

Journal Pre-proof

Community recommendations for geochemical data, services and analytical capabilities in the 21st century

Marthe Klöcking, Lesley Wyborn, Kerstin A. Lehnert, Bryant Ware, Alexander M. Prent et al.

PII: S0016-7037(23)00191-6
DOI: <https://doi.org/10.1016/j.gca.2023.04.024>
Reference: GCA 13033

To appear in: *Geochimica et Cosmochimica Acta*

Received date: 2 November 2022

Accepted date: 24 April 2023

Please cite this article as: M. Klöcking, L. Wyborn, K.A. Lehnert et al., Community recommendations for geochemical data, services and analytical capabilities in the 21st century, *Geochimica et Cosmochimica Acta*, doi: <https://doi.org/10.1016/j.gca.2023.04.024>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier.



Community recommendations for geochemical data, services and analytical capabilities in the 21st century

Marthe Klöcking^a, Lesley Wyborn^b, Kerstin A. Lehnert^c, Bryant Ware^d, Alexander M. Prent^{e,d,f}, Lucia Profeta^c, Fabian Kohlmann^g, Wayne Noble^g, Ian Bruno^h, Sarah Lambartⁱ, Halimulati Ananuer^j, Nicholas D. Barber^{k,l}, Harry Becker^m, Maurice Brodbeckⁿ, Hang Deng^o, Kai Deng^p, Kirsten Elger^q, Gabriel de Souza Franco^r, Yajie Gao^b, Khalid Mohammed Ghasera^s, Dominik C. Hezel^t, Jingyi Huang^{u,v}, Buchanan Kerswell^w, Hilde Kochⁿ, Anthony W. Lanati^{x,j}, Geertje ter Maat^f, Nadia Martínez-Villegas^y, Lucien Nana Yobo^z, Ahmad Redaa^{aa,ab}, Wiebke Schäfer^{ac}, Megan R. Swing^{ad}, Richard J. M. Taylor^{ae}, Marie Katrine Traun^{af}, Jo Whelan^{ag}, Tengfei Zhou^{ah}

^aGeoscience Centre, Georg-August-Universität, Goldschmidtstr. 1, Göttingen, 37077, Germany

^bResearch School of Earth Sciences, The Australian National University, 142 Mills Rd, Acton, 0200, ACT, Australia

^cLamont-Doherty Earth Observatory, Columbia University, 61 Rte 9W, Palisades, 10964, NY, United States

^dJohn de Laeter Centre, Curtin University, Building 301, Murdoch Ct, Bentley, 6845, WA, Australia

^eAuScope Ltd, Melbourne, Australia

^fFaculty of Geosciences, Utrecht University, Princetonlaan 8a, Utrecht, 3584, CB, Netherlands

^gLithodat Pty Ltd, Melbourne, Australia

^hCambridge Crystallographic Data Centre, 12 Union Road, Cambridge, CB2 1EZ, United Kingdom

ⁱDepartment of Geology and Geophysics, The University of Utah, Salt Lake City, UT, United States

^jAustralian Research Council Center of Excellence for Core to Crust Fluid Systems (CCFS) and GEMOC, School of Natural Sciences, Macquarie University, Wallumattagal Campus, North Ryde, 2109, NSW, Australia

^kDepartment of Earth and Planetary Science, McGill University, Montréal, Québec, Canada

^lDepartment of Earth Sciences, University of Cambridge, Cambridge, United Kingdom

^mGeological Sciences, Freie Universität Berlin, Malteserstr. 74-100, Berlin, 12249, Germany

ⁿDepartment of Geology, Trinity College Dublin, Dublin 2, Ireland

^oDepartment of Energy and Resources Engineering, College of Engineering, Peking University, Beijing, China

^pInstitute of Geochemistry and Petrology, Department of Earth Sciences, ETH Zürich, Clausiusstrasse 25, Zürich, 8092, Switzerland

^qGFZ German Research Centre for Geosciences, Telegrafenberg, Potsdam, 14473, Germany

^rSchool of Earth Ocean and Environment, University of South Carolina, Columbia, United States

^sDepartment of Geology, Aligarh Muslim University, Aligarh, 202001, India

^tInstitut für Geowissenschaften, Goethe-Universität Frankfurt, Altenhöferallee 1, Frankfurt, 60437, Germany

^uDepartment of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, 21218, United States

^vDepartment of Computer Science, University of Idaho, Moscow, 83844, United States

^wDepartment of Geology and Environmental Earth Science, Miami University, Oxford, OH, United States

Email address: marthe.kloecking@cantab.net (Marthe Klöcking)

^x*Institut für Mineralogie, Universität Münster, Corrensstrasse 24, Münster, 48149, Germany*
^y*IPICT, Instituto Potosino de Investigación Científica y Tecnológica, División de Geociencias Aplicadas, Camino a la Presa San José No. 2055, Col. Lomas 4a Sec., San Luis Potosí, 78216, SLP, Mexico*

^z*Department of Geology & Geophysics, Texas A&M University, College Station, 77843, TX, United States*

^{aa}*Department of Earth Sciences, Metal Isotope Group (MIG), Adelaide, SA, Australia*

^{ab}*Department of Mineral Resources and Rocks, Faculty of Earth Sciences, King Abdulaziz University, Jeddah, Saudi Arabia*

^{ac}*GeoZentrum Nordbayern, Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, 91054, Germany*

^{ad}*Department of Earth Sciences, University of Toronto, Toronto, ON, Canada*

^{ae}*Carl Zeiss Microscopy Ltd, CB23 6DW, Cambridgeshire, United Kingdom*

^{af}*Department of Geosciences and Natural Resource Management, University of Copenhagen, Copenhagen, Denmark*

^{ag}*Northern Territory Geological Survey, Darwin, NT, Australia*

^{ah}*School of Geosciences, China University of Petroleum (East China), Qingdao, China*

11 Abstract

The majority of geochemical and cosmochemical research is based upon observations and, in particular, upon the acquisition, processing and interpretation of analytical data from physical samples. The exponential increase in volumes and rates of data acquisition over the last century, combined with advances in instruments, analytical methods and an increasing variety of data types analysed, has necessitated the development of new ways of data curation, access and sharing. Together with novel data processing methods, these changes have enabled new scientific insights and are driving innovation in Earth and Planetary Science research. Yet, as approaches to data-intensive research develop and evolve, new challenges emerge. As large and often global data compilations increasingly form the basis for new research studies, institutional and methodological differences in data reporting are proving to be significant hurdles in synthesising data from multiple sources. Consistent data formats and data acquisition descriptions are becoming crucial to enable quality assessment, reusability and integration of results fostering confidence in available data for reuse. Here, we explore the key challenges faced by the geo- and cosmochemistry community and, by drawing comparisons from other communities, recommend possible approaches to over-

come them. The first challenge is bringing together the numerous sub-disciplines within our community under a common international initiative. One key factor for this convergence will be gaining endorsement from the international geochemical, cosmochemical and analytical societies and associations, journals and institutions. Increased education and outreach, spearheaded by ambassadors recruited from leading scientists across disciplines, will further contribute to raising awareness, and to uniting and mobilising the community. Appropriate incentives, recognition and credit for good data management as well as an improved, user-oriented technical infrastructure will be essential for achieving a cultural change towards an environment in which the effective use and real-time interchange of large datasets is common-place. Finally, the development of best practices for standardised data reporting and exchange, driven by expert committees, will be a crucial step towards making geo- and cosmochemical data more Findable, Accessible, Interoperable and Reusable by both humans and machines (FAIR).

Keywords: FAIR data, data standards, data quality

1. Introduction

Data are the backbone of geochemical and cosmochemical research, and their acquisition and use are central to many aspects of our research and education. Over the last century, an ever-increasing volume of geochemical data have been acquired and used to explore a variety of past, present and future processes in the Earth, environmental and planetary sciences (Fig. 1). The growing rate of data generation is complemented by new capabilities in storing, accessing, processing and modelling of large datasets (e.g. Morrison et al., 2017; Duke et al., 2022; He et al., 2022; Wieser et al., 2022).

The increasing need for globally standardised geochemical data has become a common subject of discussion amongst the international scientific community in the last few years (e.g. Stall et al., 2019; Chamberlain et al., 2021; Wyborn et al., 2021; Pourret and Irawan, 2022). Motivated by these developments, the three geochemical data

systems EarthChem, GEOROC and AusGeochem held a joint workshop at the Goldschmidt Conference 2022: “Earth Science meets Data Science: what are our needs for geochemical data, services and analytical capabilities in the 21st century?” (<https://conf.goldschmidt.info/goldschmidt/2022/meetingapp.cgi/Session/3301>). This workshop primarily focused on exploring the data and infrastructure requirements for addressing future scientific challenges. More information about the workshop programme, participating data systems and attendees is available in the Supplementary Material. This paper summarises the workshop outcomes and provides recommendations for a global geochemical data framework, required to tackle and accomplish the scientific challenges of the 21st century and beyond.

2. Motivation

2.1. Diversity and Fragmentation of Geochemical Data

We understand geochemistry as the discipline that integrates geology and chemistry by using the principles and tools of chemistry to develop fundamental understanding of the dynamics of geological systems, from the interior of the Earth to its surface environments on land, in the oceans, and in the air, to planetary systems and the entire galaxy. Geochemistry emerged as a discipline of its own in 1838 and, since then, acquisition and analysis of geochemical data have become pervasive in the Earth, environmental, and planetary sciences (Fairbridge, 1998). Geochemistry is exceedingly diverse with many recognised subdisciplines, including aqueous, organic, inorganic, isotope, bio- and physical geochemistry as well as cosmochemistry. Geochemical data have further applications in other disciplines such as archaeology, environmental science and technology, resource exploration and development (groundwater, minerals, energy), geohealth, oceanography, and agriculture, and are thus relevant to many United Nations Sustainable Development Goals (e.g. Bundschuh et al., 2017; Gill, 2017; Alexakis, 2021; Wyborn and Lehnert, 2021).

Geochemical data are incredibly diverse in nature and generally only have two common

