May 2021

# Current State of Advanced Research Computing in Canada

## An Update to the 2017 LCDRI Advanced Research Computing Position Paper

Authored by the Alliance's Advanced Research Computing Working Group:

Seppo Sahrakorpi, Maxime Boissonneault, Emmanuel Chateau-Dutier, Chris Loken, Catherine Lovekin, Carolyn McGregor, Felipe Pérez-Jvostov, Ghilaine Roquet, and Lisa Strug

**Digital Research Alliance** of Canada | **Alliance de recherche numérique** du Canada

Funded by the Government of Canada | Canada

# Table of Contents

# Table of Figures

# Table of Tables

# 1 Executive Summary

This report serves as an update to the 2017 Advanced Research Computing Position Paper submitted to Innovation, Science and Economic Development Canada (ISED) by the Leadership Council for Digital Research Infrastructure (LCDRI). The report summarizes the Advanced Research Computing (ARC) landscape in Canada, and documents strengths, challenges and opportunities for the current ARC ecosystem and the Digital Research Alliance of Canada (the Alliance). The report is written by the Alliance ARC Working Group and reflects the team's view and emphasis. This report is not the Alliance's formal statement of the current state of ARC in Canada.

The intent of this work is to help the Alliance to build on the current state and chart a path forward that advances ARC in coordination with other digital research infrastructure (DRI) elements to support research excellence in Canada. Findings and observations in this document, alongside the Research Data Management (RDM) and Research Software (RS) Current State Assessment publications, are meant to provide counsel and background information to the Alliance analysts and management, the Alliance Board, and the Alliance Researcher Council in order to support the needs assessment process, and the Alliance's New Service Delivery Model, Strategic Plan and Funding development processes.

ARC is critical to modern research, and the demand is growing phenomenally (e.g., in big data, and artificial intelligence (AI) etc.). It is highly competitive, and the technology and tools are continually and rapidly changing, with new disciplines rising and transforming as a result. All this demands an agile and highly-responsive ecosystem with expert highly-qualified personnel (HQP), and sustained and predictable funding if Canada is to meet 21$^{st}$ century challenges and remain competitive internationally.

Canada is a very diverse, huge country with a long history, and the Canadian research community reflects that diversity, both as a society and via geography. Canadian ARC and DRI ecosystems need to serve all members of this diverse community in order to better advance Canadian research, and benefit society and progress in Canada. ARC can help not only solve problems and scientific questions that would otherwise be difficult to solve, but also contribute to solutions that would not be possible via regular (e.g., analytic, experimental, or workstation level computing) means.

ARC is a suite of digital technologies and services that enables researchers to solve research issues that are either too large or too complex for them to undertake on their own. In addition to traditional server clusters, modern ARC also includes virtual infrastructures and cloud computing. ARC equips researchers with an advanced digital environment in much the same way that higher education physically equips itself with traditional infrastructure and facilities to host research activities, in an environment where the resources are very highly utilized and pooled. The key components of ARC are

- Computation (e.g. central processing units (CPUs) and graphical processing units (GPUs),

- Active storage and backup (e.g., runtime, nearline and temporary storage),

- Support, training and consultation from highly-qualified personnel (HQP),

- Software stewardship and support (system software and commonly used libraries and communication protocols),

- Privacy, security & authentication,

- High-speed connections to national and international R&E networks as well as between sites

- and Support of and coordination with other DRI components (networking, data management, mid- and long-term storage, research software).

Notably within Compute Canada federation (CCF), due to funding mandates, these components do not include mid-term (repository), and long-term (archival) storage hardware, systems and services, which nevertheless are an integral part of the modern research computing cycle and DRI as a whole. These storage components and considerations are a key part of the Alliance's future mandate.

## Current state summary

Researchers have varying degrees of access to ARC resources at the group, departmental, institutional, regional, national, and international levels. It should be noted that many researchers are unaware of the existence of the CCF resources, think these resources are not for them, or make use of other non-CCF systems. Unfortunately, the only level we can currently quantify and understand well is the national CCF level, although the non-CCF resources available for Canadian researchers are substantial. This current report uses CCF usage data almost exclusively.

The number of **CCF registered users** has grown significantly in the last decade, to 16,000 in 2020. Since 2014 the compounded annual growth rate (CAGR) has been roughly 12%. As of January 1st 2020, the largest user group was 'faculty' at 27%. The four largest user groups (faculty, doctoral and master's students, and post-docs) added up to three quarters of all CCF users.

The largest CCF faculty user groups are from engineering, and biological and life sciences, at 19% each. The humanities, social sciences, business, and psychology faculty account for ca. 10% of the faculty user base  while these disciplines have roughly 46% of all the full-time academic faculty in Canada, clearly indicating how these disciplines are underrepresented in the Canadian ARC in terms of number of users (keeping in mind that it is difficult to ascertain how much is due to differences in ARC needs in different disciplines) and in the general trend of academic research going digital. The number of users from these disciplines has grown roughly 40% over three years since 2017. The absolute number of humanities etc. users still has substantial room and potential to grow, presenting an opportunity for targeted DRI training and support campaigns. There is a real need for access to advanced computing resources in the humanities and social sciences. A strong transformation of these disciplines is underway which requires specific resources that are not currently well served in Canada and which also to a degree explains the under-representation of these disciplines.

**Recent substantial Canada Foundation for Innovation (CFI) investments in Canadian cyberinfrastructure** have resulted in consolidation of resources, new systems being built and a major increase in ARC capacity and capability in Canada. The five main CCF sites host five new

national systems and are affiliated with regional CCF member organizations as follows, from West to East:

- University of Victoria, Arbutus (Compute Canada/WestGrid),

- Simon Fraser University, Cedar (WestGrid),

- University of Waterloo, Graham (Compute Ontario),

- University of Toronto, Niagara (Compute Ontario),

- and McGill University/ École de technologie supérieure, Béluga (Calcul Québec).

Cedar, Graham, and Béluga are general purpose heterogeneous ARC clusters for a variety of ARC workloads. Arbutus is a general purpose ARC cloud system for hosting (mostly Linux based) virtual machines and other cloud workloads. Niagara is a massively-parallel homogeneous ARC cluster for scalable ARC workloads.

The **CPU resource usage** within CCF was roughly 200 000 CPU-years in 2019, indicating a four-fold increase in CPU resource usage since 2010, corresponding to roughly 16% compounded annual growth rate (CAGR). Keeping in mind that the usage has always been limited by available resources. i.e. historically, whatever resources CCF has had have been used by the community. Three research fields (engineering, physics & astronomy, and chemistry & biochemistry) consumed roughly three quarters of the CPU resources. Social sciences, psychology, business, and humanities increased their CPU resource usage eight-fold over the last decade, although in absolute terms the usage was modest at roughly 1250 CPU-years in 2019. The growth in the ARC CPU resource usage in these disciplines indicates the strong interest and potential these disciplines have for leveraging DRI in future. Keeping in mind that in these disciplines the CPU-year indicator is not necessarily the most relevant one for evaluating the use of resources.

Considering **GPU usage** across the CCF ecosystem, in 2019 the total usage of GPU resources was roughly 1300 GPU years, corresponding to roughly 56% CAGR since 2012. This growth was severely restricted by the available supply and is thus not indicative of actual growth rate of general-purpose GPU (GPGPU) computing need in general. Four disciplines (computer and information science, chemistry and biochemistry, biological and life sciences, and physics & astronomy) used roughly 87% of the GPU resources in 2019.

**Cloud usage** across the CCF systems is currently not monitored in a detailed fashion by CCF. In January 2021 CCF offered cloud resources on multiple systems and regions totaling roughly 7% of the total CPU compute capacity at CCF. The dedicated Arbutus cloud system contributes 87% of all cloud resources in the CCF infrastructure. The leading disciplines using Arbutus last year were Covid-19 and physics. Covid-19 research leading in cloud resource usage is an interesting indicator how cloud resources can be flexibly deployed for new research and needs.

Regarding **research software**, the Alliance will be publishing an extensive RS Current State document in the Summer 2021 discussing research software usage and trends in detail. The key take-away from research computing support perspective is that the software usage has a very long tail of myriad software products being used by relatively small number of users per each software package (partly since CFI has mandated CCF to support all research computing beyond

workstation level), clearly presenting its own challenges in the software maintenance, and support ecosystem.

Within CCF's **research support ticketing system** the number of tickets in 2019 was ca. 12,000, with large growth between 2016 and 2019 when users and support migrated from legacy systems with regional/institutional support to national systems with national support. The overwhelming majority of the CCF support tickets are related to general use of the ARC and CCF clusters and infrastructure. That is, a large majority of the support tickets are not related to e.g. domain specific scientific (programming etc.) needs.

**Training** is a very important activity at the CCF, including seminars, workshops, summer schools etc. It is critical for adoption, outreach, training of digital workforce, updating skills of researchers etc. In 2019-20 the CCF consortium delivered a total of 46,000 hours of received training (event hours times the number of attendees) to roughly 14,000 attendees in roughly 460 in-person events.

## Key strengths in the Canadian ARC ecosystem

### Strong ARC service provisioning

As part of its annual account application renewal process CCF queries its user base on their impressions on CCF resources and services. In 2020 85% of users of the ARC platforms were 'satisfied' or 'very satisfied' with the CCF offerings, while only 3% of respondents were either 'dissatisfied' or 'very dissatisfied'. Users from all research disciplines seem to be roughly equally satisfied with the CCF resources and services.

### Developed and refreshed ARC infrastructure

CCF and its predecessors have a 20+ year history of delivering and supporting extremely productive, high-end ARC systems through the ARC consortia to the Canadian research community. The current set of five national CCF systems was installed during the period of fall 2016 through spring 2019 as a result of CFI's Cyberinfrastructure Fund. Additional ISED funding was then used to expand and upgrade four of the systems in early 2020 and will result in a new system by fall 2021. Despite being originally envisioned as similar systems, the three general purpose systems vary significantly in size as a result of a 3x spread in their capital funding. In total, these systems represent just under $170M of federal and provincial capital investment.

Canada currently has 5 research related systems on the most recent (November 2020) Top500 rankings of world's most powerful supercomputers - Beluga, Cedar and Niagara, and two federal government systems which are used primarily for climate and meteorology. The currently fastest supercomputer in the world, Japan's Fugaku, is ca. 123 times faster than the top Canadian entry.

### Centralized service delivery

As a part of CCF's service modernization, the organization moved to a more national operations and support model including a more consistent and coherent computing and data environment and establishment of multiple national teams. Service model improvements included e.g., uniform access via centralized credentials, improved quality of the (centralized and bilingual) documentation, improved data transfer services, centralized approach to more uniform storage offerings (via standardized file system layout and policies), security overhaul, centralized software

stack, and a centralized application process for accounts and resource allocation. End-users now have a single point of contact for research computing support, while local on-campus support personnel are still available as needed. CCF and affiliates have also made improvements in workload portability across platforms. Using the same batch scheduler allows end-users to use similar job submission scripts on different systems with minor modifications. CCF has also improved system status reporting via its centralized resource publishing service that provides current information about available resources.

### A strong and dedicated community of HQP

The CCF network consists of roughly 200 HQP full-time equivalent (FTE) running the CCF operations and sites across Canada. These people provide a variety of critical services related to ARC systems administration, procurement, maintenance, networking, operations, management, planning, funding, support, research software development, data management, training, accounts and allocations management, communications, outreach etc. Many HQP analysts and sysadmins have advanced degrees and have research experience as ARC users. The ARC systems almost by definition are leading edge highly complex systems in nearly all aspects of their configuration, software and hardware stacks, operations, and use, requiring senior level expertise that takes several years of specialization to master. Maintaining the skillsets and retaining the CCF HQP is of critical importance to Canadian DRI ecosystem.

The CCF conducts a systematic post-ticket satisfaction survey to assess the satisfaction toward customer service. The survey responses speak to the high quality of service provided by the HQP people within CCF. Regarding timeliness of response, 94% of respondents rated timeliness good, or excellent. Similarly on the topic of 'solution provided', 91% were happy with the solution.

### Renewed funding commitment

The Canadian Government via the Ministry of Innovation, Science and Economic Development of Canada (ISED) clearly sees the value of DRI for Canadian society, as proven by the $572.5M 2018 budget commitment. On the Alliance side this translates to a total of $375M federal funding until March 2024, providing important (relatively) long term continuity to the DRI funding. Moreover, this restructuring also balances and centralizes the funding, covering the three key elements of a modern DRI ecosystem under one operation.

## Key challenges and opportunities in the Canadian ARC ecosystem

### Coordinated national strategic and operational planning

A more coordinated and centralized approach in the national strategic and operational planning is needed for DRI in Canada, to be enabled by a more sustained and predictable funding. This approach needs to be national in scope for e.g., increased synergies and efficiencies (e.g., in better optimized resource usage), improved interoperability, improved usability, and better utilization of HQP expertise across Canada. In case of a major fire, flooding or other major catastrophe at a host site, the CCF could potentially fully and completely lose a site and all of the data stored at that site. The planning also needs to explicitly consider not only traditional ARC, but also short-, mid- and long-term storage, off-site or cross-site backups, research data management, and research software needs holistically in one envelope, while putting added focus on underserved disciplines, audiences, and communities. Keeping in mind that not all ARC

equipment is the same, and different disciplines and communities require different types of ARC infrastructure and services.

Due to funding constraints the past and to a degree current ARC systems and service providers have not had the opportunity to focus on research data management and research software accessibility & usability, long-term storage needs in their service delivery, and the needs of wider audiences, and disciplines, e.g., humanities and social sciences. It should be noted that this situation is not necessarily due to the lack of vision or recognition by the ARC providers — it has more to do with restrictions within the funding mandates.

## Insufficient ARC supply

Considering the **CPU compute resource supply and demand** in the CCF from 2012 to 2020 the total available capacity has fluctuated within a relatively narrow range between ca 155,000 and 230,000 CPU years, while the Compute Canada (CC) Resource Allocations Competition (RAC) request based demand has grown from ca 100,000 CPU years to 450,000 CPU years. In CAGR terms the growth in the demand for CPU computing cycles was ca. 21%. In 2020 roughly 40% of the demand was met. The 270,000 CPU years of unmet demand corresponds to roughly 3.4 times Niagara supercomputer's worth of compute resources. The overall utilization of the systems is high - roughly 90% of all theoretically available cycles are used.

Among its G7 peers Canada is last when one considers aggregate total compute power in Top500. Looking into compute power relative to gross domestic product (TFlops/GDP) Canada is second last within G7. Our ARC capacity should be at least doubled in order to keep up with our peers in the middle of the pack of G7.

The demand for **GPU computing resources** has grown rapidly in the last decade to nearly 13,000 GPU years per the 2020 RAC request process, and the GPU resource supply was roughly 5x oversubscribed. Keeping in mind that many workflows and applications can not fully leverage GPUs, and underutilization and lack of optimization are of concern with GPU usage. In absolute terms the unmet GPU demand in 2020 was roughly 11,000 GPU years, equaling roughly eight current Cedar supercomputer's worth of GPU cycles. Or, to consider the scale from a different point of view, assuming modern NVIDIA V100 Volta GPU cards the total cost of the accelerator cards alone would be roughly $100M. The actual cost for fulfilling the 2020 unmet GPU need would be even higher once the cost of thousands of ARC servers, and other supporting infrastructure is included.

Compared to the above CPU and GPU capacity shortfall, the active **storage capacity** at CCF has been able to keep up with demand much better in recent years. This is a positive trend keeping in mind that storage needs are often non-transient and quite different from temporal CPU and GPU resource usage: Researchers don't expect their storage allocations to disappear in following years, and storage that is used can not be used by someone else. In 2020 the total storage capacity of 140 PB was ca. 30 PB larger than the total request, keeping in mind that much of this additional head room in the storage infrastructure is needed for efficient operation of the system. The ca. five-fold growth in demand has been roughly linear over the last five years, corresponding to ca. 39% CAGR.

Notably the above active storage analysis does not consider archival or long-term repository data supply or needs. Per its mandate CCF does not provide this kind of storage, even though the enterprise tape systems used for nearline storage at CCF could from the technical point of view

support such needs and CCF also has the required HQP expertise. Federal level substantial investment in long-term nearline and archival storage capacity and coordination with corresponding backup storage capacity is critically needed to support RDM initiatives in cohesive and sustainable fashion.

## Keeping pace with the technological and cultural diversity

Keeping pace with the technological and cultural diversity (including geographical and age diversity) is a challenge for the current Canadian ARC ecosystem and service providers. In addition to the methodological and digital research infrastructure (DRI) toolchain advances, the DRI ecosystem has become a valuable and emerging tool for non-traditional ARC users and disciplines, including e.g., social sciences, humanities, health sciences, indigenous studies etc. As more and more data and content is becoming available, these disciplines have recognized the enormous opportunities DRI systems can potentially provide for their research disciplines. In many cases these initiatives additionally include concerns related to sensitivity, privacy, ownership, and security of the data. In the context of humanities, it is not only data but also content that is getting digitized. The handling of digitized content and natively digital corpora (data from social networks, the web, etc.) reinforces the movement towards the use of computational methods already at work in these disciplines. For example, it is becoming more and more common for projects to mobilize computer processing for the automated analysis of text, voice, sound, images, or videos, either to perform mining or classification.

The needs of various underrepresented groups are being recognized better by the society, putting emphasis on addressing the needs and requirements of these communities, including e.g., racialized, LGBTQ+, and indigenous communities. If the researchers in these disciplines and communities are not familiar with modern ARC systems and tools, they may benefit from new DRI innovations, training, documentation, dedicated HQP, or tools to enable their access to DRI systems. Just purely on the traditional ARC technology side the pace of new emerging technologies is fast and to a degree diversifying, e.g., with the emergence and increasing adoption of GPU computing, new AI chip architectures, quantum computing, cloud computing (e.g., infrastructure-as-a-service IaaS, platform-as-a-service PaaS, software-as-a-service SaaS models) etc.

Canada's current DRI ecosystem has not been well equipped and funded to keep up with and address the above needs. The focus has been more on the needs of traditional ARC user communities, with some expansion in testing of new technologies (e.g., new use paradigms beyond command line access leveraging science gateways etc.) but without coordinated national level effort or funding.

## Researcher awareness and adoption of ARC

The limited researcher awareness and adoption of ARC and DRI is a major problem in Canada and globally. There are roughly 33,000 full and associate level university professors in Canada, while currently CCF lists roughly 5,500 principal investigator (PI) accounts, i.e. ca. 17% of full and associate professors have registered to use CCF infrastructure.

The fields of social sciences, humanities, psychology, business, health sciences, etc. that are currently underrepresented in the ARC community have huge potential for leveraging ARC and digital research infrastructure to advance their fields. Globally there are multiple interesting organizations participating in funding and operations of DRI for these disciplines that are

becoming more and more computerized, e.g. EGI Federation, an international e-Infrastructure providing advanced computing and data analytics services; Parthenos Virtual Research Environment (VRE), an online environment for Humanities integrating cloud storage with services and tools for the research data lifecycle; ARIADNEplus, offering cloud-based VREs for data-based archaeological research; DARIAH, a pan-European infrastructure for arts & humanities scholars; and IPERION HS, a European Integrated Research Infrastructure Platform for Heritage Science. In France the Huma-Num infrastructure provides Humanities and Social Sciences researchers not only ARC computing services, but a full portfolio of DRI services all through the research lifecycle. This infrastructure is considered a "Very Large Research Infrastructure (TGIR)" at the government funding level, and provides platforms and tools for processing, conservation, dissemination, and long-term preservation of digital research data. By not serving these disciplines at their full representative level, Canada is faced with substantial risk of lost opportunity, and being left behind by the global competition.

In addition to specific disciplines and communities that are not leveraging DRI there are researchers even within "traditional" disciplines who don't access CCF systems for a variety of reasons, e.g., researchers are not aware of the service portfolio, consider the user interfaces and usage too complicated, or perhaps have given up because their application was rejected, or they had a bad experience. Keeping in mind that there are researchers who meet their ARC needs in other ways and do not need CCF accounts and resources. These researchers can be experienced and use their own systems, customized to their needs, and working on problems that do not require high-end heavy-duty compute and storage resources.

To address researcher awareness, explicit concentrated efforts are needed to reach out to all communities that are not yet leveraging DRI. Any outreach for awareness needs to be coupled with DRI resources, both infrastructure, services, training, and support staffing so that any interested researchers will have the ability to start leveraging DRI. While raising awareness, improved access and usability technologies are needed to allow researchers who are not necessarily by nature interested or savvy with information technology and DRI. More resources are required to enable researchers to leverage ARC and DRI resources effectively for their research, via for example targeted training, support, documentation, and innovative new middleware and gateways to access DRI.

## Equity, diversity, and inclusion (EDI) and minority representation

Historically in the ARC field, equity, diversity, and inclusion (EDI) and representation of minorities have not been recognized as challenge areas, and lack of understanding and solutions addressing EDI is a major concern. The support and service delivery should include and consider the needs of indigenous communities, immigrants to Canada, researchers in rural and remote areas, early career as well as senior researchers, and disabled researchers with specific accessibility needs. As an example, the cloud computing model has enormous potential to reduce inequalities in access to resources. The CCF does not currently explicitly collect EDI data so the current situation within the CCF is not well understood. EDI should not be considered only as a separate item, in its own silo, and should instead be part of all discussions and decision making.

The DRI ecosystem should also accommodate non-native English speakers, in particular French communities, but also non-native English-speaking users (via e.g., clearly written and edited user documentation). All documentation and services should be available in both official languages, and so that the quality of translation is on par with natively written text and not at the level of semi-

automated translations. Key documentation and services should additionally be available in selected indigenous languages. Key events and conferences should include sign language and bidirectional (if not multidirectional) translation.

# 2 Methodology

This ARC Current State update report was written in the Winter 2020 – Spring 2021 in a multi-step process including consultations with representatives of Canadian ARC community. The report was written by the Alliance ARC Working Group, including

- Seppo Sahrakorpi (the Alliance Senior Analyst for ARC, Chair)

- Ghilaine Roquet (the Alliance Vice President Strategy and Planning)

- Felipe Pérez-Jvostov (the Alliance Senior Analyst for Researcher Outreach and Communications)

- Maxime Boissonneault (Team Lead, Research Support, Compute Canada / Calcul Québec)

- Chris Loken (Chief Technology Officer, Compute Ontario)

- Prof. Emmanuel Chateau-Dutier (Digital Museology, University de Montréal)

- Prof. Catherine Lovekin (Astronomy, Mount Allison U.)

- Prof. Carolyn McGregor (Health Informatics, Ontario Technology University)

- Prof. Lisa Strug (Biostatistics, University of Toronto).

The Alliance Senior Analysts Shahira Khair (the Alliance Senior Analyst for Research Data Management), and Qian Zhang (the Alliance Senior Analyst for Research Software) also assisted the ARC WG in the process. The ARC WG met weekly, also assisting and advising the Alliance in its Needs Assessment process.

Key data sources for the review report were the detailed historical Compute Canada Resource Allocation Competition (RAC) results and other historical internal resource usage data provided generously by Compute Canada. Due to the time and resource constraints, the ARC WG team could not conduct any new research or surveys to support the findings. Some aspects of such research will be conducted as a part of the Alliance's Needs Assessment and Environmental Scan projects during the first half of 2021.

Findings and observations in this document, alongside the Research Data Management and Research Software Current State Assessment publications, are meant to support and assist the Alliance analysts and management, the Alliance board, and the Alliance Researcher Council in the Needs Assessment process, leading to a new service delivery model, strategic plan and funding models for the Canadian DRI ecosystem until 2024.

# 3 Introduction

## 3.1 What is Advanced Research Computing?

In the 2017 ARC Position Paper the Leadership Council for Digital Research Infrastructure (LCDRI) defined ARC as follows: "ARC provides researchers with digital technology and expertise to help them solve research issues that are either too large or too complex for them to undertake on their own. It includes services, advice, hardware and software, all supported by highly qualified personnel (HQP), to enable research activities with significant data or computation requirements, including data acquisition, simulation, experimentation, analysis, and exploration."[1] In addition to traditional server clusters, modern ARC also includes virtual infrastructures. ARC needs to equip research and researchers in the digital environment in much the same way that higher education physically equips itself with traditional infrastructure and facilities to host research activities.

Compute Canada Federation (CCF) affiliate Westgrid defines ARC as follows: "Advanced Research Computing (ARC) is everything beyond a standard desktop workstation, which includes cloud, supercomputers / high performance computing (HPC), data management, and data storage, all in support of research."[2]  While another CCF affiliate, Compute Ontario, in 2019 defined ARC as an extension of HPC and research work done on supercomputers: "Modern research, in virtually all domains, often involves significant computational work which may not require supercomputers and massively parallel codes. Policymakers in Canada introduced the term "advanced research computing" or "ARC" to refer to the full-range of computing needs of researchers while using the term "high performance computing" or "HPC" to refer to the subset of those computing needs which can only be met on a supercomputer."[3] In this report we adopt this Canadian definition and relationship between ARC and HPC.

Comparing ARC to enterprise computing, ARC puts great value on performance and agility. Research (methods, techniques, software) and technology (including capability and cost-effectiveness) both evolve on very quick timescales (2 years can be considered a long time on some occasions). For ARC to remain competitive and relevant for leading edge research, it needs to keep-up on all fronts. Enterprise computing values stability and reliability, which is critical for e.g., payroll, email, and student registration systems since these systems and technologies advance on longer time scales. In general, it does not make a big difference if an enterprise system can suddenly process payroll 50% faster, but if ARC systems suddenly run 50% faster, then much more research can be done more quickly, satisfying the huge and growing demand.

In a similar fashion, a key distinction between ARC and regular campus IT services should be made recognizing that campus IT is aimed at production environment requiring stringent Service Level Agreement(s) (SLAs), while research computing infrastructure is in general aimed at getting as much throughput as possible for the available money, perhaps leveraging leading edge technologies that might not be as reliable as enterprise class IT services, and emphasizes

---

[1] Advanced Research Computing (ARC) Position Paper, by LCDRI (August 2017).

[2] Westgrid: What We Do https://www.westgrid.ca/about_westgrid/what_we_do (retrieved November 2020).

[3] Compute Ontario: Thinking Forward Through the Past:A Brief History of Supercomputing in Canada and its Emerging Future https://computeontario.ca/wp-content/uploads/2019/07/A-Brief-History-of-Supercomputing-in-Canada-and-its-Emerging-Future.pdf (June 2019).

flexibility in service delivery and configuration, and thus does not typically have SLAs that are as strict. ARC is also characterized by interconnectivity and mutual interaction of multiple leading-edge technologies that combined will deliver the custom deliverable required, all in an environment that is shared by large groups of users and utilized at very high capacity.

Depending on the audience and jurisdiction ARC is known with different names with slightly different emphasis: cyberinfrastructure in the US as a synonym for ARC particularly in the context of networked infrastructure, high performance computing (HPC) for higher-end ARC systems (often excluding single workstation or server scale systems) and operations, and supercomputing for the very-high end ARC systems (e.g. IBM Bluegene, and Cray supercomputers with custom hardware, I/O, and communications systems) operations. A key distinction between ARC and HPC globally (but not in Canada where HPC is considered a subset of ARC) is that the former emphasizes and focuses on scientific applications and research at scale and not necessarily performance, while the latter focuses on performance and can also include production environments e.g., in commercial sector and security agency usage. For a more enterprise- / market-oriented definition Intersect360 Research HPC-focused market research and consulting firm defines HPC having a wider scope than ARC, beyond just research computing to include also large scale production related computing, defining HPC  "as the use of servers, clusters, and supercomputers—plus associated software, tools, components, storage, and services—for scientific, engineering, or analytical tasks that are particularly intensive in computation, memory usage, or data management. HPC is used by scientists and engineers both in research and in production across industry, government, and academia." [4]

In addition to the above general considerations of ARC, the 2017 LCDRI ARC Position Paper organized ARC into a set of six core functions. [1] In the following we present these functions with additional updated commentary (*in italics*):[1]

- Computation: the use of compute resources such as cores, graphics processing units (GPUs) *or other accelerators*, and memory.

- Active storage: the data in regular use during the life of the project (in contrast to archival storage used for long-term preservation).

- Advice, support, and training: the leadership provided by Compute Canada and the regional consortia on behalf of the ARC community, and the staff expertise available to assist researchers *and to support and develop new digital approaches in research.*

- Software stewardship and support: system and application-level care, maintenance, evolution, and the development of software useful to many researchers and projects. *Notably the research software used and developed by the researchers is considered under the Research Software pillar / framework.*

- Privacy, security & authentication: stringent, detailed, or unusual requirements beyond those provided by default.

- Support of and coordination with other DRI components (network, data management, storage, software): the responsibility of ARC service providers to coordinate effectively with

---

[4] Intersect360 Research - https://www.intersect360.com/what-is-hpc (retrieved September 2020).

other components of the DRI ecosystem, *including commitment to open science initiatives as appropriate*.

In this Current State report we will follow the above slightly updated functional categorization for definition for ARC, while recognizing that the Alliance's future mandate will not be that restricted. That is, this functional definition does not include mid-term (repository), and long-term (archival) storage hardware, systems, and services, which are an integral part of the modern research computing cycle.

Interestingly, depending on the discipline, the emphasis on the ARC functionality can differ from the above technology oriented one. For example, the European Union's Common Language Resources and Technology Infrastructure (CLARIN) for humanities and social sciences focuses on the tools and research process: "In 2012 CLARIN ERIC was established and took up the mission to create and maintain an infrastructure to support the sharing, use and sustainability of language data and tools for research in the humanities and social sciences. Currently CLARIN provides easy and sustainable access to digital language data (in written, spoken, or multimodal form) for scholars in the social sciences and humanities, and beyond.  CLARIN also offers advanced tools to discover, explore, exploit, annotate, analyse or combine such data sets, wherever they are located."[5] While CLARIN is focused on language resources, DARIAH is focused on a pan-European infrastructure for arts & humanities. Its operations are based on four Virtual Competency Centres (VCCs) and constituent working groups. Each VCC has its own core mandate and focus, e.g., VCC1 is focused on CLARIN's technological foundations and e-infrastructure.[6]

## 3.2 Who uses HPC and ARC systems globally?

According to Intersect360 Research market research, HPC was globally a ca. $39B marketplace in 2019, growing at ca. 8.2% annual rate compared to 2018. The main vertical markets were Academia (17.1%), Government (25.4%), and Industry (57.5%). [7] Within the Government category, the main users of HPC were national security (12%), national research laboratories (10%), and national agencies (3%). On the industry side the main user groups were quite evenly distributed with financial services (13%), large product manufacturing (8%), biosciences (8%), energy (5%), consumer product manufacturing (5%), and retail (5%) being the largest sectors. On one hand the large majority of HPC investment is thus driven by non-academic consumption, so that the solutions the marketplace provides are not primarily tuned to serve research needs, but rather commercial and e.g., national security related use cases. On the other hand, the needs of academia can be seen to align with some of the governmental and industry needs, for example national research laboratories (10%), biosciences (8%), energy (5%), and chemical engineering

---

[5] CLARIN: CLARIN in a nutshell https://www.clarin.eu/content/clarin-in-a-nutshell (retrieved December 2020).

[6] DARIAH-EU: DARIAH in a nutshell https://www.dariah.eu/about/dariah-in-nutshell/ (retrieved February 2021).

[7] Intersect360 Research - Worldwide HPC Market 2019 Actuals, 2020-24 Forecast, Including Effects of COVID-19 https://www.intersect360.com/presentations/Intersect360%20WW%20HPC%202019%20market%20and%202020-24%20forecast.pdf (September 2020).

(4%) amount to 27% which in combination with academia (17%) adds up to 44% of the marketplace being driven by 'traditional' research computing / ARC needs.

The two new main changes in 2019 compared to earlier years were the emergence of government-led growth compared to industry, and major growth in cloud and cloud like deployments. The four main expense categories were servers (ca. $14B), software (ca. $9B), storage (ca. $6B) and services (ca. $4B) in this order. Notably the spending in cloud-based solutions showed major growth but was still less than ca. $2B overall.[8]

## 3.3 Who is involved in delivering ARC for Canadian researchers?

### Summary of local, regional, and national levels for ARC delivery in Canada

The Canadian ARC is delivered via a non-centralized network of local, regional, and national level organizations, a framework that has not changed fundamentally since the 2017 LCDRI ARC Current State[1] report.

At the local level the Canadian DRI ecosystem is vibrant with universities providing various research computing services and support, either as part of their central IT operations, libraries, or as separate independent university level research computing operations. The local growth has often been supported by national funding, for example CANARIE's Research Software and Data Management initiatives. The fundamental problem of fragmented and varied local digital research infrastructure (DRI) service delivery remains though in the Canadian academia. Some universities have relatively strong DRI operations and support, while some provide little DRI support.

At the regional level the coordination of DRI delivery is provided by WestGrid, Compute Ontario, Calcul Québec, and ACENET. The major funding provided by CFI for refreshing Canada's ARC infrastructure has further consolidated the role of the five Compute Canada federation (CCF) main host sites, owned by University of Victoria, Simon Fraser University, University of Waterloo, University of Toronto, and McGill University. It is important to note that three of the five sites (SFU, UW and McGill) are managed by distributed teams that include members outside of the owner institution. Moreover, many of the tasks required to operate these infrastructures – such as user support, documentation, software installation, monitoring, scheduling - are handled by national teams which include people located across the institutions that hire CCF staff members. This has left the ACENET regional consortia without its own main CCF ARC host site, although it operates its Siku high-performance cluster at Memorial University.

At the national level the key change has been the launch of the Digital Research Alliance of Canada (the Alliance) in 2019. The Alliance's mandate from ISED is to coordinate and fund Canada's emerging DRI ecosystem, including not only ARC, but also research software (RS), and data management (DM).[9] This approach will allow much more centralized, coherent, and predictable funding for the DRI ecosystem as a whole. As a part of the Alliance assuming its responsibilities, Compute Canada's operations will be absorbed and assumed by the Alliance as

---

[8] Intersect360 Research - Worldwide HPC Market 2019 Actuals, 2020-24 Forecast, Including Effects of COVID-19 https://www.intersect360.com/presentations/Intersect360%20WW%20HPC%202019%20market%20and%202020-24%20forecast.pdf (September 2020).

[9] Alliance: Background https://engagedri.ca/about-engage-dri/background (retrieved December 2020).

of April 1ˢᵗ 2022. On the data management side, CARL Portage operations and funding are already under the Alliance's umbrella, while CANARIE's data management and Research Data Canada (RDC) operations will be brought to the Alliance by April 2022. On both data management and research software sides the Alliance will incorporate key teams, policies, procedures, and operations from CANARIE's and Compute Canada's data management and research software initiatives.

