

Networking: The View from HEP

Shawn McKee / University of Michigan

With input from many colleagues in WLCG, OSG and LHCOPN/LHCONE

CHEP 2016, San Francisco, October 11, 2016

HEP Networking

- I was asked to present the networking view “from HEP” and while I don’t want to presume to speak for all of you or the broader HEP community, I do want to tell you about recent activities related to Global HEP networking, our challenges and a bit about how I see networks evolving
- In discussions about networking, I have gotten a broad range of feedback:
 - Most physicists I have spoken with are very happy with the global R&E network infrastructure supporting HEP and its positive enabling role.
 - A lot of people involved in supporting our infrastructure have very basic concerns about our ability to **effectively use the networks we have**.
 - Almost everyone would like the network to remain transparent as possible while continuing to deliver excellent performance

First: A Little Context

- As **Jim Seigrist** noted in his presentation yesterday, there is a long history of work by physicists (led in large part by **Harvey Newman**) to enable **HEP** networking going back to 1986 (and actually starting in 1981 with a Caltech-CERN modem link)
- In the late 1990s the **MONARC** team developed a model of how **LHC** experiments might construct a suitable infrastructure accounting for compute, storage and networking
 - Model assumed the network was expensive, somewhat unreliable and not very performant.
 - The hierarchy of tiered computing centers was the output
- **After the LHC turn-on the experiments found that the network was actually one of the most reliable and best performing components of our global infrastructure**
 - And that excellent wide-area networking (WAN) was generally being provided without direct cost to the experiments
- Based upon the experience in Run-1 the LHC experiments evolved their computing models to **take better advantage of the network.**
 - The hierarchical model was replaced by more egalitarian access to data and sites
 - Direct access to data across the WAN became part of the toolkits (AAA, FAX, etc)

HEP Network Summary

- HEP (and especially LHC) networking is:
 - **Global**
 - **Foundational** to our computing models and infrastructure
 - Continuing an **exponential increase in bandwidth use**
 - **Functioning well** but facing some current and future challenges
- The **HEP** community has significantly benefited from the world-wide Research & Education (R&E) networking community
- There are a number of (relatively) small efforts in HEP engaged in network-related areas which I will try to cover
- **While our wide-area networking needs are significant and have historically been the poster child for globally distributed e-Science, this may be changing over the coming years.**

Ongoing Work

HEP Networking

- Here is a quick snapshot of those involved in HEP Networking
 - The Open Science Grid
 - The WLCG Network and Transfer Metrics WG
 - Many institutions and communities supporting HEP networking
 - R&E **backbone networks** like ESnet, Internet2, GEANT,...
 - **NRENs** across the globe
 - **Communities** like LHCOPN/LHCONE, GLIF, perfSONAR Developers...
 - And all the **many institutions** around the world involved in network research relevant to HEP (**way too many to list!**)
 - **Our challenge is to incorporate this work into our infrastructure**

