

Don't learn, Ground. Image Generation for Grounded NLI

Anonymous ACL submission

Abstract

We propose a zero-shot method for Natural Language Inference (NLI) that leverages multi-modal representations by grounding language in visual contexts. Our approach generates visual representations of premises using text-to-image models and performs inference by comparing these with textual hypotheses. The pipeline achieves an accuracy comparable to text-based NLI classifiers while offering additional transparency. Our findings suggest that grounding language in vision is a viable and effective strategy for advancing robust natural language understanding.

1 Introduction

Language models trained and fine-tuned on various textual tasks exhibit impressive performance, especially with more data. At the same time, the extent to which unimodal language models can truly represent meaning has been criticized (Bender and Koller, 2020; Bisk et al., 2020). Achieving *functional linguistic competence* requires models to take into account the relationship between language and the world, e.g., through visual or other perceptual channels (but see Pavlick, 2023; Mandelkern and Linzen, 2024, for different positions on this issue). Enhancing language models with multimodal capabilities has now become common (recent examples include Deitke et al., 2025; Peng et al., 2024; Li et al., 2025; Chen et al., 2025).

Natural Language Inference (NLI; aka Textual Entailment) is a case in point. This task is typically framed in terms of the relationship between a premise text p and a hypothesis text h : whether h follows from (is entailed by) p , contradicts it, or whether the relationship is neutral (Dagan et al., 2006; MacCartney, 2009). This narrow framing of inference makes NLI classifiers susceptible to learned biases that arise in text (e.g., McCoy et al., 2019).

Modern approaches to NLI often focus on fine-tuning pretrained language models. Despite their impressive quantitative performance, such NLI models have notable drawbacks. First, successful fine-tuning requires substantial computational resources and large datasets. Second, fine-tuned models often scale poorly to unseen data due to biases present in the training datasets (Nie et al., 2020; Gururangan et al., 2018; McCoy et al., 2019; Zgreabăn et al., 2025).

We introduce a framework that operationalizes a truth-conditional view of entailment by grounding the premise in a concrete situation: we first render a visual scene consistent with the premise, and then ask a vision-language model to judge whether the hypothesis holds in that scene. This design is in line with recent evidence that generated assistive images can boost chain-of-thought (CoT) reasoning (Zhou et al., 2025). As our pipeline is fully zero-shot, it avoids task-specific fine-tuning and lets us probe whether grounding can reduce reliance on textual heuristics. Our contributions are as follows:

- We evaluate a zero-shot, visually grounded NLI pipeline that uses generated images as intermediate representations.
- On SNLI *easy/hard* subsets, we show that grounded inference can approach strong text-only baselines when paired with high-fidelity image generators (Table 1).
- We quantify hypothesis-side shortcut behavior via a premise-ablation control (RoBERTa (abl.)) and compare sensitivity to dataset artifacts across text-only and grounded settings.
- We benchmark multiple image generators and show that a smaller, faster model (FLUX.1-schnell, Labs, 2024) can yield competitive grounded NLI accuracy, suggesting a practical route to scaling the method.

2 Related work

Natural Language Inference Large-scale datasets developed to study NLI include SICK (Bentivogli et al., 2016), SNLI (Bowman et al., 2015), and MultiNLI (Williams et al., 2018). With the advent of large-scale pretraining, the field has witnessed a steady increase in NLI performance on standard benchmarks. For example, RoBERTa (Zhuang et al., 2021) achieves 90.8% accuracy on SNLI. More recently, large language models can be deployed in a zero-shot or few-shot fashion (Brown et al., 2020), although their accuracy in this case remains below that of models fine-tuned on the task. For example, on SNLI, Mistral-7B (Jiang et al., 2023) has a reported accuracy of around 90%, while SOTA performance is achieved with a few-shot version of T5 (Raffel et al., 2020) further trained with synthetic data, reaching 94.7% (Banerjee et al., 2024).

Visually grounded inference A separate line of work investigates the role of visual grounding in reasoning. This includes Visual Question Answering (Antol et al., 2015; Goyal et al., 2017; Hudson and Manning, 2019; Acharya et al., 2019) and visual commonsense reasoning (Zellers et al., 2019; Park et al., 2020). Importantly, Zhou et al. (2025) showed that generated intermediary images (e.g., blueprints) can be used to guide CoT reasoning on complex tasks; we adapt a similar approach to NLI.

