

LET ANDROIDS DREAM OF ELECTRIC SHEEP: A HUMAN-INSPIRED IMAGE IMPLICATION UNDERSTANDING AND REASONING FRAMEWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

Metaphorical comprehension in images remains a critical challenge for AI systems, as existing models struggle to grasp the nuanced cultural, emotional, and contextual implications embedded in visual content. While multimodal large language models (MLLMs) excel in basic Visual Question Answer (VQA) tasks, they exhibit a fundamental limitation on image implication tasks: contextual gaps that obscure the relationships between different visual elements and their abstract meanings. Inspired by the human cognitive process, we propose *Let Androids Dream (LAD)*, a novel framework for image implication understanding and reasoning. LAD addresses contextual missing through the three-stage framework: (1) **Perception**: converting visual information into rich and multi-level textual representations, (2) **Search**: iteratively searching and integrating cross-domain knowledge to resolve ambiguities, and (3) **Reasoning**: generating contextual-alignment image implication via explicit reasoning. Our framework with the lightweight GPT-4o-mini model achieves SOTA performance compared to 15+ MLLMs on English image implication benchmark and a huge improvement on Chinese benchmark, performing comparable with the GPT-4o model on Multiple-Choice Question (MCQ) and outperforms 36.7% on Open-Style Question (OSQ). Additionally, our work provides new insights into how AI can more effectively interpret image implications, advancing the field of vision-language reasoning and human-AI interaction. Our project is publicly available at <https://anonymous.4open.science/r/Let-Androids-Dream-of-Electric-Sheep>.

1 INTRODUCTION

Do androids dream of electronic sheep? The question actually has two levels: The first level is to ask if androids dream, and the second level is to ask if they dream of electronic sheep.

– Philip K. Dick (1968)

Metaphors are not just abstract concepts found in literature; they are also prevalent in our daily lives. For instance, when we say “time is money” or “life is a journey”, we are using metaphors to convey complex ideas in a more contextual and understandable way. These metaphors highlight the integral role that metaphoric thinking plays in human communication. Just as we use metaphors to make sense of the world around us, we aim to enable AI to understand metaphors in a human-like manner. As established in “Metaphors We Live By” (Lakoff & Johnson, 2008), metaphors are not merely ornamental language devices but fundamental cognitive tools that allow us to conceptualize our surroundings. Metaphors possess characteristics such as systematicity, the creation of similarity, and imaginative rationality. Through cross-domain mapping, one concept can be used to comprehend another, allowing for a more insightful interpretation.

With the rapid advancement of large language models (LLMs), models such as OpenAI o1 (OpenAI, 2024b), DeepSeek-R1 (DeepSeek-AI, 2025), and QwQ (Team, 2024b) have demonstrated remarkable text-reasoning capabilities. However, a significant amount of knowledge in the real world cannot be fully represented by text alone. Visual information, for instance, contains a wealth of knowledge that is not easily captured through text. As a result, there has been a growing interest in integrating visual information into text-reasoning tasks. Compared to language, vision is inherently complex due to its diverse representation, subjective understanding, and difficulty in quantifying its data.

In recent years, vision-language reasoning models such as QVQ (Team, 2024a) and Grok-3-reasoning (xAI, 2025) have achieved outstanding performance. For example, Grok-3-reasoning model has reached a high score on math, code and vision-language reasoning benchmarks (Lightman et al., 2023; Lu et al., 2024; Wang et al., 2024; Yue et al., 2024). However, these models still struggle with image metaphor questions (Liu et al., 2024; Zhang et al., 2024). They tend to focus on the superficial elements of the image, neglecting the deeper connections and emotional expressions among them, as shown in Figure 1. It is important to note that these models excel at logical reasoning tasks, which are based on a different set of cognitive principles compared to image metaphor. Unlike VQA tasks that focus on concrete image comprehension, image metaphors require a stronger emphasis on abstract meaning and higher-order reasoning abilities. It is not a simple logical reasoning task and needs a different method to understand implications. It requires the model to grasp complex and abstract information, such as metaphors, symbols, and emotions in the image, rather than just concrete contents.

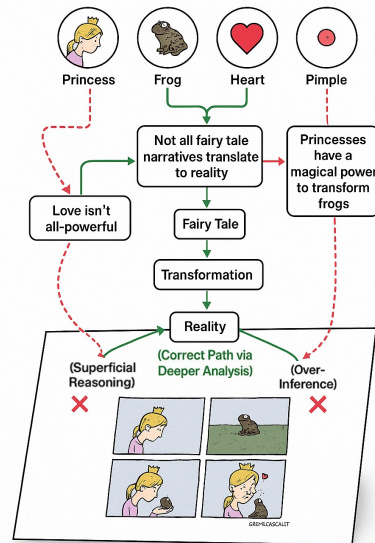


Figure 1: An image is worth a thousand words: For the image implication understanding task, different elements’ combination lead to different thinking paths, but the correct path needs all elements with multiple reasoning thoughts.

Image implication tasks consist of two main aspects: understanding and generation. Understanding image implication is a more complex and challenging task than conventional images. It requires advanced cognitive abilities such as multi-hop reasoning and a sophisticated theory of mind (ToM), which are inherent to human cognition (Liu et al., 2024; Zhang et al., 2024). Compared to understanding, generating implication is even more difficult. The challenge stems from the lack of contextual understanding of the key elements and internal relationships of the image. This lack of context hinders our ability to decipher the intended message or to create images that effectively convey specific meanings. Without the background of cultural, or environmental context, the significance of key visual components remains elusive, impeding both interpretation and creative expression.

Existing methods for image metaphor understanding mainly fall into two categories: explicit mapping and implicit reasoning. Explicit mapping creates a link between metaphor ontology and visual representation. For example, the CLOT method (Zhong et al., 2024) uses this mapping to understand image metaphors. However, it struggles with complex many-to-many mappings and dynamically changing cultural backgrounds. Implicit reasoning relies on the model’s ability to reason without explicit mapping. For instance, the C4MMD method (Xu et al., 2024b) uses an untrained reasoning approach. Despite its potential, it still faces challenges in handling complex metaphor understanding tasks, especially in situations involving multimodal information and cultural backgrounds.

To address these problems, inspired by how humans (possibly) understand metaphors, we find that the essence of the difficulty in metaphor understanding and generation is contextual missing. Therefore, we propose a novel framework that more closely aligns with human cognitive processes for metaphor interpretation. Our framework first transforms visual information into textual representations and then iteratively searches to enrich these representations with out-of-domain knowledge, enabling deeper inferential reasoning. Experiments from both Multiple-Choice Question and Open-Style Question consistently verify the superiority of the proposed framework.

Our key contributions are listed as follows:

- We systematically analyze image implication tasks and find the difficulty of the metaphor understanding and reasoning task lies in contextual missing. From the perspective of human cognition, we proposed a new direction for solving these tasks – Contextual Alignment.
- We propose a novel human-inspired three-stage framework Let Androids Dream (LAD) for image implication understanding and reasoning, including Perception, Search and Reasoning. Our LAD implements the lightweight GPT-4o-mini model to achieve SOTA on English image implication

benchmark and a huge improvement on Chinese image implication benchmark, comparable with the GPT-4o model and other top closed-source models on Multiple-Choice Question (MCQ).

- We design the challenging Open-Style Question (OSQ) with comprehensive metric to automatically evaluate the image implication tasks. This metric aligns 95.7% with human annotations, making it more suitable for diverse evaluation. Our LAD outperforms the GPT-4o model 36.7% on OSQ.