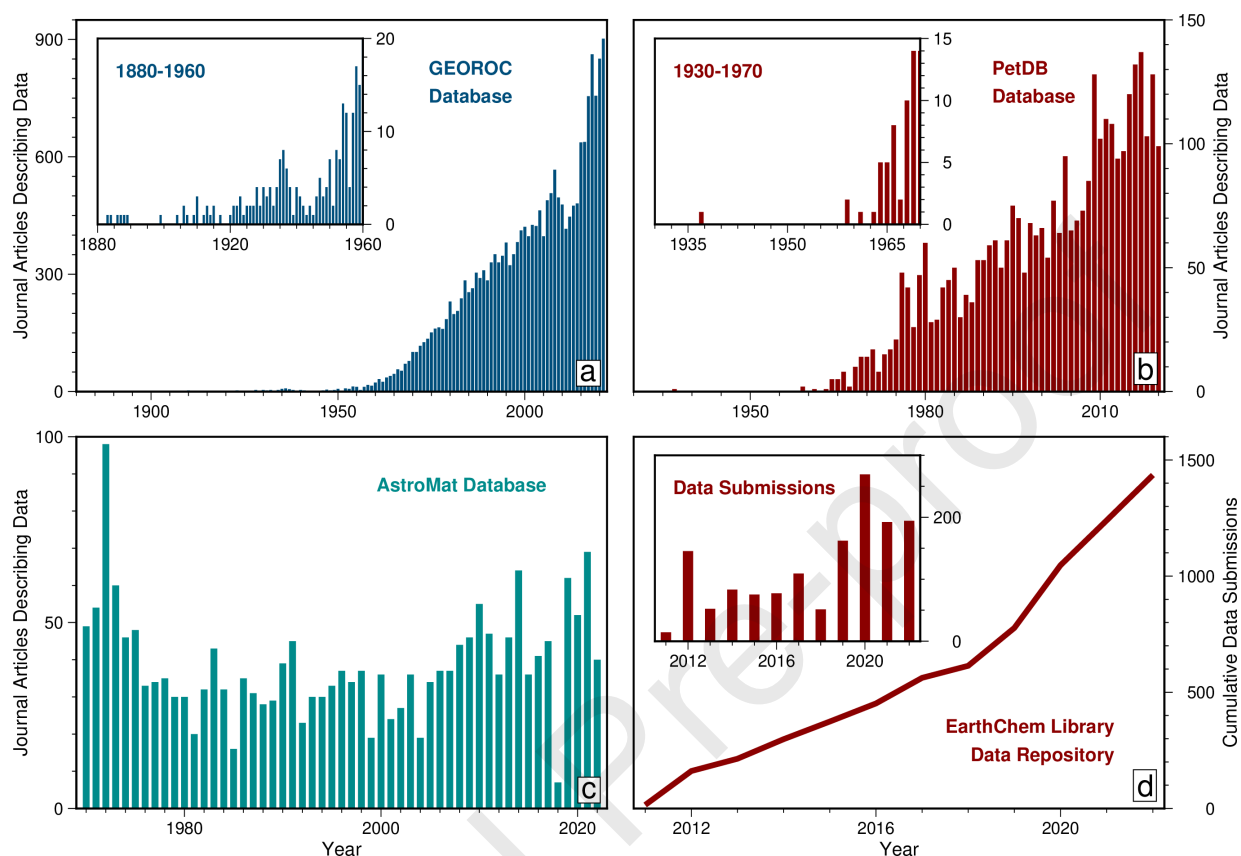


Figure 1: Increase in geochemical data published in journals and repositories since the late 19th century. **(a)** Data compiled within the GEOROC database, by publication year of the respective journal articles, as a proxy for the increase in data production within the subdiscipline of igneous geochemistry in the continental realm. Inset: Close-up of earliest publication years. **(b)** Data compiled within the Petrological Database (PetDB) which contains data complementary to GEOROC with a focus on the oceanic realm, mantle xenoliths and tephra. Inset: Close-up of earliest publication years. **(c)** Data compiled within the Astromaterials Data System, including data from the MetBase database, as a proxy for data production within cosmochemistry. **(d)** Cumulative number of data submissions to the EarthChem Library, a domain repository for all subdisciplines of geochemistry. Inset: individual number of data submissions per year.

attributes: firstly, they are “Long Tail”, i.e. highly variable and small in volume (Heidorn, 2008); and secondly, they are primarily acquired by individual investigators or small teams, often across multiple organisations and disciplines with uncertain funding sustainability. Due to this diversity, many geochemical datasets are stored in incompatible and often inaccessible silos, such as individual computers and locally developed database solutions, or they are restricted to figures without accompanying data tables. As a consequence, and despite numerous data rescue efforts, harnessing the wealth of existing geochemical data is a critical and ongoing challenge.

Although there have been many attempts to improve the aggregation, sharing and reuse of geochemical data (e.g. Wyborn and Ryburn, 1989; Carbotte and Lehnert, 2007; Geochemical Society, 2007; Goldstein et al., 2014), present-day practices tend to focus on building geochemical databases in either personal, institutional, national, or programmatic silos with a noticeable divide in approaches to data management among the sectors of academia, government and industry. Most of these databases are built for specific research projects and do not offer a long-term sustainable solution. There are very few standard practices amongst authors and publishers to make data easily shareable and interoperable. As a result, geochemical data are highly fragmented, blocked from discovery and difficult to reuse directly from the source dataset without considerable efforts in reformatting the data. Moreover, the same data are duplicated numerous times into multiple compilations and credit is rarely given to those who funded, collected, and/or analysed the original datasets. This fragmentation has a measurable financial impact: the European Commission estimated the annual direct cost of managing non-standardised research data at EUR 10.2bn, with an additional indirect cost to society of EUR 16bn per year (European Commission, 2018).

2.2. Drivers and Rationale for Connecting the Silos

A number of important resources for geochemical and cosmochemical data were established during the past 30 years, including EarthChem (<https://earthchem.org/>),

78 GEOROC (<https://georoc.eu/>), MetBase (<https://metbase.org/>), and the Astroma-
79 terials Data System (<https://www.astromat.org/>). More recent initiatives are National
80 Research Infrastructures in Germany (NFDI4Earth), Europe (EPOS), Australia (AuS-
81 cope), the US (EarthCube), or Norway (NIRD), to name a few. However, barriers around
82 individual data silos remain, hindering simple, inclusive and global access to geochemical
83 data. To overcome these silo walls, we must develop and implement common, community-
84 agreed, global standards for geochemical data and metadata. These standards are critical
85 to making geochemical data Findable, Accessible, Interoperable and Reusable to both hu-
86 mans and machines (FAIR; Wilkinson et al., 2016). Not only will FAIR data standards and
87 curation procedures increase the value of new data as they are generated and published,
88 they likewise have large potential for utilising the significant proportion of unpublished
89 geochemical data in research and public sectors from the last century.

90 Recognising that mainstream scientific journals were the most effective agents to rectify
91 problems in data reporting and implement best practices, an Editors Roundtable was held
92 in 2007 as an initiative to bring together editors, publishers, and database providers to
93 implement consistent publication practices for geochemical data. Academic societies such
94 as the Geochemical Society also adopted a policy for geochemical data publication at that
95 time (Geochemical Society, 2007). The Editors Roundtable created and signed a policy
96 statement in January 2009 (version 1.1) that laid out ‘Requirements for the Publication
97 of Geochemical Data’ (Goldstein et al., 2014). Unfortunately, even 14 years on these
98 recommendations are rarely followed.

99 Recently, the nationally-funded, global data systems Astromaterials (USA), Earth-
100 Chem (USA), GEOROC (Germany), EPOS-MSL (European Plate Observing System Mul-
101 tiScale Laboratories, Europe), MetBase (Germany) and AusGeochem (Australia) came
102 together to enable interoperability between their systems. Yet a vast amount of geo-
103 chemical data lies outside these initiatives. In response to Open Science policies and
104 demands from the scientific community, a Town Hall meeting on ‘OneGeochemistry: To-

ward a Global Network of Geochemistry Data’ was held at the AGU Fall Meeting 2019 to raise awareness of the increasingly urgent need for global standards and best practices for geochemical data— aiming towards better sharing and linking of data resources into a global network (<https://www.agu.org/Fall-Meeting-2019/Events/Data-TH23L>). The goal of this meeting was to broaden community awareness of and participation in the initiative and speakers represented relevant stakeholders such as geochemical societies, geochemical journal editors, data infrastructure providers, researchers, and funders. The OneGeochemistry initiative was launched. Since then, the OneGeochemistry initiative regularly leads and contributes to scientific sessions during Goldschmidt, EGU and AGU meetings— including a Great Debate and Webinar at EGU22 (‘Where is my data, where did it come from and how was it obtained? Improving Access to Geoanalytical Research Data’; <https://meetingorganizer.copernicus.org/EGU22/session/42788>; <https://www.youtube.com/watch?v=nqjp0ePQU0w>)— as well as international fora such as SciDataCon and the International Science Council’s Committee on Data (CODATA) meetings (e.g. Lehnert et al., 2021; Wyborn et al., 2021).

2.3. OneGeochemistry Mission

OneGeochemistry is an international collaboration between multiple national organisations that support geochemistry capability and data production. The focus of this initiative is to better coordinate global efforts in geochemical data standardisation, facilitate communication between groups and lessen duplication of efforts. OneGeochemistry is now taking action, predominantly through volunteer work of its member organisations, to collect, synthesise and promote global, community-driven data conventions and best practices. Such global best practices will enable and simplify the (re)use of geochemical data, making possible a global network of trusted geochemical data, which will accelerate the generation of new geoscientific knowledge and discoveries.

Data standardisation begins with community agreement on concepts and vocabularies used to describe analytical data. Such vocabularies are critical to organise and classify

data: they set out the common terminology. We require experts for each data type to come together to develop the required vocabularies in both human and machine readable forms, whilst building on and integrating existing definitions from the broader geoscience terminology and other related domains. The community must then agree to use these vocabularies to refer to their concepts of interest, as well as evolve and govern them as requirements change.

In line with modern informatics best practices, all geochemical data will need to comply with the FAIR principles of Wilkinson et al. (2016). OneGeochemistry seeks to make geochemical data outputs as well as related inputs (including samples, instruments, software codes):

1. **Findable (F)** through machine-actionable metadata and the systematic use of unique and persistent identifiers on inputs and outputs;
2. **Accessible (A)** using standards and internet protocols;
3. **Interoperable (I)** through common formats that incorporate authoritative and referable domain vocabularies; and
4. **Reusable (R)** through use of rich metadata that provide guidelines on provenance, quality and uncertainty, that clearly show identity, funders, and provide open licences.

It is also essential to ensure compliance with the CARE and TRUST principles. The CARE Principles for Indigenous Data Governance (Collective Benefit, Authority to Control, Responsibility, and Ethics) protect Indigenous rights and interests in Indigenous data including traditional knowledge, particularly in the sample collection phase (Carroll et al., 2020). The TRUST Principles (Transparency, Responsibility, User focus, Sustainability and Technology) ensure long-term data preservation and trustworthiness in digital repositories. (Lin et al., 2020).

Efforts have already been made to set standards for specific analytical data types: Deines et al. (2003); Demetriades et al. (2020, 2022); Boone et al. (2022); Flowers et al. (2022); Brantley et al. (2021); Abbott et al. (2022); Horstwood et al. (2016); Dutton

et al. (2017); Walker et al. (2008); Courtney Mustaphi et al. (2019); Schaen et al. (2020); Khider et al. (2019); Damerow et al. (2021); Peng et al. (2022); Wallace et al. (2022). These publications are an excellent first step, however they only cover a subset of the chemical data types and very few conform with the FAIR principles that require data to be machine readable. Hence, these standards need to be converted into the digital space (e.g., the IUPAC Digital Chemistry Initiative; <https://iupac.org/what-we-do/digital-standards/>). The next step towards standardisation of geochemical data is to follow Cox et al. (2021) and make the vocabularies, recommended within each standard to define different data types, FAIR and available from online repositories such as Research Vocabularies Australia (RVA, <https://vocabs.ardc.edu.au/>) or FAIRsharing (<https://fairsharing.org/>). Another important point often missing in existing recommendations is a governance structure that allows vocabularies and best practices to evolve.

OneGeochemistry aims to become an organisation that coordinates across all geo- and cosmochemical data types, both supporting existing community standards as well as facilitating the development of new ones where needed. Importantly, OneGeochemistry will act as the facilitator in these efforts: the initiative will neither set standards nor implement them, but rather support the community in doing so. A starting point will be to support the digitisation of existing standards to make them, and the vocabularies defined within them, fully FAIR. Fundamental to OneGeochemistry's approach is ensuring that networking common components across disciplines still enables a capacity for deeper disciplinary specialisation. This will be an ongoing, long-term project that must be continually adapted in line with new or improved developments of data acquisition and with support of, and commitment from, the global geochemical and cosmochemical communities.

