# ON THE CYCLE CONSISTENCY OF IMAGE-TEXT MAPPINGS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

Paper under double-blind review

#### ABSTRACT

We present an empirical study of cycle consistency in image-text mappings. We observe growing cycle consistency across a wide range of image-to-text and text-to-image models, i.e., images and text are becoming increasingly interchange-able in their representations. First, we investigate the factors driving this trend and identify that scaling language models and employing high-quality dataset re-captioning enhance cycle consistency. Next, we analyze the types of images and texts that are exchangeable, and find that cycle consistency strongly correlates with various desired properties such as reduced text hallucination, better descriptions, and improved compositionality and prompt-following in images. Lastly, we examine various sources of variance in cycle consistency demonstrating that text-to-image models are sensitive to specific prompt styles.

#### 022 1 INTRODUCTION

How would you convey the visual look of your hometown to a friend? One approach would 025 be to share a set of photos, showing differ-026 ent architectural elements and city scenes. An-027 other would be to give a verbal description: 028 "The roofs are made of half-cylinder terra cotta 029 tiles, layered one on top of the other." Both approaches convey visual information, even 031 though the latter is in the format of text. To wit, visual information can be communicated 033 either by images – the preferred format of the 034 computer vision scientist - or by language - the currency of NLP. The same is true in the other direction: information expressed via language 036 can be visualized in an image or infographic 037 that conveys some of the same meaning. Increasingly, multimodal models are blurring the lines between these two representational for-040



Figure 1: **Bidirectional image**  $\leftrightarrow$  **text mappings.** Image-text mappings are able to exchange text descriptions into images and vice versa. We analyze how close cycle reconstruction is to the original data.

mats (e.g., Liu et al., 2024b; Gemini, 2023), and there is interest in both the computer vision and
 NLP communities in forging links between the two modalities, so that we can apply tools from NLP
 to problems in vision (e.g., Surís et al., 2023), and vice versa (e.g., Hu et al., 2024).

This leads us to ask: are pixels and words fundamentally exchangeable formats, or are there limits 044 to how effectively text can represent images, and in how well images can convey the meaning in text? We address this question by studying the degree to which images can be translated into text 046 without losing information, and, vice versa, how faithfully can text be represented via an image. 047 We quantify this by looking in particular at the *cycle consistency* of image-text mappings. A cycle-048 consistent image mapping is one in which translating from an image to a text description, and back, results in the original image. Symmetrically, a cycle-consistent text mapping translates a starting text into an image, and back into text which matches the original input. While cycle consistency is 051 a desirable emergent property for image-text mappings, it also gives us insights about 1) what kinds of images and text lead to successful exchanges of information and 2) the models which generated 052 these images and text. Furthermore, recent models have begun to incorporate cycle consistency at training time Betker et al. (2023); Esser et al. (2024); Sharifzadeh et al. (2024); Li et al. (2024b).



Figure 3: Examples of text cycle consistency from different model combinations. Models still exhibit a degree of cycle consistency despite being trained independently. Each example shows the generated image on the left for the input text on the right, followed by the reconstructed text below. SBERT similarity (<sup>↑</sup>) between the input and output text is reported in brackets. Highlighted phrases in the descriptions are for better comparison.

We make the following findings. First, we experiment with combining off-the-shelf text-to-image and image-to-text models to create both image and text cycles. We find that current models are fairly cycle-consistent semantically, but still distant from pixel to pixel (as seen in Figures 2 and 3). Furthermore, we observe an increasing correlation between cycle consistency and model performance. We analyze several key advancements contributing to this trend: language model scale,

higher resolutions, and training data with densely captioned images. Secondly, we observe that
cycle-consistent captions are descriptive, exhibit reduced object hallucination and omission, and
dense in length. Cycle-consistent images demonstrate better prompt-following and improved compositionality across different categories. Because image-text mappings are not one-to-one, we compare sources of variance in these mappings, including forms of sampling and prompt style choices, and their effect on cycle consistency.

114 115

116

128

129

130

131

# 2 PRELIMINARIES

117 118 We examine to what degree current models display cycle-consistent properties. Similar to how au-119 to encoders calculate error between original inputs and decoded outputs to evaluate performance, we 120 measure both how well are images reconstructed through text and how well are text descriptions 120 preserved by images. We use I to denote a set of real images, and T to denote a set of text descrip-121 tions. Given an image-to-text model F and a text-to-image model G, they exhibit cycle consistency 122 if  $G(F(i)) \approx i$  for all  $i \in I$  and  $F(G(t)) \approx t$  for all  $t \in T$ . We measure *image cycle consistency* 123 and *text cycle consistency* by computing the following reconstruction losses respectively:

$$\mathcal{L}_{\text{img}} = \mathbb{E}_{i \in I}[d_{\text{img}}(G(F(i)), i)], \tag{1}$$

$$\mathcal{L}_{\text{text}} = \mathbb{E}_{t \in T}[d_{\text{text}}(F(G(t)), t)].$$
(2)

We measure the distance between image i and reconstructed image G(F(i)) with the DreamSim (Fu et al., 2023b) image distance metric  $d_{img}$ . We find that DreamSim cycle consistency best correlates with text descriptiveness and aligns with human perception of image distance. Similarly, we use SBERT (Reimers & Gurevych, 2019) to measure similarity between input and reconstructed text. See Appendix C.1 for ablations on measuring cycle consistency.

132 133 134

135

# 3 WHAT FACTORS ARE DRIVING CYCLE CONSISTENCY?

In this section, we analyze the driving factors for cycle consistency. We evaluate image and text cycle consistency for 13 image-to-text models and 5 text-to-image models (i.e., 130 cycle-consistent mappings). These models were trained with varying datasets, architecture, and scale, and were selected based on public availability and disclosure of details. See Appendix A for a complete list of models and summary of differences.