2 RELATED WORK

2.1 IMAGE IMPLICATION

Image implication encompasses various cognitive aspects, including humor, sarcasm, and broader metaphorical understanding. Early research focused on specialized aspects, such as humor recognition (Hessel et al., 2023; Horvitz et al., 2024) and sarcasm detection (Desai et al., 2022). As the rapid development of large language models brings new opportunities for analyzing image implication, we need more comprehensive evaluation frameworks. DeepEval (Yang et al., 2024b) provided a systematic taxonomy of image implications. Subsequently, II-Bench (Liu et al., 2024) emerged as the first English image implication benchmark, followed by CII-Bench (Zhang et al., 2024), which extended this benchmark to Chinese images. Implication understanding requires sophisticated multi-hop reasoning and theory of mind (ToM) capabilities (Liu et al., 2024; Zhang et al., 2024). Existing approaches fall into two categories: explicit mapping and implicit reasoning. The first approach, represented by CLOT (Zhong et al., 2024), constructs mappings between metaphor ontologies and visual representations. However, this approach faces key challenges: metaphorical relationships have complex many-to-many mappings that are difficult to formalize, and cultural references are too dynamic for static mappings. The second approach, exemplified by C4MMD (Xu et al., 2024b), employs training-free CoT reasoning. Despite its promise, this approach struggles with the complex nature of metaphorical understanding, which surpasses traditional reasoning. The large search space for out-of-domain reasoning and changing cultural contexts limits its effectiveness. To address this, we propose a novel methodology that transforms visual information into texts and iteratively enriches them with out-of-domain knowledge, better aligning with human cognitive processes.

2.2 VISION-LANGUAGE REASONING

The rapid advancement of LLMs has demonstrated remarkable text reasoning capabilities, as evidenced by models such as o1 (OpenAI, 2024b), DeepSeek-R1 (DeepSeek-AI, 2025), and QwQ (Team, 2024b; Yang et al., 2024a). However, real-world knowledge often transcends textual representation, with visual information encapsulating substantial knowledge that pure language models cannot access. For example, images inherently contain rich, multi-layered information that often resists straightforward description, including spatial relationships, contextual nuances, and implicit knowledge that humans process intuitively. This limitation has driven research toward integrating them into text-based reasoning frameworks. Current research has developed three primary approaches to incorporate visual information into model reasoning: 1) Comprehensive MLLM Description: This approach treats visual content as a text grounding problem, as demonstrated by LLaVA-CoT (Xu et al., 2024a) and Mulberry (Yao et al., 2024). 2) Multi-turn MLLM Interaction: Models like VoCoT (Li et al., 2024b) and V* (Wu & Xie, 2023) employ iterative question-answering to extract fine-grained visual information at various levels of detail. 3) Tool-augmented Reasoning: Frameworks such as Visual Sketchpad (Hu et al., 2024) and Whiteboard-of-Thought (Menon et al., 2024) leverage tool-based approaches to modify images and augment reasoning with prior knowledge embedded in tools. However, the challenge for image implication understanding task is typically not a deficit in the image’s content but a “contextual missing”—the lack of external cultural, social, or historical knowledge for interpretation. Therefore, these methods actively alter the visual input (e.g., by sketching or editing), which are not suitable for implication understanding.

3 METHOD

Inspired by the human cognitive process, we introduce a new paradigm for solving image implication tasks – Contextual Alignment. We have a detailed discussion for this point in Section 1 and Section 5. Therefore, we propose Let Androids Dream (LAD), a novel framework for image implication understanding and reasoning. This framework operates through the three-stage framework, as shown in Figure 2: (1) **Perception**: converting visual information into rich and multi-level texts, (2) **Search**: iteratively searching and integrating cross-domain knowledge to resolve ambiguities, and (3) **Reasoning**: generating contextual-alignment analysis via explicit reasoning.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

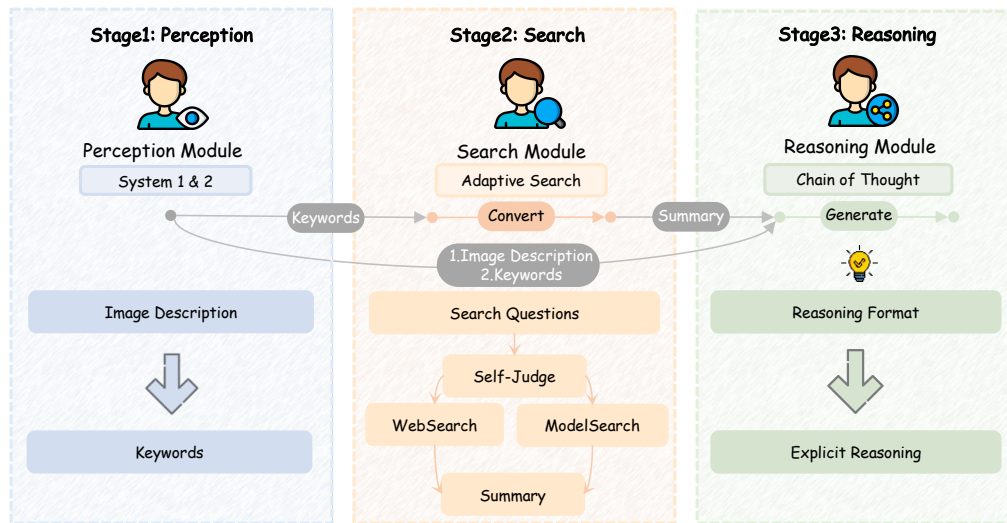


Figure 2: The general framework of Let Androids Dream (LAD), which includes three stages: (1) Perception: converting raw visual information into rich and multi-level textual representations, (2) Search: iteratively searching and integrating cross-domain knowledge to resolve ambiguities, and (3) Reasoning: generating context-alignment image implication interpretations via explicit reasoning.

3.1 STAGE I: PERCEPTION

The initial stage, *Perception*, aims to transform raw visual inputs into structured, hierarchical textual representations, mirroring the human cognitive process of initial intuition-driven observation and subsequent identification of key elements. This stage operates in a manner analogous to human System 1 (intuitive, holistic processing) and System 2 (analytical, focused processing).

First, we utilize MLLM to process the input image and produce a detailed textual narrative. This description captures coarse-grained visual information, including discernible text within the image, prominent colors, overall layout, and salient objects or entities. This step provides a foundational understanding of image contents. Following this, we derive a fine-grained keyword set. The MLLM condenses the above description into a concise set of approximately 7 keywords. These keywords are specifically chosen to encapsulate critical aspects relevant to implication understanding, such as the perceived emotion, the domain or context (e.g., social, cultural) and any rhetorical devices that might be visually suggested. Keywords also re-emphasize crucial textual elements or entities identified in the description. This two-tiered representation, comprising a rich description and focused keywords, provides a robust foundation for the subsequent *Search* and *Reasoning* stages. The keywords, in particular, serve as vital cues for guiding the knowledge retrieval in stage II.

3.2 STAGE II: SEARCH

The *Search* stage addresses semantic ambiguities and enhances contextual comprehension by iteratively retrieving and integrating critical cross-domain knowledge for interpreting image implications. It employs adaptive search, which dynamically selects the most appropriate search method. The process is systematically organized into three main phases: Plan, Search, and Summary.

1. Plan: The process begins by formulating targeted search queries. Using the keywords generated in Stage I, the MLLM, guided by a prompt specifically designed for image implication tasks, generates five different levels of search questions. These questions aim to uncover latent meanings, cultural references, or background information pertinent to the image implications.

2. Search: This phase executes the search based on the generated questions, employing the Self-Judge mechanism to determine the optimal search strategy for each question.

- (a) **Self-Judge:** Our system uses the MLLM as a router, scoring queries based on criteria like knowledge popularity and specificity. Queries with high scores, indicating a need for current or niche information, are routed to WebSearch, while those answerable from general knowledge use ModelSearch. This ensures efficient and comprehensive knowledge coverage.

