

A Picture is Worth a Thousand Words: Principled Recaptioning Improves Image Generation

Eyal Segalis
Google Research
eyalis@google.com

Dani Valevski
Google Research
daniv@google.com

Danny Lumen
Google Research
dwasserman@google.com

Yossi Matias
Google Research
yossi@google.com

Yaniv Leviathan
Google Research
leviathan@google.com

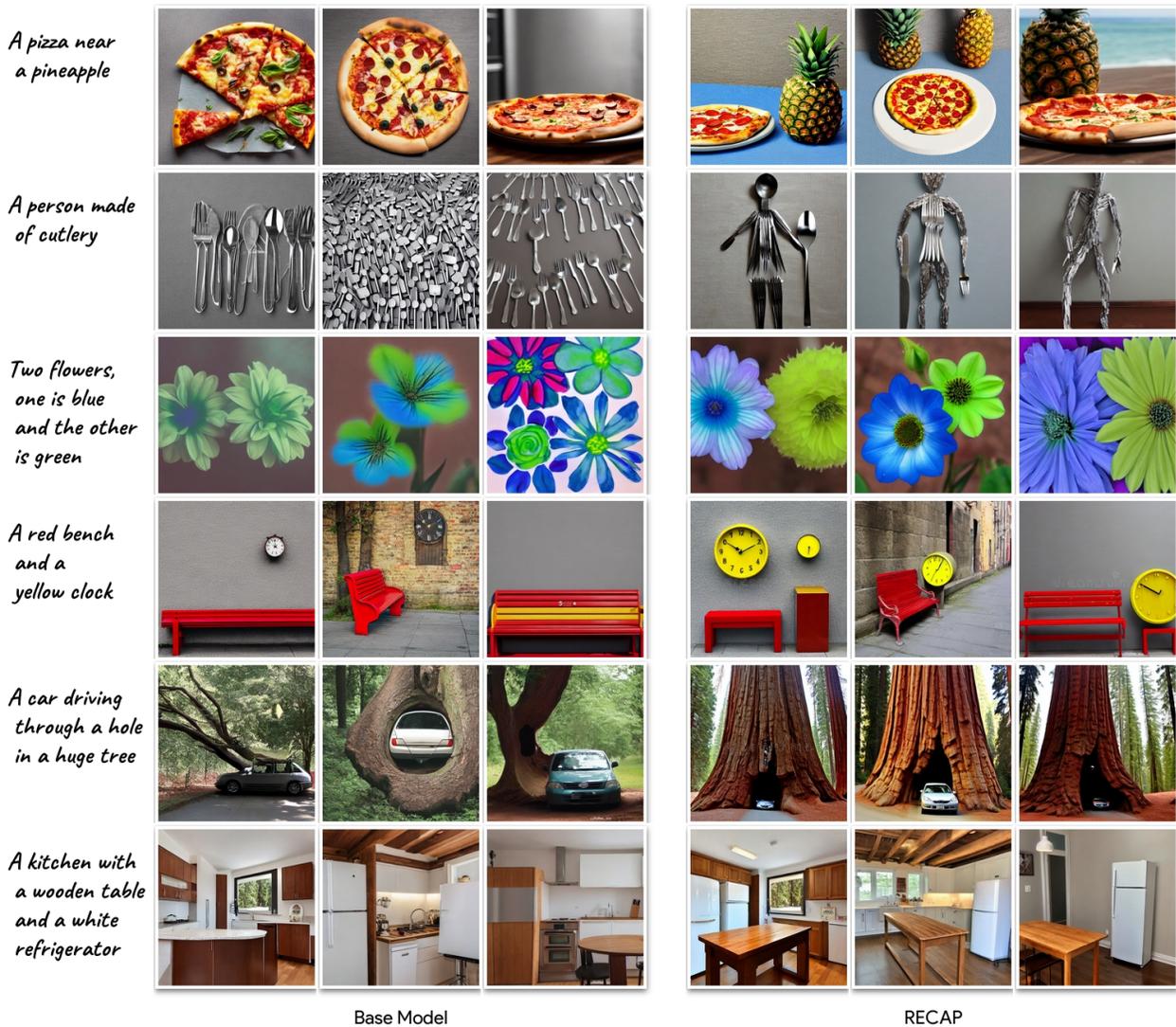


Figure 1. Examples of non-cherrypicked generations from the base Stable Diffusion model (left) and our model trained on a recaptioned dataset (right), on the same set of random seeds.

1. Abstract

Text-to-image diffusion models achieved a remarkable leap in capabilities over the last few years, enabling high-quality and diverse synthesis of images from a textual prompt. However, even the most advanced models often struggle to precisely follow all of the directions in their prompts. The vast majority of these models are trained on datasets consisting of (image, caption) pairs where the images often come from the web, and the captions are their HTML alternate text. A notable example is the LAION dataset, used by Stable Diffusion and other models. In this work we observe that these captions are often of low quality, and argue that this significantly affects the model’s capability to understand nuanced semantics in the textual prompts. We show that by relabeling the corpus with a specialized automatic captioning model and training a text-to-image model on the recaptioned dataset, the model benefits substantially across the board. First, in overall image quality: e.g. FID 14.84 vs. the baseline of 17.87, and 64.3% improvement in faithful image generation according to human evaluation. Second, in semantic alignment, e.g. semantic object accuracy 84.34 vs. 78.90, counting alignment errors 1.32 vs. 1.44 and positional alignment 62.42 vs. 57.60. We analyze various ways to relabel the corpus and provide evidence that this technique, which we call RECAP, both reduces the train-inference discrepancy and provides the model with more information per example, increasing sample efficiency and allowing the model to better understand the relations between captions and images.

2. Introduction

In recent years, text-to-image (T2I) generation models such as Imagen [1], Muse [2], Dall-E [3], Dall-E 2 [4], Parti [5], and Stable Diffusion [6] have undergone significant advancements. This progress has enabled the generation of remarkably high-quality and diverse images by conditioning on textual inputs. However, while revolutionary, even modern state-of-the-art text-to-image models may fail to generate images that fully convey the semantics and nuances from the given textual prompts. Failure modes include: missing one or more subjects from the input prompt [5, 7]; incorrect binding of entities and modifiers [5, 7, 8]; and incorrect placement and spatial composition of entities [5, 9, 10].

In this work we first observe that open-web datasets used to train open text-to-image models suffer from significant issues. For example, the captions in the LAION [11] dataset, used to train Stable Diffusion, come from alt HTML tags (Alttext). According to W3C’s web content accessibility guidelines¹, the alt attribute is used to convey the meaning and intent of the image, and not necessarily being a literal

¹<https://www.w3.org/TR/2016/NOTE-WCAG20-TECHS-20161007/H37>

description of the image itself. Indeed, we observe, that often the Alttext describes only a narrow aspect of the image, neglecting significant visual details. For example, an image of a person can have as Alttext the name of the person and the name of the photographer, but not a description of their appearance, their clothes, their position, or the background. Also, sometimes Alttext tags contain inaccuracies, mistakes and out of context information. See Fig. 4 for examples.

