

# NegVQA: Can Vision Language Models Understand Negation?

**Yuhui Zhang<sup>1</sup>**

**Yuchang Su<sup>2</sup>**  
<sup>1</sup>Stanford University

**Yiming Liu<sup>2</sup>**

**Serena Yeung-Levy<sup>1</sup>**

<sup>2</sup>Tsinghua University

Correspondence: [yuhuiz@stanford.edu](mailto:yuhuiz@stanford.edu)

## Abstract

Negation is a fundamental linguistic phenomenon that can entirely reverse the meaning of a sentence. As vision language models (VLMs) continue to advance and are deployed in high-stakes applications, assessing their ability to comprehend negation becomes essential. To address this, we introduce *NegVQA*, a visual question answering (VQA) benchmark consisting of 7,379 two-choice questions covering diverse negation scenarios and image-question distributions. We construct *NegVQA* by leveraging large language models to generate negated versions of questions from existing VQA datasets. Evaluating 20 state-of-the-art VLMs across seven model families, we find that these models struggle significantly with negation, exhibiting a substantial performance drop compared to their responses to the original questions. Furthermore, we uncover a U-shaped scaling trend, where increasing model size initially degrades performance on *NegVQA* before leading to improvements. Our benchmark reveals critical gaps in VLMs' negation understanding and offers insights into future VLM development. Project page available at <https://yuhui-zh15.github.io/NegVQA/>.

## 1 Introduction

Vision language models (VLMs) such as GPT-4o and Claude have demonstrated remarkable capabilities in understanding and reasoning about visual content through natural language interactions (OpenAI, 2023; Anthropic, 2024). These models can answer image-based questions, generate descriptions, and engage in multi-turn dialogues about visual scenes (Liu et al., 2023; Deitke et al., 2024; Wang et al., 2024b). More recently, they have been integrated into embodied AI systems and robotics, allowing direct interaction with environments and humans in high-stakes scenarios (Driess et al., 2023; Brohan et al., 2023; Kim et al., 2024a).

Despite their impressive progress, VLMs' ability to understand negation (Ackrill et al., 1975)—a fundamental linguistic phenomenon that can completely alter the meaning of a sentence—remains poorly understood. A failure to correctly interpret negation can lead to critical errors, particularly in interactive AI systems. For instance, if a user instructs a VLM not to take a certain action or asks about something that is absent, misunderstanding negation could result in actions contrary to user intent and pose serious safety risks.

To address this, we introduce *NegVQA*, a visual question answering (VQA) benchmark designed to assess VLMs' comprehension of negation. While existing VQA datasets primarily focus on affirmative questions, *NegVQA* systematically examines negation understanding across diverse scenarios. The dataset consists of 7,379 two-choice questions, covering a range of negation types, including cases where objects are absent, attributes such as colors or sizes are negated, actions are described in terms of what is not happening, and more complex forms of negation that require deeper reasoning. To construct *NegVQA*, we leverage large language models to generate natural negations of questions from existing VQA datasets, ensuring fluency while creating challenging evaluation cases that test both linguistic and visual understanding.

We evaluate 20 state-of-the-art VLMs across seven model families and find that negation remains a major challenge. Despite their strong performance on standard VQA tasks, all models struggle significantly when faced with negated questions. For instance, Qwen2-VL-72B (Wang et al., 2024b), the best-performing model, achieves 92.2% accuracy on original questions but drops nearly 20 percentage points to 72.7% on *NegVQA*. Furthermore, we observe a U-shaped scaling trend, where increasing model size initially leads to worse performance on negation before eventually improving. This finding raises important questions about how

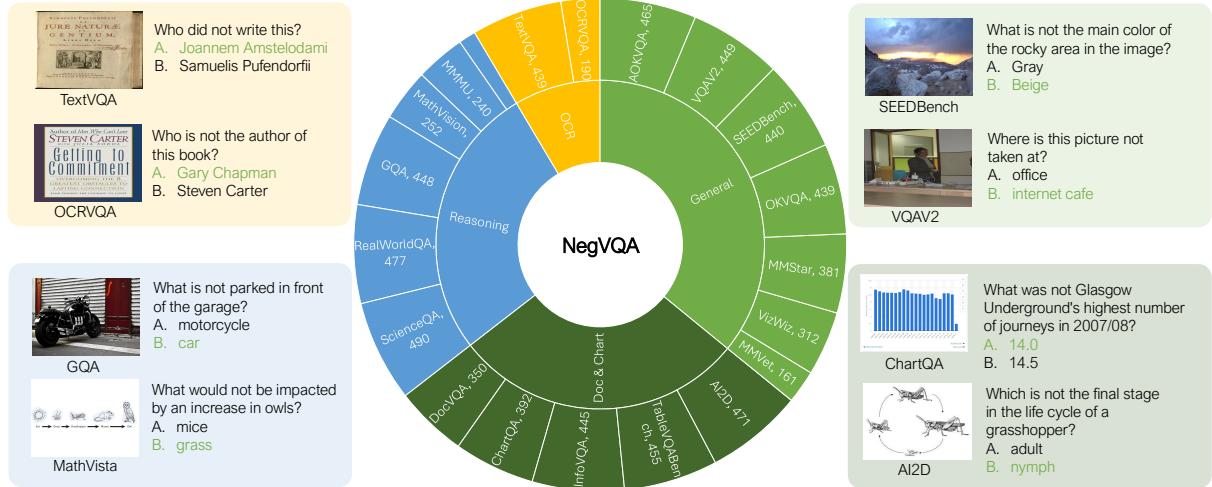


Figure 1: **NegVQA dataset overview.** (Middle) NegVQA comprises a diverse set of negated questions, totaling 7,379 instances sourced from various VQA datasets and domains (general, document/chart, reasoning, OCR). (Left/Right) Example questions from different datasets and domains, with correct answers highlighted in green.

