# tnc23

DIGITAL GENERATIONS
TIRANA, ALBANIA | 5-9 JUNE 2023

# Managed Network Services for Large Data Transfers

SENSE Orchestration

**Tom Lehman (ESnet), Chin Guok (ESnet)**

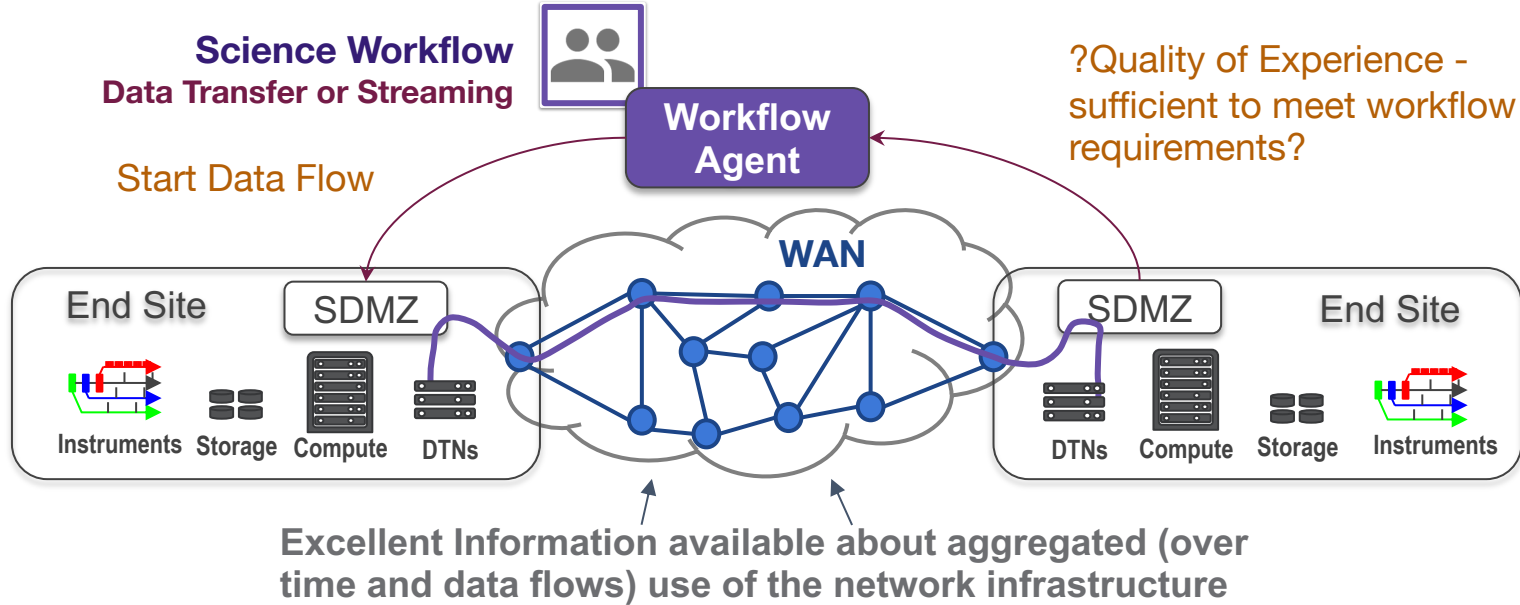TNC23

June 6, 2023

GÉANT

Co-funded by
the European Union

# Presentation Outline

- Multi-Domain, Multi-Resource Service Orchestration
  - objectives, issues, approach
- SENSE Orchestration System, Architecture, Implementation
- SENSE Orchestration services for Rucio/FTS/XRootD Data Movement and Management System
  - with a focus on LHC CMS workflows
- Next Steps

# Enable Science Workflow and Network Interaction with Deterministic "Quality of Experience"

- **No realtime per flow data available for planning or monitoring**
- **No "deterministic" network services available**
- **Start data flow, and hope for the best**

**Science Workflow**
**Data Transfer or Streaming**

**Workflow Agent**

?Quality of Experience - sufficient to meet workflow requirements?