OSG and WLCG Network Efforts

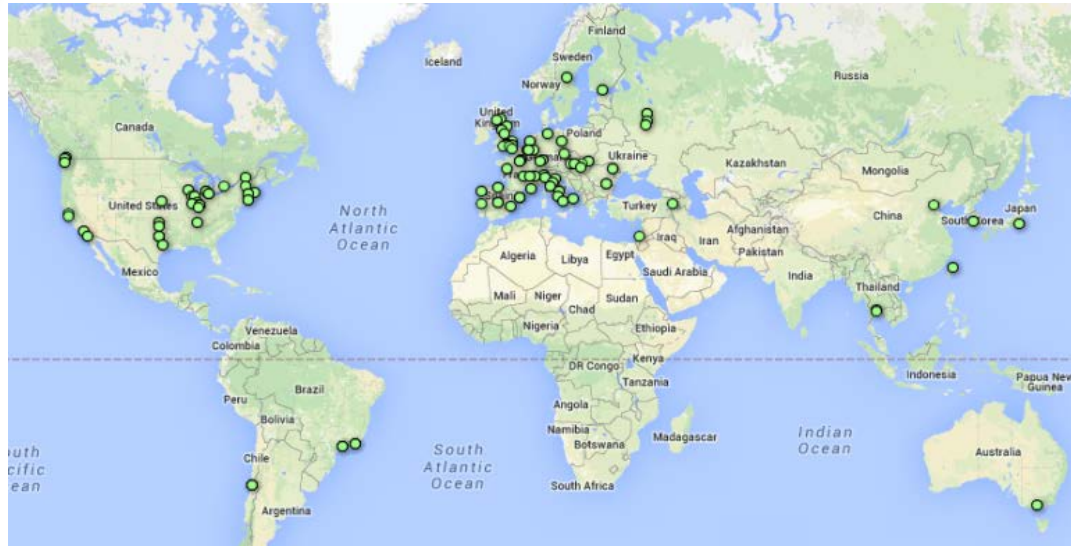
- OSG is in its fifth year supporting WLCG/OSG networking and is focused on:
 - Assisting its users and affiliates in **identifying** and **fixing network bottlenecks**
 - Improving the ability to manage and use network topology and network metrics: Analytics Platform based upon ELK in use
 - Supporting higher-level network services via network metric summarization
 - Developing effective **Alarming and Alerting** for network problems
- The WLCG Network and Transfer Metrics working group has created a **support unit** to coordinate responses to potential network issues
 - Tickets opened in the support group can be triaged to the right destination
 - Many issues are potentially resolvable within the working group
 - Network issues can be identified and directed to the appropriate network support centers
 - Documented at https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network_Performance_Incidents
 - Many issues resolved **within hours** mainly due to using perfSONAR information

Importance of Measuring Our Networks

- End-to-end network issues are frequently difficult to spot and localize
 - Network problems are multi-domain, complicating the process
 - Standardizing on specific tools and methods allows groups to focus resources more effectively and better self-support
 - Performance issues involving the network are complicated by the number of components involved end-to-end.
- **Network problems can severely impact WLCG experiment's workflows** and can take weeks, months and even years to get addressed!
- **perfSONAR** provides a number of standard metrics we can use
- **Latency measurements** provide one-way delays and packet loss metrics
 - Packet loss is almost always very bad for performance
- **Bandwidth tests** measure achievable throughput and track TCP retries (using Iperf3)
 - Provides a baseline to watch for changes; identify bottlenecks
- **Traceroute/Tracepath** track network topology
 - Measurements are only useful when we know the exact path they are taking through the network.
 - Tracepath additionally measures MTU but is frequently blocked

Current perfSONAR Deployment

http://grid-monitoring.cern.ch/perfsonar_report.txt for stats



249 Active perfSONAR instances

199 Running latest version (3.5.x)

95 sonars in latency mesh

- 8930 links measured at 10Hz
- packet-loss, one-way latency, jitter, ttl, packet-reordering

115 sonars in traceroutes mesh

- 13110 links
- hourly traceroutes, path-mtu

102 sonars in bandwidth mesh

- 10920 links (iperf3)

<https://www.google.com/fusiontables/DataSource?docid=1QT4r17HEufkvnqhJu24nIptZ66XauYEIBWWh5Kpa#map:id=3>

- Initial deployment coordinated by WLCG perfSONAR TF
- Network commissioning by WLCG Network and Transfer Metrics WG

LHCOPN/LHCONE

- The **LHCOPN** working group was established by CERN, the WLCG Tier-1 sites and the various HEP related research and education networks to define, deploy and operate the LHC Optical Private Network interconnecting the Tier-1 and the Tier-0 at CERN
- The success of **LHCOPN** for the Tier-1s led to the creation of a similar network to support the Tier-2s and their interactions with the Tier-1s: The LHC Open Network Environment (**LHCONE**)
- The **LHCOPN/LHCONE** group meets jointly 2-3 times per year to discuss policy, operations and future evolution necessary to support the LHC (and now beyond) community.
 - This mostly volunteer effort has been very beneficial for LHC
 - There is a request to increase the participation from the experiments

