
Understanding Space Is Rocket Science - Only Top Reasoning Models Can Solve Spatial Understanding Tasks

Nils Hoehing
 School of Computer Science
 University College Dublin
 Dublin 4, Ireland
 nils.hohing@ucdconnect.ie

Mayug Maniparambil
 School of Computing
 Dublin City University
 Dublin 9, Ireland
 mayugmaniparambil@gmail.com

Ellen Rushe
 School of Computing
 Dublin City University
 Dublin 9, Ireland
 ellen.rushe@dcu.ie

Noel E. O'Connor
 School of Electronic Engineering
 Dublin City University
 Dublin 9, Ireland
 noel.oconnor@dcu.ie

Anthony Ventresque
 School of Computer Science and Statistics
 Trinity College Dublin
 Dublin 2, Ireland
 anthony.ventresque@tcd.ie

Abstract

We propose RocketScience, an open-source contrastive VLM benchmark that tests for spatial relation understanding. It is comprised of entirely new real-world image-text pairs covering mostly relative spatial understanding and the order of objects. The benchmark is designed to be very easy for humans and hard for the current generation of VLMs, and this is empirically verified. Our results show a striking lack of spatial relation understanding in open source and frontier commercial VLMs and a surprisingly high performance of reasoning models. Additionally, we perform a disentanglement analysis to separate the contributions of object localization and spatial reasoning in chain-of-thought-based models and find that the performance on the benchmark is bottlenecked by spatial reasoning and not object localization capabilities. We release the dataset with a CC-BY-4.0 license and make the evaluation code available at: <https://github.com/nilshoehing/rocketscience>.

1 Introduction

Language models with vision capabilities have enabled powerful applications [35], from turning sketches into website prototypes to recommending recipes based on a photo of fridge contents. At the same time, these models still struggle with fundamental tasks that are often trivial to humans, particularly those that centre around understanding spatial relationships between objects in an image.

Several benchmarks have attempted to measure these shortcomings. However, many suffer from significant limitations: they often recycle existing datasets [17, 26, 38, 39, 44, 55, 62, 68], lack contrastive structure [10, 16, 19, 27, 56], or rely on synthetic or schematic images [3, 30, 31] (see

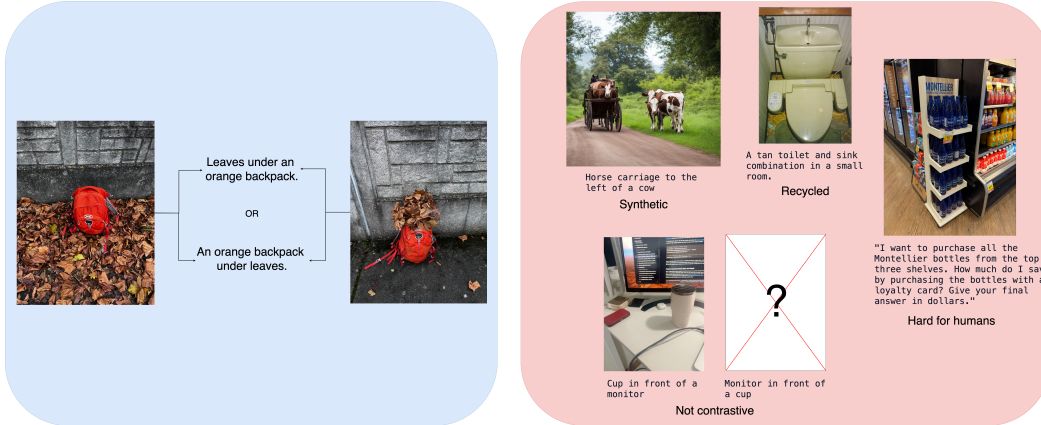


Figure 1: Left: A contrastive pair from our RocketScience benchmark, showing two images and captions differing only in object positions. Right: Examples of problematic data in other benchmarks, such as reused or synthetic data, from [3] (cc-by-sa-4.0), [17] (originally from [37] cc-by-4.0), [53] (MIT) and [56] (apache-2.0)

Figure 1). As a result, these benchmarks tend to overestimate model performance, sometimes even allowing pure language models to score highly on supposedly vision-language tasks [24, 48, 61].

The reuse of older datasets raises concerns about data contamination, as their contents may already be present in the training corpora of contemporary language models. Non-contrastive benchmarks often permit the use of unintended shortcuts, thereby failing to evaluate the specific capabilities they are designed to test. Furthermore, the inclusion of synthetic images in evaluation datasets poses challenges, as performance on such data does not reliably transfer to real-world scenarios. A notable example is the CLEVR dataset [30], where models have achieved near-saturating performance for years, despite continuing to struggle with real-image counterparts.

To address these issues, we introduce *RocketScience*, a benchmark specifically designed to rigorously evaluate spatial understanding in VLMs. The dataset is comprised of 482 manually curated, contrastive image-text pairs representing diverse, real-world scenes (indoors, outdoors, across varying lighting conditions - see Appendix C). Each example forms a question-answer pair that is trivially solvable by humans within seconds, yet proves difficult for current vision-language models.

We evaluate three major categories of models: (1) dual-encoder models such as those in the CLIP family, (2) vanilla multimodal large language models (MLLMs), both open- and closed-source, and (3) advanced reasoning-based MLLMs like o1-mini and Gemini 2.5 Pro. Our results show that all models, except those explicitly designed for multimodal reasoning, perform at chance levels. In contrast, models utilizing chain-of-thought (CoT) prompting or reinforcement learning-based reasoning approaches approach near-perfect performance on this benchmark. We further disentangle model performance along two key axes: entity localization and spatial reasoning. Our analysis reveals that poor performance on spatial understanding tasks stems primarily from limitations in reasoning capabilities, rather than failures in localizing objects.

Our contributions are summarized as follows:

- A new open-source, contrastive benchmark, RocketScience, built entirely from scratch using diverse, real-world (non-synthetic) data, specifically designed to evaluate spatial reasoning capabilities in VLMs.
- An evaluation of three classes of models on the benchmark: CLIP-like models, VLMs and reasoning VLMs
- A disentangled analysis of reasoning-based model performance along two axes: object localization and spatial reasoning. We demonstrate that chain-of-thought reasoning is the primary bottleneck for solving spatial reasoning tasks.

2 Related Work

2.1 VLM Benchmarks and VLMs

Several benchmarks have been proposed to evaluate vision language models in recent years. They span commonsense reasoning [8, 9], understanding multiple images at the same time [67], noticing small differences between images [22], counting objects [47], abductive reasoning [66], visual analogies [71] and diagram understanding [76]. Larger benchmarks like OmniBench [34], MMEvalPro [29], MMStar [12] and WildVision [40] evaluate a wide range of VLM phenomena at the same time. Recently, particularly challenging benchmarks such as ZeroBench [53] and Humanity’s last exam [49] have been released although it remains debatable whether a benchmark which is also hard for humans is even desirable, since the difficulty might be the result of poorly designed questions.

At the same time, vision-language models have also substantially improved in recent years. While traditional VLMs answer questions immediately, recently models have been trained to produce a chain-of-thought (CoT) [70], which is step-by-step reasoning, to improve their performance.

2.2 Contrastive VLM Benchmarks

A common problem with VLM benchmarks is that they can largely be solved using language models alone. [12, 25, 26, 48, 64] This is because certain compositions of objects are more likely to appear in the world and questions about them can therefore be solved with common linguistic co-occurrences instead of visual understanding. To avoid this, contrastive benchmarks [3, 23, 31, 42, 50, 61] are composed of so-called contrastive pairs, which are tuples of two images and two matching texts, that ideally only differ in the exact concept that we want to test for. This means that if the one caption within a contrastive pair contains a particularly likely scenario, like "a person on a horse" the other caption will have to be the unlikely opposite, "a person under a horse". Due to this contrastive design, a surface level statistical linguistic understanding is not enough to solve them and models need to have an integrated visio-linguistic understanding of the world to be able to solve them.

2.3 Spatial Understanding

Spatial understanding concerns reasoning about object positions and is evaluated across different modalities. For VLMs, benchmarks such as [15, 28, 52, 57, 77] test abstract spatial reasoning with simple shapes. Language-only datasets like SpartQA [43] and WorldSense [6] use question answering, while CVR [74] provides a visual outlier-detection task. Text-to-image benchmarks, including DrawBench [54], DALL-Eval [14], and VISOR [20], assess spatial understanding through compositional prompts. Others, such as SpatialRGBT-Bench [13] and Q-Spatial Bench [36], target distance reasoning. Among these diverse notions of spatial understanding, we focus specifically on spatial relations (e.g., on, under), as they form the foundation for more complex reasoning.

2.4 Spatial Relation Benchmarks

To measure progress on the understanding of spatial relations, a range of benchmarks have been developed over the past few years. We provide an extensive overview of them in Table 2 in Appendix A that also includes benchmarks with only a subset dedicated to spatial relations. A large portion of them is non-contrastive, making them potentially sensitive to shortcut solutions. Recycling data from other datasets instead of sourcing new data is common as well, which can lead to unrepresentative scores when this old data leaks into model training, inflating the scores. We will also show empirical evidence for those benchmarks being easier than RocketScience in a later section in Figure 3a. Additionally a few test sets like VisMin [3] are synthetically generated, which means we cannot directly infer models’ competence on real data from their results, especially if the synthetic samples are very schematic. We know that slight changes applied to images can already change predictive models’ outputs [18] [1], so it is prudent not to assume that this domain transfer works automatically. Additionally, we can also still observe many artifacts and unrealistic appearing objects in synthetic data, although this is likely to improve with better image generation models in the future. Some of the older benchmarks are of significant size, which was common at the time, but can actually be expensive to evaluate on today, where most top models are not open source.

