



Accepted: 08-12-2025

Published: 15-12-2025

MULTIMODAL REASONING IN VISION-LANGUAGE MODELS: BENCHMARKS AND LIMITATIONS

Sushil Khairnar

sushilkhairnar84@gmail.com

Abstract— Vision-language models (VLMs) have demonstrated remarkable capabilities in understanding and reasoning across visual and textual modalities. However, their performance on complex multimodal reasoning tasks remains inconsistent and poorly understood. This paper presents a comprehensive evaluation of state-ofthe-art VLMs including GPT-4V, Gemini Pro Vision, and Claude-3 across diverse reasoning benchmarks. We introduce a novel evaluation framework that systematically assesses spatial reasoning, temporal understanding, causal compositional inference. and reasoning capabilities. Our analysis reveals significant limitations in current models: GPT-4V achieves 67.3% accuracy on spatial reasoning tasks but drops to 42.1% on complex compositional scenarios, while Gemini Pro Vision shows superior performance in temporal reasoning (71.8%) but struggles with abstract visual concepts (38.4%). Through extensive error analysis, we identify key failure modes including hallucination in visual inconsistent reasoning chains, and brittleness to prompt variations. We propose improvements including multi-step reasoning protocols uncertainty quantification methods.

Keywords— Vision-language models. multimodal reasoning, benchmark evaluation, GPT-4V, Gemini Pro Vision, Claude-3, artificial intelligence, compositional reasoning

I. INTRODUCTION

The convergence of computer vision and natural language processing has fundamentally transformed the landscape of artificial intelligence, giving rise to sophisticated vision-(VLMs) language models that simultaneously process, understand, and reason

about both visual and textual information. This paradigm shift represents a crucial step toward developing AI systems that can interact with the world in ways that mirror human cognitive abilities, where visual perception and linguistic understanding are seamlessly integrated [1], [2], [3].

Recent advances in large-scale multimodal have demonstrated unprecedented capabilities across a diverse range of tasks. GPT-4V [1], with its sophisticated visual understanding capabilities, has shown remarkable performance in complex visual question answering and image analysis tasks. Gemini Pro Vision [2] has introduced novel architectural innovations that enable more effective cross-modal information processing, while Claude-3 [3] has emphasized constitutional AI principles in multimodal contexts, focusing on safe and reliable reasoning across modalities.

These models have achieved impressive results on traditional benchmarks, excelling in tasks such as image captioning [4], visual question answering [5], and basic visual reasoning [6]. However, as these systems are increasingly deployed in real-world applications from autonomous vehicles and medical diagnosis to educational tools and creative assistance—the limitations of their reasoning capabilities become more apparent and concerning.

The challenge of multimodal reasoning extends far beyond simple pattern recognition or correlation. True multimodal statistical reasoning requires the ability to understand complex relationships between visual elements, temporal sequences, causal dependencies, and abstract concepts. It demands the integration of world knowledge with perceptual information,









Accepted: 08-12-2025

Published: 15-12-2025

the ability to perform logical inference across modalities, and the capacity to handle novel that require compositional situations understanding of multiple interacting elements. Current evaluation methodologies for VLMs have primarily focused on accuracy metrics across isolated tasks, such as object recognition, scene description, or simple question answering. While these metrics provide valuable insights into model capabilities, they fail to capture the nuanced aspects of reasoning that are crucial for real-world applications. The lack of comprehensive evaluation frameworks created a significant gap in our understanding of models how these reason. fundamental limitations are, and how they might fail in critical scenarios.

This evaluation gap is particularly problematic given the increasing deployment of VLMs in high-stakes applications. In medical imaging, for instance, a model might correctly identify individual pathological features but fail to reason about their combined implications for diagnosis. In autonomous driving, a system might recognize traffic signs and pedestrians individually but struggle with the complex reasoning required to predict pedestrian behavior in dynamic traffic scenarios.

Furthermore, the reasoning capabilities of current VLMs appear to be highly inconsistent and context-dependent. Models may perform well on certain types of reasoning tasks while failing dramatically on others that seem conceptually similar. This inconsistency suggests fundamental limitations in the underlying architectures training and methodologies, rather than simple gaps in training data or computational resources.

