# Rucio/SENSE in DC24

Frank Würthwein, Jonathan Guiang, Aashay Arora, **Diego Davila**, John Graham, Dima Mishin, Thomas Hutton, Igor Sfiligoi, Harvey Newman, Justas Balcas, Preeti Bhat, Tom Lehman, Xi Yang, Chin Guok, Oliver Gutsche, Asif Shah, Chih-Hao Huang, Dmitry Litvinsev Phil Demar, Marcos Schwarz
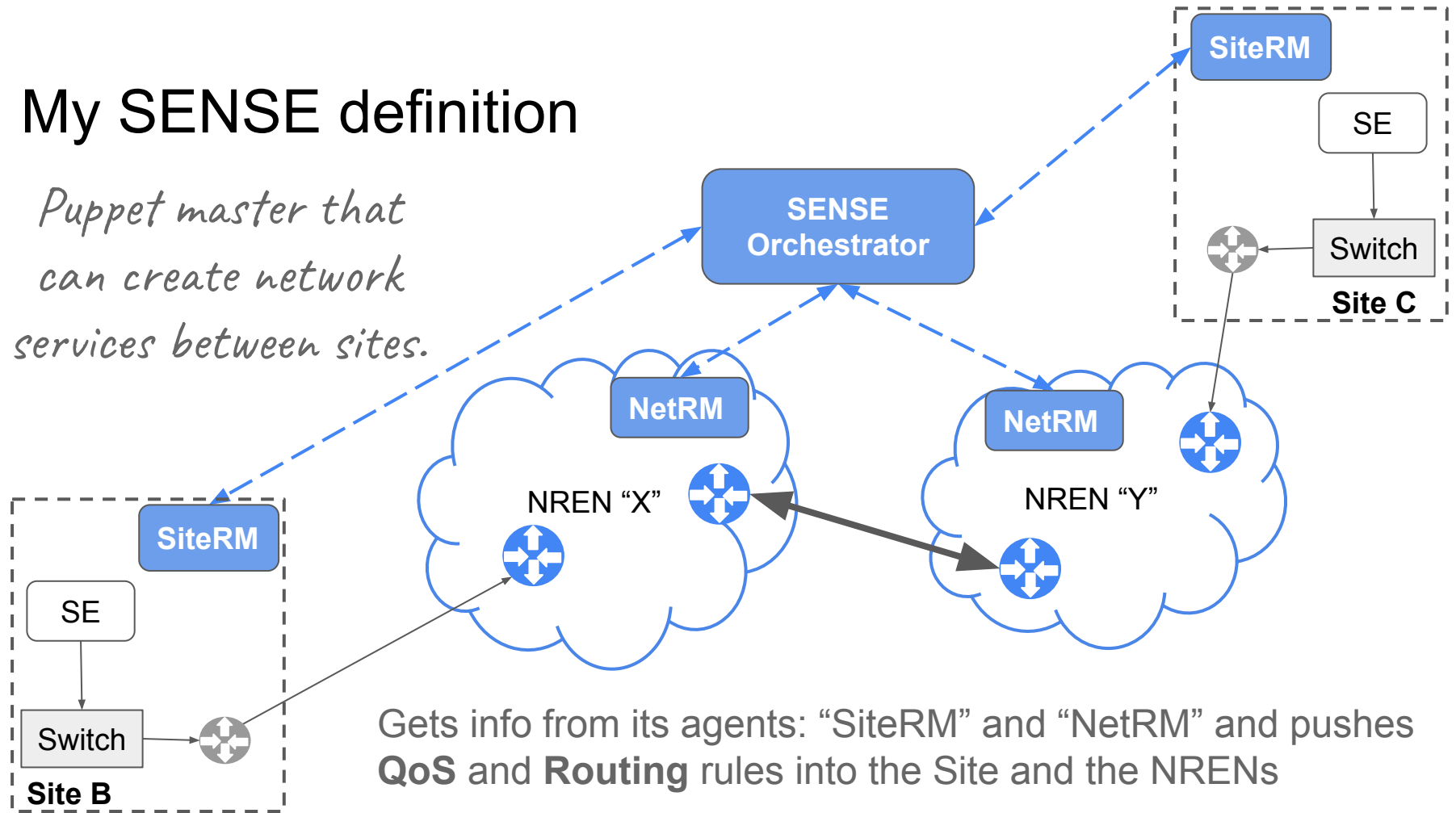
Nov 9th, 2023

# Introduction

The overall goal is to enable Rucio to negotiate and manage network capabilities with SENSE

This creates end-to-end accountability of the network utilization and allows individual stakeholders (e.g. VOs) to manage their internal priorities.