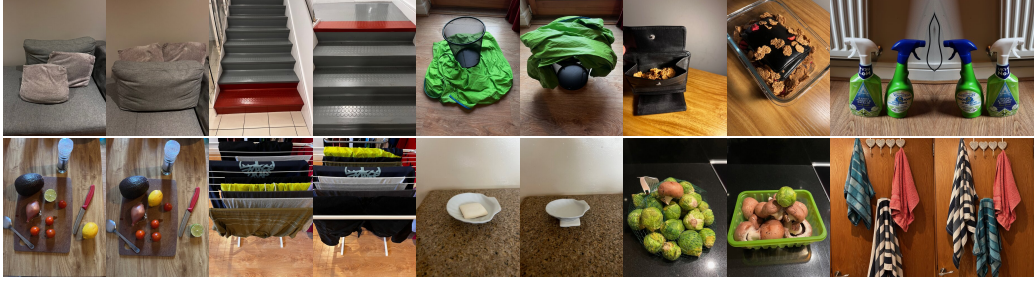


Figure 2: Overview of the contrastive RocketScience dataset.

Winoground [61], Naturalbench [32] and MMVP [63] only contain a small number of spatial understanding questions and no identifiable subset evaluating this property specifically. Therefore they are not included in Table 2. CLEVR [30] is the oldest of the benchmarks. It contains synthetic images of abstract shapes and is very schematic. VALSE [48], Visual Genome Relation [73] and SpatialEval [68] are known to be solvable to a high degree by a blind language model without any image inputs. [48, 64, 68] This likely also applies to many of the other non-contrastive benchmarks. Some of the non-contrastive benchmarks like SugarCrepe [26], ConMe [27] and SugarCrepe++ [17] put a great deal of effort into creating text foils (also known as *hard negatives*) to choose from. NuScenes-SpatialQA [62] is a huge benchmark with automatically generated captions that only covers the self driving car domain. Among the contrastive benchmarks, Rel3D [23] is impressively diverse for a synthetic dataset. FOREST [50] and VSR [38] add perspective change from the perspective of objects, introducing additional complexity.

3 RocketScience Benchmark

3.1 Benchmark Design

Based on the shortcomings outlined in Section 2, we have developed a carefully curated, real-world spatial understanding benchmark that includes a variety of different testing scenarios which avoid predefined schemas and that also follows a contrastive design to avoid shortcut solutions. (see Figure 2) In contrast to other recent benchmarks like Zerobench [53], we create image-text pairs that are understandable to humans while being challenging for machines. We prove that examples are not challenging merely due to them being inherently ambiguous with extremely high human performance. (See Section 4.4)

The main focus of this benchmark is to systematically evaluate relative spatial understanding and the order of objects, as several works have suggested that these properties are not effectively learned by many vision-language models [25, 26, 31, 38, 69]. We design the benchmark with contrastive pairs of new images and texts (each contrastive pair consists of two images and two matching captions with minimal differences). The benchmark requires two levels of understanding: object localization within the image and inference of their spatial relation. The questions are designed to require both steps, without shortcuts, by including the same objects in both parts of each contrastive pair.

3.2 Data Collection

We approximately balance the five main categories which all test for relative spatial understanding: horizontal position, vertical position, depth, proximity and order. The distribution of categories can be seen in Figure 8.

All images were collected in Europe and the USA with an iPhone 13 Mini. We intentionally exclude people and personally identifiable information for the purposes of data privacy. After collection, one author labeled the images and two additional authors checked their agreement with the labels, suggesting changes if necessary. We iterated this process until all authors agreed on the labels.

We add *Label* and *Category* tags to each contrastive pair, where the label specifies the spatial relation necessary to distinguish the two samples in the pair (e.g. "left of") and the category indicates the spatial category (e.g. horizontal position). A contrastive pair can have two categories when the

spatial relations are not polar opposites (e.g. "shoe in front of box" vs. "shoe to the right of box" would be assigned the categories `horizontal position` and `depth`).

3.2.1 Images

As some objects are impossible or extremely challenging to move in the real world, many of the pairs contrasting left and right positions consist of an original photo and its mirrored version as its opposite. Mirroring was only applied when no relevant text would be distorted through of mirroring.

To account for variation in lighting, images were captured in different lighting scenarios. These included: natural light outside, natural light inside, nighttime outside and artificial light inside. We note, however, that the lighting conditions are always consistent within a single contrastive pair.

3.2.2 Captions

Captions were designed such that the two captions for a pair of images could only differ in the position of the objects. For example:

- **Word order:** This represents the largest share of the dataset, e.g. "A chair to the left of a table" vs. "A table to the left of a chair"
- **Swapped Preposition:** This is necessary for some relations that cannot be inverted by simply changing the word order, e.g. "The scissors are close to the door" vs. "The scissors are far from the door".

In both cases, the semantics of the two captions are opposites (also termed *hard negatives*) with respect to the objects' positions. These hard negatives are required to assess whether a spatial relation has been successfully learned rather than simply the likelihood of a relation being associated with certain nouns (i.e. people on chairs, not chairs on people). With the exception of the "left" vs. "right" distinction, in many cases, hard negatives control for unlikely noun-relation cases, while a benchmark without hard negatives simply checks whether models have memorized the most likely case.

3.3 Scope and Variety

The benchmark is designed for spatial understanding, but the localization of the objects still requires additional understanding such as: counting, negation, quantifiers ("most of"), materials, size and colour. In comparison to What's Up [31], RocketScience is significantly more diverse and less schematic. We include a wide variety of scene characteristics such as indoor and outdoor environments, daytime and nighttime settings, objects at varying distances (both near and far), a range of object sizes (small and large), as well as natural and rural environments. (See Appendix C) This broad coverage makes our dataset substantially more representative of real-world conditions.

The physical objects found in the images are sourced from a wide range of objects from daily life within Europe and the US. The distribution of objects can be seen in Figure 9 in Appendix D. The most frequent adjectives are colors, but materials like "*wooden*" or "*metal*", in addition to counts and other object properties can be found in the dataset, see Figure 10 which is also in Appendix D.

3.4 Ambiguity

By design, most of the questions compare direct opposites, making ambiguity rare, however samples involving relations that are not direct opposites leave more room for interpretation. To prove solvability and to ensure unambiguous design we measure human performance. The humans obtain extremely high scores with low variance, from which we conclude that it is not ambiguous. (See Section 4.4 for details)

To ensure clarity in spatial references, we adopt the camera view as the default perspective for all annotations. Since the dataset excludes humans, we eliminate potential ambiguity arising from subjective references such as "on a person's left/right side". The camera-based viewpoint provides a consistent and semantically grounded interface between visual and linguistic modalities, reducing interpretational variability. Although most prompts contrast clearly defined opposites minimizing ambiguity, certain spatial relations, such as "near" vs. "far", or comparative attributes like "big" vs. "small", can be inherently subjective. To mitigate this subjectivity, each instance presents

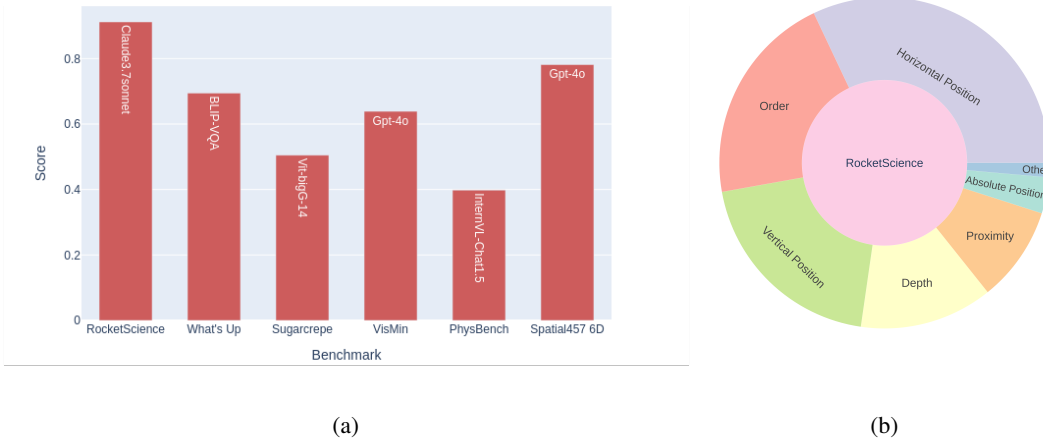


Figure 3: Overview of the RocketScience benchmark. **(a)** Benchmark quality score for popular benchmarks where higher values indicate greater challenge and headroom. **(b)** Categories contained in the RocketScience dataset. (Full breakdown in Figure 8 in Appendix D)

the model with two alternatives (e.g. a pair of captions or images), ensuring that the comparative context disambiguates the intended meaning. For instance, given two images and the caption “scissors close to the door”, the correct answer becomes the one where the scissors are closer to the door relative to the alternative. As discussed in Section 3.2.1, a subset of examples in the horizontal position category includes horizontally mirrored images. In some cases, mirrored text may be visible, but never in a way that would simplify the task.

3.5 Difficulty

Even in contrastive benchmarks, considerable care must be taken to prevent models from exploiting unintended correlations or shortcut signals (i.e. solving tasks via surface-level cues rather than the intended reasoning). In our benchmark, samples are designed to control for this by ensuring that both target entities are present in each image. For example, in the contrastive pair “beetroot in a bin” vs. “no beetroot in a bin”, both the beetroot and the bin appear in both images. The latter image still contains the beetroot, but located outside the bin. This ensures that solving the task requires understanding the spatial relation, rather than simply detecting object presence. This construction prevents the task from being reduced to simple object detection and enforces a requirement for spatial comprehension.

Additionally, a subset of image-text pairs is constructed to require fine-grained attribute localization, rather than reliance on nouns alone. For instance, in a caption such as “a crushed milk carton to the left of an intact milk carton”, both objects belong to the same category, but differ in visual attributes. Tasks of this nature increase linguistic complexity and demand comprehension of adjectives in context, further discouraging reliance on shallow pattern matching.

