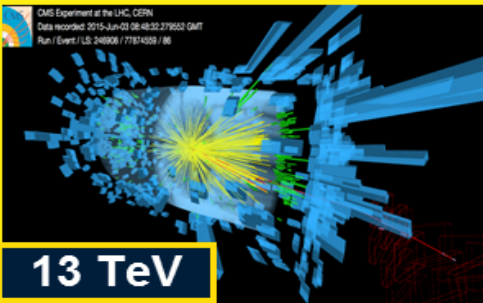# The GNA-G Missions and Working Groups
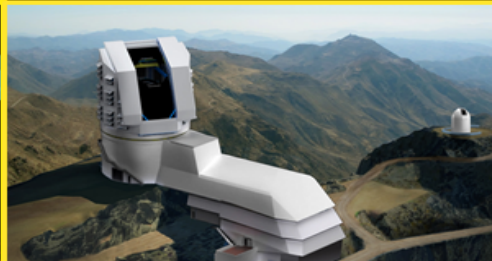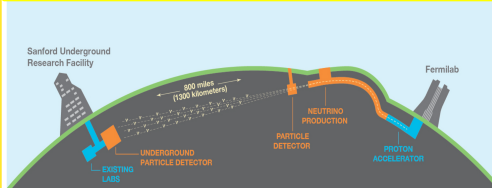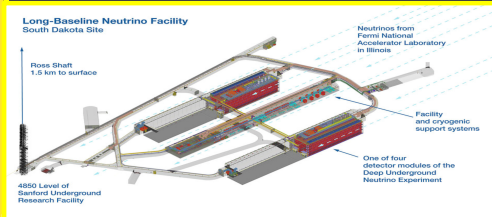## Next Generation Networks for Global Science Programs

LHC Run3 and HL-LHC

DUNE

LSST   SKA

BioInformatics

Earth Observation

*Gateways to a New Era*

13 TeV

LSST

LBNF/DUNE

LHC

SKA

**Harvey Newman, Caltech**
GNA-G Virtual Meeting
May 26, 2020

newman@hep.Caltech.edu   SENSE DOE-AC02-07CH11359   SANDIE NSF CC* 1659403

https://www.gna-g.net/

GNA-G
Global Network Advancement Group

# A New Era of Challenges: Global Exabyte Data Distribution, Processing, Access and Analysis

**GNA-G**
Global Network Advancement Group

## Exascale Data for the LHC Experiments

- **~1 Exabyte Stored by 2019; to ~ 10-50 EB during HL LHC Era**

## Network Flow: 45-60 Gbytes/sec

- **~1.5 Exabyte flowed over WLCG in 2019**

## Emergence of 400-800G in Hyper-Data Centers, 100-200G on Terrestrial WANs

- **400G in Wide Area by 2022 ?**

## Network Dilemma: Per technology generation (~10 years)

- **Capacity at same unit cost: 4X**
- **Bandwidth growth: 35-70X in Internet2, GEANT, ESnet**

## *LHC Run3: likely reach a network limit*

## Unlike the past: Optical and switch advances are evolutionary

### *Physics Limits by ~HL LHC Start*

## New Levels of Challenge

- **Global data distribution, processing, access and analysis**
- **Coordinated use of massive but still limited *diverse* compute, storage and network resources**
- **Coordinated operation and collaboration *within and among* scientific enterprises**



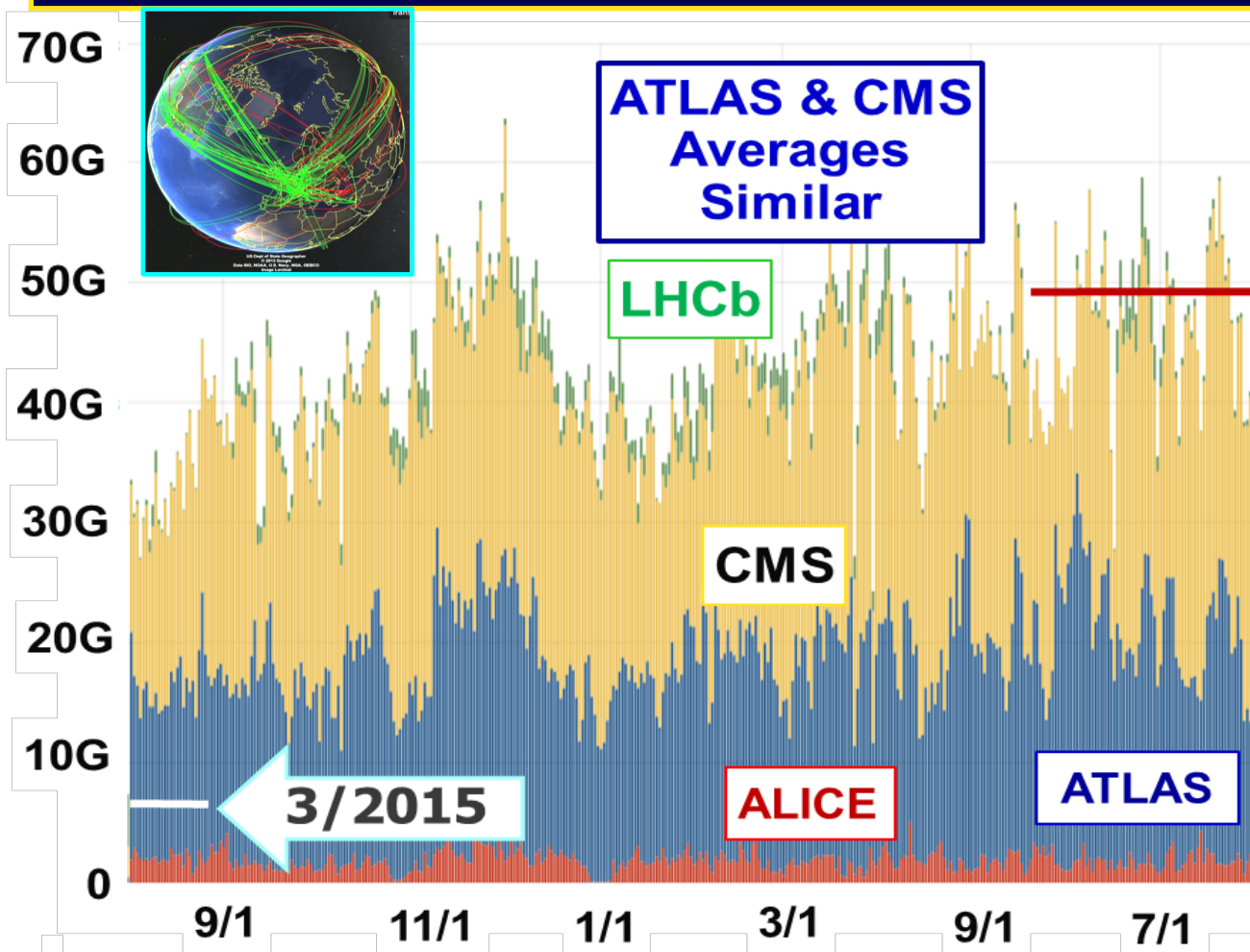Earth Observation | LCLS-II

- **HEP will experience increasing Competition from other data intensive programs**
  - **Sky Surveys: LSST, SKA**
  - **Next Gen Light Sources**
  - **Earth Observation**
  - **Genomics**

# LHC Data Flows Have *Increased* in Scale and Complexity since the start of LHC Run2 in 2015

**GNA-G**
Global Network Advancement Group

## WLCG Transfers Dashboard: Throughput Aug. 2018 – Aug. 2019

*45-50 GBytes/s Sustained*
*60+ GBytes/s Peaks*

**Complex Workflow**

- **700k jobs (threads) simultaneously**
- **Multi-TByte to Petabyte Transfers;**
- **6-17 M File Transfers/Day**
- **100ks of remote connections**