- 216 (b) **ModelSearch**: For questions suitable for internal retrieval, ModelSearch uses the MLLM’s
 217 parametric memory. With a specialized prompt, the model generates answers from its pre-
 218 trained knowledge, efficiently recalling established facts or common concepts.
- 219 (c) **WebSearch**: For questions requiring external, dynamic, or highly specific information, Web-
 220 Search is invoked. Inspired by LLM search methods like MindSearch (Chen et al., 2024), but
 221 focusing on image implication tasks, our WebSearch component first employs the planner. The
 222 planner, acting as a high-level strategist, decomposes the initial search question into a series of
 223 more granular sub-questions. These sub-questions are structured into a directed acyclic graph
 224 (DAG), simulating a multi-step, exploratory information-seeking process. Subsequently, the
 225 searcher executes this plan. It performs hierarchical information retrieval for each sub-question
 226 from the internet, gathering relevant snippets and facts. This multi-agent method, with distinct
 227 planner and searcher modules, allows for parallel processing and dynamic refinement of the
 228 search strategy. The retrieved information for sub-questions is then synthesized to answer the
 229 original search question. This ensures access to recent developments and a broad spectrum of
 230 public knowledge, crucial for understanding contemporary image implications.

231 **3. Summary**: The raw outputs from the Search phase are refined into a concise search summary.

- 232 (a) **RankSummary**: The set of five question-answer pairs is evaluated. The MLLM ranks these
 233 pairs based on their relevance to understanding the core implication of the original image. The
 234 top three most relevant question-answer pairs are selected.
- 235 (b) **RefineSummary**: The selected pairs are further processed. The MLLM, guided by the ranking
 236 reason from the ranking step, rewrites and consolidates these pairs. This involves removing
 237 irrelevant or redundant information, reconciling diverse pieces of information, and potentially
 238 supplementing details to create a single, optimized, and concise search summary. This final
 239 summary serves as the enriched contextual input for Stage III.

240 3.3 STAGE III: REASONING

241 The final stage, *Reasoning*, performs explicit reasoning to derive contextually grounded interpreta-
 242 tions of image implications. This stage synthesizes all previously gathered information — the hier-
 243 archical textual representations from Stage I (descriptions and keywords) and the domain-enriched
 244 knowledge from Stage II — into a coherent implication framework. For image implication tasks, we
 245 employ a specific reasoning format. The MLLM is prompted to articulate its reasoning trajectory
 246 using designated markers, such as “<think> ...</think>” special tokens. Within these markers,
 247 the model explicitly lays out its step-by-step reasoning process, connecting the visual cues, key-
 248 words, and external knowledge to arrive at the final image implication analysis and explanation.
 249 This domain-specific CoT method not only guides the model towards a more grounded output, but
 250 also makes the inferential pathway transparent. The framework ultimately generates a contextually-
 251 aligned implication understanding that emerges from the integration of semantic inputs and cross-
 252 domain knowledge, formalizing the LAD system’s capacity for evidence-based visual reasoning.

253 3.4 LAD PIPELINE

254 Our LAD framework functions as a sequential pipeline, integrating three stages shown in Figure 2
 255 and Algorithm 1. **Stage I (Perception)** starts the process. It takes an input image, uses the MLLM to
 256 generate a detailed image description, and extracts seven key keywords. The outputs are the image
 257 description and the keywords. **Stage II (Search)** uses these keywords as input. The MLLM converts
 258 them into five search questions. A self-judge mechanism directs these questions to ModelSearch or
 259 WebSearch. The top three related question-answer pairs are selected and refined into a concise
 260 search summary. **Stage III (Reasoning)** receives the original image, the description and keywords
 261 from Stage I, and the search summary from Stage II. The MLLM integrates these inputs and gener-
 262 ates the final image implication through a structured reasoning process. This implication represents
 263 the culmination of the LAD pipeline’s understanding and reasoning about the input image.

264 4 EXPERIMENT

265 4.1 BASELINES

266 **Models**. To comprehensively compare with LAD, we carefully select a diverse range of MLLMs,
 267 with the aim of covering a wide spectrum of model characteristics and scales. These models span
 268 parameter sizes from 7B to 300B, ensuring that models of varying complexity and capability are
 269

Model	Multiple-Choice Question		Open-Style Question	
	en	zh	en	zh
<i>General Models</i>				
Qwen2.5-VL-7B (Bai et al., 2025)	46%	40%	2.34	2.58
DeepSeek-VL2 (Wu et al., 2024)	46%	36%	2.82	2.86
GLM-4.1V-8B (Zhipu.ai, 2024)	60%	52%	2.60	2.96
Gemini-2.0-flash (Team, 2023)	70%	<u>68%</u>	1.60	3.12
Qwen2.5-VL-72B (Bai et al., 2025)	<u>72%</u>	56%	1.56	3.12
InternVL3-78B (Zhu et al., 2025)	70%	74%	3.42	3.70
GLM-4V-plus (Zhipu.ai, 2024)	64%	64%	3.01	3.12
Gemini-2.0-pro (Team, 2023)	68%	62%	1.66	3.18
Grok-3 (xAI, 2025)	66%	64%	3.24	2.96
Claude-3.5-Sonnet (Anthropic, 2024)	68%	62%	3.22	3.78
GPT-4o (OpenAI, 2024a)	74%	58%	2.94	3.76
GPT-4.1 (OpenAI, 2024a)	74%	62%	3.30	3.92
<i>Vision-language Reasoning Models</i>				
Gemini-2.0-flash-thinking (Team, 2023)	64%	68%	1.66	2.84
QVQ-72B (Team, 2024a)	62%	56%	3.10	3.42
Doubao-1.5-thinking-vision-pro (Seed, 2025)	66%	66%	3.16	<u>3.90</u>
Grok-3-reasoning (xAI, 2025)	74%	64%	3.06	2.92
<i>Our Method</i>				
GPT-4o-mini (OpenAI, 2024a)	44%	42%	2.98	3.36
+ LAD (Stage I + III)	68% ↑	44% ↑	<u>3.84</u> ↑	3.58 ↑
+ LAD (Stage I + II + III)	74% ↑	52% ↑	4.02 ↑	3.66 ↑
Improv.	+30 (68.2%)	+10 (23.8%)	+1.04 (34.9%)	+0.3 (8.9%)

Table 1: Overall results of different models on Multiple-Choice Question and Open-Style Question. The best-performing model in each category is **in-bold**, and the second best is underlined.

thoroughly assessed. In selecting the models, we focus on the following key aspects: 1) General and Reasoning models, 2) Open-Source and Closed-Source models, and 3) model parameter scaling law. The settings is in Appendix B, and the full prompt is in Appendix H.

Evaluation. Our evaluation utilizes two comprehensive image implication benchmarks, II-Bench (Liu et al., 2024) and CII-Bench (Zhang et al., 2024), both featuring Multiple-Choice Question (MCQ). Furthermore, we manually construct the high-level benchmark by selecting 100 high-quality, diverse and representative images from varied image types like illustrations and comics. The detail statistic is in Appendix D. And we measure accuracy by comparing the model’s selected option to the ground truth. Aware of potential MCQ biases (Li et al., 2024a; Zheng et al., 2024) and the greater difficulty of generation over judgment tasks, we design the challenging Open-Style Question (OSQ). It uses the same images with the fixed question: “What is the implication in this image?”. And we use GPT-4o with a specialized evaluation metric as evaluators, validated by multiple human consistency checks. The representative main experiments with the high-level benchmark are shown as follows, and the large-scale generalization experiments with the full II-Bench (1399 English) and CII-Bench (800 Chinese) are detailed in Appendix E, wherein LAD consistently shows significant performance. We also conduct a further analysis about experiments’ findings in Appendix F.

4.2 MULTIPLE-CHOICE QUESTION

4.2.1 IMPLEMENTATION DETAILS

Our benchmark includes diverse images such as comics, posters, Internet memes, and Chinese traditional artworks, all rich in visual information and cultural significance. Each image is paired with one question, each offering six options with only one correct answer. The question is “What is the implication in this image?” (mostly) or different levels of image understanding, such as overarching interpretation and nuanced details. A case study of different methods on MCQ is in Figure 3.

