

Earth and Space Science Informatics Perspectives on Integrated, Coordinated, Open, Networked (ICON) Science

D.J. Hills^{1,2}, J.E. Damerow³, B. Ahmmed⁴, N. Catolico⁵, S. Chakraborty⁶, T.Y. Chen⁷, C.M. Coward⁸, R. Crystal-Ornelas³, W.D. Duncan³, L.N. Goparaju⁹, C. Lin¹⁰, Z. Liu^{4,11}, M. K. Mudunuru¹², Y. Rao¹³, R.J. Rovetto^{14,15}, Z. Sun¹¹, B.P. Whitehead¹⁶, L. Wyborn¹⁷, T. Yao^{6,18}

¹[Geological Survey of Alabama](#), Tuscaloosa, AL, USA.

²[Ronin Institute for Independent Scholarship](#), Tuscaloosa, AL, USA.

³[Lawrence Berkeley National Laboratory](#), Berkeley, CA, USA.

⁴[Los Alamos National Laboratory](#), Los Alamos, NM, USA.

⁵[Battelle, National Ecological Observatory Network](#), CO, USA.

⁶[NASA's Goddard Space Flight Center](#), Greenbelt, MD, USA.

⁷Academy for Mathematics, Science, and Engineering, Rockaway, NJ, USA

⁸[Jet Propulsion Laboratory](#), California Institute of Technology, Pasadena, CA, USA.

⁹Vindhyan Ecology and Natural History Foundation, U.P. India.

¹⁰Atkinson Center for Sustainability and Department of Information Science, [Cornell University](#), Ithaca, NY, USA.

¹¹[George Mason University](#), Fairfax, VA, USA.

¹²[Pacific Northwest National Laboratory](#), Richland, WA, USA.

¹³North Carolina Institute for Climate Studies, [North Carolina State University](#), Asheville, NC, USA.

¹⁴Center for Orbital Debris Education & Research, [University of Maryland](#), MD, USA.

¹⁵Independent, New York, NY, USA.

¹⁶[Manaaki Whenua – Landcare Research](#), Palmerston North, New Zealand

¹⁷[Australian National University](#), Canberra, ACT, Australia

¹⁸[Science Systems and Applications, Inc.](#), Lanham, MD, USA.

Corresponding authors: Denise Hills (denise.j.hills@gmail.com); Joan Damerow (JoanDamerow@lbl.gov)

Key Points:

- **Networks** across communities, with **Coordinated** data and information modeling practices, improve scientific outcomes for all involved.
- **Integrated, Coordinated, and Open** data requires sustainable support to create and maintain infrastructure for interdisciplinary **Networks**.

34
35
36

- **Integrated** and **Coordinated** use of data in machine learning calls for **Open** benchmark datasets, shared across **Networks** for improved outcomes.

37 **Abstract**

38 This article is composed of three independent commentaries about the state of ICON principles
39 (Goldman et al., 2021) in Earth and Space Science Informatics (ESSI) and includes discussion on
40 the opportunities and challenges of adopting them. Each commentary focuses on a different
41 topic: **(Section 2)** Global collaboration, cyberinfrastructure, and data sharing; **(Section 3)**
42 Machine learning and multiscale modeling; **(Section 4)** Remote sensing for advancing Earth
43 system model development by integrating field and ancillary data. ESSI addresses data
44 management practices, computation and analysis, and hardware and software infrastructure. Our
45 role in ICON science therefore involves collaborative work to assess, design, implement, and
46 promote practices and tools that enable effective data management, discovery, integration, and
47 reuse for interdisciplinary work in Earth and space science disciplines. Networks of diverse
48 people with expertise across Earth, space, and data science disciplines are essential for efficient
49 and ethical exchanges of FAIR research products and practices. Our challenge is then to
50 coordinate the development of standards, curation practices, and tools that enable integrating and
51 reusing multiple data types, software, multi-scale models, and machine learning approaches
52 across disciplines in a way that is as open and/or FAIR as ethically possible. This is a major
53 endeavor that could greatly increase the pace and potential of interdisciplinary scientific
54 discovery.

55 **Plain Language Summary**

56 We present commentaries on the state of “ICON principles” in Earth and Space Science
57 Informatics. ICON principles (Integrated, Coordinated, Open, and Networked) are meant to
58 improve the research experience for all. Ultimately, data standardized according to community
59 conventions and formats lead to more effective and efficient collaboration, data discovery,
60 integration, and analyses. Data standards, tools, and machine learning developed using ICON
61 principles enhance our understanding of Earth processes. Using ICON principles improves
62 model results and efficacy, fosters interdisciplinary research, and provides a framework by which
63 non-experts can confidently contribute volunteered data and findings. Standardized data also
64 provides reliable common resources to help train and benchmark machine learning algorithms.
65 When networked communities work together to standardize and share data openly, the resulting
66 web of research products is more readily findable, accessible, interoperable, and reusable
67 (FAIR). Ongoing support is crucial to develop and sustain the people, systems, and tools
68 necessary to realize ICON principles in Earth and Space Science Informatics now and in the
69 future.

