

GEOGALACTICA: A SCIENTIFIC LARGE LANGUAGE MODEL IN GEOSCIENCE

Zhouhan Lin^{1,*}, Cheng Deng¹, Le Zhou¹, Tianhang Zhang¹, Yi Xu¹, Yutong Xu¹,
Zhongmou He^{1,2}, Yuanyuan Shi¹, Beiya Dai¹, Yunchong Song¹, Boyi Zeng¹, Qiyuan Chen¹,
Yuxun Miao¹, Bo Xue¹, Shu Wang³, Luoyi Fu¹, Weinan Zhang¹, Junxian He⁴,
Yunqiang Zhu³, Xinbing Wang¹, Chenghu Zhou^{1,3}

¹Shanghai Jiao Tong University ²University of Michigan

³Institute of Geographical Science and Natural Resources Research, CAS

⁴The Hong Kong University of Science and Technology

lin.zhouhan@gmail.com, davendw@sjtu.edu.cn, junxianh@cse.ust.hk, zhouch@lreis.ac.cn

ABSTRACT

Large language models (LLMs) have achieved huge success for their general knowledge and ability to solve a wide spectrum of tasks in natural language processing (NLP). Due to their impressive abilities, LLMs have shed light on potential inter-discipline applications to foster scientific discoveries of a specific domain by using artificial intelligence (AI for science, AI4S). In the meantime, utilizing NLP techniques in geoscience research and practice is wide and convoluted, contributing from knowledge extraction and document classification to question answering and knowledge discovery. In this work, we take the initial step to leverage LLM for science, through a rather straightforward approach. We try to specialize an open-sourced LLM into geoscience, by further pre-training the model with a vast amount of texts in geoscience, as well as supervised fine-tuning (SFT) the resulting model with our custom collected instruction tuning dataset. These efforts result in a model **GEOGALACTICA** consisting of **30 billion parameters**. To our best knowledge, it is the largest language model for the geoscience domain. More specifically, GEOGALACTICA is from further pre-training of Galactica – a top-performing LLM trained with a large number of scientific documents. We train GEOGALACTICA over a geoscience-related scientific text corpus containing 65 billion tokens, preserving as the largest geoscience-specific text corpus. Then we fine-tune the model with 1 million pairs of instruction-tuning data consisting of questions that demand professional geoscience knowledge to answer. We validate GEOGALACTICA on various geoscience examinations and geoscience-related open-domain questions evaluated by a group of senior geoscientists. GEOGALACTICA demonstrates the state-of-the-art performance in a diverse range of NLP tasks in geoscience, as well as revealing the potential of using geoscience-related tools. In this technical report, we will illustrate in detail all aspects of GEOGALACTICA, including data collection, data cleaning, base model selection, pre-training, SFT, and evaluation. We open-source our data curation tools and the checkpoints of GEOGALACTICA during the first 3/4 of pre-training in <https://github.com/geobrain-ai/geogalactica>[§].

Keywords Geoscience Language Model · Generative AI · Academic Language Model

*Zhouhan Lin is the corresponding author (lin.zhouhan@gmail.com).

[†]Version: v2 (major update on April 10, 2024).

[‡]For detailed author contributions, please refer to Appendix L.

[§]For all the checkpoints during the 3/4 pre-training can be accessed on [geobrain-ai/geogalactica-ckpt](https://github.com/geobrain-ai/geogalactica-ckpt). One can apply for the download links for further research and investigation, we will have a strict verification on the usage of the checkpoints for responsible AI principles.

Contents

1	Introduction	5
2	Related Work	7
2.1	Machine Learning in Geoscience	7
2.2	Natural Language Processing in Geoscience	7
2.3	Domain-specific Large Language Model	7
3	Preliminary and Vocabulary	8
4	Data Collection and Cleaning	9
4.1	The Customized Pre-training dataset: GeoCoprux	9
4.2	The Customized SFT dataset: GeoSignal Version 2	11
4.2.1	Domain General Natural Language Instruction	12
4.2.2	Restructured Knowledge-intensive Instruction	13
4.2.3	Self-Instruct	14
5	Training	14
5.1	Further Pre-training	14
5.2	Supervised Fine-Tuning	17
5.3	Tool Learning	18
6	Evaluation	20
6.1	Automatic Evaluation	20
6.1.1	GeoBench	20
6.1.2	MMLU	20
6.2	Human Evaluation	21
6.2.1	Noun Definition	23
6.2.2	Beginner Level Q&A	23
6.2.3	Intermediate Level Q&A	24
6.2.4	Advanced Level Q&A	24
6.2.5	Knowledge-based associative judgment question	24
6.2.6	Research Paper Titling Task	25
6.2.7	Geoscience Research Functionality	25
7	discussion	26
7.1	The Necessity of Pre-training	26
7.2	The Necessity of Further Pre-training	27
7.3	Carbon Emissions	27
7.4	Towards Unified Foundation Model in Geoscience	27
8	Conclusion	28

A Appendix: Progression of geoscience with AI	33
B Appendix: GeoCorpus	33
C Appendix: GeoSignal V2 Curation	34
C.1 MinDat	34
C.2 USGS	36
C.3 NGDB	37
C.4 Fossil Ontology	39
C.5 Fossil calibrations	39
D Appendix: Prompts	41
E Appendix: Training setup	46
F Appendix: Model Card	47
G Appendix: Evaluation	48
G.1 Open-ended Tasks	48
G.1.1 Noun Definition	48
G.1.2 Beginner Level Q&A	48
G.1.3 Intermediate Level Q&A	49
G.1.4 Advanced Level Q&A	49
G.2 Functional Tasks	49
G.2.1 Knowledge-based associative judgment question.	49
G.2.2 Research Paper Proposition Task.	50
G.2.3 Geoscience Research Functionality	50
H Generation Examples	51
H.1 Noun Definition	51
H.2 Beginner Level Q&A	51
H.3 Intermediate Level Q&A	52
H.4 Advanced Level Q&A	53
H.5 Knowledge-based associative judgment question	53
H.6 Research Paper Titling Task	54
H.7 Geoscience Research Functionality	54
I Appendix: Tool Learning Use cases	56
J Appendix: GEOGALACTICA Generation	58
J.1 Example Research Papers Written by GEOGALACTICA	58
J.2 Example Opinions Written by GEOGALACTICA	59
J.3 Example Summary of Scientific Articles Written by GEOGALACTICA	60

K Appendix: Lessons and Progresses	61
K.1 Phase 1: Prepare for Training on HPC	61
K.2 Phase 2: Training on HPC	62
K.3 Summary	65
L Membership and Contributions	67
L.1 Data preparation	67
L.2 Model Training	67
L.3 Model Evaluation and Application	67
L.4 Manuscript Writing	67
L.5 Project Management	67
L.6 Evaluation Team	68
L.7 Illustration in Arts	68
L.8 HPC Sponsor	68



Figure 1: The overview of the processing, construction, components, and applications of GEOGALACTICA.

1 Introduction

The rapid advancement of Large Language Models (LLMs) has ushered in a transformative era in natural language processing (NLP), where these models have exhibited remarkable capabilities across a wide spectrum of tasks and domains. These advanced AI models have demonstrated their prowess in handling diverse natural language tasks, including reading comprehension, open-ended question answering, code generation, etc. Their ability to harness vast amounts of general knowledge and apply it to solve specific challenges has sparked interest in exploring their potential applications in various scientific disciplines. In this context, the intersection of artificial intelligence (AI) and science, often referred to as AI for Science (AI4S), has emerged as a promising frontier for driving scientific discoveries and innovations.

Within the realm of AI4S, one particularly intriguing avenue is the integration of NLP techniques into geoscience research and practice. Geoscience is a comprehensive discipline encompassing fields such as geophysics, geology, meteorology, environmental science, etc., with a primary focus on unraveling the complexities of natural processes and phenomena on Earth. Traditionally, geoscientists have relied on theoretical and empirical approaches to advance their understanding of the Earth’s systems. However, the sheer volume of data generated in contemporary geoscience research necessitates new strategies and tools for knowledge discovery. The integration of computer science methodologies and AI technologies into geoscience has thus emerged as a transformative paradigm, offering the potential to accelerate scientific progress and address pressing global challenges effectively. In an era characterized by global challenges such as climate change and natural disaster mitigation, the need for efficient data acquisition, information sharing, and knowledge dissemination in geoscience has never been more critical.

In the field of geoscience, domain-specific geoscientific knowledge is usually presented in various forms of text data, such as scientific literature, textbooks, patents, industry standards, etc., which traditionally require the utilization of knowledge systems [1], knowledge graphs[2], or semantic models [3] to extract a structured form of these knowledge. More broadly, applying NLP techniques for geoscience use cases has been widely accepted [4], ranging from less complex tasks such as document classification [5], topic modeling [6], and entity recognition[7, 8], to more complex tasks such as knowledge graph construction [9], question answering [10] and summarization [11].

While general domain LLMs like Galactica [12], LLaMA [13], and GLM [14] have achieved impressive performance across various NLP tasks, they lack the domain-specific knowledge required for geoscience applications. These models have been trained on general datasets that lack authoritative geoscience-related data, limiting their adequacy in addressing the unique challenges posed by the geoscience domain. Although our recent attempt to adapt the LLaMA-7B model for geoscience using geoscience-specific data, i.e. the K2[15] model, has shown promising results, this primitive attempt is constrained by its model size and data scale, which consequently may not fully capture the complexity of geoscientific terminology and concepts. However, training a larger LLM comes with new technical challenges, since many aspects of the process become fundamentally different as the model scales up. For example, the stability of training will become more vulnerable, and the training data needs to be scaled up accordingly, resulting in a more systematic way of managing different data sources, etc.

Therefore, tailoring a general, larger LLM for the scientific domain of geoscience with a more systematically designed dataset and training pipeline is imperative in this era of LLMs. In response to these necessities, this work presents a significant step forward in the development of the model as well as the set of toolchains around it.

Leveraging the vast amount of resources of scientific literature’s meta-data, particularly the data resources collected for the OpenAlex ¹, Web of Science ², Semantic Scholar ³, and Acemap ⁴, we can create, organize, and manage a large and comprehensive geoscience dataset targeted for all stages in large language model training. In particular, we have introduced GAKG [2], Deep Literature⁵, GSO⁶, and other platforms as carriers and repositories of geoscience text knowledge. These concerted efforts have not only allowed us to accumulate a comprehensive geoscience data archive but also have served as foundations for constructing an extensive instruction-tuning dataset for geoscience-related questions, **GeoSignal-v2**, which has been employed in supervised fine-tuning (SFT). In addition, we have developed and customized a series of data-cleaning tools that allow us to automatically convert various forms of raw data, such as PDF files, forms, equations, knowledge graphs, etc., into clean texts suited as training corpus for large language models. To our best knowledge, our collected corpus has become the largest geoscience dataset.

¹<https://openalex.org/>

²<https://www.webofscience.com/wos/>

³<https://www.semanticscholar.org/>

⁴<https://www.acemap.info>

⁵<https://idea.acemap.info/>

⁶<https://gso.acemap.info/>

we have then successfully further pre-trained a language model with **30B** parameters, with Galactica-30B [12] as its base model. The resulting model is thus named as GEOGALACTICA, empowering various academic tasks in the geoscience field. With its 30 billion parameters, this model represents the culmination of further pre-training and supervised fine-tuning, making it the largest language model dedicated to the geoscience domain. Our experimental findings demonstrate that, compared to models of equivalent scale, GEOGALACTICA exhibits exceptional performance on GeoBenchmark [15]. Regarding human evaluation, our model showcases impressive competence in geoscience-related tasks when compared with 5 general language models (*ChatGPT*⁷, *Yiyan*⁸, *Qianwen*⁹, *MOSS*¹⁰, *ChatGLM*¹¹).

Moreover, since our GEOGALACTICA model provides a unified representation space and computational approach for diverse geological data described in various textual formats, it holds tremendous potential in narrowing the technological gap between different earth science tasks.

In the subsequent sections of this technical report, we will provide a detailed description of the data collection and cleaning processes, base model selection, pre-training, supervised fine-tuning, and extensive evaluations in the creation of GEOGALACTICA. Additionally, we are committed to promoting open science by making our data curation tools and pre-training checkpoints available to the research community through our GitHub repositories¹².

Broad Contribution

In addition to establishing the academic mega-model in geoscience, our goal is to contribute to a broader research community. Specifically, the experiences documented in this paper provide evidence for further community understanding of several open questions in the literature. **Warning: The model in this manuscript might produce hallucinations and reader discretion is recommended.**

1. **A Domain-specific LLM:** Our construction of GEOGALACTICA, following in the footsteps of our previous work K2 [15], represents a geoscience LLM that focuses on interacting with humans and generating contents on highly professional academic topics.
2. **A Toolchain for Data Cleaning:** A high-quality training dataset is crucial for successfully training large language models. Therefore, our contribution to the community includes developing an efficient academic data preprocessing toolchain to construct a clean training corpus from PDF documents¹³
3. **Primitive Explorations to Use Tools:** As for training GEOGALACTICA to use tools, we also construct a set of supervised data *Geotools* for training GEOGALACTICA to use tools. We also open-source the codes and data on Github.¹⁴
4. **Training Details and pre-training Checkpoints:** We conducted model training on the accelerator hardware provided by the Advanced Computing East China Sub-center. We will describe in detail the pre-training and SFT processes in the remainder of this paper. In addition, we are releasing the training checkpoints during the first 3/4 of the pre-training process on Hugging Face.¹⁵
5. **Model and data analysis process:** In building a domain-specific LLM, the model and the data should be effectively evaluated and analyzed. We provide a set of analysis and visualization methods for the SFT data and the weights of the GEOGALACTICA, open-sourced on Github.¹⁶

In summary, we aim to contribute to the research community by developing the GEOGALACTICA model and providing insights and tools related to data construction, training processes, and evaluation strategies. The organization of the paper can be seen in the contents section listed above.

⁷<https://chat.openai.com/>

⁸<https://yiyao.baidu.com/>

⁹<https://qianwen.aliyun.com/>

¹⁰<https://moss.fastnlp.top/>

¹¹<https://chatglm.cn/>

¹²The list of related tools, data, and codes can be found in <https://github.com/geobrain-ai/geogalactica>

¹³The toolchain is open-sourced on Github repos: https://github.com/Acemap/pdf_parser and <https://github.com/davendw49/sciparser>. In addition, an online demo of this toolchain can be found at <https://sciparser.acemap.info/>.

¹⁴<https://github.com/zthang/geotools>

¹⁵<https://huggingface.co/geobrain-ai/geogalactica>

¹⁶https://github.com/dbylynn/GeoGalactica_Analysis

2 Related Work

2.1 Machine Learning in Geoscience

With the advancement of artificial intelligence, utilizing machine learning, natural language processing, and recent large-scale model techniques to tackle critical problems in geoscience has become a crucial direction. Various subtasks in geoscience involve significant data collected from sensors, making them suitable for end-to-end learning using machine learning approaches. Some studies model multiple aspects of seismic signals using deep learning models to extract information relevant to earthquake prediction. Among them, [16] uses supervised learning with end-to-end training, while [17, 18] employs self-supervised learning to obtain models applied to downstream tasks. [19, 20] utilize machine learning to explore the latent correlations among different rock properties for rock type prediction. Beyond relatively straightforward classification tasks, there are numerous works applying machine learning to address more complex scenarios in geoscience, such as calculating wellhead flow rate [21], capturing and storing carbon [22], and predicting the condition of SPBM tunnels [23]. Additionally, machine learning is introduced to evaluate the real-world environment: [24] explores the use of Few-Shot Learning (FSL) methods to enhance the accuracy of high-resolution pre-stack seismic inversion, and [25] employs various machine learning techniques and ensemble methods to improve landslide hazard prediction, demonstrating their high practical value. Machine learning is also being used to aid geoscience exploration, [26] attempts to use machine learning to do data-driven modeling of solid earth science, [27] attempts to use machine learning to reveal the link between fast and slow earthquakes, [28] uses machine learning to reveal the impact of aerosols on climate impact.

2.2 Natural Language Processing in Geoscience

In addition to the diverse and heterogeneous data collected from various sensors, the field of geoscience also encompasses a significant amount of text data with standardized formats. The application of natural language processing (NLP) in earth science has witnessed remarkable progress. [29, 6] embed different sources of textual information into a unified space, [29] employs joint training of language models with text and points of interest (POI) for POI retrieval, while [6] integrates geological attribute information into the textual representation space to enable better knowledge extraction. [30, 31] enhance language models with knowledge graph techniques, where [30] constructs a knowledge graph on geological text to discover ore-forming environments, and [31] proposes an automatic entity and relation extraction approach via three-level extraction to build a geological knowledge graph from extracted information in geological reports. [32] combines retrieval techniques with language models creates an integrated solution incorporating contextual retrieval and the GeoBERT model. [33] focuses on various language irregularities encountered in natural language texts, introducing the NeuroSPE model for spatial extraction using neural networks. NLP techniques provide a unified representation space and computational approach for diverse geological data described in various textual formats, narrowing the technological gap between different earth science tasks.

2.3 Domain-specific Large Language Model

The recent emergence of large-scale language models marks a significant step towards unified information processing in geoscience. These models are pre-trained on vast amounts of text data and efficiently compress all input data. Currently, in addition to earth science, various domains have seen the development of domain-specific pre-trained models trained on domain-specific corpora. [34, 35, 36, 12, 37, 38] performs large-scale pre-training on domain-specific texts and has resulted in foundational models equipped with domain knowledge, while [39, 40, 41] fine-tuning these base models using domain-specific data, achieving models tailored to specific downstream tasks at a lower cost. These works have made significant strides in developing domain-specific LLMs through dedicated data integration and model training efforts. Recently, [42, 43, 44] explored the use of prompt engineering to unlock the potential of models without additional training, offering the possibility of unifying various geoscience tasks and further reducing the cost of employing large models in domain applications. In the field of geoscience, the exploration of large models is still in its early stages. [15] collected a substantial amount of high-quality data from earth science Wikipedia and research papers, and further fine-tuned the base model, leading to impressive scientific competence and knowledge in earth science. For the first time, our work utilizes a large corpus of earth science papers and textbooks, which were cleaned using a dedicated toolchain for constructing large-scale earth science models, ensuring data quality. Furthermore, our work completes the entire process of “*further pre-training, supervised fine-tuning, augmented learning*” for large foundation models for geoscience, bringing the largest scale and highest quality proprietary language models to the geoscience field. This will unlock tremendous possibilities for future research conducted by earth science researchers.

We have outlined the progression of geoscience research with the use of cutting-edge AI techniques, including neural network (NN), K-nearest neighbor (KNN), recurrent neural network (RNN), convolutional neural network (CNN), backpropagation (BP), reinforcement learning (RL), support vector machine (SVM), long-short term memory (LSTM),

graph convolutional neural network (GCN), Transformers, BERT, ChatGPT, and large language model (LLM). [45] The investigation reveals that the time intervals between AI technology advancements and their application in geoscience have significantly shortened, indicating an increasing reliance on advanced AI technology in the field of geoscience. The illustration is presented in Figure 2, and detailed information about the progression is shown in Appendix A.

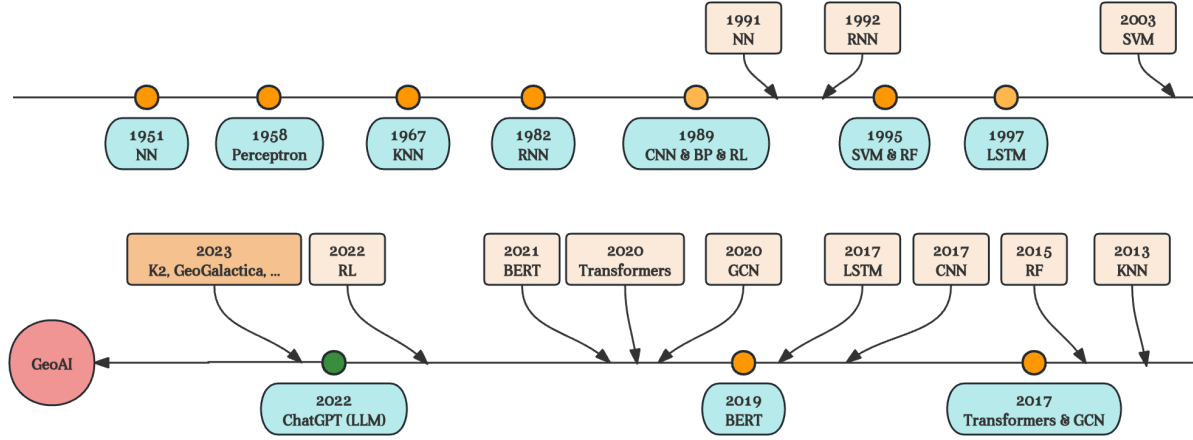


Figure 2: The progression illustration of geoscience research with the use of cutting-edge AI techniques. The textboxes in *PaleTurquoise* show the techniques from computer science, The textboxes, in *Bisque* show the research that probably the first time geoscientists used the techniques.

3 Preliminary and Vocabulary

To facilitate understanding of our work and the overview of our model, here are some key terms and their corresponding explanations that will be widely used in the narrative of this article.

Vocab	Stage and illustration
Galactica-30B	The vanilla Galactica model
GeoGalactica-FP	Checkpoint after pre-training over geoscience data (Further Pre-train)
GeoGalactica-Alpaca	Applying supervised fine-tuning with Alpaca data on top of GeoGalactica-FP
GeoGalactica-GeoSignal	Applying supervised fine-tuning with GeoSignal data on top of the first checkpoint (GeoGalactica-FP)
GeoGalactica	Applying supervised fine-tuning following training recipe of K2 on top of the first checkpoint (GeoGalactica-FT)
GeoCorpus	Geoscience text corpus for pre-training
GeoSignal	Supervised fine-tuning data for geoscience
GeoBench	Benchmarks for evaluating the performance of the geoscience LLM

Table 1: Vocabulary for this technical report.

Here we list the terms widely used in this report:

- **Sciparser**. A PDF parsing toolkit for preparing text corpus to transfer PDF to Markdown.
- **GeoTools**. A set of supervised instruction data for training GEOGALACTICA to use tools.
- **K2**. The first-ever geoscience large language model trained by firstly further pre-training LLaMA on collected and cleaned geoscience literature, including geoscience open-access papers and Wikipedia pages, and secondly fine-tuning with knowledge-intensive instruction tuning data (GeoSignal).
- **Deep Literature**. Deep Literature is a literature platform, aiming to construct a knowledge information system for geoscience scholars, which, step-by-step goes through knowledge arrangement, knowledge mining and knowledge discovery.
- **GAKG**. GAKG [2] is a multimodal Geoscience Academic Knowledge Graph (GAKG) framework by fusing papers’ illustrations, text, and bibliometric data.

- **DeepShovel**. DeepShovel [46] is an Online Collaborative Platform for Data Extraction in Geoscience Literature with AI Assistance.
- **GSO**. Similar to WordNet, **GeoScience Ontology**, (GSO) is a hierarchical tree of geological terms contains a vast amount of synonyms and word explanations, providing valuable geoscience connections between terms.
- **Acemap**. AceMap is a platform displaying relationships among academic publications and a search engine for academic literature reviews.
- **DataExpo**. A one-stop dataset service for geoscience research.

Finally, we share the model card in Appendix F.

4 Data Collection and Cleaning

The training corpus of Galactica primarily consists of literature related to computer science and biochemistry rather than earth science. This suggests that Galactica may lack sufficient knowledge in the field of geoscience. To address this, we have collected approximately **six million** research documents specifically focused on earth science. These papers were carefully selected by professional experts in the field. Furthermore, we have expanded the GeoSignal dataset based on K2 to better support natural language processing tasks in earth science research. This expanded dataset was used for fine-tuning the model after further pre-training. In the following sections, we will provide a detailed explanation of how our dataset was constructed.

4.1 The Customized Pre-training dataset: GeoCopus

According to our long-term collection efforts on geoscience papers, with the research fields subfields of geology and geography, through Acemap, we have accumulated a total of 5,980,293 papers.

During this process, we commenced our data collection in early 2020 by gathering a list of journals in geoscience from LetPub¹⁷. We identified the corresponding publishers’ websites using the journal names and ISSN to collect open-access data through web page parsing. This process continued towards 2023 when we collaborated with experts in various sub-disciplines of geoscience, we collect paper from high-quality SCI journals in the mathematical geosciences, metamorphic petrology, geochronology, geomagnetism and paleomagnetism, geomorphology, tectonics, stratigraphy, hydrogeology, geophysical, geothermics, igneous and geochemistry, surficial geochemistry, geological mapping, sedimentological, petroleum geology, paleontology, paleogeography, and mineralogy. In total, we integrated a list of **849** geoscience-related journals (Appendix B shows the distribution of the collected papers in geoscience).

We employed the journal name list to search for journal information and their publishers’ websites. Through web page scraping, we collected HTML pages and subsequently conducted data parsing to extract metadata from papers. For open-access articles, we matched them with the parsed DOI and the corresponding PDF from Sci-Hub¹⁸. If no PDF was available, we downloaded it based on the URL.

Throughout this process, we adhered to the network conventions of information acquisition. When faced with obstacles such as anti-scraping measures like 5-second shields, JavaScript encryption, IP proxy bans, and account logins, we constrained our actions to ensure the compliance of our data. Moreover, data security remained our utmost priority during this process; thus, we refrained from publicly disclosing the data obtained during this stage.

In conclusion, we obtained a total of 5,980,293 papers. Our data collection system operated through a distributed information fusion mechanism, utilizing an 8-workstation k8s cluster. Data collection was conducted using Scrapy-Redis¹⁹ framework. Additionally, we implemented compression techniques for HTML data to address challenges related to large-scale data storage.

Furthermore, we have leveraged the copyrights obtained from the publishers we have been collaborating with over the years to parse and anonymize the PDFs of these articles, creating a dataset of textual data. Additionally, referring to [47, 48], we have reason to believe that the inclusion of program code in the model’s pre-training, alongside the text, can significantly enhance the reasoning capabilities of the LLM model.

Therefore, after collecting datasets from Acemap and ArXiv, we incorporated the training dataset from Codedata. Finally, our overall training corpus is detailed in Table 2, totaling **78B**. The data from a specific source is concatenated into a single record. After tokenization, we then split it according to a block size of 2048, with each instance ending

¹⁷<https://letpub.com.cn/>

¹⁸<https://sci-hub.se/>

¹⁹<https://github.com/rmax/scrapy-redis>

with the *tokenizer.eos* token. For each training batch, the proportion of geoscience papers to the other two datasets is 8 : 1 : 1.

