

# Earth and Space Science



## COMMISSIONED MANUSCRIPT

10.1029/2020EA001562

### Key Points:

- Open science is a collaborative culture enabled by technology that empowers the open sharing of data, information, and knowledge within the scientific community and the wider public to accelerate scientific research and understanding
- This study provides a synopsis of various open science activities occurring throughout the community and synthesizes those activities around three broad open science focus areas
- Science has become increasingly data driven, and data programs now play a critical role in enabling and accelerating open science
- Data programs can set strategic policies and directions that are critical to enabling and promoting open science

### Correspondence to:

R. Ramachandran,  
[rahul.ramachandran@nasa.gov](mailto:rahul.ramachandran@nasa.gov)

### Citation:

Ramachandran, R., Bugbee, K., & Murphy, K. (2021). From open data to open science. *Earth and Space Science*, 8, e2020EA001562. <https://doi.org/10.1029/2020EA001562>

Received 16 NOV 2020  
Accepted 13 APR 2021

Published 2021. This article is a U.S. Government work and is in the public domain in the USA. *Earth and Space Science* published by Wiley Periodicals LLC on behalf of American Geophysical Union.

This is an open access article under the terms of the [Creative Commons Attribution License](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

## From Open Data to Open Science

Rahul Ramachandran<sup>1</sup> , Kaylin Bugbee<sup>1</sup> , and Kevin Murphy<sup>2</sup> 

<sup>1</sup>NASA/MSFC, Huntsville, AL, USA, <sup>2</sup>NASA/HQ, Washington, D.C., USA

**Abstract** The open science movement continues to gain momentum, attention, and discussion. However, there are a number of different interpretations, viewpoints, and perspectives as to what the term “open science” means. In this study, we define open science as a collaborative culture enabled by technology that empowers the open sharing of data, information, and knowledge within the scientific community and the wider public to accelerate scientific research and understanding. As science has become increasingly data driven, data programs now play a critical role in enabling and accelerating open science.

In this study, we describe specific actions that data programs can take to make the open science paradigm shift a reality. These actions range from implementing open data and software policies to reimaging data systems that move data out of organizational silos and into cyberinfrastructures that enable efficient research processes and accelerate knowledge dissemination. There are still a number of obstacles to be overcome by data programs which range from mitigating the risk of open data misuse to overcoming the inertia of legacy systems. Data programs need to support open science through the thoughtful development of open policies, systematic investment in innovative and collaborative infrastructures, and the promotion of cultural change. On the other hand, individual researchers play an equally important role by serving as advocates for open science principles and by adopting a number of best practices outlined in this study. By working together, a new and more open age of scientific research can be achieved to benefit science and society.

## 1. Introduction

Open science, as both a concept and a term, is increasing in popularity and usage. However, definitions, interpretations, and perceptions as to what the term “open science” means varies. Some definitions are fairly narrow and only focus on providing more open access to science as a body of knowledge. These narrow definitions place an emphasis on openly sharing scientific knowledge as early as possible in the research process (University of Cambridge, 2020). On the other hand, broader definitions of open science acknowledge that science is both a body of knowledge and a systematic method for thinking. Broad definitions place an emphasis on encouraging a culture of openness (Bartling & Friesike, 2014) that includes the entire process of conducting science (National Academies of Sciences, Engineering, & Medicine, 2018a, 2018b) and encourages open collaboration and access to knowledge (Vicente-Saez & Martinez-Fuentes, 2018). In its broadest definition, the term “open science” refers to a paradigm shift in how the methods of science are conducted. This expansive vision of open science acknowledges that rapid technology changes, primarily driven by the Internet, may enable a second scientific revolution that fundamentally changes research methods and standards across science.

To complicate matters, the term “open science” is sometimes used interchangeably to represent various principles that support the broader idea of open science itself. These principles include ideas such as open data, open source software, open journal access, and reproducibility. For example, reproducibility, or the ability to verify another scientist’s results, is enabled by the principles of open data, open code, and transparent methodologies, yet reproducibility itself is not equivalent to open science.

While open science definitions are variable and ambiguous, the value of open science as both a concept and a paradigm change is accepted by the majority of the scientific community. Open science not only benefits the scientific endeavor itself but has also been shown to benefit individual researchers through increased citations and media attention, a larger collaborative network, and exposure to new career and funding opportunities (McKiernan et al., 2016; Murphy et al., 2020). However, in order for researchers, organizations,

and programs to more effectively foster and enable open science, a workable definition of open science is needed.

Therefore, we define open science as a collaborative culture enabled by technology that empowers the open sharing of data, information, and knowledge within the scientific community and the wider public to accelerate scientific research and understanding. This definition asserts that the process of open science is inherently collaborative and functions at its best when diverse backgrounds, perspectives, and expertise are included. This collaborative process, enabled by advances in technology, spurs the goal of openly sharing data, information, and knowledge to an ever-growing audience of both the scientific community and the public at large. The success of the open science endeavor is, therefore, measured in a number of ways including the accelerated scientific research, broader scientific literacy, and increased diversity throughout the process. This vision of open science converges around three overarching dimensions: (a) increasing the accessibility to the scientific process and the corresponding body of knowledge; (b) making both the research process and knowledge sharing more efficient; and (c) understanding and assessing scientific impact through innovative new metrics.

The contributions of this study are threefold. First, the study provides a summary of the drivers behind open science and an overview of the current state of various open science activities across the community. This overview is meant to concretely define the scope of open science activities so that organizations may take actionable steps in order to operationalize open science principles. Second, we posit that science is increasingly data driven, placing data programs and the organizations in charge of overseeing the life cycle of scientific data in a critical position to enable and accelerate open science. In this study, we describe specific actions, such as developing technologically innovative data systems and evolving data stewardship practices, that data programs can take to make the open science paradigm shift reality. In addition, we detail a number of challenges that will need to be addressed by the community in order to advance open science. Third, we provide a number of best practices that may be adopted by both individual researchers and broader data programs in order to make open science a reality.

## 2. Drivers for More Open Science

While science, in some measure, have been open since the development of the journal publication system in the mid-1600s (Bartling & Friesike, 2014), three key drivers are accelerating the open science movement and redefining what it means to be open. The three factors driving this movement are technology advances, the rapid growth in data volume and variety (Xu & Yang, 2014), and the increasing interdisciplinary nature of the science questions being tackled. These three factors are not independent but rather are profoundly intertwined, with data and technology being especially synergistic.

Rapid technology advances are fundamentally changing how science is conducted and in tandem are accelerating the adoption of open science principles. Technology has enabled new workflows that not only make the process of science more efficient, but also create new mediums for sharing knowledge with other scientists and the broader community. Furthermore, technological innovations are continually evolving the scientific process itself. New collaboration technologies make sharing ideas, data, algorithms, software, and experiments easier (Friesike et al., 2015), while new software tools are now available that quickly incorporate the latest improvements in algorithms and analysis methods such as machine learning libraries (Woelfle et al., 2011). In addition, researchers have access to better, more cost-effective computational power, more substantial and affordable storage via technologies such as the cloud, and faster networks. These changes in technology also allow broader participation in the scientific process, making it possible to successfully harness the public's participation through various citizen science activities (Newman et al., 2012).

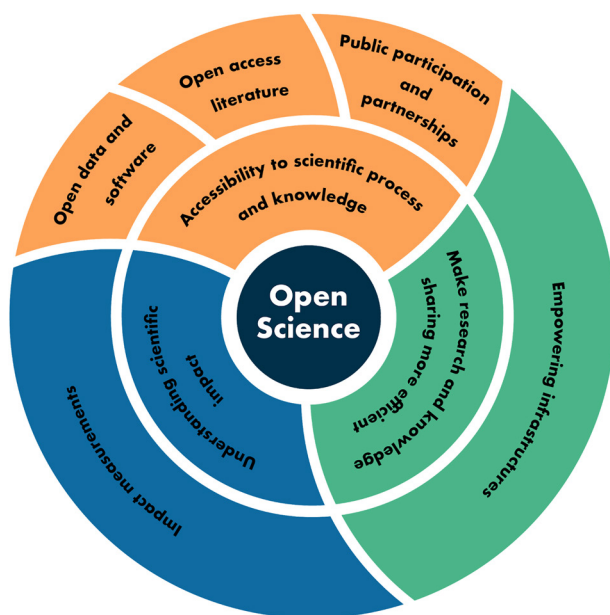
Technology has changed the way scientists communicate and access information. Until recently, scientific knowledge dissemination was controlled by the major journal publishers within each field. This centralized means of communication was built upon a 17th-century publication model that was intended to make science more open at the time. However, a major disruption to this model occurred with the advent of the Internet, where scientists publish thoughts, ideas, results, and conclusions openly and to everyone (Laakso et al., 2011). While trusted journals remain the primary medium for archiving and sharing peer-reviewed scientific knowledge, valuable information is also now available from these nonjournal sources. These “gray

literature” sources include reports, blogs, articles, and various publications produced outside of the traditional commercial and academic publication workflows (Schöpfel & Prost, 2016). Nontraditional communication sources allow scientists to share initial results and lessons learned in real time and are challenging the prevailing belief that scientific results must be completely validated and verified before publication (De Roure et al., 2010). The details of results and analysis are also augmented by these new communication channels, which allow scientists to more thoroughly describe the development of algorithms and source code used to generate results.