Ongoing work in Network Analytics

- The volume and complexity of network related data being collected by OSG and the experiments is challenging to use, but holds the promise of providing much deeper insights into our networks and hard to identify network problems.
 - To be most useful, the data requires cleaning, augmenting, transforming & correlating
- Ilija Vukotic (Univ. of Chicago) has developed ELK/jupyter stack for ATLAS Analytics and worked with Xinran Wang on [anomaly detection and advanced alerting/notifications](#) for network problems (See Track 5 talk Thursday afternoon)
 - Also looked at detection of the anomalies based on machine learning models
- Jerrod Dixon and Brian Bockelman (UNL) exploring network analytics in CMS
- Henryk Giemza (NCBJ), Federico Stagni integrating perfSONAR in DIRAC for LHCb
- Shawn McKee (Univ. of Michigan) working on real-time root cause analysis ([PuNDIT](#)) in collaboration with perfSONAR developers
- Hendrik Boras and Marian Babik (CERN) working on developing models for network cost-matrix - determine performance of network paths

Challenges

Network Operations

- Deployment of perfSONARs at all WLCG sites made it possible for us to see and **much more easily debug end-to-end network problems**
- A group focusing on helping sites and experiments with network issues using perfSONAR was formed - [WLCG Network Throughput](#)
 - Reports of non-performing links are actually quite common (almost on a weekly basis)
 - Most of the end-to-end issues are due to faulty switches or mis-configurations at sites
 - Some cases also due to link saturation (recently in LHCOPN) or issues at NRENs
- Recent network analytics of LHCOPN/LHCONE perfSONAR data also point out **some very interesting facts** about our networks:
 - **Packet loss greater than 2% for a period of 3 hours on almost 5% of all LHCONE links**
- Network telemetry (real-time network link usage) is likely to become available in the mid-term (but likely not from all NRENs at the same time)
- **It is becoming increasingly important to focus on site-based network operations**

Network Diversity

- We have a range of capacities and funding models across our global set of sites.
 - Some sites don't explicitly pay for WAN but others may need to pay for:
 - Their WAN connections
 - Any "special" services
 - Excessive bandwidth use
 - Support
 - Tier-1 and Tier-2 sites have a range of bandwidth to the WAN from 1 - 200 Gbps
- **This diversity in capability and cost leads to differences in perspectives about HEP networking planning and goals**
 - Sites with excellent networking and small costs want to see the network emphasized (to reduce other costs or improve capability)
 - Conversely, sites with lower capacity or "expensive" networking want to have the infrastructure able to conserve its networking use.
 - It can be challenging to get consensus in how much we should emphasize use of the network

Making the Most of our Networks

- Much of our WLCG infrastructure is NOT tuned to take the best advantage of the networks we currently have
 - There are a wide range of **mis-configurations**, **non-optimal tunings** and **incorrect application** and **hardware** settings that lead to inefficient use of our networks
 - As mentioned, we have a wealth of data now available and ready for analysis to **identify bottlenecks** and **poor performance**.
- As we identify bottlenecks and poor performance we need to take the next step and work to improve our end-host's ability to effectively utilize the network we have
 - Doesn't require **SDN**, **new hardware** or **new networks** but can make a huge difference in network throughput for sites
 - Should we organize a near-term workshop to share best practices, tools and tuning information?