VLMs process negation and how to scale up VLMs to enhance negation understanding abilities.

In summary, we propose *NegVQA*, a critical diagnostic tool for evaluating negation comprehension in VLMs. Our study establishes baseline performance across major VLM families, reveals their significant shortcomings and uncovers scaling behaviors. These insights highlight the need to develop more robust and trustworthy VLMs that can accurately handle negation, a fundamental aspect of natural language understanding.

## 2 Dataset: NegVQA

This section details the construction and statistics of *NegVQA*, our benchmark for evaluating vision language models’ ability to handle negation.

### 2.1 Data Curation

We construct *NegVQA* by systematically transforming questions from VMCBench (Zhang et al., 2025), a multi-choice visual question answering (VQA) benchmark spanning various datasets and domains, into negated versions using GPT-4o (OpenAI, 2023). Our curation process consists of two main steps.

First, we prompt GPT-4o to generate negated versions of the original questions while preserving their syntactic structure and meaning (see Figure 2 for prompt details). For example, the question “Who wrote this book?” is transformed into “Who did not write this book?” We exclude questions that cannot be meaningfully negated (e.g., “Find the value of x.”), as determined by GPT-4o’s assess-

ment of their negotiability. After filtering, 7,379 out of 9,018 questions were identified as negotiable and successfully transformed. To assess the accuracy of GPT-4o’s negation process, we manually verified 100 sampled negated questions and found that 97% were correctly negated—including both the question stems and the two answer choices—confirming the high reliability of the method. Three errors are provided in Appendix Figure 4.

Second, we adjust the answer choices to reflect the negation. Each original four-choice question is reduced to a two-choice format, where we select the correct answer and randomly sample an incorrect choice, then invert their correctness. For instance, in the original question “Who wrote this book?”, if the correct answer is “Samuelis Pufendorfii” and an incorrect choice is “Joannem Amstelodami”, we generate “Who did not write this book?” where “Joannem Amstelodami” becomes the correct answer, and “Samuelis Pufendorfii” the incorrect one. This ensures that the negation meaningfully impacts the answer selection.

### 2.2 Statistics and Examples

*NegVQA* incorporates questions from 20 widely-used VQA datasets within VMCBench, covering a broad range of vision language understanding tasks. It includes datasets for **general VQA capabilities** (VQAv2 (Goyal et al., 2017), OKVQA (Marino et al., 2019), MMVet (Yu et al., 2024), VizWiz (Gurari et al., 2018), A-OKVQA (Schwenk et al., 2022), MMStar (Chen et al., 2024), SEEDBench (Li et al., 2024)),

## GPT-4o Negation Prompt

### Task:

You will be given a question collected from existing visual question answering datasets. Your task is to produce a minimally modified, negated version of the question by inserting a negation (e.g., “not”, “do not”, “isn’t”, etc.) in a way that:

1. **Minimal Changes:** Alters the original question as little as possible.
2. **Answer Inversion:** Causes the original correct answer to become incorrect while making one of the originally incorrect answers correct.
3. **Linguistic Accuracy:** Adheres to proper grammar and preserves the semantic intent of the question.

### Special Case:

1. Do not negate any background that is provided along with the question (e.g., mathematical conditions, background information, etc). Only negate the question itself (usually the last sentence).
2. If it is not possible to create a valid negation that meets these criteria, return an empty string for the negated question and set the flag `is_negatable` to `false`.

### Output Format:

Your response should be an object with the following structure:

```
{  
    "negated_question": "<your negated question (with original background  
                        information) here, or an empty string if not negatable>",  
    "is_negatable": <true/false>  
}
```

Figure 2: Detailed prompts for adding the negation using GPT-4o.

**reasoning tasks** (MathVision (Wang et al., 2024a), GQA (Hudson and Manning, 2019), MMMU (Yue et al., 2024), RealWorldQA (xAI, 2024), MathVista (Lu et al., 2024b), ScienceQA (Lu et al., 2022)), **OCR-based VQA** (OCRVQA (Mishra et al., 2019), TextVQA (Singh et al., 2019)), and **document/chart comprehension** (DocVQA (Mathew et al., 2021), InfoVQA (Mathew et al., 2022), ChartQA (Masry et al., 2022), TableVQABench (Kim et al., 2024b), AI2D (Kembhavi et al., 2016)). The final dataset contains 7,379 questions distributed across these datasets and domains, with the detailed distribution and example questions visualized in Figure 1.

*NegVQA* is designed to systematically test VLMs’ ability to process negation in diverse visual scenarios. The dataset ensures diversity in negation forms, covering cases related to objects, attributes, logical reasoning, spatial relationships, and more. Additionally, all transformed questions have strong visual relevance, requiring models to understand both the image content and the linguistic negation to generate correct answers. *NegVQA* thus serves as a comprehensive benchmark that evaluates vision language models’ ability to under-

