# CS348 Notes
# BGP+IGP
# Video Numbers: 20, 21

### OjMaha

I have prepared these notes by watching the videos from Networks Playlist. The following notes may be asynchronous and irrelevant to what Prof. Vinay teaches in class (cuz I do not pay attention during lectures lol). Further, these notes might not cover *everything* as explained in the video lectures. Consider these to be a supplemental read :). If you find any errors, do notify me so they can be edited.

# Border Gateway Protocol (BGP)

It conducts inter-domain routing b/w ASes.

Customer: The one who pays for service. (Multi-home customer: pays multiple providers for service)

Provider : The one who provides service.

Tier-1 AS: pays nobody. It only acts as a provider. They are widely connected

Tier-2 AS: takes internet service from Tier-1 AS.
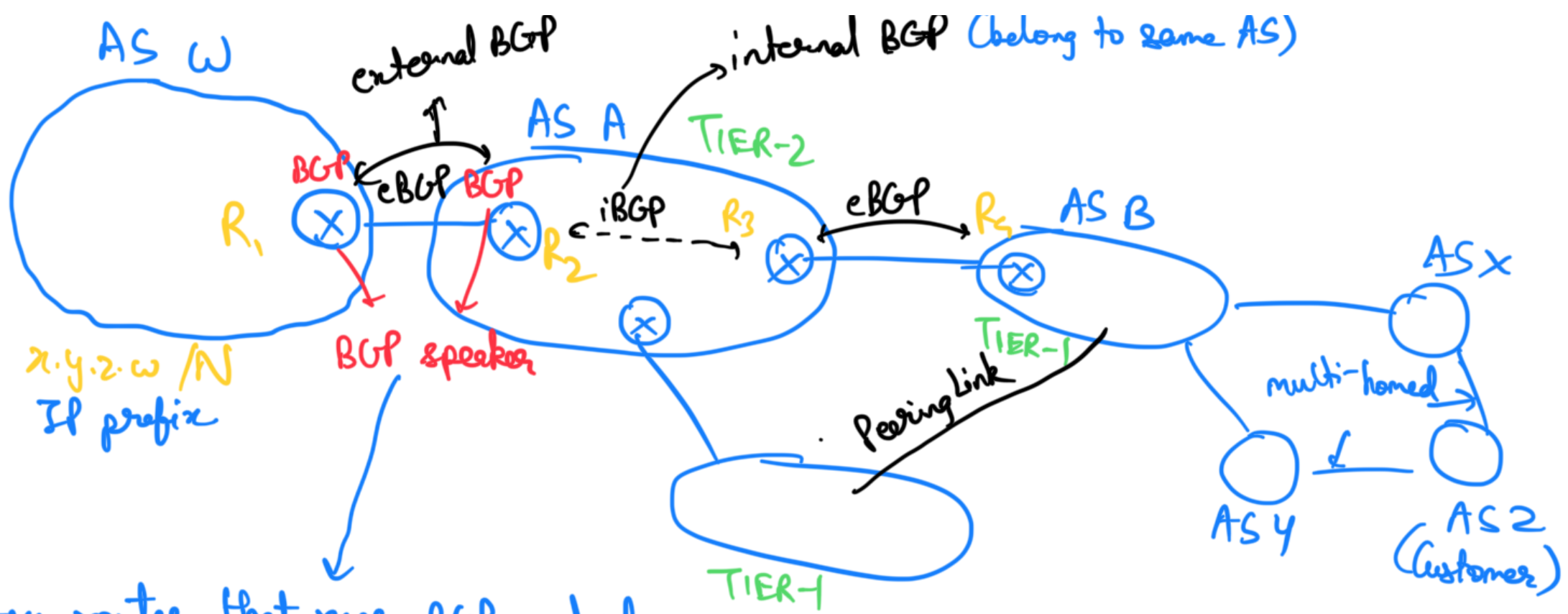
Define Tier-n likewise.

Peering Link: Connection b/w two Tier-1 AS for convenient data sharing to others. No one pays the other here. Since this is unpaid, the bandwidth is low to avoid "cheating" by sending huge amts of data to the other Tier-1. It is "under-provisional".

Service Level Agreement: A QoS guarantee b/w 2 AS abt latency, bandwidth. Note that there cannot be a full end-to-end latency/bandwidth guarantee since AS do not know the whole path. The agreement is just for within own network.

