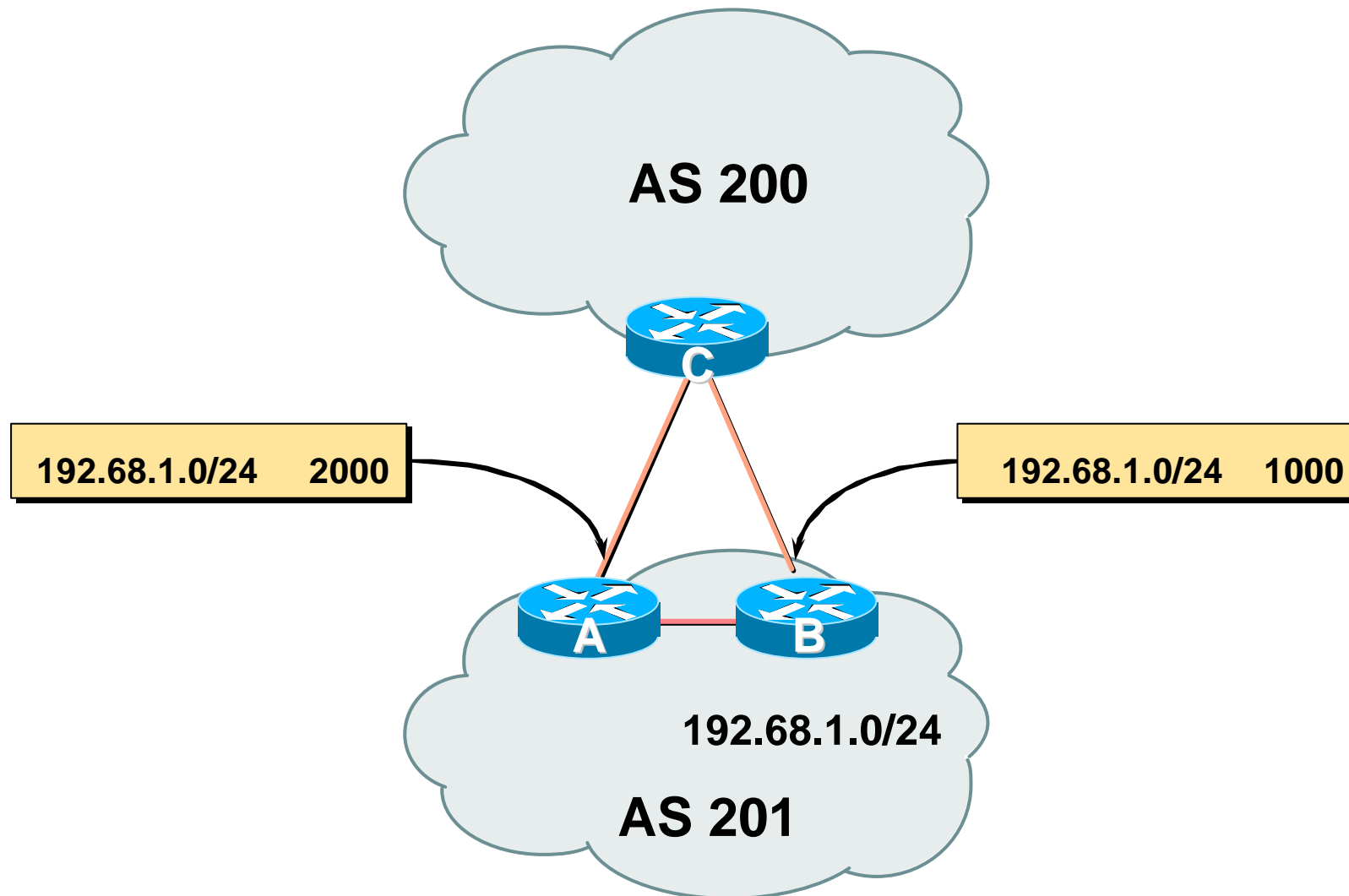
Local Preference

- **Configuration of Router B:**

```
router bgp 400
  neighbor 220.5.1.1 remote-as 300
  neighbor 220.5.1.1 route-map local-pref in
!
route-map local-pref permit 10
  match ip address prefix-list MATCH
  set local-preference 800
!
ip prefix-list MATCH permit 160.10.0.0/16
```

Multi-Exit Discriminator (MED)

Cisco.com



Multi-Exit Discriminator

Cisco.com

- **Inter-AS – non-transitive**
- **Used to convey the relative preference of entry points**
 - determines best path for *inbound* traffic
- **Comparable if paths are from same AS**
- **IGP metric can be conveyed as MED**
 - set metric-type internal* in route-map

Multi-Exit Discriminator

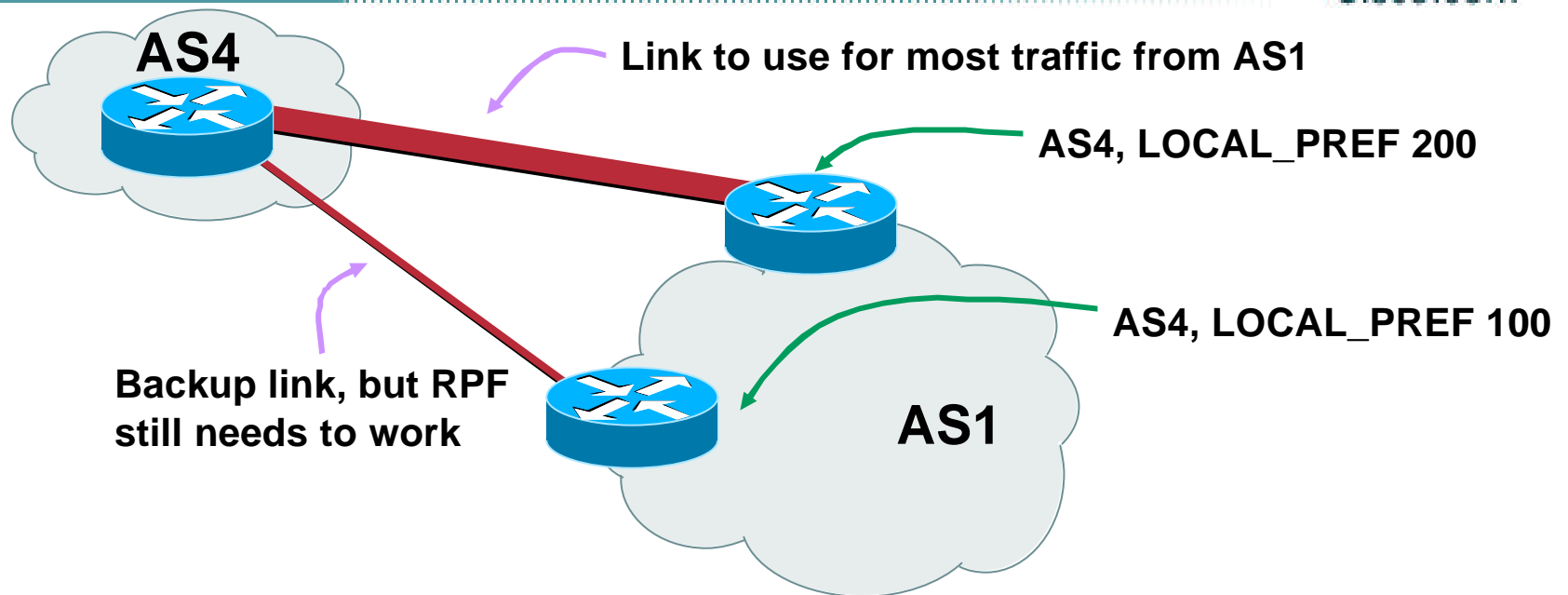
Cisco.com

- **Configuration of Router B:**

```
router bgp 400
  neighbor 220.5.1.1 remote-as 200
  neighbor 220.5.1.1 route-map set-med out
!
route-map set-med permit 10
  match ip address prefix-list MATCH
  set metric 1000
!
ip prefix-list MATCH permit 192.68.1.0/24
```

Weight – Used to Deploy RPF

Cisco.com



- Local to router on which it's configured
Not really an attribute
- route-map: **set weight**
- Highest weight wins over all valid paths
- Weight customer eBGP on edge routers to allow RPF to work correctly

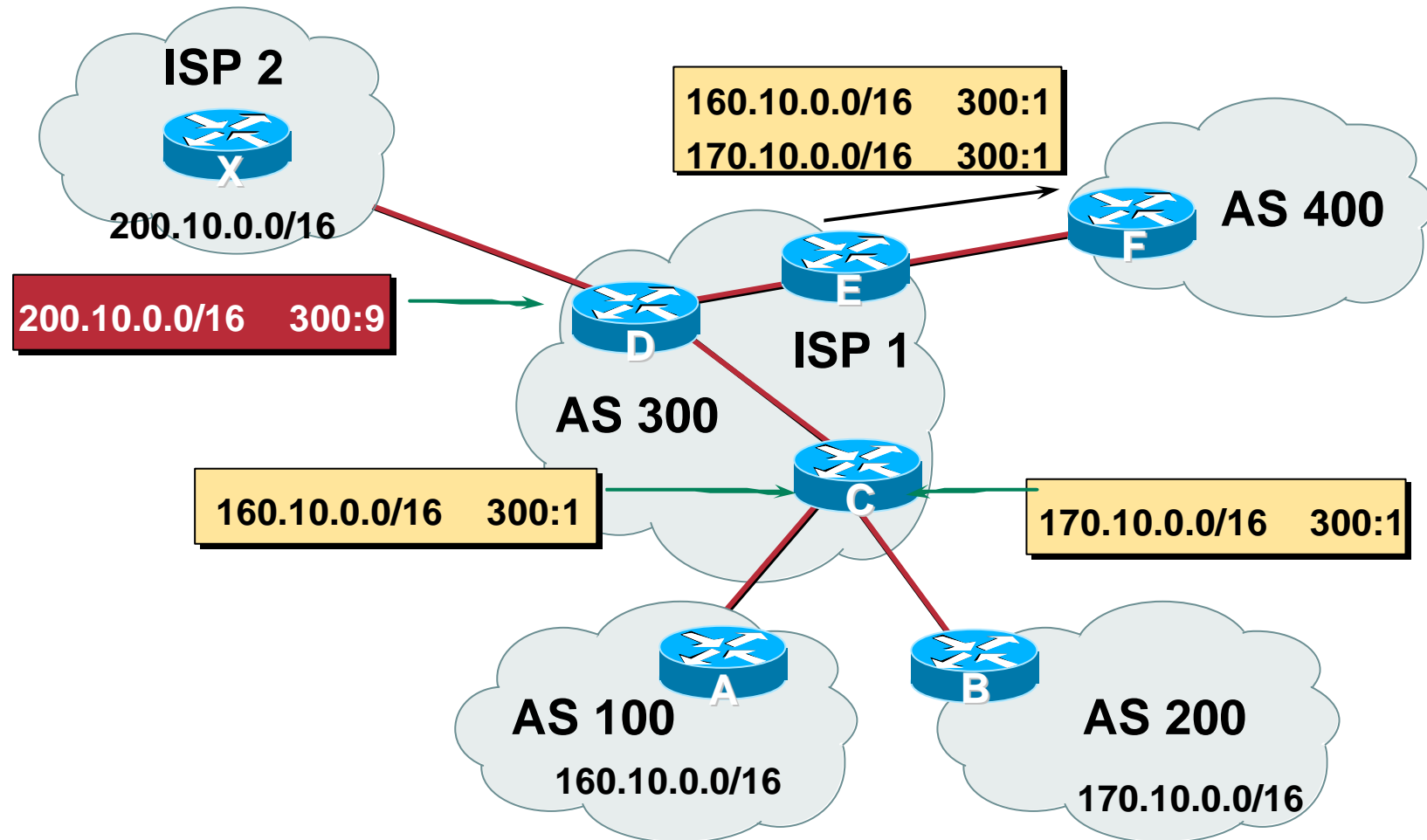
Community

Cisco.com

- **BGP attribute**
- **Described in RFC1997**
- **32 bit integer**
 - Represented as two 16 bit integers**
- **Used to group destinations**
 - Each destination could be member of multiple communities**
- **Community attribute carried across AS's**
- **Very useful in applying policies**

Community

Cisco.com



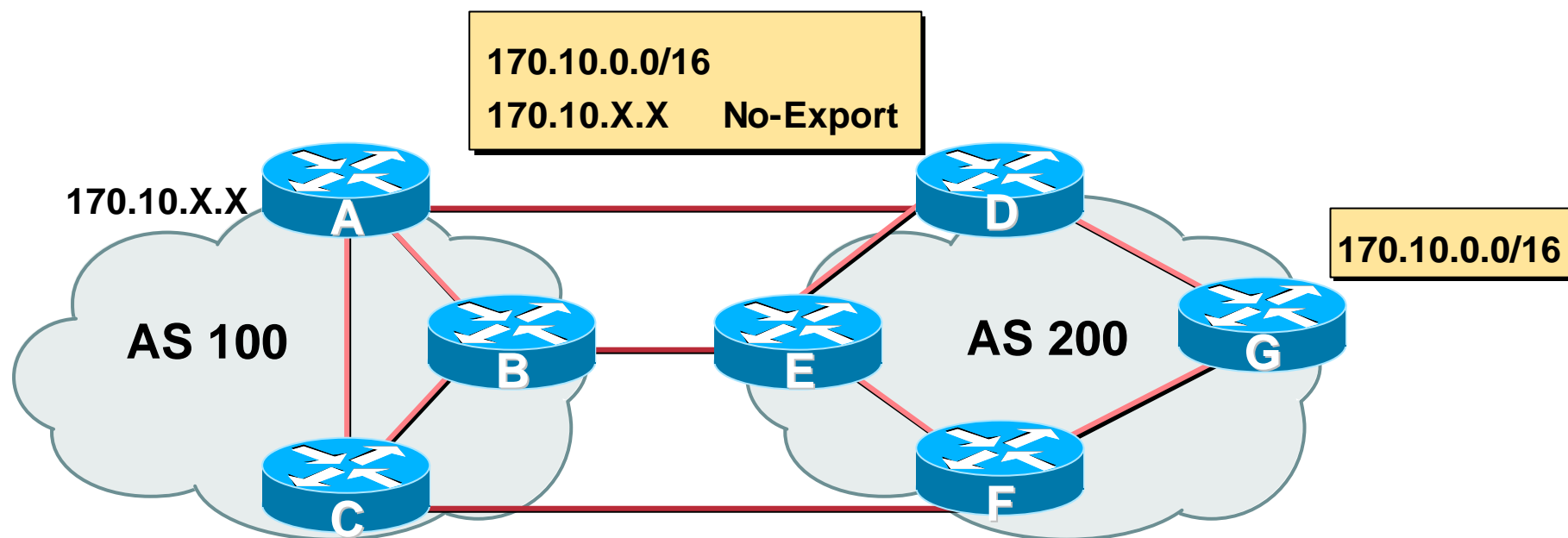
Well-Known Communities

Cisco.com

- **no-export**
do not advertise to eBGP peers
- **no-advertise**
do not advertise to any peer
- **local-AS**
do not advertise outside local AS (only used with confederations)

No-Export Community

Cisco.com



- AS100 announces aggregate and subprefixes
aim is to improve loadsharing by leaking subprefixes
- Subprefixes marked with **no-export** community
- Router G in AS200 does not announce prefixes with **no-export** community set

BGP Path Selection Algorithm

Why Is This the Best Path?

