

Federal Register Notice 86 FR 46278, <https://www.federalregister.gov/documents/2021/08/18/2021-17737/request-for-information-rfi-on-an-implementation-plan-for-a-national-artificial-intelligence>, October 1, 2021.

---

# Request for Information (RFI) on an Implementation Plan for a National Artificial Intelligence Research Resource: Responses

**DISCLAIMER:** Please note that the RFI public responses received and posted do not represent the views and/or opinions of the U.S. Government nor those of the National AI Research Resource Task Force., and/or any other Federal agencies and/or government entities. We bear no responsibility for the accuracy, legality, or content of all external links included in this document.

Akira Bell  
SVP & Chief Information Officer

**October 1, 2021**

Wendy Wigen

Re: RFI Response: National AI Research Resource

Dear Wendy Wigen:

Thank you for giving Mathematica and other members of the Artificial Intelligence (AI) research community an opportunity to provide input on issues being considered by the National Artificial Intelligence Research Resource (NAIRR) Task Force. A Research Resource is a logical next step to build on the Evidence Act recommendations and infrastructure created by the Open Data initiative. Current cloud technology capabilities can provide the ability to store, analyze, and visualize massive amounts of data. Although AI methods are becoming more accessible, they are not yet ubiquitous and this rapidly changing field has not yet established methods to address equity issues.

For more than 50 years, Mathematica has been at the forefront of assessing the effectiveness of policies and programs to improve public well-being. Our deep bench of more than 1,500 experts translates big questions into insights for our public and private sector partners. We apply our expertise at the intersection of data science and social science by leveraging data assets using advanced technologies and methods such as AI, reusable and scalable data and platforms, and high-performance secure cloud infrastructure. We have a reputation for quality and objectivity rooted in our rigor and commitment to improving well-being for all, which centers on the ethical use of data and equity issues. We offer recommendations for the NAIRR based on our expertise as end users of data and methodologists who have a deep understanding of government policies and programs.

In our experience, good AI requires strong interconnections between the data, technology, and methods. Understanding of the data content, program issues, and questions decision makers need answered are often missing links that cause AI projects to fail. Mathematica's approach to AI provides actionable recommendations by grounding our work in an understanding of the policy issues, pairing the expertise of multidisciplinary subject matter experts with methodological experts, and using data ethically with an emphasis on equity at every stage.

As such, I am pleased to represent my colleagues in our response to **Questions 1A, 1Bii, 1D-1G, 2, 3, 4, and 6**, as published in *Federal Register* no. 2021-15660 by the National Science Foundation and Office of Science and Technology Policy. We share these insights in the hope they help to inform your work and facilitate important discussions and action among federal agencies.

Please direct all follow-up questions or comments to my colleague David Roberts, Mathematica's director of Strategic Communications.

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 2

Mathematica

## **Q1: What options should the Task Force consider and why?**

### **A. Goals for establishment and sustainment of a National Artificial Intelligence Research Resource and metrics for success;**

The NAIRR seeks to increase access to computational resources, data, educational tools, and user support for AI research. As such, goals should include the successful adoption, use, and sustainability of this resource, which requires effective promotion, establishing appropriate use guidelines, and implementing metrics to track usage. To incorporate key indicators as appropriate, the Task Force should define metrics for measuring success as the resource is designed.

The Task Force should develop a strategically targeted communication and outreach strategy to raise awareness among potential users, particularly those disadvantaged historically with respect to resource access to address equity. The NAIRR should require site registration and include tracking metrics for specific resources (data sets, training, and so on) to collect data that demonstrate where uptake has been most successful. The NAIRR can in turn use this information to target communication messages to achieve the objective of supporting a broad range of AI researchers.

To use resources effectively, the NAIRR should promote an approach to AI research founded in methodological rigor, including input from topical subject matter experts (SMEs) and designed with end users in mind. SMEs serve as a critical bridge between the data, technology, methods, and good AI systems by ensuring AI researchers and system developers understand the issues the AI system aims to solve. Similarly, incorporating a human-centered design (HCD) approach will help ensure that end users can practically implement resulting AI products. The NAIRR should also provide guidelines for practitioners related to the responsible use of AI products, such as ways to evaluate equity and bias in the use of AI products.

