The Association Between Training Data and Text-to-Image Generation Capabilities

Anonymous Author(s) Affiliation Address email

Abstract

Text-to-image (T2I) models are often touted for their supposed ability to create 1 2 compositional images with many components. However, these models can fail 3 to faithfully generate images when presented with prompts containing just two or three entities. In this work, we seek an explanation for such failures with 4 respect to the training data. We introduce the *training appearance ratio*, which 5 compares the number of training images depicting specific entities vs. the number 6 of training captions mentioning those same entities, and examine how well this 7 8 measure correlates with generation success rates. We find positive and significant correlations between these ratios and successful image generations. Furthermore, 9 our proposed measure yields stronger correlations with model success rates than 10 existing training data frequency measures. These associations suggest that our 11 proposed measure (training appearance ratio) better captures the relationship 12 between training data statistics and generation success. 13

14 **1 Introduction**

When asked to generate an image of "a bicycle and a skateboard", Stable Diffusion, a popular text-to-image (T2I) model [Rombach et al., 2022], succeeds only 8% of the time. Despite "bicycle" and "skateboard" being common objects that are generated separately nearly 100% of the time, the model fails to generate both jointly. The inability of models to handle such simple cases showcase their weak compositional capabilities.

In this work, we aim to explain models' failures with respect to their training data properties. 20 Drawing from previous works that have shown that pretraining data frequencies correlate with model 21 performance [Razeghi et al., 2022, Kandpal et al., 2023, Udandarao et al., 2024], we first seek to 22 replicate such findings for our setup of generating multiple common entities. However, our results 23 indicate that simple caption frequencies correlate poorly with models' generation success rates. Upon 24 digging into the training data, we observe that captions mentioning entities may pair with images that 25 only showcase a subset of those entities, or none at all, as shown in Figure 1b. For instance, there 26 are more than 9,000 captions in LAION2B-en [Schuhmann et al., 2022] that mention both "bicycle" 27 and "skateboard", but only 9% of corresponding images actually contain both objects. These findings 28 indicate that captions alone provide an inaccurate measure of how often entities are actually depicted 29 in training images. 30

Based on these findings, we adjust our frequencies to only consider training examples for which both the captions and images contain all specified entities, similar to Udandarao et al. [2024]. While these adjusted frequencies correlate better with models' generation success rates, they do not account for how T2I models are trained and utilized in practice (i.e., images are conditioned on texts). Therefore, we consider the ratio between entities appearance in training images vs. captions, which explicitly

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.



(a) Generated images for the prompt "a **bicycle** and a **skateboard**". The model (SD1.5) mostly generates *one of the two* objects (primarily bicycles).



(b) Training images where *either* **skateboard** or **bicycle** are shown, *but not both*. Many of these images depict parks and outdoors spaces that are suitable for both skateboarding and bicycling, but only include one.

Figure 1: Examples of generated/training images where prompts/captions mention "skateboard" and "bicycle", but corresponding images do not include both.

³⁶ incorporates this conditioning, and formalize this measure to be the *training appearance ratio*. We

³⁷ find that this ratio exhibits stronger correlations with model generation capabilities across various

combinations of models, prompts, and entities ($\rho = 0.43$ vs. 0.27 for 2 entities, and $\rho = 0.31$ vs.

 $_{39}$ 0.19 for 3 entities, averaged). These stronger correlations show that our measure better associates

40 success in generating images with the training data.

In summary, our work demonstrates that models are poor at basic compositional generations, and proposes a new training data measure that correlates better with models' success rates than existing approaches. Our findings suggest that simple training appearance ratios help better understand model behavior, at least in part, and establish a foundation for future work that investigates concrete and measure that correlates are poor at basic compositional generations, and the behavior at least in part, and establish a foundation for future work that investigates concrete and

45 comprehensive explanations for model failures and successes.

