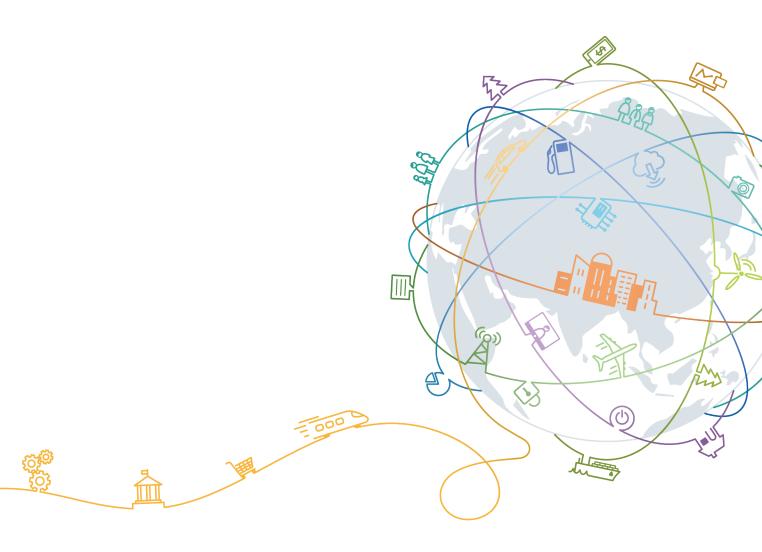
如何从堆叠切换为 M-LAG

文档版本 02

发布日期 2020-09-18





版权所有 © 华为技术有限公司 2020。 保留一切权利。

非经本公司书面许可,任何单位和个人不得擅自摘抄、复制本文档内容的部分或全部,并不得以任何形式传播。

商标声明



HUAWE和其他华为商标均为华为技术有限公司的商标。

本文档提及的其他所有商标或注册商标,由各自的所有人拥有。

注意

您购买的产品、服务或特性等应受华为公司商业合同和条款的约束,本文档中描述的全部或部分产品、服务或特性可能不在您的购买或使用范围之内。除非合同另有约定,华为公司对本文档内容不做任何明示或默示的声明或保证。

由于产品版本升级或其他原因,本文档内容会不定期进行更新。除非另有约定,本文档仅作为使用指导,本文档中的所有陈述、信息和建议不构成任何明示或暗示的担保。

华为技术有限公司

地址: 深圳市龙岗区坂田华为总部办公楼 邮编: 518129

网址: https://www.huawei.com

客户服务邮箱: support@huawei.com

客户服务电话: 4008302118

目录

1 什么是堆叠?	1
2 什么是 M-LAG?	2
3 为什么要从堆叠切换为 M-LAG?	3
4 堆叠切换 M-LAG 之前需要哪些检查项	5
4.1 检查 MSTP、Smart Link 配置是否生效	5
4.2 检查 ACL 资源是否充足	6
4.3 检查二层子接口资源是否充足	10
4.4 检查业务是否双归	11
4.5 检查 License 是否一致	11
4.6 检查动态路由协议建立的邻居数量是否超限	13
4.7 检查 IPv6 业务是否存在	14
4.8 检查跨设备的二层端口隔离配置是否生效	15
4.9 检查跨设备的单向业务是否存在	15
4.10 检查系统 MAC 是否为堆叠主设备	17
4.11 检查是否是 SDN 网络	18
4.12 记录 MAC 表项和流量情况	18
5 五步完成 堆叠切换为 M-LAG	20
5.1 删除堆叠双主检测链路	21
5.2 隔离堆叠备设备,增加 M-LAG 配置	21
5.3 切换业务	23
5.4 隔离堆叠主设备,增加 M-LAG 配置	23
5.5 查询切换后状态,切换前后业务对比	25

1 什么是堆叠?

堆叠技术是指把多个支持堆叠的设备组合在一起,逻辑上合为一台整体设备。用户可以将这多台设备看成一台单一设备进行管理和使用。这样既可以通过增加设备来扩展 端口数量和交换能力,同时也通过多台设备之间的互相备份增强了设备的可靠性。

如<mark>图1-1</mark>所示,DeviceA和DeviceB通过堆叠链路连接在一起,从逻辑上构成一台设备,并作为一个整体参与数据转发。

DeviceA和DeviceB互相备份,当DeviceA故障时,DeviceB可以接替DeviceA保证系统的正常运行。堆叠基本工作原理请参见《配置指南-虚拟化》里的<mark>堆叠配置</mark>。

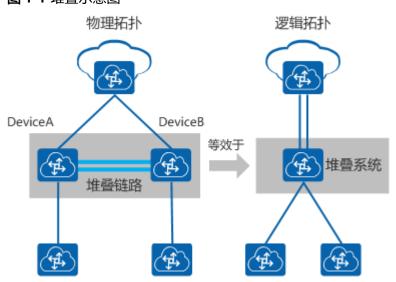


图 1-1 堆叠示意图

2 什么是 M-LAG?

M-LAG(Multichassis Link Aggregation Group)是一种新兴的跨设备链路聚合的技术。其基本思想是让两台接入交换机以同一个状态和被接入的设备进行链路聚合协商,在被接入的设备看来,就如同和一台设备建立了链路聚合关系。通过跨设备的链路聚合,可以将可靠性从单板级提高到设备级。

如<mark>图2-1</mark>所示,DeviceA和DeviceB间部署M-LAG,M-LAG设备和ServerA通过跨设备的链路聚合进行链路聚合协商,实现ServerA的双归接入。

DeviceA与DeviceB形成负载分担,共同进行流量转发。当DeviceA或DeviceB发生故障时,流量可以快速切换到另一台设备,保证业务的正常运行。M-LAG基本工作原理请参见《配置指南-以太网交换》里的M-LAG(跨设备链路聚合)配置。

M-LAG不仅解决了传统聚合链路可靠性低的问题,同时规避了堆叠在升级过程中时间长、风险高等缺点。

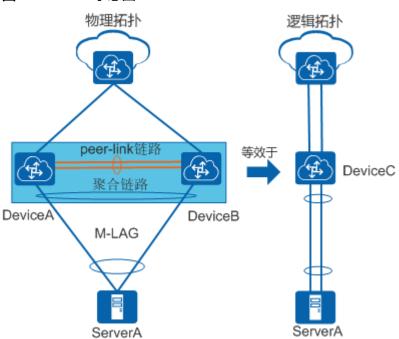


图 2-1 M-LAG 示意图

3 为什么要从堆叠切换为 M-LAG?

