



Innovator-VL: A Multimodal Large Language Model for Scientific Discovery

Zichen Wen^{1*}, Boxue Yang^{1*}, Shuang Chen¹, Yaojie Zhang¹, Yuhang Han¹, Junlong Ke¹, Cong Wang¹, Yicheng Fu¹, Jiawang Zhao¹, Jiangchao Yao¹, Xi Fang², Zhen Wang², Henxing Cai², Lin Yao², Zhifeng Gao², Yanhui Hong², Nang Yuan², Yixuan Li², Guojiang Zhao², Haoyi Tao², Nan Wang², Han Lyu², Guolin Ke², Ning Liao³, Xiaoxing Wang³, Kai Chen³, Zhiyu Li³, Feiyu Xiong³, Sihan Hu⁴, Kun Chen⁴, Yanfeng Wang¹, Weinan E^{1†}, Linfeng Zhang^{2†}, Linfeng Zhang^{1†}

¹School of Artificial Intelligence, Shanghai Jiao Tong University, ²DP Technology, ³MemTensor,

⁴Institute of Theoretical Physics, Chinese Academy of Sciences

Abstract

We present Innovator-VL, a scientific multimodal large language model (MLLM) designed to advance multimodal understanding and reasoning across diverse scientific domains while still maintaining excellent performance on general vision tasks. Contrary to the prevailing trend of relying on massive domain-specific pretraining data and opaque training pipelines, our work demonstrates that principled training design and transparent methodology can yield strong scientific multimodal intelligence with substantially reduced data requirements. (i) First, we provide a fully transparent and end-to-end reproducible training pipeline for scientific multimodal modeling, covering all stages from data collection and cleaning to preprocessing, supervised fine-tuning, reinforcement learning, and evaluation, together with detailed optimization and hyperparameter recipes. This enables faithful reproduction of our results and facilitates systematic extension and adaptation by the community. (ii) Second, Innovator-VL exhibits remarkable data efficiency, achieving competitive performance on a wide range of scientific tasks using fewer than five million carefully curated scientific training samples, despite not relying on large-scale scientific pretraining. These results highlight that effective scientific multimodal reasoning can be achieved through principled data selection and training strategies rather than indiscriminate data scaling. (iii) Third, Innovator-VL demonstrates strong generalization beyond scientific domains, achieving competitive performance among MLLMs of comparable size on general vision benchmarks, multimodal reasoning benchmarks, and scientific benchmarks. This indicates that scientific alignment can be integrated into a unified multimodal model without compromising general-purpose capabilities. Together, our practices suggest that, in the absence of large-scale scientific data, efficient, reproducible, and high-performing scientific multimodal models can be built, thereby providing a practical and transparent foundation for future research in scientific multimodal modeling.

- 🌐 **Homepage:** InnovatorLM.github.io/Innovator-VL
- 🐙 **Github:** [InnovatorLM/Innovator-VL](https://github.com/InnovatorLM/Innovator-VL)
- 🤗 **Instruct Model:** InnovatorLab/Innovator-VL-8B-Instruct
- 🤗 **Thinking Model:** InnovatorLab/Innovator-VL-8B-Thinking
- 🤗 **Instruct Data:** InnovatorLab/Innovator-VL-Instruct-46M
- 🤗 **RL Data:** InnovatorLab/Innovator-VL-RL-172K

*Equal contribution.

†Corresponding author.

Contents

1	Introduction	3
2	Model Architecture	5
2.1	Vision Encoder	5
2.2	Projector	5
2.3	Language Model	6
3	Pre-training	6
4	Post-training	7
4.1	Supervised Fine-tuning	7
4.1.1	General Multimodal Instruction Data	8
4.1.2	Chain-of-Thought and Multi-step Reasoning Data	8
4.1.3	Scientific Understanding Data	8
4.2	Reinforcement Learning	9
4.2.1	RL Data	9
4.2.2	Optimization Algorithm	10
4.2.3	Reward System	10
5	Infrastructure	11
6	Evaluation	12
6.1	Comparison Models	12
6.2	Benchmarks	12
6.3	Implementation Details	12
6.4	Performance Comparison	13
6.4.1	General Multimodal Understanding	14
6.4.2	Multimodal Reasoning	14
6.4.3	Scientific Knowledge	14
7	On Token Efficiency of Reasoning	14
8	Conclusion and Future Works	15
A	More Results	21
A.1	Qualitative Case Study	21
A.1.1	General Visual Tasks	21
A.1.2	Math & Reasoning Tasks	25
A.1.3	Scientific Tasks	30
B	Evaluation Details	46
B.1	Benchmark Details	46
B.2	Evaluation Prompts	49

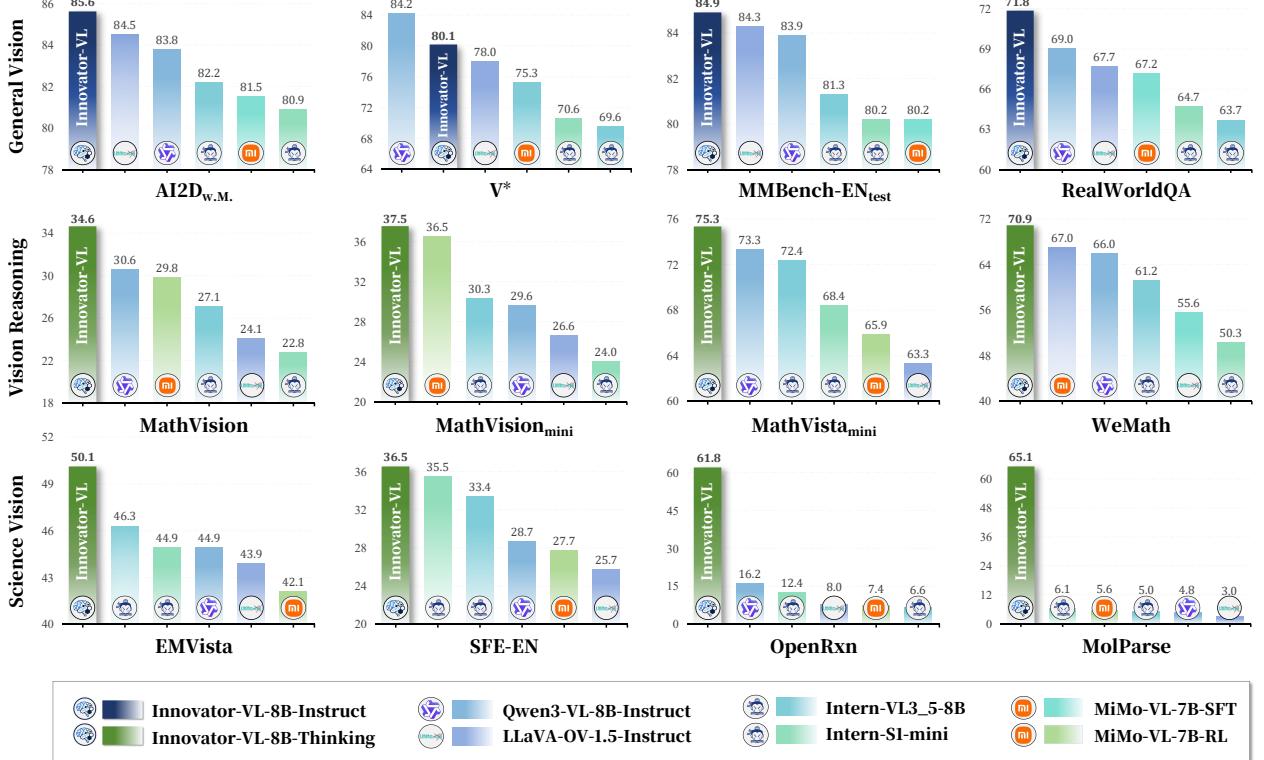


Figure 1 Innovator-VL-8B performance across different benchmarks. The first row displays the results of the Innovator-VL-8B-Instruct model on general benchmarks, while the second and third rows present the Innovator-VL-8B-Thinking model on mathematical reasoning and scientific benchmarks. The yellow dashed line represents the average score of the selected models.

1 Introduction

In recent years, multimodal large language models (MLLMs) have advanced rapidly [Bai et al., 2025a, An et al., 2025, Wang et al., 2024a, Team et al., 2025a,b, Zhu et al., 2025], driven by the scaling of model parameters [Liu et al., 2024a, Team et al., 2025c], the availability of large-scale multimodal data, and improvements in training paradigms. By jointly modeling visual and textual information, these models are able to perceive and reason, enabling a wide range of applications that were previously difficult to address with unimodal systems. As a result, MLLMs have demonstrated strong performance on diverse tasks, including visual question answering, image captioning, visual grounding, document understanding, and multimodal reasoning, highlighting their growing importance as general-purpose intelligent systems [Wen et al., 2025a].

Despite these impressive advances, scientific-domain tasks remain particularly challenging and underexplored [Zhang et al., 2023]. Scientific problems often require precise understanding, rigorous multi-step reasoning, and the ability to integrate domain knowledge across modalities, making them substantially more demanding than general multimodal tasks. Importantly, strong performance on scientific tasks is widely regarded as a crucial milestone toward building artificial general intelligence (AGI) and scientific general intelligence (SGI), as it reflects a model’s capacity for systematic reasoning and knowledge-intensive problem solving. However, existing multimodal large language models, especially open-source models, still struggle to achieve performance in scientific tasks comparable to that observed on general-purpose benchmarks, revealing a significant gap between current capabilities and the requirements of scientific reasoning.

Recognizing this challenge, researchers in both academia and industry have begun to devote increasing attention to scientific multimodal understanding and reasoning, leading to a number of recent advances. Several scientific-oriented multimodal large language models [Bai et al., 2025b, Zhong et al., 2024, Liao et al., 2025a] and benchmarks [Li et al., 2025a, Zhao et al., 2025a, Shen et al., 2024] have been proposed to

enhance performance on tasks involving mathematics, physics, chemistry, and other STEM domains. These efforts demonstrate the potential of multimodal models in scientific problem solving and represent important steps toward bridging the gap between general-purpose perception and scientific domain-specific reasoning. While these efforts have led to encouraging progress, many existing approaches rely on large-scale or highly specialized scientific datasets, making them expensive to reproduce and difficult to extend. In addition, the training pipelines of several systems are either partially opaque or tightly coupled with task-specific heuristics, limiting reproducibility and broader adoption by the community. Moreover, improvements on scientific benchmarks often come at the cost of degraded performance on general multimodal understanding and reasoning tasks, suggesting a lack of balanced generalization.

Motivated by these observations, we present **Innovator-VL**, a scientific multimodal large language model designed to advance multimodal understanding and reasoning across diverse scientific domains. In contrast to prior approaches that depend on massive domain-specific data or opaque training pipelines, Innovator-VL adopts a principled and fully transparent training framework that emphasizes data efficiency, reproducibility, and balanced generalization. Specifically, we use RICE-ViT [Xie et al., 2025] as the visual encoder for Innovator-VL. Scientific understanding in multimodal settings critically depends on the accurate perception of structured visual elements, such as symbols, annotations, and relational components. By capturing fine-grained, region-level semantics, RICE-ViT decomposes scientific images into semantically coherent visual units, enabling faithful alignment between visual entities and scientific concepts. This structured visual representation significantly reduces the burden on downstream language reasoning modules, facilitating more reliable scientific reasoning. Following the design of the QwenVL series [Bai et al., 2025a, Wang et al., 2024a], we adopt PatchMerger as the vision-language projector. This design effectively balances representational capacity and computational efficiency by merging visual patches into a compact yet semantically informative representation. Considering that Qwen3-8B-Base [Yang et al., 2025a] has been extensively pre-trained on a broad and diverse corpus, we adopt it as the language decoder for Innovator-VL.

During **Pre-training**, we focus on robust language-image alignment and high-quality knowledge injection using LLaVA-1.5 558k [Liu et al., 2024b] and LLaVA-OneVision-1.5-Mid-Training-85M [An et al., 2025]. Given that Qwen3-8B-Base has already undergone extensive pretraining and high-quality scientific-domain data is unavailable, we did not further pretrain the language model on scientific texts or include scientific data in the multimodal pretraining stage. In the **Supervised Fine-tuning (SFT)** stage, we prioritize data quality over quantity. We construct high-quality, human-in-the-loop synthetic data pipelines for multiple specialized scientific domains (e.g., in-the-wild OCSR, chemical reaction understanding, and microstructural characterization). Using less than 5 million scientific samples, Innovator-VL achieves competitive performance among models of comparable size. Meanwhile, we incorporate data with chain-of-thought and multi-step reasoning to prepare the model for solving complex scientific problems and the subsequent RL stage. Finally, to bridge the gap between the model’s latent reasoning capabilities and consistent performance, we introduce **Reinforcement Learning (RL)**. We curate a discrepancy-driven RL dataset and employ Group Sequence Policy Optimization (GSPO) [Zheng et al., 2025] to optimize long-horizon reasoning paths, guided by a hierarchical reward system that enforces both structural format and semantic accuracy.

In summary, our work and practice make three main contributions:

- We present **Innovator-VL**, a fully transparent and end-to-end reproducible training recipe for scientific multimodal modeling, facilitating further exploration of scientific multimodal large language models by the research community.
- Innovator-VL demonstrates **remarkable data efficiency**, achieving competitive performance across a wide range of scientific tasks using fewer than **five million** carefully curated high-quality scientific samples, despite not relying on large-scale scientific-domain pretraining.
- Innovator-VL exhibits **strong capability beyond scientific domains**, attaining competitive performance among multimodal large language models of comparable size on general vision benchmarks, multimodal reasoning benchmarks, and scientific benchmarks. This indicates that scientific reasoning enhancement can be effectively incorporated into multimodal models without compromising general-purpose capabilities.

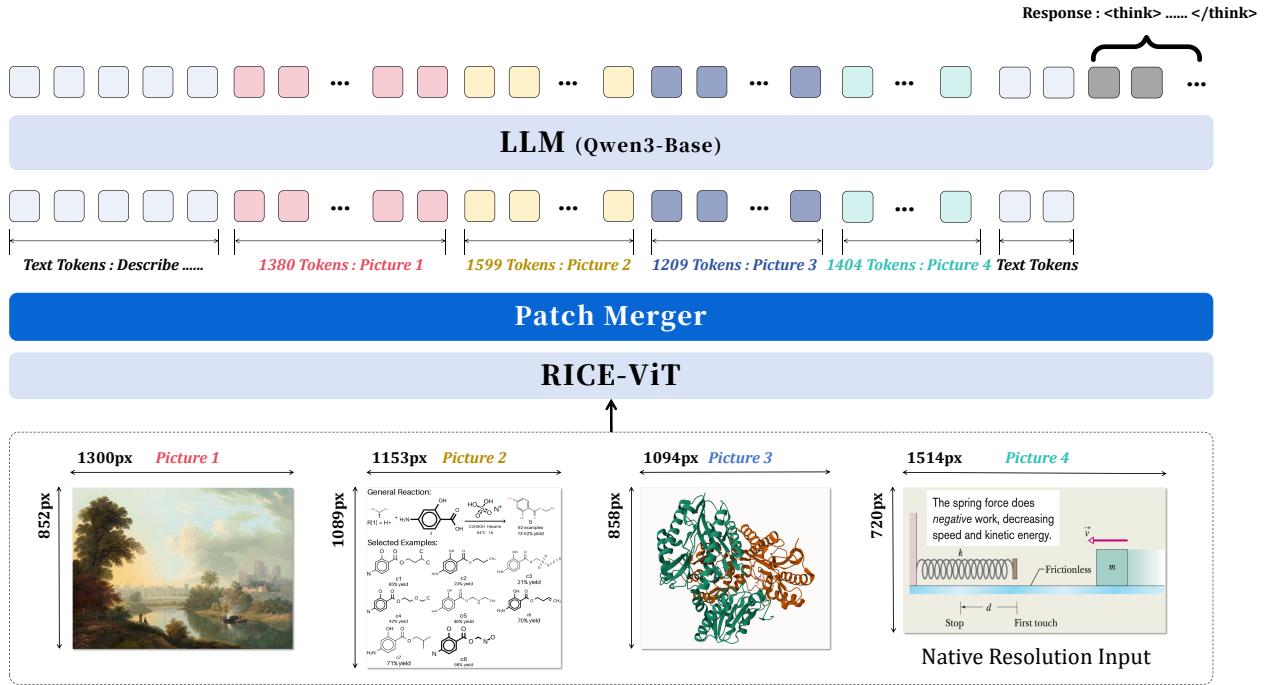


Figure 2 Overall architecture of our native-resolution multi-image reasoning model. Given a text prompt and multiple images with heterogeneous resolutions, RICE-ViT encodes each image at its native size, producing variable-length visual token sequences. **Patch Merger** then aggregates the visual tokens into a compact sequence that is concatenated with text tokens and fed into Qwen-8B-base to generate the final response with explicit reasoning.

