# Ethical recommendations for Artificial Intelligence technology in the Geological Sciences – with a focus on Language Models

*Paul H. Cleverley\**

Robert Gordon University, Garthdee, Aberdeen, United Kingdom

## Abstract

*Artificial Intelligence (AI) offers many opportunities for the geosciences to improve productivity, reduce uncertainty in models and stimulate discovery of new knowledge. There are also risks to geoscience, from the spread of obsolete, inaccurate and misinformation, to threats on fundamental human rights.*

*Whilst ethical AI frameworks exist from numerous institutions such as UNESCO, they are high level and lack practical detail in the geosciences particularly for Large Language Models (LLM). This is evidenced by the misalignment between the way current geoscience AI/LLMs are being designed, trained and deployed, with core ethical principles.*

*Using principles and frameworks from UNESCO and the International Science Council (ISC), a set of ten recommendations are proposed to bridge the gap between practice and these ethical frameworks. Critical Realism is used as an underlying philosophy which allows the potential to provide justifiable recommendations to ethical and moral questions using judgemental rationality.*

*These recommendations may help stakeholders in the international community reach conclusions on what "good looks like" for ethical AI in the geological sciences focusing on Language Models and their applications. This may inform developers, regulators, policy advisors, journal editors, geological surveys, societies, institutions and unions, publishers, funding bodies, geoscientists and decision makers.*

*This is believed to be the first research paper on AI ethics in the geological sciences with a focus on Generative AI. Understanding the nuances of our ethical choices for both the development and use of LLMs and other AI tools in the geosciences, has the potential to positively impact science integrity, and critically, ensure fairness, personal privacy, democratic norms and human rights are safeguarded.*

Keywords: Geology; Geosciences; Artificial Intelligence; Large Language Models; Ethics

# 1. Introduction

Artificial Intelligence (AI) presents an enormous opportunity for the geosciences to speed up information discovery and disrupt geoscientific practices by uncovering relationships within diverse datasets, revealing patterns that were previously hidden [Sun, 2024]. Over 950 papers have been published in 2024 mentioning "geoscience" and "large language models", three times that published in 2023[1],indicating heightened interest in this topic. Frontier AI also comes with risks to science, as Large Language Models (LLM) are designed to produce well written, convincing responses without any guarantees on accuracy or alignment with fact. This can be an issue as people often anthropomorphise the AI generated outputs, often trusting it as a human-like information source [Mittelstadt et al., 2023]. Without ethical guardrails, AI risks threatening fundamental human rights and freedoms [UNESCO, 2021]. There are also concerns in society regarding bias and disinformation in AI systems [Stall et al., 2023] which may be particularly pertinent to the geosciences as it continues to embrace a "social geoscience" agenda to better connect the geological sciences to multidisciplinary earth science, society, geoethics and human well-being [Stewart et al., 2022].

Geological knowledge supports the UN Sustainable Development Goals (SDG), critical for the discovery, use and conservation of natural resources, mitigation of natural hazards, the geotechnical support of infrastructure development, and protection of the environment [Acocella, 2015]. Geoscientists think of the earth and

---

[1] Number of hits using Google Scholar (Conducted 3 November 2024).

other planetary bodies as a complex system grounded in terms of time and space, with two key features distinguishing them from the general population: a long deep view of time, and expectation for low frequency, high impact events [Kastens et al., 2009]. In the field and elsewhere, geological observation comes first, followed by reconstructive spatiotemporal interpretation through abduction, characterising earth science distinctively from many of its sister sciences [Oh, 2023]. Having limited observational data, the non-linear nature of processes and capacity of modelling, make uncertainty inherent in the earth sciences. Where observational (AI training) data is limited, fundamental principles may be combined with LLM's towards a hybrid approach for modelling [Chen et al., 2024].

The world's geological surveys have made significant amounts of geological data publicly available [e.g. BGS 2024; EPOS 2022; USGS 2024; Vidovic et al., 2020]. Large amounts still remain non digital, which is less an issue about technology and AI, and more an issue of national investment choices and funding. It is believed that over 50% of academic published research papers are now open access and growing. The paywalled back catalogues held by commercial publishers and non-commercial geological society publishers is an ongoing subject of ethical discussion which is out of scope for this geoscience AI technology paper. However, where the use for text and data mining is purely academic only, some projects have agreements to use paywall content as well, creating a corpus of over 18 million geological full text papers to support geoscientific discovery[2]. Extensive seismic and well data has been made publicly available by some of the world's natural resources national regulators for many years [e.g. NSTA, 2024]. Common Crawl, an open and up to date dataset of large swaths of the Internet has been freely available for several years and is used by most foundation LLMs and also some domain specific ones[3]. Over recent years, the amount of freely available satellite and remote sensing data has also grown considerably[4]. In summary, whilst significant amounts of proprietary geoscientific data exist, there is also a vast amount which is open, free to use and growing rapidly. Ethical frameworks for geoscience AI have yet to be fully developed [Rivas et al., 2023]. Those that exist are high level and lack practical detail in many areas which may not be helping deliver ethical geoscience AI technology and deployments [Cleverley, 2024]. Frameworks for AI ethics and LLMs exist for some other disciplines [Corrêa et al., 2023; WHO, 2024] and Earth Observation [Kochupillai et al., 2022] but none within geological sciences especially around LLMs and Generative AI.

---

[2] https://geodeepdive.org/ (accessed 3 November 2024).
[3] https://commoncrawl.org/ (accessed 4 November 2024).
[4] https://eos.com/blog/free-satellite-imagery-sources/ (accessed 3 November 2024).

Recognising the unique aspects of the geosciences can help inform ethical design of AI in our discipline in addition to standard ethical principles. The proposed recommendations presented in this paper support the UNESCO recommendations on AI ethics [UNESCO, 2021] adopted by 194 countries, adding timely practical recommendations specifically relevant for the international geological sciences community. The proposed recommendations support existing self-assessment lists and analytical frameworks for trustworthy AI [ISC, 2024; EC, 2020].

The proposed recommendations are not intended to repeat all aspects of the existing UNESCO recommendations on AI ethics or international copyright law which must be adhered to. The proposed recommendations are also not intended to cover existing guidelines on the ethical use of AI in research and education [EU, 2022; EGU, 2024; Stall et al., 2023] or repeat definitions or arguments in the areas of open-source, open-access or open-science. The recommendations are also not intended to address ethics related to the creation of resource intensive Foundation LLMs or the ethics of their energy consumption [Martínez-Martín et al., 2024].