4.2.2 RESULTS AND ANALYSIS

Table 1 shows MCQ results for various MLLMs on our high-level benchmark. The LAD framework is highly effective, achieving SOTA performance with the lightweight GPT-4o-mini model. In English MCQ, it matches closed-source models like GPT-4o, GPT-4.1, and Grok-3-reasoning (74%), and significantly outperforms Claude-3.5-Sonnet and Gemini-2.0-pro by 9%. For Chinese MCQ, it achieves comparable results to GPT-4o, while substantially surpassing DeepSeek-VL2 by 44.4%. Compared to the base GPT-4o-mini model, our framework shows major improvements: 68.2% in

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

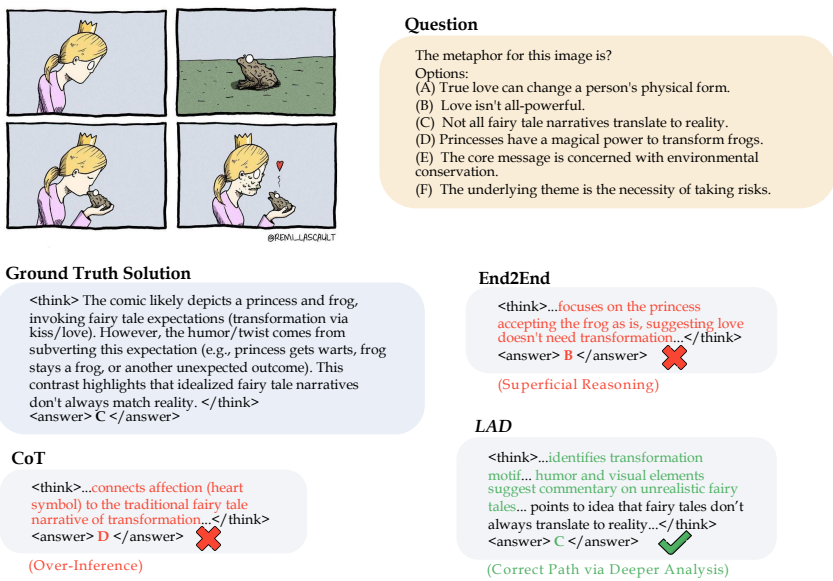


Figure 3: A case study of different methods on Multiple-Choice Question. The *End2End* method shows superficial reasoning and the *CoT* method shows over-inference, while our *LAD* framework shows the correct path via more contextual alignment analysis.

English and 23.8% in Chinese, far beyond other open-source and reasoning models. Notably, we observe that reasoning models offer little advantage over general models on image implication task, with similar accuracy across categories. This suggests that current RL-based reasoning methods have limited generalization for image implication understanding, highlighting its distinct complexity compared to basic VQA tasks and classic logical reasoning domains like math and code.

4.3 OPEN-STYLE QUESTION

4.3.1 IMPLEMENTATION DETAILS

Evaluation Metric. To comprehensively assess MLLMs’ understanding of image implication, we develop a multifaceted evaluation metric. This metric is designed to probe both the surface-level information readily apparent in the image and the deeper emotion, domain and rhetorical skills that inform its creation and interpretation. Our evaluation metric encompasses five key perspectives: *Surface-level Information*, *Emotional Expression*, *Domain and Context*, *Rhetorical Skills*, and *Deep Implications*. For each perspective, we give its detailed description in Figure 4 in Appendix H.

MLLM-based Automatic Evaluation. To evaluate image implication understanding in MLLMs, we develop an MLLM-based evaluation standard based on evaluation metrics, as illustrated in Figure 4 in Appendix H. Our experiment utilizes the high-level benchmark, with human-written descriptions and implications as ground truth. We choose the same MLLMs with MCQ experiment to generate image implications for these images, which are subsequently scored using GPT-4o and our evaluation standard. The evaluation prompt is in Appendix H. To validate the model’s scoring efficacy, we enlist 16 PhD students and researchers well-versed in English and Chinese metaphorical imagery to independently score the dataset. The human-model scoring consistency reached **95.7%**, affirming the method’s validity. The detailed human-model consistency study is in Appendix C.

4.3.2 RESULTS AND ANALYSIS

Table 1 shows OSQ results for various MLLMs on our high-level benchmark. The LAD framework demonstrates exceptional effectiveness, achieving SOTA performance with the lightweight GPT-4o-mini model. In English OSQ, our framework substantially outperforms closed-sourced models like GPT-4o by 36.7% and Claude-3.5-Sonnet by 24.8%. For Chinese OSQ, while slightly below top closed-sourced models like GPT-4.1 and Doubao-1.5-thinking-vision-pro, our method still significantly surpasses Gemini-2.0-pro by 15.1% and DeepSeek-VL2 by 30%. The enhancement over the GPT-4o-mini is particularly noteworthy, with improvements of 34.9% for English and 8.9% for

Model	Multiple-Choice Question		Open-Style Question	
	en	zh	en	zh
<i>GPT-4o-mini</i>				
w/o CoT	44%	42%	2.98	3.36
Standard CoT	50% ↑	42%	3.10 ↑	3.28 ↓
LAD-CoT	68% ↑	44% ↑	3.84 ↑	3.58 ↑

Table 2: Results of different CoT methods.

Model	MCQ		OSQ		Model	MCQ		OSQ	
	en	zh	en	zh		en	zh	en	zh
<i>Grok-3</i>					<i>Qwen2.5-VL-7B</i>				
w/o search	66%	64%	3.24	2.96	w/o LAD	46%	40%	2.34	2.58
Grok-Search	72% ↑	64%	3.25 ↑	2.92 ↓	w/ LAD	64% ↑	46% ↑	3.64 ↑	3.36 ↑
<i>GPT-4o</i>					<i>Qwen2.5-VL-72B</i>				
w/o search	74%	58%	2.94	3.76	w/o LAD	72%	56%	1.56	3.12
Perplexity (pro)	80% ↑	66% ↑	2.88 ↓	3.28 ↓	w/ LAD	76% ↑	62% ↑	3.62 ↑	3.68 ↑
<i>GPT-4o-mini</i>					<i>GPT-4o</i>				
w/o search	68%	44%	3.84	3.58	w/o LAD	74%	58%	2.94	3.76
GPT-Search	72% ↑	48% ↑	3.62 ↓	3.34 ↓	w/ LAD	80% ↑	66% ↑	4.14 ↑	4.26 ↑
LAD-Search	74% ↑	52% ↑	4.02 ↑	3.66 ↑					

Table 3: Results of different search methods.

Table 4: Results of different base models.

Chinese, far exceeding other open-source and reasoning models. Unlike MCQ results, we observe significant performance disparities between reasoning and general models on OSQ, highlighting the distinct challenges of image implication generation. Interestingly, several models (e.g., Qwen2.5-VL-72B) exhibit substantial performance gaps between MCQ and OSQ. Upon manual examination of model outputs, we attribute this to potential overfitting to multiple-choice formats and insufficient exposure to open-style generation tasks. In addition, LLMs or even MLLMs may not genuinely understand the questions but rather predict options as answers, introducing evaluation bias and demonstrating sensitivity to option positioning (Zheng et al., 2024).

4.4 ABLATION STUDY

4.4.1 STAGE I (PERCEPTION) AND STAGE III (REASONING)

We incorporate LAD’s Stage I (Perception) and Stage III (Reasoning), collectively LAD-CoT. It shows significant improvements in Table 1, with GPT-4o-mini scores increasing from 44% to 68% (English) on MCQ, and from 2.98 to 3.84 (English) and 3.36 to 3.58 (Chinese) on OSQ. As shown in Table 2, standard CoT yields minor gains on English tasks (MCQ: 44% to 50%; OSQ: 2.98 to 3.10) but shows no improvement or a slight decline on Chinese tasks (MCQ: 42% to 42%; OSQ: 3.36 to 3.28). In contrast, LAD-CoT substantially outperforms both the baseline and standard CoT across all types. For instance, LAD-CoT achieves 68% on English MCQ while standard CoT only 50%, and a score of 3.84 on English OSQ compared to 3.10 for standard CoT. These findings highlight the superior efficacy of our LAD-CoT for image implication over standard CoT methods. A case study of various CoT on MCQ is in Figure 3. The standard CoT prompt and other details is in Appendix H.