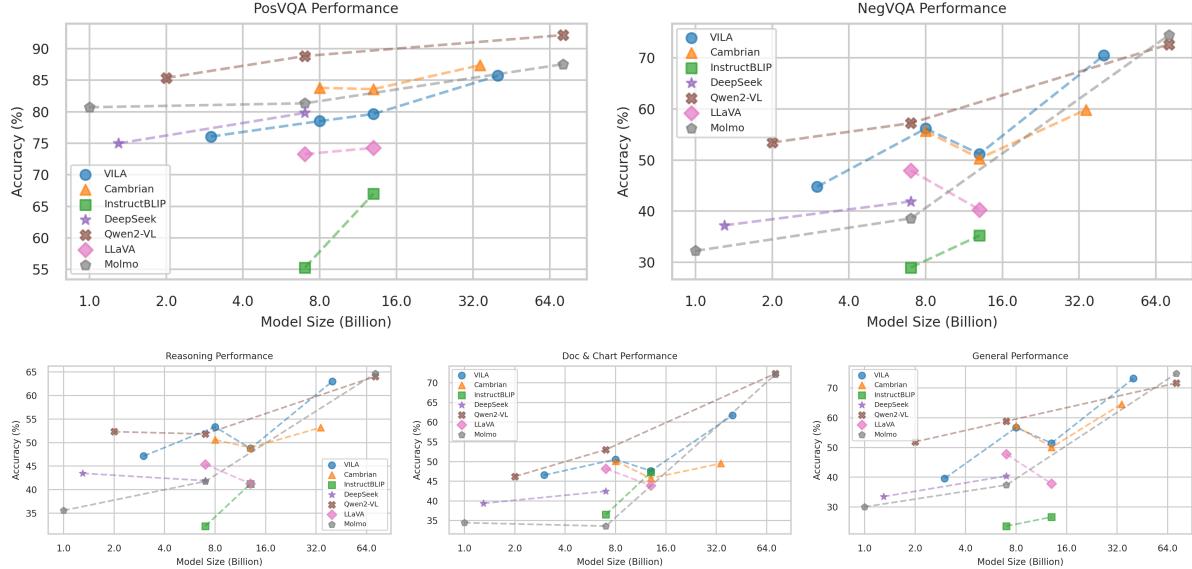
stand negation in different visual scenarios, providing critical insights into their limitations and potential improvements.

## 3 Results

In this section, we describe our experimental setup and present our findings on VLM performance on *NegVQA*. Our evaluation highlights two key insights: current VLMs exhibit significant difficulty in understanding negation, regardless of their size or architecture, and model scaling exhibits a U-shaped performance trend.

### 3.1 Experimental Setup

We selected 20 top-performing vision language models (VLMs) from 7 families based on the OpenVLM Leaderboard (Duan et al., 2024) and evaluated them on *NegVQA*, including Qwen2-VL (Wang et al., 2024b), Molmo (Deitke et al., 2024), Cambrian (Tong et al., 2024), VILA (Lin et al., 2024), DeepSeek-VL (Lu et al., 2024a), LLaVA1.5 (Liu et al., 2023), and InstructBLIP (Dai et al., 2023). These models vary in architecture, training data, and model size, providing a diverse and representative set for evaluation. For each fam-



**Figure 3: Model performance and scaling analysis on *NegVQA* across different VLM families and task categories.** (*Top left*) Performance on the original non-negated two-choice questions shows high accuracy and a positive scaling trend. (*Top right*) Performance on *NegVQA* (negated two-choice questions) is significantly lower, with models exhibiting a U-shaped scaling pattern—initially decreasing before improving as model size increases. (*Bottom*) Category-wise breakdown of *NegVQA* performance (reasoning, document/chart, general), where the U-shaped scaling effect is more pronounced in reasoning and document/chart categories.

ily, we tested multiple model sizes to analyze scaling behavior. All evaluations were conducted in a zero-shot setting using the prompt:

Question: <image> {question}

Options: A. {A} B. {B} C. {C} D. {D}

Answer with the option's letter from the given choices directly.

The results are summarized in Figure 3, with detailed performance provided in Appendix Table 1.

### 3.2 Findings

**VLMs struggle with negation understanding.** Our evaluation reveals that current VLMs consistently underperform on *NegVQA* compared to the corresponding non-negated VQA tasks (which we term *PosVQA*). As shown in Figure 3 (top left vs. top right), performance drops significantly across all model families when answering negated questions. The highest-performing model, Qwen2-VL-72B (Wang et al., 2024b), achieves only 72.7% accuracy on *NegVQA*, compared to 92.2% on non-negated questions—a gap of 19.5 percentage points. On average, model performance decreases by 29.7 points on negated questions relative to their original non-negated counterparts. This substantial decline is observed across different question types and domains, indicating a fundamental limitation in how VLMs process negation.

To contextualize model performance, we added a human baseline for the *NegVQA* benchmark. We manually answered 100 questions and found that humans achieved 89% accuracy. The remaining 11% of errors were due to two factors: 9% required domain-specific knowledge (mostly in subsets like MMMU and ScienceQA), and 2% resulted from conversion errors. This 89% human accuracy, significantly higher than the 72.7% achieved by the state-of-the-art Qwen2-VL-72B, highlights the difficulty of negation understanding for current VLMs and the room for improvement revealed by our benchmark.

The fact that VLMs struggle with negation understanding—evidenced by the performance gaps between negated and non-negated questions, as well as between VLMs and humans—underscores a critical challenge for deploying VLMs in real-world scenarios such as robotics and other embodied environments. Appendix Table 1 provides detailed numerical results.

One potential reason VLMs struggle with negated questions is their limited exposure to negation during training. For example, in the fine-tuning data of a typical VLM like LLaVA (Liu et al., 2023), only 1.1% of conversations contain the word “not.” Enhancing VLMs’ ability to understand negation through training represents a promising direction

for future research. One potential approach is to augment instruction-tuning datasets with carefully curated examples involving negation, thereby guiding models toward a deeper comprehension of such constructs.