堆叠和M-LAG作为广泛运用于数据中心网络接入层的两种横向虚拟化技术,他们都可实现终端的冗余接入,实现链路冗余备份,提高数据中心网络的可靠性和可扩展性。然而,与堆叠技术相比,M-LAG存在更高的可靠性和独立升级的优势。

如表3-1所示,对比了堆叠和M-LAG的优劣。针对升级过程业务中断时间要求高、对组网可靠性要求高的场景,我们推荐用户使用M-LAG技术替代堆叠技术,用作数据中心网络终端接入技术。

表 3-1 堆叠和 M-LAG 的对比

对比维度	堆叠	M-LAG(推荐)
可靠性	一般: 控制面集中,故障可能在成员设备上扩散设备级、单板级、链路级等都具备高可靠性	更高: 控制面独立,故障域隔离设备级、单板级、链路级等都具备高可靠性
配置复杂度	简单:逻辑上是一台设备	简单: 两台设备独立配置
成本	一般: 需要部署堆叠线缆	一般:需要部署Peer-link连线
性能	一般:主交换机控制面要控制 所有堆叠成员的转发面,CPU 负载加重	高:成员交换机独立转发,CPU负 载保持不变
升级复杂度	高:通过堆叠快速升级可以降 低业务中断时间,但升级操作 时间变长,升级风险变高	低:两台设备可分别单独升级, 升级操作简单,风险低
升级中断时间	相对较长:通过堆叠快速升级,典型配置组网下,业务中断时间在20秒~1分钟左右,与业务量强相关	短:流量秒级中断
网络设计	相对简单:堆叠设备逻辑上为 一台设备,网络结构较简单	相对复杂:M-LAG设备逻辑上仍 然是两台设备,网络结构较复杂

对比维度	堆叠	M-LAG(推荐)
适用场景	对软件版本升级中断时间无要求希望网络维护简单	对软件版本升级时业务中断时间要求较高对网络可靠性要求更高可接受增加一定程度的维护复杂度

4 堆叠切换 M-LAG 之前需要哪些检查项

- 4.1 检查MSTP、Smart Link配置是否生效
- 4.2 检查ACL资源是否充足
- 4.3 检查二层子接口资源是否充足
- 4.4 检查业务是否双归
- 4.5 检查License是否一致
- 4.6 检查动态路由协议建立的邻居数量是否超限
- 4.7 检查IPv6业务是否存在
- 4.8 检查跨设备的二层端口隔离配置是否生效
- 4.9 检查跨设备的单向业务是否存在
- 4.10 检查系统MAC是否为堆叠主设备
- 4.11 检查是否是SDN网络
- 4.12 记录MAC表项和流量情况

4.1 检查 MSTP、Smart Link 配置是否生效

检查是否存在多实例MSTP、Smart Link配置。执行命令行**display stp global**,查看 **Protocol Status**和**Mode**字段是否设备使能MSTP且功能生效。执行命令行**display smart-link group**,**查看Smart Link group**字段为**enabled**时表明设备已经使能 Smart Link。由于M-LAG系统不支持MSTP和Smart Link,这种情况下不进行切换。

显示MSTP概要信息。 <HUAWEI> display stp global

Protocol Status :Enabled Bpdu-filter default :Disabled Tc-protection :Enabled Tc-protection threshold :1 Tc-protection interval :2s Edged port default :Disabled Pathcost-standard :Dot1T Timer-factor :3 Transmit-limit :6 Bridge-diameter :7

CIST Global Information:

//使能生成树协议

Mode :MSTP //生成树协议模式是MSTP CIST Bridge :32768.0019-7459-3301 Config Times :Hello 2s MaxAge 20s FwDly 15s MaxHop 20 Active Times :Hello 2s MaxAge 20s FwDly 15s MaxHop 20 CIST Root/ERPC :32768.0019-7459-3301 / 0 (This bridge is the root) CIST RegRoot/IRPC :32768.0019-7459-3301 / 0 (This bridge is the root) CIST RootPortId :0.0 BPDU-Protection :Disabled TC or TCN received :9 TC count per hello :0 STP Converge Mode :Normal Share region-configuration :Enabled Time since last TC :0 days 1h:37m:17s Number of TC :10 Last TC occurred :10GE4/0/12 Topo Change Flag :0

查看Smart Link组的状态信息。

```
<HUAWEI> display smart-link group 1
Smart Link group 1 information:
 Smart Link group: enabled //Smart Link组已使能
 Link status: Lock
 Wtr-time is: 60 sec.
 Load-Balance Instance: 10
 Protected-VLAN reference-instance: --
 DeviceID: 0025-9e80-2494 Control-VLAN ID: 505
 Member
            Role InstanceID State Flush Count LastFlushTime
 10GE1/0/1 Master
                       0 Active
                                        0 0000/00/00 00:00:00 UTC+00:00 //成员接口为转发状态
 10GE1/0/1
            Master
                       10 Inactive
                                        0 0000/00/00 00:00:00 UTC+00:00
                       0 Inactive
 10GE1/0/2
            Slave
                                      0 0000/00/00 00:00:00 UTC+00:00
 10GE1/0/2 Slave
                      10 Active
                                       0 0000/00/00 00:00:00 UTC+00:00
```

4.2 检查 ACL 资源是否充足

查看设备上的ACL资源,判断ACL资源是否足够。M-LAG会占用ACL资源如下:

- 入方向占用CPCAR L2已有的组,同时下发M-LAG ARP和M-LAG Protocol。
- 出方向占用一个分组160位宽(单芯片出方向总共320位宽),包括下发M-LAG IPV4 UC、M-LAG IPV6 UC、M-LAG Isolate。

□ 说明

因为设备出方向ACL总资源较少,所以在检查ACL资源时需要特别关注。

以CE12800为例:

步骤1 执行display system tcam service brief命令,查看不同的业务占用的组索引和规则计数。

<huawei> (Slot: 1</huawei>	display sy	stem tcam s	service brief slot 1	
Chip Grou (FEI/FE		th Stage	ServiceName	Count
0 2/2	320Bit	Ingress	BPDU Deny	21
2/2	320Bit	Ingress	CPCAR L2	4
2/2	320Bit	Ingress	L2 Protocol Tunnel	1
3/3	320Bit	Ingress	App-Session	2

