

# MSEarth: A Multimodal Benchmark for Earth Science Phenomenon Discovery with MLLMs

Anonymous ACL submission

## Abstract

The rapid advancement of multimodal large language models (MLLMs) offers new opportunities for complex scientific challenges, yet their application in earth science—especially at the graduate level—remains underexplored due to a lack of benchmarks reflecting the depth and complexity of geoscientific reasoning. Existing datasets often rely on synthetic data or simple figure-caption pairs, failing to capture the nuanced reasoning required for real-world applications. To address this, we introduce MSEarth, a multimodal scientific dataset and benchmark curated from high-quality, open-access publications. Covering the five major spheres of Earth science—atmosphere, cryosphere, hydrosphere, lithosphere, and biosphere—MSEarth features over 289K figures with refined captions enriched by contextual discussions and reasoning from the original papers. The benchmark supports tasks such as scientific figure captioning, multiple choice questions, and open-ended reasoning, providing a scalable, high-fidelity resource for developing and evaluating MLLMs in scientific reasoning. MSEarth is publicly available to foster further research and innovation: <https://anonymous.4open.science/r/MSEarth-2B3F>.

## 1 Introduction

The advent of multimodal large language models (MLLMs) (Liu et al., 2023a; Liang et al., 2024b) has transformed AI, enabling major advances across scientific fields. Examples include ChemVLM (Li et al., 2025) in chemistry, GeoChat (Kuckreja et al., 2024) in geography, and WeatherQA (Ma et al., 2024) in atmospheric science. These models perform domain-specific visual question answering by integrating specialized knowledge: ChemVLM supports analysis of molecular structures, reactions, and chemistry exam questions, while WeatherQA enables reasoning about severe weather events in real-world settings.

Building MLLMs that understand advanced geoscientific knowledge requires rigorous datasets and benchmarks to improve performance on complex, discipline-specific problems. As shown in Table 1, existing benchmarks often use synthetic data or high school/undergraduate textbook materials (Lu et al., 2022; Yue et al., 2024), which lack the depth needed for professional, graduate-level tasks. Recent work (Li et al., 2024c; Roberts et al., 2024; Li et al., 2024b) uses academic papers to build multimodal scientific benchmarks, leveraging graduate-level complexity, but typically extracts only figures and captions and ignores key reasoning in the surrounding text. As a result, tasks are often oversimplified as basic *figure-caption matching*, offering limited insight into a model’s reasoning ability.

A further challenge is designing questions that rigorously evaluate MLLMs’ data-analysis ability to uncover Earth-science phenomena from observational imagery. In scientific papers, images are often paired with hypotheses, evidence, analyses, and conclusions that are mainly in the main text rather than captions. As Figure 1(a) shows, caption-only question generation oversimplifies tasks: question quality and difficulty are constrained by generation models, and missing paper context makes verification difficult. Existing benchmarks thus neglect the high knowledge density of scientific reasoning, raising the challenge: *How to effectively align high-value Earth science images with long-context information for phenomena uncovering?*

We address this with a new benchmark construction approach featuring two innovations. First, we introduce the *refined caption*. Observational images visualize phenomena, while deeper insights—hypotheses, supporting evidence, analytical reasoning, and conclusions—are embedded in the paper body. Raw captions are brief and lack context for complex reasoning; refined captions combine the raw caption with relevant domain-specific information extracted from the paper to

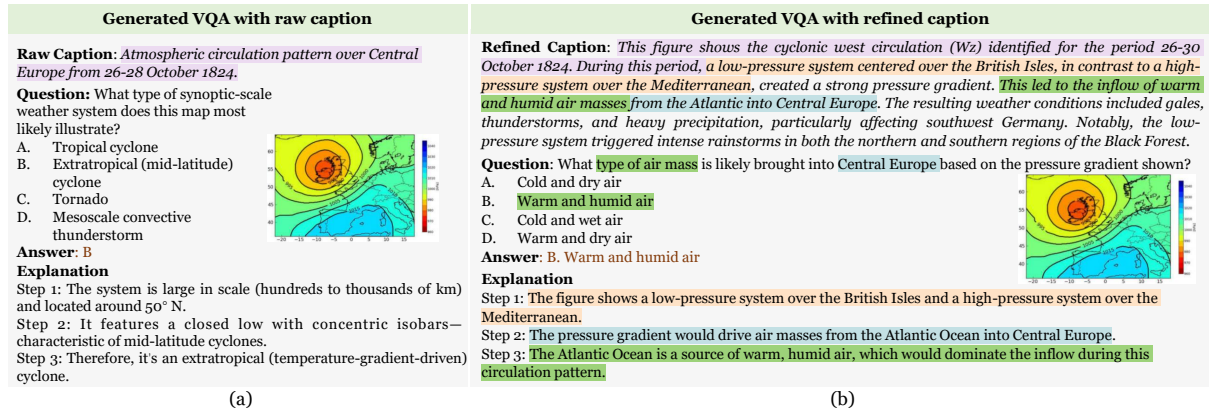


Figure 1: Illustration of VQA generation methodologies: (a) VQA relying exclusively on figure captions, and (b) VQA utilizing refined captions that integrate figure captions with content from academic papers. Highlighted areas denote questions and answers supported by evidence.

Table 1: Comparison with previous multimodal scientific benchmarks. Task types include OE (Open-ended QA), MCQ (Multiple-choice QA), and CG (Caption Generation). Originality refers to newly constructed QAs rather than those collected from existing benchmarks or textbooks. Validated indicates whether the QAs are verified by experts or supported by evidence from papers.

Benchmark Dataset	Science Topics	Tasks	Difficulty	#Ques.	Originality	Validated
<b>Multimodal Scientific benchmarks</b>						
ScienceQA (Lu et al., 2022)	General Science	MCQ	Primary	21,208	✗	✓
ScemQA (Liang et al., 2024a)	General Science	MCQ,OE	Pre-College	1,045	✗	✓
Mmmu (Yue et al., 2024)	General Science	MCQ,OE	College-level	11,550	✗	✓
OlympiadBench (He et al., 2024a)	Math, Physics	OE	Competition	8,476	✗	✓
EMMA (Hao et al., 2025)	General Science	MCQ	Mixed	2,788	✓	✓
<b>Paper-Based Multimodal Scientific benchmarks</b>						
SciFIBench (Roberts et al., 2024)	General Science	MCQ	Graduate-Level	2,000	✗	✓
ArxivCap/QA (Li et al., 2024b)	General Science	CG,MCQ	Graduate-Level	32K/100K	✓	✗
MMSci (Li et al., 2024c)	General Science	CG,MCQ	Graduate-Level	1,079,797	✗	✗
MSEarth(Ours)	Earth Science	CG,MCQ,OE	Graduate-Level	448,980	✓	✓

provide a more complete, scientifically meaningful image description. As in Figure 1(b), refined-caption-based questions are higher quality and can be supported and validated by paper content, ensuring professionalism and accuracy. Second, we apply rigorous quality control by combining multi-agent automated evaluation with expert evaluation, ensuring questions are relevant, coherent, and matched to professional-level geoscientific phenomena uncovering. This improves reliability and accuracy beyond current LLM-based scientific question generation, producing a robust and valuable dataset for advancing MLLMs in scientific domains.

Using this adaptive annotation methodology, we present MSEarth, a comprehensive multimodal benchmark for graduate-level Earth science. It is built from 64,560 open-access publications across five spheres, eight subjects, and 66 sub-subjects, from which we extract 289,891 figures. We

augment figures with refined captions averaging 136.29 tokens (vs. 37.56 originally). The test set unifies scientific figure captioning with multiple-choice and open-ended reasoning tasks for holistic, interdisciplinary evaluation, enabling rigorous assessment of MLLMs in professional geoscientific contexts and filling a key gap in graduate-level multimodal benchmarks. MSEarth also provides a scalable, high-fidelity pipeline for domain-specific scientific benchmark construction. Using this framework, we create a resource-rich training dataset spanning captioning, open-ended QA, and multiple-choice tasks for post-training (e.g., instruction tuning and GRPO-based reinforcement learning), yielding substantial gains across tasks for open-source models. Our contributions are:

**Development of a Scalable Adaptive Framework:** We introduce a semi-automated tool for automatic VQA generation and machine-assisted filtering, providing a robust, scalable method for

high-fidelity domain-specific benchmarks extensible to other scientific fields.

**High-Quality benchmark Resources for Earth Science MLLMs:** We provide an expert-annotated graduate-level Earth science benchmark and a diverse training corpus (captioning, open-ended QA, multiple-choice) to support advanced post-training of MLLMs.

**Comprehensive Evaluation and Validation of State-of-the-Art MLLMs:** We extensively evaluate MLLMs on MSEarth, offering insights into limitations and future directions, and demonstrate the effectiveness of our training data by building a state-of-the-art baseline model.

## 2 MSEarth - A Multimodal Scientific Benchmark for Earth Science

This section provides a detailed overview of the construction process for MSEarth. As illustrated in Figure 2, we outline the framework used to develop MSEarth from open-access scientific publications. The section is organized into three main parts: first, we detail the data collection and preprocessing steps. Next, we elaborate on the construction procedures for the two D&B within MSEarth, namely MSEarthCap and MSEarthQA. Finally, we describe the process of ensuring the reliability of the test data, which involves expert annotation and manual screening of the sampled test data.

### 2.1 Data Preparation

The first part of the D&B construction focuses on data collection and preprocessing. The data collection begins with more than 400K Earth science papers obtained in PDF format. These are uniformly converted into structured JSON text using the MinerU (Wang et al., 2024) parser. To classify the papers, semantic similarity is calculated between the abstracts and keywords from the five Earth spheres: hydrosphere, biosphere, lithosphere, atmosphere, and cryosphere. Based on this, the papers are assigned to respective disciplinary categories. Details are provided in Appendix F.3. We then selected papers based on the criterion of containing high-quality, Earth science-related images, resulting in a subset of around 83k papers. Specifically, Qwen-2.5-VL-72B (Bai et al., 2025b) is utilized to filter and select images, with the filtering prompts detailed in the Appendix F.4.

### 2.2 MSEarthCap

**Figure-Caption Extraction:** Figures and their corresponding captions are extracted from the JSON files processed by MinerU. As shown in Appendix F.2, MinerU has already extracted the figures along with their original captions, which can be directly utilized for subsequent processing. To ensure accurate alignment between figures and their references within the text, we employ a regex-based method to identify the labels of each figure. This approach enables precise matching between the figures and the relevant sections of the articles.

**Relevant Context Extraction:** To further enrich the captions with contextual information, we use the figure labels obtained in the previous step to perform approximate matching against the main body of the paper. Since MinerU processes papers with segmented paragraphs, we apply regular expression matching to each paragraph to extract contextual text that references the target figure. This ensures the inclusion of descriptions and reasoning associated with each figure within the paper. To guarantee that the extracted context provides sufficient detail about the target figure, only paragraphs exceeding two sentences were included in the final dataset. From this filtered subset, we selected around 64K papers that met the criteria for subsequent processing. For more details, refer to Appendix F.5.

**Refined Caption Generation:** To create professional-level figure descriptions, we employ GPT-4o for refined caption generation. The model takes as input the extracted figure, its original caption, and the contextual text from the relevant sections of the paper. Figure refinement is performed only for data that includes valid relevant contexts. The specific prompts used for this process are detailed in the Appendix F.6. After statistical analysis, we observe that the average word length of the raw captions is 37.56, while the average length of the refined captions increases to 136.29, reflecting the incorporation of richer, domain-specific content.

### 2.3 MSEarthQA

To generate high-quality multiple-choice questions (MCQs) and open-ended questions, we use a question generation pipeline that takes the figure, its original caption, and the refined caption as input. The generation prompts, detailed in the Appendix F.6, are crafted to encourage the model to highlight differences between the original and re-

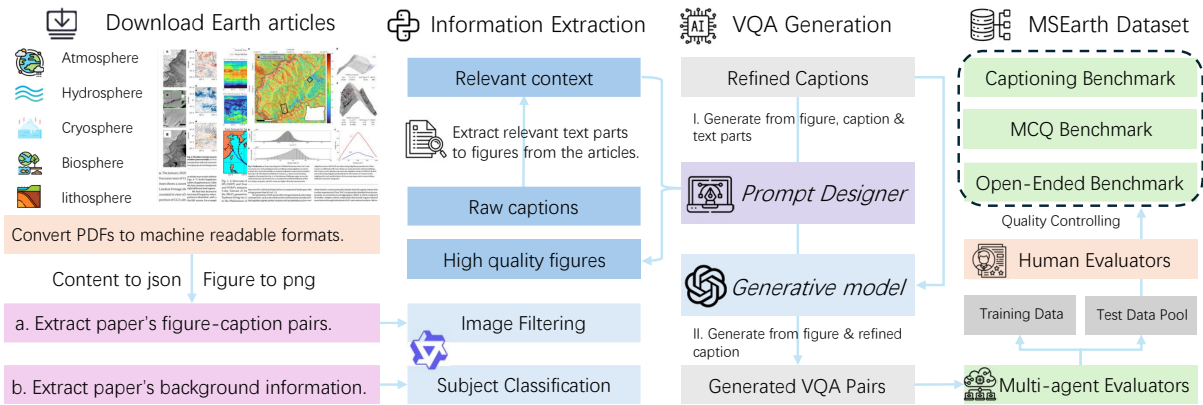


Figure 2: Data curation process for MSEarth. The two parts on the left represent data preprocessing, while the two parts on the right encompass the automated generation of VQA and expert-AI collaborative filtering.

