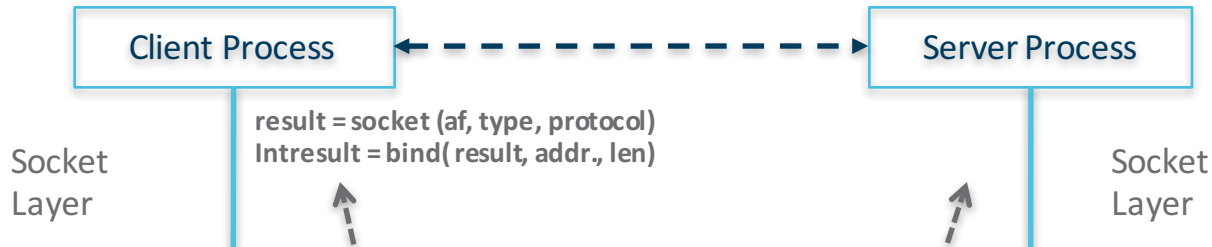# Cracking open the network 'black-box'

Inder Monga
Executive Director, Energy Sciences Network
Division Director, Scientific Networking
Lawrence Berkeley National Lab

Supercomputing Asia
March 12th, 2019
Singapore

ESnet

# Socket Interface: Most successful data plane abstraction
# Forces the network to be a <u>black box</u> to applications

Client Process ◄ - - - - - - - - - ► Server Process

result = socket (af, type, protocol)
Intresult = bind( result, addr., len)

Socket
Layer

Socket
Layer

- Gives file system like abstraction to the capabilities of the network
- Hides the complexity of the network and its operation

BGP        IS-IS

MPLS              OTN

        Carrier
OSPF    Ethernet

..and many more

Complexity:
7000+ IETF RFCs
ITU-T
IEEE
GSM
Others…

ESnet

# High-performance science network user facility
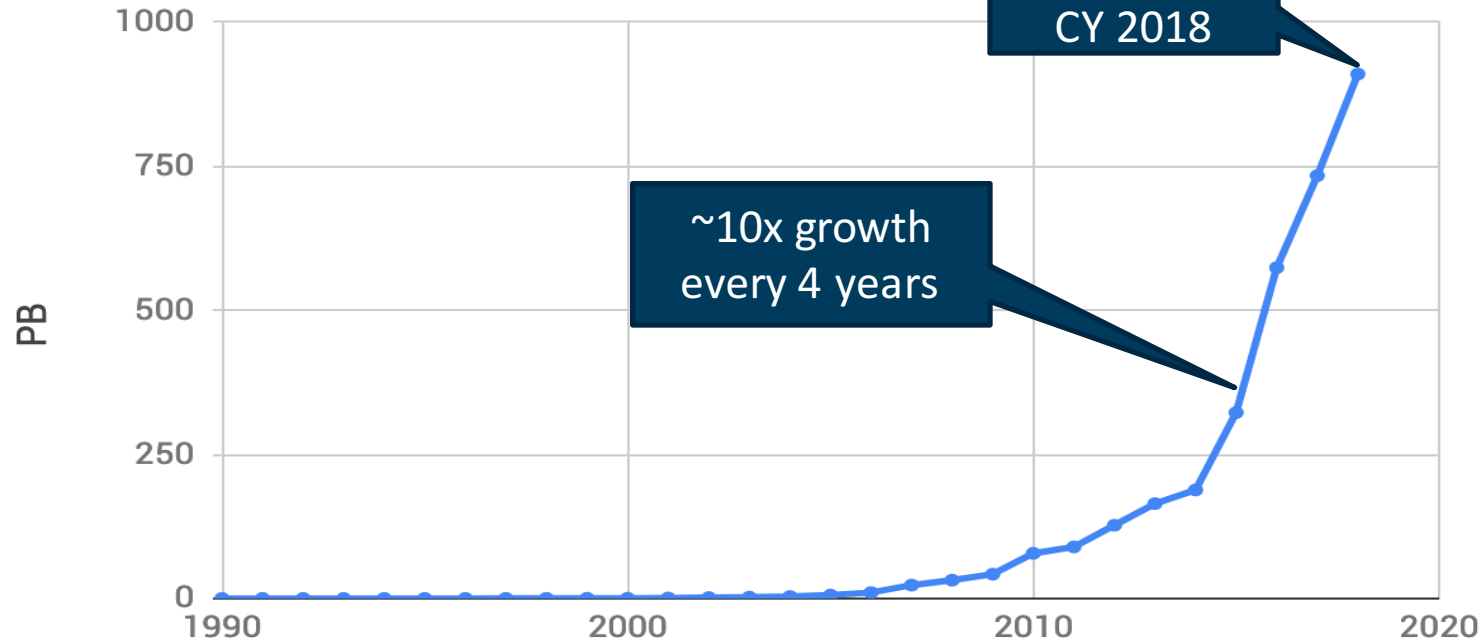# Optimized for enabling big-data science



Provides connectivity to <u>all of the DOE labs</u>,
experiment sites, & user facilities (> 34417 users)

# An ~exabyte scale network today

## Data Moved (in PB) vs. Year

910PB in CY 2018

~10x growth every 4 years

PB

1000

750

500

250

0

1990          2000          2010          2020

exponential traffic growth over past 28 years
*measures ingress or egress only, not traffic per link*