Rapid technology advances have increased the volume and variety of data by dramatically improving the instrumentation for observing and collecting data, the numerical models for simulations, and the processing abilities to efficiently analyze data. This rapid increase in the volume, velocity, and variety of data disrupts the scientist’s traditional analytic workflows and the corresponding data management practices required for working with data. The exponential growth in data volume and complexity make science more reliant on complex computational platforms called cyberinfrastructures. Cyberinfrastructure, a term used by the US National Science Foundation (NSF) (Cyberinfrastructure Council, 2007), is defined as an infrastructure consisting of computing systems, data storage systems, advanced instruments, data repositories, visualization environments, and people all linked together by high-speed networks, making scholarly innovation and discoveries possible.

The increasing accessibility and availability of these data have opened opportunities to tackle interdisciplinary science problems that span domain boundaries. Solving these interdisciplinary problems requires collaboration across traditionally siloed scientific communities and the convergence of different types of expertise, knowledge, and resources (Chesbrough, 2015). This focus on interdisciplinary research has led to a shift from individuals conducting research to a team approach with each member providing specialized expertise. Science teams are frequently including programmers and computer scientists to help conduct analyses and to optimize algorithms for efficient analysis of large volumes of data.

The three open science drivers made the Human Genome Project (HGP), an exemplar of the open science movement, possible. The HGP was a 15-year program that set out to address a complex science problem: mapping and sequencing the human genome (Hood & Rowen, 2013; Watson, 1990). In order to address such an intricate problem, the HGP was highly collaborative in nature and included international participation from 20 groups from the United States, the United Kingdom, Japan, France, Germany, and China (International Human Genome Sequencing Consortium, 2001). Not only did the HGP promote international collaboration but cross-disciplinary collaboration was also encouraged between computer scientists, mathematicians, engineers, and biologists in order to make the needed advances in computational and mathematical approaches (Collins et al., 2003; Hood & Rowen, 2013). The generation of large data volumes, distributed across organizations and needed in a timely manner, required the open sharing of data. To support this open sharing ideal, the HGP drafted the 1996 Bermuda Principles, which committed the HGP laboratories to openly sharing DNA sequencing information on a daily basis (Cook-Deegan & McGuire, 2017). This open sharing of data led to new insights for biologists including the discovery of new information on 30 disease genes (Arias et al., 2015; Toronto International Data Release Workshop Authors, 2009). The HGP also upheld open-source principles by making the software needed for analysis openly available (Hood & Rowen, 2013). Most importantly, the HGP has had lasting scientific impacts on biology and medicine. The HGP led to the emergence of proteomics, a discipline focused on identifying and quantifying the proteins present in discrete biological compartments (Hood & Rowen, 2013), has advanced scientists understanding of evolution, and initiated the comprehensive discovery and cataloging of a “parts list” of most human genes (Hood & Rowen, 2013). In addition, the HGP spawned a number of new scientific projects including the HapMap Project to catalog human genetic variation (The International HapMap Consortium, 2005) and the ENCODE project (Encyclopedia of DNA elements) to understand the functional parts of the genome (Hood & Rowen, 2013; The ENCODE Project Consortium, 2011). The HGP’s advances in open data sharing policies, open source code for analysis, technological advancements, and an effective international collaboration model provide a benchmark for open science in the modern era.



**Figure 1.** The open science concept represented as layers. The concept of open science is found in the center. The three open science focus areas are represented in the middle layer, while the outer layer represents data program-specific strategies that enable open science.

### 3. Open Science Focus Areas

This section provides an overview of the current state of various open science activities across the community, in order to explicitly scope those activities for data programs. The activities are categorized into three broad focus areas as shown in Figure 1.

Figure 1 illustrates the three broad focus areas of open science: (a) increasing the accessibility to the scientific process and the corresponding body of knowledge; (b) making both the research process and knowledge sharing more efficient; and (c) understanding and assessing scientific impact through innovative new metrics.

#### 3.1. Accessibility to Science

The accessibility to science focus area places an emphasis on providing access to science as a way of thinking but also on providing wider access to science as a body of knowledge.

##### 3.1.1. Science as a Way of Thinking

The scientific method is a disciplined way of thinking and conducting research. In the past, access to the scientific research process was perceived to be limited to scientists with specialized degrees. While there have always been amateur scientists contributing to the scientific endeavor, the digital age has made it easier for more nonscientists to participate in research. Encouraging broader and more inclusive participation in all

stages of the scientific research process is a key aspect of emerging open science activities. “Citizen science” is a commonly used term to describe the participation of nonscientists and amateurs in research (Fecher & Friesike, 2014). Citizen science activities leverage a volunteer workforce in a scientifically meaningful manner via activities such as systematically collecting data or analyzing data to discover interesting patterns. Citizen science activities are gaining momentum in use by governments and other organizations as a way to address scientific and societal challenges (Shanley et al., 2019), therefore making these activities an emerging yet significant component of open science.

Just as including amateurs in the scientific process is important, making research results understandable and comprehensible to the general public are equally important. To make science more accessible, open science efforts offer new avenues, tools, and formats of science communication beyond the traditional journal model. These efforts include science writing that targets a broader community by deconstructing interesting yet complex research results into easy-to-understand pieces of information. Science blogs, such as Dan’s Wild Wild Science Journal (Satterfield, 2020), are one effective mechanism for communicating interesting scientific results. Social media platforms are also proving to be an effective communication tool for a broader audience. Astrophysicist Neil deGrasse Tyson has around 14 million followers on Twitter, while physicist Brian Cox and biologist Richard Dawkins have around 3 million followers each. While blogs and social media are wide-reaching dissemination platforms, these tools are only as effective as the scientists who are willing to devote time and energy to creating these types of content. For this reason, there is a growing need for people, organizations or groups within the scientific community to establish themselves as credible boundary spanners, or mediators, between scientists and nonscientists (Safford et al., 2017) for effective communication.

##### 3.1.2. Science as a Body of Knowledge

Equitable access to the scientific body of knowledge is an essential dimension of open science. The scientific body of knowledge is the products of research and includes data, software, research publications, and other supporting materials (Fecher & Friesike, 2014). Equitable access to these objects is enabled through open data policies, open source software principles and open access literature.



### 3.2. Open Data

Data drive the scientific process in two ways. First, data are products of research activities and are key components in the scientific body of knowledge. Second, data are analyzed for scientific insights and results. Since data are essential to the scientific process, open science efforts have focused on making data more openly available. Open data are data that may be accessed, used, and shared for any purpose without restrictions. Open data policies, typically developed by government and commercial organizations, define what data will be shared, with whom, at what price, and under what conditions the data can be reused or redistributed (Borowitz, 2017). Data sharing policies fall on a spectrum of openness with the most open data being made fully available either free of charge or at no more than the cost of reproduction (Group on Earth Observations, 2020; Open Knowledge Foundation, 2020). On the other end of the spectrum, data may have limited or restricted access due to security concerns, the inclusion of personally identifiable information (PII), or licensing agreements often associated with commercial data purchases.

Open data benefit the open science movement in a number of ways. First, open data policies prevent duplicating the collection of data across organizations, freeing up resources to amass a more diverse array of data and making it possible to have a more comprehensive record of observations. For example, data exchange agreements between NASA and the European Space Agency (ESA) have made a virtual constellation of observations from the Landsat series and Sentinel-2 possible. Combining data from these two platforms increases the frequency of observations over land which is essential to land monitoring applications research. Second, open data policies significantly increase data use and reuse, especially when data are made freely available. The Landsat free and open data policy represents the epitome of a successful open data policy. After making Landsat data freely and openly available in 2008, the USGS saw a 20-fold increase in data downloads from 2009 to 2017 and a 4-fold increase in the use of the data in the annual number of publications (Zhu et al., 2019). Providing open access to Landsat's long-term record of observations allowed scientists to move from only analyzing single images to conducting time-series analyses, enabling advances in a number of applications including monitoring land surface changes, tracking rates of change in shoreline erosion and measuring glacial fluctuations (Kennedy et al., 2014; Roy et al., 2014; Wulder et al., 2012).

More broadly, there are a number of economic and societal benefits to open data. Remote sensing data have been of particular use to economists primarily because of its high spatial resolution, its wide geographic coverage, and its ability to offer access to information unavailable by other means (Donaldson & Storeygard, 2016). Remote sensing data have been used in a variety of economic contexts including agriculture, infrastructure investments, tourism, resource availability, and insurance (Donaldson & Storeygard, 2016). Additionally, there are a number of societal benefits to providing open data. Remote sensing data aid in disaster mitigation, response and recovery, and is also a valuable input into monitoring conflicts, illegal activities, pollution events, and the effects of policies on land use (Donaldson & Storeygard, 2016; Zhu et al., 2019).

### 3.3. Open Source Software

Software, in combination with data, acts as a tool for new knowledge discoveries and insights and often serves as a representation of knowledge in and of itself (Keyes & Taylor, 2011). However, software, unlike data, is protected by copyright, making the free use of software restricted unless the copyright owner has granted a license. Software is a broad term that applies to computer programs, applications, and source code that provide a certain level of utility to users or assist in producing a result (Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division On Engineering and Physical Sciences, 2018). Software is typically used to either conduct scientific analyses or to provide the supporting infrastructure needed to manage data. In order to be considered open source, these software must be made publicly available and include a software license that grants permissions for anyone to examine, use, change, and distribute the source code for any purpose (Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division on Engineering and Physical Sciences, 2018). Source code without a license is considered "all rights reserved", and is therefore not available for free and open use.

