

Making LVLMs Look Twice: Contrastive Decoding with Contrast Images

Anonymous ACL submission

Abstract

Large Vision-Language Models (LVLMs) are becoming increasingly popular for text-vision tasks requiring reasoning over both modalities, but often struggle with fine-grained visual discrimination. This limitation is evident in recent benchmarks like NaturalBench and D3, where closed models such as GPT-4o achieve only 39.6% accuracy, and open-source models perform below random chance (25%). We introduce Contrastive decoding with Contrast Images (CoCI), which adjusts LVLM outputs by contrasting them against outputs for similar images (Contrast Images - CIs). We first evaluate CoCI using naturally occurring CIs in benchmarks with curated image pairs, achieving improvements of up to 98.9% on NaturalBench, 69.5% on D3, and 37.6% on MMVP. For real-world applications where natural CIs are unavailable, we show that given sufficient training data, a lightweight neural classifier can effectively select CIs from similar images at inference time, improving NaturalBench performance by up to 36.8%. For scenarios lacking training data, we develop a caption-matching technique that selects CIs by comparing LVLM-generated descriptions of candidate images. Our method demonstrates the potential for improving LVLMs at inference time through different CI selection approaches, each suited to different data availability scenarios.

1 introduction

Large Vision-Language Models (LVLMs) are becoming increasingly popular for text-vision tasks that require reasoning over both modalities. However, they often struggle with fine-grained visual discrimination — that is, the ability to tell two similar yet distinct images apart — a crucial capability for real-world applications such as multimodal search, manufacturing, and robotics. Recent benchmarks have exposed this limitation: on NaturalBench (Li et al., 2024a), which tests vi-

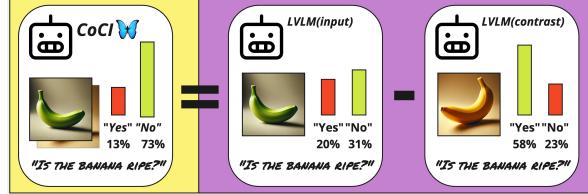


Figure 1: CoCI penalizes target image logits using those from a contrast image, weighted by hyperparameter α .

sual question answering over closely related images, state-of-the-art closed models like GPT-4o (OpenAI et al., 2024) achieve only 39.6% accuracy. Similarly, on the D3 benchmark (Gaur et al., 2024), which requires describing differences between paired images, open-source models perform below random chance (25%).

Efforts to address fine-grained visual discrimination in LVLMs are still under-explored. Current strategies addressing other LVLM shortcomings often rely on fine-tuning with specialized datasets (Wang et al., 2023; Chen et al., 2023; Liu et al., 2024a; Sarkar et al., 2024), multi-step correction pipelines (Yin et al., 2023; Zhou et al., 2023), or inference-time methods (Leng et al., 2023; Manevich and Tsarfaty, 2024; Liu et al., 2024b; Huang et al., 2023). Inference-time methods are particularly appealing as they do not require expensive model training and are less prone to compounding errors that can affect multi-step systems.

Building on the advantages of inference-time methods, we propose Contrastive decoding with Contrast Images (CoCI), an approach specifically designed to improve fine-grained visual discrimination in LVLMs. CoCI penalizes LVLM next-token probabilities with those obtained by feeding a different, contrasting image input (See Figure 1).

We evaluate CoCI across three different supervision regimes. First, using naturally occurring Contrast Images in curated benchmarks like NaturalBench, D3 and MMVP, we demonstrate improvements up to 98.9%, 69.5%, 37.6% respec-

tively. This establishes a performance ceiling for CoCI when ideal CIs are available. For applications where natural CIs are unavailable but training data exists, we show that a lightweight classifier can effectively select CIs from visually similar images at inference time, improving NaturalBench performance by up to 36.5%. In settings without training data, we propose a caption-matching technique that selects CIs at inference time by comparing LVLM-generated descriptions of candidate images.

Experiments with leading LVLMs — Qwen2-VL, LLaVA-OneVision, and Llama 3.2 (Wang et al., 2024a; Li et al., 2024b; Grattafiori et al., 2024) — establish the potential of contrastive decoding strategies with contrastive images for improved multimodal reasoning in real-world tasks.

