

# Mitigating Hallucinations in Large Vision-Language Models (LVLMs) via Language-Contrastive Decoding (LCD)

Anonymous ACL submission

## Abstract

Large Vision-Language Models (LVLMs) are an extension of Large Language Models (LLMs) that facilitate processing both image and text inputs, expanding AI capabilities. However, LVLMs struggle with object hallucinations due to their reliance on text cues and learned object co-occurrence biases. While most research quantifies these hallucinations, mitigation strategies are still lacking. Our study introduces a Language Contrastive Decoding (LCD) algorithm that adjusts LVLM outputs based on LLM distribution confidence levels, effectively reducing object hallucinations. We demonstrate the advantages of LCD in leading LVLMs, showing up to %4 improvement in POPE F1 scores and up to %36 reduction in CHAIR scores on the COCO validation set, while also improving captioning quality scores. Our method effectively improves LVLMs without needing complex post-processing or retraining, and is easily applicable to different models. Our findings highlight the potential of further exploration of LVLM-specific decoding algorithms for improved multimodal performance.

## 1 Introduction

Large Vision-Language Models (LVLMs) are a multimodal extension of Large Language Models (LLMs), transforming textual prompts and image inputs into text. However, they frequently produce object hallucinations, where absent objects are mentioned in the output (Yifan Li and Wen, 2023; Lovenia et al., 2023).

While hallucination-mitigation techniques in LLMs are actively researched, specific strategies for LVLMs are less developed. Current methods involve model-specific adjustments, additional training, or auxiliary models for post-hoc correction, and are often proven inefficient, costly, or limited by training data and model biases (Wang et al., 2023; Zhou et al., 2023; Gunjal et al., 2023; Yin



Figure 1: An illustration of LLM vs. LVLM token probabilities given an image and a text prefix mid-generation. According to the LLM, the word "dog" is much more likely to appear next. LCD dynamically contrasts these probabilities to mitigate language biases in LVLM outputs.

et al., 2023). Conversely, LVLM hallucination evaluation has seen progress with object hallucination benchmarks like NOPE (Lovenia et al., 2023) and POPE (Yifan Li and Wen, 2023), and recent works that aim for more holistic LVLM hallucination evaluation such as FaithScore (Jing et al., 2023) and HallusionBench (Guan et al., 2023).

A key reason for LVLM hallucinations is their tendency to over-rely on linguistic information, as was first observed by Guan et al. (2023). Based on this insight, we propose to intervene in the LVLM decoding phase so that model outputs are less informed by language biases. Specifically, we propose to use Contrastive Decoding (Li et al., 2023a;

O’Brien and Lewis, 2023) to alter LVLM output probabilities with respect to the internal LLM’s probabilities, guided by a dynamic weighting mechanism based on the LLM distribution’s entropy.

Our experiments show that our proposed method, Language Contrastive Decoding (LCD), improves hallucination scores on POPE (Yifan Li and Wen, 2023) and CHAIR (Rohrbach et al., 2018) on InstructBLIP variants based on Vicuna and Flan-T5 (Dai et al., 2023), and LLAVA 1.5 (Liu et al., 2023). We assess LCD’s overall generation quality by reporting captioning metrics and conducting a GPT4-V (OpenAI et al., 2023) assisted evaluation. LCD, as a decoding strategy, can be applied to other models without additional training or output modifications, emphasizing its utility for broader LVLM use.

The contributions of this paper are thus manifold. First, we introduce a novel decoding method tailored for LVLMs to mitigate object hallucinations. Next, we present a dynamic weighting strategy based on entropy which is applicable in various CD scenarios. Finally, we share our dataset and code, including a user-friendly LCD implementation via Huggingface, to encourage further research into LVLM-specific decoding strategies, a promising avenue for future research.

## 2 Motivation and Background

The integration of vision capabilities into Large Language Models has led to the development of Large Vision-Language Models, merging LLMs’ textual understanding with vision-text encoders. This trend towards multimodal systems is exemplified in commercial platforms such as GPT4-V (OpenAI et al., 2023) and Google’s Gemini (Team et al., 2023).

**Large Vision-Language Models** combine LLMs and vision-text encoders to generate text from textual prompts and visual inputs. An LVLM generally comprises three main components: a vision-text encoder like CLIP (Radford et al., 2021), an LLM such as LLAMA (Touvron et al., 2023) or Flan-T5 (Chung et al., 2022), and a cross-modal alignment module linking the vision-text encoder output with the LLM.

Initially, LVLMs were fine-tuned for specific tasks (Li et al., 2022; Wang et al., 2022). However, advancements in LLMs have led to a shift towards general-purpose, instruction-tuned LVLMs. These models are designed to handle a wide range of tasks

based on instructions, making them more versatile. Despite these advancements, LVLMs grapple with hallucinations of different types.

**LVLMs Hallucinations and their Mitigation** Hallucinations in LVLMs, particularly object hallucinations where nonexistent entities are mentioned, are often attributed to LVLMs’ reliance on spurious correlations and language biases, as demonstrated by Li et al. (2023b) and Zhou et al. (2023). Moreover, Guan et al. (2023) highlight LVLMs’ tendency to prioritize language over visual data, leading to hallucinations.

Mitigation strategies proposed by Gunjal et al. (2023) and Wang et al. (2023) involve further model training with augmented datasets or reward models. Zhou et al. (2023); Yin et al. (2023) developed auxiliary models to correct outputs post-generation. These solutions often require dataset-specific work or additional model training, potentially leading to overfitting or new biases, and are not easily transferable across LVLMs. In this work we aim to address these shortcomings with a more general, modular approach.

