

Potential LHC Networking SC24 NRE Demonstrations

Edoardo Martelli, Carmen Misa Moreira, Joe Mambretti,
Bruno Hoefft, Tom Lehman, Shawn McKee, Marian Babik,
Vitaliy Kondratenko, Tristan Sullivan, Phil Demar, Syed Asid
Shah, et al,

LHCOPN-LHCONE MEETING #52

CITTADELLA UNIVERITARIA

CANTANIA, ITALY

APRIL 9-11, 2024



Planning for SC24 (Atlanta, Georgia)

- ▶ **SCinet Sponsored Network Research Exhibition (NRE) Descriptions (June 1, 2024)**
- ▶ **NRE Submissions Define Demonstrations**
- ▶ **Lead To Assessment of Required Resources, Including WANs, Edge Devices (Enhanced Descriptions Planned for SC24)**
- ▶ **Results: Implementation of Services/Resources**
- ▶ **Assists With Pre-Conference Staging Facilities**



NREs: Verifying/Authenticating New Advanced Concepts

- ▶ Formulating New Architecture, Services, Techniques, Technologies Through Large Scale, WAN Demonstrations
- ▶ Proving Concepts With Empirical, Replicable Experiments
- ▶ Creating Prototypes
- ▶ Communicating Results To Wide Audiences
- ▶ Leveraging Large Scale Testbeds, e.g., SCinet
- ▶ Contributing To The Design and Implementation of Testbeds
- ▶ Building Blocks for Global Research Platform (GRP)

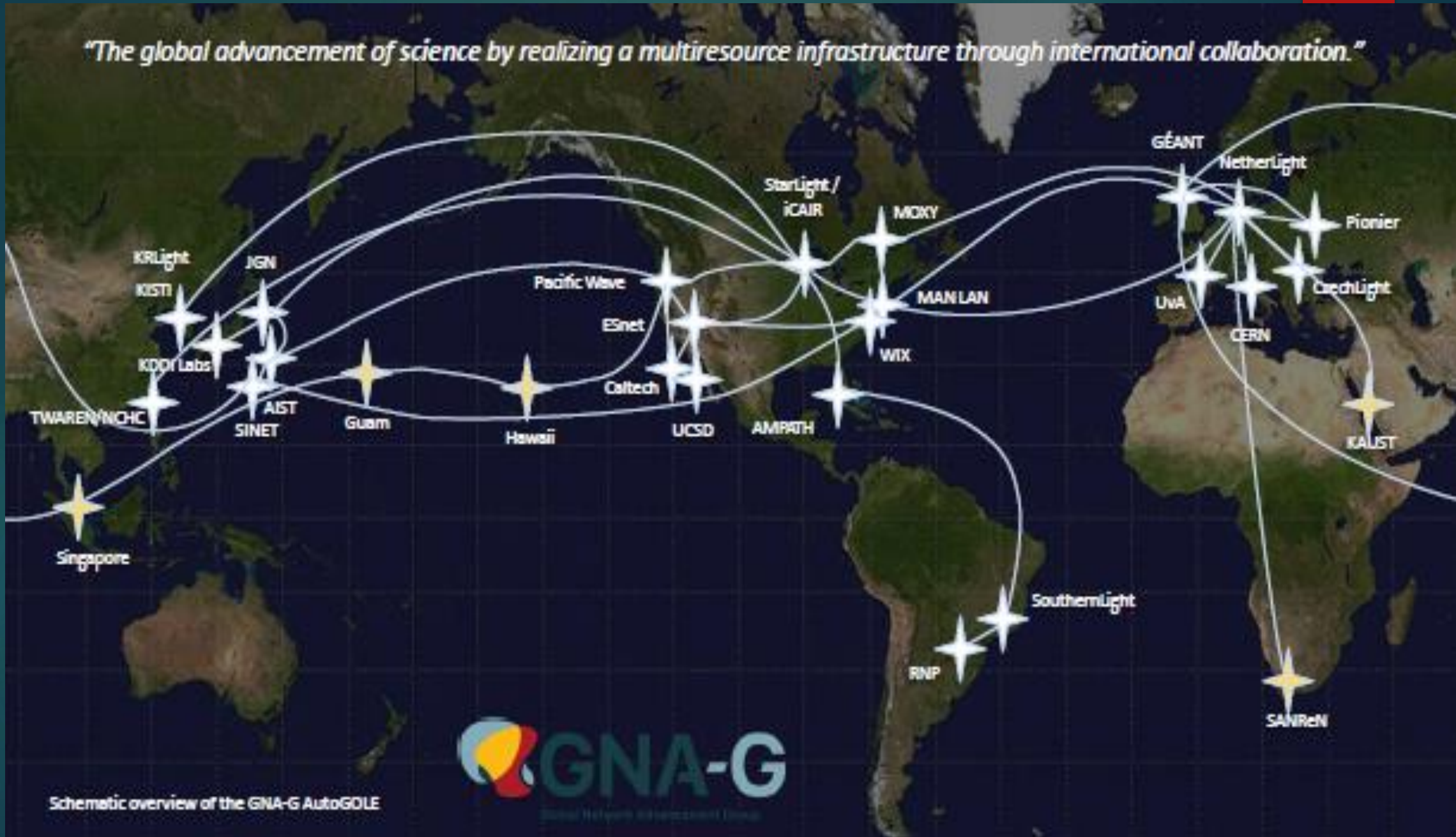
The GRP: A Platform For Global Science



GLOBAL RESEARCH PLATFORM

*A Next Generation, Software Defined,
Globally Distributed, Multi-Domain
Computational Science Environment*

"The global advancement of science by realizing a multiresource infrastructure through international collaboration."