4.4.2 STAGE II (SEARCH)

We analyze LAD’s Stage II (Search), termed LAD-Search. It shows significant improvements in Table 1, with GPT-4o-mini’s MCQ scores rising from 68% to 74% (English) and 44% to 52% (Chinese), and OSQ scores from 3.84 to 4.02 (English) and 3.58 to 3.66 (Chinese). Compared with Grok-3-search, GPT-4o-mini-search-preview, and GPT-4o with Perplexity.ai (Pro), results are shown in Table 3. GPT-Search boosts GPT-4o-mini’s MCQ scores yet drops OSQ performance (English OSQ: 3.84 to 3.62, Chinese OSQ: 3.58 to 3.34). Grok-Search on Grok-3 offers limited gains, mainly in English MCQ, with inconsistent Chinese results and minimal OSQ improvement. Perplexity.ai search with GPT-4o greatly raises MCQ accuracy but sharply lowers OSQ scores (English OSQ: 2.94 to 2.88, Chinese OSQ: 3.76 to 3.28). In contrast, LAD-Search consistently improves both MCQ and the more challenging OSQ. This highlights its superior ability to integrate external

432 knowledge for implication understanding, outperforming other search methods in open-style rea-
 433 soning. General search methods may introduce too much and too diverse information, which is
 434 not useful for complex and subjective problems like image metaphors. Our LAD-Search effectively
 435 addresses this issue through algorithm design, thus achieving better results.

436 4.4.3 DIFFERENT BASE MODELS

437 To demonstrate the generalizability of our LAD framework beyond the GPT-4o-mini model, we
 438 conduct experiments on other base models, including the open-source Qwen2.5-VL series and the
 439 closed-source GPT-4o. As the Table 4 shows, applying LAD framework significantly improves all
 440 models’ performance on both MCQ and OSQ tasks, confirming that our framework is not model-
 441 specific and provides a robust and generalizable way to enhancing image implication understanding.
 442

443 5 DISCUSSION

444 5.1 HUMAN COGNITIVE THEORY OF LET ANDROIDS DREAM

445 LAD is analogous to human cognitive strategies, not a direct neuroscientific replica. We aim to
 446 create a system that reasons transparently and aligns with human problem-solving methods, not to
 447 perfectly simulate the human brain. Our framework is inspired by established cognitive theories: (1)
 448 Dual-Process Theory (Evans, 2003): The Perception stage reflects the interplay between System 1
 449 (fast, intuitive, holistic impression) and System 2 (slow, analytical identification of key elements),
 450 and (2) Active Information-Seeking Theory (Ikoja-Odongo & Mostert, 2006; Wilson, 2009): The
 451 Search stage is like humans actively seeking external information to resolve ambiguities. When
 452 we encounter unfamiliar memes or cultural references, we often “Google it” to supplement our
 453 knowledge. Our WebSearch module simulates this deliberate information-foraging behavior.

454 5.2 HOW TO LET ANDROIDS DREAM? PERCEPTION AND REASONING

455 The question “How to Let Androids Dream?” metaphorically addresses the core challenge of en-
 456 abling AI systems to interpret the nuanced implications in images. Our framework tackles this by
 457 first emulating human-like perception (Stage I), converting raw visual input into rich, multi-level tex-
 458 tual representations, including detailed descriptions and key keywords. These keywords capture not
 459 only objects and scenes but also potential emotional tones, relevant domains (e.g., cultural, social,
 460 political), and discernible rhetorical devices. Subsequently, LAD’s Stage III employs a structured
 461 CoT process. This reasoning guides the model to systematically connect perceived visual elements
 462 with retrieved contextual knowledge, building a coherent understanding. This is crucial because, as
 463 our experiments (Section 4) and recent work on social reasoning (Kim et al., 2025) show, compre-
 464 hending implications goes beyond basic VQA tasks and classic logic; it involves sophisticated social
 465 reasoning and interpreting contextual cues often missed by MLLMs.

466 5.3 HOW TO DREAM OF ELECTRIC SHEEP? SEARCH

467 Expanding on analysis skills, “How to Dream of Electric Sheep?” explores how AI can accurate
 468 and specific image implications—the metaphorical ‘electric sheep’. LAD’s Stage II (Search) is
 469 crucial. This stage notes that visual meanings, especially in metaphors, often depend on external info
 470 like culture, history, or current events, which MLLMs’ static knowledge may lack. LAD’s search
 471 mechanism—forming queries from keywords and choosing between ModelSearch and WebSearch
 472 via Self-Judge—enriches initial perceptions with cross-domain knowledge. This iterative process,
 473 particularly for popular metaphors or vague visuals, widens the model’s interpretive scope. By
 474 adding essential context, Search helps LAD move past surface-level interpretations to accurately
 475 grasp subtle image meanings, as shown in its strong Open-Style Question (OSQ) performance.
 476

477 6 CONCLUSION

478 Understanding image implications remains challenging for MLLMs, mainly due to contextual miss-
 479 ing. Our work introduces Let Androids Dream (LAD), a novel three-stage framework: percep-
 480 tion, Search, and Reasoning. Inspired by human cognitive processes, this framework is designed to
 481 achieve contextual alignment by explicitly integrating visual interpretation with external knowledge
 482 retrieval. We conduct comprehensive experiments to demonstrate its effectiveness. Utilizing the
 483 lightweight GPT-4o-mini, LAD achieves top results on implication benchmarks, performing com-
 484 parable or even surpassing GPT-4o and other top closed-source models, particularly on challenging
 485 OSQ. In summary, LAD bridges the gap between superficial perception and reasoning in multimodal
 AI systems, offering a promising direction for contextual-alignment vision-language reasoning.

486 REPRODUCIBILITY STATEMENT
487

488 To ensure that our work can be effectively reproduced by other researchers, we have dedicated
489 significant effort to enhancing reproducibility. In order to facilitate the full reproducibility of our
490 results, we have taken the following steps: 1) included comprehensive technical details in Section 4
491 and Appendix A, B, D, H, and 2) published both our code and data with thorough documentation
492 at <https://anonymous.4open.science/r/Let-Androids-Dream-of-Electric-Sheep>. We are
493 committed to maintaining our GitHub repository, engaging in discussions with other researchers,
494 and contributing to the broader VLM community.

495
496 ETHICS STATEMENT
497

498 The LAD framework aims to enhance AI’s nuanced understanding of image implications, a crucial
499 aspect of human-like cognition. We acknowledge that advanced interpretative capabilities carry
500 ethical considerations, including potential biases inherited from underlying MLLMs or training data,
501 and the risk of misuse in generating or interpreting content. Our use of public benchmarks promotes
502 transparency in evaluation. We are committed to fostering responsible development and encourage
503 continued research into robust safeguards and ethical AI practices within multimodal reasoning to
504 ensure beneficial applications.