## 3 Language Contrastive Decoding (LCD)

Before presenting LCD, we briefly introduce the essentials of decoding in LVLMs 3.1, followed by our formal proposal 3.2 and research hypothesis 3.3.

### 3.1 Decoding Techniques and Contrastive Decoding: Essential Preliminaries

Decoding in auto-regressive generative models is the stage that transforms an input representation into a sequence of output tokens. In LVLMs, this process involves a model  $M$ , an image  $I$ , a textual prompt  $X$ , and a particular timestamp  $t$  during generation. It can be described as a series of selections from the model’s probability distribution, producing a token sequence  $T$ , as formalized in eq. (1).

$$T_t \sim P(\cdot | I, X, T_{<t}; M) \quad (1)$$

Greedy decoding, selecting the most probable token at each step (or the top  $k$  tokens in a beam search with  $k$  beams), is the simplest approach. However, high likelihood sequences do not necessarily align with human preferences, leading to the “likelihood trap” (Zhang et al., 2021). This has led to the use of sampling-based methods, such as top-k sampling, nucleus sampling (Holtzman et al.,

2020), and locally typical sampling (Meister et al., 2023), which either truncate the set of candidate tokens or adjust the model’s distribution, e.g. through temperature scaling.

Contrastive Decoding (CD) has been introduced for LLMs as a method to penalize the outputs of an expert model with those from a less powerful model (Li et al., 2023a). CD can be applied to any two probability distributions with the same support and has been adapted as a sampling strategy, improving various text generation tasks (O’Brien and Lewis, 2023; Chuang et al., 2023; Sennrich et al., 2024). CD uses both truncation and reshaping of probability distributions. The truncation phase ("adaptive plausibility") is described by eq. (2), where  $\alpha$  is a hyper-parameter,  $\mathcal{V}$  and  $\mathcal{V}'$  are the original and truncated token vocabularies at time  $t$ , and  $P$  is the conditional distribution on the prefix  $T_{<t}$ .

$$\mathcal{V}' = \{v \in \mathcal{V} : P(v|T_{<t}) \geq \alpha \max_w P(w|T_{<t})\} \quad (2)$$

Finally, the formula for CD, as suggested by O’Brien and Lewis (2023), given here generally for two conditional distributions  $P$  and  $P'$  on variable  $x$  with the same support, conditioned on  $X$  is presented in eq. (3).

$$CD_t(x, X, P, P') = \begin{cases} (1 + \beta) \log P(x|X) - \beta \log P'(x|X), & \text{if } x \in \mathcal{V}' \\ -\infty, & \text{otherwise} \end{cases} \quad (3)$$

$\beta$  is a fixed weight hyper-parameter. Our proposed method, detailed shortly, alters CD by introducing an entropy-based dynamic weighting scheme.

### 3.2 Proposed Method

Our intuition, based on previous findings by (Guan et al., 2023; Rohrbach et al., 2018; Yifan Li and Wen, 2023), is that an LVLM can be "misled" by its constituent LLM during the generation process.

Consider for example an LVLM that is describing an image (see illustration 1). Mid-generation, given the text "An image of a man walking his," it may predict "dog" due to language biases, even if it is a bear shown that is actually shown. A 'plain' LLM, without seeing the image, reinforces these biases by highly rating "dog". Our method builds on this insight to guide an LVLM towards more accurate predictions using Contrastive Decoding.

Our method operates as follows: At each generation step  $t$ , for each token  $x$ , we first determine the next-token probabilities from the LVLM,  $P_{LVLM}$ , based on the current token sequence  $T_{<t}$ , text  $X$ , and image  $I$ . We then obtain a second distribution,  $P_{LLM}$ , by inputting all data except the image into the LLM. The LLM’s conditional entropy  $H_{LLM}$  informs the dynamic weight as per eq.4. We then adjust token  $x$ ’s logits using the LCD formula in eq. 5.

$$\beta_t = \frac{\beta}{H_{LLM}(x|X, T_{<t})} \quad (4)$$

$$LCD_t(x, T_{<t}, I, P_{LVLM}, P_{LLM}) = (1 + \beta_t) \log P_{LVLM}(x|I, X, T_{<t}) - \beta_t \log P_{LLM}(x|X, T_{<t}) \quad (5)$$

We use the LCD logits for sampling, with temperature set to 1.0 in all experiments unless stated otherwise.

### 3.3 Research Hypothesis

Our hypothesis is that contrasting LVLM outputs with LLM outputs conditioned only on the textual data, can mitigate language biases, therefore reducing hallucinations in LVLMs.

## 4 Experiments and Results

We set out to assess the effect of LCD on object hallucinations in LVLM outputs against popular decoding settings. Additionally, we verify that LCD does not degrade output quality. To this end, we asses LCD on the POPE benchmark (Yifan Li and Wen, 2023), and on an image detailed-description task where we report hallucination and captioning metrics and conduct a GPT4-V assisted evaluation.

### Polling-based Object-Probing Evaluation

POPE consists of object-presence binary questions on 500 COCO dataset images (Lin et al., 2015), with questions equally divided between present and absent objects. Absent objects are chosen based on three criteria: *random*, *popular* (common in COCO), and *adversarial* (commonly co-occurring with present objects). POPE’s drawback is its one-word response structure, which limits the influence of decoding strategies and does not evaluate open-ended generation capabilities.

