# The perfSONAR Power Hour

Andy Lake – ESnet – andy@es.net

*perfSONAR is developed by a partnership of*

# What is perfSONAR?

- An **open source software collaboration** led by ESnet, GEANT, Indiana University, Internet2, RNP and the University of Michigan.

- **Goal is to provide network measurements between organizations** to help identify and troubleshoot network issues. Most commonly these include (but are not limited to):
  - Throughput
  - Packet Loss
  - One-way latency
  - Traceroute

# perf5.⊕NAR

- perfSONAR 5.0 released **April 17th, 2023.**

- **Over 50% of perfSONAR deployments** currently running 5.0

- **Enables greater visualization and analysis capabilities** through the replacement of the backend measurement storage database with OpenSearch

- **New pScheduler test plugins** to support WiFi BSSID, 802.1X authentication, DHCP response time and more

- **Ubuntu 20 support added** with additional OSes like
  - EL8 and EL9 added in early summer
  - Ubuntu 22 and Debian 11 coming soon

- **Looking ahead, 5.1 will focus on improving UI** and leverage the changes put in place by 5.0 to add new capabilities
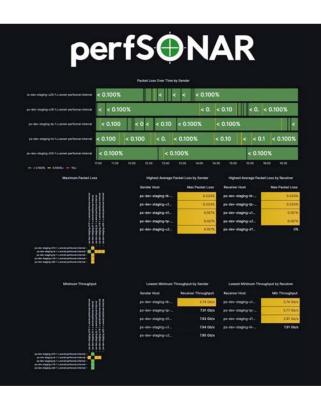


*Image: Example of Grafana dashboard users can setup in 5.0 using our guide*
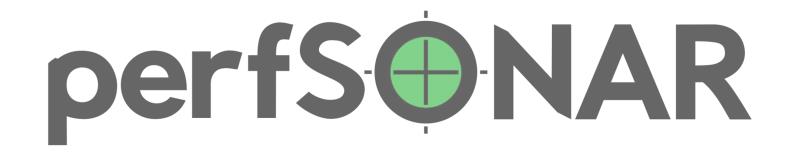
# Building on the Foundation

**UI:** Toolkit UI and MaDDash Beyond 5.0

**UI:** On-Demand Testing

**Tool Enhancements:** Multi-threaded iperf3

**Deployment Models:** perfSONAR on Internet2 Backbone

# Toolkit UI and MaDDash Beyond 5.0

Andy Lake – ESnet – andy@es.net

*perfSONAR is developed by a partnership of*

# Today's Toolkit UI

# Today's MaDDash UI

# Thinking about a new interface

- User Requirements
  - Users need to still have ability to view test results of local tests
  - Authentication required to see additional system stats, configure tests, and configure LS registration
    - Whole subset of requirements around configuring test that won't cover in this deck
  - System information needs to be available via an API (see current JSON service)
- Dev Requirements
  - Devs need to get rid of old perl cgi scripts that are increasingly difficult to maintain
  - Ditto old Javascript
  - Consolidate redundant interfaces (e.g. PWA and Toolkit Test Config)
  - Stop discovering system stats that differ across OSes, when other things have this solved (eg. node_exporter)
- Other nice to haves
  - Users have greater flexibility in building own views of data
  - Tighter integration with MaDDash and main toolkit UI
  - We don't have a lot of UI devs, if there is a way to make easier to maintain without being JS expert, that is good thing.

# What is Grafana?

- Grafana is an **open-source platform** for for exploring data from a variety of sources
- It has a few key features
  - It's **multi-data source**
  - It has a bunch of **built-in visualizations** that **don't require you to be a Javascript developer** to use
  - It has a **plugin framework** for all of the above so they can be extended and a process for becoming official plugins
- More Info: https://grafana.com/docs/grafana/latest/introduction/

Custom Plugins by NetSage and ESnet

# New UI: Focus on Measurements

# New UI: Enhancing the Fundamentals

# New UI: Instrumentation

# New UI: Data Correlation

iperf3 overlay on host metrics

# New UI: MaDDash Integration

# New UI: Customization

# Multi-Threaded iperf3

*First There was One, Now There are Many*

Bruce Mah and Sarah Larsen, ESnet

*perfSONAR is developed by a partnership of*

# What is iperf3?

- iperf3 is an open-source tool that measures network traffic performance between a client and server

- It was designed to be used in perfSONAR, but can also be used on its own

- Linux, MacOS, and FreeBSD are all officially supported

# Current State

- Before multi-threading, iperf3 was capable of 30-50 Gbps, with single stream TCP, possibly more with tuning

- Many links in ESnet6 are faster than iperf3
  - Site connections: N x 100G
  - Backbone: 400G +

- So what if we want to support connections with higher bandwidth? How do we get more throughput?

