
Geo-Expert: Towards Expert-Level Geological Reasoning via Parameter-Efficient Fine-Tuning

Anonymous Authors¹

Abstract

While general-purpose Large Language Models (LLMs) applied to Geology often hallucinate when reasoning about subsurface structures and deep-time evolution, current AI in Earth sciences predominantly targets surface remote sensing and GIS. To bridge this gap, we introduce Geo-Expert, a family of parameter-efficient geological LLMs fine-tuned on a custom-curated, high-quality instruction dataset processed using our custom instruction synthesis pipeline. We investigate the impact of model scaling and architecture by fine-tuning three base models: Qwen3-8B, Qwen3-32B, and Gemma-3-27B, with Low-Rank Adaptation (LoRA) method. Our extensive evaluation on a novel domain-specific benchmark, Geo-Eval, reveals that a domain-aligned 8B model can outperform open-weight 70B generalists and proprietary GPT-4o on specialized geological reasoning, while a 32B variant approaches frontier reasoning models. The optimized 8B model further offers a competitive cost-performance ratio for deployment. This work provides a reproducible recipe for democratizing scientific LLMs and establishes a baseline for geological artificial intelligence.

1. Introduction

Geology plays a vital role in underground resource exploration, geologic hazard mitigation, and understanding planetary evolution. Unlike general knowledge domains, geological reasoning requires interpreting complex spatio-temporal relationships and a large amount of professional data. However, LLMs often struggle with the specialized logic required in geology. General-purpose models frequently exhibit severe hallucinations when asked to perform a rigorous deduction on stratigraphy or structural geology,

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

limiting their utility in geological research and industrial applications.

Current large models in the vertical field of geography predominantly focus on Geography Knowledge Agent and Earth Observation, demonstrating proficiency in general geospatial intelligence and utilizing computer vision to analyze imagery for surface monitoring. They effectively handle tasks related to surface geography, environmental monitoring, and human-geography interactions. However, a critical gap remains between general *Geography* and specialized *Geology*. While existing models can retrieve facts about Earth’s surface, they often lack the deep-domain reasoning required to understand the solid Earth such as interpreting subsurface stratigraphic sequences, reconstructing deep-time tectonic evolution, or analyzing complex petrogeologic processes.

To bridge this gap, we introduce Geo-Expert, a suite of parameter-efficient geological LLMs. We construct a high-quality, domain-specific instruction dataset derived from five canonical geology textbooks. Utilizing the custom framework, we synthesize structured instruction-response pairs to instill rigorous domain knowledge. We investigate the scaling laws of geological reasoning by fine-tuning three base models: Qwen3-8B, Qwen3-32B, and Gemma-3-27B, employing Low-Rank Adaptation (LoRA) to improve training efficiency.

In sum, our main contributions are as follows:

- **Dataset Construction:** We present a reproducible pipeline for converting unstructured geological textbooks into high-quality instruction-tuning datasets.
- **Model Scaling Analysis:** We provide a comparative analysis of model architectures (Qwen vs. Gemma) and parameters scales (8B vs. 27B/32B), revealing that while smaller models are efficient and easy to deploy, larger parameters are crucial for complex structural reasoning.
- **Geo-Expert Benchmark:** We establish Geo-Eval, a benchmark for geological reasoning, demonstrating our fine-tuned models significantly outperform general-purpose baselines in specialized tasks.

2. Model Framework and Benchmark Construction

2.1. Data Sources

Geological reasoning is fundamentally grounded in established physical principles and deep-time evolutionary models, necessitating data sources of the highest scientific rigor. To construct a high-fidelity instruction dataset, we identified five canonical textbooks as the data source. These sources were strategically selected to cover the full spectrum of solid earth science, balancing theoretical depth, cognitive difficulty, and practical application:

- *Principles of Sequence Stratigraphy I & II* (Catuneanu) (Catuneanu, 2006) — This text serves as the standard for understanding temporal frameworks in sedimentation and basin analysis. It is renowned for its highly abstract conceptual models, challenging the model’s deep-time evolutionary reasoning.
- *Structural Geology & Advanced Structural Geology* (Fossen, 2016) — These volumes provide the theoretical foundation for rock mechanics, stress analysis, and rheology. They feature a steep learning curve with complex mathematical derivations, establishing a rigorous physical baseline.
- *Tectonics and Sedimentation: Implications for Petroleum Systems* (Gao, 2012) — A critical resource bridging the gap between plate tectonics, basin formation, and hydrocarbon exploration. It provides complex, interdisciplinary macro-scale scenarios.
- *Structural Analysis and Synthesis* (Rowland et al., 2020) — It focuses on practical application, including 3D geometric interpretation and structural mapping techniques. This volume strictly evaluates the model’s spatial imagination and geometric deduction capabilities.

2.2. Dataset Construction

Transforming static textbook content into a dynamic instruction-tuning dataset involves a pipeline designed to maximize data fidelity and reduce noise. We implemented a multi-stage process: digitization, heuristic cleaning, and instruction synthesis.