221 fined captions, ensuring that the generated ques- 257  
 222 tions are grounded in evidence from the paper. The 258  
 223 questions are constructed using GPT-4o to main- 259  
 224 tain high detail and relevance. However, due to 260  
 225 inherent challenges such as self-inconsistency and 261  
 226 uncertainty in the generation process, the gener- 262  
 227 ated questions undergo an automated and expert 263  
 228 validation process to ensure quality. 264

### 2.3.1 Automated Validation 265

230 Inspired by the LLM Voting (Yang et al., 2024b; 266  
 231 Lee et al., 2025; Kaesberg et al., 2025) method, we 267  
 232 developed a multi-agent, voting-based approach 268  
 233 to automate the validation of generated ques- 269  
 234 tions. Specifically, we employ a *Majority Voting* 270  
 235 strategy, where multiple agents independently gener- 271  
 236 ate responses, and the final decision is based on the 272  
 237 majority consensus of these agents. In our setup, we 273  
 238 utilize the following MLLMs for decision-making: 274  
 239 Qwen2.5-VL-72B, Qwen2.5-VL-7B, InternVL2.5- 275  
 240 7B, InternVL2.5-78B, and GPT-4o. A key aspect 276  
 241 of our evaluation process is the use of the *refined* 277  
 242 *caption*, which incorporates scientists’ reasoning 278  
 243 and insights about the figure extracted from the 279  
 244 paper. This refined caption provides additional 280  
 245 context and domain-specific information that goes 281  
 246 beyond the original caption. By comparing model 282  
 247 performance with and without the refined caption, 283  
 248 we can assess the quality of the questions and de- 284  
 249 termine whether they test the model’s ability to 285  
 250 grasp scientific reasoning and insights. The de- 286  
 251 tailed decision-making process is outlined below: 287

252 **Phase A:** The question and original caption are 288  
 253 provided to a suite of models  $\{M_1, M_2, \dots, M_n\}$ . 289  
 254 The correctness of the model responses is used 290  
 255 to evaluate the types and quality of the questions. 291  
 256 A threshold of 60% is defined for supermajority

257 voting. Specifically, if more than 40% of the mod- 258  
 259 els produce incorrect responses, the question is 260  
 261 flagged for further analysis. Additionally, we dis- 262  
 263 card questions that all models answer correctly, as 264  
 265 these questions do not contribute to the effective 266  
 267 testing or training of the models. Questions iden- 268  
 269 tified through this process are categorized as either 270  
 271 potentially difficult or of poor quality, with the dis- 272  
 273 tinction made in subsequent phases. 274

266 **Phase B:** In this phase, the question and refined 267  
 268 caption are provided to the same suite of models. 269  
 270 If more than 60% of the models answer the ques- 271  
 272 tion correctly with the refined caption, it indicates 273  
 274 that the question requires relatively specialized sci- 275  
 276 entific knowledge to answer. Such questions are 277  
 278 categorized as *specialized QA*, as their answers 279  
 280 rely on the model’s ability to understand and ap- 281  
 282 ply specific domain knowledge rather than simply 283  
 284 perceiving the image or relying on commonsense 285  
 286 reasoning. Questions that fail this phase proceed to 287  
 288 the next stage for further evaluation. 289

278 **Phase C:** In this phase, only models with 70B+ 279  
 280 parameters are used for voting, including GPT- 281  
 282 4o (the same model used for question genera- 283  
 284 tion), InternVL2.5-78B, and Qwen2.5-VL-72B. 285  
 286 The question and its refined caption are provided 287  
 288 to these large-scale models. If more than 60% of 289  
 290 the large models answer the question correctly, it 291  
 292 suggests that the difficulty of the question likely 293  
 294 lies in the model’s ability to perceive and interpret 295  
 296 the image content. Such questions are categorized 297  
 298 as *hard QA*. Subsequent human validation will in- 299  
 300 volve sampling and additional annotation across 301  
 302 QA filtered in all phases to ensure the overall qual- 303  
 304 ity and accuracy of the benchmark. 305

292 This pipeline identifies high-quality questions by 293  
 294 filtering out overly simplistic or poorly constructed 295

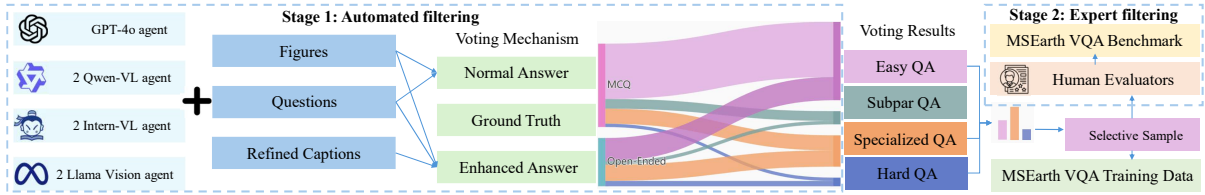


Figure 3: Overall approach of our multi-agent, voting-based approach to automate the validation of generated questions.

ones. In **Phase A**, approximately 70% of the questions were categorized as *easy*, as most models could answer them correctly without refined captions. After **Phase B**, around 20% were classified as *specialized QA*, where refined captions enabled correct answers, indicating the need for domain-specific knowledge. In **Phase C**, 5% were labeled as *hard QA*, requiring high-performing models to interpret image content accurately, suggesting that these questions test the model’s ability to perceive and interpret image content. The remaining 5% were deemed flawed and discarded. Detailed processes and examples are provided in Appendix G. For the training dataset, we sampled more than 150K VQA pairs, including both multiple-choice and open-ended questions. Of these, 20% were drawn from Phase A and 80% from Phase B, ensuring a balanced distribution of question difficulty.

Table 2: Main statistics in MSEarth-Bench. The dataset covers 66 subjects and 64,560 articles.

Statistic	Training Set	Test Set
Captioning as QA	289,891	3,000
MCQ	102,753	2,784
Open-Ended QA	49,141	1,411
<b>Total Questions</b>	<b>441,785</b>	<b>7,195</b>
Total Figures	289,891	
Avg. Caption Length	37.56 (Original) / 136.29 (Refined)	

### 2.3.2 Expert Validation

Ensuring that synthetic data closely mirrors real-world distributions is critical for evaluation tasks. To achieve this, domain experts are engaged to review and annotate the curated QA pairs for accuracy and relevance. During this process, low-quality or invalid questions are identified and filtered out to ensure the overall quality of the dataset. The annotation process is conducted from two perspectives: image types and question types. For im-

age types, we categorize the data into three groups: single-image question answering, single-image-focused question answering within multi-image figures, and multi-image relational question answering. For question types, we define two categories: scientific reasoning and perception questions. Scientific reasoning questions are constructed based on inferences or scientific discoveries presented in research papers, making them more specialized and challenging. In contrast, image perception questions focus on interpreting images and require less background knowledge of scientific concepts. This expert-AI collaborative process, combined with rigorous quality control, results in a high-quality dataset comprising 1,500 open-ended questions and 3,000 MCQs, forming the MSEarthQA benchmark. The annotated results are summarized in Table 2, with further details and analysis provided in Appendix H and Table 8.

## 3 Experiments

### 3.1 Evaluated Models

We evaluate different families of MLLMs on our benchmark. We evaluate the following closed-source models: GPT-4 series (Hurst et al., 2024), Gemini-2.5 series (Team et al., 2023) and Claude-3 series (Anthropic, 2024). We also evaluate the following open-source models: LLaVA-OneVision (Li et al., 2024a), Qwen-2.5-VL (Yang et al., 2024a), InternVL2.5/3/S1 (Chen et al., 2024; Zhu et al., 2025; Bai et al., 2025a) and Llama-3.2-Vision-Instruct (Grattafiori et al., 2024). We use chat/instruction-tuned variants of each model and compare the performance of multiple model sizes where available. Details can be found in Appendix J. To validate the effectiveness of our training data, we conducted post-training on the Interns1-mini and Qwen-2.5-VL-7b models. Specifically, we employed instruct tuning method for fine-tuning on captioning and open-ended QA tasks. For the MCQ task, we applied GRPO (Shao et al., 2024) reinforcement learning method.

### 3.2 Evaluation Metrics

Both captioning and open-ended QA tasks require generating freeform textual outputs grounded in complex scientific data and reasoning. To evaluate these tasks, we use lexical overlap-based metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), and METEOR (Banerjee & Lavie, 2005) for surface-level similarity, while BERTScore (Zhang et al., 2019) assesses deeper semantic alignment. Additionally, following G-Eval (Liu et al., 2023b), we utilize the Qwen2.5-VL-72B model with a specialized prompt to compute a factual scientific score. For the captioning task, we define a Cap-Eval score ranging from 1 to 5, where higher scores indicate better caption quality. For the open-ended QA task, we introduce OE-Eval, which evaluates the reasonableness of generated answers using a binary 0/1 scoring system.

Evaluating MCQs is relatively straightforward, as these tasks require selecting the correct answer from a predefined set of options. In our experiments, models were guided by carefully structured prompts to ensure adherence to a specific output format. Regular expression rules were employed to extract the selected choice, ensuring strict alignment with the predefined answer format. Detailed evaluation metrics are provided in Appendix K.

### 3.3 Main Results

**Current models struggle with scientific question-answering tasks, particularly on questions requiring specialized knowledge and reasoning across multiple images.** The performance of MCQs is summarized in Table 3. The results reveal that most models do not perform exceptionally well on scientific question-answering tasks, with proprietary models generally achieving better results. Further analysis of the models’ failure rates on reasoning and perception-based questions is provided in the Appendix (Figure 13). This analysis shows that models are more prone to errors on questions requiring specialized knowledge, underscoring significant room for improvement in scientific reasoning question-answering. In contrast, for relatively simpler perception-based questions, which require less domain-specific knowledge, the models tend to perform better. Similarly, when analyzing performance across different image types, we observe that most models achieve their best results on tasks involving single-image inputs. How-

ever, for tasks requiring multi-image inputs, particularly those that demand reasoning across multiple images to derive an answer, the models perform the worst. Additional experimental results can be found in Appendix L.

**Proprietary models consistently outperform open-source models in both Scientific Figure Captioning and Open-Ended VQA tasks, with LLM-based metrics providing a more reasonable evaluation.** The captioning results are presented in Table 4, where overlap-based, similarity-based, and LLM-based metrics exhibit similar trends, with no significant differences observed among the overlap-based and similarity-based metrics. The Gemini-3-Flash model achieves the best performance across most metrics. The LLM-based metric, designed to evaluate the professionalism and accuracy of generated captions, demonstrates greater variance compared to similarity-based metrics, making it more suitable for assessing the Scientific Figure Captioning task. Open-source models still show a noticeable gap compared to proprietary models, consistent with the findings from the MCQ results, suggesting a close interconnection between a model’s understanding and reasoning capabilities. Similarly, the results for open-ended question answering, presented in Table 4, show that overlap-based and similarity-based metrics tend to yield higher scores due to the shorter nature of both ground truth answers and model-generated responses. However, for open-ended questions, the focus should be on the rationality and correctness of the answers, making the LLM-based metric a more reasonable evaluation method. This metric also reveals trends consistent with the previous tasks, further highlighting the performance gap between open-source and proprietary models.

#### 3.3.1 Analysis and Discussion

As shown in Table 3, most MLLMs exhibit a pronounced gap between perception-based and reasoning-based performance. On simple visual questions—where answers depend on direct feature extraction—these models routinely exceed 75% accuracy. However, on scientific reasoning tasks that demand domain-specific knowledge, their scores drop sharply (e.g. Gemini-2.5-Pro-Thinking: 77.06% perception vs. 56.31% reasoning). This divergence suggests that while robust perception is a necessary foundation, it is not sufficient for Earth-science inference: once a model surpasses the  $\approx 75\%$  perception threshold, fur-

Table 3: Accuracies (%) of different models on multiple-choice questions. The best results are highlighted in bold, with the second-best underlined.

Model	Image-Type			Task Type		Overall ACC
	SINGLE	MULTI	CROSS	REASONING	PERCEPT	
<i>Open-source Models</i>						
LLaVA-onvision-72B	53.55	49.48	47.95	46.58	65.52	51.11
Qwen2.5-VL-7B	47.65	44.07	37.53	40.53	58.47	44.83
Qwen2.5-VL-72B	52.11	50.43	46.30	44.40	70.46	50.65
InternVL2-8B	44.86	43.99	38.36	38.97	58.47	43.64
InternVL2.5-78B	53.23	49.74	44.38	43.17	74.21	50.61
InternVL3-78B	57.53	51.37	45.48	47.00	73.61	53.38
Llama3.2-90B-Vision	45.98	40.46	38.90	38.26	56.97	42.74
DeepSeek-VL2	52.43	49.23	44.66	46.06	62.82	50.07
Qwen3-VL-32B	54.58	50.77	47.95	45.58	72.86	52.12
GLM-4.1V-Thinking	54.18	49.83	38.36	46.29	62.97	50.29
Intern-S1-mini	60.00	57.22	49.04	51.68	75.56	58.12
Intern-S1	67.01	65.62	64.11	61.22	79.61	65.63
Qwen-7B-MSEarth	<u>57.61</u>	<u>52.75</u>	<u>45.20</u>	<u>50.68</u>	<u>64.32</u>	<u>53.95</u>
Intern-S1-mini-MSEarth	<u>65.49</u>	<u>62.97</u>	<u>58.63</u>	<u>58.81</u>	78.56	<u>63.54</u>
<i>Proprietary Models</i>						
Gemini-2.5-Flash	58.33	54.55	53.42	49.98	75.56	56.11
Gemini-2.5-Flash-Thinking	60.64	54.64	53.70	51.35	75.86	57.22
Gemini-2.5-Pro-Thinking	64.78	59.36	55.34	56.31	77.06	61.28
Claude-3.5-Haiku	49.48	47.16	42.47	42.18	64.77	47.59
Claude-3.7-Sonnet	59.52	56.53	57.53	51.68	78.11	58.01
GPT-4o	63.03	55.76	47.67	50.45	<b>81.86</b>	57.97
GPT-5.2	61.99	59.54	54.52	54.98	75.86	59.99
Gemini-3-Flash	<b>70.44</b>	<b>67.70</b>	<b>66.85</b>	<b>65.14</b>	<u>80.51</u>	<b>68.82</b>

ther gains hinge on integrating specialized knowledge and enabling multi-step reasoning. By contrast, models pretrained or fine-tuned on scientific datasets, such as Intern-S1, demonstrate substantially higher accuracy in both perception and reasoning, thereby confirming that general MLLMs lack the requisite Earth-science expertise. Furthermore, we find that further training on the MSEarth training set boosts performance across the board, with the largest relative improvement appearing in reasoning tasks. Taken together, these results underscore the critical role of domain-focused data and architectures in closing the reasoning gap.

Key factors driving low reasoning-task performance include (1) insufficient coverage of Earth science content in general pretraining corpora, (2) the absence of iterative chain-of-thought reasoning modules in standard multimodal fusion backbones, and (3) the scarcity of annotated, multimodal datasets that provide step-by-step scientific rationales.

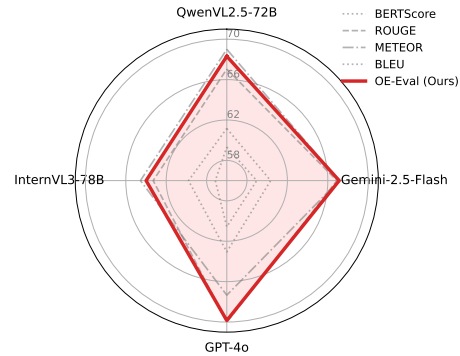


Figure 4: Comparison of Spearman correlations across different models. **OE-Eval** (red line) demonstrates consistently high correlation with human judgments compared to traditional metrics.

### 3.3.2 MLLM-BASED METRICS

Following LAVE (Mañas et al., 2024), in order to assess the validity of OE-Eval, we calculated its correlation with human judgment using Spearman’s rank correlation coefficients. To validate alignment with human judgment, four Earth Science Ph.D. candidates evaluated 160 random samples from the MSEarth Open Ended benchmark

Table 4: **MLLM’s performance on Scientific Open-ended QA and Figure Captioning.** We report ROUGEL (R-L), METEOR (Met), BLEU, BERTSCORE (BS), and the respective LLM-based evaluations (OE-EVAL and CAP-EVAL). The best results are highlighted in bold, with the second-best underlined.

Model	Open-ended QA					Figure Captioning				
	Overlap			Sim.	MLLM	Overlap			Sim.	MLLM
	R-L	MET	BLEU	BS	OE(%)	R-L	MET	BLEU	BS	CAP
<i>Open-source Models</i>										
LLaVA-onevision-72B	37.91	27.88	1.94	89.72	41.56	<b>17.82</b>	18.35	2.15	83.47	2.07
Qwen2.5-VL-7B	35.24	29.00	2.40	88.62	40.68	15.72	16.57	1.55	83.82	2.22
Qwen2.5-VL-72B	38.34	30.65	2.20	89.22	44.82	16.53	21.35	2.36	83.87	2.56
InternVL2.5-8B	35.95	28.25	2.04	89.14	39.05	16.71	19.28	1.58	83.25	1.91
InternVL2.5-78B	40.34	31.15	2.23	90.05	45.64	17.24	20.66	2.27	83.56	2.30
InternVL3-78B	40.42	31.77	2.37	89.98	47.00	16.95	20.95	2.32	83.72	2.43
Llama3.2-90B-Vision	37.53	29.08	1.99	89.16	42.72	12.98	21.21	1.49	78.89	1.82
DeepSeek-VL2	36.37	27.24	1.83	89.38	40.68	16.71	18.22	1.69	83.64	2.22
Intern-S1-mini	37.46	29.72	2.34	89.25	43.69	17.42	24.18	2.26	83.65	2.65
Qwen-7B-MSEarth	38.47	30.10	2.37	89.82	48.19	17.74	26.59	2.86	83.71	2.73
Intern-S1-mini-MSEarth	<u>41.36</u>	31.87	2.44	<b>90.42</b>	49.74	<u>17.94</u>	<u>27.61</u>	3.25	83.86	3.02
<i>Proprietary Models</i>										
Gemini-2.5-Flash	40.02	<u>32.34</u>	2.06	89.39	<u>52.00</u>	17.79	23.01	3.10	<b>83.96</b>	2.98
Gemini-2.5-Flash-Thinking	39.47	30.77	1.87	89.58	46.49	17.42	23.47	2.97	83.85	3.04
Gemini-2.5-Pro-Thinking	38.38	31.73	<u>2.50</u>	88.93	47.70	17.15	27.45	<u>3.33</u>	83.62	<u>3.35</u>
Claude-3.7-Sonnet	40.21	30.75	1.73	89.37	48.33	15.94	21.15	2.15	83.62	2.71
GPT-4o	41.03	32.70	2.04	89.78	48.55	16.33	20.19	2.15	<u>83.93</u>	2.72
Gemini-3-Flash	<b>41.48</b>	<b>32.85</b>	<b>2.46</b>	<u>90.26</u>	<b>52.61</b>	<b>17.96</b>	<b>28.32</b>	<b>3.35</b>	83.92	<b>3.40</b>

Table 5: Accuracy on MSEarth-Bench-mini across Earth science domains. Human expert scores are averages of three Ph.D.-level Earth science evaluators.