Conceptually close to the original definition of the NLI task, Vu et al. (2018) created a grounded version of SNLI by linking the premises to their original images in Flickr30k (Young et al., 2014). They found that the inclusion of visual information sometimes led to a change in the gold label for a premise-hypothesis pair, but also showed that models do not benefit significantly from the inclusion of images. Later approaches benefited more, but not by a large margin (Kiela et al., 2019; De et al., 2023). In a different vein, Suzuki et al. (2019) propose a logical formalism for entailment involving both images and texts.

Xie et al. (2019) presented a new visual-textual entailment task (VTE) and developed the SNLI-VE dataset, where the entailment relationship is defined purely between an image (which replaces the textual premise) and a hypothesis. As with Vu et al. (2018), it was observed that the grounding of image-hypothesis pairs can result in a change of label compared to the original, text-only pairs in SNLI (Do et al., 2021; Kayser et al., 2021). Reijtenbach et al. (2025) showed that images generated from SNLI premises are applicable for VTE. While our approach is similar, we do not consider images independently, but rather as intermediary representations. Their study also focused on the applicability of generated data, while we consider the *pro* and *contra* in more detail.

Our primary question is whether zero-shot NLI can be achieved by means of a visual representation of the text. To that end, we compare our implementation of this pipeline with text-only baselines. In addition, we consider two questions:

3 Experiments

1) Whether grounded NLI helps bypass some of the recognized problems in text-centric NLI, such as hypothesis-side heuristics. The latter are keywords and grammatical features in hypotheses (e.g., negation) that are spuriously correlated with NLI classes in SNLI (Gururangan et al., 2018; Geirhos et al., 2020).

2) How suitable are different models for image generation for the task of NLI grounding.

Data We evaluate the proposed approach using text and images from V-SNLI, (Vu et al., 2018). Due to the large size of the corpus, we restrict our analysis to a limited, manually analyzable subset. Specifically, to test our method’s robustness against heuristics, we draw on Gururangan et al. (2018)’s distinction between *hard* and *easy* subsets of SNLI data. For our own tests, we sample 300 hypotheses related to 100 premises from both subsets. When more than three hypotheses were associated with one premise, we discarded the surplus hypotheses. This resulted in 276 *easy* hypotheses related to 92 premises and 285 *hard* hypotheses related to 95 premises. In what follows, we refer to these as the *easy* and *hard* subsets.

Visual Representation To generate images, we use four text-to-image (TTI) models: Stable Diffusion XL (Podell et al., 2024), DALL-E 3 (Betker et al., 2024), FLUX.1-schnell (Labs, 2024), and Qwen Image (Wu et al., 2025). We produce three images from each model and report the average of three inference passes. In addition, we make use of images generated by Reijtenbach et al. (2025) using a fine-tuned version of Stable Diffusion 1.5 (Rombach et al., 2022); their dataset includes 1 image per premise *p*. As a sanity check, we also use the original Flickr images from V-SNLI.

Method		Easy	Hard
RoBERTa		98.9%	83.2% (-15.7)
RoBERTa (abl.)		52.9%	28.4% (-24.5)
gpt-4o-2024-05-13		91.2%	80.6% (-10.6)
Qwen-VL-72B-Instruct		89.0%	80.2% (-8.8)
Images+Qwen-VL	DSG		
Flickr	-	86.6%	73.3% (-13.3)
Qwen-Image	0.949	94.7%	82.1% (-12.6)
DALL-E 3	0.914	92.9%	79.8% (-13.1)
FLUX.1-schnell	0.884	91.2%	78.5% (-12.7)
SDXL	0.883	92%	78.4% (-13.6)
SD1.5	0.623	79.7%	60.7% (-19.0)

Table 1: Percentage accuracy of NLI methods on easy and hard data. All scores except Flickr and SD1.5 are the average of 3 runs. In parentheses: accuracy delta between the subsets. DSG (Cho et al., 2024) scores are reported for image generators (see Section 4)

Our choice of TTI models allows for a comparison between parameter-heavy, large language model (LLM)-enhanced models (Qwen Image, DALL-E 3) and lighter models with weaker text encoding (FLUX.1-schnell, Stable Diffusion).

Inference For visual inference, we utilize Qwen2.5-VL-72B-Instruct (Bai et al., 2025), an open-weight MLLM that offers a solid approximation of the current state of the art. In addition, we considered LLaVa-NeXT (Liu et al., 2024) as an alternative backbone, but left it out due to a high rate of formatting errors (8%). The model is queried with an image and a textual prompt that includes three hypotheses and instructs the model to produce NLI labels for all three (see Appendix A).