Ethics is not about rules, but judgements based on moral principles and is highly nuanced. The proposed recommendations are intended as a contribution to help bridge the gap between high level principles and practical implementation choices with geoscience AI. This may help the international community come to some conclusions on "what good looks like" for ethical AI in the geological sciences with a focus on Language Models. The proposals are intended to act as a reference to inform developers, regulators, policy advisors, journal editors, geological surveys, societies, institutions and unions, publishers, funding bodies, geoscientists and decision makers.

## 2. Methodology

Critical realism is chosen as the philosophy for the study as this creates the possibility for providing justifiable proposals to ethical and moral questions [Collier, 1994].

Critical realist philosophy has an interpretive fallibilistic epistemology. As researchers we are unable to separate ourselves from what we know, and this influences our research question, methods and findings. Critical realism therefore rejects naïve realism which assumes there is a close association between our knowledge of reality and reality itself. Knowledge can, however, be checked for its effectiveness; judgmental rationality is used to compare and assess competing theories on the basis of their explanatory adequacy or power.

Critical realism steers a course between agency (where people don't always follow

cultural norms) and structure, where there is a cultural habitus that shapes what people believe and how they behave [Tett, 2015].

As an axiology, the researcher is an academic and practitioner who has conducted research and worked in industry within the international digital geoscience sector for over 30 years. The researcher has freely blogged on Digital Geoscience techniques since 2015 which has a readership from over 144 countries. As stated by Mingers [2010], "It is the researcher(s) who, based on their own particular interests and pre-dispositions, carve out the object of scientific enquiry, both by defining time frames and the boundaries of the investigation". Despite these preconceptions, as stated by Malterud [2001] "Preconceptions are not the same as bias, unless the researcher fails to mention them."

The motivation of the researcher is twofold. Firstly, the development of AI which increases the productivity of geologists and helps ideation and discovery of new theories in geoscience to support society and industry. Secondly, the adherence of such AI systems to widely held democratic values and ethical frameworks, in a way that both safeguards the international geological community but also increases equity to technologically under-represented communities.

The methodology consists of combining the researchers own first hand experiences participating in the international geoscience AI sector, supplemented by published literature on geoscience LLM deployments. A scoping literature review was conducted using Google Scholar, focusing on geoscience LLM research with a Generative AI component that may pose different types of ethical questions. This was considered across the components of an AI system such as a training dataset, model, user interface, API [Riemer and Peter, 2024]. These are compared to the UNESCO AI ethics principles [UNESCO, 2021] and the International Science Council (ISC) analytical framework for LLMs [ISC, 2024]. The conceptual framework is shown in Figure 1.

On the left hand side (Figure 1) are the theoretical UNESCO principles for ethical AI such as "Transparency and Explainability" [UNESCO, 2021] and the micro-macro contexts of ethical AI such as "Individual, Data, Model, Geopolitics" [ISC, 2024]. On the right hand side (Figure 1) are the practical stakeholders involved in actual geoscience LLM deployments and the various technical components of an AI system such as "Model, API, Security, User Interface, Transaction Logging" [Riemer and Peter, 2024]. Generative mechanisms, not always visible, that give rise to these artefacts are shown by the blue boxes in Figure 1.

Judgements of misalignments between current practice for deployments and decisions on training data, models, design of user interfaces, deployment and hosting choices (right hand side Figure 1) to the ethical frameworks (left hand side Figure 1) give rise to recommendations for ethical improvement (centre arrows in Figure 1).

This allows a holistic socio-technical view of geoscience AI deployments, rather than a reductionist lens. The study places a particular focus on areas where there may be current misalignments or a risk of likely misalignments in the future, putting forward practical proposed alignment recommendations for discussion. Explanations for some generative mechanisms that may be driving behaviours are also postulated.
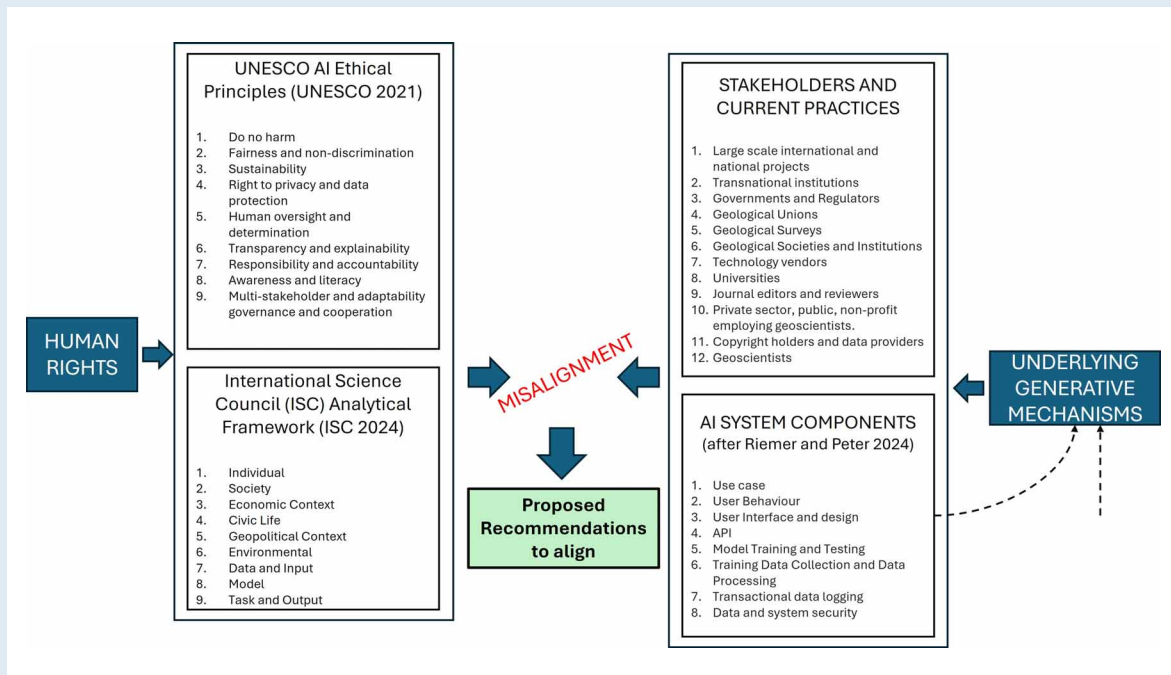


**Figure 1.** Conceptual framework for the study.

# 3. Proposed recommendations

Transparency and openness, explainability, fairness and biases, and participatory design could be considered 4 key principles for trustworthy [Albahri et al., 2024] and responsible AI in the geosciences [Stall et al., 2023; UNESCO, 2021]. Ethical communication from vendors and projects is also likely to be key to gain trust, avoiding rhetoric and hyperbole and explaining limitations with open and honest disclosure. These elements are threaded through the following ten recommendations (Table 1).

| N. | Recommendation area |
|:---:|:---:|
| 1 | Transparency of geoscience training data |
| 2 | Robustness, fairness and bias of geoscience training data |
| 3 | Openness, traceability and accuracy of geoscience models |
| 4 | Manipulation of geoscience AI models |
| 5 | Citation of sources for AI generated answers |
| 6 | Presentation of uncertainty and reasoning |
| 7 | Disclosure of appropriate digital watermarking |
| 8 | Protection of personal data |
| 9 | Treatment of uploaded geological data |
| 10 | Use case, jurisdiction, business model, ownership, governance and participatory design, and geopolitics |

**Table 1.** Ten recommendation areas.

Each recommendation will be discussed in the following sections highlighting opportunities and risks with examples where they exist.

