STRUCTURED AND ABSTRACTIVE REASONING ON MULTI-MODAL RELATIONAL KNOWLEDGE IMAGES

Anonymous authors

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

031

033

034

037 038

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Understanding and reasoning with abstractive information from the visual modality presents significant challenges for current multi-modal large language models (MLLMs). Among the various forms of abstractive information, Multi-Modal Relational Knowledge (MMRK), which represents abstract relational structures between multi-modal entities using node-edge formats, remains largely underexplored. In particular, STructured and Abstractive Reasoning (STAR) on such data has received little attention from the research community. To bridge the dual gaps in large-scale high-quality data and capability enhancement methodologies, this paper makes the following key contributions: (i). An automatic STAR data engine capable of synthesizing images with MMRK to build multi-modal instruction data with reliable chain-of-thought thinking for various STAR tasks and (ii). A comprehsive two-stage capability enhancement training framework, accompanied by a suite of evaluation protocols tailored to different STAR tasks. Based upon these contributions, we introduce STAR-64K, a dataset comprising 64K high-quality multi-modal instruction samples, and conduct experiments across 5 open-source MLLMs. Experimental results show that our two-stage enhancement framework enables smaller 3B/7B models to significantly outperform GPT-4o in STAR. Additionally, we provide in-depth analysis regarding the effectiveness of various designs, data transferability, and scalability.

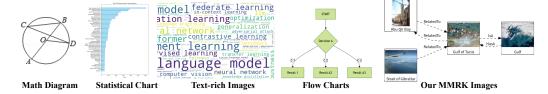


Figure 1: Different kinds of images contain abstractive information with complex semantics.

1 Introduction

Multi-modal large language models (MLLMs) (Song et al., 2023) achieve state-of-the-art understanding and reasoning capabilities across various multi-modal tasks, and are increasingly adopted in fields such as automatic driving (Cui et al., 2024), health care (Liu et al., 2023a), agriculture (Zhu et al., 2024), etc. Capability enhancement and evaluation of MLLMs are highly active areas, with a focus on advancing holistic model capabilities and charting the limits of the MLLMs.

While much of the existing research is concentrated on understanding and reasoning about real-world objects and scenes depicted within images (Fu et al., 2024), other studies have begun to explore models' abilities to interpret and reason over images that convey highly abstract semantic information. As illustrated in Figure 1, these abstractive semantic elements are highly diverse, including charts (Masry et al., 2022), mathematical diagrams (Lu et al., 2024), and more. Such abstractive semantic information is frequently presented at a conceptual level through artificially constructed visual patterns, which are defined by humans and are absent in nature. Effectively reasoning about abstractive image inputs poses an elevated challenge for MLLMs, as it demands not only basic object recognition but also a deeper understanding and interpretation of the complex information encapsulated within these human-defined abstractive visual forms.

Among the diverse array of abstractive images, an important area remains underexplored: <u>ST</u>ructured and <u>A</u>bstractive <u>Reasoning</u> (STAR) on images with <u>Multi-Modal Relational Knowledge</u> (MMRK). As illustrated in Figure 1, MMRK consists of multiple multi-modal entities and concepts that are interconnected by abstract relational edges, representing well-organized and structured factual knowledge. Unlike natural or other abstractive images, MMRK offers a flexible and structured format for encoding complex semantic relations, with broad application potential (An et al., 2025). The relational links act as higher-order human-defined abstractions, modeling intricate connections among entities, and thus place greater demands on MLLM's reasoning capabilities. To accurately perform STAR, MLLMs must understand both the entities and the underlying relational structure. However, STAR remains largely unaddressed, with only a few studies (Zhang et al., 2024a; 2025d) briefly investigating this capability, which still face two critical challenges:

- (i) Lack of large-scale data synthesis method for STAR. From the data perspective, there is a shortage of high-quality MMRK images and corresponding multi-modal instruction data. Automated pipelines for generating diverse and scalable MMRK datasets are missing, along with reliable chain-of-thought (CoT) reasoning annotations needed to improve MLLM's complex thinking and generalization ability.
- (ii) **Absence of effective enhancement and evaluation frameworks for STAR.** From a methodology perspective, a systematic training and evaluation framework for STAR is lacking. Existing work (Zhang et al., 2025d) only addresses zero-shot evaluation. Fine-tuning MLLMs on large-scale synthetic data is necessary to effectively enhance their STAR capabilities.

To tackle these challenges, we develop an automatic STAR data synthesis engine that first generates images containing MMRK and then produces instruction data accompanied by fine-grained CoT reasoning. Given the current limitations of MLLMs, our approach leverages multi-modal knowledge graphs (MMKGs) as the data source, which are structured repositories of reliable multi-modal information. We further introduce a variety of MMRK-related tasks during the synthesis process. In addition, we propose a two-stage training framework, combining supervised fine-tuning and preference alignment, to enhance MLLMs' STAR capabilities and introduce a specialized evaluation protocol. Our contribution in this paper can be summarized as follows:

- Automatic STAR Data Engine. We introduce the data synthesis engine, which examines MLLM capabilities from a novel perspective called structured and abstractive reasoning (STAR) using MMRK images. Our engine automatically generates high-quality instruction data using large-scale MMKGs with rich relational knowledge, eliminating costly manual annotation. By visualizing sampled multi-modal subgraphs and creating diverse seed STAR tasks, each data instance includes an MMRK image, a task-specific question, and a detailed CoT answer.
- Comprehensive Training and Evaluation Pipeline. We propose a systematic pipeline for enhancing and evaluating STAR capabilities in MLLMs. Our two-stage training combines instruction tuning for general competency and preference alignment for targeted optimization, utilizing the data synthesized by our engine. We also establish a dedicated protocol for STAR evaluation.
- In-depth Experimental Exploration. We conduct extensive experiment exploration on 5 famous open-source MLLM backbones from 3B to 34B, aiming to identify key factors influencing STAR enhancement. Our results demonstrate that targeted training can substantially improve MLLMs' STAR abilities with abstractive visual information, uncovering the mechanisms that enable accurate reasoning in complex multi-modal semantic contexts. Smaller MLLMs with 3B/7B parameters can outperform mainstream product-level MLLMs like GPT-4o.

2 Related Works

Multi-modal Large Language Model (MLLM) Enhancement and Evaluation MLLMs (Song et al., 2023) extend LLMs with multi-modal understanding and reasoning capabilities by incorporating multi-modal information into LLMs with different connectors (Zhu et al., 2025), which supports diverse modalities such as images (Zhu et al., 2023), audio (Huang et al., 2024), videos (Zhang et al., 2025a). A lot of datasets and benchmarks are developed to enhance and evaluate the specific capabilities of MLLMs, including multi-disciplinary knowledge (Lu et al., 2022; Yue et al., 2024b), fine-grained recognition and perception (Tong et al., 2024), structured chart understanding

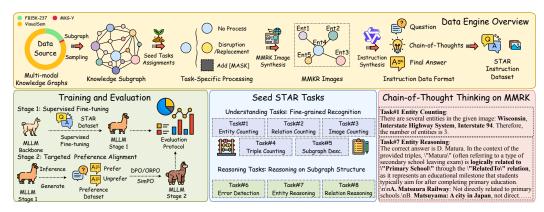


Figure 2: The overview of our data engine, the training pipeline, the seed tasks, and the CoT prompts.

(Masry et al., 2022; Wang et al., 2024), visual mathematical reasoning (Lu et al., 2024; Zhang et al., 2025b), etc. Automatic data synthesis pipelines (Zhang et al., 2024a) are usually designed for these benchmarks to generate high-quality multi-modal instruction data from complex data sources.

Multi-modal knowledge graphs (MMKGs) (Chen et al., 2024) consists of structured triple knowledge with multi-modal contents such as entity images (Liu et al., 2019), text descriptions (Yao et al., 2019), and knowledge-grounded audio/videos (Pan et al., 2022). These multi-modal contents enhance traditional triple-based KG with rich semantic information to serve diverse application scenarios by providing multi-modal factual knowledge. In the age of LLM, the combination of LLM and MMKG (Zhang et al., 2024b; Wan & Yu, 2025) attracts widespread attention from both academia and industry, which focuses on leveraging the high-quality multi-modal knowledge to reduce LLM's hallucination (Zhang et al., 2023). This work provides a new perspective to incorporate the abstractive reasoning ability of MLLMs by synthesizing multi-modal instruction data with MMKGs to provide reliable MMRK.