ATLAS & CMS Averages Similar

LHCb

CMS

ALICE

ATLAS

3/2015

70G
60G
50G
40G
30G
20G
10G
0

9/1   11/1   1/1   3/1   9/1   7/1

*7X Growth in Sustained Throughput in 4.3 Years: +60%/Yr; ~100X per Decade*
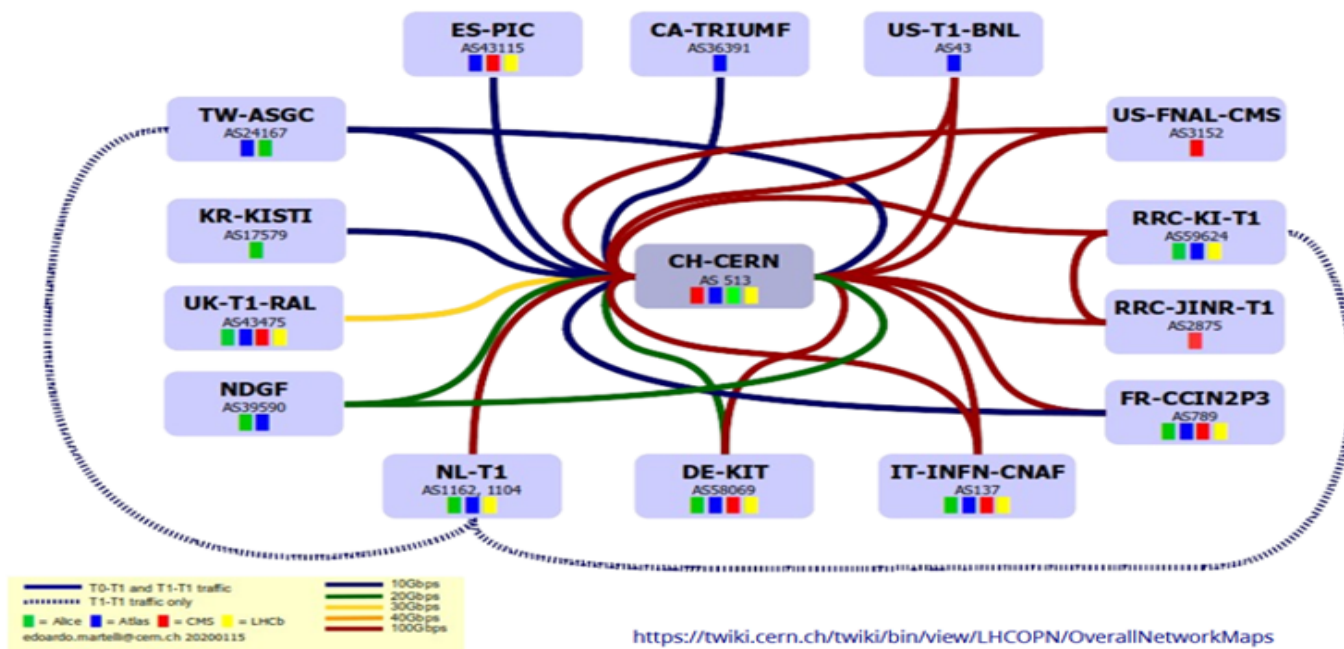
**Edoardo Martelli**
**CERN/IT**
**LHCONE Meeting May 2020**

## Numbers

- 14 Tier1s + 1 Tier0

- 12 countries in 3 continents

- Dual stack IPv4-IPv6

- 1.1Tbps to the Tier0
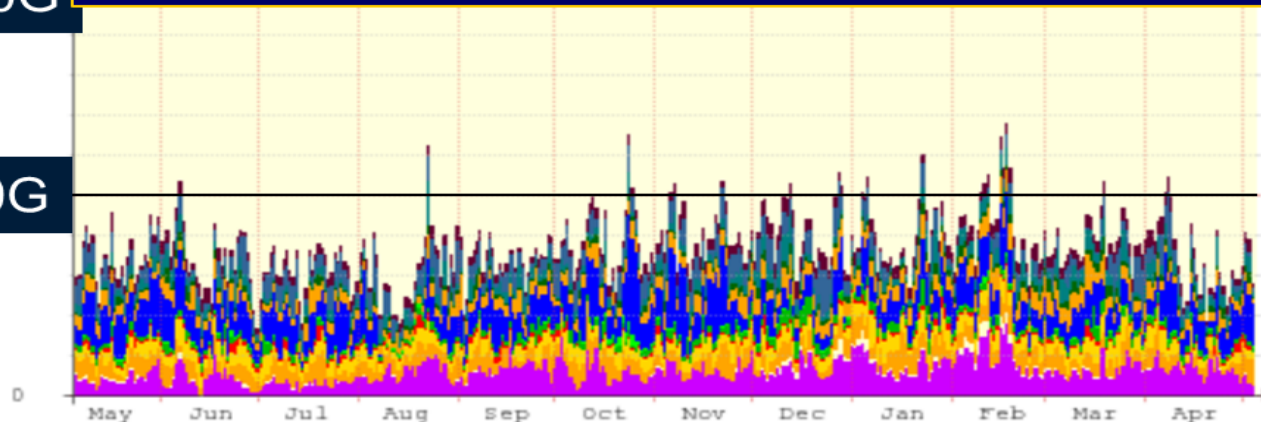
- Moved ~288 PB in the last 12 months

**LHCOPN**

**CERN ⟺ Tier1s**

ES-PIC AS43115
CA-TRIUMF AS36391
US-T1-BNL AS43
TW-ASGC AS24167
US-FNAL-CMS AS3152
KR-KISTI AS17579
RRC-KI-T1 AS59624
UK-T1-RAL AS43475
CH-CERN AS 513
RRC-JINR-T1 AS2875
NDGF AS39590
FR-CCIN2P3 AS789
NL-T1 AS1162, 1104
DE-KIT AS58069
IT-INFN-CNAF AS137

T0-T1 and T1-T1 traffic
T1-T1 traffic only
= Alice  = Atlas  = CMS  = LHCb
edoardo.martelli@cern.ch 20200115
10Gbps
20Gbps
30Gbps
40Gbps
100Gbps

https://twiki.cern.ch/twiki/bin/view/LHCOPN/OverallNetworkMaps

**LHCOPN 5/19 to 4/20: Peaks over 100G**

200G

100G

May Jun Jul Aug Sep Oct Nov Dec Jan Feb Mar Apr

# LHCONE: a Virtual Routing and Forwarding (VRF) Fabric
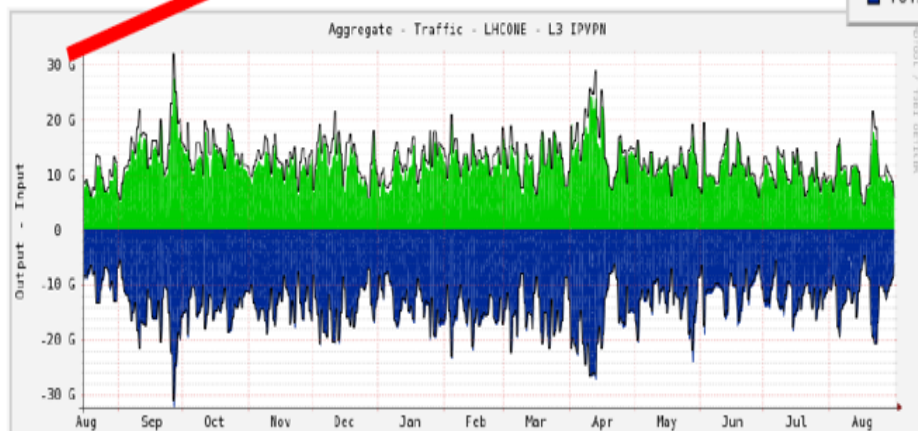## Global infrastructure for *HEP (LHC, Belle II, NOvA, Auger, Xenon)* data flows