General layout of current government funded DRI delivery mechanisms in Canada

**CURRENT NATIONAL STRUCTURE OF DRI ECOSYSTEM**



Innovation, Science, and Economic Development Canada, 2020 with adjustments by the Digital Research Alliance of Canada (the Alliance)

**Figure 1:** Current National Structure of DRI Ecosystem

The fundamental structure of DRI funding in Canada has not changed since 2017. The Government of Canada's ISED is still the main funder of the DRI ecosystem and operations in Canada, via CFI's Innovation Fund (IF) and Major Science Initiatives (MSI) funds, or via CANARIE. The funding for ARC systems and operations is relatively well established, with CFI IF and MSI funds flowing to individual host sites and universities with guidance from CC. The matching formula is still 40/60, i.e., 40% of the funding comes from CFI, while the remaining 60% match comes from sources such as provinces, universities and in-kind donations. Any local ARC infrastructure is in general not eligible for CFI funding unless a strong case is made for e.g., real-time, data sensitivity or edge computing needs. In general ARC infrastructure should be located on main CCF sites and be e.g., a contributed system addition to the existing infrastructure. The funding model is further complicated by the differences in capital (IF) and operations (MSI) funding formulas. This separation is also not well-equipped to address the emerging cloud computing trend that shifts the expenses from capital to operational expenses.

Going forward, one of the key mandates of the Alliance is to bring all three key DRI components, ARC, RS, and DM under one funding and planning umbrella, including potentially new or improved funding models. The Alliance's mandate from ISED also includes major new five-year funding for DRI, totaling $375M until March 2024.

## Compute Canada Federation affiliated organizations

Compute Canada Federation consists of four regional partner organizations ACENET, Calcul Québec, Compute Ontario and WestGrid. The key characteristics of these organizations are as follows:

### ACENET

Established in 2003, ACENET is the Digital Research Infrastructure organization in Atlantic Canada; a partnership of 14 universities and colleges in the region representing almost all post-secondary institutions. It provides advanced research computing (ARC) infrastructure, technical support, and digital skills development to over 1000 academic researchers, post-secondary students, and industry R&D in Atlantic Canada.

It houses regional ARC infrastructure at Memorial University, which holds a special national system designation from CFI that enables it to incorporate researcher contributed systems. These are systems dedicated to individual research groups, where in return for ACENET operating, maintaining, and managing them, excess compute cycles are returned to the shared resource pool, thereby maximizing the use of ARC resources and funding. There are currently four contributed systems attached to ACENET's infrastructure, with an additional five in the procurement process.

ACENET is a distributed organization with 20 staff. It is governed by a Board of Directors made up of the Vice Presidents Research (or designate) at its six host institutions: Dalhousie University, Memorial University, St. Francis Xavier University, Saint Mary's University, the University of New Brunswick, and the University of Prince Edward Island. ACENET's Executive Team is advised by a 10-member Research Directorate made up of active cross-disciplinary researchers from several Atlantic Canada's member institutions. Funding for ACENET is provided by CFI, ACOA and the four Atlantic provinces. [10]

### Calcul Québec

In Québec, Compute Canada's regional partner is Calcul Québec, an incorporated not-for-profit consortium composed of eleven Québec universities. The member universities have pooled together their local ARC investment and resources to form the coalition. More than 550 research groups and around 1925 users take advantage of these resources.[11] Funding for Calcul Québec is from province of Québec and CFI. Network connectivity is provided by Réseau d'informations scientifiques du Québec (RISQ) and CANARIE.[12] Calcul Québec enlists over 40 HQP staff, and

---

[10] Ines Hessler, CTO of Acenet, private communication (July 2021).

[11] Calcul Québec: Who are we? https://www.calculquebec.ca/en/about-us/who-are-we/ (retrieved November 2020).

[12] Calcul Québec: Partenaires https://www.calculquebec.ca/en/about-us/partenaires/ (retrieved November 2020).

the mature governance model includes a13-member Board, 10-member Scientific Council, Operations Council, and Technology development, operations, and research support council.[13]

Calcul Québec hosts one of CCF's main national systems[170], Béluga[175], a general-purpose ARC cluster owned by McGill University, located at École des technologies supérieures (ETS) and operated by a team distributed across the consortium. As of January 2021, Calcul Québec also operates an earlier generation Helios[14] supercomputer located at Université Laval, and MP2, located at Université de Sherbrooke.

**Compute Ontario**

Compute Ontario was incorporated in 2014 as a not-for-profit corporation, with a mandate to coordinate advanced computing in Ontario. [15] Its foundation is built on two decades of prior work by the province and the existing high performance computing consortia in Ontario. Compute Ontario is funded by the Ministry of Colleges and Universities (MCU) of Ontario and regional network connectivity is provided by ORION in collaboration with CANARIE.

Compute Ontario's partner consortia are the Centre for Advanced Computing, SHARCNET, SciNet, and HPC4Health.[16] Compute Ontario works with the consortia to centralize strategy and planning for advanced computing assets, including hardware, software, data management, storage, security, connectivity, and Highly Qualified Personnel.

SHARCNET is a consortium composed of 18 colleges, universities (Western University is the lead university), and research institutes operating a network of high-performance computer clusters across southwestern, central and northern Ontario. SHARCNET hosts Graham, a general-purpose heterogeneous supercomputer located at one of CCF's main national sites at University of Waterloo campus.

SciNet consortia is led by University of Toronto. SciNet hosts Niagara, a massively parallel homogenous supercomputer at University of Toronto on one of CCF's main national sites.

**Westgrid**

WestGrid is a coalition of seven member institutions from British Columbia, Alberta, Saskatchewan, and Manitoba. The member institutions receive funding for WestGrid's operations and maintenance needs through CFI's MSI funding. The provincial partners and academic institutions provide matching funding for all CFI grants.[17]

WestGrid's facilities are connected through a dedicated WestGrid core network leveraging CANARIE's national network infrastructure and regionally National Research and Education Network (NREN) partner networks (British Columbia (BCNET), Alberta (Cybera), Saskatchewan

---

[13] Calcul Québec: Gouvernance https://www.calculquebec.ca/a-propos/gouvernance/ (retrieved November 2020).

[14] Compute Canada: Hélios https://docs.computecanada.ca/wiki/H%C3%A9lios/en (retrieved November 2020).

[15] Compute Ontario: About Compute Ontario https://computeontario.ca/about-compute-ontario/ (retrieved November 2020).

[16] Compute Ontario: Partners https://computeontario.ca/partners/ (retrieved November 2020).

[17] WestGrid: What we do https://www.westgrid.ca/about_westgrid/what_we_do (retrieved November 2020).

(SRNET) and Manitoba (MRnet)).[18] WestGrid affiliate Simon Fraser University hosts and manages Cedar, a general purpose heterogenous supercomputer located at one of CCF's main national sites at SFU campus. Additionally, considering key systems within the WestGrid's regional envelope, the University of Victoria hosts Arbutus, a general-purpose ARC cloud system as part of the CCF ARC resource offering.[170]

## Academic non-CCF organizations

The Canadian DRI ecosystem is very vibrant with existing organizations in constant change and innovation and new endeavours being launched with varied scope and service offerings. In the following we provide **a non-comprehensive sampling** of such Canadian non-CCF institutions categorized under 'infrastructure', 'service and platform', and 'research and training' providers. These categories are not mutually exclusive as some of the larger institutions provide services in all categories. Moreover, some of the operations below are provided by CCF affiliated organizations, but are not part of the main CCF national, publicly available, offering located on the five main host sites.

### *Infrastructure Providers*

Infrastructure providers operate and focus on physical DRI infrastructure and additionally often provide DRI services, platforms, or training or conduct research.

### ARC UBC

The University of British Columbia's (UBC) Advanced Research Computing (ARC) is a major university-owned non-CCF ARC infrastructure operation. It consists of Sockeye compute cluster, and Chinook object storage. Sockeye's original funding of $7.9M materialized in 2018, and in early 2020 the system was further enhanced with an additional investment of $10.1M.[19] Currently it consists of over 16,000 CPU cores, 200 GPUs, InfiniBand EDR interconnect,[20] and up to 20 PB of storage, distributed between Vancouver and Okanagan campuses. The system also has 192TB of flash storage for fast temporary I/O needs. Looking at only CPU core counts, the system thus is roughly half the size of CCF's Beluga general purpose compute cluster located in Montreal. UBC's ARC services are available to researchers with faculty appointment at UBC or UBC Principal Investigators with particular emphasis to support new faculty and researchers who require their data to be hosted locally and not on shared CCF storage servers.[21]

### CAC

The Centre for Advanced Computing consortia lists Queen's University (the lead institution), Carleton University, University of Ottawa, and the Royal Military College of Canada as its members. It particularly specializes in secure, advanced computing resources and support for academic and medical researchers. CAC supports over 400 Canadian research groups totaling

---

[18] WestGrid: Partners https://www.westgrid.ca/about_westgrid/members-partners (retrieved November 2020).

[19] UBC ARC: Enhancing Support for Advanced Research Computing https://arc.ubc.ca/enhancing-support-advanced-research-computing (retrieved January 2021).

[20] UBC ARC: Sockeye - Detailed Technical Specifications https://arc.ubc.ca/sockeye-techspecs (retrieved January 2021).

[21] UBC ARC Sockeye https://arc.ubc.ca/ubc-arc-sockeye (retrieved January 2021).

some 2100 researchers working in a variety of fields.[22] CAC is not one of the main CCF national sites. Its Frontenac and Katarokwi platforms are available for CCF affiliated researchers but not free of charge as the main CCF systems.

**CRDCN**

The Canadian Research Data Centre Network (CRDCN) is a partnership between a consortium of Canadian universities and Statistics Canada, headquartered at McMaster University. Its mission is, through its Research Data Centre (RDC) Program, to provide researchers access to social, economic and health confidential microdata. Currently such access is provided in secure office spaces with strictly controlled workstation and server facilities, located on university campuses across the country. CRDCN's core funding comes from a mix of SSHRC/CIHR directed grant and from a CFI MSI award. Host universities and Statistics Canada also provide material cash and in-kind support to CRDCN .[23] In 2017, it secured substantial funding via CFI IF mechanism to create a national level centralized HPC infrastructure to meet the projected growing data processing and storage needs in the network.[24] More recently the HPC / virtual Research Datacenter (vRDC) platform design specification has been expanded to include remote access capability,  i.e. access to the central resources from outside the secure RDC office spaces.[25] This involves a detailed security design and review of the system in order to meet the strict cyber-security requirements that apply to the Protected B micro-data files[26] (as defined by e.g. Government of Canada Treasury Board, Communications Security Establishment, and Statistics Canada IT Security and Microdata requirements).

**CYBERA**

Cybera is a not-for-profit CANARIE affiliate in Alberta responsible for running the CYBERANET, part of Canada's National Research and Education Network (NREN).[27] It also provides ARC infrastructure and services for CANARIE's Digital Accelerator for Innovation and Research (DAIR) program that provides small and medium sized enterprises (SME) access to cloud testbeds and technology. It also hosts and provides access to Rapid Access Cloud, a free cloud service for Alberta academics and SMEs who are not eligible for DAIR. CYBERA also collaborates with Pacific Institute for Mathematical Sciences (PIMS) in providing the very popular Syzygy Jupyter Science Gateway to Canadian researchers. Cybera's revenue in FY2018-19 was ca. $5.9M and had 39 staff.[28]

---

[22] Centre for Advanced Computing: What is CAC? https://cac.queensu.ca/about_us/ (retrieved November 2020).

[23] CRDCN: About the CRDCN https://crdcn.org/about-crdcn (retrieved November 2020).

[24] McMaster University: McMaster Data Center receives $2.7M to support research infrastructure development for economic, social & health data https://www.economics.mcmaster.ca/news/mcmaster-data-center-recieves-2-7m-for-study-of-economic-social-and-health-data (October 2017).

[25] CRDCN 2019-24 Strategic Plan https://crdcn.org/sites/default/files/strategic_plan_0.pdf (June 2019).

[26] CRDCN September 2020 Newsletter https://us4.campaign-archive.com/?u=c3b811df1cf083f6ae6fb612b&id=5a4556d80e (retrieved November 2020).

[27] CYBERA: Services https://www.cybera.ca/services/ (retrieved January 2021).

[28] CYBERA: 2018-2019 Annual Report https://www.cybera.ca/wp-content/uploads/2020/03/Cybera_Annual_Report_2018-19.pdf (October 2019).

**HPC4Health**

HPC4Health is a consortium of SickKids and University Health Network in the Toronto area, building a cloud based secure compute engine for clinical research. The services are primarily available to members of these two institutions, while outside organizations can access services on a cost-recovery basis.[29] HPC4Health hosts a 7000 CPU-core OpenStack based cloud infrastructure so that each health care institution can access their own, fully private cloud while enjoying the benefits of resource pooling. Each participating institution is guaranteed a minimum amount of CPU cores when they need it (80% of their contribution), allowing for critical and time sensitive computing needs. The remaining 20% is shared allowing all users to leverage underutilized capacity.[30]

**NRC**

Government of Canada's National Research Council Canada (NRC) is Canada's largest federal research and development organization. It both runs its own research operations, collaborates with Canadian academia and research institutions, and funds research operations and small- and medium sized enterprises (SMEs) and industries in Canada. In fiscal year 2019 NRC had 4109 FTE and showed $184M in revenue with $1214M of total expenditures.[31] NRC teams operate and leverage multiple ARC resources, both internally and in collaboration with Shared Services Canada. But even an organization of NRC's scale can suffer from insufficient ARC services, for example in 2018 NRC's Security and Disruptive Technologies (SDTech) Research Centre had insufficient and outdated ARC capacity, with plans to potentially collaborate with Compute Canada for improvements.[32] Unfortunately, public details of NRCs ARC operations are not readily available.

**Ouranos**

Ouranos is a not-for-profit domiciled in Montreal, Quebec, focusing on climate change and its impacts, as well as relevant socio-economic and environmental vulnerabilities to drive policy and adaption strategies. It employs over 50 people and is involved in 13 scientific programs and over 100 projects. Key members of Ouranos are Province of Quebec, Hydro Quebec, UQAM, McGill University, Universite Laval, and Environment and Climate Change Canada. Annual revenue is between $8M and $12M, and the funding comes from a variety of sources including the Province of Quebec.[33] In 2015 Ouranos had three Cray supercomputers and collaborated with Calcul Quebec for ARC resources.[34]

---

[29] HPC4Health: Accessing Our Services http://www.hpc4health.ca/services.html (retrieved November 2020).

[30] HPC4Health: Overview http://www.hpc4health.ca/overview.html (retrieved November 2020).

[31] National Research Council Canada: Annual Report 2019-20 https://nrc.canada.ca/sites/default/files/2020-08/annual-report-2019-2020.pdf (August 2020).

[32] NRC Office of Audit and Evaluation: Evaluation of NRC's Security and Disruptive Technologies Research Centre - Final Report https://nrc.canada.ca/sites/default/files/2019-03/sdtech_report_2018_e.pdf (January 2018).

[33] Ouranos https://www.ouranos.ca/en/ouranos/ (retrieved November 2020).

[34] Ouranos: PLAN STRATÉGIQUE 2014-2020 (December 2014).

**SciNet4Health**

In September 2020, the University of Toronto and SciNet announced[35] a new SciNet4Health initiative that "will allow researchers and clinician scientists at U of T and its partner hospitals to access and analyze massive databases of patient health information – in a secure way that protects patients' privacy – using technologies such as machine learning." The core system has a theoretical performance peak of one petaflop and consists of 20 compute nodes, each with AMD EPYC processors and eight AMD Radeon Instinct GPU accelerators, donated by AMD. The system will be hosted at the main SciNet data center next to the Niagara supercomputer. The initiative will leverage experiences from HPC4Health for procedures and protocols. The two organizations are planning to work together for delivering health care related ARC in the Toronto area.

**Siku**

Acenet hosts Siku, a 2300 CPU core high-performance computer cluster commissioned in 2019 and located at Memorial University in St. John's, Newfoundland. It is not one of the CCF national systems, and the location is not one of CCF's main sites. Siku is funded in large part by ACOA with the goal of generating regional economic benefits through engaging the local industries, while also supporting academic research in the Atlantic region. The system is only accessible to selected clients, including both industrial researchers and academic research groups. Priority is given to industrial users, allowing academic users to use the remaining resource free-of-charge. Thanks to this industry aspect, the funding model is self-sustaining. It has both traditional ARC batch system access, and cloud computing interface.[36]

**SOSCIP**

Southern Ontario Smart Computing Innovation Platform, headquartered at Toronto's MaRS Discovery District and established in 2012, is a coalition between 15 Ontario post-secondary institutions, IBM Canada, and a variety of Ontario based small- and medium-sized businesses. It provides eligible projects fee-for-service based access to GPU, massively parallel, and cloud computing resources.[37] The eligible projects need to be an industry-academic collaboration, have advanced computing needs, and have 'clear and realizable commercialization objectives'. The academic collaborator must be a principal investigator in a SOSCIP member institution, while the commercial partner has to be based in southern Ontario.[38] By 2020 SOSCIP has worked with over 120 Ontario SMEs and has supported over 195 R&D projects.[39] SOSCIP's current GPU platform 'Mist' was launched in 2020 and consists of 54 IBM Power9 nodes with four Nvidia V100 GPU cards each, connected via Mellanox InfiniBand EDR. It is located next to SciNet's Niagara

---

[35] University of Toronto: U of T and AMD launch supercomputing program dedicated to big-data health research https://www.utoronto.ca/news/u-t-and-amd-launch-supercomputing-program-dedicated-big-data-health-research (retrieved September 2020).

[36] ACENET: Siku https://www.ace-net.ca/wiki/Siku (retrieved November 2020).

[37] SOSCIP https://www.soscip.org/ (retrieved November 2020).

[38] SOSCIP: Project Requirements https://www.soscip.org/project-guide/ (retrieved November 2020).

[39] SOSCIP: SOSCIP By the Numbers https://www.soscip.org/soscip-by-the-numbers/ (retrieved November 2020).

massively parallel supercomputer and shares the user file system with Niagara.[40] The massively parallel resource is a 2880-core equivalent allocation of SciNet's Niagara supercomputer. The cloud analytics OpenStack-based platform is designed to be a self-service big data analytics system. It consists of over 4600 CPU cores in a mixed x86 and PowerPC cluster, and over 70 Nvidia GPUs.[41]

**SSC/ECCC**

Shared Services Canada's (SSC) hosts Environment and Climate Change Canada's (ECCC) main ARC computing systems. Currently SSC operates two Cray supercomputers for ECCC, Banting and Daley. Banting was commissioned in 2017, while Daley came online in 2020. [42] The machines support ECCC's weather modelling and forecasting services, [43] and do not seem to be available to Canadian academic researchers in general.

*Service and Platform Providers*

Service and platform providers focus on providing DRI services, and operating platforms, and additionally often also provide training or conduct research, but differentiating from the infrastructure providers do not primarily own or operate their own ARC infrastructure.

**CADC**

NRC's Canadian Astronomy Data Centre (CADC) was established in 1986 and is located at NRC's Herzberg Astronomy and Astrophysics (HAA) Research Centre in Victoria, BC.[44] Its mandate is to host Canadian telescope data and to operate its science platform for data-intensive astronomy. CADC offers cloud computing, user-managed storage, group management, data publication services, and permanent storage for major data collections. CADC does not have its own major ARC infrastructure, and rather uses services provided in collaboration with Shared Services Canada, Compute Canada, CANARIE, and universities via CFI funding. A 2016 survey[45] of HAA's operations found that CADC operations were hindered by outdated IT infrastructure, and limited network capabilities, and that "transfer of HAA's IT infrastructure to Shared Services Canada (SSC) has had major impacts on the Portfolio's ability to plan, implement and acquire new IT equipment and networking capacity." CADC hosts the Canadian Advanced Network for Astronomical Research (CANFAR), and the same survey also found that "CANFAR's cloud

---

[40] SciNet: Mist GPU cluster https://www.scinethpc.ca/mist/ (retrieved November 2020).

[41] SOSCIP: SOSCIP's Advanced Computing Platforms https://www.soscip.org/platforms/ (retrieved November 2020).

[42] Shared Services Canada: High Performance Computing https://www.canada.ca/en/shared-services/corporate/data-centre-consolidation/high-performance-computing.html ; and High Performance Computing environment upgraded to support digital government https://www.canada.ca/en/shared-services/campaigns/stories/hpc-upgrade.html (retrieved September 2020).

[43] ECCC: Weather Analyses and Modelling https://weather.gc.ca/mainmenu/modelling_menu_e.html (retrieved November 2020).

[44] Canadian Astronomy Data Centre: About the CADC https://www.cadc-ccda.hia-iha.nrc-cnrc.gc.ca/en/about.html (retrieved January 2021).

[45] NRC: Evaluation of NRC Herzberg Astronomy and Astrophysics (HAA) Portfolio https://nrc.canada.ca/sites/default/files/2019-03/haa_evaluation_report_2016_e.pdf (November 2016).

computing solution is a key evolution in how the CADC serves the Canadian astronomy community."

## CANFAR

Canadian Advanced Network for Astronomical Research (CANFAR) was established to provide Canadian astronomers with a variety of ARC based services supporting their data-intensive research. The integrated suite of services includes research data management, user-managed storage, cloud processing, and specialized visualization and analytics services.[46] CANFAR does not have its own ARC infrastructure, and rather uses Compute Canada's Openstack servers, in particular the Arbutus cluster.[47]

## CMC

CMC Microsystems is a not-for-profit organization managing Canada's National Design Network (CNDN). CNDN is a national network of 10,000 academic participants and 1,000 companies focusing on research and innovation in micro-nanotechnologies.[48] The Network has received funding from CFI's Major Science Initiatives (MSI) grant. CMC offers over 50 computer-aided-design (CAD) software tools via its CADpass licensing platform, and additional supporting compute cluster and cloud platforms for Canadian academics, all under CAD branding.[49] CMC provides CCF users with a variety of commercial software licenses that are used on CCF infrastructure. In addition to CAD, CMC also provides academic partners multi-project wafer services and related fabrication services, and related tools for testing and demonstration needs.

## GenAP

The Genetics and Genomics Analysis Platform (GenAP) is a computing infrastructure and software environment for life science researchers. It was established in 2015 with funding from CANARIE, GenomeQuebec, CFI, and NSERC. GenAP offers turn-key Web applications running on Compute Canada Cloud (Arbutus) and HPC infrastructure.[50]

## LINCS

Linked Infrastructure for Networked Cultural Scholarship (LINCS) project was established in April 2020 with the goal of creating semantic web infrastructure to convert large datasets into an organized, interconnected, machine-processable set of resources for Canadian cultural research.[51] The key features will be the subsystems for 1) converting and interlinking the varied data sources, 2) storing the data, and 3) accessing the data with the ability to filter, analyze, annotate, and edit automatically created semantic results.[52] The project has dozens of university, private sector, and institutional partners from both Canada and the US. Primary funding source is

---

[46] CANFAR Portfolio https://www.canfar.net/assets/CANFAR_portfolio.pdf (May 2016).

[47] CANFAR: OpenStack Cloud https://www.canfar.net/en/docs/openstack_cloud_portal/ (retrieved January 2021).

[48] CMC Microsystems: About Us https://www.cmc.ca/about-us/ (retrieved February 2021).

[49] CMC Microsystems: CAD https://www.cmc.ca/cad/ (retrieved February 2021).

[50] GenAP: Genetics & Genomics Analysis Platform: Introduction to GenAP https://genap.ca/p/help/introduction (retrieved January 2021).

[51] LINCS https://lincsproject.ca/ (retrieved January 2021).

[52] LINCS: Research Data Infrastructure https://lincsproject.ca/development/ (retrieved January 2021).

CFI Cyberinfrastructure Initiative's $2M grant in 2019,[53] while the backend infrastructure will be provided by Compute Canada Federation host sites.

**OHDP**

In the summer and fall 2020 Ontario launched the Ontario Health Data Platform, providing researchers integrated access to Covid-19 related data. The platform focuses on health-related data and addresses the related security and privacy concerns. Key objectives are to provide insights in detection of Covid-19 in populations, discovering risk factors, predicting outbreaks, evaluating effectiveness of treatment measures, and optimize resource allocation.[54] On the methodology side OHDP emphasizes artificial intelligence and machine learning, providing access to dedicated ARC capability.[55] The platform is a collaboration between Ontario Ministry of Health, Compute Ontario, Queen's University, Vector Institute and other Ontario research institutions.[56] The funding comes from Ontario MOH.

**Syzygy**

A major recent Canadian example of cloud-based platform-as-a-service (PaaS) services is Syzygy, a collaboration between Pacific Institute for the Mathematical Sciences (PIMS), Compute Canada, and Cybera, initiated in 2017, providing Canadian researchers Jupyter Notebook based computing resources free of charge.[57] End-users can login using their university credentials if their university collaborates with Syzygy or Google accounts, and do code development and run light production and test runs. The platform has been a 'catastrophic success' in its founder's words[58], and in late 2020 had over 34,000 users.

*Research and Training Providers*

Research and training providers focus on conducting research or providing training on their specific disciplines. They can also potentially operate platforms or infrastructure specific to their discipline.

**Amii**

Alberta Machine Intelligence Institute (Amii) was established in 2002 and focuses on artificial intelligence and machine learning. It is one of three main AI institutions in Canada. Amii is located

---

[53] University of Guelph: U of G-Led Network Gets $2 Million to Link Cultural Researchers https://news.uoguelph.ca/2019/02/university-of-guelph-led-network-gets-2-million-to-link-cultural-researchers/ (February 2019, retrieved January 2021).

[54] Ontario News Release: Province Developing New Health Data Platform to Help Defeat COVID-19 https://news.ontario.ca/en/release/56659/province-developing-new-health-data-platform-to-help-defeat-covid-19 (April 2020).

[55] OHDP: About https://computeontario.ca/covid-19-health/about-ohdp/overview/ (retrieved November 2020).

[56] OHDP: Project Team https://computeontario.ca/covid-19-health/about-ohdp/project-team/ (retrieved November 2020).

[57] Syzygy.ca https://syzygy.ca/# (retrieved December 2020).

[58] James Colliander at Berkeley Computing, Data Science, and Society's 2020 National Workshop on Data Science Education – National Scale Interactive Computing https://data.berkeley.edu/academics/resources/data-science-education-resources/2020-national-workshop-data-science-education (retrieved December 2020).

at the University of Alberta and was incorporated as a not-for-profit in 2017. It currently lists over 25 researcher Fellows and over two dozen staff.[59] Amii is funded by Alberta Innovates, CIFAR, Province of Alberta, and University of Alberta.[60]

## IQC

University of Waterloo's Institute for Quantum Computing (IQC) is a research institute established in 2002 with 32 affiliated faculty members and over 300 researchers working on developing new quantum technologies.[61] It is funded and supported by Mike and Ophelia Lazaridis, the Government of Canada, the Government of Ontario and the University of Waterloo, and attracted over $30M in research funding in fiscal year 2019-2020.[62] IQC does not host or operate any client facing quantum computing system, and rather focuses on researching and advancing quantum computing technologies.

## IVADO

The mission of the Institut de Valorisation des Données (IVADO) is to expand scientific and industry-based talent in digital intelligence (including data science, artificial intelligence and operations research), and to accelerate adoption of digital intelligence.[63] It was established in 2016 when it received a major $93.6M grant from the Canada First Research Excellence Fund.[64] In 2019 it awarded $4.5M to research projects supporting over 40 projects and employing over 350 people.[65] It has 95 active industrial members and over 1400 members in its scientific community. IVADO is a major provider of training in the use of digital intelligence tools, with over 750 people trained in 2019.

## MILA

Montreal Institute for Learning Algorithms (MILA) is a not-for-profit partnership between Université de Montréal and McGill University in Montreal, Quebec. Polytechnique Montréal and HEC Montréal are also closely linked with MILA. It was established in 1993 and incorporated in 2017. MILA focuses on artificial intelligence and machine and deep learning, bringing together over 500 researchers.[66] It is one of the three main AI research institutions in Canada. In Fiscal Year 2018-19 MILA's revenue was ca. $7M, with most of the income coming from government grants ($6M), funded by CIFAR, and Quebec's Ministry of Economy and Innovation.[67]

---

[59] Amii: Our People https://www.amii.ca/about/our-people/ (retrieved November 2020).

[60] Amii: Our Story https://www.amii.ca/about/our-story/ (retrieved November 2020).

[61] University of Waterloo: About Institute for Quantum Computing https://uwaterloo.ca/institute-for-quantum-computing/about (retrieved January 2021).

[62] IQC: Annual Report April 1, 2019 – March 31, 2020 https://uwaterloo.ca/institute-for-quantum-computing/sites/ca.institute-for-quantum-computing/files/uploads/files/iqc_report_to_ised_2019-2020-eng_aug_2020.pdf (July 2020).

[63] IVADO https://ivado.ca/en/ivado/ (retrieved January 2021).

[64] HEC Montreal: IVADO receives a $93.6 M grant from Canada First https://www.hec.ca/en/news/2016/IVADO-receives-93-6-M-grant-from-Canada-First.html (September 2016, retrieved January 2021),

[65] IVADO: 2019 Activity Report https://ivado.ca/rapport-activites-2019-EN/IVADO_RapportActivites2019_ENG_v3_web.pdf (February 2020).

[66] MILA: About MILA https://mila.quebec/en/mila/ (retrieved November 2020).

[67] MILA: Annual Report from April 1, 2018 to March 31, 2019 https://mila.quebec/wp-content/uploads/2020/01/Mila-Annual-Report-2018-2019.pdf (January 2020).

**OICR**

Ontario Institute for Cancer Research (OICR) is a not-for-profit research collaborative located at MaRS centre in Toronto, Ontario. It focuses on cross-disciplinary cancer research in fields such as genomics, immuno-oncology, informatics, drug discovery, and molecular pathology. OICR has partners in health care, research, government, and the private sector.[68] It employs over 300 people and supports close to 2000 HQP researchers in Ontario. It does not have large ARC systems on-site, but has several ARC initiatives, including e.g., Computational Biology Program focusing on genomics related development of new algorithms, software, and visualization tools for large datasets.[69] OICR is the second largest funding agency for cancer research in Canada. Its revenue in fiscal year 2019-20 was ca. $85M, funded primarily by Ontario Ministry of Colleges and Universities.[70]

**Vector Institute**

Vector Institute focuses on AI, ML and DL fields. It was incorporated in 2017 as a not-for-profit with the help of University of Toronto. It encompasses a community of over 500 researchers, and over 1000 Master's students, representing over two dozen Ontario academic institutions.[71] In fiscal year 2019-20 Vector's revenue was roughly $27M, out of which province of Ontario provided ca. $10M and the federal government ca. $5.5M. In addition to the government funding, forty-seven companies sponsored Vector's operations for close to $9M. Notably Vector's funding has large annual variations due to the province of Ontario providing its funding on a front-loaded basis. On the expense side Vector's total expenses in fiscal year 2019-20 were roughly $18M, reflecting the general size of annual operations more accurately.[72] Vector Institute operates its own $6M AI computing infrastructure, consisting of 1163 GPUs distributed among nearly 200 servers, including 11 large-memory CPU only compute nodes.[73]

## Commercial organizations

A multitude of commercial organizations provide ARC computing services to Canadian researchers. Cloud computing has emerged as the primary mechanism for accessing external ARC resources thanks to its relative ease of use, low cost of entry, and pay-as-you go economic model. While on the risk side uncertainties related to budgeting, and the current concentration of the sector in the hands of a few players exposes users to the risk of uncontrolled price changes. The three big cloud providers are Amazon's AWS, Microsoft's Azure, and Google's GCP, capturing ca. 88% of global HPC deployments in 2019 according to InsideHPC.[74] All three of

---

[68] OICR: About Us https://oicr.on.ca/about-us/ (retrieved November 2020).

[69] OICR: Computational Biology https://oicr.on.ca/research-portfolio/computational-biology/ (retrieved November 2020).

[70] OICR: Ontario Institute for Cancer Research Statement of Revenue and Expenses and Changes in Net Assets https://oicr.on.ca/wp-content/uploads/2020/09/OICR-Financials-1920.pdf (September 2020).

[71] Vector Institute: April 2019 – March 2020 Annual Report https://vectorinstitute.ai/wp-content/uploads/2020/11/vector_annual-reportv8.pdf (November 2020).

[72] Vector Institute: Financial Statements for the Year Ended March 31, 2020 https://vectorinstitute.ai/wp-content/uploads/2020/11/2020-fs-final-vector-institute.pdf (November 2020).

[73] Vector Institute Position Paper submission to NDRIO: Canada's Future DRI Ecosystem: AI Research Needs https://engagedri.ca/wp-content/uploads/2020/12/Canada%E2%80%99s-Future-DRI-Ecosystem_-AI-Research-Needs.pdf December 2020).

[74] InsideHPC White Paper: Cloud Adoption for HPC: Trends and Opportunities https://insidehpc.com/white-paper/cloud-adoption-for-hpc-trends-and-opportunities/ (November 2019).

these cloud providers are also certified by the Government of Canada to host Protected B classified data in selected service offerings.[75]

## Amazon AWS

Amazon is the original and quintessential cloud service provider, both in general and for ARC. It launched its Amazon Web Service (AWS) offering initially in 2002 and moreover the Elastic Compute Cloud (EC2) cloud instance offering in 2006. Since then, service portfolio has grown astronomically and now includes full suite of services[76] for networking, routing, storage, and ARC including high-speed interconnects, large memory nodes, GPU computing, FPGA, and even quantum computing (Amazon Braket)[77]. Amazon captured ca. 46% of the global HPC cloud deployments in 2019.[72] The scale of Canadian academic use is not known at the moment but can be assumed to be substantial thanks to the relative ease of use and dynamic nature of this service.

## Google Cloud Platform

Like AWS, Google Cloud Platform (GCP) is a very mature cloud service offering a variety of Infrastructure-as-a-Service (IaaS), Platform-as-a-Service (PaaS), and serverless services.[78] Google's equivalent of Amazon's EC2 is Google Compute Engine (GCE) that provides CPU, GPU, and large memory compute nodes for ARC use cases, but does not provide FPGA based computing. In 2019 GCP captured ca. 18% of the global HPC cloud deployments.[72] Google has strong focus on AI and ML by not only providing a dedicated Google AI Platform,[79] but also by developing in-house a Google TPU ML processor that is available for general use via Google Cloud TPU cloud service.[144]

## Microsoft Azure

Microsoft's Azure is the third major ARC cloud provider with ca. 24% of the global ARC cloud deployments in 2019.[72] Similar to AWS, Azure's ARC IaaS service include CPU, FPGA, GPU, and high speed interconnect optimized compute instances.[80] Quantum computing development and testing full-stack environment is also available as Azure Quantum as of February 2021.[81] Notably for high-end ARC needs, Azure also has Cray XC and CS supercomputers available as part of its cloud offering. Although branded as part of Azure cloud, the arrangement is more like

---

[75] Government of Canada Cloud Brokering Service: GC Cloud Providers (Protected B) https://cloud-broker.canada.ca/s/central-provider-page-v2?language=en_CA (retrieved November 2020).