Open source software culture is slowly gaining momentum within the sciences as the community recognizes the benefits and impacts of making software openly available. First, open source software encourages software reuse, the benefits of which include reducing the time working with data, lowering duplication of effort, enhancing the use of open data, and ensuring the longevity of code (Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division on Engineering and Physical Sciences, 2018). Open source software also makes teamwork easier across organizations and allows scientists to have a larger collaborative network. Lastly, open code enables scientific and computational reproducibility and transparency. Scientific reproducibility is elusive without access to open code (Gil et al., 2016), and publications with only natural-language descriptions of methods, algorithms, and code implementation are often insufficient for reproducibility (Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division on Engineering and Physical Sciences, 2018). Computational reproducibility is even more challenging to achieve, but new, forward-thinking approaches such as reusable research objects, executable publications, and containers are making computational reproducibility easier (Chard et al., 2020; Gil et al., 2016; Koop et al., 2011; Ton That et al., 2017).

### 3.4. Open Access Literature

Journal articles have been and continue to be a key mechanism for sharing scientific results with the broader community. However, access to journal articles has been limited by the prohibitive cost of accessing a large number of journals needed for research and by copyright restrictions that limit the freedom of sharing. Open access literature removes some of these limitations by making articles available digitally online, free of charge and free of most copyright and licensing restrictions (Suber, 2012). There are two types of open access literature: gold open access and green open access. Gold open access articles are delivered and distributed by journals. To be designated a gold article, authors simply submit a manuscript to an open-access journal (Suber, 2012). Green open access articles, however, are distributed by a repository. An author delivers a manuscript to an open access repository, an act also known as self-archiving, and the manuscript is required to be published in a journal that allows self-archiving (Suber, 2012). The flexibility of gold versus green open access makes it possible for authors to share manuscripts via a green repository, especially when a high prestige open access journal is not available for a given domain. There are a number of green repositories available including EarthArXiv (EarthArXiv, 2020), OSFPreprints (Center for Open Science, 2020), the Earth and Space Science Open Archive (ESSOAr, 2020), and ArXiv (Cornell University, 2020).

Journal publishers have adopted different publication approaches in order to make research articles openly accessible. For example, five of the American Geophysical Union's (AGU) journals are inherently open access while 16 other journals within AGU employ a hybrid model, where a study within a journal can be made open access if the article processing cost is borne by the author. American Mathematical Society (AMS), Society of Photographic Instrumentation Engineers (SPIE) and Institute of Electrical and Electronics Engineers (IEEE) Geoscience, and Remote Sensing Society (GRSS) publications have also adopted different variations of this publication model to address the need for increasing access to research articles. In addition, new policies at some journals allow authors to self-archive in green repositories. These policies allow authors to offer manuscripts and other information openly to others so long as the repository is not-for-profit and encourages scientific engagement. Some publishers even provide a green repository for self-archiving, the most notable being AGU's Earth and Space Science Open Archive (ESSOAr).

### 3.5. Efficient Research Processes and Knowledge Dissemination

The second open science focus area considers how to make both the research process and scientific collaborations more efficient.

#### 3.5.1. Cyberinfrastructure to Support Science

The shift in science toward data-intensive scientific discovery (Gray, 2009) has necessitated the need for new and better computation infrastructures to support science at scale. de La Beaujardière (2019) posed the question on how to enable "science at scale," such that researchers and other users can work with large, multisource data sets. Similarly, Robinson et al. (2020) stated the need for infrastructures that facilitate the construction of a robust, scalable, and adaptable data analysis pipeline. These infrastructures enable

researchers to cope with the volume of data, provide effective user/data interfaces and visualizations, and utilize more powerful algorithms to extract more information from these datasets. These enabling infrastructures move researchers away from the drudgery of data management and wrangling and back to intuitive and productive workflows, thus accelerating scientific progress. While there have been numerous efforts to build cyberinfrastructures (Droegemeier et al., 2005; Ludäscher et al., 2006; Nemani, 2012), renewed efforts are needed to couple infrastructures with research software and practices to enable more efficient Earth system science research (Bandaragoda et al., 2019).

### 3.5.2. Collaborations

Large collaborative teams composed of individuals with different types of expertise are needed to solve increasingly complex and interdisciplinary scientific problems. These contributors play a number of roles in the scientific process including performing experiments, curating data, conducting analyses, developing software, validating results, and conducting critical reviews (CASRAI, 2020). Innovative new platforms are needed that seamlessly enable scientific collaboration from a number of geographically distributed contributors who are performing a variety of tasks. These effective collaborative platforms need to offer a number of key capabilities (Bartling & Friesike, 2014; Roure et al., 2008) including the ability to easily manage research objects, incentivized sharing of those research objects, capacity to be open and extensible to future technology changes, ability to peer review scientific research objects (Himmelstein et al., 2019), and the means to support actionable research beyond simply serving as an object repository. A number of collaborative open science platforms, such as myExperiment (De Roure et al., 2009), JetStream (Jetstream, 2020) and GeneLab (Berrios et al., 2020; NASA, 2020b), support online collaborations in novel ways, along with science specific social networks, including ResearchGate and Mendeley (Nentwich & König, 2014), that allow researchers to connect and share journal articles.

### 3.6. Understanding Scientific Impact

The third open science focus area seeks to understand the broader implications of scientific research. As scientific research is shared more openly across a variety of platforms and to a number of different audiences, there is a need to comprehensively understand the impact of scientific contributions both to academia and to the broader public. Quantitative metrics, also known as impact measurements, offer one method for assessing scientific impact in the digital era. In the past 25 years, impact measurements have been limited to citation analysis of academic journal articles to assess scientific contributions (Fenner, 2014). While this constraint was initially necessary due to the physical printing of journal articles, the movement toward electronic publishing has eliminated this limitation. Citation analysis has offered some insights into scientific impact but the slow adoption of citations for some content like data and software (Fenner, 2014), means that many first-class research objects are not considered. Additionally, as scholarly workflows migrate to the web, other scientific impact measurements need to be considered such as article-level metrics, social network sharing metrics, and usage metrics (Fenner, 2014).

To address some of these impact measurement needs, the new discipline of altmetrics was established in 2010 to measure attention on the social web (Bar-Ilan et al., 2019). Altmetrics are collected for individual scientific outputs, such as a study or data set, and consider a variety of input variables including Mendeley document additions, social media shares, or blog post references (Bar-Ilan et al., 2019). The input variables vary from organization to organization as does the weight each input variable receives when calculating an altmetric (Crotty, 2017), making altmetric values different from service to service. Altmetrics are calculated by a number of service providers including Altmetric.com (Altmetric, 2020), PlumX (Plum Analytics, 2020), ImpactStory (Our Research, 2020), and Scienceopen.com (ScienceOpen, 2020). However, some altmetric calculation methods, such as those from Altmetric.com are proprietary, rendering these metrics fundamentally opposed to open science and open source principles.

## 4. Data Program's Role to Enable Open Science

The three open science focus areas can be encouraged and accelerated through thoughtful design (National Academies of Sciences, Engineering, & Medicine, 2018a, 2018b). Since modern science is primarily data-driven, data programs are in a unique position to design policies and systems that support and promote

open science. In this section, we describe specific actions that data programs can take to make the open science paradigm shift reality and we structure those actions around the elements described in Section 3. These actions represent an expansion in what some may consider the fundamental scope of a data program in that an emphasis is placed not only on data but on other key research aspects, such as software, documentation, and collaboration, that are essential to open science. For each action, we also provide specific examples of steps that NASA's Earth Science Data Systems (ESDS) is taking to support open science.

#### 4.1. Improving Accessibility to Science

##### 4.1.1. Encouraging Science as a Way of Thinking

Data programs can support data collection activities that involve public participation. For example, ESDS's Citizen Science for Earth Systems Program (CSESP) focuses on developing and implementing projects that encourage the general public's contributions to advance Earth system science. NASA's citizen science projects have the same scientifically rigorous standards as any other science mission directorate project (SMD Science Management Council, 2018), yet offer citizens an opportunity to participate in the process. Data systems also support the legitimacy of citizen science efforts by ensuring these data are subject to key stewardship activities, including standardized documentation, metadata, file formats, and quality assessments (Earth Science Data Systems, 2020a). The application of data stewardship processes ensures that these data are discoverable and usable to the broader community.

Data programs may also engage the public by systematically supporting scientific challenges, hackathons, and other open events. These events benefit open science by harnessing the public's creativity and innovation and by enabling the broader public to learn about and engage with science data. Several examples of these types of events include NASA's International Space Apps Challenge (NASA, 2020c), the Copernicus hackathons (Copernicus, 2020), and labeling events in platforms like Zooniverse (Zooniverse, 2020) and the Sentinel-Hub Classification App (Sentinel Hub, 2020). One archetypal example of a practical and collaborative approach to creating data challenges is the Climate Data Initiative (CDI) Innovation challenges. The Climate Data Initiative was an essential aspect of President Obama's Climate Action Plan (Office of the Press Secretary, 2014) that leveraged the federal government's extensive open data catalog to spur innovation and advance resilience to the impacts of climate change. The CDI Innovation challenges catalyzed new, data-driven solutions to help communities build resilience to climate change and included participants such as NASA, NOAA, USGS, USDA, Esri, Microsoft, and Research Data Alliance. More recently, NASA hosted the Space Apps COVID-19 Challenge to solve problems surrounding the COVID-19 pandemic. Over 15,000 people from 150 countries used Earth observation data from NASA and its partner space agencies to show how satellite information can aid in the understanding of the COVID-19 outbreak on both global and local scales (Landau, 2020).