3 THE STAR DATA ENGINE

In this section, we introduce our STAR data engine, designed to synthesize images paired with MMRK and corresponding text instructions for constructing high-quality multi-modal instruction datasets tailored to STAR tasks. With this engine, we generate STAR-64K with diverse task types.

3.1 Data Engine Overview

Figure 2 presents the overview of our data engine and synthesis pipeline. With an input subgraph sampled from MMKG with MMRK, the engine generates multi-modal instruction data, comprising pairs of MMRK and text instructions. Note that MMRK is a multi-modal knowledge subgraph, which is the visual modality input for MLLMs. In this work, we define 8 different seed tasks, focusing on MMRK for both structured and abstractive relational reasoning. These seed STAR tasks designed by us are divided into two categories:

Understanding the MMRK Data. Before MLLMs can perform complex reasoning, they must accurately identify and describe various components within MMRK. Task types in this category include **Entity Counting** (EC, Task #1), **Relation Counting** (RC, Task #2), **Image Counting** (IC, Task #3), **Triple Counting** (TC, Task #4), and **Subgraph Description** (SD, Task #5). EC, RC, IC, and TC, respectively, require the model to count entities, relations, images, and triples present in MMRK, demonstrating its understanding of fundamental elements. SD asks the MLLM to briefly describe the given visual MMRK, requiring it to grasp the global context of the data. These understanding tasks are inspired by classic visual recognition and understanding benchmarks; however, in the context of MMRK images, recognizing and describing such complex semantic networks in the visual modality becomes significantly more challenging.

Reasoning on the MMRK Data. Upon their understanding of MMRK data, MLLMs are expected to integrate information from MMRK with their own knowledge to perform advanced reasoning and predictions, which would be an advanced capability for MLLMs. Therefore, we propose **Error**

Detection (ED, Task #6), **Entity Reasoning** (ER, Task #7), and **Relation Reasoning** (RR, Task #8). ED requires MLLMs to detect the anomalous entity that causes a factual error in the given MMRK while ER and RR ask MLLMs to make a choice for a certain missing entity/relation in the given MMRK. The design of these tasks is motivated by classic reasoning tasks like knowledge graph completion (Zhang et al., 2025c) on KGs, and we hope that MLLMs can demonstrate similar reasoning abilities on the visual modality information containing MMRK.

_ _

3.2 Detailed Synthesis Pipeline

Based on the eight seed tasks discussed previously, we devise a five-step pipeline to synthesize MMRK images and text prompts, thereby constructing multi-modal instruction data for STAR tasks.

- **Step 1. Data Source.** We select three public MMKGs as our source data: VisualSem (Alberts et al., 2020), FB15K-237 (Liu et al., 2019), and MKG-Y (Xu et al., 2022), which contain million-scale encyclopedic common-sense knowledge with images and entity descriptions as multi-modal contents. The statistical information of the three MMKGs is presented in Appendix A.1. We can denote one MMKG as $\mathcal{KG} = (\mathcal{E}, \mathcal{R}, \mathcal{T}, \mathcal{I}, \mathcal{D})$, where $\mathcal{E}, \mathcal{R}, \mathcal{T}$ represent the entity, relation, and triple sets, respectively. \mathcal{I}, \mathcal{D} are the image and text description sets for entities in the MMKG.
- **Step 2. Subgraph Sampling.** Next, we sample knowledge subgraphs $\mathcal{KG}' \subseteq \mathcal{KG}$ where its entity/relation/triple sets are the subsets of the full MMKG. For each sample, an entity is selected as the starting point, followed by a random walk that combines depth-first and breadth-first search until a specified number of entities and relations are collected. To control data complexity, we limit each subgraph to a maximum of 9 entities. After sampling the subgraphs, we also sample images from \mathcal{I} for each entity in \mathcal{KG}' for visualization.
- Step 3. Task-specific Processing. Subgraph instances are then assigned to each seed task for further task-specific processing. EC, RC, TC, and SD require no additional modification. For IC, we randomly remove some of the images of entities from \mathcal{I}' , introducing missing image information and differentiating it from EC. In ED, a single entity is randomly replaced with another from the global entity set to introduce an error. For ER and RR, a target entity or relation is masked in the subgraph by replacing its text and image with a [MASK] mark, challenging the model to infer the missing information. After processing, we obtain a modified subgraph \mathcal{KG}'' for each subgraph \mathcal{KG}' .
- Step 4. MMRK Image Data Synthesis. We visualize each processed subgraph, converting it into the image modality using a KG visualization tool such as GraphViz (Ellson et al., 2004), formally represented as: $\mathcal{V} \leftarrow \texttt{GraphViz}(\mathcal{KG''},\mathcal{I'},\mathcal{D'})$ where $\mathcal{I'},\mathcal{D'}$ are the image set and text description set of $\mathcal{KG'}$ respectively. Finally, the MMRK image \mathcal{V} integrates all entities in $\mathcal{KG''}$ with their images and textual descriptions, connected via directional relations. These relational edges provide structured, abstractive information that links entities across multiple modalities, making full understanding and reasoning over such structured abstractions a significant challenge for MLLMs.
- **Step 5. Instruction Data Synthesis.** We then synthesize the corresponding instruction data for the synthesized MMRK images. For each \mathcal{V} , we prepare an input question \mathcal{Q} and answer \mathcal{A} to form one instruction instance $(\mathcal{V}, \mathcal{Q}, \mathcal{A})$. Distinct question and answer templates are crafted for each seed task type. For a given task, \mathcal{Q} is generated using a fixed template, while the answer is determined by the specific context of \mathcal{V} . Meanwhile, the answer \mathcal{A} is divided into two segments: a chain-of-thought (CoT) thinking process and the final answer. Different tasks are associated with specific question and answer templates, as well as distinct methods for generating chain-of-thought (CoT) reasoning. We present the technical details in Appendix A.2 and the instruction templates in Appendix B.
- **Information of the Synthetic Data.** With the mentioned pipeline in our data engine, we finally generate 8000 data instances for each task and split them into train/valid/test sets with 8:1:1. Therefore, the full training set consists of 51200 instances, while the validation and test sets consist of 6400 instances of data, respectively. Therefore, we name the 64K data synthesised by us as STAR-64K. Next, we would present the training and evaluation pipeline.

4 Training and Evaluation Protocol

In this section, we detail our training and evaluation protocol designed to enhance the STAR capability of MLLMs using the large-scale benchmark synthesized by our data engine. Given that

existing MLLMs lack robust STAR abilities, we first employ a two-stage fine-tuning strategy to imbue the models with these capabilities, and then comprehensively assess their performance using our evaluation pipeline.

4.1 TWO-STAGE TRAINING FOR CAPABILITY ENHANCEMENT

To strengthen the STAR capability of MLLMs, we propose a two-stage training pipeline. In Stage 1, we perform supervised fine-tuning (SFT) for general capability enhancement; in Stage 2, we apply preference alignment (PA) methods to target specific optimization of failure cases.

Stage 1: Supervised Fine-tuning for General Enhancement. We first fine-tune MLLMs with visual instruction data $\mathcal{D}_{sft} = \{(\mathcal{V}_i, \mathcal{Q}_i, \mathcal{A}_i)\}_{i=1}^{N_1}$ synthesised by our data engine. By doing this, MLLMs can learn the basic STAR ability and the basic output format (structured CoT reasoning plus final answer) by fitting on the training data.

Stage 2: Preference Alignment for Targeted Optimization. Upon completion of Stage 1, the MLLMs demonstrate baseline competency for simple understanding and reasoning over structured information in visual knowledge graphs. However, we find that a single round of SFT is insufficient to fully unlock the model's potential, especially in complex or error-prone scenarios where hallucinations persist. To address this, we introduce PA methods in Stage 2 for targeted performance improvement on these challenging cases.

