

# SciEvalKit: An Open-source Evaluation Toolkit for Scientific General Intelligence

SciPrismaX Team and Community Contributors\*, Shanghai Artificial Intelligence Laboratory

## Abstract

We introduce SciEvalKit, a unified benchmarking toolkit designed to evaluate AI models for science across a broad range of scientific disciplines and task capabilities. Unlike general-purpose evaluation platforms, SciEvalKit focuses on the core competencies of scientific intelligence, including Scientific Multimodal Understanding, Scientific Symbolic Reasoning, Scientific Code Generation, Science Hypothesis Generation and Scientific Knowledge Understanding. It supports six major scientific domains, spanning from physics and chemistry to astronomy and materials science. SciEvalKit builds a foundation of expert-grade scientific benchmarks, curated from real-world, domain-specific datasets, ensuring that tasks reflect authentic scientific challenges. The toolkit features a flexible, extensible evaluation pipeline that enables batch evaluation across models and datasets, supports custom model and dataset integration, and provides transparent, reproducible, and comparable results. By bridging capability-based evaluation and disciplinary diversity, SciEvalKit offers a standardized yet customizable infrastructure to benchmark the next generation of scientific foundation models and intelligent agents. The toolkit is open-sourced and actively maintained to foster community-driven development and progress in AI4Science.

Page <https://prismax.opencompass.org.cn/>  
 Code <https://github.com/InternScience/SciEvalKit>

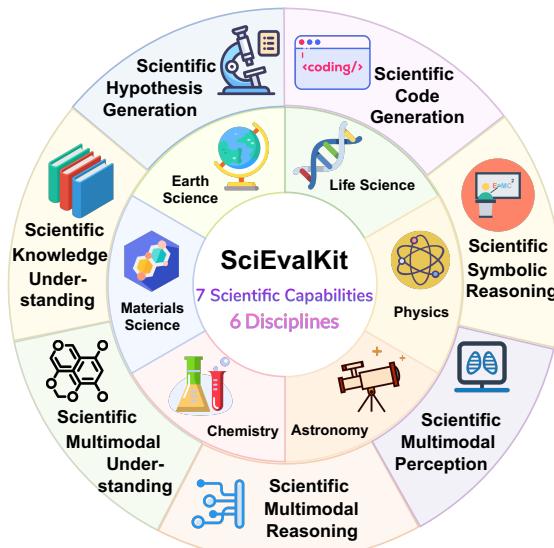


Figure 1 | Overview of the SciEvalKit scientific intelligence evaluation framework.

\*SciEvalKit contributors can join the author list of the report based on their contribution to the repository. Specifically, it requires 3 major contributions (implement a new benchmark, foundation model, or contribute a major feature). We will update the report quarterly and an additional section that details each developer's contribution will be appended in the next update.

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Benchmark Suite</b>	<b>4</b>
2.1	Core Competencies Taxonomy of Scientific Intelligence	4
2.2	Scientific Discipline Coverage	6
2.3	Expert-Aligned Benchmark Construction	6
2.3.1	Principles of Expert-Aligned Benchmark Design	6
2.3.2	Benchmark Overview	7
<b>3</b>	<b>Evaluation Framework</b>	<b>7</b>
3.1	Abstraction Layer	7
3.2	Unified interface for prompt construction and prediction	9
3.3	Capability-Oriented Evaluation	9
3.4	Evaluation Modes	10
<b>4</b>	<b>Evaluation Results</b>	<b>11</b>
4.1	Capability-Wise Leaderboard	13
4.2	Discipline-Wise Leaderboard	14
<b>References</b>		<b>15</b>
<b>A</b>	<b>Appendix</b>	<b>17</b>
A.1	Authors	17
A.2	Full Evaluation Results Across Core Benchmarks	17
<b>B</b>	<b>Benchmark Description</b>	<b>19</b>
<b>C</b>	<b>Representative Task Cases</b>	<b>21</b>
C.1	MaScQA	21
C.2	Chembench	22
C.3	LLM4Chem	23
C.4	SciCode	23
C.5	PHYSICS	24
C.6	CMPPhysBench	24
C.7	ClimaQA	26
C.8	EarthSE	26
C.9	ProteinLMBench	27
C.10	TRQA	27
C.11	ResearchBench	28
C.12	MSEarth	29
C.13	AstroVisBench	29
C.14	SLAKE	33
C.15	SFE	33

## 1. Introduction

The advances in large language models (LLMs) have demonstrated remarkable general-purpose reasoning and broad knowledge retrieval. Recently, researchers are increasingly interested in probing whether these models exhibit elements of scientific intelligence including conceptual understanding, symbolic reasoning, procedural modeling and hypothesis-driven exploration. Despite encouraging progress on individual benchmarks, current evaluations largely focus on surface-level correctness or narrow task-specific metrics, and therefore fail to assess whether LLMs can truly operate across the full spectrum of scientific reasoning. Real-world scientific problem solving fundamentally differs from generic reasoning: it requires conceptual abstraction, symbolic manipulation, hypothesis formation and testing, multi-step procedural thinking, and the ability to interpret structured visual representations such as chemical diagrams, protein structures. Yet, existing benchmarks rarely capture this holistic view, nor do they systematically evaluate these capabilities across scientific disciplines, modalities and cognitive dimensions.

From a cognitive perspective, scientific reasoning is inherently structural, relational, and multi-representational. The famous DSRP Theory which represents Distinctions, Systems, Relationships, Perspectives respectively posits that all complex reasoning emerges from these four fundamental cognitive patterns. Moreover, he emphasizes that individuals can improve their reasoning capacities by explicitly engaging with these four elements. This theoretical perspective provides a principled foundation for the core capabilities required to model scientific intelligence, prompting us to move beyond factual memorization or pattern recognition toward evaluating large language models' abilities in relational reasoning, structured reasoning, and representational alignment across textual, symbolic and visual modalities. Grounded in this framework, we propose a taxonomy of seven core dimensions of scientific intelligence as shown in Fig. 1 that reflect essential capabilities in modern scientific practice: (1) **Scientific Knowledge Understanding** which assesses the models' grasp of domain-specific concepts and factual relationships. (2) **Scientific Code Generation** which captures the ability to translate scientific descriptions and algorithmic procedures into executable logic. (3) **Scientific Symbolic Reasoning** which evaluates the manipulation of equations, physical laws, symbolic expressions, and structured notation. (4) **Scientific Hypothesis Generation** which measures the capabilities to propose plausible hypotheses and research directions. (5) **Scientific Multimodal Perception** which focuses on the entity localisation and grounding in paired visual–text inputs. (6) **Scientific Multimodal Reasoning** which involves chain-of-thought inference with domain scientific image data typically found in research papers. (7) **Scientific Multimodal Understanding** which probes rigorous interpretation of raw scientific data.

Existing benchmarks typically assess isolated abilities such as factual question answering, code completion and visual recognition. However, they seldom capture the reasoning processes for realistic scientific workflows. Most multimodal benchmarks are confined to generic image caption tasks, lacking expert-level semantic alignment capabilities for scientific diagrams. Furthermore, mainstream leaderboards often emphasize overall scores and overlook the detailed differences between scientific dimensions.

Modern frontier models are predominantly engineered and tuned for general-purpose utility—handling dialogue, broad-domain retrieval or generic reasoning tasks. Yet science tasks impose specific requirements: precise symbolic manipulation, code–execution fidelity and the ability to align dense textual arguments with highly specialised diagrams or experimental data. To quantify whether current general models meet these standards, we compare each model's score on general benchmarks against its score on corresponding scientific domain benchmarks (shown in Fig. 2). The strongest model Gemini-3 Pro already approaches 90 score on general tasks, but it and every other model falls below 60 score once they are probed under rigorous scientific scenario. This systematic gap highlights

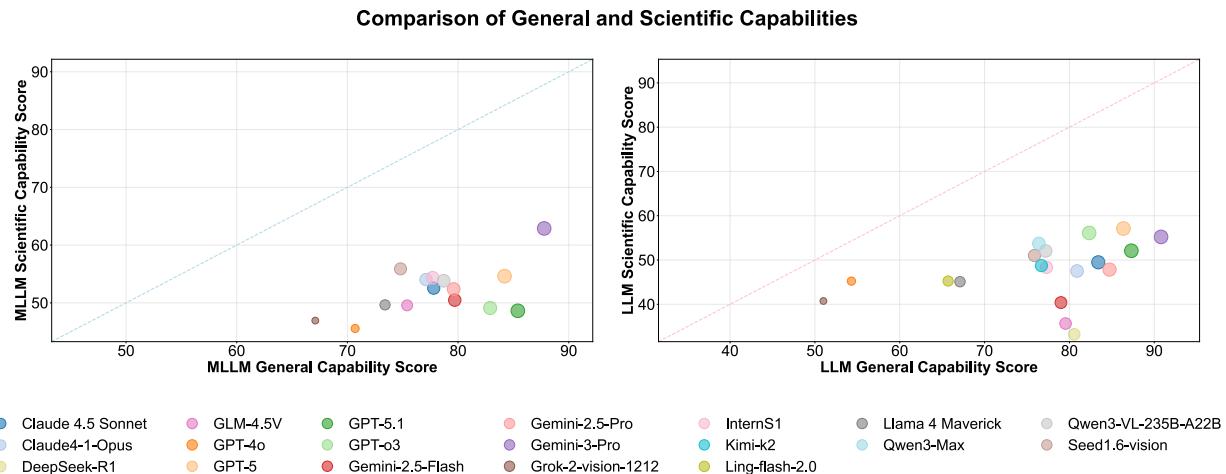


Figure 2 | Comparison of model performance on scientific versus general tasks.

the need for an integration of general and specialized capabilities that broad instruction tuning must be fused with expert-level skills in coding, symbolic reasoning and diagram understanding.

To fill this critical gap, we introduce SciEvalKit, an open-source scientific intelligence evaluation toolkit and leaderboard system designed to assess LLMs across the five core dimensions of scientific intelligence. SciEvalKit unifies 14 expert-curated benchmarks spanning six major scientific disciplines including life sciences, chemistry, earth sciences, materials science, physics, and astronomy. It integrates language-based and execution-based scoring, multimodal inputs, semantic LLM-as-a-judge validation, and expert-aligned evaluation criteria. Through this unified framework, SciEvalKit establishes a transparent, cognitively grounded, and scientifically credible evaluation paradigm for next-generation scientific AI systems.

Using SciEvalKit, we evaluate the scientific intelligence of cutting-edge LLMs, spanning proprietary, open-source, and multimodal architectures. Our findings reveal substantial disparities among models: while most achieve moderate to strong performance in knowledge understanding, capabilities such as symbolic reasoning and code generation remain underdeveloped. Notably, these shortcomings persist even among vision-enabled or instruction-tuned models, highlighting the need for capability-oriented evaluation.

In summary, this work contributes: (1) a seven-dimensional capability taxonomy grounded in expert-defined reasoning demands, (2) SciEvalKit: an open-source, multimodal, execution-aware, and expert-aligned evaluation toolkit, and (3) a comprehensive benchmark analysis of leading LLMs, uncovering critical gaps in their readiness for real scientific problem solving.

## 2. Benchmark Suite

### 2.1. Core Competencies Taxonomy of Scientific Intelligence

To comprehensively evaluate scientific intelligence in large language models, we construct a modality-aware taxonomy of seven core dimensions, classified into multimodal and text-only categories based on their input format and cognitive demands. Multimodal benchmarks emphasize the integration and alignment of visual and text contents to meet the reasoning demands of real-world scientific workflows, while text-only benchmarks probe symbolic, conceptual, and generative reasoning capabilities solely through language.

**Scientific Multimodal Perception** captures a model's ability to detect and localize scientifically meaningful entities from multimodal input. Unlike general visual perception, scientific perception requires identifying scientific structures such as organs in CT or MRI scans or chemically relevant patterns in the images. **Scientific Multimodal Understanding** represents the multimodal capability to extract and interpret structured scientific information from visual elements where the images themselves carry a high degree of scientific specificity. Benchmarks like SFE fall into this category where they require the model to align scientific symbols, notations and visual encodings with domain knowledge, making this capability essential for scientific workflows. **Scientific Multimodal Reasoning** refers to the model's ability to integrate visual and textual modalities to support coherent scientific inference. This capability goes beyond recognizing modality-specific patterns. It emphasizes cross-modal grounding, multi-step inference and domain-aware reasoning strategies. A key facet involves Chain-of-Thought (CoT) style reasoning, where models articulate intermediate steps when answering complex science questions. Within text-only benchmarks, we categorize benchmarks according to four key scientific capabilities. The first is **Scientific Knowledge Understanding**, which assesses a model's grasp of domain-specific concepts and factual relationships across disciplines such as Chemistry, Earth Science and Life Science. The second concerns **Scientific Code Generation** involving algorithmic comprehension and code generation tasks that demand precise mapping from scientific descriptions to executable logic. The third focuses on **Scientific Symbolic Reasoning** which targets a model's ability to manipulate equations, units, and structured scientific notations. This ability is particularly critical in disciplines such as physics, where symbolic representation plays a central role in modeling physical systems and deriving formal solutions. Finally, we include **Science Hypothesis Generation** where models engage in abductive inference and explanatory synthesis under open-ended or minimally structured prompts.