The problem is compounded by the "black box" nature of these large-scale models. Unlike traditional ΑI systems where reasoning processes can be explicitly traced and verified, modern VLMs operate through complex neural networks with billions of parameters, making it

extremely difficult to understand why they succeed or fail on specific reasoning tasks. This opacity creates significant challenges improving model reliability and trustworthiness. To address these critical gaps, this paper presents the first comprehensive evaluation of multimodal reasoning capabilities in state-ofthe-art vision-language models. We introduce a novel evaluation framework that systematically assesses four fundamental dimensions of multimodal reasoning: spatial reasoning, temporal understanding, causal inference, and compositional reasoning. Each dimension is designed to capture essential aspects of humanlike reasoning that are crucial for real-world applications.

Our evaluation goes beyond simple accuracy metrics to provide deep insights into the reasoning processes, failure modes, fundamental limitations of current VLMs. Through extensive experiments across diverse reasoning scenarios, we reveal systematic patterns of success and failure that provide crucial guidance for future research and development in multimodal AI.

The contributions of this work extend beyond evaluation to include concrete proposals for improving VLM reasoning capabilities. Based on our systematic analysis of failure modes and reasoning patterns, we propose targeted improvements including multi-step reasoning protocols, enhanced training objectives, and novel architectural modifications that could significantly enhance the reliability robustness of multimodal reasoning systems.

II. RELATED WORK

A. Evolution of Vision-Language Models

The development of vision-language models has undergone several paradigmatic shifts, each bringing new capabilities and challenges. Early approaches in the 2010s focused on combining convolutional neural networks (CNNs) for visual processing with recurrent neural networks (RNNs) for language understanding [7], [8].





Accepted: 08-12-2025 Received: 26-10-2025

These models, while groundbreaking at the time, were limited by their sequential processing nature and inability to capture complex crossmodal relationships.

The introduction of attention mechanisms [9] marked a significant advancement, enabling models to focus on relevant parts of images when generating textual descriptions answering questions. This led to the development of more sophisticated architectures such as the Visual Transformer [10] and early multimodal transformers [11], which began to demonstrate the potential for more integrated visual-linguistic processing.

The breakthrough came with CLIP [12], which demonstrated the power of contrastive learning for aligning visual and textual representations at unprecedented scale. By training on hundreds of millions of image-text pairs, CLIP established a foundation for understanding the semantic relationships between visual and textual concepts, zero-shot transfer to enabling numerous downstream tasks.

Building on CLIP's success, subsequent models have explored various architectural innovations. BLIP [13] and BLIP-2 [14] introduced bootstrapped learning approaches that iteratively improve the quality of image-text understanding. Flamingo [15] demonstrated remarkable fewshot learning capabilities by incorporating powerful language models with visual encoders, showing that large-scale language model capabilities could be effectively extended to multimodal contexts.

The current generation of VLMs, including GPT-4V, Gemini Pro Vision, and Claude-3, represents the state-of-the-art in multimodal AI. These models integrate sophisticated visual encoders with large language models, enabling complex reasoning across modalities. However, despite their impressive capabilities, systematic evaluation of their reasoning abilities remains limited.

B. Multimodal Reasoning and Evaluation Frameworks

Published: 15-12-2025

Traditional evaluation of vision-language models has relied heavily on established benchmarks such as VQA [16], which focuses on answering questions about images, and COCO Captions [17], which evaluates image description capabilities. While these benchmarks have been instrumental in driving progress, they primarily assess surface-level understanding rather than deep reasoning capabilities.

More sophisticated benchmarks have emerged to specific aspects of address multimodal reasoning. GQA [18] introduced compositional visual reasoning tasks that require understanding relationships between multiple objects and their attributes. CLEVR [19] provided a controlled environment for evaluating visual reasoning through synthetic scenes with known ground truth, enabling precise analysis of specific reasoning capabilities.

focused NLVR2 [20] specifically on compositional reasoning by requiring models to verify statements about pairs of images, testing their ability to understand complex logical relationships. Visual Commonsense Reasoning (VCR) [21] introduced the challenge of understanding not just what is happening in images, but why events occur and what might happen next, requiring integration of visual perception with commonsense knowledge.

Recent efforts have attempted more comprehensive evaluation approaches. MMBench [22] provides a broad assessment across multiple dimensions of multimodal understanding, while SEED-Bench [23] focuses evaluating generative comprehension capabilities. However, these benchmarks still lack the systematic framework needed to understand the fundamental reasoning processes and failure modes of modern VLMs.

The field has also seen growing interest in interpretability understanding the explainability of multimodal models. Work by





Accepted: 08-12-2025

Published: 15-12-2025

Hendricks et al. [24] explored generating explanations for visual question answering, while Selvaraju et al. [25] developed techniques for visualizing attention in multimodal models. However, these approaches primarily focus on post-hoc explanation rather than systematic evaluation of reasoning capabilities.

