NOT HOW YOU THINK, IT'S WHAT YOU SEE: DECOU-PLING PERCEPTION FROM REASONING

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

018

019

021

024

025

026

027

028

029

031

034

037

040

041

042

043

044

045

046

047

048

Paper under double-blind review

ABSTRACT

The ability of Vision-Language Models (VLMs) to reason depends on a complex interplay between visual perception and abstract cognition. While it is widely recognized that perception is a significant bottleneck, systematically diagnosing how it fails and developing methods to unlock latent reasoning capabilities remains a key challenge. To address this, we introduce a cognitively-inspired framework that decomposes VLM behavior through four distinct paradigms: 1) Direct Visual Rule Learning (holistic processing), 2) Deductive Rule Learning (explicit rule extraction), 3) Componential Analysis (CA), which decouples perception by reasoning over task-agnostic textual descriptions, and 4) Interactive Componential Analysis (ICA), which introduces a feedback loop for targeted visual probing. Our framework's emphasis on task-agnostic decomposition and cognitive parallels provides a unique lens for analysis compared to prior decoupling efforts. Applying this framework across an expanded suite of benchmarks, we conduct a comprehensive evaluation on both proprietary and open-source multi-image VLMs. Our results confirm that perception is a primary bottleneck and show that our CA and ICA paradigms yield substantial performance gains, unlocking the latent reasoning abilities of powerful LLMs. Crucially, ICA demonstrates that an interactive loop can resolve fine-grained visual ambiguities that static descriptions cannot, outperforming the non-interactive CA approach. Our work provides a robust diagnostic toolkit for the community and offers concrete architectural insights, demonstrating that interactive, decoupled systems are a promising path toward more general and capable visual intelligence.

1 Introduction

Human cognition adeptly integrates visual perception with abstract reasoning to navigate the world (Kunda, 2020; Lake et al., 2017). While Vision-Language Models (VLMs) have made remarkable progress (Zhao et al., 2023; Radford et al., 2021), their ability to perform complex visual reasoning remains brittle. It is widely recognized that a primary failure point is the model's visual perception, but we lack systematic tools to diagnose how these perceptual systems fail and frameworks to mitigate these failures to unlock latent reasoning.

To investigate these questions, we test models on two challenging task families that probe the limits of perception and reasoning. First, Bongard Problems (BPs) (Bongard, 1968), a classic test requiring few-shot discovery of an abstract visual rule. We use natural image variants, Bongard-OW (Wu et al., 2024) and Bongard-HOI (Jiang et al., 2022). Second, Winoground (Thrush et al., 2022), which tests visio-linguistic compositional reasoning through minimally contrastive image-caption pairs.

This paper introduces an evaluation framework designed to dissect the cognitive processes of VLMs on these tasks. Our core contribution is a framework grounded in cognitive science that, unlike prior work that inventories static capabilities, uses dynamic paradigms to model problem-solving strategies:

- Direct Visual Rule Learning (DVRL): Simulates holistic processing (Biederman, 1987), where the model analyzes all images simultaneously.
- Deductive Rule Learning (DRL): Mimics explicit, rule-based deduction (Rips, 1994), separating rule extraction from application.

068

069

071 072 073

074 075

080

081

082 084

085 087

096

090

098 099 100

101

102 103

104 105

106 107

- 3. Componential Analysis (CA): Parallels analytical decomposition (Gluck et al., 2008), reasoning over structured, task-agnostic textual descriptions.
- 4. **Interactive Componential Analysis (ICA):** Extends CA with a feedback loop, allowing the reasoning module to actively probe the perception module for targeted details.

This framework allows for systematic analysis of VLM behavior, identifying specific processing bottlenecks. By generating comprehensive, task-agnostic image descriptions, our componential paradigms allow us to disentangle perception from reasoning. ICA enhances this by enabling a dynamic perceptual process guided by reasoning needs. This approach facilitates multi-image reasoning on single-image architectures and even allows us to evaluate text-only LLMs by providing high-quality descriptions (Section 7.2).

Applying this framework across a broadened set of benchmarks and models, we find that our CA and ICA paradigms yield substantial performance gains. These methods achieve highly competitive results on Bongard-OW, Bongard-HOI, and Winoground, primarily by pairing high-fidelity descriptions with powerful reasoning models. This success across diverse tasks highlights the robustness of decoupling perception from reasoning. Concurrently, our analysis confirms a significant perception bottleneck, as models' performance drastically improves when their perceptual front-end is bypassed or guided.

Our contributions are:

- 1. A novel, cognitively-inspired framework with four distinct paradigms (including interactive reasoning) for the diagnostic evaluation of VLMs.
- 2. A componential method (CA and ICA) to disentangle perception from reasoning, enabling multi-image task evaluation for diverse architectures and unlocking latent reasoning in LLMs.
- 3. Comprehensive empirical results on multiple benchmarks and models, confirming the perception bottleneck and demonstrating that our interactive, decoupled methods can significantly mitigate it.
- 4. A demonstration of the effectiveness of this approach, achieving strong, competitive performance on challenging visual reasoning tasks.

RELATED WORK

VLM Benchmarks. The evaluation of VLMs has rapidly evolved from foundational tasks like VQA (Antol et al., 2015) to benchmarks testing more complex reasoning. Recent efforts focus on multi-image understanding through interleaved corpora (Laurençon et al., 2024) and dedicated benchmarks such as MuirBench (Wang et al., 2024) or low-level perception tests like BLINK (Fu et al., 2024). Our work complements these by focusing on benchmarks specifically designed to probe core cognitive abilities that resist simple linguistic mediation: abstract few-shot rule discovery using natural image Bongard Problems (BPs) (Wu et al., 2024; Jiang et al., 2022) and fine-grained compositional reasoning via Winoground (Thrush et al., 2022).

Cognitive Science Grounding. Our framework is grounded in cognitive science perspectives on human problem-solving (Newell et al., 1972). Our paradigms model distinct cognitive strategies: Direct Visual Rule Learning (DVRL) mirrors rapid, holistic processing (Biederman, 1987); Deductive Rule Learning (DRL) reflects explicit, rule-based deduction (Rips, 1994); and our Componential Analysis (CA) paradigms parallel analytical decomposition, where problems are broken into constituent parts for systematic reasoning (Gluck et al., 2008).

Decoupling Frameworks and Chain-of-Thought. Our approach is related to a growing body of work on decoupling perception and reasoning in VLMs. This includes Multimodal Chain-of-Thought (CoT) prompting (Zhang et al., 2023; Zheng et al., 2023) and dedicated evaluation frameworks. For instance, Prism (Qiao et al., 2024) provides a valuable framework for assessing a static inventory of fine-grained VLM skills like object recognition and counting.

Our contribution is distinct in three key ways. First, our paradigms are grounded in cognitive processes (e.g., holistic vs. deductive), analyzing how a model solves a problem, not just what skills it possesses. Second, our CA paradigm deliberately uses task-agnostic descriptions, creating a

clean separation between raw perceptual capability and downstream reasoning, unlike many CoT methods that generate task-conditioned descriptions. Finally, our new Interactive CA (ICA) paradigm introduces a novel dynamic feedback loop, where the reasoning module actively probes the perception module. This moves beyond the static, one-pass evaluations common in prior work to model a more realistic, iterative reasoning process.

3 Models

We evaluated a diverse suite of VLMs, distinguishing between models based on their multi-image context capacity. State-of-the-art proprietary systems, including **GPT-4o** (OpenAI, 2024) and **Gemini 2.0** (Google, 2024), can natively process the large number of images (13) required for our DVRL and DRL paradigms on Bongard Problems. In contrast, while some contemporary open-source models like **Pixtral-12B** (Agrawal et al., 2024), **Llama-Vision-3.2** (Meta, 2024), and **LLaVA** variants (Liu et al., 2023; XTuner, 2025) accept multiple images, none currently support the large context required for a direct evaluation in these paradigms. This technical constraint underscores the necessity of our Componential Analysis (CA) paradigm as the primary method for assessing complex, multi-image reasoning on these powerful open-source architectures. For ablation studies, we also used text-only LLMs (**Llama3** (Grattafiori et al., 2024), **Phi-4** (Abdin et al., 2024), etc.). All evaluations used few-shot prompting at zero temperature. Further details are in Appendix A.4.

4 Dataset and Task

We test our framework on a diverse suite of benchmarks chosen to probe distinct cognitive abilities. Our primary testbed is Bongard-OW (Wu et al., 2024), a 500-case subset testing few-shot abstract rule discovery on natural images (see Figure 1 for an example). To assess generalization, we use Bongard-HOI (Jiang et al., 2022) (400 samples) to evaluate reasoning about human-object interactions, and Winoground (Thrush et al., 2022) (400 samples) for fine-grained compositional grounding. Together, these tasks provide a robust testbed for analyzing high-level visual reasoning, from abstraction to compositionality, across our different cognitive paradigms.