Model	Atmospheric	Solid Earth Geophysics	Geography	Ecology	Geology	Hydrology	Oceanography	Polar	All
InternVL3-78B	50.70%	29.73%	28.57%	47.06%	25.00%	51.02%	25.00%	30.00%	47.33%
gemini-2.5-pro	46.48%	58.11%	35.00%	35.71%	47.06%	50.00%	59.18%	50.00%	51.33%
o4-mini	50.00%	45.00%	55.88%	71.43%	49.30%	48.98%	43.33%	62.50%	53.00%
Expert	86.49%	85.00%	85.29%	92.86%	87.32%	85.71%	86.67%	87.50%	87.00%

493 across four MLLMs, as detailed in Appendix K.1.  
 494 With an inter-annotator agreement (Krippendorff’s  
 495  $\alpha$ ) of 0.695, the results in Figure 4 demonstrate that  
 496 OE-Eval achieves higher consistency with human  
 497 judgment than all baselines, effectively capturing  
 498 the scientific nuance required by our benchmark.

#### 499 4 Human Performance Baseline

500 To clarify benchmark difficulty and further justify  
 501 its educational relevance, we additionally report  
 502 human expert scores on the MSEarth-Bench-mini  
 503 set. MSEarth-Bench-mini is constructed by sam-  
 504 pling 300 specialized questions from the original  
 505 MCQ dataset, enabling us to examine whether ex-  
 506 isting methods can improve model performance on  
 507 these domain-focused problems. We hired three  
 508 Ph.D. students with backgrounds in Earth sciences  
 509 to evaluate the tasks, and report their average scores

510 as a human-performance baseline in Table 5. The  
 511 results show that human experts consistently out-  
 512 perform current MLLMs across all Earth science  
 513 domains.

#### 514 5 Conclusion

515 We introduce MSEarth, a graduate-level multi-  
 516 modal dataset designed for MLLMs in geoscientific  
 517 applications. MSEarth not only serves as a robust  
 518 test set but also includes rich training resources  
 519 aimed at enhancing the geoscientific understand-  
 520 ing and reasoning capabilities of existing MLLMs.  
 521 Our evaluation reveals significant gaps in current  
 522 MLLMs’ ability to handle complex, graduate-level  
 523 geoscientific reasoning, highlighting opportunities  
 524 for improvement. We believe MSEarth will serve  
 525 as a valuable resource for advancing MLLMs in  
 526 scientific reasoning.

527  
528  
529  
530  
531  
532  
  
533  
534  
  
535  
536  
537  
538  
539  
  
540  
541  
542  
543  
  
544  
545  
546  
547  
548  
549  
  
550  
551  
552  
553  
554  
555  
  
556  
557  
558  
559  
560  
  
561  
562  
563  
564  
565  
  
566  
567  
568  
569  
570  
571  
  
572  
573  
574  
575  
  
576  
577  
578  
579  
580

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 1:1, 2024.

Lei Bai, Zhongrui Cai, Maosong Cao, Weihao Cao, Chiyu Chen, Haojiong Chen, Kai Chen, Pengcheng Chen, Ying Chen, Yongkang Chen, et al. Intern-s1: A scientific multimodal foundation model. *arXiv preprint arXiv:2508.15763*, 2025a.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025b.

Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pp. 65–72, 2005.

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.

Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024a.

Conghui He, Wei Li, Zhenjiang Jin, Chao Xu, Bin Wang, and Dahua Lin. Opendatalab: Empowering general artificial intelligence with open datasets. *arXiv preprint arXiv:2407.13773*, 2024b.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

Lars Benedikt Kaesberg, Jonas Becker, Jan Philip Wahle, Terry Ruas, and Bela Gipp. Voting or consensus? decision-making in multi-agent debate. *arXiv preprint arXiv:2502.19130*, 2025.

Kartik Kuckreja, Muhammad Sohail Danish, Muza-mmam Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. Geochat: Grounded large vision-language model for remote sensing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27831–27840, 2024.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023.

Xian Yeow Lee, Shunichi Akatsuka, Lasitha Vidyaratne, Aman Kumar, Ahmed Farahat, and Chetan Gupta. Reliable decision-making for multi-agent llm systems. 2025.

Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024a.

Junxian Li, Di Zhang, Xunzhi Wang, Zeying Hao, Jingdi Lei, Qian Tan, Cai Zhou, Wei Liu, Yaotian Yang, Xinrui Xiong, et al. Chemvlm: Exploring the power of multimodal large language models in chemistry area. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 415–423, 2025.

Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv preprint arXiv:2403.00231*, 2024b.

Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungy-oung Ji, Byungju Lee, Xifeng Yan, et al. Mmsci: A multimodal multi-discipline dataset for phd-level scientific comprehension. In *AI for Accelerated Materials Design-Vienna 2024*, 2024c.

Zhenwen Liang, Kehan Guo, Gang Liu, Taicheng Guo, Yujun Zhou, Tianyu Yang, Jiajun Jiao, Renjie Pi, Jipeng Zhang, and Xiangliang Zhang. Scemqa: A scientific college entrance level multimodal question answering benchmark. *arXiv preprint arXiv:2402.05138*, 2024a.

Zijing Liang, Yanjie Xu, Yifan Hong, Penghui Shang, Qi Wang, Qiang Fu, and Ke Liu. A survey of multi-model large language models. In *Proceedings of the 3rd International Conference on Computer, Artificial Intelligence and Control Engineering*, pp. 405–409, 2024b.

636	Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pp. 74–81, 2004.	
637		
638		
639	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36:34892–34916, 2023a.	
640		
641		
642		
643	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. G-eval: Nlg evaluation using gpt-4 with better human alignment. <i>arXiv preprint arXiv:2303.16634</i> , 2023b.	
644		
645		
646		
647	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. <i>Advances in Neural Information Processing Systems</i> , 35:2507–2521, 2022.	
648		
649		
650		
651		
652		
653	Chengqian Ma, Zhanxiang Hua, Alexandra Anderson-Frey, Vikram Iyer, Xin Liu, and Lianhui Qin. Weatherqa: Can multimodal language models reason about severe weather? <i>arXiv preprint arXiv:2406.11217</i> , 2024.	
654		
655		
656		
657		
658	Oscar Mañas, Benno Krojer, and Aishwarya Agrawal. Improving automatic vqa evaluation using large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 4171–4179, 2024.	
659		
660		
661		
662		
663	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pp. 311–318, 2002.	
664		
665		
666		
667		
668	Jonathan Roberts, Kai Han, Neil Houlsby, and Samuel Albanie. Scifibench: Benchmarking large multimodal models for scientific figure interpretation. <i>arXiv preprint arXiv:2405.08807</i> , 2024.	
669		
670		
671		
672	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>arXiv preprint arXiv:2402.03300</i> , 2024.	
673		
674		
675		
676		
677		
678	Noah Siegel, Zachary Horvitz, Roie Levin, Santosh Divvala, and Ali Farhadi. Figureseer: Parsing result-figures in research papers. In <i>Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14</i> , pp. 664–680. Springer, 2016.	
679		
680		
681		
682		
683		
684	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.	
685		
686		
687		
688		
689		
	Bin Wang, Chao Xu, Xiaomeng Zhao, Linke Ouyang, Fan Wu, Zhiyuan Zhao, Rui Xu, Kaiwen Liu, Yuan Qu, Fukai Shang, et al. Mineru: An open-source solution for precise document content extraction. <i>arXiv preprint arXiv:2409.18839</i> , 2024.	690
		691
		692
		693
		694
	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. <i>Advances in neural information processing systems</i> , 33:5776–5788, 2020.	695
		696
		697
		698
		699
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. <i>arXiv preprint arXiv:2412.15115</i> , 2024a.	700
		701
		702
		703
	Joshua C Yang, Damian Dalisan, Marcin Korecki, Carina I Hausladen, and Dirk Helbing. Llm voting: Human choices and ai collective decision-making. In <i>Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society</i> , volume 7, pp. 1696–1708, 2024b.	704
		705
		706
		707
		708
		709
	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 9556–9567, 2024.	710
		711
		712
		713
		714
		715
		716
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> , 2019.	717
		718
		719
		720
	Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. <i>arXiv preprint arXiv:2504.10479</i> , 2025.	721
		722
		723
		724
		725
		726

727  
728  
729  
730  
731  
732  
733  
  
734  
  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
  
745  
  
746  
747  
748  
749  
750  
  
751  
  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773

## A Limitations

The limitation of this work lies in the vastness and diversity of geosciences as a discipline. While we have made efforts to cover a wide range of topics, it is inevitable that certain niche or highly specialized subfields may not be adequately represented in MSEarth.

## B Potential Risks

All papers used were obtained from OpenDataLab (He et al., 2024b) under the CC BY 4.0 license, which permits adaptation and redistribution with attribution. We strictly adhered to all licensing terms and usage requirements specified by OpenDataLab. We have not identified any potential risks or negative societal impacts associated with this work. The dataset is constructed from publicly available scientific literature and is intended solely to advance research in multimodal scientific reasoning.

## C Usage of Language Models

We utilized a large language model (LLM) to aid in the preparation of this manuscript. Its use was limited to editorial tasks, including proofreading for typographical errors, correcting grammar, and improving the clarity and readability of the text.

## D Related Works

**Multimodal Scientific Datasets and Benchmarks.** Numerous multimodal benchmarks have been developed to evaluate scientific understanding across various domains. These benchmarks often integrate text, images, and other modalities to assess models’ reasoning and cross-modal capabilities. However, their creation typically requires significant manual effort in data collection and validation. ScienceQA (Lu et al., 2022) is an early multimodal benchmark that features multiple-choice questions (MCQs) collected from online resources and manually filtered for quality. It covers general science topics such as physics, chemistry, and biology, with a focus on elementary and high school-level reasoning. SceMQA (Liang et al., 2024a) and Mmmu (Liang et al., 2024a) extended this by incorporating both MCQs and open-ended questions (OE) from textbooks and online resources, targeting pre-college and college-level difficulty. OlympiadBench (He et al., 2024a) introduced competition-level problems in mathematics and physics, offering open-ended tasks sourced

from Olympiad exams. These problems are highly challenging but limited to specific domains. More recently, EMMA (Hao et al., 2025) combined manually designed questions with existing benchmarks, covering a broader range of topics with mixed difficulty levels. In contrast, our objective is to enhance models’ ability to comprehend multimodal, complex scientific problems—drawn from high-quality research papers in earth science—that demand graduate-level, domain-specific expertise.

**Paper-Based Multimodal Scientific Datasets and Benchmarks.** Benchmarks based on academic papers aim to leverage the rich, domain-specific content found in scientific literature. FigureSeer (Siegel et al., 2016) first extracts figures from academic papers, focusing on chart figures to evaluate the understanding of chart figures. SciFiBench (Roberts et al., 2024) extended this by introducing figure-to-caption and caption-to-figure matching tasks, while MMSci (Li et al., 2024c) further advanced this approach using figures from Nature papers. However, these benchmarks lack original questions, limiting their ability to assess advanced reasoning and contextual understanding. ArxivQA/Cap (Li et al., 2024b) expanded the scope by generating new questions for figures from 32 subjects on arXiv. However, these questions were generated solely using the inherent capabilities of GPT-4V (Achiam et al., 2023) and did not have contextual support from the relevant text in the papers, raising concerns about their scientific validity. In contrast, our proposed benchmark, MSEarth, addresses these limitations by introducing original, evidence-supported questions grounded in refined captions. This approach enables a rigorous evaluation of MLLMs in professional-level geoscientific applications.

## E Reproducibility

All papers used were obtained from OpenDataLab (He et al., 2024b) under the CC BY 4.0 license, which permits adaptation and redistribution with attribution. We strictly adhered to all licensing terms and usage requirements specified by OpenDataLab. This work establishes a benchmark for evaluating the multimodal Earth scientific exploration capabilities of MLLMs in the field of Earth sciences. It has broader positive impacts, including promoting the responsible use of AI in scientific research and enhancing public understanding of Earth sciences. We believe MSEarth will serve as a valuable re-

774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
  
811  
  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823

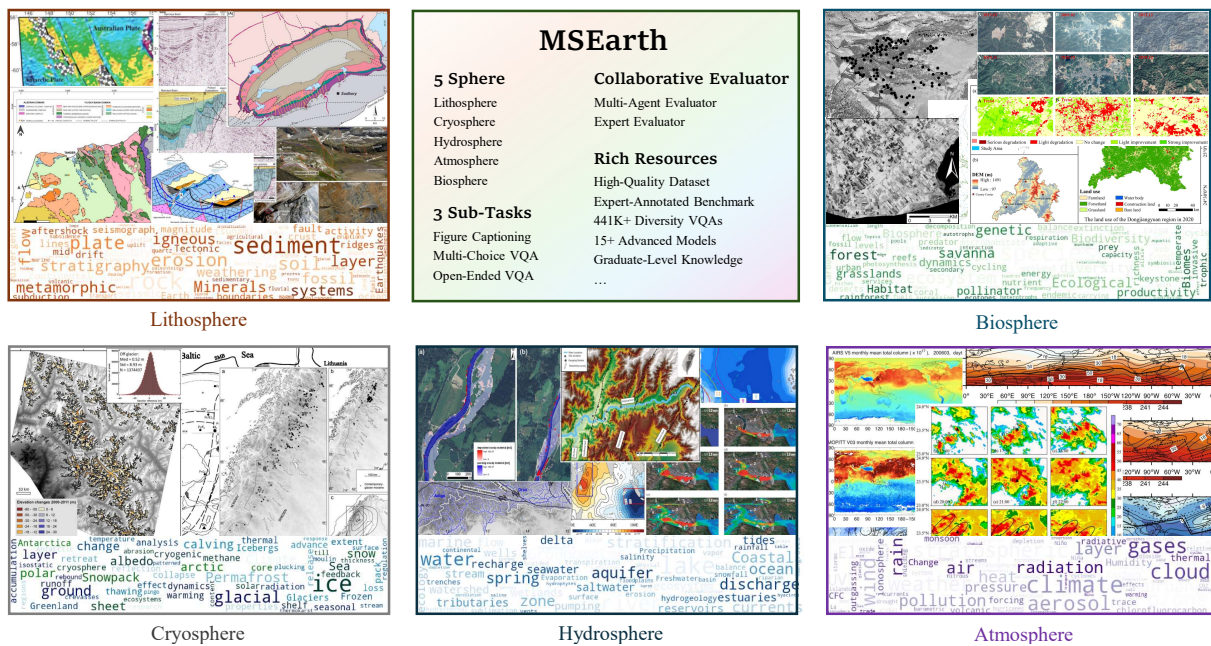


Figure 5: Illustrative examples of the diverse types of scientific figures in MSEarth, sourced from open-access articles available from website.