Last, data programs should make a concerted effort to increase public awareness of the value of the data collected to advance science and humanity. Similar to science blogs, systematic communication channels should be employed to craft data stories for the public. These stories create a connection between scientific data and how they impact people's lives. NASA's Earth Observatory website (NASA, 2020a) and the Land Processes Distributed Active Archive Center (LP DAAC)'s Data In Action (Land Processes Distributed Active Archive Center, 2020a) articles are two examples of content which communicate the broader importance of data to society. Data programs should also consider participating in working groups, such as the GEOValue community (GEOValue, 2020), to better understand and promote Earth observation data's benefits to the broader community.

##### 4.1.2. Enabling Access to Science as a Body of Knowledge

Data system programs support open access to knowledge by developing and implementing open data and open source software policies. These policies, such as the 25-year-old ESDS open data policy, enable not only scientific research for a broad community of users, but also make international collaboration possible. Open data policies empower partnerships, such as the Group on Earth Observations (GEO)'s partnership of 100 national governments, to occur. In turn, these open data policies enable the creation of new and innovative products that benefit the wider Earth observation community. For example, the Harmonized Landsat Sentinel (HLS) data product (Claverie et al., 2018; Land Processes Distributed Active Archive Center, 2020b) is



being produced by NASA to support land applications and uses USGS's Landsat data and ESA's Sentinel 2 data. Data products like HLS are only possible when open data policies are in place.

Open data and software policies should be straightforward and easy to understand in order to benefit the user community and the scientists who are creating the data and software. Clear and unambiguous policies help users who need to understand any use or sharing constraints associated with data and software. Data programs can provide clarity to users by leveraging standard licenses for data and software instead of custom licenses, which may be interpreted in a number of ways. Straightforward, simple, and permissive policies also reduce the compliance burden on scientists who are not subject matter experts in open data and open source software and licensing.

Successful implementation of open policies by data programs may be achieved by providing clear and consistent guidelines for each step of the scientific data life cycle and by gaining the support of science program managers. Clear communication on policy expectations is essential and begins with solicitations and new research project requirements. Data programs can drive policies in the project formulation phase by including these policies in the data management plan (DMP) requirements for any project producing data of value to the broader community. Data programs should also consider requiring software management plans (SMPs) either as a stand-alone document or as a significant component within the DMP. Clear guidelines should also be provided regarding what data, software, and code are expected to be openly shared. Educational resources and well-formulated examples provided by the data program can also ensure that DMPs, SMPs, and the created research objects comply with open policies. While NASA's ESDS has well-documented guidelines for data management plans (Earth Science Data Systems, 2020d), SMPs are not currently required and are not a major component of DMPs. While ESDS does require that proposals include a plan for committing their software as open source (Earth Science Data Systems, 2020c), an opportunity exists to more formally support SMPs within the program.

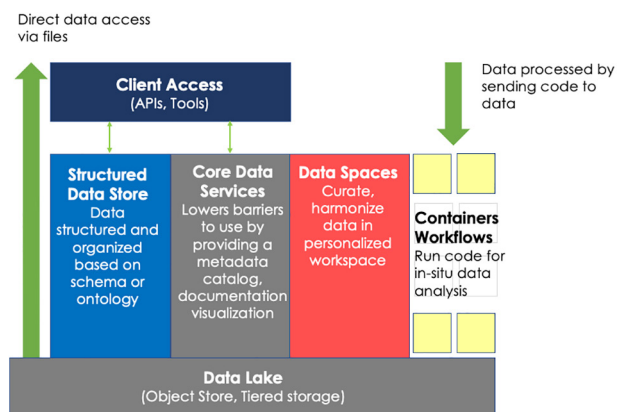
Once data and software are created, data programs can support open science by improving the discovery and use of crucial research objects through following strong data stewardship best practices and adopting the FAIR principles which state that data should be findable, accessible, interoperable, and reusable (Mons et al., 2017). These best practices include creating comprehensive metadata and assigning persistent identifiers and citations to all data and software in order to support understanding and reuse. While progress has been made in embracing data citation standards, there is a need for data programs to adopt software citation principles (Stodden et al., 2018). Data programs can support software citation by leveraging existing software citation standards, by generating persistent identifiers for software and by encouraging and supporting authors to cite software in journal articles. While NASA's ESDS open source policy requires that all software source code developed through ESDS-funded research solicitations be designated, developed, and distributed to the public as open source, the program still has a need to develop software metadata, citations, and digital object identifiers (DOIs) to enable broader discovery.

Finally, data programs should support access to research publications related to data and software. Digital library portals, such as the Astrophysics Data System (ADS) (Smithsonian Astrophysical Observatory, 2020), can serve not only as a primary entry point to relevant scientific publications, but also as a linkage between journal articles and related data and code. Data programs should consider developing and maintaining infrastructures like the ADS which allows for the discovery of journal articles and also serves as a green repository for authors to upload preprints or journal articles that are free and open. The success of a data program-sponsored green repository can be magnified by incentivizing authors to participate in the open journal publication process. Authors should be encouraged to either publish in an open-access journal or to a journal that allows publishing to a green repository.

## 4.2. Leveraging Technology to Enable Efficient Research Processes and Knowledge Dissemination

### 4.2.1. From Data Systems to Enabling Collaborative Infrastructures

In current data system architectures, most data archives are separate from computing resources. Any analysis requires data movement to either a user's machine or some computing resource. Cloud platforms are quickly becoming a viable building block for designing collaborative infrastructures that move data out



**Figure 2.** A conceptual data system architecture in the cloud (Bugbee et al., 2020).

of organizational silos and instead collocate them with compute (de La Beaujardière, 2019). As a new computing paradigm, cloud computing delivers scalable, on-demand, pay-as-you-go access to a pool of computing resources. These cloud technologies make it easier, more efficient, and economical to configure data analysis platforms for large scale computing. New data systems need to be reimagined to be cloud-native and integrated with data analysis platforms.

A conceptual cloud-native data system framework that supports such new cyberinfrastructures is depicted in Figure 2 below (Bugbee et al., 2020). The data lake serves as a central repository of different data. The core set of data services in this framework cover scientific data stewardship functions such as metadata generation, data ingest, reformatting, documentation, and data publication, and services required by a user to discover, visualize, and access the data. This framework allows users to have multiple modes of accessing data. While users have direct access to the data files, structured data stores can be instantiated when needed, as envisioned in an analytics-optimized data store (Ramachandran et al., 2019). Technol-

ogy for implementing structured data stores can range from new databases such as Rasdaman (Baumann et al., 2013) to data cubes (Giuliani et al., 2019) to software frameworks like Pangeo (Guillaume et al., 2019), which utilize cloud-optimized formats like Zarr and Dask for parallelizing processing tasks. These structured data stores minimize the data wrangling burden on the end-users and allow fast reads on the data. Virtualization enables users to package their code as containers and run it on the data residing in the data lake. The ability to perform analysis either using a structured data store or by running containers on the data lake minimizes the need to download data. The data spaces component provides a private personalized workspace where users can curate and integrate different data. Collaboration within data spaces allows the sharing of both code and data.

Data programs should invest in cloud-native data systems that follow this conceptual framework. Data programs should also invest in cyberinfrastructures that build on such cloud-native data systems to enable collaboration and allow users to facilitate new science. The Joint ESA-NASA multimission algorithm and analysis platform (MAAP) is one such example. The MAAP is a collaborative, cloud-based, open science platform dedicated to addressing the biomass community's unique research needs (Bugbee et al., 2020). The MAAP provides a new approach to accessing, sharing, analyzing, and processing data. The MAAP's users have seamless access to curated airborne, spaceborne, and field from both ESA and NASA.

Data programs should also invest in collaborative, open-source data stewardship tools that support open science principles. Tools such as the Algorithm Publication Tool (APT) make it easier for scientists to collaboratively write Algorithm Theoretical Basis Documents (ATBDs) (Bugbee et al., 2020). ATBDs describe the physical theory, mathematical procedures, and assumptions made for the algorithms that create higher-level data products. While these documents reinforce the data program's commitment to supporting reproducibility and transparency, providing a tool like the APT makes it easier for scientists to collaborate on these documents and, in the end, makes the data program more open. The availability, adoption, and use of such tools, and new enabling cyberinfrastructures are needed to make research processes efficient and to accelerate knowledge dissemination.

