# CS 31006: Computer Networks – Iinternet Routing
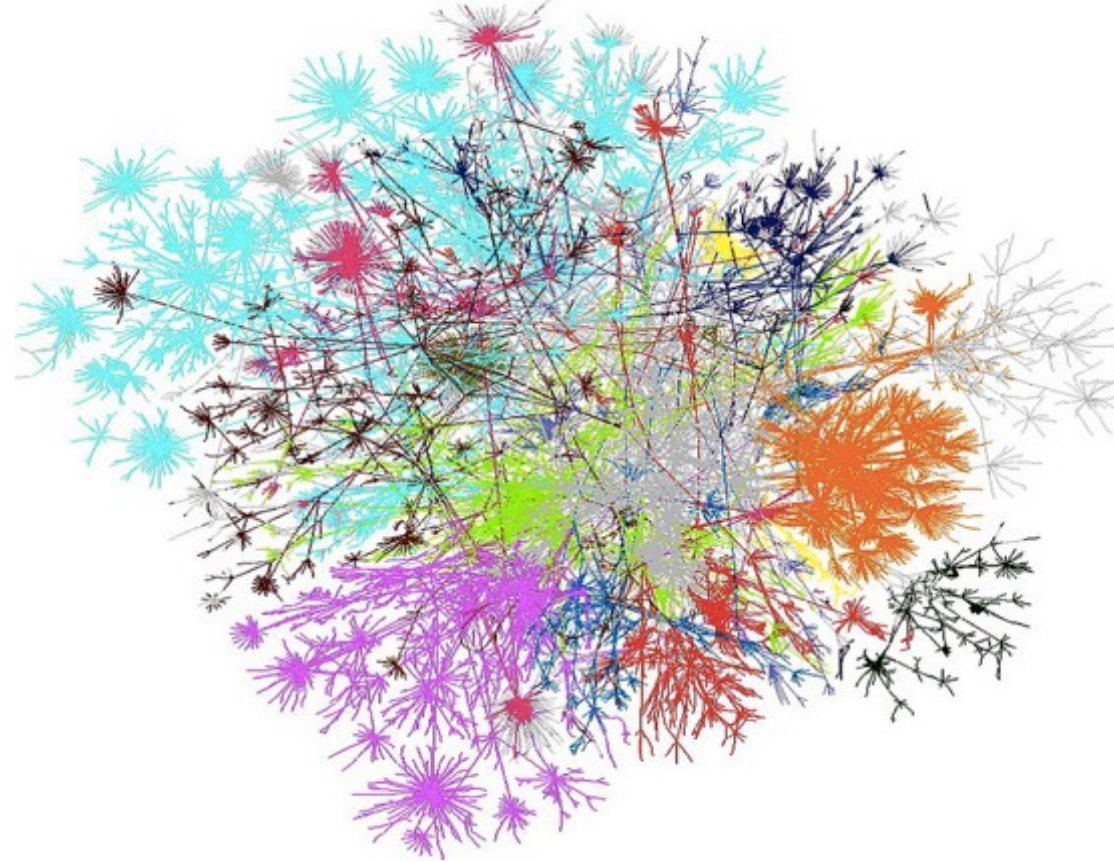
**Department of Computer Science and Engineering**

INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

**Rajat Subhra Chakraborty**
rschakraborty@cse.iitkgp.ac.in

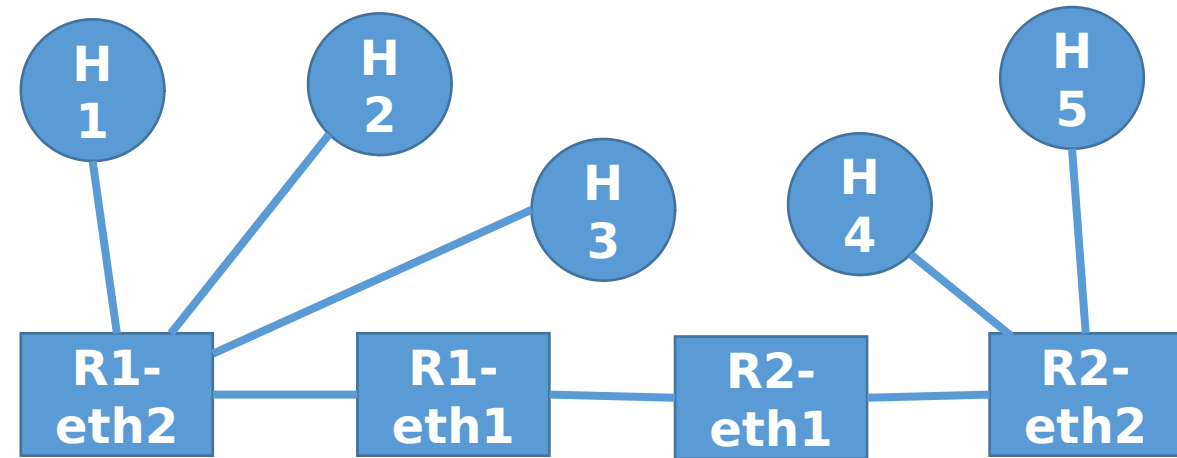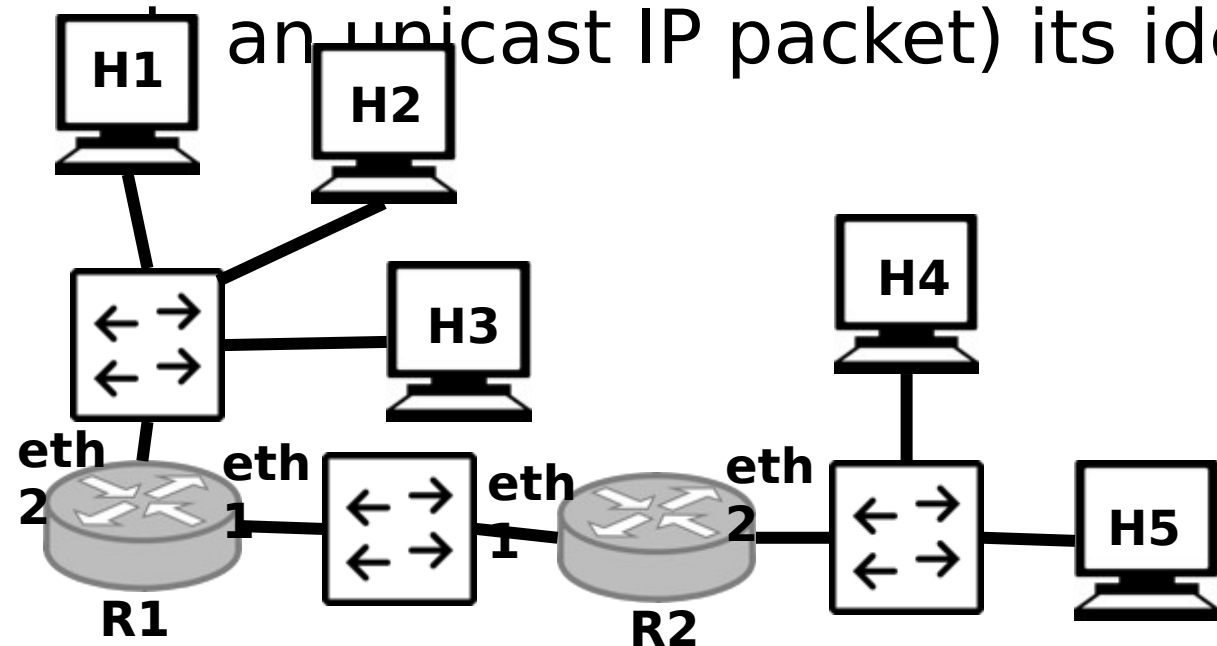**Sandip Chakraborty**
sandipc@cse.iitkgp.ac.in