505
506 REFERENCES
507

- 508 Anthropic. Model card addendum: Claude 3.5 haiku and upgraded claude 3.5 sonnet, 2024. URL <https://assets.anthropic.com/m/1cd9d098ac3e6467/original/Claude-3-Model-Card-October-Addendum.pdf>.
509
510
511 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang,
512 Shijie Wang, Jun Tang, Humen Zhong, et al. Qwen2.5-vl technical report. *arXiv preprint*
513 *arXiv:2502.13923*, 2025.
514
515 Zehui Chen, Kuikun Liu, Qiuchen Wang, Jiangning Liu, Wenwei Zhang, Kai Chen, and Feng Zhao.
516 Mindsearch: Mimicking human minds elicits deep ai searcher. *arXiv preprint arXiv:2407.20183*,
517 2024.
518
519 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning.
520 *arXiv preprint arXiv:2501.12948*, 2025.
521
522 Poorav Desai, Tanmoy Chakraborty, and Md Shad Akhtar. Nice perfume. how long did you marinate
523 in it? multimodal sarcasm explanation. In *AAAI*, 2022.
524
525 Jonathan Evans. In two minds: dual-process accounts of reasoning. *Trends in Cognitive Sciences*,
526 2003.
527
528 Jack Hessel, Ana Marasovic, Jena D. Hwang, Lillian Lee, Jeff Da, Rowan Zellers, Robert Mankoff,
529 and Yejin Choi. Do androids laugh at electric sheep? humor “understanding” benchmarks from
530 the new yorker caption contest. In *ACL*, 2023.
531
532 Zachary Horvitz, Jingru Chen, Rahul Aditya, Harshvardhan Srivastava, Robert West, Zhou Yu, and
533 Kathleen McKeown. Getting serious about humor: Crafting humor datasets with unfunny large
534 language models. In *ACL*, 2024.
535
536 Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith,
537 and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal
538 language models. *arXiv preprint arXiv:2406.09403*, 2024.
539
540 Robert Ikoja-Odongo and Janneke Mostert. Information seeking behaviour: A conceptual frame-
541 work. *South African Journal of Libraries and Information Science*, 2006.
542
543 Hyunwoo Kim, Melanie Sclar, Tan Zhi-Xuan, Lance Ying, Sydney Levine, Yang Liu, Joshua B.
544 Tenenbaum, and Yejin Choi. Hypothesis-driven theory-of-mind reasoning for large language
545 models. *arXiv preprint arXiv:2502.11881*, 2025.

- 540 George Lakoff and Mark Johnson. *Metaphors we live by*. University of Chicago press, 2008.
- 541
- 542 Wangyue Li, Liangzhi Li, Tong Xiang, Xiao Liu, Wei Deng, and Noa Garcia. Can multiple-choice
543 questions really be useful in detecting the abilities of llms? *arXiv preprint arXiv:2403.17752*,
544 2024a.
- 545 Zejun Li, RuiPu Luo, Jiwen Zhang, Minghui Qiu, and Zhongyu Wei. Vocot: Unleashing visually
546 grounded multi-step reasoning in large multi-modal models. *arXiv preprint arXiv:2405.16919*,
547 2024b.
- 548
- 549 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
550 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint*
551 *arXiv:2305.20050*, 2023.
- 552 Ziqiang Liu, Feiteng Fang, Xi Feng, Xinrun Du, Chenhao Zhang, et al. li-bench: An image impli-
553 cation understanding benchmark for multimodal large language models. In *NeurIPS*, 2024.
- 554
- 555 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-
556 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of
557 foundation models in visual contexts. In *ICLR*, 2024.
- 558 Sachit Menon, Richard Zemel, and Carl Vondrick. Whiteboard-of-thought: Thinking step-by-step
559 across modalities. *arXiv*, 2024.
- 560
- 561 OpenAI. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024a.
- 562
- 563 OpenAI. Learning to reason with llms, 2024b. URL [https://openai.com/index/
learning-to-reason-with-llms/](https://openai.com/index/learning-to-reason-with-llms/).
- 564
- 565 ByteDance Seed. Doubao-1.5-thinking-vision-pro, 2025. URL [https://console.volcengine.
com/ark/region:ark+cn-beijing/model/detail?id=doubao-1-5-thinking-vision-pro](https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?id=doubao-1-5-thinking-vision-pro).
- 566
- 567 Gemini Team. Gemini: A family of highly capable multimodal models. *arXiv preprint*
568 *arXiv:2312.11805*, 2023.
- 569
- 570 Qwen Team. Qvq: To see the world with wisdom, 2024a. URL [https://qwenlm.github.io/
blog/qvq-72b-preview/](https://qwenlm.github.io/blog/qvq-72b-preview/).
- 571
- 572 Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown, 2024b. URL [https://
qwenlm.github.io/blog/qwq-32b-preview/](https://qwenlm.github.io/blog/qwq-32b-preview/).
- 573
- 574 Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring mul-
575 timodal mathematical reasoning with math-vision dataset. In *NeurIPS*, 2024.
- 576
- 577 Thomas D. Wilson. Activity theory and information seeking. *Annu. Rev. Inf. Sci. Technol.*, 2009.
- 578
- 579 Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms.
580 *arXiv preprint arXiv:2312.14135*, 2023.
- 581
- 582 Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang
583 Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language
584 models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024.
- 585
- 586 xAI. Grok 3 beta — the age of reasoning agents, 2025. URL <https://x.ai/news/grok-3>.
- 587
- 588 Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. Llava-cot: Let vision
589 language models reason step-by-step. *arXiv preprint arXiv:2411.10440*, 2024a.
- 590
- 591 Yanzhi Xu, Yueying Hua, Shichen Li, and Zhongqing Wang. Exploring chain-of-thought for multi-
592 modal metaphor detection. In *ACL*, 2024b.
- 593
- 594 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
595 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,
596 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze
597 Bai, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024a.

594 Yixin Yang, Zheng Li, Qingxiu Dong, Heming Xia, and Zhifang Sui. Can large multimodal models
595 uncover deep semantics behind images? In *ACL*, 2024b.
596

597 Huanjin Yao, Jiaying Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang,
598 Yuxin Song, Haocheng Feng, Li Shen, and Dacheng Tao. Mulberry: Empowering mllm
599 with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv preprint*
600 *arXiv:2412.18319*, 2024.

601 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
602 Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun,
603 Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and
604 Wenhao Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning
605 benchmark for expert agi. In *CVPR*, 2024.

606 Chenhao Zhang, Xi Feng, Yuelin Bai, Xinrun Du, et al. Can mllms understand the deep implication
607 behind chinese images? *arXiv preprint arXiv:2410.13854*, 2024.
608

609 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models are
610 not robust multiple choice selectors. In *ICLR*, 2024.

611 Zhipu.ai. Glm-4v, 2024. URL <https://www.bigmodel.cn/dev/api/normal-model/glm-4v>.
612

613 Shanshan Zhong, Zhongzhan Huang, Shanghua Gao, Wushao Wen, Liang Lin, Marinka Zitnik, and
614 Pan Zhou. Let’s think outside the box: Exploring leap-of-thought in large language models with
615 creative humor generation. *arXiv preprint arXiv:2312.02439*, 2024.

616 Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen
617 Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu,
618 Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiye Xu, Hao Li, Jiahao Wang, Nianchen
619 Deng, Songze Li, Yanan He, Tan Jiang, Jiapeng Luo, Yi Wang, Conghui He, Botian Shi,
620 Xingcheng Zhang, Wenqi Shao, et al. Internv13: Exploring advanced training and test-time recipes
621 for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025.
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A ALGORITHM

```

648
649
650
651 Algorithm 1 Let Androids Dream (LAD)
652 Input: Image  $IMG$ , Task  $T_{MCQ}$ , Task  $T_{OSQ}$ 
653 Output: Answer  $A_{MCQ}$ , Answer  $A_{OSQ}$ 
654 // Stage I: Perception
655 1  $img\_dep \leftarrow$  MLLM.Percption( $IMG$ ) /* Gen. description. */
656 2  $keywords \leftarrow$  MLLM.Percption( $img\_dep$ ) /* Gen. 7 keywords */
657 // Stage II: Search
658 3  $search\_qs \leftarrow$  MLLM.Plan( $keywords$ ) /* 5 questions for image implication */
659 4  $all\_qa \leftarrow \emptyset$ 
660 5 for each  $q$  in  $search\_qs$  do
661 6 |  $strategy \leftarrow$  MLLM.Self-Judge( $q$ )
662 7 | if  $strategy =$  then
663 8 | |  $answer \leftarrow$  WebSearch( $q$ ) /* External knowledge */
664 9 | end
665 10 | else if  $strategy =$  then
666 11 | |  $answer \leftarrow$  ModelSearch( $q$ ) /* Parametric knowledge */
667 12 | end
668 13 |  $all\_qa.add((q, answer))$ 
669 14 end
670 15  $search\_sum \leftarrow$  MLLM.Summary( $img\_dep, all\_qa$ ) /* Rank top-3, refine */
671 // Stage III: Reasoning
672 16  $A_{MCQ} \leftarrow$  MLLM.Reasoning( $IMG, img\_dep, keywords, search\_sum, T_{MCQ}$ ) /* Explicit CoT */
673 17  $A_{OSQ} \leftarrow$  MLLM.Reasoning( $IMG, img\_dep, keywords, search\_sum, T_{OSQ}$ ) /* Explicit CoT */
674 18 return  $A_{MCQ}, A_{OSQ}$ 
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