Model	Method	METEOR $\uparrow$	WMD $\uparrow$	ROUGE_L $\uparrow$	Acc $\uparrow$	Det $\uparrow$	CHAIRs $\downarrow$	CHAIRi $\downarrow$
InstructBLIP $_F$	Baseline	.157	.367	.161	4.92	4.02	.662	.146
	LCD	<b>.159</b>	<b>.370</b>	<b>.168</b>	<b>5.4</b>	4.01	<b>.566</b>	<b>.131</b>
InstructBLIP $_V$	Baseline	.178	.423	.291	3.7	3.51	.274	.126
	LCD	<b>.199</b>	<b>.48</b>	<b>.38</b>	<b>4.59</b>	<b>3.83</b>	<b>.174</b>	<b>.107</b>
LLAVA 1.5	Baseline	.163	<b>.357</b>	.169	4.77	4.56	.672	.182
	LCD	<b>.171</b>	.352	<b>.181</b>	<b>5.39</b>	4.54	<b>.610</b>	<b>.161</b>

Table 1: Image Description results.  $F$  and  $V$  stand for the Flan-T5 and Vicuna. Acc and Det are mean GPT4-V scores for Accuracy and Detailedness. METEOR, WMD and ROUGE $_L$  are popular captioning metrics (Kusner et al., 2015; Banerjee and Lavie, 2005; Lin, 2004).

**Image Detailed-Descriptions** To complement POPE, we introduce a long-form text generation task, inspired by findings from Zhou et al. (2023), that more extensive context increases hallucinations. The task uses the same COCO images as POPE, for which we generate detailed descriptions. The prompts we use are detailed in appendix A.1. This task tests LCD in a scenario where the propensity to hallucinate is high, and its longer outputs facilitate qualitative evaluation.

**Baselines and Metrics** For POPE, we use sampling as the baseline and report F1 scores.<sup>1</sup> For the descriptions task, we use the popular nucleus sampling algorithm<sup>2</sup> and report CHAIR metrics (Rohrbach et al., 2018). To assess description quality, we use captioning metrics against COCO’s gold captions, which serve as an approximation considering length differences. Additionally, following Yin et al. (2023), we use GPT4-V to evaluate the descriptions for Detailedness and Accuracy (details provided in Appendix A.1).

**Models** We conduct our experiments with leading LVLMS: two versions of the InstructBLIP model (with Flan-T5 and Vicuna LLMs) and LLAVA 1.5. The complete experimental details, such as exact model variants and generation hyperparameters, are given in the appendix.

## 5 Results and Discussion

In POPE, LCD improves F1 scores across 8 of 9 configurations (table 2). In the image description task (table 1), it significantly reduces hallucinations at both sentence and instance levels in all three models. However, CHAIR numbers are still high (especially in the Vicuna based models, InstructBLIP $_V$  and LLAVA), demonstrating the

<sup>1</sup>Complete POPE results are in the appendix, table 4

<sup>2</sup>We find that nucleus-sampling gives better results than vanilla sampling (see table 3 in the appendix for ablations)

POPE	Model	Baseline F1	LCD F1
Random	InstructBLIP $_V$	83.95	<b>87.55</b>
Popular		82.80	<b>84.34</b>
Adversarial		80.25	<b>81.64</b>
Random	InstructBLIP $_F$	84.05	<b>84.27</b>
Popular		80.74	<b>82.81</b>
Adversarial		78.87	<b>80.69</b>
Random	LLAVA 1.5	<b>84.17</b>	83.76
Popular		83.10	<b>83.47</b>
Adversarial		81.34	<b>81.62</b>

Table 2: POPE results for different models and methods.

prevalence of object hallucinations in long form LVLMS outputs. Notably, InstructBLIP $_V$  outperforms LLAVA 1.5 in most assessed metrics, in spite of being considered a weaker model. LCD shows particular effectiveness in InstructBLIP models probably due to their LLMs being frozen during training, which lends them a stronger language bias. In terms of overall generation quality as reflected by the captioning metrics, LCD is better in all cases but one (WMD in LLAVA 1.5, where the reduction is  $\sim 1\%$ ). In the GPT4-V eval, LCD improves Accuracy in all cases, and is as detailed as the baseline, suggesting it reduces hallucinations but does not increase the granularity of the descriptions.

## 6 Conclusion

In this paper we present Language Contrastive Decoding, a novel method to reduce hallucinations in LVLMS. By dynamically adjusting output probabilities using the LVLMS’s internal LLM, LCD significantly improves hallucination metrics across different LVLMS architectures, enhancing the quality and reliability of generated content without necessitating retraining or auxiliary models and post-processing. This work highlights the potential of specialized decoding strategies in enhancing multimodal AI models and lays the groundwork for further exploration into more sophisticated LVLMS decoding methods.

## 7 Limitations

Firstly, while LCD shows promise in reducing hallucinations, it only targets hallucinations caused by language biases, but hallucinations can arise from other sources. For instance, previous work has shown that some hallucinations are caused by poor visual understanding (Guan et al., 2023). We believe LCD can be used as a platform to craft LVLM-specific decoding algorithms that would mitigate hallucinations stemming from different factors, and leave this pursuit for future work.

Secondly, our evaluation method primarily addresses object hallucinations, which are only one form of hallucination that LVLMs may exhibit. Preliminary results signal that LCD mitigates more complex manifestations of language-induced hallucinations as assessed by recent benchmarks such as FAITHSCORE (Jing et al., 2023) and HallusionBench (Guan et al., 2023), but further work is required to establish this.

Moreover, LCD relies on current LVLM architectures that combine an LLM and a text-vision encoder, and requires access to an LLM that emits output probabilities on the same set of tokens as the LVLM. It is possible that the future generation of multimodal AI systems will have a different architecture that will make LCD obsolete. Additionally, LCD requires an LLM forward pass for each LVLM decoding step. The added latency could be mitigated with efficient inference techniques, and also by using a smaller LLM as the contrasting model. The effectiveness of LCD in this scenario is left for future work.

