# Notes from LHCONE Point-to-Point session at the Helsinki LHCONE meeting

Oct. 11, 2016
W. E. Johnston, wej@es.net
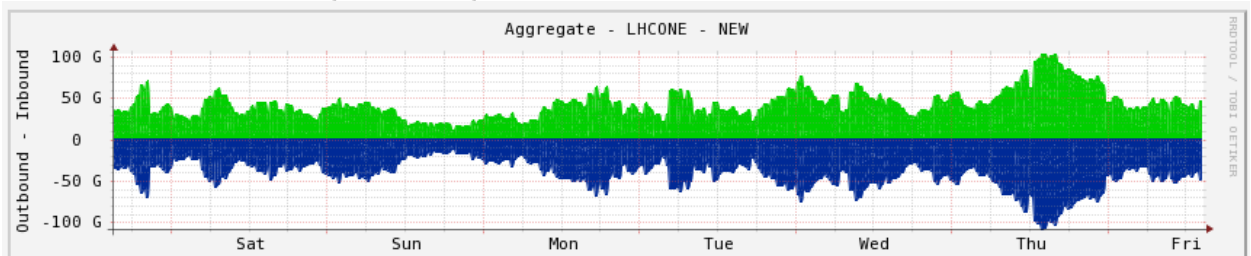
Meeting web site with all talks: https://indico.cern.ch/event/527372/

**Please take note**: The comments that I make here about the talks are NOT complete summaries of the talks. I have abstracted what I consider to be highlights or of special interest. In almost all cases much has been omitted and I recommend that you look at the talks on the Indico site.

## 1 First a few comments about the LHCONE overlay network.

*Mian Usman, GEANT*
- Traffic within LHCONE is steadily growing
- GÉANT has seen peaks of over 100Gbps
- Growth of over 65% from Q2 2015 to Q2 2016



- **https://indico.cern.ch/event/527372/contributions/2158737/attachments/1338766/2016155/LHCONE_L3VPN_Update.pdf**

*Mike O'Connor, ESnet*
- Two US Tier1 LHC centers (ATLAS and CMS) are regularly exceeding 50Gbps each in LHCONE
- ESnet LHCONE traffic to and from CERN regularly exceeds 40 Gb/s
- ESnet LHCONE traffic has increased 118% in the past year

**https://indico.cern.ch/event/527372/contributions/2236890/attachments/1338767/2015141/LHCONE-OperationsUpdate-Helsinki.pdf**

The global infrastructure of LHCONE has been a resounding success. [wej]

## 2 Point-to-Point session

The term "Point-to-Point" session is a bit misleading.

At the 2011 Amsterdam LHCONE meeting, the Architecture Working Group was charged to investigate five areas:
1) VRF (L3 VPN)
2) L2 multipath using 802.1.aq or TRILL
3) Openflow

4) Point to point circuit pilot
5) Diagnostic infrastructure
6) Determine LHCONE impact on / interaction with LHC software stacks

The results have been:
- The LHCONE overlay (L3 VPN) network is now in full production and some time ago all related issues have been moved to the LHCONE Operations Group.

- L2 multipath investigation is quiescent [as far as I know – wej]

- Openflow (now SDN) is rapidly moving to center stage for next generation networks.

- The point-to-point pilot is making slow progress

- Diagnostic infrastructure has mostly moved to its own working group.
  (See Shawn Mckee's talk
  https://indico.cern.ch/event/527372/contributions/2210680/attachments/1338706/2016165/LHCONE_perfSONAR_update-Helsinki-2016.pdf )

- Software stack interaction is an ongoing exercise.

The P2P experiment has been the primary LHCONE Architecture Group focus for the past several years. However, I think that there ia increasing interest in SDN networking and how it might be used in LHCONE.

While there have not been any LHCONE SDN experiments proposed yet, SDN and SDX (SDN Exchanges) are being worked on in the GLIF community and at some point this approach will likely provide a more flexible, lower overhead way to configure overlay networks like LHCONE.

Further, as the P2P use cases frequently include many P2P circuits between well identified groups of sites, it may be that an SDN configured overlay network with performance guarantees will supplant P2P circuits. (However, I am sure the same multi-domain authorization and authentication issues that have shown up in the P2P work will, show up in SDN.)

In this regard we can expect LHCONE Architecture Group will be involved with SDN experiments, especially as the R&E networks start to deploy SDN in the WAN.


## 2.1   LHC Networking And SDN/SDX Services

*Joe Mambretti, iCARI and StarLight*

SDN allows network designers to create a wider range of services than those provided by traditional networks.