Start Data Flow

**WAN**

End Site

SDMZ

Instruments    Storage    Compute    DTNs

SDMZ

DTNs    Compute    Storage    Instruments

End Site

**Excellent Information available about aggregated (over time and data flows) use of the network infrastructure**

# Elevate Network to First Class Resource
# API driven Automation and Orchestration

**Science Workflow**
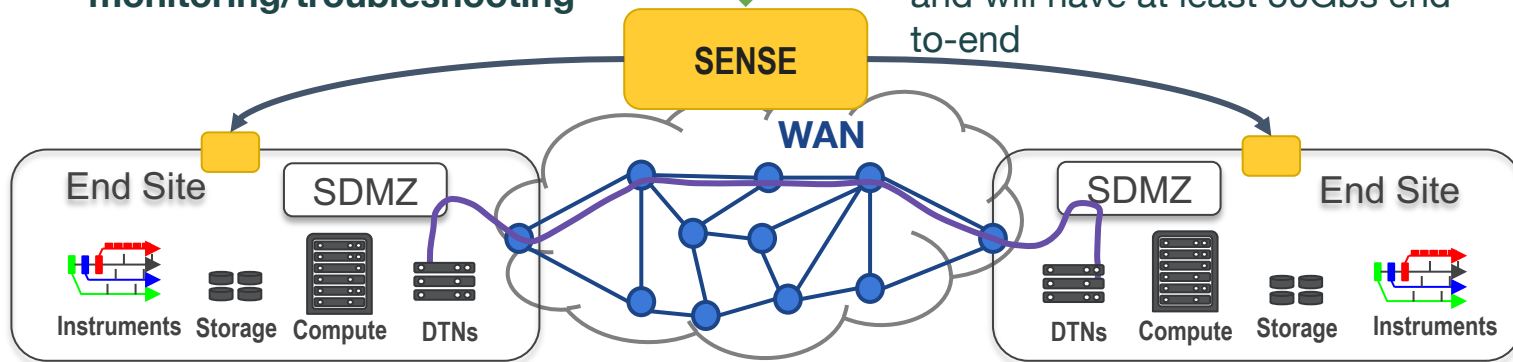**Data Transfer or Streaming**

**Workflow Agent**

SENSE operates between science workflow and the distributed cyberinfrastructure

**Workflow and Network can interact for planning, resource discovery, negotiation, and full life cycle monitoring/troubleshooting**
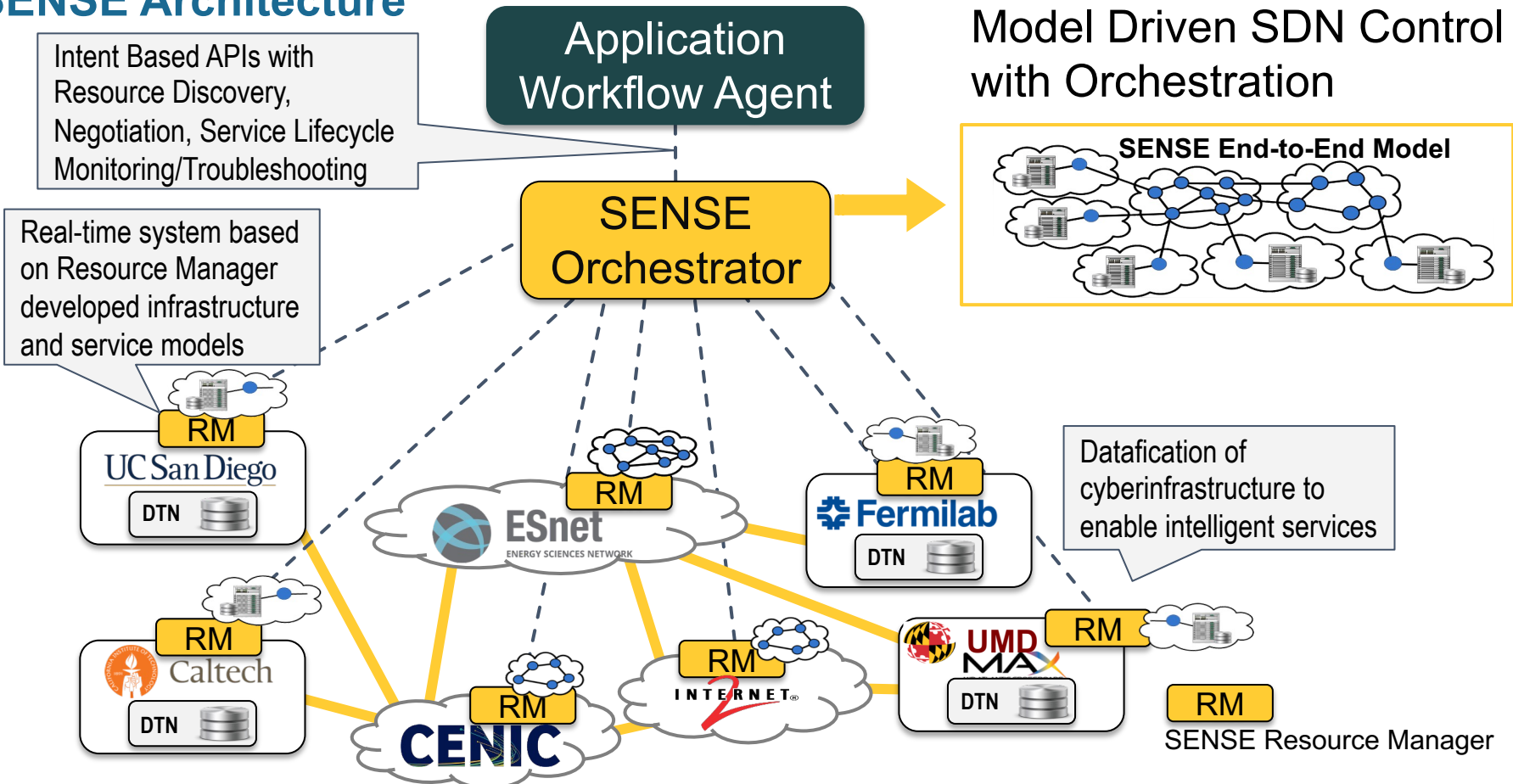
Workflow: Would like to move 1TB anytime in the next 24 hours

Network: You can start in 2 hours, and will have at least 50Gbs end-to-end

**SENSE**

**WAN**

End Site

SDMZ

Instruments  Storage  Compute  DTNs

SDMZ

DTNs  Compute  Storage  Instruments

End Site

- Allows workflows to identify data flows which are higher priority
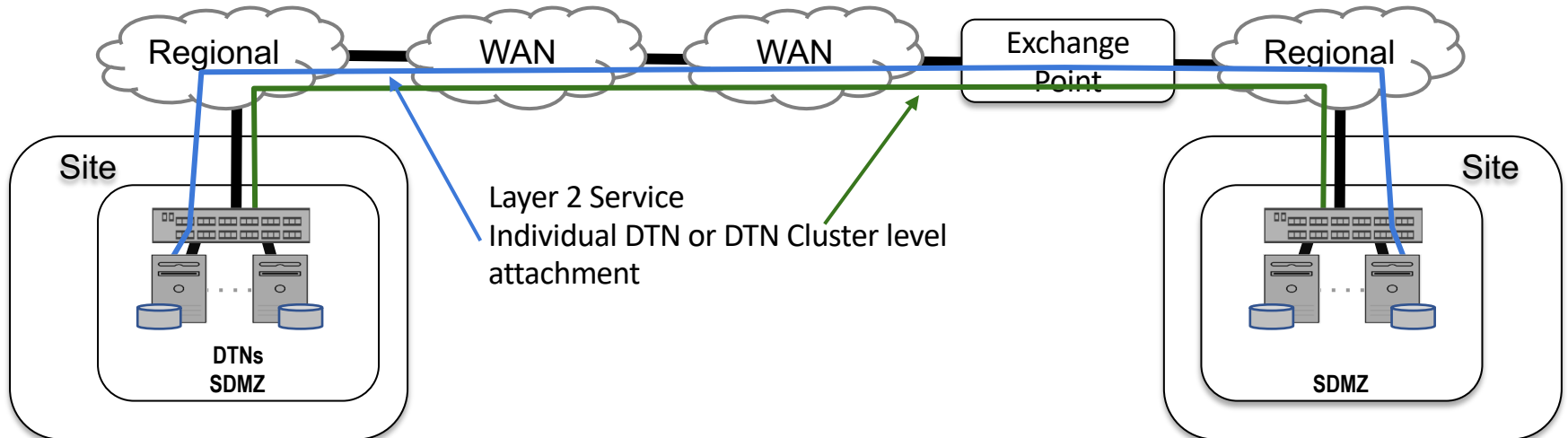- Allows the network to traffic engineer to fully utilize all network paths

# SENSE Architecture

Model Driven SDN Control with Orchestration



Application Workflow Agent

SENSE Orchestrator

Intent Based APIs with Resource Discovery, Negotiation, Service Lifecycle Monitoring/Troubleshooting

Real-time system based on Resource Manager developed infrastructure and service models

Datafication of cyberinfrastructure to enable intelligent services

SENSE End-to-End Model

UC San Diego
DTN

Caltech
DTN

CENIC

ESnet
ENERGY SCIENCES NETWORK

INTERNET 2

Fermilab
DTN

UMD MAX
DTN

RM   SENSE Resource Manager

# SENSE Solution Approach – Application Interactions

- **Intent Based** – Abstract requests and questions in the context of the application objectives.