Text Extraction and Sanitization. We first used MinerU (OpenDataLab, 2024), a high-precision PDF-to-Markdown tool, to extract textual content while preserving the structure of the original textbooks. Following the extraction, we developed a custom Python-based cleaning module to process the raw Markdown. This pipeline first segments the text into discrete blocks based on paragraph lines containing

metadata keywords like page numbers and copy notices. Furthermore, to prevent redundancy caused by layout overlaps, we implemented a global deduplication mechanism using a hash-based method to filter out identical text blocks across the entire corpus.

Instruction Synthesis. To transform the cleaned Markdown segments into a high-quality, geoscientific instruction-tuning corpus, we engineered a custom, domain-specific instruction synthesis pipeline. This pipeline operates on a project-based workflow that deeply integrates earth science document processing, geological domain modeling, and reasoning-enhanced generation. The synthesis process consists of three key stages:

1. **Chapter-Aware Recursive Chunking:** Unlike simple fixed-length splitting, the framework utilizes a semantic segmentation algorithm based on the Markdown header structure. This ensures that each text chunk retains complete semantic context within its hierarchical scope, preventing the fragmentation of complex geological descriptions.
2. **Domain-Structured Question Generation:** An LLM first constructs a hierarchical domain tree (extracting outlines and concepts) to bind specific tags to text segments. Questions are then dynamically generated based on these domain tags and character density, minimizing redundancy and ensuring broad coverage of geological concepts.
3. **Chain-of-Thought (CoT) Answer Construction:** To capture the rigorous logic required for geological reasoning, the system employs reasoning-oriented models to generate answers. This process produces not only the final ground truth but also the intermediate reasoning steps (CoT), ensuring strict alignment with the source material while explicitly detailing the deduction process.

This automated pipeline enabled us to efficiently scale the dataset to 11,518 high-quality instruction pairs while preserving the authoritative accuracy of the original canonical sources.

2.3. Model Fine-Tuning

We performed instruction tuning using the ms-swift framework (ModelScope, 2024), applying distinct hyperparameter strategies tailored to the scale of each model. We targeted all linear layers (all-linear) for Low-Rank Adaptation (LoRA) (Hu et al., 2021).

For the compact Qwen3-8B model (Yang et al., 2024), we adopted a standard efficiency configuration with LoRA rank

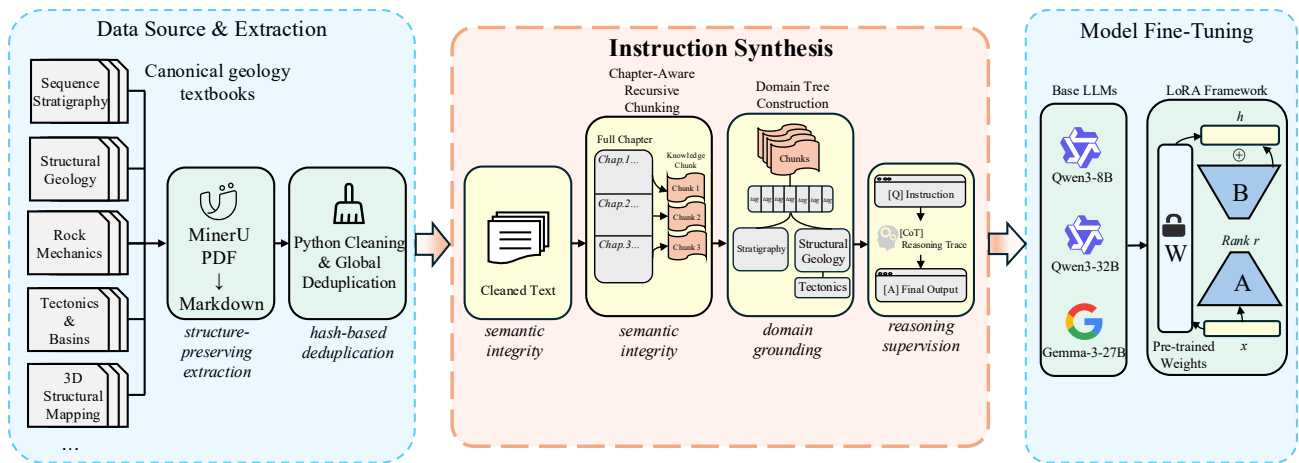


Figure 1. **Overview of the Geo-Expert framework.** The pipeline consists of three main stages: (1) extracting and sanitizing text from canonical geology textbooks; (2) synthesizing high-quality, CoT-enhanced instruction pairs via a domain-structured generation pipeline; (3) applying parameter-efficient fine-tuning on base LLMs.

$r=32$, $a=32$, and a learning rate of $2e-5$ using FP16 precision. This lightweight setup allowed for rapid iteration on a single NVIDIA RTX 5090.

In contrast, to effectively adapt the deeper architectures of Gemma-3-27B (Team, 2024) and Qwen3-32B, we scaled up the adaptation capacity, increasing the LoRA rank to $r=64$ and alpha to $a=128$. To mitigate the high memory footprint and ensure numerical stability during training, we employed BF16 (bfloat16) precision and activated gradient checkpointing with a gradient accumulation step of 4. This optimized configuration enabled the fine-tuning of these 27B+ parameter models on a $4 \times$ RTX 5090 setup, demonstrating a feasible recipe for fine-tuning 27B+ models on consumer-grade GPUs.

