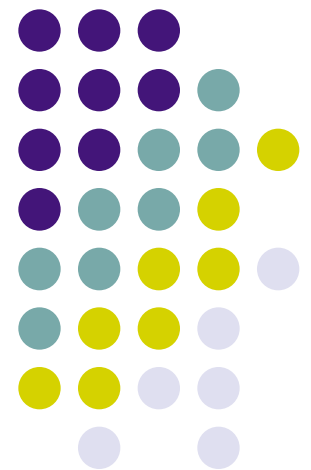
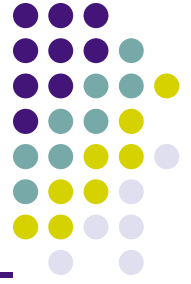


BGP

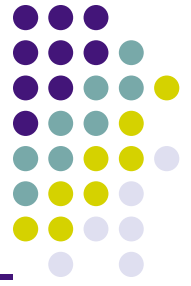
BGP protocol
iBGP configuration
BGP convergence



References



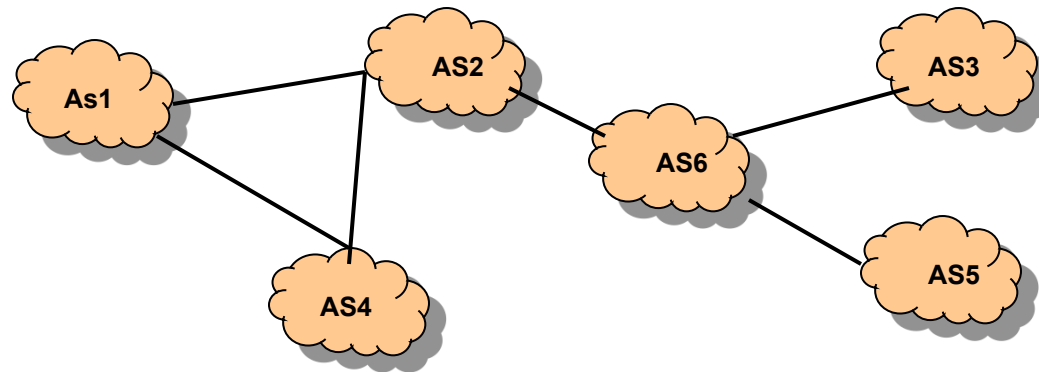
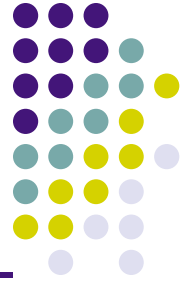
1. BGP tutorial - BPG4 case studies by Sam Halabi
2. BGP routing policies in ISP networks by Mathew Caesar and Jennifer Rexford



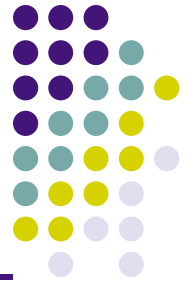
Autonomous System

- Definition: internet is network of networks glued by IP
- Within a network (intra-domain) any routing policy can be chosen
- A common routing policy is needed when routing between networks or domains
- A Domain is a network that has unified administrative routing policy
- Autonomous System (domain) or AS: Has a number assigned to it and provides routing information to other ASes

Internet structure



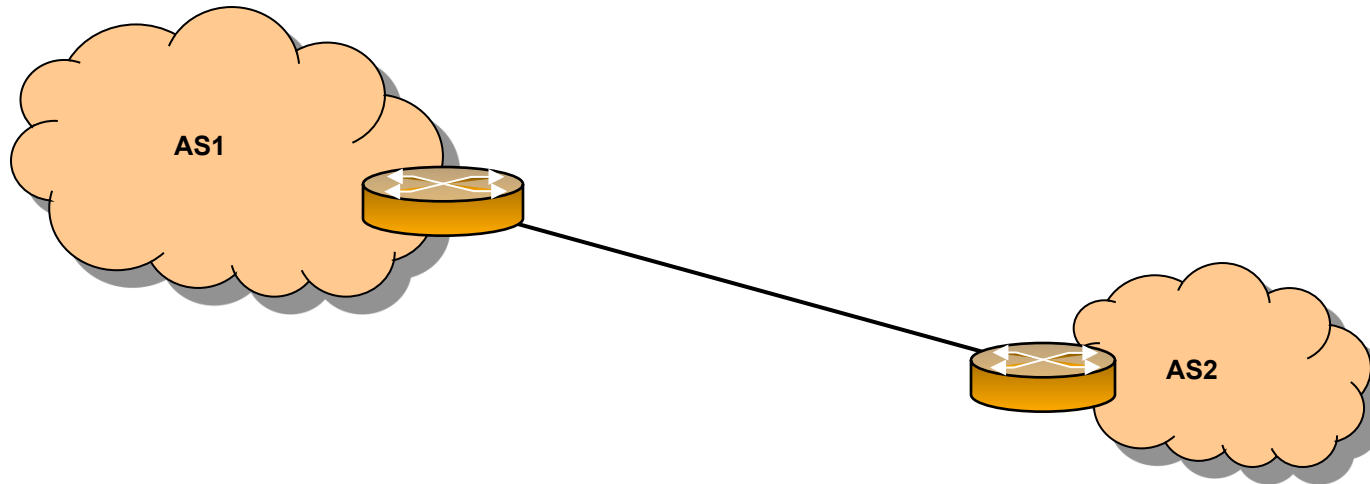
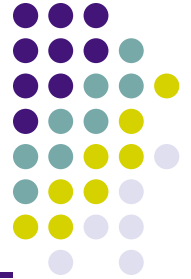
- AS provide reachability information to other ASes
- Within AS, local routing protocols used (optimize path metric)
- Inter-AS concerned with reachability and policy implementation
 - Usually \$\$ involved with relationships



Autonomous system

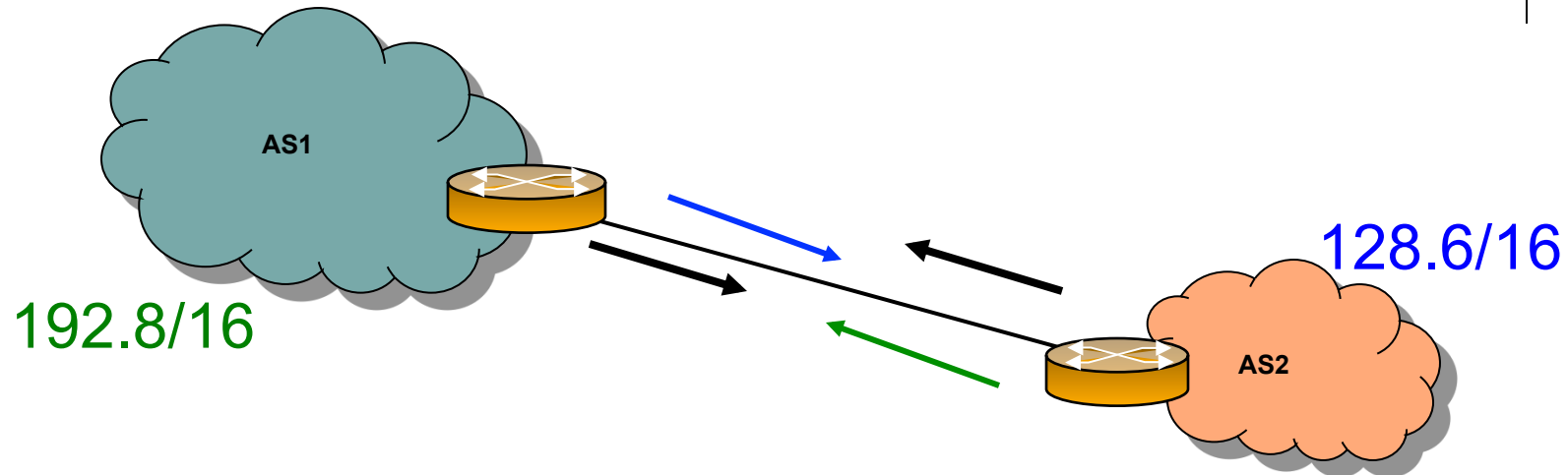
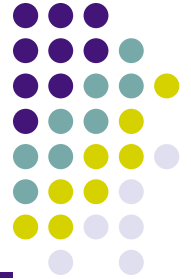
- The actual entity that participates in interdomain routing
- Has a unique 16 bit number assigned
- Examples:
 - RUTGERS: 46, STANFORD;32, MIT: 3, CMU: 9
 - AT&T: 6431, ...
 - Quest: 209, ...
 - Sprint: 1239, ...
- How do ASes interconnect to provide global connectivity?
- How does routing information get exchanged?
- How is policy specified and implemented?

Internet routing



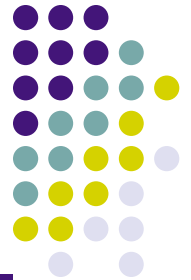
- Intra domain
 - OSPF, RIP
 - Route on IP addresses
 - Path metrics
- Inter domain
 - BGP
 - Route on AS numbers
 - Policy and business relations based

BGP: basic idea

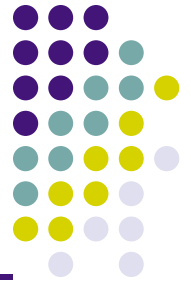


- AS1 needs to inform AS2 that it can route to 192.8/16 and AS2 needs to inform AS1 that it can route to 128.16/16
- After this, what else
 - Route updates/changes
 - Policy: what is AS1 does not want to route to anyone else but its own domain?
 - What paths should be preferred?
- This is essentially BGP

