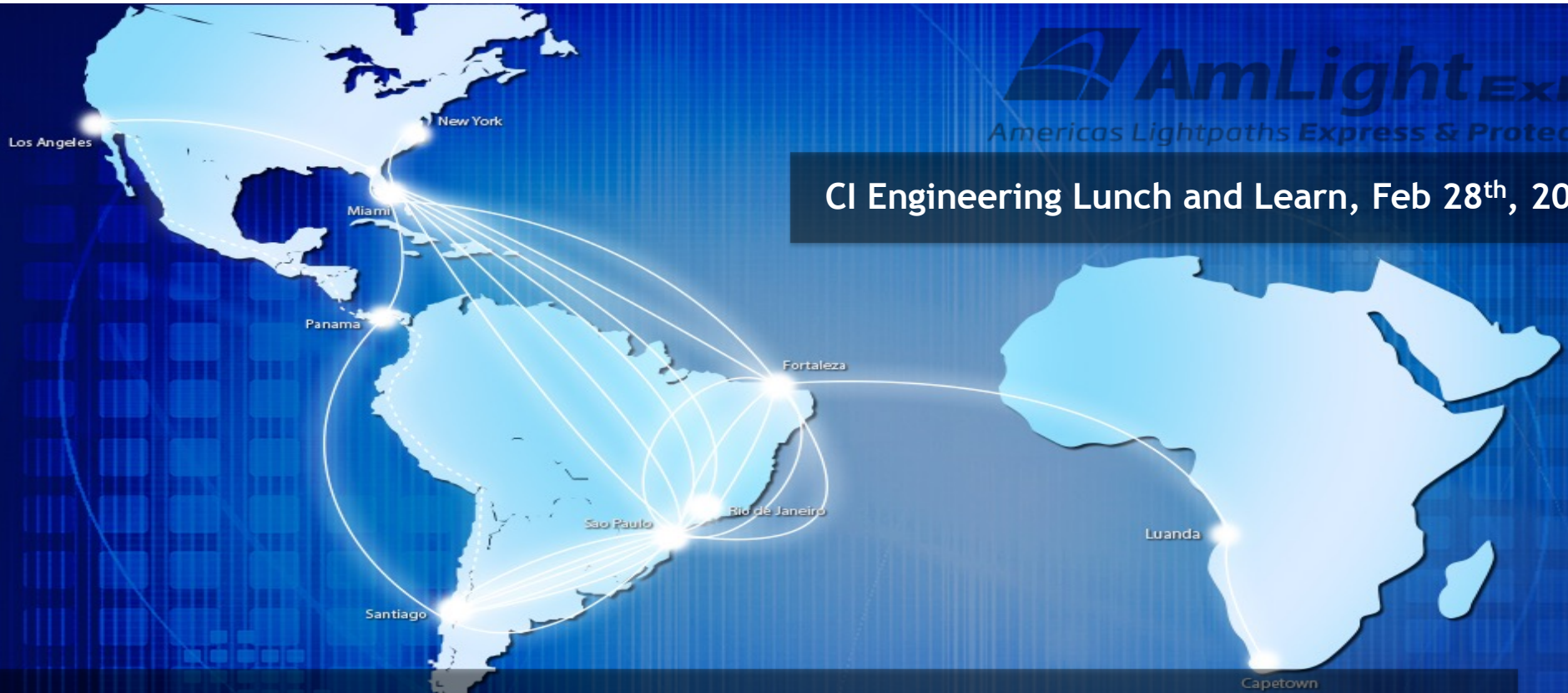AmLight ExP
Americas Lightpaths Express & Protect

CI Engineering Lunch and Learn, Feb 28th, 2025

**BERToD: An automated BER testing framework to search for packet loss at AmLight**

Jeronimo Bezerra - FIU/AmLight

Renata Frez – RNP/AmLight

Italo Valcy – FIU/AmLight

# Outline

- Motivation
- How BERToD works
- Lessons Learned
- Next Steps
- Conclusion

**AmLight** ExP
Americas Lightpaths Express & Protect

# Disclaimer

- Packet vs Frame
  - Interchangeable in this presentation

- Packet Loss vs Packet Drop:
  - Drop: We drop our packets
    - Tail drop, QoS, blocking topologies, traffic engineering, small buffers
  - Loss: Someone/thing loses/corrupts our packets
    - Fiber cuts, power outages, damage components

- For this talk, our focus is on packet loss!

