

BGP 基础

BGP (Border Gateway Protocol) 边界网关协议 BGP 知识点

BGP 知识点：

BGP 基础配置，BGP 5 种报文，6 种邻居状态，4 大类 细分 10 种属性，IBGP EBGP (环回口 物理接口) 建立邻居，BGP 认证，fake-as，路由传递原则，IBGP 防环，EBGP 防环，RR 防环，BGP 路由自动聚合，手工聚合 (detail-suppressed，suppress-policy，attribute-policy，origin-policy，as-set)，BGP 5 种 community 属性，BGP 选路，BGP 联盟，路由反射器，BGP 路由过滤，引入，下放默认路由

BGP 概述

边界网关协议 BGP (Border Gateway Protocol) 是一种实现自治系统 AS (Autonomous System) 之间的路由可达，并选择最佳路由的 **高级路径矢量路由协议**。

早期发布的三个版本分别是 BGP-1 (RFC1105)、BGP-2 (RFC1163) 和 BGP-3 (RFC1267)，1994 年开始使用 BGP-4 (RFC1771)，2006 年之后单播 IPv4 网络使用的版本是 BGP-4 (RFC4271)，其他网络 (如 IPv6 等) 使用的版本是 MP-BGP (RFC4760)。

BGP 版本号：如果两端版本号不一致，则协调至相同为止，版本号高的服从版本号低的。假设一端是 BGP3，另一端是 BGP 4，则最后协商的结果为 BGP3

MP-BGP 是对 BGP-4 进行了扩展，来达到在不同网络中应用

的目的，BGP-4 原有的消息机制和路由机制并没有改变。M P-BGP 在 IPv6 单播网络上的应用称为 BGP4+，在 IPv4 组播网络上的应用称为 MBGP (Multicast BGP)。为方便管理规模不断扩大的网络，网络被分成了不同的自治系统。1982 年，外部网关协议 EGP (Exterior Gateway Protocol) 被用于实现在 AS 之间动态交换路由信息。

但是 EGP 设计得比较简单，只发布网络可达的路由信息，而不对路由信息进行优选，同时也没有考虑环路避免等问题，很快就无法满足网络管理的要求。

BGP 是为取代最初的 EGP 而设计的另一种外部网关协议。不同于最初的 EGP，BGP 能够进行路由优选、避免路由环路、更高效率的传递路由和维护大量的路由信息。

虽然 BGP 用于在 AS 之间传递路由信息，但并不是所有 AS 之间传递路由信息都需要运行 BGP。比如在数据中心上行的连入 Internet 的出口上，为了避免 Internet 海量路由对数据中心内部网络的影响，设备采用静态路由代替 BGP 与外部网络通信。

BGP 从多方面保证了网络的安全性、灵活性、稳定性、可靠性和高效性：

- 1、BGP 采用认证和 GTSM 的方式，保证了网络的安全性。
- 2、BGP 提供了丰富的路由策略，能够灵活的进行路由选路。
- 3、BGP 提供了路由聚合和路由衰减功能用于防止路由振荡，有效提高了网络的稳定性。
- 4、BGP 使用 TCP 作为其传输层协议（端口号为 179），并支持 BGP 与 BFD 联动，提高了网络的可靠性。

BGP 按照运行方式分为 EBGp (External/Exterior BGP) 和 I

BGP (Internal/Interior BGP) 。

1、EBGP：

运行于不同 AS 之间的 BGP 称为 EBGP。为了防止 AS 间产生环路，当 BGP 设备接收 EBGP 对等体发送的路由时，会将带有本地 AS 号的路由丢弃。

2、IBGP：

运行于同一 AS 内部的 BGP 称为 IBGP。为了防止 AS 内产生环路，BGP 设备不将从 IBGP 对等体学到的路由通告给其他 IBGP 对等体，并与所有 IBGP 对等体建立全连接。为了解决 IBGP 对等体的连接数量太多的问题，BGP 设计了路由反射器和 BGP 联盟（详情见后面）。

如果在 AS 内一台 BGP 设备收到 EBGP 邻居发送的路由后，需要通过另一台 BGP 设备将该路由传输给其他 AS，此时推荐使用 IBGP。

同步规则 BGP synchronization

当一台路由器从自己的 IBGP 对等体学习到一条 BGP 路由时，它将不能使用该条路由或把这条路由通告给自己的 EBGP 对等体，除非它又从 IGP 协议学习到这条路由，也就是要求 IBGP 路由与 IGP 路由同步。

同步规则主要用于规避 BGP 路由黑洞。

解决方法：

- 1 所有设备都运行 BGP
- 2 IGP 与 BGP 进行引入
- 3 采用 MPLS

华为默认 BGP 同步是被关闭的

Router ID

BGP 的 Router ID 是一个用于标识 BGP 设备的 32 位值，通常是 IPv4 地址的形式，在 BGP 会话建立时发送的 Open 报文

中携带。

对等体之间建立 BGP 会话时，每个 BGP 设备都必须有唯一的 Router ID，否则对等体之间不能建立 BGP 连接。

BGP 的 Router ID 在 BGP 网络中必须是唯一的，可以采用手工配置，也可以让设备自动选取。

缺省情况下，BGP 选择设备上的 Loopback 接口的 IPv4 地址作为 BGP 的 Router ID。如果设备上没有配置 Loopback 接口，系统会选择接口中最大的 IPv4 地址作为 BGP 的 Router ID

BGP 工作原理

BGP 对等体的建立、更新和删除等交互过程主要有 5 种报文、6 种状态机、4 类属性和 5 个原则。

BGP 报文

BGP 对等体间通过以下 5 种报文进行交互，其中 Keepalive 报文为周期性发送，其余报文为触发式发送：

1、Open 报文：

用于建立 BGP 对等体连接。

2、Update 报文：

用于在对等体之间交换路由信息。需要在 BGP 中 network 才会有 Update 报文

3、Notification (通告) 报文：

用于中断 BGP 连接。

4、Keepalive 报文：

用于保持 BGP 连接。

5、Route-refresh (刷新) 报文：

用于在改变路由策略后请求对等体重新发送路由信息。只有支持路由刷新 (Route-refresh) 能力的 BGP 设备会发送和响应此报文。

可以抓取到 route-refresh 报文
<>refresh bgp all import

重置 BGP

<>reset bgp all

修改计时器，默认 60，180

bgp 100

timer keepalive 5 hold 15

BGP 邻居建立状态：

idle:初始状态

connect:BGP 等待 TCP 连接的建立

active:TCP 连接失败，重新建立 TCP 连接

opensent : TCP 建立成功，发送 open 报文

openconfirm:收到正确的 OPEN 报文

established:BGP 邻居建立成功



1、Idle 状态是 BGP 初始状态。

在 Idle 状态下，BGP 拒绝邻居发送的连接请求。只有在收到本设备的 Start 事件后，BGP 才开始尝试和其它 BGP 对等体进行 TCP 连接，并转至 Connect (连接) 状态。Start 事件是由一个操作者配置一个 BGP 过程，或者重置一个已经存在的过程或者路由器软件重置 BGP 过程引起的。

任何状态中收到 Notification (通告) 报文或 TCP 拆链通知等

Error 事件后，BGP 都会转至 Idle 状态。

2、在 Connect (连接) 状态下，BGP 启动连接重传定时器 (Connect Retry)，等待 TCP 完成连接。

如果 TCP 连接成功，那么 BGP 向对等体发送 Open 报文，并转至 OpenSent 状态。

如果 TCP 连接失败，那么 BGP 转至 Active (活跃) 状态。

如果连接重传定时器超时，BGP 仍没有收到 BGP 对等体的响应，那么 BGP 继续尝试和其它 BGP 对等体进行 TCP 连接，停留在 Connect 状态。

3、在 Active 状态下，BGP 总是在试图建立 TCP 连接。

如果 TCP 连接成功，那么 BGP 向对等体发送 Open 报文，关闭重传定时器，并转至 OpenSent 状态。

如果 TCP 连接失败，那么 BGP 停留在 Active 状态。

如果连接重传定时器超时，BGP 仍没有收到 BGP 对等体的响应，那么 BGP 转至 Connect 状态。

4、在 OpenSent 状态下，BGP 等待对等体的 Open 报文，并对收到的 Open 报文中的 AS 号、版本号、认证码等进行检查。如果收到的 Open 报文正确，那么 BGP 发送 Keepalive 报文，并转至 OpenConfirm 状态。

如果发现收到的 Open 报文有错误，那么 BGP 发送 Notification 报文给对等体，并转至 Idle 状态。

5、在 OpenConfirm 状态下，BGP 等待 Keepalive 或 Notification 报文。如果收到 Keepalive 报文，则转至 Established 状态，如果收到 Notification 报文，则转至 Idle 状态。

6、在 Established 状态下，BGP 可以和对等体交换 Update、

Keepalive、Route-refresh 报文和 Notification 报文。

如果收到正确的 Update 或 Keepalive 报文，BGP 就认为对端处于正常运行状态，将保持 BGP 连接。

如果收到错误的 Update 或 Keepalive 报文，BGP 发送 Notification 报文通知对端，并转至 Idle 状态。

Route-refresh 报文不会改变 BGP 状态。

如果收到 Notification 报文，那么 BGP 转至 Idle 状态。

如果收到 TCP 拆链通知，那么 BGP 断开连接，转至 Idle 状态。

常见的 三种状态 Idle，Active，Established

AS	MsgRcvd	MsgSent	OutQ	Up/Down	State	Pre
100	2	2	0	00:00:22	Established	
100	2	2	0	00:00:17	Established	
200	0	1	0	00:00:37	Active	
200	0	0	0	00:00:26	Idle	