2 Contrastive Decoding with Contrast Images (CoCI)

We present CoCI, a method to improve LVLM outputs by penalizing token probabilities that are likely under a contrast image. The choice of contrast image is crucial: e.g., when querying about fruit ripeness with an input image of an unripe banana, contrasting against an image of a ripe banana provides strong contrastive signal, while an image of a ripe pear offers weaker contrast and an image of a bus provides no useful signal and may degrade performance. This intuition guides our CI selection strategies across different scenarios. Before formalizing this intuition, we first review key concepts in LVLM text generation.

2.1 Preliminaries: Text Generation in LVLMs

LVLMs extend LLMs’ next-token prediction capability by conditioning on both text and images.¹ Generation (also called decoding) proceeds by sampling from next-token distributions, concatenating selected tokens with the context, and feeding this back into the model until an EOS token or length limit is reached. The LVLM next-token prediction formula is:

$$\text{LVLM}_t(y_{<t}, I) = P(y|y_{<t}, I) \quad \forall y \in \mathcal{V} \quad (1)$$

where $y_{<t}$ is the token prefix including both input and generated text up to position t , I is the input image, and \mathcal{V} is the model’s vocabulary.

¹In this work, we focus on the single image input case.

2.2 Contrastive Decoding

Since Li et al. (2023) introduced Contrastive Decoding, multiple variants have been explored (Sennrich et al., 2024; Jin et al., 2024; Phan et al., 2024), varying in weighting mechanisms, truncation strategies, and probability space operations. We implement CoCI based on Sennrich et al. (2024)’s minimal-hyperparameter variant as follows:

$$\begin{aligned} \text{CoCI}_t(y_{<t}, I, I') = \\ \log \left(P(y|y_{<t}, I) - \alpha P(y|y_{<t}, I') \right) \quad \forall y \in \mathcal{V} \end{aligned} \quad (2)$$

At each timestep t , CoCI penalizes token probabilities from distribution $P(y|y_{<t}, I)$ with those from $P(y|y_{<t}, I')$, where I is the target image and I' is the contrast image. The hyperparameter α controls the contrastive penalty strength.²

2.3 Obtaining Contrast Images

We propose three approaches for obtaining CIs:

Naturally occurring CIs. Many tasks naturally provide pairs of images that can serve as contrast images (CIs). For instance, a home assistant robot searching for “the blue ceramic mug with a chip on the handle” needs to distinguish between similar cups to find the exact match. We evaluate this scenario using LVLM benchmarks with curated image pairs designed to test fine-grained discrimination capabilities. These paired images serve as natural CIs in our experiments.

Classifier-obtained CIs. For cases without natural CIs, we train a lightweight classifier to select them at inference time. Given an LVLM L and training triplets $\langle q, I, I' \rangle$ where q is a binary question and I, I' are images with different ground-truth answers, we (a) Extract LVLM hidden states $h_{q,i} \in R^{d_L}$ for each image-question pair. (b) Concatenate states for image pairs to get $h_{q,i,i'} \in R^{2*d_L}$. (c) Create negative samples using the j least similar images from the top- k similar images to I in dataset D ³. (d) Train a three-layer MLP as a CI classifier⁴

²Throughout this work, we set $\alpha = 0.5$ for VQA tasks, and $\alpha = 0.8$ for open-ended generation tasks, without tuning.

³We set $j = 5$, $k = 100$ without tuning. For D , we use flickr30k (Young et al., 2014). We use open-clip with checkpoint "laion/CLIP-ViT-L-14-DataComp.XL-s13B-b90K" as the image encoder (Ilharco et al., 2021; Cherti et al., 2023; Radford et al., 2021a; Schuhmann et al., 2022). We use cosine-similarity as the similarity score throughout this work.

⁴See appendix A.1 and A.3 for PyTorch code, training setup and ablations.

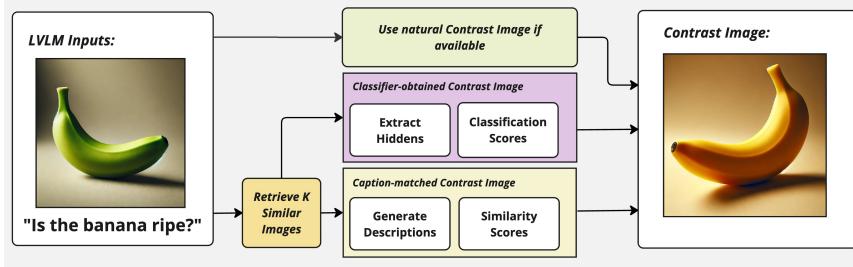


Figure 2: Illustration of the approaches we explore for obtaining a Contrast Image (CI).