- The resolution to the counting to infinity problem enforces a maximum cost for a network path (generally 15 in RIP). This limits the diameter of a AS to a maximum of 15 hops.

- High signaling overhead - Periodic broadcasting of the distance vector table can result in increased utilization of the network resources for signaling.

- The algorithm is relatively slow to converge; you require information from all the nodes in the AS.

# Link State Routing

- 1979: The ARPANET routing protocol was replaced by link state routing, as an impact to count-to-infinity problem (convergence become slow)

- The routing protocol – **Open Shortest Path First (OSPF)**

- The protocol is fairly simple
  - Discover neighbors and learn their network addresses
  - Set the distance or cost metric to each of the neighbors
  - Construct a packet telling all it has learned
  - Broadcast this packet – every router periodically learns the **link state** of the network graph
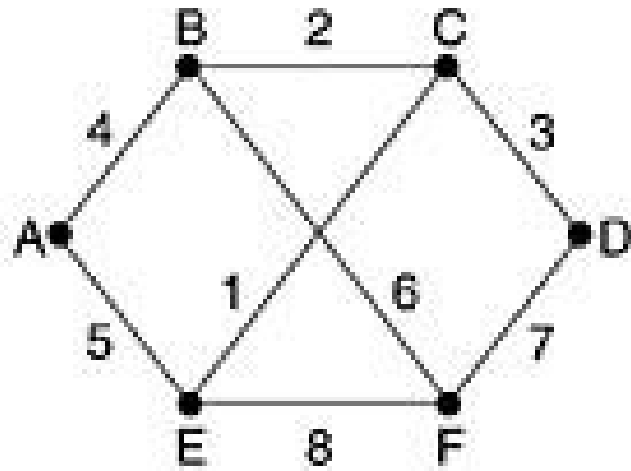  - Compute the shortest path to every other routers

- When a router is booted, it first learn the neighbors – broadcast a HELLO packet on each point to point line **– note the use of broadcast IP address**

- Once a router receives a HELLO message, it sends back (this ~~b~~ an unicast IP packet) its identity (IP address)

# LSR– Setting Link Costs

- Each link is assigned with a link cost or distance (hop count, delay) which is used as the routing metric to find out the shortest path

- A standard approach – inverse of the link bandwidth  - higher capacity paths are better choices (minimize the routing cost)

- Some networks use link delay - computed through a ICMP ECHO packet from the IP layer
  - ICMP – **Internet Control Message Protocol –** a set of message suites for IP layer management functionalities
  - ICMP Echo Request and ICMP Echo Reply

**Image: Computer Networks, Andrew S. Tanenbaum, David J. Wetherall**

- All of the routers must get all of the link state packets quickly and reliably
  - If different routers have different information, the routing inconsistency may occur

- Use flooding to distribute the link state packets to all the

s a sequence number that is
ket sent – **used to identify stale**



**Image:** http://www.danzig.jct.ac.il/tcp-ip-lab/ibm-tutorial/3376c33.html

**Image:** http://cnp3book.info.ucl.ac.be/1st/html/network/network.html

- Once a Link State (LS) packet is received, the sequence number is checked
  - If the sequence number is higher than the last observed LS packet, then it is accepted; otherwise it is discarded

- **What if the sequence number wraps around?**
  - Use a 32 bit sequence number, and 1 LS packet per second – it would take 137 years to wrap around

- Every entry in the router is associated with an **age** – denotes the lifetime of an entry
  - Deletes the old entries from the routing table