BGP 邻居建立不成功的原因

- 1.AS 号或 peer 邻居地址出错；
- 2.BGP 的 router ID 是否有冲突；
- 3.BGP 对等体两端是否均采用环回口创建邻居；
- 4.物理上非直连的 EBGP 邻居是否配置多跳；
- 5.用于创建底层 TCP 的路由是否可达；
- 6.创建 BGP 对等体两端认证配置是否一致；

- 7.BGP 对等体是否配置了 peer x.x.x.x ignore ;
- 8.是否配置了禁止 TCP 端口 179 的 ACL。

ignore 的应用

需要短暂中断邻居会话且该邻居配置量较大时，通过执行命令 peer ignore 可以避免重新配置的工作量。例如，在一段时间内，对端升级或调整链路导致邻居频繁建立连接时，为了避免路由或邻居关系频繁震荡，需要暂时中断 BGP 邻居，则可以在较稳定的一端使用该命令。

使用该命令可以停止与指定对等体（组）之间的会话，并且清除所有相关路由信息。对于一个对等体组，这就意味着大量与对端的会话突然终止。

R1:

```
bgp 100
```

```
peer 192.168.12.2 ignore
```

邻居关系会变成 idle

Peer	V	AS	MsgRcvd	MsgSent
OutQ	Up/Down	State	PrefRcv	
192.168.12.2		4	100	0
0	0 00:00:03	Idle(Admin)		

BGP 对等体之间的交互原则

BGP 设备将最优路由加入 BGP 路由表，形成 BGP 路由。

BGP 设备与对等体建立邻居关系后，采取以下交互原则：

- 1、从 IBGP 对等体获得的 BGP 路由，BGP 设备只发布给它的 EBGP 对等体。

- 2、从EBGP对等体获得的BGP路由，BGP设备发布给它所有EBGP和IBGP对等体。
- 3、当存在多条到达同一目的地址的有效路由时，BGP设备只将最优路由发布给对等体。
- 4、路由更新时，BGP设备只发送更新的BGP路由。
- 5、所有对等体发送的路由，BGP设备都会接收。

BGP 属性

4 类属性， 10 种

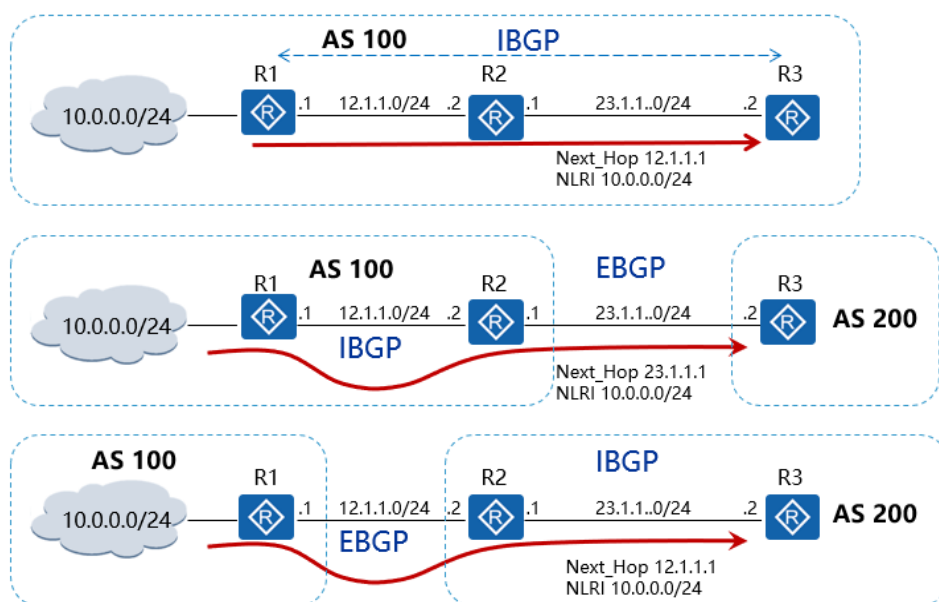
属性名称	类别
ORIGIN	公认必须遵循
AS_PATH	公认必须遵循
NEXT_HOP	公认必须遵循
LOCAL_PREF	公认可选
ATOMIC_AGGREGATE	公认可选
AGGREGATOR	可选过渡
COMMUNITY	可选过渡
MULTI_EXIT_DISC (MED)	可选非过渡
ORIGINATOR_ID	可选非过渡
CLUSTER_LIST	可选非过渡

Atomic_aggregate 是一个公认可选属性，它只相当于一种预警标识，而并不承载任何信息。当路由器收到一条 BGP 路由更新时，发现该条路由携带 Atomic_aggregate 属性时，它便知道这条路由可能出现了路径属性的丢失，此时该路由器把这条路由通告给其他对等体时，需保留路由的 Atomic_aggregate 属性。另外，收到该路由更新的路由器不通将这条路由再度明细化。

Aggregator 是一个可选过渡属性，用于标记路由汇总行为发生在哪个 AS 及哪台 BGP 路由器上。

Next_Hop

Next_Hop 属性记录了路由的下一跳信息。BGP 的下一跳属性和 IGP 的有所不同，不一定是邻居设备的 IP 地址。通常情况下，Next_Hop 属性遵循下面的规则：



1 BGP Speaker 将本地始发路由发布给 IBGP 对等体时，会把该路由信息的下一跳属性设置为本地与对端建立 BGP 邻居关系的接口地址。

2 BGP Speaker 在向 EBGP 对等体发布某条路由时，会把该路由信息的下一跳属性设置为本地与对端建立 BGP 邻居关系的接口地址。

从 EBGP 邻居收到的路由，传给 EBGP 邻居时一定会改变下一跳吗？

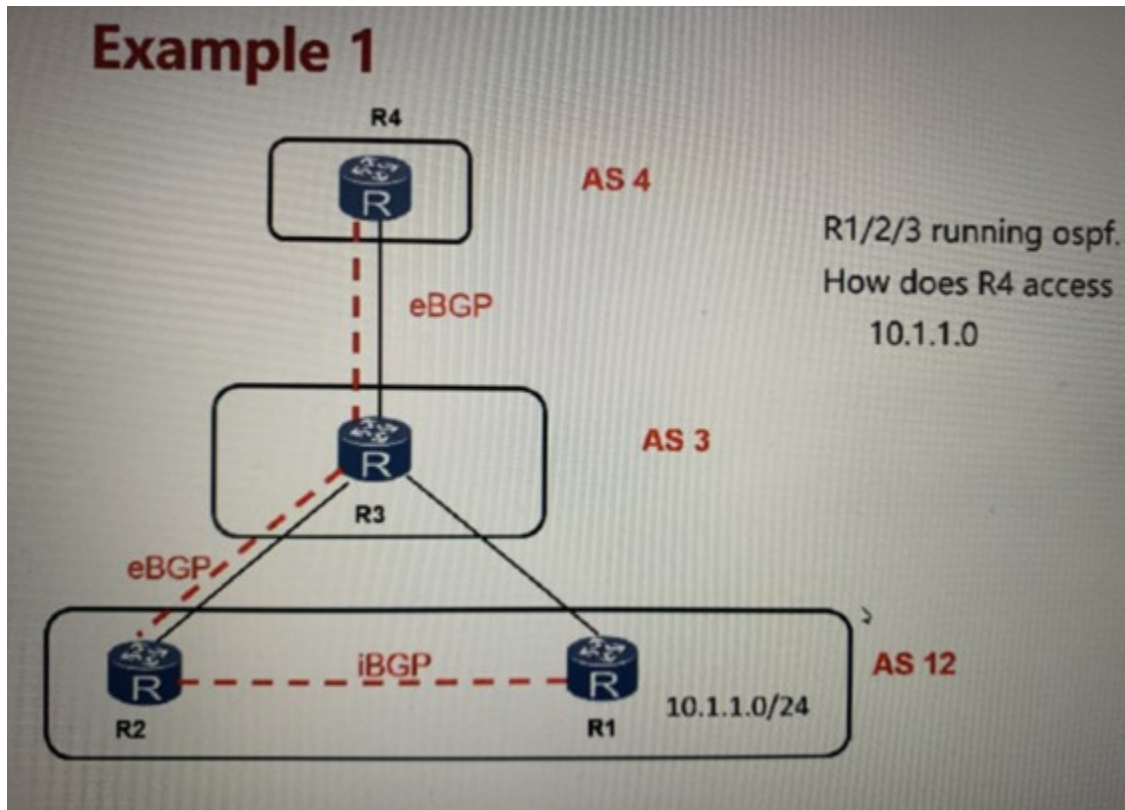
默认情况下是会改变的，但是可以通过配置命令使下一跳不改变，或者在一些特殊场景下，发出路由给 EBGP 邻居时下一

跳是不改变的，主要是为了防止出现次优路径的问题。

R2：

```
bgp 12
```

```
peer 192.168.23.3 next-hop-invariable
```



peer next-hop-invariable 命令配置不同 AS 域的 PE 向 EBGP 对等体发布路由时不改变下一跳；向 IBGP 对等体发布引入的 IGP 路由时使用 IGP 路由的下一跳地址。保证对端 PE 可以在流量传输时迭代到通往本端 PE 的 BGP LSP。