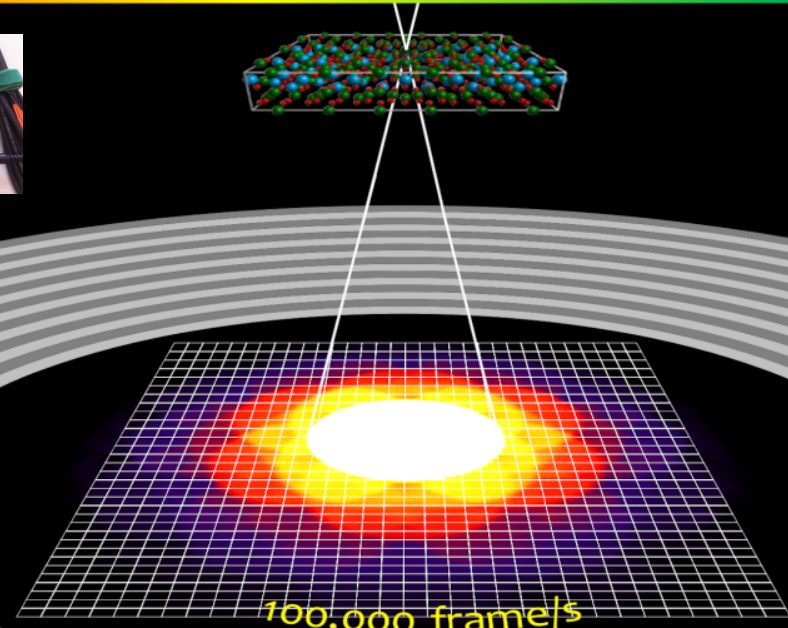ESnet

# New instruments, more data: NCEM 4D-Stem



FPGA based readout system

400 –

1 Tb/s

100,000 frame/s Pixilated Detector

Segmented HAADF Detector

MOLECULAR FOUNDRY

# New instruments, more data : LCLS-II



LU34 experiment: Taking Snapshots of O-O Bond Formation in Photosynthetic Water-Splitting Using Simultaneous X-ray Emission Spectroscopy and Crystallography – Y. Vital (LCLS PI)

Diffraction pattern from LU34

Peak Throughput (prior to data reduction)

~1 Tb/s

# Science DMZ architecture [ESnet] has been impactful around the world and adopted by {Pacific, National, Asia} Research platforms

Petascale DTN Project

November 2017
L380 Data Set

Gigabits per second (min/avg/max), three transfers
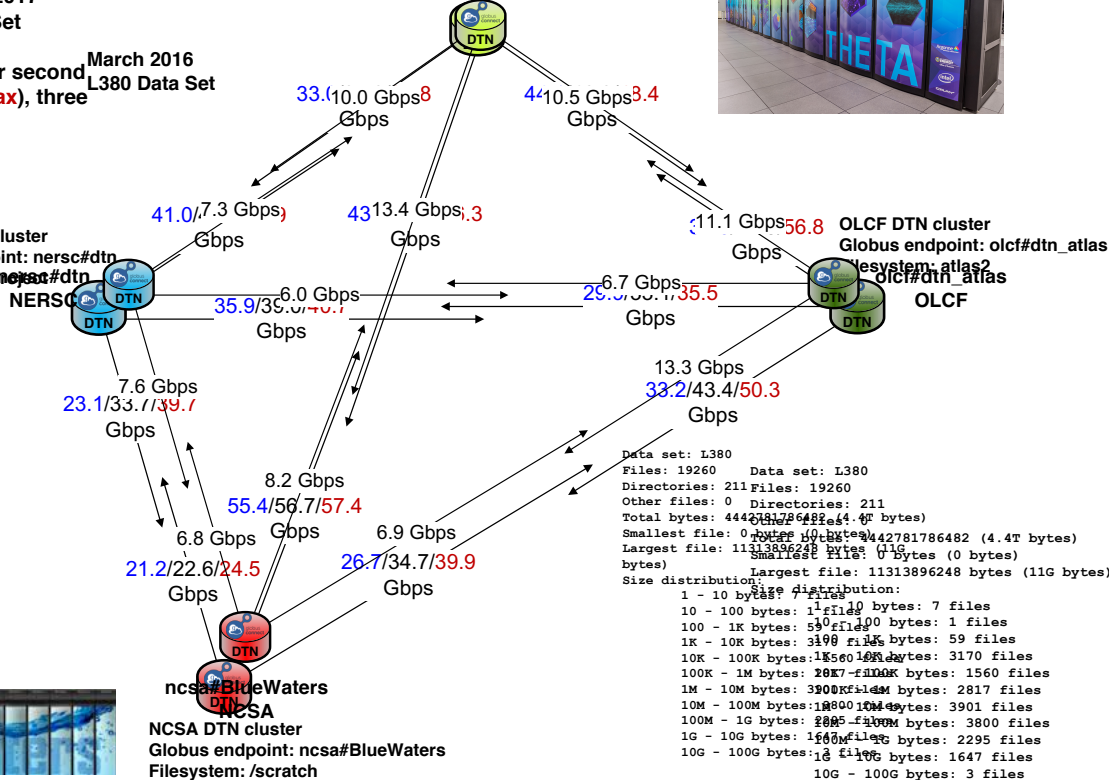
March 2016
L380 Data Set