141 We use Densely Captioned Images (DCI) dataset (Urbanek et al., 2024), which features high-142 resolution images annotated with dense captions, compared to other datasets (e.g., 480×640 pixels, 143 13.54 tokens for MSCOCO (Lin et al., 2014)). Due to limited prompt length of text-to-image mod-144 els, we use sDCI which summarizes DCI captions to fit 77 tokens (1500×2250 pixels, 49.21 tokens). 145 We sample 1K examples from the train split. We report the average cycle consistency for each textto-image model, computed across 13 image-to-text models and 3 random seeds. We follow the same 146 procedure with image-to-text models. Image and text cycle consistency calculations for all possible 147 model combinations are shown in Figures 14, 15, Figure 16 provides a baseline comparison. 148

- 149
- 150 3.1 CYCLE CONSISTENCY IMPROVES WITH LLM SCALE

151 An image-to-text model consists of a vision encoder, a projector, and a large language model (LLM). 152 Scaling the vision transformer (ViT) for the vision encoder is reported to enhance performance (Li 153 et al., 2023b), yet a simple MLP projection remains the dominant approach (Liu et al., 2023b) 154 OpenGVLab, 2024; Li et al., 2024a). As no model offers open-sourced weights with varying vision 155 encoder scales while keeping other parameters fixed, we focus our analysis on ablating the LLM 156 size. Figure 4 demonstrates that scaling the LLM enhances both image and text cycle consistency 157 across all image-to-text model families. Figure 5 highlights the effect of LLM size on image cycle 158 consistency, comparing the InternVL2 model family trained on the same architecture and dataset but 159 with varying LLM scales. We observe that scaling the LLM improves caption descriptiveness, e.g., InternVL2-40B is the only model capable of accurately describing both the color and the presence 160 of a corner turret. In contrast, models with smaller LLMs fail to capture such fine-grained details, 161 leading to reduced image cycle consistency.



3.2 RE-CAPTIONED DATASET QUALITY

202 Current image-to-text models are predominantly trained on re-captioned datasets where real images 203 are annotated with detailed descriptions generated by large language models (e.g., GPT-4) or vision-204 language models (e.g., GPT-4V). Similarly, recent works demonstrate that training text-to-image 205 models with descriptive captions generated by high-quality captioning models significantly enhances 206 their prompt-following ability (Betker et al.) 2023; Esser et al., 2024). As most text-to-image models 207 do not disclose information on training data, our analysis primarily focuses on image-to-text models and their re-captioned datasets. Ideally, the analysis would involve the same model trained with 208 and without re-captioned datasets, or with datasets of varying quality; however, such open-sourced 209 weights are not available. Consequently, we compare different models of similar sizes. To limit the 210 number of uncontrollable variables, we compare models with similar number of parameters. Note 211 that other factors, such as architecture, pre-trained backbones, and model training still differ between 212 the models. 213

Table I demonstrates that the quality of the re-captioned dataset (e.g., dataset re-captioned by GPT 4V, LLaVA1.6-34B) aligns with improved image cycle consistency. Models trained on such datasets often exhibit better consistency than those trained on larger datasets annotated by less-performant

216		Re-	captioned Dataset	Cycle Consistency		
217	Model	Size	Re-captioning Model	Image $(\downarrow)$	Text $(\uparrow)$	
218	BLIP2-6.7B	244M	BLIP	0.5936	0.5215	
220	LLaVA1.5-7B	23K	GPT-4*	0.5022	0.6290	
221	LLaVA1.6-7B	112K	GPT-4V	0.4833	0.6173	
222	LLaVA-OV-7B	3.5M	LLaVA1.6-34B	0.4742	0.6259	

Table 1: **Re-captioned dataset quality and cycle consistency.** The quality of the re-captioned dataset (i.e., generated by a high-performing model) aligns with improved image cycle consistency. Models trained on such datasets often exhibit better consistency than those trained on larger datasets annotated by less-performant models. In contrast, text cycle consistency shows little difference between the LLaVA models due to limited descriptiveness of the input text (sDCI), resulting in diminishing improvements for longer, more detailed captions, such as those reconstructed by LLaVA1.6 and LLaVA-OV. Image cycle consistency is measured by DreamSim (lower is better), and text cycle consistency by SBERT (higher is better).



Figure 6: Cycle consistency correlates with better prompt-following images. We measure cycle consistency (averaged across all image-to-text models) as a function of prompt-following quality in generated images. We observe a strong correlations for both kinds of cycle consistency and prompt-following quality in images.

models (e.g., BLIP). Our findings align with existing literature (Li et al., 2024a) that highlights "quality over quantity" for training multimodal models. On the other hand, text cycle consistency shows little difference between the LLaVA models, as the input text from sDCI often lacks fine-grained detail (evidenced in Figure 16) compared to longer and more descriptive synthetic captions, such as those produced by LLaVA1.6 and LLaVA-OV. We believe higher-quality human annotations and text-to-image models with longer context would enhance the analysis of text cycle consistency. Note that for LLaVA1.5, GPT-4 (i.e., language model) is prompted with captions and bounding boxes to generate detailed captions. We exclude InternVL2 from this analysis as information regarding its pre-training dataset is not disclosed.

#### 4 PROPERTIES OF CYCLE-CONSISTENT IMAGES AND TEXTS

In this section, we explore the quality of cycle-consistent texts and images with respect to various
 properties. We find more descriptive and less-hallucinated captions and better prompt-following
 images generally align with higher cycle consistency.

265 4.1 COMPOSITIONALITY AND PROMPT-FOLLOWING IN IMAGES

We investigate how cycle consistency varies with compositionality and prompt-following in text-to-image generation. We measure these qualities on two benchmarks: T2I-Compbench (Huang et al., 2023), which focuses on image compositionality, and Drawbench which includes a variety of categories for general text-to-image synthesis (Saharia et al., 2022).



Figure 7: Cycle consistency strongly correlates with descriptiveness in text. While both cycles exhibit a strong correlation, we observe a better alignment between text cycle consistency and descriptiveness in the generated captions.



Figure 8: Text descriptiveness and cycle consistency. From left to right: Original photograph, and generated 300 image reconstructions made by different captions under each image with the DreamSim ( $\downarrow$ ) reconstruction in brackets. As captions include more specific and correct visual details, the reconstruction quality increases. All images are generated with SD3 from the same random seed.