We further observe that while trained mainly on similar datasets of open (image, caption) pairs, recent automatic captioning systems, such as PaLI [12], produce highly accurate captions. See examples in Fig. 3. This may be due to the fact that the inverse problem of image-to-text (I2T) is easier, or to the fact that these captioning models are larger, are trained longer than the T2I models, or leverage large pre-trained language components.

With these observations, we suggest a new method for horizontally improving T2I models by training them on improved captions, auto-generated by a custom I2T model. We call our method RECAP, and show that applying it to Stable Diffusion results in a model that is better than the baseline across the board, with a battery of standard metrics substantially improving, e.g. FID 17.87→14.84, as well as in human evaluation of successful image generation 29.25% → 48.06% (see Section 5).

In Sec. 4 we provide the details of our method, RECAP. Sec. 5 discusses our results and shows the horizontal improvements both in image quality as well as in semantic fidelity to the prompts. In Sec. 6 we analyze the issues with the original captions, demonstrate that the improvements are indeed due to the new captions, and that they arise as a result of both minimizing the train-test skew as well as increasing sample efficiency.

To summarize, our main contributions are:

- A new method we call RECAP, that leverages automatic captioning to improve the quality of a text-to-image model in a substantial way horizontally, both in fidelity and semantics, measured on a set of 7 standard metrics as well as with human evaluations.
- An analysis showing how Alttext captions used by current training methods suffer from train-inference skew and lack in semantic details, resulting in text-to-image models that often fail in fidelity and semantic alignment, and how different captions mitigate both issues.

3. Related Work

Text-to-Image Models. Deep generative models for image generation from text have shown notable progress in recent years, transitioning from using GAN-based methods [13] to using autoregressive transformers [2, 3, 5] and diffusion models [1, 4, 14]. An important area of enhancement is in improving a model’s capability to align with the input text effectively. Methods condition the image on the output

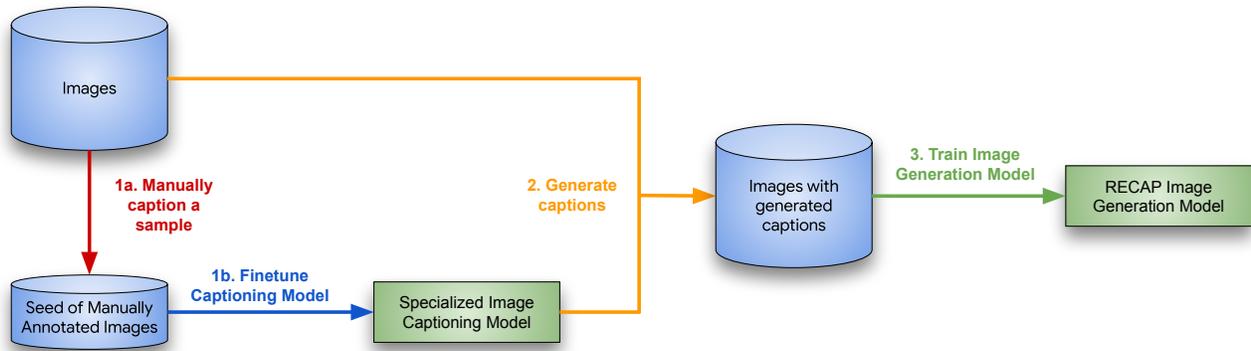


Figure 2. Schematic diagram of our method RECAP. In steps (1a) and (1b) we fine-tune an image-to-text captioning model on a small set of detailed human captions. In step (2) we use this fine-tuned model to recaption the images in the training dataset of a text-to-image model, and with this dataset, in step (3) we train an image generation model with the recaptioned dataset.

of a pre-trained text embedder, usually CLIP [15]. Imagen [1] shows that using a strong T5 [16] encoder significantly improves text-image alignment. Parti [5] uses a SimVLM model [17] to annotate some of the training set images.

Image captioning. Image captioning is a fundamental problem in computer graphics. Recent models [12, 18–21] have made significant progress on this task, thanks to high scale training data and utilizing pre-trained image and text models, allowing them to address a variety of multi-modal tasks. In this work we fine-tune PaLI [12] to recaption the training dataset of a text-to-image model.

Synthetic multimodal datasets. Some concurrent works augment multimodal dataset via automatic means to improve the capabilities of image captioning and embedding models. Nguyen et al. [22] show that training CLIP on data with generated captions improves its performance. Li et al. [23] show how to improve the BLIP model by iteratively removing captions with high loss, replacing them with captions from the previous epoch or generating a new image for them using Stable Diffusion. Ma et al. [24] generate a synthetic text-to-image model using Stable Diffusion and then uses it to train a captioning algorithm.

Improving text-alignment of diffusion models. An important line of work attempts to rectify text-to-image alignment issues in diffusion models. This is usually done by lexical analysis of the input prompt modifying the attention maps throughout sampling. Chefer et al. [7] try to prevent the diffusion models from ignoring certain tokens by guiding the sample process to reweigh the token attention maps. Rassin et al. [8] analyze the prompt to identify modifiers, and then uses a sample-time guidance to bind their attention maps to those of their entities. Wu et al. [9] isolates phrases that describe individual entities and attends to them. Phung et al. [10] allows placing objects in certain region by boosting their attention in those regions.

Concurrently with our work, Dall-E 3 [25] proposes to use

an automatic captioning system to regenerate the captions used to train a T2I model. Our work uses an open model (Stable Diffusion) and we provide more details and focus more on analysis and evaluation, but otherwise the main ideas are very similar.

4. Method

Our method, RECAP consists of 3 steps: (1) fine-tune an automatic captioning system to produce desirable labels; (2) relabel the images from our text-to-image training set with this automatic recaptioning system; and (3) train the text-to-image model on the dataset consisting of the images and new captions. Fig. 2 visualizes the overall method.

4.1. Training Dataset

We selected a subset of 10M photos from *LAION-2B-en improved Aesthetics* dataset. We followed the data filters used to train Stable Diffusion versions 1.2-1.4², as follows: *aesthetics score* ≥ 5.0 , *pwatermark* < 0.5 , *nsfw* is *UNLIKELY* and both *height* & *width* ≥ 512 . We note that this filtering operation might amplify biases in the dataset [26]. We further excluded an arbitrary subset of 10K photos from the training set to be used as an internal validation set.

4.2. Captioning Model

We used a pre-trained I2T captioning model (PaLI [12]). As the model outputs are relatively terse and lack in detail, we first collected a small set of 100 manual captions from human raters and fine-tuned the captioning model on that set.

In order to experiment with the effect of different captioning distributions, the raters were asked to provide two types

²<https://huggingface.co/CompVis/stable-diffusion>

of captions. First, a detailed caption for each image, with these instructions: “Describe what you see in each image using 1-2 detailed sentences”. Note that the instructions limited the length to only 1-2 detailed sentences at most, due to CLIP’s (the downstream text encoder) context size of only 77 tokens, which longer captions exceeded. Second, we collected a short and less detailed caption for each image with this instruction: “Describe what you see in each image using a single short sentence”. We did not iterate on the quality of these manual captions. Fig. 3 provides example captions given by the human raters, as well as those produced by the non fine-tuned PaLI model.

