

Lawrence Berkeley National Laboratory

LBL Publications

Title

Arizona State University Requirements Analysis Report

Permalink

<https://escholarship.org/uc/item/03r3h2d5>

Authors

Zurawski, Jason

Southworth, Douglas

Meade, Brenna

Publication Date

2022-03-04

Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial-NoDerivatives License, available at <https://creativecommons.org/licenses/by-nc-nd/4.0/>

Peer reviewed



Arizona State University Requirements Analysis Report

March 4, 2022



U.S. DEPARTMENT OF
ENERGY
Office of Science



ESnet

ENERGY SCIENCES NETWORK



INDIANA UNIVERSITY

Disclaimer

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor the Trustees of Indiana University, nor any of their employees, makes any warranty, express or implied, or assumes any legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California or the Trustees of Indiana University. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California, or the Trustees of Indiana University.

Arizona State University Requirements Analysis Report

March 4, 2022

The Engagement and Performance Operations Center (EPOC) is supported by the National Science Foundation under Grant No. 1826994.

ESnet is funded by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research. Carol Hawk is the ESnet Program Manager.

ESnet is operated by Lawrence Berkeley National Laboratory, which is operated by the University of California for the U.S. Department of Energy under contract DE-AC02-05CH11231.

This is a University of California, Publication Management System report number LBNL-2001502¹.

¹<https://escholarship.org/uc/item/03r3h2d5>

Participants & Contributors

Rebecca Belshe, Arizona State University
Marisa Brazil, Arizona State University
Alan Chapman, Arizona State University
Po-Lin Chiu, Arizona State University
William Dizon, Arizona State University
Douglas Jennewein, Arizona State University
Chris Kurtz, Arizona State University
John McCutcheon, Arizona State University
Johnathan Lee, Arizona State University
Natalie Mason, Arizona State University
Brenna Meade, Indiana University
Ken Miller, ESnet
Anna Muldoon, Arizona State University
Kristy Roschke, Arizona State University
Ian Shaeffer, Arizona State University
Douglas Shepherd, Arizona State University
Michael Simeone, Arizona State University
Doug Southworth, Indiana University
Gil Speyer, Arizona State University
Philip Tarrant, Arizona State University
Thomas Taylor, Arizona State University
Kevin Tinnin, Arizona State University
Shawn Walker, Arizona State University
Dewight Williams, Arizona State University
Jason Yalim, Arizona State University
Jason Zurawski, ESnet

Report Editors

Brenna Meade, Indiana University: Brenna Meade meadeb@iu.edu
Doug Southworth, Indiana University: dojosout@iu.edu
Jason Zurawski, ESnet: zurawski@es.net

Contents

1 Executive Summary	9
Deep Dive Review Purpose and Process	9
This Review	9
This review includes case studies from the following campus stakeholder groups:	9
The review produced several important findings and recommendations from the case studies and subsequent virtual conversations:	11
2 Deep Dive Findings & Recommendations	12
2.1 Findings	12
2.2 Recommendations	13
3 Process Overview and Summary	17
3.1 Campus-Wide Deep Dive Background	17
3.2 Campus-Wide Deep Dive Structure	18
3.3 Arizona State University Deep Dive Background	20
3.4 Organizations Involved	21
4 Arizona State University Case Studies	22
4.1 Departments of Biophysics and Biochemistry: Structural Biology Research	23
4.1.1 Use Case Summary	23
4.1.2 Collaboration Space	23
4.1.3 Instruments & Facilities	23
4.1.4 Data Narrative	23
4.1.4.1 Data Volume & Frequency Analysis	23
4.1.4.2 Data Sensitivity	23
4.1.4.3 Future Data Volume & Frequency Analysis	23
4.1.5 Technology Support	24
4.1.5.1 Software Infrastructure	24
4.1.5.2 Network Infrastructure	24
4.1.5.3 Computation and Storage Infrastructure	24
4.1.5.4 Data Transfer Capabilities	24
4.1.6 Internal & External Funding Sources	25
4.1.7 Resource Constraints	25
4.1.8 Ideal Data Architecture	25
4.1.9 Outstanding Issues	25
	5

4.2 School of Life Sciences: Cellular Mechanisms of Evolution	26
4.2.1 Use Case Summary	26
4.2.2 Collaboration Space	26
4.2.3 Instruments & Facilities	26
4.2.4 Data Narrative	26
4.2.4.1 Data Volume & Frequency Analysis	26
4.2.4.2 Data Sensitivity	26
4.2.4.3 Future Data Volume & Frequency Analysis	26
4.2.5 Technology Support	27
4.2.5.1 Software Infrastructure	27
4.2.5.2 Network Infrastructure	27
4.2.5.3 Computation and Storage Infrastructure	27
4.2.5.4 Data Transfer Capabilities	27
4.2.6 Internal & External Funding Sources	27
4.2.7 Resource Constraints	27
4.2.8 Ideal Data Architecture	27
4.2.9 Outstanding Issues	28
4.3 Department of Physics: Optical Biophysics and Spatial Transcriptomics	29
4.3.1 Use Case Summary	29
4.3.2 Collaboration Space	29
4.3.3 Instruments & Facilities	29
4.3.4 Data Narrative	30
4.3.4.1 Data Volume & Frequency Analysis	30
4.3.4.2 Data Sensitivity	30
4.3.4.3 Future Data Volume & Frequency Analysis	30
4.3.5 Technology Support	31
4.3.5.1 Software Infrastructure	31
4.3.5.2 Network Infrastructure	31
4.3.5.3 Computation and Storage Infrastructure	31
4.3.5.4 Data Transfer Capabilities	31
4.3.6 Internal & External Funding Sources	31
4.3.7 Resource Constraints	31
4.3.8 Ideal Data Architecture	31
4.3.9 Outstanding Issues	32

4.4 School of Molecular Sciences: Structural Biophysics and Cryo-EM	33
4.4.1 Use Case Summary	33
4.4.2 Collaboration Space	33
4.4.3 Instruments & Facilities	33
4.4.4 Data Narrative	33
4.4.4.1 Data Volume & Frequency Analysis	33
4.4.4.2 Data Sensitivity	33
4.4.4.3 Future Data Volume & Frequency Analysis	33
4.4.5 Technology Support	33
4.4.5.1 Software Infrastructure	33
4.4.5.2 Network Infrastructure	33
4.4.5.3 Computation and Storage Infrastructure	33
4.4.5.4 Data Transfer Capabilities	34
4.4.6 Internal & External Funding Sources	34
4.4.7 Resource Constraints	34
4.4.8 Ideal Data Architecture	34
4.4.9 Outstanding Issues	34
4.5 School of Social and Behavioral Sciences: Combating Mis/Disinformation via Critical Data Studies	35
4.5.1 Use Case Summary	35
4.5.2 Collaboration Space	35
4.5.3 Instruments & Facilities	36
4.5.4 Data Narrative	36
4.5.4.1 Data Volume & Frequency Analysis	38
4.5.4.2 Data Sensitivity	38
4.5.4.3 Future Data Volume & Frequency Analysis	38
4.5.5 Technology Support	38
4.5.5.1 Software Infrastructure	38
4.5.5.2 Network Infrastructure	39
4.5.5.3 Computation and Storage Infrastructure	39
4.5.5.4 Data Transfer Capabilities	39
4.5.6 Internal & External Funding Sources	39
4.5.7 Resource Constraints	39
4.5.8 Ideal Data Architecture	40

4.5.9 Outstanding Issues	40
4.6 Arizona State University Research Technology Office & Research Computing	41
4.6.1 Use Case Summary	41
4.6.2 Collaboration Space	41
4.6.3 Capabilities & Special Facilities	41
4.6.4 Technology Narrative	41
4.6.4.1 Network Infrastructure	41
4.6.4.2 Computation and Storage Infrastructure	44
4.6.4.3 Network & Information Security	45
4.6.4.4 Monitoring Infrastructure	45
4.6.4.5 Software Infrastructure	45
4.6.5 Organizational Structures & Engagement Strategies	46
4.6.5.1 Organizational Structure	46
4.6.5.2 Engagement Strategies	46
4.6.6 Internal & External Funding Sources	47
4.6.7 Resource Constraints	47
4.6.8 Outstanding Issues	47
Appendix A - Research Computing Facilities Statement	48
Personnel	48
Advanced Computing	48
Advanced Computing Systems	49
Dell Center of Excellence for HPC and Artificial Intelligence	49
Open Science Grid	50
Data Center	50
Network	51
Science DMZ	51
Data Storage	52

1 Executive Summary

Deep Dive Review Purpose and Process

EPOC uses the Deep Dive process to discuss and analyze current and planned science, research, or education activities and the anticipated data output of a particular use case, site, or project to help inform the strategic planning of a campus or regional networking environment. This includes understanding future needs related to network operations, network capacity upgrades, and other technological service investments. A Deep Dive comprehensively surveys major research stakeholders' plans and processes in order to investigate data management requirements over the next 5–10 years. Questions crafted to explore this space include the following:

- How, and where, will new data be analyzed and used?
- How will the process of doing science change over the next 5–10 years?
- How will changes to the underlying hardware and software technologies influence scientific discovery?

Deep Dives help ensure that key stakeholders have a common understanding of the issues and the actions that a campus or regional network may need to undertake to offer solutions. The EPOC team leads the effort and relies on collaboration with the hosting site or network, and other affiliated entities that participate in the process. EPOC organizes, convenes, executes, and shares the outcomes of the review with all stakeholders.

This Review

Between October 2021 and February 2022 staff members from the Engagement and Performance Operations Center (EPOC) met with researchers and staff from Arizona State University (ASU) for the purpose of a Deep Dive into scientific and research drivers. The goal of this activity was to help characterize the requirements for a number of campus use cases, and to enable cyberinfrastructure support staff to better understand the needs of the researchers within the community.

This review includes case studies from the following campus stakeholder groups:

- Departments of Biophysics and Biochemistry: Structural Biology Research
- School of Life Sciences: Cellular Mechanisms of Evolution
- Department of Physics: Optical Biophysics and Spatial Transcriptomics
- School of Molecular Sciences: Structural Biophysics and Cryo-EM
- School of Social and Behavioral Sciences: Combating Mis/Disinformation via Critical Data Studies
- Arizona State University Research Technology Office & Research Computing

Material for this event included the written documentation from each of the profiled research areas, documentation about the current state of technology support, and a write-up of the discussion that took place via e-mail and video conferencing.

The case studies highlighted the ongoing challenges and opportunities that ASU has in supporting a cross-section of established and emerging research use cases. Each case study mentioned unique challenges which were summarized into common needs.