Finally, there are ethical considerations related to the mitigation of hallucinations in LVLMs. As these models become more reliable, it is crucial to continue evaluating the potential impacts of their use, ensuring they do not perpetuate or exacerbate biases present in their training data. LCD indeed mitigates some biases, but it is important to keep in mind that it might amplify other biases, unknown to us. Responsible deployment of these models requires ongoing vigilance and a commitment to transparency and fairness.

## References

Satanjeev Banerjee and Alon Lavie. 2005. **METEOR: An automatic metric for MT evaluation with improved correlation with human judgments**. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Transla-*

*tion and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. 2023. **Dola: Decoding by contrasting layers improves factuality in large language models**.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. **Scaling instruction-finetuned language models**.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. **Instructblip: Towards general-purpose vision-language models with instruction tuning**.

Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023. **Hallusionbench: An advanced diagnostic suite for entangled language hallucination visual illusion in large vision-language models**.

Anisha Gunjal, Jihan Yin, and Erhan Bas. 2023. **Detecting and preventing hallucinations in large vision language models**.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. **The curious case of neural text degeneration**.

Liqliang Jing, Ruosen Li, Yunmo Chen, Mengzhao Jia, and Xinya Du. 2023. **Faithscore: Evaluating hallucinations in large vision-language models**.

M. J. Kusner, Y. Sun, N. I. Kolkin, and K. Q. Weinberger. 2015. From word embeddings to document distances. In *ICML*.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. **Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation**.

Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. 2023a. **Contrastive decoding: Open-ended text generation as optimization**.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. **Evaluating object hallucination in large vision-language models**.

403	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	462
404		463
405		464
406		465
407	Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. <a href="#">Microsoft coco: Common objects in context</a> .	466
408		467
409		468
410		469
411		470
412	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. <a href="#">Improved baselines with visual instruction tuning</a> .	471
413		472
414		473
415	Holy Lovenia, Wenliang Dai, Samuel Cahyawijaya, Ziwei Ji, and Pascale Fung. 2023. <a href="#">Negative object presence evaluation (nope) to measure object hallucination in vision-language models</a> .	474
416		475
417		476
418		477
419	Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2023. <a href="#">Locally typical sampling</a> .	478
420		479
421	Sean O’Brien and Mike Lewis. 2023. <a href="#">Contrastive decoding improves reasoning in large language models</a> .	480
422		481
423	OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madeleine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal	482
424		483
425		484
426		485
427		486
428		487
429		488
430		489
431		490
432		491
433		492
434		493
435		494
436		495
437		496
438		497
439		498
440		499
441		500
442		501
443		502
444		503
445		504
446		505
447		506
448		507
449		508
450		509
451		510
452		511
453		512
454		513
455		514
456		515
457		516
458		517
459		518
460		519
461		520
		521
		522
		523