70 **1 Introduction**

71 Integrated, Coordinated, Open, Networked (ICON) science aims to enhance synthesis,
72 increase resource efficiency, and create transferable knowledge (Goldman et al., 2021). This
73 article belongs to a collection of commentaries ([Leadership Team et al., 2021](#)) spanning
74 geoscience on the state and future of ICON science. Earth and Space Science Informatics (ESSI)
75 encompasses a broad field that addresses data management practices, computation and analysis,
76 and hardware and software infrastructure. ESSI’s role in ICON science therefore involves
77 collaborative work to assess, design, implement, and promote practices and tools that enable
78 effective data management, discovery, integration, and reuse for interdisciplinary work in Earth
79 and space science (ESS) disciplines. In this series of commentaries, we examine the current state,
80 challenges, and opportunities of ICON science through the lenses of global collaboration,

81 [cyberinfrastructure](#), and data sharing (Section 2); machine learning and multiscale modeling
82 (Section 3); and remote sensing for advancing Earth system models (ESM) development by
83 integrating field and ancillary data (Section 4).

84 **2 Global collaboration, cyberinfrastructure, and data sharing**

85 2.1 Current state and challenges

86 Global collaboration across disciplines is essential to the development and
87 implementation of data/metadata standards and cyberinfrastructures. Thus, many organizations
88 have emerged to facilitate such collaboration, e.g., [Research Data Alliance](#), [World Data System](#),
89 [OneGeology](#), [Earth Science Information Partners](#). These organizations have produced numerous
90 active [groups involved in Earth, space and environmental science data and research](#), and
91 developed many data tools and services, e.g. [Earth, Space and Environmental Sciences Data](#)
92 [Vocabulary Repositories](#). Research is more efficient with **Networked** data practices and
93 cyberinfrastructures that support scientific discovery. Yet, there is still a large disconnect and
94 lack of **Coordination** across many informatics communities and the broader communities we
95 aim to support.

96 Research teams often lack sufficient resources (e.g., appropriate cyberinfrastructure,
97 expert data/software personnel, financial allotment) to effectively manage, standardize, and
98 publish high-quality data (Mons, 2020). This hinders data from being **Open and/or Findable**,
99 **Accessible, Interoperable, and Reusable** (FAIR; Wilkinson et al., 2016). Further, specific
100 criteria to make data FAIR (Gries et al., 2019; Jones et al., 2019) inevitably vary across
101 disciplines and data types. Because there are no widely accepted standards to evaluate FAIR-
102 ness, data may be miscaterogized (e.g., Kinkade & Shepherd, 2021; Mons et al., 2017; Stall et
103 al., 2019). Importantly, FAIR does not mean **Open**; data can be **Open** without being FAIR, and
104 *vice versa* (see [What is the difference between “FAIR data” and “Open data” if there is one?](#)).

105 Supporting ESS research requires assessing, designing, building, and maintaining
106 cyberinfrastructures (e.g., data repositories/archives, application programming interfaces (APIs),
107 visualization tools, search interfaces) that are often organized around a particular data type,
108 discipline, or organization. Interoperability issues are then minimized using bespoke or *ad hoc*
109 conventions within that particular community (e.g., [Deep Carbon Observatory](#), [HydroShare](#),
110 [Long-Term Ecological Research Network](#), [National Ecological Observatory Network](#)). However,
111 most cyberinfrastructures lack the resources for **Integration** and **Coordination** necessary for
112 interdisciplinary work, including guidance and leading practices; domain semantics; technical,
113 data, methodological, and instrumentation standards; workflow management; training; and
114 sustainable technical and financial support. These deficits hinder **Open** data that fosters machine
115 actionable, interdisciplinary scientific discovery.

116 While existing standards and practices may address similar concepts, they are not fully
117 interoperable or **Integrated** within and across relevant disciplines. Valuable resources are spent
118 developing/updating translators, or disciplinary standards are simply disconnected and inefficient
119 for interdisciplinary users. **Coordination** is needed to implement standards for effective
120 interdisciplinary data discovery and exchange. A major limitation to **Coordination** involves a
121 lack of consistent and transparent protocols (e.g., data and code production, processing methods)
122 across interdisciplinary teams that limits reuse and replication. These combined factors create
123 barriers to **Open and FAIR** data.

124 Ever-increasing volumes of open data and tools now allow us to ask science questions
125 that synthesize data and knowledge across scientific disciplines from globally distributed
126 resources, thus expanding the impact of funded research (e.g., Michener, 2015; Rosenberg et al.,
127 2019). More successful **Networked** data sharing efforts (e.g., [Global Biodiversity Information](#)
128 [Facility](#), [Ameriflux](#), [Consortium of Universities for the Advancement of Hydrologic Science,](#)
129 [Inc.](#)) have been driven by 1) demand for a specific data type (Barrett et al., 2012; Novick et al.,
130 2018; Robertson et al., 2014); 2) reporting standards that enable global data search and
131 integration (e.g., Wieczorek et al., 2012; Yilmaz et al., 2011); and 3) associated user-friendly
132 tools (Clark et al., 2016; Robertson et al., 2014).