The review produced several important findings and recommendations from the case studies and subsequent virtual conversations:

- The ASU campus has invested considerable time, effort, and funding to develop the Science DMZ architecture used by the research community.
- Data storage is a critical part of the research workflow. A number of research use cases have noted that they will require access to more storage capacity (e.g. multiple TBs, approaching PB) on an ongoing basis in the years to come.
- External collaborators to ASU research staff have the ability to request access to ASU resources, but can run into problems due to the consistency of how access is managed, and the length of time required to complete the review and account creation process.
- ASU Research Technology Office (RTO) provides a number of hardware and software support approaches that enhance and facilitate ASU research use cases, and has had success in building technology support solutions at multiple layers (networking, computation, storage, and data mobility software).
- To scale the ASU RTO services for future support requirements, it is recommended that ASU increase the ASU support team size, services offered, and available documentation to better scale the engagement, integration, and operational services that are currently available.
- It is recommended that RTO technology service requests be integrated with the Research Advancement support office, to better identify technology needs and use cases at the time of application, and work with teams that are funded.
- It is recommended that ASU RTO pursue a holistic strategy to addressing the data storage needs of campus researchers through several critical upgrades and service offerings of different sizes and speeds to integrate in with the research data lifecycle.
- It is recommended that ASU RTO perform regular reviews of usage, usage cases, and technology that use the Science DMZ infrastructure. This should be coupled with a routine review of CI financial and sustainability approaches.
- The RTO data mobility infrastructure should be upgraded to integrate other Globus endpoints, as well as adopting approaches to portal applications for some users.

2 Deep Dive Findings & Recommendations

The deep dive process helps to identify important facts and opportunities from the profiled use cases. The following outlines a set of findings and recommendations from the ASU Deep Dive that summarize important information gathered during the discussions surrounding case studies, and possible ways that could improve the CI support posture for the campus.

2.1 Findings

- The ASU campus has invested considerable time, effort, and funding to develop the Science DMZ architecture used by the research community. This design is fully featured, and has shown success for early use cases that have adopted it.
- Data storage is a critical part of the research workflow, and several technical factors matter significantly when gauging how it will integrate with use cases: storage size, access speed for reads and writes, and locality to the users. A number of research use cases identified through this activity have noted that they will require access to more storage capacity (e.g. multiple TBs, approaching PB) on an ongoing basis in the years to come; compounding this are:
 - Technical requirements on the aforementioned storage to have low latency (e.g. located on premises)
 - Maintain access times that can match base to instrumentation and computation
 - Concerns that when cloud-based storage is utilized, it may not be as permanent as one would expect (e.g., change of vendors which could force a migration of data between clouds)
- The ASU network and information security policy on the Science DMZ is stable for the early use cases and users, and receives regular reviews to ensure that it is delivering on stated goals and expectations.
- External collaborators to ASU research staff have the ability to request access to ASU resources via their ASURITE user IDs, or apply for other forms of university affiliation. This helps to facilitate collaborative activities. This process is not without downsides:
 - The speed of acquiring this access is not consistent, however, and it can take an extended period of time to grant access which harms research progress.
 - The amount of information required to gain access is not consistent with the research use case, e.g., treating temporal users the same as a permanent student or staff member. A middle ground that enables access, without requiring as much personal information, is desirable.

- ASU RTO provides a number of hardware and software support approaches that enhance and facilitate ASU research use cases. Without these resources, many groups would not have the expertise or ability to deliver on funded milestones.
- ASU RTO has had success in building technology support solutions at multiple layers (networking, computation, storage, and data mobility software) for a number of use cases, and continues to investigate ways they can support existing and emerging requirements.
- The ASU RTO environment supports a wide range of research with computation, storage, networking, software, and engagement support. To date this commitment can be seen in the number of users, and successful research outcomes, that are using RTO facilities instead of directly operating bespoke infrastructure. The dedicated RTO staff that are available to manage these activities is necessary, and often quite rare in the university system. This commitment is commendable and the successful outcomes justify the investment.
- ASU RTO currently maintains a set of hardware and software to support data mobility activities (e.g. sharing with collaborators, transfer to national and international facilities, etc.). A dedicated data transfer capability, consisting of a number of DTNs running Globus software, with a subscription to critical services that enable secure management of sensitive data and connectivity to cloud environments, forms the backbone of this support. This capability is critical for modern research use cases, and will see increased use in the coming years.

2.2 Recommendations

- To scale the ASU RTO services for future support requirements, it is recommend to consider the following:
 - Increase the ASU support team size to better scale the engagement, integration, and operational services that are currently available.
 - Increase the services that are offered, and devices ways to advertise and champion these to existing, and potential, users through a structured portal/request system.
 - Improve the quality and quantity of service documentation, to better scale the number of available staff in RTO.
 - Continue to work with ASU KE Communications to create a strategic communication plan which may include updating documentation given to new and potential faculty, and working at researcher-facing events to raise awareness about and explain RTO service offerings
 - Working to integrate RTO technology service requests with the Research Advancement support office, to better identify technology needs and use cases at the time of application, and work with teams that are funded.

- Continuing to partner with affiliated ASU IT organizations to offer services, and centralize support for complicated use cases that require alternative arrangements for networking, computation, storage, or software support.
- It is recommended that ASU RTO pursue a holistic strategy to addressing the data storage needs of campus researchers through several critical upgrades and service offerings:
 - Create a “tiered” approach to campus storage that recognizes 4 key layers related to the research lifecycle:
 - “Small and Fast” storage, provided through upgraded and expanded Data Transfer Nodes. These will serve as critical “data capacitors” that will share data both within, and external to ASU.
 - “Medium and Average” storage, provided through the existing RTO Storage Area Network. These resources form a critical backbone for the university, accepting data from research groups that require access to computation, as well as those that need a location to store valuable research data sets and lack a local storage solution.
 - “Large and Slow” storage, in the form of a new archival system for long-term storage requirements. This system will remove unused, but still valuable, data sets from the previous layer freeing up valuable space, and still offering assurance and security for use cases that are required by grant guidelines to main research data after projects end.
 - “Large, Slow, and Remote”, via existing relationships with cloud providers. This off-site storage can be a critical backup to the on-campus archival system, and can be leveraged as a last resort due to the slow access speeds and times.
- It is recommended that ASU RTO perform regular reviews of usage, usage cases, and technology that use the Science DMZ infrastructure. Through this evaluation process, the Science DMZ will remain secure, performance oriented, and useful to the ASU research community.
- It is recommended that ASU RTO work with EPOC, and TrustedCI, to evaluate the information and network security readiness of the Science DMZ infrastructure. This activity will help to make critical choices regarding exposed services, and expansion possibilities. A special focus on preparation for long-term management of sensitive data (e.g. HIPPA or other forms of CUI controls) is being performed correctly for the ASU research community.
- ASU RTO should investigate alternative methods to grant access to a subset of research technology components that does not require the full access, and

often cumbersome application process, for ASURITE accounting. A number of users have identified more narrow use cases (e.g. sharing a data set on a non-routine basis) that could benefit from a light-weight way to exchange account information, and access resources that may reside in the Science DMZ.

- It is recommended that ASU RTO schedule routine CI financial and sustainability meetings to better understand the ongoing costs related to supporting campus science. ASU RTO will begin to catalog equipment and service cost ratios to the overall institutional research awards, expenditures, and F&A in order to provide overall CI costs. As ASU RTO continues the process of researcher engagement potential use cases requiring additional CI resources can be addressed and planned for. This effort would consist of a regular reviewing and planning of:
 - Current CI services, and their costs
 - Potential CI services, and their costs
 - Review of equipment lifecycle
 - Review of budgets (near and long term)
 - Personnel requirements
 - CI grant opportunities

- It is recommended that ASU RTO explore adding different classes of operational assurance to the research community on the technology hardware and software they support. This comes from the observation that certain research activities can occasionally reach a level of maturity that requires more stable operational support. For instance, a project that routinely accesses remote network resources, and runs continuously, can be harmed if there is scheduled or unscheduled downtime of the underlying technical stack. It is recommended that RTO consider implementing a second class of service that offers higher levels of operational assurance for these services that are moving toward production.

- It is recommended that ASU RTO continue to work with the research community to incorporate additional use cases (e.g. facilities, laboratories, instruments, experiments, etc.) that could benefit from core service offerings in the form of pilot workflow activities. These may consist of creating direct access to the Science DMZ infrastructure to remote portions of campus by extending the network, integrating centrally managed computational and storage components, helping with scientific software, or providing other critical research services to advance research outcomes.

- The RTO data mobility infrastructure should be upgraded to:
 - Integrate other Globus endpoints that are not under the current subscription to Globus, which will allow a transfer of benefits and

services (e.g. single sign on, the ability to make and manage groups, etc.).

- Expand DTN capabilities in the form of more servers as the demand for data mobility increases across campus.
- Consider adopting components of the Globus Modern Research Data Portal ²(MRDP) to facilitate a more user friendly experience to browsing research data sets
- Participate in the EPOC/ESnet Data Mobility Exhibition³ (DME) measure and improve data transfer performance.

² <https://docs.globus.org/modern-research-data-portal/>

³ <https://fasterdata.es.net/performance-testing/2019-2020-data-mobility-workshop-and-exhibition/2019-2020-data-mobility-exhibition/>

3 Process Overview and Summary

3.1 Campus-Wide Deep Dive Background

Over the last decade, the scientific community has experienced an unprecedented shift in the way research is performed and how discoveries are made. Highly sophisticated experimental instruments are creating massive datasets for diverse scientific communities and hold the potential for new insights that will have long-lasting impacts on society. However, scientists cannot make effective use of this data if they are unable to move, store, and analyze it.

The Engagement and Performance Operations Center (EPOC) uses the Deep Dives process as an essential tool as part of a holistic approach to understand end-to-end research data use. By considering the full end-to-end research data movement pipeline, EPOC is uniquely able to support collaborative science, allowing researchers to make the most effective use of shared data, computing, and storage resources to accelerate the discovery process.

EPOC supports five main activities

- Roadside Assistance via a coordinated Operations Center to resolve network performance problems with end-to-end data transfers reactively;
- Application Deep Dives to work more closely with application communities to understand full workflows for diverse research teams in order to evaluate bottlenecks and potential capacity issues;
- Network Analysis enabled by the NetSage monitoring suite to proactively discover and resolve performance issues;
- Provision of managed services via support through the Indiana University (IU) GlobalNOC and our Regional Network Partners; and
- Coordinated Training to ensure effective use of network tools and science support.

