

Solución de problemas de las pausas TX en Nexus 2232

Contenido

[Introducción](#)

[Prerequisites](#)

[Requirements](#)

[Componentes Utilizados](#)

[Búfers de ingreso](#)

[Configuración de control de flujo](#)

[Causas de la pausa TX en Nexus 2232](#)

[Casos de prueba de laboratorio](#)

[Diagrama de la red](#)

[Prueba 1. Tráfico en ráfagas con control de flujo no habilitado en host](#)

[Prueba 2. Tráfico en ráfagas con control de flujo habilitado en host](#)

[Prueba 3. Colisión de Hash de EtherChannel](#)

[Remediación](#)

[Conclusiones y prácticas recomendadas](#)

Introducción

Este documento describe la información para ayudar a resolver problemas de las pausas de transmisión (TX) en los puertos de la interfaz de host (HIF) Nexus 2232. Se centra en el tráfico en la dirección de host a red (H2N) (tráfico que llega de los servidores a la red, de sur a norte). No cubre escenarios relacionados con los flujos de tráfico de red a host (N2H).

Este documento se ha creado principalmente para Nexus 2232 Fabric Extender (FEX), pero el concepto se aplica a B22 y 2248UPQ FEX.

Prerequisites

Requirements

Cisco recomienda que tenga conocimiento sobre estos temas

- Configuración de Cisco Nexus serie 2000
- Configuración de Cisco Nexus serie 6000

Componentes Utilizados

La información que contiene este documento se basa en las siguientes versiones de software y hardware.

- Cisco Nexus N2K-C2232PP-10GE

- Cisco Nexus 6001
- 7.1(1)N1(1)

The information in this document was created from the devices in a specific lab environment. All of the devices used in this document started with a cleared (default) configuration. If your network is live, make sure that you understand the potential impact of any command.

Búfers de ingreso

El Nexus 2232 cuenta con 32 puertos de 1/10 G de interfaz de host (HIF) y 8 puertos de 10 G de interfaz de red (NIF).

Antes de profundizar en el problema de la pausa TX, debe entender los búfers disponibles en las interfaces FEX. Los buffers asignados a la interfaz/qos-group se pueden verificar a través de este comando en el switch primario:

```
esc-6001# show queuing interface ethernet 147/1/1
if_slot 79, ifidx 0x1f920000
Ethernet147/1/1 queuing information:
Input buffer allocation:
Qos-group: 0
frh: 8
drop-type: drop
cos: 0 1 2 3 4 5 6
xon      xoff      buffer-size
-----+-----+-----
0        126720    151040
<snip>
```

Como se ve, con la calidad de servicio (QoS) predeterminada, para el tráfico de clase de descarte (qos-group 0), el FEX HIF tiene 151040 bytes para almacenar en búfer el tráfico H2N y el umbral XOFF es 126720 bytes.

Configuración de control de flujo

El Nexus 2232 está sobresuscrito a 8:1. Para evitar caídas de paquetes en la dirección H2N debido a una suscripción excesiva y a desbordamientos de búfer, Nexus 2232 tiene HIF flow-control send on de forma predeterminada:

```
esc-6001# show run int ethernet 147/1/1 all | inc flow
priority-flow-control mode auto
flowcontrol receive off
flowcontrol send on
esc-6001# show interface ethernet 147/1/1 flowcontrol
```

```
-----
Port          Send FlowControl  Receive FlowControl  RxPause  TxPause
              admin   oper    admin   oper
-----
Eth147/1/1   on     on     off     off     0         0
```

Causas de la pausa TX en Nexus 2232

Cuando se alcanza el umbral XOFF de 126720 bytes, el Nexus 2232 envía una pausa TX hacia el

host en el HIF. Las causas comunes para esto son:

1. El tráfico H2N que entra en el FEX está muy saturado, lo que hace que las memorias intermedias de ingreso estén llenas y alcancen el umbral XOFF.
2. La mayoría de las implementaciones de FEX utilizan canales de puerto para agregar varios NIF. TX Pause también se ve debido a las memorias intermedias de ingreso, que se llenan debido a la colisión de hash de etherchannel en FEX. Esto sucede cuando varios puertos HIF intentan salir de un solo NIF debido a los resultados de etherchannel.