- **Interactive** – What is possible? What is recommended? Let's negotiate.

- **Real-time** – Resource availability, provisioning options, service status, troubleshooting.

- **End-to-End** – Multi-domain networks, end sites, and the network stack inside the end systems.

- **Full Service Lifecycle Interactions** – Continuous conversation between application and network for the service duration.

# SENSE Services

- **Orchestration** (of other domain owned systems)
- **Multi-Resource** (networks, end systems, instruments, clouds)
- **Multi-Domain** (Sites, Regionals, WANs, Exchange Points)
- **Multi-Service** (L2 Point-to-Point, L2 MultiPoint, L3VPN, QoS, Traffic engineered paths)
- **Intelligent Services** (realtime interaction, full-lifecycle monitoring)



Layer 2 Service
Individual DTN or DTN Cluster level attachment

# SENSE - Model based Resource Descriptions

# SENSE - Model based Resource Descriptions

- Read only and optionally with user editable parameters

- Allows user to run with one time "ticket" or multiple time-use allocations

# SENSE - Northbound API

# Multi-Resource Orchestration

- Networks, End-Systems, Cloud Resources, Instruments

- No need to manage/orchestrate all of the resources end-to-end, just the ones that matter
  - congestion, performance, or policy reasons



**Cloud provider resources and connections**

**Traffic Engineered End-to-End Paths**

# SENSE and Rucio/FTS/XRootD Interoperation

- **Rucio identifies groups of data flows (IPv6 subnets) which are "high priority"**

Scientific Data Management and Movement Suite

Primary system for LHC and others

**Rucio**

**FTS**

**SENSE Orchestrator**

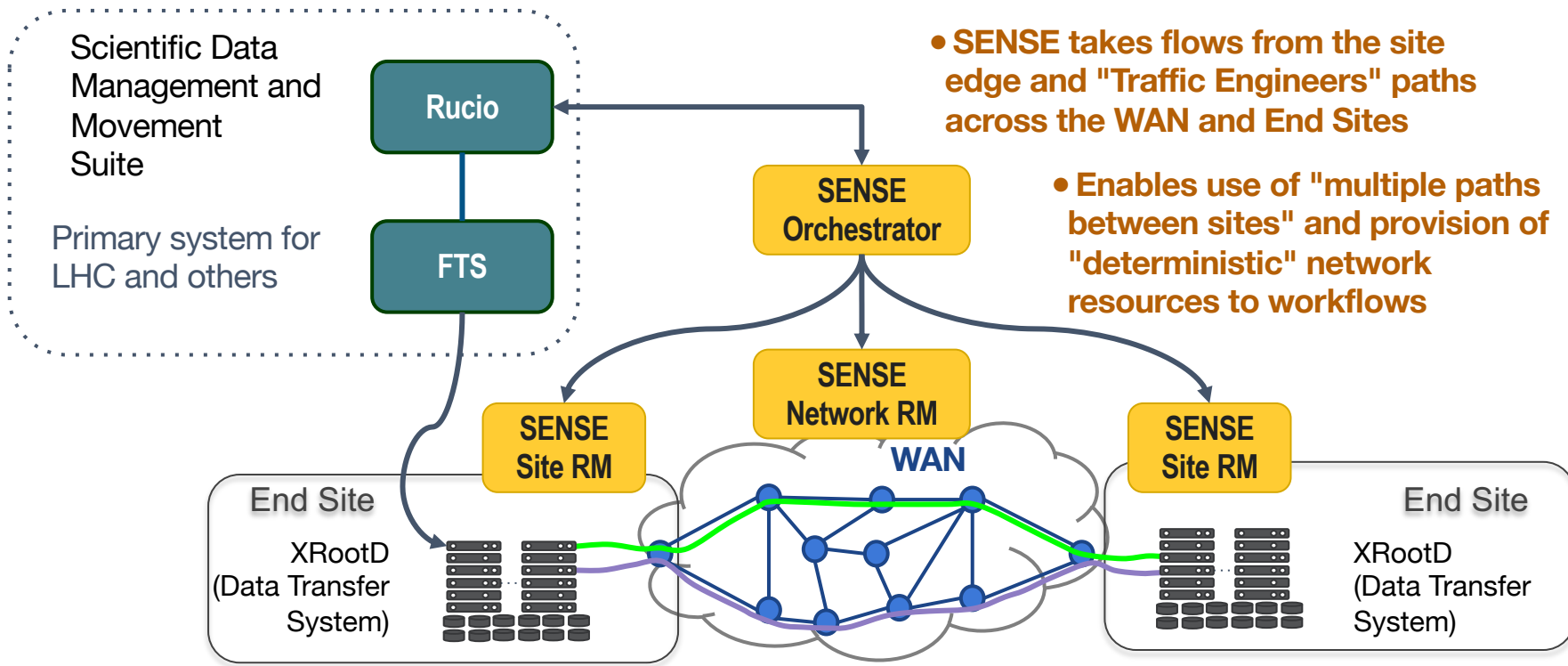**SENSE Network RM**

**SENSE Site RM**

**SENSE Site RM**

**WAN**

End Site

XRootD (Data Transfer System)

End Site

XRootD (Data Transfer System)

- **SENSE takes flows from the site edge and "Traffic Engineers" paths across the WAN and End Sites**

- **Enables use of "multiple paths between sites" and provision of "deterministic" network resources to workflows**