- The problem: Adding more parallel connections doesn't increase throughput

# What are we doing about it?

- Adding multi-threading to iperf3
  - Goal: Use multiple CPU cores

Example command: iperf3 –parallel 3

- Note: currently this refers to multiple parallel streams, but in the new multi-threaded iperf, this will represent the number of threads. This is because in the new version each stream will get its own thread.

# Performance Testing

- Test Setup:
  - Computers:
    - Back-to-back, no routers
    - 200 Gbps link
    - 16 cores, 2 packages, 32 cores total
    - More cores than parallel streams
    - Mellanox ConnectX6
    - ~3 GHz AMD CPU

200Gbps

# Effects on Performance

iperf3-mt has significants improvements in throughput performance over iperf3.

# How does it compare to iperf2?

Multi-threaded iperf3
has similar throughput
results to iperf2

# When?

"Soon"

But seriously Q4 2023/Q1 2024

We're going to test it on ESnet perfSONAR hosts before releasing the final version.

Any interest in testing or want the current working version:
https://github.com/esnet/iperf/tree/mt

# Summary

The new multi-threaded iperf3 can better test performance of faster paths.

GitHub: https://github.com/esnet/iperf

Multi-threaded: https://github.com/esnet/iperf/tree/mt

Contact email: iperf@es.net

# Questions?

# perfSONAR on the Internet2 Backbone

Mark Feit – Internet2 / perfSONAR Development Team
mfeit@internet2.edu

2023 INTERNET2 TECHNOLOGY exchange

# WE DID A THING

# A Plethora of Public perfSONAR Points in PoPs



PNWGP Seattle/Starlight

University of Montana Missoula

Link Oregon Portland/Boise

UETN Salt Lake City/Las Vegas

Nevada System of Higher Education Las Vegas/Reno

Pacific Wave Denver

CENIC Los Angeles/Sunnyvale/Sacramento

MREN Starlight/Starlight

University of Illinois- Urbana-Champaign Chicago/Starlight

Big Ten Academic Alliance OmniPoP Chicago/Starlight

University Of Memphis Memphis

GPN Kansas City/Tulsa

Sun Corridor Network Tucson/Phoenix

MissiON Jackson/Atlanta

LONI Baton Rouge

LEARN Dallas/Houston

OARnet Cincinnati/Cleveland

Merit Network, Inc. Toledo/Chicago

Northern Lights GigaPOP Minneapolis

Indiana GigaPOP Chicago/Indianapolis

NoX Albany/New York City

CEN Hartford/New York City

NYSERNet Buffalo/New York City

KINBER Philadelphia

Drexel University Philadelphia

OSHEAN New York City

NJEdge Philadelphia/New York City

MAGPI Philadelphia/Pittsburgh

University of Pittsburgh Pittsburgh

CAAREN Ashburn

MAX Ashburn//Washington

KyRON Louisville

MARIA Ashburn

MCNC/C-Light Charlotte/Raleigh

SoX Atlanta/Nashville

FLR Jacksonville/Pensacola

## Using the Internet2 Public perfSONAR Nodes

# *PoP*`.ps.internet2.edu`

- IPv4 and IPv6 available

- Log into a perfSONAR system (yours, not ours) and run a task with pScheduler

```
pscheduler task throughput --dest PoP.ps.internet2.edu
```

# PoP Directory

| | | | | | |
|---|---|---|---|---|---|
| **Albany** | alba | **Hartford** | hart2 | **Pensacola*** | pens |
| **Ashburn** | ashb | **Houston** | houh | **Philadelphia** | phil |
| **Atlanta** | atla | **Houston** | hous | **Phoenix** | phoe |
| **Boise*** | bois | **Indianapolis** | indi | **Pittsburgh** | pitt |
| **Boston**** | bost | **Jackson** | jcsn | **Portland** | port |
| **Charlotte** | char | **Jacksonville** | jack | **Raleigh** | rale |
| **Chicago** | chic | **Kansas City** | kans | **Reno** | reno |
| **Chicago** | eqch | **Las Vegas** | lasv | **Sacramento** | sacr |
| **Chicago** | star | **Los Angeles** | losa | **Salt Lake City** | salt |
| **Cincinnati** | cinc | **Los Angeles** | losa2 | **San Jose** | sanj |
| **Cleveland** | clev | **Lousville** | loui | **Seattle*** | seat |
| **Dallas** | dall | **Minneapolis** | minn | **Sunnyvale** | sunn |
| **Dallas** | dall3 | **Missoula** | miss2 | **Toledo** | tole2 |
| **Denver** | denv | **Nashville** | nash | **Tucson** | tucs |
| **El Paso** | elpa | **New York** | newy2 | **Tulsa** | tuls |
| **Fargo**** | farg | **New York** | newy32aoa | **Washington** | wash |