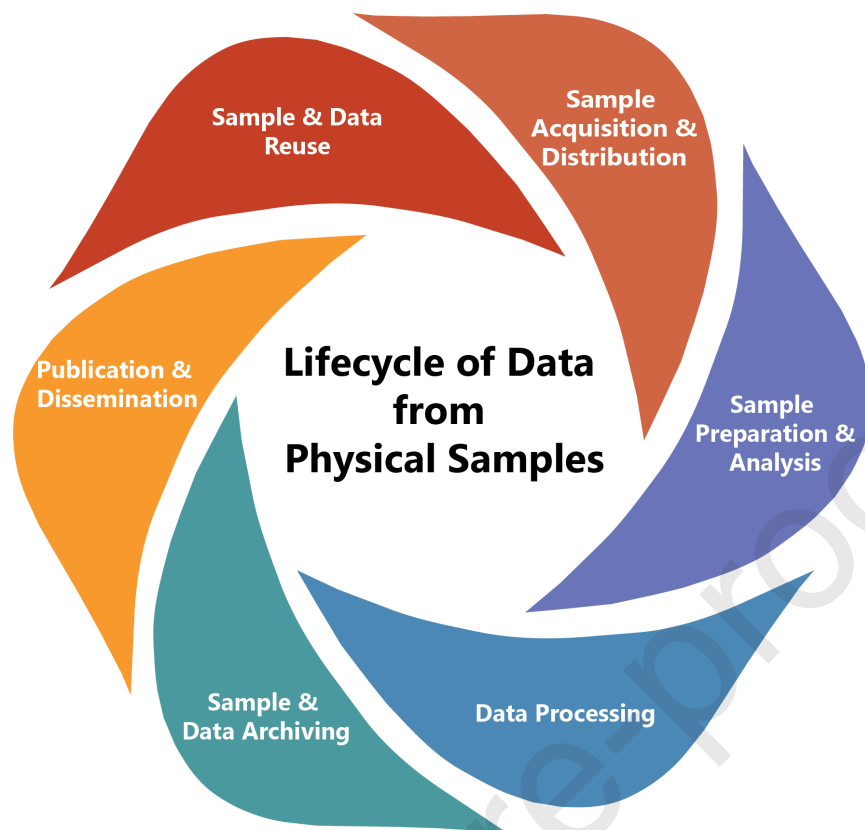


Figure 2: The sample and data life cycle from acquisition to publication to reuse (adapted from Ramdeen et al., 2022). Tools that support researchers throughout this process include SESAR, a registry for physical samples. AusGeochem, StraboSpot and Sparrow are examples of systems that support researchers from field acquisition of samples through sample preparation and analysis to publication in a domain repository. Repositories such as the EarthChem Library serve the Archiving and Publication of Data, while synthesis databases such as the Astromaterials Data Synthesis, PetDB, GEOROC or MetBase facilitate dissemination and data reuse.

3. Challenges for the Community

This paper tackles challenges faced by both the active research community (predominantly at academic and government institutions) and the curated data systems that support this community throughout the research data lifecycle. These data systems can be grouped into four types: 1) Laboratory Information Management Systems, 2) Repositories, 3) Data Portals, and 4) Synthesis Databases. Firstly, Laboratory Information Management

189 Systems focus on physical samples and cover the first half of the research data lifecycle
 190 from sample collection or generation to processing and analysis (Fig. 2). Examples of
 191 such systems include AusGeochem (<https://www.auscope.org.au/ausgeochem>), Stra-
 192 boSpot (<https://www.strabospot.org/>) and Sparrow (<https://sparrow-data.org/>).
 193 Secondly, the final data products derived from samples might then be published in Repos-
 194 itories as well as cited in journal publications. Generalist repositories, such as Figshare
 195 (<https://figshare.com/>), Dryad (<https://datadryad.org/>) or Zenodo (<https://zenodo.org/>), publish research outputs irrespective of academic discipline and without review. Do-
 196 main repositories, in contrast, cater to specific disciplines or subdisciplines and therefore
 197 offer data services targeted to the particular requirements of these domains. PANGAEA
 198 (<https://www.pangaea.de/>) and GFZ Data Services (<https://bib.telegrafenberg.de/dataservices/>) are examples of domain repositories for the Earth Sciences, whilst the
 199 Astromaterials Data Repository (<https://repo.astromat.org/>), the EarthChem Library
 200 (<https://earthchem.org/ec1/>) or the GEOROC Data Repository (<https://georoc.eu/>) are domain repositories specifically for geochemical data. Thirdly, Data Portals
 201 offer a catalogue of datasets hosted by different repositories. For example, DataONE
 202 (<https://dataone.org/>) searches across 44 data repositories of all disciplines operated by
 203 research centres, universities, libraries, scientific consortia, non-profit organisations, citizen
 204 science initiatives, corporate divisions, governmental and non-governmental organisations.
 205 Such data portals greatly increase the discoverability of data products stored in the respec-
 206 tive systems by searching through their metadata catalogues, including the title, abstract or
 207 keywords of individual datasets. Finally, Synthesis Databases compile individual data pub-
 208 lications and harvest data from the scientific literature to enable data discovery and reuse
 209 across multiple datasets. In contrast to data portals, synthesis databases do not only sup-
 210 port searches across the metadata of datasets in multiple repositories (e.g. title, keywords,
 211 etc), they further compile the actual data held in each of these records and allow download
 212 of single, combined datasets. Similar to domain repositories, synthesis databases usually

specialise in a particular subdiscipline or have a geographical focus. However, in contrast to repositories they do not serve as a data publisher but instead only focus on synthesising and compiling previously published data. Note that we do not consider research datasets derived from literature compilations as databases here as they usually are ephemeral, one-off research products that are not continuously curated and more importantly, rarely uniquely identify each analysis so that the author and funder can track citations and measure impact. The Astromaterials Data Synthesis, GEOROC, LEPR (<https://lepr.earthchem.org/>), MetBase and PetDB (<https://search.earthchem.org/>) are all examples of synthesis databases. These databases provide valuable resources not only for further research but also for teaching. Both repositories and synthesis databases also play an important role in data rescue efforts. Figure 3 shows an example of the flow of geochemical data from natural samples through the IEDA2 (Interdisciplinary Earth Data Alliance) and affiliated data systems.

In an ideal world, all analytical data produced in a laboratory and subsequently published in the scientific literature, would eventually be made available in a federated, global data system that makes it easy for others to find, access and reuse these data. Features of such an ideal data system include:

- 1. Relevance & Findability:** A variety of data types are available for all types of sample material (natural and synthetic). It is easy to combine multiple databases to search, capture and organise all existing data. These databases contain minimal redundancy and the use of globally unique, persistent and resolvable identifiers (e.g. digital object identifiers, DOIs, and the international generic sample number, IGSN) allows compilation of analyses from the same sample or publication. Database versioning allows reproducibility of previous searches.
- 2. Accessibility:** User access is facilitated by optimised complex queries, for example through a customisable search engine, visualisation, data analysis and export options. Access through standard programming languages guarantees machine-readability.

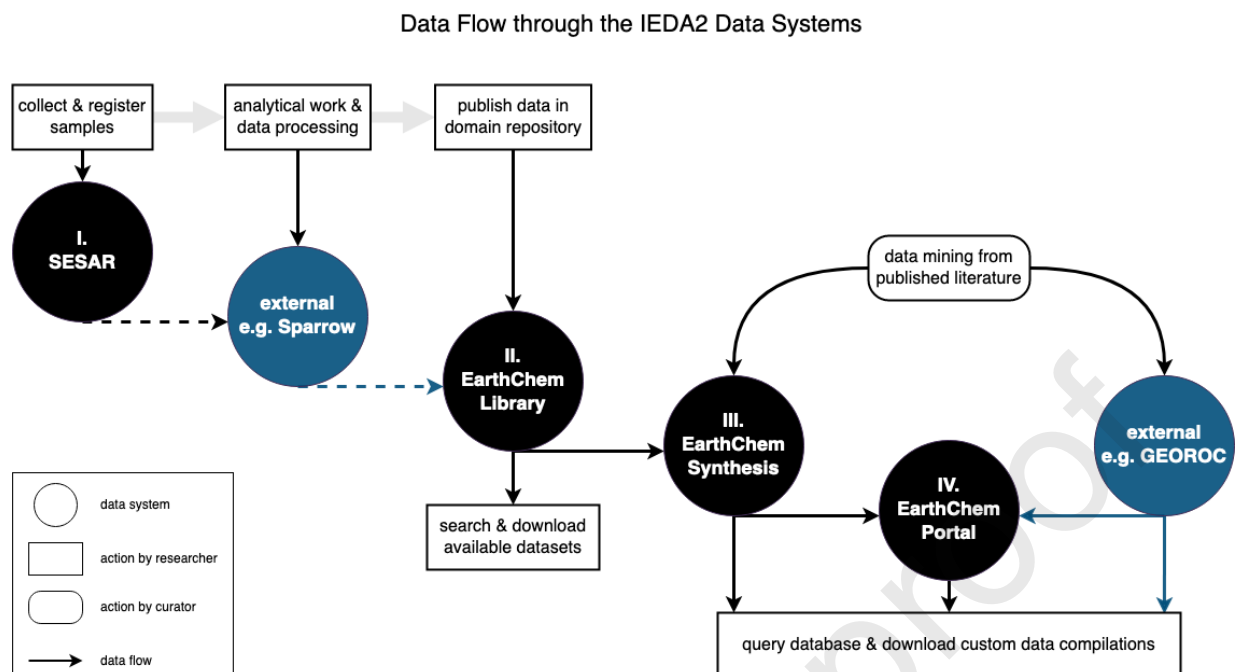


Figure 3: An example of the flow of geochemical data from natural samples through the IEDA2 (black) and partner (blue) data systems. Together, these data systems cover the entire research data lifecycle as shown in Fig. 2. Note that the EarthChem Portal enables data searches across distinct synthesis databases, in contrast to the data portals described in the text that facilitate metadata searches across different repositories. Not included in this schema is the Library of Experimental Phase Relations (LEPR) for experimental and synthetic materials. For comparison, the AusGeochem system covers stages I to III of this diagram for data produced by Australian geochemistry laboratories.

Furthermore, access is free and open to all: there should be no cost to the researcher in either publishing or accessing data.

3. **Data Quality:** Data are reliable and their quality is straightforward to assess, i.e. they follow a common standard that ensures availability of rich sample and analytical metadata (e.g. provenance, description of method and analysis conditions). Completeness of metadata allows assessment of accuracy and precision, and ensures reproducibility. Both data providers and data users perform QA/QC; any data quality issues are reported and promptly resolved.
4. **Attribution:** Appropriate citation of the people, laboratories, organisations, fun-

ders, research artefacts and data is ensured through use of globally unique, persistent and resolvable identifiers and compliance with international metadata standards (e.g. the IGSN for samples, the Open Researcher and Contributor Identifier, ORCID, for authors, the Research Organization Registry, ROR, for institutions; or the DataCite metadata standard).

Many of the data systems mentioned above strive to provide such a comprehensive data infrastructure. It is now increasingly recognised that data and metadata capture should start with the collection/production of the sample itself, and not only after data publication (e.g. Damerow et al., 2021). However, there are many challenges along the path towards FAIR geochemical data, many of which have been introduced above. One of the goals of the Goldschmidt 2022 workshop was to investigate these challenges in more detail, so that appropriate solutions for each of them might be developed. These challenges are rooted in the current research culture around geoanalytical data, as well as the limitations of the existing data systems and their often precarious funding situation.

3.1. Challenges for Researchers

The current research culture in geochemistry means that few researchers are willing to share their data (Chamberlain et al., 2021). Although the recent push for open science has benefited the open data landscape, community understanding and adoption are still centred around individuals. The majority of data producers remain reluctant to share their data unless forced by journal or funding requirements: the EarthChem Library reported an increase in data submissions after several of the AGU journals enforced data publications in trusted domain repositories in 2019 (Fig. 1d; https://www.agu.org/Share-and-Advocate/Share/Polymakers/Position-Statements/Position_Data). Nevertheless, there is still a widespread lack of adoption of these policies by the research community. Common barriers to data sharing include the additional effort of organising and formatting of data, distrust and protection of personal interests, e.g. with additional work in progress, insecurity about copyright and licensing, lack of knowledge about the most appropriate

repository, lack of time, as well as the costs of sharing data (Stuart et al., 2018; Science et al., 2021; Tedersoo et al., 2021). Yet even those researchers who are willing to share their data are faced with a number of considerable challenges that we discuss in the following.