- Construct the network graph from the link state packets



**Image:** http://cnp3book.info.ucl.ac.be/1st/html/network/network.html
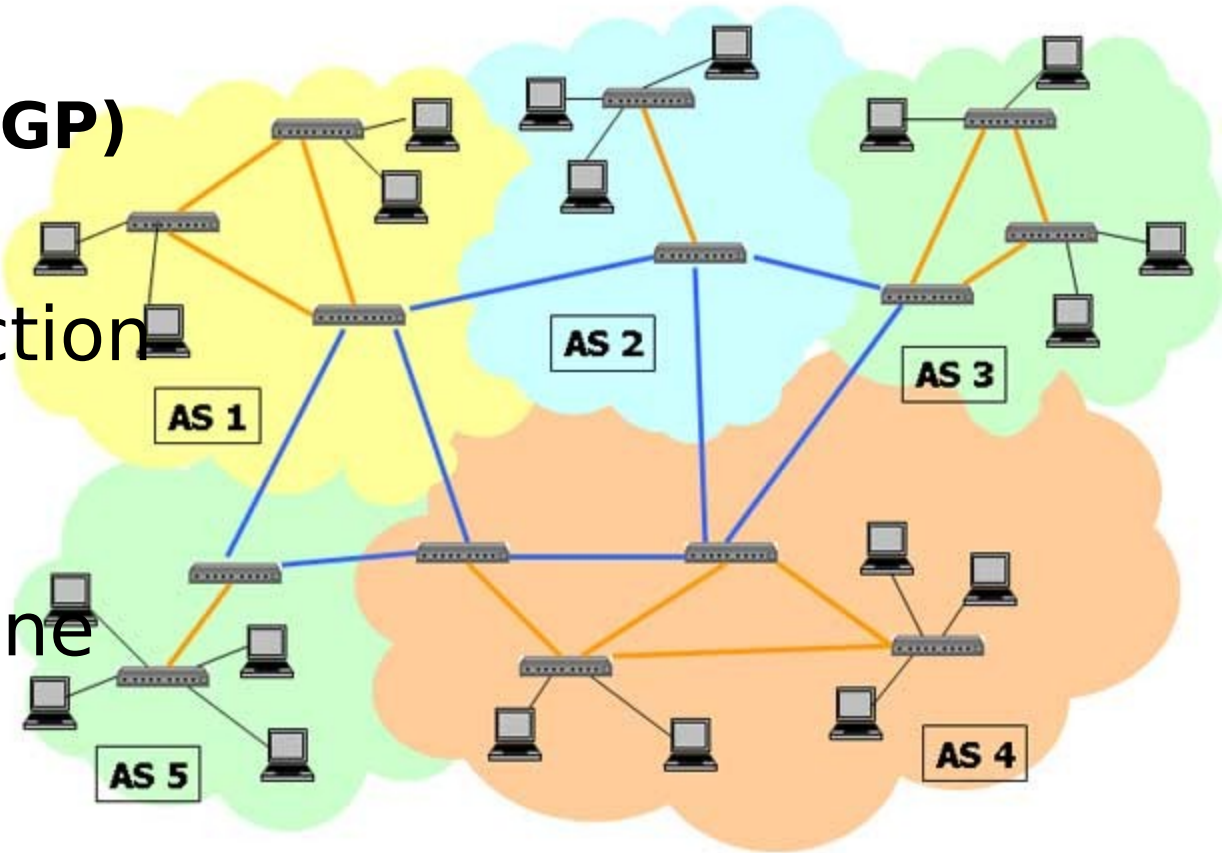
- Link State Messages – **Overhead for the LSR**

- **OSPF Messages**
  - HELLO: Used to establish neighborhood
  - Database Descriptor (DD or DBD): Broadcast the local routing database among neighbors – check consistency of database information among routers
  - Link State Requests (LSR): Explicitly requests for link state information based on the database comparison
  - Link State Updates (LSU): Forward link state information
  - Link State Acknowledgements (LSAck): Acknowledges the receipt of link state information
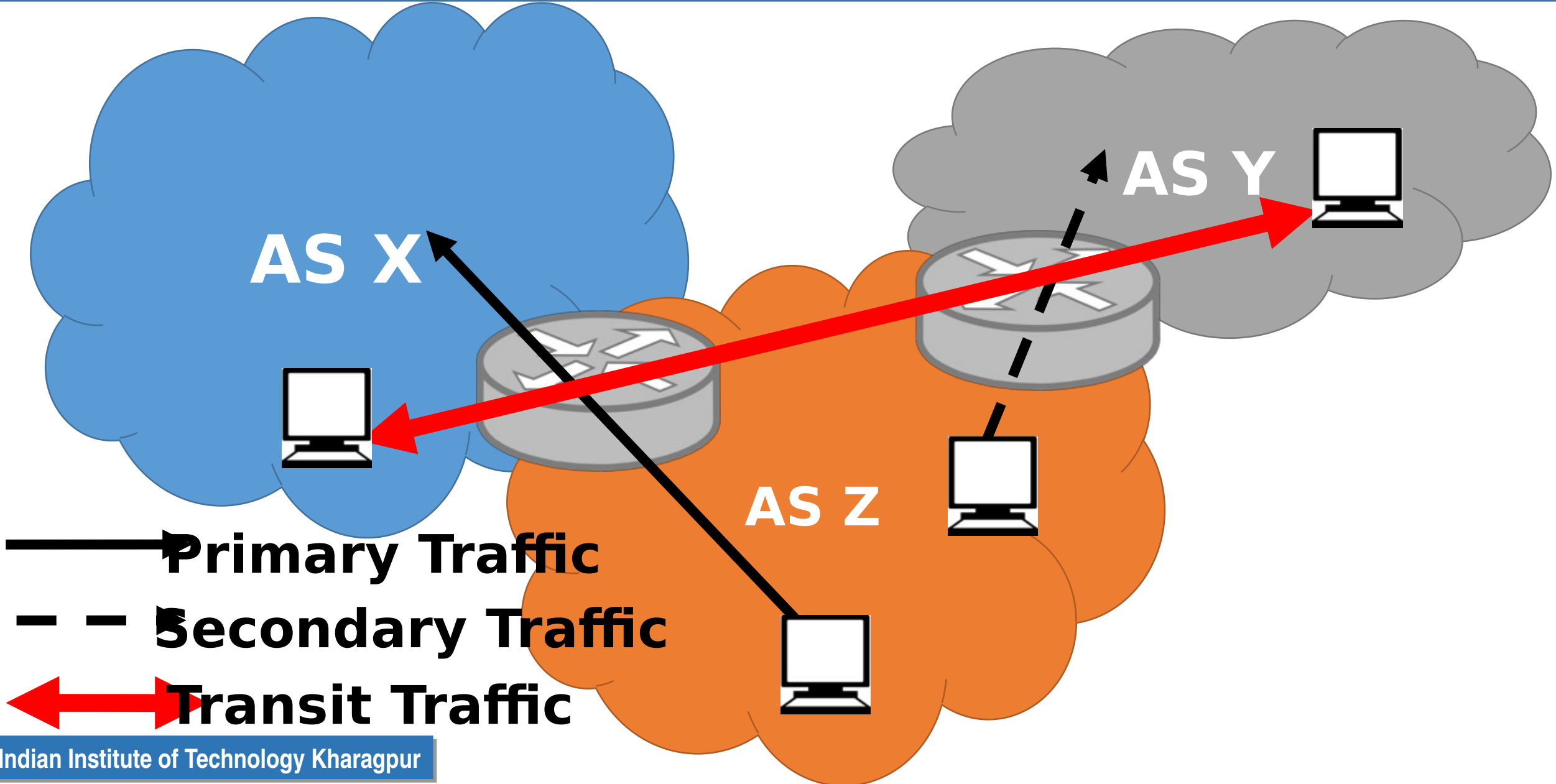
# Inter-domain Routing Protocols