Schematic overview of the GNA-G AutoGOLE

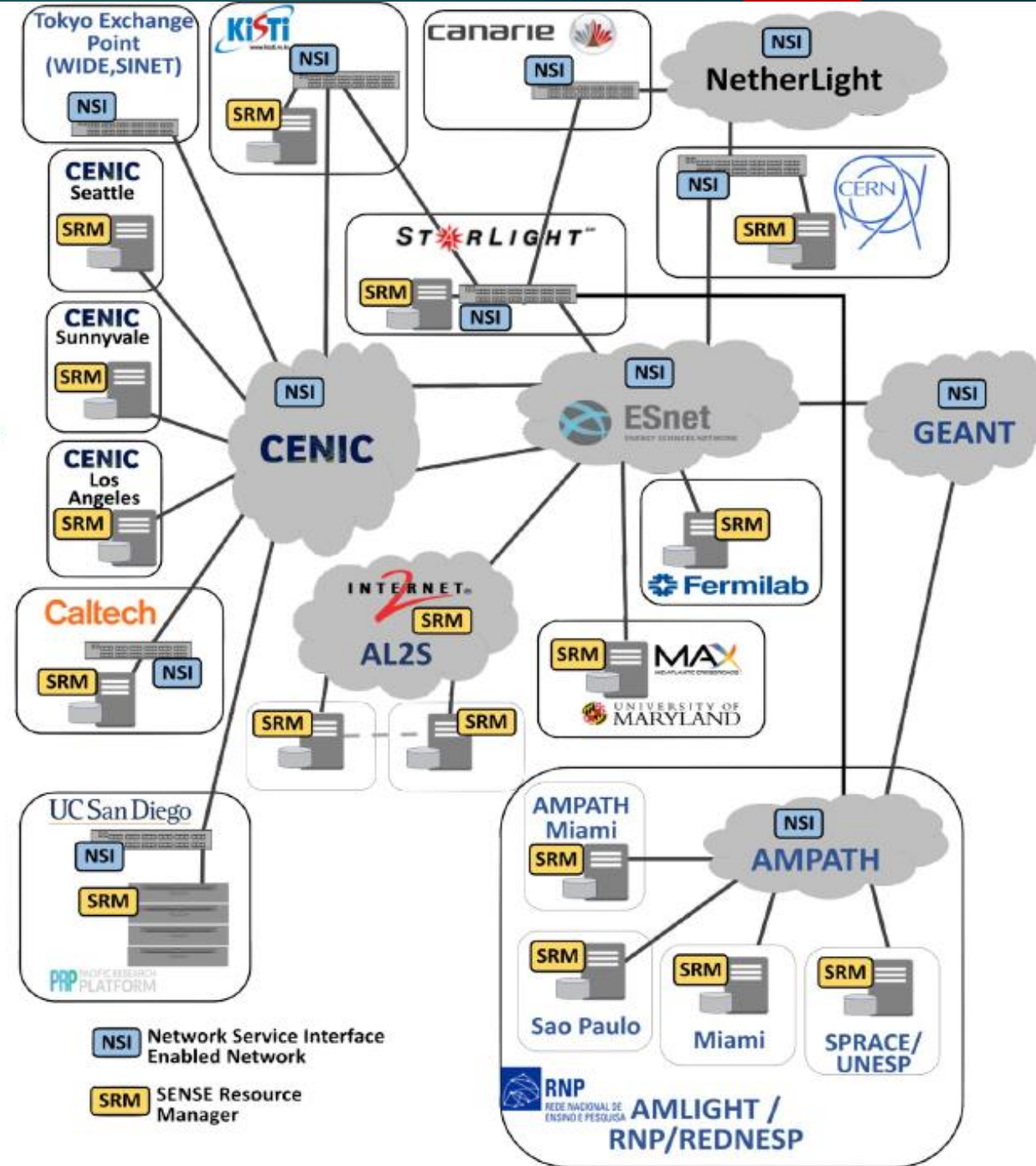


Global Research Platform/AutoGOLE Open R&E Exchanges

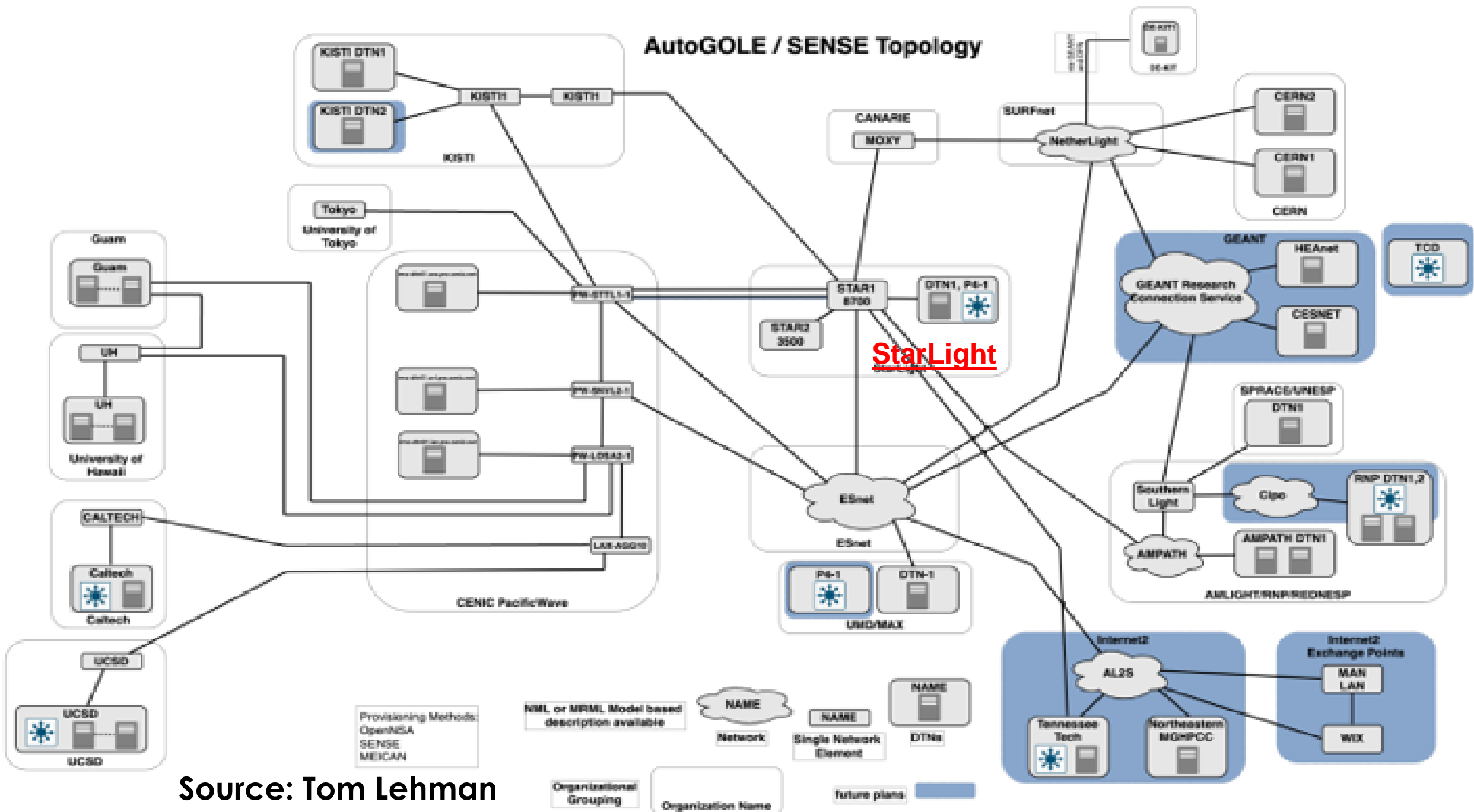
SENSE/AutoGole

- AutoGOLE, NSI, and SENSE working together provide the mechanisms for complete end-to-end services which includes the network and the attached End Systems (DTNs).