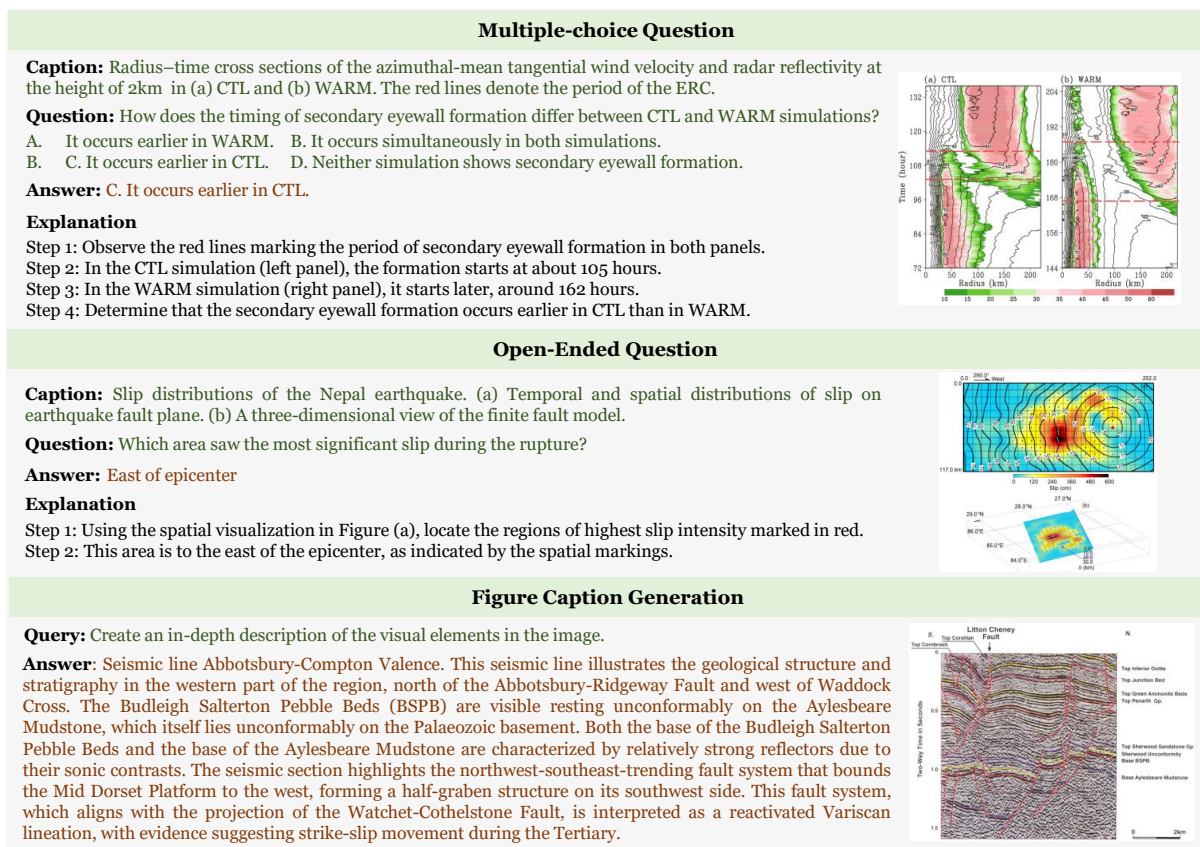


Figure 6: Examples of the three types of scientific question-answering tasks presented in our benchmark.

824 source for advancing multimodal language models  
 825 (MLLMs) in scientific reasoning, and we plan to ex-  
 826 pand its scope to other scientific domains in future  
 827 work. The benchmark is publicly available to foster

further research and innovation in this field. All  
 data in MSEarth are released anonymously, includ-  
 ing the complete dataset on HuggingFace (<https://huggingface.co/MSEarth-Data>) and all data-

828  
 829  
 830  
 831

processing, training, and evaluation code on Anonymous Github (<https://anonymous.4open.science/r/MSEarth-2B3F>)

## F Benchmark Details

### F.1 Field Explanation of MSEarth

In Table 6 and Figure 6, we provide an explanation of each field for the three tasks in MSEarth.

### F.2 Format conversion

Specifically, our data source is a collection of papers gathered by OpenDataLab (He et al., 2024b) from online resources. These papers were processed using MinerU, which converted the textual content of the PDFs into JSON format and saved the images as PNG files. In Figure 7, we present a portion of the content list from a processed PDF paper, highlighting the original caption (raw caption) of the image and the corresponding discussion section. It is evident that the discussion of figures in the paper contains substantial scientific reasoning, which is crucial for a comprehensive understanding of scientific figures.

### F.3 Paper Filtering

To classify scientific papers into relevant Earth system categories, we employ a similarity-based approach using pre-trained sentence embeddings and cosine similarity. The process begins by generating embeddings for both the paper’s title and predefined keywords using the pre-trained all-MiniLM-L6-v2 (Wang et al., 2020) model, which captures the semantic meaning of textual data. First, we calculate the similarity between the paper’s title and a set of general positive keywords, such as “Earth,” “Earth system,” “hydrosphere,” “biosphere,” “lithosphere,” “atmosphere,” and “cryosphere.” The cosine similarity is computed between the embedding of the paper’s title and the embedding of the general positive keywords. If the similarity score is below a predefined threshold (0.2), the paper is excluded from further analysis, as it is deemed irrelevant to the sciences of the Earth system.

To further filter out irrelevant papers, we calculate the similarity between the paper’s title and a set of general negative keywords, such as “cell biology,” “virus,” “pharmaceuticals,” “chemistry,” “physics,” and “astronomy.” If the similarity score exceeds a predefined threshold (0.1), the paper is excluded, as it is likely to belong to unrelated disciplines. For papers that pass the initial filtering,

we calculate their similarity with predefined positive classification keywords for each Earth system category (e.g., hydrosphere, biosphere, lithosphere, atmosphere, cryosphere). Each category contains a list of domain-specific keywords. For example, the hydrosphere category includes keywords such as “water cycle,” “ocean,” “rivers,” “lakes,” and “groundwater,” while the biosphere category includes “ecosystem,” “biodiversity,” “habitat,” and “species.” The cosine similarity is computed between the paper’s title and each keyword within a category, and the average similarity score for each category is calculated.

The category with the highest average similarity score is selected as the most relevant classification for the paper, provided the score exceeds a predefined threshold (0.15). To ensure robustness, we also calculate the similarity between the paper’s title and negative classification keywords for each category. For example, the hydrosphere negative keywords include “chemistry,” “universe,” “planets,” and “astronomy,” while the biosphere negative keywords include “cell biology,” “medicine,” and “pharmacology.” The final relevance score for each category is computed as the difference between the positive and negative similarity scores, ensuring that papers with high relevance to unrelated fields are excluded. The final classification of a paper is determined by passing the general positive and negative keyword thresholds, identifying the category with the highest positive relevance score adjusted by subtracting the negative relevance score, and ensuring the adjusted relevance score exceeds the classification threshold (0.15). This approach allows us to systematically classify papers into Earth system categories while filtering out irrelevant content, leveraging semantic embeddings and cosine similarity to ensure that the classification is both accurate and interpretable.

### F.4 Image Filtering

Next, we further filtered the images in these 103,108 papers. Our goal was to retain Earth observation images, as our task focuses on evaluating the model’s ability to understand and reason about scientific phenomena in Earth sciences. These images include various types of visual data, such as those representing geophysical processes, atmospheric phenomena, geographic features, weather patterns, and cartographic representations.

To achieve this, we employed a systematic filtering pipeline based on the Qwen2-VL-7B-Instruct

Field Name	Input	Description
<b><i>Multiple-choice Question</i></b>		
question_id	✗	The unique identifier for the question.
query	✓	Contains the original caption, question, and options.
response	✗	The correct answer to the question.
images	✓	The file path to the associated image(s).
refined_caption	✗	The enhanced image description based on the paper content.
classification	✗	The classification of the question, including its domain and discipline.
reasoning_chain	✗	The reasoning steps to arrive at the answer.
<b><i>Open-Ended Question</i></b>		
question_id	✗	The unique identifier for the question.
query	✓	Contains the original caption and the question.
response	✗	The correct answer to the question.
images	✓	The file path to the associated image(s).
refined_caption	✗	The enhanced image description based on the paper content.
classification	✗	The classification of the question, including its domain and discipline.
reasoning_chain	✗	The reasoning steps to arrive at the answer.
<b><i>Caption Generation</i></b>		
question_id	✗	The unique identifier for the question.
query	✓	The question.
response	✗	The correct answer (refined caption).
images	✓	The file path to the associated image(s).
context	✗	The text from the paper that describes the image.
original_caption	✗	The original caption of the image.
classification	✗	The classification of the question, including its domain and discipline.

Table 6: **Field Descriptions for Different Tasks.** The table provides details about each field, whether it is used as input, and its description. Fields are grouped by task type: MCQ, OE, and Caption Generation.

931 model. The filtering process was guided by a carefully  
932 designed prompt, which instructed the model  
933 to classify each image as either an Earth obser-  
934 vation image or not. Specifically, the prompt de-  
935 fined Earth observation images as those depicting  
936 remote sensing imagery, atmospheric data visual-  
937 izations, aerial views of geographical features (e.g.,  
938 rivers, urban landscapes), weather-related images  
939 (e.g., precipitation maps, typhoon tracks), and car-  
940 tographic representations. Conversely, the prompt  
941 explicitly excluded images containing biological  
942 entities (e.g., humans, plants, animals), artificial ob-  
943 jects (e.g., vehicles, device structures), data visual-  
944 izations (e.g., statistical charts, line graphs, scatter  
945 plots), text-based content, or blank images.

946 The filtering process was implemented as fol-

947 lows: for each image, the model was provided  
948 with both the image and the prompt, and it gen-  
949 erated a binary output (“1” for Earth observation  
950 images and “0” otherwise). To ensure robustness,  
951 the model’s output was validated through multiple  
952 sampling attempts with slight variations in gener-  
953 ation parameters (e.g., temperature). If the model  
954 consistently classified an image as “1,” it was re-  
955 tained; otherwise, it was discarded. This iterative  
956 and robust classification approach allowed us to  
957 minimize false positives and negatives in the fil-  
958 tering process. After this step, we retained around  
959 83K papers, which contained images classified as  
960 Earth observation images. These filtered images  
961 form the basis for subsequent analysis and eval-  
962 uation of the model’s capabilities in understand-

```

Structure of Content List Field

"content_list": [
{
"type": "text",
"text": "LUMINESCENCE STUDIES ON NEOTECTONIC EVENTS IN SOUTH-CENTRAL KUMAUN HIMALAYA -- A FEASIBILITY STUDY",
"text_level": 1
},
.....
{
"type": "text",
"text": "The continued northward movement of the Indian plate has caused accumulation of stresses which get periodically released. resulting in earthquakes and neotectonic activity along faults. Geophysical and structural studies suggest that seismicity in the Himalaya is related to movements along three major longitudinal thrusts/faults (Fig. 1) viz. the Himalayan Frontal Fault (HFF), the Main Boundary Thrust (MBT) and the Main Central Thrust (MCT) (Valdiya, 1986, 1988; Nakata, 1989). These thrusts divide the Himalaya into three distinct lithotectonic zones. These zones are further dissected by numerous transverse faults (Valdiya, 1976; Khattri and Tyagi, 1983a). A concept of locked segments has been propounded to assess the seismogenic potential of various sectors and is related to accumulation of stresses and their eventual release along the faulted zones (Khattri and Tyagi 1983b). The repeat frequency of this process of locking of stresses and release of energy is not yet well established due to lack of dating methods. Needless to say, this is an aspect cardinal to the estimation of seismogenic hazards involved in planning large scale engineering and societal projects in the Himalaya."
"text_level": 0
},
{
"type": "image",
"img_path": "s3://llm-pipeline-media/pdf-imgs/f5def6d72e7dbfe47ada87d0bc9084997c67f989011db77a255d1817a33fe86b.png",
"img_caption": "FIG. 1. Geological map indicating fault zones and locked segments in Himalaya"
},
.....
]

```

Figure 7: Examples of the content list field in a paper.

Analyze the provided image and classify it as an Earth observation image or not. Earth observation images include, but are not limited to:

- Remote sensing imagery,
- Atmospheric data visualizations,
- Aerial views of geographical features (e.g., rivers, urban landscapes),
- Weather-related images (e.g., precipitation maps, typhoon tracks),
- Cartographic representations.

Exclude images depicting:

- Biological entities (humans, plants, animals),
- Artificial objects (vehicles, device structures),
- Data visualizations (statistical charts, line graphs, scatter plots),
- Text-based content or blank images.

Output format:

- Return "1" if the image is an Earth observation image.
- Return "0" if the image does not meet the Earth observation criteria.

Provide only the numerical output (1 or 0) without any additional text or explanation.

Figure 8: Prompt for retaining Earth observation images.

ing and reasoning about Earth science phenomena. Prompt for retaining Earth observation images are shown in Figure 8.

### F.5 Content Filtering

To construct our benchmark, which requires generating VQA tasks supported by the content of the papers, we ensured that the selected figures not only had captions but were also discussed in detail within the text of the papers. Using regular expressions, we extracted the figure numbers and iden-

tified corresponding discussions in the main body of the papers. Figures with discussions exceeding two sentences were included in the final dataset. Finally, we selected 64,560 papers, resulting in a total of 289,891 figures for further processing.

### F.6 Prompt Designer for MSEarth

The prompt was used to generate a refined caption are shown in Figure 9.

The prompt was used to generate diverse VQAs are shown in Figure 10.

963  
964  
965  
966  
967  
968  
969  
970  
971  
972

973  
974  
975  
976  
977  
978  
979  
980  
981  
982

Category	Type	Keywords
General	Positive	Earth, Earth system, hydrosphere, biosphere, lithosphere, atmosphere, cryosphere
General	Negative	cell biology, virus, pharmaceuticals, chemistry, physics, astronomy, food science, proteins, microbiology
Hydrosphere	Positive	water cycle, ocean, rivers, lakes, groundwater, ice caps, aquifers, precipitation, evaporation, humidity
Hydrosphere	Negative	chemistry, universe, planets, astronomy, astrophysics, space, stars, galaxy, cosmology
Biosphere	Positive	ecosystem, biodiversity, habitat, species, biomes, ecological balance, carbon cycle
Biosphere	Negative	cell biology, chemistry, medicine, pharmacology, microbiology, biochemistry, toxicology, pathology, clinical
Lithosphere	Positive	earthquake, tectonic plates, earth's crust, minerals, rocks, soil, sediments, mountains, volcanoes, landforms, geological processes
Lithosphere	Negative	ancient texts, archaeology, culture, history, artifacts, civilization, prehistoric, mythology, anthropology
Atmosphere	Positive	stratosphere, troposphere, weather, climate, greenhouse gases, ozone layer, air pressure, humidity, winds, carbon dioxide, temperature
Atmosphere	Negative	universe, galaxy, astronomy, astrophysics, space, stars, planets, cosmology, black holes, nebula, solar system
Cryosphere	Positive	glaciers, ice sheets, sea ice, permafrost, snowpack, icebergs, frozen ground, climate change, albedo effect, polar regions
Cryosphere	Negative	frozen food, ice cream, refrigeration, freezing, cold storage, ice cubes, food preservation, chilling, frost

Table 7: Keywords for positive and negative classifications across different Earth system categories. The table includes general keywords as well as specific keywords for hydrosphere, biosphere, lithosphere, atmosphere, and cryosphere.