[76] Amazon White Paper: Overview of Amazon Web Services https://d0.awsstatic.com/whitepapers/aws-overview.pdf (August 2020).

[77] Amazon: Quantum computing is now available on AWS through Amazon Braket https://aws.amazon.com/about-aws/whats-new/2020/08/quantum-computing-available-aws-through-amazon-braket/ (retrieved November 2020).

[78] Google Cloud Platform Services Summary https://cloud.google.com/terms/services (retrieved November 2020).

[79] Google AI Platform https://cloud.google.com/ai-platform/ (retrieved November 2020).

[80] Microsoft Azure High-Performance Computing https://azure.microsoft.com/en-us/solutions/high-performance-computing/ (retrieved November 2020).

[81] InsideHPC: Azure Quantum Now in Public Preview https://insidehpc.com/2021/02/azure-quantum-now-in-public-preview/ (retrieved February 2021).

a short-term hourly lease of dedicated Cray hardware (with access to standard Azure based storage infrastructure) and not the standard cloud IaaS virtual instance offering.[82]

**OVHCloud**

Even though the big three cloud providers capture nearly 90% of ARC cloud deployments, multiple niche or more targeted cloud providers are being leveraged by the ARC community. For example OVH is a privately-owned company established in 1999 and is currently Europe's largest cloud hosting provider.[83] It presents itself as a transparent and secure alternative to the big three, working on for example sensitive storage and cloud solutions.[84] It reported ca 500M Euro sales in 2018, with aggressive multi-billion Euro investment in 2021-26 to catch up to the rivals.[85] The company announced in January 2021 a major Storage-as-a-Service initiative, in collaboration with IBM and Atempo, providing secure, sovereign and resilient storage for European enterprises and public institutions.[86]

**IBM Canada Watson & Bluegene**

Since 2012 IBM Canada has collaborated with and provided substantial computing resources to Canadian researchers, particularly in collaboration with SOSCIP and Queen University's CAC. Earlier in the decade researchers in Ontario had for example access to a IBM Watson analytics platform (at CAC), IBM Bluegene /Q supercomputer (at SciNet on behalf of SOSCIP), FPGA platform, GPU-accelerated platform, and IBM Cloud services.[87] The three former services seem to have been discontinued while the latter two services can be accessed via SOSCIP's fee-for-service offering.[88, 89]

International Organizations available for Canadians

Canadian researchers have access to a variety of international research computing resources, usually through collaborations with foreign researchers. The foreign PI submits the primary application for access and the Canadian researcher will get access as a foreign collaborator of that research team. In some cases, there could be limitations regarding accessing compute systems, and for example foreigner access to Oak Ridge National Lab's Summit and Frontier

---

[82] Microsoft Azure: Cray in Azure https://azure.microsoft.com/en-us/solutions/high-performance-computing/cray/ (retrieved November 2020).

[83] IT World Canada: Canadian customers' heads are still in the clouds, and so is VMware's https://www.itworldcanada.com/article/canadian-customers-heads-are-still-in-the-clouds-and-so-is-vmwares/421294 (retrieved January 2021).

[84] Reuters: France's OVH partners with Google for European cloud computing push https://www.reuters.com/article/ctech-us-ovh-google-cloud-idCAKBN27Q0OP-OCATC (November 10, 2020, retrieved January 2021.)

[85] Reuters: France's OVH to triple spending to take on Google, Amazon in cloud computing https://www.reuters.com/article/us-ovh-strategy-idUSKCN1MS17L (October 18, 2018, retrieved January 2021).

[86] InsideHPC: OVHcloud Teams with IBM and Atempo for Cloud Storage https://insidehpc.com/2021/01/ovhcloud-teams-with-ibm-and-atempo-for-cloud-storage/ (retrieved January 2021).

[87] IBM Canada: Why investing in Canadian R&D matters https://www.ibm.com/ibm/ca/en/ibmcanada100/investing-in-canadian-rd.html (retrieved November 2020).

[88] SOSCIP https://www.soscip.org/ (retrieved November 2020).

[89] CAC: CAC Services https://cac.queensu.ca/services/ (retrieved November 2020).

supercomputer is vetted via Oak Ridge National Lab (ORNL) Personnel Access System (PAS) mechanism[90] and might not be available for foreigners at all, or only under limited circumstances. If access to international systems is limited to collaborators of foreign PIs, this might disenfranchise some Canadian researchers who might not have established such international collaborations. In the following we discuss a few representative examples from the US, EU, and Australia.

In the US, at the national level, the Department of Energy (DOE) and the National Science Foundation (NSF) are the two major funders of cyberinfrastructure via a myriad of different mechanisms. Key DOE national programs hosting supercomputing systems[91] are National Nuclear Security Administration (NNSA, Trinity supercomputer hosted at Los Alamos National Labs, and Sierra hosted at Lawrence Livermore National Labs)[92], Oak Ridge Leadership Computing Facility (OLCF, hosting e.g. Summit and Titan)[93], Argonne Leadership Computing Facility (hosting e.g. forthcoming Aurora and Theta)[94], and National Energy Research Scientific Computing Center (NERSC, hosting e.g. the forthcoming Perlmutter, and current Cori supercomputers)[95]. Particularly the NNSA hosted systems are not available for general science use (even by US citizens) since they run highly sensitive nuclear stockpile simulations. The primary public access and resource allocation to ORNL and ALCF systems is via DOE Office of Science's Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program.[96] These labs also have Director's Discretionary programs and other mechanisms for providing smaller and shorter-term allocations. For example, the Advanced Scientific Computing Research (ASCR) Leadership Computing Challenge (ALCC) allocation program provides large resource allocation at OLCF, ALCF, and NERSC for high-risk, high-payoff simulations.[97] NERSC's primary mechanism for resource allocation is its own Energy Research Computing Allocations Process (ERCAP).[98]

NSF's primary mechanism for providing ARC resources is its Office of Advanced Cyberinfrastructure (OAC). NSF provides e.g., major funding to the National Center for Supercomputing Applications (NCSA).[99] NCSA, located at University of Illinois, Urbana-Champaign, hosts the Blue Waters supercomputer and leads the Extreme Science and

---

[90] Oak Ridge National Laboratory: Applying for a user account https://docs.olcf.ornl.gov/accounts/accounts_and_projects.html#applying-for-a-user-account (retrieved November 2020).

[91] US Department of Energy: Supercomputing and Exascale https://www.energy.gov/science-innovation/science-technology/computing (retrieved November 2020).

[92] National Nuclear Security Administration: Maintaining the Stockpile https://www.energy.gov/nnsa/missions/maintaining-stockpile (retrieved November 2020).

[93] Oak Ridge Leadership Computing Facility: Compute Systems https://www.olcf.ornl.gov/olcf-resources/compute-systems/ (retrieved November 2020).

[94] Argonne Leadership Computing Facility: ALCF Resources https://www.alcf.anl.gov/alcf-resources (retrieved November 2020).

[95] NERSC: Systems https://www.nersc.gov/systems/ (retrieved November 2020).

[96] INCITE Leadership Computing: INCITE Program https://www.doeleadershipcomputing.org/about/ (retrieved November 2020).

[97] Argonne Leadership Computing Facility: ALCC Allocation Program https://www.alcf.anl.gov/science/alcc-allocation-program (retrieved November 2020).

[98] NERSC: Allocations of computer time and storage https://www.nersc.gov/users/accounts/allocations/ (retrieved November 2020).

[99] NCSA: About NCSA http://www.ncsa.illinois.edu/about (retrieved December 2020).

Engineering Discovery Environment (XSEDE) project. XSEDE provides US researchers centralized access to multiple supercomputing resources, including those located at other NSF funded supercomputing centers, e.g. Indiana University (IU/TACC, Jetstream), Pittsburgh Supercomputing Center (PSC, Bridges, Bridges-2, Anton 2), San Diego Supercomputer Center (SDSC, Comet, Expanse), and Texas Advanced Computing Center (TACC, Ranch, Stampede2).[100] Access to various NSF funded resources is not straightforward, for example TACC's Frontera supercomputer is not available via XSEDE allocation mechanism, rather than TACC has a dedicated allocation application framework in place for that purpose[101], and the National Center for Atmospheric Research (NCAR)'s Cheyenne supercomputer is also not part of the XSEDE program.[102]

Besides the above discussed examples of DOE and NSF funded major supercomputing 'leadership class' operations ARC infrastructure in the US is characterized by a myriad of smaller scale systems, located, and managed by individual universities, states, or coalitions, many of which are potentially available for Canadian researchers via international collaborations. For example, Harvard University pooled its resources in 2011 with Massachusetts Institute of Technology and other local universities, state and private sector establishing the $168M USD Massachusetts Green High Performance Computing Center facility.[103] The center provides facilities and infrastructure for member universities for hosting their ARC infrastructure. For example, Harvard University operates main components of its 100,000 CPU core Cannon hybrid cluster in the facility.[104]

In Europe the historically main organization coordinating ARC resources is Partnership for Advanced Computing (PRACE) that was established in 2010, with a total funding of 125M EUR until 2019.[105] The European ARC ecosystem under PRACE is split into three tiers where Tier-0 covers the supercomputers (petaflop systems), Tier-1 is for the national level systems, and Tier-2 contains the regional and university level systems.[106] The hosting member countries (Germany, France, Italy, Spain, and Switzerland) have committed to fund and deliver research infrastructure services to the 26 country PRACE member coalition. Currently PRACE's ARC infrastructure consists of seven Tier-0 level systems, some of which are hosted in multiple locations or have multiple functionally different segments. The most recent is HAWK, located at High-Performance Computing Centre Stuttgart (HLRS) that came online in 2020 with 26 peak Pflops performance.[107] There are nineteen Tier-1 systems, totaling 16 PFlops peak compute power. The PRACE systems are in principle available for Canadian researchers with possible restrictions imposed by the host

---

[100] XSEDE: XSEDE Resource Information https://portal.xsede.org/allocations/resource-info (retrieved December 2020).

[101] TACC: FRONTERA ALLOCATION SUBMISSION GUIDELINES https://frontera-portal.tacc.utexas.edu/allocations/policy/ (retrieved December 2020).

[102] NCAR:Allocations https://www2.cisl.ucar.edu/user-support/allocations (retrieved December 2020).

[103] The Harvard Crimson: Harvard Helps Build $168M Supercomputing Facility https://www.thecrimson.com/article/2011/10/31/supercomputers-research-facility-holyoke/ (retrieved December 2020).

[104] Harvard University Faculty of Arts and Sciences Research Computing: Cluster Architecture https://www.rc.fas.harvard.edu/about/cluster-architecture/ (retrieved December 2020).

[105] PRACE: Introduction https://prace-ri.eu/about/introduction/ (retrieved December 2020).

[106] ARCHER Training: HPC in Europe https://www.archer.ac.uk/training/course-material/2017/11/intro-epcc/slides/L12-PRACE.pdf (retrieved December 2020).

[107] PRACE: HPC Systems https://prace-ri.eu/hpc-access/hpc-systems/ (retrieved December 2020).

institutions. Additionally, the terms of reference for the latest resource call for example state that collaboration with a PI from a European country that contributes to PRACE will improve the chances of approval for resources.[108] PRACE is also encouraging international collaboration via calls for collaboration with XSEDE and RIST.[109] Since the European ARC ecosystem is quite complex, PRACE has recently taken the initiative to launch HPC-in-Europe portal in order to coordinate European ARC services via a bottom-up approach, providing a one-stop shop for ARC users. This portal, located at hpc-portal.eu is still under development.[110]

In 2018 European Union formed the European High-Performance Computing Joint Undertaking (EuroHPC JU) with the goal of coordinating efforts and funding exascale computers, and funding of EUR 1.1B in FY 2019-20. The three main hosting sites for pre-exascale systems will be Barcelona Supercomputing Centre (Spain), CSC (Finland), and CINECA (Italy). It currently has five supercomputers under construction, largest of which (LUMI) would have a peak 552 PFlops performance, planned to come online in 2021.[111] According to Oriol Pineda of PRACE roughly 5% of PRACE resources were used by foreigners, while the policies for eligibility of foreign researchers to use EuroHPC resources are still under development.[107]

There are two main Tier-1 supercomputing centers in Australia, both funded by Australian government's Department of Education's National Collaborative Research Infrastructure Strategy (NCRIS).[112] The National Computational Infrastructure (NCI) organization hosts Australia's fastest supercomputer, 9 peak PFlops Gadi,[113] while Pawsey Supercomputing Centre (PSC) operates the Magnus petascale supercomputer among others.[114] In October 2020 PSC announced that their next supercomputer will be a peak 50 PFlops system from HP/Cray.[115] The current systems at Pawsey are in general available for foreign researchers if they collaborate with an eligible PI.[116]

## 3.4 How is ARC delivered and funded in other jurisdictions?

The 2017 LCDRI ARC Current State Paper discussed the global ARC service delivery in chapter 4.3, and in Appendix C. Since then, the general global funding landscape has remained relatively unchanged, with supranational (i.e. EU with PRACE and more recently EuroHPC), complex national (i.e. the US with e.g. DOE and NSF at federal level mixing with myriad regional and local

---

[108] PRACE: PRACE Project Access Terms of Reference –22nd Call for Proposals https://prace-ri.eu/wp-content/uploads/Terms_of_Reference_Call22.pdf (retrieved December 2020).

[109] PRACE: Collaborative Calls https://prace-ri.eu/hpc-access/collaborative-calls/ (retrieved December 2020).

[110] SC20: European HPC Ecosystem Summit presentation by Oriol Pineda  https://cdmcd.co/Qqz7Eq (see Q/A at 37:55 mark, November 2020).

[111] EuroHPC: Discover EuroHPC https://eurohpc-ju.europa.eu/discover-eurohpc (retrieved December 2020).

[112] NCRIS Network: Infrastructure Projects Funded by NCRIS https://www.ncris-network.org.au/capabilities (retrieved December 2020).

[113] NCI Australia: HPC Systems https://nci.org.au/our-systems/hpc-systems (retrieved December 2020).

[114] PSC: Magnus https://pawsey.org.au/systems/magnus/ (retrieved December 2020).

[115] PSC: Powering the next generation of Australian research with HPE https://pawsey.org.au/powering-the-next-generation-of-australian-research-with-hpe/ (retrieved December 2020).

[116] Pawsey: Application Process https://support.pawsey.org.au/documentation/display/US/Application+Process (retrieved December 2020).

level funding sources), centralized national (e.g. Japan), and collaborative national (e.g. Australia) delivery and funding models still applying to characterize ARC funding on global scale. As the Alliance will propose to ISED a Strategic Plan and a New DRI Funding Model in the Fall 2021, a thorough analysis of the global DRI delivery and funding landscape is required to provide potential ideas for future Canadian DRI delivery. Due to the major scope and resources required for such a survey, a detailed global DRI, including not only ARC but also RS and RDM, delivery and funding model review and analysis will be conducted as a separate Environmental Scan project in early to mid 2021.

As discussed in the previous section, a major development in the ARC international landscape in recent years has been the creation of EuroHPC to deliver world-leading supercomputing resources within the European Union, including major funding increases beyond and independent of PRACE's budget envelope. On the technical advancements side the previous chapter also touches on various new foreign supercomputing systems that have become available or will soon be available for Canadian researchers through potential individual level international collaborations.

As discussed in Chapter 4.1, Canadian ARC is still largely consumed by traditional ARC disciplines of hard sciences like physics, engineering, and computer science. Diversity in usage of ARC has been increasing at a very fast rate in Canada, including demand for not just ARC, but for a more comprehensive and encompassing DRI service portfolio. Globally some of these underserved research disciplines already have substantial DRI support. For example, in France the Huma-Num infrastructure provides Humanities and Social Sciences researchers not only ARC computing services, but a full portfolio of DRI services all through the research lifecycle, while Prodego focuses on social science data This infrastructure is considered a "Very Large Research Infrastructure (TGIR)" at the government funding level, and provides platforms and tools for processing, conservation, dissemination and long-term preservation of digital research data.[117]

## 3.5 Future trends in ARC and AI architectures, markets and needs

In the following we will discuss general, mostly global ARC and AI market trends and needs. In the Winter 2020 and Spring 2021 the Alliance will conduct three more detailed and focused assessments of future trends 1) a Needs Assessment survey and analysis of the needs and trends in the Canadian DRI community, 2) a global Environmental Scan, and 3) a Technological Advancements Review.

Despite Covid-19 and related dire economic circumstances in 2020 (and likely beyond) ARC and AI markets are predicted to continue to grow substantially in the next five years, at the estimated compounded annual growth rate (CAGR) of 7.1% to a $55B USD market in 2024 according to Intersect360 Research. This growth will have temporal undulating variations so that the spending will drop in 2020, but the pent-up demand will increase the spending substantially in 2021. Such undulations are predicted to result in an overall unchanged CAGR (compared to pre-covid-19 estimates) over the long term. [118] Hyperion Research forecasted ca. 8.7% CAGR pre-covid and

---

[117] Huma-Num: About us https://www.huma-num.fr/about-us (retrieved December 2020).

[118] Addison Snell of Intersect360 Research at HPC-AI Advisory Council 2020 Australia Conference: Supercomputing to the Rescue: HPC/AI Market Update http://www.hpcadvisorycouncil.com/events/2020/australia-conference/pdf/HPCAIMarketUpdate_020920_ASnell.pdf (September 2020).

is more cautious than Intersect Research regarding long-term spending post-covid, citing the likely need for budget tightening in the government sector in 2023 and beyond.[119] In their SC20 conference update in November 2020, Hyperion Research estimated the effect of Covid-19 to be short but sharp, with ARC market revenue dropping 11.5% in the first half of 2020.[120]

On the client/customer side government spending is forecast to increase over the next five years both in absolute and relative terms due to weakening spending on commercial sector (energy, retail, and large product manufacturing).[115] Biosciences, defense and government laboratory spending are predicted to grow over the next few years with biosciences getting an extra boost in 2020-22 thanks Covid-19 related research.[116]

## Cloud

The ARC cloud landscape is very dynamic with constant innovation, growing quickly, and would justify its own research paper. Cloud computing can be defined as "...the delivery of computing services—including servers, storage, databases, networking, software, analytics, and intelligence—over the Internet ("the cloud") ...".[121]

In the following we will touch upon a few interesting new service offerings and review the projected growth rates. On the service offering front the various platform-as-a-service (PaaS) offerings, located between the more traditional infrastructure-as-a-service (IaaS, e.g., Amazon EC2 instances), and software-as-a-service (SaaS, e.g. Microsoft Office365 cloud collaboration platform), are particularly interesting. A major recent Canadian example of PaaS services is the above-mentioned Syzygy, a collaboration between Pacific Institute for the Mathematical Sciences (PIMS), Compute Canada, and Cybera, initiated in 2017, providing Canadian researchers Jupyter Notebook based computing resources free of charge.[122]

Another interesting project is PanGeo.io community platform, developing open-source python tools for Big Data geosciences including scalability for petabytes of data in on-premise and cloud ARC environments.[123] The platform has global funding and support from government agencies (NSF, NASA, UK Met Office), academia (University of Washington), and private sector (Anaconda). Pangeo developers are for example helping NASA to process massive datasets in its Earthdata Cloud by leveraging metadata to load and process only relevant parts of massive datasets.[124] The platform is not limited to geosciences, for example McGill University's Neuro institute is currently doing a PanNeuro pilot leveraging PanGeo. The Pangeo Cloud part of the

---

[119] InsideHPC: Hyperion Research Forecasts Widespread Covid-19 Disruption to HPC Market https://insidehpc.com/2020/04/hyperion-research-forecasts-widespread-covid-19-disruption-to-hpc-market/ (retrieved September 2020).

[120] InsideHPC: At SC20: Hyperion Sees COVID HPC Impact Sharp but Short; HPE in Server Lead; Aurora 12+/- Months Late; Cloud HPC Heating Up https://insidehpc.com/2020/11/at-sc20-hyperion-sees-covid-hpc-impact-sharp-but-short-hpe-in-server-lead-aurora-12-months-late-cloud-hpc-heating-up/ (retrieved November 2020).

[121] Microsoft Azure: What is cloud computing? https://azure.microsoft.com/en-us/overview/what-is-cloud-computing/ (retreived May 2021).

[122] Syzygy.ca https://syzygy.ca/# (retrieved December 2020).

[123] Pangeo: About Pangeo http://pangeo.io/about.html (retrieved December 2020).

[124] NASA Earthdata: The Pangeo Project: Developing Community Tools For a New Era of Data Analysis https://earthdata.nasa.gov/learn/articles/pangeo-project (retrieved December 2020).

platform is a AWS or Google Cloud based data-science environment and is currently in the experimental development phase with no direct long-term funding.[125]

Even though cloud services are relatively easy to launch for anyone with a credit card, understanding the true cost of these services is very complicated and depends on the exact nature of computing done. To address these issues, with support from NSF, University of California San Diego has collaborated with multiple California-based institutions to form CloudBank, "a cloud access entity that will help the computer science community access and use public clouds for research and education by delivering a set of managed services designed to simplify access to public clouds".[126] At a practical level the service provides NSF funded researchers with need for cloud resources, a user portal and related training to access cloud resources from multiple vendors at favorable rates thanks to avoiding indirect costs, and pooled purchasing power.

Looking at future growth and changes in products and services categories, the cloud sector is estimated to grow at a much higher rate than the rest of the market, at over 20% CAGR according to Intersect360 Research.[115]. Hyperion Research foresees a CAGR of nearly 25% until 2023, driven by movement of workgroup class (systems priced less than $100k USD) workloads to the cloud. [127] The total size of the ARC in the cloud market is estimated to have been ca. $2.2B to $2.8B USD in 2019.[128]

According to InsideHPC in 2019 roughly 30% of the ARC community is already using cloud for production workloads, while over 90% overall are interested in using cloud. Of the organizations using cloud, over 60% are seeing positive results or consider the cloud as a 'real winner'. [116] Hyperion Research reports that the proportion of ARC workloads in the cloud grew from 10% to 20% just between 2018 and 2019. This recent substantial increase in ARC cloud adoption is seen to be driven by improved software application offerings, easier access models for cloud services, and increased compute capabilities. Earlier in the decade the key drivers for ARC cloud adoption were the ability to test and develop, 'experiment' on leading edge technologies in the cloud, and to run embarrassingly parallel workloads.[124]

### On-premise solutions

On the traditional on-premise server side the purchases are expected to focus increasingly on the high-end and supercomputing systems side of the market since this segment is not yet well served by commercial cloud providers due to highly advanced and often customized technology these systems use. There is also a wide range push for funding exascale computing efforts, a need that can not (yet) be served by cloud systems.[115,116] At the lower performance end of ARC, i.e. the less than $100k USD workstation or server side, the purchases dropped 25% year-on-year from Q2

[125] Pangeo: Pangeo Cloud http://pangeo.io/cloud.html (retrieved December 2020).

[126] CloudBank: About CloudBank https://www.cloudbank.org/about (retrieved December 2020).

[127] Hyperion Research white paper: Bringing HPC Expertise to Cloud Computing https://www.dellemc.com/resources/en-us/asset/analyst-reports/products/ready-solutions/hyperion-dell-cloud-hpc-ai.pdf (April 2020).

[128] Addison Snell of Intersect360 Research: Pre-SC20 Market Update http://www.intersect360.com/LiteratureRetrieve.aspx?ID=158848 (retrieved November 2020).

2019 to Q2 2020, indicating pressure on this segment (beyond short-term Covid-19 related issues) as these workloads (continue to) move to the more flexible and variable cloud offerings.[129]

## Changes in vendor profiles and ARC workflows

Considering the vendor and service provider side, the number of major vendors is predicted to decline, while there will be growth in smaller niche / start-up vendor solutions.[130] There will also be increased competition between chip makers (e.g., x86 v. ARM v. growing custom chip market) and increased demand for GPU accelerator resources.[127] Custom chip designs for e.g., AI solutions can potentially be 10 to 100 times faster than standard x86 or GPU based solutions. The venture capital funding in the US on AI acceleration market is estimated to be over $4B USD.[131] Hyperion Research estimates the quantum computing market to be of the order of $320M USD in 2020, with substantial part going to cloud based quantum computing. They estimate this market to grow at ca. 27% CAGR.[132]

A major future technical and operational trend will be improved workflow and efficiencies in the whole DRI ecosystem. There will be a need for integrated workflows with networked pre- and post-edge and cloud computing components.[127] Embedding edge computing at various workflow stages will allow reduction in massive data volumes and efficiencies in data transfer and storage needs. Such technologies include e.g., on-sensor / field-deployable-processing, near sensor and real time processing, smart ARC interconnects, and specialized accelerators as discussed above.[128] "AI for AI" will be a key component in bringing individual components together to a 'seamlessly integrated facilities' that execute automated experiments and analysis. The key benefit being that integrated facilities will improve productivity and reproducibility in science.[128]

On the research domain side, the central National Computational Infrastructure organization in Australia lists increased demand in ARC services, growing scale of ARC jobs, increasing number of disciplines, and the above-mentioned need for integrated workflows as key trends.[127]

## Storage and data management

Storage and data management will be more and more important in established (traditional ARC) and emerging (e.g., big data, AI, ML, deep learning etc.) ARC fields. The annual growth of these markets will be substantial. In the life sciences alone, the data analysis, storage & management market is predicted to grow at 17.1% CAGR between 2018 and 2024, to $41.1B USD globally.[133]

---

[129] InsideHPC: Hyperion: Covid-19 Driving Down HPC Server Revenues, But Impact May Be Moderating https://insidehpc.com/2020/09/hyperion-research-covid-19-driving-decline-in-hpc-server-revenues-but-may-be-moderating/ (retrieved September 2020).

[130] Allan Williams of NCI Australia at HPC-AI Advisory Council 2020 Australia Conference: HPC Impact: Future of Scientific Computing http://www.hpcadvisorycouncil.com/events/2020/australia-conference/pdf/FutureofSciComp_020920_AWilliams.pdf (September 2020).

[131] Kathy Yelick of Lawrence Berkley National Laboratories at HPC-AI Advisory Council 2020 Australia Conference: AI for Science https://www.youtube.com/watch?v=sLjI9p3u7Mo&list=PLafs-cr09EuW71NepWOlQ98Ht8K40VnJx&index=7 (September 2020).

[132] Hyperion Research at High Performance Computing at AWS conference: HPC in the cloud with Hyperion Research https://hpcaws.splashthat.com/ (November 2020).

[133] MarketsandMarkets: HPC, Data Analysis, Storage & Management Market in Life Sciences By Products & Services (Data Analysis, Cloud Computing), Applications (NGS, Microscopy, Chromatography), End User (Pharmaceutical & Biotechnology, Hospitals) - Global Forecast to 2024

Hyperion Research forecasts that specifically the global ARC storage market would grow at ca. 7% CAGR from $5.5B USD in 2018 to $7.0B USD in 2023. [134] Iterative simulation methods and growing Big Data and AI workloads generating massive amounts of files and volumes of data are driving the need for improved data management for more efficient use of storage capacity and better performance. Most ARC storage users estimate their annual capacity growth rate to be less than 50%, with emphasis in the 25-50% annual growth rate. While the following paragraphs describe technical challenges, notably the key inhibitor for storage operations is recruiting and hiring HQP personnel.[131]

AI, ML and deep learning are emphasizing existing ARC storage characteristics, while also introducing new storage workflows that are different from traditional ARC. Some of the key requirements to consider are 1) need for high-performance networking, 2) need for shared storage, 3) need for multi-tiered storage, 4) need for parallel access, 5) need to support multiple protocols, and 6) advanced metadata handling.[135]

High-performance networking is of course the bread and butter of traditional ARC systems, but in the case of storage such networks can be designed for e.g., low-latency memory based networked storage solutions like NVMe Over Fabrics (NVMe-oF). Shared storage where all loads have access to all storage is a common feature in ARC storage, but for AI this is very important in order to distribute loads more evenly, and for keeping the expensive GPU units well-utilized. Multi-tiered storage at the age of AI/ML must take into account in particular the (conflicting) characteristics that the training datasets are often huge, need to be accessed with very fast I/O when utilized, and need to be kept for long times since (re)collecting the data can be costly (or perhaps even impossible). On one hand all-Flash solutions can be prohibitively expensive, so a multi-tiered hybrid solution that includes persistent memory, Flash, and a pool of disks is required.[136] On the other hand such tiered storage systems can be confusing and inconvenient for the end-users. To address this, NERSC's forthcoming Perlmutter supercomputer system is aiming to have only two storage tiers presented to the end-users, while tiering is handled automatically in the backend by software-based solutions.[137] Parallel access is another common ARC storage feature, but this need is again emphasized by the sheer scale of compute processes in typical AI/ML workloads that need simultaneous access to the (same) data. Support for multiple protocols, for example native parallel FS, NFS, SMB, S3 etc., is required since AI/ML data can be collected from devices with their own protocols, e.g., Internet-of-things network connected devices, while on the other hand the data is consumed within the ARC massively-parallel storage

https://www.marketsandmarkets.com/Market-Reports/hpc-data-analysis-storage-management-market-47829739.html (April 2019, public summary retrieved September 2020).

[134] Hyperion Research White Paper (sponsored by Panasas): New Study Details Importance of TCO for HPC Storage Buyers https://www.panasas.com/wp-content/uploads/2020/04/Hyperion_Importance-of-TCO-for-HPC-Storage-Buyers_Q1-20_FINAL_2020-04-22.pdf (April 2020).

[135] Storage Switzerland LLC White Paper (sponsored by Panasas): Is Your Storage Infrastructure Ready for the Coming AI Wave? https://insidehpc.com/white-paper/is-your-storage-infrastructure-ready-for-the-coming-ai-wave/ (January 2020).

[136] The Next Platfom: Divide Deepens Between HPC and Enterprise Storage https://www.nextplatform.com/2020/10/20/divide-deepens-between-hpc-and-enterprise-storage/?mc_cid=61a88daab6&mc_eid=79a1266800 (retrieved October 2020).

[137] HPCWire: SC20 Panel – OK, You Hate Storage Tiering. What's Next Then? https://www.hpcwire.com/2020/11/25/sc20-panel-ok-you-hate-storage-tiering-whats-next-then/ (retrieved November 2020).

system. Advanced metadata handling is not unique to AI/ML, but the sheer volume of billions of files with metadata attached puts additional I/O pressure on the storage system.[132]

## Towards and past exascale

Globally the leading edge of the ARC community is working towards exascale performance, even though this is not a key or realistic objective in the mid-term in Canada. Currently very large scale, massively-parallel jobs in the CCF systems are not that frequent. 96% of workloads run on CCF systems are less than 2048 cores, and 72% are less than 1024 cores. Keeping in mind that most of the CCF clusters are in theory capable of running 30k+ core jobs. Whether this situation is simply due to lack of access to resources at larger scale or due to lack of massively parallel scalable codes that can leverage such scale needs further investigation. For Canadian researchers to gain access to and leverage massively parallel exascale computing it might be beneficial for Canada to establish a partnership with US-based providers to be able to run those exceptionally large cases.

In the words of HPE/Cray's Chief Technologist Nic Dube "Nothing is easy in Exascale".[138] The key issues will be power, software and system resiliency requirements. The power requirements for an exascale system are estimated to range at 30-40 MW, and upwards. This will require major infrastructure updates at the datacenter level. On the software front the sheer scale of parallelism will require focused code development efforts in order to be able to run real life applications (leveraging e.g., CUDA libraries that abstract away the parallelism to a degree) and not just 'trophy' type Linpack codes for Top500. On the system resiliency side, the system software work done in the US supercomputing centers has successfully provided "an appearance of resiliency" to the end-users even though the underlying components are not resilient (or don't even need to be).[135]

Reaching for the exascale, the two prevailing architectural approaches seem to be either going CPU cores only, or CPU cores + GPU accelerator combination.[139] The current Top500 #1 machine, Japanese ARM-based 'Fugaku' is an example of the former, although Fugaku's custom designed processor is both many-core ARM CPU and accelerated GPU-like processor, while still maintaining wide support for applications via "application first" and being "the leader in all HPC benchmarks" design principles.[140]

The recent US leadership class machines (e.g., Summit and Frontera) are examples of the latter (CPU+GPU) approach. Also some of the future European pre-exascale machines fall into this latter category, for example EuroHPC JU's over 200 peak PF 'Leonardo' at CINECA will be based on leveraging Intel CPUs and Nvidia A100 GPUs in one to four ratio.[141] Interestingly the new Leonardo machine is from Atos/Bull, indicating its rapidly growing presence in the ARC

---

[138] InsideHPC: Getting to Exascale: Nothing Is Easy https://insidehpc.com/2020/10/getting-to-exascale-nothing-is-easy/ (retrieved October 2020).

[139] Prof. Mark Parsons of EPCC at HPC-AI Advisory Council 2020 UK Conference: Exascaling AI http://www.hpcadvisorycouncil.com/events/2020/uk-conference/pdf/day-one/M_Parsons_ExascalingAI_131020.pdf (October 2020).

[140] SC20 conference presentation by Satoshi Matsuoka: Fugaku: 'Exascale' and Applications First' https://cdmcd.co/Kra6dr (November 2020).