Figure 1: Example Bongard-OW task. *Left*: Positive examples. *Center*: Negative examples. *Right*: Query. Rule: *A group photo at a wedding reception*. Query is negative. (3 of 6 examples shown per set).

5 COGNITIVELY-INSPIRED EVALUATION PARADIGMS

We evaluate VLMs using four paradigms designed to probe different facets of visual reasoning and assess performance under systematically varied cognitive demands, inspired by human cognitive strategies. All paradigms require the model to output a structured response including analysis, the derived rule, query description, and classification (positive/negative). Figure 2 provides a schematic overview. Specific prompts are detailed in Appendix A.5.

5.1 DIRECT VISUAL RULE LEARNING (DVRL)

This paradigm assesses holistic reasoning by presenting all 13 images (6 positive, 6 negative, 1 query) simultaneously to the VLM. It demands the model integrate information across the entire set to identify the rule and classify the query in one step. This mirrors the human ability to quickly grasp the 'gist' of a visual scene or problem. Due to requiring simultaneous multi-image input, only models like Gemini 2.0 and GPT-40 were tested under this paradigm.

5.2 DEDUCTIVE RULE LEARNING (DRL)

Mimicking deliberative, rule-based deduction, DRL involves two stages:

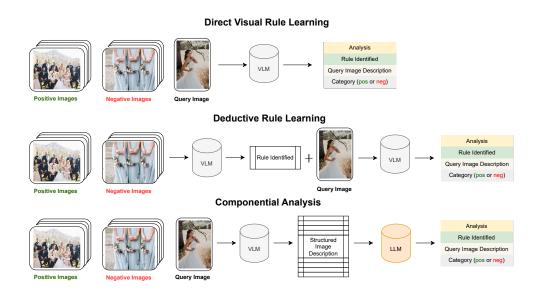


Figure 2: Cognitively-Inspired Evaluation Paradigms. DVRL (Direct Visual Rule Learning): Concurrent processing of all images, mimicking holistic perception. Requires multi-image input capability. DRL (Deductive Rule Learning): Two-stage process separating rule extraction from application, mimicking explicit deduction. CA (Componential Analysis): Multi-stage process involving individual image description followed by reasoning over text, mimicking analytical decomposition and enabling perception-reasoning separation.

- 1. **Rule Extraction:** The VLM analyzes the 12 context images (positive/negative sets) to identify and concisely summarize (max 20 words) the distinguishing rule.
- Rule Application: The VLM receives the previously generated rule summary and the query image, classifying the query based solely on the provided rule.

This separation allows examining the fidelity of both rule formation and rule application processes.

5.3 COMPONENTIAL ANALYSIS (CA)

Reflecting analytical problem decomposition, CA proceeds in stages based on textual representations:

- 1. **Image Description:** The VLM generates a detailed, structured, and ideally *task-agnostic* JSON description for each of the 13 images *individually*.
- 2. **Text-Based Reasoning:** A powerful LLM receives the collection of 13 JSON descriptions (labeled positive/negative/query) and performs rule extraction and query classification based *only* on this textual input.

This paradigm is crucial as it (a) allows evaluating models lacking direct multi-image input, (b) enables assessing reasoning largely independent of perceptual errors, and (c) facilitates evaluation of text-only LLMs on visual reasoning tasks.

5.4 INTERACTIVE COMPONENTIAL ANALYSIS (ICA)

To mitigate the limitations of static perception revealed by CA, this paradigm extends it with a dynamic, multi-step feedback loop that emulates a "look again" strategy.

 Initial Description: Same as in CA, the VLM generates initial, task-agnostic descriptions for each image.

- Ambiguity Identification & Question Formulation: A reasoning LLM analyzes the initial
 descriptions and the task goal (e.g., the Winoground captions) to identify the most critical,
 ambiguous visual detail needed for a confident decision. It then formulates a specific,
 targeted question about this detail.
- 3. **Focused Re-Perception:** The VLM is shown the relevant image again, but this time is asked to answer only the targeted question from the previous step.
- 4. **Synthesized Reasoning:** The LLM integrates the initial descriptions with the new, high-precision information from the Q&A step to make its final classification.

This interactive process allows the model to actively resolve perceptual ambiguities, moving beyond a single static "glance" to perform more robust, human-like visual verification.

6 RESULTS AND ANALYSIS

This section details the performance of the evaluated VLMs, beginning with the primary Bongard-OW benchmark and then examining generalizability.

6.1 Performance on Bongard-OW

Table 1 presents the core results on our 500-sample Bongard-OW subset. Under **Direct Visual Rule Learning** (**DVRL**), applicable only to GPT-40 and Gemini 2.0, performance was strong but below optimal (Gemini 2.0: 82.2%, GPT-40: 80.0%), suggesting limitations in purely holistic, simultaneous multimage reasoning for this complex task.

Performance improved markedly under **DRL** for both models (GPT-40: 88.0%, Gemini 2.0: 86.8%). The explicit separation of rule extraction and application stages appears beneficial, aligning with the idea that breaking down complex cognitive tasks can improve performance.

Model	DVRL	DRL	CA
GPT-40	80.0	88.0	92.8
Gemini 2.0	82.2	86.8	93.6
Pixtral-12B	-	-	87.2
Llama-Vision-11B	-	-	53.4
Llama-Vision-90B	-	-	55.1
Llava-7B	-	-	66.2
Llava-Llama3-8B	-	-	53.2
Prior SOTA	(GPT-4 +	- InstructBLIP)	63.8
Human Average	(acro	ss samples)	91.0

Table 1: Classification accuracy (%) across evaluation paradigms on the Bongard-OW subset. Paradigms abbreviated: DVRL, DRL, CA. Dashes (-) indicate non-applicability due to model input limitations.

Componential Analysis (CA), reasoning over textual descriptions, yielded the

highest accuracies for the top models (GPT-40: 92.8%, Gemini 2.0: 93.6%). Notably, these results establish a **new state-of-the-art (SOTA)** on the Bongard-OW Text-score benchmark, surpassing the reported human average (91.0% Wu et al. (2024)). The previous best machine performance Wu et al. (2024) involved using GPT-4 to reason over captions generated by models such as InstructBLIP, achieving a maximum accuracy of 63.8%. Our significant improvement with the CA paradigm underscores the efficacy of its comprehensive, task-agnostic description generation (Stage 1) coupled with advanced reasoning engines (Stage 2). Pixtral-12B also achieved strong CA performance (87.2%). However, a significant gap emerged with other open-source models. Models like Llama-Vision and LLaVA variants exhibited much lower CA accuracy, often with dramatic imbalances between positive and negative sample performance (e.g., Llava-Llama3-8B: 53.2% overall, heavily biased towards negative samples). This pattern strongly suggests that the bottleneck for these models is not necessarily the abstract reasoning itself, but rather the fidelity of their internal *visual perception* and subsequent translation into usable representations. (here, text descriptions).

The consistent trend of accuracy increasing from DVRL through DRL to CA for GPT-40 and Gemini 2.0 further reinforces the value of structured reasoning and, particularly, the effectiveness of the component-based textual reasoning approach for this task when perception is adequate.

6.2 Performance on Bongard HOI

On Bongard-HOI, we evaluated GPT-40 and Gemini 2.0 across the four standard test splits (sosa, soua, uosa, uoua; N=100 each from balanced sampling). The results, shown in Table 2, largely replicated the trends observed on Bongard-OW. Performance systematically improved with increased paradigm structure (DVRL < DRL < CA) for both models. Our Compo-nential Analysis (CA) paradigm, particularly with GPT-40 as the reasoning engine, achieved an average accuracy of 77.3% across the splits (with individual splits like sosa and soua reaching 83%), establishing a new state-of-the-art for VLM-based approaches on this benchmark. This surpasses prior SOTA Raghuraman et al. (2024) results from non-VLM spe-cialized methods, such as the reported 76.4% average from a CLIP fine-tuned via PMF approach Raghuraman et al. (2024). Gemini 2.0 with CA also demonstrated strong competitive per-formance with an average of 74.5%. This consistency validates

Model		Avg
Gemini 2.0	DVRL DRL CA	50.8 61.3 74.5
GPT-40	DVRL DRL CA	68.5 71.8 77.3
Prior SOTA	PMF	76.4
Human Avg.	_	91.4

Table 2: Average performance (%) on Bongard-HOI four test-splits across three paradigms.

uation across different complex natural image reasoning datasets. Notably, overall model performance on HOI is lower than on OpenWorld, and a significant gap remains to the high human average scores (avg. 91.4% Jiang et al. (2022)), suggesting HOI's unique challenges in discerning subtle interaction-based rules.