Prioritization between different stakeholders are managed at the SENSE layer

# My SENSE definition

*Puppet master that can create network services between sites.*

**SiteRM**

SE

Switch

**Site C**

**SENSE Orchestrator**

**SiteRM**

SE

Switch

**Site B**

**NetRM**

NREN "X"

**NetRM**

NREN "Y"

Gets info from its agents: "SiteRM" and "NetRM" and pushes **QoS** and **Routing** rules into the Site and the NRENs

# Multi subnet Storage System

- Sense services are created based on subnets
- **Problem:** current Storage Systems live in a single subnet
- **Solution:** expose our Storage Systems over multiple subnets
- We managed to do this in XRootD by adding a bunch of configuration
  - No extra hardware needed

More details here:
https://indico.cern.ch/event/1185600/contributions/5109192/attachments/2545788/4383989/Automated%20Network%20Services%20for%20Exascale%20Data%20Movement%20(1).pdf

# Rucio & SENSE

Rucio: Data Management System used by CMS and ATLAS, it knows the data workflows, how big, where they have to go, how important are.

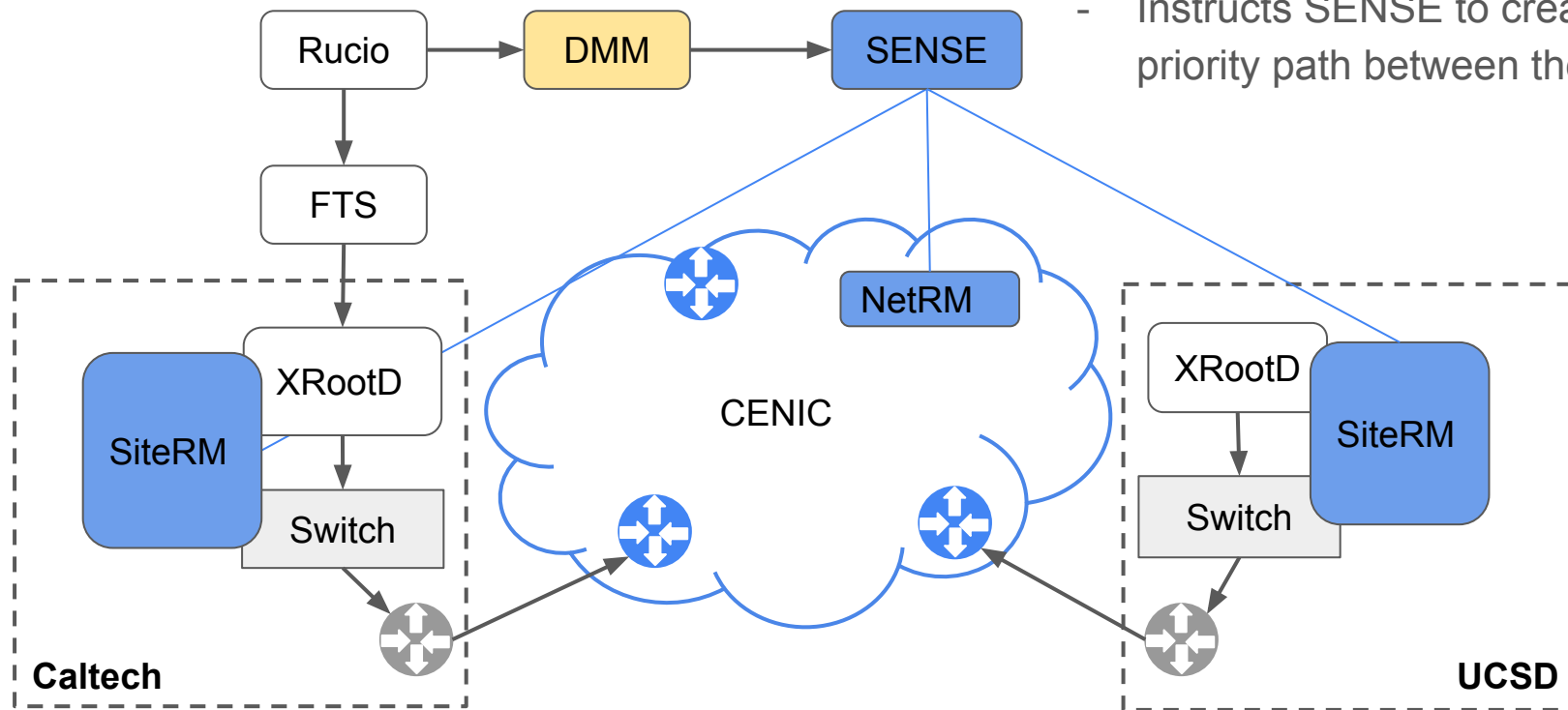By joining forces we can let Rucio leverage SENSE capabilities to:
- Isolate   => different data workflows travel on different subnets
- QoS        => Allocate bandwidth
- Routing  => Select different paths