Notice: All the data employed in this manuscript is derived from web pages that are publicly available, thereby introducing a potential bias in the reliability of the data. Users are cautioned to be mindful of potential hallucination problems that may occur when utilizing large-scale models. Furthermore, this paper adheres to all copyright concerns. Should any issues arise, stakeholders are encouraged to notify the authors.

Dataset	#blockNum	#tokenNum	#itemNum	#tokenSize	#batchRatio
<i>GeoCorpus</i>	25,743,070	52,721,798,004	5,548,479	98.21G	80%
<i>ArXiv</i>	6,691,886	13,704,981,558	742,835	25.53G	10%
<i>Codedata</i>	6,066,725	12,424,652,670	3,456,887	23.14G	10%
Total	38,501,681	78,851,432,232	9,748,201	146.88G	-

Table 2: Data distribution of the corpus used for training GEOGALACTICA

We utilized data processing and enhancement tools based on DeepShovel and K2 during this process. With the help of Grobid [49] and pdffigure2 [50], we provided a comprehensive parsing solution for extracting text, images, tables, formulas, and other data types from research documents. This was further enhanced by DeepShovel for parsing tables and formulas, resulting in the development of the SciParser tool. We plan to open-source this tool and share it on GitHub.

Within PDF documents, there are various types of data, including text, images, tables, and formulas, all organized within the articles’ hierarchical structure and page layout. Data preprocessing is necessary to extract and ensure the readability of such content. It entails utilizing a PDF parsing tool to perform an initial parsing of the PDF document, resulting in a parsing file that contains various information from the document. However, the readability of this file is often poor, and it may have a significant amount of redundant information. Subsequently, the parsing file needs to undergo data cleansing, extracting the desired text, images, tables, formulas, and other data, and converting it into Markdown format for further processing, analysis, or display purposes.

Currently, we are utilizing Grobid as our PDF parsing tool. Grobid can accurately extract text from PDF documents and perform structured extraction of articles. It provides an outline of the text, forming an XML structure that allows for restoring the original PDF layout. Additionally, Grobid enables precise localization of images, tables, formulas, and other data types. With the provided bounding boxes, we can obtain the corresponding images using the PyMuPDF tool²⁰. Further leveraging the OCR recognition integrated into DeepShovel [46], we can convert tables, formulas, and other elements into Markdown format. The parsing process is completed by writing all the parsed content into a markdown file for further use. Throughout the entire process, for tables, we utilize the DeepShovel PDF Table Parser²¹. This tool ensures the completeness and accuracy of the table content while preserving the table structure, making it convenient to reconstruct tables using Markdown. As for formulas, we employ an improved version of Latex-OCR²² for the recognition, converting the parsing results into the string format. We open-source our PDF parsing solution on GitHub²³.

Tokenization is a crucial component of text corpus construction. To aid language models in comprehending academic papers, we utilize dedicated tokens for different types of special data. Finally, we use special tokens similar to the original Galactica paper [12] to unify various forms of text extracted from various sources into one standard protocol. Below is an explanation of our special tokens.

- **Figures:** We use the special tokens [START_FIGURE] and [END_FIGURE] to mark the captions of figures in the paper.
- **Tables:** The special tokens [START_TABLE] and [END_TABLE] are employed to identify the position of tables within paragraphs. During this process, we convert tables from the PDF into Markdown format.
- **References:** We use the special tokens [START_REF] and [END_REF] to annotate citations. The title of the article is placed within these special tokens to enhance readability.

²⁰<https://github.com/pymupdf/PyMuPDF>

²¹<https://github.com/ShaoZhang0115/Table-Extraction-for-Geoscience-Literature>

²²<https://github.com/lukas-blecher/LaTeX-OCR>

²³https://github.com/Acemap/pdf_parser

- **Formulas:** For mathematical content or formulas, we employ regular expressions and rule-based methods to filter and clean irregular formulas parsed from the PDF. Additionally, we use the special tokens [START_FORMULA] and [END_FORMULA] to capture them.

And the dedicated tokens for these different types of special data are shown in Figure 3 (We use the one in [15])



[START_FIGURE]

Fig. 2. SEBLUP estimates of perceived disorder in Manchester (division in 6 quantiles).

[END_FIGURE]

A. Figure Processed Text.

Simulation is an established practice in the imaging radar community that serves many purposes, including inSAR accuracy assessment [1], [2], analysis of new algorithms [3], [4] and new sensor designs [5]–[7], and improvement of image interpretation algorithms [8]–[10].

Simulation is an established practice in the imaging radar community that serves many purposes, including inSAR accuracy assessment [START_REF]Terrain height measurement accuracy of interferometric synthetic aperture radars [END_REF], [START_REF]A time-domain raw signal simulator for interferometric SAR[END_REF], analysis of new algorithms [START_REF]Application of angular correlation function of clutter scattering and correlation imaging in target detection[END_REF], [START_REF]Clutter effects on ground moving target velocity estimation with SAR along-track interferometry[END_REF] and new sensor designs [START_REF]Multi-frequency and multipolarization SAR system analysis with simulation software developed at CSA[END_REF]– [START_REF]Spatial considerations in SAR speckle simulation[END_REF], and improvement of image interpretation algorithms [START_REF]Radar image simulation[END_REF] – [START_REF]Enhanced simulation of radar backscatter from forests using LiDAR and optical data[END_REF].

C. Citation Processed Text.

Table 5 Summary of covariates and coefficients of correlation of each variable with direct estimates of perceived disorder.									
	min	1st quartile	Median	Mean	3rd quartile	Max	Spearman coeff		
Proportion BME	0.03	0.15	0.25	0.31	0.42	0.89	0.31	a	
Proportion unemployed	0.04	0.11	0.16	0.17	0.23	0.64	0.29	a	
Income deprivation	0.01	0.12	0.24	0.34	0.52	0.93	0.35	a	
Population churn	0.21	0.33	0.39	0.39	0.44	0.68	0.13	a	
Mixed land-uses	0.00	0.19	0.31	0.31	0.44	0.88	0.13	a	

[START_TABLE]

Table 5

Summary of covariates and coefficients of correlation of each variable with direct estimates of perceived disorder.

[None][Min][First quartile][Median][Mean][Third quartile][Max][Spearman coeff]
[a][b][c][d][e][f][g][h][i][j]
[Proportion BME][0.03][0.15][0.25][0.31][0.42][0.89][0.31 a]
[Proportion unemployed][0.04][0.11][0.16][0.17][0.23][0.64][0.29 a]
[Income deprivation][0.01][0.12][0.24][0.34][0.52][0.93][0.35 a]
[Population churn][0.21][0.33][0.39][0.39][0.44][0.68][0.13 a]
[Mixed land-uses][0.00][0.19][0.31][0.31][0.44][0.88][0.13 a]
[a p-value < 0.01, *p-value < 0.05, [None][None][None][None][None][None][None][None]

[END_TABLE]

B. Table Processed Text.

$$I = sn \quad (2)$$

with a speckle pdf given by

$$P(n) = \exp\{-n\}. \quad (3)$$

[START_FORMULA]I = sn (2)[END_FORMULA]

with a speckle pdf given by

[START_FORMULA]P (n) = exp{-n}. [END_FORMULA]

D. Formula Processed Text.

Figure 3: Tokenization processed text. A. shows an example of a figure marker, we only choose to preserve the captions; B. shows an example of a table marker, we transfer the tables into the form of Markdown; C. shows the tokenization of the citations, we replace the reference numbers into reference papers’ title to preserve the readability of the text corpus; D. shows an example of the special tokens for formulas.

4.2 The Customized SFT dataset: GeoSignal Version 2

Through extensive investigation and research, we have thoroughly explored natural language processing tasks specifically tailored to geoscience. In this endeavor, we have identified a set of tasks that cater to the unique requirements of geoscience applications. However, during this process, we have observed numerous unsupervised signals within these tasks that have yet to be fully harnessed and summarized.

- **Geoscience Knowledge Graph:** Named entity recognition (NER) for temporal scales, rock types, etc., relation extraction (RE) for linking knowledge points, text-to-graph transformation, and knowledge discovery through reasoning
- **Academic Applications:** Keyword extraction, summarization, and information retrieval.
- **General Applications:** Question and Answering (Q&A), conversations related to geoscience education, and text classification.
- **Geographical Applications:** Point of Interest (POI) queries and multimodal Q&A.

However, the supervised signals for these tasks can be reconstructed using professional geoscience data websites. Based on the data scheme provided by K2, we further elaborate on the entire data construction process. In this process, we have built three categories of data:

1. Literature-related data can be used to construct general natural language instructions, enabling the model to possess basic semantic reasoning abilities.
2. Geoscience-related data, which is used to build a knowledge-intensive instruction dataset, allowing the model to understand and comprehend the specific forms of natural language tasks in the field of geoscience.
3. Self-instruction-related data, following the examples of Alpaca [51] and Baize [52], we have distilled geoscience-related data from ChatGPT and invited geoscience experts to annotate it. This data is used to construct high-quality geoscience question-answering datasets.

4.2.1 Domain General Natural Language Instruction

For the general instruction learning data, we have integrated four platforms constructed by Acemap, and reconstructed the data accordingly.



Figure 4: Four platforms that contribute most to our GeoSignal.

Referring to RST [53] and K2 [15], we restructure the signals from various geoscience-related platforms. The following paragraphs will provide a detailed explanation for each platform and the illustrations for restructured domain-general natural language instruction.

Deep Literature and DataExpo. These two platforms can be understood as collections of papers and datasets. Therefore, the Related Paper (with abstract) and Reference resolution of Deep Literature, as well as the Reference resolution of DataExpo, serve as excellent datasets for establishing referential relationships.

Using the text processing tool mentioned earlier, we explicitly employ a multi-threaded Grobid to process all documents and convert them into an intermediate XML format. Within the converted XML, we identify the `bibl_id` of in-text citations and then locate the corresponding reference paper titles in the XML’s reference section.

GSO. Similar to WordNet, the hierarchical tree of geological terms contains a vast amount of synonyms and word explanations, providing valuable supervised signals. As a result, we traverse all the nouns in GSO, extract all the synonyms for each term, and combine them with the term itself to create a set. We then construct all possible pairs of (term, synonym_term) and add them to a list of results.

For the word description, we traverse all the nouns of GSO, extract the definition of the respective noun to serve as the description and create signal pairs (word, description). Additionally, there is also a specialized geology dictionary, which includes a dataset of categorized geology terms. The original data is in PDF format, and we convert it into JSON format through PDF parsing. In this process, we first use a parsing tool to convert the PDF into a docx format, and then use a data processing script to convert its content into JSON format. Subsequently, we proceed with content processing, removing hyphens at line breaks, and merging multiple definitions of a single term. GSO uses two geoscience dictionaries. For the geology dictionary, each entry consists of a "name" and "description". For the geography knowledge dictionary, it includes one more "attribute" field.

GAKG. GAKG is rich in images, tables, and other elements from geology papers. Meanwhile, the text describing these images and tables, as well as their captions, can serve as excellent sequence-to-sequence (seq2seq) supervised data. Regarding the papers and their graphical information, four types of binary pairs can be generated. During this process,

we transform the original text of the paper, tables, and illustrations in PNG format along with their corresponding captions, including table numbers and contents, into the target data format: (*illustration caption, illustration content*), (*illustration caption, referring sentence*), (*table caption, table content*), (*table caption, referring sentence*). For detailed information regarding this specific aspect, please refer to the Appendix.

Our approach to handling this is as follows:

1. The captions and contents of tables and illustrations are stored in separate JSON files within their respective folders and can be extracted from there.
2. The referring sentences, on the other hand, need to be retrieved from the original text of the paper by referencing the table/illustration numbers mentioned in their captions.

Specifically, we search for the keywords “fig” (or variations like “Fig” and “FIG”) and “table” (or “Table” and “TABLE”) in the original text and identify the associated numbers (i.e., “i”) immediately following them. We then search for complete sentences between two periods preceding and following these numbers.

Our program handles some unexpected scenarios, such as excluding cases like “Fig11” or “Fig12” when searching for “Fig1,” and partially excluding cases where the confusion in numbering arises from referring to tables/illustrations from other papers. We also consider disorders caused by the dot used in English sentences and abbreviations, among other cases.

However, there are still a few limitations to this method:

1. When the keywords “fig” or “table” appear at the end of a sentence, our program includes both that sentence and the subsequent one as the corresponding referring sentence.
2. There might be instances where figures/tables from other papers are referenced. Our program can identify such cases if:
 1. The figure/table numbers are more significant than the current paper’s total number of figures/tables.
 2. The word “of” appears close after “Fig” in the text.

In scenarios where it is difficult to discern whether a referenced figure/table belongs to another paper, we prioritize data quality. If we encounter any unmatched or garbled text, or the text is concise, we will discard that particular supervisory signal.

Wikipedia. Wikipedia contains a lot of crowd-sourcing information for geoscience. Consequently, we have also incorporated geoscience data from the Wikipedia page. To retrieve the information, we utilized web scraping techniques and relevant libraries.

For the article’s title, we used the Wikipedia library in Python²⁴, which supports accessing sections of a Wikipedia page. Each section’s title, text, and sub-sections form a nested structure. By recursively traversing each page, we obtained a list of triplets comprising each section’s level, title, and paragraph. The triplets are structured as (level, title, paragraph), where the level indicates the depth of nesting, the title represents the section’s title, and the paragraph contains the corresponding text content.

To retrieve the “Summary & Abstract” of the article, we utilize the Wikipedia library in Python to access the abstract of the corresponding Wikipedia page directly. We then concatenate the paragraphs from the abovementioned sections to form the full text. Finally, we output the tuple (full text, abstract).

To extract the Entity mentioned in the article, we use the requests library and the BeautifulSoup²⁵ library to scrape the Wikipedia page directly. We retrieve the text from all tags labeled “p” and “ul” and treat them as paragraphs. Next, within these paragraph tags, we search for tags labeled “a” with a href attribute starting with “/wiki/”. These represent the highlighted blue hyperlinked sections within the text. We collect these entities and output the tuple (*paragraph, entities*).

4.2.2 Restructured Knowledge-intensive Instruction

In our work of building restructured knowledge-intensive instruction data, we begin by searching for authoritative websites related to paleontology, dinosaurs, fossils, rocks, and other fields within geoscience. We then filter these websites, specifically selecting those with structured data available for extraction.

²⁴<https://github.com/goldsmith/Wikipedia>

²⁵<https://www.crummy.com/software/BeautifulSoup/>

Disciplines	Websites	Websites Intro.
Dinosaur	https://dinoanimals.com/dinosaurdatabase/	A comprehensive Dinosaur Database, offering a detailed catalog of dinosaurs.
Fossil	https://fossilcalibrations.org/	A specialized resource offering a curated collection of fossil calibrations.
Fossil	http://fossil-ontology.com/	A multi-dimensional scientific database of fossil specimens.
Mineral	https://ruff.info/	A dedicated resource for the study and identification of minerals.
Mineral	https://zh.mindat.org/	A comprehensive online mineralogical database.
Sedimentary	https://mrdata.usgs.gov/	A system with interactive maps and data for analyzing mineral resources on a regional and global scale.
Earthquake	https://www.usgs.gov/	A website collecting all the earthquake world wide.
Hazard	https://public.opendatasoft.com/explore/	A platform for exploring various datasets sorted by their modification date.

Table 3: Knowledge Intensive Data Sources.

For the websites that can be structured, we perform corresponding restructured processing like K2 [15]. Taking the provided image as an example, we match the structured data on the website using Key-Value pairs and create natural Instruction and Response pairs.

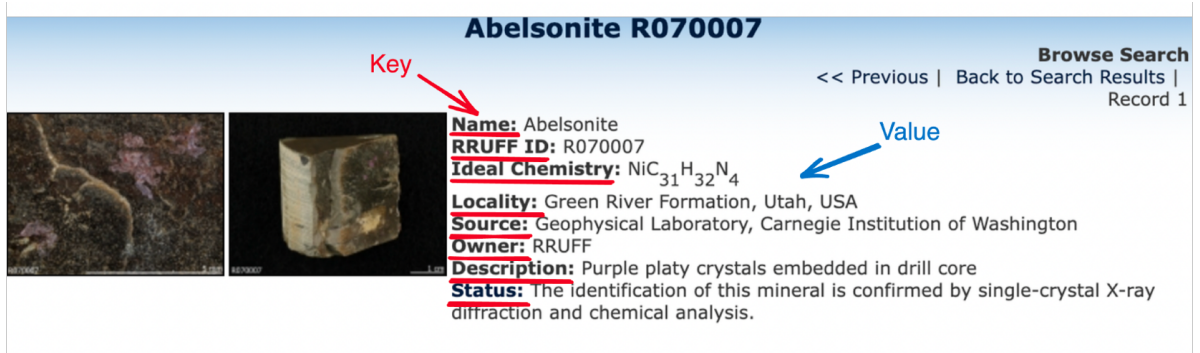


Figure 5: An example for illustrating the construction of restructured knowledge-intensive instruction data.

4.2.3 Self-Instruct

According to Alpaca [51] and Baize [52], using problem seeds to generate answer from ChatGPT²⁶ is an appropriate way to build instruction tuning data. In geoscience scenarios, we generate 1000 questions per subject under the geoscience, and we put the problem seeds on GEOGALACTICA’s Github Repo.

In terms of overall data collection, the total amount is as follows. And we select a certain proportion of data to be included in our supervised fine-tuning process. In the final version, after further manual verification and cleaning, we choose to use a dataset of **100K** samples as GeoSignal Version 2 for instructional data during the supervised fine-tuning. The detailed statistic of the instruction tuning data is shown in Table 4.

5 Training

Taking the lessons from GLM-130B [14], we design the frameworks and plans of the GEOGALACTICA. The following are the details of our progress.

5.1 Further Pre-training

After the initial pre-training by Meta AI, the model Galactica can undergo additional training on a geoscience-specific dataset. We hope this fine-tunes the model’s understanding and generation capabilities in particular domains or styles.

We utilize a supercomputing cluster based on the Hygon DCU architecture, combined with the Megatron-LM framework [54], to further pre-train our models. The computing cluster consists of 512 nodes, with each node equipped with a 32-core CPU, 128GB of memory, and 4 DCU acceleration cards, each with 16GB of memory, resulting in a total of 2048 acceleration cards, where each acceleration card is equivalent to approximately 0.2 times the computing power of an NVIDIA A100 GPU. The Megatron-LM framework employs 3D parallelism strategies, including pipeline-parallel,

²⁶In the data curation process and experiments throughout this paper, we use the 2023 March version of ChatGPT and 2023 March version of GPT-4 unless otherwise specified

		Signals	tuples	#NumofSamples
Scholar		Title (with Abstract)	(abstract; title)	2,690,569
		Abstract (with Publications Fulltext)	(fulltext; abstract)	2,601,879
		Category (with abstract)	(abstract; category)	12,321,212
		Related Paper (with abstract)	(source abstract; target abstract; reference sentence)	40,047,777
		One Sentence Summary (with abstract)	(abstract; question; answer)	2,690,569
		Reference resolution	(sentence; pronoun.; reference item) [including citation]	2,329,820
DataExpo		Title	(abstract; title)	216,036
		Summary & Abstract	(fulltext; abstract)	216,036
GAKG	GAKG	Principal Concepts	(sentence; entity; types)	3,892,102
		Relations	(abstract; sentence; head entity; relation; tail entity)	30,123
		Paper table caption	(table caption; referring sentence)	2,772,166
		Paper illustration caption	(illustration caption; referring sentence)	9,128,604
		Paper table content	(table caption; table content)	2,772,166
		Paper illustration content	(illustration caption; illustration content)	9,128,604
	GSO	Factual knowledge	(sentence; facts; improper statement)	114,392
		Taxonomy	(upper term; term)	112,298
		Synonyms	(term; synonym term)	23,018
		Word description	(word; description; source)	110,209
	GA-Dialogue	Future content and Previous content	(corrupted text; corrupted positions; target spans)	5,434
GeoOpenData	dinosaur	Factual knowledge	(property; property value)	11,348
	fossilcalibrations	Factual knowledge	(property; property value)	1,749
	fossilontology	Factual knowledge	(property; property value)	3,210
	mindat	Factual knowledge	(property; property value)	51,291
	ngdb	Factual knowledge	(property; property value)	148,212
	opendatasoft	Factual knowledge	(property; property value)	37,823
	rruff	Factual knowledge	(property; property value)	32,778
	usgsearthquake	Factual knowledge	(property; property value)	37,284
WordNet		Synonyms	(term; synonym term)	6,408
		Word description	(word; description; source)	27,123
Wikipedia		Title	(term; abstract)	3,033,595
		Summary & Abstract	(fulltext; abstract)	753,920
		Entity mentions	(paragraph; entities)	3,688,926
		Relation	(text; subject; property; object)	630,210
IODP		Title	(abstract; title)	2,839
		Summary & Abstract	(fulltext; abstract)	2,638

Table 4: GeoSignal Statistics Table.

model-parallel, and data-parallel, to maximize GPU performance while reducing communication overhead. Given the four acceleration cards per node, we set the model parallel size to 4 for optimal model-parallel efficiency. Additionally, in the case of a mini-batch size of 1, we set the pipeline-parallel size to 16 to fully utilize the memory resources.

We preprocess all training text data by performing tokenization. The tokenized results of each document are then concatenated using an *end-of-sentence* (*eos*) marker. Subsequently, we crop the concatenated sequences into fixed lengths of 2048, resulting in 30 million training samples, corresponding to 7324 training steps. Before formally starting the training, we conduct a preliminary experimental analysis of node failures and save checkpoints at intervals of 100 steps. We initiate the pre-training process after transforming the initial checkpoint format into the format Megatron-LM requires. Ultimately, after running for 16 days, the computing cluster completes the further pre-training of the model at a speed of 3 minutes per step. Due to the frequent occurrence of node failures, the actual training takes nearly a month to complete. After the pre-training, we convert the checkpoints into the Hugging Face format for subsequent applications.

Challenge in further pre-training

1. **Over-fitting:** Further pre-training may increase the risk of overfitting, especially when the training data is relatively limited compared to the original Galactica pre-training data (refer to Section 4).
2. **Catastrophic forgetting:** In Further pre-train, ensuring that the training on the initial pre-training data is not forgotten is crucial. Sudden increases in the loss of new data sources can lead to the loss of knowledge acquired from the Galactica pre-training. It is essential to address how to effectively transfer higher-level language abilities to specific tasks and prevent the loss of the model’s generality obtained during the initial pre-training during the fine-tuning process.
3. **Stability and Convergence:** Further pre-training models may be more prone to training instability and convergence difficulties. During the training process, more sophisticated optimization techniques and strategies may be required to ensure that the model converges smoothly to an appropriate state.

Parameters transformation from Galactica to Megatron GPT-2 Since Galactica belongs to the OPT model, we referred to and modified the code available on Hugging Face for converting HF GPT-2 to Megatron GPT-2. The conversion parameters can be adjusted based on the actual scale of pipeline parallelism (PP), model parallelism (MP), and data parallelism (DP) during runtime.

Training detail

- * **Training Setup:** In this study, we utilized a supercomputing cluster based on the hygon DCU architecture and combined it with the Megatron-LM framework for further pre-training of the model. The computing cluster consisted of 512 nodes, with each node equipped with a 32-core CPU, 128GB of memory, and 4 DCU accelerator cards with 16GB of VRAM, totaling 2048 accelerator cards, each of which is equivalent to approximately 0.2 times the computational power of an NVIDIA A100.
- * **Parallel Configuration:** The Megatron-LM framework employed 3D parallelism techniques, including pipeline parallelism, model parallelism, and data parallelism, to maximize GPU performance and minimize communication overhead. Since each node had 4 accelerator cards, we set the model parallel size to 4 to achieve optimal parallel efficiency. Additionally, in cases where the mini-batch size was 1, we set the pipeline-parallel size to 16 to fully utilize the VRAM resources.
- * **Data Preprocessing:** We performed tokenization on all training text data, and the tokenized results of each document were concatenated using the `<eos>` marker. Subsequently, we cropped the concatenated tokens into fixed lengths of 2048, resulting in 30 million training samples, corresponding to 7324 training steps.
- * **Checkpoints:** Before formally starting the training process, we analyzed the node failure patterns through preliminary experiments and saved checkpoints at intervals of 100 steps.
- * **Hyperparameter Selection:** We conducted extensive experiments for hyperparameter selection in Further pre-train. Regarding learning rate scheduling, initial experiments showed that directly adopting a maximum learning rate of $1e-4$ from the Galactica-30B model led to a rapid increase in loss after a certain number of steps, resulting in training failure. Hence, we observed the gradient norm curve during the training warm-up phase and selected the learning rate corresponding to the minimum gradient norm, which was $1e-5$, as the actual maximum learning rate for training, which remained constant throughout the entire training process. For the training warm-up, we employed a linear training warm-up strategy and tested different training warm-up steps, and the optimum result was achieved with 100 training warm-up steps. Regarding other hyperparameters, we opted for the Adam optimizer with a β_1 of 0.9, a β_2 of 0.95, a weight decay rate of 0.1, and epsilon of $1e-8$. To balance effectiveness and efficiency, we set the global batch size to 4096 and utilized checkpoint activations to save VRAM. Additionally, we set the gradient clip threshold to 1.0 and a dropout rate of 0.1.

For a better understanding of our training, we list the hyperparameters of the model and the configured setting of the training in Appendix E.

Training curves We share the curves of the training loss and gradient normalization as Figure 6 and Figure 7. We observed that the training loss quickly dropped from about 1.60 to 1.40 during the first 300 steps and then smoothly decreased from 1.40 to 1.32 in the subsequent steps. Although the gradient normalization showed several spikes, sharply increasing from 0.1 to approximately 0.3~4.8, the model exhibited no signs of saturation after further pre-training on 60 billion tokens. This demonstrates the stability of the entire further pre-training process.