Recommendation #1: Transparency of geoscience training data

All reasonable efforts should be made to publish and make openly available full datasets used to train any geoscience AI through FAIR (Findable, Accessible, Interoperable and Reusable) principles. These datasets might be used for machine learning using literature/data/images, machine learning using geological Question & Answer (Q&A) pairs, Retrieval Augmented Generation (RAG) vector chunked input data, Knowledge Graphs (KG), databases or a combination of these elements. This applies equally for unsupervised machine learning, as well as supervised machine learning such as labelling of images for training (machine vision), text or other data-structures. This would support Open Science principles for reproducibility and the ability to identify undesirable biases used during "training" by other researchers and actors.

Where training data is not publicly released by the institution or vendor, it should be encouraged to disclose why this is the case. This might be because proprietary Intellectual Property (IP) is being developed by the institution or vendor through its training datasets. Where training data is not released, institutions and vendors are encouraged to describe the training dataset in as much detail as possible, including provenance and percentage of non-peer review content used. There have been examples over the past year or so, where large volumes of geoscience specific training data were illegally obtained by institutions building LLMs [Cleverley et al., 2024] which categorically breaks the UNESCO recommendations on AI ethics and may evidence lack of governance.

One should not claim to be open-source or Open science in communications for any geoscience AI where the training datasets are not publicly released in full, enabling re-use. Releasing these data differentiates an open academic project, from one which is more closed and proprietary/commercial in nature. It is important for transparency that the international geoscience community understands where projects and AI technology are positioned. A lack of transparency can fracture trust and cooperation.

Release of training data is encouraged to foster international cooperation, development of digital skills and equity in geoscience AI. The for-profit sector is also encouraged to publicly release as much which is reasonably possible given their business model.

## Recommendation #2: Robustness, fairness and bias of geoscience training data

It is desirable that training data should be collated so it is as broad as possible given the use case. For example, avoiding overt bias towards training data from any single country (especially if the audience is an international one), significant non-peer reviewed content, consideration of under-represented communities, and any historic biases that perpetuate discrimination. Heterogeneous and multidisciplinary training datasets are encouraged to stimulate new connections, and supporting multilingual training data is important for equity and diversity.

Some datasets are of course inherently biased. For example, georeferenced remote sensing data to track anthropomorphic climate change is biased towards developed nations within North America, Europe and China due to the number of satellite missions owned by these countries, which is independent of the degree to which they are affected by climate change [Kochupillai et al., 2022].

Keeping geoscientific information up to date within an AI system is crucial. For example, one well known AI chatbot gave the incorrect date for the base of the Barremian for over a year, because this was updated by the International Commission on Stratigraphy (ICS) after the AI model had been trained. For geoscience chatbots in particular, how current their training data is, should be clearly communicated to geoscientists in the user interface to mitigate the use and further dissemination of obsolete/out of date information.

These geoscience training data will be filtered and manipulated according to certain criteria by institutions and vendor data scientists. This may not always be apparent from the outputs of recommendation #1 even if they are publicly released. Any criteria which seek to manipulate or bias geoscience related (including social geoscience) training data to ensure any government and its policies, political party, organisation or individual is seen in an uncritical viewpoint, may infringe on democratic values and should be avoided. Any criteria to manipulate or omit geoscience related data (including social geoscience) from information collections, or promotion of certain content for motives other than purely geoscience aims should also be avoided. This may be akin to aspects of "Dark Knowledge" constructs introduced by Burnett and Lloyd [2020]. To respect national sovereignty for some countries, training data may need to be manipulated in this way, but the promotion of such to the international geoscience community in the name of a geological union should be avoided.

Methodological procedures for filtering training data are encouraged to be published. There are examples of LLMs that are truly open source, where institutions have openly released their training data and methods used for reproducibility [AllenAI, 2024]. Openly releasing these procedures for collation and

filtering of training data differentiates an open academic project, from one which is more closed and proprietary/commercial in nature. It is important for transparency that the international geoscience community understands where projects and AI technology are positioned in this spectrum.

The release of data laid out in recommendations #1 and #2 would allow geoscience researchers to critique and debate the advantages and disadvantages of such training datasets and techniques, and combine datasets appropriately, producing better science and AI.

## Recommendation #3: Openness, traceability and accuracy of geoscience models

One key element to support trustworthy AI is published documentation of the processes, algorithms and parameters used to generate models, supporting explainability and reproducibility. This is of paramount performance for geoscience AI where it is used to support prediction of natural disasters such as landslides, earthquakes and flooding as it directly impacts people's lives [Albahri et al., 2024]. Benchmarking of AI against a test set which may include calibrated questions, can replace marketing hyperbole, helping geoscientists ascertain the accuracy of such AI systems, comparing it with others already available. Ensuring these questions and test sets are published under FAIR principles can help geoscientists understand the types of questions used to test geoscience AI systems, and the AI's performance against those questions and test sets. Test sets of questions could be amalgamated to form a single global benchmark for geoscience LLM type AI solutions.

Recent European research papers published in 2024 on geoscience LLMs have not always disclosed the geoscience questions used to benchmark the AI [Baucon and de Carvalho, 2024], something that geoscience journal editors might wish to improve on with respect to ethical disclosure guidelines. Journal editors of non-technical journals (in an AI sense) may wish to consult with new Task Group on AI in Geosciences of the IUGS Commission on Geoethics[5] that have been established, where they feel their peer review process does not have the necessary expertise to review geoscience AI aspects of papers.

