

# SciMDR: Advancing Scientific Multimodal Document Reasoning

Ziyu Chen<sup>C\*</sup>    Yilun Zhao<sup>Y\*</sup>    Chengye Wang<sup>Y</sup>    Rilyn Han<sup>Y</sup>  
Manasi Patwardhan<sup>T</sup>    Arman Cohan<sup>Y</sup>

<sup>Y</sup>Yale University    <sup>C</sup>University of Chicago    <sup>T</sup>TCS Research

## Abstract

Constructing scientific multimodal document reasoning datasets for foundation model training involves an inherent trade-off among scale, faithfulness, and realism. To address this challenge, we introduce the *synthesize-and-reground* framework, a two-stage pipeline comprising: (1) *Claim-Centric QA Synthesis*, which generates faithful, isolated QA pairs and reasoning on focused segments, and (2) *Document-Scale Regrounding*, which programmatically re-embeds these pairs into full-document tasks to ensure realistic complexity. Using this framework, we construct **SciMDR**, a large-scale training dataset for cross-modal comprehension, comprising 300K QA pairs with explicit reasoning chains across 20K scientific papers. We further construct **SciMDR-EVAL**, an expert-annotated benchmark to evaluate multimodal comprehension within full-length scientific workflows. Experiments demonstrate that models fine-tuned on **SciMDR** achieve significant improvements across multiple scientific QA benchmarks, particularly in those tasks requiring complex document-level reasoning.

 SciMDR     SciMDR

## 1 Introduction

While rapid publication accelerates the spread of ideas, it also makes it harder to locate the most consequential results and to integrate them into coherent understanding (Bornmann and Mutz, 2015; Kusumegi et al., 2025). LLMs and their multimodal counterparts (i.e., MLLMs) offer a promising way to navigate this flood of information, providing tools to quickly summarize, synthesize, and query scientific knowledge (Taylor et al., 2022; Luo et al., 2025). However, scientific papers remain difficult for general-purpose models because evidence is distributed across long, multimodal documents

(text, figures, and tables) and often requires domain expertise to interpret specialized terminology and connect claims to supporting context (Song et al., 2025; Wang et al., 2025; Zhao et al., 2025a). As a result, current models still struggle to provide reliable assistance in real scientific workflows (Zhao et al., 2025c; Tang et al., 2025; Xu et al., 2025).

A primary reason for this limitation is a deficit in high-quality training data that mirrors the complexity of real-world scientific inquiry. This data gap is reflected in the existing Scientific QA (SciQA) datasets. Early efforts rely on costly human annotation and remained small-scale and often text-only (Dasigi et al., 2021; Malaviya et al., 2024). Subsequent work turned to visual elements but adopted a *sanitized-context* approach, focusing on isolated figures or tables (Masry et al., 2022; Kahou et al., 2017). Recent work have begun to incorporate full-document contexts, presenting models with more realistic, *in-the-wild* tasks (Pramanick et al., 2024). This shift, however, has exposed a deeper, unresolved methodological challenge: a fundamental trade-off between *faithfulness* and *realism* in synthetic data. Specifically, to achieve *faithfulness*, QA generators can be prompted with concise, atomic contexts, which simplifies the task to yield verifiable outputs. However, this setup sacrifices realism as it leaves the generation pipeline underexposed to the full-length, complex documents. Conversely, to achieve *realism*, querying with lengthy, unprocessed documents can more closely mirror practical use cases. However, this long-context approach leads to attention dilution, increasing the likelihood of hallucinations (Ji et al., 2023) and undermining faithfulness in the generated ground-truth answers (Liu et al., 2024b; Bai et al., 2024).

As illustrated in Figure 1, to resolve this *faithfulness-realism dilemma*, we propose a new data synthesis paradigm that decouples faithfulness and realism across two stages. The first stage deliberately reduces data synthesis difficulty by structur-

\*Equal contributions. Correspondence to: Yilun Zhao (yilun.zhao@yale.edu)

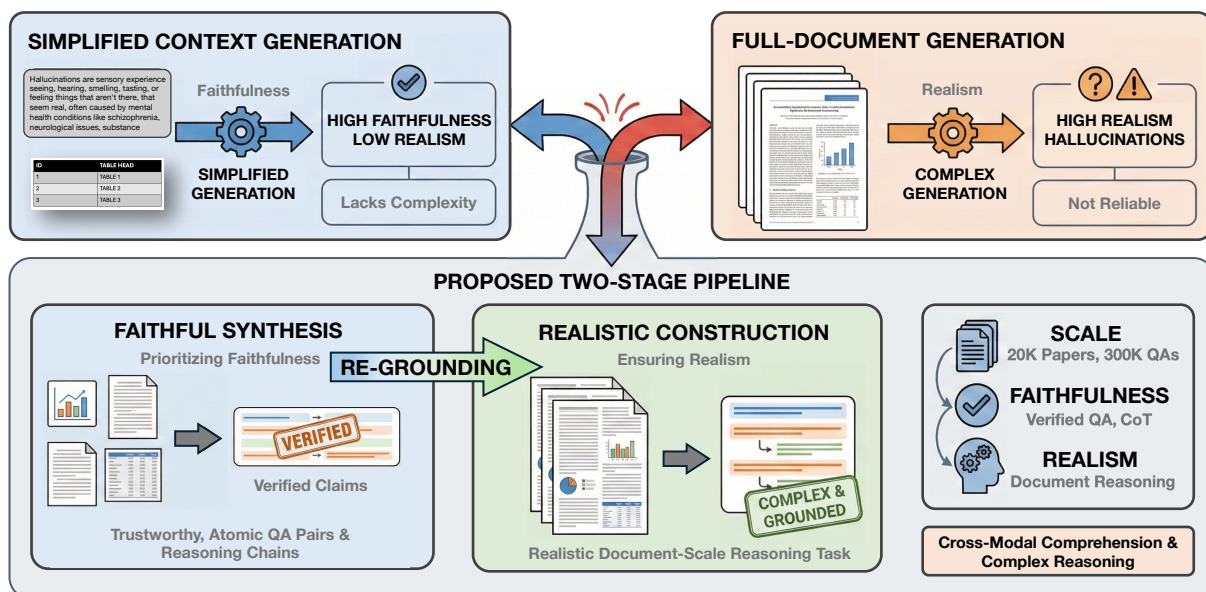


Figure 1: **The Faithfulness-Realism Dilemma in scientific data synthesis and our proposed solution.** Existing approaches face an inherent trade-off: simplifying context ensures *faithfulness* but lacks real-world complexity, while generating directly from full documents ensures *realism* but risks hallucination. We resolve this by decoupling the objectives into a two-stage *synthesize-and-reground* framework. By first generating verified QA pairs on atomic contexts and subsequently re-embedding them into full-document tasks, we achieve a dataset that simultaneously satisfies *Scale*, *Faithfulness*, and *Realism*.

ing synthesis around isolated, claim-centric units and a backward construction to ensure *faithfulness*, while the second stage reintroduces full-document complexity during training instance construction to achieve *realism*. Specifically, our approach first prioritizes *faithfulness* through synthesis stage. By operating on small, verifiable, and atomic contexts, this stage allows a generator to reliably produce grounded QA pairs and their detailed Chain-of-Thought (CoT) rationales (Wei et al., 2022). By constraining the core task and minimizing auxiliary demands, the generator is better positioned to produce trustworthy outputs. Second, we address *realism* via a training instance construction stage. We re-embed this golden QA-CoT pair within its original, full-document context. This design is the key to our solution: the model is presented with a realistic, *in-the-wild* task, but is simultaneously equipped with the precise CoT as ground truth. This demonstration teaches the model both *how to find the evidence* and *how to use evidence to answer questions*, bridging the gap between faithful synthesis and realistic application. Using this pipeline, we construct **SCIMDR**, a new, large-scale (300K QA pairs from 20K papers) dataset for multimodal scientific document reasoning, enabling models to be trained to help users understand central claims, supporting evidence, mechanisms, and comparisons

under realistic full-document conditions.

To comprehensively evaluate model performance in real-world scientific scenarios, we construct **SCIMDR-EVAL**, a benchmark comprising 907 human annotated QA pairs requiring evidence localization within lengthy, noisy documents, which further enables us to investigate the impact of long-context noise on model robustness. To validate our approach, we fine-tune Qwen2.5-VL-7B and LLaVA-1.5-7B on **SCIMDR**. Our empirical evaluation shows that this model significantly outperforms baselines across a comprehensive suite of three established benchmarks (i.e., ChartQA (Masry et al., 2022), CharXiv (Wang et al., 2024b) and SPIQA (Pramanick et al., 2024)) and **SCIMDR-EVAL**. Ablation studies confirm the value of the high-quality reasoning chains within our generated data, and experimental results validate that such data effectively teaches models the skills required for real-world scientific QA. Our main contributions are summarized below:

- We introduce **SCIMDR-EVAL**, an expert-annotated benchmark designed to evaluate model performance in realistic, in-the-wild scientific QA scenarios (§3).
- We propose a novel *synthesize-and-reground* paradigm that resolves *faithfulness-realism*

Category	Data	CoT	Q-Gen	Num QA	Source	Domain	Full Text	Visual
<b>Bench.</b>	QASPER (Dasigi et al., 2021)		human	5K	1.5K papers	NLP	×	×
	QASA (Lee et al., 2023)		human	1.8K	112 papers	AI/ML	✓	×
	ArgSciChat (Ruggeri et al., 2023)		human	41	20 papers	NLP	✓	×
	MMLongBench-Doc (Ma et al., 2024)	-	human + 11ms	2.5K	1612 charts	STEM	✓	✓
	CharXiv (Wang et al., 2024b)		human	11.5K	2.3K charts	STEM	×	✓
	ChartQAPro (Masry et al., 2025)		human + 11ms	1.9K	1.3K charts	STEM	×	✓
	DomainCQA (Zhong et al., 2025)		11ms	1.7K	482 charts	STEM	×	✓
<b>SciMDR-EVAL</b>	-	human	907	200 papers	STEM	✓	✓	
<b>DataSet.</b>	ChartQA (Masry et al., 2022)	×	human + 11ms	23K	28K charts	STEM	×	✓
	ArXivQA (Li et al., 2024b)	✓	GPT-4	100K	32K charts	STEM	×	✓
	MMSci (Li et al., 2024c)	×	GPT-4	1M	128K papers	STEM	×	✓
	SPIQA (Pramanick et al., 2024)	✓	human + 11ms	270K	25.5K papers	CS	×	✓
<b>SciMDR</b>	✓	GPT-5.1	300K	20K papers	STEM	✓	✓	

Table 1: **Comparison of Scientific QA Benchmarks & Datasets.** Unlike prior works that rely on sanitized contexts or lack reasoning annotations, **SciMDR** integrates *Full-Text* understanding, *Visual* modality, and explicit *chain-of-thought* reasoning at *scale*, bridging the gap between faithful synthesis and realistic document complexity.