[141] HPCWire: Nvidia and EuroHPC Team for Four Supercomputers, Including Massive 'Leonardo' System https://www.hpcwire.com/2020/10/15/nvidia-and-eurohpc-team-for-four-supercomputers-including-massive-leonardo-system/ (retrieved October 2020).

marketplace, taking market share from the usual big players (HPE/Cray, Dell, Lenovo, IBM).[142] The cores only route is characterized by longer system lifetimes, easier SW implementation for traditional simulation codes, poorer AI performance, and larger physical space and power requirements compared to the cores+GPU approach.[136] In order to balance the needs of different audiences EuroHPC is taking a hybrid approach with its forthcoming LUMI pre-exascale supercomputer. LUMI will have a CPU-nodes-only cluster with various memory allocations and then a very large cluster of CPU+GPU nodes, all connected via shared Ethernet based high-speed interconnect. The system is expected to deliver 552 PF peak capacity and is expected to be built in 2021.[143] AMD's Epyc processor family in combination with its Instinct GPUs is gaining a lot of traction primarily thanks to the compute density and cost-efficiency of the Epyc's compared to Intel Xeon's, for example Frontier and El Capitan systems in the US and the new supercomputer at Pawsey Supercomputing Center in Perth, Australia will leverage these chips.[144]

Looking at longer-term technology trends that are needed to make 'post-exascale' computing possible, the prevailing view in particular in the US is that the era of scale based improvements must transfer more to development based improvements: Slowing down of Moore's law in combination with the decrease in scale of improvements in transistor density, thread performance, clock frequency, power efficiency, and number of cores per socket will require new innovations beyond just these 'traditional' fronts.[145]

Novel architectural development efforts can be organized in to three (overlapping) categories, 1) building of special purpose machines, 2) designing chip systems with heterogeneous integration of myriad custom accelerators for specific 'microtasks', and 3) improving the workload adaptability of the ARC system via resource disaggregation.[142] The first category includes custom built machines similar to D. E. Shaw Research's Anton series for molecular dynamics[146], Google Tensor Processing Unit (TPU) ASIC for machine learning,[147] custom neuromorphic chips for spiking neural networking based AI, or potentially a custom supercomputer for density-functional-theory (DFT) based computations that currently comprise 25% of NERSC workload.[142]

The second category (heterogeneous integration of custom accelerators) involves aggregation of highly customized accelerator chiplets located immediately next to the main CPU chip (beyond

---

[142] The Next Platform: With Another Key Supercomputer Win, Atos Looks Stronger Than Ever https://www.nextplatform.com/2020/10/15/with-another-key-supercomputer-win-atos-looks-stronger-than-ever/?mc_cid=ee478a79f8&mc_eid=79a1266800 (retrieved October 2020).

[143] The Next Platform: The Resurrection Of Cray And AMD In A Trifurcating HPC Space https://www.nextplatform.com/2020/10/22/the-resurrection-of-cray-and-amd-in-a-trifurcating-hpc-space/?mc_cid=61a88daab6&mc_eid=79a1266800 (retrieved October 2020).

[144] The Next Platform: HPE And AMD Bag The Big Supercomputer Deal Down Under https://www.nextplatform.com/2020/10/19/hpe-and-amd-bag-the-big-supercomputer-deal-down-under/?mc_cid=ee478a79f8&mc_eid=79a1266800 (retrieved October 2020).

[145] John Shalf of Lawrence Berkeley National Laboratories at Oklahoma Supercomputing Symposium 2020: Pathfinding for Post-Exascale HPC http://www.oscer.ou.edu/Symposium2020/oksupercompsymp2020_talk_shalf_20200930.pdf (September 2020).

[146] Pittsburgh Supercomputing Center: Anton https://psc.edu/resources/computing/anton (retrieved October 2020).

[147] Google Cloud: Cloud TPU https://cloud.google.com/tpu/ (retrieved October 2020).

the usual integrated general purpose GPU accelerators).[148] Already existing commercial examples include the highly specialized Apple's Bionic chips for e.g. machine learning and motion analysis in smartphones,[149,150] or Amazon's AWS Graviton customized ARM chip line for cloud workloads.[151] The Networking and Information Technology Research and Development (NITRD) Program in The US is investigating such opportunities for ARC applications in their Project 38.[152] A related interesting technology is Data Processing Unit (DPU), a silicon-on-chip (SoC) unit that resides next to the traditional CPU and GPU units, and combines ARM processor cores w/ integrated networking and GPU capabilities for off-loading some of the compute and networking management tasks from the CPU to the DPU.[153]

The third category (resource disaggregation) aims to create reconfigurable compute nodes with very high-speed interconnect between key components. An internal interconnect at the speed of photonic bandwidth would allow flexible allocation and low-level (re)-configuration of memory, CPU, GPU, I/O, networking etc. resources. Allowing a single system to execute efficiently a variety of workloads (e.g., AI training, AI inference, data mining or graph analytics) with very different memory, compute, networking and I/O needs. Such a design could be based on e.g., photonic multi-chip-modules (MCMs) as proposed by the Photonic Integrated Networked Energy efficient datacenter (PINE) project at Columbia University.[154] Another benefit of hardware disaggregation is the flexibility for adapting to various workloads, e.g., between traditional ARC, high-throughput (HTP) and AI workloads. In the US some recent DoD supercomputing procurement decisions have gone to disaggregated solutions emphasizing flexibility over pure (and more traditional) FLOPs per dollar metric.[155]

---

[148] Prof. Simon McIntosh-Smith of University of Bristol at HPC-AI Advisory Council 2020 UK Conference: Exascale Research and Development Opportunities http://www.hpcadvisorycouncil.com/events/2020/uk-conference/pdf/day-one/S_McIntoshSmith_ExaRandDOpps_131020.pdf (October 2020).

[149] Apple: Apple unveils all-new iPad Air with A14 Bionic, Apple's most advanced chip https://www.apple.com/newsroom/2020/09/apple-unveils-all-new-ipad-air-with-a14-bionic-apples-most-advanced-chip/ (retrieved October 2020).

[150] Wired: An Exclusive Look Inside Apple's A13 Bionic Chip https://www.wired.com/story/apple-a13-bionic-chip-iphone/ (retrieved October 2020).

[151] ZDNet: AWS Graviton2: What it means for Arm in the data center, cloud, enterprise, AWS https://www.zdnet.com/article/aws-graviton2-what-it-means-for-arm-in-the-data-center-cloud-enterprise-aws/ (retrieved October 2020).

[152] NITRD Project 38 Technical Report: HPC Performance Improvements Through Innovative Architecture https://www.nitrd.gov/Presentations/files/HPC-Performance-Improvements-Project-38.pdf (October 2019).

[153] EnterpriseAI: Nvidia Expands Its DPU Family, Unveils New Datacenter on Chip Architecture https://www.enterpriseai.news/2020/10/05/nvidia-expands-its-dpu-family-unveils-new-datacenter-on-chip-architecture/ (retrieved October 2020).

[154] Columbia University in the City of New York: Photonic Integrated Networked Energy efficient datacenter (PINE) https://lightwave.ee.columbia.edu/research-projects/photonic-integrated-networked-energy-efficient-datacenter-pine (retrieved October 2020).

[155] TheNextPlatform: For HPC And AI, Composability Might Trump Cheap Flops https://www.nextplatform.com/2020/10/27/for-hpc-and-ai-composability-might-trump-cheap-flops/?mc_cid=95e0f6bf8a&mc_eid=79a1266800 (retrieved November 2020).

## 3.6 Covid-19

Covid-19 substantially increased interest in cloud computing, both in the general and ARC marketplaces. Customers that already had cloud (some, perhaps e.g., experimental) presence in the cloud increased their cloud usage due to increased general computing demand, directing the new computing demand to cloud at an accelerated adoption rate compared to their earlier timelines. An additional motivation for accelerated cloud adoption was the ability to manage increased volatility in computing needs. The scientific community has put substantial focus on Covid-19 related research, leveraging specialized ARC and AI cloud resources, and special grants from commercial providers. On the infrastructure side, many corporations and institutions that previously considered looking into ARC cloud(s) as a multi-year project accelerated their cloud testing and adoption due to changing workplace demands, and difficulties and dependencies related to managing on-prem ARC datacenters.[156]

On the Compute Canada Federation front main hosting sites had minimal downtime due to Covid-19 and have been able to provide services all through the pandemic. CCF has also provided dedicated Covid-19 support for Canadian researchers, including "Providing access to cloud resources, high-performance clusters, storage or boosting job priority; Consulting in high-performance computing (HPC), data management, data analysis, machine learning, and visualization; and connecting Canadian scientists from other research areas and institutions to further collaborations".[157] Compute Canada has also contributed resources to the international folding@home distributed computing project by optimizing code to GPUS, and by running SARS-CoV-2 protein structure simulations on Arbutus cloud and Cedar supercomputer.[158] By the end of September 2020 CCF had granted 27 allocations for Covid-19 projects that had requested a total of ca. 4600 CPU core years of priority compute cycle access. At that point the research projects had used 1300 CPU years of these resources. This usage is roughly 0.7% of the total available CPU year resource available in 2020 at CCF.[159] Notably this number does not consider all Covid-19 research done on CCF systems but considers only if the research group had requested additional resources or prioritization beyond their standard allocation. As an indication of the scope of such research activity, CCF community and researchers were mentioned in at least 70 Covid-19 related news stories and social media posts between March and October 2020.[160] Two special use cases at CCF that directly impacted decision making and were not exploratory in nature were: All modeling done for INSPQ (Institut national de santé publique du Québec), which drive the decisions about public health restrictions, have been done by a Université Laval group,

---

[156] Altair HPC Virtual Summit 2020 - Cloud Roundtable: Is Cloud Officially Inevitable? Experts from Azure, Oracle, Advania and Google get candid about 2020's biggest cloud computing trends and challenges https://player.vimeo.com/video/455986574 (September 2020).

[157] Compute Canada: Support for COVID-19 research projects https://www.computecanada.ca/featured/support-for-covid-19-research-projects/ (retrieved September 2020).

[158] Compute Canada: Harnessing the Power of Scientific Cloud Computing to Fight COVID-19 https://www.computecanada.ca/featured/harnessing-the-power-of-scientific-cloud-computing-to-fight-covid-19/ (retrieved September 2020).

[159] Source: Compute Canada Database, provided by Maxime Boissonneault (October 2020).

[160] Per Compute Canada communications team analysis provided by Maxime Boissonneault (November 2020).

with support from Charles Coulombe, a CQ Analyst, and used Graham,[161] and the sequencing to identify variants of COVID-19 are currently being done on Béluga, by a group from McGill, also with the help of a CQ analyst from McGill.[162]

At global scale, individual institutions and corporations have also increased their investment in ARC due to Covid-19. For example, AMD donated 5 Petaflops worth of computing power for Covid-19 research in academia in September 2020,[163] including the substantial hardware donation for SciNet's new SciNet4Health initiative.[164] In November 2020, the US DOE procured a dedicated supercomputer funded by the Coronavirus Aid, Relief and Economic Security (CARES) Act. The 'Mammoth' large-memory system will be located at Lawrence Livermore National Labs (LLNL) and consists of 64 twin AMD Epyc based server nodes with 2TB of RAM and 4TB of non-volatile memory each, connected via Omnipath high speed interconnect. The machine is designed for COVID-19 research, e.g., genomics analysis, non-traditional ARC simulations and graph analytics.[165] The world's fastest supercomputer, Japan's Fugaku at RIKEN, was brought online nearly a year in advance in order to tackle Covid-19 research. The supercomputer has been used at full scale to study e.g., droplet transmission, airflow, filtering performance of face masks, and effects of humidity on the viability of the Coronavirus, driving Japanese policy making directly.[166] In the US the capacity of the world's ninth fastest supercomputer, Frontera, will increase by roughly 5% in January 2021 via special grant from National Science Foundation (NSF) and donation from Dell.[167]

## 3.7 Academic Return-on-Investment for ARC

ARC and DRI have a myriad of benefits to the individual researchers and science, societies, and industries. ARC can help not only solve problems and scientific questions that would otherwise be difficult to solve, but also contribute to solutions that would not be possible via regular (e.g., analytic, experimental, or workstation level computing) means. The problems solved with ARC can range from nanoscale like drug discovery to macroscopic like severe weather simulations or

---

[161] Quebec INSPQ: Épidémiologie et modélisation de l'évolution de la COVID-19 9 avril 2021 - Mise à jour des projections du 18 mars https://www.inspq.qc.ca/covid-19/donnees/projections/9-avril-2021 (retrieved April 2021).

[162] Calcul Quebec: COVID-19: Opening Up the Data https://www.calculquebec.ca/en/recherche/covid-19-opening-up-the-data/ (retrieved April 2021).

[163] AMD: AMD COVID-19 HPC Fund Adds 18 Institutions and Five Petaflops of Supercomputer Processing Power to Assist Researchers Fighting COVID-19 Pandemic https://www.amd.com/en/press-releases/2020-09-14-amd-covid-19-hpc-fund-adds-18-institutions-and-five-petaflops (retrieved September 2020).

[164] University of Toronto: U of T and AMD launch supercomputing program dedicated to big-data health research https://www.utoronto.ca/news/u-t-and-amd-launch-supercomputing-program-dedicated-big-data-health-research (retrieved September 2020).

[165] HPCWire: Lawrence Livermore Announces Mammoth Cluster to Fight COVID-19 https://www.hpcwire.com/2020/11/04/lawrence-livermore-announces-mammoth-cluster-to-fight-covid-19/ (retrieved November 2020).

[166] HPCWire: It's Fugaku vs. COVID-19: How the World's Top Supercomputer Is Shaping Our New Normal https://www.hpcwire.com/2020/11/09/its-fugaku-vs-covid-19-how-the-worlds-top-supercomputer-is-shaping-our-new-normal/ (retrieved November 2020).

[167] InsideHPC: TACC's Frontera HPC System Expansion for 'Urgent Computing' – COVID-19, Hurricanes, Earthquakes https://insidehpc.com/2020/11/taccs-frontera-hpc-system-expansion-for-urgent-computing-covid-19-hurricanes-earthquakes/ (retrieved November 2020).

climate change, and often with direct consequences to people and society. ARC is also emerging as a vital component for e.g., social sciences and digital humanities. The competitiveness of Canadian industry depends on the efficiency and value-add of its products, the design of which is often aided by ARC systems. A company can for example test a myriad of product variations using ARC, and then only build the most interesting product candidates, reducing cost and shortening time-to-market.

A more rigorous exercise in studying the benefits of ARC should also consider the related costs, in terms of looking into return-on-investment (ROI) of ARC, and in particular in academic settings. This in turn can be a very difficult exercise, since the return in monetary and non-monetary terms is very difficult to quantify, and even the seemingly more straightforward 'investment' side is difficult to estimate accurately in practice.

The US based Coalition for Academic Scientific Computation (CASC) has recently studied academic ROI in terms of total-cost-of-ownership (TCO) and both financial and non-financial benefits. On the denominator, i.e. cost side, one has to consider not only the immediate capital costs (e.g. hardware and software procurement, warranties, licensing, depreciation etc.), but also the on-going operational costs (training, staffing, licensing, power, cooling, networking, maintenance, security, monitoring and billing etc.), facility construction and infrastructure costs and related depreciation.[168] Beyond just the explicit cost, one also needs to consider and define the scope, e.g. if for example the facilities are jointly operated or funded.

On the numerator side, regarding potential benefits of ARC in financial terms, one should consider at least:

1. the benefit to end-user of ARC facilities in research (time saved)

2. ARC system resources (cost savings compared to alternative solutions)

3. personnel resources (value of support from the ARC provider)

4. value of training

5. grant income (monetary value of grant income received v. lost opportunity)

6. products and patents (monetary income), and

7. economic impact (indirect regional financial benefits, e.g., jobs and tax income).[162]

One could though argue that the true benefits of ARC and DRI in general are on the non-financial side, i.e., impacts and outcomes, where indirect and long-term benefits can potentially be almost incalculably beneficial, e.g., considering development of new life saving vaccines, or the earlier mentioned use of modeling to inform decision makers in the context of the COVID19 etc. On a bit more concrete terms such non-financial benefits include:

1. new discoveries reported in publications (improved quality of life for people)

---

[168] Craig E. Stewart et al.: Assessment of financial returns on investments in cyberinfrastructure facilities: A survey of current methods - PEARC '19: Proceedings of the Practice and Experience in Advanced Research Computing on Rise of the Machines (learning) July 2019 Article No.: 33 Pages 1–8 https://doi.org/10.1145/3332186.3332228 (July 2019).

2. people trained in new areas (a better trained workforce for the economy)

3. awards, press notices (reputational benefits to people and organizations), and

4. patents (products improving quality of life, or e.g., the sustainability of human life on earth).[169]

On the business side the US Council on Competitiveness conducted a study on benefits of ARC within its membership. According to the so-called Solve report the best metrics for justifying ARC investment in the industry were 'time to solution', 'inability to solve the problem by any other means', 'ROI', and 'reduced costs compared to physical methods'.[170] Going beyond just the benefits of ARC, there are widely ranging estimates for ARC ROI in the commercial setting, e.g. Hyperion Research reports an average $44 USD return on every dollar invested (i.e. an unrealistic ca. 4300% ROI),[171] while Intersect360 Research hedges stating that the true ROI must be much lower and is very complex to estimate.[172]

# 4. Current State

ARC and DRI are critical to a significant and ever-growing number of researchers because of both technological advances and new research paradigms. These trends are not emerging just in the traditional hard sciences, but also in other disciplines, such as artificial intelligence, natural language processing, social media analysis, large-scale qualitative and quantitative survey analysis, and gene sequencing.

It is difficult to measure the full scale and scope of ARC activities in Canada since these endeavors are so prevalent - there are ARC resources used by Canadian researchers at the research group, departmental, institutional (universities, colleges, research hospitals), research institute (e.g., OICR, OBI etc.), provincial, national, and international levels. Some of these resources are restricted to one group, institute or discipline while many are shared at some level. Some of these resources are primarily for government agencies (e.g., SSC hosted resources used by ECCC), but also include academic usage. This variety should be considered a strength of the Canadian ARC ecosystem, serving the varied and complex geographical, technical, and domain specific needs.

CCF resources are the only ARC resources which are available to literally all Canadian academic researchers. Academic ARC usage in Canada is relatively well known within the main CCF

---

[169] Craig E. Stewart et al.: Assessment of non-financial returns on cyberinfrastructure: A survey of current methods - HARC '19: Proceedings of the Humans in the Loop: Enabling and Facilitating Research on Cloud Computing July 2019 Article No.: 2 Pages 1–10 https://doi.org/10.1145/3355738.3355749 (July 2019).

[170] US Council on Competitiveness: Solve. The Exascale Effect: the Benefits of Supercomputing Investment for U.S. Industry https://www.compete.org/reports/all/2695 (October 2014).

[171] Hyperion Research: HPC Investments Bring High Returns https://www.dellemc.com/resources/en-us/asset/analyst-reports/products/ready-solutions/hyperion-hpc-investment-brings-high-returns.pdf (July 2020).

[172] HPCWire: ROI: Is HPC Worth It? What Can We Actually Measure? By Addison Snell, Intersect360 Research https://www.hpcwire.com/2020/10/15/roi-is-hpc-worth-it-what-can-we-actually-measure/ (retrieved October 2020).

ecosystem, and the current analysis below is based on usage of CCF facilities only. It should be noted that many researchers are unaware of the existence of the CCF resources, think these resources are not for them, or make use of the other non-CCF systems discussed above. Such usage is not captured in the following data.

## 4.1 Registered users of CCF systems

Distribution by position



**Figure 2:** Position of registered CCF users

The Figure 2 above shows the variety of positions the CCF registered users held as of January 1st 2020. The positions are self-identified by the end-users from a list provided by CCF in their central user account database during CCF's annual account renewal. They are also validated by a CCF team member. While most are self-explanatory, there is for example no distinction between different faculty positions (e.g., assistant, tenured, Canada Research Chair etc.), or the term 'researcher' could apply widely and is not strictly defined in this context per se.

In total there were nearly 16 000 'registered' users listed in the CCF database. Registered users are defined as users who, during the annual account renewal process, have registered in the Compute Canada Database as wanting to keep their account active regardless of whether they have logged into any systems (CCF calls such users sometimes 'active users'). The percentages represent the number of positions in the user base and do not indicate how much (if at all) these users used the resources.

The CCDB user pool has significant year-to-year turnover, indicating that there are a lot of people who get exposure to ARC even if they are leaving academia. The high turnover rate requires on-going training on an annual basis or even more frequently. This exposure is contributing to the development of a workforce that is more knowledgeable about and skilled in using DRI tools. The turn-over is estimated to be 20-50% annually, indicating that the cumulative total of unique people who get exposure to ARC is large.

The largest user group is 'faculty' at 27% , followed by 'doctoral students' and master's students, at 23% and 15%, respectively. It is to be noted that most non-PI positions must be sponsored by a PI, which explains in part the high percentage of faculty positions. Including the fourth largest user group, post-docs at 9%, these four largest user groups comprise 74% of all CCF users. Considering the more 'senior' researcher positions, including faculty, adjunct faculty, researcher, and external collaborator positions, these senior researcher type positions amount to roughly 43% of all CCF users. If one were to include post-docs into this group, senior researchers would amount to more than half of all users. The more 'junior' or early career researcher positions, ranging from post-docs to undergraduate students add up to 55% of all users. Regardless of how one considers the seniority of the post-docs, roughly half the users of CCF systems are senior researchers and the other half early career researchers.

The level of one's ARC literacy can vary widely regardless of the career state. There is an ongoing need for training in DRI, due to huge growth in use of AI within many research fields or simply due to an exponential growth in data or computing requirements. ARC itself is continually evolving because technology, software, and research methods continually improve and change so that the researchers need to update their skill sets continually. While at the same time researcher's individual scientific needs and interests will evolve and scale-up requiring more ambitious and evolved, and targeted, DRI tools and support.

Faculty by research area



**Figure 3:** Faculty by research area

The pie chart in the Figure 3 above shows the distribution of faculty positions amongst CCF users across different research disciplines. In total there were roughly 4200 users who identified as faculty in January 2020 (out of the total of ca. 16 000 users are discussed above). Again, these statistics do not indicate actual usage (which will be discussed in Section 4.3 below), only the research discipline distribution with CCF faculty users. The largest user groups are from 'engineering', and 'biological and life sciences', at 19% each. There are multiple disciplines represented by roughly half the number of faculty among the CCF faculty users, namely 'mathematics and statistics', 'computer and information science', 'medical sciences', 'physics', and 'chemistry and biochemistry', all ranging between 7% and 11% of the user base. Looking at groups of disciplines, engineering, mathematics and statistics, and computer science add up to 37% of the faculty user base (i.e., totaling 1530 faculty). Comparing to the observations in the 2017 LCDRI report, in 2017 this cohort was at 35%,[1] so that the relative percentage of faculty users stayed roughly constant between 2017 and 2020. On the life sciences side biological and life sciences, and medical sciences add up to 28% compared to 26% in 2017. The traditional ARC fields of physics, astronomy, chemistry and biochemistry, and environmental and earth sciences add up to 25% compared to ca. 30% in 2017. The 5% point drop in three years is due to a few

percentage point drops in the relative number of both chemistry and biochemistry, and physics faculty users. It should be noted that some disciplines, e.g., medical sciences, and social sciences, often necessitate using alternative systems such as highly secure or personal health information (PHI) compliant infrastructure.

The humanities, social sciences, business, and psychology faculty account for roughly 10% of the faculty user base at CCF, while the faculty in these disciplines (humanities; social and behavioural sciences and law, and business, management, and public administration) represented roughly 46% of all the full-time academic faculty in Canada,[173] clearly indicating how the number of ARC users in these disciplines is not representative of their overall relative representation in Canadian academia, keeping in mind that one should not expect all disciplines to need equal access to ARC resources. These disciplines have been growing their usage of ARC since 2017 when they accounted for roughly 7% of the user base. The growth of 3% percentage points corresponds to roughly 40% growth for this user group over three years, while the absolute number of users in these disciplines still has substantial room and potential to grow. In addition to the number of faculty, the number of sponsored users per faculty member is smaller in social sciences and other underrepresented fields. There is a distinction between underrepresentation and the needs not served, with the latter being just one potential reason for the former. But in general, like all other disciplines, social sciences and humanities are more and more impacted by digital technology. The increase in the amount of available data and digitized content is accompanied by a growing use of computational methods. As an example, Canada stands out internationally for its research in the field of digital humanities. In this field, the current ARC/DRI offer does not cover the specific needs of users. There is every reason to believe that the evolution of services could facilitate the adoption by many other users in the humanities and social sciences whose needs are insufficiently addressed.

---

[173] Canadian Association of University Teachers (CAUT): Age Distribution of Full-time University Teachers by Age, Sex and Major Discipline, 2017-2018 https://www.caut.ca/sites/default/files/3.16_age_distribution_of_full-time_university_teachers_by_sex_and_major_discipline_2017-2018.xlsx (retrieved April 2021).

## 4.2 What are the main CCF systems for ARC delivery?



**Figure 4:** Canada's national ARC platform

The Figure 4 above shows the five national hosting sites across Canada. Recent substantial CFI investments in Canadian Cyberinfrastructure have resulted in new systems being deployed and a major increase in ARC capacity and capability in Canada. This process also involved consolidating the number of CCF facilities to five main data center sites.[174] These five sites host five national systems and are affiliated with regional CCF member organizations as follows, from West to East:

- University of Victoria, Arbutus (WestGrid);

- Simon Fraser University, Cedar (WestGrid),

- University of Waterloo, Graham (Compute Ontario),

- University of Toronto, Niagara (Compute Ontario); and

- McGill University/ École de technologie supérieure, Béluga (Calcul Québec).[175]

---

[174] Compute Canada: Renewing Canada's Advanced Research Computing Platform https://www.computecanada.ca/techrenewal/ (retrieved September 2020).

[175] Compute Canada: Available Resources https://www.computecanada.ca/research-portal/accessing-resources/available-resources/ (retrieved September 2020).

The general characteristics of these national systems are:

- **Arbutus is a general-purpose ARC cloud system** for hosting (mostly Linux based) virtual machines and other cloud workloads.[176] It is based on an open-source OpenStack cloud infrastructure and has over 16,000 Intel CPU cores, spanning over 450 nodes with a total of 140 TB of memory (i.e., 300 GB per node or 10 GB per core). The communication backend is Ethernet based ranging from 10 to 25 gigabit Ethernet. The main storage capacity is 17 PB. The system has minimal GPU accelerator capacity.[177]

- **Cedar is a general-purpose heterogeneous ARC cluster** for a variety of ARC workloads.[170] It has nearly 95 000 Intel CPU cores, spanning over 2470 nodes with available memory per node ranging from 125 GB to 3 TB (i.e., from 4 GB to 90 GB per core). The communication backend is a high-speed low-latency 100 Gbit/s Intel Omni-Path fabric. The tiered storage system ranges from small persistent 'home' (at over 500 TB total capacity), and persistent 'project' (at 23 000 TB), to high-speed non-persistent 'scratch' (at 5400 TB). The system has a total of 1350 Nvidia GPU cards available as accelerator capability.[178]

- **Graham is a general-purpose heterogeneous ARC cluster** for a variety of ARC workloads.[170] It has nearly 42 000 Intel CPU cores, spanning over 1185 nodes with available memory per node ranging from 124 GB to 3 TB (i.e., from 4 GB to 50 GB per core). The communication backend is high-speed low-latency 56 Gbit/s (100 Gbit/s) Mellanox FDR (EDR) InfiniBand fabric. The tiered storage system ranges from small persistent 'home' (at over 130 TB total capacity), and persistent 'project' (at 16000TB), to high-speed non-persistent 'scratch' (at 3600 TB). The system has a total of 520 Nvidia GPU cards available as accelerator capability. Roughly one fourth of the GPU cards are the latest generation Turing T4 cards that are designed for deep learning workloads.[179]

- **Niagara is a massively-parallel homogeneous ARC cluster** for scalable ARC workloads.[170] It has nearly 81 000 Intel CPU cores, spanning over 2016 nodes with available memory per node fixed at 200 GB (i.e. 5 GB per core). The communication backend is a high-speed low-latency 100 Gbit/s Mellanox EDR InfiniBand fabric leveraging leading edge Dragonfly+ topology. The tiered storage system ranges from small persistent 'home' (at 200 TB total capacity), and persistent 'project' (at 2000 TB), to high-speed non-persistent 'scratch' (at 7000 TB), and very-high speed non-persistent 'burst' (at 230 TB). The system used to have no GPU accelerator capability, but thanks to the recent expansion funding it currently has 64 GPUs.[180]

- **Béluga is a general-purpose heterogeneous ARC cluster** for a variety of ARC workloads.[170] It has nearly 35 000 Intel CPU cores, spanning over 872 nodes with available

[176] Compute Canada: National Systems https://www.computecanada.ca/techrenewal/national-systems/ (retrieved September 2020).

[177] Compute Canada: Arbutus cloud https://docs.computecanada.ca/wiki/Cloud_resources#Arbutus_cloud_.28arbutus.cloud.computecanada.ca.29 (retrieved September 2020).

[178] Compute Canada: Cedar https://docs.computecanada.ca/wiki/Cedar (retrieved September 2020).

[179] Compute Canada: Graham https://docs.computecanada.ca/wiki/Graham (retrieved September 2020).

[180] Compute Canada: Niagara https://docs.computecanada.ca/wiki/Niagara (retrieved September 2020).

memory per node ranging from 90 GB to 750 GB (i.e., from 2 GB to roughly 20 GB per core). The communication backend is a high-speed low-latency 100 Gbit/s Mellanox EDR InfiniBand fabric. The tiered storage system ranges from small persistent 'home' (at over 100 TB total capacity), and persistent 'project' (at 25000 TB), to high-speed non-persistent 'scratch' (at 2600 TB). The system has a total of ca. 688 Nvidia GPU cards available as accelerator capability.[181]

## 4.3 What is the current and past usage of ARC in CCF facilities?

Registered Users:



**Figure 5:** CCF registered users

In the Figure 5 above the number of registered users in the CCF is plotted as a function of time, as of January 1st for each year. 'Registered users' refers to users in the CCF user account database who have kept their account active during CCF's annual account renewal (CCF often refers to such users as 'active users'). It should be noted that many of these people do not necessarily run jobs during the year, or even access the systems. Some users just need to or want to keep their accounts active if needed in future, or potentially the research project was delayed or got redirected. All Principal Investigators (PIs) need accounts to sponsor their students but might not access the ARC systems or run jobs themselves. Such faculty and PIs are critical as facilitators and funders of the research even though they are not strictly speaking explicit day-to-day users of the ARC systems. Indirect evidence based on software downloads as discussed

---

[181] Compute Canada: Béluga https://docs.computecanada.ca/wiki/B%C3%A9luga/en (retrieved September 2020).

later in this document indicates that roughly 10,000 users were actively using the main systems in January-October 2020, compared to roughly 18,000 registered users in October 2020.

The number of registered users has grown significantly in the last decade, from 829 in 2010 to 15 994 in 2020, corresponding to over 1700% growth over ten years. Notably the growth has been relatively linear over the years but split in two segments. Between 2010 and 2013 the growth was very fast, then almost static from 2013 to 2014, and then steady and slower growth from 2014 until the present day. The overall compounded annual growth rate (CAGR) has been roughly 34% year-on-year over the last decade. Until 2013 the CAGR was even higher, at ca 110%, while since 2014 the CAGR has been roughly 12%.

As of late October 2020, the total number of users with a registered CCF account in the CCF systems was roughly 18,000 users, indicating an increase by roughly 2000 users during the first ten months of the year, i.e. slightly above the annual CAGR. The regional distribution of these users was as follows:

- Westgrid: 5100 affiliated registered users,

- Compute Ontario: 4800 (SharcNet: 2765, CAC: 1019, SciNet: 1865)

- Calcul Québec: 4600, and

- ACENET: 1000,

totaling 16,400 registered users who were affiliated with universities that in turn were members of their regional CCF organizations. On the flip side roughly 900 users were not affiliated with any of the main CCF affiliate organizations (ca. 700 users did not report any affiliation in the database).

The number of institutions using CCF resources in late October 2020 was over 600, far exceeding the number of Canadian institutions and indicating the long tail of international collaboration. The institutions with more than 500 CCF users were, in descending order:

- McGill University (1700 registered users),

- University of Toronto (1500),

- University of British Columbia (1500),

- University of Alberta (1100),

- University of Waterloo (800),

- Université de Montreal (800),

- Simon Fraser University (600), and

- University of Western Ontario (600).

The total number of registered users from these eight institutions was roughly 8500 users, i.e., ca. 47% of all registered CCF users. On the other hand, over 600 institutions had ten or fewer users per institution, indicating a very long tail in the diversity of users and institutions.

CPU usage:



**Figure 6:** Historical ARC CPU usage (aggregate and per research discipline)

The Figure 6 above shows the historical total and discipline specific CPU resource usage in CCF aggregate systems. The CPU resource usage is measured in CPU core years, i.e., the compute resource used by running a program on a single CPU core for one full year. Notably this metric does not consider the changes in computational power of the CPU, i.e., the advances in CPU architecture and technology as new generations of CPUs are introduced. The time is indicated in full calendar years.

The thick solid black line shows the total absolute CPU core year usage over the last decade (see the vertical axis on the right for units), while the dashed blue line indicates the CPU supply in core-years per RAC application process (highlighting that the demand is limited by supply). In 2010 the usage was ca. 50 000 CPU-years, while in 2019 the usage was ca. 200 000 CPU-years, indicating a four-fold increase in CPU resource usage. This corresponds to roughly 16% CAGR. Notably the growth was quick from 2010 to 2013 and then stagnated at ca. 150 000 CPU-year level between 2013 and 2017 due to the 9-year funding gap for ARC systems between NPF in 2006, and CI in 2015, and related budget finalizations and spacing of expenditures. Since 2017 the usage has grown rapidly again and in nearly linear fashion as the new CFI modernization funded systems came online. Historically as well as currently the usage in ARC (both in Canada and globally) has been limited by supply so that the usage doesn't go up since users decide to use more resources, i.e., due to increases in demand per se. The usage increases annually since the researchers have access to more resources.

The thin colored solid lines and corresponding shaded areas show the relative distribution of the CPU-year usage between different disciplines (see the vertical axis on the left for units). In 2010 the largest user group was Biological and Life Sciences at ca. 30% followed by Physics and Astronomy at ca. 25% (astronomy was not tracked separately in 2010), Chemistry and Biochemistry at ca. 18%, and Engineering at ca. 12% of the total usage. All in all, these four disciplines used roughly 85% of the resources.

In 2019 the largest user group was:

- Engineering at ca. 28%,

- Physics (ca. 20%) & Astronomy (ca. 5%) at combined ca. 25%, and

- Chemistry and Biochemistry at ca. 20%.

These three research fields consumed roughly three quarters of the resources, while roughly 25% was used by other fields, including:

- Biological and Life Sciences (at ca. 8%),

- Environmental and Earth Sciences (at ca. 6%),

- Computer and Information Science (at ca. 4%),

- Mathematics and Statistics (at ca. 2%), and

- Medical Sciences (at ca. 2%).


These five disciplines consumed a total of roughly 22% of the resources, leaving only 3% to 'Other' (at ca 1%), followed by Social Sciences, Psychology, Business, and Humanities.