- Routing between two autonomous systems (AS) – we call them as **routing domains**
  - **Exterior Gateway Protocol (EGP)**

- **Multi-homed AS:** has connection to more than one AS

- **Stub AS:** Connected to only one other AS

- **Transit AS:** Provides connection to other AS



**Photo courtesy:** **http://www.web3.lu/**

- Each AS can run its own intra-domain routing protocols, we call them as **interior gateway protocols (IGP)**
  - Open Shortest Path First (OSPF)
  - Routing Information Protocol (RIP)
  - Can even use static routing or a mixed of IGPs at different subnets

- Inter-domain routing problem – **the AS shares *reachability information –*** description of the set of IP addresses that can be reached via a given AS

- **Challenge –** Each AS has to determine its own *routing policies* (can be complex)
  - *Whenever possible, I prefer to send traffic via AS X than via AS Y, but I'll use AS Y if it is the only path; and I never want to carry traffic*
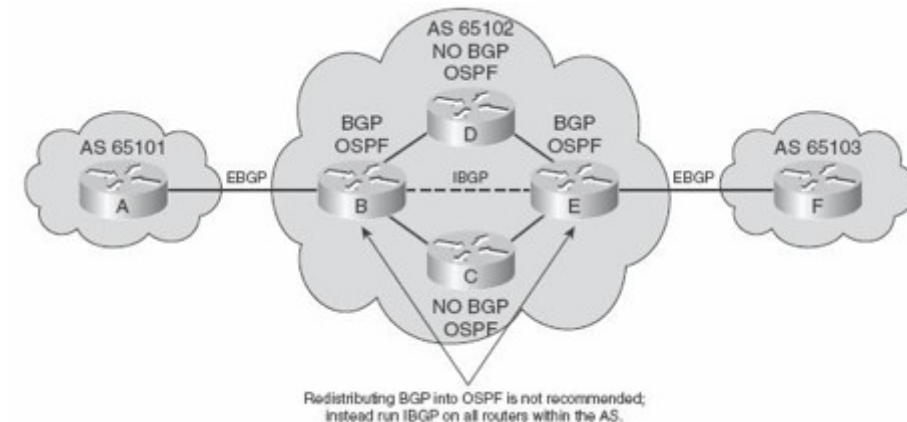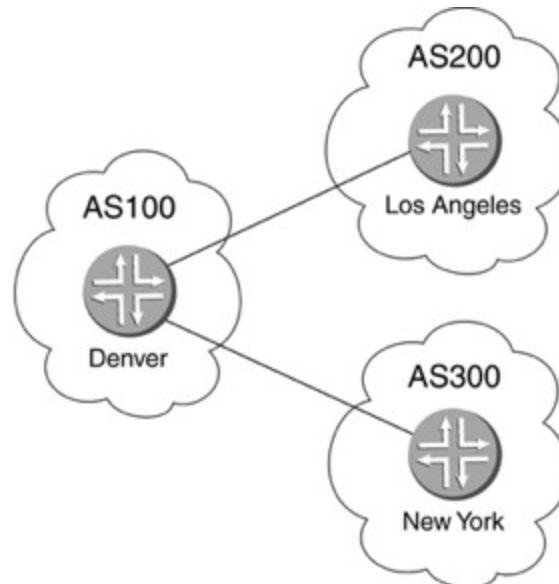
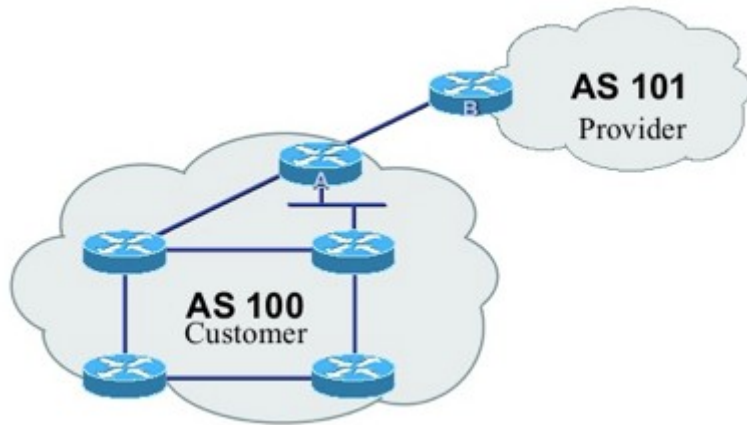# Routing Policies in the Internet



**AS X**

**AS Y**

**AS Z**

— Primary Traffic

- - Secondary Traffic

←→ Transit Traffic

- The initial protocol (called EGP) was designed for specialized topology, such as a tree topology.

- BGP replaces EGP – generalizes the topology structure of the Internet.

- BGP assumes that the Internet is an arbitrary interconnected set of ASs.

- **Local Traffic:** Originates at or terminates on nodes within an AS

- **Stub AS:** Only carry local traffic

- **Multi-homed AS:** Only carry local traffic, refuses to carry transit traffic

- **Transit AS:** Carry both transit and local traffic

- **Best non-looping policy-complaint path**
  - Loop free path through the ASs
  - Complaint with the policies of the various ASs along the path

- **Scaling** – CIDR at the Internet scale may not be scalable – You need to store IP/netmask information for thousands of subnets
  - Define paths by AS numbers, not the IP; Example AS12-AS14-AS76-AS132-AS45-AS61