## Where were we?

**LHCONE in Europe**
**GEANT**

GÉANT

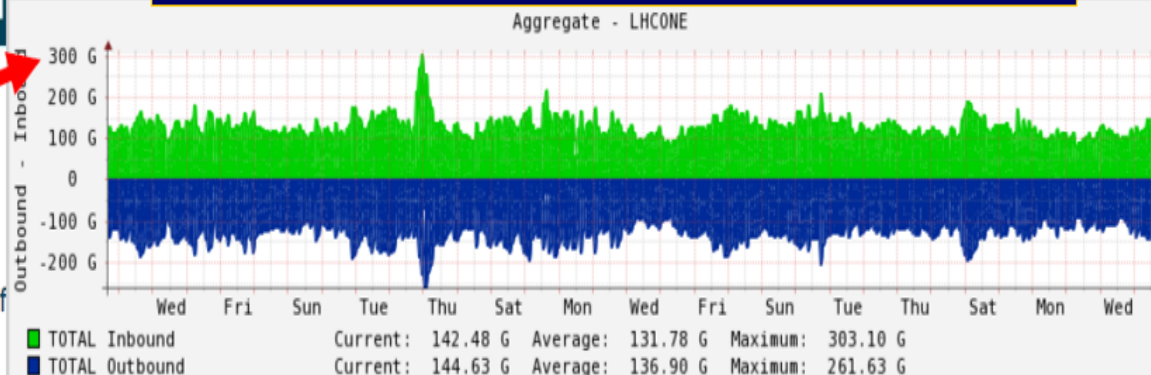**Now: LHCONE: Peaks to 300G**

- Aggregate LHCONE traffic from all the NRENs and Peers
  - Average traffic ~25Gbps
  - Sustained Peaks ~35Gbps
  - Trans-Atlantic Traffic ~ 20Gbps (Peak)
- Graphs shows 1 day average traffic over last 12 months the peak traffic is much higher

**10x**

Aggregate - LHCONE

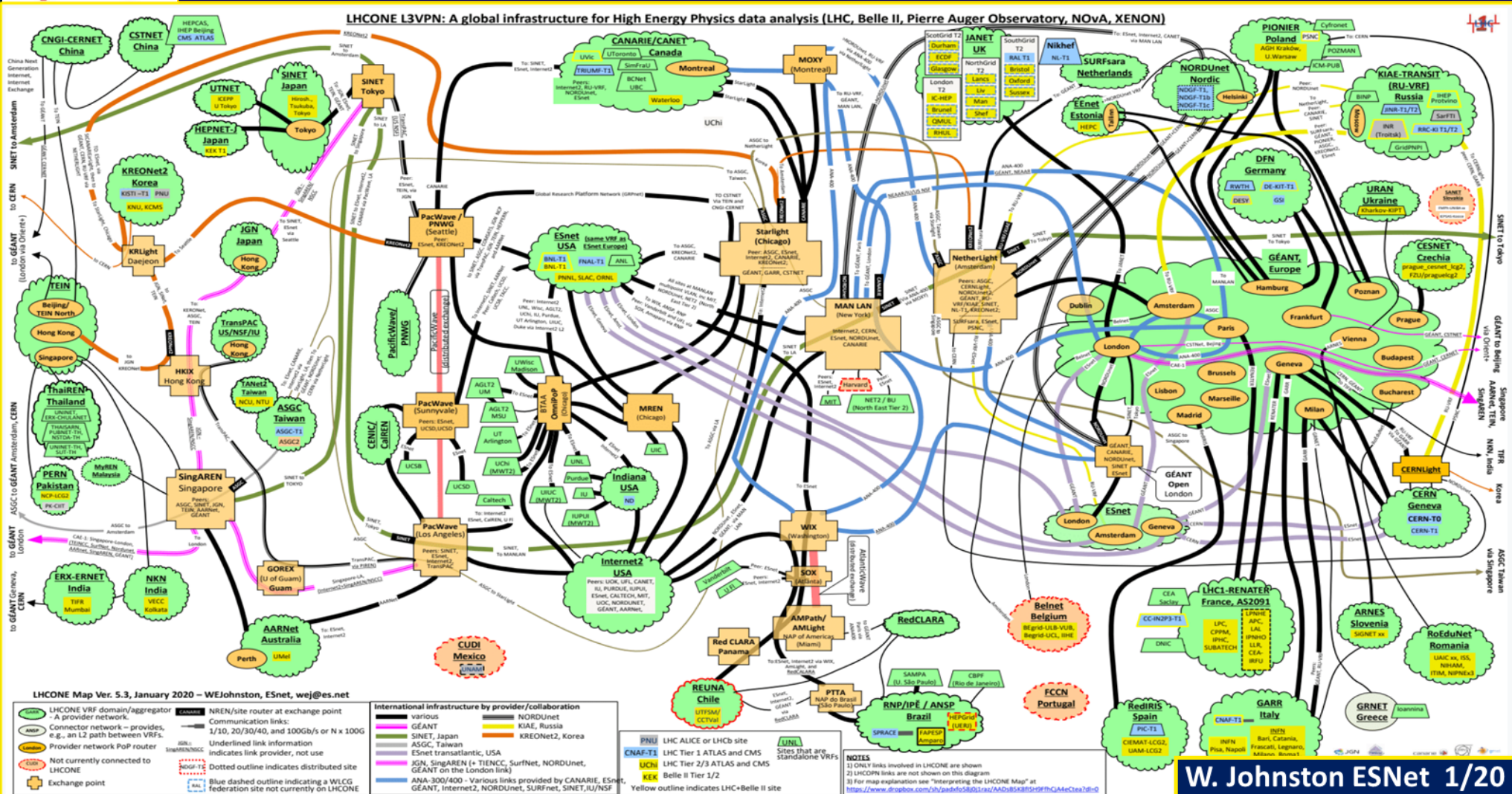| | | | |
|---|---|---|---|
| TOTAL Inbound | Current: 142.48 G | Average: 131.78 G | Maximum: 303.10 G |
| TOTAL Outbound | Current: 144.63 G | Average: 136.90 G | Maximum: 261.63 G |

Aggregate - Traffic - LHCONE - L3 IPVPN

**Before Versus After LHC Run 3:**
**5 – 6X Average, 10X Peak Growth**

**Good News:** The Major R&E Networks Have Mobilized on behalf of HEP
LHCONE traffic growing by 60-70%/Yr: **a challenge** already in LHC Run3 (2021-4)

# LHCONE: a Virtual Routing and Forwarding (VRF) Fabric

## Global infrastructure for *HEP (LHC, Belle II, NOvA, Auger, Xenon)* data flows



LHCONE L3VPN: A global infrastructure for High Energy Physics data analysis (LHC, Belle II, Pierre Auger Observatory, NOvA, XENON)

LHCONE Map Ver. 5.3, January 2020 – WEJohnston, ESnet, wej@es.net

W. Johnston ESNet 1/20

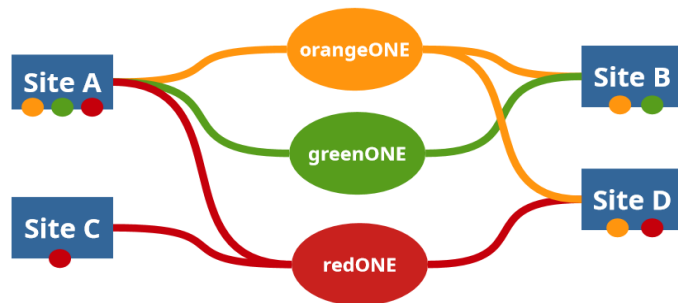**Good News: The Major R&E Networks Have Mobilized on behalf of HEP**
**Challenge: A complex system with limited scaling properties.**
**Response: New Mode of Sharing ? Multi-One ?**