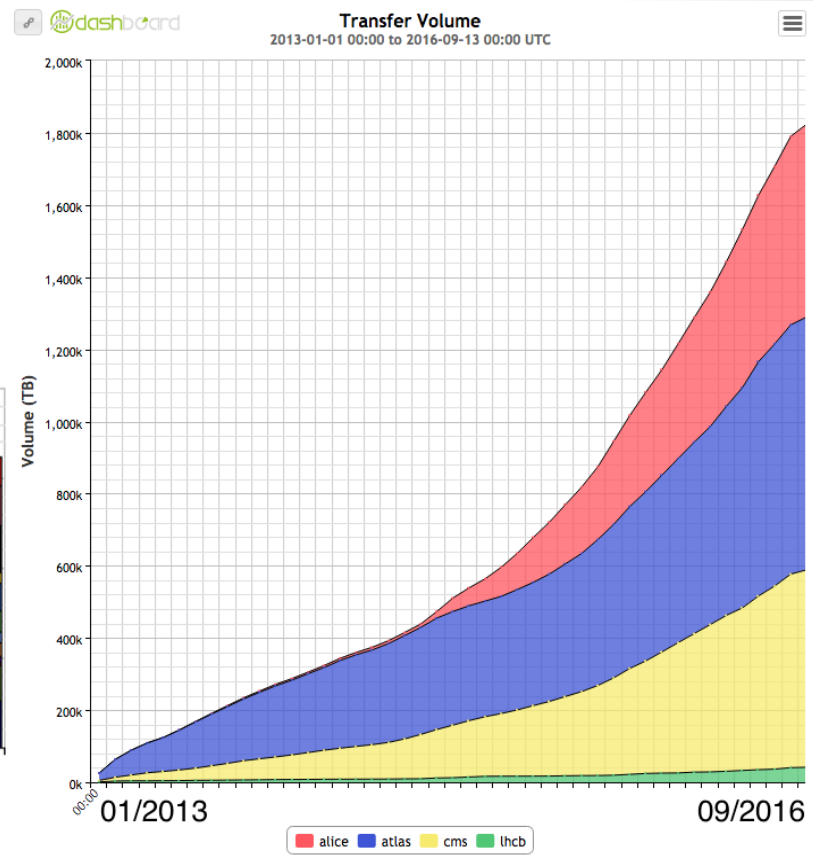
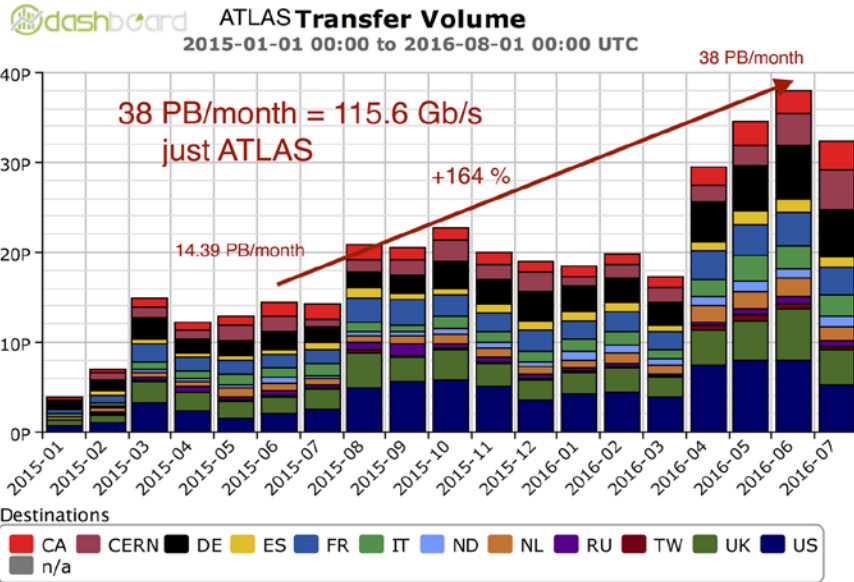
Planning for the Future

Planning for Future HEP Networking

- As a community we need to think about what we want to do regarding networking and at what timescales
 - There is a vision of a long-term evolution producing “Smart Nets” ...what characteristics with they have and how much work will we need to do to take best advantage of them?
 - What things should we worry about in the **near-term**? The **mid-term**? The **long-term**?
- Much of this will be part of the HSF community white-paper effort that was discussed this last weekend.
- I have some thoughts I wanted to share here.