- **Path Cost** – the autonomous systems are "*autonomous*" – every AS has their own interior routing protocol and own routing metrics – how to define a path metric for a BGP path?
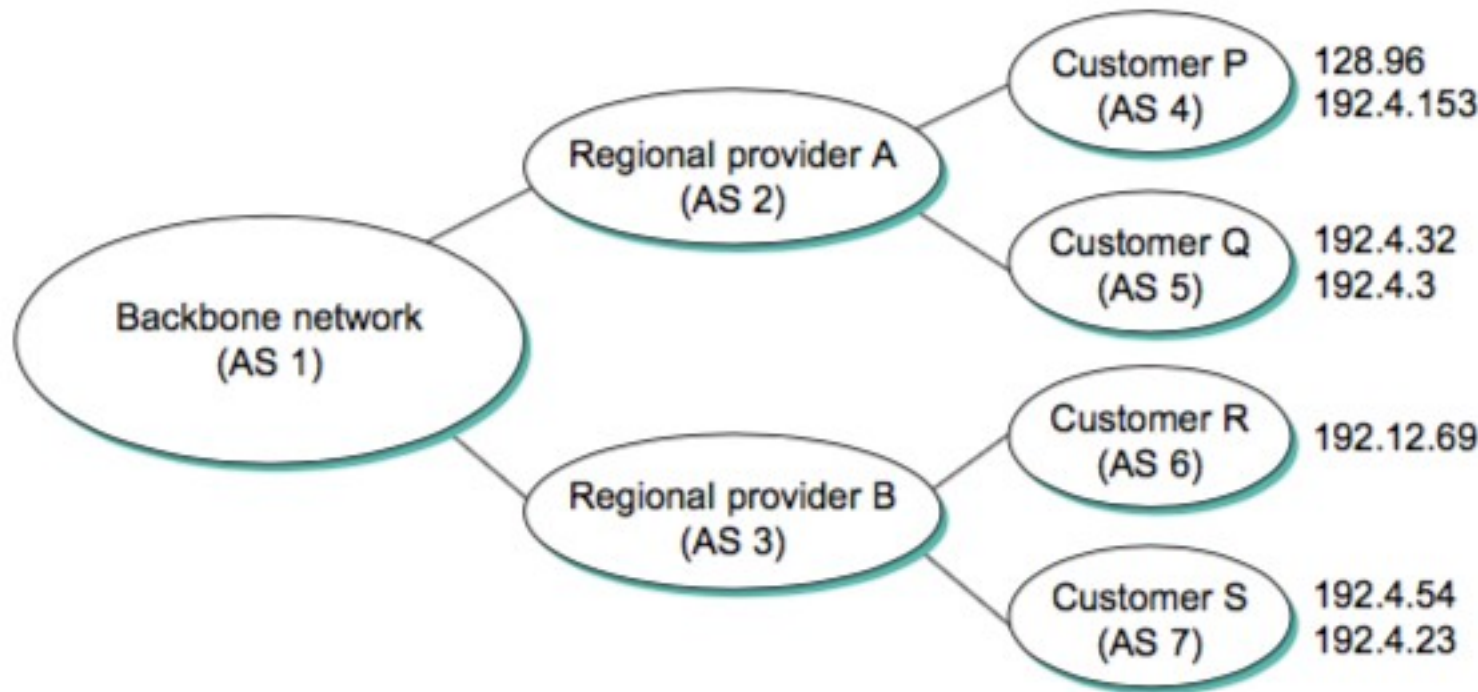  - Advertise only the ***reachability*** information, no path costs are

- AS Reachability information shared by the BGP supported routers (called the BGP peers)

- Shared between BGP peers through BGP UPDATE messages

- An IPv4 prefix (with IPv4 protocol) with the corresponding network IP
  - Example: 110.12, /16

- Used to find out the AS reachability as well as route aggregation through CIDR
  - Combine 110.12/16 and 110.12.8/24 together to form 110.12/16 – You do not need to maintain duplicate paths

- **BGP Speakers** – A router configured with BGP – a spoke-person for the entire AS
  - Advertises the reachability information for this AS

- Once initialized, uses the well known BGP port (TCP port 179) to connect to other configured BGP peers in the Internet, and share the AS reachability information

- BGP speakers advertised the path information with the BGP speakers in the peer ASs.
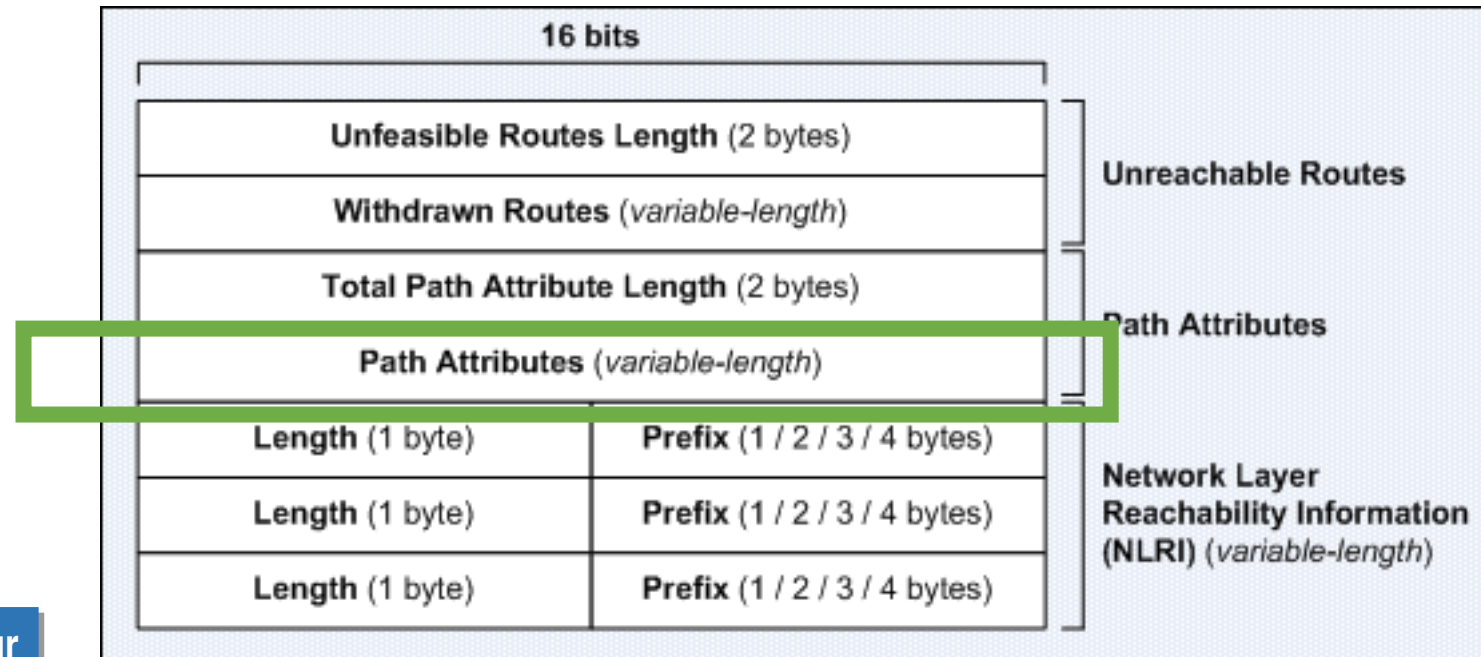
- BGP advertises **complete path** information as enumerated list of ASs to reach a particular network
    - Necessary to enable the sorts of policy decisions
    - Enables routing loops to be readily detected