*dilemma* in synthetic data generation by decoupling data generation from training instance construction, ensuring both atomic precision and holistic realism (§4).

- We release **SciMDR**, a large-scale high-faithfulness training dataset, by using the designed data synthesis pipeline (§4).
- Experiments show that fine-tuning on **SciMDR** improves scientific QA performance, and analyses further confirm that our data provides strong training signals for robust, in-the-wild multi-modal reasoning under long-context noise (§5).

## 2 Related Work

Crafting datasets to benchmark and enhance the scientific reasoning capabilities of LLMs necessitates a balance of three critical attributes: *scale*, *faithfulness*, and *realism*. However, achieving this balance presents a fundamental dilemma for prior work. As the general capabilities of LLMs have advanced, their expanding knowledge base offers opportunities for large-scale data synthesis. Yet, existing approaches often compromise one attribute to optimize the others, as summarized in Table 1.

**Human-Annotated SciQA.** Early scientific QA datasets relied on manual annotation to overcome the challenge of generating diverse, open-ended and domain-specific questions. Initial efforts like PubMedQA (Jin et al., 2019), BioASQ (Krithara et al., 2023), and QASPER (Dasigi et al., 2021) yielded thousands of examples but were often limited to abstracts or fixed formats. Subsequent work, such as QASA (Lee et al., 2023)

and Covid-QA (Möller et al., 2020), utilized full-text annotation for free-form questions, while ExpertQA (Malaviya et al., 2024), SCIDQA (Singh et al., 2024), and MISS-QA (Zhao et al., 2025b) further enhanced question complexity. While human annotation typically ensures quality, it faces a bottleneck in *scale*. The expensive nature of expert annotation limits these datasets’ size, making them insufficient for training modern foundation models that require vast quantities of data.

**Sanitized-Context SciQA.** With the development of visual capabilities in LLMs, attention has increasingly turned to the visual context within scientific documents, such as figures and tables. Datasets such as DVQA (Kafle et al., 2018), FigureQA (Kahou et al., 2017), PlotQA (Methani et al., 2020), ChartQA (Masry et al., 2022), and ChartQAPro (Masry et al., 2025) were proposed to benchmark with QA centered on visual contexts, placing new demands on the models’ visual understanding and reasoning. More recently, MathVista (Lu et al., 2023) and ArXivQA (Li et al., 2024b) have further broadened this task’s scope by incorporating more charts and diagrams. However, these datasets typically operate on *sanitized contexts*, isolating visual elements from their surrounding textual analysis. This approach creates a discrepancy between the benchmark task and the real-world challenge of navigating noisy, long-form documents. By simplifying the information retrieval process to isolated snippets, these methods compromise *realism*, failing to reflect the complexity of holistic scientific reasoning.

**Long-Context SciQA.** In real-world cases, users frequently query with long, complex documents.

Driven by the extension of context windows in LLMs (Team et al., 2024; Liu et al., 2025), many datasets have begun to focus on models’ ability to process and answer questions based on long-context. For instance, SciREX (Jain et al., 2020) is a document-level information extraction dataset, QuALITY (Pang et al., 2022) involves annotated QA over complete passages, and MMLongBench-Doc (Ma et al., 2024) and M3SciQA (Li et al., 2024a) incorporate visual information and multi-document reasoning through expert curation. The reliance on human annotators constrains the *scale* of these datasets. To address scalability, benchmarks like SPIQA (Pramanick et al., 2024), Loong (Wang et al., 2024a) and LongReason (Ling et al., 2025) typically synthesize questions based on short contexts, introducing extended noise documents during the evaluation. While providing final answers suffices for *benchmarking*, effective *training* demands explicit **reasoning** that guide models to locate evidence and filter noise. Originating from sanitized contexts, existing synthetic data inherently lacks these global traces, limiting its utility in enhancing *needle-in-a-haystack* reasoning capabilities.

### 3 SCIMDR-EVAL Benchmark

We focus on document-level scientific QA, where models must comprehend lengthy, multimodal documents in realistic scenarios. However, existing benchmarks mainly evaluate models on sanitized contexts—isolated figures, tables, or short passages. To bridge this gap and provide an evaluation of models’ capabilities in *in-the-wild* scientific reasoning, we construct **SCIMDR-EVAL**, an expert-annotated benchmark specifically designed to evaluate document-level multimodal QA performance. This benchmark serves dual purposes: (1) demonstrates the difficulty of *in-the-wild* scientific reasoning, and (2) provide a general, reliable testbed for evaluating multimodal document understanding in real-world scientific scenarios.

#### 3.1 Benchmark Construction

**SCIMDR-EVAL** is constructed through human annotation to ensure the quality and accuracy. We recruited three annotators (graduate students in computer science) to manually craft QA pairs from 300 scientific papers sourced from arXiv. To ensure coverage of scientific reasoning capabilities, we define five question types based on established

practices in scientific inquiry and our analysis of real-world SciQA requirements:

- *Evidence-Based Explanation & Quantification*: Explaining *how* and *why* visual element supports textual claim, often with quantitative analysis.
- *Concept-to-Instance Mapping*: Linking abstract concepts, architectures, or processes described in text to their concrete visual representations.
- *Hypothesis Validation & Inferential Reasoning*: Using textual and visual evidence to validate hypotheses, infer conclusions, or predict outcomes.
- *Critical Analysis & Consistency Check*: Critically evaluate whether textual claims are accurately supported by visual data, identifying potential inconsistencies or mischaracterizations.
- *Argumentative Role & Synthesis*: Synthesizing the overall scientific contribution and understanding the role of visual evidence in main argument.

For each assigned paper, the annotator was instructed to read the paper and formulate questions that necessitate synthesizing information across both textual content and visual elements distributed throughout the paper. Each entry was authored by one annotator and verified by the other two. Annotators were instructed to balance the questions across all types and provided with detailed guidelines and examples to ensure consistency and quality. Annotators also marked key points in each answer to facilitate fine-grained evaluation. This process yielded 907 high-quality QA pairs with detailed reasoning chains and answer key points for evaluation.

#### 3.2 Evaluation Protocol

Given the open-ended nature of our questions, exact-match and binary score might be inappropriate. Instead, we employ GPT-5-mini as an LLM judge to evaluate model responses. LLM-assisted evaluations are commonly used in many benchmarks (Lu et al., 2023; Yu et al., 2023; Wang et al., 2024b). The judge is provided with the question, annotated answer with key points, and response with reasoning chain. It assigns scores based on factual correctness, reasoning quality, and coverage of key points. We provide the implementation details in Appendix A.

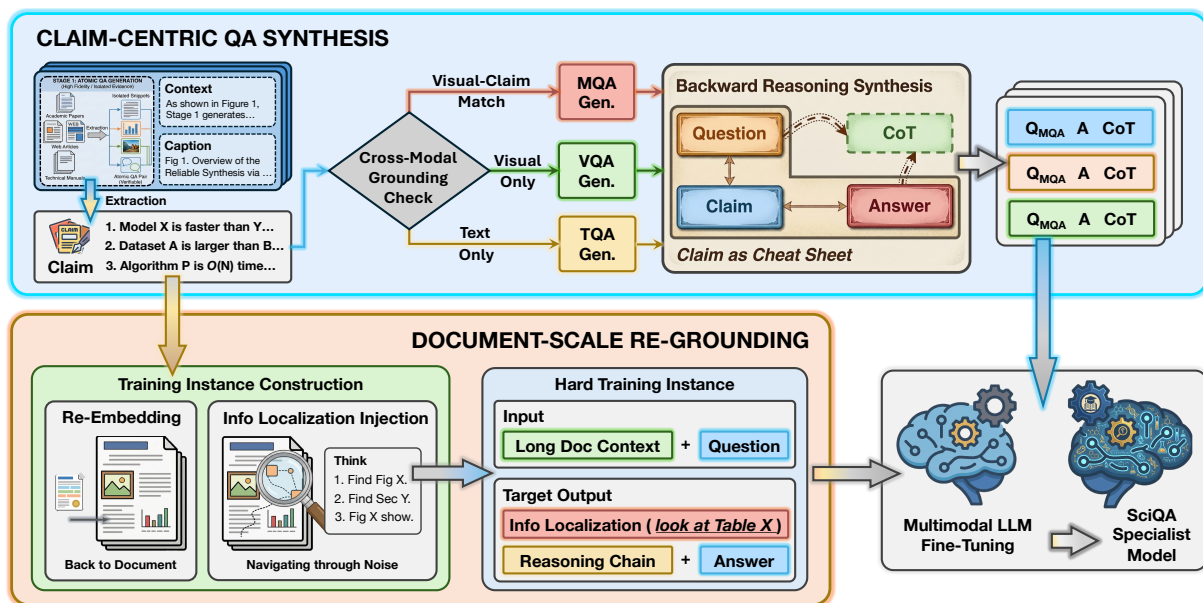


Figure 2: **Overview of the synthesize-and-reground framework.** The pipeline operates in two stages: *Claim-Centric QA Synthesis* ensures faithfulness by extracting atomic claims and employing backward reasoning to generate QA pairs with chain-of-thought; *Document-Scale Re-grounding* ensures realism by re-embedding these pairs into full-document contexts and injecting information localization steps to create hard training instances.

## 4 Training Data Synthesis Pipeline

To resolve the aforementioned *faithfulness-realism dilemma*, we introduce a two-stage paradigm that decouples data synthesis process from training instance construction, as outlined in Figure 2:

- **Claim-Centric QA Synthesis:** We first generate high-quality, trustworthy data by reducing the task difficulty for the generator model to ensure correctness and traceability.
- **Document-Scale Regrounding:** Then use this data to construct complex, realistic training instances for full-document comprehension.

This approach allows us to achieve all three goals: generated at scale, high-faithfulness content, and formatted for realistic, complex training.