Baseline We employed three baselines to represent zero-shot and fine-tuned NLI classifiers. We use gpt-4o-2024-05-13 (Hurst et al., 2024) as a strong zero-shot text-only baseline. We also run Qwen2.5-VL-72B-Instruct on textual premises and hypotheses. The prompt is adjusted to accommodate texts. Lastly, we test fine-tuned RoBERTa (Zhuang et al., 2021) from sentence-transformers (Reimers and Gurevych, 2019).

Even though RoBERTa can compete with state-of-the-art models on SNLI, it is known to be prone to bias as a consequence of fine-tuning (Gururangan et al., 2018). To control for how much this affects RoBERTa, we additionally report an ablation: we substitute textual premises with uninformative text ("something is happening") that cannot entail or contradict most hypotheses and report the results. If these are notably above random, this may indicate that the model is using textual shortcuts.

3.1 Results

We report the results in Table 1. Overall, grounded zero-shot NLI attains reasonably high accuracy, showing that the method is generally applicable. With high-fidelity image generators, grounded inference approaches the performance of strong baselines: Qwen-Image+VL reaches 94.7% accuracy on the *easy* subset and 82.1% on the *hard* subset, comparable to GPT-4o and close to fine-tuned RoBERTa.

Concerning robustness to surface heuristics, the results do not provide definite evidence that grounding eliminates sensitivity to dataset artifacts. All models, independent of grounding, exhibit lower performance on the *hard* subset. Fine-tuned RoBERTa shows the largest accuracy drop (15.7%), suggesting a strong reliance on hypothesis-side heuristics; this is supported by the ablation experiment, in which RoBERTa remains above chance on the *easy* subset (52.9% vs. 33%). At the same time, GPT-4o and image-conditioned models show comparable *easy/hard* deltas ($\approx 12\%$), indicating that the *easy/hard* split reflects general dataset difficulty rather than hypothesis-only reasoning alone.

Nevertheless, the results indicate that visual representations consistently influence inference. Performance varies systematically with the source and fidelity of the images: higher-quality generators (Qwen-Image, DALL-E 3) yield stronger results than lower-fidelity ones (SD1.5), despite identical textual inputs. Notably, using original Flickr images results in lower accuracy than most generated images on both subsets, suggesting that visual properties such as resolution and salience affect grounded inference. Further evidence comes from error overlap analysis (Table 2). The limited overlap between errors made by textual models and those made by image-assisted inference suggests that grounded models are not driven by the same hypothesis-level features alone. Remarkably, this is more pronounced on the *hard* subset, where the ablated RoBERTa performs below chance level.

Finally, the comparison across image generators reveals that LLM-backed diffusion models produce images that support more accurate grounded inference; particularly, Qwen-Image achieves the strongest overall performance. At the same time, FLUX.1-schnell performs competitively despite its smaller size and reduced computational requirements, indicating that effective grounded NLI does not necessarily require the most resource-intensive

Method	Easy	Hard
RoBERTa	0% (3)	35.4% (48)
gpt-4o-2024-05-13	24% (25)	37.5% (55)
Qwen-VL-72B-Instruct	29% (31)	42.8% (56)

Table 2: Baselines: relative error overlap with Qwen-Image + Qwen-VL; in parentheses: total error count.

image generation models. With FLUX.1-schnell requiring the fewest diffusion steps (16 vs. ≈ 50), it could have sufficient throughput for real-time production of assistive images.

4 Discussion

The fact that Qwen2.5-VL-72B-Instruct yields higher accuracy on artificial data than on Flickr images, contrary to expectations, may challenge its adequacy for our task. We consider the explanation that the low resolution of Flickr images (500×500) may be too coarse for the model. Indeed, similarly sized SD1.5 images (512×512) also lead to low accuracy; by contrast, 1024×1024 images systematically yield higher accuracy. Our analysis of predictions further reinforces this view: we find that Flickr+VL predicts the Neutral label disproportionately often: its ratio in *hard* and *easy* subsets is 36.8% and 38.9%. In contrast, the proportions of Neutral in the gold labels are 24.6% and 31.7% respectively. For SD1.5+VL, the ratio is also above 33%. We supply the confusion matrices in Appendix B. We see this as an indication that low resolution prevents the model from distinguishing the supporting details, although the VL model appears adequate overall. Still, the resolution issue may have implications for further work with V-SNLI and its descendants, as it could explain why these datasets continue to be challenging for the most modern VL models (Pitta et al., 2025); likewise, artificial VTE data production (Reijtenbach et al., 2025) has to account for different resolutions being easy or difficult for vision models.

