# Big Data in the Campus Landscape

## Basic Infrastructure Support

This paper is a publication of the ECAR Campus Cyberinfrastructure (ECAR-CCI) Working Group. It is a part of a series focusing on Big Data in the Campus Landscape. ECAR working groups bring together higher education IT leaders to address core technology challenges. Individuals at EDUCAUSE member institutions are invited to collaborate on projects that advance emerging technologies important to colleges and universities. ECAR-CCI helps educational institutions develop institutional strategies and plan resource deployment and helps users harness and optimize the power and capabilities of these new integrated IT tools and systems for educational and research applications in higher education.

## Introduction

Today, virtually every important research breakthrough—be it designing a new drug, developing new materials, analyzing literary texts, or forecasting climate change—depends on computing resources. Computing resources have become as important to scientific discovery as theory and experimentation, and with the rise of big data across all academic disciplines—including the humanities, social sciences, and even the arts—research computing infrastructure is critical to the ability of research groups to attract funding and retain the best talent. Yet the reality today is that most campus data infrastructure is developed to support administrative data requirements, which are very different from those of research big data.

In this paper, we discuss issues that big data research raises for basic campus infrastructure support, including network design, server hosting strategies, backup and disaster recovery, and campus bridging. Research computing infrastructure support services have often been specialized and at times isolated. The size and scale of research big data will, by default, impact campus IT. A well-coordinated strategy is essential.

Network architecture must be reconsidered because conventional campus network architectures that seek to minimize security threats and manage costs are not adequate to meet the performance needs for efficiently transporting research big data. Computation for big data research places high and continuous workloads on servers, outstripping most data center capacities. Standard backup and disaster recovery strategies are insufficient to address big data requirements. Finally, scientists need to use the cyberinfrastructure resources of their lab, their campus, other campuses, and regional or national (e.g., National Science Foundation or Department of Energy) centers in a seamlessly integrated way, necessitating novel campus bridging models.

We will explore emerging approaches to address these issues. There may not be a common solution for all the issues presented by big data research, but it is becoming clear that the easy approach—or worse, no approach—will negatively impact a university's competitiveness in the era of big data.

# Network Design

There is an increasingly well-understood need to maximize the network as a means for efficiently moving big data, particularly but not only in support of the campus research mission. But this need must be balanced with the continuing need to minimize the likelihood that the network could become a security threat, not only for the research big data (including, for instance, big data in health and genomics, where there are special security and privacy complications) but also for other campus data that may be exposed by the network. Another ongoing consideration is to minimize the cost of the campus network as an ordinary utility shared by the entire institution.

Several factors combine to challenge our ability to rapidly move big scientific data sets, necessitating additional requirements for network design:

- First, the size of data—both individual data sets and the aggregate data volumes used in research projects—continues to grow, driven by research needs, by the plunging costs per byte of storage, and by the rapidly increasing precision and sample rates of instruments.

- Second, unlike administrative data, for instance, these data sets are often the subject of collaborations between geographically distant scientists.

- Third, these data sets must often be processed by different high-performance computing systems and related computers in different locations. For example, raw input might be processed by large-scale MPI-based models at an XSEDE[1] site and then subjected to analysis by campus clusters.

- Finally, while the speed of computer networks does steadily increase, the rate of increase during the early 21st century has not kept pace with the increase in the size of data sets. Simultaneously, there are pressures on campus network engineers to support highly functional firewalls, deep packet inspection devices, and bandwidth limiters to protect data. These two interests—for greater performance to meet big data needs and for more packet-forwarding complexity to meet greater security needs—are in contention with each other.

Add to these factors the common state of limited funding for campus infrastructure, and it is clear that there is a need for new approaches to network design.

These new approaches include several themes. One approach, the Science DMZ, focuses on carving out a small part of the campus network to address the needs of big data movement.[2] As a portion of the campus network topology, the Science DMZ is connected to the campus's wide-area routers in such a way that high-speed, wide-area flows are not burdened by firewalls and other "packet disruption devices." Because it separates the network for use specifically by research big data, the Science DMZ is also an ideal location for file transfer nodes carefully designed to support high-speed, wide-area flows to and from high-performance file systems where key scientific data sets are stored. This approach also reduces pressures on the design and operation of the rest of the campus network.

While the Science DMZ innovates in network topology and organization, software-defined networking (SDN) innovates in switch design and protocols.[3] The key technical idea is to separate the *control* plane (which is often proprietary, with consequently high costs and obstacles to interoperability) from the *data*

plane (where vendor-specific innovation and performance are welcome). Separately, this same technology permits innovative network architectures to be implemented because SDN can be used to implement such architectures, which can then be run on a collection of switches that might come from various vendors but that all support SDN.