BGP Path Selection Algorithm

Cisco.com

- **Do not consider path if no route to next hop**
- **Do not consider iBGP path if not synchronised (Cisco IOS)**
- **Highest weight (local to router)**
- **Highest local preference (global within AS)**
- **Prefer locally originated route**
- **Shortest AS path**

BGP Path Selection Algorithm (continued)

Cisco.com

- **Lowest origin code**

IGP < EGP < incomplete

- **Lowest Multi-Exit Discriminator (MED)**

If *bgp deterministic-med*, order the paths before comparing

If *bgp always-compare-med*, then compare for all paths

otherwise MED only considered if paths are from the same AS (default)

BGP Path Selection Algorithm (continued)

Cisco.com

- Prefer eBGP path over iBGP path
- Path with lowest IGP metric to next-hop
- Lowest router-id (originator-id for reflected routes)
- Shortest Cluster-List

Client **must** be aware of Route Reflector attributes!

- Lowest neighbour IP address

Applying Policy with BGP

Control!

Applying Policy with BGP

Cisco.com

- **Applying Policy**

Decisions based on AS path, community or the prefix

Rejecting/accepting selected routes

Set attributes to influence path selection

- **Tools:**

Prefix-list (filter prefixes)

Filter-list (filter ASes)

Route-maps and communities

Policy Control

Prefix List

Cisco.com

- Filter routes based on prefix
- Inbound and Outbound

```
router bgp 200
  neighbor 220.200.1.1 remote-as 210
  neighbor 220.200.1.1 prefix-list PEER-IN in
  neighbor 220.200.1.1 prefix-list PEER-OUT out
!
ip prefix-list PEER-IN deny 218.10.0.0/16
ip prefix-list PEER-IN permit 0.0.0.0/0 le 32
ip prefix-list PEER-OUT permit 215.7.0.0/16
```

Policy Control

Filter List

Cisco.com

- Filter routes based on AS path
- Inbound and Outbound

```
router bgp 100
  neighbor 220.200.1.1 remote-as 210
  neighbor 220.200.1.1 filter-list 5 out
  neighbor 220.200.1.1 filter-list 6 in
!
ip as-path access-list 5 permit ^200$
ip as-path access-list 6 permit ^150$
```

Policy Control

Regular Expressions

Cisco.com

- **Like Unix regular expressions**
 - .** Match one character
 - *** Match any number of preceding expression
 - +** Match at least one of preceding expression
 - ^** Beginning of line
 - \$** End of line
 - _** Beginning, end, white-space, brace
 - |** Or
 - ()** brackets to contain expression

Policy Control

Regular Expressions

Cisco.com

- **Simple Examples**

| | |
|----------------------|---|
| .* | Match anything |
| .+ | Match at least one character |
| ^\$ | Match routes local to this AS |
| _1800\$ | Originated by 1800 |
| ^1800_ | Received from 1800 |
| _1800_ | Via 1800 |
| _790_1800_ | Passing through 1800 then 790 |
| _(1800_)+ | Match at least one of 1800 in sequence |
| _\\(65350\\)_ | Via 65350 (confederation AS) |

Policy Control

Route Maps

Cisco.com

- A route-map is like a “programme” for IOS
- Has “line” numbers, like programmes
- Each line is a separate condition/action
- Concept is basically:
 - if *match* then do *expression* and *exit*
 - else
 - if *match* then do *expression* and *exit*
 - else *etc*

Policy Control

Route Maps

Cisco.com

- Example using prefix-lists

```
router bgp 100
  neighbor 1.1.1.1 route-map infilter in
  !
  route-map infilter permit 10
    match ip address prefix-list HIGH-PREF
    set local-preference 120
  !
  route-map infilter permit 20
    match ip address prefix-list LOW-PREF
    set local-preference 80
  !
  route-map infilter permit 30
  !
  ip prefix-list HIGH-PREF permit 10.0.0.0/8
  ip prefix-list LOW-PREF permit 20.0.0.0/8
```

Policy Control

Route Maps

Cisco.com

- Example using filter lists

```
router bgp 100
  neighbor 220.200.1.2 route-map filter-on-as-path in
  !
route-map filter-on-as-path permit 10
  match as-path 1
  set local-preference 80
  !
route-map filter-on-as-path permit 20
  match as-path 2
  set local-preference 200
  !
route-map filter-on-as-path permit 30
  !
ip as-path access-list 1 permit _150$
ip as-path access-list 2 permit _210_
```

Policy Control

Route Maps

Cisco.com

- **Example configuration of AS-PATH prepend**

```
router bgp 300
  network 215.7.0.0
  neighbor 2.2.2.2 remote-as 100
  neighbor 2.2.2.2 route-map SETPATH out
!
route-map SETPATH permit 10
  set as-path prepend 300 300
```

- **Use your own AS number when prepending**

Otherwise BGP loop detection may cause disconnects

Policy Control

Setting Communities

Cisco.com

- **Example Configuration**

```
router bgp 100
  neighbor 220.200.1.1 remote-as 200
  neighbor 220.200.1.1 send-community
  neighbor 220.200.1.1 route-map set-community out
!
route-map set-community permit 10
  match ip address prefix-list NO-ANNOUNCE
  set community no-export
!
route-map set-community permit 20
!
ip prefix-list NO-ANNOUNCE permit 172.168.0.0/16 ge 17
```

Policy Control

Matching Communities

Cisco.com

- Example Configuration

```
router bgp 100
  neighbor 220.200.1.2 remote-as 200
  neighbor 220.200.1.2 route-map filter-on-community in
!
route-map filter-on-community permit 10
  match community 1
  set local-preference 50
!
route-map filter-on-community permit 20
  match community 2 exact-match
  set local-preference 200
!
ip community-list 1 permit 150:3 200:5
ip community-list 2 permit 88:6
```

BGP for Internet Service Providers

Cisco.com

- **BGP Basics (quick recap)**
- **Scaling BGP**
- **Deploying BGP in an ISP network**
- **Multihoming Examples**

BGP Scaling Techniques

BGP Scaling Techniques

Cisco.com

- **How to scale iBGP mesh beyond a few peers?**
- **How to implement new policy without causing flaps and route churning?**
- **How to reduce the overhead on the routers?**
- **How to keep the network stable, scalable, as well as simple?**

BGP Scaling Techniques

Cisco.com

- **Dynamic Reconfiguration**
- **Peer groups**
- **Route flap damping**

Dynamic Reconfiguration

Soft Reconfiguration and Route Refresh

Soft Reconfiguration

Cisco.com

Problem:

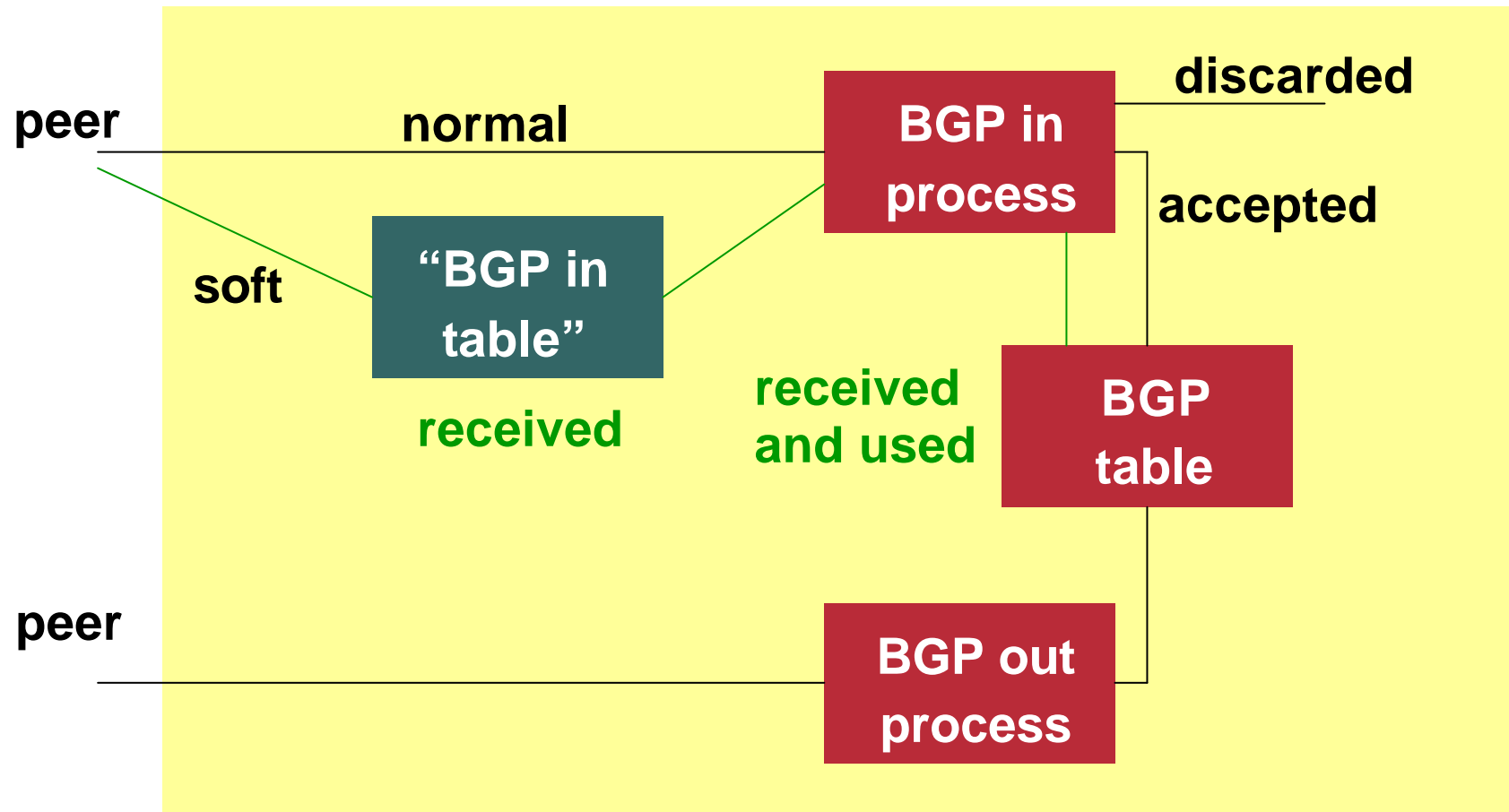
- **Hard BGP peer clear required after every policy change because the router does not store prefixes that are denied by a filter**
- **Hard BGP peer clearing consumes CPU and affects connectivity for all networks**

Solution:

- **Soft-reconfiguration**

Soft Reconfiguration

Cisco.com



Soft Reconfiguration

Cisco.com

- **New policy is activated without tearing down and restarting the peering session**
- **Per-neighbour basis**
- **Use more memory to keep prefixes whose attributes have been changed or have not been accepted**

Configuring Soft Reconfiguration

Cisco.com

```
router bgp 100
  neighbor 1.1.1.1 remote-as 101
  neighbor 1.1.1.1 route-map infilter in
  neighbor 1.1.1.1 soft-reconfiguration inbound
```

! Outbound does not need to be configured !

Then when we change the policy, we issue an exec command

```
clear ip bgp 1.1.1.1 soft [in | out]
```

Route Refresh Capability

Cisco.com

- Facilitates non-disruptive policy changes
- No configuration is needed
- No additional memory is used
- Requires peering routers to support “route refresh capability” – RFC2918
- **clear ip bgp x.x.x.x in** tells peer to resend full BGP announcement

Soft Reconfiguration vs. Route Refresh

Cisco.com

- **Use Route Refresh capability if supported**
find out from “show ip bgp neighbor”
uses much less memory
- **Otherwise use Soft Reconfiguration**
- **Only hard-reset a BGP peering as a last resort**

Peer Groups

Peer Groups

Cisco.com

Without peer groups

- **iBGP neighbours receive same update**
- **Large iBGP mesh slow to build**
- **Router CPU wasted on repeat calculations**

Solution – peer groups!

- **Group peers with same outbound policy**
- **Updates are generated once per group**

Peer Groups – Advantages

Cisco.com

- **Makes configuration easier**
- **Makes configuration less prone to error**
- **Makes configuration more readable**
- **Lower router CPU load**
- **iBGP mesh builds more quickly**
- **Members can have different inbound policy**
- **Can be used for eBGP neighbours too!**

Configuring Peer Group

Cisco.com

```
router bgp 100
  neighbor ibgp-peer peer-group
  neighbor ibgp-peer remote-as 100
  neighbor ibgp-peer update-source loopback 0
  neighbor ibgp-peer send-community
  neighbor ibgp-peer route-map outfilter out
  neighbor 1.1.1.1 peer-group ibgp-peer
  neighbor 2.2.2.2 peer-group ibgp-peer
  neighbor 2.2.2.2 route-map infilter in
  neighbor 3.3.3.3 peer-group ibgp-peer
```

! note how 2.2.2.2 has different inbound filter from peer-group !