3 BGP Speaker 在向 IBGP 对等体发布从 EBGP 对等体学来的路由时，并不改变该路由信息的下一跳属性。要设置 peer 1.1.1.1 next-hop-local, 改变下一跳

BGP 防环机制

IBGP 防环：

路由器从它的一个 BGP 对等体那里接收到的路由条目不会将该路由器再传递给其它 IBGP 对等体，这个原则被称为 BGP 水平分割

路由反射器的防环：Originator_id, Cluster_list

Originator_id 可选非过渡属性，由 RR 产生，封装在 Update 消息中，使用 router-id 值标识路由的始发者，用于防止集群内路由环路。

Cluster_list 可选非过渡属性，记录路由经过的每个集群的 Cluster_id，用来在集群间避免环路。

EBGP 防环：

当路由器从 EBGP 邻居收到 BGP 路由时，如果该路由的 AS_Path 中包含了自己的 AS 编号，则该路由将会直接丢弃。

=====

BGP 是目前 Internet 骨干网上运行的核心路由协议，也是部署最广泛的路由协议之一。在过去的几十年里，Internet 的发展日新月异，新兴应用的不断涌现，对 Internet 网络的可靠性、扩展性提出了更高的要求。作为整个 Internet 稳定运行的基础，BGP 为了适应 Internet 的发展趋势，也推出了许多高级特性。

fake-as

RTB:

bgp 2000

peer 1.1.1.1 fake-as 200

=====

配置 BGP 负载分担

在大型网路中，到达同一目的地通常会存在多条有效路由，但是 BGP 只将最优路由发布给对等体，这一特点往往会造成很多流量负载不均衡的情况。通过配置 BGP 负载分担，可以流量负载均衡，减少网络拥塞。

一般情况下，只有“BGP 选择路由的策略”所描述的前 8 个属性完全相同，BGP 路由之间才能相互等价，实现 BGP 的负载分担。

```
bgp 100
maximum load-balancing 2
```

配置完成后，查看全局 ip 路由表
同一个 BGP 路由条目，是有两个下一条

44.44.44.0/24	IBGP	255	0	RD	1.1.1.1	GigabitEthernet
/0/0						
	IBGP	255	0	RD	2.2.2.2	GigabitEthernet
/0/1						
55.55.55.0/24	IBGP	255	0	RD	1.1.1.1	GigabitEthernet
/0/0						
	IBGP	255	0	RD	2.2.2.2	GigabitEthernet
/0/1						

在 BGP 路由表中，还是只优选一个条目

connect-interface

用来指定发送 BGP 报文的源接口，并可指定发起连接时使用的源地址。

在两台设备通过多链路建立多个对等体时，使用 peer connect-interface 命令。

使用环回口建立邻居时，把更新源接口由物理接口改变环回口

next-hop-local

命令一般在 ASBR 上配置。当设备通过 EBGP 邻居学到路由再转发给其他 IBGP 邻居时，默认不修改下一跳，但其 EBGP 邻居发来的路由的下一跳都是其 EBGP 邻居的 Peer 地址，本端对等体所属 AS 域内的 IBGP 邻居收到这样的路由后，由于下一跳不可达导致路由无法活跃。因此，需要在 ASBR 上对 IBGP 邻居配置 `peer next-hop-local` 命令，使得发给 IBGP 邻居的路由的下一跳是其自身的地址，IBGP 邻居收到这样的路由后（由于域内都配置了 IGP）发现下一跳可达，路由即为活跃路由。

ebgp-max-hop

通常情况下，EBGP 对等体之间必须具有直连的物理链路，如果不满足这一要求，则必须使用 `peer ebgp-max-hop` 命令允许它们之间经过多跳建立 TCP 连接。BGP 使用 Loopback 口建立 EBGP 邻居时，必须配置命令 `peer ebgp-max-hop`（其中 `hop-count ≥ 2`），否则邻居无法建立。

路由反射器与联盟对比

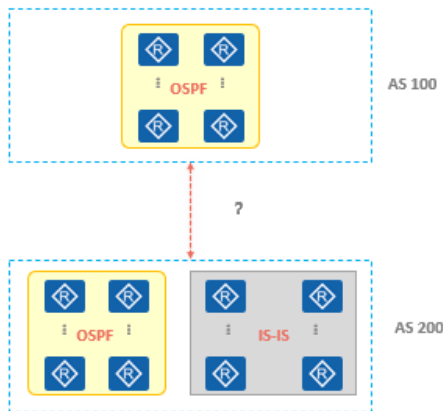
表 6-8

路由反射器和联盟对比

路由反射器	联盟
不需要更改现有的网络拓扑，兼容性好	需要改变逻辑拓扑
配置方便，只需要在反射器上配置	所有设备需要重新配置
集群之间需要全互联	联盟的 AS 之间由特殊的互联
适用于中大型网络	适用于特大规模网络

- 为方便管理规模不断扩大的网络，网络被分成了不同的 AS (Autonomous System，自治系统)。早期，EGP (Exterior Gateway Protocol，外部网关协议) 被用于实现在 AS 之间动态交换路由信息。但是 EGP 设计得比较简单，只发布网络可达的路由信息，而不对路由信息进行优选，同时也没有考虑环路避免等问题，很快就无法满足网络管理的要求。
- BGP 是为取代最初的 EGP 而设计的另一种外部网关协议。不同于最初的 EGP，BGP 能够进行路由优选、避免路由环路、更高效率的传递路由和维护大量的路由信息。
- 本章节将介绍 BGP 的基本概念。

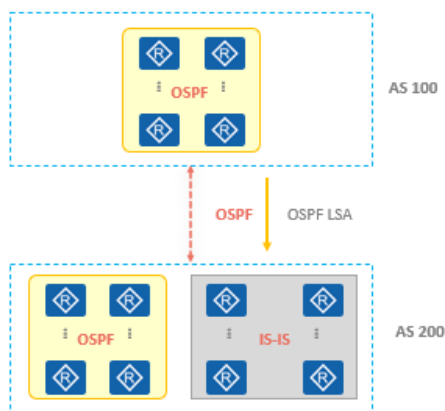
AS



- OSPF、IS-IS等IGP路由协议在组织机构网络内部广泛应用，随着网络规模扩大，网络中路由数量不断增长，IGP已无法管理大规模网络，AS的概念由此诞生。
- AS指的是在同一个组织管理下，使用统一选路策略的设备集合。
- 不同AS通过AS号区分，AS号存在16bit、32bit两种表示方式。IANA负责AS号的分发。
- 当不同AS之间需要进行通信时，在AS之间应使用何种路由协议进行路由的传递？

- IANA (Internet Assigned Numbers Authority ，因特网地址分配组织) ： IAB (Internet Architecture Board ，因特网体系委员会) 的下设组织。IANA 授权 NIC (Network Information Center ，网络信息中心) 和其他组织负责 IP 地址和域名分配，同时，IANA 负责维护 TCP/IP 协议族所采用的协议标识符数据库，包括自治系统号。
- 在长度为 16bit 的 AS 号表示方式中：64512-65534 为私有 AS 号，在长度为 32bit 的 AS 号表示方式中：4200000000-4294967294 为私有 AS 号。

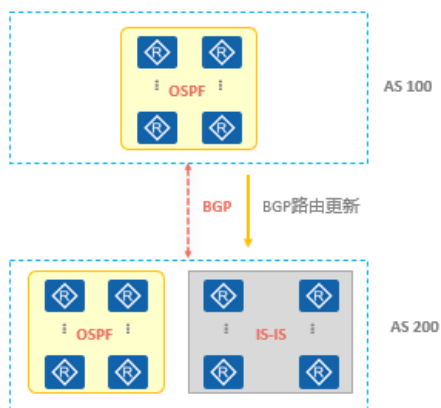
使用IGP传递路由



- AS之间需要直连链路，或通过VPN协议构造逻辑直连（例如GRE Tunnel）进行邻居建立。
- AS之间可能是不同的机构、公司，相互之间无法完全信任，使用IGP可能存在暴露AS内部的网络信息的风险。
- 整个网络规模扩大，路由数量进一步增加，路由表规模变大，路由收敛变慢，设备性能消耗加大。

- VPN (virtual private network , 虚拟专用网) ：使用虚拟专业网络技术可以从逻辑上建立一个直接连接的网络。

使用BGP传递路由



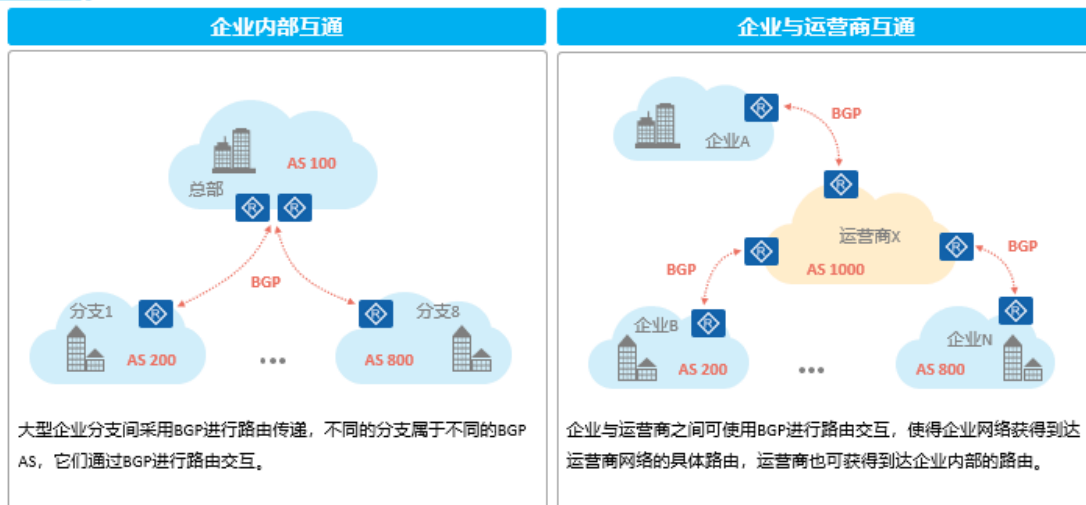
- 为此在AS之间专门使用BGP (Border Gateway Protocol, 边界网关协议) 协议进行路由传递，相较于传统的IGP协议：
 - BGP基于TCP，只要能够建立TCP连接即可建立BGP。
 - 只传递路由信息，不会暴露AS内的拓扑信息。
 - 触发式更新，而不是进行周期性更新。