The bottleneck of the training

- The embedding layer is not treated as a stand-alone component in the training process. Instead, it is combined with the first transformer layer. As a result, the VRAM usage on certain cards is 60% higher than on others, leading to

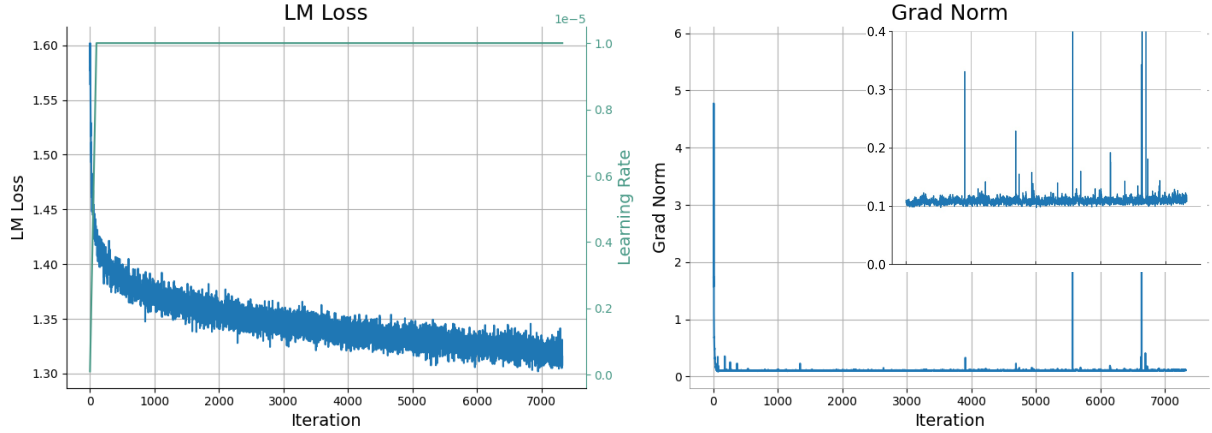


Figure 6: Training curve during the further pre-training.

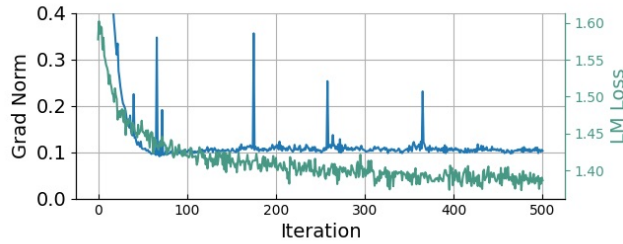


Figure 7: Training curve of the first 500 steps during the further pre-training.

decreased training efficiency. This is because a larger PP value is required to accommodate the entire model, which increases communication overhead.

- Due to some bugs in Megatron, the "continue pre-train" function cannot utilize distributed optimizers. This results in each DP group or model replica having a complete copy of the optimizer state, significantly increasing VRAM usage.

5.2 Supervised Fine-Tuning

After pre-training, LLMs can be supervised fine-tuning (SFT) on a smaller, more targeted dataset under human supervision. This process adapts the model to specific tasks or improves its performance in certain areas.

We employed SFT to enhance the geoscientific reasoning performance of our large-scale models on specific geoscientific tasks. This process is essential to effectively transfer advanced language capabilities to geoscientific-specific tasks and preserve the model’s generalization acquired during pre-training.

We utilized two major frameworks, Huggingface and DeepSpeed, during this stage to facilitate our training work. This aimed to accomplish instruction fine-tuning and model prediction tasks. In the training process, the Hygon DCU cluster remained our primary resource. Compared to the pre-training stage, SFT truncation only took advantage of 128 nodes and their accompanying 512 DCUs. We continued to employ the learning rate schedule used during pre-training, where the maximum learning rate was set to $1e - 5$, combined with a linear warm-up consisting of 100 warm-up steps. For the optimizer, we still selected the Adam optimizer, with β_1 and β_2 set to 0.9 and 0.999, respectively. Additionally, a weight decay of 0.05 and ϵ value of $1e-8$ was chosen to better adapt to the required fine-tuning tasks.

Considering the enormous scale of the model, we utilized the DeepSpeed ZeRO3 technique for memory optimization, along with the gradient checkpoint method, to further reduce memory pressure. The maximum input sequence length was limited to 512 in this process to avoid unnecessary computational overhead. However, due to the limitations of DeepSpeed, the global batch size had to be no smaller than the number of accelerator cards. Therefore, we opted

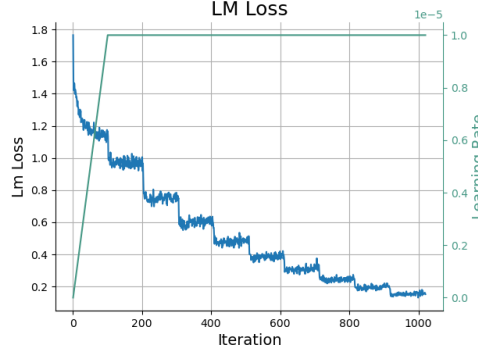


Figure 8: Training curve during the SFT on dataset Alpaca.

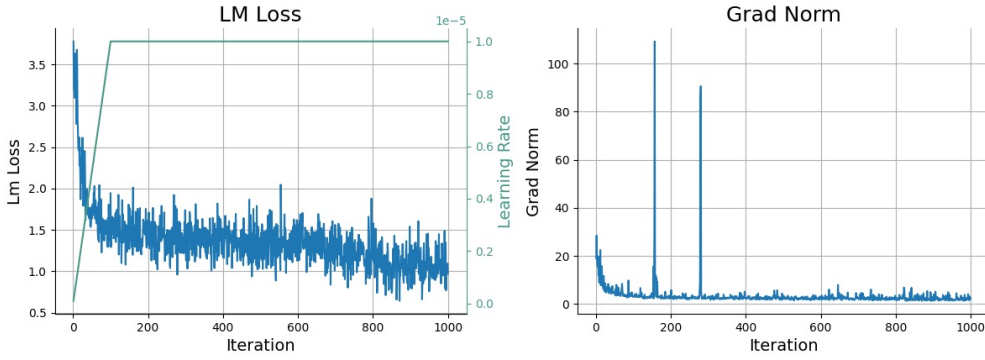


Figure 9: Training curve during the SFT on Geosignal.

for a larger global batch size of 512. Regarding the settings of other parameters, we followed the default values of the Huggingface trainer framework. For the subsequent training, we used the Alpaca dataset and conducted training for three epochs, which only took about one day to obtain the final SFT model. This training process, supported by Megatron-LM, supported our research work.

Following the recipe proposed by K2 [15] and a similar experience in RadiologyGPT [55], we did the SFT in two stages. For the first stage, we aligned the model with humans via Alpaca instruction tuning data, while using the GeoSignal v2 in the second stage. The training curve of SFT on Alpaca is Figure 8 while the SFT on GeoSignal is Figure 9

Moreover, we compare the variety of the instruction tuning data of Dolly and GeoSignal in Figure 11, showing that the general domain instructions dataset has less variety than the knowledge-intensive instructions dataset.

5.3 Tool Learning

In addition, LLMs can be designed to interact with and learn from various tools, such as browsers, databases, or other software interfaces. This allows the model to perform more complex tasks that require external information or specific functionalities.

We leveraged the ToolBench dataset [56], an open-source resource, to enable geoscientific large-scale models to leverage tool API capabilities. We sampled five types of tool QA data from ToolBench, namely *arxiv*, *bing_search*, *database*, *weather*, and *wolframalpha*, and supplemented it with our collected **geo_search** data, resulting in approximately 10k training samples. We open-source this dataset on <https://github.com/zthang/geotools>.

During the SFT stage, we trained the models together with training data such as alpaca. As the training samples from the tool data tend to be longer, we set the max_length to be 2048. During training, we only calculate the loss and backpropagate the gradients for the part of the API call, specifically the thought, action, action input, and the corresponding tokens for the final answer.

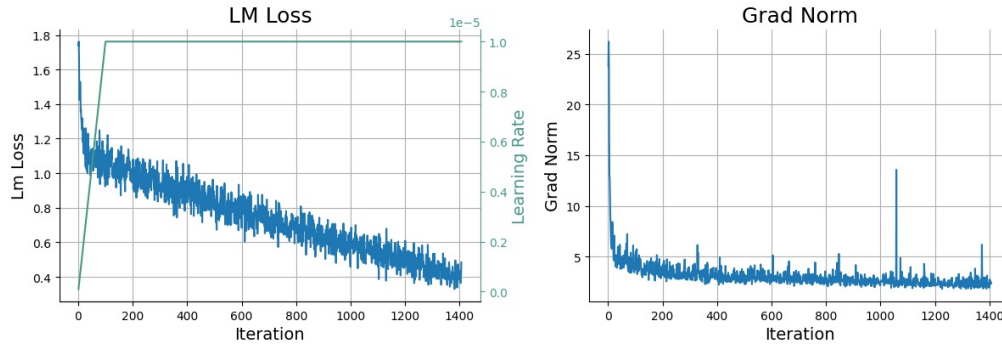


Figure 10: Training curve during the tools SFT.

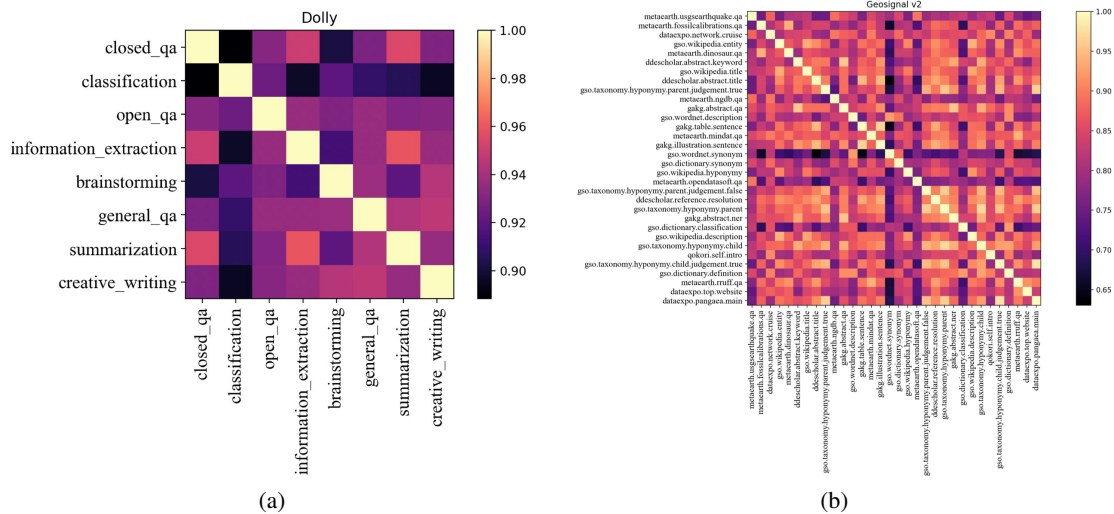


Figure 11: Variety of the instruction tuning data in Dolly and GeoSignal.

Once the model is trained, we specify the tool’s description and corresponding API parameter instructions in the prompt. For a given question, we first let the model output the related API call (thought, action, action input) to obtain the results returned by the external tool. These results are then used as observations and fed back into the model, generating a new set of thought, action, and action input for the next iteration (if further tool calls are required). This process continues until the model gathers enough information and outputs the final answer.

Here are two naive examples of how the Galactica-30b models use the tool. The detailed examples are shown in the Appendix:

- **Example 1:**

Question: "What is the weather in New York 3M years ago?"

Thought: "weather"

Action: "geo_search"

Action Input: "New York, Weather, 3M years"

- **Example 2**

Question: "What is the definition of plate tectonics?"

Thought: "arxiv"

Action: "search"

Action Input: "query: plate tectonics"

6 Evaluation

Once we have completed the model’s training, we proceed to examine its grasp of scientific and geoscientific knowledge. We have divided the evaluation into two parts.

- The first part involves automated evaluation using the GeoBench provided by K2. This enables us to assess the model’s performance in handling geoscientific tasks. Additionally, to examine if the newly learned knowledge has affected the pre-existing ability, we conducted MMLU (Minimal Meaningful Learning Units) tests. These tests are compared against the original Galactica model.
- The second part encompasses manual evaluation, where we carefully selected several subtasks from geoscience. For this evaluation, we invited 10 researchers specializing in geoscience to participate in voting and scoring. Ultimately, we compare the model’s performance with five other large-scale platforms in open testing.

By conducting these evaluations, we aim to comprehensively assess the model’s abilities and compare its performance against automated benchmarks and human assessments, ensuring its competence in scientific and geoscientific domains.

6.1 Automatic Evaluation

6.1.1 GeoBench

GeoBench, proposed by [15] is a benchmarking tool specifically designed to evaluate and test the geoscientific understanding and capabilities of LLMs. It focuses on assessing how well LLMs can process and generate responses involving geographic and geological information.

Baselines	NPEE	APTest
Random	27.1	20.0
Gal-6.7B	25.7	29.9
LLaMA-7B	21.6	27.6
K2-7B	39.9	29.3
ChatGPT	48.8	20.0
Gal-30B	41.2	<u>38.5</u>
GalAlp-30B	42.6	44.1
GEOGALACTICA-30B	<u>46.6</u>	36.9

Table 5: comparison among baselines on Objective tasks in GeoBench.

Through testing our models on GeoBench, we have observed that larger and more academic models outperform benchmarks like NPEE, which are inclined toward academic research. However, they do not perform well in benchmarks like AP Study, which lean more towards foundational education. This difference may be caused by training materials that guide the model to contemplate more advanced knowledge. The training data consists of academic research achievements, namely papers, which may result in a deviation from and lack of basic knowledge. This is an area we intend to focus on for improvement in the future.

It is worth noting that Galactica, with 30 billion parameters, often fails to outperform Llama, with 7 billion parameters, in general benchmark tasks. However, in our GeoBench, we have successfully developed GEOGALACTICA, which builds upon Galactica, surpassing K2, built upon Llama.

6.1.2 MMLU

The MMLU has been divided into math and non-math sections by Galactica, and we have been following their reports closely. From the results (Shown in Table 6), it is evident that after processing 6 million geoscience-related literature documents, specific skills of the model, such as algebra, biology, chemistry, and mathematics, have shown improvement. This phenomenon appears to be linked to papers focusing on mathematical geology, biological geoscience, and chemical geology, highlighting the interdisciplinary nature of geoscience. Surprisingly, machine learning has experienced significant enhancement, likely due to the inclusion of GitHub code in our corpus. In summary, subjects closely related to geoscience, including those logically connected to geology and its subfields, have shown notable progress. However, disciplines like physics indicate that the original Galactica outperforms our GEOGALACTICA and subjects unrelated to geosciences, such as medical genetics, medicine, and electrical engineering, have shown a decline in performance. It is noteworthy that GEOGALACTICA and the original Galactica are generally at a similar stage regarding average performance in math-related subjects within the MMLU.

Subject	GEOGALACTICA30B	GAL 30B	GalAlp 30B
Abstract Algebra	0.300	0.250	0.320
Astronomy	0.461	0.500	0.474
College Biology	0.576	0.576	0.514
College Chemistry	0.370	0.320	0.350
College Computer Science	0.400	0.410	0.370
College Mathematics	0.320	0.350	0.350
College Medicine	0.480	0.520	0.445
College Physics	0.284	0.333	0.294
Econometrics	0.377	0.368	0.368
Electrical Engineering	0.538	0.579	0.503
Elementary Mathematics	0.328	0.310	0.288
Formal Logic	0.302	0.270	0.278
High School Biology	0.565	0.561	0.535
High School Chemistry	0.360	0.399	0.355
High School Computer Science	0.500	0.480	0.510
High School Mathematics	0.311	0.256	0.304
High School Physics	0.298	0.364	0.325
High School Statistics	0.333	0.352	0.319
Machine Learning	0.411	0.339	0.366
Medical Genetics	0.550	0.580	0.520
Average	0.4032	0.40585	0.3894

Table 6: We report the results of the three models in math.

After assessing the mathematical subject, we examined the results of the subjects that were excluded. Overall, GEOGALACTICA performs slightly better than the original Galactica in the average of non-math-related subjects in MMLU. Interestingly, subjects like global facts, US History, and World History have significantly improved compared to the original Galactica. This phenomenon can be attributed to the fact that many aspects of history, such as significant discoveries and political knowledge, are closely intertwined with geoscience. This underscores the significance of geoscience, which can profoundly influence global progress.

Furthermore, in conceptual physics, learning from geoscience papers has led to a better understanding of the model. This suggests that several concepts in geoscience do not align with the knowledge taught in colleges and high schools. Consequently, models struggle to apply this related knowledge when solving problems at the college and high school levels.

Observation on ablation Fortunately, we came across Galpaca-30B on Hugging Face ²⁷, which significantly reduced the carbon emissions from our finetuning experiments. This model utilized Alpaca’s instructions to learn from the dataset and was applied to SFT on Galactica-30B. Upon horizontal comparison, Galpaca-30B performed notably worse than the original Galactica and GEOGALACTICA in the majority of disciplines. This indicates that instruction learning in the general domain can significantly impact the performance of specialized domain models during practical evaluations.

6.2 Human Evaluation

In this part, we have selected five open models to evaluate together with our GEOGALACTICA model. These models include:

1. MOSS, an open-source tool-augmented conversational language model, was released by Qiu Xipeng’s team from the School of Computer Science at Fudan University as a ChatGPT-like model.
2. Qwen is a chatbot developed by Alibaba Cloud, a technology company under the Alibaba Group. Alibaba announced its intention to open Tongyi Qianwen to the public, indicating its readiness for the market and reflecting China’s growing focus on AI technology.
3. ChatGPT is an AI language model developed by OpenAI, known for its ability to generate human-like text based on prompts, facilitate engaging conversations, answer questions, and perform a wide range of language-related tasks. ²⁸

²⁷<https://huggingface.co/GeorgiaTechResearchInstitute/galpaca-30b>

²⁸We use the 2023 March version of ChatGPT.

Subject	GeoGal 30B	Gal 30B	GalAlp 30B
Anatomy	0.496	0.541	0.533
Business Ethics	0.430	0.420	0.470
Clinical Knowledge	0.532	0.555	0.491
Computer Security	0.600	0.650	0.620
Conceptual Physics	0.481	0.434	0.417
Global Facts	0.390	0.300	0.340
High School European History	0.533	0.606	0.491
High School Geography	0.581	0.540	0.515
High: School Gov & Politis	0.534	0.565	0.461
High School Macroeconomics	0.408	0.405	0.367
High School Microeconomics	0.424	0.458	0.424
High School Psychology	0.613	0.628	0.556
High School US History	0.436	0.352	0.319
High School World History	0.620	0.456	0.446
Human Aging	0.552	0.552	0.511
Human Sexuality	0.511	0.565	0.481
International Law	0.612	0.644	0.554
Jurisprudence	0.491	0.472	0.444
Logical Fallacies	0.423	0.472	0.442
Management	0.573	0.602	0.515
Marketing	0.641	0.705	0.607
Miscellaneous	0.522	0.501	0.470
Moral Disputes	0.480	0.462	0.468
Moral Scenarios	0.238	0.244	0.245
Nutrition	0.536	0.520	0.448
Philosophy	0.444	0.492	0.431
Prehistory	0.503	0.522	0.435
Professional Accounting	0.344	0.312	0.319
Professional Iaw	0.326	0.326	0.327
Professional Medicine	0.438	0.449	0.379
Professional Psychology	0.472	0.505	0.449
Public Relations	0.473	0.445	0.455
Security Studies	0.424	0.408	0.322
Sociology	0.537	0.547	0.483
US Foreign Policy	0.550	0.510	0.540
Virology	0.434	0.422	0.410
World Religion	0.421	0.427	0.380
Average	0.487	0.486	0.448

Table 7: We report the results of the three models in social sciences.

4. Yiyan, also known as Ernie Bot, is an AI chatbot service product developed by Baidu. It has been under development since 2019 and is based on a large language model named "Ernie 4.0", which was announced on October 17, 2023
5. ChatGLM is an open bilingual language model developed by Tsinghua University. It is optimized for Chinese conversation and is based on the General Language Model architecture.

For the selected projects, our evaluation is designed as follows:

We refer to K2’s Human Evaluation and define the evaluation metrics for open-ended questions: scientificity, correctness, and coherence (score range is [1, 2, 3]). The specific explanations are as follows:

- **Scientificity:** It represents whether the generated content appears as something that a geoscience professional would say. A score of 1 indicates not good, 2 indicates acceptable, and 3 indicates very good.
- **Correctness:** From the perspective of a geoscience expert, whether the model convinces you and if the information obtained is correct. A score of 1 indicates incorrect, 2 indicates possibly right, and 3 shows correct.
- **Coherence:** This metric is used to evaluate the consistency and coherence of the model, i.e., whether the text consistently discusses a specific topic and reads smoothly. A score of 1 indicates not good, 2 indicates acceptable, and 3 indicates very good.

Category	Problem	Prompt	Skills
Open-ended	Noun Definition	What is carbonate rock?	Knowledge
Open-ended	Beginner Level Q&A	How many continents are there in the world?	Knowledge
Open-ended	Intermediate Level Q&A	Would the Ohio train derailment leading to vinyl chloride leakage affect the ecological environment around the Great Lakes based on ocean currents or air dispersion?	Analysis
Open-ended	Advanced Level Q&A	How did dinosaurs become extinct?	Discovery
Functional	Confirmation of Geoscience Knowledge System	Is carbonate rock a type of limestone?	Judgment
Functional	Geoscience Paper Titling	This is the abstract of my paper. Can you help me come up with a title?	Summarization
Functional	Paper Summary	This is my passage. Can you help me summarize the passage?	Summarization
Functional	Speech Writing	Please help me write a speech based on my topic.	Writing
Functional	Pre-requisite Knowledge Recommendations	This is my article. Can you recommend some prerequisite knowledge points?	Information extraction

Table 8: Tasks we designed in human evaluation parts.

Based on this, the cumulative score can be calculated. Additionally, for the functional questions of the large model, the evaluation metric is relative ranking. Participants in the evaluation will receive replies from all six models on the same input, and our expert judges will rate these models in the order of 1, 2, 3, 4, 5, and 6. Finally, the total ranking of each model will be calculated. In this part, we invite **10** geoscience practical people, including 6 students and 4 teachers. (The contribution of the human evaluation is shown in Appendix L).

Open-ended Tasks For open-ended questions, the general large-scale model uses the interface output provided by ChatALL²⁹ for consistency. Our large model interacts through our UI interface, where higher metric scores are preferred.

6.2.1 Noun Definition

In our geoscience entrance exam question set, we randomly selected 20 geoscience vocabulary terms to evaluate the model’s understanding of domain-specific terminology, the whole terms are in subsubsection G.1.1.

	Scientificity	Correctness	Coherence
MOSS	291	302	351
Qianwen	419	435	435
ChatGPT	337	351	357
Yiyan	236	276	305
ChatGLM	278	291	347
GEOGALACTICA	339	361	393

Table 9: Comparison on Noun Definition tasks.

In this task, our model has demonstrated remarkable vocabulary proficiency, but what indeed astonishes us is its exceptional ability to handle scientific questions and professional respond beyond what other models can achieve. One example is shown in subsection H.1.

6.2.2 Beginner Level Q&A

We selected 10 easy questions from the high school geoscience Olympiad in China and had them translated into English by professional translators. The questions are shown in subsubsection G.1.2. One example is shown in subsection H.2.

Overall, our model ranks third when considering all aspects, but ChatGPT outperforms other models significantly in this category of questions.

²⁹<https://github.com/sunner/ChatALL>

	Scientificity	Correctness	Coherence
MOSS	116	120	147
Qianwen	191	177	207
ChatGPT	219	214	225
Yiyan	176	174	187
ChatGLM	160	156	184
GEOGALACTICA	176	173	202

Table 10: Comparison on Beginner Level Q&A.

6.2.3 Intermediate Level Q&A

We have selected 10 moderately difficult questions from Chegg and SaveMyExam, which require a certain level of geo-science knowledge training. The questions are shown in subsection G.1.3. One example is shown in subsection H.3.

	Scientificity	Correctness	Coherence
MOSS	143	154	178
Qianwen	178	180	193
ChatGPT	210	206	207
Yiyan	180	186	189
ChatGLM	161	163	179
GEOGALACTICA	162	169	171

Table 11: Comparison on Intermediate Level Q&A.

Overall, our model ranks tied for fourth place, but ChatGPT outperforms other models significantly in this category of questions. The results are shown in Table 11.

6.2.4 Advanced Level Q&A

We have selected 9 highly difficult questions from the urgent geoscience problems proposed by the Institute of Geography, Chinese Academy of Sciences. These questions require extensive training in geoscience knowledge as well as the ability to reason through scientific research. The questions are shown in subsection G.1.4. One example is shown in subsection H.4.

	Scientificity	Correctness	Coherence
MOSS	166	173	194
Qianwen	202	199	209
ChatGPT	137	133	181
Yiyan	190	192	200
ChatGLM	172	171	194
GEOGALACTICA	185	187	206

Table 12: Comparison on Advanced Level Q&A.

Overall, according to Table 12 our model ranks third, but ChatGPT seems to lack sufficient capability to handle these types of questions.

Functional Tasks When it comes to the evaluation of functional questions, we have chosen to apply GEOGALACTICA to scientific research literature. GEOGALACTICA is dedicated to facilitating the comprehension and interpretation of scientific research literature. When external information input is not required, we utilize the consistent output provided by the interface of ChatALL. In terms of overall evaluation, since it involves ranking, lower scores are preferred.

6.2.5 Knowledge-based associative judgment question

To determine the presence or absence of knowledge system relationships, the questions are derived from the Knowledge trees in GSO. The questions are shown in subsection G.2.1. One example is shown in subsection H.5.

Models	Sum of Rank
MOSS	579
Qianwen	557
ChatGPT	600
Yiyan	570
ChatGLM	752
GEOGALACTICA	725

Table 13: Comparison on knowledge-based associative judgment question.

Overall, our model ranks fifth, indicating that there is still significant room for improvement in handling these logical questions. Further advancements can be achieved by constructing CoT-type data and injecting more expertise into the model.

6.2.6 Research Paper Titling Task

In this phase, we randomly selected abstracts from 20 geoscience research papers and inputted them into the model, asking it to generate a title. This task showcases the model’s understanding of knowledge points and familiarity with the field. The questions are shown in subsubsection G.2.2. One example is shown in subsection H.6.

Models	Sum of Rank
MOSS	805
Qianwen	426
ChatGPT	326
Yiyan	561
ChatGLM	440
GEOGALACTICA	451

Table 14: Comparison on research paper titling task.