To account for the fact that benchmark datasets often leak into model training over time, we argue that benchmark difficulty and informativeness are best evaluated at the time of their release. To quantify this, we introduce a normalized scoring metric:

$$\text{Score} = \frac{1 - \text{SOTA}_{\text{AtRelease}}}{1 - \text{Random}} \quad (1)$$

where $\text{SOTA}_{\text{AtRelease}}$ is the best-performing *non-CoT* model at the time of dataset release, and Random represents the expected score from uniform guessing. This metric captures the benchmark’s potential to differentiate between random and competent models at release time, rather than its current saturation level. It also allows for fairer comparisons across datasets of varying ages. Under this metric, our benchmark RocketScience outperforms several contemporaneous datasets as can be observed in figure 3a, even though the metric favors older benchmarks as they do not have to compare to current SOTA models.

4 Experimental Setup & Evaluation

4.1 Models

We evaluate three main types of models: CLIP-like (contrastive) embedding models, language models with vision capabilities (VLMs) and reasoning VLMs with a hidden chain-of-thought (reflective models). The most interesting models in our context are those that have been trained for more fine-grained visual understanding like NegCLIP [73], Paligemma [7], SpaceOm [11] and GLM4.1 [60]. We also select a range of VLMs from top frontier VLMs like GPT-4o [46] and Llama4 [41]. Additionally, we also evaluate Gemini 2.5 [21] and o4-mini [45] as examples of reflective models. We only include commercial models that can be run for less than ten US Dollars on the benchmark and only run them once due to their cost.

4.2 Evaluation

Preprocessing: We resize all images to 1024×1024 . Most of the CLIP models and Paligemma include custom preprocessing that downsizes images further. Closed VLMs like GPT-4o and Claude Sonnet allow for large image inputs. For all API models we use .png to maintain image quality.

Inputs: Each contrastive pair is split into four questions:

- Q1: First image + both captions \rightarrow Which is the correct caption?
- Q2: Second image + both captions \rightarrow Which is the correct caption?
- Q3: First text + both images \rightarrow Which is the correct image?
- Q4: Second text + both images \rightarrow Which is the correct image?

The correct answer for Q1 and Q3 is always 1 and for Q2 and Q4 it's always 2. This automatic alternation means it is not necessary to shuffle the answers as is necessary for non-contrastive benchmarks to avoid models simply always choosing the first answer.

CLIP-like models receive the images and captions individually and compute the similarity between them. The exact inputs and outputs can be seen in Appendix B.

Decoding: For API models we set temperature = 0 where possible, to minimize variance. Please note that this still does not lead to fully deterministic according to their documentation [21, 41, 46]. CLIP models are deterministic during evaluation. For local VLMs, we use greedy decoding in order to ensure results are reproducible.

Hardware and Runtime: The local models were evaluated on a T4-GPU. Most CLIP models evaluate within a few minutes. The API models take between 1 and 2.5 hours to run, with the reflective models taking the most time.

4.3 Metrics

We use text score, image score and group score, as introduced by Winoground [61], for the CLIP-like models and modify the scores to be suitable for VLMs as follows: for text score, one image and two captions are used as input with the model then choosing the matching caption. For image score, two images and one caption act as the input and the model must select the correct image. This approach is not perfectly comparable to the original scores, as they are based on individually encoded texts and images, making them slightly harder than our newly defined VLM scores. The adapted image score, text score and group score are defined in Equations 2, 3 and 4 as follows:

$$f(I_0, I_1, T_0, T_1) = \begin{cases} 1, & \text{if } \hat{y}(I_0, I_1, T_0) = \text{"choose } I_0\text{" and } \hat{y}(I_0, I_1, T_1) = \text{"choose } I_1\text{"} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

$$g(I_0, I_1, T_0, T_1) = \begin{cases} 1, & \text{if } \hat{y}(I_0, T_0, T_1) = \text{"choose } T_0\text{" and } \hat{y}(I_1, T_0, T_1) = \text{"choose } T_1\text{"} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

$$h(I_0, I_1, T_0, T_1) = \begin{cases} 1, & \text{if } f(I_0, I_1, T_0, T_1) = 1 \text{ and } g(I_0, I_1, T_0, T_1) = 1 \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where f, g, h are the equivalents of image score, text score and group score and I_0, I_1, T_0, T_1 are the images and texts from a contrastive pair.

4.4 Human Evaluation

To assess the solvability of the benchmark we present the questions to humans with no formal linguistics education. (n=4) The evaluation scheme is the same as for VLMs, so they either get a caption and have to choose between two images or get an image and have to choose the right of two captions. However because humans can remember their previous responses, we only present one question from each contrastive pair to each human. The scores are then computed across all testers. For humans the two images or the two captions are shuffled - otherwise the correct one would always be in the same spot. The image resolution for the human evaluation is 600 by 600 pixels so that two images fit onto a laptop screen comfortably. (This is lower than the resolution the models receive) The exact instructions and interface are available in Appendix G. Our participants score a mean accuracy of 0.985 with a standard deviation of 0.008. We conclude that RocketScience has an extremely low level of ambiguity.

5 Results

We present the results in Table 1. Table 3 in Appendix E includes additional CLIP-like models. Model names that end with *_cot* have been prompted to do chain-of-thought reasoning on top of their usual system prompt. Gemini 2.5 and o4-mini perform reasoning internally before responding. We find that all open vision-language models, even those trained for spatial understanding perform very poorly and often below random chance. This phenomenon has also been observed in other contrastive vision-language benchmarks, such as Winoground [61] and WhatsUp [31]. Only reasoning models come close to human performance.

Model Name	Horiz	Vert	Depth	Prox	Order	Abs Pos	Total
Random chance	0.17	0.17	0.17	0.17	0.17	0.17	0.17
Human	0.96	0.98	0.95	1.00	0.92	0.80	0.95
ViT-B-32negCLIP [73]	0.00	0.02	0.00	0.00	0.03	0.00	0.01
CoCa_ViT-L [72]	0.00	0.04	0.00	0.00	0.02	0.00	0.01
PaliGemma-3b-mix-448 [7]	0.01	0.04	0.00	0.00	0.02	0.20	0.02
Qwen-2.5-vl-72b-instruct [5]	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Qwen-vl-max [4]	0.01	0.02	0.00	0.00	0.02	0.00	0.01
Claude-3-7-sonnet-20250219 [2]	0.17	0.38	0.11	0.52	0.18	0.00	0.24
Llama-4-maverick [41]	0.17	0.23	0.05	0.52	0.15	0.40	0.20
GPT-4o-2024-08-06 [46]	0.08	0.40	0.24	0.40	0.07	0.00	0.19
Llama-4-maverick_cot [41]	0.53	0.42	0.21	0.64	0.45	0.20	0.44
GPT-4o-2024-08-06_cot [46]	0.44	0.68	0.39	0.56	0.52	0.60	0.51
SpaceOm [11]	0.01	0.02	0.03	0.04	0.00	0.00	0.01
Glm-4.1v-9b-thinking [60]	0.77	0.68	0.42	0.4	0.68	1.00	0.64
Gemini-2.5-pro-preview-03-25 [21]	0.91	0.87	0.76	0.76	0.80	0.80	0.83
o4-mini [45]	0.97	0.91	0.89	0.68	0.88	1.00	0.89

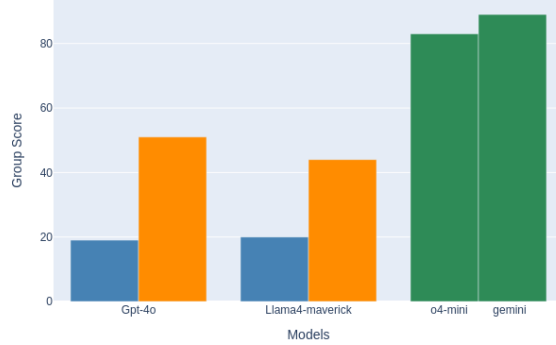
Table 1: Group score on RocketScience by category of the task (Horizontal Position, Vertical Position, Depth, Proximity, Order, Absolute Position, Total Score). Models run with chain-of-thought prompting are marked in orange and reflective models are marked in green. The highest score in each category is bold.

5.1 Why Are Chain-of-Thought Models Better?

We hypothesize that there are two main steps necessary to detect a spatial relation: localization of objects and inference of spatial relation. We examine both stages to determine why reasoning models like o4-mini perform much better than their non-reasoning counterparts like gpt-4o:

Model	Acc lf	Acc rf
gpt-4o	96.11 ± 0.96	90.55 ± 0.96
o4-mini	96.66 ± 0.00	95.56 ± 1.92

(a) Mean accuracy and standard deviation over 3 runs on the localization task (horizontal position subset). The difference in localization performance between non-CoT and CoT models is minimal. A small performance gap is also observed based on object order in the prompt: accuracy is slightly higher when the first-mentioned object appears on the left (Acc lf) compared to when it appears on the right (Acc rf), for both model types.



(b) Group score comparison of models without chain-of-thought (blue), with explicit chain-of-thought prompting (orange), and with implicit chain-of-thought reasoning (green).

Figure 4: Disentanglement experiments of (a) localisation and (b) CoT reasoning

- **Localization of objects:** We test whether o4-mini is better at localizing the objects than gpt-4o. We use the horizontal position subset for this (on which gpt-4o performs very poorly). The two models are prompted to provide bounding boxes for both objects in each image. We then check whether the coordinates of both objects are in the correct spatial configuration. (not exact location) Figure 4a shows the results. We find that gpt-4o is very close to o4-mini’s performance, indicating that reasoning does not help with the localisation stage, but only with concluding the correct spatial relation.
- **Inference of spatial relation:** We use two sets of prompts for non-reflective models: The first prompts each model to output the answer alone immediately and nothing else (non-CoT). The second prompts each model to first reason and then output the results (CoT). We find that the second prompt significantly improves performance over the first one, indicating that chain-of-thought reasoning plays the biggest role in the improved performance. (See Figure 4b)

6 Discussion

6.1 Limitations

First, models accessed via APIs were only evaluated once due to cost constraints. Although we set the sampling temperature to zero where possible to ensure determinism, outputs may still vary slightly between runs. Second, while our dataset aims to reflect real-world complexity, some scenes remain less cluttered with objects than typical real-world environments. Introducing additional clutter while maintaining clear, unambiguous relations remains a significant challenge. Third, we observe notable performance improvements using chain-of-thought prompting; however, this benefit may be specific to top-tier commercial models and might not generalize to smaller or open-source models. Finally, although we strive to minimize changes within each contrastive pair, slight variations in camera angle may occur. This could introduce a potential shortcut where models exploit angle differences to infer spatial relations rather than relying solely on object configurations.