BGP发展历史



- 目前关于 BGP-4 最新的 RFC 是 4271, 相比较于 RFC1771, 对于一些细节进行了进一步说明, 如事件、状态机以及 BGP 路由决策流程等。

BGP在企业中的应用



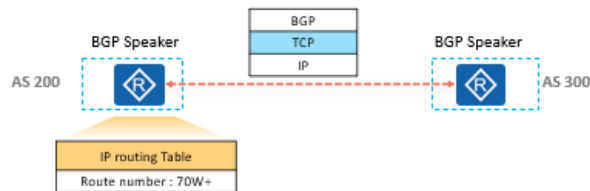


BGP概述

- BGP是一种实现自治系统AS之间的路由可达，并选择最佳路由的矢量性协议。早期发布的三个版本分别是BGP-1（RFC1105）、BGP-2（RFC1163）和BGP-3（RFC1267），1994年开始使用BGP-4（RFC1771），2006年之后单播IPv4网络使用的版本是BGP-4（RFC4271），其他网络（如IPv6等）使用的版本是MP-BGP（RFC4760）。
- BGP的特点：
 - BGP使用TCP作为其传输层协议（端口号为179），使用触发式路由更新，而不是周期性路由更新。
 - BGP能够承载大批量的路由信息，能够支撑大规模网络。
 - BGP提供了丰富的路由策略，能够灵活的进行路由选路，并能指导对等体按策略发布路由。
 - BGP能够支撑MPLS/VPN的应用，传递客户VPN路由。
 - BGP提供了路由聚合和路由衰减功能用于防止路由振荡，通过这两项功能有效地提高了网络稳定性。



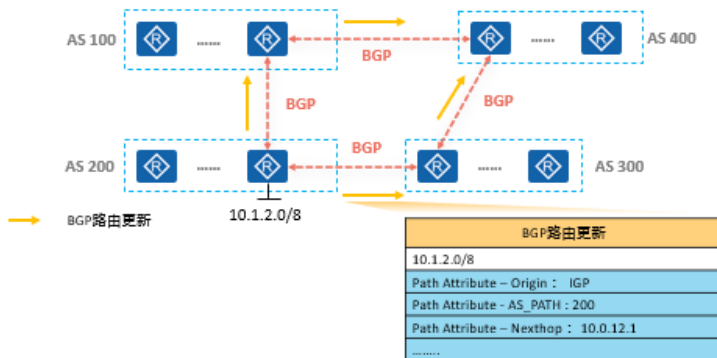
BGP特征 (1)



- BGP使用TCP为传输层协议，TCP端口号179。路由器之间的BGP会话基于TCP连接而建立。
- 运行BGP的路由器被称为BGP发言者（BGP Speaker），或BGP路由器。
- 两个建立BGP会话的路由器互为对等体（Peer），BGP对等体之间交换BGP路由表。
- BGP路由器只发送增量的BGP路由更新，或进行触发式更新（不会周期性更新）。
- BGP能够承载大批量的路由前缀，可在大规模网络中应用。



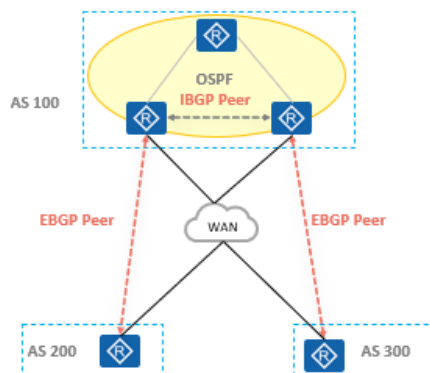
BGP特征 (2)



- BGP通常被称为路径矢量路由协议 (Path-Vector Routing Protocol)。
- 每条BGP路由都携带多种路径属性 (Path attribute)，BGP可以通过这些路径属性控制路径选择，而不像IS-IS、OSPF只能通过Cost控制路径选择，因此在路径选择上，BGP具有丰富的可操作性，可以在不同场景下选择最合适的路径控制方式。



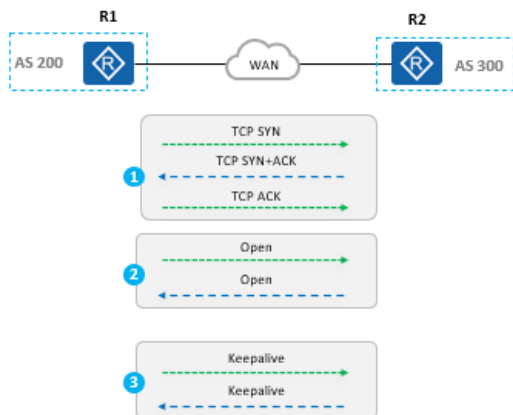
BGP对等体关系



- 与OSPF、IS-IS等协议不同，BGP的会话是基于TCP建立的。建立BGP对等体关系的两台路由器并不要求必须直连。
- BGP存在两种对等体关系类型：EBGP及IBGP：
 - EBGP (External BGP)：位于不同自治系统的BGP路由器之间的BGP对等体关系。两台路由器之间要建立EBGP对等体关系，必须满足两个条件：
 - 两个路由器所属AS不同 (即AS号不同)。
 - 在配置EBGP时，Peer命令所指定的对等体IP地址要求路由可达，并且TCP连接能够正确建立。
 - IBGP (Internal BGP)：位于相同自治系统的BGP路由器之间的BGP邻接关系。



BGP对等体关系建立 (1)

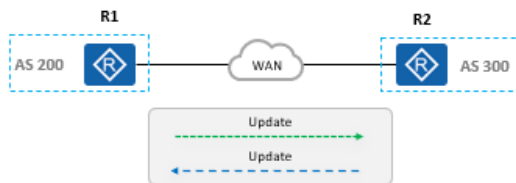


- 先启动BGP的一端先发起TCP连接，如左图所示，R1先启动BGP，R1使用随机端口号向R2的179端口发起TCP连接，完成TCP连接的建立。
- 三次握手建立完成之后，R1、R2之间相互发送Open报文，携带参数用于对等体建立，参数协商正常之后双方相互发送Keepalive报文，收到对端发送的Keepalive报文之后对等体建立成功，同时双方定期发送Keepalive报文用于保持连接。
- 其中Open报文中携带：
 - My Autonomous System：自身AS号
 - Hold Time：用于协商后续Keepalive报文发送时间
 - BGP Identifier：自身Router ID

- BGP 建立对等体的对等体都会发起 TCP 三次握手，所以会建立两个 TCP 连接，但是实际 BGP 只会保留其中一个 TCP 连接，从 Open 报文中获取对端 BGP Identifier 之后 BGP 对等体会比较本端的 Router ID 和对端的 Router ID 大小，如果本端 Router ID 小于对端 Router ID，则会关闭本地建立的 TCP 连接，使用由对端主动发起创建的 TCP 连接进行后续的 BGP 报文交互。

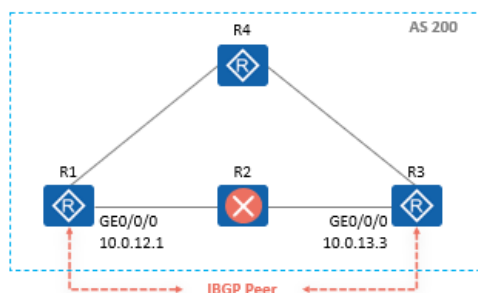


BGP对等体关系建立 (2)



BGP对等体关系建立之后，BGP路由器发送BGP Update (更新) 报文通告路由到对等体。

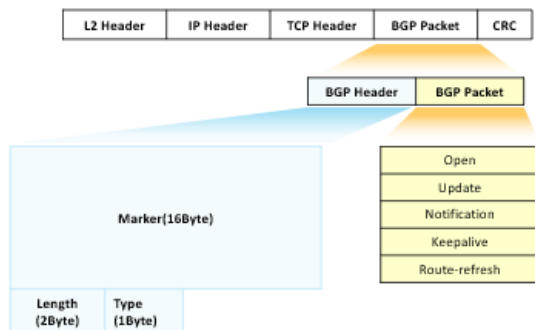
TCP连接源地址



- 缺省情况下，BGP使用报文出接口作为TCP连接的本地接口。
- 在部署IBGP对等体关系时，建议使用Loopback地址作为更新源地址。Loopback接口非常稳定，而且可以借助AS内的IGP和冗余拓扑来保证可靠性。
- 在部署EBGP对等体关系时，通常使用直连接口的IP地址作为源地址，如若使用Loopback接口建立EBGP对等体关系，则应注意EBGP多跳问题。