304 For T2I-Compbench, we evaluate on color, shape, texture, and spatial fine-grained categories. Sim-305 ilarly to Section 3 we calculate cycle consistency for text-to-image models by averaging across all 13 image-to-text models. Figure 6 plots image and text cycle consistency against text-to-image per-306 formance for all 5 text-to-image models. Cycle consistency highly correlates with text-to-image 307 generation quality. Intuitively, images which are more faithful to visual details represented in text 308 preserve more information and therefore facilitate better image and text reconstructions. Ablations 309 for measuring cycle consistency are discussed in Appendix C.1. 310

311 312

287

289

290 291

293

295

296

297

298

299

301

302

303

#### 4.2 TEXT DESCRIPTIVENESS

313 Similarly to Section 4.1, we analyze the relationship between the descriptiveness of captions and cy-314 cle consistency. Typically, image captions are evaluated using the CIDEr score on the COCO Karpa-315 thy split (Karpathy & Fei-Fei, 2015). This dataset covers a limited distribution of images annotated 316 with short captions which poses a challenge in evaluating modern image-to-text models which gener-317 ate long, descriptive captions. Inspired by recent text-to-image benchmarks (Huang et al., 2023; Sa-318 haria et al. 2022), we instead conduct visual question answering without the image component, i.e., 319 "VQA without V". Given a VQA dataset  $\{v, q, a\}_{i=1}^{N}$ , an image-to-text model F, we first generate 320 synthetic text F(v) for images in the VQA dataset. Then we prompt a large language model (LLM) 321 to answer the question based on the generated caption. This allows us to measure whether synthetic text accurately describes fine-grained details of the image, as the LLM must answer a diverse range 322 of questions based solely on the description. We use Meta-Llama-3.1-8B-Instruct (Dubey) 323 et al., 2024) as the LLM evaluator in all experiments. Exact prompts are detailed in Appendix A.3.

325

326

327

328

330

331

332

333

334

335

336 337

338

339

340

341

355



Figure 10: Cycle consistency strongly correlates with reduced hallucination in text. We observe that both cycles strongly correlates with reduced hallucination in text, with image cycle consistency being more prominent.

To assess how informative synthetic text is across fine-grained attributes (e.g., counting, position), we perform "VQA without V" on three VQA benchmarks: SEEDBench (Li et al.) 2023a), MME (Fu et al.) 2023a), and MMStar (Chen et al.) 2024a). As we evaluate how well synthetic text describes images, we focus on questions from the "perception" categories across all datasets. Cycle consistency is calculated as in Section 3.

347 Figure 7 compares image-to-text model captioning performance 348 scores against image and text cycle consistency scores (averaged 349 across all text-to-image models). This demonstrates a strong cor-350 relation between both image and text cycle consistency and cap-351 tioning performance. Generally, more informative captions lead to better image reconstructions through text, and also better re-352 covery of input text details such as in Figure 8. For a discussion 353 of ablations measuring cycle consistency see Appendix C.1. 354

**Sometimes less is more.** Although generally descriptive cap-356 tions exhibit better cycle consistency, we observe examples of 357 high cycle consistency using short, undescriptive captions. Such 358 instances are not uncommon, and mainly occur when the input 359 image displays a very typical scene. In Figure 9 (Top Row), im-360 ages of large golden statues are often taken by people at ground 361 level looking up at the statue. Perhaps such images are "com-362 mon" enough that a few keywords are sufficient to fully describe the scene. See Appendix Figure 22 for more examples. 364



[0.274] Anemone, Clownfish

Figure 9: **Short text captions** can yield high image cycle consistency despite lack of detail. Strong bias between certain images and text captions leads to easy reconstructions.

365
 366
 4.3 OBJECT HALLUCINATION IN TEXT

To investigate object hallucination in the generated text, we utilize the POPE benchmark (Li et al., 2023c), using their annotated COCO dataset. POPE constructs a set of triples consisting of an image, multiple questions and their answers  $\langle x, \{q(o_i), a_i\}_{i=1}^l \rangle$ , where x is the image,  $q(o_i)$  is a question probing object  $o_i$  based on a template "Is there a/an <object> in the image?",  $o_i$  is the *i*-th object to be probed, and  $a_i$  is the answer to the question ("Yes" or "No"), and l denotes the number of questions per image. While POPE focuses on evaluating hallucination at the model level, our goal is to measure hallucination in the generated text.

To achieve this, we perform the VQA without V analysis on POPE. Specifically, for a given LLM M, we evaluate  $M(f(x), q(o_i)) = a_i$ , which measures how many questions can be answered based on the generated caption. This contrasts with evaluating  $f(x, q(o_i)) = a_i$ , which measures how accurately an image-to-text model f can answer questions directly. The experimental setup is identical to that described in Section [4.2]. Figure [10, [11] demonstrates that cycle consistency strongly



385

386

387

388

389 390

391

392 393

394

378



Figure 11: Reduced object hallucination (POPE  $\uparrow$ ) correlates with better cycle consistency (DreamSim  $\downarrow$ ).

correlates with texts with reduced hallucination. Notably, we observe that *image* cycle consistency shows a stronger correlation with reduced hallucination.

#### 4.4 TEXT DENSITY

Given that more informative captions correlate with more cycle consistent images as seen in Section 4.2. 396 we now study how densely information should be 397 packed into a caption - i.e., to get the best image re-398 construction, what is the ideal caption length? To 399 properly control both the level of detail and length 400 of captions, we use summarized captions for im-401 ages in the Densely Captioned Images dataset (Ur-402 banek et al., 2024) created by Huh et al. (2024). In 403 this dataset, each image is accompanied by captions 404 of different lengths: 5, 10, 20, 30, and 50 word 405 summaries of the original fully detailed DCI cap-406 tion. For each captions of each summary length, we 407 use text-to-image models to generate images over 10 different random seeds and then report the av-408 erage DreamSim score between the generated im-409 ages and the original image described by the cap-410 tions. Figure 12 plots the relationship between cap-411 tion density and image reconstruction. Similarly to 412 Section 4.2, we find that increasing the amount of 413 granularity of captions improves reconstruction er-414 ror for all text-to-image models, although with high 415 variance. Furthermore, for models SD1.5, SDXL,



Figure 12: **Image cycle consistency vs. caption density.** We measure cycle consistency across captions for the same image summarized into varying lengths. Reconstruction scores increase with more descriptive captions, with reduced benefit after about 30 words. Error bars show standard deviation across 10 different random seeds.

and SDXL-Turbo image reconstruction sees little benefit beyond 30 tokens, whereas SD3 and
 FLUX-Time continue to show improvement. Figure 13 provides examples of summary captions
 and their corresponding synthetic images.