# Objectives

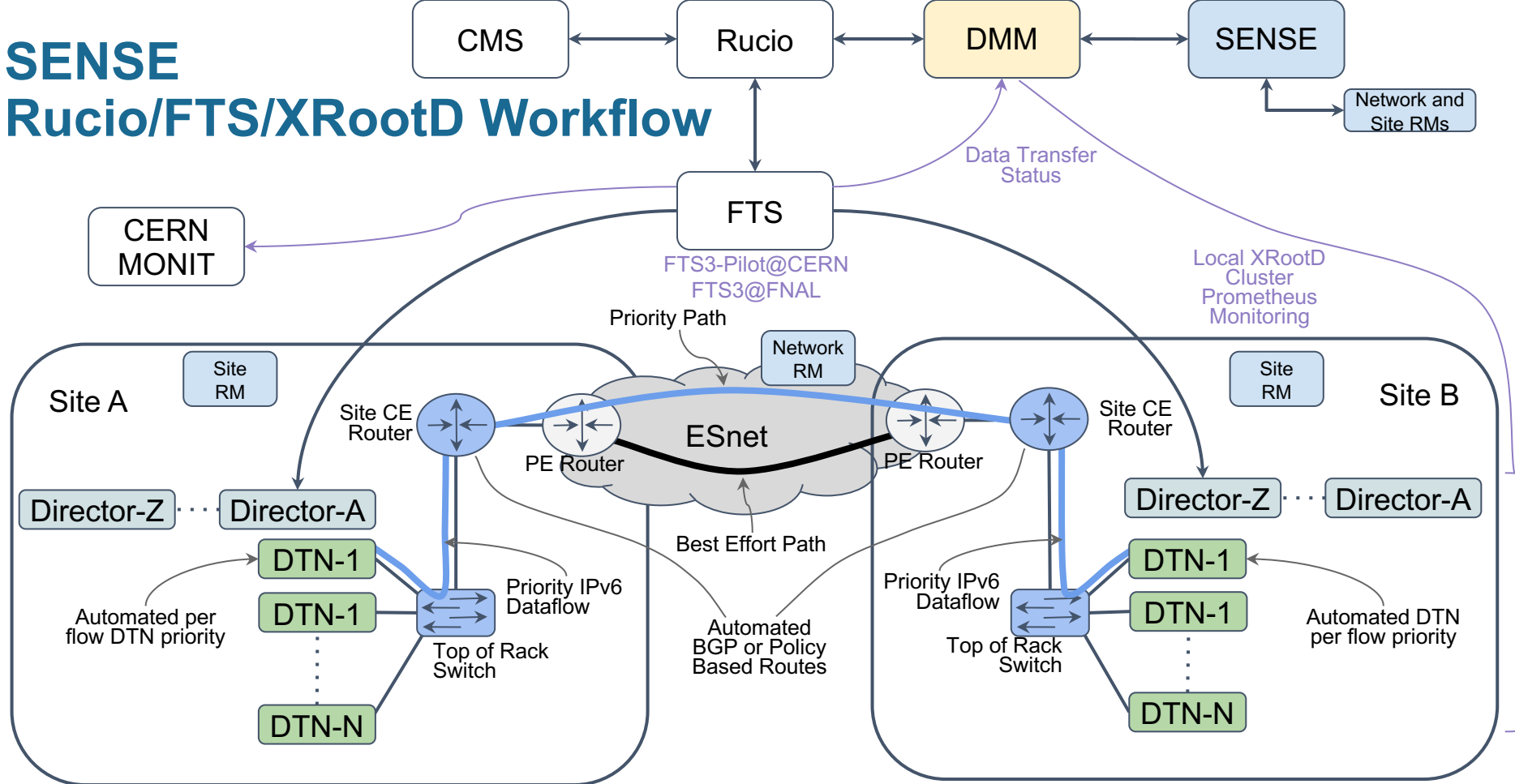*Overall objective is to develop a better way to manage CMS transfers*

*Accountability:  determine where the issues are and develop a process to correct*
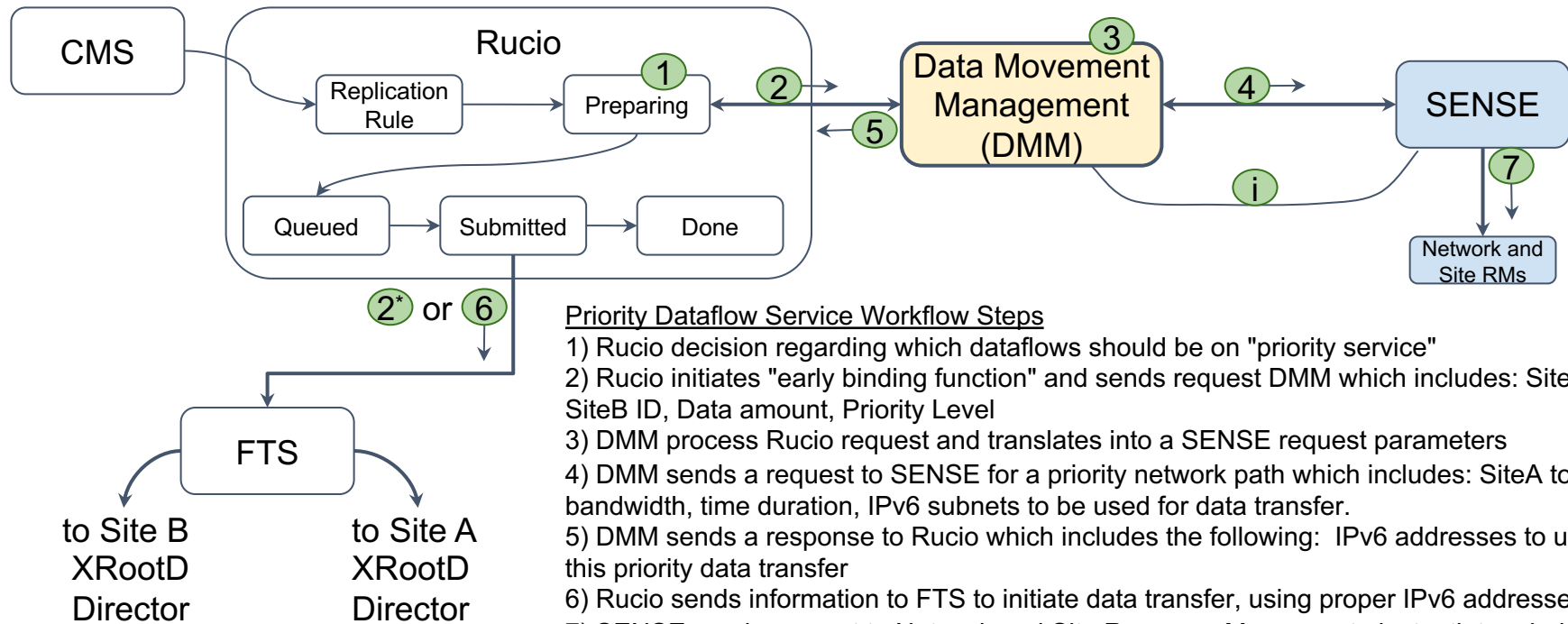
*Focus on the largest flows (not ALL transfers)*