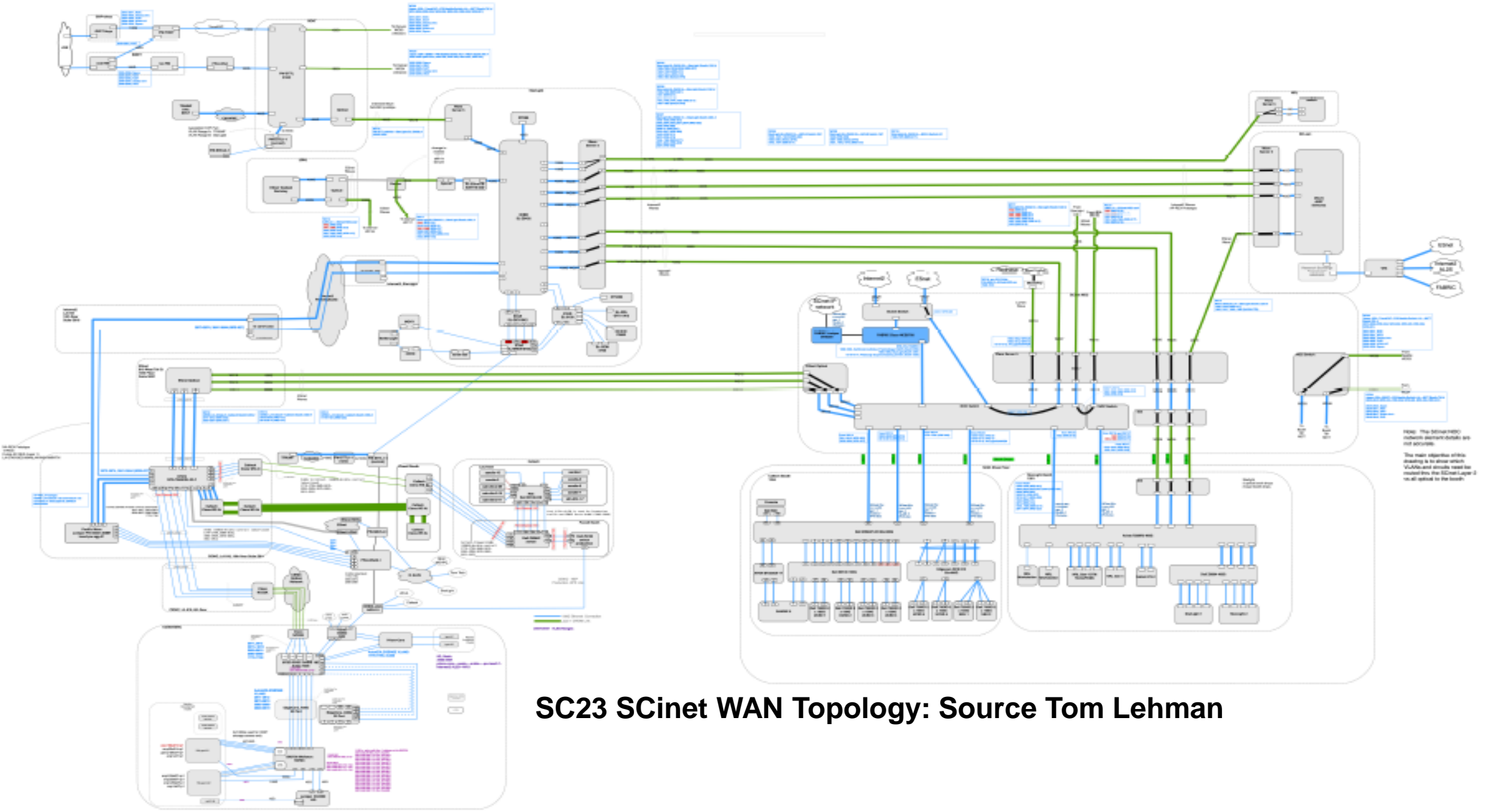
Source: Tom Lehman



AutoGOLE / SENSE Topology

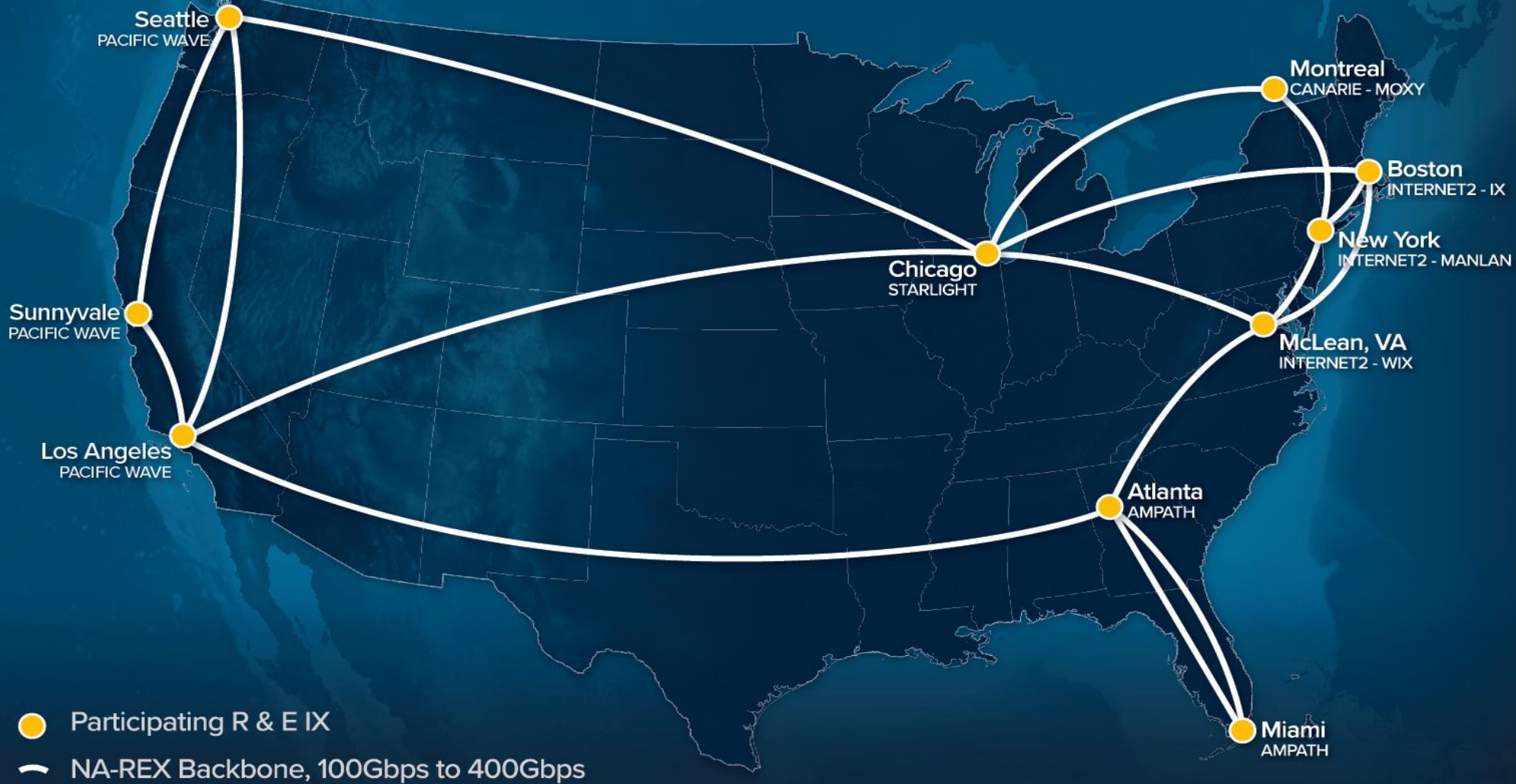


Source: Tom Lehman



SC23 SCinet WAN Topology: Source Tom Lehman

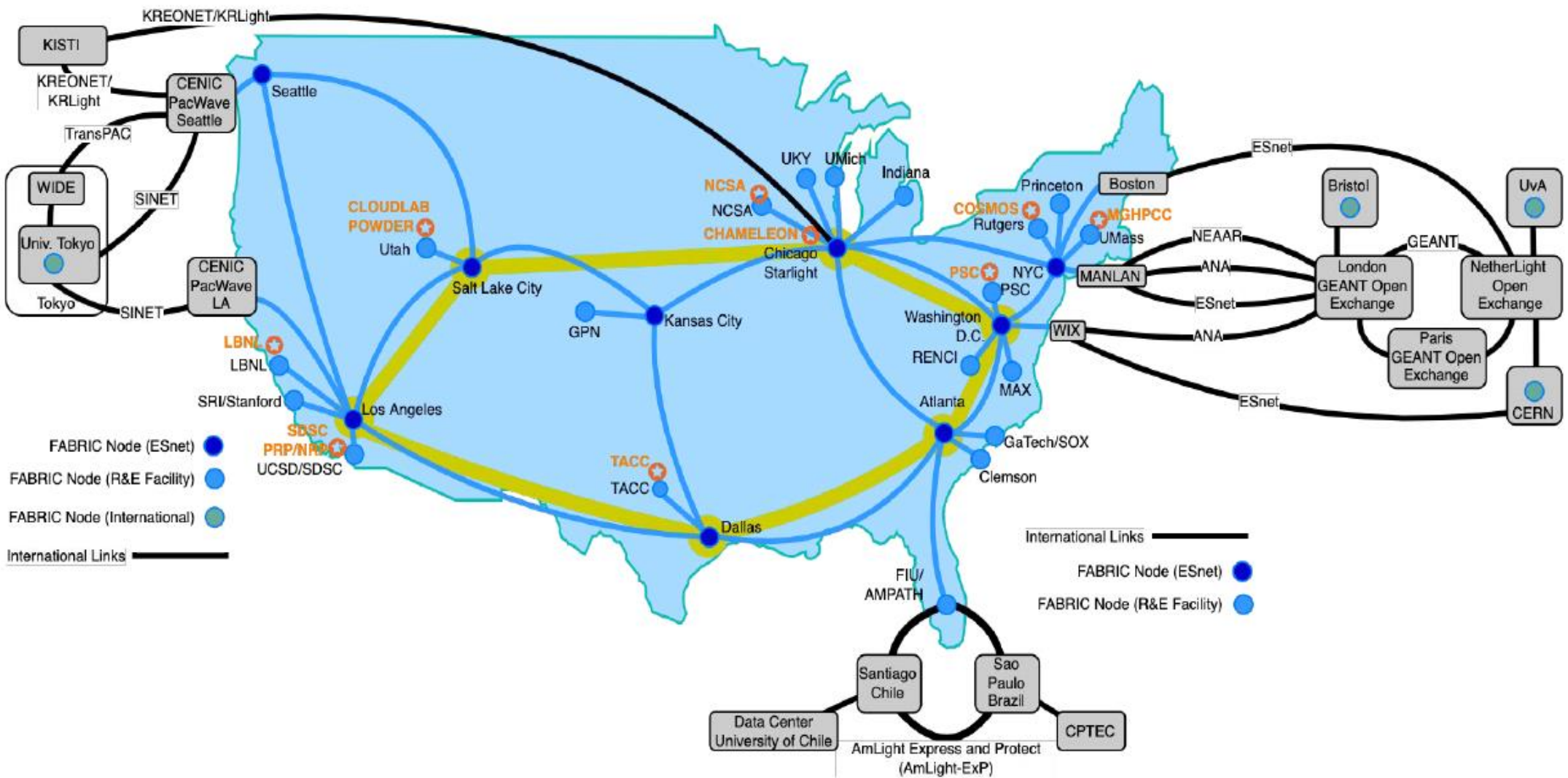
NA-REX North America Research & Education Exchange Collaboration



International Networks
at Indiana University

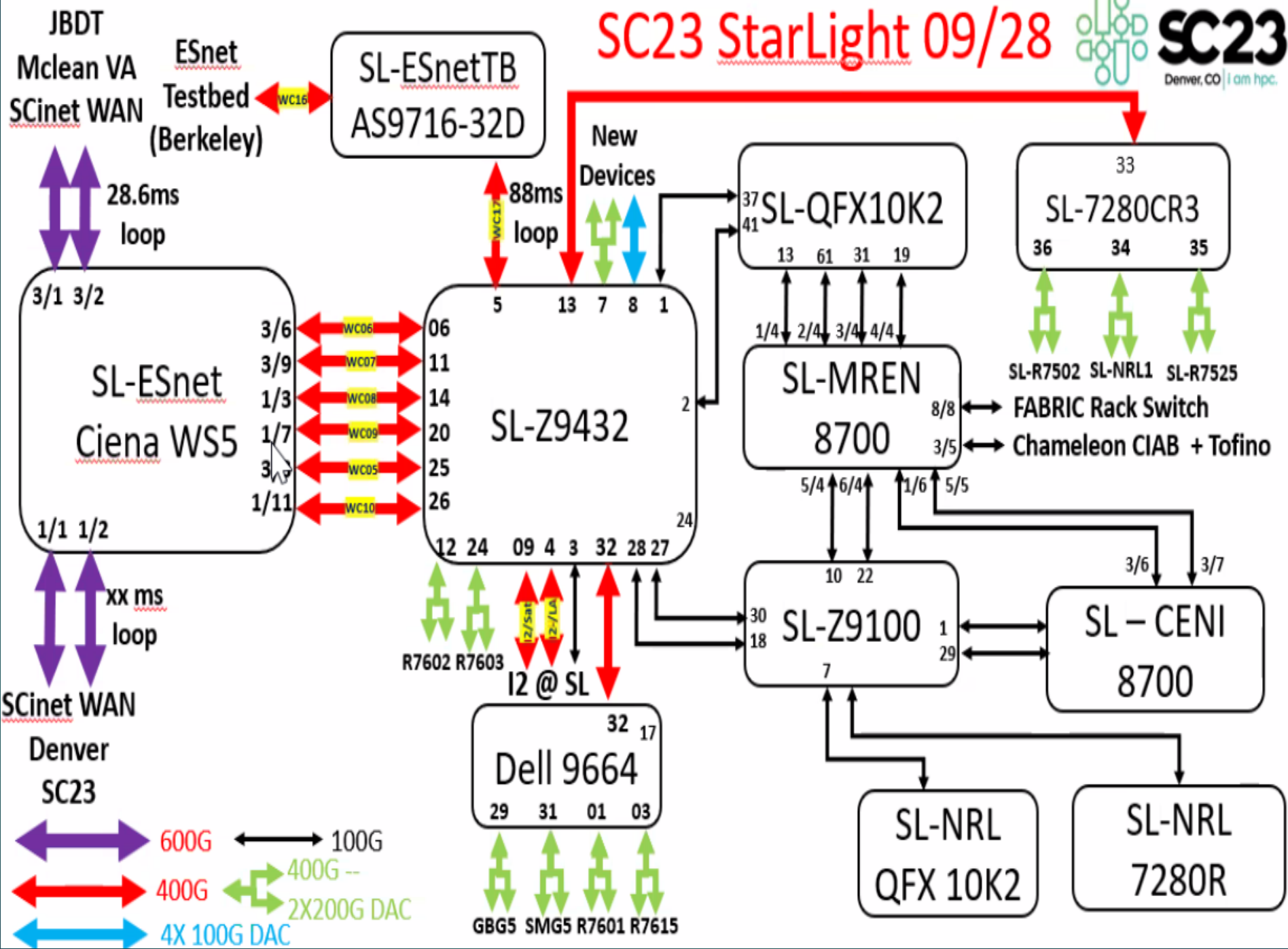


FABRIC Testbed (+FAB)



FABRIC Topology - with FAB Sites

SC23 StarLight 09/28



SENSE provisioning system

SENSE (SDN for E2E Networked Science at the Exascale): provision system that dynamically builds end-to-end virtual guaranteed networks across administrative domains without manual intervention.

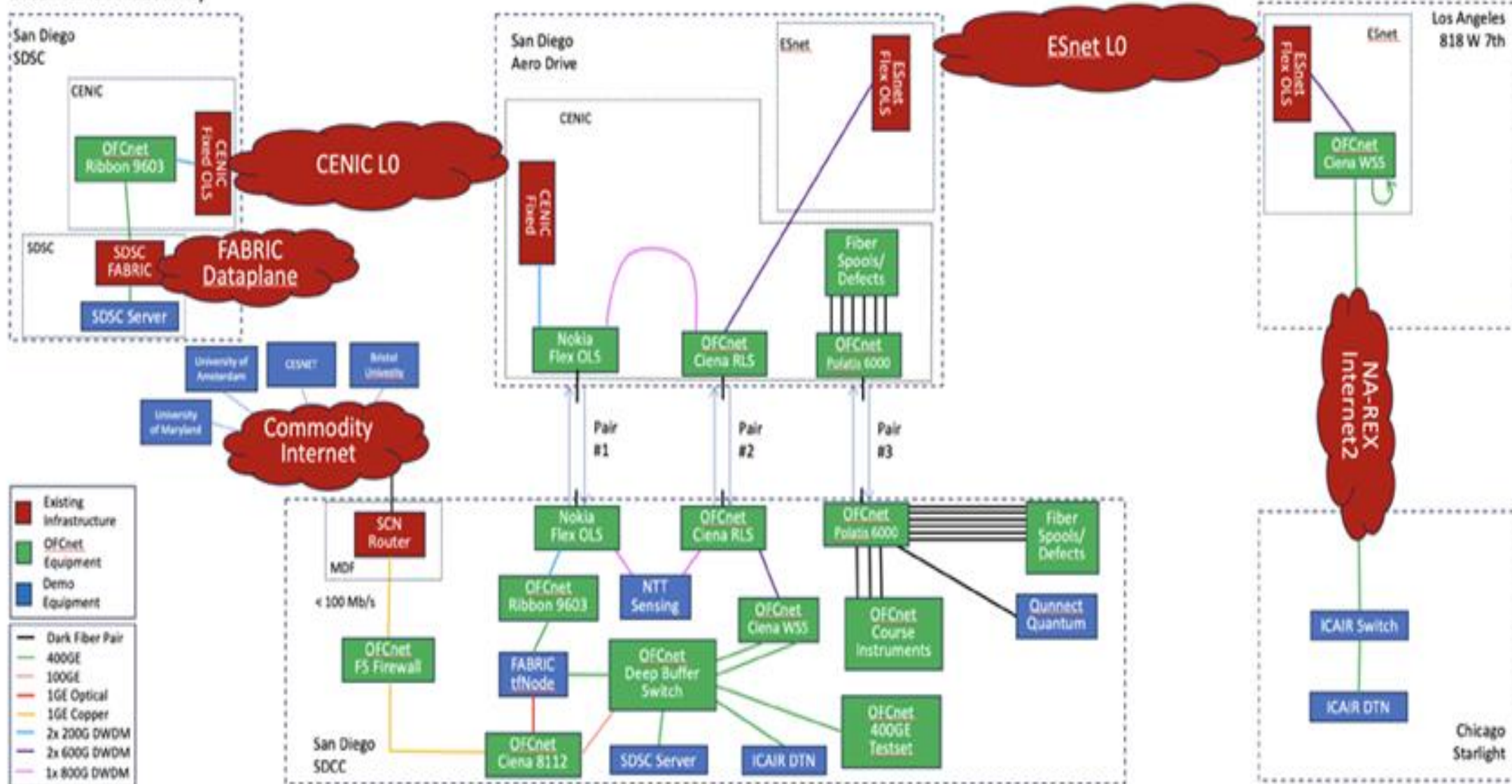
- ❑ Provisioning automation: bring-up and management of services without human involvement.
- ❑ Multi-domain: multiple administrative domains, independent policies and AUP (Acceptable Use Policy).
- ❑ Resource orchestration: allocation and reservation of resources including compute, storage and network.
- ❑ End-to-end: DTN NIC to DTN NIC, across Science DMZ (Demilitarized zone), WANs, Open exchange points...

