CISCO *Live!*

Let's go

#CiscoLive

# Cisco Webex App

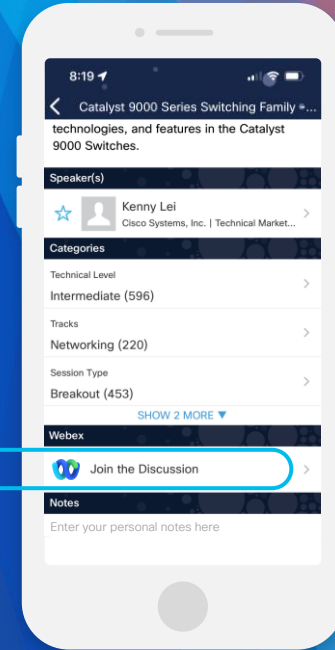## Questions?
Use Cisco Webex App to chat
with the speaker after the session

## How

1. Find this session in the Cisco Live Mobile App

2. Click "Join the Discussion"

3. Install the Webex App or go directly to the Webex space

4. Enter messages/questions in the Webex space

Webex spaces will be moderated
by the speaker until June 9, 2023.

https://ciscolive.ciscoevents.com/ciscolivebot/#BRKDCN-3641

# Agenda

- Overview

- Understanding SAN Congestion

- Detecting SAN Congestion

- Troubleshooting SAN Congestion

- Proactively preventing SAN Congestion

# Overview

What is this 'SAN Congestion' thing?

- Why am I referring to 'SAN Congestion' instead of 'Slow Drain'?

- Everyone knows 'Slow Drain', so why 'SAN Congestion'?

- Why should I be concerned?

# Understanding SAN Congestion

# Understanding SAN Congestion

## Fibre Channel Buffer-to-Buffer Flow Control – The Basics

- Fibre Channel is a 'lossless' network protocol
- Sender does not send a frame unless the receiver has a buffer
- 'Fibre Channel utilizes Buffer-to-Buffer(B2B) Credit based flow control
- Each side of link informs adjacent side of the number of buffers/credits
- Each frame sent requires a B2B credit to be returned
- B2B credits are also called 'R_RDYs'
- Frame receivers can slow rate of ingress traffic by 'withholding' credits
- If a sender runs out of credits it must stop sending until it receives one

R_RDY

FC Frame

# Understanding SAN Congestion

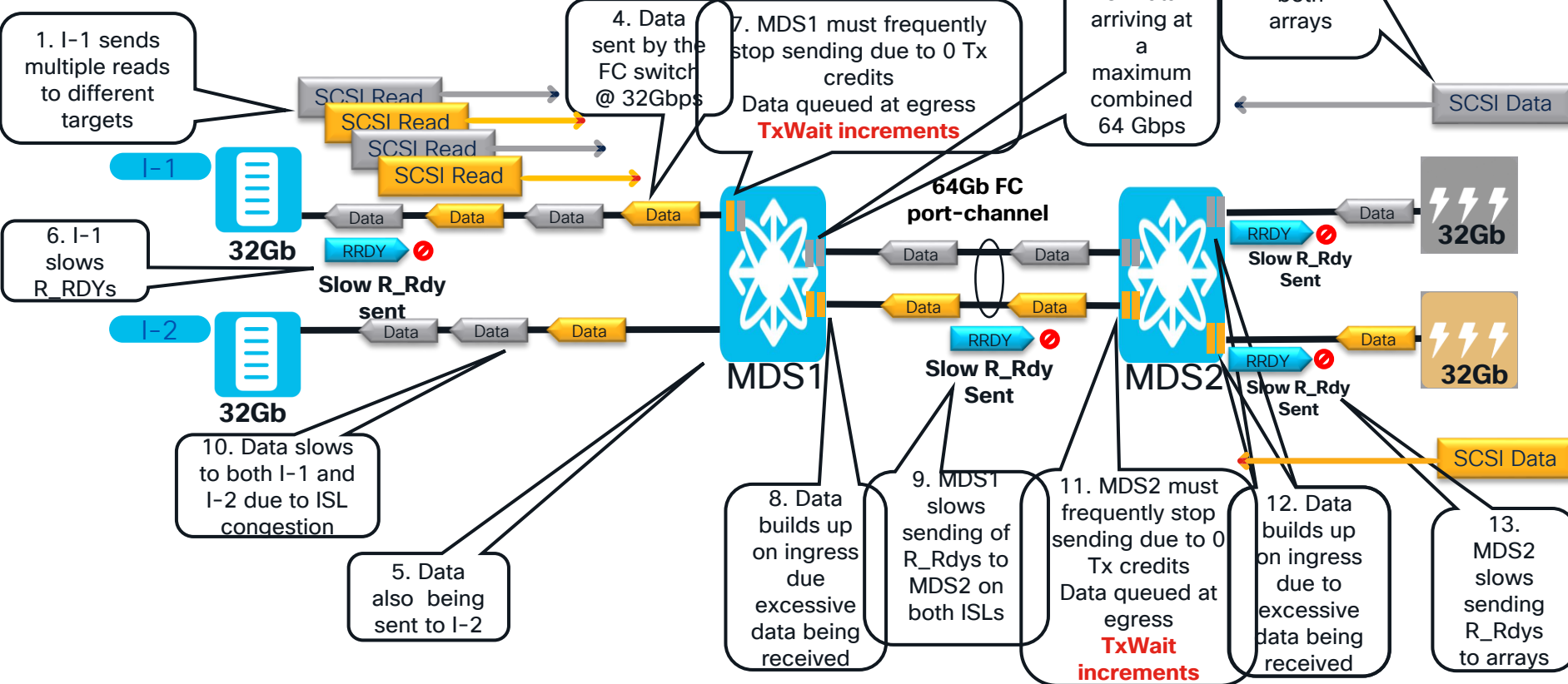There are 4 reasons for congestion in a Fibre Channel SAN