6.2 Ethics

People and personal data are explicitly excluded as subjects within this dataset in an effort to minimize the unnecessary use of human data in experiments. Each sample within the dataset was meticulously reviewed by three of the authors of this paper in an effort to maintain quality standards by minimizing errors and omissions. It is necessary to emphasize that the geographical locations of images are limited to the US and Europe, making the environment and objects specific to these locations. It is therefore crucial to stress that a model that performs well on this benchmark will not necessarily perform well in all geographic locations or with objects specific to them. Additionally, though this benchmark has successfully revealed that most vision-language models fail to effectively model spatial relations, it should not be relied upon alone as a form of evaluation for these relations, as benchmark performance should not act as a replacement for application and location-specific quality

control and testing. We caution that though benchmarks can reveal *some* model shortcomings, they are never exhaustive and we caution against their use alone to rank algorithms in real-world applications/production as this can lead to unexpected societal outcomes.

6.3 Conclusion

We introduce RocketScience, a challenging new benchmark for evaluating spatial relation understanding in vision language models. Built from scratch using real-world, contrastive image-text pairs, RocketScience reveals that most open-source and commercial models perform surprisingly poorly. Our analysis shows that chain-of-thought prompting significantly improves performance. By disentangling the effects of object localization and spatial reasoning, we find that the primary limitation lies in models’ ability to perform structured reasoning about spatial relations, rather than in their visual perception. We hope RocketScience serves as a diagnostic and development tool for future VLMs, encouraging research toward models with more robust spatial understanding.

Acknowledgements and Funding

This publication has emanated from research conducted with the financial support of Science Foundation Ireland under Grant number 18/CRT/6183. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

References

- [1] Ahmad Mustafa Anis, Hasnain Ali, and Saquib Sarfraz. On the limitations of vision-language models in understanding image transforms, 2025. URL <https://arxiv.org/abs/2503.09837>.
- [2] Anthropic. Claude sonnet 3.7. <https://www.anthropic.com/news/claude-3-7-sonnet>, 2025. Accessed: May 2025.
- [3] Rabiul Awal, Saba Ahmadi, Le Zhang, and Aishwarya Agrawal. Vismin: Visual minimal-change understanding, 2025. URL <https://arxiv.org/abs/2407.16772>.
- [4] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. URL <https://arxiv.org/abs/2308.12966>.
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL <https://arxiv.org/abs/2502.13923>.
- [6] Youssef Bencheikroun, Megi Dervishi, Mark Ibrahim, Jean-Baptiste Gaya, Xavier Martinet, Grégoire Mialon, Thomas Scialom, Emmanuel Dupoux, Dieuwke Hupkes, and Pascal Vincent. Worldsense: A synthetic benchmark for grounded reasoning in large language models, 2023.
- [7] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer, 2024. URL <https://arxiv.org/abs/2407.07726>.
- [8] Nitzan Bitton-Guetta, Yonatan Bitton, Jack Hessel, Ludwig Schmidt, Yuval Elovici, Gabriel Stanovsky, and Roy Schwartz. Breaking common sense: Whoops! a vision-and-language benchmark of synthetic and compositional images, 2023. URL <https://arxiv.org/abs/2303.07274>.

- [9] Nitzan Bitton-Guetta, Aviv Slobodkin, Aviya Maimon, Eliya Habba, Royi Rassin, Yonatan Bitton, Idan Szpektor, Amir Globerson, and Yuval Elovici. Visual riddles: a commonsense and world knowledge challenge for large vision and language models, 2024. URL <https://arxiv.org/abs/2407.19474>.
- [10] Wenxiao Cai, Iaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and Bo Zhao. Spatialbot: Precise spatial understanding with vision language models, 2025. URL <https://arxiv.org/abs/2406.13642>.
- [11] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities, 2024. URL <https://arxiv.org/abs/2401.12168>.
- [12] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models?, 2024. URL <https://arxiv.org/abs/2403.20330>.
- [13] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision language models, 2024. URL <https://arxiv.org/abs/2406.01584>.
- [14] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *CoRR*, abs/2202.04053, 2022. URL <https://arxiv.org/abs/2202.04053>.
- [15] Keng Ji Chow, Samson Tan, and Min-Yen Kan. TraVLR: Now you see it, now you don’t! a bimodal dataset for evaluating visio-linguistic reasoning. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 3322–3347, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.
- [16] Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. Embspatial-bench: Benchmarking spatial understanding for embodied tasks with large vision-language models, 2024. URL <https://arxiv.org/abs/2406.05756>.
- [17] Sri Harsha Dumpala, Aman Jaiswal, Chandramouli Sastry, Evangelos Milios, Sageev Oore, and Hassan Sajjad. Sugarcrepe++ dataset: Vision-language model sensitivity to semantic and lexical alterations, 2024. URL <https://arxiv.org/abs/2406.11171>.
- [18] Matias Duran, Thomas Laurent, Ellen Rushe, and Anthony Ventresque. Metamorphic testing for pose estimation systems, 2025. URL <https://arxiv.org/abs/2502.09460>.
- [19] 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, 2024. URL <https://arxiv.org/abs/2404.12390>.
- [20] Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vibhav Vineet, Eric Horvitz, Ece Kamar, Chitta Baral, and Yezhou Yang. Benchmarking spatial relationships in text-to-image generation, 2022.
- [21] Google. Gemini 2.5. <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>, 2025. Accessed: May 2025.
- [22] Brian Gordon, Yonatan Bitton, Yonatan Shafir, Roopal Garg, Xi Chen, Dani Lischinski, Daniel Cohen-Or, and Idan Szpektor. Mismatch quest: Visual and textual feedback for image-text misalignment, 2024. URL <https://arxiv.org/abs/2312.03766>.
- [23] Ankit Goyal, Kaiyu Yang, Dawei Yang, and Jia Deng. Rel3d: A minimally contrastive benchmark for grounding spatial relations in 3d, 2020. URL <https://arxiv.org/abs/2012.01634>.
- [24] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–6334, 2017. doi: 10.1109/CVPR.2017.670.

- [25] Nils Hoehing, Ellen Rushe, and Anthony Ventresque. What’s left can’t be right – the remaining positional incompetence of contrastive vision-language models, 2023. URL <https://arxiv.org/abs/2311.11477>.
- [26] Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. Sugar-crepe: Fixing hackable benchmarks for vision-language compositionality, 2023.
- [27] Irene Huang, Wei Lin, M. Jehanzeb Mirza, Jacob A. Hansen, Sivan Doveh, Victor Ion Butoi, Roei Herzig, Assaf Arbelle, Hilde Kuehne, Trevor Darrell, Chuang Gan, Aude Oliva, Rogerio Feris, and Leonid Karlinsky. Conme: Rethinking evaluation of compositional reasoning for modern vlms, 2024. URL <https://arxiv.org/abs/2406.08164>.
- [28] Jen-Tse Huang, Dasen Dai, Jen-Yuan Huang, Youliang Yuan, Xiaoyuan Liu, Wenxuan Wang, Wenxiang Jiao, Pinjia He, and Zhaopeng Tu. Visfactor: Benchmarking fundamental visual cognition in multimodal large language models, 2025. URL <https://arxiv.org/abs/2502.16435>.
- [29] Jinsheng Huang, Liang Chen, Taian Guo, Fu Zeng, Yusheng Zhao, Bohan Wu, Ye Yuan, Haozhe Zhao, Zhihui Guo, Yichi Zhang, Jingyang Yuan, Wei Ju, Luchen Liu, Tianyu Liu, Baobao Chang, and Ming Zhang. Mmevalpro: Calibrating multimodal benchmarks towards trustworthy and efficient evaluation, 2025. URL <https://arxiv.org/abs/2407.00468>.
- [30] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [31] Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s “up” with vision-language models? investigating their struggle with spatial reasoning. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9161–9175, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.568. URL <https://aclanthology.org/2023.emnlp-main.568>.
- [32] Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. Naturalbench: Evaluating vision-language models on natural adversarial samples, 2024. URL <https://arxiv.org/abs/2410.14669>.
- [33] Xianhang Li, Zeyu Wang, and Cihang Xie. An inverse scaling law for clip training, 2023. URL <https://arxiv.org/abs/2305.07017>.
- [34] Yizhi Li, Ge Zhang, Yinghao Ma, Ruibin Yuan, Kang Zhu, Hangyu Guo, Yiming Liang, Jiaheng Liu, Zekun Wang, Jian Yang, Siwei Wu, Xingwei Qu, Jinjie Shi, Xinyue Zhang, Zhenzhu Yang, Xiangzhou Wang, Zhaoxiang Zhang, Zachary Liu, Emmanouil Benetos, Wenhao Huang, and Chenghua Lin. Omnibench: Towards the future of universal omni-language models, 2024. URL <https://arxiv.org/abs/2409.15272>.
- [35] Zongxia Li, Xiyang Wu, Hongyang Du, Fuxiao Liu, Huy Nghiem, and Guangyao Shi. A survey of state of the art large vision language models: Alignment, benchmark, evaluations and challenges, 2025. URL <https://arxiv.org/abs/2501.02189>.
- [36] Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. Reasoning paths with reference objects elicit quantitative spatial reasoning in large vision-language models, 2024. URL <https://arxiv.org/abs/2409.09788>.
- [37] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015.
- [38] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning, 2023.

