Can Multimodal LLMs See Materials Clearly? A Multimodal Benchmark on Materials Characterization

Zhengzhao Lai¹, Youbin Zheng², Zhenyang Cai¹, Haonan Lyu³, Jingpu Yang², Hongqing Liang³, Yan Hu^{1*}, Benyou Wang¹

¹The Chinese University of Hong Kong, Shenzhen

²Northeastern University ³Zhejiang University

{zhengzhaolai, huyan}@cuhk.edu.cn

Abstract

Materials characterization is fundamental to acquiring materials information, revealing the processing-microstructure-property relationships that guide material design and optimization. While multimodal large language models (MLLMs) have recently shown promise in generative and predictive tasks within materials science, their capacity to understand real-world characterization imaging data remains underexplored. To bridge this gap, we present MatCha, the first benchmark for materials characterization image understanding, comprising 1,500 questions that demand expertlevel domain expertise. MatCha encompasses four key stages of materials research comprising 21 distinct tasks, each designed to reflect authentic challenges faced by materials scientists. Our evaluation of state-of-the-art MLLMs on MatCha reveals a significant performance gap compared to human experts. These models exhibit degradation when addressing questions requiring higher-level expertise and sophisticated visual perception. Simple few-shot and chain-of-thought prompting struggle to alleviate these limitations. These findings highlight that existing MLLMs still exhibit limited adaptability to real-world materials characterization scenarios. We hope MatCha will facilitate future research in areas such as new material discovery and autonomous scientific agents. MatCha is available at https://github.com/ FreedomIntelligence/MatCha.

1 Introduction

Materials characterization serves as a critical means of obtaining information about the physical world (Leng, 2013), offering insights that transcend the limitations of human sensory perception. It enables multi-scale analysis, providing rich data on material morphology, composition, and structure. This information, in turn, reveals the physical, chemical, and mechanical properties essential

for guiding new material design and optimization (Robertson et al., 2011). For instance, scanning electron microscopy (SEM) and transmission electron microscopy (TEM) images are instrumental in determining the underlying mechanisms of steel bar fracture in buildings (Inkson, 2016). Despite its importance, interpreting diverse and complex imaging data generated by various characterization techniques remains a significant challenge that demands extensive domain expertise. Typically, even experienced materials researchers invest considerable time analyzing these multifaceted results. This analytical burden becomes particularly acute when handling high-throughput data, where efficiency bottlenecks arise. While convolutional neural networks (CNNs) have been utilized for various materials characterization tasks (Madsen et al., 2018; Maksov et al., 2019; Zaloga et al., 2020; Warmuzek et al., 2021; Leitherer et al., 2023), providing valuable insights, their application faces notable limitations. Predominantly, these CNNbased approaches are task-specific, exhibiting limited cross-task generalization. Furthermore, prior work has largely focused on morphological perception, resulting in shallow image content understanding that often falls short of the nuanced interpretations required by materials scientists in real-world scenarios. These constraints highlight the need for more versatile and deeply understanding models.

Leveraging MLLMs (Zhu et al., 2023; Liu et al., 2023; Li et al., 2023; Achiam et al., 2023; Wang et al., 2024b; Liu et al., 2024) for these challenges offers distinct advantages. MLLMs have demonstrated strong performance and generalization in both natural and domain-specific image understanding (Li et al., 2025), and have spurred revolutionary changes in materials science, including property prediction (Rubungo et al., 2023; Antunes et al., 2024; Xie et al., 2023), new material design (Tang et al., 2025; Mishra et al., 2024; Xie et al., 2024), and autonomous scientific agents (Zhang

^{*}Corresponding author.

et al., 2024; Ding et al., 2024). These models show great potential to assist materials scientists with diverse materials characterization tasks via natural language interaction, thereby facilitating new material development and boosting scientific productivity. To realize this potential, MLLMs must first be capable of accurately interpreting diverse materials characterization images, recognizing fundamental visual content, and performing reasoning based on visual cues. However, existing evaluations of MLLMs on scientific imaging data primarily focus on biomedical domains (He et al., 2020; Huang et al., 2023; Lozano et al., 2024), are confined to relatively simplistic figures and charts in limited scientific fields (Yue et al., 2024a,b; Chen et al., 2024a; Roberts et al., 2024), or lacking sufficient depth to represent the complexity of authentic materials research scenarios (Alampara et al., 2024; Li et al., 2024b; Verma et al., 2024). Consequently, a comprehensive, expert knowledge-anchored multimodal benchmark specifically designed for materials science to rigorously assess current models is notably absent. This gap impedes the progress toward AI-assisted research and autonomous scientific discovery agents.

To bridge this gap, we present MatCha, a challenging multimodal benchmark for materials characterization imaging data understanding. The core strengths of MatCha are threefold: (1) practical and realistic task design: The design philosophy of MatCha originates from real-world scientific workflows. The tasks are derived directly from the research processes of materials scientists and are designed to reflect authentic challenges in practice. (2) task diversity and broad coverage: MatCha incorporates 21 sub-tasks, each representing a concrete step within the scientific workflow. These tasks collectively cover a wide range of characterization methods and corresponding problems. (3) expert-level difficulty: MatCha includes 1,500 multiple-choice questions of varying complexity, each requiring visual understanding and expertlevel scientific expertise.

We first benchmark state-of-the-art MLLMs on MatCha under the zero-shot setting, observing a substantial performance gap between models and human experts as well as noticeable performance degradation across different task stages due to limited generalization capability. Next, we further investigate whether these performance gaps can be bridged by in-context learning or by guiding the model through a chain-of-thought (CoT) process.

The results show that while some models do benefit from these strategies, others exhibit unstable or even degraded performance, and a significant gap to human expert performance persists.

In summary, our contributions are as follows:

- We introduce MatCha, the first multimodal perception and understanding benchmark for materials characterization, which comprises 21 expert-defined tasks reflecting real scientific challenges.
- We conduct extensive experiments on various proprietary and open-source models under different settings.
- We reveal their current limitations and directions for future enhancements through detailed quantitative and qualitative analyses.

2 Related Works

Materials Characterization Analysis. Computer vision has revolutionized the extraction of visual information and quantitative analysis from material microscopy images. Applications include image-based structure recognition (Abouelatta, 2013; DeCost and Holm, 2015; Chowdhury et al., 2016; DeCost et al., 2017), detection of individual atomic sites and defects (Madsen et al., 2018; Li et al., 2018; Yang et al., 2021; Shen et al., 2021), and segmentation of microstructures or particles (DeCost et al., 2019; Roberts et al., 2019; Baskaran et al., 2020; Bals and Epple, 2023). However, existing studies largely focus on perception in specific electron microscopes, neglecting deeper imaging content and cross-modal data analysis, which results in a gap with actual scientific research problems. Additionally, material science characterization data significantly differ from natural images, thus limiting generalization between tasks. Furthermore, the scarcity of large labeled datasets impedes supervised deep learning model training (Holm et al., 2020). Our research seeks to leverage and investigate the capabilities of MLLMs to address the diverse challenges in characterization data analysis and elucidation, thereby tackling real-world problems faced by materials scientists.

Multimodal Benchmarks in Science. Recent advancements in MLLMs (Liu et al., 2023; Zhu et al., 2023; Li et al., 2023; Dai et al., 2023) have spurred the development of benchmarks to evaluate their scientific problem-solving capabilities. For example, MMMU (Yue et al., 2024a)

and MMMU-Pro (Yue et al., 2024b) present extensive multi-discipline college-level tasks. In materials and chemistry, MaScQA (Zaki et al., 2024) offers a text-based question dataset, while MaCBench (Alampara et al., 2024) proposes a multimodal benchmark. However, MaCBench (Alampara et al., 2024) concentrates on chemistry and general laboratory knowledge, with limited content on crystalline materials. SciFIBench (Roberts et al., 2024) targets scientific figure interpretation using arXiv papers, which do not cover materials science or chemistry extensively (Hsu et al., 2021; Li et al., 2024a) and, being non-peer-reviewed, may have quality concerns. While MMSci (Li et al., 2024b) sources its content from Nature Communications, its task formats-captioning and figurecaption matching-are not representative of the queries posed in actual scientific research, which limits its practical applicability. MatCha, in contrast, is designed with problem formats that authentically mirror the challenges scientists face during the process of scientific discovery.

3 The MatCha benchmark

We introduce MatCha, a comprehensive and challenging benchmark for advancing multimodal materials characterization visual analysis and understanding. MatCha comprises 1,500 questions across 21 tasks, reflecting expert-level difficulty and grounded in real-world scientific scenarios faced by scientists. In the following sections, we elaborate on how we construct tasks (§3.1), collect data (§3.2), generate questions (§3.3), and analyze its composition (§3.4). The holistic construction pipeline is shown in Fig. 1.

3.1 Task Construction

To cover a wide range of domain-specific knowledge for a comprehensive evaluation, we collaborate with experienced researchers in materials science. Drawing from the typical workflow of materials science research—"Processing" \rightarrow "Morphology" \rightarrow "Structure" \rightarrow "Property"—we designed a corresponding chain of stages that reflects this logical progression. The details are as follows:

Processing Correlation (PC). As a foundational step in materials characterization, it involves identifying the characterization technique and its intended purpose. It evaluates the ability to accurately awareness characterization methods and their appropriate application contexts.

Morphology Analysis (MA). This stage focuses on observing and assessing the surface or cross-sectional of materials to gather morphological information. It measures the capability of models to perceive both macro- and micro-scale visual characteristics in electron microscopy images.

Structure Analysis (SA). It targets the interpretation of material structure at the micro- or atomic-scale, which is essential for understanding the underlying mechanisms of various properties, such as mechanical behavior. It further assesses the ability to integrate and link cross-modal knowledge, for instance, linking spectral peaks to chemical bonds or functional groups.