1. Slow Drain – Receiver purposely slowing down traffic by withholding R_RDYs

2. Over-Utilization – Receiver requesting more data than can be transmitted

3. Insufficient B2B credits for the link's distance(latency), speed and frame size

4. B2B credits lost due to bit errors or Invalid Transmission Words(ITW)

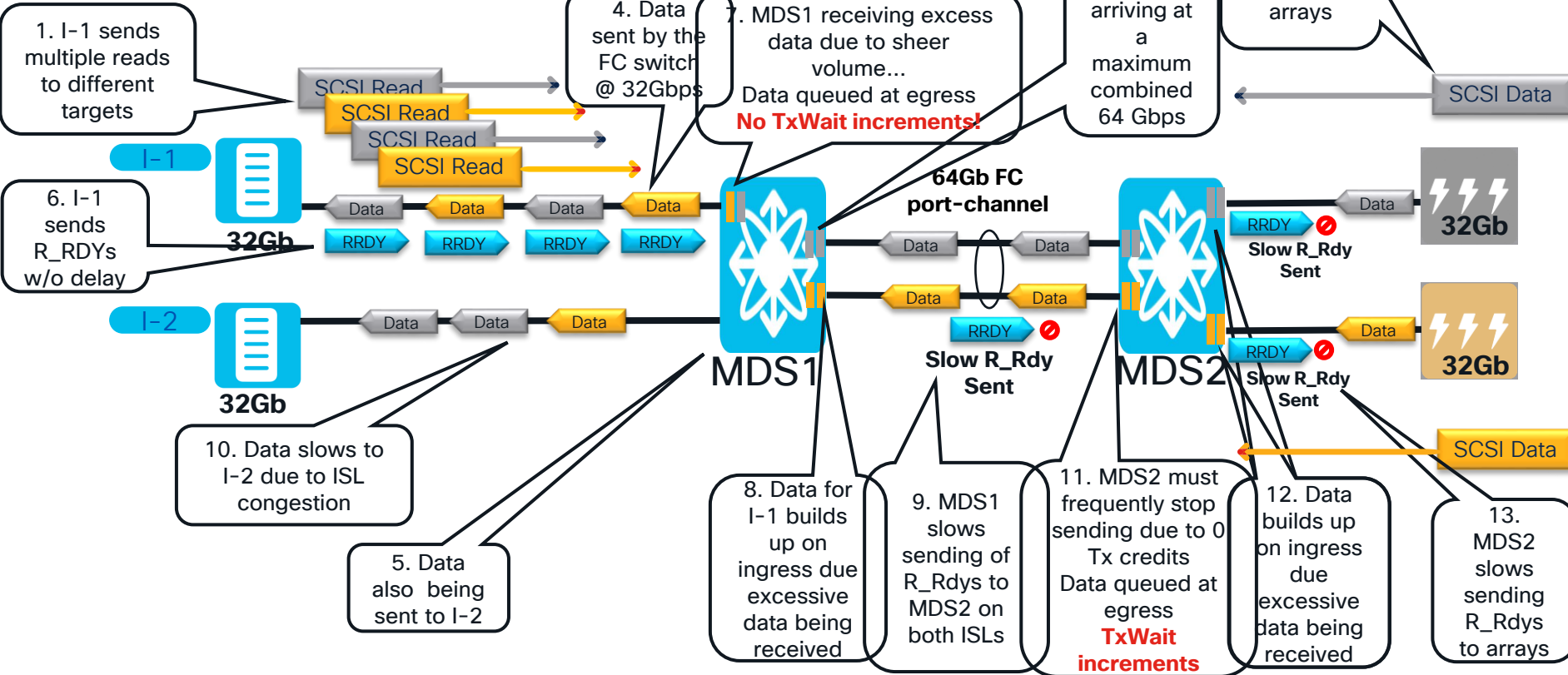#1, #2 and #4 are the focus of this presentation

# Slow Drain – Example

## "Typical" Slow Drain causing ISL and array congestion

1. I-1 sends multiple reads to different targets

2. Data sent by both arrays

3. Data arriving at a maximum combined 64 Gbps

4. Data sent by the FC switch @ 32Gbps

5. Data also being sent to I-2

6. I-1 slows R_RDYs

7. MDS1 must frequently stop sending due to 0 Tx credits
Data queued at egress
**TxWait increments**

8. Data builds up on ingress due excessive data being received

9. MDS1 slows sending of R_Rdys to MDS2 on both ISLs

10. Data slows to both I-1 and I-2 due to ISL congestion

11. MDS2 must frequently stop sending due to 0 Tx credits
Data queued at egress
**TxWait increments**

12. Data builds up on ingress due to excessive data being received

13. MDS2 slows sending R_Rdys to arrays

SCSI Read
SCSI Read
SCSI Read
SCSI Read

I-1
32Gb

I-2
32Gb

Data
RRDY
Slow R_Rdy sent

64Gb FC port-channel

Data
RRDY
Slow R_Rdy Sent

MDS1

MDS2

SCSI Data

32Gb

32Gb

RRDY
Slow R_Rdy Sent

RRDY
Slow R_Rdy Sent

SCSI Data

**Both arrays and all devices utilizing ISLs are affected!**

# Over-Utilization – Example

## Multiple Reads causing ISL and array congestion

1. I-1 sends multiple reads to different targets

2. Data sent by both arrays

3. Data arriving at a maximum combined 64 Gbps

4. Data sent by the FC switch @ 32Gbps

7. MDS1 receiving excess data due to sheer volume... Data queued at egress **No TxWait increments!**

6. I-1 sends R_RDYs w/o delay

SCSI Read
SCSI Read
SCSI Read
SCSI Read

I-1

32Gb

Data | Data | Data | Data
RRDY | RRDY | RRDY | RRDY

I-2

32Gb

Data | Data | Data

SCSI Data

64Gb FC port-channel

Data | Data

Data | Data

RRDY **Slow R_Rdy Sent**

MDS1

MDS2

RRDY **Slow R_Rdy Sent**

Data | 32Gb

RRDY **Slow R_Rdy Sent**

Data | 32Gb

SCSI Data

10. Data slows to I-2 due to ISL congestion

5. Data also being sent to I-2

8. Data for I-1 builds up on ingress due excessive data being received

9. MDS1 slows sending of R_Rdys to MDS2 on both ISLs

11. MDS2 must frequently stop sending due to 0 Tx credits Data queued at egress **TxWait increments**

12. Data builds up on ingress due excessive data being received

13. MDS2 slows sending R_Rdys to arrays

**Not strictly "slow drain" but the effects are exactly the same!**

# Comparison of Slow Drain vs. Over-Utilization

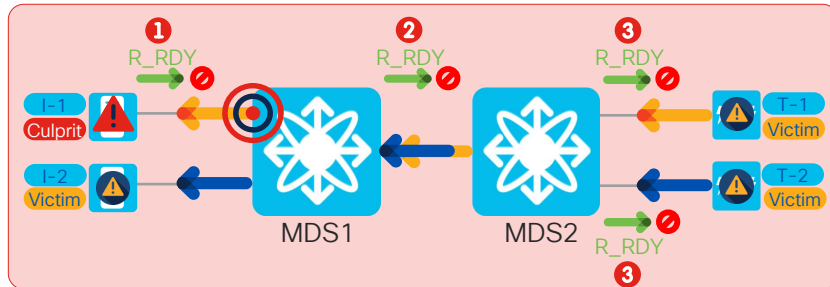| Slow Drain | Over-utilization |
|---|---|
| Tx B2B credit starvation | Receive data rate on ISL port is faster than the host port speed |