- [39] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player?, 2024. URL <https://arxiv.org/abs/2307.06281>.
- [40] Yujie Lu, Dongfu Jiang, Wenhui Chen, William Yang Wang, Yejin Choi, and Bill Yuchen Lin. Wildvision: Evaluating vision-language models in the wild with human preferences, 2024. URL <https://arxiv.org/abs/2406.11069>.
- [41] Meta. Llama 4. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, 2025. Accessed: May 2025.
- [42] Imanol Miranda, Ander Salaberria, Eneko Agirre, and Gorka Azkune. Bivlc: Extending vision-language compositionality evaluation with text-to-image retrieval, 2024. URL <https://arxiv.org/abs/2406.09952>.
- [43] Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang Ning, and Parisa Kordjmeshidi. Spartqa: : A textual question answering benchmark for spatial reasoning, 2021.
- [44] NVIDIA, :, Alisson Azzolini, Hannah Brandon, Prithvijit Chattopadhyay, Huayu Chen, Jinju Chu, Yin Cui, Jenna Diamond, Yifan Ding, Francesco Ferroni, Rama Govindaraju, Jinwei Gu, Siddharth Gururani, Imad El Hanafi, Zekun Hao, Jacob Huffman, Jingyi Jin, Brendan Johnson, Rizwan Khan, George Kurian, Elena Lantz, Nayeon Lee, Zhaoshuo Li, Xuan Li, Tsung-Yi Lin, Yen-Chen Lin, Ming-Yu Liu, Alice Luo, Andrew Mathau, Yun Ni, Lindsey Pavao, Wei Ping, David W. Romero, Misha Smelyanskiy, Shuran Song, Lyne Tchapmi, Andrew Z. Wang, Boxin Wang, Haoxiang Wang, Fangyin Wei, Jiashu Xu, Yao Xu, Xiaodong Yang, Zhuolin Yang, Xiaohui Zeng, and Zhe Zhang. Cosmos-reason1: From physical common sense to embodied reasoning, 2025. URL <https://arxiv.org/abs/2503.15558>.
- [45] o4 mini. o4-mini system card. <https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf>, 2025. Accessed: May 2025.
- [46] OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoonchian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braundstein, Andrew Cann, Andrew Codisopoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Gierler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian O’Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang,

Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavín Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljube, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Pene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunningham, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiye Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.

- [47] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten, 2023. URL <https://arxiv.org/abs/2302.12066>.
- [48] Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena, 2022.
- [49] Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnav Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin, Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehringer, Jiaqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu

Ren, Usman Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li, Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy, Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori, Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, Joseph Marvin Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao, Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever, Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Anmol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoun, Alvin Jin, Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Antonenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nurdin Kaparov, Allen Zang, Iliia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Poritski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe, Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Victor Efren Guadarrama Vilchis, Immo Klose, Ujjwala Ananthaswaran, Adam Zweiger, Kaivalya Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans, Haline Heidinger, Maksim Radionov, Václav Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Goldfarb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall, Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla, Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayeaz Aziz, Mark H Inlow, Hao He, Ling Zhang, Younesse Kaddar, Ivar Ångquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakrishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avishy Carmi, Ethan D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw, JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam, Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu, Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hendrik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou Nie, Anna Szytber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani, Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson, Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu, Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bosio, Fereshteh Kazemi, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes, Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Durand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff, Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William Merrill, Moritz Firsching, Carter Harris, Stefan Ciobăcă, Jason Gross, Rohan Pandey, Ilya Gusev, Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piperski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Dersch, Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Zebaze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler, Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy

Zou, Muthu Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasouqi, Alexander Shen, Bitá Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya, Nick Winter, Miguel Orbegoza Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hossain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflor, Haile Kassahun, Alena Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani, Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristyy, Stephen Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday, Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mariana Costa, Hubeyb Gurdogan, Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek, Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchhof, Linjie Dai, Evan Kim, Johan Ferret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca, Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob Votava, Vladimir Vinnikov, Ethan Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin, Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li, Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou, Julien Degorre, Yigit Yalin, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Bosca, Oleg Shumar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie, Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara, Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Christoph Demian, Peyman Kasani, Xinyu Zhang, David Avagian, Eshawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Liakhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja Somrak, Eric Vergo, Juehang Qin, Benjámín Borbás, Eric Chu, Jack Lindsey, Antoine Jallon, I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer, Tran Duc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie, Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long, Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasilios Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselynn Grace Montecillo, Jakub Łucki, Russell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deepakkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane, Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Mündler, Sören Möller, Luca Arnaboldi, Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff, Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, Innocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahaloohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy,

Darling Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh, Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ragavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun, Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang, Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbbar, Lorenzo Vaquero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia, Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le, Mickaël Noyé, Michał Perełkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Antonio A. W. L. Wong, Maria del Rio-Chanona, Dániel Kondor, Pieter Francois, Ed Chilstrey, Jakob Zsambok, Dan Hoyer, Jenny Reddish, Jakob Hauser, Francisco-Javier Rodrigo-Ginés, Suchandra Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto, Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bartomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger, Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li, Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang, Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park, Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam Perlitz, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra, Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao, Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial, Jiachen Liu, Xingyu Qu, Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda Chen, Mengze Tang, Jiaqi Deng, Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, Yushun Chen, Sara Vera Marjanović, Junwoo Ha, Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, Cedegao E. Zhang, Zhun Wang, Gaël Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, Zhe Ye, Stefano Ermon, Ignacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi Wei, Hongzheng Chen, Kunal Pai, Ahmed Elkhany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra, Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosing, Timothy Kang, Barbara Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel Villanueva, Ivan Rannev, Igor Chernyavsky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao Dong, Jianxin Wang, Laila Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bohdal, Atharv Singh Patlan, Shehzaad Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal, Bingsen Chen, Kutay Tire, Atak Talay Yücel, Brandon Christof, Veerupaksh Singla, Zijian Song, Sanxing Chen, Jiaxin Ge, Kaustubh Ponkshe, Isaac Park, Tianneng Shi, Martin Q. Ma, Joshua Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda Zhu, Ziyi Zhang, Vaidehi Patil, Ketan Jha, Qiutong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling Vieira, Alham Fikri Aji, Jae-Won Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu, Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, James Bailey, Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong Yeo, Eeshaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan Zhang, Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gulati, Songyang Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M. Cao, Aakaash Nattanmai, Gordon McKellips, Anish Cheraku, Asim Suhail, Ethan Luo, Marvin Deng, Jason Luo, Ashley Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael Liu, Advait Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua Lass, Hubert Yang, Surya Sunkari, Vishruth Bharath, Violet Ai, James Leung, Rishit Agrawal, Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, Aryan Singh, Ethan Luo, Clark Peng, Tyler Osbey,

- Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel, Vidhi Kulkarni, Vijaykaarti Sundarapandiyani, Ashley Zhang, Andrew Le, Zafir Nasim, Srikar Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet Saha, Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo, Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity’s last exam, 2025. URL <https://arxiv.org/abs/2501.14249>.
- [50] Tanawan Premisri and Parisa Kordjamshidi. Forest: Frame of reference evaluation in spatial reasoning tasks, 2025. URL <https://arxiv.org/abs/2502.17775>.
- [51] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 18–24 Jul 2021.
- [52] Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind, 2024. URL <https://arxiv.org/abs/2407.06581>.
- [53] Jonathan Roberts, Mohammad Reza Taesiri, Ansh Sharma, Akash Gupta, Samuel Roberts, Ioana Croitoru, Simion-Vlad Bogolin, Jialu Tang, Florian Langer, Vyas Raina, Vatsal Raina, Hanyi Xiong, Vishaal Udandara, Jingyi Lu, Shiyang Chen, Sam Purkis, Tianshuo Yan, Wenye Lin, Gyungin Shin, Qiaochu Yang, Anh Totti Nguyen, David I. Atkinson, Aaditya Baranwal, Alexandru Coca, Mikah Dang, Sebastian Dziadzio, Jakob D. Kunz, Kaiqu Liang, Alexander Lo, Brian Pulfer, Steven Walton, Charig Yang, Kai Han, and Samuel Albanie. Zerobench: An impossible visual benchmark for contemporary large multimodal models, 2025. URL <https://arxiv.org/abs/2502.09696>.
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 36479–36494. Curran Associates, Inc., 2022.
- [55] Fatemeh Shiri, Xiao-Yu Guo, Mona Golestan Far, Xin Yu, Gholamreza Haffari, and Yuan-Fang Li. An empirical analysis on spatial reasoning capabilities of large multimodal models, 2024. URL <https://arxiv.org/abs/2411.06048>.
- [56] Chan Hee Song, Valts Blukis, Jonathan Tremblay, Stephen Tyree, Yu Su, and Stan Birchfield. Robospacial: Teaching spatial understanding to 2d and 3d vision-language models for robotics, 2025. URL <https://arxiv.org/abs/2411.16537>.
- [57] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs, 2019.
- [58] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale, 2023.
- [59] Emilia Szymanska, Mihai Dusmanu, Jan-Willem Buurlage, Mahdi Rad, and Marc Pollefeys. Space3d-bench: Spatial 3d question answering benchmark, 2024. URL <https://arxiv.org/abs/2408.16662>.
- [60] V Team, Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, Lihang Pan, Shuaiqi Duan, Weihang Wang, Yan Wang, Yean Cheng, Zehai He, Zhe Su, Zhen Yang, Ziyang Pan, Aohan Zeng, Baoxu Wang, Bin Chen, Boyan Shi, Changyu Pang, Chenhui Zhang, Da Yin, Fan Yang, Guoqing Chen, Jiazheng Xu, Jiale Zhu, Jiali Chen, Jing Chen, Jinhao Chen, Jinghao Lin, Jinjiang Wang, Junjie Chen, Leqi Lei, Letian Gong, Leyi Pan, Mingdao Liu, Mingde Xu, Mingzhi Zhang, Qinkai Zheng, Sheng Yang, Shi Zhong, Shiyu Huang, Shuyuan Zhao, Siyan Xue, Shangqin Tu, Shengbiao Meng, Tianshu Zhang, Tianwei Luo, Tianxiang Hao, Tianyu Tong, Wenkai Li, Wei Jia, Xiao Liu, Xiaohan Zhang,

