

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/365485244>

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

Preprint · November 2022

DOI: 10.31223/X5H07Q

CITATIONS

0

READS

157

35 authors, including:



Marthe Klöcking

Georg-August-Universität Göttingen

29 PUBLICATIONS 179 CITATIONS

[SEE PROFILE](#)



L.A.I. Wyborn

Australian National University

255 PUBLICATIONS 3,658 CITATIONS

[SEE PROFILE](#)



Kerstin Annette Lehnert

Columbia University

150 PUBLICATIONS 1,958 CITATIONS

[SEE PROFILE](#)



Bryant Ware

Curtin University

34 PUBLICATIONS 138 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



PARSEC: Project Building New Tools for Data Sharing and Reuse through a Transnational Investigation of the Socioeconomic Impacts of Protected Areas [View project](#)



Capacity development through institutional linkage for a green and sustainable arsenic remediation process for safe food and water provision in San Luis Potosí [View project](#)

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

^aGeoscience Centre, Georg-August-Universität, Göttingen, Germany; ^bResearch School of Earth Sciences, The Australian National University, Acton, Australia; ^cLamont-Doherty Earth Observatory, Columbia University, Palisades, United States; ^dJohn de Laeter Centre, Curtin University, Bentley, Australia; ^eLithodat Pty Ltd, Melbourne, Australia; ^fCambridge Crystallographic Data Centre, Cambridge, United Kingdom; ^gDepartment of Geology and Geophysics, The University of Utah, Salt Lake City, United States; ^hSchool of Natural Sciences, Macquarie University, North Ryde, Australia; ⁱDepartment of Earth and Planetary Science, McGill University, Montréal, Canada; ^jDepartment of Earth Sciences, University of Cambridge, Cambridge, United Kingdom; ^kGeological Sciences, Freie Universität Berlin, Berlin, Germany; ^lDepartment of Geology, Trinity College Dublin, Dublin 2, Ireland; ^mDepartment of Energy and Resources Engineering, College of Engineering, Peking University, Beijing, China; ⁿInstitute of Geochemistry and Petrology, Department of Earth Sciences, ETH Zürich, Zürich, Switzerland; ^oGFZ German Research Centre for Geosciences, Potsdam, Germany; ^pSchool of Earth Ocean and Environment, University of South Carolina, Columbia, United States; ^qDepartment of Geology, Aligarh Muslim University, Aligarh, India; ^rInstitut für Geowissenschaften, Goethe-Universität Frankfurt, Frankfurt, Germany; ^sDepartment of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, United States; ^tDepartment of Computer Science, University of Idaho, Moscow, United States; ^uDepartment of Geology and Environmental Earth Science, Miami University, Oxford, United States; ^vInstitut für Mineralogie, Universität Münster, Münster, Germany; ^wFaculty of Geosciences, Utrecht University, Utrecht, Netherlands; ^xIPICYT, Instituto Potosino de Investigación Científica y Tecnológica, División de Geociencias Aplicadas, San Luis Potosí, Mexico; ^yDepartment of Geology & Geophysics, Texas A&M University, College Station, United States; ^zDepartment of Earth Sciences, Metal Isotope Group (MIG), Adelaide, Australia; ^{aa}Department of Mineral Resources and Rocks, Faculty of Earth Sciences, King Abdulaziz University, Jeddah, Saudi Arabia; ^{ab}GeoZentrum Nordbayern, Friedrich-Alexander-Universität Erlangen-Nürnberg, Erlangen, Germany; ^{ac}Department of Earth Sciences, University of Toronto, Toronto, Canada; ^{ad}Carl Zeiss Microscopy Ltd, Cambridgeshire, United Kingdom; ^{ae}Department of Geosciences and Natural Resource Management, University of Copenhagen, Copenhagen, Denmark; ^{af}Northern Territory Geological Survey, Darwin, Australia; ^{ag}School of Geosciences, China University of Petroleum (East China), Qingdao, China

This preprint has been submitted for publication in *Geochimica et Cosmochimica Acta*. Please note, this preprint has not yet been peer-reviewed. The final published version of this paper may, therefore, have slightly different content. If accepted, the final version of this manuscript will be available via the 'Peer-reviewed Publication DOI' link on the right-hand side of the webpage. Please feel free to contact the authors; we welcome feedback. Thank you.

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

^a*Geoscience Centre, Georg-August-Universität, Goldschmidtstr. 1, Göttingen, 37077, Germany*

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

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

^d*John de Laeter Centre, Curtin University, Building 301, Murdoch Ct, Bentley, 6845, WA, Australia*

^e*Lithodat Pty Ltd, Melbourne, Australia*

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

^g*Department of Geology and Geophysics, The University of Utah, Salt Lake City, UT, United States*

^h*CCFS/GEMOC, School of Natural Sciences, Macquarie University, Balaclava Road, North Ryde, 2109, NSW, Australia*

ⁱ*Department of Earth and Planetary Science, McGill University, Montréal, Québec, Canada*

^j*Department of Earth Sciences, University of Cambridge, Cambridge, United Kingdom*

^k*Geological Sciences, Freie Universität Berlin, Malteserstr. 74-100, Berlin, 12249, Germany*

^l*Department of Geology, Trinity College Dublin, Dublin 2, Ireland*

^m*Department of Energy and Resources Engineering, College of Engineering, Peking University, Beijing, China*

ⁿ*Institute of Geochemistry and Petrology, Department of Earth Sciences, ETH Zürich, Clausiusstrasse 25, Zürich, 8092, Switzerland*

^o*GFZ German Research Centre for Geosciences, Telegrafenberg, Potsdam, 14473, Germany*

^p*School of Earth Ocean and Environment, University of South Carolina, Columbia, United States*

^q*Department of Geology, Aligarh Muslim University, Aligarh, 202001, India*

^r*Institut für Geowissenschaften, Goethe-Universität Frankfurt, Altenhöferallee 1, Frankfurt, 60437, Germany*

^s*Department of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, 21218, United States*

^t*Department of Computer Science, University of Idaho, Moscow, 83844, United States*

^u*Department of Geology and Environmental Earth Science, Miami University, Oxford, OH, United States*

^v*Institut für Mineralogie, Universität Münster, Corrensstrasse 24, Münster, 48149, Germany*

^w*Faculty of Geosciences, Utrecht University, Princetonlaan 8a, Utrecht, 3584, CB, Netherlands*

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

^x*IPICYT, 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*

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

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

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

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

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

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

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

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

^{ag}*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 descriptions as well as appropriate information on data quality are becoming crucial to enabling reproducibility and integration of results and fostering confidence for data 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. 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 working groups, 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 has been acquired and used to explore a large 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.

Motivated by a growing need for globally standardised geochemical data, 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 ad-

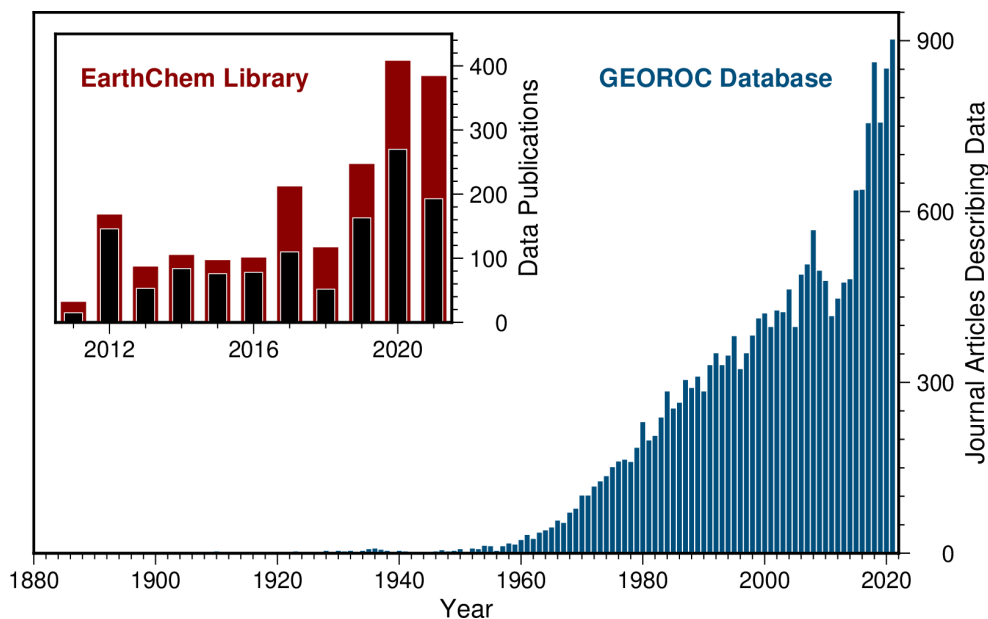


Figure 1: Increase in geochemical data published in journals and repositories since the late 19th century. Blue histogram shows data compiled within the GEOROC database, by publication year of the respective journal articles, as a proxy for the increase in data production within the subdiscipline of igneous geochemistry. The oldest article contributing to the GEOROC compilations was published in 1883. Inset: Number of data submissions to the EarthChem Library (ECL), a domain repository for all subdisciplines of geochemistry. Red = files submitted to the ECL each year; black = datasets published. Note the significant increase in submissions following a change in publisher requirements in 2019.