**Slow Drain side:**

I-1 is busy

I-1 slows down its ingress traffic rate by slowing down sending of R_RDY to MDS1

Tx Utilization % (tx-datarate)
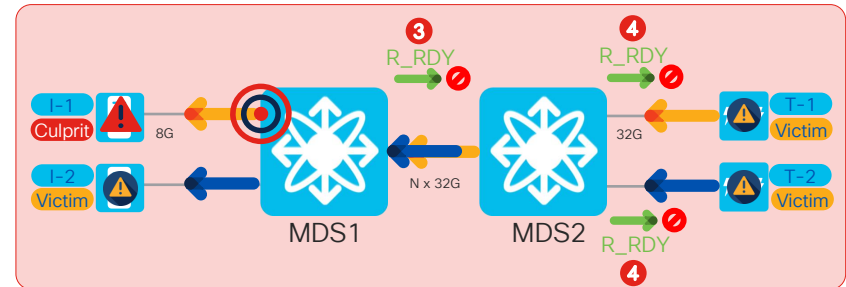
Tx Congestion % (txwait)

**Over-utilization side:**

I-1 is receiving at full capacity

MDS1 slows down traffic rate to I-1 by slowing down sending of R_RDY to MDS2

Tx Utilization % (tx-datarate)

Tx Congestion % (txwait)



Frames are not dropped in FC fabric. Rather, they consume switch buffers causing a fabric-wide congestion spreading

# Detecting SAN Congestion

# Understanding TxWait

- TxWait is the basic metric for determining/quantifying **Slow Drain**

- TxWait is an ASIC counter that increments by 1 as a port is unable to transmit a queued frame for 2.5 microseconds due to Tx B2B credit unavailability

```
mds9710# show interface fc1/1 counters | include ignore-case wait
 26009409536 2.5us TxWait due to lack of transmit credits
 Percentage TxWait for last 1s/1m/1h/72h: 0%/50%/22%/6%
```

- Convert TxWait to seconds by (TxWait * 2.5) / 1000000
  - In the above output, **26009409536** * 2.5/1000000 = 65,023 seconds
  - MDS was not able to transmit for 65,023 seconds since the counter was last cleared

- MDS enriches the raw TxWait counter:
  - For storing on switch OBFL (On-board Failure Logging (Buffer)) for troubleshooting
  - TxWait History graphs
  - For automated alerting and actions by port-monitor (PMon)
  - Export via SNMP or NX-API to remote systems like NDFC/DCNM slow drain analysis

# TxWait OBFL on MDS

- TxWait delta value is logged periodically(20 seconds) into OBFL, if delta value >=100ms.
- Displays TxWait time in 2.5µs ticks as well as in seconds.
- Timestamp of event occurrence also recorded.

> **Logged individually per module**

> **Congestion percentage is calculated over the 20 second interval**

```
MDS9706-C# show logging onboard txwait
----------------------------------
Module: 10 txwait
----------------------------------
Notes:
    - Sampling period is 20 seconds
    - Only txwait delta >= 100 ms are logged


-------------------------------------------------------------------------------
| Interface | Virtual Link | Delta TxWait Time  | Congestion | Timestamp        |
|           |              | 2.5us ticks | seconds |          |                  |
-------------------------------------------------------------------------------
|   fc1/15  | None         |    86510    |    0    |    1%    | Thu Feb 10 15:11:42 2022 |
|   fc1/15  | None         |    46459    |    0    |    0%    | Thu Feb 10 15:11:22 2022 |
|   fc1/15  | None         |   1129160   |    2    |   14%    | Sat Oct 16 00:09:52 2021 |
|   fc1/15  | None         |    658894   |    1    |    8%    | Tue Oct 12 02:18:50 2021 |
```

# Understanding Tx-datarate – Port Utilization

- Tx-Datarate is the basic metric for determining **Over-Utilization**

- Port-monitor on MDS measures datarate in percent utilization.

- Two available methods:
  - Tx-datarate: tx utilization >= 80% (*) continuously for 10 seconds (*)
  - Tx-datarate-burst: 5 (*) times in 10 seconds (*) tx utilization > 90% (*) continuously for 1 second

- An event is recorded when the high threshold(rising-threshold) is reached

- An event is recorded when the low threshold(falling-threshold is reached

- The interface was highly utilized for the time between those events

For all practical purposes, due to longer polling intervals in production environments, treat any occurrence of high utilization the same as over-utilization, which may cause congestion

(*) = user configurable

# Tx-datarate OBFL in MDS

- High-utilization events are stored in the switch
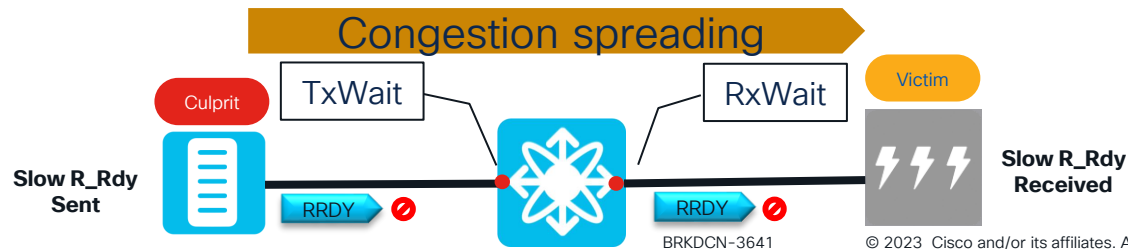
```
MDS9706-C# show logging onboard datarate
------------------------------------------------------------------------------------
| Interface |  Speed |       Alarm-types        |   Rate   |        Timestamp          |
------------------------------------------------------------------------------------
|  fc1/13   |   4G   |  TX_DATARATE_BURST_FALLING |   0@0%   | Fri Apr 29 16:41:06 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_FALLING       |   63%    | Fri Apr 29 16:40:56 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_RISING        |   98%    | Fri Apr 29 16:34:03 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_BURST_RISING  |   6@98%  | Fri Apr 29 16:34:00 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_BURST_FALLING |   0@0%   | Fri Apr 29 16:33:04 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_FALLING       |   54%    | Fri Apr 29 16:32:53 2022 |
|  fc1/13   |   4G   |  TX_DATARATE_RISING        |   98%    | Fri Apr 29 16:25:41 2022 |
```

TX_DATARATE_RISING it started at 10 seconds prior to when it was recorded 16:25:31and ended 10 seconds prior to when the TX_DATARATE_FALLING was recorded 16:32:43. There was high utilization for 7 min 12 seconds.

