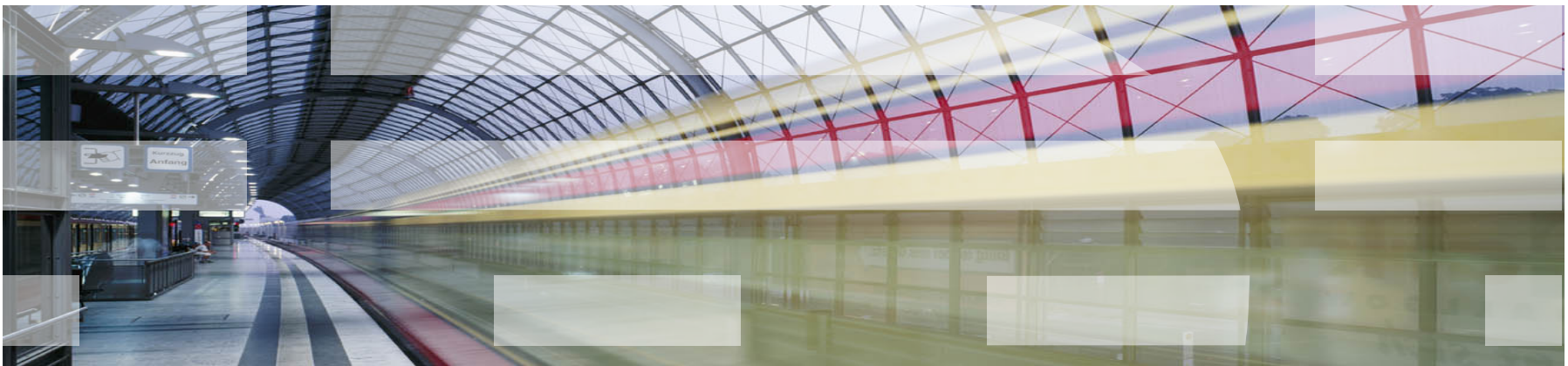


Connecting HPC and High Performance Networks for Scientists and Researchers: Final Survey Results



***SC15 Austin, Texas, November 18, 2015
Final Survey Results as of January 14, 2016***



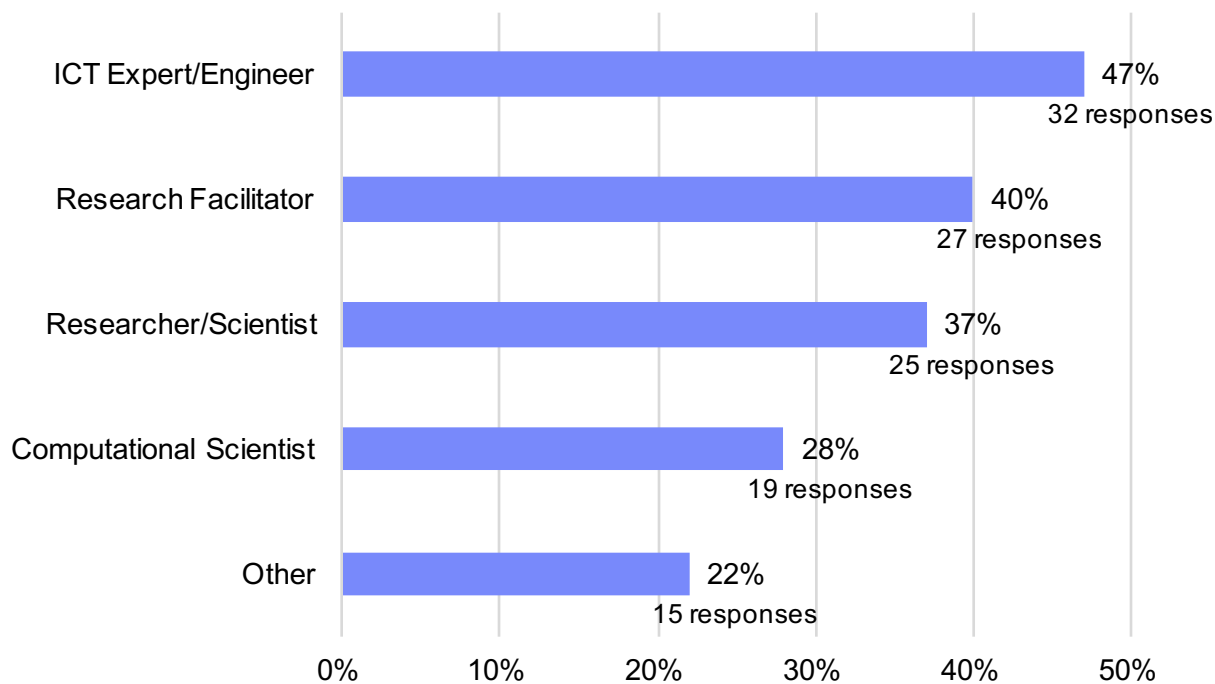
Prominent themes from survey participants:

“High Performance Computing / High Throughput Computing resource needs, challenges and best practices” (68 responses)

1. Energy & Environment and Healthcare & Life Sciences are emerging users/consumers of HPC and networking
2. Most respondents are using local HPC/HTC resources
3. IT departments are heavily involved with HPC/HTC
4. HPC data size is expected to double within next 2 years – currently working with petabytes or terabytes of data
5. Difficulty with data transfer speeds amongst multiple locations and organizations
6. Speed and resources (storage, technical, analytical, tools, and human) are problematic
7. Globus is a popular choice in leveraging high performance networks for research computing

Nearly half of respondents classify their role as ICT Expert/Engineer, followed by Research Facilitator and Researcher/Scientist

Q1: What is your role? Choose all that apply.

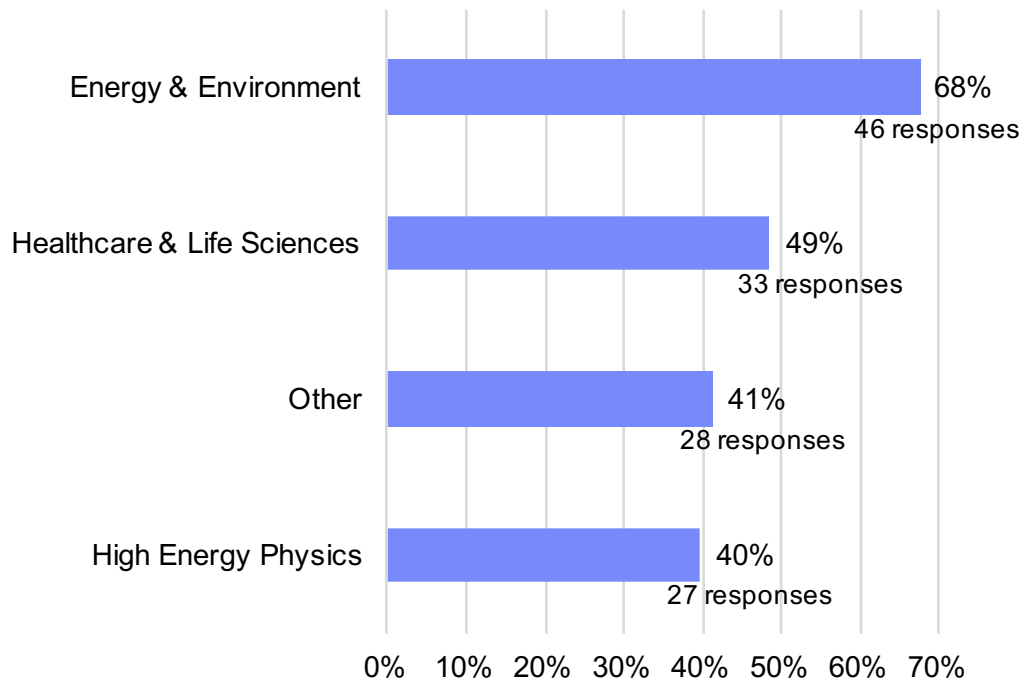


Other responses included:

- Director (2)
- Director of Advanced Research Computing
- Research Computing Director
- Director of High Performance Computing
- IT senior director for university research computing
- Deputy CIO for Research... interface between central IT resources and our researchers
- CIO
- Instructor
- State coordinator for research cyberinfrastructure in higher education
- IT Project Manager
- Research support
- HPC user support
- Research computing service provider
- Data scientist

Energy & Environment and Healthcare & Life Sciences are the leading disciplines for HPC/HTC research areas

Q2: Please specify your area of research or areas of research that you support.