The case for SDN:
- A network technology that can both take advantage of and enable modern IT technology

- High level of routine customization of the network – including virtualization -  that is just a regular part of the normal operational milieu

Key attributes for SDXs = Open services, architecture, connectivity

SDN also enables:
- Highly granulated views into network capabilities and resources, including individual flow data
- Many options for control over SDN resources, including distributed control of managed network services by edge applications
- Dynamic provisioning and adjustment options, including those that are automated and implemented in real time
- Faster implementation of new services

SDX – Software Defined Network Exchanges
- SDN technology is local to a domain, so SDXs are required to interconnect SDN "islands"
- SDXs provide highly granulated views into and control over all flows within the exchange
- Democratization of exchange facilities – options for edge control over exchanges

SDX architectural components (examples)
- Hybrid network services (multiple services, multi-layer, multi-domain, integration of Open Flow and non-OF paths)
- Multi-domain resource advertisement, discovery, and signaling (including edge signaling)
- Support for multi-domain, federated path controllers of different types
- Topology exchange services
- Control and network resource APIs
- Network programming languages (e.g. P4 and Frenetic)
  - E.g. see "high-level language for programming protocol-independent packet processors" - http://onrc.stanford.edu/p4.html
- Abstraction definitions
- Service signaling and policy bundling and distribution
- BGP extension and substitutes
- Network Description Language (NDL) schemas
  - See, e.g. https://ivi.fnwi.uva.nl/sne/ndl/
- Orchestration processes
  - E.g. Using orchestrated SDX services to implement and control WAN "superchannels," in part enabled by DTNs, which demonstrates highly scalable dynamic provisioning – a scalability not possible on today's networks
- Etc.

**https://indico.cern.ch/event/527372/contributions/2159195/attachments/1339843/2017213/IRNC_SDX_LHCONE-LHCOPN_Presentation_September_2016.pdf**

## 2.2   NSI and Automated GOLE update

*John MaAuley, ESnet*

**NSI-Connection Service 2.0 published**

- The Modify operations are added.
- Reservations are now versioned.
- Added a generic notification message that carries a number of new notifications for activation. Folded forcedEnd in to this as well.
- Updated PathList and STP definitions to follow approved proposal.
  - Changed path to ERO (explicit routing object).
  - Included new optional "symmetric" element to PathType to constrain bidirectional connections.
  - STP is not composed of networkId, localId, an optional label, and an optional orientation to indicated ingress or egress of bidirectional connections.
- Removed "nilable" from globalReservationId and made it optional to follow the design pattern of all other types.
- Collapsed ServiceParametersType - moved schedule and serviceAttributes into ReservationInfoGroup instead of having them in a subtype.

Public Comment
- NSI Signaling and path finding v2
- NSI Policy (completed)
- NSI NSA description (completed)
- NSI AA (completed)
- NSI Topology (in process)

**AutoGOLE fabric delivers dynamic network services between GOLEs and networks:**
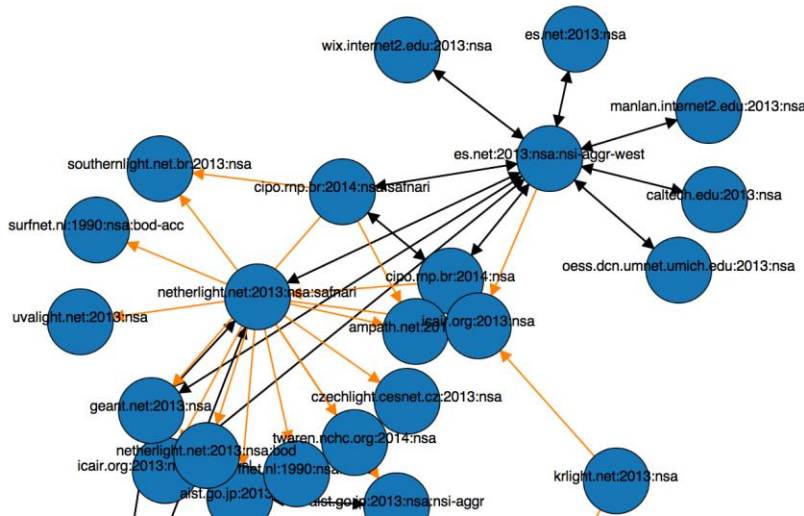


- Based on NSI Connection Service v2.0
  - Redundant Aggregator backbone with a leaf uPA per network (hub and spoke architecture)

- o 29 Network Service Agents (6 aggregators, 23 uPA) advertising 30 networks
- Using DDS service for NSA discovery and document propagation between aggregators
- Introduction of monitoring, troubleshooting, and provisioning tools
  - o Dashboard, MEICAN, DDS Portal, etc.
- Advancing capabilities
  - o Experimenting with new path finding and signaling algorithms
  - o Additional network modeling for optimizations

About eight projects currently use AutoGOLE.