## Findings

- Conventional approaches to campus networking, optimized for cost-effective and secure networking for the routine work of the campus community, will not likely meet the big data needs of researchers.

- The Science DMZ approach greatly enhances the ability of campuses to address big data networking needs. Several implementations of this approach, some quite simple and some more sophisticated, have succeeded.

- Software-defined networking is emerging as a very promising approach to cost-effectively moving big data.

- The Science DMZ and SDN are complementary approaches.

## Recommendations

- Each research university should implement a Science DMZ in a manner using an implementation approach that is in line with the campus's network strategy and is appropriate to the big data needs of its researchers.

- Every research university should keep itself aware of the emerging SDN technology and consider applying it as the technology matures and as needs suggest.

# Server Hosting Strategies

Most campus data infrastructure was developed to support administrative data requirements. The conversation about server hosting strategies today usually includes a review on the use of public versus private clouds. Cloud services are maturing, and there is an abundance of cloud-related information, services, and vendors available in the market today. It is likely that most campus CIOs will have to evaluate and determine whether to use specific cloud services. It is important to note that cloud services may be used to provision computing and/or to store data. But putting data in the cloud and doing the computing elsewhere can be problematic, given that getting in and out of the cloud can be expensive and a bottleneck. A case can be made, for example, for putting administrative computing (payroll, human resources, web, etc.) in the cloud—the data are structured and the volume is relatively small. Research big data computing is on a different scale in more ways than just size and complexity. It has different funding models, different cost structures, different support requirements, and—maybe most importantly— different technical requirements and patterns of use. For example, a frequent dilemma for cloud providers is the issue of latency. This is currently a significant stumbling block where data generation is integrated with computation. Workflows for instruments such as nucleic acid sequencers and MRIs, for example, are difficult to integrate reliably.

Another significant consideration for research computing is the fact that many of the computing technologies that are used are not well supported by current cloud service providers. This will change over time, but cloud services will probably always lag when it comes to new tools because it does not

make economic sense for the services to support emerging and specialized technologies until high demand requires them. The value of control and flexibility is another important and often ignored issue in discussions about cloud services. Although administrative data and applications are structured and well defined, research data are almost always more complex and the algorithms used are continually refined—hence the need for flexibility and control.

For many institutions, data storage, sharing, and access for sensitive data (i.e., protected health information) are problems that could be addressed by thoughtful partnerships with cloud vendors. For example, a vendor may attest to be HIPAA compliant to support a university's clinical research, but due diligence by the university requires that it conduct its own NIST-based evaluation to document the vendor's ability to secure its data. Cloud vendors have not generally provided this level of access but are now responding to market demand. The university is then able to satisfy dependencies and ensure end-to-end security using its existing standards to offer cloud-based storage services for sensitive data.[4]

Cloud services for research computing are proving useful to satisfy "bursty" computing needs—when the need for resources varies greatly day to day. The typical research computing workload, however, is 24/7 and spans weeks, months, or longer. With this pattern of usage, the economics of cloud services are not favorable. The experience at the Massachusetts Green High Performance Computing Center and at universities has shown that the cost of delivering on-premises research computing cycles for the typical continuous workload can be three to five times lower than discounted rates from cloud services. A key to gaining the advantages of privately operated research computing resources is the ability to operate at a large scale and in a location where power costs are competitive. On a campus, designing the architecture and services to meet both administrative and research computing needs can help achieve an economy of scale.

One area where cloud providers have an advantage today is in automated management of computing resources, which allows them to manage more computers with fewer people, hence at less cost. This gap could be reduced by initiatives such as the Massachusetts Open Cloud project, which seeks to develop openly accessible versions of cloud-provider proprietary automated management systems.[5]

## Collaborative On-Premises Models to Support Research

### MGHPCC

The Massachusetts Green High Performance Computing Consortium was formed in 2010 to improve access to research computing resources and to encourage collaboration between institutions in Massachusetts. The founding members of the consortium are Boston University, Harvard University, MIT, Northeastern University, and the University of Massachusetts, with additional support from private companies and the state government. It is a novel collaboration that continues to evolve and grow. The facility is designed specifically to support research computing and is catalyzing a number of joint research initiatives.

### Princeton University

Princeton University developed and implemented an institutional strategy to support research and administrative computing in one facility. Changing economic and organizational factors influenced strategy development and the critical elements behind the successful collaboration between research faculty, the IT organization, and senior administration. This collaboration informed Princeton's decision to build a high-performance computing center to meet the university's research and administrative computing needs now and in the future.