*Plan to use this system as part mini-Data Challenges in 2023 and official Data Challenge in 2024*

# SENSE Rucio/FTS/XRootD Workflow

# Rucio, DMM, SENSE Workflow



**Priority Dataflow Service Workflow Steps**

1) Rucio decision regarding which dataflows should be on "priority service"

2) Rucio initiates "early binding function" and sends request DMM which includes: SiteA ID, SiteB ID, Data amount, Priority Level

3) DMM process Rucio request and translates into a SENSE request parameters

4) DMM sends a request to SENSE for a priority network path which includes: SiteA to SiteB, bandwidth, time duration, IPv6 subnets to be used for data transfer.

5) DMM sends a response to Rucio which includes the following: IPv6 addresses to use for this priority data transfer

6) Rucio sends information to FTS to initiate data transfer, using proper IPv6 addresses

7) SENSE sends request to Network and Site Resource Managers to instantiate priority network service

i) DMM to SENSE "discovery services" (one time at DMM startup)
This is the mechanism for DMM to discover information about sites which includes: sites available for service, IPv6 subnets available, site network connection speed

*Rucio to FTS and DMM interactions can be asynchronous

# DMM  - Data Movement Manager

- React to and process Rucio's "priority" data flow request

- Translate that into actionable information
  - Network provisioning (via SENSE)
  - Data Transfer initiation (identify the proper IPv6 subnet for Rucio-FTS-XRootD to use for a data flow)

- Longer term Focus:  Designing effective policies for how "priority" should be established, who decides, what is the proper mix between priority services and best effort
  - Eventually DMM functions may be distributed between Rucio, SENSE, and/or other parts of the Domain

# Rucio, DMM, SENSE Workflow

- A "priority" data flow is a flexible concept, and could be:
  - all data between Site A and Site B for a specific time period
  - all data between Site A and Site B on a specific IPv6 subnet
  - almost anything based on Site and IPv6/subnet parameters
- End-to-End Data Transfer monitoring
  - Performance evaluation (was the performance as expected?)
  - If not, analysis of why?  (network?, congestion?  where? end-system config/tuning? data movement protocols?  other?)
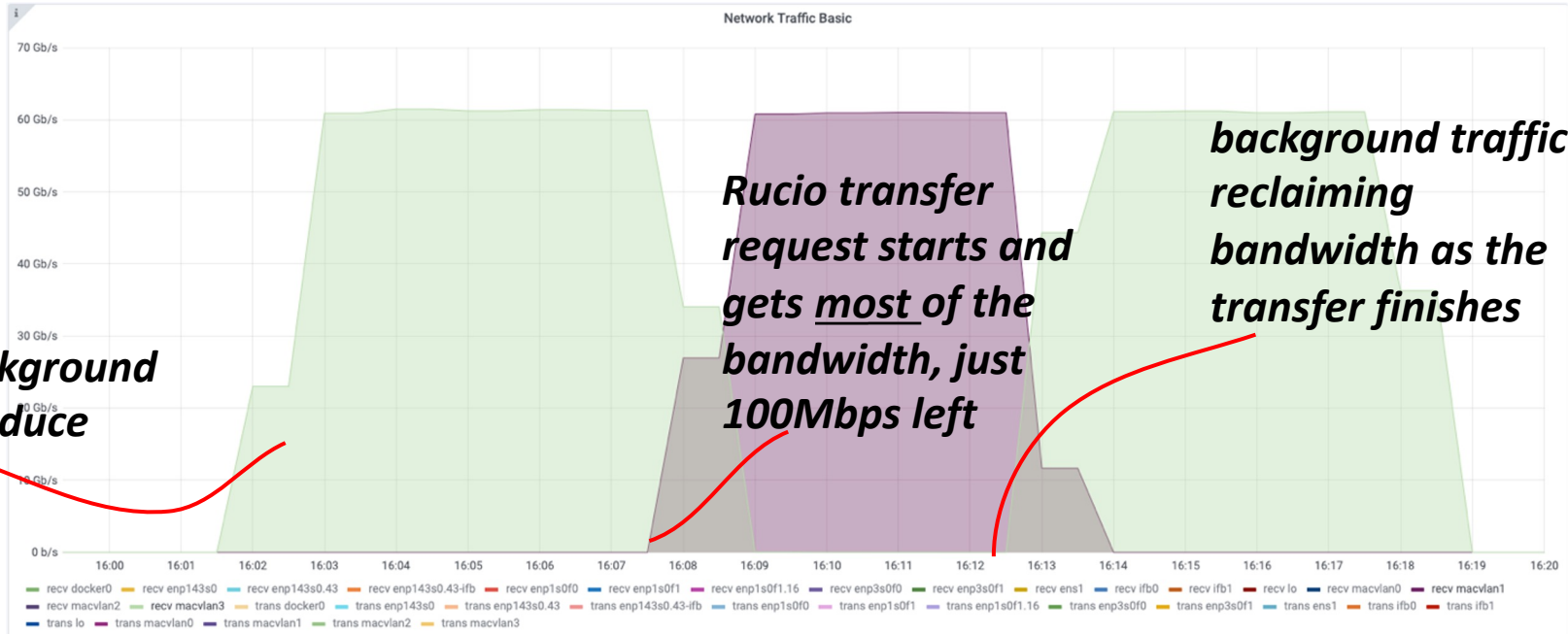
# End-to-End Performance Monitoring