BGP protocol

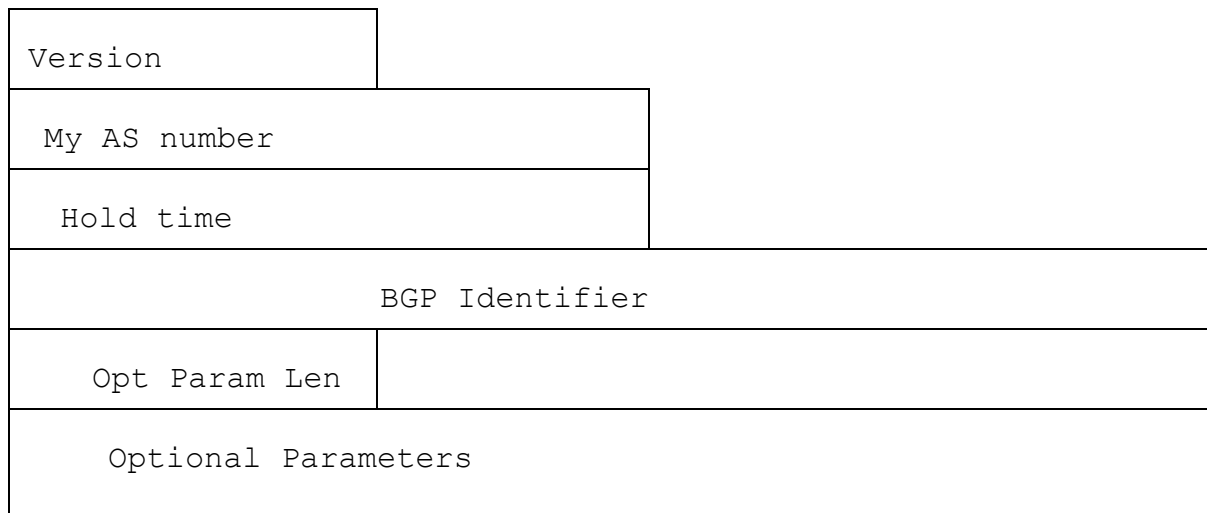


- BGP uses TCP as its transport protocol, on port 179. On connection start, BGP peers exchange complete copies of their routing tables, which can be quite large. However, only changes (deltas) are then exchanged, which makes long running BGP sessions more efficient than shorter ones.
- Four Basic messages:
 - *Open*:
Establishes BGP session (uses TCP port #179)
 - *Notification*:
Report unusual conditions
 - *Update*:
Inform neighbor of new routes that become active
Inform neighbor of old routes that become inactive
 - *Keepalive*:
Inform neighbor that connection is still viable

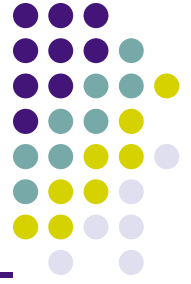


OPEN Message

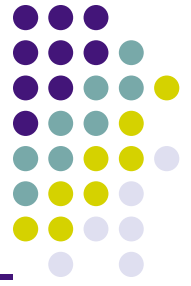
- During session establishment, two BGP speakers exchange their
 - AS numbers
 - BGP identifiers (usually one of the router's IP addresses)
 - Select hold timer : max time before declaring peer is down
- A BGP speaker has option to refuse a session
- authentication information (optional)



NOTIFICATION and KEEPALIVE Messages



- NOTIFICATION
 - Indicates an error
 - terminates the TCP session
 - gives receiver an indication of why BGP session terminated
 - Examples: header errors, hold timer expiry, bad peer AS, bad BGP identifier, malformed attribute list, missing required attribute, AS routing loop, etc.
- KEEPALIVE
 - protocol requires some data to be sent periodically.
If no UPDATE to send within the specified time period,
then send KEEPALIVE message
to assure partner that connection is still alive

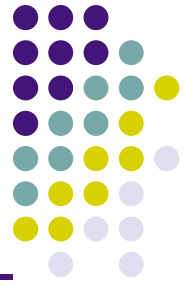


UPDATE Message

- used to either advertise and/or withdraw previously announced prefixes
- path attributes: list of attributes that pertain to ALL the prefixes in the Reachability Info field

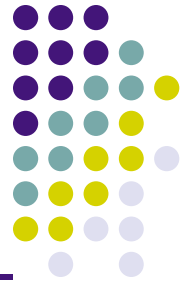
FORMAT:

Withdrawn routes length (2 octets)	
Withdrawn routes (variable length)	
Total path attributes length (2 octets)	
Path Attributes (variable length)	
(NLRI) Reachability Information (variable length)	



BGP update message

- Withdrawn Routes: Length field 2 Bytes
- Withdrawn route list
- Path attributes: Length field 2 bytes
- Path attributes list
- NLRI list : a list of entries
 - Length field (1 byte), Prefix (variable length)
- Path attributes apply to all the prefixes in the NLRI list



Advertising a prefix

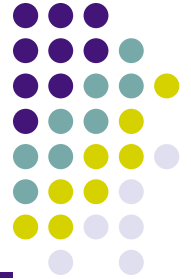
- When a router advertises a prefix to one of its BGP neighbors:
 - information is valid until first router explicitly advertises that the information is no longer valid
 - BGP does not require routing information to be refreshed
 - if node A advertises a path for a prefix to node B, then node B can be sure node A is using that path itself to reach the destination.



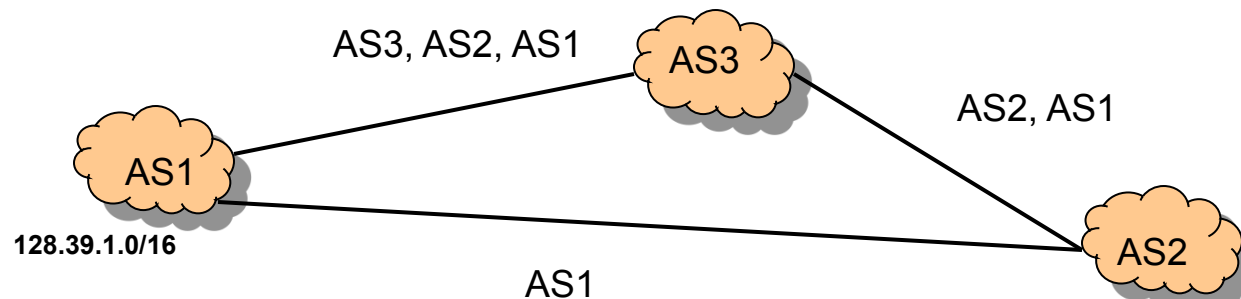
BGP attributes

- BGP protocol announcements carries with it several attributes
- Attribute describes characteristics of a prefix
- BGP chooses a single path for a given prefix based on attributes (can choose to ignore!)
- BGP always announces the best path to neighbors
- Attributes
 - 1 ORIGIN
 - 2 AS_PATH
 - 3 NEXT_HOP
 - 4 MED
 - 5 LOCAL_PREF
 - 6 WEIGHT
 - 7 COMMUNITY
 - 8 AGGREGATOR

PATH ATTRIBUTES



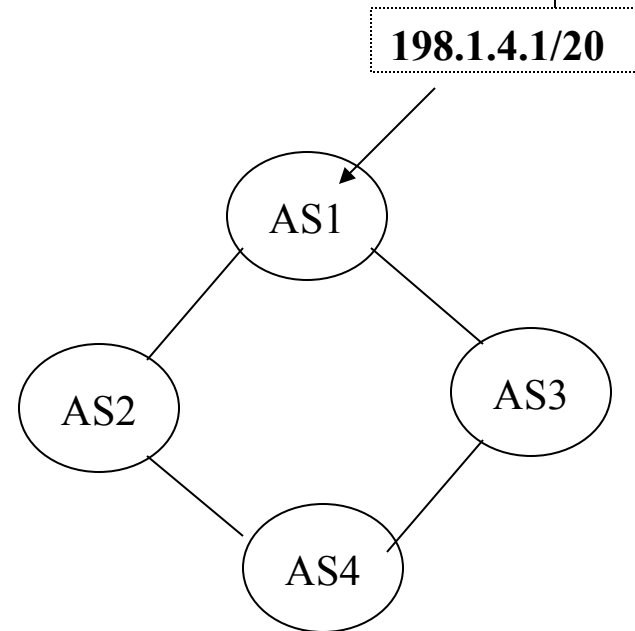
- **ORIGIN(TYPE CODE=1):**
 - Who originated the announcement? Where was a prefix *injected* into BGP?
 - Manually configured, directly connected, by other intra-routing protocols
 - IGP, EGP, default – incomplete (learnt from some other means)
- **AS-PATH (TYPE CODE =2)**
 - a list of AS' s through which the announcement for a prefix has passed
 - each AS prepends its AS # to the AS-PATH attribute when forwarding an announcement
 - useful to detect and prevent loops
 - AS length can be used to select among routes unless a LOCAL PREF attribute overrides



Attribute: Local Preference (type code = 5)



- Used to indicate preference among multiple paths for the same prefix *anywhere* in the internet.
- The higher the value the more it is preferred
- Default value is 100
- Local to the AS (non-transitive)
- Often used to select a specific exit point for **outbound** traffic
- Override influence of AS path length



BGP table at AS4:

Destination	AS Path	Local Pref
198.1.4/20	AS3 AS1	300
198.1.4/20	AS2 AS1	100

Use of local pref

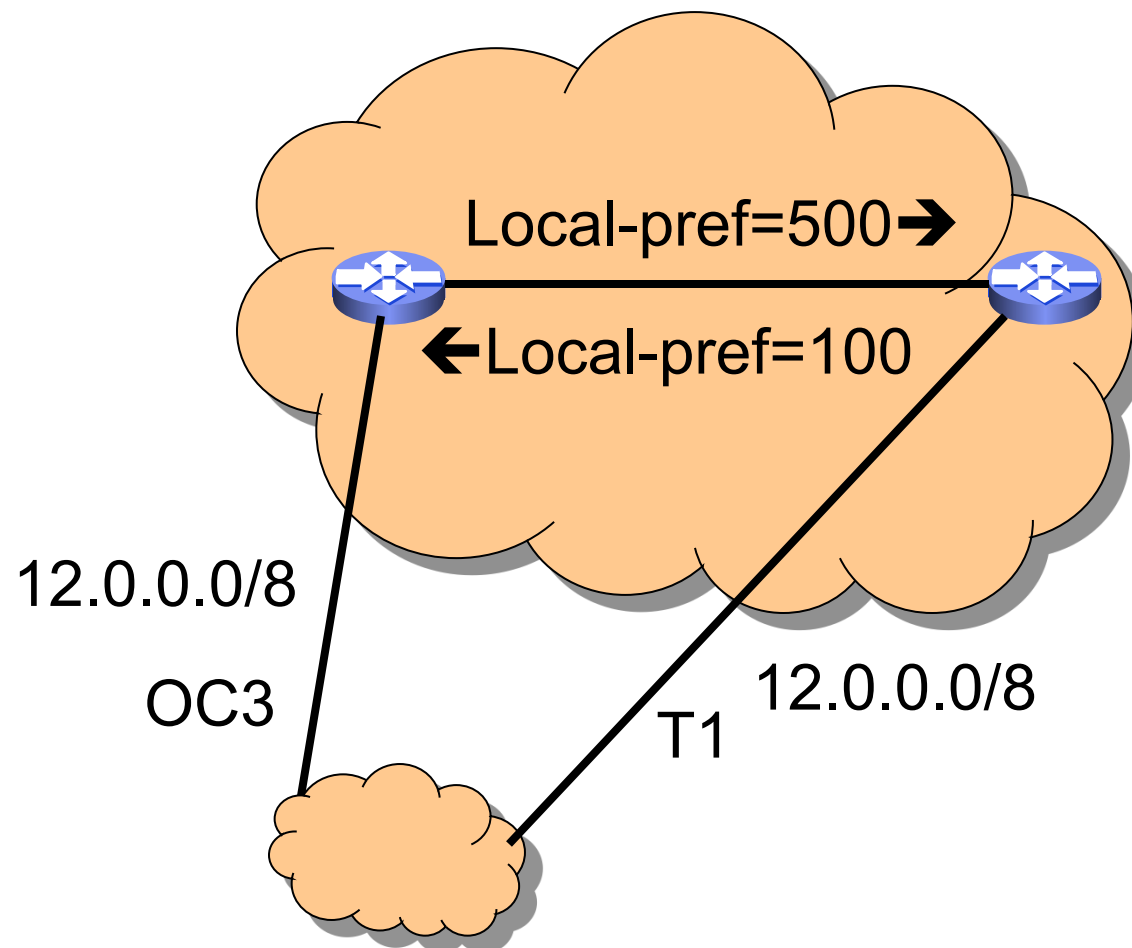


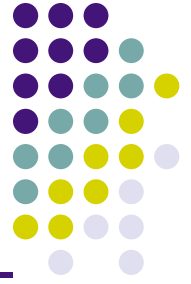
Diagram illustrating a BGP routing scenario. A central cloud contains three routers labeled A, B, and C. Router A has IP 128.64.1.1, Router B has IP 128.64.3.1, and Router C has IP 128.64.8.2. All three routers are connected via iBGP. Router A is connected to an external router B? (labeled B?) via a yellow arrow. Router C is connected to an external router 198.6.4.0/22 via a black arrow. Router A is also connected to an external router 198.6.4.0/22 via a blue arrow labeled eBGP. A table titled "BGP Table at R" shows the destination 198.6.4.0/22.

Destination	N
198.6.4.0/22	

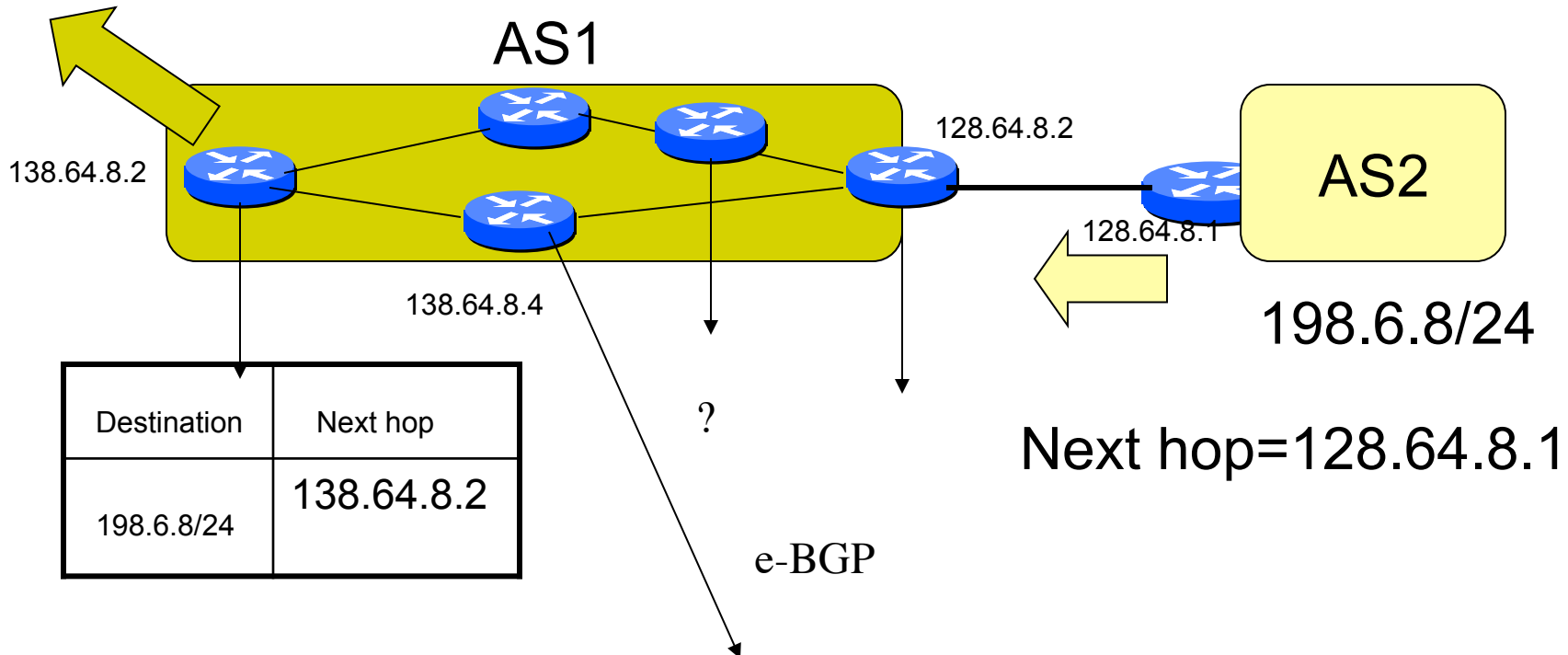
Destination	Nexthop
198.6.4.0/22	128.64.8.1
BGP Table at Router B:	

Destination	Nextthop
198.6.4.0/22	

Use of next hop



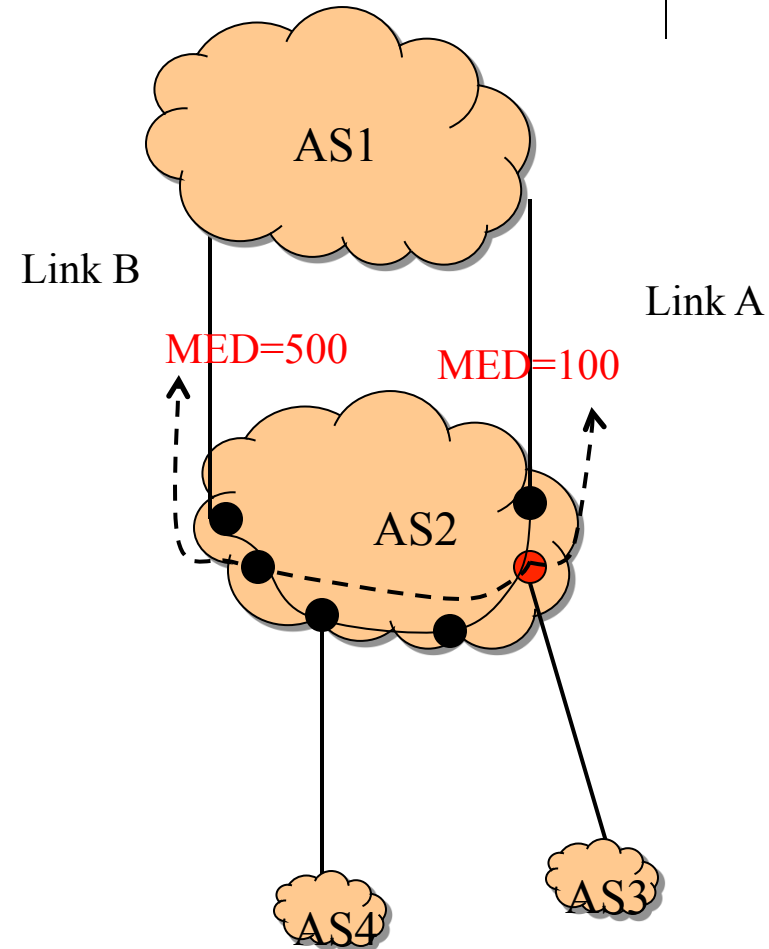
Next hop=138.64.8.2 for 198.6.8/24

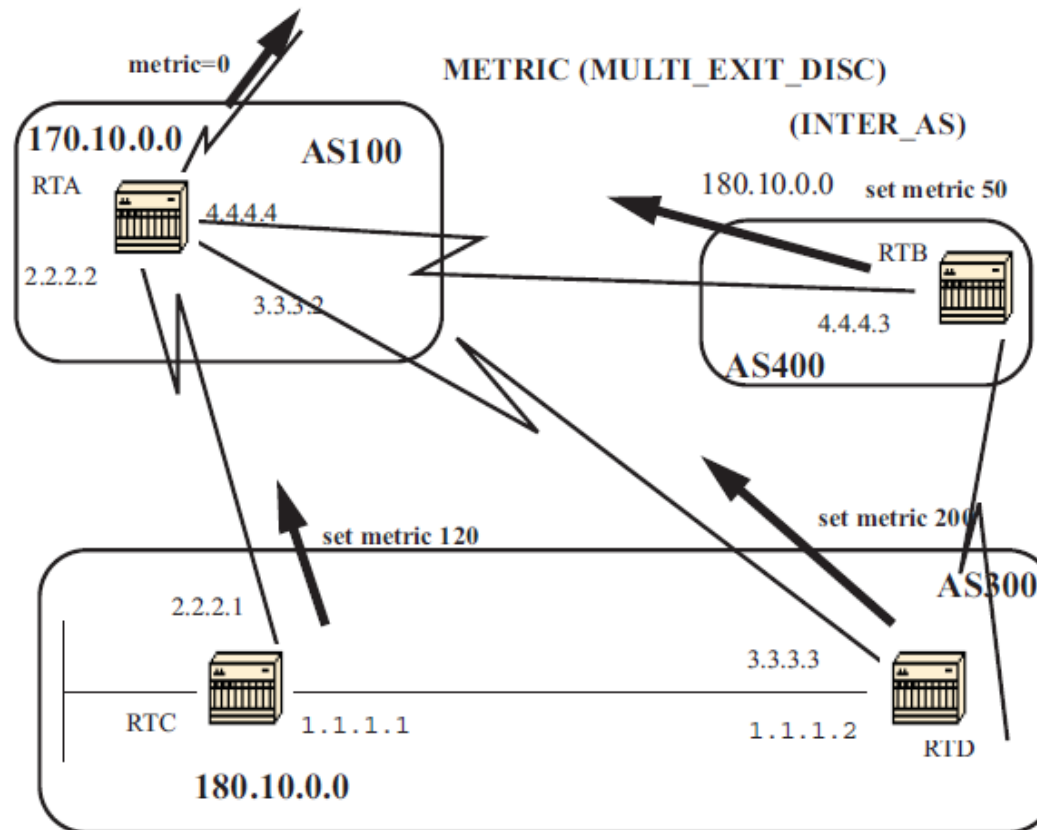
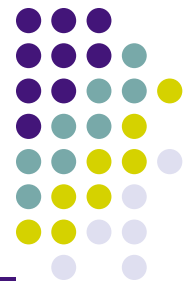