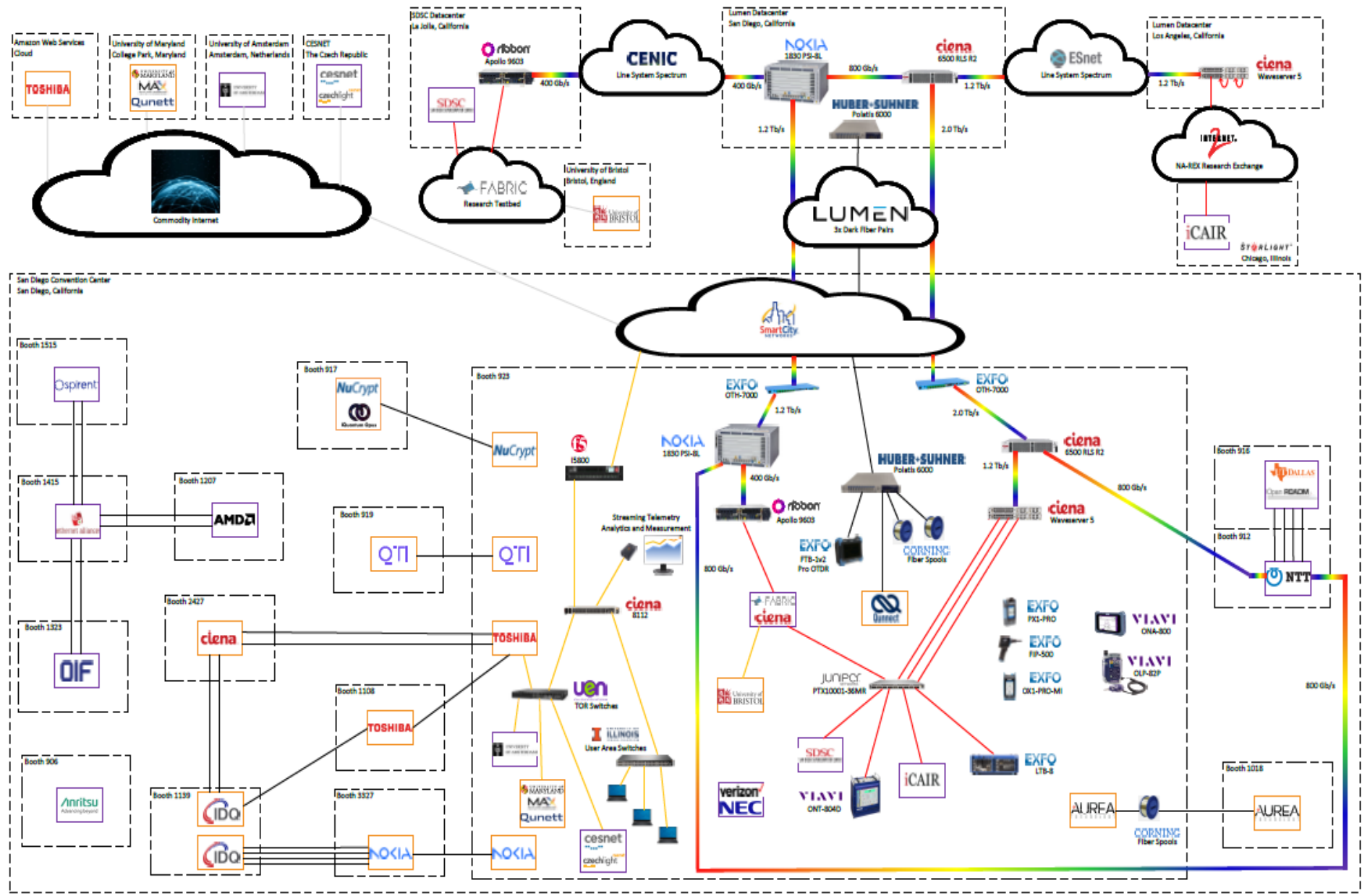


ESnet

ENERGY SCIENCES NETWORK

External Connectivity

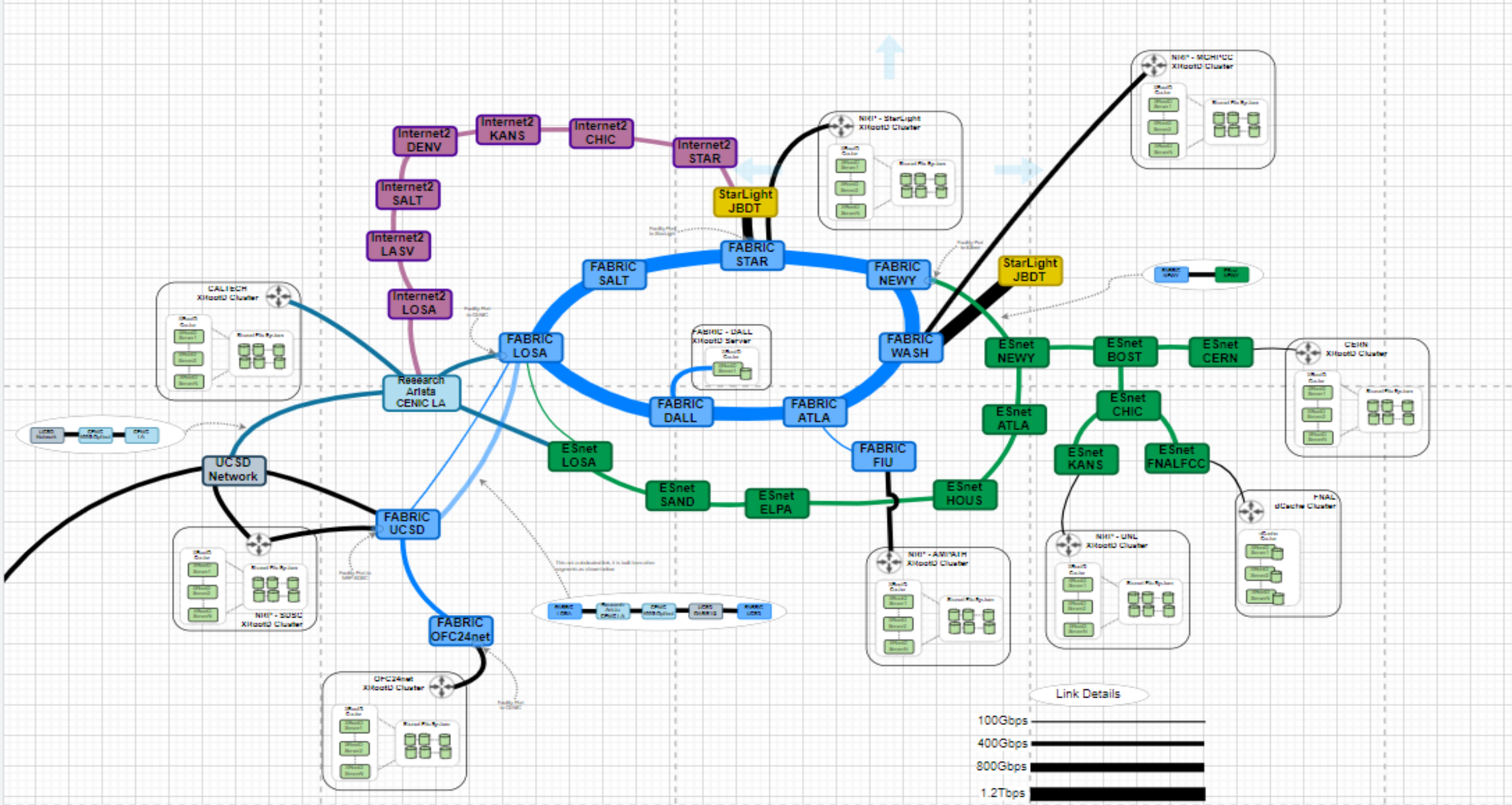




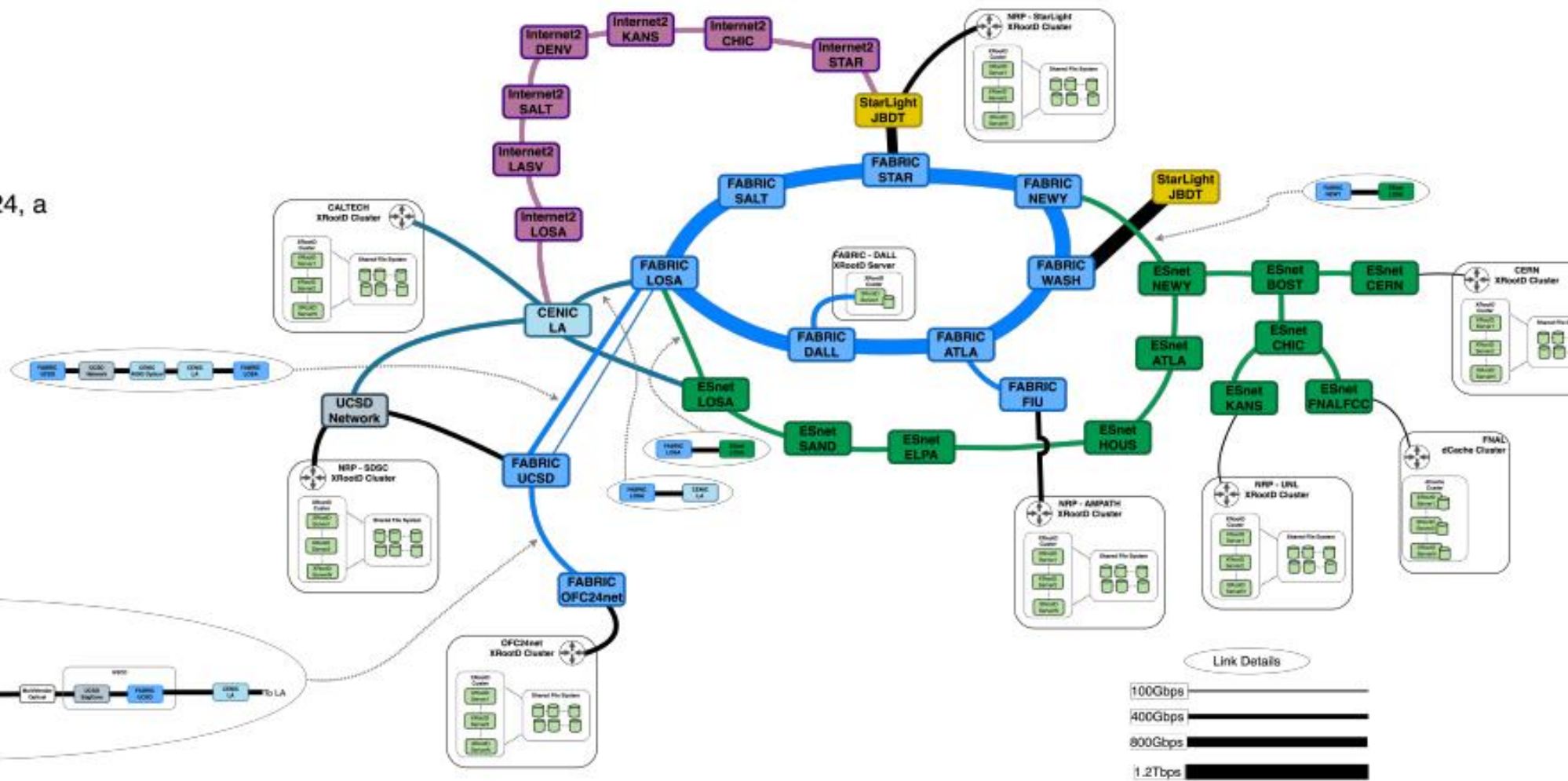
- Abstracted Connectivity
- Dark Fiber
- 1 Gb/s Ethernet
- 400 Gb/s Ethernet
- DWDM
- Quantum Demonstration
- Classical Demonstration