- From local XRootD cluster Prometheus
  - Allocated vs achieved bandwidth
  - Total data transferred vs total transfer size
  - DMM summarizes when a transfer finishes
- FTS records in monIT
  - Data transfer performance from FTS/XRootD perspective
- Correlate data transfer layer throughput with network utilization
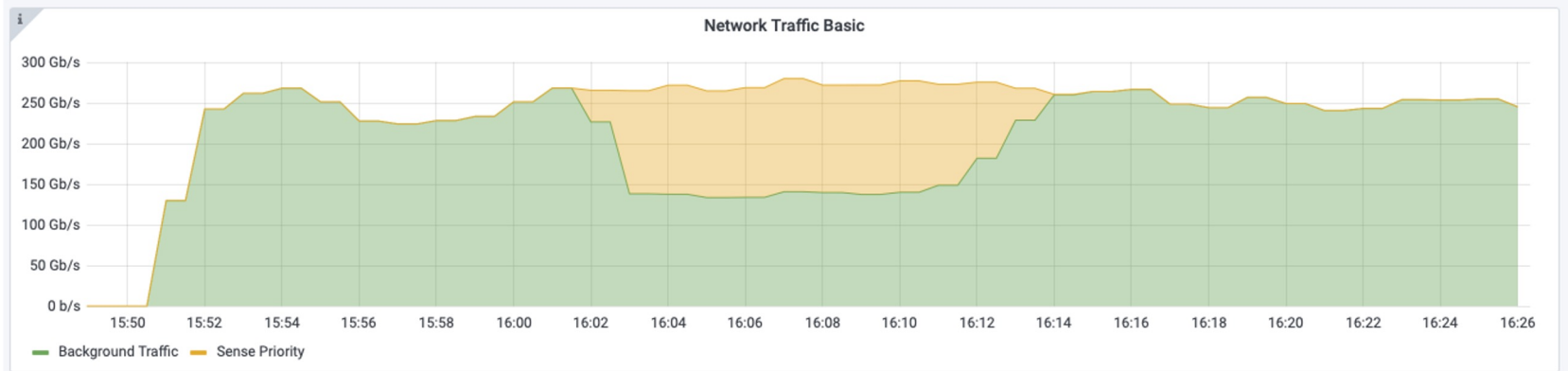  - Still working on the details of data collection, storage, and correlation/analysis

# Proof of Concept Testing

Currently working toward ~400 Gbps site-to-site. Only a few hosts needed for these rates.

**60 Gbps**

*Rucio transfer request starts and gets <u>most</u> of the bandwidth, just 100Mbps left*

*background traffic reclaiming bandwidth as the transfer finishes*

*Artificial background traffic to produce congestion*

# UCSD to Caltech Testing at higher speeds



- Using FDT (Not FTS/XRootD)
- Green – background traffic, Yellow – Priority path requested via SENSE
- Total Capacity between UCSD-Caltech (300gbps). Background 200G, Priority 100G.
- Host level QoS uses Linux TC, Kubernetes/Multus.   Also evaluating use of BPF and Smart NICs for end-system options.

# 400Gbps Benchmark of XRootD- HTTP third-party-copy Transfers

- We can reach ~400* Gbps and sustain it for hours! (345 Gbps over a network path capable of doing 350 Gbps).  Using 40 streams of 1 GB files for each of the 13 servers with Caltech as sink, i.e. 520 streams coming out of UCSD

- XRootD-HTTP is capable of supporting the high throughputs required for the HL-LHC era

- Systematically running transfers can enable us to parameterize by number of CPU cores, number of streams, etc.  Need at least $\mathcal{O}(10)$ streams per XRootD instance for ideal throughput.

- Use of redirectors does not affect performance. Choice of transfer tool does affect throughput.

- Reference UCSD, Caltech team presentation for more details:
  - CHEP23, https://indico.jlab.org/event/459/contributions/11303/

```
sense@sn3700:~$ show interfaces counters -i Ethernet4,Ethernet16,Ethernet20,Ethernet124
     IFACE    STATE         RX_OK        RX_BPS    RX_UTIL   RX_ERR    RX_DRP     RX_OVR            TX_OK          TX_BPS     TX_UTIL    TX_ERR
TX_DRP     TX_OVR
----------   -------   -------------   ----------  ---------  --------  --------   --------    -------------   -------------   ---------   --------
----------  --------
  Ethernet4      U    8,982,229,719   36.80 MB/s     0.29%        0         0          0   60,127,506,940    9958.94 MB/s     79.67%         0
40,750,689       0                                                                                                    +
  Ethernet16     U    9,002,491,671   38.91 MB/s     0.31%        0         0          0   58,470,633,925   11048.23 MB/s     88.39%         0
24,316,754       0                                                                                                    +
  Ethernet20     U    8,847,434,599   31.74 MB/s     0.25%        0         0          0   60,157,515,021    9912.32 MB/s     79.30%         0
29,518,679       0                                                                                                    +
Ethernet124      U    8,845,126,430   36.77 MB/s     0.29%        0         0          0   58,698,987,555   11940.03 MB/s     95.52%         0
26,414,746       0
```
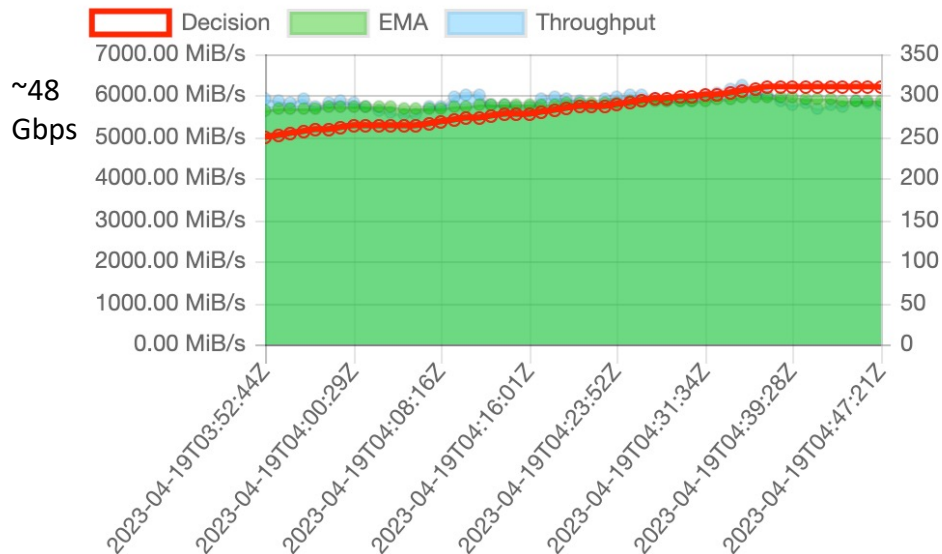
*=345 Gbps*