With this small dataset, we fine-tuned PaLI for 300 steps, using a learning rate of $4e-5$, dropout rate of 0.1 and a batch size of 64, mixing 50% short captions and 50% long captions for multiple copies of the 100 images, using a different fixed conditioning prefix for the short vs. long captions. When generating captions from the fine-tuned model, we use the terms *RECAP Short* and *RECAP Long* to refer to generations conditioned on the short and long prefixes respectively. Example outputs, as well as a comparison to the original Alttext captions can be found in Fig. 4. See Appendix C for additional dozen random examples of the generated RECAP captions.

Overall, the captions produced by this bespoke model improve both of the issues above - their distribution better matches inference time prompts and they contain much more detail to improve sample efficiency (see Sec. 6.1.1).

4.3. Image Generation Model

Next, we fine-tuned Stable Diffusion v1.4 for an additional $250k^3$ steps with a learning rate of $1e-5$, batch size of 512, and prompt dropout rate of 0.1^4 . We fine-tuned the model training both UNet and CLIP weights and used a 50%-50% mixture of RECAP Short and RECAP Long captions (RECAP Mix) as this performed best. The results in the main text are all from this configuration.

RECAP is independent of the sampling method, so we can use any sampling method with it. That said, for all of the experiments in this work, we used DDIM sampling with 50 inference steps and a guidance scale of 7.5.

5. Results

We compare the RECAP model to two models: Baseline and Alttext. Baseline is the Stable Diffusion v1.4 model. Alttext is the baseline model fine-tuned for the same number of steps and on the same set of images as RECAP, but with the

³When further fine-tuning the model for 1M steps, we observe better visual results along with diminishing returns in the auto computed metrics. All results in the paper are for 250K steps except Fig. 1 and the human evaluation, where we used a model fine-tuned for 1M steps.

⁴Prompts with more than 77 tokens (CLIP’s upper limit) were dropped. This happened for <1% of the data in each training set.

original captions (Alttext) instead of the RECAP captions. The Alttext model resolves contamination concerns, as it includes the exact set of images. In all the comparisons we used the same random seeds across models.

We compare the models using a variety of automated metrics, human evaluation and a qualitative evaluation of examples. We observe improvements in all metrics (see Tabs. 1 and 2).

5.1. Automated metrics

We evaluated the performance and semantic capabilities of the RECAP model using a battery of metrics suggested by [27] (using their publicly available code) on the MS-COCO validation dataset. Tab. 1 contains a summary of the results.

To assess overall generation quality we use the standard FID metric and observe that images generated with the RECAP model have a significantly better score ($17.87 \rightarrow 14.84$).

In addition, we assess the semantic capabilities of our model: to check that the model generates faithfully the requested objects we use Semantic Object Accuracy ($80.80 \rightarrow 86.17$), to check the number of generated objects we use Counting Alignment errors ($1.44 \rightarrow 1.32$), to check that the locations of objects are correct we use Positional Alignment ($57.60 \rightarrow 62.42$), and finally to check the overall adherence to the prompt we measure Clip score ($92.78 \rightarrow 93.80$).

In all the metrics, we see no improvement in the Alttext model compared to the baseline, proving the improvements stem from the captions themselves and not from the additional training.

Note that following [27], we omit the IS and O-IS metrics, as all Stable Diffusion models yield a higher score than real images. As noted by [27], IS is better suited for datasets with a single object and does not perform well for the MS-COCO dataset, containing multiple objects per each image.

Additional details on the automated metrics are in Appendices A.1 to A.4.

5.2. Human evaluation

For complementary evaluation of model performance, we used human raters. Results are summarized in Tab. 2.

Raters were asked to select images generated from each model, only if they successfully follow a given prompt. We evaluated once on 200 random prompts from the MS-COCO validation set, and separately on the challenging DrawBench dataset [1]. We presented four images (from different seeds) for each prompt, using the same seeds across models.

We calculated two metrics: percentage of successful image generation across all prompts and seeds (i.e. given a prompt and a seed, the chance of a successful image generation); and percentage of at least one successful image generation for a given prompt, out of four seeds (i.e. given a prompt, the chance of successfully generating an image



Short Caption A blue willys gasser car

Long Caption A willys gasser car in blue color. it is placed near the car shed on the floor. in the back, there is a man seated in the chair

PaLI a blue coupe with a big engine in it .

Short Caption A glass of iced tea

Long Caption A glass of iced tea placed on a saucer decorated with mint leaves. it is located on the wood table

PaLI a glass of iced tea with a straw and mint .

Figure 3. Examples of captions given by human raters, and the automatically generated caption from the non fine-tuned PaLI model. Photos taken from LAION.

for it, in four attempts). We see a relative 64.3% improvement in successful image generation on MS-COCO, and a 41.7% improvement on DrawBench. We also see a relative improvement of 42.1% in successful prompt generation on MS-COCO and 37.5% improvement on DrawBench. The Alttext model showed minor improvement on the MS-COCO dataset (12%-13%) and did not improve the DrawBench dataset. Further details can be found in Appendix A.5.

5.3. Qualitative Results

Fig. 1 provides representative examples where RECAP outperforms the base Stable Diffusion model (and sometimes also larger models, see Appendix B).

Generally, we observe that RECAP can better interpret relations between entities. Prepositions like "near", "through", "made off" are often ignored by the base Stable Diffusion model. The RECAP model applies the prepositions correctly, while the base model often resorts to the most common relation between the entities (based on the training data distribution). For example, in the prompt "A pizza near a pineapple", the base model places the pineapple on the pizza, as the probable relation between the two, while RECAP generates a pineapple near it, as requested.

RECAP also better handles cases where different modi-

fiers are applied to multiple entities (e.g. "A red bench and a yellow clock"). The base model will treat the sentence as a bag of words, applying all modifiers to all entities or ignore some of them. RECAP is also able to interpret complex modifiers like anaphors. For examples, in the prompt "Two flowers, one is blue and the other one is green" it understands that "one" refers to a flower.

6. Analysis

We hypothesize that the underlying improvement in image generation quality, as measured in the results above, stems from two improvements in the training captions: (1) reducing the discrepancy between the train and inference prompts, and (2) giving the model more textual information per image, thus improving the training sample efficiency. Below we provide an ablation analysis showing that RECAP captions achieve both properties, and that the resulting model benefits from both.

6.1. Comparing Different Caption Types

6.1.1 Generated Captions

Tab. 3 compares the generated captions in the training set, to the caption in the MS-COCO validation set, on a variety



Alttext 2013 ducati monster 1100 evo diesel motorcycle photos and specifications. Black Bedroom Furniture Sets. Home Design Ideas

PaLI a motorcycle with a black background

RECAP Short A green benelli motorcycle

RECAP Long A modern motorcycle with a combination of black, grey, and green colors. it is placed on a black background

Cover of Oregon Wine Press February 2019

a table with a dinner plate and a glass of wine .