The data format for stage 2 can be denoted as $\mathcal{D}_{pa} = \{(\mathcal{V}_i, \mathcal{Q}_i, \mathcal{A}_i^{(p)}, \mathcal{A}_i^{(p)})\}_{i=1}^{N_2}$ where $\mathcal{A}_i^{(p)}, \mathcal{A}_i^{(p)}$ represent a preferred answer and an unpreferred answer for current input $\mathcal{V}_i, \mathcal{Q}_i$. Specifically, we run inference on the training data after Stage 1 and select instances where the model fails to generate correct outputs. For these instances, the gold answers are treated as preferred, and the incorrect, model-generated answers as unpreferred, thus forming the PA dataset \mathcal{D}_{pa} . We then adopt PA methods such as DPO (Rafailov et al., 2023) to further optimize the MLLMs, explicitly improving performance on hard cases. By maximizing the likelihood of preferred answers and minimizing that of unpreferred ones, the loss function refines the model's output distribution, thereby boosting accuracy on challenging data. More details of the two-stage pipeline are provided in Appendix A.3.

4.2 EVALUATION PROTOCOL

Following the two-stage training, we evaluate MLLM performance on all eight STAR tasks using the protocol below. For all tasks except Task #5, we define ground-truth answers: counting tasks require a numerical value as the answer, while detection/reasoning tasks require a response or selecting a specific entity or relation within the image with MMRK. Accuracy is calculated by comparing predictions against the standard answers.

To assess the quality and correctness of model-generated CoT reasoning, we follow the LLM-as-a-Judge paradigm (Li et al., 2024), leveraging a stronger LLM as an evaluator that scores the generated CoT relative to our gold labels. As Task 5 is open-ended, we assess subgraph descriptions using the same scoring approach as for CoT evaluation. This comprehensive protocol enables holistic assessment of MLLM STAR capabilities using a diverse set of metrics.

5 EXPERIMENTS AND ANALYSIS

In this section, we introduce the detailed experiment settings and present our experiment results and further analysis, focusing on the following **research questions** (**RQ**):

- **RQ1**: Does the two-stage training pipeline enhance the STAR capabilities of MLLM?
- **RQ2**: Can the different tasks we design influence each other with positive or negative transfer?
- **RQ3**: What scale of data is required to incorporate STAR capabilities?
- **RQ4**: Does the CoT process have a positive effect on the final performance?
- **RQ5**: How much do the entity images and texts influence the final model performance?
- **RQ6**: Are there any intuitive cases to show the performance of MLLMs after two-stage training?

Table 1: The main experiment results on two-stage training on 5 open-source MLLMs. For stage 1(S1), we conduct two groups of experiments S1(single) and S1(Full), representing SFT on single task/full data. For stage 2(S2), we employ three classic PA methods including DPO/OPRO/SimPO.

Experiment Settings		Tas	k#1	Task#2 Task#3		k#3	Task#4		Task#5			Task#7		Task#8		AVG		
		ACC	CoT	ACC	CoT	ACC	CoT	ACC	CoT	Score	ACC	CoT	ACC	CoT	ACC	CoT	AVG	
QVQ-72B Qwen2.5-72B GPT-4v GPT-4o-mini GPT-4o		30.75 38.25 37.75 67.50 43.75	- - - -	8.25 65.38 41.25 72.25 56.33	- - - -	5.50 8.63 14.00 29.88 17.38	- - - -	42.50 39.13 40.00 31.25 34.50	- - - -	50.13 65.00 59.25 69.13 82.38	4.63 0.25 3.63 3.50 2.73	- - - -	25.38 40.38 29.83 29.25 53.88	- - - -	16.00 52.88 39.13 23.00 40.00	- - - -	22.89 38.74 33.11 40.72 41.37	
Qwen2.5-VL	3В	Zero-shot S1(Single) S1(Full) S2(DPO) S2(ORPO) S2(SimPO)	18.25 51.00 42.75 55.50 39.00 71.00	62.07 52.67 73.28 64.79 79.03	20.13 56.63 67.00 89.25 84.62 89.38	74.75 79.74 95.67 94.06 96.82	3.50 10.38 57.13 66.88 59.00 37.25	28.38 29.17 65.37 60.29 52.31	12.75 20.13 23.50 26.13 17.88 28.13	31.90 31.87 51.77 43.81 53.23	57.71 58.31 59.94 66.64 66.81 67.92	6.25 20.37 37.25 37.50 37.75 37.62	32.34 32.00 38.42 39.00 40.06	47.63 52.75 61.13 60.00 59.63 59.13	53.76 55.88 60.94 61.21 60.89	38.25 64.50 77.25 76.85 77.63 78.25	54.00 56.00 64.77 65.79 66.78	25.56 41.76 53.24 59.84 55.29 58.59
	7B	Zero-shot S1(Single) S1(Full) S2(DPO) S2(ORPO) S2(SimPO)	6.13 77.13 64.88 66.50 65.75 69.63	82.47 81.79 82.11 81.96 82.96	12.25 91.13 92.75 94.00 93.38 93.75	95.85 97.38 97.79 98.89 97.76	0.13 65.88 71.37 73.50 71.88 75.38	71.10 76.70 77.32 77.09 78.55	13.13 24.75 27.62 30.25 27.38 29.00	54.96 54.07 58.67 54.45 56.77	68.62 74.85 75.71 76.44 76.65 76.32	0.75 52.75 55.87 58.63 56.38 57.75	42.93 45.23 46.10 47.11 46.52	26.00 64.38 67.50 69.37 70.00 68.50	65.84 69.40 68.79 68.96 68.13	42.88 77.63 80.13 82.00 79.75 81.50	68.24 71.52 72.35 71.20 71.95	21.24 66.06 66.98 68.84 67.65 68.98
	32B	Zero-shot S1(Single) S1(Full) S2(DPO) S2(ORPO) S2(SimPO)	55.50 77.25 67.75 73.25 60.38 76.37	81.29 79.84 82.12 78.65 82.01	70.25 88.00 93.63 94.50 93.75 89.25	83.07 99.70 97.73 97.39 92.25	14.88 57.75 63.13 65.75 60.75 66.63	64.35 70.21 72.46 69.98 73.23	1.63 23.13 27.50 32.50 25.25 32.75	46.87 53.93 62.61 55.44 60.29	72.08 69.98 75.07 72.98 73.36 74.16	5.38 41.50 54.00 52.38 52.75 54.38	41.05 44.16 44.67 46.07 45.44	37.50 64.63 73.50 71.25 69.63 71.38	65.11 68.90 69.95 68.92 69.78	40.38 77.00 81.75 81.63 80.62 82.25	65.24 70.59 70.75 69.35 71.59	37.20 62.41 67.04 68.03 64.56 68.40
LLaVA-1.5/NEXT	7B	Zero-shot S1(Single) S1(Full) S2(DPO) S2(ORPO) S2(SimPO)	11.75 34.13 70.38 71.25 71.25 67.38	28.91 39.02 49.41 49.72 48.86	2.13 43.25 66.75 86.13 86.13	66.18 83.89 92.02 92.27 92.07	12.88 23.38 60.75 62.63 63.25 63.50	28.99 34.50 46.40 46.01 46.75	4.38 11.25 22.25 44.75 45.50 39.00	22.27 27.30 41.84 42.22 41.69	20.27 25.77 33.24 42.34 42.87 42.41	1.13 2.25 6.50 19.37 19.75 20.13	18.68 22.93 29.23 28.99 28.97	36.50 33.63 54.50 60.63 60.00 59.50	37.75 41.76 59.50 59.45 57.65	59.88 61.35 79.62 80.25 79.88 80.13	41.32 48.61 58.59 59.78 58.46	18.62 29.38 49.25 58.42 58.58 57.27
	34B	Zero-shot S1(Single) S1(Full) S2(DPO) S2(ORPO) S2(SimPO)	12.25 57.13 97.50 97.88 97.88 98.25	75.12 86.05 93.05 92.47 92.96	7.63 81.50 96.75 99.25 99.25 99.25	90.20 98.04 99.62 99.43 99.29	8.25 68.38 85.00 98.88 98.75 98.75	52.77 80.79 89.98 89.71 90.09	19.63 84.00 68.13 66.87 70.13 62.25	40.78 64.35 83.30 82.89 82.31	59.31 65.23 72.02 79.43 79.86 79.42	0.63 58.50 66.63 66.38 66.87 67.00	31.33 38.67 46.94 46.98 45.46	45.88 52.38 75.25 74.75 76.00 73.88	56.75 60.90 72.32 73.00 71.86	54.88 67.63 85.00 84.62 85.25 85.87	58.91 66.39 75.79 75.55 76.20	26.06 59.03 80.79 83.51 84.25 83.08