**AmLight** ExP
Americas Lightpaths **Express & Protect**

# Motivation

- ***Data transfers over long-haul links suffer extra with packet loss***
  - 125 ms round-trip time (RTT) from Chile or Brazil to Jacksonville, FL
  - 131 ms RTT from Chile or Brazil to Atlanta
  - *A packet loss rate of $1\times10^{-3}$ is enough to disrupt data movement workflows over 100+ ms RTTs.*
    - Reference: July 7th, 2023, CI Engineering Lunch and Learn *Handling Microbursts @ AmLight – Part 2 of 2*
    - *https://www.es.net/science-engagement/ci-engineering-lunch-and-learn-series/*

- Science applications are <u>expecting</u> better network performance
  - SLA-driven science drivers are <u>demanding</u> more granular loss detection measurement ($1\times10^{-9}$ or *1 out of 1,000,000,000* packets)
  - Such granularity is hard to achieve by just using standard hardware and software

- AmLight has grown in complexity in the last 5 years
  - Next slide

- Current solutions for packet loss detection have <u>scalability</u>, <u>accuracy</u>, or <u>granularity</u> limitations

**AmLight** ExP
Americas Lightpaths **Express & Protect**

# NSF IRNC: AmLight Network – 2020-2025

- A distributed academic exchange point built to enable collaboration among Latin America, Africa, and the U.S.

- Supported by NSF, OAC, and the IRNC program under award # OAC-2029283 for the 2021-2025

- Partnerships with R&E networks in the U.S., Latin America, Caribbean and Africa, built upon layers of trust and openness by sharing:
  - Infrastructure resources
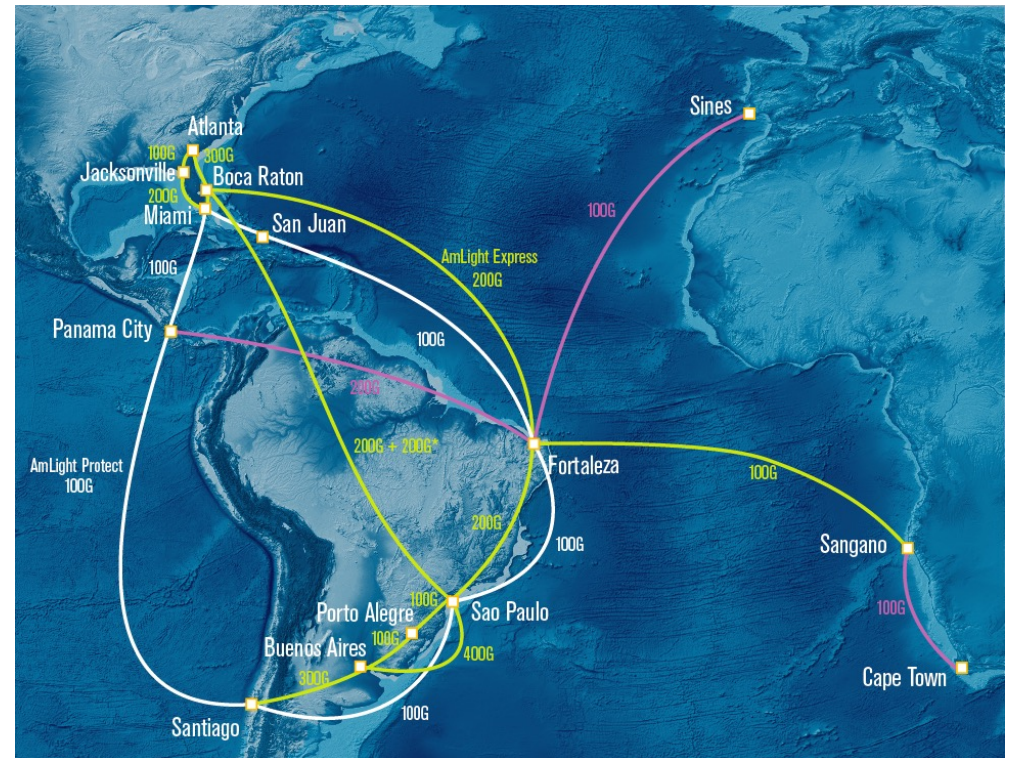  - Human resources

(NSF Award # OAC-2029283)

# NSF IRNC: AmLight Network – 2020-2025

- **39x 100G links:**
  - 2.1+ Tbps of <u>international</u> connectivity
  - AmLight will reach 5+ Tbps of <u>total</u> capacity[2025]
  - Dark fiber, spectrum, waves, and lit services

- **9x Sites / 19x racks:**
  - Miami, Boca Raton, Jacksonville, Sao Paulo, Fortaleza, Santiago, San Juan, Panama City, Cape Town, Atlanta, and Buenos Aires