Property Analysis (PA). Since the structure of a material determines its properties, this stage poses a more difficult challenge, by evaluating the logical reasoning capabilities in connecting structural features with properties in physical and chemical.

To concretize these research stages, we define several scientifically meaningful sub-tasks for each, as detailed in App. A. These stages represent the essential steps that materials scientists typically undertake, while the sub-tasks capture a diverse range of real-world challenges encountered during new materials development. Evaluating MLLMs on these tasks is crucial to reveal their potential and limitations in scientific research, providing critical insights into their capability boundaries and guiding future advancement for applications.

3.2 Data Curation

Data Collection. Given the diversity of characterization data, firstly, we define a set of advanced search terms and their synonyms that exhaustively cover the types of data required across the various sub-tasks. With these terms, we employ Exsclaim (Schwenker et al., 2021) to search for and retrieve publicly accessible articles from the Nature platform under CC BY-4.0 license. The search results are sorted by relevance. For each article, we download the HTML file along with all associated figures and their corresponding captions. We crawl 340 articles containing 2,165 figures in total.

Data Processing. First, we query GPT-40 with prompts presented in App. H.1 to segment the full caption of each figure into sub-captions corresponding to sub-figures. Second, Exsclaim (Schwenker et al., 2021) is applied to split each figure into its constituent sub-figures and assign the appropri-

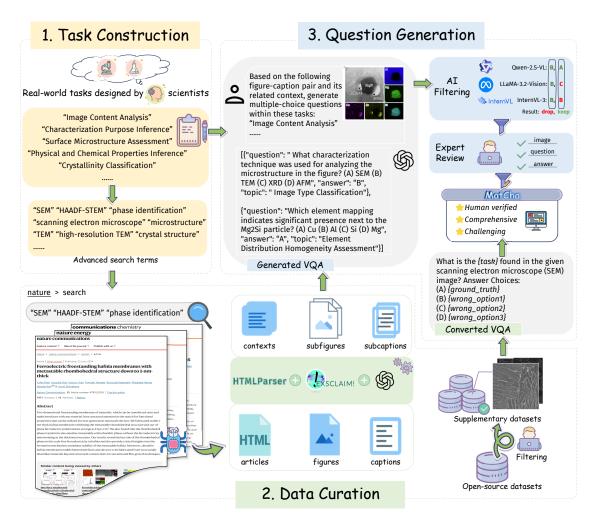


Figure 1: MatCha construction pipeline. First, experts define scientifically meaningful, practical tasks and extract key terms. Second, data is collected and processed using these terms. Third, GPT-based generation and template-based conversion are employed to construct samples from the gathered data, followed by quality filtering and review.

ate sub-captions to each one. The resulting sub-figures are categorized. To ensure the authenticity of the dataset, we retain categories that accurately reflect real-world materials characterization scenarios and exclude other simulated types of sub-figures, such as illustrations. To compensate for insufficient information in the subsequent VQA generation caused by the overly short sub-captions, we use a parser from (Toland et al., 2023) to parse HTML files of the articles and extract the main body text. We then develop a regular expression matching function to retrieve the context relevant to each sub-figure from the main content. Finally, these remaining article contents and images form the basis for the subsequent question generation.

Data Supplementation. Figures in published papers often inevitably contain annotations or markings that may hint at microstructural characteristics. Moreover, although we download high-resolution figures directly, their visual clarity can

still be inferior to images captured directly from characterization instruments. To alleviate this and enhance benchmark diversity and challenge, we source additional non-simulated, human-annotated datasets. Following rigorous expert review and filtering, three datasets (Hecht et al., 2017; Baskaran et al., 2020; Dennler et al., 2021) are selected as supplementary data sources. These datasets consist of high-quality, real-world electron microscopy images, and are used to construct supplementary sub-tasks: surface microstructure analysis (Suppl. SMA), defect type classification (Suppl. DTC), and image content analysis (Suppl. ICA). Specifically, Suppl. SMA requires analyzing the microstructural features of Ti-6Al-4V alloy images; Suppl. DTC examines defect identification and distinguishment between different defect types; and Suppl. ICA assesses the analysis of primary microstructural components in low-carbon steel. These supplementary tasks focus on common microstructural

analyses in practical yet fundamental materials scenarios, evaluating the domain-specific knowledge and visual perception capabilities of MLLMs.

3.3 Question Generation

We formulate our benchmark in a closed-ended VQA format. This approach facilitates easier analysis compared to open-ended VQA and obviates the need for LLMs in automatic evaluation. Consequently, it eliminates subjectivity in answer assessment, ensuring that evaluation results more accurately reflect the true performance of MLLMs.

Generated VQA. Equipped with collected (subfigure, sub-caption, context) triplets, we employ GPT-40 to generate multiple-choice questions. Considering the unique nature of images in each article and to prevent overly image-specific or divergent questions lacking generalizability, we carefully design a prompt that constrains VQA generation within our predefined sub-task scopes. However, we do not restrict the number of answer choices, allowing the model to fully utilize the unique information in each data triplet, as detailed in App. H.2. Each generated VQA sample includes a sub-figure, a question within a sub-task, and a set of answer choices, where one is correct and the others are distractors generated by GPT-40. This process finally yields 26,891 samples, which subsequently undergo data filtering and expert review.

Converted VQA. For the three supplementary datasets, we convert their metadata and labels into multiple-choice questions using a QA template. Similarly, each converted VQA sample consists of an image, a question designed under a certain sub-task, and several answer choices. One choice is the correct answer derived from the original label, while the remaining options are distractors sampled from the label set, excluding the correct one.

Data filtering. For the generated VQA samples, we perform a coarse-grained filtering using AI experts as the first step. Concretely, we use Qwen2.5-VL-7B (Wang et al., 2024a), InternVL3-8B (Chen et al., 2024c), and LLaMA-3.2-11B-Vision (Grattafiori et al., 2024) to each attempt the question three times. If all AI experts answer the question correctly in all attempts, the question is deemed too simple and removed. The remaining questions are retained. This process effectively filters out easy questions and preserves those both challenging and discriminative, ensuring quality.

Experts review. Following coarse filtering by AI experts, two materials science experts-Ph.D. candidates from the material science and engineering department-conduct a manual review to ensure: 1) each figure is a real photograph or data-generated plot, not a simulation or illustration, ensuring data authenticity; 2) each question is grounded in visual content, answerable solely through visual cues and reasoning with intrinsic domain knowledge, without requiring external contextual information, ensuring question validity; 3) overly simple optical character recognition (OCR)-style questions lacking domain-specific expertise are removed to emphasize professionalism and challenge, reflecting realistic scenarios for materials researchers. After this, 994 samples remain, forming the core of MatCha. Additionally, as the three supplementary datasets have already undergone manual validation and annotation, we randomly select 506 samples from their converted VQA set. Together, these comprise a total of 1,500 samples that constitute the complete benchmark. In App. G.1, we provide some examples as illustrations.

3.4 Analysis of MatCha Benchmark

Scientific discovery in materials science demands multiple and complex characterizations beyond basic natural image perception and shallow domain expertise. To investigate the diversity and representativeness of our benchmark, we count and analyze the distribution of sub-tasks, characterization techniques, and material types. A detailed statistical breakdown is provided in App. B. The quantitative statistics show that the dataset encompasses a wide array of methods. The material types include metallic materials, inorganic non-metallic materials, composites, and organic polymers. Furthermore, the collected source articles are retrieved from 14 different journals under the Nature platform, with publication dates ranging from 2015 to 2025, mitigating potential biases towards specific research topics or time periods. This broad scope ensures MatCha reflects a wide range of real-world scientific challenges. The statistics of different subtask samples are shown in Fig. 2 and Fig. 3.

4 Experiments

4.1 Experimental Setup

Evaluation models. We use random selection as a baseline, where a randomly chosen option is treated as the answer. We evaluate a diverse set of

Statistics	Value
MatCha Instances	
Total Images	1,260
Total QA samples	1,500
Average size (px)	559×660
Maximum size (px)	1876×1064
Processing Correlation	
- questions	153
- average length	160
- options	4':153
Morphology Analysis	
- questions	795
- average length	236
options	7':206,5':4
- *	4': 308, 3': 275, 2': 2
- # Among them,	
for the supplementary	
- questions	506
- average length	271
- options	7': 206, 4': 33, 3': 267
Structure Analysis	
- questions	370
- average length	187
options	6':1,5':6
-	4': 345, 3': 15, 2': 3
Property Analysis	
- questions	182
- average length	194
- options	5': 4, 4': 163, 3': 12, 2': 3

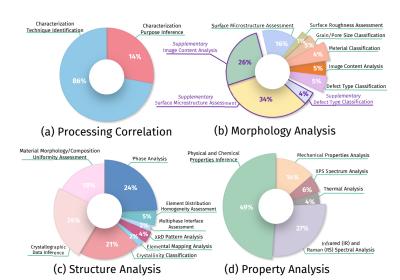


Figure 2: Statistical overview of Figure 3: Composition of the MatCha benchmark, illustrating the proportion MatCha samples. Superscripts and distribution statistics of its 21 sub-tasks across the four progressive research indicate the number of multiple- stages: Processing Correlation (PC), Morphology Analysis (MA), Structure choice options for those questions. Analysis (SA), and Property Analysis (PA).