### **B. A plan for ownership and administration of the National Artificial Intelligence Research Resource, including: ii. A governance structure for the Research Resource, including oversight and decision-making authorities;**

Governance of the NAIRR should include government and private sector organizations to leverage knowledge of government processes and structures with expertise in cutting-edge applications of AI to social science problems. This mix of participants, along with time limits for governance committee members, will help ensure committees remain current on AI industry trends and represent a range of perspectives. Selection criteria for private sector organizations should include companies engaged in AI with public and private sector lines of work across various industries. The NAIRR should also establish safeguards, such as term limits or conflict of interest agreements, to ensure private organizations do not participate for their own gain.

### **D. Capabilities required to create and maintain a shared computing infrastructure to facilitate access to advanced computing resources for researchers across the country, including provision of curated data sets, compute resources, educational tools and services, a user-interface portal, secure access control, resident expertise, and scalability of such infrastructure;**

To maximize access and use among AI researchers, the NAIRR should be a cloud-based portal with a range of computational and analytic resources curated by topic and tailored by experience level—from beginner to more highly trained AI practitioner. Mathematica found this approach was successful for increasing access to our [COVID-19 Curated Data, Modeling, and Policy Resources](#).

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 3

Mathematica

Organizing the breadth of resources and providing ongoing support will require extensive curation and coordination. The NAIRR should establish robust and scalable metadata standards to help organize and maintain resources, such as [Document, Discover and Interoperate](#) standards. Using open-source tools can help control cost, and they have the benefit of substantial community support.

At a minimum, the NAIRR should include several core capabilities:

- Ingesting, cataloging, storing, and archiving data sets
- Version control for data sets and other resources
- Computational resources and visualization tools for working with data
- Configurable environments with the ability to preserve them for future use or save documentation about the configuration
- Monitoring and analysis of system access, resource use, and system security
- User access controls
- Metrics to track system access and use of specific resources by individual users
- Communication tools researchers and SMEs can use to discuss specific data sets or resources either privately or publicly (with discussions accessible to the NAIRR community)

#### **E. An assessment of, and recommended solutions to, barriers to the dissemination and use of high-quality government data sets as part of the National Artificial Intelligence Research Resource;**

Researcher knowledge, access to appropriate tools and resources, and agency publication of high quality data sets can all be significant barriers to disseminating and using NAIRR data sets. To understand and address these issues, the NAIRR should provide a variety of resources:

**Access to SMEs.** SMEs provide critical context for data sets and how to practically address problems using AI. Providing access to SMEs who are knowledgeable about data sets hosted in the NAIRR and related policy issues, as well as supports for communicating with SMEs, will help ensure researchers understand and properly interpret data sets used in AI research.

**Support for multiple data formats.** Providing data sets in multiple formats or conversion tools can increase accessibility among researchers who have limited experience with data manipulation tools or programming languages. AI researchers can also have a broad interest in using government records that are publicly available, but stored as unstructured data sets (for example, emails, *Congressional Record*, and other text-based data). To facilitate the use of unstructured data sets, the NAIRR should provide support for NoSQL databases or data lakes and resources to provision unstructured data sets or structured databases in third normal form (that is, tidy format) whenever possible. The NAIRR could also include semi-automated text analytic tools to facilitate exploring these resources or reference existing supports for the use of unstructured data.<sup>1</sup>

**Tools for managing data access.** Some NAIRR data sets might be static, but others could reside elsewhere or be available only through live data streams. Providing support to access live data streams or data sets housed outside the NAIRR, particularly using application programming

---

<sup>1</sup> Existing examples of resources to support the use of unstructured data include <https://www.ibm.com/cloud/blog/structured-vs-unstructured-data> and <https://www.essentialsql.com/database-normalization/>.

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 4

Mathematica

interfaces (APIs), is critical and can provide the portal an advantage over other methods of accessing data. The NAIRR can apply best practices for data sharing as outlined by the [HHS CDO initiative](#).

**Data matching guidance.** Different data sets might code similar constructs inconsistently and provide hard-to-follow data documentation, making it difficult for users to combine data sources. Promoting standards in data use guides for NAIRR data sets and providing guidelines for combining data sets can help address usage barriers related to users' knowledge. Encouraging standard, detailed data set documentation can also alert AI researchers to potential issues related to responsible use and equity associated with data quality issues often introduced during processing.

**A data operations framework.** A framework (such as DataOps) can define an approach for designing and implementing the NAIRR portal in addition to outlining best practices, recommended workflows, and use of guidelines for AI researchers. Providing a framework such as this can support access to data for use in AI research by removing the requirement for researchers to have the knowledge to develop and implement workflows and follow best practices on their own.