Geoscience specific AI models derived from training data (recommendation #1) should be placed in freely available model repositories with standard open-source licenses under FAIR principles so the entire community can download and use the models in their own secure information environments. Exemplars in the earth

---

[5] https://www.geoethics.org/tg-cg-ai-geosciences (accessed 2 December 2024).

sciences would be the NASA Indus Language Models that are openly available [Bhattacharjee et al., 2024].

It is important to note that if an existing "open weight" foundation LLM is used by an institution or vendor and further trained in some way, this does not mean by inference their resulting AI model or tool is also "open". This is only the case if that derived model itself is placed in a freely available model repository for download and the source code of their tools placed in a freely available code library with an open-source license. Otherwise, the institution or vendor has a proprietary (potentially commercial) AI model and tool.

One recent North American published research paper made the training data publicly available but not the geoscience AI model [Lawley et al., 2022] and vice versa for another North American research paper [Bhattacharjee et al., 2024]. In both cases all data may have been openly released, had the journal editor guided authors according to these proposed ethical recommendations.

If models are not made open and only made available as part of web hosted tools or APIs by an institution or vendor, this should be made clear. Openly releasing these models differentiates an open academic project and technology, from one which is more closed and proprietary/commercial in nature. Open release of models and training data is to be encouraged as this has the potential to increase international cooperation and equity across the geosciences as projects from different countries build on each other's research and work.

There are a number of free or cheap API's available to large LLMs making it quite easy to create a web application which acts as an AI conversational assistant with some prompt engineering. In this situation it is important to disclose to what extent the web application adds a significant intellectual contribution over and above that API. There have been situations where a web application has been created, branded and marketed as geology specific to the international geological community, but for some aspects it may virtually be 99% OpenAI's ChatGPT API. This has not been disclosed and this risks misleading the user in terms of what they think they are using. It also does not give full credit to the authors of the work creating that underlying API; this is not explainable AI, and this opaqueness should be avoided.

## Recommendation #4: Manipulation of geoscience models

Feedback mechanisms to ensure geoscience expert knowledge and geological principles are incorporated into models are already being implemented and will likely be more important in the future. This may mitigate some "Hallucination" effects for obviously incorrect or physically unrealistic answers and outputs. Any

alignment of geoscience AI models such as Reinforcement Learning through Human Feedback (RLHF), trigger words or structured guiding of answers should be published and declared. Any algorithmic method or tagging which affects RAG chunk retrieval should also be published to aid transparency.

Any criteria which seek to manipulate or bias geoscience related (including social geoscience) training data to ensure any government and its policies, political party, organisation, culture or individual is seen in an uncritical viewpoint, may infringe on democratic values and should be avoided for the wider international geoscience community. Any criteria to manipulate or omit geoscience related data or promote certain content for motives other than purely geoscience aims to an international community should be avoided. It is recognised that to respect national sovereignty for some countries, LLM models may need to be "aligned" driven by a generative mechanism of technological sovereignty, however promotion of such to the international geoscience community should be avoided.

Openly releasing these methodological data differentiates an open academic project and technology, from one which is more closed and proprietary/commercial in nature. It is important for transparency that the international geoscience community understands where projects and AI technology are positioned on this spectrum.

### Recommendation #5: Citation of sources for AI generated answers

Any AI generated output must ensure sources are cited in its user interface from which the multimodal output (text, data-files, images, tables, graphics and video) and assertions were generated. No answers should be provided without a link to their source(s) to support traceability. This will allow the Geoscientist to check the validity of the answer. User interface scaffolding (functionality) may be beneficial, to present the AI generated output visually alongside original source text and images to aid verification for geoscientists and help explainability for answers.

This proposed recommendation also recognises the authors and copyright holders used to train the AI models, as we all do when we cite sources in papers ourselves as researchers and is also often a legal condition of the licensing agreement even for open-access CC-BY-4.0 content. RAG is the most likely arrangement to deliver this in practice for an LLM. Exemplars exist in the geosciences from mid-2023. Evidenced by the publicly available LLM digital conversational assistant that references answers from geological reports held by the Norwegian Petroleum Directorate[6]. Despite this, at least two geological LLM chatbots were publicly

---

[6] https://npd.fabriqai.com/ (accessed 7 December 2024).

released to the international geological community in 2024 without referencing the sources of the answers produced.

So called "Hallucinations" are inevitable with non-deterministic LLM's. Presenting user interface scaffolding to show original source data in innovative ways may be a way to support the scientific inquiry process where all source information needs to be cited.

## Recommendation #6: Presentation of uncertainty and reasoning

It is important geoscience AI generated outputs are grounded in real world geological principles and factual knowledge, and a measure of uncertainty provided for answers. Human agency should be paramount in design, ensuring geologists are in the driving seat of AI tools, with the ability to view underlying source data and change underlying algorithmic parameters.

Automated fact checking may be desirable, against authoritative public domain geological reference data and fundamental principles. For example, in 2023 OpenAI's ChatGPT convincingly "Hallucinated" an approved mineral with diacritical marks, named "Eötvösite" which does not exist [Ralph, 2023]. It is believed this was "generated" from a Reddit Dungeons and Dragons forum discussing an imaginary world of made-up minerals, a common source on Common Crawl used by foundation LLM's. If the AI system had cross referenced the AI generated answer to authoritative peer reviewed public geological data of approved minerals, this could have highlighted this to the geoscientist through a low probability confidence measure that the answer may be suspect. Note, at that time OpenAI's ChatGPT did not cite its sources as standard, so it was not easy to check where the answer came from, highlighting the importance of proposed recommendation #5 for ethical AI in science use cases.

For factual answers, Wei et al. [2024] found a positive correlation between the stated confidence of an LLM and its accuracy, although it currently appears to be below a straight line of confidence v accuracy. Nevertheless, this could become a best practice in prompt engineering behind the scenes, where every answer provided by an LLM to a geoscientist comes with an assessment percentage of confidence combined with a heuristic of how often that model tends to get factual answers correct. If an LLM for example has a heuristic of 60% accuracy and it has a stated confidence of 80% in an answer, the overall confidence is 0.48. The key element proposed, is the construct of providing *some measure of uncertainty* to the user.

Currently, there are little to no examples of this in user interfaces of geoscience focused LLM driven digital assistants, despite it being straightforward to do using prompt engineering. As a generalised indicative example, in the background, a

context can be set by the designer such as *"You are a geological assistant. Always answer with a percentage of how confident you are in the accuracy of your answer"*. In the user interface when a geoscientist now asks a question such as *"What minerals are found in granites?"* along with the answer a confidence percentage can be displayed as standard with a narrative.