Lack of consistent guidelines: Policies on data management vary widely amongst the different funding agencies, institutions, publishers and journals. Funders often require a data management plan at the proposal stage, yet few enforce these requirements once grants are approved. Researchers are neither penalised nor rewarded in response to how they manage their data, prompting the question as to why this requirement exists in the first instance if there is no mechanism for ensuring compliance. In addition, institutional open access policies often do not extend to include research data or a requirement for machine-readable formats— a PDF-copy of published journal articles in the institutional repositories is usually enough to fulfil these guidelines. This effect is compounded by many institutions lacking the resources to support their researchers in appropriate data management. Finally, the publishing landscape is as diverse as the journals available. Each publisher has defined their own policies on data management, and often these guidelines differ for each journal even with the same publisher. The publishers Springer Nature, American Association for the Advancement of Science (AAAS) and American Geophysical Union (AGU) are proponents of consistent data management practices, requiring data publication in domain repositories prior to manuscript acceptance across many of their journals, yet each have developed their own— differing— guidelines on how to comply with this policy. Dedicated data journals, such as Data in Brief (Elsevier) and Scientific Data (Springer Nature), perhaps present a good alternative in requiring data submission to (domain) repositories and, in addition, providing a platform for publishing and describing data that might otherwise never be made public— for example, data from unfinished or abandoned thesis projects or those transcribed from old, non-digital formats. However, most other journals still accept data tables in formats ranging from tabular (CSV or XLS) to text (DOC, PDF) and even image files (JPEG, PNG) as part of supplementary mate-

rials or they encourage submission to generalist repositories, such as Figshare, Zenodo or
 Dryad, where there is no quality control or agreed reporting standards on geochemical data.
 Researchers, therefore, are faced with the impossible task of navigating these conflicting
 guidelines, and will generally follow the policy of the journal or publisher they submit to
 out of fear that their manuscript might otherwise be rejected. When faced with the com-
 plexity of submission to domain repositories (see below), often the publishing option with
 the lowest workload is chosen. This behaviour naturally leads to highly heterogeneous data
 published following very different standards, if any, in very different formats across a wide
 range of repositories or other data publishers. In addition to the many different formats
 that prevent data from being easily combined and compared, many datasets remain behind
 a journal paywall and are very hard to access in the first place. Data availability “upon
 request” also remains a popular option, even though it has been shown to be burdensome
 and ineffective as a means for data sharing (Vines et al., 2014; Tedersoo et al., 2021). Even
 for Science, a journal that adopted an open data policy in 2016, 30% of articles do not
 publish their data at all, and only for about a quarter of articles can research findings be
 accurately reproduced (Stodden et al., 2018; Yeston, 2021).

Complexity of data submission: Good data management takes time. The assem-
 bling and submission of data tables and related information require time and additional
 effort outside of the primary process of manuscript submission. Usually, substantial pro-
 cessing is performed on raw data coming from an analytical instrument. While this process-
 ing is a common research practice, information on data reduction and reference materials
 used are often not reported, or only a simplified version is included in the methods or sup-
 plementary information. Yet, reporting this information is crucial for the reproducibility
 of data and, therefore, a prerequisite for data submission to domain repositories. This
 considerable, additional investment of research time and resources is often voluntary, and
 not appropriately rewarded within the current academic structure (Piwowar et al., 2007;
 Kim and Stanton, 2012). Even though data publications are increasingly visible via (au-

333 tomatic) indexing in ORCID profiles, for example, they are rarely counted towards the
334 research track record or valued by recruiting and promotion committees. Whilst assigning
335 DOIs to datasets helps to emphasise the value of data publications, the lack of awareness in
336 the broader research community means that these publications are often not appropriately
337 cited. In addition, researchers who consider submitting to domain repositories are often
338 deterred by the additional processing time before the final data publication. The Earth-
339 Chem Library, for example, that specialises in geochemical data, advises a turnaround
340 time ranging from a few days to up to two weeks. PANGAEA, a domain repository for
341 all disciplines within the Earth Sciences, has a data publication timeline of three months.
342 Even though there are good reasons for these timelines— mostly centred around curation
343 as discussed below—, they discourage even more researchers from publishing their data.

344 **Variable quality of the available published data:** A direct result of the lack of
345 guidelines combined with the complexity of data submission is the highly variable quality
346 of the available datasets. The lack of enforced standard formats for publishing geochem-
347 ical data often precludes any quality assessment and, therefore, reuse of published data.
348 Common issues include: dead links or non-existent supplementary material; errors in data
349 reporting; lack of reproducibility due to missing analytical information; and the use of unde-
350 fined abbreviations only understood by the owner of the dataset. Data quality assessment
351 is often impossible due to a lack of analytical details or measures of uncertainties, includ-
352 ing inconsistent units on uncertainty reporting (e.g. standard deviations, standard errors,
353 confidence interval, 1σ vs. 2σ errors, etc.). When compiling data from multiple sources,
354 additional challenges include inconsistent, non-standardised terminology (e.g. eclogite vs
355 arclogite) and missing units of measurement. Finally, the original owner, funder, and/or
356 creator of the data are rarely credited in compiled datasets.

357 **Complexity of citation for data compilation work:** The inclusion of all refer-
358 ences to the original data sources in published data tables, which is common standard for
359 data collections, does not automatically provide credit in measurable form. In order for

these citations to be tracked, references must be included in the ‘References’ sections of scholarly literature. Unfortunately, journals commonly limit the total number of citations allowed (often between 40–70) and ask authors to move any additional references into the supplemental information. Yet, references in supplemental information are not properly indexed, not linked to the manuscript, nor tracked accurately— all of which is essential to enable reproducible research and for researchers and institutions to trace data usage and receive appropriate credit for their work. The new “Complex Citations Working Group” of the Research Data Alliance (RDA; <https://www.rd-alliance.org/groups/complex-citations-working-group>) is currently developing a method for handling the citation of large numbers of objects— particularly datasets, software, and physical samples— in scholarly work (Agarwal et al., 2021). They propose the term ‘reliquary’ to describe a collection/package of aggregated individual datasets that make up a data compilation used within a specific article. By citing this ‘data reliquary’, all component datasets would also receive a citation without needing to be included in the article reference list. Work by the RDA group now focuses on (1) the development of a scalable solution and the infrastructure to enable credit for each individual element of this ‘reliquary’, and (2) its acknowledgement and implementation by journals.

Sensitive data: Finally, an important consideration within both the FAIR and CARE principles is how to handle sensitive data that should only be discoverable by certain, authorised persons or only available after an embargo period. This access control is particularly important for geochemical data produced or funded by industry and for agencies that deal with classified information. Fortunately, good technical solutions already exist, simply requiring clear licensing of datasets and the ability of repositories to handle management of temporary embargo periods during the publication phase. Such solutions are already implemented in many geochemical repositories, including, for example, CUAHSI HydroShare (<https://www.hydroshare.org/>) or the EarthChem Library.

3.2. Challenges for Data Systems

Some of the challenges for researchers detailed above are related to current limitations of data repositories and synthesis databases. One major issue lies with the resources available to these data systems and the sustainability of funding. Long-term staffing solutions for data curators that assist researchers with data submissions are vital for data systems. The advantage of publishing data in domain repositories is that the research data are documented in a format specific to the discipline and the respective data type, which ensures that data quality can be easily assessed and data users have greater trust in individual datasets. By collecting data in domain repositories, they are also more visible and easier to discover for others in the field. Even though data sharing practices vary widely between scientific disciplines, the greater discoverability of datasets published in curated domain repositories often leads to greater reuse—and ultimately citation—of these data and the associated publications (e.g. Piwowar et al., 2007; Science et al., 2021). Yet in order to consistently provide this service, domain repositories need to employ curators with domain expertise who carefully review each data submission. Many researchers of today are not familiar with all intricacies of data management, and hence data submissions are often not consistent. While it takes the researchers a considerable amount of time to collate this information, repository curators then need to invest further time to convert submissions to their internal standard and ensure all data and metadata are transparent and easy to understand by third parties.

More often than not, repositories are not funded for this additional work and are struggling with staffing issues. These problems arise because many of the data systems catering to a specific domain were born out of research projects that succeeded in attracting additional funding to further develop their infrastructure. However, this funding is usually temporary and restricted to the development of new technologies or services—system maintenance and curation are rarely funded by national science foundations. What is more, these data systems compete for funding with researchers within their domain. Although

it has long been recognised that the benefits of open data infrastructure, and the measurable resources saved by their existence, far outweigh the costs of building and maintaining this infrastructure (e.g. Ball et al., 2004), most data systems still struggle for long-term survival. Far too often, data systems that are widely used by the research community are orphaned because of discontinued funding: MetPetDB and SedDB are pertinent examples of such systems that are no longer maintained, and at worst are no longer available to the community.

The availability of resources is intricately linked with community-uptake of domain repository services. For many data systems, it is an ongoing struggle to entice more researchers to submit their data, something which they require as an indicator for their success and continued funding. With additional resources, data systems could better raise awareness within the community, as well as expand their user support, in turn increasing the number of datasets submitted by researchers. Ideally, resources would also be allocated to provide training materials and build guided workflows that operate across repositories and other publication platforms to make it easy for researchers to follow best practices.

4. Approaches to similar challenges in other communities

Despite the various challenges outlined in the previous section, this topic is not new and other disciplines have successfully begun adopting FAIR data practices. In analytical science, particularly where the same data type is collected by multiple laboratories and institutions, informed decisions on whether or how to (re)use any digital analytical dataset is dependent on a consideration of what practices have been used to obtain the data and the provision of information about the quality specifications (Peng et al., 2022). The following summarises successful approaches to data standardisation and quality assurance in other communities that the geochemistry community can learn from.

4.1. Chemistry

The International Union of Pure and Applied Chemistry (IUPAC) has a record of over 100 years in fostering a global consensus to define and develop a common and systematic nomenclature for chemistry. IUPAC has developed the International Chemical Identifier (InChI; Heller et al., 2013), a non-proprietary identifier for chemical substances that provides a standard way to encode molecular information. IUPAC has also produced a series of colour books that are regarded as the world's authoritative resource for chemical nomenclature, terminology, and symbols. International committees of experts in the relevant sub-disciplines of chemistry draft the recommendations that are then ratified by IUPAC's Interdivisional Committee on Terminology, Nomenclature and Symbols (ICTNS; <https://iupac.org/body/027/>). The Terminology definitions are published by IUPAC and include books for

1. Naming Chemical Structures

- Blue Book: Nomenclature of Organic Chemistry
- Red Book: Nomenclature of Inorganic Chemistry
- White Book: Biochemical Nomenclature

2. Describing Chemistry Concepts:

- Orange Book: Terminology for Analytical Methods
- Purple Book: Polymer Terminology and Nomenclature
- Silver Book: Properties in Clinical Laboratory Sciences
- Green Book: Quantities, Units and Symbols in Physical Chemistry

Other IUPAC initiatives include the Gold Book Compendium of Chemical Terminology (<https://goldbook.iupac.org/>), the Commission on Isotopic Abundances and Atomic Weights (<https://www.ciaaw.org/>) and the Machine Actionable Periodic Table (<https://pubchem.ncbi.nlm.nih.gov/ptable/>). Advancement of digital activities and strategy

within IUPAC largely sits with the Committee on Publications and Cheminformatics Data Standards. IUPAC is currently transforming from a Centre of Excellence for Chemistry Standards to a Centre of Excellence for Digital Chemistry Standards. Many of their digital standards could be leveraged by the global geochemistry community (Stall et al., 2020).

IUPAC is primarily a volunteer-based organisation with a modest amount of project funding primarily supported through subventions paid by its member bodies (chemical societies or national academies, and some publications income). A small staff office supports the organisation generally, but most volunteers utilise basic infrastructure of their organisations while they work on projects. After the life of the projects, standard specifications are generally available as open access publications. Further development and ongoing support are primarily coordinated through partnerships with external and affiliated organisations. For example, the InChI Trust is a member-supported charity organisation affiliated with IUPAC who develops and maintains the code-base that encapsulates the IUPAC InChI standard specification. Organisations contributing to the InChI Trust include journal publishers, chemical societies, government organisations, software vendors and academic organisations.

4.2. Crystallography

Crystallography has a long history of discipline standardisation starting with development of the Crystallographic Information Framework (CIF) in 1991 under the auspices of the International Union of Crystallography (IUCr). The CIF standard is a general, flexible and easily extensible free-format archive file that was designed to be a machine-readable standard for submissions to *Acta Crystallographica* and to crystallographic databases (Hall et al., 1991). A CIF dictionary also stores the name, version and time of update, thus enabling precise citation of the standards used to support a particular data set (Hall and Cook, 1995; Hall and McMahon, 2016). Domain repositories ensure the long term preservation and access to derived results and processed data published in standard formats (Bruno et al., 2017; Groom et al., 2016; Bergerhoff and Brown, 1987; Berman et al., 2003). These

crystallographic repositories also support joint workflows with journal publishers that lower technical barriers to data publication by researchers. Further, domain repositories provide services that enable the discovery and reuse of both data and derived knowledge across domains in academia and industry (Taylor and Wood, 2019). For example, the IUCr is taking a lead in ensuring that the preservation of raw diffraction data is viable at a number of distributed and centralised data archives, each of which registers a dataset and uniquely identifies it with a persistent identifier (Kroon-Batenburg et al., 2022). The IUCr provides tools with online validation checks and validation of the data is part of the peer review process for journals (Spek, 2020). Some journals that publish papers on crystallography also sponsor the development of validation tools.