### 4.1 Scientific Paper Collection and Processing

We collected raw academic papers from two primary sources to construct our dataset: CoRR in arXiv and Nature Communications. Papers from arXiv focus on the Computer Science, comprising a total of 9,847 papers ranging from 2017 to 2025. To ensure our dataset reflects the most recent research advancements, we prioritized papers from the last three years (2023–2025), which constitute over 97% of our arXiv subset. We also gathered

9,273 General Science articles from Nature Communications published between 2018 and 2025, ensuring a broad coverage of high-quality scientific content. To parse the multimodal content of each paper, we use the MinerU2.5 OCR model (Niu et al., 2025) with a vLLM backend. Given a downloaded PDF, our adopted OCR pipeline extracts the full body text, section boundaries, figures, tables, and associated captions. We serialize these outputs into JSON files, which are then used by the subsequent data synthesis pipeline. For each paper, we then use GPT-5.1 to assess whether it reports an original, experiment-driven study, filtering out surveys, position papers, tutorials, and purely conceptual work. Table 2 presents a detailed breakdown of the resulting dataset statistics.

### 4.2 Claim-Centric QA Synthesis

The objective of this stage is to produce a corpus of trustworthy, atomic QA pairs and their corresponding reasoning chains, all grounded in the source document. We achieve this quality by operating on small and isolated contexts, and employing a claim-centric mechanism. QA pairs can be classified into three types based on the information source required for an answer: **VQA** (*Vision-Only QA*), answerable solely from visual information (figures and tables); **TQA** (*Text-Only QA*), answerable solely from textual context; and **MQA** (*Multi-*

*modal QA*), which requires synthesizing information from both text and visuals. Each category is further defined by specific sub-types to balance generation diversity with controllability.

Our synthesis process begins with a multi-modal context unit, each comprising a segment of raw text, an associated visual (figure or table), and its caption. The core of this process is a claim-centric mechanism. We first perform a context-aware pre-processing step to identify all sentences within the text that reference the associated visual (e.g., *As shown in Figure X...*). We then feed the processed text into the LLM generator. At this time, the visual information is temporarily withheld to ensure a purely text-based analysis. Our prompt marks the previously identified referencing sentences, prioritizing segments most likely to contain arguments for later visual grounding. Following this guidance, the LLM generator breaks down the text into discrete, declarative claims, each representing a core finding or conclusion.

These extracted claims then serve as the unified *blueprint* for both QA and reasoning synthesis. First, there is a cross-modal grounding step, the LLM generator revises its claims by checking each one against the previously withheld visual information to determine whether a direct visual correlate exists. Claims with visual correlates are routed for MQA generation, text-only claims are routed for TQA, and VQA pairs are generated in parallel by focusing the LLM exclusively on the visual. Besides, for each QA pair, we guide the generation of its reasoning chain. We reframe this from an inference task to a low-risk, constrained articulation task. The claim is the key to this shift, acting as a cheating sheet with the ground-truth conclusion. By providing this answer upfront, we transform the task of LLM generator from finding an answer, to articulate a step-by-step rationale that logically connects a newly generated question to the supplied claim. This backward construction paradigm makes the synthesis easy by offloading tasks of evidence retrieval and open-ended inference, yielding reasoning chains both trustworthy and controllable.

### 4.3 Document-Scale Regrounding

The atomic QA pairs and reasoning chains generated on small and isolated contexts are suited for benchmarking a model’s capabilities, but they are suboptimal as training data. This is because, in realistic application scenarios, users rarely filter relevant paragraphs before posing a query. Instead,

the more common use case involves interrogating the entire noisy, complex document. Simply training on the atomic QA pairs would fail to prepare the model for this full-context challenge.

We bridge this gap by re-purposing claims. The claim, which served as a generation blueprint in synthesis, now functions as a ground-truth evidence map for the training stage. Because each QA pair is bound to a claim, which records the precise location of its textual and visual evidence, we can programmatically construct an ideal *Information Localization* step. This is achieved by populating a pre-defined template with the specific identifiers (e.g., *Section X, Table Y*) stored in the claim. This content, which explicitly states how to find the necessary information, is then prepended to the synthetic reasoning chains. For example: *To answer this question, I need to first consult Section X, and then cross-reference the results in Table Y...*

This deterministic synthesis of CoT rationales provides the downstream model with an accurate, verifiable, and imitable reasoning demonstration. This creates the hard training instance: the task is no longer a simple query on a filtered easy context, but a realistic, hard challenge that requires finding evidence within the full document. Critically, while the task difficulty is high, the solutions we provide via demonstrations are detailed and well-structured. With such data, the model is not just learning *what* the answer is; it is learning *how to find the answer* within a complex context. The final training data format is structured as: *(Full Document Context, Question) → (Information Localization + Reasoning + Final Answer)*. This format compels the model to first practice localizing related information and then execute grounded reasoning, thereby enhancing its practical utility in real-world scientific QA applications.

### 4.4 Human Evaluation of Synthetic Data

To directly assess the faithfulness and reliability of the synthesized training data with our pipeline, we conducted a manual evaluation on a random sample of 300 QA pairs from *SciMDR*. The sample consists of 100 instances from each category: VQA, TQA, and MQA. The sampled QA pairs were manually reviewed using two criteria: *correctness*, which measures whether the answer is factually accurate, and *relevance*, which captures whether the QA is properly grounded in the source document and associated visual evidence.

For VQA and TQA, all reviewed instances were

Part I: SciMDR (Training Dataset)		
Category	Description	Count
TQA	Answerable solely from textual context	47,389
VQA	Answerable solely from figures/tables	125,052
MQA	Requires synthesis of text and visuals	132,020
<b>Total</b>		<b>304,461</b>
Part II: SciMDR-EVAL (Benchmark)		
Type	Focus	Count
EEQ	Explanation & quantitative analysis	205
CIM	Linking abstract concepts to visuals	240
HVI	Inferential reasoning & prediction	244
CAC	Consistency check & critical evaluation	97
ARS	Synthesis of argument & visual role	121
<b>Total</b>		<b>907</b>

Table 2: Details about the SciMDR training dataset and the SciMDR-EVAL benchmark. SciMDR is categorized by modality, while SciMDR-EVAL is categorized by reasoning type, a taxonomy that also guides the synthesis of multi-modal samples in SciMDR.

factually correct and well grounded in the provided visual or textual context. For MQA, 91 out of the 100 sampled instances were judged to be high-quality and fully accurate. The remaining 9 instances were still factually correct, but were occasionally more verbose or focused on high-level concepts, resulting in longer and more complex reasoning chains. These results validate the effectiveness of our pipeline in generating faithful and well-grounded QA data.

## 5 Experiments

We conduct experiments to verify the effectiveness of our proposed data construction pipeline and SciMDR, addressing two research questions:

- **RQ1:** Does fine-tuning on SciMDR enhance model performance on scientific reasoning?
- **RQ2:** Does our synthetic data pipeline possess the capability to produce useful training data that improves model scientific reasoning?

### 5.1 Experimental Setup

**Dataset.** Our dataset SciMDR comprises three categories based on information sources: VQA, TQA, and MQA. The dataset was constructed following the pipeline in Section 4, generating approximately 300K QA pairs with claim-centric reasoning chains from 20K research papers with gpt-5.1.

**Training Configuration.** We employ a two-stage training, using Qwen2.5-VL-7B (Bai et al., 2025b) as our primary base model. In Stage 1, we train on VQA and TQA data for 1 epoch with a peak learning rate  $1 \times 10^{-5}$  and batch size 64. In Stage 2, we continue training on MQA data for 1 epoch with learning rate  $1 \times 10^{-6}$ . In fine-tuning with SPIQA, we train the language model for 1 epoch with a learning rate of  $1 \times 10^{-5}$  and batch size 64. We fine-tune the language model while keeping the visual encoder and projector frozen.

**Evaluation Benchmarks.** We evaluate models on four benchmarks: (1) **ChartQA** (Masry et al., 2022), a foundational chart QA benchmark to evaluate logical and visual reasoning over standard real-world charts; (2) **CharXiv** (Wang et al., 2024b), a benchmark for scientific QA that uses expert-curated charts from research papers to assess both *Descriptive examination* and complex *Reasoning capabilities*; (3) **SPIQA** (Pramanick et al., 2024), a benchmark with 3 subsets designed to assess multimodal comprehension of academic content, which requires a holistic understanding of complex figures and tables within full-text papers; and (4) **SciMDR-EVAL**, our annotated benchmark for full-document scientific reasoning.

**Baselines.** We benchmark our method against the base model Qwen2.5-VL-7B to measure relative gains, and reproduce SPIQA, a recent synthetic baseline, by fine-tuning the same base model to isolate data quality effects. We also include several strong open-source multimodal models Qwen-3-VL-8B (Bai et al., 2025a), LLaVA-OV-1.5-8B (An et al., 2025), and InternVL-3-8B (Zhu et al., 2025) as competitive references. In addition, we evaluate some advanced models GPT-4o (OpenAI, 2024), GPT-5.1 (OpenAI, 2025a), and GPT-5.2 (OpenAI, 2025b) on SciMDR-EVAL to establish a performance upper bound and analyze the development of scientific multimodal document reasoning capability.

### 5.2 Main Results

Table 3 presents the comparative performance of model fine-tuned with SciMDR, against the baselines across all four benchmarks. The results substantiate the efficacy of our approach (RQ1). *Model fine-tuned with SciMDR achieves substantial improvements over the base model across the board, effectively transforming a general-purpose multimodal model into a specialized scientific assis-*

Model	ChartQA	CharXiv-D	CharXiv-R	SPIQA-A	SPIQA-B	SPIQA-C	SciMDR-EVAL
GPT-5.1	-	90.9	58.3	79.4	79.8	71.6	47.2
GPT-5.2	-	95.2	73.1	79.9	75.4	74.0	49.9
Qwen-3-VL-8B	87.4	74.2	40.1	73.2	64.0	62.3	34.2
LLaVA-OV-1.5-8B	85.9	66.3	32.9	66.0	62.7	51.1	15.5
InternVL-3-8B	86.2	66.7	34.6	59.6	46.9	40.8	16.8
Qwen2.5-VL-7B	84.6	65.0	37.7	66.4	56.6	48.9	19.8
+ SPIQA	81.8 <sup>-2.8</sup>	50.9 <sup>-14.1</sup>	33.3 <sup>-4.4</sup>	62.7 <sup>-3.7</sup>	44.7 <sup>-11.9</sup>	40.0 <sup>-8.9</sup>	5.6 <sup>-14.2</sup>
+ SciMDR	86.3 <sup>+1.7</sup>	75.6 <sup>+10.6</sup>	37.9 <sup>+0.2</sup>	68.6 <sup>+2.2</sup>	58.8 <sup>+2.2</sup>	47.3 <sup>-1.6</sup>	49.1 <sup>+29.3</sup>

Table 3: **Main results on scientific QA benchmarks.** Fine-tuning with SciMDR outperforms the base model and the recent synthetic dataset across most metrics, particularly on complex reasoning tasks.