Configuring Peer Group

Cisco.com

```
router bgp 109
  neighbor external-peer peer-group
  neighbor external-peer send-community
  neighbor external-peer route-map set-metric out
  neighbor 160.89.1.2 remote-as 200
  neighbor 160.89.1.2 peer-group external-peer
  neighbor 160.89.1.4 remote-as 300
  neighbor 160.89.1.4 peer-group external-peer
  neighbor 160.89.1.6 remote-as 400
  neighbor 160.89.1.6 peer-group external-peer
  neighbor 160.89.1.6 filter-list infilter in
```

Route Flap Damping

Stabilising the Network

Route Flap Damping

Cisco.com

- **Route flap**

Going up and down of path or change in attribute

BGP WITHDRAW followed by UPDATE = 1 flap

eBGP neighbour going down/up is NOT a flap

Ripples through the entire Internet

Wastes CPU

- **Damping aims to reduce scope of route flap propagation**

Route Flap Damping (continued)

Cisco.com

- **Requirements**

Fast convergence for normal route changes

History predicts future behaviour

Suppress oscillating routes

Advertise stable routes

- **Documented in RFC2439**

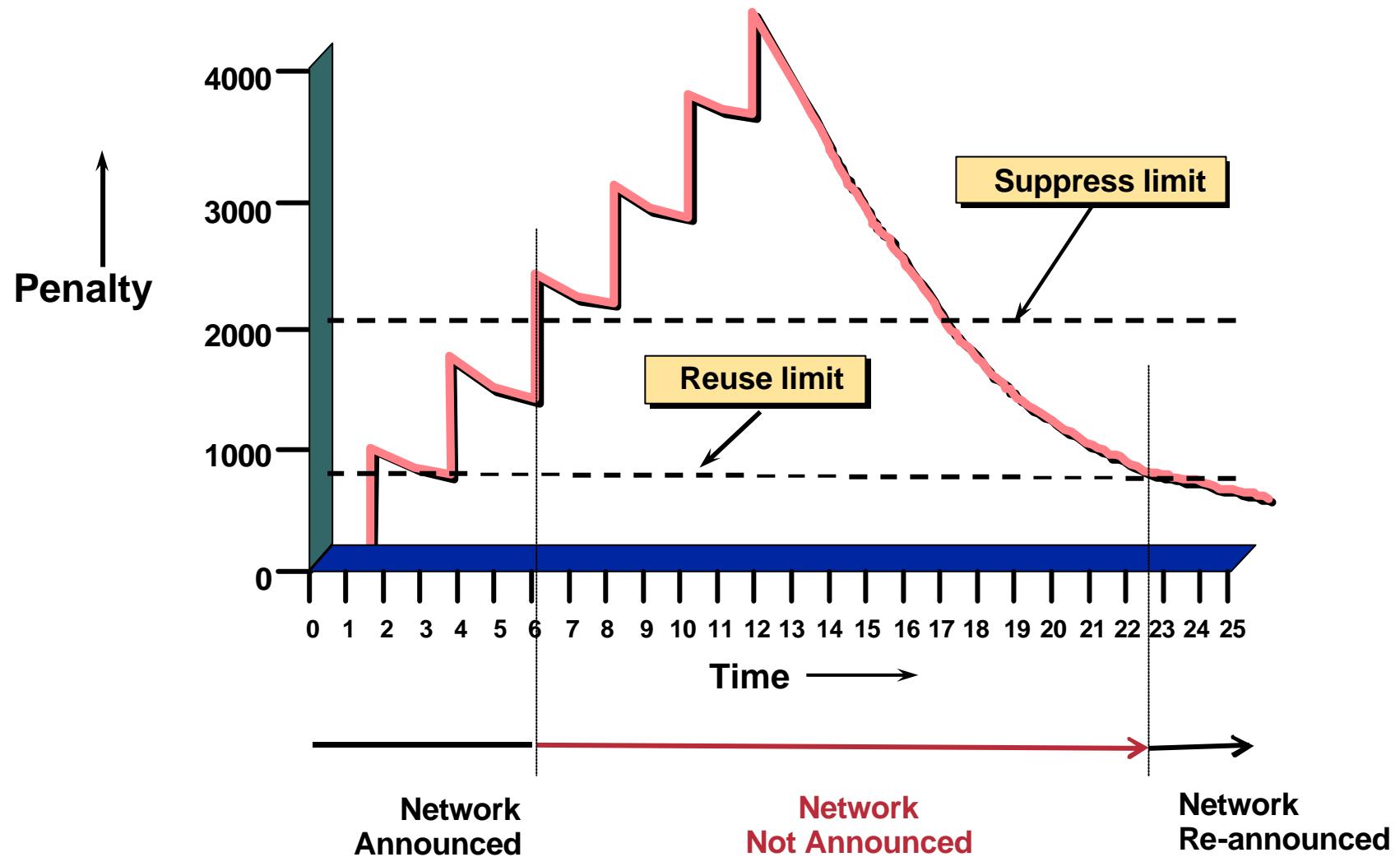
Operation

Cisco.com

- **Add penalty (1000) for each flap**
Change in attribute gets penalty of 500
- **Exponentially decay penalty**
half life determines decay rate
- **Penalty above suppress-limit**
do not advertise route to BGP peers
- **Penalty decayed below reuse-limit**
re-advertise route to BGP peers
penalty reset to zero when it is half of reuse-limit

Operation

Cisco.com



Operation

Cisco.com

- **Only applied to inbound announcements from eBGP peers**
- **Alternate paths still usable**
- **Controlled by:**
 - Half-life (default 15 minutes)**
 - reuse-limit (default 750)**
 - suppress-limit (default 2000)**
 - maximum suppress time (default 60 minutes)**

Configuration

Cisco.com

Fixed damping

```
router bgp 100
```

```
  bgp dampening [<half-life> <reuse-value> <suppress-  
    penalty> <maximum suppress time>]
```

Selective and variable damping

```
  bgp dampening [route-map <name>]
```

Variable damping

recommendations for ISPs

<http://www.ripe.net/docs/ripe-229.html>

BGP Scaling Techniques

Cisco.com

- **These 3 techniques should be core requirements in all ISP networks**

Soft reconfiguration/Route Refresh

Peer groups

Route flap damping

BGP for Internet Service Providers

Cisco.com

- **BGP Basics (quick recap)**
- **Scaling BGP**
- **Deploying BGP in an ISP network**
- **Multihoming Examples**

Deploying BGP in an ISP Network

Current Practices

BGP versus OSPF/ISIS

Cisco.com

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

BGP versus OSPF/ISIS

Cisco.com

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

BGP versus OSPF/ISIS

Configuration Example

Cisco.com

```
router bgp 34567
  neighbor core-ibgp peer-group
  neighbor core-ibgp remote-as 34567
  neighbor core-ibgp update-source Loopback0
  neighbor core-ibgp send-community
  neighbor core-ibgp-partial peer-group
  neighbor core-ibgp-partial remote-as 34567
  neighbor core-ibgp-partial update-source Loopback0
  neighbor core-ibgp-partial send-community
  neighbor core-ibgp-partial prefix-list network-ibgp out
  neighbor 222.1.9.10 peer-group core-ibgp
  neighbor 222.1.9.13 peer-group core-ibgp-partial
  neighbor 222.1.9.14 peer-group core-ibgp-partial
```

BGP versus OSPF/ISIS

Cisco.com

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

Aggregation

Quality or Quantity?

Aggregation

Cisco.com

- ISPs receive address block from Regional Registry or upstream provider
- **Aggregation** means announcing the **address block** only, not subprefixes
 - Subprefixes should only be announced in special cases – see later.
- **Aggregate should be generated internally**
 - Not on the network borders!**

Configuring Aggregation Method One

- **ISP has 221.10.0.0/19 address block**
- **To put into BGP as an aggregate:**

```
router bgp 100
```

```
network 221.10.0.0 mask 255.255.224.0
```

```
ip route 221.10.0.0 255.255.224.0 null0
```

- **The static route is a “pull up” route**
more specific prefixes within this address block ensure connectivity to ISP’s customers
“longest match lookup”

Configuring Aggregation Method Two

- **Configuration Example**

```
router bgp 109
  network 221.10.0.0 mask 255.255.252.0
  aggregate-address 221.10.0.0 255.255.224.0 [summary-only]
```

- **Requires more specific prefix in routing table before aggregate is announced**

- **{summary-only} keyword**

ensures that only the summary is announced if a more specific prefix exists in the routing table

- **Sets “aggregator” attribute**

Useful for debugging

Announcing Aggregate – Cisco IOS

Cisco.com

- **Configuration Example**

```
router bgp 100
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 101
  neighbor 222.222.10.1 prefix-list out-filter out
!
ip route 221.10.0.0 255.255.224.0 null0
!
ip prefix-list out-filter permit 221.10.0.0/19
```

Announcing an Aggregate

Cisco.com

- **ISPs who don't and won't aggregate are held in poor regard by community**
- **Registries' minimum allocation size is now a /20**

no real reason to see subprefixes of allocated blocks in the Internet

BUT there are currently >62000 /24s!

The Internet Today

Cisco.com

- **Current Internet Routing Table Statistics**

BGP Routing Table Entries 111947

Prefixes after maximum aggregation 73017

Unique prefixes in Internet 53184

Prefixes larger than registry alloc 45107

/24s announced 62487

only 5471 /24s are from 192.0.0.0/8

ASes in use 13045

Receiving Prefixes

Receiving Prefixes from downstream peers

Cisco.com

- **ISPs should only accept prefixes which have been assigned or allocated to their downstream peer**
- **For example**
 - downstream has 220.50.0.0/20 block**
 - should only announce this to peers**
 - peers should only accept this from them**

Receiving Prefixes: Cisco IOS

Cisco.com

- **Configuration Example on upstream**

```
router bgp 100
```

```
neighbor 222.222.10.1 remote-as 101
```

```
neighbor 222.222.10.1 prefix-list customer in
```

```
!
```

```
ip prefix-list customer permit 220.50.0.0/20
```

Receiving Prefixes from upstream peers

Cisco.com

- **Not desirable unless really necessary
special circumstances – see later**
- **Ask upstream to either:
originate a default-route
-or-
announce one prefix you can use as default**

Receiving Prefixes from upstream peers

Cisco.com

- **Downstream Router Configuration**

```
router bgp 100
  network 221.10.0.0 mask 255.255.224.0
  neighbor 221.5.7.1 remote-as 101
  neighbor 221.5.7.1 prefix-list infilter in
  neighbor 221.5.7.1 prefix-list outfilter out
!
ip prefix-list infilter permit 0.0.0.0/0
!
ip prefix-list outfilter permit 221.10.0.0/19
```

Receiving Prefixes from upstream peers

Cisco.com

- **Upstream Router Configuration**

```
router bgp 101
  neighbor 221.5.7.2 remote-as 100
  neighbor 221.5.7.2 default-originate
  neighbor 221.5.7.2 prefix-list cust-in in
  neighbor 221.5.7.2 prefix-list cust-out out
!
ip prefix-list cust-in permit 221.10.0.0/19
!
ip prefix-list cust-out permit 0.0.0.0/0
```

Receiving Prefixes from upstream peers

Cisco.com

- **If necessary to receive prefixes from upstream provider, care is required**

don't accept RFC1918 etc prefixes

<http://www.ietf.org/internet-drafts/draft-manning-dsua-07.txt>

don't accept your own prefix

don't accept default (unless you need it)

don't accept prefixes longer than /24

This guideline may change “soon”

Receiving Prefixes

```
router bgp 100
  network 221.10.0.0 mask 255.255.224.0
  neighbor 221.5.7.1 remote-as 101
  neighbor 221.5.7.1 prefix-list in-filter in
!
ip prefix-list in-filter deny 0.0.0.0/0                ! Block default
ip prefix-list in-filter deny 0.0.0.0/8 le 32
ip prefix-list in-filter deny 10.0.0.0/8 le 32
ip prefix-list in-filter deny 127.0.0.0/8 le 32
ip prefix-list in-filter deny 169.254.0.0/16 le 32
ip prefix-list in-filter deny 172.16.0.0/12 le 32
ip prefix-list in-filter deny 192.0.2.0/24 le 32
ip prefix-list in-filter deny 192.168.0.0/16 le 32
ip prefix-list in-filter deny 221.10.0.0/19 le 32 ! Block local prefix
ip prefix-list in-filter deny 224.0.0.0/3 le 32    ! Block multicast
ip prefix-list in-filter deny 0.0.0.0/0 ge 25      ! Block prefixes >/24
ip prefix-list in-filter permit 0.0.0.0/0 le 32
```

Prefixes into iBGP

Injecting prefixes into iBGP

Cisco.com

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

Router Configuration network statement

Cisco.com

- **Example:**

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

Injecting prefixes into iBGP

Cisco.com

- **interface flap will result in prefix withdraw and re-announce**

use “ip route...permanent”

**Static route always exists, even if interface is down
® prefix announced in iBGP**

- **many ISPs use redistribute static rather than network statement**

only use this if you understand why

Inserting prefixes into BGP: redistribute static

Cisco.com

- Care required with **redistribute!**

redistribute <routing-protocol> means everything in the <routing-protocol> will be transferred into the current routing protocol

Does not scale if uncontrolled

Best avoided if at all possible

redistribute normally used with “route-maps” and under tight administrative control

Router Configuration: redistribute static

Cisco.com

- **Example:**

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

Injecting prefixes into iBGP

Cisco.com

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

Configuration Tips

iBGP and IGPs

Cisco.com

- **Make sure loopback is configured on router**
iBGP between loopbacks, **NOT** real interfaces
- **Make sure IGP carries loopback /32 address**
- **Make sure IGP carries DMZ nets**
Use ip-unnumbered where possible
Or use next-hop-self on iBGP neighbours
neighbor x.x.x.x next-hop-self