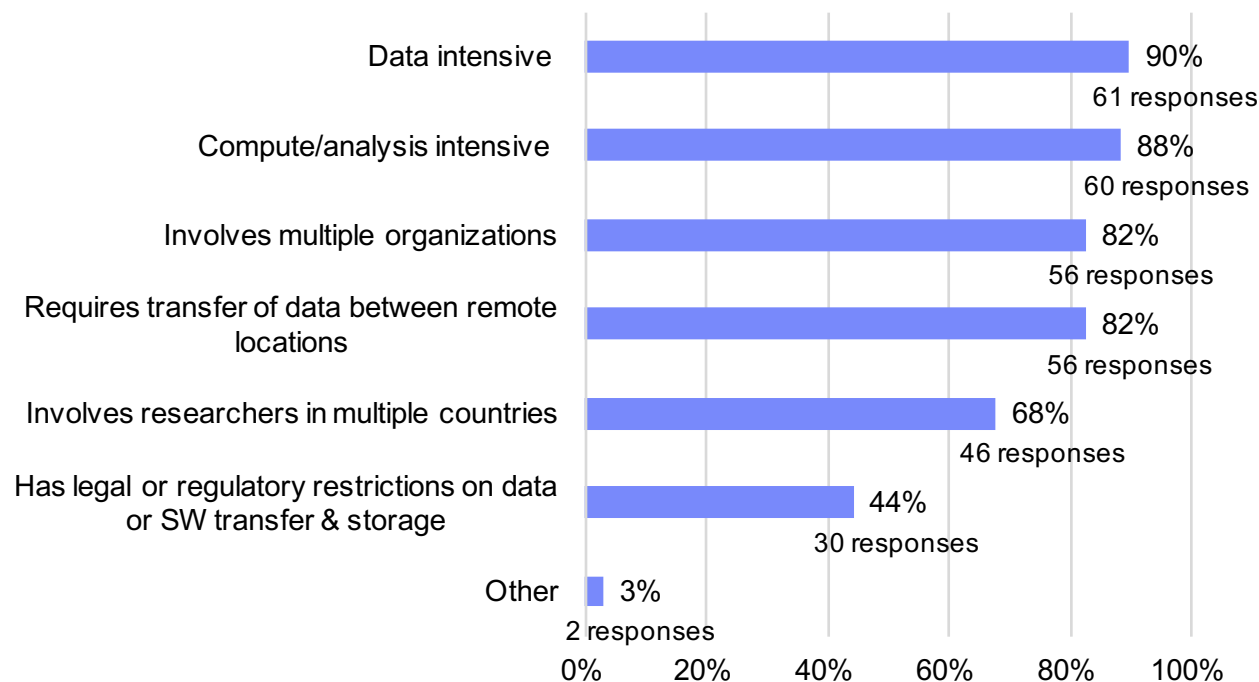


Other responses included:

- Statistical computing, geostatistics, time series
- Condensed matter physics
- Materials Science / Condensed Matter Physics / Statistical Physics
- All disciplines across campus
- All sciences and engineering
- Computational chemistry & material science, virtual training & collaboration environments, IoT/CPS for freshwater research
- Broad Academic
- All areas of modern computational research from political science to CFD to high-energy physics
- Humanity
- Geophysics
- All engineering disciplines, all science disciplines, social sciences, humanities, architecture and business
- Weather, Computational Chemistry, Material Science
- Hydrodynamic / offshore engineering; marketing; economics
- Engineering
- Chemistry, materials, macromolecular science, mechanical engineering, biomedical engineering
- Computer science
- HPC
- Marine Science, Extreme Computing, Desalination
- All computational science
- All
- Big Data/Analytics, Internet of Things/Sensor Nets, Water Quality/Environment
- Modelling
- Financial services
- Materials science
- Social science, Humanities, Physical Sciences, and Engineering generally
- Transportation
- All science domains across all campuses
- Materials sciences

HPC/HTC needs revolve around sharing compute intensive data amongst multiple locations and organizations