*Not yet in service     **Future

# Topology

| PoP Type | Router Connection |
|---|---|
| Distributed | Core* |
| Multi-Degree | |
| Interconnect | |
| Split Interconnect | Aggregation* |

*First where more than one is installed

**Network Reachability**

<u>Now</u>
R&E + I2PX

<u>Later</u>
Elsewhere

# Initial Administrative Limitations

- Caps on **throughput** test bandwidth
  - R&E*        10 Gb/s
  - Elsewhere 1 Gb/s
  - Higher bandwidth considered on a case-by-case basis

- No, **disk-to-disk**, **s3throughput**, **idleex** or **wifibssid** tests

- Tests with a **duration** parameter are limited to one minute

- Repeating tests
  - Not more-frequently than once per hour
  - Must wrap up within 24 hours
  - **repeat-until** parameter not allowed

- Testing priorities
  - Internet2 Internal
  - R&E Networks
  - Everyone Else

- These limits will be refined periodically to make sure the community's needs are being met.

*Determined using ESnet's R&E network list: `http://stats.es.net/sample_configs/pscheduler/ren`

# Beta Period

- Now through January

- Feature set close to production

- Feedback is appreciated
  ### pas@internet2.edu

# Beta Will Be Beta

- Working through some teething with the NICs in some PoPs
  - Systems disappear from the network

- Ongoing experimentation on systems in **CLEV**, **PHOE** and **TUSC** PoPs

β

# Performance Tuning

### New OS

### + New Kernel

### + New NIC Driver

### = New System Tunings

- Higher-speed throughput requires additional test parameters to run well

# THE SYSTEMS

# A ~~Chicken~~ System in Every ~~Pot~~ PoP

- Dell R6515 (NGI Buildout) and R6615 (Later)
- AMD EPYC 7402P 24 Cores / 48 Threads at 2.8 GHz
- 128 GiB RAM

- Broadcom 2x 10 GbE
- Mellanox Connect-X 5 2x 100 GbE

- AlmaLinux 9

*Connections to routers vary by PoP type.*

# Host Architecture

# ~~Special~~ Secret Sauce: The `macvlan` Network Driver

- Binds a host interface directly into a container

- Bypasses additional container networking code

- Negligible performance difference vs. bare metal

- No address assigned on the host
  - Prevents external access to the OS

**Big MACVLAN**

# Resources for Performance

- ## 100 GbE Interface Containers
  - ### 12 Dedicated CPU Cores (Threads)
  - ### 32 GiB Dedicated RAM
  - ### IRQ Affinity (Tuning)
  - ### 93+ Gb/s

- ## 10 GbE Interface Containers
  - ### Shared CPU Cores
  - ### Shared Memory
  - ### No special performance tuning

Cores

Memory

IRQ Affinity

**perfSONAR Container**

100 GbE Interface

**perfSONAR Container**

10 GbE Interface

# Deployment Technology Stack(s)

- Single data set with detailed information on PoPs, systems, interfaces and networks
  - System kickstart files
  - Ansible
  - Salt
  - Assorted shell scripts
  - DNS records
  - Internal proxy configuration and ACLs

- Ultimate goal is to use Salt

# Why so many different technologies?

- The tools are buggy.

- IPv4 /31s used for point-to-point connections
  - Halves address use compared to /30s
  - See RFC 3021

- Docker couldn't handle those at all
  - Patch submitted in 2021, released in Docker 23.0.0 (2022).

- Podman is fine with them but its web API isn't.

# The Long and Winding Road

- Began with Ansible for expediency
  - Problems with networks being re-created with each run
  - Destroyed/rebuilt the container. *No bueno.*

- Tried Salt
  - Stymied by the /31 problem

- Ended up with a set of shell scripts
  - Tied together for single-command provisioning of the entire system
  - Small bites that can be easily converted to Ansible or Salt

# Future Development: Deployment Kit

- Will be derived from deployments at Internet2 and elsewhere
  - Based around Ansible
  - Minimal manual perfSONAR host configuration
  - Driven by data: Configure and go
  - Container-per-interface model (Plain or VLANs)