A poster of food served on plates

A black plate filled with a variety of indian food items along with a glass of white wine. the caption reads "spice it up"



Alttext You Want To Learn About What Human Anatomy The Skeletal System

PaLI a girl drawing a skeleton made out of cardboard .

RECAP Short A girl making a human skeleton model

RECAP Long A girl who is kneeling down and drawing a human skeleton model. the skeleton model is made of cardboard

GraphicRiver Toucan Pattern 7750027

tropical birds and flowers seamless pattern .

Floral and toucans seamless pattern

A seamless vector pattern with toucan birds and tropical flowers. the background color is beige

Figure 4. Examples of captions generated by the RECAP model conditioned on the short or long prefixes, the original PaLI model, and the original Alttext captions. Photos taken from LAION.

	FID↓	O-FID↓	SOA-C↑	SOA-I↑	CA↓	PA↑	RP↑
Baseline	17.87	8.19	78.90	80.80	1.44	57.60	92.78
Alttext	17.53	8.90	78.99	80.85	1.47	57.40	91.32
RECAP	14.84	6.23	84.34	86.17	1.32	62.42	93.80
Real Images	2.62	0.00	90.02	91.19	1.05	100.0	83.54

Table 1. Results for the automated metrics for RECAP model vs. baseline and Alttext models. RECAP model improves across all metrics.

	MS-COCO Successful Images	MS-COCO Successful Prompts	DrawBench Successful Images	DrawBench Successful Prompts
Baseline	29.3%	53.5%	15.6%	33.1%
Alttext	33.4%	60.0%	13.2%	33.8%
RECAP	48.1%	76.0%	22.1%	45.5%

Table 2. Human evaluation results comparing RECAP, Alttext and Baseline models, on two benchmarks (MS-COCO and DrawBench) for two metrics: percentage of images generated that fully follow the prompt, and percentage of prompts with at least one generated image (out of four seeds) that fully followed the prompt.

of language metrics. We used MS-COCO to represent the prompts we expect to see in inference time. We observe that our generated captions are closer in distribution to MS-COCO in several senses.

First, we compare them by using standard language metrics from the *textstat* python package⁵ calculated over the 10M examples in each train dataset. *Flesch Reading Ease* score is a standard measure for how difficult it is to read a text in a given language. We observe that RECAP is able to generate easy to read sentences, while the original Alttext is often difficult to read. Similarly, the *text_standard* score is a consensus score based on a battery of metrics to estimate the smallest grade of a native speaking reader to be able to read the text. RECAP generates texts which a 4th grader (the lowest possible) can read, while the original Alttext is estimated to require 8th graders to fully understand.

To measure more directly the reduction in train-inference skew, we compare the distribution of the embedded texts. Since Stable Diffusion is using CLIP to encode the text, we calculate the Fréchet distance between the CLIP embeddings of the various datasets⁶. Results are summarized in Tab. 3. As suspected, RECAP generated captions are closer in distribution to the MS-COCO captions than the Alttext captions. Furthermore, the RECAP Short captions are closer to MS-COCO than RECAP Long captions. These results are in line with the improvement in FID of the images from models trained with the corresponding datasets, as detailed in Sec. 6.1.2.

Next, we measured how well the automatically generated captions describe the images, using a human evaluation of 100 random images, asking raters to score each caption con-

sidering both faithfulness to the image, and completeness in describing the image, on a scale of 1-5. RECAP generated captions were rated as more faithful and complete (with an average score of 3.58 for RECAP Short and 4.3 for RECAP Long) than Alttext captions (scoring on average 2.9). This supports our hypothesis that automated captions make training more efficient by providing more textual information.

6.1.2 Generated Images

Next, we compared the results of fine-tuning Stable Diffusion on our two caption types, RECAP Long and RECAP Short. The results are summarized in Tab. 4. We observe that training on RECAP Short captions achieves better FID scores, and faster, but with little semantic improvement, while the RECAP Long captions exhibit significant semantic improvement (see representative metrics in Fig. 5). Mixing the caption sets (RECAP Mix) provides the best of both worlds.

Note that the improvement in FID on the MS-COCO validation set, is correlated to the improvement in the Fréchet distance of the generated captions, as detailed in Sec. 6.1.1. This indicates that the semantic improvement in the RECAP Long model stems from improved training sample efficiency, and not only from reducing the train-inference skew.

6.2. Training Different Model Weights

To further explore the contribution of the improved examples to each part of the Stable Diffusion model, we compared training only the UNet weights, vs. the CLIP weights, vs. both, using the same training and evaluation procedures. For simplicity, we only report the results for RECAP Mix vs. Alttext and baseline models.

The results are summarized in Tab. 5. Overall, as expected, training more weights achieves better performance.

⁵<https://pypi.org/project/textstat>

⁶Similarly to the FID metric used to compare between image distributions, we take the Fréchet distance of the CLIP embedding distributions, modeled as a Gaussian.

	Alttext	RECAP Long	RECAP Short	MS-COCO
# sentences	1.15	2.15	1.00	1.00
# words	11.38	21.29	5.89	10.45
# letters	61.08	87.96	24.39	41.67
Flesch Reading Ease score	45.71 (= Difficult)	88.35 (= Easy)	86.61 (= Easy)	86.93 (= Easy)
<i>text_standard</i> score	8.78 (8th grader)	4.50 (4th grader)	3.11 (4th grader)	4.59 (4th grader)
Distance to MS-COCO	0.45	0.24	0.18	0.00

Table 3. Comparison of average statistics on the captions generated by RECAP vs. the original Alttext over the 10m examples in the training set, and the 10k examples in the MS-COCO validation set. Distance to MS-COCO is the Fréchet distance of the CLIP embeddings of the RECAP caption sets (calculated on 250k samples) vs. the MS-COCO validation set. Note that since the embeddings are normalized, the distance is between 0-1.

	FID↓	O-FID↓	SOA-C↑	SOA-I↑	CA↓	PA↑	RP↑
Baseline	17.87	8.19	78.90	80.80	1.44	57.60	92.78
Alttext	17.53	8.90	78.99	80.85	1.47	57.40	91.32
RECAP Short	14.85	5.81	79.81	81.95	1.43	60.00	93.30
RECAP Long	15.61	7.48	84.59	86.16	1.32	63.27	93.26
RECAP Mix	14.84	6.23	84.34	86.17	1.32	62.42	93.80
Real Images	2.62	0.00	90.02	91.19	1.05	100.0	83.54

Table 4. Results for the automated metrics for RECAP models vs. baseline and Alttext models, comparing models that trained on different sets of captions.

Interestingly, training only CLIP weights (which are $\sim 12\%$ of the total weights) achieves better FID, with less training steps, but with little semantic improvement. Training both CLIP and UNet weights results in significantly higher improvement to semantic scores than training only on one of them. See Appendix A.6 for more details.