**Q3: How would you describe your research or the research you support?
Choose all that apply.**

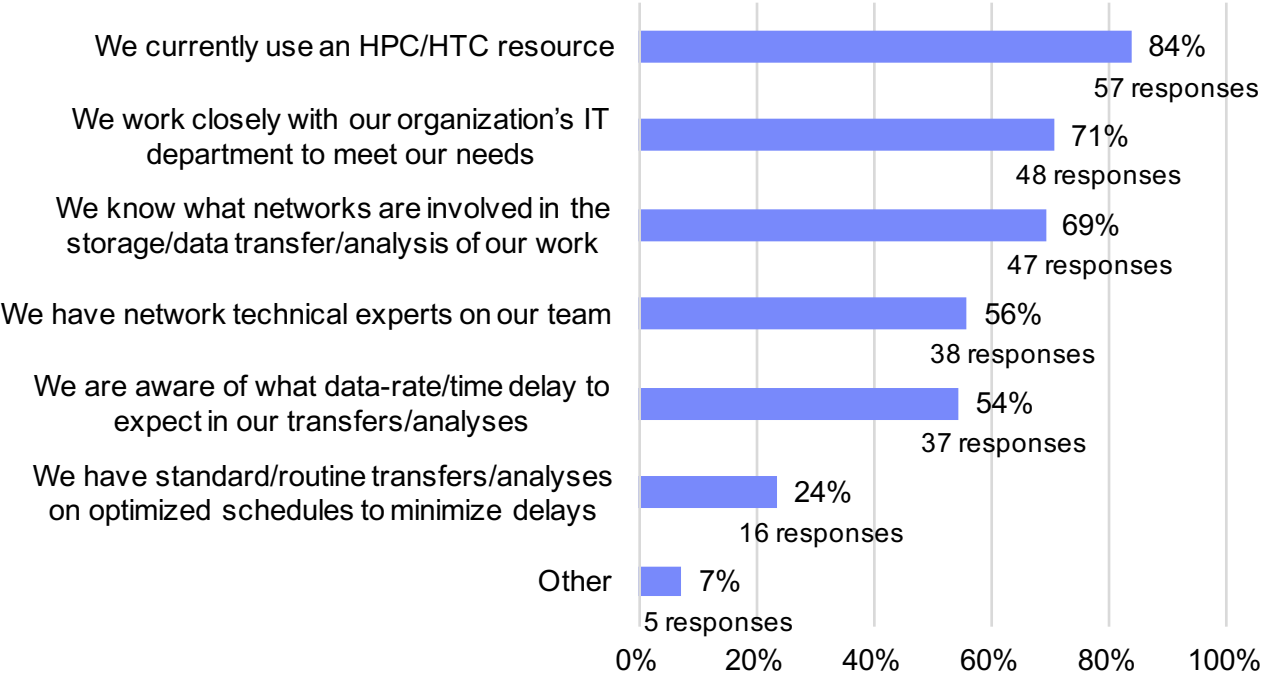


Other responses included:

- N/A
- Requires service provider carrier grade technology for the networks

Most respondents are currently using an HPC/HTC resource, understand the networks involved with their data, and are closely tied to IT

Q4: How would you describe your HPC/HTC resource needs? Choose all that apply.

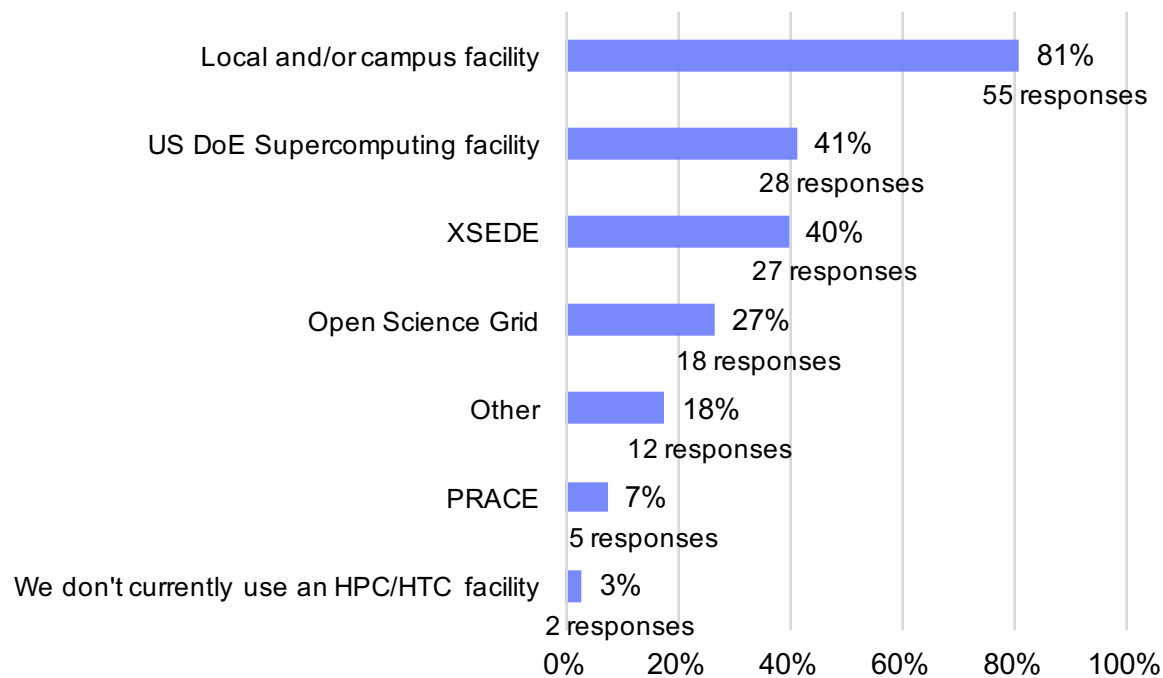


Other responses included:

- We are part of the IT department
- We make extensive use of Globus services
- We have our own HPC including the 7th fastest computer in the world
- I represent the needs of a large research university. The needs are fundamentally insatiable in some areas. There is an urgent need for a truly competitive market for research computing services to emerge.
- We run HPC and HTC resources

Most respondents are currently using a local or campus HPC/HTC facility

Q5: Which HPC/HTC resource do you currently use? Choose all that apply.