**Lessons learned from data-sharing initiatives.** The NAIRR should consider lessons learned from existing data-sharing and open data initiatives about what works well and potential pitfalls. For example, incorporating incentives for sharing data, providing support for agencies to share data in requested ways, developing channels for data users to provide feedback (such as ratings and data use popularities) and suggestions on shared data, and establishing clear lines of responsibility and ownership are all important.

#### **F. An assessment of security requirements associated with the National Artificial Intelligence Research Resource and its management of access controls;**

The Task Force must consider requirements such as security of the NAIRR portal, integrity of the resources housed there, and user-level access controls. The Task Force should tailor management strategies based on the level of risk to balance security with accessibility. These strategies should align with [federal thinking on cybersecurity](#) while providing flexibility and ease of access for the broad range of AI researchers the resource is intended to support. A [Zero Trust Architecture](#) approach best suits this type of resource: the focus must be on the data and the level of protection required based on the risk profile of the data, and the purposes the data serve. Real-time, contextual details that explicitly describe each user should be the basis for granting access. For example, the potential risk associated with access to public data is low and should not require substantial resources. By contrast, the integrity of the data sets and guidance documentation, as well as availability of computational resources, pose a higher level of risk that warrants greater investment.

#### **G. An assessment of privacy and civil rights and civil liberties requirements associated with the National Artificial Intelligence Research Resource and its research;**

Privacy, civil rights, and civil liberties inherently intertwine in AI research. Multiple stages of the research can introduce privacy and bias issues—from selecting to cleaning data and developing and implementing the model. Public use datasets are reviewed for risk of individual disclosure and often remove personal identifiers such as race, gender, and ethnicity. However, combining data sets from disparate sources could result in reidentification of individuals based on new field combinations. The NAIRR can holistically assess the data sets it hosts to ensure possible end users cannot maliciously or inadvertently reidentify individuals by combining data sets. [Standards for deidentifying data sets](#) provided for AI research might need to be more stringent. The NAIRR should also establish a disclosure risk committee to evaluate potential risk based on the combination of data sources

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 5

Mathematica

available. It might also consider a record identification system that can facilitate matching across sources and allow for removing fields that could result in reidentification.

To address bias that infringes on civil rights or civil liberties, the NAIRR should provide guidance to researchers on how to address algorithmic bias, such as a process to evaluate how resources might reinforce inequity and provide responsible use guidance to researchers. This guidance should focus on transparency in algorithms so researchers can understand what the models do and how they reach conclusions. If data sets available on the NAIRR include sensitive data such as race or ethnicity that could potentially lead to infringement of civil liberties, the NAIRR should implement more stringent requirements for adhering to provided guidelines.

## **Q2: Which capabilities and services provided through the NAIRR should be prioritized?**

The goal of establishing a national computing and collaboration environment for AI researchers and students is ambitious and timely. To realize this vision, the NAIRR should prioritize several factors:

**Scalable, simple, and low-cost infrastructure.** There are many options for relatively cheap, secure, easily provisioned cloud-based infrastructure with free tiers for students and researchers to use. The NAIRR could provide access to Amazon Web Services (AWS), Azure, and/or Google Cloud Platform resources directly through the portal, with subsidies or expanded free tiers for students and researchers whose research serves a public mission. This would enable the NAIRR to be stood up more quickly and allow AI researchers to take full advantage of cloud-based services.

**Access to large data sets.** The NAIRR should provide general access to curated public use files (PUF), as well as access to restricted use files (RUF) with secure controls. It should also support standardized APIs that authorized end users can query programmatically so they can seamlessly integrate their services and data sets into existing data pipelines. Although the NAIRR should incentivize open access to data, proper vetting of third-party API providers is essential to reduce the risk of malicious or unauthorized use or provision of data. Potential users should present compelling evidence for why they need access to APIs with sensitive data. In addition, the NAIRR could provide access to larger real-world and so-called toy data sets researchers could use to follow along with training materials to help students and aspiring AI researchers learn big data processing and analysis techniques.

**Coordination with SMEs.** The NAIRR should provide access to knowledgeable SMEs who can help students and researchers understand and properly interpret data sets. These SMEs should include agency personnel with deep expertise in NAIRR data sets, as well as outside volunteers with domain expertise. Experienced AI researchers can provide insight on the analytical approaches, and SMEs can help aspiring and seasoned data scientists identify domain-specific problems, understand the intricacies of complex or idiosyncratic data sets, and ensure models are developed with end users in mind. The NAIRR should also provide communication channels—through video conference software, Slack channels, communities of practice, or other virtual communication mechanisms—to enable users of the platform to reach out to SMEs and confirm they properly understand and interpret their findings. Unfortunately, curating data sets and making SMEs available can be costly, time consuming, and resource intensive. To address these resources constraints, the NAIRR should give certain users of the platform—such as students and researchers working on high-priority government contracts—privileged access to SMEs and expedited access to curated data sets.