Model	SciMDR-EVAL
GPT-5.2	<b>49.9</b>
GPT-5.1	47.2
GPT-4o	24.7
Qwen2.5-VL-7B	19.8
+ SciMDR	49.1 <sup>+29.3</sup>

Table 4: **Performance comparison on SciMDR-EVAL against advanced models.** Despite having only 7B parameters, our model matches the performance of GPT-5.2 and GPT-5.1 on this domain-specific task.

tant. To further contextualize the difficulty of our proposed benchmark and the effectiveness of our method, we compare our fine-tuned model against advanced proprietary models on SciMDR-EVAL. Despite its smaller parameter size 7B, model with SciMDR exhibits competitive performance on this scientific reasoning task.

### 5.3 Pipeline Effectiveness and Analysis

Having established the performance gains, we address RQ2 by analyzing the quality of our synthetic data and deconstructing the contributions of our pipeline components.

#### 5.3.1 Data Quality Comparison

To assess the quality of our synthetic data independent of the base model’s intrinsic capabilities, we conduct a controlled comparison using LLaVA-1.5-7B (Liu et al., 2024a). We chose LLaVA-1.5 as our probing model for two strategic reasons: its fully transparent training data ensures no prior exposure to our evaluation benchmarks, and as a more modest baseline, it is more sensitive to data quality, allowing us to clearly observe the marginal gains from different instruction-tuning datasets. We fine-tune LLaVA-1.5-7B on three configurations: (1) 50K samples from SPIQA, (2) 50K VQA samples from SciMDR, and (3) 50K samples from SPIQA re-annotated using our claim-centric

Method	ChartQA	CharXiv	SPIQA-A
LLaVA-1.5-7B	19.6	27.8	31.5
+ SPIQA (50k)	26.3 <sup>+6.7</sup>	13.5 <sup>-14.3</sup>	35.7 <sup>+4.2</sup>
+ SciMDR (50k)	<b>26.8<sup>+7.2</sup></b>	<b>28.5<sup>+0.7</sup></b>	<b>36.7<sup>+5.2</sup></b>
+ SPIQA (re-annotated)	25.5 <sup>+5.9</sup>	28.1 <sup>+0.3</sup>	<b>39.8<sup>+8.3</sup></b>

Table 5: **Controlled data quality comparison.** The results base on LLaVA-1.5-7B show that re-annotating SPIQA’s source documents with our pipeline yields improvements over the original labels, confirming the superiority of our generation methodology.

pipeline. All models are trained for 2 epochs and evaluated on on single-image benchmarks to match the model’s input constraint.

Table 5 confirms that re-annotating SPIQA with our pipeline outperforms the original labels (39.8 vs. 35.7) using identical source documents. This isolates the gains to our methodology rather than data selection. We attribute this improvement to the rich reasoning signals in our data: notably, the model trained on our re-annotated SPIQA generates responses on CharXiv that are 5× longer on average than the original data, reflecting a substantial enhancement in reasoning depth and details.

#### 5.3.2 Ablation Study on Reasoning Chains

We further investigate which components of our training data contribute to full-document comprehension. Using the Stage 1 checkpoint, we evaluate three variants on SciMDR-EVAL: (1) full data with explicit information localization and reasoning chains, (2) removing localization, and (3) removing reasoning chains (QA pairs only).

Table 6 reveals that removing reasoning chains leads to a significant drop in performance (49.1 → 16.9), underscoring that simple QA pairs are insufficient for teaching complex scientific logic. Removing information localization also causes a drop, indicating that explicit guidance on *where* to look is important for helping models navigate the

Info Loc	Reasoning	SciMDR-EVAL
✓	✓	49.1
×	✓	22.8 <sup>-26.3</sup>
×	×	16.9 <sup>-32.2</sup>

Table 6: **Ablation study of training data components.** Both explicit information localization and step-by-step reasoning are critical for successful fine-tuning.

noise in full-text documents.

### 5.3.3 Impact of Long-Context Noise

Our pipeline is motivated by the observation that generating data directly from long, noisy contexts reduces faithfulness. To empirically quantify the impact of noise, we evaluate Qwen2.5-VL-7B on SciMDR-EVAL under three input settings: (1) **Oracle Context**, which provides only the ground-truth visual and referencing text with zero distractors; (2) **Standard Setting**, SciMDR-EVAL default which simulates realistic retrieval by including limited noise (maximum 8 images and 6 paragraphs); and (3) **Full-Paper**, which supplies the entire document content to maximize distractor density.

Table 7 reveals a clear performance degradation as noise increases. The gap between **Oracle Context** (32.9) and **Full-Paper** (12.8) confirms that long-context distractors are a source of error; even when the information is present, the model struggles to localize evidence within dense content.

### 5.3.4 Failure Analysis

We conduct a failure analysis of SciMDR-EVAL and comparing predictions from the base model and its fine-tuned counterpart on SciMDR. We observe four main error types:

- *Incorrect Evidence Localization*: Selecting the wrong visual element or paragraph instead of the true supporting context.
- *Reasoning / Logic Error*: Correctly locating relevant evidence but failing in multi-step deduction or computation.
- *Hallucination of Context*: Fabricating numbers, visual features, or statements not present in the document.
- *Incomplete Synthesis*: Identifying correct evidence but missing key annotated answer points.

Overall, the fine-tuned model shows clear improvements in grounding and evidence localization,

Input	SciMDR-EVAL
Standard	19.8
Oracle	32.9
Full-Paper	12.8

Table 7: **Challenge of Attention Dilution.** Effect of context noise on accuracy. Performance degrades as the amount of irrelevant context increases.

suggesting that the structured reasoning signals in SciMDR effectively can reduce hallucination and improve document-level reasoning. Details can be found in Appendix A.4.

## 6 Conclusion and Discussion

In this work, we addressed the *faithfulness-realism dilemma* in constructing synthetic datasets for multimodal scientific document reasoning. We introduced the *synthesize-and-reground* framework, which decouples atomic reasoning synthesis from full-document training. With SciMDR and SciMDR-EVAL, we demonstrate that our approach enables open-source models to bridge the performance gap with proprietary systems in complex multimodal document reasoning. Given reliance on proprietary models and STEM focus, future work will explore distilling synthesis into open-source models and expanding domains.

### Limitations

While our *synthesize-and-reground* framework effectively enhances scientific multimodal reasoning, several limitations remain. The fidelity of our training data is intrinsically bounded by the capabilities of the teacher model used for atomic synthesis. We assume that breaking the task into atomic claims minimizes hallucinations, yet any subtle factual errors or reasoning flaws generated at this stage are coded into the training signal. In practice, if the teacher model exhibits specific biases or misconceptions regarding niche scientific domains, these will inevitably propagate to the student model. Regarding the scope of our claims, our empirical validation is concentrated on papers within STEM disciplines (primarily Computer Science and General Science). This focus partly reflects the current scarcity of data resources outside of the hard sciences. Consequently, our results have not yet been validated in fields with distinct reasoning paradigms, where scientific discourse may follow different structures.

## References

- Xiang An, Yin Xie, Kaicheng Yang, Wenkang Zhang, Xiuwei Zhao, Zheng Cheng, Yirui Wang, Songcen Xu, Changrui Chen, Chun Yat Wu, Huajie Tan, Chunyuan Li, Jing Yang, Jiecao Yu, Xiyao Wang, Bin Qin, Yumeng Wang, Zizhen Yan, Ziyong Feng, and 3 others. 2025. *Llava-onevision-1.5: Fully open framework for democratized multimodal training*. *ArXiv*, abs/2509.23661.
- Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xiong-Hui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Rongyao Fang, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Qidong Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, and 46 others. 2025a. *Qwen3-vl technical report*. *ArXiv*, abs/2511.21631.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025b. *Qwen2.5-vl technical report*. *arXiv preprint arXiv:2502.13923*.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, and 1 others. 2024. Longbench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd annual meeting of the association for computational linguistics (volume 1: Long papers)*, pages 3119–3137.
- Lutz Bornmann and Rüdiger Mutz. 2015. Growth rates of modern science: A bibliometric analysis based on the number of publications and cited references. *Journal of the association for information science and technology*, 66(11):2215–2222.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. *arXiv preprint arXiv:2105.03011*.
- Sarthak Jain, Madeleine Van Zuylen, Hannaneh Hajishirzi, and Iz Beltagy. 2020. Scirex: A challenge dataset for document-level information extraction. *arXiv preprint arXiv:2005.00512*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):1–38.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pages 2567–2577.
- Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656.
- Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2017. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*.
- Anastasia Krithara, Anastasios Nentidis, Konstantinos Bougiatiotis, and Georgios Paliouras. 2023. Bioasqqa: A manually curated corpus for biomedical question answering. *Scientific Data*, 10(1):170.
- Keigo Kusumegi, Xinyu Yang, Paul Ginsparg, Mathijs de Vaan, Toby Stuart, and Yian Yin. 2025. Scientific production in the era of large language models. *Science*, 390(6779):1240–1243.
- Yoonjoo Lee, Kyungjae Lee, Sunghyun Park, Dasol Hwang, Jaehyeon Kim, Hong-in Lee, and Moontae Lee. 2023. Qasa: advanced question answering on scientific articles. In *International Conference on Machine Learning*, pages 19036–19052. PMLR.
- Chuhan Li, Ziyao Shangguan, Yilun Zhao, Deyuan Li, Yixin Liu, and Arman Cohan. 2024a. M3sciqqa: A multi-modal multi-document scientific qa benchmark for evaluating foundation models. *arXiv preprint arXiv:2411.04075*.
- Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024b. Multi-modal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv preprint arXiv:2403.00231*.
- Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoung Ji, Byungju Lee, Xifeng Yan, and 1 others. 2024c. Mmsci: A multimodal multi-discipline dataset for phd-level scientific comprehension. In *AI for Accelerated Materials Design-Vienna 2024*.
- Zhan Ling, Kang Liu, Kai Yan, Yifan Yang, Weijian Lin, Ting-Han Fan, Lingfeng Shen, Zhengyin Du, and Jiecao Chen. 2025. Longreason: A synthetic long-context reasoning benchmark via context expansion. *arXiv preprint arXiv:2501.15089*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306.
- Jiaheng Liu, Dawei Zhu, Zhiqi Bai, Yancheng He, Huanxuan Liao, Haoran Que, Zekun Wang, Chenchen Zhang, Ge Zhang, Jiebin Zhang, and 1 others. 2025. A comprehensive survey on long context language modeling. *arXiv preprint arXiv:2503.17407*.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024b. Lost in the middle: How language