**ALCF DTN cluster**
Globus endpoint: alcf#dtn_mira
Filesystem: /projects

ALCF

THETA

33.0 (10.0 Gbps) 8 Gbps

44 10.5 Gbps 3.4 Gbps

41.0/47.3 Gbps Gbps

43 13.4 Gbps 1.3 Gbps

11.1 Gbps 56.8 Gbps

**NERSC DTN cluster**
Globus endpoint: nersc#dtn
Filesystem: /project

NERSC

Cori

35.9/39.0/46.7 6.0 Gbps Gbps

29.0/33.1/35.5 6.7 Gbps Gbps

**OLCF DTN cluster**
Globus endpoint: olcf#dtn_atlas
Filesystem: atlas2

olcf#dtn_atlas
OLCF

summit

23.1/33.7/39.7 7.6 Gbps Gbps

13.3 Gbps
33.2/43.4/50.3 Gbps

55.4/56.7/57.4 8.2 Gbps Gbps

21.2/22.6/24.5 6.8 Gbps Gbps

26.7/34.7/39.9 6.9 Gbps Gbps

ncsa#BlueWaters
NCSA

**NCSA DTN cluster**
Globus endpoint: ncsa#BlueWaters
Filesystem: /scratch

Data set: L380
Files: 19260
Directories: 211
Other files: 0
Total bytes: 4442810786482 (4.4T bytes)
Smallest file: 0 bytes (0 bytes)
Largest file: 11313896248 bytes (11G bytes)
Size distribution:
    1 - 10 bytes: 7 files
    10 - 100 bytes: 1 files
    100 - 1K bytes: 59 files
    1K - 10K bytes: 3170 files
    10K - 100K bytes: 1560 files
    100K - 1M bytes: 2817 files
    1M - 10M bytes: 3901 files
    10M - 100M bytes: 3800 files
    100M - 1G bytes: 2295 files
    1G - 10G bytes: 1647 files
    10G - 100G bytes: 3 files

Data set: L380
Files: 19260
Directories: 211
Other files: 0
Total bytes: 4442810786482 (4.4T bytes)
Smallest file: 0 bytes (0 bytes)
Largest file: 11313896248 bytes (11G bytes)
Size distribution:
    1 - 10 bytes: 7 files
    10 - 100 bytes: 1 files
    100 - 1K bytes: 59 files
    1K - 10K bytes: 3170 files
    10K - 100K bytes: 1560 files
    100K - 1M bytes: 2817 files
    1M - 10M bytes: 3901 files
    10M - 100M bytes: 3800 files
    100M - 1G bytes: 2295 files
    1G - 10G bytes: 1647 files
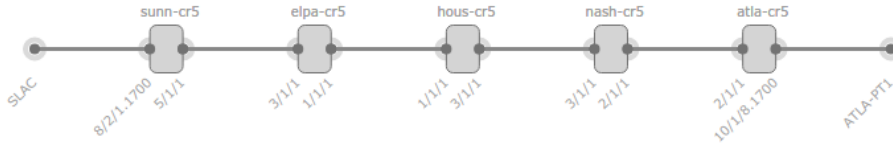    10G - 100G bytes: 3 files

Data courtesy of Eli Dart

ESnet

# Data movement software keeps on improving: from 1 PB/week to 1 PB/day (approx.)

HOME › OSCARS »

## SLAC latency loop - 1 of 2 - OVERRIDE - VLAN 1700

sunn-cr5    elpa-cr5    hous-cr5    nash-cr5    atla-cr5

SLAC    8/2/1.1700    5/1/1    3/1/1    1/1/1    1/1/1    3/1/1    3/1/1    2/1/1    2/1/1    10/1/8.1700    ATLA-PT1

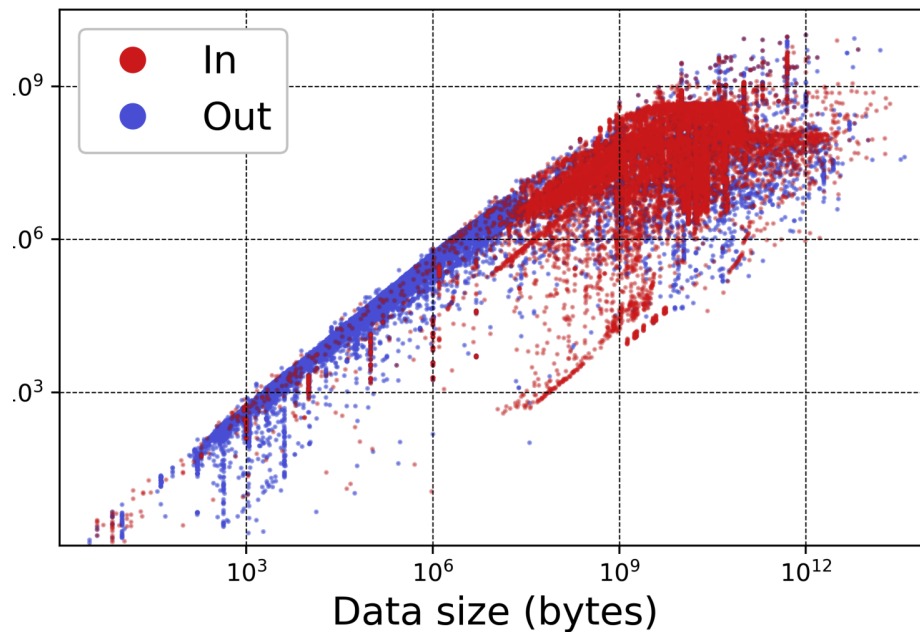**ESnet's Network, Software Help SLAC Researchers in Record-Setting Transfer of 1 Petabyte of Data**

Using a 5,000-mile network loop operated by ESnet, researchers at the SLAC National Accelerator Laboratory (SLAC) and Zettar Inc. (Zettar) recently transferred 1 petabyte in 29 hours, with encryption and checksumming, beating last year's record by 5 hours, almost a 15 percent improvement.