Whilst this may cater for simple factual answers, exploratory search goals are an area for further research, where there is no right answer. Techniques in this area may include provision of multiple answers from different sources (with perhaps multiple underlying models created to reflect different sources or parameters) to a user's prompts, questions and intent. Another technique may be to use the level of dissonance in the literature (disagreement between authors and sources) which could be presented as one type of uncertainty or confidence measure within the user interface.

AI systems curate what we see so therefore influence what we know, in some regards they have become an epistemology. This algorithmic curation is therefore of ethical significance in AI designed for scientists. Providing scaffolding in the user interface to allow a geoscientist to change algorithmic variables may avoid perceptions of a "black box" for geoscience AI, improving transparency which can build trust. For example, when using RAG, most LLM answers often only come from the top ranked "chunks" which may be just from a few papers. For exploratory search goals, what is statistically ranked top in Information Retrieval (IR) of text chunks can be somewhat arbitrary in a large corpus with many results of relevant content of which there is no right answer.

Allowing the geoscientist a sense of control to change the retrieval algorithm or recommendations presented by AI systems may support transparency and explainability, influencing dimensions such as the importance of publish date, citation rank of the authors, paper or journal, country of origin, native language, source/collection, peer/non peer reviewed etc. This may help surface different answers, responses and recommendations, mitigating algorithmic biases embedded by the designers of the system and providing agency to the geoscientist to understand aspects of uncertainty in the information space.

Designers of such AI systems should, where possible, present to the user "the workings" of answers generated to show the geological reasoning rather than just the final answer. To take a trivial example, asking what is older, the "Barremian" or the "Bathonian", to an AI system, an answer that just returns *"Bathonian"* is simply reciting. Whereas an answer with reasoning may say *"Bathonian is a geological stage between 168.3 to 166.1 million years, Barremian is a geological stage 129.4 to 125 million years, so the Bathonian is around 40 million years older than the Barremian. The answer is the Bathonian"*. If such techniques are embedded within

multi-modal geoscience AI tools, it may aid explainability and transparency of the geoscience AI system.

## Recommendation #7: Disclosure of appropriate digital watermarking

Hidden cryptographic digital watermarking of text, images and video is already implemented within some generic AI systems. This can support copyright holders and owners of AI systems trace provenance of data, plagiarism, misuse, abuse or illegal use of outputs from AI systems. However, it also has the potential for user surveillance impacting data privacy. All reasonable measures should be taken to make geoscientists aware of hidden cryptographic digital watermarking of outputs from geoscience specific AI tools. The geoscientist should be informed, perhaps via a hyperlink in the user interface, what data has been encrypted in the output including whether any of their own personal data was used. It is proposed that for data privacy reasons, geoscientists should ideally be allowed to choose to opt-out of any encryption of their personal details (such as name, email, location, IP address, prompts used etc.) in the outputs of geoscience AI tools.

Governance of blockchain / distributed ledgers for digital watermarks should be approved by all stakeholders, be openly disclosed, with as transparent mechanisms as reasonably possible. It is proposed that images, text or data should not be digitally watermarked from an AI system if it is from the training data or not generative, following the European AI Act (section 133) for Digital Watermarking [EU, 2024]. This includes proportionate use of digital watermarking, so should not be used where the AI is assistive only. For example, AI tools that allow natural language interaction with structured data are generally assistive, such as interacting with a geospatial shapefile to create a geological map, so no digital watermarking should be used.

## Recommendation #8: Protection of personal data

For web hosted AI tools, personal data, such as email addresses, affiliations, IP addresses, location data, prompts and usage tracking patterns should be safeguarded and anonymised if used to improve performance of AI models. Personal details should not be passed to third parties or used for any purpose other than that consented to by the user. Geoscientists should have optionality, the right to withdraw consent, the right to object and the right to be forgotten for any geoscience AI tool they sign up for or use. For example, the European Union AI Act [EU, 2024] provides significant laws that need to followed regarding the storage and

use of personal data of EU citizens. Any user should be informed and provided with the full extent of any surveillance of their online activities through any geoscience AI system. National legislation on data privacy for citizens based on nationality and location must be respected and complied with by vendors of geoscience AI tools.

## Recommendation #9: Treatment of uploaded geological data

For web hosted AI tools, all reasonable measures should be taken to inform geoscientists before they upload documents, maps, datafiles etc., that owners of the AI are likely to have the rights in their terms and conditions to use and share these data for their own purposes. Generative mechanisms operating here are likely to be a need for an organisation or nation state to mine geoscientific data, from wherever it can be obtained, to support areas such as natural resources security (critical minerals, fossil fuels and sustainable energy), and/or natural disaster response and early warning, and/or the furthering of science.

Confidential or copyrighted content should therefore not be uploaded unless legal agreements and architectures have been setup in participation with geoscientists and their institutions according to their information security policies. Geoscientists should be made aware that the cost of many free web hosted AI tools and APIs may be their (or their institutions) data. The AI owners should implement appropriate measures to make geoscientists aware of these terms with clear disclosure and messaging. Whilst legal notices are included within the lengthy terms and conditions of use for some geoscience LLMs, it is argued this should be made more obvious. When presented with an "upload document" button in the user interface, some text indicating what rights the user is giving away by uploading this document should be stated. A system that informs the user in this way, is more ethically safe, than one that relies on a user to read lengthy legal jargon on a terms and conditions link that many geoscientists will never read.As part of governance processes and good practice due diligence, geoscience AI systems should ideally have automated proactive measures to identify users who are uploading copyrighted or obviously confidential information, rather than wait for such issues to be reported by third parties.

## Recommendation #10: Use case, jurisdiction, business model, ownership, governance and geopolitics

Use cases should be fully disclosed, minimising unintentional or unexpected harms or preventing them entirely if possible. There should be a focus on *"what is the*

*problem we are trying to solve"* by funding agencies, researchers, designers and project managers of geoscience AI systems.

In Earth Observation for example, labelling (for AI prediction) dwellings on satellite imagery and tagging them as "slums" could pose ethical risks of stigmatization, and the identification of small-scale artisanal mining may cause some to lose their livelihoods posing ethical dilemmas [Kochupillai et al., 2022].

It is proposed that a cautionary approach is taken regarding any use case where LLMs are promoted to be used to mark geoscientists work. Some early grading systems are being promoted to geology students and teachers [Baucon and de Carvalho, 2024]. Whilst this stimulates possibilities, there is little transparency on how the marking works, what biases may exist and what negative consequences this may have. We should be careful not to "speed things up" to the detriment of human input and accuracy of results [Martínez-Martín et al., 2024]. Participatory design from stakeholders will likely lead to more thoughtful and ethical solutions, rushing out tools quickly just because we can, without ethical safety assessments, is arguably not what "good looks like" for ethical AI in the geosciences.