Data infrastructure in crystallography is funded through a variety of mechanisms including research grants, subscription and licensing, and governmental support (Bruno et al., 2017). The development of standards in crystallography is supported by IUCr, with the checkCIF service being supported by sponsorship from publishing organisations. Standard activities also rely heavily on volunteer effort as the scientific unions are limited in the level of support and coordination they can provide. The work of the Worldwide Protein Data Bank (wwPDB) in structural biology is primarily supported by direct funding from government. Conversely, data organisations supporting chemical crystallography do not receive direct public funding and must generate their own revenue, which is typically done by charging industry and academia for access to value-added software and services.

4.3. Seismology

Another example in the development of global community standards for a geoscience data type has been the International Federation of Digital Seismograph Networks (FDSN; <https://www.fdsn.org/>) which is a commission of the International Association for Seismology and Physics of the Earth's Interior (IASPEI) of the International Union of Geodesy and Geophysics (IUGG). The FDSN began in 1984 when multiple countries agreed to create a global network around those scientists using broadband instrumentation compatible with

community developed specifications (Dziewonski, 1994). In 1987 expert groups within the FDSN were instrumental in the development of a universal standard for the distribution of broadband waveform data and related parametric information, the SEED format (Standard for Exchange of Earthquake Data). The SEED format was adopted by instrument manufacturers and has since gone through several evolutions. The FDSN also developed a specification that defines RESTful web service interfaces for accessing common FDSN data types online and publishes a list of Federated Data Centres that provide FDSN-compliant web services (<https://www.fdsn.org/webservices/datacenters/>). Network operators can apply for FDSN Network codes through the FDSN website to provide unique identifiers for seismological data streams, which are required in publications to uniquely identify and attribute the networks that generated the data (Evans et al., 2015). FDSN is an international non-governmental organisation with volunteer membership (Suárez et al., 2008). All funding is derived from voluntary contributions by member institutions.

4.4. Geological Map Data

In 2003, the GeoSciML (Geoscience Markup Language) project was initiated under the auspices of the Commission for Geoscience Information (CGI) working group on Data Model Collaboration and endorsed by the International Union of Geological Sciences. GeoSciML is an XML-based data transfer standard for the exchange of digital geoscientific information, which is mainly focussed on the representation and description of features found on geological maps, but is extensible to other geoscience data such as drilling, sampling and analytical data (Sen and Duffy, 2005). In 2007, GeoSciML was adopted by the OneGeology initiative to underpin and improve the accessibility of global, regional and national geological map data (Jackson and Wyborn, 2008).

4.5. The Oceans Best Practice System and IODP

The Ocean Best Practices System (OBPS, www.oceanbestpractices.org), is an initiative of the global Intergovernmental Oceanographic Commission (IOC) of UNESCO,

supported by the International Oceanographic Data and Information Exchange (IODE) and the Global Oceans Observing System (GOOS). The OBPS site supports technological solutions and community approaches to ensure FAIR methods and associated data and to facilitate the development, documentation and sharing of ocean best practices. As of 1 March 2023, the OBPS site contains 1787 best practice documents from 52 institutions/organisations: as new documents are submitted, they are reviewed and endorsed by expert teams (Przeslawski et al., 2022). OBPS further runs an ambassador programme to promote equitable access to ocean best practices across communities, disciplines, and regions.

Each institution/organisation can submit their best practice documents including quality documents specific to their data acquisition programmes. The Australian Integrated Marine Observing System (IMOS), for example, operates a wide range of observing equipment throughout Australia's coastal and open oceans and makes all of its data openly and freely accessible. Documents related to the quality of their datasets, including quality specifications, quality evaluation, execution and dissemination are published by IMOS on the international OBPS site (Ruth and Atkins, 2022, <https://repository.oceanbestpractices.org/handle/11329/556>). Publication of best practice documents in a single site from so many organisations leads to convergence and ultimately globalisation of best practices, meaning that a practice can be accessible and usable in multiple regions. At the same time, best practices can be adapted to match regional infrastructure capabilities (Przeslawski et al., 2022).

The International Oceans Discovery Program (IODP, the successor of the Ocean Drilling Program, ODP; <https://www.iodp.org/>) further requires that samples collected on their cruises are archived in one of three recommended repositories. Access to samples is open and transparent to scientists, educators, museums and outreach officers, but regulated by strict policies that ensure their appropriate use and specify the reporting of any research outcomes derived from these samples (<https://www.iodp.org/top-resources/program->

documents/policies-and-guidelines/519-iodp-sample-data-and-obligations-policy-
implementation-guidelines-may-2018-for-expeditions-starting-october-2018-and-
later/file). These outcomes are made available through the integrated data and publi-
cation portal SEDIS (Scientific Earth Drilling Information System; <http://sedis.iodp.org/>).

Core funding for OBPS is provided jointly by co-sponsors IODE and GOOS (both in
turn funded through the International Oceanographic Commission, IOC). Any technologi-
cal developments and implementation of the OBPS objectives and community recommen-
dations has to be supplemented by external project funding, such as IMOS. The work of
OBPS is overseen by a UNESCO-funded project manager and 24 volunteer steering group
members.

4.6. What can be learned from these initiatives?

The examples from crystallography, chemistry, seismology, geology and oceanography
demonstrate that it is indeed possible to unite community efforts and together define,
implement and enforce best practices and standards for data reporting at an international
level. The geochemical and cosmochemical communities can benefit by implementing many
common threads outlined in the above initiatives, including:

1. Securing endorsements from recognised, authoritative groups that are connected to
leading International Science Unions/organisations; in some cases, these groups also
provide limited funding;
2. Establishing expert committees for developing data standards and regularly updating
these standards as additional requirements emerge;
3. Publishing community-agreed, time-stamped standards and vocabularies online in
both human and machine-readable formats in governed, sustainable repositories;
4. Connecting with funding agencies to adopt commonly defined standards and enforce
research data management plans and data submissions;

5. Connecting with publishers and editors to enforce compliance with data standards within publications;
6. Developing and implementing tools that validate data standards compliance;
7. Enforcing data submission to domain repositories that work with publishers to implement standards and ensure long-term preservation and increased discoverability of data;
8. Adoption of standard data and file formats by instrument manufacturers;
9. Developing education and outreach programs to teach data management and disseminate existing standards and best practices for data users and contributors.

5. The Path Forward: OneGeochemistry

During the workshop at Goldschmidt 2022, organisers and participants discussed possible solutions to the aforementioned challenges and towards the goal of a standardised network of geochemical data resources. The options promising the highest short-term impact are: official endorsement of the OneGeochemistry initiative; establishment of expert committees to collect and define best practices for each data type; and a broad education and outreach programme that highlights the benefits of community engagement in this issue. Each of these strategies is discussed in detail below.

5.1. Endorsement

Standards and data management should be developed bottom-up but need to be enforced top-down. As a consequence, OneGeochemistry is pursuing endorsement from (i) societies, (ii) publishers, (iii) funders and (iv) instrument manufacturers to gain authority for the initiative and thus increase community participation.

5.1.1. Societies and Unions

The heterogeneity of geochemical data and the multiple purposes that geochemistry can be used for, has resulted in geochemistry being a part of at least four International Science Council (ISC) Science Unions and tens, if not hundreds, of geochemical associations,

societies, and commissions at both international and national level. The four main unions that are relevant to geochemical and cosmochemical data include the International Union of Geological Sciences (IUGS), International Union of Geodesy and Geophysics (IUGG), International Union of Crystallography (IUCr) and the International Union of Pure and Applied Chemistry (IUPAC).

As of December 2022, the OneGeochemistry initiative is acting as the OneGeochemistry CODATA Working Group under the International Science Council to bring together all the disparate initiatives that are happening in geochemistry across Scientific Unions, Associations, Societies and Commissions (<https://codata.org/initiatives/decadal-programme2/worldfair/onegeochemistry-wg/>). Over the next two years, this Working Group will be utilised to recruit a larger membership base to the initiative that will then be able to vote on a long-term governance structure for OneGeochemistry. The OneGeochemistry interim board has so far secured endorsement from the following six international geochemical societies and associations: the Geochemical Society, the European Association of Geochemistry, the Association of Applied Geochemists, the International Association of Geochemistry, the Meteoritical Society and the IUGS commission on Global Geochemical Baselines. A final decision is pending from the International Association of Geoanalysts and the International Association of Geochemists. These developments lend authority to OneGeochemistry as the trusted international initiative tasked with bringing together the community and coordinate global efforts in geochemical data standardisation. Society endorsement will further help disseminate the goals and activities of OneGeochemistry to a broad membership throughout the geochemical sub-disciplines, and increase participation in the initiative. Additional national and/or sub-disciplinary societies will be contacted in the future and the OneGeochemistry board invites suggestions and recommendations from the community.

5.1.2. *Publishers*

OneGeochemistry will continue the discussion with journal publishers and editors to raise awareness for the need for data standards in geochemistry to be enforced. The Commitment Statement developed by the Coalition for Publishing Data in the Earth and Space Sciences (COPDESS; <https://copdess.org/enabling-fair-data-project/commitment-statement-in-the-earth-space-and-environmental-sciences/>) has united many of the repositories, publishers, societies, institutions and infrastructure in an agreement to uphold minimum standards. OneGeochemistry will build upon this commitment and, through town halls and other meetings at international conferences, will work towards establishing domain repositories as trusted data publishers that collaborate with journals and publishers to ensure that data submitted to a journal comply with agreed community standards and the FAIR principles.

5.1.3. *Funders*

As a community we need to communicate with the national and regional funding agencies to alert them to our requirements for data management. Many funders have FAIR data policies but most do not yet enforce them or check compliance. In addition, funders play an important role in guiding the academic credit system. For example, the German Research Foundation (DFG) recently changed their rules to recognise article preprints, data sets or software packages as research outcomes, which is an important and positive signal to the scientific community (https://www.dfg.de/en/research_funding/announcements_proposals/2022/info_wissenschaft_22_61/index.html).

5.1.4. *Instrument Manufacturers*

At Goldschmidt 2022, members of the OneGeochemistry interim board connected with some of the geochemical instrument manufacturers, who were very supportive of the initiative and committed to implementing community-agreed data, metadata and formatting standards once they were developed and accepted. As shown by the example from the seis-

mological community, support and adoption by instrument manufacturers of community-agreed data standards, aided by common file formats, is crucial to their widespread implementation within laboratories. The increasing adoption of electronic laboratory notebooks, for example, could be exploited to implement data standards and provide a direct data pipeline into certified domain repositories.

5.2. *Expert Committees*

Multiple best practices and recommendations for specific data types, analytical techniques or sub-disciplines have already been defined and are variably adhered to across the globe. A growing number of publications aim to establish agreement on minimum variables and vocabularies for various geochemical data types Deines et al. (2003); Demetriades et al. (2020, 2022); Boone et al. (2022); Flowers et al. (2022); Brantley et al. (2021); Abbott et al. (2022); Horstwood et al. (2016); Dutton et al. (2017); Walker et al. (2008); Courtney Mustaphi et al. (2019); Schaen et al. (2020); Khider et al. (2019); Damerow et al. (2021). Effective development of scientific standards requires a participatory framework with a need for ongoing, open dialogue within and across research communities (Yarmey and Baker, 2013). The larger the size of the community that agrees and commits to a particular standard, the larger the community that can share and reuse data, particularly in machine-to-machine environments. Hence, to enable global data exchange, we need to harmonise and curate these existing standards through a number of expert committees that are endorsed and/or recognised by authoritative, international geochemical societies and unions. The task of these expert committees would be to compile and further develop standards for each distinct analytical technique or related groups of analytical methods. A committee would be made up of experts within a specific method that are representative of the diversity of users for each data type, including geographical regions, institutions and career levels.