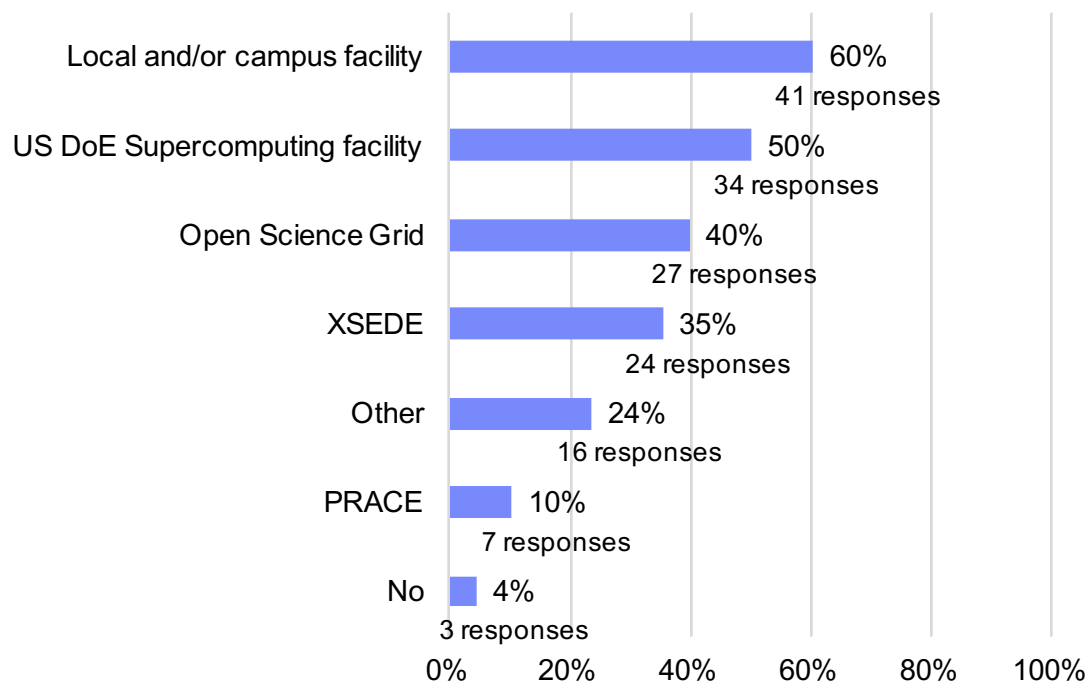


Other responses included:

- National e-infra resources (compute, storage, network)
- AWS
- TACC
- Statewide HPC resource
- European grid, national research cluster, national research cloud infrastructure
- Users we support use these resources when they can
- Network Infrastructure is in pale for up to 4x10 I2 connection to National Labs
- Compute Canada
- Other federal government facilities (NASA, NOAA), cloud providers (AWS, Softlayer)
- Yellowstone (NSF Climate Science Facility)
- NCAR, NASA
- HPCI (Japan)

Most respondents would like to continue using a local facility or a US DoE supercomputing facility

Q6: Would you like to use an HPC/HTC resource? Choose all that apply.

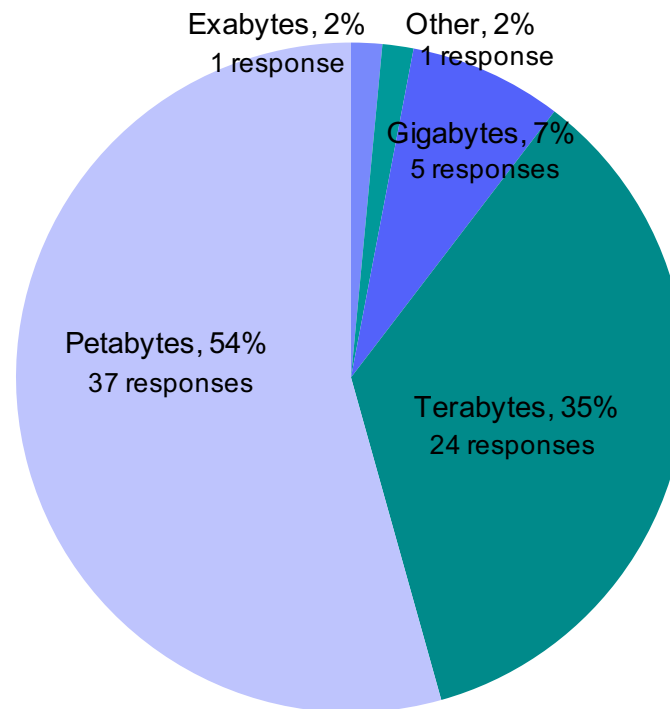


Other responses included:

- Private cloud/Boinc Volunteer computing
- We would like to assist our researchers in utilizing whatever HPC/HTC resource they need
- AWS
- AWS cloud
- I would like to help our researchers gain access to any/all of the resources available to them. The facilities checked above are ones of which I'm aware are being used
- Yes
- Network Infrastructure is in place for up to 4x10¹² connection to National Labs
- Compute Canada
- The community I represent already uses all the listed resources, as well as commercial services
- We run HPC and HTC resources
- Already using per Q5
- Globus
- CloudLab
- Yellowstone
- Would like to continue using the resources
- HCPI (Japan)

**54% of respondents are working with petabytes of data.
89% are working with data measured in at least terabytes.**

Q7: What is the total volume of data you use and/or generate for your research (or the research you support)?