一般而言在AS内部，网络具备一定的冗余性。在R1与R3之间，如果采用直连接口建IBGP邻居关系，那么一旦接口或者直连链路发生故障，BGP会话也就断了，但是事实上，由于冗余链路的存在，R1与R3之间的IP连通性其实并没有DOWN（仍然可以通过R4到达彼此）。

BGP报文类型 (1)



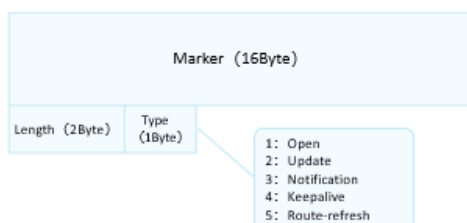
BGP存在5种类型的报文，不同类型的报文拥有相同的头部（header）。

- 不同于常见的IGP协议，BGP使用TCP作为传输层协议，端口号179，这使得BGP支持在非直连的路由器之间建立对等体关系。

BGP报文类型 (2)

报文名称	作用	发送时刻
Open	协商BGP对等体参数，建立对等体关系	BGP TCP连接建立成功之后
Update	发送BGP路由更新	BGP对等体关系建立之后有路由需要发送或路由变化时向对等体发送Update报文
Notification	报告错误信息，中止对等体关系	当BGP在运行中发现错误时，发送Notification报文将错误通告给BGP对等体
Keepalive	标志对等体建立，维持BGP对等体关系	BGP路由器收到对端发送的Keepalive报文，将对等体状态置为已建立，同时后续定期发送keepalive报文用于保持连接
Route-refresh	用于在改变路由策略后请求对等体重新发送路由信息。只有支持路由刷新能力的BGP设备会发送和响应此报文	当路由策略发生变化时，触发请求对等体重新通告路由

BGP报文格式 - 报文头格式



- BGP五种报文都拥有相同的报文头，格式如左侧所示，主要字段解释如下：
 - Marker: 16Byte，用于标明BGP报文边界，所有bit均为“1”。
 - Length: 2Byte，BGP报文总长度（包括报文头在内），以Byte为单位。
 - Type: 1Byte，BGP报文的类型。其取值从1到5，分别表示Open、Update、Notification、Keepalive和Route-refresh 报文。

-
- Opt Parm Len：Optional parameters 的长度。
- Optional parameters：宣告自身对于一些可选功能的支持，比如认证、多协议支持。
- 除了 IPv4 单播路由信息，BGP4+还支持多种网络层协议（如 IPv6、组播），在协商时 BGP 对等体之间会通过 Optional parameters 字段协商对网络层协议的支持能力。



BGP报文格式 - Update

Unfeasible routes length (2Byte)
Withdrawn routes (NByte)
Total path attribute length (2Byte)
Path attributes (NByte)
NLRI (NByte)

- Update报文用于在对等体之间传递路由信息，可以用于发布、撤销路由。
- 一个Update报文可以通告具有相同路径属性的多条路由，这些路由保存在NLRI (Network Layer Reachable Information, 网络层可达信息) 中。同时Update还可以携带多条不可达路由，用于告知对方撤销路由，这些保存在Withdrawn Routes字段中。
- 报文格式如左侧所示，主要字段解释如下：
 - Withdrawn routes: 不可达路由的列表。
 - Path attributes: 与NLRI相关的所有路径属性列表，每个路径属性由一个TLV (Type-Length-Value) 三元组构成。
 - NLRI: 可达路由的前缀和前缀长度二元组。
- Unfeasible routes length: 不可达路由字段的长度，以 Byte 为单位。如果为 0 则说明没有 Withdrawn Routes 字段。
- Withdrawn Routes Length: 标明 Withdrawn Routes 部分的长度。其值为零时，表示没有撤销的路由。
- Total path attribute length: 路径属性字段的长度，以 Byte 为单位。如果为 0 则说明没有 Path Attributes 字段。



BGP报文格式 - Notification

Error code (8bit)	Error subcode (8bit)	
Data (可变长度)		

- 当BGP检测到错误状态时（对等体关系建立时、建立之后都可能发生），就会向对等体发送Notification，告知对端错误原因。之后BGP连接将会立即中断。
 - Error Code、Error subcode: 差错码、差错子码，用于告知对端具体的错误类型。
 - Data: 用于辅助描述详细的错误内容，长度并不固定。



BGP报文格式 - Keepalive



- BGP路由器收到对端发送的Keepalive报文，将对等体状态置为已建立，同时后续定期发送keepalive报文用于保持连接。
- Keepalive报文格式中只包含报文头，没有附加任何其他字段。



BGP报文格式 - Route-refresh



- Route-refresh报文用来要求对等体重新发送指定地址族的路由信息，一般为本端修改了相关路由策略之后让对方重新发送Update报文，本端执行新的路由策略重新计算BGP路由。
- 相关字段内容如下：
 - AFI: Address Family Identifier, 地址族标识, 如IPv4。
 - Res.: 保留, 8个bit必须置0。
 - SAFI: Subsequent Address Family Identifier, 子地址族标识。

- 在 Open 报文协商时会协商是否支持 Route-refresh，如果对等体支持 Route-refresh 能力，则可以通过 **refresh bgp** 命令手工对 BGP 连接进行软复位，BGP 软复位可以在不中断 BGP 连接的情况下重新刷新 BGP 路由表，并应用新的策略。
- 对于不支持 Route-Refresh 能力的 BGP 对等体，可以配置 **keep-all-routes** 命令，保留该对等体的所有原始路由，这样不需要复位 BGP 连接即可完成路由表的刷新。
- 缺省情况下未开启 **keep-all-routes**。

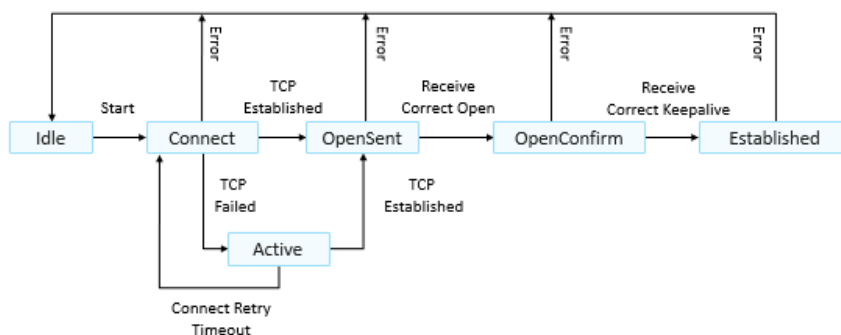


BGP状态机 (1)

Peer状态名称	用途
Idle	开始准备TCP的连接并监视远程对等体, 启用BGP时, 要准备足够的资源
Connect	正在进行TCP连接, 等待完成中, 认证都是在TCP建立期间完成的。如果TCP连接建立失败则进入Active状态, 反复尝试连接
Active	TCP连接没建立成功, 反复尝试TCP连接
OpenSent	TCP连接已经建立成功, 开始发送Open包, Open包携带参数协商对等体的建立
OpenConfirm	参数、能力特性协商成功, 自己发送Keepalive包, 等待对方的Keepalive包
Established	已经收到对方的Keepalive包, 双方能力特性协商发现一致, 开始使用Update通告路由信息



BGP状态机 (2)



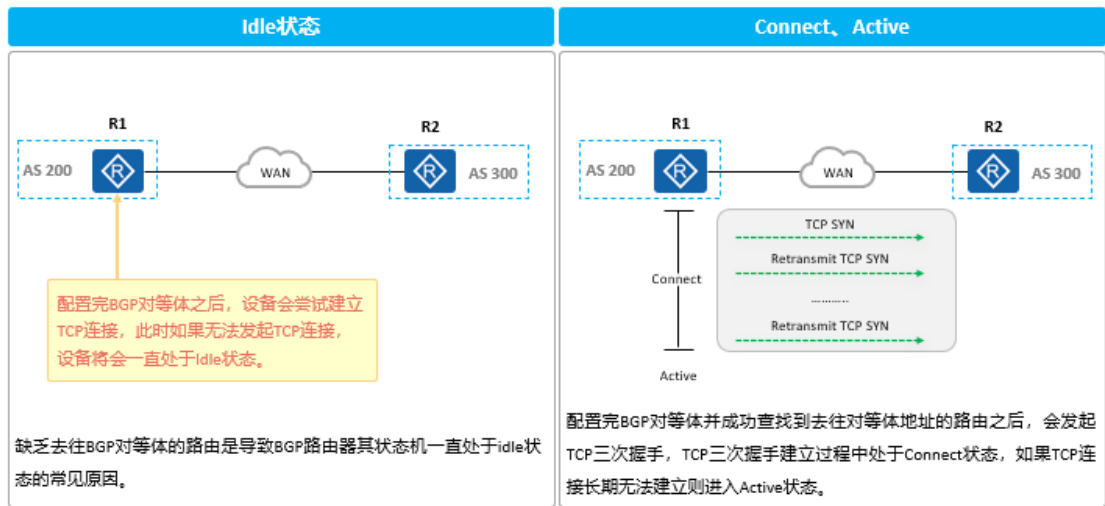
- Idle 状态是 BGP 初始状态。在 Idle 状态下, BGP 拒绝对等体发送的连接请求。只有在收到本设备的 Start 事件后, BGP 才开始尝试和其它 BGP 对等体进行 TCP 连接, 并转至 Connect 状态。
- Start 事件是由一个操作者配置一个 BGP 过程, 或者重置一个已经存在的过程或者路由器软件重置 BGP 过程引起的。
- 任何状态中收到 Notification 报文或 TCP 拆链通知等 Error 事件后, BGP 都会转至 Idle 状态。
- 在 Connect 状态下, BGP 启动连接重传定时器 (Connect Retry), 等待 TCP 完成连接。
- 如果 TCP 连接成功, 那么 BGP 向对等体发送 Open 报