5.1 EXPERIMENT SETTINGS

Baselines. In our experiments, we utilize Qwen2.5-VL-3B/7B/32B (Bai et al., 2025), LLaVA-1.5-7B (Liu et al., 2023b), and LLaVA-NEXT-34B (Liu et al., 2024) as MLLM backbones. For each backbone, we report three groups of results: (1) zero-shot performance without any finetuning, (2) results after stage 1 training (Vanilla SFT), and (3) results after both stage 1 and stage 2 training (SFT + PA). In the SFT stage, we experiment with two settings: training on single-task data and on the full STAR-64K dataset. For PA, we employ three mainstream PA methods: DPO (Rafailov et al., 2023), ORPO (Hong et al., 2024), and SimPO (Meng et al., 2024). We include the zero-shot results of larger models, including QVQ-72B (Team, 2024), Qwen2.5-VL-72B (Bai et al., 2025), GPT-4V (OpenAI, 2023), GPT-4o-mini, and GPT-4o (OpenAI, 2024), for comprehensive comparison.

Hyper-parameter Settings We implement our two-stage training (SFT&PA) and inference process with two famous open-source projects: LLaMA-Factory (Zheng et al., 2024) and vLLM (Kwon et al., 2023). We conduct our experiments on $8 \times$ NVIDIA A100 GPUs. The max sequence length is set to 8192 and the global batch size to 8 with BF16 precision. We train MLLMs with LoRA (Hu et al., 2022) and search the rank in $\{8, 16\}$. AdamW (Loshchilov & Hutter, 2019) optimizer is used for both training stages with a cosine scheduler. For stage 1, we set the training epoch to 3. The learning rate is searched in $\{1e^{-5}, 1e^{-4}, 3e^{-4}\}$. For stage 2, we train MLLMs with further 1 epoch on the checkpoints of stage 1 and set the learning rate to $1e^{-6}$. The PA data scale $N_2 = 16663$ according to the bad cases we collect from the training set of MMRK-64K after stage 1.

For evaluation, we employ Qwen2.5-VL-72B (Bai et al., 2025) as an LLM judger to score model predictions (CoT and unstructured zero-shot results) against the golden labels, providing a more objective and scalable assessment. The evaluation prompt templates used are detailed in Appendix B.

5.2 MAIN EXPERIMENT RESULTS (RQ1)

We summarize the main experimental results in Table 1, which reports the performance of five MLLM backbones on STAR tasks before and after the proposed two-stage training pipeline. For

each backbone, we evaluate six experimental settings: (1) zero-shot performance, (2–3) stage 1 (SFT) trained with either single-task data or the full STAR-64K dataset, and (4–6) stage 2 (PA) using three alternative methods, DPO, ORPO, and SimPO. Based on these results, we draw the following key observations:

Existing mainstream MLLMs fail on the STAR tasks. From the zero-shot results of OpenAI's GPT models and other open-source MLLMs, it is evident that current leading MLLMs struggle with STAR tasks, indicating that their generalization capabilities do not readily extend to synthetic images and MMKR scenarios. Notably, after SFT with single-task data, Qwen2.5-VL-3B achieves a comparable or even slightly better overall accuracy than much larger models such as Qwen2.5-VL-72B and GPT-40 (41.76% vs. 38.74% / 41.37%). This suggests that the limited performance of current MLLMs on STAR tasks is mainly due to insufficient relevant data during their training phases. Simply applying SFT on single-task data can partially unlock this latent capability, but still requires further fine-grained optimization.

Two-stage pipeline progressively improves the STAR capabilities. Comparing results across the full STAR-64K dataset, we observe that stage 1 SFT leads to substantial improvements over zero-shot performance as models adapt to synthetic multimodal instructions and learn to solve diverse task types. Stage 2 (PA) delivers additional performance gains, albeit smaller than those achieved in stage 1, highlighting the complementary effect of preference alignment. The three PA methods exhibit strong generality and consistently improve results across different backbones. In terms of backbone comparison, LLaVA-1.5-7B consistently underperforms relative to Qwen2.5-VL-7B, whereas LLaVA-NEXT-34B demonstrates clear superiority over Qwen2.5-VL-32B, particularly in counting-related tasks such as EC, RC, IC, and TC. This suggests that both architecture design and scale, alongside our training paradigms, are crucial for advancing STAR performance.

5.3 Transferability Experiments (RQ2)

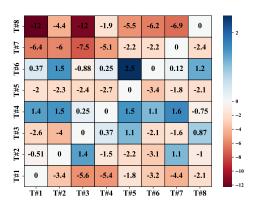


Figure 3: The task-wise transferability experiments. The y-axis is the basic task.

In addition to the main experiments, we conduct a series of supplementary SFT trials using singletask datasets, aiming to investigate the transferability among different STAR tasks. Compared to joint training on the full dataset, single-task training generally results in diminished performance across most tasks, with Task #1 being a notable exception. This suggests that mixed training with diverse multi-task instructions promotes knowledge transfer across tasks and collectively enhances overall model performance. The unique case of Task #1, which relies predominantly on basic entity recognition abilities, indicates that this fundamental capability is not further improved by subsequent training on more complex recognition or reasoning tasks. Complex tasks, on the other hand, present greater learning

challenges for MLLMs, while simpler recognition tasks enable models to better capture underlying patterns in MMKR images. Additionally, to further probe task transferability, we conduct supplementary SFT experiments by pairing tasks during SFT. As illustrated in Figure 3, pairwise task combinations reveal more nuanced mutual enhancement effects, with Tasks 4 and 6 benefiting especially from being trained alongside other tasks. This observation is consistent with findings from the main experiments. Overall, these results suggest that MLLMs can develop emergent STAR capabilities through training on a broader and more complex set of tasks, gradually generalizing to new or related tasks. However, the emergence and effectiveness of such generalization critically depend on the diversity and richness of the training data provided.

5.4 SCALABILITY EXPERIMENTS (RQ3)

We further investigate the scalability of the STAR data, aiming to determine the data volume required to instill fundamental STAR capabilities in MLLMs. In Figure 4, we present the answer and CoT

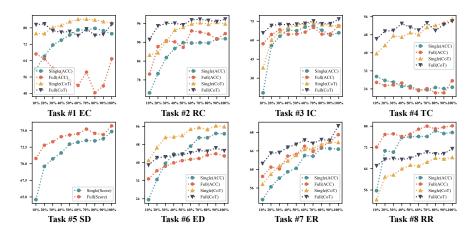


Figure 4: The scalability experiments Qwen2.5-VL-7B for 8 STAR tasks.

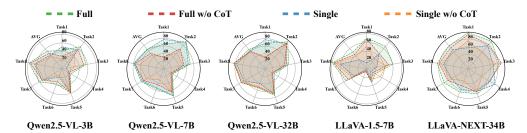


Figure 5: Ablation study on the effectiveness of CoT prompts in the instruction data.

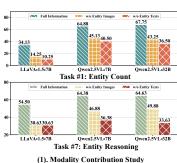
quality results of 8 STAR tasks trained under single-task and full-data settings from 10% to 100% data. The experimental result indicates that, for all other tasks except for Task #1 and Task #4, most tasks exhibit a clear trend of increasing STAR capability as the training data scale grows, consistent with established scale laws in data-driven learning.