**Port-monitor tx-datarate *must* be configured to log to OBFL!**

# Introducing... RxWait!

- RxWait is the basic metric for determining ingress congestion

- RxWait is *new* in 64G modules and switches starting in NX-OS 9.3(2)

- RxWait measures the amount of time the switchport is preventing ingress frames

- When a switch is experiencing Tx congestion it withholds B2B credits on ports sending to the Tx congested port causing ingress congestion

- RxWait is an ASIC counter that increments by 1 as a port is at 0 Rx B2B credits for 2.5µs

- RxWait indicates ports affected by congestion not those causing congestion

- Previous generations used a software derived counter indicating 100ms of zero Rx credits

- Convert RxWait to seconds by (RxWait * 2.5) / 1000000 (just like TxWait)

# Introducing RxWait

```
MDS9710# show interface fc10/1 counters detailed
fc10/1
…
 Congestion Stats:
  Tx Timeout discards:                                0
  Tx Credit loss:                                     0
  TxWait 2.5us due to lack of transmit credits:       0
  Percentage TxWait for last 1s/1m/1h/72h:    0%/0%/0%/0%
  RxWait 2.5us due to lack of receive credits:    12345
  Percentage RxWait for last 1s/1m/1h/72h:    0%/0%/0%/0%
  Rx B2B credit remaining:                         1000
  Tx B2B credit remaining:                         1000
  Tx Low Priority B2B credit remaining:            1000
  Rx B2B credit transitions to zero:                  2
  Tx B2B credit transitions to zero:                  3
```

- In the above output, 12345 * 2.5/1000000 = 0.0308625 seconds

  - MDS was not able to receive for 0.0308625 seconds since the counter was last cleared

- MDS enriches the raw RxWait counter:
  - Graphical display – **show interface rxwait-history**
  - Last 1 second, 1 minute, 1 hour, 72 hours – **show interface <counters detailed>**
  - Historical logging every 20 seconds in OBFL(On-board Failure Logging) – **show logging onboard rxwait**

# RxWait OBFL on MDS

- RxWait delta value is logged periodically(20 seconds) into OBFL, if delta value >=100ms.
- Displays RxWait time in 2.5µs ticks as well as in seconds.
- Timestamp of event occurrence also recorded.

Logged individually per module

Ingress Congestion **percentage** is calculated over the 20 second interval

```
MDS9710# show logging onboard rxwait module 10
---------------------------------
Module: 10 rxwait
---------------------------------
Notes:
    - Sampling period is 20 seconds
    - Only rxwait delta >= 100 ms are logged


-----------------------------------------------------------------------------
| Interface | Virtual Link | Delta RxWait Time   | Congestion | Timestamp               |
|           |              | 2.5us ticks | seconds | (Ingress)  |                         |
-----------------------------------------------------------------------------
|    fc10/1 | None         | 6242188     | 15      | 78%        | Thu Jan 12 13:44:34 2023 |
|    fc10/1 | None         | 6211282     | 15      | 77%        | Thu Jan 12 13:44:14 2023 |
|    fc10/1 | None         | 6240818     | 15      | 78%        | Thu Jan 12 13:43:54 2023 |
|    fc10/1 | None         | 6229296     | 15      | 77%        | Thu Jan 12 13:43:34 2023 |
```

# Other Congestion Indications

- **Timeout-drops** - Frames dropped due to age in the switch
  - Each frame is time stamped when received on an interface
  - If age of frame exceeds 500ms(default) when it reaches egress interface it is dropped
  - Dropped frames(for any reason) cause IO errors, aborted IOs, application errors

- **Credit-Loss-Recovery** - 1/1.5 seconds of zero Tx credits
  - Occurs when an interface is at zero Tx B2B credits continuously for 1/1.5 seconds
    - 1 second for F/NP ports and 1.5 seconds for E (ISL) ports
  - Link Reset protocol is performed resulting in recovery of credits
  - Most severe indication of congestion in a Fibre Channel SAN
  - Can be caused by bit errors or severe congestion

# OBFL error-stats

```
MDS9710# show logging onboard module 10 error-stats
--------------------------------
Module: 10 error-stats
--------------------------------
Notes:
    - Sampling period is 20 seconds

------------------------------------------------------------------------
 ERROR STATISTICS INFORMATION FOR DEVICE DEVICE: FCMAC
------------------------------------------------------------------------
   Interface |                                       |        |   Time Stamp
     Range   |        Error Stat Counter Name        | Count  |MM/DD/YY HH:MM:SS
             |                                       |        |
------------------------------------------------------------------------
 fc10/1     |F64_MAC_KLM_CNTR_RX_FEC_UNCORRECTED BLOCKS |1316  |11/11/22 05:12:13
 fc10/48    |F64_CMON_CREDIT_LOSS_CH0_TMR2_HIT      |5      |07/26/22 17:39:00
 fc10/48    |F64_CMON_TX_WT_100MS_CH0_TMR1_HIT      |763    |07/26/22 17:39:00
 fc10/48    |F64_TMM_PORT_FRAME_DROP                |78876  |07/26/22 17:39:00
 fc10/48    |F64_TMM_PORT_OFFLINE                   |75408  |07/26/22 17:39:00
 fc10/48    |F64_TMM_PORT_TIMEOUT_DROP              |3477   |07/26/22 17:39:00
 fc10/48    |F64_CMON_CREDIT_LOSS_CH0_TMR2_HIT      |4      |07/26/22 17:38:20
 fc10/48    |F64_CMON_TX_WT_100MS_CH0_TMR1_HIT      |748    |07/26/22 17:38:20
 fc10/48    |F64_TMM_PORT_FRAME_DROP                |55050  |07/26/22 17:38:20
 fc10/48    |F64_TMM_PORT_OFFLINE                   |51829  |07/26/22 17:38:20
 fc10/48    |F64_TMM_PORT_TIMEOUT_DROP              |3229   |07/26/22 17:38:20
```

**Delta Credit-Loss**

F64_CMON_CREDIT_LOSS_CH0_TMR2_HIT 5 - 4 = 1 credit-loss

**Delta timeout-drops**

F64_TMM_PORT_TIMEOUT_DROP 3477 - 3229 = 248 drops

**Time intervals**

Credit-loss and timeout-drops occurred in 20 second interval ending in 17:39:00

**Other counters**

error-stats includes many other types of error counters

**Count is total – Must subtract from previous to get delta value**

# Timeout-drop S_ID / D_ID Identification

Identifies specific S_ID/D_ID of dropped frames – Useful when multiple logins

show hardware internal fcmac port x tmm_timeout_stat_buffer

Shows FC frame header info including S_ID, D_ID to identify victims

Module command (either 'attach module x' or 'slot x' prefix)

```
`slot 1 show hardware internal fcmac port 79 tmm_timeout_stat_buffer`
+----------------------------------------------------------------------------+
| PORT:78 ASIC PORT: 5 PG:1 PG PORT:1 START: 4 END: 7 WR:4 RD:0 NUM PKTS:4 |
+------+--------+--------+---+----+------+------+----+----+----+----+-+----+
|Delay | Chip  |Vegashdr|TS | FC | Src  | Dest |RCTL| CTL| SI | DI |A|OFF |
|(msec)|time(0x)|time(0x)|VLD|TYPE| ID   | ID   |(0x)|(0x)|(0x)|(0x)|T|LINE|
+------+--------+--------+---+----+------+------+----+----+----+----+-+----+
|   630|    6ff8|    6fb9|  1|   8|220340|6c0a40|   1|1800|  32|   5|0|   0|
|   630|    6ff8|    6fb9|  1|   8|220340|6c0a40|   1|1800|  32|   5|0|   0|
|   630|    6ff8|    6fb9|  1|   8|2203a0|6c0a40|   1|1800|  31|   5|0|   0|
|   630|    6ff8|    6fb9|  1|   8|2203a0|6c0a40|   1|1800|  31|   5|0|   0|
+------+--------+--------+---+----+------+------+----+----+----+----+-+----+
```

**Slot 1... port 79**
Interface fc1/79

**Src ID**
FCID of sender

**Dest ID**
FCID of destination

**Delay (msec)**
Age of frame before it was dropped

**Captures the last 4 packets dropped due to timeout per port**

# Troubleshooting
# SAN Congestion

CISCO Live!