419 420

421

# 5 VARIANCE IN CYCLE CONSISTENCY

422 In the previous sections, we mainly observe cycle consistency by greedy sampling from image-to-423 text models and averaging over three random seeds for text-to-image models. In this section, we 424 observe how stochasticity can affect cycle consistency. Text-to-image and image-to-text models can 425 exhibit stochastic behavior due to factors such as random seed initialization, temperature sampling, 426 and differences in prompt wording. While random seed and temperature sampling are only relevant 427 to one mapping direction, prompt style applies to both. Different prompt styles for text-to-image 428 generation come from changing the text while maintaining its meaning. One such example is editing word choice and syntax. For image-to-text models, input prompts can be used to query the model 429 and request different kinds of text descriptions. These sources of variability can lead to different 430 results given identical inputs. In this section, we analyze the extent to which each factor causes 431 variance in cycle consistency.

432 Input Image 5 Word Summary 10 Word Summary 20 Word Summary 30 Word Summary 50 Word Summary 433 434 435 436 [0.405] A man and woman walk towards a KFC [0.457] A man 437 [0.528] Man and [0.633] Man walks to KFC oman walking to a pizza under a umbrella, 438 KFC with bright blue restaurant. bright blue sky with white clouds and a red umbrella. sky and red u ıbrella 439 440 441 442 443 444 [0.625] White [0.552] Coconut ice cream with topped ice cream cups, baskets ... signs with price colorful vendor 445 cups with umbrellas coconut on a colorful stand. umbrellas. s, with signs listing and appl 446

Figure 13: Effect of caption density on image cycle consistency. From the left to right: Original photograph, and synthetic reconstructions made by different summarized versions of the same caption. Under each synthetic image is the DreamSim reconstruction followed by the summary caption. Generally as the captions become longer and more descriptive, they better match the original image.

Source of Variance	Image Cycle Consistency	Text Cycle Consistency
Text-to-image Models		
Random seed	0.0415	0.0634
Caption style	0.0709	0.0719
Image-to-text Models		
Temperature sampling	0.0642	0.0588
Prompt style	0.0371	0.0427

 Table 2: Sources of variance in cycle consistency. We compare how random seed, prompt style, and temperature sampling affect variance in cycle consistency. For each source of variance we report the average standard deviation using DreamSim and SBERT for image and text cycle consistency respectively.

463 For each factor, we generate N = 10 variations and calculate the average standard deviation using 464 DreamSim for image cycle consistency and SBERT for text cycle consistency. For temperature sam-465 pling, we set the temperature to 0.7. For text-to-image models, we modify prompt style by using 466 Meta-Llama-3.1-8B-Instruct (Dubey et al., 2024) to rewrite a given prompt while main-467 taining its original meaning and number of words. The choice of prompts for the models is included 468 in the Appendix A.1 We use the DCI dataset and sample 100 examples with N = 10 variations 469 each across 60 model combinations, resulting in variance calculations over 6,000 examples. Note that we excluded BLIP2 models from measuring prompt style variance as they are not instruction-470 tuned and often produce captions of less than 3 words, making it challenging to change the style 471 without changing its meaning. 472

Table 2 demonstrates that image-to-text models exhibit higher variance due to temperature sampling but remain relatively robust to changes in prompt style. In contrast, text-to-image models are significantly more sensitive to prompt style than random seed sampling.

476 477 478

479

447

448

449

#### 6 RELATED WORK

Large multimodal models are rapidly improving, particularly in vision and language. Image-to-text models are capable of producing comprehensive image descriptions (Gemini, 2023; OpenAI) by scaling the language model (Liu et al., 2023b;a) and training on semantically-rich synthetic captions (Li et al., 2022; 2023b; Sharifzadeh et al., 2024; Liu et al., 2023b;a).

Concurrently, text-to-image models can generate images that follow a wide range of prompts (Podell
 et al., 2023; Sauer et al., 2023; Esser et al., 2024; BlackForestLabs). Recent works (Betker et al., 2023; Brooks et al., 2024) further enhance the prompt-following ability by learning a descriptive

image captioner and generating pseudo image-text pairs. The models trained on this dataset are capable of generating images faithful to long, descriptive captions.

While several works implicitly encourage image-text cycle consistency during sampling, recent work (Li et al., 2024b) explicitly enforces cycle consistency by leveraging unpaired image and text data, ensuring alignment between the original samples and their cycle-generated counterparts. Despite large image-text models becoming increasingly cycle-consistent, this property has been surprisingly little studied. In this work, we provide an in-depth analysis of cycle-consistent properties across a wide range of off-the-shelf image-to-text and text-to-image models.

494 495 496

506 507

508

512

513

514

527

528

529

# 7 CONCLUSION

497 This paper studies cycle consistency in current image-text mappings. We find that existing image-498 to-text and text-to-image models achieve a certain level of cycle consistency, even without explicit 499 training for it. We observe that image-to-text models with larger LLMs trained on high-quality 500 re-captioned data are associated with higher cycle consistency. Moreover, we show that cycle con-501 sistency improves with the quality of synthetic text and image generations. Generated images that 502 follow the compositionality and details provided by the input text model tend to be more cycleconsistent, and similarly for more detailed, informative, and accurate captions. Lastly, we highlight 504 stochastic factors that may affect cycle consistency, and find that prompt style for text-to-image 505 models contributes the most variance to cycle consistency.

## References

- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang
   Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3):8, 2023.
  - BlackForestLabs. Announcing black forest labs. https://blackforestlabs.ai/ announcing-black-forest-labs/, Accessed: 2024-09-24.
- Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. URL https://openai.com/research/video-generation-models-as-world-simulators.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
  Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language
  models? *arXiv preprint arXiv:2403.20330*, 2024a.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qing-long Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *arXiv* preprint arXiv:2312.14238, 2023.
  - Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024b.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
  Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- 538 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
  539 Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023a.