6.3 PERFORMANCE ON WINOGROUND WITH STATIC CA

our framework's applicability and the benefit of structured eval-

We applied our static CA paradigm to Winoground, generating task-agnostic descriptions and then using an LLM for matching. As shown in Table 3, this approach achieves new state-of-the-art results across all three metrics. Using GPT-40 as the reasoning engine yields Text: 75.5%, Image: 58.5%, and Group: 52.0%, significantly surpassing prior SOTA. This success demonstrates that our task-agnostic, decoupled strategy is highly effective not just for rule-discovery, but also for fine-grained compositional reasoning.

6.4 Interactive CA on Winoground: Mitigating Perceptual Gaps

While CA is powerful, we hypothesized its static, one-pass descriptions might miss the single, subtle visual detail that differentiates Winoground pairs. To address this, we applied our Interactive CA (ICA) paradigm, allowing the reasoner to "look again" by asking a targeted question to the perception module.

The results, presented in Table 3, show a significant and consistent performance uplift over the static CA approach for both models tested.

Model	Text Score	Image Score	Group Score
GPT-40 + CA Gemini 2.0 + CA	75.50 71.00	58.50 48.75	52.00 42.00
GPT-40 + ICA	78.00	62.75	55.25
Gemini 2.0 + ICA	72.50	55.25	46.75
Llama3.3-70B + CA	68.25	49.25	41.75
Qwen2.5-32B + CA Phi-4-14B + CA	67.00 65.25	46.25 46.00	40.00 37.75
Qwen2.5-14B + CA	59.25	34.50	27.25
MMICL + CoCoT (Zhang et al., 2024)	64.25	52.5	50.75

Table 3: State-of-the-art performance on the **Winoground benchmark** achieved using our **Componential Analysis** (**CA**) paradigm. Scores reported are the standard Winoground metrics: *Text Score* (correct caption selection per image description), *Image Score* (correct image selection per caption), and *Group Score* (all selections correct per sample), averaged over 400 samples.

Critically, the largest gains are on the Image and Group scores, which are most sensitive to fine-grained visual details. For instance, Gemini's Image Score improved by a remarkable 6.5 points. This demonstrates that the interactive feedback loop is highly effective at resolving the exact visual ambiguities that Winoground is designed to test. This finding confirms that the perception bottleneck is not immutable; it can be actively mitigated by a dynamic, multi-pass reasoning process that guides perception.

In summary, across diverse reasoning tasks, our Componential Analysis paradigms consistently achieve high and often state-of-the-art performance. The success of the interactive ICA variant further highlights that dynamic, decoupled approaches—where reasoning can actively probe perception—are a powerful and promising direction for building more robust and accurate VLMs.

7 ABLATION STUDIES: ISOLATING PERCEPTION AND REASONING

To further investigate the interplay between visual perception, rule representation, and reasoning, we conducted targeted ablation studies. Both studies presented below serve to underscore the critical role of the initial representation derived from visual input – whether it's applying a rule *to* a perceived query image (Section 7.1) or reasoning *from* perceived context images (Section 7.2).

7.1 RULE APPLICATION FIDELITY

How well can models apply an abstract rule once it's formulated? To isolate rule application from rule extraction, we provided models with high-quality rule summaries (generated by GPT-40) and the query image, tasking them solely with classification based on the given rule. This tests the model's ability to ground the symbolic rule in the visual input of the query image.

Table 4 shows performance for several open-source models under this condition. Models like Pixtral-12B demonstrate relatively strong and balanced rule application. Comparing this to Table 1, the generally higher scores here than in CA (where models generate their own descriptions) support the idea that rule application itself is less challenging for these models than the initial perception/description phase.

Model	Acc
LLaVA-7B + DRL	72.0
Llama-vision-11B + DRL	68.2
Pixtral-12B + DRL	83.8
LLaVA-13B + DRL	70.0
LLaVA-34B + DRL	74.8
Llama-vision-90B + DRL	74.2

Table 4: **Rule Application Accuracy:** Accuracy (Acc) in % when classifying query images based on externally provided rules in DRL paradigm.

7.2 IMPACT OF DESCRIPTION QUALITY ON REASONING

Complementing the previous ablation, we investigated how reasoning performance changes when the initial perceptual stage (description generation) is standardized using a high-fidelity source. We generated descriptions for all context and query images using GPT-40 and then used these descriptions as input to the reasoning stage (Stage 2) of the Componential Analysis paradigm for various target models, including weaker VLMs and even text-only LLMs.

The results in Table 5 were revealing. Providing high-quality descriptions dramatically improved the reasoning accuracy of VLMs that struggled when using their own descriptions. Llama-Vision-11B, for example, improved from 53.4% (Table 1) to 84.17%, and Llama-Vision-90B from 55.1% to 90.98%. This provides strong evidence that the reasoning capabilities of these models are significantly underestimated

Model	Acc
Llava:7b + CA	80.56
Llava: $34b + CA$	81.56
Llama-vision:11b + CA	84.17
Llama-vision:90b + CA	90.98
Deepseek-r1:14b + CA	87.98
Gemma2:27b + CA	88.98
Qwen 2.5:7b + CA	90.38
Phi4:14b + CA	91.98
Qwen2.5:32b + CA	92.79
Qwen 2.5:14b + CA	92.99

Table 5: **Impact of High-Quality Descriptions**: Accuracy (Acc) in % using Componential Analysis (Stage 2 reasoning) with image descriptions generated externally by GPT-4o. Includes VLMs and LLMs models respectively separated by line.

by end-to-end evaluations; their primary limi-

tation lies in generating accurate perceptual representations. Further illustrating this sensitivity to description source quality, Table A.6 in the Appendix details a comparison using components generated by Pixtral-12B.

Remarkably, this approach also enabled text-only LLMs to perform the visual reasoning task effectively. Models such as Phi-4 (14B) achieved 91.98% accuracy, while several Qwen models also exceeded 90%. This demonstrates that: (1) High-quality textual descriptions can serve as effective surrogates for visual input, enabling modality transfer for reasoning tasks. (2) The CA paradigm, particularly when coupled with controlled descriptive input, serves as a powerful tool for isolating and evaluating the core symbolic reasoning abilities of both VLMs and LLMs, independent of their integrated perceptual systems. These findings strongly reinforce the conclusion that improving visual perception is paramount for enhancing end-to-end visual reasoning in many current models.

7.3 Additional Analysis

Semantic Similarity Analysis during DRL (Table A.3) confirmed that derived rules generally aligned well with query descriptions, particularly for positive samples. The relatively high similarity for negative samples highlights the challenge of the dataset's near-miss counterexamples.

Qualitative Error Analysis Examining samples misclassified by both top models (GPT-40, Gemini 2.0) under CA revealed recurring error patterns (details in Appendix A.7.7, Table A.9). Frequent issues involved over-generalizing rules, missing critical objects/properties present in positive examples, focusing on spurious correlations, or failing to consistently apply derived rules. These qualitative examples underscore that even highly capable models exhibit fragility in nuanced visual detail processing and robust symbolic rule manipulation.

8 DISCUSSION

This research leveraged a cognitively-inspired framework to dissect the mechanisms of visual reasoning in VLMs. By evaluating performance across paradigms mirroring human cognitive strategies (holistic-*DVRL*, deductive-*DRL*, analytical-*CA*), we moved beyond aggregate scores to probe how these models process complex information. A central finding emerges: a critical perception bottleneck limits many contemporary VLMs, masking strong downstream reasoning capabilities. While advanced models (GPT-40, Gemini 2.0) excel when reasoning over high-fidelity textual descriptions, many open-source models falter, pointing to failures in reliably extracting relevant visual information.

The success of the CA paradigm is particularly insightful. Its strength lies in effectively decoupling perception from reasoning via task-agnostic descriptions. Unlike context-dependent CoT approaches, CA first builds a comprehensive, independent textual "world model" of each image. This allows powerful LLMs to apply their sophisticated reasoning abilities on a clean, symbolic representation. The resulting robust performance across different reasoning types (abstraction in BPs, compositionality in Winoground) suggests that modular architectures—featuring specialized perception modules that output rich symbolic data for general reasoning engines—are a highly promising direction.

The success of our Interactive CA (ICA) paradigm takes this insight a crucial step further. While static CA bypasses the bottleneck, ICA demonstrates how to actively mitigate it. On a fine-grained task like Winoground, where a single detail is critical, the ability for the reasoner to formulate a targeted question and prompt a "second look" from the perception module led to significant performance gains (Table 3). This result is pivotal: it suggests the perception bottleneck is not an immutable property but often a failure of one-pass processing. By modeling a more realistic, iterative process of verification—akin to human visual scanning—ICA shows that the connection between perception and reasoning modules should not be a one-way street, but a dynamic, bidirectional dialogue.

Our ablation studies provide strong converging evidence. Isolating reasoning by providing high-quality external descriptions (Section 7.2) resulted in dramatic performance improvements for bottle-necked models and enabled high performance even from text-only LLMs. This clearly demonstrates that a model's latent reasoning capability often far exceeds what its end-to-end performance suggests.