## G Multi-Agent Voting

### G.1 Prompt

The prompt was used to generate a normal answer for MCQ are shown in Figure 11.

The prompt was used to generate a enhanced answer for MCQ are shown in Figure 12.

### G.2 Example of Different Levels of Questions

Figure 14 illustrates an example of a simple problem in multi-agent voting, while Figure 15 presents an example of a domain-specific problem, and Figure 16 demonstrates an example of a challenging problem. The most notable distinction lies in the varying levels of perceptual ability required by the model: simple and challenging problems primarily differ in the model's ability to perceive and inter-

pret images, whereas domain-specific problems emphasize the model's knowledge in specialized fields. Additionally, the answers to domain-specific questions are often supported by evidence found in the "refined caption" field provided in the paper.

To construct the benchmark datasets, we sampled data from the multi-agent automated filtering process as follows: 900 questions from Phase A, 1800 questions from Phase B, and 300 questions from Phase C were selected to form the multiple-choice question (MCQ) set, while 500 questions from Phase A and 1000 questions from Phase B were selected to form the open-ended question set. All sampled data were subsequently validated by experts to ensure accuracy and quality.

You are an expert assistant in scientific image analysis and caption generation. Your task is to rewrite or generate a new, detailed caption for the provided figure using the original caption and only the sentences or information from the Relevant Content that are directly associated with this figure.

**Please strictly follow these guidelines:**

- Assume the figure does not reference or depend on other figures in the document.
- Exclude any mention of other figures, their content, or references in the caption.
- If subfigures are present, provide specific descriptions for each subfigure accordingly. Otherwise, assume it represents a single figure.
- The new caption must be detailed, precise, and include only the relevant details from the provided content.

**Inputs for caption generation:**

- Original Caption: {caption}
- Relevant Content: {content}

Now write a detailed, high-quality caption for this figure below:

Figure 9: Prompt for generating refined captions.

## H Expert Validation

### H.1 Details

We recruited annotators with a background in Earth sciences and a master’s degree through an annotation company to label the data. The annotated dataset consists of 3,000 MCQs and 1,500 open-ended QAs. We provided the annotators with figures, queries, reasoning chains, and our processed refined captions to assist them in evaluating whether the provided answers were reasonable. For questions where the answers could not be found in the refined captions, the annotators were required to use their own knowledge to determine the correctness of the answers. If they were unable to make a judgment, such questions were discarded to ensure that the filtered dataset contained only accurate and complete questions. The tasks assigned to the annotators are described below:

The evaluation framework categorizes questions based on several criteria. First, the **Image Type of Reasoning Required** distinguishes between questions involving a single image, where the input consists of just one image, and those with a **Single-image focus**, where multiple images are present but the question pertains to one specific image. Additionally, **Multi-image reasoning** questions require comparing or reasoning across multiple images.

Next, the **Type of Scientific Question** is considered. **Perception Questions** are those where

answers can be derived through basic observation, such as identifying position or color. These questions do not have answers in the refined captions and require manual evaluation of their validity. In contrast, **Reasoning Questions** necessitate domain-specific knowledge for answering, and annotators must verify if the answer can be derived from the refined caption field.

The **Completeness of Questions** is another criterion, where questions are classified as **Complete** if all necessary information is provided in the question or image, and **Incomplete** if missing information makes it difficult or impossible to answer.

Finally, the **Correctness of Questions** assesses whether the provided answer is accurate, categorizing them as **Correct** or **Incorrect** based on the accuracy of the answer. After manual screening, 216 invalid entries were identified in the MCQ task, and 89 invalid entries were found in the open-ended task. To evaluate the effectiveness of our multi-agent filtering process, we conducted a statistical analysis of the three phases of data. In Phase A, 59 out of 900 sampled questions were deemed invalid after manual review; in Phase B, 80 out of 1800 questions were invalid; and in Phase C, 77 out of 300 questions were invalid. These results demonstrate the utility of the initial model-based filtering: questions supported by refined captions and correctly answered by most models (Phases A and B) tend to be of higher quality, while questions

You are an advanced AI model specialized in generating high-quality Visual Question Answering (VQA) tasks. Your role is to generate a diverse set of VQA questions, answers, and reasoning chains based on the provided visual input (a figure) and its captions.

**Definitions:**

1. **Figure:** A scientific or illustrative figure provided as the primary visual input. Test-takers will analyze this image to answer the questions.
2. **Caption:** A concise summary describing key aspects of the Figure.
3. **Supplementary:** In-depth information (e.g., summarized expert insight, detailed analysis, or background knowledge) that you can use to assist in designing advanced and meaningful questions. However, test-takers cannot access this information.

**Input Information Provided:**

- **Caption:** {raw caption}
- **Supplementary:** {refined caption}

**Task Instructions:**

**1. Use of Input Sources:**

- Ensure that no question can be answered entirely using Caption without observations.
- **Supplementary Usage:** The correct answers are encouraged to be derived from...

**2. Question Types:**

- **Multiple Choice Questions (MCQs):** At least 2 questions must be of this type, with 4 distinct options (A-D) and one correct answer.
- **Open-Ended Questions:** At least 2 questions must be open-ended, requiring concise and precise answers (no more than 4 words).

**3. Reasoning Chains:**

- For every question, you must include a reasoning chain. The chain explains the logical process by which the correct answer can be determined.
- The reasoning chain must:

**4. Output Structure:**

The output must be written in **JSON format**

**5. Task Guidelines:**

1. Questions that are grounded in the Supplementary context are highly encouraged. These questions should require the test-taker to refer to in-depth knowledge and insights not immediately visible in the Figure or Caption.
2. Avoid referencing the Supplementary in any question and reasoning\_chain (e.g., "According to the Supplementary" or "The Supplementary states").

Provide your response below:

Figure 10: Prompt for generating VQAs.

You are tasked with answering a multiple-choice question about the given input image.

**Instructions:**

1. Carefully analyze the input image and the provided query.
2. Based on the image, select the correct option (e.g., 'A', 'B', 'C') or directly state the correct option content.
3. Provide reasoning explaining how to derive the correct answer.

**Input:**

- **Query:** {query}

**Output Format:**

The output must be written in **JSON format** using the structure below:

```
{
  "answer": "Correct option or short answer",
  "Explanation": "Explaining how to derive correct answer."
}
```

Figure 11: Prompt for generating normal answers for MCQs.

Table 8: Main statistics in MSEarth-Bench.

Statistic	Number
Total questions	7,195
MCQ	2,784
Questions with single images	1,255 (45.1%)
Questions with multiple images	1,529 (54.9%)
* Single-image focus	≈1,164 (41.8%)
* Multi-image relational	≈365 (13.1%)
Reasoning Question	2,117 (76.0%)
Perception Question	667 (24.0%)
Open-Ended	1,411
Questions with single images	679 (45.1%)
Questions with multiple images	832 (54.9%)
* Single-image focus	≈619 (41.8%)
* Multi-image relational	≈113 (13.1%)
Captioning	3,000
Average caption length	37.71
Average refined caption length	137.47

## I Classification of Research Problems in Earth Sciences

Under the framework of the five major spheres, we further categorized the generated research problems into specific academic disciplines according to a standardized classification system. Within the broad category of Earth Sciences, we refined the classification into detailed sub-disciplines or sub-fields. The classification process involves three main steps: first, identifying the primary sphere to which the research problem belongs, selecting from eight major disciplines (referred to as primary

spheres), including Atmospheric Sciences, Ecology and Biosciences, Hydrology, Oceanography, Geology, Geography, Solid Earth Geophysics, and Polar Science. Second, the classification is further refined by selecting the most appropriate sub-discipline or sub-field from a detailed hierarchy. Third, for interdisciplinary problems, the primary classification is clearly stated, and any relevant secondary classifications are noted. This hierarchical approach ensures a systematic and precise categorization of research problems, enabling a deeper understanding of their academic and scientific context.

**Summary of Classification:** The classification system includes a total of 8 first-level disciplines and 66 second-level disciplines. Each research problem is assigned to one of the primary disciplines and further refined into a specific sub-discipline based on its characteristics and context.

### Classification Hierarchy

**Atmospheric Sciences:** Atmospheric Chemistry, Meteorology, Climatology, Hydrometeorology, Paleoclimatology, Atmospheric Physics, Numerical Weather Prediction and Simulation, Atmospheric Remote Sensing.

**Ecology and Biosciences:** Regional Ecology, Population Ecology, Community Ecology,

1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111

You are tasked with answering a multiple-choice question about the given input image.

**Input:**

- **Question:** {question}
- **Refined Caption:** {caption}

**Instructions:**

1. Carefully analyze the input image and its caption.
2. Based on the image and caption, select the correct option (e.g., 'A', 'B', 'C') or directly state the correct option content.

**Output Format:**

The output must be written in **JSON format** using the structure below:

```
{
  "answer": "Correct option or short answer",
  "Explanation": "Explaining how to derive correct answer."
}
```

Figure 12: Prompt for generating enhanced answers for MCQs.

1112	Ecosystem Ecology, Ecological Engineering,	<b>Solid Earth Geophysics:</b> Geodynamics, Seismol-	1139
1113	Restoration Ecology, Landscape Ecology,	ogy, Geomagnetism, Gravimetry, Goelec-	1140
1114	Aquatic Ecology and Limnological Ecology,	tricity, Geothermal Science, Tectonophysics,	1141
1115	Biogeochemistry, Biogeography.	Exploration Geophysics, Computational Geo-	1142
1116	<b>Hydrology:</b> Hydrology, Hydrogeology, Limnol-	physics, Experimental Geophysics, Other Dis-	1143
1117	ogy, River Hydrology and Estuarine Hydrol-	ciplines in Solid Earth Geophysics.	1144
1118	ogy, Groundwater Hydrology, Regional Hy-	<b>Polar Science:</b> Polar Ecology, Polar Oceanogra-	1145
1119	drology, Ecohydrology, Hydrological Physics,	phy, Glaciology, Permafrost Science, Polar	1146
1120	Hydrological Geography, Hydrological Mete-	Climate Science.	1147
1121	orology, Hydrological Measurement, Hydro-		
1122	logical Cartography.	<b>J MLLMs' versions</b>	1148
1123	<b>Oceanography:</b> Ocean Chemistry, Ocean	For open-source models, we use vllm (Kwon et al.,	1149
1124	Physics, Ocean Biology, Ocean Geology, Re-	2023) for local testing; for proprietary models, we	1150
1125	remote Sensing Oceanography, Environmental	conduct tests via API calls. The download paths	1151
1126	Oceanography, Marine Resources Science.	for specific models and the versions of models ac-	1152
1127	<b>Geology:</b> Economic Geology, Engineering Ge-	cessed via API are provided in Figure 10.	1153
1128	ology, Environmental Geology, Quaternary	<b>K Evaluation Metrics</b>	1154
1129	Geology, Sedimentology, Stratigraphy, Pale-	<b>K.1 MLLM-based Metrics</b>	1155
1130	ogeography, Volcanology, Mineralogy and	Following G-Eval (Liu et al., 2023b), we utilize	1156
1131	Petrology, Regional Geology, Remote Sens-	MLLM (Qwen2.5-VL-72B) with a specialized	1157
1132	ing Geology.	prompt to compute a factual scientific score. For	1158
1133	<b>Geography:</b> Physical Geography, Human Geog-	the captioning task, we define a Cap-Eval score	1159
1134	raphy, Regional Geography, Urban Geogra-	ranging from 1 to 5, where higher scores indicate	1160
1135	phy, Tourism Geography, World Geography,	better caption quality. For the open-ended QA task,	1161
1136	Historical Geography, Geomorphology, Bio-	we introduce OE-Eval, which evaluates the reason-	1162
1137	geography, Chemical Geography, Other Dis-		
1138	ciplines in Geography.		

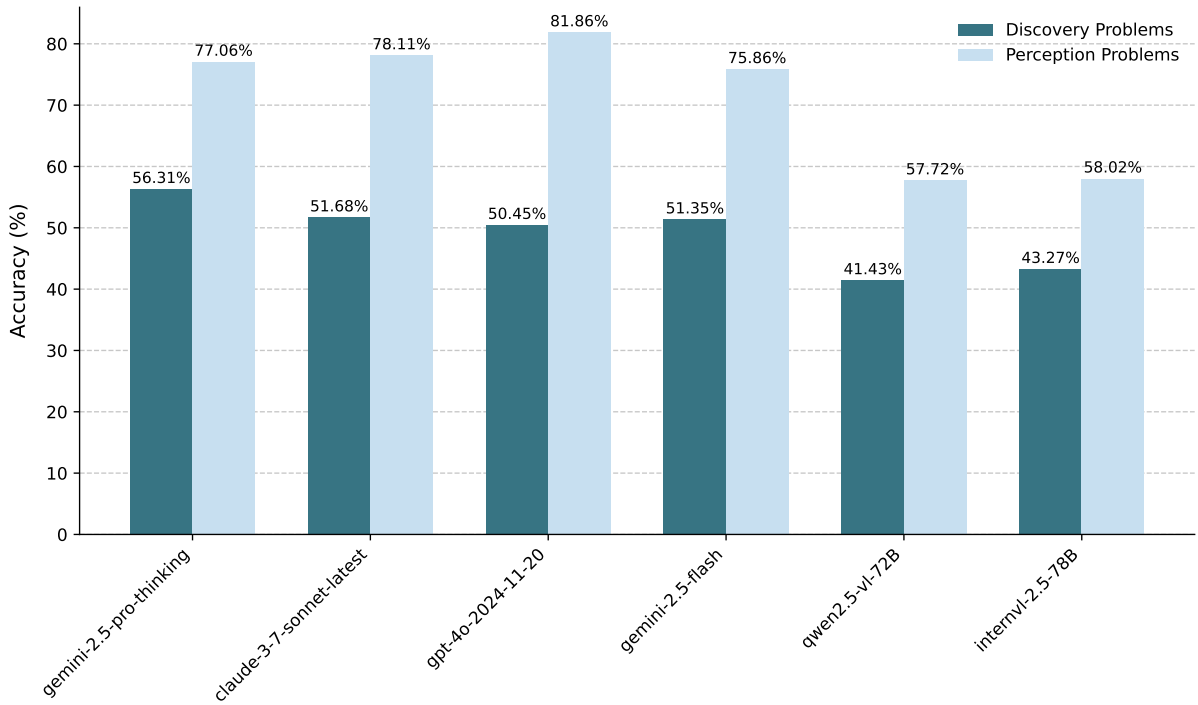


Figure 13: Models' accuracy on reasoning and perception problems.

Table 9: **Top Sub-disciplines in Various Scientific Subjects.** The table lists the top three sub-disciplines by count within each major scientific subject.

Subject	Top 1 Sub-subject		Top 2 Sub-subject		Top 3 Sub-subject	
Hydrology	River Hydrology and Estuarine Hydrology	805	Groundwater Hydrology	790	Limnology	439
Ecology and Biosciences	Aquatic Ecology and Limnological Ecology	562	Landscape Ecology	298	Ecosystem Ecology	280
Geology	Sedimentology	1068	Quaternary Geology	298	Structural Geology	215
Solid Earth Geophysics	Seismology	845	Tectonophysics	343	Exploration Geophysics	74
Geography	Physical Geography	1575	Urban Geography	76	Geomorphology	40
Polar Science	Glaciology	352	Polar Climate Science	46	Permafrost Science	31
Atmospheric Sciences	Meteorology	920	Climatology	619	Atmospheric Remote Sensing	159
Oceanography	Ocean Physics	698	Ocean Geology	163	Environmental Oceanography	104

1163 ableness of generated answers using a binary 0/1  
1164 scoring system.

