

H3C交换机技术专题篇

交换那些事儿M-LAG技术专题新系列
M-LAG典型案例及FAQ
每期推送2个案例+4个FAQ
让你在遇到类似问题时不迷茫不发愁
解决方案手到擒来~

本期内容

典型案例:

- Case1 M-LAG+VRRP服务器主备网卡接入流量泛洪
- Case2 M-LAG peer-link故障时主设备被MAD DOWN

FAQ:

- Q1 M-LAG组网需要修改mac老化时间吗?
- Q2 园区核心设备M-LAG组网,单挂接口是否只能配置为trunk接口?
- Q3 M-LAG设备健康值要检查哪些?
- Q4 M-LAG设备间可以同步静态mac表项吗?

DRNI更名为M-LAG

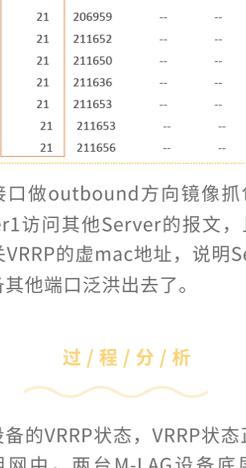
原DR口>>M-LAG接口
原DR group>>M-LAG group
原IPL链路>>peer-link链路
原Ipp口>>peer-link接口
本文涉及命令行暂不做修改

Case1

M-LAG+VRRP服务器主备网卡接入流量泛洪

组 / 网 / 说 / 明

设备及版本: S6850 R6635
组网: 两台S6850配置M-LAG+VRRP作为服务器的网关,服务器采取主备网卡方式分别接在两台M-LAG设备上, M-LAG设备上没有配置M-LAG聚合口。



问 / 题 / 描 / 述

现场 M-LAG+VRRP 组网, 两台 S6850 起 M-LAG+VRRP 作为服务器网关, 服务器采取主备网卡方式分别连接两台 M-LAG 设备, M-LAG 设备上不起 M-LAG 聚合口, 相当于服务器单挂接入。现场开局时发现服务器使用主网卡访问时流量正常, 使用备网卡访问时, Server1 访问其他 Server 不通, 且 M-LAG 从设备上其他还没上业务的接口出方向有大量泛洪流量。

查看从设备接口流量, 发现其他接口 WGE1/0/1-WGE1/0/17 入方向几乎没有流量:

```
=====display counters rate inbound interface=====
Usage: Bandwidth utilization in percentage
Interface Usage (%) Total (pps) Broadcast (pps) Multicast (pps)
WGE1/0/1 0 0 -- --
WGE1/0/2 0 0 -- --
WGE1/0/3 0 0 -- --
WGE1/0/4 0 0 -- --
WGE1/0/5 0 0 -- --
WGE1/0/6 0 0 -- --
WGE1/0/7 0 0 -- --
WGE1/0/8 0 516 -- --
WGE1/0/9 0 0 -- --
WGE1/0/10 0 0 -- --
WGE1/0/17 0 0 -- --
```

但是这些接口的出方向流量都达到了21%:

```
=====display counters rate outbound interface=====
Usage: Bandwidth utilization in percentage
Interface Usage (%) Total (pps) Broadcast (pps) Multicast (pps)
WGE1/0/1 21 206896 -- --
WGE1/0/2 21 206903 -- --
WGE1/0/3 21 206854 -- --
WGE1/0/4 21 206876 -- --
WGE1/0/5 21 206959 -- --
WGE1/0/6 21 211652 -- --
WGE1/0/7 21 211650 -- --
WGE1/0/8 21 211636 -- --
WGE1/0/9 21 211653 -- --
WGE1/0/10 21 211653 -- --
WGE1/0/17 21 211656 -- --
```

在这些接口做outbound方向镜像抓包, 发现报文均为Server1访问其他Server的报文, 且报文目的mac均为网关VRRP的虚mac地址, 说明Server1发出的报文从设备其他端口泛洪出去了。

过 / 程 / 分 / 析

查看从设备的VRRP状态, VRRP状态正常, 在M-LAG+VRRP组网中, 两台M-LAG设备底层都会下发VRRP虚mac, 也就是说从设备收到目的mac为VRRP虚mac地址的报文时应该直接走三层转发, 不应该泛洪:

```
=====display vrrp verbose=====
IPv4 Virtual Router Information:
Running mode: Standard
Total number of virtual routers: 5
Interface Vlan-interface100
VRID : 100 Adver Timer : 100
Admin Status : Up State : Backup
Config Pri : 100 Running Pri : 100
Preempt Mode : Yes Delay Time : 0
Become Master : 2940ms left
Auth Type : None
Virtual IP : 10.10.1.254
Virtual MAC : 0000-5e00-0001
Master IP : 10.10.1.1
```