Table 1 | Taxonomy of Benchmarks for Scientific Intelligence capabilities.

Scientific Capability	Subject	Benchmark	Modality	Source
Scientific Multimodal Perception	Life Science	<a href="#">SLAKE</a>	Image	ISBI 21
Scientific Multimodal Reasoning	Earth Science	<a href="#">MSEarth</a>	Image	arXiv 25
Scientific Multimodal Understanding	Multidisciplinary	<a href="#">SFE</a>	Image	NeurIPS 25
	Earth Science	<a href="#">OmniEarth</a>	Image	arXiv 25
	Life Science	<a href="#">OmniMedVQA</a>	Image	CVPR 24
	Physics	<a href="#">PhyX</a>	Image	arXiv 25
Scientific Knowledge Understanding	Chemistry	<a href="#">ChemBench</a>	Text	Nat Chem 25
	Chemistry	<a href="#">ChemBench4K</a>	Text	arXiv 25
	Chemistry	<a href="#">LLM4Chem</a>	Image	COLM 24
	Earth Science	<a href="#">ClimaQA</a>	Text	ICLR 25
	Earth Science	<a href="#">EarthSE</a>	Text	arXiv 25
	Life Science	<a href="#">ProteinLMBench</a>	Text	BIBM 24
	Life Science	<a href="#">BioProbench</a>	Text	arXiv 25
	Materials Science	<a href="#">MaScQA</a>	Text	Digit Discov 24
	Life Science	<a href="#">TRQA</a>	Text	arXiv 25
	Life Science	<a href="#">Biology-Instructions</a>	Text	ACL 25
	Life Science	<a href="#">Mol-Instructions</a>	Text	ICLR 24
	Life Science	<a href="#">PEER</a>	Text	NeurIPS 22
Scientific Code Generation	Multidisciplinary	<a href="#">SciCode</a>	Text	NeurIPS 24
	Astronomy	<a href="#">AstroVisBench</a>	Image	NeurIPS 25
Scientific Symbolic Reasoning	Physics	<a href="#">CMPhysBench</a>	Text	arXiv 25
	Physics	<a href="#">PHYSICS</a>	Text	arXiv 25

Scientific Capability	Subject	Benchmark	Modality	Venue
Science Hypothesis Generation	Multidisciplinary	<a href="#">ResearchBench</a>	Text	ICML 25

## 2.2. Scientific Discipline Coverage

A key design principle of our benchmark suite is to ensure comprehensive coverage of the major scientific disciplines where large language models are expected to demonstrate expert-level reasoning. Rather than evaluating models through isolated subject-specific tasks, our framework spans the full landscape of natural sciences, including **Life Science, Chemistry, Earth and Environmental Science, Physics, and Materials Science**.

Each discipline is represented through benchmarks that reflect not only factual or textbook-level knowledge, but also procedural reasoning, mechanistic interpretation, and context-dependent application. For example, ProteinLMBench and TRQA-lit capture biomolecular and biomedical reasoning, requiring the integration of protein sequence understanding with biological function and therapeutic context. ChemBench and MaScQA evaluate higher-order chemical and materials reasoning, including thermodynamics, phase transitions, and structure–property analysis. ClimaQA and EarthSE focus on earth system, climate science, and geospatial interpretation, while PHYSICS and CMPHysBench emphasize fundamental physical laws, mathematical modeling, and symbolic derivations.

This broad disciplinary coverage enables the evaluation of not only knowledge-centric tasks, but also the cognitive heterogeneity across scientific domains where reasoning formats, cognition levels, modality dependencies, and knowledge structures differ substantially. As such, our benchmark suite offers a more faithful and holistic view of how LLMs generalize across scientific fields, rather than excelling in narrow domains.

## 2.3. Expert-Aligned Benchmark Construction

### 2.3.1. Principles of Expert-Aligned Benchmark Design

To ensure the benchmark suite faithfully evaluates the scientific capabilities of LLMs, we adopt a construction paradigm grounded in domain expertise, cognitive coverage and procedural transparency. The benchmarks aims to represent core scientific workflows through rigorously selected tasks that reflect pressing and high-impact scientific questions of contemporary importance.

We conduct a multi-round consultation with domain experts from diverse scientific fields including chemistry, earth science, life science, materials science, and physics. These experts are invited to propose benchmark tasks reflective of authentic research challenges they encounter ranging from climate assessment, protein function inference, thermodynamic reasoning in material design and scientific code generation. Proposals are evaluated based on their fidelity to real-world scientific reasoning and alignment with high-priority questions in specific scientific domains.

Following the open-ended proposal phase, a second-stage selection process is implemented to identify core benchmarks that satisfied below criteria:

1. **Scientific Validity:** Tasks must be grounded in real scientific content and reasoning, avoiding mechanical factual memory. For example, questions like chemical reaction pathways or protein structure interpretation are preferred over superficial concept definitions.
2. **Expert Calibration:** Each benchmark undergoes manual verification and calibration by domain experts, who validate the correctness of task formulations, solution rationales, and scoring criteria.
3. **Capability Coverage:** Selected benchmarks must collectively span the five core dimensions of scientific intelligence—Scientific Knowledge Understanding, Scientific Code Generation, Symbolic Reasoning, Hypothesis Generation, and Diagram Understanding—ensuring that the evaluation reflects both analytical depth and reasoning breadth.
4. **Modality and Task Diversity:** The suite covers multiple modalities (text, diagrams, molecular structures, protein sequences, scientific plots, radiological imagery, etc.) and multiple task formats (multiple-choice, free-form generation, code execution, document analysis), thereby capturing the multimodal and procedural nature of real scientific workflows.

5. **Community Recognition:** Benchmarks are preferably endorsed by the broader scientific or industrial community at domain-leading conferences, or released by reputable research groups. This ensures the benchmarks' credibility, relevance and alignment with community-validated standards.

### 2.3.2. Benchmark Overview

To assess the scientific intelligence across diverse modalities and tasks, we curated a suite of expert-aligned benchmarks that reflect real-world scientific workflows. Each benchmark is designed to evaluate specific dimension of scientific intelligence. The benchmarks discussed below are only those have been evaluated in the current release and appear on our leaderboard. A broader benchmarks can be found on the SciEvalKit homepage. Additionally, these benchmarks span a broad spectrum of scientific disciplines including chemistry, earth sciences, life sciences, materials science, and astrophysics

We provide a systematic evaluation of scientific intelligence from two dimensions: text-only and multimodal evaluation. The text-only benchmarks focus on evaluating a model's ability to comprehend, reason over, and generate scientific content using purely textual input.

In contrast, the multimodal benchmarks introduce scientific problems where text and visual inputs are both required. Such settings simulate real-world scientific scenarios where visual and textual reasoning must be integrated, enabling a more complete evaluation of the model's ability to engage with complex scientific artifacts. We provide detailed descriptions of each benchmark including their disciplinary coverage, task design, and alignment with scientific capabilities below.

Our first release evaluates models on a curated subset of SciEvalKit benchmark pool. Detailed dataset description are referred to Appendix B. At a glance, the present suite comprises:

- *Text-only evaluation:* nine knowledge-centric benchmarks (e.g. ChemBench, MaScQA, ProteinLMBench), two specialised reasoning sets for code generation (SciCode, AstroVisBench) and symbolic manipulation (CMPhysBench, PHYSICS).
- *Multimodal evaluation:* three vision-language benchmarks which are MSEarth, SLAKE and SFE, requiring joint reasoning over scientific figures and textual context.

Together these tasks span chemistry, earth and climate science, life science, materials science, physics and astronomy. The suite was selected through our expert-aligned construction workflow to ensure (i) high scientific validity, (ii) coverage of all seven capability dimensions, and (iii) diversity of modalities and question types.

Table 1 summarizes the core benchmarks included in this release of SciEvalKit. These benchmarks are selected based on their coverage of key scientific intelligence capabilities across diverse disciplines and modalities. Specifically, this table highlights the most representative and capability-aligned benchmarks that form the evaluation backbone for our first leaderboard release. And we will continuously expand the coverage in future releases. This includes incorporating additional modalities, disciplines, and newly proposed agent tasks from the community.

## 3. Evaluation Framework

### 3.1. Abstraction Layer

The framework of SciEvalKit is organised as four cooperating layers which are Dataset, Model Inference, Evaluation & Testing, and Report & Storage that together deliver an end-to-end, reproducible pipeline for multimodal scientific benchmarking. Each layer's scope is deliberately kept concise with clearly defined boundaries so that researchers can extend one part of the stack without impacting the others.

**Dataset Layer.** The Dataset Layer serves as the entry point for data ingestion and task specification. Dataset construction is handled through the `build_dataset` routine, which maps a dataset identifier to its respective dataset class using centralized registries (e.g., `supported_video_datasets`, `supported_text_datasets`). Each dataset class inherits from either `TextBaseDataset`, `ImageBaseDataset`, or `VideoBaseDataset`, which provide unified interfaces for TSV or metadata loading, index normalization, and modality-specific data caching. Each dataset implements a custom `build_prompt()` method that encapsulates raw task data such as questions, images, video frame paths, code

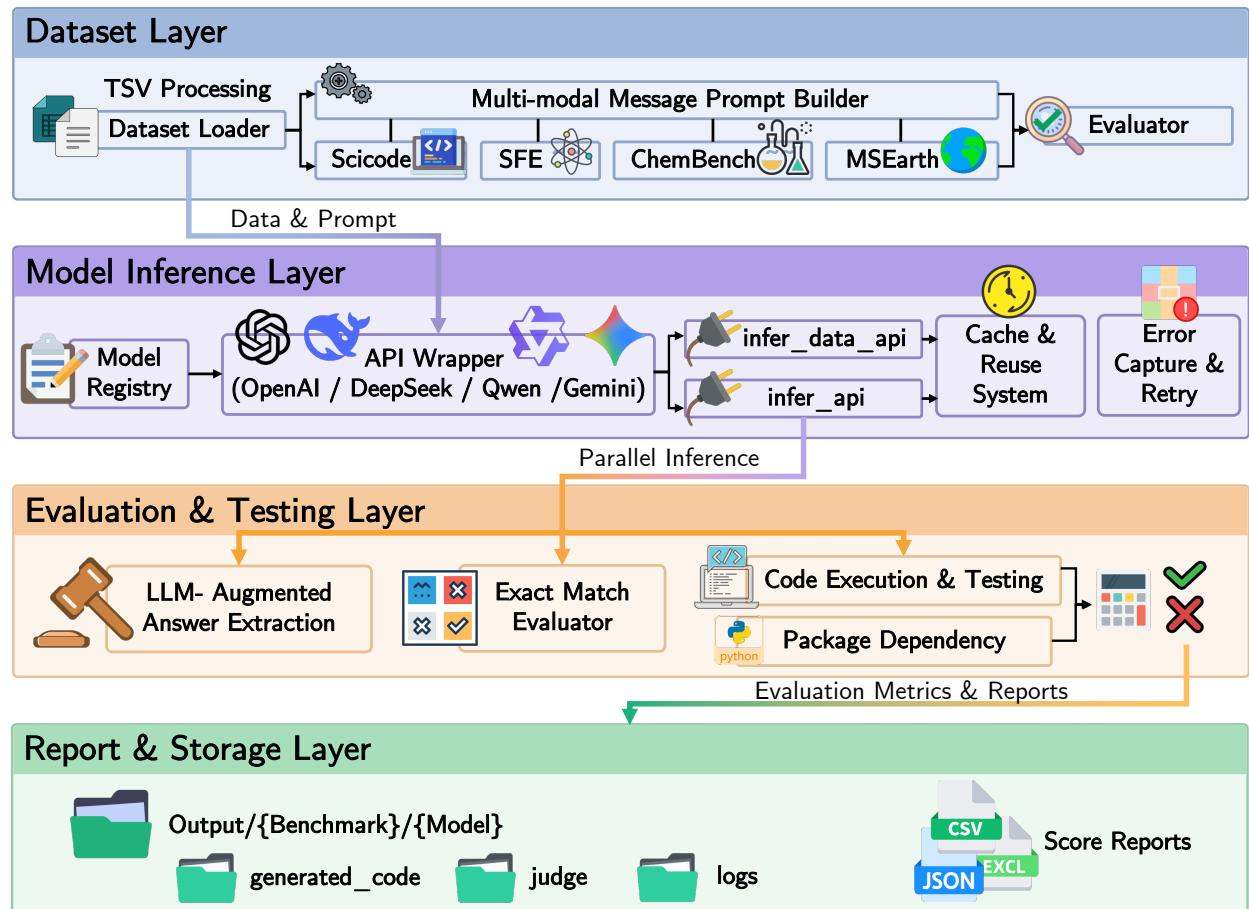


Figure 3 | Evaluation pipeline used in SciEvalKit.

snippets, and answer options into interleaved multi-modal messages. This representation constitutes the atomic unit for model-level inference.