133 2.2. Opportunities and moving forward

134 Replicable and transparent research that reflects ICON principles requires sustainable
135 investment in cyberinfrastructure to improve interoperability and **Integration**. Global high-level
136 **Coordination** across organizations is needed to bridge siloed efforts across disciplines,
137 organizations, and/or countries. A commitment to community engagement is needed to bring
138 together input across disciplines, understand data management challenges and needs, and
139 promote the adoption of shared practices. Making data as **Open and/or FAIR** as ethically
140 possible requires key advocates who facilitate **Networked** collaboration.

141 Data users, code contributors, and tool developers should align with established standards
142 or community practices. We can encourage practices that promote ICON principles, such as
143 **Open** publication of study plans (e.g., [PLOS ONE study proposals](#)), data production and
144 processing protocols (e.g., [Common Workflow Language](#)), and software code. We must
145 continually evaluate how to **Coordinate** and **Integrate** across existing cyberinfrastructure from
146 local to global scales, which involves iterative rounds of engagement; education and outreach;
147 and feedback across data providers, tool and service creators, and scientists who use ESS data
148 and services. **Coordinating Networks** across disciplines will involve technical approaches to
149 connect related data (e.g., PIDs, APIs, ontologies, geospatial standards) and promoting
150 widespread adoption of community standards that improve scientific outcomes and benefit all
151 participants in the network.

152 3 Machine learning and multiscale modeling

153 3.1 Success and current status of AI/ML

154 Over the past decade, artificial intelligence approaches, including machine learning
155 (AI/ML), have revolutionized scientific discovery across disciplines, including ESS (Maskey et
156 al., 2020a). The AI/ML revolution, driven by a wealth of **Open** data and rapid technological
157 development in computational cyberinfrastructure, has led to more processing power and greater
158 **Networking** which allows unprecedented resource and data sharing. There are many success
159 stories demonstrating how AI/ML has been used to address challenging issues in ESS, e.g.,
160 extreme weather prediction (Maskey et al., 2020b; Pradhan et al., 2018; Wimmers et al., 2019),
161 land use/land cover change monitoring (Hansen et al., 2013), earth system modeling (Reichstein
162 et al., 2019), endangered species identification (Allen et al., 2021), spatial downscaling of
163 climate models and satellite observations (López López et al., 2018; Vandal et al., 2019), space
164 weather forecasting (Wintoft et al., 2017), and lunar and planetary landform classification
165 (Palafox et al., 2017; Silburt et al., 2019). Various funding agencies worldwide have recently

166 released their strategic plans and guidelines to expand the investment in AI/ML research which
167 will further its adoption within ESSI for at least the next decade.

168 3.2 Common challenges in AI/ML

169 To accelerate this adoption, the ESS community needs to collectively address three key
170 challenges. First, most AI/ML applications in ESS are *ad hoc* research that lacks system-wide
171 **Coordination** and is time-consuming. There are little AI-ready data (e.g., cleaned, harmonized,
172 formatted, well understood) that can be efficiently **Integrated** across domains or applications
173 and few recommended practices on proper model development and documentation (Maskey et
174 al., 2020a). Thus, amplifying the value of AI/ML in ESS requires an ecosystem including AI-
175 ready training datasets and standardized model development practices. This ecosystem would
176 enable the ESS community to collaboratively develop open AI/ML applications at scale. A
177 second challenge is related to the wealth of **Open** data in ESS. Currently, there are no
178 community-recommended practices on how to properly develop, document, and share the AI/ML
179 applications that track provenance and enable reproducibility (Sun et al., 2020). Third, the
180 explainability and generalizability of AI/ML models are also major concerns for the ESS
181 community (McGovern et al., 2019; Toms et al., 2020). To address complex questions in ESS
182 systems, we need to better understand why AI/ML models perform in a certain way, their
183 consistency with domain knowledge, and how models developed using a specific set of data can
184 adjust dynamically to shifts in ESS data. Additionally, ethical awareness, conduct, and
185 responsibility in AI/ML and related activities are essential to the practice of principled research.

186 3.3 Opportunities and moving forward

187 We identify five opportunities where researchers may focus their efforts to make ESS
188 AI/ML more efficient. One opportunity relates to big data in ESS. Because the capacity and
189 application scope of AI/ML heavily depends on patterns in training data, it should be as
190 representative as possible. The requirements for big training datasets have led to calls for
191 libraries of **Open** and FAIR benchmark datasets ([WILDS](#), Koh et al., 2020; [Radiant Earth](#)
192 [Foundation](#); Rasp et al., 2020) related to questions within ESS (Crystal-Ornelas et al., 2021). A
193 second opportunity is increased **Networking** through cloud computing (Gorelick et al., 2017;
194 Mayer-Schönberger & Cukier, 2013). By sharing data and models in the cloud, researchers
195 around the world can access these resources without being limited by local computing power.
196 More work needs to be done to make cloud computing more accessible for ESS despite recent
197 progress. Increased **Openness** in the exchange of data handling practices to allow sharing
198 common workflows while handling large datasets is a third opportunity. A fourth opportunity is
199 to improve interpretability through **Integration** across disciplines by: (1) including physics in
200 ML models (Jia et al., 2019; Raissi et al., 2019), (2) leveraging machine learning exploratory
201 tools (Montavon et al., 2017; Ying et al., 2019), and (3) involving domain experts into AI/ML
202 pipelines. A final opportunity for growth is to automate workflows to improve the development
203 efficiency (e.g., auto-sklearn, AutoKeras) (He et al., 2021). To improve AI engineering
204 efficiency and reduce data collection and processing costs, modelers may also use data
205 augmentation methods such as mixup (Zhang et al., 2017) to fill in the missing data and enhance
206 data quality (Alexandrov & Vesselinov, 2014; Vesselinov et al., 2018). We emphasize that these
207 opportunities for ESS to inform and apply AI/ML models is not exhaustive; rather it is a starting
208 point for exploring how ICON science can benefit the future of this rapidly growing field within
209 ESS.