For training, we use 60% of NaturalBench data, augmented with GPT-4o-generated question paraphrases. At inference, we select the CI that maximizes the classifier score among the k most similar images to the input.⁵

Caption-matched CIs. For scenarios without training data, we select CIs by comparing LVLM-generated image descriptions. Given an input image, we (a) Retrieve k similar images⁶. (b) Generate LVLM descriptions for all $k + 1$ images. (c) Embed descriptions using a text encoder. (d) Select the image whose description is most similar to the input image's description

2.4 Research Hypothesis

We test whether: (a) Contrastive decoding with CIs improves LVLM fine-grained reasoning, (b) A lightweight classifier trained on LVLM hidden states can effectively select CIs, and (c) Images with similar LVLM descriptions can serve as CIs.

3 Experiments

We evaluate CoCI with three leading LVLMs⁷, on three benchmarks testing fine-grained visual discrimination:

NaturalBench (Li et al., 2024a) tests discrimination between similar images through yes/no and multiple-choice questions, where answers differ between paired images. The benchmark contains 1900 image pairs with two questions per pair. We split the data into train (60%), dev (20%), and test (20%) sets. Metrics include image accuracy (both questions correct per image), question accuracy (per-question correctness), and group accuracy (G-Acc, requiring correct answers for all four image-question combinations).

⁵See table 2 for a comparison of using different k values at inference time. Our inference-time retrieval setup is identical to the one used in training (i.e. dataset, embedding model).

⁶We set $k = 5$ without tuning.

⁷See appendix A.2 for details on the checkpoints we used.

Model	Method	D3 (self-ret.)	MMVP (acc.)	NB (g-acc.)
Qwen2-VL	Baseline	30.8	46.0	30.8
	CoCI _{CAP}	34.8	48.7	31.3
	CoCI _{NAT}	52.2	63.3	46.6
LLaVA-OV	Baseline	25.1	52.7	28.2
	CoCI _{CAP}	31.6	57.3	31.6
	CoCI _{NAT}	38.1	66.7	56.1
Llama 3.2	Baseline	28.7	39.3	21.1
	CoCI _{CAP}	33.6	41.3	22.4
	CoCI _{NAT}	35.6	43.3	29.2

Table 1: CoCI performance comparison with provided CIs across three benchmarks, with natural CIs (CoCI_{NAT}) and caption-matched CIs (CoCI_{CAP}). Random baseline score for all metrics is 25.

MMVP (Multimodal Visual Patterns) (Tong et al., 2024) evaluates visual difference detection through multiple-choice questions on image pairs. Each pair differs in at least one visual aspect (e.g., object state, position, or orientation), with questions targeting these differences. The dataset comprises 150 image pairs with one question per pair. A model is considered correct only if it answers correctly on both images in a pair.

D3 (Detect, Describe, Discriminate) (Gaur et al., 2024) evaluates LVLMs' ability to generate discriminative descriptions that uniquely distinguish between similar images and consists of 247 image pairs. We adapt D3 to evaluate CoCI by reformulating it as a single-input task, where the model describes each image separately. We use the original self-retrieval evaluation protocol, which tests whether an image-text encoder can match the descriptions back to their respective images.

4 Results and Discussion

In Table 1 we can see that using natural CIs yields substantial improvements: up to 21.4 points on D3 (Qwen), 17.3 points on MMVP (LLaVA), and 27.9 points on NaturalBench (LLaVA). Caption-matched CIs show moderate but consistent gains,