Generally, we distinguish two types of errors pertaining mostly to text-to-image: factual errors (1) and neutrality errors (2). The former include incorrect rendering of object counts, concept bleeding (Podell et al., 2024), etc. We estimate the impact of these errors on our pipeline using two measures. First, the average score on Davidsonian Scene Graphs (Cho et al., 2024), i.e., the proportion of correctly rendered atomic statements from 0 to 1 (see Appendix D). The score is computed with Qwen-2.5-VL-Instruct and verified manu-

ally. Second, the proportion of Contradiction labels. As factual errors often lead to the Entailment label flipping, the latter score can be a coarse proxy. Our conclusions are as follows: we find that faithfulness and contradiction ratio are negatively correlated (Pearson’s $r = -0.97$, $n = 10$, $p < 0.05$); more importantly, we also observe that the impact of factual errors is substantial in SD1.5, but much less so in the modern diffusion models with $\approx 5\%$ inaccuracies for Qwen-Image (see Table 1). This calls for caution when using TTI models, but also shows that some of them are fit for the task.

Neutrality errors occur when the image inadvertently overspecifies initially ambiguous details, an issue that complicated turning SNLI into a VTE dataset (Do et al., 2021; Kayser et al., 2021). Oftentimes, this problem is unavoidable: consider the premise "The dog about to catch a frisbee", for which image generators are certain to produce an outside scene, and the originally neutral hypothesis "The dog is outside". We argue that predicting Entailment here may be meaningful, as it can reflect a possible implicature: e.g., this example was rated as Entailment by 2 out of 5 SNLI labelers. With Qwen images, the pipeline agrees with 1+ annotators on 8 out of 10 misclassified Neutrals. In that sense, some of these predictions may be justified as a form of fuzzy entailment.

5 Conclusion

This paper explored the use of assistive images as intermediate representations for zero-shot Natural Language Inference. Our results show that visually grounded NLI is feasible when premise information can be rendered visually, achieving accuracy comparable to strong text-only baselines. By grounding inference in concrete visual situations, the proposed pipeline aligns with a truth-conditional view of entailment and offers increased transparency.

Our comparison of text-to-image models shows that while none of them is error-free, they can support accurate NLI judgments, and the impact of factual rendering errors appears limited for the strongest of them. Importantly, we find that grounded inference does not require the most computationally expensive models: lighter-weight systems such as FLUX.1-schnell yield competitive results, suggesting that image-assisted inference may be viable beyond small-scale or offline settings.

6 Limitations

Several limitations of the present study should be noted. First, text-to-image generation models are known to exhibit social, demographic, and representational biases, and they may fail to visualize certain statements due to safety filters or training artifacts. Although we encountered only one such case in our experimental data, these issues remain a general concern for image-assisted reasoning methods and warrant careful consideration.

Second, while we evaluated a diverse set of image generators and a strong VL model, our coverage is not exhaustive. Other architectures, prompting strategies, or vision–language models may yield different results, and a more systematic exploration of the design space could provide further insights.

Finally, scalability remains a practical consideration. Generating high-resolution images incurs considerable computational cost, which may limit throughput in real-time or large-scale deployments. Although our results indicate that cheaper and faster generators can still support accurate grounded inference, alternative approaches such as image retrieval or hybrid generation–retrieval pipelines may offer more efficient solutions. We leave these directions for future investigation.

In the present work, we used generative AI tools for spelling and grammar checks.

References

Manoj Acharya, Kushal Kafle, and Christopher Kanan. 2019. [TallyQA: Answering complex counting questions](#). In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*, Honolulu, Hawaii. Association for the Advancement of Artificial Intelligence. ArXiv: 1810.12440 ISSN: 2159-5399.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. [VQA: Visual Question Answering](#). In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV'15)*, pages 2425–2433, Santiago, Chile. IEEE. ArXiv: 1505.00468v1.

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, and 8 others. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Sourav Banerjee, Anush Mahajan, Ayushi Agarwal, and Eishkaran Singh. 2024. [First Train to Generate, then](#)

[Generate to Train: UnitedSynT5 for Few-Shot NLI](#). ArXiv:2412.09263 [cs] version: 2.

Emily M. Bender and Alexander Koller. 2020. [Climbing towards NLU: On meaning, form, and understanding in the age of data](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.