Overall, our model ranks fourth, with no clear distinction between the ChatGLM, Qianwen, and our model in terms of performance on this task.

6.2.7 Geoscience Research Functionality

To ensure fairness when incorporating external research papers for evaluation, we employ our own PDF parsing solution to interpret the papers. We then use the consistent output provided by the ChatALL interface. As for our GEOGALACTICA, we utilize our UI interface for interactions and obtain outputs accordingly.

More specifically, when it comes to interpreting scientific literature, we often inquire about the following aspects:

1. Can you help me write a speech based on the content of the article?
2. Can you help me summarize the article?
3. Could you please recommend some prerequisite knowledge points?

These three scenarios are all closely related to us as researchers. We have assessed five papers, which cover various domains of Earth sciences and are written in different styles. The papers are listed in subsubsection G.2.3. One example is shown in subsection H.7. After the evaluation, our ranking is as follows:

	Writing	Summary	Extraction
MOSS	114	164	115
Qianwen	135	185	232
ChatGPT	62	86	51
Yiyan	178	139	160
ChatGLM	106	168	169
GEOGALACTICA	135	100	212

Table 15: Comparison on geoscience research functionality.

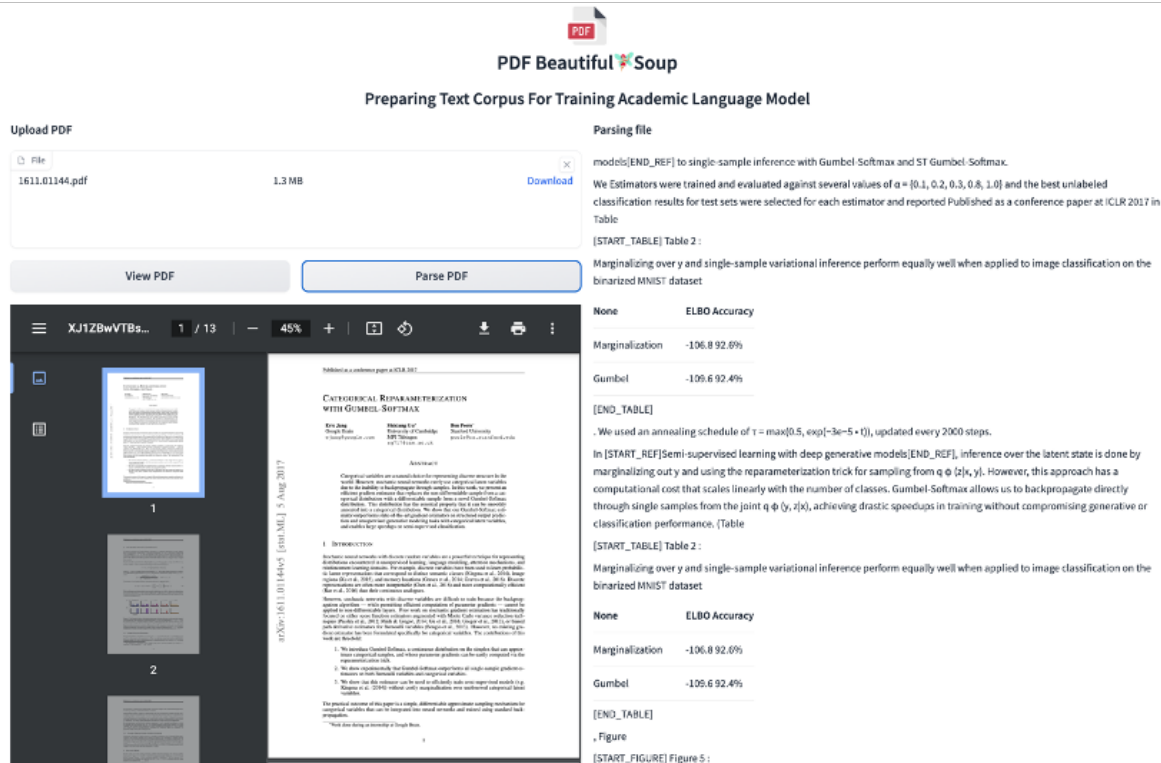


Figure 12: The user interface for evaluating the GEOGALACTICA.

Overall, considering the summation of the three scores, ChatGPT and MOSS are better than GEOGALACTICA, and GEOGALACTICA is tied third with ChatGLM. Our model demonstrates excellent summarization skills, thanks to its comprehensive incorporation of knowledge in the field of Earth sciences. The same principle applies to the task of extracting key information, for which we still need to gather a certain amount of expert thinking data.

7 discussion

7.1 The Necessity of Pre-training

Initiating training for a domain-specific language model from scratch in the field of geoscience is a complex decision that requires careful consideration of multiple factors. Here, we discuss some thoughts on why we did not consider training from scratch:

1. Geoscience data is relatively limited, and training a high-quality model from scratch requires sufficient data support. The availability of geoscience data is constrained, potentially leading to the issue of data scarcity when training from scratch.
2. Training a large-scale language model from scratch demands substantial computational resources and time. As we lack sufficient resources and time, employing pre-trained models and conducting transfer learning may yield more cost-effective results.
3. Using pre-trained models and conducting transfer learning has shown promising results within a relatively short timeframe. Thus, further training can be practical. Refer to K2. Therefore, while training from scratch would enable models to comprehend and capture domain-specific terminology, concepts, and relationships in geoscience, as well as train from the outset on geoscience data to better adapt to domain-specific language and knowledge, we opted for a strategy that involves further pre-training due to cost, time, and data considerations.

7.2 The Necessity of Further Pre-training

Thus, in our perspective, employing more general-purpose models for transfer learning and further pre-training in the field of geoscience can be meaningful because:

1. Geoscience encompasses various specialized domains, such as geology, meteorology, and environmental science. Further pre-training can enhance the model’s understanding and capture of these domain-specific concepts, terms, and relationships through training on geoscience-related textual data. Additionally, geoscience often involves many domain-specific terms and contextual information that may not be commonly found in everyday language. The model can better comprehend and contextualize these terms through further pre-training, thereby improving performance in geoscience texts.
2. Geoscience texts may rely on specific geographical backgrounds, spatial and temporal relationships, and regional information. Further pre-training can assist the model in better understanding these contextual dependencies, leading to more accurate information processing and generation. Further pre-training can enhance model performance for tasks such as text classification, information extraction, and generating geological reports. The model can learn more task-specific feature representations by training on relevant domain-specific data.
3. During this process, we observed the alleviation of data scarcity in geoscience. In certain domains within geoscience, scarce data may limit training samples. However, further pre-training enables the model to learn general language abilities from a larger-scale dataset and subsequently fine-tune on a smaller amount of domain-specific data, mitigating the impact of data scarcity.

In conclusion, further pre-training can enable large language models to better adapt to the characteristics and requirements of the geoscience field, leading to enhanced performance in geoscience text processing and task execution.

7.3 Carbon Emissions

During our cumulative training, **1,488,137.26 DCU hours** were consumed, resulting in cumulative carbon emissions of $212\ tCO_2eq$ calculated by Equation 1. Our work provides a foundational model for subsequent geoscience researchers to fine-tune their smaller models, potentially reducing carbon emissions in their future work.

$$C_{Emission}(kg) = DCU\ hours * TDP\ (kW) * C_{Intensity}\ (kg/kWh) \quad (1)$$

7.4 Towards Unified Foundation Model in Geoscience

The application of artificial intelligence in geosciences demonstrates vast prospects. In terms of geoscientific literature analysis, AGI systems especially the unified foundation model can assist researchers in identifying the frequency of specific vocabulary while addressing any ambiguities, thereby enhancing the accuracy of literature comprehension. Furthermore, AGI can integrate dispersed geoscientific knowledge by analyzing extensive literature uncovering novel correlations and trends, thus providing new perspectives and directions for geoscience research. Additionally, AGI systems can aid in geoscience education, offering personalized content and teaching methods to facilitate students’ ease of learning and understanding of geoscientific knowledge.

On the other hand, the application of a unified foundation model in the geoscience domain extends beyond academic research to practical uses such as geological hazard warnings, resource exploration, and environmental protection. Regarding geological hazard warnings, AGI can utilize big data and models to provide accurate predictions and assessments, helping to mitigate the damages caused by natural disasters. Concurrently, a unified foundation model plays a crucial role in underground resource exploration, improving efficiency and accuracy in exploration endeavors. In the realm of environmental protection, a unified foundation model aids in the real-time monitoring of environmental conditions through the analysis of remote sensing data, supporting decision-making processes for environmental conservation.

In the future, with sufficient abundant data and computing powers, and other feasibility of achieving unified foundation models in the field of geoscience, the future of AGI in Geoscience can be expected. In the future, the unified foundation model will continue to play a role in advancing frontiers in geoscience research. It can assist scientists in conducting large-scale data analysis, unraveling complex phenomena such as internal Earth structures, plate tectonics, and crustal evolution, thereby providing deeper scientific comprehension. AGI also contributes to environmental monitoring and protection, aeromagnetic data interpretation, water resource management, carbon capture, and other domains. By analyzing hydrological data and geological information, a unified foundation model can predict groundwater resources’ distribution and sustainable utilization, providing scientific foundations for water resource management. In carbon capture, a unified foundation model can assist researchers in selecting suitable geological storage layers and sealing

rocks, thereby driving the development of carbon reduction technologies. Overall, AGI accelerates the accumulation of scientific knowledge in geosciences and offers unprecedented support in addressing global challenges, providing robust intelligent assistance for humanity’s future sustainable development.

8 Conclusion

In conclusion, the utility of NLP in geoscience research and practice is vast, and large language models (LLMs) have shown great success in various NLP domains. However, specialized LLMs in geoscience are scarce. We introduce GEOGALACTICA, a 30B parameters language model designed explicitly for geoscience applications. Through training on a comprehensive geoscience academic dataset and fine-tuning with geoscience-knowledge intensive instruction pairs, GEOGALACTICA outperforms existing models in geoscience NLP tasks. Our validation with senior geoscientists confirms its effectiveness. The release of GEOGALACTICA and our training experience aims to contribute to the advancement of unified foundation models in geoscience.

Acknowledgement

The computation resource was supported by the Advanced Computing East China Sub-center. This work is supported by NSF China (No.62020106005, 61960206002, 42050105, 62061146002, 62106143), National Key Technologies R&D Program (No. 2022YFB3904201), Shanghai Pilot Program for Basic Research - Shanghai Jiao Tong University. The second author would like to thank Wu Wen Jun Honorary Doctoral Scholarship, AI Institute, Shanghai Jiao Tong University.

References

- [1] Shu Wang, Yunqiang Zhu, Yanmin Qi, Zhiwei Hou, Kai Sun, Weirong Li, Lei Hu, Jie Yang, and Hairong Lv. A unified framework of temporal information expression in geosciences knowledge system. *Geoscience Frontiers*, 2022.
- [2] Cheng Deng, Yuting Jia, Hui Xu, Chong Zhang, Jingyao Tang, Luoyi Fu, Weinan Zhang, Haisong Zhang, Xinbing Wang, and Cheng Zhou. Gakg: A multimodal geoscience academic knowledge graph. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021.
- [3] Rahul Ramachandran, Muthukumaran Ramasubramanian, Pravesh Koirala, Iksha Gurung, and Manil Maskey. Language model for earth science: Exploring potential downstream applications as well as current challenges. *IGARSS 2022 - 2022 IEEE International Geoscience and Remote Sensing Symposium*, pages 4015–4018, 2022.
- [4] Hao Zhang and Jin-Jian Xu. When geoscience meets foundation models: Towards general geoscience artificial intelligence system. *arXiv preprint arXiv:2309.06799*, 2023.
- [5] Qinjun Qiu, Zhong Xie, Liang Wu, and Wenjia Li. Geoscience keyphrase extraction algorithm using enhanced word embedding. *Expert Systems with Applications*, 125:157–169, 2019.
- [6] Christopher JM Lawley, Michael G Gadd, Mohammad Parsa, Graham W Lederer, Garth E Graham, and Arianne Ford. Applications of natural language processing to geoscience text data and prospectivity modeling. *Natural Resources Research*, pages 1–25, 2023.
- [7] Qinjun Qiu, Zhong Xie, Liang Wu, and Liufeng Tao. Automatic spatiotemporal and semantic information extraction from unstructured geoscience reports using text mining techniques. *Earth Science Informatics*, 13:1393–1410, 2020.
- [8] Qinjun Qiu, Zhong Xie, Liang Wu, and Wenjia Li. Dgeosegmenter: A dictionary-based chinese word segmenter for the geoscience domain. *Computers & geosciences*, 121:1–11, 2018.
- [9] Chengbin Wang, Xiaogang Ma, Jianguo Chen, and Jingwen Chen. Information extraction and knowledge graph construction from geoscience literature. *Computers & geosciences*, 112:112–120, 2018.
- [10] Cheng Deng, Bo Tong, Luoyi Fu, Jiabin Ding, Dexing Cao, Xinbing Wang, and Chenghu Zhou. Pk-chat: Pointer network guided knowledge driven generative dialogue model. *arXiv preprint arXiv:2304.00592*, 2023.
- [11] Kai Ma, Miao Tian, Yongjian Tan, Xuejing Xie, and Qinjun Qiu. What is this article about? generative summarization with the bert model in the geosciences domain. *Earth Science Informatics*, pages 1–16, 2022.
- [12] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony S. Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *ArXiv*, abs/2211.09085, 2022.
- [13] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023.
- [14] Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, P. Zhang, Yuxiao Dong, and Jie Tang. Glm-130b: An open bilingual pre-trained model. *ArXiv*, abs/2210.02414, 2022.
- [15] Cheng Deng, Tianhang Zhang, Zhongmou He, Qiyuan Chen, Yuanyuan Shi, Le Zhou, Luoyi Fu, Weinan Zhang, Xinbing Wang, Cheng Zhou, Zhouhan Lin, and Junxian He. Learning a foundation language model for geoscience knowledge understanding and utilization. *ArXiv*, abs/2306.05064, 2023.
- [16] Fanchun Meng, Tao Ren, Zhenxian Liu, and Zhida Zhong. Toward earthquake early warning: A convolutional neural network for repaid earthquake magnitude estimation. *Artificial Intelligence in Geosciences*, 2023.
- [17] Yifeng Fei, Hanpeng Cai, Junhui Yang, Jiandong Liang, and Guang Hu. Unsupervised pre-stack seismic facies analysis constrained by spatial continuity. *Artificial Intelligence in Geosciences*, 2023.
- [18] Qingkai Kong, Andrea Chiang, Ana C Aguiar, M Giselle Fernández-Godino, Stephen C Myers, and Donald D Lucas. Deep convolutional autoencoders as generic feature extractors in seismological applications. *Artificial intelligence in geosciences*, 2:96–106, 2021.
- [19] Priyadarshi Chinmoy Kumar and Kalachand Sain. Machine learning elucidates the anatomy of buried carbonate reef from seismic reflection data. *Artificial Intelligence in Geosciences*, 4:59–67, 2023.
- [20] Mazahir Hussain, Shuang Liu, Umar Ashraf, Muhammad Ali, Wakeel Hussain, Nafees Ali, and Aqsa Anees. Application of machine learning for lithofacies prediction and cluster analysis approach to identify rock type. *Energies*, 15(12):4501, 2022.

- [21] Reda Abdel Azim. A new correlation for calculating wellhead oil flow rate using artificial neural network. *Artificial Intelligence in Geosciences*, 3:1–7, 2022.
- [22] Peiyi Yao, Ziwang Yu, Yanjun Zhang, and Tianfu Xu. Application of machine learning in carbon capture and storage: An in-depth insight from the perspective of geoscience. *Fuel*, 333:126296, 2023.
- [23] Deming Xu, Yusheng Wang, Jingqi Huang, Sijin Liu, Shujun Xu, and Kun Zhou. Prediction of geology condition for slurry pressure balanced shield tunnel with super-large diameter by machine learning algorithms. *Tunnelling and Underground Space Technology*, 131:104852, 2023.
- [24] Ting Chen, Yaojun Wang, Hanpeng Cai, Gang Yu, and Guangmin Hu. High resolution pre-stack seismic inversion using few-shot learning. *Artificial Intelligence in Geosciences*, 3:203–208, 2022.
- [25] Anik Saha and Sunil Saha. Integrating the artificial intelligence and hybrid machine learning algorithms for improving the accuracy of spatial prediction of landslide hazards in kurseong himalayan region. *Artificial Intelligence in Geosciences*, 3:14–27, 2022.
- [26] Karianne J. Bergen, Paul A. Johnson, Maarten V. de Hoop, and Gregory C. Beroza. Machine learning for data-driven discovery in solid earth geoscience. *Science*, 363(6433), mar 2019.
- [27] Claudia Hulbert, Bertrand Rouet-Leduc, Paul A. Johnson, Christopher X. Ren, Jacques Rivière, David C. Bolton, and Chris Marone. Similarity of fast and slow earthquakes illuminated by machine learning. *Nature Geoscience*, 12(1):69–74, dec 2018.
- [28] Ying Chen, Jim Haywood, Yu Wang, Florent Malavelle, George Jordan, Daniel Partridge, Jonathan Fieldsend, Johannes De Leeuw, Anja Schmidt, Nayeong Cho, Lazaros Oreopoulos, Steven Platnick, Daniel Grosvenor, Paul Field, and Ulrike Lohmann. Machine learning reveals climate forcing from aerosols is dominated by increased cloud cover. *Nature Geoscience*, 15(8):609–614, aug 2022.
- [29] Xiao Liu, Juan Hu, Qi Shen, and Huan Chen. Geo-bert pre-training model for query rewriting in poi search. In *Conference on Empirical Methods in Natural Language Processing*, 2021.
- [30] Qinjun Qiu, Kai Ma, Hairong Lv, Liufeng Tao, and Zhong Xie. Construction and application of a knowledge graph for iron deposits using text mining analytics and a deep learning algorithm. *Mathematical Geosciences*, 55(3):423–456, 2023.
- [31] Bin Wang, Liang Wu, Zhong Xie, Qinjun Qiu, Yuan Zhou, Kai Ma, and Liufeng Tao. Understanding geological reports based on knowledge graphs using a deep learning approach. *Computers & Geosciences*, 168:105229, 2022.
- [32] Huseyin Denli, HassanJaved Chughtai, Brian Hughes, Robert Gistri, and Peng Xu. Geoscience language processing for exploration. *Day 3 Wed, November 17, 2021*, 2021.
- [33] Qinjun Qiu, Zhong Xie, Kai Ma, Liufeng Tao, and Shiyu Zheng. Neurospe: A neuro-net spatial relation extractor for natural language text fusing gazetteers and pretrained models. *Transactions in GIS*.
- [34] Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. In *Conference on Empirical Methods in Natural Language Processing*, 2019.
- [35] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(1):1–23, 2021.
- [36] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *ArXiv*, abs/2303.17564, 2023.
- [37] Renqian Luo, Liai Sun, Yingce Xia, Tao Qin, Sheng Zhang, Hoifung Poon, and Tie-Yan Liu. Biogpt: Generative pre-trained transformer for biomedical text generation and mining. *Briefings in bioinformatics*, 2022.
- [38] Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, W. Zhang, Imran Razzak, and Bram Hoex. Darwin series: Domain specific large language models for natural science. *ArXiv*, abs/2308.13565, 2023.
- [39] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36:1234 – 1240, 2019.
- [40] Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *ArXiv*, abs/1904.05342, 2019.
- [41] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. Legalbert: The muppets straight out of law school. *ArXiv*, abs/2010.02559, 2020.

- [42] Yifan Zhang, Cheng Wei, Shangyou Wu, Zhengting He, and Wenhao Yu. Geogpt: Understanding and processing geospatial tasks through an autonomous gpt. *arXiv preprint arXiv:2307.07930*, 2023.
- [43] Chong Ma, Zihao Wu, Jiaqi Wang, Shaochen Xu, Yaonai Wei, Zhengliang Liu, Lei Guo, Xiaoyan Cai, Shu Zhang, Tuo Zhang, et al. Impressiongpt: an iterative optimizing framework for radiology report summarization with chatgpt. *arXiv preprint arXiv:2304.08448*, 2023.
- [44] Zhiyuan Peng, Xuyang Wu, and Yi Fang. Soft prompt tuning for augmenting dense retrieval with large language models. *arXiv preprint arXiv:2307.08303*, 2023.
- [45] wikipedia. History of artificial neural networks. 2023.
- [46] Shao Zhang, Yuting Jia, Hui Xu, Ying Wen, Dakuo Wang, and Xinbing Wang. Deepshovel: An online collaborative platform for data extraction in geoscience literature with ai assistance. *arXiv preprint arXiv:2202.10163*, 2022.
- [47] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- [48] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- [49] Grobid: A machine learning software for extracting information from scholarly documents. <https://github.com/kermitt2/grobid>, 2008–2023.
- [50] Christopher Clark and Santosh Kumar Divvala. Pdffigures 2.0: Mining figures from research papers (jcdl’16). 143–152. *Google Scholar Google Scholar Digital Library Digital Library*, 2016.
- [51] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [52] Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *ArXiv*, abs/2304.01196, 2023.
- [53] Weizhe Yuan and Pengfei Liu. restructured pre-training. *ArXiv*, abs/2206.11147, 2022.
- [54] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- [55] Zhengliang Liu, Aoxiao Zhong, Yiwei Li, Longtao Yang, Chao Ju, Zihao Wu, Chong Ma, Peng Shu, Cheng Chen, Sekeun Kim, et al. Radiology-gpt: A large language model for radiology. *arXiv preprint arXiv:2306.08666*, 2023.
- [56] Yujia Qin, Shi Liang, Yining Ye, Kunlun Zhu, Lan Yan, Ya-Ting Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Runchu Tian, Ruobing Xie, Jie Zhou, Marc H. Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis. *ArXiv*, abs/2307.16789, 2023.
- [57] Michael D. McCormack. Neural computing in geophysics. *Geophysics*, 10:11–15, 1991.
- [58] Dave Hale. Methods to compute fault images, extract fault surfaces, and estimate fault throws from 3d seismic images. *Geophysics*, 78, 2013.
- [59] Anders U. Waldebrand and Anne H. Schistad Solberg. Salt classification using deep learning. 2017.
- [60] Li-Xin Wang and Jerry M. Mendel. Adaptive minimum prediction-error deconvolution and source wavelet estimation using hopfield neural networks. *Geophysics*, 57:670–679, 1992.
- [61] Yusuf Nasir and Louis J. Durlofsky. Deep reinforcement learning for optimal well control in subsurface systems with uncertain geology. *J. Comput. Phys.*, 477:111945, 2022.
- [62] Pablo Guillén, Germán Larrazábal, Gladys Gonzalez, Dainis Boumber, and Ricardo Vilalta. Supervised learning to detect salt body. *Seg Technical Program Expanded Abstracts*, 2015.
- [63] Heidi Anderson Kuzma. A support vector machine for avo interpretation. *Seg Technical Program Expanded Abstracts*, pages 181–184, 2003.
- [64] Henri Blondelle, A. Juneja, J. Micaelli, and Philip Neri. Machine learning can extract the information needed for modelling and data analysing from unstructured documents. 2017.
- [65] L. Y. Zheng, A. Albayrak, W. L. Teng, M. G. Khayat, and L. Pham. Using Transformer Networks and Knowledge Graphs in Earth Science Literature to Synthesize Mass Information for Transdisciplinary Research. In *AGU Fall Meeting Abstracts*, volume 2020, pages IN030–04, December 2020.

- [66] Danfeng Hong, Lianru Gao, Jing Yao, Bing Zhang, Antonio J. Plaza, and Jocelyn Chanussot. Graph convolutional networks for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 59:5966–5978, 2020.
- [67] Prasanna Koirala, Muthukumaran Ramasubramanian, Iksha Gurung, Manil Maskey, and Rahul Ramachandran. BERT-E: An Earth Science Specific Language Model for Domain-Specific Downstream Tasks. In *AGU Fall Meeting Abstracts*, volume 2021, pages IN15B–06, December 2021.
- [68] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Anand Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei A. Zaharia. Efficient large-scale language model training on gpu clusters using megatron-lm. *SC21: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–14, 2021.

A Appendix: Progression of geoscience with AI

Here we show the progression of geoscience research with the use of cutting-edge AI techniques summarized by [45].

Methods	In CS	In Geoscience	Gap
NN	1951	1991 Neural computing in geophysics [57]	40
Perceptron	1958	-	
KNN	1967	2013 Methods to compute fault images, extract fault surfaces, and estimate fault throws from 3D seismic images [58]	46
CNN/BP	1980/1989	2017 Salt classification using deep learning [59]	37
RNN	1982	1992 Adaptive minimum prediction-error deconvolution and source wavelet estimation using Hopfield neural networks [60]	10
BP	1986	-	
RL	1989	2022 Deep reinforcement learning for optimal well control in subsurface systems with uncertain geology [61]	33
Random Forest	1995	2015 Supervised learning to detect salt body [62]	20
SVM	1995	2003 A support vector machine for avo interpretation [63]	8
LSTM	1997	2017 Machine learning can extract the information needed for modelling and data analysing from unstructured documents [64]	20
Transformer	2017	2020 Using Transformer Networks and Knowledge Graphs in Earth Science Literature to Synthesize Mass Information for Transdisciplinary Research [65]	3
GCN	2017	2020 Graph Convolutional Networks for Hyperspectral Image Classification [66]	3
BERT	2018	2021 BERT-E: An Earth Science Specific Language Model for Domain-Specific Downstream Tasks [67]	3
ChatGPT/LLM	2022	2023 K2: A Foundation Language Model for Geoscience Knowledge Understanding and Utilization [15]	1

Table 16: The progression of geoscience research with the use of cutting-edge AI techniques.

B Appendix: GeoCorpus

Here we show the distribution of the collected papers from top-10 amounts journals in geoscience.