## AutoGOLE dashboard (e.g. NSI control plane peerings):

The control plane graph shows the NSI control plane peerings. On the graph it is possible to see control plane peering mismatches, NSA host reachability and Unknown NSAs. Alive NSA hosts marked as unreachable might need to allow ICMP traffic. More information for each NSA can be seen by clicking on a node and by looking at the tables below.
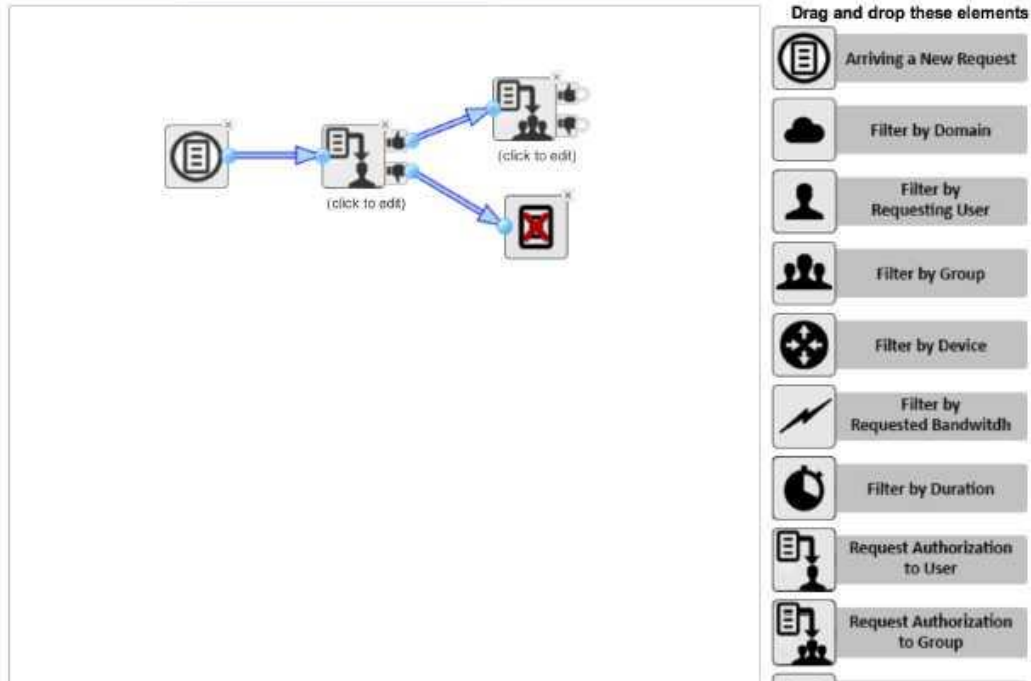


## Management Environment of Inter-domain Circuits for Advanced Networks (MEICAN) (RPN)

- Topology view
- Monitoring view
- Workflows

- Discovery
- Device inventory
- Port inventory
- Automated tests

**Work items, 2016**

- AutoGOLE Dashboard
    - Overview of both control plane and data plane of the AutoGOLE
- Supporting LHC Sites
    - Supporting LHC sites that want to connect to the AutoGOLE (Brookhaven and NLT1 tested last year)
- Connecting Data Transfer Nodes
    - Kick-off by StarLight, Caltech, RNP, University of Amsterdam this fall
- AutoGOLE MEICAN Pilot
    - Run a pilot of the RNP's Cipó Service front-end interface – the MEICAN – being used by participant research and education networks (RENs) and exchange point (IXP) operators to configure multi-domain point to point circuits. The participants will evaluate the MEICAN as the main interface for AutoGOLE GLIF Project.
    - https://wiki.rnp.br/display/secipo/AutoGOLE+MEICAN+Pilot

Under discussion

- Using the AutoGOLE for automated interconnects with service providers:
    - GÉANT-Microsoft Azure ExpressRoute connections are now being setup using GÉANT and NetherLight's automated provisioning systems to prevent manual configuration for each service request.

https://indico.cern.ch/event/527372/contributions/2288042/attachments/1339832/2017195/LHCONE.NSI.AGOLE.update.pdf

## *2.3* High Performance Data Transfer Node - Design to Reality
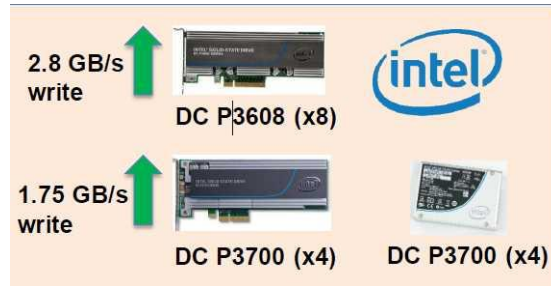
*Azher Mughal, Caltech*