OFC 2024 - OFCnet Architecture





January 25, 2024, a

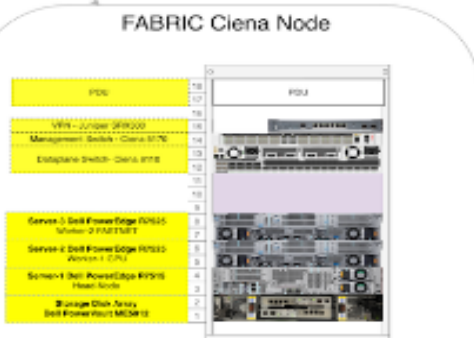
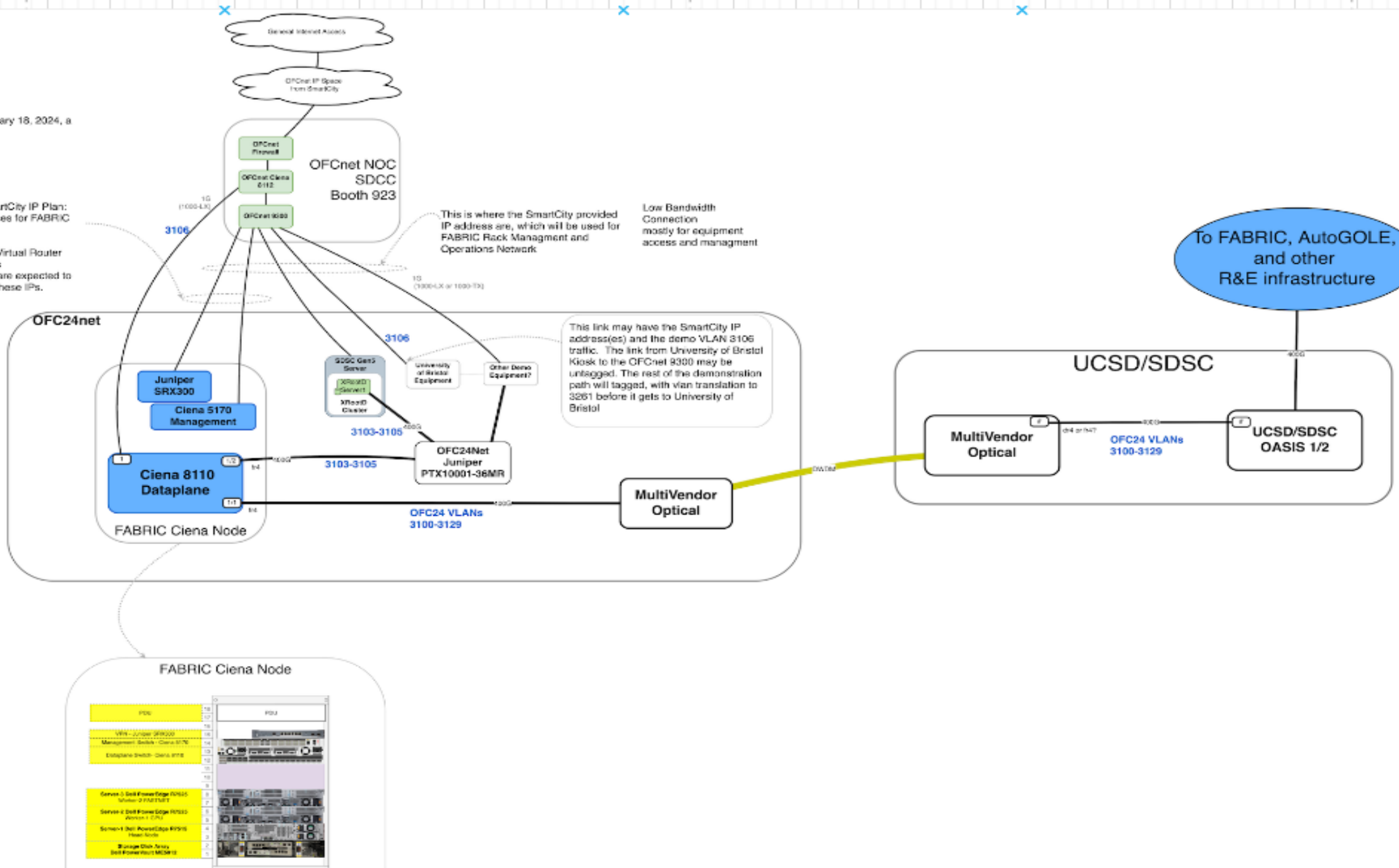


FABRIC utilizes a OFC2net for inter-campus connectivity
FABRIC also utilizes a 1.2Tbps link between UCSD and LA

February 18, 2024, a

FABRIC SmartCity IP Plan:
-5 IP addresses for FABRIC SRX300
Headnode
Open Stack Virtual Router
2 IPs for VMs
-no firewalls are expected to be inline for these IPs.

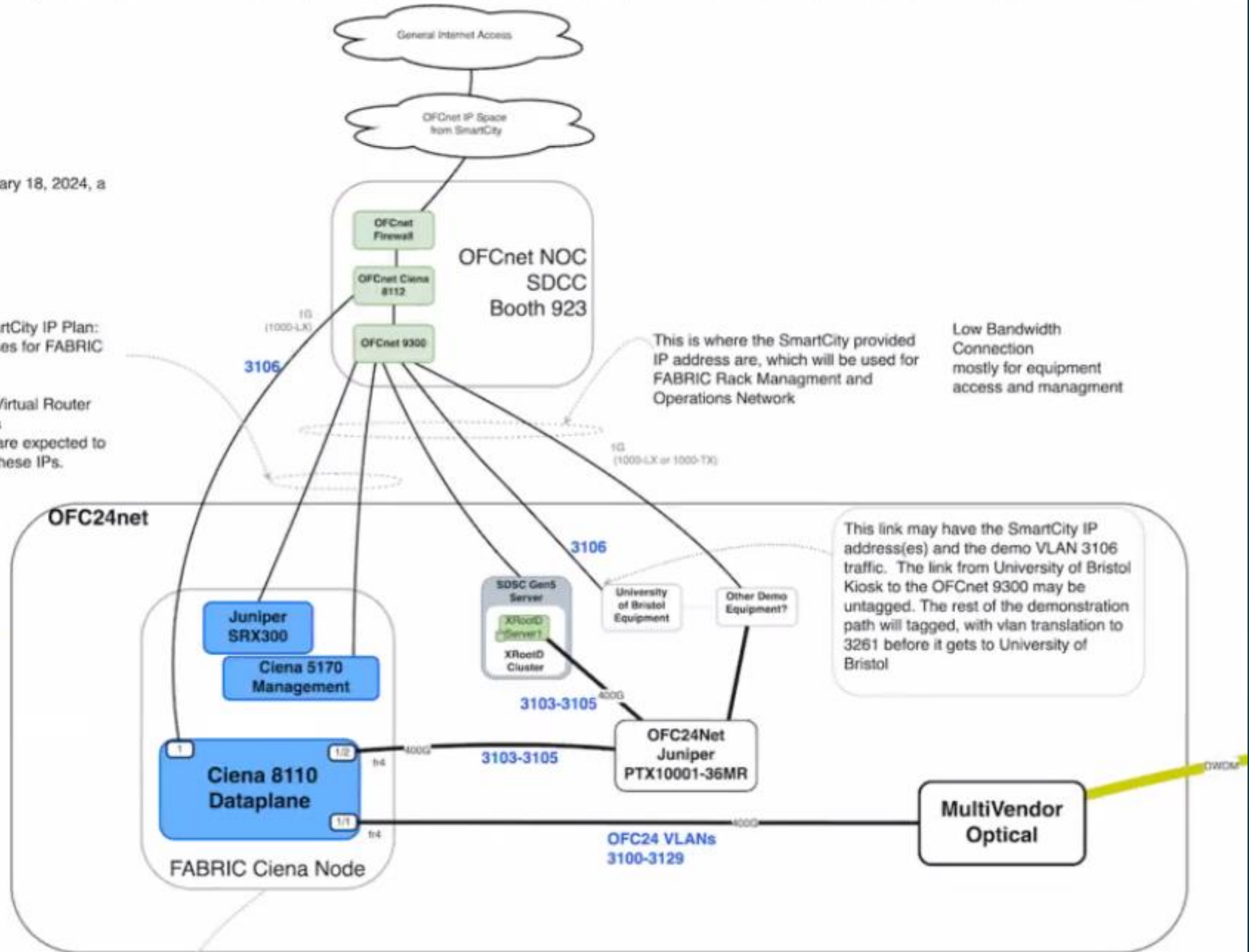
VLAN Demo Assignments:
3100-3102: FABRIC Infrastructure
3103-3105: SDSC Gen5 Server
3106: Bristol



February 18, 2024, a

FABRIC SmartCity IP Plan:
-5 IP addresses for FABRIC SRX300
Headnode
Open Stack Virtual Router
2 IPs for VMs
-no firewalls are expected to be inline for these IPs.

VLAN Demo Assignments:
3100-3102: FABRIC Infrastructure
3103-3105: SDSC Gen5 Server
3106: Bristol



OFCnet NOC SDCC Booth 923

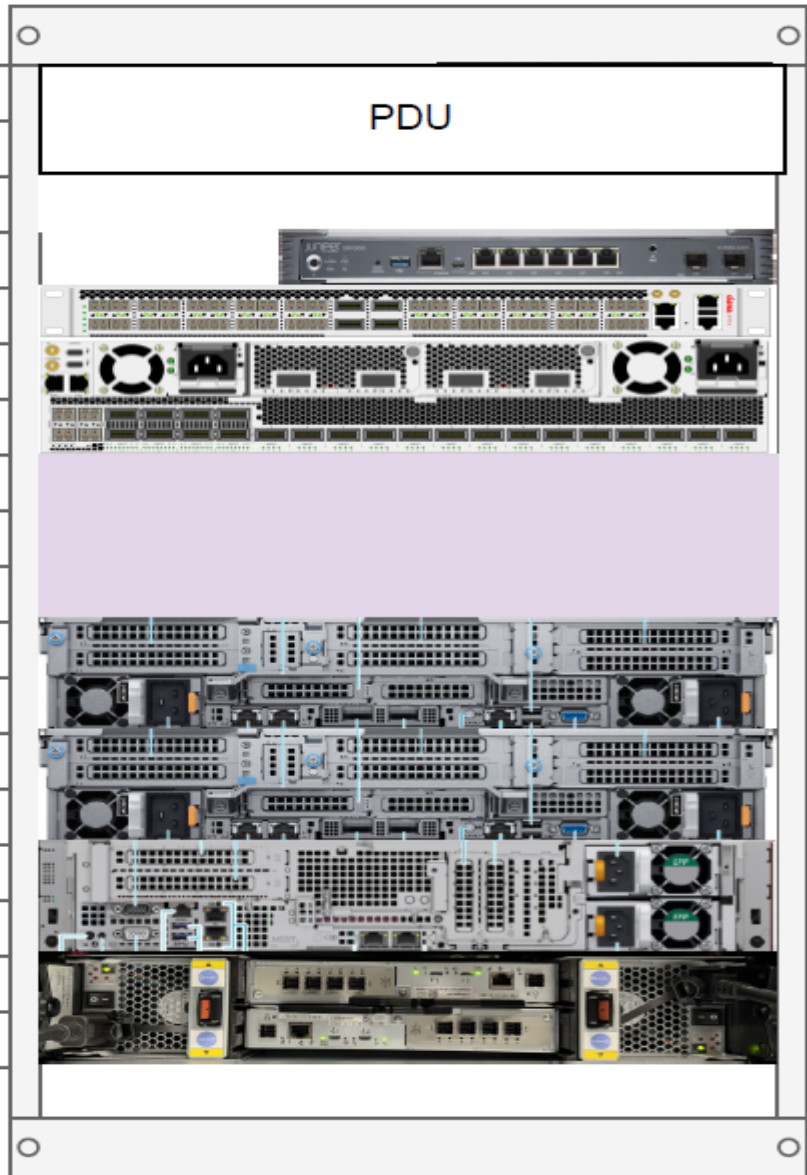
This is where the SmartCity provided IP address are, which will be used for FABRIC Rack Management and Operations Network