26 dressing future scientific challenges. More information about the workshop programme,
 27 participating data systems and attendees is available in the Supplementary Material.

28 This paper summarises the workshop outcomes and provides recommendations for a
 29 global geochemical data framework, required to accomplish the scientific challenges of the
 30 21st century and beyond.

31 2. Motivation

32 2.1. Diversity and Fragmentation of Geochemical Data

33 We understand geochemistry as the discipline that integrates geology and chemistry
 34 by using the principles and tools of chemistry to discover and develop fundamental under-

standing 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 many other disciplines including, but not limited to, archaeology, environmental science and technology, resource exploration and development (groundwater, minerals, energy), geohealth, oceanography, and agriculture and is 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 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, e.g., individual computers, locally developed database solutions, or 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 Ryborn, 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 sectors (academia, government, 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/>), GEOROC (<https://georoc.eu/>), MetBase (<https://metbase.org/>), and the Astromaterials Data System (<https://www.astromat.org/>). More recent initiatives are National Research Infrastructures in Germany (NFDI4Earth), Europe (EPOS), Australia (AuScope), the US (EarthCube), or Norway (NIRD), to name a few. However, walls around individual data silos remain, hindering simple, inclusive and global access to geochemical data. To overcome these walls, we must develop common, community-agreed, global standards for geochemical data and metadata. These standards are critical to making geochemical data Findable, Accessible, Interoperable and Reusable to both humans and machines (FAIR; Wilkinson et al., 2016). Not only will FAIR data standards and curation procedures increase the value of new data as they are generated and published, they likewise have large potential for utilising the significant proportion of unpublished geochemical data in research and public sectors from the last century.

Recognising that mainstream scientific journals were the most effective agents to rectify problems in data reporting and implement best practices, an Editors Roundtable

was launched in 2007 as an initiative to bring together editors, publishers, and database providers to implement consistent publication practices for geochemical data. Academic societies such as the Geochemical Society also adopted a policy for geochemical data publication at that time (Geochemical Society, 2007). The Editors Roundtable created and signed a policy statement in January 2009 (version 1.1) that laid out ‘Requirements for the Publication of Geochemical Data’ (Goldstein et al., 2014).

Recently, the nationally-funded, global data systems EarthChem (USA), GEOROC (Germany), EPOS-MSL (European Plate Observing System MultiScale Laboratories, Europe) and AusGeochem (Australia) came together to enable interoperability between their systems. Yet a vast amount of geochemical data lies outside these initiatives. In response to Open Science policies and demands from the scientific community, a Town Hall meeting on ‘OneGeochemistry: Toward a Global Network of Geochemistry Data’ (<https://www.agu.org/Fall-Meeting-2019/Events/Data-TH23L>) 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. The goal of the 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 CODATA meetings (e.g. Lehnert et al., 2021; Wyborn et al., 2021).

2.3. OneGeochemistry Mission

OneGeochemistry is taking action to develop and promote global, community-driven data conventions and best practices necessary to build a global network of trusted geochemical data. These actions will enable and simplify the (re)use of geochemical data and 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 also 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. \(2016\)](#) and be readable by both humans and machines. 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 (Collective Benefit, Authority to Control, Responsibility, and Ethics) Principles for Indigenous Data Governance to protect Indigenous rights and interests in Indigenous data (including traditional knowledge),

particularly in the sample collection phase (Carroll 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. (2020); Abbott et al. (2022); Horstwood et al. (2016); Dutton et al. (2017); Walker et al. (2008); Mustaphi et al. (2019); Schaen et al. (2020); Khider et al. (2019); Damerow et al. (2021). 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/>). A further common sticking point is that the vocabularies recommended to define each data type are not FAIR and are not available from online repositories such as Research Vocabularies Australia (<https://vocabs.ardc.edu.au/>) or FAIRsharing (<https://fairsharing.org/>). There is also no evidence in most of these papers that the recommended vocabularies have a governance structure in place that allows them to evolve.

OneGeochemistry aims to become an organisation that would coordinate across all geo- and cosmochemical data types. Fundamental to its 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.

3. Challenges for the Community

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

3) Synthesis Databases. Laboratory Information Management Systems focus on physical samples and cover the first half of the research data lifecycle from sample collection or generation to processing and analysis (Fig. 2). Examples of such systems include AusGeochem (<https://www.auscope.org.au/ausgeochem>) and Sparrow (<https://sparrow-data.org/>). The final data products derived from samples are then published in Repositories. Generalist repositories, such as Figshare (<https://figshare.com/>), Dryad (<https://datadryad.org/>) or Zenodo (<https://zenodo.org/>), publish research outputs irrespective of academic discipline. Domain repositories, in contrast, cater to specific disciplines or subdisciplines and therefore offer data services targeted to the particular requirements of these domains. PANGAEA (<https://www.pangaea.de/>) and GFZ Data Services (<https://bib.teagrafenberg.de/dataservices/>) are examples of domain repositories for the Earth Sciences, whilst the EarthChem Library (<https://earthchem.org/ec1/>) or the GEOROC Data Repository (<https://georoc.eu/>) are domain repositories specifically for geochemical data. Synthesis Databases compile these individual data publications, as well as harvesting data from the scientific literature, to enable data discovery and reuse across multiple datasets. Similar to domain repositories, synthesis databases usually specialise in a particular subdiscipline or have a geographical focus. AstroMat, GEOROC, MetBase and PetDB 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.

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

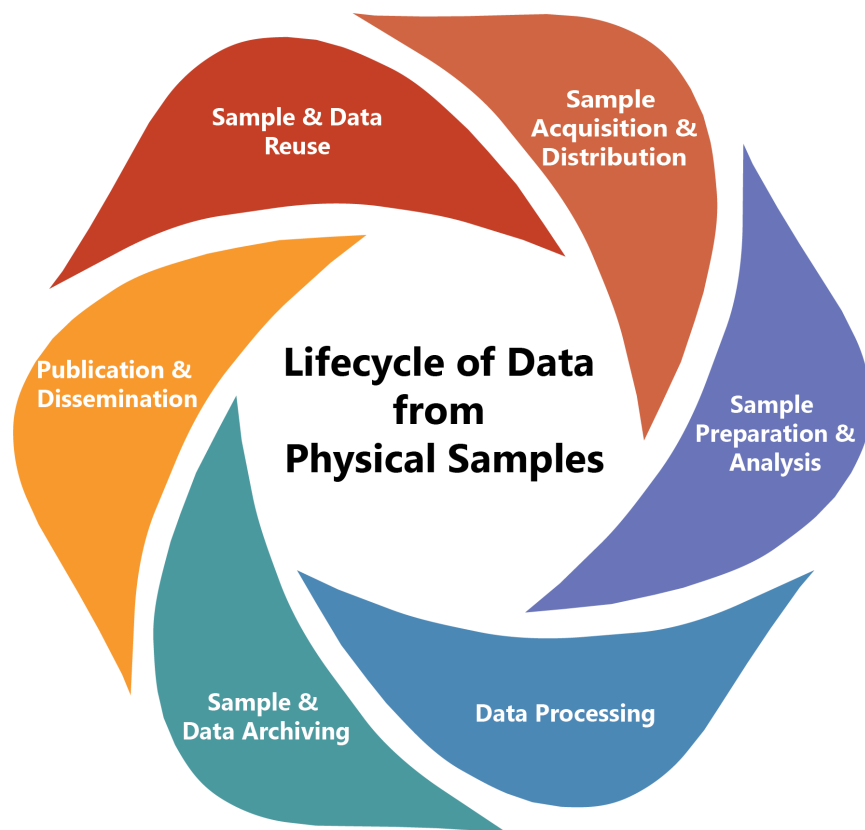


Figure 2: The sample and data life cycle from field to publication to reuse (adapted from [Ramdeen et al., 2022](#)). Tools that support researchers throughout this process include SESAR, a registry for physical samples. AusGeochem or the StraboSpot-Selenocene-Sparrow system of EarthCube support researchers 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 AstroMat, EarthChem, GEOROC or MetBase facilitate dissemination and data reuse.

redundancy and the use of unique identifiers (e.g. DOI, 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.