**A summary of the Caltech SC15 demo and the SC16 goals.**
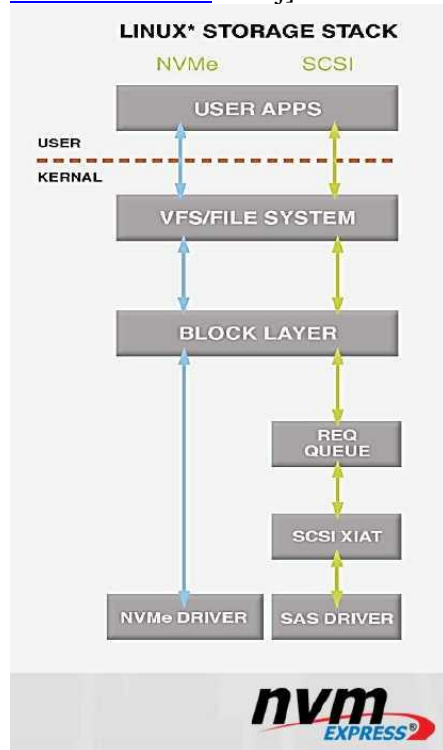
SC16:
- SDN Traffic Flows
    - Network should solely be controlled by the SDN application.
    - Relying mostly on the North Bound interface
        - Install flows among a pair of DTN nodes
        - Re-engineer flows crossing alternate routes across the ring (shortest or with more bandwidth)
- High Speed DTN Transfers
    - 100G to 100G (Network to Network)
    - 100G to 100G (Disk to Disk)
    - 1 Tbps to 1Tbps (Network to Network using RoCE and TCP)
- ODL (OpenDaylight) (PhEDEx + ALTO (RSA + SPCE)
    - Substantially extended OpenDaylight controller using a unified multilevel control plane programming framework to drive the new network paradigm
    - Advanced integration functions with the data management applications of the CMS experiment

**DTN Design Considerations**

- Different Choices and Opinions …
    - How many rack units are needed / available.
    - Single socket vs dual socket systems
    - Many cores vs fewer cores at high clock rates
    - SATA 3 RAID Controllers vs HBA Adapters vs NVME
    - White box servers vs servers from traditional vendors (design flexibility ?)
    - How many PCIe slots are needed (I/O + network). What should be the slot width (x16 for 100GE)
    - Onboard network cards vs add-on cards
    - Airflow for heat load inside the chassis for different workloads (enough fans ?)
    - Processor upgradeable motherboard
    - Remote BMC / ILOM / IPMI connectivity
    - BIOS Tweaking
    - Choice of Operating system, kernel flavors
    - Redundant power supply

- Results of SC15 hardware testing are shown
- NVME (Non-Volatile Memory Express - communications interface/protocol developed specially for SSDs) advantages
    - NVMExpress introduced in 2011
    - Bypasses all the AHCI / SCSI layers
    - Extremely fast (dedicated FPGA) (Seagate recently announced 10GB/sec drive

- o
- o Low latency
- o Supported by large number of vendors
- o Generally available in both 2.5" or PCIe cards form factor (PCIe Gen3 x4/x8/x16)
- o Prices are getting low:
  - - Sata3 SSDs are about 24 - 40 cents per GB
  - - NVME are about $2 per GB (expensive)
- o [NVME intro: http://www.pcworld.com/article/2899351/everything-you-need-to-know-about-nvme.html  - wej]
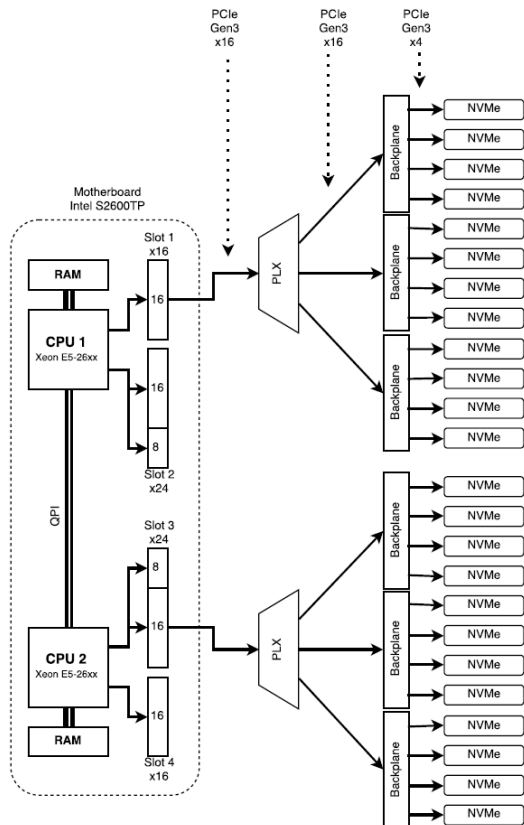