In the last decade the major trends have been the strong growth in Engineering, the drop in relative position of Biological and Life Sciences resource consumption, and emergence and strong growth of new disciplines leveraging ARC. Among the top three, Engineering grew from 12% allocation to 28% over the last decade, while Physics & Astronomy and Chemistry and

Biochemistry stayed roughly the same. The gains by Engineering were to a degree mirrored by the drop in Biological and Life Sciences that 'only' consumed ca. 8% of all the resources in 2019, while compared to 2010, when the consumption was ca. 30%. In absolute terms the resource usage by Biological and Life Sciences has been constant at ca. 16 000 CPU-years while the relative allocation has dropped four-fold over the last decade. In the 'middle-of-the pack' resource consumer range (between 2-6% per discipline) Mathematics and Statistics (over seven-fold absolute growth) and Environmental and Earth Sciences (nearly seven-fold absolute growth) increased their resource usage relative to others over the last decade, i.e., at faster rate than the overall four-fold growth. Medical Sciences usage grew at roughly 4.5 growth rate, while Computer and Information Science usage 'only' grew at 2.5 rate over the decade.

Interestingly Social Sciences, Psychology, Business, and Humanities increased their CPU resource usage substantially over the last decade. In 2010 this group consumed roughly 0.3% of the CPU resources, while in 2019 they consumed roughly 0.6% of the ARC resources, corresponding to overall two-fold growth in usage compared to other disciplines. This growth is even more impressive in absolute terms, roughly eight-fold growth over the decade from 157 CPU-years in 2010 to roughly 1250 CPU-years in 2019. The growth in the ARC CPU resource usage in these disciplines indicates strong interest, and potential these disciplines have for leveraging DRI in future. Keeping in mind that given the type of tasks performed in these disciplines, the raw CPU cycle measurement is not necessarily relevant for reporting on the use of computer resources.

**CPU usage per position type**



**Figure 7:** CPU usage per position type

The Figure 7 above shows the relative CPU usage per position type across the CCF systems, for the period ranging from March 2020 to February 2021. The data is not normalized based on the number of users in each group. It shows the total CPU time used by users in each position type. Doctoral students as a whole is the largest user group at 37%, followed by Postdoctoral Fellows at 16%, and Master's students at 13%, adding up to two thirds of the total CPU resource usage. Users in the Faculty position directly use 9% of the CCF CPU resources. The remaining quarter of the resources is used by a variety of external and internal groups, with undergraduate students at 5% of the total. GPU usage in general is similar to CPU usage per position type, except that Master's students use 22% of the GPU resources, i.e. almost double from their relative CPU usage, and Faculty uses only 2% of the GPU resources compared to 9% of the CPU uses.

Comparing the CPU resource usage to the number of users in various positions (as discussed earlier in this document) provides interesting insights: Faculty represents 27% of the CCF user base, while using 9% of the CPU resources. Postdoctoral fellows, and doctoral students comprise 32% of the user base, while consuming 53% of the CPU resources. While Master's students leverage the resources at roughly the same level as their numbers would imply (13% usage compared to 15% of the user base). In general, this kind of under- and over-representation in resource usage is not surprising given that faculty and senior researchers often need to focus on guiding their research teams, writing grant applications and public outreach, leaving less time for running workloads on the ARC systems.

**CPU usage per host site v. PI's region**

| Regional resource | PI's regional affiliation | | | | |
|---|---|---|---|---|---|
| | | ACENET region | Compute Ontario region | Calcul Quebec region | Westgrid region | Grand Total (CPU years) |
| | beluga-compute | 3.68% | 11.14% | **68.02%** | 17.16% | 23,052 |
| | cedar-compute | 9.29% | 22.11% | 17.03% | **51.56%** | 74,780 |
| | graham-compute | 14.34% | **47.49%** | 14.17% | 24.00% | 30,573 |
| | mp2 | 3.77% | 21.51% | **65.72%** | 9.00% | 15,723 |
| | niagara-compute | 0.82% | **73.07%** | 7.77% | 18.33% | 71,774 |
| | **Grand Total** | 6.19% | 41.43% | 22.54% | 29.84% | 215,903 |

**Table 1:** Regional CPU resource usage allocation per PI's regional affiliation

The Table 1 above shows the regional CPU resource usage allocation per CCF region based on the principal investigator's home institution. The data is for a 12-month period, from April 2020. Each line corresponds to one of the 4 national sites, in addition to Mammouth Parallel 2 (mp2), a legacy cluster still allocated for 2020. The columns three to six then show the regions PIs come from. The values for each row / resource add up to one hundred percent (not shown as a percentage, but the absolute total value is shown in the last column) and show the regional usage distribution of the resource. For example, users affiliated with Calcul Québec use 68% of the CPU cycles on Beluga, while only 8% of all CPU cycles on Niagara. The usage that is local to the infrastructure, i.e., usage by users from that region is highlighted in bold for clarity. To assign proper absolute context to the relative distribution on each row, the last column lists the total CPU years for each resource respectively. The primary resources are Niagara at 72k CPU years and Cedar at 75k CPU years, while Beluga and Graham supply 23k and 31k CPU years respectively. The last row indicates the relative distribution of the total 216k CPU year allocation between users in different regions. The relative percentages in each column (rows three to six) are not meant to add up to the 'grand total' percentage value at the bottom row.

Focusing on the last row and the last column, users in Compute Ontario region (4800 registered users) consume roughly 41% of all CPU cycles, while Compute Ontario provides roughly 47% of the compute cycles (i.e., 102k of 215k). The next largest user group originates from the Westgrid

region (5100 registered users), consuming roughly 30% of the CPU cycles, while Westgrid provides roughly 35% of the compute resources. Users originating from regions affiliated with Calcul Quebec (4600 registered users) use roughly 23% of the total, while Calcul Quebec's Beluga cluster provides roughly 11% of the compute resources. Acenet regional users (100 registered users) consume roughly 6% of the total CPU cycles. The rest of the compute capacity is provided by the legacy `Mammouth Parallel 2 cluster (mp2), operated by Calcul Québec.

Looking at rows three to six, i.e., regional distribution of users for each service queue, the general trend is that each 'local' resource is primarily used by the users from that same region. At the more locally focused end, 73% of the CPU resources on Niagara are used by Compute Ontario affiliated users. Similarly, roughly 68% of the usage on Beluga is by Calcul Quebec affiliates. The two other main clusters, Graham and Cedar, on the other hand show locally originated usage at roughly 50% level. The remaining half on Westgrid's Cedar is roughly evenly distributed between users from Compute Ontario and Calcul Quebec regions. In the case of Graham, users from the Westgrid region use roughly 25% of the resources, while the remaining quarter is split roughly evenly between Acenet and Calcul Quebec affiliates.

Independent analysis of regional usage versus PI's regional affiliation has been done by Compute Ontario, adding some color to the above data. The results indicate that large majority of PIs make use of their regional system but also that at least 50% of the PIs in any region also computed outside their region (and 18% of PIs used 3 of the 5 systems). The analysis also investigated the role of the RAC allocation and process, showing that e.g., Ontario based PIs used 40% of all the cycles while also receiving 40% of the RAC allocations. The regional usage also aligns well with population data, i.e., users from different regions use ARC resources corresponding to the population size in that region.

We note that the above correlation between the location of an infrastructure and that of the researchers is mostly historical and complicated. All CCDB-registered PIs can access all systems (whether or not they have an allocation). The choice of which system(s) they use for computation can be driven by many factors including technical as well as more individual reasons. Some obvious factors include: RAC allocation (awards are made for specific systems), availability of hardware (e.g., large-memory nodes, a particular type of GPU), scheduling policy (e.g., allowing jobs with long wall-clock times or large core-counts), training, recommendations by colleagues, personal perception of system performance, and other personal preferences. The need to have data available locally for efficient computation (as well as data management issues in general) may tend to minimize the number of systems a PI uses but, on the other hand, using more systems increases the amount of "default" cycles available to a PI and can help reduce wait times. The distribution of PIs and their usage can also relate to the timing of when systems went into production. Globally, we can calculate that about 60% of the CPU cores are being used by users from the region that hosts the infrastructure, while 40% of the CPU cores are being used by users from other regions.

**Figure 8:** Historical ARC GPU usage (aggregate and per research discipline)

The Figure 8 above discusses the historical GPU accelerator usage at CCF systems over the last decade. The absolute units on the right-hand-side vertical axis are GPU years, i.e., the equivalent of executing a program on one GPU processing unit for a full calendar year. The compute efficiency of this unit does not stay constant over time due to architectural advancements, and the number of 'micro' compute cores within the massively parallel GPU can change over time, for example the latest generation Nvidia Ampere GPU card can have over 8000 CUDA compute cores per GPU while the previous generation Nvidia Turing GPU processing unit can have around

3000 CUDA cores per GPU.[182] Compute efficiency increases is also true for compute nodes (core counts increase) and compute cores (due to microprocessor improvements though sadly not to the same degree any longer).

The thick black solid line indicates the growth in GPU usage in absolute terms. In 2010 the usage of GPUs at CCF was practically zero, as GPU computing was emerging as an important ARC trend at that point. The first year the CC GPUs were in production and being monitored and reported was 2012, with 58 GPU years being used. In 2019 the total usage of GPU resources was ca. 1300 GPU years, corresponding to roughly 56% CAGR since 2012. As discussed elsewhere in this document this growth was severely restricted by the available supply and is thus not indicative of actual growth rate of GPU computing in general.

The thin colored solid lines and corresponding shaded areas show the relative distribution of the GPU-year usage between different research disciplines (see the vertical axis on the left for percentages). In 2019 the largest user group was Computer and Information Science at 40%, followed by Chemistry and Biochemistry (24%), Biological and Life Sciences (12%), and Physics & Astronomy (11%), All in all these four disciplines used roughly 87% of the GPU resources in 2019. Historically, GPU usage has first and foremost been associated with molecular dynamics (Chemistry and Biochemistry), with GROMACS and NAMD being the primary applications. More recently, and comparing to their respective CPU usage, Computer and Information Science has become much more prominent as a GPU user at CCF.

---

[182] NVIDIA White Paper: NVIDIA AMPERE GA102 GPU ARCHITECTURE https://www.nvidia.com/content/dam/en-zz/Solutions/geforce/ampere/pdf/NVIDIA-ampere-GA102-GPU-Architecture-Whitepaper-V1.pdf (September 2020).

**GPU usage per host site v. PI's region**

| Regional resource | | PI's regional affiliation | | | | |
|---|---|---|---|---|---|---|
| | | ACENET region | Compute Ontario region | Calcul Quebec region | Westgrid region | Grand Total (GPU years) |
| | beluga-gpu | 3.37% | 17.46% | 69.65% | 9.52% | 547.47 |
| | cedar-gpu | 2.21% | 29.99% | 27.02% | 40.78% | 989.13 |
| | graham-gpu | 8.29% | 50.08% | 24.92% | 16.71% | 306.49 |
| | helios-gpu | 0.00% | 52.43% | 46.12% | 1.45% | 87.40 |
| | **Grand Total** | 3.41% | 30.64% | 39.64% | 26.31% | 1930.49 |

**Table 2** Regional GPU resource usage allocation per PI's regional affiliation

The Table 2 above shows the regional GPU resource usage allocation per CCF region based on the principal investigator's home institution. The data is for a 12-month period, from April 2020. Each line corresponds to one of the 3 national sites that have GPUs, in addition to Helios (helios), a legacy cluster still allocated for 2020 and 2021. The columns three to six then show the different regions PIs originate from. The values for each row / resource add up to hundred percent (not shown as a percentage, but the absolute total value is shown in the last column) and show the regional usage distribution of the resource. For example, users affiliated to the Calcul Québec region use 70% of total available GPU cycles on Beluga, and 27% of all available GPU cycles on Cedar. To assign proper absolute context to the relative distribution on each row, the last column lists the total GPU years for each resource respectively. The primary resources are Cedar at 990 GPU years, and Beluga at 550 GPU years, while Graham supplies 310 GPU years. The legacy Helios system, operated by Calcul Québec, provided 90 GPU years since April 2020. The last row indicates the relative distribution of the total 1930 GPU year allocation between users in different regions. The relative percentages in each column (rows three to six) are not meant to add up to the 'grand total' percentage value at the bottom row.

Focusing on the last row and the last column, users in Calcul Quebec region consume roughly 40% of all GPU cycles, while Calcul Quebec provides roughly 33% of the compute cycles (i.e., 634 out of 1930 GPU-years). The next largest user group originates from the Compute Ontario region, consuming roughly 31% of the GPU cycles, while Compute Ontario provides roughly 16% of the compute resources. Users originating from regions affiliated with Westgrid use roughly 26%

of the total, while Westgrid's Cedar cluster provides roughly 51% of the GPU compute resources at CCF. Acenet regional users consume roughly 3% of the total GPU cycles.

Looking at rows three to six, i.e., regional distribution of users for each service queue, the general trend is that each 'local' resource is primarily used by the users from that same region (excluding the small usage of Helios). At the more locally focused end, 70% of the GPU resources on Beluga are used by Calcul Quebec affiliated users. Compute Ontario affiliated users then use 17% of Beluga's GPU resources, while Westgrid affiliates use ca. 10%. The two other main clusters, Graham and Cedar, on the other hand show locally originated usage at roughly 50% and 41% level, respectively. In the case of Graham, users from Calcul Quebec region use roughly 25% of its GPU resources, while the remaining quarter of the GPU resource is split between Westgrid affiliated users at 17% and Acenet affiliates at 8%. The remaining 61% of Westgrid Cedar's GPU capacity is roughly evenly distributed between users from Compute Ontario and Calcul Quebec regions, at 30% and 27% respectively. We note that this correlation between the location of an infrastructure and that of the researchers is mostly historical and accidental. Indeed, users tend to migrate to a new cluster mostly when older ones are being retired. Technical considerations are considered when CCF staff recommend a specific site to a research group. Globally, we can calculate that about 50% of the CPU cores are being used by users from the region that hosts the infrastructure, while 50% of the CPU cores are being used by users from other regions.

## Software usage

CCF ARC infrastructure supports and provides a wide array of research software to the user community. The software distribution mechanisms and availability vary even between main systems, the details of which are discussed in chapter 4.5 of this document. One of the key software distribution mechanisms centrally managed by CCF staff leverages 'module' packaging technology and is available to users at all sites.

**Figure 9:** Number of distinct users of software modules

The Figure 9 above shows the number of distinct users (y-axis, on a logarithmic scale) that have loaded specific software packages via CCF's software modules packaging system. Each software module package is indicated with a numbered label (x-axis). In the first ten months in 2020, over 700 different software packages were loaded using CCF's central 'module' mechanism. 15 SW modules were loaded by over 1000 distinct users, while over 600 SW modules were loaded by less than 100 distinct users each. Moreover, of such 600 SW modules, 450 SW modules were loaded by fewer than 20 distinct users each. That is, the SW usage has a very **'long tail'** with a myriad of pieces of SW used by a relatively small number of users, clearly putting stress in the SW maintenance and support ecosystem. In addition, most of these software packages are built for multiple versions of the package, as well as multiple generations of CPU architectures by the CCF team, to ensure optimal performance on a given infrastructure.

The most popular, default, packages were loaded by close to 10,000 users, indicating the number of active users who executed jobs in the CCF systems during the first ten months of 2020 (excluding many jobs on 'niagara'). Since many software modules are loaded as part of interconnected and interdependent sets of software packages it is very difficult to interpret actual

explicit software usage and popularity from the individual module loads. One can nevertheless deduce some information for software packages that are relatively independent and need to be explicitly loaded by the end-user. Looking at such cases, the GCC compilers were loaded roughly 3000 times, indicating that ca. 30% of users wanted to switch from the default Intel compilers to GNU compilers. The module loads also indicate that roughly 14% of active users were running or at least were interested in GPU computing (since CUDA libraries were loaded by ca. 1400 users). The R scientific package is a relatively independent scientific computing package, so that the roughly 1100 downloads of the R module package indicate that it was used by roughly one tenth of the active users. Since Python is used by many software packages, one can not draw similar conclusions based on the corresponding ca. 4100 loads of the Python module package. For example, in the CCF software stack, 38 modules depend explicitly on Python, while only 3 modules depend explicitly on R.

## Cloud usage

In January 2021 CCF offered cloud resources on multiple systems and regions, including Arbutus (ca. 16,000 CPU cores), East (ca. 600 cores), Cedar (ca. 1000 cores) and Graham (ca. 800 cores), totaling ca. 18,400 CPU cores. As the numbers indicate the dedicated Arbutus cloud system contributes 87% of all cloud resources in the CCF infrastructure. The CCF supercomputing infrastructure has roughly 268k CPU cores available, so that the cloud offering is roughly 7% of the total CPU compute capacity. Looking at the GPU resources on the cloud side Arbutus offers ca. 100 GPUs, while CCF's traditional supercomputers have combined ca. 2500 GPUs. That is, in relative terms CCF cloud systems have roughly half of the GPU (ca. 4%) capacity compared to available CPU capacity. There are two main types of users: those who launch virtual machines (VMs) as persistent instances, and those who launch compute / temporal platform-as-a-service instances or software-as-a-service services via dedicated middleware platforms e.g., CANFAR, Syzygy etc. Some CCF users might not adhere to these categories as savvy users are known to take and leverage whatever resources they can if technically available.

| Research Area | Total vCPU Core Years | Total Projects |
|---|---|---|
| Anthropology | 9.45 | 1 |
| Astronomy | 2,990.68 | 51 |
| Biological and Life Sciences | 2,785.39 | 52 |
| Business | 97.55 | 12 |
| CCF Staff | 279.25 | 36 |
| Chemistry and Biochemistry | 887.56 | 8 |
| Computer and Information Science | 2,093.78 | 69 |
| COVID-19 | 7,815.38 | 4 |
| Engineering | 767.37 | 35 |
| Environmental and Earth Science | 525.54 | 20 |
| FRDR | 143.26 | 4 |
| History | 9.02 | 1 |
| Humanities | 573.52 | 37 |
| Mathematics and Statistics | 359.19 | 11 |
| Medical Science | 428.37 | 23 |
| National Team or Service | 350.46 | 24 |
| Physics | 7,069.72 | 18 |
| Psychology | 160.79 | 13 |
| Research Data Management | 0.18 | 1 |
| Social Science | 348.84 | 35 |
| Training | 232.44 | 11 |
| **Grand Total** | **27,927.76** | **466** |

**Table 3:** Cloud usage per research area in 2020 on Arbutus

Cloud usage across the above systems is currently not monitored in a detailed fashion by CCF. The Table 3 above lists the cloud usage per research area on CCF's primary cloud provider system, the Arbutus cluster, in calendar year 2020. The usage in the second column is reported as vCPU core years with hyperthreading enabled, i.e., each physical CPU core is presented as two virtual cores by the operating system and for the purposes of resource allocation, effectively doubling the available compute resource presented to the end-users. The total utilization in 2020 was ca. 28,000 CPU-years which is roughly 13,000 CPU years more than the 15,000 strict non-hyperthreaded CPU-year capacity of this cluster. With hyperthreading based oversubscription

Arbutus is presented as a 30,000 CPU core resource. Oversubscribing is a standard feature of well managed cloud systems, particularly with server workloads that do not always utilize the resources at 100% as traditional ARC workloads often do. The oversubscription ratio in Arbutus is roughly 1.9. An optimal oversubscription ratio depends strongly on the type of workloads and can range from 1.0 (i.e., no oversubscription, e.g., traditional CPU compute intensive ARC workloads), to over 5 (e.g., for lightweight office applications, or lightly used web portals).

The leading disciplines using Arbutus last year were Covid-19 (ca. 7.8k CPU years), Physics (7.1k), Astronomy (3.0k), Biological and Life Sciences (2.8k), and Computer and Information Sciences (2.1k), totaling ca. 22.8k CPU years, i.e., 81% of the total usage. Covid-19 research leading in resource usage is an interesting indicator how cloud resources can be flexibly deployed for new research and needs. Interestingly the Covid-19 research was done by only four projects, compared to e.g., Physics, where roughly the same amount of compute cycles was distributed over 18 projects.

Another interesting comparison can be made between Physics and Astronomy (ca. 10k CPU-years, 69 projects), and Social Sciences and Humanities (0.9k CPU years, 72 projects), indicating positive uptake of cloud computing by the underrepresented disciplines when it comes to number of projects (with much smaller compute power requirements than the comparison group).

| Research Area | Total Users |
|---|---|
| Astronomy | 41 |
| Biological and Life Sciences | 189 |
| Business | 20 |
| CCF Staff | 236 |
| Chemistry and Biochemistry | 20 |
| Computer and Information Science | 221 |
| Engineering | 87 |
| Environmental and Earth Science | 76 |
| Humanities | 57 |
| Mathematics and Statistics | 28 |
| Medical Science | 70 |
| Physics | 37 |
| Psychology | 29 |
| Social Science | 63 |
| **Grand Total** | **1174** |

**Table 4:** Arbutus users by Research Area in 2020

The Table 4 above lists the number of cloud users per research discipline (using slightly different nomenclature as in the previous table per CCF accounting) in 2020 on Arbutus. The largest user group is Computer and Information Sciences (221 users), corresponding to the largest number of cloud projects (69 projects). Interestingly there are less Astronomy and Physics users (78) than Social Sciences and Humanities users (120).

The above statistics are nevertheless difficult to interpret and can be misleading. For the statistics on the number of cloud users, it should be noted that several groups are running large portals that serve large communities. For example, the Syzygy project by James Colliander will look like a single user in the above stats, while hosting on average 1000 unique users every day and have a user base in tens of thousands and growing. Other platforms like CANFAR[183], GenAp[184], Magic Castle[185], or iReceptor[186] that are serving large communities are not properly captured by the current statistics either.

It should be further noted that the above data does not include other CCF cloud resources, or commercial cloud usage, the latter of which could be substantial. Further targeted research is clearly needed to assess the scope of DRI cloud usage in Canada.

## Storage usage

Due to resource constraints the ARC WG did not have the opportunity to analyse the actual storage usage patterns across the CCF host sites and storage systems. Instead, the Alliance has formed a separate Storage WG that will review both the current storage usage and future storage needs across the CCF infrastructure. This work will not only consider active storage as covered under CCF's current mandate from CFI, but also nearline, repository, and archival storage needs and policies. The Storage WG's findings will be available in the Spring-Summer 2021 and will contribute to the Alliance's Strategic Plan and New Service Models in the Fall 2021.

Instead of looking at actual usage of storage, the CCF RAC gives a window to current and historical active storage requests and allocations. Such RAC aggregate active storage data is discussed in Chapter 4.5. It should be noted that the end-user facing RAC process data does not consider backup, repository, and archival data needs.

[183] CANFAR: CADC https://www.canfar.net/en/nodes/cadc/ (retrieved April 2021).
[184] Compute Canada: Canadians lead in transforming genomic data into knowledge to drive medical innovations https://www.computecanada.ca/news/canadians-lead-in-transforming-genomic-data-into-knowledge-to-drive-medical-innovations/ (retrieved April 2021).
[185] InsideHPC: Compute Canada's Magic Castle: Terraforming the Cloud for HPC https://insidehpc.com/2020/02/compute-canadas-magic-castle-terraforming-the-cloud-for-hpc/ (retrieved April 2021).
[186] SFU: iReceptor Architecture http://ireceptor.irmacs.sfu.ca/architecture (retrieved April 2021).
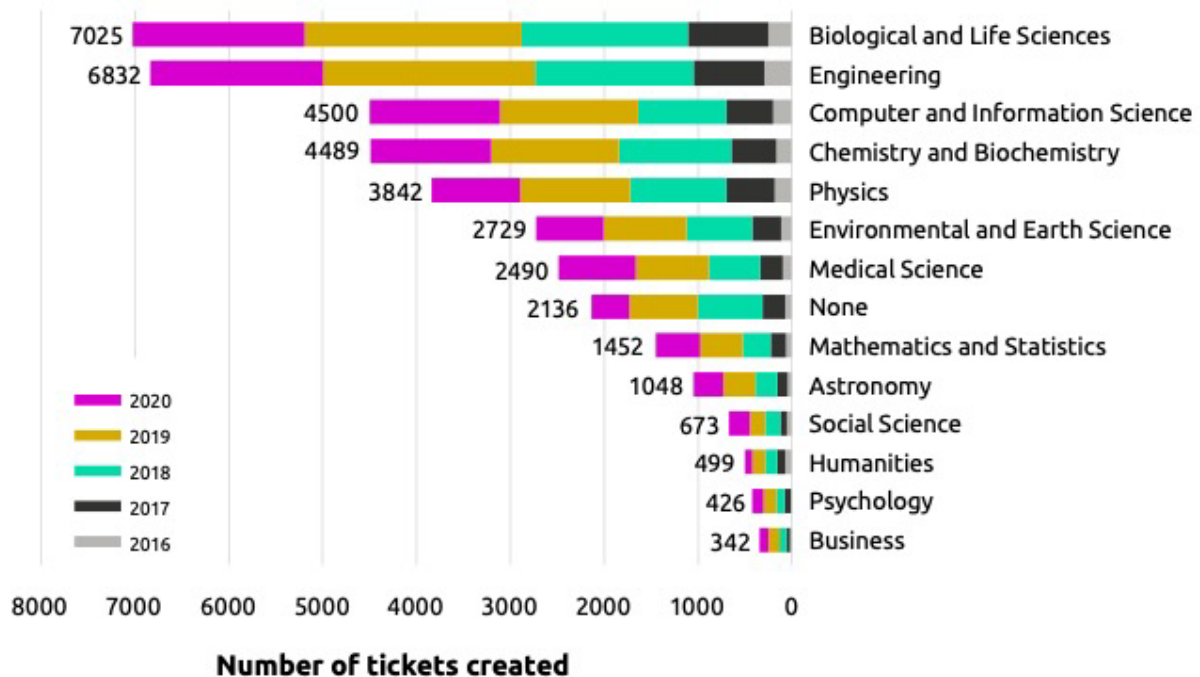
IT tech support usage



**Figure 10:** Support ticket distribution by the research area

The Figure 10 above shows the number of support tickets submitted to the CCF support team between 2016 and late October 2020. Each horizontal bar corresponds to different research disciplines while each color corresponds to different years. The total length of each bar is the cumulative number of support tickets for that discipline. The total number of tickets is roughly 38,500 over the five years. Historically the number of tickets was lower in 2016 and 2017, at ca. 1600 and 4400 respectively, but since then the growth in total number of tickets has slowed down with ca. 9600 in 2018, and ca. 12300 in 2019. Much of the early increase was due to users migrating from legacy systems with regional/institutional support to national systems with national support. The end of October total for 2020 is ca 10,500 indicating that in 2020 the total number of tickets will be roughly the same as in 2019. The vast majority of these CCF support tickets are related to general use of the ARC and CCF clusters and infrastructure, that is, a large majority of the support tickets are not related to domain specific scientific needs.

The growth in number of tickets has been much larger since 2016 than the corresponding increase in number of users (as discussed earlier, the CAGR for the CCF user growth has been ca. 12% since 2014), indicating how CCF central team has grown the support service and brought it up to speed to serve the community. The stabilization in 2019 and 2020 in the number of tickets processed is due to a few factors. As discussed above, much of the early growth (2016, 2017, 2018) was due to the transition from legacy systems to national systems, with users also transitioning from regional helpdesks to national helpdesk. The growth in the national helpdesk was likely matched with a decrease in regional helpdesks (although the ARC WG does not have regional data to support this assumption). Another interesting aspect is that the number of tickets is more closely related to the number of *new* users, rather than to the total number of users.

Experienced users tend to ask fewer questions than beginners. We therefore expect the number of tickets to follow the growth rate rather than the absolute number of users.

The Table 5 below lists the percentage of support tickets by research discipline in 2019 (second column). For comparison purposes, the third column lists the user distribution, while the fourth and fifth columns show the relative CPU and GPU resource usage (as discussed earlier in this chapter).

| Discipline | Support tickets in 2019 | Users by discipline | CPU usage in 2019 | GPU usage in 2020 |
|---|---|---|---|---|
| Biological and life sciences | 19% | 18% | 8% | 12% |
| Engineering | 18% | 19% | 28% | 7% |
| Computer and information science | 12% | 14% | 4% | 41% |
| Chemistry and biochemistry | 11% | 9% | 20% | 24% |
| Physics | 10% | 10% | 20% | 9% |
| Environmental and earth science | 7% | 6% | 6% | 0.4% |
| Medical science | 6% | 8% | 2% | 2.0% |
| Mathematics and statistics | 4% | 5% | 2% | 0.8% |
| Astronomy | 3% | 3% | 6% | 1.4% |
| **Social science** | **1.3%** | **2.2%** | **0.3%** | 0.2% |
| **Humanities** | **1.1%** | **1.2%** | **0.03%** | 0.0% |
| **Psychology** | **1.1.%** | **1.5%** | **0.2%** | 0.5% |
| **Business** | **0.9%** | **1.1%** | **0.2%** | 0.2% |

**Table 5:** Support tickets by research discipline

Looking at disciplines with more than 3% of the tickets submitted, the number of tickets seems to correlate with the number of users per discipline at CCF. The CPU usage varies by discipline, so that biological and life sciences need support relative to their representation but leverage only half of the CPU-years in relative terms. On the other hand, physics and astronomy file a similarly representative number of tickets, but (not surprisingly) consume double the CPU resources compared to the size of the user base. With respect to Computer and Information Sciences, it should be noted that they consume relatively few CPU resources (4% of the total) and much more of the GPU resources (40% of the total), so likely their support requests are mostly related to that resource.

The disciplines in the four categories that have submitted the least number of tickets (social sciences, humanities, psychology, and business) contribute roughly 4.4% of all the support tickets, while representing 6% of the users by discipline. Looking at the corresponding CPU resource usage, these disciplines consume roughly 1% of both the CPU and GPU resources. The researchers in these fields file support tickets at roughly the rate corresponding to their representation in user numbers but consume only a quarter of the CPU resources in relative terms compared to the traditional heavy users. This discrepancy could indicate a need for targeting these disciplines, not only in the terms of getting researchers accounts, but moreover enabling the researchers to leverage ARC and DRI resources effectively for their research, via for example targeted training, support, documentation, and innovative new middleware and gateways to access DRI.

## Training

Training is a very important activity in the CCF federation, including seminars, workshops, summer schools etc. It is critical for adoption, outreach, training of digital workforce, updating skills of researchers (many have learned about ML techniques for the first time) etc. It is also one of the metrics CC reports on to CFI every year.

| Region | Total # of in-person events delivered within region | Total # sites/locations where in-person events were delivered | Total # of attendees at all in-person events | Total # training hours delivered at in-person events | Total # of online events delivered within region |
|---|---|---|---|---|---|
| ACENET | 75 | 7 | 1,298 | 3,522 | 2 |
| Calcul Québec | 44 | 11 | 905 | 4,453 | 2 |
| Compute Ontario totals | 254 | 18 | 10,132 | 30,114 | 85 |
| WestGrid | 85 | 7 | 1,883 | 7,898 | 14 |
| **TOTALS** | **458** | **43** | **14,218** | **45,987** | **103** |

**Table 6:** Training within the CCF consortium

The Table 6 above lists the training within the CCF consortium from April 2019 to March 2020. The rows in the top part show the training given by regions, while the lower part (in italics) lists the training at Compute Ontario given by local affiliates. Training hours in the fifth column are calculated as total training hours received per person, i.e., a one hour training session with ten attendees is counted as ten training hours delivered.

In 2019-20 the CCF consortium delivered total of 46,000 training hours (equaling 5.2 years) to roughly 14,000 attendees in roughly 460 in-person events. Regionally two thirds of the hours were delivered by Compute Ontario, followed by Westgrid at roughly 8000 hours delivered, and then Acenet and Calcul Quebec roughly 4000 hours delivered. Considering the number of in-person events held, Compute Ontario hosted roughly half of all 458 events in 2019-20, followed by Westgrid with 19% of the events, ACENET (16%), and Calcul Quebec (10%). The average number of attendees per event was roughly 31, indicating the demand for and how popular such training has been. Looking at average number of attendees per region we notice that Compute Ontario had most attendees per event (40 per event), while the events in rest of the regions were attended at roughly equal rate of around 20 people per event.

The composition of training has changed over the last decade. At SciNet the 'traditional' ARC (MPI and OpenMP) training hours delivered have stayed roughly the same since 2012, at around 1000 hours. Data sciences (python,R, ML  etc.) has emerged as the largest training category so that in 2018 SciNet delivered roughly 5000 hours of training in data sciences. Another recent emerging training field has been scientific computing, where roughly 2000 hours of training was delivered in 2018.

On the international ARC training front CCF and SciNet have been involved in the International HPC Summer School, which was hosted in Toronto in 2015, and will be hosted by SciNet virtually in July 2021.[187] This event typically has a total of 10 Canadian grad students from across the country participating in the competition each year, and the participants are given training, mentorship etc. from expert instructors, including Canadian HQP.

## 4.4 What are the current strengths of Canada's ARC platform?

The 2017 LCDRI ARC Current State Assessment listed multiple strengths in the Canadian ARC ecosystem, as follows:

- Strong ARC service provision.

- Developed and refreshed ARC infrastructure.

- Improved service delivery to researchers through better coordination.

- A strong and dedicated community of highly qualified personnel who are committed to delivering high quality infrastructure and services to Canada's researchers.

- A strong culture and well-developed centres for innovation.

---

[187] International HPC Summer School https://www.ihpcss.org/index.html (retrieved May 2021).

- Strong track record of adaptability and diversity.

- Stable and well-developed regulatory environment.

Since 2017 the ARC ecosystem has maintained or even improved on these strengths. Moreover, the renewed and centralized long-term funding from ISED for DRI, to be managed by the Alliance should be added as a new strength for the Canadian DRI ecosystem going forward. In the following we will touch on some of the items above, particularly if there have been recent important changes or emphasis is warranted. For a detailed review of these historical strengths, please see Chapter 4.1 of the 2017 LCDRI ARC Paper.

### Strong ARC service provision.

As part of its annual account application renewal process CCF queries its user base on their impressions on CCF resources and services. In 2020 CCF received 10,900 responses to the survey with 85% of users of the ARC platforms being 'satisfied' or 'very satisfied' with the CCF offerings in general. On the flip side only 3% of respondents were either 'dissatisfied' or 'very dissatisfied'. The level of satisfaction has stayed high and static over the last few years, e.g., in 2017 the percentage of satisfied users was the same 85% as in 2020.

Users from all research disciplines seem to be equally satisfied with the CCF resources and services. The variations between disciplines are small, in 2020 total average was 4.32 on the scale of 1 to 5. Differentiating the user satisfaction survey results by position, region, or RAC award status also does not show any large statistically clear variations.