524	Gemini Team, Rohan Anil, Sebastian Borgeaud,	Anirudh Baddepudi, Alex Goldin, Adnan Ozturk,	588
525	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu	Albin Cassirer, Yunhan Xu, Daniel Sohn, Deven-	589
526	Soricut, Johan Schalkwyk, Andrew M. Dai, Anja	dra Sachan, Reinald Kim Amplayo, Craig Swan-	590
527	Hauth, Katie Millican, David Silver, Slav Petrov,	son, Dessie Petrova, Shashi Narayan, Arthur Guez,	591
528	Melvin Johnson, Ioannis Antonoglou, Julian Schrit-	Siddhartha Brahma, Jessica Landon, Miteyan Patel,	592
529	twieser, Amelia Glaese, Jilin Chen, Emily Pitler,	Ruizhe Zhao, Kevin Vilella, Luyu Wang, Wenhao	593
530	Timothy Lillicrap, Angeliki Lazaridou, Orhan Fir-	Jia, Matthew Rahtz, Mai Giménez, Legg Yeung,	594
531	rat, James Molloy, Michael Isard, Paul R. Barham,	Hanzhao Lin, James Keeling, Petko Georgiev, Di-	595
532	Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm	ana Mincu, Boxi Wu, Salem Haykal, Rachel Sapu-	596
533	Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins,	tro, Kiran Vodrahalli, James Qin, Zeynep Cankara,	597
534	Clemens Meyer, Eliza Rutherford, Erica Moreira,	Abhanshu Sharma, Nick Fernando, Will Hawkins,	598
535	Kareem Ayoub, Megha Goel, George Tucker, En-	Behnam Neyshabur, Solomon Kim, Adrian Hut-	599
536	rique Piqueras, Maxim Krikun, Iain Barr, Nikolay	ter, Priyanka Agrawal, Alex Castro-Ros, George	600
537	Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White,	van den Driessche, Tao Wang, Fan Yang, Shuo yiin	601
538	Anders Andreassen, Tamara von Glehn, Lakshman	Chang, Paul Komarek, Ross McIlroy, Mario Lučić,	602
539	Yagati, Mehran Kazemi, Lucas Gonzalez, Misha	Guodong Zhang, Wael Farhan, Michael Sharman,	603
540	Khalman, Jakub Sygnowski, Alexandre Frechette,	Paul Natsev, Paul Michel, Yong Cheng, Yamini	604
541	Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan,	Bansal, Siyuan Qiao, Kris Cao, Siamak Shakeri,	605
542	Xi Chen, James Lottes, Nathan Schucher, Federico	Christina Butterfield, Justin Chung, Paul Kishan	606
543	Lebron, Alban Rustemi, Natalie Clay, Phil Crone,	Rubenstein, Shivani Agrawal, Arthur Mensch, Kedar	607
544	Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu,	Soparkar, Karel Lenc, Timothy Chung, Aedan Pope,	608
545	Heidi Howard, Adam Bloniarz, Jack W. Rae, Han	Loren Maggiore, Jackie Kay, Priya Jhakra, Shibo	609
546	Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober,	Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,	610
547	Dan Garrette, Megan Barnes, Shantanu Thakoor, Ja-	Andrea Tacchetti, Maja Trebacz, Kevin Robinson,	611
548	cob Austin, Gabriel Barth-Maron, William Wong,	Yash Katariya, Sebastian Riedel, Paige Bailey, Ke-	612
549	Rishabh Joshi, Rahma Chaabouni, Deeni Fatiha,	fan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose	613
550	Arun Ahuja, Ruibo Liu, Yunxuan Li, Sarah Cogan,	Slone, Neil Houlsby, Xuehan Xiong, Zhen Yang,	614
551	Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang,	Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa	615
552	Jordan Grimstad, Ale Jakse Hartman, Martin Chad-	Lee, Music Li, Thais Kagohara, Jay Pavagadhi, So-	616
553	wick, Gaurav Singh Tomar, Xavier Garcia, Evan	ophie Bridgers, Anna Bortsova, Sanjay Ghemawat,	617
554	Senter, Emanuel Taropa, Thanumalayan Sankara-	Zafarali Ahmed, Tianqi Liu, Richard Powell, Vijay	618
555	narayana Pillai, Jacob Devlin, Michael Laskin, Diego	Bolina, Mariko Iinuma, Polina Zablotskaia, James	619
556	de Las Casas, Dasha Valter, Connie Tao, Lorenzo	Besley, Da-Woon Chung, Timothy Dozat, Ramona	620
557	Blanco, Adrià Puigdomènech Badia, David Reitter,	Comanescu, Xiance Si, Jeremy Greer, Guolong Su,	621
558	Mianna Chen, Jenny Brennan, Clara Rivera, Sergey	Martin Polacek, Raphaël Lopez Kaufman, Simon	622
559	Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,	Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie	623
560	Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yim-	Miao, Mohamed Elhawaty, Aditya Siddhant, Nenad	624
561	ing Gu, Kate Olszewska, Yujing Zhang, Ravi Ad-	Tomasev, Jinwei Xing, Christina Greer, Helen Miller,	625
562	danki, Antoine Miech, Annie Louis, Laurent El	Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma,	626
563	Shafey, Denis Teplyashin, Geoff Brown, Elliot Catt,	Angelos Filos, Milos Besta, Rory Blevins, Ted Kli-	627
564	Nithya Attaluri, Jan Balaguer, Jackie Xiang, Pi-	menko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi	628
565	dong Wang, Zoe Ashwood, Anton Briukhov, Al-	Mu, Oscar Chang, Mantas Pajarskas, Carrie Muir,	629
566	bert Webson, Sanjay Ganapathy, Smit Sanghavi,	Vered Cohen, Charline Le Lan, Krishna Haridasan,	630
567	Ajay Kannan, Ming-Wei Chang, Axel Stjerngren,	Amit Marathe, Steven Hansen, Sholto Douglas, Ra-	631
568	Josip Djolonga, Yuting Sun, Ankur Bapna, Matthew	j Kumar Samuel, Mingqiu Wang, Sophia Austin,	632
569	Aitchison, Pedram Pejman, Henryk Michalewski,	Chang Lan, Jiepu Jiang, Justin Chiu, Jaime Alonso	633
570	Tianhe Yu, Cindy Wang, Juliette Love, Junwhan Ahn,	Lorenzo, Lars Lowe Sjösund, Sébastien Cevey,	634
571	Dawn Bloxwich, Kehang Han, Peter Humphreys,	Zach Gleicher, Thi Avrahami, Anudhyan Boral,	635
572	Thibault Sellam, James Bradbury, Varun Godbole,	Hansa Srinivasan, Vittorio Selo, Rhys May, Kon-	636
573	Sina Samangooei, Bogdan Damoc, Alex Kaskasoli,	stantinos Aisopos, Léonard Hussenot, Livio Baldini	637
574	Sébastien M. R. Arnold, Vijay Vasudevan, Shubham	Soares, Kate Baumli, Michael B. Chang, Adrià Rec-	638
575	Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tan-	casens, Ben Caine, Alexander Pritzel, Filip Pavetic,	639
576	burn, Srivatsan Srinivasan, Hyeontaek Lim, Sarah	Fabio Pardo, Anita Gergely, Justin Frye, Vinay	640
577	Hodkinson, Pranav Shyam, Johan Ferret, Steven	Ramasesh, Dan Horgan, Kartikeya Badola, Nora	641
578	Hand, Ankush Garg, Tom Le Paine, Jian Li, Yu-	Kassner, Subhrajit Roy, Ethan Dyer, Víctor Cam-	642
579	jia Li, Minh Giang, Alexander Neitz, Zaheer Abbas,	pos, Alex Tomala, Yunhao Tang, Dalia El Badawy,	643
580	Sarah York, Machel Reid, Elizabeth Cole, Aakanksha	Elspeth White, Basil Mustafa, Oran Lang, Ab-	644
581	Chowdhery, Dipanjan Das, Dominika Rogozińska,	hishek Jindal, Sharad Vikram, Zhitao Gong, Sergi	645
582	Vitaly Nikolaev, Pablo Sprechmann, Zachary Nado,	Caelles, Ross Hemsley, Gregory Thornton, Fangxi-	646
583	Lukas Zilka, Flavien Prost, Luheng He, Marianne	aoyu Feng, Wojciech Stokowiec, Ce Zheng, Phoebe	647
584	Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan,	Thacker, Çağlar Ünlü, Zhishuai Zhang, Moham-	648
585	Dawei Jia, Miltiadis Allamanis, Clara Huiyi Hu,	mad Saleh, James Svensson, Max Bileschi, Piyush	649
586	Raoul de Liedekerke, Justin Gilmer, Carl Saroufim,	Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias,	650
587	Shruti Rijhwani, Shaobo Hou, Disha Shrivastava,	Arpi Vezer, Marco Selvi, Toby Shevlane, Mikel Ro-	651