As put by researchers regarding the use of generative AI at UNESCO Global GeoParks, *"due to the current state of generative AI and considering all its development and workflow, they neither provide reliability, nor can they be aligned with the [UNESCO Global GeoPark] core values"* [Martínez-Martín et al., 2024] who go on to say the need for at least one professional to analyse results obtained from the AI is essential. Every use case should undergo an ethical assessment by those skilled in the art and ideally include people outside the technology/research project concerned. Where conflict may exist between core values of an area and the use of generative AI, this requires careful discussion. At the moment in the geosciences, evidence suggests this is not happening.

If the audience is international, non-native English speakers will be disadvantaged by an AI system that requires interrogation only in English. Multi-lingual support and other user interface scaffolding techniques is desirable.

Arguably, significant data driven scientific discovery in the geosciences involving Language Models is more likely to come from dense vectors/embeddings models [Lawley et al., 2022], rather than the "consumer" style chatbot AI experience many are focusing on. This may be important in terms of progressing international cooperation. For example, copyright holders that may be reticent to have their IP used for LLM training in chatbot style systems that are not just used by academia (they may also be used by state actors and private sector). This may be because LLMs can regurgitate this IP (text, data and images) sometimes verbatim. Therefore copyright holders may be far more amenable to provide or publicly license their IP for embeddings use only. There may be good ethical reasons to do so in the name

of equity, balancing copyright protection with the needs of the international geoscience community to accelerate data driven discovery. A lack of knowledge and an overly simplistic view of Language Models may be blinding some in the geoscience community to opportunities in this area.

For hosted solutions, the national laws by which any geoscience AI is governed should be disclosed. Potential for infringement of democratic values such as freedom of expression and information should be avoided. Any national laws embedded in AI which have the affordance to refuse to answer certain geoscience related questions (including social geoscience), or bias geoscience related outputs to ensure a government and its policies, political party, organisation or individual is seen in an uncritical viewpoint, should be avoided for the international geoscience community [Cleverley et al., 2024]. Geoscience AI that infringes such fundamental democratic values cannot be AI for international public good.

The business model (private, public, non-profit) for the institution or vendor of the geoscience AI should be disclosed along with the academic and industrial sector(s) targeted.

Many private sector and state-owned organisations employing geoscientists have, and continue to use, LLM's applied to their own proprietary information. Some of the world's largest companies such as Aramco [Malin, 2024] and ExxonMobil [Denli et al., 2021] have fine-tuned their own LLM models and there are many innovative startups. For competitive reasons these training data, models and tools are unlikely to be publicly released in full. However, sharing of experiences and findings applying AI to the geosciences with the wider community are encouraged to increase information sharing between industry and academia. These proposed recommendations on ethics apply equally to the "for-profit" sector as much as they do for the "non-profit" and public institutions.

Funding and legal ownership should be disclosed to the geoscience community. Mechanisms should be put in place to ensure clear accountability for the Geoscience AI, responsibilities for the development, deployment and/or use of geoscience AI systems, identifying and mitigating risks in a transparent way. Ethical impact assessments should be conducted and approved by stakeholders before the piloting and public release of major AI tools that involve LLM's in the geosciences, with regular auditing systems. All reasonable measures must be taken to mitigate the misuse of geoscientific AI systems.

The International Science Council (ISC), which the International Union of Geological Sciences (IUGS) adheres to, has developed an analytical framework for LLM's [ISC, 2024]. This checklist has items on geopolitics. Those items of particular relevance to the geosciences include:

- *"Geopolitics: Is a desire for technological sovereignty driving behaviours?"*

- *"Digital Colonialism: Could state or non-state actors harness systems and data to understand and control other countries' populations, ecosystems or undermine jurisdictional control".*
- *"Geopolitics: Could the system stir competition between nations over harnessing individual and group data for economic, medical and security interests?"*
- *"Digital Divide: Are existing digital inequalities exacerbated, or new ones created?"*
- *"Dual Use: Is there a possibility for both military application as well as civilian use?"*

These questions require sensitive consideration by the world's geological unions. It is important large scale projects do not unintentionally exacerbate or create new digital inequalities, introduce conflicts, unintentional harms and infringe democratic values. Some of these points are discussed further in the next section.

# 4. Discussion

There are benefits to openly release geoscience AI training data, resultant models and tools to support ethical deployment of geoscience AI for the public good. However, it may not always be possible for these data to be released due to IP and commercial business models. This should always be made clear to the geological community where an AI project or technology is positioned.

Geoscience AI vendors and projects can publicly declare they are adhering to all UNESCO ethics guidelines, however under closer scrutiny they may not be ethical to UNESCO principles, based on comparing against the practice based proposed recommendations in this study.

Many aspects of search engines and LLM Chatbot design are likely to be driven in 2024 by a "Google Habitus" and "ChatGPT Habitus" respectively as generative mechanisms. This may explain why some AI tools have been developed the way they have been, without linking to sources for example. But we don't have to follow the habitus.

Cognitive biases, such as emotional attachment, self-interest and social biases such as groupthink may be causal factors for some the behaviours observed in the literature. These observations include decision makers on some major geoscience AI projects not identifying or dismissing/ignoring ethical aspects that may seem serious to many subject matter experts outside their project. Lack of experience and deep expertise in this subject matter, a desire to not "rock the boat" in projects where they have no real control and ambiguous accountability, may be contributing causal factors. These factors can all lead good people to make or participate in bad ethical decisions.

Geoscience AI models and tools could be assessed against these 10 recommendations where appropriate. With an assessment of transparency, openness, explainability, fairness & bias, and participatory design according to each of these 10 proposed recommendations. Institutions or vendors could self-assess their AI. International geoscientific bodies may also wish to play an oversight role. Journal editors may also wish to use the recommendations to enhance review and disclosure processes for submitted papers on the topic of AI in the geosciences. It's important to note that some Language Models might just be trained for a specific narrow task like classification, or extraction such as Named Entity Recognition (NER) or provide embeddings/dense vectors whose purpose is not generative, so all 10 recommendations won't always apply. For example, in the previous examples, proposed recommendation #5 (citing source of answers) and #7 (disclosure of digital watermarking) are not applicable. Taking another example, for a remote sensing AI model detecting sedimentary copper, or an AI model which detects mineral grains in a rock thin section, released publicly in a model repository available to download, recommendations #8 (Protection of personal data) and #9 (Treatment of uploaded data) will not necessarily apply as models are within the user's own firewall/network. These proposed recommendations can therefore be scaled appropriately.