46 2 Explaining Successes Through Training Data Statistics

T2I models often fail to generate images following simple prompts with multiple common entities. 47 Our main goal in this study is to investigate to what extent models' ability to faithfully generate 48 images from prompts can be attributed to statistics from their training data. To address this objective, 49 we need to first define how we measure and compare training data statistics and image generation 50 success. Consider a training dataset $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), ..., (x_N, y_N)\}$ consisting of N (image, 51 caption) pairs. We also assume a prompt p that instructs the model to generate some entities 52 $e = \{e_1, e_2, \dots, e_k\}$, where $\forall i, e_i \in p$. To identify relevant examples from \mathcal{D} we select training 53 captions that mention the entities e specified in p. For example, for the prompt "a bicycle and a 54 skateboard", we query from \mathcal{D} and choose image-caption pairs whose captions include the entities 55 "bicycle" and "skateboard". 56

Note that while entities e may appear in a caption y_i , the image x_i corresponding to that caption may not contain all entities (sometimes even none), as depicted in Figure 1, and as was observed in Udandarao et al. [2024].¹ Since raw counts provide a biased estimation of entity occurrences in images, we instead propose measuring the proportion of captions whose images also contain all specified entities.

¹Table 5 (Appendix) shows example image-caption pairs that fall under these cases.

⁶² We define this quantity to be the training appearance ratio $(tar_{e,ic})$:

$$tar_{e,ic} = \frac{|\mathcal{D}_{e,i}|}{|\mathcal{D}_{e,c}|}$$

where $\mathcal{D}_{e,c}$ is the subset of \mathcal{D} whose captions contain entities e, and $\mathcal{D}_{e,i}$ is the subset of \mathcal{D} whose captions and images contain entities e. A higher value of $tar_{e,ic}$ indicates that image-caption pairs that mention a set of entities in captions also tend to include those entities in images.

After computing $tar_{e,ic}$, we generate images for prompt p using a T2I model to obtain generated images $\mathcal{G}_{e,p}$. We calculate the proportion of images that depict all entities, which we call the generation appearance ratio $(gar_{e,ip})$.

$$gar_{e,ip} = \frac{|\mathcal{G}_{e,i}|}{|\mathcal{G}_{e,p}|}$$

Similar to above, $\mathcal{G}_{e,i}$ is the subset of generated images whose prompts and images contain entities 69 e. We then examine whether the generation appearance ratio of generated entities that are explicitly 70 specified in prompts $(gar_{e,ip})$ correlates with corresponding ratios from the training data $(tar_{e,ic})$. 71 While previous works highlight correlations between model behavior and frequencies in the data 72 [Razeghi et al., 2022, Kandpal et al., 2023, Udandarao et al., 2024], we hypothesize that training 73 appearance ratios exhibit stronger associations with model generation capabilities, since $tar_{e,ic}$ 74 75 directly captures discrepancies in how often entities occur in training images vs. texts (similar to 76 how $gar_{e,ip}$ captures discrepancies in how often entities occur in generated images vs. prompts). In other words, we argue that $tar_{e,ic}$ more closely matches what we measure at generation, resulting in 77 stronger correlations as we show in Section 4. 78

79 **3** Experimental Setup

Entities We select entities from the MS COCO dataset [Lin et al., 2014] classes in addition to
 manually added entities (e.g., fruits, vegetables) as shown in Table 3 (Appendix), resulting in 84
 entities. We intentionally focus on frequent entities that models succeed in generating individually,
 and whose presence or absence are easy to evaluate.

Automated Image Evaluation To determine whether an image contains specified entities, we
utilize an automated approach. We use visual question answering (VQA) and employ PaliGemma
[Google, 2024] as our VQA model. More specifically, we ask the model whether an image contains a
given entity, which is done for all entities in the prompt, and consider an image to contain all entities
if the model answers "yes" for every entity. Note that PaliGemma achieves 91% on human annotated
images, as discussed in Appendix A.5.