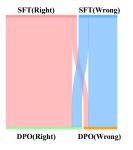
For Task #1, ACC fluctuates with increasing dataset size, whereas CoT quality steadily improves, indicating that while the model's overall counting accuracy does not consistently increase, its precision in entity recognition within CoT becomes progressively better. This can be attributed to the fact that MLLMs possess an inherent counting ability that does not markedly improve with further training. In contrast, their aptitude for recognizing and distinguishing objects in MMKR images advances noticeably. Task #7 follows a similar pattern: ACC remains mostly unchanged, yet CoT quality continues to climb, suggesting that the model is refining its content identification in MMKR images even as its aggregate counting capability remains limited by architectural constraints. Considering the main experiments, the upper limit of this capability on the Qwen2.5-VL-7B model aligns closely with these results. To surpass current limitations, it is necessary to utilize more powerful backbones, as further scaling of training data alone yields diminishing returns for certain task types.

5.5 ABLATION STUDY (RQ4&RQ5)

To further investigate the key factors contributing to the STAR performance, we conduct two ablation studies for the CoT prompts and the modality information incorporated in the MMKR images.

The effectiveness of CoT. As mentioned before, we construct CoT prompts for different tasks to guide MLLMs in identifying and reasoning over the relevant elements within the given MMRK image. To further explore its effectiveness, we conduct several experiments that remove the CoT prompts in the training data. As shown in Figure 5, the STAR performance of all 5 different MLLMs consistently degrades when the CoT prompts are omitted, regardless of whether models are trained on single-task or full multi-task data. These results demonstrate that CoT contributes to performance improvement, providing effective guidance for models to think and solve STAR tasks.







(2). Corrected Cases by Stage 2

(3). Commonsense Knowledge Retention

Figure 6: The modality contribution experiments and case studies.

Table 2: Study of modality contribution on full datasets.

Modality contribution. To synthesize the MMRK images, we incorporate the entity images and texts in the MMKG to construct semantic-rich visualized subgraphs. To assess their impact on them, we conduct SFT experiments for Task #1

Se	Task1	Task2	Task3	Task4	Task5	Task6	Task7	Task8	
Qwen2.5-VL 7B	w/o ent. images w/o ent. texts full dataset	55.50 59.13 64.88	75.88 74.62 92.75	48.62 47.88 71.37	26.63 25.37 27.62	67.99 67.90 75.71	32.00 34.87 55.87	52.63 41.50 67.50	65.75 68.12 80.13
Qwen2.5-VL 32B	w/o ent. images w/o ent. texts full dataset	49.75 58.25 67.75	83.25 82.25 93.63	42.25 41.00 63.13	29.88 25.88 27.50	66.05 65.61 75.07	29.63 28.63 54.00	42.50 46.25 73.50	68.00 66.88 81.75
LLaVA-1.5 7B	w/o ent. images w/o ent. texts full dataset	33.87 37.13 70.38	68.13 67.50 66.75	38.38 42.13 60.75	20.50 21.25 22.25	33.20 33.09 33.24	6.00 6.50 6.50	31.50 33.38 54.50	70.13 69.62 79.62

and Task #7 by re-synthesizing MMKR images without entity images or without texts. These two tasks are entity-centric and are greatly affected by the completeness of entity information. As shown in Figure 6(1), performance drops noticeably on both tasks when either modality is removed, underscoring the value of both visual and textual entity information for effective MLLM reasoning. Notably, omitting entity texts leads to a greater decline, suggesting that textual information is particularly critical. Besides, we present the full task results in Figure 2, which highlight the importance of multi-modal entity information, with textual content playing a dominant role in the STAR tasks.

5.6 CASE STUDY (RQ6)

To provide a more intuitive understanding of the effectiveness of our two-stage training pipeline, we present a case study in this section. Before SFT (stage 1), the model's STAR performance is notably poor, but undergoes marked improvement following SFT. The main experiments further reveal that the second-stage PA process delivers additional gains in STAR accuracy. To identify where these improvements occur, we analyze the distribution of the model's inference results after SFT and subsequent DPO training, as shown in Figure 6(2). Our analysis shows that PA training in the second stage corrects a substantial proportion of erroneous outputs generated after SFT, although some errors persist in a small number of cases. Overall, the higher rate of corrected test predictions confirms a net improvement in model performance. Moreover, expanded case studies presented in Appendix C demonstrate that the stage 2 PA not only increases answer accuracy, but also significantly reduces hallucinations in the CoT reasoning process. Meanwhile, as illustrated in Figure 6(3), we assess the retention of Qwen2.5-VL-7B's commonsense knowledge at various stages using MMMU (Yue et al., 2024a). The results show that two-stage training with STAR-64K not only preserves but, in some domains—such as arts and business—even enhances commonsense knowledge. This demonstrates that STAR capabilities can be effectively integrated into existing MLLMs, improving their performance while maintaining their commonsense reasoning abilities, which would be a win-win strategy. Further details of the commonsense knowledge retention are in Appendix C.

6 Conclusion

In this paper, we investigate structured and abstractive reasoning on images enriched with multimodal relational knowledge for MLLMs. To address this research gap, we design a data engine that synthesizes STAR instruction data and introduces STAR capabilities to models through a customized training and evaluation pipeline. We systematically assess model performance and thoroughly validate the extent of current MLLMs' STAR capabilities, as well as the improvements enabled by our pipeline. Furthermore, we conduct comprehensive analyses of task transferability, data scalability, and design reasonability with intuitive cases to show the effectiveness of our design.

ETHICS STATEMENT

In this paper, we utilize three open-source knowledge graphs (KGs) as our data sources, which we then modify to generate new datasets. Additionally, the primary MLLM backbones we employ are mainstream open-source models. We did not collect data or conduct computational experiments in ways that violated scientific ethics. Therefore, our work does not involve any ethical issues.

REPRODUCIBILITY STATEMENT

We detail the entire pipeline in our methodology section and elaborate on the hyperparameters involved in the experimental settings. Additionally, we provide the relevant pipeline code in the supplementary materials to ensure the reproducibility of this work.

REFERENCES

- Houda Alberts, Teresa Huang, Yash Deshpande, Yibo Liu, Kyunghyun Cho, Clara Vania, and Iacer Calixto. Visualsem: a high-quality knowledge graph for vision and language. *CoRR*, abs/2008.09150, 2020.
- Shuowen An, Si Zhang, Tongyu Guo, Shuang Lu, Wenying Zhang, and Zhihui Cai. Impacts of generative AI on student teachers' task performance and collaborative knowledge construction process in mind mapping-based collaborative environment. *Comput. Educ.*, 227:105227, 2025.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv* preprint arXiv:2502.13923, 2025.
- Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD Conference*, pp. 1247–1250. ACM, 2008.
- Zhuo Chen, Yichi Zhang, Yin Fang, Yuxia Geng, Lingbing Guo, Xiang Chen, Qian Li, Wen Zhang, Jiaoyan Chen, Yushan Zhu, Jiaqi Li, Xiaoze Liu, Jeff Z. Pan, Ningyu Zhang, and Huajun Chen. Knowledge graphs meet multi-modal learning: A comprehensive survey. *CoRR*, abs/2402.05391, 2024.
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, Tianren Gao, Erlong Li, Kun Tang, Zhipeng Cao, Tong Zhou, Ao Liu, Xinrui Yan, Shuqi Mei, Jianguo Cao, Ziran Wang, and Chao Zheng. A survey on multimodal large language models for autonomous driving. In *WACV (Workshops)*, pp. 958–979. IEEE, 2024.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pp. 248–255. IEEE Computer Society, 2009.
- John Ellson, Emden R. Gansner, Eleftherios Koutsofios, Stephen C. North, and Gordon Woodhull. Graphviz and dynagraph static and dynamic graph drawing tools. In *Graph Drawing Software*, pp. 127–148. Springer, 2004.
- Chaoyou Fu, Yifan Zhang, Shukang Yin, Bo Li, Xinyu Fang, Sirui Zhao, Haodong Duan, Xing Sun, Ziwei Liu, Liang Wang, Caifeng Shan, and Ran He. Mme-survey: A comprehensive survey on evaluation of multimodal llms. *CoRR*, abs/2411.15296, 2024.
- Jiwoo Hong, Noah Lee, and James Thorne. ORPO: monolithic preference optimization without reference model. In *EMNLP*, pp. 11170–11189. Association for Computational Linguistics, 2024.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*. OpenReview.net, 2022.

Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, Yi Ren, Yuexian Zou, Zhou Zhao, and Shinji Watanabe. Audiogpt: Understanding and generating speech, music, sound, and talking head. In *AAAI*, pp. 23802–23804. AAAI Press, 2024.