MLLMs, including the following proprietary models: GPT-{40 (Achiam et al., 2023), 40-mini (Hurst et al., 2024)}, Gemini-1.5-{Pro (Team et al., 2024), Flash (Team et al., 2024)}, Claude-3.5-Sonnet (Anthropic, 2024), Llama-4-Maverick (Meta, 2025). We also evaluate the following popular open-source models: LLaVa-1.5 series (Liu et al., 2024), Qwen-2.5-VL series (Bai et al., 2025), InternVL-3 series (Chen et al., 2024b), Llama-3.2-Vision (Grattafiori et al., 2024), Janus-Pro (Chen et al., 2025), and Gemma-3 (Team et al., 2025). Notably, we choose their chat or instruction-tuned version for each model for better capability of instruction following.

Implementation details. Since all questions are in multiple-choice format, we instruct the models to constrain their responses to the provided option letters. A detailed prompt can be found in App. H.3. For proprietary models, inference is conducted via the API platform. For open-source models, we utilize the Transformers library (Wolf et al., 2020) and the LLaMA-Factory toolkit (Zheng et al., 2024) to perform inference on NVIDIA RTX 6000 GPUs. More details are presented in App. C.

To better understand the performance gap between models and humans, we invite doctoral researchers familiar with various characterization techniques from the material science and engineering department to participate in tests as a human baseline. They are presented with the same format of questions and instructions as the models to

ensure a fair comparison.

4.2 Experimental Results

We adopt a generic zero-shot strategy for evaluation, and the quantitative results of 6 proprietary models and 9 open-source models with varying sizes and architectures are in Tab. 1. Considering the different sources of data, we report the Generated VQA subset and the Converted VQA subset, respectively. The main results are as follows:

MatCha represents a challenging benchmark.

On the generated VQA subset, the best-performing model, GPT-40, achieves only 62.58% accuracy, exhibiting a 26.29% gap from human performance (88.87%). When facing the more difficult converted VQA subset, the top-performing model, LLaMA-4-Maverick, scores only 57.71%, significantly lower than the 88.93% accuracy attained by human experts. These results demonstrate that MatCha provides a well-calibrated level of difficulty and highlight the considerable performance gap between current MLLMs and human experts in materials characterization.

A no-image ablation study is also conducted in App. E to validate that the questions in MatCha are strongly grounded in visual content.

Proprietary models outperform open-source models. Considering the overall tasks in the generated VQA subset, the leading proprietary model outperforms its open-source counterpart by 9.96%

M 11		Gen	erated V	/QA		Converted VQA				MatCha
Model	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All
Baselines										
Random Choice	19.61	26.64	23.24	24.72	15.79	34.83	24.24	15.53	23.91	24.73
Human	90.26	89.31	87.57	88.59	88.87	94.76	87.88	81.55	88.93	88.89
			Propi	ietary M	odels					
GPT-40 (Achiam et al., 2023)	85.62	70.59	71.62	66.48	62.58	60.30	63.64	39.81	52.17	59.07
GPT-4o-mini (Hurst et al., 2024)	66.67	61.59	57.57	51.65	45.27	54.31	48.48	18.93	39.53	43.33
Gemini-1.5-Flash (Team et al., 2024)	75.16	62.63	56.76	54.95	48.39	46.82	48.48	31.07	40.51	45.73
Gemini-1.5-Pro (Team et al., 2024)	86.27	70.24	67.03	62.64	59.76	55.43	66.67	39.81	49.80	56.40
Claude-3.5-Sonnet (Anthropic, 2024)	83.66	71.28	66.49	68.68	61.37	50.56	75.76	46.60	50.59	<u>57.73</u>
LlaMA-4-Maverick (Meta, 2025)	79.74	62.98	60.81	59.89	53.32	69.66	48.48	<u>43.69</u>	57.71	54.80
			Open-	source M	10dels					
Qwen2.5-VL-7B (Bai et al., 2025)	66.01	57.44	54.05	53.85	43.46	45.69	54.55	20.87	36.17	41.00
Qwen2.5-VL-32B (Bai et al., 2025)	69.28	66.44	60.81	66.48	52.62	<u>58.43</u>	54.55	<u>25.24</u>	44.66	49.93
InternVL3-8B (Chen et al., 2024b)	41.83	50.87	48.38	48.90	36.12	46.82	54.55	21.36	36.96	36.40
InternVL3-38B (Chen et al., 2024b)	67.97	66.44	58.38	60.99	49.70	63.30	54.55	23.79	46.64	48.67
LLaVA-1.5-7B (Liu et al., 2024)	22.22	26.99	20.00	29.12	15.90	41.57	24.24	14.56	29.45	20.47
LLaVA-1.5-13B (Liu et al., 2024)	43.14	40.14	38.92	37.36	29.18	33.71	27.27	14.56	25.49	27.93
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	60.13	40.48	40.27	41.21	31.39	38.95	12.12	18.45	28.85	30.53
Janus-Pro-7B (Chen et al., 2025)	48.37	49.13	51.08	53.85	39.54	25.09	48.48	19.90	24.51	34.47
Gemma-3-4b-it (Team et al., 2025)	60.13	47.06	43.24	41.21	34.41	39.33	45.45	25.73	34.19	34.33

Table 1: Evaluation results of model performance on MatCha. Generated VQA: PC (processing correlation), MA (morphology analysis), SA (structure analysis), PA (property analysis). Converted VQA: Suppl. SMA (supplementary surface microstructure analysis), Suppl. ICA (supplementary image content analysis), Suppl. DTC (supplementary defect type classification). **Bolded** values signify the optimal in-class outcomes (open-source or proprietary) and <u>underlined</u> values indicate the suboptimal performance.

(62.58% to 52.62%), indicating a considerable performance gap. This disparity widens to 11.07% (57.71% to 46.64%) in the converted VQA subset. Overall, open-source models explicitly exhibit low performance across both generated and converted VQA subsets, with most models failing to correctly answer more than 40% of the questions.

Performance and generalization capability disparities exist across models. Models show considerable differences in performance across task dimensions, and even within the same model, there are marked fluctuations across different research stages (*e.g.*, PC vs. PA). These results suggest that current models struggle with task generalization and knowledge transfer within the materials science domain.

4.3 Analysis

Most models still struggle with relatively simple perceptual tasks. Some proprietary models, such as Gemini-1.5-Pro, GPT-4o, and Claude-3.5-Sonnet, lead in tasks from the PC and MA stages. Gemini-1.5-Pro, for instance, trails human performance by only 3.99% on the PC stage, indicating its reasonable ability to handle basic characterization techniques. However, even the best-performing model on the MA stage-Claude-3.5-Sonnet-still lags behind human experts by approximately 18.03%. This highlights that, despite recent

advancements in understanding natural images, current cutting-edge MLLMs remain insufficient for achieving expert-level pattern recognition in fine-grained morphological analysis of materials imaging. We attribute this to the current lack of high-quality scientific training corpora, which hinders models from generalizing the ability developed on natural images to these specialized tasks.

Nevertheless, beyond these three models, the performance of most other proprietary and opensource models remains suboptimal. For example, most open-source models failed to exceed 50% accuracy for tasks in the MA stage. This suggests that many current open-source models still struggle to extract and interpret morphological features from electron microscopy images, particularly when recognizing multiscale structures. However, Qwen2.5-VL-32B and InternVL3-38B achieve a notable 66.44% on the MA stage, outperforming certain proprietary models such as GPT-4o-mini (61.91%). This indicates the potential for opensource models to match or even surpass generalpurpose proprietary models in specialized scientific domains when properly optimized.

Most models degrade on tasks requiring more expertise and reasoning ability. Benefiting from our stage-based task design, we can conduct a fine-grained analysis of both proprietary and open-source models across different materials research

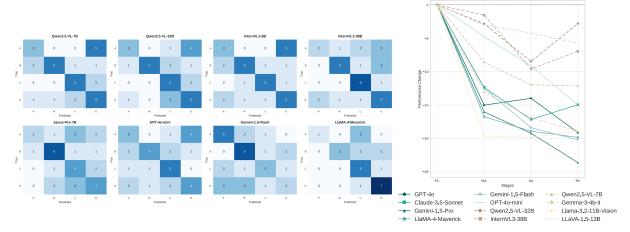


Figure 4: Performance analysis. Left: Confusion matrices of several models on the Suppl. DTC task. Right: Performance degradation trends of models across progressive stages.

stages, as this design aligns well with the research workflow of materials scientists. As shown in the right figure of Fig. 4, we observe a clear trend: as tasks progress from morphology analysis to structure interpretation and properties elucidationstages requiring deeper materials domain knowledge and reasoning ability-most models exhibit a significant decline in performance. Interestingly, this performance deterioration is more pronounced among proprietary MLLMs. On average, proprietary models show a 15.96% accuracy degradation across the MA, SA, and PA stages. In contrast, open-source models demonstrate a smaller average decline of 10.29%. Consequently, several open-source models, including Qwen2.5-VL-7B, Qwen2.5-VL-32B, and InternVL3-38B, have outperformed certain proprietary models in the MA, SA, and PA stages, in some cases approaching the accuracy of the best proprietary models. We hypothesize that this may be related to the incorporation of a larger proportion of scientific imagery or domain-specific textual data during training, particularly in scientific discourse.

These findings suggest that current mainstream MLLMs still lack sufficient materials science domain-specific understanding and multi-level reasoning capabilities, rendering them inadequate for research tasks that demand deep materials expertise and complex logical analysis. This limitation constrains their potential for broader application in scientific research support.