- AS 2 can advertise NLRI for the subnets given to Customer P and Customer Q



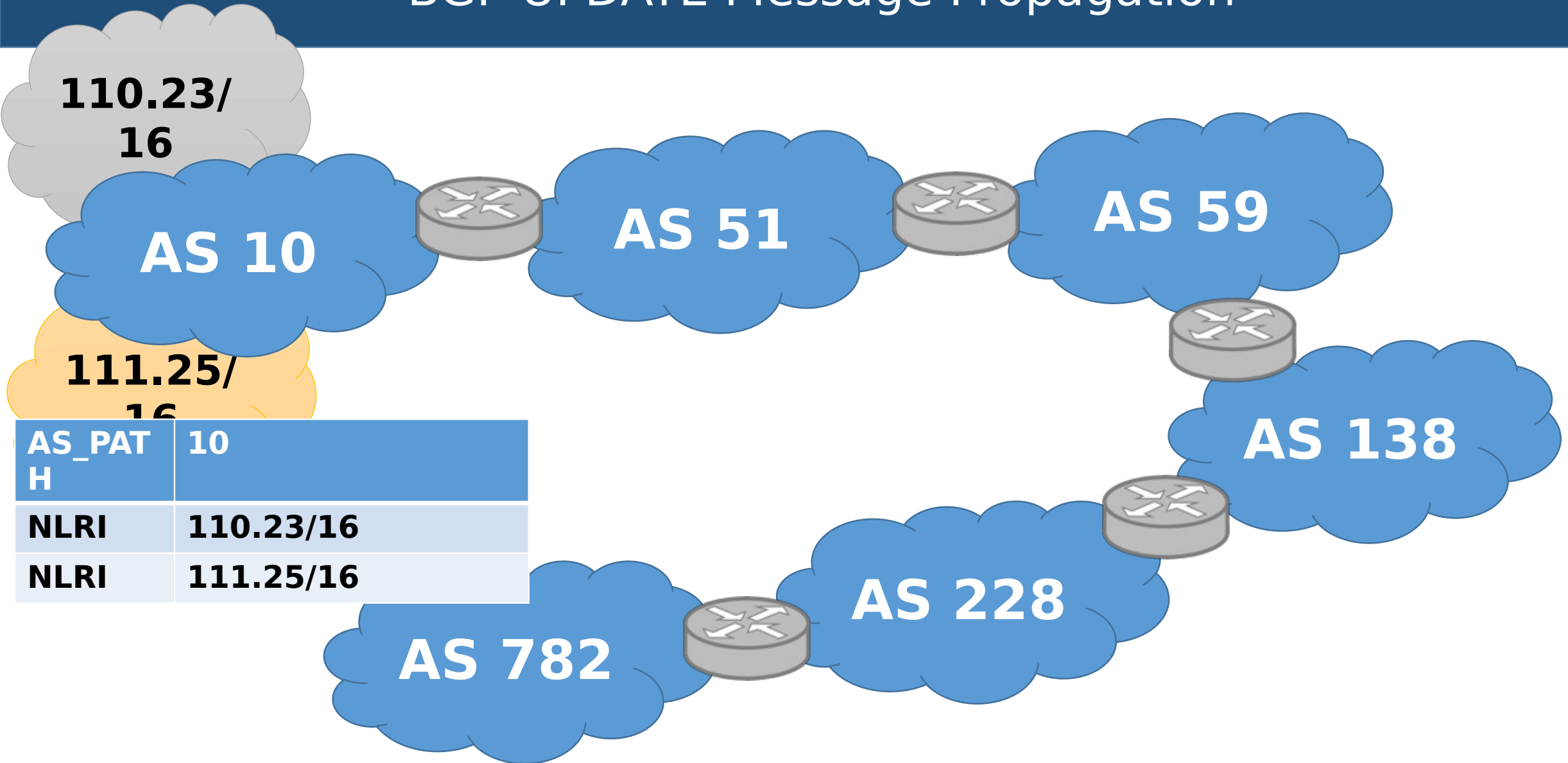- AS 3 can reach to the subnet 128.96 via the AS Path AS3-AS1-AS2-AS4

**Image courtesy: Computer Networks, Larry L Peterson and Bruce S Davie**

- Stores all the paths across various ASs through which a BGP UPDATE message has passed