40G
20G
0.0
20G
40G
60G
80G
100G

09 AM    12 PM    03 PM    06 PM    09 PM    Thu 27    03 AM    06 AM    09 AM    12 PM    03 PM    06 PM

3/21/19

ESnet

# But, from an application perspective, what happens when your network data transfer fails?



ESnet

# Even well tuned infrastructure does not get consistent service



**(d)** Petrel

Nine orders of variability

Chard K, Dart E, Foster I, Shifflett D, Tuecke S, Williams J. (2018) The Modern Research Data Portal: a design pattern for networked, data-intensive science. *PeerJ Computer Science* 4:e144  https://doi.org/10.7717/peerj-cs.144

# How to 'crack open the network black box' without destroying the power of the abstraction?

1. High-precision telemetry

2. Scalable analytics infrastructure

3. Model-based approach to request network services

4. Network prediction using machine learning techniques

ESnet

# 1. High-precision telemetry: deep insight into flows



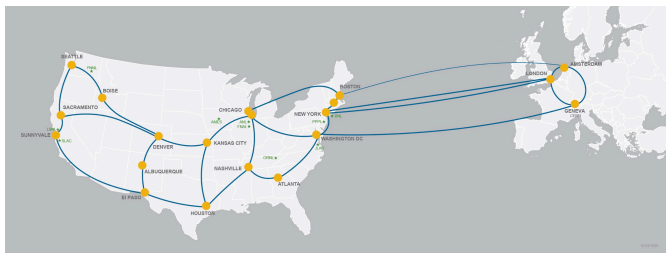Jupiter with the naked eye



Jupiter Close Up

**Per flow, high-precision telemetry**

- Per packet-metadata tracking (e.g. timestamp, ingress location, etc)

- 10 ns precision in timing

Use high-fidelity data to get better insights and analytics:
- Packet Microbursts
- Path deviations ( RTT and Delay )
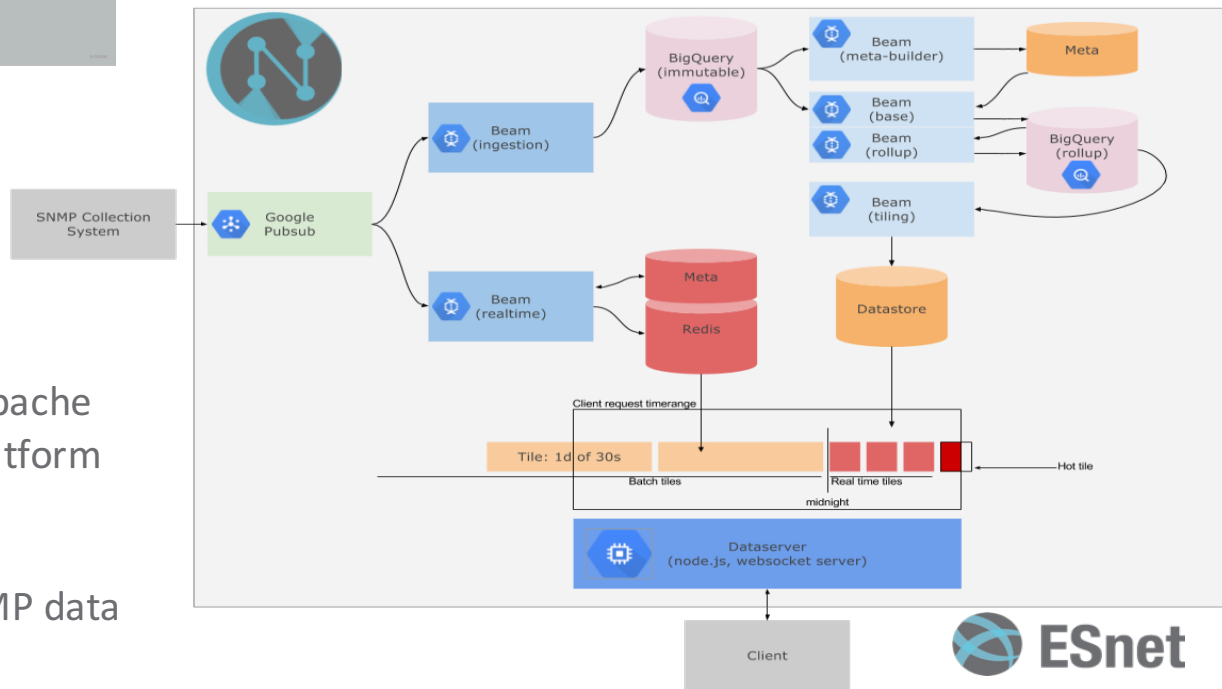- Security / anomaly detection
- Head of Queue Blocking
- Many others...