Model	Method	Q-acc	I-acc	Acc	G-acc
Qwen2-VL	Baseline	55.3	59.3	76.8	30.8
	$Cl_{sk=4}$	55.5	58.8	76.4	32.1
	$Cl_{sk=8}$	56.3	58.9	76.7	32.4
	$Cl_{sk=16}$	57.4	60.1	77.2	33.7
	$Cl_{sk=32}$	57.8	60.1	77.4	34.2
	$Cl_{sk=64}$	58.2	60.8	77.9	33.9
LLaVA-OV	Baseline	53.8	56.1	74.6	28.2
	$Cl_{sk=4}$	59.2	59.6	77.6	35.3
	$Cl_{sk=8}$	57.8	60.1	77.5	34.5
	$Cl_{sk=16}$	57.6	58.7	77.0	33.4
	$Cl_{sk=32}$	60.3	62.1	78.5	38.4
	$Cl_{sk=64}$	59.7	62.1	78.2	37.6
Llama 3.2	Baseline	46.3	50.5	71.8	21.1
	$Cl_{sk=4}$	49.2	52.8	73.2	23.2
	$Cl_{sk=8}$	49.1	52.2	73.1	21.8
	$Cl_{sk=16}$	48.8	52.4	73.1	22.4
	$Cl_{sk=32}$	49.9	52.5	73.7	22.1
	$Cl_{sk=64}$	49.7	52.5	73.6	22.1

Table 2: CoCI accuracy metrics on the NaturalBench test set with CIs chosen using a lightweight classifier. $k = j$ denotes the classifier ran on the j most similar images to the input image.

particularly on D3 where LLaVA improves from 25.1% to 31.6%, suggesting that contrasting against images with similar captions effectively guides visual discrimination.

Table 2 shows that for Qwen and LLaVA, performance improves with larger candidate pools (k), peaking around $k=32$. Llama performs best with small pools ($k=4$), perhaps due to differences in its architecture affecting classifier effectiveness on its hidden states.

In NaturalBench, G-Acc shows particularly strong improvement with natural CIs (e.g., from 28.2% to 56.1% for LLaVA-OV) as it requires consistency across all image-question combinations. This pattern persists with classifier-selected CIs, where G-Acc improves by up to 10.2 points while other metrics show more modest gains.

The substantial performance gap between natural CIs and other methods indicates that while classifier-selected and caption-matched CIs provide improvements, they don't yet capture all aspects that make natural pairs effective.⁸

5 Related Work

Inference-time methods for enhancing multimodal reasoning. Recent work has focused on hallucination reduction through confidence-based probability adjustments (Huo et al., 2024), compar-

⁸See appendix A.3 for ablation tests with different CI selection strategies.

ison with semantic references (Yang et al., 2024), and contrastive decoding with perturbed inputs (Leng et al., 2023; Manevich and Tsarfaty, 2024). Our work extends these approaches to address fine-grained visual discrimination.

Alignment and grounding in LVLMs. Prior work has enhanced visual-textual alignment through object-level information synthesis (Wang et al., 2024b), targeted fine-tuning (Lu et al., 2024), and specialized dataset construction (Li et al., 2024c). While these methods improve foundational capabilities, they don't directly address fine-grained discrimination.

Contrastive examples in multimodal models. CLIP (Radford et al., 2021b) established contrastive learning between images and text for modality alignment. Recent works focus on leveraging contrast pairs: (Le et al., 2023) and (Zhang et al., 2024) generate synthetic contrast datasets using text-to-image models, while (Abbasnejad et al., 2020) and (Zhou et al., 2024) use contrastive examples to address dataset biases. Unlike these approaches that require data generation or model training, our method operates at inference time.⁹

6 Conclusion

We introduced Contrastive decoding with Contrast Images (CoCI), showing its effectiveness in improving LVLMs' fine-grained visual discrimination capabilities on multiple benchmarks. While naturally occurring contrast pairs yielded the strongest performance gains, both classifier-based and caption-matching approaches can provide meaningful improvements without requiring dataset curation, model training or extensive computational resources.

The effectiveness of our approach across different supervision regimes establishes CoCI as a practical solution for improving LVLM performance at inference time. Our results suggest that contrastive decoding strategies, when combined with appropriate contrast image selection, can enhance LVLMs' ability to make fine-grained visual distinctions, opening new avenues for improving multimodal reasoning through inference-time techniques.

⁹Classifier-selected CIs require minimal preprocessing compared to model finetuning or dataset curation.

7 Limitations

CoCI has several limitations worth noting. While we demonstrate its effectiveness with classifier-based and caption-matching approaches, the substantial performance gap between natural and automatically selected CIs indicates significant headroom for finding more effective contrast images. We tested simple selection methods to establish the viability of the approach, leaving the exploration of more sophisticated CI selection strategies to future work. Additionally, our evaluation focuses primarily on VQA and self-retrieval protocols; exploring additional evaluation methods could reveal other aspects of how CoCI affects LVLM generations.