- 201 3. **Data Quality:** Data are reliable, i.e. they follow a common standard that ensures
202 availability of rich sample and analytical metadata (e.g. provenance, description of
203 method and analysis conditions). Completeness of metadata allows assessment of
204 accuracy and ensures reproducibility.
- 205 4. **Attribution:** Appropriate citation of the people, laboratories, organisations, fun-
206 ders, research artefacts and data is ensured through use of globally unique, persistent
207 and resolvable identifiers and compliance with international metadata standards (e.g.
208 ORCID, ROR, DataCite, Crossref, IGSN).

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

218 3.1. Challenges for Researchers

219 The current research culture in geochemistry means that only few researchers are
220 willing to share their data (Chamberlain et al., 2021). Although the recent push for
221 open science has also benefited the open data landscape, community understanding and
222 adoption are still centred around individuals. The majority of data producers remain
223 reluctant to share their data unless forced by journal or funding requirements: the in-
224 set in Fig. 1 shows the rapid increase in submissions to the EarthChem Library after
225 several of the AGU journals enforced data publications in trusted domain repositories
226 in 2019 ([https://www.agu.org/Share-and-Advocate/Share/Polymakers/Position-
227 Statements/Position_Data](https://www.agu.org/Share-and-Advocate/Share/Polymakers/Position-Statements/Position_Data)). The lack of adoption by the research community is, in part,

caused by a number of considerable challenges facing those researchers that are willing to share their data.

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 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. Springer Nature, Science and AGU are leaders in this respect, requiring data publication in domain repositories prior to manuscript acceptance, yet each have developed their own—differing—guidelines on how to comply with this policy. Data journals such as Data in Brief and Scientific Data also require data submission to (domain) repositories and, in addition, provide 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 CSV or XLS to PDF and even JPEG as part of supplementary materials or they encourage submission to generalist repositories, such as Figshare, Zenodo or Dryad, where there is no quality control or reporting standard. 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 lest their manuscript be rejected. When faced with the complexity 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, in very different formats across a wide range of repositories. In addition to the many different formats that make data hard to combine and compare, 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 for Science, a journal that adopted an open data policy in 2016, 30% of articles do not publish their data at all (Yeston, 2021).

Complexity of data submission: Good data management takes time. The assembling and submission of data tables and related information require time and additional effort outside of the primary process of manuscript submission. Usually, substantial processing is performed on raw data coming from an analytical instrument. While this processing 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 supplementary 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. Even though data publications are increasingly visible via (automatic) indexing in ORCID profiles, for example, they are rarely counted towards the research track record or valued by recruiting and promotion committees. Whilst assigning DOIs to datasets helps to emphasise the value of data publications, the lack of awareness in the broader research community means that these publications are often not appropriately cited. In addition, researchers who consider submitting to domain repositories are often deterred by the additional processing time before the final data publication. The EarthChem Library, for example, that specialises in geochemical data, advises a turnaround time ranging from a few days to up to two weeks. PANGAEA, a domain repository for all disciplines within the Earth Sciences, has a data publication

282 timeline of three months. Even though there are good reasons for these timelines— mostly
283 centred around curation as discussed below—, they discourage even more researchers from
284 publishing their data.

285 Sensitive data: An important consideration within both the FAIR and CARE princi-
286 ples is how to handle sensitive data that should only be discoverable by certain, authorised
287 persons or only available after an embargo period. This access control is particularly
288 important for data produced or funded by industry and for agencies that deal with clas-
289 sified information. Good technical solutions already exist, simply requiring clear licensing
290 of datasets and the ability of repositories to handle management of temporary embargo
291 periods during the publication phase.

292 Variable quality of the available published data: The final issue to be highlighted here
293 are the considerable challenges caused by a lack of standard formats for publishing geo-
294 chemical data, which often precludes quality assessment and, therefore, reuse of published
295 data. Common issues include: dead links or non-existent supplementary material; errors
296 in data reporting; lack of reproducibility due to missing analytical information; and the
297 use of abbreviations only understood by the owner of the dataset. Data quality assessment
298 is often impossible due to a lack of analytical details or measures of uncertainties, includ-
299 ing inconsistent units on uncertainty reporting (e.g. standard deviations, standard errors,
300 confidence interval, 1σ vs. 2σ errors, etc.). When compiling data from multiple sources,
301 additional challenges include inconsistent, non-standardised terminology (e.g. eclogite vs
302 arclogite) and missing units of measurement. Finally, the original owner, funder, and/or
303 creator of the data are rarely credited in synthesis or compiled datasets.

304 *3.2. Challenges for Data Systems*

305 Some of the challenges for researchers detailed above are related to current limitations
306 of data repositories and synthesis databases. One major issue lies with the resources avail-
307 able to these data systems and the sustainability of funding. Long-term staffing solutions
308 for data curators that assist researchers with data submissions are vital for data systems.

309 The advantage of publishing data in domain repositories is that the research data are doc-
310 umented in a format specific to the discipline and the respective data type, which ensures
311 that data quality can be easily assessed and data users have greater trust in individual
312 datasets. By collecting data in domain repositories, they are also more visible and easier
313 to discover for others in the field, leading to greater reuse— and ultimately citation— of
314 these data.

315 Yet in order to consistently provide this service, domain repositories need to employ cu-
316 rators with domain expertise who carefully review each data submission. Many researchers
317 of today are not familiar with proper data management, and hence data submissions are
318 not consistent: column headers are not standardised, there are spelling errors, inconsistent
319 text, and widespread use of non-standard abbreviations. While it takes the researchers a
320 considerable amount of time to collate this information, repository curators then need to
321 invest further time to convert submissions to their internal standard and ensure all data
322 and metadata are transparent and easy to understand by third parties.

323 More often than not, repositories are not funded for this additional work and are strug-
324 gling with staffing issues. This issue arises because many of the data systems catering to a
325 specific domain were born out of research projects that succeeded in attracting additional
326 funding to further develop their infrastructure. However, this funding is usually temporary
327 and restricted to the development of new technologies or services— system maintenance
328 and curation are rarely funded by national science foundations. What is more, these data
329 systems compete for funding with researchers within their domain. Far too often, data sys-
330 tems that are widely used by the research community are orphaned because of discontinued
331 funding: MetPetDB, SedDB and NAVDAT are all pertinent examples of such systems that
332 are no longer maintained, and at worst are no longer available to the community.

333 The availability of resources is intricately linked with community-uptake of domain
334 repository services. For many data systems, it is an ongoing struggle to entice more
335 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

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.

4.1. 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 (Bruno et al., 2017; Groom et al., 2016; Bergerhoff and Brown, 1987; Berman et al., 2003) ensure the long term preservation and access to derived results and processed data published in standard formats, and 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). 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 (Spek, 2020) and validation of the data is part of the peer review process for journals. Some journals that publish papers on crystallography also sponsor the development of validation tools.

4.2. 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 standardise nomenclature, including 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).

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/>). The FDSN began in 1984 when multiple countries agreed to create a global network around those 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 (Standard for Exchange of Earthquake Data) format. 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).