Entity Caption Occurrences We use WIMBD [Elazar et al., 2024] to retrieve counts of entities from the training data. Specifically, we extract captions that mention a set of entities ($\mathcal{D}_{e,c}$), and randomly sample up to 1,000 image-caption pairs. Based on the corresponding images, we calculate the proportion of images that depict the specified entities to measure $tar_{e,ic}$. We multiply the number of captions ($|\mathcal{D}_{e,c}|$) by the ratios computed previously, $tar_{e,ic}$, to estimate the number of training examples that both mention entities in captions and include them in images.

Prompts We prompt the model to generate images with one, two, and three entities using the prompts shown in Table 4 in Appendix A.1. For each prompt, we generate 50 images using different random seeds, resulting in 100 images total for single entity prompts and 200 images total for double and triple entity prompts.

Data & Models We focus on Stable Diffusion [Rombach et al., 2022], a popular set of text-to-image models. Specifically, we use SD1.1 and SD1.5, which are both trained on 2.3 billion image-caption pairs filtered to contain only English captions (LAION2B-en). Additionally we use SD2.1, which is trained on LAION-5B [Schuhmann et al., 2022], a dataset of 5.9 billion multilingual image-captions pairs (including LAION2B-en). Notably, we use Stable Diffusion models for our study because they are the only models with publicly available training datasets, which is necessary for deriving connections between training dataset properties and model behavior.



Figure 2: Correlations between generation appearance ratios $(gar_{e,ip})$ and training appearance ratios $(tar_{e,ic})$ for 1, 2, and 3 entities, shown for SD1.1 and prompt 1. We bin examples into 10 equallysized groups or deciles based on $tar_{e,ic}$ and compute median $tar_{e,ic}$ and $gar_{e,ip}$ values for each bin, which correspond to the navy blue points.

Model	1 Entity	2 Entities	3 Entities
SD1.1	0.98	0.44	0.18
SD1.5	0.99	0.50	0.21
SD2.1	0.96	0.66	0.32

Table 1: Generation appearance ratios $(gar_{e,ip})$ for different models and # of entities, averaged across prompts.

107 4 Results

Generation Appearance Ratios How good are models at compositional generation? To answer 108 this question, we examine generation appearance ratios $(gar_{e,ip})$, which capture the success rate of 109 generating images with all specified entities, for different models and number of entities (Table 1). 110 We find that all models successfully generate single entities > 96% of the time, validating that models 111 are capable of generating common individual entities. However, models exhibit massive drops when 112 generating two and three entities – for example, both SD1.1 and SD1.5 models generate two entities 113 <= 50% of the time. Although SD2.1 is notably better at generating two entities (at nearly 66%), it 114 still struggles in this simple compositional setting. In summary, we see that models fail increasingly 115 as prompts depict more entities. We do not go beyond 3 entities, since Stable Diffusion generates 116 four entities < 5% of the time. 117

Correlations between Model Behavior and Training Data Statistics We wish to explain model success rates in generating various entities with respect to the training data. To do so, we first analyze frequency-based approaches, building on related work that explores the impact of training data in different settings [Razeghi et al., 2022, Udandarao et al., 2024]. We then show that our proposed measure ($tar_{e,ic}$) is more strongly correlated with model behavior.

Baselines: Frequency-based Approaches As baselines, we compute Pearson's correlation between $gar_{e,ip}$ and (1) frequencies of entities in captions and (2) estimated frequencies of entities in images (counts multiplied by $tar_{e,ic}$). Following Udandarao et al. [2024], we compute the log_{10} of frequencies to capture log-linear associations, and refer to the resulting correlations as ρ_{cap} and ρ_{im} . Results are presented in the first two sections of Table 2 for various models and number of entities, averaged across prompts.