In conclusion, our cognitively-inspired framework serves a dual purpose. It is both a valuable diagnostic tool for pinpointing the prevalent perception bottleneck and a proof-of-concept for a more powerful architectural class. The success of our componential, and especially our interactive, paradigms reveals the significant latent reasoning potential within today's models. More importantly, it offers a clear path forward: building modular, interactive systems where reasoning can dynamically guide perception is a key step towards more robust and general visual intelligence.

Limitations Our work has several important limitations that define avenues for future research.

The primary limitation of our componential paradigms (CA and ICA) is their reliance on language as an intermediate representation. Their effectiveness is likely highest on tasks where critical visual properties are readily verbalizable. For challenges that hinge on non-verbalizable or geometric reasoning—such as the fine-grained correspondence tasks in benchmarks like BLINK (Fu et al., 2024)—the utility of a purely text-mediated approach may be reduced. While our interactive ICA shows that a dynamic dialogue can resolve ambiguities missed by static descriptions, its scope is still bounded by what can be effectively queried and described in text.

Second, while our study covers abstract, compositional, and interactive reasoning, the framework's applicability to other complex visual domains, such as scientific chart interpretation or mathematical reasoning, requires further investigation.

Finally, we acknowledge several practical scope limitations. This work did not conduct a systematic analysis of prompt sensitivity, a known factor in VLM performance. A deeper investigation into the computational costs and latency trade-offs of our multi-stage paradigms, especially the interactive ICA, is also warranted for practical application.

9 Conclusion

This paper introduced a cognitively-inspired framework to dissect the perception-reasoning interface in VLMs. Through four distinct paradigms—including our novel Interactive Componential Analysis (ICA)—we systematically analyzed VLM problem-solving strategies, revealing two key insights. First, our diagnostic approach confirms that a critical perception bottleneck limits many contemporary VLMs, masking significant latent reasoning abilities.

Second, and more importantly, we demonstrate a powerful architectural solution. Our componential paradigms, which decouple perception from reasoning via task-agnostic textual descriptions, achieve highly competitive performance across diverse benchmarks testing abstraction (Bongard-OW), interaction (Bongard-HOI), and compositionality (Winoground). The success of the interactive ICA paradigm, which allows the reasoning module to actively probe and guide perception, is particularly significant. It shows that the perception bottleneck is not an immutable barrier but can be dynamically mitigated.

Ultimately, our work suggests that the path toward more robust visual intelligence lies not just in scaling monolithic models, but in developing modular, interactive architectures. By providing both a diagnostic toolkit and a proof-of-concept for this interactive approach, we offer a blueprint for a new class of systems capable of more deliberate, verifiable, and human-like visual reasoning.

REFERENCES

Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.

Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.

- Alan Baddeley. Working memory: Theories, models, and controversies. *Annual Review of Psychology*, 63(1):1–29, January 2012. ISSN 1545-2085. doi: 10.1146/annurev-psych-120710-100422. URL http://dx.doi.org/10.1146/annurev-psych-120710-100422.
 - Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate, 2016. URL https://arxiv.org/abs/1409.0473.
 - I. Biederman. Recognition-by-components: A theory of human image understanding. *Psychological Review*, 94:115–147, 1987.
 - MM Bongard. The recognition problem. tech. rep. 1968.
 - Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024.
 - Mark A. Gluck, Russell A. Poldrack, and Szabolcs Kéri. The cognitive neuroscience of category learning. *Neuroscience & Biobehavioral Reviews*, 32(2):193–196, 2008. ISSN 0149-7634. doi: https://doi.org/10.1016/j.neubiorev.2007.11.002. URL https://www.sciencedirect.com/science/article/pii/S0149763407001364. The Cognitive Neuroscience of Category Learning.
 - Google. Gemini 2.0 flash-exp. Large language model, 2024. URL https://gemini.google.com/. Accessed on 2024-08-06.
 - Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
 - Huaizu Jiang, Xiaojian Ma, Weili Nie, Zhiding Yu, Yuke Zhu, and Anima Anandkumar. Bongard-hoi: Benchmarking few-shot visual reasoning for human-object interactions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19056–19065, 2022.
 - Maithilee Kunda. Ai, visual imagery, and a case study on the challenges posed by human intelligence tests. *Proceedings of the National Academy of Sciences*, 117(47):29390–29397, 2020.
 - Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017.
 - Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Advances in Neural Information Processing Systems*, 36, 2024.
 - Fei Fei Li, Rufin VanRullen, Christof Koch, and Pietro Perona. Rapid natural scene categorization in the near absence of attention. *Proceedings of the National Academy of Sciences*, 99(14): 9596–9601, June 2002. ISSN 1091-6490. doi: 10.1073/pnas.092277599. URL http://dx.doi.org/10.1073/pnas.092277599.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
 - AI Meta. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI Blog. Retrieved December*, 20:2024, 2024.
- Allen Newell, Herbert Alexander Simon, et al. *Human problem solving*, volume 104. Prentice-hall Englewood Cliffs, NJ, 1972.
 - OpenAI. Chatgpt 4o. Large language model, 2024. URL https://chat.openai.com/chat. Accessed version 2024-08-06.

- Yuxuan Qiao, Haodong Duan, Xinyu Fang, Junming Yang, Lin Chen, Songyang Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Prism: A framework for decoupling and assessing the capabilities of vlms. *Advances in Neural Information Processing Systems*, 37:111863–111898, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- Nikhil Raghuraman, Adam W Harley, and Leonidas Guibas. Support-set context matters for bongard problems. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=kUuPUIPvJ6.
- Lance J Rips. The psychology of proof: Deductive reasoning in human thinking. Mit Press, 1994.
- Julian Risch, Timo Möller, Julian Gutsch, and Malte Pietsch. Semantic answer similarity for evaluating question answering models. In Adam Fisch, Alon Talmor, Danqi Chen, Eunsol Choi, Minjoon Seo, Patrick Lewis, Robin Jia, and Sewon Min, editors, *Proceedings of the 3rd Workshop on Machine Reading for Question Answering*, pages 149–157, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.mrqa-1.15. URL https://aclanthology.org/2021.mrqa-1.15/.
- Larry R. Squire. Memory and the hippocampus: A synthesis from findings with rats, monkeys, and humans. *Psychological Review*, 99(2):195–231, 1992. ISSN 0033-295X. doi: 10.1037/0033-295x. 99.2.195. URL http://dx.doi.org/10.1037/0033-295X.99.2.195.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248, 2022.
- Mohit Vaishnav and Thomas Serre. GAMR: A guided attention model for (visual) reasoning. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=iLMgk2IGNyv.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- Fei Wang, Xingyu Fu, James Y Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, Wenxuan Zhou, Kai Zhang, et al. Muirbench: A comprehensive benchmark for robust multi-image understanding. *arXiv preprint arXiv:2406.09411*, 2024.
- Rujie Wu, Xiaojian Ma, Zhenliang Zhang, Wei Wang, Qing Li, Song-Chun Zhu, and Yizhou Wang. Bongard-openworld: Few-shot reasoning for free-form visual concepts in the real world. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=hWS4MueyzC.
- XTuner. Llava-llama3: A llava model fine-tuned from llama 3 instruct and clip-vit-large-patch14-336 with sharegpt4v-pt and internvl-sft. https://github.com/XTuner/llava-llama3, 2025.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia,

Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. arXiv preprint arXiv:2412.15115, 2024. Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs. arXiv preprint arXiv:2401.02582, 2024. Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. arXiv preprint arXiv:2302.00923, 2023. Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv preprint arXiv:2303.18223, 2023. Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. Advances in Neural Information Processing Systems, 36:5168–5191, 2023.

APPENDIX

A.1 Broader Relevance

This study offers insights with broader implications for developing more robust and human-like AI systems. Our cognitively-inspired evaluation paradigms provide valuable tools for assessing and understanding the strengths and limitations of Vision-Language Models (VLMs) on complex visual reasoning tasks. The insights gained extend beyond Bongard problems, contributing to the development of VLMs capable of advanced reasoning in real-world applications. Our key finding regarding the visual processing bottleneck in many models has significant implications for future research aimed at bridging the performance gap and unlocking the full potential of accessible models. The demonstration of high performance by advanced VLMs underscores the potential for sophisticated visual understanding, reinforcing the importance of architectures integrating robust perception and reasoning. Finally, our comparative evaluation contributes to discussions about AI accessibility and transparency, identifying specific areas for improvement and paving the way for more reliable AI.

A.2 ATTENTION AND MEMORY IN VISUAL REASONING

While our study primarily focuses on the interplay between perception and reasoning, the roles of attention and memory are also implicitly present in our paradigms. The DVRL paradigm likely engages VLM "visual attention" mechanisms (Bahdanau et al., 2016) to identify salient features across the image set, akin to human holistic processing (Biederman, 1987; Li et al., 2002). DRL relies on the model's ability to "memorize" the extracted rule, involving processes related to working memory (Baddeley, 2012) and internal representation storage (Squire, 1992). Although not directly measured, their involvement is inherent. Future work could explore these aspects more explicitly, perhaps via attention map analysis (Vaswani et al., 2017) or probing memory representations (Vaishnav and Serre, 2023).