Luisa Bentivogli, Raffaella Bernardi, Marco Marelli, Stefano Menini, Marco Baroni, and Roberto Zamparelli. 2016. [SICK Through the SemEval Glasses](#). *Language Resources and Evaluation*, 50(1):95–124.

James Betker, Gabriel Goh, Li Jing, † TimBrooks, Jianfeng Wang, Linjie Li, † LongOuyang, † JuntangZhuang, † JoyceLee, † YufeiGuo, † WesamManassra, † PrafullaDhariwal, † CaseyChu, † YunxinJiao, and Aditya Ramesh. 2024. [Improving image generation with better captions](#).

Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. [Experience grounds language](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. 2025. Janus-Pro: Unified Multimodal Understanding and Generation with Data and Model Scaling.

Jaemin Cho, Yushi Hu, Jason Baldridge, Roopal Garg, Peter Anderson, Ranjay Krishna, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. 2024. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation. In *ICLR*.

Ido Dagan, Oren Glickman, Bernardo Magnini, and Ramat Gan. 2006. [The PASCAL Recognising Textual Entailment.pdf](#). In J. Quinonero-Candela, I. Dagan, B. Magnini, and F D’Alche-Buc, editors, *Machine Learning Challenges*, pages 177–190. Springer, Berlin and Heidelberg.

Arkadipta De, Maunendra Sankar Desarkar, and Asif Ekbal. 2023. [Towards improvement of grounded cross-lingual natural language inference with visio-textual attention](#). *Nat. Lang. Process. J.*, 4:100023.

466	Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, and 31 others. 2025. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 91–104.	
477	Virginie Do, Oana-Maria Camburu, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-SNLI-VE: Corrected Visual-Textual Entailment with Natural Language Explanations . ArXiv:2004.03744 [cs].	
481	Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. 2020. Shortcut learning in deep neural networks . <i>Nature Machine Intelligence</i> , 2(11):665–673. Number: 11 Publisher: Nature Publishing Group.	
487	Yash Goyal, Tejas Khot, Aishwarya Agrawal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering . In <i>Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR’17)</i> , pages 6904–6913. ArXiv: 1612.00837 ISSN: 15731405.	
495	Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.	
504	Drew A. Hudson and Christopher D. Manning. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering . <i>Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition</i> , 2019-June:6693–6702. ArXiv: 1902.09506v3 ISBN: 9781728132938.	
511	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. <i>arXiv preprint arXiv:2410.21276</i> .	
516	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L�lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth�e Lacroix, and William El Sayed. 2023. Mistral 7B . ArXiv:2310.06825 [cs].	
	Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-vil: A dataset and benchmark for natural language explanations in vision-language tasks . In <i>IEEE/CVF International Conference on Computer Vision</i> , pages 1224–1234. IEEE.	524 525 526 527 528 529 530
	Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, and Davide Testuggine. 2019. Supervised multimodal bitransformers for classifying images and text . <i>ArXiv</i> , abs/1909.02950.	531 532 533 534
	Black Forest Labs. 2024. Flux. https://github.com/black-forest-labs/flux .	535 536
	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2025. Llava-onevision: Easy visual task transfer . <i>Transactions on Machine Learning Research</i> . Accepted by TMLR.	537 538 539 540 541
	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. Llava-next: Improved reasoning, ocr, and world knowledge .	542 543 544
	Bill MacCartney. 2009. <i>Natural language inference</i> . Ph.D. thesis, Stanford University.	545 546
	Matthew Mandelkern and Tal Linzen. 2024. Do Language Models’ Words Refer? <i>Computational Linguistics</i> , 50(3):1191–1200.	547 548 549
	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	550 551 552 553 554 555
	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4885–4901, Online. Association for Computational Linguistics.	556 557 558 559 560 561 562
	Jae Sung Park, Chandra Bhagavatula, Roozbeh Motlaghi, Ali Farhadi, and Yejin Choi. 2020. Visual-COMET: Reasoning About the Dynamic Context of a Still Image . In <i>Proceedings of the European Conference on Computer Vision</i> , pages 508–524, Berlin and Heidelberg. Springer. ArXiv: 2004.10796 ISSN: 16113349.	563 564 565 566 567 568 569
	Ellie Pavlick. 2023. Symbols and grounding in large language models . <i>Philosophical Transactions of the Royal Society A</i> .	570 571 572
	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shao-han Huang, Shuming Ma, Qixiang Ye, and Furu Wei. 2024. Grounding multimodal large language models to the world . In <i>International Conference on Learning Representations (ICLR)</i> . ICLR 2024 poster.	573 574 575 576 577