Models fail to recognize microstructural details in more realistic and complex perceptual scenarios. In the converted VQA subset, which involves real-world materials electron microscopy characterization images, current MLLMs encounter substantial challenges despite these tasks being relatively basic for human experts. In the Suppl. SMA task, LLaMA-4-Maverick achieves the best performance with an accuracy of 69.66%, still approximately 25.1% lower than human performance. Open-source models struggle to identify the microstructural hierarchy in alloy images, with most accuracies falling between 30% and 50%. In the Suppl. DTC task, which requires the detection and classification of subtle structural defects, both proprietary and open-source models exhibited significant difficulties. Many failed to determine the presence of defects or distinguish between defects from different origins, indicating a breakdown in their recognition capabilities, as shown in the left figure of Fig. 4. The Suppl. ICA task poses an even greater challenge, demanding fine-grained differentiation of multiple microstructural components in low-carbon steel. This task places high demands on both visual discrimination and metallography knowledge. All open-source models achieve accuracies below 30%, and proprietary models also perform suboptimally, revealing severe limitations in the understanding and reasoning abilities of current models regarding materials microstructure.

Notably, GPT-40 consistently ranked among the top three performers across all three supplementary tasks, suggesting a relatively stronger generalization capability. This may be attributed to its more advanced multimodal alignment mechanisms and broader pretraining data coverage.

4.4 Impact of Few-shot and Chain-of-Thought

To investigate whether the performance limitations observed in the zero-shot setting can be mitigated,

we conduct further experiments using few-shot incontext learning and CoT prompting. The detailed results of performance across various few-shot settings and CoT are presented in App. D.

Few-shot Learning Our few-shot experiments reveal that the ability to leverage in-context examples varies dramatically across models. On the one hand, some models show significant positive scaling. GPT-40, for instance, demonstrates a remarkable improvement on the challenging Converted VQA subset, with its accuracy jumping from 52.17% (zero-shot) to 73.52% (16-shot). Other models like Claude-3.5-Sonnet and GPT-4o-mini also show a clear upward trend. This suggests that for some models, in-context examples can effectively elicit domain-specific knowledge and improve pattern recognition. On the other hand, this benefit is not universal and can even be detrimental. For Gemini-1.5-Pro, performance is inconsistent and can even degrade at higher shot counts. The effect is more severe for Qwen2.5-VL-32B, which suffers a notable performance drop when provided with examples. These cases indicate that for certain models, in-context examples may introduce noise or conflict with their internal knowledge, making the effective utilization of such prompts a non-trivial challenge.

Chain-of-Thought Prompting Our evaluation of CoT prompting reveals mixed results. On one hand, most models show slight improvement with step-by-step reasoning. Janus-Pro, for instance, sees the most significant benefit, with its accuracy increasing by nearly 6 percentage points to 40.27%. On the other hand, for top-performing models like Gemini-1.5-Pro and Llama-4-Maverick, CoT prompting actually hindered performance compared to the direct zero-shot approach.

These contradictory outcomes suggest that CoT is not a universally effective strategy. While it can help some models, it may conflict with the internal reasoning pathways of others. This reinforces our finding that simple prompting techniques, such as in-context learning and CoT, are insufficient to overcome fundamental gaps in intrinsic domain knowledge and visual perception. Therefore, other techniques, such as the retrieval mentioned in App. F, remain for future exploration.

4.5 Error Analysis

To gain deeper insights into the nature of model failures, we sample 100 incorrect predictions from

Table 2: Distribution of error types for three models.

Model	Visual Perception Error	Lack of Material Knowledge	Language & Logic Failure	Image-Text Alignment Error
GPT-40 (Achiam et al., 2023)	27%	71%	1%	1%
GPT-4o-mini (Hurst et al., 2024)	34%	59%	5%	2%
Gemini-1.5-Pro (Team et al., 2024)	32%	64%	2%	2%

GPT-40 and categorize the primary cause of error using GPT-40. We predefine four main error types: 1) Visual Perception Error, where the model fails to correctly identify visual features, such as structures, boundaries, or grains; 2) Lack of Material Knowledge, where the model misunderstands or lacks necessary materials domain concepts; 3) Language and Logic Failure, involving errors in explanation or the semantic understanding of the question and options; and 4) Image-Text Alignment Misunderstanding, where the model incorrectly links image content to the question or options text.

As shown in Tab. 2, the predominant error category across all models is *Lack of Material Knowledge*, accounting for over 60-70% of failures, confirming that a core limitation of current MLLMs is their deficiency in specialized scientific domain knowledge. *Visual Perception Error* is the second most common failure type, highlighting that even leading models struggle with the fine-grained and complex visual patterns in scientific imagery, a conclusion that aligns with our previous analysis. We provide error cases for illustration in App. G.2.

5 Conclusion

Automated interpretation of complex materials characterization imaging data remains a significant bottleneck. In this work, we introduce MatCha, the first benchmark designed to assess the capability of current MLLMs to understand materials characterization imagery. We construct a suite of realistic and scientifically meaningful tasks to evaluate state-of-the-art MLLMs. Our comprehensive evaluations pinpoint their inadequacy for generalization, deep domain knowledge, complex morphology perception, and nuanced materials analysis. MatCha thus serves as a critical tool for diagnosing these core deficiencies, aiming to guide the development of MLLMs that can truly accelerate materials research and enable autonomous scientific discovery.

Limitations

Although MatCha aims for broad coverage with 21 sub-tasks, the vast field of materials science means it cannot encompass every material, characteriza-

tion technique, experimental nuance, or specific research question. Certain emerging areas or highly specialized analyses might be underrepresented. Meanwhile, the dataset of 1,500 questions, while reviewed by experts, might inadvertently contain biases or not fully capture the statistical distribution of all real-world scenarios.

Potential Risks

While MatCha is designed to rigorously assess the specialized knowledge and visual understanding capabilities of MLLMs within the domain of materials characterization, a potential risk or limitation lies in its direct applicability to all professional scenarios without further consideration. High performance on MatCha indicates strong foundational capabilities, but evaluating MLLMs in specific operational laboratory or industrial settings may necessitate adaptation or fine-tuning to those unique professional contexts. Real-world applications often involve distinct instrument variations, proprietary data intricacies, evolving experimental procedures, or specific analytical goals not exhaustively covered by any standardized benchmark. Therefore, there is a risk that strong MatCha scores might be overgeneralized, and deploying these MLLMs effectively into diverse, practical workflows will likely still require careful contextual adjustments and validation to ensure optimal and reliable performance in those specialized environments.

Acknowledgments

work was supported by This Shenthe zhen Science and Technology Program (JCYJ20220818103001002), Shenzhen Doctoral Startup Funding (RCBS20221008093330065), Tianyuan Fund for Mathematics of National Natural Science Foundation of China (NSFC) (12326608), Shenzhen Science and Technology Program (Shenzhen Key Laboratory Grant No. ZDSYS20230626091302006), and Shenzhen Stability Science Program 2023, Shenzhen Key Lab of Multi-Modal Cognitive Computing. The authors thank the financial support from the National Natural Science Foundation of China (Grant no. 22302174) and the Natural Science Foundation of Zhejiang Province (Grant no. LZ25E030005). Dr. Hong-Qing Liang acknowledges gratefully the research startup package from Zhejiang University.

References

- Ossama B Abouelatta. 2013. Classification of copper alloys microstructure using image processing and neural network. *Journal of American Science*, 9(6):213–223.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Nawaf Alampara, Indrajeet Mandal, Pranav Khetarpal, Hargun Singh Grover, Mara Schilling-Wilhelmi, NM Anoop Krishnan, and Kevin Maik Jablonka. 2024. Macbench: a multimodal chemistry and materials science benchmark. In *Proceedings of the 38th Conference on Neural Information Processing Systems (NeurIPS 2024)*.
- Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.
- Luis M Antunes, Keith T Butler, and Ricardo Grau-Crespo. 2024. Crystal structure generation with autoregressive large language modeling. *Nature Communications*, 15(1):10570.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report.
- Jonas Bals and Matthias Epple. 2023. Deep learning for automated size and shape analysis of nanoparticles in scanning electron microscopy. *RSC advances*, 13(5):2795–2802.
- Arun Baskaran, Genevieve Kane, Krista Biggs, Robert Hull, and Daniel Lewis. 2020. Adaptive characterization of microstructure dataset using a two stage machine learning approach. *Computational Materials Science*, 177:109593.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and 1 others. 2024a. Are we on the right way for evaluating large vision-language models? *Advances in Neural Information Processing Systems*, 37:27056–27087.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. 2025. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv* preprint arXiv:2501.17811.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024b. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.

- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024c. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198.
- Aritra Chowdhury, Elizabeth Kautz, Bülent Yener, and Daniel Lewis. 2016. Image driven machine learning methods for microstructure recognition. *Computational Materials Science*, 123:176–187.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. Advances in neural information processing systems, 36:49250–49267.
- Brian L DeCost, Toby Francis, and Elizabeth A Holm. 2017. Exploring the microstructure manifold: image texture representations applied to ultrahigh carbon steel microstructures. *Acta Materialia*, 133:30–40.
- Brian L DeCost and Elizabeth A Holm. 2015. A computer vision approach for automated analysis and classification of microstructural image data. *Computational materials science*, 110:126–133.
- Brian L DeCost, Bo Lei, Toby Francis, and Elizabeth A Holm. 2019. High throughput quantitative metallography for complex microstructures using deep learning: A case study in ultrahigh carbon steel. *Microscopy and Microanalysis*, 25(1):21–29.
- Nik Dennler, Antonio Foncubierta-Rodriguez, Titus Neupert, and Marilyne Sousa. 2021. Learning-based defect recognition for quasi-periodic hrstem images. *Micron.* 146:103069.
- Qianggang Ding, Santiago Miret, and Bang Liu. 2024. Matexpert: Decomposing materials discovery by mimicking human experts. *arXiv preprint arXiv:2410.21317*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Xuehai He, Yichen Zhang, Luntian Mou, Eric Xing, and Pengtao Xie. 2020. Pathvqa: 30000+ questions for medical visual question answering. *arXiv* preprint *arXiv*:2003.10286.
- Matthew D Hecht, Brian L DeCost, Toby Francis, Elizabeth A Holm, Yoosuf N Picard, Bryan A Webler, and 1 others. 2017. Ultrahigh carbon steel micrographs.