# FTS Transfers via SENSE Path logged in MONIT (using CERN FTS3@Pilot Instance)
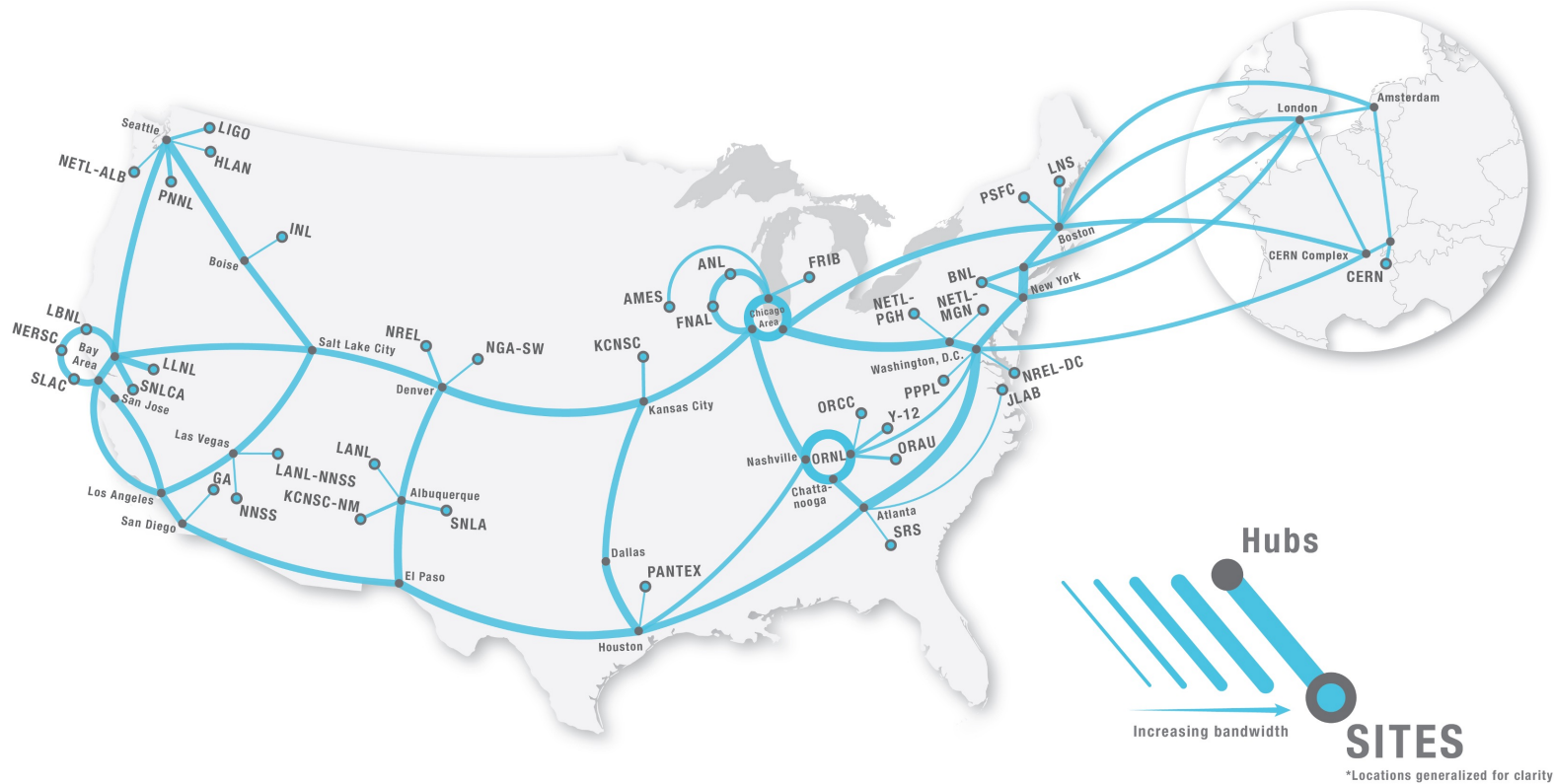
# Higher Speed transfers using FTS/XRootD

- Once FTS Transfers are submitted, FTS Slowly increase number of active transfers (see red line).

- Due to this, XRootD endpoints do not get enough streams to reach >200gbps.

- Working to increase transfer rates

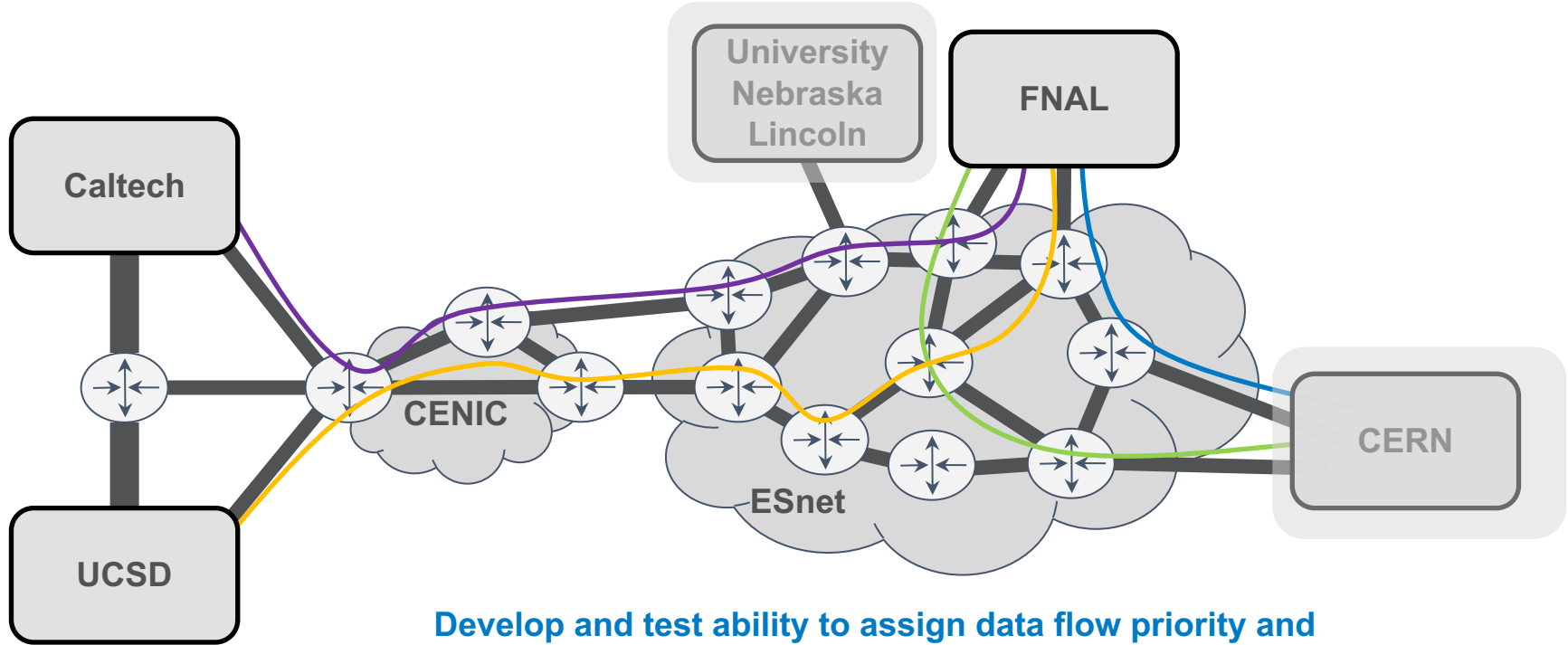- Including a dynamic way to control submission rate (FTS to XRootD)