We believe that training CLIP weights mainly reduces the skew in the distribution of texts between the training set and the evaluation set, while training the UNet weights mostly improves the alignment of the text to the image.

7. Summary and Discussion

In this paper, we show how text-to-image models can be improved across the board by training on synthetically generated captions. We performed an in-depth analysis demonstrating that short descriptions which narrow the train-inference gap are helpful, as are long and detailed descriptions that improve sample efficiency despite being different from the inference set. We further demonstrated that mixing these descriptions in the training set improves all fronts simultaneously.

There are several interesting directions for future work. It would be interesting to check if by tuning the captioning model to produce ample detail in narrow areas, the same recipe can be used to improve semantic capabilities in new domains (for example, could we create models that can accurately generate hair styles, room designs, facial expressions, clothing, etc. based on detailed descriptions?). Similarly, it is possible to use RECAP to train T2I models in domains

that lack textual captions altogether (e.g., a personal photo album or screencaps of TV shows). We conducted initial experiments here, and they show a lot of promise.

We experimented with fine-tuning a model with the RECAP captions, but it would be interesting to compare that to a model that was pre-trained on the RECAP captions from scratch. Relatedly, it would be interesting to further experiment with different mixtures of the three caption types we have (RECAP Short, RECAP Long, and Alttext). Even more generally, we could imagine creating and mixing together several more flavors of recaptioning models. This could allow us to circumvent the token limit (by training on the same image multiple times with a different subset of the caption each time). Regardless of token limit, it could be interesting to explore training on several shorter captions per image instead of a single long one. It would also be interesting to explore the effects of RECAP on larger models trained on larger datasets.

Finally, RECAP shows the importance of high quality datasets, and that it can be improved with synthetic data, we hope that this provides yet another encouragement to apply such techniques even beyond the T2I domain.

8. Acknowledgments

We would like to extend our gratitude to Eyal Molad, Matan Kalman, Jason Baldrige, and the Theta Labs team at Google, for great reviews, suggestions, and support to this paper.

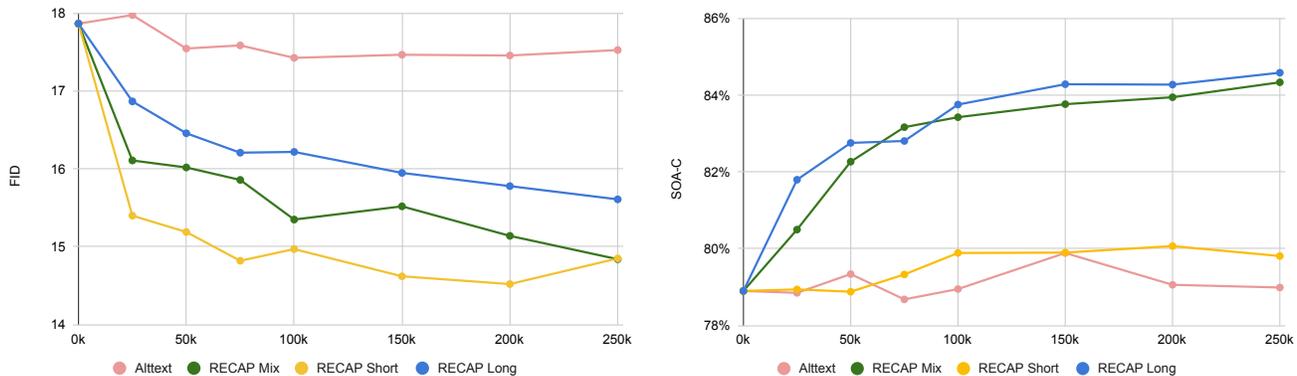


Figure 5. FID (left) and SOA (right) scores for various checkpoints of the Alltext model vs. RECAP models, comparing models that trained on different sets of captions. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. Lower FID is better. We see that RECAP Short achieves better FID and faster, but with no semantic improvement. RECAP Long achieves some FID improvement with significant semantic improvement, and RECAP Mix achieves both.

	FID↓	O-FID↓	SOA-C↑	SOA-I↑	CA↓	PA↑	RP↑
Baseline	17.87	8.19	78.90	80.80	1.44	57.60	92.78
Alltext UNet	16.83	8.20	76.90	78.72	1.50	56.05	91.89
Alltext CLIP	17.04	8.32	79.08	81.10	1.52	58.72	91.34
Alltext UNet+CLIP	17.53	8.90	78.99	80.85	1.47	57.40	91.32
RECAP Mix UNet	15.49	6.47	82.15	84.04	1.40	59.77	92.76
RECAP Mix CLIP	14.60	6.09	80.19	82.04	1.37	60.44	92.48
RECAP Mix UNet+CLIP	14.84	6.23	84.34	86.17	1.32	62.42	93.80
Real Images	2.62	0.00	90.02	91.19	1.05	100.0	83.54

Table 5. Results for the automated metrics for RECAP models vs. baseline and Alltext models, comparing models that trained different set of weights.

References

- [1] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. [2](#), [3](#), [4](#), [12](#)
- [2] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T. Freeman, Michael Rubinstein, Yuanzhen Li, and Dilip Krishnan. Muse: Text-to-image generation via masked generative transformers, 2023. [2](#)
- [3] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021. [2](#)
- [4] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. [2](#)
- [5] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. [2](#), [3](#)
- [6] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. [2](#)
- [7] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models, 2023. [2](#), [3](#)
- [8] Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map alignment, 2023. [2](#), [3](#)
- [9] Qiucheng Wu, Yujian Liu, Handong Zhao, Trung Bui, Zhe Lin, Yang Zhang, and Shiyu Chang. Harnessing the spatial-temporal attention of diffusion models for high-fidelity text-to-image synthesis, 2023. [2](#), [3](#)
- [10] Quynh Phung, Songwei Ge, and Jia-Bin Huang. Grounded text-to-image synthesis with attention refocusing, 2023. [2](#), [3](#)
- [11] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. [2](#)
- [12] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, Alexander Kolesnikov, Joan Puigcerver, Nan Ding, Keran Rong, Hassan Akbari, Gaurav Mishra, Linting Xue, Ashish Thapliyal, James Bradbury, Weicheng Kuo, Mojtaba Seyedhosseini, Chao Jia, Burcu Karagol Ayan, Carlos Riquelme, Andreas Steiner, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali: A jointly-scaled multilingual language-image model, 2023. [2](#), [3](#)
- [13] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. [2](#)
- [14] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. [2](#)
- [15] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. [3](#)
- [16] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. [3](#)
- [17] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision, 2022. [3](#)
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. [3](#)
- [19] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022.
- [20] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models, 2022.