OneGeochemistry's role will be to facilitate and support these expert committees, as well as to disseminate best practice recommendations and invite feedback from the wider

community. In addition, OneGeochemistry will set up a technical committee that converts existing standards into machine-readable format. Overall, the focus of the OneGeochemistry initiative is to coordinate global efforts in geochemical data standardisation, facilitate communication amongst distributed groups and thus minimise duplication and redundancy. In a first step, OneGeochemistry will work with the wider community to compile existing standards, determine which additional data types require standards/vocabularies and which analytical methods are currently in use or have been used in the past for each data type. The role of the expert committees would then be to:

1. Compile lists of existing standards or best practices (including data models and vocabularies) and ensure they are in the public domain, accessible online in a repository or vocabulary service, such as OBPS and RVA, respectively;
2. Review neighbouring fields and disciplines that have already defined data standards to ensure interoperability (e.g. IUPAC terminologies, government agencies or industry standards);
3. Provide governance to existing standards and harmonise where possible;
4. Monitor and update each agreed upon standard as needed;
5. Develop new data standards where required.

The technical committee led by OneGeochemistry would then work with the expert committees to digitise these standards and make them FAIR. A timeframe of two years per thematic expert committee is envisaged, culminating in a formal publication of the recommended standard and its presentation to the community at one of the annual workshops facilitated by OneGeochemistry.

All community-agreed standards are to be published through the ‘Brown Book’, part of the IUPAC Colour Books Series described in Section 4 above which has been offered to OneGeochemistry. With this Brown Book the geochemistry community will be able to publish any nomenclature, terminology or standards that are not already covered in the geochemistry literature. This resource will be invaluable not only in documenting nomenclatures

defined by the geochemical expert committees, but also in ensuring that relevant, existing digital chemical standards are leveraged wherever possible (e.g., the Machine-Accessible Periodic Table).

A successful example of an existing expert committee in geochemistry is the Tephra Community that has developed data submission templates for the EarthChem Library (Wallace et al., 2022). EarthChem has further recently started a working group to develop a method directory. Whilst we acknowledge the risk that this modular approach might further divide the community, we propose that it is the most viable solution to: 1) Involve the community in the process of developing data standards; 2) Provide well-defined, feasible work packages with clear credit/reward/outcome that will motivate community-participation; and 3) Give authority to the standards developed to ensure they are accepted by the wider community. To contribute to or join the OneGeochemistry initiative please visit www.onegeochemistry.org for more information and contact onegeochemistry@codata.org.

5.3. Incentives, Education & Outreach

We recognise that a critical component for the success of OneGeochemistry is increasing outreach and dissemination while establishing appropriate incentives that invite more community members to join. An unexpected outcome of the Goldschmidt 2022 workshop was the observation how poorly known the existing data systems are, especially among early career researchers. Through the OneGeochemistry initiative we hope to achieve greater community engagement via (i) passive advertising of data efforts within research presentations and publications; (ii) virtual campaigns and the open sharing of resources; and (iii) active training through workshops and data mentoring programmes. Whilst this active training can be primarily facilitated by members of the OneGeochemistry board, passive advertising and sharing of resources rely on community participation. For example, passive advertising may include the proper attribution of data systems in publications, following citation guidelines and templates provided by the systems, or the addition of data sys-

tem logos to presentation materials (e.g. conference slides, posters, graphical abstracts). Virtual campaigns include a broad social media presence (e.g. on Twitter, LinkedIn), blog posts, webinars and a dedicated YouTube channel to disseminate tutorials and teach data management skills. All of these activities would greatly benefit from the participation of a broad group of active community members and ‘*OneGeochemistry ambassadors*’ could drive these initiatives. Ambassadors are envisaged as early to mid-career, cutting-edge researchers that promote good data management following current best practices and standards. Assisted by the OneGeochemistry board members, ambassadors will spread awareness in the communities of the importance of data management in geo- and cosmochemistry, the existing landscape of data systems, and inspire new and future generations to contribute. In parallel, OneGeochemistry and its participating data systems would continue to host workshops at scientific conferences, organise data hackathons, contribute to the Data Help Desks coordinated by ESIP at major Earth Science conferences such as the AGU Fall Meeting, the EGU General Assembly and the Geological Society of America meeting (<https://www.esipfed.org/data-help-desk>) and hold Data FAIR workshops (<https://data.agu.org/datafair/>). In addition, data management could be integrated into mentoring schemes at these conferences and inter-institution and international data mentoring programs could focus on available resources in the communities.

While communicating and advertising OneGeochemistry, we must always be aware of motivations and incentives (or disincentives) to contribute to standard development, data publication and global databases for each stakeholder. Options to increase community uptake of data sharing practices have been discussed at length in other communities and center around a balance between the perceived cost *versus* benefit of data sharing (e.g. Kim and Stanton, 2012; Kidwell et al., 2016). Yet, the precise incentives will differ widely between different groups in the community (Fig. 4). For OneGeochemistry, the focus is on engaging:

- **Publishers and editors** who ensure peer review, storage and release of datasets in

certified domain repositories prior to publication.

- **Funding agencies** who require compliance with certified standards, and provide necessary funds for data curation and staff.
- **Data repositories** who are key to storing, curating and making geoanalytical data FAIR.
- **Government surveys/agencies** who have a long history of generating and archiving publicly funded research data as well as industry data.
- **Professional societies/science unions/associations** who can both endorse and help to promote the standards/best practices.
- **Instrument manufacturers** who can ensure any data generated with their instruments and output by their software are compliant with standards.
- **Laboratory managers** and other geoanalytical data producers to ensure consistency and quality of geochemical data at the point of generation.
- **Researchers** who generate, (re)use and publish geochemical data.

For *researchers*, the main incentive for engaging in good data management practices is credit received towards their scientific track record. As more funding, recruitment and promotion bodies start considering more than journal publications as a measurable research output, data publications in domain repositories will gain importance. OneGeochemistry and/or its member data systems will further strive to support researchers through acknowledging the number and quality of individual contributions on their websites or, as is common practice with software, through regular version releases. Tracking of citations to data publications independently of a related research paper will provide an additional measure of impact of specific research outputs. Tracking data citations is also a convenient way for *funders*, *institutions* and *laboratories* to measure their impact. Both ‘data



Figure 4: The place of OneGeochemistry within the broader research data landscape (adapted from OECD, 2017). Each group of stakeholders has different needs and motives for contributing to or enforcing FAIR data practices. Blue circles symbolise the role of OneGeochemistry in coordinating expert committees and facilitating education and ambassadorship.

reliquaries' and the new 'smart citation' frameworks, such as scite_, are promising developments that will aid this cause. For *instrument manufacturers*, clear guidance for data and file formats through community-agreed standards would significantly reduce the resources spent on developing custom data formats for each analytical instrument. At the same time, proprietary file formats need not be forfeited as long as final data outputs follow the community-agreed standards.

Industry, such as mining or environmental companies, have been omitted from the

list since this initiative is born out of the academic (and governmental research) domain. However, we acknowledge that these companies produce large data volumes and we would welcome future contact and participation with industry representatives. Some countries, such as Australia, already require that all industry data be made available to local geological surveys after a certain time period— providing an incentive for companies to comply with common data standards to facilitate data sharing, whilst still ensuring a competitive advantage through time-limited, confidential agreements.

6. Conclusions

There is an urgent need in the geochemistry and cosmochemistry communities to define data-type specific best practices and standards for reporting geoanalytical data. Only once these best practices exist, are implemented in research workflows and are consistently followed will geoanalytical data become easy to find, trust and reuse for education or further data-driven research that is increasingly employed to tackle the next big, data intensive and complex scientific questions. We propose that the international OneGeochemistry initiative enacts this change, driven and supported by the community, through facilitating a global, online network of machine-readable data that is persistent, interoperable and reusable, and above all minimises duplication. Once the community has adopted and fully integrated a culture of standardised data and metadata reporting practices, such a framework will also ensure reliable attribution of those who collected, analysed, curated and made accessible any geochemical and cosmochemical data. Endorsement by societies, publishers and funders will give the OneGeochemistry initiative authority to establish expert committees that develop and promote best practices and standards for specific data types. Community engagement and participation at all stages of the process will be pursued through active outreach and dissemination.

Acknowledgements

This manuscript is the result of a workshop held at the Goldschmidt 2022 conference hosted by the Geochemical Society. We thank Jerry Carter (IRIS) and Rob Casey (EarthScope) for helpful comments on the history of the development of data standards in seismology. Jay Pearlman, Rachel Przeslawski, Pauline Simpson provided valuable background information on OBPS. Richard Hartshorn and Leah McEwen provided detailed feedback on the practices in chemistry and crystallography. Michael Badawi, Jieun Kim and Nicolas Randazzo are thanked for their contributions to the Goldschmidt 2022 workshop. We thank Olivier Pourret for helpful comments on the preprint and are grateful to Tao Wen, Penny Wieser and Jamie Farquharson for their thoughtful and constructive reviews and to Jeff Catalano for expert editorial handling of the manuscript. MK is supported by the German Research Foundation (DFG grant 437919684). KL and LP acknowledge funding from US NSF Award Number 2148939 and NASA Grant Number 80NSSC19K1102. BW, AMP, and the AuScope Geochemistry Network are supported by AuScope and the National Collaborative Research Infrastructure Strategy (NCRIS). AMP is supported through AuScope which is a beneficiary in the WorldFAIR project, coordinated by CODATA, and funded by the European Union's Horizon Europe Framework Programme (grant agreement 101058393). SL was supported by the NSF (EAR grant 1946346). NDB acknowledges funding from the NERC Centre for the Observation and Modelling of Earthquakes, Volcanoes and Tectonics (COMET), the Bill & Melinda Gates Foundation (Grant Number OPP1144), and the Gates Cambridge Trust. HB was supported by DFG CRC TRR 170 (Project-ID 263649064). MB and HK are supported by Science Foundation Ireland (SFI) grant 13/RC/2092 and co-funded by iCRAG industry partners. HK is further supported by SFI grant 16/RP/3849. KD thanks the support by the ETH Zurich Postdoctoral Fellowship 20-1 FEL-24. DCH is supported through the NFDI4Earth funded by the DFG (project number 460036893). BK was supported by NSF grant OIA1545903. AWL was funded by the Geological Society of Australia - Victorian

Division, the German Research Exchange Service (DAAD Grant No. 57507869) and the Australian Government Research Training Program (Allocation No. 2018177). WS was funded by DFG grant KE 2395/3-1 (project number 447528294). MKT was supported by Geocenter Danmark (grant nr. GC4-2019). This is contribution 1760 from the ARC Centre of Excellence for Core to Crust Fluid Systems and 1529 in the GEMOC Key Centre.

Appendix A. Supplementary Material

The Supplementary Material contains additional information on the Goldschmidt 2022 workshop “Earth Science meets Data Science: what are our needs for geochemical data, services and analytical capabilities in the 21st century?”, including the workshop programme, details on the participating data systems and the complete list of contributors.

References

- Abbott, P., Bonadonna, C., Bursik, M., Cashman, K., Davies, S., Jensen, B., Kuehn, S., Kurbatov, A., Lane, C., Plunkett, G., Smith, V., Thomlinson, E., Thordarsson, T., Walker, J.D., Wallace, K., 2022. Community established best practice recommendations for tephra studies-from collection through analysis. doi:10.5281/ZENODO.3866266.
- Agarwal, D.A., Damerow, J., Varadharajan, C., Christianson, D.S., Pastorello, G.Z., Cheah, Y.W., Ramakrishnan, L., 2021. Balancing the needs of consumers and producers for scientific data collections. *Ecological Informatics* 62, 101251. doi:10.1016/J.ECOINF.2021.101251.
- Alexakis, D.E., 2021. Linking DPSIR model and water quality indices to achieve sustainable development goals in groundwater resources. *Hydrology* 8, 90. doi:10.3390/hydrology8020090.
- Ball, C.A., Sherlock, G., Brazma, A., 2004. Funding high-throughput data sharing. *Nature Biotechnology* 2004 22:9 22, 1179–1183. doi:10.1038/nbt0904-1179.