2 Model Architecture

Innovator-VL adopts a three-stage architecture following the widely adopted paradigm of *vision encoder-projector-language model*, which has been shown to effectively bridge visual and textual modalities. Our model architecture is illustrated in Figure 2.

2.1 Vision Encoder

We adopt **RICE-ViT** [Xie et al., 2025], a Vision Transformer variant designed for region-aware representation learning. Unlike CLIP [Radford et al., 2021] or SigLIP [Zhai et al., 2023, Tschannen et al., 2025] that primarily focus on global patch interactions, RICE-ViT integrates region-level semantics through a specialized Region Transformer layer, enabling the model to capture both holistic and localized visual cues within a single forward pass. This is particularly important for scientific imagery, where fine-grained structures, dense labels, and spatially localized patterns frequently occur. The region-aware cluster discrimination mechanism used in RICE-ViT facilitates richer semantic embedding by jointly optimizing object and OCR region representations, which improves downstream tasks such as grounding, dense prediction, and multimodal reasoning. Existing studies have shown that RICE-ViT consistently outperforms traditional visual backbones across diverse visual tasks, including segmentation and dense detection. By integrating RICE-ViT as our vision encoder, Innovator-VL benefits from enhanced visual representation capacity that is well aligned with the demands of scientific multimodal understanding and reasoning.

2.2 Projector

To bridge the visual encoder and language model, Innovator-VL employs a learned token compression module, namely **PatchMerger**, following the design used in recent high-performance multimodal language models [Bai et al., 2025a, Wang et al., 2024a]. The PatchMerger module is specifically introduced to address

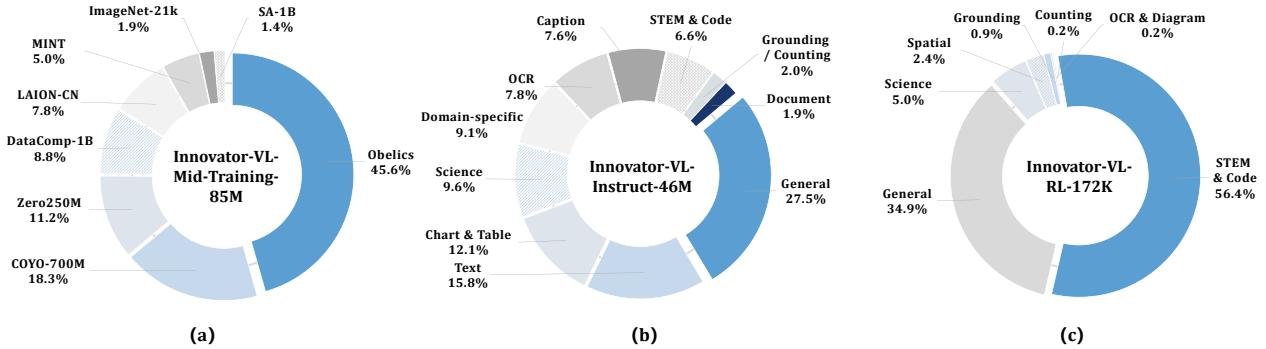


Figure 3 Data distribution across different training stages. (a) Distribution of data sources within the Innovator-VL-Mid-Training-85M dataset. (b) Distribution of data sources within the Innovator-VL-Instruct-46M dataset. (c) Distribution of data sources within the Innovator-VL-RL-172K dataset.

the quadratic computational cost associated with processing dense visual tokens from the Vision Transformer. By learning to merge a larger set of input patch embeddings into a smaller set of representative tokens, PatchMerger significantly reduces the token sequence length that downstream components must process, leading to substantial reductions in computational complexity and memory usage. Empirical studies on token compression [Liao et al., 2025b, Wen et al., 2025b] have demonstrated that such learned merger modules can achieve significant efficiency gains with negligible accuracy degradation on both upstream and downstream tasks [Zhang et al., 2025a, Liu et al., 2025a, Han et al., 2024, Wen et al., 2025c,d, Liu et al., 2025b]. In the context of multimodal integration, this compression facilitates more efficient cross-modal interaction with the language model without compromising the richness of visual representations, which is particularly important when training on limited scientific datasets. By adopting PatchMerger as the projector, Innovator-VL achieves a favorable balance between computational efficiency and representational fidelity, enabling scalable multimodal reasoning and fine-grained understanding.

2.3 Language Model

For the language modeling and reasoning component of Innovator-VL, we adopt **Qwen3-8B-Base** [Yang et al., 2025a]. Pretrained on a large and diverse corpus spanning multiple domains and languages, Qwen3-8B-Base exhibits strong performance in STEM, logical reasoning, and long-context understanding, making it particularly suitable for scientific multimodal tasks that require integrating textual and visual information. Furthermore, Qwen3-8B-Base is open-source and supported by a mature tooling ecosystem, which aligns with our goal of a fully transparent and reproducible scientific multimodal large language model. By integrating Qwen3-8B-Base as the language backbone, Innovator-VL is able to perform advanced scientific reasoning, multimodal question answering, and descriptive generation.

3 Pre-training

Our pre-training consists of two stages.

- **Language-Image Alignment.** We pretrain the projector to align visual features with the word embedding space of large language models (LLMs).
- **High-quality Mid-Training.** We introduce a high-quality knowledge learning (Mid-training) stage to balance computational efficiency with the effective injection of new knowledge into multimodal large language models (MLLMs). In this stage, we transition to full-parameter training across all modules.

Given the scarcity and annotation difficulty of large-scale, high-quality scientific training data, and considering that Qwen3-8B-Base has already undergone extensive pretraining on a comprehensive and diverse corpus, we chose not to continue pretraining the LLM on scientific text data*. We reason that Qwen3-8B-Base

*We have also developed an enhanced internal variant based on a Qwen3-8B-Base that underwent continued pretraining on

Table 1 Overview of training stages, trainable components, datasets, and data scales used to train Innovator-VL, where Innovator-VL-Mid-Training-85M is derived from LLaVA-OneVision-1.5-Mid-Training-85M.

Stage	Training	Dataset	Samples
Language-Image Alignment	Projector	LLaVA-1.5 558k	~558K
High-quality Mid-Training	All	Innovator-VL-Mid-Training-85M	~85M
Supervised Fine-tuning	All	Innovator-VL-Instruct-46M	~46M
Reinforcement Learning	All	Innovator-VL-RL-172K	~172K

has already absorbed substantial scientific knowledge during its pretraining. Without access to high-quality and diverse scientific data, further pretraining on scientific texts could introduce bias and overfitting to some domains. At the same time, our goal is for Innovator-VL to excel in scientific tasks while maintaining its general capabilities, with a robust and reliable foundation in visual understanding serving as the cornerstone for higher-level scientific reasoning.

This philosophy extends naturally to the multimodal pretraining stage. In the first stage, **Language-Image Alignment**, we train our projector using the LLaVA-1.5 558k [Liu et al., 2024b] dataset, aiming to align visual features with the word embedding space of the LLM.

Building upon this alignment stage, we introduce the **High-quality Mid-Training**. In this stage, we transition to full-parameter training across all modules. We employ the **LLaVA-OneVision-1.5-Mid-Training** dataset [An et al., 2025], a large-scale multimodal corpus specifically designed to enhance the foundation ability of multimodal large language models (MLLMs). This dataset comprises approximately **85 million** high-quality image-text pairs, including 65 million English and 20 million Chinese samples. The data are curated from a diverse array of prominent sources, including COYO-700M [Byeon et al., 2022], Obelics [Laurençon et al., 2023], DataComp-1B [Gadre et al., 2023], LAION-CN [Zhang et al., 2022], ImageNet-21K [Russakovsky et al., 2015], SAM-1B [Kirillov et al., 2023], MINT [Wang et al., 2023], and Zero250M [Xie et al., 2023]. A key characteristic of this dataset is the integration of a **feature-based concept-balanced sampling strategy**. To address the challenges of missing annotations or suboptimal caption quality (e.g., the brief and incomplete labels in COYO-700M or caption-free nature of SAM-1B), the dataset construction process moves beyond traditional raw text matching. Instead, it utilizes the pre-trained **MetaCLIP-H/14-Full-CC2.5B** [Xu et al., 2023] encoders to project both images and a taxonomy of 500,000 concepts into a shared embedding space. By retrieving the top- K nearest concept embeddings for each image, refined **pseudo-captions** are generated. This methodology ensures semantic diversity and balance across the massive-scale data, ultimately facilitating a more robust and fine-grained visual-language alignment during pre-training.

4 Post-training

Following multimodal pretraining, Innovator-VL undergoes staged post-training in order to further refine the model’s capabilities, including supervised fine-tuning and reinforcement learning. This stage aims to enhance the model’s ability to follow instructions, perform multi-step reasoning for complex problems, and conduct multimodal reasoning in scientific scenarios.

4.1 Supervised Fine-tuning

To further empower the model to address a wide spectrum of visual tasks with precise and appropriate responses, we continue to perform full-parameter visual instruction fine-tuning on the model. At this stage, we pursue three objectives. First, we endow the model with fundamental visual instruction-following capability, enabling it to leverage the knowledge acquired during pretraining to solve general multimodal understanding tasks. Second, to prepare for the subsequent reinforcement learning stage, we perform a cold start by equipping the model with preliminary chain-of-thought and multi-step reasoning abilities. Third, and most importantly, we strengthen the model’s scientific understanding, training it to effectively

high-quality scientific corpora. Due to ongoing data compliance review procedures, this version is not yet available for public release.

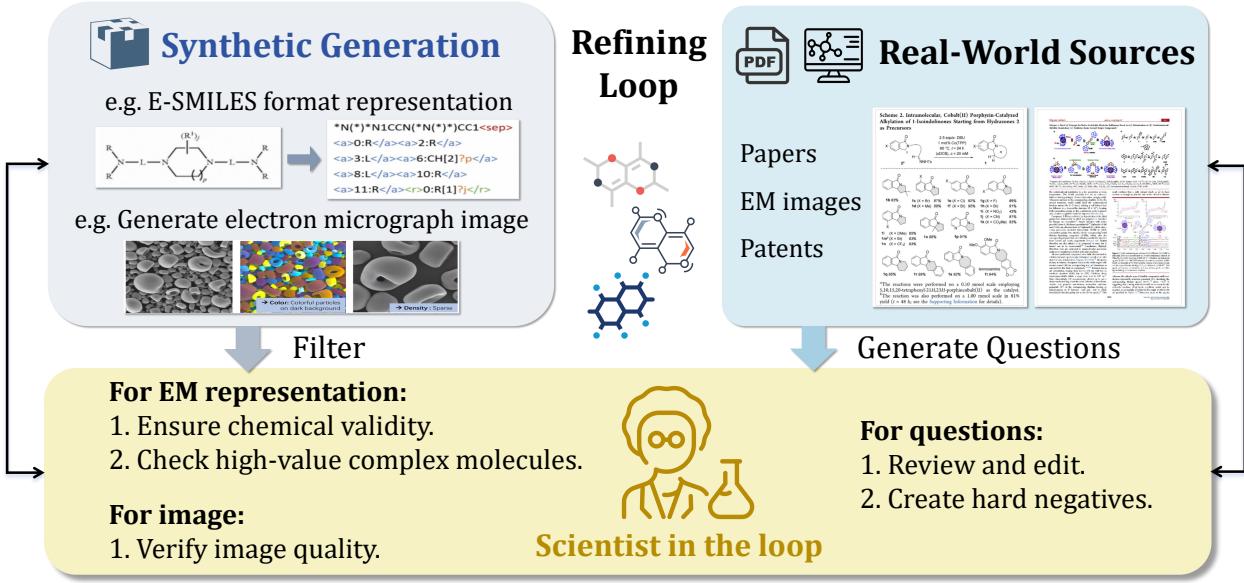


Figure 4 Data Construction Pipeline. The raw data are obtained from both synthetic generation and real-world sources. For each data modality (EM representations, images, and questions), domain experts apply modality-specific inspection and refinement strategies. Through repeated iterative optimization, a final high-quality dataset is produced.

solve problems in scientific domains. To this end, we curate three categories of training data for supervised fine-tuning (SFT), including general multimodal instruction data, chain-of-thought and multi-step reasoning data, and scientific understanding data.

4.1.1 General Multimodal Instruction Data

We employ the LLaVA-OneVision-1.5-Instruct dataset [An et al., 2025], which consolidates a diverse collection of instruction-tuning data from multiple sources. The dataset is carefully curated to ensure balanced coverage across a wide range of categories, including Caption, Chart & Table, Code & Math, General VQA, Grounding & Counting, OCR. This dataset contains approximately 22 million samples. The strong performance of LLaVA-OneVision-1.5 on general vision benchmarks empirically demonstrates the effectiveness of this data.

4.1.2 Chain-of-Thought and Multi-step Reasoning Data

Compared to general multimodal tasks, scientific tasks are inherently more challenging. Therefore, in the visual instruction fine-tuning stage, we not only employ general multimodal instruction data but also introduce data with chain-of-thought and multi-step reasoning annotations. This design enables the model to acquire reasoning patterns at this stage, laying a solid foundation for tackling complex scientific reasoning tasks as well as for the subsequent reinforcement learning stage. Concretely, we adopt a large-scale reasoning-oriented instruction dataset Honey-Data-15M [Zhang et al., 2025b] containing approximately 15 million samples. The dataset spans a wide range of domains and is specifically designed to elicit chain-of-thought and multi-step reasoning behaviors. Different from the original dataset, we further refine the data by removing a large number of explicit think tags. We find that such heavily templated markers may introduce undesirable noise during training; their removal helps preserve the underlying reasoning structure while encouraging more natural and robust reasoning patterns.

4.1.3 Scientific Understanding Data

To enhance the model’s generalizable understanding and cross-domain robustness on scientific tasks, we synthesize high-quality datasets across diverse scientific domains. Our data construction philosophy em-

phasizes principled methodology, systematic quality control, and domain-specific refinement strategies. To illustrate this approach, we present three representative scientific subfields as exemplars: in-the-wild optical chemical structure recognition (OCSR), chemical reaction understanding from scientific literature, and microstructural characterization from electron micrographs (EM). For each subfield, we develop a dedicated data construction and quality control pipeline that systematically integrates data generation, filtering, and consistency verification, demonstrating our generalizable framework for scientific data synthesis. The overall data processing pipeline is illustrated in Figure 4.

Optical chemical structure recognition (in-the-wild OCSR). A low-cost, large-scale paired training corpus is developed via a human-in-the-loop data engine with two stages: (i) synthetic bootstrapping; (ii) active-learning-driven expansion on real-world data. To represent patent-prevalent yet SMILES-inexpressible structures (e.g., Markush patterns and connectivity constraints), E-SMILES is adopted as a unified annotation format, written as SMILES<sep>EXTENSION, where EXTENSION uses compact XML-like tokens to encode R-groups and connectivity information while remaining RDKit-compatible and sequence-model-friendly. The initial model is trained on ~7M synthesized and rendered paired molecular structure images–E-SMILES samples, and subsequently expanded with real patent/paper PDFs through detection, cropping, and deduplication. Confidence is estimated via a 5-fold ensemble and SMILES-similarity scoring, and mid-confidence samples (0.6–0.9) are prioritized for expert correction. The pipeline iterates in a closed loop of “model pre-annotation–human verification–periodic retraining”, enabling scalable acquisition of high-quality real annotations and progressive alignment to real-world distributions.

Chemical reaction understanding from scientific literature. To systematically evaluate multimodal large language models (MLLMs) in realistic PDF-based reaction-reading scenarios, a hierarchical benchmark is curated to cover both single-figure fine-grained perception and full-document cross-modal reasoning. High-impact open-access papers from the past five years are selected to ensure broad reaction diversity and authentic document complexity. Reaction schemes are extracted via PDF layout parsing and automatic cropping (retaining titles and label regions). QA pairs are generated in a human-in-the-loop manner: candidates are proposed by a large model and subsequently reviewed and refined by domain experts to ensure chemical consistency and verifiability. Expert-driven adversarial distractor design further yields chemically plausible, visually dependent hard negatives, while verification mechanisms such as “None of the Above” are incorporated to suppress guess-based hallucinations and improve benchmark robustness.

Microstructural characterization from electron micrographs (EM). The EM dataset is assembled at scale with dense instance-level segmentation annotations and structured textual descriptions to support microstructure segmentation and generation. Multi-source real EM data are aggregated and cleaned via deduplication and quality filtering, while non-structural regions (e.g., watermarks and instrument text) are cropped to obtain a cleaner and more consistent visual distribution. A nine-dimensional, attribute-disentangled description schema is introduced, leveraging multi-model independent extraction with cross-validation followed by expert adjudication for accuracy and interpretability. Instance segmentation annotations are produced through a high-performance annotation platform and a human-in-the-loop iterative workflow, forming a closed loop of “expert seed labeling–model pseudo-label generation–human refinement–multi-stage quality control” to enable scalable construction of high-confidence annotations.