Attribute: Multi-Exit Discriminator (MED) (code=4)

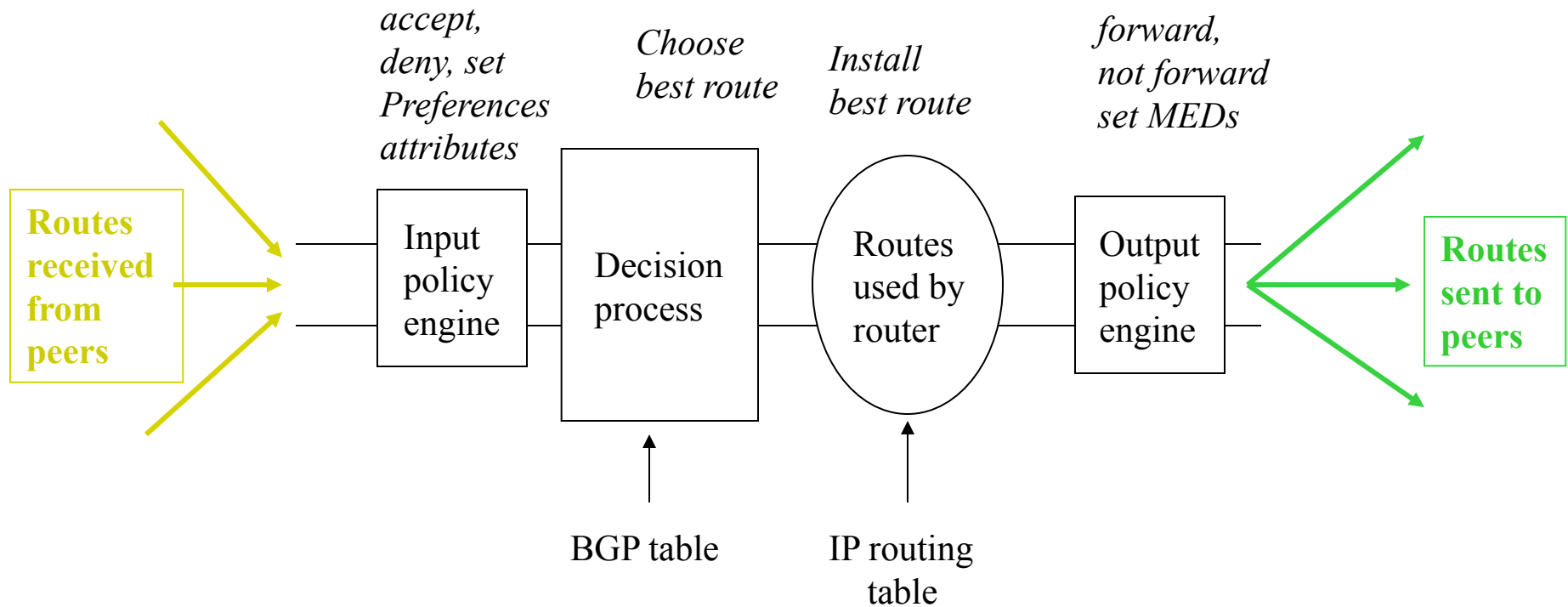
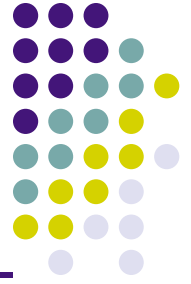
- when AS' s interconnected via 2 or more links
- AS path length are same
- AS announcing prefix, sets MED value
- enables AS2 to indicate its preference (lower MED is better)
- AS receiving prefix uses MED to select link
- a way to specify how close a prefix is to the link it is announced on



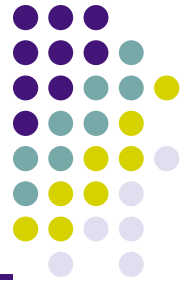


- Med values from the same AS are compared
- A lower MED value is preferred
- MED values exchanged between ASs- non-transitive

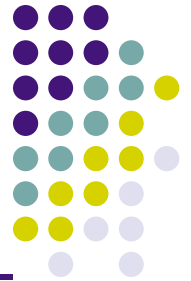
Routing Process Overview



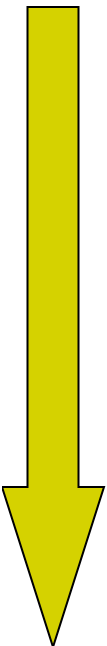
BGP Decision Process: Path Selection on a Router



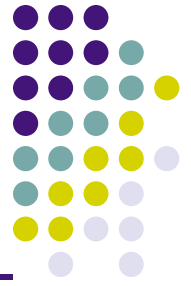
- Routing Information Base
 - Store all BGP routes for each destination prefix
 - Withdrawal message: remove the route entry
 - Announcement message: update the route entry
- Selecting the best route
 - Consider all BGP routes for the prefix
 - Apply rules for comparing the routes
 - Select the one best route
 - Use this route in the forwarding table
 - Send this route to neighbors



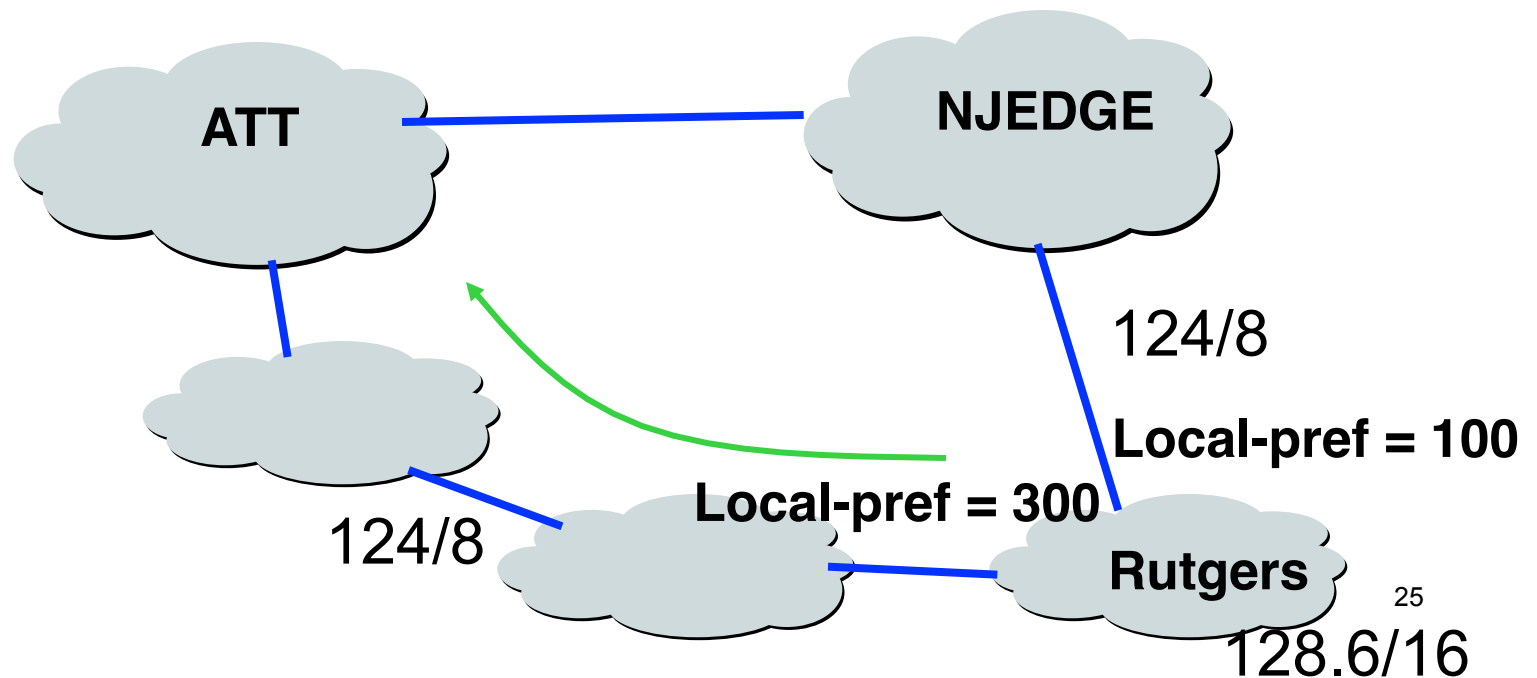
BGP Decision Process

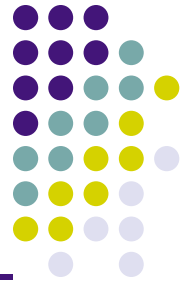
- 
1. Choose route with highest LOCAL-PREF
 2. If have more than 1 route, select route with shortest AS-PATH
 3. If have more than 1 route, select according to lowest ORIGIN type where $IGP < E\text{-}BGP < \text{default}$
 4. If have more than 1 route, select route with lowest MED value
 5. Select e-BGP learned over i-BGP learned path
 6. Select min cost path to NEXT HOP using IGP metrics (lowest IGP cost to BGP egress)
 7. If have multiple internal paths, use lowest BGP Router ID to break the tie.
- See: BGP routing policies in ISP networks by Caesar and Rexford