### 4.3. Understanding Scientific Impact Through Data System Measurements

Data programs should adopt or develop impact measurements for all first-class research objects within the data system, including data, software, documentation, and services. Impact measurements help data programs understand the value and use of data, software, services, and information not just within the scientific community but more broadly. A greater understanding of these objects' impact helps informing future decisions about new projects, new technological infrastructure needs, and new data stewardship requirements. Insights from data program impact measurements also help scientists quantitatively demonstrate the value of data and research. A secondary benefit to these measurements is a better understanding of the effectiveness of internal data system processes including data stewardship best practices. Last, data

program impact measurements may enhance the open science process by serving as a metric or weight for discovering new research objects. Combining impact metrics with traditional search techniques enables discovery by highlighting research objects of interest.

Data program impact measurements should include data usage metrics and citation metrics along with emerging metrics techniques such as altmetrics. Most data systems, including NASA's ESDS, collect data usage metrics (Earth Science Data Systems, 2020b) such as the number of views of a data set, the number of downloads or how many times a data set has been accessed. While most data systems collect these metrics, quantifying data and software impact metrics in the research community via citations has been difficult due to inconsistent or nonexistent citation practices. Emerging methods, such as altmetrics and machine learning techniques, make it possible to overcome the lack of citations in articles in order to better understand the impacts of data, software, and other information. For example, leveraging machine learning techniques across open access repositories and other unstructured data assets provides an alternative way to assess data and software impacts. These techniques may make it possible to better understand relationships between journal articles and data, along with the provenance of data.

Similarly, a new type of altmetrics for data systems that combine data, software, code, and documentation usage may provide another means of assessing impact. While traditional altmetrics have focused on journal articles as the primary research object of interest, the same approach may not be effective for data and software. Data system altmetrics should include already established input variables such as social media shares but should also explore other variables of interest, including measures of data stewardship such as metadata quality scores or data stewardship maturity matrix assessment scores (Peng et al., 2019). Similarly, software usage should be assessed by how widely a piece of code is supported by the community through measures such as forks on Github and incorporation into other tools and workflows. Last, documentation, such as ATBDs, data guides, data recipes, and blogs should be treated as first-class research objects within the data system. Assigning DOIs and citations to these key documents will make it possible to track these objects' impact and inclusion as a variable for data system altmetrics. Establishing linkages between data, software, and documentation, similar to those seen in Scholix Link Information Packages (Lowenberg et al., 2019), are a comprehensive way of measuring impact across the data system.

Most importantly, data system altmetrics should conform to open science and open-source principles. The methodology used to generate altmetrics should be transparent and reproducible, while the code used to create altmetrics should be open-sourced. Ensuring data system altmetrics are open also makes community buy-in and agreement possible. More broadly, data programs should ensure that all impact measurements are open to the scientists and researchers who create these valuable research objects. Scientists may use these measurements to support career opportunities and to demonstrate the value of data to peers. This openness, in turn, may motivate scientists to cite data and software in articles more consistently.

## 5. Challenges for Data Programs

There are a number of challenges facing data programs that must be considered when supporting the open science paradigm. These challenges require creative solutions along with community engagement and discussion and are described below by specific focus areas.

### 5.1. Accessibility to Science

Greater engagement with the public is an important aspect of open science, yet there are still hurdles to be overcome to support this engagement. For example, while there is a growing interest in supporting citizen scientist activities, there is some reluctance to accept citizen science data as scientifically legitimate from both within the science community and even by some data centers. This reluctance is due to the perception that citizen science data is of lower quality and may not be of use to either science or decision making (Shanley et al., 2019). Overcoming these biases will be critical to ensure the long term adoption and viability of citizen science data. Data programs, in collaboration with scientists, may need to develop new policies and best practices to overcome these challenges. These policies and best practices may focus on making citizen science data more trustworthy and usable by requiring data management plans for these data, by designing

these data to be easily interoperable with existing Earth observation data, by providing standard metadata and by ensuring the data are geolocated for easy use in decision making tools (Newman et al., 2017).

Similarly, many scientists and data centers are averse to committing limited resources in order to effectively communicate data value to the wider public. To encourage more scientists to become these types of boundary spanners, the scientific community will need to evolve to value these contributions to science. Data programs and scientific programs can foster acceptance by providing funding for these types of roles and creating solicitations that emphasize science communication to the broader public.

Providing open access to data makes it possible for data to be reused for scientific benefit but also poses some issues for data programs. First, making data openly available introduces conditions for data misuse. Individual users may unwittingly use data in inappropriate ways or, more insidiously, use data to misrepresent facts. For example, Earth observation data was recently misused or misrepresented in maps and images generated during the Australian fire crisis (BBC, 2020). Similarly, some third-party organizations bulk download open data in order to redistribute the data. The level of data stewardship and understanding provided by these third party distributors is typically lower than those of the original data providers and may lead to data misunderstanding and misuse by users who do not always understand the need to use quality flags or the caveats surrounding data use. Last, requiring researchers to make data openly available is sometimes seen as burdensome and a distraction from research itself. Many scientists acknowledge that they are holding “dark data, or data that has never been published or otherwise made available to the rest of the scientific community” (Heidorn, 2008). With little professional rewards for sharing these data and a lack of time and resources, scientists are often unwilling to put any effort into making these data openly available. In order to increase open data sharing, data programs need to develop solutions that minimize the burden on researchers to share data and also foster incentives and credits for sharing.

Open source software presents similar problems for both researchers and data programs. Similar to open data requirements, a move toward open source software introduces extra requirements for scientists that is often viewed as distracting (Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division On Engineering and Physical Sciences, 2018). Scientists are typically not experts in the open source process and are confused by the lack of consensus as to where to share code and which license to select for code. Assigning an appropriate license to software is confusing and intimidating to a novice code developer and is often overwhelming for scientists who want to spend time pursuing science and not learning copyright laws. This confusion is compounded by the fact that the use of standard open-source licenses may be limited by the constraints imposed by government mandates. In addition to the licensing, there is currently minimal guidance on which code repositories should be used for sharing software and code. Some source code is individually shared using code repositories such as Github, while some organizations leverage an organizational Github/Gitlab repository instance. Official repositories, such as the IEEE Remote Sensing Code Library (IEEE, 2020), may alleviate this issue but options are currently limited. Data programs may need to consider providing clear and easy-to-understand guidelines on the open source software process to promote success.

While open access publication primarily poses challenges for researchers, data programs have a need to provide free and open access to journal articles related to data and software. Data programs may address this need by providing a green repository for open access to relevant journals and by working with scientists to ensure data articles are published in journals which support self-publication.

## 5.2. Enabling Efficient Research Processes and Knowledge Dissemination

Developing enabling infrastructures that support both collaboration and analysis at scale requires a fundamental transformation in existing data systems. Data programs must invest in and evolve data systems toward infrastructures that facilitate analysis. This data system evolution will require data centers to move beyond simply functioning as a data archive to instead serving as knowledge centers for the scientific community. This evolution will not be easy due to the sunk costs of existing processing, ingest, archive and distribution systems along with the need to retrain existing staff for new roles and responsibilities. Yet, data programs may view adopting cloud technologies as an opportunity to design and build new data analysis



platforms that function as a cohesive and interdisciplinary ecosystem rather than individual, stand-alone data silos.

For data programs, there is a risk of building new infrastructure solutions that do not meet the scientific community's needs or have low adoption rates that do not justify the cost of implementation. The “Build it, and they will come” type of tool development mentality always runs a risk of low adoption if the potential users are not engaged at the beginning of a project (Cutcher-Gershenfeld et al., 2016). This risk is especially true for collaboration tools which need to be designed to help scientists do what they are already doing, not what the tool developer feels scientists should be doing. Data programs may mitigate this risk by including scientific stakeholders early and throughout the infrastructure development process.

Similarly, while cloud computing has the potential to be a transformational technology for improving research processes and collaboration, cloud computing is not a panacea. For data programs, developing cloud-native tools and learning to utilize the cloud platform effectively in a cost-optimal manner is a non-trivial task. Advances in cloud platforms happen very rapidly, with new services added every week. These rapid advances place an extra burden on code maintenance and refactoring. Clearly, the cost of managing and sustaining these new types of cloud-based infrastructures will require the development of new sustainable business models and best practices. Cloud-based collaborative infrastructures may also require advances in data stewardship practices. As cloud computing enables different groups to build data sharing and analysis platforms quickly, there is a risk that these collaboration platforms will lead to data lakes turning into swamps filled with data of dubious quality (Sicular, 2016). These platforms will require systematic, semi-automated data governance, and management plans in order to prevent the degradation of the data lake into a data swamp.

Cloud computing poses some challenges for researchers as well. Cloud computing still does not solve equity issues due to the disparity in bandwidth and the lack of needed technological infrastructures. The lack of access to the internet, high bandwidths, or necessary computing resources means that not everyone has equal access to open data and information. Similarly, while cloud computing reduces the initial cost needed to run analyses at scale, funding is still required to run cloud computing. The need for funding is particularly pressing as science adopts artificial intelligence (AI) and machine learning (ML) techniques to build models which require computing at a scale that only a few organizations can afford. Finally, cloud capacity building is needed to assist researchers and end users to effectively utilize cloud platform capabilities.

### 5.3. Impact Measurements

There are still several obstacles to be addressed in order for data programs to calculate data system impact metrics effectively. The biggest hurdle centers around the need to accurately understand and calculate data usage beyond the immediate data system itself. Data programs collect download metrics, and data, and software citations in journal articles may serve as one measure of usage when data programs provide unique identifiers for citation. However, the research community has not widely adopted data and software citation best practices and citations are often lost during the journal submission process (Lowenberg et al., 2019). Similarly, data and code may be used for applications and in decision-making processes that are never published in a peer-reviewed journal. Last, as more third-party vendors harvest open data and make it available for use on alternative platforms, data programs will need to decide how to account for the varied usage of valuable data on these platforms.