- **Network and Monitoring Devices:**
  - 20x programmable switches and 5x Juniper routers
  - 10x 10G perfSonar nodes
  - 4x 100G servers
  - 4x In-band Network Telemetry (INT) collectors
    - ~10Mpps & 96TB of telemetry data per day

**AmLight** ExP
Americas Lightpaths **Express & Protect**

# Current approaches to detect/isolate packet loss

- AmLight monitoring:
  - SDN counters, ICMP, SNMP, traceroute, optical telemetry metrics (polling-based), INT and JTI reports (streaming-based), perfSONAR measurements, and dedicated 100G nodes

- Even with all of those, fault isolation and mitigation are still challenging and have a high OPEX:
  - Evaluate results, correlate data, run extra tests, send field technicians to clean/replace suspicious components, steer traffic, and run again with different outcomes
    - Days/Weeks of work

- Existing packet generators/network testers are used in an ad-hoc fashion
  - Hardware-based granularity,
  - But manual configuration and with a learning curve to configure and read results

**AmLight** **ExP**
Americas Lightpaths **Express & Protect**

# BERToD - Bit Error Rate Testing on Demand

- An automated packet loss detection framework that uses granular per-packet network telemetry (INT), SDN, a customized networking pipeline, and hardware-based packet generator to detect bit error rates as low as $1\times10^{-12}$

- BERToD leverages recent developments at AmLight:
  - Flexible forwarding rules provided by the SDN switches
  - Link and buffer utilization monitoring provided by In-band Network Telemetry (INT)
  - Topological data and dynamic service instantiation provided by the Kytos-ng SDN Controller

- *Near* deterministic results due to specialized network hardware being used end-to-end:
  - Highly accurate with granular results

**AmLight** ExP
Americas Lightpaths **Express & Protect**

Before we go any deeper, let's go through some concepts and technologies

AmLight ExP
Americas Lightpaths **Express & Protect**

# Traffic Generators/Network Testers

- Also known as <u>network performance testers</u>, they are appliances with specialized application and hardware focused on benchmarking network performance and reliability, as well as protocols, devices, and applications.

- Traffic Generators perform packet creation and processing entirely on specialized ASIC/FPGA to achieve deterministic results.

- Highly flexible in terms of packet creation: types of packets, size, headers, number of packets, packet rate/bandwidth, and even customizable payloads.

- Traffic generators have APIs to support remote integration.

- During SC23 and SC24, SCinet had access to EXFO, Viavi, and KeySight solutions to test the WAN links.

**AmLight** ExP
Americas Lightpaths **Express & Protect**

# Traffic Generators/Network Testers



- **AmLight has an EXFO FTB-1 NetBlazer with 4x100G interfaces**
  - One 2x100G module for experimentation/testbed
  - One 2x100G module for BERToD/production

- **AmLight created Python wrapper to use EXFO's SCPI API**
  - SCPI (Standard Commands for Programmable Instruments) is no fun!

- **BERToD uses two applications: EtherBERT and MonGen.**
  - For EtherBERT, PRBS31 is supported (high accuracy)
  - Others are supported, such as RFC2544 and RFC6349

# In-band Network Telemetry (INT)

- INT is a streaming telemetry solution based on P4 that records network telemetry data in the packet while the packet traverses a path between two points in the network

- Telemetry is exported directly from the Data Plane and the Control Plane is not affected:
  - Translation: you can track/monitor/evaluate EVERY single packet at line rate and in real time.

**AmLight** ExP
Americas Lightpaths Express & Protect

# In-band Network Telemetry (INT)

# SDN: Building paths over AmLight

~~Before we go deeper, let's go through some concepts and technologies~~

Ok, let's move on.

AmLight ExP
Americas Lightpaths Express & Protect

# BERToD – Components Explained

- **Reads testing policies**
  - Test duration, application, packet length, TCP/IP headers, number of packets, maximum bandwidth, scheduling, etc.

- **Kytos-ng:**
  - Generates the current network topology and instantiates testing paths over the network, from the packet generator to remote loops.

- **BAPM (Behavior, Anomaly, and Performance Manager)**
  - Generates the network state based on the telemetry sources available. For instance:
    - Identify the bandwidth available and buffer utilization for each interface based on the last 30 seconds of utilization.

- **Packet Generator**
  - Sends packets based on the policy (next slide)

# BERToD - Bit Error Rate Test on Demand [2]

- Test every possible link every 30 min:
  - Latency, jitter, frame loss, and out-of-sequence tests
  - Multiple frame sizes: 68, 256, 512, 1024, 1518, 9000 bytes
  - Each test runs for up to 10 seconds, and we send up to 500,000 frames
    - In case a test fails, run again with a multiplier metric (for instance, 3)
  - Choice for max bandwidth comes from BAPM
    - Up to 50% of the available bandwidth based on the last 30 seconds (and up to 40 Gbps)

- Displaying results:
  - Last hour, Last 7 days, heatmap, and text outputs

- Grafana Annotations are used to document known topology events and actions to help correlate events.