Import Policy: Local Preference



- Favor one path over another based on local policy
 - Override the influence of AS path length
 - Local admin policy given priority
 - Apply local policies to prefer a path

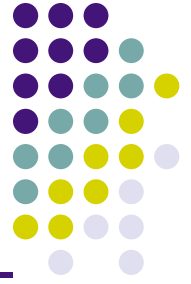




Import Policy: Filtering

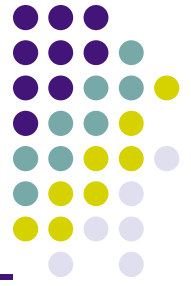
- Discard some route announcements
 - Detect configuration mistakes and attacks
- Examples on session to a customer
 - Discard route if prefix not owned by the customer
 - Does not want routing for that prefix via the peer

Export Policy: Filtering

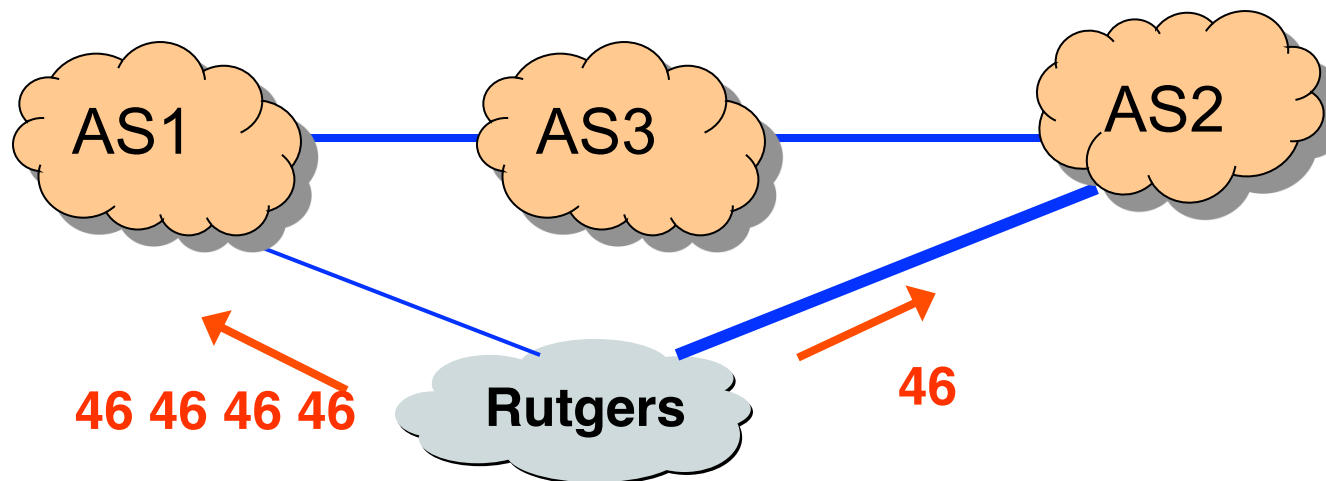


- Discard some route announcements
 - Limit propagation of routing information
- Not forwarding prefixes
 - Do not want others to use you as an intermediary for that prefix

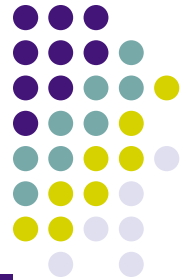
Export Policy: Attribute Manipulation



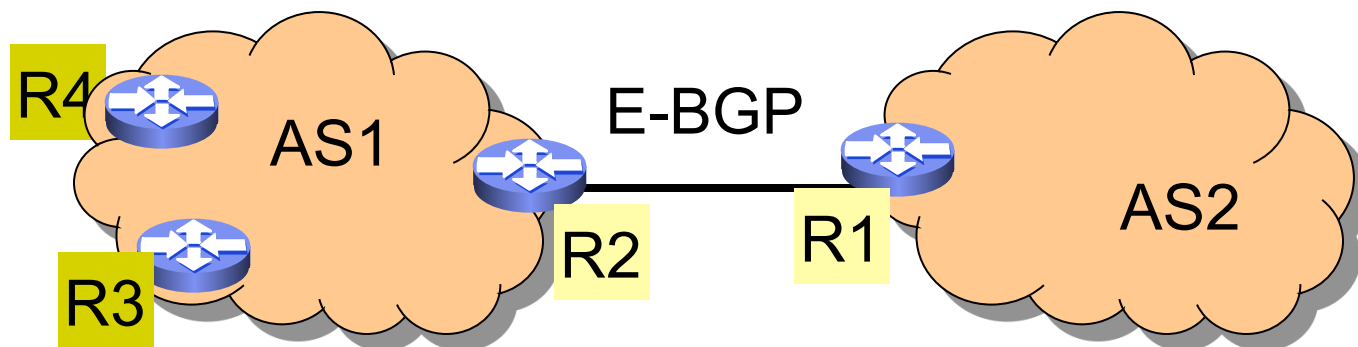
- Modify attributes of the active route
 - To influence the way other ASes behave
- Example: AS_PATH padding
 - Artificially inflate AS path length seen by others
 - Convince some ASes to send traffic another way
 - May not work always: AS2 may have a higher LOCAL_PREFERENCE



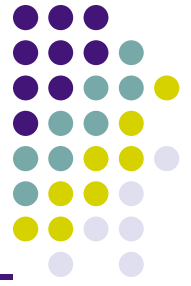
Internal vs. External BGP



- Internal-BGP or i-BGP used to distribute routes within AS
- Egress routers use E-BGP or BGP
- R4 and R3 learn routes from R2 using i-BGP
- R1 and R2 talk e-BGP (different AS)
- R2, R4 and R2, R3 and R3, R4 taal i-BGP (same AS)

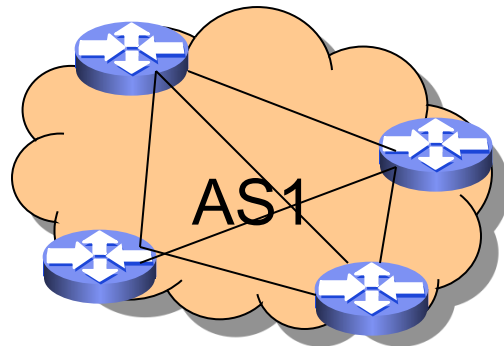
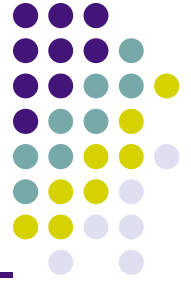


Internal BGP (I-BGP)

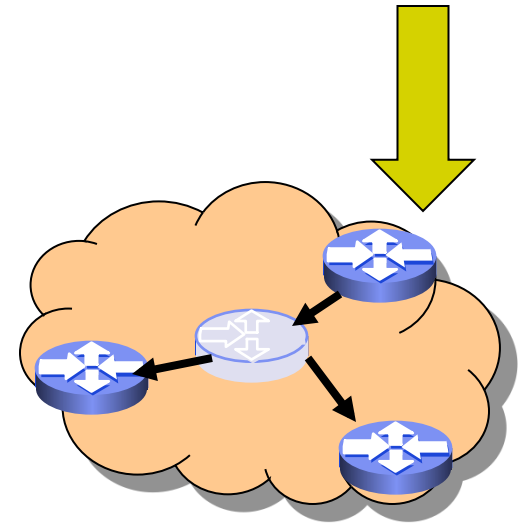


- Same messages as E-BGP
- Different rules about re-advertising prefixes:
 - Rule #1: Prefix learned from E-BGP can be advertised to I-BGP neighbor and vice-versa, but
 - Rule #2: Prefix learned from one I-BGP neighbor cannot be advertised to another I-BGP neighbor
 - Reason: no AS PATH within the same AS and thus danger of looping.
 - Means each I-BGP speaker must be connected directly with every other I-BGP within the same AS
 - Full MESH!!!