3,	/3	320Bit	Ingress	CPCAR L3	19

设备共有12个组,GroupID字段显示已有两个组被占用,还剩下10个组,足够M-LAG ARP和M-LAG Protocol两个组使用。

步骤2 执行display system tcam bank resource命令,查看各业务使用资源的情况。

				_					
Bankld	Entry	Entry Free	Entry Used	Stage (Fl	Grou _l EI/FE)		ВТуре	KBId Service	Name
0,1	320Bit	957 6 1 4 31 2	21	3/3	L2	2,3 2,3 2,3 2,3	CPCAR L2 Prote App-Se	ocol Tunnel ession LL3	
2,3	320Bit	1003 2 12	5	6/1 6/1 6/1 6/1		IPv6 4,5 4,5 4,5 4,5	4,5 CPCAR EVN Pa	CPCAR Ipv6 Dci	
4	160Bit	1020	3	Ingress 5/6 5/6	5/6 L2 MPLS IPv6	IPv4 1		CPCAR Termina	ited v4
5 6 7 8 9 10 11 12 13	160Bit	1011	12	Ingress	8/7	MPLS		MPLS PHP R Vxlan Ipv6	
Usage	otal Used								
L2 8	2(1,4) 5(1,2,3,4,5 3(1,4,5) 7(1,2,3,4,	5) 5.6.7)	3(0	0,2,3,5,6,7) 0,6,7) 1,2,3,6,7) (0)					

□ 说明

回显**Stage**字段没有Egress资源被占用,说明出方向有足够的资源给M-LAG使用。步骤3~步骤9里的命令行是当步骤1~步骤2的资源不足时进一步判断占用ACL资源较多的业务,从而通过删除业务的方式释放资源。

步骤3 执行display system tcam resource命令,查看外扩TCAM的资源信息。

[~HUAWEI] Resource D		-	ım resource s	lot 3		
Slot Chip	TCAM	Service	Banks	Total	Used	Free

3	0	internal All	12		24576	566	24010
3	0	internal - ACL				560	
3	0	internal - UCv6Route				6	
3	0	internal - MCv4Route				0	
3	0	internal - MCv6Route				0	
3	1	internal All	12		24576	566	24010
3	1	internal - ACL				560	
3	1	internal - UCv6Route				6	
3	1	internal - MCv4Route				0	
3	1	internal - MCv6Route				0	
3	2	internal All	12		24576	566	24010
3	2	internal - ACL				560	
3	2	internal - UCv6Route				6	
3	2	internal - MCv4Route				0	
3	2	internal - MCv6Route				0	
3	3	internal All	12		24576	566	24010
3	3	internal - ACL				560	
3	3	internal - UCv6Route				6	
3	3	internal - MCv4Route				0	
3	3	internal - MCv6Route				0	
		o Tomplato Information:					
		e Template Information:					
Slo	t	Type Running	gTem	pla	te Next	Template	
3		CE-L24LQ-EA					

步骤4 执行display system tcam resource acl命令,查看三重内容寻址内存TCAM的资源信息。

Slot	t Ch	ip TCAM Resource Stage	Total	Used	Limited	Free
1	0	Internal Banks Ingress+Egress	12	2	2 10	
1	0	Internal Rules Ingress+Egress	24576	128	3968 2	20480
1	0	Internal Meters Ingress+Egress	65536	0	0 65	536
1	0	Internal Counters Ingress	16384	0	0 1638	4
1	0	Internal Counters Egress	2816	0	0 2816	j
1	1	Internal Banks Ingress+Egress	12	2	2 10	
1	1	Internal Rules Ingress+Egress	24576	128	3968 2	20480
1	1	Internal Meters Ingress+Egress	65536	0	0 65	536
1	1	Internal Counters Ingress	16384	0	0 1638	4
1	1	Internal Counters Egress	2816	0	0 2816	;

步骤5 执行display system tcam acl group resource命令,查看各业务使用资源的情况。

```
      <HUAWEI> display system tcam acl group resource slot 1

      STG: Stage
      KCP: Key Construction Program

      ING: Ingress
      EGR: Egress

      CYC: Cycle
      PTYPE: PortType

      FRT: Front Ports
      RCY: Recycle Ports

      16-L: 16bit-LSB Copy Engine
      16-M: 16bit-MSB Copy Engine

      32-L: 32bit-LSB Copy Engine
      32-M: 32bit-MSB Copy Engine

      F: Free
      T: Total

      Slot: 1 Chip: 0 UseRate:Normal

      STG KCP PacketType
      PTYPE
      CYC Group UsedKey 16-L 32-L 16-M 32-M

      F|T F|T F|T

      ING 1 L2 FRT 0 2 2,3 2|8 5|8 7|8 6|8

      ING 2 IPV4 FRT 0 3 2,3 0|8 4|8 4|8 6|8

      ING 3 TRILL FRT 0 1 2,3 6|8 5|8 7|8 6|8

      ING 4 IPV6 FRT 0 4 2,3 2|8 5|8 1|8 6|8
```

步骤6 执行display system tcam fail-record命令,查看TCAM下发失败的信息。

<huaw< th=""><th>'EI> display syster</th><th>m tcam fail-record</th><th></th></huaw<>	'EI> display syster	m tcam fail-record	
Slot Ch	ip Time	Service	ErrInfo
1 1 full	2019-03-24 06:40	D:11 Traffic Policy V	LAN Group resource
Total: 1			

步骤7 执行display system forwarding resource命令,查看芯片关键资源的占用情况。

[~HUAWEI] display system forwarding resource slot 4 Local Common Hardware Forwarding Tables: Slot Chip Name Total Remain Used[%] 4 LEM 262144 262142 2[0%] 1[0%] 4 0 - MAC - IP host 0[0%] 4 0 1[0%] 4 0 - ILM 262144 262136 8[0%] 4 LEM 1 - MAC 7[0%] 4 - IP host 0[0%] 1 4 - ILM 1[0%] 4 0 LPM 32768 32768 0[0%] 4 0 - IPv4 UC 0[0%] - IPv4 MC 0[0%] 4 0 - IPv6 UC 0[0%] 4 0 - IPv6 MC 0[0%] 0[0%] 4 LPM 32768 32768 - IPv4 UC 0[0%] 0[0%] 4 - IPv4 MC 1 - IPv6 UC 0[0%] - IPv6 MC 0[0%] 4 1

步骤8 执行display traffic-policy applied-record命令,查看流策略的应用记录。

<huawei> display tra Total records : 4</huawei>	affic-policy applied-re	cord	
Policy Type/Name	Apply Param	eter	Slot State
dsc	Global(IN)	1 f 2 fail(3) 4 fail(3)	
n4 p1	10GE4/0/2(IN) 10GE4/0/5(IN)	. 4	1 fail(4) 1 fail(4)

Fail reason:

- 3 -- The numbers of matched conditions and actions in the traffic policy exceed the limit.
- 4 -- Insufficient ACL resources.