A.3 DATASET DETAILS

A.3.1 BONGARD OPENWORLD DATASET

We utilize a subset of 500 test cases from the Bongard OpenWorld dataset (Wu et al., 2024). The full dataset contains 1001 samples, each with 7 positive and 7 negative real-world images distinguished by a "commonsense" rule. Our evaluation set was created by taking the first 250 samples and generating two test cases from each (one positive query, one negative query), resulting in 500 balanced test cases. Specific sample IDs used will be released.

A.3.1.1 COMMONSENSE VALUE CATEGORIES

Table A.1 summarizes the rule categories.

ID	Concept Category	Example
0	Anything else	Animals are running.
1	Human-Object Interaction (HOI)	A person playing the guitar.
2	Taste / Nutrition / Food	A plate of high-calorie food.
3	Color / Material / Shape	A wooden floor in the living room.
4	Functionality / Status / Affordance	An animal capable of flying in the tree.
5	And / Or / Not	A man without beard.
6	Factual Knowledge	A building in US capital.
7	Meta Class	Felidae animals.
8	Relationship	A bench near trees.
9	Unusual Observations	Refraction of light on a glass cup.

Table A.1: Commonsense ID Categories and Examples in Bongard OpenWorld dataset (Wu et al., 2024).

A.3.1.2 COMMONSENSE VALUE DISTRIBUTION IN OUR SUBSET

Table A.2 shows the distribution in our subset. Category '0' is predominant.

ID	Count	Percentage (%)
0	365	73.0
3	36	7.2
9	26	5.2
1	15	3.0
2	14	2.8
4	12	2.4
5	10	2.0
6	10	2.0
8	8	1.6
7	4	0.8
Total	500	100.0

Table A.2: Distribution of Commonsense ID Categories in the 500 Bongard-OW test cases used in our evaluation subset.

A.3.2 BONGARD-HOI DATASET

To assess generalizability on natural images with a different reasoning focus (human-object interactions), we used the Bongard-HOI dataset (Jiang et al., 2022). We evaluated performance on its four standard test splits, defined by object/action novelty:

- sosa: seen object, seen action
- soua: seen object, unseen action
- uosa: unseen object, seen action
- · uoua: unseen object, unseen action

The original splits vary significantly in size and balance (e.g., sosa: 200 pos/200 neg queries; soua: 2236 pos/1348 neg; uosa: 660 pos/660 neg; uoua: 695 pos/695 neg). For consistent cross-split evaluation in this work, we created balanced subsets by sampling 100 test cases from each of the four splits, ensuring an equal distribution of 50 positive and 50 negative query images per split. This resulted in a total evaluation set of 400 samples for Bongard-HOI (100 per split), used for the results reported in Table 2.

A.3.3 WINOGROUND DATASET

To test performance on fine-grained visio-linguistic compositional reasoning, we utilized the Winoground dataset (Thrush et al., 2022). This dataset comprises 400 samples specifically designed to challenge compositional understanding. Each sample contains a pair of minimally contrastive images (I_0, I_1) and a corresponding pair of minimally contrastive captions (C_0, C_1) , requiring models to correctly match image I_0 to caption C_0 and image I_1 to caption C_1 . We used all 400 samples provided in the standard dataset release for our Winoground evaluations reported in Section 6.3 and Table 3.

A.3.4 DATASET AVAILABILITY

Bongard OpenWorld: https://rujiewu.github.io/Bongard-OW.github.io/.

Bongard-HOI: https://github.com/NVlabs/Bongard-HOI/blob/master/assets/dataset.md.

Winoground: https://huggingface.co/datasets/facebook/winoground

Details on the specific subsets and samples used in our evaluations will be released upon publication.

A.4 MODEL AND EXPERIMENT DETAILS

A.4.1 MODEL DETAILS

VLMs: GPT-40; Gemini 2.0; Pixtral-12B; Llama-Vision-3.2 (11B, 90B); LLaVA (Llama-2 based; 7B, 13B, 34B); LLaVA-Llama3-8B. **Text-Only LLMs (for Ablation 7.2):** Phi-4 (14B) (Abdin et al., 2024); Qwen2.5 (7B, 14B, 32B) (Yang et al., 2024); Deepseek-r1 (32B, 70B) (Guo et al., 2025); Gemma2 (27B) (Team et al., 2024).

A.4.2 EXPERIMENT CONFIGURATION

- Access: APIs for closed models; Ollama for open models.
- **Input:** Base64 images in prompts (see Appendix A.5).
- **Image Handling:** API defaults or max 1024px (Ollama). Multi-image calls for DVRL where supported.
- **Decoding:** Temperature 0.
- Fine-tuning: None.
- Hardware: NVIDIA GPUs (2080Ti, 3090, 6000 Ada).

A.4.3 EVALUATION METRICS

- Classification Accuracy: Primary metric (% correct).
- Semantic Similarity: Cosine similarity of OpenAI embeddings ('text-embedding-3-large') between descriptions/rules. Inspired by (Risch et al., 2021).

A.4.4 WINOGROUND SCORE CALCULATION USING COMPONENTIAL ANALYSIS

This section defines the calculation of Winoground (Thrush et al., 2022) scores (text_score, image_score, group_score) within our Componential Analysis (CA) paradigm (Section 6.2).

In the standard Winoground task, a sample i consists of two images $I_{0,i}$, $I_{1,i}$ and two captions $C_{0,i}$, $C_{1,i}$, where $(C_{0,i}, I_{0,i})$ and $(C_{1,i}, I_{1,i})$ are the ground truth correct pairs. Models are typically evaluated based on a scoring function s(C, I) indicating the match between a caption and an image.

In our CA paradigm, the reasoning model (Stage 2) does not access images $I_{0,i}$, $I_{1,i}$ directly. Instead, it operates on textual image descriptions $D_{0,i}$, $D_{1,i}$ generated in Stage 1. The model is prompted to make explicit choices about the best match between descriptions and captions. Let $Choice_C(D_k, \{C_0, C_1\})$ denote the caption $(C_0 \text{ or } C_1)$ chosen by the model as the best match for description D_k . Similarly, let $Choice_D(C_k, \{D_0, D_1\})$ denote the description $(D_0 \text{ or } D_1)$ chosen for caption C_k .

The scores for each sample i in the dataset W (where N = |W| = 400) are calculated as follows:

1. Text Score (f_{CA}): This measures if the correct caption is selected for each image description. We use an indicator function $\mathbb{I}[\cdot]$ which is 1 if the condition inside is true, and 0 otherwise.

$$Choice_{C}(D_{0,i}, \{C_{0,i}, C_{1,i}\}) = C_{0,i}$$

$$f_{CA}(i) = \mathbb{I}\begin{bmatrix} and \\ Choice_{C}(D_{1,i}, \{C_{0,i}, C_{1,i}\}) = C_{1,i} \end{bmatrix}$$
(1)

This score is 1 only if the model correctly identifies the caption for both description $D_{0,i}$ and description $D_{1,i}$.

2. Image Score (g_{CA}) : This measures if the correct image description is selected for each caption.

$$Choice_{D}(C_{0,i}, \{D_{0,i}, D_{1,i}\}) = D_{0,i}$$

$$g_{CA}(i) = \mathbb{I}\begin{bmatrix} and \\ Choice_{D}(C_{1,i}, \{D_{0,i}, D_{1,i}\}) = D_{1,i} \end{bmatrix}$$
(2)

This score is 1 only if the model correctly identifies the description for *both* caption $C_{0,i}$ and caption $C_{1,i}$.

3. Group Score (h_{CA}) : This requires all associations within the sample to be correct.

$$h_{CA}(i) = f_{CA}(i) \land g_{CA}(i) \tag{3}$$

Equivalently, $h_{CA}(i) = 1$ if and only if $f_{CA}(i) = 1$ and $g_{CA}(i) = 1$.