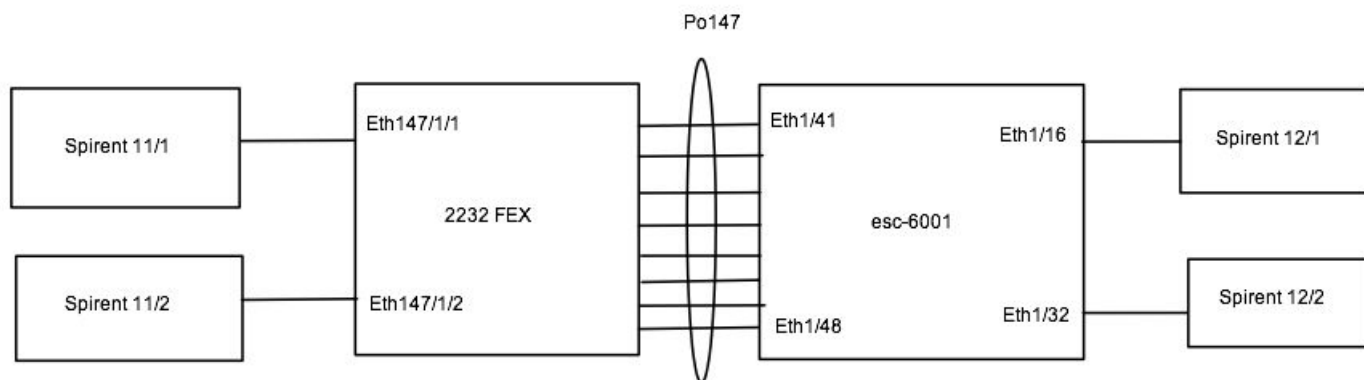
Caídas en dirección H2N

Para evitar caídas de paquetes, TX Pause se envía una vez que se alcanza el umbral XOFF. Sin embargo, se pueden ver caídas del tráfico H2N si:

1. los servidores no cumplen con la pausa o,
2. tenga un retraso para cumplir la pausa que hace que se alcance el umbral de caída de 151 KB.

Casos de prueba de laboratorio

Diagrama de la red



Para esta prueba en el laboratorio, hay cuatro puertos de espirente de 10 G que actúan como hosts, dos están en el FEX y dos en el Nexus principal 6001. Todos los puertos están en la VLAN 50. No hay otro puerto activo en el FEX o el primario:

```
esc-6001# show port-channel summary
Flags:  D - Down          P - Up in port-channel (members)
        I - Individual    H - Hot-standby (LACP only)
        s - Suspended     r - Module-removed
        S - Switched      R - Routed
        U - Up (port-channel)
        M - Not in use. Min-links not met
-----
Group Port-      Type      Protocol  Member Ports
Channel
-----
147  Po147(SU)  Eth       NONE      Eth1/41(P)  Eth1/42(P)  Eth1/43(P)
                                           Eth1/44(P)  Eth1/45(P)  Eth1/46(P)
                                           Eth1/47(P)  Eth1/48(P)
```

```
esc-6001# show fex 147 detail | exc Down
```

```

FEX: 147 Description: FEX0147   state: Online
FEX version: 7.1(1)N1(1) [Switch version: 7.1(1)N1(1)]
FEX Interim version: 7.1(1)N1(1)
Switch Interim version: 7.1(1)N1(1)
Extender Serial: FOT1635R003
Extender Model: N2K-C2232PP-10GE, Part No: 73-12533-05
Card Id: 82, Mac Addr: 20:3a:07:34:5b:02, Num Macs: 64
Module Sw Gen: 21 [Switch Sw Gen: 21]
post level: complete
Pinning-mode: static   Max-links: 1
Fabric port for control traffic: Eth1/47
FCoE Admin: false
FCoE Oper: true
FCoE FEX AA Configured: false
Fabric interface state:
  Po147 - Interface Up. State: Active
  Eth1/41 - Interface Up. State: Active
  Eth1/42 - Interface Up. State: Active
  Eth1/43 - Interface Up. State: Active
  Eth1/44 - Interface Up. State: Active
  Eth1/45 - Interface Up. State: Active
  Eth1/46 - Interface Up. State: Active
  Eth1/47 - Interface Up. State: Active
  Eth1/48 - Interface Up. State: Active
Fex Port      State  Fabric Port
  Eth147/1/1   Up     Po147
  Eth147/1/9   Up     Po147

```

Logs:

```

04/21/2015 21:58:30.162193: Module register received
04/21/2015 21:58:30.164611: Registration response sent
04/21/2015 21:58:30.196708: create module inserted event.
04/21/2015 21:58:30.197425: Module Online Sequence
04/21/2015 21:58:35.051474: Module Online

```

Prueba 1. Tráfico en ráfagas con control de flujo no habilitado en host

Cuando envía ráfagas de unidifusión de velocidad de línea de 100K 1500 bytes desde el host en Eth147/1/1(a Eth1/16) y Eth147/1/9(do Eth1/32). Cada flujo es un flujo único. El control de flujo está desactivado en el host (Spirent).