**Our target is ONLY the LARGEST and/or TIME-SENSITIVE data workflows**

# How it looks

DMM - Data Movement Manager:
- Picks a free subnet at each site
- Calculates bandwidth allocation
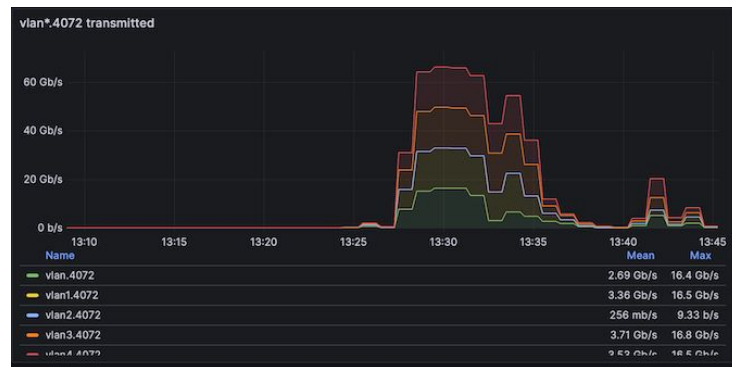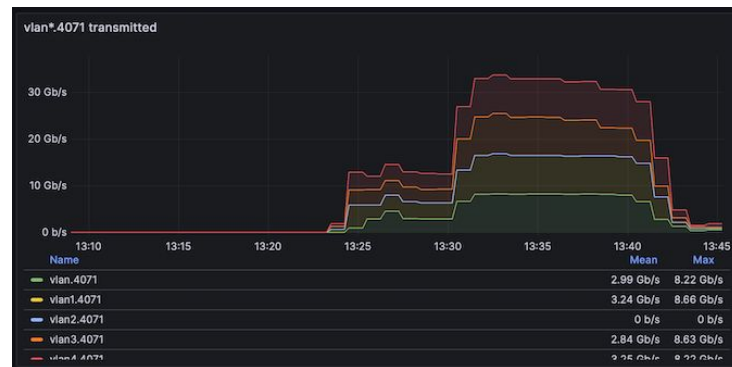- Instructs SENSE to create a priority path between them

# Project status

- Last year: PofC at 10Gbps with 1 managed allocation and background traffic.
- Last week:  Managed 100Gbps between 2 allocations (more on this later)
- Working on new site deployments at UNL and FNAL
- Agreed with the Rucio team on initial plan for integration of DMM into Rucio
- Demonstrated 350Gbps transfers within our testbed (no SENSE managed)
    - Presented in CHEP23
- New 400Gbps link from LA to ESnet allows us to expand our testbed beyond UCSD and Caltech
    - Planning to do high-throughput O(100Gbps) tests at higher latencies

# Mini-challenges status

**#1 (COMPLETED).** High bandwidth demonstration of **multiple** Rucio initiated priority data transfers between UCSD and Caltech

In the image on the right we can see a 100Gbps link between UCSD and Caltech being shared by 2 Priority Paths created by SENSE, each of them using only its allocated share 33/66 Gbps (top/bottom)

**#2 (ONGOING).** Demonstration of 3 priority paths between 3 different pairs of sites. We are currently testing FNAL's deployment.