We find that ρ_{cap} is not statistically significant (significance level < 0.01) across all combinations of 129 models, prompts, and number of entities except for SD1.1 with one entity. For the overwhelming 130 majority of cases, raw captions counts do not correlate with $gar_{e,ip}$. These results are unsurprising, 131 since raw caption counts are poor indicators of how often entities actually occur in training images. 132 We observe negative correlations for ρ_{cap} in the three entity case, which is somewhat surprising, but 133 these values are not statistically significant. In contrast, ρ_{im} exhibits consistently positive correlations 134 for two and three entities, and is statistically significant across all prompts and models in the two 135 entity case. When comparing ρ_{im} values for two and three entities, we observe a clear reduction 136 in ρ_{im} across models (0.08 absolute decrease). This reduction may be due to models exhibiting 137

		Number of Entities		
Corr	Model	1	2	3
	SD1.1	**0.37	0.06	-0.12
ρ_{cap}	SD1.5	0.12	0.07	-0.06
	SD2.1	0.20	0.02	-0.06
	SD1.1	**0.40	**0.31	0.20
$ ho_{im}$	SD1.5	0.18	**0.28	0.17
	SD2.1	0.26	**0.23	0.21
	SD1.1	0.17	**0.47	**0.34
ρ_{ratio}	SD1.5	0.29	**0.42	**0.28
(ours)	SD2.1	0.23	**0.40	**0.30

Table 2: Pearson's correlation coefficients between generation appearance ratios and various training data measures: (1) frequency of entities in captions (ρ_{cap}) as a baseline, (2) estimated frequency of entities in images (ρ_{im}) as another baseline, and (3) our proposed measure (ρ_{ratio}), averaged across prompts. We compute the \log_{10} of frequencies for (1) and (2) to capture log-linear associations. ** indicates correlations are statistically significant (significance level < 0.01) for all prompts.

poor generation capabilities as a whole for three entities. Overall, these findings indicate that frequency-based measures may not be effective in capturing generation success for multiple entities.

Proposed Measure: Training Appearance Ratios We present correlation results between $gar_{e,ip}$ and $tar_{e,ic}$ in the last section of Table 2 (ρ_{ratio}). We find that all models exhibit positive, but not statistically significant correlations for single entities. Since we select frequently occurring entities by design, we can expect models to generate them successfully irrespective of $tar_{e,ic}$.

For prompts with two and three entities, we observe positive and statistically significant correlations 144 across all models, prompts, and number of entities. Both Figures 2b (two entities, and 2c (three enti-145 ties) show linear associations between generation and training appearance ratios. These associations 146 become much clearer when data points are binned into deciles based on $tar_{e.ic}$, with $\rho_{ratio}=0.95$ for 147 2 entities and $\rho_{ratio}=0.90$ for 3 entities. We observe some variability across prompts with $\sigma \leq 0.07$ 148 for two entities and $\sigma \leq 0.06$ for three entities. Similar to ρ_{im} , we see a decrease in ρ_{ratio} going 149 from two to three entities (0.12 absolute decrease). That being said, ρ_{ratio} consistently exhibits 150 statistical significance and higher values relative to ρ_{im} . Overall, these results suggest ρ_{ratio} is a 151 stronger indicator of successful generations for compositional prompts depicting multiple entities. 152

153 5 Conclusion

This work studies the connection between models' generation success and training appearance 154 ratios. Although numerous studies have shown that model performance strongly correlates with 155 the frequency of entities [Razeghi et al., 2022, Kandpal et al., 2023, Udandarao et al., 2024], we 156 show that for image generation, successful generations correlate better with the proportion of training 157 captions that depict specified entities. Our findings are complemented by Seshadri et al. [2024], who 158 also show that model generations are associated with ratios from the training data in the context of 159 gender-occupation biases. Our results emphasize the need for improving data quality by limiting 160 image-caption mismatches and further necessitate open access to pretraining corpuses to be able to 161 characterize model behaviors and their flaws. 162

163 **References**

J. Betker, G. Goh, L. Jing, T. Brooks, J. Wang, L. Li, L. Ouyang, J. Zhuang, J. Lee, Y. Guo, et al.
 Improving image generation with better captions. 2023. URL https://cdn.openai.com/
 papers/dall-e-3.pdf.