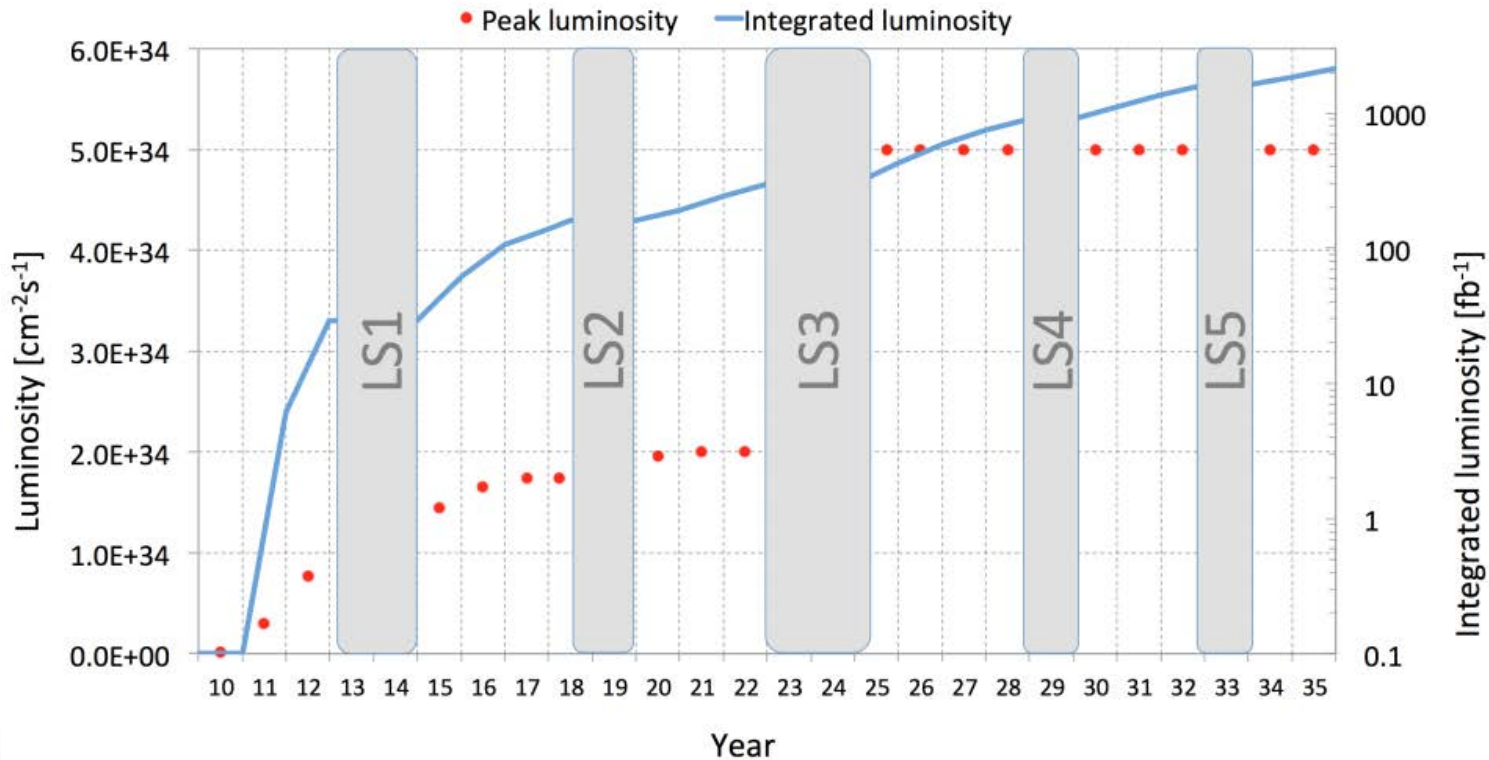
LHC Data Growth

Experiments have been transferring exponentially increasing amount of data since startup. **This trend is likely to continue as it's driven by increasing data volumes, more capable infrastructure and excellent networks.**



LHC schedule

We will see significant pressure on network resources, which will likely accelerate in HL-LHC (x10). Major increases in funding are not expected and **will likely remain flat**.