# MultiOne and DUNEOne
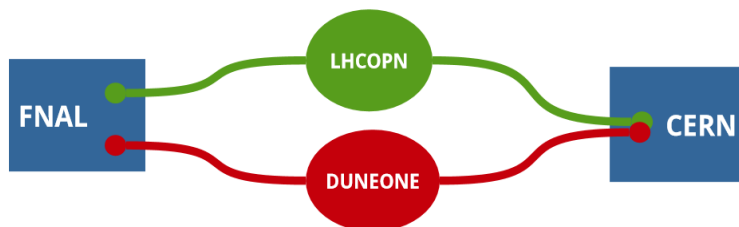

GNA-G
Global Network Advancement Group

## Recap: multiple "ONEs"

- Each site joins only the VPNs it is collaborating with, to reduce the exposure of their data-centre/Science-DMZ
- If doable, each Collaboration funds its own VPN



CERN | IT Information Technology Department

## DUNEONE prototype

ProtoDUNE and DUNE identified as possible use case to build a multiONE prototype
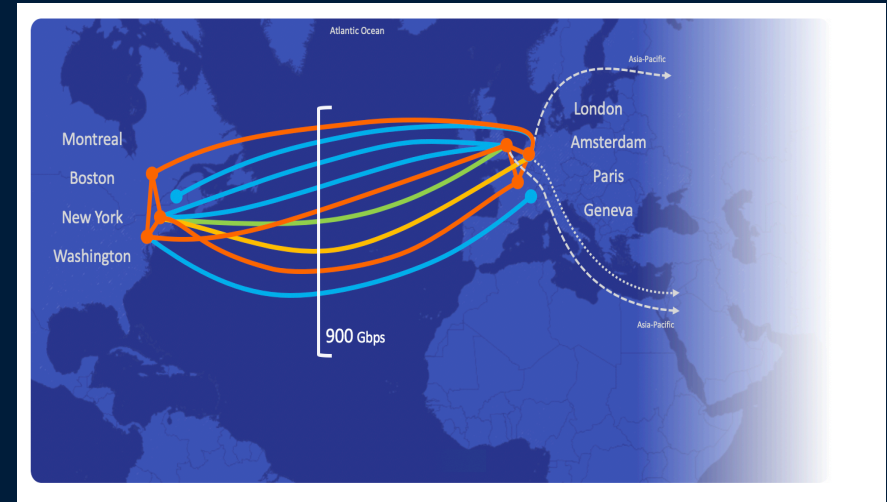


## Status

- Not identified a solution to easily separate traffic, yet

- Explored traffic marking for policy routing with router vendor. Not possible with existing network processor, but it may be possible with upcoming ones

- ESnet is ready to implement a L2 circuit between CERN and FNAL. L3VPN will be considered when necessary at a later stage

- Analysing protoDUNE traffic to check if it could be identified by src and destination addresses

Edoardo Martelli at LHCOPN/LHCONE
Meeting May 13, 2020

# Update: Advanced North Atlantic (ANA) Collaboration

- **Currently: 9x 100 Gbit/s lambdas between GXP**
  - **7x Internet2, NORDUnet, ESnet, SURFnet, CANARIE, and GÉANT**
  - **1x NSF-funded NEAAR Project**
  - **1x Japan's NII/SINET**
- **Started in 2012**
- **First light in 2013**



- **Possible Future Directions from Late 2020 or 2021**
  - **ANAv2: Long-term commitments on bandwidth or spectrum**
  - **ANAv3: At the table with new cable builds, anchor tenantship?**

**Aim: Rightsized, upgradable, resilient bandwidth for less money across the North Atlantic Ocean**

# Annual CMS Data Volume

|  | # of collisions | # of events simulated | RAW event size [MB] | AOD event size [MB] | Total per year [PB] |
|---|---|---|---|---|---|
| Today | 9 Billion | 22 Billion | 0.9 | 0.35 | ~20 |
| HL-LHC | 56 Billion | 64 Billion | 6.5 | 2 | ~600 |

**The beams get "brighter" by x6**
**Data taking rate goes up by x6**
**Simulations go up by x3**

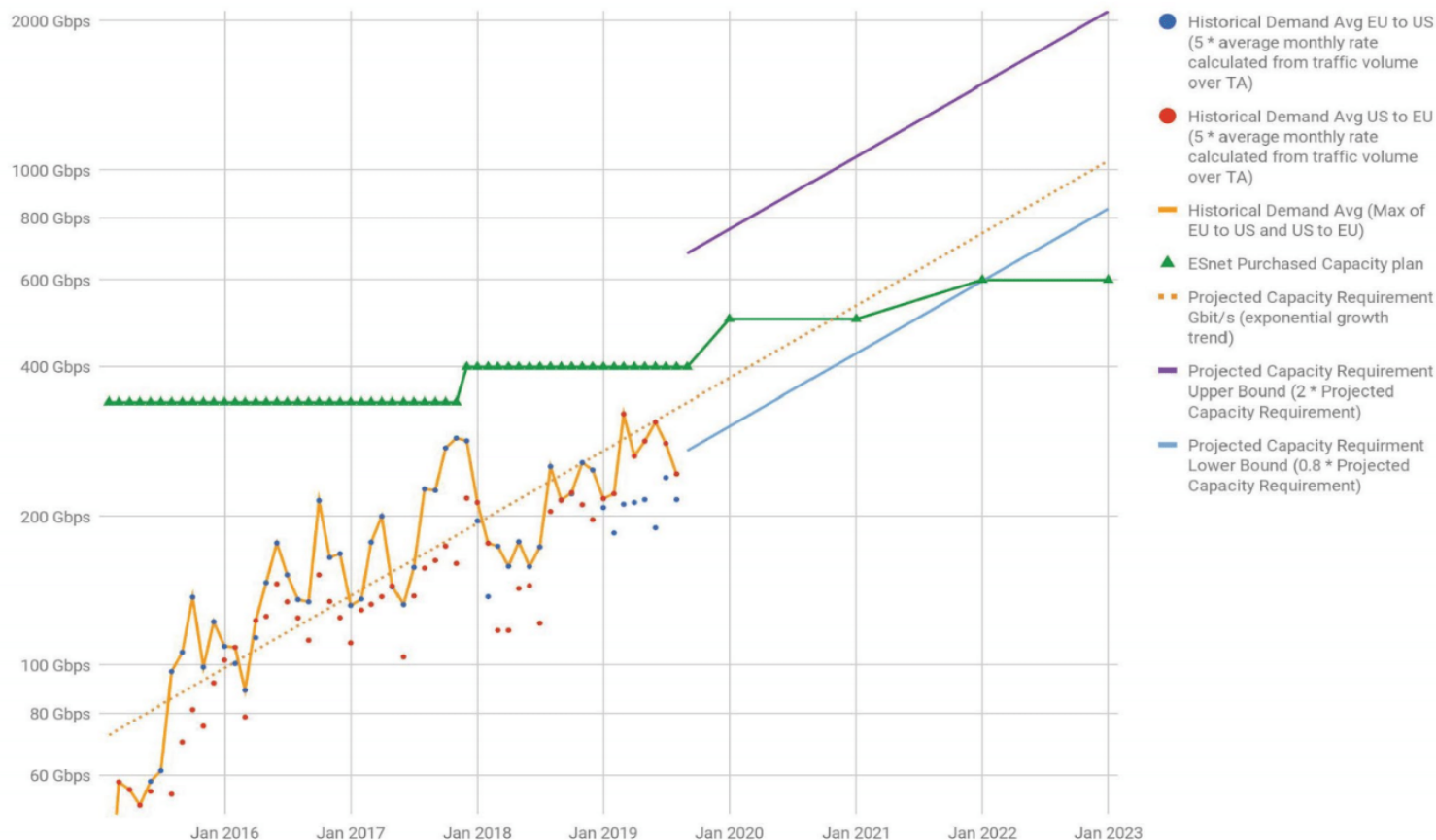**Primary Data volume
per year goes up by x30**

**This talk is about R&D strategies to keep the cost the ~same
despite a x30 increase in data volume per year.**

Will motivate the R&D via a detour on how science is done.