2.4. Benchmark Construction: Geo-Eval

Standard benchmarks in Earth sciences predominantly focus on remote sensing tasks or general geography trivia, failing to capture the rigorous deductive logic required for solid earth geology. To bridge this evaluation gap and comprehensively evaluate the true reasoning capabilities of LLMs in the vertical domain of geology, we constructed Geo-Eval, a novel domain-specific benchmark consisting of high-fidelity, open-ended reasoning questions. Unlike traditional static datasets manually extracted from exams, Geo-Eval was constructed through an automated adversarial mining pipeline coupled with expert human-in-the-loop verification. The construction process proceeded in four distinct phases:

1. **Large-Scale Generative Extraction:** We initially deployed DeepSeek-R1 (DeepSeek-AI et al., 2025) to perform a comprehensive semantic analysis of the source

geology textbooks. Through this process, we extracted a foundational pool of 2,591 complex geological questions and synthesized their corresponding preliminary reference answers.

2. **Comparative Inference:** To evaluate the baseline domain capability, both our fine-tuned Qwen3-8B-Geo and DeepSeek-R1 were tasked with generating independent responses to the entire pool of 2,591 questions.
3. **Difficulty-Aware Evaluation Paradigm (LLM-as-a-Judge):** To move beyond static, undifferentiated testing, we introduce a *difficulty-aware evaluation paradigm* designed to identify the true cognitive frontier of geological reasoning. We utilized GLM-4.5 (Zeng et al., 2024) as an impartial evaluator to blindly score both models' responses on a 10-point scale. Rather than evaluating on a random distribution, we systematically isolated a subset of 387 "hard boundary" questions where the score differential between our Qwen3-8B-geo and the frontier generalist DeepSeek-R1 was marginal (score difference ≤ 4). This dynamic calibration of difficulty yields exceptionally valuable boundary samples. They represent complex cognitive tasks where general-purpose giants exhaust their zero-shot deductive capacities and lose their decisive advantage. Consequently, this paradigm shifts the focus from basic factual recall to a much stricter, boundary-aware stress test of true expert-level geological reasoning.
4. **Expert Vetting and Ground Truth Finalization:** To ensure absolute scientific validity, the isolated 387 candidate questions were submitted to geology professors for rigorous manual screening, refinement, and calibra-

tion. Then we get the test dataset with the reference answers.

The resulting finalized dataset serves as a highly targeted, expert-verified ground truth. By deliberately focusing on these high-complexity boundary problems, Geo-Eval provides a robust metric for assessing whether a model has truly internalized deep geological logic, moving far beyond superficial pattern matching.

2.5. Mitigating Selection Bias

We acknowledge that the boundary mining process uses Qwen3-8B-Geo as one comparator, which could in principle introduce selection bias. We mitigate this in three ways: (i) the 387 candidate questions were independently re-vetted, refined, and—where necessary—rewritten by geology professors, and the reference answers used for scoring are expert-authored ground truth, not model outputs; (ii) the final evaluation employs GPT-4o as a reference-guided judge against this expert-verified ground truth, decoupling scoring from the mining process; (iii) two of our three reported models (Qwen3-32B-Geo and Gemma-3-27B-Geo) were not used in the boundary-mining loop, yet exhibit consistent gains, suggesting the observed improvements generalize beyond the mining comparator.

3. Benchmark and Evaluation: Geo-Eval

3.1. Benchmark Content

To avoid one-dimensional evaluation, we rigorously partition the Geo-Eval test dataset into three cognitive levels:

- **Concept:** Focuses on the identification and definition of core terminology, such as mineralogy, petrology, and stratigraphy.
- **Process:** Focuses on the logical deduction of dynamic geological mechanisms, including tectonic movements and diagenesis.
- **Engineering:** Focuses on the comprehensive application analysis of geological conditions in practical scenarios, such as engineering site selection and geological hazard prevention.

3.2. Benchmark Baselines

After applying Geo-Eval to different models, we adopt the widely recognized LLM-as-a-Judge paradigm, utilizing GPT-4o as a reference-guided LLM evaluator comparing against expert-verified ground truth. We impose strict domain-fact-checking constraints within the prompt to generate a quantitative score ranging from 0 to 10. We select 11

highly representative models for a comprehensive comparison, categorized into three groups: proprietary state-of-the-art models (e.g., GPT-5.4, DeepSeek-V3.2), open-weight generalist models of varying parameter scales (e.g., Llama-3.1-70B, Qwen3-32B), and our domain-specific models fine-tuned via the proposed Geo-Expert framework.