210 **4 Remote sensing for advancing Earth system model development by integrating field and** 211 **ancillary data**

212 4.1 Current Status

213 Remote sensing technology combined with field and ancillary data (e.g., field
214 measurements, other imagery; Acton, 1996) has transformed the development of ESMs as they
215 have advanced from aerial imagery of the early nineteenth century (Necsoiu et al., 2013) to the
216 present-day's Google Earth Engine (Gorelick et al., 2017) and Unmanned Aerial Vehicles (Singh
217 & Frazier, 2018). Most publicly-funded remote sensing datasets are **Open** and hosted on public
218 repositories (e.g., government-sponsored repositories, Github, Zenodo). In addition, this data is
219 collected through **Coordinated** standards between government agencies across the globe
220 (Alameh, 2020). **Integration** of remote sensing technology with independent field measurements
221 and high spatial resolution satellite imagery has been essential for ESM validation. This also
222 includes estimating derived data products (e.g., from satellites) accuracy and quantifying
223 uncertainty (Strahler et al., 2006). Crowdsourcing and citizen science have further advanced the
224 integration of remote sensing with field data (e.g., [RaspberryShake](#), Khan et al., 2018; Saralioglu
225 & Gungor, 2020; Worldwide Hydrobiogeochemistry Observation Network for Dynamic River
226 Systems [[WHONDRS](#)], Stegen & Goldman, 2018), resulting in broader **Networked** efforts that
227 benefit researchers and a wide variety of data users. We note that agencies in the US and Europe
228 have open-sourced their data to all users internationally. Some popular open data sources,
229 associated cyberinfrastructure, and tools are included in [an associated github repository](#).

230 4.2 Challenges and call to action

231 Two primary challenges which the ESSI community faces are limited global data
232 collection and inadequate cyberinfrastructure. Despite advances in sensors, crowdsourcing, and
233 citizen science (e.g., RaspberryShake, WHONDRS), collecting and hosting high-quality global
234 data present immense challenges. For example, RaspberryShake has collected more than 30TB
235 of seismographic data over the past decade but lacks the necessary cyberinfrastructure to reliably
236 and sustainably store it.

237 Recent progress in AI/ML has improved the representation of Earth system processes
238 (e.g., thermal, land physics and hydrology, radiation, atmospheric ocean circulation) in ESMs
239 (Rasp et al., 2018). ML, in particular, requires massive datasets to represent processes at both
240 normal and extreme events (e.g., hurricanes, wildfires); however, extreme event data are rare due
241 to the unique challenges faced during collection. Thus, the concept of crowdsourcing data
242 collection, using **Coordinated** methods (e.g., RaspberryShake, WHONDRS) on extreme events,
243 is an attractive option that improves **Networked** research.

244 There has been a **Coordinated** effort from US and European agencies to develop
245 cyberinfrastructure that improves and increases access to data to enhance predictions and
246 understanding of various Earth system processes. For example, the European Space Agency
247 Sentinel data products are recently available in the [Copernicus Data and Information Access](#)
248 [Service](#) cloud environments. In addition, the US Geological Survey Landsat satellite data
249 inventory has been open to the public since 2008 and has been in the cloud since 2020 (U.S.
250 Geological Survey, 2008). Furthermore, the National Aeronautics and Space Administration
251 (NASA) and the National Oceanic and Atmospheric Administration (NOAA) have adopted a
252 strategic vision to leverage cloud computing and operate multiple components of their data

253 systems in a retail cloud environment. This calls for action to identify the opportunities to
254 improve policy and strategy planning across various countries to make satellite data products
255 accessible to all users in open data portals. In addition, automated quality assurance of satellite
256 observations is needed to support global, regional, or local data services. **Coordinated** across
257 international agencies, a standard open data cyberinfrastructure will help to assure ESM data
258 from multiple sources (national, regional, governments, academia, and the private sector) are
259 available and easily **Integrated** into open-source platforms and networks.

260 4.3 Opportunities and moving forward

261 First, close coordination would help international agencies and organizations build a
262 standard open data cyberinfrastructure to ensure that earth science data are free, open, and easily
263 integrated into ESMs. Second, we need next-generation sensors and satellites which provide
264 more fine resolution data to increase the accuracy of ESMs. For example, the joint NASA-Indian
265 Space Research Organization (ISRO) Synthetic Aperture Radar (SAR) ([NISAR](#)) mission is
266 anticipated to provide fine-scale resolution radar data with a spatial resolution of less than a
267 centimeter to study the earth's features and processes. Third, the role of AI/ML needs to be
268 expanded to plug in the gaps of remote sensing data.