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

- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. From generation to judgment: Opportunities and challenges of llm-as-a-judge. *arXiv preprint arXiv:* 2411.16594, 2024.
- Fenglin Liu, Tingting Zhu, Xian Wu, Bang Yang, Chenyu You, Chenyang Wang, Lei Lu, Zhangdaihong Liu, Yefeng Zheng, Xu Sun, Yang Yang, Lei A. Clifton, and David A. Clifton. A medical multimodal large language model for future pandemics. *npj Digit. Medicine*, 6, 2023a.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023b.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
- Ye Liu, Hui Li, Alberto García-Durán, Mathias Niepert, Daniel Oñoro-Rubio, and David S. Rosenblum. MMKG: multi-modal knowledge graphs. In *ESWC*, volume 11503 of *Lecture Notes in Computer Science*, pp. 459–474. Springer, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR (Poster)*. Open-Review.net, 2019.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *NeurIPS*, 2022.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *ICLR*. OpenReview.net, 2024.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq R. Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *ACL* (*Findings*), pp. 2263–2279. Association for Computational Linguistics, 2022.
- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. In *NeurIPS*, 2024.
- Roberto Navigli and Simone Paolo Ponzetto. Babelnet: Building a very large multilingual semantic network. In *ACL*, pp. 216–225. The Association for Computer Linguistics, 2010.
- OpenAI. GPT-4 technical report. CoRR, abs/2303.08774, 2023.
- OpenAI. Gpt-4o system card, 2024. URL https://arxiv.org/abs/2410.21276.
- Haojie Pan, Yuzhou Zhang, Zepeng Zhai, Ruiji Fu, Ming Liu, Yangqiu Song, Zhongyuan Wang, and Bing Qin. Kuaipedia: a large-scale multi-modal short-video encyclopedia. *CoRR*, abs/2211.00732, 2022.
 - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*, 2023.

- Shezheng Song, Xiaopeng Li, and Shasha Li. How to bridge the gap between modalities: A comprehensive survey on multimodal large language model. *CoRR*, abs/2311.07594, 2023.
- Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: a core of semantic knowledge. In *WWW*, pp. 697–706. ACM, 2007.
 - Qwen Team. Qvq: To see the world with wisdom, December 2024. URL https://qwenlm.github.io/blog/qvq-72b-preview/.
 - Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *CVPR*, pp. 9568–9578. IEEE, 2024.
 - Denny Vrandecic and Markus Krötzsch. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85, 2014.
 - Xueyao Wan and Hang Yu. Mmgraphrag: Bridging vision and language with interpretable multimodal knowledge graphs, 2025. URL https://arxiv.org/abs/2507.20804.
 - Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. In *NeurIPS*, 2024.
 - Derong Xu, Tong Xu, Shiwei Wu, Jingbo Zhou, and Enhong Chen. Relation-enhanced Negative Sampling for Multimodal Knowledge Graph Completion. In *ACM Multimedia*, pp. 3857–3866. ACM, 2022.
 - Liang Yao, Chengsheng Mao, and Yuan Luo. KG-BERT: BERT for knowledge graph completion. *CoRR*, abs/1909.03193, 2019.
 - Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024a.
 - Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. In *CVPR*, pp. 9556–9567. IEEE, 2024b.
 - Boqiang Zhang, Kehan Li, Zesen Cheng, Zhiqiang Hu, Yuqian Yuan, Guanzheng Chen, Sicong Leng, Yuming Jiang, Hang Zhang, Xin Li, Peng Jin, Wenqi Zhang, Fan Wang, Lidong Bing, and Deli Zhao. Videollama 3: Frontier multimodal foundation models for image and video understanding. *CoRR*, abs/2501.13106, 2025a.
 - Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Yichi Zhang, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Shanghang Zhang, Peng Gao, and Hongsheng Li. MAVIS: mathematical visual instruction tuning with an automatic data engine. In *ICLR*. OpenReview.net, 2025b.
 - Wenqi Zhang, Zhenglin Cheng, Yuanyu He, Mengna Wang, Yongliang Shen, Zeqi Tan, Guiyang Hou, Mingqian He, Yanna Ma, Weiming Lu, and Yueting Zhuang. Multimodal self-instruct: Synthetic abstract image and visual reasoning instruction using language model. In *EMNLP*, pp. 19228–19252. Association for Computational Linguistics, 2024a.
 - Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu, Wen Zhang, and Huajun Chen. Making large language models perform better in knowledge graph completion. In *ACM Multimedia*, pp. 233–242. ACM, 2024b.
 - Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu, Binbin Hu, Ziqi Liu, Wen Zhang, and Huajun Chen. Multiple heads are better than one: Mixture of modality knowledge experts for entity representation learning. In *ICLR*. OpenReview.net, 2025c.

Yichi Zhang, Zhuo Chen, Lingbing Guo, Yajing Xu, Min Zhang, Wen Zhang, and Huajun Chen. Abstractive visual understanding of multi-modal structured knowledge: A new perspective for mllm evaluation, 2025d. URL https://arxiv.org/abs/2506.01293.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren's song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219, 2023.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL http://arxiv.org/abs/2403.13372.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models, 2023. URL https://arxiv.org/abs/2304.10592.

Hongyan Zhu, Shuai Qin, Min Su, Chengzhi Lin, Anjie Li, and Junfeng Gao. Harnessing large vision and language models in agriculture: A review. *CoRR*, abs/2407.19679, 2024.

Xun Zhu, Zheng Zhang, Xi Chen, Yiming Shi, Miao Li, and Ji Wu. Connector-s: A survey of connectors in multi-modal large language models. *CoRR*, abs/2502.11453, 2025.

Table 3: Statistical information about the MMKG data source used in our data engine.

Dataset	Entity	Relation	Triple	Data Source
FB15K-237	14541	237	310116	FreeBase
MKG-Y	15000	16	26638	YAGO
VisualSem	89896	13	1481007	Wikipedia, ImageNet, BabelNet

A DETAILS OF THE DATA ENGINE AND TRAINING PIPELINE

A.1 DETAILS OF OUR DATA SOURCE

We present the detailed information of the MMKGs used in our data engine in Table 3, which includes FB15K-237 (Bollacker et al., 2008), MKG-Y (Xu et al., 2022), and VisualSem (Alberts et al., 2020). They are constructed from heterogeneous knowledge bases like FreeBase (Bollacker et al., 2008), YAGO (Suchanek et al., 2007), WikiPedia (Vrandecic & Krötzsch, 2014), ImageNet (Deng et al., 2009), and BabelNet (Navigli & Ponzetto, 2010). These MMKGs encompass diverse entities and relation types, with the knowledge triples they form containing a wide range of encyclopedic and commonsense knowledge. We utilized these MMKGs to construct a large-scale STAR dataset based on our data engine.

A.2 DETAILS OF INSTRUCTION SYNTHESIS

For EC/RC/IC/TC, their CoT prompts are based on several CoT templates that guide MLLMs to recognize the detailed information in V, and the final answer is the proper number of elements.

For SD, its CoT and final answer are combined. We generate a paragraph of words to describe the current MMRK by prompting a strong LLM with the detailed texts of the subgraph.

For ED, the final answer is the disrupted entity, and we generate a CoT for error analysis based on the given knowledge contexts. For ER and RR, the final answer is the option where the entity/relation is masked, and we generate CoT to guide MLLMs' thinking and reasoning.

The overarching principle for constructing CoT is to utilize information from the subgraph before visualization as prompts for the LLM, thereby guiding the LLM to generate the corresponding reasoning process. Compared to presenting MLLMs with a synthesized MMRK image, this approach

yields higher-quality CoT data with fewer hallucinations. Current MLLMs lack sufficient understanding of MMRK images, leading to numerous errors. However, when provided with accurate text prompts of subgraphs, LLMs can generate appropriate results.