**Model Inference Layer.** The Model Inference Layer mediates between structured prompts and model outputs. Model instantiation is handled by `build_model_from_config`, which resolves model metadata from `supported_VLM`. Each model object exposes a unified `.generate(message, dataset)` interface, abstracting the distinction between local inference (via vLLM or `torch.distributed`) and API-based cloud models (OpenAI, DeepSeek, Gemini, Anthropic). Inference workflows are orchestrated via `infer_data()`, `infer_data_api()` and `infer_data_job_video()`, which provide transparent support for batching, parallel token generation, retry or error tolerance, and partial progress restoration through `_supp.pkl`.

**Evaluation & Testing Layer.** Once predictions are generated, dataset-specific evaluate methods perform capability-aligned scoring through a combination of exact matching, semantic retrieval, numerical scoring, code execution, and LLM-based judging function `evaluate()`. Matching utilities, multiple-choice answer extraction and general-purpose judging `build_judge()` provide deterministic and LLM-augmented evaluation paths. Code-execution tasks (e.g., Scicode) invoke sandboxed Python environments to verify visual output fidelity and computational correctness.

**Report & Storage Layer.** The Report & Storage Layer ensures reproducibility and transparent logging. All predictions, logs, reasoning traces, metadata, and evaluation results follow a structured file convention. Helper functions (e.g., `get_pred_file_path()`, `prepare_reuse_files()`, `get_intermediate_file_path()`) ensure consistency across model runs. Final evaluation metrics are serialized in CSV, JSON, or XLSX formats based on benchmark requirements, facilitating both longitudinal comparison and leaderboard hosting.

### 3.2. Unified interface for prompt construction and prediction

SciEvalKit is developed on the top of the VLMEvalKit [1] code base, preserving its modular abstractions while introducing extensions for scientific multimodal inputs, discipline-aware prompt construction and capability-oriented evaluation. A central design goal of SciEvalKit is to provide a unified interface through which datasets, models, and evaluators interact in a modality-agnostic manner. To achieve this, the framework establishes a standardized prompt construction and prediction interface that applies consistently to text-only tasks, image-based visual reasoning, and multi-modal scientific problems involving diagrams, code snippets, symbolic expressions, molecular structures, and geospatial imagery. This interface is implemented at both the dataset and model abstraction layers, enabling fully unified end-to-end execution without requiring per-dataset or per-model procedural handling.

On the dataset side, every dataset class inherits from a base abstraction such as `TextBaseDataset`, `ImageBaseDataset`, or `VideoBaseDataset`, each of which exposes a common set of required functions. Among them, the `.build_prompt()` method is the core entrypoint responsible for converting a structured sample from the dataset (e.g., TSV row containing question text, answer options, image encoding, or code references) into a standardized multi-modal message representation. A multi-modal message is represented as an ordered list of typed content segments, where each segment explicitly declares both its modality and payload.

```
dict(type='text', value=text)
dict(type='image', value=tgt_path)
```

This explicit specification ensures that models can interpret consistently. If a dataset supports custom instruction formats, such as SFE's discipline-aware prompt templates for multiple-choice, exact-match, or open-ended questions, these are applied within `.build_prompt()` while still adhering to the unified output message schema. The interface also accommodates advanced message packing for video or sequential images when activated in `infer_data_api()`.

Additionally, datasets provide a `.evaluate()` method as a standardized evaluator interface, which get model predictions from the storage layer and applies deterministic or judge-assisted scoring. Depending on the dataset type, the evaluation pipeline may involve exact matching, choice extraction, code execution for scientific programming tasks, or LLM-based scoring via `build_judge()`. Despite these internal variations, the signature of the `evaluate()` function remains constant across datasets, ensuring compatible execution. Supporting utility functions such as `display()` and `dump_image()` further contribute to this unified interface by enabling consistent visual inspection, debugging, and handling of base64-encoded or remote images, ensuring that dataset parsing and integrity checks are conducted through shared mechanisms.

On the model side, the framework enforces a single inference interface through `.generate()`, regardless of whether the model is accessed via API endpoints (e.g., GPT-4o, Gemini, DeepSeek) or runs locally through vLLM or PyTorch-based implementations. All inference functions including `infer_data()`, `infer_data_api()` and `infer_data_job_video()` construct messages from dataset prompts and invoke `model.generate(message=..., dataset=..., **kwargs)` in a unified manner. This uniform invocation mechanism abstracts away differences in backend execution, request formatting, batching, retry handling, or temperature sampling. Moreover, models may optionally override prompt formatting by checking `use_custom_prompt(dataset_name)`, while still maintaining conformity to the unified interface contract.

This unified interface design enables powerful decoupling between dataset logic, model execution, and evaluation strategy. Researchers can incorporate new datasets by providing only `.build_prompt()` and `.evaluate()`, without modifying inference or scoring pipelines. Models can be swapped freely by exposing a `.generate()` method that accepts standardized multi-modal messages. Inference scripts can construct prompts, dispatch model generation, transfer results, and trigger evaluation without conditional branching or dataset-specific logic. As a result, the entire pipeline is fully modular, reproducible and extensible, highlighting the architectural robustness of SciEvalKit.

### 3.3. Capability-Oriented Evaluation

To move beyond aggregate accuracy and better reflect the multifaceted nature of scientific intelligence, we adopt a capability-oriented evaluation paradigm. Rather than treating all benchmarks uniformly, we classify them into five core competency dimensions based on the underlying cognitive demands, task structure, and the

knowledge representations required for successful problem solving. This allows each model to be evaluated not only on performance, but on what kind of scientific reasoning it is capable of. For each capability dimension, a model's score is computed as the average performance across all benchmarks belonging to that capability, ensuring that evaluation is both comprehensive and domain-balanced.

**Scientific Image Understanding.** This dimension assesses a model's ability to interpret visually encoded scientific information and align it with textual. It includes *MSEarth*, *SLAKE*, and *SFE*, which involve high-fidelity scientific diagrams such as satellite environmental maps, CT/MRI scans, medical annotation images and molecular diagrams. These tasks require entity localization, structural interpretation, and diagram-text alignment—capabilities that differ fundamentally from generic visual understanding. The final capability score is the mean across these benchmarks.

**Scientific Code Generation.** This dimension evaluates a model's capacity to translate scientific intent into executable computational procedures. *AstroVisBench* and *SciCode* are assigned to this category because their tasks require not only generating syntactically correct code, but also reasoning over domain-specific computational logic (e.g., astrophysics data processing, numerical simulation and scientific plotting). The score for this capability is computed as the average over these benchmarks, reflecting both semantic correctness and the ability to produce executable code that aligns with real-world scientific workflows and engineering practices.

**Scientific Symbolic Reasoning.** This dimension focuses on symbolic knowledge representation, formula manipulation, unit reasoning, and quantitative derivation. *CMPHysBench* and *PHYSICS* both evaluate these capacities using tasks that require algebraic transformation, dimensional consistency, and symbolic inference rather than mere textual recall. The capability score is derived from the average results across these symbolic reasoning benchmarks. These tasks are particularly representative of the unique reasoning challenges found in physics where symbolic representations are fundamental to problem solving.

**Science Hypothesis Generation.** *ResearchBench* is the only benchmark dedicated to this capability and uniquely represents open-ended hypothesis formulation, scientific discourse planning, and research proposal synthesis. Unlike conventional QA tasks, it evaluates abductive reasoning, conceptual integration, novelty generation, and literature-grounded justification.

**Scientific Knowledge Understanding.** This dimension evaluates a model's capacity to comprehend and reason over domain-specific scientific knowledge across a broad spectrum of disciplines. Benchmarks such as ProteinLMBench, MaScQA, ClimaQA, TRQA-lit, Earth-Silver, and ChemBench focus on concept interpretation, factual consistency, mechanistic understanding, property relations, and application-oriented reasoning grounded in real scientific contexts. These tasks typically involve integrating procedural knowledge, conceptual hierarchy, and disciplinary logic rather than symbolic derivation or general linguistic inference. The score for this capability dimension is computed as the average performance across all benchmarks categorized under Scientific Knowledge Understanding.

### 3.4. Evaluation Modes

To ensure that model performance is assessed in a manner faithful to scientific rigor, our evaluation framework adopts a hybrid scoring paradigm that integrates **deterministic rule-based matching**, **semantic LLM-based judging**, and **execution-based verification**. This design accommodates the diverse answer formats and reasoning modalities present across scientific tasks ranging from symbolic problem solving and diagram interpretation to code synthesis and explanatory reasoning.

**Evaluation Methods.** We employ two primary evaluation pipelines:

- Natural-language matching: This is used for most knowledge, reasoning, and visualization oriented tasks. Model predictions are extracted and normalized via rule-based processing (e.g., option inference, unit normalization, numerical span extraction, or text canonicalization), and then compared against ground truth using dataset-specific metrics. For different task types, we adopt appropriate scoring schemes such as accuracy (MCQ) or relaxed numeric matching (scientific problem solving).
- Code-execution-based evaluation: For benchmarks such as *SciCode* that explicitly evaluate scientific code generation and algorithmic reasoning, the model's output is interpreted as executable Python programs. The predicted code is stitched into functional scripts, dependencies are resolved, and official unit tests are executed to verify correctness. Scores are computed based on the number of passed test

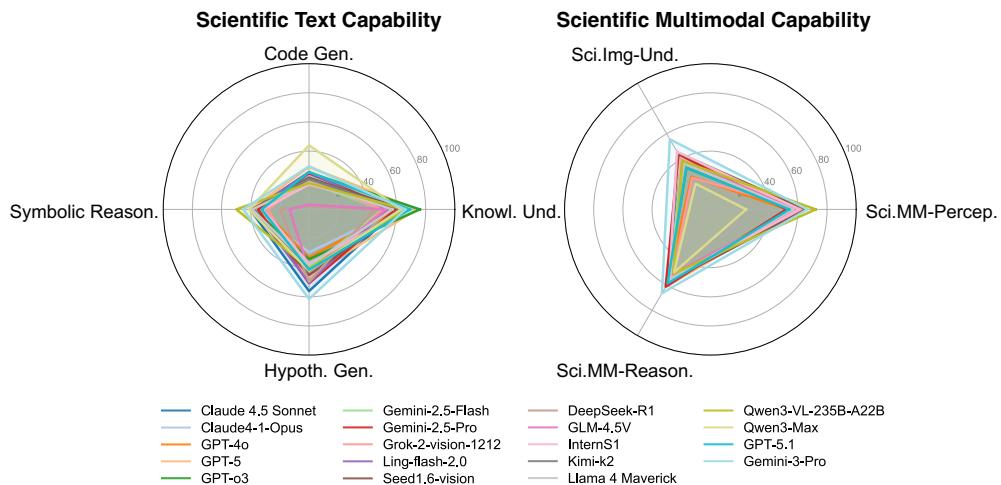


Figure 4 | Large-language-model (LLM) scientific capabilities (left) versus multimodal-language-model (MLLM) scientific capabilities (right) comparing the evaluated models on the SciEval leaderboard. Each axis reports the score (0 – 100) for one capability or scientific field; concentric rings mark 25 intervals up to the outer 100 score.

cases, reflecting not only syntactic validity but also functional correctness and engineering reliability.

**Question Format Handling.** To accommodate heterogeneous formats across benchmarks, we adopt differentiated evaluation strategies:

- Multiple-choice questions (MCQ): Predictions are mapped to option labels through rule-based extraction. If the model outputs free-form text, a secondary semantic alignment step infers the intended option (A/B/C/D). Accuracy is reported, optionally with per-category results.
- Fill-in-the-blank / cloze questions: Extracted responses are reduced to core information (numbers, chemical formulas, single phrases) using handcrafted parsing rules. Answers are scored based on exact match or relaxed matching with semantic normalization.
- Open-ended / free-form questions: For explanation questions, descriptive captioning, or scientific hypothesis formulation, we compute string- and semantics-based metrics (e.g., BLEU, VQA score, or semantic overlap). In certain benchmarks, optional qualitative grading (e.g., correctness, reasoning soundness) is also permitted using LLM-as-judge.

**Judgment Strategies.** To accurately evaluate answers requiring semantic equivalence, contextual reasoning, or scientific justification, we adopt two scoring strategies:

- Rule-based strategy: Deterministic string normalization, regular expression extraction, unit conversion, and option inference handle most cases involving structured answers or symbolic formats.
- LLM Evaluation Strategy: When rule-based approaches fail—particularly for free-form responses, complex reasoning explanations, or ambiguous mappings, we optionally invoke verification models (primarily closed-source SOTA models and SciEvalKit predominantly using GPT series) to assess semantic equivalence or correctness.