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 November 2022, the OBPS site contains 1728 best practice documents from 52 institutions/organisations: as new documents are submitted, they are reviewed and endorsed by expert teams (Przeslawski et al., 2022).

Each institution/organisation can submit their best practice documents including quality documents specific to their data acquisition programs. The Australian Integrated Marine Observing System (IMOS), 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, while at the same time, best practices can be adapted to match regional infrastructure capabilities

(Przeslawski et al., 2022).

The International Oceans Drilling 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 publication portal SEDIS (Scientific Earth Drilling Information System; <http://sedis.iodp.org/>).

4.6. What can be learned from these initiatives?

The examples from crystallography, chemistry, seismology, geology and oceanography show that it is indeed possible to unite a community 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 sources;
2. Establishing expert working groups 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 disseminate existing standards and best practices for data users and contributors;
10. Incorporating data management into the undergraduate curriculum.

5. The Path Forward: OneGeochemistry

During the workshop at Goldschmidt 2022, organisers and participants discussed possible solutions to the aforementioned challenges. The options promising the highest short-term impact are: official endorsement of the OneGeochemistry initiative; establishment of expert working groups 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).

OneGeochemistry is proposing to form a CODATA Working Group to bring together all the disparate initiatives that are happening in geochemistry across Scientific Unions, Associations, Societies and Commissions. The OneGeochemistry interim board has applied to the following seven international geochemical societies and associations for endorsement of this work: Geochemical Society, European Association of Geochemistry, International Association of Geoanalysts, International Association of Geochemists, Association of Applied Geochemists, IUGS Commission on Global Geochemical Baselines and Meteoritical Society. Further national and/or sub-disciplinary societies will be contacted in the future and the OneGeochemistry board is open to additional 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 geochemistry data standards 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 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 seismological 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 Working Groups

There are multiple different standards currently in use and a growing number of publications aim to establish agreement on minimum variables and vocabularies for various geochemical data types (Table 1). 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 working groups that are endorsed and/or recognised by authoritative, international geochemical societies and unions. The task of these expert working groups would be to develop standards for each distinct analytical technique or related groups of analytical methods. A working group 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 could be to facilitate regular workshops or hackathons to develop, improve and disseminate best practice recommendations and invite feedback from the wider community. In a first step, OneGeochemistry would work with the wider community to determine which 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 working groups 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;
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. Continuously monitor and update each agreed upon standard;
5. Develop new data standards where required.

A successful example of an expert working group in geochemistry is the Tephra Community that has developed data submission templates for the EarthChem Library. 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.

OneGeochemistry has been offered the Brown Book by IUPAC as part of their Colour Books Series described above. This resource will be invaluable not only in documenting nomenclature defined by the geochemical expert working groups but also in ensuring that relevant, existing digital chemical standards are leveraged wherever possible (e.g., the Machine-Accessible Periodic Table).

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. A surprising outcome of the Goldschmidt 2022 workshop was the observation how poorly known the existing data systems are, especially among early career researchers. Options for increased community engagement include:

1. Attribution and advertising of OneGeochemistry and the existing geochemical data systems in research outputs through:

- citation of data systems in publications following citation guidelines and templates provided by the systems
- Encouraging the addition of data system logos to presentation materials (e.g. conference slides, poster, graphical abstract)
- Where tools for plotting or analysis are provided by data systems, resulting figures should be watermarked

2. Virtual activities and resources

- Maintaining and increasing a social media presence (e.g. Twitter, LinkedIn);
- Using blog posts, webinars and a dedicated YouTube channel to disseminate tutorials and teach data management skills

- Organising data hackathons
- Developing ready-to-use teaching lessons, materials, slides and exercises to increase use of databases in teaching (e.g. ask lecturers globally to “publish” and share their materials, and promote the teaching materials to lecturers). Seek collaborations with dedicated educational initiatives, such as NAGT (<http://www.nagt.org>) and SERC (<https://serc.carleton.edu/index.html>).

3. Workshops and Data-Mentoring

- Continue hosting workshops at scientific conferences (e.g., Goldschmidt, GSA, AGU, EGU), where future expert working groups present their progress and liaise with the wider community
- Contributing to the Data Help Desks coordinated by ESIP at major Earth Science conferences such as AGU, EGU, Geological Society of America (<https://www.esipfed.org/data-help-desk>) and holding Data FAIR workshops (<https://data.agu.org/datafair/>)
- Integrating data management into mentoring schemes at these conferences
- Implementing inter-institution and international data mentoring programs that also focus on available resources in the communities

OneGeochemistry ambassadors will be recruited to assist with these activities and initiatives. Ambassadors are envisaged as 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.

While communicating and advertising OneGeochemistry, we must also be aware of motivations and incentives (or disincentives) to contribute to standard development, data publication and global databases for each stakeholder. These incentives will differ between

different groups in the community (Fig. 3). 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 could further 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



Figure 3: 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 working groups and facilitating education and ambassadorship.

of impact of specific research outputs. Tracking data citations is also a convenient way for *funders*, *institutions* and *laboratories* to measure their impact. 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

above even though they likely produce far larger data volumes than both academic and governmental communities and might have the human and financial resources to quickly develop and implement data standards. However, some countries like Australia require that all data be made available to local geological surveys after a certain time period— providing an incentive to comply with the 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 and are 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 scientific questions. We propose OneGeochemistry as a community-driven initiative that can enact this change by building a global, online network of machine-readable data that is persistent, interoperable and reusable and above all, minimises duplication of the same data. Such a system will also ensure reliable citation for those who collected, analysed, curated and made accessible any geochemical and cosmochemical data. Endorsement by societies, publishers and funders will give OneGeochemistry the authority to establish expert working groups that develop and promote best practices and standards for specific data types. We will seek to increase community engagement 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) for helpful comments on the history of the development of data standards in seismology. Michael Badawi, Jieun Kim and Nicolas Randazzo are thanked for their contributions to the Goldschmidt 2022

workshop. MK is supported by the German Research Foundation (DFG grant 437919684). 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. HB was supported by DFG CRC TRR 170 (Project-ID 263649064). DCH is supported through the NFDI4Earth funded by the DFG (project number: 460036893). AWL is funded by the Deutscher Akademischer Austauschdienst (German Research Exchange Service, Grant No. 57507869) and an Australian Government Research Training Program (RTP) Stipend and RTP Fee-Offset Scholarship through Macquarie University (Allocation No. 2018177).

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. URL: <https://zenodo.org/record/3866266>, doi:10.5281/ZENODO.3866266.
- Alexakis, D.E., 2021. Linking DPSIR model and water quality indices to achieve sustainable development goals in groundwater resources. *Hydrology* 8, 90. URL: <https://doi.org/10.3390/hydrology8020090>, doi:10.3390/hydrology8020090.
- Bergerhoff, G., Brown, I., 1987. International Union of Crystallography, Chester.

702 Berman, H., Henrick, K., Nakamura, H., 2003. Announcing the worldwide protein data
 703 bank. *Nature Structural & Molecular Biology* 10, 980–980. URL: [https://doi.org/](https://doi.org/10.1038/nsb1203-980)
 704 [10.1038/nsb1203-980](https://doi.org/10.1038/nsb1203-980), doi:[10.1038/nsb1203-980](https://doi.org/10.1038/nsb1203-980).

705 Boone, S.C., Dalton, H., Prent, A., Kohlmann, F., Theile, M., Gréau, Y., Florin, G.,
 706 Noble, W., Hodgekiss, S.A., Ware, B., Phillips, D., Kohn, B., O'Reilly, S., Gleadow, A.,
 707 McInnes, B., Rawling, T., 2022. AusGeochem: An open platform for geochemical data
 708 preservation, dissemination and synthesis. *Geostandards and Geoanalytical Research*
 709 46, 245–259. URL: <https://doi.org/10.1111/ggr.12419>, doi:[10.1111/ggr.12419](https://doi.org/10.1111/ggr.12419).