步骤9 执行display system tcam acl resourcekey-buffer命令,查看当前ACL的KB资源使用信息。

<huawei> dis KB : Key Buffer Slot: 1</huawei>	splay system tcam ac	:l resou	ırce key	y-buffer verbo
Chip Directio UsedKBID	n ServiceName	(Group	КВТуре
0 Ingress Ingress Ingress Ingress Ingress	BPDU Deny CPCAR L2 Protocol Tunnel App-Session CPCAR	2 2 2 3 3	L2 L2 L2 IPv4 IPv4	2,3 2,3 2,3 2,3 2,3 2,3

Ingress	ECMP Hash	105	IPv4	4		
Ingress	LAG Hash	119	IPv4	4		
Ingress	Traffic Policy VLAN	294	IPv4	5		

----结束

□ 说明

框式设备用到的命令行参考上述CE12800例子,具体的支持情况参见CloudEngine系列交换机的产品文档。

盒式设备用到的命令行是: display system tcam resource acl、display system tcam fail-record、display system tcam bank resource、display system tcam service brief、display traffic-policy applied-record。

4.3 检查二层子接口资源是否充足

在M-LAG接入VXLAN场景下,M-LAG系统在BD(Bridge Domain)绑定VNI后,会在Peer-Link口创建隐式二层子接口。由于隐式二层子接口的数量与BD绑定VNI配置的数量相等,所以在堆叠切换M-LAG之前,需要查看设备上绑定VNI的BD数量和二层子接口的数量之和是否超过二层子接口规格,详细产品规格查询工具链接:https://support.huawei.com/onlinetoolweb/sqt/index。

若超过则要判断这些接口是否可以删除:

- 可以删除:通过删除接口释放规格保证设备上绑定VNI的BD数量和二层子接口的数量不超过规格。
- 不可以删除:不能进行堆叠到M-LAG的切换。

查看设备上绑定VNI的BD数量。

```
<HUAWEI>display bridge-domain binding-info
BDID VNI VSI
                         EVPN
2
     2
3
     3
4
5
6
7
     7
8
     8
9
10
     10
11
      11
12
```

查看设备上二层子接口的数量。

```
<HUAWEI>display current-configuration | in mode l2 interface Eth-Trunk10.1 mode l2 interface Eth-Trunk10.2 mode l2 interface Eth-Trunk10.3 mode l2 interface Eth-Trunk12.1 mode l2 interface Eth-Trunk200.1 mode l2
```

4.4 检查业务是否双归

如<mark>图4-1</mark>所示,堆叠系统正常工作时,现网中的业务分为双归和单归接入。双归业务流量通过堆叠两个成员设备转发,单归业务流量通过一台堆叠成员设备转发。两种业务在堆叠切换M-LAG时处理方式如下:

- 双归业务:若单个成员设备带宽能够承载整个网络业务,则可以进行堆叠切换M-LAG。
- 单归业务:若单归业务不能中断,需要先将单归业务移到备份设备,保证业务正常情况下再切换。

若单归业务可以中断,在切换之前可以先将业务中断,待切换M-LAG成功后再叠加业务。

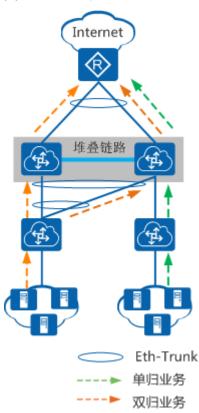


图 4-1 堆叠业务流量

4.5 检查 License 是否一致

框式设备通过命令行**display license** [**verbose**] **chassis** *chassis-number*分别查看堆叠成员设备的License。

盒式设备通过命令行**display license** [**verbose**] **slot** *slot-id*分别查看堆叠成员设备的 License,需要保证堆叠成员设备都已经加载License成功。

● 若其中一台设备无License,需要激活License文件,避免堆叠分裂后造成相应的业务无法下发的问题。

- 将License文件上传到设备上。 License申请和上传具体步骤可以参见License使用指南。
- 激活License文件,获取相应授权。 license active filename
- 若堆叠成员设备都没有License,则不需要关注此检查项。

#显示盒式设备的License文件信息。查看所有成员设备的License state是否都为 Normal_o

<HUAWEI> display license slot 0

slot 0:

Active License : flash:/CloudEngine7800.dat

License state : Normal Revoke ticket : No ticket

RD of Huawei Technologies Co., Ltd. Product name : CloudEngine 7800 Product version: V100R006

License Serial No: LIC201411261KSH50 Creator : Huawei Technologies Co., Ltd.
Created Time : 2014-11-26 09:09:51
Feature name : CELIC

Authorize type : demo Expired date : 2015-02-20 Trial days : -

Item name Item type Value Description CE-LIC-VXLAN Function YES CE-LIC-VXLAN

License state: Demo. The license for the current configuration will expire in 86 day(s).

Apply for authentic license before the current license expires.

显示框式设备的License文件信息。查看所有成员设备的License state是否都为 Normal_o

<HUAWEI> display license chassis 1/2

Active License : flash:/LICCloudEngine12800_V200R019_20190725UMKG5T.xml

License state : Normal Revoke ticket : No ticket

RD of Huawei Technologies Co., Ltd. Product name : CloudEngine 12800

Product version: V200R019

License Serial No: LIC20190725UMKG5T : Huawei Technologies Co., Ltd. Creator Created Time : 2019-07-25 14:36:20 Feature name : CELIC

Authorize type : comm Expired date : PERMANENT

Trial days : 60

Item name Item type Value Description

CE128-LIC-IPV6 -- 1 CloudEngine 12800 IPv6 Function

CE128-LIC-IPv6 Function YES CE128-LIC-IPv6

CE12800 Telemetry Function CE128-LIC-TLM -- 1

CE-LIC-TLM CE128-LIC-VS Function YES -- 1 CE-LIC-TLM

CloudEngine 12800 Virtual System Function

DE0S0000VS01 Function YES CE128-LIC-VS

Master board license state: Trial. The trial days remains 60 day(s). Apply for authentic license before the current license expires.