## Findings

- Big data research computing workloads are typically high and continuous.
- Many needed technologies are not currently well supported by cloud service providers.
- Cloud services that support "bursty" applications are maturing.

## Recommendations

- Universities should evaluate cloud offerings for server hosting, at least as a hybrid model, to accommodate sporadic burst compute needs before building or expanding a data center.
- Universities that are expanding or building new data centers should integrate research computing needs in their strategy.

# Backup and Disaster Recovery

Backup and disaster recovery strategies for research big data focus on two needs:

1. While generating data, a researcher needs to be protected from disaster (data loss).
2. When the project is completed, a researcher needs to make enough data accessible so others can reproduce the science or disprove it (data sharing).

Both of these support the principle that research is conducted with reproducibility as a fundamental attribute. One way to accomplish this is by beginning projects with well-documented data management plans that point to clear protocols for backup and disaster recovery strategies. Of course, this presumes that supporting information technologies are in place.

Standard backup and disaster recovery strategies, however, have proven insufficient to address big data requirements. Research big data are typically very large in volume, diverse (structured and unstructured), rapidly acquired, and often variable in quality. As such, traditional backup approaches usually take a long time, must be done more frequently, and cost more. The result is that research data are not always backed up with the same rigor as administrative data. Often we see that research computing home directories and some network shares are backed up but that the bulk of the data are not backed up or are backed up locally by individual researchers because IT cannot do it.

One can say that by identifying research as "different" to justify bypassing administrative computing restraints, we are partly responsible for this current situation. Well, research is different, but in a big data world—where the data are our greatest asset—we all need to work differently. Funding agencies are recognizing this and are including regulatory requirements in grants for data management. (This will be explored further in the Curation paper of this series.) There are many issues to consider that are both technology and process based. What should be backed up? Who should back up the data? How should they be backed up? How are they made available to other investigators? For how long? What uses are best served by private clouds and which ones by public cloud services? Developing a campus strategy for big data backup and recovery requires an interdisciplinary team (i.e., IT, library, research, and finance) to design and deliver.

## Findings

- Typical campus IT services are ill equipped to effectively back up and archive research big data due to their volume, complexity, velocity, and variety.

- The issue is not just about technology. Policy and regulatory requirements must be addressed as well.

## Recommendations

- Enterprise backup and disaster recovery strategies should include plans for research big data because they are an important university asset.

- Universities should create multidisciplinary groups—including IT, research IT, researchers, library, finance, and policy—to develop integrated strategies to protect research big data.

- EDUCAUSE should support universities by organizing workshops and resource libraries for documenting and sharing good practices.

# Campus Bridging

"Campus bridging" refers to enabling scientists to use the cyberinfrastructure resources of their lab, their campus, other campuses, and regional/national (e.g., NSF or DOE) centers in a seamlessly integrated way.[6] Accomplishing this enhances the effectiveness of the scientist's personal research and, particularly, of collaborative research involving colleagues scattered geographically but united (perhaps over time or for a brief period) in a common intellectual pursuit. Big data are relevant to campus bridging because large, complex, and/or rapidly changing data sets are often used in research, including in cases where participants in the research are scattered and/or where the cyberinfrastructure resources needed to work on those data sets may themselves be geographically scattered. We note that if campus bridging is successfully achieved, then the prospects for our scientists to engage in cutting-edge collaborative research are significantly enhanced. Thus, there is a strong incentive for each university to ensure that campus bridging becomes a successful and even routine experience for our faculty, students, and staff.

Anecdotally, many obstacles stand in the way, including inadequate network performance, lack of data and metadata standards, and lack of interoperability in how data management is structured at different sites. Thus, in a sense, successful campus bridging requires achieving success in each of the three areas that we've already touched on in this paper:

1. Success in *network infrastructure* is needed to make possible the rapid sharing of the big data objects that are often used during the course of research, particularly in collaborative research. Similarly, success in networking infrastructure is needed to support the nimble use by campus researchers of local and remote computing resources needed for use with big data objects.

2. Success in *server hosting strategies* is needed to support the local computing and storage resources that are necessary to complement remote cyberinfrastructure resources. Otherwise, researchers depend solely on remote resources, thus weakening their impact on scientific communities they care about while also weakening the real and perceived attractiveness of the local campus as a great place to be to do research. Stated positively, successful server hosting

strategies permit local resources to be interesting complements to national and other remote resources.

3. Success in *backup and disaster recovery* is needed to ensure that data-intensive research—whether conducted with local, remote, or combined resources—can be relied on over a long period of time.