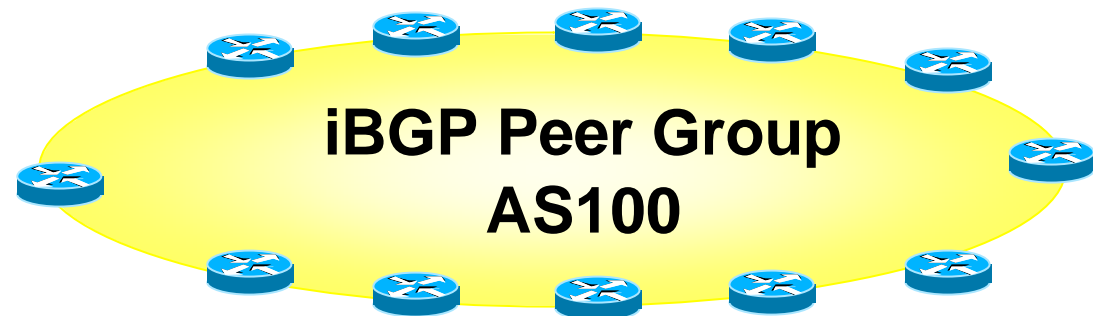
Next-hop-self

Cisco.com

- **Used by many ISPs on edge routers**
 - Preferable to carrying DMZ /30 addresses in the IGP**
 - Reduces size of IGP to just core infrastructure**
 - Alternative to using `ip unnumbered`**
 - Helps scale network**
 - BGP speaker announces external network using local address (loopback) as next-hop**

BGP Template – iBGP peers

Cisco.com



```
router bgp 100
neighbor internal peer-group
neighbor internal description ibgp peers
neighbor internal remote-as 100
neighbor internal update-source Loopback0
neighbor internal next-hop-self
neighbor internal send-community
neighbor internal version 4
neighbor internal password 7 03085A09
neighbor 1.0.0.1 peer-group internal
neighbor 1.0.0.2 peer-group internal
```

BGP Template – iBGP peers

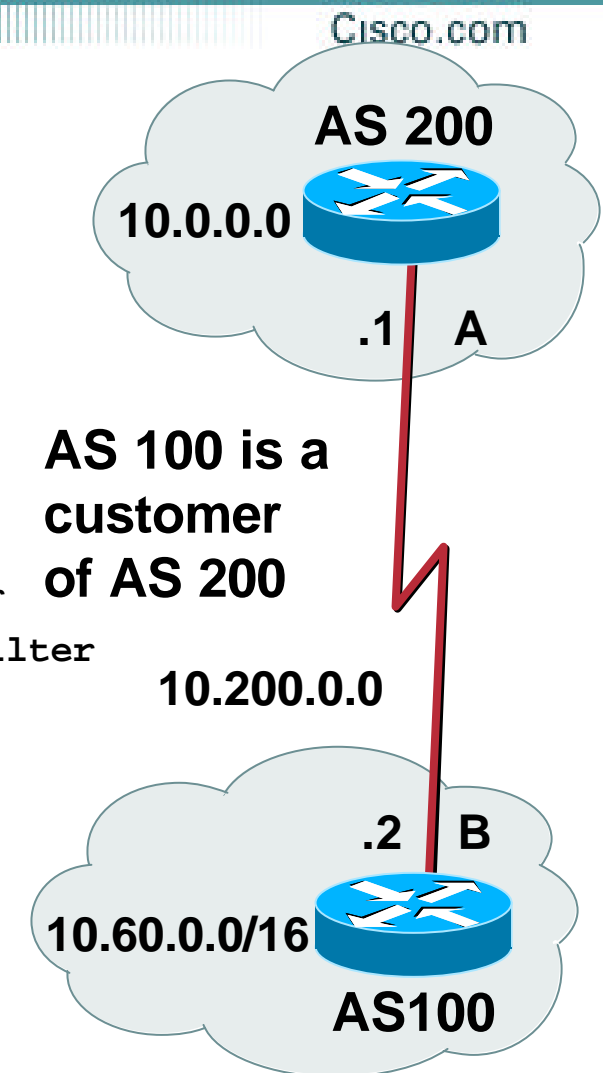
Cisco.com

- **Use peer-groups**
- **iBGP between loopbacks!**
- **Next-hop-self**
Keep DMZ and point-to-point out of IGP
- **Always send communities in iBGP**
Otherwise accidents will happen
- **Hardwire BGP to version 4**
Yes, this is being paranoid!
- **Use passwords on iBGP session**
Not being paranoid, **VERY** necessary

BGP Template – eBGP peers

Router B:

```
router bgp 100
bgp dampening route-map RIPE229-flap
network 10.60.0.0 mask 255.255.0.0
neighbor external peer-group
neighbor external remote-as 200
neighbor external description ISP connection
neighbor external remove-private-AS
neighbor external version 4
neighbor external prefix-list ispout out ! "real" filter
neighbor external filter-list 1 out      ! "accident" filter
neighbor external route-map ispout out
neighbor external prefix-list ispin in
neighbor external filter-list 2 in
neighbor external route-map ispin in
neighbor external password 7 020A0559
neighbor external maximum-prefix 120000 [warning-only]
neighbor 10.200.0.1 peer-group external
!
ip route 10.60.0.0 255.255.0.0 null0 254
```



BGP Template – eBGP peers

Cisco.com

- **BGP damping – use RIPE-229 parameters**
- **Remove private ASes from announcements**
Common omission today
- **Use extensive filters, with “backup”**
Use as-path filters to backup prefix-lists
Use route-maps for policy
- **Use password agreed between you and peer on eBGP session**
- **Use maximum-prefix tracking**
Router will warn you if there are sudden changes in BGP table size, bringing down eBGP if necessary

More BGP “defaults”

Cisco.com

- Log neighbour changes

bgp log-neighbor-changes

- Enable deterministic MED

bgp deterministic-med

Otherwise bestpath could be different every time BGP session is reset

- Make BGP admin distance higher than any IGP

distance bgp 200 200 200

BGP for Internet Service Providers

Cisco.com

- **BGP Basics (quick recap)**
- **Scaling BGP**
- **Deploying BGP in an ISP network**
- **Multihoming Examples**

Multihoming

Multihoming Definition

Cisco.com

- **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

AS Numbers

Cisco.com

- **An Autonomous System Number is required by BGP**
- **Obtained from upstream ISP or Regional Registry**
- **Necessary when you have links to more than one ISP or exchange point**

Configuring Policy

Cisco.com

- **Three BASIC Principles**
 - prefix-lists** to filter **prefixes**
 - filter-lists** to filter **ASNs**
 - route-maps** to apply **policy**
- **Avoids confusion!**

Originating Prefixes

Cisco.com

- **Basic Assumptions**

MUST announce assigned address block to Internet

MAY also announce subprefixes – reachability is not guaranteed

RIR minimum allocation is /20

**several ISPs filter RIR blocks on this boundary
called “Net Police” by some**

Part

```
!! APNIC
ip prefix-list FILTER permit 61.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 202.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 210.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 218.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 220.0.0.0/8 ge 9 le 20
!! ARIN
ip prefix-list FILTER permit 24.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 63.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 64.0.0.0/6 ge 9 le 20
ip prefix-list FILTER permit 68.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 199.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 200.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 204.0.0.0/6 ge 9 le 20
ip prefix-list FILTER permit 208.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 216.0.0.0/8 ge 9 le 20
!! RIPE NCC
ip prefix-list FILTER permit 62.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 80.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 193.0.0.0/8 ge 9 le 20
ip prefix-list FILTER permit 194.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 212.0.0.0/7 ge 9 le 20
ip prefix-list FILTER permit 217.0.0.0/8 ge 9 le 20
```

“Net Police” prefix list issues

Cisco.com

- meant to “punish” ISPs who won’t and don’t aggregate
- impacts legitimate multihoming
- 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 it current**

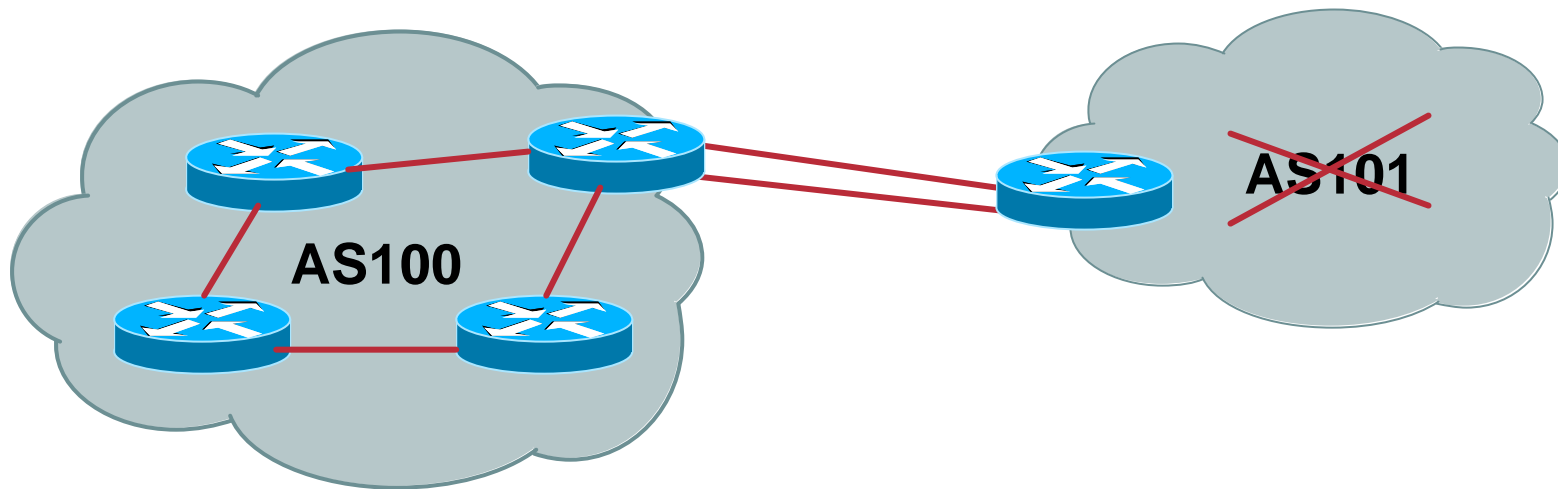
Multihoming Options

Multihoming Scenarios

Cisco.com

- **Stub network**
- **Multi-homed stub network**
- **Multi-homed network**
- **Configuration Options**