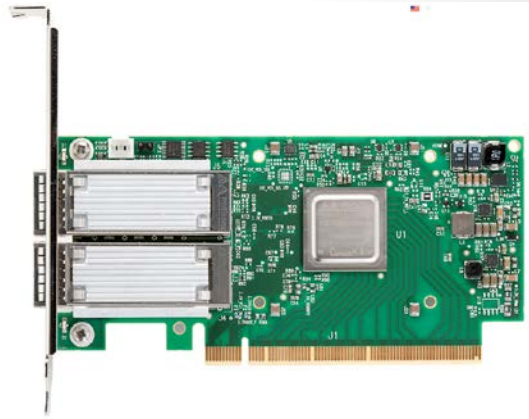


Improving End-host Networking

- **New operating systems and associated end-host improvements in hardware are making it easier to get high-performance on our wide-area networks**
- TCP more stable in CC7, throughput ramp ups much quicker
 - Detailed [report](#) available from Brian Tierney / ESN
- Fair Queueing Scheduler (FQ) available from kernel 3.11+
 - Even more stable, works better with small buffers
- Best single flow tests show TCP LAN at 79Gbps, WAN (RTT 92ms) at 49Gbps
 - IPv6 slightly faster on the WAN, slightly slower on the LAN
- New TCP congestion algorithm ([TCP BBR](#)) from Google
 - Google reports 2-4x performance improvement on path with 1% loss (100ms RTT)
 - Early testing from ESN less conclusive, there is also question how tolerant BBR will be with other congestion algorithms on the same link.

WAN vs LAN capacity

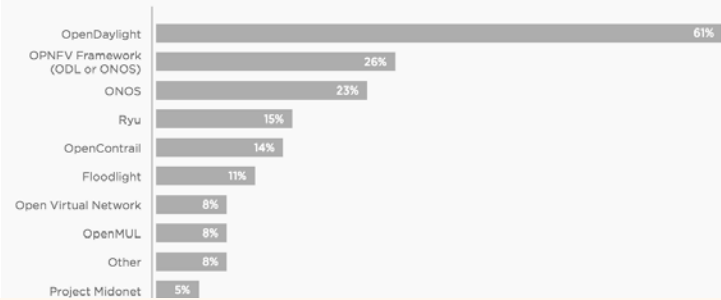
- Historically WAN capacity has not always had a stable relationship compared to data-centre
 - In recent history WAN technologies grew rapidly and for a while outpaced LAN or even local computing bus capacities
 - Today 100Gbps WAN links are the typical high-performance network backbone link speeds, but LANs are also in the same range
 - List price for 100Gbit dual port card is ~ \$1000, but significant discounts can be found (as low as \$400), list price for 16 port 100Gbit switch is \$9000
- Today it is easy to over-subscribe WAN links
 - **in terms of \$ of local hardware at many sites**
- **Will WAN be able to keep up ? Likely yes**, however:
 - We did benefit from the fact that 100Gbit WAN was deployed on time for Run2, might not be the case for Run3 and 4
 - By 2020 800 Gbps waves likely available, but at significant cost since those can be only deployed at proportionally shorter distances (thus more repeaters are needed)
- Planning of the capacities and upgrades (NREN vs sites) will be needed



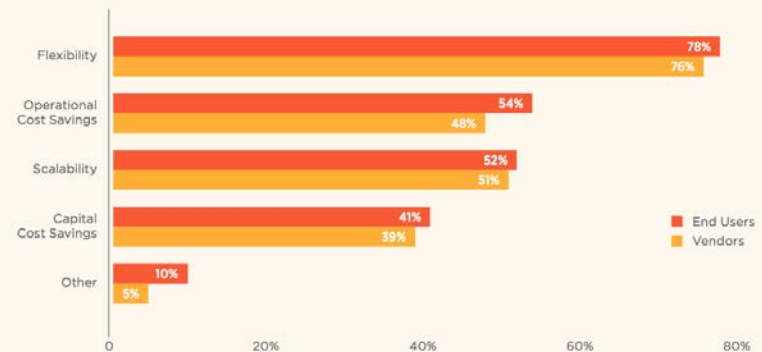
Software Defined Networks (SDN)

- SDN is a set of technologies offering solutions for many of the future challenges
 - Current links might handle ~ 6x more traffic if we could avoid peaks and be more efficient
- Many different point-to-point efforts and successes reported within LHCOPN/LHCONE
 - The challenge remains getting this end-to-end
- While it's still unclear which technologies will become mainstream, it's already clear that software will play major role in networks in the mid-term (commercially driven)
 - Will experiments have effort to engage in the existing SDN testbeds to determine what impact it will have on their data management and operations ?