The relationship, however, is bidirectional: Not only does campus bridging require success in these three areas, but it also motivates and in some cases becomes part of the solution. For example, the Compact Muon Solenoid (CMS) high-energy physics community has worked to store its high-value big data objects in a geographically distributed way.[7] By solving problems of metadata, cataloging, and high-speed data movement between geographically distributed repositories, CMS enables a university's scientists to use local and remote compute resources to access and create these data sets, thereby enhancing the value of local storage and compute resources. It also helps solve a difficult problem in backup and recovery in that the cataloging and data-movement parts of the system ensure that key data sets always have multiple geographically distributed replicas.

Fundamental to success in campus bridging is organizing campus support staff to identify and provide helpful support to campus researchers who use national resources such as XSEDE and DoE National Laboratory resources. These staff should be capable of both facilitating the effective use of these remote resources and understanding how relevant local infrastructure does or could strengthen research efforts. In some cases, this might mean removing network bottlenecks. In other cases, it might mean identifying local compute, storage, networking, or visualization resources that could dovetail with remote resources to help campus researchers accomplish things not practical with only local or only remote resources.

## Findings

- Most research universities have researchers who benefit from powerful remote compute, storage, and/or data resources. Often, there is only partial understanding between local IT leadership and these researchers.

- Local IT leadership can break down barriers to the effective use of these remote resources.

- More interestingly, if more complete understanding emerges between IT leadership and researchers, opportunities to combine the use of local and remote resources can be identified.

## Recommendations

- Every research university should organize a means to engage in national resources such as the XSEDE Campus Champions program,[8] including the strengthening of the understanding of the needs of local researchers and the strengths and weaknesses of local campus infrastructure.

- Every research university should organize parallel efforts targeted at the needs of researchers that use other remote resources such as those of the DoE National Labs.

- Where practical, local campus cyberinfrastructure should be designed to dovetail with these remote resources.

# Conclusion

By its nature, research computing is on the leading edge of many technology innovations. The IT needs of campus research eventually become the norm for other area computing needs. For example, we are beginning to see big data approaches applied to student recruitment, retention, and success. Addressing these campus cyberinfrastructure issues now as a campus strategic initiative for big data research will likely pay off for all campus data in the long run.

## Authors

**Guy T. Almes**
Director, Academy for Advanced
   Telecommunications and Learning Technologies
Texas A&M University

**Ralph J. Zottola**
Chief Technology Officer, Research Computing
University of Massachusetts Central Office

## Notes

1. XSEDE is "a single virtual system that scientists can use to interactively share computing resources, data, and expertise" that allows researchers to establish "private, secure environments" that enable collaboration and provide access to computing resources and services.

2. The Science DMZ was pioneered by network engineers at ESnet and is best described in "Science DMZ: A Scalable Network Design Model for Optimizing Science Data Transfers."

3. For more on software-defined networking and the current status of SDN implementations, see the upcoming brief from the ECAR Campus Cyberinfrastructure (ECAR-CCI) and Communications Infrastructure and Applications (ECAR-CIA) working groups. This document, *The Promise and Reality of SDN*, aims to help clear up some of the confusion, explain specific reasons why cyberinfrastructure leaders on our campuses might want to take an interest in SDN, and what campuses can do to prepare for SDN in the near term. Available from https://www.educause.edu/ecar/ecar-working-groups.

4. See, for example, how two universities have satisfied their due-diligence requirements using different approaches, thereby providing needed services to their research communities: William Barnett, Robert Flynn, and Anurag Shankar, "Bringing Box into HIPAA Alignment," and Ruth Marinshaw, "Stanford Medicine Box Discussion" (both presented at the fall 2014 Coalition for Advanced Scientific Computing meeting, Arlington, D.C.). The first presentation was also given at the 2014 Internet2 Global Summit and is available as a netcast.

5. Learn more about the Massachusetts Open Cloud.

6. To learn more about the role of the campus cyberinfrastructure in campus bridging, see the NSF Advisory Committee for Cyberinfrastructure Task Force on Campus Bridging, *Final Report*, March 2011, and the 2012 response to that report from this working group, "What's Next for Campus Cyberinfrastructure? ACTI Responds to the NSF ACCI Reports."

7. Learn more about the Compact Muon Solenoid, a particle detector in CERN's Large Hadron Collider. Information about how the data have been managed can be found at J. Adelman-McCarthy et al., "CMS Computing Operations During Run 1," *Journal of Physics: Conference Series* 513, track 3 (2014).

8. The XSEDE Campus Champions program "supports campus representatives as a local source of knowledge about high-performance and high-throughput computing and other digital services, opportunities, and resources."