那为什么从设备会将这部分流量泛洪呢? 我们在两台M-LAG设备上分别查看设备底层虚mac下发情况, 发现M-LAG主设备上底层下发了VRRP虚mac地址, 而M-LAG从设备的底层则没有该硬件表项:

```
M-LAG主设备表项:
[probe]debug l3intf-drv show virtualmac 0 slot 1
*****
-L3INTF Info by index Slot 1
*****
- IntfId: 100
- VlanId: 100
- VmacStatus: 1
- RefCount: 1
- Vmac: 0000-5e00-0001
M-LAG从设备表项:
[probe]debug l3intf-drv show virtualmac 0 slot 1
*****
-L3INTF Info by index Slot 1
*****
- IntfId: 0
- VlanId: 0
- VmacStatus: 0
- Vmac: 0x00-00-00-00-00-00
*****
```

从设备底层没有VRRP虚mac的硬件表项, 当从设备收到Server1备网卡发过来的目的mac为VRRP虚mac的报文时, 不认为目的mac是本地mac, 因此没有走三层转发流程, 而是走了二层转发; 由于查不到该mac表项, 就二层泛洪了。

那为什么从设备上底层没有下发VRRP虚mac表项呢? 现场服务器为主备网卡接入, 相当于纯单挂接入, 设备上也没有配置任何M-LAG聚合口。当前设备在纯单挂场景下存在如下限制:

服务器主备接入 (bond1), 并且配置VRRP网关的VLAN没有配置M-LAG接口时, 需要创建一个M-LAG聚合组, 并配置M-LAG聚合接口允许该VLAN通过。

例如: VLAN接口100配置了VRRP, 但是没有任何M-LAG聚合接口加入VLAN 100, 此时VLAN 100会出现预期外的未知单播流量泛洪, 泛洪流量太大则可能导致丢包。

解 / 决 / 方 / 法

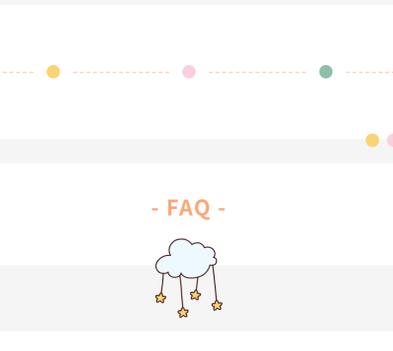
根据设备限制, 创建一个M-LAG聚合接口, 并配置聚合接口允许VLAN 100通过。后续我们计划开发一条专用命令来优化这个M-LAG纯单挂场景的限制。

Case2

M-LAG peer-link故障时主设备被MAD DOWN

组 / 网 / 说 / 明

设备及版本: S12504G-AF R7624P15
组网: 两台S12504G-AF起M-LAG, 服务器主备网卡分别接M-LAG两台设备, 没有配置M-LAG接口 (和上文一样属于纯单挂场景)。下图为局部组网图:



问 / 题 / 描 / 述

现场 M-LAG 组网, 未起 Vlan 双活或者 VRRP 网关, 服务器主备网卡分别接在 M-LAG 两台设备, 相当于纯单挂接入场景。现场 M-LAG-1 为配置的主设备, M-LAG-2 为配置从设备。但由于 M-LAG-1 之前重启过, 发生过一次主备倒换, 实际生效的角色 M-LAG-2 为主设备, M-LAG-1 为从设备, 因此实际流量到达服务器网卡。



过 / 程 / 分 / 析

由于服务器主备网卡发生了倒换, 与现场工程师确认, 服务器主备倒换时自身需要对业务表项进行恢复, 会导致业务中断。因此现场业务报障的直接原因是服务器主备网卡切换, 主备网卡切换的原因又是因为 M-LAG-2 的业务口被 MAD DOWN 了。

此时有一个疑问, 现场明明 M-LAG-2 是实际主设备, 为什么 peer-link 故障时 MAD DOWN 的却是 M-LAG-2 设备?

这里我们再来详细研究一下 M-LAG 的故障处理机制, 当 peer-link 故障而 keepalive 链路正常时, 确实应该 MAD DOWN 从设备的接口, 但这个从设备是当时计算出的从设备, 而不是故障发生之前的从设备。

我们仔细去查阅官网配置指导, 可以看到如下说明:

```
M-LAG角色计算触发条件包括:
• M-LAG设备在系统初始化时(包括新配置M-LAG或带M-LAG配置重启设备)。
• peer-link链路UP时,设备角色通过peer-link计算。
• peer-link故障,Keepalive正常工作,设备角色通过Keepalive链路计算。
• peer-link和Keepalive链路均故障,根据本端M-LAG设备上M-LAG接口状态决定设备角色。
```

划重点了朋友们, peer-link故障, keepalive 正常工作时, 是会触发 M-LAG 角色计算的。也就是说, 现场 peer-link 故障后, M-LAG-1 和 M-LAG-2 设备的主从关系要通过 keepalive 链路重新计算, 角色计算规则在配置手册中也有详细说明:

- 当通过 peer-link 或 Keepalive 链路交互报文计算设备角色时, 依次比较如下因素:
- (1) 比较设备所有 M-LAG 接口的状态, 有可工作 M-LAG 接口的一端为优;
 - (2) 比较计算前角色, 若有一端为 Primary, 另一端为 None, 则 Primary 端优;
 - (3) 比较 MAD DOWN 状态, 若一端存在处于 MAD DOWN 状态的接口, 另一端不存在处于 M-LAG MAD DOWN 状态的接口, 则不存在处于 M-LAG MAD DOWN 状态的接口的一端优;
 - (4) 比较设备健康状态, 健康值越小越优。设备的健康值可通过 display system health 命令查看, 健康值越小设备越健康, 设备无故障运行时, 健康值为 0;
 - (5) 比较设备角色优先级, 越高越优;
 - (6) 比较设备桥 MAC, 越小越优。

根据这个计算规则, 我们依次来看一下: 全单挂场景, 两台设备均无可用 M-LAG 口, 所以第 1 条 pass;

在故障前不存在为 none 角色的设备和被 MAD DOWN 的设备, 因此第 2 条、第 3 条规则也不满足;

由于故障前没有查看过设备健康度, 所以第 4 条存疑, 但故障前整个系统是正常运行的, 大概率两台设备都没有问题, 即健康度都为 0;

M-LAG-1 是配置的主设备, 角色优先级配置得比 M-LAG-2 高, 因此第 5 条 M-LAG-1 胜出;

由于第 5 条已经可以计算出两台设备的角色, 第 6 条就没有比较的必要了。

从上述分析来看, 现场 peer-link 故障后重新进行角色计算, 很大可能 M-LAG-1 会被计算为主设备, 那么 M-LAG-2 作为重新计算出的备设备被 MAD DOWN 就是正常的。

从以上信息可知, M-LAG 的 MAD DOWN 实际会按照角色计算规则选举出一个从设备进行 MAD DOWN, 而非简单的将故障前实际生效的角色为从设备的接口给 MAD DOWN。

解 / 决 / 方 / 法

M-LAG 本身就是跨设备链路聚合, 若业务侧切换时会存在恢复表象等导致业务中断的情况, 建议及时修改组网为 M-LAG 聚合的形式; 如果确有 M-LAG 全单挂的业务需求, 建议尽可能保证业务主备与 M-LAG 配置的主从设备保持一致, 同时注意前文所说的限制---创建一个 M-LAG 聚合接口, 并配置聚合接口允许相应 VLAN 通过。

- FAQ -

Q1: M-LAG组网需要修改mac老化时间吗?

A1: 在M-LAG组网环境中, 如果设备存在大量MAC地址表项, 请通过 mac-address timer aging 命令增加MAC地址老化时间, 建议配置MAC地址老化时间在20分钟以上。

Q2: 园区核心设备M-LAG组网,单挂接口是否只能配置为trunk接口?

A2: 原本单挂接口需要与peer-link接口类型保持一致, peer-link一般为trunk口, 因此单挂接口也只能配置为trunk接口。但当前园区核心设备M-LAG推荐的R7624Pxx及之后的版本, 单挂接口可以使用access接口, 无需与peer-link保持一致。

Q3: M-LAG设备健康值要检查哪些?

A3: M-LAG角色计算时需要考虑设备健康值, 健康值是指设备会定期执行一系列检查, 根据检查项状态, 判断设备的健康状态。健康值检查项包括: 内存检测、进程异常检测、芯片堵塞检测、CPU死循环检测、MMU表项错误检测、单板状态检测、单板间HG检测等等, 可以通过display system health查看具体检查项:

表1-18 display system health命令显示健康值列表

| 字段 | 描述 |
|--------------------|---------------|
| Health Normal | 健康值正常, 取值为: 0 |
| Memory | 内存使用 |
| Communication | 通信链路 |
| OSN monitor | OSN板卡健康值 |
| Chip fan | 芯片风扇故障检测 |
| CPU load | CPU负载检测 |
| Forwarding channel | 转发通道健康值 |
| HG | 单板间HG检测 |
| External TCM | 板间外部流量控制检测 |
| LSW CPU | LSW CPU健康值 |
| MSU CPU | MSU CPU健康值 |
| LSW port | LSW端口健康值 |
| Board status | 单板健康值 |
| Fan status | 风扇健康值 |
| Temperature status | 温度健康值 |

Q4: M-LAG设备间可以同步静态mac表项吗?

A4: 两端M-LAG设备之间不会同步静态MAC地址表项和黑洞MAC地址表项。

每多处理一个问题
对M-LAG的理解也更多一层
经验案例多看看
典型配置多学习
大家都可以成为大佬~
有问题欢迎留言哦~

END

扫码关注我们哦