**Technical documentation, training, and tutorials.** The NAIRR should pair access to data sets and SMEs with technical documentation, tutorials, and standardized metadata to empower users. The portal should include (English and non-English) materials and tutorials written for multiple

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 6

Mathematica

free-to-use programming languages—including Python, R, Julia, C#, and Java. In addition, the portal's design should integrate with popular version control and containerization services (such as GitHub, Bitbucket, Docker, and Kubernetes). Adoption of the NAIRR is likely to be higher the more it prioritizes popular open-source programming languages and tools. Using legacy languages or not integrating with tools commonly used in the industry will lower credibility and push people toward other platforms.

**Human-centered design.** Researchers often view AI projects as a technology or data science effort aimed at achieving a certain level of accuracy or technical rigor, rather than an effort to improve human decision making. However, to be valuable, practitioners must adopt and use AI to drive improvement and change. As part of its mission, the NAIRR should promote incorporating HCD at the earliest stages of an AI project to properly frame analytic problems; develop solutions that deliver results to the right people at the right time; and provide wraparound support in the form of explainability, transparency, and user support. Mathematica's work on the [Centers for Medicare & Medicaid Services Artificial Intelligence Health Outcomes Challenge](#) reflects this mindset. Our team implemented an HCD process with multiple rounds of discovery, research, and user testing with patients, doctors, and other clinicians to accomplish these goals. The result was analytic output that could be explained and interpreted easily by clinical end users, as well as actionability by integrating with established workflows to create a seamless user experience that turned outputs into solutions.

### **Q3: How can the NAIRR and its components reinforce principles of ethical and responsible research and development of AI, such as those concerning issues of racial and gender equity, fairness, bias, civil rights, transparency, and accountability?**

The NAIRR and its components should establish clear and actionable ethics standards, document potential bias and nuance in resource manuals, set guidelines to enhance appropriate use of the resource, and protect data privacy. Specifically, the NAIRR should consider the following aspects to reinforce ethical and responsible AI research and development:

**Create measurable and reportable ethics standards and rubrics.** When developing ethics standards and rubrics, the NAIRR can increase awareness and discussion of work done by advocacy groups (such as the [Algorithmic Justice League](#)), media outlets (for example, [ProPublica](#)), academic institutions, and other organizations focusing on responsible use of AI and algorithmic bias. Equipping interactive AI ethics tools in the ethics rubrics will make the standards more actionable. For example, Mathematica developed an [open-source tool](#) that enables users to assess fairness in algorithms predicting binary outcomes and quantify the uncertainty of each estimate under a Bayesian framework.

**Assess potential bias in NAIRR data sets and models and provide considerations in use documentation.** The NAIRR should assess the data sets it makes available through the portal for underlying biases, and it should clearly articulate any potential biases to potential end users. Sometimes it might be difficult to discern the types of biases that underly a given data set, especially when working with large or unstructured data. For instance, natural language models trained on text data might reflect the unconscious biases of the original authors. Furthermore, the NAIRR should assess predictive models developed using NAIRR for their proxy discrimination before applying to other applications within the portal. Complex predictive models can often obfuscate prejudiced conclusions by using features that might not seem to correlate with race, ethnicity, and so on, but actually do. AI developers should assess their models for proxy discrimination when validating their results. The NAIRR should require documentation of such nuances in the data user manual and

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 7

Mathematica

keep updating the document to incorporate feedback from data users on underlying data biases through a crowdsourcing mechanism (such as a user-reported rating and feedback system).

**Require proposals to address and account for potential bias when developing AI tools using NAIRR resources.** The proposal should require researchers to assess biases that might affect the data sets they want to use to develop a new AI tool and outline their plan to account for these biases when developing the new tool. For example, if a developer planned to use Medicare claims data to develop an AI prediction tool, it would be important to consider regional variations in health care service use and the underlying racial biases in health care delivery, both of which could introduce substantial bias into prediction tools. AI developers should also consider whether implementing or using the tool they propose to develop could adversely affect historically disadvantaged priority groups. For tools with the potential to adversely affect certain groups, it will be important to assess the actual impact on those subgroups when they implement the tool. When validating a new AI tool, developers should provide information on the tool's performance for relevant subgroups. For example, developers should provide AI model performance metrics, such as the Area Under Curve (AUC) score for subgroups based on gender, race and ethnicity, and age, among other individual characteristics.