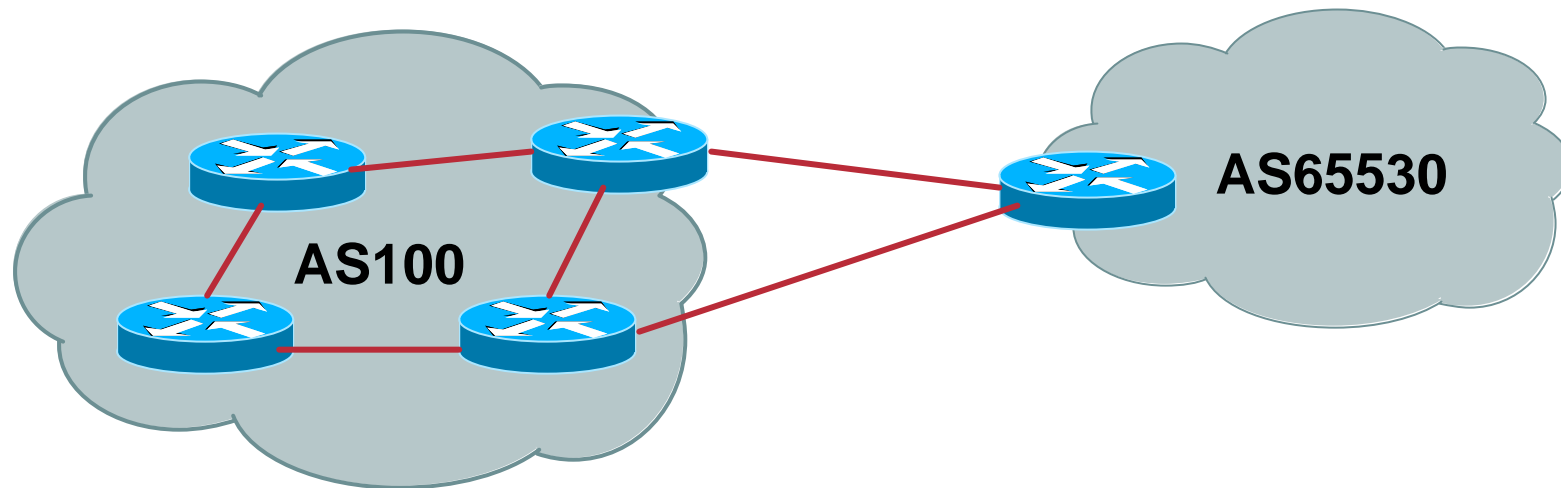
Stub Network



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

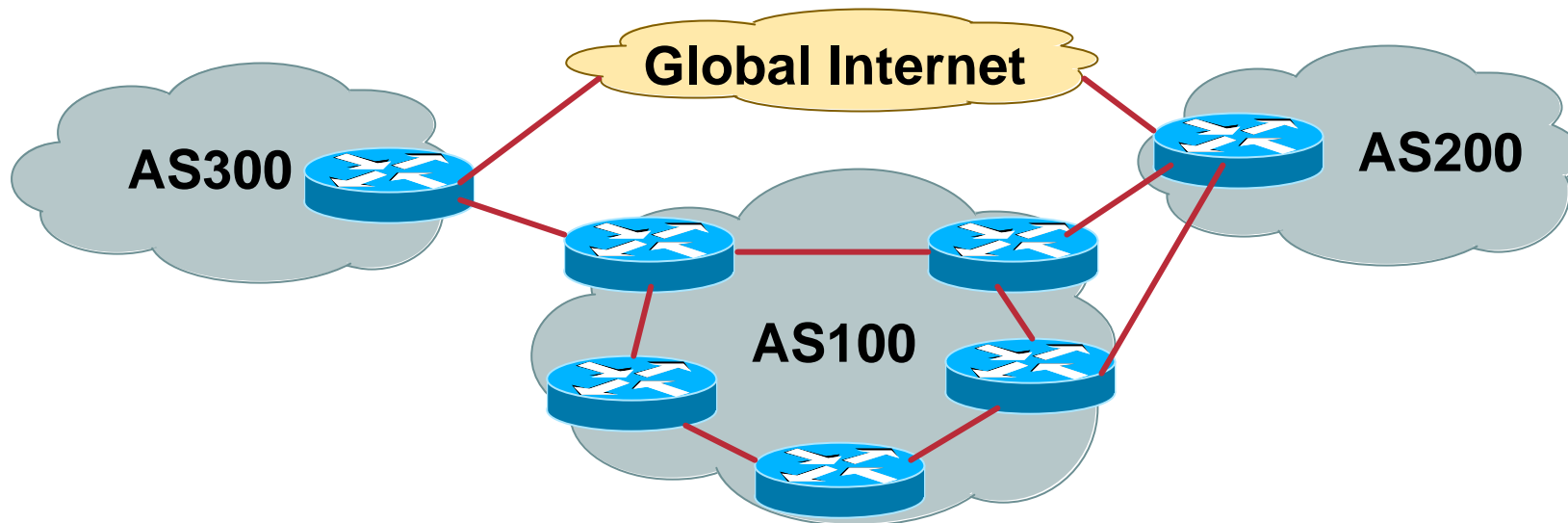
Multi-homed Stub Network

Cisco.com



- 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

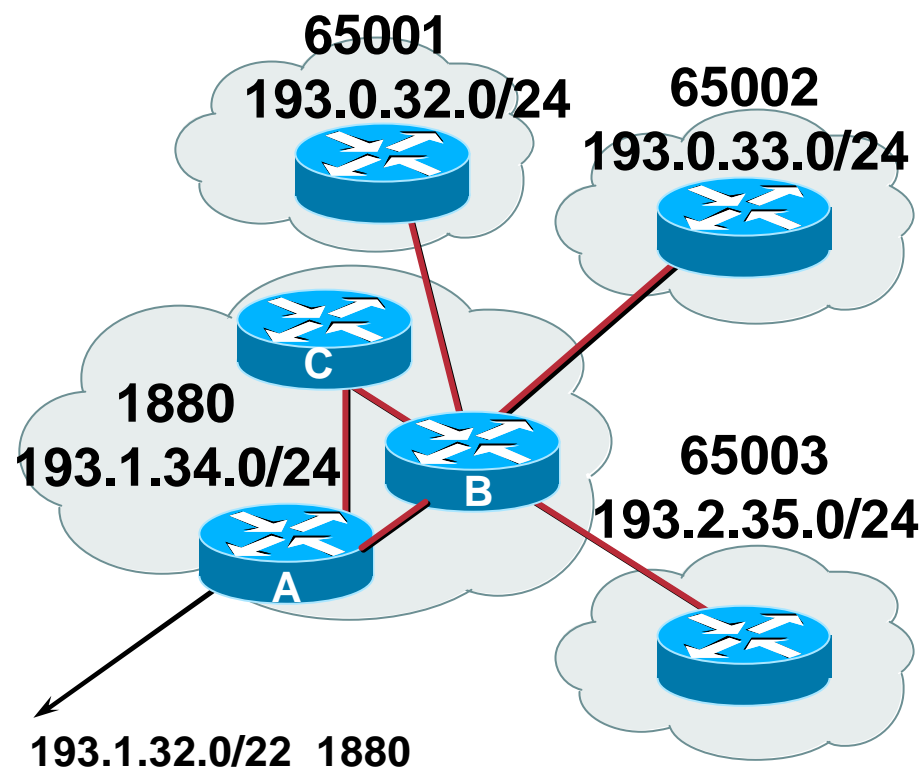
Private-AS – Application

Cisco.com

- **Applications**

ISP with single-homed customers (RFC2270)

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



Private-AS Removal

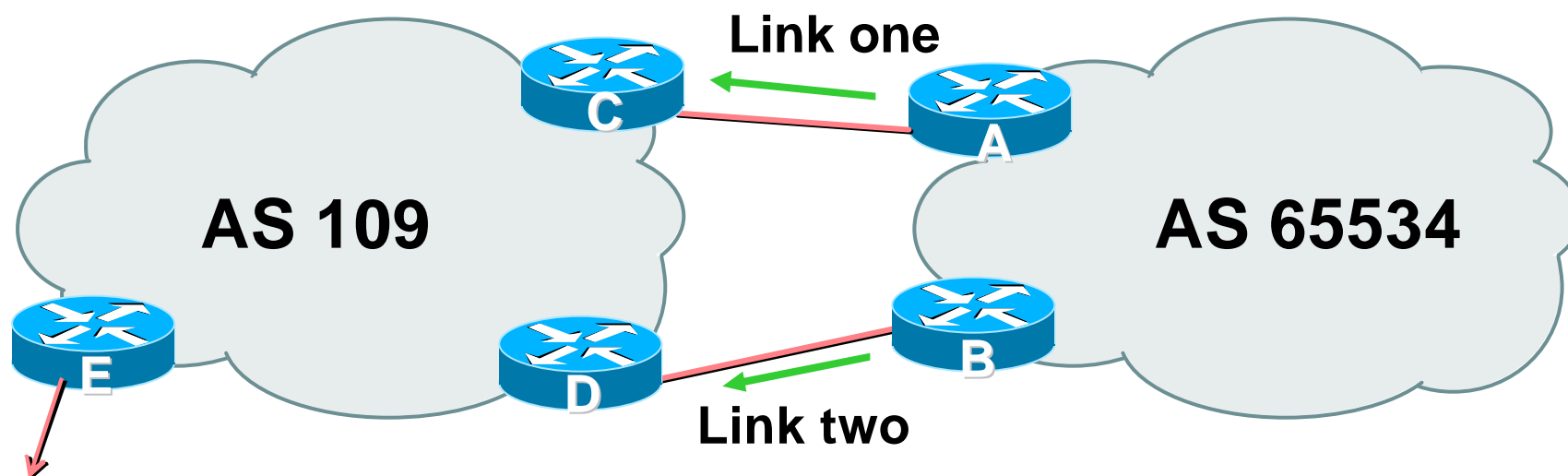
- **neighbor x.x.x.x remove-private-AS**
- **Rules:**
 - available for eBGP neighbors only**
 - if the update has AS_PATH made up of private-AS numbers, the private-AS will be dropped**
 - if the AS_PATH includes private and public AS numbers, private AS number will not be removed...it is a configuration error!**
 - if AS_PATH contains the AS number of the eBGP neighbor, the private-AS numbers will not be removed**
 - if used with confederations, it will work as long as the private AS numbers are after the confederation portion of the AS_PATH**

Two links to the same ISP

With Redundancy and Loadsharing

Two links to the same ISP (with redundancy)

Cisco.com



- **AS109 removes private AS and any customer subprefixes from Internet announcement**

Loadsharing to the same ISP

Cisco.com

- **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 approx as most ISP and end-site traffic is inbound

Two links to the same ISP

Cisco.com

- **Router A Configuration**

```
router bgp 65534
  network 221.10.0.0 mask 255.255.224.0
  network 221.10.0.0 mask 255.255.240.0
  neighbor 222.222.10.2 remote-as 109
  neighbor 222.222.10.2 prefix-list routerC out
  neighbor 222.222.10.2 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list routerC permit 221.10.0.0/20
ip prefix-list routerC permit 221.10.0.0/19
!
ip route 221.10.0.0 255.255.240.0 null0
ip route 221.10.0.0 255.255.224.0 null0
```

Router B configuration is similar but with the other /20

Two links to the same ISP

- **Router C Configuration**

```
router bgp 109
```

```
neighbor 222.222.10.1 remote-as 65534
```

```
neighbor 222.222.10.1 default-originate
```

```
neighbor 222.222.10.1 prefix-list Customer in
```

```
neighbor 222.222.10.1 prefix-list default out
```

```
!
```

```
ip prefix-list Customer permit 221.10.0.0/19 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

Cisco.com

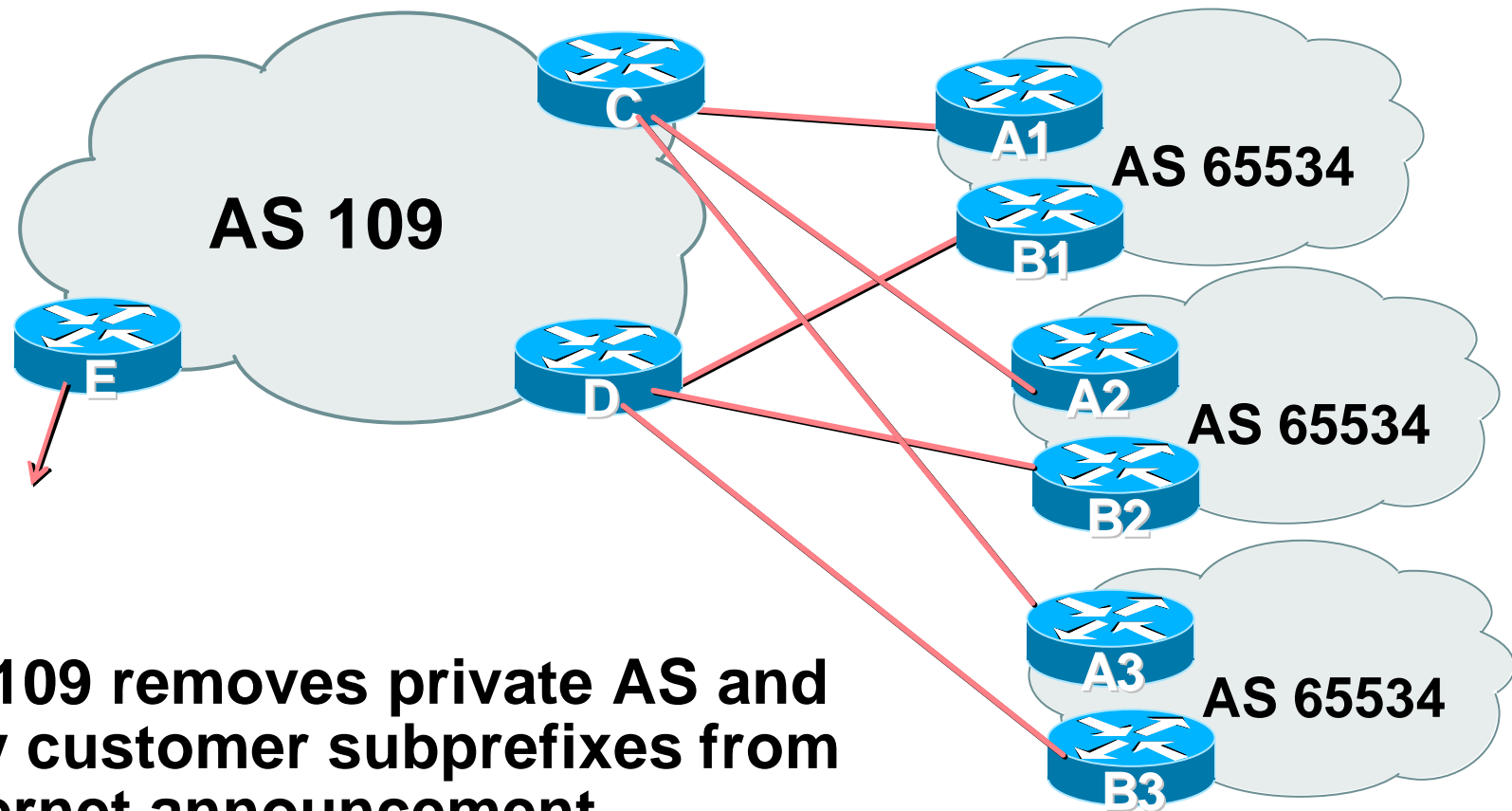
- **Loadsharing configuration is only on customer router**
- **Upstream ISP has to**
 - remove customer subprefixes from external announcements**
 - remove private AS from external announcements**
- **Could also use BGP communities**

Two links to the same ISP

**Multiple Dualhomed Customers
(RFC2270)**

Multiple Dualhomed Customers (RFC2270)

Cisco.com



- **AS109 removes private AS and any customer subprefixes from Internet announcement**

Multiple Dualhomed Customers

Cisco.com

- **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**

Two links to the same ISP

Cisco.com

- **Router A1 Configuration**

```
router bgp 65534
  network 221.10.0.0 mask 255.255.224.0
  network 221.10.0.0 mask 255.255.240.0
  neighbor 222.222.10.2 remote-as 109
  neighbor 222.222.10.2 prefix-list routerC out
  neighbor 222.222.10.2 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
ip prefix-list routerC permit 221.10.0.0/20
ip prefix-list routerC permit 221.10.0.0/19
!
ip route 221.10.0.0 255.255.240.0 null0
ip route 221.10.0.0 255.255.224.0 null0
```

Router B1 configuration is similar but for the other /20

Multiple Dualhomed Customers

Cisco.com

- Router C Configuration

```
router bgp 109
  neighbor bgp-customers peer-group
  neighbor bgp-customers remote-as 65534
  neighbor bgp-customers default-originate
  neighbor bgp-customers prefix-list default out
  neighbor 222.222.10.1 peer-group bgp-customers
  neighbor 222.222.10.1 description Customer One
  neighbor 222.222.10.1 prefix-list Customer1 in
  neighbor 222.222.10.9 peer-group bgp-customers
  neighbor 222.222.10.9 description Customer Two
  neighbor 222.222.10.9 prefix-list Customer2 in
```


Multiple Dualhomed Customers

Cisco.com

```
neighbor 222.222.10.17 peer-group bgp-customers
neighbor 222.222.10.17 description Customer Three
neighbor 222.222.10.17 prefix-list Customer3 in
!
ip prefix-list Customer1 permit 221.10.0.0/19 le 20
ip prefix-list Customer2 permit 221.16.64.0/19 le 20
ip prefix-list Customer3 permit 221.14.192.0/19 le 20
ip prefix-list default permit 0.0.0.0/0
```

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

Multiple Dualhomed Customers

Cisco.com

- **Router E Configuration**

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

```
router bgp 109
  neighbor 222.222.10.17 remote-as 110
  neighbor 222.222.10.17 remove-private-AS
  neighbor 222.222.10.17 prefix-list Customers out
!
ip prefix-list Customers permit 221.10.0.0/19
ip prefix-list Customers permit 221.16.64.0/19
ip prefix-list Customers permit 221.14.192.0/19
```

- **Private AS still visible inside AS109**

Multiple Dualhomed Customers

Cisco.com

- 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 109
  neighbor 222.222.10.17 remote-as 110
  neighbor 222.222.10.17 prefix-list my-aggregate out
!
ip prefix-list my-aggregate permit 221.8.0.0/13
```