The method introduces additional computation at inference time, running the LVLM twice per generation step and requiring CI selection overhead. While this aligns with the growing trend of leveraging test-time compute for improved performance, the current implementation could be optimized. Future work could explore more efficient implementations of contrastive decoding and investigate fusing operations like hidden state extraction with the generation procedure to reduce computational overhead.

Our implementation uses Flickr30k as the image database for CI selection - using larger, more diverse image collections could improve performance. Alternative image retrieval models and similarity scoring methods could also enhance CI selection. Additionally, our approach does not address cases where multiple contrasts might be informative - we only use a single contrast image, while some scenarios might benefit from multiple contrasting viewpoints.

The experiments use a fixed contrastive weight (α) across tasks within each category (VQA/generation). A more nuanced approach to setting this parameter, dynamically per sample or per token, based on the specific input or task, could yield better results.

While CoCI improves visual discrimination, it could potentially amplify biases present in contrast image databases or introduce new failure modes when inappropriate contrast images are selected. These risks should be carefully evaluated before deployment in sensitive applications.

Finally, our experiments focus exclusively on English-language benchmarks. Extending CoCI to multilingual settings and investigating how contrastive decoding approaches perform across differ-

ent languages represents an important direction for future research.

References

Ehsan Abbasnejad, Damien Teney, Amin Parvaneh, Javen Shi, and Anton van den Hengel. 2020. *Counterfactual vision and language learning*. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10041–10051.

Zhiyang Chen, Yousong Zhu, Yufei Zhan, Zhaowen Li, Chaoyang Zhao, Jinqiao Wang, and Ming Tang. 2023. *Mitigating hallucination in visual language models with visual supervision*.

Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829.

Manu Gaur, Darshan Singh S, and Makarand Tapaswi. 2024. *Detect, describe, discriminate: Moving beyond vqa for mllm evaluation*.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Isahan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-teng Jia, Kalyan Vasudevan Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal

394	Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazéri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madiam Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Prithish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	458
395		459
396		460
397		461
398		462
399		463
400		464
401		465
402		466
403		467
404		468
405		469
406		470
407		471
408		472
409		473
410		474
411		475
412		476
413		477
414		478
415		479
416		480
417		481
418		482
419		483
420		484
421		485
422		486
423		487
424		488
425		489
426		490
427		491
428		492
429		493
430		494
431		495
432		496
433		497
434		498
435		499
436		500
437		501
438		502
439		503
440		504
441		505
442		506
443		507
444		508
445		509
446		510
447		511
448		512
449		513
450		514
451		515
452		516
453		517
454		518
455		519
456		520
457		521

522	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wencheng Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaoqian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. <i>The llama 3 herd of models</i> .	579
523		580
524		581
525		582
526		583
527		
528	Zhaowei Li, Qi Xu, Dong Zhang, Hang Song, Yiqing Cai, Qi Qi, Ran Zhou, Junting Pan, Zefeng Li, Van Tu Vu, Zhida Huang, and Tao Wang. 2024c. <i>Groundinggpt:language enhanced multi-modal grounding model</i> .	584
529		585
530	Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2024a. <i>Mitigating hallucination in large multi-modal models via robust instruction tuning</i> .	586
531		587
532		
533	Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2024b. <i>Reducing hallucinations in vision-language models via latent space steering</i> .	588
534		589
535		590
536	Ilya Loshchilov and Frank Hutter. 2019. <i>Decoupled weight decay regularization</i> .	591
537		592
538	Junyu Lu, Dixin Zhang, Songxin Zhang, Zejian Xie, Zhuoyang Song, Cong Lin, Jiaxing Zhang, Bingyi Jing, and Pingjian Zhang. 2024. <i>Lyrics: Boosting fine-grained language-vision alignment and comprehension via semantic-aware visual objects</i> .	593
539		594
540	Avshalom Manevich and Reut Tsarfaty. 2024. <i>Mitigating hallucinations in large vision-language models (LVLMs) via language-contrastive decoding (LCD)</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 6008–6022, Bangkok, Thailand. Association for Computational Linguistics.	595
541		600
542		601
543		602
544		603
545	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 Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrej 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 Giertler, 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	604
546		605
547		606
548		607
549		608
550		609
551		610
552		611
553		612
554		613
555		614
556		615
557		616
558		617
559		618
560		619
561		620
562		621
563		622
564		623
565		624
566		625
567		626
568		627
569		628
570		629
571		630
572		631
573		632
574		633
575		634
576		635
577		636
578		637