While data system altmetrics may be another method to help data programs better understand usage, considerable effort will be needed in order to define and test data set and software altmetrics. Some initial work around data altmetrics is being accomplished by groups such as ImpactStory (Fenner, 2014), but there is still considerable work to be done to define what altmetrics mean to data programs, what input metrics should be considered for these altmetrics and whether data program altmetrics can be implemented in a viable manner.

More broadly, the acceptance and adoption of data system altmetrics is not assured. While data programs can take steps to ensure wider adoption, such as leveraging these metrics when reviewing new proposals, there is no guarantee that altmetrics will be deemed of value by the community. Even if altmetrics are adopted, there are still limitations associated with collecting these metrics. Data, software, and research

may receive attention for the wrong reasons, thus misrepresenting the metrics and, in the end, the value of an object. In addition, there is a risk that these metrics are subject to self-promotion, gaming, and the constantly changing nature of the web (Fenner, 2014). Variables such as the career stage of an object creator, the specific discipline in which an object is created, the number and mix of coauthors and whether an object is cross-disciplinary in nature makes metrics comparisons problematic at best (Kurtz & Henneken, 2017). Effectively calculating impact metrics will require careful thought by data programs. However, if fairly and mindfully adopted, these metrics have the potential to provide greater insights into the impacts of data more broadly.

## 6. Conclusions

The open science movement continues to gain momentum as both scientists and organizations adopt many of the principles and ideas described in this study. Open science beckons with the promise of more collaborative and efficient research, a more educated and engaged public, and scientific results that are reproducible and easier to understand. No doubt, there will be issues as the open science paradigm shift continues to expand. Data and compute equity will always be a challenge until reliable internet access is available to the majority. In addition, there are risks associated with both guaranteeing data authenticity and increasing data misuse. However, these issues should not keep us, as a community, from moving open science forward.

In order to move the open science paradigm forward, active participation is needed from both supporting organizations such as data programs and journal publishers but also individuals from the research community. As highlighted in this study, data programs need to acknowledge the major role they play in supporting open science through the development of open policies, investment in innovative, and collaborative infrastructures, and the promotion of cultural change.

Data programs should support data collection activities that involve public participation and help legitimize such citizen science efforts by enforcing standard scientific stewardship activities. The application of scientific data stewardship processes will ensure both usability and trust of the collected data. Engaging the public via scientific challenges and hackathons will enable data programs to harness the public's creative and innovative engagement with science data and increase public awareness of the scientific data's value. Data system programs should support open access to knowledge by developing and implementing open data and open source software policies and improving the discovery and use of crucial research objects. Knowledge is open only when data, software, documentation, and publications are linked and discoverable. Data programs should invest in developing the next generation of enabling cyberinfrastructures that remove the drudgery of data management and wrangling and allow the use of open science tools to transform science. Finally, data programs should adopt or develop impact measurements for all first-class research objects within the data system, including data, software, documentation, services, and users. Such measures will help guide future investments both in science as well as the data and information systems.

Individual researchers, on the other hand, can be advocates for open science by adopting a number of best practices. First, researchers should make their data available in an open repository and in a nonproprietary, standardized format. Whenever possible, they should create a DOI for their data and be sure to provide clear licensing information along with use constraints about the data. Second, software and code should be made open source through a community adopted license that is as permissive as possible to encourage reuse. In addition, any libraries used should also be open sourced and follow software development best practices, such as rigorous version controls. Third, researchers should contribute back by supporting the community development of open source software, libraries, code, and tools that are used by the wider community. When appropriate, code should be open to community development and feedback as well. Fourth, peer-reviewed articles should be published to gold journals when possible. If a reputable gold journal is not available for a particular domain, researchers should make sure their publishers allow self-publishing in a green repository. Fifth, the public should be actively engaged by starting or supporting science blogs, citizen science projects, hack-a-thons, or sharing results on social media. Finally, researchers should cite data, software, and documentation whenever possible but especially in journal articles.

Together individual researchers, data programs and organizations will usher in a new and more open age of scientific research that will not only benefit science itself but will bring further advances to society as a whole.

## Data Availability Statement

Data sharing is not applicable to this article as no new data were created or analyzed in this study.

## Acknowledgments

Authors would like to thank Derek Koehl for his assistance in editing and improving the manuscript. They also thank all our colleagues within the Earth Science Data Systems Program NASA who provided insight and expertise that greatly assisted this research.

## References

- Altmetric. (2020). *Discover the attention surrounding your research*. Altmetric. Retrieved from <https://www.altmetric.com/>
- Arias, J. J., Pham-Kanter, G., & Campbell, E. G. (2015). The growth and gaps of genetic data sharing policies in the United States. *Journal of Law and the Biosciences*, 2, 56–68. <https://doi.org/10.1093/jlb/lso032>
- Bandaragoda, C., Castronova, A., Istanbuluoglu, E., Strauch, R., Nudurupati, S. S., Phuong, J., et al. (2019). Enabling collaborative numerical modeling in earth sciences using knowledge infrastructure. *Environmental Modelling and Software*, 120, 104424. <https://doi.org/10.1016/j.envsoft.2019.03.020>
- Bar-Ilan, J., Halevi, G., & Milojević, S. (2019). Differences between altmetric data sources: A case study. *Journal of Altmetrics*, 2, 1. <https://doi.org/10.29024/joa.4>
- Bartling, S., & Friesike, S. (2014). Towards another scientific revolution. In S. Bartling, & S. Friesike (Eds.), *Opening science* (pp. 3–15). Springer International Publishing. [https://doi.org/10.1007/978-3-319-00026-8\\_1](https://doi.org/10.1007/978-3-319-00026-8_1)
- Baumann, P., Dumitru, A. M., & Merticariu, V. (2013). The array database that is not a database: File based array query answering in radsaman. In M. A. Nascimento, T. Sellis, R. Cheng, J. Sander, Y. Zheng, H.-P. Kriegel, M. Renz, & C. Sengstock (Eds.), *Advances in spatial and temporal databases* (pp. 478–483). Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-642-40235-7\\_32](https://doi.org/10.1007/978-3-642-40235-7_32)
- BBC. (2020). *Australia fires: Misleading maps and pictures go viral*. BBC News. Retrieved from <https://www.bbc.com/news/blogs-trending-51020564>
- Berrios, D. C., Galazka, J., Grigorev, K., Gebre, S., & Costes, S. V. (2020). NASA GeneLab: Interfaces for the exploration of space omics data. *Nucleic Acids Research*, 49, D1515–D1522. <https://doi.org/10.1093/nar/gkaa887>
- Borowitz, M. (2017). *Open space: The global effort for open access to environmental satellite data*. The MIT Press. <https://doi.org/10.7551/mitpress/10659.001.0001>
- Bugbee, K., Lynnes, C., Ramachandran, R., Maskey, M., Barciauskas, A., Kaulfus, A., et al. (2020). Advancing open science through innovative data system solution: The joint ESA-NASA multi-mission algorithm and analysis platform (MAAP)'S data ecosystem. In *IGARSS 2020 - 2020 IEEE international geoscience and Remote sensing symposium*. <https://doi.org/10.1109/igarss39084.2020.9323731>
- CASRAI. (2020). *CRedit: Contributor roles taxonomy*. CASRAI. Retrieved from <https://casrai.org/credit/>
- Center for Open Science. (2020). *OSF. OSFHome*. Retrieved from <https://osf.io/>
- Chard, K., Gaffney, N., Hategan, M., Kowalik, K., Ludäscher, B., McPhillips, T., et al. (2020). Toward enabling reproducibility for data-intensive research using the whole tale platform. *Advances in Parallel Computing*, 36, 766–778. <https://doi.org/10.3233/APC200107>
- Chesbrough, H. (2015). *From open science to open innovation* (pp. 51–66). Institute for Innovation and Knowledge Management, ESADE.
- Claverie, M., Ju, J., Masek, J. G., Dungan, J. L., Vermote, E. F., Roger, J.-C., et al. (2018). The harmonized landsat and sentinel-2 surface reflectance data set. *Remote Sensing of Environment*, 219, 145–161. <https://doi.org/10.1016/j.rse.2018.09.002>
- Collins, F. S., Morgan, M., & Patrinos, A. (2003). The human genome project: lessons from large-scale biology. *Science*, 300(5617), 286–290. <https://doi.org/10.1126/science.1084564>
- Committee on Best Practices for a Future Open Code Policy for NASA Space Science, Space Studies Board, & Division on Engineering and Physical Sciences. (2018). *Open source software policy options for NASA Earth and space sciences* (p. 25217). National Academies of Sciences, Engineering, and Medicine. <https://doi.org/10.17226/25217>
- Cook-Deegan, R., & McGuire, A. L. (2017). Moving beyond Bermuda: Sharing data to build a medical information commons. *Genome Research*, 27, 897–901. <https://doi.org/10.1101/gr.216911.116>
- Copernicus. (2020). *Copernicus hackathons*. EU. Retrieved from <https://hackathons.copernicus.eu/>
- Cornell University. (2020). *ArXiv.org e-Print archive*. Retrieved from [www.arxiv.org](http://www.arxiv.org)
- Crotty, D. (2017). Altmetrics. *European Heart Journal*, 38, 2647–2648. <https://doi.org/10.1093/eurheartj/ehx447>
- Cutcher-Gershenfeld, J., Baker, K. S., Berente, N., Carter, D. R., DeChurch, L. A., Flint, C. C., et al. (2016). Build it, but will they come? A geoscience cyberinfrastructure baseline analysis. *Codata*, 15, 8. <https://doi.org/10.5334/dsj-2016-008>
- Cyberinfrastructure Council. (2007). *Cyberinfrastructure vision for 21st century discovery*. National Science Foundation.
- de La Beaujardière, J. (2019). A geodata fabric for the 21st century. *Eos*, 100. <https://doi.org/10.1029/2019EO136386>
- De Roure, D., Goble, C., Alekseyevs, S., Bechhofer, S., Bhagat, J., Cruickshank, D., et al. (2010). Towards open science: The myExperiment approach. *Concurrency and Computation: Practice and Experience*, 22, 2335–2353. <https://doi.org/10.1002/cpe.1601>
- De Roure, D., Goble, C., & Stevens, R. (2009). The design and realisation of the virtual research environment for social sharing of workflows. *Future Generation Computer Systems*, 25, 561–567. <https://doi.org/10.1016/j.future.2008.06.010>
- Donaldson, D., & Storeygard, A. (2016). The view from above: Applications of satellite data in economics. *Journal of Economic Perspectives*, 30, 171–198. <https://doi.org/10.1257/jep.30.4.171>
- Droegemeier, K. K., Gannon, D., Reed, D., Plale, B., Alameda, J., Baltzer, T., et al. (2005). Service-oriented environments in research and education for dynamically interacting with mesoscale weather. *IEEE Computing in Science & Engineering*, 7, 24–32. <https://doi.org/10.1109/mcse.2005.124>
- EarthArXiv. (2020). *EarthArXiv*. EarthArXiv. Retrieved from <https://eartharxiv.org/>
- Earth Science Data Systems. (2020a). *NASA ESDS citizen science data working group white paper (1.0)*. NASA. Retrieved from <https://cdn.earthdata.nasa.gov/conduit/upload/14273/CSDWG-White-Paper.pdf>
- Earth Science Data Systems. (2020b). *EOSDIS annual metrics Reports NASA Earth science data systems program*. Retrieved from <https://earthdata.nasa.gov/eosdis/system-performance/eosdis-annual-metrics-reports>
- Earth Science Data Systems. (2020c). *ESDS open source software policy. NASA Earth science data systems program*. Retrieved from <https://earthdata.nasa.gov/collaborate/open-data-services-and-software/esds-open-source-policy>