652	driguez, Tom Kwiatkowski, Samira Daruki, Keran Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks, Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi Hashemi, Richard Ives, Yana Hasson, YaGuang Li, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal, Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević, Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot, Matthew Lamm, Nicola De Cao, Charlie Chen, Gamaleldin Elsayed, Ed Chi, Mahdis Mahdieh, Ian Tenney, Nan Hua, Ivan Petrychenko, Patrick Kane, Dylan Scandinaro, Rishub Jain, Jonathan Uesato, Romina Datta, Adam Sadovsky, Oskar Bunyan, Dominik Rabiej, Shimu Wu, John Zhang, Gautam Vasudevan, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu, Nan Wei, Ivy Zheng, Betty Chan, Pam G Rabinovitch, Piotr Stanczyk, Ye Zhang, David Steiner, Subhajt Naskar, Michael Azzam, Matthew Johnson, Adam Paszke, Chung-Cheng Chiu, Jaime Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin Miao, Andrew Lee, Nino Vieillard, Sahitya Potluri, Jane Park, Elnaz Davoodi, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luwei Zhou, Jonathan Evens, William Isaac, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Chris Gorgolewski, Peter Grabowski, Yu Mao, Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Paul Suganthan, Evan Palmer, Geoffrey Irving, Edward Loper, Manaal Faruqui, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Michael Fink, Alfonso Castaño, Irene Gianoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian LIN, Marin Georgiev, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Alena Repina, Xihui Wu, Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse, Fan Yang, Jeff Piper, Nathan Ie, Minnie Lui, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Lam Nguyen Thiet, Daniel Andor, Pedro Valenzuela, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Sarmishta Velury, Sebastian Krause, Cassidy Hardin, Lucas Dixon, Lili Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Tejasi Latkar, Mingyang Zhang, Quoc Le, Elena Allica Abellan, Dayou Du, Dan McKinnon, Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Sid Lall, Ken Franko, Egor Filonov, Anna Bulanova, Rémi Leblond, Vikas Yadav, Shirley	716
653	Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu, Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Hao Zhou, Alek Dimitriev, Hannah Forbes, Dylan Banarse, Zora Tung, Jeremiah Liu, Mark Omernick, Colton Bishop, Chintu Kumar, Rachel Sterneck, Ryan Foley, Rohan Jain, Swaroop Mishra, Jiawei Xia, Taylor Bos, Geoffrey Cideron, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Petru Gurita, Hila Noga, Premal Shah, Daniel J. Mankowitz, Alex Polozov, Nate Kushman, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu, Meghana Thotakuri, Tom Natan, Anhad Mohananey, Matthieu Geist, Sidharth Mudgal, Sertan Girgin, Hui Li, Jiayu Ye, Ofir Roval, Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Quan Yuan, Sumit Bagri, Danila Sinopalnikov, Sabela Ramos, John Mellor, Abhishek Sharma, Aliaksei Severyn, Jonathan Lai, Kathy Wu, Heng-Tze Cheng, David Miller, Nicolas Sonnerat, Denis Vnukov, Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy, Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang, Rui Zhu, Mark Geller, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Andrei Sozanschi, Daniel Toyama, Evan Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Abhimanyu Goyal, Diane Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana, Jing Li, Sabaer Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Tao Zhu, Yeongil Ko, Laura Knight, Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Dustin Tran, Yeqing Li, Nir Levine, Ariel Stolovich, Norbert Kalb, Rebecca Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Balaji Lakshminarayanan, Charlie Deck, Shyam Upadhyay, Hyo Lee, Mike Dusenberry, Zonglin Li, Xuezhi Wang, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem, Summer Yue, Sho Arora, Eric Malmi, Daniil Mirylenka, Qijun Tan, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Steven Zheng, Francesco Pongetti, Mukarram Tariq, Yanhua Sun, Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Ragha Kotikalapudi, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan, Aaron Parisi, Joe Stanton, Chenkai Kuang, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu, Abe Ittycheriah, Prakash Shroff, Pei Sun, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David Gaddy, Ishita Dasgupta, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht, Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivière, Alanna Walton, Clément Crepy, Alicia Parrish, Yuan Liu, Zongwei Zhou, Clement Farabet, Carey Radebaugh, Praveen Srinivasan, Claudia van der Salm, Andreas Fidje-land, Salvatore Scellato, Eri Latorre-Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Ce-	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		730
667		731
668		732
669		733
670		734
671		735
672		736
673		737
674		738
675		739
676		740
677		741
678		742
679		743
680		744
681		745
682		746
683		747
684		748
685		749
686		750
687		751
688		752
689		753
690		754
691		755
692		756
693		757
694		758
695		759
696		760
697		761
698		762
699		763
700		764
701		765
702		766
703		767
704		768
705		769
706		770
707		771
708		772
709		773
710		774
711		775
712		776
713		777
714		778
715		779