***Conclusion:*** *CMS Data ~Exabyte/Year by ~2028 at HL-LHC*

*Preparation in progress: a 10 Petbyte Data Challenge*

9

European Demand and Capacity Forecasts (updated Sept 2019)



- ● Historical Demand Avg EU to US (5 * average monthly rate calculated from traffic volume over TA)
- ● Historical Demand Avg US to EU (5 * average monthly rate calculated from traffic volume over TA)
- ▬ Historical Demand Avg (Max of EU to US and US to EU)
- ▲ ESnet Purchased Capacity plan
- ▪▪ Projected Capacity Requirement Gbit/s (exponential growth trend)
- ▬ Projected Capacity Requirement Upper Bound (2 * Projected Capacity Requirement)
- ▬ Projected Capacity Requirment Lower Bound (0.8 * Projected Capacity Requirement)

- **Recommendation from ESnet6 technical review:**

**ESnet should consider spectrum acquisition as an option for the non-OLS footprint to serve the science community that depends upon capacity growth of this connectivity.**

**᛭ Fermilab**

9    3/17/2020    Phil DeMar | US-CMS BluePrint WG meeting

https://www.dropbox.com/s/yi9b1gc8v5q8jke/DeMar-US-CMS-BluePrint_3-17-20.pdf?dl=0

# Network Requirements Update for the HL-LHC Era
## LHC Experiments Awaken

* **In January, at the 43ʳᵈ LHCOPN/LHCONE meeting at CERN https://indico.cern.ch/event/828520/, the LHC experiments expressed the need for Terabit/sec links by the start of HL-LHC operations in 2027-28, preceded by the usual Computing and Storage (and Network) challenges starting during LHC Run3 (2021-4)**

* **This was reinforced by the requirements presented by the DOMA project which *"foresees requiring 1 Tbps links by HL-LHC (ballpark) to support WLCG needs. This is for the network backbones and larger sites…"***

  * **References: (1) E. Martelli, S. McKee LHCOPN-LHCONE Report to the Grid Deployment Board, (2) DOMA project presentation at the LHCONE meeting https://indico.cern.ch/event/828520/contributions/3570904/attachments/1968554/3274036/LHCONE-DOMA-01-2020.pdf**

* **NB: The quoted network capacity requirements are an order of magnitude greater than what is currently available through the national and transoceanic networks based on 100GE links.**

  * **As discussed at the LHCONE meeting, in the GNA-G Leadership group meetings, and in the HEPIX Techwatch technology tracking group, these requirements cannot be accommodated solely through the exploitation of technology evolution within a constant budget.**

# Capacity Requirements Analysis, Using
## ESnet Transatlantic Network Traffic Projections

- **Current Requirements: 0.35 – 0.85 Tbps**
  **[0.8 to 2X the 2016-19 traffic projection]**

- **Growth Rate 1.4X per year**

- **Hence *16X capacity requirement in 2028* = 5.6 to 13.6 Tbps;**
  **Since this is an Esnet only, and not a global projection,**
  **the upper limit may be the better requirements metric**

- **Traditional long-term capacity per unit cost rate: +15 – 20% per year;**
  **Hence 3.1 to 4.3 times affordable capacity by 2028**

- **Implied Shortfall: 3.7 to 5.2X**

- **Naïve Implementation Outlook by 2028: 68 200G links across the Atlantic**
  **(for example 17 links on each along 4 disjoint paths);**
  **compare the ANA consortium today: 9 100G links at present**

- **Ways to bring down the costs: Acquire spectrum IRUs on undersea cables;**
  **Move towards co-ownership on undersea cables if and where possible**

- **Outlook: This will get us part of the way there (within a factor of 2?)**

- **Bottom Line: Need to develop a new system that comprehensively monitors,**
  **tracks, manages and controls use, coordinated with compute and storage use**

# Hierarchical Storage via Data Lakes
## Regional Caches

- **Store most data on "active archive" on inexpensive,** high latency media **(e.g. Tape).**

- **Keep a "golden copy" on redundant high availability disk [fewer copies].**

  - **This defines the working set allowed to be accessed.**

  - **Jobs requesting data not in working set will queue up** until data is recalled from archive

- **Regional Caches at processing centers (e.g. Tier1s & 2s; ~1 petabyte)**

  - **Size of region determined by** latency tolerance of application

  - **Cost trade-off: between cache size vs network use**

- **Useful distance metric: 10% IO penalty among merged caches**

- **EU example: ~500-1000 km**

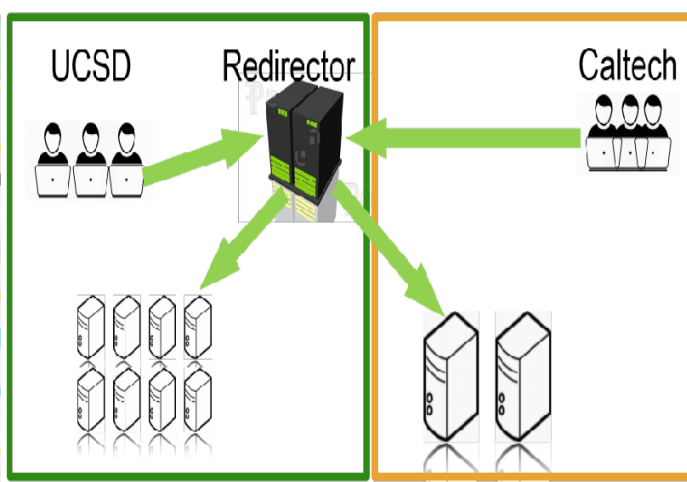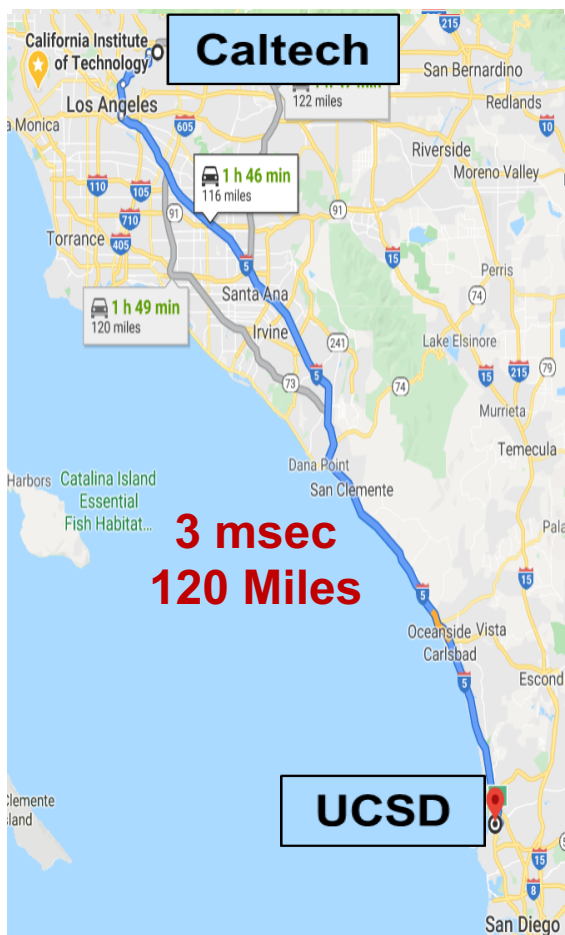- **Advanced protocol, caching methods: could extend distance**