4.6 检查动态路由协议建立的邻居数量是否超限

首先检查堆叠系统是否存在动态路由协议,其次检查已经存在的动态路由协议建立的动态邻居数量是否超限。

- 由于V200R019C10版本及其之前版本不支持动态路由协议接入M-LAG口,因此需要检查堆叠的所有跨堆叠成员的Eth-trunk聚合口是否存在动态路由协议。
 - a. 通过命令行查询动态路由信息,目前动态路由协议包括RIP、OSPF、BGP和IS-IS。

执行**display rip neighbor**查看RIP的邻居信息。其中Number of RIP routes表示RIP的邻居个数。

执行display ospf peer查看OSPF中各区域邻居的信息。其中Total number of peer(s)表示OSPF的邻居个数。

执行display ospfv3 peer查看OSPFv3邻居信息。其中Total number of peer(s)表示OSPFv3的邻居个数。

```
      <HUAWEI> display ospfv3 1 peer vlanif 10

      OSPFv3 Process (1)

      Total number of peer(s): 1
      //OSPFv3邻居个数

      Peer(s) in full state: 1

      OSPFv3 Area (0.0.0.0)

      Neighbor ID
      Pri State
      Dead Time Interface
      Instance ID

      10.1.1.1
      1 Full/ -
      00:00:30 Vlanif10
      0
```

执行display isis peer用来查看IS-IS的邻居信息。其中Total Peer(s)表示IS-IS的邻居个数。

```
| Comparison | Co
```

执行display bgp peer用来查看BGP对等体信息。其中Total number of peers表示BGP的邻居个数。

```
**HUAWEI> display bgp peer

Status codes: * - Dynamic

BGP local router ID : 10.2.3.4

Local AS number : 10

Total number of peers : 2 //BGP邻居个数

Peers in established state: 1

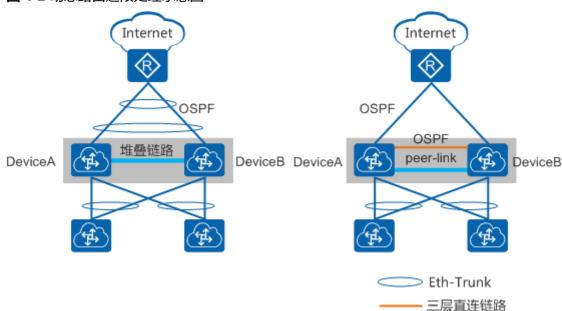
Total number of dynamic peers: 0

Peer V AS MsgRcvd MsgSent OutQ Up/Down State PrefRcv
```

```
10.1.1.1 4 100 0 0 0.00:00:07 Idle 0
10.2.5.6 4 200 32 35 0.00:17:49 Established 0
```

- b. 如果存在上述任何一种动态路由,则不能直接继承到M-LAG场景,可以修改动态路由为静态路由(ip route-static/ipv6 route-static)方式。静态路由具体配置参见《配置指南-IP单播路由》里的静态路由配置。
- 如图4-2所示,假设堆叠设备和上行网络建立了OSPF邻居关系且上行网络ospf邻居满规格。将堆叠切换为M-LAG相当于将一台设备切换成两台设备,因此M-LAG的两台设备和上行网络间就需要建立超规格的邻居。为了避免因为超规格造成邻居建立失败,需要通过建立一条三层直连链路,将动态路由邻居搬迁到M-LAG成员设备间来实现。同时,新建的链路带宽需要满足容量。

图 4-2 动态路由超限处理示意图



4.7 检查 IPv6 业务是否存在

检查堆叠设备的版本和IPv6业务,若设备版本在V200R005C10之前且设备存在IPv6业务,如DHCPv6业务等,需要先将堆叠设备升级到V200R005C10或之后版本(堆叠升级步骤参见《配置指南-虚拟化》里的升级堆叠),再进行堆叠切换M-LAG的步骤,否则IPv6业务在M-LAG组网里不会生效。

查看设备的版本信息

<HUAWEI> display version
Huawei Versatile Routing Platform Software
VRP (R) software, Version 8.200 (CloudEngine 6800 V200R005C10)
Copyright (C) 2012-2020 Huawei Technologies Co., Ltd.
......

查看DHCPv6 Server业务。此处仅以DHCPv6业务为例,具体业务可以查看<mark>配置文</mark> 档。

<HUAWEI> display dhcpv6 pool pool1
DHCPv6 pool: pool1
Address prefix: FC00:2::/64
Lifetime 172800 seconds, preferred 86400 seconds
100 in use, 0 conflicts //存在DHCPv6 Server业务
Information refresh time: 86400
DNS server address: FC00:2::3

DNS server domain name: huiwei.com Conflict-address expire-time: 172800

renew-time-percent: 50 rebind-time-percent: 80 Active normal clients: 0

4.8 检查跨设备的二层端口隔离配置是否生效

通过命令行display port-isolate group查看堆叠跨设备的二层端口隔离组的配置信息。

查看所有二层端口隔离组的配置。

<HUAWEI> display port-isolate group all

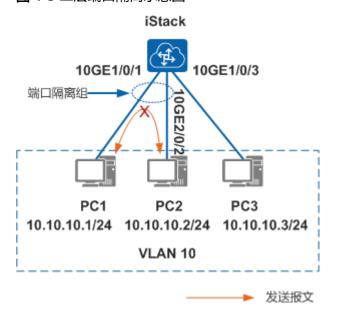
The ports in isolate group 1: 10GE**1**/0/1 10GE**2**/0/2

二层端口隔离是指相同VLAN内的二层端口不能互访,如<mark>图4-3</mark>所示,PC1、PC2和PC3 同属于VLAN10,将PC1与PC2对应的端口10GE1/0/1和10GE2/0/2加入端口隔离组后,PC1与PC2在VLAN10内不能互相访问。