Route reflectors

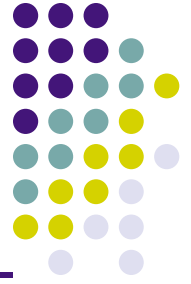


Mesh does not Scale
 $O(N^2)$ sessions

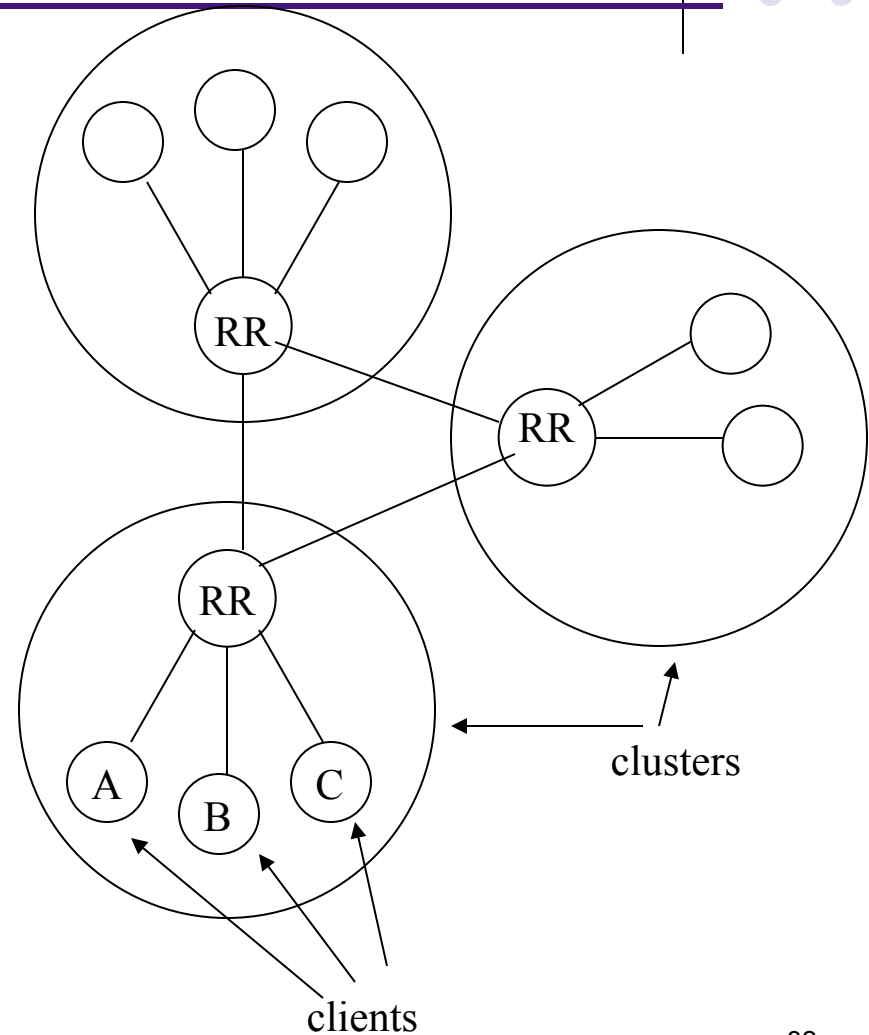


Only $N-1$ sessions
The RR only advertises best routes

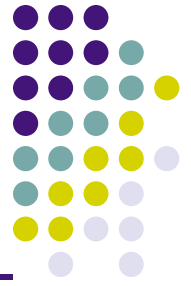
Route Reflectors



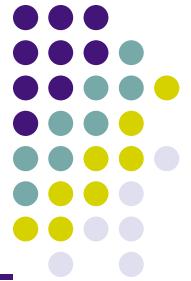
- Problem: requiring a full mesh of I-BGP sessions between all pairs of routers is hard to manage for large AS' s.
- Solution:
 - group routers into **clusters**.
 - Assign a leader to each cluster, called a **route reflector** (RR).
 - Members of a cluster are called **clients** of the RR
- I-BGP Peering
 - clients peer only with their RR
 - RR' s must be fully meshed



Route Reflectors: Rule on Announcements



- If received from RR, reflect to clients
- If received from a client, reflect to RRs and clients
- If received from E-BGP, reflect to all - RRs and clients
- RR' s reflect only the best route to a given prefix, not all announcements they receive.
 - helps size of routing table
 - sometimes clients don' t need to carry full table



Announcement loop

CISCO manual on BGP configuration



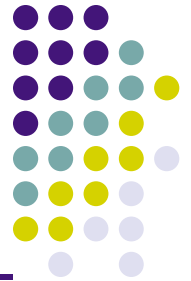
Caution Incorrectly setting BGP attributes for a route reflector can cause inconsistent routing, routing loops, or a loss of connectivity. Setting BGP attributes for a route reflector should be attempted only by an experienced network operator.

Command	Purpose
Router(config-router)# no bgp client-to-client reflection	Disables client-to-client route reflection.

RFC 4456- BGP Route Reflectors

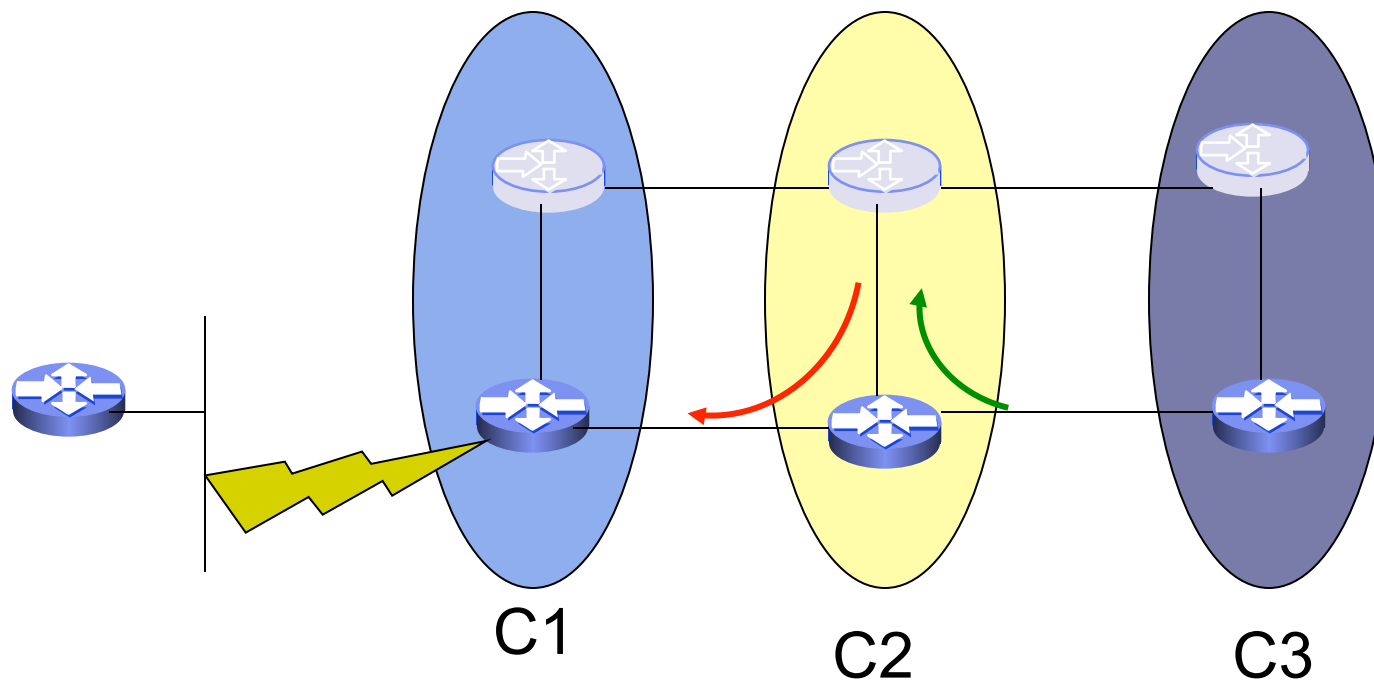
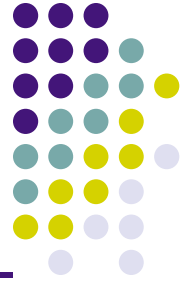
When a route is reflected, it is possible through misconfiguration to form route re-distribution loops.

Avoiding Loops with Route Reflectors



- Loops cannot be detected by traditional approach using AS-PATH because AS-PATH not modified within an AS.
- Announcements could leave a cluster and re-enter it.
- Two new attributes introduced:
 - ORIGINATOR_ID: router id of route's originator in AS
rule: announcement discarded if returns to originator
 - CLUSTER_LIST: a sequence of cluster id's. set by RRs.
rule: if an RR receives an update and the cluster list contains its cluster id, then update is discarded.

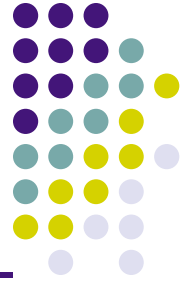
Announcement loops prevention



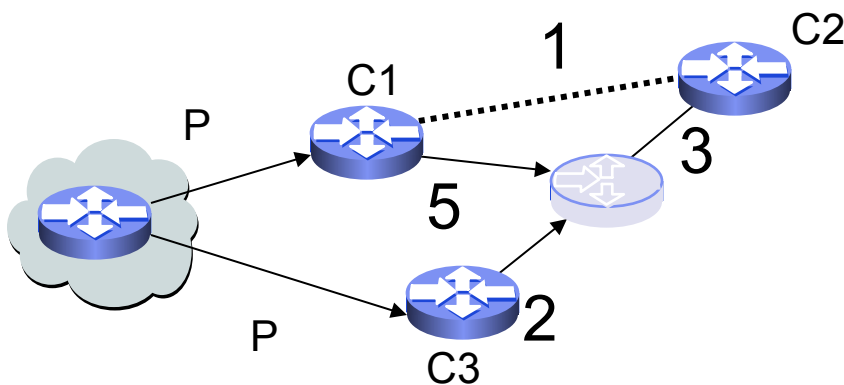
Drop if Originator ID= router id

Drop if Cluster List contains ClusterID

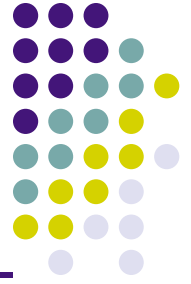
Route reflector vs Full mesh



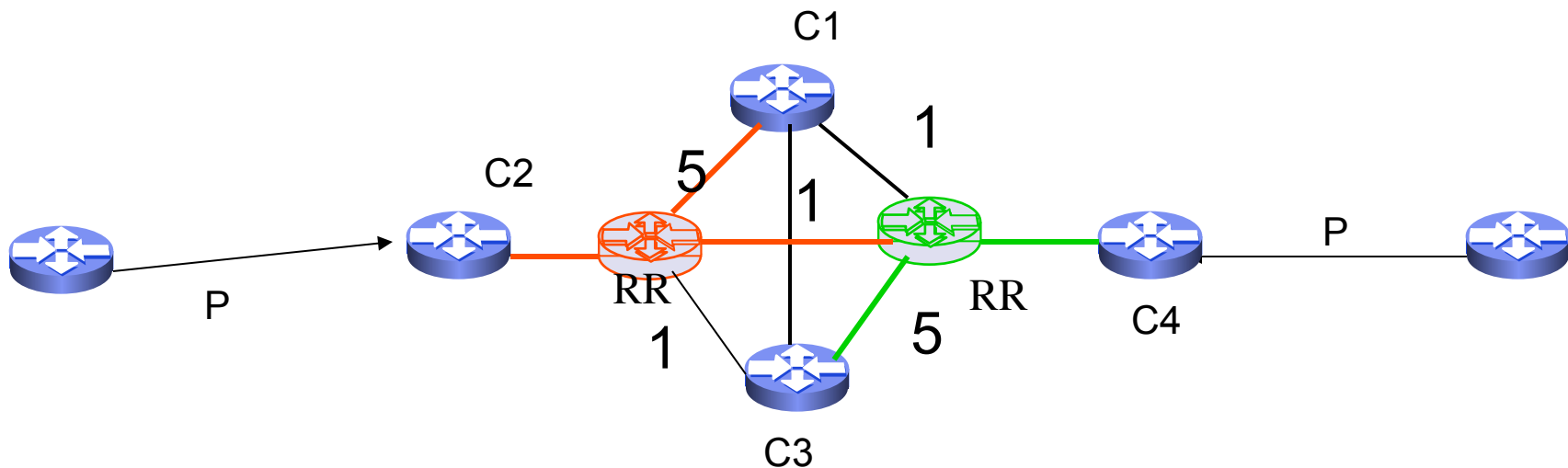
- In a full mesh, every router hears about every egress announcement
- Has complete visibility ; each router picks the shortest IGP path (among all routers announcing a prefix)
- Not so with RRs
- Who does C2 choose as the egress point?