540 Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and 541 Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic 542 data. arXiv preprint arXiv:2306.09344, 2023b. 543 Gemini. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805, 544 2023. 546 Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, 547 and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal 548 language models. arXiv preprint arXiv:2406.09403, 2024. 549 Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A compre-550 hensive benchmark for open-world compositional text-to-image generation. Advances in Neural 551 Information Processing Systems, 36:78723–78747, 2023. 552 553 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation 554 hypothesis. In International Conference on Machine Learning, 2024. 555 Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descrip-556 tions. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3128-3137, 2015. 558 559 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. arXiv preprint 561 arXiv:2408.03326, 2024a. 562 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Bench-563 marking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125, 564 2023a. 565 566 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-567 training for unified vision-language understanding and generation. In International conference on 568 machine learning, pp. 12888–12900. PMLR, 2022. 569 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image 570 pre-training with frozen image encoders and large language models. In International conference 571 on machine learning, pp. 19730–19742. PMLR, 2023b. 572 573 Tianhong Li, Sangnie Bhardwaj, Yonglong Tian, Han Zhang, Jarred Barber, Dina Katabi, Guil-574 laume Lajoie, Huiwen Chang, and Dilip Krishnan. Leveraging unpaired data for vision-language generative models via cycle consistency. In The Twelfth International Conference on Learning 575 Representations, 2024b. 576 577 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating 578 object hallucination in large vision-language models. arXiv preprint arXiv:2305.10355, 2023c. 579 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 580 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer 581 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, 582 Proceedings, Part V 13, pp. 740–755. Springer, 2014. 583 584 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 585 tuning, 2023a. 586 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In NeurIPS, 587 2023b. 588 589 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 590 tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-591 *tion*, pp. 26296–26306, 2024a. 592 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances 593

in neural information processing systems, 36, 2024b.

594 595 596	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> , 2023c.
598 599	OpenAI. Hello gpt-40. https://openai.com/index/hello-gpt-40/. Accessed: 2024-09-24.
600	OpenGVLab. Internvl-2.0. 2024. URL https://internvl.github.io/blog/
601	2024-07-02-InternVL-2.0/
603 604 605	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. <i>arXiv preprint arXiv:2307.01952</i> , 2023.
606 607 608 609	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021.
610 611 612 613 614	Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert- networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language</i> <i>Processing</i> . Association for Computational Linguistics, 11 2019. URL <a href="https://arxiv.&lt;br&gt;org/abs/1908.10084">https://arxiv.</a>
615 616 617	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF confer-</i> <i>ence on computer vision and pattern recognition</i> , pp. 10684–10695, 2022.
618 619 620 621	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in neural information processing systems</i> , 35:36479–36494, 2022.
623 624	Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion dis- tillation. <i>arXiv preprint arXiv:2311.17042</i> , 2023.
625 626 627	Sahand Sharifzadeh, Christos Kaplanis, Shreya Pathak, Dharshan Kumaran, Anastasija Ilic, Jovana Mitrovic, Charles Blundell, and Andrea Banino. Synth2: Boosting visual-language models with synthetic captions and image embeddings. <i>arXiv preprint arXiv:2403.07750</i> , 2024.
628 629 630 631	Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 11888–11898, 2023.
632 633 634 635	Jack Urbanek, Florian Bordes, Pietro Astolfi, Mary Williamson, Vasu Sharma, and Adriana Romero- Soriano. A picture is worth more than 77 text tokens: Evaluating clip-style models on dense cap- tions. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 26700–26709, 2024.
636 637 638 639	Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
640 641	Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text generation. Advances in Neural Information Processing Systems, 34:27263–27277, 2021.
642 643 644	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 586–595, 2018.
645 646 647	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluat- ing text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> , 2019.

# 648 A MODEL DETAILS

650 We use 13 different image-to-text models and 5 text-to-image models to study cycle consistency. 651 Models are chosen based on public availability and disclosure of details such as architecture, scale, 652 training dataset and method etc. We use the following models for image-to-text mappings: BLIP-653 2.7B, BLIP-6.7B, BLIP-2-Flan T5-XXL Li et al. (2023b), LLaVA 1.5-7B, LLaVA 1.5-13B Liu et al. 654 (2023a), LLaVA OneVision-Qwen2-0.5B, LLaVA OneVision-Qwen2-7B Li et al. (2024a), LLaVA 1.6 Mistral-7B, LLaVA 1.6-34B Liu et al. (2024a), InternVL2-2B, InternVL2-8B, InternVL2-26B, 655 and InternVL2-40B Chen et al. (2023; 2024b). Table 5 provides a summary of model differences in 656 terms of scale and architecture. 657

We use the following models for text-image mappings: Stable Diffusion 1.5 (Rombach et al.) 2022),
Stable Diffusion XL (Podell et al.) 2023), Stable Diffusion XL Turbo (Sauer et al.) 2023), Stable Diffusion 3 (Esser et al.) 2024), and FLUX (Timestep-distilled) (BlackForestLabs). Tables 6 provides a summary of model differences.

662 663

677 678 679

680

692 693 694

#### A.1 HYPERPARAMETERS FOR IMAGE-TO-TEXT MODELS

To ensure that all image-to-text models can produce image descriptions to the best of their ability, we use the prompt recommended by the model distributor, as shown in Table 3. We use greedy search for all experiments (except for temperature sampling in Section 5), and 77 maximum tokens, i.e., maximum prompt length supported by text-to-image models.

Model	Prompt
BLIP2	"this is a picture of"
LLaVA1.5	"Write a detailed description of the given image."
LLaVA1.6	"Write a detailed description of the given image."
LLaVA-OV	"Write a detailed description of the given image."
InternVL2	"Please describe the image in detail."

Table 3: Prompts used for generating image descriptions for image-to-text models.

### A.2 HYPERPARAMETERS FOR TEXT-TO-IMAGE MODELS

To ensure that all text-to-image models can produce outputs to the best of their ability, for each model we use the settings recommended by the model distributor. Hyperparameters for each model are reported in Table 4. We use random seeds 0, 123, and 324229 in our experiments.