然而, M-LAG两台主备设备默认二层端口是不隔离的, 可以通过ACL或VLAN隔离。

- ACL隔离: PC1流量通过ACL隔离无法转发到PC2,此时会占用ACL资源。
- VLAN隔离: PC1与PC2端口加入不同的VLAN进行隔离。但当组网环境复杂,隔离的端口庞大时,VLAN隔离的成本将上升。

图 4-3 二层端口隔离示意图



4.9 检查跨设备的单向业务是否存在

跨设备的单向业务主要包括跨设备重定向、跨设备镜像,如果存在这两种业务,需要 修改流量模型,让流量从单设备出。

- 查看跨设备重定向配置
 - a. 查看基于MQC的跨设备重定向配置 # 查看流策略的应用记录

[~SwitchA] display traffic-policy applied-record

Total records : 1

Policy Type/Name Apply Parameter Slot State

10GE1/0/1(IN) 1 success

查看流策略的配置信息

<SwitchA> display traffic policy
Traffic Policy Information:
Policy: p1
Classifier: c1
Type: OR
Behavior: b1
Redirect:
Redirect interface
10GE2/0/3

Total policy number is 1

□ 说明

当Apply Parameter和Redirect interface在不同的堆叠成员设备上,则认为存在跨设备的重定向。

b. 查看基于ACL的跨设备重定向配置

查看流策略的应用记录

查看流策略的配置信息

[~SwitchA] interface 10ge 1/0/1 [~SwitchA-10GE1/0/1] display this | in traffic-redirect

[~SwitchA-10GE1/0/1] traffic-redirect acl 4001 interface 10ge 2/0/3 inbound

□ 说明

当Apply Parameter和traffic-redirect的重定向接口在不同的堆叠成员设备上,则认为存在跨设备的重定向。

• 查看跨设备镜像配置

a. 查看设备上配置的观察端口。

查看设备上配置的观察端口

<HUAWEI> display observe-port
-----Index : 1
Slot: 1

Interface: 100GE1/0/3

b. 查看设备上镜像的配置信息。

查看设备上镜像的配置信息

<HUAWEI> display port-mirroring
Observe port mirroring:

MirroringPort Direction ObservePort

100GE4/0/10 Inbound 1

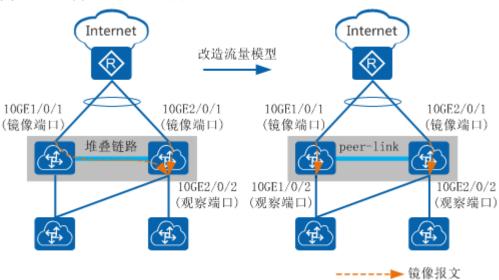
Traffic mirroring:

TrafficBehavior	ObservePort
b	1

当观察端口和镜像端口在不同的堆叠成员设备上,则认为存在跨设备的镜像。

以跨设备镜像为例,如<mark>图4-4</mark>所示,当观察口和镜像端口不在同一台设备上时,通过下发命令行**observe-port**将观察端口拆开的形式,解除跨设备的镜像流量,避免在堆叠切换M-LAG后跨设备的镜像流量丢失。

图 4-4 跨设备镜像改造示意图



4.10 检查系统 MAC 是否为堆叠主设备

执行命令display system mac-address,通过Stack MAC Information字段检查系统 MAC是否为堆叠主设备。如果有Stack MAC Information字段,认为系统MAC不是堆叠主设备,需要通过命令行undo set system mac-address恢复设备的MAC地址为出厂值。目的是防止在系统MAC不是堆叠主设备的情况下,堆叠分裂会造成堆叠两台设备的系统MAC地址冲突。

杏麦惟叠设备系统MAC.

		MAC.			
[HUA\	WEI]display system	mac-addre	S		
Currer	nt System MAC add	ress: 8446-fe	1-dd20(Used Stack	MAC)	
Currer	nt System MAC num	nber : 16			
User-c	configured MAC add	dress:			
User-c	configured MAC nur	mber :			
Systen	n MAC Switch :	Enable			
Systen	n MAC Switch-delay	y: 10(minute:)		
Systen	n MAC Inconsistenc	e-alarm :	Enable		
Systen	n MAC Inconsistenc	e-alarm Dela	y: 10(minutes)		
Manu	facture MAC Inform	ation:			
Slot	MAC	Number			
1	8446-fea1-dd20	 16			
		 16			
1	8446-fea1-dd20 8446-fea1-e480	 16			
1	8446-fea1-dd20 8446-fea1-e480	16 16			

Slot	MAC	Number
	8446-fea1-dd20	16
-		
2	8446-fea1-dd20	16

修改堆叠的系统MAC地址,不需要重启,系统MAC会更新,业务会有1~2个报文丢包,并且触 发对端设备的ARP更新。

4.11 检查是否是 SDN 网络

当AC控制器对设备进行控制和管理时(软件定义网络(Software Defined Network, SDN)),若直接修改配置将一台堆叠设备拆成两台M-LAG成员,AC控制器将无法识 别这两台M-LAG设备,即无法控制和管理M-LAG系统。需按照以下步骤进行改造:

- 1. 将全部业务迁移到备份设备上。
- 堆叠切换M-LAG后将M-LAG成员设备都加入AC控制器。
- 再重新通过AC控制器下发业务。

4.12 记录 MAC 表项和流量情况

堆叠切换M-LAG前需要检查MAC表项、端口流量和路由信息,与切换后的信息进行对 比。路由查询可以参见4.6 检查动态路由协议建立的邻居数量是否超限章节。通过以下 命令查看MAC表项和端口流量信息。

查看MAC地址表项。

```
<HUAWEI> display mac-address
Flags: * - Backup
BD : bridge-domain Age : dynamic MAC learned time in seconds
MAC Address VLAN/VSI/BD Learned-From
0000-0000-0033 100/-/- 10GE1/0/1
                                       dynamic 4294367295
0000-0000-0001 200/-/- 10GE1/0/2
                                       static
Total items: 2
```

# 端口流量信息。	,					
<huawei> display ir</huawei>	iterface brief					
PHY: Physical						
*down: administrative	ly down					
^down: standby						
(l): loopback						
(s): spoofing						
(b): BFD down						
(e): ETHOAM down						
(d): Dampening Supp	ressed					
(p): port alarm down						
(dl): DLDP down						
(c): CFM down						
(sd): STP instance dis-	carding					
InUti/OutUti: input ut	ility rate/output u	itility ra	te			
Interface	PHY Protocol In	Uti Out	Uti in[Frrors out	Errors	
10GE1/0/1	down down	0%	0%	0	0	
10GE1/0/2	down down	0%	0%	0	0	
10GE1/0/3	down down	0%	0%	0	0	
10GE1/0/4	down down	0%	0%	0	0	
40GE1/0/1	down down	0%	0%	0	0	