3.3. Main Results and Insights

The experimental results fully demonstrate the significant advantages of the Geo-Expert framework in breaking through the knowledge barriers of generalist models. We analyze the results from three core dimensions:

1. **Remarkable Performance Breakthrough** As shown in Table 1, general open-weight base models struggle with high-depth geological questions. However, after fine-tuning with the Geo-Expert framework, all models achieve a substantial leap in performance. For instance, the average score of Qwen3-32B surges from 5.00 to 6.82, and Gemma-3-27B jumps from 5.16 to 6.59. Notably, in the “Engineering” dimension—which most rigorously tests comprehensive deductive reasoning—the score of Qwen3-8B rises from 4.41 to 6.44 post-fine-tuning. This strongly suggests that in highly specialized STEM branches, relying solely on general knowledge injected during the pre-training phase is severely limited. The high-quality data synthesis and alignment strategies provided by Geo-Expert effectively reconstruct the models’ domain cognitive networks. Furthermore, a paired t-test on the 387 samples reveals that the performance gains of the Geo-Expert models are statistically extremely significant ($p = 3.7121 \times 10^{-106} < 0.001$), proving that the improvements stem from deep domain alignment rather than random variance.
2. **Extreme Parameter Efficiency** One of the most striking findings of this study is the extreme parameter efficiency exhibited by the fine-tuned models. The Qwen3-8B-geo model, with merely 8B parameters, comprehensively outperforms the internationally leading open-weight benchmark, Llama-3.1-70B-Instruct, across all three dimensions (overall score 6.27 vs. 4.12), and even surpasses the proprietary industry benchmark, GPT-4o (5.93). The suboptimal performance of Llama-3.1-70B profoundly reveals the limitations of the Scaling Law in extremely vertical domains: large parameter models lacking deep alignment with specific scientific corpora are highly prone to superficial term usage and common-sense hallucinations. Conversely, a “small” model with refined domain injection can replace bloated generalist giants in professional scenarios.
3. **Approaching Frontier Performance** Our best-

Table 1. Performance comparison of various models on the Geo-Eval Benchmark. The best results in each column are highlighted in bold, and the second-best are underlined. The Δ column indicates the absolute performance gain achieved by our Geo-Expert fine-tuning framework compared to the respective base model. *** indicates statistical significance ($p = 3.7121 \times 10^{-106} < 0.001$) compared to the corresponding base models via paired t-test.

Model	Size	Concept	Process	Engineering	Average	Δ
<i>Proprietary SOTA</i>						
GPT-5.4 (Frontier)	-	7.35	7.10	7.00	7.15	-
DeepSeek-V3.2-Instruct	-	<u>6.80</u>	6.75	6.67	6.74	-
GPT-4o	-	6.10	5.90	5.80	5.93	-
<i>Open-weight Generalists</i>						
Gemma-3-27B-IT	27B	5.30	5.10	5.08	5.16	-
Qwen3-32B	32B	5.20	4.90	4.90	5.00	-
Qwen3-8B	8B	4.80	4.68	4.41	4.63	-
GLM-4-9B	9B	4.70	4.50	4.51	4.57	-
Llama-3.1-70B-Instruct	70B	4.30	4.10	3.96	4.12	-
Geo-Expert (Ours)						
Qwen3-32B-geo	32B	6.78	<u>6.79</u>	<u>6.90</u>	<u>6.82</u> ***	+1.82
Gemma-3-27B-geo	27B	6.70	6.60	6.47	6.59***	+1.43
Qwen3-8B-geo	8B	6.10	6.27	6.44	6.27***	+1.64

performing model, Qwen3-32B-geo, achieves an outstanding score of 6.82, exhibiting remarkable balance across the Concept (6.78), Process (6.79), and Engineering (6.90) dimensions. This performance not only significantly leads GPT-4o but also trails only marginally behind the current top-tier proprietary model, GPT-5.4 (7.15), which relies on massive computational power and reinforcement learning-based Chain-of-Thought (CoT). Considering its substantial advantages in inference cost and parameter scale, Geo-Expert models demonstrate exceptional value for practical engineering deployment and academic application.

3.4. Qualitative Analysis

To intuitively demonstrate the effectiveness of the Geo-Expert framework in mitigating catastrophic domain-shifting errors, we present a qualitative comparison between GPT-4o and our fine-tuned model (Qwen3-8B-geo) on a highly specialized structural geology question, as shown in Table 2.

As illustrated in Table 2, the base model, lacking deep geological corpus alignment, exhibits a severe cross-domain hallucination. It incorrectly interprets the polysemous term “wedge” as a mechanical or civil engineering component, generating completely irrelevant advice such as using carbon fiber cloth for concrete reinforcement. This literal interpretation leads to a complete failure in answering the geological prompt.

In stark contrast, the Geo-Expert fine-tuned model successfully anchors the semantic context to structural geology. It accurately explains the mechanical responses (stress concentration, fault reactivation) and correctly identifies the complex tectonic accommodation mechanisms (thrust sliding, lateral material flow). This striking contrast (0 vs. 9) highlights a critical vulnerability of generalist LLMs: when encountering polysemous terms across disciplines, they tend to default to high-frequency general concepts. The Geo-Expert framework effectively recalibrates the model’s prior distribution, ensuring precise semantic anchoring and rigorous factual reasoning in the target scientific domain.