1165 The prompt was used for Cap-Eval are shown in  
1166 Figure 19.

1167 The prompt was used for OE-Eval are shown in  
1168 Figure 20.

1169 To further establish the correlation between  
1170 LLMs and human judgment specifically in the do-  
1171 main of Earth Science VQA, we conducted a hu-  
1172 man evaluation with four Ph.D. candidates special-  
1173 izing in Earth sciences. They scored a random  
1174 sample of 160 questions from our MSEarth Open  
1175 Ended benchmark. The models evaluated included  
1176 Gemini-2.5-Flash, GPT-4o, InternVL3-78B and  
1177 QwenVL2.5-72B. Our inter-annotator agreement,  
1178 measured by Krippendorff's alpha, is 69.5. Fol-  
1179 lowing LAVE (Mañas et al., 2024), in order to  
1180 assess the validity of OE-Eval, we calculated its  
1181 correlation with human judgment using Spearman's

1182 rank correlation coefficients. We derive a single  
1183 "quality" score from the 4 binary ratings (correc-  
1184 t/incorrect) per answer as follows: 1.0 if at least 3  
1185 annotators rate the answer as correct, 0.5 if only 2  
1186 did so, and 0.0 otherwise.

1187 From the table's results, OE-Eval demonstrates a  
1188 higher consistency with human judgment compared  
1189 to all the considered baselines.

## 1190 K.2 Similarity-based Metrics

1191 In cases where some models fail to strictly fol-  
1192 low instructions and only output the correct an-  
1193 swer, resulting in regular expression matching fail-  
1194 ures, we use the all-MiniLM-L6-v2 model (Wang  
1195 et al., 2020) to calculate the similarity between the  
1196 model's output and each option. The option with  
1197 the highest similarity is then selected as the model's  
1198 answer.

Easy Multiple-choice Question

**Image Caption:**

Coastline totally covered by seawall and concrete armor units on Shimizu coast (July 1995).

**Question:**

How has the natural aesthetic of the coastal environment likely changed due to the structures visible in the figure?

**Options:**

- A. Enhanced beauty due to added greenery
- B. Preservation of the sandy beach's appearance
- C. Loss of natural shoreline aesthetics
- D. Creation of diverse wildlife habitats

**Response:**

C. Loss of natural shoreline aesthetics

**Reasoning Chain:**

Step 1: Use visual cues from the figure to observe the dominance of artificial structures, such as seawalls and concrete armor units.

Step 2: Reflect on the visual impact of manmade elements completely covering the coastline.

Step 3: Identify that no sandy beaches or natural aesthetics remain visible.

Step 4: Determine that "Loss of natural shoreline aesthetics" best describes the impact based on both image and caption.



**Image Caption:**

Hydrogeological cross section in the Grombalia basin

**Question:**

Which layer lies directly below the shallow aquifer in most parts of the Grombalia basin?

**Options:**

- A. Topsoil
- B. Clayey sands aquicluds
- C. Semi-deep aquifer
- D. Sandstone

**Response:**

B. Clayey sands aquicluds

**Reasoning Chain:**

Step 1: The visual information from the figure shows a cross-sectional view of the stratigraphy.

Step 2: Observing the labeling in the legend and the diagram, clayey sands aquicluds are consistently shown directly below the shallow aquifer.

Step 3: Correlating these observations, the layer below the shallow aquifer is identified as the clayey sands aquicluds.

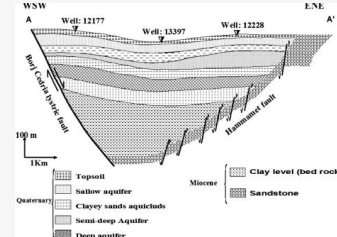


Figure 14: An example of easy multiple-choice VQA.

**L Detailed MSEarth-MCQ Results**

The inputs to our model, MSEarth, consist of images, questions, and the original captions. The original captions provide contextual information about the images, such as the meanings of specific symbols. Therefore, we conducted tests on different models to evaluate their performance with and without the original captions.

For open source models, we performed experiments ino settings: with and without the original caption. The results show that providing the original caption improves performance in all tasks. Notably, the improvement is more significant for perception tasks compared to reasoning tasks. This is likely because perception tasks rely more heavily on understanding the image content, and the original caption provides helpful contextual information for interpreting the image.

We have compiled several case studies to illustrate the necessity of the original caption when

answering questions in certain situations. In example 22, if the original caption is not provided, InternVL3-78B will be unable to accurately determine that the geographical location is in Germany, resulting in an incorrect answer. In contrast, some proprietary models may possess stronger perceptual capabilities and can correctly identify the location as Germany even without the original caption. Similarly, in example 23, providing the original caption aids the model in understanding the image, thereby facilitating task completion. Both scenarios are prevalent in scientific question-answering contexts. To address this, we conducted separate experiments and explicitly integrated these settings into the design of the MSEarth-MCQ task.

**M Results with Compute Scaling**

From the main experiments, it is evident that the performance of various models declines significantly on questions requiring specialized knowl-

1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237

Moderate Multiple-choice Question

**Image Caption:**

Change in sea level pressure between different pairs of sensitivity experiments with the sea surface temperature gradient altered in the Atlantic Ocean, Pacific Ocean, or in combinations in both oceans.

....

Bottom Panel: Difference between experiments AdPi and AiPd. Most changes outside of the range shown ( $-\$2.2\$$  to  $\$2.2\sim\mathrm{mb}\$$ ) are significant at the 95% level.

**Question:**

What mechanism best explains the Southern Hemisphere pressure changes observed in the bottom panel (AdPi-AiPd)?

**Options:**

- A. Pacific tropical warming amplifies wave refraction.
- B. High-latitude cooling destabilizes atmospheric conditions.
- C. Sea ice variations influence atmospheric circulation.
- D. Increased sensible heat transports planetary momentum.

**Response:**

C. Sea ice variations influence atmospheric circulation.

**Reasoning Chain:**

Step 1: Look at the bottom panel (AdPi-AiPd), which demonstrates significant pressure responses in the Southern Hemisphere.

Step 2: Identify pressure anomalies during Southern Hemisphere summer that align with locations of seasonal sea ice.

Step 3: Use context from the caption and figure dynamics to deduce influence from sea ice variations.

**Refined Caption:**

Sea level pressure differences between sensitivity experiments with altered sea surface temperature (SST) gradients in the Atlantic and Pacific Oceans:

.....

•Bottom Panel: The difference between Experiment AdPi (Atlantic high-latitude warming and tropical cooling combined with Pacific tropical warming and high-latitude cooling) and Experiment AiPd (the reverse gradient alterations). This configuration produces a large positive NAO change, with warm tropical Pacific SSTs driving equatorward wave refraction and poleward angular momentum transport, while warm northern North Atlantic SSTs reduce low-altitude northward sensible heat transport and destabilize the local atmosphere. Notable responses are also observed in the Southern Hemisphere, even during summer, which are linked to sea ice variations. Most changes outside the range of  $-\$2.2\$$  to  $\$2.2\sim\mathrm{mb}\$$  are statistically significant at the 95% confidence level.

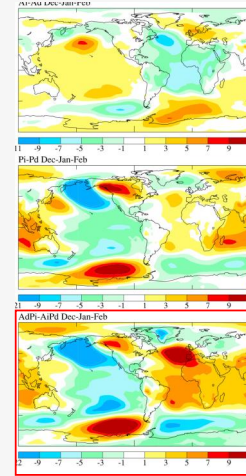


Figure 15: An example of specialized multiple-choice VQA.

1238 edge. To explore whether existing methods can  
 1239 enhance model performance on such questions, we  
 1240 sampled 300 specialized questions from the MCQ  
 1241 dataset to create the MSEarth-mini set. We then  
 1242 evaluated the effectiveness of Chain-of-Thought  
 1243 (CoT) reasoning and majority voting mechanism,  
 1244 which selects the most frequent response among  
 1245 (N) candidate responses; in the case of a tie, one of  
 1246 the most frequent answers is randomly chosen. The  
 1247 results are presented in figure 21. Notably, for the  
 1248 Gemini-Pro-thinking model, which inherently incor-  
 1249 porates a thinking mechanism, introducing CoT  
 1250 reasoning led to a decline in performance. Simi-  
 1251 larly, for some open-source models, such as Qwen  
 1252 and InternVL, the addition of CoT reasoning also  
 1253 resulted in performance degradation. However, the  
 1254 majority voting mechanism proved effective for  
 1255 most models.

1256 **N Impact of Explicit Reasoning**

1257 We assess the impact of explicit chain-of-thought  
 1258 (CoT) prompting on MSEarth-Bench-mini across  
 1259 both open-source and proprietary LVLMs. For

open-source models, we compare InternVL3 and  
 QwenVL2.5; for proprietary models, we examine  
 the Gemini-2.5-Flash series. Within the proprietary  
 family, variants with dedicated “thinking” capa-  
 bilities (e.g., Gemini-2.5-Flash-Thinking) gener-  
 ally outperform counterparts without such capa-  
 bilities (e.g., Gemini-2.5-Flash). In contrast, for  
 open-source models, adding explicit CoT some-  
 times leads to performance declines, which we hy-  
 pothesize stems from limited training for explicit  
 reasoning behaviors (e.g., GRPO-style preference  
 optimization).

To further probe the role of explicit CoT, we  
 include GPT-o4-mini, which exposes configuration  
 options that control reasoning depth (low, medium,  
 high), roughly corresponding to the length of the  
 reasoning chain. Results are shown in Table 12.

Overall, we observe the following:

Models explicitly equipped and trained for  
 “thinking” benefit from enabling CoT (e.g., Gemini-  
 2.5-Flash-Thinking). When a model already ex-  
 hibits strong inherent reasoning, additional explicit  
 CoT can reduce performance, as seen in o4-mini

Hard Multiple-choice Question

**Image Caption:**

MERIS RGB image of scene over Mediterranean Sea from 2 May 2003 (left) and zoom into daily composite of TCWV (middle) and uncertainty of TCWV (right) from MERIS and SSM/I measurements.

**Question:**

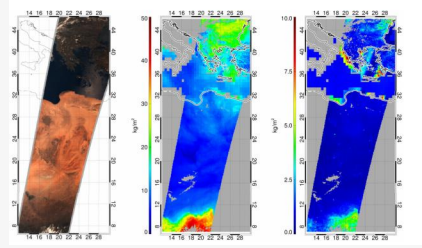
In which regions is TCWV density the highest based on the middle panel?

**Options:**

- A. Over the Mediterranean Sea
- B. Over land near coastal areas
- C. Over desert areas in the bottom left
- D. In the southernmost water regions

**Response:**

D. In the southernmost water regions



**Reasoning Chain:**

- Step 1: Observe the middle panel, which shows TCWV values represented by a color gradient.
- Step 2: The highest values correspond to the red and yellow sections in the southernmost part of the image.
- Step 3: These sections lie over the near-equatorial regions of the southern water zones where higher water vapour levels are observed.

**Refined Caption:**

MERIS RGB image of the Mediterranean Sea area from 2 May 2003 (left), zoomed-in view showing the daily composite of Total Column Water Vapour (TCWV) (middle), and the uncertainty of TCWV measurements (right). The middle panel illustrates the smooth transition of the water vapour field between land and ocean, with increased uncertainty in coastal areas. This uncertainty is linked to the use of MERIS data to fill gaps in the SSM/I measurements. The right panel highlights the regions of elevated uncertainty, particularly along the coast. Additionally, the figure emphasizes MERIS's high sensitivity to small-scale variations in the water vapour field, as seen over Western Turkey.

Figure 16: An example of hard multiple-choice VQA.

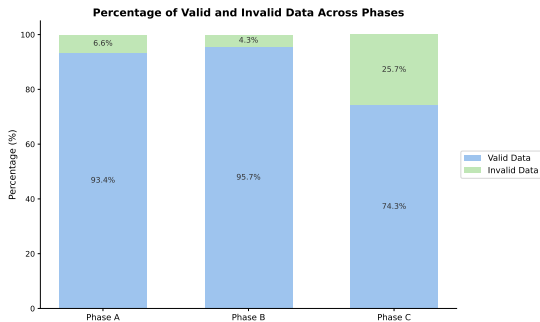


Figure 17: Proportion of valid and invalid data after manual screening across different phases. Phase B, or Specialized VQA, exhibits the highest quality.

at medium/high reasoning depths. Open-source models do not consistently benefit from CoT without targeted training for reasoning behaviors, suggesting a direction for future supervised and RL post-training.

**O More Results**

For MCQ and OE questions, we used radar charts to illustrate the performance of various models across different disciplines. We also give some case studies in Figure 25 and Figure 26. We also present detailed performance breakdowns of all models across every sub-discipline of Earth science in Tables 16.

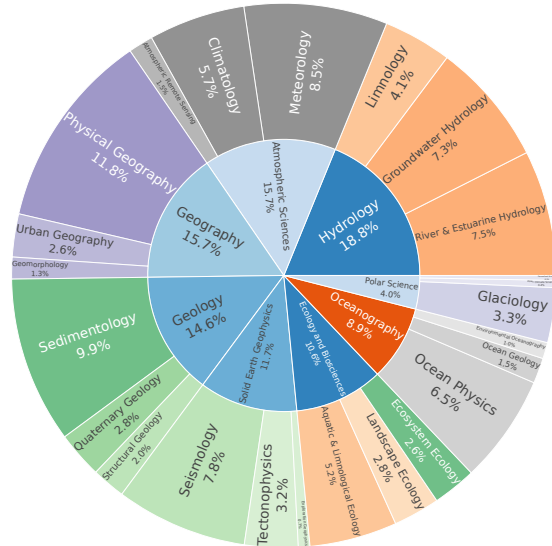


Figure 18: Subjects distribution in MSEarth.

Table 10: Evaluated MLLMs in our experiments with their versions or Huggingface model paths.

<i>Open-source Models</i>	
Model	Model path
Qwen2.5-VL-7B	<a href="https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct">https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct</a>
Qwen2.5-VL-32B	<a href="https://huggingface.co/Qwen/Qwen2.5-VL-32B-Instruct">https://huggingface.co/Qwen/Qwen2.5-VL-32B-Instruct</a>
Qwen2.5-VL-72B	<a href="https://huggingface.co/Qwen/Qwen2.5-VL-72B-Instruct">https://huggingface.co/Qwen/Qwen2.5-VL-72B-Instruct</a>
InternVL2.5-8B	<a href="https://huggingface.co/OpenGVLab/InternVL2_5-8B">https://huggingface.co/OpenGVLab/InternVL2_5-8B</a>
InternVL2.5-78B	<a href="https://huggingface.co/OpenGVLab/InternVL2_5-78B">https://huggingface.co/OpenGVLab/InternVL2_5-78B</a>
InternVL3-78B	<a href="https://huggingface.co/OpenGVLab/InternVL3-78B">https://huggingface.co/OpenGVLab/InternVL3-78B</a>
LLaVA-onvision-72B	<a href="https://huggingface.co/llava-hf/llava-onevision-qwen2-72b-ov-hf">https://huggingface.co/llava-hf/llava-onevision-qwen2-72b-ov-hf</a>
Llama3.2-90B-Vision	<a href="https://huggingface.co/meta-llama/Llama-3.2-90B-Vision">https://huggingface.co/meta-llama/Llama-3.2-90B-Vision</a>
DeepSeek-VL2	<a href="https://huggingface.co/deepseek-ai/deepseek-vl2">https://huggingface.co/deepseek-ai/deepseek-vl2</a>
Intern-S1-mini	<a href="https://huggingface.co/internlm/Intern-S1-mini">https://huggingface.co/internlm/Intern-S1-mini</a>
<i>Proprietary Models</i>	
Model	Model versioning
GPT-4o	GPT-4o-20
Gemini-2.5-Pro-Thinking	gemini-2.5-pro-preview-05-06
Gemini-2.5-Flash	Gemini-2.5-Flash-preview-04-17
Gemini-2.5-Flash-Thinking	Gemini-2.5-Flash-preview-04-17
Claude-3.7-Sonnet	Claude-3.7-Sonnet-20250219
Claude-3.5-Haiku	claude-3-5-haiku-20241022
GPT-4o-mini	GPT-4o-mini-2024-07-18

Table 11: Accuracies (%) of different models on multiple-choice questions. The best results are highlighted in bold, with the second-best underlined. OC: original caption.