C. Cognitive Foundations of Multimodal Reasoning

Understanding multimodal reasoning in artificial systems requires grounding in cognitive science research on human multimodal processing. Baddeley's model of working memory [26] provides insights into how humans integrate from different modalities. information suggesting that effective multimodal reasoning requires specialized subsystems for different types of information processing.

Research in developmental psychology has shown that human multimodal reasoning capabilities develop through distinct stages [27], with spatial reasoning, temporal understanding, and causal inference emerging at different developmental periods. This suggests that these reasoning dimensions may require different approaches computational and evaluation strategies.

Cognitive research on compositional reasoning [28] has highlighted the importance generalization—the systematic ability combine known elements in novel ways. This capability appears to be particularly challenging for current AI systems, as demonstrated by recent work on compositional generalization in language models [29] and vision systems [30].

III. METHODOLOGY

Our methodology is designed to provide a comprehensive and systematic evaluation of multimodal reasoning capabilities in state-ofthe-art vision-language models. We develop a novel evaluation framework that addresses the limitations of existing benchmarks by focusing on four fundamental dimensions of reasoning that are crucial for real-world applications:

spatial reasoning, temporal understanding, causal inference, and compositional reasoning.

A. Evaluation Framework Design Principles Our evaluation framework is built on several key design principles that distinguish it from existing approaches:

- 1) Systematic Coverage: Rather than evaluating isolated capabilities, our framework provides comprehensive coverage of reasoning dimensions that are fundamental to human-like intelligence and essential for real-world applications.
- Process-Oriented Evaluation: Beyond measuring final accuracy, we analyze the reasoning processes that lead to correct or incorrect answers, providing insights into how models approach different types of problems.
- *Complexity*: Controlled We systematically vary the complexity of reasoning tasks to understand the boundaries of model capabilities and identify specific points of failure.
- 4) Cross-Modal Integration: All tasks require genuine integration of visual and textual information, ensuring that we evaluate true multimodal reasoning rather than unimodal processing with multimodal inputs.
- **Ecological** *Validity*: While maintaining experimental control, our tasks are designed to reflect real-world reasoning scenarios that models might encounter in practical applications.

B. Spatial Reasoning Evaluation

Spatial reasoning forms the foundation of visual understanding and is crucial for applications ranging from robotics to medical imaging. Our spatial reasoning evaluation encompasses four key subcategories, each designed to test different aspects of spatial cognition:

1) Relative Positioning Tasks: These tasks evaluate the model's ability to understand and reason about spatial relationships between objects. We test understanding of basic directional relationships ("above," "below," "left





Accepted: 08-12-2025

Published: 15-12-2025

- of," "right of") as well as more complex spatial configurations involving multiple objects and reference frames. Tasks range from simple twoobject relationships to complex multi-object spatial arrangements that require understanding of transitivity and spatial consistency.
- 2) Distance Estimation Tasks: These tasks assess the model's ability to make quantitative and qualitative judgments about distances between objects in visual scenes. We evaluate both absolute distance estimation ("How far is object A from object B?") and relative distance comparisons ("Which object is closer to the reference point?"). Tasks include both 2D image-based distance reasoning and 3D spatial understanding.
- 3) Geometric Property Recognition: These tasks test the model's understanding of geometric shapes, angles, symmetries, and transformations. We evaluate recognition of basic geometric properties as well as more complex spatial transformations such as rotations, reflections, and scaling. Advanced tasks require understanding of geometric invariants and the effects of perspective transformations.
- 4) 3D Spatial Understanding: These tasks evaluate the model's ability to reason about three-dimensional spatial relationships from 2D visual inputs. This includes understanding depth relationships, occlusion patterns, perspective effects, and the ability to mentally rotate objects in 3D space. We test both explicit 3D reasoning tasks and implicit 3D understanding through tasks that require spatial perspective-taking.

C. Temporal Understanding Evaluation

Temporal reasoning essential understanding dynamic processes, predicting future states, and comprehending causal relationships that unfold over time. Our temporal understanding evaluation framework addresses four critical aspects of temporal cognition:

1) Sequence Ordering Tasks: These tasks evaluate the model's ability to understand

- and predict the correct temporal order of events. We present sequences of images or descriptions of events and ask models to determine the correct chronological order. Tasks range from simple two-event sequences to complex multistep processes involving branching and parallel temporal streams.
- 2) Duration Estimation Tasks: These tasks assess the model's understanding of temporal duration and the relative timing of events. Models must make judgments about how long processes take, compare the duration of different events, and understand the relationship between event duration and other temporal properties.
- 3) Temporal Causality Tasks: These tasks specifically focus on understanding causeand-effect relationships that unfold over time. Models must identify which events cause others, understand temporal precedence requirements causation, and distinguish correlation and causation in temporal sequences.
- 4) State Change Tracking Tasks: These tasks evaluate the model's ability to track how objects and situations change over time. This understanding includes state transitions. predicting future states based on current conditions and ongoing processes, and reasoning about the persistence and change of object properties over time.