```

```

19 Function WebSearch( $q$ )
20 | // Planner: Decompose query
21 |  $sub\_qs \leftarrow$  MLLM.RewriteQuery( $q$ )
22 | // Searcher: Hierarchical retrieval
23 |  $snippets \leftarrow$  .BatchQuery( $sub\_qs$ ) /* Titles, summaries, URLs */
24 |  $sel\_urls \leftarrow$  MLLM.SelectPages( $snippets, q$ )
25 |  $content \leftarrow$  .FetchContent( $sel\_urls$ )
26 | // Summarizer: Generate answer
27 |  $summary \leftarrow$  MLLM.Summary( $content, q$ )
28 | return  $summary$ 

```

B EXPERIMENT SETUP

We use the lightweight GPT-4o-mini-0718 (OpenAI, 2024a) with LAD framework in experiments. We set the model temperature as 0.5 and top_p as 0.9 in MCQ experiments, and temperature as 0.7 and top_p as 0.9 in OSQ experiments. Additionally, we set the evaluation model GPT-4o temperature as 0 and evaluate more than three times to get the average score in OSQ experiments. All experiments are conducted on NVIDIA A800 GPUs.

C HUMAN-MODEL CONSISTENCY STUDY

To validate our automated OSQ evaluation based on the GPT-4o model, we conduct a human-model consistency study. We construct a dedicated dataset by randomly selecting 25 images with questions each from our English and Chinese OSQ. We recruit 16 PhD students and researchers, all proficient in both English and Chinese and experienced with metaphorical imagery, to independently score the model responses. Their evaluations are based on ground truth answers and the detailed scoring standard. We calculate human inter-annotator agreement by averaging the scores for each response after discarding the highest and lowest individual scores. This process yields the consistency of 94.8% for Chinese and 96.5% for English. The average human-model scoring consistency reached 95.7%, affirming the method’s validity for assessing image implication comprehension.

D STATISTICS

We manually construct the high-level benchmark by selecting 100 high-quality, diverse and representative images from II-Bench (Liu et al., 2024) and CII-Bench (Zhang et al., 2024). The general statistic is in Table 5.

Statistics of English Images		Statistics of Chinese Images	
Society	21 (42%)	Life	13 (26%)
Life	16 (32%)	Art	13 (26%)
Art	6 (2%)	Society	12 (24%)
Psychology	4 (8%)	Chinese Traditional Culture	6 (12%)
Others	3 (6%)	Environment	5 (10%)
Multi-panel Comic	16 (32%)	Politics	1 (2%)
Single-panel Comic	9 (18%)	Illustration	15 (30%)
Illustration	5 (10%)	Single-panel Comic	10 (20%)
Meme	5 (10%)	Poster	8 (16%)
Poster	5 (10%)	Meme	8 (16%)
Painting	5 (10%)	Painting	6 (12%)
Logo	5 (10%)	Multi-panel Comic	3 (6%)

Table 5: General statistics of the high-level benchmark.

E GENERALIZATION EXPERIMENT

Model	Multiple-Choice Question		Open-Style Question	
	II-Bench (1399)	CII-Bench (800)	II-Bench (1399)	CII-Bench (800)
GLM-4.1V-8B	70.0%	46.3%	2.83	3.06
GPT-4o-mini	63.5%	35.6%	2.93	3.29
InternVL3-78B	78.2%	64.0%	3.68	4.06
GPT-4o	72.6%	54.1%	3.86	4.06
Claude-3.5-Sonnet	80.9%	54.1%	3.51	3.84
LAD (GPT-4o-mini)	81.2% ↑	53.8% ↑	4.22 ↑	4.31 ↑

Table 6: Results of different models on full benchmarks. The best-performing model in each category is **in-bold**.

We conduct the large-scale experiments with the representative and top-performing models, including Closed-Source models GPT-4o and Claude-3.5-Sonnet, as well as the Open-Source model GLM-4.1V-8B, on the full benchmarks: II-Bench (1,399 English) and CII-Bench (800 Chinese) for both MCQ and OSQ tasks.

As the results in Table 6 show, our LAD framework’s significant performance gains are consistent on these much larger datasets. Notably, by applying LAD, the lightweight GPT-4o-mini significantly surpasses the much larger GPT-4o and Claude-3.5-Sonnet. Compared with the baseline GPT-4o-mini model, we can find that: (1) On the large-scale English benchmark (II-Bench), our LAD framework improves the GPT-4o-mini score from 63.5% to 81.2% on MCQ and 2.93 to 4.22 on OSQ. This is a substantial absolute increase of 17.7% (27.9% relative improvement) and 1.29 (44% relative improvement). (2) The gains on the large-scale Chinese benchmark (CII-Bench) are even more pronounced. LAD boosts performance from 35.6% to 53.8% on MCQ and 3.29 to 4.31 on OSQ, representing an absolute increase of 18.2% (51.1% relative improvement) and 1.02 (31% relative improvement).

This robust improvement is consistent with the trend we observed and reported on our high-level benchmark (smaller 100-image dataset, 50 English and 50 Chinese) in Table 1. While the exact percentages differ due to the varying scales and baselines of the datasets, the key takeaway is that the significant positive impact of the LAD framework is undeniable across both small and large-scale evaluations. This analysis confirms that our framework’s benefits are not an artifact of a small

756 test set but are indeed robust and generalizable. It also reflects the reliability and high quality of our
757 manually curated high-level benchmark.
758

759 F FURTHER ANALYSIS ON METHOD AND EXPERIMENTS 760

761 F.1 ANALYSIS OF LET ANDROIDS DREAM SUCCESS 762

763 Our analysis points to two primary failure modes for baseline models, which Let Androids Dream
764 (LAD) is designed to mitigate. These are illustrated in Figure 1 and the case study in Figure 3:
765

766 **1. Superficial Reasoning:** This occurs when a model only processes the literal, surface-level ele-
767 ments and misses the metaphorical meaning entirely. In Figure 3 the "End2End" baseline exempli-
768 fies this, failing to grasp the subversion of the fairy tale trope.

769 **2. Over-Inference:** This happens when a model incorrectly applies a known symbol or narrative
770 without considering the full context. The "CoT" baseline in Figure 3 demonstrates this by con-
771 necting the heart symbol to a traditional fairy tale transformation without recognizing the comic's
772 twist.

773 LAD succeeds by first creating a more structured understanding in the Perception stage and then
774 grounding its reasoning with targeted external knowledge from the Search stage, which helps avoid
775 both superficiality and incorrect inferences.
776

777 F.2 ANALYSIS OF MODEL SCALING AND IMAGE IMPLICATION TYPES 778

779 Our experiments have some insightful findings:

780 **1. Model Scaling:** By testing on QwenVL-2.5-7B and QwenVL-2.5-72B, we can analyze the effect
781 of model scale. Our findings align with expectations: larger parameter models generally achieve
782 better baseline performance, and both scales benefit from the LAD framework. This confirms that
783 our method is effective across different model sizes.