4.2 Reinforcement Learning

To further strengthen the model’s multimodal reasoning and its capacity to tackle complex scientific problems, we apply reinforcement learning (RL) to Innovator-VL-8B-Instruct.

4.2.1 RL Data

We curate a specialized RL training dataset, **Innovator-VL-RL-172K**. Our construction pipeline integrates discrepancy-driven selection, reward-based filtering, and rigorous format standardization.

Discrepancy-Driven Selection and Filtering. We view RL primarily as an *elicitation* mechanism rather than knowledge injection. Consequently, our data selection targets the model’s “effective learnable boundary.” We

measure the divergence between *Pass@N* and *Pass@1* performance across candidate samples; a significant gap indicates that while the model possesses the potential to generate correct solutions (high *Pass@N*), its policy fails to reliably prioritize them (low *Pass@1*). To ensure high training efficiency, we further employ a **reward-based sampling** strategy. We generate multiple response candidates for each instance and filter based on average reward scores, retaining only medium-difficulty instances. This effectively removes trivial samples the model has already mastered and unsolvable cases beyond its current grasp.

Standardization and Composition. Raw data sourced from diverse public datasources [An et al., 2025, Zhang et al., 2025c, Xu et al., 2025, Deng et al., 2025, Qiao et al., 2025, Wang et al., 2025a, Li et al., 2025b] often suffers from inconsistent answer styles and reasoning structures. To mitigate training instability, we standardize all collected samples into a unified reasoning format, normalizing textual structures and ensuring consistent step-wise outputs. The resulting curated corpus contains approximately 172K instances. As illustrated in Figure 3, the dataset is heavily weighted towards complex reasoning tasks. Specifically, **STEM & Code** represents the majority share at 56.4%, followed by **General** multimodal tasks at 34.9%. To ensure comprehensive coverage of fine-grained capabilities, the dataset also includes specific subsets for **Science** (5.0%), **Spatial** reasoning (2.4%), **Grounding** (0.9%), **Counting** (0.2%), and **OCR & Diagram** understanding (0.2%).

4.2.2 Optimization Algorithm

To effectively train Innovator-VL on complex reasoning tasks requiring long Chain-of-Thought (CoT) paths, we employ **Group Sequence Policy Optimization (GSPO)** [Zheng et al., 2025]. While previous group-based methods like GRPO apply token-level importance weights to sequence-level rewards, this mismatch can introduce high-variance noise and instability during training. GSPO addresses this by strictly aligning the optimization unit with the reward unit, performing importance sampling and clipping at the sequence level.

GSPO optimizes the policy π_θ by maximizing the expected reward of generated response groups while constraining the update using a sequence-level trust region. The objective function is defined as:

$$\mathcal{J}_{\text{GSPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, \{y_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|x)} \left[\frac{1}{G} \sum_{i=1}^G \min \left(s_i(\theta) \hat{A}_i, \text{clip}(s_i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right) \right], \quad (1)$$

where G is the group size, and \hat{A}_i is the normalized advantage for the i -th response derived from the reward $r(x, y_i)$. A key innovation of GSPO is the definition of the importance ratio $s_i(\theta)$, which is based on the length-normalized sequence likelihood:

$$s_i(\theta) = \left(\frac{\pi_\theta(y_i|x)}{\pi_{\theta_{\text{old}}}(y_i|x)} \right)^{\frac{1}{|y_i|}}. \quad (2)$$

By utilizing this sequence-level objective, GSPO eliminates the inconsistency of token-level updates in previous algorithms, ensuring robust convergence and superior efficiency when scaling RL to the long-context reasoning tasks present in our dataset.

4.2.3 Reward System

To facilitate deep reasoning capabilities, we explicitly instruct the model to separate its internal thought process from the final conclusion. All RL training instances are formatted with the following instruction:

Instruction: Think and solve the following question step by step. Please put your thinking and analysis procedure within <think>...</think>. Put ONLY your final answer within <answer>...</answer>.

To ensure the generated responses are both structurally compliant and factually correct, we design a comprehensive reward function $R(x, y)$ that integrates format adherence with a hierarchical accuracy verification process. The total reward is calculated as a weighted sum:

$$R(x, y) = \alpha \cdot r_{\text{format}} + (1 - \alpha) \cdot r_{\text{accuracy}}. \quad (3)$$

where we empirically set the weighting factor $\alpha = 0.1$ to prioritize reasoning correctness while maintaining structural discipline.

Format Reward (r_{format}). To facilitate stable Chain-of-Thought (CoT) and multi-step reasoning, we enforce a strict output structure. The model receives a format reward of 1.0 if the response follows the specific XML-like template `<think>...</think>`, `<answer>...</answer>` or employs standard mathematical formatting like `\boxed{}`. Partial credit (e.g., 0.8) is assigned if the response contains mathematical reasoning and a discernible answer but lacks precise formatting tags, while completely unstructured outputs receive a score of 0.0.

Accuracy Reward (r_{accuracy}). Correctness verification is the core of our reward system. Given the diversity of tasks (e.g., multiple-choice, math, open-ended QA), we employ a **hierarchical verification strategy** that balances computational efficiency with semantic understanding. The process proceeds in three cascading stages:

1. **Heuristic & Rule-Based Matching:** We first extract the answer candidate using regular expressions (targeting `<answer>` tags or `\boxed{}`). For Multiple-Choice Questions (MCQs), we utilize an extensive pattern-matching algorithm to handle various output styles (e.g., "(A)", "The answer is A", "Option A"). For general text, we apply relaxed string matching to account for minor variations in punctuation or spacing.
2. **Symbolic Verification:** If exact matching fails for mathematical problems, we employ a symbolic verification tool (using `math_verify`) to parse both the ground truth and the prediction. This allows us to verify mathematical equivalence (e.g., recognizing that $\frac{1}{2}$ is equivalent to 0.5) which simple string matching would miss.
3. **LLM-as-a-Judge:** For open-ended questions or cases where deterministic methods yield a negative result, we utilize a strong external LLM (e.g., Qwen3-VL-235B-A22B-Instruct [Yang et al., 2025a, Chen et al., 2024a]) as a final judge. The judge is prompted with a strict evaluation protocol to output a binary score (0 or 1) based on semantic alignment with the ground truth, ignoring trivial formatting differences.

This cascaded approach ensures that clear, deterministic answers are rewarded cheaply and quickly, while complex semantic matches are accurately captured by the model-based judge.

5 Infrastructure

During the pre-training and supervised fine-tuning (SFT) stages of Innovator-VL, our training pipeline is built upon **AIAK-Training-LLM**[†], a robust distributed training framework derived from Megatron-LM [Shoeybi et al., 2019]. To maximize computational throughput and memory efficiency, we employ a suite of advanced optimization strategies, including the **Transformer Engine** for accelerating mixed-precision computation, distributed optimizer parallelism, and selective activation recomputation to significantly reduce memory overhead. This high-performance training infrastructure serves as the backbone of our pre-training and SFT stages, enabling efficient and scalable learning over massive multimodal datasets.

Following LLaVA-OneVision-1.5 [An et al., 2025], we adopt a data packing strategy to further improve training efficiency. In large-scale multimodal training, substantial variation in sequence lengths leads to significant computational redundancy caused by excessive padding. To address this issue and improve GPU utilization, we employ a carefully designed offline data packing strategy. Specifically, instead of dynamically collating samples at runtime, we pre-optimize the data layout through a parallelized, bucket-based algorithm that groups and concatenates multiple short instances into dense, contiguous sequences. This

[†]An enterprise-grade optimization of Megatron-LM provided by Baidu Cloud.

global preprocessing step enables more effective utilization of the context window while minimizing padding overhead. As a result, our approach substantially improves training throughput and overall computational efficiency.

To overcome the efficiency bottlenecks inherent in synchronous RL training, where the optimization process is frequently blocked by time-consuming inference, we utilize **AReAL** [Fu et al., 2025], a high-throughput asynchronous reinforcement learning framework. By fundamentally decoupling the *Rollout* phase from the *Learning* phase, AReAL transforms the training pipeline into a non-blocking producer-consumer system. This architecture allows dedicated inference workers to continuously produce long-horizon reasoning trajectories, filling the replay buffer without stalling the trainer nodes. Consequently, the learner can perform gradient updates in parallel with data generation, ensuring near-100% GPU utilization and drastically reducing the wall-clock time required for convergence on complex reasoning tasks.

6 Evaluation

6.1 Comparison Models

To evaluate **Innovator-VL**, we compare it against a range of state-of-the-art multimodal models with similar parameter scales (7B–9B), including: **Qwen3-VL-8B** [Bai et al., 2025c], **InternVL3.5-8B** [Zhu et al., 2025], **Intern-S1-mini (9B)** [Bai et al., 2025b], **LLaVA-OneVision (LLaVA-OV) 1.5-8B** [An et al., 2025], **MiMo-VL-7B** [Xiaomi, 2025] (SFT and RL versions), and **MiniCPM-V 4.5 (8B)** [Yao et al., 2024]. These models provide a robust baseline for assessing our model’s performance in general perception, mathematical reasoning, and scientific domains.

6.2 Benchmarks

To comprehensively evaluate our model **Innovator-VL**, we select a set of influential and representative benchmarks and systematically categorize them into three capability dimensions: **general**, **math & reasoning**, and **science** evaluations. Specifically, we use:

- **General Vision:** AI2D [Kembhavi et al., 2016], OCRBench [Liu et al., 2024c], ChartQA [Masry et al., 2022], MMMU(Val) [Yue et al., 2024], MMMU-Pro (Standard) [Wang et al., 2024b], MMStar [Chen et al., 2024b], VStar-Bench [Wu and Xie, 2023], MMBench-EN (Dev/Test) [Liu et al., 2024d], MME-RealWorld (EN/CN) [Zhang et al., 2025d], DocVQA(Val) [Mathew et al., 2021a], InfoVQA(Val) [Mathew et al., 2021b], SEED-Bench [Li et al., 2023], SEED-Bench-2-plus [Li et al., 2024], and RealWorldQA [Dhakal, 2023].
- **Math & Reasoning:** MathVision (Test/Mini) [Wang et al., 2024c], MathVerse (Mini) [Zhang et al., 2024a], MathVista (Mini) [Lu et al., 2024], and WeMath [Qiao et al., 2025].
- **Science:** ScienceQA [Lu et al., 2022], RxnBench (EN/ZH) [Li et al., 2025a], MolParse[†], OpenRxn[§], EMVista[¶], SuperChem (EN/CN) [Zhao et al., 2025a], SmollInstruct [Yu et al., 2024], ProteinLM-Bench [Shen et al., 2024], SFE (EN/ZH) [Zhou et al., 2025], MicroVQA [Burgess et al., 2025], MSEarth-MCQ [Zhao et al., 2025b], and XLRS-Bench-lite [Wang et al., 2025b].

These benchmarks are primarily **multimodal** in nature, enabling an objective and thorough comparison between **Innovator-VL** and state-of-the-art models across diverse scenarios and capability aspects. More details about these datasets can be found in the Appendix B.1.

6.3 Implementation Details

We run all evaluations using the `lmms-eval` framework [Zhang et al., 2024b]. For Innovator-VL-8B-Instruct and Innovator-VL-8B-Thinking, we use deterministic decoding (temperature 0.0, top-p 1.0) across all benchmarks. Additional evaluation details are provided in Appendix B.

[†]<https://huggingface.co/datasets/InnovatorLab/MolParse>

[§]<https://huggingface.co/datasets/InnovatorLab/OpenRxn>

[¶]<https://huggingface.co/datasets/InnovatorLab/EMVista>

Table 2 Comparison with state-of-the-art multimodal models across general, mathematical, and scientific benchmarks. The **best** and second-best results are highlighted in bold and underline, respectively.

Benchmark	Innovator-VL		Qwen3-VL	InternVL3.5	Intern-S1	LLaVA-OV	MiMo-VL	MiniCPM-V	
	8B-Instruct	8B-Thinking	8B	8B	mini (9B)	1.5-8B	7B-SFT	7B-RL	4.5 (8B)
General									
AI2D	85.56	<u>85.07</u>	83.78	82.19	80.86	84.49	81.54	81.51	83.42
AI2D (no mask)	94.46	<u>94.11</u>	93.04	93.10	92.36	93.91	92.03	92.20	93.88
OCR Bench	80.00	79.30	84.90	80.00	80.20	<u>82.50</u>	76.90	78.00	81.50
ChartQA	86.80	86.24	85.12	<u>86.88</u>	87.12	86.64	71.16	71.52	72.84
MMMU(Val)	55.22	52.22	51.44	<u>56.89</u>	57.00	55.44	48.11	49.44	48.44
MMMU-Pro (Standard)	37.40	36.76	36.47	<u>40.17</u>	41.50	37.11	35.49	35.43	36.19
MMStar	64.96	66.15	60.45	<u>66.16</u>	63.50	68.01	58.24	59.25	50.73
VStar-Bench	80.10	<u>81.68</u>	84.29	69.63	70.68	78.01	75.39	75.92	73.30
MBBench-EN (Dev)	85.31	84.88	84.79	83.16	82.99	<u>85.48</u>	80.58	81.01	86.34
MBBench-EN (Test)	<u>84.87</u>	85.82	83.91	81.28	80.21	84.30	80.16	79.65	84.19
MME-RealWorld (EN)	63.70	63.04	<u>63.29</u>	62.54	56.39	61.57	57.47	59.55	60.94
MME-RealWorld (CN)	57.99	<u>62.26</u>	64.53	60.57	52.75	54.64	44.18	46.68	55.37
DocVQA(Val)	94.94	94.56	<u>95.64</u>	91.20	92.59	97.85	86.51	87.05	94.94
InfoVQA(Val)	79.37	78.11	<u>83.13</u>	74.71	75.16	79.03	<u>83.04</u>	82.92	72.14
SEED-Bench	73.84	72.61	75.06	<u>74.41</u>	73.90	72.79	73.06	72.70	74.00
SEED-Bench-2-Plus	70.44	69.83	71.19	<u>71.15</u>	70.44	68.99	69.35	69.39	69.26
RealWorldQA	71.50	<u>70.98</u>	69.02	64.71	63.66	67.71	67.19	66.01	69.28
Average	74.50	74.33	74.71	72.87	71.84	74.03	69.44	69.90	70.99
Math & Reasoning									
MathVision (Test)	31.32	34.64	30.56	27.11	22.76	24.11	<u>32.99</u>	29.77	23.72
MathVision (Mini)	28.95	37.50	29.61	30.26	24.01	26.64	<u>33.88</u>	<u>36.51</u>	18.75
MathVerse (Mini)	54.77	58.73	<u>57.28</u>	48.38	39.97	43.88	54.14	50.28	42.92
MathVista (Mini)	<u>74.30</u>	75.30	73.30	72.40	68.40	63.30	65.40	65.90	61.10
WeMath	65.00	<u>70.86</u>	66.03	61.21	50.34	55.63	73.79	67.01	53.97
Average	50.87	55.41	51.36	47.87	41.10	42.71	<u>52.04</u>	49.89	40.09
Science									
ScienceQA	91.49	91.65	92.95	96.09	<u>96.04</u>	94.86	88.89	89.15	95.07
RxnBench (EN)	82.36	79.80	83.41	78.43	80.79	81.64	80.00	79.61	<u>82.49</u>
RxnBench (ZH)	82.75	81.11	83.61	80.26	82.16	81.70	80.26	80.59	<u>83.34</u>
MolParse	<u>64.90</u>	65.10	4.75	5.00	6.10	2.95	4.55	5.55	3.80
OpenRxn	<u>57.05</u>	61.75	16.15	6.55	12.35	8.00	8.15	7.35	9.15
EMVista	45.95	50.05	44.86	46.32	44.92	43.89	39.62	42.05	<u>49.46</u>
SuperChem (EN)	12.00	11.40	10.20	<u>12.60</u>	13.20	8.60	10.40	8.00	11.40
SuperChem (CN)	10.20	9.40	10.40	<u>11.40</u>	13.60	7.60	10.00	8.60	<u>11.40</u>
SmolInstruct	22.72	24.89	20.44	<u>62.41</u>	63.54	12.21	20.27	18.63	11.99
ProteinLMBench	62.39	60.06	59.32	<u>62.82</u>	64.41	61.55	58.58	58.90	61.97
SFE (EN)	<u>35.90</u>	36.49	28.71	33.39	35.50	25.66	29.37	27.71	30.30
SFE (ZH)	<u>35.69</u>	<u>34.58</u>	28.85	31.47	33.45	24.98	27.18	26.68	27.77
MicroVQA	<u>51.06</u>	40.79	49.23	50.67	50.00	48.17	48.75	51.82	29.65
MSEarth-MCQ	52.59	52.87	48.56	<u>56.50</u>	56.93	54.81	46.23	46.66	50.50
XLRS-Bench-lite	44.83	46.98	49.98	<u>50.50</u>	51.12	44.00	44.40	45.04	47.38
Average	50.13	<u>49.79</u>	42.09	45.63	46.94	40.04	39.78	39.76	40.38
Overall	<u>61.42</u>	61.83	58.33	58.45	57.59	56.02	55.06	54.97	54.40

6.4 Performance Comparison

As summarized in Table 2, we evaluate **Innovator-VL** across 37 diverse benchmarks spanning general, mathematical, and scientific domains. Overall, **Innovator-VL-8B-Thinking** achieves a state-of-the-art (SOTA) average score of **61.83%**, surpassing all compared models of similar scale. These results highlight the remarkable performance of our model, which delivers excellent general multimodal capabilities while achieving significant breakthroughs in specialized scientific domains.