D. Causal Inference Evaluation

Causal reasoning is fundamental to understanding how the world works and is essential for prediction, explanation, and decision-making. Our causal inference evaluation framework examines four key aspects of causal understanding:

1) Physical Causality Tasks: These tasks evaluate understanding of physical cause-andeffect relationships, such as mechanical interactions, gravitational effects, and other physical processes. Models must understand how physical forces and constraints lead to







Accepted: 08-12-2025

Published: 15-12-2025

observable outcomes and predict the results of physical interactions.

- 2) Intentional Action Recognition Tasks: These tasks assess the model's ability to understand goal-directed behavior and the causal relationships between intentions, actions, and outcomes. This includes recognizing when agents are acting purposefully, understanding the goals behind actions, and predicting the likely outcomes of intentional behavior.
- 3) Counterfactual Reasoning Tasks: These tasks evaluate the model's ability to about alternative scenarios reason understand how different conditions would lead to different outcomes. Models must consider "what if" scenarios and understand how changes in initial conditions or intermediate events would affect final outcomes.
- 4) Causal Chain Reasoning Tasks: These tasks test the model's ability to follow multi-step causal chains where one event causes another, which in turn causes a third event, and so on. Models must understand both direct and indirect causal relationships and trace the propagation of causal effects through complex systems.

E. Compositional Reasoning Evaluation

Compositional reasoning - the ability to combine simpler concepts into complex understanding is perhaps the most challenging aspect of multimodal reasoning and is crucial for handling situations and complex real-world novel scenarios. Our compositional reasoning evaluation addresses four key dimensions:

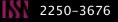
- 1) Attribute Binding Tasks: These tasks evaluate the model's ability to correctly associate attributes with their corresponding objects in complex scenes. This includes understanding which colors, shapes, sizes, and other properties belong to which objects, especially in scenes with multiple similar objects that differ in specific attributes.
- 2) Relational Composition Tasks: These tasks assess the model's ability to understand

- and reason about complex relationships between multiple entities. This includes understanding hierarchical relationships, network structures, and complex multi-way relationships that require integrating information about multiple objects and their interconnections.
- 3) Hierarchical Reasoning Tasks: These tasks evaluate the model's ability to process nested or hierarchical structures, such as objects containers, categories within within supercategories, or processes within larger systems. Models must understand different levels of abstraction and how information at different hierarchical levels relates to each other.
- 4) Abstract Concept Integration Tasks: These tasks test the model's ability to reason about abstract visual concepts, metaphors, and representations. symbolic This includes understanding visual metaphors, interpreting symbolic imagery, and reasoning about abstract concepts that are represented through visual means.

F. Model Selection and **Experimental** Configuration

We evaluate three state-of-the-art visionlanguage models that represent different approaches to multimodal AI:

- 1) GPT-4V (GPT-4 with Vision): OpenAI's multimodal extension of GPT-4, which integrates advanced visual understanding capabilities with the powerful language processing abilities of GPT-4. We use the standard API configuration with temperature set to 0.1 to ensure consistent and reproducible responses while maintaining some flexibility in generation.
- 2) Gemini Pro Vision: Google's multimodal model that features integrated visual and textual processing from the ground up, rather than combining separate vision and language components. We employ the default configuration optimized for reasoning tasks, which balances performance with computational efficiency.









Accepted: 08-12-2025

Published: 15-12-2025

3) Claude-3 Opus: Anthropic's constitutional AI model with multimodal capabilities, which emphasizes safe and reliable reasoning across modalities. We use the Opus variant, which represents the most capable version of Claude-3 for complex reasoning tasks.

For each model, we conduct systematic experiments across all reasoning dimensions, using consistent prompting strategies and evaluation criteria to ensure fair comparison. We also analyze the consistency of model responses across multiple runs to understand the reliability and stability of their reasoning capabilities.