- [21] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. 3
- [22] Thao Nguyen, Samir Yitzhak Gadre, Gabriel Ilharco, Sewoong Oh, and Ludwig Schmidt. Improving multi-modal datasets with image captioning, 2023. 3
- [23] Wenyan Li, Jonas F. Lotz, Chen Qiu, and Desmond Elliott. Data curation for image captioning with text-to-image generative models, 2023. 3
- [24] Feipeng Ma, Yizhou Zhou, Fengyun Rao, Yueyi Zhang, and Xiaoyan Sun. Text-only image captioning with multi-context data generation, 2023. 3
- [25] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Lia, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2023. 3
- [26] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes, 2021. 3
- [27] Tan M. Dinh, Rang Nguyen, and Binh-Son Hua. Tise: Bag of metrics for text-to-image synthesis evaluation, 2022. 4, 12
- [28] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 12
- [29] Tobias Hinz, Stefan Heinrich, and Stefan Wermter. Semantic object accuracy for generative text-to-image synthesis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3):1552–1565, mar 2022. 12
- [30] T. Xu, P. Zhang, Q. Huang, H. Zhang, Z. Gan, X. Huang, and X. He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1316–1324, Los Alamitos, CA, USA, jun 2018. IEEE Computer Society. 12

A. Detailed Results

The following sub-sections explain in more detail the metrics used to evaluate our model, as well as provides plots of several metrics over the number of training steps.

A.1. Image Realism

FID [28] is a measure of how close two image datasets are in terms of the distribution of the semantic content across a large set of photos in each dataset. It is commonly used for evaluation of text-to-image models, by taking a paired set of texts and images, and comparing the set of original images to a set of newly generated images from the texts. A lower score means the distributions are more alike.

O-FID [27] is a variant of FID where off-the-shelf object detector crops out objects from the images first, and FID is calculated on the cropped image sets, giving somewhat more granular measurement of the distribution of objects.

Fig. 6 shows the calculated FID and O-FID scores for several checkpoints throughout training. It is clear that the RECAP model improves significantly the FID and O-FID scores over the MS-COCO dataset, while fine-tuning the same amount of steps on the same set of images, but with the original Alttext captions, produces only a minor improvement.

A.2. Semantic Object Accuracy

The SOA metric was suggested in [29] and is used to evaluate the accuracy of a text-to-image generative model, by measuring how well it follows the instructions to generate specific objects as part of a prompt. To do so, it uses off-the-shelf specialized object detectors of 80 different classes (e.g. a motorcycle or a keyboard) on labeled data. There are two variants to this metric, one that averages across images (SOA-I) and another that averages across classes (SOA-C).

Results can be found in Fig. 7. We see a very significant improvement for the RECAP model vs. the baseline, while the Alttext model does not show any improvement.

A.3. Counting and Positional Alignments

Generative text to image models are known to struggle counting objects, i.e. if a prompt specifies a specific number of objects (e.g. "3 birds") the model often generates a different number of objects. In the Counting Alignment (CA) metric [27], MS-COCO was filtered to prompts which contain a specific instruction to generate a known number of objects. An off-the-shelf counting model was used to count the number of expected objects in the generated image (per object type). The lower the score the better.

Similarly, generative text to image models often struggle to follow positional cues, e.g. "a girl in front of a boy" or "a plate of avocado under the table". In the Positional Alignment (PA) metric [27], MS-COCO was filtered to prompts

with specific positional cues. Each image generated by such prompt gets a CLIP score vs. the original prompt, and also vs. each replacement of the positional cue in the original prompt with a different (wrong) one (e.g. "under" instead of "in front of").

Results for both metrics can be found in Fig. 8. Once again, the RECAP model improves the baseline while the Alttext model does not.

A.4. Text Alignment

R-Precision metric (RP), also known as CLIP score, [30] is a popular measure to how close a prompt is to an image generated from it, by using the CLIP embeddings distance between each prompt and the image generated from it. However, CLIP is also used by Stable Diffusion as the text encoder, producing a bias towards images generated by Stable Diffusion model, and in particular yielding higher RP score for the generated images vs. the original real images. Still, we see that relative to the base model (scoring 92.78), the Alttext model achieves lower score (91.31, probably due to the small dataset size which reduces the training data variance), while the RECAP model improves it (93.8, despite the small dataset size).

A.5. Human Evaluation

We sent samples from the base, Alttext and RECAP models for human evaluation. Raters were presented with four images for each model along with the prompt used to generate the image. The instructions were to only select images that strictly followed the prompt and did not contain any major deformities (minor deformities are fairly common with Stable Diffusion 1.4). Presented results are averaged across raters.

We evaluate all of these models on 200 randomly sampled prompts from MS-COCO and the DrawBench dataset [1]⁷, which Stable Diffusion 1.4 is known to have difficulty with. The RECAP model generates 64.3% (29.25% → 48.06%) more valid images, and is 42.1% (53.5% → 76%) more likely to be able to generate at least one valid image (out of four seeds) over the base model. This indicates that RECAP is both better at generating images as well as being capable of following more difficult prompts (Tab. 2).

A.6. Model Weights

We provide additional plots for the various models trained with different subset of unfrozen weights, as described in Sec. 6.2, in Figs. 9 and 10.

⁷2 prompts which were >77 CLIP tokens were dropped

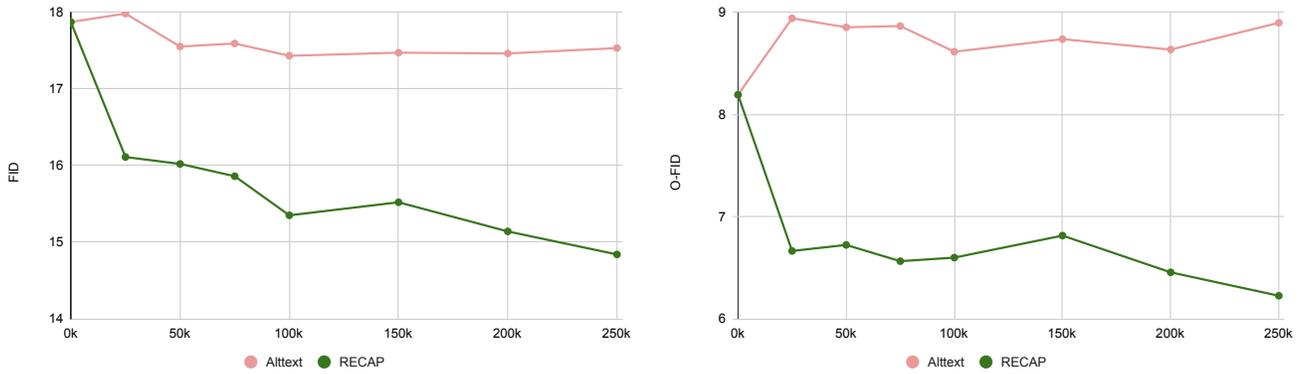


Figure 6. FID (left) and O-FID (right) scores for various checkpoints of the Alltext model vs. the RECAP model. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. Lower score is better.