## 4. Evaluation Results

Table 2 | Evaluation of large language models across scientific capability benchmarks: Scientific Knowledge Understanding (Knowl. Und.), Scientific Code Generation (Code Gen.), Scientific Symbolic Reasoning (Symbolic Reason.), and Scientific Hypothesis Generation (Hypoth. Gen.).

Model	Knowl. Und.	Code Gen.	Symbolic Reason.	Hypothesis Gen.	Overall
<b>Closed-Weight MLLMs</b>					
Claude 4.5 Sonnet	60.67	21.73	40.36	56.10	44.72
Claude4-1-Opus	60.87	25.32	38.69	29.47	38.58
GPT-5.1	69.23	25.63	32.44	41.45	42.19
GPT-5	74.05	29.21	39.91	45.67	47.21
GPT-4o	60.84	17.67	32.09	33.04	35.91
GPT-o3	76.05	25.26	38.14	34.14	43.40
Gemini-3-Pro	66.06	29.57	45.19	61.51	50.58
Gemini-2.5-Flash	50.46	18.28	32.07	40.86	35.42
Gemini-2.5-Pro	59.34	24.77	34.96	50.73	42.45
Grok-2-vision-1212	50.14	20.60	28.21	49.63	37.14
Seed1.6-vision	65.78	21.49	39.24	45.00	42.88
<b>Open-Weight MLLMs</b>					
GLM-4.5V	52.78	3.24	13.43	42.23	27.92
InternS1	66.14	17.08	31.62	37.45	38.07
Ling-flash-2.0	53.39	25.60	37.98	50.29	41.81
DeepSeek-R1	45.17	0.06	20.00	49.73	28.74
Kimi-k2	62.49	20.86	38.59	42.28	41.06
Llama 4 Maverick	57.22	18.26	38.97	38.31	38.19
Qwen3-VL-235B-A22B	65.98	18.00	49.93	40.62	43.63
Qwen3-Max	63.14	43.97	41.04	42.12	47.57

Fig. 4 reveals several systematic patterns that align with the quantitative leaderboards while providing additional insights. First, the Scientific Knowledge Understanding spoke is consistently the strongest for nearly every model while the Code Generation and Symbolic Reasoning axes lag considerably in comparison. This imbalance underscores that contemporary LLMs have largely solved factual recall but still struggle with executable logic and equation-level manipulation, two skills that remain essential for modern computational science. Second, three proprietary models which are Gemini 3 Pro, GPT-5 and GPT-o3 form the widest and most balanced polygons, reflecting the highest average scores and the most even performance. Notably, the open-source Qwen3-VL-235B-A22B achieves scores that approach those of the leading proprietary models, indicating that some open-source systems now perform competitively across multiple scientific competencies.

A closer comparison between GPT-5 and its follow-up GPT-5.1 indicates that scores decrease across nearly every axis, indicating that the latest iteration prioritises incremental alignment refinements over new scientific competence. From a disciplinary perspective, the radar chart on the right reveals the weakest links in Physics and Life sciences. This contraction contrasts sharply with the relatively broad coverage in earth sciences and chemistry, highlighting areas where benchmark designers and model developers should concentrate their efforts. Finally, the highest scores are still achieved by general-purpose models such as Gemini 3 Pro and GPT-o3 rather than by any domain-specific, science-tuned models.

Table 3 | Evaluation of LLMs on scientific multimodal capabilities, including Scientific Multimodal Perception (Sci.MM-Percep.), Scientific Multimodal Understanding (Sci.MM-Und.), and Scientific Multimodal Reasoning (Sci.MM-Reason.).

Model	Sci.MM-Percep.	Sci.MM-Und.	Sci.MM-Reason.	Overall
<b>Closed-Weight MLLMs</b>				
Claude 4.5 Sonnet	57.87	43.64	56.11	52.54
Claude4-1-Opus	58.25	③ 45.19	58.66	54.03
GPT-5.1	54.10	33.05	58.73	48.63
GPT-5	59.94	42.44	② 61.46	③ 54.61
GPT-4o	52.78	25.93	57.97	45.56
GPT-o3	55.23	32.84	59.27	49.11
Gemini-3-Pro	② 66.54	① 55.62	① 66.49	① 62.88
Gemini-2.5-Flash	55.98	38.20	57.22	50.47
Gemini-2.5-Pro	52.12	43.76	③ 61.28	52.39
Grok-2-vision-1212	64.00	25.04	51.76	46.93
Seed1.6-vision	③ 65.79	44.75	57.11	② 55.88
<b>Open-Weight MLLMs</b>				
GLM-4.5V	59.10	38.57	51.04	49.57
InternS1	60.89	② 45.73	56.47	54.36
Llama 4 Maverick	56.74	36.83	55.39	49.65
Qwen3-VL-235B-A22B	① 72.29	38.35	50.83	53.82
Qwen3-Max	24.51	20.40	49.86	31.59

#### 4.1. Capability-Wise Leaderboard

We first assess the models based on five scientific core capabilities: Knowledge Understanding (Knowl. Und.), Code Generation (Code Gen.), Symbolic Reasoning (Symbolic Reasoning), Hypothesis Generation (Hypoth. Gen.), and Image Understanding (Image Und.). Table ?? summarizes model performance across five scientific capability dimensions. Across the board, Gemini 3 Pro demonstrates the strongest overall capability, ranking first in four of five dimensions. It achieves the best performance in Hypothesis Generation (61.51), Symbolic Reasoning (45.19), Code Generation (16.92), and exhibits strong multimodal capability with the highest score in Scientific Multimodal Understanding (62.88). GPT-5, GPT-o3, and Claude 4.5 Sonnet form a competitive group with balanced capability profiles. GPT-5 attains the second-highest scores in Knowledge Understanding (74.05) and third-highest in Image Understanding (54.61), while GPT-o3 shows strongest performance in Knowledge Understanding (76.05). Claude models, although slightly behind GPT-series in higher-order reasoning, demonstrate stable performance in Hypothesis Generation (56.10) and Symbolic Reasoning (40.36).

Scientific Code Generation remains the weakest competency for all evaluated models, underscoring the gap between text-based reasoning and executable scientific problem-solving. Even top-performing models achieve relatively low scores (Gemini-3-Pro: 16.92, GPT-5: 13.85, Qwen3-Max: 12.31), suggesting that current LLMs struggle with implementing algorithmic structures and translating scientific logic into runnable code.

In contrast, Scientific Multimodal Understanding has already shown notably stronger performance for multimodal models. For example, Gemini-3-Pro (62.88), Seed1.6-vision (55.88) and GPT-5 (54.61), all demonstrate solid competency in extracting spatial, structural, and quantitative information from scientific imagery (e.g., SLAKE, SFE). InternS1 demonstrates strong and stable competencies in Scientific Multimodal Understanding, achieving a competitive score of 54.36, surpassing Claude models (52.54 for Claude 4.5 Sonnet) and closely trailing GPT-5 (54.61).

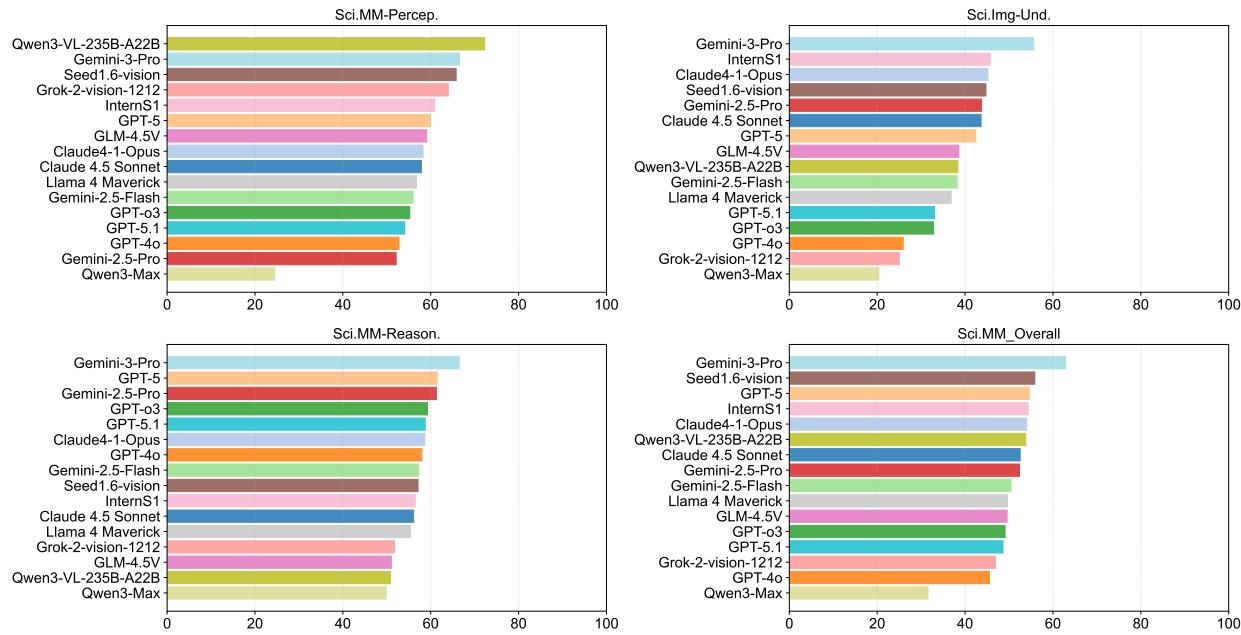


Figure 5 | Model scores on three multimodal competencies: Scientific MM Perception (Sci. MM Perception), Scientific Multimodal Understanding (Sci.MM-Und.), and Scientific MM Reasoning (Sci.MM-Reason.) together with their average (Sci.MM Overall).

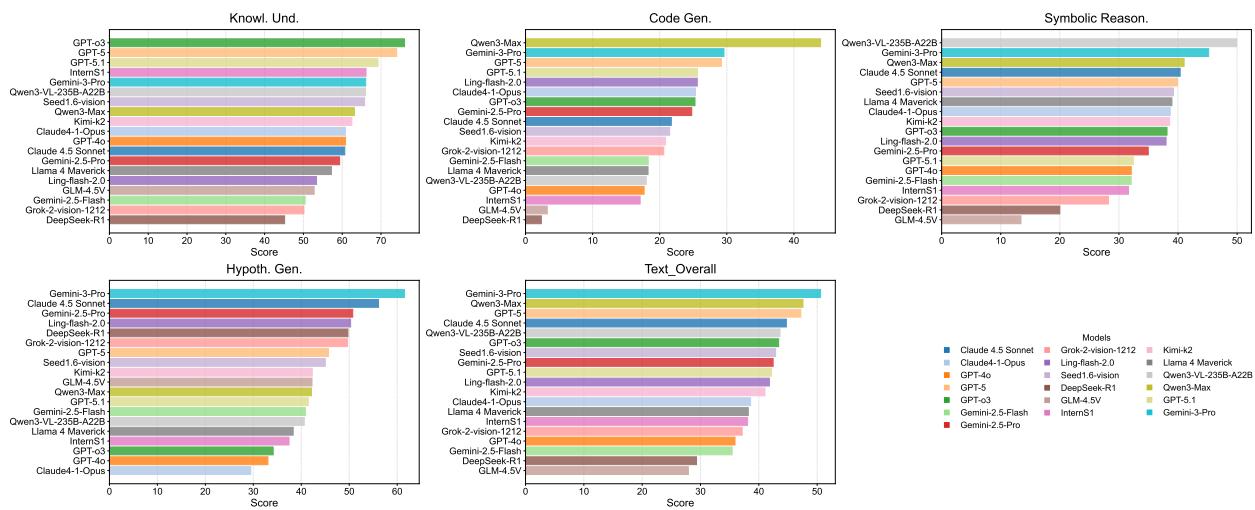


Figure 6 | Model scores on five text-only capacities: Scientific Knowledge Understanding (Knowl. Und.), Scientific Code Generation (Code Gen.), Scientific Symbolic Reasoning (Symbolic Reason.), Scientific Hypothesis Generation (Hypoth. Gen.), and their mean (Text Overall).

#### 4.2. Discipline-Wise Leaderboard

We aggregate results across six scientific disciplines: Life Science, Astronomy, Earth Science, Chemistry, Materials Science, and Physics as shown in Table 4. Chemistry and Materials Science appears to be the most "solved" domains for current frontier models. Performance in Chemistry, Earth Science and Materials Science tends to be substantially higher than in Physics and Life Science. This reflects the varied cognitive demands of the disciplines: while Chemistry and Materials Science rely heavily on structured factual knowledge and diagrammatic interpretation, Physics requires symbolic manipulation, mathematical derivation, and application of physical laws whose capabilities that current LLMs do not yet master reliably.

Across disciplines, Gemini-3-Pro leads with the highest average disciplinary performance, particularly excelling in Chemistry (83.08). GPT-5 closely follows in these domains, achieving 81.62 in Chemistry and 93.54 in Materials Science indicating strong material property reasoning. In contrast, Physics emerges as a difficult frontier across nearly all models. Even the leading systems, such as Qwen3-VL-235B-A22B (49.93), Gemini-3-Pro (45.19) and GPT-5 (39.91) achieve significantly lower scores relative to their performance in factual disciplines, highlighting persistent limitations in symbolic reasoning, derivation, and equation-grounded problem solving.