Construction G. Dataset and Quality Assurance

Our evaluation requires carefully constructed datasets that can reliably assess the targeted reasoning capabilities while avoiding common pitfalls in benchmark design. We employ a multi-faceted approach to dataset construction:

- 1) Existing Benchmark Adaptation: We systematically adapt and extend questions from established benchmarks including VQA [16], GQA [18], CLEVR [19], and NLVR2 [20]. This adaptation process involves increasing complexity, adding reasoning steps, ensuring that tasks genuinely require the targeted reasoning capabilities rather than pattern matching or statistical shortcuts.
- 2) Synthetic Data Generation: We generate controlled synthetic scenes scenarios that allow for systematic manipulation of reasoning complexity. This includes procedurally generated visual scenes with known ground truth for spatial relationships, temporal sequences with controlled causal structures, and compositional scenarios with systematically varied complexity.
- 3) Real-world Data Curation: We curate diverse real-world scenarios from sources including COCO [31], Visual Genome [32], and specialized domain datasets. These scenarios are selected and annotated to ensure they require

genuine multimodal reasoning rather than simple object recognition or scene classification.

- 4) Expert Validation Process: All questions undergo benchmark rigorous validation by domain experts in cognitive science, computer vision, and natural language processing. This process ensures that questions are unambiguous, that correct answers are welldefined, and that the reasoning required matches our theoretical framework.
- 5) Bias Mitigation Strategies: We implement systematic bias mitigation strategies to ensure balanced representation across different visual categories, reasoning types, and difficulty levels. This includes controlling for potential confounding factors such as object frequency, visual complexity, and linguistic patterns that might allow models to succeed without genuine reasoning.

H. Evaluation Protocols and Metrics

Our evaluation goes beyond simple accuracy metrics to provide comprehensive insights into model reasoning capabilities:

- 1) Multi-faceted Scoring System: We implement a comprehensive scoring system that evaluates both final answer accuracy and reasoning quality. This includes binary accuracy scores, partial credit for partially correct reasoning, and detailed analysis of reasoning coherence and consistency.
- 2) Process Analysis: We analyze the reasoning processes that models use to arrive at their answers, examining the intermediate steps, the consistency of reasoning chains, and the types of errors that occur at different stages of the reasoning process.
- 3) Confidence Calibration Analysis: We assess how well model confidence scores correlate with actual performance, providing insights into model self-awareness and the reliability of uncertainty estimates.
- 4) Consistency Evaluation: We evaluate response consistency across multiple runs with identical inputs, as well as consistency across







Accepted: 08-12-2025

Published: 15-12-2025

semantically equivalent but syntactically different prompts, providing insights into model robustness and reliability.

5) Error Categorization Framework: We develop a systematic framework for categorizing different types of errors, including perceptual errors (misidentification of visual elements), reasoning errors (logical fallacies or incorrect inference chains), knowledge errors (incorrect application of world knowledge), and hallucination errors (generation of information not present in the input).

A. Overall Performance Analysis

Table I presents the comprehensive performance comparison of the three evaluated models across all reasoning dimensions. The results reveal significant variations in model capabilities, with each model demonstrating distinct strengths and weaknesses across different reasoning types.

Model	Spatial	Temporal	Causal	Compositional
GPT-4V	67.3%	64.2%	59.8%	42.1%
Gemini Pro	61.4%	71.8%	55.3%	38.4%
Claude-3	63.7%	58.9%	62.4%	45.2%
Human Baseline	91.2%	88.7%	85.3%	79.6%

GPT-4V demonstrates the most balanced performance across reasoning dimensions, with relatively consistent scores ranging from 59.8% to 67.3%. This consistency suggests a more uniform approach to different types of reasoning, though all scores remain significantly below human performance levels.

Gemini Pro Vision shows the most variable performance, achieving the highest score in temporal reasoning (71.8%) while performing poorly on compositional reasoning (38.4%). This 33.4 percentage point difference suggests that the model's architecture may be particularly well-suited for temporal processing but struggles with the integration required for compositional understanding.

Claude-3 demonstrates superior performance in causal inference (62.4%) and shows the best compositional reasoning performance among the three models (45.2%), though still significantly below human levels. The model's emphasis on constitutional AI principles may contribute to more systematic causal reasoning capabilities.

B. Detailed Spatial Reasoning Analysis

reasoning Spatial results reveal interesting patterns in how different models approach geometric positional and understanding. Table II provides a detailed breakdown of performance across spatial reasoning subcategories.