578	Elena Pitta, Tom Kouwenhoven, and Tessa Verhoef.	Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin,	635
579	2025. Probing vision-language understanding	Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai Bai,	636
580	through the visual entailment task: promises and pit-	Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang,	637
581	falls . In <i>Proceedings of the 2nd LUHME Workshop</i> ,	Zekai Zhang, Zhengyi Wang, An Yang, Bowen Yu,	638
582	pages 74–83, Bologna, Italy.	Chen Cheng, Dayiheng Liu, Deqing Li, and 20 oth-	639
		ers. 2025. Qwen-image technical report . <i>Preprint</i> ,	640
		arXiv:2508.02324.	641
583	Dustin Podell, Zion English, Kyle Lacey, Andreas	Ning Xie, Farley Lai, Derek Doran, and Asim Kadav.	642
584	Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna,	2019. Visual Entailment: A Novel Task for Fine-	643
585	and Robin Rombach. 2024. Sdxl: Improving latent	Grained Image Understanding. <i>arXiv</i> , 1901.06706.	644
586	diffusion models for high-resolution image synthesis .	ArXiv: 1901.06706v1.	645
587	In <i>International Conference on Learning Representa-</i>		
588	<i>tions (ICLR)</i> . ICLR 2024 spotlight.		
589	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-	646
590	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	enmaier. 2014. From image descriptions to visual	647
591	Wei Li, and Peter J. Liu. 2020. Exploring the limits	denotations: New similarity metrics for semantic in-	648
592	of transfer learning with a unified text-to-text trans-	ference over event descriptions. <i>Transactions of the</i>	649
593	former. <i>Journal of Machine Learning Research</i> , 21:1–	<i>Association for Computational Linguistics</i> , 2:67–78.	650
594	67. ArXiv: 1910.10683.		
595	Rob Reijtenbach, Suzan Verberne, and Gijs Wijnholds.	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin	651
596	2025. Dataset creation for visual entailment using	Choi. 2019. From recognition to cognition: Visual	652
597	generative AI . In <i>Proceedings of the 5th Work-</i>	commonsense reasoning . In <i>Proceedings of the IEEE</i>	653
598	<i>shop on Natural Logic Meets Machine Learning</i>	<i>Computer Society Conference on Computer Vision</i>	654
599	<i>(NALOMA)</i> , pages 8–17, Bochum, Germany. Associ-	<i>and Pattern Recognition (CVPR'19)</i> , pages 6713–	655
600	ation for Computational Linguistics.	6724. ArXiv: 1811.10830 ISSN: 10636919.	656
601	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	Mădălina Zgreabăn, Tejaswini Deoskar, and Lasha	657
602	Sentence embeddings using siamese bert-networks .	Abzianidze. 2025. Merge: Minimal expression-	658
603	In <i>Proceedings of the 2019 Conference on Empirical</i>	replacement generalization test for natural language	659
604	<i>Methods in Natural Language Processing</i> . Associa-	inference. <i>arXiv preprint arXiv:2510.24295</i> .	660
605	tion for Computational Linguistics.		
606	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	Yiyang Zhou, Haoqin Tu, Zijun Wang, Zeyu Wang,	661
607	Patrick Esser, and Björn Ommer. 2022. High-	Niklas Muennighoff, Fan Nie, Chaorui Deng, Shen	662
608	resolution image synthesis with latent diffusion mod-	Yan, Haoqi Fan, Yejin Choi, and 1 others. 2025.	663
609	els. In <i>Proceedings of the IEEE/CVF conference</i>	When visualizing is the first step to reasoning: Mira,	664
610	<i>on computer vision and pattern recognition</i> , pages	a benchmark for visual chain-of-thought. In <i>NeurIPS</i>	665
611	10684–10695.	<i>2025 Workshop on Efficient Reasoning</i> .	666
612	Riko Suzuki, Hitomi Yanaka, Masashi Yoshikawa, Koji	Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. 2021. A	667
613	Mineshima, and Daisuke Bekki. 2019. Multimodal	robustly optimized BERT pre-training approach with	668
614	Logical Inference System for Visual-Textual Entail-	post-training . In <i>Proceedings of the 20th Chinese</i>	669
615	ment . In <i>Proceedings of the 57th Annual Meeting of</i>	<i>National Conference on Computational Linguistics</i> ,	670
616	<i>the Association for Computational Linguistics: Stu-</i>	pages 1218–1227, Huhhot, China. Chinese Informa-	671
617	<i>dent Research Workshop</i> , pages 386–392, Florence,	tion Processing Society of China.	672
618	Italy. Association for Computational Linguistics.		
619	Hoa Trong Vu, Claudio Greco, Aliia Erofeeva, Somayeh	A Prompts	673
620	Jafaritazehjan, Guido Linders, Marc Tanti, Alberto	Our final prompts result from systematic varying	674
621	Testoni, Raffaella Bernardi, and Albert Gatt. 2018.	of the initial prompt across several dimensions. We	675
622	Grounded textual entailment . In <i>Proceedings of the</i>	picked the combination yielding the highest accu-	676
623	<i>27th International Conference on Computational Lin-</i>	racy overall. The dimensions are as follows:	677
624	<i>guistics</i> , pages 2354–2368, Santa Fe, New Mexico,	– terms 'premise' and 'hypothesis' kept/replaced	678
625	USA. Association for Computational Linguistics.	with 'statement';	679
626	Adina Williams, Nikita Nangia, and Samuel Bowman.	– terms 'entailment', 'contradiction', 'neither'	680
627	2018. A Broad-Coverage Challenge Corpus for Sen-	kept/replaced with 'accurate', 'contradicting', 'nei-	681
628	tence Understanding through Inference . In <i>Proceed-</i>	ther';	682
629	<i>ings of the 2018 Conference of the North American</i>	– provided/omitted explanations of NLI classes;	683
630	<i>Chapter of the Association for Computational Lin-</i>	– provided/omitted detailed explanations of the	684
631	<i>guistics: Human Language Technologies, Volume</i>	Neutral class.	685
632	<i>1 (Long Papers)</i> , pages 1112–1122, New Orleans,		
633	Louisiana. Association for Computational Linguis-		
634	tics.		