A.3 DETAILS OF TRAINING

702

703

704

705 706

707

708

709

710 711 712

713 714

715

716 717

718

719

720

721 722

723 724

725

726

727

728

729 730

731 732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750 751

752 753

754

755

Stage 1. Supervised Fine-tuning The SFT process, following the next token prediction paradigm, can be denoted as:

$$\mathcal{L}_{sft} = -\mathbb{E}_{(\mathcal{V}_i, \mathcal{Q}_i, \mathcal{A}_i) \sim \mathcal{D}_{sft}} \left[\log P_{\mathcal{M}}(\mathcal{A}_i \mid \mathcal{V}_i, \mathcal{Q}_i) \right] \tag{1}$$

 $\mathcal{L}_{sft} = -\mathbb{E}_{(\mathcal{V}_i,\mathcal{Q}_i,\mathcal{A}_i) \sim \mathcal{D}_{sft}} \left[\log P_{\mathcal{M}}(\mathcal{A}_i \mid \mathcal{V}_i,\mathcal{Q}_i) \right]$ where $P_{\mathcal{M}}$ represents the conditional probability of the current answer given by the MLLM \mathcal{M} .

Stage 2. Preference Alignment The PA process with DPO can be denoted as:

$$\mathcal{L}_{pa} = -\mathbb{E}\left[\log \sigma \beta \left(\log \frac{P_{\mathcal{M}}(\mathcal{A}_{i}^{(p)} \mid \mathcal{V}_{i}, \mathcal{Q}_{i})}{P_{\mathcal{M}_{ref}}(\mathcal{A}_{i}^{(p)} \mid \mathcal{V}_{i}, \mathcal{Q}_{i})} - \log \frac{P_{\mathcal{M}}(\mathcal{A}_{i}^{(u)} \mid \mathcal{V}_{i}, \mathcal{Q}_{i})}{P_{\mathcal{M}_{ref}}(\mathcal{A}_{i}^{(u)} \mid \mathcal{V}_{i}, \mathcal{Q}_{i})}\right)\right]$$
(2)

where σ is the sigmoid function and β is the temperature hyper-parameter. \mathcal{M}_{ref} represents the reference model, which is the MLLM trained after stage 1 in practice. Note that in our experiments, other improved versions of DPO, such as ORPO (Hong et al., 2024) and SimPO (Meng et al., 2024), are also employed in stage 2.

INSTRUCTION TEMPLATES

This section presents the instruction templates used in our MMRK data synthesis and performance evaluation process. The instruction templates we presented in this section include: Figure 7: the question templates for 8 STAR tasks; Figure 8: the answer templates (w/ CoT) for 8 STAR tasks; Figure 9: the instruction template used for subgraph description (Task #5) quality evaluation; Figure 10: the instruction template used for CoT quality evaluation for other tasks. Qwen2.5-72B.

CASE STUDY

In this section, we present more detailed case studies to illustrate the effectiveness of the two-stage training pipeline. We present three cases in Figure 11, 12, 13. From these cases, we observe a common pattern: both pure zero-shot results and those trained solely on S1 exhibit severe hallucinations. MLLM generates numerous entities and relations in CoT that are entirely unrelated to the MMRK within the current image, leading to erroneous final answers. However, through targeted optimization in S2, MLLM's hallucinations are suppressed, and its accuracy is evidenced by statistical results. This indicates that our two-stage design has indeed functioned as we expected.

For the commonsense knowledge retention experiments, we employ MMMU (Yue et al., 2024a), which is one of the most popular MLLM benchmarks for commonsense knowledge evaluation. MMMU consists of diverse subjects which can be categorized into arts & designs, business, science, health & medical, humanities & social science, and tech & engineering. We evaluated Qwen2.5-VL-7B's performance across these six domains on its validation set before and after two-stage training. To enable clearer comparison, we applied max-min normalization to the results. The findings reveal that MLLM models trained through the second phase of the STAR task demonstrate improved performance across all domains except science. This contrasts with the common observation that models lose generalizability after instruction-based fine-tuning. This demonstrates that training on the STAR task can activate or enhance the common-sense knowledge of MLLMs to a certain extent without causing catastrophic forgetting. Consequently, it can be integrated as a new capability into existing MLLMs, which underscores the significance of our research.

THE USE OF LARGE LANGUAGE MODELS

The primary research subject of this paper is LLM & MLLM. Additionally, LLMs are employed as a general assistant for code debugging and polishing certain paragraphs. Core idea conception, experimental design, and paper writing are completed by human authors.

1	J	U	
7	5	7	
7	5	8	
7	5	9	
7	6	0	
7	6	1	
	6		
	6		
	6		
	6		
	6		
	6		
	6		
	6		
7			
7	7		
7		2	
	7		
	7		
	7		
7		6	
7		7	
7	7	8	
7	7	9	
7	8	0	
7	8	1	
7	8	2	
7	8	3	
7	8	4	
7	8	5	
7	8	6	
7	8	7	
7	8	8	
7	8	9	
7	9	0	
7	9	1	
7	9	2	
	9		
7	9	4	
	9		
	9		
	9		
	9		
	9		
8			
8			
8			
8			
8			
8			
g	0	0	
24		-//	

756

Question Templates for 8 STAR Tasks

Task #1: Entity Counting

<image>Given the multi-modal knowledge graph. Please count the number of entities in it.

Task #2: Relation Counting

<image>Given the multi-modal knowledge graph. Please count the number of different relations in it.

Task #3: Image Counting

<image>Given the multi-modal knowledge graph. Please count the number of entities that have image information in the given knowledge graph.

Task #4: Triple Counting

<image>Given the multi-modal knowledge graph. Please count the number of knowledge triples in it.

Task #5: Subgraph Description

<image>Given the multi-modal knowledge graph. Please describe the knowledge presented by it.

Task #6: Error Detection

<image>Given the multi-modal knowledge graph. Please point out the wrong entity in it.

Task #7: Entity Reasoning

<image>Given the multi-modal knowledge graph. One entity in it is replaced by [MASK]. Please select one correct entity from the options.

Task #8: Relation Reasoning

<image>Given the multi-modal knowledge graph. One relation in it is replaced by [MASK]. Please select one correct relation from the options.

Figure 7: The question templates for STAR tasks.

Answer (w/ CoT) Templates for 8 STAR Tasks

Task #1: Entity Counting

<think> There are several entities in the given multi-modal knowledge graph: {ENT1, ENT2,, ENT K} Therefore, the number of entities is {ENTITY NUMBER} <answer>{ENTITY NUMBER}</answer>

Task #2: Relation Counting

<think> There are several different relations in the given multi-modal knowledge graph: {REL1, REL2,, REL K} Therefore, the number of different relations is {RELATION NUMBER} NUMBER} Are the several different relations in the given multi-modal knowledge graph: {RELATION NUMBER}

Task #3: Image Counting

<think> There are several entities with images in the given multi-modal knowledge graph: $\{ENT1, ENT2,, ENT M\}$ Other entities without images are: $\{ENT1, ENT2,, ENT N\}$ Therefore, the number of entities is $\{IMAGE \ NUMBER\} < think> <answer> \{IMAGE \ NUMBER\} < think> <answer> \}$

Task #4: Triple Counting

<think> There are several knowledge triples in the given multi-modal knowledge graph:
Therefore, the number of triples is
<think> <answer>

Task #5: Subgraph Description

Description of the subgraph.

Task #6: Error Detection & Task #7: Entity Reasoning & Task #8: Relation Reasoning https://doi.org/10.1007/j.c/ https://doi.o

Figure 8: The answer templates for STAR tasks.

Instruction Template for Task5 Quality Evaluation

As an automated answer-scoring system, please evaluate the similarity between the model's generated responses and the correct answers.

Both of the standard answer and model generated answer are describing a knowledge graph with several sentences.

You must determine whether the key entities, relations, and knowledge mentioned in the model's generated response align with the standard answer.

Ultimately, output an integer between 0 and 100, where a higher number indicates greater similarity. Below are our defined basic scoring rules:

- 0 points: No similarity at all
- 1 to 40 points: Minor information overlap
- 40 to 60 points: Moderate information overlap
- 60 to 90 points: Substantial and detailed information overlap
- Above 90 points: Virtually identical, with only minor syntactic variations

Standard Answer:

Model Generated Answer:

Please response a number for the score directly. Do not provide any other text in the response.