Low Bandwidth Connection mostly for equipment access and management

This link may have the SmartCity IP address(es) and the demo VLAN 3106 traffic. The link from University of Bristol Kiosk to the OFCnet 9300 may be untagged. The rest of the demonstration path will be tagged, with vian translation to 3261 before it gets to University of Bristol

MultiVendor Optical

	PDU	18
		17
		16
Serial: CV0923AN0416	VPN - Juniper SRX300	15
Serial:	Management Switch - Ciena 5170	14
Serial:	Dataplane Switch- Ciena 8110	13
		12
		11
		10
		9
Service Tag: J1FB0R3	Server-3 Dell PowerEdge R7525 Worker-2 FASTNET	8
		7
Service Tag: DKRT7Y3	Server-2 Dell PowerEdge R7525 Worker-1 GPU	6
		5
Service Tag: 7MLG7Y3	Server-1 Dell PowerEdge R7515 Head-Node	4
		3
Service Tag: 51MQ7Y3	Storage Disk Array Dell PowerVault ME5012	2
		1

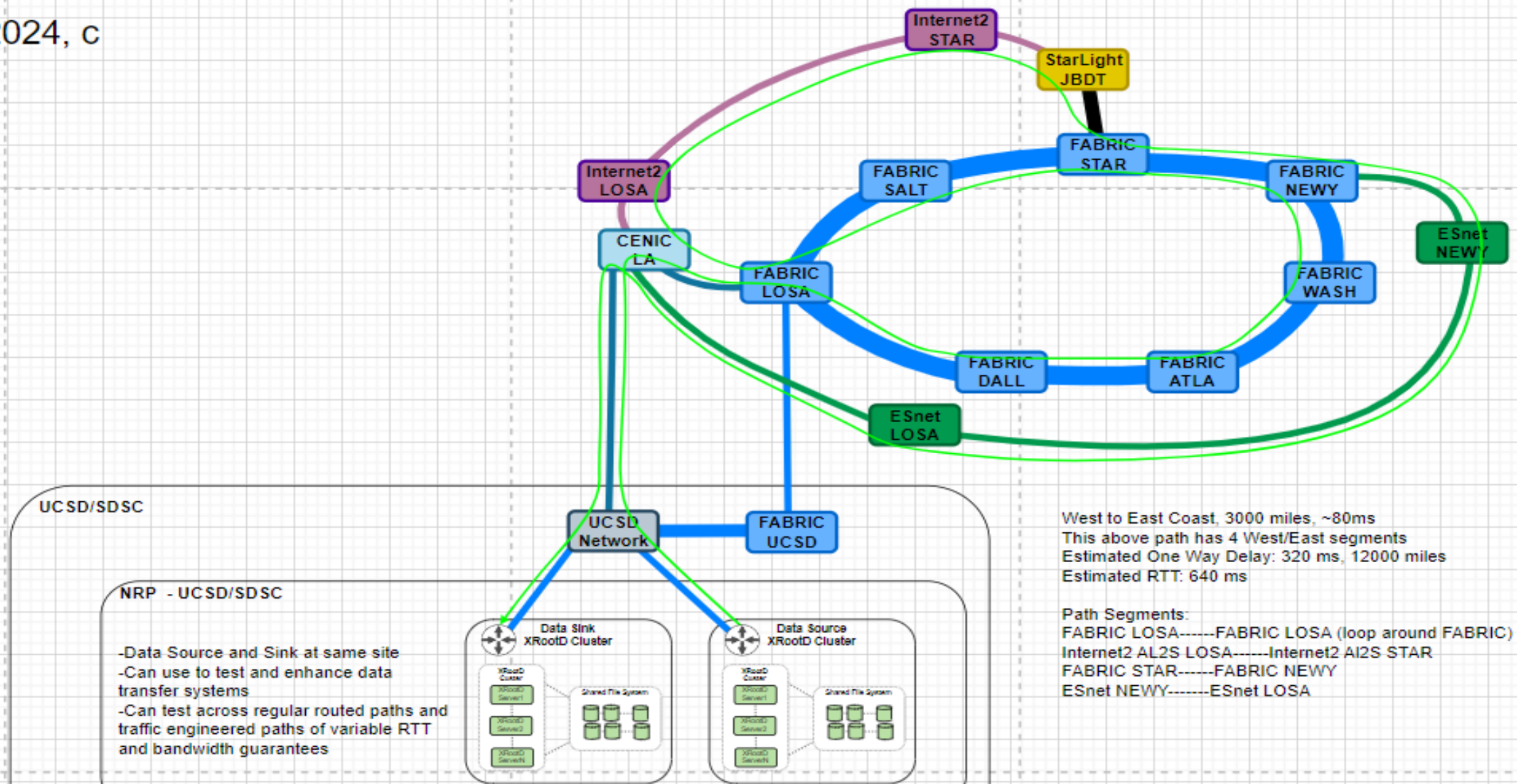


Rack Enclosure
PCC-18U30
(See Note-1)

PDU-1 Serial:
PDU-2 Serial:

y 25, 2024, c

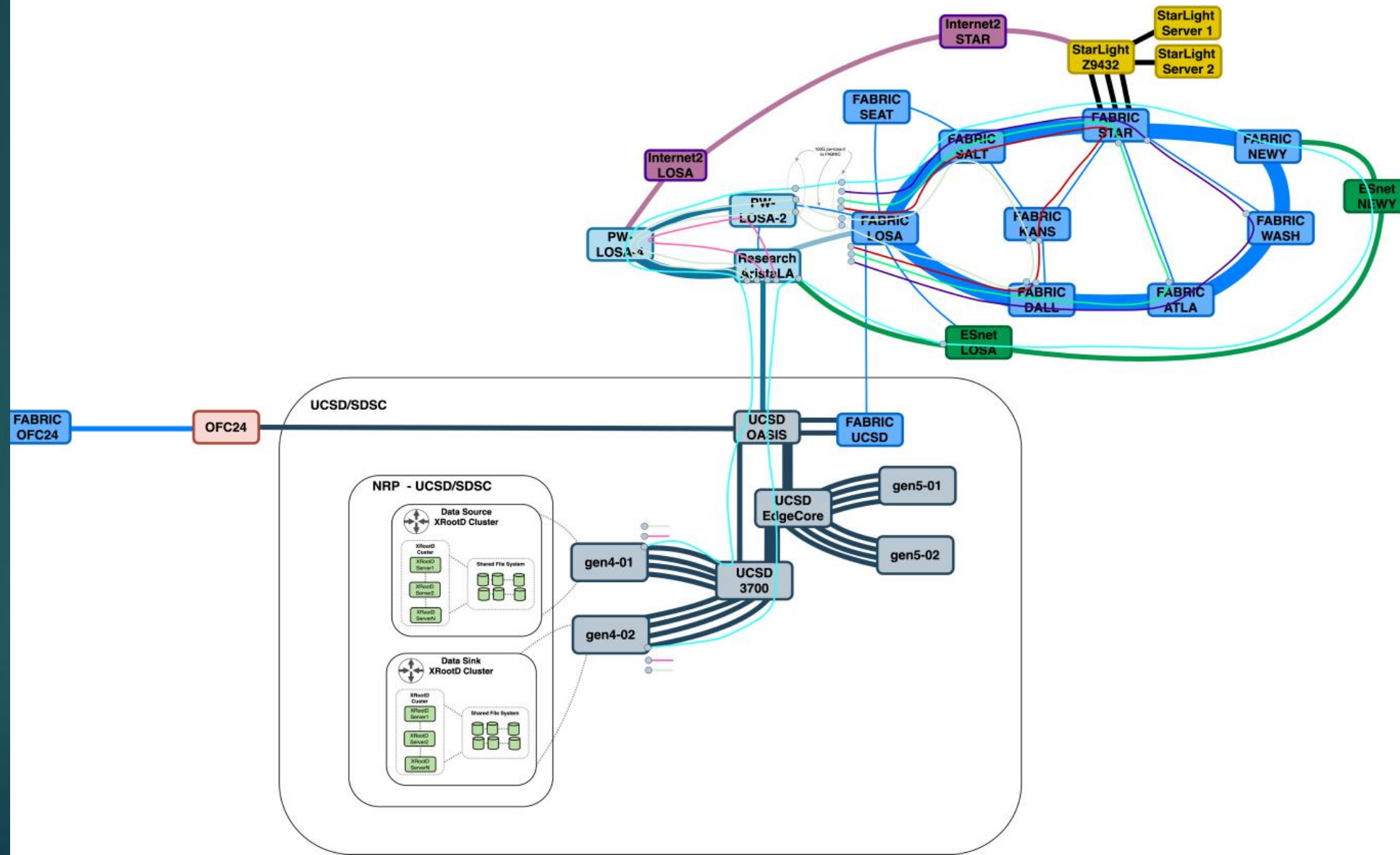
Example Topology for Testing, Research, Development Data Management and Movement Systems



West to East Coast, 3000 miles, ~80ms
This above path has 4 West/East segments
Estimated One Way Delay: 320 ms, 12000 miles
Estimated RTT: 640 ms

Path Segments:
FABRIC LOSA-----FABRIC LOSA (loop around FABRIC)
Internet2 AL2S LOSA-----Internet2 AI2S STAR
FABRIC STAR-----FABRIC NEWY
ESnet NEWY-----ESnet LOSA

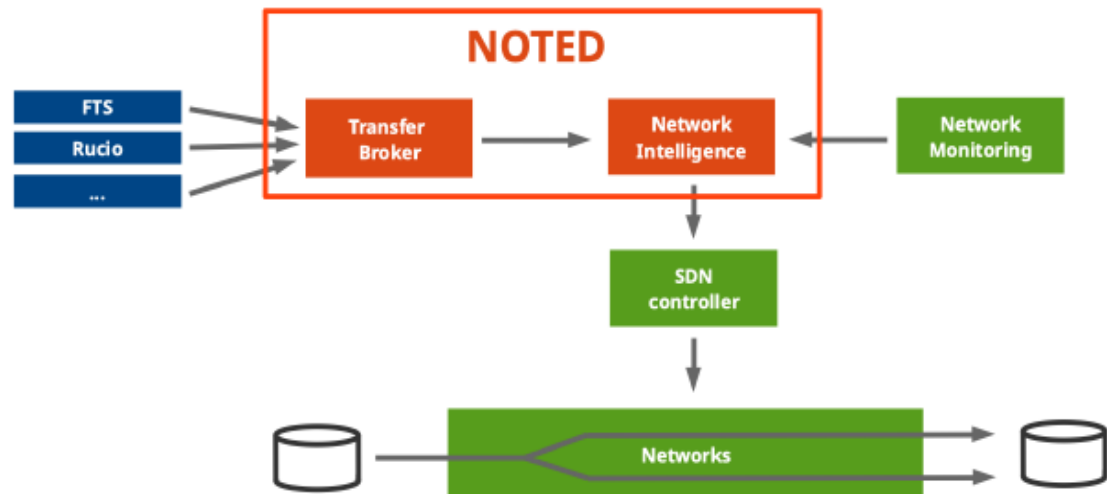
136 ms RTT
120 ms RTT
114 ms RTT
98 ms RTT
80 ms RTT



Example SC24 SCinet Network Research Exhibitions

- ▶ **Global Research Platform (GRP)**
- ▶ **SDX 1.2 Tbps WAN Services**
- ▶ **SDX E2E 400 Gbps WAN Services**
- ▶ **400 Gbps DTNs & Smart NICs**
- ▶ **Network Optimized Transport for Experimental Data (NOTED) – With AI/ML Driven WAN Network Orchestration**
- ▶ **SDX International Testbed Integration**
- ▶ **StarLight SDX for Petascale Science**
- ▶ **DTN-as-a-Service For Data Intensive Science**
- ▶ **P4 Integration With Kubernetes**
- ▶ **PetaTrans Services Based on NVMe-Over-Fabric**
- ▶ **NASA Goddard Space Flight Center HP WAN Transport Services**
- ▶ **Resilient Distributed Processing & Rapid Data Transfer**
- ▶ **PRP/NRP Demonstrations**
- ▶ **Open Science Grid Demonstrations**
- ▶ **N-DISE Named Data Networking for Data Intensive Science**
- ▶ **Orchestration With Packet Marking (SciTags)**
- ▶ **Data Tsunami**

SKELTON AND ELEMENTS OF NOTED



FTS (File Transfer Service):

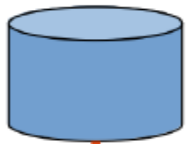
- Inspect and analyse data transfers to estimate if an action can be applied to optimise the network utilization → get on-going and queued transfers.