**Establish data report agreement to enhance transparency in AI development.** If AI developers receive access to a data set or multiple data sets for the purpose of developing an AI tool, the developers should have to provide information regarding how representative their training, testing, and validation data sets are of the population to which they hope to apply the AI tool. They should present such information for data sets directly used for AI model training, testing, and validation, rather than the original data sets from the NAIRR. For example, if a developer excludes certain variables because of data quality issues or other reasons in the process of cleaning the data to build an AI tool, the final cleaned data might not be representative even if the initial data were.

**Deidentify data when possible to ensure disclosure avoidance.** The governance and oversight structure for the NAIRR must account for different types of data that impose different levels of risk to people or groups represented by those data. Because it is unlikely obtaining consent for use of particularly sensitive data—such as health or financial data—will be feasible, the NAIRR should deidentify these data when possible. Moreover, the NAIRR should carefully consider data privacy issues that are introduced when providing access to many data sets. For example, when using multiple linked data sets, researchers must still protect the confidentiality of deidentified information, checking by differentially private algorithms to ensure disclosure avoidance.

**Provide disclosure guidelines or requirements for AI tools developed using NAIRR resources.** AI developers should make enough information available about an AI tool to enable a potential user to determine whether it is appropriate to use in a specific population or setting, to understand the performance of the tool, and to understand how to interpret prediction results. Such information includes the data used to develop the AI tool, the approach to develop and validate the tool, how the algorithm reaches its results, and information regarding overall AI model performance and its performance for specific subpopulations. Transparency does not always require that AI developers make an AI algorithm itself available though. The level of transparency should depend to some extent on the level of risk associated with the use of the AI tool, the level of precision of the predictions, the clarity of the recommended actions to end users, and the potential for legal liability ([NAM AI report](#), pp. 192, 219–220).

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 8

Mathematica

#### **Q4: What building blocks already exist for the NAIRR, in terms of government, academic, or private-sector activities, resources, and services?**

The NAIRR can benefit from government open data initiatives, evolving interoperable and digital information, technology offerings from government contractors and cloud service providers, and established private–academic partnerships. The NAIRR can use several types of building blocks:

**Existing open data initiatives from the federal government and agencies.** For example, [Data.gov](#) includes more than 300,000 data sets and rich data tools, data incubator sources, and skills development resources that enable citizen participation in government, create opportunities for economic development, and inform decision making in both private and government sectors. The [Data Optimization Initiative](#), led by the Office of the Chief Technology Officer within the U.S. Department of Health and Human Services (HHS), has released more than 4,500 data sets collected by HHS for public use by researchers and entrepreneurs.

**Health-related AI research that incorporates evolving interoperable and digital information.** HHS is in the process of advancing the connectivity of electronic health information and interoperability of health information. For example, the electronic health information exchange enables doctors, nurses, pharmacists, other health care providers and patients to appropriately access and securely share a patient’s vital medical information electronically, improving the speed, quality, safety, and cost of caring for patients. Future health-related AI research developed using the NAIRR could apply and incorporate information from such data resources.

**Existing technology offerings and data sources from government contractors.** Many government contractors and research institutes have invested in technology offerings and curated data sources that could support specific needs of NAIRR users. Specifically, many of them operate virtual data enclaves (VDEs) that provide secure access to restricted data through a virtual private network connection to a portal on researchers’ computers. For example, the data library recently developed by Mathematica is an AWS cloud solution that enables users to search for, discover, and interact with ingested data sets in a secure and compliant environment. NORC offers a data enclave that provides data storage and management, computational resources, and reporting tools. The NAIRR could also benefit from existing, curated data sources. For example, Mathematica surveyed and organized a [repository of publicly available data, modeling, and policy resources on COVID-19](#) that enables states, health care decision makers, providers, and others to predict need and direct resources based on the best available evidence during the pandemic.