# Troubleshooting SAN Congestion

3 Step Process

1. Understand goals

2. Classify problem

3. Follow methodology

# Troubleshooting SAN Congestion
## Goals

Two main troubleshooting goals

1. Primary - Determine the culprit

2. Secondary – Determine the various victims

# Culprits and Victims

New terminology to describe devices causing problems and those affected

- Culprits
  - Those devices causing congestion

- Victims
  - Those devices affected by the congestion
  - Three types
  1. **Direct** – Devices zoned with the culprit
  2. **Indirect** – Devices zoned with the "direct victim"
  3. **Same-path** – Devices utilizing the congested network path

Understanding culprits and victims explains the scope of the congestion

# Culprit/Victim Identification

# Troubleshooting SAN Congestion

## Victim Identification

To identify the victims(and there will be many) first understand culprit zoning

Zone members will be 'Direct Victims'

```
zone name I-1-T-1
member I-1 init
member T-1 target
```

I-1 is Culprit

T-1 is Direct Victim

Zone members of 'Direct Victims' will be 'Indirect Victims'

```
zone name I-3-T-1
member I-3 init
member T-1 target
```

I-3 is Indirect Victim of T-1

### Culprit/Victim Identification



## Identify congested path(s)

- All devices utilizing congested paths (e.g. All devices on MDS1) are potential 'Same-Path Victims'

# Troubleshooting SAN Congestion

## Victim Identification

Next look for congestion indications

- ## Culprits
  - TxWait or Tx-Datarate
- ## Direct Victims
  - RxWait
- ## Indirect Victims
  - None
- ## Same-Path Victims
  - Sender – RxWait
  - Receiver - None



Culprit/Victim Identification

# Troubleshooting SAN Congestion

Classifying Congestion Symptoms

| Level | Host Symptoms | Switch Behavior | Indications | Applicable Commands |
|---|---|---|---|---|
| 1 | Latency | Frame queuing | TxWait < 30% – Culprits, ISLs<br>RxWait – Victims, ISLs | show interface <counters><br>show logging onboard txwait<br>show logging onboard rxwait |
| 1.5 | Severe latency | Frame queuing | TxWait >= 30% – Culprits, ISLs<br>RxWait – Victims, ISLs | Same as Level 1 |
| Over-Utilization | Latency | Frame queuing | High Tx-Datarate - Culprits<br>RxWait – Victims, ISLs<br>TxWait - ISLs | Same as level 1/1.5 +<br>show logging onboard datarate |
| 2 | SCSI errors / retransmissions | Frame dropping | TxWait – Culprits, ISLs<br>RxWait – Victims, ISLs<br>Timeout-drops – Culprits, ISLs | Same as level 1/1.5 +<br>Show logging onboard error-stats |
| 3 | Extreme Delay / Application Failures | Links failing/reset (FC only) | TxWait – Culprits, ISLs<br>RxWait – Victims, ISLs<br>Timeout-drops, Culprits, ISLs<br>Credit-Loss-Recovery, Culprits, ISLs<br>Link Failures due to LR failures | Same as level 2 +<br>show logging onboard credit-loss |

Note: Each level includes all the symptoms of the previous levels

# Troubleshooting SAN Congestion
## Methodology

- Cisco recommends troubleshooting congestion in the following order:

# Troubleshooting SAN Congestion

## Methodology – Follow Congestion to Culprit

- If Rx congestion, then find ports communicating with this port that have Tx congestion

  - Zoning defines which devices communicate with this port

  - Understand topology

- If port communicating with port showing Rx congestion is FCIP

  - Check for TCP retransmits

  - Check for overutilization of FCIP



Victim

Data    Data    Data    Data

F    E

R_RDY      R_RDY

| Ingress(Rx) Congestion |
| --- |
| RxWait |
| or |
| RX_WT_AVG_B2B(100ms) |

**Congestion**

| Egress(Tx) Congestion |
| --- |
| TxWait |
| or |
| High Tx-Datarate |

# Troubleshooting SAN Congestion

## Methodology – Follow Congestion to Culprit

- ## If Tx congestion found

  - If TxWait on F port then device attached is slow drain device

  - If High Tx-Utilization on F port then attached device has Over-Utilization    Note: No TxWait in this case!

  - If E port then go to adjacent switch and continue troubleshooting

  - Continue to track through the fabric until destination F-port is discovered



| Victim | Data | Data | Data | Data | | Data | Culprit |

| F | E | E | F |

R_RDY    R_RDY    R_RDY

Congestion

| Ingress(Rx) Congestion | Egress(Tx) Congestion | Ingress(Rx) Congestion | Egress(Tx) Congestion |
|---|---|---|---|
| RxWait | TxWait | RxWait | TxWait |
| or | or | or | or |
| RX_WT_AVG_B2B(100ms) | High Tx-Datarate | RX_WT_AVG_B2B(100ms) | High Tx-Datarate |

# Troubleshooting SAN Congestion

Now that I've located the culprit, what's next?

- For Level (1) (1.5) (2) problems:
  - Investigate end device for internal bottlenecks causing the TxWait
  - Go to Prevention section

- For (OU) Over-utilization problems:
  - Increase speed of HBA(e.g. 16Gbps to 32Gbps)
  - Increase number of HBAs
  - Implement storage array based initiator rate limiting
  - Go to Prevention section

- For Level (3) problems
  - Determine if due to severe congestion on end device
  - Determine if due to lost B2B credits due to physical(bit) errors
  - Consider port-monitor to error-disable

# Troubleshooting SAN Congestion

## Determining cause of level **3** congestion

- Physical errors – Look for evidence of bit errors (switch side)

```
MDS9132T# show interface fc1/13 counters detailed
fc1/13
...
 Link Stats:
...
  Rx Primitive sequence protocol errors:          0
  Rx Invalid transmission words:                   0
  Rx Invalid CRCs:                                 0
  Rx Delimiter errors:                             0
  Rx fragmented frames:                            0
  Rx frames with EOF aborts:                       0
  Rx unknown class frames:                         0
  Rx Runt frames:                                  0
  Rx Jabber frames:                                0
  Rx too long:                                     0
  Rx too short:                                    0
  Rx FEC corrected blocks:                         0
  Rx FEC uncorrected blocks:                       0
...
  BB_SCs credit resend actions:                    0
  BB_SCr Tx credit increment actions:              0
```