710 Brantley, S., Wen, T., Agarwal, D., Catalano, J., Schroeder, P., Lehnert, K., Varadharajan,
 711 C., Pett-Ridge, J., Engle, M., Castronova, A., Hooper, R., Ma, X., Jin, L., McHenry, K.,
 712 Aronson, E., Shaughnessy, A., Derry, L., Richardson, J., Bales, J., Pierce, E., 2020. A
 713 vision for the future low-temperature geochemical data-scape URL: [https://doi.org/](https://doi.org/10.31223/x5zp5w)
 714 [10.31223/x5zp5w](https://doi.org/10.31223/x5zp5w), doi:[10.31223/x5zp5w](https://doi.org/10.31223/x5zp5w).

715 Bruno, I., Gražulis, S., Helliwell, J.R., Kabekkodu, S.N., McMahon, B., Westbrook, J.,
 716 2017. Crystallography and databases. *Data Science Journal* 16. URL: [https://doi.](https://doi.org/10.5334/dsj-2017-038)
 717 [org/10.5334/dsj-2017-038](https://doi.org/10.5334/dsj-2017-038), doi:[10.5334/dsj-2017-038](https://doi.org/10.5334/dsj-2017-038).

718 Bundschuh, J., Maity, J.P., Mushtaq, S., Vithanage, M., Seneweera, S., Schneider, J., Bhat-
 719 tacharya, P., Khan, N.I., Hamawand, I., Guilherme, L.R., Reardon-Smith, K., Parvez,
 720 F., Morales-Simfors, N., Ghaze, S., Pudmenzky, C., Kouadio, L., Chen, C.Y., 2017.
 721 Medical geology in the framework of the sustainable development goals. *Science of The*
 722 *Total Environment* 581-582, 87–104. URL: [https://doi.org/10.1016/j.scitotenv.](https://doi.org/10.1016/j.scitotenv.2016.11.208)
 723 [2016.11.208](https://doi.org/10.1016/j.scitotenv.2016.11.208), doi:[10.1016/j.scitotenv.2016.11.208](https://doi.org/10.1016/j.scitotenv.2016.11.208).

724 Carbotte, S., Lehnert, K., 2007. WORKSHOP REPORT | building a global data network
 725 for studies of earth processes at the world's plate boundaries. *Oceanography* 20, 124–125.
 726 URL: <https://doi.org/10.5670/oceanog.2007.38>, doi:[10.5670/oceanog.2007.38](https://doi.org/10.5670/oceanog.2007.38).

- Carroll, S.R., Garba, I., Figueroa-Rodríguez, O.L., Holbrook, J., Lovett, R., Materechera, S., Parsons, M., Raseroka, K., Rodriguez-Lonebear, D., Rowe, R., Sara, R., Walker, J.D., Anderson, J., Hudson, M., 2020. The CARE principles for indigenous data governance. *Data Science Journal* 19. URL: <https://doi.org/10.5334/dsj-2020-043>, doi:10.5334/dsj-2020-043.
- Chamberlain, K.J., Lehnert, K.A., McIntosh, I.M., Morgan, D.J., Wörner, G., 2021. Time to change the data culture in geochemistry. *Nature Reviews Earth & Environment* 2, 737–739. URL: <https://doi.org/10.1038/s43017-021-00237-w>, doi:10.1038/s43017-021-00237-w.
- Damerow, J.E., Varadharajan, C., Boye, K., Brodie, E.L., Burrus, M., Chadwick, K.D., Crystal-Ornelas, R., Elbashandy, H., Alves, R.J.E., Ely, K.S., Goldman, A.E., Haberman, T., Hendrix, V., Kakalia, Z., Kemner, K.M., Kersting, A.B., Merino, N., O'Brien, F., Perzan, Z., Robles, E., Sorensen, P., Stegen, J.C., Walls, R.L., Weisenhorn, P., Zavarin, M., Agarwal, D., 2021. Sample identifiers and metadata to support data management and reuse in multidisciplinary ecosystem sciences. *Data Science Journal* 20, 11. URL: <https://doi.org/10.5334/dsj-2021-011>, doi:10.5334/dsj-2021-011.
- Deines, P., Goldstein, S.L., Oelkers, E.H., Rudnick, R.L., Walter, L.M., 2003. Standards for publication of isotope ratio and chemical data in chemical geology. *Chemical Geology* 202, 1–4. URL: <https://doi.org/10.1016/j.chemgeo.2003.08.003>, doi:10.1016/j.chemgeo.2003.08.003.
- 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. URL: <https://zenodo.org/record/7267967>, 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.

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. URL: <https://doi.org/10.1016/j.quageo.2017.03.001>, doi:10.1016/j.quageo.2017.03.001.

Dziewonski, A.M., 1994. The FDSN: history and objectives. *Annals of Geophysics* 37.
URL: <https://doi.org/10.4401/ag-4191>, 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. URL: <https://data.europa.eu/doi/10.2777/02999>, 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.
URL: <https://doi.org/10.1029/2015eo036971>, doi:10.1029/2015eo036971.

Fairbridge, R.W., 1998. History of geochemistry, in: *Encyclopedia of Earth Science*. Kluwer
Academic Publishers, pp. 315–322. URL: https://doi.org/10.1007/1-4020-4496-8_156, doi:10.1007/1-4020-4496-8_156.

Flowers, R., Zeitler, P., Danišík, 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* URL: <https://doi.org/10.1130/b36266.1>,
doi:10.1130/b36266.1.

777 Geochemical Society, 2007. Geochemical society policy on geochemical databases. URL:
778 <https://www.geochemsoc.org/about/positionstatements/datapolicy>.

779 Gill, J.C., 2017. Geology and the sustainable development goals. Episodes 40, 70–
780 76. URL: <https://doi.org/10.18814/epiiugs/2017/v40i1/017010>, doi:[10.18814/epiiugs/2017/v40i1/017010](https://doi.org/10.18814/epiiugs/2017/v40i1/017010).

782 Goldstein, S., Lehnert, K., Hofmann, A., 2014. Requirements for the publication of geo-
783 chemical data. URL: <https://ecl.earthchem.org/view.php?id=735>, doi:[10.1594/IEDA/100426](https://doi.org/10.1594/IEDA/100426).

785 Groom, C.R., Bruno, I.J., Lightfoot, M.P., Ward, S.C., 2016. The cambridge structural
786 database. Acta Crystallographica Section B Structural Science, Crystal Engineering and
787 Materials 72, 171–179. URL: <https://doi.org/10.1107/s2052520616003954>, doi:[10.1107/s2052520616003954](https://doi.org/10.1107/s2052520616003954).

789 Hall, S.R., Cook, A.P.F., 1995. STAR dictionary definition language: Initial specification.
790 Journal of Chemical Information and Computer Sciences 35, 819–825. URL: <https://doi.org/10.1021/ci00027a005>, doi:[10.1021/ci00027a005](https://doi.org/10.1021/ci00027a005).

792 Hall, S.R., McMahon, B., 2016. The implementation and evolution of STAR/CIF ontolo-
793 gies: Interoperability and preservation of structured data. Data Science Journal 15.
794 URL: <https://doi.org/10.5334/dsj-2016-003>, doi:[10.5334/dsj-2016-003](https://doi.org/10.5334/dsj-2016-003).

795 Heidorn, P.B., 2008. Shedding light on the dark data in the long tail of science. Library
796 Trends 57, 280–299. URL: <https://doi.org/10.1353/lib.0.0036>, doi:[10.1353/lib.0.0036](https://doi.org/10.1353/lib.0.0036).

798 Heller, S., McNaught, A., Stein, S., Tchekhovskoi, D., Pletnev, I., 2013. InChI - the
799 worldwide chemical structure identifier standard. Journal of Cheminformatics 5. URL:
800 <https://doi.org/10.1186/1758-2946-5-7>, doi:[10.1186/1758-2946-5-7](https://doi.org/10.1186/1758-2946-5-7).