Model	Resolution	Steps	Guidance Scale
SD1.5	512	50	7.5
SDXL-Turbo	512	4	0
SDXL	1024	50	7.5
SD3	1024	50	7.5
FLUX-Time	1024	4	0

Table 4: Hyperparameters for text-to-image models.

### A.3 VQA WITHOUT V

We use the following prompt for the LLM judge in the VQA without V experiment in Section 4.2

For a finite formula formula

To generate the image description, we use the prompt "*Write a caption for this image*." for all image-to-text models.

Model	# Params	Vision Encoder	Projector	LLM
BLIP2-2.7B	3.8B	EVA-CLIP ViT-g (1.1B)	QFormer	OPT (2.7B)
BLIP2-6.7B	7.8B	EVA-CLIP ViT-g (1.1B)	QFormer	OPT (6.7B)
LLaVA1.5-7B	7.1B	CLIP ViT-L (304M)	MLP	Vicuna-1.5 (7B)
LLaVA1.5-13B	13.4B	CLIP ViT-L (304M)	MLP	Vicuna-1.5 (13B)
LLaVA1.6-7B	7.6B	CLIP ViT-L (304M)	MLP	Mistral (7B)
LLaVA1.6-34B	34.8B	CLIP ViT-L (304M)	MLP	Nous-Hermes-2-Yi (34B)
InternVL2-2B	2.5B	InternViT (304M)	MLP	InternML (2.2B)
InternVL2-8B	8.1B	InternViT (304M)	MLP	InternML (7.7B)
InternVL2-26B	25.5B	InternViT (5.5B)	MLP	InternML (19.9B)
InternVL2-40B	40B	InternViT (5.5B)	MLP	InternML (34.4B)
LLaVA-OV-0.5B	0.9B	SigLIP ViT-L/14 (307M)	MLP	Qwen-2 (0.5B)
LLaVA-OV-7B	8B	SigLIP ViT-L/14 (307M)	MLP	Qwen-2 (7B)

Table 5: Summary of image-to-text models on model architecture and scale.

Model	# Params	Image Generator	Context Dim.	Text Encoder	Dataset Re-captioning
SD1.5	983M	UNet (860M)	768	CLIP ViT-L	×
SDXL	3.5B	UNet (2.6B)	2048	CLIP ViT-L & OpenCLIP ViT-G (817M)	×
SDXL-Turbo	3.5B	UNet (2.6B)	2048	CLIP ViT-L & OpenCLIP ViT-G (817M)	×
FLUX-Time	12B	MMDiT	2816	CLIP VIT-L & T5 XXL	-
SD3	2B	MMDiT	4096	CLIP ViT-L & OpenCLIP ViT-G & T5 XXL	50% real 50% CogVLM cap

Table 6: Summary of text-to-image models on model architecture and scale.

### B CYCLE CONSISTENCY FOR ALL MODEL COMBINATIONS

We report cycle reconstruction scores across all different model combinations (13 image-to-text models  $\times$  5 text to image models) for both text and image cycles. Figures 14 and 15 display heatmaps of scores for image and text cycle consistency respectively.



Figure 14: Image cycle consistency on DCI dataset. We report the average score across 3 random seeds.

### C ABLATIONS ON MEASURING CYCLE CONSISTENCY

#### C.1 METRICS

For image cycle consistency, we measure DreamSim (Fu et al., 2023b), LPIPS (Zhang et al., 2018),
CLIP (Radford et al., 2021), and MSE between input and reconstructed images. Correspondingly,
we measure the distance between input text and reconstructed text with various text similarity metrics: BertScore (Zhang et al., 2019), BartScore (Yuan et al., 2021), SBERT (Reimers & Gurevych,
2019), and CLIP (Radford et al., 2021). For CLIP, we measure the cosine similarity between features
from the CLIP text encoder. Image-text alignment is measuring using CLIP and ImageReward Xu
et al. (2024).



Figure 16: Human text and real image as baselines. We compare human text as a baseline to synthetic captions generated by image-to-text models, and real image as a baseline to synthetic images generated by text-to-image models. We detail why real image performs worse in Figure 17.

Similarly to Figure 7, we investigate correlation between cycle consistency and captioning/text-to-image generation quality by ablating different similarity metrics. Tables 7, 8 report the Pearson correlation coefficient between cycle consistency and captioning or image generation performance respectively for different reconstruction metrics. Perceptual similarity metrics, i.e., DreamSim and LPIPS, align best with captioning performance, followed by CLIP and MSE. We additionally com-pare image-text similarity metrics and find that image reconstruction scores are more predictive of image captioning performance. For MSCOCO images and captions, we also report compare cycle consistency computed with various metrics and POPE in Table 9. Again, DreamSim has the strongest correlation. 

#### C.2 MODEL ABLATIONS FOR MEASURING CYCLE CONSISTENCY

In this paper, cycle consistency calculations for image-to-text models are averaged across all 5 text-to-image models. Another option would be to choose a text-to-image model to fix, and then use this one fixed model to calculate cycle consistency for all image-to-text models. This section investigates the how this choice of fixed model (or averaging across all fixed models) affects cycle consistency correlations in Figure 7. 

We report the Pearson correlation coefficient per model. As shown in Figure xxx, the correlation is consistently strong for most models ( $R^2 > 0.65$ ), except for BLIP2-2.7B and LLaVA-OV-0.5B with lower coefficients of 0.349 and 0.241, respectively. We attribute the low correlation to their use of small-scale, less-performant language models (OPT-2.7B, Qwen2-0.5B) as pre-trained backbones which may cause poorer text reconstruction. 

#### DOWNSTREAM PERFORMANCE D

Input Text

A group of diverse tents are set up on the ground, surrounded by people, cars, and buildings. Orange cones and logs are scattered on the sidewalk and ground. The buildings in the background are a mix of group beck

background are a mix of gray, brown, and trees, with a palm tree on the left.

The tents include a red, white, and neon green tent tops.