**These counters indicate various types of bit errors**

**These two counters indicate B2B credits have been lost due to bit errors and recovered**

**Bit errors can cause a loss of B2B credits**

# Troubleshooting SAN Congestion

Determining cause of level ③ congestion – Continued

Physical errors – Check transceiver(SFP) power levels (switch side)

```
MDS9710# show interface fc10/1 transceiver details
fc10/1 sfp is present
 ...
    Cisco pid is DS-SFP-FC64G-SW
    Firmware version is 0.149

    No tx fault, no rx loss, in sync state, diagnostic monitoring type is 0x68
    SFP Diagnostics Information:
    --------------------------------------------------------------------------
                                    Alarms                    Warnings
                              High          Low          High          Low
    --------------------------------------------------------------------------
    Temperature   30.60 C     75.00 C     -5.00 C      70.00 C       0.00 C
    Voltage        3.29 V      3.63 V      2.97 V       3.46 V       3.13 V
    Current        7.00 mA    12.00 mA     3.00 mA     11.20 mA      3.60 mA
    Tx Power       0.58 dBm    5.00 dBm   -12.20 dBm    4.00 dBm     -8.20 dBm
    Rx Power     -12.36 dBm -  5.00 dBm   -15.20 dBm    2.00 dBm    -11.20 dBm
    Transmit Fault Count = 0
    --------------------------------------------------------------------------
    Note: ++  high-alarm; +  high-warning; --  low-alarm; -  low-warning
```

**fc10/1 has Rx Power(low light level) below the "Low Warning" threshold**

# Troubleshooting SAN Congestion

Determining cause of level **3** congestion

- Physical errors – Look for evidence of bit errors (adjacent side)

```
MDS# show rdp fcid 0xc90280 vsan 200
------------------------------------------------------------
                    RDP frame details
------------------------------------------------------------
…
Link Error Status:
-----------------------------
VN PHY port type               : FC
Link failure count             : 2
Loss of sync count             : 3
Loss of signal count           : 3
Primitive sequence proto error : 0
Invalid Transmission word      : 0
Invalid CRC count              : 0
...
FEC Status:
-----------------------------
Corrected blocks   : 0
Uncorrected blocks : 0
```

These counters indicate various types of bit errors on the adjacent device

These counters indicate Forward Error Correction(FEC) errors on the adjacent device

**RDP – Read Diagnostic Parameters – Queries stats from adjacent device**

# Troubleshooting SAN Congestion

Determining cause of level **3** congestion – Continued

## Physical errors – Check transceiver(SFP) power levels on adjacent side

```
MDS# show rdp fcid 0xc90280 vsan 200
-------------------------------------------------------------
                      RDP frame details
-------------------------------------------------------------
…
Optical Product Data:
-----------------------------
Vendor Name     : AVAGO
…
        -----------------------------------------------------------------------
                Current              Alarms                   Warnings
                Measurement     High        Low          High          Low
        -----------------------------------------------------------------------
Temperature    49.01 C        75.00 C     -5.00 C      70.00 C       0.00 C
Voltage         3.36 V         3.61 V      2.97 V       3.46 V       3.10 V
Current         7.50 mA        9.73 mA     1.54 mA      8.19 mA      2.56 mA
Tx Power        0.73 dBm       5.39 dBm  -15.92 dBm     2.34 dBm    -9.90 dBm
Rx Power       -1.09 dBm       3.38 dBm  -12.91 dBm     2.34 dBm   -11.15 dBm
        -----------------------------------------------------------------------
Note: ++  high-alarm; +  high-warning; --  low-alarm; -  low-warning
```

**RDP shows adjacent SFP Rx and Tx power values look OK**

# Troubleshooting SAN Congestion

Determining cause of level **3** congestion - Continued

Physical errors – Ensure B2B State Change(BB_SC) is functional
BB_SC is a B2B credit recovery mechanism to recover 'lost' credits

```
MDS9710# show interface fc9/2
fc9/2 is up
    Hardware is Fibre Channel, SFP is short wave laser w/o OFC (SN)
…
    Port mode is F
    Port vsan is 1
    Admin Speed is auto
    Operating Speed is 32 Gbps
    Rate mode is dedicated
    Port flow-control is R_RDY

    Transmit B2B Credit is 32
    Receive B2B Credit is 32
    B2B State Change: Admin(on) Oper(up), Negotiated Value(14)
```

**BB_SC is operational**

**BB_SC is configured on**

If operational state is down check HBA settings

**BB_SC can recover B2B credits prior to a total loss**

# Troubleshooting SAN Congestion

Determining cause of level ⬤3 congestion - Continued

If **any** physical errors are found or SFP levels are low, check and/or replace:
- SFP (switch side)
- SFP (adjacent side)
- Cable(s)
- Ensure cables do not exceed length for cable type, SFP type and speed
- Patch panels(if any)

If **no** physical errors are found:
- Investigate end device for reasons for severe congestion

Consider using port-monitor to error-disable on counter credit-loss-reco

Tip: Credit-loss-recovery only on a single fabric's connection usually means problems with the physical connection

Credit-loss-recovery on both fabrics' connections usually means severe congestion in the end device(initiator or target)

# SAN Congestion Alerting

# Automated Alerting and Congestion Prevention
## Port-monitor (PMon) on Cisco MDS

PMon monitors each switchport at a low granularity (as low as 1 second).

When a threshold exceed, PMon automatically takes actions like generating alerts, shutting down (errdisable) ports, flapping the port, isolating the port, or Dynamic Ingress Rate Limiting(DIRL).

Port-monitor has 23 counters that can be monitored

Port-monitor has 9 congestion related counters that can be monitored

# Automated Alerting and Congestion Prevention

## Available Port-monitor Counters

| | |
|---|---|
| **credit-loss-reco** | **Monitor credit loss recovery counter** |
| err-pkt-from-xbar | Monitor err-pkt-from-xbar counter |
| err-pkt-to-xbar | Monitor err-pkt-to-xbar counter |
| input-errors | Monitor input-errors counter |
| invalid-crc | Monitor invalid-crc counter |
| invalid-words | Monitor invalid-words counter |
| link-loss | Monitor link-failure counter |
| **lr-rx** | **Monitor the number of link resets received by the fc-port** |
| **lr-tx** | **Monitor the number of link resets transmitted by the fc-port** |
| rx-datarate | Monitor rx performance counter |
| rx-datarate-burst | Monitor rx-datarate-burst counter |
| sfp-rx-power-low-warn | Monitor sfp receive power low warning |
| sfp-tx-power-low-warn | Monitor sfp transmit power low warning |
| signal-loss | Monitor signal-loss counter |
| state-change | Monitor state-change counter |
| sync-loss | Monitor sync-loss counter |
| **timeout-discards** | **Monitor timeout discards counter** |
| **tx-credit-not-available** | **Monitor credit not available counter** |
| **tx-datarate** | **Monitor tx performance counter** |
| **tx-datarate-burst** | **Monitor tx-datarate-burst counter** |
| tx-discards | Monitor tx discards counter |
| **tx-slowport-oper-delay** | **Monitor tx slow port operation delay** |
| **txwait** | **Monitor tx total wait counter** |

**Congestion**

Counters in orange are congestion related Level 1, 1.5, 2 and 3

**Over-Utilization**

tx-datarate and tx-datarate-burst **are needed** for detecting Over-Utilization **On** by default in 8.5(1), 9.2(1) and later

# Automated Alerting and Congestion Prevention
## Port-monitor (PMon) on Cisco MDS

How to configure Port-Monitor?