- Every BGP speaker, when receives an UPDATE message, appends its own AS number and advertise that to the BGP peers
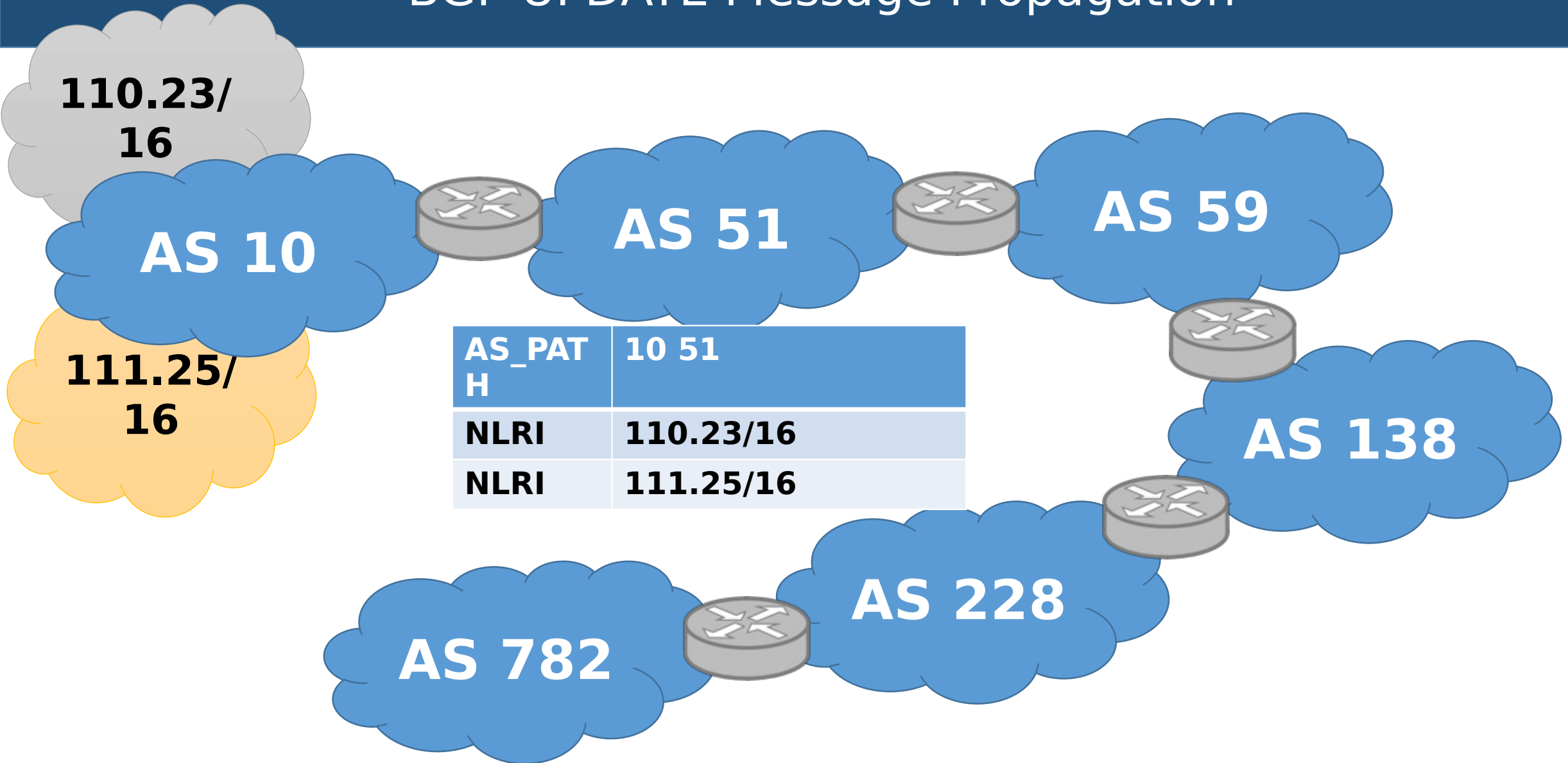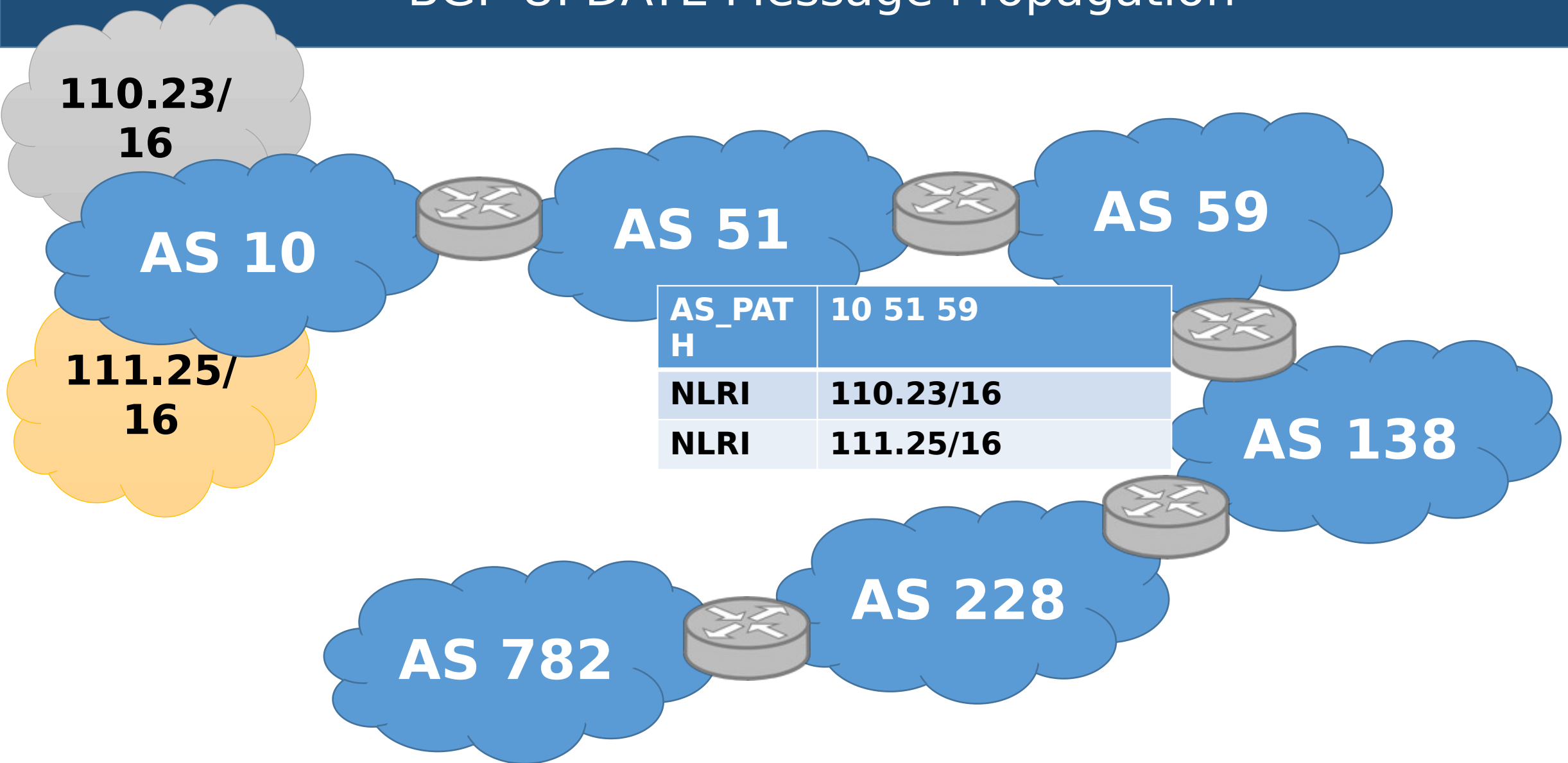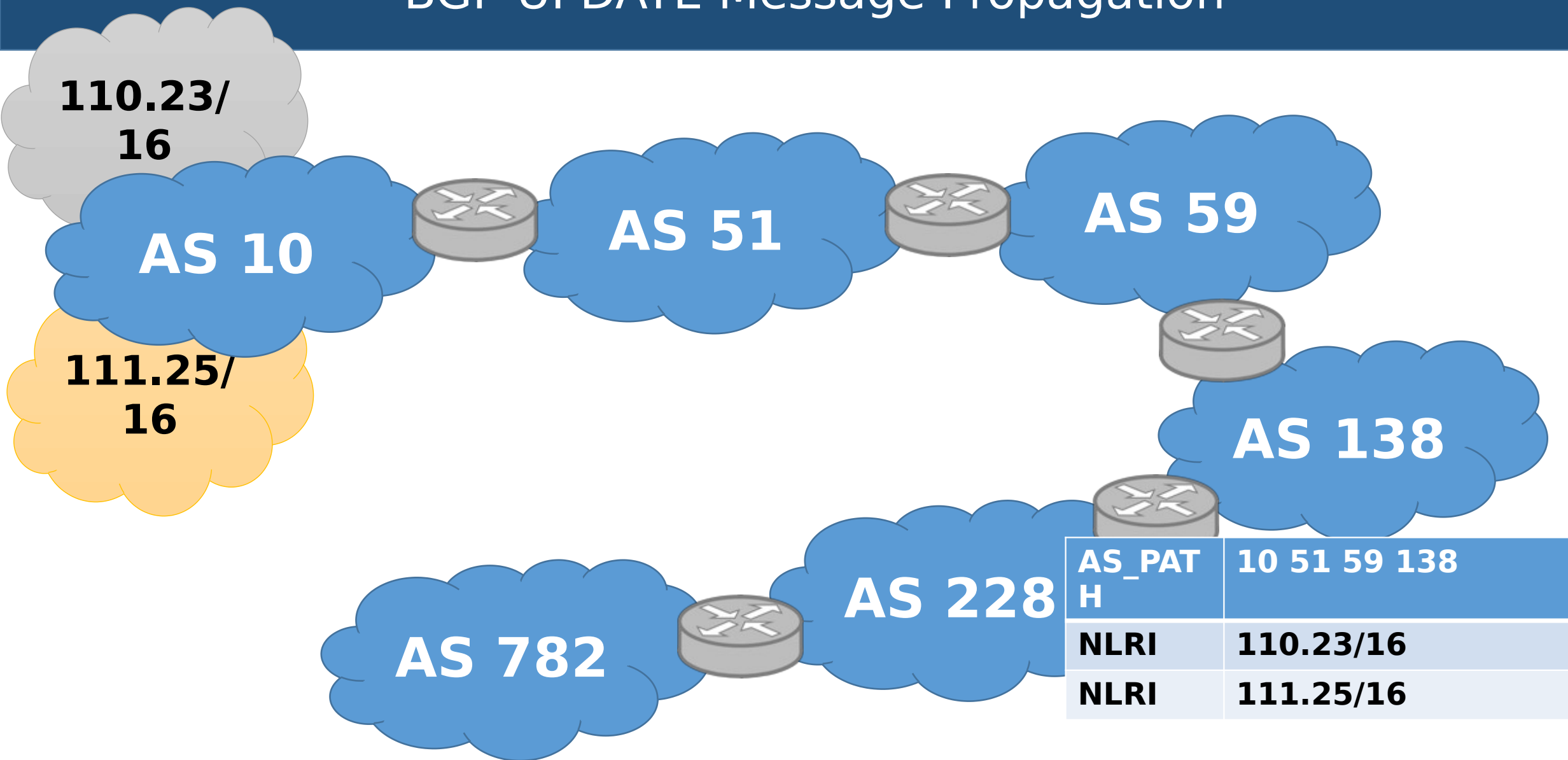
# BGP UPDATE Message Propagation