AS W

external BGP

internal BGP (belong to same AS)

AS A    TIER-2

eBGP BGP    iBGP    R₃    eBGP    R₅    AS B

BGP    cBGP    BGP

R₁    ×    ×    R₂    ×    ×    AS×

x.y.z.w /N    ×    ×

IP prefix    BGP speaker    TIER-1    multi-homed

Peering Link    AS Y    f    AS Z

TIER-1    (customer)

Any router that runs BGP protocol.

cBGP  R₁ tells R₂ :    x.y.z.w /N    W    ⏜    ⇒ Advertisement.
                       IP prefix   AS Path   attributes

iBGP

cBGP  R₃ tells R₄ :    x.y.z.w /N    A-W    ⏜
                                            attributed

R₅ tells R₆ .    — —    B-A-W

The AS path is the path from me to reach the IP prefix. It also means that if R₆ gives R₅ an IP pkt with Dst. IP matching prefix:

it will be able to route it along B-A-W.

⚡ B does not necessarily have to tell C that there is a path B-A or B-A-W.
   B & C can decide to keep some of their connections quiet. Because both are
   Tier-1 AS, they do not want the other to send its traffic to it (B).
   (say B)                                      (say C)

Similarly, Z is a small network and doesn't want packets from X to Y
to be routed through it. It doesn't advertise the connection.

BGP does not compel any router to share its ads.

e-BGP: protocol b/w speakers of diff. ASes. They use TCP port 179.

no. of routers in AS.

i-BGP:              same AS. $\exists$ pairwise TCP connections. $O(n^2)$ connections

IGP: Interior Gateway Protocol (Intra-domain routing) ⇒ LSR, DV, etc.

Steps to decide e-BGP paths:

(i) e-BGP speakers learn AS-paths from neighboring routers in other ASes.
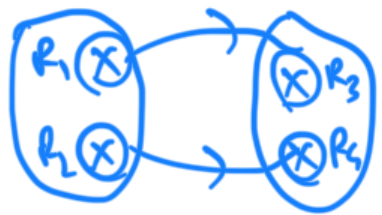(ii) e-BGP nodes share learned info via iBGP with other BGP speakers in own AS

(iii) BGP speakers select routes to various IP prefixes.

(iv) Insert chosen routes into IGP. We do this cuz there may exist nodes & speakers that do not follow BGP but only IGP. Thus this allows all routes to fwd packets to a destination IP.

(v) eBGP speakers can now advertise newly created routes to neighboring ASes.

## BGP attributes:

(i) Local-preferences

Suppose AS2 receives two ads for same prefix. The network administrator determines what the local pref for these routes. Higher local pref ⇒ more preferred path. This pref. gets forwarded.
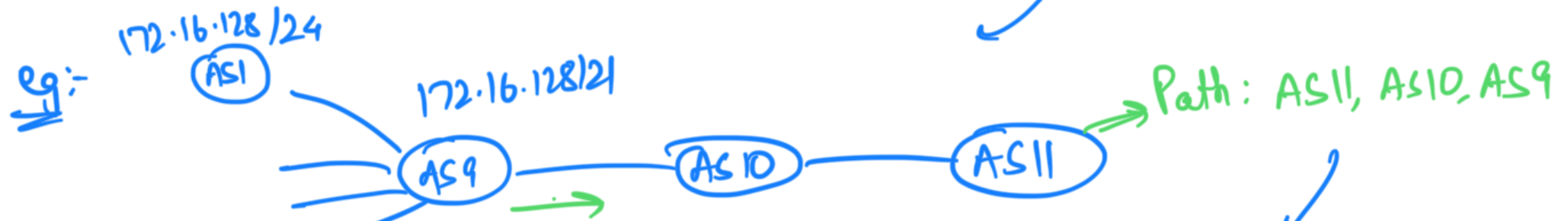
(ii) Multi-exit discriminator (MED)



Suppose $R_1 \to R_3$ & $R_2 \to R_4$ advertise abt. same prefix but diff.
MED: 2000    MED: 1000    ∥ MEDs.
say P

This means; AS1 is telling AS2 that it prefers packets with dest. P to be routed via $R_4 \to R_2$. Lower MED ⇒ higher pref. The value may be the IGP distance

to the next BGP speaker.

**(iii) AS Path:** List of all AS nos all the way to the) dest. with that IP prefix.

whoever first advertises that prefix.