**AmLight** ExP
Americas Lightpaths **Express & Protect**

# BERToD – Granular Individual Results

- Using Grafana to plot each test's loss, jitter, latency, and out-of-sequence
- Great way to understand the last 24 hours
- Filters available to visualize test results based on frame size and individual paths
- Not great for correlating events

# BERToD – Historical Results

- Using Grafana to plot each test's loss per day

- Great way to correlate events and identify patterns

- Filters available to visualize test results based on frame size and individual paths

- Used with annotations to add context

# Lessons Learned

- **Testing infrastructure vs testing user experience**
  - To achieve deterministic results, network resources must be fully available.
  - Even with buffer occupancy monitoring, traffic engineering at AmLight had to be enhanced to avoid new drops ➔

- How to monitor the actual user experience?
  - Using perfSONAR and BERToD in the same queue as users

- "*There is no such innocent maintenance at the datacenter*".
  - Mishandled patch cords are the main reason for sudden spikes of errors ➔

- Dirty fiber/connector is the main reason for discrete errors (<0.0001%)

- Some vendors have weird policies, and small frames are delivered out of sequence (under investigation)

| AmLight Traffic Prioritization Policy | |
|---|---|
| Queue 7 | Reserved for future use. |
| Queue 6 | Reserved for management traffic |
| Queue 5 | Reserved for future use. |
| Queue 4 | Reserved for "deterministic" monitoring (BERT). |
| Queue 3 | Vera Rubin Observatory over shared links. |
| Queue 2 | Reserved for more than best effort. Not in use. |
| Queue 1 | (Default) Best effort traffic & BERToD for users |
| Queue 0 | Less than Best Effort. Experiments/Microbursts |

% of Frame Loss Grouped by Day



< 1e-7%   1e-7%+   0.00000100%+   0.0000100%+   0.000100%+   0.00100%+   0.0100%+   0.100%+

**AmLight** ExP
*Americas Lightpaths* **Express & Protect**

Tests using Priority Queue BE    Tests using Priority Queue 4

# Next steps: Re-Testing after a link flap

- Integration with Kytos-ng SDN Controller to test links after **each link flap**:
  - The goal is to evaluate if the link is clean after a maintenance/repair before using it again!
    - After a link flaps, the SDN controller waits up to 2 min to confirm the link is stable and then initiates the quarantine mode
    - BERToD is notified of the quarantine and starts a 5 min test
    - If results are clean, BERToD sets the link as operational/ready
    - The SDN controller then makes the link available to all applications

AmLight ExP
Americas Lightpaths **Express & Protect**

# Next steps: Automate fault isolation

Automate the fault isolation process using all data sources available

- SDN logs, topology changes, EVC optimizations, events/demos, optical monitoring, and visits to the data center.

- An AI/ML researcher is playing with our testing dataset and log entries to isolate issues.



% of Frame Loss Grouped by Day

# Next steps: Automate fault isolation [2]

How to create a correlation between a, for instance, high utilization queue and the BERToD test results?

# Next Steps: Comparing results with perfSONAR

- We are still learning how to compare the results provided by BERToD and perfSONAR.

- An example to share:
  - BERToD started reporting consistent errors when testing the leased link between Fortaleza, Brazil and Miami.
  - The issue started on Jan 14th, 2025. After a maintenance on the Jan 24th, the errors disappeared (unknown cause).
  - Although the errors reported were around 0.0002%, we believe perfSONAR should have detected something before the 24th.

# Conclusion

- BERToD is a fantastic addition to the network monitoring portfolio thanks to the hardware-based traffic generator and enhanced network telemetry provided by the AmLight SDN solution.

  - Production since September 2024. Used daily by AmLight OPS.

- Having a hardware-based traffic generator enables quick testing with extreme accuracy

  - Helps us follow the demands of our SLA-driven science drivers

- BERToD is a great complement to perfSONAR @ AmLight.

  - While perfSONAR allows AmLight to test applications and protocols with excellent per-direction visibility, BERToD provides extreme performance visibility for applications over ultra-long paths where any packet loss causes damage.

**AmLight** ExP
Americas Lightpaths *Express & Protect*

# BERToD - Bit Error Rate Test on Demand [5]

- **Per-Hour Heatmap visualization created to help identify patterns across tests**



- **Command-line to access full results and test configuration**

**AmLight** ExP
Americas Lightpaths **Express & Protect**