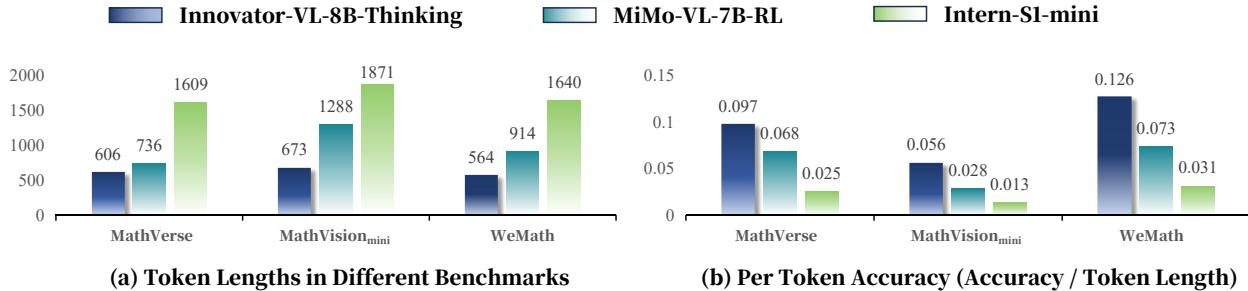


Figure 5 Token efficiency comparison across vision reasoning benchmarks. (a) Average token lengths, showing that Innovator-VL-8B-Thinking generates significantly shorter reasoning chains compared to the other models. (b) Accuracy-to-token ratio, which measures the reasoning efficiency, demonstrating that Innovator-VL-8B-Thinking achieves 1.4x to 2x higher accuracy-to-token ratio than MiMo-VL-7B-RL and 3.9x to 4.3x higher than Intern-S1-mini.

6.4.1 General Multimodal Understanding

As shown in Table 2, **Innovator-VL-8B-Instruct** demonstrates excellent general multimodal capabilities. With an average score of **74.50%**, it is on par with the state-of-the-art Qwen3-VL-8B (74.71%) while significantly outperforming other representative models such as InternVL3.5-8B and LLaVA-OV-1.5-8B. It also demonstrates superior scene understanding on **MME-RealWorld** and achieves the best performance on **AI2D** and **RealWorldQA**. These results indicate that Innovator-VL possesses exceptionally robust visual recognition and instruction-following abilities.

6.4.2 Multimodal Reasoning

The effectiveness of our reinforcement learning (RL) strategy is most prominently reflected in the **Math & Reasoning** category. **Innovator-VL-8B-Thinking** achieves a state-of-the-art average score of **55.41%**, representing a significant **4.54%** absolute improvement over **Innovator-VL-8B-Instruct**. Notably, **Innovator-VL-8B-Thinking** attains the highest performance in this category across all compared models. These results demonstrate that RL enhances the model’s performance in complex reasoning.

6.4.3 Scientific Knowledge

As a specialized scientific multimodal model, **Innovator-VL** exhibits a dominant position in the **Science** category, with its versions securing the top two average scores (**50.13%** and **49.79%**), significantly exceeding all general models. The most striking advantage is observed in specialized chemistry tasks. For instance, on the **OpenRxn** and **MolParse** benchmarks, Innovator-VL achieves scores exceeding **57%** and **64%** respectively, while all other baselines fail to exceed 17%. Such a substantial performance gap demonstrates that Innovator-VL has successfully internalized complex scientific knowledge.

7 On Token Efficiency of Reasoning

Beyond the accuracy improvements demonstrated in Section 6, Innovator-VL exhibits remarkable efficiency in reasoning token consumption. As shown in Figure 5 (a), Innovator-VL-8B-Thinking generates significantly shorter reasoning chains compared to Intern-S1-mini [Bai et al., 2025b] and MiMo-VL-7B-RL [Xiaomi, 2025]. Specifically, our model consumes approximately 62% to 66% fewer tokens than Intern-S1-mini and 18% to 48% fewer than MiMo-VL across the three benchmarks. On WeMath, for instance, our model requires merely 564 average tokens, whereas Intern-S1-mini needs 1,640 tokens.

More importantly, as illustrated in Figure 5 (b), Innovator-VL-8B-Thinking achieves substantially higher accuracy-to-token ratios across all benchmarks, indicating superior reasoning efficiency. Our model delivers approximately 1.4x to 2x higher accuracy-to-token ratio than MiMo-VL-7B-RL and 3.9x to 4.3x higher than Intern-S1-mini. This metric directly captures the effectiveness of each token consumed during reasoning;

higher values indicate that the model extracts more useful information per token, rather than generating redundant or irrelevant content.

This efficiency advantage stems from our reinforcement learning stage that explicitly optimizes for both correctness and conciseness. During RL training, the model learns to identify critical reasoning steps and bypass redundant computations, forming compact yet effective reasoning pathways. This capability is particularly valuable for scientific applications where lengthy inference chains not only incur computational costs but also risk accumulating errors. Compared to recent trends that pursue performance through extensive test-time computation scaling, which often incurs prohibitive token costs, Innovator-VL demonstrates that superior accuracy-efficiency trade-offs can be achieved through principled post-training recipes. Such efficient reasoning ability makes our model more practical for real-world deployment, especially in resource-constrained research environments [Chen et al., 2025] and latency-sensitive applications [Xiong et al., 2025, Yang et al., 2025b, He et al., 2025].

8 Conclusion and Future Works

In this work, we present Innovator-VL, a scientific multimodal large language model that achieves strong performance across diverse scientific domains while maintaining excellent general vision capabilities through a fully transparent and reproducible training pipeline. Our work demonstrates three key principles for building scientific multimodal intelligence: first, that competitive scientific reasoning can be achieved with fewer than five million carefully curated samples, challenging the prevailing reliance on massive domain-specific pretraining; second, that scientific alignment and general-purpose capabilities can coexist within a unified framework without mutual degradation; and third, that principled reinforcement learning training can elicit both accurate and concise reasoning, as evidenced by our model’s superior token efficiency compared to peer models. Looking forward, we aim to extend Innovator-VL toward real-world scientific workflows such as automated experimental design and cross-modal hypothesis generation, expand its modalities to encompass video, 3D molecular structures, and time-series data, develop lightweight variants through compression techniques to broaden accessibility, and enhance its reasoning through integration with external scientific tools and knowledge bases. We hope that our transparent methodology and the practical foundation established by Innovator-VL will enable the research community to advance scientific multimodal modeling and contribute meaningfully to scientific discovery and innovation.

References

- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. [arXiv preprint arXiv:2502.13923](#), 2025a.
- Xiang An, Yin Xie, Kaicheng Yang, Wenkang Zhang, Xiuwei Zhao, Zheng Cheng, Yirui Wang, Songcen Xu, Changrui Chen, Didi Zhu, et al. Llava-onevision-1.5: Fully open framework for democratized multimodal training. [arXiv preprint arXiv:2509.23661](#), 2025.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. [arXiv preprint arXiv:2409.12191](#), 2024a.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. [arXiv preprint arXiv:2501.12599](#), 2025a.
- Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. Kimi-vl technical report. [arXiv preprint arXiv:2504.07491](#), 2025b.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. [arXiv preprint arXiv:2504.10479](#), 2025.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. [arXiv preprint arXiv:2412.19437](#), 2024a.

- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. [arXiv preprint arXiv:2507.20534](#), 2025c.
- Zichen Wen, Yiyu Wang, Chenfei Liao, Boxue Yang, Junxian Li, Weifeng Liu, Haocong He, Bolong Feng, Xuyang Liu, Yuanhuiyi Lyu, et al. Ai for service: Proactive assistance with ai glasses. [arXiv preprint arXiv:2510.14359](#), 2025a.
- Xuan Zhang, Limei Wang, Jacob Helwig, Youzhi Luo, Cong Fu, Yaochen Xie, Meng Liu, Yuchao Lin, Zhao Xu, Keqiang Yan, et al. Artificial intelligence for science in quantum, atomistic, and continuum systems. [arXiv preprint arXiv:2307.08423](#), 2023.
- Lei Bai, Zhongrui Cai, Yuhang Cao, Maosong Cao, Weihan Cao, Chiyu Chen, Haojiong Chen, Kai Chen, Pengcheng Chen, Ying Chen, et al. Intern-s1: A scientific multimodal foundation model. [arXiv preprint arXiv:2508.15763](#), 2025b.
- Xiaohui Zhong, Lei Chen, Xu Fan, Wenzhu Qian, Jun Liu, and Hao Li. Fuxi-2.0: Advancing machine learning weather forecasting model for practical applications. [arXiv preprint arXiv:2409.07188](#), 2024.
- Ning Liao, Xiaoxing Wang, Zehao Lin, Weiyang Guo, Feng Hong, Shixiang Song, Geng Yu, Zihua Zhao, Sitao Xie, Longxuan Wei, et al. Innovator: Scientific continued pretraining with fine-grained moe upcycling. [arXiv preprint arXiv:2507.18671](#), 2025a.
- Hanzheng Li, Xi Fang, Yixuan Li, Chaozheng Huang, Junjie Wang, Xi Wang, Hongzhe Bai, Bojun Hao, Shenyu Lin, Huiqi Liang, et al. Rxnbench: A multimodal benchmark for evaluating large language models on chemical reaction understanding from scientific literature. [arXiv preprint arXiv:2512.23565](#), 2025a.
- Zehua Zhao, Zhixian Huang, Junren Li, Siyu Lin, Junting Zhou, Fengqi Cao, Kun Zhou, Rui Ge, Tingting Long, Yuexiang Zhu, et al. Superchem: A multimodal reasoning benchmark in chemistry. [arXiv preprint arXiv:2512.01274](#), 2025a.
- Yiqing Shen, Zan Chen, Michail Mamalakis, Luhan He, Haiyang Xia, Tianbin Li, Yanzhou Su, Junjun He, and Yu Guang Wang. A fine-tuning dataset and benchmark for large language models for protein understanding. In [2024 IEEE International Conference on Bioinformatics and Biomedicine \(BIBM\)](#), pages 2390–2395. IEEE, 2024.
- Yin Xie, Kaicheng Yang, Xiang An, Kun Wu, Yongle Zhao, Weimo Deng, Zimin Ran, Yumeng Wang, Ziyong Feng, Roy Miles, et al. Region-based cluster discrimination for visual representation learning. In [Proceedings of the IEEE/CVF International Conference on Computer Vision](#), pages 1793–1803, 2025.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. [arXiv preprint arXiv:2505.09388](#), 2025a.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In [Proceedings of the IEEE/CVF conference on computer vision and pattern recognition](#), pages 26296–26306, 2024b.
- Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, et al. Group sequence policy optimization. [arXiv preprint arXiv:2507.18071](#), 2025.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In [International conference on machine learning](#), pages 8748–8763. PMLR, 2021.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In [Proceedings of the IEEE/CVF international conference on computer vision](#), pages 11975–11986, 2023.
- Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. [arXiv preprint arXiv:2502.14786](#), 2025.
- Chenfei Liao, Wensong Wang, Zichen Wen, Xu Zheng, Yiyu Wang, Haocong He, Yuanhuiyi Lyu, Lutao Jiang, Xin Zou, Yuqian Fu, et al. Are we using the right benchmark: An evaluation framework for visual token compression methods. [arXiv preprint arXiv:2510.07143](#), 2025b.
- Zichen Wen, Yifeng Gao, Weijia Li, Conghui He, and Linfeng Zhang. Token pruning in multimodal large language models: Are we solving the right problem? [arXiv preprint arXiv:2502.11501](#), 2025b.
- Evelyn Zhang, Fufu Yu, Aoqi Wu, Zichen Wen, Ke Yan, Shouhong Ding, Binqing Qi, and Linfeng Zhang. D2pruner: Debiased importance and structural diversity for mllm token pruning. [arXiv preprint arXiv:2512.19443](#), 2025a.

- Xuyang Liu, Zichen Wen, Shaobo Wang, Junjie Chen, Zhishan Tao, Yubo Wang, Tailai Chen, Xiangqi Jin, Chang Zou, Yiyu Wang, et al. Shifting ai efficiency from model-centric to data-centric compression. [arXiv preprint arXiv:2505.19147](#), 2025a.
- Yuhang Han, Xuyang Liu, Zihan Zhang, Pengxiang Ding, Junjie Chen, Donglin Wang, Honggang Chen, Qingsen Yan, and Siteng Huang. Filter, correlate, compress: Training-free token reduction for mllm acceleration. [arXiv preprint arXiv:2411.17686](#), 2024.
- Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang, Qintong Zhang, Weijia Li, Conghui He, and Linfeng Zhang. Stop looking for important tokens in multimodal language models: Duplication matters more. [arXiv preprint arXiv:2502.11494](#), 2025c.
- Zichen Wen, Shaobo Wang, Yufa Zhou, Junyuan Zhang, Qintong Zhang, Yifeng Gao, Zhaorun Chen, Bin Wang, Weijia Li, Conghui He, et al. Efficient multi-modal large language models via progressive consistency distillation. [arXiv preprint arXiv:2510.00515](#), 2025d.
- Xuyang Liu, Ziming Wang, Junjie Chen, Yuhang Han, Yingyao Wang, Jiale Yuan, Jun Song, Linfeng Zhang, Siteng Huang, and Honggang Chen. Global compression commander: Plug-and-play inference acceleration for high-resolution large vision-language models. [arXiv preprint arXiv:2501.05179](#), 2025b.
- Minwoo Byeon, Beomhee Park, Haecheon Kim, Sungjun Lee, Woonhyuk Baek, and Saehoon Kim. Coyo-700m: Image-text pair dataset. <https://github.com/kakaobrain/coyo-dataset>, 2022.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. [Advances in Neural Information Processing Systems](#), 36:71683–71702, 2023.
- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. [Advances in Neural Information Processing Systems](#), 36:27092–27112, 2023.
- Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, et al. Fengshenbang 1.0: Being the foundation of chinese cognitive intelligence. [arXiv preprint arXiv:2209.02970](#), 2022.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. [International journal of computer vision](#), 115(3):211–252, 2015.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In [Proceedings of the IEEE/CVF international conference on computer vision](#), pages 4015–4026, 2023.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. [arXiv preprint arXiv:2309.10691](#), 2023.
- Chunyu Xie, Heng Cai, Jincheng Li, Fanjing Kong, Xiaoyu Wu, Jianfei Song, Henrique Morimitsu, Lin Yao, Dexin Wang, Xiangzheng Zhang, et al. Ccmb: A large-scale chinese cross-modal benchmark. In [Proceedings of the 31st ACM International Conference on Multimedia](#), pages 4219–4227, 2023.
- Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data. [arXiv preprint arXiv:2309.16671](#), 2023.
- Yi Zhang, Bolin Ni, Xin-Sheng Chen, Heng-Rui Zhang, Yongming Rao, Houwen Peng, Qinglin Lu, Han Hu, Meng-Hao Guo, and Shi-Min Hu. Bee: A high-quality corpus and full-stack suite to unlock advanced fully open mllms, 2025b. URL <https://arxiv.org/abs/2510.13795>.
- Kaichen Zhang, Keming Wu, Zuhao Yang, Bo Li, Kairui Hu, Bin Wang, Ziwei Liu, Xingxuan Li, and Lidong Bing. Openmmreasoner: Pushing the frontiers for multimodal reasoning with an open and general recipe. [arXiv preprint arXiv:2511.16334](#), 2025c.
- Guowei Xu, Peng Jin, Ziang Wu, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. Llava-cot: Let vision language models reason step-by-step. In [Proceedings of the IEEE/CVF International Conference on Computer Vision](#), pages 2087–2098, 2025.

Yihe Deng, Hritik Bansal, Fan Yin, Nanyun Peng, Wei Wang, and Kai-Wei Chang. Openvlthinker: Complex vision-language reasoning via iterative sft-rl cycles. In The Thirty-ninth Annual Conference on Neural Information Processing Systems, 2025.

Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Jiapeng Wang, Zhuoma Gongque, Shanglin Lei, Yifan Zhang, et al. We-math: Does your large multimodal model achieve human-like mathematical reasoning? In Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 20023–20070, 2025.

Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, Fangzhen Lin, and Wenhua Chen. Vl-rethinker: Incentivizing self-reflection of vision-language models with reinforcement learning. arXiv preprint arXiv:2504.08837, 2025a.