Two links to different ISPs

With Redundancy

Two links to different ISPs (with redundancy)

Cisco.com

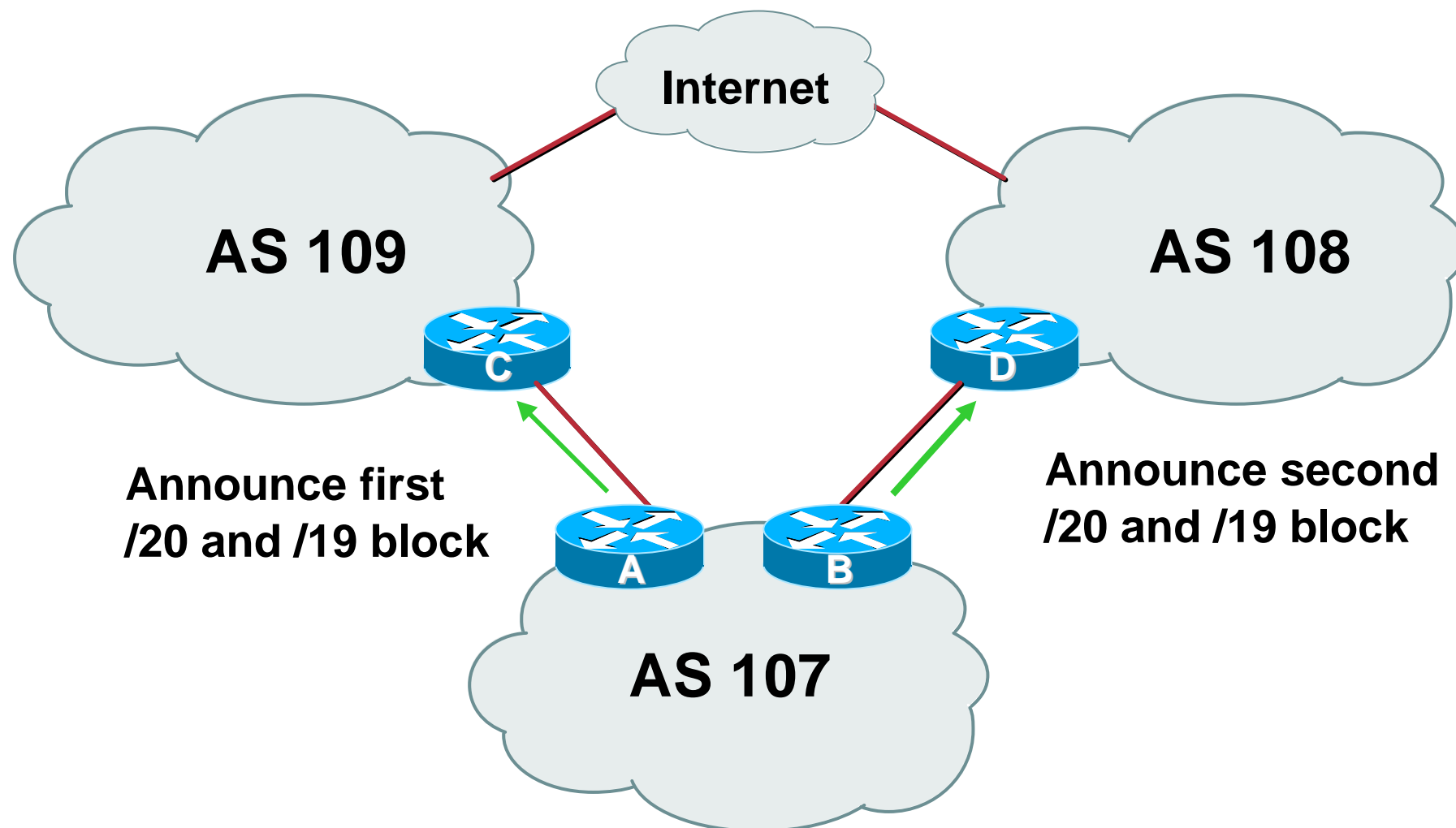
- **Announce /19 aggregate on each link**
- **Split /19 and announce as two /20s, one on each link**

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 redundancy)

Cisco.com



Two links to different ISPs (with redundancy)

Cisco.com

- Router A Configuration

```
router bgp 107
  network 221.10.0.0 mask 255.255.224.0
  network 221.10.0.0 mask 255.255.240.0
  neighbor 222.222.10.1 remote-as 109
  neighbor 222.222.10.1 prefix-list firstblock out
  neighbor 222.222.10.1 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
!
ip prefix-list firstblock permit 221.10.0.0/20
ip prefix-list firstblock permit 221.10.0.0/19
```

Two links to different ISPs (with redundancy)

Cisco.com

- Router B Configuration

```
router bgp 107
  network 221.10.0.0 mask 255.255.224.0
  network 221.10.16.0 mask 255.255.240.0
  neighbor 220.1.5.1 remote-as 108
  neighbor 220.1.5.1 prefix-list secondblock out
  neighbor 220.1.5.1 prefix-list default in
!
ip prefix-list default permit 0.0.0.0/0
!
ip prefix-list secondblock permit 221.10.16.0/20
ip prefix-list secondblock permit 221.10.0.0/19
```


Two links to different ISPs

More Controlled Loadsharing

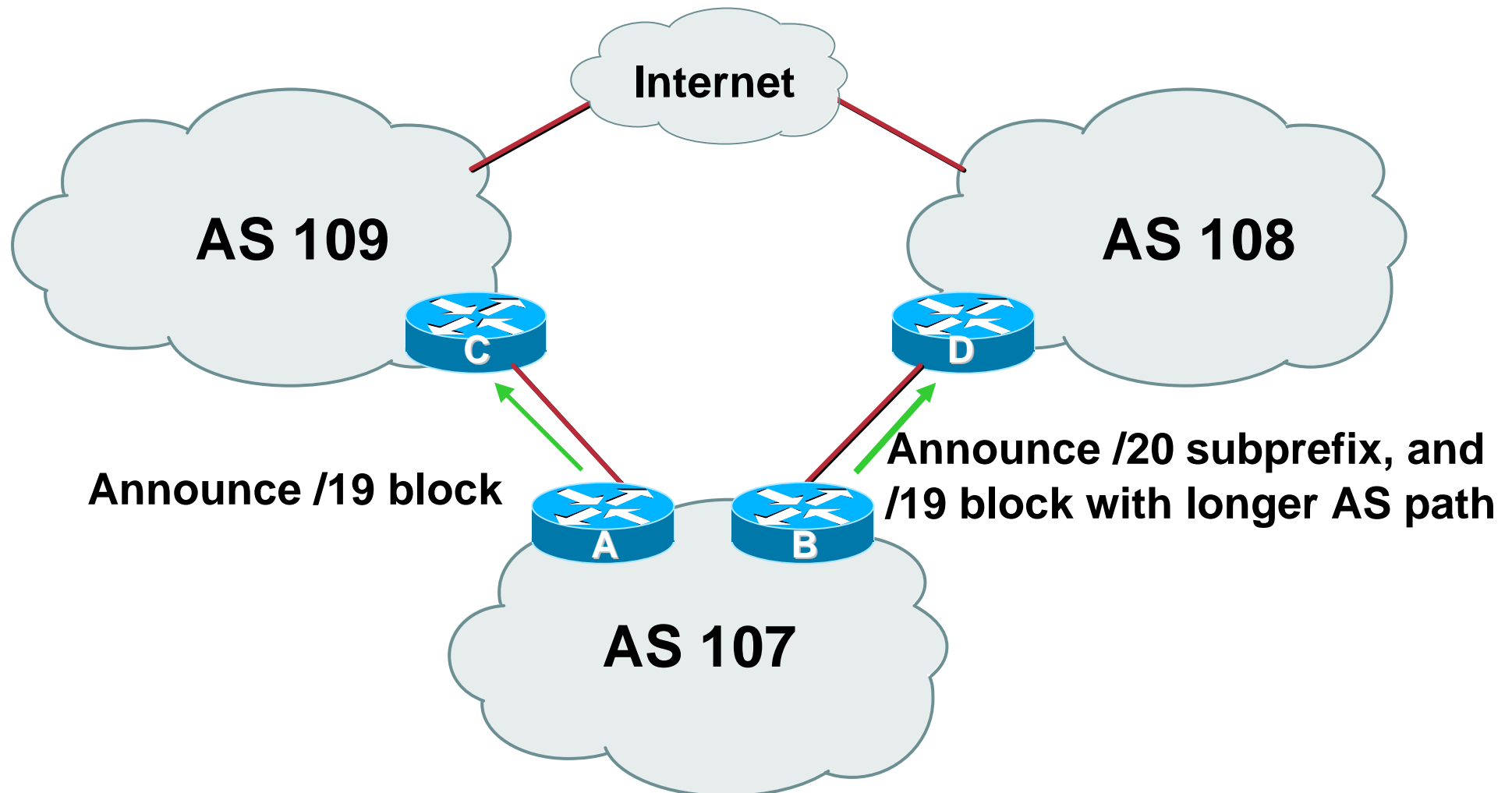
Loadsharing with different ISPs

Cisco.com

- **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

Cisco.com



Loadsharing with different ISPs

Cisco.com

- **Router A Configuration**

```
router bgp 107
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 109
  neighbor 222.222.10.1 prefix-list default in
  neighbor 222.222.10.1 prefix-list aggregate out
!
ip prefix-list aggregate permit 221.10.0.0/19
```

Loadsharing with different ISPs

Cisco.com

- Router B Configuration

```
router bgp 107
  network 221.10.0.0 mask 255.255.224.0
  network 221.10.16.0 mask 255.255.240.0
  neighbor 220.1.5.1 remote-as 108
  neighbor 220.1.5.1 prefix-list default in
  neighbor 220.1.5.1 prefix-list subblocks out
  neighbor 220.1.5.1 route-map routerD out
!
..next slide..
```

Loadsharing with different ISPs

Cisco.com

```
route-map routerD permit 10
  match ip address prefix-list aggregate
  set as-path prepend 107 107
route-map routerD permit 20
!
ip prefix-list subblocks permit 221.10.0.0/19 le 20
ip prefix-list aggregate permit 221.10.0.0/19
```

Service Provider Multihoming

Service Provider Multihoming

Cisco.com

- **Previous examples dealt with loadsharing inbound traffic**

What about outbound?

- **ISPs strive to balance traffic flows in both directions**

Balance link utilisation

Try and keep most traffic flows symmetric

Service Provider Multihoming

Cisco.com

- **Balancing outbound traffic requires inbound routing information**

Common solution is “full routing table”

Rarely necessary – the “routing mallet” to try solve loadsharing problems

Keep It Simple (KISS) is often easier (and \$\$\$ cheaper) than carrying n-copies of the full routing table

Service Provider Multihoming

Cisco.com

- **Examples**

One upstream, one local peer

One upstream, local exchange point

Two upstreams, one local peer

- **All examples require BGP and a public ASN**

Service Provider Multihoming

One Upstream, One local peer

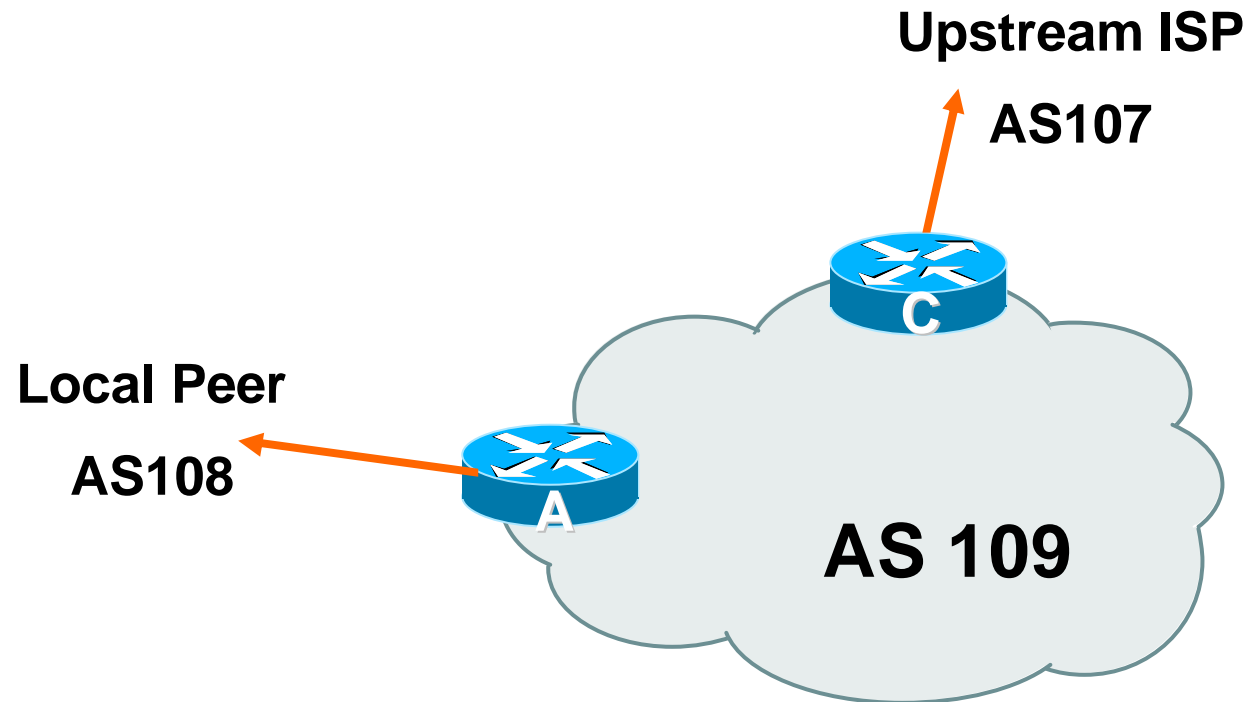
One Upstream, One Local Peer

Cisco.com

- **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

Cisco.com



One Upstream, One Local Peer

Cisco.com

- Router A Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.2 remote-as 108
  neighbor 222.222.10.2 prefix-list my-block out
  neighbor 222.222.10.2 prefix-list AS108-peer in
!
ip prefix-list AS108-peer permit 222.5.16.0/19
ip prefix-list AS108-peer permit 221.240.0.0/20
ip prefix-list my-block permit 221.10.0.0/19
!
ip route 221.10.0.0 255.255.224.0 null0
```

One Upstream, One Local Peer

Cisco.com

- **Router A – Alternative Configuration**

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.2 remote-as 108
  neighbor 222.222.10.2 prefix-list my-block out
  neighbor 222.222.10.2 filter-list 10 in
!
ip as-path access-list 10 permit ^(108_)+$
!
ip prefix-list my-block permit 221.10.0.0/19
!
ip route 221.10.0.0 255.255.224.0 null0
```

