



**Hewlett Packard**  
Enterprise

# **TROUBLESHOOTING LNET MULTI-RAIL NETWORKS - DEMO**

---

Chris Horn, Lustre Software Engineer

May, 2022

# OUTLINE

---

- Part 1 Recap
- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



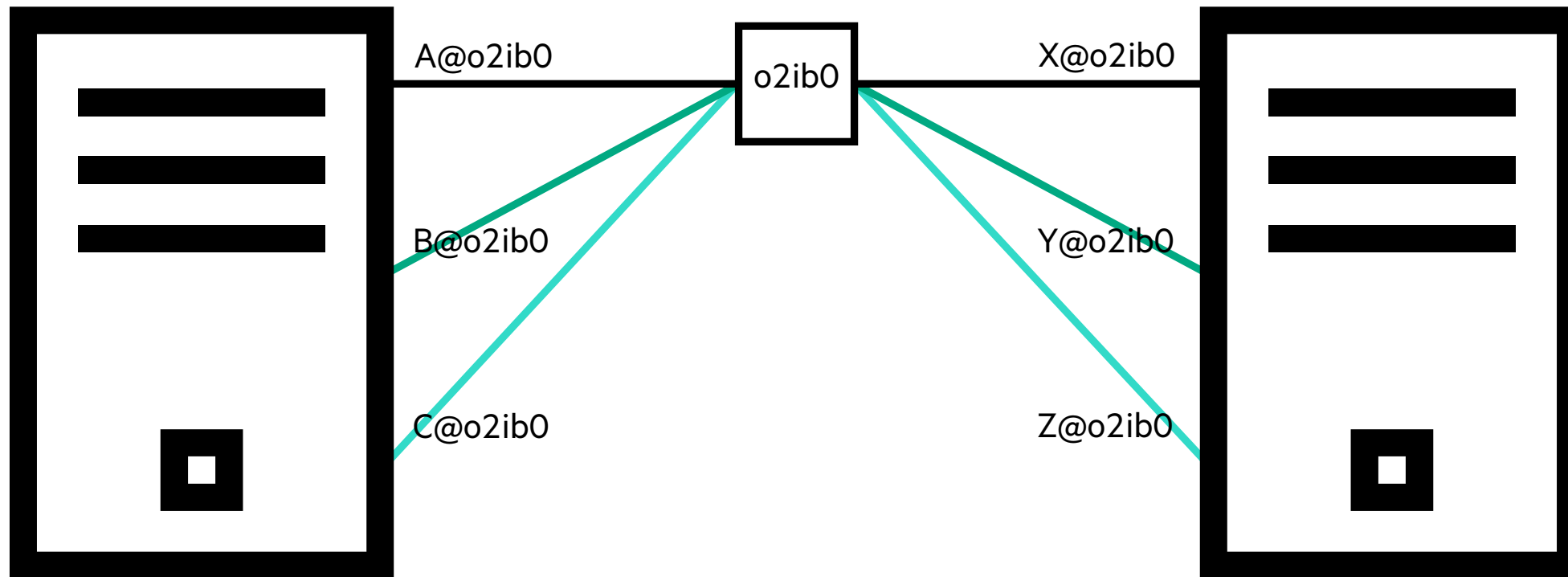
# PART 1 RECAP

---

- Links
  - Event Page: <https://www.eofs.eu/events/lad21>
  - Video: <https://youtu.be/j3m-mznUdac>
  - Slides: [https://www.eofs.eu/\\_media/events/lad21/lnet\\_multi-rail\\_troubleshooting.pdf](https://www.eofs.eu/_media/events/lad21/lnet_multi-rail_troubleshooting.pdf)
- LNet Multi-Rail Overview
  - Basics
  - Role of Primary NID
- Important Statistics
  - Local and Peer NI send and receive counts
- Validating Expected Behavior
- A Closer Look at LNet Health



# LNET MULTI-RAIL OVERVIEW



In Lustre 2.10, the LNet Multi-Rail feature allows for multiple interfaces in the same LNet network.