- models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Ziming Luo, Zonglin Yang, Zexin Xu, Wei Yang, and Xinya Du. 2025. Llm4sr: A survey on large language models for scientific research. *arXiv preprint arXiv:2501.04306*.
- Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, and 1 others. 2024. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. *Advances in Neural Information Processing Systems*, 37:95963–96010.
- Chaitanya Malaviya, Subin Lee, Sihao Chen, Elizabeth Sieber, Mark Yatskar, and Dan Roth. 2024. Expertqa: Expert-curated questions and attributed answers. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3025–3045.
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Findings of the association for computational linguistics: ACL 2022*, pages 2263–2279.
- Ahmed Masry, Mohammed Saidul Islam, Mahir Ahmed, Aayush Bajaj, Firoz Kabir, Aaryaman Kartha, Md Tahmid Rahman Laskar, Mizanur Rahman, Shadikur Rahman, Mehrad Shahmohammadi, and 1 others. 2025. Chartqapro: A more diverse and challenging benchmark for chart question answering. *arXiv preprint arXiv:2504.05506*.
- Nitesh Methani, Pritha Ganguly, Mitesh M Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1527–1536.
- Timo Möller, Anthony Reina, Raghavan Jayakumar, and Malte Pietsch. 2020. Covid-qa: A question answering dataset for covid-19. In *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*.
- Junbo Niu, Zheng Liu, Zhuangcheng Gu, Bin Wang, Linke Ouyang, Zhiyuan Zhao, Tao Chu, Tianyao He, Fan Wu, Qintong Zhang, Zhenjiang Jin, Guang Liang, Rui Zhang, Wenzheng Zhang, Yuan Qu, Zhifei Ren, Yuefeng Sun, Yuanhong Zheng, Dongsheng Ma, and 42 others. 2025. [Mineru2.5: A decoupled vision-language model for efficient high-resolution document parsing](#). *Preprint*, arXiv:2509.22186.
- OpenAI. 2024. [Hello gpt-4o](#).
- OpenAI. 2025a. [GPT-5.1: A smarter, more conversational ChatGPT](#).
- OpenAI. 2025b. [Introducing GPT-5.2](#).
- Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen, Vishakh Padmakumar, Johnny Ma, Jana Thompson, He He, and 1 others. 2022. Quality: Question answering with long input texts, yes! In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5336–5358.
- Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. 2024. Spiga: A dataset for multimodal question answering on scientific papers. *Advances in Neural Information Processing Systems*, 37:118807–118833.
- Federico Ruggeri, Mohsen Mesgar, and Iryna Gurevych. 2023. A dataset of argumentative dialogues on scientific papers. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7684–7699.
- Shruti Singh, Nandan Sarkar, and Arman Cohan. 2024. Scidqa: A deep reading comprehension dataset over scientific papers. *arXiv preprint arXiv:2411.05338*.
- Zhangde Song, Jieyu Lu, Yuanqi Du, Botao Yu, Thomas M Pruyn, Yue Huang, Kehan Guo, Xiuzhe Luo, Yuanhao Qu, Yi Qu, and 1 others. 2025. Evaluating large language models in scientific discovery. *arXiv preprint arXiv:2512.15567*.
- Xiangru Tang, Zhuoyun Yu, Jiapeng Chen, Yan Cui, Daniel Shao, Weixu Wang, Fang Wu, Yuchen Zhuang, Wenqi Shi, Zhi Huang, and 1 others. 2025. Cellforge: agentic design of virtual cell models. *arXiv preprint arXiv:2508.02276*.
- Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Chengye Wang, Yifei Shen, Zexi Kuang, Arman Cohan, and Yilun Zhao. 2025. [SciVer: Evaluating foundation models for multimodal scientific claim verification](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8562–8579, Vienna, Austria. Association for Computational Linguistics.
- Minzheng Wang, Longze Chen, Fu Cheng, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan Xu, Lei Zhang, Run Luo, and 1 others. 2024a. Leave

- no document behind: Benchmarking long-context llms with extended multi-doc qa. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5627–5646.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, and 1 others. 2024b. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Zhijian Xu, Yilun Zhao, Manasi Patwardhan, Lovekesh Vig, and Arman Cohan. 2025. [Can LLMs identify critical limitations within scientific research? a systematic evaluation on AI research papers](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 20652–20706, Vienna, Austria. Association for Computational Linguistics.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.
- Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2024. [Lmms-eval: Reality check on the evaluation of large multimodal models](#). *Preprint*, arXiv:2407.12772.
- Yilun Zhao, Weiyuan Chen, Zhijian Xu, Manasi Patwardhan, Chengye Wang, Yixin Liu, Lovekesh Vig, and Arman Cohan. 2025a. [AbGen: Evaluating large language models in ablation study design and evaluation for scientific research](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12479–12491, Vienna, Austria. Association for Computational Linguistics.
- Yilun Zhao, Chengye Wang, Chuhan Li, and Arman Cohan. 2025b. [Can multimodal foundation models understand schematic diagrams? an empirical study on information-seeking QA over scientific papers](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 18598–18631, Vienna, Austria. Association for Computational Linguistics.
- Yilun Zhao, Kaiyan Zhang, Tiansheng Hu, Sihong Wu, Ronan Le Bras, Yixin Liu, Xiangru Tang, Joseph Chee Chang, Jesse Dodge, Jonathan Bragg, Chen Zhao, Hannaneh Hajishirzi, Doug Downey, and Arman Cohan. 2025c. [Sciarena: An open evaluation platform for non-verifiable scientific literature-grounded tasks](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. [Llamafactory: Unified efficient fine-tuning of 100+ language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.
- Ling Zhong, Yujing Lu, Jing Yang, Weiming Li, Peng Wei, Yongheng Wang, Manni Duan, and Qing Zhang. 2025. [Domaincqa: Crafting expert-level qa from domain-specific charts](#). *arXiv preprint arXiv:2503.19498*.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Yue Cao, Yangzhou Liu, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Han Lv, and 29 others. 2025. [InternV3: Exploring advanced training and test-time recipes for open-source multimodal models](#). *ArXiv*, abs/2504.10479.

## A Data and Experimental Details

### A.1 Configuration

**Qwen2.5-VL-7B.** We fine-tuned Qwen2.5-VL-7B using LLaMA-Factory (Zheng et al., 2024) with the following configurations. The maximum sequence length was set to 16K tokens (including both visual and language tokens) to accommodate long-context scientific documents. For image inputs, we set a maximum of 8 images per instance with  $\text{max\_pixels} = 512 \times 512$ . Images are automatically resized to maintain their aspect ratio within the specified pixel range.

For *VQA + TQA*. We trained on visual-only and text-only QA pairs for 1 epoch with learning rate  $1 \times 10^{-5}$  and batch size 64. Only the language model was trained while the visual encoder and projector remained frozen.

For *MQA*. We continued training on multimodal QA pairs for 1 epoch with learning rate  $1 \times 10^{-6}$  and batch size 64, maintaining the same freeze strategy.

**LLaVA-1.5-7B.** For data quality comparison experiments, we fine-tuned LLaVA-1.5-7B using LLaMA-Factory for 2 epochs with learning rate  $1 \times 10^{-5}$ , batch size 64, and warmup ratio 0.1. Unlike the Qwen experiments, all model components (vision encoder, projector, and language model) were trained without freezing.

### A.2 Evaluation Framework

All evaluations were conducted using `lmms-eval` (Zhang et al., 2024), which provides standardized evaluation protocols for large multimodal models. We implemented a custom evaluation module for *SciMDR-EVAL* to ensure consistency with existing benchmarks.

### A.3 LLM Judge Configuration

Given the open-ended nature of questions in *SciMDR-EVAL*, we employed GPT-5-mini as an LLM judge to evaluate model responses. The judge assesses each response based on factual correctness, reasoning quality, and coverage of annotated key points.

**Binary Scoring.** For main results shown in Table 3, we use strict binary scoring: a response receives score 1 only if it correctly addresses all key points with accurate reasoning; otherwise it receives 0. The accuracy is computed as the percentage of fully correct responses.

**Fine-grained Metrics.** For detailed analysis, we also report text correctness rate (percentage correctly interpreting textual evidence), visual correctness rate (percentage correctly interpreting visual evidence), and partial credit score (average proportion of key points addressed). These fine-grained metrics provide additional insight but are not used for main benchmark comparison. The complete judge prompt is provided in Figure 4.

### A.4 Failure Mode Analysis

**Setup.** To analyze failure patterns, we randomly sampled 100 questions from *SciMDR-EVAL* and compared outputs from the base model Qwen2.5-VL-7B and the fine-tuned on *SciMDR*. We manually categorized incorrect predictions into four error types.

**Failure Categories.** We define the following error categories:

- *Incorrect Evidence Localization*: Selecting the wrong visual element or paragraph instead of the true supporting context.
- *Reasoning / Logic Error*: Correctly locating relevant evidence but failing in multi-step deduction or computation.
- *Hallucination of Context*: Fabricating numbers, visual features, or statements not present in the document.
- *Incomplete Synthesis*: Identifying correct evidence but missing key annotated answer points.

**Findings.** Both quantitative error analysis and qualitative inspection demonstrate that the structured reasoning signals in *SciMDR* are important for improving multimodal document-level scientific QA. The fine-tuned model benefits from explicit localization supervision and exhibits stronger grounding behavior compared to the base model.

Failure Type	Qwen	SciMDR
Incorrect Evidence Localization	18	5
Reasoning / Logic Error	6	9
Hallucination of Context	11	3
Incomplete Synthesis	8	7
Total Errors	43	24

Table 8: Failure type comparison on 100 randomly sampled *SciMDR-EVAL* questions.

## B Annotator and Data Usage

**Annotator Recruitment.** For constructing *SciMDR-EVAL*, we recruited three graduate students in Computer Science with at least one year of experience in machine learning research and scientific paper analysis. Annotators were compensated above local minimum wage, consistent with standard research assistant rates. All annotators provided written informed consent before participating.