- o

## NVME over fabrics

- Goals of NVMe over Fabrics is to extend the low-latency efficient NVMe block storage protocol with no more additional 10uSec.
- NVMe over Fabrics maintains the architecture and software consistency of the NVMe protocol across different fabric types, providing the benefits of NVMe regardless of the fabric type or the type of non-volatile memory used in the storage target.
- Initial code is available starting with Linux kernel 4.7
- Current NVME transports:
  - o Fiber Channel
  - o NVME fabric is built over exiting
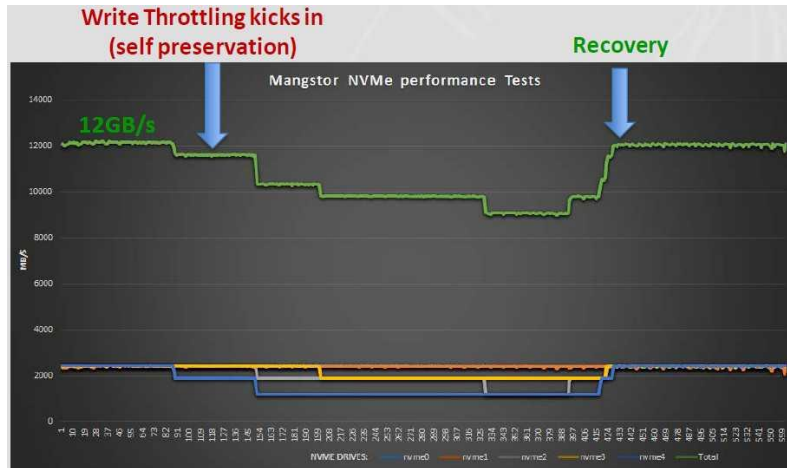
## 2CRSI / SuperMicro 2U -NVME Servers

- PCIe Switching Chipset for NVME



- 
- 2CRSI Server with 24 NVMe drives



- 
- Max throughput reached at 14 drives (7 drives per processor)
- A limitation due to combination of single PCIe x16 bus (128Gbps), processor utilization and application overheads.
- Temperature effects on SSD drives

- 
- Lesson: Follow the manufacturer's guidelines for the minimum airflow requirement.
- Interesting notes on "Build a low cost NVME storage"


## System Design Considerations for200GE / 400GE and beyond … 1Tbps

Server Readiness:
1) Current PCIe Bus limitations
- PCIeGen 3.0(x16can reach 128GbsFull Duplex)
- PCIe Gen 4.0(x16can reach double the capacity, i.e. 256Gbps
- PCIeGen 4.0(x32can reach double the capacity, i.e. 512Gbps
2) Increased number of PCIe lanes within processor
- Haswell/Broadwell (2015/2016)
  o PCIelanes per processor = 40
  o Supports PCIeGen 3.0 (8GT/sec) [8 giga transfers/sec]
  o Up to DDR4 2400MHz memory
- Skylake (2017)
  o PCIelanes per processor = 48
  o Supports PCIe Gen 4.0 (16GT/sec)
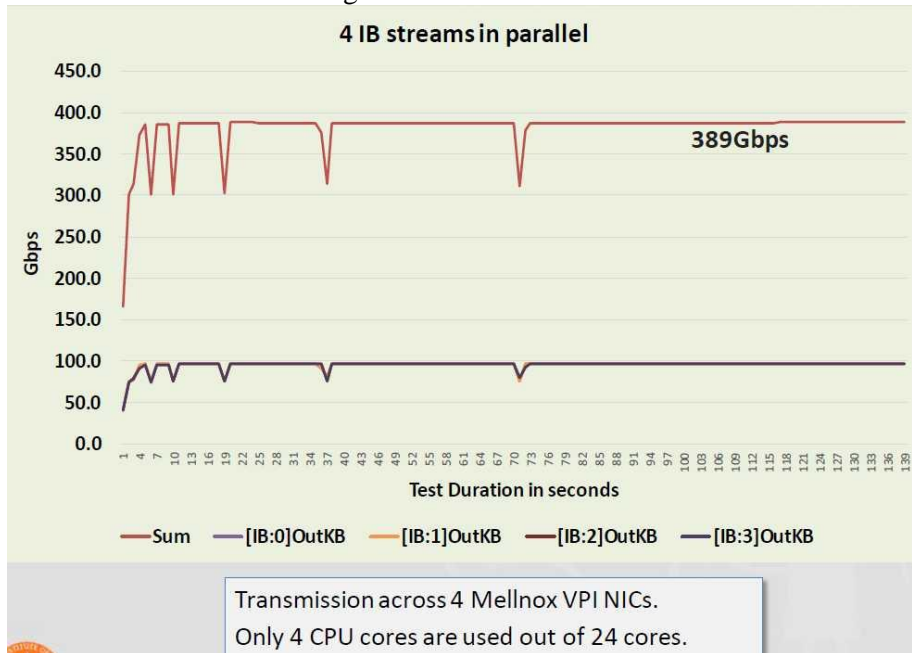3) Faster core rates, or Over clocking (what's best for production systems)
4) Increased memory controllers at higher clock rate reaching 3000MHz
5) TCP / UDP / RDMA over Ethernet