A.5 MODEL PROMPTS

810

811

812

813

814 815

816 817

818 819

820 821

822

823

824

825826827

828

829 830

831

832 833

834

835

836

837

838

839

840

841

843

845

847

848

849 850

851

852

853 854

855

856

861

862

863

A.5.1 DIRECT VISUAL RULE LEARNING

The prompt used for the Direct Visual Rule Learning paradigm is designed to elicit a holistic analysis of the provided images, encouraging the model to identify a distinguishing rule and apply it to the query image. The prompt emphasizes the distinction between positive (cat_2) and negative (cat_1) examples and guides the model to provide a structured output containing its analysis, the identified rule, details about the query image, and the final classification.

```
def visual_concept_test_prompt(m, n):
    Generates a visual analysis prompt.
    Args:
        m (int): Number of positive samples.
        n (int): Number of negative samples.
    Returns:
        str: The formatted prompt string.
    return f"""
    You are provided with \{m + n + 1\} images: the first \{m\} samples are
        'cat_2', the next {n} samples are 'cat_1', and the last image is
       the 'query image'.
    Analyze the common characteristics or patterns found in the 'cat_2'
       samples (positive samples: following 1 common rule) that
       distinctly separate them from the 'cat_1' samples (negative
       samples: it might not follow any possible rule).
    Your task is to:
    1. Determine the rule or criterion that distinguishes the 'cat_2'
       samples from the 'cat_1' ones.
    2. Analyse the 'query image' (last image).
    3. Provide your conclusion for the 'query image' if it can be
       categorized as either 'cat_1' or 'cat_2' based on the analysis
       and the rule.
    Ensure that the output is clear, well-formatted, and free of
       unnecessary explanations.
    Omit the ''' tags at the beginning and end of the page. The format
       of your output should be as follows:
    - **Analysis**: (Your analysis here)
    - **Rule**: (The distinguishing rule here)
    - **Query Image**: (Query image details)
    - **Conclusion**: (cat_1 or cat_2)
```

A.5.2 DEDUCTIVE RULE LEARNING

The Deductive Rule Learning paradigm employs a two-stage prompting strategy. The first stage focuses on rule extraction from positive and negative examples, while the second stage applies the extracted rule to classify a query image. The prompts for each stage are detailed below.

865 866

867

868

870

871

872 873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

290

891

892

893

894

895

896

897

898 899

900

902

903 904 905

906

908

909 910

911

912

913

914

915

916

917

A.5.2.1 FIRST-STAGE PROMPT (RULE EXTRACTION)

This prompt guides the model to identify and summarize a distinguishing rule based on provided positive and negative examples. It emphasizes conciseness in the rule summary.

```
def visual_concept_prompt(m, n):
try:
    if m < 0 or n < 0:
        raise ValueError(f"Invalid input: m and n must be
           non-negative. Received m={m}, n={n}.")
    if m > 0 and n > 0:
        prompt = f"""
            You are provided with {m + n} images: the first {m}
                samples are cat_2, the next {n} samples are cat_1.
                Analyze the common characteristics or patterns found
                in the cat_2 samples (positive samples: following 1
                common rule) that distinctly separate them from the \,
                cat_1 samples (negative samples: it might not follow
                any possible rule).
            Your task is to provide the rules that defines cat_2
                samples. At the end, write "summary" of the rule
                identified in less than 20 words.
            Ensure that the output is clear, well-formatted, and
                free of unnecessary explanations. Omit the ''' tags
                at the beginning and end of the page.
    if n == 0:
        prompt = f"""
            You are provided with {m} images: {m} samples are cat_2.
                Analyze the common characteristics or patterns found
                in the cat_2 samples (positive samples: following 1
                common rule) that distinctly separate them from
                negative samples which might not follow any possible
                rule.
            Your task is to provide the rules that defines cat_2
                samples. At the end, write "summary" of the rule
                identified in less than 20 words.
            Ensure that the output is clear, well-formatted, and
                free of unnecessary explanations. Omit the ''' tags
                at the beginning and end of the page.
    return prompt
except ValueError as e:
    print(f"Error: {e}")
    raise
```

A.5.2.2 SECOND-STAGE PROMPT (RULE APPLICATION)

This prompt presents the previously extracted rule summary and a query image, prompting the model to classify the image based on the rule. It reinforces the Bongard problem context and requests a structured output.

```
# Define the visual analysis prompt
def visual_concept_test_prompt(m, n, summary):
    return f"""
    We are working with Bongard dataset where there are {m} image in the
        cat_2 and {n} images in the cat_1. Summary of the common
        characteristics or patterns found in the cat_2 samples (positive
        samples: following 1 common rule) that distinctly separate them
        from the cat_1 samples (negative samples: it might not follow
        any possible rule) is as follows: \n {summary}.
```

```
918
           Your task is to ponder over the rule and provide your conclusion for
919
               the 'query image' if it can be categorized as either "cat_1" or
920
               "cat 2".
921
           Ensure that the output is clear, well-formatted, and free of
              unnecessary explanations.
923
           Omit the ''' tags at the beginning and end of the page. The format
924
              of your output should be as follows:
925
926
           - **Analysis**: (Your analysis here)
           - **Rule**: (The distinguishing rule here)
927
           - **Query Image**: (Query image details)
928
           - **Conclusion**: (cat_1 or cat_2)
929
930
```

A.5.3 COMPONENTIAL ANALYSIS

Define the visual analysis prompt

931 932

933

934

935

936 937

938 939

940

941

942

943

944

The Componential Analysis paradigm also uses a two-stage prompting strategy. The first stage generates detailed image descriptions, while the second stage derives a rule from these descriptions and applies it to a query image. The specific prompts for each stage are presented below.

A.5.3.1 FIRST-STAGE PROMPT (IMAGE DESCRIPTION GENERATION)

This prompt instructs the model to generate a comprehensive, hierarchical description of a given image in JSON format. It guides the model to cover various aspects of the image, from scene and objects to activities and contextual elements, facilitating detailed comparative analysis in the subsequent stage.

```
def visual_concept_prompt():
945
946
           Generates a visual analysis prompt.
947
           Args:
948
949
           Returns:
950
               str: The formatted prompt string.
951
           return """
952
                   Carefully examine the provided image and identify all
953
                       possible visual elements, organizing them into a
954
                       detailed hierarchical structure. Start with broad
955
                       categories and progress to more specific subcategories.
                       This should cover everything visible in the image,
956
                       ensuring no detail is overlooked. Structure your
957
                       findings in a JSON format to enable easy comparison and
958
                       synthesis of data from other images. This will help
959
                       discern patterns, contexts, and rules valuable for
960
                       identifying or understanding query images.
961
                   Your hierarchy might encompass the following elements:
962
963
                   1. **Scene/Environment**: Description of the overall setting
964
                       depicted, such as urban, natural, indoor, or outdoor
965
                       scenes.
                   2. **Objects**: Define distinct items or entities present in
966
                       the scene.
967
                   - **Living Beings**: Animals, humans, or other biological
968
                       entities.
969
                       - Species or classification (e.g., dog, bird, human).
970
                       - Characteristics (e.g., color, posture, movement).
                   - **Inanimate Objects**: Both synthetic and natural elements.
971
                       - Categories (e.g., vehicle, building, trees).
```

```
972
                       - Properties (e.g., color, size, material, shape).
973
                   3. **Activities**: Observable actions or interactions
974
                       involving any objects or beings.
                   - Specific descriptions of actions (e.g., walking, flying).
975
                   - Participants involved in these actions.
                   4. **Contextual Elements**: Environmental conditions and
977
                       time markers, such as time of day or weather.
978
                    - Detailed characteristics (e.g., cloudy, night, winter).
979
                   5. **Visual Patterns**: Prominent colors, textures, and
980
                       patterns that are visually significant.
                   6. **Emotional Undertones**: Any emotional presence or
981
                       expressions evident in the image.
982
                   7. **Textual Information**: Any visible text within the
983
                       image, including what it says and its visual style.
984
                   8. **Summary**: A concise narrative summarizing the overall
                       content and context of the image.
985
986
                   Ensure that every aspect from the image is represented under
987
                       these categories. The information should be presented in
988
                       the following JSON format:
989
990
                    "Scene": {
991
                        "Description": "..."
992
993
                   "Objects": {
994
                        "Living Beings": [...],
                       "Inanimate Objects": [...]
995
996
                   "Activities": [...],
997
                   "Contextual Elements": {
998
                        "Time of Day": "...",
999
                        "Weather": "..."
1000
                    "Visual Patterns": {
1001
                        "Dominant Colors": [...],
1002
                        "Textures": [...]
1003
1004
                   "Emotional Undertones": "..."
                   "Textual Information": "..."
1005
                   "Summary": "..."
1006
1007
                   Ensure that the JSON output is clear, well-formatted, and
1008
                       free of unnecessary explanations. Omit the '''json tags
1009
                       at the beginning and end of the page.
1010
1011
1012
```

A.5.3.2 SECOND-STAGE PROMPT (RULE DERIVATION INSTRUCTION)