Y. Elazar, A. Bhagia, I. H. Magnusson, A. Ravichander, D. Schwenk, A. Suhr, E. P. Walsh, D. Groen eveld, L. Soldaini, S. Singh, H. Hajishirzi, N. A. Smith, and J. Dodge. What's in my big
 data? In *The Twelfth International Conference on Learning Representations*, 2024. URL
 https://openreview.net/forum?id=RvfPn0kPV4.

D. Ghosh, H. Hajishirzi, and L. Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment, 2023.

Google. Big vision: Paligemma project configurations. https://github.com/
 google-research/big_vision/tree/main/big_vision/configs/proj/paligemma,
 2024. Accessed: 2024-06-14.

 J. Hessel, A. Holtzman, M. Forbes, R. Le Bras, and Y. Choi. CLIPScore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.595.
 URL https://aclanthology.org/2021.emnlp-main.595.

 Y. Hu, B. Liu, J. Kasai, Y. Wang, M. Ostendorf, R. Krishna, and N. A. Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 20406–20417, October 2023. URL https://openaccess.thecvf.com/content/ICCV2023/papers/ Hu_TIFA_Accurate_and_Interpretable_Text-to-Image_Faithfulness_Evaluation_ with_Question_Answering_ICCV_2023_paper.pdf.

K. Huang, K. Sun, E. Xie, Z. Li, and X. Liu. T2i-compbench: A comprehensive benchmark for
 open-world compositional text-to-image generation. *arXiv preprint arXiv: 2307.06350*, 2023.
 URL https://arxiv.org/pdf/2307.06350.

N. Kandpal, H. Deng, A. Roberts, E. Wallace, and C. Raffel. Large language models struggle
 to learn long-tail knowledge. In *International Conference on Machine Learning*, 2023. URL
 https://proceedings.mlr.press/v202/kandpal23a/kandpal23a.pdf.

 T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick.
 Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755.
 Springer, 2014. URL https://api.semanticscholar.org/CorpusID:14113767.

A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
 J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language
 supervision. In *International Conference on Machine Learning*, 2021. URL https://api.
 semanticscholar.org/CorpusID:231591445.

 A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. URL https://api.semanticscholar.
 org/CorpusID:248097655.

R. Rassin, E. Hirsch, D. Glickman, S. Ravfogel, Y. Goldberg, and G. Chechik. Linguistic binding
 in diffusion models: Enhancing attribute correspondence through attention map alignment. In
 Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https://
 openreview.net/forum?id=AOKU4nRw1W.

Y. Razeghi, R. L. Logan IV, M. Gardner, and S. Singh. Impact of pretraining term frequencies on few-shot numerical reasoning. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 840–854, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
findings-emnlp.59. URL https://aclanthology.org/2022.findings-emnlp.59.

R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
 with latent diffusion models. In 2022 IEEE/CVF Conference on Computer Vision and Pattern
 Recognition (CVPR), pages 10674–10685. IEEE Computer Society, jun 2022. URL https:
 //doi.ieeecomputersociety.org/10.1109/CVPR52688.2022.01042.

C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta,
C. Mullis, M. Wortsman, P. Schramowski, S. Kundurthy, K. Crowson, L. Schmidt, R. Kaczmarczyk, and J. Jitsev. Laion-5b: An open large-scale dataset for training next generation
image-text models. In *Advances in Neural Information Processing Systems*, volume 35, pages
25278-25294, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/
file/a1859debfb3b59d094f3504d5ebb6c25-Paper-Datasets_and_Benchmarks.pdf.

P. Seshadri, S. Singh, and Y. Elazar. The bias amplification paradox in text-to-image generation.
In K. Duh, H. Gomez, and S. Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6367–6384, Mexico City, Mexico, June 2024.
Association for Computational Linguistics. doi: 10.18653/v1/2024.naacl-long.353. URL https:
//aclanthology.org/2024.naacl-long.353.