OPEN SOURCE NV/SDN SOLUTIONS DEPLOYED



BIGGEST BENEFITS OF NV



Respondents could choose multiple benefits of network virtualization.
Results from 116 end-users and 85 technology vendors

SD-WAN (Software Defined Wide-Area Networking)

- Large **Network as a Service** providers include several well established CSPs such as Amazon, Rackspace, AT&T, Telefonica, etc.
- Recently more niche **NaaS** providers have appeared offering **SD-WAN** solutions
 - Aryaka, Cloudgenix, Pertino, VeloCloud, etc.
 - Their offering is currently limited and not suitable for high throughput, but evolving fast
- **SD-WAN** market is estimated to grow to **\$6 billion in 2020** (sdxcentral)
- Will low cost WAN become available in a similar manner we are now buying cloud compute and storage services ?
 - **Unlikely**, our networks are shared and global, not easy to support LHC requirements
 - Transit within major cloud providers such as Amazon currently not possible and unlikely in the future, limited by regional business model - but great [opportunity for NRENs](#)

R&E Networking

- R&E network providers have been working in mutually beneficial ways with HEP for a long time because:
 - HEP (especially LHC) has been representative of future data intensive science domains
 - Can serve as a testbed environment for early prototyping of evolving capabilities
- Big data analytics requiring high throughput no longer limited to HEP
 - SKA (Square Kilometer Array) plans to operate at data volumes **200x current LHC scale**
 - Besides Astronomy there are MANY science domains anticipating data scales beyond LHC, cf. [ESRFI 2016 roadmap](#)
- **What does n more HEP-scale science domains competing for the same network resources imply?**
 - Will HEP continue to enjoy “unlimited” bandwidth and prioritised attention or will we need to compete for the networks with other data intensive science domains ?
 - Will there be **AstroONE** or **BioONE** soon? **Will they bring additional network funding?**

Draft Perspective on Needed Effort

- **Short-term (1-2 years):** Focus on network monitoring, debugging and analytics. Find and fix network problems, improving our ability to utilize the networks we have.
- **Medium-term (3-7 years):** Plan for and evaluate the use of SDN for our infrastructures. Work on integration of those aspects deemed beneficial. Estimate the impact of other data-intensive science domains on our R&E networks and collaborate with them on their ramp-up to our scale.
- **Long-term(8-12 years):** Plan for and deal with the R&E network environment: **sharing, orchestration, automation** and the **implementation of smart networks**. Ensure our software can interact with smart network capabilities and agilely respond to dynamically changing infrastructure capacities and problems.

Upcoming Network Meetings

- Pre-GDB on Networking, January 10, 2017
<https://indico.cern.ch/event/571501/>
- LHCOPN/LHCONE Meeting at BNL end of March / beginning of April 2017. Being scheduled soon

Conclusion and Summary

- HEP Networking has been a reliable, high-performing infrastructure component
 - Still work to do in **finding / localizing problems** (new or existing) and **fixing bottlenecks**
- New technologies and our own work will make it easier to increase data transfers
- We must track network capacities and technology changes
 - HEP will continue to exponentially increase its use of the network for the foreseeable future and, considering only HEP, **this seems sustainable by R&E networks**
 - Site vs NREN capacity upgrades, **HEP computing model evolution** and their relative timing, needs to be watched
 - But increasingly, HEP likely no longer the only domain using global R&E networking
- Sharing the future capacity will require greater interaction with networks
 - While unclear on what technologies will become mainstream, we know that software will play a major role in the networks of the future and we need to be ready to use it

Questions or Comments?

References

- WLCG network Use-cases document for experiments and middleware
<https://docs.google.com/document/d/1ceiNITUJCwSuOuvbEHZnZp0XkWkwdkPQTQic0VbH1mc/edit>
- Harvey's slides from Nordunet covering HEP networking history and ongoing work
https://www.dropbox.com/s/at2ky4rdc6szkmq/NGenIAGlobalNetworks_hbn091916.pptx?dl=0
- OSG Network Documentation
<https://www.opensciencegrid.org/bin/view/Documentation/NetworkingInOSG>
- WLCG Network and Transfer Metrics Working Group
<https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics>
- perfSONAR deployment documentation for OSG and WLCG
<https://twiki.opensciencegrid.org/bin/view/Documentation/DeployperfSONAR>
- WLCG workshop October 8, 2016 networking session presentations
<https://indico.cern.ch/event/555063/sessions/203482/#20161008>