Input Text

812 813 814

810

811

815 816 817

818

819

820 821

822 823

824 825

826

827

837

838

A group of diverse tents are set up on the ground, surrounded by people, cars, and buildings. Orange cones and logs are scattered on the sidewalk and ground. The buildings in the background are a mix of gray, brown, and grown with a pale trave on the loft and trees, with a palm tree on the left. The tents include a red, white, and neon green tent tops.

colorful tents and consist, will inductors colorful tents and consists at up in the foreground, indicating a festival or market. The tents are of various colors, including red, green blue, and white, and are arranged in a somewhat organized manner, suggesting different stalls or booths for Reconstructed Text

Reconstructed Text

[SBERT 0.683] The image depicts a bustling outdoor event taking place in a park-like setting

The scene is lively and vibrant, with numerous



Real Image

[SBERT 0.833] The image depicts an outdoor scene where veral tents are set up on wh appears to be a paved area, possibly a parking lot or a similar open space. The tents are of various colors, including red, white, and orderly fashion. The tents are pitched on the ground, and some are supported by wooden s placed

Figure 17: Why do synthetic images achieve better text cycle consistency compared to real images? We visualize text cycle consistency from a real image vs. synthetic image. Compared to real images containing more complex detail, synthetic images only generate details *described in the input text* which occupy larger areas of the generated image. Therefore, such details are easier to reconstruct for the image-to-text model, resulting in better text reconstruction.

	Image Cycle Consistency				icy Image-Text Similarity	
Dataset	DreamSim	LPIPS	CLIP	MSE	CLIP	Image Reward
MME	0.705	0.475	0.806	0.247	0.748	0.427
SEEDBench	0.728	0.794	0.621	0.431	0.587	0.067
MMStar	0.640	0.599	0.544	0.477	0.617	0.228
Average	0.833	0.820	0.788	0.692	0.831	0.265

Table 7: Pearson correlation coefficient between similarity metrics and captioning performance. Perceptual similarity metrics, such as DreamSim and LPIPS, are the best predictors of captioning performance, followed by CLIP and MSE. DreamSim reconstruction generally shows a higher correlation than both imagetext similarity metrics.

	Text Cycle Consistency				Text-I	mage Similarity
Dataset	BARTScore	BERTScore	CLIP	SBERT	CLIP	Image Reward
T2I-CompBench	0.940	0.795	0.966	0.992	0.720	0.961
DrawBench	0.945	0.894	0.832	0.861	0.465	0.814
Average	0.954	0.825	0.952	0.979	0.675	0.944

Table 8: Pearson correlation coefficient between similarity metrics and text-to-image performance. Overall, text cycle consistency strongly correlates with text-to-image quality across all metrics. SBERT shows the highest correlation, followed by BARTScore and CLIP, all of which outperform image-text similarity metrics.

Ц С

847 848

846

849

850 Section 4.2, discusses how captioning performance mea-851 sured by VQA without V is strongly associated with both 852 image and text cycle consistency. Now, we examine if the 853 same correlation holds for model VQA and downstream 854 performance. For image-to-text models, we examine the 855 relationship between cycle consistency and VQA performance in the table below for both images and text. Note 856 that this is an evaluation of model performance, while 857 VQA without V evaluates the quality of the text gener-858 ated by the model. 859



860 For VQA performance, we use reported scores on benchmarks MMBench Liu et al. (2023c) and MME Fu et al. 861 (2023a). MME is split into perception and cognition 862 categories. Cycle consistency is computed on the sDCI 863



dataset and averaged across five different text-to-image models with 3 different random seeds. Fig-

864	Metric	Accuracy	Precision	Recall	F1 Score
865		riccuracy	Treeston	Itecuii	1150010
866	Image Cycle Consistency				
867	DreamSim	79.10	99.32	58.66	73.76
868	LPIPS	77.47	99.29	55.33	71.06
000	CLIP	77.99	99.18	56.44	71.94
809	MSE	74.64	99.07	49.74	66.23
870					
871	Image-text Similarity				
872	CLIP	78.04	99.30	56.49	72.01
873	Image Reward	75.92	99.13	52.30	68.48

874 Table 9: Top-1 hallucination and descriptiveness on MSCOCO. For a given image and a set of corresponding 875 captions, we select the top-1 caption based on each metric. We observe that DreamSim cycle consistency favors 876 captions with less hallucination and better descriptiveness compared to other metrics. Higher precision indicates 877 reduced hallucination, while higher recall reflects increased descriptiveness.

879	Fixed I2T Model	ICC $R^2$	TCC $R^2$
880	BLIP2-2.7B	0.836	0.349
881	BLIP2-6.7B	0.834	0.657
882	BLIP2-FlanT5-XXL	0.879	0.871
883	LLaVA1.5-7B	0.915	0.966
884	LLaVA1.5-13B	0.916	0.964
885	LLaVAOV-0.5B	0.932	0.201
886	LLaVAOV-7B	0.954	0.910
887	LLaVA1.6-7B	0.942	0.963
000	LLaVA1.6-34B	0.953	0.952
000	InternVL2-2B	0.904	0.935
889	InternVL2-8B	0.903	0.904
890	InternVL2-26B	0.880	0.879
891	InternVL2-40B	0.902	0.913
892	All Models	0.902	0.950

Fixed T2I Model	ICC $R^2$	TCC $R^2$
SD1.5	0.741	0.875
SDXL	0.759	0.845
SDXL-Turbo	0.731	0.879
SD3	0.790	0.870
FLUX-Time	0.794	0.861
All Models	0.766	0.864

878

Table 10: Pearson correlation coefficients between text-to-image performance and image cycle consistency (ICC) and text cycle consistency (TCC) for different fixed image-to-text models (Left) and different fixed text-to-image models (Right). Correlations are generally strong except when fixing BLIP2-2.7B or LLaVAOV-0.5B to calculate text cycle consistency. Higher  $R^2$  value indicates stronger correlation.

895

896

899 900

901

902

903

904

905

906

907

ures 18 shows plots of image and text cycle consistency vs VQA performance on each benchmark, with points representing different image-to-text models.

We find that image cycle consistency best correlates with VQA performance on MMBench and MME (cognition), with weaker association for other benchmarks. This may be somewhat surprising because MME (perception) in the VQA without V setting has a strong correlation with cycle consistency. Text cycle consistency did not have as strong as a correlation across all VOA benchmarks. It is important to note that VQA scores examine if the model can answer diverse questions about an image, while VQA without V examines if the text caption can answer diverse questions about the image. While these two evaluations are somewhat related, the difference between these tasks could account for the difference in correlation.