Horstwood, M.S.A., Košler, J., Gehrels, G., Jackson, S.E., McLean, N.M., Paton, C.,
 Pearson, N.J., Sircombe, K., Sylvester, P., Vermeesch, P., Bowring, J.F., Condon, D.J.,
 Schoene, B., 2016. Community-derived standards for LA-ICP-MS u-(th)pb geochronol-
 ogy – uncertainty propagation, age interpretation and data reporting. *Geostandards*
 and *Geoanalytical Research* 40, 311–332. URL: [https://doi.org/10.1111/j.1751-](https://doi.org/10.1111/j.1751-908x.2016.00379.x)
[908x.2016.00379.x](https://doi.org/10.1111/j.1751-908x.2016.00379.x), doi:[10.1111/j.1751-908x.2016.00379.x](https://doi.org/10.1111/j.1751-908x.2016.00379.x).

Jackson, I., Wyborn, L., 2008. International viewpoint and news. *Environmental Geology*
 53, 1377–1380. URL: <https://doi.org/10.1007/s00254-007-1085-z>, doi:[10.1007/](https://doi.org/10.1007/s00254-007-1085-z)
[s00254-007-1085-z](https://doi.org/10.1007/s00254-007-1085-z).

Khider, D., Emile-Geay, J., McKay, N.P., Gil, Y., Garijo, D., Ratnakar, V., Alonso-Garcia,
 M., Bertrand, S., Bothe, O., Brewer, P., Bunn, A., Chevalier, M., Comas-Bru, L., Csank,
 A., Dassié, E., DeLong, K., Felis, T., Francus, P., Frappier, A., Gray, W., Goring, S.,
 Jonkers, L., Kahle, M., Kaufman, D., Kehrwald, N.M., Martrat, B., McGregor, H.,
 Richey, J., Schmittner, A., Scroxton, N., Sutherland, E., Thirumalai, K., Allen, K.,
 Arnaud, F., Axford, Y., Barrows, T., Bazin, L., Birch, S.E.P., Bradley, E., Bregy, J.,
 Capron, E., Cartapanis, O., Chiang, H.W., Cobb, K.M., Debret, M., Dommain, R., Du,
 J., Dyez, K., Emerick, S., Erb, M.P., Falster, G., Finsinger, W., Fortier, D., Gauthier,
 N., George, S., Grimm, E., Hertzberg, J., Hibbert, F., Hillman, A., Hobbs, 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. URL: <https://doi.org/10.1029/2019pa003632>,
 doi:[10.1029/2019pa003632](https://doi.org/10.1029/2019pa003632).

827 Kroon-Batenburg, L.M.J., Helliwell, J.R., Hester, J.R., 2022. *IUCrData* launches raw data
 828 letters. *IUCrData* 7. URL: <https://doi.org/10.1107/s2414314622008215>, doi:[10.](https://doi.org/10.1107/s2414314622008215)
 829 [1107/s2414314622008215](https://doi.org/10.1107/s2414314622008215).

830 Lehnert, K., Wyborn, L., Bennett, V.C., Hezel, D., McInnes, B.I.A., Plank, T., Rubin, K.,
 831 2021. Onegeochemistry: Towards an interoperable global network of fair geochemical
 832 data URL: <https://zenodo.org/record/5767950>, doi:[10.5281/ZENODO.5767950](https://doi.org/10.5281/ZENODO.5767950).

833 Mustaphi, C.J.C., Brahney, J., Aquino-López, M.A., Goring, S., Orton, K., Noronha, A.,
 834 Czaplewski, J., Asena, Q., Paton, S., Brushworth, J.P., 2019. Guidelines for reporting
 835 and archiving 210pb sediment chronologies to improve fidelity and extend data lifecycle.
 836 *Quaternary Geochronology* 52, 77–87. URL: [https://doi.org/10.1016/j.quageo.](https://doi.org/10.1016/j.quageo.2019.04.003)
 837 [2019.04.003](https://doi.org/10.1016/j.quageo.2019.04.003), doi:[10.1016/j.quageo.2019.04.003](https://doi.org/10.1016/j.quageo.2019.04.003).

838 OECD, 2017. Co-ordination and support of international research data networks. URL:
 839 <https://doi.org/10.1787/e92fa89e-en>, doi:[10.1787/e92fa89e-en](https://doi.org/10.1787/e92fa89e-en).

840 Peng, G., Lacagnina, C., Downs, R.R., Ganske, A., Ramapriyan, H.K., Ivánová, I.,
 841 Wyborn, L., Jones, D., Bastin, L., Lin Shie, C., Moroni, D.F., 2022. Global community
 842 guidelines for documenting, sharing, and reusing quality information of individual digital
 843 datasets. *Data Science Journal* 21. URL: <https://doi.org/10.5334/dsj-2022-008>,
 844 doi:[10.5334/dsj-2022-008](https://doi.org/10.5334/dsj-2022-008).

845 Przeslawski, R., Pearlman, J., Karstensen, J., 2022. Dataset quality information in
 846 australia’s integrated marine observing system, in: *SciDataCon 2022*. URL: [https:](https://www.scidatacon.org/IDW-2022/sessions/431/paper/969/)
 847 [//www.scidatacon.org/IDW-2022/sessions/431/paper/969/](https://www.scidatacon.org/IDW-2022/sessions/431/paper/969/).

848 Ramdeen, S., Wyborn, L.A.I., Lehnert, K.A., Klump, J., 2022. The role of unique iden-
 849 tifiers in tracing the life cycle of a sample and any data derived from it, in: *Gold-*
 850 *schmidt2022 abstracts*, Geochemical Society. URL: [https://conf.goldschmidt.info/](https://conf.goldschmidt.info/goldschmidt/2022/meetingapp.cgi/Paper/12644)
 851 [goldschmidt/2022/meetingapp.cgi/Paper/12644](https://conf.goldschmidt.info/goldschmidt/2022/meetingapp.cgi/Paper/12644).

852 Ruth, P.D.V., Atkins, N., 2022. Dataset quality information in australia’s integrated
 853 marine observing system, in: SciDataCon 2022. URL: [https://www.scidatacon.org/](https://www.scidatacon.org/IDW-2022/sessions/431/paper/1052/)
 854 [IDW-2022/sessions/431/paper/1052/](https://www.scidatacon.org/IDW-2022/sessions/431/paper/1052/).

855 Schaen, A.J., Jicha, B.R., Hodges, K.V., Vermeesch, P., Stelten, M.E., Mercer, C.M.,
 856 Phillips, D., Rivera, T.A., Jourdan, F., Matchan, E.L., Hemming, S.R., Morgan, L.E.,
 857 Kelley, S.P., Cassata, W.S., Heizler, M.T., Vasconcelos, P.M., Benowitz, J.A., Koppers,
 858 A.A., Mark, D.F., Niespolo, E.M., Sprain, C.J., Hames, W.E., Kuiper, K.F.,
 859 Turrin, B.D., Renne, P.R., Ross, J., Nomade, S., Guillou, H., Webb, L.E., Cohen,
 860 B.A., Calvert, A.T., Joyce, N., Ganerød, M., Wijbrans, J., Ishizuka, O., He, H.,
 861 Ramirez, A., Pfänder, J.A., Lopez-Martínez, M., Qiu, H., Singer, B.S., 2020. Interpreting
 862 and reporting 40ar/39ar geochronologic data. GSA Bulletin 133, 461–487. URL:
 863 <https://doi.org/10.1130/b35560.1>, doi:[10.1130/b35560.1](https://doi.org/10.1130/b35560.1).

864 Sen, M., Duffy, T., 2005. GeoSciML: Development of a generic GeoScience markup language.
 865 Computers & Geosciences 31, 1095–1103. URL: [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.cageo.2004.12.003)
 866 [cageo.2004.12.003](https://doi.org/10.1016/j.cageo.2004.12.003), doi:[10.1016/j.cageo.2004.12.003](https://doi.org/10.1016/j.cageo.2004.12.003).

867 Spek, A.L., 2020. *checkCIF* validation ALERTS: what they mean and how to respond.
 868 Acta Crystallographica Section E Crystallographic Communications 76, 1–11. URL:
 869 <https://doi.org/10.1107/s2056989019016244>, doi:[10.1107/s2056989019016244](https://doi.org/10.1107/s2056989019016244).