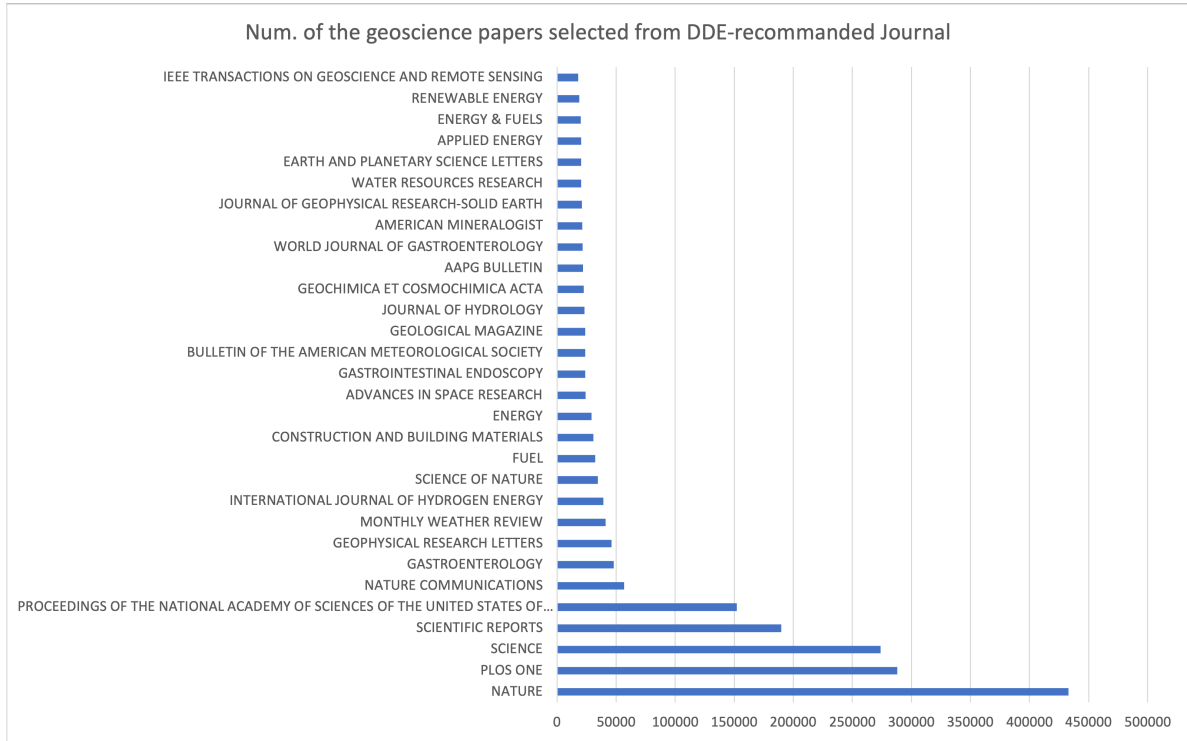


Figure 13: Distribution of the collected papers from top-10 amounts journals in geoscience.

C Appendix: GeoSignal V2 Curation

Below, we will provide a detailed explanation of how we obtain useful supervision signals for geoscience tasks from websites like MinDat, USGS, NGDB, Fossil Ontology, and Fossil calibrations.

C.1 MinDat

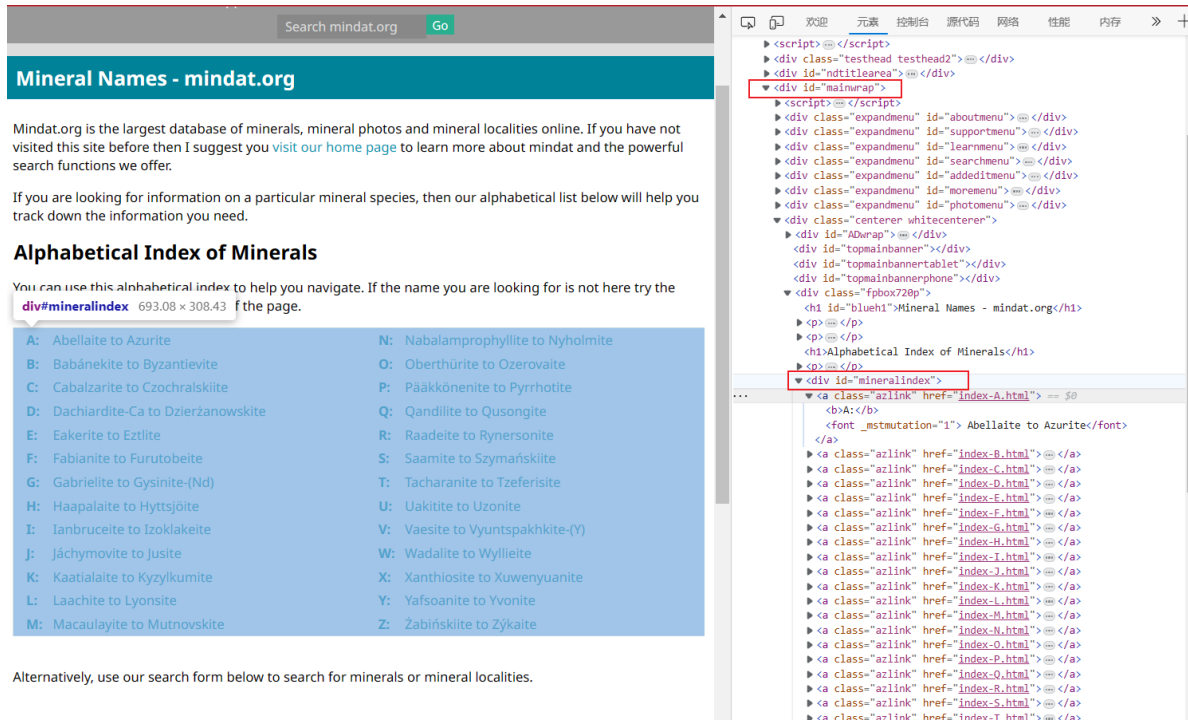


Figure 14: Mines in MinDat.

This website provides a list of all minerals in alphabetical order, and it is possible to obtain the collection of all minerals starting with a certain letter from this page. (As shown in Figure 14) After scraping the HTML page, the following Python spider logic can be used to retrieve the URLs of all minerals starting with letters A-Z: (as shown in Figure 15) `find(id="mainwrap") → find(id="mineralindex")`

To extract the relevant information for each mineral, we can create a dictionary for each mineral containing the following information:

- **Mineral Name.** Key: Mineral, Value: Name of Mineral.
- **Physical Properties.** Extract information on the physical properties of the mineral, including Colour, Lustre, Specific Gravity, Crystal System, Hardness, Name.
- **Chemical Properties.** Extract the chemical elements present in the mineral as Chemical Element and flatten them. For example, Abellaite should be Na, Pb, C, O, H.
- **Type and Occurrence.** Extract information on the Type and Occurrence of the mineral, including Type locality (which will be a separate URL in the format shown below), General Appearance, Place of Conservation.
- **References.** Extract and include the References in the dictionary, with value as a list.

Not all minerals will have all the above information. For minerals with missing information, we will still create a dictionary entry but the value will be "No corresponding information". For example, the dictionary for Abellaite will be illustrated in Figure 16.

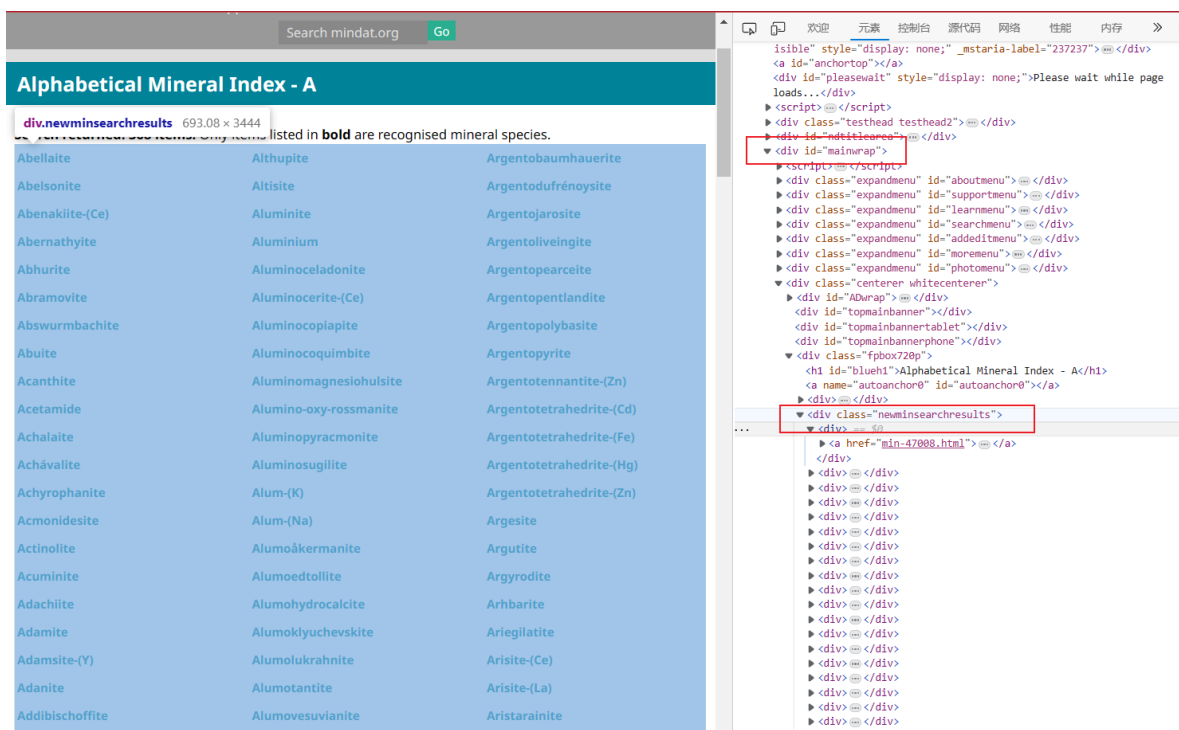


Figure 15: Mines in MinDat.

```
{'Mineral': 'Abellaite',
'Colour': 'Colourless to white',
'Lustre': 'Vitreous',
'Specific Gravity': '5.93 (Calculated)',
'Crystal System': 'Hexagonal',
'Hardness': 'No corresponding information',
'Name': 'Named in honour of Catalan gemmologist Joan Abella i Creus (born 13 December 1968, Sabadell, Catalonia, Spain) who has long studied the minerals of the Eureka mine and also found the mineral.',
'Chemical Element': 'Na,Pb,C,O,H',
'Type Locality': 'Eureka mine, Castell-estao, La Torre de Capdella (La Torre de Cabdella), La Vall Fosca, El Pallars Jussà, Lleida, Catalonia, Spain. The website is https://www.mindat.org/loc-53316.html',
'General Appearance': 'As sparse coatings consisting of subhedral crystals not larger than 10 µm as well as larger hexagonal platelets up to ~30 µm',
'Place of Conservation': 'Cotype material is deposited in the collections of the Natural History Museum of Barcelona, Barcelona, Spain, specimen number MGB 26.350',
'References': ['Ibáñez-Insa, J., Elvira, J.J., Oriols, N., Llovet, X., Viñals, J. (2016) Abellaite, IMA 2014-111. CNMC Newsletter No. 29, February 2016, page 200. Mineralogical Magazine: 80: 199-205.',
'Ibáñez-Insa J., Elvira J.J., Llovet X., Pérez-Cano J., Oriols N., Busquets-Masó M., Hernández S. (2017): Abellaite, NaPb2(CO3)2(OH), a new supergene mineral from the Eureka mine, Lleida province, Catalonia, Spain. European Journal of Mineralogy: 29: 915-922.',
'Siidra O., Nekrasova D., Depmeier W., Chukanov N., Zaitsev A., Turner R. (2018): Hydrocerussite-related minerals and materials: structural principles, chemical variations and infrared spectroscopy. Acta Crystallographica: B74: 182-195.']}

```

Figure 16: Signals in MinDat.

C.2 USGS

	Details	Phases	Magnitudes
Did You Feel It?			
ShakeMap	Magnitude	4.41 mw	
PAGER	Location uncertainty	36.801°N 121.323°W ± 0.1 km	
Technical	Depth uncertainty	10.2 km ± 0.3	
Origin	Origin Time	2023-04-04 22:23:17.820 UTC	
Moment Tensor	Number of Stations	111	
Focal Mechanism	Number of Phases	126	
Waveforms	Minimum Distance	8.5 km (0.08°)	
ShakeAlert®	Travel Time Residual	0.17 s	
Download Event KML	Azimuthal Gap	48°	
View Nearby Seismicity	EE Region	Central California (39)	
Earthquakes	Review Status	REVIEWED	
Hazards	Catalog	NC (nc73866925)	
Data & Products			

<https://earthquake.usgs.gov/earthquakes/eventpage/nc73866925/technical>

Figure 17: USGS collections.

```
{'Name': 'M 4.4 - 1km NNW of Tres Pinos, CA',
'Location': '36.801°N 121.323°W',
'Origin Time': '2023-04-04 22:23:17.820 UTC',
'Minimum Distance': '8.5 km ( 0.08° )',
'Azimuthal Gap': '48°',
'Moment': '5.084e+15 N-m',
'Magnitude': '4.40 Mw',
'Depth': '8.0 km',
'Percent DC': '91%',
'Messages Issued': 'Initial: 5.1 s, Peak: 8.2 s, Final: 31.1 s',
'Magnitude Estimates': 'Initial: M 4.6, Peak: M 4.7, Final: M 4.7, ANSS report: 230 s after origin, M4.7',
'Nearby Cities': ['Hollister, 9 km (6 mi)',
'Salinas, 33 km (20 mi)',
'San Jose, 78 km (49 mi)',
'San Francisco, 145 km (90 mi)']}
```

Figure 18: USGS signals.

When the webpage is opened, it displays important earthquake information from 1900 to 2023 around the world, and the information can be organized by year. For a given year "x", the corresponding webpage for significant earthquakes in that year can be accessed through the URL: <https://earthquake.usgs.gov/earthquakes/browse/significant.php?year=x>

ID number

Search

Results: 335026 records

Lab ID	Source	Map location	State	County	Submitter	Date submitted
62M110	stream/river	Map	AK		Detterman, Robert L.	12/6/1962
62M111	stream/river	Map	AK		Detterman, Robert L.	12/6/1962
62M112	stream/river	Map	AK		Detterman, Robert L.	12/6/1962
62M113	stream/river	Map	AK		Detterman, Robert L.	12/6/1962
63M1132	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1133	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1134	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1135	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1136	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1137	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1138	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1139	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1140	stream/river	Map	AK		Detterman, Robert L.	4/8/1964
63M1141	stream/river	Map	AK		Detterman, Robert L.	4/8/1964

https://mrdata.usgs.gov/ngdb/sediment/show-ngdbsed.php?lab_id=62M112

Figure 19: NGDB collections.

For example, the webpage for significant earthquakes in the year 2023 can be accessed through the URL: <https://earthquake.usgs.gov/earthquakes/browse/significant.php?year=2023>

The webpage for the year 2002 can be accessed through the URL: <https://earthquake.usgs.gov/earthquakes/browse/significant.php?year=2002>

Therefore, using this method, all significant earthquake information pages from 1900 to 2023 can be obtained.

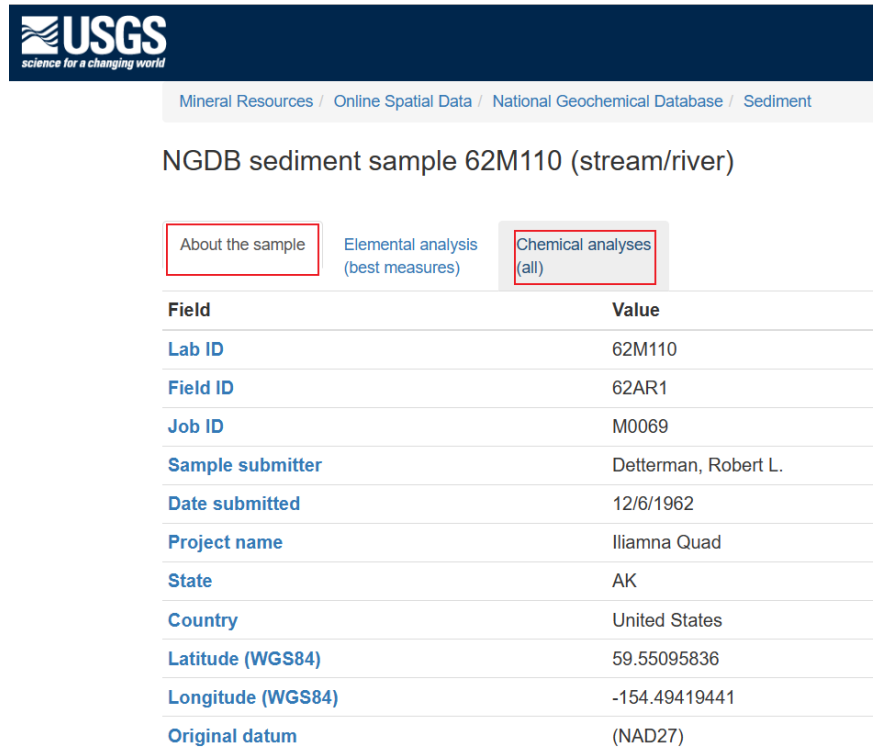
For each specific webpage corresponding to a given year, the list of earthquake URLs can be obtained using the function "find(class='eqitem')".

- **Basic Information.** The basic information that needs to be obtained is the name, with the key value set as "Name" and the value as the name obtained from the webpage title.
- **Origin.** The information needed for this section includes: Location, Origin Time, Minimum Distance, and Azimuthal Gap.
- **Moment Tensor.** The information needed for this section includes: Moment, Magnitude, Depth, and Percent DC.
- **Post ShakeAlert.** The information needed for this section includes: Messages Issued, Magnitude Estimates, and Nearby Cities. Since there may be multiple nearby cities, a list will be created to store them as the value for this specific dictionary entry.

Finally, we obtain the signals in USGS, specifying which signals are retrieved and providing details on how they are processed.

C.3 NGDB

The webpage provides a list of URLs containing information about all the sedimentary rocks. Figure 19



Field	Value
Lab ID	62M110
Field ID	62AR1
Job ID	M0069
Sample submitter	Detterman, Robert L.
Date submitted	12/6/1962
Project name	Iliamna Quad
State	AK
Country	United States
Latitude (WGS84)	59.55095836
Longitude (WGS84)	-154.49419441
Original datum	(NAD27)

Figure 20: Mines in NGDB.

```
{'Lab ID': '62M110',
 'Submitter': 'Detterman, Robert L.',
 'Date submitted': '12/6/1962',
 'State': 'AK',
 'Country': 'United States',
 'Location': '59.55163005, -154.49200195',
 'Location Precision': '1 second',
 'Source': 'stream/river',
 'Chemical': 'Sodium, Magnesium, Potassium, Calcium, Titanium, Iron, '}
```

Figure 21: NGDB final signals.

Information Acquisition:

For each type of sedimentary rock, corresponding information needs to be extracted. The extracted information for each type of sedimentary rock will be combined into one dictionary, and each dictionary will be a list for each type of sedimentary rock. In practice, it is possible that some sedimentary rocks may lack information. In this case, a corresponding dictionary should still be created, with the value set to "No corresponding information." The final data will be stored in a JSON file. The information that needs to be extracted for each type of sedimentary rock is as follows, taking 62M110 as an example.

The information needs to be extracted from "About the sample" and "Chemical analysis" sections (marked in red in the figure).

- **About the sample.** The following information needs to be extracted from this section: Lab ID, Submitter, Date submitted, State, Country, Location (this part is obtained by synthesizing the data from Original Latitude and Original Longitude), Location Precision, and Source.
- **Chemical analysis** This section requires obtaining the major fossil elements of the sedimentary rock, by only recording chemical elements with a composition greater than or equal to 1%.

Finally, we obtain the signals in NGDB, with details on the nature of these signals and the processing steps involved.

C.4 Fossil Ontology

```
{'Lineage': 'Plant,Angiospermae,Magnoliopsida,Vitales,Vitaceae,Tetrastigma,Tetrastigma shantungensis',
'Age': '23.03~5.333',
'Locality': 'Shandong province,China',
'GPS': '(36.412,118.594)',
'House': '中国科学院南京地质古生物研究所',
'References': ['Range-expansion effects on the belowground plant microbiome',
'Mycorrhizal fungi influence global plant biogeography',
'Influences on plant nutritional variation and their potential effects on hominin diet selection',
'Fungi-plant-arthropods interactions in a new conifer wood from the uppermost Permian of China reveal complex ecological relationships and trophic networks',
'Global plant-symbiont organization and emergence of biogeochemical cycles resolved by evolution-based trait modelling']}
```

Figure 22: Fossil ontology signals.

The webpage displays a list of all fossils, so the list can be directly organized. The final data format will be a list, where each element in the list is a dictionary representing a type of fossil, and the data will be saved in a JSON file. Taking *Tetrastigma shantungensis* as an example, below shows the data structure:

- **Basic Information of the Fossil.** Basic Information of the Fossil This part only requires the information of the biological lineage to which the fossil belongs, with key value "Lineage".
- **Geological Information of the Fossil.** This part requires information about the age of the fossil, the location of the fossil, and GPS information. The key values are "Age", "Locality", and "GPS", respectively.
- **Relevant Research Information of the Fossil.** This part requires information about the collection points and reference sources of the fossil, with key values "House" and "References", respectively. Among them, there are many References, which are organized as a list as the value.

Finally, we obtain the signals in Fossil Ontology.

C.5 Fossil calibrations

Starting from the aforementioned URLs, all the specific webpage information about 220 fossil types is included. After grabbing the HTML, please organize and save it accordingly.

To search, we use "find(class='listed-calibrations')", as shown in the following image:

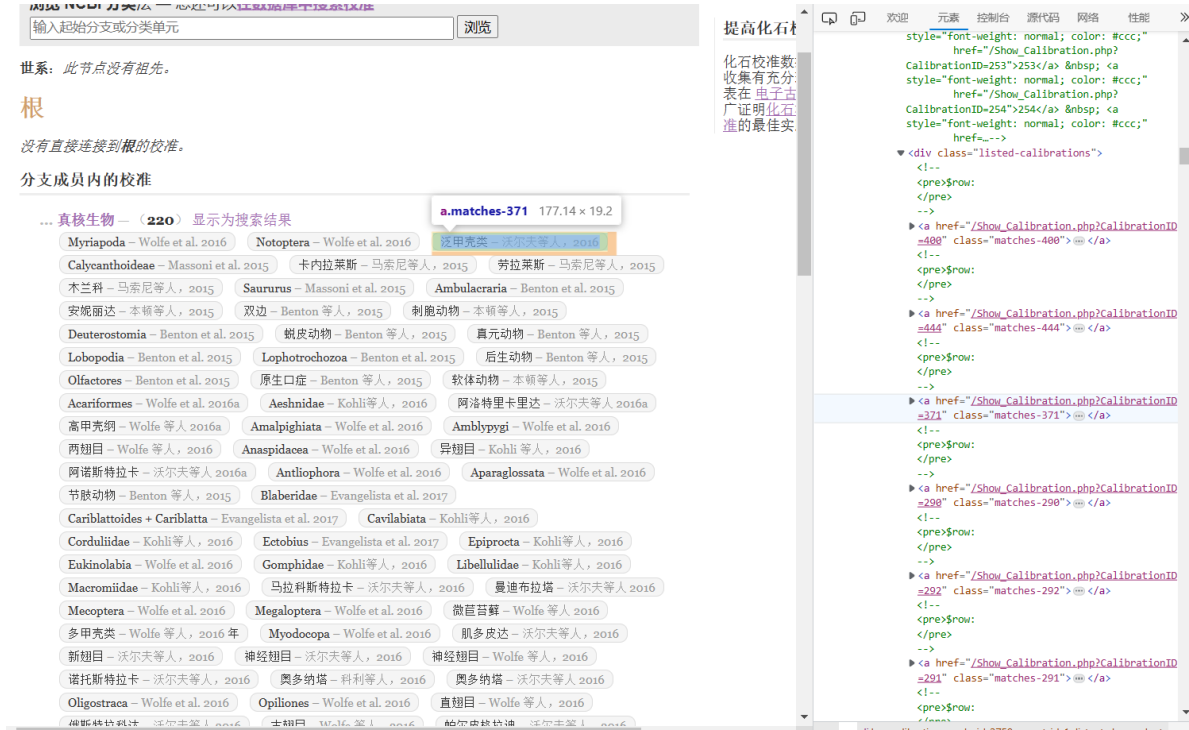


Figure 23: Fossil calibrations collections.

After performing the above search, it is possible to obtain the URLs for all the fossils, and then extract information about each of them one by one.

The final data format will be a list, where each element in the list is a dictionary representing a type of fossil, and the data will be saved in a JSON file. Taking the fossil Sessilia as an example, below shows the information to be extracted and the format:

- **Basic Information.** This section contains two parts: the species to which the fossil belongs, and the biological lineage of the species. The key value for species name is "Name", and the value is the name of the species to which the fossil belongs; The key value for biological lineage is "Lineage", and the value is the biological lineage of the species.
- **Age Information.** This part includes the earliest and latest time of the fossil, with key values "Minimum age" and "Maximum age", respectively.
- **Key Fossil Information.** This part mainly concerns the key fossils used to determine the node period, including the location and geological age, with key values "Locality" and "Geological age" respectively.
- **References.** This part mainly includes the calibration and reference sources of the fossil information, with key values "Calibration" and "Reference", respectively.

Finally, we obtain the signals in Fossil Calibrations.

```
{'Name': 'Sessilia',
 'Lineage': 'Eukaryota,Opisthokonta,Metazoa,Eumetazoa,Bilateria,Coelomata,Protostomia,Ecdysozoa,Panarthropoda,Arthropoda,Mandibulata,Pancrustacea,Crustacea,Maxillopoda,Thecostraca,Cirripedia,Thoracica,Sessilia',
 'Minimum age': '89.5Ma',
 'Maximum age': '521Ma',
 'Locality': 'Pit No. 125, Froxfield, Hampshire, England',
 'Geological age': 'Cretaceous, Mesozoic',
 'Calibration': 'Wolfe, J.M., Daley, A.C., Legg, D.A., and Edgecombe, G.D. 2016. Fossil calibrations for the arthropod tree of life. Earth Science Reviews 160: 43-110.',
 'Reference': 'Gale, A.S., Sørensen, A.M. 2015. Origin of the balanomorph barnacles (Crustacea, Cirripedia, Thoracica): new evidence from the Late Cretaceous (Campanian) of Sweden. J. Syst. Palaeontol. 13, 791-824.'}
```

Figure 24: Fossil calibrations signals.