### Developed and refreshed ARC infrastructure.

CCF and its predecessors have a 20+ year history of building, supporting, and delivering high-end ARC systems through its consortia to the Canadian research community. Compute Ontario's recent "Thinking Forward Through the Past: A Brief History of Supercomputing in Canada and its Emerging Future" has a good summary of the developments since the beginning of the 1950's.[188] A testament to the health of the underlying fundamentals is the robust and high user satisfaction across CCF's user base.

In its 2018 budget, the Government of Canada decided to make a major $575.5M investment in Canadian Digital Research Infrastructure. While the bulk of this investment was to fund longer term new DRI organization and what would become the Alliance ($375M), and for CANARIE to support academic networking ($145M), the budget also included $50M to support immediate ARC infrastructure expansion and upgrades.[189] Including the matching by provinces and partners the total investment in Canadian ARC improvements was $94M, allocated to all main host sites: McGill University (total of $28.1M), Simon Fraser University ($39.7M), University of Victoria ($9.6M), University of Toronto ($11M), and University of Waterloo ($5.6M),[190] contributing towards

---

[188] Compute Ontario: Thinking Forward Through the Past:A Brief History of Supercomputing in Canada and its Emerging Future https://computeontario.ca/wp-content/uploads/2019/07/A-Brief-History-of-Supercomputing-in-Canada-and-its-Emerging-Future.pdf  (June 2019).

[189] Government of Canada, Innovation, Science, and Economic Development Canada – Digital Research Infrastructure http://www.ic.gc.ca/eic/site/136.nsf/eng/home (retrieved January 2021).

[190] Government of Canada, Innovation, Science, and Economic Development Canada – Digital Research Infrastructure Questions and Answers: http://www.ic.gc.ca/eic/site/136.nsf/eng/00003.html (retrieved January 2021).

the major system upgrades as detailed in section 4.2. above. Four sites are expected to have capacity online by Winter 2020, and the last one by Winter 2021.

In addition to the ARC upgrade funding, ISED has also expressed interest in supporting any potential short-term immediate ARC storage needs, for the fiscal year 2021-22. The Alliance has formed a Storage Working Group that will present an evidence-based storage upgrade proposal to CFI in early 2021.

**Top500**

There are currently 12 Canadian systems (8 of which are for ARC) listed in the most recent (November 2020) Top500 rankings of world's most powerful supercomputers. Two of these entries are old versions of Cedar, and two are CPU and GPU versions of current Cedar, so strictly speaking Canada has 9 distinct systems on Top500. All recently commissioned CCF non-cloud systems except Graham are on the list: Niagara at #82, Cedar (GPU) at #89, and Beluga at #188. In June 2020 listing Graham was at #500 but dropped from the list as of November 2020. The top Canadian system (Niagara) clocks in ca. 3.6 Petaflops of measured computing performance, while the slowest, Beluga, on the list performs at ca. 2.3 PFlops level. The slowest systems on the November list clock 1.3 PFlops so that Graham's 1.2 PFlops is close but not anymore enough to be on the Top 500. Notably the currently fastest supercomputer in the world, Japan's Fugaku has ca. 442 PFlops performance, i.e. it is ca. 123 times faster than the top Canadian entry.[191]

Going beyond CCF affiliated systems, the Top500 has multiple additional Canadian entries: Shared Services Canada's (SSC) Banting at #128, and Daley at #139. These Cray XC50 supercomputers are hosted by SSC, and primarily used by Environment and Climate Change Canada (ECCC).[192] On the Top500 list there are also four Lenovo based Canadian systems provided by "Cloud Provider", ranked from #333 to #336. The owner/operator of these (seemingly identical) commercial systems is not known publicly.

In the Table 7 below we compare Canada's November 2020 Top500 entries to other G7 countries on an aggregate basis, including both CCF and non-CCF systems. It should be noted that these statistics are based on the Top500 results as is, without trying to parse multiple entries for the same or similar systems etc. As such, for the sake of consistency and comparability, the number of entries for Canada is 12 and not 9 as discussed above.

Looking at the number of entries, Canada is ranked 5th ahead of Italy, and tied with the U.K. Considering aggregate total compute power (Rmax Peta Flops) Canada is last, while The US is the clear leader in absolute compute power, although Japan is not trailing too far behind, by nearly three Canada's. Notably Italy, which has only half the entries compared to Canada, has three times as much aggregate compute power. Comparing individual countries like Canada with E.U. members is nevertheless not straightforward due to pooled resources at E.U. level. The second to last column tries to provide a more representative measure of ARC investment by comparing

---

[191] Top500: November 2020 https://www.top500.org/lists/top500/2020/11/download/TOP500_202011.xlsx (retrieved January 2021).

[192] Shared Services Canada: High Performance Computing https://www.canada.ca/en/shared-services/corporate/data-centre-consolidation/high-performance-computing.html ; and High Performance Computing environment upgraded to support digital government https://www.canada.ca/en/shared-services/campaigns/stories/hpc-upgrade.html (retrieved September 2020).

aggregate compute power to national GDP.[193] Notably Japan takes a clear lead on this metric, with Italy, Germany, France, and The US in the middle of the pack. Canada is second to last in this metric among G7 (at aggregate 16 Rmax TFlops per GDP in USD), indicating that based on our national wealth, our ARC capacity should be at least doubled to keep up with our peers (that are at ca 31 to 39 aggregate Rmax TFlops per GDP in USD).

| Country | Top500 entries | Aggregate Rmax PFlops | GDP (2019, Billion USD) | Aggregate Rmax TFlops per GDP (TFlops per Billion USD) | Ranking within G7 based on Top500 compute power per GDP |
|---|---|---|---|---|---|
| Canada | 12 | 27 | 1736 | 16 | 6 |
| France | 18 | 89 | 2715 | 33 | 3 (tie) |
| Germany | 17 | 131 | 3861 | 34 | 3 (tie) |
| Italy | 6 | 79 | 2003 | 39 | 2 |
| Japan | 34 | 594 | 5081 | 117 | 1 |
| United Kingdom | 12 | 34 | 2829 | 12 | 7 |
| Unites States | 113 | 669 | 21433 | 31 | 3 (tie) |

**Table 7:** Comparison of Top500 rankings and GDP among G7 countries

Per November 2013 Top500 listings, Canada's aggregate compute power per GDP was roughly 1369 GFlops per GDP Billion, i.e., ca. 1.4 TFlops per GDP.[194] In other words in the last seven years Canada (and in general all other countries thanks to advances in and value of HPC technology) has been able to get nearly eleven times more compute power per GDP dollar, while improving its G7 ranking by one step, ascending from #7 to #6.

Centralized service delivery to researchers through better coordination.

As a part of CCF's service modernization, the organization moved to a more national operations and support model including a more consistent and coherent computing and data environment.

---

[193] The Worldbank: GDP (current US$) https://data.worldbank.org/indicator/NY.GDP.MKTP.CD?view=chart (retrieved January 2021).

[194] Top500: Who are the Top500 list countries? https://www.top500.org/news/who-are-the-top500-list-countries/ (June 13 2014, retrieved January 2021).

Service model improvements included e.g., uniform access via centralized credentials (based on LDAP), improved quality of the (centralized and bilingual) documentation, improved data transfer services, centralized approach to more uniform storage offerings (via standardized file system layout and policies), and a centralized application process for accounts and resource allocation. End-users now have a single point of contact for research computing support, while local on-campus support personnel are still available as needed. Centralized support provides multiple benefits e.g., via access to deeper expertise across Canada, access to bilingual support across Canada, and via better distribution of available staff and resources across multiple time zones.[195] Besides the centralized support, at least a third of support tickets for SciNet's Niagara system are handled outside the central CCF support ticketing system.

CCF and affiliates have made improvements in workload portability across platforms: "The new systems will allow Compute Canada users to more easily shift their workload between the different systems, to make optimal use of available resources. This will be facilitated by deploying a single high-performance computing batch system (SLURM), having a common naming scheme for software, modules, and file system mount points, and incorporating mechanisms for data movement with the workload manager."[196] Using the same batch scheduler allows end-users to use the same job submission scripts on different systems with minor modifications. Unfortunately, the users still need to login to individual systems to submit jobs, that is, the individual schedulers are not part of one single centrally accessible scheduling system.

The software stacks on general purpose clusters Graham, Cedar and Béluga are identical, while on the massively parallel Niagara the SW system is split so that both the main CCF SW stack and Niagara specific stacks are available with the latter being the default.[197] The main CCF provided SW stack environment is documented, tracked, and publicized in Github.[198] The SW stack is optimized for the CPU architecture, portable and scalable. It is publicly available for both CCF affiliates and for the general public globally via CERN Virtual Machine File System (CVMFS) technology.[199] In Canada, it is used by various software groups and institutions, as well as by NRC.

In addition to the above-mentioned improvements in identification and authorization service, and software distribution service, the CCF also has improved system status reporting via its centralized resource publishing service that provides current information about available resources.

[195] NDRIO Position Paper submission by Maxime Boissonneault: https://engagedri.ca/successes-and-shortfalls-of-the-current-canadian-arc-platform-and-ideas-to-improve-it-further (retrieved January 2021).

[196] Compute Canada: Usage and Capabilities https://www.computecanada.ca/techrenewal/usage-and-capabilities/ (retrieved September 2020).

[197] Compute Canada: Available Software https://docs.computecanada.ca/wiki/Available_software (retrieved September 2020).

[198] Github: Compute Canada Software Management https://github.com/ComputeCanada/software-stack/blob/master/doc/INDEX.md (retrieved January 2021).

[199] Compute Canada: Accessing CVMFS https://docs.computecanada.ca/wiki/Accessing_CVMFS (retrieved January 2021).

A strong and dedicated community of highly qualified personnel (HQP) who are committed to delivering high quality infrastructure and services to Canada's researchers.

The CCF network consists of roughly 250 highly qualified personnel (HQP) running the CCF operations and sites across Canada. These people provide a variety of critical services related to ARC systems administration, procurement, maintenance, networking, operations, management, planning, funding, support, research software development, data management, training, accounts, and allocations management, communications, and outreach. The ARC systems almost by definition are leading edge highly complex systems in nearly all aspects of their configuration, software and hardware stacks, operations, and use, requiring senior level expertise that takes several years of specialization to master. Maintaining the skillsets and retaining the CCF HQP is of critical importance to Canadian DRI ecosystem.

Some of the staff need to be local to the host sites for quick physical access to the hardware and networking in the datacenters, while most other functions can be and should be provided remotely. Many ARC systems are by design meant to be accessed remotely, so that the only limit for many operations is the bandwidth and latency of the network connection, and other external factors like potential additional security requirements, human resources needs and requirements, or team building preferences. In any case, this built-in remote work enabling nature of ARC allows Canadian DRI providers to potentially leverage a wide pool of talent regardless of their geographical location. Off-site positioning of the staff is even preferred for e.g., support services where support personnel can be located at the local university near the user community.

| | Westgrid | Compute Ontario | Calcul Quebec | Acenet | Total |
|---|---|---|---|---|---|
| Budgeted FTEs | 62,1 | 67 | 53 | 19 | 201,1 |
| Number of employees | 98 | 76 | 55 | 20 | 249,0 |
| Number of institutions | 7 | 14 | 8 | 6 | 35 |

*Employees of Compute Canada Central are not included*

*14 institutions in Ontario include 2 research hospitals*

**Table 8:** Staff distribution – ARC budget 2020-2021

The Table 8 discusses the budgeted HQP staffing levels within CCF in fiscal year (FY) 2020-21 including regional distribution of the personnel, excluding central CC employees. Canada's ARC community is composed of ca. 250 highly qualified personnel across the country. The formal full-time equivalent (FTE) resources in FY20-21 were ca 200.

The 200 FTE strong HQP workforce translates to roughly 6 FTE per contributing institution, and to roughly 1:80 ratio between number of HQP FTEs and CCF registered users. Interestingly this ratio is more favorable at e.g., Texas Advanced Computing Center (TACC), the host site for

Frontera, one of the world's fastest supercomputers, that serves 'several thousand users'[200] with a staffing of roughly 190 people[201], corresponding to a staff-to-user ratio in the range of 1:16 to 1:55, i.e., from five times more favorable to nearly two times more favorable at TACC. On the upside, it is remarkable how the CCF HQP team is consistently delivering 85% user satisfaction with much less resources per user than e.g., at TACC.

Regionally, Westgrid and Compute Ontario contribute the largest number of FTEs, at 62 and 67 respectively while Calcul Québec contributed 53 and Acenet 19 FTEs. Comparing these FTE resources to number of users from the same regions, at Westgid the ratio is 82 registered users to a HQP FTE, at Compute Ontario the ratio is 84, at Calcul Québec 87, and Acenet ca 55. In other words, the ratio of HQP resources provided per registered user is roughly the same in the three largest regions while the situation is better at Acenet. It should be noted that this comparison does not consider the fact that not all HQP resources are dedicated to user facing services such as support and training. For example, the three larger CCF regional affiliates need substantial systems administrator staffing to run their main hosting sites (which Acenet does not currently have). Furthermore, the support and training within CCF is not locked-in regionally so that users can get support from HQP that are not local or located in their home region.

---

[200] TACC: Texas boosts U.S. science with fastest academic supercomputer in the world https://www.tacc.utexas.edu/-/texas-boosts-u-s-science-with-fastest-academic-supercomputer-in-the-world (retrieved December 2020).

[201] TACC: Staff Directory https://www.tacc.utexas.edu/about/directory (retrieved December 2020).
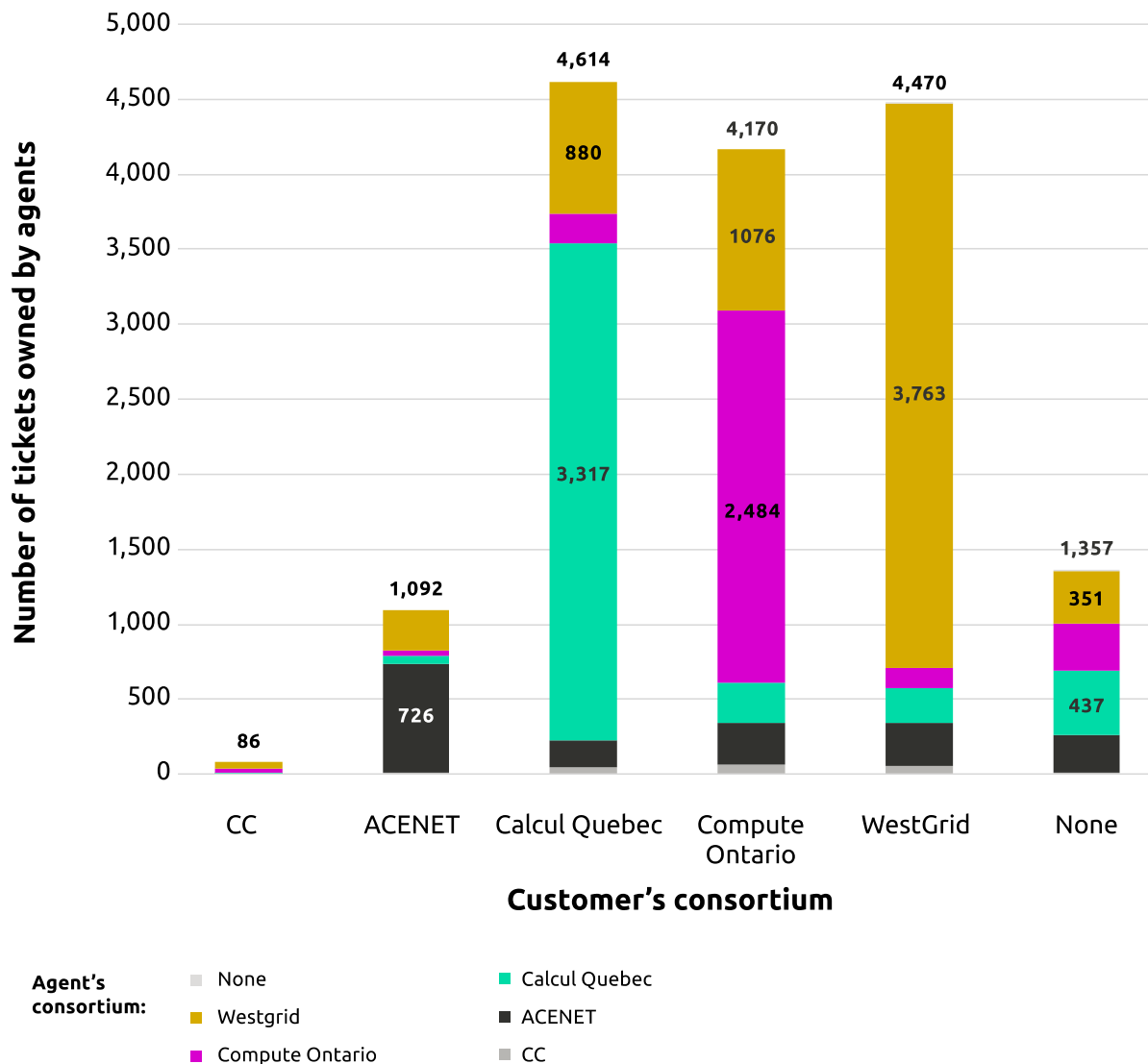
**Figure 11:** Support ticket count by customer and agent consortium

The Figure 11 above shows the number of support tickets that originated from customers located in different CCF regional consortia. The data is for the calendar year 2020. The total number of tickets from different regions is indicated at the top of each bar. Users in three regions, Calcul Quebec, Westgrid, and Compute Ontario submitted a similar number of tickets, ca. 4600, 4500, and 4200 tickets, respectively. Users affiliated with Acenet region submitted roughly 1100 tickets. The colored bars inside every main bar indicate the consortia that handled the ticket. For example, Calcul Quebec support staff handled ca. 3300 tickets that originated from users affiliated with Calcul Quebec, and similarly Compute Ontario handled ca. 2500 support tickets that originated from users affiliated with Compute Ontario. The general trend for support tickets submitted by users in all regions is that the tickets are handled primarily by the corresponding regional CCF affiliate. It should be kept in mind that considering only ticket counts can be misleading. For example, this type of analysis does not account for all the people involved in resolution of a given ticket, nor does it account for the effort required for tickets (which can range from a minute to

hours) or the fact that experts for specific issues may not be distributed evenly across regions. This current analysis also does not include the consideration of the topic being addressed in the ticket.

A closer look in the data reveals interesting trends. Keeping in mind that BC, Ontario, and Quebec host main CCF systems and need to have corresponding administrative staff so that these regions will have in relative terms less staffing available for research computing support ticket handling. The total volume of tickets from users in Calcul Quebec and Compute Ontario regions is more than the number of tickets these regions handled respectively. The ca. 4600 registered Calcul Quebec users submitted ca. 4600 tickets, while Calcul Quebec staff (ca. 53 FTEs) handled ca. 4100 tickets (as indicated by adding up all green areas). The ca. 5600 registered Compute Ontario users submitted ca. 4200 tickets, while Compute Ontario staff (ca. 67 FTEs) handled ca.3100 tickets (as indicated by adding up all purple areas).  In other words, if Calcul Quebec and Compute Ontario support staff would handle all tickets originating from their regions additional local tickets would remain to be handled by other regions, that is, 500 and 1100 tickets, respectively.

The additional tickets from Calcul Quebec and Compute Ontario need to be picked up by other regions, namely Acenet and Westgrid. In the case of Acenet the local staff (ca. 19 FTEs) handles more than half of the locally originated tickets (ca. 700 of 1100 tickets) submitted by the ca.1050 local registered users, and then supports tickets submitted by users in different regions in roughly equal share per region, adding up the total of ca 1400 tickets handled by Acenet staff (as indicated by adding up all black areas). In other words, Acenet support staff handled roughly 300 more tickets than were submitted by all users in Acenet region. In the case of Westgrid the local staff (ca. 62 FTEs) handles more than 80% of the locally originated tickets (ca. 3800 of the total of ca. 4500 tickets) submitted by the ca. 5050 local registered users, and then supports tickets submitted by users in other regions substantially, adding up the total of ca 6300 tickets handled by Westgrid staff (as indicated by adding up all brown areas). In other words, Westgrid support staff handled roughly 1800 more tickets than were submitted by all users from the Westgrid region. Westgrid staff handled ca. 1100 tickets originating from Compute Ontario region, and ca. 900 tickets originating from Calcul Quebec region.
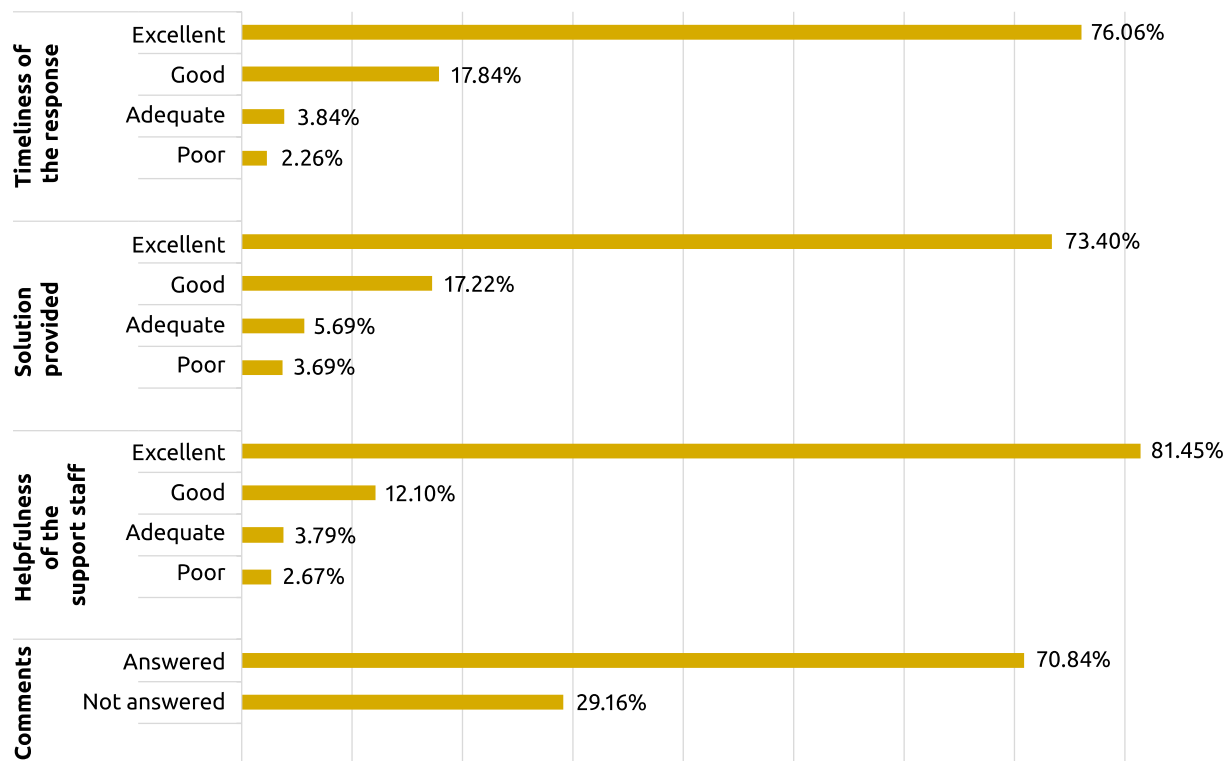
**Figure 12:** CCF support ticket system post-ticket satisfaction survey responses

Recent CCF support ticket system post-ticket satisfaction survey responses speak to the high quality of service provided by the HQP people within CCF. The above Figure 12 shows the distribution of responses collected from September 2019 (when such data collection started at CCF) until early October 2020. Regarding timeliness of response, 94% of respondents rated timeliness good, or excellent. Similarly on the topic of 'solution provided', 91% were happy with the solution. The same trend persisted on the question of staff helpfulness, where 94% of respondents rated the staff good or excellent. These results indicate a high level of satisfaction with the CCF support team. The survey did not sample the users of SciNet's 'Niagara' supercomputer sufficiently, since many support tickets are still being sent to SciNet's local helpdesk. For more reliable analysis and conclusions, a more comprehensive data collection and survey would be required.

### A strong culture and well-developed centres for innovation.

The main CCF host sites are operated by very highly qualified personnel providing solutions and infrastructure that is both available nationally, but also tailored as needed to local university and regional level needs. These centers attract talent and act as training grounds for new HQP generation, while also attracting academics and researchers who appreciate local access to support and resources. Collaboration between host sites has led to national level service delivery (in most cases), leveraging individual innovation to improve the whole. As indicated by the Top500 rankings, the main CCF compute systems are truly world class.

### Strong track record of adaptability and diversity.

The main CCF host sites must satisfy a myriad of different user requirements and needs while maintaining a steady baseline of ARC compute cycles and storage services for volume users. The CCF coalition continues to rise to this challenge by providing CPU, and GPU computing services, various cloud-based offerings for adaptability. This desire to serve the end-user community needs to be balanced with the current funding restrictions and possibilities related to e.g., contributed systems, tape storage, ability to provide physically separated secure storage, targeted programs for the EDI communities, and lack of flexibility regarding separation between capital and operational funding.

### Stable and well-developed regulatory environment.

Canada is a stable progressive safe democracy with well-functioning government branches, excellent public education, and universal healthcare systems, all providing the necessary framework and stability, directly or indirectly for advanced digital research operations. Investors and funders, corporations, academia, HQP and researchers can rely on data and science driven policy and investment. Such a secure environment fosters innovation as people can trust that services and resources will not vanish overnight. For example, the reorganizing of funding and operations of Canada's DRI ecosystem was announced as a part of the 2018 federal budget so that the Alliance will take over the daily operations in April 2022, leaving all parties enough time to plan.

Regulatory frameworks are necessary to protect people's privacy and security, while open access to science and data can unleash vast new knowledge, research, and innovation. Balancing these two competing aspects for the benefit of society is one of the big challenges for Canadian DRI community in the coming years.

### Renewed funding commitment from the government for DRI

The Canadian government via the Ministry of Innovation, Science and Economic Development of Canada (ISED) clearly sees the value of DRI for Canadian society, as proven by the $572.5M 2018 budget commitment as discussed above. On the Alliance side this translates to a total of $375M federal funding until March 2024, providing important (relatively) long term continuity to the DRI funding. Moreover, this funding restructuring also balances the focus away from only ARC to also RS and RDM, covering the three key pillars of a modern DRI ecosystem under one operation.

## 4.5 What are the current challenges and opportunities facing Canada's ARC platform?

The 2017 LCDRI ARC Current State Assessment listed multiple challenges in the Canadian ARC ecosystem, as follows:

- Insufficient ARC supply to meet current and future demand.

- Lack of sustained and predictable funding.

- National platform development and the current funding model.

- Coordinated national strategic and operational planning

- Attraction and retention of HQP.

- International collaboration and competitiveness.

- Coordination of federal science investment and the provision of ARC.

- Keeping pace with technological and market changes.

- Leveraging cross-sectoral ARC resources.

- Researcher awareness and adoption of ARC.

- Environmental impact.

- Securing the national platform.

Since 2017 many of these challenges still stand unresolved, even though there has been substantial material investment to address resource constraint issues. Moreover, ISED has recently mandated the Alliance to address and mitigate many of these issues. In the following we will touch on some of the items above, particularly if there has been recent important changes or emphasis is warranted. For a detailed review of these historical challenges, please see Chapter 4.2 of the 2017 LCDRI ARC current state document.
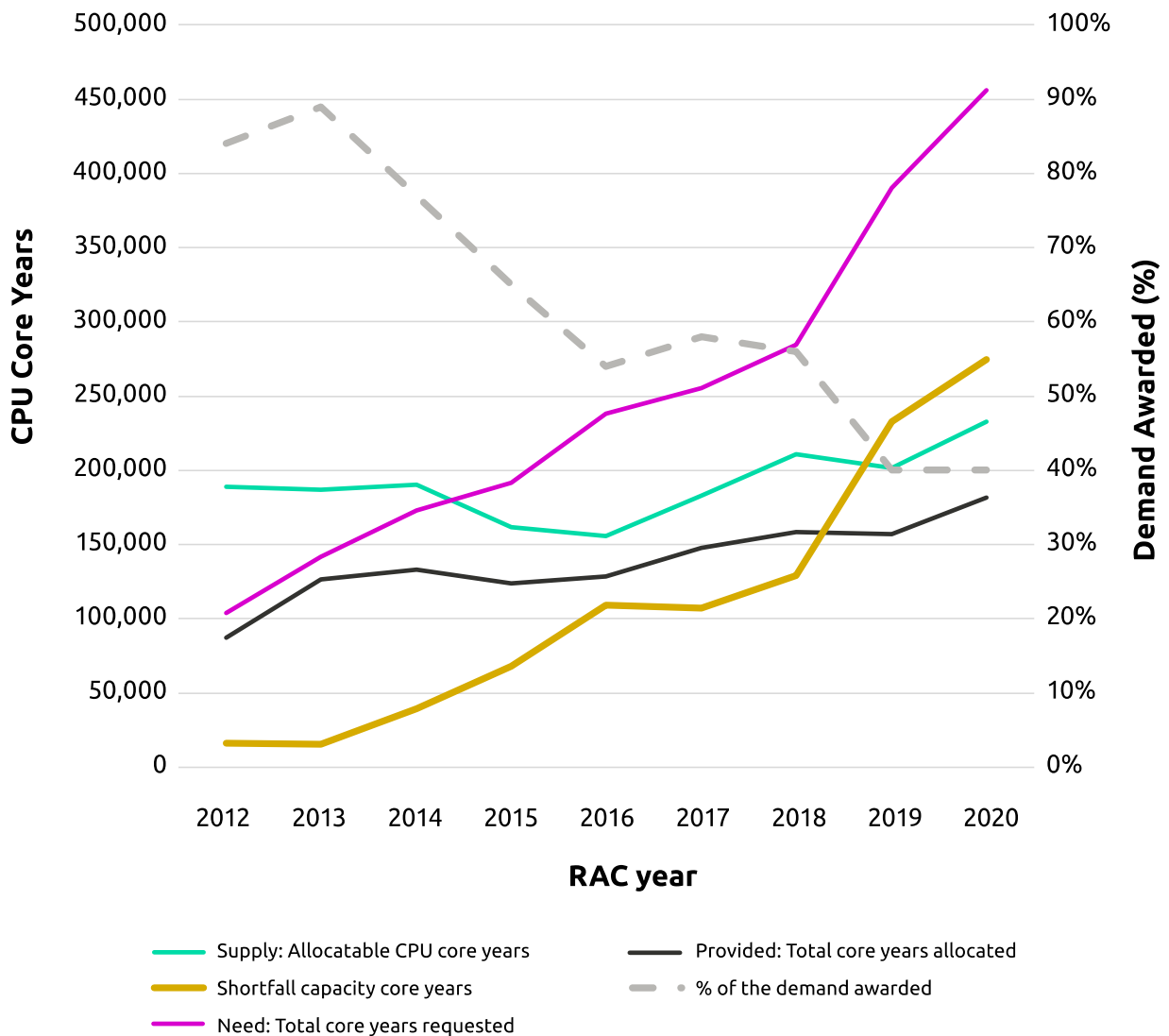
**CPU compute resource supply and demand**



**Figure 13:** CCF historical CPU allocation and demand

The Figure 13 above compares CPU compute resource supply and demand in the CCF consortium from 2012 to 2020. The supply is based on actual available CPU capacity in CCF systems, and the demand is based on the submitted resource allocation requests in CCF's annual Resource Allocations Competition (RAC). The horizontal axis corresponds to RAC allocation years, while the units on the vertical axis correspond to an equivalent of running a single CPU core at 100% capacity for a year (i.e., CPU core years). The green line is the total annual raw compute capacity available in the main CCF systems. The total available capacity has fluctuated within a relatively narrow range between ca 155k and 232k CPU core years in the last eight years.

In 2014 there was significant capacity that went offline due to age and to CCF decisions made about viability supporting older equipment. This is indicated by the dip in the capacity (green line) after 2014 that did not rise significantly above 2012 levels until 2020 when ISED expansion was deployed. Notably the CPU core year metric does not consider the actual compute power of each cycle; that is, the increase in processor compute capability thanks to architecture developments.

The purple line is the total capacity requested in the RAC competitions. In the last eight years the demand has grown from ca 100k CPU years to 450k CPU years. The growth in demand is very rapid and semi-linear but does not seem to be exponential. In CAGR terms the growth in the demand for CPU computing cycles was ca. 21% per year.

The black line indicates the actual RAC allocated capacity. It should be noted that the total capacity is provided to end-users via the RAC application process (deliberately capped at about 80% of the capacity), while the remaining capacity (roughly 20%) is available for use by any user on as needed basis without a need for a formal application. In 2020 the available total resource of ca. 232 000 CPU years is distributed between the RAC (ca. 182 000 CPU years), and unallocated/non-competitive use. The unallocated 50 000 CPU year capacity is utilized by the rapid access RAS users, 'opportunistic' use. Just as an indication of the scale, in 2020 this total unallocated capacity across the national sites corresponds to ca. 60% of Niagara's total annual compute capacity. It should be noted that overall utilization of the systems is high - roughly 90% of all theoretically available cycles are used (the remainder is accounted for by downtime of individual components of systems, planned and unplanned outages and the fact that job scheduling on shared systems is never perfect).

The allocated capacity (black line) has closely followed the available capacity (green line), leaving the above-mentioned ca. 20% margin for the rapid access services. Comparing the supply (green line) and demand (purple line) we see that the CPU computing capacity is falling behind with the rapidly growing need and the ARC infrastructure development has not kept up with the demand.

The thick brown solid line shows the unmet demand in absolute terms. In 2020 this is roughly 274k CPU years. This unmet demand corresponds to roughly 3.4 times the capacity of the Niagara supercomputer. The grey dashed thick line highlights the scale of the problem by showing what percentage of the computing demand was actually allocated. Historically this can be seen to decrease from roughly 80% of the demand being satisfied in 2012 to only 40% in 2020.

The infrastructure modernization in 2016-18 was able to stabilize the shortfall temporarily (see the thick brown line) but was not able to reduce the shortfall in absolute or relative terms. In 2019 the decommissioning of MP2 again reduced the available compute capacity, while the more recent additions were able to moderately increase the total available capacity by some 30 000 CPU years in 2020. All the while as the rapidly growing demand for resources further exacerbated the absolute gap between supply and demand. In summary, in the last decade the CPU compute capacity shortfall has increased significantly in both absolute and relative terms.