Table 4 | Evaluation of LLMs across six scientific disciplines.

Model	Life Sci.	Astronomy	Earth Sci.	Chemistry	Mat. Sci.	Physics
<b>Closed-Weight MLLMs</b>						
Claude 4.5 Sonnet	43.86	34.13	64.66	71.27	83.23	40.36
Claude4-1-Opus	42.49	39.87	67.47	71.94	83.23	38.69
GPT-5.1	② 60.54	42.02	73.88	76.45	61.08	32.44
GPT-5	③ 59.49	44.57	② 74.43	② 81.62	② 93.54	39.91
GPT-4o	54.46	30.73	63.08	72.37	61.54	32.09
GPT-o3	① 61.57	42.82	③ 74.15	81.77	① 93.85	38.14
Gemini-3-Pro	49.20	42.23	72.49	⑤ 83.08	③ 91.23	② 45.19
Gemini-2.5-Flash	35.61	31.94	68.21	75.12	62.15	32.07
Gemini-2.5-Pro	34.85	③ 43.40	70.52	③ 78.29	83.23	34.96
Grok-2-vision-1212	43.44	33.51	58.80	63.04	49.08	28.21
Seed1.6-vision	58.53	32.21	69.46	75.48	68.92	39.24
<b>Open-Weight MLLMs</b>						
Ling-flash-2.0	28.83	② 48.12	69.55	66.80	63.69	37.98
Kimi-k2	41.10	35.66	① 77.77	77.04	71.38	38.59
DeepSeek-R1	24.13	0.11	63.14	77.28	44.46	20.00
GLM-4.5V	54.11	0.32	52.68	71.38	57.08	13.43
InternS1	53.70	29.55	65.87	75.28	81.69	28.60
Llama 4 Maverick	42.92	31.91	59.82	68.06	80.77	38.97
Qwen3-VL-235B-A22B	56.69	29.84	67.81	77.39	81.85	① 49.93
Qwen3-Max	38.95	① 75.64	67.38	76.38	69.54	③ 41.04

## References

- [1] Haodong Duan et al. “Vlmevalkit: An open-source toolkit for evaluating large multi-modality models”. In: *Proceedings of the 32nd ACM international conference on multimedia*. 2024, pp. 11198–11201.
- [2] Xiangyu Zhao et al. “MSEarth: A Benchmark for Multimodal Scientific Comprehension of Earth Science”. In: *arXiv preprint arXiv:2505.20740* (2025).
- [3] Bo Liu et al. “Slake: A semantically-labeled knowledge-enhanced dataset for medical visual question answering”. In: *2021 IEEE 18th international symposium on biomedical imaging (ISBI)*. IEEE. 2021, pp. 1650–1654.
- [4] Yuhao Zhou et al. “Scientists’ First Exam: Probing Cognitive Abilities of MLLM via Perception, Understanding, and Reasoning”. In: *arXiv preprint arXiv:2506.10521* (2025).

- [5] Sebastian Antony Joseph et al. “Astrovisbench: A code benchmark for scientific computing and visualization in astronomy”. In: *arXiv preprint arXiv:2505.20538* (2025).
- [6] Minyang Tian et al. “Scicode: A research coding benchmark curated by scientists”. In: *Advances in Neural Information Processing Systems* 37 (2024), pp. 30624–30650.
- [7] Adrian Mirza et al. “A framework for evaluating the chemical knowledge and reasoning abilities of large language models against the expertise of chemists”. In: *Nature Chemistry* (2025). ISSN: 1755-4349. DOI: [10.1038/s41557-025-01815-x](https://doi.org/10.1038/s41557-025-01815-x). URL: <http://dx.doi.org/10.1038/s41557-025-01815-x>.
- [8] Veeramakali Vignesh Manivannan et al. “ClimaQA: An Automated Evaluation Framework for Climate Question Answering Models”. In: *arXiv preprint arXiv:2410.16701* (2024).
- [9] Yiqing Shen et al. “A fine-tuning dataset and benchmark for large language models for protein understanding”. In: *2024 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. IEEE. 2024, pp. 2390–2395.
- [10] Mohd Zaki, NM Anoop Krishnan, et al. “MaScQA: investigating materials science knowledge of large language models”. In: *Digital Discovery* 3.2 (2024), pp. 313–327. DOI: [10.1039/D3DD00188A](https://doi.org/10.1039/D3DD00188A). URL: <https://pubs.rsc.org/en/content/articlehtml/2024/dd/d3dd00188a>.
- [11] Zhongyue Zhang et al. “OriGene: A Self-Evolving Virtual Disease Biologist Automating Therapeutic Target Discovery”. In: *bioRxiv* (2025), pp. 2025–06.
- [12] Weida Wang et al. “CMPhysBench: A Benchmark for Evaluating Large Language Models in Condensed Matter Physics”. In: *arXiv preprint arXiv:2508.18124* (2025).
- [13] Shenghe Zheng et al. “Scaling physical reasoning with the physics dataset”. In: *arXiv preprint arXiv:2506.00022* (2025).
- [14] Haofei Yu et al. “Researchtown: Simulator of human research community”. In: *arXiv preprint arXiv:2412.17767* (2024).

## A. Appendix

### A.1. Authors

#### Leading Author

Yiheng Wang, Yixin Chen

#### Project Contributor

Shuo Li, Yifan Zhou, Bo Liu, Jiakang Yuan, Jia Bu, Wanghan Xu, Yuhao Zhou, Xiangyu Zhao, Zhiwang Zhou, Fengxiang Wang, Jianxin Dong, Jiaye Ge, Yinghui Zhang, Meng Li, Haodong Duan, Songyang Zhang, Boming Shi

#### Community Contributor

Jun Yao, Han Deng, Yizhou Wang, Jiabei Xiao, Jiaqi Liu, Encheng Su, Yujie Liu, Weida Wang, Junchi Yao, Haoran Sun, Runmin Ma, Bo Zhang, Dongzhan Zhou, Shufei Zhang, Peng Ye, Xiaosong Wang, Shixiang Tang

#### Corresponding Authors

Wenlong Zhang, Lei Bai

### A.2. Full Evaluation Results Across Core Benchmarks

Table 5 | Evaluation of LLMs on scientific multimodal benchmarks: SLAKE, SFE, and MSEarth.

Model	SLAKE	SFE	MSEarth
<b>Closed-Weight MLLMs</b>			
Claude 4.5 Sonnet	57.87	43.64	56.11
Claude4-1-Opus	58.25	45.19	58.66
GPT-5.1	54.10	33.05	58.73
GPT-5	59.94	42.44	61.46
GPT-4o	52.78	25.93	57.97
GPT-o3	55.23	32.84	59.27
Gemini-3-Pro	66.54	55.62	66.49
Gemini-2.5-Flash	55.98	38.20	57.22
Gemini-2.5-Pro	52.12	43.76	61.28
Grok-2-vision-1212	64.00	25.04	51.76
Seed1.6-vision	65.79	44.75	57.11
<b>Open-Weight MLLMs</b>			
GLM-4.5V	59.10	38.57	51.04
InternS1	60.89	45.73	56.47
Llama 4 Maverick	56.74	36.83	55.39
Qwen3-VL-235B-A22B	72.29	38.35	50.83
Qwen3-Max	24.51	20.40	49.86

Table 6 | Evaluation of large language models across five scientific benchmarks.

Model	ChemBench	ClimaQA	EarthSE	ProteinLM	MaScQA
<b>Closed-Weight MLLMs</b>					
Claude 4.5 Sonnet	69.20	71.27	66.60	15.57	83.23
Claude4-1-Opus	69.02	71.94	71.80	11.65	83.23
GPT-5.1	63.90	76.45	86.44	68.22	61.08
GPT-5	70.42	81.62	80.20	59.22	93.54
GPT-4o	61.67	72.37	58.90	59.43	61.54
GPT-o3	69.81	81.77	81.40	66.10	93.85
Gemini-3-Pro	73.08	83.08	67.89	21.19	91.23
Gemini-2.5-Flash	42.32	75.12	72.30	4.34	62.15
Gemini-2.5-Pro	70.10	78.29	72.00	0.11	83.23
Grok-2-vision-1212	60.77	63.04	61.60	20.97	49.08
Seed1.6-vision	64.64	75.48	75.80	63.88	68.92
<b>Open-Weight MLLMs</b>					
Ling-flash-2.0	59.87	66.80	72.30	22.78	63.69
DeepSeek-R1	45.97	77.28	49.00	0.00	44.46
GLM-4.5V	49.35	71.38	35.60	60.81	57.08
InternS1	65.87	75.28	73.80	46.72	81.69
Kimi-k2	65.83	77.04	78.50	28.71	71.38
Llama 4 Maverick	66.44	68.06	56.00	19.70	80.77
Qwen3-VL-235B-A22B	63.68	77.39	75.20	50.11	81.85
Qwen3-Max	64.72	76.38	75.90	40.00	69.54

Table 7 | Evaluation of large language models across five scientific benchmarks.

Model	TRQA	CMPHYSBENCH	PHYSICS	ResearchBench	SciCode	AstroVisBench
<b>Closed-Weight MLLMs</b>						
Claude 4.5 Sonnet	58.14	43.83	36.90	56.10	9.23	34.23
Claude4-1-Opus	57.56	43.57	33.80	29.47	10.77	39.87
GPT-5.1	59.30	26.92	37.96	41.45	9.23	42.02
GPT-5	59.30	45.52	34.30	45.67	13.85	44.57
GPT-4o	51.16	35.37	28.80	33.04	4.62	30.73
GPT-o3	63.37	40.98	35.30	34.14	7.69	42.82
Gemini-3-Pro	59.88	55.97	34.40	61.51	16.92	42.23
Gemini-2.5-Flash	46.51	38.69	25.45	40.86	4.62	31.94
Gemini-2.5-Pro	52.33	38.17	31.75	50.73	6.15	43.40
Grok-2-vision-1212	45.35	28.27	28.15	49.63	7.69	33.51
Seed1.6-vision	45.93	43.18	35.30	45.00	10.77	32.21
<b>Open-Weight MLLMs</b>						

Model	TRQA	CMPHYSBENCH	PHYSICS	ResearchBench	SciCode	AstroVisBench
Ling-flash-2.0	34.88	46.35	29.60	50.29	3.08	48.12
DeepSeek-R1	48.26	0.00	33.50	49.73	0.00	0.11
GLM-4.5V	42.44	0.00	26.85	42.23	6.15	0.32
InternS1	53.49	28.60	34.65	37.45	4.62	29.55
Kimi-k2	53.49	44.78	32.40	42.28	6.15	35.57
Llama 4 Maverick	52.33	44.29	33.65	38.31	4.62	31.91
Qwen3-VL-235B-A22B	47.67	64.75	35.10	40.62	6.15	29.84
Qwen3-Max	52.33	46.13	35.95	42.12	12.31	75.64

To provide a more complete view of model capabilities, we present the full quantitative results of all evaluated models across both multimodal and text-only scientific tasks.

Table 5 summarizes the performance of 19 representative models on three key multimodal scientific benchmarks: SLAKE, SFE, and MSEarth, covering visual reasoning on medical and scientific diagrams, code-based visualization, and entity localization.

Table 7 reports model-level performance on non-multimodal scientific reasoning tasks, such as causal inference, knowledge retrieval, mathematical derivation, and scientific QA, providing a comprehensive comparison of large language models across diverse text-only settings.

## B. Benchmark Description

**MSEarth.** MSEarth [2] is a multimodal benchmark designed to assess scientific understanding in Earth science by integrating visual figures with textual reasoning. Its questions are drawn from curated figures and refined captions across atmosphere, cryosphere, hydrosphere, lithosphere, and biosphere, emphasizing not only perceptual recognition but also domain-informed inference. It serves as a representative task for evaluating multimodal scientific understanding and geophysical causal reasoning and reaches graduate-level.

**SLAKE.** SLAKE [3] is a medical visual question-answering benchmark designed for multimodal scientific reasoning in clinical imaging. Each instance is annotated with semantic segmentation masks and bounding-boxes for key organs or structures, and is further structured around a built knowledge graph relational triples mapping organs, functions and disease entities. The dataset comprises real-world medical images drawn from three primary imaging modalities—CT, X-ray, and MRI—spanning anatomical regions such as the brain, neck, chest, abdomen, and pelvic cavity.

**SFE.** The SFE [4] benchmark is a multimodal, multilingual evaluation suite designed to assess the scientific cognitive capacities of advanced models across perception, understanding, and reasoning. Drawing from authentic scientific raw data formats, SFE spans five high-value disciplines (astronomy, chemistry, earth science, life science, and materials science) and comprises 66 expert-curated tasks and 830 verified visual question-answer pairs. The tasks are structured across three hierarchical cognitive levels—signal perception, attribute understanding, and comparative reasoning. By requiring models to process real scientific imagery and textual context, and to reason about them at an expert level, SFE pushes beyond superficial knowledge retrieval toward genuine scientific reasoning.