# PRIMARY NIDS

- Primary NID uniquely identifies a multi-rail peer
- The first NID configured on a node is designated the primary NID.
- Primary NID used as a key for modifying peer with Inetctl CLI



# IMPORTANT STATISTICS

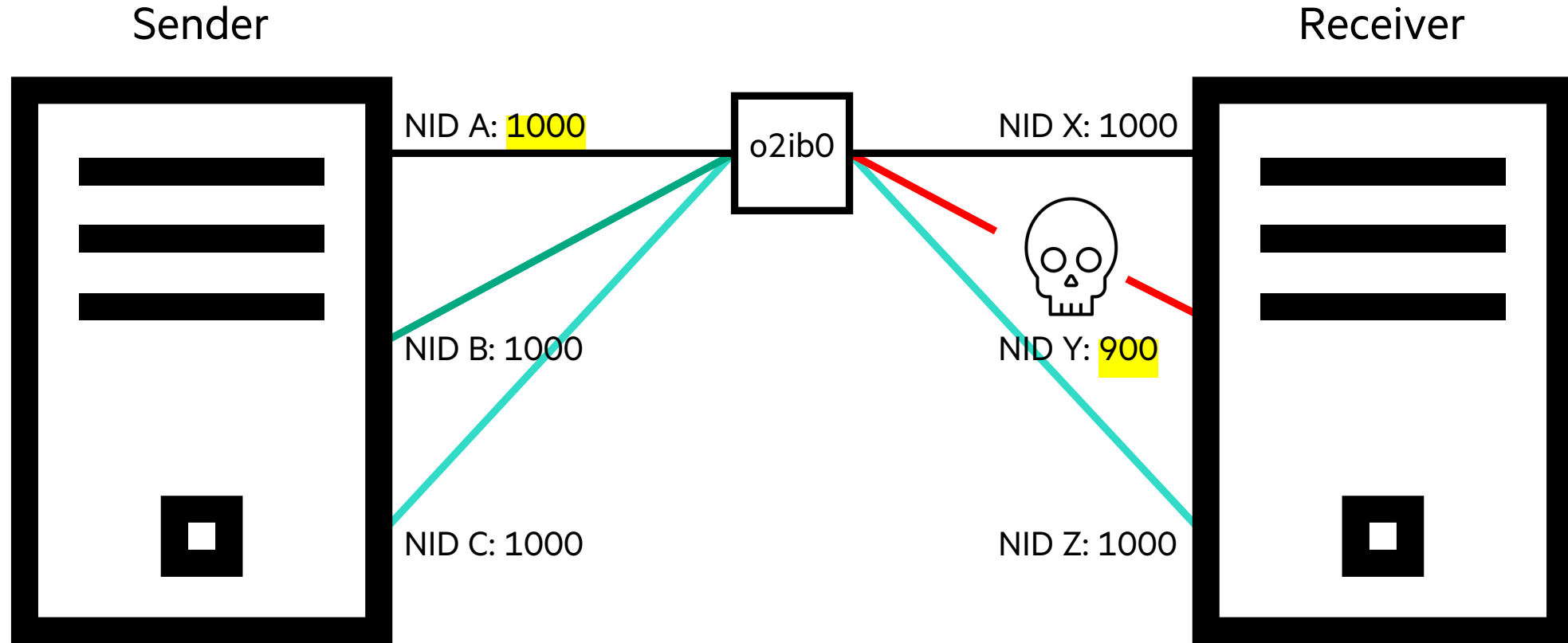
---

```
# lnetctl net show -v --net tcp | egrep -e 'nid|send_count|recv_count'
  - nid: 192.168.2.30@tcp
    send_count: 27
    recv_count: 27
  - nid: 192.168.2.31@tcp
    send_count: 25
    recv_count: 25
# lnetctl peer show --nid 192.168.2.38@tcp -v | egrep -e 'nid|send_count|recv_count'
  - primary nid: 192.168.2.38@tcp
    - nid: 192.168.2.38@tcp
      send_count: 26
      recv_count: 26
    - nid: 192.168.2.39@tcp
      send_count: 26
      recv_count: 26
```

#



# LNET HEALTH



NID Y has decremented health.  
Future sends will avoid NID Y.



# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths





# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



# DISCOVERY DEMO RECAP

---

- Discovery controls whether MR peers are created as a result of traffic
  - MR peers can always be created via CLI
- Lustre 2.12
  - If discovery enabled locally then MR peer is created if the peer supports MR (even if peer has discovery disabled)
  - If discovery disabled locally then MR peer is not created
- Lustre 2.15
  - Discovery must be enabled on both local host and remote peer for MR peer to be created



# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- **Show how to view LNet health state and activities**
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



# LNET HEALTH DEMO RECAP

---

- Health values tracked for local and remote (peer) interfaces
  - View local with `Inetctl net show -v 2 | egrep -e nid -e 'health value'`
  - View remote with `Inetctl peer show -v 2 | egrep -e nid -e 'health value'`
- Health values modified based on how LNDs classify send failure (message status)
  - LOCAL\_\* message status tells LNet problem was with local (source) interface
    - Health value for source interface decremented
  - REMOTE\_\* message status tells LNet problem was with remote (destination) interface
    - Health value for destination interface decremented
  - NETWORK\_TIMEOUT message status specified when LND does not know where problem was
    - Health values for both source and destination interfaces are decremented
- Interface with decremented health is placed into recovery mode
  - View local interfaces in recovery with `Inetctl debug recovery -l`
  - View remote interfaces in recovery with `Inetctl debug recovery -p`



# LNET HEALTH DEMO RECAP (CONT)

---

- Lustre 2.12
  - Interfaces in recovery ping'd every recovery\_interval seconds (lnetctl global show | grep recovery\_interval)
  - Interfaces in recovery are stuck if the interface is removed from cluster (client reboot/LNet configuration changed, etc.)
- Lustre 2.15
  - Interfaces in recovery ping'd using exponential backoff
    - next\_ping = current\_time + 2^pings\_sent ; 1<sup>st</sup> ping after 1 seconds; 2<sup>nd</sup> ping after 2 seconds, 3<sup>rd</sup> ping after 4 seconds, etc.
    - 15 minute max interval (configurable with <https://jira.whamcloud.com/browse/LU-14979>)
  - We must receive a message from an interface before it is eligible for recovery
  - Remote (peer) interfaces can age out of recovery
    - Configurable with recovery\_limit parameter (default to 0: remote interfaces do not age)
    - After recovery\_limit (seconds) the peer interface is removed from recovery
    - Peer interface becomes eligible again once we receive a message from it



# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



# PRIMARY NID DEMO RECAP

---

- Primary NID identifies a multi-rail peer
- Lustre log messages will reference the primary NID
  - The primary NID may be on a network that is `_not_` used for Lustre
  - This does not mean Lustre is using that NID for traffic. It is only identifying a peer
- Modifying the primary NID of a mounted client (or server) is a bad idea
  - Lustre log messages can still reference the old primary NID





# DEMO

---

- Demo Environment
- Show how discovery feature affects the Multi-Rail feature
- Show how to view LNet health state and activities
- Show how primary NID is referenced by Lustre
- Show how Inetctl ping can and cannot find broken paths



# LNETCTL PING DEMO RECAP

---

- The destination argument to lnetctl ping command specifies a `_peer_`, not a specific endpoint
- If the nid of a MR peer is specified LNet may send the ping to any peer interface (not the specified NID)
- Lustre 2.15 adds `--source` option to lnetctl ping command
  - Specifies a local interface to use for the ping
  - This fixes the destination NID so LNet will not select a different peer interface



# Q & A

Chris Horn  
email: [chris.horn@hpe.com](mailto:chris.horn@hpe.com)



# TICKETS

---

- Inetctl --source flag
  - <https://jira.whamcloud.com/browse/LU-14939> (Landed for 2.15)
- Restore round robin when NI returned to service
  - <https://jira.whamcloud.com/browse/LU-13575> (Landed for 2.15)
- NIs stuck in recovery
  - <https://jira.whamcloud.com/browse/LU-13569> (Landed for 2.15)
- Correct classification of send errors
  - <https://jira.whamcloud.com/browse/LU-13571> (Landed for 2.14)
  - <https://jira.whamcloud.com/browse/LU-14540> (Landed for 2.15)