D Appendix: Prompts

Below, we will list the prompt words we use for constructing Instruction tuning data.

deepliterature.abstract.title

1. **Input:** What is a suitable title for this geoscience paragraph: [input]
Output: The title can be '[output]'
2. **Input:** Provide a fitting title for this geoscience text: [input]
Output: One possible title can be '[output]'
3. **Input:** Suggest a title for this geoscience passage: [input]
Output: '[output]' can be a suitable title.
4. **Input:** What would be an appropriate title for this geoscience paragraph: [input]
Output: Here is a title for you to consider: [output].
5. **Input:** What should be the title of this geoscience text considering its content: [input]
Output: '[output]' is one option.

deepliterature.abstract.keyword

1. **Input:** Extract the main keywords from this geoscience paragraph: [input]
Output: The main keywords in the paragraph can be listed as follows: [output]
2. **Input:** Identify the key terms in this geoscience passage: [input]
Output: The key terms in the passage can be identified as: [output]
3. **Input:** What are the most important keywords in this geoscience text: [input]
Output: The most important keywords in the text are: [output]
4. **Input:** What are the fields of study connected to this geoscience text: [input]
Output: The fields of study connected to this text include: [output]
5. **Input:** Determine the fields of research related to this geoscience paragraph: [input]
Output: The related fields of research for this paragraph are: [output]

deepliterature.reference.resolution

1. **Input:** When you see the sentence like [input], can you recommend some related papers for me?
Output: Alright, I have found a relevant paper for you: \n[output]
2. **Input:** Based on the content of this sentence [input], could you suggest a few relevant papers?
Output: Okay, I've got the related paper right here: \n[output]
3. **Input:** Considering the information in this sentence [input], what papers would you recommend?
Output: Sure, here's the associated paper you're looking for: \n[output]
4. **Input:** Given the context of this sentence [input], can you provide a list of related papers?
Output: No problem, I have the corresponding paper for you: \n[output]
5. **Input:** After reading this sentence [input], what are some papers you would recommend?
Output: Certainly, I'm showing you the related paper now: \n[output]

gso.wikipedia.title

1. **Input:** Can you explain the concept of [title] in relation to [keyword]?
Output: Certainly, when discussing [title] in the context of [keyword], [paragraph]
2. **Input:** Please share information about [title] concerning [keyword].
Output: Of course, regarding [title] and its connection to [keyword], [paragraph]
3. **Input:** What does the [title] section cover about the topic [keyword]?
Output: In the [title] section about [keyword], we have: [paragraph]
4. **Input:** Please describe [title] as it pertains to [keyword].
Output: Certainly, in terms of [title] and its relationship to [keyword], [paragraph]
5. **Input:** Can you elaborate on [title] in connection with [keyword]?
Output: Of course, when discussing [title] and its connection to [keyword], [paragraph]

gso.wikipedia.entity

1. **Input:** Can you identify the key terms related to earth science mentioned in this passage? The passage: [input]
Output: Sure, the terms related to earth science mentioned in the passage include: [output]
2. **Input:** What are the terms connected to earth science present in this passage? The passage: [input]
Output: The terms connected to earth science in the passage are: [output]
3. **Input:** Please list the terms associated with earth science found in the given passage. The passage: [input]
Output: The terms associated with earth science in the passage are as follows: [output]
4. **Input:** Can you help me find the terms relevant to earth science in this passage? The passage: [input]
Output: Sure, the terms relevant to earth science in the passage are: [output]
5. **Input:** Identify the terms in the passage that are related to earth science. The passage: [input]
Output: The terms in the passage related to earth science are: [output]

gso.wordnet.description

1. **Input:** Can you provide a description for the word '[input]'?
Output: Sure, the description for the word '[input]' is: [output]
2. **Input:** What does the word '[input]' mean?
Output: The meaning of the word '[input]' is: [output]
3. **Input:** Please explain the meaning of the term '[input]'.
Output: The term '[input]' means: [output]
4. **Input:** What can you tell me about the meaning of the word '[input]'?
Output: The word '[input]' has the following meaning: [output]
5. **Input:** Help me understand the definition of the term '[input]'.
Output: The definition of the term '[input]' is: [output]

gso.wordnet.synonym

1. **Input:** Can you provide a synonym for the word '[input]'?
Output: Sure, a synonym for the word '[input]' is: [output]
2. **Input:** What is a synonym of the word '[input]'?
Output: A synonym of the word '[input]' is: [output]
3. **Input:** Please give me a synonym for the term '[input]'.
Output: A synonym for the term '[input]' is: [output]
4. **Input:** What would be a synonym for the word '[input]'?
Output: A synonym for the word '[input]' is: [output]
5. **Input:** Help me find a synonym for the term '[input]'.
Output: A synonym for the term '[input]' is: [output]

gso.dictionary.definition

1. **Input:** Can you provide the specialized definition of the term '[input]'?
Output: Sure, in the context of geoscience, the term '[input]' is defined as: [output]
2. **Input:** What does the term '[input]' mean in a specialized context?
Output: In the field of geoscience, the term '[input]' means: [output]
3. **Input:** Please explain the meaning of the specialized term '[input]'.
Output: The specialized term '[input]', in the context of geoscience, means: [output]
4. **Input:** Can you suggest the meaning of the term '[input]' in a specific field?
Output: In the field of geoscience, the meaning of the term '[input]' is: [output]
5. **Input:** What would be the definition of the word '[input]' in a specialized context?
Output: In the context of geoscience, the definition of the word '[input]' is: [output]

gso.dictionary.synonym

1. **Input:** Can you provide a synonym for the word '[input]'?
Output: Sure, you might [output]

2. **Input:** What is a synonym of the word '[input]'?
Output: You can [output]
3. **Input:** Please give me a synonym for the term '[input]'.
Output: Of course, you should [output]
4. **Input:** What would be a synonym for the word '[input]'?
Output: One can [output]
5. **Input:** What's a synonym for the word '[input]'?
Output: You can [output]

gso.dictionary.classification

1. **Input:** Identify the scientific discipline associated with the term '[input]':
Output: [output]
2. **Input:** In which scientific discipline would you find the term '[input]'?
Output: [output]
3. **Input:** The term '[input]' is primarily associated with which of these disciplines?
Output: [output]
4. **Input:** In the realm of Earth sciences, which discipline is most closely linked to the term '[input]'?
Output: [output]
5. **Input:** Which scientific discipline is most relevant when discussing the term '[input]'?
Output: [output]

gso.taxonomy.hyponymy.child

1. **Input:** What are some subfields or specific topics that fall under the broader concept of [parent]?
Output: Some subfields or specific topics that fall under the broader concept of [parent] in geoscience include [child].
2. **Input:** What are some specific examples of [parent] within geoscience?
Output: Some specific examples of [parent] within geoscience include [child].
3. **Input:** What are some narrower categories or subdisciplines within [parent]?
Output: Some narrower categories or subdisciplines within [parent] in geoscience include [child].
4. **Input:** What are some specific types or varieties of [parent] in geoscience?
Output: Some specific types or varieties of [parent] in geoscience include [child].
5. **Input:** What are some specific techniques or methodologies used to study [parent] in geoscience?
Output: Some specific techniques or methodologies used to study [parent] in geoscience include [child].

gso.taxonomy.hyponymy.parent

1. **Input:** What is the overarching category that [child] belongs to?
Output: It can be classified under [parent], which is a broad field of study within geoscience.
2. **Input:** What are the subfields that fall under [parent]?
Output: Some subfields that are part of [parent] include [child], among others.
3. **Input:** What are some related concepts to [parent] in geoscience?
Output: Some related concepts to [parent] in geoscience are [child], among others.
4. **Input:** What is the specific category that [child] belongs to within geoscience?
Output: The specific category that [child] belongs to within geoscience is [parent].
5. **Input:** What are the different branches of geoscience that [child] is part of?
Output: Some branches of geoscience that [child] is part of include [parent], among others.

gso.taxonomy.hyponymy.judgment.parent

1. **Input:** Is [child] a child concept of [parent] in geoscience?
Output: [Answer]
2. **Input:** Does [child] fall under the broader concept of [parent] in geoscience?
Output: [Answer]

3. **Input:** Is [child] a narrower category or subdiscipline within [parent] in geoscience?
Output: [Answer]
4. **Input:** Are [child] specific types or varieties of [parent] in geoscience?
Output: [Answer]
5. **Input:** Do scientists use [child] as a technique or methodology to study [parent] in geoscience?
Output: [Answer]

gso.taxonomy.hyponymy.judgment.child

1. **Input:** Is [parent] the broader concept that encompasses [output] in geoscience?
Output: [Answer]
2. **Input:** Does [parent] serve as the parent category for [output] in geoscience?
Output: [Answer]
3. **Input:** Is [parent] a higher-level concept that includes [output] in geoscience?
Output: [Answer]
4. **Input:** Does [parent] have a broader field of study than [child] within geoscience?
Output: [Answer]
5. **Input:** Is [parent] a common parent concept of [child] within geoscience?
Output: [Answer]

metaearth

1. **Input:** What's the [key] of the [object]?
Output: The [object]'s [key] is [value].
2. **Input:** Can you tell me the [key] of the [object]?
Output: Sure, the [object]'s [key] is [value].
3. **Input:** What is the [object]'s [key]?
Output: The [key] of the [object] is [value].
4. **Input:** What would be the [key] of the [object]?
Output: The [object]'s [key] is [value].
5. **Input:** I'd like to know the [key] of the [object].
Output: The [object] has a [key] of [value].

gakg.qa

1. **Input:** What's the [key] of the paper according to the abstract: \n[object]?
Output: The paper's [key] is [value].
2. **Input:** Can you tell me the [key] of the paper according to the abstract: \n[object]?
Output: Sure, the paper's [key] is [value].
3. **Input:** What is the paper's [key]? According to the abstract: \n[object]
Output: The [key] of the paper is [value].
4. **Input:** Please inform me of the paper's [key]. According to the abstract: \n[object]
Output: The [key] associated with the paper is [value].
5. **Input:** I'd like to know the [key] of the paper According to the abstract: \n[object].
Output: The paper has a [key] of [value].

ner

1. **Input:** Please extract the named entities in: [input].
Output: The named entities in the passage can be listed as follows: \n[output]
2. **Input:** What are the named entities mentioned in: [input]?
Output: The named entities mentioned in the passage are: \n[output]
3. **Input:** Can you identify the named entities in: [input]?
Output: Sure, the named entities in the passage are: \n[output]

4. **Input:** What are the geological terms or concepts mentioned in: [input]?
Output: The geological terms or concepts mentioned in the passage are: \n [output]
5. **Input:** What are the names of any geological formations mentioned in: [input]?
Output: The names of the geological formations mentioned in the passage are: \n [output]

gakg.illustration

1. **Input:** What kind of illustration when you see these content? [input]
Output: [output]
2. **Input:** What type of illustration comes to mind when you observe the following content? [input]
Output: [output]
3. **Input:** "What sort of illustration is associated with this content? [input]
Output: [output]
4. **Input:** What kind of illustration can you imagine when presented with this content? [input]
Output: [output]
5. **Input:** When examining the following content, what type of illustration would you associate with it? [input]
Output: [output]

gakg.table

1. **Input:** What kind of table when you see these elements? [input]
Output: [output]
2. **Input:** What type of table comes to mind when you observe the following content? [input]
Output: [output]
3. **Input:** What sort of table is associated with this elements? [input]
Output: [output]
4. **Input:** What kind of table can you imagine when presented with this content? [input]
Output: [output]
5. **Input:** Upon seeing the content, what kind of table do you think of? [input]
Output: [output]

E Appendix: Training setup

- Model setup (30B parameters)
 - num layers: 48
 - num attention heads: 56
 - hidden size: 7168
 - max position embeddings: 2048
 - layernorm epsilon: $1e - 5$
- Regularization Setup
 - optimizer: Adam
 - attention dropout: 0.1
 - hidden dropout: 0.1
 - weight decay: 0.1
 - clip-grad: 1.0
 - adam β_1 : 0.9
 - adam β_2 : 0.95
 - adam ϵ : $1e - 8$
- Training Setup
 - micro-batch-size: 1
 - global-batch-size: 4096
 - recompute-activations: True (gradient checkpointing)
 - train-samples: 30M (60B token)
 - disable-bias-linear: True (turn off the bias of *nn.linear*)
 - seed: 42
 - save-interval: 100
- Learning Rate Setup
 - lr-decay-style: linear
 - lr-warmup-steps: 100
 - lr: $1e - 5$
 - min-lr: $1e - 7$
- Mixed Precision Setup
 - FP16: False
 - BF16: False
- Parallel Configuration
 - tensor-model-parallel-size: 4
 - pipeline-model-parallel-size: 16
 - distributed-backend: NCCL
 - sequence-parallel: True

Notices: With a model parallel size (TP) of 4 and a pipeline parallel size (PP) of 16, it can be considered that a 30B model with 48 layers is divided into 16 parts with 3 layers in each part, and each of the 4 accelerator cards in a node is responsible for processing 3 continuous layer.

F Appendix: Model Card

Our model is based on Galactica, a standard GPT2 structure with 48 transformer blocks and 56 attention heads per layer, with a hidden dim of 7,168. The parameters and tokenizer of our model are initialized from the hugging face release checkpoint of Galactica-30B, and the maximum input length is 2048 and no bias settings are followed. Although Galactica-30B is an fp16 model, to ensure the stability of the training process, we used fp32 to train and save model parameters.

Model Card - GEOGALACTICA	
Model Details	
• Developed by: Shanghai Jiao Tong University and Deep-time Digital Earth Science Center.	
• Shared by: Shanghai Jiao Tong University and GeoBRAIN.ai.	
• Model type: Further pre-train and Supervised Fine-tuning.	
• Language(s) (NLP): English.	
• License: Apache License 2.0.	
• Further pre-train from model: Galactica [12].	
Model Sources	
• Repository: https://github.com/geobrain-ai/geogalactica	
• Paper: GEOGALACTICA: A Scientific Large Language Model in Geoscience	
Intended Use	
• Research Assistance: Providing support in academic and industrial research by summarizing current scientific literature, suggesting hypotheses, and identifying gaps in existing research.	
• Educational Tool: Serving as an educational resource for students and professionals in geosciences, offering explanations of complex concepts and providing interactive learning experiences.	
• Collaboration and Communication: Facilitating collaboration among geoscientists by providing a platform for sharing data and insights, and helping in communicating complex geoscientific information to non-experts.	
Ethical Considerations	
• This model inherits from Galactica [12], and in the training corpus, we have conducted sufficient data governance to ensure that the training data embodies geographical community, transparency, inclusiveness, respect for privacy, and topic neutrality.	
Training Data	
• Further pre-train: A geoscience-related text corpus containing 65 billion tokens , preserving as the largest geoscience-specific text corpus.	
• Supervised Fine-tuning: daven3/geosignal.	
• Tool-Augmented Learning: zthang/geotools.	
Evaluation Data	
• MMLU: Massive Multitask Language Understanding, a large-scale research initiative aimed at improving language models' understanding and reasoning abilities across a diverse range of subjects and tasks.	
• GeoBench: The benchmark mentioned in K2 [15]. The data can be access on [daven3/geobench](https://huggingface.co/datasets/daven3/geobench).	
• Human Evaluation: Selected questions.	
Model Card Contact	
GEOGALACTICA is a research preview intended for non-commercial use only. Please contact us if you find any issues. For details, you can email via davendw@sjtu.edu.cn .	

Figure 25: Model Card for GEOGALACTICA.

G Appendix: Evaluation

G.1 Open-ended Tasks

G.1.1 Noun Definition

For a better understanding of the words, we added the definitions from the Geology Dictionary.

- 1 Physical weathering and chemical weathering are processes that break down rocks and minerals through mechanical and chemical means, respectively.
- 2 Sedimentary differentiation (*the sorting and separation of sediments during the formation of sedimentary rocks*).
- 3 A continental margin (*the boundary between a continent and an ocean, characterized by various geological features*).
- 4 Seafloor spreading (*the process where new oceanic crust is formed at mid-ocean ridges as tectonic plates move apart*).
- 5 Stratum occurrence (*relates to the presence and distribution of rock layers in a particular geological context*).
- 6 Normal fault (*a type of fault where the hanging wall moves downward relative to the footwall*).
- 7 Plate tectonics (*the theory that describes the movement and interaction of Earth's lithospheric plates*).
- 8 Continental margin (*the boundary between a continent and an ocean, characterized by various geological features*).
- 9 The Yanshan Movement (*refers to a tectonic event in China that resulted in significant geological changes*).
- 10 Continental margin (*the boundary between a continent and an ocean, characterized by various geological features*).
- 11 A united paleocontinent (*a reconstructed ancient landmass formed by merging the continents' positions from the distant past*).
- 12 Earth resources (*natural materials and substances that are valuable to humans for various purposes*).
- 13 Reverse fault (*a type of fault where the hanging wall moves upward relative to the footwall*).
- 14 Ediacara fauna (*refers to a group of early, soft-bodied, and mostly extinct organisms from the Ediacaran Period*).
- 15 Fingerfacies fossil (*represents a specific type of fossil that provides information about sedimentary environments and conditions*).
- 16 Walther's Law (*describes the vertical succession of sedimentary rock layers that were originally deposited in lateral proximity*).
- 17 Vertical accumulation (*is the process of sediment layers accumulating on top of each other over time*).
- 18 The original levelness principle (*suggests that sedimentary layers were initially deposited horizontally*).
- 19 The stratum overlap principle (*relates to the idea that younger sedimentary layers can cover or overlap older ones*).
- 20 Lateral accumulation (*refers to the process of sedimentary material accumulating horizontally, typically in a depositional environment*).

G.1.2 Beginner Level Q&A

- 1 How does a cloud fill up with water?
- 2 How does diffraction make a tree's shadow blurry?
- 3 How does trash in the ocean disappear?
- 4 How does water dowsing work?
- 5 How does wind create all the ocean currents?
- 6 If I jump, will the entire earth move a little bit?
- 7 If I were able to dig a hole from the U.S. through the center of the earth, what part of China would I end up in?
- 8 Is a quadruple rainbow possible?
- 9 What causes the water going down a drain to swirl clockwise in the northern hemisphere and counter-clockwise in the southern hemisphere?
- 10 What keeps the continents floating on a sea of molten rock?

G.1.3 Intermediate Level Q&A

- 1 How does the movement of tectonic plates contribute to the formation of earthquakes and volcanic activity?
- 2 What are the main factors that influence the formation and intensity of hurricanes in the Atlantic Ocean?
- 3 How does the process of erosion shape the landscape and contribute to the formation of features such as canyons and valleys?
- 4 What are the primary mechanisms responsible for the formation and movement of glaciers?
- 5 How do ocean currents influence the distribution of marine organisms and impact the productivity of marine ecosystems?
- 6 What factors contribute to the formation and intensity of tornadoes in regions prone to severe weather events?
- 7 How do geological processes such as weathering and sedimentation contribute to the formation and transformation of soil?
- 8 What role does the Earth's magnetic field play in protecting the planet from harmful solar radiation?
- 9 How do variations in atmospheric pressure and temperature contribute to the formation and behavior of weather systems?
- 10 What are the primary processes responsible for the formation and transformation of different types of rocks, such as igneous, sedimentary, and metamorphic rocks?

G.1.4 Advanced Level Q&A

- 1 How did Earth and other planets form? Were planets formed in situ?
- 2 Was there ever a collision of the Earth with another planet Theia, giving birth to our satellite?
- 3 What is the long-term heat balance of Earth?
- 4 What made plate tectonics a dominant process only on Earth?
- 5 How inherent to planetary evolution is the development of life conditions?
- 6 As planets age and cool off, their internal and surface processes coevolve, chemically and mechanically, shaping the atmospheric composition. What are the chemical composition and mechanical properties of rocks in the Earth's mantle at the extreme pressure and temperature they undergo?
- 7 What are the dynamic processes in the Earth's interior that accommodate and fuel plate tectonics?
- 8 How does the geomagnetic field link to the iron convection properties at the deep Earth?
- 9 Are intraplate hotspots made by deep sources of uprising materials (mantle plumes) coming from the deepest Earth's mantle?"

G.2 Functional Tasks

G.2.1 Knowledge-based associative judgment question.

- 1 What is the specific category that magnetostratigraphy belongs to within geoscience?
- 2 What is the overarching category that magnetic polarity stratigraphy belongs to?
- 3 What are the subfields that fall under magnetic polarity stratigraphy?
- 4 What are some related concepts to geomagnetic polarity in geoscience?
- 5 What are some related concepts to geomagnetic polarity in geoscience?
- 6 What is the specific category that transitional polarity belongs to within geoscience?
- 7 What are the subfields that fall under magnetic polarity stratigraphy?
- 8 What are the subfields that fall under magnetostratigraphic polarity units?
- 9 What are the subfields that fall under magnetostratigraphic polarity units?
- 10 What are the subfields that fall under magnetostratigraphic polarity units?

G.2.2 Research Paper Proposition Task.

Here are two examples, and we will release the whole data on Github:

- **Abstract:** The Wenchuan Earthquake on 12 May 2008 triggered a large number of geo-hazards including landslides, slope collapses and debris flows. Field investigations and remote-sensing interpretation identified 11,308 geo-hazards in 16 seriously damaged counties in Sichuan Province, southwest China. The paper reports an analysis of the distribution of these geo-hazards, particularly the earthquake-triggered landslides. Not surprisingly, the most significant geo-hazards were related to the main fault and on the hanging-wall side, although some occurred in deeply incised river gorges further away from the main rupture zone. Due to the high seismic intensity of the earthquake, most of the large landslides moved at high speed and for considerable distances.

Title: Analysis of the geo-hazards triggered by the 12 May 2008 Wenchuan Earthquake, China

- **Abstract:** The Modern-Era Retrospective Analysis for Research and Applications-2 (MERRA2) version of the Goddard Earth Observing System-5 (GEOS-5) atmospheric general circulation model (AGCM) is currently in use in the NASA Global Modeling and Assimilation Office (GMAO) at a wide range of resolutions for a variety of applications. Details of the changes in parameterizations after the version in the original MERRA reanalysis are presented here. Results of a series of atmosphere-only sensitivity studies are shown to demonstrate changes in simulated climate associated with specific changes in physical parameterizations, and the impact of the newly implemented resolution-aware behavior on simulations at different resolutions is demonstrated. The GEOS-5 AGCM presented here is the model used as part of the GMAO MERRA2 reanalysis, global mesoscale simulations at 10 km resolution through 1.5 km resolution, the real-time numerical weather prediction system, and for atmosphere-only, coupled ocean-atmosphere and coupled atmosphere-chemistry simulations. The seasonal mean climate of the MERRA2 version of the GEOS-5 AGCM represents a substantial improvement over the simulated climate of the MERRA version at all resolutions and for all applications. Fundamental improvements in simulated climate are associated with the increased re-evaporation of frozen precipitation and cloud condensate, resulting in a wetter atmosphere. Improvements in simulated climate are also shown to be attributable to changes in the background gravity wave drag, and to upgrades in the relationship between the ocean surface stress and the ocean roughness. The series of resolution-aware parameters related to the moist physics was shown to result in improvements at higher resolutions and result in AGCM simulations that exhibit seamless behavior across different resolutions and applications..

Title: GMD - Development of the GEOS-5 atmospheric general circulation model: evolution from MERRA to MERRA2

G.2.3 Geoscience Research Functionality

Here we share the five papers used to evaluate the geoscience research functionality of the LLMs in the citation form of MLA:

- 1 Zheng, Yadong, et al. "A challenge to the concept of slip-lines in extrusion tectonics." *Geoscience Frontiers* 2.1 (2011): 23-34.
- 2 Fedo, Christopher M., H. Wayne Nesbitt, and Grant M. Young. "Unraveling the effects of potassium metasomatism in sedimentary rocks and paleosols, with implications for paleo weathering conditions and provenance." *Geology* 23.10 (1995): 921-924.
- 3 Owens, Jeremy D., et al. "Iron isotope and trace metal records of iron cycling in the proto-North Atlantic during the Cenomanian-Turonian oceanic anoxic event (OAE-2)." *Paleoceanography* 27.3 (2012).
- 4 Percival, L. M. E., et al. "Determining the style and provenance of magmatic activity during the Early Aptian Oceanic Anoxic Event (OAE 1a)." *Global and Planetary Change* 200 (2021): 103461.
- 5 Jenkyns, Hugh C. "Transient cooling episodes during Cretaceous Oceanic Anoxic Events with special reference to OAE 1a (Early Aptian)." *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376.2130 (2018): 20170073.

H Generation Examples

H.1 Noun Definition

We present an example with prompt "What is physical weathering and chemical weathering?" as follows:

- **ChatGPT:**

Physical weathering and chemical weathering are two processes that contribute to the breakdown and alteration of rocks and minerals on Earth's surface. Let's explore each of them:

1. **Physical Weathering:** Physical weathering, also known as mechanical weathering, refers to the disintegration of rocks and minerals into smaller fragments without changing their chemical composition. It is primarily caused by physical forces and environmental factors.

2. **Chemical Weathering:** Chemical weathering involves the alteration or decomposition of rocks and minerals through chemical reactions. It occurs when minerals in rocks come into contact with water, atmospheric gases, or other substances that cause them to dissolve, decompose, or undergo chemical changes.

It's important to note that physical and chemical weathering often occur simultaneously and can complement each other, accelerating the overall weathering process. Additionally, the type and rate of weathering are influenced by factors such as climate, rock composition, topography, and the presence of vegetation or organisms.

- **GEOGALACTICA**

Physical weathering and chemical weathering are two of the main processes that cause rocks and minerals to break down and change over time.

Physical weathering refers to the breakdown of rocks and minerals due to the action of natural forces such as wind, water, ice, and gravity. These forces can cause rocks to crack, crumble, and disintegrate into smaller pieces. For example, when water flows over the surface of a rock, it can erode and wear away the rock's surface, causing it to break apart.

Chemical weathering, on the other hand, refers to the breakdown of rocks and minerals due to the action of chemical reactions. These reactions occur when water, air, and other chemicals interact with the rocks and minerals, causing them to dissolve, decompose, or react with each other. For example, when carbon dioxide in the air dissolves in water, it forms carbonic acid, which can dissolve calcium carbonate in limestone, causing it to break down into calcium and bicarbonate ions.

In summary, physical weathering refers to the breakdown of rocks and minerals due to the action of natural forces, while chemical weathering refers to the breakdown of rocks and minerals due to the action of chemical reactions. Both processes play an important role in shaping the Earth's landscape over time.