269 5 Concluding remarks

270 ESSI science that utilizes ICON principles enables data synthesis, increases resource
271 efficiency, and creates knowledge that transcends individual systems (Goldman et al., 2021).
272 ESSI can work to ensure that diverse scientists have user-friendly resources to contribute
273 and use data that follows community conventions. Such collections of **Open and/or FAIR** data,
274 shared across **Networks** for mutual benefit, are critical to appropriately train AI/ML, which
275 furthers **Integration** and **Coordination** in ESSI science. Cross-community **Networks** improve
276 scientific outcomes for all involved. Communities must work together to share data openly using
277 community standards, to produce **Open and/or FAIR** data that enables data synthesis and can
278 revolutionize fields of research (e.g., Kelling et al., 2009). Ongoing, sustainable support is vital
279 to create and maintain the cyberinfrastructure and human resources necessary for **Integrated**,
280 **Coordinated**, and **Open and/or FAIR** data (as much ethically as possible) for interdisciplinary
281 **Networks**.

282 Acknowledgments

283 DJH, JED, NC, CC, WDD, ZL, RJR, BPW, and LW authored section 2 'Global
284 collaboration, cyberinfrastructure, and data sharing.' RCO, SC, BA, CL, YR, TYC, and ZS
285 authored section 3 'Machine learning and multiscale modeling.' LNG, MKM, and TY authored
286 section 4 'Remote sensing for advancing Earth system model development by integrating field
287 and ancillary data.'

288 Sky Bristol (USGS) was instrumental in early discussions, particularly of cost-benefit
289 analysis.

290 JED and RCO were funded by the ESS-DIVE repository and WDD by the National
291 Microbiome Data Collaborative, both by the U.S. DOE's Office of Science Biological and
292 Environmental Research under contract number DE-AC02-05CH11231. NC was supported by
293 NEON, a program sponsored by the NSF and operated under cooperative agreement by Battelle

294 Memorial Institute. SC was supported by an appointment to the NASA Postdoctoral Program at
295 NASA Goddard Space Flight Center, administered by Universities Space Research Association
296 under contract with NASA. CL was supported by an appointment as a postdoctoral fellow at the
297 Cornell Atkinson Center for Sustainability, and an affiliation with the Department of Information
298 Science. ZL was supported by Cooperative Geoinformation Research with the NASA GSFC
299 Earth Sciences Data and Information Services Center (GES DISC). MKM was supported by the
300 U.S. DOE-SC, SFA at PNNL. YR was supported by NOAA through the Cooperative Institute for
301 Satellite Earth System Studies under Cooperative Agreement NA19NES4320002. BPW was
302 supported by the Ministry of Business Innovation and Employment (MBIE) Infrastructure
303 Platform. TY was supported by [NASA Applied Sciences Disasters Program](#) and NASA's
304 [LANCE](#) system, part of NASA's EOSDIS.

305 The views and opinions of authors expressed herein do not necessarily state or reflect
306 those of the US Government or any international agency thereof.

307 **References**

- 308 Acton, C. H. (1996). Ancillary data services of NASA's Navigation and Ancillary Information
309 Facility. *Planetary and Space Science*, 44(1), 65–70. [https://doi.org/10.1016/0032-](https://doi.org/10.1016/0032-0633(95)00107-7)
310 [0633\(95\)00107-7](https://doi.org/10.1016/0032-0633(95)00107-7)
- 311 Alameh, N. (2020). A future of location data integration. *Geo: GeoConnexion International*
312 *Magazine*, 19(6), 18–19. Retrieved from [https://www.geoconnexion.com/publication-articles/a-](https://www.geoconnexion.com/publication-articles/a-future-of-location-data-integration)
313 [future-of-location-data-integration](https://www.geoconnexion.com/publication-articles/a-future-of-location-data-integration)
- 314 Alexandrov, B. S., & Vesselinov, V. V. (2014). Blind source separation for groundwater
315 pressure analysis based on nonnegative matrix factorization. *Water Resources Research*, 50(9),
316 7332–7347. <https://doi.org/10.1002/2013wr015037>
- 317 Allen, A. N., Harvey, M., Harrell, L., Jansen, A., Merkens, K. P., Wall, C. C., et al. (2021). A
318 Convolutional Neural Network for Automated Detection of Humpback Whale Song in a Diverse,
319 Long-Term Passive Acoustic Dataset. *Frontiers in Marine Science*, 8, 165.
320 <https://doi.org/10.3389/fmars.2021.607321>
- 321 Barrett, T., Clark, K., Gevorgyan, R., Gorelenkov, V., Gribov, E., Karsch-Mizrachi, I., et al.
322 (2012). BioProject and BioSample databases at NCBI: facilitating capture and organization of
323 metadata. *Nucleic Acids Research*, 40(D1), D57–D63. <https://doi.org/10.1093/nar/gkr1163>

324 Clark, K., Karsch-Mizrachi, I., Lipman, D. J., Ostell, J., & Sayers, E. W. (2016). GenBank.
325 *Nucleic Acids Research*, 44(D1), D67–D72. <https://doi.org/10.1093/nar/gkv1276>