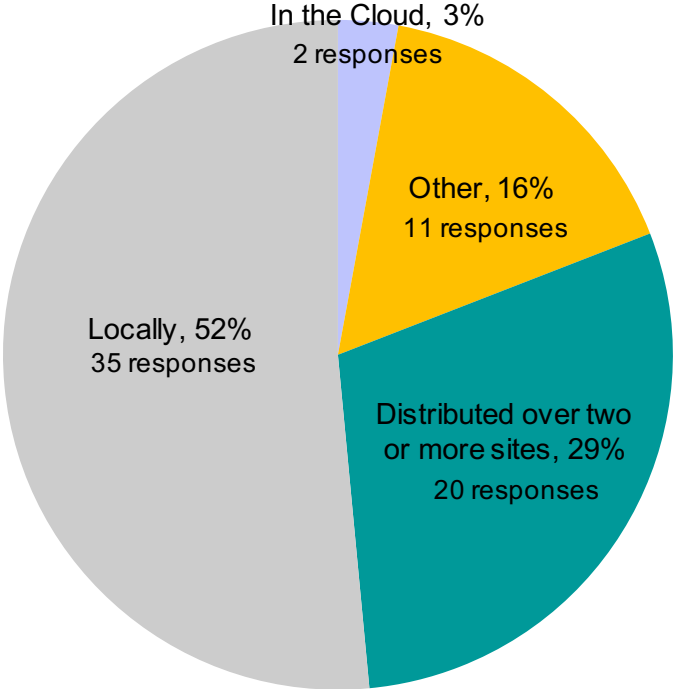


Other responses include:

- We run HPC and HTC resources

Most respondents store their data locally, followed closely by distributed over multiple sites

Q8: How do you store your data?



Other responses included:

- National e-infra resources which include storage
- All of the above (5)
- We provide cloud storage
- Hadoop HDFS cluster
- We run HPC and HTC resources
- Multiple platforms
- PSU GPFS storage

Issues identified that impede research center around speed and proper resources

Q9: Is there any known problem that impedes your research activities? (e.g., data transfer to/from remote site is too slow, inconsistent or unpredictable, etc.)

Speed

1. Slow transfer, sometimes trouble accessing files between team members
2. Low bandwidth, low performance last mile, and networking
3. Inconsistent and unpredictable at times
4. Transfer to/from remote sites slow, no reliable latency on internal network
5. Quality of Service for data transfers; conflicts between research and campus security needs.
6. Inter-campus data transfer and transfer to/from XSEDE
7. Need more aggregate bandwidth 10->100Gbps
8. Data transfer to/from remote site is too slow, inconsistent
9. Data mobility between facilities remains a challenge. There are challenges around basic authentication and authorization as well as challenges around efficient services that use available bandwidth effectively.
10. Always better to be faster
11. Too slow
12. Fileserver I/O speed are unpredictable due to million of files being updated. ZFS with SSDs alleviating this.
13. Data transfer too slow, transfers inconsistent
14. "Last mile" problem me of limited data rates from labs to campus core network.
15. Robust data transfer
16. Local infrastructure is slow/obsolete
17. Network speed constraints; in the process of upgrading to 100Gb or better connections
18. It could always be faster
19. Transfer speed fluctuation
20. There is variability in transfer speed between sites. The data management work should be automated (data movement, archive to HPSS, standard analysis, documentation). While some sites have this capability and some are in progress, others have not started. There is a limitation in the amount of open data storage that can be utilized.
21. Data transfer too slow
22. Data transfer to/from remote site is too slow

Storage

23. Storage performance
24. Local storage limitations
25. Access to high speed storage

Other

26. Yes

Proper Resources

27. Inefficiency when transferring a large number of small files
28. Need to establish policy for storage, access, and transport of Research Data
29. Existing 'legacy' software that expects data to be available on a locally mounted filesystem
30. Various "costs" associated with transient data placement and aggregation that would allow HTC solutions vs. simplicity of existing HPC solution
31. Complex user administration; complex system administration
32. Users are still developing data management strategies.
33. Lack of local network expertise to facilitate data movement
34. Data share
35. High cost of data communications in and out of Saudi Arabia
36. Campus firewall
37. Non-standard network configurations are unstable even with the greatest experts on the problem
38. Accessibility due to business firewalls
39. Mostly on campus networking issues in specific buildings here specific researchers are located
40. Data storage, transfer and archiving is not yet seamless
41. Regulatory and data custody issues
42. Relative immaturity of master data management practice, lack of accurate global time
43. synchronization, nascent internal state of semantic technologies and huge web of conflicting regulatory (domestic and international) requirements and constraints.
44. Works very well for large institutions; small sites have firewall/bandwidth/expertise limitations
45. Lack of super computing facility
46. Try to make it easier for new partners
47. The usual tuning issues
48. Globus is not installed at all sites, or DTNs not configured properly at all sites
49. Workflow orchestration tools need to be made more mature
50. Tools for interacting with remote data
51. Having the right people that can support the computational needs of the disciplines
52. Security is likely the biggest issue
53. Technical support/staffing available to assist in optimal use of resources
54. Data Management between distributed locations
55. Checkpoint capabilities incompatible with maximum allowed wall time
56. Yes, all of the above. Plus lack of security on administrative networks, would prefer private networks.