Model	Input OC	Image-Type			Task Type		Overall ACC
		SINGLE	MULTI	CROSS	REASONING	PERCEPT	
<i>Open-source Models</i>							
LLaVA-onvision-72B	✗	49.40	45.52	41.10	41.92	61.86	46.69
Qwen2.5-VL-7B	✗	39.12	35.65	39.18	37.27	38.98	37.68
Qwen2.5-VL-32B	✗	42.07	39.78	40.00	37.03	52.92	40.84
Qwen2.5-VL-72B	✗	47.65	43.30	43.84	41.43	57.72	45.33
InternVL2-8B	✗	35.94	34.11	32.33	34.25	36.13	34.70
InternVL2.5-78B	✗	48.13	45.88	45.21	43.27	58.02	46.80
InternVL3-78B	✗	51.95	44.85	45.75	44.54	59.67	48.17
Llama3.2-90B-Vision	✗	44.30	40.64	36.16	38.64	51.42	41.70
DeepSeek-VL2	✗	45.42	42.70	46.85	43.74	46.78	44.47
LLaVA-onvision-72B	✓	53.55	49.48	47.95	46.58	65.52	51.11
Qwen2.5-VL-7B	✓	47.65	44.07	37.53	40.53	58.47	44.83
Qwen2.5-VL-32B	✓	52.59	46.99	43.84	42.47	70.16	49.10
Qwen2.5-VL-72B	✓	52.11	50.43	46.30	44.40	70.46	50.65
InternVL2-8B	✓	44.86	43.99	38.36	38.97	58.47	43.64
InternVL2.5-78B	✓	53.23	49.74	44.38	43.17	74.21	50.61
InternVL3-78B	✓	57.53	51.37	45.48	47.00	73.61	53.38
Llama3.2-90B-Vision	✓	45.98	40.46	38.90	38.26	56.97	42.74
DeepSeek-VL2	✓	52.43	49.23	44.66	46.06	62.82	50.07
<i>Proprietary Models</i>							
Gemini-2.5-Flash	✓	58.33	54.55	<u>53.42</u>	49.98	75.56	56.11
Gemini-2.5-Flash-Thinking	✓	60.64	54.64	53.70	<u>51.35</u>	75.86	57.22
Gemini-2.5-Pro-Thinking	✓	<b>64.78</b>	<b>59.36</b>	55.34	<b>56.31</b>	77.06	<b>61.28</b>
Claude-3.5-Haiku	✓	49.48	47.16	42.47	42.18	64.77	47.59
Claude-3.7-Sonnet	✓	59.52	<u>56.53</u>	<b>57.53</b>	51.68	<u>78.11</u>	<u>58.01</u>
GPT-4o-mini	✓	52.51	48.63	43.01	43.65	<u>68.67</u>	49.64
GPT-4o	✓	<u>63.03</u>	55.76	47.67	50.45	<b>81.86</b>	57.97

Evaluate the quality of a generated caption for a geoscience research paper figure or image.

**Evaluation Criteria:**

1. **Scientific Accuracy:** Does the generated caption accurately describe the scientific content of the figure or image?
2. **Clarity and Coherence:** Is the caption well-structured, logically organized, and easy to understand?
3. **Relevance and Completeness:** Does the caption provide all necessary information to understand the figure or image?

**Evaluation Steps:**

1. Compare the **Generated Caption** to the **Standard Caption**. Assess whether the generated caption aligns with the scientific content and intent of the standard caption.
2. Assign a score for coherence on a scale of 1 to 5, where 1 is the lowest and 5 is the highest, based on the Evaluation Criteria.

**Input:**

- **Standard Caption:** {response}
- **Generated Caption:** {generated\_caption}

**Important Instructions:**

- Only output the score in the specified JSON format.
- Do not provide any explanations, comments, or additional text.

**Output Format:**

The output must be written in **JSON format** using the structure below:

```
{  
  "score": 1-5  
}
```

Figure 19: Prompt for evaluating the quality of generated captions.

Model	CoT (Accuracy %)	Non-CoT (Accuracy %)
Gemini-2.5-Pro	50.67%	52.33%
Gemini-2.5-Flash-no-think	42.00%	40.00%
Gemini-2.5-Flash-Thinking	52.00%	46.00%
o4-mini (low)	52.00%	51.00%
o4-mini (medium)	50.67%	53.00%
o4-mini (high)	50.33%	54.33%

Table 12: Effect of explicit chain-of-thought (CoT) prompting on MSEarth-Bench-mini. Higher is better; values are accuracy (%). For o4-mini, low/medium/high denote shorter-to-longer reasoning traces.

You are tasked with evaluating the correctness of a generated answer to an open-ended question about a given input image.

**Input:**

- **Question:** {query}
- **Refined Caption:** {refined caption}
- **Standard Answer:** {response}
- **Generated Answer:** {generated\_answer}

**Instructions:**

1. Based on the refined caption, question, and standard answer, determine if the generated answer is correct.
2. Only output the determination in the specified JSON format.
3. Do not provide any explanations, comments, or additional text.

**Output Format:**

The output must be written in **JSON format** using the structure below:

```
{
  "is_correct": true or false
}
```

Figure 20: Prompt for evaluating the quality of generated answers to open-ended questions.

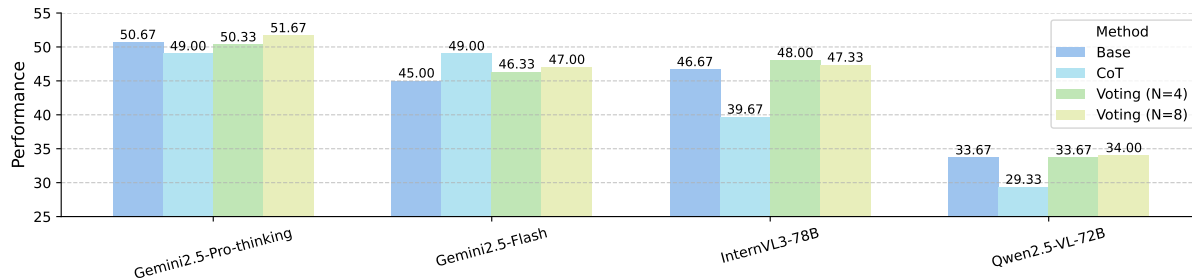


Figure 21: Performance of different strategies on MSEarth-mini.

Model	Atmospheric Sciences			Ecology and Biosciences		
	Meteor.	Climat.	Atmos. RS	Ecosys. Ecol.	Landsc. Ecol.	Aquat. Ecol.
InternVL-8B	0.4039	0.3643	0.5238	0.3833	0.4792	0.3030
InternVL-78B	0.4624	0.3643	0.5476	0.6333	0.5208	0.6061
InternVL3-78B	0.4847	0.3857	0.4762	0.6333	0.5417	0.6364
Qwen2.5-VL-72B	0.4903	0.4071	0.4524	0.4500	0.4375	0.5455
Claude-3.7-Sonnet	0.5599	0.4714	0.6429	0.6500	0.6042	0.5758
Gemini-2.5-Pro-Thinking	0.5877	0.5071	0.5714	0.7000	0.5833	0.6667
GPT-4o	0.5097	0.4429	0.5714	0.5667	0.5208	0.6667
GPT-4o-mini	0.4457	0.3857	0.5000	0.6333	0.4375	0.5455
Gemini-2.5-Flash-Thinking	0.5265	0.3929	0.5476	0.7500	0.5833	0.6364
Gemini-2.5-Flash	0.5153	0.4000	0.5952	0.6167	0.5625	0.5455
Intern-S1	0.6320	0.5870	0.6750	0.7350	0.6580	0.7020
Intern-S1-mini-MSEarth	0.6080	0.5630	0.6420	0.7120	0.6350	0.6780
Intern-S1-mini	0.5850	0.5310	0.6180	0.6870	0.6090	0.6450

Table 13: Model Performance on Primary and Sub-Disciplines of Earth Science (Accuracy)

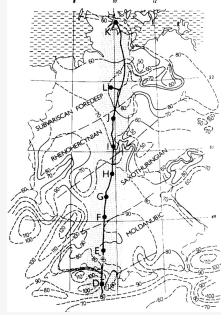
**Raw Caption:** Heat-flow pattern in Germany and in adjacent areas (after Hurlig et al., 1992) together with the position of the Central Segment of the European Geotraverse and shotpoints \$D-K\$

**Question:** Which region along the EGT segment has the most sparse heat-flow data?

- A. North German Lowland
- B. Molasse Basin
- C. Saxothuringian
- D. Rhenohercynian

**Answer:** A

**Refined Caption:** Heat-flow pattern in Germany and adjacent areas, illustrating the position of the Central Segment of the European Geotraverse (EGT). The EGT extends from the southern rim of the Molasse Basin near Chur, northward to the Baltic Sea near Kiel. This section traverses significant geological formations, including the Alpine Molasse Basin, Hercynian (Variscan) units of the Swabian Jura, South German Basin, Mid-German Crystalline Rise, Hessian Depression, and Leine Graben, terminating in the North German Lowland. The heat-flow data, particularly sparse in the North German Basin, range from values north of Hamburg to east of Hannover. The pattern depicted is based on previously obtained temperature data. Shotpoints D-K are marked along the traverse.



**InternVL3-78B w original caption:** A - The North German Lowland region shows fewer heat-flow contour lines compared to other regions, indicating sparser data. The contour lines are less dense, suggesting fewer measurements or less variation recorded, which implies sparse data collection in this area. ✓

**InternVL3-78B w/o original caption:** B - The Molasse Basin region shows fewer heat-flow data points compared to other regions along the EGT segment. The contour lines are less dense, indicating sparser data collection in this area. ✗

**gemini-2.5-flash w/o caption:** A - The image shows heat flow contours (dashed lines) and geological regions. The EGT segment is represented by the thick black line with points D through K. We need to identify the region along the EGT segment with the most sparse heat-flow data, which is indicated by the density of heat flow contours. Examining the map, the northernmost part of the EGT segment, around point K and extending further north, is likely within the North German Lowland. In this area, the heat flow contours are relatively sparse (e.g., 60, 70). In contrast, the Rhenohercynian region along the EGT (from G to J) shows a higher density of contours (60, 70, 80, 90). The Saxothuringian region along the EGT (around I and J) also has several contours (50, 60, 70, 80). The Moldanubic region (D to E) has a high density of contours (80, 90, 100, 110, 120). The Molasse Basin is not shown along the EGT segment. Comparing the density of contours along the EGT in the given options, the North German Lowland (northernmost part) exhibits the most sparse heat-flow data. ✓

**gpt-4o-2024-11-20 w/o caption:** A - The North German Lowland region along the EGT segment has the most sparse heat-flow data, as indicated by the lack of detailed contour lines and data points in this area compared to the other regions. The other regions, such as the Molasse Basin, Saxothuringian, and Rhenohercynian, show more detailed heat-flow contour lines and data points, indicating denser data coverage. ✓

Figure 22: Performance comparison of different models under two settings: with and without the original caption. The results indicate that more powerful models exhibit less dependence on the original caption, highlighting their ability to interpret image content independently.

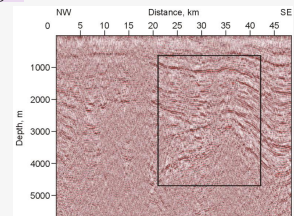
**Raw Caption:** Defining the mud volcano boundary, layers, and faults on the CDS stack migrated section. Black rectangle shows part of the section showing in subsection in Fig. 14 b.

**Question:** What distinguishes the NW section of the Figure from the SE section?

- A. The NW section shows more layered reflectors.
- B. The SE section contains more faults and disrupted layers
- C. The NW section has deeper visible structures
- D. The SE section is closer to the surface.

**Answer:** B

**Refined Caption:** Migrated section of the common diffraction surface (CDS) stacked data, illustrating the boundary, internal layers, and faults of the mud volcanoes. The left portion of the section, which is of particular interest, shows an improved imaging of the mud volcano boundary. Reflectors and other subsurface structures responsible for diffraction events in the stacked section are clearly visible, providing enhanced detail and clarity of the geological features.



**InternVL3-78B w original caption:** B - The SE section within the black rectangle shows more disrupted layers and faults compared to the NW section, which appears to have more continuous and less disturbed reflectors. This indicates that the SE section contains more structural disturbances. ✓

**InternVL3-78B w/o original caption:** A - The NW section of the figure shows more distinct and continuous layered reflectors compared to the SE section, which has more disrupted and less continuous layers. This indicates that the NW section has more pronounced stratification. ✗

**gemini-2.5-flash w/o caption:** A - The image displays a seismic section. The NW section, particularly in the upper part, shows more continuous and distinct layered reflectors compared to the SE section, which appears more complex with folding and potentially faulting within the boxed area. Option B is incorrect as the SE section, while complex, doesn't necessarily show more faults and disrupted layers than the NW section, and the NW section also exhibits some disruption. Option C is incorrect as both sections show structures down to similar depths. Option D is incorrect as the depth scale is the same for both sections, and the SE section is not closer to the surface than the NW section. ✗

**gemini-2.5-pro w/o caption:** B - Upon examining the seismic reflection profile, the SE section (roughly from 20 km to 45 km distance) displays prominent folding of the geological layers. Associated with this folding, there are also clear indications of faults, such as offsets and abrupt terminations of reflectors. This intense deformation results in disrupted layers. In contrast, the NW section (roughly from 0 km to 20 km distance) shows layers that are generally flatter or gently dipping, with less evidence of significant folding or faulting. Therefore, the SE section is distinguished by containing more faults and disrupted layers compared to the NW section. ✓

**gpt-4o-2024-11-20 w/o caption:** B - In the image, the NW section shows relatively continuous and layered reflectors, while the SE section displays more disrupted and faulted layers, as evidenced by the irregular and broken patterns in the SE region. This indicates the presence of faults and structural deformation in the SE section. ✓

Figure 23: Performance comparison of different models under two settings: with and without the original caption.