## SC16 –1 Tbps network server demo

- The server design and the network details are discussed in the paper
- 400GE network server testing



Transmission across 4 Mellnox VPI NICs.
Only 4 CPU cores are used out of 24 cores.

- 

https://indico.cern.ch/event/527372/contributions/2267011/attachments/1339980/2017493/Helsinki_DTN_Mughal.pdf

## 2.4 Global Networks for High Energy Physics and Data Intensive Sciences: A New Network Paradigm for LHC and HL LHC, and Exascale HPC Systems in HEP's Data Intensive Ecosystem [NORDUnet conference talk as part of LHCONE meeting]

*Harvey Newman, Caltech*
- Networking for High Energy Physics and Global Science: A 30 year retrospective
  - Network Trends in 2015-16 – 100G WAN networks well under way
  - We are midway in the 7-8 Year generational cycle of 100G networks
  - Issue: Will next generation 400G networks proliferate in time for Run3 ?
- A new era of challenges and opportunity
  - Complex workflows: The flow patterns have increased in scale and complexity, even at

the start of LHC Run2
- o WLCG Dashboard for April-May indicates
    - 28 gbytes/s average,75 gbytes/s peak rates
    - Complex workflows
    - Multi-tbyte dataset transfers
    - Transfers of up to 60 million files daily
    - Access to tens of millions of object collections/day
    - >100k of remote connections simultaneously (e.g. using AAA "Any data, Anytime, Anywhere" - a CMS XRootD federated storage approach that aims to: Provide low-latency access to any single event, Reduce data access error rate, Overflow jobs from busy sites to less busy ones, Use opportunistic resources, Make life at T3s easier. See, e,g, https://indico.egi.eu/indico/event/1417/session/53/contribution/245/material/slides/0.pdf - wej]

- o 2.7x traffic growth (+166%) in last 12 months; +60% in April alone

- Addressing a new era of challenges as we move to exascale data and computing
    - o The largest science datasets under management today, from the LHC program, are >400 petabytes (PB)
        - Exabyte datasets are on the horizon, by the end of Run2 in 2018
        - 850 PB flowed Across the WLCG, 300 PB over Esnet in last 12 months
        - These datasets will grow by ~100X, to the ~50-100 Exabyte range, during the HL LHC era from 2026
    - o Reliance on high performance networks will continue to grow as many Exabytes are distributed, processed & analyzed at 100s of sites
    - o As needs of other fields continue to grow, HEP will face stiff competition for use of limited network resources.
    - o Location independent access: blurring the boundaries among sites + analysis vs computing
        - Once the archival functions are separated from the Tier-1 sites, the functional difference between Tier-1 and Tier-2 sites becomes small [and the analysis/computing-ops boundary blurs]
        - Connections and functions of sites are defined by their capability, including the network!!
    - o Domains of Big Data in 2025. In each, the projected annual and storage needs are presented, across the data lifecycle

| Data Phase | SKA | Twitter | YOU TUBE | GENOMICS | HL LHC |
|---|---|---|---|---|---|
| Acquisition | 25 ZB/Yr | 0.5–15 billion tweets/year | 500–900 million hours/year | 1 Zetta-bases/Yr | |
| Storage | 1.5 EB/Yr | 1–17 PB/year | 1–2 EB/year | 2-40 EB/Yr | 2-10 EB/Yr |
| Analysis | In situ data Reduction | Topic and sentiment mining | Limited requirements | | |
| | Real-time processing | Metadata analysis | | Variant Calling $2 \times 10^{12}$ CPU-h | |
| | Massive Volumes | | | All-pairs genome alignment $10^{16}$CPU-h | 0.065 to 0.2 X $10^{12}$ CPU Hrs |
| Distribution | DAQ 600 TB/s | Small units of distribution | Major component of modern user's bandwidth (10 MB/s) | Many at 10 MBps Fewer at 10 TB/sec | DAQ to 10 TB/s Offline ~0.1 TB/s |