**Consent and Usage Rights.** Prior to annotation, participants received detailed consent forms explaining the research purpose, public data release, withdrawal rights, confidentiality measures, and compensation structure.

For source papers in *SciMDR* and *SciMDR-EVAL*, we exclusively used open-access publications from arXiv (various Creative Commons licenses) and Nature Communications (CC-BY license). These licenses permit text and data mining for research purposes, requiring no additional consent from paper authors.

**Quality Control.** To ensure annotation quality, annotators underwent training with detailed guidelines and examples. Each QA pair was authored by one annotator and verified by the other two. Weekly meetings addressed challenging cases and maintained consistency.

**Annotation Cost.** For pre-annotation, the total setup time for the annotation team was approximately 5 hours, which included designing guidelines, creating samples, and conducting a training session to align the annotators with the protocol. For annotation, the average time for reading and annotating a single paper was approximately 10 minutes.

## C Data Synthesis Prompts

This section presents the complete prompts used in our data synthesis pipeline, corresponding to the stages described in Section 4.

**Claim Extraction.** Figure 5 shows the prompt that guides the LLM to distill paragraphs into structured, verifiable claims serving as blueprints for QA generation.

**Visual Grounding.** Figure 6 presents the prompt for matching textual claims with visual evidence and determining their relationship types.

**Multimodal QA Generation.** Figure 7 details the prompt for generating questions requiring synthesis of textual and visual information across five reasoning types (EEQ, CIM, HVI, CAC, ARS).

**Visual-Only QA Generation.** Figure 8 provides the prompt for generating questions answerable solely from visual information across eight reasoning categories.

**Text-Only QA Generation.** Figure 3 shows the prompt for generating questions testing deep understanding of scientific content without visual evidence.

## D MQA Examples

This section presents examples of multimodal QA pairs across the five question types.

**Evidence-Based Explanation & Quantification (EEQ).** Figure 9 illustrates an EEQ-type question requiring quantitative analysis of visual evidence to support textual claims.

**Concept-to-Instance Mapping (CIM).** Figure 10 shows a CIM-type question that links abstract architectural concepts described in text to their concrete visual representations in diagrams.

**Hypothesis Validation & Inferential Reasoning (HVI).** Figure 11 presents an HVI-type question demonstrating inferential reasoning by synthesizing visual patterns and textual explanations to draw conclusions.

**Critical Analysis & Consistency Check (CAC).** Figure 12 provides a CAC-type question that critically evaluates the consistency between textual characterizations and visual data.

**Argumentative Role & Synthesis (ARS).** Figure 13 displays an ARS-type question requiring synthesis of visual evidence and textual arguments to understand the overall scientific contribution.

**System:** You are an expert AI Research Assistant specializing in scientific text analysis and reasoning. Your goal is to generate challenging, high-quality questions based on textual claims that do not have direct visual evidence.

**Task**

For each claim, generate a question that requires deep understanding of the textual content and the scientific claims presented. The question should test comprehension, analysis, and reasoning about the scientific concepts, methodologies, or findings described in the text.

**Inputs**

1. **Textual Claims:** A list of claims without visual grounding that serve as the basis for question generation. {claims\_json}

2. **Context:** {context}

**Methodology: 3-Step Chain of Thought**

**Step 1: Claims Analysis**

Review all textual claims and identify the key scientific concepts, methodologies, or findings for each claim. Determine what makes each claim scientifically significant or interesting. Identify potential areas of complexity or nuance that would make good question topics.

**Step 2: Question Strategy Selection**

Choose the most appropriate approach for each question:

- **Conceptual Understanding:** Test deep understanding of scientific concepts or theories
- **Methodological Analysis:** Focus on experimental design, procedures, or analytical approaches
- **Critical Evaluation:** Assess the validity, limitations, or implications of findings
- **Comparative Analysis:** Compare different approaches, results, or interpretations
- **Causal Reasoning:** Explore cause-and-effect relationships described in the text

**Step 3: Question and Answer Generation**

Formulate a clear, specific question that requires synthesis of information from multiple claims for each claim. Ensure the question cannot be answered by simply restating a single claim. Do not give implications in the problem description: Your question should not imply that it is related to specific information. Let the respondent judge which information is relevant based on the questions. You can only ask one specific question for each claim and it cannot contain multiple sub-questions. Avoid using commas to connect multiple different questions in a single sentence (like "What ..., and how ... ?", should be avoided).

Generate a comprehensive answer with step-by-step reasoning. When you answer questions, remember you are the questioner, and you need to build a good example as a reference answer. Therefore, you can use Verified Textual Claims as shortcuts, which will help you understand accurate way of thinking when answering questions. But as you are building a reference answer for respondent, so when you are organizing the language of your answer, please do not reveal shortcuts. You should simulate a clever respondent who is good at answering questions, only using the original text in the answering process.

**Output Format**

```
[
  {
    "question_type": "The strategy you selected
      from the list",
    "question": "The final, challenging question
      based on the textual claims",
    "answer": {
      "understand_question": "Hint: think and plan",
      "chain_of_thought_answer": [
        {
          "step": 1,
          "reasoning": "Hint: Analysis of relevant claims
            and their scientific significance"
        },
        {
          "step": 2,
          "reasoning": "Hint: Deeper exploration of
            the concepts, methods, or findings"
        },
        {
          "step": 3,
          "reasoning": "Hint: Synthesis and conclusion
            that directly answers the question"
        }
      ],
      "conclusion": "Hint: The final, concise answer"
    }
  }
]
```

Figure 3: **TQA generation prompt.** This prompt generates questions testing deep understanding of scientific content without visual evidence.

**System:** You are an expert evaluator for multi-modal question-answering tasks. Evaluate the model's prediction based on the following three-component rubric.

**Question**

{question}

**Ground Truth Answer**

{ground\_truth\_answer}

**Model Prediction**

{prediction}

**Scoring Rubric (Total: 1.0 point)**

Evaluate the model's response across three dimensions:

**1. Text Citation Score (0.30 points)**

Evaluate whether the model accurately found and cited relevant textual content:

- **0.30 points:** The model accurately identified and cited all relevant text passages that fully support the answer
- **0.20 points:** The model identified and cited most relevant text passages, with minor omissions
- **0.10 points:** The model cited some relevant text but missed many key passages or included significant irrelevant text
- **0.0 points:** The model failed to identify or cite relevant textual content, or only cited irrelevant text

**2. Image Citation Score (0.30 points)**

Evaluate whether the model accurately identified and referenced relevant images:

- **0.30 points:** The model accurately identified and referenced all relevant images needed to answer the question
- **0.20 points:** The model identified and referenced most relevant images, with minor omissions
- **0.10 points:** The model referenced some relevant images but missed many or included significant irrelevant images
- **0.0 points:** The model failed to identify or reference relevant images, or only referenced irrelevant images

**3. Answer Accuracy Score (0.40 points)**

Evaluate whether the model correctly answered the key points of the question:

- **0.40 points:** The model's answer correctly addresses all key points and matches the ground truth
- **0.20 points:** The model's answer partially addresses the question but misses some key points or contains minor errors
- **0.0 points:** The model's answer is incorrect or fails to address the key points of the question

Figure 4: **LLM judge prompt.** This prompt evaluates model responses based on text citation (0.30), image citation (0.30), and answer accuracy (0.40).

**System:** You are a senior scientific editor with exceptional skills in logical analysis and information synthesis. Your task is to read a paragraph from a research paper and distill its core information into structured **CLAIMS**.

#### Core Principles

- **Summarization Over Detail:** Extract significant, high-level assertions. A claim should be a complete, conclusive statement.
- **Filter Non-Essentials:** Ignore introductory phrases, transitional sentences, or overly granular details.
- **Fidelity to Source:** Every claim must be strongly and directly supported by the original text.

#### Chain of Thought

Follow these steps internally before generating the final output:

- **Read for Gist:** Read the entire paragraph to understand its overall purpose and main argument.
- **Identify the Core:** Locate the key sentences that form the "skeleton" of the argument.
- **Synthesize and Refine:** Merge and summarize sentences into a single, concise content statement.
- **Classify and Populate:** Choose the most appropriate `claim_type` and fill in the required JSON fields.

#### Hints

- When selecting a `claim_type`, refer to these definitions:
  - **Conclusion/Finding:** The central takeaway or primary discovery.
  - **Comparative:** Compares advantages, disadvantages, or performance of entities.
  - **Causal/Explanatory:** Describes the reason for a phenomenon or mechanism of a process.
  - **Methodological:** Describes experimental setup, architecture, or operational steps.
  - **Descriptive/Quantitative:** Provides objective description or specific data point/metric.
- Sentences that reference a Figure or Table often contain important claims.

#### Input

**Context:** {context}

#### Output Requirements

Provide output in the following JSON format. You can generate multiple claims if there are multiple claims in the paragraph.

```
[
  {
    "id": "T_claim_01",
    "claim_type": "Choose from ['Conclusion/Finding',
      'Comparative', 'Causal/Explanatory',
      'Methodological', 'Descriptive/Quantitative']",
    "content": "Concise, clear, and complete summary.",
    "key_entities": ["List 1-3 core entities"],
    "source_sentences": ["Original supporting sentences"],
    "potential_question_focus": "Central question
      this claim is best suited to answer."
  }
]
```

Figure 5: **Claim extraction prompt.** This prompt guides the LLM to distill paragraphs into structured, verifiable claims serving as blueprints for QA generation.

**System:** You are a meticulous Multimodal Analyst and Fact-Checker. Your expertise is in precisely matching textual claims with visual evidence.

**Task**

Your mission is to receive structured textual claims (JSON format) and, for **each claim**, systematically find corresponding visual evidence within a given image and caption. You will annotate findings by adding a `visual_grounding` field to each object.

**Inputs**

1. **Structured Textual Claims (JSON):** {text\_json}
2. **Caption:** {caption}

**Process**

For **each claim object**, strictly follow this procedure:

1. **Analyze the Claim:** Read the claim's content and `key_entities` to understand its core assertion.
2. **Locate Visual Elements:** Search the image and caption for relevant visual elements (e.g., lines in a graph, bars in a chart, labels, specific regions, or text).
3. **Critically Evaluate the Evidence:** Is the visual element **directly** and **explicitly** related to the claim? Thematic relevance alone is insufficient. If evidence exists, determine its nature: Does it support, quantify, illustrate, or contradict the claim?
4. **Construct the visual\_grounding Object:** If you **cannot find** direct visual evidence, construct {"exists\_in\_visual": false}. If you **find** direct evidence, construct a complete object with "exists\_in\_visual": true along with `relationship_type`, `visual_element_description`, and `justification`.