- Elizabeth A Holm, Ryan Cohn, Nan Gao, Andrew R Kitahara, Thomas P Matson, Bo Lei, and Srujana Rao Yarasi. 2020. Overview: Computer vision and machine learning for microstructural characterization and analysis. *Metallurgical and Materials Transactions A*, 51(12):5985–5999.
- Ting-Yao Hsu, C Lee Giles, and Ting-Hao'Kenneth' Huang. 2021. Scicap: Generating captions for scientific figures. *arXiv* preprint arXiv:2110.11624.
- Zhi Huang, Federico Bianchi, Mert Yuksekgonul, Thomas J Montine, and James Zou. 2023. A visual–language foundation model for pathology image analysis using medical twitter. *Nature medicine*, 29(9):2307–2316.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Beverley J Inkson. 2016. Scanning electron microscopy (sem) and transmission electron microscopy (tem) for materials characterization. In *Materials characterization using nondestructive evaluation (NDE) methods*, pages 17–43. Elsevier.
- Jakub Lála, Odhran O'Donoghue, Aleksandar Shtedritski, Sam Cox, Samuel G Rodriques, and Andrew D White. 2023. Paperqa: Retrieval-augmented generative agent for scientific research. arXiv preprint arXiv:2312.07559.
- Andreas Leitherer, Byung Chul Yeo, Christian H Liebscher, and Luca M Ghiringhelli. 2023. Automatic identification of crystal structures and interfaces via artificial-intelligence-based electron microscopy. *npj Computational Materials*, 9(1):179.
- Yang Leng. 2013. Materials characterization: introduction to microscopic and spectroscopic methods. John Wiley & Sons.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024a. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language models. *arXiv* preprint arXiv:2403.00231.
- Wei Li, Kevin G Field, and Dane Morgan. 2018. Automated defect analysis in electron microscopic images. *npj Computational Materials*, 4(1):36.
- Yifan Li, Zhixin Lai, Wentao Bao, Zhen Tan, Anh Dao, Kewei Sui, Jiayi Shen, Dong Liu, Huan Liu, and Yu Kong. 2025. Visual large language models for generalized and specialized applications. *arXiv* preprint arXiv:2501.02765.

- Zekun Li, Xianjun Yang, Kyuri Choi, Wanrong Zhu, Ryan Hsieh, HyeonJung Kim, Jin Hyuk Lim, Sungyoung Ji, Byungju Lee, Xifeng Yan, and 1 others. 2024b. Mmsci: A dataset for graduate-level multidiscipline multimodal scientific understanding.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. Advances in neural information processing systems, 36:34892– 34916.
- Alejandro Lozano, Jeffrey Nirschl, James Burgess, Sanket Rajan Gupte, Yuhui Zhang, Alyssa Unell, and Serena Yeung. 2024. Micro-bench: A microscopy benchmark for vision-language understanding. *Advances in Neural Information Processing Systems*, 37:30670–30685.
- Jacob Madsen, Pei Liu, Jens Kling, Jakob Birkedal Wagner, Thomas Willum Hansen, Ole Winther, and Jakob Schiøtz. 2018. A deep learning approach to identify local structures in atomic-resolution transmission electron microscopy images. Advanced Theory and Simulations, 1(8):1800037.
- Artem Maksov, Ondrej Dyck, Kai Wang, Kai Xiao, David B Geohegan, Bobby G Sumpter, Rama K Vasudevan, Stephen Jesse, Sergei V Kalinin, and Maxim Ziatdinov. 2019. Deep learning analysis of defect and phase evolution during electron beaminduced transformations in ws2. *npj Computational Materials*, 5(1):12.
- AI Meta. 2025. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation.
- Vaibhav Mishra, Somaditya Singh, Dhruv Ahlawat, Mohd Zaki, Vaibhav Bihani, Hargun Singh Grover, Biswajit Mishra, Santiago Miret, NM Krishnan, and 1 others. 2024. Foundational large language models for materials research. arXiv preprint arXiv:2412.09560.
- Graham Roberts, Simon Y Haile, Rajat Sainju, Danny J Edwards, Brian Hutchinson, and Yuanyuan Zhu. 2019. Deep learning for semantic segmentation of defects in advanced stem images of steels. *Scientific reports*, 9(1):12744.
- Jonathan Roberts, Kai Han, Neil Houlsby, and Samuel Albanie. 2024. Scifibench: Benchmarking large multimodal models for scientific figure interpretation. *Advances in Neural Information Processing Systems*, 37:18695–18728.
- Ian M Robertson, Christopher A Schuh, John S Vetrano, Nigel D Browning, David P Field, Dorte Juul Jensen, Michael K Miller, Ian Baker, David C Dunand, Rafal Dunin-Borkowski, and 1 others. 2011. Towards an integrated materials characterization toolbox. *Journal of Materials Research*, 26(11):1341–1383.

- Andre Niyongabo Rubungo, Craig Arnold, Barry P Rand, and Adji Bousso Dieng. 2023. Llm-prop: Predicting physical and electronic properties of crystalline solids from their text descriptions. *arXiv* preprint arXiv:2310.14029.
- Eric Schwenker, Weixin Jiang, Trevor Spreadbury, Nicola Ferrier, Oliver Cossairt, and Maria KY Chan. 2021. Exsclaim!—an automated pipeline for the construction of labeled materials imaging datasets from literature. *arXiv preprint arXiv:2103.10631*.
- Mingren Shen, Guanzhao Li, Dongxia Wu, Yuhan Liu, Jacob RC Greaves, Wei Hao, Nathaniel J Krakauer, Leah Krudy, Jacob Perez, Varun Sreenivasan, and 1 others. 2021. Multi defect detection and analysis of electron microscopy images with deep learning. *Computational Materials Science*, 199:110576.
- Yingheng Tang, Wenbin Xu, Jie Cao, Weilu Gao, Steve Farrell, Benjamin Erichson, Michael W Mahoney, Andy Nonaka, and Zhi Yao. 2025. Matterchat: A multi-modal llm for material science. *arXiv preprint arXiv:2502.13107*.
- 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.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, and 1 others. 2025. Gemma 3 technical report. arXiv preprint arXiv:2503.19786.
- Aubrey Toland, Huan Tran, Lihua Chen, Yinghao Li, Chao Zhang, Will Gutekunst, and Rampi Ramprasad. 2023. Accelerated scheme to predict ring-opening polymerization enthalpy: simulation-experimental data fusion and multitask machine learning. *The Journal of Physical Chemistry A*, 127(50):10709–10716.
- Prateek Verma, Minh-Hao Van, and Xintao Wu. 2024. Beyond human vision: The role of large vision language models in microscope image analysis. In 2024 IEEE International Conference on Big Data (Big-Data), pages 1700–1705. IEEE.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, and 1 others. 2024a. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Song XiXuan, and 1 others. 2024b. Cogvlm: Visual expert for pretrained language models. *Advances in Neural Information Processing Systems*, 37:121475–121499.

Małgorzata Warmuzek, Marcin Żelawski, and Tomasz Jałocha. 2021. Application of the convolutional neural network for recognition of the metal alloys microstructure constituents based on their morphological characteristics. *Computational Materials Science*, 199:110722.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, and 1 others. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.

Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, Wenjie Zhang, and 1 others. 2023. Darwin series: Domain specific large language models for natural science. *arXiv preprint arXiv:2308.13565*.

Tong Xie, Yuwei Wan, Yixuan Liu, Yuchen Zeng, Shaozhou Wang, Wenjie Zhang, Clara Grazian, Chunyu Kit, Wanli Ouyang, Dongzhan Zhou, and 1 others. 2024. Darwin 1.5: Large language models as materials science adapted learners. *arXiv preprint arXiv:2412.11970*.

Sang-Hyeok Yang, Wooseon Choi, Byeong Wook Cho, Frederick Osei-Tutu Agyapong-Fordjour, Sehwan Park, Seok Joon Yun, Hyung-Jin Kim, Young-Kyu Han, Young Hee Lee, Ki Kang Kim, and 1 others. 2021. Deep learning-assisted quantification of atomic dopants and defects in 2d materials. *Advanced Science*, 8(16):2101099.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.

Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, and 1 others. 2024b. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*.

Mohd Zaki, NM Anoop Krishnan, and 1 others. 2024. Mascqa: investigating materials science knowledge of large language models. *Digital Discovery*, 3(2):313–327.

Alexander N Zaloga, Vladimir V Stanovov, Oksana E Bezrukova, Petr S Dubinin, and Igor S Yakimov. 2020. Crystal symmetry classification from powder x-ray diffraction patterns using a convolutional neural network. *Materials Today Communications*, 25:101662.

Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. 2024. Honeycomb: A flexible llm-based agent system for materials science. *arXiv preprint arXiv:2409.00135*.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient finetuning of 100+ language models. *arXiv preprint arXiv:2403.13372*.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

A Sub-task Details

We present each stage along with its corresponding sub-tasks and explanation in Tab. 3, all defined by materials scientists. These sub-tasks are designed to comprehensively encompass common challenges and characterization techniques encountered in materials characterization.