The presentation of uncertainty and reasoning (recommendation #6) is probably the most open ended proposed recommendation. This is more a call for designers of geoscience AI systems to reflect and think in terms of uncertainty of answers/outputs at many levels and how best to implement an assessment of uncertainty within geoscience AI tools.

Providing free web hosted geoscience AI solutions from technologically advanced countries can help deliver new capabilities and increase equity for geoscientists in countries with lower technological maturity, incomes and fewer funding options particularly in parts of the Global South. A competing theory on the ethics is that this could perpetuate inequalities, by consolidating know-how, IP and geoscientific data into just a few technologically advanced large countries. Any solution which is proprietary in nature, which provides free access to its web hosted AI geoscience solutions targeting low and middle income countries in return for their geoscientific data (which may be exploited for natural resources and energy) could arguably be viewed as a form of Data Colonialism in the geosciences, as per the International Science Council (ISC) analytical framework. Respect for the self-governance of indigenous people's data is part of the UNESCO ethical principle on governance. Recommendations from this research propose that geoscience LLM training data, models and source code for tools should be openly released, especially from academic non-profit geoscience initiatives supported by geological unions. This

would allow the potential for National Geological Surveys around the world, and other geological organisations, to take what has been done, modify as appropriate, and host derived AI/LLMs within their own respective jurisdictions that they control and can further develop. This would allow benefits to be realised, without geologists and institutions being forced to relinquish sovereignty over their data assets to technologically advanced large countries, in exchange to use such capabilities.

It is argued this is a more ethical approach. Where individuals and institutions wish to publicly share particular geoscience data they can (and are doing so). As well as it being more ethical, more international cooperations, projects, institutions and individuals can take what has been openly released, and develop it further, providing greater diversity of thought and brainpower, and release that derivative work openly. It is argued this "open" and "federated" approach towards international geoscience AI, is likely to be more ethical and more successful than any "centrally" controlled approach dominated by a single large country. This can accelerate how the geosciences tackle the UN Sustainable Development Goals (SDGs).

Any geoscience AI which is intentionally configured through manipulation of training data, and/or its LLM models and/or its legal terms and conditions, to protect the image of a certain government (its policies), any political party, any organisation or any individual, is unsuitable for promotion to the wider international scientific community in the name of an international geological institution. Democratic norms of freedom of expression and information, and human rights within the international community should be protected.

# 5. Conclusion

As stated previously, ethics is not about rules, it's about judgements made against principles which can give several answers. However, some explanations may carry more ethical weight than others when tested against detailed ethical frameworks. The proposed recommendations are intended as a contribution to help bridge the gap between high level principles and practical implementation choices with geoscience AI and stimulate further research. This may help the international community come to some conclusions on "what good looks like" for ethical AI in the geological sciences with a focus on Language Models.

These proposed recommendations sit alongside comprehensive recommendations on AI Ethics from UNESCO. The aim of these proposed recommendations is to connect ethics to practice for geoscience AI especially around LLMs and promote debate and discussion on this topic within the earth sciences.

The opportunity for AI in the geosciences continues to be potentially

transformational for productivity and scientific discovery. Thoughtful consideration on ethical design can help secure safe and strong foundations to realise the benefits of AI in the geosciences for current and future generations.

# References

Acocella V., (2015). *Grand challenges in Earth science: research toward a sustainable environment.* Frontiers in Earth Sciences, 3. https://doi.org/10.3389/feart.2015.00068

Albahri A.S., Khaleel Y.L., Habeeb M.A., Ismael R.D., Hameed Q.A. et al., (2024). *A systematic review of trustworthy artificial intelligence applications in natural disasters.* Computers and Electrical Engineering, 118(B). https://doi.org/10.1016/j.compeleceng.2024.109409

AllenAI, (2024). *OLMo 2 is a family of fully-open language models, developed start-to-finish with open and accessible training data, open-source training code, reproducible training recipes, transparent evaluations, intermediate checkpoints, and more.* https://allenai.org/olmo (accessed 7 December 2024).

Baucon A. and de Carvalho C.N., (2024). *Can AI Get a Degree in Geoscience? Performance Analysis of a GPT-Based Artificial Intelligence System Trained for Earth Science (GeologyOracle).* Geoheritage, 16, 121. https://doi.org/10.1007/s12371-024-01011-2

BGS, British Geological Survey, (2024). Website: https://www.bgs.ac.uk/geological-data/opengeoscience/ (accessed 3 November 2024).

Bhattacharjee B., Trivedi A., Muraoka M., Ramasubramanian M., Udagawa T. et al., (2024). *INDUS: Effective and Efficient Language Models for Scientific Applications.* https://arxiv.org/pdf/2405.10725 (accessed 3 November 2024).

Burnett S. and Lloyd A., (2020). *Hidden and forbidden conceptualising Dark Knowledge.* Journal of Documentation, 76(6), 1341-1358. https://doi.org/10.1108/JD-12-2019-0234

Chen M., Qian Z., Boers N., Creutzig F., Camps-Valls G. et al., (2024). *Collaboration between artificial intelligence and Earth science communities for mutual benefit.* Nature Geoscience, 17, 949-952. https://doi.org/10.1038/s41561-024-01550-x

Cleverley P., (2024). *Geoscience AI in crisis?* Geoscientist, 34(3), 22-25. https://doi.org/10.1144/geosci2024-024

Cleverley P., Peppoloni S., Baily C. and Thompson S., (2024). *Advancing transparent and ethical AI.* Geoscientist online. https://geoscientist.online/sections/viewpoint/advancing-transparent-and-ethical-ai/ (accessed 2 November 2024).

Collier A., (1994). *Critical Realism: An Introduction to Roy Bhaskar's Philosophy*. London, UK: Verso.

Corrêa N.K., Galvão C., Santos J.W., Del Pino C., Pontes Pinto E. et al., (2023). *Worldwide AI ethics: A review of 200 guidelines and recommendations for AI governance*. Patterns, 4(10). https://doi.org/10.1016/j.patter.2023.100857

Denli H., Chughtai H.A., Hughes B., Gistri R. and Xu P., (2021). *Geoscience Language Processing for Exploration*. Abu Dhabi Petroleum Exhibition & Conference, Abu Dhabi, UAE, November 15-18, 2021. https://doi.org/10.2118/207766-MS

EC, (2020). *European Commission Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self-assessment.* https://digital-strategy.ec.europa.eu/en/library/assessment-list-trustworthy-artificial-intelligence-altai-self-assessment (accessed 3 November 2024).