**Output Requirements**

Output a single, complete, augmented JSON object. Ensure the original JSON's structure is fully preserved. The only modification should be the addition of the `visual_grounding` field to every original claim object.

**visual\_grounding Field Definitions**

- `exists_in_visual`: (Boolean) Required. true if direct visual evidence is found, otherwise false.
- `relationship_type`: (String) Required if `exists_in_visual` is true. Must be one of: ["Supports", "Quantifies", "Illustrates", "Elaborates", "Contradicts"].
- `visual_element_description`: (String) Required if `exists_in_visual` is true. A clear, textual description of the relevant visual element(s).
- `justification`: (String) Required if `exists_in_visual` is true. A detailed explanation of how the visual element is linked to the textual claim.

**relationship\_type Definitions**

- **Supports:** Visual evidence directly confirms or strengthens the assertion.
- **Quantifies:** Visual evidence provides specific numerical data or measurements.
- **Illustrates:** Visual evidence provides a visual example or representation.
- **Elaborates:** Visual evidence provides additional details or context.
- **Contradicts:** Visual evidence directly refutes or opposes the assertion.

Figure 6: **Visual grounding prompt.** This prompt matches textual claims with visual evidence, determining relationship types (Supports, Quantifies, Illustrates, Elaborates, Contradicts).

**System:** You are an expert AI Research Assistant specializing in multimodal scientific reasoning. Your goal is to generate challenging, high-quality multimodal questions for each claim with visual grounding.

### Question Categories

#### 1. Evidence-Based Explanation & Quantification (EEQ)

Core: Explain HOW and WHY a visual element supports a textual claim, and quantify that support. Example: "The authors claim that [statement]. How exactly does the data in [Figure/Table X] support this claim, and can you quantify the effect?"

#### 2. Concept-to-Instance Mapping (CIM)

Core: Link an abstract concept, architecture, or process described in text to its concrete visual representation. Example: "The paper defines '[concept]' in Section X. Identify the corresponding components in [Figure Y] and explain how they match the description."

#### 3. Hypothesis Validation & Inferential Reasoning (HVI)

Core: Use combined evidence from text and visuals to validate a hypothesis, infer conclusions, or predict outcomes. Example: "The hypothesis is that [hypothesis]. How do the results in [Figure X], combined with the text's interpretation, validate this hypothesis?"

#### 4. Critical Analysis & Consistency Check (CAC)

Core: Critically evaluate whether textual claims are accurately supported by visual data. Example: "The text describes the improvement in [Figure X] as 'significant'. Based on the visual evidence and scale, is this characterization accurate?"

#### 5. Argumentative Role & Synthesis (ARS)

Core: Summarize the overall scientific takeaway and the specific role of visual evidence in the paper's main argument. Example: "What is the core scientific takeaway from the combination of [Figure X] and its description in the text?"

### Task

For each claim, generate one question requiring deep, integrated understanding of THREE sources: (1) Figure/Table image, (2) Caption, (3) Context text. The question MUST NOT be answerable by looking at only one or two sources in isolation.

### Inputs

1. **Verified Textual Claims:** {claims\_json}
2. **Context:** {context}
3. **Caption:** {caption}

### Guiding Philosophy

Think Like a Scientist, Not a Grader. Your primary goal is to frame questions from the perspective of a curious scientist trying to understand the *phenomenon*, not a teacher asking a student to interpret a chart. The question should be about the 'what', 'why', or 'how' of the scientific concept itself. The text and visual are simply the *resources* the user needs to construct the answer.

### Methodology: 4-Step Chain of Thought

#### Step 1: Deep Analysis & Synthesis

Review all claims. What is the central scientific story they tell with the visual? Analyze the visual, caption, and context. Identify key relationships, trends, and mechanisms. Synthesize sources to form a high-level understanding.

#### Step 2: Strategic Question Type Selection

Based on synthesis, choose the MOST appropriate question sub-type (EEQ, CIM, HVI, CAC, or ARS) that would lead to a challenging and insightful question.

#### Step 3: Draft the Final Multimodal Question

**Core Principle:** Ask about the concept, not the evidence. Frame the question around the scientific idea, process, or implication. The user will naturally have to find the evidence in the provided materials.

#### Crucially, re-examine your question using these examples:

- **BAD:** "How does the diagram illustrate that you must visit all target regions?"
- **GOOD:** "Why must all potential target regions be searched before the target's precise location can be confirmed?"

#### Use Diverse Framing Strategies:

- Causal Reasoning (Why?): "What is the primary reason for the observed drop in efficiency when temperature exceeds 50°C?"
- Procedural Understanding (How?): "If a new data point were introduced, walk through the steps the algorithm would take to classify it."
- Hypothetical Scenarios (What if?): "What would be the likely impact if the 'attention' module were removed?"
- Defining Conditions (Under what...): "For which demographic does the interface show the most significant usability improvement?"
- Integrative Synthesis (What implication?): "What is the overarching conclusion regarding system stability when combining error rates in Table 1 with the risk assessment in the conclusion?"
- Critical Discrepancy (What inconsistency?): "What specific data behaviors contradict the author's claim of a 'linear relationship'?"

**Final Check:** Read aloud. Does it sound natural? It must not contain words like "figure," "table," "diagram," "graph," or "according to the text."

#### Step 4: Generate a Conclusive Answer

Provide a comprehensive, step-by-step answer. Your reasoning must explicitly cite evidence from the visual, caption, and context. When describing visual information, meticulously detail all relevant visual elements and explicitly connect them to the question.

Figure 7: **MQA generation prompt.** This prompt generates questions requiring synthesis of textual and visual information across five reasoning types (EEQ, CIM, HVI, CAC, ARS).

**System:** You are an expert Visual Question-Answering pair generator. Your task is to analyze a provided image context (e.g., chart, graph, table) and a specified question category to create a structured JSON object.

**Question Categories**

1. **DR, Data Retrieval (Lookup):** Asks for a specific data point that can be directly read from the image.
2. **EI, Extremum Identification (Max/Min):** Asks to find the highest, lowest, largest, or smallest value, or the entity associated with it.
3. **CO, Computation:** Requires a mathematical calculation (e.g., sum, difference, average, percentage change) based on data points from the image.
4. **CT, Counting:** Requires counting the number of elements that meet a specific numerical criterion.
5. **CR, Comparison & Ranking:** Requires comparing two or more data points or finding an entity with a specific rank.
6. **TP, Trend & Pattern Analysis:** Focuses on overall behavior of data over time, correlations, or specific patterns.
7. **IP, Inference & Prediction:** Asks for a projection, estimation based on a trend, or hypothetical outcome.
8. **MS, Compositional Reasoning (Multi-Step):** A complex question that requires combining two or more of the above types.

**Task**

Based on the provided image context and the specified question category, generate one QA pair.

**Question Category:** Choose the MOST appropriate question sub-type that would lead to a challenging and insightful question. {VISUAL\_ONLY\_QUESTION\_CATEGORY}

**Caption:** {caption}

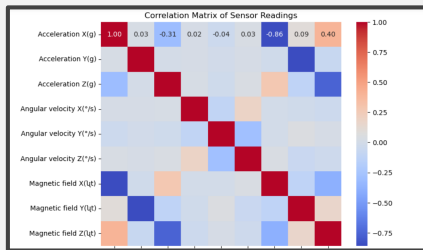
**Rules**

1. **Question Generation:** The generated question must be relevant to the specified category and must be answerable *solely* by analyzing the visual information in the image context.
2. **Global Image Description:** First, give a comprehensive and detailed description of what you see in the image. Describe the type of visualization, its main components, labels, colors, layout, values, the magnitude and positional relationships of values of each element, and any important visual elements.
3. **Relevant Parts of Image:** Connect the image description to the specific question being asked. Identify which parts of the image are relevant to answering the question.
4. **Step-by-Step Reasoning:** Provide step-by-step reasoning to find the answer. Each step should build on the previous one.
5. **Answer:** State the final answer clearly in a single, complete sentence.
6. **Short Form Answer:** Provide a concise version of the answer, typically a number, word, or short phrase, suitable for automated evaluation.
7. **JSON Structure:** Your final output **MUST** be a single, raw JSON object strictly adhering to the following structure.

**Output Format**

```
[
  {
    "question_type": "Select from [DR, EI, CO, CT,
      CR, TP, IP, MS]",
    "question": "The question you generated",
    "global_image_description": "...",
    "relevant_parts_of_image": "...",
    "step_by_step_reasoning": "...",
    "answer": "A full-sentence answer",
    "short_form_answer": "A concise, short-form answer"
  }
]
```

Figure 8: **VQA generation prompt.** This prompt generates questions answerable solely from visual information across eight reasoning categories.



How does the pattern of low correlation among different sensor axes imply that each sensor can provide unique and complementary information for identifying activities? From both statistical interpretation and the empirical evidence.

**From Visual**, most off-diagonal squares are pale blue or pale red, corresponding to correlation values near zero (between -0.3 and 0.4). Only a few variable pairs show moderate correlations, such as Acceleration X and Magnetic field X (around -0.86) or Magnetic field X and Magnetic field Z (about 0.40). The caption confirms it is a correlation matrix, thus these colors quantitatively represent correlation coefficients between -1 and 1. The general light tones indicate weak relationships.

**The context text states that** low correlation values between variables indicate minimal redundancy, meaning no single sensor axis is overly predictable from another. This property allows each variable to possibly capture distinctive aspects of the motion or environment. The text also emphasizes that these sensor readings could provide unique information for activity recognition.

Integrating the heatmap patterns (mostly low correlations), the caption labeling the plot as a correlation matrix, and the textual description of the implications, we can conclude that the low correlations indicate each axis measures relatively independent aspects of movement or field disturbances. Low redundancy means that combining them could enrich an activity recognition model since each sensor's data adds separate, complementary insight into the physical activity pattern.

**Conclusion:** Because most sensor axes show correlations close to zero, as visualized in the correlation matrix, they measure largely independent aspects of movement and environment. This minimal redundancy, confirmed by the statistical interpretation in the text, means each axis can contribute unique information, making them valuable complementary features for activity recognition.