4. Related Work

4.1. Large Language Models in Earth Sciences

The application of AI in Earth sciences has grown rapidly, spanning medium-range global weather forecasting (Lam et al., 2023) and scientific language modeling (Taylor et al., 2022). Within geoscience specifically, two complementary lines of work have emerged. The first builds *foundation models* via continued pre-training on broad geoscience corpora: K2 (Deng et al., 2024) adapts LLaMA on geoscience literature and pairs it with a curated knowledge graph, while GeoGalactica (Lin et al., 2024) extends Galactica with over 65B geoscience-specific tokens. These models improve factual recall across general geoscience topics but do not explicitly target multi-step deductive reasoning. The second line develops *geospatial agents and benchmarks*: GeoGPT (Zhang et al., 2023) demonstrates autonomous spatial query

Table 2. Qualitative comparison of model outputs. The base model suffers from cross-domain hallucination, whereas Geo-Expert accurately anchors to geology.

Prompt:

What phenomena does the local thickening of a wedge cause, and how is this deformation accommodated?

GPT-4o:**[Score: 0/10]***Cross-domain Hallucination:*

...Local thickening may alter the structural stress distribution, leading to **concrete cracking** or gear fatigue. To adjust this deformation, it is recommended to use **carbon fiber cloth for reinforcement...**

Geo-Expert: Qwen3-8B-geo**[Score: 9/10]***Accurate Geological Reasoning:*

...Local thickening of a wedge leads to **internal stress concentration** and an increase in thrust... This deformation is primarily accommodated through **thrust fault sliding, extensional collapse, and lateral material flow...**

execution and facility siting from natural language, and Uni-EARTH (Kao et al., 2025) evaluates LLM agents on Earth observation tasks using NASA satellite data.

Despite this progress, a critical gap persists in *solid-earth* geology. Foundation models like K2 and GeoGalactica are trained predominantly on surface-oriented literature (remote sensing, hydrology, atmospheric science), while agent-based systems operate on 2D geospatial workflows. Neither paradigm addresses the core challenges of structural geology and stratigraphy—cross-cutting relationship inference, polyphase tectonic restoration, and deep-time evolutionary reasoning grounded in abstract physical principles. GEO-EXPERT targets this neglected subdomain, shifting the adaptation objective from factual retrieval over surface phenomena to structured reasoning about subsurface geological processes.

4.2. Domain-Specific Adaptation and PEFT

General-purpose LLMs remain prone to hallucination in specialized STEM domains, motivating a growing body of domain-adapted models in medicine (Li et al., 2023), finance (Wu et al., 2023), and law (Cui et al., 2023). A common limitation of these efforts is their emphasis on *factual recall*—e.g., symptom lookup or terminology matching—rather than multi-step deductive reasoning. Recent work on reasoning-oriented distillation (Mukherjee et al., 2023; Yu et al., 2023) suggests that explicit chain-of-thought supervision is essential for eliciting logical inference, yet this insight has rarely been applied to domain-specific scientific adaptation. Parameter-Efficient Fine-Tuning (PEFT), particularly Low-Rank Adaptation (LoRA) (Hu et al., 2021), offers a practical path to internalize domain knowledge while preserving base-model capabilities. However, its effectiveness is fundamentally bounded by the quality of the

supervision signal: tuning on raw domain text teaches the model *what to say* but not *how to reason*. In the geological setting this gap is especially acute, as correct answers often require chaining stratigraphic principles, cross-cutting relationships, and tectonic context across multiple inference steps. GEO-EXPERT addresses this by coupling LoRA with a CoT-enhanced instruction synthesis pipeline, shifting the adaptation objective from surface-form memorization to structured geological reasoning.

4.3. Scientific Benchmarks and LLM-as-a-Judge

Evaluating LLMs in specialized scientific domains poses a significant challenge. Traditional benchmarks often rely on static, multiple-choice questions extracted from standardized exams (Hendrycks et al., 2021). While useful for assessing factual recall, these static datasets fail to capture the multi-step deductive logic required in real-world scientific research. To overcome the limitations of human evaluation scalability, the “LLM-as-a-Judge” paradigm has been increasingly adopted to assess open-ended, complex reasoning tasks (Zheng et al., 2024). Geo-Eval advances this paradigm by introducing an adversarial mining pipeline. By isolating “hard boundary” samples where general-purpose models fail but domain-aligned models succeed, Geo-Eval provides a much more rigorous and discriminative metric for assessing true expert-level geological reasoning.

5. Conclusion and Future Work

In this work, we present Geo-Expert, a framework that addresses the critical reasoning bottlenecks of generalist LLMs in highly specialized scientific domains. At the heart of our contribution is a robust, scalable dataset construction and instruction synthesis pipeline. Rather than passively fine-tuning, this pipeline achieves deep domain alignment by systematically recalibrating the model’s semantic representations. By doing so, Geo-Expert successfully transforms text generation into well-rounded scientific reasoning.