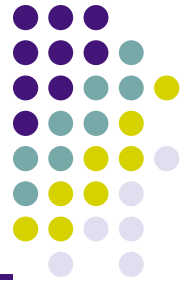
Forwarding loops



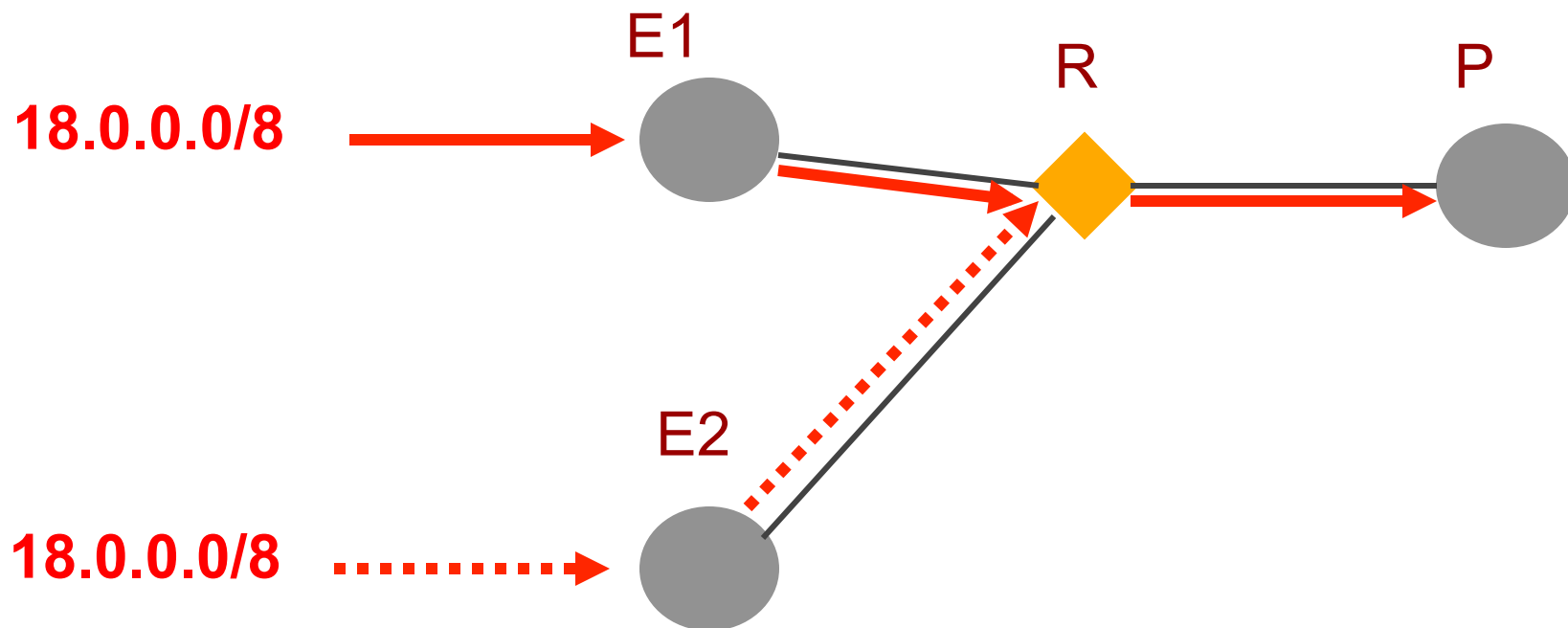
- Each router picks the shortest path among the routes it has heard
- Two different routers can consider each other has the intermediary to the shortest path
- Forwarding loop!!



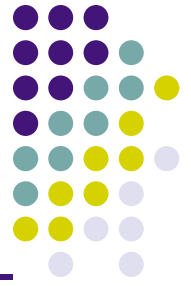
Key insight for emulating full-mesh



- For every BGP router P, every egress E
 - P and E have iBGP session, OR \rightarrow if true for all P, what do we have?
 - P should be the client of a **route reflector on the shortest path** between P and E

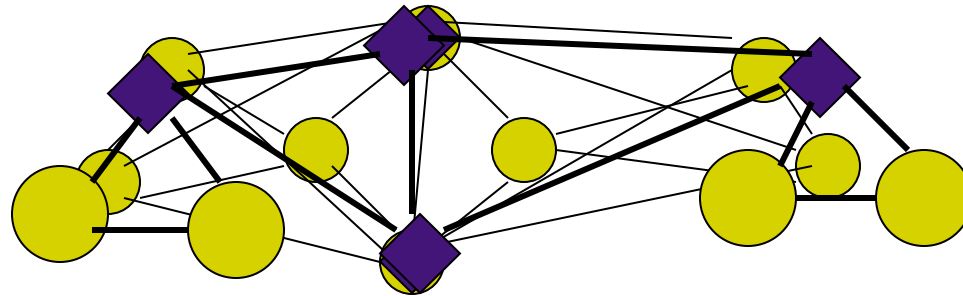
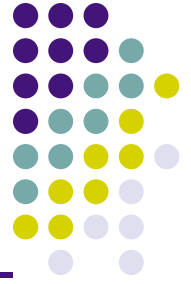


Can we do better than full mesh?

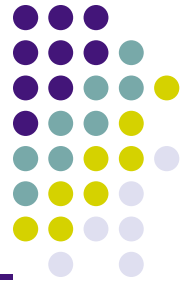


- Graph separator
- Choose a separator and make them route reflectors
- Connect RRs into a mesh
- Make members of connected components clients of all RRs in the separator
- Connect all members in each connected component into a mesh
- Recursively apply to the components

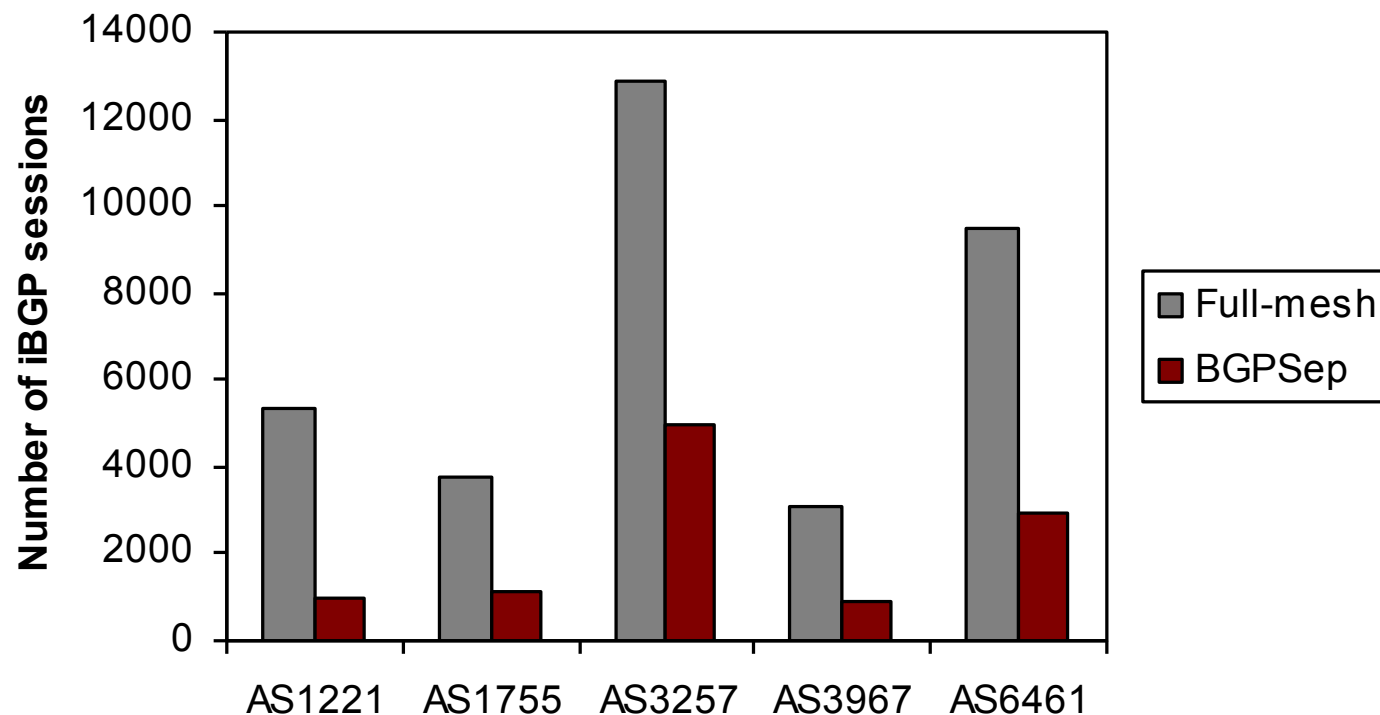
Separator Algorithm

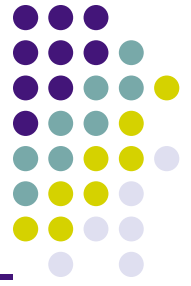


Evaluation



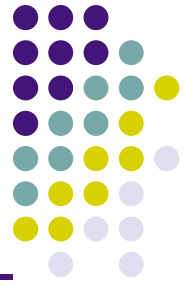
- 2.5 to 5X fewer iBGP sessions on ISP topologies [Source: Rocketfuel]





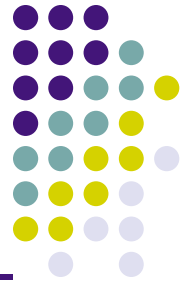
BGP convergence delay

- How long before a route change converges in the network
- Two Time factors
- Time to detect a failure
 - Keep-alive 60 seconds
 - Hold timer: 180 seconds
- On failure detection, throw away peer routes and announce changes