- Xin Lyu, Xinyue Fan, Xuancheng Huang, Yanling Wang, Yadong Xue, Yanfeng Wang, Yanzi Wang, Yifan An, Yifan Du, Yiming Shi, Yiheng Huang, Yilin Niu, Yuan Wang, Yuanchang Yue, Yuchen Li, Yutao Zhang, Yuting Wang, Yu Wang, Yuxuan Zhang, Zhao Xue, Zhenyu Hou, Zhengxiao Du, Zihan Wang, Peng Zhang, Debing Liu, Bin Xu, Juanzi Li, Minlie Huang, Yuxiao Dong, and Jie Tang. Glm-4.5v and glm-4.1v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning, 2025. URL <https://arxiv.org/abs/2507.01006>.
- [61] 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 (CVPR)*, pages 5238–5248, June 2022.
 - [62] Kexin Tian, Jingrui Mao, Yunlong Zhang, Jiwan Jiang, Yang Zhou, and Zhengzhong Tu. Nuscenes-spatialqa: A spatial understanding and reasoning benchmark for vision-language models in autonomous driving, 2025. URL <https://arxiv.org/abs/2504.03164>.
 - [63] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms, 2024. URL <https://arxiv.org/abs/2401.06209>.
 - [64] Michael Tschannen, Manoj Kumar, Andreas Steiner, Xiaohua Zhai, Neil Houlsby, and Lucas Beyer. Image captioners are scalable vision learners too, 2023.
 - [65] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, Olivier Hénaff, Jeremiah Harmsen, Andreas Steiner, and Xiaohua Zhai. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features, 2025. URL <https://arxiv.org/abs/2502.14786>.
 - [66] Mor Ventura, Michael Toker, Nitay Calderon, Zorik Gekhman, Yonatan Bitton, and Roi Reichart. NI-eye: Abductive nli for images, 2024. URL <https://arxiv.org/abs/2410.02613>.
 - [67] Fei Wang, Xingyu Fu, James Y. Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, Wenxuan Zhou, Kai Zhang, Tianyi Lorena Yan, Wenjie Jacky Mo, Hsiang-Hui Liu, Pan Lu, Chunyuan Li, Chaowei Xiao, Kai-Wei Chang, Dan Roth, Sheng Zhang, Hoifung Poon, and Muhao Chen. Muirbench: A comprehensive benchmark for robust multi-image understanding, 2024. URL <https://arxiv.org/abs/2406.09411>.
 - [68] Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Yixuan Li, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models, 2024. URL <https://arxiv.org/abs/2406.14852>.
 - [69] Xingrui Wang, Wufei Ma, Tiezheng Zhang, Celso M de Melo, Jieneng Chen, and Alan Yuille. Spatial457: A diagnostic benchmark for 6d spatial reasoning of large multimodal models, 2025. URL <https://arxiv.org/abs/2502.08636>.
 - [70] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
 - [71] Eunice Yiu, Maan Qraitem, Charlie Wong, Anisa Noor Majhi, Yutong Bai, Shiry Ginosar, Alison Gopnik, and Kate Saenko. Kiva: Kid-inspired visual analogies for testing large multimodal models, 2024. URL <https://arxiv.org/abs/2407.17773>.
 - [72] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=Ee277P3AYC>.
 - [73] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it?, 2022. URL <https://arxiv.org/abs/2210.01936>.

- [74] Aimen Zerroug, Mohit Vaishnav, Julien Colin, Sebastian Musslick, and Thomas Serre. A benchmark for compositional visual reasoning, 2022.
- [75] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023.
- [76] Fengji Zhang, Linqun Wu, Huiyu Bai, Guancheng Lin, Xiao Li, Xiao Yu, Yue Wang, Bei Chen, and Jacky Keung. Humaneval-v: Benchmarking high-level visual reasoning with complex diagrams in coding tasks, 2025. URL <https://arxiv.org/abs/2410.12381>.
- [77] Yizhe Zhang, He Bai, Ruixiang Zhang, Jiatao Gu, Shuangfei Zhai, Josh Susskind, and Navdeep Jaitly. How far are we from intelligent visual deductive reasoning?, 2024. URL <https://arxiv.org/abs/2403.04732>.

A Appendix: Related Work

Benchmark	contrastive	new data	real data	non-schematic	size
MMBench[39]	✗	✗	✓	✓	125
SpatialEval-Real [68]	✗	✗	✓	✓	1000
VSR [38]	✗	✗	✓	✓	2195
CLEVR [30]	✗	✓	✗	✗	15,000
VALSE [48]	✗	✗	✓	✓	535
SugarCrepe [26]	✗	✗	✓	✓	1406
ConMe [27]	✗	✗	✓	✓	6793
SC++ [17]	✗	✗	✓	✓	1406
VGR (ARO) [73]	✗	✗	✓	✓	23,937
RoboSpatial-Home [56]	✗	✓	✓	✓	123
BLINK[19]	✗	✗	✓	✓	286
SpatialBench [10]	✗	✓	✓	✓	35
Space3D-bench [59]	✗	✗	✓	✓	188
Spatial-MM [55]	✗	✗	✓	✓	2,000
EmbSpatial Bench [16]	✗	✗	✓	✓	3,640
NuScenes-SpatialQA [62]	✗	✗	✓	✓	2,500,000
Cosmos1 [44]	✗	✗	✓	✓	292
Spatial457 [69]	✗	✓	✗	✓	9,990
FOREST [50]	✓	✓	✗	✓	4,352
BiVLC [42]	✓	✓	✗	✓	1,400
Rel3D [23]	✓	✓	✗	✓	27,336
VisMin [3]	✓	✓	✗	✗	622
What's Up (A+B) [31]	✓	✓	✓	✗	820
RocketScience	✓	✓	✓	✓	482

Table 2: Overview of Spatial Relations benchmarks (or with a subset for that). Contrastive means fully contrastive, new data means entirely new and no recycled parts, real data means not synthetic, non-schematic means different scenes and objects not always in the same positions, size is only the spatial relations subset and measured as number of image-text pairs.

B Appendix: Inputs and Outputs

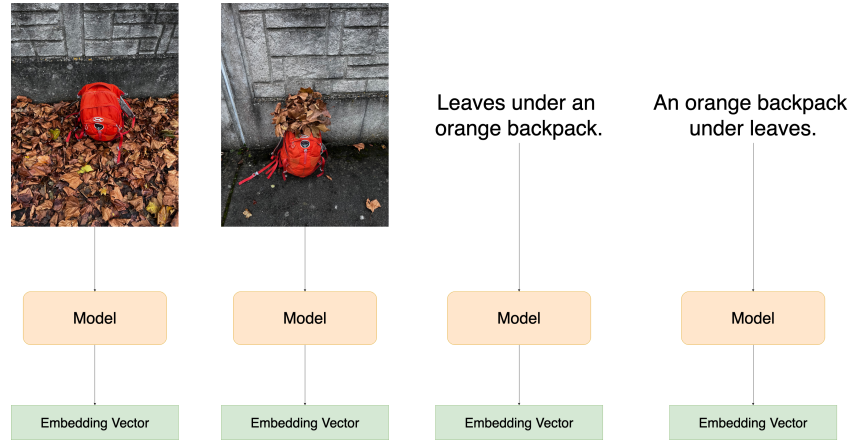


Figure 5: Function of CLIP-like models: they embed each image and text into a vector independently without having access to the other inputs at the same time

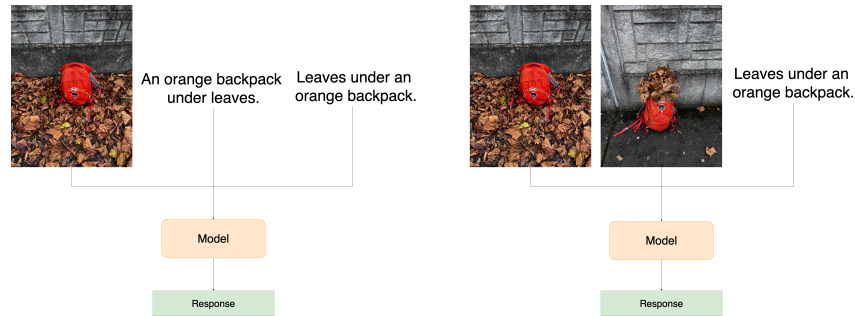


Figure 6: Function of VLMs: they receive either an image and two captions and have to choose a caption or they receive two images and a caption and have to choose an image

The model prompts are as follows: CLIP models simply receive the captions from the dataset without additional prompting. The other VLMs (both normal and reflective) receive two types of prompts:

- "Which caption fits the image best? Reply only with the number 1 or 2, nothing else. 1.) [CAPTION1] 2.) [CAPTION2]"
- "Which image fits the caption best? Reply only with the number 1 or 2, nothing else. Caption: [CAPTION1]"

For our additional experiments we test models' explicit chain-of-thought capabilities. These models are denoted with `_cot` after their name. They receive different prompts to make them reason:

- "Which caption fits the image best? Reason about it and at the end write RESPONSE and reply only with the number 1 or 2. 1.) [CAPTION1] 2.) [CAPTION2]"
- "Which image fits the caption best? Reason about it and at the end write RESPONSE and reply only with the number 1 or 2. Caption: [CAPTION1]"

Example reasoning trace from gpt-4o for a case with two images and one caption: The caption describes "A grey bin on a white towel." Image 1 shows a grey bin placed directly on a white towel. Image 2 shows a towel covering the grey bin. The best fit for the caption is image 1, as it correctly shows the bin on the towel. RESPONSE 1

C Appendix: Selected examples from the dataset



Figure 7: Examples for scene diversity in RocketScience.