326 Crystal-Ornelas, R., Varadharajan, C., Christianson, D., Damerow, J., Weierbach, H., Robles, E.,
327 et al. (2021). *A library of AI-assisted FAIR water cycle and related disturbance datasets to*
328 *enable model training, parameterization and validation*. Office of Scientific and Technical
329 Information (OSTI). <https://doi.org/10.2172/1769646>

330 Goldman, A. E., Emani, S. R., Pérez-Angel, L. C., Rodríguez-Ramos, J. A., Stegen, J. C., & Fox,
331 P. (2021). Special collection on open collaboration across geosciences. *Eos*, 102.
332 <https://doi.org/10.1029/2021EO153180>

333 Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., & Moore, R. (2017). Google
334 Earth Engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment*,
335 202, 18–27. <https://doi.org/10.1016/j.rse.2017.06.031>

336 Gries, C., Servilla, M., O’Brien, M., Vanderbilt, K., Smith, C., Costa, D., & Grossman-Clarke, S.
337 (2019). Achieving FAIR Data Principles at the Environmental Data Initiative, the US-LTER
338 Data Repository. *Biodiversity Information Science and Standards*, 3, e37047.
339 <https://doi.org/10.3897/biss.3.37047>

340 Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., et al.
341 (2013). High-resolution global maps of 21st-century forest cover change. *Science*, 342(6160),
342 850–853. <https://doi.org/10.1126/science.1244693>

343 He, X., Zhao, K., & Chu, X. (2021). AutoML: A Survey of the state-of-the-art. *Knowledge-*
344 *Based Systems*, 212, 106622. <https://doi.org/10.1016/j.knosys.2020.106622>

345 Jia, X., Willard, J., Karpatne, A., Read, J., Zwart, J., Steinbach, M., & Kumar, V. (2019). Physics
346 Guided RNNs for Modeling Dynamical Systems: A Case Study in Simulating Lake Temperature

- 347 Profiles. In *Proceedings of the 2019 SIAM International Conference on Data Mining (SDM)* (pp.
348 558–566). Society for Industrial and Applied Mathematics.
349 <https://doi.org/10.1137/1.9781611975673.63>
- 350 Jones, M. B., Slaughter, P., & Habermann, T. (2019). *Quantifying FAIR: automated metadata*
351 *improvement and guidance in the DataONE repository network*.
352 <https://doi.org/10.5281/zenodo.3408466>
- 353 Kelling, S., Hochachka, W. M., Fink, D., Riedewald, M., Caruana, R., Ballard, G., & Hooker, G.
354 (2009). Data-intensive Science: A New Paradigm for Biodiversity Studies. *Bioscience*, 59(7),
355 613–620. <https://doi.org/10.1525/bio.2009.59.7.12>
- 356 Khan, A., Denton, P., Stevenson, J., & Bossu, R. (2018). Engaging citizen seismologists
357 worldwide. *Astronomy & Geophysics*, 59(4), 4.15–4.18. <https://doi.org/10.1093/astrogeo/aty190>
- 358 Kinkade, D., & Shepherd, A. (2021). Geoscience data publication: Practices and perspectives on
359 enabling the FAIR guiding principles. *Geoscience Data Journal*, (gdj3.120).
360 <https://doi.org/10.1002/gdj3.120>
- 361 Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., et al. (2020).
362 WILDS: A Benchmark of in-the-Wild Distribution Shifts. *arXiv [cs.LG]*. Retrieved from
363 <https://arxiv.org/abs/2012.07421>
- 364 López López, P., Immerzeel, W. W., Rodríguez Sandoval, E. A., Sterk, G., & Schellekens, J.
365 (2018). Spatial Downscaling of Satellite-Based Precipitation and Its Impact on Discharge
366 Simulations in the Magdalena River Basin in Colombia. *Frontiers of Earth Science in China*, 6,
367 68. <https://doi.org/10.3389/feart.2018.00068>
- 368 Maskey, M., Alemohammad, H., Murphy, K. J., & Ramachandran, R. (2020a). Advancing AI for
369 Earth Science: A data systems perspective. *Eos*, 101. <https://doi.org/10.1029/2020EO151245>

- 370 Maskey, M., Ramachandran, R., Ramasubramanian, M., Gurung, I., Freitag, B., Kaulfus, A., et
371 al. (2020b). Deepti: Deep-Learning-Based Tropical Cyclone Intensity Estimation System. *IEEE*
372 *Journal of Selected Topics in Applied Earth Observations and Remote Sensing*.
373 <https://doi.org/10.1109/jstars.2020.3011907>
- 374 Mayer-Schönberger, V., & Cukier, K. (2013). *Big Data: A Revolution that Will Transform how*
375 *We Live, Work, and Think*. Houghton Mifflin Harcourt. Retrieved from
376 <https://play.google.com/store/books/details?id=uy4lh-WEhhIC>
- 377 McGovern, A., Lagerquist, R., Gagne, D. J., Eli Jergensen, G., Elmore, K. L., Homeyer, C. R., &
378 Smith, T. (2019). Making the Black Box More Transparent: Understanding the Physical
379 Implications of Machine Learning. *Bulletin of the American Meteorological Society*, 100(11),
380 2175–2199. <https://doi.org/10.1175/BAMS-D-18-0195.1>
- 381 Michener, W. K. (2015). Ecological data sharing. *Ecological Informatics*, 29, 33–44.
382 <https://doi.org/10.1016/j.ecoinf.2015.06.010>
- 383 Mons, B. (2020). Invest 5% of research funds in ensuring data are reusable. *Nature*, 578(7796),
384 491–491. <https://doi.org/10.1038/d41586-020-00505-7>
- 385 Mons, B., Neylon, C., Velterop, J., Dumontier, M., da Silva Santos, L. O. B., & Wilkinson, M.
386 D. (2017). Cloudy, increasingly FAIR; revisiting the FAIR Data guiding principles for the
387 European Open Science Cloud. *Information Services & Use*, 37(1), 49–56.
388 <https://doi.org/10.3233/isu-170824>
- 389 Montavon, G., Lapuschkin, S., Binder, A., Samek, W., & Müller, K.-R. (2017). Explaining
390 nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition*, 65,
391 211–222. <https://doi.org/10.1016/j.patcog.2016.11.008>