**AstroVisBench.** AstroVisBench [5] is a code-centric benchmark designed to assess large language models' capabilities in implementing scientific workflows and generating research-quality visualizations within the domain of astronomy. Drawing from 110 publicly available Jupyter notebooks curated for astronomy research workflows, AstroVisBench require models to generate executable code for data processing and produce scientific visualizations conforming to domain standards.

**SciCode.** SciCode [6] is a scientist-curated benchmark that evaluates a model's ability to translate natural-language research problems into executable Python solutions. Drawn from 80 challenging research tasks that span 16 sub-fields including mathematics, physics, chemistry, biology, and materials science, the benchmark decomposes each task into fine-grained sub-problems and each accompanied by gold-standard reference

implementations and unit-test suites. SciCode emphasises realistic scientific workflows where problems often require domain-specific knowledge recall, multi-step reasoning, and calls to external scientific libraries.

**ChemBench.** ChemBench [7] focuses on assessment of LLM's chemical knowledge comprehension and discipline-based reasoning across chemistry and materials science. Its subdomains ranges from general chemistry to more specialized fields such as inorganic, analytical and physical chemistry. Serving as a comprehensive probe of chemical intelligence, ChemBench offers fine-grained insights into models' strength and weakness, making it an

**ClimaQA.** ClimaQA [8] transforms graduate-level climate science textbooks into scientifically grounded questions with domain-expert refinement. The benchmark contains two complementary subsets which are ClimaQA-Gold, manually curated and validated by experts, and ClimaQA-Silver, programmatically generated but aligned with the same scientific rigor, covering multiple QA formats such as multiple-choice, cloze-style, and free-form reasoning. By grounding tasks in authentic scientific content and emphasizing conceptual understanding, causal reasoning, and domain-specific inference, ClimaQA provides a more faithful lens to assess a model's capability for climate science reasoning beyond general QA performance.

**EarthSE.** EarthSE is developed to systematically probe LLMs' competencies across the full breadth of Earth-science disciplines, covering five major spheres of Earth systems and 114 subfields, with diverse task formats tailored to evaluate foundational knowledge and domain-specific reasoning. Earth-Silver in particular is curated to represent professional-level difficulty, intended to test models' depth of Earth-science knowledge and capability for scientific exploration.

**ProteinLMBench.** ProteinLMBench [9] aggregates a curated set of tasks drawn from widely used protein-analysis datasets covering protein-based property prediction, protein descriptions, and protein sequence understanding and comprises 944 six-choice questions. Every item interleaves natural-language context with an amino-acid sequence span, forcing the model to align textual clues with residue patterns rather than rely on surface keyword cues. By sequence questions, ProteinLMBench offers a fine-grained, domain-specific complementary to broader biomolecular benchmarks.

**MaScQA.** MaScQA [10] is a specialized question-answering benchmark designed to evaluate large language models' understanding and reasoning capabilities in materials science and metallurgical engineering. The questions are categorized into 14 domains: thermodynamics, atomic structure, mechanical behaviour, materials manufacturing, material applications, phase transition, electrical properties, material processing, transport phenomenon, magnetic properties, material characterization, fluid mechanics, material testing, and miscellaneous. This fine-grained categorization enables targeted evaluation of LLMs' competence across diverse subfields, reflecting the interdisciplinary nature of modern materials science.

**TRQA.** TRQA [11] is a benchmark designed to evaluate large language models on biomedical reasoning and literature-based inference, designed to evaluate biomedical reasoning across literature evidence and real-world drug pipeline data. The benchmark targets a broad set of core capabilities: scientific planning, literature-grounded information retrieval, tool selection, reasoning toward biological conclusions, and critical self-evaluation.

**CMPHysBench.** CMPHysBench [12] is designed to evaluate large language models' scientific reasoning abilities in condensed matter physics, with a particular focus on symbolic derivation, algebraic manipulation and physical interpretation. Its questions span core subfields of condensed matter physics including magnetism, superconductivity, strongly correlated systems, etc. CMPHysBench serves as an evaluation to benchmark the depth of scientific reasoning in LLMs, particularly their ability to operate with formal representations, physical principles, and structured scientific logic.

**PHYSICS.** PHYSICS [13] evaluates the ability of language models to perform undergraduate-level physical reasoning spanning Mechanics, Electromagnetism, Thermodynamics, Optics, and Modern Physics. Each problem demands a combination of domain knowledge and symbolic reasoning ranging from conceptual understanding and detailed calculations.

**ResearchBench.** ResearchBench [14] aims to evaluate the capability of language models to simulate research workflows, particularly focusing on paper writing and peer-review generation. The benchmark consists of 1,000 paper-writing tasks and 200 review-writing tasks, all curated from recent top-tier machine learning venues such as NeurIPS 2024 and ICLR 2024. This benchmark assesses a model's ability to generate research ideas, understand technical contributions and provide coherent scientific writing aligned with the realistic academic standards.

## C. Representative Task Cases

### C.1. MaScQA

MaScQA — Thermodynamic Reasoning in Materials Science

**Question:**

A student aims to deposit a thin metallic film on  $SiO_2$  substrate, with an adhesion layer between the metal film and substrate, in a contiguous planar fashion. Island type of growth must be avoided. The student performs an extensive optimization exercise. Which one of the following steps is in the right direction? (A) Choose a metallic adhesion layer with very low interfacial energy with the deposited thin film (B) Choose a metallic adhesion layer with very low interfacial energy with  $SiO_2$ , irrespective of its interaction with metal film to be deposited (C) Increase the substrate temperature and decrease the deposition rate (D) Use intermittent stages of deposition followed by annealing

**Question Type:** Multi-Choice Question

**Answer:** A

**Response from Gemini 3 Pro:**

To achieve the growth of a contiguous planar thin metallic film and avoid island (Volmer-Weber) growth, one must ensure that the deposited film "wets" the surface beneath it.

\*\*Theoretical Analysis:\*\* The growth mode is determined by the balance of surface and interfacial energies, described by Young's equation and the spreading parameter  $S$ . For a film ( $f$ ) depositing on a substrate ( $s$ ), the condition for layer-by-layer (planar) growth is:

$$\gamma_s \geq \gamma_f + \gamma_{int}$$

where: \*  $\gamma_s$  is the surface energy of the substrate (in this case, the adhesion layer). \*  $\gamma_f$  is the surface energy of the depositing metal film. \*  $\gamma_{int}$  is the interfacial energy between the film and the substrate.

\*\*Evaluating the Options:\*\*

\* \*\*(A) Choose a metallic adhesion layer with very low interfacial energy with the deposited thin film:\*\* In this scenario, the "substrate" for the top metal film is the adhesion layer. To satisfy the wetting condition ( $\gamma_{AdhesionLayer} \geq \gamma_{MetalFilm} + \gamma_{Interface}$ ), it is crucial to minimize  $\gamma_{Interface}$ . A very low interfacial energy implies strong bonding (affinity) between the adhesion layer and the metal film, which promotes wetting and planar growth. This is the correct direction.

\* \*\*(B) Choose a metallic adhesion layer with very low interfacial energy with  $SiO_2$ , irrespective of its interaction with metal film to be deposited:\*\* While the adhesion layer must adhere to the  $SiO_2$ , ignoring the interaction with the top metal film is incorrect. If the interfacial energy between the adhesion layer and the top metal film is high (poor bonding), the inequality  $\gamma_{AdhesionLayer} \geq \gamma_{MetalFilm} + \gamma_{Interface}$  will likely not be met, leading to island growth of the metal film on top of the adhesion layer.

\* \*\*(C) Increase the substrate temperature and decrease the deposition rate:\*\* High substrate temperatures increase surface diffusivity. In systems that are not perfectly wetting (which is common for metals on dissimilar surfaces), high mobility allows atoms to aggregate into thermodynamically stable 3D islands (agglomeration). To force a continuous film in a system that might tend toward islanding, one typically uses \*lower\* substrate temperatures and \*higher\* deposition rates to increase nucleation density and freeze the atoms in a continuous layer (kinetic control). Increasing temperature generally promotes islanding and dewetting for thin metal films.

\* \*\*(D) Use intermittent stages of deposition followed by annealing:\*\* Annealing provides thermal energy which drives the system toward thermodynamic equilibrium. For thin metallic films, the equilibrium state is often a dewetted set of droplets (islands) to minimize total surface energy (Rayleigh instability). Therefore, annealing usually destroys the continuity of a thin film rather than promoting it.

\*\*Conclusion:\*\* Only option (A) correctly addresses the thermodynamic requirements to prevent island growth and promote a contiguous planar film.

[A]

**Judgment:** Correct.

**Note:** Gemini 3 Pro is the only model that answers this question correctly.

## C.2. Chembench

Chembench — Analytical Chemistry (Molecular Vibrations)

**Question:**

What is the total number of normal vibrational modes for a C<sub>60</sub> molecule?

- (A) 174 (B) 180 (C) 54 (D) 60

**Question Type:** Multi-Choice Question

**Answer:** A

**Response from Gemini 3 Pro:**

To determine the total number of normal vibrational modes for a C<sub>60</sub> molecule, we follow these steps:

- Identify the number of atoms (N):** The molecule is C<sub>60</sub>, which consists of 60 carbon atoms. So, N = 60.
- Calculate the total degrees of freedom:** The total degrees of freedom for a molecule with N atoms is given by 3N.

$$3 \times 60 = 180$$

**3. Subtract translational and rotational degrees of freedom:**

- There are always 3 translational degrees of freedom.
- Since C<sub>60</sub> is a non-linear molecule (it has a spherical cage-like structure), it has 3 rotational degrees of freedom.
- The formula for the number of vibrational modes for a non-linear molecule is 3N – 6.

**4. Perform the calculation:**

$$\text{Vibrational modes} = 3(60) - 6 = 180 - 6 = 174$$

Therefore, the total number of normal vibrational modes is 174.

[ANSWER] A [/ANSWER]

**Judgment:** Correct.

**Note:** Gemini 3 Pro accurately applies the degrees of freedom formula (3N – 6) for non-linear molecules to solve this structural chemistry problem.

Chembench — Analytical Chemistry (NMR Spectroscopy)

**Question:**

Determine the number of peaks in the 1H NMR spectrum of the molecule with the SMILES OCC1=CNC(=O)C1=O.

**Question Type:** Numeric Question

**Answer:** 4

**Response from Gemini 3 Pro:**

To determine the number of peaks in the 1H NMR spectrum of the molecule with the SMILES

OCC1=CNC(=O)C1=O, we first analyze its structure.

### 1. SMILES Parsing:

- **OCC**: A hydroxymethyl group (HO-CH2-) attached to the ring anchor.
- **Ring Structure**: The connectivity corresponds to a **1H-pyrrole-2,3-dione** core substituted at position 4.
- **Substituent**: A hydroxymethyl group (-CH2OH) at position 4.

### 2. Identify Proton Environments:

- **Hydroxyl Proton (-OH)**: The proton on the oxygen of the hydroxymethyl group. This represents **1** distinct chemical environment.
- **Methylene Protons (-CH<sub>2</sub>-)**: The two protons on the carbon of the hydroxymethyl group. Since the molecule has no chiral center and the ring is planar, these two protons are enantiotopic (equivalent). They represent **1** distinct chemical environment.
- **Vinylic Proton (=CH-)**: The proton on the ring carbon at position 5 (adjacent to the nitrogen). This represents **1** distinct chemical environment.
- **Amide/Imide Proton (-NH-)**: The proton on the ring nitrogen. This represents **1** distinct chemical environment.

3. Count the Peaks: There are 4 distinct sets of equivalent protons: (1) -OH, (2) -CH<sub>2</sub>-, (3) =CH-, (4) -NH-. Assuming a standard solvent where exchangeable protons are visible, we expect 4 signals.

[ANSWER]4[/ANSWER]

Judgment: **Correct**.

Note: Gemini 3 Pro correctly parses the SMILES string, identifies equivalent protons based on molecular symmetry, and provides the exact integer count required.

## C.3. LLM4Chem

### LLM4Chem — Property Prediction

#### Question:

Predict the octanol/water distribution coefficient logD under the circumstance of pH 7.4 for <SMILES> NC(=O)C1=CC=CC=C1O </SMILES> .

Answer: **1.090**

## C.4. SciCode

### SciCode

#### Question:

Create a function to solve the linear system  $\mathbf{Ax} = \mathbf{b}$  using the conjugate gradient method. This function takes a matrix  $\mathbf{A}$  and a vector  $\mathbf{b}$  as inputs.