**Model exhibits a U-shaped trend scaling.** Intriguingly, a hint of a U-shaped scaling trend (Wei et al., 2022; Zhang et al., 2023) is observed: as models grow larger, their performance on *NegVQA* initially declines before improving at the highest scales. This U-shaped trend is evident in model families such as Cambrian (Tong et al., 2024) and VILA (Lin et al., 2024) (Figure 3, top right), and is especially pronounced in reasoning and document/chart-based tasks (Figure 3, bottom left). Appendix Figure 5 provides a detailed breakdown of performance across individual datasets.

Conceptually, this U-shaped trend can be understood as the composition of two underlying capabilities: original question answering, which tends to improve steadily with model scale, and negation understanding, which follows a tanh-like activation curve. Smaller models with limited reasoning ability often treat negated questions as if they were non-negated, ignoring the negation and selecting answers accordingly. As models scale up, their performance on non-negated questions improves, but their misunderstanding of negation becomes more detrimental, leading to a dip in performance on negated questions. Only when models reach a sufficient level of sophistication to handle negation properly does their performance on negated questions recover, completing the U-shaped trajectory.

Overall, these results underscore the persistent challenges VLMs face in handling negation and highlight the intriguing scaling behavior of VLMs.

## 4 Related Work

**Vision language models (VLMs).** VLMs enable multimodal understanding by modeling  $p(y_t|y_{<t}, x)$  in an auto-regressive manner, where  $y_i$  represents text tokens and  $x$  represents visual input. Modern VLMs typically comprise three key components: a visual encoder (often CLIP (Radford et al., 2021)), a language model, and a linear or MLP projector connecting them. Notable examples include proprietary models such as GPT-4o (OpenAI, 2023) and Claude (Anthropic, 2024), as well as open-source models like LLaVA (Liu et al., 2023) and BLIP (Li et al., 2023). These models are generally trained on image-text pairs and

instruction-tuning datasets, leveraging pre-trained vision and language components. While they exhibit strong performance on various image understanding tasks (Liu et al., 2023; Deitke et al., 2024; Wang et al., 2024b) and have been applied in embodied AI and robotics (Driess et al., 2023; Brohan et al., 2023; Kim et al., 2024a), their ability to handle negation remains largely unexplored.

**Negation understanding.** Negation plays a fundamental role in language comprehension (Ackrill et al., 1975). Most prior research has focused on evaluating language models’ ability to understand negation (Hossain et al., 2020; Fancellu and Webber, 2015; Kassner and Schütze, 2020; Zhang et al., 2023). More recently, studies have begun assessing CLIP (Radford et al., 2021)’s understanding of negation (Alhamoud et al., 2025; Singh et al., 2024; Quantmeyer et al., 2024). However, to the best of our knowledge, no prior work has systematically evaluated negation comprehension in generative VLMs. In this work, we introduce *NegVQA*, the first benchmark designed to assess VLMs’ ability to handle negation. Given the increasing deployment of VLMs in real-world embodied AI systems, understanding their limitations in processing negation is crucial, as failures in user intent interpretation could lead to unintended and risky scenarios.

**Scaling trends.** Scaling up models has been a dominant approach in advancing foundation models. However, most scaling studies have focused on language models (Kaplan et al., 2020; Brown et al., 2020; Ruan et al., 2024). While many tasks benefit from scaling, some exhibit inverse scaling (Lin et al., 2022; McKenzie et al., 2023) or U-shaped scaling (Wei et al., 2022; Zhang et al., 2023). In this work, we analyze scaling effects in vision language models on the negation task and reveal a similar U-shaped scaling pattern.

## 5 Conclusion

In this work, we present *NegVQA*, a benchmark designed to evaluate vision language models’ ability to comprehend negation. Our analysis of 20 VLMs highlights their significant limitations in handling negation and uncovers a U-shaped scaling pattern in performance. We envision *NegVQA* as a valuable resource for advancing linguistically competent, safe, and trustworthy vision language models.

**Acknowledgments.** S.Y. is a Chan Zuckerberg Biohub — San Francisco Investigator.

## Limitations

Our study has three limitations: First, while our multiple-choice format enables controlled experimentation and easy evaluation metrics, it may not fully capture how VLMs handle negation in more open-ended or real-world scenarios where models cannot rely on predefined answer choices. Second, we focus exclusively on zero-shot evaluation, due to current VLMs’ architectural constraint of accepting only single image inputs, leaving unexplored how few-shot prompting might affect negation understanding and performance scaling. Finally, this work primarily investigates how vision-language models (VLMs) handle negation. Enhancing their ability to understand and process negation during training is a promising direction for future research. One potential approach is to augment instruction-tuning datasets with carefully curated examples involving negation, thereby guiding models toward a deeper comprehension of such constructs. Despite these limitations, our work provides the first comprehensive analysis of how VLMs process negation, uncovering both their current limitations and a U-shaped scaling pattern. The *NegVQA* benchmark establishes a foundation for systematically evaluating and improving how future vision language models handle this fundamental linguistic operation.

## References

- John L Ackrill et al. 1975. *Categories and De interpretatione*. Clarendon Press.
- Kumail Alhamoud, Shaden Alshammari, Yonglong Tian, Guohao Li, Philip Torr, Yoon Kim, and Marzyeh Ghassemi. 2025. Vision-language models do not understand negation. *arXiv preprint arXiv:2501.09425*.
- Anthropic. 2024. Introducing the next generation of claudie.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. 2023. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. In *NeurIPS*.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. 2024. Are we on the right way for evaluating large vision-language models? In *NeurIPS*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *NeurIPS*.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. 2023. Palm-e: An embodied multimodal language model. In *ICML*.
- Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, et al. 2024. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. In *MM*.
- Federico Fancellu and Bonnie Webber. 2015. Translating negation: A manual error analysis. In *Proceedings of the Second Workshop on Extra-Propositional Aspects of Meaning in Computational Semantics (Ex-ProM 2015)*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *CVPR*.
- Md Mosharaf Hossain, Venelin Kovatchev, Pranoy Dutta, Tiffany Kao, Elizabeth Wei, and Eduardo Blanco. 2020. An analysis of natural language inference benchmarks through the lens of negation. In *EMNLP*.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. In *ACL*.

- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Min-joon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In *ECCV*.
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. 2024a. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*.
- Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. 2024b. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2024. Seed-bench: Benchmarking multimodal large language models. In *CVPR*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.
- Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024. Vila: On pre-training for visual language models. In *CVPR*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods. In *ACL*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *NeurIPS*.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. 2024a. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024b. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *ICLR*.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *NeurIPS*.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *CVPR*.
- Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *ACL Findings*.
- Minesh Mathew, Viraj Bagal, Ruben Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. 2022. Infographicvqa. In *WACV*.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *WACV*.
- Ian R McKenzie, Alexander Lyzhov, Michael Martin Pieler, Alicia Parrish, Aaron Mueller, Ameya Prabhu, Euan McLean, Xudong Shen, Joe Cavanagh, Andrew George Gritsevskiy, et al. 2023. Inverse scaling: When bigger isn't better. *TMLR*.
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. 2019. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*.
- OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Vincent Quantmeyer, Pablo Mosteiro, and Albert Gatt. 2024. How and where does clip process negation? *arXiv preprint arXiv:2407.10488*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*.
- Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. 2024. Observational scaling laws and the predictability of language model performance. *arXiv preprint arXiv:2405.10938*.
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvqa: A benchmark for visual question answering using world knowledge. In *ECCV*.
- Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *CVPR*.
- Jaisidh Singh, Ishaan Shrivastava, Mayank Vatsa, Richa Singh, and Aparna Bharati. 2024. Learn "no" to say "yes" better: Improving vision-language models via negations. *arXiv preprint arXiv:2403.20312*.
- Shengbang Tong, Ellis L Brown II, Penghao Wu, Sanghyun Woo, ADITHYA JAIRAM IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. In *NeurIPS*.
- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

Jason Wei, Yi Tay, and Quoc V Le. 2022. Inverse scaling can become u-shaped. *arXiv preprint arXiv:2211.02011*.

xAI. 2024. [Realworldqa dataset](#).

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2024. Mm-vet: Evaluating large multimodal models for integrated capabilities. In *ICML*.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *CVPR*.

Yuhui Zhang, Yuchang Su, Yiming Liu, Xiaohan Wang, James Burgess, Elaine Sui, Chenyu Wang, Josiah Akllilu, Alejandro Lozano, Anjiang Wei, et al. 2025. Automated generation of challenging multiple-choice questions for vision language model evaluation. In *CVPR*.

Yuhui Zhang, Michihiro Yasunaga, Zhengping Zhou, Jeff Z HaoChen, James Zou, Percy Liang, and Serena Yeung. 2023. Beyond positive scaling: How negation impacts scaling trends of language models. In *ACL 2023*.