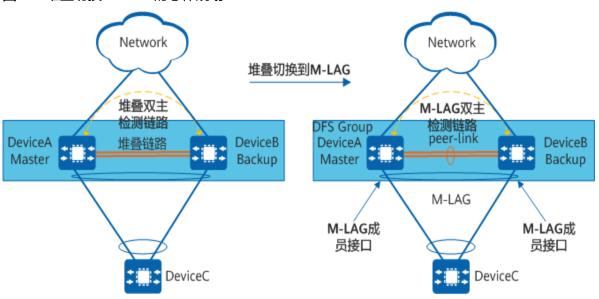
40GE1/0/2	down down	0%	0%	0	0
Eth-Trunk0	down down	0%	0%	0	0
GE1/0/1	up up 0.0	1% 0.0	1%	0	0
GE1/0/2	up up 0.0	1% 0.0	1%	0	0
GE1/0/3	down down	0%	0%	0	0
GE1/0/4	down down	0%	0%	0	0
GE1/0/5	down down	0%	0%	0	0
GE1/0/6	up up 0.0	1% 0.0	1%	0	0
GE1/0/7	up up 0.0	1% 0.0	1%	0	0
GE1/0/8	down down	0%	0%	0	0
GE1/0/9	down down	0%	0%	0	0
GE1/0/10	down down	0%	0%	0	0
More					

5 五步完成堆叠切换为 M-LAG

堆叠切换M-LAG的总体规划如图5-1所示。

- 1. 删除堆叠双主检测链路。
- 2. 将堆叠备设备从原堆叠中隔离,给堆叠备设备配置M-LAG。
- 3. 将堆叠主设备的业务切换到原堆叠备设备,堆叠主设备隔离并配置M-LAG。
- 4. 两台设备组建M-LAG系统。
- 5. 检查切换后状态及业务恢复情况。

图 5-1 堆叠切换 M-LAG 的总体规划



- 5.1 删除堆叠双主检测链路
- 5.2 隔离堆叠备设备,增加M-LAG配置
- 5.3 切换业务
- 5.4 隔离堆叠主设备,增加M-LAG配置
- 5.5 查询切换后状态,切换前后业务对比

5.1 删除堆叠双主检测链路

堆叠的双主检测链路是用来检测和处理堆叠分裂的协议,可以实现堆叠分裂的检测、冲突处理和故障恢复,降低堆叠分裂对业务的影响。由于在堆叠正常工作时双主检测链路故障不会对业务造成影响,因此堆叠切换M-LAG时,可先将堆叠双主检测链路删除。删除双主检测链路可以防止堆叠分裂后双主检测功能生效造成备设备的业务口被Error-down,另外还可以防止在拆除堆叠时产生的告警。

堆叠的双主检测链路可以配置在业务口、Eth-trunk口、管理网口或堆叠端口,删除双主检测链路的步骤如下。

• 删除业务口直连方式的双主检测

system-view interface interface-type interface-number undo dual-active detect mode direct quit

● 删除Eth-Trunk口代理方式双主检测

system-view interface eth-trunk *trunk-id* undo dual-active detect mode relay quit

• 删除管理网口方式双主检测

system-view interface meth 0/0/0 undo dual-active detect enable quit

删除堆叠端口方式双主检测

system-view interface stack-port member-id/port-id undo dual-active detect mode direct auit

5.2 隔离堆叠备设备,增加 M-LAG 配置

本节详细描述隔离堆叠的备设备后在备设备上的配置。隔离堆叠备设备,通过在所有业务端口执行shutdown命令行实现。按照如下步骤增加M-LAG配置,建议直接清空重新配置。

1. shutdown堆叠备所有业务端口

system-view interface interface-number shutdown

2. 增加DFS-GROUP配置。

system-view dfs-group dfs-group-id source ip ip-address [vpn-instance vpn-instance-name] [peer peer-ip-address [udp-port port-number]] [timeout seconds] //配置双主检测链路。quit

□ 说明

选择合适的链路作为M-LAG的双主检测链路,可以是管理口、直连线、上行互通链路等。

3. 配置STP模式

system-view stp mode rstp stp v-stp enable stp root primary stp bridge-address mac-address

当设备的角色是根桥时,需要配置stp root primary和stp bridge-address *macaddress*。

4. 配置peer-link端口

- 若堆叠线缆组建的堆叠,需要拆除堆叠线缆,用业务口普通连线作为peer-link端口。

system-view

interface eth-trunk trunk-id

peer-link peer-link-id //配置peer-link

auit

- 复用之前的堆叠链路作为peer-link端口。

system-view

interface eth-trunk trunk-id

peer-link peer-link-id //配置peer-link

quit

5. 配置原堆叠跨设备的Eth-Trunk口为M-LAG口。

system-view

interface eth-trunk trunk-id

mode { lacp-static | lacp-dynamic }

dfs-group *dfs-group-id* **m-lag** *m-lag-id* //绑定DFS Group和用户侧Eth-Trunk接口,即配置为M-LAG成员接口。

quit

6. 配置双活网关的虚拟IP地址和虚拟MAC地址。

配置IPv4双活网关:

system-view

interface { vlanif vlan-id | vbdif bd-id }

ip address ip-address { mask | mask-length } [sub]

mac-address mac-address

quit

配置IPv6双活网关:

system-view

interface { vlanif vlan-id | vbdif bd-id }

ipv6 enable

ipv6 address { ipv6-address prefix-length | ipv6-address | prefix-length } [eui-64]

mac-address mac-address

auit

□ 说明

M-LAG接入VXLAN网络或三层网络时该步骤是必须的。

7. 配置Monitor Link

system-view

monitor-link group group-id

port interface-type interface-number { downlink [downlink-id] | uplink }

quit

8. 修改OSPF、BGP的router-id

system-view

ospf [process-id | router-id | vpn-instance vpn-instance-name] *

□ 说明

原堆叠系统的OSPF、BGP协议的**router-id**只有一个,堆叠切换M-LAG后,需要在M-LAG两个成员设备上配置不同的router-id。同样,在网络侧的**router-id**也需要修改为两个不同的值,分别和M-LAG两个成员建立动态路由。