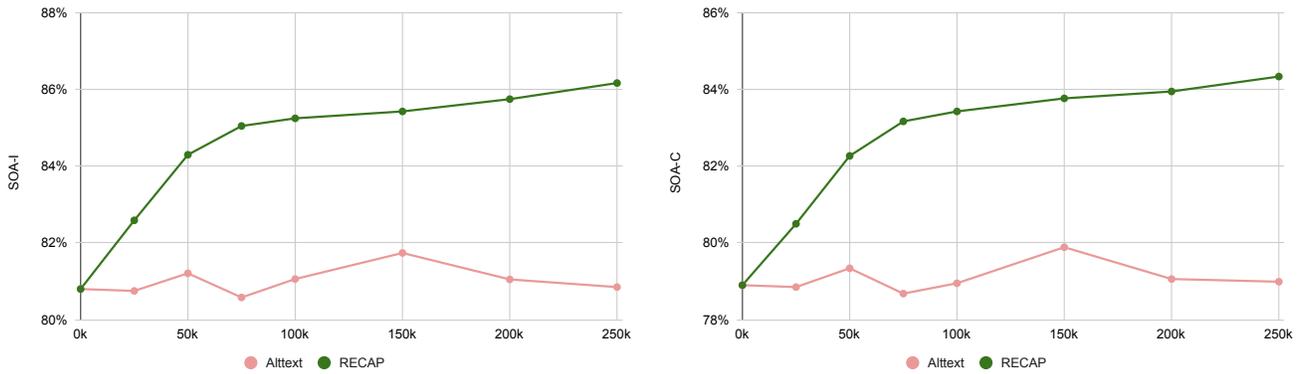


Figure 7. Semantic Object Accuracy (SOA) scores for various checkpoints of the Alltext model vs. the RECAP model. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. SOA-I (left) averages across images, while SOA-C (right) averages across classes.

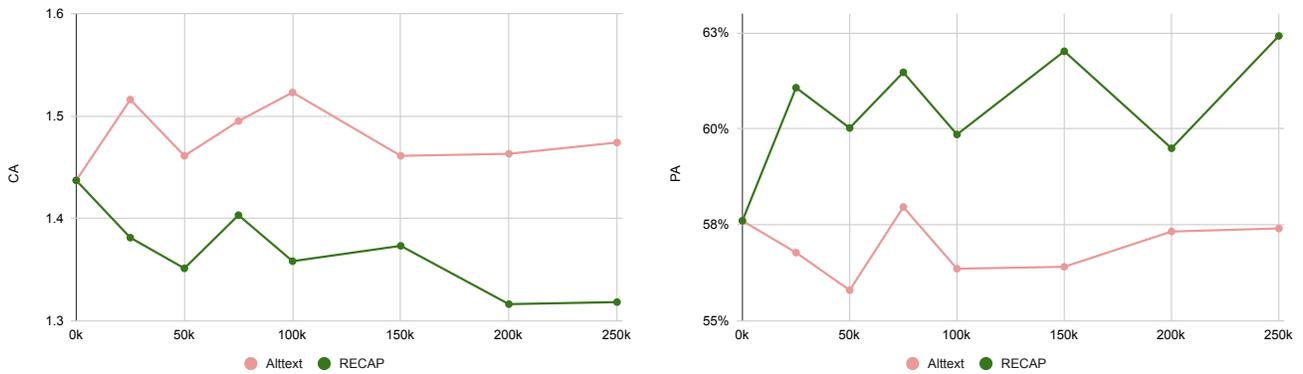


Figure 8. Counting Alignment (CA) (left) and Positional Alignment (PA) (right) scores for various checkpoints of the Alltext model vs. the RECAP model. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. Lower CA score is better.

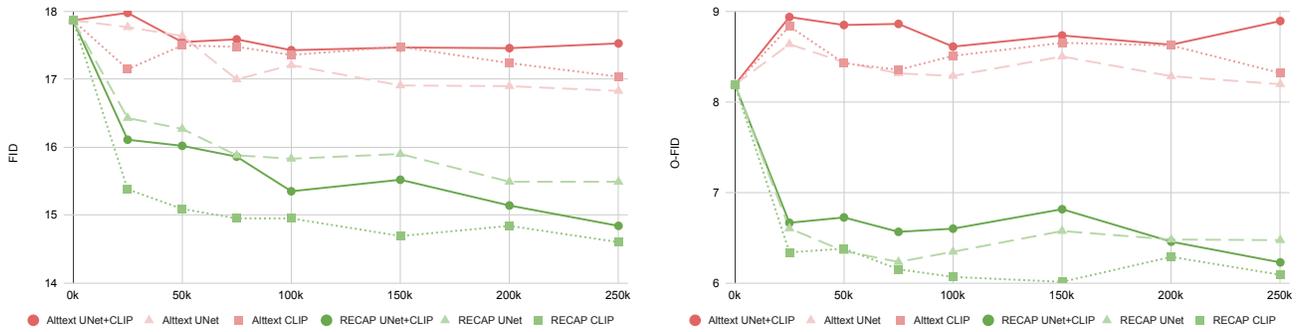


Figure 9. FID (left) and O-FID (right) scores for various checkpoints of the Alltext model vs. RECAP models, comparing models that trained on different set of weights. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. Lower is better.

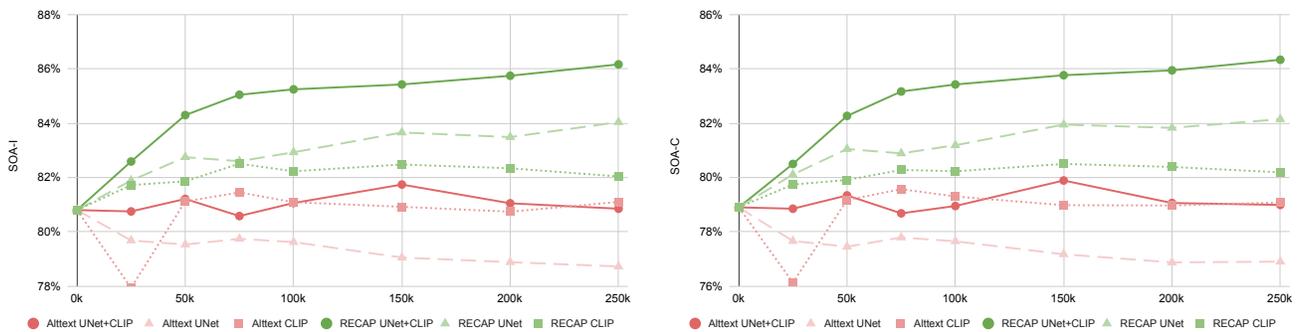


Figure 10. Semantic Object Accuracy (SOA) scores for various checkpoints of the Alltext model vs. RECAP models, comparing models that trained on different set of weights. 0 fine-tuning steps is vanilla Stable Diffusion 1.4. SOA-I (left) averages across images, while SOA-C (right) averages across classes.

B. Other Image Generation Models

Fig. 11 shows example images generated by SDXL and Midjourney for the prompts in the top figure.