文，并转至 OpenSent 状态。

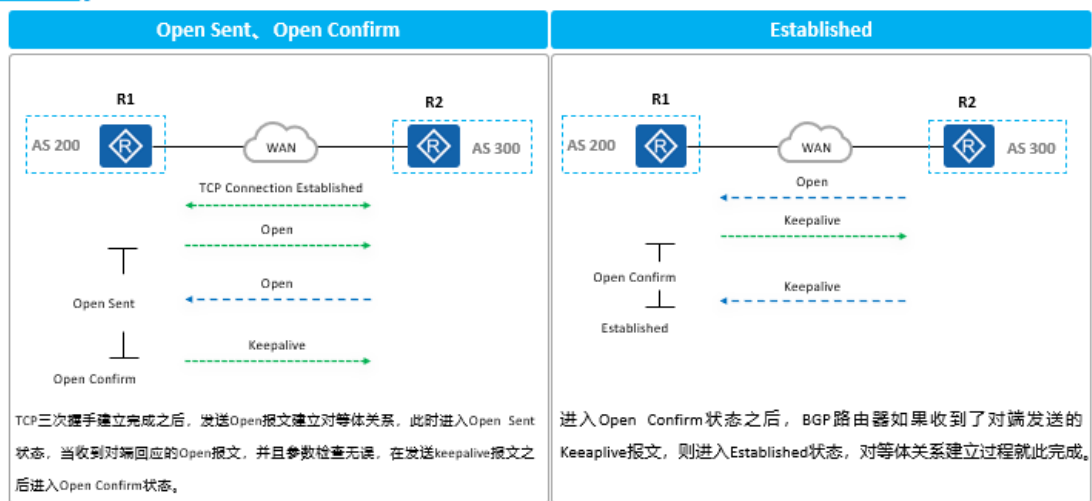
- 如果 TCP 连接失败，那么 BGP 转至 Active 状态。
- 如果连接重传定时器超时，BGP 仍没有收到 BGP 对等体的响应，那么 BGP 继续尝试和其它 BGP 对等体进行 TCP 连接，停留在 Connect 状态。
- 在 Active 状态下，BGP 总是在试图建立 TCP 连接。
- 如果 TCP 连接成功，那么 BGP 向对等体发送 Open 报文，关闭连接重传定时器，并转至 OpenSent 状态。
- 如果 TCP 连接失败，那么 BGP 停留在 Active 状态。
- 如果连接重传定时器超时，BGP 仍没有收到 BGP 对等体的响应，那么 BGP 转至 Connect 状态。
- 在 OpenSent 状态下，BGP 等待对等体的 Open 报文，并对收到的 Open 报文中的 AS 号、版本号、认证码等进行检查。
- 如果收到的 Open 报文正确，那么 BGP 发送 Keepalive 报文，并转至 OpenConfirm 状态。
- 如果发现收到的 Open 报文有错误，那么 BGP 发送 Notification 报文给对等体，并转至 Idle 状态。
- 在 OpenConfirm 状态下，BGP 等待 Keepalive 或 Notification 报文。如果收到 Keepalive 报文，则转至 Established 状态，如果收到 Notification 报文，则转至 Idle 状态。
- 在 Established 状态下，BGP 可以和对等体交换 Update、Keepalive、Route-refresh 报文和 Notification 报文。
- 如果收到正确的 Update 或 Keepalive 报文，那么 BGP 就认为对端处于正常运行状态，将保持 BGP 连接。
- 如果收到错误的 Update 或 Keepalive 报文，那么 BGP 发送 Notification 报文通知对端，并转至 Idle 状态。
- Route-refresh 报文不会改变 BGP 状态。
- 如果收到 Notification 报文，那么 BGP 转至 Idle 状态。
- 如果收到 TCP 拆链通知，那么 BGP 断开连接，转至 Idle

状态。

BGP状态机详解 (1)



BGP状态机详解 (2)





BGP对等体表

```
<R1>display bgp peer
BGP local router ID : 10.0.1.1
Local AS number : 100
Total number of peers : 1      Peers in established state : 1
```

Peer	V	AS	MsgRcvd	MsgSent	OutQ	Up/Down	State	PrefRcv
10.0.12.2	4	100	25719	25714	0	0428h32m	Established	1

- 在设备上通过**display bgp peer**命令查看BGP对等体表，其中主要参数含义：

- Peer：对等体地址
- V：version，版本号
- AS：对等体AS号
- Up/Down：该对等体已经存在up或者down的时间
- State：对等体状态，这里显示的为BGP状态机的状态
- PrefRcv：prefix received，从该对等体收到的路由前缀数目

- BGP 对等体表的作用为列出本设备的 BGP 对等体，以及对等体的状态等信息。
- MsgRcvd、MsgSent：从对等体收到的报文个数，向对等体发送的报文个数。
- OutQ：out queue，对外发送报文队列中排队的个数，一般为 0。



BGP路由表 (1)

```
<R1>display bgp routing-table
BGP Local router ID is 10.0.1.1
Status codes: * - valid, > - best, d - damped,
              h - history, i - internal, s - suppressed, S - Stale
Origin : i - IGP, e - EGP, ? - incomplete
Total Number of Routes: 2
```

	Network	NextHop	MED	LocPrf	PrefVal	Path/Ogn
*>i	10.0.45.0/24	10.0.4.4	0	100	0	?
* i		10.0.4.4	0	100	0	?

- 在设备上通过**display bgp routing-table**查看BGP路由表：

- Network：路由的目的网络地址以及网络掩码
- NextHop：下一跳地址

- 如果想要查看某条路由更加详细的信息，可以通过**display bgp routing-table ipv4-address { mask / mask-length}**查看，该命令会将匹配的BGP路由信息详细展示。

- 列出本设备发现的所有 BGP 路由，如果到达同一个目的

地存在多条路由，则将路由都进行罗列，但每个目的地只会优选一条路由。

概览 对等体关系 报文及状态机 协议表项 路由生成 通告原则

BGP路由表 (2)

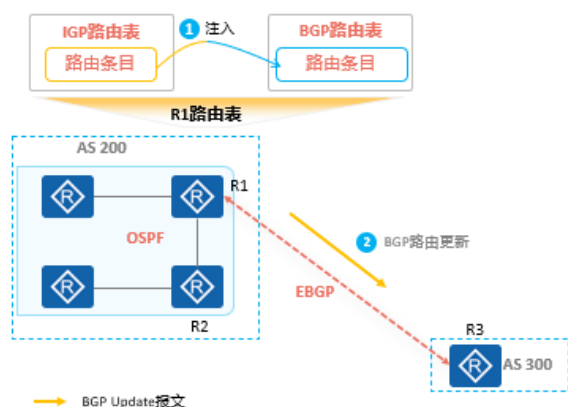
```
<R1>display bgp routing-table 10.0.45.0 24
BGP local router ID : 10.0.1.1
Local AS number : 100
Paths: 2 available, 1 best, 1 select
BGP routing table entry information of 10.0.45.0/24:
From: 10.0.2.2 (10.0.2.2)          #标明路由来源
Route Duration: 06h19m44s
Relay IP Nexthop: 10.0.12.2
Relay IP Out-Interface: GigabitEthernet0/0/0
Original nexthop: 10.0.4.4          #路由下一跳地址
Qos information : 0x0
AS-path Nil, origin incomplete, MED 0, localpref 100, pref-val 0, valid, internal, best,
select, active, pre 255, IGP cost 2 #路径属性、是否被优选
Originator: 10.0.4.4
Cluster list: 10.0.2.2
Not advertised to any peer yet
```

```
BGP routing table entry information of 10.0.45.0/24:
From: 10.0.3.3 (10.0.3.3)
Route Duration: 05h17m56s
Relay IP Nexthop: 10.0.12.2
Relay IP Out-Interface: GigabitEthernet0/0/0
Original nexthop: 10.0.4.4
Qos information : 0x0
AS-path Nil, origin incomplete, MED 0, localpref 100, pref-val 0, valid, internal, pre
255, IGP cost 2, not preferred for peer address
Originator: 10.0.4.4
Cluster list: 10.0.3.3
Not advertised to any peer yet
```

- 通过 `display bgp routing-table ipv4-address { mask / mask-length }` 可以显示指定 IP 地址/掩码长度的路由信息，在其中有关于该 BGP 路由的详细信息，如：路由始发者、下一跳地址、路由的路径属性等。