1. Start by enabling Port-Monitor for sending alerts

2. Refine the thresholds over weeks/months. Solve the real culprits. Avoid too many alerts.

3. Finally, enable actions, such as congestion prevention using DIRL

4. Ensure tx-datarate and/or tx-datarate-burst are on for Over-Utilization!

5. Go to step 2

Sample PMon policies: https://www.cisco.com/c/en/us/support/docs/storage-networking/mds-9000-nx-os-software-release-62/200102-Sample-MDS-port-monitor-policy-for-alert.html

# PMon Policy on MDS

```
#
port-monitor name fabricmon_edge_policy
  logical-type edge
  counter txwait poll-interval 1 delta rising-threshold 30 event 4 falling-threshold 10 event 4 alerts syslog rmon portguard DIRL
  counter tx-datarate poll-interval 10 delta rising-threshold 80 event 4 falling-threshold 70 event 4 alerts syslog rmon obfl portguard DIRL
  counter tx-datarate-burst poll-interval 10 delta rising-threshold 5 event 4 falling-threshold 1 event 4 alerts syslog rmon obfl datarate 90


# Show port-monitor
Policy Name  : fabricmon_edge_policy
Admin status : Not Active
Oper status  : Not Active
Port type    : All Edge Ports
```

| Counter | Threshold Type | Interval (Secs) | Warning | | Thresholds | | | Rising/Falling actions | | Congestion-signal | |
| | | | Threshold | Alerts | Rising | Falling | Event | Alerts | PortGuard | Warning | Alarm |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Link Loss | Delta | 30 | none | n/a | 5 | 1 | 4 | syslog,rmon | FPIN | n/a | n/a |
| Sync Loss | Delta | 30 | none | n/a | 5 | 1 | 4 | syslog,rmon | FPIN | n/a | n/a |
| Signal Loss | Delta | 30 | none | n/a | 5 | 1 | 4 | syslog,rmon | FPIN | n/a | n/a |
| Invalid Words | Delta | 30 | none | n/a | 1 | 0 | 4 | syslog,rmon | FPIN | n/a | n/a |
| Invalid CRC's | Delta | 30 | none | n/a | 5 | 1 | 4 | syslog,rmon | FPIN | n/a | n/a |
| State Change | Delta | 60 | none | n/a | 5 | 0 | 4 | syslog,rmon | none | n/a | n/a |
| TX Discards | Delta | 60 | none | n/a | 200 | 10 | 4 | syslog,rmon | none | n/a | n/a |
| LR RX | Delta | 60 | none | n/a | 5 | 1 | 4 | syslog,rmon | none | n/a | n/a |
| LR TX | Delta | 60 | none | n/a | 5 | 1 | 4 | syslog,rmon | none | n/a | n/a |
| Timeout Discards | Delta | 60 | none | n/a | 200 | 10 | 4 | syslog,rmon | none | n/a | n/a |
| Credit Loss Reco | Delta | 1 | none | n/a | 1 | 0 | 4 | syslog,rmon | none | n/a | n/a |
| TX Credit Not Available | Delta | 1 | none | n/a | 10% | 0% | 4 | syslog,rmon | none | n/a | n/a |
| RX Datarate | Delta | 10 | none | n/a | 80% | 70% | 4 | syslog,rmon,obfl | none | n/a | n/a |
| TX Datarate | Delta | 10 | none | n/a | 80% | 70% | 4 | syslog,rmon,obfl | DIRL | n/a | n/a |
| TX-Slowport-Oper-Delay | Absolute | 1 | none | n/a | 50ms | 0ms | 4 | syslog,rmon | none | n/a | n/a |
| TXWait | Delta | 1 | none | n/a | 30% | 10% | 4 | syslog,rmon | DIRL | 40% | 60% |
| RX Datarate Burst | Delta | 10 | none | n/a | 5@90% | 1@90% | 4 | syslog,rmon,obfl | none | n/a | n/a |
| Input Errors | Delta | 60 | none | n/a | 5 | 1 | 4 | syslog,rmon | none | n/a | n/a |

```
On falling threshold portguard actions FPIN, DIRL, Cong-Isolate-Recover will initiate auto recovery of ports.
```

# NDFC Congestion/Congestion Analysis

## Best Practice – Run in always-on mode.

- Slow-drain analysis is not enabled by default

- After adding a new fabric:
  - Enable performance monitoring
  - Schedule to run slow drain analysis daily for 24 hou



*Slow-drain Analysis is renamed to Congestion Analysis in NDFC 12.1.1e

- DCNM/NDFC slow-drain analysis has minimal/negligible effect on the switches

# NDFC Congestion/Slow-drain Analysis

## Always-on, historical view with trending and seasonality

- fc1/33 is congested in Tx direction
- TxWait increases but not all the time. Only two spikes in last 12 hours.
- Next Steps –
  - Correlate with host and app. Does it correlate with a cron job on the host?
  - Look at SAN Insights metrics to find the root cause.

# SAN Congestion Preventing

# SAN Congestion

Including Slow Drain and Over-Utilization

We talked about

Understanding,

Detection,

Troubleshooting

and

Alerting

Now, let's talk about Prevention

# SAN Congestion Innovation on Cisco MDS

**DIRL**
MDS limits the I/O from the congested devices – The ultimate solution to SAN Congestion that works today.
Cisco Patented

**FPIN**
Notify the end devices about congestion. Action is dependent on the end devices
T11 Standard

**Traffic Segregation**
VSAN | Virtual Links
Segregate traffic on ISL resulting in reduced spread of congestion

**Congestion Isolation**
Isolate traffic going to a congested device into a dedicated Virtual Link on an ISL resulting in reduced spread of congestion

**No-credit-drop**
Drop frames going to a congested device resulting in easing up fabric-wide congestion

2015    2017    2017    2021    2021

CISCO *Live!*

Not an exhaustive list

# Common Causes of SAN Congestion

**Slow Drain**

**Over-Utilization**

**Different detection**
On the switch port connected to the culprit device

**Similar effects**
On fabric and victim devices

# Common Causes of SAN Congestion

Slow Drain

Over–Utilization

**Different detection**
On the switch port connected to the culprit device

**Similar effects**
On fabric and victim devices

**Same Root Cause**
The culprit device is receiving more than it can ingest

# The Root Cause of SAN Congestion

**The Root Cause**

I-1 is receiving more than it can ingest

# The Root Cause of SAN Congestion

**The Root Cause**

I-1 is receiving more than it can ingest

Why is I-1 receiving more than it can ingest?