- Earth Science Data Systems. (2020d). *Data management plan Guidance NASA Earth science data systems program*. Retrieved from <https://earthdata.nasa.gov/esds>
- ESSOAr. (2020). *Earth and space science open archive*. ESSOAr. Retrieved from <https://www.essoar.org/>
- Fecher, B., & Friesike, S. (2014). Open Science: One term, five school of thoughts. *Opening science*. Springer.
- Fenner, M. (2014). Altmetrics and other novel measures for scientific impact. In S. Bartling, & S. Friesike (Eds.), *Opening science* (pp. 179–189). Springer. [https://doi.org/10.1007/978-3-319-00026-8\\_12](https://doi.org/10.1007/978-3-319-00026-8_12)
- Friesike, S., Widenmayer, B., Gassmann, O., & Schildhauer, T. (2015). Opening science: Towards an agenda of open science in academia and industry. *The Journal of Technology Transfer*, 40, 581–601. <https://doi.org/10.1007/s10961-014-9375-6>
- GEOValue. (2020). *Geo Value*. Retrieved from [www.geovalue.org](http://www.geovalue.org)
- Gil, Y., David, C. H., Demir, I., Essawy, B. T., Fulweiler, R. W., Goodall, J. L., et al. (2016). Toward the geoscience paper of the future: Best practices for documenting and sharing research from data to software to provenance. *Earth and Space Science*, 3, 388–415. <https://doi.org/10.1002/2015EA000136>
- Giuliani, G., Camara, G., Killough, B., & Minchin, S. (2019). Earth observation open science: Enhancing reproducible science using data cubes. *Data*, 4, 147. <https://doi.org/10.3390/data4040147>
- Gray, J. (2009). A transformed scientific method. In T. Hey, S. Tansley, & K. Tolle (Eds.), *The fourth paradigm: Data-intensive scientific discovery*. Microsoft Research.
- Group on Earth Observations. (2020). *GEO data sharing principles implementation*. Group on Earth Observations. Retrieved from [https://www.earthobservations.org/geoss\\_dsp.shtml](https://www.earthobservations.org/geoss_dsp.shtml)
- Guillaume, E.-B., Abernathey, R., Hamman, J., Ponte, A., & Rath, W. (2019). The pangeo big data ecosystem and its use at CNES. In *Big data from space (BiDS'19) turning data into insights*. Retrieved from <https://archimer.ifremer.fr/doc/00503/61441/>
- HeidornBryan, P. (2008). Shedding light on the dark data in the long tail of science. *Library Trends*, 57, 280–299. <https://doi.org/10.1353/lib.0.0036>
- Himmelstein, D. S., Rubinetti, V., Slochower, D. R., Hu, D., Malladi, V. S., Greene, C. S., & Gitter, A. (2019). Open collaborative writing with Manubot. *PLoS Computational Biology*, 15, e1007128. <https://doi.org/10.1371/journal.pcbi.1007128>
- Hood, L., & Rowen, L. (2013). The human genome project: Big science transforms biology and medicine. *Genome Medicine*, 5, 79. <https://doi.org/10.1186/gm483>
- IEEE. (2020). *Remote sensing code library*. IEEE. Retrieved from <http://www.grss-ieee.org/publication-category/rscl/>
- International Human Genome Sequencing Consortium. (2001). Initial sequencing and analysis of the human genome. *Nature*, 409, 860–921. <https://doi.org/10.1038/35057062>
- Jetstream. (2020). *Jetstream: A national Science and engineering cloud*. Jetstream. Retrieved from <https://jetstream-cloud.org>
- Kennedy, R. E., Andréfouët, S., Cohen, W. B., Gómez, C., Griffiths, P., Hais, M., et al. (2014). Bringing an ecological view of change to Landsat-based remote sensing. *Frontiers in Ecology and the Environment*, 12, 339–346. <https://doi.org/10.1890/130066>
- Keyes, D., & Taylor, V. (2011). *National science foundation advisory committee for CyberInfrastructure task force on software for science and engineering [final report]*. National Science Foundation. Retrieved from [https://www.nsf.gov/cise/oac/taskforces/TaskForceReport\\_Software.pdf](https://www.nsf.gov/cise/oac/taskforces/TaskForceReport_Software.pdf)
- Koop, D., Santos, E., Mates, P., Vo, H. T., Bonnet, P., Bauer, B., et al. (2011). A provenance-based infrastructure to support the life cycle of executable papers. *Procedia Computer Science*, 4, 648–657. <https://doi.org/10.1016/j.procs.2011.04.068>
- Kurtz, M. J., & Henneken, E. A. (2017). *Measuring metrics: A forty year longitudinal cross-validation of citations, downloads, and peer review in Astrophysics*. <https://doi.org/10.1002/asi.23689>
- Laakso, M., Welling, P., Bukvova, H., Nyman, L., Björk, B.-C., & Hedlund, T. (2011). The development of open access journal publishing from 1993 to 2009. *PloS One*, 6, e20961. <https://doi.org/10.1371/journal.pone.0020961>
- Landau, E. (2020). *NASA space Apps COVID-19 challenge winners share stories of innovation*. NASA. Retrieved from <https://www.nasa.gov/feature/nasa-space-apps-covid-19-challenge-winners-share-stories-of-innovation>
- Land Processes Distributed Active Archive Center. (2020a). *Data in action*. LPDAAC. Retrieved from <https://lpdaac.usgs.gov/resources/data-action/>
- Land Processes Distributed Active Archive Center. (2020b). *Harmonized landsat-sentinel 2 (HLS) overview*. LPDAAC. Retrieved from <https://lpdaac.usgs.gov/data/get-started-data/collection-overview/missions/harmonized-landsat-sentinel-2-hls-overview/>
- Lowenberg, D., Chodacki, J., Fenner, M., Kemp, J., & Jones, M. B. (2019). *Open data metrics: Lighting the fire (version 1) [computer software]*. Zenodo. <https://doi.org/10.5281/ZENODO.3525349>
- Ludäscher, B., Altintas, I., Berkley, C., Higgins, D., Jaeger, E., Jones, M., et al. (2006). Scientific workflow management and the Kepler system. *Concurrency and Computation: Practice and Experience*, 18, 1039–1065. <https://doi.org/10.1002/cpe.994>
- McKiernan, E. C., Bourne, P. E., Brown, C. T., Buck, S., Kenall, A., Lin, J., et al. (2016). How open science helps researchers succeed. *ELife*, 5, e16800. <https://doi.org/10.7554/eLife.16800>
- Mons, B., Neylon, C., Velterop, J., Dumontier, M., da Silva Santos, L. O. B., & Wilkinson, M. D. (2017). *Cloudy, increasingly FAIR: revisiting the FAIR Data guiding principles for the European open science cloud* (37, pp. 49–56). Information Services & Use. <https://doi.org/10.3233/ISU-170824>
- Murphy, M. C., Mejia, A. F., Mejia, J., Yan, X., Cheryan, S., Dasgupta, N., et al. (2020). Open science, communal culture, and women's participation in the movement to improve science. *Proceedings of the National Academy of Sciences*, 117, 24154–24164. <https://doi.org/10.1073/pnas.1921320117>
- NASA. (2020a). *Earth observatory*. NASA. Retrieved from <https://earthobservatory.nasa.gov/>
- NASA. (2020b). *NASA GeneLab: Open science for life in space*. GeneLab. Retrieved from <https://genelab.nasa.gov>
- NASA. (2020c). *NASA international space apps challenge*. Space Apps Challenge. Retrieved from <https://www.spaceappschallenge.org/>
- National Academies of Sciences, Engineering, and Medicine. (2018a). *Open science by design: Realizing a vision for 21st century research*. The National Academies Press.
- National Academies of Sciences, Engineering, & Medicine. (2018b). *Open source software policy options for NASA earth and space sciences*. In *National Academies of sciences, engineering, and medicine 2018*. The National Academies Press. <https://doi.org/10.17226/25217>
- Nemani, R. (2012). NASA earth exchange: Next generation earth science collaborative. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, XXXVIII. <https://doi.org/10.5194/isprsarchives-XXXVIII-8-W20-17-2011>
- Nentwich, M., & König, R. (2014). Academia goes facebook? the potential of social network sites in the scholarly realm. In S. Bartling, & S. Friesike (Eds.), *Opening science* (pp. 107–124). Springer International Publishing. [https://doi.org/10.1007/978-3-319-00026-8\\_7](https://doi.org/10.1007/978-3-319-00026-8_7)
- Newman, G., Chandler, M., Clyde, M., McGreavy, B., Haklay, M., Ballard, H., et al. (2017). Leveraging the power of place in citizen science for effective conservation decision making. *Biological Conservation*, 208, 55–64. <https://doi.org/10.1016/j.biocon.2016.07.019>