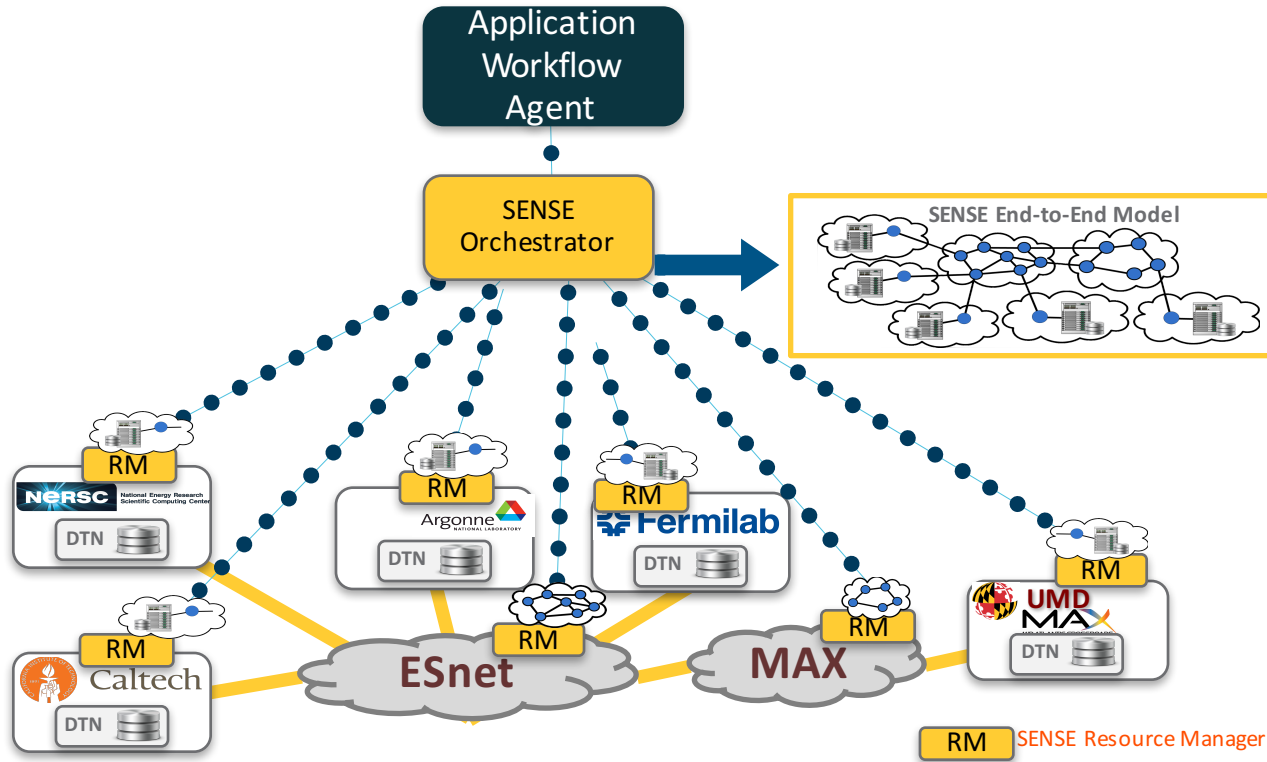ESnet

# 2. Scalable analytics infrastructure



Streams telemetry ➤

On-demand analytics infrastructure



- Real-time telemetry from the network
- **netbeam** platform: Using Apache BEAM and Google Cloud Platform
- Both Batch and Stream processing in parallel
- In production for ESnet SNMP data
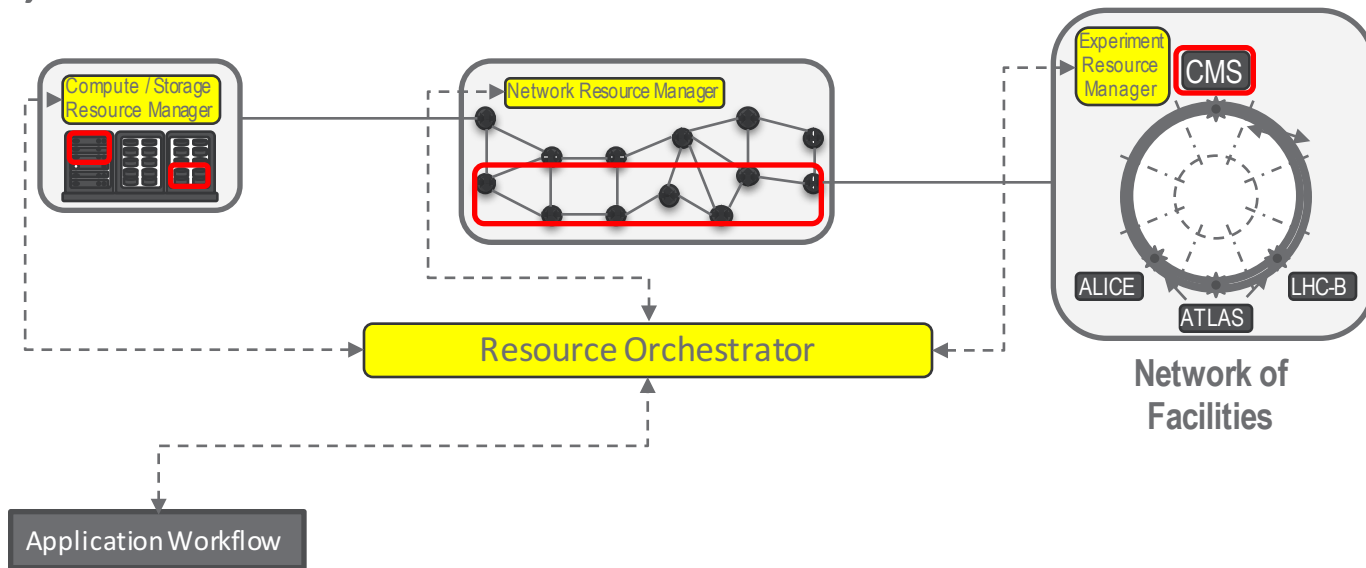
# 3. Model based approach to request network services