eg:-



172.16.128/24
AS1

172.16.128/21

AS9

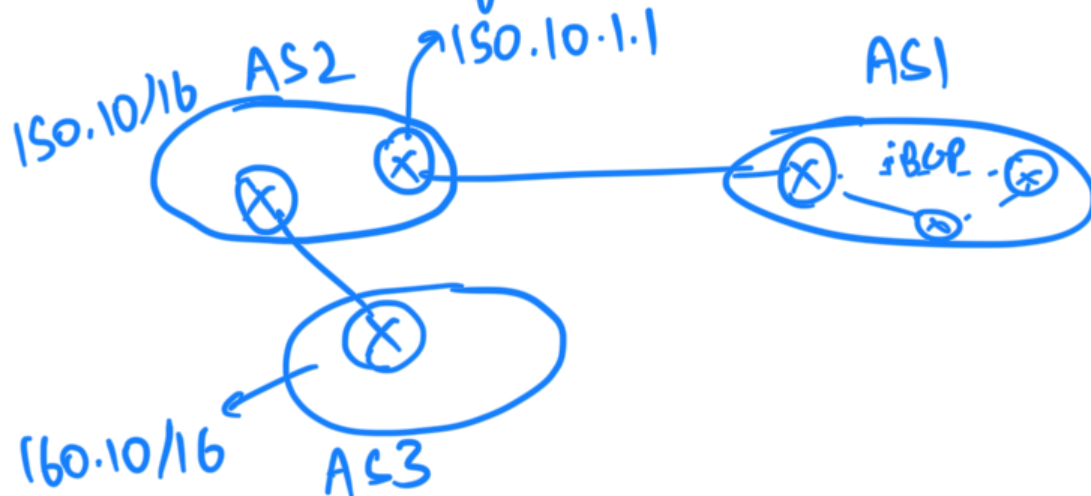AS 10

AS11

→ Path: AS11, AS10, AS9

AS8

172.16.135/24

if the dest. is AS1 but AS9 is the first who advertised, the path will not contain AS1 but end at AS9.

After AS9 receives info; it forwards to AS1.

**(iv) Next-HOP:**

IP address of the first router is the next AS.



150.10/16   AS2   ↗150.10.1.1

150.10/16

AS1

iBGP

info that AS1 has:   Next-HOP

150.10/16   AS2   150.10.1.1
⤷ dest^n.        ↓ ↳Path
160.10/16   AS2,AS3   150.10.1.1

160.10/16   AS3

**Rules to choose Routes:**

Each BGP speaker must decide which ~~route to use among many~~ for the same prefix. (tie-breaks along the way)

→ if no local pref was set

tie-breaker for $a_1 == a_2$
⌃ && $b_1 == b_2$

a) Use largest local-pref.    b) Use shortest AS-path.    c) choose lowest MED path

d) Choose path learned by eBGP over path learned by iBGP.

e) Choose path with lowest IGP metric to next-hop. ⇒ HOT-Potato routing

$A \rightarrow B$
means A advertises to B



Here, $R_2$ chooses to route over $N_2$ rather than $R_1$ in case #AS, and all
(eBGP)                    (iBGP)
other metrics are the same.

Here, $R_3$ recd only iBGP from $R_1$ & $R_2$. Say dist.$(R_3, N_2) >$ dist $(R_3, N_1)$.
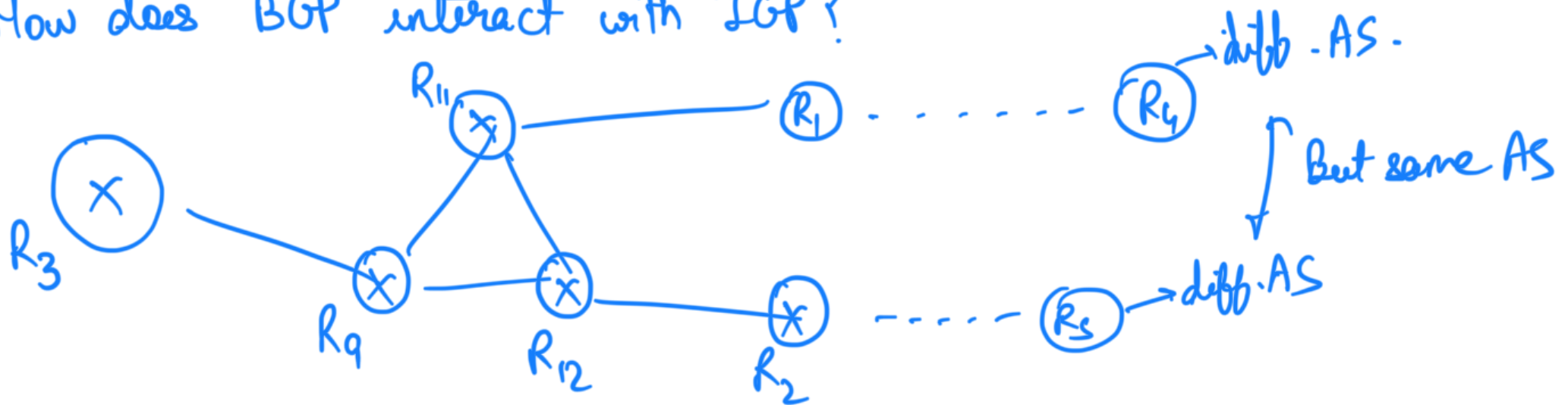Then; $R_3$ chooses to route via $R_2$. (nearest exit) (warm aboo routing)

b) Choose lowest router ID among all BGP speakers who have sent their ads.

∫ router ID: highest IP on router. This is the final tie-breaker.