Whereas the Roadside Assistance portion of EPOC can be likened to calling someone for help when a car breaks down, the Deep Dive process offers an opportunity for broader understanding of the longer term needs of a researcher. The Deep Dive process aims to understand the full science pipeline for research teams and suggest alternative approaches for the scientists, local IT support, and national networking partners as relevant to achieve the long-term research goals via workflow analysis, storage/computational tuning, identification of network bottlenecks, etc.

The Deep Dive process is based on an almost 15-year practice used by ESnet to understand the growth requirements of Department of Energy (DOE) facilities⁴. The EPOC team adapted this approach to work with individual science groups through a set of structured data-centric conversations and questionnaires.

⁴ <https://fasterdata.es.net/science-dmz/science-and-network-requirements-review>

3.2 Campus-Wide Deep Dive Structure

The Deep Dive process involves structured conversations between a research group and relevant IT professionals to understand at a broad level the goals of the research team and how their infrastructure needs are changing over time.

The researcher team representatives are asked to communicate and document their requirements in a case-study format that includes a data-centric narrative describing the science, instruments, and facilities currently used or anticipated for future programs; the advanced technology services needed; and how they can be used. Participants considered three timescales on the topics enumerated below: the near-term (immediately and up to two years in the future); the medium-term (two to five years in the future); and the long-term (greater than five years in the future).

The case study process tries to answer essential questions about the following aspects of a workflow:

- **Research & Scientific Background**—an overview description of the site, facility, or collaboration described in the Case Study.
- **Collaborators**—a list or description of key collaborators for the science or facility described in the Case Study (the list need not be exhaustive).
- **Instruments and Facilities: Local & Non-Local**—a description of the network, compute, instruments, and storage resources used for the science collaboration/program/project, or a description of the resources made available to the facility users, or resources that users deploy at the facility or use at partner facilities.
- **Process of Science**—a description of the way the instruments and facilities are used for knowledge discovery. Examples might include workflows, data analysis, data reduction, integration of experimental data with simulation data, etc.
- **Computation & Storage Infrastructure: Local & Non-Local**—The infrastructure that is used to support analysis of research workflow needs: this may be local storage and computation, it may be private, it may be shared, or it may be public (commercial or non—commercial).
- **Software Infrastructure**—a discussion focused on the software used in daily activities of the scientific process including tools that are used locally or remotely to manage data resources, facilitate the transfer of data sets from or to remote collaborators, or process the raw results into final and intermediate formats.
- **Network and Data Architecture**—description of the network and/or data architecture for the science or facility. This is meant to understand how data moves in and out of the facility or laboratory focusing on local infrastructure configuration, bandwidth speed(s), hardware, etc.
- **Resource Constraints**—non-exhaustive list of factors (external or internal) that will constrain scientific progress. This can be related to funding, personnel, technology, or process.
- **Outstanding Issues**—Listing of any additional problems, questions, concerns, or comments not addressed in the aforementioned sections.

At a physical or virtual meeting, this documentation is walked through with the research team (and usually cyberinfrastructure or IT representatives for the organization or region), and an additional discussion takes place that may range beyond the scope of the original document. At the end of the interaction with the research team, the goal is to ensure that EPOC and the associated CI/IT staff have a solid understanding of the research, data movement, who's using what pieces, dependencies, and time frames involved in the Case Study, as well as additional related cyberinfrastructure needs and concerns at the organization. This enables the teams to identify possible bottlenecks or areas that may not scale in the coming years, and to pair research teams with existing resources that can be leveraged to more effectively reach their goals.

3.3 Arizona State University Deep Dive Background

Between October 2021 and February 2022, EPOC organized a Deep Dive in collaboration with ASU to characterize the requirements for several key science drivers. The representatives from each use case were asked to communicate and document their requirements in a case-study format. These included:

- Departments of Biophysics and Biochemistry: Structural Biology Research
- School of Life Sciences: Cellular Mechanisms of Evolution
- Department of Physics: Optical Biophysics and Spatial Transcriptomics
- School of Molecular Sciences: Structural Biophysics and Cryo-EM
- School of Social and Behavioral Sciences: Combating Mis/Disinformation via Critical Data Studies
- Arizona State University Research Technology Office & Research Computing

3.4 Organizations Involved

The Engagement and Performance Operations Center (EPOC) was established in 2018 as a collaborative focal point for operational expertise and analysis and is jointly led by Indiana University (IU) and the Energy Sciences Network (ESnet). EPOC provides researchers with a holistic set of tools and services needed to debug performance issues and enable reliable and robust data transfers. By considering the full end-to-end data movement pipeline, EPOC is uniquely able to support collaborative science, allowing researchers to make the most effective use of shared data, computing, and storage resources to accelerate the discovery process.

The Energy Sciences Network (ESnet) is the primary provider of network connectivity for the U.S. Department of Energy (DOE) Office of Science (SC), the single largest supporter of basic research in the physical sciences in the United States. In support of the Office of Science programs, ESnet regularly updates and refreshes its understanding of the networking requirements of the instruments, facilities, scientists, and science programs that it serves. This focus has helped ESnet to be a highly successful enabler of scientific discovery for over 25 years.

Indiana University (IU) was founded in 1820 and is one of the state's leading research and educational institutions. Indiana University includes two main research campuses and six regional (primarily teaching) campuses. The Indiana University Office of the Vice President for Information Technology (OVPIT) and University Information Technology Services (UITS) are responsible for delivery of core information technology and cyberinfrastructure services and support.

Arizona State University (ASU) is a public research university, founded in 1885. ASU is a member of the Universities Research Association and classified among "R1: Doctoral Universities – Very High Research Activity". ASU has nearly 150,000 students attending classes, with more than 38,000 students attending online, and 90,000 undergraduates and more nearly 20,000 postgraduates across its five campuses and four regional learning centers throughout Arizona. ASU offers 350 degree options from its 17 colleges and more than 170 cross-discipline centers and institutes for undergraduate students, as well as more than 400 graduate degree and certificate programs.

4 Arizona State University Case Studies

ASU presented a number use cases during this review. These are as follows:

- Departments of Biophysics and Biochemistry: Structural Biology Research
- School of Life Sciences: Cellular Mechanisms of Evolution
- Department of Physics: Optical Biophysics and Spatial Transcriptomics
- School of Molecular Sciences: Structural Biophysics and Cryo-EM
- School of Social and Behavioral Sciences: Combating Mis/Disinformation via Critical Data Studies
- Arizona State University Research Technology Office & Research Computing

Each of these Case Studies provides a glance at research activities, the use of experimental methods and devices, the reliance on technology, and the scope of collaborations. It is important to note that these views are primarily limited to current needs, with only occasional views into the event horizon for specific projects and needs into the future. Estimates on data volumes, technology needs, and external drivers are discussed where relevant.

4.1 Departments of Biophysics and Biochemistry: Structural Biology Research

Content in this section authored by Dewight Williams, Director of Electron Microscopy Resources Laboratory for the Biophysics and Biochemistry Department

4.1.1 Use Case Summary

Cryogenic Transmission Electron Microscopes (TEM) are used to produce images of biological material in a frozen hydrated state. Subsequent images are processed computationally, to provide three dimensional volumes of these macromolecules or cellular components at near atomic resolution such that the assembly state of biomolecules in living systems can be determined. Nearly all biological researchers and even some material scientists working on organic-inorganic assemblies have become stake holders to the use of this technology at ASU.

4.1.2 Collaboration Space

Po-Lin Chiu, from ASU, is a primary collaborator for this work. Many other ASU professors and research staff leverage the facilities that are provided.

4.1.3 Instruments & Facilities

This use case leverages ASU resources primarily:

- Krios cryoTEM
- Titan eTEM
- SCOB annex with an 40Gbps fiber path that is now available

4.1.4 Data Narrative

3000-5000 images are collected daily on the cryoTEM using the direct electron detector. This data is generally about 1-2 TB in size and needs to be moved to either laboratory computers or onto a computational cluster for image processing and three dimensional structure determination. The instruments are operated for about $\frac{3}{4}$ of a month on average.

Data is retained for approximately a year, as space allows. Currently have the ability to store up to 200TBs, and on average 20TB of data is available for use. Data is migrated or deleted when space grows small.

4.1.4.1 Data Volume & Frequency Analysis

The use case produces Terabytes (TB) amounts of data volume, on a daily frequency.

4.1.4.2 Data Sensitivity

There are no sensitive aspects to the use case's data to report.

4.1.4.3 Future Data Volume & Frequency Analysis

It is expected that the future versions of this use case will produce Terabyte (TB) amounts of data volume, on a daily frequency. Increases in data output from instruments will increase this number, and adding more instruments will also result in a higher frequency of use.

4.1.5 Technology Support

4.1.5.1 Software Infrastructure

The following software packages are used during this research:

- Relion (REGularised Likelihood Optimisation)⁵: software package that employs an empirical Bayesian approach for electron cryo-EM structure determination
- Cryosparc⁶: scientific software platform for cryo-EM used in research and drug discovery pipelines
- IMOD⁷: set of image processing, modeling and display programs used for tomographic reconstruction and for 3D reconstruction of EM serial sections and optical sections
- Cistem⁸: software to process cryo-EM images of macromolecular complexes and obtain high-resolution 3D reconstructions from them
- MDFF⁹: provides commands for setting up and analyzing molecular dynamics flexible fitting simulations, i.e., simulations in which an atomic structure is flexibly fitted into a density map
- CCP4¹⁰: suite of programs that allows researchers to determine macromolecular structures by X-ray crystallography, and other biophysical techniques
- Phenix¹¹: software package for macromolecular structure determination using crystallographic (X-ray, neutron and electron) and electron cryo-microscopy data

4.1.5.2 Network Infrastructure

This facility features 10Gbps to the workstations that are involved in data collection, and that connects to a 40Gbps-capable switch to that links to ASU RTS resources.

4.1.5.3 Computation and Storage Infrastructure

Currently most data are processed on individual workstations. Storage in facility limited to 200TB, and in individual labs to 30TB.

Future operations would benefit to use locally available HPC resources at RTS, or within the sequencing facility.

4.1.5.4 Data Transfer Capabilities

Typically data sets in the 2-TB range are moved from the data acquisition machines via external media. Alternatively, rsync across university network has been used. Both are successful but slower than desired for the amount of data that is used during this research.

⁵ <https://relion.readthedocs.io/en/release-4.0/>

⁶ <https://cryosparc.com>

⁷ <https://bio3d.colorado.edu/imod/>

⁸ <https://cistem.org>

⁹ <https://www.ks.uiuc.edu/Research/vmd/plugins/mdff/>

¹⁰ <https://www.ccp4.ac.uk>

¹¹ <https://phenix-online.org>

There have been historical issues with data transfer both on and off campus:

- The data transfer rate for Agave (ASU HPC cluster) are a major bottleneck – particularly when sending of the /scratch filesystem, which is used on the compute resources
- Data transfer to external collaborators is also slow, but this may be because the endpoint isn't optimized for wide area transfer

The facility does try to use Globus when possible, but does not have a paid subscription currently. As a result, users that have access to globus at their home institution will “pull” data when they can, others will resort to using rsync or removable media, and still others may use cloud connections (e.g., Dropbox, google drive). ASU would like to pursue having a unified Globus subscription for all endpoints.