Resultados: Los puertos receptores informaron que se descartaban unos 563 paquetes por cada flujo. Dado que el control de flujo está inhabilitado en el host, puede ver mucha más pausa TX y también alta latencia (unos 100 microsegundos):

```

esc-6001# clear counters ; clear qos stat
esc-6001#
esc-6001# show interface ethernet 147/1/1, ethernet 147/1/9 | egrep Ethernet147|pause|unicast
Ethernet147/1/1 is up
  100000 unicast packets  0 multicast packets  0 broadcast packets
  0 Rx pause
  0 unicast packets  578269 multicast packets  0 broadcast packets
  578267 Tx pause
Ethernet147/1/9 is up
  100000 unicast packets  0 multicast packets  0 broadcast packets
  0 Rx pause
  0 unicast packets  578275 multicast packets  0 broadcast packets
  578273 Tx pause
esc-6001# show interface ethernet 147/1/1, eth147/1/9 flowcontrol

```

```

-----
Port          Send FlowControl  Receive FlowControl  RxPause  TxPause

```

	admin	oper	admin	oper		
Eth147/1/1	on	on	off	off	0	578267
Eth147/1/9	on	on	off	off	0	578273

Las caídas notificadas por el receptor se descartan en el FEX mismo. Hay comandos de hardware internos que pueden mostrar las caídas, pero requiere que entienda completamente la arquitectura interna de FEX que está fuera del alcance de este documento. Si necesita verificar estos contadores, póngase en contacto con el TAC para conocer este aspecto de la resolución de problemas.

Prueba 2. Tráfico en ráfagas con control de flujo habilitado en host

Cuando envía ráfagas de unidifusión de velocidad de línea de 1500 K de flujo único desde el host en Eth147/1/1(a Eth1/16) y Eth147/1/9(do Eth1/32). Cada flujo es un flujo único. El control de flujo está activado en el host (Spirent).

Resultados:

Los puertos del receptor no informan de pérdidas. La pausa mínima de TX y la latencia media son de unos 19 microsegundos:

```
esc-6001# clear counters ; clear qos stat
esc-6001# show interface ethernet 147/1/1, ethernet 147/1/9 | egrep Ethernet147|pause|unicast
Ethernet147/1/1 is up
  100000 unicast packets  0 multicast packets  0 broadcast packets
  0 Rx pause
  0 unicast packets  4743 multicast packets  0 broadcast packets
  4739 Tx pause
Ethernet147/1/9 is up
  100000 unicast packets  0 multicast packets  0 broadcast packets
  0 Rx pause
  0 unicast packets  4703 multicast packets  0 broadcast packets
  4700 Tx pause
```

```
esc-6001# show interface ethernet 147/1/1, eth147/1/9 flowcontrol
```

Port	Send FlowControl		Receive FlowControl		RxPause	TxPause
	admin	oper	admin	oper		
Eth147/1/1	on	on	off	off	0	4739
Eth147/1/9	on	on	off	off	0	4700

Caídas:

No hay caídas ya que el host honra el control de flujo enviado desde el FEX.

Prueba 3. Colisión de Hash de EtherChannel

El link ascendente entre FEX y parent es un canal de puerto. Mientras que depende del miembro en el canal de puerto que se escoja y de lo ocupado que esté, TX Pause se puede ver en los FEX HIF. En el laboratorio, sólo hay dos puertos activos en el FEX y los 8 enlaces ascendentes utilizados en el canal de puerto.

Pero para esta prueba, con el hashing predeterminado, el tráfico del host en Ethernet 147/1/1 y Ethernet 147/1/9 se envía a NIF0, que se conecta a Eth1/41 en el 6001. Si envía el 98% de tráfico

de velocidad de línea de los hosts, TX Pause se envía en ambos HIF.