**Strong private-academic partnerships.** Partnerships with academic institutions that foster AI research can expand the reach of the NAIRR and build a pipeline of future contributors. For example, the [Howard-Mathematica SICSS partnership](#) is an instructional program hosted at Howard University, a historically black college and university (HBCU), designed to promote learning and support development of expertise in computational science for graduate students, postdoctoral researchers, and beginning faculty from HBCUs and underrepresented communities. This program highlights how innovative partnerships can help counter anti-Black racism and inequity. The program includes lectures, group problem sets, participant-led research projects, office hour sessions led by industry professionals, and invited outside speakers who conduct computational social science research in a variety of settings. Topics covered included text as data, website scraping, digital field experiments, machine learning, and ethics. Partnerships such as this introduce new views on existing applications of AI that those who are not members of underrepresented communities might not always clearly understand, such as limitations of natural language processing for African American

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 9

Mathematica

vernacular English, AI tools for health care practitioners to eliminate bias for rural or unrepresented populations, and AI studies of network effects for immigrants.

**Cloud service provider collaboration.** NAIRR can work with cloud computing providers for flexible and scalable resources, while benefitting from the volume of services offered through these providers. Many organizations already partner with these cloud providers to develop AI solutions to address social needs and improve public well-being. For example, [Mathematica is collaborating with Google Cloud](#) on Google Cloud's new Healthcare Data Engine to provide health care measures, data analytics, and data science capabilities to health care teams and leaders across the country.

#### **Q6: Where do you see limitations in the ability of the NAIRR to democratize access to AI R&D? And how could these limitations be overcome?**

The past decade has seen a dramatic increase in private sector investment in AI research and development, both in the United States and throughout the world.<sup>2</sup> These commercial investments have often had a democratizing effect on AI research, particularly through corporate participation in open-source and open data communities. However, many challenges remain to ensure AI resources are widely available and researchers with varying experience levels have the capabilities to develop and deploy AI. Among the core obstacles to broader democratization of AI resources are access to (1) adequate training and expertise; (2) large data sets, especially private or industry-specific data sets;<sup>3</sup> and (3) high-performance computing infrastructure necessary to use many state-of-the-art AI algorithms and data processing techniques.

**Access to training and expertise.** AI research often requires advanced knowledge of computer science, as well as a strong foundation in mathematics such as calculus, probability and statistics, and linear algebra. This can create barriers to entry for aspiring AI researchers or those coming from fields other than science, technology, engineering, or math (STEM) who are less well versed in these concepts. If not paired with adequate training, educational resources, and subject matter expertise, AI democratization can lead to inaccurate statistical interpretations, developing biased models, and programming or technical errors that produce incorrect results. The NAIRR should provide training resources that emphasize developing data science and computer science skill sets, proper data preparation and modeling procedures, and accurate interpretation of model performance and results. A central goal of the NAIRR's training efforts should be to reach a broader and more diverse talent pool of prospective AI researchers coming from colleges and universities, technical institutes, and other non-traditional learning settings. To advance this goal, the NAIRR could also consider establishing certificate programs for aspiring technologists that would be more easily attainable for those with fewer resources at their disposal. Finally, the NAIRR should provide channels of communication with SMEs who can support students and researchers on an as-needed basis.

**Access to data sets.** Beyond training and expertise, access to data—particularly large volumes of data needed to train highly accurate machine learning models—is essential. Although private sector investment has helped to advance many AI tools and procedures, there are often limits to sharing data in the private sector. A key to democratizing AI research for the NAIRR will be providing easy access to large data sets while maintaining proper data security and role-based authentication

---

<sup>2</sup> Zachary Arnold. "Tracking AI Investment Initial Findings from the Private Markets." September 2020. Available at <https://cset.georgetown.edu/publication/tracking-ai-investment/>.

<sup>3</sup> Ahmed, Nur, and Muntasir Wahed. "The De-democratization of AI: Deep Learning and the Compute Divide in Artificial Intelligence Research." *ArXiv*, October 2020. Available at <https://arxiv.org/abs/2010.15581>.

**To:** Wendy Wigen  
**From:** Akira Bell  
**Date:** October 1, 2021  
**Page:** 10

Mathematica

procedures. Ideally, the NAIRR platform would provide access to both government and private sector data and provide tools for greater integration between the two.

**Access to computing infrastructure.** Another core challenge to democratizing AI research is unequal access to the computing resources required to process large volumes of data. NAIRR can address this issue by providing tools and infrastructure tailored to experience levels, along with clear guidance and user documentation. Tools should be available at little or no cost to ensure accessibility, but the NAIRR must develop a prioritization and approval process to ensure specific users do not hoard resources or use them for unlawful purposes. In addition, the NAIRR will have to provide a help desk support function to respond to users' issues and requests.