1013

1014

1015

1016

1017 1018

1019

1020

1021

1022

1023

1024

1025

This prompt guides the model to analyze the JSON descriptions generated in the first stage, derive a distinguishing rule, and apply it to classify a query image. It emphasizes the use of the provided JSON format and requests a structured output.

```
def user_eval_prompt(all_image_specs, m, n):
    return f"""
        We are working with the Bongard dataset, which contains {m}
            images in cat_2 (positive samples) and {n} images in cat_1
            (negative samples). These categories are defined as follows:
        - Cat_2: Positive samples that follow a single common rule.
        - Cat_1: Negative samples that may not follow any specific rule.

The image descriptions for the positive samples, negative samples, and the test image are provided in JSON format.
```

```
1026
                   Analyze the common patterns or characteristics in the cat_2
1027
                   samples that distinguish them from cat_1 samples.
1028
               Your task is to:
               1. Derive the rule that defines the cat_2 samples.
1030
               2. Apply this rule to categorize the test image.
1031
1032
               Here are the image descriptions:
1033
               ### Positive Samples (cat_2):
1034
               {all_image_specs[:m]}
1035
1036
               ### Negative Samples (cat_1):
1037
               {all_image_specs[m:m+n]}
1038
               ### Test Image:
1039
               {all_image_specs[-1]}
1040
1041
               Provide your output in the following format:
1042
1043
               - **Analysis**: (Your analysis here)
               - **Rule**: (The distinguishing rule here)
1044
               - **Test Image**: (Test image details)
1045
               - **Conclusion**: (cat_1 or cat_2)
1046
1047
```

A.6 PROMPTS FOR WINOGROUND COMPONENTIAL ANALYSIS (STAGE 2 REASONING)

For the Winoground benchmark (Thrush et al., 2022), Stage 2 of our Componential Analysis (CA) paradigm requires a reasoning model to evaluate matches between image descriptions (JSON strings generated in CA Stage 1 from the Winoground images) and the provided captions. The following prompts were used to guide the reasoning LLM in selecting the best match, forming the basis for calculating the Text Score and Image Score components as detailed in Appendix A.4.4. Both prompts instruct the model to perform a systematic, step-by-step comparison and to return its analysis and final categorization in a structured JSON format.

A.6.1 PROMPT FOR TEXT SCORE COMPONENT DECISION

1048 1049

1050

1051 1052

1053

1054

1055

1056

1057

1058

1059 1060

1061

1063

1064

1066

The following Python function defines the prompt presented to the reasoning LLM. Given one image's detailed JSON description and two candidate captions (Caption 0 and Caption 1), the model is tasked to determine which caption has a higher possibility of matching the image description. This process is repeated for the second image description in the Winoground pair to gather the necessary data points for the Text Score.

```
def text_score_prompt(image_description, caption_0, caption_1):
1068
           Generates a prompt for an LLM to determine if an image description
1069
              has a higher possibility
1070
           of matching caption_0 or caption_1 by evaluating each match
1071
              individually and comparing them,
1072
           using a detailed JSON description and commonsense reasoning.
1073
1074
           Args:
               image_description (str): A JSON string description of an image.
1075
               caption_0 (str): The first candidate caption.
1076
               caption_1 (str): The second candidate caption.
1077
1078
           Returns:
               str: The formatted prompt string.
```

```
1080
           prompt = f"""You are provided with a detailed JSON description of a
1081
               single image and two different captions (Caption 0 and Caption
1082
               1). Your task is to evaluate how well the image description
               matches *each* caption individually, determine which caption
1083
               provides a stronger match (higher possibility), and explain why.
1084
               Apply commonsense reasoning where needed.
1085
1086
           **Image Description (JSON):**
1087
           '''ison
1088
           {image_description}'''
1089
           **Caption 0:** "{caption_0}"
1090
           **Caption 1:** "{caption_1}"
1091
1092
           **Instructions:**
           1. **Deconstruct Image Description:** Identify the main entities
1093
               (using 'id's), actions ('Activities'), attributes
1094
               ('characteristics', 'properties'), and relationships ('Spatial
1095
               Relationships') detailed in the JSON description. Use
1096
               commonsense to understand the full context implied by the
1097
               description.
           2. **Evaluate Match with Caption 0:** Systematically check how well
1098
               the key elements identified in Caption 0 (entities, actions,
1099
               attributes, relationships) are supported by the details in the
1100
               'Image Description' JSON.
1101
               * Look for specific 'id's, 'characteristics', 'actor_ids', 
 'target_ids', 'action' descriptions, 'relationship' types,
1102
                   etc., in the JSON that align with Caption 0's elements.
1103
               * Use commonsense reasoning to map JSON details to caption terms
1104
                   (e.g., 'characteristics' like "elderly" might correspond to
1105
                   "old person").
1106
               * Assess the overall strength of the match (e.g., "strong
1107
                   support", "partial support", "weak support",
                   "contradiction"). Note any discrepancies.
1108
           3. **Evaluate Match with Caption 1:** Perform the same systematic
1109
               check and assessment against Caption 1.
1110
               * Look for specific JSON details supporting or contradicting
1111
                   Caption 1's elements.
1112
               * Use commonsense reasoning.
               * Assess the overall strength of the match for Caption 1. Note
1113
                   any discrepancies.
1114
               **Compare Matches and Conclude:** Compare the strength of the
1115
               match assessed for Caption 0 versus Caption 1. Explain *why* the
1116
               image description represents one caption with a higher
1117
               possibility or accuracy than the other. Highlight the specific
               JSON details (or lack thereof) that lead to this conclusion.
1118
               Explicitly mention where commonsense was applied during the
1119
               evaluation or comparison.
1120
               **Categorize:** Assign 'cat_0' if the image description has a
1121
               higher possibility of matching Caption 0, or 'cat_1' if it has a
1122
               higher possibility of matching Caption 1.
1123
           Return your response strictly in the following JSON format:
1124
1125
               "analysis": (Your detailed analysis comparing the match strength
1126
                   for each caption against the image description, explaining
1127
                   why one is a better fit, and noting the use of commonsense),
               "category": ('cat_0' or 'cat_1')
1128
           }}
1129
1130
           Do not include any text outside of the JSON structure. Your decision
1131
               must be based on evaluating the match between the image
1132
               description and each caption, then comparing those evaluations.
1133
           return prompt
```

```
1134
       A.6.2 PROMPT FOR IMAGE SCORE COMPONENT DECISION
1135
1136
       Similarly, the following Python function defines the prompt used for the Image Score component.
1137
       Given one caption and two candidate image descriptions (Image 0 Description and Image 1 Descrip-
       tion, both JSON strings), the model is tasked to determine which image description has a higher
1138
       possibility of matching the caption. This is repeated for the second caption in the Winoground pair.
1139
1140
       def image_score_prompt(caption, image_0_description,
1141
           image_1_description):
1142
            Generates a prompt for an LLM to determine if a caption has a higher
1143
                possibility
1144
            of matching image_0_description or image_1_description by evaluating
1145
                each match
1146
            individually and comparing them, using detailed JSON descriptions
1147
                and commonsense reasoning.
1148
            Args:
1149
                caption (str): The caption to evaluate.
1150
                image_0_description (str): The JSON string description of the
1151
                    first image.
1152
                image_1_description (str): The JSON string description of the
                    second image.
1153
1154
            Returns:
1155
                str: The formatted prompt string.
1156
            prompt = f"""You are provided with a single caption and detailed
1157
```

image description individually, determine which description
provides a stronger match (higher possibility), and explain why.
Apply commonsense reasoning where needed.

Caption: "{caption}"

Image 0 Description (JSON):
'''json
{image_0_description}'''

Your task is to evaluate how well the caption matches *each*

JSON descriptions of two different images (Image 0 and Image 1).

{image_1_description}'''

'''json

Image 1 Description (JSON):

1158

1159

1160

1161

1162

1163 1164

1165

1166

1167

1168

1169

1170 1171

1172

1173

1174

1175

1176

1177

1178 1179

1180

1181 1182

1183

1184

1185

1186

1187

Instructions:
1. **Deconstruct Caption:** Identify the main entities, actions,
 attributes, and relationships mentioned in the caption (e.g.,
 "old person", "kisses", "young person"). Use commonsense to
 understand the full context implied by the caption.

- **Evaluate Match with Image 0:** Systematically check how well the key elements identified in the caption are supported by the details in 'Image 0 Description'.
 - * Look for specific 'id's, 'characteristics', 'actor_ids', 'target_ids', 'action' descriptions, 'relationship' types, etc., in the JSON that align with the caption's elements.
 - * Use commonsense reasoning to map caption terms to JSON details (e.g., "old person" might correspond to 'characteristics' like "elderly").
 - * Assess the overall strength of the match (e.g., "strong support", "partial support", "weak support", "contradiction"). Note any discrepancies.
- **Evaluate Match with Image 1:** Perform the same systematic check and assessment against 'Image 1 Description'.
 - * Look for specific JSON details supporting or contradicting the caption's elements.

```
1188
               * Use commonsense reasoning.
1189
               st Assess the overall strength of the match for Image 1. Note any
1190
                   discrepancies.
              **Compare Matches and Conclude: ** Compare the strength of the
1191
              match assessed for Image 0 versus Image 1. Explain *why* one
1192
              description represents the caption with a higher possibility or
1193
              accuracy than the other. Highlight the specific JSON details (or
1194
              lack thereof) from *both* descriptions that lead to this
1195
              conclusion. Explicitly mention where commonsense was applied
              during the evaluation or comparison.
1196
              **Categorize:** Assign 'cat_0' if the caption has a higher
1197
              possibility of matching Image 0 Description, or 'cat_1' if it
1198
              has a higher possibility of matching Image 1 Description.
1199
          Return your response strictly in the following JSON format:
1201
           {{
               "analysis": (Your detailed analysis comparing the match strength
1202
                   for each description against the caption, explaining why one
1203
                   is a better fit, and noting the use of commonsense),
1204
               "category": ('cat_0' or 'cat_1')
1205
          }}
          Do not include any text outside of the JSON structure. Your decision
1207
              must be based on evaluating the match between the caption and
1208
              each description, then comparing those evaluations.
1209
1210
           return prompt
1211
```

A.7 RESULTS AND EXTENDED ANALYSIS

A.7.1 Performance on Bongard Openworld

	Gemini 2.0		emini 2.0 GPT-4o	
Category	Mean	Std Dev	Mean	Std Dev
Positive	0.915	0.02	0.902	0.02
Negative	0.868	0.02	0.866	0.02

Table A.3: Semantic Similarity (Cosine) between query descriptions and rules derived during Deductive Rule Learning.