Figure 9: The instruction template used for subgraph description (Task #5) quality evaluation with Qwen2.5-72B.

Instruction Template for Chain-of-thought Quality Evaluation

As an automated answer-scoring system, please evaluate the similarity between the model's generated thought process and the golden label for thought process.

You must determine whether the key entities, relations, and knowledge mentioned in the model's generated thought process align with the standard answer.

Ultimately, output an integer between 0 and 100, where a higher number indicates greater similarity. Below are our defined basic scoring rules:

- 0 points: No similarity at all
- 1 to 30 points: Minor information overlap
- 30 to 60 points: Moderate information overlap
- 60 to 90 points: Substantial and detailed information overlap
- Above 90 points: Virtually identical, with only minor syntactic variations
- If both the thoght process and the final answer match the golden label, full score is awarded.
- If the reasoning process is incorrect but the final answer is correct, partial score may be given.
- If neither the reasoning process nor the final answer is correct, a lower score is assigned.

Standard Thought Process: {}

Model Generated Thought Process: {}

Please response a number for the score directly. Do not provide any other text in the response.

Figure 10: The instruction template used for CoT quality evaluation with Qwen2.5-72B.

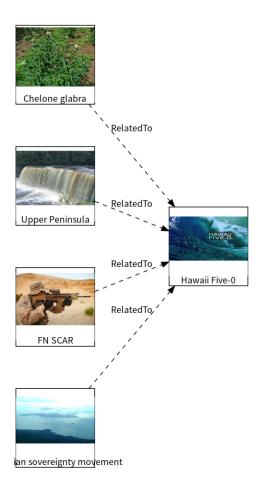


Figure 11: The MMKR image case for Task #1.

Case for Task #1

Zero-shot Answer: The multi-modal knowledge graph contains 10 entities. **Stage 1 Answer:**

\n<think>\nThere are several entities in the given multi-modal knowledge graph:\nUpper Peninsula\nChelone glabra\nChelone glabra\nHawaii Five-0\nFN SCAR\nIndigenous sovereignty movement\nHawaii\nTherefore, the number of entities is 7\n

Stage 2(DPO) Answer:

Golden Answer:

\n<think>\nThere are several entities in the given multi-modal knowledge graph:\nChelone glabra\nHawaii Five-0\nUpper Peninsula\nFN SCAR\nHawaiian sovereignty movement\nTherefore, the number of entities is 5\n

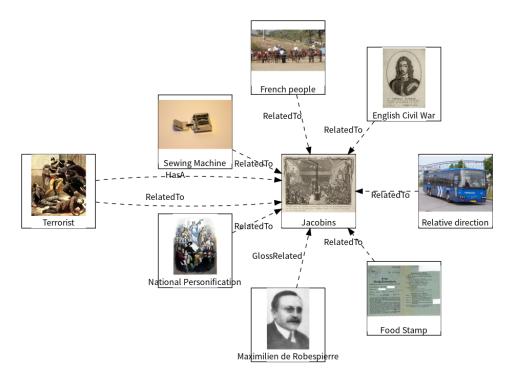


Figure 12: The MMKR image case for Task #2.

Case for Task #2

Zero-shot Answer:

The multi-modal knowledge graph in the image contains 10 different relations. These relations are: $\n = \frac{n}{1}$. Related To $\n = \frac{n}{2}$. Submachine gun $\n = \frac{n}{3}$. Machine Pistol $\n = \frac{n}{3}$. Light Machine Gun $\n = \frac{n}{3}$. Automatic Rifle $\n = \frac{n}{3}$. Firing Port Weapon $\n = \frac{n}{3}$. Colt AR 15 $\n = \frac{n}{3}$. M201 Firing Port Weapon $\n = \frac{n}{3}$. Heckler & Koch G36"

Stage 1 Answer:

 \n think>\nThere are several different relations in the given multi-modal knowledge graph:\nRelatedTo\nGlossRelated\nTherefore, the number of different relations is $2\n<\frac{1}{n}$ is $2\n<\frac{1}{n}$.

Stage 2(DPO) Answer:

Golden Answer:

 \n think>\nThere are several different relations in the given multi-modal knowledge graph:\nRelatedTo\nIsA\nGlossRelated\nTherefore, the number of different relations is $3\n<\t$ hink>\n<answer>3</answer>\n\n

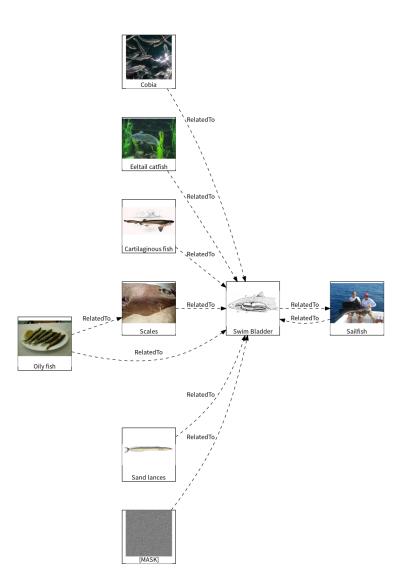


Figure 13: The MMKR image case for Task #7.

Case for Task #7

Zero-shot Answer:

To determine the correct entity that replaces the [MASK] in the given multi-modal knowledge graph, we need to analyze the relationships and the context provided by the other entities and their connections.\n\n1. The graph shows a hierarchy of relationships between different types of fish and their features.\n2. The [MASK] is connected to Scales and Swim Bladder, which are features of fish.\n3. The options provided are:\n - A. Tropical Africa\n - B. Carnivorous Fish\n - C. Bony Fish\n - D. Tropical fish\n - E. Philippine butterflyfish\n\nGiven that the [MASK] is related to Scales and Swim Bladder, it is likely to be a type of fish. Among the options, Carnivorous Fish and Tropical fish are types of fish, but Carnivorous Fish is a specific type of fish, while Tropical fish is a broader category that could include various types of fish.\n\nSince the graph is showing a hierarchy and the [MASK] is connected to Scales and Swim Bladder, it is more likely to be a specific type of fish rather than a broader category like Tropical fish. Therefore, the most appropriate choice is:\n\nB. Carnivorous Fish\n\nSo, the correct answer is B. Carnivorous Fish.

Stage 1 Answer:

\n<think>\nThe correct answer is Bony Fish because it is the most general and inclusive category that can logically relate to Swim Bladder across all the provided triples. The other options are either too specific (e.g., Tropical Africa, Tropical fish) or not directly relevant (e.g., Carnivorous Fish, Philippine butterflyfish). Bony Fish fits well with the general relationship to Swim Bladder and aligns with the context of the other fish-related entities in the triples.\n</hr>

Stage 2(DPO) Answer:

\n<think>\nThe correct entity Tropical fish is more suitable because it directly relates to the biological and ecological context of the Swim Bladder, which is a common feature in many fish species, including tropical fish. The other options are less suitable:\n\n- **A. Tropical Africa**: This is a geographical region and not a biological classification, making it irrelevant to the relationship with Swim Bladder.\n- **B. Carnivorous Fish**: While some fish are carnivorous, this option is too specific and does not encompass all fish that have a swim bladder.\n- **C. Bony Fish**: This is a broader classification but less specific than Tropical fish, which is a subset of bony fish that are commonly found in tropical waters.\n- **E. Philippine butterflyfish**: This is a specific species and not a general category, making it less suitable for the context of the triples.\n

Golden Answer:

\n<think>\nThe correct answer is Tropical fish (D) because it fits well within the context of the triples, which involve various types of fish and their relationship to the swim bladder. The other options are less suitable:\n\n- A. Tropical Africa: This is a region, not a type of fish, and does not fit the pattern of the triples.\n- B. Carnivorous Fish: While some carnivorous fish may have swim bladders, this term is too specific and doesn't encompass the variety of fish mentioned.\n- C. Bony Fish: This is a broader category that includes many fish with swim bladders, but Tropical fish is more specific and aligns better with the other fish types listed.\n- E. Philippine butterflyfish: This is a specific species, which is too narrow compared to the general category of Tropical fish.\n\nTropical fish is the most appropriate choice as it represents a diverse group of fish that often have swim bladders and fits well with the other fish types mentioned in the triples.\n