500 Miles is an interesting distance for merging caches !!!

**Examples in Production:**
**"SoCal" (UCSD + Caltech); INFN**

**F. Wuerthwein (UCSD) et al**

# (Southern) California ((So)Cal) Cache

**(Roughly 20,000 cores across Caltech & UCSD … half typically used for analysis)**
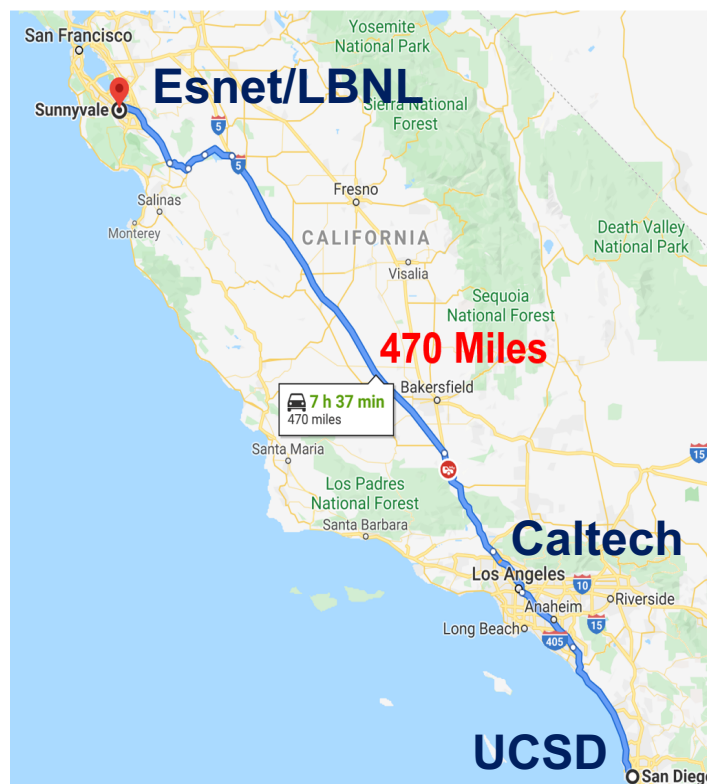


**3 msec
120 Miles**

**CPU in both places can
access storage in both places.**

**How much disk space is enough?**

**Cache MINI and measure
working set accessed:**

**0.45 Petabytes in October 2019**

**470 Miles**

**In early May, we added a cache
at the ESNet POP in Sunnyvale
to the SoCal cache.**

- **Top Line Message**

   **A comprehensive R&D program to develop the architecture, design, prototyping, scaling and optimization of the HL-LHC Computing Model is required**

   * **Including the worldwide network as a first class resource** coordinated with distributed computing and storage

   * **Including innovative approaches** in several areas

   * **Leveraging, coordinating and pushing forward** several key developments: from ML methods to computing to regional caches to SDN networks

   * **Integrating or mediating among regional developments** to form a **worldwide fabric** supporting HEP workflow

* **A leading activity to help drive these developments is the emerging** High Luminosity LHC Computing Model

# SDN Enabled Networks for Science at the Exascale
## SENSE: https://arxiv.org/abs/2004.05953

**Model-based Site and Network Resource Managers**

**Designed to Adapt to Available SDN Systems**

**SENSE Native RMs are Available if no current automation layer**

**Application Workflow Agents**

**SENSE**

**SENSE operates between the SDN Layer controlling the individual networks/end-sites, and science workflow agents/middleware**

**Intent-Based APIS with Resource Discovery, Negotiation, Service Lifecycle Monitoring/Troubleshooting**
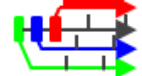
**SDN Layer**

Regional — WAN — WAN — SDX — Regional
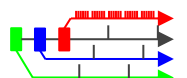
End Site
SDMZ
Instruments   Storage   Compute   DTNs
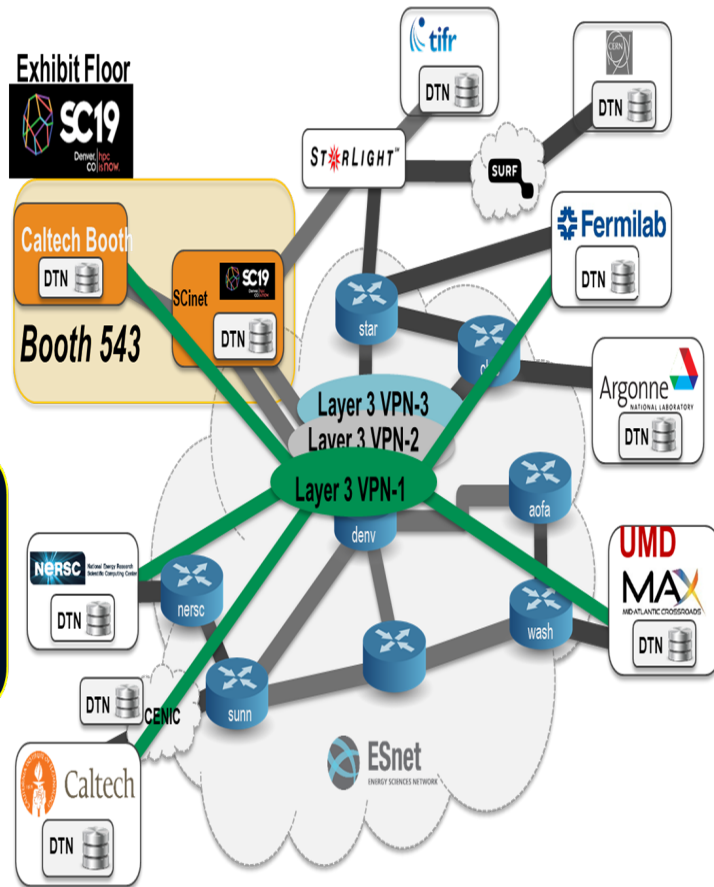
SDMZ   End Site
DTNs   Compute   Storage   Instruments
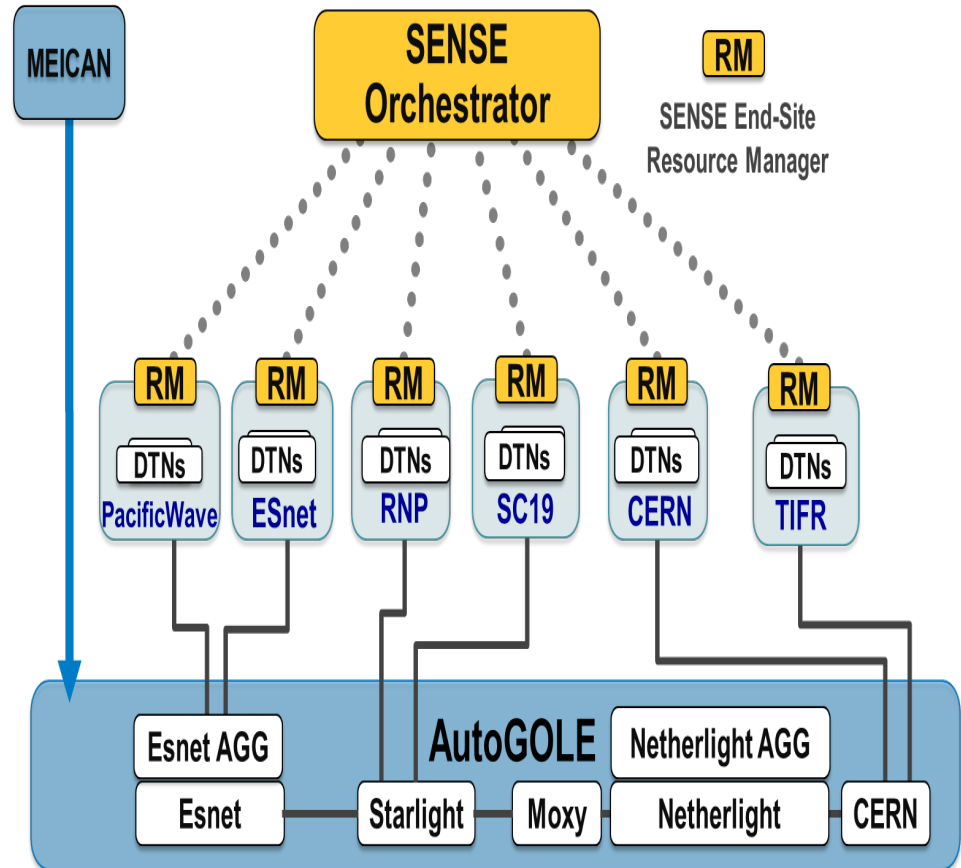
# SENSE SC19 Demonstration Topology