Beyond achieving state-of-the-art domain performance, this research exposes a profound paradigm shift regarding the Scaling Law in vertical disciplines. Our empirical results demonstrate extreme parameter efficiency: a domain-aligned 8B model can comprehensively eclipse a 70B open-weight giant. This unequivocally proves that high-quality data curation and targeted logical alignment profoundly outweigh sheer parameter volume. Consequently, Geo-Expert provides a highly cost-effective and easily deployable solution. It enables localized, flexible deployment for real-world geological engineering and on-site exploration, freeing researchers from the reliance on prohibitive computational clusters.

Our work has several limitations. First, the training cor-

pus is built from five canonical textbooks, biased toward structural geology, stratigraphy, and tectonics; mineralogy, geochemistry, and geophysics are under-represented. Second, the Geo-Eval benchmark contains 387 expert-vetted questions—sufficient for boundary-aware stress testing but smaller than general-purpose benchmarks. Finally, the current framework is text-only and does not yet handle the inherently multimodal nature of geological data (cross-sections, well logs, field photos).

While Geo-Expert establishes a formidable foundation for text-based geological AI, we acknowledge that Earth sciences are inherently multimodal. Future work will extend this parameter-efficient framework by integrating Vision-Language Models (VLMs) to parse stratigraphic profiles and remote sensing imagery. Coupled with Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) to dynamically track evolving literature, we aim to evolve Geo-Expert from a static reasoning engine into a dynamic, autonomous agent for geoscience discovery.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Catuneanu, O. *Principles of sequence stratigraphy*. Elsevier, 2006.
- Cui, J., Li, Z., Yan, Y., Chen, B., and Yuan, L. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *arXiv preprint arXiv:2306.16092*, 2023.
- DeepSeek-AI et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Deng, C., Zhang, T., He, Z., Xu, Y., Chen, Q., Shi, Y., Fu, L., Zhang, W., Wang, X., Zhou, C., Zhuang, L., et al. K2: A foundation language model for geoscience knowledge understanding and utilization. *Proceedings of the 17th ACM International Conference on Web Search and Data Mining (WSDM)*, pp. 161–170, 2024.
- Fossen, H. *Structural geology*. Cambridge university press, 2016.
- Gao, D. (ed.). *Tectonics and Sedimentation: Implications for Petroleum Systems*. AAPG Memoir 100. American Association of Petroleum Geologists, 2012.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Kao, C. H., Zhao, W., Revankar, S., Speas, S., Bhagat, S., Datta, R., Phoo, C. P., Mall, U., Vondrick, C., Bala, K., and Hariharan, B. Towards llm agents for earth observation. *arXiv preprint arXiv:2504.12110*, 2025. Accepted at ICML 2025 Workshop TerraBytes.
- Lam, R., Sanchez-Gonzalez, A., Willson, M., Wiratunga, P., Ewalds, T., Ferraro, A., Pritzel, A., et al. Learning skillful medium-range global weather forecasting. *Science*, 382(6677):1416–1421, 2023.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, 2020.
- Li, Y., Li, Z., Zhang, K., Dan, R., Jiang, S., and Zhang, Y. Chatdoctor: A medical chat model fine-tuned on a large language model meta-llama. *arXiv preprint arXiv:2303.14070*, 2023.
- Lin, Z., Ding, C., Yao, H., Hao, Z., Zheng, J., et al. Geogalactica: A scientific large language model in geoscience. *arXiv preprint arXiv:2401.00434*, 2024.
- ModelScope. ms-swift: Scalable lightweight infrastructure for fine-tuning. <https://github.com/modelscope/swift>, 2024.
- Mukherjee, S., Mitra, A., Jawahar, G., Agarwal, S., Palangi, H., and Awadallah, A. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.
- OpenDataLab. Mineru: A high-precision pdf to markdown extraction tool. <https://github.com/opendatalab/MinerU>, 2024.
- Rowland, S. M., Duebendorfer, E. M., and Schiefelbein, I. M. *Structural analysis and synthesis: a laboratory course in structural geology*. John Wiley & Sons, 2020.
- Taylor, R., Kardas, M., Cucurull, G., Scialom, T., Hartshorn, A., Elvis, E., Pumperla, M., Toyama, Y., Guillemotte, T., and Jones, L. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.