870 Stall, S., McEwen, L., Wyborn, L., Hoebelheinrich, N., Bruno, I., 2020. Growing the
 871 FAIR community at the intersection of the geosciences and pure and applied chemistry.
 872 Data Intelligence 2, 139–150. URL: https://doi.org/10.1162/dint_a_00036, doi:[10.](https://doi.org/10.1162/dint_a_00036)
 873 [1162/dint_a_00036](https://doi.org/10.1162/dint_a_00036).

874 Taylor, R., Wood, P.A., 2019. A million crystal structures: The whole is greater than the
 875 sum of its parts. Chemical Reviews 119, 9427–9477. URL: [https://doi.org/10.1021/](https://doi.org/10.1021/acs.chemrev.9b00155)
 876 [acs.chemrev.9b00155](https://doi.org/10.1021/acs.chemrev.9b00155), doi:[10.1021/acs.chemrev.9b00155](https://doi.org/10.1021/acs.chemrev.9b00155).

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. URL: <https://zenodo.org/record/4313859>, doi:[10.5281/ZENODO.4313859](https://doi.org/10.5281/ZENODO.4313859).

Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., Bouwman, J., Brookes, A.J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers, R., Gonzalez-Beltran, A., Gray, A.J., Groth, P., Goble, C., Grethe, J.S., Heringa, J., 't Hoen, P.A., Hoof, R., Kuhn, T., Kok, R., Kok, J., Lusher, S.J., Martone, M.E., Mons, A., Packer, A.L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone, S.A., Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M.A., Thompson, M., van der Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., Mons, B., 2016. The FAIR guiding principles for scientific data management and stewardship. Scientific Data 3. URL: <https://doi.org/10.1038/sdata.2016.18>, doi:[10.1038/sdata.2016.18](https://doi.org/10.1038/sdata.2016.18).

Wyborn, L., Elger, K., Prent, A., Lehnert, K., Bruno, I., Klöcking, M., Klump, J., Profeta, L., Quinn, D.P., Ramdeen, S., ter Maat, G., 2021. The onegeochemistry initiative: Mobilising a global network of fair geochemical data to support research into the grand challenge of an environmentally sustainable future URL: <https://zenodo.org/record/5765464>, doi:[10.5281/ZENODO.5765464](https://doi.org/10.5281/ZENODO.5765464).

Wyborn, L., Lehnert, K., 2021. OneGeochemistry: Creating a global network of geochemical data to support the 17 united nations sustainable development goals, in: Goldschmidt2021 abstracts, European Association of Geochemistry. URL: <https://doi.org/10.7185/gold2021.6562>, doi:[10.7185/gold2021.6562](https://doi.org/10.7185/gold2021.6562).

Wyborn, L.A.I., Ryborn, R.J., 1989. PETCHEM data set : Australia and Antarctica -

902 , documentation. Record 1989/019. Geoscience Australia, Canberra. URL: [http://pid.](http://pid.geoscience.gov.au/dataset/ga/14256)
903 [geoscience.gov.au/dataset/ga/14256](http://pid.geoscience.gov.au/dataset/ga/14256).

904 Yarmey, L., Baker, K.S., 2013. Towards standardization: A participatory framework for
905 scientific standard-making. *International Journal of Digital Curation* 8, 157–172. URL:
906 <https://doi.org/10.2218/ijdc.v8i1.252>, doi:[10.2218/ijdc.v8i1.252](https://doi.org/10.2218/ijdc.v8i1.252).

907 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/)
908 [//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/)
909 [code-deposition/](https://blogs.sciencemag.org/editors-blog/2021/07/15/progress-in-data-and-code-deposition/).

Supplementary Material for manuscript

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

M. Klöcking, L. Wyborn, K.A. Lehnert, B. Ware, A.M. Prent, L. Profeta, F. Kohlmann, W. Noble, I. Bruno, S. Lambart, H. Ananuer, N.D. Barber, H. Becker, M. Brodbeck, H. Deng, K. Deng, K. Elger, G. de Souza Franco, Y. Gao, K.M. Ghasera, D.C. Hezel, J. Huang, B. Kerswell, H. Koch, A.W. Lanati, G. ter Maat, N. Martínez-Villegas, L. Nana Yobo, A. Redaa, W. Schäfer, M.R. Swing, R.J.M. Taylor, M.K. Traun, J. Whelan, T. Zhou

S1. Goldschmidt 2022 Workshop

The workshop [“Earth Science meets Data Science: what are our needs for geochemical data, services and analytical capabilities in the 21st century?”](#) was held at the Goldschmidt Conference 2022 (hybrid format). The goals of this workshop were to:

1. Explore scientific challenges in geochemistry;
2. Showcase examples of existing data solutions/infrastructures/services; and
3. Discuss recommendations, best practices and essential features of a globally standardised geochemical data framework;

The primary focus of the workshop lay on exploring the data and infrastructure requirements for addressing future scientific challenges through keynote seminars, breakout working groups, panel and open discussions. In addition, two dedicated tutorials demonstrated the features and capabilities of the EarthChem, GEOROC (DIGIS) and AusGeochem (AuScope Geochemistry Network) data systems, inviting feedback from participants.

Contribution to this workshop gave participants the opportunity to voice their needs directly to the platform and repository creators/managers. Participants benefited from discussions around data transparency, best practices, big data synthesis and inter-laboratory analytical comparisons, actively contributing to bring about a cultural change in the geochemistry community. All workshop organisers and participants have been invited to contribute to this paper (see [Section S1.3](#) for a detailed list of contributors and their roles).

S1.1 Workshop Programme

Day 0 (8 July, 6:00-10:00 HADT): EarthChem & GEOROC tutorial	
06:00-06:15	Welcome & Introductions
06:15-06:45	Overview of Systems
06:45-08:15	Accessing Data
08:15-08:45	Coffee Break
08:45-09:00	Sample Management & Identification
09:00-10:00	Publishing Data
Day 1 (9 July, 10:00-17:00 HADT): Main Workshop	
10:00-10:15	Welcome & Introductions
10:15-11:00	Keynote 1: Sarah Lambart (slides)
11:00-11:15	Coffee Break
11:15-12:00	Keynote 2: Ian Bruno (slides)
12:00-12:45	Round-the-Table Discussion: expectations for this workshop
12:45-13:45	Lunch Break
13:45-14:45	Break-out Session 1: what are requirements for geochemical data?
14:45-15:15	Introduction to OneGeochemistry and the landscape of existing data resources and standards (Lesley Wyborn, slides)
15:15-15:45	Coffee Break
15:45-16:45	Break-out Session 2: how can we achieve these requirements?
16:45-17:00	Final Words & Way Forward

Day 2 (10 July, 12:00-17:00 HADT): AusGeochem tutorial	
12:00-13:00	Rooftop Lunch (in-person participants)
13:00 - 13:45	Introduction to Day 3 and hosts; What is the AGN, Goals of the AGN and The AusGeochem Platform
13:45 - 14:30	Technical Creation of the Platform
14:30 - 16:45	Exercise 1: AusGeochem Landing Page, Registration and Demonstration
	Exercise 2: Adding and Minting a Sample - Single sample upload to a created package - attach data to that sample - visualize data in map and interrogate data through various drop downs.
	Coffee Break
	Exercise 3: Adding Geochemical Data to Created Sample: Single sample upload to - created package - attach data to that sample - visualize data in map and interrogate data through various drop downs.
16:45 - 17:00	AusGeochem into the Future - Developments/Outlook

S1.2 Participating Data Systems

EarthChem is a disciplinary data facility established in 2003 that curates and provides open access to geochemical and petrological observations in digital data collections. Between 2010 and 2020, EarthChem has operated as part of the Interdisciplinary Earth Data Alliance (IEDA), and since 2022 within the reimagined IEDA2, both supported by the US National Science Foundation. EarthChem provides data stewardship services for a broad community of Earth and environmental scientists, including data publishing through the trusted repository service of the EarthChem Library, and data access and mining through synthesis databases (PetDB, EarthChem Portal, Library of Experimental Phase Relations, and the Decade Volcano Portal). EarthChem also established best practices for geochemical and sample-based data through the Editors Roundtable (Goldstein et al. 2014), which led to the Coalition for Publishing Data in the Earth & Space Sciences COPDESS (Hanson et al. 2015, Lehnert & Hsu 2015) that continues to advance best practices for data in scholarly publications and provided the foundation for the AGU project “Enabling FAIR Data” (Stall et al. 2017). EarthChem applications and data holdings can be accessed at <https://earthchem.org>.