**Mini-challenge #1.  Two Priority Paths created between SDSC and Caltech**

# Proposal for DC24

1. **Any day**: 400Gbps of SENSE managed throughput **SDSC => Caltech**
   a. completely orthogonal: via CENIC i.e. no prod links used
   b. <u>do we want to monitor/account this for as DC?</u>
2. **Day 4** ("reprocessing"): **FNAL => Caltech & SDSC => Caltech** share of 100Gbps of SENSE managed throughput. The share 66/33 gets swapped in intervals of 1 hr
   a. Will use a non-Prod path from FNAL -> ESnet
   b. within ESnet level the link is shared with Prod
3. **Day 11 or 12** (Contingency): **FNAL => Caltech** at max throughput available via SENSE
   a. Same as #2

For 2 and 3: we would ask for **green-light from FNAL** and **synch up with DUNE** activities**.** Also we **need to agree on how to monitor**: "RSE names" and "Activity"

For all 3 cases NO Production Storage would be used

# How Rucio/SENSE can improve DC in the future?

Monitoring is a very important part of DC but obtaining precise throughput between a given pair of sites is not trivial.

Rucio/SENSE can create priority paths between a given pair of sites for the "Data Challenge" injected data and provide not only a bandwidth guaranteed path but also instantaneous throughput monitoring for such path.
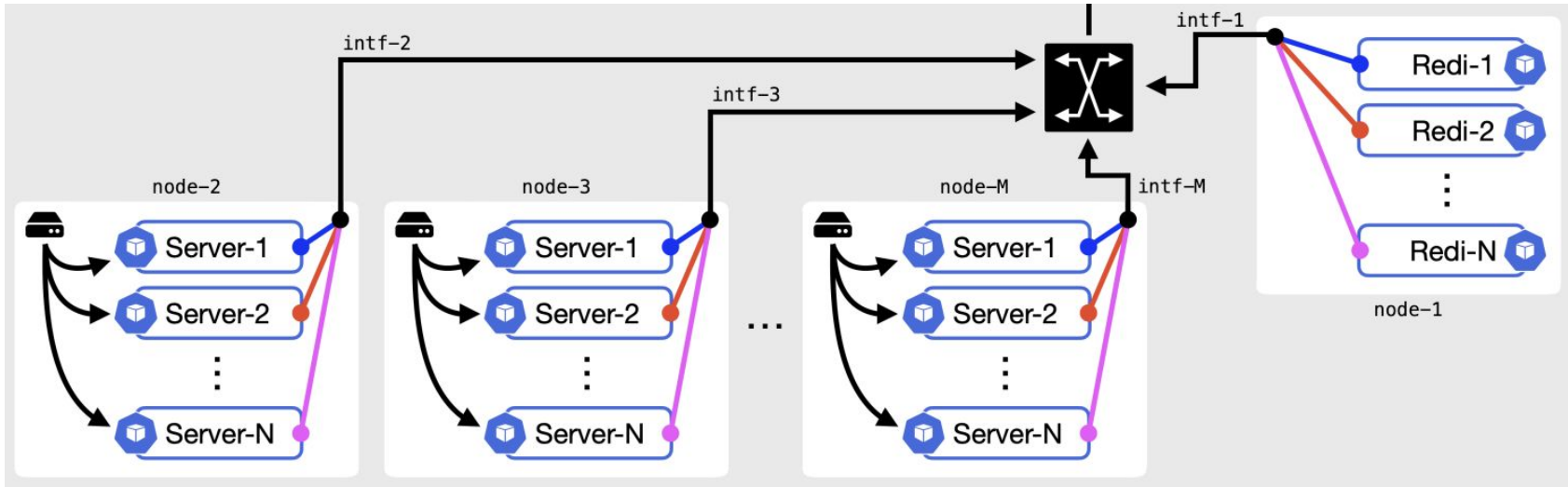
# Questions?

**ACKNOWLEDGMENTS**

# Backup slides

# Isolation using XRootD multi-endpoint

- A single data server is configured to listen in N different IPv6 addresses.
- We use IPv6 because we need many IP addresses



XRootD cluster with M servers and N subnets, Every color represents a different subnet

# New 400Gbps