H.2 Beginner Level Q&A

We present an example with prompt "How does a cloud fill up with water?" as follows:

- **ChatGPT:**

"A cloud forms and "fills up" with water through a process known as condensation. It involves the transformation of water vapor into liquid water droplets or ice crystals, which are the visible components of clouds.

Here's a simplified explanation of the process:

1. **Evaporation:** The sun's heat causes water from various sources like oceans, lakes, and rivers to evaporate, turning it into water vapor. This water vapor is essentially invisible.

2. **Moisture in the air:** As the water evaporates, it rises into the atmosphere and mixes with the surrounding air. The capacity of the air to hold water vapor depends on factors such as temperature and humidity.

3. **Cooling and saturation:** As the warm, moist air rises higher into the atmosphere, it encounters cooler temperatures. Cooler air cannot hold as much moisture as warm air, so the water vapor begins to condense into tiny water droplets or ice crystals around microscopic particles in the air, such as dust, salt, or pollution. This process is called saturation.

4. Cloud formation: The condensed water droplets or ice crystals cluster together to form visible clouds. Clouds consist of countless tiny droplets or ice crystals suspended in the air, which give them their characteristic appearance.

The process of cloud formation is influenced by various factors such as air temperature, humidity levels, and the presence of condensation nuclei (particles that water vapor can condense onto). Different types of clouds, such as cumulus, stratus, or cirrus, form under specific atmospheric conditions and have distinct appearances.

It's important to note that while clouds contain water, they don't "fill up" like a container. The amount of water vapor in the atmosphere continuously changes as a result of various meteorological processes, including evaporation, condensation, and precipitation. Clouds form, evolve, and dissipate as these processes occur within the atmosphere."

- **GEOGALACTICA**

A cloud fills up with water when water vapor in the air cools and condenses into tiny droplets or ice crystals. This happens when the temperature of the air drops below the dew point or frost point, which is the temperature at which the air can hold the maximum amount of water vapor.

When the temperature drops, the water molecules in the air slow down and collide with each other more frequently. This causes the water vapor to condense into tiny droplets or ice crystals, forming clouds. The more water vapor in the air, the more likely it is that the temperature will drop low enough for condensation to occur.

Clouds can also form when warm, moist air rises and cools. As the air rises, it expands and cools, causing the water vapor in the air to condense into tiny droplets or ice crystals. This process is known as adiabatic cooling.

In summary, a cloud fills up with water when water vapor in the air cools and condenses into tiny droplets or ice crystals. This can happen when the temperature of the air drops below the dew point or frost point, or when warm, moist air rises and cools.

H.3 Intermediate Level Q&A

We present an example with prompt "How does the movement of tectonic plates contribute to the formation of earthquakes and volcanic activity," as follows:

- **ChatGPT:**

"The movement of tectonic plates is a key factor in the formation of earthquakes and volcanic activity. The Earth's lithosphere, which consists of several large and numerous smaller tectonic plates, is divided into rigid pieces that float and move on the underlying semi-fluid asthenosphere. There are three main types of plate boundaries where interactions between these plates occur: divergent boundaries, convergent boundaries, and transform boundaries. At divergent boundaries, tectonic plates move away from each other, creating a gap or rift. This movement results in the upwelling of magma from the mantle, forming new crust. As the magma rises to the surface, it can lead to volcanic activity, with molten rock erupting to form volcanoes. However, the volcanic activity at divergent boundaries is generally less explosive compared to other types of volcanic activity. At convergent boundaries, tectonic plates collide with each other. There are three types of convergent boundaries: oceanic-oceanic, oceanic-continental, and continental-continental. When an oceanic plate converges with another oceanic plate or a continental plate, the denser oceanic plate is usually subducted (forced beneath) into the mantle. The subduction process can generate intense pressure and friction, resulting in the release of seismic energy and causing earthquakes. Additionally, as the subducted plate sinks into the mantle, it can melt, generating magma that rises to the surface and forms volcanoes. These volcanoes can be explosive and often create volcanic arcs, such as the Pacific Ring of Fire. At transform boundaries, tectonic plates slide past each other horizontally. The friction between the plates can cause them to become locked, preventing smooth movement. As the plates continue to push against each other, stress builds up along the locked boundary. Eventually, the stress overcomes the friction, leading to a sudden release of energy, resulting in earthquakes. Transform boundaries are primarily characterized by seismic activity rather than volcanic activity. In summary, the movement of tectonic plates at divergent boundaries leads to volcanic activity, while at convergent boundaries, it results in both earthquakes and volcanic activity. Transform boundaries primarily contribute to seismic activity. These interactions between tectonic plates play a fundamental role in shaping Earth's surface and are responsible for the occurrence of earthquakes and the formation of volcanic features around the world."

- **GEOGALACTICA**

"The movement of tectonic plates is the driving force behind the formation of earthquakes and volcanic activity on Earth. Tectonic plates are large, rigid sections of the Earth's crust that move slowly over the surface of the planet. The movement of these plates is caused by convection currents in the Earth's mantle, which is the layer beneath the crust and above the core.

When two tectonic plates move towards each other, they may collide, causing one plate to be forced under the other. This process is known as subduction, and it is responsible for the formation of many of the world's mountain ranges, such as the Himalayas and the Andes. When two plates collide, the pressure can build up until it is released in a sudden, violent event known as an earthquake.

On the other hand, when two plates move in opposite directions, one may be forced upwards, causing the Earth's crust to stretch and thin. This process is known as rifting, and it is responsible for the formation of many of the world's rift valleys, such as the Great Rift Valley in Africa. When the crust is stretched too far, it can break apart, leading to the formation of new ocean basins.

Volcanic activity is also closely linked to the movement of tectonic plates. When two plates move towards each other, the pressure can build up until it is released in a sudden, violent event known as an earthquake. This can cause the Earth's crust to crack."

H.4 Advanced Level Q&A

We present an example with prompt "How did Earth and other planets form? Were planets formed in situ," as follows:

- **ChatGPT:**

"The formation of planets depends on several factors, including the gravitational pull of nearby celestial bodies, the temperature and pressure of the surrounding environment, and the availability of necessary resources for planet formation. Planets may form in situ (i.e., within their own gravitational field) or through the interaction with other celestial bodies. Some planets may form from dust particles in space, while others may form around existing celestial bodies such as comets or asteroids. The process of planet formation involves the accumulation of material over time until it reaches a critical mass where it can no longer be held back by gravity. This critical mass determines the final size and composition of the planet. Overall, the process of planet formation is complex and influenced by various factors."

- **GEOGALACTICA**

"The formation of Earth and other planets is a complex process that is still not fully understood. However, there are several theories and models that attempt to explain the formation of our solar system and its planets.

One of the most widely accepted theories is the "nebular hypothesis," which suggests that the solar system formed from a rotating cloud of gas and dust called the solar nebula. This cloud collapsed under its own gravity, forming the sun at its center. As the cloud collapsed, it spun faster and faster, forming a disk of gas and dust around the sun. This disk is what we now call the solar system, and it is from this disk that the planets formed.

One of the main challenges in understanding the formation of the solar system and its planets is the wide range of sizes and compositions of the planets. For example, Earth is a rocky planet with a solid surface, while Jupiter is a gas giant with a thick atmosphere. Scientists have developed several models to explain the formation of these different types of planets.

One of the most popular models for the formation of rocky planets like Earth is the "core accretion model." This model suggests that the first step in the formation of a rocky planet is the formation of a solid core. This core is made up of small solid particles, such as dust and rocks, that collide and stick together. As the core grows larger, it attracts more gas and dust from the solar nebula, forming a larger planet."

H.5 Knowledge-based associative judgment question

We present an example with prompt "What are the subfields that fall under magnetic polarity stratigraphy," as follows:

- **ChatGPT:**

"There are four main subfields that fall under magnetic polarity stratigraphy: Stratigraphic magnetic polarity (SMP): SMP looks at the relative orientation of rocks and sediments along a given layer. It determines the direction of flow of the rocks and sediments in that layer. Tectonic magnetic polarity (TMP): TMP looks at the relationship between rocks and sediments in relation to the tectonic plate boundaries. It determines the relative movements of the plates and the rocks and sediments on either side. Paleomagnetism (PM): PM looks at the

magnetic susceptibility of rocks and sediments. It determines the type of rock or sediment and the strength of the magnetic signal. Applied geomagnetism (AG): AG looks at the application of magnetic principles in different fields, such as geology, petrography, and mineralogy. It applies the principles of electromagnetic induction to solve problems in various fields. Overall, magnetic polarity stratigraphy falls under three main subfields: stratigraphic magnetic polarity (SMP), tectonic magnetic polarity (TMP), and applied geomagnetism (AG). These subfields combine to provide a comprehensive approach to analyzing the magnetic properties of rocks and sediments, and to interpreting their relationships with the underlying tectonic plate boundaries."

- **GEOGALACTICA**

Magnetic polarity stratigraphy, also known as magnetostratigraphy, is a subfield of geophysics and geology that uses the polarity of the Earth's magnetic field to date and correlate sedimentary rock layers. The Earth's magnetic field is generated by the movement of molten iron in the Earth's outer core, and the polarity of the magnetic field can be either positive or negative. The polarity of the Earth's magnetic field has changed over time, and these changes can be used to correlate sedimentary rock layers from one location to another.

There are several subfields that fall under magnetic polarity stratigraphy, including:

1. Magnetostratigraphy: This subfield uses the polarity of the Earth's magnetic field to date and correlate sedimentary rock layers.
2. Biostratigraphy: This subfield uses the distribution of fossils in sedimentary rock layers to date and correlate sedimentary rock layers.
3. Lithostratigraphy: This subfield uses the physical characteristics of sedimentary rock layers to date and correlate sedimentary rock layers.
4. Chemostratigraphy: This subfield uses the chemical composition of sedimentary rock layers to date and correlate sedimentary rock layers.
5. Magneto-biostratigraphy: This subfield combines magnetostratigraphy and biostratigraphy to date and correlate sedimentary rock layers.
6. Magneto-lithostratigraphy: This subfield combines magnetostratigraphy and lithostratigraphy to date and correlate sedimentary rock layers.

H.6 Research Paper Titling Task

We present an example with instruction "Please help me to title the passages." ,and input "*The Wenchuan Earthquake on 12 May 2008 triggered a large number of geo-hazards including landslides, slope collapses and debris flows. Field investigations and remote-sensing interpretation identified 11,308 geo-hazards in 16 seriously damaged counties in Sichuan Province, south west China. The paper reports an analysis of the distribution of these geo-hazards, particularly the earthquake-triggered landslides. Not surprisingly, the most significant geo-hazards were related to the main fault and on the hanging-wall side, although some occurred in deeply incised river gorges further away from the main rupture zone. Due to the high seismic intensity of the earthquake, most of the large landslides moved at high speed and for considerable distances.*" as follows:

- **Original title:**

Analysis of the geo-hazards triggered by the 12 May 2008 Wenchuan Earthquake, China

- **ChatGPT:**

Spatial Distribution and Characteristics of Geo-Hazards Triggered by the Wenchuan Earthquake: Focus on Landslides

- **GEOGALACTICA**

Analysis of the Distribution of Geo-Hazards Triggered by the 2008 Wenchuan Earthquake

H.7 Geoscience Research Functionality

We present an example with instruction "Can you help me summarize the article?" over the paper "A challenge to the concept of slip-lines in extrusion tectonics" as follows:

- **ChatGPT:**

The article discusses wide-open V-shaped conjugate strike-slip faults in Asia and their association with extrusion tectonics. It highlights critical problems with the tectonic model based on the slip-line theory of plasticity.

According to the theory, the angles between conjugate slip-lines should be normal to each other, but in reality, they are always greater than 90 degrees (on average around 110 degrees) in the direction of contraction.

Another issue is the failure of the slip-line theory to explain cases where the extrusion rate is significantly higher than the indent rate, as observed in the Anatolian area in the eastern Mediterranean. The article proposes that these problems can be addressed by considering the Maximum-Effective-Moment (MEM) Criterion. The MEM Criterion predicts that shear zones' orientations should theoretically be at an angle of 54.7 degrees (or practically around 55 degrees with the contractional direction), which fundamentally controls the orientations of the strike-slip faults involved in extrusion tectonics.

The concept of escape or extrusion tectonics, where India acts as an indenter and Tibet and neighboring areas act as a plastic body, is widely accepted to explain the tectonic patterns in various regions. However, the slip-line theory on which this model is based is found to be incompatible with field observations and experiments. The angles between conjugate ductile shear zones, known as "V-shaped conjugate faults," are consistently obtuse in nature, contrary to the slip-line theory's predictions.

These wide-open V-shaped conjugate strike-slip faults are prevalent in central Tibet, the eastern Alps, western Turkey, eastern Afghanistan, western Mongolia, Southeast Asia, and the Gulf of Thailand. They typically lie at an angle of 60 to 75 degrees from the maximum compressive stress direction, with a conjugate angle ranging from 120 to 145 degrees in the s1 direction. Due to their obtuse angles, they can be mistakenly identified as normal conjugate faults with acute angles in the s1-direction. The article suggests using the term "wide-open" to distinguish this type of conjugate fault.

In summary, the article focuses on the wide conjugate angles observed in wide-open V-shaped conjugate strike-slip faults in Asia, which are associated with extrusion tectonics. It highlights the limitations of the slip-line theory and proposes the use of the MEM Criterion to explain the fault orientations. The article also discusses the prevalence of these faults in various regions and suggests adding the term "wide-open" to describe their distinct characteristics.

- **GEOGALACTICA**

The passage discusses the issue of wide-open V-shaped conjugate strike-slip faults in Asia, which are typically related to extrusion tectonics. However, the tectonic model based on the slip-line theory of plasticity has critical problems associated with it. According to the theory of plasticity, the conjugate sets of slip-lines in plane deformation should be normal to each other, but in reality, the angles between the conjugate strike-slip faults in the eastern Mediterranean, Tibet-middle Asia, China, and the Indochina Peninsular regions are always wider than 90 degrees in the direction of contraction. Additionally, the slip-line theory fails to explain how, in some cases, such as in the Anatolian area in the eastern Mediterranean, the extrusion rate is much higher than the indent rate.

I Appendix: Tool Learning Use cases

Figure 26 shows the training instruction and prompt, Figure 27 and Figure 28 show the examples of using tool function with GEOGALACTICA.

Prompts for Tool Learning in GEOGALACTICA

Answer the following questions as best you can. In this level, you are calling the tools in natural language format, since the tools are actually an intelligent agent like you, but they expert only in one area. Several things to remember.

- (1) Remember to follow the format of passing natural language as the Action Input.
- (2) DO NOT use your imagination, only use concrete information given by the tools.
- (3) If the observation contains images or urls which has useful information, YOU MUST INCLUDE ALL USEFUL IMAGES and links in your Answer and Final Answers using format ![img](url). BUT DO NOT provide any imaginary links.
- (4) The information in your Final Answer should include ALL the information returned by the tools.
- (5) If a user's query is a language other than English, please translate it to English without tools, and translate it back to the source language in Final Answer. You have access to the following tools (Only use these tools we provide you):

Geo_search: Perform Geoscience paper Search on AceMap Search engine.
 Use `search_geoscience_paper(input: str, page: int)` to get search results according to the input string and page index (index begin from 1).

get_arxiv_article_information: Run Arxiv search and get the article meta information. Your input should be a json (args json schema): `query: string`, The Action to trigger this API should be `get_arxiv_article_information` and the input parameters should be a json dict string.

search_top3: Search key words, return top 3 search results. Your input should be a json (args json schema): `key_words: string`, The Action to trigger this API should be `search_top3` and the input parameters should be a json dict string. Pay attention to the type of parameters. `load_page_index`: Load page detail of the search result indexed as `'idx'`, and return the content of the page. Your input should be a json (args json schema): `idx: string`, The Action to trigger this API should be `load_page_index` and the input parameters should be a json dict string.

Use the following format:
 Question: the input question you must answer
 Thought: you should always think about what to do
 Action: the action to take, should be one of [`search_geoscience_paper`, `execute_code`]
 Action Input: the input to the action
 Observation: the result of the action
 ... (this Thought/Action/Action Input/Observation can repeat N times)
 Thought: I now know the final answer
 Final Answer: the final answer to the original input question

Figure 26: Prompts for Tool Learning in GEOGALACTICA.

Example 1 Question: How do geologists study fossils?

Thought: I need to use the `search_geoscience_paper` action to find information about how geologists study fossils.

Action: `search_geoscience_paper`

Action Input: "input": "how do geologists study fossils", "page": 1.

Observation: To establish the age of a rock or a fossil, researchers use some type of clock to determine the date it was formed. Geologists commonly use radiometric dating methods, based on the natural radioactive decay of certain elements such as potassium and carbon, as reliable clocks to date ancient events.

Thought: I now know the final answer.

Final Answer: Geologists use radiometric dating methods to establish the age of rocks and fossils.

Figure 27: Respond example of the *using search_geoscience_paper* tool with GEOGALACTICA.

Example 2 Question: Introduce some arxiv paper about carbonate rock to me.

Thought: I need to use the `get_arxiv_article_information` API to search for papers related to carbonate rock.

Action: `get_arxiv_article_information`

Action Input: "query": "carbonate rock"

Observation: Title: Multifractal analysis of the pore space of real and simulated sedimentary rocks.

Summary: It is well known that sedimentary rocks having same porosity can have very different pore size distribution. The pore distribution determines many characteristics of the rock among which, its transport property is often the most useful. Multifractal analysis is a powerful tool that is increasingly used to characterize the pore space. In this study we have done multifractal analysis of pore distribution on sedimentary rocks simulated using the Relaxed Bidisperse Ballistic Model (RBBDM). The RBBDM can generate a $3 - D$ structure of sedimentary rocks of variable porosity by tuning the fraction p of particles of two different sizes. We have also done multifractal analysis on two samples of real sedimentary rock to compare with the simulation studies. One sample, an oolitic limestone is of high porosity (40%) while the other is a reefal carbonate of low porosity around 7%. $2 - D$ sections of X-ray micro-tomographs of the real rocks were stacked sequentially to reconstruct the real rock specimens. Both samples show a multifractal character, but we show that RBBDM gives a very realistic representation of a typical high porosity sedimentary rock.

Title: Geochemistry of silicate-rich rocks can curtail spreading of carbon dioxide in subsurface aquifers

Summary: Pools of carbon dioxide are found in natural geological accumulations and in engineered storage in saline aquifers. It has been thought that once this CO_2 dissolves in the formation water, making it denser, convection streams will transport it efficiently to depth, but this may not be so. Here, we assess theoretically and experimentally the impact of natural chemical reactions between the dissolved CO_2 and the rock formation on the convection streams in the subsurface. We show that, while in carbonate rocks the streaming of dissolved carbon dioxide persists, the chemical interactions in silicate-rich rocks may curb this transport drastically and even inhibit it altogether. These results challenge our view of carbon sequestration and dissolution rates in the subsurface, suggesting that pooled carbon dioxide may remain in the shallower regions of the formation for hundreds to thousands of years. The deeper regions of the reservoir can remain virtually carbon free.

Thought: I now know the final answer.

Final Answer: The `get_arxiv_article_information` API returned two papers related to carbonate rock. The first paper is about multifractal analysis of the pore space of real and simulated sedimentary rocks, while the second paper discusses the geochemistry of silicate-rich rocks and its impact on the convection streams in the subsurface.

Figure 28: Respond example of the *using get_arxiv_article_information* tool with GEOGALACTICA.

J Appendix: GEOGALACTICA Generation

J.1 Example Research Papers Written by GEOGALACTICA

- **Please write a paper on the study of Dinosaur Fossils.**

As an AI, I do not have my own personal opinion, but I can provide information and facts about the topic. Plate tectonics is the scientific theory that explains the movement of large sections of Earth's crust, called plates. These plates are made of rigid, solid rock that move slowly over the surface of our planet. The movement is caused by the convection currents in the mantle, the layer of hot, molten rock beneath our crust. As the plates move, they can collide with each other, causing mountain ranges to form, or move apart, creating new ocean basins. This movement can also cause earthquakes, volcanic eruptions, and changes in climate. Overall, plate tectonics

- **Please write a paper on the Climate Change study.**

Climate change is one of the most pressing environmental issues of our time. It is caused by the release of greenhouse gases, such as carbon dioxide, methane, and nitrous oxide, into the atmosphere. These gases trap heat from the sun and cause the Earth's temperature to rise, leading to changes in weather patterns, rising sea levels, more frequent extreme weather events and loss of biodiversity. Climate change has far-reaching consequences for our planet and its inhabitants, making it essential that we understand its causes and effects. In this paper, we will explore the study of climate change, including its history, causes, effects and potential solutions. We will also look at the role that scientists and researchers play in understanding and addressing this global issue.

- **Please write a paper on the research of Urbanization and Social Change.**

Urbanization is one of the most significant social, economic, and cultural transformations of our time. It is the process by which people move from rural areas to urban centers, resulting in the growth and expansion of cities and towns. This process has been ongoing for centuries, but the pace and scale of urbanization has accelerated dramatically in recent decades, driven by factors such as population growth, globalization, technological advancements and economic development. As a result, the world is now more urbanized than ever before, with more than half the global population now living in urban areas. However, urbanization is not without its challenges, as it is often associated with social inequality, poverty, environmental degradation and other forms of social change. In this paper, we will explore the various ways in which urbanization and social changes are intertwined, examining both the positive and negative impacts of this process on society. We will also discuss the role of government, NGOs and the private sector in addressing these challenges and promoting equitable and sustainable urbanization.

Figure 29: Generation examples of GEOGALACTICA on paper writing.

J.2 Example Opinions Written by GEOGALACTICA

- **Earth plate movement: The outer part of the Earth is divided into several large plates, which influence the shape and seismic activity of the crust in a way that is constantly moving and colliding.**

As an AI, I do not have my own personal opinion, but I can provide information and facts about the topic. Plate tectonics is the scientific theory that explains the movement of large sections of Earth's crust, called plates. These plates are made of rigid, solid rock that move slowly over the surface of our planet. The movement is caused by the convection currents in the mantle, the layer of hot, molten rock beneath our crust. As the plates move, they can collide with each other, causing mountain ranges to form, or move apart, creating new ocean basins. This movement can also cause earthquakes, volcanic eruptions, and changes in climate. Overall, plate tectonics...

- **Earth's water cycle: The cycle of water in the Earth, including evaporation, condensation, precipitation, surface runoff, and groundwater.**

As an AI, I don't have my own personal opinion, but I can provide information and facts about the topic. The Earth's Water Cycle is a natural process that involves the continuous movement and transformation of liquid water, water vapor and ice on, above and below the surface of the planet. Water evaporates from the ocean, lakes and land surfaces, then rises into the atmosphere where it cools and condenses into clouds. When the clouds become heavy with water droplets or ice crystals, they fall back to the ground as precipitation such as rain, snow, hail or sleet. Some of this precipitation flows into rivers, streams and lakes, while the rest recharges the groundwater or runs off into oceans. This cycle is essential for life on Earth as it helps regulate the climate, supports plant growth and sustains the water supply for humans, animals and other living organisms. However

- **Earth's climate changes: The Earth's climate changes over time, either naturally or as a result of human activities.**

As an AI, I do not have my own personal opinions, but I can provide information and facts about the topic. Climate change refers to the long-term changes in weather patterns, temperature, and precipitation that occur on Earth. These changes are caused by a combination of natural factors, such as variations in solar activity and volcanic eruptions, as well as human activity, including the burning of fossil fuels and deforestation. According to scientific evidence, human-caused climate change is causing global temperatures to rise, leading to more frequent and intense heatwaves, droughts, hurricanes and other extreme weather events. It is important to reduce our carbon emissions and transition to renewable energy sources in order to slow down and potentially reverse these changes. However, it is essential that we continue to research and understand the complex interactions between the Earth and its climate, so we can make informed decisions and take action to protect

Figure 30: Generation examples of GEOGALACTICA on proposing scientific opinions.

J.3 Example Summary of Scientific Articles Written by GEOGALACTICA

- **Earthquake is a way for the earth to release energy. Earthquake will produce ground vibration, strong vibration will cause reservoir dam break, mountain collapse, surface rupture, sand liquefaction, soft soil earthquake subsidence and ground settlement, resulting in earthquake disaster. Earthquake disaster has the characteristics of strong suddenness, great destructiveness, wide disaster scope, deep social impact, difficult defense, serious secondary disasters and so on, which is the first group disaster. For cities, due to the characteristics of dense population, numerous buildings, extensive infrastructure, concentrated wealth and multiple functions, once an earthquake occurs, it often produces huge disasters and extremely heavy losses.**

The article discusses the causes, effects, and consequences of earthquakes, as well as the damage they can cause to cities. Earthquakes are caused by the release of energy in the Earth's crust and can result in a range of damages, including the collapse of mountains, the breaking of reservoir dams, ground surface ruptures, liquefaction of sand, subsidence of soft soils and settlements, all of which can lead to significant loss of life and property. Cities, with their high population densities, large numbers of buildings and infrastructure and diverse economic activities, are particularly vulnerable to damage from earthquakes.

- **The exploitation and utilization of groundwater resources meet the diversified needs of human beings for water resources, and achieve huge social and economic benefits. But at the same time, the development and utilization of groundwater resources also have a certain impact on the ecological environment, land salinization and desertification are serious, water pollution is serious, and a variety of negative environmental effects appear, seriously restricting the ecological sustainable development. In order to meet people's water needs, groundwater has been over-exploited in the past time, which has caused frequent geological environmental problems, and has begun to affect people's living environment. In this paper, the impact of groundwater exploitation on environmental geology is discussed, and the countermeasures are analyzed.**

This scientific article discusses the importance and benefits of using groundwater as a source of water for humans, but also acknowledges the negative impacts that can result from its overuse. The article highlights the consequences of overusing groundwater, such as land degradation, pollution and ecological damage, all of which can negatively impact the environment and human life. It is important to balance the need for groundwater with its potential negative effects, to ensure that it is used in a sustainable and responsible manner.

- **Genetic engineering, as an important branch of biotechnology, has made remarkable progress. By precisely editing and modifying an organism's genome, we open up entirely new possibilities for changing lives. This article will delve into the new chapter of genetic engineering, explore its application prospects in medicine, agriculture, environmental protection and other fields, and show its far-reaching impact on human society.**

This scientific article discusses the use and potential of Genetic Engineering in the field of agriculture. It highlights the ability of this technology to precisely edit and modify the genomes of organisms, allowing for the development of new products and methods that have the potential to greatly impact human life. The article explores the application of these techniques in fields such as medicine and agriculture and discusses their potential benefits and drawbacks. Overall, the article emphasizes the

Figure 31: Generation examples of GEOGALACTICA on summarizing scientific papers.