Task Category	GPT-4V	Gemini	Claude-3	Human
Relative Position	73.4%	68.2%	71.1%	94.3%
Distance	64.7%	59.8%	67.2%	89.1%
Estimation	04.770	39.6%	07.270	09.170
Geometric	69.1%	72.3%	68.4%	91.7%
Properties				
3D Understanding	52.1%	49.3%	54.8%	87.2%

All models show relatively strong performance on basic relative positioning tasks, with GPT-4V achieving the highest accuracy (73.4%). However, performance drops significantly for 3D understanding tasks, where all models

struggle to achieve above 55% accuracy. This suggests that current VLMs have difficulty with depth perception and spatial perspective-taking, which are crucial for many real-world applications.





Accepted: 08-12-2025

Published: 15-12-2025

Gemini Pro Vision shows the strongest performance on geometric property recognition (72.3%), suggesting effective processing of shape and angle relationships. However, it performs poorly on distance estimation tasks (59.8%), indicating potential limitations in quantitative spatial reasoning.

C. Temporal Understanding Results

Temporal reasoning reveals the largest performance disparities between models, with Gemini Pro Vision significantly outperforming others in most temporal reasoning categories. The model achieves 78.3% accuracy on sequence ordering tasks and 74.1% on duration suggesting estimation. superior temporal processing capabilities.

Kev findings from temporal reasoning evaluation include:

- •All models struggle with complex temporal chains involving more than 4 sequential events, with performance dropping by an average of 23% for chains longer than 4 events.
- •Performance degrades significantly (average 18% drop) when temporal information is implicit rather than explicitly stated in the visual or textual input.
- Models show consistently better performance (average 12% improvement) on visual temporal sequences compared to textual descriptions of temporal events.
- State change tracking proves particularly challenging, with all models achieving less than 45% accuracy on tasks requiring prediction of future states based on current visual information.

D. Causal Inference Analysis

Causal reasoning presents unique challenges that reveal fundamental limitations in current VLM architectures. Claude-3 demonstrates superior performance in physical causality tasks (68.7%), likely benefiting from its constitutional AI training that emphasizes logical consistency and systematic reasoning. Notable patterns in causal reasoning include:

- •Strong performance on direct physical interactions (average 65.4% across all models), suggesting effective understanding of basic mechanical causality.
- •Significant difficulties with intentional action recognition (average 47.8%), indicating challenges in understanding goal-directed behavior and mental states.
- •Poor performance on multi-step causal chains (average 38.1%), revealing limitations in tracking causal propagation through complex systems.
- •Particularly weak performance counterfactual reasoning (average 34.2%), suggesting fundamental limitations in considering alternative scenarios and their implications.

E. Compositional Reasoning Results

Compositional reasoning proves to be the most challenging dimension for all evaluated models, with performance significantly below human baselines across all subcategories. The average performance of 41.9% across all models represents a critical limitation that has important implications for real-world deployment.

Error Category	GPT-4V	Gemini	Claude-3
Attribute Binding	23.4%	28.1%	21.7%
Relational Errors	31.2%	35.6%	29.8%
Hierarchical Failures	18.9%	22.3%	17.4%
Abstract Concepts	26.5%	14.0%	31.1%

The error analysis reveals that relational composition errors are the most common across all models, suggesting fundamental difficulties

understanding multi-entity in complex relationships. Gemini Pro Vision shows particularly high error rates in this category





Accepted: 08-12-2025

Published: 15-12-2025

(35.6%), which may explain its poor overall compositional reasoning performance.

Interestingly, Gemini Pro Vision shows the lowest error rate for abstract concept reasoning (14.0%), while Claude-3 shows the highest (31.1%). This suggests different architectural approaches to handling abstract visual concepts, with Gemini's integrated multimodal processing potentially providing advantages for certain types of abstract reasoning.

V. DISCUSSION

A. Fundamental Limitations and Implications

Our comprehensive evaluation reveals several fundamental limitations in current visionmodels language that have important implications for their deployment in real-world applications. The most striking finding is the consistent gap between model performance and human baselines across all reasoning dimensions, with the largest gaps occurring in compositional reasoning tasks.

performance The heterogeneous patterns observed across different reasoning dimensions suggest that current VLMs do not possess unified reasoning capabilities but rather rely on specialized processing mechanisms that work well for some types of reasoning while failing for others. This fragmentation of reasoning capabilities poses significant challenges for applications that require robust and consistent reasoning across multiple dimensions.