780	sare, Tom Hudson, Piermaria Mendolicchio, Lexi	and Enhong Chen. 2023. <a href="#">Woodpecker: Hallucination</a>	840
781	Walker, Alex Morris, Ivo Penchev, Matthew Mauger,	<a href="#">correction for multimodal large language models.</a>	841
782	Alexey Guseynov, Alison Reid, Seth Odoom, Lucia		
783	Loher, Victor Cotruta, Madhavi Yenugula, Dominik	Hugh Zhang, Daniel Duckworth, Daphne Ippolito, and	842
784	Grewe, Anastasia Petrushkina, Tom Duerig, Antonio	Arvind Neelakantan. 2021. <a href="#">Trading off diversity and</a>	843
785	Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson,	<a href="#">quality in natural language generation.</a> In <i>Proceed-</i>	844
786	Adam Kurzrok, Lynette Webb, Sahil Dua, Dong	<i>ings of the Workshop on Human Evaluation of NLP</i>	845
787	Li, Preethi Lahoti, Surya Bhupatiraju, Dan Hurt, Ha-	<i>Systems (HumEval)</i> , pages 25–33, Online. Associa-	846
788	roon Qureshi, Ananth Agarwal, Tomer Shani, Matan	tion for Computational Linguistics.	847
789	Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei		
790	Wang, Chetan Tekur, Mihir Sanjay Kale, Jinliang	Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun	848
791	Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty,	Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and	849
792	Yi Sun, Yao Zhao, Stephan Lee, Pandu Nayak, Doug	Huaxiu Yao. 2023. <a href="#">Analyzing and mitigating object</a>	850
793	Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi	<a href="#">hallucination in large vision-language models.</a>	851
794	Vyas, Martin Wicke, Xiao Ma, Taylan Bilal, Ev-		
795	genii Eltyshev, Daniel Balle, Nina Martin, Hardie	<b>A Appendix</b>	852
796	Cate, James Manyika, Keyvan Amiri, Yelin Kim,	<b>A.1 Detailed Experimental Setup</b>	853
797	Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripu-	For POPE and the descriptions experiment, we use	854
798	raneni, David Madras, Mandy Guo, Austin Waters,	the following LCD parameters $\beta = 3.0$ , $\alpha = 0.1$ .	855
799	Oliver Wang, Joshua Ainslie, Jason Baldridge, Han	We set the temperature to 0.5 in POPE and 1.0 in	856
800	Zhang, Garima Pruthi, Jakob Bauer, Feng Yang, Ri-	the descriptions experiment. We limit the descrip-	857
801	ham Mansour, Jason Gelman, Yang Xu, George	tions length to 250 tokens in all models we tested.	858
802	Polovets, Ji Liu, Honglong Cai, Warren Chen, Xi-	We don’t tune any of these parameters. The prompt	859
803	angHai Sheng, Emily Xue, Sherjil Ozair, Adams Yu,	we use for the descriptions experiment is <i>"Describe</i>	860
804	Christof Angermueller, Xiaowei Li, Weiren Wang, Ju-	<i>this image in detail:"</i> . The models we use have the	861
805	lia Wiesinger, Emmanouil Koukoumidis, Yuan Tian,	following Huggingface identifiers:	862
806	Anand Iyer, Madhu Gurusurthy, Mark Goldenson,		
807	Parashar Shah, MK Blake, Hongkun Yu, Anthony	• Salesforce/instructblip-vicuna-7b	863
808	Urbanowicz, Jennimaria Palomaki, Chrisantha Fer-		
809	nando, Kevin Brooks, Ken Durden, Harsh Mehta,	• Salesforce/instructblip-flan-t5-xl	864
810	Nikola Momchev, Elahe Rahimtoroghi, Maria Geor-		
811	gaki, Amit Raul, Sebastian Ruder, Morgan Red-	• llava-hf/llava-1.5-7b-hf	865
812	shaw, Jinhyuk Lee, Komal Jalan, Dinghua Li, Ginger	<b>GPT4-V Assisted Evaluation</b> We follow the	866
813	Perng, Blake Hechtman, Parker Schuh, Milad Nasr,	evaluation protocol given in Yin et al. (2023),	867
814	Mia Chen, Kieran Milan, Vladimir Mikulik, Trevor	where an image and two descriptions are given	868
815	Strohman, Juliana Franco, Tim Green, Demis Has-	to the model, formatted with the prompt in figure 2.	869
816	sabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol	The model outputs scores in two dimensions: Ac-	870
817	Vinyals. 2023. <a href="#">Gemini: A family of highly capable</a>	curacy and Detailedness. We use the <i>gpt-4-vision-</i>	871
818	<a href="#">multimodal models.</a>	<i>preview</i> model on February 2024.	872
819	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier		
820	Martinet, Marie-Anne Lachaux, Timothée Lacroix,		
821	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal		
822	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard		
823	Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open</a>		
824	<a href="#">and efficient foundation language models.</a>		
825	Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie		
826	Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and		
827	Lijuan Wang. 2022. <a href="#">Git: A generative image-to-text</a>		
828	<a href="#">transformer for vision and language.</a>		
829	Lei Wang, Jiabang He, Shenshen Li, Ning Liu, and		
830	Ee-Peng Lim. 2023. <a href="#">Mitigating fine-grained halluci-</a>		
831	<a href="#">nation by fine-tuning large vision-language models</a>		
832	<a href="#">with caption rewrites.</a>		
833	Kun Zhou Jinpeng Wang Wayne Xin Zhao Yifan Li,		
834	Yifan Du and Ji-Rong Wen. 2023. <a href="#">Evaluating object</a>		
835	<a href="#">hallucination in large vision-language models.</a> In <i>The</i>		
836	<i>2023 Conference on Empirical Methods in Natural</i>		
837	<i>Language Processing.</i>		
838	Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao		
839	Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun,		

```

prompt = Lambda A, B: f"""
You are required to score the performance of two AI assistants in describing a given image. You should pay extra
attention to the hallucination, which refers to the part of descriptions that are inconsistent with the image content,
such as claiming the existence of something not present in the image or describing incorrectly in terms of the counts,
positions, or colors of objects in the image. Note that the descriptions may be accompanied by bounding boxes,
indicating the position of objects in the image, which are represented as [x1, y1, x2, y2] with floating numbers ranging
from 0 to 1. These values correspond to the top left x1, top left y1, bottom right x2, and bottom right y2.
Please rate the responses of the assistants on a scale of 1 to 10, where a higher score indicates better performance,
according to the following criteria:
1: Accuracy: whether the response is accurate with respect to the image content. Responses with fewer hallucinations
should be given higher scores.
2: Detailedness: whether the response is rich in necessary details. Note that hallucinated descriptions should not count
as necessary details.
Please output a single line for each criterion, containing only two values indicating the scores for Assistant 1 and 2,
respectively. The two scores are separated by a space. Following the scores, please provide an explanation of your
evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not
affect your judgment.

[Assistant 1]
{A}
[End of Assistant 1]

[Assistant 2]
{B}
[End of Assistant 2]

Output format:

Accuracy:
Scores of the two answers:
Reason:

Detailedness:
Scores of the two answers:
Reason:
"""