No Known Issues

57. No (7)
65. Not really (2)
67. NA (2)

Most respondents expect their data to at least double within the next 2 years

Q10: What is the expected growth rate of your data? (e.g., will double annually for 5 years)

At Least Double Within Next 3 Years

1. Double next 3 years
2. Expect to double every year or so
3. Annual doubling is a reasonable guess (but just a guess)
4. Increases by a third each year.
5. Doubles every 3-5 years
6. Could double over next 2-3 years.
7. 35% year over year growth
8. Double every two years
9. Double every 3 years
10. Double in 2 years
11. Double every two years
12. Doubling around every 12-18 months
13. Double each second year
14. Roughly double in three years
15. Double every 2.5 years
16. Will double over next three years
17. Double within two years ... not know farther out

Other

18. Linear growth
19. Hard to say, but exponential
20. Order of magnitude growth as new 5-year multi-institutional project kicks off with both computational and data-intensive components
21. Haven't thought about that.---!!!!!!!
22. 150 TB / year for the next 3 years
23. 50% annual increase.
24. Petabytes per month in 5 years
25. Exponential
26. Lots
27. 150Gb a year
28. 20% growth per year
29. 50%/year
30. 20% per year
31. Ten to hundreds of petabytes
32. Higher than is comfortable
33. Not expected

12

More Than Double Within Next 3 Years

34. More than doubling annually
35. Double every 2 years
36. Doubles appx every year
37. We expect doubling every 3 years.
38. Double every three years
39. Will grow about 5X
40. Move from 1 Ped/yr to 1 Ped/mo in 3 yr, w/ 2 Ped/mo in 5 yr.
41. Factors of 5 to 10
42. Expect raw data to double every year but archived data storage will be less
43. Double each year ; no end in sight
44. Unknown but likely to double every year or so
45. Quadruple every 2 to 3 years
46. Will triple annually for the next 5 years
47. It is more than doubling every year
48. > doubling annually
49. Expected to double annually for foreseeable future
50. Circa 30% CAGR
51. Triple in 5 years
52. Doubling every 1-2 years
53. Double every year
54. Every 3-4 years the data amounts will quadruple. If significant amounts of storage are available we can start examining higher temporal and spatial resolution data, and thus, our data rates can grow
55. Will double every year
56. 35% annual
57. Double every six months -----???

Unknown

58. Unknown (11)

Globus is a popular choice in leveraging high performance networks for research computing

Q11: What success stories and/or best practices can you share for research leverage of high performance networks? (Optional)

1. Have successfully been using Globus as a solution for internal data transfers and transfers to/from other institutions, including implementation of dedicated endpoints on campus and at data centers.
2. Best Practice: Globus Online combined with a high-speed Science Network and Science DMZ
3. PerfSonar, DMZ and Globus deployments
4. Globus helps
5. Globus is a great tool for transferring large data sets
6. Globus online, with sharing, and better protocol to utilize network
7. When Globus transfers are not limited by the requirement of OSG certificates, they work great even if there is variability.
8. Science DMZ/perfSonar very helpful in moving large data sets
9. <http://www.internet2.edu/research-solutions/case-studies/accelerating-genomic-research-advanced-networking-collaborations>
10. Reprocessing all CrIS sounder data from mission start, making this generally available, statistical analysis and intercomparison of AIRS, CrIS, IASI sounder data with ECMWF model data
11. Combining the resources of the RON and UT System provides leverage to access I2, ESnet, and Community Networks to reduce payments to LEC/CLEC sources.
12. Use automated movement of data
13. Faculty led governance model.
14. CC-NIE award was instrumental but we still need a cyberinfrastructure engineer
15. Multiple virtual machines doesn't provide network performance advantage over bare-metal servers.
16. We have several projects that share observational, simulation datasets broadly. These span
17. Contact the experts supporting the infrastructure you use in case of problems or questions
18. Standardization and consolidation of resources
19. Leveraged high speed SAN with GridFTP and 40GbE connections to 100GbE uplink
20. Send the analytics to the data, not vice versa.
21. I like the new technologies, we need it to make it for the next generation to excel.
22. Part of calculation (eigenvector) is done in remote site on capability machine on one national Lab and transferred to capacity machine on another lab to use it to maximize the efficiency of each site.
23. We are piloting a staging service optimized for transfer of data (using gridftp)
24. We don't have one yet
25. Don't understand the question

Other challenges center around automation, proper tools, and human resources to handle data

Q12: What else do you want to share as a need, challenge or solution for leveraging advanced technologies in research and science endeavors? (Optional)