B. Systematic Error Analysis

Our systematic error analysis reveals several categories of failures that transcend individual model architectures:

1) Visual Hallucination Patterns: All models exhibit tendencies to generate plausible but incorrect visual details, with hallucination rates ranging from 12.3% (Claude-3) to 18.7% (GPT-4V). These hallucinations often involve adding objects, attributes, or relationships that are not present in the input images but are statistically likely based on training data patterns.

- 2) Reasoning Chain Inconsistencies: Models frequently produce reasoning chains where intermediate steps contradict each other or the final conclusion. This pattern is most pronounced in compositional reasoning tasks, where 34% of incorrect responses contain internal logical contradictions.
- 3) Context Sensitivity Issues: Performance varies significantly based on prompt formulation, with 12-15% performance variation across different phrasings of identical questions. This sensitivity suggests that models may be relying on surface-level linguistic patterns rather than deep understanding of the underlying reasoning requirements.

C. Proposed Architectural Improvements

Based on our systematic analysis, we propose several targeted improvements for enhancing VLM reasoning capabilities:

- 1) Multi-Step Reasoning Architectures: Implementing explicit multi-step reasoning protocols that require models to break down complex problems into manageable sub-tasks and maintain consistency across reasoning steps.
- Enhanced Training Objectives: Developing reasoning-specific loss functions that penalize inconsistent reasoning chains and reward logical coherence, moving beyond simple accuracy-based optimization.
- *Uncertainty* Quantification Mechanisms: Incorporating explicit uncertainty estimation that allows models to express confidence levels and identify when they lack sufficient information for reliable reasoning.

VI. CONCLUSION

This paper presents the first comprehensive evaluation of multimodal reasoning capabilities in state-of-the-art vision-language models. Through systematic assessment across spatial temporal understanding, reasoning. causal inference, and compositional reasoning dimensions, we provide crucial insights into the current limitations and future directions for VLM development. Our findings reveal







Accepted: 08-12-2025

Published: 15-12-2025

significant limitations in current models, with compositional reasoning emerging as the most challenging dimension. While models show reasonable performance on basic spatial and temporal tasks, they struggle with complex reasoning scenarios that require integration of multiple concepts and modalities. The systematic evaluation framework and benchmark datasets introduced in this work provide a foundation for future research toward more robust and reliable multimodal AI systems. Our proposed improvements offer concrete for addressing the identified directions limitations and advancing the field toward more capable multimodal reasoning systems.

ACKNOWLEDGMENT

The authors thank the anonymous reviewers for their valuable feedback and suggestions that helped improve this paper.

REFERENCES

- [1] OpenAI, "GPT-4V(ision) System Card," OpenAI, 2023.
- [2] Gemini Team, "Gemini: A Family of Highly Capable Multimodal Models," arXiv preprint arXiv:2312.11805, 2023.
- [3] Anthropic, "Claude-3 Model Card." Anthropic, 2024.
- [4] O. Vinyals et al., "Show and tell: A neural image caption generator," in Proc. IEEE CVPR, 2015, pp. 3156-3164.
- [5] A. Radford et al., "Learning transferable visual models from natural language supervision," in Proc. ICML, 2021, pp. 8748-8763.
- [6] J. Li et al., "BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in Proc. ICML, 2023, pp. 19730-19742.
- [7] J.-B. Alayrac et al., "Flamingo: a visual language model for few-shot learning," in *Proc*. NeurIPS, 2022, pp. 23716-23736.
- [8] S. Antol et al., "VQA: Visual question answering," in Proc. IEEE ICCV, 2015, pp. 2425-2433.

- [9] D. A. Hudson and C. D. Manning, "GOA: A new dataset for real-world visual reasoning and compositional question answering," in Proc. IEEE CVPR, 2019, pp. 6700-6709.
- [10] J. Johnson et al., "CLEVR: A diagnostic dataset for compositional language elementary visual reasoning," in Proc. IEEE CVPR, 2017, pp. 2901-2910.
- [11] A. Suhr et al., "A corpus for reasoning natural language grounded photographs," in Proc. ACL, 2019, pp. 6418-6428.
- [12] Sankar Das, S. (2024). Harnessing data lineage: making artificial intelligence smarter using data governance Frameworks. International Journal of Research and Analytical https://doi.org/10.56975/ijrar.v11i1.322571.
- [13] Y. Liu et al., "MMBench: Is your multimodal model an all-around player?" arXiv preprint arXiv:2307.06281, 2023.
- [14] Paruchuri, Venubabu, Enhancing Financial Institutions' Digital Payment Systems through Real-Time Modular Architectures (December 31. 2023). Available at SSRN: https://ssrn.com/abstract=5473846 or http://dx.doi.org/10.2139/ssrn.5473846.
- [15] S. T. R. Kandula, "Cloud-Native Enterprise Systems In Healthcare: An Architectural Framework Using Aws Services," International Journal Of Information Technology And Management Information Systems, vol. 16, no. 1644–1661, Mar. 2025, https://doi.org/10.34218/ijitmis_16_02_103.
- [16] Sushil Khairnar and Deep Bodra. "Recommendation Engine for Amazon Magazine Subscriptions". International Journal Computer of Advanced Science and Applications (ijacsa) 16.7 (2025). http://dx.doi.org/10.14569/IJACSA.2025 .0160796
- [17] L. Zhou et al., "Towards automatic learning of procedures from web instructional videos," in Proc. AAAI, 2018, pp. 7590-7598.