B Dataset Statistics and Diversity

Based on our statistics, the papers used in MatCha are sourced from the following journals: Communications Earth & Environment, Communications Materials, Light: Science & Applications, NPG Asia Materials, Nature Biotechnology, Nature Communications, Nature Materials, Nature Photonics, Nature Synthesis, Polymer Journal, Scientific Data, Scientific Reports, npj Computational Materials, and npj Heritage Science. The distribution statistics below show that the dataset encompasses a wide range of mainstream characterization techniques and material types. These types are annotated by instructing GPT-40 and are then classified and merged by human experts accordingly. Previous datasets in the materials characterization domain are largely limited to techniques such as SEM and XRD, while MatCha integrates data across the spectrum of mainstream methods, making it the most diverse dataset of its kind to date.

B.1 Distribution of Characterization Types

The data in MatCha cover a wide range of characterization techniques, categorized as follows:

• **Microscopy**: Includes transmission electron microscopy (147), scanning electron microscopy (550), scanning transmission electron microscopy (83), optical microscopy (42), X-ray imaging/tomography (1), scanning

Stage	Sub-task	Explanation
Dun and distribution	Characterization Technique Identification	Determine the characterization technique used (e.g., SEM, TEM, XRD, AFM).
Processing Correlation	Characterization Purpose Inference	Deduce the scientific purpose of using a particular characterization technique.
	Material Classification	Infer the general category of the material (e.g., metal, ceramic, polymer, composite).
	Image Content Analysis	Analyze microstructural content from electron microscope images
		to extract material characteristics.
	Surface Microstructure Assessment	Determine structural features such as the presence of surface defects,
Morphology Analysis		the order of the crystal structure, and the existence of layered structures, etc.
Worphology Analysis	Surface Roughness Assessment	Evaluate whether the material surface appears smooth or rough,
		or determine the level of surface roughness.
	Defect Type Classification	Recognize and classify defect types,
	such as dislocations, vacancies, stacking faults, grain boundaries, etc.	
	Grain/Pore Size Classification	Categorize the size scale of grains or pores (e.g., nanometer, micrometer, or millimeter range).
	Crystallographic Data Inference	Utilize unit cell parameters and lattice spacings to determine structural features,
		such as symmetry, space group, and overall lattice architecture.
	Crystallinity Classification	Assess whether the material is amorphous, polycrystalline, or single crystalline
		based on the image.
	Multiphase Interface Assessment	Examine the presence of multiple phases or interfaces within the image,
		and analyze their structural and compositional features.
	X-ray diffraction (XRD) Pattern Analysis	Extract and analyze key information from XRD spectra,
		including peak positions and other characteristic features.
Structure Analysis	Phase Analysis	Include phase identification and classification, interpretation of phase composition, assess-
		ment of phase homogeneity, determination of crystal structure and polymorphic forms, etc.
	Elemental Mapping Analysis	Identify elements represented by different colors or regions in elemental mapping images ,
		(e.g., EDS or EELS maps).
	Element Distribution Homogeneity Assessment	Analyze the image to assess whether elements are uniformly distributed across the material.
	Material Morphology/Composition Uniformity Assessment	Assess the uniformity of material morphology and composition, such as
		component ratios, particle size distribution, the homogeneity of internal microstructures, etc.
	Physical and Chemical Properties Inference	Predict potential physical or chemical properties of materials.
	Mechanical Properties Analysis	Extract key parameters such as yield strength, ultimate tensile strength,
		and ductility from the stress-strain curve to assess material performance.
Property Analysis	Thermal Analysis	Extract critical information from various thermal analysis methods (e.g., TGA, DTA, DSC).
1 topolty Allalysis	Infrared (IR) and Raman (RS) Spectral Analysis	Elucidate the molecular structure and chemical composition, identify functional groups,
		and bond types associated with specific spectral peaks, etc.
	X-ray Photoelectron Spectroscopy (XPS) Spectrum Analysis	Analyze XPS spectra to identify peak positions, determine elemental composition,
		and chemical states, and infer chemical structures.

Table 3: MatCha taxonomy of sub-tasks.

probe microscopy (11), atom probe tomography (5), focused ion beam-SEM (6), and X-ray photoemission electron microscopy (2).

- **Spectroscopy**: Includes Raman spectroscopy (42), photoluminescence/fluorescence spectroscopy (22), infrared spectroscopy (23), energy-dispersive X-ray spectroscopy (21), Xray photoelectron spectroscopy (9), nuclear magnetic resonance (10), time-resolved spectroscopy (8), X-ray absorption spectroscopy (6), ultraviolet-visible spectroscopy (5), mass spectrometry (3), extended X-ray absorption fine structure (3), electroluminescence (3), electron probe micro-Analysis (2), Fouriertransform spectroscopy (1), atomic absorption spectroscopy (1), Auger electron spectroscopy (1), electron energy loss spectroscopy (1), electron paramagnetic resonance (1), and cathodoluminescence (1).
- Diffraction and Scattering: Includes electron diffraction (68), X-ray diffraction (58), electron backscatter diffraction (18), smallangle X-ray scattering (8), transmission Kikuchi diffraction (3), grazing-incidence X-ray diffraction (2), neutron diffraction (2), and reciprocal space mapping (1).
- Electrochemical Analysis: Includes general

- electrochemical tests (15), cyclic voltammetry (9), electrochemical impedance spectroscopy (3), voltammetry (2), cycling stability test (2), performance evaluation (Faradaic/Coulombic efficiency, 2), galvanostatic intermittent titration technique (1), and galvanostatic charge-discharge (1).
- Computation and Simulation: Includes microscopy/diffraction simulation (7), general simulation/computation (7), quantum chemistry calculation (4), first-principles calculation (2), optical simulation (1), phase diagram/thermodynamic calculation (1), finite element method (1), and molecular dynamics (1).
- **Mechanical Testing**: Includes stress-strain test (8), general mechanical testing (5), fracture toughness (1), scratch test (1), and nanoindentation (1).
- Thermal Analysis: Includes differential scanning calorimetry (4), thermal imaging (1), and thermal conductivity measurement (1).
- Magnetic Characterization: Includes magnetic measurement (3), and X-ray magnetic circular dichroism (2).

• Other and Performance Evaluation: Includes data/image analysis (65), electrical/optoelectronic device performance (18), elemental/compositional mapping (10), and physical property measurement (9).

B.2 Distribution of Material Types

The benchmark covers four major categories of materials: metallic materials (683), inorganic non-metallic materials (435), composite materials (167), and organic polymer materials (98).

C Parameters Settings

To reduce randomness, the temperature is fixed at 0 for models using API interface. For models executed with the Transformers library, the default setting is retained. Furthermore, in zero-shot and few-shot settings where models are required to directly output multiple-choice answers, the maximum generation length (max_new_tokens) is set to 32. In contrast, for the CoT experiments, the generation length is set to 8192.

D Detailed Few-shot and Chain-of-Thought Results

We evaluate model performance under both few-shot and CoT prompting settings. For the few-shot experiments, we randomly sample 2, 4, 8, and 16 in-context examples for each test instance from its corresponding data source. For CoT prompting, we append the phrase "Let's think step by step" to the end of the original instruction. The results are shown in Tab. 4, Tab. 5, Tab. 6, Tab. 7, and Tab. 8, respectively.

In the 8- and 16-shot settings, several models, including LLaVA-1.5 (7B, 13B) and Janus-Pro-7B, fail to output a valid option, likely due to challenges in processing the extended context.

E Ablations Experiments

Followed by previous research (Goyal et al., 2017; Chen et al., 2024a), we conduct a no-image ablation study to investigate the contribution of visual information, shown in Tab. 9.

The notably low score of LLaMA-4-Maverick in the no-image condition is primarily due to its tendency to bypass the question and directly generate answers without adequately adhering to instructions in the absence of visual input. Other models also exhibited substantial performance degradation—over 20%—with accuracy only marginally (around

10%) above random guessing. This underscores the significant role of visual information in MatCha and highlights the rigorous demands on the visual understanding capabilities of models.

Model	MatCha All	Ablation no-image drop
GPT-40 (Achiam et al., 2023)	34.33	-24.74
Gemini-1.5-Pro (Team et al., 2024)	29.53	-26.87
Claude-3.5-Sonnet (Anthropic, 2024)	34.60	-23.40
LlaMA-4-Maverick (Meta, 2025)	5.93	-20.47

Table 9: No-image ablation study on MatCha.

F Future Direction: The Potential Role of Retrieval-Augmented Generation

A primary conclusion from our benchmark is that MLLMs are significantly constrained by a lack of specialized domain knowledge, a common challenge in vertical domains. A promising solution is Retrieval-Augmented Generation (RAG), which allows models to access external knowledge bases. RAG has already shown success in scientific domains, such as PaperQA (Lála et al., 2023), which answers complex questions by retrieving information from scientific articles.

Our few-shot experiments indicate that providing in-context examples can, to some extent, elicit domain knowledge of a model. Similarly, RAG is another powerful paradigm providing dynamic, external knowledge to a model before it generates a response. Given that both methods function by supplying context, we believe RAG can also further improve performance in materials characterization scenarios. This makes MatCha an ideal testbed for evaluating future multimodal RAG systems in materials science.