Why do routers have diff. IP addresses? Because the router may be a part of several networks, each of which requiring a diff. prefix. To satisfy all the cond^n for those nets, it has diff. IPs allowing flexibility.

Watch vid lec 21 till 22 mins for complete example.

How does BGP interact with IGP?



(i) <u>Encapsulation:</u>

Suppose $R_a$, $R_{11}$, $R_{12}$ have no BGP info. It only has info on how

to route pkt to routers of own AS.

in $R_a$ table (IGP)

| Dest. | Next Hop |
|-------|----------|
| $R_1$ | $R_{11}$ |
| $R_2$ | $R_{12}$ |
| $R_3$ | $R_3$ |
| $R_{11}$ | $R_{11}$ |
| $R_{12}$ | $R_{12}$ |

Sidenote: each interface in a router has diff. IP. (Here, we ignore that)

Suppose dest is some router D with IP 142.32.6.21. It is not in routing table of $R_3$. When $R_3$ receives this IP packet, it creates another IP pkt where the recd packet becomes the payload.



payload of new IP pkt

og IP pkt

new IP pkt.

$R_1$ receives this pkt, strips off the outer IP layer and forwards the internal pkt to $R_4$. ($R_1$'s table has info of $R_4$ via eBGP)

**(ii)** <u>Pervasive BGP:</u>

All routers must speak BGP.

Suppose there is a unique exit for 142.32.6 /24

<u>BGP table:</u> Prefix     Exit/Gateway       <u>IGP table:</u> IP addr    Next Hop

142.32.6/24      $R_1$          :         :

true for all

142.32.6/24

→ Recursive Lookup.      (say $R_a$)

On receiving a packet, it finds out dest. prefix; looks for gateway in BGP table ($R_1$) & then looks for next hop in the IGP table ($R_{11}$). $R_{11}$ does a similar thing and so on.

**(iii)** <u>Tagged IGP ;</u>

Internal routers may not be BGP speakers here. But IGP allows add^n of "tags" in the ads.

$R_1$ inserts into its own IGP table :

142.32.6/24    $R_1$ → gateway router.
prefix          tag

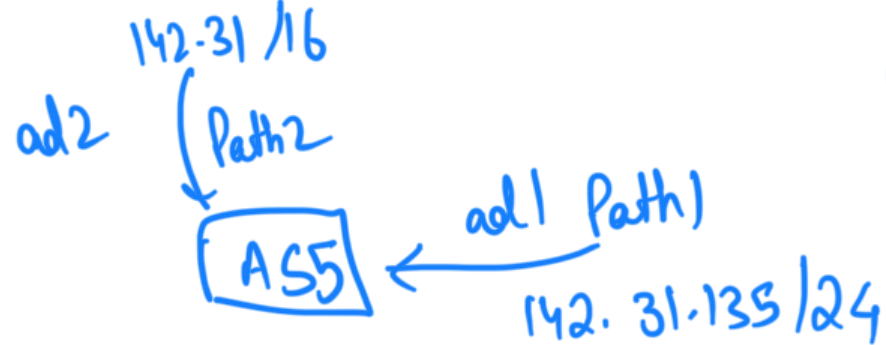This info is propagated to all routers in AS using IGP

Now @ $R_q$:

IGP table:

| Dest | Next | Tag | Cost |
|------|------|-----|------|
| $R_1$ | | | 5 |
| $R_2$ | | | 10 |
| ⋮ | ⋮ | ⋮ | ⋮ |

announced by $R_1$  ↰ 142.32.6/24          $R_1$

announced by $R_2$  ↲ 142.32.6/24          $R_2$

                                                    → hot-potato routing.
Thus, $R_q$ looks up for the prefix & finds closest match. (cost-wise)

# Longest Prefix Matching :

We do supernetting to decrease the size of routing tables.

Now say there is an AS that receives 2 ads in which one IP is the subset of the other. eg:

142.31/16

ad2 ( Path2

AS5 ← ad1 Path1

142.31.135/24

AS5 decides to use Path1 for 142.31.135/24 & Path2 for 142.31/16 excluding 142.31.135/24

Now say dest. IP of a packet rcd @ AS5 is 142.31.135.20

This particular addr matches both the prefixes as mentioned earlier in the routing table. So, **longest prefix** is chosen to route.

why? Because longer prefix match would mean that that particular path is shorter. If we used shorter prefix, there woulda been more number of hops in order to aggregate/supernet into the one with shorter prefix.

Further, longest prefix ensures pkt reaches dest?. If we used shorter prefix. there is a possibility that a pkt keeps hopping b/w 2 routers who each had advertised to each other $IP_1$ and $IP_2$ s.t. either $IP_1 \subseteq IP_2$ or $IP_2 \subseteq IP_1$.

CAM: Content addressable memory. Allows quick look up in the table.