Automation

1. The need is to automate the handling of data to reduce project cost and to have scalable data handling. This would allow handling 10-100x more data. Our project gets ~10-13% of total cycles on major machines. Given the last two items we are unlikely to need significantly more cycles so there is little ability to grow in simulation length. Where the growth can come from is having access to significantly more scratch and archival storage, we can output more data and get finer detailed view of our simulation

Human Resources

2. Challenge: education/training for users of hpc(s)
3. Institutional resources in IT Network are necessary and not sufficient to support Research Lab requirements and Data Infrastructures. Must address HR/sourcing issues.
4. Research support resources are scarce, in particular providing last mile connectivity on campus.
5. More non-HPC researchers using HPC!!!
6. We do not have sufficient support staff to help researchers maximize use of HPC, GPU, Hadoop/Spark environments. I suspect this is a problem for many institutions.
7. Data management expertise.
8. Research data management services
9. Human resource with both domain science and computer skills

Tools

10. SDN looks promising, ScienceDMZ is needed
11. Need for tools to help federate regional resources and inter-federate with national resources
12. Virtual machines and head nodes for intensive data analysis and customized solution provisioning and sharing (such as Galaxy instances or Globus instances)
13. Need additional identity management federation
14. There is a need investment in a hierarchy of resources. Over emphasis on large monopolistic enterprises (whether tax-payer or commercially funded) are troubling and a challenge to research needs. A healthy market place for computing probably needs the challenge of data lock-in to be solved more effectively. The current market for research computing services has many problems that distort the economics and inhibit competition. Networking can play an important role in creating a functioning market.
15. Cluster tools like slurm are usable but need refinement

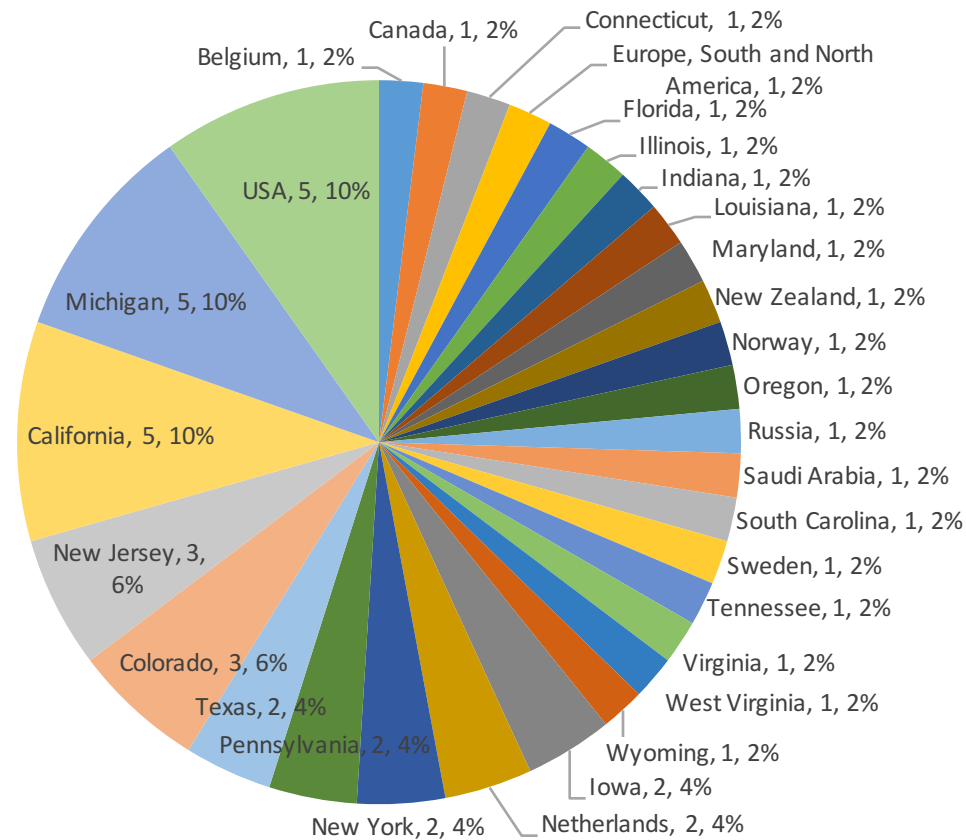
Other

16. Challenge: Cost-effective preservation for large-scale data
17. Local high speed science network
18. Challenge: getting expected performance from all the components of an emerging technology high speed network (e.g.. 40 Gb)
19. Help the next Frontier
20. Mid-tier computing using new architectures.
21. Need for movement of compute to where data rests

SC15 Survey Respondents primarily (80%) from the U.S.

More than half of respondents requested final results

Q13: What country/state do you work in? (Optional)



***Thank you very much !
Domo arigato gozaimasu !
Muchas gracias !
Merci beaucoup !
Grazie mille !***



ESnet

ENERGY SCIENCES NETWORK



engage@es.net

cino@internet2.edu

businessdevelopment@geant.org