V. Udandarao, A. Prabhu, A. Ghosh, Y. Sharma, P. H. Torr, A. Bibi, S. Albanie, and M. Bethge.
 No "zero-shot" without exponential data: Pretraining concept frequency determines multimodal
 model performance. *ArXiv*, abs/2404.04125, 2024. URL https://api.semanticscholar.
 org/CorpusID:268987507.

X. Wu, D. Yu, Y. Huang, O. Russakovsky, and S. Arora. Conceptmix: A compositional image
 generation benchmark with controllable difficulty, 2024. URL https://arxiv.org/abs/2408.
 14339.

M. Yuksekgonul, F. Bianchi, P. Kalluri, D. Jurafsky, and J. Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International*

238 Conference on Learning Representations, 2023. URL https://openreview.net/forum?id=

239 KRLUvxh8uaX.

240 Limitations

We compare properties in the training data with model behavior using correlational analysis and observe clear trends: higher training appearance ratios are associated with higher generation successes. However, we cannot assert that our measure explains or definitively impacts model behavior without employing a causal approach, and leave this important direction to future work.

Our results suggest that different entity combinations with similar training appearance ratios can have variable generation success rates. Although correlations between training appearance ratios and model success rates are consistently positive and significant in the two and three entity settings, they are weakly to moderately positive. These results suggest that simple training appearance ratios offer some insights into models' generation capabilities, but do not provide the full story. Perhaps there are more nuanced training data measures to consider, or other factors beyond the data such as model scale, architecture, and training.

Along these lines, it is worth noting that closed models such as DALL-E 2 [Ramesh et al., 2022], and especially DALL-E 3 [Betker et al., 2023], are much better at handling compositional prompts. While we do not know the exact factors that contribute to this improvement, we speculate that training data quality and curation play a huge role. Perhaps the image-caption pairs used to train such models were filtered or augmented to have much higher training appearance ratios as a whole. However, without access to such datasets, it is unclear to what extent training appearance ratios are a driving force behind more capable models.

In addition, we focus on the specific setup of generating between 1-3 entities, which is a fundamental 259 aspect of compositional understanding. As we show, models fail considerably even in this simple 260 setting. However, there are other well-known failure modes [Ghosh et al., 2023, Huang et al., 2023, 261 Rassin et al., 2023] in text-to-image generation that should be considered, as well as more complex 262 263 compositions [Wu et al., 2024]. Furthermore, our study focuses exclusively on English prompts. We encourage researchers to study the association between training data and text-to-image generation 264 for other languages. This study is among the first to investigate text-to-image failure modes with 265 respect to training data, and we hope that it motivates future work to further probe and expand on 266 these findings. 267

		Entities		
airplane	apple	asparagus	backpack	banana
bear	bed	beet	bench	bicycle
bird	boat	book	bottle	bowl
broccoli	bus	cake	car	carrot
cat	chair	clock	coconut	corn
couch	cow	cup	daisy	dog
donut	elephant	fork	garlic	giraffe
grapes	handbag	horse	hydrangea	iris
kale	keyboard	kite	knife	laptop
lily	lime	mango	microwave	motorcycle
onion	orchid	oven	peony	pineapple
pizza	pomegranate	refrigerator	remote	rose
sandwich	sheep	sink	skateboard	skis
snowboard	spoon	strawberry	suitcase	sunflower
surfboard	tie	toaster	toilet	tomato
toothbrush	train	truck	tulip	tv
umbrella	vase	watermelon	zebra	

Table 3: List of 84 common entities used to study models' ability to generate multiple entities.

268 A Appendix

269 A.1 Prompts

The prompts used for generating images are presented in Table 4. For each prompt, we have the following number of instances (i.e., entity combinations after filling in [E1], [E2], [E3]): we have 84 instances for 1 entity, 440 instances for 2 entities, and 440 instances for 3 entities.