Model	Original Non-negated Questions					Negated Questions (NegVQA)				
	General	Reason	OCR	Doc&Cht	Average	General	Reason	OCR	Doc&Cht	Average
Cambrian-8B	87.6 <sup>88.9</sup> <sub>86.3</sub>	74.0 <sup>75.9</sup> <sub>72.1</sub>	93.4 <sup>95.3</sup> <sub>91.5</sub>	80.9 <sup>82.6</sup> <sub>79.2</sub>	83.8 <sup>84.6</sup> <sub>83.0</sub>	57.2 <sup>59.1</sup> <sub>55.3</sub>	50.6 <sup>52.8</sup> <sub>48.4</sub>	71.8 <sup>75.3</sup> <sub>68.3</sub>	50.1 <sup>52.2</sup> <sub>48.0</sub>	55.7 <sup>56.8</sup> <sub>54.6</sub>
Cambrian-13B	87.7 <sup>89.0</sup> <sub>86.4</sub>	73.5 <sup>75.4</sup> <sub>71.6</sub>	95.9 <sup>97.4</sup> <sub>94.4</sub>	80.5 <sup>82.2</sup> <sub>78.8</sub>	83.6 <sup>84.4</sup> <sub>82.8</sub>	50.1 <sup>52.0</sup> <sub>48.2</sub>	48.9 <sup>51.1</sup> <sub>46.7</sub>	69.2 <sup>72.8</sup> <sub>65.6</sub>	45.7 <sup>47.8</sup> <sub>43.6</sub>	50.3 <sup>51.4</sup> <sub>49.2</sub>
Cambrian-34B	90.0 <sup>91.1</sup> <sub>88.9</sub>	80.3 <sup>82.0</sup> <sub>78.6</sub>	96.6 <sup>98.0</sup> <sub>95.2</sub>	85.0 <sup>86.5</sup> <sub>83.5</sub>	87.4 <sup>88.2</sup> <sub>86.6</sub>	64.6 <sup>66.4</sup> <sub>62.8</sub>	53.2 <sup>55.4</sup> <sub>51.0</sub>	81.5 <sup>84.5</sup> <sub>78.5</sub>	49.5 <sup>51.6</sup> <sub>47.4</sub>	59.9 <sup>61.0</sup> <sub>58.8</sub>
InstructBLIP-7B	58.5 <sup>60.4</sup> <sub>56.6</sub>	53.9 <sup>56.1</sup> <sub>51.7</sub>	70.0 <sup>73.6</sup> <sub>66.4</sub>	48.4 <sup>50.5</sup> <sub>46.3</sub>	55.3 <sup>56.4</sup> <sub>54.2</sub>	23.6 <sup>25.2</sup> <sub>22.0</sub>	32.2 <sup>34.3</sup> <sub>30.1</sub>	20.7 <sup>23.9</sup> <sub>17.5</sub>	36.5 <sup>38.6</sup> <sub>34.4</sub>	28.9 <sup>29.9</sup> <sub>27.9</sub>
InstructBLIP-13B	75.8 <sup>77.4</sup> <sub>74.2</sub>	62.5 <sup>64.6</sup> <sub>60.4</sub>	68.1 <sup>71.7</sup> <sub>64.5</sub>	53.9 <sup>56.0</sup> <sub>51.8</sub>	67.0 <sup>68.1</sup> <sub>65.9</sub>	26.6 <sup>28.3</sup> <sub>24.9</sub>	41.2 <sup>43.4</sup> <sub>39.0</sub>	20.3 <sup>23.4</sup> <sub>17.2</sub>	47.3 <sup>49.4</sup> <sub>45.2</sub>	35.2 <sup>36.3</sup> <sub>34.1</sub>
DeepSeek-VL-1.3B	81.5 <sup>83.0</sup> <sub>80.0</sub>	66.6 <sup>68.7</sup> <sub>64.5</sub>	88.6 <sup>91.1</sup> <sub>86.1</sub>	65.9 <sup>67.9</sup> <sub>63.9</sub>	75.0 <sup>76.0</sup> <sub>74.0</sub>	33.5 <sup>35.3</sup> <sub>31.7</sub>	43.5 <sup>45.7</sup> <sub>41.3</sub>	34.6 <sup>38.3</sup> <sub>30.9</sub>	39.4 <sup>41.5</sup> <sub>37.3</sub>	37.2 <sup>38.3</sup> <sub>36.1</sub>
DeepSeek-VL-7B	84.7 <sup>86.1</sup> <sub>83.3</sub>	71.2 <sup>73.2</sup> <sub>69.2</sub>	91.3 <sup>93.5</sup> <sub>89.1</sub>	73.1 <sup>75.0</sup> <sub>71.2</sub>	79.8 <sup>80.7</sup> <sub>78.9</sub>	40.4 <sup>42.3</sup> <sub>38.5</sub>	41.9 <sup>44.1</sup> <sub>39.7</sub>	53.7 <sup>57.6</sup> <sub>49.8</sub>	42.4 <sup>44.5</sup> <sub>40.3</sub>	41.9 <sup>43.0</sup> <sub>40.8</sub>
LLaVA-1.5-7B	81.0 <sup>82.5</sup> <sub>79.5</sub>	67.7 <sup>69.8</sup> <sub>65.6</sub>	85.5 <sup>88.3</sup> <sub>82.7</sub>	61.1 <sup>63.2</sup> <sub>59.0</sub>	73.3 <sup>74.3</sup> <sub>72.3</sub>	47.7 <sup>49.6</sup> <sub>45.8</sub>	45.4 <sup>47.6</sup> <sub>43.2</sub>	49.7 <sup>53.6</sup> <sub>45.8</sub>	48.2 <sup>50.3</sup> <sub>46.1</sub>	47.9 <sup>49.0</sup> <sub>46.8</sub>
LLaVA-1.5-13B	82.8 <sup>84.2</sup> <sub>81.4</sub>	66.5 <sup>68.6</sup> <sub>64.4</sub>	86.4 <sup>89.1</sup> <sub>83.7</sub>	62.3 <sup>64.4</sup> <sub>60.2</sub>	74.3 <sup>75.3</sup> <sub>73.3</sub>	37.8 <sup>39.6</sup> <sub>36.0</sub>	41.2 <sup>43.4</sup> <sub>39.0</sub>	40.4 <sup>44.2</sup> <sub>36.6</sub>	43.9 <sup>46.0</sup> <sub>41.8</sub>	40.3 <sup>41.4</sup> <sub>39.2</sub>
Molmo-1B	83.6 <sup>85.0</sup> <sub>82.2</sub>	71.7 <sup>73.7</sup> <sub>69.7</sub>	92.0 <sup>94.1</sup> <sub>89.9</sub>	77.7 <sup>79.5</sup> <sub>75.9</sub>	80.7 <sup>81.6</sup> <sub>79.8</sub>	30.0 <sup>31.7</sup> <sub>28.3</sub>	35.6 <sup>37.7</sup> <sub>33.5</sub>	30.4 <sup>34.0</sup> <sub>26.8</sub>	34.5 <sup>36.5</sup> <sub>32.5</sub>	32.2 <sup>33.3</sup> <sub>31.1</sub>