- Towards a next generation network-integrated system: SDN systems for exascale science
  - o Vision: Distributed environments where resources can be deployed flexibly to meet the demands
    - - SDN is a natural path to this vision:
      - • Separating the functions that control the flow of traffic, from the switching infra-structure that forwards the traffic
      - • Through open deeply programmable "controllers".
  - o With many benefits:
    - - Replacing stovepiped vendor HW/SW solutions by open platform-independent software services
    - - Virtualizing services and networks: lowering cost and energy, with greater simplicity
    - - Adding intelligent dynamics to system operations
    - - A major direction of Research networks + Industry
    - - A Sea Change that is still emerging and maturing
  - o Prerequisites: Dynamic circuits
    - - Generalized to: Multicircuit, multisite, SDN driven systems:
      - • In LHCONE
      - • For LSST in the future
  - o SDN Demonstration at the FTW Workshop (Caltech, Amlight/FIU, ESnet, Internet2, Michigan, Sao Paolo)
    - - Dynamic path creation via DYNES, FDT agent, OSCARS, OESS for OpenFlow data plane provisioning over Internet2 AL2S, MonALISA agents at the end-sites provide detailed monitoring information
  - o SC15 demonstration: SDN-driven large flow steering, load balancing, site orchestration, over terabit/sec LANs and over global networks
  - o [Progress is being made in usable SDN technology – wej]
- Bringing the leadership HPC facilities into the data intensive echo systems of the LHC and other major science programs
  - o [considerable progress is being made – see slides -wej]
- **Summary**:
  - o Advanced networks will continue to be a key to the discoveries in HEP and other data intensive fields of science
  - o Near Term and Decadal Challenges must be addressed:  Greater scale, complexity and scope; challenging the available capacity

- New approaches: a new class of deeply programmable software driven networked systems to handle globally distributed Exabyte-scale data are required, and being developed
- NGenIA: New paradigm - Consistent SDN-driven end-to-end ops with stable, load balanced, high throughput managed flows
  - A new horizon in the way networks are operated and managed
- Adapting Exascale Computing Facilities to meet the needs of data intensive science, with high energy physics as the first use case (followed by others) will have multiple benefits
  - Short Term: Enable Rapid Responses, including full reprocessing
  - Medium Term: Paving the Way to the next LHC Computing Model, within the bounds of networking and storage
  - Long Term: Empower the HEP and other communities to make the next rounds of discoveries in science

**https://indico.cern.ch/event/527372/contributions/2270877/attachments/1352263/2041805/NGenIAGlobalNetworks_hbn091916.pdf**

## 2.5 Network Traffic Optimization

*Yatish Kumar, Corsa*

Suggestions for an approach, or rather a set of steps of increasing sophistication to an approach for traffic optimization
- Typical WANs: Multiple paths. Not all are heavily loaded
- Simple Assumption
  - Using this topology and link rate information a system like Panda can compute the offered load onto all links in the network.
  - Ignore other sources of traffic.
  - Benefits
    1) Panda can avoid competition with itself for various data transfers
    2) It can potentially select non competing paths, as long as there is a transfer opportunity to use multiple path segments
- Better Situation - Other Traffic
  - Using topology and link rate information a system like Panda can compute the offered load onto all links in the network.
  - Incorporate path utilization updates to account for other traffic not under Panda's control
  - Benefits
    1) Panda can avoid competition with itself for various data transfers
    2) It can potentially select non competing paths, as long as there is a transfer opportunity to use multiple path segments
    3) Be somewhat adaptive about scheduling transfers, based on network activity\
- Even Better Situation - Other Traffic
  - Using this topology and link rate information a system like Panda can compute the offered load onto all links in the network.
  - Publish scheduled data transfers to a calendar.
  - Benefits
    1) Panda can avoid competition with itself for various data transfers
    2) It can potentially select non competing paths, as long as there is a transfer opportunity to use multiple path segments
    3) Be somewhat adaptive about scheduling transfers, based on network activity