```

Figure 2: Prompt used to evaluate descriptions with GPT4-V, taken from Yin et al. (2023)

## A.2 Detailed Experimental Results

Model	Method	METEOR $\uparrow$	WMD $\uparrow$	ROUGE_L $\uparrow$	CHAIRs $\downarrow$	CHAIRi $\downarrow$
InstructBLIP <sub>F</sub>	Baseline	0.151	0.361	0.156	0.666	0.174
	Baseline <sub>N</sub>	0.157	<u>0.367</u>	0.161	0.662	0.146
	LCD <sub>-dw</sub>	<u>0.159</u>	0.364	<u>0.163</u>	<u>0.594</u>	<u>0.133</u>
	LCD	<b>0.163</b>	<b>0.370</b>	<b>0.168</b>	<b>0.566</b>	<b>0.131</b>
InstructBLIP <sub>V</sub>	Baseline	0.171	0.408	0.274	0.308	0.138
	Baseline <sub>N</sub>	0.178	0.423	0.291	0.274	0.126
	LCD <sub>-dw</sub>	<b>0.202</b>	<u>0.474</u>	<u>0.366</u>	<u>0.23</u>	<u>0.116</u>
	LCD	<u>0.199</u>	<b>0.48</b>	<b>0.38</b>	<b>0.174</b>	<b>0.107</b>
LLAVA 1.5	Baseline	0.160	<u>0.353</u>	0.167	0.632	0.183
	Baseline <sub>N</sub>	0.163	<b>0.357</b>	0.169	0.672	0.182
	LCD <sub>-dw</sub>	<u>0.169</u>	0.352	<u>0.179</u>	<u>0.624</u>	<b>0.157</b>
	LCD	<b>0.171</b>	0.352	<b>0.181</b>	<b>0.610</b>	<u>0.161</u>

Table 3: Image Description ablations. *-dw* is an LCD variant without dynamic weighting, with  $\beta = 0.5$ . Baseline<sub>N</sub> is using nucleus sampling with  $p = 0.95$ , Baseline is vanilla sampling.

POPE	method	model	accuracy	precision	recall	f1	yes_ratio
random	Baseline	InstructBLIP Vicuna	84.90%	89.57%	79.00%	83.95%	44.10%
random	LCD	InstructBLIP Vicuna	87.53%	87.43%	87.67%	87.55%	50.13%
popular	Baseline	InstructBLIP Vicuna	83.30%	85.35%	80.40%	82.80%	47.10%
popular	LCD	InstructBLIP Vicuna	83.73%	81.31%	87.60%	84.34%	53.87%
adversarial	Baseline	InstructBLIP Vicuna	80.23%	80.17%	80.33%	80.25%	50.10%
adversarial	LCD	InstructBLIP Vicuna	80.27%	76.33%	87.73%	81.64%	57.47%
random	Baseline	InstructBLIP FlanT5	85.63%	94.43%	75.73%	84.05%	40.10%
random	LCD	InstructBLIP FlanT5	86.03%	96.47%	74.80%	84.27%	38.77%
popular	Baseline	InstructBLIP FlanT5	82.07%	87.17%	75.20%	80.74%	43.13%
popular	LCD	InstructBLIP FlanT5	84.43%	92.44%	75.00%	82.81%	40.57%
adversarial	Baseline	InstructBLIP FlanT5	79.83%	82.83%	75.27%	78.87%	45.43%
adversarial	LCD	InstructBLIP FlanT5	82.03%	87.22%	75.07%	80.69%	43.03%
random	Baseline	LLAVA 1.5	85.87%	95.67%	75.13%	84.17%	39.27%
random	LCD	LLAVA 1.5	85.73%	97.18%	73.60%	83.76%	37.87%
popular	Baseline	LLAVA 1.5	84.80%	93.57%	74.73%	83.10%	39.93%
popular	LCD	LLAVA 1.5	85.40%	96.17%	73.73%	83.47%	38.33%
adversarial	Baseline	LLAVA 1.5	82.77%	88.67%	75.13%	81.34%	42.37%
adversarial	LCD	LLAVA 1.5	83.33%	90.98%	74.00%	81.62%	40.67%

Table 4: Complete POPE results.