It should be noted that assessing the ARC supply and demand situation based on looking at user requests and allocation as in the RAC above can provide an incomplete or misleading picture of the underlying conditions and situation. ARC administrators are aware of situations where researchers might request too many resources leaving allocations underutilized. However, the CCF does not currently have the staffing commitment required to help researchers in these situations and properly assess the efficiency of codes and optimize them for better resource

utilization. This is an avenue that the Alliance could consider investing into to ensure that DRI is used to its best capacity.

It should be kept in mind that ARC resources by their nature are always in short supply due to the constantly growing number of disciplines leveraging DRI, increasing resolution of experimental or observatory instruments, and the need for higher resolution and accuracy simulations, often adjusted to the available compute allocation. If the available compute resource increases, the researchers quickly switch to leveraging that capacity for new science.

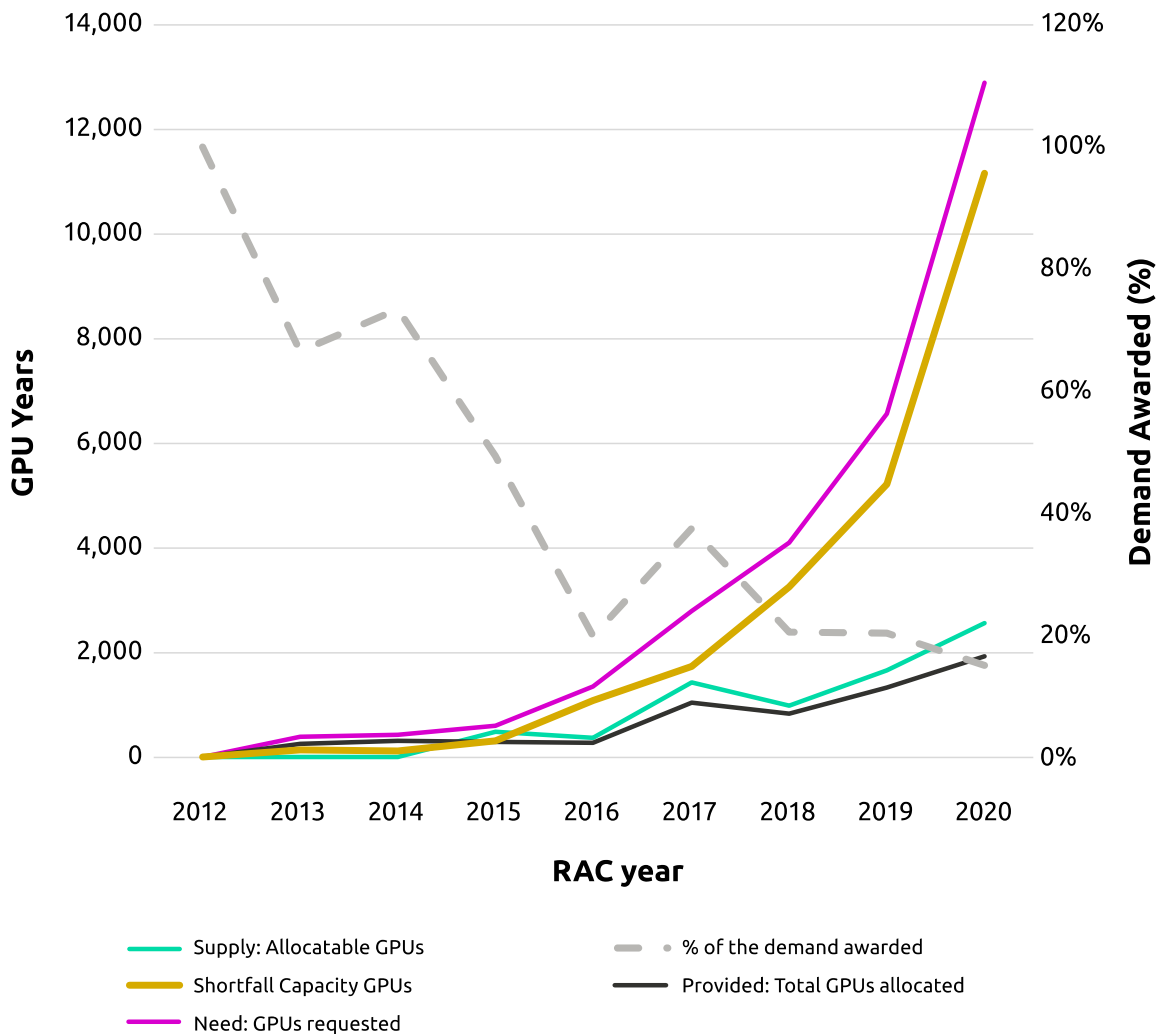**GPU accelerator supply and demand**



**Figure 14:** CCF historical GPU allocation and demand

The Figure 14 above compares GPU accelerator supply and demand in the CCF consortium from 2012 to 2020. The supply is based on actual available GPU capacity in CCF systems, and the demand is based on submitted resource allocation requests in CC's annual Resource Allocations

Competition (RAC).[202] The horizontal axis corresponds to RAC allocation years, while the units on the vertical axis correspond to an equivalent of running a single GPU accelerator at 100% capacity for a year (i.e., GPU years). The green line is the total annual allocatable capacity while the purple line is the RAC requested capacity. The black line is the capacity allocated to end-users via the RAC mechanism.

The demand for GPU computing resources has grown strongly between 2012 and 2020. In 2012 the requests were minimal, at 10 GPU years, while in 2020 RAC round the total request was nearly 13 000 GPU years. The growth has been non-linear, growing exponentially year over year. In CAGR terms the growth was ca 67% since 2017 (when the demand was ca 2800 GPU years), indicating the very rapidly increasing demand for this resource. Such demand fits the global trends, e.g., the current Top500 list has a record 149 systems with accelerators (146 of those using Nvidia GPUs) and 6 of the top10 use GPUs. - including #2 Summit with over 24,000.

GPU computing capability is falling behind with the exponentially growing needs. On one hand this is positive showing significant increasing interest for accelerator technologies, while on the other hand the ARC infrastructure has not kept up with the demand. The infrastructure modernization in 2016-17 was able to catch up and keep up with the demand temporarily, but since then the substantially increased desire for adopting these accelerators has clearly outpaced the supply.

In absolute terms the unmet GPU capacity in 2020 is roughly 11 100 GPU years, as shown by the solid brown thick line in the Figure above. This gap between GPU supply and demand equals roughly eight current Cedar supercomputers. Or, to consider the scale from a different point of view, considering modern NVIDIA V100 Volta GPU cards (list price of ca. $9500 CAD),[203] the total cost for purchasing ca. 11 000 accelerator cards alone would be roughly $100M. The actual cost of catching up to the 2020 unmet GPU need would be even higher once the cost of thousands of ARC servers, and other supporting infrastructure, is included.

In relative terms the demand that has been fulfilled (the dashed thick grey line) has dropped from ca 100% in 2012 to roughly 20% in 2020. Notably the relative shortfall has stabilized at this level in the last three years thanks to infrastructure updates, indicating that the supply and demand are following roughly the same growth trajectories, but at different absolute levels as discussed above.

Interestingly, similar trends could be seen at HPC centers worldwide in the ascendancy of CPU based supercomputing decades ago. Usage of GPUs for HPC started to become a serious "trend" around 2010 but using them was difficult due to complexity of programming even with earlier scientific-programming-oriented GPU programming paradigms like CUDA. Molecular dynamics (MD) codes were some of the earliest to be ported and uptake in that community was particularly strong. Emergence of AI in the last decade has more recently exploded the demand for GPU resources. With such complexities - new technology, new compute methods and paradigms – there is a possibility that the demand as measured by RAC proposals may be inflated even more

---

[202] Compute Canada: 2020 Resource Allocations Competition Results
https://www.computecanada.ca/research-portal/accessing-resources/resource-allocation-competitions/rac-2020-results/ (retrieved August 2020).

[203] CDW online store: NVIDIA Tesla V100 - GPU computing processor - Tesla V100 - 16 GB
https://www.cdw.ca/product/NVIDIA-Tesla-V100-GPU-computing-processor-Tesla-V100-16-GB/4939179 (retrieved September 2020).

so than for CPUs. With rapid emergence and adoption of GPUs researchers do not have good baselines for code performance, or for the number of training runs required or the human-time needed to manage and interpret results. This results in overestimation of the GPU computing need. Keeping in mind also that e.g., porting the whole scientific problem to GPUs can be very difficult or impossible. As a result, there is potentially a built-in inefficiency in GPU resource utilization per Amdahl's Law, when parts of the execution are CPU bound only while keeping GPUs reserved.

**Storage allocation supply and demand**



**Figure 15:** CCF historical storage allocation and demand

The Figure 15 above shows the historical storage supply and demand at CCF. The aggregate available supply of the various available storage types is shown by the solid green line. In 2015 the total capacity was ca. 15 PB, increasing nearly tenfold to ca. 143 PB in 2020 with substantial annual variations as old systems have been retired and new ones brought online. This total storage is distributed among functionally different storage types including Project, dCache, Cloud,

and Nearline storage systems.  Project (57 PB supply in 2020)[204] is the large main storage for active research data and files, dCache (15 PB) is object file storage system for large datasets (in particular in high-energy physics), Cloud (4 PB) is for cloud instances, and Nearline (68 PB) is a disk-tape hybrid filesystem for less active data.

On the demand side the purple line indicates the historical aggregate storage demand, while the brown line indicates the storage provided. The demand has increased roughly five-fold from ca. 21 PB in 2015 to ca. 110 PB in 2020. Notably in 2020 the total storage capacity (green line) was ca. 34 PB larger than the total request (purple line). Much of the additional head room in the storage infrastructure is needed for efficient operation of the system. Since the storage system capacity has increased even more over this time period, the supply has been able to keep up with the demand so that the provided storage has grown at the same pace as the demand. In absolute terms the unmet demand in 2020 was ca. 9 PB. The ca. five-fold growth in demand has been roughly linear over the last five years, corresponding to ca. 39% CAGR.

The black dashed line indicates the awarded/allocated demand as a percentage of the annual storage requests. This fulfilled demand has ranged from 72% in 2011 to 91% in 2020, while dipping to 64% in 2016. Thanks to increases in supply, the storage system on aggregate has been able to keep up with the demand.

Notably per CCF's data retention policies Nearline is not an archival or backup file system and is available only for active projects.[205]  Archival or long-term storage is commonly understood as multi-year (i.e., 5, 10, 20 years) storage, requiring predictable long-term funding. Per its mandate CCF does not provide this kind of storage, even though the enterprise tape systems used for Nearline storage at CCF can from technical point of view support such needs. CCF also has the required HQP expertise for providing such storage thanks to running tape systems for years if funding and policies would be in place. The Alliance has setup a dedicated Storage Working Group in late 2020 with the mandate to holistically consider both short- and long-term storage needs in the Canadian DRI ecosystem.

**Projected needs in astronomy and astrophysics**

As discussed above, astronomy and astrophysics are major users of DRI resources, both in Canada and globally. These disciplines also have very mature and longstanding global DRI ecosystems e.g., for processing massive amounts of observational instrument data, and for running astronomical simulations. As new instruments are launched and resolutions are increasing, the needs of these 'traditional ARC disciplines' are growing rapidly as well. In December 2020 Canadian Astronomical Society (CASCA) published its LRP2020 report, the Long-Range Plan for Canadian astronomy and astrophysics for the period from 2020 to 2030.[206] The estimated general need in 2025 is 100 PF years of CPU capacity, 100 PF years of GPU capacity, and 75 PB of online storage. Comparing these requirements to current capacity across

---

[204] Compute Canada: 2020 Resource Allocations Competition Results https://www.computecanada.ca/research-portal/accessing-resources/resource-allocation-competitions/rac-2020-results/ (retrieved November 2020).

[205] Compute Canada documentation:  Technical glossary for the resource allocation competitions https://docs.computecanada.ca/wiki/Technical_glossary_for_the_resource_allocation_competitions (retrieved November 2020).

[206] Canadian Astronomical Society – LRP2020 report https://casca.ca/wp-content/uploads/2020/12/LRP2020_December2020-1.pdf  (December 2020).

all CCF installations, according to CASCA they correspond to roughly ten times of the current total capacity for CPUs, and 25 times of the current capacity for GPUs. The storage needs in 2025 corresponds to the total volume of all current project storage in CCF systems. In addition to general supercomputing needs, the Square Kilometer Array (SKA1) project will need additional 10 PF years of CPU compute capacity, and close to 900 PB of storage. The general computing needs are planned to be funded through the Alliance while the SKA1 computing needs will be funded through the SKA project. Clearly these projected supercomputing and DRI needs will require special and dedicated consideration and funding in the next decade.

## Lack of sustained and predictable funding

Lack of sustained and predictable funding continues to be an issue in the Canadian DRI landscape. This has had major effects on long-term planning for the DRI ecosystem. Lack of clarity on the timing or size of the next round of funding has prevented effective planning for growth, and new technologies and capabilities – resulting in a cycle of boom-and-bust. The substantial funding commitment in the 2018 budget in connection with the formation of the Alliance will provide important predictability and continuity until March 2024. As discussed in the 2017 LCDRI ARC report the issues with funding continuity can negatively affect achieving system efficiencies, can hinder researchers' ability to plan for long-term multi-year research projects, and can negatively affect the retention of HQP due to short term contracts.[1] The misalignment of capital and operational funding due to funding rules and policies is also a major issue. The separate funding mechanisms (and time periods) for capital and operational expenses not only affect the traditional datacenter-based ARC operations, but also hinder the adoption of modern cloud technologies due to 'artificially' limited funding options. The need for establishing a predictable and sustainable funding for ARC was the number one recommendation by Hyperion Research in their report for Compute Ontario in November 2019.[207]

Figure 16 illustrates the highly complex flow of funds in the current ARC funding model. Due to legacy reasons, the CFI funding for ARC infrastructure and operations is funneled through one university, Western University, and not through the central authority, i.e., Compute Canada. Additionally, the main site funding ultimately flows to individual universities hosting the datacenter, and not through the regional CCF affiliate organization. The more recent $50M ISED ARC Expansion injection of funds by-passed Western University and was given directly to main host sites. Without going into further details, and not discussing the operational funding, it is clear that the funding and flow of funds for ARC is unnecessarily complicated in Canada. The formation of

---

[207] Hyperion Research: Summary Report of Phase 1 Study to Support Compute Ontario ARC Planning https://computeontario.ca/wp-content/uploads/2019/11/Hyperion_Summary.pdf (November 2019).
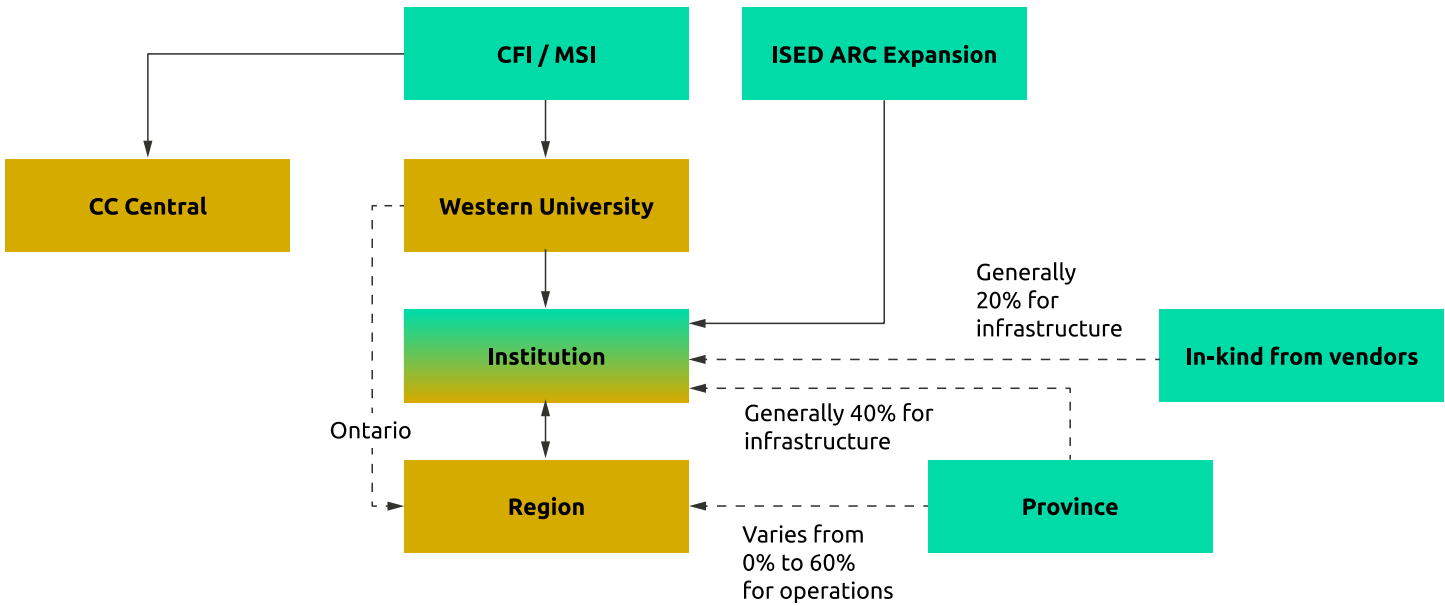
**FUNDING FLOW – ARC**



**Figure 16:** ARC funding flow

the Alliance and recognition from ISED and CFI regarding these complexities is a welcome improvement going forward.

Notably, since the funding often comes in sporadic fashion, addressing e.g., an immediate need every 4-5 years, it results in large purchases with multi-year gaps in between. Since the computing technologies advance on an on-going basis any purchase that can be delayed will result in more value for money. As discussed earlier in this document in the context of Top500, just over the last 8 years roughly the same relative investment (in terms of GDP) has yielded over ten times more raw compute power in Canada.

National platform development and the current funding model

ARC funding in Canada can in general be categorized in three buckets, first for capital funding to procure hardware and supporting software and warranties, second for operational funding to operate the hardware including utility, rental and insurance costs and systems administrator staffing, and the third for operational funding for supporting the users including the related staffing costs. The first stream to fund infrastructure comes in waves with the possibility of multi-year gaps and inherent uncertainty about long-term funding. The second stream is provided on an annual basis through CFI's MSI funding channel and is tied to operating the capital infrastructure on the host sites on the first stream. The third stream is independent from the second stream and is tied to the home institutions of the researchers leveraging ARC resources across Canada. Currently the funding matching formula is in general 40/60, i.e. for 40% of federal funding there needs to be a 60% matching from universities, provinces, vendor discounts etc. In the capital funding stream, the 40/40/20 ratio in practice translates to roughly 50/50 thanks to generous vendor 'CFI discounts' (i.e., the 20% vendor discount is nearly given and in practice the federal and remaining

contributions even out), while in the second and third streams the 40/60 ratio holds more strictly. Notably the recent $50M ARC Expansion Program capital funding injected funding to the first stream, capital funding, but did not include corresponding increases in operational funding streams.

As discussed in the 2017 LCDRI report, the 40% CFI, 60% matching funding model continues to be problematic for universities and regions in the light of trying to achieve a fair distribution of capital and operational costs. Institutions that host main CCF infrastructure might have committed to funding operations and on-going costs beyond their fair local or regional use. While at the same time some universities and regions are using the same resources at no cost, or are only contributing to infrastructure costs, but not on operational costs, or vice versa. Only a small portion of universities and colleges (35) contribute to fund ARC costs. Some regions provide less funding or support to CCF, while in other cases the funding match is available. Funding can also be dependent on other emerging budget priorities and thus due to its sporadic nature can not be (fully) relied on, keeping in mind for example the fiscal consequences of current Covid-19 pandemic to provincial budgets.

The current funding model and award processes also by their nature do not induce collaboration between regions and host sites since it is for example difficult and unlikely (but not impossible) to get provinces to fund things outside their borders. Regarding the systems on main host sites, the sites own and operate these systems beyond just hosting so that there are multiple drivers with local systems administrators operating the local systems per local needs and expertise. On the user support side, the operations are more tied to the researchers in the local institutions, and whether the institutions will spend time for national initiatives or not varies widely. The inherent lack of incentive for collaboration makes it more difficult for the federal government and the funding agencies to coordinate, homogenize, and optimize the national level platform delivery and service model to cohesive federal DRI infrastructure.

Due to the highly complex flow of funding in the current funding model, the roles and responsibilities of various local, regional, and national level entities are not clear, or easily enforced. A more straightforward operational model would have a tighter coupling of funding and operational and management responsibilities.

The Alliance should in particular look to the EU for how DRI for underrepresented fields e.g., social sciences and humanities etc. are funded and operated in a federated model. There are multiple interesting organizations participating in funding and operations of DRI for these disciplines, e.g. EGI Federation, an international e-Infrastructure providing advanced computing and data analytics services[208]; Parthenos Virtual Research Environment (VRE), an online environment for Humanities integrating cloud storage with services and tools for the research data lifecycle[209]; ARIADNEplus, offering cloud-based VREs for data-based archaeological research[210]; DARIAH, a pan-European infrastructure for arts & humanities scholars[211]; and IPERION HS, a

---

[208] EGI: About https://www.egi.eu/about/ (retrieved February 2021).

[209] Parthenos: About VRE https://parthenos.d4science.org/web/parthenos_vre/about-parthenos-vre (retrieved February 2021).

[210] ARIADNEplus: About https://ariadne-infrastructure.eu/about-ariadne/ (retrieved February 2021).

[211] DARIAH-EU: DARIAH in a nutshell https://www.dariah.eu/about/dariah-in-nutshell/ (retrieved February 2021).

European Integrated Research Infrastructure Platform for Heritage Science[212]. ESFRI, the European Strategy Forum on Research Infrastructures is a major research funding agency in Europe, funding for example the above mentioned DARIAH, and the PRACE partnership for ARC.[213]

In its 2018 major funding announcement for Canadian DRI ecosystem the Government of Canada recognized some of the above-mentioned issues. The new funding for DRI is longer-term, for 5 years until March 2024, while a new entity, the Alliance, was later formed to centrally plan and allocate most of the funds. This continuity allows the Alliance to plan purchases ahead of time rather than be subject to cycles. Infrastructure purchases can instead be distributed over several years, in particular over the last three Fiscal Years (FY) 2021-22, 2022-23, and 2023-24. Part of the Alliance's mandate is also to design the national DRI platforms holistically, considering ARC, research software, and research data management needs and national level efficiencies. Another major change is in the matching funds requirement so that going forward ISED/Alliance will be the majority funder with (in general) 60% of the capital and operations costs while the remaining 40% would come from matching contributions.[214] The details of the matching ratio are currently being discussed and negotiated by the key parties.

### Coordinated national strategic and operational planning

The current ARC infrastructure, and moreover the past mandate for CCF operations and funding has emphasized traditional ARC needs, models of operation, user disciplines, and audiences. In practice the ARC delivery systems have focused on high-performance and massively parallel systems and clusters, and active runtime or project storage systems, with some cloud systems coming online more recently. This approach has to a degree missed the needs of wider audiences, and disciplines, e.g. humanities and social sciences, and needs of more diverse audiences. These systems have also relied on traditional technologies for accessing and using the systems, e.g., using the command line instead of graphical user-interface-based solutions, or using traditional batch schedulers instead of cloud based solutions that often hide the scheduler from the end-user. Considering interactive jobs (graphical user interfaces are by definition interactive) vs batch jobs, the primary trade-off is whether to keep resources idle (wasted) vs usable at once (interactive). The focus has mostly been on not keeping resources idle. Moreover, the past and to a degree current ARC systems and service providers do not include RDM, RS, and long-term storage needs in their service delivery, and planning. It should be emphasized that often this situation is not due to the lack of vision or recognition by the ARC providers, rather than due to lack of available (in particular FTE) funding and restrictions within the funding mandates.

A more coordinated and centralized approach in the strategic and operational planning for DRI in Canada is needed. This approach needs to be national in scope for e.g., increased synergies and efficiencies (e.g., in better optimized resource usage), improved interoperability, improved usability, and better utilization of HQP expertise across Canada. The lack of high availability and

---

[212] IPERION HS: About http://www.iperionhs.eu/about/ (retrieved February 2021).

[213] ESFRI Roadmap 2018: ESFRI Projects http://roadmap2018.esfri.eu/media/1044/part1-project-landmarks-list.pdf (retrieved February 2021).

[214] Government of Canada ISED: DIGITAL RESEARCH INFRASTRUCTURE CONTRIBUTION PROGRAM - PROGRAM GUIDE https://www.ic.gc.ca/eic/site/136.nsf/vwapj/DRIContributionProgram_ProgramGuide.pdf/$file/DRIContributionProgram_ProgramGuide.pdf (retrieved March 2021).

high redundancy in CCF infrastructure is also of concern. In case of a major fire a host site could potentially fully and completely lose a site and all of the data stored at that site. Currently there is no off-site replication, and typically, data on tape is hosted next to the cluster. Having nationally coordinated operation, and funding for backup and archival storage systems should include off-site replication to mitigate potential for data loss.

The more central and coordinated planning also needs to explicitly consider not only traditional ARC, but also short-, mid- and long-term storage, RDM, and RS needs holistically in one envelope, while putting added focus on underserved disciplines, audiences, and communities. Since 2018, the GC and ISED have recognized this critical need so that one of the key mandates for the Alliance will be to take this more coordinated approach to funding and operating Canadian DRI ecosystem going forward. Currently the Alliance is engaged in its Needs Assessment process including ARC, RDM, and RS Current State assessments, Call for Position Papers and DRI Documentation, Researcher Needs Survey, Townhalls etc. This process will result in the Alliance's holistic Strategic Plan for Canadian DRI, New Service Delivery Model, and ultimately the comprehensive DRI Funding Proposal, to be submitted to ISED in late 2021.

To unify the delivery of ARC services federally the Alliance needs to create consistent Participation Agreement and Service Level Agreement (SLA) policies across all service providers, regions, and host sites. Currently services are provided either without any SLAs or using SLAs that are mostly 'best effort' with few measurable Key Performance Indicators (KPIs). In order to optimize federal DRI investment and manage service delivery expectations, going forward the SLAs need to include clear KPIs with corresponding enforcement framework. This process will also need to include researcher consultation regarding evolving expectations, e.g., if the user community would require 24/7 support, or if zero-data loss for storage.

## Attraction and retention of HQP

ARC systems are almost by definition highly complex, state-of-the-art systems combining a variety of leading-edge technologies as required by a myriad of distinct use cases. Working on these systems requires highly qualified senior level staff that will take years to train. Such staff will also require time and resources for on-going training and learning on the job to keep up with constantly evolving trends. These skills of these HQP staff are also often highly sought after by the commercial operators, particularly in recent years as cloud computing, AI, and quantum computing are either emerging or already being adopted by more mainstream business operations.

CCF has a pool of roughly 250 very highly skilled and motivated employees. Keeping these employees at the Alliance should be a very high priority for the Alliance so that it can continue providing the high quality of resources and services in the years ahead. Monetary compensation is important no matter how motivated the person is. Unfortunately, the academic and governmental salary rates can not directly compete with the private sector. Where CCF and the Alliance ARC community has an edge, are the less directly tangible factors provided largely by the universities, like benefits, job security, work environment, work-life balance etc. The Alliance needs to collaborate with all DRI providers to strengthen these factors to keep public sector DRI HQP positions competitive and to retain the HQP.

A particularly clear problem is the long-term job security of the HQP. Many full-time ARC HQP positions are still limited term. In Ontario ARC staff have been hired explicitly to be dedicated to

the national systems, while the funding comes in 5-yr tranches and can be reviewed and changed annually by CFI and to some extent CC. Universities are unlikely to make these permanent positions and it is quite common for a position to be contingent on grant funding. This situation could potentially be a problem, although there are many HQP who are 10 year "veterans" even under these circumstances. Anecdotally, the biggest immediate concern of CCF staff seems to be the uncertainty caused by the transition to the Alliance – and CCF staff employment past April 2022.

The ARC WG does not currently have data on the number of limited term contracts within the HQP FTEs beyond a general observation that 'many people work on yearly contracts'. Any short-term contracts (particularly less than 12 months) are discouraging to the employees due to longer term uncertainty and can easily result in employees looking for other opportunities with permanent positions. If term contracts are necessary due to budgeting reasons, these reasons need to be articulated clearly to the employees with the emphasis on long term continuity. Preferably the funding institutions should commit to funding staffing beyond the explicit project funding term. Moreover, the fundamental funding agency, e.g., ISED and Government of Canada should commit to funding core DRI operations and staffing on an on-going, non-term basis.

Compute Ontario commissioned a report with recommendations on HQP in Ontario in 2018. The Malatest report[215] recommended twelve focus areas to attract, retain, and develop HQP. In the category of attracting talent, the recommendations were to

1. promote an image of ARC as a meaningful career path, similar to e.g., the Research Software Engineer (RSE) as a career path initiative that started in the UK and has now become an international initiative[216],

2. support women working with ARC,

3. increase recognition for HQP working with ARC, and

4. promote ARC in social sciences and humanities.

To retain HQP talent, expanding to all of Canada and not just Ontario, the report recommended

1. support longer-term employment in academia,

2. nurture Canada's HQP community,

3. establish Canada as a globally recognized hub for ARC, and

4. market Canada as a leading place to live.

Lastly, the report emphasized the need for developing HQP expertise and skills via

1. providing more educational opportunities to develop computing technical skills and computational thinking skills,

---

[215] Malatest, commissioned by Compute Ontario: Highly Qualified Personnel Study https://computeontario.ca/publications/reports/highly-qualified-personnel-study/ (April 2018).

[216] Society of Research Software Engineering: History https://society-rse.org/about/history/ (retrieved March 2021).

2.  establishing campus champions,

3.  supporting development of ARC curricula, and

4.  entrenching teamwork and communication skills.

The ARC WG emphasizes #5 the need for long-term funding, #2 the support for diversity and women working with DRI, and #4 the need for promoting DRI in social sciences and humanities.

## International collaboration and competitiveness

International collaboration is front and center on many fields. For Canada to stay competitive and relevant in sciences globally, it needs to have a competitive domestic DRI ecosystem. Lack of sustained and predictable funding and insufficient domestic ARC compute resources continue to be a problem and the situation has not changed for the better in any fundamental way since the 2017 LCDRI ARC report identified this issue. On the contrary, the need for global collaboration has been increasing via global initiatives in sharing knowledge, funding, and data related for example to climate change and COVID-19. Canadians already collaborate internationally on multiple specific ARC intensive fields, e.g., CERN's ATLAS-TRIUMF project in particle physics and the corresponding integrated Tier-2s, and Tier-3 analysis platforms[217], or Canadian Astronomy Data Centre (CADC) in astronomy[218], leveraging dedicated and substantial Canadian ARC resources. Going forward the DRI ecosystem needs to expand beyond these well-established discipline- and project-based solutions to embrace new international collaboration and projects, both by providing basic facilities infrastructure to support new CADC type ARC systems and by providing more general DRI capacity and services to attract global research and collaborators in new and underserved disciplines.

Canadian researchers might not be considered as attractive international collaborators if they are not able to carry their own weight when it comes to ability to conduct research leveraging domestic ARC resources, forcing the researchers to use their international collaborators' ARC allocations, located in e.g., the US, Australia, EU, or China. The scope of this problem, Canadian academics having to go abroad for ARC resources, is not well quantified at the moment. The Alliance is currently conducting a researcher needs survey as a part of its Needs Assessment process. This survey will also gauge the use of international ARC resources, and while not strictly quantitative the results will provide insights to the scope of this issue.

Additionally, attracting top academic talent globally could be a problem. One can imagine a situation in the highly competitive top academic faculty marketplace where Canadian academic institutions would have more difficulty hiring top researchers when the peer foreign institutions can potentially provide better local, regional, and national level ARC services and resources. An indicator of such poor global competitiveness is Canada's ranking as second last among the G7 when it comes to Top500 listed compute power per GDP, as discussed earlier in this document.

In addition to international collaboration, recent years have seen increasing collaboration at regional and interprovincial levels. As with the international collaboration, these initiatives face

---

[217] TRIUMF: ATLAS Tier-1 Data Centre https://www.triumf.ca/atlas-tier-1-data-centre (retrieved February 2021).

[218] CFI: Canadian Astronomy Data Centre https://navigator.innovation.ca/en/facility/national-research-council-canada/canadian-astronomy-data-centre (retrieved February 2021).

material legal, policy, interoperability, and funding challenges that need substantial resources and research support to solve and mitigate. Going forward the Alliance can and should play a key role in establishing connections and collaborations with its international peer organizations, provinces, and regional organizations to reduce and remove barriers to collaborative research initiatives.

## Coordination of federal science investment and the provision of ARC

In their 2017 ARC Current State report LCDRI addressed the still persistent issue of federal level science funding and allocation of ARC resources as follows: "Federally funded research is increasingly reliant upon access to ARC resources. However, there is currently no formal coordination between Compute Canada and the three federal granting councils, nor with specific, high impact programs such as the Canada Research Chairs (CRC), the National Centres of Excellence (NCEs), or the Canada First Research Excellence Fund (CFREF). This can mean that researchers and the ARC administrators responsible for serving them are often caught by surprise in terms of finding access and responding to resource requirements. Stronger, more formalized, collaborative planning among these parties would ensure that efficiencies are realized, investments are exploited maximally, and that more streamlined processes are available to researchers, reducing the risk that they could find themselves with grant money, but without access to the ARC resources they need."[1]

Currently Compute Canada does not have access to tri-agency grant applications, preventing them from preparing for future resource needs, commenting on the technical viability of the application, or assessing general future computing trends. Sometimes this disconnect results in the PIs having to first apply and get research funding from one of the tri-agencies via a peer-reviewed process, only to then having to apply in a separate independent, peer-reviewed, process for compute cycles from CCF. This process is inefficient, cumbersome, and potential rejection or reduction of available compute resources in the second part can derail otherwise funded research.

In some fields, e.g., social sciences and humanities, the funding agencies might not allow funding for IT systems administrative and operational costs. Such policies exacerbate the issues related to DRI adoption in these fields since they already historically are underrepresented in DRI usage, and do not necessarily have internal IT and ARC resources to support any prospective ARC systems, compared to fields with long history in leveraging ARC and building corresponding ARC support operations. For example, for social sciences funding in the EU, the technical aspects and operational support needs of a DRI leveraging application are assessed as part of the application process.

In short, a more holistic grant review process is needed in Canada for projects that require DRI resources. This process needs to include both research funding agencies and relevant service provider institutions, incorporating the technical feasibility assessment, and consideration of funding the operational costs. The Alliance and tri-council are working on addressing these major issues and are for example considering a 'pre-award evaluation' process where the DRI needs of an application would be reviewed before approval. Fundamentally, such a holistic, coordinated planning and cooperation between main federal funding agencies is needed to build a viable DRI ecosystem that incorporates ARC, RDM and RS needs for current and future DRI user disciplines.