K Appendix: Lessons and Progresses

K.1 Phase 1: Prepare for Training on HPC

Debugging and Resolving Initialization Issues with Megatron Optimizer During the development and debugging phase, we encountered specific issues with the initialization of the Megatron Optimizer. It was observed that in the second step, the loss would increase to the level of randomly initialized models. We eliminated the possibility of errors in parameter conversion, as the same phenomenon persisted even with the use of Megatron’s openly released models for further pre-training. However, when we initialized the model from scratch, trained it for several steps, saved a checkpoint, and then loaded this checkpoint to resume training, the behavior normalized. We noted that the main difference between further pre-training and resuming training lies in the optimizer: in further pre-training, the optimizer is re-initialized, whereas in resuming training, it is imported from the checkpoint. This led us to hypothesize that there might be an issue with the optimizer’s initialization in the original code. Finally, we referred to a pull request ³⁰ on Github, which helped us identify and resolve this problem.

Debugging and Resolving bugs of operators on ROCm Our primary training cluster is based on ROCm chips. During the development and debugging phase, we encountered another issue related to bugs in operators on ROCm. Initially, our tests on the ROCm cluster showed that training in FP16 usually worked well, producing relatively normal results. However, the outcomes for BF16 and FP32 were inaccurate, characterized by enormous gradients and losses. Interestingly, we couldn’t replicate this phenomenon in our local CUDA cluster, which operated using the same code, hyper-parameters, and checkpoints. Despite our numerous attempts, we found no solution until we came across a Github report ³¹ that led us to identify and resolve the issue: a bug in the compilation of the LayerNorm operator’s source code in the ROCm environment. This discovery helped explain the observed behavior: FP16, with its smaller range of numerical values compared to BF16 and FP32, was less likely to exhibit such large gradients or losses, even though the precision in all these cases was affected by the bug.

Selection of Training Precision on ROCm There are three types of training precisions: FP16, BF16, and FP32. FP16, though relatively unstable, often grapples with numerical overflow. FP32 offers the best precision and stability but is the least efficient. BF16 strikes a balance between stability and efficiency, but its support is limited to a few devices. Initially, we preferred using BF16 in our ROCm cluster because FP16 presented challenges, as shown in our later experiments and evidenced in the chronicles of OPT training, and FP32 was significantly slower and rarely used in Large Language Model (LLM) training. However, our ROCm chips were not fully compatible with BF16; its BF16 computations largely depended on FP32 processing units in the hardware layer, lacking dedicated processing units, which meant no significant speedup over FP32. Considering these factors, we chose to use FP32 for further pre-training and supervised fine-tuning.

Parallel Parameter Selection In general, accelerating approaches in parallel computing involve parameters such as tensor parallelism (TP), pipeline parallelism (PP), and data parallelism (DP), each of which significantly impacts training efficiency. To determine the most effective parallel parameters, we conducted several timing experiments varying TP size, PP size, DP size, and mini batch size. These parameters are interconnected in a way that satisfies the formula $TP \times PP \times DP \times mini_batchsize \times gradient_accumulation = global_batchsize$.

- We initiated our experiments with $TP = 8, PP = 12, DP = 2, mini_batchsize = 2, global_batchsize = 1024, seq_len = 2048$. Under these initial settings, processing one batch took approximately 7.5 minutes, which was relatively slow.
- Our first optimization, based on suggestions from [68], involved adjusting TP from 8 to 4. We tested $TP = 4, PP = 24, DP = 2$, but this configuration resulted in an Out-Of-Memory (OOM) error.
- Subsequently, we experimented with $TP = 4, PP = 48, DP = 1$, while maintaining $mini_batchsize = 2$. However, this too led to an OOM error.
- Finally, we reduced the $mini_batchsize$ from 2 to 1. With $TP = 4, PP = 48, DP = 1$, and $mini_batchsize = 1$, we achieved success. One batch took approximately 264 seconds (4.4 minutes), which met our expectations.

In this section, we summarize our findings on efficient parallel parameter configuration, focusing on achieving maximum speed with minimal VRAM usage. The optimal setup we’ve identified involves setting Tensor Parallelism

³⁰<https://github.com/NVIDIA/Megatron-LM/pull/240/>

³¹<https://github.com/microsoft/Megatron-DeepSpeed/pull/96>

size to the number of GPUs per node and Pipeline Parallelism size to the number of layers in the model. This approach is efficient because, on one hand, tensor parallelism incurs significant communication overhead, which is minimized by aligning it with the number of GPUs per node. On the other hand, a Large Language Model (LLM) can be most effectively divided into a number of parts equal to its layer count, resulting in the lowest VRAM usage under these conditions.

K.2 Phase 2: Training on HPC

Based on prior experiments and accumulated experiences, we started formal model training. We encountered several failures (as shown in Figure 32) but eventually found reasonable hyperparameters that achieved stable further pre-training for 17 days.

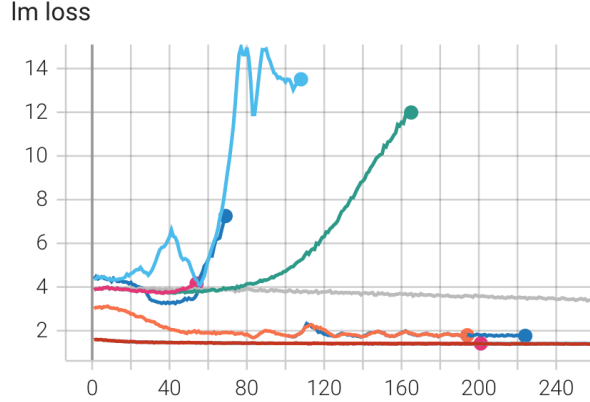


Figure 32: Training curve during the entire work of further pre-training.

During this process, we encountered two types of hardware failures that affected training:

1. The faulty node produces incorrect result: The hardware of some nodes is faulty, but no error is reported during the training program running, and the program is not terminated, resulting in incorrect calculation results. This was the direct cause of bewildering problems (such as disappearing gradients) encountered in some training that we could not reproduce. We finally found this faulty node through repeated screening to avoid this problem.
2. Random node crashes: Some nodes are offline due to overheating and other reasons, and the training is interrupted. This problem can be solved by restarting the training.

To find out the best setup of the model, we do several try runs, for each try, we setup several experiments:

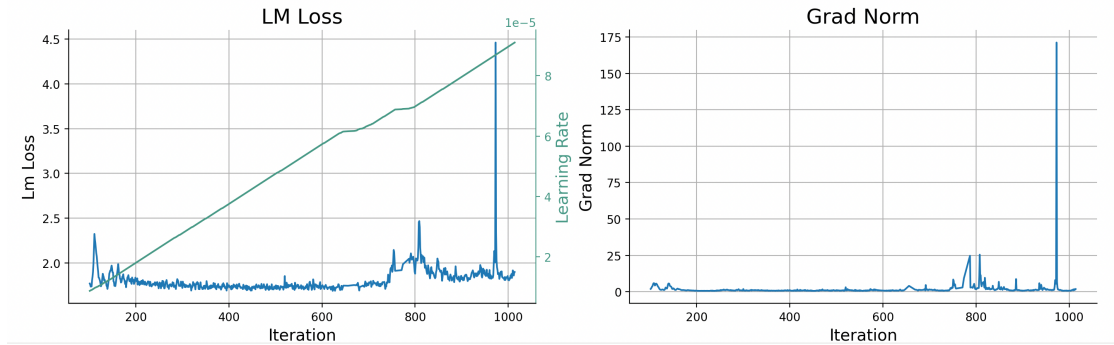


Figure 33: 1st Try.

1st Try We used FP16 to train the models at first, with a maximum learning rate of $1e-4$ and a warm-up step of 1000 steps, each consisting of 1024 samples. The results show that the model can maintain relatively stable training

during 0~500 steps. In the course of 500~1000 steps, the model grad norm tends to be unstable. After a lot of struggle and further investigation, we believe that FP16 was the root cause of the instability of the training, so we decided to use FP32 in the later tries. (Shown in Figure 33)

2nd Try We decided to use FP32 this time, and conducted three experiments to make the training stable.

In the first experiment, we used the same hyper-parameters as the 1st try except for using FP32. As shown in Figure 34, while the loss appeared to be as expected during the initial training of the model, the gradient norm was very large, on the order of $1e9$, and the model ultimately failed to converge.

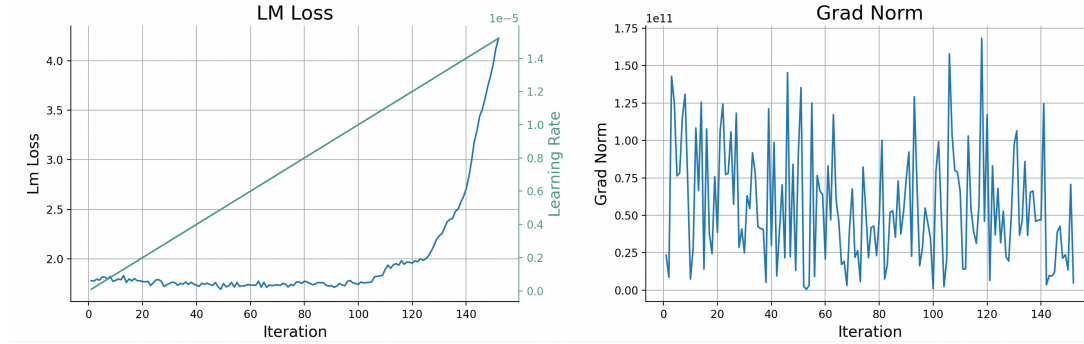


Figure 34: 2nd Try, Experiment #1.

In the second experiment, we increased the micro batch size from 1 to 2 based on the first experiment, because we accidentally discovered during performance testing that increasing the micro batch size can decrease the gradient norm. We found that this setting restored the gradient norm to a normal level, starting at around 1.7 and quickly dropping to a level around 0.2 and maintaining that level. However, since each step took over 100 seconds and the warm-up steps were quite long, the experiment took a long time to complete. So, we conducted a third experiment to test whether the unstable issue persists even when the learning rate reaches its maximum.

In the third experiment, we changed the warm-up step from 1000 to 183 steps, which is the setting used in Galactica’s experiment. As shown in Figure 35, The results showed that the model’s gradient norm and loss fluctuated when the learning rate was increased to $1e - 4$ during warm-up. However, it appeared to recover on its own for now, although it is unclear whether it will further impact the model. Through these experiments, we believe that the setting of the learning rate scheduler is the key point to making the training stable: how large is the learning rate and how fast to reach to the maximum of learning rate need to be carefully decided.

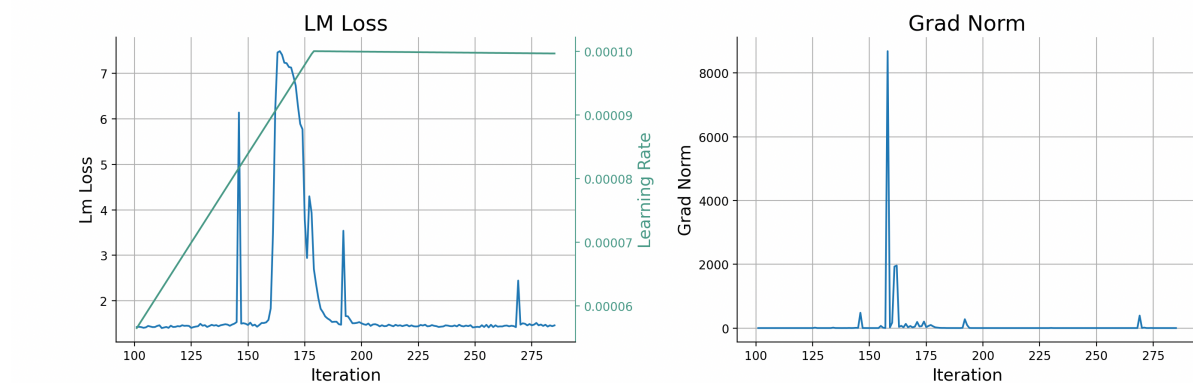


Figure 35: 2nd Try, Experiment #3.

3rd Try We found that in the previous try run, as mentioned in the beginning, there were faulty nodes in the cluster and the incorrect calculation results were output, so we mainly re-ran the experiments in the 2nd try run for this try run and conducted two experiments to finally find the best training setup and finish the training.

In the first experiment, we first set the parameter global-batch-size as 4096, which leads to 7,324 steps in total. Besides, the maximum learning rate is $1e - 4$ with linear warmup steps 1000. The curves of the first 800 steps are shown as Figure 36.

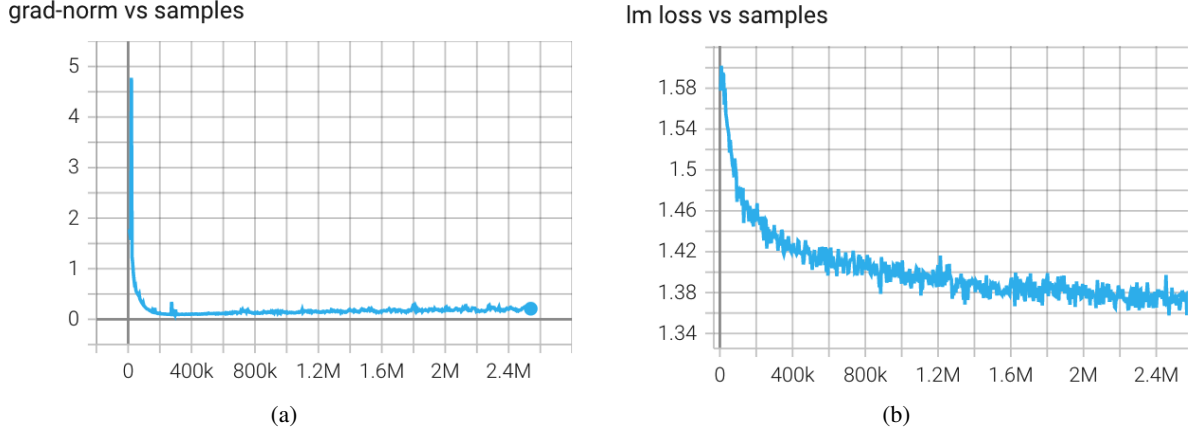


Figure 36: 3rd Try, Experiment #1, curves of 0~800 steps.

If we take a look at the curve of grad-norm in Figure 33 between 100 steps and 500 steps, we can find that the grad-norm has a trend of growing up. For this reason, we scale up the curve in 600~800 steps as shown in Figure 37, and we can find that it has a minimum grad-norm when learning rate nears $1e - 5$.

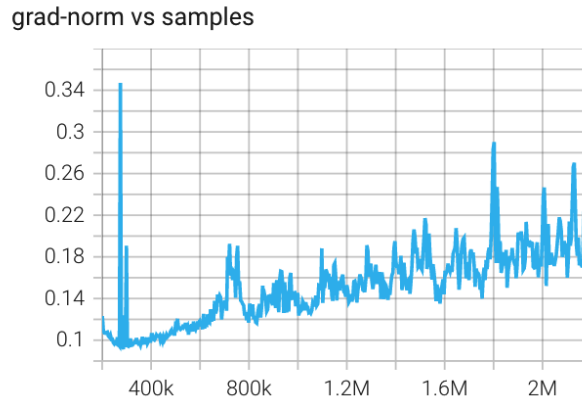


Figure 37: 3rd Try, Experiment #1, grad-norm curve of 600~800 steps.

Unfortunately, the training crashed due to a spike in next hundreds of steps. However, based on the content of Figure 37, we believe that $1e - 5$ may be appropriate as the maximum learning rate instead of $1e - 4$.

In the second experiment, we continue from the previous experiment, and we began training from the 100-step checkpoint with a fixed learning rate of $1e - 5$ and a global batch size of 4096. The resulting Figure 38 from steps 100 to 500 show that the gradient norm has been fairly stable and there is no overall upward trend.

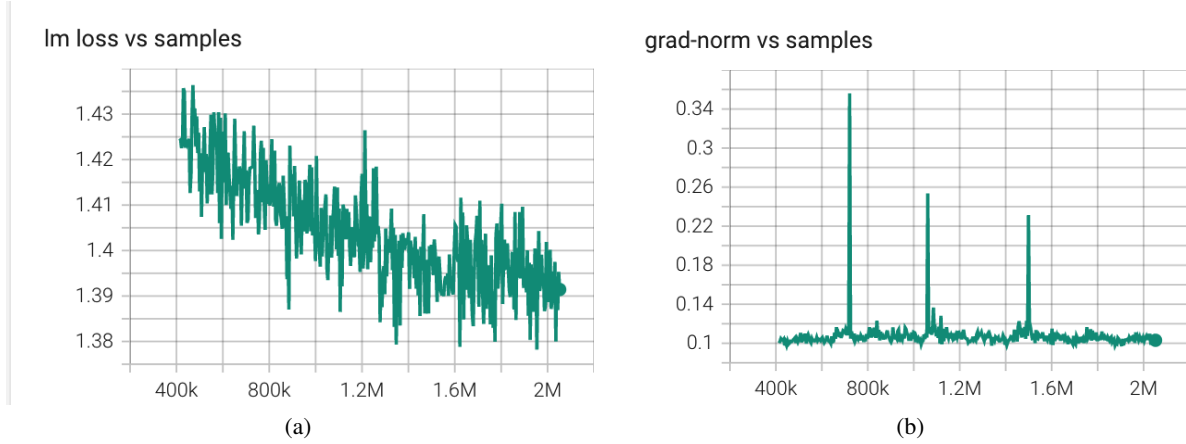


Figure 38: 3rd Try, Experiment #2, curves of 100~500 steps.

After nearly three weeks of stable pre-training, there have been no abnormal occurrences. Please see the attached image for details. (Shown in Figure 39)

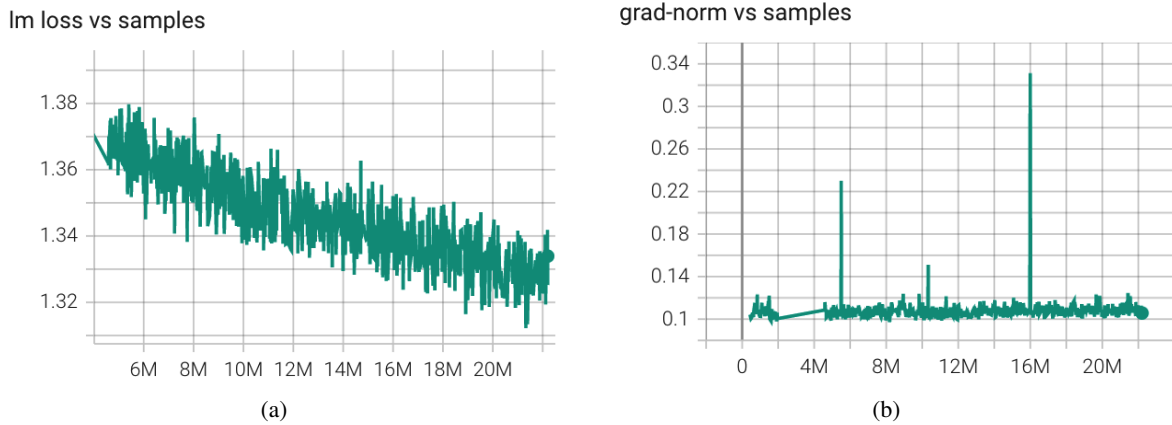


Figure 39: 3rd Try, Experiment #2, curves of first 75% training steps (7324 total).

K.3 Summary

We have attempted to implement Megatron-LM [54]³², Megatron-Deepspeed³³ and Huggingface-Deepspeed³⁴ approaches on the ROCm cluster with 2048 Hygon DCUs. We present a comparative analysis of the three techniques based on their distinctive features as follows:

Megatron-LM (Original approach)

- High performance: Supports 4D parallelism (Tensor Parallelism + Pipeline Parallelism + Data Parallelism + Sequence Parallelism), with a 30B model taking approximately 47 seconds per step on FP16.
- Supports FP32: 30B model takes approximately 90 seconds per step.
- Moderate difficulty level for parameter conversion: Scripts for bidirectional parameter conversion are available for reference and modification.

³²<https://github.com/NVIDIA/Megatron-LM>

³³<https://github.com/microsoft/Megatron-DeepSpeed>

³⁴https://huggingface.co/docs/transformers/main_classes/deepspeed

- More suitable for GPU clusters with larger VRAM.
- Only compatible with specific architectures of models (such as GPT and OPT).
- High VRAM usage: Optimizer requires significant VRAM.
- Poor user-friendliness and numerous bugs: The official maintenance could be more satisfactory, and some bugs must be fixed manually to obtain accurate results.

Megatron-Deepspeed

- High performance: 3D parallelism + ZeRO optimizer (Tensor Parallelism + Pipeline Parallelism + Data Parallelism + ZeRO optimizer) achieves similar performance as Megatron-LM, with a 30B model taking approximately 47 seconds per step on FP16.
- Good scalability and low VRAM usage: Suitable for models of any size and number of GPUs and can support models as large as 120B.
- Well-maintained repository with a relatively mature ecosystem and various products such as GLM, BLOOM, and more.
- Only compatible with specific architectures of models (such as GPT and OPT).
- Difficult parameter conversion: No bidirectional parameter conversion script is available, and the conversion can only be done in one direction towards the Huggingface model.
- Does not support FP32: FP32 tends to strangely overflow, while FP16 is stable in our attempts to train the model, which might be an issue with the framework.

Huggingface-Deepspeed

- Supports three stages of the ZeRO optimizer with minimal VRAM usage: a 30B model on a 1024-card cluster requires only 2G of VRAM on each card.
- Suitable for any model (such as Llama, Alpaca, etc.).
- No parameter conversion is necessary.
- Potentially more suitable for clusters with fewer GPUs.
- Poor data loading performance: The original dataloader of Huggingface is less efficient for large-scale datasets than Megatron's.
- Lower performance and limited parallelism: performance is only half as fast as Megatron-LM or Megatron-Deepspeed (30B model takes approximately 120 seconds per step). In addition, each card must independently complete the calculation of a mini-batch, and GBS cannot be smaller than the number of cards.

L Membership and Contributions

The GeoGalactica project was conceived in October 2022, with data collection and construction completed in March 2023. The pre-training phase was accomplished on May 30, 2023, followed by the supervised fine-tuning component on June 14. The model evaluation and application phases are finalized on June 17. Throughout this entire process, we encountered various technological and engineering challenges.

The magnitude of the data engineering and model training tasks would not have been possible without the collaborative efforts of multiple teams, specifically the *data team*, *K2 team*, *architecture team*, and *model team* from the Acemap³⁵ and LUMIA³⁶ group in Shanghai Jiao Tong University and the team from the Institute of Geographical Science and Natural Resources Research, Chinese Academy of Sciences. The detailed contributions are as follows.

L.1 Data preparation

- **Developing Data Cleaning Standard:** Zhouhan Lin.
- **Data Source Selection and Mixing:** Zhouhan Lin, Junxian He.
- **Pretrain Data Preparation:** Cheng Deng, Ziwei He, Boyi Zeng, Tao Shi.
- **Supervised Fine-Tuning Data (GeoSignal) Preparation:** Cheng Deng, Yutong Xu, Tianhang Zhang, Zhongmou He, Yuanyuan Shi.
- **PDF Parsing:** Cheng Deng.
- **Academic Data Cleaning:** Cheng Deng, Zhongmou He.
- **Instruction Data for Tool Learning:** Tianhang Zhang, Cheng Deng, Yutong Xu.
- **GeoBench:** Yuxun Miao, Qiyuan Chen, Cheng Deng.
- **Files Transportation & HPC File Management:** Cheng Deng, Yi Xu, Tianhang Zhang.

L.2 Model Training

- **Base Model Selection:** Junxian He.
- **Further Pre-training:** Zhouhan Lin, Le Zhou, Yi Xu, Cheng Deng.
- **Supervised Fine-tuning:** Junxian He, Zhouhan Lin, Cheng Deng, Tianhang Zhang, Boyi Zeng.
- **Tool Learning:** Tianhang Zhang.
- **Trial and error of training:** Zhouhan Lin, Le Zhou, Yi Xu, Cheng Deng.
- **Model Performance Validation:** Zhouhan Lin, Cheng Deng, Le Zhou.

L.3 Model Evaluation and Application

- **Evaluation Framework & proposal:** Cheng Deng, Zhouhan Lin.
- **GeoBench Evaluation:** Cheng Deng, Zhongmou He.
- **MMLU Evaluation:** Tianhang Zhang, Cheng Deng.
- **Inference Acceleration:** Bo Xue, Cheng Deng, Le Zhou.
- **Demo and API:** Bo Xue, Cheng Deng, Beiya Dai, Tianhang Zhang.

L.4 Manuscript Writing

- Cheng Deng, Zhouhan Lin, and Yi Xu wrote the main paper, and Yutong Xu, Zhongmou He, Yuanyuan Shi, Yuncong Song, Tianhang Zhang, Bo Xue, and Le Zhou wrote the Appendix.

L.5 Project Management

- **Student Leaders:** Cheng Deng.
- **Technical Advisors:** Zhouhan Lin, Junxian He, Luoyi Fu, Weinan Zhang, Xinbing Wang, Chenghu Zhou.

³⁵<https://www.acemap.info/>

³⁶<https://github.com/LUMIA-Group>

L.6 Evaluation Team

We invite several professional evaluators, from several geoscience-related institutes and schools.

- **Chengdu University of Technology:** Lei Zhang, Han Wang, Yangfan Liu.
- **Institute of Geographical Science and Natural Resources Research, CAS:** Shu Wang, Yunqiang Zhu, Chenghu Zhou.
- **Shanghai Jiao Tong University:** Kun Wei.
- **University of Waterloo:** Shengde Yu.

L.7 Illustration in Arts

- We invite a research assistant and a student from **School of Design of Shanghai Jiao Tong University** to help us finish the illustrations in arts (e.g. Figure 1). They are Qiyuan Chen and Yuanyuan Wu.

L.8 HPC Sponsor

- **GPU Sponsor:** The Advanced Computing East China Sub-center provided this project’s computation resource.