4.1.6 Internal & External Funding Sources

This research has external funding from:

- Department of Defense (e.g. US Army)
- National Institutes of Health (NIH)
- National Science Foundation (NSF)

4.1.7 Resource Constraints

Throughput is always an issue as expectations are always increasing based on national and international competitors capabilities.

4.1.8 Ideal Data Architecture

Ideally HPC resources would be available to support this work:

- 2-4 TB of SSD cache space for data processing
- 1-3 GPU
- 20-40 CPUs
- 512 to 1 TB of RAM
- 10-40 GB network connectivity for data transfer
- At least 10-20 TB per project data storage space that is accessible for 3-6 months duration.

4.1.9 Outstanding Issues

Currently, the available data storage space is limiting on the Agave cluster for cryoTEM data. Further transfer from scratch server to nodes is very slow, thus limiting processing during allowed wall-times. This has frustrated 90% of cryoTEM users into focusing on individual workstations to ensure throughput can be maintained and projects are competitive.

4.2 School of Life Sciences: Cellular Mechanisms of Evolution

Content in this section authored by John McCutcheon, Professor in the School of Life Sciences and associate director of the Biodesign Center for Mechanisms of Evolution

4.2.1 Use Case Summary

The McCutcheon lab is primarily focused on determining the cellular mechanisms of evolution. Our specific system is the symbiotic relationship between bacteria and the mealybug. The genome of the symbiotic bacteria in the mealybug has been reduced much like those found in the presumed symbiotic mitochondria and chloroplasts. We hope to use cellular tomography at increasing resolutions to begin to understand the symbiotic exchange between the gut bacteria and the mealybug in terms of biomolecule and energy exchange. We work closely with the bioEM group through the EMC to generate volume images of the special organelle hosting the bacteria in the mealybug gut.

4.2.2 Collaboration Space

The following people are core collaborators:

- Elizabeth Villa (UCSD)
- Ke Hu (ASU)
- Michael Lynch (ASU)
- Dewight Williams (ASU)
- Po-Lin Chiu (ASU)

4.2.3 Instruments & Facilities

This research will leverage ASU resources Helios FIB/SEM, Krios, and eTEM Titan. We anticipate performing to some extent imaging at UCSD in collaboration with Elizabeth Villa. Large volume tomograms will need to be either reconstructed locally from data collected there or large volume tomograms transferred locally for volume analysis.

4.2.4 Data Narrative

Typically tomograms covering the gut organelle will require 100 GB file sizes, but there will be hundreds to thousands of these tomograms. Individual regions will require averaging that means 1000's to 10,000 of volumes of specific regions within tomograms will need to be averaged. High memory nodes and large data storage will be required to perform this work.

4.2.4.1 Data Volume & Frequency Analysis

The use case is still being developed, and has not produced data yet. It will eventually produce Terabytes over the course of a year.

4.2.4.2 Data Sensitivity

There are no sensitive aspects to the use case's data to report.

4.2.4.3 Future Data Volume & Frequency Analysis

It is expected that the future versions of this use case will produce Petabyte (PB) amounts of data volume, on a monthly frequency.

4.2.5 Technology Support

4.2.5.1 Software Infrastructure

The following software packages are used during this research:

- IMOD¹²: set of image processing, modeling and display programs used for tomographic reconstruction and for 3D reconstruction of EM serial sections and optical sections
- Amira¹³: 2D–5D solution for visualizing, analyzing and understanding life science and biomedical research data from many image modalities, including Optical and Electron Microscopy, CT, MRI and other imaging techniques
- motioncor2¹⁴: multi-GPU program that corrects beam-induced sample motion recorded on dose fractionated movie stacks
- Relion (REGularised Likelihood Optimisation)¹⁵: software package that employs an empirical Bayesian approach for electron cryo-EM structure determination

4.2.5.2 Network Infrastructure

There are no local networking requirements or configurations to report. This research uses existing ASU components that are documented in Section 4.6.

4.2.5.3 Computation and Storage Infrastructure

There are no local computing or storage requirements or configurations to report. This research uses existing ASU components that are documented in Section 4.6.

4.2.5.4 Data Transfer Capabilities

The research has just started, and cannot report on experiences with data transfer.

4.2.6 Internal & External Funding Sources

This research has no external funding from:

- National Science Foundation (NSF)
- Howard Hughes Medical Institute (HHMI)

4.2.7 Resource Constraints

The need for large data storage and processing of large dataset on an HPC.

4.2.8 Ideal Data Architecture

This research would benefit from nodes that contain larger main memory resources (e.g. 512G or beyond), GPU availability to support back-projection of large data volumes, and large data storage space approaching TB to PB.

¹² <https://bio3d.colorado.edu/imod/>

¹³ <https://www.thermofisher.com/us/en/home/electron-microscopy/products/software-em-3d-vis/amira-software.html>

¹⁴ <https://emcore.ucsf.edu/ucsf-software>

¹⁵ <https://relion.readthedocs.io/en/release-4.0/>

4.2.9 Outstanding Issues

Researchers in this lab are requesting training in how to better utilize HPC resources on campus.

Being able to support remote visualizations (e.g., X-Windows over a network that is linking a remote HPC resource) is a critical use case to support.

4.3 Department of Physics: Optical Biophysics and Spatial Transcriptomics

Content in this section authored by Douglas Shepherd, Assistant Professor from the Department of Physics

4.3.1 Use Case Summary

The goal of this research project is to generate nanoscale, 3D maps of RNA expression in the human olfactory system. One way to think about this experiment is imaging-based, 3D, targeted RNA sequencing. By building a healthy atlas of RNA expression across multiple human samples, we aim to understand, for the first time, how nerves are wired from the nose (olfactory epithelium) to the brain (olfactory bulb). This effort is led by Shepherd group at Arizona State University. The collaborating labs are the Presse group at Arizona State University and the Restrepo group at University of Colorado Anschutz Medical Campus.

The research is funded by the NIH BRAIN initiative and is expected to deposit both raw and processed data in centralized databases. Data is generated on a microscope control computer at a rate of 10-50 TB/week, transferred over a private fiber network in the Shepherd lab to a NAS system, and processed using 1 of 2 private Linux servers. The results are then shared with collaborators in summary tables and highly down-sampled forms. The data does not leave one single laboratory at ASU until it is ready for collaborators to view, because the network transfers rates out of the laboratory room often average 10 megabytes/second.

We have a second spatial transcriptomics project in the human lung, funded by the Chan Zuckerberg Initiative. The collaborators on this project are at Northwestern University, Duke University, and the Chan Zuckerberg Biohub. This uses the same infrastructure as above and generates the same amount of data in parallel. This data also does not leave the lab until summary results are ready for collaborators.

4.3.2 Collaboration Space

The following are collaborators in spatial transcriptomics:

- Diego Restrepo, University of Colorado Anschutz Medical Campus
- Alexander Misharin, Northwestern University
- Purushothama Rao Tata, Duke University
- Elizabeth Duong, UCSD
- Nicholas Banovich, TGen
- Roy Wollman, UCLA

The following are collaborators in optics:

- Reto Fiolka, UT Southwestern
- Andrew York, Calico Labs

4.3.3 Instruments & Facilities

The research is performed in the laboratory space located in Building ISTB5, room 171:

- **Acquisition computer:** custom AMD threadripper platform with NVMe RAID card (16 TB) and Nvidia Titan RTX GPU

- **NAS:** Synology with expansion bays, currently with 1PB of raw storage and configured at RAID 10.
- **Server:** custom dual Intel Xeon (12 cores each), 1 TB ram, 32 TB of NVMe RAID, and 2x Nvidia Titan RTX GPU
- **Network:** custom local 10Gbps fiber network with direct connections between:
 - acquisition <-> NAS
 - NAS<-> server

Acquisition computer controls a one-of-kind high numerical aperture oblique plane microscope with fluidics unit¹⁶.

A new server is on order, and will feature:

- AMD Epyc (64 cores)
- 1 TB ram
- 32 TB of NVMe RAID
- 1x Nvidia

A new NAS is also on order, and will double storage capability in the lab.

4.3.4 Data Narrative

We get tissue from collaborators, prepare the libraries and tissue for imaging, run the imaging experiment, process raw data, and store imaging data and results in compressed Zarr files. We write all of our own experimental control, simulation, and data analysis code.

4.3.4.1 Data Volume & Frequency Analysis

The use case produces a Terabytes (TB) amount of data volume, on a weekly frequency.

4.3.4.2 Data Sensitivity

There are sensitive aspects to the use case's data to report, related to the provenance of the samples that are used.

All verified raw data must eventually be deposited in either CZI or NIH funded repositories with sufficient metadata for re-analysis by other groups. Our analysis will be deposited and published.

4.3.4.3 Future Data Volume & Frequency Analysis

It is expected that the future versions of this use case will produce Petabyte (PB) amounts of data volume, on a weekly frequency.

¹⁶ <https://elifesciences.org/articles/57681>

4.3.5 Technology Support

4.3.5.1 Software Infrastructure

We develop all tools for our experiment in-house. We rely on scRNAseq data and analysis tools from other groups, including:

- Kallisto¹⁷: program for quantifying abundances of transcripts from bulk and single-cell RNA-Seq data, or more generally of target sequences using high-throughput sequencing reads
- Scanpy¹⁸: scalable toolkit for analyzing single-cell gene expression data built jointly with anndata
- Scvelo¹⁹: scalable toolkit for RNA velocity analysis in single cells

4.3.5.2 Network Infrastructure

The networking infrastructure, and connections to campus computing, are in need of upgrade. The building infrastructure is known to only support 100Mbps speeds. In particular the bottlenecks described in Section 4.1 and 4.2 also apply to this use case: there are several performance abnormalities when transferring to the local HPC resources as well as sharing with collaborators. It has been observed that on and off campus transfers can perform at around 10Mbps.

4.3.5.3 Computation and Storage Infrastructure

We maintain our own Linux servers with high memory and GPU capability. We can launch Apache instances as necessary. See Section 4.3.3.

4.3.5.4 Data Transfer Capabilities

The use case cannot transfer data, and must rely on removable media.