- Initial version targets Debian
  - EL to follow

**Thanks!**

Public perfSONAR

# *PoP*`.ps.internet2.edu`

Feedback

# `pas@internet2.edu`

# Multiple perfSONAR activites in GÉANT

- Lookup Service dashboards
- perfSONAR deployments
- Microdep integration
- On-demand perfSONAR Graphical User Interface (psGUI)

## Lookup Service dashboards

- Display, filter and search the content of the Lookup Service

- Grafana 8 based
  - Filtering on text, domains, communities
  - Stats on hosts and services, maps

- https://stats.perfsonar.net
  - Replaces ESnet Service Directory

- Next steps:
  - Port to Grafana 9
  - Filter on multiple values

# perfSONAR deployments in the GÉANT network (1/2)

- 10 public deployments on the core network: https://network.geant.org/perfsonar/



**Amsterdam**
psmp-gn-bw-ams-nl.geant.org
psmp-gn-owd-ams-nl.geant.org

**Budapest**
psmp-gn-bw-bud-hu.geant.org
psmp-gn-owd-bud-hu.geant.org

**Frankfurt**
psmp-lhc-bw-fra-de.geant.org
psmp-lhc-owd-fra-de.geant.org

**Geneva**
psmp-lhc-bw-gen-ch.geant.org
psmp-lhc-owd-gen-ch.geant.org

**Lisbon**
psmp-gn-bw-lis-pt.geant.org
psmp-gn-owd-lis-pt.geant.org

**London**
psmp-gn-bw-lon-uk.geant.org
psmp-gn-owd-lon-uk.geant.org

psmp-lhc-bw-lon-uk.geant.org
psmp-lhc-owd-lon-uk.geant.org

**London2**
psmp-gn-bw-lon2-uk.geant.org
psmp-gn-owd-lon2-uk.geant.org

**Paris**
psmp-gn-bw-par-fr.geant.org
psmp-gn-owd-par-fr.geant.org

psmp-lhc-bw-par-fr.geant.org
psmp-lhc-owd-par-fr.geant.org

**Poznan**
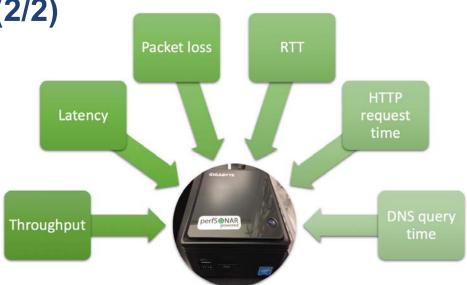psmp-gn-bw-poz-pl.geant.org
psmp-gn-owd-poz-pl.geant.org

**Vienna**
psmp-gn-bw-vie-at.geant.org
psmp-gn-owd-vie-at.geant.org
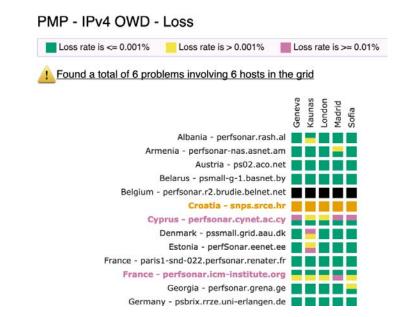
R&E Nodes
LHCONE Nodes

# perfSONAR deployments in the GÉANT network (2/2)

- Performance Measurement Platform (PMP)
  - Small nodes (Intel NUC) and VM
  - Deployed in GÉANT partners organisations
- Measurements
  - Diverse set of measurements
  - Regularly to GÉANT core network
  - Verify GÉANT access links
  - International connections (ESnet, Internet2, RNP, …)
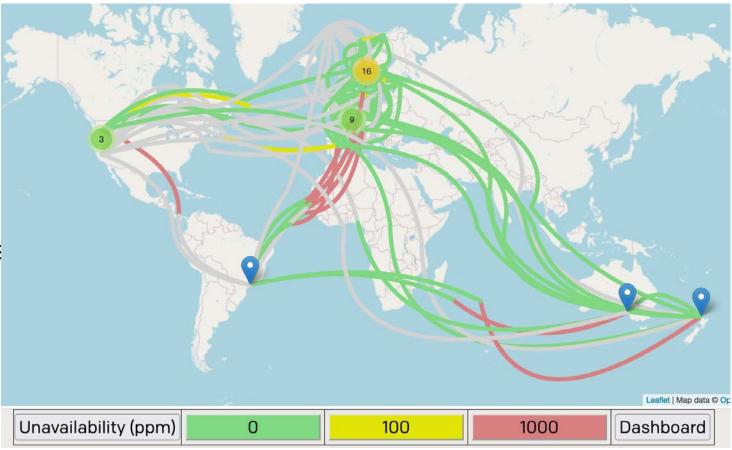  - 2nd tiers: University networks
- https://pmp-central.geant.org