In addition to the above, the on-going operational costs related to the so-called 'contributed systems' is problematic. Currently CFI requires that systems funded through their Innovation programs are installed on CCF main host sites so that they can contribute resources to the

common computing pool when not in use by the primary award recipient. If the applicant argues that the requested ARC system cannot be part of the main CCF systems, such requests need to be approved on a case-by-case basis by Compute Canada, and by the Alliance as of April 2022. If a system is permitted to not be a contributed system and if the system is not co-located on CCF datacenter, the operational costs of such systems naturally fall on the applicant or applicant's home institution. In the case of contributed systems, the situation is currently complicated since there is no link between the installation of the contributed system and operations funding. The operational funding of a system that is partially used by an individual research group, and partially used by CCF community in general is not addressed by the funding agencies at the time of funding approval and is left for the host institution to sort out with the applicant, often after the system infrastructure funding has been approved. ISED in collaboration with CFI and the Alliance has recognized this issue and has set up a working group with the mandate to create terms and policies for funding contributed systems operations. It could be potentially easier for researchers if they could purchase computing services directly from the Alliance, rather than researchers having to purchase the hardware themselves.

## Keeping pace with technological and cultural diversity

Historically, the technological advances in DRI were focused on ARC. Recent advances have brought to the forefront the relevance and benefits of RS and RDM to modern research, so that the DRI needs to involve all three aspects in its design. In addition to the methodological and DRI toolchain advances, the DRI ecosystem has become a valuable and emerging tool for non-traditional ARC users and disciplines, including for example social sciences, health sciences or indigenous studies. As more and more data is becoming available, these disciplines have recognized the enormous opportunities DRI systems can potentially provide for their research. In many cases these initiatives additionally include concerns related to sensitivity, privacy, ownership, and security of the data. The needs of various underrepresented groups are being better recognized by society, putting emphasis on addressing the needs and requirements of these communities, including for example racialized, LGBTQ+, and indigenous communities. If the researchers in these disciplines and communities may not be familiar with modern ARC systems and tools. They may require new DRI innovations, training, documentation, dedicated HQP or tools to enable their access to DRI systems. The challenges related to EDI are discussed in more detail in the "Equity, diversity, and inclusion and indigenous access and minority representation" section below.

Just purely on the traditional ARC technology side the landscape is changing at a quick pace with, for example, the emergence and increasing adoption of GPU computing, new AI chip architectures, quantum computing or cloud computing. Currently the Alliance, CC, and CCF host and regional institutions have experts following recent technological developments in the ARC field.

All of the above groups and trends require new dedicated focus and resources in addressing their needs and related technologies and technological advances that can serve them. Canada's current DRI ecosystem has not been well equipped and funded to keep up with and address these needs. The focus has been on the needs of traditional ARC user communities, with some expansion in testing of new technologies (e.g., new use paradigms beyond command line access leveraging science gateways etc.) but without coordinated national level effort or funding.

As discussed earlier the 2018 GC DRI budget announcement and formation of the Alliance included a recognition of some of the above technological and cultural diversity issues. Going forward the Alliance will include the above concerns in its central, national level DRI planning and funding. The Alliance DRI team will be actively keeping up with technological trends via e.g., following international technical publications and participating in (currently virtual) international DRI conferences. The Alliance will be conducting a Current Needs Assessment in the Winter-Spring 2021 that will review what future DRI technologies Canadian researchers would have a need for. These inputs will then drive a gap analysis, and holistic strategic planning, resulting in a funding proposal to ISED that will be submitted in late 2021.

### Leveraging cross-sectoral ARC resources

The authors are not aware of any ARC resources that would be available to Canadian researchers in general at SSC (ECCC), Communication Security Establishment's (CSE) and Canadian Security Intelligence Service's (CSIS) hosted ARC facilities. The CSE and CSIS ARC operations understandably are limited to providing (classified and unclassified) services only to their core audience(s) and are not integrated with the CCF infrastructure. Regarding ECCC's ARC resources, climate and weather academics have historically had mechanisms and connections to get access to ECCC's SSC hosted systems.  Meteorological services typically buy two systems at a time - one dedicated to production work and one "failover" system that is used for research (unless the other ones fail).

In the United States the vast and varied ecosystem of ARC resources is more widely available to academic researchers. Moreover, resources that are primarily funded by e.g. defense funding are also (partially) available for general science initiatives per their mandates: the world's third fastest supercomputer, Sierra, hosted by Lawrence Livermore National Laboratory (LLNL) Livermore Computing Centre and primarily funded by US Department of Energy (DOE) and National Nuclear Security Administration (NNSA), is used both for classified computer simulations related to e.g. management of the United States' nuclear weapons stockpile, and for unclassified simulations in the fields of computational physics, climate change and global security.[219] As a recent example of cross-sectoral ARC collaboration in the US, in March 2020 The White House announced Covid-19 partnership that provided international researchers access to supercomputing resources at multiple DOE supercomputing facilities in addition to commercial cloud providers.[220]

There is no public information available regarding Communication Security Establishment's (CSE) and Canadian Security Intelligence Service's (CSIS) ARC capabilities. Nor are any such systems listed on the current Top500 list,[221] suggesting that CSE and CSIS have not submitted their systems to the list, likely due to national security reasons. Due to the computationally intensive work these institutions conduct, they are likely to have very powerful supercomputers. As the above example from the US shows, such systems could potentially be of interest to

---

[219] Lawrence Livermore National Laboratories Livermore Computing Centre: Mission Support https://hpc.llnl.gov/about-us/mission-support (retrieved September 2020).

[220] Lawrence Livermore National Laboratories: New partnership to unleash U.S. supercomputing resources in the fight against COVID-19 https://www.llnl.gov/news/new-partnership-unleash-us-supercomputing-resources-fight-against-covid-19 (retrieved September 2020).

[221] Top500: November 2020 https://www.top500.org/lists/top500/2020/11/download/TOP500_202011.xlsx (retrieved January 2021).

Canadian researchers, although understandably such collaboration can only happen if mandates are fundamentally changed (similar to e.g., DOE's split mandate).

Having access to classified (CSE, CSIS), or unclassified (SSC/ECCC) government owned and funded ARC resources, even at limited scale and scope, would benefit Canadian academics and the country as a whole. One could imagine that selected research projects of critical national relevance (e.g., Covid-19 vaccine development) that leverage highly scalable codes and require leading edge supercomputing capability (i.e., not simply capacity computing) could be run on non-traditional Canadian publicly funded ARC resources, provided by SSC/ECCC. A central DRI organization, e.g., the Alliance, could have a role in initiating, negotiating and managing cross-sectoral DRI collaboration(s) and related frameworks.

## Researcher awareness and adoption of ARC

The limited researcher awareness and adoption of ARC and DRI is a major problem in Canada and globally. There are roughly 33000 full and associate level university professors in Canada,[222] while currently CCF lists roughly 5500 PI accounts related to the RAC process, i.e., ca. 17% of full and associate professors are potentially leveraging CCF resources (directly). Out of these 5500 PI accounts and corresponding research groups, 3177 had CPU usage over the past 12 months, i.e., ca. 10% of the potential PI led research groups in Canada were leveraging CCF resources. Setting the CPU usage threshold at more than 4 core-years (i.e., the rough equivalent of a semi-modern laptop running 365/24/7), this figure drops to 1545 PI led groups, and to 578 PI led research groups if we consider the RAC threshold of 50 core-years (i.e., rough equivalent of 5-10 workstations running at 365/24/7). This last group, the "heavy" users of CCF facilities, correspond to roughly 2% of eligible PI led research groups in Canada.

Adoption of DRI can be viewed also from the point of view of compute cycle usage in different disciplines. Researchers in "underrepresented" disciplines (e.g., social sciences, humanities, psychology, and business) represent 17% of all the faculty level users in the CCF user database but use only 1% of all the CPU resources. Part of this imbalance can be attributed to smaller research group sizes in these fields in the CCF database (i.e., user to PI ratio) and potentially attributed to the CPU heavy simulations that e.g., physics and astrophysics researchers have developed and utilize in their research compared to other disciplines that are not necessarily as compute intensive. But nothing should in principle stop these underrepresented disciplines from utilizing much more ARC resources once one considers the size and variety of potentially available data heavy resources, e.g., databases, and longitudinal and transactional data owned by private sector entities, provincial and federal governments, Statistics Canada, regional health authorities etc. Combining and cross-referencing such information sources would require substantial storage and computing resources. Potential AI and ML based analysis of these resources would additionally provide new insights in these disciplines. That is, the fields of social sciences, humanities, psychology, business etc. that are currently underrepresented in the ARC community have huge potential for leveraging ARC and DRI to advance their fields. By not serving these disciplines at their full representative level, Canada is faced with substantial risk of lost opportunity, and being left behind by the global competition.

---

[222] Statistics Canada: Number of full-time teaching staff at Canadian universities, by rank, sex https://doi.org/10.25318/3710007601-eng (retrieved April 2021).

In addition to specific disciplines and communities that are not leveraging DRI, there are researchers even within "traditional" disciplines who don't access CC systems for a variety of reasons, e.g., researchers are not aware of the service portfolio, consider the user interfaces and usage too complicated, or perhaps have given up because their application was rejected, or they had a bad experience. To address researcher awareness, explicit concentrated efforts are needed to reach out to all these communities and people. Such outreach should not only discuss currently available solutions in these fields but should include forward looking brainstorming advocacy where senior level DRI experts would meet with experts in these disciplines to hash out potential new approaches and technologies to advance science. Advocacy should also highlight relevant existing use cases in the respective fields to show how DRI tools can benefit these disciplines. Any outreach for awareness needs to be coupled with DRI resources, including infrastructure, services, training, support staffing, and new improved access and usability technologies so that any interested researchers will have the ability to start leveraging DRI. Without this coupling, the outreach can potentially just result in hypotheticals without actual actionables. Interesting usability solutions could e.g. be Virtual Desktop Infrastructure (VDI) solutions, where the end-users can access ARC resources in a familiar desktop environment, or web browser based access, such as notebooks in Jupyter or web portals. Another important factor for further adoption of ARC and DRI is the need for improved software and library stacks and corresponding integration. All of the above should also include substantial targeted training component.

As part of its mandate going forward, the Alliance will try to address the above issues and provide solutions related to researcher awareness and adoption in underrepresented fields and disciplines. At a very fundamental level, the expansion of the mandate from ARC to RS and RDM will allow for more holistically designed and comprehensive DRI solutions that should by design be more accessible and relevant to wider audiences. The Alliance is already involved in an explicit and concentrated effort to reach out to underrepresented research communities in its Needs Assessment process, where the target audience has not been only the traditional ARC users (e.g., based on CCF user databases), but so that the outreach has explicitly included major discipline specific organizations, VPRs of all Canadian universities. The resulting new Strategic Plan, Service Delivery Model, and Funding will explicitly address and solve issues related to the current situation of underrepresentation.

## Environmental impact

ARC data centers are dense with servers and supporting infrastructure that by its nature is being run at close to full compute utilization, consuming a lot of power (estimated 1% of global electricity consumption). The current estimate for the power consumption of CC data centers is of the order of 3-5 MW. Canada's overall power supply is 67% from renewables and 82% from non-GHG emitting sources so that the direct greenhouse gas emission footprint is reduced.

The compute power efficiency (flops per W) can be used as an indicative proxy for certain aspects of the environmental friendliness of an ARC system. Although this metric only considers the power consumed by the ARC system itself and does not consider additional factors, e.g., that all power consumed must be cooled off, adding to the total energy and environmental impact. The current cooling overhead is estimated to be of order 10-30% for CC national sites. Going beyond the cluster efficiencies, reuse of heat generated by the cluster is probably one of the key ingredients that many CC sites are missing. Notably this is a function of the data center infrastructure, not a

function of the cluster itself. Improving compute power efficiency and heat reuse are key components in improving the total environmental impact of a supercomputer system.

Top500 and related Green500 are well-known and well-established ARC benchmarks that have their substantial limitations but can be used as proxies for ARC computing performance and efficiency. These benchmarks represent semi-theoretical performance of hardware by running a micro-benchmark which is often not representative of the actual workloads being run by the research community. Additionally, GPUs might be efficient in theory, but they can also be very poorly used or not be usable at all by various applications. A much more ideal approach would be to focus on code performance of real research code on one side, and on data center efficiency (energy reuse and such) on the other side but expanding such metrics across the global ecosystem for comparison purposes would be difficult.

Within the above caveats the recently installed CCF systems are positioned relatively well internationally in Green500 benchmark. Canada has five entries in top 100 of Green500 rankings with Cedar (two entries with different CPU+GPU combinations: #16, ca. 11 GFlops/W and #33, ca. 8 GFlops/W), Beluga (#21, ca. 9.5 GFlops/W), Cedar (pure CPU, #33, ca. 8 GFlops/W), and Niagara (#67, ca. 3.9 GFlops/W) .[223] It is encouraging that the recently acquired Canadian systems are well positioned on the compute power efficiency front. On the other hand, looking at raw GFlops performance Canada has only two systems in top100 of Top500 with Niagara (#70) as the highest ranked Canadian ARC system and Cedar (GPU accelerated, at #74) trailing closely behind.[224] In other words, Canadian ARC systems are internationally better positioned on the power efficiency side compared to rankings on the raw compute power side.

For some applications and use cases, when applicable, the efficiency of accelerators is impressive both when it comes to raw compute power and compute power efficiency. The top ranked system on the Green500 list, Japanese MN-3 system at ca. 21 GFlops per W, uses custom designed MN-Core ASIC accelerator chips that are engineered for one specific task, the training phase in deep learning workloads.[225] Six machines on top 10 of the Top500 list use accelerators. Compute power efficiency of GPU accelerators is demonstrated by comparing purely CPU based Niagara (#70 on Top500) to Nvidia V100 GPU accelerated Cedar (#74 on Top500). Niagara achieves roughly the same compute power as Cedar while consuming nearly three times more electricity (920 kW v. 310 kW).  It should be emphasized that many problems and applications can not be easily (or at all) recast or ported to GPUs. For this reason, TACC's Frontera, an NSF-funded system for general academic research, is primarily a pure CPU system with additional GPU capability.

It is important to note that looking at compute power efficiency alone does not consider other important environmental factors, e.g., environmental sustainability and friendliness of electricity production, total lifetime environmental impact of the systems and materials used, power efficiency of the A/C and cooling systems (if even required) or potential secondary use of waste heat. A key practical implementation is the effective reuse of any excess heat,y for example by

---

[223] Green500: June 2020 listing https://www.top500.org/files/green500/green500_top_202006.xls (retrieved September 2020).

[224] Top500: November 2020 https://www.top500.org/lists/top500/2020/11/download/TOP500_202011.xlsx (retrieved January 2021).

[225] InsideHPC White Paper - Supermicro Contributes to the MN-3 Supercomputer Earning #1 on Green500 list https://insidehpc.com/white-paper/supermicro-contributes-to-the-mn-3-supercomputer-earning-1-on-green500-list/ (retrieved September 2020).

heating of the adjacent buildings or even communities. For example, the heat produced by Mammouth at Université de Sherbrooke and Colosse at Université Laval is used to heat part of their respective campuses.[226,227]

Repurposing of used equipment extends the lifetime of equipment and avoids unnecessary landfill. It must be kept in mind that old equipment is not as energy and compute power efficient as newer hardware and fails more often, increasing electricity costs, reducing space usage efficiency, and increasing staffing costs.

## Securing the national platform

According to the Canadian Security Intelligence Service (CSIS) in recent years, there has been an increase in frequency and sophistication of cyber threats against Canadian research interests, particularly in the fields of biopharma and health sectors, artificial intelligence, quantum computing, ocean technology and aerospace.[228] The security breaches are often related to human factors (e.g., poor passwords, or social engineering) but can also include technological solutions. For example, the US is concerned about the security of Huawei 5G networking equipment and is encouraging its closest allies to not use such equipment.[229] The Canadian Government has not yet decided on this matter, while such issues are current and of concern for Canadian DRI ecosystem since some CCF equipment (e.g., Graham) was made by Huawei. Such public concerns and pressure about some particular vendors need to be balanced with the formal procurement policies and rules for e.g., Requests for Proposal (RFPs).

As the DRI infrastructure is becoming more national and centralized, any actions to secure the infrastructure need to be coordinated between all participants, at local host sites, regional and federal levels. These efforts need to protect sensitive data, personal information, intellectual property, and digital and strategic research assets.

ARC solutions were traditionally built for and operated in isolated networks so that external threats were not a primary design concern, rather the focus was traditional ARC performance characteristics of raw compute power, latency, bandwidth, or energy efficiency. As such, ARC systems were explicitly not designed to consider prevention of malicious attacks or attempts. As new opportunities, sensitive datasets, new paradigms like edge computing, and new research disciplines are starting to leverage ARC, the systems are becoming more accessible and nationwide, and thus also more exposed to unauthorized external access attempts and threats. The threat landscape has evolved significantly in the last couple of years with widespread and well-publicized attacks on many HPC systems as in the crypto-mining related hacking on

---

[226] U. Sherbrooke: L'ordinateur Mammouth au premier rang au Canada https://www.usherbrooke.ca/sciences/accueil/nouvelles/nouvelles-details/article/17844/ (retrieved May 2021).

[227] CBC Radio-Canada: Des serveurs informatiques pour chauffer l'Université Laval https://ici.radio-canada.ca/nouvelle/1143499/serveurs-informatiques-chauffage-universite-laval-recuperation-energie (retrieved May 2021).

[228] Global News: China and Russia 'aggressively' targeting Canadians, CSIS director warns https://globalnews.ca/news/7629494/china-and-russia-targeting-canadians-csis-director/ (retrieved February 2021).

[229] CBC: Biden team sees Huawei as a threat and wants to talk to allies https://www.cbc.ca/news/world/biden-huawei-canada-1.5900991 (retrieved February 2021).

European supercomputing infrastructure in May 2020.[230] CCF ARC sites have also experienced various types of attacks over the years. As a response to the emerging threat landscape, CCF affiliated ARC sites have increasingly added resources and operations to address security concerns, including improved security culture, funding, and dedicated group meetings. CC has had a security officer for several years, and all sites have hired full-time dedicated security staff as of 2020. Host sites have also installed more network monitoring capabilities and are working much more closely with their institutional Chief Information Security Officers (CISOs) and each other than in the past.

Going forward, balancing the ARC performance, accessibility, and security will continue to be a challenge as DRI systems need to be hardened against cyberattacks, both technologically, operationally, and at policy level. To address the cybersecurity in DRI is part of the Alliance's mandate, and the Alliance will need to collaborate with and leverage the Government of Canada's Canadian Center for Cyber Security's (CCCS) expertise and resources in planning and securing DRI operations and services within its funding envelope. The Alliance will also be forming a Cybersecurity working group to involve DRI security considerations in its planning, operation and funding decision making.

## Heterogeneous service portfolio and delivery across national ARC computing systems

ARC and DRI services and requirements are complex, and the needs and technologies of different user disciplines and groups are not homogeneous. Even individual systems with similar functional goals are not identical (e.g., cedar and graham have multiple CPU and GPU types in each due to different funding and procurement cycles). Niagara is architected specifically for the performance needs of large-scale parallel computations. Arbutus cloud is naturally very different by design from the general-purpose compute systems. Systems for sensitive data will need to be customized to comply with local, regional, and national requirements, in addition to the discipline and function specific needs. On the commercial side Amazon's AWS offers myriad instances (40+) and services (300+) serving its global customer base. Within DRI, homogeneity would result in less innovation and lower competitive position for researchers. The key is to apply common standards where applicable but not just for the sake of uniformity.

On the other hand, lack of homogeneity in services and setup can potentially increase support staffing costs, operating costs, and introduce duplication of effort costs within the CCF, as well as being confusing for end users. Having the same service delivered multiple times by different groups potentially introduces not only differences in the end-user facing service, but also potential security risks as the seniority and expertise of different implementing teams vary.

With the above in mind, the current CCF systems have multiple differences in system setup, access to support and documentation, and service portfolio. Some of the notable differences between the main CCF systems include: internet access availability from the compute nodes (blocked in some systems, while allowed on others); front / login nodes are not uniformly configured (e.g. on some systems one can not setup 'crontab' based scheduled execution of routine tasks, system memory and CPU limits differences, different access policies for data transfers); scheduling policies (e.g. differences in max runtime, number of jobs policies); and Globus file transfer system (fundamental differences in authentication, e.g. CC v. 3rd party main

---

[230] BBC News: Europe's supercomputers hijacked by attackers for crypto mining
https://www.bbc.com/news/technology-52709660 (retrieved March 2021).

Globus). System level differences can potentially lead to problems and inconveniences for end-users as they migrate their workloads between systems. Although data locality is a main concern for workload transferability as moving large or complex datasets across systems might not be feasible.

Different systems also offer materially different service portfolios, e.g., Virtual Deskop Infrastructure (VDIs) for remote desktop-based work, NextCloud for cloud-based storage, and Jupyter Notebook science gateway for usability and collaboration are available on some CCF sites. If the same service is available on multiple sites, it is often developed and maintained by different, local teams, potentially resulting in differences in service delivery and configuration between sites for the said service. Documentation and helpdesk delivery also differs significantly for one of the sites, compared to the other sites, with a unilingual English wiki being the primary documentation source, and a parallel helpdesk separate from the CCF national helpdesk being used.

Whenever possible, a unified service offering should be favored. When not possible, having a centralized catalogue of all CCF services, a service portfolio or portal, would be helpful for the researchers to find centrally 'where and what' services are available. Such service portals could include not only service listings and descriptions, but also links to key documentation, and service provider and contact information for further information and queries. A central service catalogue could potentially provide a high-level single point of entry for new and existing CCF users.

## Equity, diversity, and inclusion and indigenous access and minority representation

Equity, diversity, and inclusion (EDI) are key moral imperatives, the importance of which are widely recognized in modern liberal democracies. Historically, in ARC delivery EDI have not been recognized as challenge areas, and for example the 2017 LCDRI ARC Current State report does not explicitly raise EDI concerns. The CCF does not currently collect EDI data so the current situation within the CCF is not well understood. CCF does collect institutional affiliations, indicating regional diversity, but this data is not sufficient for proper consideration of how EDI are presented in the CCF community.

Major ARC resources in Canada and globally are centrally located on dedicated host sites and are accessed and used remotely (e.g., by using SSH terminal access or Jupyter web portals). In principle the remote access nature of these systems should benefit equity of access from more remote areas. In practice, the rural and Northern regions in Canada have problems with access to reliable and high-speed internet. In August 2020 CANARIE announced an important milestone for equity in academic internet networking with Nunavut Arctic College joining Canada's National Research and Education Network (NREN).[231] With this addition all thirteen provinces and territories have access to high-speed internet for academic and research use. Close collaboration between ARC and DRI providers and networking providers is required for effective and equitable access to DRI resources.

The DRI ecosystem must also accommodate non-native English speakers, in particular French communities, but also non-native English-speaking users (via e.g., clearly written, and edited user documentation). All documentation and services must be available in both official languages, and

---

[231] CANARIE Press Release: Nunavut Joins Canada's National Research and Education Network to Enhance Nunavummiut Access to Colleagues, Data, and Scientific Facilities https://www.canarie.ca/nunavut-joins-canadas-nren/ (Retrieved February 2021).

so that the quality of translation is on par with natively written text and not at the level of automated translations. Key documentation and services should additionally be available in selected indigenous languages. Key events and conferences should include sign language translation. As a federal organization, the Alliance should equip itself to provide bidirectional (if not multidirectional) translation for major events.

Going forward, the Alliance should systematically collect EDI related data regarding its operations and service delivery so that the current situation can be assessed, problem areas identified, and solutions provided, beyond just 'lip service' regarding general importance of EDI. As an indicator of the Alliance's commitment to EDI, the Researcher Needs Assessment survey that the Alliance submitted to all Canadian researchers in February 2021 was praised for its leading edge and comprehensive EDI related questions. Such data will provide critical insights to understanding and solving EDI related issues in Canadian DRI ecosystem.

EDI should not be considered only as a separate item, in its own silo, but rather should be part of all discussions and decision making. As an example of this, in forming its Researcher Council advisory body, the Alliance focused particular attention on EDI. As a result, the current RC membership in general well reflects Canadian society, the indigenous representation notwithstanding (due to difficulties finding indigenous academics who could dedicate the time needed to serve in the Council). The Alliance's RDM WG currently has two members of indigenous background, bringing valuable data rights and ownership expertise to the research data management discussion and planning.

## Lack of Facilities and Services for Sensitive Data

The potential in analyzing sensitive data sources has emerged as a key trend in recent years in Canada and globally. In discussions with the Canadian research community, the Needs Assessment position paper submission, and the related researcher survey all indicate the urgent need for sensitive data related DRI in Canada. Sensitive datasets can include personal health, indigenous, statistical census, financial, tax, social media, and transactional commercial, municipal, provincial, and federal level data. On the upside, analyzing and potentially cross-referencing such data can provide huge benefits for advancing social sciences, economics, humanities, and human health sector, all ultimately benefiting Canadian society, policy, and decision making. On the downside, these datasets often involve highly sensitive data at individual or corporate level with strict requirements for privacy and sensitivity. Balancing the benefits and risks for accessing and researching such data is an on-going challenge involving multiple non-trivial aspects including

- public perception,

- communications,

- expectation management,

- municipal, provincial, federal, and international legislation, requirements, and related mandated policies and procedures,

- questions of data ownership, and open access

The related supporting technological solutions include:

- at-rest and in-transfer encryption,

- cybersecurity,

- application-level security, and

- policy implementation enforcement technologies

while still maintaining usability and accessibility for the end-users, the researchers. All of the above factors need to be considered holistically when designing and creating secure research data platform(s) processing sensitive data.

Currently there are no fully national, centrally managed and provided, publicly available CCF DRI platforms for sensitive data in Canada. The closest to a truly national level non-CCF platform is probably Canadian Research Data Center Network (CRDCN) that in collaboration with Statistics Canada (SC), Shared Services Canada, CANARIE, and local universities provides access to primarily SC owned data for social science researchers via secure Research Data Center (RDC) offices located in over 30 university campuses across the country. Access to these resources is available to all eligible Canadian researchers in participating universities. This system is not yet fully centralized, but CRDCN is currently in the process of building a centralized remotely accessible ARC cluster (aka virtual-RDC "vRDC") to replace the current local to each RDC workstation and server infrastructure. This system will be built to adhere to highly demanding Government of Canada Protected B data protection requirements, including additional Statistics Canada and potentially local and regional requirements (in order to e.g., cross-reference SC owned data with provincial level health care transactional data (e.g., Ontario Health Insurance Plan (OHIP)). At regional and local levels there are multiple facilities for sensitive data, e.g., the earlier discussed, and highly successful HPC4Health initiative.

The challenge for the Alliance and Canadian DRI ecosystem going forward is to provide DRI solutions for sensitive data processing at national level, and at scale while adhering to all local, regional, and federal level requirements. As legal and other requirements for different datasets, e.g., health data compared to data from social media, are often very different, the design of DRI systems that could handle a wide variety of use cases at the national level will be challenging. Achieving this requires national operation, coordination, and cooperation. From a legal perspective, the "chain of command" must be clearly articulated, which will be difficult in a federated model, while at the same time considering all local, regional, and federal level requirements. Going forward the Alliance is planning on establishing sensitive data working groups to work on providing solutions for sensitive and secure data.

### Funding of data center construction, maintenance, and operations

Building a data center and related infrastructure to support ARC systems is very capital intensive, including procurement of real estate and facilities, and proper electricity, HVAC, and networking infrastructure. ARC facilities generally have more demanding requirements for power, cooling, floor-space and even floor-loading than can be met by many enterprise data centers. These investments need to be designed and built for long term projected needs involving substantial planning and estimation of future trends and needs, preferably including long-term commitments from all stakeholders to justify the costs and related amortization. Data center infrastructure costs

are not covered at the moment by the CFI or other federal funding agencies. Instead building such data center facilities is often covered by local hosting sites in connection with provincial level funding, leading to potential inefficiencies at national level, and inequalities between Canadian higher-ed institutions when some institutions acquire substantial costs while many institutions are not sharing the infrastructure costs.

ARC operating costs are not part of the CFI IF funding envelope, rather they are funded through a different CFI program (MSI). As such the procurement of new ARC systems and infrastructure is disconnected from the funding of corresponding operational costs and concerns, even though these two are closely connected. The operational costs of running these data centers, beyond just power and staffing costs, should be included in all decision making at national level. Currently CFI operating costs do not reflect true operating costs since data center space, physical security (alarms, monitoring etc.) and infrastructure maintenance (chillers, pumps etc.)  are generally not eligible costs.

In other words, modern ARC infrastructure involves substantial capital and operational costs, and staffing needs that should all be considered and secured at the time of the infrastructure investment. This should include solving the major problem with non-aligned timelines. The staggered times for commissioning and decommissioning ARC systems, while beneficial from continuity of service delivery and technical updates point of view, results to difficulties in aligning required capital and operating budgets  The Government of Canada and ISED have recognized this problem and part of the Alliance's mandate going forward is to 'replace' CFI's MSI operational funding envelope for ARC operational costs so that going forward major DRI capital and operational funding decisions will be done holistically. The operational funding planning needs to also consider costs related to contributed systems.

## Lack of long-term planning due to resource constraints

As discussed earlier in this document, the CCF HQP staff operate the ARC infrastructure and related services at relatively high users to staff ratio compared to e.g., TACC. As such the CCF staff has limited time and resources to focus on long-term future trends and developments, having to focus on short term operational needs. Having HQP staffing that has time and resources to follow future trends and e.g., test and develop new technologies will benefit the whole DRI ecosystem via better leading service delivery, more informed decision making at all levels, and better retention of HQP. More operational funding and HQP hiring is potentially needed. Going forward the Alliance should aspire to reduce day-to-day operational workloads within the staff and allow for more resources towards new service and product development.  Anecdotal evidence indicates that many jobs are not utilizing CPU and GPU resources effectively. Currently the CCF staff does not have time or commitment to work with users to optimize those inefficiencies. Reducing the day-to-day operational workloads for the HQP staff would allow for better system resource utilization and efficiency.

# 5 Relationships between ARC, DM, and RS

Historically, ARC sites have wrestled with data management and used research software. CC has had tiered storage systems for at least 10-15 years because of data management issues and has deployed research software and middleware (e.g., for high-energy physics (HEP) community). Continuous development and improvement due to advances in technology and the evolving requirements of researchers are the story of ARC and supercomputing. Disciplines with long histories of ARC usage like climate, astronomy, and HEP have been building and using RS and RDM for decades (though they didn't necessarily use those terms). RDM and RS communities have become more sophisticated, organized, mature and important over the past decade or so thanks to advances in technology and software as well as the evolving demands of research disciplines and e.g., the big data explosion. New disciplines have huge data needs so that more general and easier-to-use ARC, RS, and RDM solutions are sought after, required, and developed.

Traditionally the relationship between ARC and software development was focused on migrating existing RS codes to new classes of supercomputers, which were more and more parallel, or perhaps even had new types of accelerators. Such efforts included the US funding massively parallel commercial debugger development to improve the usability and ease the porting of research software to the new class of supercomputers. More recently ARC communities have started to put further emphasis on maturing RS and RDM ecosystems and support as they are recognized to be critical components of a viable DRI ecosystem.

A mature RS ecosystem includes professional coding practices, version controls via source code repositories, variety of compiler, driver and library offerings, SW quality assurance, reusability, containerized environments for duplication, documentation and training, accessible science gateways etc. A mature RDM ecosystem enables the Findable, Accessible, Interoperable, and Reusable (FAIR) principles for research data, and Transparency, Responsibility, User focus, Sustainability, and Technology (TRUST) principles for data repositories. Delivering and enabling such RS and RDM systems requires purpose-designed and -built ARC delivery systems, including middleware, web portals and hosting, high-speed networking, and custom mid- and long-term storage solutions for nearline and archival data in addition to corresponding backup storage. Notably many of these components are present in some form today. ARC has always been about continuous improvement in response to advances in technology and evolving user needs and needs within appropriate funding. The needs for RS and RDM are not new per se, but more of a function of available funding historically.

Going forward the Alliance will need to go beyond traditional ARC operations and will include, per its mandate, ARC, RS and RDM holistically in its planning, operations, and funding decisions. Such an approach will provide Canadian researchers a more mature and cohesive DRI ecosystem, including long term efficiencies via open science. Ultimately this will result in direct benefits to all Canadians via new research, innovation and insights driving Canadian economy, competitiveness and policy making.

# Appendix A: Community Feedback on Summer 2021

The ARC Current State report was presented to stakeholder groups and community in the summer 2021 for feedback and overall comments, including CCF HQP professional staff community, CCF regional CEOs, Research Software and Research Data Management working groups, and the Alliance Researcher Council. Most feedback was general in nature and did not uncover any major factual concerns, resulting to only minor wording changes in the final report.

In their feedback the CCF community, RS and RDM WGs, and the Alliance Researcher Council highlighted the following topics:

- Additional data for differentiating between cloud, high-throughput, and HPC (large scale massively parallel jobs with substantial inter-process communication needs) type workloads would be important for future decision making on resource allocation. For example, correlating the research disciplines with the scale of actual job submissions would be an interesting exercise.

- Breakdown of HQP FTE roles would be relevant additional information, e.g., how many of the HQP are in systems administration, v. project management v. research computing support v. cloud computing support etc. roles?

- International comparisons to similar jurisdictions, e.g., Australia would provide valuable insights.

- It would be interesting to do an international comparison to see how the 17% general pick-up rate of CCF ARC systems among Canadian academics, out of which 10% are HHS, compares globally.

- The supply and demand fulfillment rates need to be understood in a larger context. On the face of it e.g., the 40% fulfillment for CPU resources does not sound out-of-line compared to any general approval rate for grant applications. It should be noted that the data presented in the ARC report does not explicitly measure the resource needs that were never even articulated to CCF via RAC due to e.g., acknowledgement by the researchers about the futility or difficulty of getting sufficient resources, researchers already utilizing other non-CFF and international ARC resources etc. In other words, the true demand is much higher than reported by RAC applications alone.

- One must be careful in the scope and validity of projections when using past data to predict the future.

- The authorship, purpose, intended audience, and issuing organization of the document should be stated clearly

- The linkages between data and declarative statements could be made clearer

- The value and importance of CCF HQP should be highlighted more

- Moving forward the Alliance should collect support ticket information regarding which clusters the users are having problems with.

- Diversity of research disciplines should be discussed separately from the issues like gender equity, and regional equity.