CRIC (Computing Resource Information Catalog):

- Enrichment to get an overview and knowledge of the network topology → get IPv4/IPv6 addresses, endpoints, rcsite and federation.

FLOWCHART AND DATASET STRUCTURE

- Input parameters: configuration given by the user
 - In noted/config/config.yaml → define a list of {src_rcsite, dst_rcsite}, maximum and minimum throughput threshold, SENSE/AutoGOLE VLANs UUID and user-defined email notification among others.
- Enrich NOTED with the topology of the network:
 - Query CRIC database → get endpoints that could be involved in the data transfers for the given {src_rcsite, dst_rcsite} pairs.
- Analyse on-going and upcoming data transfers:
 - Query FTS recursively → get on-going data transfers for each set of source and destination endpoints.
 - The total utilization of the network is the sum of on-going and upcoming individual data transfers for each source and destination endpoints for the given {src_rcsite, dst_rcsite} pairs.
- Network decision:
 - If NOTED interprets that the link will be congested → provides a dynamic circuit via SENSE/AutoGOLE.
 - If NOTED interprets that the link will not be congested anymore → cancel the dynamic circuit and the traffic is routed back.



Rucio

NOTED at KIT

FTS

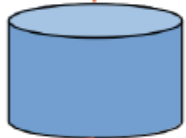
NOTED at CERN

AutoGOLE SENSE

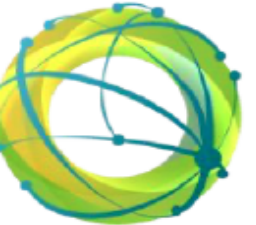
Direct Dynamic Circuit
LHCOPN default path via CERN



Dynamic Circuits
LHCOPN default path



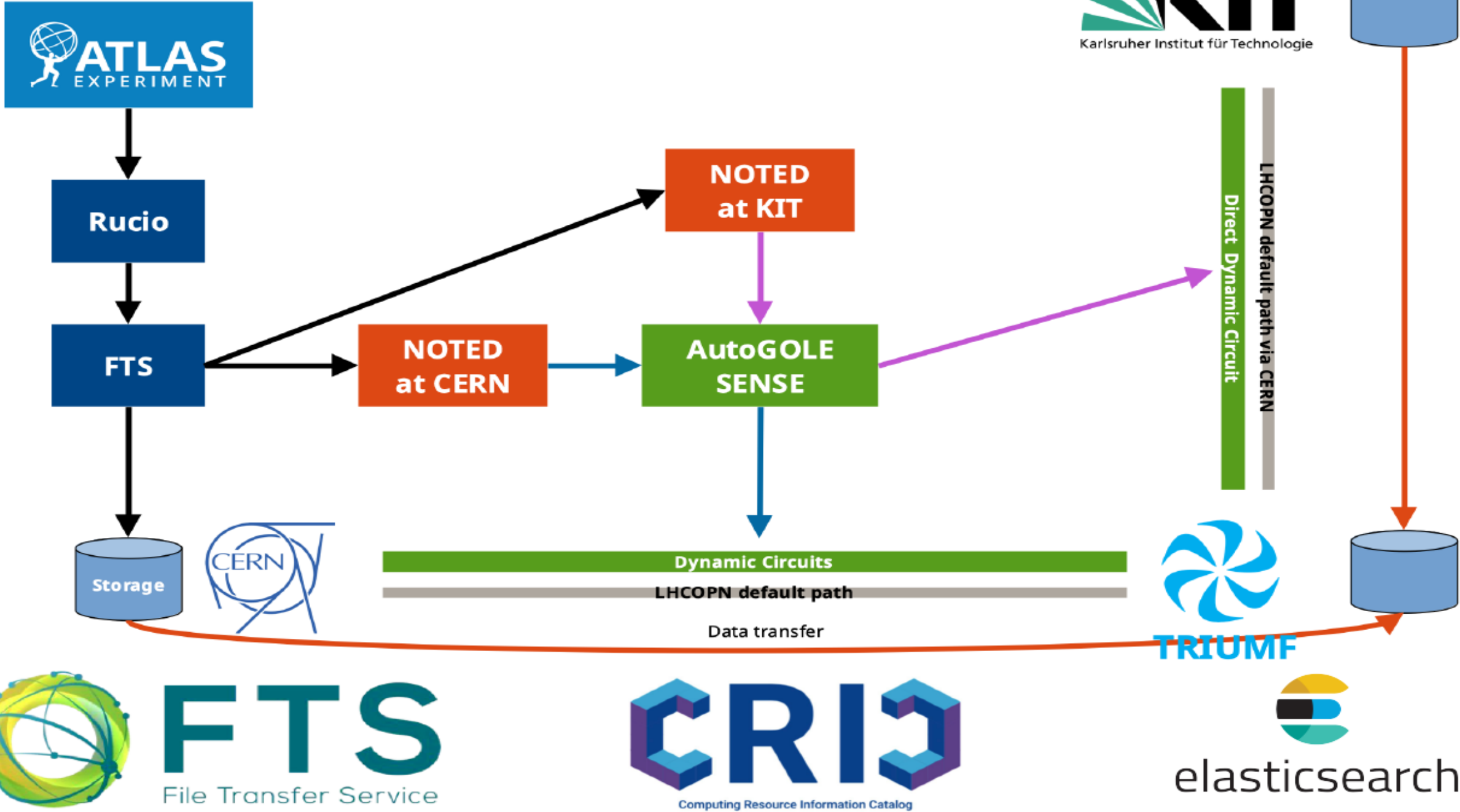
Data transfer



FTS
File Transfer Service

CRIQ
Computing Resource Information Catalog

elasticsearch



Scitags Initiative

Leads= Shawn McKee, Marian Babik

- **Scientific Network Tags** (scitags) is an initiative promoting identification of the science domains and their high-level activities at the network level.



- Enable tracking and correlation of our transfers with Research and Education Network Providers (R&Es) network flow monitoring
- Experiments can better understand how their network flows perform along the path
 - Improve visibility into how network flows perform (per activity) within R&E segments
 - Get insights into how experiment is using the networks, get additional data from R&Es on behaviour of our transfers (traffic, paths, etc.)
- Sites can get visibility into how different network flows perform
 - Network monitoring per flow (with experiment/activity information)
 - E.g. RTT, retransmits, segment size, congestion window, [etc.](#) all per flow

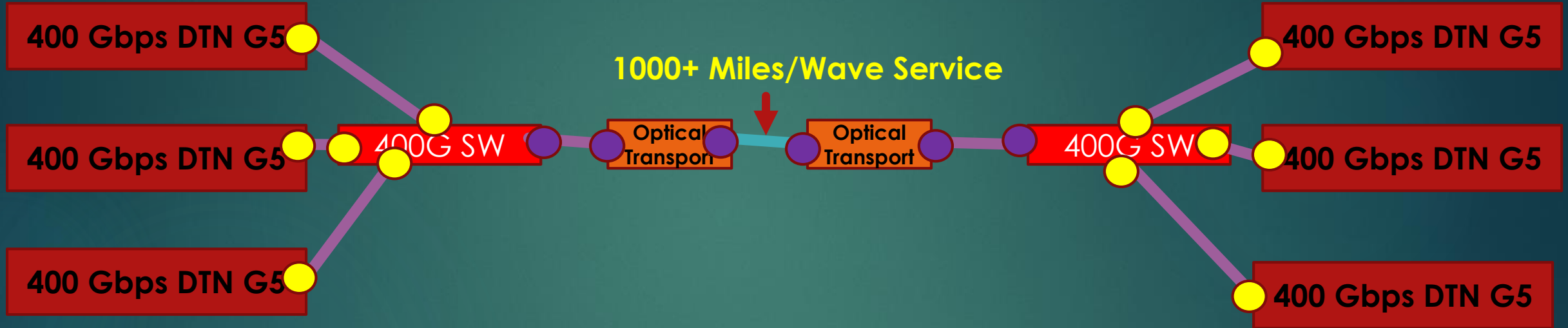
SC24 Packet/Flow Marking NRE

- ▶ **Concept: The Goals of the SC24 Packet and Flow Marking NRE Demonstrations Will Build On the SC23 Demonstrations To Showcase The Capabilities of The Scitags Architecture And Methods For Optimizing Data Intensive Science**
- ▶ **Five Demonstrations Will Be Staged**
 - ▶ IPv6 Packet Marking With eBPF-TC (100 Gbps)
 - ▶ XRootD Packet Marking with Flowd+eBPF-TC
 - ▶ Accounting For Flow Labeled Packets Using a P4 Programmable Switch
 - ▶ Measurements via Esnet High-Touch Processes
 - ▶ Scitags Integration With DTN-as-a-Service.
- ▶ **Participants:**
 - ▶ CERN, University of Victoria, KIT, ESnet, StarLight, CANARIE, Fermi National Accelerator Laboratory, SCInet, Digital Alliance, etc