One Upstream, One Local Peer

Cisco.com

- Router C Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list default in
  neighbor 222.222.10.1 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 221.10.0.0 255.255.224.0 null0
```


One Upstream, One Local Peer

Cisco.com

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

Service Provider Multihoming

One Upstream, Local Exchange Point

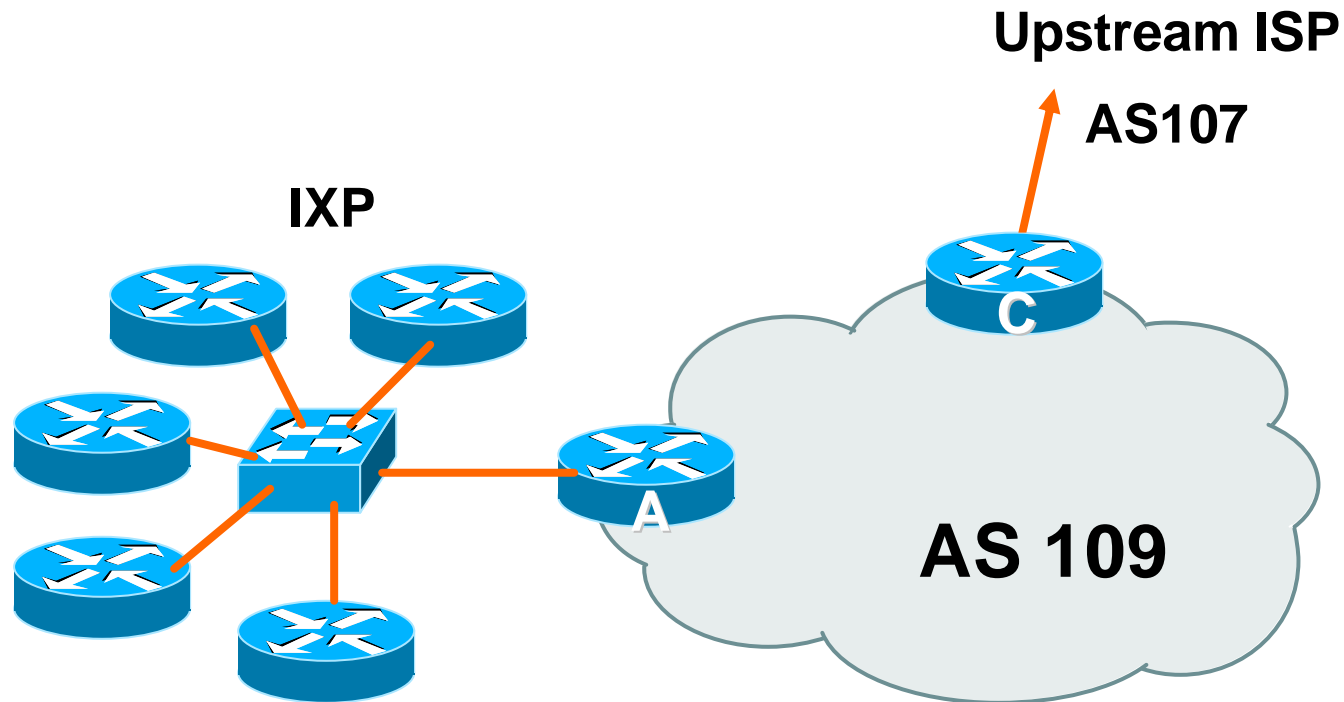
One Upstream, Local Exchange Point

Cisco.com

- **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 from IXP peers**

One Upstream, Local Exchange Point

Cisco.com



One Upstream, Local Exchange Point

Cisco.com

- **Router A Configuration**

```
interface fastethernet 0/0
  description Exchange Point LAN
  ip address 220.5.10.1 mask 255.255.255.224
  ip verify unicast reverse-path
  no ip directed-broadcast
  no ip proxy-arp
  no ip redirects
!
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor ixp-peers peer-group
  neighbor ixp-peers soft-reconfiguration in
  neighbor ixp-peers prefix-list my-block out
..next slide
```

One Upstream, Local Exchange Point

Cisco.com

```
neighbor 220.5.10.2 remote-as 100
neighbor 222.5.10.2 peer-group ixp-peers
neighbor 222.5.10.2 prefix-list peer100 in
neighbor 220.5.10.3 remote-as 101
neighbor 222.5.10.3 peer-group ixp-peers
neighbor 222.5.10.3 prefix-list peer101 in
neighbor 220.5.10.4 remote-as 102
neighbor 222.5.10.4 peer-group ixp-peers
neighbor 222.5.10.4 prefix-list peer102 in
neighbor 220.5.10.5 remote-as 103
neighbor 222.5.10.5 peer-group ixp-peers
neighbor 222.5.10.5 prefix-list peer103 in
..next slide
```

One Upstream, Local Exchange Point

Cisco.com

```
ip route 221.10.0.0 255.255.224.0 null0
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list peer100 permit 222.0.0.0/19
ip prefix-list peer101 permit 222.30.0.0/19
ip prefix-list peer102 permit 222.12.0.0/19
ip prefix-list peer103 permit 222.18.128.0/19
!
```

One Upstream, Local Exchange Point

Cisco.com

- Router C Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list default in
  neighbor 222.222.10.1 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 221.10.0.0 255.255.224.0 null0
```


One Upstream, Local Exchange Point

Cisco.com

- **Note Router A configuration**
Prefix-list higher maintenance, but safer
uRPF on the FastEthernet interface
- **IXP traffic goes to and from local IXP,**
everything else goes to upstream

Service Provider Multihoming

Two Upstreams, One local peer

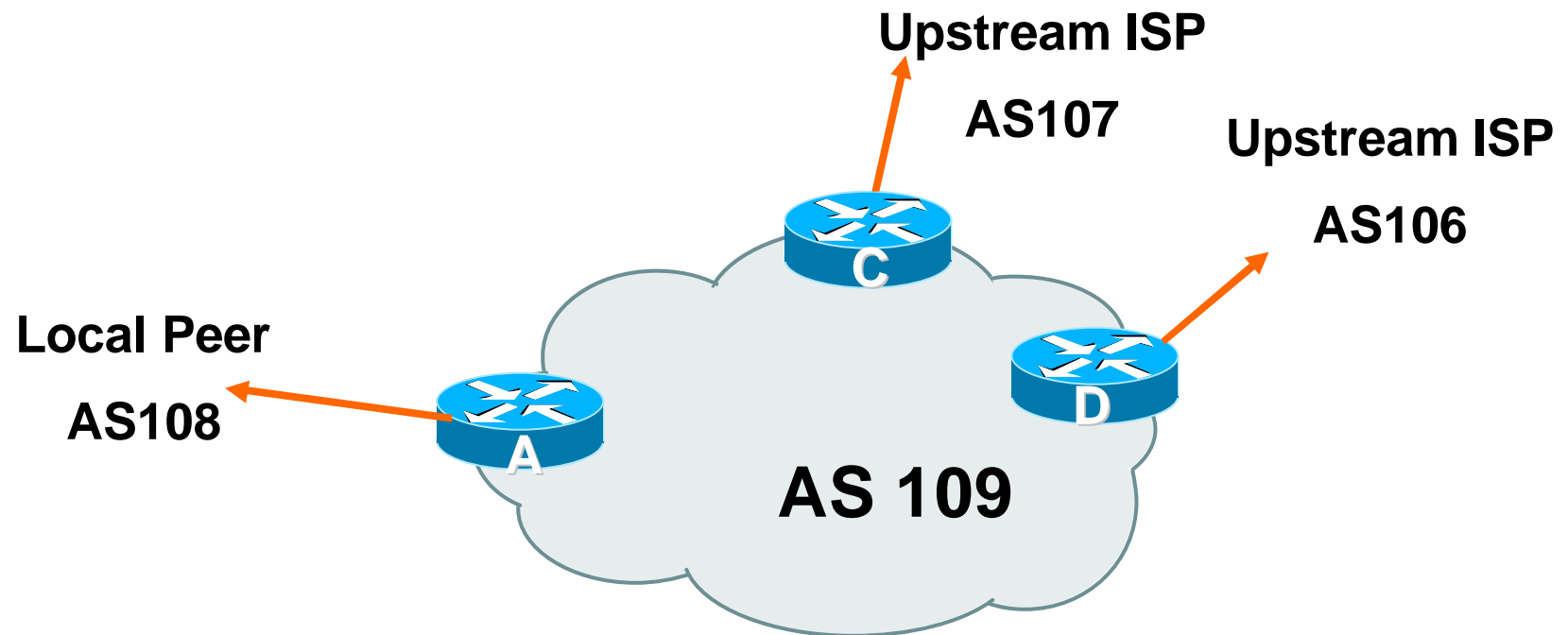
Two Upstreams, One Local Peer

Cisco.com

- **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**

Two Upstreams, One Local Peer

Cisco.com



Two Upstreams, One Local Peer

Cisco.com

- **Router A**

Same routing configuration as in example with one upstream and one local peer

Same hardware configuration

Two Upstreams, One Local Peer

Cisco.com

- Router C Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list default in
  neighbor 222.222.10.1 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 221.10.0.0 255.255.224.0 null0
```

Two Upstreams, One Local Peer

Cisco.com

- Router D Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.5 remote-as 106
  neighbor 222.222.10.5 prefix-list default in
  neighbor 222.222.10.5 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 221.10.0.0 255.255.224.0 null0
```

Two Upstreams, One Local Peer

Cisco.com

- **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

Cisco.com

- **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!

Two Upstreams, One Local Peer: Full Routes

Cisco.com

- Router C Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list rfc1918-deny in
  neighbor 222.222.10.1 prefix-list my-block out
  neighbor 222.222.10.1 route-map AS107-loadshare in
!
ip prefix-list my-block permit 221.10.0.0/19
! See earlier in presentation for RFC1918 list
..next slide
```

Two Upstreams, One Local Peer: Full Routes

Cisco.com

```
ip route 221.10.0.0 255.255.224.0 null0
!
ip as-path access-list 10 permit ^(107_)+$
ip as-path access-list 10 permit ^(107_)+_[0-9]+$
!
route-map AS107-loadshare permit 10
    match ip as-path 10
    set local-preference 120
route-map AS107-loadshare permit 20
    set local-preference 80
!
```

Two Upstreams, One Local Peer: Full Routes

Cisco.com

- **Router D Configuration**

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.5 remote-as 106
  neighbor 222.222.10.5 prefix-list rfc1918-deny in
  neighbor 222.222.10.5 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
! See earlier in presentation for RFC1918 list
```

Two Upstreams, One Local Peer: Full Routes

Cisco.com

- **Router C configuration:**
 - Accept full routes from AS107**
 - Tag prefixes originated by AS107 and AS107's neighbouring ASes with local preference 120**
 - Traffic to those ASes will go over AS107 link**
 - Remaining prefixes tagged with local preference of 80**
 - Traffic to other all other ASes will go over the link to AS106**
- **Router D configuration same as Router C without the route-map**

Two Upstreams, One Local Peer: Full Routes

Cisco.com

- **Full routes from upstreams**

Expensive – needs 128Mbytes RAM today

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 presentation for examples

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- Router C Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list rfc1918-nodef-deny in
  neighbor 222.222.10.1 prefix-list my-block out
  neighbor 222.222.10.1 filter-list 10 in
  neighbor 222.222.10.1 route-map tag-default-low in
  !
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
..next slide
```

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

```
! See earlier presentation for RFC1918 list
!
ip route 221.10.0.0 255.255.224.0 null0
!
ip as-path access-list 10 permit ^(107_)+$
ip as-path access-list 10 permit ^(107_)+_[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

Cisco.com

- Router D Configuration

```
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.5 remote-as 106
  neighbor 222.222.10.5 prefix-list default in
  neighbor 222.222.10.5 prefix-list my-block out
!
ip prefix-list my-block permit 221.10.0.0/19
ip prefix-list default permit 0.0.0.0/0
!
ip route 221.10.0.0 255.255.224.0 null0
```

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **Router C configuration:**

Accept full routes from AS107

(or get them to send less)

Filter ASNs so only AS107 and AS107's neighbouring ASes are accepted

Allow default, and set it to local preference 80

Traffic to those ASes will go over AS107 link

Traffic to other all other ASes will go over the link to AS106

If AS106 link fails, backup via AS107 – and vice-versa

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **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 presentation for examples

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **When upstreams cannot or will not announce default route**

Because of operational policy against using “default-originate” on BGP peering

Solution is to use IGP to propagate default from the edge/peering routers

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **Router C Configuration**

```
router ospf 109
  default-information originate metric 30
  passive-interface Serial 0/0
!
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.1 remote-as 107
  neighbor 222.222.10.1 prefix-list rfc1918-deny in
  neighbor 222.222.10.1 prefix-list my-block out
  neighbor 222.222.10.1 filter-list 10 in
!
..next slide
```

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

```
ip prefix-list my-block permit 221.10.0.0/19
! See earlier presentation for RFC1918 list
!
ip route 221.10.0.0 255.255.224.0 null0
ip route 0.0.0.0 0.0.0.0 serial 0/0 254
!
ip as-path access-list 10 permit ^(107_)+$
ip as-path access-list 10 permit ^(107_)+_[0-9]+$
!
```

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **Router D Configuration**

```
router ospf 109
  default-information originate metric 10
  passive-interface Serial 0/0
!
router bgp 109
  network 221.10.0.0 mask 255.255.224.0
  neighbor 222.222.10.5 remote-as 106
  neighbor 222.222.10.5 prefix-list deny-all in
  neighbor 222.222.10.5 prefix-list my-block out
!
..next slide
```