Jinke Li, Jiarui Yu, Chenxing Wei, Hande Dong, Qiang Lin, Liangjing Yang, Zhicai Wang, and Yanbin Hao. Unisvg: A unified dataset for vector graphic understanding and generation with multimodal large language models. In Proceedings of the 33rd ACM International Conference on Multimedia, pages 13156–13163, 2025b.

Zhaorun Chen, Yichao Du, Zichen Wen, Yiyang Zhou, Chenhang Cui, Zhenzhen Weng, Haoqin Tu, Chaoqi Wang, Zhengwei Tong, Qinglan Huang, et al. Mj-bench: Is your multimodal reward model really a good judge for text-to-image generation? arXiv preprint arXiv:2407.04842, 2024a.

Mohammad Shoeybi, Mostafa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. arXiv preprint arXiv:1909.08053, 2019.

Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun Mei, Jiashu Wang, et al. Areal: A large-scale asynchronous reinforcement learning system for language reasoning. arXiv preprint arXiv:2505.24298, 2025.

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibo Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report. arXiv preprint arXiv:2511.21631, 2025c.

LLM-Core-Team Xiaomi. Mimo-vl technical report, 2025. URL <https://arxiv.org/abs/2506.03569>.

Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. arXiv preprint arXiv:2408.01800, 2024.

Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images, 2016. URL <https://arxiv.org/abs/1603.07396>.

Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. Science China Information Sciences, 67(12):220102, 2024c.

Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning, 2022. URL <https://arxiv.org/abs/2203.10244>.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoxi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2024. URL <https://arxiv.org/abs/2311.16502>.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. Advances in Neural Information Processing Systems, 37:95266–95290, 2024b.

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models?, 2024b. URL <https://arxiv.org/abs/2403.20330>.

Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms, 2023. URL <https://arxiv.org/abs/2312.14135>.

- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player?, 2024d. URL <https://arxiv.org/abs/2307.06281>.
- Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun Wang, Qingsong Wen, Zhang Zhang, Liang Wang, Rong Jin, and Tieniu Tan. Mme-realworld: Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for humans?, 2025d. URL <https://arxiv.org/abs/2408.13257>.
- Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images, 2021a. URL <https://arxiv.org/abs/2007.00398>.
- Minesh Mathew, Viraj Bagal, Rubén Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C. V Jawahar. Infographicvqa, 2021b. URL <https://arxiv.org/abs/2104.12756>.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension, 2023. URL <https://arxiv.org/abs/2307.16125>.
- Bohao Li, Yuying Ge, Yi Chen, Yixiao Ge, Ruimao Zhang, and Ying Shan. Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension, 2024. URL <https://arxiv.org/abs/2404.16790>.
- Nirajan Dhakal. Realworldqa dataset. <https://huggingface.co/datasets/nirajandhakal/realworldqa>, 2023. Hugging Face Datasets.
- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset, 2024c. URL <https://arxiv.org/abs/2402.14804>.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems?, 2024a. URL <https://arxiv.org/abs/2403.14624>.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024. URL <https://arxiv.org/abs/2310.02255>.
- Pan Lu, Swaroop Mishra, Tony 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, 2022. URL <https://arxiv.org/abs/2209.09513>.
- Botao Yu, Frazier N. Baker, Ziqi Chen, Xia Ning, and Huan Sun. Llasmol: Advancing large language models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset, 2024. URL <https://arxiv.org/abs/2402.09391>.
- Yuhao Zhou, Yiheng Wang, Xuming He, Ao Shen, Ruoyao Xiao, Zhiwei Li, Qiantai Feng, Zijie Guo, Yuejin Yang, Hao Wu, et al. Scientists' first exam: Probing cognitive abilities of mllm via perception, understanding, and reasoning. [arXiv preprint arXiv:2506.10521](https://arxiv.org/abs/2506.10521), 2025.
- James Burgess, Jeffrey J Nirschl, Laura Bravo-Sánchez, Alejandro Lozano, Sanket Rajan Gupte, Jesus G. Galaz-Montoya, Yuhui Zhang, Yuchang Su, Disha Bhowmik, Zachary Coman, Sarina M. Hasan, Alexandra Johannesson, William D. Leineweber, Malvika G Nair, Ridhi Yarlagadda, Connor Zuraski, Wah Chiu, Sarah Cohen, Jan N. Hansen, Manuel D Leonetti, Chad Liu, Emma Lundberg, and Serena Yeung-Levy. Microvqa: A multimodal reasoning benchmark for microscopy-based scientific research, 2025. URL <https://arxiv.org/abs/2503.13399>.
- Xiangyu Zhao, Wanghan Xu, Bo Liu, Yuhao Zhou, Fenghua Ling, Ben Fei, Xiaoyu Yue, Lei Bai, Wenlong Zhang, and Xiao-Ming Wu. Msearch: A multimodal scientific dataset and benchmark for phenomena uncovering in earth science, 2025b. URL <https://arxiv.org/abs/2505.20740>.
- Fengxiang Wang, Hongzhen Wang, Mingshuo Chen, Di Wang, Yulin Wang, Zonghao Guo, Qiang Ma, Long Lan, Wenjing Yang, Jing Zhang, Zhiyuan Liu, and Maosong Sun. Xlrs-bench: Could your multimodal llms understand extremely large ultra-high-resolution remote sensing imagery?, 2025b. URL <https://arxiv.org/abs/2503.23771>.
- Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Lmms-eval: Reality check on the evaluation of large multimodal models, 2024b. URL <https://arxiv.org/abs/2407.12772>.
- Yuan Chen, Zichen Wen, Yuzhou Wu, Xuyang Liu, Shuang Chen, Junpeng Ma, Weijia Li, Conghui He, and Linfeng Zhang. Ipcv: Information-preserving compression for mllm visual encoders. [arXiv preprint arXiv:2512.18747](https://arxiv.org/abs/2512.18747), 2025.

Minhao Xiong, Zichen Wen, Zhuangcheng Gu, Xuyang Liu, Rui Zhang, Hengrui Kang, Jiabing Yang, Junyuan Zhang, Weijia Li, Conghui He, et al. Prune2drive: A plug-and-play framework for accelerating vision-language models in autonomous driving. [arXiv preprint arXiv:2508.13305](#), 2025.

Yantai Yang, Yuhao Wang, Zichen Wen, Luo Zhongwei, Chang Zou, Zhipeng Zhang, Chuan Wen, and Linfeng Zhang. Efficientvla: Training-free acceleration and compression for vision-language-action models. [arXiv preprint arXiv:2506.10100](#), 2025b.

Peize He, Zichen Wen, Yubo Wang, Yuxuan Wang, Xiaoqian Liu, Jiajie Huang, Zehui Lei, Zhuangcheng Gu, Xiangqi Jin, Jiabing Yang, et al. Audiomarathon: A comprehensive benchmark for long-context audio understanding and efficiency in audio llms. [arXiv preprint arXiv:2510.07293](#), 2025.

A More Results

A.1 Qualitative Case Study

To better illustrate the practical strengths of our approach beyond aggregate benchmark scores, we present a curated set of representative cases sampled from the evaluation benchmarks. For each case, our trained model produces the correct answer, while other baseline models fail under the same evaluation setting. We report the original question (and associated image/figure when applicable), together with the full model responses, to highlight where our model's ability aligns with the task requirements and where competing methods break down. These examples are not meant to be exhaustive; rather, they provide qualitative evidence of improved robustness and generalization in challenging scenarios.

A.1.1 General Visual Tasks

Case 1

Question



What is the color of the handbag?

Choices:

- A. brown
- B. black
- C. red
- D. white

Answer with the option's letter from the given choices directly.

Answer

D

Innovator-VL-8B-Thinking

Model Output

D

Qwen3-VL-8B

Model Output

B

MiMo-VL-7B-RL

Model Output

B

Case 2

Question



What is the color of the clock?

Choices:

- A. black
- B. yellow
- C. green
- D. red

Answer with the option's letter from the given choices directly.

Answer

C

Innovator-VL-8B-Instruct

Model Output

C

InternVL3.5-8B

Model Output

A

Intern-S1-mini (9B)

Model Output

A

Case 3

Question



What direction can we go?

Choices:

- A. Straight
- B. Right
- C. Left

Please answer directly with only the letter of the correct option and nothing else.

Answer

C

[Innovator-VL-8B-Thinking](#)

Model Output

C

[MiniCPM-V4.5 \(8B\)](#)

Model Output

A

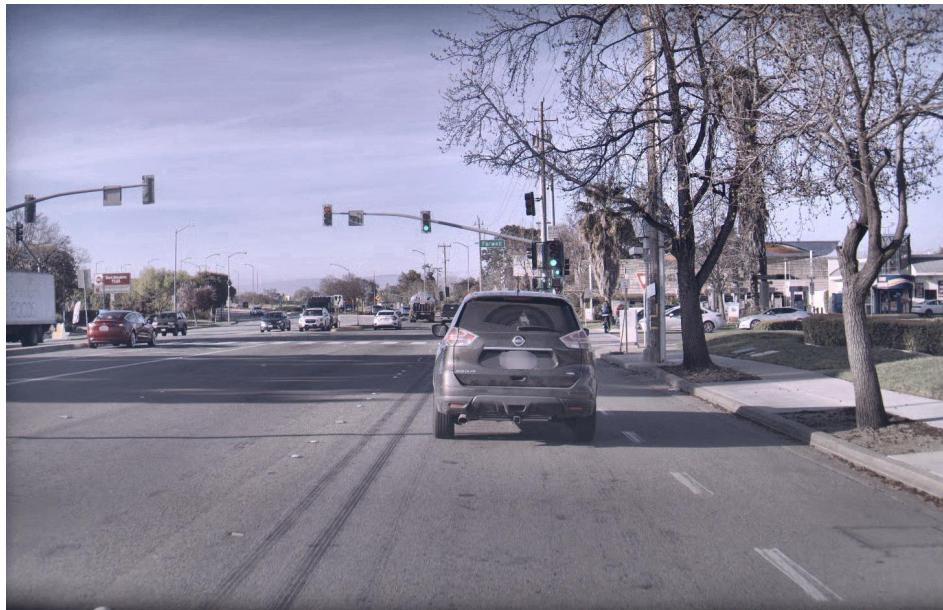
[Qwen3-VL-8B](#)

Model Output

B

[Case 4](#)

Question



The bike rider has a blue backpack.
Please answer directly with a single word or number.

Answer

Yes

Innovator-VL-8B-Instruct

Model Output

Yes

MiMo-VL-7B-RL

Model Output

No

Intern-S1-mini (9B)

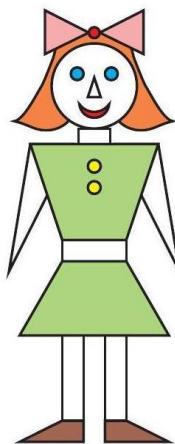
Model Output

No

A.1.2 Math & Reasoning Tasks

Case 1

Question



How many triangles can you find in the picture?

Answer

5

Innovator-VL-8B-Thinking

Model Output

To determine the number of triangles in the given image, let's break it down step by step:

Step 1: Analyze the overall structure

The image depicts a simple stick figure with a dress. The dress has a triangular skirt, and there are other triangular shapes in the figure's head and body.

Step 2: Identify individual triangles

1. Head:

- The bow on top of the head is composed of two triangles.
- The face itself does not contain any triangles, but the eyes and mouth are simple shapes.

2. Dress:

- The skirt of the dress is a large triangle.

3. Body:

- There are no additional triangles in the body beyond the dress.

4. Feet:

- Each foot is a triangle.

Step 3: Count the triangles

- Bow: 2 triangles (one for each half of the bow).
- Skirt: 1 triangle.
- Feet: 2 triangles (one for each foot).

Step 4: Sum up all the triangles

$$\text{Total} = 2 + 1 + 2 = 5.$$

Final Answer:

5

MiniCPM-V4.5 (8B)

Model Output

Let's analyze the image to find all the triangles:

1. First, look at the girl's bow on her head. The bow is a triangle shape.
2. Next, consider the dress. The main part of the dress (the green section without the belt) has a triangular lower part that forms an upside-down triangle.

Now, let's count each distinct triangle:

- The pink bow on top of her head is one triangle.
- The bottom part of the green dress (the skirt area) is another triangle.

Are there any other triangles? Let's check each part again:

- Her face features: eyes and mouth are circles and lines, not triangles.
- Arms and legs are straight lines (rectangles or trapezoids).
- Belt is a rectangle.
- Hair is curved but does not form a triangle in this stylized drawing.

So, only the bow and the bottom part of the dress are triangles.

InternVL3.5-8B

Model Output

So, let's count the triangles in the picture step by step. First, I'll look at the obvious large triangles and then check for smaller ones.

1. **The bow on the head:** There are two large triangles forming the bow.
2. **The dress:** The dress has a large triangle shape on the top part.
3. **The skirt:** The skirt has a large triangle shape.
4. **The arms:** Each arm is a triangle, and there are two arms.
5. **The legs:** Each leg is a triangle, and there are two legs.
6. **The feet:** Each foot is a triangle, and there are two feet.
7. **The buttons:** The buttons on the dress are circles, so they don't count as triangles.

Wait, let me recount to make sure I didn't miss any:

- Bow: 2 triangles
- Dress top: 1 triangle
- Skirt: 1 triangle
- Arms: 2 triangles
- Legs: 2 triangles
- Feet: 2 triangles

Adding them up:

$$2 + 1 + 1 + 2 + 2 + 2 = 10.$$

Hmm, maybe I missed some smaller triangles. Let's check again:

- The bow has 2 large triangles.
- The dress top is 1 large triangle.
- The skirt is 1 large triangle.
- Each arm is 1 triangle, so 2 triangles.
- Each leg is 1 triangle, so 2 triangles.
- Each foot is 1 triangle, so 2 triangles.

That still totals 10. But wait, maybe there are more small triangles in the design. Let's see:

- The bow might have smaller triangles inside, but they seem to be part of the larger shape.
- The dress might have some hidden triangles, but I don't see any.

I think 10 is the correct count. Let me confirm once more:

- Bow: 2
- Dress top: 1
- Skirt: 1
- Arms: 2
- Legs: 2
- Feet: 2

Total:

$$2 + 1 + 1 + 2 + 2 + 2 = 10.$$

So, the total number of triangles should be 10.

Answer: 10

Case 2

Question



Five sparrows are sitting on a rope (see picture). Some of them are looking to the left, some of them are looking to the right. Every sparrow whistles as many times as the number of sparrows he can see sitting in front of him. For example, the third sparrow whistles exactly twice. How often do all sparrows whistle altogether?

Answer

10

Innovator-VL-8B-Thinking

Model Output

To solve this problem, we need to determine how many times each sparrow whistles based on the number of sparrows it can see in front of it. Let's analyze the situation step by step.

Step 1: Understand the Problem

- There are 5 sparrows sitting on a rope.
- Each sparrow whistles as many times as the number of sparrows it can see in front of it.
- We need to calculate the total number of whistles from all sparrows.

Step 2: Analyze Each Sparrow

Let's label the sparrows from left to right as S_1, S_2, S_3, S_4 , and S_5 .

- **Sparrow S_1 :** It can see 0 sparrows in front of it, so it whistles 0 times.
- **Sparrow S_2 :** It can see 1 sparrow (S_1) in front of it, so it whistles 1 time.
- **Sparrow S_3 :** It can see 2 sparrows (S_1 and S_2) in front of it, so it whistles 2 times.
- **Sparrow S_4 :** It can see 3 sparrows (S_1, S_2 , and S_3) in front of it, so it whistles 3 times.
- **Sparrow S_5 :** It can see 4 sparrows (S_1, S_2, S_3 , and S_4) in front of it, so it whistles 4 times.

Step 3: Calculate the Total Number of Whistles

$$S_1 : 0 \quad S_2 : 1 \quad S_3 : 2 \quad S_4 : 3 \quad S_5 : 4$$

$$\text{Total whistles} = 0 + 1 + 2 + 3 + 4 = 10.$$

10

Qwen3-VL-8B

Model Output

Let's analyze the problem step by step.

We are given five sparrows sitting on a rope, and each sparrow whistles as many times as the number of sparrows it can see sitting *in front of him*. The key is the interpretation of "in front of him", which depends on the direction each sparrow is facing.

Looking at the image:

- The first sparrow (leftmost) is waving with its left hand, suggesting it is facing right.
- The second sparrow is waving with its right hand, suggesting it is facing left.
- The third sparrow is waving with its left hand, suggesting it is facing right.
- The fourth sparrow is waving with its right hand, suggesting it is facing left.
- The fifth sparrow (rightmost) is waving with its left hand, suggesting it is facing right.

So the sparrows alternate in direction:

right, left, right, left, right.

Odd-numbered sparrows (S_1, S_3, S_5) face right; even-numbered sparrows (S_2, S_4) face left.