D Appendix: Dataset Analysis

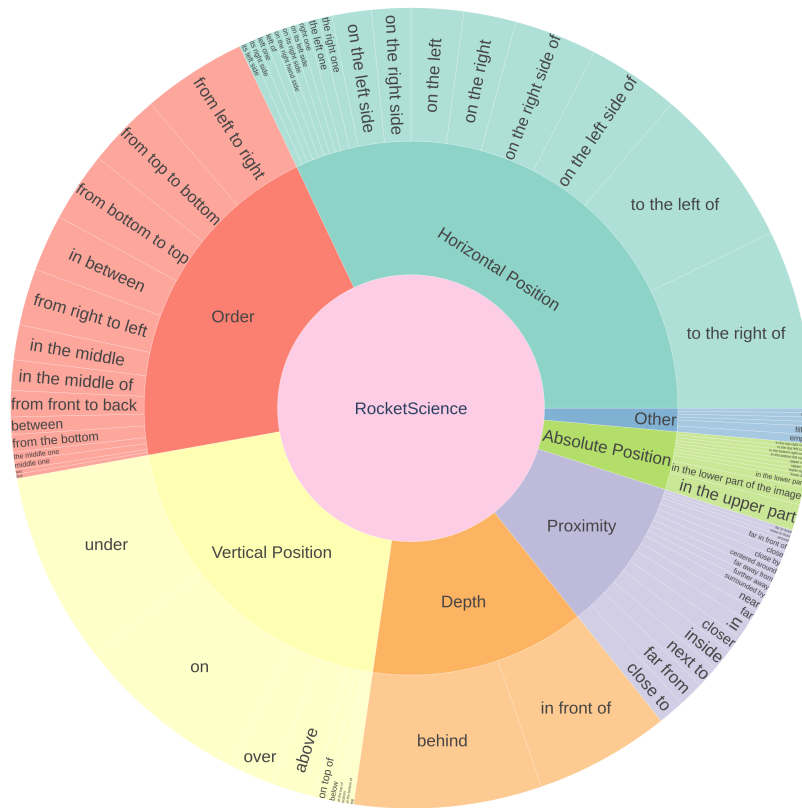


Figure 8: Dataset distribution, relative proportions

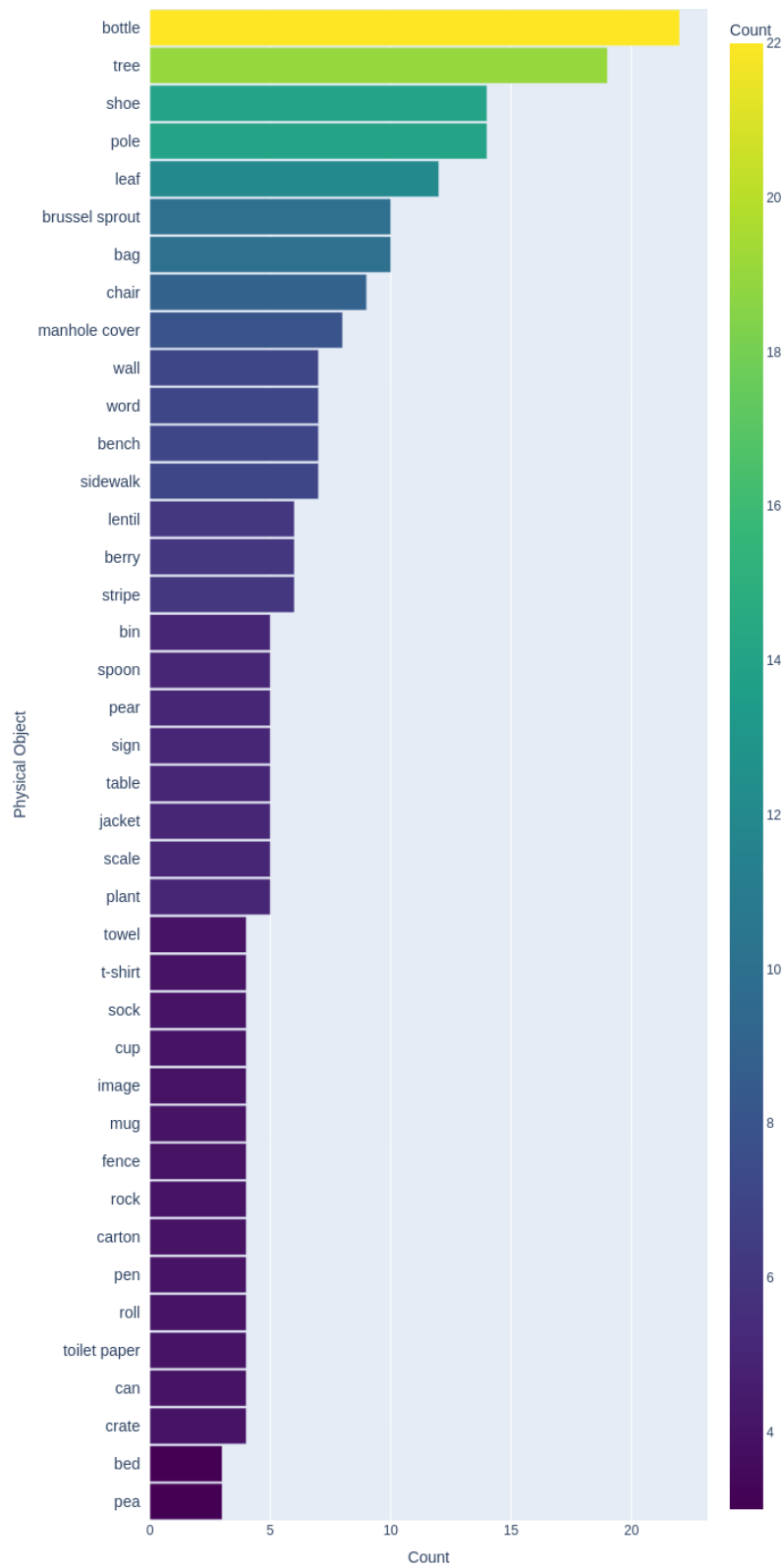


Figure 9: Physical Object Counts

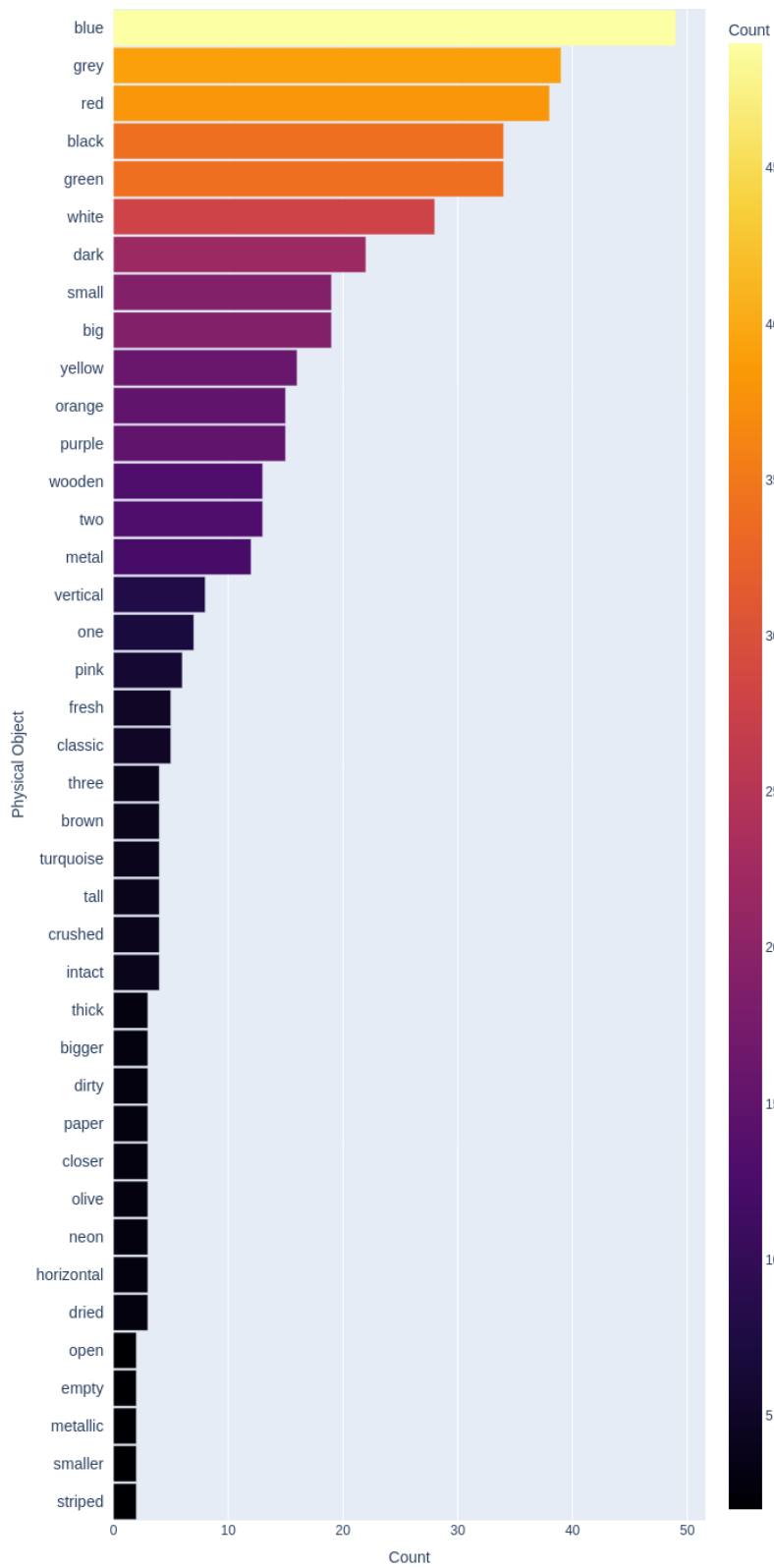


Figure 10: Adjective Counts

E Appendix: Detailed Results

modelName	textscore	imagescore	groupscore
random	0.25	0.25	0.17
human	0.97	0.98	0.95
ViT-B-32negCLIP [73]	0.08	0.04	0.01
EVA02-B-16merged2b_s8b_b131k [58]	0.13	0.04	0.02
EVA02-L-14-336merged2b_s6b_b61k [58]	0.12	0.05	0.02
ViT-B-16-SigLIPwebli [75]	0.13	0.04	0.02
ViT-L-16-SigLIP-384webli [75]	0.10	0.07	0.02
ViT-L-14-CLIPAdatcomp1b [33]	0.07	0.05	0.00
ViT-L-16-SigLIP2-512webli [65]	0.12	0.07	0.01
coca_ViT-B-32laion2b_s13b_b90k [72]	0.14	0.05	0.02
coca_ViT-L-14laion2b_s13b_b90k [72]	0.10	0.04	0.01
ViT-B-16-SigLIP2-512webli [65]	0.10	0.05	0.02
ViT-B-16openai [51]	0.11	0.05	0.02
ViT-B-32openai [51]	0.11	0.02	0.01
paligemma-3b-mix-448 [7]	0.13	0.10	0.02
qwen-2.5-vl-72b-instruct [5]	0.47	0.01	0.00
qwen-vl-max [4]	0.29	0.03	0.01
claude-3-7-sonnet-20250219 [2]	0.53	0.37	0.24
llama-4-maverick [41]	0.38	0.37	0.20
gpt-4o-2024-08-06 [46]	0.38	0.39	0.19
llama-4-maverick_cot [41]	0.59	0.66	0.44
gpt-4o-2024-08-06_cot [46]	0.73	0.66	0.51
SpaceOm [11]	0.08	0.14	0.01
glm-4.1v-9b-thinking [60]	0.84	0.72	0.64
gemini-2.5-pro-preview-03-25 [21]	0.94	0.89	0.83
o4-mini (medium) [45]	0.91	0.94	0.89

Table 3: Results on the RocketScience benchmark, the second division is CLIP-like models, the third regular VLMs, the fourth regular vlms with explicit chain-of-thought and the last VLMs with implicit chain-of-thought. All CLIP-like models and basic VLMs fail drastically while some reasoning models come very close to human performance.

Model	Horizontal			Vertical			Depth		
	ts	is	gs	ts	is	gs	ts	is	gs
paligemma-3b-mix-448	0.08	0.13	0.01	0.21	0.11	0.04	0.16	0.13	0.00
qwen-2.5-vl-72b-instruct	0.51	0.00	0.00	0.62	0.02	0.00	0.39	0.03	0.00
qwen-vl-max	0.23	0.07	0.01	0.47	0.02	0.02	0.24	0.00	0.00
claude-3-7-sonnet-20250219	0.43	0.35	0.17	0.64	0.49	0.38	0.47	0.24	0.11
llama-4-maverick	0.40	0.33	0.17	0.45	0.36	0.23	0.18	0.29	0.05
gpt-4o	0.35	0.27	0.08	0.62	0.58	0.40	0.39	0.42	0.24
llama-4-maverick_cot	0.67	0.72	0.53	0.55	0.66	0.42	0.37	0.50	0.21
gpt-4o_cot	0.72	0.63	0.44	0.83	0.77	0.68	0.74	0.55	0.39
gemini-2.5-pro-preview-03-25	0.99	0.92	0.91	0.96	0.91	0.87	0.89	0.84	0.76
o4-mini	0.97	0.97	0.97	0.92	0.94	0.91	0.95	0.92	0.89

Table 4: Text score, image score and group score for each category in the dataset.

Model	Proximity			Order			Absolute Position		
	ts	is	gs	ts	is	gs	ts	is	gs
paligemma-3b-mix-448	0.04	0.00	0.00	0.13	0.07	0.02	0.40	0.60	0.20
qwen-2.5-vl-72b-instruct	0.44	0.04	0.00	0.37	0.00	0.00	0.60	0.00	0.00
qwen-vl-max	0.44	0.00	0.00	0.23	0.03	0.02	0.20	0.00	0.00
claude-3-7-sonnet-20250219	0.72	0.68	0.52	0.53	0.28	0.18	0.40	0.00	0.00
llama-4-maverick	0.64	0.64	0.52	0.25	0.37	0.15	0.80	0.40	0.40
gpt-4o	0.52	0.60	0.40	0.18	0.22	0.07	0.20	0.40	0.00
llama-4-maverick_cot	0.68	0.88	0.64	0.68	0.57	0.45	0.40	0.80	0.20
gpt-4o_cot	0.68	0.76	0.56	0.72	0.60	0.52	0.60	1.00	0.60
gemini-2.5-pro-preview-03-25	0.88	0.88	0.76	0.92	0.87	0.80	1.00	0.80	0.80
o4-mini	0.72	0.88	0.68	0.90	0.92	0.88	1.00	1.00	1.00

Table 5: Textscore, imagescore and groupscore for each category in the dataset

F Appendix: Evaluation Stability

To prove that the size of our benchmark is appropriate, we test the standard deviation of one model with poor performance (gpt4o without chain-of-thought) and one model with good performance (gemini 2.5 pro). We run each model three times and then randomly sample subsets of size 0.5 to 1.0 of the dataset and provide their mean and standard deviation below. RocketScience yields stable evaluation results and would even do so if it were much smaller.

Share	Gpt-4o (Mean \pm Std)	Gemini 2.5 pro (Mean \pm Std)
0.5	0.21 \pm 0.03	0.86 \pm 0.01
0.6	0.20 \pm 0.03	0.85 \pm 0.01
0.7	0.20 \pm 0.03	0.85 \pm 0.01
0.8	0.19 \pm 0.02	0.84 \pm 0.01
0.9	0.19 \pm 0.02	0.85 \pm 0.02
1.0	0.18 \pm 0.02	0.86 \pm 0.02