- Newman, G., Wiggins, A., Crall, A., Graham, E., Newman, S., & Crowston, K. (2012). The future of citizen science: Emerging technologies and shifting paradigms. *Frontiers in Ecology and the Environment*, 10, 298–304. <https://doi.org/10.1890/110294>
- Office of the Press Secretary. (2014). *FACT SHEET: The president's climate data initiative: Empowering America's communities to prepare for the effects of climate change*. The White House. Retrieved from <https://obamawhitehouse.archives.gov/the-press-office/2014/03/19/fact-sheet-president-s-climate-data-initiative-empowering-america-s-comm>
- Open Knowledge Foundation. (2020). *What is open? Open knowledge foundation*. Retrieved from <https://okfn.org/opendata/>
- Our Research. (2020). *ImpactStory*. ImpactStory. Retrieved from [www.impactstory.org](http://www.impactstory.org)
- Peng, G., Milan, A., Ritchey, N. A., Partee, II, R. P., Zinn, S., McQuinn, E., et al. (2019). Practical application of a data stewardship maturity matrix for the NOAA onestop project. *Codata*, 18, 41. <https://doi.org/10.5334/dsj-2019-041>
- Plum Analytics. (2020). *About PlumX metrics*. Plum Analytics. Retrieved from <https://plumanalytics.com/learn/about-metrics/>
- Ramachandran, R., Bugbee, K., Maskey, M., & Lynnes, C. (2019). From ARDS to AODS: Future of analytics for Earth observations. In IGARSS 2019 - 2019 IEEE international geoscience and Remote sensing symposium. IGARSS.
- Robinson, N. H., Hamman, J., & Abernathy, R. (2020). *Seven principles for effective scientific big-data systems*. ArXiv. <http://arxiv.org/abs/1908.03356>
- Roure, D. D., Goble, C., Bhagat, J., Cruickshank, D., Goderis, A., Michaelides, D., & Newman, D. (2008). myExperiment: Defining the social virtual research environment. In *2008 IEEE fourth international conference on EScience* (pp. 182–189). IEEE. <https://doi.org/10.1109/eScience.2008.86>
- Roy, D. P., Wulder, M. A., Loveland, T. R., Woodcock, C. E., Allan, R. G., Anderson, M., et al. (2014). Landsat-8: Science and product vision for terrestrial global change research. *Remote Sensing of Environment*, 145, 154–172. <https://doi.org/10.1016/j.rse.2014.02.001>
- Safford, H. D., Sawyer, S. C., Kocher, S. D., Hiers, J. K., & Cross, M. (2017). Linking knowledge to action: The role of boundary spanners in translating ecology. *Frontiers in Ecology and the Environment*, 15, 560–568. <https://doi.org/10.1002/fee.1731>
- Satterfield, D. (2020). *Dan's wild wild science journal*. AGU Blogosphere. Retrieved from <https://blogs.agu.org/wildwildscience/>
- Schöpfel, J., & Prost, H. (2016). Altmetrics and grey literature: Perspectives and challenges. *GL18 international Conference on grey literature*. Retrieved from <https://hal.univ-lille.fr/hal-01405443>
- Science Open. (2020). *Science Open*. Retrieved from <https://www.scienceopen.com/>
- Sentinel Hub. (2020). *Classification App*. Retrieved from <https://apps.sentinel-hub.com/classificationApp/#/>
- Shanley, L. A., Parker, A., Schade, S., & Bonn, A. (2019). Policy perspectives on citizen science and crowdsourcing. *Citizen Science: Theory and Practice*, 4, 30. <https://doi.org/10.5334/cstp.293>
- Sicular, S. (2016). *Three architecture styles for a useful data lake* (No. G00303817) (pp. 1–32). Gartner. Retrieved from <https://www.gartner.com/doc/3380017/architecture-styles-useful-data-lake>
- SMD Science Management Council. (2018). *Science mission directorate policy: Citizen science (SMD policy document SPD-33)*. NASA. <https://science.nasa.gov/science-pink/s3fs-public/atoms/files/SPD.33Citizen.Science.pdf>
- Smithsonian Astrophysical Observatory. (2020). *Astrophysics data system*. Retrieved from <https://ui.adsabs.harvard.edu/>
- Stodden, V., Seiler, J., & Ma, Z. (2018). An empirical analysis of journal policy effectiveness for computational reproducibility. *Proceedings of the National Academy of Sciences of the United States of America*, 115, 2584–2589. <https://doi.org/10.1073/pnas.1708290115>
- Suber, P. (2012). What is open access? In *Open access*. The MIT Press. <https://doi.org/10.7551/mitpress/9286.003.000310.7551/mitpress/9286.001.0001>
- The ENCODE Project Consortium. (2011). A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biology*, 9, e1001046. <https://doi.org/10.1371/journal.pbio.1001046>
- The International HapMap Consortium. (2005). A haplotype map of the human genome. *Nature*, 437, 1299–1320. <https://doi.org/10.1038/nature04226>
- Ton That, D. H., Fils, G., Yuan, Z., & Malik, T. (2017). Sciunits: Reusable research objects. In *2017 IEEE 13th international conference on E-science (e-Science)* (pp. 374–383). <https://doi.org/10.1109/eScience.2017.51>
- Toronto International Data Release Workshop Authors. (2009). Prepublication data sharing. *Nature*, 461, 168–170. <https://doi.org/10.1038/461168a>
- University of Cambridge. (2020). *Open research. Scholarly communication*. Retrieved from [osc.cam.ac.uk](http://osc.cam.ac.uk)
- Vicente-Saez, R., & Martinez-Fuentes, C. (2018). Open Science now: A systematic literature review for an integrated definition. *Journal of Business Research*, 88, 428–436. <https://doi.org/10.1016/j.jbusres.2017.12.043>
- Watson, J. (1990). The human genome project: Past, present, and future. *Science*, 248, 44–49. <https://doi.org/10.1126/science.2181665>
- Woelfle, M., Oliaro, P., & Todd, M. H. (2011). Open science is a research accelerator. *Nature Chem*, 3, 745–748. <https://doi.org/10.1038/nchem.1149>
- Wulder, M. A., Masek, J. G., Cohen, W. B., Loveland, T. R., & Woodcock, C. E. (2012). Opening the archive: How free data has enabled the science and monitoring promise of Landsat. *Remote Sensing of Environment*, 122, 2–10. <https://doi.org/10.1016/j.rse.2012.01.010>
- Xu, C., & Yang, C. (2014). Introduction to big geospatial data research. *Annals of GIS*, 20, 227–232. <https://doi.org/10.1080/19475683.2014.938775>
- Zhu, Z., Wulder, M. A., Roy, D. P., Woodcock, C. E., Hansen, M. C., Radeloff, V. C., et al. (2019). Benefits of the free and open Landsat data policy. *Remote Sensing of Environment*, 224, 382–385. <https://doi.org/10.1016/j.rse.2019.02.016>
- Zooniverse. (2020). *Welcome to the Zooniverse: People-powered research*. Retrieved from <https://www.zooniverse.org/>