

组网及说明

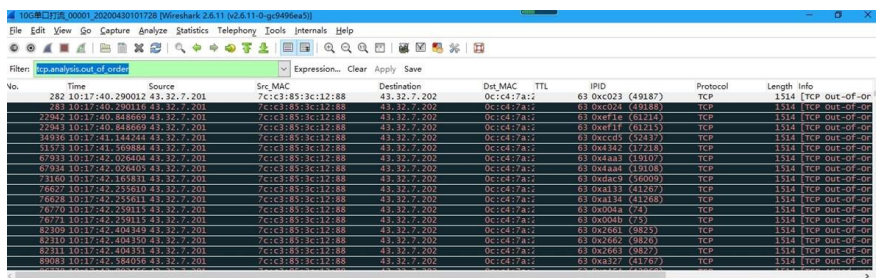
服务器===交换机===服务器 交换机侧4条链路静态聚合，服务器侧bond0绑定

问题描述

两台服务器彼此相互打流程测试，交换机二层透传vlan 3007的流量始终卡在 500Mbps左右，聚合端口是40Gbps的。服务器侧为判断问题原因在哪，于是甩开S6800交换机，服务器的4个网卡彼此直接连接。打流测试能跑到24Gbps左右。

过程分析

打流的时候抓包，看到有大量重传报文



确认是bond模式造成吞吐量下降，下面从两个角度来分析一下交换机链路聚合与服务器对接mode=0模式为什么会造成网络吞吐量下降

一、从理论角度分析

服务器和交换机做链路聚合采用mode=0模式都可能会存在这个问题，因为轮询发包的方式从两个端口发包经过的链路不同，终端收到的包可能就会出现乱序，这时候终端就会发送要求重传的报文的信号。下面是mode=0模式的详细说明， mode0 平衡负载模式：平时两块网卡工作，但需要在服务器本机网卡相连的交换机设备上进行端口聚合来支持绑定技术。

2.1 mode 0 (平衡轮询策略)

• 特点

传输数据包顺序为依次传输：即第1个包走ens33/eth0，下一个包就走ens36/eth1...一直循环下去，直到最后一个传输完毕

• 实现功能

- 实现了负载均衡
- 提供容错能力

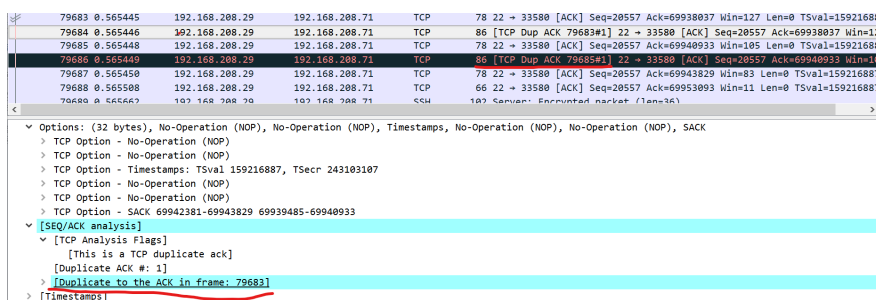
• 缺点

如果一个连接或者会话的数据包从不同的接口发出的话，中途再经过不同的链路，在客户端很有可能会出现数据包无序到达的问题，而无序到达的数据包需要重新要求被发送，这样网络的吞吐量就会下降

2.2 mode 1 (主备策略)

二、从实际情况来分析

下图是在192.168.208.71这台服务器上往192.168.208.29服务器传送文件时抓到的包，从抓到的包来看服务器经过轮询发包的形式发送数据，由于中间链路的原因确实给网络吞吐量带来了影响



当29收到乱序的包时，会回应一个DUP ACK报文，可以看出frame为79683的报文是79683的DUP ack要求重传报文，在抓到的包里通过搜索tcp.analysis.out_of_order可以看到后到的报文，所以确实有乱

序的包

The image shows a Wireshark packet capture window titled 'tcp.analysis.out_of_order'. The main pane displays a list of network packets. The columns are: No., Time, Source, Destination, Protocol, Length, and Info. The packets are all TCP segments from source 192.168.208.71 to destination 192.168.208.29. The 'Info' column for each packet indicates it is 'Out-Of-Order' and provides details like sequence number, acknowledgment number, and window size. The sequence numbers are not in order, starting at 33580 and jumping to 33588, 33590, 33592, etc.