Molmo-7B-O	83.1 <sup>84.5</sup> <sub>81.7</sub>	69.9 <sup>71.9</sup> <sub>67.9</sub>	91.2 <sup>93.4</sup> <sub>89.0</sub>	81.4 <sup>83.1</sup> <sub>79.7</sub>	81.3 <sup>82.2</sup> <sub>80.4</sub>	37.4 <sup>39.2</sup> <sub>35.6</sub>	41.7 <sup>43.9</sup> <sub>39.5</sub>	49.4 <sup>53.3</sup> <sub>45.5</sub>	33.6 <sup>35.6</sup> <sub>31.6</sub>	38.6 <sup>39.7</sup> <sub>37.5</sub>
Molmo-7B-D	85.6 <sup>86.9</sup> <sub>84.3</sub>	67.8 <sup>69.9</sup> <sub>65.7</sub>	94.8 <sup>96.5</sup> <sub>93.1</sub>	84.3 <sup>85.9</sup> <sub>82.7</sub>	83.0 <sup>83.9</sup> <sub>82.1</sub>	55.9 <sup>57.8</sup> <sub>54.0</sub>	48.6 <sup>50.8</sup> <sub>46.4</sub>	75.3 <sup>78.7</sup> <sub>71.9</sub>	49.7 <sup>51.8</sup> <sub>47.6</sub>	55.3 <sup>56.4</sup> <sub>54.2</sub>
Molmo-72B	89.4 <sup>90.6</sup> <sub>88.2</sub>	78.2 <sup>80.0</sup> <sub>76.4</sub>	96.7 <sup>98.1</sup> <sub>95.3</sub>	89.0 <sup>90.3</sup> <sub>87.7</sub>	87.5 <sup>88.3</sup> <sub>86.7</sub>	74.8 <sup>76.5</sup> <sub>73.1</sub>	64.7 <sup>66.8</sup> <sub>62.6</sub>	93.9 <sup>95.8</sup> <sub>92.0</sub>	72.1 <sup>74.0</sup> <sub>70.2</sub>	74.5 <sup>75.5</sup> <sub>73.5</sub>
Qwen2-VL-2B	88.6 <sup>89.8</sup> <sub>87.4</sub>	74.7 <sup>76.6</sup> <sub>72.8</sub>	96.1 <sup>97.6</sup> <sub>94.6</sub>	84.8 <sup>86.3</sup> <sub>83.3</sub>	85.4 <sup>86.2</sup> <sub>84.6</sub>	51.9 <sup>53.8</sup> <sub>50.0</sub>	52.3 <sup>54.5</sup> <sub>50.1</sub>	78.0 <sup>81.2</sup> <sub>74.8</sub>	46.2 <sup>48.3</sup> <sub>44.1</sub>	53.4 <sup>54.5</sup> <sub>52.3</sub>
Qwen2-VL-7B	91.3 <sup>92.4</sup> <sub>90.2</sub>	79.8 <sup>81.6</sup> <sub>78.0</sub>	97.2 <sup>98.5</sup> <sub>95.9</sub>	89.4 <sup>90.7</sup> <sub>88.1</sub>	88.8 <sup>89.5</sup> <sub>88.1</sub>	58.8 <sup>60.7</sup> <sub>56.9</sub>	51.8 <sup>54.0</sup> <sub>49.6</sub>	82.0 <sup>85.0</sup> <sub>79.0</sub>	53.0 <sup>55.1</sup> <sub>50.9</sub>	57.2 <sup>58.3</sup> <sub>56.1</sub>
Qwen2-VL-72B	93.6 <sup>94.5</sup> <sub>92.7</sub>	83.4 <sup>85.0</sup> <sub>81.8</sub>	99.0 <sup>99.8</sup> <sub>98.2</sub>	94.8 <sup>95.7</sup> <sub>93.2</sub>	92.2 <sup>92.8</sup> <sub>91.6</sub>	71.7 <sup>73.4</sup> <sub>70.0</sub>	64.1 <sup>66.2</sup> <sub>62.0</sub>	91.8 <sup>93.9</sup> <sub>89.7</sub>	72.4 <sup>74.3</sup> <sub>70.5</sub>	72.7 <sup>73.7</sup> <sub>71.7</sub>
VILA1.5-3B	83.9 <sup>85.3</sup> <sub>82.5</sub>	68.0 <sup>70.0</sup> <sub>66.0</sub>	88.2 <sup>90.7</sup> <sub>85.7</sub>	66.0 <sup>68.0</sup> <sub>64.0</sub>	76.1 <sup>77.1</sup> <sub>75.1</sub>	39.6 <sup>41.5</sup> <sub>37.7</sub>	47.1 <sup>49.3</sup> <sub>44.9</sub>	51.9 <sup>55.8</sup> <sub>48.0</sub>	46.6 <sup>48.7</sup> <sub>44.5</sub>	44.8 <sup>45.9</sup> <sub>43.7</sub>
VILA1.5-8B	85.3 <sup>86.6</sup> <sub>84.0</sub>	71.2 <sup>73.2</sup> <sub>69.2</sub>	91.0 <sup>93.2</sup> <sub>88.8</sub>	69.4 <sup>71.4</sup> <sub>67.4</sub>	78.5 <sup>79.4</sup> <sub>77.6</sub>	56.7 <sup>58.6</sup> <sub>54.8</sub>	53.3 <sup>55.5</sup> <sub>51.1</sub>	68.4 <sup>72.0</sup> <sub>64.8</sub>	50.5 <sup>52.6</sup> <sub>48.4</sub>	56.2 <sup>57.3</sup> <sub>55.1</sub>
VILA1.5-13B	85.7 <sup>87.0</sup> <sub>84.4</sub>	73.7 <sup>75.6</sup> <sub>71.8</sub>	91.6 <sup>93.8</sup> <sub>89.4</sub>	70.3 <sup>72.2</sup> <sub>68.4</sub>	79.6 <sup>80.5</sup> <sub>78.7</sub>	51.4 <sup>53.3</sup> <sub>49.5</sub>	48.7 <sup>50.9</sup> <sub>46.5</sub>	62.5 <sup>66.3</sup> <sub>58.7</sub>	47.6 <sup>49.7</sup> <sub>45.5</sub>	51.2 <sup>52.3</sup> <sub>50.1</sub>
VILA1.5-40B	89.4 <sup>90.6</sup> <sub>88.2</sub>	78.6 <sup>80.4</sup> <sub>76.8</sub>	96.3 <sup>97.8</sup> <sub>94.8</sub>	81.5 <sup>83.2</sup> <sub>79.8</sub>	85.7 <sup>86.5</sup> <sub>84.9</sub>	73.2 <sup>74.9</sup> <sub>71.5</sub>	63.0 <sup>65.1</sup> <sub>60.9</sub>	90.3 <sup>92.6</sup> <sub>88.0</sub>	61.8 <sup>63.9</sup> <sub>59.7</sub>	70.5 <sup>71.5</sup> <sub>69.5</sub>