...because I-1 is asking for it.

# The Solution

**Cisco Dynamic Ingress Rate Limiting**



DIRL

- I-1 is asking for more than it can ingest
- DIRL limits I-1's asking rate to reduce its receiving rate
- DIRL dynamically changes I-1's *asking* rate to adapt to its traffic profile

# The Solution

**Cisco Dynamic Ingress Rate Limiting**



DIRL

**I-1 is asking** for more than it can ingest

DIRL **limits I-1's** *asking* **rate** to reduce its receiving rate

DIRL **dynamically changes** I-1's *asking* rate to **adapt** to its traffic profile



DIRL prevents SAN Congestion due to slow-drain **and** over-utilization.

# Cisco Dynamic Ingress Rate Limiting

## End-device independent
Upgrading of end-devices is not needed

## Adaptive
DIRL dynamically adjusts as per the traffic profile of the host

## No side effects
Rate limits congested hosts only. Other non-congested hosts and storage ports are not impacted

## Easy adoption
DIRL is available on MDS switches after a software-only upgrade.

## Gradual Rollout
DIRL can be implemented one switch at a time

## Affordable
No additional license needed

## Topology independent
DIRL works in edge-core, edge-core-edge, or collapsed core (single switch fabric) topologies

**Without Cisco DIRL**



**With Cisco DIRL**



**Back pressure alleviated!**

# Notifications and Congestion Signals in Fibre Channel



Available on Cisco MDS in NX-OS 8.5(1) onwards

Exchange Diagnostic Capabilities (EDC) (for Congestion Signals)
Register Diagnostic Functions (RDF) (for FPINs)
Congestion Signals (Primitives)
Fabric Performance Impact Notifications (FPIN)

# DIRL vs FPIN

- DIRL helps today. FPIN readiness will take a few years.
  - DIRL is available on existing MDS switch after a software-only upgrade, without any dependency on end devices
  - Although FPIN is supported on MDS switches, action is dependent on the end devices

- DIRL is affordable
  - DIRL and FPIN work on existing MDS switches and don't need an additional license
  - Must upgrade end-devices to benefit from FPIN

- In the future, when you are ready for FPIN, DIRL will continue to be a complementary technology
  - What if a few devices don't react to FPIN and still cause congestion? DIRL within MDS switches will be the protection

# SAN Congestion Management - Recommendations

**Reactive**

- Gather 'show tech-support slowdrain' from all switches
- Use OBFLand other commands to identify culprit and victims
- TAC can help!

**Proactive**

- Schedule NDFC/DCNM Congestion Analysis to run daily for 24 hours.
- Important for troubleshooting
- Configure MDS port-monitor (PMon) for automated alerts and actions.
- Important for congestion prevention using DIRL.

**Predictive**

- Enable SAN Analytics and SAN Insights for getting visibility into application I/O traffic patterns.
- Important for finding the underlying root cause and predicting congestion

*Slow-drain Analysis is renamed to Congestion Analysis in NDFC 12.1(1e)

# Upcoming Book Available For Pre-order

# Related sessions

| Session ID | Title | Time and Venue | Speaker |
|---|---|---|---|
| BRKDCN-2945 | IP Fabric for Storage Networks Best Practice and Design | Monday, June 5 8:30 AM– 10:00 AM | Paresh Gupta |
| LTRCRT-2821 | Hands-on preparation for CCNP Data Center Certification with Cisco SAN labs – | Sunday, June 4 9:00 AM – 1:00 PM | Somit Maloo Iskren Nikolov |
| BRKDCN-3645 | DCNM SAN Insights – Real-time and always-on NVMe visibility at scale | Wednesday, June 7 1:00 PM – 2:00 PM | Paresh Gupta |
| CCP-1411 | Data Center Switching Hardware Platform Roadmap update | Thursday, June 8 8:30 AM – 9:30 AM | Becky Marques |
| BRKDCN-3677 | Dos and Don't of Deploying NVMe Over Fabrics | Thursday, June 8, 1:00 PM – 2:00 PM | Kamal Bakshi |

# Fill out your session surveys!

Attendees who fill out a minimum of four session surveys and the overall event survey will get **Cisco Live-branded socks** (while supplies last)!

Attendees will also earn 100 points in the **Cisco Live Challenge** for every survey completed.

**These points** help you get on the leaderboard and increase your chances of winning daily and grand prizes
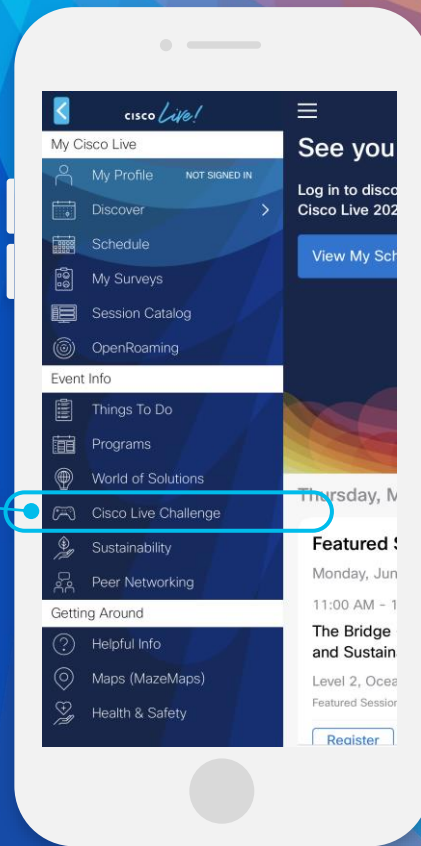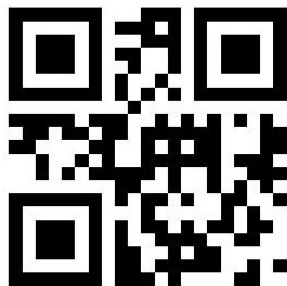
# Continue
# your education

- Visit the Cisco Showcase for related demos

- Book your one-on-one Meet the Engineer meeting

- Attend the interactive education with DevNet, Capture the Flag, and Walk-in Labs

- Visit the On-Demand Library for more sessions at www.CiscoLive.com/on-demand

# Cisco Live
# **Challenge**

## Gamify your Cisco Live experience!
Get points for attending this session!

## How:

**1** Open the Cisco Events App.

**2** Click on 'Cisco Live Challenge' in the side menu.

**3** Click on View Your Badges at the top.

**4** Click the + at the bottom of the screen and scan the QR code:

CISCO *Live!*

Let's go

#CiscoLive