		Gen	erated V	/QA		Converted VQA				MatCha
Model	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All
Proprietary Models										
GPT-40 (Achiam et al., 2023)	82.35	71.97	71.35	67.58	62.47	87.64	60.61	46.12	68.97	64.67
GPT-4o-mini (Hurst et al., 2024)	73.86	66.09	58.65	53.30	48.59	52.06	63.64	13.59	37.15	44.73
Gemini-1.5-Flash (Team et al., 2024)	76.47	63.67	58.11	58.24	51.31	63.67	60.61	33.01	50.99	51.20
Gemini-1.5-Pro (Team et al., 2024)	80.39	66.78	58.92	58.24	52.72	59.18	60.61	50.97	55.93	53.80
Claude-3.5-Sonnet (Anthropic, 2024)	74.51	70.59	67.03	68.68	59.66	26.59	78.79	57.77	42.69	53.93
LlaMA-4-Maverick (Meta, 2025)	87.58	75.78	70.81	67.58	64.59	83.52	66.67	56.80	71.54	66.93
			Open-	source M	10dels					
Qwen2.5-VL-7B (Bai et al., 2025)	64.71	58.82	53.24	59.34	43.86	38.58	57.58	24.76	34.19	40.60
Qwen2.5-VL-32B (Bai et al., 2025)	64.71	65.05	60.81	57.69	48.29	39.70	51.52	24.27	34.19	43.53
InternVL3-8B (Chen et al., 2024b)	52.29	55.71	51.89	50.00	39.74	65.92	57.58	35.92	53.16	44.27
InternVL3-38B (Chen et al., 2024b)	66.67	68.51	57.57	57.14	50.50	74.91	63.64	29.13	55.53	52.20
LLaVA-1.5-7B (Liu et al., 2024)	26.14	35.99	34.32	31.32	23.94	32.58	33.33	12.62	24.51	24.13
LLaVA-1.5-13B (Liu et al., 2024)	32.68	48.44	42.70	43.41	33.00	32.21	18.18	12.62	23.32	29.73
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	47.06	48.10	40.00	45.05	30.99	38.95	12.12	9.71	25.30	29.07
Janus-Pro-7B (Chen et al., 2025)	50.33	53.98	57.57	55.49	42.66	32.21	18.18	15.05	24.31	36.47
Gemma-3-4b-it (Team et al., 2025)	54.90	49.13	42.70	43.41	34.41	30.34	57.58	21.84	28.66	32.47

Table 4: 2-shot results of model performance on MatCha. **Bolded** values signify the optimal in-class outcomes (open-source or proprietary).

		Gen	erated V	/QA		Converted VQA				MatCha	
Model	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All	
Proprietary Models											
GPT-40 (Achiam et al., 2023)	84.31	73.01	74.05	65.93	63.88	85.02	57.58	53.88	70.55	66.13	
GPT-4o-mini (Hurst et al., 2024)	77.78	66.44	60.27	55.49	51.11	47.57	54.55	16.50	35.38	45.80	
Gemini-1.5-Flash (Team et al., 2024)	78.43	62.63	57.84	56.04	50.30	35.21	57.58	33.01	35.77	45.40	
Gemini-1.5-Pro (Team et al., 2024)	83.01	70.59	65.41	63.19	58.15	59.55	66.67	47.09	54.94	57.07	
Claude-3.5-Sonnet (Anthropic, 2024)	87.58	74.39	69.46	69.23	64.59	81.65	75.76	68.93	76.09	68.47	
LlaMA-4-Maverick (Meta, 2025)	88.24	71.28	72.70	67.03	63.98	85.39	66.67	50.00	69.76	65.93	
			Open-	source M	10dels						
Qwen2.5-VL-7B (Bai et al., 2025)	62.75	62.63	51.35	58.79	43.76	46.82	51.52	20.87	36.56	41.33	
Qwen2.5-VL-32B (Bai et al., 2025)	63.40	67.13	59.46	57.14	48.79	48.69	48.48	23.30	38.34	45.27	
InternVL3-8B (Chen et al., 2024b)	50.33	56.75	51.89	49.45	41.35	50.19	51.52	33.98	43.68	42.13	
InternVL3-38B (Chen et al., 2024b)	70.59	69.90	58.92	58.79	51.41	60.67	57.58	40.78	52.37	51.73	
LLaVA-1.5-7B (Liu et al., 2024)	21.57	23.88	31.62	28.02	17.81	32.21	21.21	13.59	23.91	19.87	
LLaVA-1.5-13B (Liu et al., 2024)	32.03	24.22	30.54	25.27	18.91	32.21	24.24	13.59	24.11	20.67	
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	47.06	54.67	44.32	45.60	35.21	23.97	18.18	16.99	20.75	30.33	
Janus-Pro-7B (Chen et al., 2025)	52.29	54.33	52.43	51.10	39.64	33.33	27.27	15.05	25.49	34.87	
Gemma-3-4b-it (Team et al., 2025)	50.98	47.75	44.32	43.96	32.80	33.33	42.42	24.76	30.43	32.00	

Table 5: 4-shot results of model performance on MatCha. **Bolded** values signify the optimal in-class outcomes (open-source or proprietary).

M-1-1		Gen	erated V	/QA		Converted VQA				MatCha	
Model	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All	
Proprietary Models											
GPT-4o (Achiam et al., 2023)	84.97	73.36	72.70	68.13	64.19	86.52	60.61	46.12	68.38	65.60	
GPT-4o-mini (Hurst et al., 2024)	77.12	63.32	59.73	57.14	49.70	63.67	54.55	18.93	44.86	48.07	
Gemini-1.5-Flash (Team et al., 2024)	81.05	64.71	57.03	57.14	50.70	52.81	33.33	43.20	47.63	49.67	
Gemini-1.5-Pro (Team et al., 2024)	84.31	69.20	60.00	66.48	56.94	64.42	66.67	43.69	56.13	56.67	
Claude-3.5-Sonnet (Anthropic, 2024)	85.62	75.43	67.84	69.23	63.68	89.89	84.85	62.14	78.26	68.60	
LlaMA-4-Maverick (Meta, 2025)	84.31	73.36	71.89	66.48	63.38	90.26	60.61	52.43	72.92	66.60	
			Open-	source M	10dels						
Qwen2.5-VL-7B (Bai et al., 2025)	71.24	61.94	53.78	57.14	45.47	32.96	48.48	20.87	29.05	39.93	
Qwen2.5-VL-32B (Bai et al., 2025)	62.09	67.47	60.27	60.44	50.00	38.58	57.58	20.39	32.41	44.07	
InternVL3-8B (Chen et al., 2024b)	56.86	57.44	52.16	50.55	42.05	60.30	48.48	35.92	49.60	44.60	
InternVL3-38B (Chen et al., 2024b)	73.20	67.47	59.46	55.49	51.51	58.43	54.55	33.98	48.22	50.40	
LLaVA-1.5-7B (Liu et al., 2024)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
LLaVA-1.5-13B (Liu et al., 2024)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	51.63	50.87	45.14	42.31	32.49	36.70	9.09	16.02	26.48	30.47	
Janus-Pro-7B (Chen et al., 2025)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
Gemma-3-4b-it (Team et al., 2025)	58.17	49.48	47.03	46.70	37.02	36.33	39.39	23.79	31.42	35.13	

Table 6: 8-shot results of model performance on MatCha. **Bolded** values signify the optimal in-class outcomes (open-source or proprietary).

Model	D.C.		erated V	•	4.11	Converted VQA				MatCha	
	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All	
Proprietary Models											
GPT-40 (Achiam et al., 2023)	83.66	73.36	73.51	73.08	65.19	89.14	66.67	54.37	73.52	68.00	
GPT-4o-mini (Hurst et al., 2024)	78.43	68.17	60.81	57.14	52.62	64.79	60.61	16.99	45.06	50.07	
Gemini-1.5-Flash (Team et al., 2024)	77.12	63.67	57.84	58.24	51.21	61.80	57.58	44.66	54.55	52.33	
Gemini-1.5-Pro (Team et al., 2024)	84.97	65.40	60.54	61.54	55.53	59.18	66.67	34.47	49.60	53.53	
Claude-3.5-Sonnet (Anthropic, 2024)	83.01	73.70	67.57	68.13	62.07	75.28	75.76	41.26	61.46	61.87	
LlaMA-4-Maverick (Meta, 2025)	84.97	75.09	71.08	68.68	64.29	89.89	42.42	62.62	75.69	68.13	
			Open-	source M	10dels						
Qwen2.5-VL-7B (Bai et al., 2025)	70.59	62.63	54.05	58.24	46.08	33.33	48.48	10.19	24.90	38.93	
Qwen2.5-VL-32B (Bai et al., 2025)	64.05	66.09	62.97	60.44	50.50	40.07	51.52	24.76	34.58	45.13	
InternVL3-8B (Chen et al., 2024b)	62.09	61.94	52.97	53.85	45.37	47.94	45.45	33.01	41.70	44.13	
InternVL3-38B (Chen et al., 2024b)	75.16	69.55	61.08	59.34	53.22	49.81	54.55	36.41	44.66	50.33	
LLaVA-1.5-7B (Liu et al., 2024)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
LLaVA-1.5-13B (Liu et al., 2024)	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	57.52	58.13	45.68	40.11	37.12	36.33	18.18	10.68	24.70	32.93	
Janus-Pro-7B (Chen et al., 2025)	55.56	55.36	61.35	57.14	46.38	31.09	45.45	21.36	28.06	40.20	
Gemma-3-4b-it (Team et al., 2025)	64.71	50.87	49.46	42.86	38.53	28.09	39.39	24.76	27.47	34.80	

Table 7: 16-shot results of model performance on MatCha. **Bolded** values signify the optimal in-class outcomes (open-source or proprietary).