4.3.6 Internal & External Funding Sources

This research has external funding from:

- NIH MH128867-01 expires 2024
- CZI Human Lung Cell Atlas v1.0 expires 2023

4.3.7 Resource Constraints

The inability to transfer data locally, as well as to remote locations, is hampering research progress. It is impossible to share GB of data with collaborators, and mechanisms that use removable media are regularly used.

On campus-storage is limited, and will reduce the ability to be productive.

4.3.8 Ideal Data Architecture

The research would benefit from an acquisition computer with high speed NVMe raid for storage. It would also be able to transfer data on-the-fly via fiber to central storage.

¹⁷ <https://pachterlab.github.io/kallisto/>

¹⁸ <https://scanpy.readthedocs.io/en/stable/>

¹⁹ <https://scvelo.readthedocs.io/en/stable/>

Being able to run pre-processing and analysis code on remote cluster would be vastly preferred to running things locally. This can be done via an x-windows session, so that we can interactively explore imaging data to set and verify code parameters on random region of interests before batch processing all data.

Lastly, data hosting such that we can share raw data with collaborators without using removeable media

4.3.9 Outstanding Issues

This use case has well documented data transfer issues with RTS, and can benefit from additional support to build a better research network to facilitate local and remote data transfers.

4.4 School of Molecular Sciences: Structural Biophysics and Cryo-EM

Content in this section authored by Po-Lin Chiu, Assistant Professor, School of Molecular Sciences, and member of the Biodesign Center for Applied Structural Discovery.

4.4.1 Use Case Summary

I am in the School of Molecular Sciences at Arizona State University. My research focuses on structural biophysics and cryo-EM to understand the molecular mechanism of biomolecules.

4.4.2 Collaboration Space

Collaborators are mostly in universities, medical schools, or pharmaceuticals. There will be significant work with the Mayo Clinic in 2022 and beyond.

4.4.3 Instruments & Facilities

The research uses the ASU electron microscopy facility and research computing facility.

4.4.4 Data Narrative

We collect mass of electron microscopic image data and perform data processing to reconstruct the 3D molecular structure from the 2D images. The data size is about several hundred thousand to several millions of images for one high-resolution 3D density map.

4.4.4.1 Data Volume & Frequency Analysis

The use case produces Terabyte (TB) amount of data volume, on a monthly frequency.

4.4.4.2 Data Sensitivity

The extent of the sensitivity varies from sample to sample, but many samples are sensitive.

4.4.4.3 Future Data Volume & Frequency Analysis

It is expected that the future versions of this use case will produce Terabyte (TB) amounts of data volume, on a monthly frequency.

4.4.5 Technology Support

4.4.5.1 Software Infrastructure

We mostly use single-particle cryo-EM package to calculate the 3D reconstruction, such as Relion, cisTEM, cryoSPARC, and so on. For visualizing the 3D reconstruction, we often use UCSF Chimera or UCSF ChimeraX.

4.4.5.2 Network Infrastructure

None to report at this time, research uses ASU institutional resources.

4.4.5.3 Computation and Storage Infrastructure

None to report at this time, research uses ASU institutional resources.

4.4.5.4 Data Transfer Capabilities

It takes very long time transferring the data (4 TB for about one week with Globus).

4.4.6 Internal & External Funding Sources

This research has external funding from:

- Army Research Office
- Department of Energy

4.4.7 Resource Constraints

None to report at this time.

4.4.8 Ideal Data Architecture

Nothing to report at this time.

4.4.9 Outstanding Issues

None to report at this time.

4.5 School of Social and Behavioral Sciences: Combating Mis/Disinformation via Critical Data Studies

Content in this section authored by Shawn Walker, Assistant Professor of Critical Data Studies Arizona State University School of Social & Behavioral Sciences

4.5.1 Use Case Summary

These system primarily support research to better understand and combat mis/disinformation on a range of topics. Scholars participating in the project are in the fields of information science, journalism, math, computer science, Barrett Honors, and communication.

On a high level, these systems support multiple research projects focusing on studying mis/disinformation as well as work surrounding social media methods and data archiving. The goals are to support a variety of activities focused on understanding the spread of misinformation in multiple contexts (COVID-19, COVID-19 dashboards, social media platforms such as Facebook and Twitter, alt-tech platforms such as Parler and GAB, social movements, elections, etc.). We also develop methods to collect and analyze data from these platforms at scale and in real-time (collection and some analysis). Methods to characterize the change in this content over time and how to archive this content for continued and future use (data archiving). This includes content from social media platforms, embedded content, and high-fidelity web archiving.

The generalized data lifecycle for most projects includes the collection of data from social media platforms or URLs in real-time (some collections use found data) -> archiving of original data from APIs/platforms/web archives -> preprocessing of data to create derivatives for analysis -> analysis of data or insertion of data into analytical tools such as ElasticSearch/Kibana -> refined analysis for reports or publications -> long-term storage.

Collaborators are from the School of Social and Behavioral Sciences, New America Foundation (external), University College Dublin (external), Cronkite School of Journalism, ASU Library Data & Analytics Lab, Math, computer science (students), and Barrett Honors College (students). Funders include the University of Waterloo/Mellon Foundation and Facebook/SSRC.

A DOD funded project not related to misinformation also uses the infrastructure to collect network survey data and analyze data.

4.5.2 Collaboration Space

Collaborators are located in AZ at multiple ASU campuses (West, Downtown, Tempe), DC (New America), NYC (New America), Stanford (MD), Bristol (Exeter, UK), and Dublin (University College Dublin).

Datasets are stored in project storage at ASU RC and VMs which are accessed via the on-campus network or via the ASU VPN. Smaller datasets and derivatives are stored in

shared Google Drive spaces -- especially when sharing with more temporary external collaborators.

Datasets and processed derivatives related to the Waterloo project come from the Internet Archive.

There are many opportunities for future collaboration and data sharing, but the datasets are too large or contain too many files to share using the services ASU already offers such as Dropbox or Google Drive. For example, a number of potential collaborators have requested access to the Parler data we have inserted into the RC ElasticSearch cluster, but that requires access to the ASU VPN.

4.5.3 Instruments & Facilities

Primarily we use 5 VMs provided by ASU RC. Two of the VMs mount a 50TB project storage which is primarily used by the New America. The other two VMs mount a shared 40TB storage with local 2TB mounts for temporary storage.

These two VMs support all other projects. For some types of data analysis (network analysis, topic modeling, etc.), data is transferred to the HPC cluster for analysis.

Depending on the need, workflows can be transferred between the two machine groups.

For the Waterloo project, some web archive data (WARC) is processed using the Archived Unleashed tools at the Internet Archive. Derivative data is downloaded from IA and further analyzed using our VMs at ASU.

Derivatives will be exported to Google Sheets, Drive, etc. as needed for less technical collaborators to work with the data.

We also use the ElasticSearch cluster and Kabana interface to make some social media data searchable to multiple teams.

4.5.4 Data Narrative

In general, we ingest data from the following sources:

- Twitter COVID-19 stream. This is a real-time stream from Twitter providing all tweets tagged by Twitter as related to COVID-19. This can consist of hundreds of thousands to tens of millions of tweets per data. Six connections to the Twitter API endpoint must be maintained at all times or data is dropped and cannot be collected again. This data is in Twitter JSON (text) format. Data files are rotated and compressed on an hourly basis.
- Twitter compliance stream. This is a real-time stream providing lists of tweets deleted or hidden as well as Twitter accounts deleted, suspended, or hidden. This requires 4 real-time connections to the API endpoint and must be maintained at all times as any drops in that connection will result in a loss of data that cannot be collected again.

- Multiple Twitter Streaming API connections - Real-time connections to the public Twitter streaming API for various projects. Data is received in real-time and any drop in the connection results in lost data.
- COVID dashboard web crawls - This is a set of docker containers which archive websites using web recorder in real-time and at high-fidelity. We use the system to make daily archives of state COVID-19 dashboards. Crawls are kicked off via corn jobs each day to archive state dashboards. The data is outputted in WARC format contains HTML, JavaScript, images, videos, etc. -- every component of a website that the browser loads.
- Archive Twitter and Facebook data - we have a number of Twitter and Facebook datasets which were capture via searches or web crawls. This data is collected one time and used for multiple types of analysis. We also have archives from Facebook data collected via their CrowdTangle service. We also use some data from Facebook's link dataset (FB/SSRC/Social Science One), but this data is usually accessed and analyzed in FB's cluster environment. All of this is text data.
- Associated crawls of URLs in social media data. Depending on the project, we will create web archives of embedded images, videos, and links in the above mentioned social media data.
- Found Parler Data - Parler 1.0 was de-platformed on Jan 10. We have multiple crawls of Parler data that we conducted ourselves using various open source tools. We also have multiple captures of found data from Parler other researchers and hackers have posted on the web. This includes over 200M posts, 25TB of Parler CDN video data, and 236GB of Parler CDN images.
- Internet Archive crawls of Geocities and COVID-19 dashboards. This data was downloaded from the Internet Archive. We have received access to their crawls for the Waterloo project.

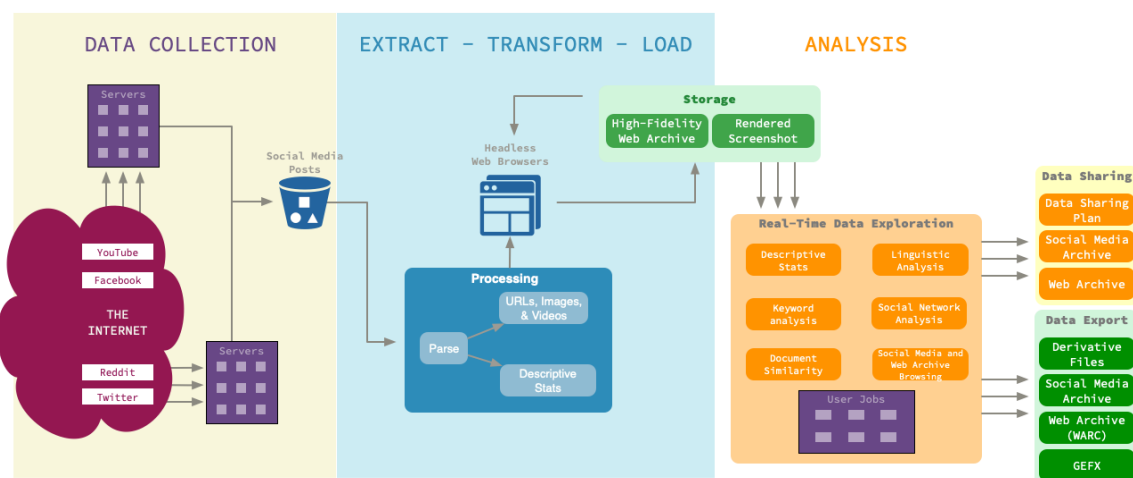


Figure 4.5.1: Workflow Diagram

A general workflow involves:

- Question development

- matching any available data to the research question via API/web crawls/found data
- archiving or collecting of that data
- various methods of data exploration and making the data available to all team members via ES, derivatives, derivative CSV files, visualizations
- often this process involves multiple steps of iteration
- report output and possible export of the data to share with the public or other researchers

Sometimes our workflow is a bit different when it comes to found or released data, we instead capture released/leaked data ASAP, stabilize it, and determine methods to analyze it to understand the possibilities. Collaborators from New America, will often combine this data analysis with OSINT methods to contextualize the data to other actors for final reporting.