Multi-resource models abstracts the network specifics away and allows for higher-level service request

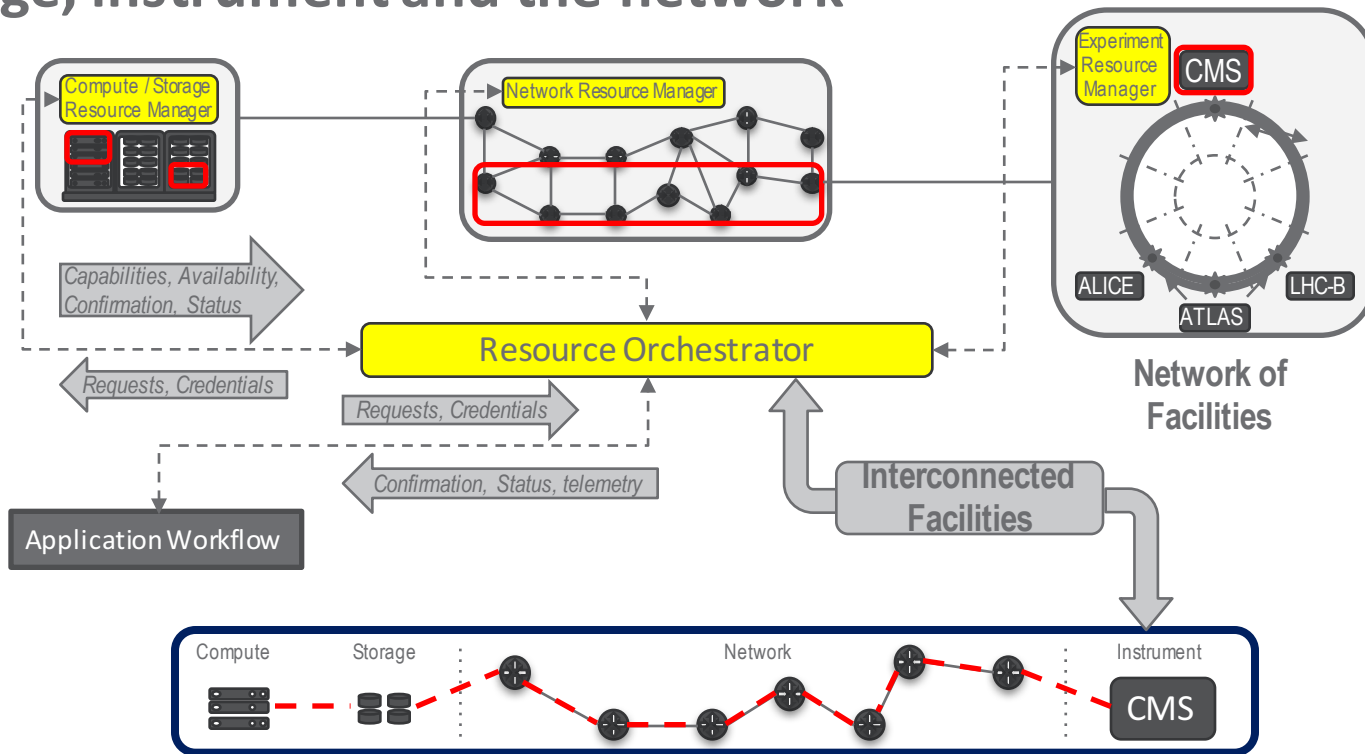*SENSE project is a research project funded by ASCR, DOE, US*

*Realtime system based on Resource Manager developed infrastructure and service models*