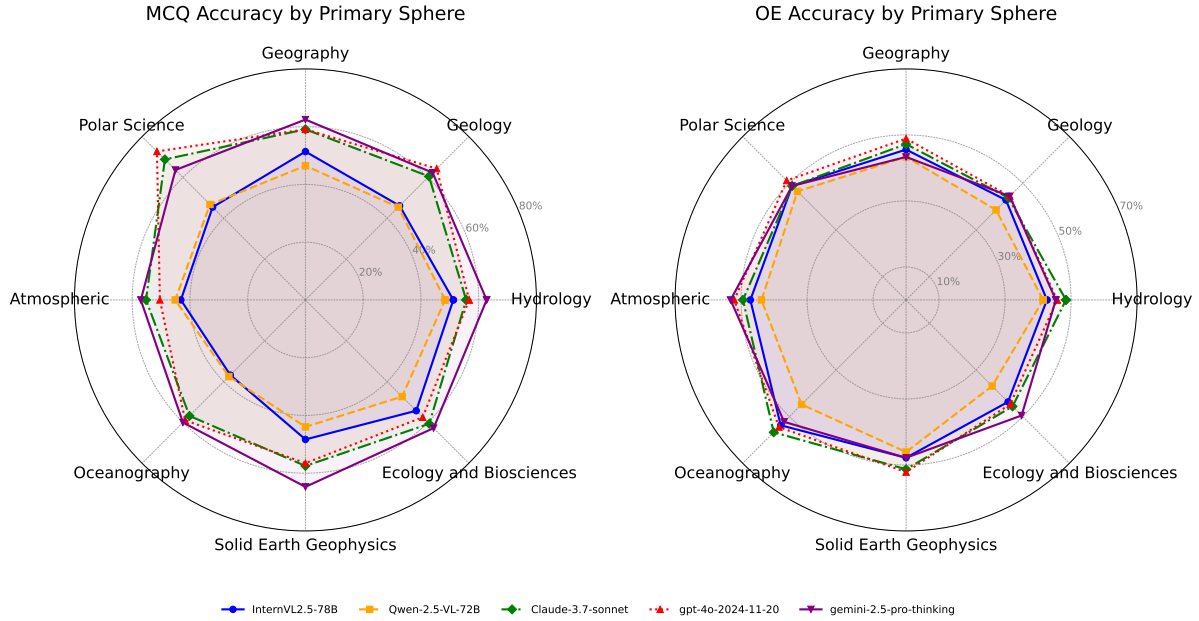


Figure 24: Performance comparison of different models across various subjects.

Model	Geography			Geology		
	Phys. Geog.	Urban Geog.	Reg. Geog.	Sediment.	Struct. Geol.	Quat. Geol.
InternVL-8B	0.4765	0.4444	0.4444	0.5090	0.4268	0.4878
InternVL-78B	0.5451	0.4815	0.5556	0.5663	0.4634	0.5610
InternVL3-78B	0.5884	0.3333	0.6667	0.5986	0.5366	0.6585
Qwen2.5-VL-72B	0.5307	0.3704	0.6667	0.5305	0.5000	0.6098
Claude-3.7-Sonnet	0.5993	0.4074	0.8889	0.6057	0.6341	0.4878
Gemini-2.5-Pro-Thinking	0.6354	0.5556	0.5556	0.6487	0.5854	0.6341
GPT-4o	0.5884	0.5556	0.7778	0.6703	0.6341	0.6829
GPT-4o-mini	0.5090	0.4815	0.6667	0.5269	0.5488	0.5122
Gemini-2.5-Flash-Thinking	0.6426	0.5556	0.7778	0.5986	0.5976	0.5122
Gemini-2.5-Flash	0.5921	0.5185	0.7778	0.5771	0.5366	0.5854
Intern-S1	0.6830	0.6250	0.9120	0.6970	0.6780	0.7250
Intern-S1-mini-MSEarth	0.6590	0.5980	0.8850	0.6730	0.6540	0.6980
Intern-S1-mini	0.6320	0.5640	0.8530	0.6450	0.6210	0.6670

Table 14: Model Performance on Primary and Sub-Disciplines of Earth Science (Accuracy, Continued)

Model	Hydrology			Oceanography		
	River Hydrol.	Groundw. Hydrol.	Limnol.	Ocean Phys.	Ocean Geol.	Env. Oceanogr.
InternVL-8B	0.4550	0.4297	0.4348	0.3800	0.4762	0.2941
InternVL-78B	0.5500	0.4688	0.5652	0.4800	0.5714	0.3529
InternVL3-78B	0.5400	0.4766	0.5652	0.5250	0.6190	0.5882
Qwen2.5-VL-72B	0.5850	0.5078	0.5435	0.4900	0.6667	0.3529
Claude-3.7-Sonnet	0.6150	0.5391	0.5435	0.5650	0.7143	0.3529
Gemini-2.5-Pro-Thinking	0.6700	0.6172	0.6739	0.5750	0.6667	0.5882
GPT-4o	0.6200	0.5625	0.5217	0.5800	0.7619	0.4706
GPT-4o-mini	0.5900	0.4766	0.5435	0.4650	0.5238	0.4706
Gemini-2.5-Flash-Thinking	0.6250	0.5625	0.5435	0.5000	0.6667	0.4706
Gemini-2.5-Flash	0.5700	0.4922	0.6304	0.4800	0.5238	0.3529
Intern-S1	0.7050	0.6680	0.6970	0.6320	0.7950	0.6230
Intern-S1-mini-MSEarth	0.6820	0.6430	0.6720	0.6080	0.7680	0.5950
Intern-S1-mini	0.6570	0.6150	0.6450	0.5830	0.7320	0.5680

Table 15: Model Performance on Primary and Sub-Disciplines of Earth Science (Accuracy, Continued)

Model	Polar Science			Solid Earth Geophysics		
	Glaciol.	Permafrost Sci.	Polar Ocean	Seismol.	Tectonophys.	Geomagn.
InternVL2.5-8B	0.4571	0.7500	0.0000	0.4248	0.5625	0.5455
InternVL2.5-78B	0.4571	0.5000	1.0000	0.4902	0.4375	0.5455
InternVL3-78B	0.5286	0.5000	1.0000	0.5033	0.5000	0.6364
Qwen2.5-VL-72B	0.6286	0.5000	1.0000	0.4706	0.3750	0.7273
Claude-3.7-Sonnet	0.7000	0.5000	1.0000	0.5752	0.5625	0.6364
Gemini-2.5-Pro-Thinking	0.6286	1.0000	1.0000	0.6863	0.5000	0.8182
GPT-4o	0.7286	0.7500	1.0000	0.5163	0.6250	0.9091
GPT-4o-mini	0.5571	0.7500	1.0000	0.4379	0.4375	0.6364
Gemini-2.5-Flash-Thinking	0.6571	0.7500	1.0000	0.5359	0.6250	0.8182
Gemini-2.5-Flash	0.6812	0.7500	1.0000	0.6013	0.5000	0.5455
Intern-S1	0.7650	1.0000	1.0000	0.7230	0.6850	0.9320
Intern-S1-mini-MSEarth	0.7380	0.9500	1.0000	0.6970	0.6580	0.9050
Intern-S1-mini	0.7120	0.9000	1.0000	0.6650	0.6230	0.8780

Note: 1. Sub-disciplines listed are the top 3 with the largest sample size in each primary discipline;

2. Abbreviations: Meteor.=Meteorology, Climat.=Climatology, Atmos. RS=Atmospheric Remote Sensing,

Ecosys. Ecol.=Ecosystem Ecology, Landsc. Ecol.=Landscape Ecology, Aquat. Ecol.=Aquatic & Limnological Ecology,

Phys. Geog.=Physical Geography, Reg. Geog.=Regional Geography, Sediment.=Sedimentology,

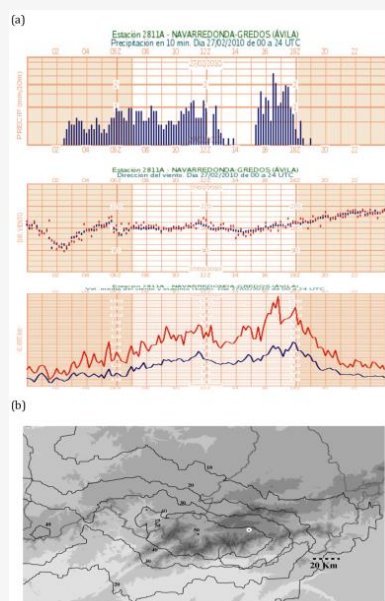
Struct. Geol.=Structural Geology, Quat. Geol.=Quaternary Geology, River Hydrol.=River & Estuarine Hydrology,

Groundw. Hydrol.=Groundwater Hydrology, Limnol.=Limnology, Env. Oceanogr.=Environmental Oceanography,

Glaciol.=Glaciology, Seismol.=Seismology, Tectonophys.=Tectonophysics, Geomagn.=Geomagnetism.

Table 16: Model Performance on Primary and Sub-Disciplines of Earth Science (Accuracy, Continued)

## Multiple Choice VQA Case Study



**Original Caption:** (a) Temporal distribution of precipitation, wind direction, and wind intensity (maximum and mean) on 27 February observed at the weather station at Navarredonda (see Fig. 1b); (b) spatial pattern of precipitation (isolines in mm) observed on 27 February over the study area (circle, Navarredonda automatic weather station).

**Question:** What does the temporal distribution of precipitation in Figure 9(a) indicate about the morning hours of February 27?

**Options:**

- A. Minimal precipitation occurred.
- B. A brief period of heavy precipitation occurred.
- C. Precipitation was evenly distributed throughout the morning.
- D. Precipitation was absent in the morning.

**Response:** B. A brief period of heavy precipitation occurred.

**Refined caption:** (a) Temporal distribution of precipitation, wind direction, and wind intensity (maximum and mean) on 27 February as observed at the Navarredonda weather station. The data highlights a brief period of increased instability indices, Froude number, and meridional moisture flux during the morning and central hours, attributed to the short-lived flow perpendicular to the mountain range. (b) Spatial pattern of precipitation on 27 February over the study area, with isolines indicating precipitation levels in millimeters. The pattern confirms a broad area of precipitation with measured values between 30 mm and higher, with maximum values observed around the Navarredonda automatic weather station.

**Reasoning chain:**

Step 1: Observe the precipitation graph in Figure (a), which shows bar heights representing precipitation recorded in 10-minute intervals.

Step 2: Analyze the morning hours (e.g., 06Z to 09Z); the bar graph shows a significant increase in precipitation during this time.

Step 3: Conclude that there was a brief period of heavy precipitation in the morning.

qwen-2.5-VL-instruct

Answer: A

Explanation: The top graph in Figure (a) shows the temporal distribution of precipitation in 10-minute intervals. In the morning hours (from 00:00 UTC to around 08:00 UTC), the bars representing precipitation are relatively low, indicating minimal precipitation during this time frame. There is no evidence of heavy precipitation or even distribution throughout the morning, nor is there an indication that precipitation was completely absent.

gemini-2.5-pro-thinking

Answer: B

Explanation: Figure (a) shows the temporal distribution of precipitation. The x-axis represents time in UTC, and the y-axis represents precipitation in mm/10min. The morning hours of February 27 can be considered as the period from approximately 06:00 UTC to 12:00 UTC. Observing the precipitation graph (top panel) during these hours: \n1. Precipitation is clearly present throughout this period, so option D ('Precipitation was absent in the morning') is incorrect. \n2. The amount of precipitation is not minimal; there are consistent bars, many of which are around 0.5 mm/10min to 1 mm/10min, and some even higher. Thus, option A ('Minimal precipitation occurred') is incorrect. \n3. The height of the precipitation bars varies significantly throughout the morning, indicating that the precipitation was not evenly distributed. For instance, there are peaks and troughs in intensity. Therefore, option C ('Precipitation was evenly distributed throughout the morning') is incorrect. \n4. Between 06:00 UTC and 07:00 UTC, there is a noticeable increase in precipitation intensity, with several 10-minute intervals showing precipitation rates of 1 mm/10min or more, and one peak reaching approximately 1.5 mm/10min around 06:30 UTC. This can be characterized as a period of heavier precipitation compared to other parts of the morning (e.g., 08:00-10:00 UTC where precipitation is generally lower). This period of higher intensity is relatively brief in the context of the entire morning. Another similar period of increased intensity occurs around 11:00-12:00 UTC. Therefore, option B ('A brief period of heavy precipitation occurred') accurately describes the observed pattern.

claude-3-7-sonnet-thinking

Answer: C. Precipitation was evenly distributed throughout the morning.

Explanation: Looking at Figure (a), the top panel shows precipitation data for February 27, 2010, with the x-axis representing hours from 00 to 24 UTC. During the morning hours (approximately 06-12 UTC), the blue precipitation bars display a relatively consistent pattern with similar heights. There are no major gaps in precipitation or extremely intense brief episodes during this period. Instead, the precipitation appears to be occurring steadily throughout the morning hours with a fairly even distribution, making option C the correct answer.

gpt-4o-2024-11-20

Answer: B. A brief period of heavy precipitation occurred.

Explanation: Figure (a) shows the temporal distribution of precipitation with vertical bars indicating precipitation levels over 10-minute intervals. During the morning hours of February 27 (especially between 06:00 UTC and 08:00 UTC), there is a noticeable spike in precipitation levels, indicating a brief period of heavy precipitation compared to other time intervals.

Figure 25: Case Study of Multiple Choice VQA.

### Open-Ended VQA Case Study 1



**Original Caption:** Schematic illustration of locations of the photometric stations.

**Question:** Which station lies in a semi-arid steppe zone?

**Response:** Baruun-Urt

**Reasoning chain:**

Step 1: From the Figure, Baruun-Urt is outside the Gobi Desert's shaded area, suggesting it is not in the arid region.

Step 2: The semi-arid steppe zone is consistent with areas having denser vegetation, implied by Baruun-Urt's location in the Figure.

**Refined caption:** Schematic illustration of the locations of the photometric stations used for atmospheric AOD measurements in the Eastern Gobi region of Mongolia. The stations are situated at approximately the same latitude and include Sainshand (located in an arid zone), Dalanzadgad (also in an arid zone), and Baruun-Urt (located in a semi-arid steppe zone with denser vegetation cover).

qwen-2.5-VL-instruct

Answer: Dalanzadgad  
Explanation: Located in semi-arid area



gemini-2.5-pro-thinking

Answer: Baruun-Urt  
Explanation: The Gobi Desert is an arid zone. Baruun-Urt is located north of the Gobi Desert, in an area likely to be a semi-arid steppe transition zone.



claude-3-7-sonnet

Answer: Baruun-Urt  
Explanation: Located at northern edge of Gobi Desert region where desert transitions to semi-arid steppe ecological zone.

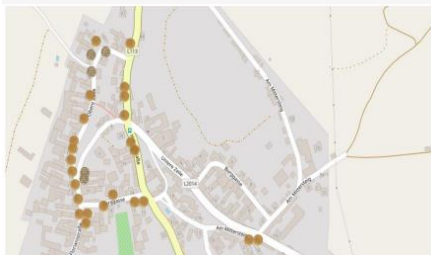


gpt-4o-2024-11-20

Answer: Sainshand,  
Explanation: The map shows Sainshand located within the semi-arid steppe zone of the Gobi Desert.



### Open-Ended VQA Case Study 2



**Original Caption:** Map of Gösing with marked earth buildings from the cadastre.

**Question:** Which natural feature primarily impacts settlement patterns visible in the map?

**Response:** Terrain slope

**Reasoning chain:**

Step 1: Observing the map shows that areas near slopes have higher earth building clusters.

Step 2: Caption reveals that terrain influences settlement and construction choices.

Step 3: Terrain slope emerges as the most significant natural factor influencing patterns.

**Refined caption:** Detailed map of Gösing illustrating the distribution of earth buildings as recorded in the earth building cadastre. The map highlights how the topography of the terrain influences the settlement patterns of these structures, providing insight into the construction techniques and methods that shape the appearance of entire streets within the village.

gemini-2.5-pro-thinking

Answer: Terrain or topography  
Explanation: The map shows buildings (settlements) concentrated along roads, and the presence of dashed lines, likely representing contour lines, suggests varying elevation. Settlements often follow flatter areas or valleys in hilly terrain for easier construction and access. Therefore, the terrain/topography is the primary natural feature influencing where settlements are located.



qwen-2.5-VL-instruct

Answer: River proximity  
Explanation: Settlements cluster near river.



gpt-4o-2024-11-20

Answer: Terrain elevation differences  
Explanation: The map shows settlement patterns concentrated along roads and flatter areas, avoiding steeper or elevated terrain, which is typical in areas with significant elevation differences.



Figure 26: Case Study of Open-Ended VQA.