Now, for a sparrow facing right, "in front" means to the right; for a sparrow facing left, "in front" means to the left. Therefore:

- Sparrow S_1 (facing right) sees S_2, S_3, S_4, S_5 : 4 sparrows \Rightarrow whistles 4 times.
- Sparrow S_2 (facing left) sees S_1 : 1 sparrow \Rightarrow whistles 1 time.
- Sparrow S_3 (facing right) sees S_4, S_5 : 2 sparrows \Rightarrow whistles 2 times.
- Sparrow S_4 (facing left) sees S_3 : 1 sparrow \Rightarrow whistles 1 time.
- Sparrow S_5 (facing right) sees none: 0 sparrows \Rightarrow whistles 0 times.

The statement "For example, the third sparrow whistles exactly twice." matches this calculation for S_3 . Thus the total number of whistles is:

$$4 + 1 + 2 + 1 + 0 = 8.$$

Answer: 8

MiMo-VL-7B-RL

Model Output

Let's analyze this step-by-step:

1. **Identify which sparrows are looking left and which are looking right.**

From the image:

- Sparrows 1 and 3 are looking left.
- Sparrows 2, 4, and 5 are looking right.

2. Calculate how many times each sparrow whistles.

- Sparrow 1 (looking left): can see 2 sparrows in front of him (2 and 3) \Rightarrow whistles 2 times.
- Sparrow 2 (looking right): can see 3 sparrows in front of him (3, 4, and 5) \Rightarrow whistles 3 times.
- Sparrow 3 (looking left): can see 2 sparrows in front of him (4 and 5) \Rightarrow whistles 2 times.
- Sparrow 4 (looking right): can see 1 sparrow in front of him (5) \Rightarrow whistles 1 time.
- Sparrow 5 (looking right): can see 0 sparrows in front of him \Rightarrow whistles 0 times.

3. Sum up the whistles.

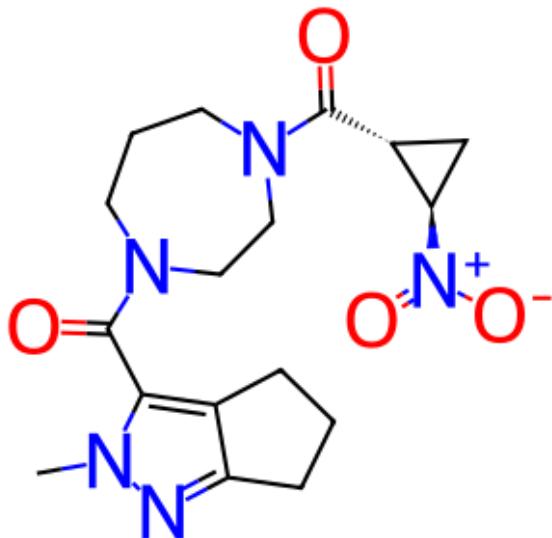
$$2 + 3 + 2 + 1 + 0 = 8.$$

Therefore, all sparrows whistle a total of 8 times.

A.1.3 Scientific Tasks

Case 1

Question



Please write the Extended SMILES (E-SMILES) of the given molecule image.

Below is the definition of Extended SMILES (E-SMILES) Format Writing Rules:

The Extended SMILES (E-SMILES) format used in the MolParser paper is designed to augment a standard RDKit-compatible SMILES string with additional machine-readable annotations, especially useful for representing Markush structures, abstract rings, repeating groups, and other extensions commonly seen in chemical patents. The general form of an E-SMILES string is:

SMILES<sep>EXTENSION

Components:

1. SMILES

This is a valid RDKit-compatible SMILES string.

To ensure consistency, the atom indexing starts from 0. You can generate this using RDKit with `rootedAtAtom=0`.

2. <sep>

A special delimiter separating the base SMILES string from the extension part.

3. EXTENSION

The extension part contains structured XML-like tokens providing additional information. These tokens annotate atom positions, ring positions, and other special features.

There are three types of tags:

- <a>ATOM_INDEX : [GROUP_NAME] — attaches a group to a specific atom.
- <r>RING_INDEX : [GROUP_NAME] </r> — attaches a group to any position on a ring.
- <c>CIRCLE_INDEX : [CIRCLE_NAME] </c> — defines abstract repeating rings.

A special token <dum> can be used to represent a connection point.

Definitions:

ATOM_INDEX: Index of the atom in the SMILES string (starting from 0).

RING_INDEX: Index of the ring in the molecule (starting from 0).

CIRCLE_INDEX: Index of the abstract repeating ring (starting from 0).

GROUP_NAME: The name of the substituent or group, which may be:

- A standard group: R, X, Y, Z, Ph, Me, OMe, CF₃, etc.
- An arbitrary or custom-defined group with superscripts or subscripts, e.g., R[1], X[2], R₁, R₃.

Examples:

1. *c1ccccc1<sep><a>0:R[1]
→ Adds substituent R[1] to atom 0.
2. clccecc1<sep><r>0:R[1]</r><r>0:R[2]</r>
→ Adds two substituents R[1], R[2] to ring 0.
3. *C1*C(*)C(C(*)C2=CC=N=C2)N=1<sep><a>0:R[4]<a>1:X<r>1:R[5]</r>
→ Complex structure with atom-level and ring-level substitutions.
4. *c(NC(*)C(*)C(*)C(*)C(=O)N(*))<sep>
<a>0:R[1]<a>2:R[3]<a>5:R[2]
<a>7:R[5]<a>8:R[4]<c>0:B</c>
<a>13:R[7]<a>14:R[6]
→ Shows the usage of abstract ring <c>, with multiple atom-level substitutions.

Answer

Cn1nc2c(c1C(=O)N1CCCN(C(=O)[C@H]3C[C@H]3[N+](=O)[O-])CC1)CCC2<sep>

Innovator-VL-8B-Instruct

Model Output

Cn1nc2c(c1C(=O)N1CCCN(C(=O)[C@H]3C[C@H]3[N+](=O)[O-])CC1)CCC2<sep>

InternVL3.5-8B

Model Output

To write the E-SMILES for the given molecule, we need to identify the SMILES string and then add any necessary extensions. Let's break down the structure:

SMILES String

The SMILES representation of the molecule is:

CN1N=C2CCCC2=C1C(=O)N1CCN(C(=O)[C@H]2C[C@H]2[N+](=O)[O-])CC1

Extensions

In this molecule, there are no explicit substituents or abstract rings that require additional annotations beyond the standard SMILES. Therefore, the E-SMILES will only include the SMILES string.

E-SMILES

CN1N=C2CCCC2=C1C(=O)N1CCN(C(=O)[C@H]2C[C@H]2[N+](=O)[O-])CC1<sep>

This E-SMILES string follows the format by providing the SMILES and an empty extension part, as no further annotations are needed for this structure.

Intern-S1-mini (9B)

Model Output

To write the Extended SMILES (E-SMILES) of the given molecule, we need to follow the rules provided. Here's a step-by-step breakdown:

1. SMILES Representation:

First, we need to write the SMILES string for the given molecule. The SMILES string for the molecule is:

CN1N=C2CCCC2=C1C(=O)N1CCN(C(=O)[C@H]2C[C@H]2[N+](=O)[O-])CC1

2. Extension Part:

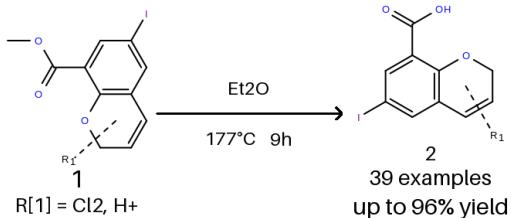
Next, we need to add the extension part to the SMILES string. The extension part contains structured XML-like tokens providing additional information.

- **Atom Indexing:** We need to identify the atom indices where we want to attach groups. In this case, we can attach groups to the nitrogen atom in the pyrazole ring (atom index 0) and the nitrogen atom in the piperazine ring (atom index 1).
- **Group Names:** We can use standard group names like R, X, Y, Z, Ph, Me, OMe, CF₃, etc. For this example, let's use R for the groups.
- **Delimiter:** We use the special delimiter <sep> to separate the base SMILES

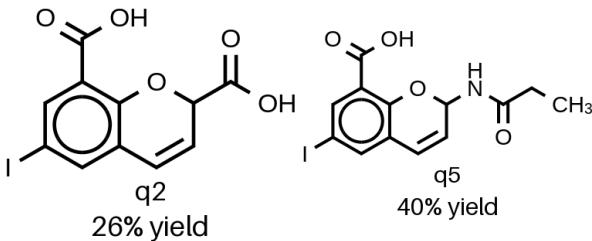
Case 2

Question

General Reaction:



Selected Examples:



Return the ESMILES of the select example with label q2.

Answer

```
<E-SMILES>OC(=O)c1cc(I)cc2c1OC(C(=O)O)C=C2</E-SMILES>
```

Innovator-VL-8B-Thinking

Model Output

```
<E-SMILES>OC(=O)c1cc(I)cc2c1OC(C(=O)O)C=C2</E-SMILES>
```

Qwen3-VL-8B

Model Output

```
O=C(O)c1cc2c(cc1I)OC=C2C(=O)O
```

Intern-S1-mini (9B)

Model Output

The ESMILES (Extended SMILES) representation of the selected example with label q2 is:

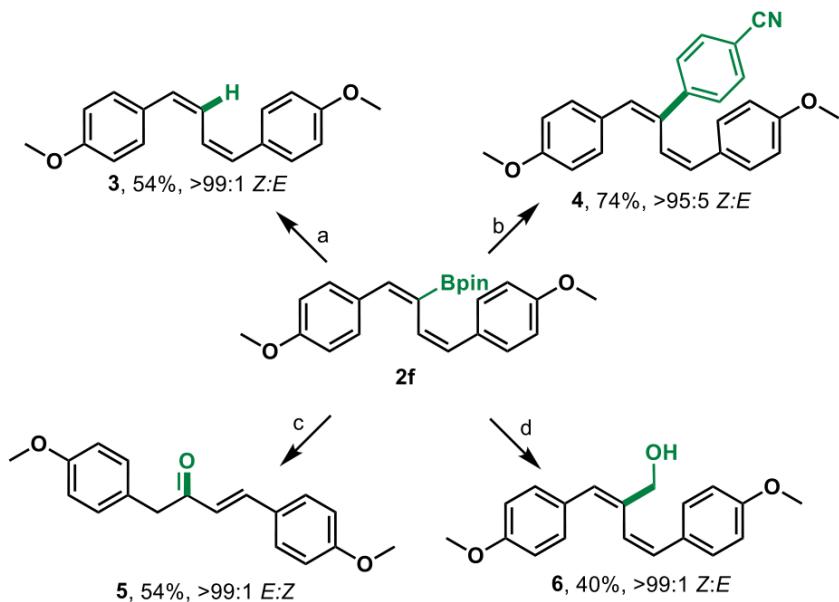
COC(=O)C1=C(I)C=CC2=C1OC(C(=O)O)C=C2

This SMILES string represents the structure of the compound shown in the image labeled q2.

Case 3

Question

Scheme 4. Synthetic Applications of 2-Boryl-1,3-Dienes^a



^aReaction conditions: a: KHF₂ (3.0 equiv), acetic acid (0.5 mL). b: Pd₂(dba)₃ (4 mol %), SPhos (5 mol %), 4-iodobenzonitrile (1.5 equiv), THF/3 M NaOH (3:1), 70 °C. c: H₂O₂ (30 equiv), 3 M NaOH (30 equiv), THF (0.5 mL). d: *n*BuLi (3.0 equiv), dibromomethane (4.0 equiv), THF (1.0 mL) –78 to 25 °C, 4 h then H₂O₂ (5.0 equiv), 3 M NaOH (5.0 equiv).

Question: Please summarize the four derivatization reactions (a–d) centered on 2-boryl-1,3-diene 2f in Scheme 4 and their corresponding products (3–6).

Choices:

- A. Scheme 4 shows four reaction pathways starting from 2-boryl-1,3-diene 2f:
 - a) KHF₂/AcOH effects removal of Bpin to give the deborylated product 3 (yield 54%, >99:1 Z:E);
 - b) Pd₂(dba)₃/SPhos-catalyzed Suzuki coupling with 4-iodobenzonitrile affords the arylated product 4 (yield 74%, >95:5 Z:E);
 - c) H₂O₂/NaOH oxidation of the boronate gives the $\alpha,\beta,\gamma,\delta$ -unsaturated ketone product 5 (yield 54%, >99:1 Z:E);
 - d) *n*BuLi/CH₂Br₂ insertion followed by H₂O₂ oxidation provides the allylic alcohol product 6 (yield 40%, >99:1 Z:E).
- B. Scheme 4 shows four reaction pathways starting from 2-boryl-1,3-diene 2f:
 - a) KHF₂/AcOH effects removal of Bpin to give the deborylated product 3 (yield 54%, 99:1 Z:E);
 - b) Pd₂(dba)₃/SPhos-catalyzed Suzuki coupling with 4-iodobenzonitrile affords the arylated product 4 (yield 74%, >95:5 Z:E);
 - c) H₂O₂/NaOH oxidation of the boronate gives the $\alpha,\beta,\gamma,\delta$ -unsaturated ketone product 5 (yield 54%, >99:1 E:Z);
 - d) *n*BuLi/CH₂Br₂ insertion followed by H₂O₂ oxidation provides the allylic alcohol product 6 (yield 40%, >99:1 Z:E).

- C. Scheme 4 shows four reaction pathways starting from 2-boryl-1,3-diene 2f:
- a) KHF₂/AcOH effects removal of Bpin to give the deborylated product 3 (yield 54%, >99:1 Z:E);
 - b) Pd₂(dba)₃/SPhos-catalyzed Suzuki coupling with 4-iodobenzonitrile affords the arylated product 4 (yield 74%, >95:5 Z:E);
 - c) H₂O₂/NaOH oxidation of the boronate gives the $\alpha, \beta, \gamma, \delta$ -unsaturated ketone product 5 (yield 54%, >99:1 E:Z);
 - d) *n*BuLi/CH₂Br₂ insertion followed by H₂O₂ oxidation provides the allylic alcohol product 6 (yield 40%, >99:1 Z:E).
- D. Scheme 4 shows four reaction pathways starting from 2-boryl-1,3-diene 2f:
- a) KHF₂/AcOH effects removal of Bpin to give the deborylated product 3 (yield 54%, >99:1 Z:E);
 - b) Pd₂(dba)₃/SPhos-catalyzed Suzuki coupling with 4-iodobenzonitrile affords the arylated product 4 (yield 74%, >95:5 Z:E);
 - c) H₂O₂/NaOH oxidation of the boronate gives the $\alpha, \beta, \gamma, \delta$ -unsaturated ketone product 5 (yield 54%, >99:1 E:Z);
 - d) *n*BuLi/CH₂Br₂ insertion followed by O₂ oxidation provides the allylic alcohol product 6 (yield 40%, >99:1 Z:E).

Based on the image and the question, choose the most appropriate answer.

Only output a single letter (A, B, C, or D). Do NOT output any other text or explanation.

Answer

C

Innovator-VL-8B-Instruct

Model Output

C

InternVL3.5-8B

Model Output

B

Intern-S1-mini (9B)

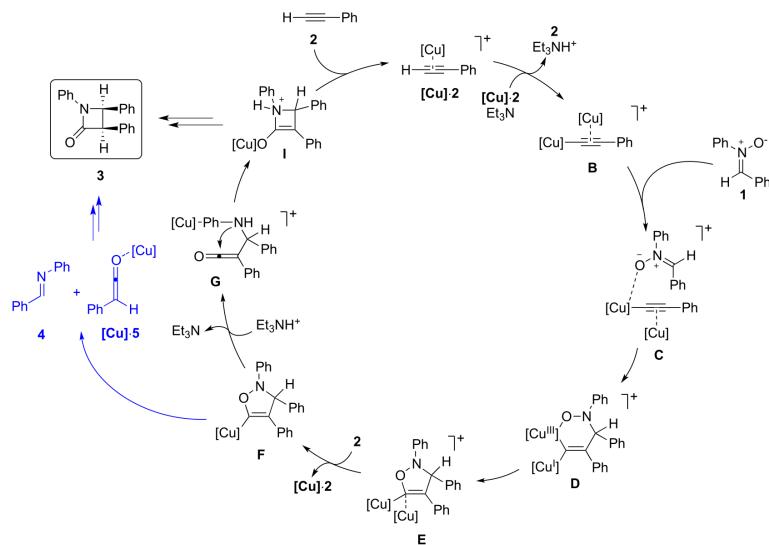
Model Output

A

Case 4

Question

Scheme 2. Previously Proposed Reaction Mechanism Based on DFT Calculations (Black),⁸ with the Alternative Pathway Derived from Experiments (Blue)^{a9}



^aNote that the naming of the species follows the nomenclature of the previous calculations reported in ref 8.