A.7.2 Performance on Bongard-HOI

(Refer to Table A.4 in main text)

1212 1213

1214 1215

1222

1223

1224 1225 1226

1227 1228

1229 1230

1231 1232

1233

1234

1236

1237

1239

1240

1241

A.7.3 WINOGROUND PERFORMANCE CONTEXT

To contextualize the performance of our Componential Analysis (CA) paradigm applied to Gemini 2.0 on Winoground (reported in Section 6.2), we also ran evaluations using Gemini Pro Vision with several prompting strategies. Table A.5 shows these comparative results on the 400-sample Winoground set used. While advanced CoT methods like DDCoT and CoCoT improve over the baseline for Gemini Pro Vision, the CA paradigm applied to Gemini 2.0 achieves competitive scores, particularly on the text metric, demonstrating its effectiveness.

A.7.4 COMPARISON OF DESCRIPTION SOURCES (PIXTRAL-12B VS. GPT-40)

The results, detailed in Table A.6, consistently show that using image components described by GPT-40 yielded higher downstream reasoning accuracy compared to using components described

Model	Paradigm	sosa	soua	uosa	uoua	Avg
	DVRL	50	54	49	50	50.8
Gemini 2.0	DRL	63	62	55	65	61.3
	CA	77	74	70	77	74.5
	DVRL	68	75	61	70	68.5
GPT-4o	DRL	73	77	64	73	71.8
	CA	83	83	66	77	77.3
Human Avg.	_	87.2	90.0	93.6	94.9	91.4

Table A.4: Performance (%) on Bongard-HOI splits across paradigms. Human average taken from (Jiang et al., 2022) **Splits:** sosa: seen_obj_seen_act, soua: seen_obj_unseen_act, uosa: unseen_obj_seen_act, uoua: unseen_obj_unseen_act. Human average from cited source.

Model / Strategy	Text	Image	Group
Gemini (Baseline)	30.75	26.00	25.00
Gemini + DDCoT	45.00	25.00	23.75
Gemini + CCoT	22.50	33.00	20.75
Gemini + CoCoT	40.00	32.50	27.75
Gemini 2.0 + CA (Ours)	71.91	48.71	42.01

Table A.5: Performance comparison on Winoground (400 samples). CA refers to our Componential Analysis paradigm. Other results use Gemini Pro Vision with different prompting strategies.

by Pixtral-12B across all tested reasoning models. While both description sources enabled strong performance, the advantage conferred by GPT-4o's descriptions (ranging from approximately 2% to over 11% improvement depending on the reasoning model) further underscores the critical dependence of reasoning outcomes on the fidelity, richness, and potentially the alignment of the initial perceptual descriptions with the concepts required by the reasoning task. This reinforces the significance of the VLM's front-end visual processing and description capabilities as a key factor influencing overall visual reasoning performance.

	Components (%)		
Model	Pixtral-12B	GPT-40	
Deepseek-R1-14B	83.21	87.98	
Llama3.2-vision-90B	89.05	90.98	
Phi-4-14B	86.86	91.98	
Owen2.5-14B	90.51	92.99	
LLaVA-7B	68.61	80.56	
Llama3.2-vision-11B	80.29	84.17	
LLaVA-34B	79.56	81.56	
Phi-3-14B	84.67	86.97	

Table A.6: Performance comparison using Componential Analysis (Stage 2) with image descriptions generated by either Pixtral-12B or GPT-4o. Evaluated across various reasoning models.

A.7.5 COMPONENTIAL ANALYSIS RESULTS BY COMMONSENSE CATEGORY

Analysis of GPT-40 and Gemini 2.0 performance in CA across commonsense categories (Appendix Table A.7) showed generally strong performance, indicating robustness to varied conceptual rules. Minor variations suggested potential differences in handling specific types of context or attributes, possibly reflecting training data nuances.

ID	Concept Category	GPT-40 (%)	Gemini 2.0 (%)
0	Anything else	92.88	94.23
1	Human-Object Interaction	86.67	92.86
	(HOI)		
2	Taste / Nutrition / Food	100.00	85.71
3	Color / Material / Shape	88.89	91.67
4	Functionality / Status /	100.00	100.00
	Affordance		
5	And / Or / Not	90.00	80.00
6	Factual Knowledge	90.00	90.00
7	Meta Class	100.00	100.00
8	Relationship	100.00	100.00
9	Unusual Observations	92.31	92.31

Table A.7: Overall accuracy (%) of GPT-40 and Gemini 2.0 on the Bongard-OW test set using Componential Analysis, broken down by Commonsense ID category. Performance variations highlight differing model strengths on specific concept types.

A.7.6 IMPACT OF COT-LIKE STRUCTURE

(Refer to Table A.8 below)

Prompt Type	Accuracy (%)		
	Overall	neg	pos
Minimal (No CoT)	61.6	39.2	84.0
Structured (CoT-like)	80.0	66.4	93.6

Table A.8: Impact of Structured Prompting on DVRL accuracy (GPT-4o).

A.7.7 DETAILED ERROR ANALYSIS EXAMPLES

(Refer to Table A.9 below)

No.	Test ID	Caption (Rule)	Reason for Error (Based on GPT-4o o/p)	
1	0021_neg_0	Cars on the city streets at night	Weak reasoning (similarity): Rule requires vehicles, test image (painting) lacks them explicitly, though context implies city.	
2	0014_neg_0	A person playing a guitar.	Rule extraction error: Rule too general (e.g., "person with instrument"), misses specific object (guitar) mentioned in analysis.	
3	0033_neg_0	A bicycle is placed in the corner	Rule extraction error: Misses key property (in a corner / specific placement context). Test image (collage) lacks this context.	
4	0037_neg_0	The girl has long and thin braids on her head.	Rule extraction error: Rule too general (e.g., "girl with braids"), misses specific property (long and thin).	
5	0076_pos_0	Various kinds of rings	Rule extraction error: Rule misses specific object (ring), focuses on property (intricate design) absent in query.	
6	0076_neg_0	Various kinds of rings	Rule extraction error: Rule misses specific object (ring), too general.	
Continued on next page				

No.	Test ID	Caption (Rule)	Reason for Error (Based on GPT-4o Output)	
7	0082_neg_0	Live coral on the sea floor.	Weak reasoning (similarity): Rule identifies 'coral', but test image description fails to mention it. Perceptual description error.	
8	0084_neg_0	A wooden fence surrounding a grassy field. Rule extraction error: Rule misses specific object (grass), uses broader term (greenery). Test image has greenery but not clearly grass.		
9	0112_neg_0	A wooden floor in the living room. Rule extraction error: Misses key objects (living room, floor), focuses only on 'wooden' and general 'indoor'.		
10	0117_neg_0	Colorful ribbons. Rule extraction error: Rule too general, misses specific object (ribbons).		
11	0122_neg_0	A satellite view of Earth.	Rule extraction error: Misses specific viewpoint (top-down satellite), uses more general 'aerial'.	
12	0136_pos_0	Spectator seats view in the stadium.	Weak reasoning/Rule Application error: Rule mentions "sports or spectators", query image description lacks both, leading to incorrect negative classification despite being stadium seats.	
13	0213_neg_0	Checkerboard pattern Rule extraction error: Misses specific object context (fabric), although pattern is identified.		
14	0234_neg_0	A beautiful stone sculpture Rule extraction error: Focuses on wrong property ('prominent' obelisk) instead of the intended rule property ('tall' obelisk).		
15	0247_pos_0	Small river filled with reeds	Rule extraction error: Misses key object (reeds), while focusing on negative constraints (no industrial presence) which are weakly present.	

Table A.9: Error Analysis: Examples of Bongard-OW cases misclassified by both GPT-40 and Gemini 2.0 in Componential Analysis. Captions indicate the ground truth rule (Wu et al., 2024). Reasoning based on analyzing GPT-40's generated analysis, rule, and query description.