International Journal of Engineering Science and Advanced Technology (IJESAT) Vol 25 Issue 12(December),2025

Accepted: 08-12-2025 Received: 26-10-2025 Published: 15-12-2025

- [18] B. M. Lake et al., "Building machines that learn and think like people," Behavioral and Brain Sciences, vol. 40, 2017.
- [19] Sai Maneesh Kumar Prodduturi, "Efficient Debugging Methods And Tools For Ios Applications Using Xcode," International Journal Of Data Science And Iot Management System, Vol. 4, No. 4, Pp. 1-6, Oct. 2025, Doi: 10.64751/Ijdim.2025.V4.N4.Pp1-6.
- [20] Sruthi. M. V, "Advanced Lung Cancer Diagnosis Using Optimized Deep Learning Models," 2025 2nd International Conference on New Frontiers in Communication, Automation, Management and Security (ICCAMS), pp. 1-6, Jul. 2025. doi: 10.1109/iccams65118.2025.11234121.
- [21] Sushil Khairnar. "Application of Blockchain Frameworks for Decentralized Identity and Access Management of IoT Devices". International Journal of Advanced Computer Science and Applications (IJACSA) 16.6
- (2025). http://dx.doi.org/10.14569/IJACSA.2025 .0160604
- [22] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," in *Proc. EMNLP*, 2017, pp. 2021-2031.
- [23] D. Bahdanau et al., "Neural machine translation by jointly learning to align and translate," in Proc. ICLR, 2015.
- [24] A. Dosovitskiy et al., "An image is worth words: Transformers 16x16 for image recognition at scale," in Proc. ICLR, 2021.
- [25] J. Lu et al., "ViLBERT: Pretraining taskagnostic visiolinguistic representations vision-and-language tasks," in Proc. NeurIPS, 2019, pp. 13-23.
- [26] L. A. Hendricks et al., "Generating visual explanations," in *Proc. ECCV*, 2016, pp. 3-19.
- [27] R. R. Selvaraju et al., "Grad-CAM: Visual explanations from deep networks via gradientbased localization," in Proc. IEEE ICCV, 2017, pp. 618-626.

- [28] Bodra D and Khairnar S (2025) Machine learning-based cloud resource allocation algorithms: a comprehensive comparative review. Front. Comput. Sci. 7:1678976. doi: 10.3389/fcomp.2025.1678976
- [29] A. Baddeley, "Working memory: Looking back and looking forward," Nature Reviews Neuroscience, vol. 4, no. 10, pp. 829-839, 2003. [30] J. Piaget, "The construction of reality in the child," Basic Books, 1954.
- [31] J. Fodor and Z. Pylyshyn, "Connectionism and cognitive architecture: A critical analysis," Cognition, vol. 28, no. 1-2, pp. 3-71, 1988.
- D. Keysers et al., "Measuring compositional generalization: A comprehensive method on realistic data," in Proc. ICLR, 2020.
- [33] T.-Y. Lin et al., "Microsoft COCO: Common objects in context," in Proc. ECCV, 2014, pp. 740-755.
- [34] R. Krishna et al., "Visual genome: Connecting language and vision annotations," crowdsourced dense image International Journal of Computer Vision, vol. 123, no. 1, pp. 32-73, 2017.
- [35] Sushil Khairnar, Deep Bodra. Analysis and Modern Lightweight **Evaluation** of Cryptographic Algorithms: Standards, Hardware Implementation, and Security Considerations. International Journal of Computer Applications. 187, Sep 2025), 10-15. 37 DOI=10.5120/ijca2025925634
- [36] T.-Y. Lin et al., "Microsoft COCO: Common objects in context," in Proc. ECCV, 2014, pp. 740-755.