4.5.4.1 Data Volume & Frequency Analysis

The use case uses Terabyte (TB) amounts of data volume. Frequency is in real-time for web archives and Twitter data. Found data, closed platforms (Parler), and archived data would be collected at a single point in time. Some social media and web archive data is re-crawled at regular intervals as appropriate for the project (i.e. daily, weekly, or monthly)

4.5.4.2 Data Sensitivity

Yes, there are sensitive aspects to the use case's data and they can be listed or explained, While social media data does not fall under HIPPA or other compliance requirements, the data contains PII and our analysis of the data focuses attention on specific users which could have ethical implications. Some data contains offensive, racist, or threatening content (i.e. hate speech, white supremacy, violence at demonstrations, threats, etc.).

4.5.4.3 Future Data Volume & Frequency Analysis

It is expected that the future versions of this use case will produce Terabyte (TB) amounts of data volume. Some datasets will be archives, but many datasets must still be collected in real-time

4.5.5 Technology Support

4.5.5.1 Software Infrastructure

The vast majority of our everyday data collection is carried out via customer Python scripts, Jupyter notebooks, and a few R scripts. The mentioned "found datasets" are often downloaded from AWS S3 buckets. Small portions of data and derivatives are transferred to Google Drive for sharing with team members or even external collaborators.

Desktop analysis tools include Gephi and ORA for network analysis and viz. We have used DataWrapper for embedding visualizations in final reports.

When using the cluster for network analysis, we use RAPIDS.ai libraries.

For some web archives analysis, we use the ARCH as an analysis tool and to produce derivatives files from large WARC²⁰ (archive) files.

Research Computing provides access to their ElasticSearch and Kibana cluster for some data analysis. This is hosted within the ASU VPN.

4.5.5.2 Network Infrastructure

This mostly includes network connectivity to the wider Internet as well as AWS S3 and Google Drive. Connectivity between VMs, the VPN, and on-campus buildings is used to transfer data for local storage or exploration. The VPN performance is not ideal, and is a serious bottleneck.

4.5.5.3 Computation and Storage Infrastructure

For our project, this includes just basic VM access, storage as well as occasional access to the HPC cluster and scratch storage.

4.5.5.4 Data Transfer Capabilities

We downloaded 36TB from AWS S3. This process took a few days to complete and cost \$3,000 in transfer charges.

We also transfer multiple GB files on a regular basis between ASU and my college at University College Dublin.

We transfer multi-TB data files from the Internet Archive and ASU a few times a year.

4.5.6 Internal & External Funding Sources

This research has external funding from:

- SSRC/Facebook (2021)
- Waterloo/Mellon Foundation (2022)
- Department of Defense (DoD)

4.5.7 Resource Constraints

The speed of connectivity (most likely on both sides) between ASU, the Internet Archive, and sometimes AWS can be an issue when downloading large datasets. This requires more advance planning to ensure that a dataset is completely downloaded before needed.

In general, we face a host of issues when attempting to acquire social media and web data. Web crawlers are ill-equipped to capture dynamic web pages (like COVID-19 dashboards). Social media companies either do not offer way to access data or provide little access to that data. I'd imagine that this is outside of the scope of EPOC -- a central clearinghouse where researchers could share, compare, and archive their social media datasets would be awesome. Something like the HathiTrust for social media data would

²⁰ <https://archivesunleashed.org/arch/>

be amazing. Would be happy to discuss something like this, but I also recognize the logistical and legal issues involved.

4.5.8 Ideal Data Architecture

Ideally, we would like to develop a workflow that take social media data we collect --> runs all URLs and web content into a web archiver --> automatically runs some analysis on the data for presentation --> provides a dashboard with basic exploratory stats about the data and interactive network graphs --> provides suggestions on additional data to collect, preemptively collects that data until a researchers confirms or removes that data form the collection --> checks the status of archived data and the arability of that in the live web on a regular basis.

Besides the very complex workflow development, this would require large amounts of storage (multi-TB per day) and a number of crawlers to collect this data in real-time.

4.5.9 Outstanding Issues

Data protection, security, and providing access are an issue for our group. Often we would like to provide temporary access to data or tools, but the effort required to request an affiliate ASURITE and provide access is too high. Since the majority of our VMs have internal IPs as a security measure, it is difficult to facilitate quick exchanges or access to tools.

4.6 Arizona State University Research Technology Office & Research Computing

Content in this section authored by Marisa Brazil, Douglas Jennewein, Gil Speyer, & Jason Yalim, from the ASU Research Technology Office and Research Computing

4.6.1 Use Case Summary

Research Computing provides advanced computing and data services to researchers across ASU on all four main campuses. Services include training, academic course support, high performance computing, virtual machines, bulk data storage, server co-location, and HIPAA-aligned secure VMs and storage.

4.6.2 Collaboration Space

We work closely with XSEDE through Campus Champions and CaRCC through the People Network. We host a regional instance of XSEDE Jetstream2 and run several OSG compute servers. We have several researchers using XSEDE systems, OSG, and some commercial cloud (GCP). We are a Software Carpentry platinum sponsor and host many workshops each year. We are a Dell Technologies center of excellence in HPC and AI and work closely with Dell on software and platform development including the Omnia provisioning stack.

4.6.3 Capabilities & Special Facilities

User-facing service overview²¹ describes all offerings

We operate three data center spaces:

- small space in the GWC Engineering building for low-risk systems requiring physical access from students or non-RC personnel
- the main Tempe campus data center in the ISTB1 building which houses project data storage and the legacy HPC cluster and VM infrastructure as well as Science DMZ gear
- the new FISMA high data center at Iron Mountain Phoenix where the new HPC cluster and VM infrastructure is currently being deployed.

4.6.4 Technology Narrative

4.6.4.1 Network Infrastructure

LAN

Research Computing local network currently consists of more than 100 switches across three separate data centers. The switches we manage range in speeds from 1G to 400G and vary in brands, which include Dell Networking, Arista, and Penguin Arctica. Last December, we worked with our campus network provider, HyeTech, to redesign the border and upgrade the Science network. This redesign migrated management of the border switches to Hyetech, which allowed us to take better advantage of the campus infrastructure and utilize VXLAN between sites.

²¹ <https://asurc.atlassian.net/wiki/spaces/RC/pages/60915741/Services+and+Pricing+Structure>

At each of Research Computing's main data centers, we manage a pair of Fortigate Firewalls. The first set, at our on-campus data center, has two bonded pairs of 40G interfaces (2x80G aggregated) for inbound and outbound traffic. The second pair is currently being installed and will feature a similar setup but with 100G interfaces (2x200G aggregated) for inbound and outbound traffic. These high-speed firewalls allow us to keep most of our equipment behind them without compromising security or performance. Having the ability to individually pick which systems are monitored with the IPS rules and which ones are only sitting behind ACLs and logged to elastic, makes sure systems like data transfer nodes have traffic logged but are not slowed down and systems like researcher-managed VM's are still protected.

Most of the network is layer 2 with all layer 3 taking place at the firewalls. We currently host about 200 VM's, running across 14 Xen nodes. Many of the VM's are researcher-managed and are on private VLANs with their own private IP space. Outside of offering the VM's we also occasionally set users up with the Fortinet SSL VPN to connect to internal services and manage hosted hardware. Beyond hosting systems, we offer some researchers and departments private ""virtual"" firewalls where they can manage their own rules and have high-speed access to the HPC cluster and internet to do their research.

MAN

The Metro Area Network is another place that has had many new changes over the last few years. In July of 2020, the campus border networks were upgraded to provide ASU with redundant 100G connections between all campuses, across different providers and data centers. The connectivity helped pave the way for the new Research Computing space that is currently being built out at the Iron Mountain data center.

ASU Metro Ring

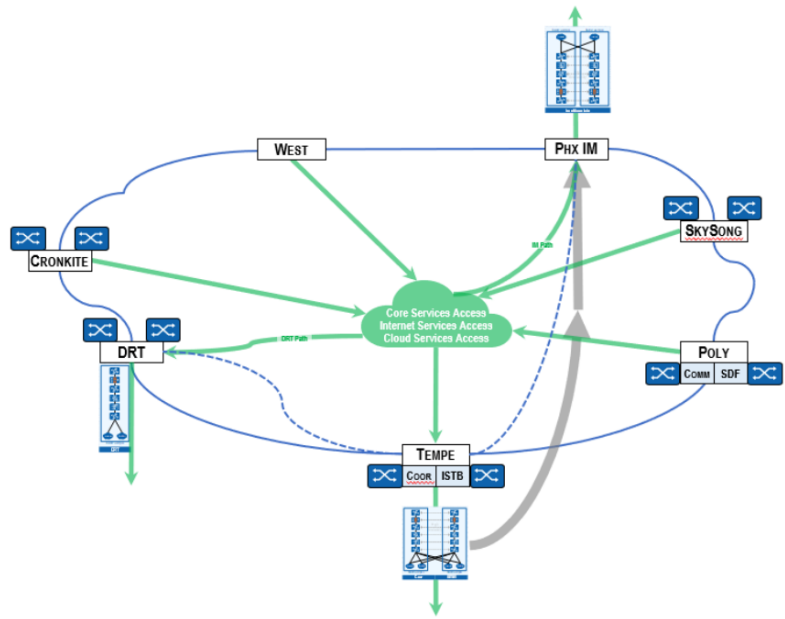


Figure 4.6.1: MAN Diagram

WAN

The WAN at ASU consists of a redundant 100G commodity connection through COX communications and a redundant 100G connection to Internet2 through Sun-corridor. The ASU WAN connections are managed by Hyetech, who provide 100G connectivity to our border along with jumbo-frame and ipv6 support.