Route change propagation

- New route announcement requires path exploration
- Path Path exploration is expensive
 - Large number of possible paths
 - Might have to explore (nearly) all of them
- Minimum Route Advertisement Interval
 - Minimum time between advertisement of routes for a given destination to a given neighbor
 - allows for combining multiple messages in one
 - Typical value of 30 seconds
- Convergence delay
 - $(30 \text{ seconds}) * (\# \text{ of paths}) + 180 \text{ seconds}$



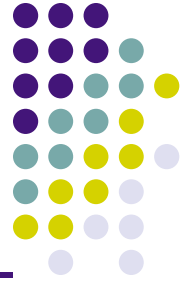
Route changes

- Link failures
- Reachability issues (router reboot)
- Session resets
- Lots of path changes
- How to deal with transient changes
- Do not want to be too quick
- At the same time not wait too long
- Strike a balance
- Route Flap damping (RFD)

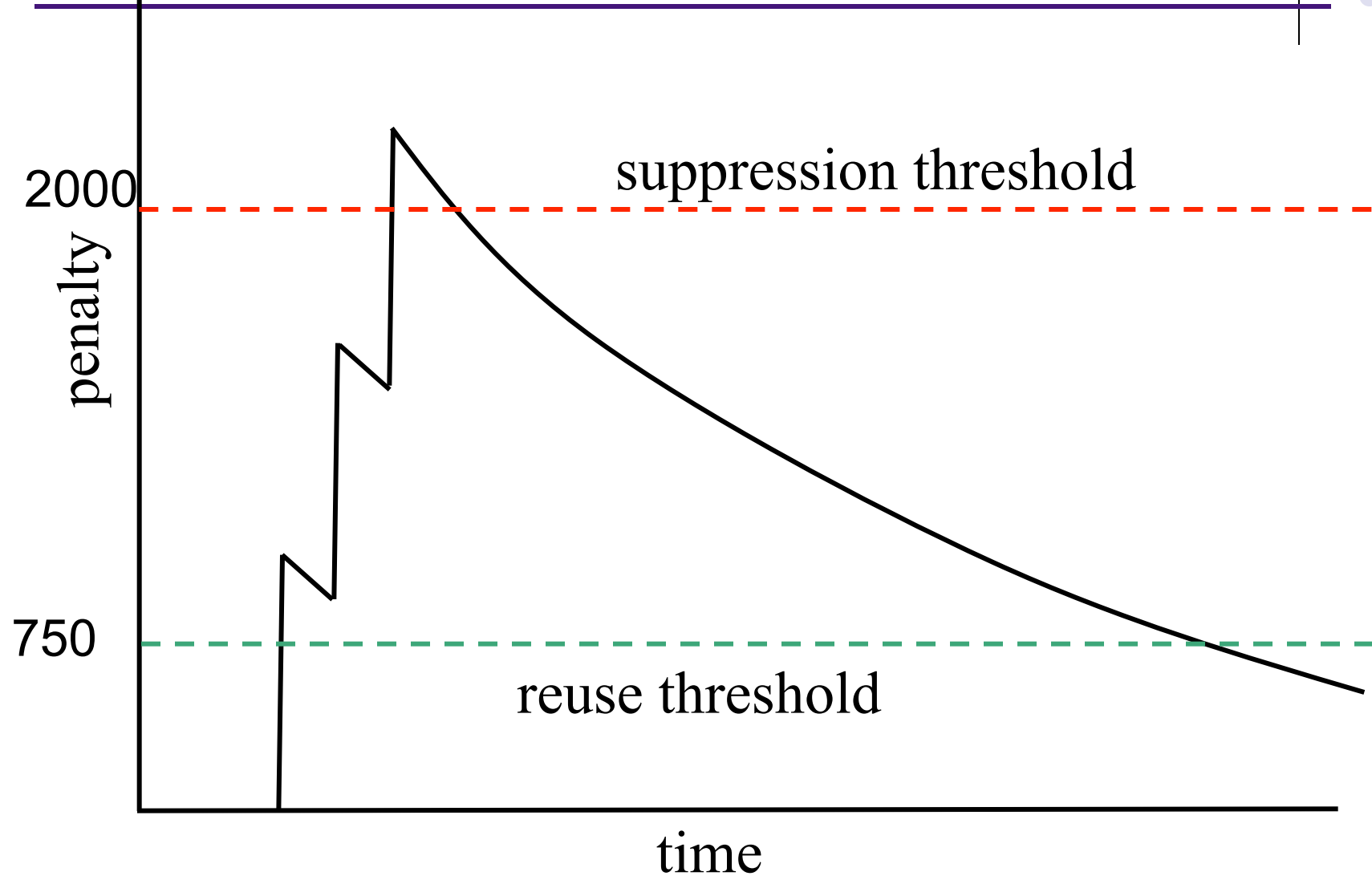


Route flap damping

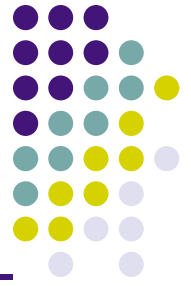
- First flap on a route (withdrawn & readvertised), asses penalty (1000 points), put the prefix in historical category
- Second flap (another 1000 points), do not advertise this route to others
- Penalty is decayed if it does not flap further
- Once the penalty falls below 750, the route is removed from dampened state



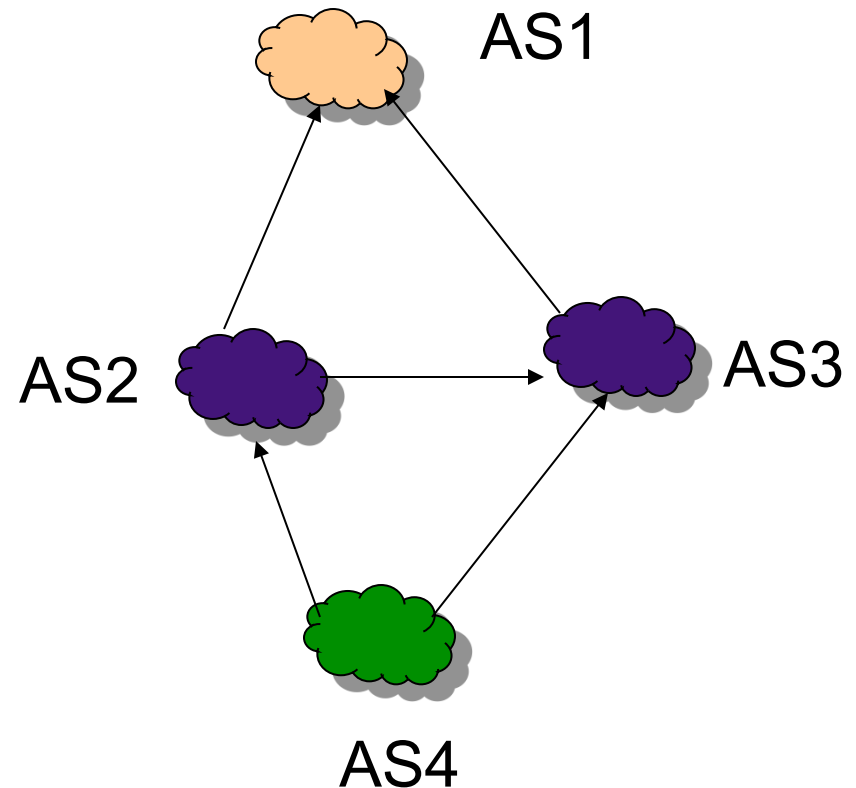
Damping Penalty Function



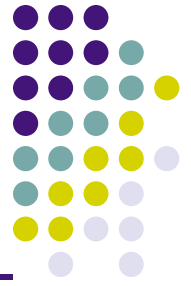
Effect of damping



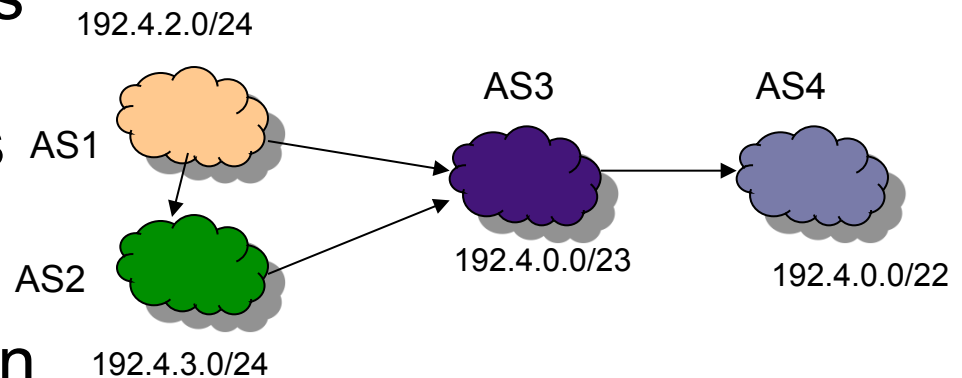
- If Route to AS1 flaps beyond suppression penalty
- AS2 will not advertise routes to AS1 via AS2
- AS4 will then stick to route AS1 via AS3



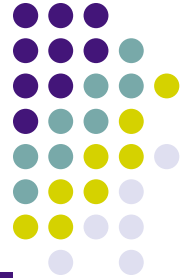
Aggregation can help route flapping



- Specific-route changes can result in flapping
- But aggregated routes may not exhibit route flapping
- Hence aggregation can mask route flapping and reduce instability because it reduces the number of networks visible in the core Internet.



Current Research



1. **A Measurement Study on the Impact of Routing Events on End-to-End Internet Path Performance** F. Wang, Z. M. Mao, J. Wang, L. Gao, R. Bush
ACM SIGCOMM 2006
2. **Rationality and Traffic Attraction: Incentives for Honest Path Announcements in BGP** Sharon Goldberg (Princeton University); Shai Halevi (IBM T. J. Watson Research); Aaron D. Jaggard (Rutgers University); Vijay Ramachandran (Colgate University); Rebecca N. Wright (Rutgers University) ACM SIGCOMM 2008
3. **How secure are inter domain Routing Protocols?** Sharon Goldberg , et.al, SIGCOMM 2010