The **GEOROC** (*Geochemistry of Rocks of the Oceans and Continents*) database contains published geochemical analyses of whole rocks, glasses, minerals and inclusions from eleven different geological settings across the world. It was set up by the Max Planck Institute for Chemistry, Mainz, Germany, and has been available online since the end of 1999. The database provides free access to >32 million individual values of major and trace element concentrations, radiogenic and nonradiogenic isotope ratios as well as analytical ages, compiled from >20,600 published scientific articles. In addition to the chemical analyses, extensive metadata describing the publication, sample and analytical method are stored. Together with PetDB and EarthChem, GEOROC developed a data and metadata schema for geochemical analyses (Lehnert et al., 2000). GEOROC is a key data contributor to the EarthChem Portal, hosted by IEDA2. Since 2021, GEOROC also provides a domain data repository that enables data submissions by the community. The GEOROC database is currently maintained by the Digital Geochemical Data Infrastructure ([DIGIS](#)) initiative at the University of Göttingen. It can be accessed at <https://georoc.eu>.

The **AusGeochem** platform was developed by the AuScope Geochemistry Network (AGN) and collaborator Lithodat to facilitate better organisation, coordination and ability to share data produced by Australian geochemistry laboratories. It differs from EarthChem and GEOROC in that it focuses on collecting data directly from the laboratories, although users can upload existing datasets. The AGN is funded by AuScope, Australia’s provider of infrastructure to the Earth and Geospatial Sciences in Australia, aiming to create wide and open access to earth and geospatial science infrastructure to drive research across government, institutions and industry. AuScope is funded by the Australian Government under the National Collaborative Infrastructure Strategy (NCRIS). As more and more data are produced in laboratories each day, the amount of data that becomes available for the scientific community through publications represents only the tip of the iceberg, with an appreciable amount of data abandoned on USB or hard disk drives. Therefore, the AGN aimed to build a platform that caters to laboratories, laboratory users and technical staff to make it easier to upload data directly from the instrument into a relational publicly accessible

database. When users upload sample metadata to AusGeochem, laboratory staff performing the analyses can upload the finished data directly into the platform, simultaneously linking the analyses to the sample metadata. With the ability to give samples unique codes and labels through an IGSN minting service, perform statistical analyses, novel capabilities to visualise and synthesise data within the context of large volumes of laboratory generated publicly funded geochemical data, the database simultaneously performs the function of repository and acts as a place for collaboration, interrogation and dissemination of high value geochemistry datasets. The platform links the private, collaborative and public domains. The AusGeochem platform can be accessed at <https://ausgeochem.auscope.org.au>.

S1.3 Workshop Contributors

Name	Affiliation	Workshop Role	Interest in this topic
Marthe Klöcking	Göttingen University; GEOROC	Organiser	
Kerstin Lehnert	Columbia University; IEDA2, EarthChem	Organiser	
Alexander Prent	Curtin University, AusGeochem	Organiser	
Lucia Profeta	Columbia University; IEDA2, EarthChem	Organiser	
Bryant Ware	Curtin University, AusGeochem	Organiser	
Fabian Kohlmann	Lithodat, AusGeochem	Tutorial organiser	
Wayne Noble	Lithodat, AusGeochem	Tutorial organiser	
Lesley Wyborn	Australian National University	Organiser; invited keynote speaker	Where to start on developing machine-actionable standards and vocabularies for geochemical data
Ian Bruno	CCDC	Participant; invited keynote speaker	To share experiences from chemistry and crystallography and learn about common challenges
Sarah Lambart	University of Utah	Participant; invited keynote speaker	Know how I can become a contributor - learn about new developments
Halimulati Ananuer	Macquarie University	Participant	Geochemical database and management
Michael Badawi	University of Lorraine	Participant	
Nicholas Barber	University of Cambridge	Participant	Interested in statistical treatments of large geochemical datasets
Harry Becker	Freie Universität Berlin	Participant	Metadata in data repositories

Maurice Brodbeck	Trinity College Dublin	Participant	What is good geochemical data? Reporting standards. Access and use of platforms.
Hang Deng	Peking University	Participant	Learn more about open access geochemical database
Kai Deng	ETH Zurich	Participant	Geochemical data compilation and reuse
Gabriel Franco	University of South Carolina	Participant	Getting and insider's perspective on the main geochemistry databases
Yajie Gao	Australian National University	Participant	Learn more about open database
Khalid Mohammed Ghasera	Aligarh Muslim University	Participant	Big data resources
Jingyi Huang	Johns Hopkins University	Participant	learn about how this data system works and how data managed
Buchanan Kerswell	Miami University	Participant	Student research
Jieun Kim	Northwestern University	Participant	Available geochemical databases
Hilde Koch	University College Dublin	Participant	Reliability of Geochemical Data of Database. Is it possible to get access to raw data? How can you make existing databases more visible?
Anthony Lanati	University Münster and Macquarie University	Participant	Ethical use of the databases, as well as contributing data and ensuring my data is FAIR
Nadia Martínez-Villegas	IPICYT	Participant	Learn about opendatabase and how to better report and manage data
Nicolas Randazzo	McMaster University	Participant	Interested in Big Data and wanted to learn more about data repositories
Ahmad Redaa	King Abdulaziz University	Participant	Learn about open access sources of geological datasets

Wiebke Schäfer	Friedrich-Alexander University	Participant	How to make my data accessible and pros and cons of making data accessible to everyone
Megan Swing	University of Toronto/ Royal Ontario Museum	Participant	Learning more about the the types of databases other researchers use
Marie Katrine Traun	University of Copenhagen	Participant	Integrating better data practices in data collection and laboratories and exploring the potential of big data geochemical studies with advanced statistics
Jo Whelan	Northern Territory Geological Survey	Participant	Learn about the data models, vocabularies and standards that are being applied with a view to ensure that the Survey is taking a consistent approach with data
Lucien Nana Yobo	Texas A&M University	Participant	Learn about the various geochemical database
Tengfei Zhou	China University of Petroleum (East China)	Participant	Big data and Geochemical data

References

- Bruno, I. (2022). Defining Standards in Crystallography and Chemistry. Zenodo. <https://doi.org/10.5281/zenodo.7120798>.
- Goldstein, S.L., Hofmann, A.W., & Lehnert, K.A. (2014). *Requirements for the Publication of Geochemical Data*. Interdisciplinary Earth Data Alliance (IEDA). <https://doi.org/10.1594/IEDA/100426>.
- Hanson, B., Lehnert, K.A., & Cutcher-Gershenfeld, J. (2015). Committing to publishing data in the Earth and space sciences. *EOS*, 96. <https://doi.org/10.1029/2015EO022207>.
- Lambart, S. (2022). Databases for petrologists: a user POV. Goldschmidt 2022, Honolulu, Hawaii. Zenodo. <https://doi.org/10.5281/zenodo.7221958>.
- Lehnert, K., & Hsu, L. (2015). The new paradigm of data publication. *Elements*, 11, 368-369.
- Lehnert, K., Su, Y., Langmuir, C. H., Sarbas, B., & Nohl, U. (2000), A global geochemical database structure for rocks, *Geochem. Geophys. Geosyst.*, 1, 1012, doi:[10.1029/1999GC000026](https://doi.org/10.1029/1999GC000026).
- Stall, S., Robinson, E., Wyborn, L., Yarmey, L.R., Parsons, M.A., Lehnert, L., Cutcher-Gershenfeld, J., Nosek, B., & Hanson, B. (2017). Enabling FAIR data across the Earth and space sciences. *Eos*, 98, <https://doi.org/10.1029/2017EO088425>.
- Wyborn, L., Lehnert, K., Prent, A., Klöcking, M., Profeta, L., Elger, K., & ter Maat, G. (2022). Introducing OneGeochemistry, landscape of data resources and standards: showcases growing list of publications on data reporting standards etc. Goldschmidt 2022, Honolulu, Hawai'i, USA. Zenodo. <https://doi.org/10.5281/zenodo.7212407>.