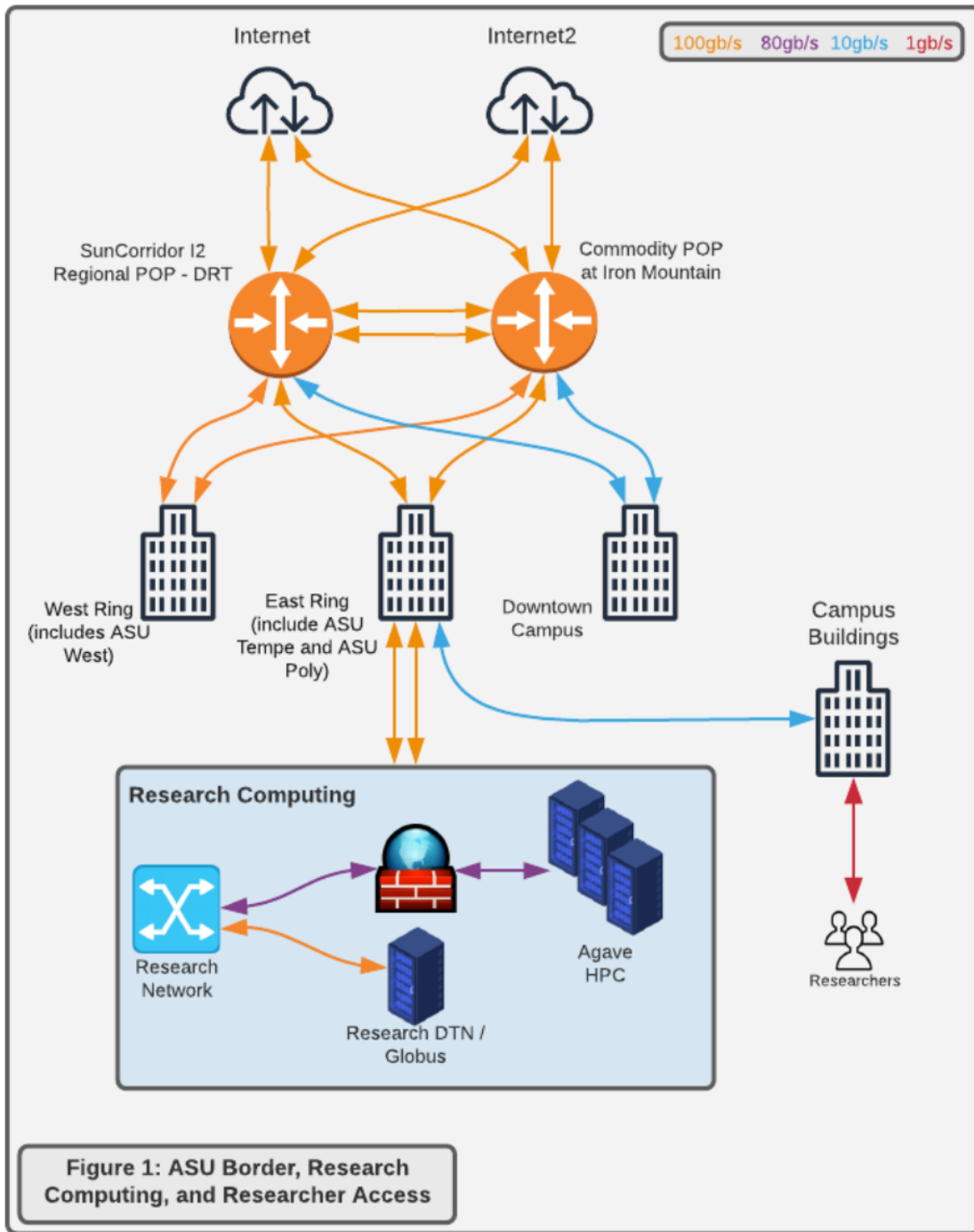


Figure 4.6.1: WAN Diagram

4.6.4.2 Computation and Storage Infrastructure

Research Computing currently manages a large 22K core, mainly intel-based, HPC cluster with more than 360 GPU's. This cluster, Agave, lives on the Tempe campus in the ISTB1 Datacenter. Our new cluster, which is currently being built at a remote data center, will be AMD EPYC based with a mixture of Nvidia A100's and A40's. Both of these clusters sit behind the Research Computing firewalls and have dedicated login nodes,

scratch, and home storage. Our Scratch storage at both sites is Beegfs, our home and project storage is Qumulo and Isilon. All of the systems are connected with 10g, 25g, or 100g interfaces with 56G or 100G IB/Omnipath. Outside of the two main clusters, we do have 20 nodes dedicated to OSG. This will likely be decommissioned when we start to offer OSG overflow on the current HPC clusters.

4.6.4.3 Network & Information Security

Research Computing as stated previously, maintains a set of high speed, high throughput firewalls, which most of the equipment we host sits behind. The logs from this are shipped off to Elasticsearch where we monitor for suspicious activity using elastiflow to then generate an IP block list which is sent back to the firewall. Outside of this, we use elasticsearch with tools like Wazuh to gather information from individual systems which is then shipped back to elasticsearch and logged. We also have alerts set up on the firewalls to catch and report suspicious login activity and compromised hosts.

4.6.4.4 Monitoring Infrastructure

We operate a perfSONAR node for the campus²² that handles wide area and regional network monitoring.

For system monitoring, we are currently using a mixture of Icinga, Nagios, elasticsearch, Grafana, and Prometheus. These tools monitor local systems and then either alert via email, Slack, or victorops. Initially we had decided on Icinga for monitoring with Elasticsearch on logging but Icinga does not seem to be as HPC friendly as we had hoped. We started using Prometheus with Grafana and it seems like a much better tool for monitoring HPC nodes. We will likely continue using Icinga for monitoring VM's and hosted systems and just using Prometheus to monitor and report on the cluster nodes.

For network performance and monitoring, we are using perfSONAR and elastiflow. We have recently started to explore using Consul to measure the performance of apps and services, many of which live on different subnets, this looks like it will be a good indicator of network performance going forward. Currently, all of this data is used internally only, to troubleshoot issues when they arise.

4.6.4.5 Software Infrastructure

Currently the main tools and software packages that we use to support HPC users on the clusters are:

- Open on Demand²³: "one-stop shop" for access to High Performance Computing resources
- Slurm²⁴: job scheduler for Linux and Unix-like kernels, used by many of the world's supercomputers and computer clusters
- BeeGFS²⁵: hardware-independent POSIX parallel file system

²² <http://perfsonar.rc.asu.edu/toolkit/>

²³ <https://openondemand.org>

²⁴ <https://slurm.schedmd.com/documentation.html>

²⁵ <https://www.beegfs.io/c/>

- Elasticsearch, Kibana, Grafana²⁶: a search engine and visualization platform that provides a distributed, multitenant-capable full-text search engine with an HTTP web interface and schema-free JSON documents
- Nagios²⁷, Icinga²⁸ and Prometheus²⁹: monitoring software
- OpenLDAP³⁰: Lightweight Directory Access Protocol
- Bind³¹: DNS management
- Ansible³² and Salt³³: automation platforms
- Lmod³⁴ and Spack³⁵: package management

This list will change as we build out the software stack for the new cluster. The new cluster will phase out some of the older tools in lieu of newer, better supported tools.

4.6.5 Organizational Structures & Engagement Strategies

4.6.5.1 Organizational Structure

ASU Knowledge Enterprise is a matrix organization with a rough hierarchy.

The Research Technology Office is led by Sean Dudley, AVP and Chief Research Information Officer. RTO units include:

- Research Computing - Doug Jennewein
- Research Engagement - Marisa Brazil
- Research Accelerator - Gil SPeyer
- Research Data Management - Philip Tarrant
- Research Enterprise Architecture - Chris Kurtz
- Research Information Security - Michael Hacker
- Research Technology Support - Tim Remmington, deskside support and ticketing

4.6.5.2 Engagement Strategies

Marisa Brazil leads the Research Engagement team in the Research Technology Office at ASU. This unit conducts first contact meetings with new faculty members as well as researchers new to research computing. It also conducts over 50 regular training workshops each year, reaching hundreds of faculty, staff, and students. We also host an annual weeklong research computing expo combining training events with research presentations. RC also participates in an annual Grants Research and Sponsored Projects conference sponsored by the VPR for attracting new researchers.

²⁶ <https://www.elastic.co>

²⁷ <https://www.nagios.org>

²⁸ <https://icinga.com>

²⁹ <https://prometheus.io>

³⁰ <https://www.openldap.org>

³¹ <https://www.isc.org/bind/>

³² <https://www.ansible.com>

³³ <https://saltproject.io>

³⁴ <https://lmod.readthedocs.io/en/latest/>

³⁵ <https://spack.io>

4.6.6 Internal & External Funding Sources

The majority of our funding comes from central administration, with additional annual support from Engineering and the College of Arts and Sciences. We also generate revenue by charging for some of our services.

Current grants include:

- CC* Compute: The Arizona Federated Open Research Computing Enclave (AFORCE), an Advanced Computing Platform for Science, Engineering, and Health³⁶
- CC* Networking Infrastructure: Science DMZ for Data-enabled Science, Engineering, and Health³⁷
- Frameworks: Collaborative Research: Integrative Cyberinfrastructure for Next-Generation Modeling Science³⁸
- We have a pending NIH SIG HEI proposal for a GPU cluster for the Cryo-EM core facility.

4.6.7 Resource Constraints

We are recovering from budget reductions put in place at the beginning of the COVID-19 Pandemic. ASU has weathered the pandemic well, financially speaking, and budgets are returning to their pre-pandemic levels. In general RC enjoys strong executive support and funding.

4.6.8 Outstanding Issues

Our lead HPC engineer departed in August and the position remains unfilled. We are in the middle of deploying a new cluster while continuing to operate a heritage HPC system.

³⁶ https://www.nsf.gov/awardsearch/showAward?AWD_ID=2126303&HistoricalAwards=false

³⁷ https://www.nsf.gov/awardsearch/showAward?AWD_ID=2018886&HistoricalAwards=false

³⁸ https://www.nsf.gov/awardsearch/showAward?AWD_ID=2103905&HistoricalAwards=false

Appendix A - Research Computing Facilities Statement

Personnel

Arizona State University (ASU) is served by both the central University Technology Office (UTO) and the Research Technology Office (RTO). UTO is the central IT organization with over 540 FTEs across multiple service areas including desktop support, wired and wireless networking, public and private cloud, identity management, information security, and web application development. UTO oversees core campus IT services such as payroll, email, instant messaging, user file storage, and document creation/collaboration. UTO also handles ASU policies regarding IT services, data governance, and information security. RTO focuses on IT services directly supporting research and researchers. Specifically, RTO comprises 65 FTEs covering Research Computing, Scientific Software Engineering, Research Data Management, Business Intelligence, and Web Services. RTO is overseen by the Chief Research Information Officer who reports to the University's Executive Vice President of Research. The Research Computing staff consists of computational scientists, programmers, engineers, and database administrators with expertise in all areas of computing, including scientific and parallel computing, big data analytics (in memory), custom software development, database engineering, and scientific visualization.