Answer:

```
def cg(A, b, x, tol):
    """
    Inputs:
    A : Matrix, 2d array size M * M
    b : Vector, 1d array size M
    x : Initial guess vector, 1d array size M
    tol : tolerance, float
    Outputs:
```

```

x : solution vector, 1d array size M
"""
# Initialize residual vector
res = b - np.dot(A, x)
# Initialize search direction vector
search_direction = res.copy()
# Compute initial squared residual norm
old_res_norm = np.linalg.norm(res)
itern = 0
# Iterate until convergence
while old_res_norm > tol:
    A_search_direction = np.dot(A, search_direction)
    step_size = old_res_norm**2 / np.dot(search_direction,
                                         A_search_direction)

    # Update solution
    x += step_size * search_direction
    # Update residual
    res -= step_size * A_search_direction
    new_res_norm = np.linalg.norm(res)

    # Update search direction vector
    search_direction = res + (new_res_norm / old_res_norm)**2 *
                                         search_direction

    # Update squared residual norm for next iteration
    old_res_norm = new_res_norm
    itern += 1
return x

```

## C.5. PHYSICS

### PHYSICS

#### Question:

To make a flat stone skip across the water surface when thrown quickly, the stone may bounce and fly towards the distance, commonly known as "stone skipping." To achieve the "stone skipping" effect, the angle between the direction of the stone's velocity and the water surface at the point of contact must not exceed  $\theta$ . To observe "stone skipping," a student throws a stone horizontally from a height  $h$  above the water surface. What is the minimum launch velocity required? (Neglect air resistance during the stone's flight, and the acceleration due to gravity is  $g$ .)

Answer:  $\frac{\sqrt{2gh}}{\tan \theta}$

## C.6. CMPhysBench

### CMPhysBench - Theoretical Foundations

#### Question:

A particle of mass  $m$  is in the ground state of a one-dimensional harmonic oscillator potential

$$V_1(x) = \frac{1}{2}kx^2, \quad k > 0$$

When the spring constant  $k$  suddenly changes to  $2k$ , the potential then becomes

$$V_2(x) = kx^2$$

Immediately measure the energy of the particle, and find the expression for the probability of the particle being in the ground state of the new potential  $V_2$ .

**Answer:**

(a) The wave function of the particle  $\psi(x, t)$  should satisfy the time-dependent Schrödinger equation

$$i\hbar \frac{\partial}{\partial t} \psi = -\frac{\hbar^2}{2m} \frac{\partial^2}{\partial x^2} \psi + V\psi \quad (3)$$

When  $V$  undergoes a sudden change (from  $V_1 \rightarrow V_2$ ) but with a finite change quantity,  $\psi$  remains a continuous function of  $t$ , implying that  $\psi$  does not change when  $V$  changes abruptly.

Denote  $\psi_0(x)$  and  $\phi_0(x)$  as the ground state wave functions of the potential  $V_1$  and  $V_2$ , respectively. After the potential suddenly changes from  $V_1$  to  $V_2$ , the wave function of the particle remains  $\psi_0$ . The probability of measuring the particle in the state  $\phi_0$  is  $|\langle \psi_0 | \phi_0 \rangle|^2$ .

Rewrite  $V_1$  and  $V_2$  in standard form:

$$V_1(x) = \frac{1}{2}kx^2 = \frac{1}{2}m\omega_1^2x^2 \quad (4)$$

$$V_2(x) = kx^2 = \frac{1}{2}m\omega_2^2x^2 \quad (5)$$

It is clear that

$$\omega_2 = \sqrt{2}\omega_1 \quad (4)$$

$\psi_0$  and  $\phi_0$  can be expressed as in formulas (3) and (5) from problem 3.2, namely

$$\begin{aligned} \psi_0(x) &= \left( \frac{\alpha}{\sqrt{\pi}} \right)^{\frac{1}{2}} e^{-\alpha^2 x^2 / 2}, & \alpha^2 &= m\omega_1/\hbar \\ \phi_0(x) &= \left( \frac{\beta}{\sqrt{\pi}} \right)^{\frac{1}{2}} e^{-\beta^2 x^2 / 2}, & \beta^2 &= m\omega_2/\hbar \end{aligned} \quad (6)$$

where

$$\beta^2/\alpha^2 = \omega_2/\omega_1 = \sqrt{2} \quad (7)$$

Thus

$$\begin{aligned} \langle \psi_0 | \phi_0 \rangle &= \sqrt{\frac{\alpha\beta}{\pi}} \int_{-\infty}^{+\infty} e^{-\frac{1}{2}(\alpha^2+\beta^2)x^2} dx = \left( \frac{2\alpha\beta}{\alpha^2+\beta^2} \right)^{\frac{1}{2}} \\ |\langle \psi_0 | \phi_0 \rangle|^2 &= \frac{2\alpha\beta}{\alpha^2+\beta^2} = \frac{2\beta/\alpha}{1+\beta^2/\alpha^2} = \frac{2^{5/4}}{1+\sqrt{2}} = 0.9852 \end{aligned} \quad (8)$$

This is the required probability. (b) Consider the time when the potential changes for the first time ( $V_1 \rightarrow V_2$ ) as  $t = 0$ , then the wave function is

$$\psi(x, 0) = \psi_0(x) \quad (9)$$

Let  $\phi_n(x)$  denote the energy eigenstates of the potential  $V_2$ , corresponding to energy levels

$$E_n = (n + \frac{1}{2})\hbar\omega_2$$

Expand  $\psi_0$  as a linear combination of  $\phi_n$ ,

$$\psi_0(x) = \sum_n C_n \phi_n(x), \quad (n \text{ can only take even values}) \quad (10)$$

For  $0 < t < \tau$ , in Schrödinger equation (3)  $V = V_2(x)$ , its solution is

$$\begin{aligned} \psi(x, t) &= \sum_n C_n \phi_n(x) e^{-iE_n/\hbar} \\ &= e^{-i\omega_2 t \cdot 2} \sum_n C_n \phi_n(x) e^{-i\omega_{\omega_2} t} \end{aligned} \quad (11)$$

To have  $\psi(x, \tau) = A\psi_0(x)$ , it must hold that

$$e^{-in\omega_2\tau} = 1, \quad n = 0, 2, 4, \dots \quad (12)$$

That is

$$e^{i\omega_2\tau} = \pm 1 \quad (12')$$

The  $\tau$  that satisfies this condition is

$$\tau = l\pi/\omega_2 = l\pi\sqrt{\frac{m}{2k}}, \quad l = 1, 2, 3, \dots \quad (13)$$

When  $t = \tau$ , after the potential changes from  $V_2$  back to  $V_1$ , the particle remains in the state  $\psi_0$ , with energy  $E = \hbar\omega_1/2$ .

## C.7. ClimaQA

### ClimaQA

**Question:**

Why does the reflectivity increase between points 3 and 4 in the vertical profile of radar data for stratiform precipitation?

- A) Due to the reduction in relative humidity below cloud, causing evaporation of raindrops.
- B) As a result of the snowflakes becoming more compact and collapsing due to increased terminal velocity.
- C) Due to the reduction in size of snowflakes as they melt and acquire a liquid layer.
- D) Because of the increase in the concentration of ice crystals in that layer.

**Question Type:** Multi-Choice Question

**Answer:** C

## C.8. EarthSE

### EarthSE

**Question:**

How do varying fire regimes and browsing intensities influence vegetation dynamics and species composition in Mediterranean ecosystems, particularly in scenarios modeled near Lago di Massaciuccoli during the Holocene, and what implications do interacting disturbances have on drought-sensitive species such as *Abies alba* and the dominance of *Quercus ilex* under different climatic conditions?

**Question Type:** Analysis

**Answer:**

Varying fire regimes and browsing intensities serve as critical factors in shaping vegetation dynamics and species composition in Mediterranean ecosystems. In simulations near Lago di Massaciucoli, small, infrequent fires coupled with low browsing intensity promote the coexistence of *Abies alba* and *Quercus ilex* in mixed forests. However, large, infrequent fires decrease overall biomass and favor shade-tolerant species such as evergreen oaks and deciduous hardwoods. High fire frequency paired with warmer-than-present summer temperatures significantly reduces *A. alba* biomass, suggesting its vulnerability to such disturbances. Browsing intensities further exacerbate this effect, with moderate browsing eliminating *A. alba* regeneration entirely, favoring species like *Q. ilex* and shade-intolerant evergreen and deciduous competitors. Interactions between fire, browsing, and climate highlight complex dynamics, where increased light availability and reduced competition from *A. alba* under higher disturbance levels allow *Q. ilex* to dominate. These findings emphasize the combined impact of anthropogenic activity and climatic changes on biodiversity and forest compositional shifts in Mediterranean regions.

## C.9. ProteinLMBench

### ProteinLMBench

**Question:**

Analyze the sequence of the following enzyme protein, and step by step deduce which reaction(s) this enzyme is involved in?

```
<seq> M S I P S S Q Y G F V F N K Q S G L K L R N D L P V H K P K A G Q L L L K V D A V G L C H S D
L H V I Y E G L D C G D N Y V M G H E I A G T V A A V G D D V I N Y K V G D R V A C V G P N G C
G G C K Y C R G A I D N V C K N A F G D W F G L G Y D G G Y Q Q Y L L V T R P R N L S R I P D
N V S A D V A A A S T D A V L T P Y H A I K M A Q V S P T S N I L L I G A G G L G G N A I Q V A
K A F G A K V T V L D K K K E A R D Q A K K L G A D A V Y E T L P E S I S P G S F S A C F D F V
S V Q A T F D V C Q K Y V E P K G V I M P V G L G A P N L S F N L G D L A L R E I R I L G S F W
G T T N D L D D V L K L V S E G K V K P V V R S A K L K E L P E Y I E K L R N N A Y E G R V V F N P </seq>
```

**Question Type:** Analysis

**Answer:**

Based on the analysis of the provided protein sequence, the function of this protein is described as follows: A zinc protein. Acts on primary or secondary alcohols or hemi-acetals with very broad specificity; however the enzyme oxidizes methanol much more poorly than ethanol. The animal, but not the yeast, enzyme acts also on cyclic secondary alcohols. So, this enzyme should belong to the alcohol:NAD<sup>+</sup> oxidoreductase. Therefore, by synthesizing the analysis of this enzyme, we can determine the reaction(s) it is involved in as follows: (1) a primary alcohol + NAD<sup>+</sup> = an aldehyde + NADH + H<sup>+</sup> (2) a secondary alcohol + NAD<sup>+</sup> = a ketone + NADH + H<sup>+</sup>.

## C.10. TRQA

### TRQA

**Question:**

What is the highest global development stage achieved by drugs targeting PDE2A for cognitive impairment in schizophrenia up to April 2025, and which drugs are in this stage?

**Answer:** Phase I; BI 474121

## C.11. ResearchBench

### ResearchBench

**Question:**

You are helping with the scientific hypotheses generation process. We in general split the period of conducting research into four steps. Firstly it's about finding a good and specific background research question, and an introduction of the previous methods under the same topic; Secondly its about finding inspirations (mostly from literatures), which combined with the background research question, can lead to a impactful research hypothesis; Thirdly it's hypothesis generation based on the background research question and found inspirations; Finally it's about designing and conducting experiments to verify hypothesis. An example is the backpropagation of neural networks. In backpropagation, the research question is how to use data to automatically improve the parameters of a multi-layer logistic regression, the inspiration is the chain rule in mathematics, and the research hypothesis is the backpropagation itself. In their paper, the authors have conducted experiments to verify their hypothesis. Now we have identified a good research question, and we have found a core inspiration in a literature for this research question. But one inspiration might not have enough information to support a finding of novel, valid, and significant research hypothesis. Therefore we also have found a series of inspiration candidates, which might provide additional useful information to assist the core inspiration for the next step of hypothesis generation. Please help us generate a novel, valid, and significant research hypothesis based on the background research question and the inspirations. The background research question is: What are the processes and sources responsible for the formation and evolution of dust in the early universe, particularly at redshifts greater than 4?

The introduction of the previous methods is: Survey not provided. Please overlook the survey.

The core inspiration is: Title: Dust formation in primordial Type II supernovae; Abstract: We have investigated the formation of dust in the ejecta of Type II supernovae (SNe), mostly of primordial composition, to answer the question of where are the first solid particles formed in the universe. However, we have also considered non-zero progenitor's metallicity values up to  $Z = Z_{\odot}$ . The calculations are based on standard nucleation theory and the scheme has been first tested on the well studied case of SN1987A, yielding results that are in agreement with the available data. We find that: i) the first dust grains are predominantly made of silicates, amorphous carbon (AC), magnetite, and corundum; ii) the largest grains are the AC ones, with sizes around  $300\text{ }\mu\text{m}$ , whereas other grain types have smaller radii, around  $10\text{--}20\text{ }\mu\text{m}$ . The grain size distribution depends somewhat on the thermodynamics of the ejecta expansion and variations in the results by a factor 2 might occur within reasonable estimates of the relevant parameters. Also, and for the same reason, the grain size distribution, is essentially unaffected by metallicity changes. The predictions on the amount of dust formed are very robust: for  $Z = 0$ , we find that SNe with masses in the range  $(12\text{--}35)\text{M}_{\odot}$  produce about  $0.08\text{M}_{\odot} < M_{\text{dust}} < 0.3\text{M}_{\odot}$  of dust/SN. The above range increases by roughly 3 times as the metallicity is increased to solar values. We discuss the implications and the cosmological consequences of the results..