385 Team, G. Gemma: Open models based on gemini research
386 and technology. *arXiv preprint arXiv:2403.08295*, 2024.
387
388 Wu, S., Irsoy, O., Lu, S., Dabravolski, V., Dredze, M.,
389 Gehrmann, S., Kambadur, P., Rosenberg, D., and Mann,
390 G. Bloomberggpt: A large language model for finance.
391 *arXiv preprint arXiv:2303.17564*, 2023.
392
393 Yang, A., Yang, B., Zhang, B., et al. Qwen3 technical report.
394 *arXiv preprint arXiv:2505.09388*, 2024.
395
396 Yu, L., Jiang, W., Shi, H., Yu, J., Liu, Z., Zhang, Y., Kwok,
397 J. T., Li, Z., Weller, A., and Liu, W. Metamath: Boot-
398 strap your own mathematical questions for large language
399 models. *arXiv preprint arXiv:2309.12284*, 2023.
400
401 Zeng, A., Liu, H., Du, Z., et al. Glm-4: An open language
402 model family. *arXiv preprint arXiv:2406.12793*, 2024.
403
404 Zhang, Y., Wei, C., Wu, S., He, Z., and Yu, W. Geogpt:
405 Understanding and processing geospatial tasks through
406 an autonomous gpt. *arXiv preprint arXiv:2307.07930*,
2023.
407
408 Zheng, L., Chiang, W.-L., Sheng, Y., Hao, S., Wu, Z., Ba, S.,
409 Zhuang, Y., Lin, Z., Li, Z., Xing, E., et al. Judging llm-
410 as-a-judge with mt-bench and chatbot arena. In *Advances*
411 *in Neural Information Processing Systems*, 2024.
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439

440 A. Implementation Details

441 In this section, we provide the detailed hyperparameter configurations and hardware setups used for fine-tuning the Geo-
442 Expert model family. All experiments were conducted using the `ms-swift` framework to ensure standardized PEFT across
443 different model architectures.

445 To balance training efficiency and memory footprint, we customized the LoRA settings based on the base models’ parameter
446 scales. For the compact Qwen3-8B model, we utilized a single NVIDIA RTX 5090 GPU with FP16 precision. For the larger
447 architectures (Gemma-3-27B and Qwen3-32B), we scaled the hardware to a 4×RTX 5090. To mitigate the OOM issues
448 on these larger models, we activated Gradient Checkpointing and employed BF16 mixed-precision training. The detailed
449 hyperparameter settings are summarized in Table 3.

450
451 *Table 3.* Hyperparameter configurations for fine-tuning the Geo-Expert model family.

453 HYPERPARAMETER	454 8B MODEL	455 27B/32B MODELS
456 LORA TARGET MODULES	457 ALL LINEAR	458 ALL LINEAR
459 LORA RANK (r)	460 32	461 64
462 LORA ALPHA (α)	463 32	464 128
465 LORA DROPOUT	466 0.05	467 0.05
468 LEARNING RATE	469 2E-5	470 1E-5
471 OPTIMIZER	472 ADAMW	473 ADAMW
474 PRECISION	475 FP16	476 BF16
477 GRADIENT CHECKPOINTING	478 FALSE	479 TRUE
480 BATCH SIZE	481 4	482 2 (PER DEVICE)
483 GRADIENT ACCUMULATION	484 1	485 4
486 HARDWARE	487 1×RTX 5090	488 4×RTX 5090

466 B. Prompt Templates

467 B.1. LLM-as-a-Judge Prompt for Geo-Eval

469 To conduct a rigorous and impartial evaluation during the hard sample mining phase of Geo-Eval, we utilized GLM-4.5
470 as the evaluator. We designed a zero-shot prompt template that constrains the judge to focus purely on geological factual
471 accuracy and logical deduction, penalizing superficial or hallucinatory responses. The specific prompt used is as follows:

472
473 *“You are an expert professor in solid earth geology. Your task is to evaluate the responses of two AI models to a
474 complex geological question.*

475 ***Question:** [Insert Question]*

476 ***Model A Response:** [Insert Model A Answer]*

477 ***Model B Response:** [Insert Model B Answer]*

478 ***Evaluation Criteria:** 1. Factual Accuracy (40%): Does the response align strictly with established geological
479 principles (e.g., Andersonian faulting, sequence stratigraphy)? 2. Deductive Logic (40%): Is the chain of thought
480 coherent? 3. Professional Terminology (20%): Are the geological terms used accurately without cross-domain
481 confusion?*

482 ***Task:** Score each model from 0 to 10. You must provide a brief justification for your score, explicitly pointing out
483 any hallucinations or logical flaws. Output your final scores in the format: Score A: [X], Score B: [Y].”*

485 B.2. Reference-Based Evaluator Prompt for Final Scoring

487 To conduct the final quantitative evaluation, we utilized GPT-4o as a reference-guided evaluator. Unlike zero-shot blind
488 testing, this phase employs a strict Reference-Based Evaluation paradigm, where the model compares candidate responses
489 against expert-verified ground truth. The prompt explicitly enforces terminology flexibility to accurately accommodate the
490 nuances of geological nomenclature.