Advanced Computing

Research Computing is an academic supercomputing facility providing high performance computing (HPC) environments, a data intensive ecosystem, connectivity to the Internet2 research and education network, and large-scale data storage with elastic capacity to the public cloud. Research Computing provides a variety of HPC (both physical and virtual), cloud, storage, development, implementation, and consulting services. Research Computing consulting services and support for computational investigations, including data analysis, simulation, modeling, visualization, and other high-performance approaches include:

- Identifying optimal systems and software platforms
- Training in computational and/or graphics algorithms, tools, and packages
- Developing custom post-processing graphics tools
- Creating virtual environments for scientific research and fine arts
- Tuning applications for peak performance and implementing parallel algorithms and programs
- Purchase consultation for server, HPC, and storage acquisitions
- Virtual server provisioning (local, cloud)
- Physical and virtual server management
- Providing state-of-the-art interfaces to HPC systems
- Recharge of non-preemptable computing time and data storage solutions
- Accessing extensive external, government-funded compute resources (XSEDE, OSG)

Advanced Computing Systems

The Agave supercomputer is ASU's flagship high performance computing (HPC) cluster. Agave is a heterogeneous Intel-based HPC cluster containing over 15,000 CPU cores. Each node is equipped with a solid-state drive and system memory ranging (depending on the node) from 128GB to 256GB of DDR4 2400 RAM. For large memory applications, the cluster also contains three nodes with 1TB, 1.5TB, and 2TB of DDR4 2666 RAM. GPU computing capabilities include access to over 330 NVIDIA A100, V100, K80, GTX 1080, and RTX 2080 GPU accelerators. A dedicated pool of 1.2PB high performance BeeGFS provides fast scratch storage for compute jobs. Compute nodes are accessible through four login nodes, one of which hosts the NSF-funded Open OnDemand web interface. Compute jobs are managed with the SchedMD Slurm scheduler. Agave is supported by a 100Gbps InfiniBand network fabric. It is connected to the campus Science DMZ, Internet2, and Data Transfer Nodes by a 100/40GE core network.

A dedicated pool of 1.2PB high performance BeeGFS fast scratch storage is presented to the Agave cluster via dual interconnected networks (InfiniBand and Omni-path), and a 1.8PB network attached storage array provides HPC home directory storage. For general purpose research data, a 4PB network attached storage array provides project storage.

Researchers may also purchase their own compute nodes and incorporate them into the Agave cluster, with Research Computing supplying all necessary rack space, power, cooling, networking, and software maintenance. Compute nodes are 52-CPU-core Intel servers manufactured and supported by Dell Technologies. Researchers and their delegates have priority access to their purchased nodes, and any idle capacity is made available to computing jobs for the general ASU research computing community. Such jobs will be guaranteed to run without preemption for at least four hours, after which jobs submitted by the node's owner will preempt them. The owner may also reserve their nodes exclusively for up to three one-week periods per year. Purchasing computing capacity in this manner allows researchers guaranteed access to the necessary computing power without needing to operate and maintain their own servers.

Research Computing will support researcher-purchased nodes for as long as feasible. However, beyond the hardware warranty period, the faculty is responsible for any hardware and labor costs necessary to maintain the hardware. Once the warranty period has expired, Research Computing may remove the node from the cluster if it is no longer technically feasible to support it.

Dell Center of Excellence for HPC and Artificial Intelligence

ASU has been designated a Dell Center of Excellence for HPC and Artificial Intelligence, the third such center in the United States, and the sixth globally. This distinction has grown out of a close collaboration with Dell HPC experts on system architecture, design, and innovation.

Through this new partnership ASU is deploying a new high performance computing system in the second half of 2021. Anticipated characteristics are a mixed CPU/GPU

environment of approximately 15,000 CPU cores, 200 GPU devices, and 6PB of high performance storage.

Open Science Grid

Research Computing runs a 20-node Open Science Grid (OSG) site for the research community at large and is investigating using spare cycles on idle lab workstations to significantly augment this OSG infrastructure.

Data Center

ASU Research Computing maintains an on-campus data center in Interdisciplinary Science and Technology Building 1 (ISTB1), in the center of the ASU Tempe campus. The ISTB1 facility was built in 2012 consists of a 5,000 square foot primary data center for critical systems, networking, and computational resources as well as a 3,000 square foot secondary data center for non-critical systems, development, and individual research development equipment. Both data centers employ a standard “hot and cool aisle” layout with computer room air conditioning units totaling 200 tons, supported by campus chilled water systems. Room-dedicated FM-200 fire suppression systems protect the facility. Power for the facility is from on-campus natural gas turbines fed by Utility natural gas. Data center power supports dual power feeds, protected by two uninterruptible power supply units totaling 1MW, with an onsite diesel generator providing 2MW of power and a 500-ton emergency chiller in the event of a loss of utility power. Access to the data center requires a keycard and PIN, and the facility is monitored 24x7 from a dedicated operations center, and physical access is controlled and maintained by the UTO operations center.

ASU Research Computing is currently building out a new data center at the Iron Mountain Phoenix facility. This facility will provide ASU with more than four times the capacity of the existing infrastructure in a commercial Tier III+ data center with advanced power, cooling, and network capabilities. ASU Research Computing will be a core tenant of the facility and will have the capabilities to provide secure research (up to FISMA³⁹ High) security. A private fiber ring will connect the facility to the ASU Tempe Campus. Internet2 and Commodity Internet circuits will be available at the facility as secondary connectivity, as well as private point-to-point circuits as Research requires.

³⁹ <https://www.dhs.gov/cisa/federal-information-security-modernization-act>

Network

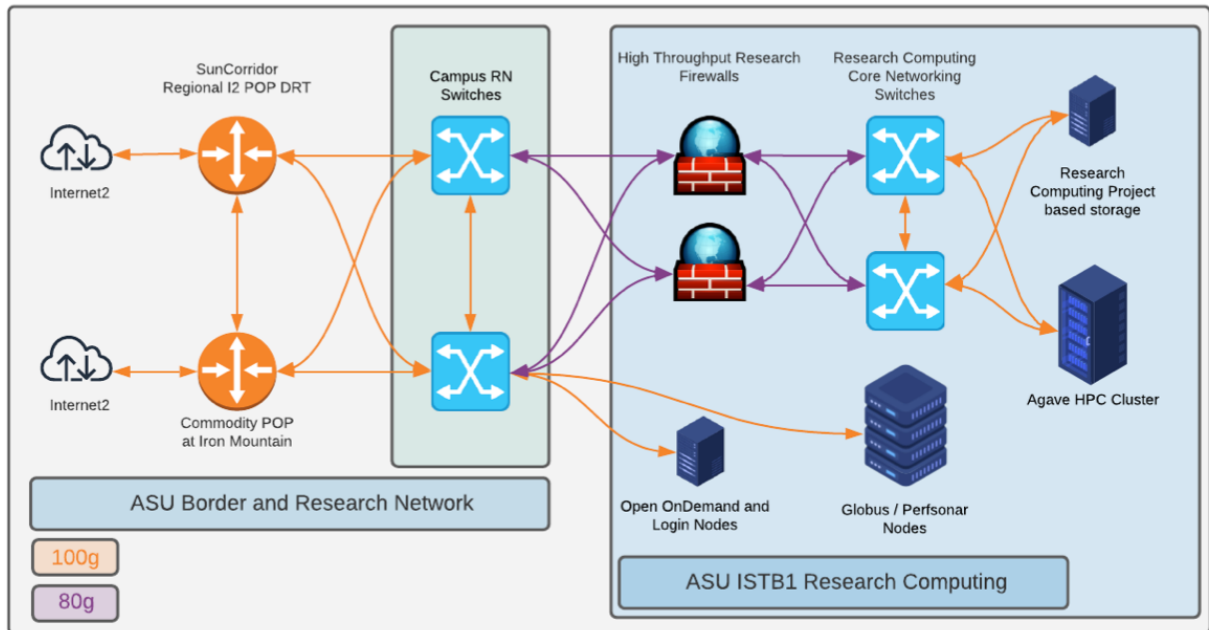


Figure 1: Research Computing Network Connectivity

Figure B-1 illustrates the current ASU network border topology. Primary network access is via a 100 Gb Internet2 circuit to the Tempe campus. Secondary 10 Gb commodity Internet circuits provide additional and partially redundant network access directly from the Tempe campus. The ASU Polytechnic campus (25 miles southeast of Tempe), the ASU West campus (25 miles northwest of Tempe), and the ASU Downtown campus (10 miles west of Tempe) connect to the Tempe campus via redundant 10Gb circuits on a commodity fiber ring. The ASU network is monitored 24x7x365 by a commercial network provider as well as by the University Technology Office.

Buildings on all ASU campuses are connected in a hub-and-spoke model, with most buildings served by redundant 10Gb links and 1Gb to end users. The campus network employs an advanced security complex consisting of a layered defense-in-depth deployment of security controls that include DDoS and IP reputation, a variety of specialized network firewalls, and anti-phishing protections. The ASU cybersecurity program also includes mandatory security education and awareness training, and the UTO Governance, Risk, and Compliance Team conducts continuous assessments evaluating risk and vulnerabilities.

Science DMZ

The ASU Science DMZ is a network enclave that bypasses the network security complex. The Science DMZ is explicitly designed for high-throughput data movement, incorporating 100/40 Gigabit Ethernet, virtual circuits, and software-defined networking capabilities as well as dedicated systems for large data movement requiring a friction-free path, with security policies and enforcement mechanisms tailored for high performance science environments.

Data Storage

The ASU University Technology Office (UTO) supports cloud storage using a variety of cloud-based storage offerings, including Enterprise Dropbox (for Staff/Faculty, 1TB limit), Microsoft OneDrive (available to all Staff, Faculty, and Students via Office 365, 1TB limit), and Google Drive (available to all Staff, Faculty, and Students via G Suite for Education, no storage limit). UTO provides storage to Business Units via SMB/CIFS on Enterprise NetApp Network Appliances.

ASU Research Computing provides 100GB of home directory storage for users of the Agave Cluster, as well as access to the 1.3PB BeeGFS high-speed short-duration scratch environment for cluster computing jobs. Research Computing also operates 4PB of network-attached project term storage. This storage is accessible to the Agave HPC cluster and individual researcher workstations via traditional network shares. The Globus data movement platform provides resilient high-speed access to data stored on Research Computing systems and allows for transfer to user's University provided Google Drive accounts.