概览 对等体关系 报文及状态机 协议表项 路由生成 通告原则

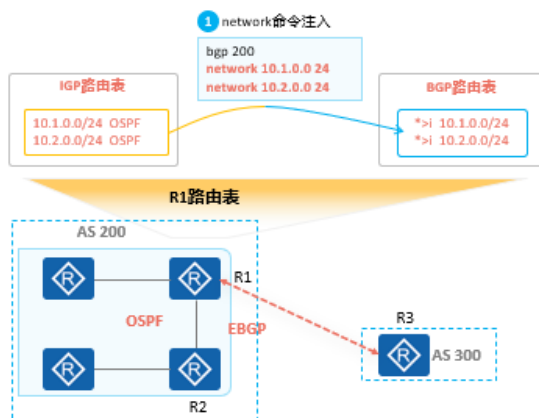
BGP路由的生成



- 不同于IGP路由协议，BGP自身并不会发现并计算产生路由，BGP将IGP路由表中的路由注入到BGP路由表中，并通过Update报文传递给BGP对等体。
- BGP注入路由的方式有两种：
 - Network
 - import-route
- 与IGP协议相同，BGP支持根据已有的路由条目进行聚合，生成聚合路由。



Network注入路由 (1)

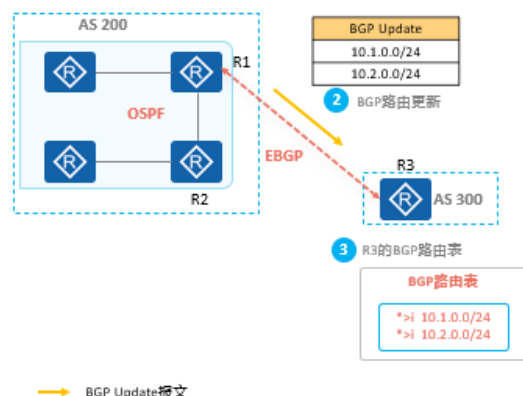


通过Network方式注入路由：

1. AS200内的BGP路由器已经通过IGP协议OSPF学习到了两条路由：10.1.0.0/24和10.2.0.0/24，在BGP进程内通过network命令注入这两条路由，这两条路由将会出现在本地的BGP路由表中。
2. AS200内的BGP路由器通过Update报文将路由传递给AS300内的BGP路由器。
3. AS300内的BGP路由器收到路由后，将这两条路由加入到本地的BGP路由表中。



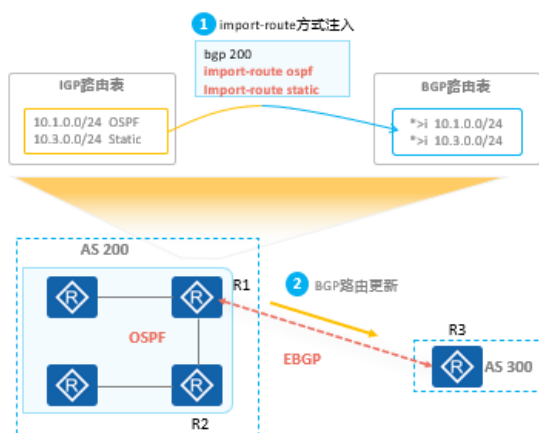
Network注入路由 (2)



通过Network方式注入路由：

1. AS200内的BGP路由器已经通过IGP协议OSPF学习到了两条路由：10.1.0.0/24和10.2.0.0/24，在BGP进程内通过network命令注入这两条路由，这两条路由将会出现在本地的BGP路由表中。
2. AS200内的BGP路由器通过Update报文将路由传递给AS300内的BGP路由器。
3. AS300内的BGP路由器收到路由后，将这两条路由加入到本地的BGP路由表中。

import-route方式注入路由

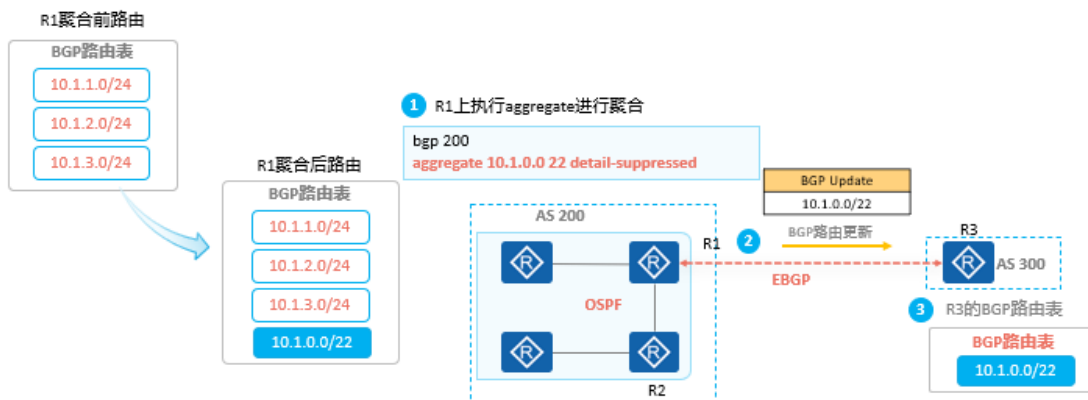


• Network方式注入路由虽然是精确注入，但是只能一条条配置逐条注入IP路由表中的路由，如果注入的路由条目很多配置命令将会非常复杂，为此可以使用import-route方式，将：

1. 直连路由
2. 静态路由
3. OSPF路由
4. IS-IS路由

等协议的路由注入到BGP路由表中。

BGP聚合路由



与众多IGP协议相同，BGP同样支持路由的手工聚合，在BGP配置视图使用aggregate命令可以执行BGP路由手工聚合，在BGP已经学习到相应的明细路由情况下，设备会向BGP注入指定的聚合路由。

- 执行聚合之后，在本地的 BGP 路由表中除了原本的明细路由条目之外，还会多出一条聚合的路由条目。
- 如果在执行聚合时指定了 **detail-suppressed**，则 BGP 只会向对等体通告聚合后的路由，而不通告聚合前的明细路由。
- 在聚合时配置了抑制明细路由的参数，R3 上查看路由表，将只能看到 BGP 路由：10.1.0.0/22，无法看到聚合前的明细路由。

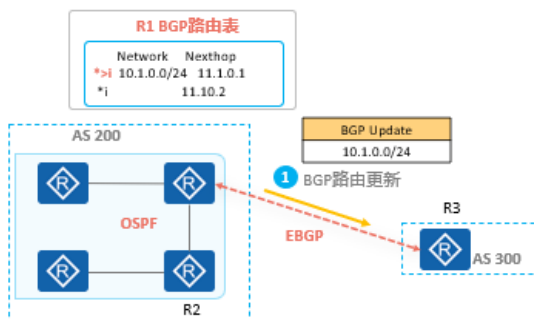


通告原则

- BGP通过network、import-route、aggregate聚合方式生成BGP路由后，通过Update报文将BGP路由传递给对等体。
- BGP通告遵循以下原则：
 - 只发布最优且有效路由。
 - 从EBGP对等体获取的路由，会发布给所有对等体。
 - IBGP水平分割：从IBGP对等体获取的路由，不会发送给IBGP对等体。
 - BGP同步规则指的是：当一台路由器从自己的IBGP对等体学习到一条BGP路由时（这类路由被称为IBGP路由），它将不能使用该条路由或把这条路由通告给自己的EBGP对等体，除非它又从IGP协议（例如OSPF等，此处也包含静态路由）学习到这条路由，也就是要求IBGP路由与IGP路由同步。同步规则主要用于规避BGP路由黑洞问题。



BGP路由通告原则一



- 第一条原则：只发布最优且有效（即下一跳地址可达）路由。
- 通过display bgp routing-table命令可以查看BGP路由表。

Total Number of Routes: 2

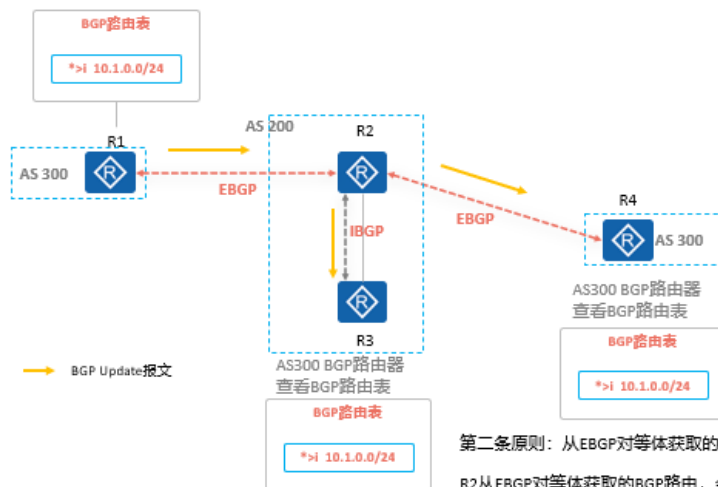
Network	NextHop	MED	LocPrf	PrefVal	Path/Ogn
*> 10.1.0.0/24	11.1.0.1	0	100	0	?
*i 11.10.1	11.1.0.2	0	100	0	?