EGU, European Geosciences Union, (2024). *Statement on the use of AI-based tools for the presentation and publication of research results in Earth, planetary and space science.* https://www.egu.eu/news/1031/statement-on-the-use-of-ai-based-tools-for-the-presentation-and-publication-of-research-results-in-earth-planetary-and-space-science/ (accessed 3 November 2024).

EPOS, European Plate Observing System, (2022). [Ref: Cocco M., Freda C., Atakan K., Bailo D., Saleh-Contell K., Lange O. and Michalek J. (2022). *The EPOS Research Infrastructure: a federated approach to integrate solid Earth science data and services.* Annals of Geophysics, 65(2). https://doi.org/10.4401/ag-8756].

EU, European Union, (2022). *Directorate-General for Education, Youth, Sport and Culture, Ethical guidelines on the use of artificial intelligence (AI) and data in teaching and learning for educators, Publications Office of the European Union.* https://data.europa.eu/doi/10.2766/153756 (accessed 3 November 2024).

EU, European Union, (2024). *Artificial Intelligence Act.* https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32024R1689 (accessed 2 November 2024).

ISC, International Science Council, (2024). *A guide for policymakers: Evaluating rapidly developing technologies including AI, large language models and beyond.* https://council.science/wp-content/uploads/2024/04/A-guide-for-policy-makers_AI.pdf (accessed 2 November 2024).

Kastens K.A., Manduca C.A., Cervato C., Frodeman R., Goodwin C. et al., (2009). *How Geoscientists Think and Learn.* Eos, 90(31), 265-266. https://doi.org/10.1029/2009EO310001

Kochupillai M., Kahl M., Schmitt M., Taubenbock H. and Zhu X.X., (2022). *Earth Observation and Artificial Intelligence: Understanding emerging ethical issues and opportunities.* IEEE Geoscience and Remote Sensing Magazine, 10(4). https://doi.org/10.1109/MGRS.2022.3208357

Lawley J.M., Raimondo S., Chen T., Brin L., Zakharov A. et al., (2022). *Geoscience*

*language models and their intrinsic evaluation.* Applied Computing and Geosciences, 14. https://doi.org/10.1016/j.acags.2022.100084

Malin C., (2024). *World's Largest Industrial LLM revealed.* Middle East AI News, https://www.middleeastainews.com/p/aramco-launches-largest-industrial-llm (accessed 4 November 2024).

Malterud K., (2001). *Qualitative research: standards, challenges and guidelines.* The Lancet, 358(9280), 483-488. https://doi.org/10.1016/S0140-6736(01)05627-6

Martínez-Martín J.E., Rosado-González E.M., Martínez-Martín B., Sá A.A., (2024). *UNESCO Global Geoparks vs. Generative AI: Challenges for Best Practices in Sustainability and Education.* Geosciences. 14(275). https://doi.org/10.3390/geosciences14100275

Mingers J., (2011). *The Contribution of Systemic Thought to Critical Realism. University of Kent Business School.* Working Paper no.232. https://kar.kent.ac.uk/28404/ (accessed 4 November 2024).

Mittelstadt B., Wachter S. and Russell C., (2023). *To protect science, we must use LLMs as zero-shot translators.* Nature Human Behaviour, 7, 1830–1832. https://doi.org/10.1038/s41562-023-01744-0

NSTA, North Sea Transition Authority, (2024). Website: https://ndr.nstauthority.co.uk/ (accessed 4 November 2024).

Oh P.S., (2023). *Abduction in Earth Science Education.* In Magnani L. (ed.), Handbook of Abductive Cognition, Springer, Cham. https://doi.org/10.1007/978-3-031-10135-9_48

Ralph J., (2023). *ChatGPT – A system too clever and creative for its own good?* https://www.mindat.org/mesg-619454.html (accessed 3 November 2024).

Riemer K. and Peter S., (2024). *Conceptualising generative AI as style engines: Application archetypes and implications.* International Journal of Information Management, 79. https://doi.org/10.1016/j.ijinfomgt.2024.102824

Rivas P., Thompson C., Tafur B., Khanal B., Ayoade O. et al., (2023). *Artificial Intelligence in Earth Sciences.* Chapter 15 – AI Ethics for Earth Sciences, pp. 379-396. https://doi.org/10.1016/B978-0-323-91737-7.00007-4

Stall S., Cervone G., Coward C. and Cutcher-Gershenfeld J., (2023). *Ethical and Responsible Use of AI/ML in the Earth, Space, and Environmental Sciences.* American Geophysical Union (AGU). 10.22541/essoar.168132856.66485758/v1

Stewart I., Capello M.A., Mouri H., Mhopjeni K. and Raji M., (2023). *Three Horizons for Future Geoscience.* Earth Science, Systems and Society, 3. https://doi.org/10.3389/esss.2023.10079

Sun Z., Brink T., Carande W., Koren G., Cristea N. et al., (2024). *Towards practical artificial intelligence in the earth sciences.* Computational Geosciences. https://doi.org/10.1007/s10596-024-10317-7

Tett G., (2016). *The Silo Effect.* UK: Abacus.

UNESCO, (2021). *Recommendations on the Ethics of Artificial Intelligence.* https://unesdoc.unesco.org/ark:/48223/pf0000381137 (accessed 3 November 2024).

USGS, (2024). *Earthquake Hazards Program: Lists, Maps and Statistics.* https://www.usgs.gov/programs/earthquake-hazards/lists-maps-and-statistics (accessed 7 December 2024).

Vidovic J., Schavemaker Y., Witteman T., Tulstrup J., van Gessel S. et al., (2020). *EuroGeoSurveys: from a non-profit association to a geological service for Europe.* In Hill P.R., Lebel D., Hitzman M., Smelror M., Thorleifson H. (eds.), The Changing Role of Geological Surveys, Geological Society, London, Special Publications 499, pp. 129-137. https://doi.org/10.1144/SP499-2019-47

Wei J., Nguyen K., Chung H.W., Jiao J., Papay S. et al., (2024). Introducing SimpleQA, Website: https://openai.com/index/introducing-simpleqa/ (accessed 2 December 2024).

WHO, World Health Organization (2024). *Ethics and governance of artificial intelligence for health. Guidance on large multi-modal models.* https://iris.who.int/bitstream/handle/10665/375579/9789240084759-eng.pdf;jsessionid=9E1B3B5DDDFADE971D7795EB4F3FB7C6?sequence=1 (accessed 7 December 2024).

*Corresponding author: **Paul H. Cleverley**

e-mail: paulhughcleverley@gmail.com