1.2 Tbps WAN Service Prototype for Data Intensive Science

StarLight International/National
Communications Exchange Facility, Chicago, IL

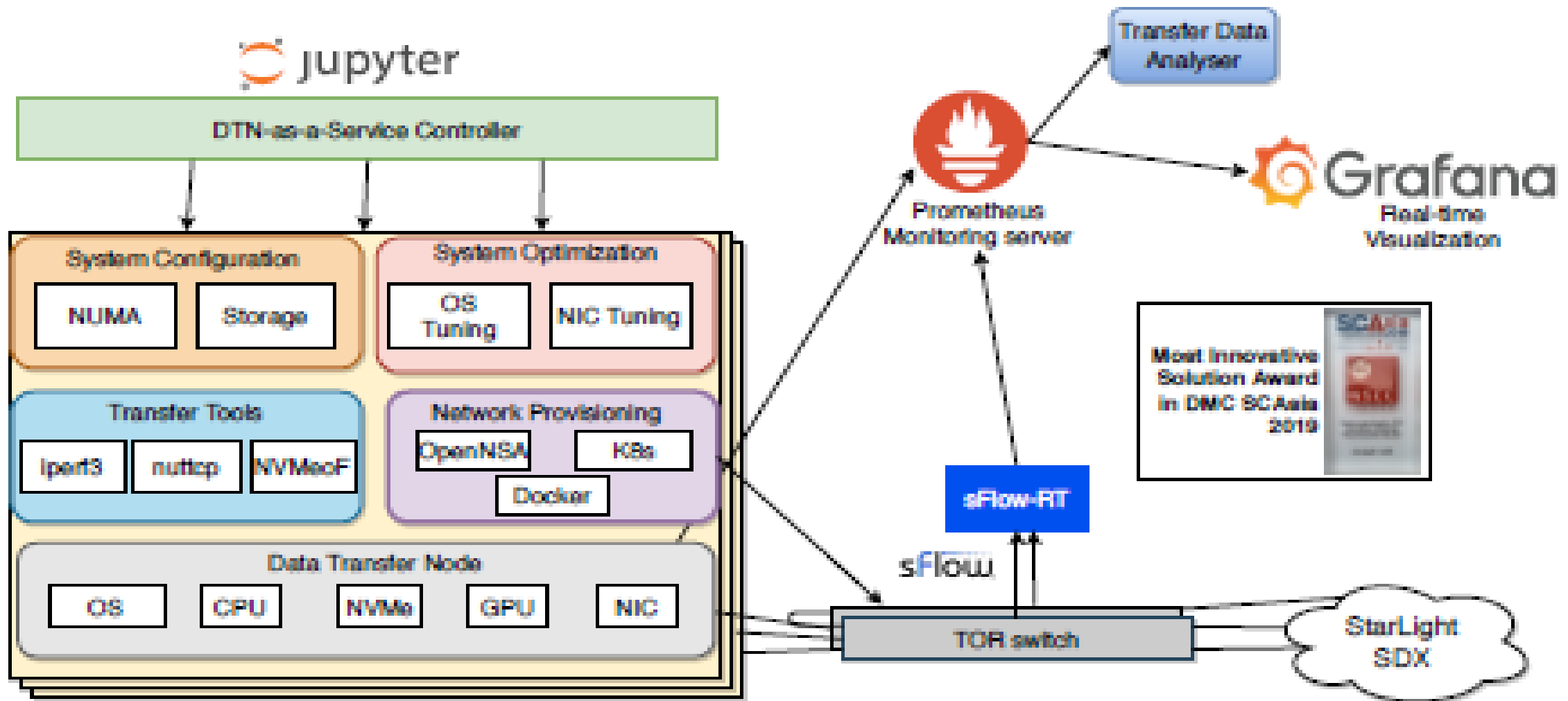
Joint Big Data Testbed McLean, Va



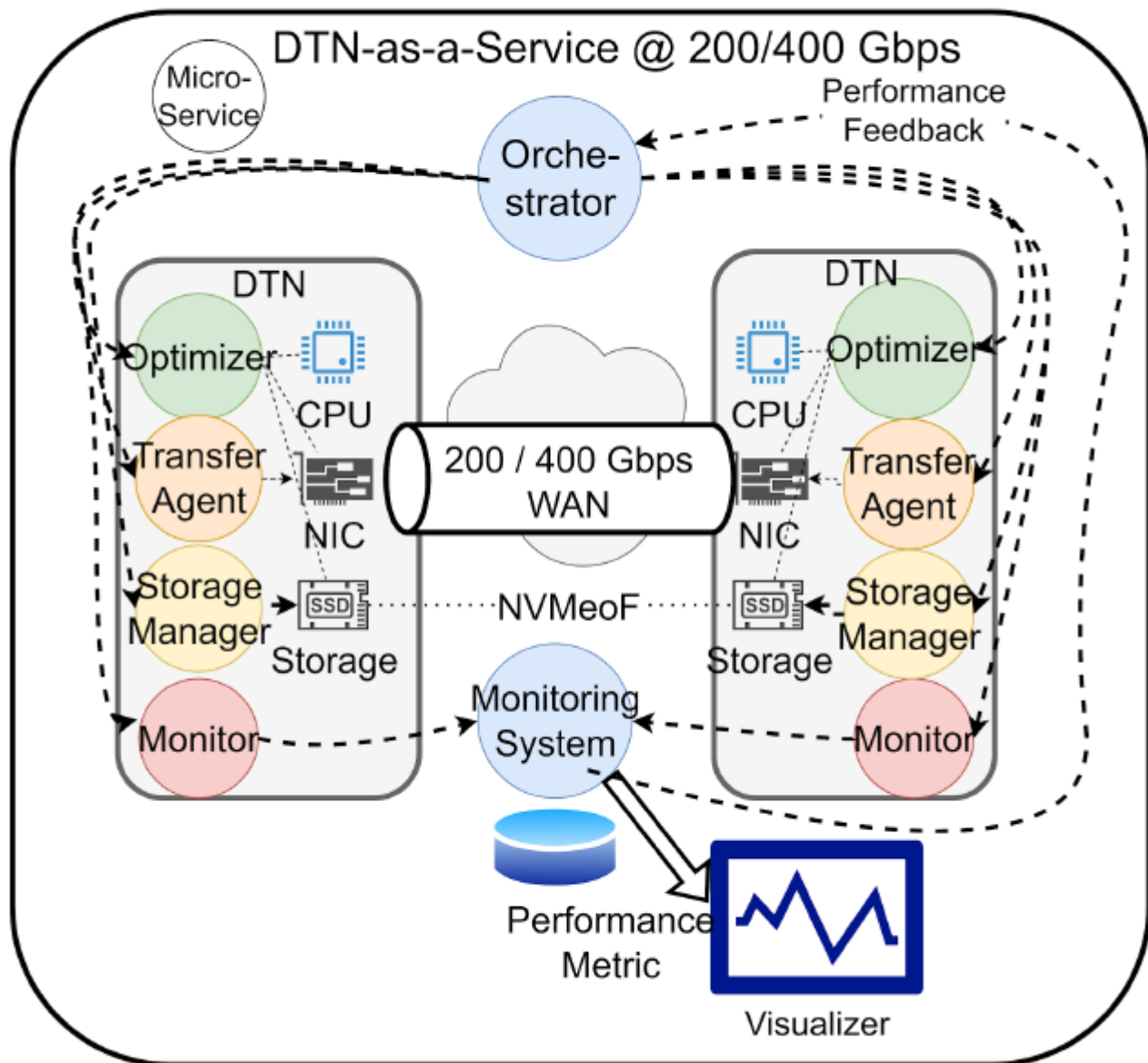
● LR4 Transceiver + Smart NIC

● X Transceiver

STARLIGHT SDX DTN-as-a-Service



200/400 Gbps DTN-as-a-Service in High-Performance Research Platform



- 200/400 Gbps end-to-end high-performance data transfer over WAN
- DTN-as-a-Service with microservice architecture, optimizing and transferring using containers
- NVMeoF with streaming support
- Performance monitoring and visualization using opensource platforms (Prometheus, Grafana, and sFlow)

Peripheral Component Interconnect Express (PCIe)

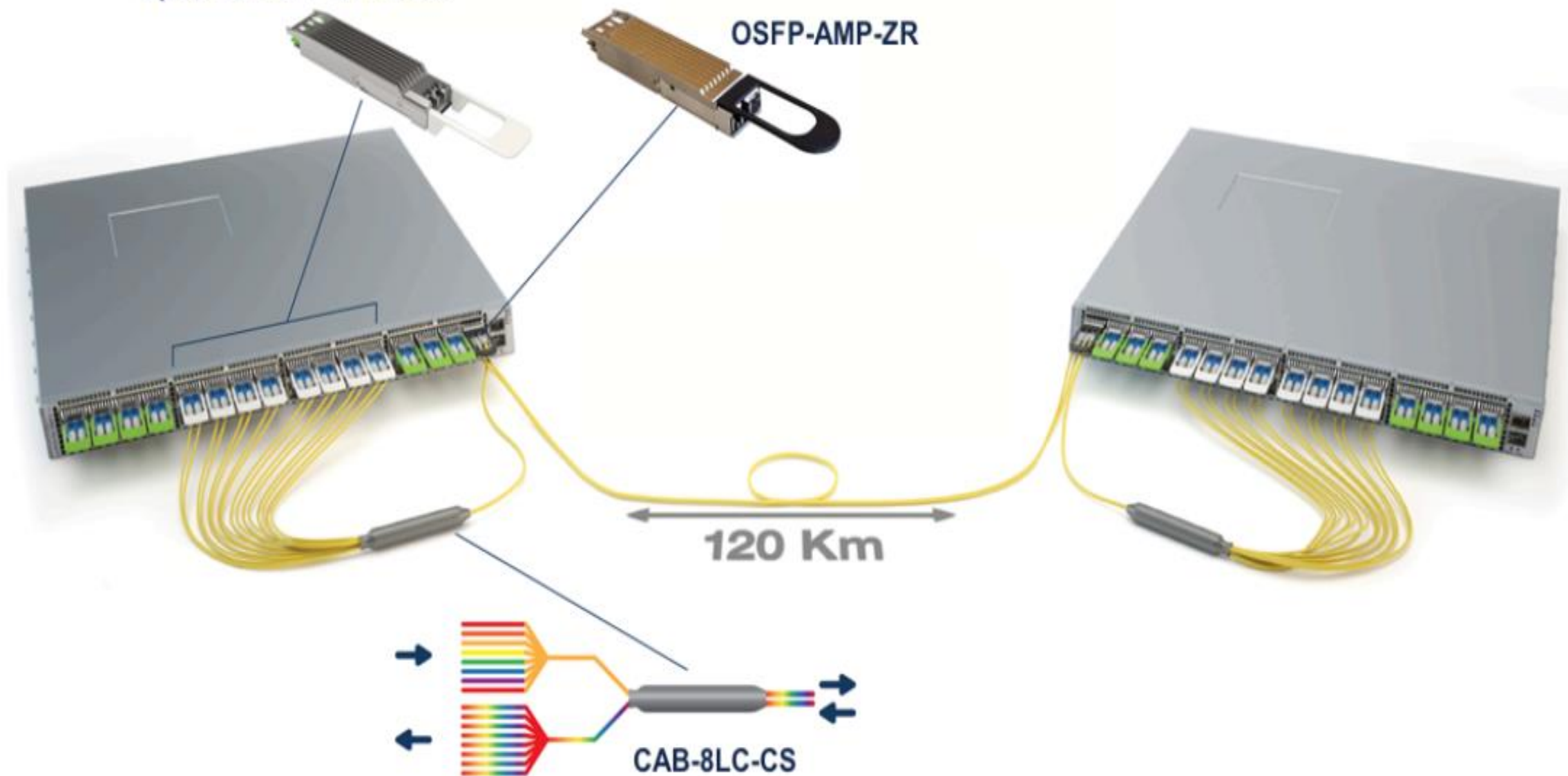
- ▶ PCIe 5.0 server with 32 Lanes = 128GB/s bandwidth 32 GHz
- ▶ 3.94 Gbps/Lane
- ▶ For Optimization Requires G5 Components
- ▶ Foundation For Non-Volatile Memory Express (NVMe) Interfaces, (e.g., SSDs)
- ▶ PCIe 6.0 = 2025-2026

Up to 8x OSFP-400G-ZR

OSFP-AMP-ZR

120 Km

CAB-8LC-CS



WAN Optical Transport

- ▶ Ciena WaveServer 100GbE and 400 GbE (Soon 800 GbE) applications, scaling capabilities to 12.8 Tb/s client and 12.8 Tb/s line capacity in 2RU



Thanks!

▶ Questions?