# Entities	Prompt
1	1. a/an [E1] 2. a photo of a/an [E1]
2	1. a/an [E1] and a/an [E2] 2. a photo of a/an [E1] and a/an [E2] 3. [E1], [E2] 4. a/an [E1] next to a/an [E2]
3	 a/an [E1] and a/an [E2] and a/an [E3] a photo of a/an [E1] and a/an [E2] and a/an [E3] [E1], [E2], [E3] a/an [E1] next to a/an [E2] and a/an [E3]

Table 4: Image generation prompts for single, double, and triple entities. [E1], [E2], and [E3] are replaced with various entities (e.g., elephant, zebra, and giraffe).

272

273 A.2 Image Generation

This work uses 3 Stable Diffusion versions: SD1.1 and SD1.5 (trained on LAION2B-en) and SD2.1 (trained on LAION-5B). We use the default generation parameters of 50 inference steps and a guidance scale of 7.5. We specify a batch size of 4. For a given instance of a prompt (i.e., filled in with entities) and model version, we generate 50 images using different random seeds. In total, our generations have taken ~600 hours in total on a single TITAN RTX GPU.



Figure 3: Correlations between training appearance ratios $(tar_{e,ic})$ for LAION2B-en and LAION-5B for 1, 2, and 3 entities. We observe strong correlations for all three.

279 A.3 Entities

The entities used to fill in prompts are presented in Table 3. We include 84 entities in total. The minimum count in in the dataset is for the word "beet" with 123,134 caption mentions for LAION2Ben and 194,530 caption mention for LAION5B. The maximum count is for the word "book" with 21,353,659 caption mentions in LAION2B-en and 28,379,268 for LAION5B.

284 A.4 VQA

For performing automated image evaluation, a common choice is to use CLIPScore [Hessel et al., 2021]. However, CLIP [Radford et al., 2021], its underlying model, struggles with compositional understanding [Hu et al., 2023, Yuksekgonul et al., 2023] and performs poorly for such prompts. As a result, we turn to Visual Question Answering (VQA). We ask a separate question for each entity using the following format: "Is there a/an [entity] in this image, yes or no?", which is then asked for all entities in the prompt. If the model responds "yes" to each of the questions, we consider the image to contain all specified entities. This approach is used for both training and generated images.

292 A.5 Human Evaluation

We perform human evaluation to assess whether our VQA approach is appropriate and effective for evaluating the presence of entities in images. The authors of this paper labeled 400 randomly selected generated images in the two entity setting, providing annotations for entity1 and entity2. We find that PaliGemma predictions match human annotations in 90.88% of cases, which indicates strong performance. The biggest disagreements between human annotations and model predictions tend to be cases for which entities are similar in appearance and use cases (e.g., backpack and handbag), as well as large size differences (e.g., toothbrush and snowboard).

300 A.6 Comparing Training Appearance Ratios

As shown in Figure 3, training appearance ratios calculated using LAION2B-en and LAION5B are highly correlated. While this is perhaps not surprising given that we focus exclusively on English and LAION2B-en is a subset of LAION5B, it is worth noting that these ratios are preserved across both datasets for the entity combinations we consider.



(a) Generated images using SD2.1 with the prompt "a toothbrush and a sink" ($gar_{e,ip}$ =0.44).



(b) Training images whose captions mention both "sink" and "toothbrush" ($tar_{e,ic}$ =0.44).

Figure 4: We sample generated and training images for the prompt "a toothbrush and a sink". Both the generation and training appearance ratios are the same. We see that generated images depicting one entity tend to show sinks, while training images depicting one entity show both toothbrush and sink individually.



(a) Generated images using SD2.1 with the prompt "a watermelon and a handbag" ($gar_{e,ip}=0.48$).

(b) Training images whose captions mention both "watermelon" and "handbag" (*tare*,*ic*=0.46).

Figure 5: We sample generated and training images for the prompt "a watermelon and a handbag". Both the generation and training appearance ratios are very similar. We see that generated images seem to always depict watermelons, and sometimes handbags (with appearances similar to a watermelon). While some training images are watermelon handbags, other examples may depict accessories or watermelon-colored handbags.