The additional inspiration candidates are: - Composition and quantities of dust produced by AGB-stars and returned to the interstellar medium - The evolution of refractory interstellar grains in the solar neighborhood

Now you have seen the background research question, the core inspiration, and many potential additional inspiration candidates. Please try to generate a novel, valid, and significant research hypothesis based on the background research question and the inspirations. (response format: 'Hypothesis: Reasoning Process: ')

**Answer:** The authors propose that multiple stellar sources and processes, particularly supernovae and asymptotic giant branch (AGB) stars, significantly contribute to the formation and evolution of dust in the early universe. They suggest that the combination of these sources, along with dust growth in the interstellar medium, can account for the observed dust masses in early galaxies.

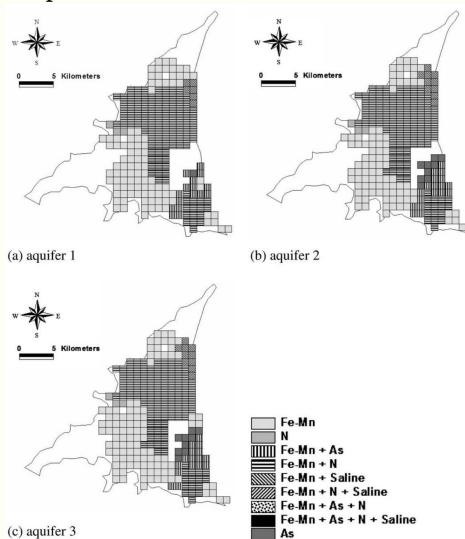
## C.12. MSEarth

### MSEarth

#### Question:

Delineation of hazardous regions for the nine classifications for  $P(u; z_p | (n)) \geq 0.5$  Fe–Mn is the Fe–Mn hazard; N is the nitrogen hazard; As is the As hazard; saline is the saline hazard.

Which aquifer shows the highest spread of the Fe–Mn hazard?



Answer: C

## C.13. AstroVisBench

### AstroVisBench — Data Processing and Visualization in Astrophysics

#### Question:

**Setup:** To set up the environment for our analysis, import the necessary astronomical and data handling packages that will allow us to perform sky coordinate calculations, access and handle data tables, and make HTTP requests to retrieve data. Additionally, configure the environment to properly display tables and set up appropriate input and output stream handling capabilities for image processing. Make sure that each imported module serves its purpose, whether that's astrometric calculations, pretty printing, image handling or data visualization.

**Processing:** For this analysis, utilize a series of functions to query an astronomical database using an API. These functions should enable you to perform a precise HSC catalog search using positional queries with an angular search radius, while also being able to retrieve metadata about the tables and columns available in the database. Implement the necessary error checking to ensure the compatibility of request parameters, and validate that query format and table specification adhere to the accepted values. Then, use the astropy name resolver to acquire the celestial coordinates of a specified astronomical target and filter the data to select only those entries that meet certain criteria, such as having a minimum number of measurements in specified bands.

**Visualization:** Plot the celestial coordinates of the extracted objects on a 2D visualization to represent their positions in relation to the center of the galaxy. Use distinct markers to denote both the location of the galaxy center and the objects obtained from the catalog search, ensuring that their sizes and colors are visibly distinct to differentiate the two. Set appropriate axes labels for Right Ascension and Declination, with an inverted x-axis to align with conventional celestial plotting practices. Add a legend to facilitate the understanding of the symbols used in the plot, and ensure aspect ratio maintenance for

accurate representation.

**Question Type:** Code Generation (Multi-step: Processing & Visualization)

**Answer:** `Code`

**Response from Intern-S1:**

**Processing Code:**

```

def get_target_coordinates(target_name):
    """Resolve the target name to RA and Dec in degrees."""
    try:
        coord = SkyCoord.from_name(target_name)
        return coord.ra.deg, coord.dec.deg
    except Exception as e:
        print(f"Error resolving target name: {e}")
        raise

def query_hsc_catalog(ra, dec, radius, constraints):
    """Query the HSC catalog with a cone search and
    constraints."""
    if radius > 2.0:
        raise ValueError("Radius must be <= 2 degrees for HSC
        cone search.")

    base_url =
        "https://catalogs.mast.stsci.edu/api/v0.1/hsc/cone.csv"
    params = {
        'ra': ra,
        'dec': dec,
        'radius': radius,
        'where': constraints,
        'format': 'csv'
    }

    response = requests.get(base_url, params=params)
    if response.status_code == 200:
        table = Table.read(BytesIO(response.content),
                           format='csv')
        return table
    else:
        print(f"Error in API request: {response.status_code}")
        print(response.text)
        response.raise_for_status()

def get_hsc_metadata():
    """Retrieve metadata about HSC tables and columns."""
    metadata_url =
        "https://catalogs.mast.stsci.edu/api/v0.1/hsc/metadata"
    response = requests.get(metadata_url)
    if response.status_code == 200:
        return response.json()
    else:
        print(f"Error retrieving metadata:")

```

```

    {response.status_code})
response.raise_for_status()

# Example usage for IC 1613 with 0.5 degree radius and
   measurement constraints
target_name = "IC 1613"
radius_deg = 0.5
constraints = "A_F475W_N >= 10 AND A_F814W_N >= 10"

try:
    ra, dec = get_target_coordinates(target_name)
    print(f"Coordinates for {target_name}: RA={ra:.4f},
          Dec={dec:.4f}")

    result_table = query_hsc_catalog(ra, dec, radius_deg,
                                     constraints)
    print(f"\nRetrieved {len(result_table)} entries matching
          the criteria.")
    print("First 5 entries:")
    result_table[:5].pprint(max_width=150)

except Exception as e:
    print(f"An error occurred: {e}")

```

#### Processing Execution Output (API Error):

```

Coordinates for IC 1613: RA=16.2017, Dec=2.1195
Error in API request: 404
{"detail":"Not Found"}
An error occurred: 404 Client Error: Not Found for url: ...

```

#### Visualization Code (Assuming data table 'tab' is available):

```

# Visualization

# Create a figure and axis
plt.figure(figsize=(10, 8))

# Plot the catalog objects as blue points
plt.scatter(tab['MatchRA'], tab['MatchDec'], s=20,
            c='skyblue', edgecolor='k', alpha=0.6, label='HSC Objects')

# Plot the galaxy center as a large red star
plt.scatter(ra_ic1613, dec_ic1613, s=200, c='red', marker='*',
            label='IC 1613 Center')

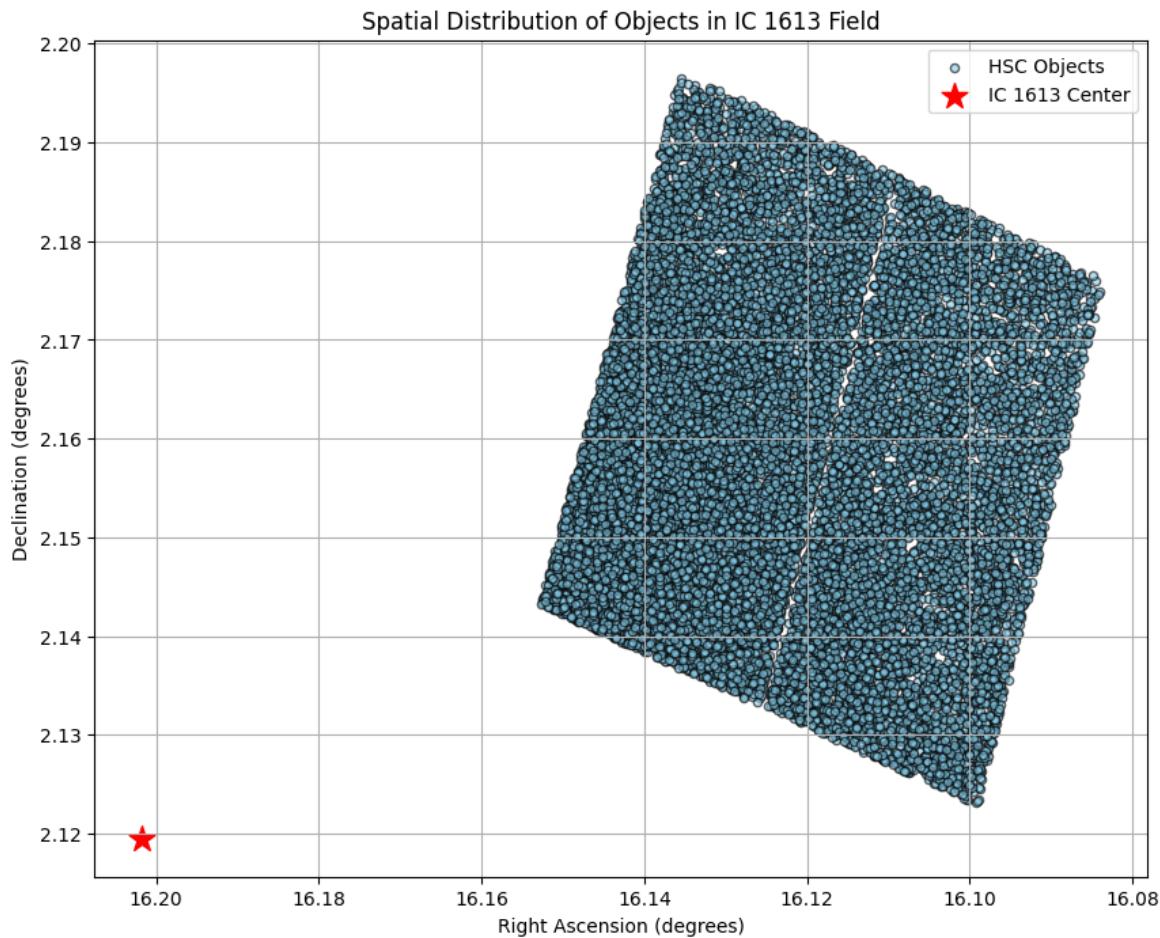
# Invert the x-axis to align with celestial plotting
   conventions
plt.gca().invert_xaxis()

# Add labels, title, and legend
plt.xlabel('Right Ascension (degrees)')
plt.ylabel('Declination (degrees)')
plt.title('Spatial Distribution of Objects in IC 1613 Field')
plt.legend()
plt.grid(True)

```

```
# Display the plot
plt.show()
```

#### Visualization Result:



#### Judgment Rationale for Visualization (V2):

Comparing the under-test image with the ground truth image, several styling differences are observed: 1) The marker style for HSC objects differs (filled circles vs. simple dots). 2) The marker size differs ( $s=20$  vs.  $markersize=1$ ). 3) The galaxy center marker differs (filled red star vs. red 'x' marker). 4) The model includes a grid and title which are absent in the ground truth. Despite these aesthetic differences, both visualizations convey the same key information: the spatial distribution of HSC measurements in relation to the IC 1613 galaxy center, with an inverted x-axis, proper axis labels, distinct markers for objects vs. center, and a legend. The overall pattern, data positions, and coordinate ranges are identical. The differences are primarily aesthetic rather than informational.

#### Overall Judgment (Processing & Visualization):

Processing Code is functionally correct but execution fails due to external API error (404 Not Found). Visualization Code executes successfully but has minor aesthetic differences compared to the target visualization.

Visualization Error Type: **Minor Error**

## C.14. SLAKE

SLAKE

**Question:**

What modality is used to take this image?

**Image:**



**Answer:** CT

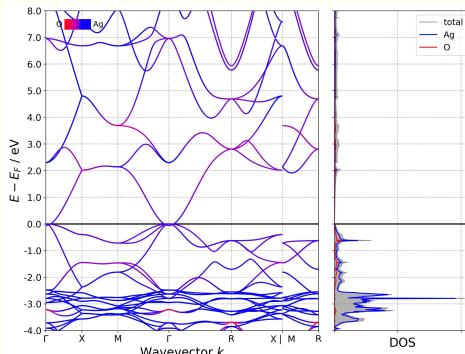
## C.15. SFE

SFE

**Question:**

Given an electronic band structure plot, analyze the positions of the valence band maximum (VBM) and conduction band minimum (CBM) to decide if the material has a direct bandgap (VBM and CBM at the same high-symmetry k-point, which you must specify), an indirect bandgap (VBM and CBM at different high-symmetry k-points, which you must specify both), or if band crossings at the Fermi level indicate a metallic state with no bandgap. Reply with a single concise narrative paragraph stating your conclusion, including the relevant high-symmetry k-points for direct or indirect cases, or for a metallic case briefly explaining that band crossings preclude any bandgap, without bullet points, lists, or additional commentary.

Based on the provided band structure <image>, determine whether the material exhibits a direct or indirect bandgap.



**Answer:** The material is metallic, doesn't exhibit a direct or indirect bandgap.