Two Upstreams, One Local Peer: Partial Routes

Cisco.com

```
ip prefix-list deny-all deny 0.0.0.0/0 le 32
ip prefix-list my-block permit 221.10.0.0/19
! See earlier presentation for RFC1918 list
!
ip route 221.10.0.0 255.255.224.0 null0
ip route 0.0.0.0 0.0.0.0 serial 0/0 254
!
```


Two Upstreams, One Local Peer: Partial Routes

Cisco.com

- **Partial routes from upstreams**

Use OSPF to determine outbound path

Router D default has metric 10 – primary outbound path

Router C default has metric 30 – backup outbound path

Serial interface goes down, static default is removed from routing table, OSPF default withdrawn

Service Provider Multihoming

Case Study

Case Study

Requirements (1)

Cisco.com

- **ISP needs to multihome:**
 - To AS5400 in Europe**
 - To AS2516 in Japan**
 - /19 allocated by APNIC**
 - AS 17660 assigned by APNIC**
 - 1Mbps circuits to both upstreams**

Case Study

Requirements (2)

Cisco.com

- **ISP wants:**

- Symmetric routing and equal link utilisation in and out
(as close as possible)**

- international circuits are expensive**

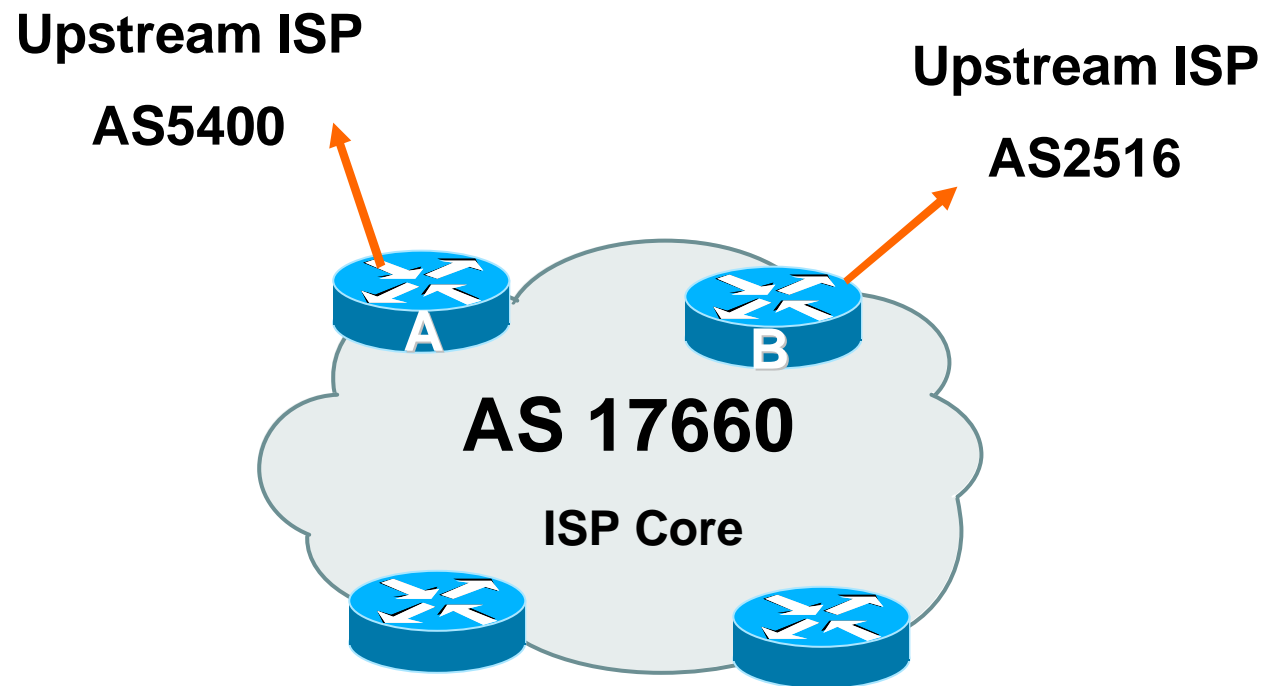
- Has two 2600 border routers with 64Mbytes memory**

- Cannot afford to upgrade memory or hardware on border
routers or internal routers**

- **“Philip, make it work, please”**

Case Study

Cisco.com



Allocated /19 from APNIC

Circuit to AS5400 is 1Mbps, circuit to AS2516 is 1Mbps

Case Study

Cisco.com

- Both providers stated that routers with 128Mbytes memory required for AS17660 to multihome

Wrong!

Full routing table is rarely required or desired

- **Solution:**

Accept default from one upstream

Accept partial prefixes from the other

Case Study

Inbound Loadsharing

Cisco.com

- **First cut: Went to a few US Looking Glasses**

Checked the AS path to AS5400

Checked the AS path to AS2516

AS2516 was one hop “closer”

Sent AS-PATH prepend of one AS on AS2516 peering

Case Study

Inbound Loadsharing

Cisco.com

- **Refinement**

Did not need any

First cut worked, seeing on average 600kbps inbound on each circuit

Does vary according to time of day, but this is as balanced as it can get, given customer profile



Case Study

Outbound Loadsharing

Cisco.com

- **First cut:**
 - Requested default from AS2516**
 - Requested full routes from AS5400**
- **Then looked at my Routing Report**
 - Picked the top 5 ASNs and created a filter-list**
 - If 701, 1, 7018, 1239 or 7046 are in AS-PATH, prefixes are discarded**
 - Allowed prefixes originated by AS5400 and up to two AS hops away**
 - Resulted in 32000 prefixes being accepted in AS17660**

Case Study

Outbound Loadsharing

Cisco.com

- **Refinement**

32000 prefixes quite a lot, seeing more outbound traffic on the AS5400 path

Traffic was very asymmetric

out through AS5400, in through AS2516

Added the next 3 ASNs from the Top 20 list

209, 2914 and 3549

Now seeing 14000 prefixes

Traffic is now evenly loadshared outbound

Around 200kbps on average

Mostly symmetric

Case Study

Configuration Router A

Cisco.com

```
router ospf 100
  log-adjacency-changes
  passive-interface default
  no passive-interface Ethernet0/0
  default-information originate metric 20
!
router bgp 17660
  no synchronization
  no bgp fast-external-fallover
  bgp log-neighbor-changes
  bgp deterministic-med
...next slide
```

Case Study

Configuration Router A

Cisco.com

```
neighbor 166.49.165.13 remote-as 5400
neighbor 166.49.165.13 description eBGP multihop to AS5400
neighbor 166.49.165.13 ebgp-multihop 5
neighbor 166.49.165.13 update-source Loopback0
neighbor 166.49.165.13 prefix-list in-filter in
neighbor 166.49.165.13 prefix-list out-filter out
neighbor 166.49.165.13 filter-list 1 in
neighbor 166.49.165.13 filter-list 3 out
!
prefix-list in-filter deny rfc1918etc in
prefix-list out-filter permit 202.144.128.0/19
!
ip route 0.0.0.0 0.0.0.0 serial 0/0 254
...next slide
```

Case Study

Configuration Router A

Cisco.com

```
ip as-path access-list 1 deny _701_  
ip as-path access-list 1 deny _1_  
ip as-path access-list 1 deny _7018_  
ip as-path access-list 1 deny _1239_  
ip as-path access-list 1 deny _7046_  
ip as-path access-list 1 deny _209_  
ip as-path access-list 1 deny _2914_  
ip as-path access-list 1 deny _3549_  
ip as-path access-list 1 permit _5400$  
ip as-path access-list 1 permit _5400_[0-9]+$  
ip as-path access-list 1 permit _5400_[0-9]+_[0-9]+$  
ip as-path access-list 1 deny .*  
ip as-path access-list 3 permit ^$  
!
```

Case Study

Configuration Router B

Cisco.com

```
router ospf 100
  log-adjacency-changes
  passive-interface default
  no passive-interface Ethernet0/0
  default-information originate
!
router bgp 17660
  no synchronization
  no auto-summary
  no bgp fast-external-fallover
...next slide
```

Case Study

Configuration Router B

Cisco.com

```
bgp log-neighbor-changes
bgp deterministic-med
  neighbor 210.132.92.165 remote-as 2516
  neighbor 210.132.92.165 description eBGP peering
  neighbor 210.132.92.165 soft-reconfiguration inbound
  neighbor 210.132.92.165 prefix-list default-route in
  neighbor 210.132.92.165 prefix-list out-filter out
  neighbor 210.132.92.165 route-map as2516-out out
  neighbor 210.132.92.165 maximum-prefix 100
  neighbor 210.132.92.165 filter-list 2 in
  neighbor 210.132.92.165 filter-list 3 out
!
```

...next slide

Case Study

Configuration Router B

Cisco.com

```
!  
prefix-list default-route permit 0.0.0.0/0  
prefix-list out-filter permit 202.144.128.0/19  
!  
ip as-path access-list 2 permit _2516$  
ip as-path access-list 2 deny .*  
ip as-path access-list 3 permit ^$  
!  
route-map as2516-out permit 10  
    set as-path prepend 17660  
!
```


Configuration Summary

Cisco.com

- **Router A**

Hears full routing table – throws away most of it

AS5400 BGP options are all or nothing

Static default pointing to serial interface – if link goes down, OSPF default removed

- **Router B**

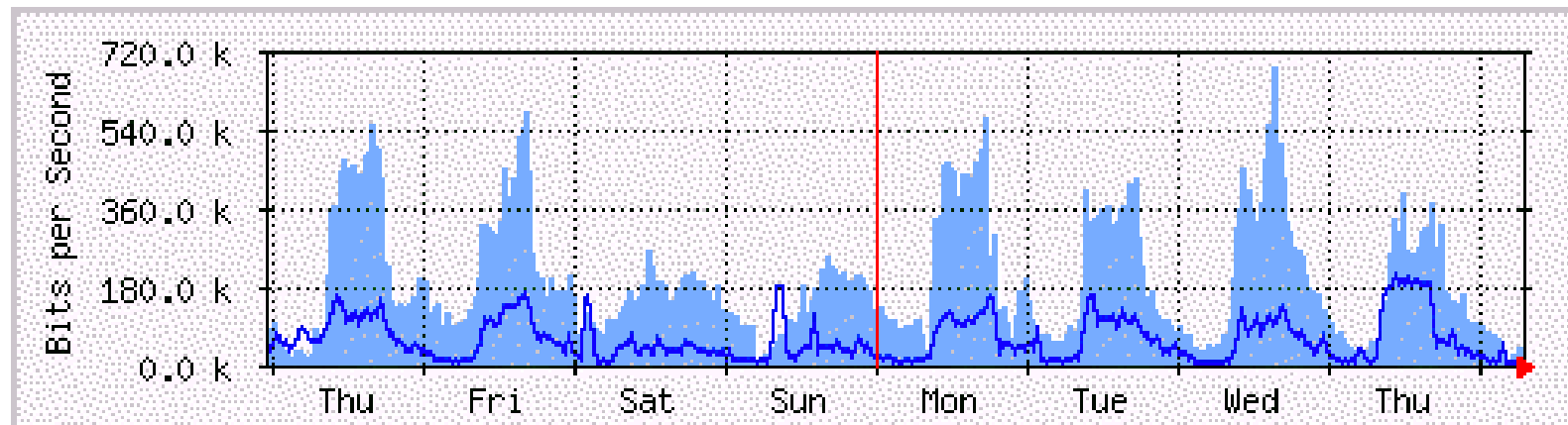
Hears default from AS2516

If default disappears (BGP goes down or link goes down), OSPF default is removed

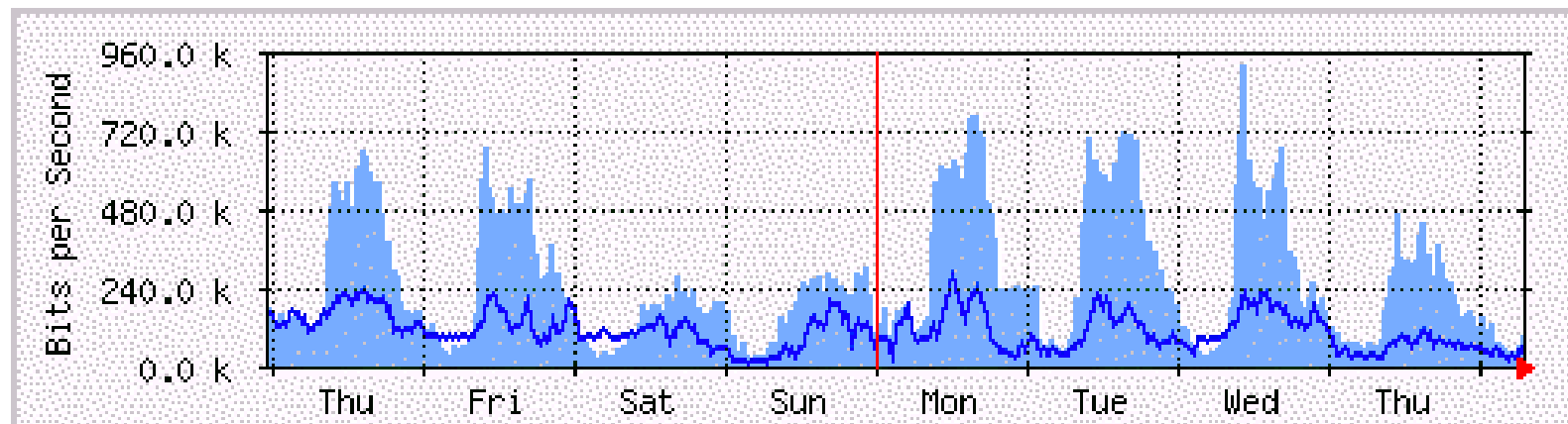
Case Study

MRTG Graphs

Cisco.com



Router A to AS5400



Router B to AS2516

Case Study Summary

Cisco.com

- **Multihoming is not hard, really!**
 - Needs a bit of thought, a bit of planning**
 - Use this case study as an example strategy**
 - Does not require sophisticated equipment, big memory, fast CPUs...**

BGP for Internet Service Providers

Cisco.com

- **BGP Basics (quick recap)**
- **Scaling BGP**
- **Deploying BGP in an ISP network**
- **Multihoming Examples**

BGP for Internet Service Providers

End of Tutorial