# 3. Enables an end-to-end service model that includes compute, storage, instrument and the network



Network of Facilities

Compute / Storage Resource Manager

Network Resource Manager

Experiment Resource Manager

CMS

ALICE

ATLAS

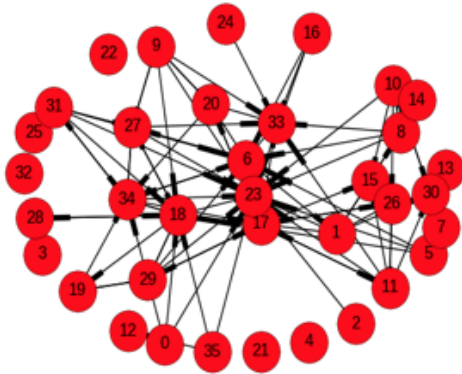LHC-B

Resource Orchestrator

Application Workflow

ESnet

# 3. Enables an end-to-end service model that includes compute, storage, instrument and the network
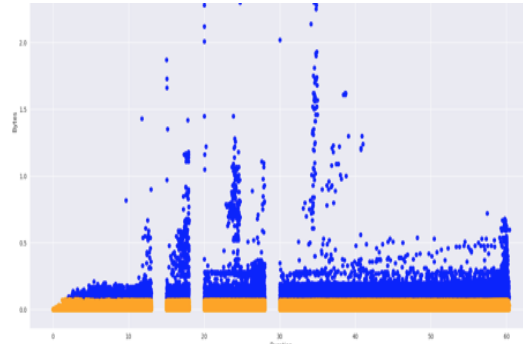


Slide courtesy of Chin Guok

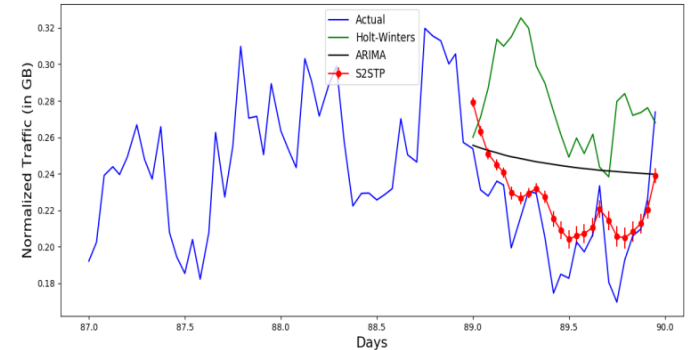# 4. Network prediction using Machine Learning techniques

*Understanding which sites are busiest at different times*

*High-Speed classifying of big and small flows to redirect packet routes*
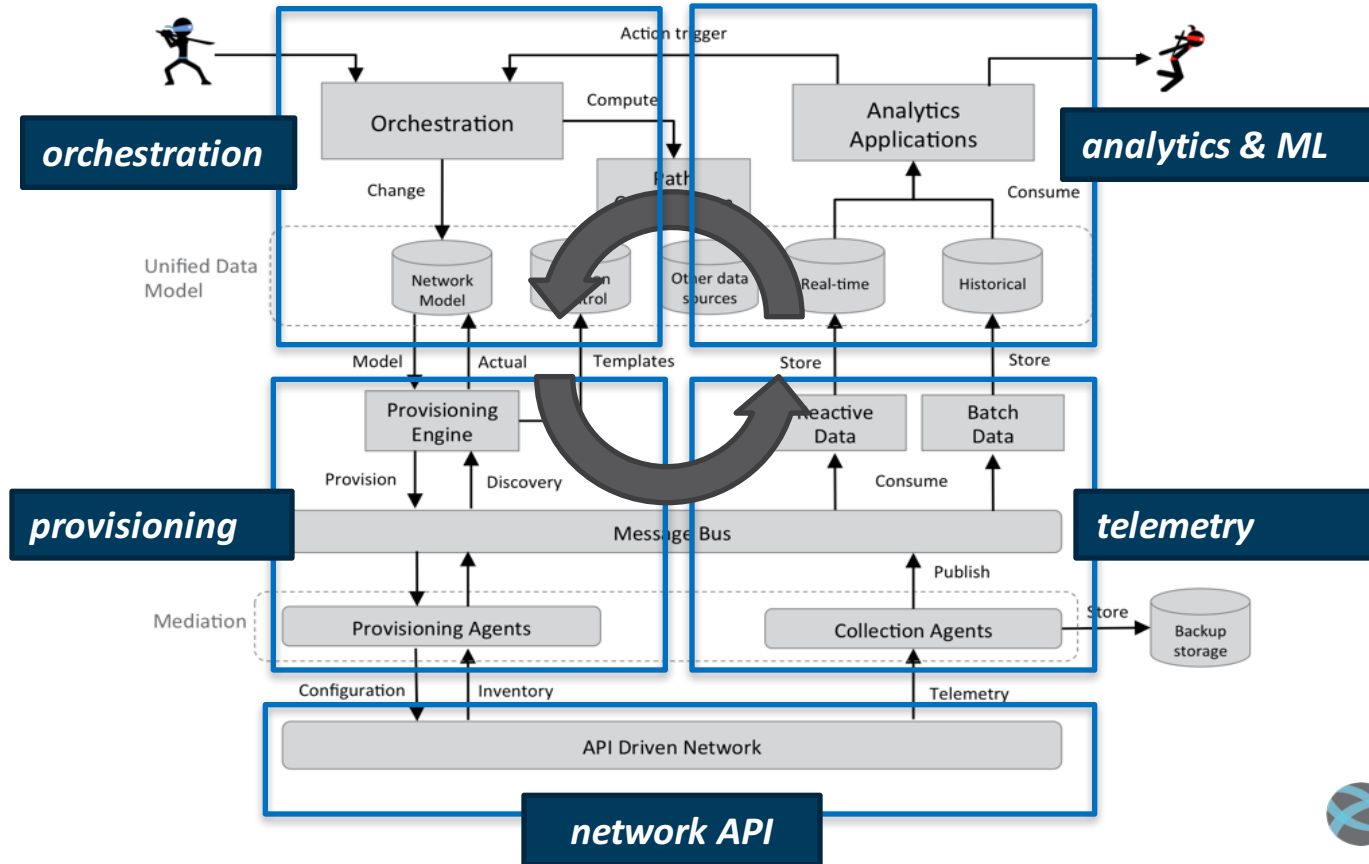
*Prevent congestion and links failures by anticipating traffic 24 hours in advance*