908 909 910

#### E DIVERGENCE IN GENERATED IMAGES

911 912

913 We use the DCI summarized captions dataset Huh et al. (2024) Urbanek et al. (2024) detailed in 914 Section 4.4 to compare the diversity of synthetic images generated from text based on the caption 915 density use to create them. For each image in the summarized DCI dataset, there are summary captions of different length. For each image, we generate generate 10 different images using random 916 seed sampling for every summary length. We then calculate the mean-pairwise distance between all 917 of the 10 generated images from the same summary caption. For each text-to-image model, we plot

<sup>893</sup> 894

<sup>897</sup> 898



Figure 18: Cycle consistency vs. Model Downstream Performance (VQA): Top Row (Left to Right): Image cycle consistency association with VQA scores on MMBench, MME (Perception), and MME (cognition). Bottom Row (Left to Right): Text cycle consistency association with VQA scores on MMBench, MME (Perception), and MME (cognition).

the mean-pairwise distance vs. the caption length seen in Figure 19. Considering all models, there is mixed consensus on how caption length affects generated image diversity.

## F ADDITIONAL RESULTS

We show additional "plug and play" examples of image and text cycle consistency in Figures 2021 respectively. Figure 22 provides more examples of increased text details hurting image reconstruction scores.

### G FAILURE CASES

We provide examples of failure cases of cycle consistency in Figures 23 and 24. Failures include: synthetic images with artifacts or implausible generations but little effect on captions, descriptions of non-existent objects, endpoint model failures (i.e. the intermediate image or text representation is reasonable but the endpoint model creates inaccuracies which affect reconstruction). Many of these mistakes can be attributed to model error and usually affect text cycle consistency much more than image, mainly because images generated from incorrect captions often have lower cycle consistency, whereas image-to-text models do not always notice inaccuracies in synthetic images.



Figure 20: Examples of image cycle consistency using different image-to-text, text-to-image combinations. The model combination used to generate the caption  $\rightarrow$  image is shown at the top of the image. DreamSim( $\downarrow$ ) distance between reconstructed images and the original are reported in brackets in front of the text captions.

A waterfall in a green body of water, surrounded by white stone statues and lush vegetation, including yellow-green willow tree leaves, light green leaves, and fluffy green tree foliage. The sky is blue with a few white clouds.

1010 1011 A restaurant interior with 1012 a colorful flower 1013 arrangement on a table, a buffet tray of appetizers 1014 on skewers (including 1015 cherry tomatoes and cantaloupe) surrounded 1016 by green lettuce, and 1017 tables of diners enjoying a meal in the background. 1018

[0.643] The The ceiling has white diving table lights illuminating the border and small circular lights throughout.

SD1.5

![](_page_18_Picture_6.jpeg)

[0.599] a

InternVL2-26B

 $SD1.5 \rightarrow InternVL2-26B$ 

[0.848] a waterfall with statues and SDXL-Turbo → BLIP2-Flan-T5

 $SDXL\text{-}Turbo \rightarrow BLIP2\text{-}Flan\text{-}T5$ 

![](_page_18_Picture_9.jpeg)

 $SD3 \rightarrow LLaVA1.6-34B$ 

[0.473]

![](_page_18_Picture_11.jpeg)

[0.777] The image

![](_page_18_Picture_14.jpeg)

[0.751]

SDXL

→ LLaVA OV-7B

[0.715] The image cap usly set v ning and

 $SDXL \rightarrow LLaVA OV-7B$ 

FLUX-Tin 5-13B

FLUX-Time  $\rightarrow$  LLaVA 1.5-13B

![](_page_18_Picture_17.jpeg)

[0.654] The image featu dining room with a long with various food items." adorned with a " larg including app as well oranges

Figure 21: Examples of text cycle consistency from various model combinations. The model combination 1022 used to generate the image  $\rightarrow$  text caption is shown at the top of each example's image. To make comparing 1023 the text easier, we highlight relevant phrases in the input and reconstructed descriptions. 1024

1025

1019

1020

1021

995

996

1003

1004

1005

1006

1007

1008

![](_page_19_Figure_1.jpeg)

measured with  $SBERT(\uparrow)$  due varying image-to-text model quality. While the synthetic image generated from the text seems reasonable, InternVL2-8B interprets the image correctly while the text from LLaVA 1.6-34B contains several errors leading to poor text cycle consistency.

**Input text**: A waterfall in a green body of water, surrounded by white stone statues and lush vegetation, including yellow-green willow tree leaves, light green leaves, and fluffy green tree foliage. The sky is blue with a few white clouds.

![](_page_20_Picture_2.jpeg)

**LLaVA 1.5-13B** [0.579]: The image features a beautiful garden with a waterfall and a pond. There are several statues of angels and other figures scattered throughout the garden, adding to its charm. A woman is sitting on a bench near the waterfall, enjoying the serene atmosphere. In addition to the woman, there are two other people in the scene, one standing near the...

**InternVL2-2B** [0.756]: The image depicts a serene and picturesque garden scene featuring a waterfall and several statues. The garden is lush and verdant, with a variety of greenery and plants surrounding the waterfall. The waterfall is a central feature, cascading down a rocky ledge into a pool below. The water appears to be a vibrant green, suggesting it might be a natural

Generated Image

**Input text**: A **red-roofed**, yellow-white **building** stands in the background, with a **balcony** and **rectangular windows**. In the foreground, a **tree** and **grassy area** with a small **fence** are visible, along with a **black car** parked in a **space**, its doors open and a man nearby. The **sky is light blue** with white clouds.

![](_page_20_Picture_7.jpeg)

#### BLIP2-FlanT5-XXL [0.628]: a yellow building with a red roof

**LLaVA 1.6-34B** [0.521]: The image features a **black car** parked in front of a **hotel**, with a man standing outside the vehicle. The car is parked in a parking space, and the man appears to be either getting out of the **car** or walking towards it. The hotel is a two-story **building** with a **red roof**, and there are several **windows** visible on both floors.

Generated Image

Figure 24: **Failure cases for text cycle consistency.** Text reconstruction is often affected by image captioning errors (top), and hallucinations and failures of the metric to address this (bottom).