- 886 Bergerhoff, G., Brown, I., 1987. International Union of Crystallography, Chester.
- 887 Berman, H., Henrick, K., Nakamura, H., 2003. Announcing the worldwide protein data
888 bank. *Nature Structural & Molecular Biology* 10, 980–980. doi:10.1038/nsb1203-980.
- 889 Boone, S.C., Dalton, H., Prent, A., Kohlmann, F., Theile, M., Gréau, Y., Florin, G.,
890 Noble, W., Hodgekiss, S.A., Ware, B., Phillips, D., Kohn, B., O'Reilly, S., Gleadow, A.,
891 McInnes, B., Rawling, T., 2022. AusGeochem: An open platform for geochemical data
892 preservation, dissemination and synthesis. *Geostandards and Geoanalytical Research*
893 46, 245–259. doi:10.1111/ggr.12419.
- 894 Brantley, S.L., Wen, T., Agarwal, D.A., Catalano, J.G., Schroeder, P.A., Lehnert, K.,
895 Varadharajan, C., Pett-Ridge, J., Engle, M., Castronova, A.M., Hooper, R.P., Ma, X.,
896 Jin, L., McHenry, K., Aronson, E., Shaughnessy, A.R., Derry, L.A., Richardson, J.,
897 Bales, J., Pierce, E.M., 2021. The future low-temperature geochemical data-scape as
898 envisioned by the U.S. geochemical community. *Computers & Geosciences* 157, 104933.
899 doi:10.1016/j.cageo.2021.104933.
- 900 Bruno, I., Gražulis, S., Helliwell, J.R., Kabekkodu, S.N., McMahon, B., Westbrook, J.,
901 2017. Crystallography and databases. *Data Science Journal* 16. doi:10.5334/dsj-
902 2017-038.
- 903 Bundschuh, J., Maity, J.P., Mushtaq, S., Vithanage, M., Seneweera, S., Schneider, J., Bhat-
904 tacharya, P., Khan, N.I., Hamawand, I., Guilherme, L.R., Reardon-Smith, K., Parvez,
905 F., Morales-Simfors, N., Ghaze, S., Pudmenzky, C., Kouadio, L., Chen, C.Y., 2017.
906 Medical geology in the framework of the sustainable development goals. *Science of The*
907 *Total Environment* 581-582, 87–104. doi:10.1016/j.scitotenv.2016.11.208.
- 908 Carbotte, S., Lehnert, K., 2007. WORKSHOP REPORT | building a global data network
909 for studies of earth processes at the worlds plate boundaries. *Oceanography* 20, 124–125.
910 doi:10.5670/oceanog.2007.38.

- 911 Carroll, S.R., Garba, I., Figueroa-Rodríguez, O.L., Holbrook, J., Lovett, R., Materechera,
912 S., Parsons, M., Raseroka, K., Rodriguez-Lonebear, D., Rowe, R., Sara, R., Walker, J.D.,
913 Anderson, J., Hudson, M., 2020. The CARE principles for indigenous data governance.
914 Data Science Journal 19. doi:10.5334/dsj-2020-043.
- 915 Chamberlain, K.J., Lehnert, K.A., McIntosh, I.M., Morgan, D.J., Wörner, G., 2021. Time
916 to change the data culture in geochemistry. Nature Reviews Earth & Environment 2,
917 737–739. doi:10.1038/s43017-021-00237-w.
- 918 Courtney Mustaphi, C.J., Brahney, J., Aquino-López, M.A., Goring, S., Orton, K.,
919 Noronha, A., Czaplewski, J., Asena, Q., Paton, S., Brushworth, J.P., 2019. Guidelines for
920 reporting and archiving 210Pb sediment chronologies to improve fidelity and extend data
921 lifecycle. Quaternary Geochronology 52, 77–87. doi:10.1016/j.quageo.2019.04.003.
- 922 Cox, S.J.D., Gonzalez-Beltran, A.N., Magagna, B., Marinescu, M.C., 2021. Ten simple
923 rules for making a vocabulary FAIR. PLOS Computational Biology 17, e1009041. doi:10.
924 1371/JOURNAL.PCBI.1009041.
- 925 Damerow, J.E., Varadharajan, C., Boye, K., Brodie, E.L., Burrus, M., Chadwick, K.D.,
926 Crystal-Ornelas, R., Elbashandy, H., Alves, R.J.E., Ely, K.S., Goldman, A.E., Haber-
927 man, T., Hendrix, V., Kakalia, Z., Kemner, K.M., Kersting, A.B., Merino, N., OBrien,
928 F., Perzan, Z., Robles, E., Sorensen, P., Stegen, J.C., Walls, R.L., Weisenhorn, P.,
929 Zavarin, M., Agarwal, D., 2021. Sample identifiers and metadata to support data man-
930 agement and reuse in multidisciplinary ecosystem sciences. Data Science Journal 20, 11.
931 doi:10.5334/dsj-2021-011.
- 932 Deines, P., Goldstein, S.L., Oelkers, E.H., Rudnick, R.L., Walter, L.M., 2003. Standards
933 for publication of isotope ratio and chemical data in chemical geology. Chemical Geology
934 202, 1–4. doi:10.1016/j.chemgeo.2003.08.003.
- 935 Demetriades, A., Huimin, D., Kai, L., Savin, I., Birke, M., Johnson, C.C., Argyraki, A.,

2020. International union of geological sciences manual of standard geochemical methods
for the global black soil project. doi:10.5281/ZENODO.7267967.
- Demetriades, A., Johnson, C.C., Smith, D.B., Ladenberger, A., Sanjuan, P.A., Argyraki,
A., Stouraiti, C., de Caritat, P., Knights, K.V., Rincón, G.P., Simubali, G.N., 2022.
International Union of Geological Sciences Manual of Standard Methods for Establishing
the Global Geochemical Reference Network. IUGS Commission on Global Geochemical
Baselines, Athens, Hellenic Republic. doi:10.5281/ZENODO.7307696.
- Duke, R., Bhat, V., Risko, C., 2022. Data storage architectures to accelerate chemical
discovery: data accessibility for individual laboratories and the community. *Chemical
Science* 13, 13646–13656. doi:10.1039/D2SC05142G.
- Dutton, A., Rubin, K., McLean, N., Bowring, J., Bard, E., Edwards, R., Henderson, G.,
Reid, M., Richards, D., Sims, K., Walker, J., Yokoyama, Y., 2017. Data reporting
standards for publication of U-series data for geochronology and timescale assessment
in the earth sciences. *Quaternary Geochronology* 39, 142–149. doi:10.1016/j.quageo.
2017.03.001.
- Dziewonski, A.M., 1994. The FDSN: history and objectives. *Annals of Geophysics* 37.
doi:10.4401/ag-4191.
- European Commission, 2018. Cost-benefit analysis for FAIR research data: cost of not
having FAIR research data. European Commission, Directorate General for Research
and Innovation and PwC EU Services. doi:10.2777/02999.
- Evans, P., Strollo, A., Clark, A., Ahern, T., Newman, R., Clinton, J., Pedersen, H.,
Pequegnat, C., 2015. Why seismic networks need digital object identifiers. *Eos* 96.
doi:10.1029/2015eo036971.
- Fairbridge, R.W., 1998. History of geochemistry, in: *Encyclopedia of Earth Science*. Kluwer
Academic Publishers, pp. 315–322. doi:10.1007/1-4020-4496-8_156.

- Flowers, R., Zeitler, P., Danišik, M., Reiners, P., Gautheron, C., Ketcham, R., Metcalf, J., Stockli, D., Enkelmann, E., Brown, R., 2022. (U-Th)/He chronology: Part 1. data, uncertainty, and reporting. *GSA Bulletin* 135, 104–136. doi:10.1130/b36266.1.
- Geochemical Society, 2007. Geochemical society policy on geochemical databases. URL: <https://www.geochemsoc.org/about/positionstatements/datapolicy>.
- Gill, J.C., 2017. Geology and the sustainable development goals. *Episodes* 40, 70–76. doi:10.18814/epiiugs/2017/v40i1/017010.
- Goldstein, S., Lehnert, K., Hofmann, A., 2014. Requirements for the publication of geochemical data. doi:10.1594/IEDA/100426.
- Groom, C.R., Bruno, I.J., Lightfoot, M.P., Ward, S.C., 2016. The Cambridge Structural Database. *Acta Crystallographica Section B Structural Science, Crystal Engineering and Materials* 72, 171–179. doi:10.1107/s2052520616003954.
- Hall, S.R., Allen, F.H., Brown, I.D., 1991. The crystallographic information file (CIF): a new standard archive file for crystallography. *Acta Crystallographica Section A Foundations of Crystallography* 47, 655–685. doi:10.1107/s010876739101067x.
- Hall, S.R., Cook, A.P.F., 1995. STAR dictionary definition language: Initial specification. *Journal of Chemical Information and Computer Sciences* 35, 819–825. doi:10.1021/ci00027a005.
- Hall, S.R., McMahon, B., 2016. The implementation and evolution of STAR/CIF ontologies: Interoperability and preservation of structured data. *Data Science Journal* 15. doi:10.5334/dsj-2016-003.
- He, Y., Zhou, Y., Wen, T., Zhang, S., Huang, F., Zou, X., Ma, X., Zhu, Y., 2022. A review of machine learning in geochemistry and cosmochemistry: Method improvements

- 984 and applications. *Applied Geochemistry* 140, 105273. doi:10.1016/J.APGEOCHEM.2022.
985 105273.
- 986 Heidorn, P.B., 2008. Shedding light on the dark data in the long tail of science. *Library*
987 *Trends* 57, 280–299. doi:10.1353/lib.0.0036.
- 988 Heller, S., McNaught, A., Stein, S., Tchekhovskoi, D., Pletnev, I., 2013. InChI - the
989 worldwide chemical structure identifier standard. *Journal of Cheminformatics* 5. doi:10.
990 1186/1758-2946-5-7.
- 991 Horstwood, M.S.A., Košler, J., Gehrels, G., Jackson, S.E., McLean, N.M., Paton, C.,
992 Pearson, N.J., Sircombe, K., Sylvester, P., Vermeesch, P., Bowring, J.F., Condon,
993 D.J., Schoene, B., 2016. Community-derived standards for LA-ICP-MS U-(Th)Pb
994 geochronology – uncertainty propagation, age interpretation and data reporting. *Geo-*
995 *standards and Geoanalytical Research* 40, 311–332. doi:10.1111/j.1751-908x.2016.
996 00379.x.
- 997 Jackson, I., Wyborn, L., 2008. International viewpoint and news. *Environmental Geology*
998 53, 1377–1380. doi:10.1007/s00254-007-1085-z.
- 999 Khider, D., Emile-Geay, J., McKay, N.P., Gil, Y., Garijo, D., Ratnakar, V., Alonso-Garcia,
1000 M., Bertrand, S., Bothe, O., Brewer, P., Bunn, A., Chevalier, M., Comas-Bru, L., Csank,
1001 A., Dassié, E., DeLong, K., Felis, T., Francus, P., Frappier, A., Gray, W., Goring,
1002 S., Jonkers, L., Kahle, M., Kaufman, D., Kehrwald, N.M., Martrat, B., McGregor,
1003 H., Richey, J., Schmittner, A., Scroxton, N., Sutherland, E., Thirumalai, K., Allen,
1004 K., Arnaud, F., Axford, Y., Barrows, T., Bazin, L., Birch, S.E.P., Bradley, E., Bregy,
1005 J., Capron, E., Cartapanis, O., Chiang, H.W., Cobb, K.M., Debret, M., Dommain,
1006 R., Du, J., Dyez, K., Emerick, S., Erb, M.P., Falster, G., Finsinger, W., Fortier, D.,
1007 Gauthier, N., George, S., Grimm, E., Hertzberg, J., Hibbert, F., Hillman, A., Hobbs,
1008 W., Huber, M., Hughes, A.L.C., Jaccard, S., Ruan, J., Kienast, M., Konecky, B., Roux,

G.L., Lyubchich, V., Novello, V.F., Olaka, L., Partin, J.W., Pearce, C., Phipps, S.J., Pignol, C., Piotrowska, N., Poli, M.S., Prokopenko, A., Schwanck, F., Stepanek, C., Swann, G.E.A., Telford, R., Thomas, E., Thomas, Z., Truebe, S., Gunten, L., Waite, A., Weitzel, N., Wilhelm, B., Williams, J., Williams, J.J., Winstrup, M., Zhao, N., Zhou, Y., 2019. PaCTS 1.0: A crowdsourced reporting standard for paleoclimate data. *Paleoceanography and Paleoclimatology* 34, 1570–1596. doi:10.1029/2019pa003632.