Figure 9: **Example of EEQ (Evidence-Based Explanation & Quantification) type question.** This example demonstrates how the model must explain how visual patterns (correlation matrix) support textual claims with quantitative analysis, integrating statistical interpretation from the figure with conceptual explanations from the text.

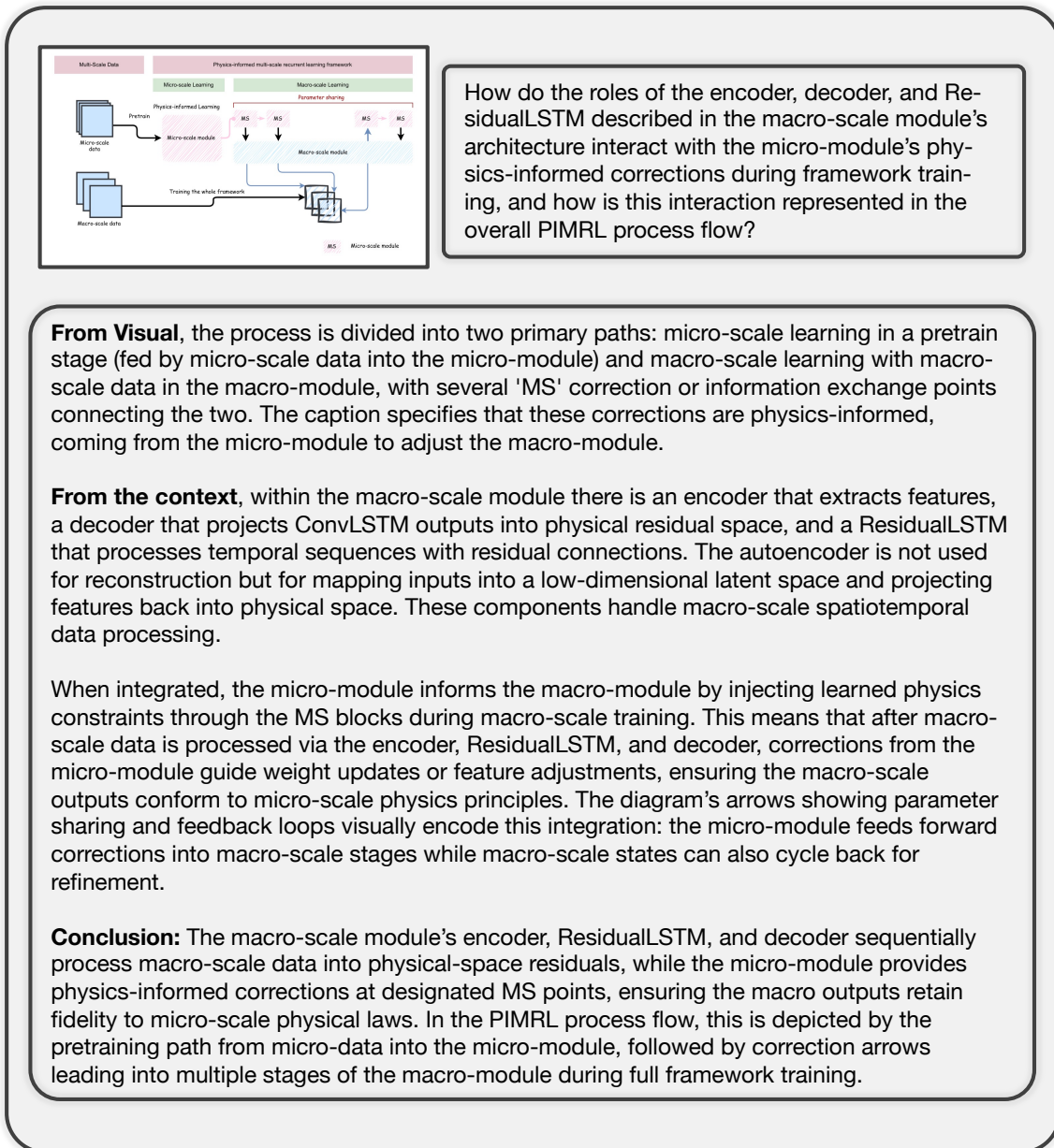
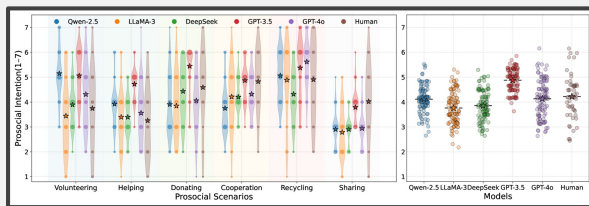


Figure 10: **Example of CIM (Concept-to-Instance Mapping) type question.** This example shows how the model links abstract architectural components (encoder, decoder, ResidualLSTM) described in text to their concrete visual representations in the system diagram, tracing information flow across modules.



What factors might explain why GPT-3.5 demonstrates the strongest prosociality in helping, donating, and cooperation tasks compared to other LLMs and humans, while models like Qwen-2.5 and LLaMA-3 display more polarized or consistently conservative behaviors across scenarios?

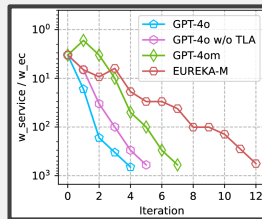
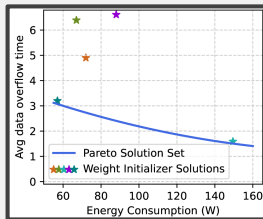
From the left panel violin plots, GPT-3.5's average scores (star markers) are consistently high in Helping, Donating, and Cooperation—above both human and other LLM averages. In these tasks, the distribution of GPT-3.5's responses often skews toward the upper Likert scale (scores 5–7). In contrast, Qwen-2.5 shows a high average only in certain tasks (e.g., Recycling, Volunteering), but dips in interpersonal ones like Sharing; LLaMA-3's stars are near the lower-middle range across all tasks, indicating hesitancy or conservative ratings.

The caption clarifies that the stars indicate average intention scores and that the six scenarios include interpersonal and public-good contexts like Volunteering, Helping, Donating, Cooperation, Recycling, and Sharing. The high GPT-3.5 averages in Helping, Donating, and Cooperation thus represent strong prosocial intentions across both interpersonal and collaborative public tasks.

The text in Section 3.1 explains that GPT-3.5 exhibits top overall prosociality, particularly in helping, donating, and cooperation. It also notes that model behavior is sensitive to scenario type: Qwen-2.5 is more polarized, excelling in low-cost/public-good contexts (e.g., Recycling) but scoring lower in interpersonal contexts (e.g., Sharing). LLaMA-3 is consistently conservative, likely due to more restrained decision-making or a cautious alignment style. Additionally, the text states LLMs are more prosocial in high-norm scenarios and less so when personal cost is high, suggesting that structural features of the scenario interact with each model's alignment and training data biases. GPT-3.5's high scores across diverse cooperative and altruistic contexts may indicate that its alignment tuning has generalized well to both low- and high-cost prosocial acts, unlike models that overfit toward certain structural contexts.

**Conclusion:** GPT-3.5's consistent top-tier performance in helping, donating, and cooperation tasks likely reflects alignment training that generalizes prosocial norms across diverse contexts, enabling high prosociality even in potentially costly interpersonal scenarios. In contrast, Qwen-2.5's polarized pattern suggests scenario-specific sensitivity, and LLaMA-3's uniformly low responses indicate a conservative approach that limits positive ratings regardless of context.

Figure 11: **Example of HVI (Hypothesis Validation & Inferential Reasoning) type question.** This example illustrates inferential reasoning where the model analyzes distributional patterns in violin plots alongside textual explanations to infer underlying factors explaining behavioral differences across models.



Does the evidence suggest that GPT-4om's slower convergence in reward weight search is fully explained by its weaker reasoning and numerical abilities, or might other factors also be contributing?

From the second plot, GPT-4om's curve (green) starts declining slower than GPT-4o's, achieving improvement later, and not reaching the low ratios as quickly. The caption confirms this is tracking the change of maximum  $w\_service/w\_ec$  during iteration.

The context attributes this to GPT-4om's weaker reasoning and numerical analysis capabilities, leading to less flexible step sizes and thus requiring more iterations. This plausibly explains the slower decline.

However, the visual shows GPT-4om's curve to be smoother and more regular than EUREKA-M's, suggesting that while reasoning limitations slow adaptation, they are not the sole cause; step size policy design and possibly input complexities also likely contribute. Its curve does decline steadily, indicating moderate optimization ability despite limitations.

**Conclusion:** GPT-4om's slower convergence is largely consistent with reduced flexibility from weaker reasoning and numerical skills, but the steadiness of its decline compared to EUREKA-M suggests that suboptimal step size strategy and other design factors may also play roles in its performance gap.

Figure 12: **Example of CAC (Critical Analysis & Consistency Check) type question.** This example demonstrates critical evaluation of whether textual claims are accurately supported by visual data, requiring careful assessment of evidence strength and potential discrepancies.

How does the breadth of research topics on underwater debris reflect the urgency of detecting and managing marine debris to protect ecosystems, especially given the diverse challenges in harsh marine environments?

The diagram visually categorizes research about underwater debris into multiple branches: Deep Learning Techniques, Environmental Conservation, Alternative Solutions, Datasets, and Other research. Each branch contains numerous references, indicating a broad and multidisciplinary engagement with the problem, as clarified in the caption describing this as the 'literature available for the subsequent topics.'

The context specifically states that detecting and managing marine debris is critical for preserving ecosystems, while also detailing particular difficulties—such as the debris' small size, submerged nature, concealment on the sea bottom, and threats to wildlife. These challenges necessitate a variety of specialized approaches.

By integrating the categories shown in the diagram with the context's problems, it becomes clear that each research branch addresses different facets of these challenges: conservation research targets ecological impacts, deep learning methods tackle detection difficulties, alternative solutions explore prevention or cleanup strategies, while datasets provide foundational resources. This broad coverage underscores the urgency and complexity of the management task, as multiple interconnected avenues are pursued to mitigate ecosystem risks.

**Conclusion:** The diversity of research topics—from conservation-focused studies to technical detection methods—demonstrates that managing marine debris is a multi-pronged effort driven by the urgency of addressing numerous detection difficulties and environmental risks in harsh marine environments.

Figure 13: **Example of ARS (Argumentative Role & Synthesis) type question.** This example shows how the model synthesizes visual evidence and textual arguments to articulate the overall scientific contribution and understand the role of visual elements in supporting the main thesis.