- 392 Necsoiu, M., Dinwiddie, C. L., Walter, G. R., Larsen, A., & Stothoff, S. A. (2013). Multi-
393 temporal image analysis of historical aerial photographs and recent satellite imagery reveals
394 evolution of water body surface area and polygonal terrain morphology in Kobuk Valley
395 National Park, Alaska. *Environmental Research Letters*, 8(2), 025007.
396 <https://doi.org/10.1088/1748-9326/8/2/025007>
- 397 Novick, K. A., Biederman, J. A., Desai, A. R., Litvak, M. E., Moore, D. J. P., Scott, R. L., &
398 Torn, M. S. (2018). The AmeriFlux network: A coalition of the willing. *Agricultural and Forest*
399 *Meteorology*, 249, 444–456. <https://doi.org/10.1016/j.agrformet.2017.10.009>
- 400 Palafox, L. F., Hamilton, C. W., Scheidt, S. P., & Alvarez, A. M. (2017). Automated detection of
401 geological landforms on Mars using Convolutional Neural Networks. *Computers & Geosciences*,
402 101, 48–56. <https://doi.org/10.1016/j.cageo.2016.12.015>
- 403 Pradhan, R., Aygun, R. S., Maskey, M., Ramachandran, R., & Cecil, D. J. (2018). Tropical
404 Cyclone Intensity Estimation Using a Deep Convolutional Neural Network. *IEEE Transactions*
405 *on Image Processing*. <https://doi.org/10.1109/tip.2017.2766358>
- 406 Raissi, M., Perdikaris, P., & Karniadakis, G. E. (2019). Physics-informed neural networks: A
407 deep learning framework for solving forward and inverse problems involving nonlinear partial
408 differential equations. *Journal of Computational Physics*, 378, 686–707.
409 <https://doi.org/10.1016/j.jcp.2018.10.045>
- 410 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in
411 climate models. *Proceedings of the National Academy of Sciences of the United States of*
412 *America*, 115(39), 9684–9689. <https://doi.org/10.1073/pnas.1810286115>

- 413 Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., & Thuerey, N. (2020).
414 WeatherBench: A benchmark data set for data-driven weather forecasting. *Journal of Advances*
415 *in Modeling Earth Systems*, 12(11). <https://doi.org/10.1029/2020ms002203>
- 416 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat.
417 (2019). Deep learning and process understanding for data-driven Earth system science. *Nature*,
418 566(7743), 195–204. <https://doi.org/10.1038/s41586-019-0912-1>
- 419 Robertson, T., Döring, M., Guralnick, R., Bloom, D., Wieczorek, J., Braak, K., et al. (2014). The
420 GBIF integrated publishing toolkit: facilitating the efficient publishing of biodiversity data on the
421 internet. *PloS One*, 9(8), e102623. <https://doi.org/10.1371/journal.pone.0102623>
- 422 Rosenberg, K. V., Dokter, A. M., Blancher, P. J., Sauer, J. R., Smith, A. C., Smith, P. A., et al.
423 (2019). Decline of the North American avifauna. *Science*, eaaw1313.
424 <https://doi.org/10.1126/science.aaw1313>
- 425 Saralioglu, E., & Gungor, O. (2020). Crowdsourcing in Remote Sensing: A Review of
426 Applications and Future Directions. *IEEE Geoscience and Remote Sensing Magazine*, 8(4), 89–
427 110. <https://doi.org/10.1109/MGRS.2020.2975132>
- 428 Silburt, A., Ali-Dib, M., Zhu, C., Jackson, A., Valencia, D., Kissin, Y., et al. (2019). Lunar crater
429 identification via deep learning. *Icarus*, 317, 27–38. <https://doi.org/10.1016/j.icarus.2018.06.022>
- 430 Singh, K. K., & Frazier, A. E. (2018). A meta-analysis and review of unmanned aircraft system
431 (UAS) imagery for terrestrial applications. *International Journal of Remote Sensing*, 39(15-16),
432 5078–5098. <https://doi.org/10.1080/01431161.2017.1420941>
- 433 Strahler, A. H., Boschetti, L., Foody, G. M., Friedl, M. A., Hansen, M. C., Herold, M., et al.
434 (2006). *Global land cover validation: Recommendations for evaluation and accuracy assessment*
435 *of global land cover maps* (Publication EUR 22156 EN). European Commission, Joint Research