Kidwell, M.C., Lazarević, L.B., Baranski, E., Hardwicke, T.E., Piechowski, S., Falkenberg, L.S., Kennett, C., Slowik, A., Sonnleitner, C., Hess-Holden, C., Errington, T.M., Fiedler, S., Nosek, B.A., 2016. Badges to Acknowledge Open Practices: A Simple, Low-Cost, Effective Method for Increasing Transparency. *PLOS Biology* 14, e1002456. doi:10.1371/JOURNAL.PBIO.1002456.

Kim, Y., Stanton, J.M., 2012. Institutional and individual influences on scientists' data sharing practices. *The Journal of Computational Science Education* 3, 47–56. doi:https://doi.org/10.22369/issn.2153-4136/3/1/6.

Kroon-Batenburg, L.M.J., Helliwell, J.R., Hester, J.R., 2022. *IUCrData* launches raw data letters. *IUCrData* 7. doi:10.1107/s2414314622008215.

Lehnert, K., Wyborn, L., Bennett, V.C., Hezel, D., McInnes, B.I.A., Plank, T., Rubin, K., 2021. Onegeochemistry: Towards an interoperable global network of FAIR geochemical data doi:10.5281/ZENODO.5767950.

Lin, D., Crabtree, J., Dillo, I., Downs, R.R., Edmunds, R., Giaretta, D., Giusti, M.D., L'Hours, H., Hugo, W., Jenkyns, R., Khodiyar, V., Martone, M.E., Mokrane, M., Navale, V., Petters, J., Sierman, B., Sokolova, D.V., Stockhause, M., Westbrook, J., 2020. The TRUST principles for digital repositories. *Scientific Data* 7. doi:10.1038/s41597-020-0486-7.

Morrison, S.M., Liu, C., Eleish, A., Prabhu, A., Li, C., Ralph, J., Downs, R.T., Golden,

J.J., Fox, P., Hummer, D.R., Meyer, M.B., Hazen, R.M., 2017. Network analysis of mineralogical systems. *American Mineralogist* 102, 1588–1596. doi:10.2138/AM-2017-6104CCBYNCND.

OECD, 2017. Co-ordination and support of international research data networks. doi:10.1787/e92fa89e-en.

Peng, G., Lacagnina, C., Downs, R.R., Ganske, A., Ramapriyan, H.K., Ivánová, I., Wyborn, L., Jones, D., Bastin, L., Lin Shie, C., Moroni, D.F., 2022. Global community guidelines for documenting, sharing, and reusing quality information of individual digital datasets. *Data Science Journal* 21, 8. doi:10.5334/dsj-2022-008.

Piwowar, H.A., Day, R.S., Fridsma, D.B., 2007. Sharing detailed research data is associated with increased citation rate. *PLoS ONE* 2. doi:10.1371/journal.pone.0000308.

Pourret, O., Irawan, D.E., 2022. Open Access in Geochemistry from Preprints to Data Sharing: Past, Present, and Future. *Publications* 10, 3. doi:10.3390/PUBLICATIONS10010003.

Przeslawski, R., Pearlman, J., Karstensen, J., 2022. Dataset quality information in australia's integrated marine observing system, in: *SciDataCon 2022*. URL: <https://www.scidatacon.org/IDW-2022/sessions/431/paper/969/>.

Ramdeen, S., Wyborn, L.A.I., Lehnert, K.A., Klump, J., 2022. The role of unique identifiers in tracing the life cycle of a sample and any data derived from it, in: *Goldschmidt2022 abstracts*, Geochemical Society. URL: <https://conf.goldschmidt.info/goldschmidt/2022/meetingapp.cgi/Paper/12644>.

Ruth, P.D.V., Atkins, N., 2022. Dataset quality information in Australia's Integrated Marine Observing System, in: *SciDataCon 2022*. URL: <https://www.scidatacon.org/IDW-2022/sessions/431/paper/1052/>.

- Schaen, A.J., Jicha, B.R., Hodges, K.V., Vermeesch, P., Stelten, M.E., Mercer, C.M., Phillips, D., Rivera, T.A., Jourdan, F., Matchan, E.L., Hemming, S.R., Morgan, L.E., Kelley, S.P., Cassata, W.S., Heizler, M.T., Vasconcelos, P.M., Benowitz, J.A., Koppers, A.A., Mark, D.F., Niespolo, E.M., Sprain, C.J., Hames, W.E., Kuiper, K.F., Turrin, B.D., Renne, P.R., Ross, J., Nomade, S., Guillou, H., Webb, L.E., Cohen, B.A., Calvert, A.T., Joyce, N., Ganerød, M., Wijbrans, J., Ishizuka, O., He, H., Ramirez, A., Pfänder, J.A., Lopez-Martínez, M., Qiu, H., Singer, B.S., 2020. Interpreting and reporting $^{40}\text{Ar}/^{39}\text{Ar}$ geochronologic data. *GSA Bulletin* 133, 461–487. doi:10.1130/b35560.1.
- Science, D., Simons, N., Goodey, G., Hardeman, M., Clare, C., Gonzales, S., Strange, D., Smith, G., Kipnis, D., Iida, K., Miyairi, N., Tshetsha, V., Ramokgola, R., Makhera, P., Barbour, G., 2021. The State of Open Data 2021. Technical Report. doi:10.6084/M9.FIGSHARE.17061347.V1.
- Sen, M., Duffy, T., 2005. GeoSciML: Development of a generic GeoScience markup language. *Computers & Geosciences* 31, 1095–1103. doi:10.1016/j.cageo.2004.12.003.
- Spek, A.L., 2020. *checkCIF* validation ALERTS: what they mean and how to respond. *Acta Crystallographica Section E Crystallographic Communications* 76, 1–11. doi:10.1107/s2056989019016244.
- Stall, S., McEwen, L., Wyborn, L., Hoebelheinrich, N., Bruno, I., 2020. Growing the FAIR community at the intersection of the geosciences and pure and applied chemistry. *Data Intelligence* 2, 139–150. doi:10.1162/dint_a_00036.
- Stall, S., Yarmey, L., Cutcher-Gershenfeld, J., Hanson, B., Lehnert, K., Nosek, B., Parsons, M., Robinson, E., Wyborn, L., 2019. Make scientific data FAIR. *Nature* 201 570:7759–570, 27–29. doi:10.1038/d41586-019-01720-7.
- Stodden, V., Seiler, J., Ma, Z., 2018. An empirical analysis of journal policy effectiveness

for computational reproducibility. *Proceedings of the National Academy of Sciences* 115,
2584–2589. doi:10.1073/PNAS.1708290115.

Stuart, D., Baynes, G., Hrynaskiewicz, I., Allin, K., Penny, D., Lucraft, M., Astell, M.,
2018. Whitepaper: Practical challenges for researchers in data sharing doi:10.6084/m9.
figshare.5975011.v1.

Suárez, G., Van Eck, T., Giardini, D., Ahern, T., Butler, R., Tsuboi, S., Suárez, G., 2008.
The International Federation of Digital Seismograph Networks (FDSN): An Integrated
System of Seismological Observatories. *IEEE SYSTEMS JOURNAL* 2. doi:10.1109/
JSYST.2008.2003294.

Taylor, R., Wood, P.A., 2019. A million crystal structures: The whole is greater than
the sum of its parts. *Chemical Reviews* 119, 9427–9477. doi:10.1021/acs.chemrev.
9b00155.

Tedersoo, L., Küngas, R., Oras, E., Köster, K., Eenmaa, H., Leijen, Ä., Pedaste, M., Raju,
M., Astapova, A., Lukner, H., Kogermann, K., Sepp, T., 2021. Data sharing practices
and data availability upon request differ across scientific disciplines. *Scientific Data* 8,
1–11. doi:10.1038/s41597-021-00981-0.

Vines, T.H., Albert, A.Y., Andrew, R.L., Débarre, F., Bock, D.G., Franklin, M.T., Gilbert,
K.J., Moore, J.S., Renaut, S., Rennison, D.J., 2014. The availability of research data
declines rapidly with article age. *Current Biology* 24, 94–97. doi:10.1016/j.cub.2013.
11.014.

Walker, D.J., Condon, D., Thompson, W., Renne, P., Koppers, A., Hodges, K., Reiners,
P., Stockli, D., Schmitz, M., Bowring, S., Gehrels, G., 2008. Geochron workshop reports
sponsored by EarthChem and EARTHTIME. doi:10.5281/ZENODO.4313859.

Wallace, K.L., Bursik, M.I., Kuehn, S., Kurbatov, A.V., Abbott, P., Bonadonna, C., Cash-
man, K., Davies, S.M., Jensen, B., Lane, C., Plunkett, G., Smith, V.C., Tomlinson, E.,

- 1107 Thordarsson, T., Walker, J.D., 2022. Community established best practice recommen-
 1108 dations for tephra studies—from collection through analysis. *Scientific Data* 9, 1–11.
 1109 doi:10.1038/s41597-022-01515-y.
- 1110 Wieser, P.E., Petrelli, M., Lubbers, J., Wieser, E., Özaydın, S., Kent, A.J., Till, C.B.,
 1111 2022. Thermobar: An open-source Python3 tool for thermobarometry and hygrometry.
 1112 *Volcanica* 5, 349–384. doi:10.30909/VOL.05.02.349384.
- 1113 Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A.,
 1114 Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., Bouwman, J., Brookes,
 1115 A.J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers,
 1116 R., Gonzalez-Beltran, A., Gray, A.J., Groth, P., Goble, C., Grethe, J.S., Heringa, J.,
 1117 't Hoen, P.A., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S.J., Martone, M.E., Mons,
 1118 A., Packer, A.L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone, S.A.,
 1119 Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M.A., Thompson, M., van der
 1120 Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft,
 1121 K., Zhao, J., Mons, B., 2016. The FAIR guiding principles for scientific data management
 1122 and stewardship. *Scientific Data* 3. doi:10.1038/sdata.2016.18.
- 1123 Wyborn, L., Elger, K., Prent, A., Lehnert, K., Bruno, I., Klöcking, M., Klump, J., Profeta,
 1124 L., Quinn, D.P., Ramdeen, S., ter Maat, G., 2021. The OneGeochemistry initiative:
 1125 Mobilising a global network of FAIR geochemical data to support research into the
 1126 grand challenge of an environmentally sustainable future doi:10.5281/ZENODO.5765464.
- 1127 Wyborn, L., Lehnert, K., 2021. OneGeochemistry: Creating a global network of geochem-
 1128 ical data to support the 17 United Nations sustainable development goals, in: Gold-
 1129 schmidt2021 abstracts, European Association of Geochemistry. doi:10.7185/gold2021.
 1130 6562.
- 1131 Wyborn, L.A.I., Ryburn, R.J., 1989. PETCHEM data set : Australia and Antarctica -

- 1132 documentation. Record 1989/019. Geoscience Australia, Canberra. URL: [http://pid.](http://pid.geoscience.gov.au/dataset/ga/14256)
1133 [geoscience.gov.au/dataset/ga/14256](http://pid.geoscience.gov.au/dataset/ga/14256).
- 1134 Yarmey, L., Baker, K.S., 2013. Towards standardization: A participatory framework for
1135 scientific standard-making. *International Journal of Digital Curation* 8, 157–172. doi:10.
1136 2218/ijdc.v8i1.252.
- 1137 Yeston, J.S., 2021. Progress in data and code deposition. URL: [https:](https://blogs.sciencemag.org/editors-blog/2021/07/15/progress-in-data-and-code-deposition/)
1138 [//blogs.sciencemag.org/editors-blog/2021/07/15/progress-in-data-and-](https://blogs.sciencemag.org/editors-blog/2021/07/15/progress-in-data-and-code-deposition/)
1139 [code-deposition/](https://blogs.sciencemag.org/editors-blog/2021/07/15/progress-in-data-and-code-deposition/).

Declaration of interests

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

--