9. 修改设备本地SNMP实体的引擎ID

system-view

undo snmp-agent local-engineid

为了在切换M-LAG后网管分别对接上M-LAG成员设备,可以删除设备本地SNMP实体的引擎ID,设备自动生成新的local-engineid。

5.3 切换业务

在上一节中已经改造完成M-LAG备设备的配置,现需要将业务从堆叠主设备切到原堆叠备设备,具体步骤如下:

1. 查看并记录堆叠主设备的MAC表项,端口流量信息和路由信息。

display mac-address display interface brief display rip neighbor display ospf peer display ospfv3 peer display isis peer display bgp peer

2. shutdown堆叠主的所有端口。

system-view interface interface-type interface-number shutdown quit

3. 执行完成步骤2后,立即执行命令undo shutdown,恢复原堆叠备设备的端口。

system-view

interface interface-type interface-number undo shutdown quit

□ 说明

如果堆叠设备是根桥,则整网存在STP收敛。

4. 验证切换后的业务是否正常

display mac-address display interface brief display rip neighbor display ospf peer display ospfv3 peer display isis peer display bgp peer

查看上述命令行,和步骤1的结果比较。若业务不正常,执行shutdown原堆叠备设备,undo shutdown堆叠主设备,流量恢复到切换之前的状态并排查配置。

5.4 隔离堆叠主设备,增加 M-LAG 配置

隔离堆叠主设备。需要执行命令行**shutdown**堆叠主设备所有业务端口。按照如下步骤增加M-LAG配置,建议直接清空重新配置。M-LAG的具体配置请参见**CloudEngine系**列交换机的产品文档。

1. shutdown堆叠主所有业务端口

system-view interface interface-type interface-number shutdown

2. 增加DFS-GROUP配置。

system-view dfs-group dfs-group-id source ip ip-address [vpn-instance vpn-instance-name] [peer peer-ip-address [udp-port port-number]] [timeout seconds] //配置双主检测链路。quit

选择合适的链路作为M-LAG的双主检测链路,可以是管理口、直连线、上行互通链路等。

3. 配置STP模式

system-view stp mode rstp stp v-stp enable stp root primary stp bridge-address mac-address

□ 说明

当设备的角色是根桥时,需要配置stp root primary和stp bridge-address *macaddress*。

- 4. 配置peer-link端口
 - 若堆叠线缆组建的堆叠,需要拆除堆叠线缆,用业务口普通连线作为peerlink端口。

system-view interface eth-trunk trunk-id peer-link peer-link-id //配置peer-link quit

- 复用之前的堆叠链路作为peer-link端口。

system-view interface eth-trunk trunk-id peer-link peer-link-id //配置peer-link quit

5. 配置原堆叠跨设备的Eth-Trunk口为M-LAG口。

system-view interface eth-trunk trunk-id mode { lacp-static | lacp-dynamic } dfs-group dfs-group-id m-lag m-lag-id |/绑定DFS Group和用户侧Eth-Trunk接口,即配置为M-LAG成员接口。 quit

6. 配置双活网关的虚拟IP地址和虚拟MAC地址。

配置IPv4双活网关:

system-view interface { vlanif vlan-id | vbdif bd-id } ip address ip-address { mask | mask-length } [sub] mac-address mac-address quit

配置IPv6双活网关:

system-view interface { vlanif vlan-id | vbdif bd-id } ipv6 enable ipv6 address { ipv6-address prefix-length | ipv6-address prefix-length } [eui-64] mac-address mac-address quit

□ 说明

M-LAG接入VXLAN网络或三层网络时该步骤是必须的。

7. 配置Monitor Link

system-view
monitor-link group group-id
port interface-type interface-number { downlink [downlink-id] | uplink }
quit

8. 修改OSPF、BGP的router-id

system-view

ospf [process-id | router-id | vpn-instance vpn-instance-name] *

□说明

原堆叠系统的OSPF、BGP协议的router-id只有一个,堆叠切换M-LAG后,需要在M-LAG 两个成员设备上配置不同的router-id。同样,在网络侧的router-id也需要修改为两个不同 的值,分别和M-LAG两个成员建立动态路由。

修改设备本地SNMP实体的引擎ID

system-view

undo snmp-agent local-engineid

□ 说明

为了在切换M-LAG后网管分别对接上M-LAG成员设备,可以删除设备本地SNMP实体的引 擎ID,设备自动生成新的local-engineid。

10. M-LAG主备成员的配置完成后,undo shutdown已经配置完M-LAG主设备的端 口,M-LAG系统组建完成。

system-view

interface interface-type interface-number

undo shutdown

quit

5.5 查询切换后状态,切换前后业务对比

为确认M-LAG系统是否正常工作,需要查询M-LAG的心跳状态和主备状态以及M-LAG 口的状态。

查询M-LAG心跳状态及主备状态。

[~HUAWEI] display dfs-group 1 m-lag

: Local node Heart beat state: OK

Node 1 *

Dfs-Group ID : 1 Priority: 150

Address : ip address 10.3.3.4

State : Master

Causation

System ID : 0025-9e95-7c31 SysName : HUAWEIA : V200R019C10 Version Device Type : CE6850EI

Node 2

Dfs-Group ID : 1 Priority: 120

: ip address 10.3.3.3 Address

State : Backup

Causation : -

System ID : 0025-9e95-7c11 SysName : HUAWEIB Version : V200R019C10 Device Type : CE6850EI

查看设备上的M-LAG信息。

[~HUAWEI] display dfs-group 1 node 1 m-lag brief

* - Local node

M-Lag ID Interface Port State Status Consistency-chec 1 Eth-Trunk 10 **Up** active(*)-active

Failed reason:

- 1 -- Relationship between vlan and port is inconsistent
- 2 -- STP configuration under the port is inconsistent
- 3 -- STP port priority configuration is inconsistent
- 4 -- LACP mode of M-LAG is inconsistent

- 5 -- M-LAG configuration is inconsistent
- 6 -- The number of M-LAG members is inconsistent

除查看M-LAG状态外,还需要查看业务流量是否恢复,可以通过查询MAC、网络侧路由和总流量,同4.12 记录MAC表项和流量情况章节记录的信息比较。