No.	Time	Source	Destination	Protocol	Length	Info
2019	0.011116	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33580 → 22 [ACK] Seq=1537197 Ack=541 Win=329 Len=...
2021	0.011119	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33580 → 22 [ACK] Seq=1540093 Ack=541 Win=329 Len=...
4750	0.024046	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=3394185 Ack=1117 Win=329 Len=...
4752	0.024049	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=3397001 Ack=1117 Win=329 Len=...
5561	0.028168	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=3056937 Ack=1297 Win=329 Len=...
5563	0.028171	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=3953833 Ack=1297 Win=329 Len=...
6002	0.030252	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4245049 Ack=1405 Win=329 Len=...
6003	0.030255	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4247945 Ack=1405 Win=329 Len=...
6005	0.030258	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4250841 Ack=1405 Win=329 Len=...
6224	0.031241	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4389933 Ack=1441 Win=329 Len=...
6226	0.031245	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4392829 Ack=1441 Win=329 Len=...
6227	0.031246	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4397725 Ack=1441 Win=329 Len=...
6229	0.031251	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4398621 Ack=1441 Win=329 Len=...
6475	0.032391	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4557029 Ack=1477 Win=329 Len=...
6477	0.032394	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4559925 Ack=1477 Win=329 Len=...
6478	0.032396	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4562821 Ack=1477 Win=329 Len=...
6649	0.033189	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4671969 Ack=1513 Win=329 Len=...
6651	0.033192	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4674865 Ack=1513 Win=329 Len=...
6652	0.033194	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4677761 Ack=1513 Win=329 Len=...
6753	0.033627	192.168.208.71	192.168.208.29	TCP	1514	[TCP Out-Of-Order] 33588 → 22 [ACK] Seq=4737649 Ack=1513 Win=329 Len=...

例如服务器1口发出了100号包到交换机1口，服务器2口发出了101号包到交换机2口，终端先收到了交换机4口发来的101号包，后收到了交换机3口发来的100号包，收到101号包的时候发现时序不对就要求重传，然后在wireshark里经过筛选可以看到重传的报文如下图所示，

The image shows a Wireshark packet capture window with a filter applied. The main pane displays a list of network packets. The columns are: Time, Source, Destination, Protocol, Length, and Info. The packets are all SSH segments from source 192.168.208.71 to destination 192.168.208.29. The 'Info' column for each packet indicates it is a retransmission of a previous packet. The sequence numbers are 2018, 4749, 5560, 6223, 6474, 6648, 6752, 6808, 13314, 14912, 15216, 16347, 17057, 19225, 19496, and 20186. The 'Info' column for each packet indicates it is a retransmission of a previous packet.

Time	Source	Destination	Protocol	Length	Info
2018	0.011114	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
4749	0.024043	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
5560	0.028165	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
6223	0.031238	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
6474	0.032388	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
6648	0.033187	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
6752	0.033625	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
6808	0.033850	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)
13314	0.081424	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=3097 Ack=10642673 Win=...
14912	0.091194	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=3493 Ack=12077641 Win=...
15216	0.093147	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=3565 Ack=12338281 Win=...
16347	0.100615	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=3853 Ack=13379393 Win=...
17057	0.104997	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=4033 Ack=14004929 Win=...
19225	0.118833	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=4645 Ack=16117561 Win=...
19496	0.120377	192.168.208.29	TCP	102	[TCP Retransmission] 22 → 33580 [PSH, ACK] Seq=4717 Ack=16352137 Win=...
20186	0.124023	192.168.208.71	SSH	1514	Client: [TCP Fast Retransmission] , Encrypted packet (len=1448)

综上所述，linux服务器网卡聚合mode=0模式与交换机对接会使得网络的吞吐量下降，解决方案：Linux服务器网卡聚合与交换机对接更换为动态链路聚合方式(即mode=4)。

解决方法

服务器侧更换其他bond模式和交换机对接，调整合适模式参数