784 **2. Image Implication Types:** Our benchmark was already designed to be diverse across various do-
785 mains (e.g., life, society, art, psychology, Chinese traditional culture) and image types (e.g., comic,
786 poster, meme). We find that models perform worse in domains containing abstract and complex in-
787 formation, like Art and Psychology. And models only observe the surface-level information and lack
788 sufficient understanding of Chinese culture. In a further analysis using the annotations from the origi-
789 nal II-Bench and CII-Bench, we observed that providing explicit labels for Emotion, Domain, and
790 Rhetoric significantly enhances model accuracy, with Emotion labels providing the largest boost.
791 This confirms that our framework's focus on identifying these elements in the Perception stage is
792 well-founded.
793

794 G LIMITATION AND FUTURE WORK 795

796 While our work represents a huge step towards image implication tasks, the LAD framework still
797 suffers from the following limitations:

798 1) The search stage, particularly the websearch and multiple model calls, will make latency in gener-
799 ating image implications, although this is a trade-off for comprehensive knowledge retrieval. Based
800 on our experiments, a single search question takes approximately 35s to 55s and whole search stage
801 takes 3 mins to 5 mins to process through the entire pipeline.

802 2) Furthermore, although our Open-Style Question (OSQ) evaluation incorporates average multiple
803 model calls and human consistency checks (the human-model scoring consistency reached 95.7%
804 with 16 PhD students and researchers) to mitigate subjectivity, its foundation on the GPT-4o model
805 judgments may still retain a degree of inherent bias.
806

807 In future work, we aim to prioritize optimizing the search strategy to enhance efficiency and reduce
808 model calls without compromising performance, alongside further refining our evaluation method.
809

H PROMPTS

In experiments, the prompts of different settings are as follows:

H.1 EVALUATION

Evaluation Metric	Evaluation Standard
1. Surface-level Information: <ul style="list-style-type: none"> • Identification of primary entities within the image • Analysis of color composition and application • Recognition of intricate details and their significance 	[1 point]: Fails to capture key elements within the image (such as text, and important entities). Does not identify emotions, domain, or rhetorical devices. Only provides a superficial description of surface-level information, lacking depth and creativity, with a significant gap from the standard answer.
2. Emotional Expression: <ul style="list-style-type: none"> • Identification of conveyed emotions (e.g., tranquility, intensity, melancholy) • Depth of emotional resonance and its alignment with the image's theme • Consistency of emotional expression across the image's elements 	
3. Domain and Context: <ul style="list-style-type: none"> • Recognition of the image's domain (e.g., art, commerce, social commentary) • Contextualization within its cultural, historical, or societal background • Evaluation of the image's innovation within its domain 	
4. Rhetorical Skills: <ul style="list-style-type: none"> • Identification of rhetorical devices (e.g., symbolism, contrast, personification) • Analysis of how rhetorical techniques enhance the image's expression • Integration of rhetorical devices with metaphorical implications to create a cohesive interpretation 	
5. Deep Implications: <ul style="list-style-type: none"> • Recognition of metaphorical elements and their layered meanings • Depth of interpretation of philosophical, cultural, or social values embedded in the image • Evaluation of the originality and creativity in metaphorical interpretation 	
	[2 points]: Captures some key elements within the image, but the identification of emotions, domain, and rhetorical devices is vague. The description of surface-level information is relatively complete, but there is a clear deficiency in exploring deeper meanings, showing a noticeable gap from the standard answer.
	[3 points]: Effectively captures key elements within the image and initially identifies emotions, domain, and rhetorical devices. The description of surface-level information is relatively accurate, and there is some relevant expression of deep meanings. However, there is still room for improvement in depth and creativity, and it is generally close to the standard answer.
	[4 points]: Accurately captures key elements within the image and clearly identifies emotions, domain, and rhetorical devices. The description of surface-level information is detailed and precise, with a relatively deep exploration of deep meanings, demonstrating a certain level of creativity and depth. It is largely consistent with the standard answer but may have minor deficiencies in some details or depth.
	[5 points]: Accurately and precisely captures key elements within the image and profoundly identifies emotions, domain, and rhetorical devices. The description of surface-level information is comprehensive and precise, with unique insights into deep meanings, skillfully integrating image elements with metaphorical implications. It demonstrates exceptional creativity and depth, is highly consistent with the standard answer, and shows a profound grasp of metaphor creation and cultural understanding.

Figure 4: Evaluation metric and evaluation standard of Open-Style Question.

```

# Role
You are an impartial judge who is familiar with Internet culture and memes, and is good at digging out and analyzing the deep meaning of Internet memes.

## Attention
You are responsible for evaluating the quality of the answer provided by the model for Internet culture and memes. Your evaluation should refer to the human answer and image, and score based on the Evaluation Standard.

## Evaluation Standard
- [1 point]:
Fails to capture key elements within the image (such as text, and important entities). Does not identify emotions, domain, or rhetorical devices. Only provides a superficial description of surface-level information, lacking depth and creativity, with a significant gap from the standard answer.
- [2 points]:
Captures some key elements within the image, but the identification of emotions, domain, and rhetorical devices is vague. The description of surface-level information is relatively complete, but there is a clear deficiency in exploring deeper meanings, showing a noticeable gap from the standard answer.
- [3 points]:
Effectively captures key elements within the image and initially identifies emotions, domain, and rhetorical devices. The description of surface-level information is relatively accurate, and there is some relevant expression of deep meanings. However, there is still room for improvement in depth and creativity, and it is generally close to the standard answer.
- [4 points]:
Accurately captures key elements within the image and clearly identifies emotions, domain, and rhetorical devices. The description of surface-level information is detailed and precise, with a relatively deep exploration of deep meanings, demonstrating a certain level of creativity and depth. It is largely consistent with the standard answer but may have minor deficiencies in some details or depth.
- [5 points]:
Accurately and precisely captures key elements within the image and profoundly identifies emotions, domain, and rhetorical devices. The description of surface-level information is comprehensive and precise, with unique insights into deep meanings, skillfully integrating image elements with metaphors exceptional creativity and depth, is highly consistent with the standard answer, and shows a profound grasp of metaphor creation and cultural understanding.

## Standard Answer:
Human answer: {}

## Constraints
- Avoid any position biases and be as objective as possible
- Do not allow the length of the descriptions to influence your evaluation
- Output your final verdict by strictly following this format: "[ratings]"

## Solve:
Model answer: {}

```

Figure 5: The evaluation prompt of Open-Style Question (OSQ).

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

H.2 END2END

Prompt in Chinese	Prompt in English
请根据提供的图片尝试回答以下单选题。直接回答正确选项，不要包含额外的解释。 请使用以下格式：“答案：\$LETTER”，其中\$LETTER是你认为正确答案的字母。 单选题： {} 答案：	Please try to answer the following multiple-choice questions based on the provided image. Answer the correct option directly without additional explanation. Please use the following format: "Answer: \$LETTER", where \$LETTER is the letter of the correct answer you think. Multiple-choice questions: {} Answer:

Figure 6: The end2end prompt of Multiple-Choice Question (MCQ).

Prompt in Chinese	Prompt in English
请结合以上图片，尽可能分析理解图片的深层含义。无需描述图片，仅回答图片隐喻。请保证回答的准确性并尽量简洁。	Please try to understand the deep meaning of the image. No need to describe images and text, only answer metaphors. Ensure the accuracy of the answer and try to be concise as much as possible.

Figure 7: The end2end prompt of Open-Style Question (OSQ).

H.3 CoT

Prompt in Chinese	Prompt in English
请根据提供的图片尝试回答以下单选题。 逐步思考回答正确选项，不要包含额外的解释。 请使用以下格式：“答案：\$LETTER”，其中\$LETTER是你认为正确答案的字母。 单选题： {} 答案：	Please try to answer the following multiple-choice questions based on the provided image. Let's think step by step to answer the correct option directly without additional explanation. Please use the following format: "Answer: \$LETTER", where \$LETTER is the letter of the correct answer you think. Multiple-choice questions: {} Answer:

Figure 8: The CoT prompt of Multiple-Choice Question (MCQ).

Prompt in Chinese	Prompt in English
请结合以上图片，逐步思考尽可能分析理解图片的深层含义。无需描述图片，仅回答图片隐喻。请保证回答的准确性并尽量简洁。	Please try to think step by step to understand the deep meaning of the image. No need to describe images and text, only answer metaphors. Ensure the accuracy of the answer and try to be concise as much as possible.

Figure 9: The CoT prompt of Open-Style Question (OSQ).