Slide courtesy of Mariam Kiran

Network traffic prediction will help us bring route and engineer flows appropriately and on-the-fly

ESnet

# ESnet's next-generation (ESnet6) software architecture



**orchestration**

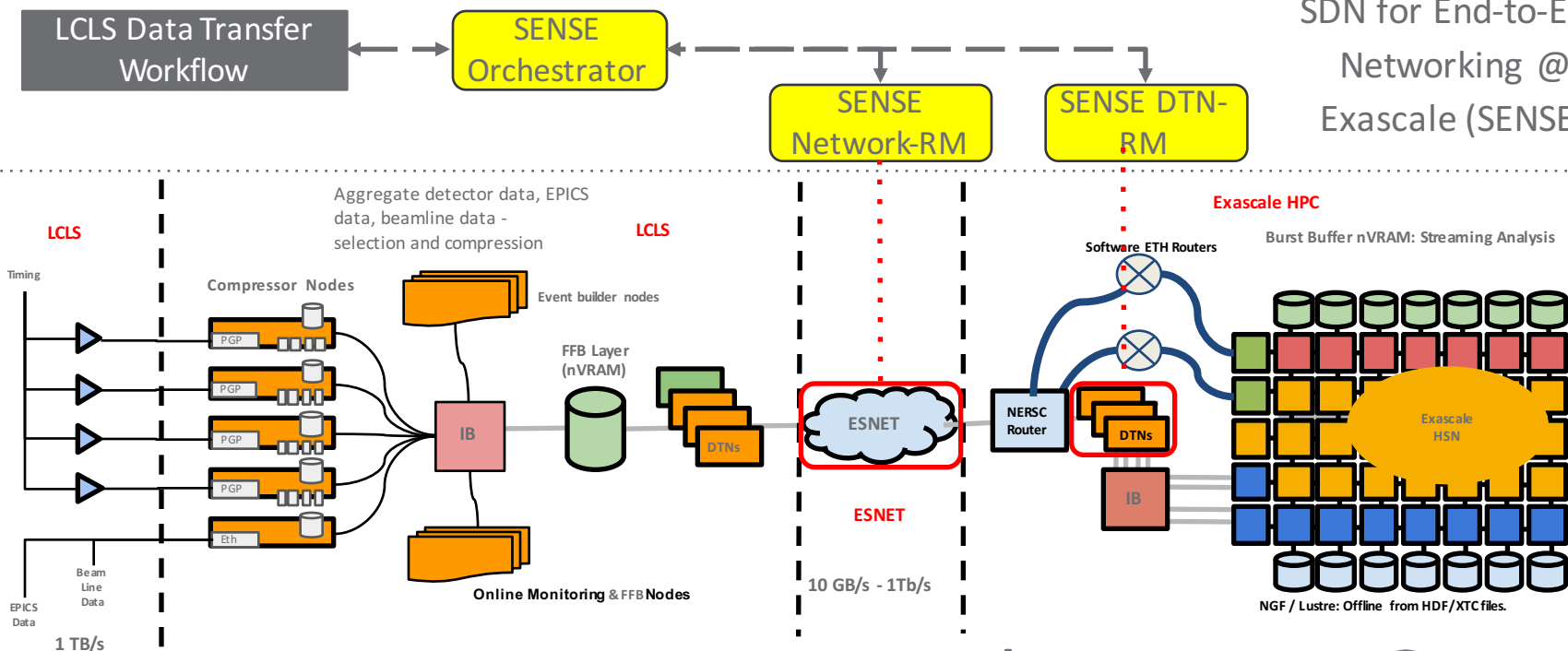**analytics & ML**

**provisioning**

**telemetry**

**network API**

Action trigger

Orchestration

Compute

Change

Analytics
Applications

Consume

Path

Unified Data
Model

Network
Model

Other data
sources

Real-time

Historical

Model

Actual

Templates

Store

Store

Provisioning
Engine

Reactive
Data

Batch
Data

Provision

Discovery

Consume

Message Bus

Publish

Mediation

Provisioning Agents

Collection Agents

Store

Backup
storage

Configuration

Inventory

Telemetry

API Driven Network

ESnet

# Superfacility model for productive, reproducible science



Light Sources

Sequencers

ESnet

Computing and
Data Facilities

Interconnected facilities where data is
acquired, stored, analyzed and
served

Telescopes

Experimental
Facilities

Expertise

CAMERA
Applied Math

Particle
Detectors

Jupyter

VisIt

Microscopes

User Community

Environmental Sensors

ESnet

# ExaFEL: A science example of the Superfacility model

# The circulatory system for DOE collaborative science



- **Networking tailored to meeting science demands**
  - Bandwidth reservations, performance monitoring, Science "DMZ" model for the last mile,..
- An effective **network design and application interaction** is extremely important to meet the needs of next-generation of data science

ESnet

# Thank you.

imonga@es.net

ESnet