639	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 Kvlichan, 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 Quinonero 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, Kavin 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, Lillian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Afak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, 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, Shiron Wu, Shuaiqi Xia, Sonia Phene, 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, Weiyi 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. 2024. <i>Gpt-4o system card</i> .	703
640		704
641		705
642		706
643		707
644		708
645		709
646		710
647		711
648		712
649		713
650		714
651		715
652		716
653		717
654		718
655		719
656		720
657		721
658		722
659		723
660		724
661		725
662		726
663		727
664		728
665		729
666		
667		
668		
669		
670		
671		
672		
673		
674		
675		
676		
677		
678		
679		
680		
681		
682		
683		
684		
685		
686		
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697		
698		
699		
700		
701		
702		
703	Phuc Phan, Hieu Tran, and Long Phan. 2024. <i>Distillation contrastive decoding: Improving llms reasoning with contrastive decoding and distillation</i> .	730
704		731
705		732
706		
707	Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021a. Learning transferable visual models from natural language supervision. In <i>ICML</i> .	733
708		734
709		735
710		736
711		737
712		738
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733	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. 2021b. Learning transferable visual models from natural language supervision.	739
734		740
735		741
736		742
737		743
738		744
739		
740		
741		
742		
743		
744		
745	Pritam Sarkar, Sayna Ebrahimi, Ali Etemad, Ahmad Beirami, Sercan Ö. Arik, and Tomas Pfister. 2024. <i>Data-augmented phrase-level alignment for mitigating object hallucination</i> .	745
746		746
747		747
748		748
749	Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. <i>LAION-5b: An open large-scale dataset for training next generation image-text models</i> . In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	749
750		750
751		751
752		752
753		753
754		754
755		755
756		756
757		757
758		758
759	Rico Sennrich, Jannis Vamvas, and Alireza Mohammadmashahi. 2024. <i>Mitigating hallucinations and off-target machine translation with source-contrastive and language-contrastive decoding</i> .	759
760		760
761		761
762		762

763 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma,
764 Yann LeCun, and Saining Xie. 2024. Eyes wide
765 shut? exploring the visual shortcomings of multi-
766 modal llms.

767 Lei Wang, Jiabang He, Shenshen Li, Ning Liu, and
768 Ee-Peng Lim. 2023. Mitigating fine-grained hallucin-
769 ation by fine-tuning large vision-language models
770 with caption rewrites.

771 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-
772 hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin
773 Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei
774 Du, Xuancheng Ren, Rui Men, Dayiheng Liu,
775 Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a.
776 Qwen2-vl: Enhancing vision-language model's per-
777 ception of the world at any resolution.

778 Wei Wang, Zhaowei Li, Qi Xu, Linfeng Li, YiQing
779 Cai, Botian Jiang, Hang Song, Xingcan Hu, Pengyu
780 Wang, and Li Xiao. 2024b. Advancing fine-grained
781 visual understanding with multi-scale alignment in
782 multi-modal models.

783 Dingchen Yang, Bowen Cao, Guang Chen, and
784 Changjun Jiang. 2024. Pensieve: Retrospect-then-
785 compare mitigates visual hallucination.

786 Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao
787 Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun,
788 and Enhong Chen. 2023. Woodpecker: Hallucination
789 correction for multimodal large language models.

790 Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-
791 enmaier. 2014. From image descriptions to visual
792 denotations: New similarity metrics for semantic in-
793 ference over event descriptions. *Transactions of the*
794 *Association for Computational Linguistics*, 2:67–78.

795 Jianrui Zhang, Mu Cai, Tengyang Xie, and Yong Jae
796 Lee. 2024. Countercurate: Enhancing physical and
797 semantic visio-linguistic compositional reasoning via
798 counterfactual examples.