~48 Gbps

| Source | Destination | VO | Submitted | Active | Staging | S.Active | Archiving | Finished | Failed | Cancel | Rate (last 1h) | Thr. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| + davs://xrootd-sense-ucsd-redirector.sdsc | davs://sense-redir-01.ultralight.or | cms | 1284 | 190 | – | – | – | 14117 | – | – | 100.00 % | 5223.57 MiB/s |
| | | | 1284 | 190 | 0 | 0 | 0 | 14117 | 0 | 0 | 100.00 % | – |

# ESnet Network Topo + US-CMS Sites



Hubs

Increasing bandwidth

SITES
*Locations generalized for clarity

# SENSE Rucio/FTS/XRootD Interoperation System Deployment



**Develop and test ability to assign data flow priority and traffic engineer different end-to-end paths**

May add other sites: CERN, UNL, Vanderbilt, SPRACE

Deployment underway

# Next Steps

- Development Goals:
  - DMM Development and policies. Allow it be adaptable – and define importance of data transfer.
  - Add more sites – US (Fermilab (T1), Nebraska (T2), Vanderbilt (T2)), Brazil - Sprace (T2), CERN (T2). Looking for more European site(s).
  - More NOS (Network Operating Systems) support in SiteRM (Dell OS 10, FreeRTR, Juniper)
  - Quality of Service (Hard QoS, Soft QoS) What to do once underutilized/oversubscribed?
  - Link weights on WAN:
    - Caltech-LasVegas-CERN (130ms, 10gbit max); Caltech-SFO-CERN (163ms, 20gbit max)
  - Policy for fair-share between experiments. Who gets how much and what?
  - Automated End-to-End troubleshooting, monitoring, alarming. (pin-point exact hop failing, alerting)
  - Other experiment use cases and support in SENSE.

- Participate in the WLCG Data Challenge 2024

- Software-Defined Network for End-to-end Networked Science at the Exascale, Elsevier Future Generation Computer Systems, Volume 110, September 2020, Pages 181-201, https://doi.org/10.1016/j.future.2020.04.018

- SENSE Northbound API Program
  - https://app.swaggerhub.com/apis/xi-yang/SENSE-O-Intent-API

- Contacts
  - Xi Yang, xiyang@es.net
  - Tom Lehman, tlehman@es.net
  - SENSE Information, sense-info@es.net

- SENSE Website: sense.es.net

# Acknowledgements

**ESnet** — Chin Guok, Tom Lehman, Inder Monga, Xi Yang

**California Institute of Technology** — Harvey Newman, Justas Balcas, Preeti Bhat

**UCSD / CMS** — Frank Würthwein, Jonathan Guiang, Aashay Arora, Diego Davila, John Graham, Dima Mishin, Thomas Hutton, Igor Sfiligoi

**Fermilab** — Oliver Gutsche, Phil Demar

# Thank you
Any questions?

Tom Lehman <tlehman@es.net>

# Extra Slides

# Key Themes

- Today, science workflows view the network as an opaque infrastructure - inject data and hope for an acceptable Quality of Experience

- We should allow workflow agents to interact with the network - ask questions, see what is possible, get flow specific data and resources

- Science workflow planning should be able to include the networks as a first-class resource (alongside compute, storage, instruments)

- This requires collaborative cross-discipline teams for workflow co-design

- The same mechanisms that allow the above can also be used by individual networks to distribute traffic more efficiently across entire infrastructure

## Objectives

- Provide mechanisms for domain science workflows and middleware (Rucio) to identify "priority" data flows
- Realtime integration of site data flows and wide area traffic engineering
  - in response to "priority" request
  - and/or just allow better overall network (link) utilization via traffic distribution/optimization
- Traffic engineering may include paths with QoS, or to traverse lightly loaded links

# SENSE - Site Layer 3 Flow to WAN Traffic Engineered Path Service

# Objectives

- Make Rucio capable to **schedule transfers on the network.**

- **Improve accountability**.

- **Predetermined transfer speed and quality of service (time to completion).**

- **Fine-grain managed transfers can be also fine-grain monitored** since they travel alone within a well-identified network channel.

- Comparing **Achieved V.S. Allocated bandwidth** will make network & endpoint issues evident.

# Important Link Management

- There are **multiple transatlantic and transpacific links**, operated by multiple organizations

- Goal is to more flexibly control how these are utilized on a per flow, group, or use basis

- Do not want to manage "**every**" flow in the network; but we should be able to manage "**any**" flow in the network

- **An equally important goal** is to understand the load vs capacity and leave room for other traffic

- **Remain compatible** with other network operations

- **Two timescales:** SENSE overlay network of virtual circuits with BW guarantees is relatively stable; IPv6 subnets and Directors provide more dynamic flow mapping to various traffic engineered paths