Para esta prueba, el control de flujo se inhabilita en los hosts:

```
esc-6001# show interface ethernet 147/1/1, ethernet 147/1/9 | inc Ethernet14|rate|pause
Ethernet147/1/1 is up
 30 seconds input rate 9836009128 bits/sec, 819667 packets/sec
 30 seconds output rate 2516922296 bits/sec, 4915863 packets/sec
   input rate 9.84 Gbps, 819.67 Kpps; output rate 2.52 Gbps, 4.91 Mpps
 0 Rx pause
 98376923 Tx pause
Ethernet147/1/9 is up
 30 seconds input rate 9836252112 bits/sec, 819687 packets/sec
 30 seconds output rate 2516980960 bits/sec, 4915978 packets/sec
   input rate 9.84 Gbps, 819.69 Kpps; output rate 2.52 Gbps, 4.91 Mpps
 0 Rx pause
 98376916 Tx pause
```

```
esc-6001# show port-channel traffic interface port-channel 147
ChanId      Port Rx-Ucst Tx-Ucst Rx-Mcst Tx-Mcst Rx-Bcst Tx-Bcst
-----
 147  Eth1/41 99.99%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/42 0.0%    0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/43 0.0%    0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/44 0.0%    0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/45 0.0%    0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/46 0.0%    0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/47 0.00%  99.00% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/48 0.0%    1.00% 12.50% 12.50%  0.0%  0.0%
```

```
esc-6001# attach fex 147
Attaching to FEX 147 ...
To exit type 'exit', to abort type '$.'
```

```
fex-147# dbgexec w
woo> rate
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+
| Port  || Tx Packets | Tx Rate | Tx Bit  || Rx Packets | Rx Rate | Rx Bit  | Avg
Pkt|Avg Pkt|   |          | (pkts/s) | Rate    ||           | (pkts/s) | Rate    | (Tx) |
|      ||           |          |         ||           |          |         | (Tx) |
(Rx) |Err|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+
| 0-NI8 ||           | 24 | 4 | 11.23Kbps || 22 | 4 | 16.49Kbps | 272
| 448 | |
| 0-NI7 ||           | 15 | 3 | 4.17Kbps  || 17 | 3 | 3.81Kbps  | 154
| 120 | |
| 0-NI6 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI5 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI4 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI3 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI2 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI1 ||           | 1 | 0 | 1.76Kbps  || 1 | 0 | 1.08Kbps  | 1080
| 656 | |
| 0-NI0 || 4108297 | 821659 | 10.05Gbps || 1 | 0 | 1.08Kbps | 1509
| 656 | |
| 0-HI31 ||           | 1 | 0 | 2.28Kbps  || 1 | 0 | 2.28Kbps  | 1412
| 1412 | |
| 0-HI30 ||           | 1 | 0 | 2.28Kbps  || 1 | 0 | 2.28Kbps  | 1412
```

1412									
0-HI29	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI28	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI27	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI26	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI25	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI24	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI23	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI22	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI21	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI20	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI19	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI18	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI17	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI16	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI14	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI13	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI12	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI11	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI10	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI9	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI8	24556087	4911217	3.30Gbps	4094470	818894	9.95Gbps	64		
1500									
0-HI6	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI5	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI4	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI3	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI2	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI1	1	0	2.28Kbps	1	0	2.28Kbps	1412		
1412									
0-HI0	24560241	4912048	3.30Gbps	4095156	819031	9.95Gbps	64		
1500									

Caídas:

Hay caídas porque el host no está configurado para el control de flujo.

Con el control de flujo habilitado en los hosts, los hosts hacen la pausa de honor y las devoluciones del acelerador:

```

esc-6001# clear counters ; clear qos stat
esc-6001#
esc-6001# show interface ethernet 147/1/1, ethernet 147/1/9 | inc Ethernet14|rate|pause
Ethernet147/1/1 is up
 30 seconds input rate 4926871976 bits/sec, 410572 packets/sec
 30 seconds output rate 1288637816 bits/sec, 2516870 packets/sec
   input rate 4.93 Gbps, 410.57 Kpps; output rate 1.29 Gbps, 2.52 Mpps
   0 Rx pause
   88129183 Tx pause
Ethernet147/1/9 is up
 30 seconds input rate 4924820632 bits/sec, 410401 packets/sec
 30 seconds output rate 1287225224 bits/sec, 2514111 packets/sec
   input rate 4.92 Gbps, 410.40 Kpps; output rate 1.29 Gbps, 2.51 Mpps
   0 Rx pause
   88069874 Tx pause