Table 6: Performance of Gpt-4o and Gemini 2.5 pro over three runs each on random subsets of RocketScience. The standard deviation stays low at both half the dataset size and the full dataset.

G Appendix: Human Baseline

The full set of instructions given to the participants (apart from the consent form) is: "You will be asked to answer several questions. Each question will consist of two images and a caption, and you will need to click on the image that best matches the caption."

The testing interface can be seen in Figure 11 and Figure 12.

We do not see any significant risks for the study participants and we obtained permission for the human evaluation from University College Dublin's Human Research Ethics Committee.

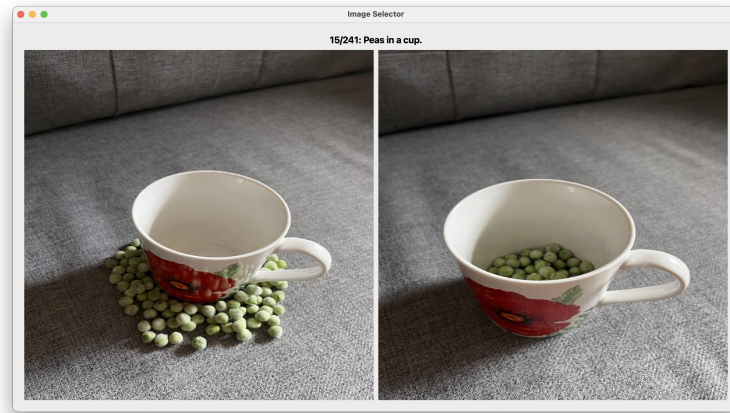


Figure 11: Human baseline interface with two images

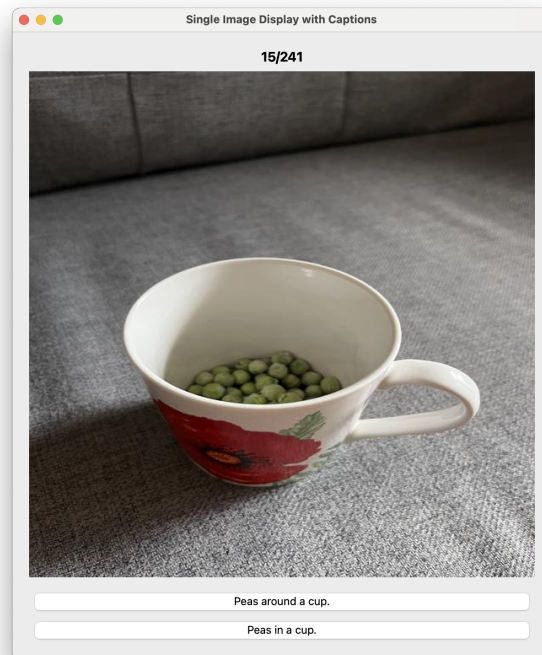


Figure 12: Human baseline interface with two captions

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: 1) we provide an overview of existing spatial relations benchmarks and their shortcomings in the related work section, 2) we propose RocketScience, a manually curated benchmark that addresses the shortcomings in the main section 3) We examine whether the localization or reasoning is more important for the success of CoT models in the results section.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We include a section on limitations and also discuss limitations throughout the rest of the text when necessary.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: We do not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The dataset is published on huggingface and the evaluation script is available on github via the link in the abstract, including instructions for how to run it. The model's full names are available in the code, but also stated in the results tables for full reproducibility.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The dataset is published on huggingface and the evaluation script is available on github via the link in the abstract, including instructions for how to run it. The model's full names are available in the code, but also stated in the results tables for full reproducibility.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We describe all hyperparameters used for evaluation of the models and they are also available in the code.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report mean and standard deviation in our experiment in table 3a). For our main results table we explain how we achieve reproducibility on open models and that we can only run API models once due to their cost.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We point out the type of GPU and rough execution times.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We have a section on ethics where we discuss how we avoid potential issues.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss societal impacts in the ethics section.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our dataset does not contain personal information and we are not aware of high-risk scenarios posed by the release.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We cite all models used and for images that are not ours we cite the source and state the license of the image.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: We document the details of our dataset in the paper as well as on [huggingface](https://huggingface.co) and also submit a croissant file.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[Yes\]](#)

Justification: We provide full documentation of the human performance experiment in the appendix.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[Yes\]](#)

Justification: We obtained permission for this experiment by University College Dublin's Human Research Ethics Committee as a low-risk study.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. **Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core contribution is a benchmark which was manually created.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.