问题：该反应的总体转化通式是什么？请用结构简式描述主要底物和产物。

选项：

- A. 底物苯基炔（ $\text{Ph}-\text{C}\equiv\text{CH}$ ）与N-氧化亚胺1在Cu催化和 Et_3N 作用下，最终生成五元环化合物3。
- B. 底物苯基炔（ $\text{Ph}-\text{C}\equiv\text{CH}$ ）与N-氧化亚胺1在Cu催化和 Et_3N 作用下，最终生成四元环化合物3。
- C. 底物苯基炔（ $\text{Ph}-\text{C}\equiv\text{CH}$ ）与N-氧化亚胺1在Cu催化和吡啶作用下，最终生成四元环化合物3。
- D. 底物苯基炔（ $\text{Ph}-\text{C}\equiv\text{CH}$ ）与N-氧化亚胺1在Cu催化和 Et_3N 作用下，最终生成桥环化合物3。

请根据图像和问题，从以上四个选项中选择最合适答案。

只输出单个字母（A, B, C 或 D），不要输出选项内容，也不要输出任何解释。

Answer

B

Innovator-VL-8B-Thinking

Model Output

B

MiMo-VL-7B-RL

Model Output

A

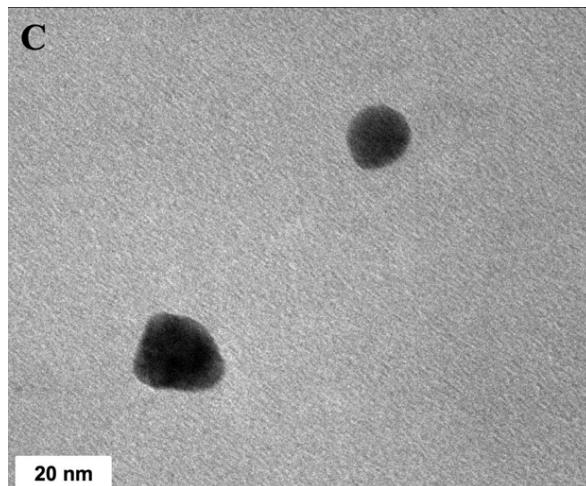
MiniCPM-V4.5 (8B)

Model Output

A

Case 5

Question



Carefully examine the electron microscopy image, select the most accurate diagnosis based on the visual evidence observed in the microstructure, and answer the question: Is the particle-area distribution uniform? (Define non-uniform as $\frac{\text{max/min diameter}}{\text{mean}} > 30\%$) by choosing the most suitable option from:

- A. Yes
- B. No
- C. Uncertain

Answer

A

Innovator-VL-8B-Instruct

Model Output

A

Qwen3-VL-8B

Model Outpu

B

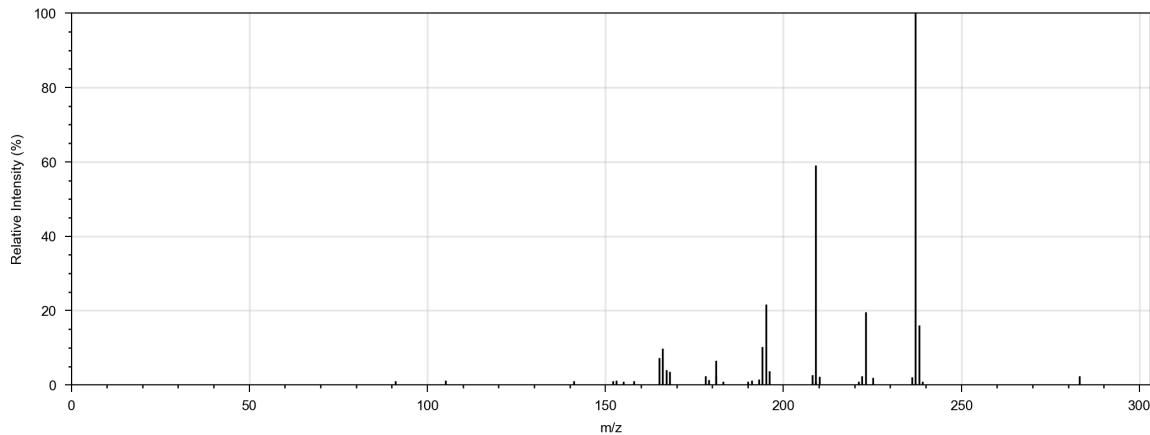
MiMo-VL-7B-RL

Model Output

B

Case 6

Question



You are an expert in life and need to solve the following question. The question is an exact match question. Answer the question using a single word or phrase.

From the compound's MS/MS spectrum, determine the number of sulfur atoms based on isotope pattern and other compound information in the spectrum.

Answer

0

Innovator-VL-8B-Thinking

Model Output

0

Qwen3-VL-8B

Model Output

2

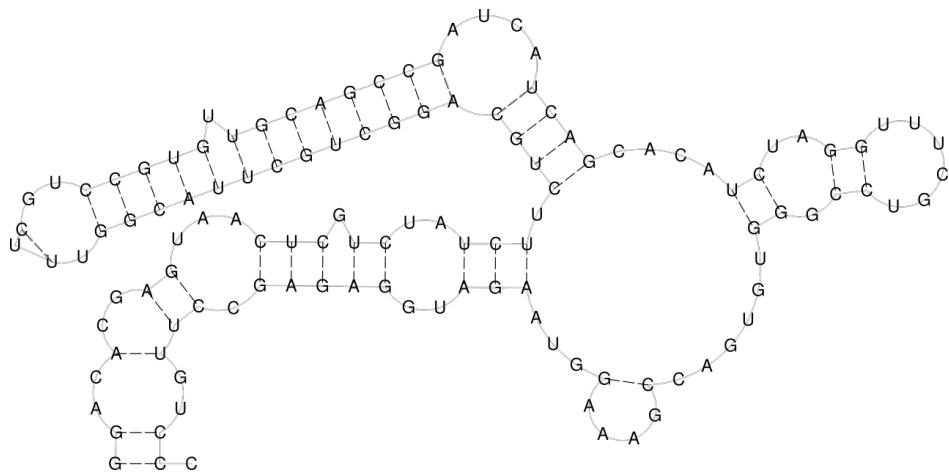
MiMo-VL-7B-RL

Model Output

1

Case 7

Question



You are an expert in life and need to solve the following question. The question is a multiple-choice question. Answer with the option letter from the given choices.

Based on the RNA secondary structure diagram, which of the following structures might the RNA contain? (Multiple choices)

- (A) stem
- (B) Hairpin loop
- (C) Interior loop
- (D) bulge loop
- (E) Multi-branched loop
- (F) Pseudoknot

Answer

ABCD

Innovator-VL-8B-Instruct

Model Output

ABCD

Qwen3-VL-8B

Model Output

ABCDE

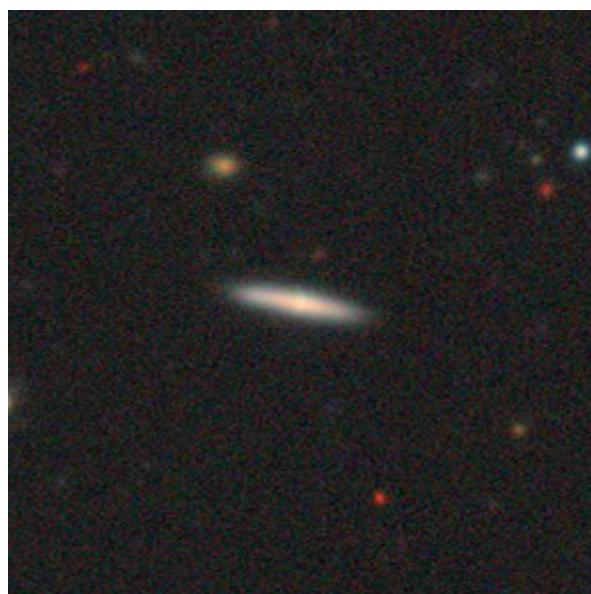
MiMo-VL-7B-RL

Model Output

ABCDEF

Case 8

Question



You are an expert in astronomy and need to solve the following question. The question is a multiple-choice question. Answer with the option letter from the given choices. What is the structure of galaxy in the image?

- (A) Unbarred Tight Spiral Galaxies
- (B) Unbarred Loose Spiral Galaxies
- (C) Cigar Shaped Smooth Galaxies
- (D) Edge-on Galaxies with out Bulge

Answer

D

Innovator-VL-8B-Instruct

Model Output

D

Qwen3-VL-8B

Model Output

C

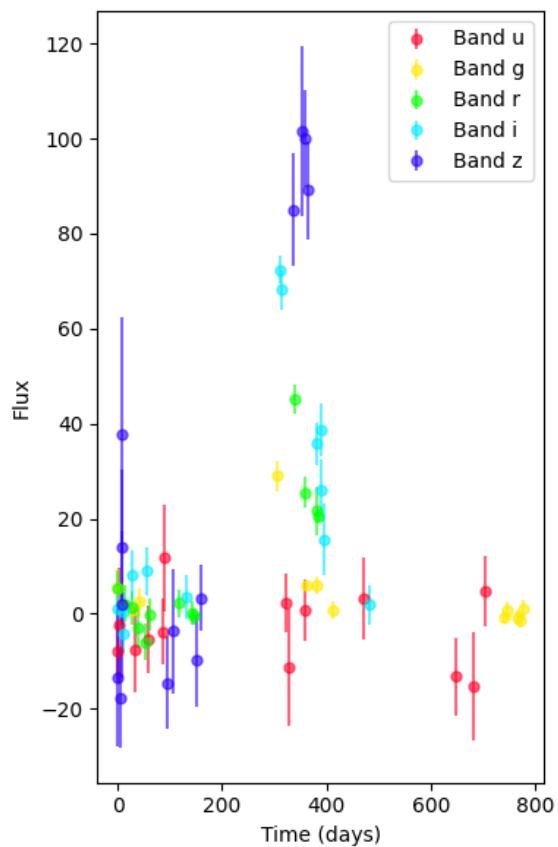
MiMo-VL-7B-RL

Model Output

C

Case 9

Question



根据所提供的不同波段的图像，这是什么类型的光变曲线？

- (A) EB
 (B) KN
 (C) SNII
 (D) SNIbc

Answer

C

Innovator-VL-8B-Thinking

Model Output

C

InternVL3.5-8B

Model Output

B

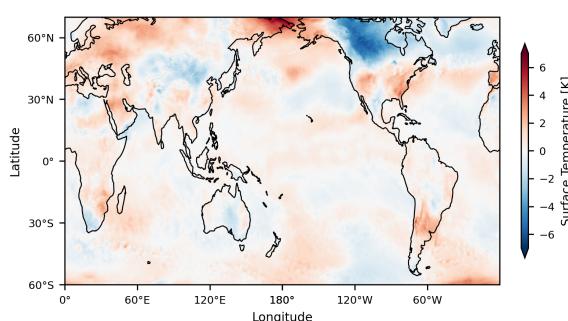
Intern-S1-mini (9B)

Model Output

D

Case 10

Question



下图显示了全球地表温度异常。请分别选择温度异常最高和最低的大致区域。

- 白令海
- 澳大利亚
- 中亚
- 北美北部

Answer

AD

Innovator-VL-8B-instruct

Model Output

AD

InternVL3.5-8B

Model Output

CD

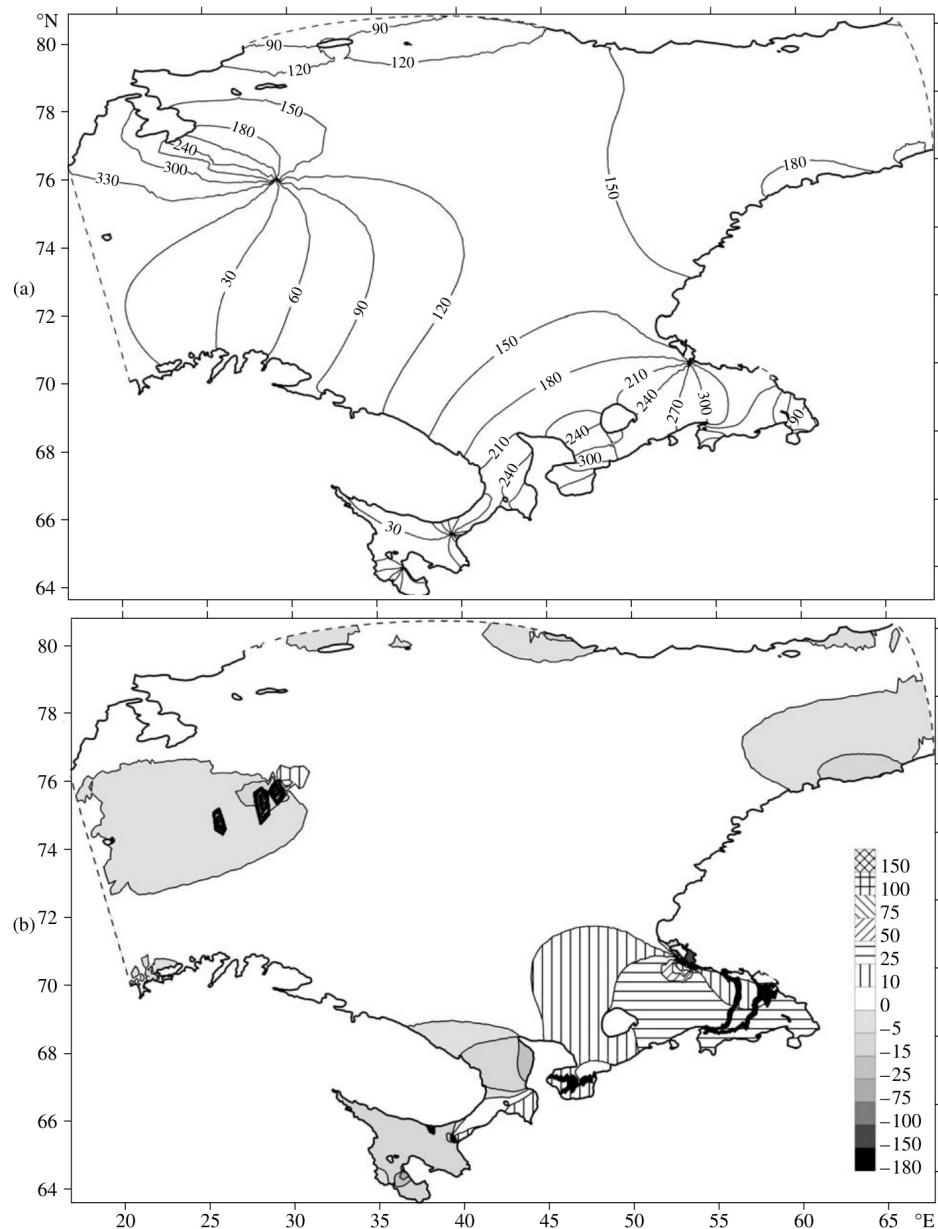
MiMo-VL-7B-RL

Model Output

A

Case 11

Question



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

Caption: (a) Phases (in degrees) of tidal sea surface level elevations in the ice-free period and (b) their variations (in degrees) under ice or ice-free conditions in the White Sea.

Question: Which region exhibits degenerate amphidromy characteristics during the ice-free period based on Figure (a)?

Options:

- A. Near the throat outlet
- B. Near Cape Letni Orlov
- C. In the northern section of the White Sea

D. Along the southern coastline

Based on the image, select the correct option (e.g., A, B, C) or directly state the correct option content.

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

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

Answer

B. Near Cape Letni Orlov

Innovator-VL-8B-Instruct

Model Output

```
{  
    "answer": "B",  
    "Explanation": "The image shows that near Cape Letni Orlov, the  
    tidal sea surface level elevations exhibit degenerate amphidromy  
    characteristics during the ice-free period."  
}
```

Qwen3-VL-8B

Model Output

```
{  
    "answer": "C. In the northern section of the White Sea",  
    "Explanation": "In Figure (a), the contour lines representing  
    tidal phase (in degrees) show a pattern where the phase values  
    increase radially outward from a central point in the northern  
    section of the White Sea, indicating a degenerate amphidromy. This  
    is characteristic of a single amphidromic point with a radial phase  
    progression, which is typical in the northern part of the White Sea  
    during the ice-free period."  
}
```

MiMo-VL-7B-RL

Model Output

```
{  
    "answer": "C",  
    "Explanation": "In Figure (a), the northern section of the  
    White Sea shows degenerate amphidromy characteristics as indicated  
    by the phase lines converging towards a central point, which is  
    characteristic of a single amphidromic point."}
```

```
    typical of degenerate amphidromy."  
}
```

B Evaluation Details

B.1 Benchmark Details.