Table 1: **Performance of 20 vision language models from 7 families on NegVQA and the original non-negated dataset.** Each reported accuracy is accompanied by a 95% binomial confidence interval, with the lower bound shown as a subscript and the upper bound as a superscript.

### GPT-4o Negation Errors

**Example 1: Improper Negation**

**Original Question:** how many total singles does he have?

**Negated Question:** how many total singles does he **not have**?

**Example 2: Condition Negation Error**

**Original Question:** As shown in the figure, points A, B, and C are three points on O, and the straight line CD and O are tangent to point C. If DCB = 40.0, then the degree of CAB is ()

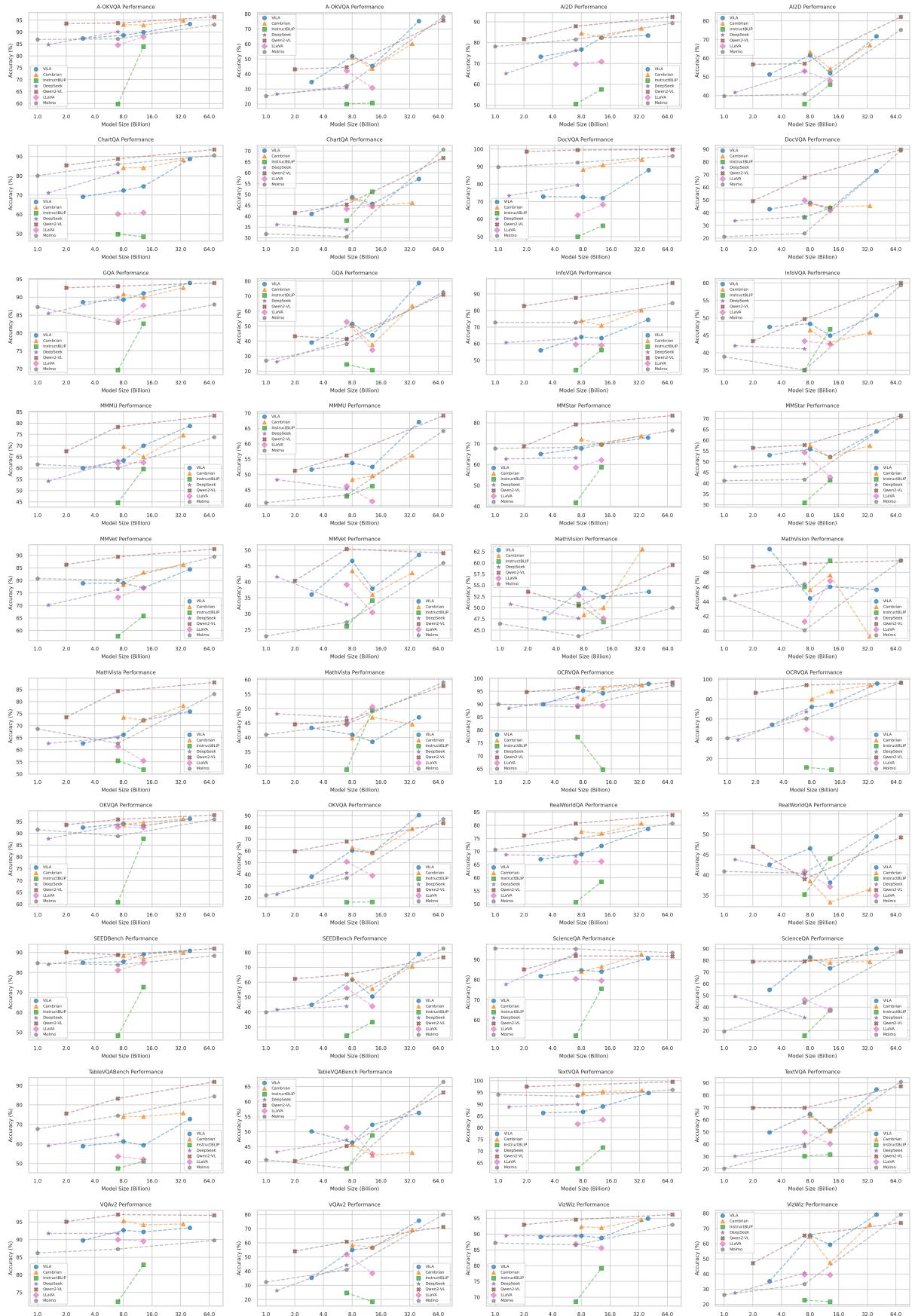
**Negated Question:** As shown in the figure, points A, B, and C are three points on O, and the straight line CD and O are **not tangent** to point C. If DCB = 40.0, then the degree of CAB is ()

**Example 3: Condition Negation Error**

**Original Question:** If cricket was removed from the food web, there would be

**Negated Question:** If cricket was **not removed** from the food web, there would be

Figure 4: **Errors in negated questions generated by GPT-4o.** The first question cannot be negated, while the second and third questions are negated in the condition, whereas the negation should apply to the main question.



**Figure 5: Model performance and scaling analysis on NegVQA across different VLM families and datasets.** For each of the 20 subsets in NegVQA, we present scaling curves for both the original non-negated dataset and the negated dataset from left to right, resulting in a total of 40 figures.