## SENSE Testbed and L3 VPN Service

**SENSE enabled resources at DOE Labs, Universities, Research Facilities, + SC19**

**Dynamic attachment of End Site resources to L3VPNs advertised by ESnet**



# SC19-NRE-020 Intercontinental Demonstration
## Multi-Resource Orchestration via AutoGole and SENSE



**MEICAN**

**SENSE Orchestrator**

**RM** — SENSE End-Site Resource Manager

**RM** DTNs PacificWave | **RM** DTNs ESnet | **RM** DTNs RNP | **RM** DTNs SC19 | **RM** DTNs CERN | **RM** DTNs TIFR

**AutoGOLE**

Esnet AGG | Netherlight AGG

Esnet | Starlight | Moxy | Netherlight | CERN

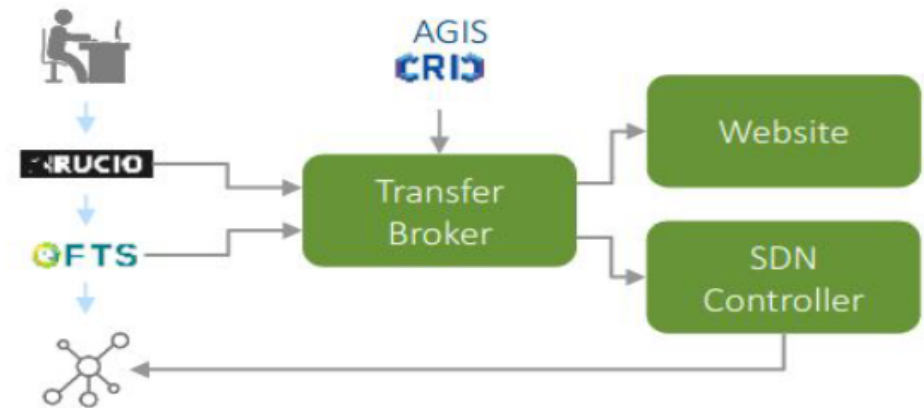**SENSE – AutoGOLE** Joint Interworking Demo ➔ **Candidate Inter-regional Mediation Layer for Global Workflows** (as discussed in GNA-G)

**For a global fabric, including Australia and Africa we would need to include genomics, LSST, SKA, and others in the overall concept along with HEP**
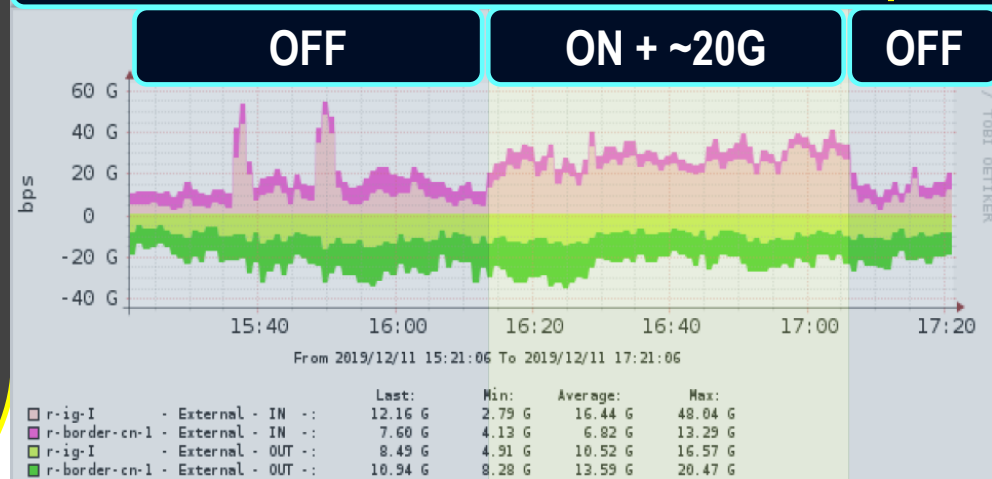
# NOTED: Network Optimized Transfer of Experimental Data CERN/IT Project (C. Busse-Grawitz)

- **NOTED publishes network aware information on on-going massive data transfers,** that can be used to provide additional capacity by orchestrating the network behavior **(e.g. more effective use of existing network paths; finding alternates; load balancing).**

- **The advantage of starting with NOTED is that its Transfer Broker, as shown, can already interpret Rucio and FTS queues** and translate them into network aware information with the help of the WLCG's database.

- **While still in the prototyping stage, NOTED has already demonstrated the full chain** with transfers between CERN and the Tier1s in Germany (DE-KIT) and the Netherlands (NLT1).

## Transfer Broker Interfaces to Job Queues, SDN Controller, WLCG Database



## Switch some traffic to DE-KIT LHCOPN path

| OFF | ON + ~20G | OFF |



From 2019/12/11 15:21:06 To 2019/12/11 17:21:06

| | | Last: | Min: | Average: | Max: |
|---|---|---|---|---|---|
| □ r-ig-I | - External - IN -: | 12.16 G | 2.79 G | 16.44 G | 48.04 G |
| ■ r-border-cn-1 | - External - IN -: | 7.60 G | 4.13 G | 6.82 G | 13.29 G |
| □ r-ig-I | - External - OUT -: | 8.49 G | 4.91 G | 10.52 G | 16.57 G |
| ■ r-border-cn-1 | - External - OUT -: | 10.94 G | 8.28 G | 13.59 G | 20.47 G |

# Missions and Working Groups

**Work areas**
- **Architecture**
- **Tools**
- **Research Support**
- **Operations**
- **Technical Policy**
- **Futures**

**Active work groups**
- **AutoGOLE / SENSE** (Gerben van Malenstein)
- **GREN Mapping** (Ryan & Thomas Fryer)
- **Virtualization** (Jerry Sobieski)
- **GXP Architecture and Services** (Mian Usman)
- **Telemetry (Dale Finkelson)**

**Ideas for WG**
- NSI interoperability & operations coordination
- Intradomain Automation
- GRP
- Routing Anomalies

Erik-Jan Bos at Internet2 Tech-X 12/19

# Missions and Working Groups

- **Establish a plan for operation and evolution of the infrastructure & services; present & next generation roadmap**
  - **WG: Management and Forward Planning (Leadership Team)**
  - **Activities:**
    - **Mission Goals and Program**
    - **Policy/AUP**
    - **Interorganization interfaces: R&E Networks; Science Programs**
    - **Cost & funding evolution**
    - **Technology and market tracking**
    - **Architecture evolution & choices**
    - **Service requirements & choices**
    - **R&D directions and Goals**
- **Define, operate, evolve and extend a global production infrastructure of GXPs and shared links**
  - **WGs: Architecture, Operations, GREN Mapping**

# Missions and Working Groups

- **Establish reliable high throughput across the global footprint**

    - **WGs: Tools, Telemetry, Engagement, GRP, AutoGOLE/SENSE**

- **Define and evolve an advanced services architecture (ASA); Establish and Grow a Global Persistent Testbed; Establish Major Use Cases (LHC, LSST et al)**