M-4-1		Gen	erated V	/QA		Converted VQA				MatCha
Model	PC	MA	SA	PA	All	Suppl. SMA	Suppl. DTC	Suppl. ICA	All	All
GPT-4o (Achiam et al., 2023)	84.97	73.36	68.65	70.88	63.08	55.43	69.70	46.12	52.57	59.53
GPT-4o-mini (Hurst et al., 2024)	71.90	57.09	61.35	58.24	48.39	53.93	54.55	24.76	42.09	46.27
Gemini-1.5-Flash (Team et al., 2024)	78.43	61.94	56.49	59.89	49.70	47.19	42.42	30.58	40.12	46.47
Gemini-1.5-Pro (Team et al., 2024)	83.66	69.55	60.81	64.29	55.63	47.19	72.73	42.23	46.84	52.67
Claude-3.5-Sonnet (Anthropic, 2024)	84.31	71.97	68.11	69.78	61.77	49.06	81.82	51.46	52.17	58.53
LlaMA-4-Maverick (Meta, 2025)	60.78	54.33	51.62	54.95	40.85	55.06	54.55	55.34	55.14	45.67
			Open-	source M	10dels					
Qwen2.5-VL-7B (Bai et al., 2025)	62.75	56.40	50.81	57.69	44.57	49.06	60.61	22.82	39.13	42.73
Qwen2.5-VL-32B (Bai et al., 2025)	67.32	67.82	61.08	64.29	52.52	53.18	60.61	32.52	45.26	50.07
InternVL3-8B (Chen et al., 2024b)	45.10	54.67	52.16	56.04	40.44	53.93	51.52	29.61	43.87	41.60
InternVL3-38B (Chen et al., 2024b)	60.78	65.40	60.27	61.54	50.30	59.18	57.58	41.75	51.98	50.87
LLaVA-1.5-7B (Liu et al., 2024)	5.23	15.57	15.95	22.53	9.26	32.58	24.24	10.68	23.12	13.93
LLaVA-1.5-13B (Liu et al., 2024)	5.88	16.96	16.49	21.98	9.46	35.58	24.24	15.05	26.48	15.20
Llama-3.2-11B-Vision (Grattafiori et al., 2024)	60.13	45.67	47.03	46.70	37.32	35.96	24.24	22.82	29.84	34.80
Janus-Pro-7B (Chen et al., 2025)	55.56	55.71	61.35	57.69	46.48	31.09	45.45	21.36	28.06	40.27
Gemma-3-4b-it (Team et al., 2025)	58.17	49.13	49.46	45.05	36.82	28.84	57.58	24.27	28.85	34.13

Table 8: CoT results of model performance on MatCha. **Bolded** values signify the optimal in-class outcomes (open-source or proprietary).

G Cases

G.1 VQA Cases

Physical and Chemical Properties Inference

Question: What does the red circle in the 230 °C frame indicate regarding the nanorods' crystallization?

Choices: (A) The maximum diffraction intensity (B) Onset of the first diffraction spot

(C) Completion of crystallization (D) Absence of any crystallization

Answer: B

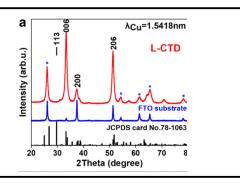


Crystallographic Data Inference

Question: What is the space group designation associated with the phase identified in the L-CTD XRD pattern?

Choices: (A) Fm-3m (B) Pbcn (C) P63/mmc

(D) Ia3d **Answer**: B



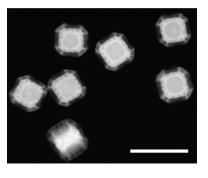
Surface Microstructure Assessment

Question: What structural feature is clearly visible in the HAADF-STEM image of type 2 p-ANHs?

Choices: (A) Cylindrical canyons on the 100 facets (B) Entirely smooth surface (C) Spherical nanoparticle morphology (D) Randomly

Answer: A

structured pores

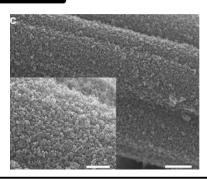


Material Morphology and Composition Uniformity Assessment

Question: What type of material morphology is observed in the SEM image?

Choices: (A) Rod-like structures (B) Uniformly distributed spherical nanoparticles (C) Thin film coating (D) Porous sponge-like structures

Answer: B



Mechanical Properties Analysis

Question: What is the approximate value of the yield stress observed in the engineering stress-strain curve?

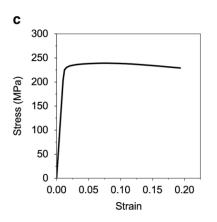
Choices: (A) 230 MPa (B) 300 MPa (C) 200

MPa (D) 50 MPa Answer: A

Question: What is the approximate strain at which stress begins to decrease in the engineering stress-strain curve?

Choices: (A) 0.08 (B) 0.20 (C) 0.01 (D) 0.50

Answer: A

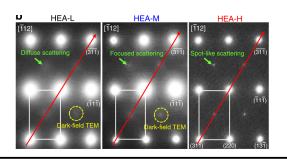


Characterization Technique Identification

Question: What characterization technique was used to capture the diffraction patterns shown in Figure 5b?

Choices: (A) SEM (B) XRD (C) TEM (D)

STEM **Answer**: C



G.2 Error Cases

Lack of Material Knowledge

Question: How many distinct lattice structures are identified from the FFT pattern inset? (A) One (B) Two (C) Three (D) Four

Correct Answer: (C) Model Answer: (B)

Model Prediction Rationable:

To determine how many distinct lattice structures are identified from the FFT (Fast Fourier Transform) pattern inset, let's analyze the image step by step:

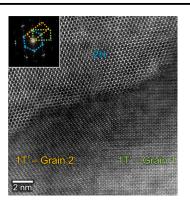
1. Observation of the FFT Pattern:

The inset FFT pattern shows three distinct sets of diffraction spots, each marked with different colored dashed lines (blue, orange, and green). These correspond to different lattice structures or orientations.

- 2. Identification of Lattice Structures in the Image:
- The main image shows three regions labeled as "2H," "1T' Grain 1," and "1T' Grain 2."
- The "2H" region corresponds to one lattice structure.
- The "1T' Grain 1" and "1T' Grain 2" regions correspond to the same lattice structure (1T') but with different grain orientations.
- Conclusion:

From the FFT pattern, we can identify two distinct lattice structures:

- The "2H" lattice structure.



- The "1T'" lattice structure (with two grains but the same structure). Thus, the correct answer is: (B) Two

Error Analysis:

The model incorrectly concluded that there are only two distinct lattice structures based on the FFT pattern and the labeled regions in the image. The mistake lies in the interpretation of the FFT pattern inset, which clearly shows three distinct sets of diffraction spots marked by different colored dashed lines (blue, orange, and green). These correspond to three distinct lattice structures or orientations. The model failed to recognize that the FFT pattern indicates three distinct lattice structures, not just two, despite the labeling of '1T' grains in the main image. This error stems from a lack of material knowledge and misunderstanding of how FFT patterns correspond to lattice structures.

Visual Perception Error

Question: What are the main microstructure in the given scanning electron microscope (SEM) images of Ti-6Al-4V alloy? (A) acicular: acicular morphology (B) lamellar: α lamellae (C) duplex: duplex microstructure, a bimodal morphology consisting of equiaxed grains and α -lamellae

Correct Answer: (C) Model Answer: (B)

Model Prediction Rationable:

To identify the main microstructure in the given SEM image of the Ti-6Al-4V alloy, let's analyze the features step by step:

- 1. Observation of the Image:
- The image shows a clear pattern of elongated, parallel structures within larger grains.
- These structures resemble lamellae, which are characteristic of a lamellar microstructure.
- 2. Comparison with Descriptions:
- Acicular: Acicular morphology typically consists of needle-like structures, which are not observed in this image.
- Lamellar: The image shows α lamellae within the grains, matching the description of a lamellar microstructure.
- Duplex: A duplex microstructure would show a combination of equiaxed grains and α lamellae. However, the image does not show distinct equiaxed grains alongside lamellae.
- 3. Conclusion: The microstructure in the SEM image corresponds to (B) lamellar: α lamellae. **Error Analysis**:

The model incorrectly identified the microstructure as lamellar (B) instead of duplex (C). The SEM image clearly shows a bimodal morphology consisting of equiaxed grains and α -lamellae, which is characteristic of a duplex microstructure. The model failed to recognize the presence of equiaxed grains in the image, focusing only on the lamellae structures. This indicates a visual perception error, as the model did not correctly identify the equiaxed grains in the image.

H Prompt Setting

H.1 Sub-caption Segmentation Prompt

```
System prompt:

Subfigure labels are letters referring to individual subfigures within a larger figure. Please separate the given full caption into the exact subcaptions and format as a syntactically valid JSON format with keys the letter of each subcaption. If there is no full caption then return an empty JSON.

User prompt:

Caption: {caption}
```

H.2 QA Generation Prompt

```
System prompt:
You are a scientific expert. Based on the following figure-caption pair and its related context
from the article in the field of materials science and characterization, generate 1 to 4 visual
question answering (VQA)-style question-answer pairs, depending on the amount of information
provided. The question should be related to both the figure-caption pair and the context, but
the answer should be able to be inferred and analyzed ONLY from the figure.
{sub-tasks with explanation}
You should generate multiple-choice questions, ensure that the answers are concise and clear.
For multiple-choice questions, please generate plausible but incorrect options. The number
of options is not limited, and enclose all options in parentheses (e.g., (A)) as part of the
{\it question.} \ {\it After providing the question and answer, also include the topic of this question, and}
output in JSON format: { "vqas": [ "question": ..., "answer": ..., "topic": ... ] } Example template: { "vqas": [ "question": ... (A) xx (B) xx (C) xx (D) xx, "answer": "B", "topic": ...
] }
User prompt:
Figure {sub-figure}
Caption: {sub-caption}
Related Context: {related context}
```

H.3 Evaluation Prompt

```
System prompt:

You are a helpful materials science assistant. Based on the figure, please answer the following question. The answer could be inferred from the figure and must be concise and clear. Answer directly without any explanation.

User prompt:

Question: {question with options}

Answer with the option's letter from the given choices directly:
```