```

```

esc-6001# show port-channel traffic interface port-channel 147
ChanId      Port Rx-Ucst Tx-Ucst Rx-Mcst Tx-Mcst Rx-Bcst Tx-Bcst
-----
 147  Eth1/41 99.99%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/42  0.0%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/43  0.0%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/44  0.0%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/45  0.0%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/46  0.0%  0.0% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/47  0.00% 99.00% 12.50% 12.50%  0.0%  0.0%
 147  Eth1/48  0.0%  1.00% 12.50% 12.50%  0.0%  0.0%

```

```

esc-6001# attach fex 147
Attaching to FEX 147 ...
To exit type 'exit', to abort type '$.'
fex-147# dbgexec w
woo> rate

```

Port	Tx Packets	Tx Rate	Tx Bit	Rx Packets	Rx Rate	Rx Bit	Avg
Pkt Avg Pkt		(pkts/s)	Rate		(pkts/s)	Rate	(Tx)
(Rx) Err							
0-NI8	32	6	19.76Kbps	19	3	16.01Kbps	366
506							
0-NI7	13	2	3.85Kbps	20	4	5.14Kbps	165
140							
0-NI6	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI5	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI4	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI3	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI2	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI1	1	0	1.76Kbps	2	0	2.16Kbps	1080
656							
0-NI0	4105292	821058	10.04Gbps	2	0	2.16Kbps	1509
656							
0-HI31	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							

0-HI30	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI29	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI28	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI27	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI26	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI25	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI24	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI23	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI22	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI21	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI20	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI19	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI18	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI17	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI16	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI14	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI13	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI12	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI11	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI10	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI9	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI8	12556848	2511369	1.68Gbps	2049754	409950	4.98Gbps	63
1500							
0-HI6	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI5	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI4	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI3	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI2	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI1	1	0	2.28Kbps	1	0	2.28Kbps	1412
1412							
0-HI0	12573036	2514607	1.68Gbps	2051092	410218	4.98Gbps	64
1499							

+-----+-----+

Remediación

De forma predeterminada, para el tráfico IP, el balance de carga de FEX se basa en la dirección MAC/IP de destino de origen. Para problemas como este, cambie el algoritmo de hash para obtener una mejor distribución del tráfico sobre el canal de puerto de entramado. Utilice este método si observa un equilibrio de carga desigual. Esta opción no es una solución absoluta:

```
esc-6001# show port-channel load-balance
```

```
Port Channel Load-Balancing Configuration:  
System: source-dest-ip
```

```
Port Channel Load-Balancing Addresses Used Per-Protocol:  
Non-IP: source-dest-mac  
IP: source-dest-ip source-dest-mac
```

Which hashing algorithm to choose depends on traffic profile. Here are the options available.

```
esc-6001# conf t
```

```
Enter configuration commands, one per line. End with CNTL/Z.
```

```
esc-6001(config)# port-channel load-balance ethernet ?  
destination-ip      Destination IP address  
destination-mac     Destination MAC address  
destination-port    Destination TCP/UDP port  
source-dest-ip      Source & Destination IP address (includes 12)  
source-dest-ip-only Source & Destination IP addresses only  
source-dest-mac     Source & Destination MAC address  
source-dest-port    Source & Destination TCP/UDP port (includes 12 and 13)  
source-dest-port-only Source & Destination TCP/UDP port only  
source-ip           Source IP address  
source-mac          Source MAC address  
source-port         Source TCP/UDP port
```

Conclusiones y prácticas recomendadas

1. TX Pause es un mecanismo operativo normal para evitar caídas de paquetes en 2232/2248UPQ/B22 FEX.
2. Maximice el número de enlaces ascendentes entre 2232/2248UPQ/B22 FEX y parent. Para poder tener más trayectos hacia la red y también ayuda a tener búfers máximos para el tráfico N2H.
3. Si los links ascendentes entre FEX y parent y no se utilizan de forma uniforme, el cambio del hashing de canal de puerto puede ayudar.
4. Dado que no hay switching local en FEX, evite tener perfiles de flujo de tráfico horizontal en los hosts en FEX.
5. Evite los dispositivos con ráfagas, como dispositivos NAS, chasis de servidor blade en FEX. Estos deben estar en el padre.
6. FEX 2348UPQ más reciente con búfer compartido de 32 millones, tiene un búfer compartido de 1 MB por HIF para el tráfico H2N para una mejor absorción de ráfagas. Además, con los enlaces ascendentes 40G NIF, las posibilidades de colisiones de hash/congestión de enlace ascendente se minimizan en gran medida.