686

A.1 Vision-language prompt

687

Question: With respect to the objects in the image, is the statement in square brackets a) entailment, b) contradiction, c) neutral?

688

689

Answer 'entailment' if the statements accurately describe the objects in the image.

691

Answer 'contradiction' if the statement contradicts the image.

692

693

Answer 'neutral' if the statement isn't a contradiction, but adds unverifiable details (hidden/off-screen facts, motifs, intentions, identities/relationships etc.).

694

Put your answers into angled brackets.

695

Statement 1: []

699

Statement 2: []

700

Statement 3: []

701

Answer 1 (entailment/contradiction/neutral): <...>

702

Answer 2 (entailment/contradiction/neutral): <...>

703

Answer 3 (entailment/contradiction/neutral): <...>

704

705

A.2 Language prompt

706

Question: With respect to the premise, is the statement in square brackets a) entailment, b) contradiction, c) neutral?

707

Answer 'entailment' if the statements accurately describe the premise.

708

Answer 'contradiction' if the statement contradicts the premise.

709

Answer 'neutral' if the statement isn't a contradiction, but adds unverifiable details (hidden/off-screen facts, motifs, intentions, identities/relationships etc.).

710

Put your answers into angled brackets.

711

Statement 1: []

712

Statement 2: []

713

Statement 3: []

714

Premise:

715

Answer 1 (entailment/contradiction/neutral): <...>

716

Answer 2 (entailment/contradiction/neutral): <...>

717

Answer 3 (entailment/contradiction/neutral): <...>

718

719

720

721

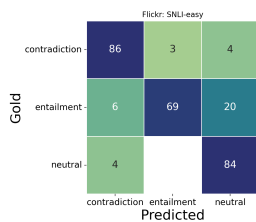
722

723

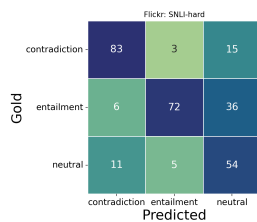
724

725

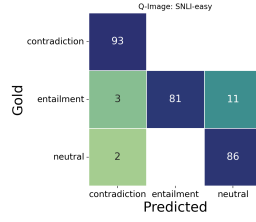
B Confusion Matrices



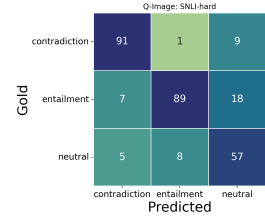
(a) Flickr + Qwen2.5-VL-72B-Instruct: *easy*