**110.23/16**

**AS 10**

**AS 51**

**AS 59**

**AS 138**

**111.25/16**

| AS_PATH | 10 |
|---------|-----|
| NLRI | 110.23/16 |
| NLRI | 111.25/16 |

**AS 228**

**AS 782**

# BGP UPDATE Message Propagation

110.23/16

AS 10

AS 51

AS 59

111.25/16

| AS_PATH | 10 51 |
|---------|-------|
| NLRI | 110.23/16 |
| NLRI | 111.25/16 |

AS 138

AS 228

AS 782

Indian Institute of Technology Kharagpur

# BGP UPDATE Message Propagation



| AS_PATH | 10 51 59 138 |
|---|---|
| NLRI | 110.23/16 |
| NLRI | 111.25/16 |

110.23/16

111.25/16

AS 10

AS 51

AS 59

AS 138

AS 228

AS 782

# BGP UPDATE Message Propagation

110.23/16

AS 10

111.25/16

AS 51

AS 59

AS 138

AS 228

AS 782

| AS_PATH | 10 51 59 138 228 |
|---------|------------------|
| NLRI | 110.23/16 |
| NLRI | 111.25/16 |

# BGP UPDATE Message Propagation



| AS_PATH | 10 51 59 138 228 782 |
|---------|----------------------|
| NLRI | 110.23/16 |
| NLRI | 111.25/16 |

110.23/16

111.25/16

AS 10

AS 51

AS 59

AS 138

AS 228

AS 782

- Based on the BGP UPDATE message, a BGP speaker may have multiple paths to a subnet

- The BGP speaker chooses the best one according to its own local policies -> It advertises this route in the next BGP UPDATE message
  - Check the set of rules that are followed for BGP path establishment algorithm: https://www.cisco.com/c/en toc ol-bgp/13753-25.html

**age courtesy: https://www.cisco.com**



**Indian Institute of Technology Kharagpur**

# BGP Looking Glass



**https://stat.ripe.net/widget/looking-glass**