(a) Generated images using SD2.1 with the prompt "a giraffe and a bear" ($gar_{e,ip}$ =0.46).



(b) Training images whose captions mention both "giraffe" and "bear" ($tar_{e,ic}$ =0.43).

Figure 6: We sample generated and training images for the prompt "a giraffe and a bear". We observe that while the generation and training appearance ratios are highly similar, the ways in which entities are depicted at generation and training differ quite noticeably (e.g., training images mostly show toys or cartoons).



(a) Generated images using SD1.5 with the prompt "a motorcycle and a bench" ($gar_{e,ip}=0.08$).



(b) Training images whose captions mention both "motorcycle" and "bench" ($tar_{e,ic}$ =0.08).

Figure 7: We sample generated and training images for the prompt "a motorcycle and a bench". The generation and training appearance ratios are identical. At generation, the model generates images of motorcycles individually a clear majority of the time. The training data, however, also includes images of benches individually as well as images without either entity.



(a) Generated images using SD1.5 with the prompt "a photo of a bus and a horse" ($gar_{e,ip}$ =0.18).



(b) Training images whose captions mention both "motorcycle" and "bench" ($tar_{e,ic}$ =0.21).

Figure 8: We sample generated and training images for the prompt "a photo of a bus and a horse". The generation and training appearance ratios are very close. At generation, the model often generates buses individually, specifically red buses. While training images also depict buses individually in several cases, they seem to capture a more diverse set of buses.



(a) Generated images using SD1.5 with the prompt "elephant, daisy" ($gar_{e,ip}$ =0.24).

(b) Training images whose captions mention both "elephant" and "daisy" ($tar_{e,ic}=0.30$).

Figure 9: We sample generated and training images for the prompt "elephant, daisy". The generation and training appearance ratios are fairly close. At generation, the model mostly depicts elephants individually, and they look reasonably realistic. In training images, we mainly see artistic renditions of elephants.



(a) Generated images using SD1.5 with the prompt "boat, chair" ($gar_{e,ip}$ =0.16).



(b) Training images whose captions mention both "boat" and "chair" ($tar_{e,ic}$ =0.19).

Figure 10: We sample generated and training images for the prompt "boat, chair". The generation and training appearance ratios are fairly close. At generation, the model primarily depicts a boat or chair, often individually, in an outdoor setting. In training images, while we see some entities in outdoor setting, many just depict a chair in a staged setting.

Image	Caption	VQA Predictions
	How To Make An Asparagus Bed	asparagus: yes, bed: no
	Bluetooth Speaker Panda with Remote Shutter Release White 4.3x4.5cm	panda: yes, remote: no
	Candy apple Red Volkswagen bus for couple and bridal party at waterfront wedding	apple: no, bus: yes
	Sweet potato, coconut and tomato lentil dahl in a bowl beside a bowl of cherry tomatoes	coconut: no, tomato: yes
	Extreme BMX Bicycle Riding in Concrete Skateboard Park - Bar spin to tire tap Stock Footage	bicycle: yes, skate- board: no
K	Lily the Borzoi chasing other dog	lily: no, dog: yes
	LED Waterproof RGB Colorful Wedding Party Vase Base Light Submersible+ Remote	vase: no, remote: yes
	An elephant cow taking a dust bath with her calf (Kruger National Park, South Africa).	elephant: yes, cow: no
	Collapsible Chair From Skis Ski Woodcraft Pinterest	chair: yes, skis: no
	Jungle Animal Shapes - Cake Toppers or Party Decorations monkey giraffe lion elephant tiger zebra snake hippo baby shower birthday party	cake: no, giraffe: yes

Table 5: Example training images and captions for which captions mention two specified entities (captions may mention other entities as well), but images only depict one of the specified entities clearly. Specified entities are in **bold**. One potential explanation for such occurrences is the ambiguity of words (e.g., "Lily" is both a name and a flower). Another explanation is that a combination of entities may have their own meaning (e.g., "asparagus bed" is not the same as "asparagus" + "bed").