- 在BGP路由表中同时存在以下两个标志的路由为最优、有效：
 - *:代表有效
 - >:代表最优

→ BGP Update报文



BGP路由通告原则二

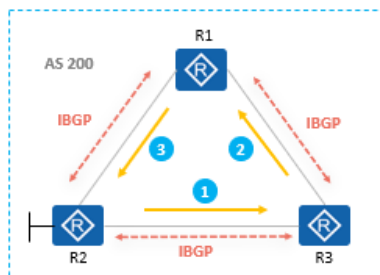


第二条原则：从EBGP对等体获取的路由，会发布给所有对等体。

R2从EBGP对等体获取的BGP路由，会发布给所有EBGP、IBGP对等体。

BGP路由通告原则三 (1)

→ BGP Update报文

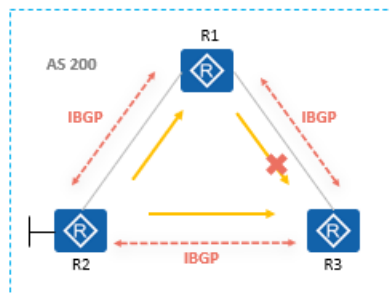


- 第三条原则：从IBGP对等体获取的BGP路由，不会再发送给其他IBGP对等体。
- 该条原则也被称为“IBGP水平分割”。
- 如图所示，如果IBGP对等体学习到的路由会继续传递给其他的IBGP对等体：
 - R2将一条路由传递给了IBGP对等体R3
 - R3收到路由之后传递给IBGP对等体R1
 - R1继续传递给IBGP对等体R2
 路由环路形成。

BGP路由通告原则三 (2)

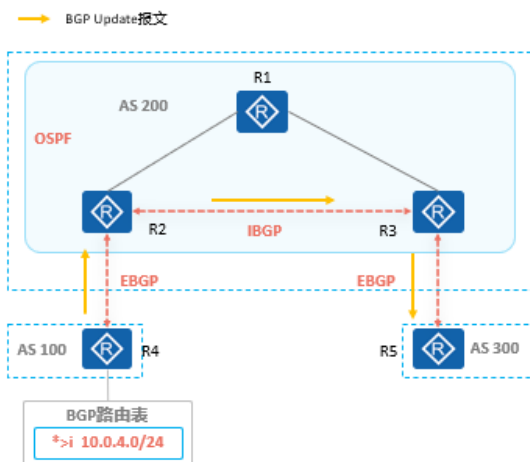
IBGP全互联

→ BGP Update报文



- 第三条原则可能会带来新的问题，如左侧所示，当BGP路由器R2将路由传递给BGP路由器R1时，由于第三条原则限制，R1无法将BGP路由传递给R3，R3将无法学习到路由。
- 为解决该问题可以采用AS内IBGP全互联的方式，即：R2、R3之间建立非直连的IBGP对等体关系，以此让BGP路由器R2将路由传递给BGP路由器R3。

BGP路由通告原则四 (1)

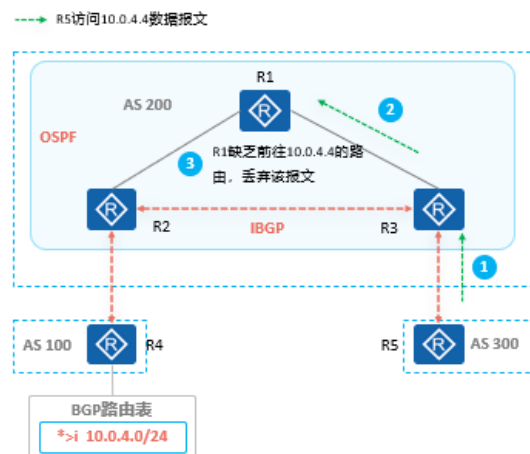


- 第四条原则：当一台路由器从自己的IBGP对等体学习到一条BGP路由时（这类路由被称为IBGP路由），它将不能使用该路由或把这条路由通告给自己的EBGP对等体，除非它又从IGP协议（例如OSPF等，此处也包含静态路由）学习到这条路由，该条规则也被称为BGP同步原则。

如图所示：

1. BGP路由器R4上存在一条路由由10.0.4.0/24，R4将其传递给了R2。
2. R2将路由传递给非直连IBGP对等体R3。
3. R3将路由传递给R5。
4. 之后R5向10.0.4.4发起访问。

BGP路由通告原则四 (2)



R5访问10.0.4.4：

1. R5查找路由表，将报文发送给R3。
2. R3收到报文后查找路由表，匹配到一条BGP路由，其下一跳为R2，但是R2为非直连下一跳，需要进行路由迭代，通过IGP学习到的路由迭代出下一跳为R1。R3将报文发送给R1。
3. R1收到报文后查找路由表，因为R1并非BGP路由器，未与R2建立IBGP对等体关系，因此R1上并无BGP路由10.0.4.0/24，路由查找失败，R1将报文丢弃。

- 产生该问题根本原因为 AS200 域内未运行 BGP 的路由器并无从 BGP 学习到的路由条目，查找路由失败，导致 R1 丢弃报文。为此制定了 BGP 同步原则：

当 BGP 的路由条目也存在于 IGP 路由表时才对外发送，以图中场景为例，当 R3 查看 IGP 路由表，OSPF 路由表中并无路由 10.0.4.0/24，因此并不会向 R5 发送该路由，自然也不会产生后续的访问失败问题。

- 解决该问题的方式有：
- 将 BGP 路由重分发到 IGP 中，基本不会使用该方式。

- 建立全互联的 IBGP 对等体关系，让全网所有路由器都拥有 BGP 路由。

配置介绍

1. 启动BGP进程

```
[Huawei] bgp { as-number-plain | as-number-dot }  
[Huawei-bgp] router-id ipv4-address
```

启动BGP，指定本地AS编号，并进入BGP视图。使用router-id命令配置BGP的Router ID，建议将BGP Router ID配置为设备Loopback接口的地址。

2. 配置BGP对等体

```
[Huawei-bgp] peer { ipv4-address | ipv6-address } as-number { as-number-plain | as-number-dot }
```

创建BGP对等体，指定对等体地址以及AS号。

3. 配置建立对等体使用的源地址、EBGP对等体最大跳数

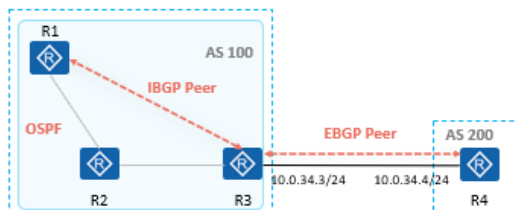
```
[Huawei-bgp] peer ipv4-address connect-interface interface-type interface-number [ ipv4-source-address ]  
[Huawei-bgp] peer ipv4-address ebgp-max-hop [ hop-count ]
```

指定发送BGP报文的源接口，并可指定发起连接时使用的源地址。

指定建立EBGP连接允许的最大跳数。缺省情况下，EBGP连接允许的最大跳数为1，即只能在物理直连链路上建立EBGP连接。

- 如果没有配置 Router ID，则 BGP 会自动选取系统视图下的 Router ID 作为 BGP 协议的 Router ID。系统视图下的 Router ID 选择规则，请参见命令 router-id 中的描述。

配置案例 (1)



R1的配置如下：

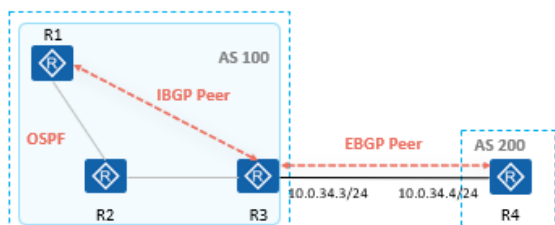
```
[R1] bgp 100  
[R1-bgp] router-id 10.0.1.1  
[R1-bgp] peer 10.0.3.3 as-number 100  
[R1-bgp] peer 10.0.3.3 connect-interface LoopBack1
```

R3的配置如下：

```
[R3] bgp 100  
[R3-bgp] router-id 10.0.3.3  
[R3-bgp] peer 10.0.1.1 as-number 100  
[R3-bgp] peer 10.0.1.1 connect-interface LoopBack1  
[R3-bgp] peer 10.0.34.4 as-number 200
```

- BGP对等体关系、AS号、设备互联地址如图所示。
- 所有设备的Loopback1接口地址为10.0.x.x/32，其中x为设备编号，所有设备都使用Loopback1地址作为Router ID。
- R1、R3之间使用Loopback1地址作为更新源地址建立IBGP对等体关系，R3、R4之间使用互联接口地址作为更新源地址建立EBGP对等体关系。

配置案例 (2)



R4的配置如下:

```
[R4] bgp 200
[R4-bgp] router-id 10.0.4.4
[R4-bgp] peer 10.0.34.3 as-number 100
```

- BGP对等体关系、AS号、设备互联地址如图所示。
- 所有设备的Loopback1接口地址为10.0.x.x/32，其中x为设备编号，所有设备都使用Loopback1地址作为Router ID。
- R1、R3之间使用Loopback1地址作为更新源地址建立IBGP对等体关系，R3、R4之间使用互联接口地址作为更新源地址建立EBGP对等体关系。

配置案例 (3)

在R3上查看BGP对等体状态:

```
<R3> display bgp peer
BGP Local router ID : 10.0.3.3
Local AS number : 100
Total number of peers : 2
Peers in established state : 2
```

Peer	V	AS	MsgRcvd	MsgSent	OutQ	Up/Down	State	PreRcv
10.0.1.1	4	100	0	0	0	00:00:07	Established	0
10.0.34.4	4	200	32	35	0	00:17:49	Established	0

思考题：

- (简答题) BGP 使用的 TCP 目的端口号是多少？
- (简答题) BGP 对等体关系有哪几种？划分的依据是什么？
- (多选题) BGP 对等体关系建立、更新路由分别使用 ()、() 报文。
- Route-refresh
- Open
- Notification
- Update

答案：

- 179
- IBGP、EBGP 对等体关系，对等体是否和自身处于一个 AS 内。
- B、D