4) Allow other users to cooperatively avoid competition with Panda for network resources
- Further refinement
  - Using this topology and link rate information a system like Panda can compute the offered load onto all links in the network.
  - Incorporate Bandwidth SLAs
  - Benefits
    1) Panda can avoid competition with itself for various data transfers
    2) It can potentially select non competing paths, as long as there is a transfer opportunity to use multiple path segments
    3) Be somewhat adaptive about scheduling transfers, based on network activity
    4) Allow other users to cooperatively avoid competition with Panda for network resources
    5) Create more predictable behaviours when multiple users are involved
- Getting Carried Away
  - Using this topology and link rate information a system like Panda can compute the offered load onto all links in the network.
  - Install Data Staging in quiet parts of the network, or at key junctions
  - Benefits
    1) Panda can avoid competition with itself for various data transfers
    2) It can potentially select non competing paths, as long as there is a transfer opportunity to use multiple path segments
    3) Be somewhat adaptive about scheduling transfers, based on network activity
    4) Allow other users to cooperatively avoid competition with Panda for network resources
    5) Create more predictable behaviours when multiple users are involved
- Summary
  - Easy Steps: Don't ask anything of the service provider. Just document topology.
  - Next Level : Ask for a bandwidth calendar and traffic steering
  - Next Level : Ask for SLAs / bandwidth guarantees

  - All of the above are relatively easy requests possible on existing networks
  - Note: No BOD, no RSVP, no MPLS-TE, no L2 circuit setups

https://indico.cern.ch/event/527372/contributions/2309006/attachments/1339734/2017007/LHCONE_Helsinki_Corsa.pdf


## 2.6  BGP Route Server Proof of Concept [for P2P circuits]

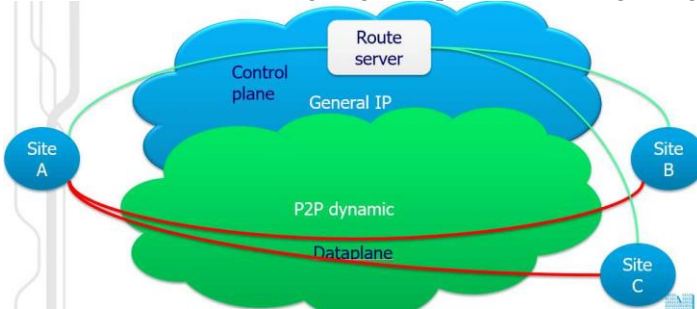*Magnus Bergroth, NORDUnet and Bruno Hoeft , DE-KIT*


**The Problem**

- Dynamic P2P links has two end points that normally terminates in a aggregation router at each site.
- On logical interface per destination site.
- eBGP are configured over the logical interface to each site.
- Reachability is advertised after the P2P link is up and BGP is established.
- Full mesh of BGP sessions.
- Extensive amount of configuration.

- BGP sessions over short lived P2P links are most of the time down and causes alarms.

**Use of a route server**

- Simplify the BGP setup
- Only one BGP session per site
- Route server with one outgoing RIB per site, steering using communities



- 
- Edoardo found:
  - Unfortunately, been quiet, and what I know no implementation ready for use.
  - A new draft: BFD for Multipoint Networks draft-ietf-bfd-multipoint-08 lead by Juniper + Cisco means that maybe this will be eventually move forward
- Experience with Bird route server
  - [See "The BIRD Internet Routing Daemon," http://bird.network.cz
  - Connecting NGDF and DE-KIT via GTS
    - First approach
    - Static vlan to GTS (Géant Test system)
    - Synchronize and deploy dynamic BGP via
      - Route server
      - BFD
    - Second approach
      - MultiNREN GTS (GTS@NORDUNet+Géant+DFN)
  - Third approach
    - Include transatlantic end sites (FNAL?)
  - [see slides for test setup details]

https://indico.cern.ch/event/527372/contributions/2296779/attachments/1339877/2017635/P2P_BGP_Route_Server_Proof_of_Concept_BrH.pdf


# 3   What next?

- There are still important use cases for dynamic, guaranteed bandwidth, circuits
- However, end users who consider the service potentially important are needed to participate in and to drive testing
- DTNs are an enabling technology for large-scale data movement
  - Some reference implementations for high-speed DTNs need to be documented and put in a publicly assessable location (perhaps the LHCONE web site at CERN)
    - A know working implementation with part numbers for everything (motherboard, CPUs, storage elements, storage controllers, etc.), software and firmware version numbers, chassis layout for cooling and exact slot usage, etc. Sufficient information that when such a blueprint is followed that result will just work.
- Shawn McKee reports that they have had some successes and some failures with Open vSwitch

managing parallel production and experimental network environments. He will report his experiences at the next meeting.

- Azher Mughal will also report on Caltech's experience with Open vSwitch, as will Justas Balcas of CERN
- There is a WLCG meeting on Jan 10, 2017 where the LHC software developers will discuss their plans (at CERN, I believe). Some of the LHCONE Architecture group folks should try and attend.
- The next LHCONE meeting will be hosted by Brookhaven National Lab (US ATLAS Tier 1) and probably physically be held at Columbia University in New York City.