491
492 You are an authoritative Professor of Geology and an expert academic evaluator.
493 Your task is to grade a Model’s Answer to a complex geological question against
494

495 a provided professional Reference Answer.
496
497 [Grading Criteria]
498 1. Score from 0 to 10.
499 2. Focus on SEMANTIC CORRECTNESS and GEOLOGICAL ACCURACY, not verbatim matching.
500 3. Terminology Flexibility: Be lenient with accepted geological synonyms (e.g.,
501 "igneous" vs "magmatic", "thrust fault" vs "reverse fault" depending on context).
502 4. Completeness: A 10/10 answer must contain all core geological facts mentioned
503 in the reference.
504 5. Penalty: Deduct heavily for hallucinations, factually incorrect statements,
505 or mixing up crucial concepts (e.g., misidentifying rock types or eras).
506

507 [Output Format]
508 Output ONLY a JSON object with this exact structure:
509 {"score": 8, "reason": "brief professional explanation"}
510

511 Question: {question}
512 Reference Answer: {reference}
513 Model Answer: {model_ans}
514

515 516 **B.3. Instruction-Tuning Data Example** 517

518 As described in Section 2.2, our instruction synthesis pipeline generates high-quality, reasoning-enhanced QA pairs. Below
519 is a translated JSON snippet extracted from our training corpus, demonstrating the required <think> trajectory before
520 final generation.
521

```
522  
523 [  
524   {  
525     "messages": [  
526       {  
527         "role": "user",  
528         "content": "How is the method of measuring vorticity using rotated  
529           porphyroblasts applied to the study of rocks near the  
530           Alpine Fault?"  
531       },  
532       {  
533         "role": "assistant",  
534         "content": "<think>\n1. Identify core concepts: vorticity measurement,  
535           rotated porphyroblasts, Alpine Fault.\n2. Recall geological  
536           context: The Alpine Fault is a major plate boundary fault  
537           with intense ductile deformation.\n3. Formulate explanation:  
538           Porphyroblasts record kinematic history during shear  
539           deformation...\n</think>\n\nThe method of measuring  
540           vorticity using rotated porphyroblasts is a sensitive  
541           kinematic technique used to study the behavior of rocks  
542           during shear deformation. By analyzing the rotational  
543           dynamics of these crystals within the ductile shear zones..."  
544       }  
545     ]  
546   }  
547 ]
```

550 **B.4. Geo-Eval Hard Boundary Sample Analysis**

551 In Section 2.4, we detailed the automated mining of “hard boundary” samples where general-purpose models fail but
552 domain-expert models succeed. The following real evaluation cases illustrate severe cross-domain hallucinations by the base
553 model and precise domain anchoring by Geo-Expert.
554

555 **Case 1: Semantic Disambiguation of Polyphase Deformation**

556 **Question:** What is polyphase deformation?

557 **Ground Truth:** In structural geology, Polyphase Deformation (or Superimposed Deformation) refers to the complex
558 process where the same geological body undergoes two or more distinct regional tectonic stress fields over a long period of
559 geological history...
560

561 **Base Model Response (Score: 0.0 / 10):**

562 **[Exposed Thinking Process]:** *Hmm, the user is asking about “polyphase deformation”. First, I need to determine the*
563 *context of this term. It could appear in materials science, engineering, or computer graphics. For instance, in computer*
564 *graphics, it might involve multiple steps of deformation algorithms...*

565 **[Final Output]:** *“Polyphase deformation” is an interdisciplinary concept. In materials science, it refers to different stages*
566 *of deformation under varying temperatures or stress... In computer graphics, it refers to the combination of multiple*
567 *animation steps...*
568

569 **Geo-Expert Response (Score: 9.0 / 10):**

570 **[Final Output]:** *In structural geology, Polyphase Deformation refers to the complex superimposed deformation process*
571 *generated by two or more distinct tectonic stress fields over geological time. Each independent deformation event leaves*
572 *specific structural traces in the rock, such as folds, faults, or foliations. Studying polyphase deformation is of decisive*
573 *significance for reconstructing paleostress fields and understanding tectonic evolution...*
574

575 **Case 2: Disciplinary Anchoring of “Deformation” Terminology**

576 **Question:** What is the academic definition of deformation terminology?

577 **Ground Truth:** The academic definition of deformation refers to the process and result of relative positional changes among
578 internal particles of an object under external factors (e.g., stress, temperature), leading to changes in shape, size, orientation,
579 or volume. It is a core concept in structural geology and solid mechanics...
580

581 **Base Model Response (Score: 0.0 / 10):**

582 **[Exposed Thinking Process]:** *Hmm, the user is asking about the academic definition of “deformation terminology”. This*
583 *could have different interpretations in linguistics, philosophy, or sociology. For instance, in linguistics, it might refer to*
584 *inflection or word morphing. Or it could mean a term that is morphed or transformed across contexts...*

585 **[Final Output]:** *“Deformation terminology” does not have a universally accepted definition. Its meaning varies by*
586 *discipline: In linguistics, it may refer to word inflection or semantic shifts. In philosophy, it may refer to conceptual*
587 *evolution. In sociology, it implies the contextualization of a term across different social groups...*
588

589 **Geo-Expert Response (Score: 9.0 / 10):**

590 **[Final Output]:** *The academic definition of deformation terminology refers to the scientific concepts used to describe the*
591 *shape changes of rocks under stress, encompassing major types such as elastic, plastic, and brittle deformation. Elastic*
592 *deformation refers to reversible shape changes... Furthermore, it includes key concepts like stress, strain, and strain rate, as*
593 *well as deformation mechanisms (e.g., dislocation creep). These terms collectively form a systematic description of rock*
594 *deformation processes, providing a theoretical foundation for studying geological structures.*
595
596
597
598
599
600
601
602
603
604