436 Center. Retrieved from [https://op.europa.eu/en/publication-detail/-/publication/52730469-6bc9-](https://op.europa.eu/en/publication-detail/-/publication/52730469-6bc9-47a9-b486-5e2662629976)
437 [47a9-b486-5e2662629976](https://op.europa.eu/en/publication-detail/-/publication/52730469-6bc9-47a9-b486-5e2662629976)

438 Stall, S., Yarmey, L., Cutcher-Gershenfeld, J., Hanson, B., Lehnert, K., Nosek, B., et al. (2019).
439 Make scientific data FAIR. *Nature*, 570(7759), 27. <https://doi.org/10.1038/d41586-019-01720-7>

440 Stegen, J. C., & Goldman, A. E. (2018). WHONDRS: a Community Resource for Studying
441 Dynamic River Corridors. *mSystems*, 3(5), e00151–18. [https://doi.org/10.1128/mSystems.00151-](https://doi.org/10.1128/mSystems.00151-18)
442 [18](https://doi.org/10.1128/mSystems.00151-18)

443 Sun, Z., Di, L., Burgess, A., Tullis, J. A., & Magill, A. B. (2020). Geoweaver: Advanced
444 Cyberinfrastructure for Managing Hybrid Geoscientific AI Workflows. *ISPRS International*
445 *Journal of Geo-Information*, 9(2), 119. <https://doi.org/10.3390/ijgi9020119>

446 Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. (2020). Physically interpretable neural networks
447 for the geosciences: Applications to earth system variability. *Journal of Advances in Modeling*
448 *Earth Systems*, 12(9). <https://doi.org/10.1029/2019ms002002>

449 U.S. Geological Survey. (2008). *Imagery for Everyone: Timeline Set to Release Entire USGS*
450 *Landsat Archive at No Charge*. Retrieved from [https://prd-wret.s3.us-west-](https://prd-wret.s3.us-west-2.amazonaws.com/assets/palladium/production/s3fs-public/atoms/files/USGSStechann-20080421-landsat-imagery-release.pdf)
451 [2.amazonaws.com/assets/palladium/production/s3fs-public/atoms/files/USGSStechann-20080421-](https://prd-wret.s3.us-west-2.amazonaws.com/assets/palladium/production/s3fs-public/atoms/files/USGSStechann-20080421-landsat-imagery-release.pdf)
452 [landsat-imagery-release.pdf](https://prd-wret.s3.us-west-2.amazonaws.com/assets/palladium/production/s3fs-public/atoms/files/USGSStechann-20080421-landsat-imagery-release.pdf)

453 Vandal, T., Kodra, E., & Ganguly, A. R. (2019). Intercomparison of machine learning methods
454 for statistical downscaling: the case of daily and extreme precipitation. *Theoretical and Applied*
455 *Climatology*, 137(1-2), 557–570. <https://doi.org/10.1007/s00704-018-2613-3>

456 Vesselinov, V. V., Alexandrov, B. S., & O'Malley, D. (2018). Contaminant source identification
457 using semi-supervised machine learning. *Journal of Contaminant Hydrology*, 212, 134–142.
458 <https://doi.org/10.1016/j.jconhyd.2017.11.002>

- 459 Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., Giovanni, R., et al. (2012).
460 Darwin Core: An Evolving Community-Developed Biodiversity Data Standard. *PloS One*, 7(1),
461 e29715. <https://doi.org/10.1371/journal.pone.0029715>
- 462 Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., et al.
463 (2016). The FAIR Guiding Principles for scientific data management and stewardship. *Scientific*
464 *Data*, 3, 160018. <https://doi.org/10.1038/sdata.2016.18>
- 465 Wimmers, A., Velden, C., & Cossuth, J. H. (2019). Using Deep Learning to Estimate Tropical
466 Cyclone Intensity from Satellite Passive Microwave Imagery. *Monthly Weather Review*, 147(6),
467 2261–2282. <https://doi.org/10.1175/MWR-D-18-0391.1>
- 468 Wintoft, P., Wik, M., Matzka, J., & Shprits, Y. (2017). Forecasting Kp from solar wind data:
469 input parameter study using 3-hour averages and 3-hour range values. *Journal of Space Weather*
470 *and Space Climate*, 7, A29. <https://doi.org/10.1051/swsc/2017027>
- 471 Yilmaz, P., Kottmann, R., Field, D., Knight, R., Cole, J. R., Amaral-Zettler, L., et al. (2011).
472 Minimum information about a marker gene sequence (MIMARKS) and minimum information
473 about any (x) sequence (MIxS) specifications. *Nature Biotechnology*, 29(5), 415–420.
474 <https://doi.org/10.1038/nbt.1823>
- 475 Ying, R., Bourgeois, D., You, J., Zitnik, M., & Leskovec, J. (2019). GNNExplainer: Generating
476 Explanations for Graph Neural Networks. *Advances in Neural Information Processing Systems*,
477 32, 9240–9251. Retrieved from <https://www.ncbi.nlm.nih.gov/pubmed/32265580>
- 478 Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2017). mixup: Beyond Empirical Risk
479 Minimization. *arXiv [cs.LG]*. Retrieved from <http://arxiv.org/abs/1710.09412>