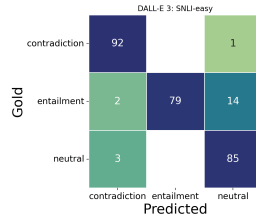
(b) Flickr + Qwen2.5-VL-72B-Instruct: *hard*



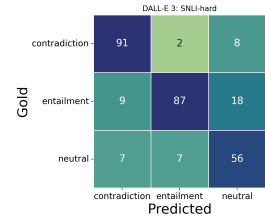
(a) Qwen-Image + Qwen2.5-VL-72B-Instruct: *easy*



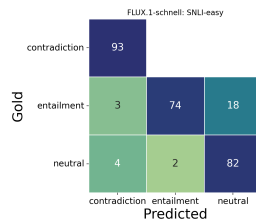
(b) Qwen-Image + Qwen2.5-VL-72B-Instruct: *hard*



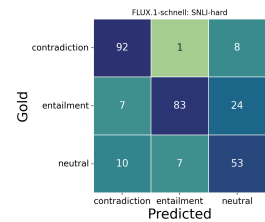
(a) DALL-E 3 + Qwen2.5-VL-72B-Instruct: *easy*



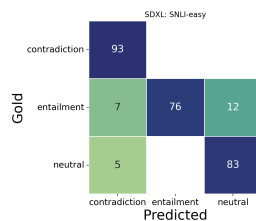
(b) DALL-E 3 + Qwen2.5-VL-72B-Instruct: *hard*



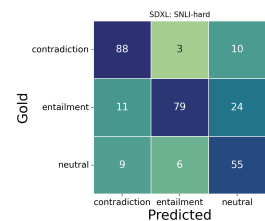
(a) FLUX.1-schnell + Qwen2.5-VL-72B-Instruct: *easy*



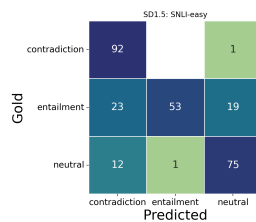
(b) FLUX.1-schnell + Qwen2.5-VL-72B-Instruct: *hard*



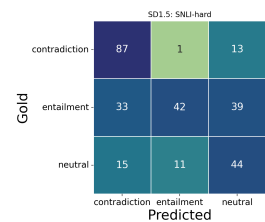
(a) SDXL + Qwen2.5-VL-72B-Instruct: *easy*



(b) SDXL + Qwen2.5-VL-72B-Instruct: *hard*



(a) SD1.5 + Qwen2.5-VL-72B-Instruct: *easy*



(b) SD1.5 + Qwen2.5-VL-72B-Instruct: *hard*

C Misclassified neutrals

Method	Agr. cases
Qwen-Image	8/10 (80%)
DALL-E 3	10/11 (90%)
FLUX.1-schnell	8/11 (73%)
SDXL	11/13 (85%)

Table 3: Proportion of neutral instances misclassified in all 3 inference runs, where TTI+VL pipeline agrees with 1+ SNLI annotators. We consider these instances as possible implicational readings.

D Davidsonian Scene Graphs

TTI evaluation with DSGs was proposed in [Cho et al. \(2024\)](#). The pipeline consists of three steps: 1) breaking the source text down into atomic statements (DSG tuples); 2) rewriting the tuples as questions; 3) answering the questions based on the generated image. Steps 1 and 2 can be performed with a text-only LLM, while 3 requires a VQA model. The answers can then be manually verified. Here, we provide the used prompts (same as in the original paper) as well as an evaluation example.

Tuple generation: Task: given input prompts, describe each scene with skill-specific tuples. Do not generate same tuples again. Do not generate tuples that are not explicitly described in the prompts.
output format: id | tuple

Question generation: Task: given input prompts and skill-specific tuples, re-write tuple each in natural language question.
output format: id | question

Question answering: Task: Answer all questions by number with 'yes' or 'no' based on the image; reply using tuple syntax. e.g., [(1, yes), (2, no)].

Example: Text: A dog jumps over the pole.

Tuples:

1 | entity - whole (dog)

2 | entity - whole (pole)

3 | action - (dog, jump)

4 | relation - spatial (dog, pole, over)

Questions:

1 | Is there a dog?

2 | Is there a pole?

3 | Is the dog jumping?

4 | Is the dog jumping over the pole?



Figure 7: Image: FLUX.1-schnell

Answers:

[(1, yes), (2, yes), (3, yes), (4, yes)] → score=1.0