**PMP IPv4 Dashboard**

PMP - IPv4 OWD - Loss

Loss rate is <= 0.001%  Loss rate is > 0.001%  Loss rate is >= 0.01%

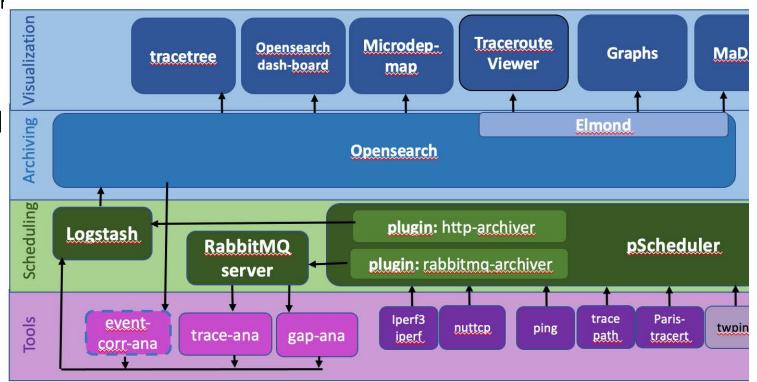⚠ Found a total of 6 problems involving 6 hosts in the grid

# Microdep integration with perfSONAR (1/2)

- Microdep is a packet loss analysis and visualisation tool
  - Spotting packet gaps, micro failures, ~10 packets loss
  - Using 100 packet/sec probes
  - Traceroutes and
    ICMP response monitoring
- Realtime event analysis:
  - Packet-loss (gaps)
  - Queues (jitter)
  - Route failures and changes
    (traceroute)
  - Joint event anomality and alarms
    (ELK and ML)



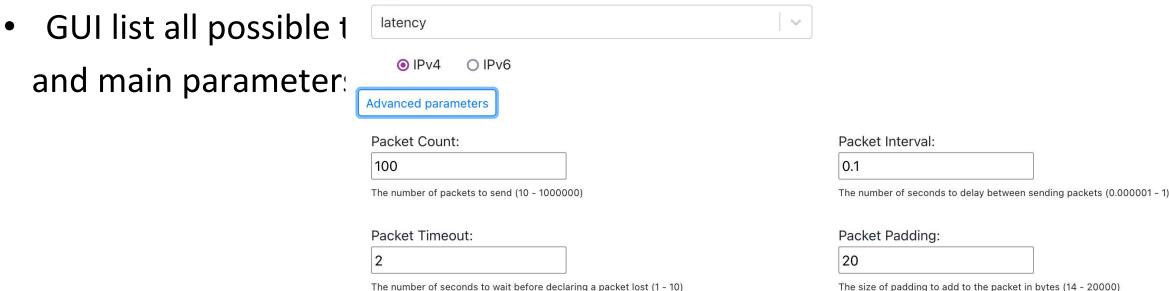| Unavailability (ppm) | 0 | 100 | 1000 | Dashboard |

# Microdep integration with perfSONAR (2/2)

- Using perfSONAR to generate probes
  - OWAMP for paced packets
  - Traceroute
  - Rely on 2000+ public perfSONAR hosts
  - Use pSConfig and pScheduler
- Adding a data pipeline to
  - Analyse packet gaps
  - Store history for further anal
- Next steps:
  - Package and bundle with pS

# On-demand perfSONAR Graphical User Interface (psGUI) (1/2)

- GUI to drive perfSONAR / pScheduler

- Use case:
  - MaDDash setup, grids, regular measurements
  - Want to do a one off, on-demand additional test
  - List of pS nodes coming from pSConfig file, MaDDash grids

- GUI list all possible t
  and main parameter

Test:

| latency | ⌄ |
|---------|---|

⦿ IPv4    ◯ IPv6

Advanced parameters

Packet Count:

100

The number of packets to send (10 - 1000000)

Packet Interval:

0.1

The number of seconds to delay between sending packets (0.000001 - 1)

Packet Timeout:

2

The number of seconds to wait before declaring a packet lost (1 - 10)

Packet Padding:

20

The size of padding to add to the packet in bytes (14 - 20000)

# On-demand perfSONAR Graphical User Interface (psGUI) (2/2)

- Results:

- Packaged as a Docker Image to be built

- https://github.com/perfsonar/psgui/

perfSONAR

THANKS FOR JOINING US!

Cartoon text courtesy of textstudio.com

69

ESnet  GÉANT  INDIANA UNIVERSITY  INTERNET2  RNP  UNIVERSITY OF MICHIGAN