799 Baohang Zhou, Ying Zhang, Kehui Song, Hongru Wang,
800 Yu Zhao, Xuhui Sui, and Xiaojie Yuan. 2024. MCIL:
801 Multimodal counterfactual instance learning for low-
802 resource entity-based multimodal information extrac-
803 tion. In *Proceedings of the 2024 Joint International*
804 *Conference on Computational Linguistics, Language*
805 *Resources and Evaluation (LREC-COLING 2024)*,
806 pages 11101–11110, Torino, Italia. ELRA and ICCL.

807 Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun
808 Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and
809 Huaxiu Yao. 2023. Analyzing and mitigating object
810 hallucination in large vision-language models. *ArXiv*,
811 abs/2310.00754.

812 A Appendix

813 A.1 Lightweight Classifier Implementation Details

814 Below is the PyTorch code of the lightweight classifier.

```
815 class Classifier(torch.nn.Module):
816     def __init__(self, input_dim: int):
817         super(Classifier, self).__init__()
818         # factor of 2 due to concatenation of target and candidate features
819         self.linear1 = torch.nn.Linear(input_dim * 2, input_dim)
820         self.linear2 = torch.nn.Linear(input_dim, input_dim)
821         self.linear3 = torch.nn.Linear(input_dim, 1)
822         self.dropout = torch.nn.Dropout(p=0.3)
823
824     def forward(self, x) -> torch.Tensor:
825         x = self.dropout(self.linear1(x))
826         x = F.relu(x)
827         x = self.dropout(self.linear2(x))
828         x = F.relu(x)
829         x = self.linear3(x)
830         return x
```

831 We trained a classifier per tested LVLM, all with the following parameters, using the AdamW ([Loshchilov and Hutter, 2019](#)) optimizer.

```
833 batch_size=256
834 num_epochs=13
835 learning_rate=3e-4
836 weight_decay=1e-6
```

837 A.2 LVLM Checkpoints Tested

838 The following are the LVLM checkpoints we tested CoCI with:

```
839 Qwen/Qwen2-VL-7B-Instruct
840 llava-hf/llava-onevision-qwen2-7b-ov-hf
841 meta-llama/Llama-3.2-11B-Vision-Instruct
```

A.3 Effect of Choosing a Contrast Image on NaturalBench Performance

842

Method	Setting	Q-acc	I-acc	Acc	G-acc
CoCI ablations	Baseline	51.6	55.4	75.1	25.6
	CI \leftarrow Random (out of top-5 most similar to input)	49.6	52.1	73.8	23.2
	CI \leftarrow Natural	71.8	70.8	84.3	51.6
	CI \leftarrow Most similar to input	49.7	52.5	73.6	23.9
	CI \leftarrow Most similar to Natural	60.3	60.7	78.9	35.0
	CI \leftarrow Least similar to Natural	46.7	48.9	72.6	21.8
Classifier	$k = 4$	51.7	54.3	74.5	26.6
	$k = 8$	53.0	55.4	75.3	26.6
	$k = 16$	54.3	56.8	76.1	29.2
	$k = 32$	52.2	54.6	75.1	25.8
	$k = 64$	51.8	53.9	74.7	26.3
	$k = 100$	52.1	54.1	74.8	25.5
$\text{Classifier}_{\text{+augmentations}}$	$k = 4$	52.0	54.3	74.6	27.1
	$k = 8$	52.8	55.9	75.0	27.9
	$k = 16$	54.5	57.8	76.1	29.2
	$k = 32$	54.9	58.2	75.9	30.0
	$k = 64$	54.7	57.9	76.1	30.3
	$k = 100$	54.7	58.0	76.1	30.0

Table 3: CoCI performance on the NaturalBench dev set with different CI selection methods, using Qwen2-VL. $\text{Classifier}_{\text{+augmentations}}$ indicates training data augmentation with GPT-4o paraphrased questions and standard image augmentations. Using natural CIs provides the strongest performance gains, with a 26-point improvement in group accuracy over baseline (51.6% vs 25.6%). Selecting CIs by similarity to natural CIs improves performance significantly (35.0% G-acc), while using the least similar images performs worse than baseline (21.8%), validating the importance of CI selection strategy. Random CI selection hurts performance (23.2% G-acc) even when restricted to similar images, highlighting that similarity alone is insufficient. Training with augmented data provides modest but consistent improvements across all metrics, with G-acc increasing by about 4 points compared to the non-augmented classifier. The augmented classifier also demonstrates more robust performance, maintaining consistent scores across different k values compared to the higher variance seen in the non-augmented version.