    - **WGs: GRP, Architecture, Telemetry, Engagement, AutoGOLE/SENSE, Virtualization, …**

    - **Activities: Deployment, Prototyping, Workflow Integration**

    - **Trajectory: Towards preproduction then limited scale production trials**

- **To discuss, develop: Working group structures, groupings and/or taskforces to refine and carry out the mission**

- *To meet the challenges of globally distributed Exascale data and computation faced by the major science programs*
- *Mission: Coordinate provisioning the maximum feasible capacity across a global footprint, and enable best use of that infrastructure*
- *Beyond capacity alone, enable the science. An Approach:*
  - *A new "Consistent Operations" paradigm:* **goal-oriented, policy-driven**
  - *Stable, resilient high throughput flows*
  - *Controls at the network edges, and in the core*
  - *Real-time dynamic, adaptive operations among the sites and networks; Increasing negotiation, adaptation, with built-in intelligence*
  - *Coordination among VO and Network Orchestrators*
- **Bringing Exascale, pre-Exascale HPC and Web-scale cloud facilities, into the data intensive ecosystems of global science programs**
  - **Petabyte transactions and caching using state of the art + emerging network and server technology generations; Tbit/sec demonstrators**
- **Engaging with the full range of technologies and many partners**
- **We require a *comprehensive, forward looking* global R&D program**
- *The GNA-G can have a key role in this important endeavor*

22

# Extra Slides

# Follow

# Global Network Advancement Group (GNA-G) Leadership Team: Since September 2019

leadershipteam@lists.gna-g.net



**Erik-Jan Bos NorduNet**    **Buseung Cho KISTI**    **Dale Finkelson Internet2**    **Gerben van Malenstein SURFnet**    **Harvey Newman Caltech**    **David Wilde Aarnet**

- **The GNA-G is an open volunteer group devoted to developing the blueprint** to make using the Global R&E networks both simpler and more effective, operating under GNA-G.

- **Its primary mission is to support global research and education** using the technology, infrastructures and investments of its participants.

- **The GNA-G needs to be a data intensive research & science** engager that facilitates and accelerates global-scale projects by (1) **enabling high-performance data transfer,** and (2) **acting as a partner in the development of next generation intelligent network systems** that support the workflow of data intensive programs

See https://www.dropbox.com/s/qsh2vn00f6n247a/GNA-G%20Meeting%20slides%20-%20TechEX19%20v0.8.pptx?dl=0

# GNA-G Mission

- **The primary mission of the Global Network Advancement Group (GNA-G) is to support global research and education using the technology, infrastructures and investments of its participants.**

- **The GNA-G exists to bring together researchers, National Research and Education Networks (NRENs), Global eXchange Point (GXP) operators, Regionals, and other R&E Network providers in developing a common global infrastructure to support of these needs.**

# GNA-G Mission (2)

- **The GNA-G exists to further the understanding about how to effectively use the rich networking infrastructure the NREN community has created:**
  - **The GNA-G will suggest best practices, efficient ways for interaction among NREN's, be a means to encourage NREN's to try new techniques and technologies.**
  - **In doing these things the GNA-G will help to make the investment the NRENs have made in building networks have a greater value to the community those investments are meant to serve, the Education and Scientific/Research communities.**

# GNA-G Challenges (1/3)

- **Cooperation between NRENs to create resilient infrastructure in all regions**
  - We want everyone to make packets move fast globally
  - Truly fast end to end transfers maximizing line speed capacity

- **Although bandwidth keeps increasing, in the long term there is scarce capacity on the array of intercontinental circuits, in specific trans-oceanic links**
  - How to effectively use this capacity for research?
  - Intelligent inter-operation

- **The GNA-G has a role in promoting collaboration within regions** to foster the growth of resilient regional infrastructure connectivity as well as **collaboration on region to region interconnectivity**

- **The GNA-G should also be a forum for the discussion and trialing of new technologies** aimed at furthering research

- **The GNA-G could also be a forum for discussion about mission- oriented networking,** e.g. Service for data intensive science and remote campuses

- **The GNA-G needs to be a data intensive research & science engager that facilitates and accelerates global-scale projects by enabling end-to-end high-performance data transfer as well as a partner in the development of next generation intelligent network systems that support the workflow of data intensive programs**

- **We are mindful not to duplicate any efforts in areas we work in, and we work with other groups to find optimal synergies**

- **Areas that naturally fall upon us:**
  - **Bringing regional efforts together on intelligent networks**
  - **What happens across oceans**

* The LHCONE/LHCOPN meetings at CERN are focused on bringing the WLCG experiments, NRENs and researchers together to discuss network requirements, and possible future work https://indico.cern.ch/event/828520/
  * The January meeting was summarized at the Grid Deployment Board (GDB: https://indico.cern.ch/event/813757/
* There have been significant activities in the networking R&E community and the WLCG experiments, but not yet focused on targeting joint production goals
* A survey of existing efforts and possible directions was provided by the **HEPiX Network Function Virtualization WG**
  * Their report (April 2020) is available at https://zenodo.org/record/3741402#.XsFU1MBID8A

* The NFV WG identified three work areas of interest to the NRENs and experiments:
  * Marking R&E network packets
  * Shaping and pacing traffic
  * Network Orchestration
* The consensus at the meeting resulted in an effort to instantiate a technical networking working group to help define, prototype and guide the work outlined above.
* The evolving draft charter is available at

https://docs.google.com/document/d/1I4U5dpH556kCnolHzyRpBl74lPc0gpgAG3VPUp98lo0/edit?usp=sharing

30

# Towards a Global Network Fabric

✳ **Following discussions and presentations at the Americas and Global Research Platform workshops in September and the Internet2 Tech Exchange in December, it became clear that the services in the SENSE project could be further developed to serve as a mediator among the intelligent network software systems being developed in the various world regions including Europe (AutoGOLE), Latin America (AmLight), and Asia (Virtual Dedicated Networks). A living example of this was demonstrated at SC19 [*] where interoperation of the SENSE and AutoGOLE network service frameworks, and integral control of the DTN systems and network systems by the SENSE Resource Managers developed by Caltech and ESnet were shown.**

✳ **This led to plans for a persistent national and global R&D testbed as a venue for ongoing and future network developments in the context of the HL-LHC Computing Model. These developments are also planned to leverage NSF's major investment in FABRIC, "a unique national research infrastructure to enable cutting-edge and exploratory research at-scale in networking, cybersecurity, distributed computing and storage systems, machine learning, and science applications".**

✳ **[*] SC19 Network Research Exhibition: "LHC Multi-Resource, Multi-Domain Orchestration via AutoGOLE and SENSE", https://sc19.supercomputing.org/app/uploads/2019/11/SC19-NRE-020.pdf**

* It was agreed in subsequent discussions that the HEPIX Technology Watch WG and/or the Global Network Advancement (GNA-G) leadership group that was formed in the fall of 2019 [*], can help define how much of it can be satisfied through technology evolution by 2027, and by 2024 in the preparatory phase [Evolution to 400G links; nearing technology limits across oceans]

* The rest will involve a change in paradigm including a system composed of end-to-end services involving coordinated operation among sites and networks, and orchestration:

  * We can leverage developments underway in projects such as SENSE, NOTED and SANDIE.

  * Ongoing discussions should continue to define what the new services and classes required entail.

    * Solutions will vary by region and by network.

* An important part of this is the persistent testbed being deployed by SENSE in collaboration with AutoGOLE and other projects.

  * This is proceeding: starting with the current SENSE testbed sites, plus extensions to CERN, Starlight in Chicago, SURFnet in Amsterdam, UCSD, and several other sites in the US, Europe, Latin America and Asia