- **General Vision:**

- **AI2D:** AI2D [Kembhavi et al., 2016] contains over 5,000 charts and 15,000 questions and answers, focusing on the syntactic parsing and question-answering tasks for charts. It uses Diagram Parse Graphs (DPG) to support the understanding and reasoning of the structure of charts and their constituent relationships.
- **OCRBench:** OCRBench [Liu et al., 2024c] is a comprehensive evaluation benchmark consisting of 29 datasets, aimed at assessing the performance of large multimodal models in text-related visual tasks. It covers a wide range of text recognition tasks, including text recognition, scene text-based question answering, document-based question answering, key information extraction, and handwritten mathematical expression recognition, providing diverse scenarios and tasks to comprehensively evaluate OCR capabilities.
- **ChartQA:** ChartQA [Masry et al., 2022] includes 9.6K human-written questions and 23.1K questions generated from human-written chart summaries, focusing on complex reasoning tasks for charts. Unlike existing template-based question sets, the tasks in this dataset involve multiple logical and arithmetic operations, and the questions often refer to the visual features of the charts. These questions assess the visual and logical reasoning abilities of models and challenge their capacity to answer complex reasoning questions.
- **MMMU:** MMMU [Yue et al., 2024] is a large-scale multimodal benchmark with 11.5K questions spanning 6 core disciplines, 30 subjects, 183 subfields, and 30 image types. It is designed to test advanced perception and domain-specific reasoning abilities.
- **MMMU-Pro:** MMMU-Pro [Yue et al., 2024] is an enhanced version of the MMMU benchmark, which rigorously evaluates the understanding and reasoning capabilities of multimodal models through a three-step evaluation process: (1) filtering out questions answerable by text-only models, (2) augmenting candidate options, and (3) introducing a vision-only input setting where questions are embedded within images. This setting tests models on their ability to process both visual and textual information simultaneously, and results show that models perform significantly worse on MMMU-Pro compared to MMMU.
- **MMStar:** MMStar [Chen et al., 2024b] is a vision-indispensable benchmark with 1,500 human-curated items balanced over six core capabilities and 18 axes. It explicitly filters out samples that can be solved without images and proposes metrics to quantify the true multimodal gains.
- **VStar-Bench:** VStar-Bench [Wu and Xie, 2023] is a benchmark specifically designed to evaluate large multimodal language models (MLLMs) in their ability to process high-resolution images and focus on visual details. The benchmark introduces a V-visual search mechanism to test MLLMs' performance in visual querying, collaborative reasoning, and contextual understanding, emphasizing the importance of incorporating visual search capabilities into multimodal systems.
- **MBBench:** MBBench [Liu et al., 2024d] is a systematically designed bilingual multimodal evaluation benchmark that covers a wide range of fine-grained vision–language capabilities. Through rigorous data quality control and a cyclic evaluation strategy, it provides robust and objective model assessment. The benchmark supports both Chinese and English multiple-choice evaluations and introduces large language model-assisted grading mechanisms, enabling fair and scalable comparison across models with different instruction-following capabilities.

- **MME-RealWorld:** MME-RealWorld [Zhang et al., 2025d] is a large-scale, human-annotated multimodal benchmark focusing on high-resolution images and complex real-world scenarios. It systematically characterizes the core bottlenecks of current MLLMs in perception and reasoning. By covering diverse scenes and subtasks with challenging, manually constructed questions, the benchmark reveals significant shortcomings even in state-of-the-art models under realistic application settings.
- **DocVQA:** DocVQA [Mathew et al., 2021a] is a document image-based visual question answering dataset containing over 12,000 document images and 50,000 questions. It is designed to systematically evaluate models' capabilities in document understanding, layout reasoning, and structural perception.
- **InfoVQA:** InfoVQA [Mathew et al., 2021b] targets infographic understanding and visual question answering, covering diverse information graphics annotated with natural language questions and answers. The dataset requires joint reasoning over document layout, textual content, graphical elements, and data visualizations, emphasizing foundational reasoning and arithmetic skills in complex infographic comprehension.
- **SEED-Bench:** SEED-Bench [Li et al., 2023] is a benchmark for evaluating the understanding capabilities of generative multimodal large language models. It consists of 19,000 human-annotated multiple-choice questions spanning 12 key understanding dimensions across both image and video modalities.
- **SEED-Bench-2-Plus:** SEED-Bench-2-Plus [Li et al., 2024] is a benchmark for evaluating text-rich visual comprehension of multimodal large language models, consisting of 2.3K human-annotated multiple-choice questions. The dataset covers three categories—Charts, Maps, and Webs—to assess model performance in real-world scenarios involving dense embedded text in images.
- **RealWorldQA:** RealWorldQA [Dhakal, 2023] contains 765 images annotated with question–answer pairs, covering anonymized vehicle-captured data and a wide range of real-world scenes. It aims to evaluate multimodal models under realistic and unconstrained visual environments.

- **Math & Reasoning:**

- **MathVista:** MathVista [Lu et al., 2024] targets mathematical reasoning in visual contexts, aggregating 6,141 examples from 28 existing datasets and 3 newly created ones, testing models' abilities in fine-grained visual understanding and compositional reasoning.
- **MathVision:** MathVision [Wang et al., 2024c] curates 3,040 competition-style math problems with visual contexts across 16 disciplines and five difficulty levels, highlighting gaps between current models and human performance in multimodal mathematical reasoning.
- **MathVerse:** MathVerse [Zhang et al., 2024a] is a comprehensive visual mathematics benchmark for multimodal large language models, consisting of 2,612 image-based math problems. It introduces multiple modality variants to fairly assess whether models truly leverage visual information rather than relying solely on textual reasoning. The benchmark further adopts fine-grained chain-of-thought-based evaluation strategies to analyze intermediate reasoning processes in depth.
- **WeMath:** WeMath [Qiao et al., 2025] is the first benchmark focusing on knowledge acquisition and generalization mechanisms in visual mathematical reasoning. It contains 6.5K problems covering 67 hierarchical mathematical knowledge concepts. By decomposing questions into subproblems and introducing four-dimensional evaluation metrics, WE-MATH systematically reveals the fundamental capabilities and limitations of multimodal models in understanding and generalizing mathematical knowledge.

- **Science:**

- **ScienceQA:** ScienceQA [Lu et al., 2022] is a large-scale multimodal science question-answering benchmark containing approximately 21,000 multiple-choice questions across diverse scientific

disciplines. Each question is accompanied by detailed explanations and reasoning chains, explicitly designed to characterize models' chain-of-thought reasoning, interpretability, and generalization in multi-hop multimodal scientific reasoning.

- **RxnBench:** RxnBench [Li et al., 2025a] is a multimodal benchmark for chemical reaction understanding, constructed from real scientific paper PDFs. It evaluates models' joint understanding of reaction diagrams and textual descriptions through both single-image QA and full-document QA tasks, systematically exposing limitations in chemical structure recognition and mechanistic reasoning.
- **MolParse:** Mol QA & Parser is an evaluation benchmark for molecular understanding, containing 2,000 question–answer pairs based on single-molecule images. It is used to evaluate molecular question answering and structure parsing tasks, including functional group recognition, property-related queries, and SMILES or extended SMILES generation from molecular images.
- **OpenRxn:** Reaction is an evaluation benchmark for chemical reaction understanding, containing 2,000 reaction-based question–answer pairs derived from organic reaction images and SMILES representations. It is used to evaluate chemical knowledge and reaction-level reasoning in large language models.
- **EMVista:** A benchmark of 1,850 EM microstructure samples evaluating instance-level perception, microstructural attribute understanding (morphology, density, spatial distribution, layering, scale variation), and robustness under dense, overlapping, multi-scale scenes.
- **SMolInstruct:** SMolInstruct [Yu et al., 2024] is a large-scale benchmark for evaluating chemistry-oriented language models, covering 14 small-molecule tasks and approximately 3.3M question–answer pairs with standardized evaluation splits.
- **SuperChem:** SUPERChem [Zhao et al., 2025a] is a 500-problem, expert-curated chemistry reasoning benchmark (multimodal + text-only) with expert solution paths for process-level evaluation via Reasoning Path Fidelity (RPF).
- **ProteinLMBench:** ProteinLMBench [Shen et al., 2024] consists of 944 manually verified multiple-choice questions covering protein sequence, structure, and function. It is specifically designed for evaluating protein understanding in LLMs.
- **SFE:** Scientists' First Exam (SFE) [Zhou et al., 2025] probes scientific cognition of MLLMs across perception, attribute understanding, and comparative reasoning, comprising 830 expert-verified VQA items across five disciplines.
- **MicroVQA:** MicroVQA [Burgess et al., 2025] is a microscopy-centric multimodal benchmark with 1,042 multiple-choice questions spanning diverse imaging modalities and biological topics, targeting scientific analysis and reasoning.
- **MSEarth-MCQ:** MSEarth-MCQ [Zhao et al., 2025b] is the multiple-choice subset of the MSEarth Earth-science benchmark, providing approximately 2.78K expert-derived figure-grounded questions across the atmosphere, cryosphere, hydrosphere, and biosphere.
- **XLRS-Bench:** XLRS-Bench [Wang et al., 2025b] evaluates MLLMs on ultra-high-resolution remote-sensing imagery (average resolution 8.5k × 8.5k pixels), with human-verified annotations across diverse tasks reflecting real-world RS scenarios.

B.2 Evaluation Prompts

AI2D

```
<image>
{question}
A. {choice_1}
B. {choice_2}
C. {choice_3}
D. {choice_4}
```

Answer with the option's letter from the given choices directly.

AI2D (no mask)

```
<image>
{question}
A. {choice_1}
B. {choice_2}
C. {choice_3}
D. {choice_4}
```

Answer with the option's letter from the given choices directly.

OCRBench

```
<image>
{question}
```

ChartQA

```
<image>
{question}
```

Answer the question with a single word.

MMMU (Val)

multiple_choice_prompt:

```
<image>
{question}
A. {option_1}
B. {option_2}
C. {option_3}
D. {option_4}
```

Answer with the option's letter from the given choices directly.

open-ended_prompt:

```
<image>
{question}
```

Answer the question using a single word or phrase.

MMMU-Pro (Standard)

<image>
{question}
A. {choice_1}
B. {choice_2}
...
J. {choice_10}

Answer with the option letter from the given choices directly.

MMStar

<image>
{question}

Answer with the option's letter from the given choices directly.

VStar-Bench

<image>
{question}
A. {option_1}
B. {option_2}
C. {option_3}
D. {option_4}

Answer with the option's letter from the given choices directly.

MMBench-EN (Dev/Test)

<image>

{hint} {question} {options}

Answer with the option's letter from the given choices directly.

MME-RealWorld (EN)

<image>
{question} The choices are listed below:
{options}

Select the best answer to the above multiple-choice question based on the image. Respond with only the letter (A, B, C, D, or E) of the correct option.

The best answer is:

MME-RealWorld (CN)

<image>
{question} 选项如下所示:
{options}

根据图像选择上述多项选择题的最佳答案。只需回答正确选项的字母（A, B, C, D 或 E）。
最佳答案为：

DocVQA(Val)

<image>
{question}
Answer the question using a single word or phrase.

InfoVQA(Val)

<image>
{question}
Answer the question using a single word or phrase.

SEED-Bench / SEED-Bench-2-Plus

<image>
{question}
A. {choice_A} ... D. {choice_D}
Answer with the option's letter from the given choices directly.

RealWorldQA

<image>
{question}

MathVision / MathVision-mini

Think and solve the following question step by step. Please put your thinking and analysis procedure within <think></think>. Put ONLY your final answer within <answer></answer>.

<image>
{question}
Choices: {choices}

MathVerse

Think and solve the following question step by step. Please put your thinking and analysis procedure within <think></think>. Put ONLY your final answer within <answer></answer>.

<image>
{question}

MathVista

Think and solve the following question step by step. Please put your thinking and analysis procedure within <think></think>. Put ONLY your final answer within <answer></answer>.

Question: {question} (Unit: {unit})
Choices: {choices}
Hint: {hint}

WeMath

Think and solve the following question step by step. Please put your thinking and analysis procedure within <think></think>. Put ONLY your final answer within <answer></answer>.

<image>
{question}
{option}

ScienceQA

<image>
Context: {hint}
{question}
{choices}

Answer with the option's letter from the given choices directly.

RxnBench (EN)

<image>
Question: {question}
Choices:
A. {choice_A}
B. {choice_B}
C. {choice_C}
D. {choice_D}

Based on the image and the question, choose the most appropriate answer.

Only output a single letter (A, B, C, or D). Do NOT output any other text or explanation.

RxnBench (ZH)

<image>
问题: {question}
选项:
A. {choice_A}
B. {choice_B}
C. {choice_C}
D. {choice_D}

请根据图像和问题，从以上四个选项中选择最合适答案。

只输出单个字母（A, B, C 或 D），不要输出选项内容，也不要输出任何解释。

MolParse

<image>
{question}

OpenRxn

<image>
{question}

EMVista

<image>
{problem}

Answer with the option's letter from the given choices directly.

SuperChem (EN)

<image>
Question:
{question}
{option}

Select the best answer to the above multiple-choice question based on the image. Respond with only the letter of the correct option.

SuperChem (CN)

<image>
问题:
{question}
{option}

根据图像选择上述单选题的最佳答案。请仅输出正确选项的字母。

ProteinLMBench

Answer the multiple-choice question based solely on the provided context. If you are still unsure about the answer, output option 7.

Select only ONE correct option by its number. Start your response with 'The correct option is' followed by the option number ONLY.

Question: {question}

Options: {options}

The correct option is:

SFE-EN / SFE-ZH

multiple_choice_prompt:

<image>

You are an expert in {field} and need to solve the following question.

The question is a multiple-choice question. Answer with the option letter from the given choices.

{question}

{options}

exact_match_prompt:

<image>

You are an expert in {field} and need to solve the following question.

The question is an exact match question. Answer the question using a single word or phrase.

{question}

open-ended_prompt:

<image>

You are an expert in {field} and need to solve the following question.

The question is an open-ended question. Answer the question using a phrase.

{question}

MicroVQA

<image>

The following is a multiple choice question (with answers).

Think step by step and then output the answer in the format of "The answer is (X)" at the end, where X is the correct letter choice.

{question}

Options: {options}

MSEarth-MCQ

<image>

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

{question}

Based on the image, select the correct option (e.g., 'A', 'B', 'C') or directly state the correct option content.

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

```
  {"json": {  
    "answer": "Correct option or short answer",  
    "Explanation": "Reasoning explaining how to derive the correct answer."  
  }}
```

SmolInstruct

HIV

Analyze the SMILES string and predict its HIV activity. Answer strictly with “Yes” or “No”.

Input: {input}

Output:

BBBP

Does this molecule penetrate the Blood-Brain Barrier (BBBP)? Answer strictly with “Yes” or “No”.

Input: {input}

Output:

ClinTox

Predict the clinical toxicity of this molecule. Answer strictly with “Yes” or “No”.

Input: {input}

Output:

Side Effect

Does this drug cause the specific side effect? Answer strictly with “Yes” or “No”.

Input: {input}

Output:

Aqueous Solubility

As a specialized chemist, calculate the logSol (aqueous solubility) for the molecule below. Consider the molecular weight, number of rotatable bonds, and aromatic proportion. Provide the final logSol value as a floating point number. Output: [value]

Input: {input}

Output:

Lipophilicity

Analyze the lipophilicity (logP) of this molecular structure. Focus on the hydrophobic and hydrophilic balance. Provide the numerical logP value only. Output: [value]

Input: {input}

Output:

Forward Synthesis

Predict the major product(s) for these reactants. Output the product(s) in SMILES format.

Input: {input}

Output:

SmollInstruct

Retrosynthesis

Suggest the starting materials (reactants) for the given product. If multiple reactants are needed, separate them with a period (.). Output only the SMILES strings.

Input: {input}

Output:

SMILES to IUPAC

Convert this SMILES structure to its IUPAC name. Output only the IUPAC name, without any explanation or extra text.

Input: {input}

Output:

IUPAC to SMILES

Convert this IUPAC name to its SMILES structure. Output only the SMILES string, without any explanation or extra text.

Input: {input}

Output:

SMILES to Molecular Formula

Determine the molecular formula of this SMILES string. Output only the molecular formula, without any explanation or extra text.

Input: {input}

Output:

IUPAC to Molecular Formula

Determine the molecular formula of this IUPAC name. Output only the molecular formula, without any explanation or extra text.

Input: {input}

Output:

Captioning

Describe the following molecule's structure and chemical classification in detail.

Input: {input}

Output:

Generation

Generate the SMILES string for a molecule that fits this description.

Input: {input}

Output:

w/o Land Use Classification:

<image>

{question} The choices are listed below:

{option_A}

{option_B}

{option_C}

{option_D}

Select the best answer for the multiple-choice question based on the image. Only respond with the letter corresponding to the correct answer (A, B, C, D).

The answer is:

w/ Land Use Classification:

<image>

{question} The choices are listed below:

{option_A}

{option_B}

{option_C}

{option_D}

Select the best answer(s) for the multiple-choice question based on the image. There may be more than one correct option. Only respond with the letter(s) corresponding to the correct answer(s) (A, B, C, D), with multiple choices separated by spaces. The answer(s) is(are):