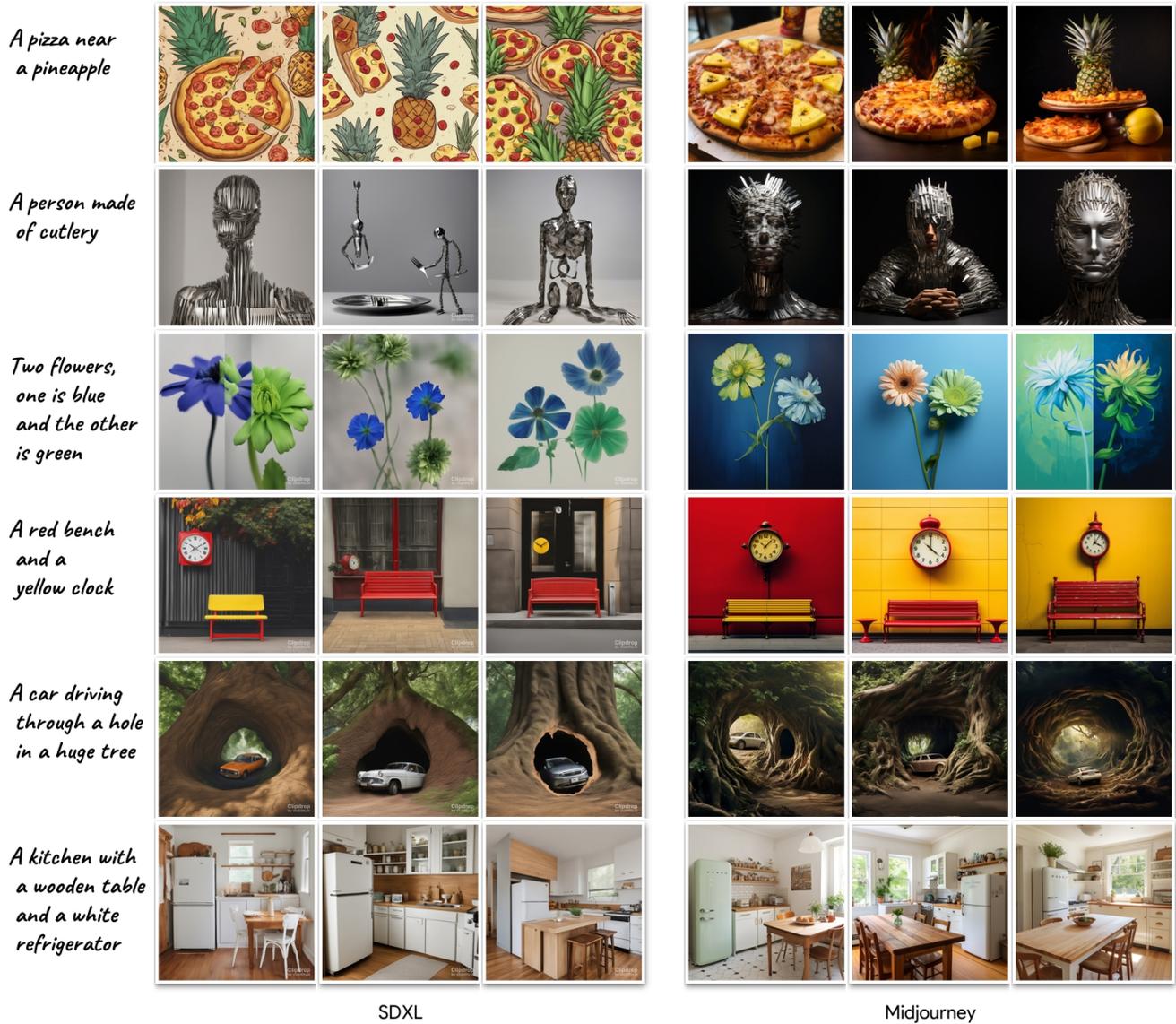


Figure 11. Examples of images generated by SDXL 1.0 and Midjourney 5.1 for the prompts in the top figure.

C. Additional Captioning Examples

Fig. 12 and Fig. 13 show additional example captions generated by RECAP versus the original Alttext.

			
Alttext	iPad Wallpaper - Tribute To Michael Jackson	Simply creative minimalist photography by peechaya burroughs for Minimal artistte	Siberian Rubythroat
PaLI	a silhouette of a man in a suit and hat .	a person holding a balloon in front of a paper egg	a bird standing on one leg on a concrete surface .
RECAP Short	A poster of michael jackson	A hand holding an orange balloon	A bird standing on one leg
RECAP Long	A painting of michael jackson dancing with his arms wide open, the background of the painting is in dark blue color	A hand holding an orange balloon in front of a paper egg, it is made like an egg	A bird with a red beak and white breast, it is standing on one leg
			
Alttext	5 biblioth ques originales faire soi m me - Fabrication etagere garage ...	koala wallpapers pets cute and docile	Black-headed Bunting, male
PaLI	a shelf with a row of shelves and a row of shelf holders .	a koala bear sitting in a tree looking at the camera .	a bird perched on a tree branch against a blue sky
RECAP Short	A shelf with different colour brackets	A koala bear hanging from a tree	A bird perched on a branch
RECAP Long	A row of wooden shelves connected by colorful metal brackets, these shelves are placed against the wall	A koala bear sitting on top of a tree branch, it is looking at the camera	A bird with a black beak and a brown and yellow feather combination perched on a tree branch, it is looking straight ahead
			
Alttext	The Path of the Missing Malaysian Airliner: What We Know, and How ...	Ant Hill with Clover	000 00000 00000 0000 0000000 2013 - 0000 000000 00000 00000 0000 0000000 2013 - Volkswagen Tiguan Photos
PaLI	a map showing the route of the missing passenger jet .	a small ant crawling next to a small plant .	a red car parked on a beach next to a bridge .
RECAP Short	A map of the mh370 route	An ant crawling next to a plant	A poster of the volkswagen tiguan
RECAP Long	A map of the world with countries and a plane in it, the missing malaysia airlines flight mh370 is flying over the bay of bengal, last radar contact with plane	An ant is crawling near a small plant in the dirt, it is surrounded by blurry soil	A red car with a roof rack is parked on the beach, in the background, we can see the golden gate bridge

Figure 12. Examples of captions generated by RECAP variants vs. the original Alttext. Photos taken from LAION.



Alttext GraphicRiver Red Carpet 4026867

PaLI a red carpet leading to a bright light .

RECAP Short A red carpet leading towards a light

RECAP Long A red carpet with golden barriers and red ropes. the end of the carpet is lit by bright light

Alttext valentine's day food art

PaLI a collage of pictures of various foods .

RECAP Short A collage of valentine's day food recipes

RECAP Long A collage of different pictures of valentine's day food. these pictures are all made to look like heart-shaped cakes, chocolates, and sweets

Alttext MALE WHINCHAT

PaLI a small bird perched on top of a rock .

RECAP Short A bird perched on a rock

RECAP Long A bird is perched on a rock with green background. it is looking for something in the rock



Alttext art at Country Music Hall of Fame

PaLI a photo taken at a museum by person

RECAP Short A poster of willie nelson

RECAP Long A poster of willie nelson standing by playing guitar. in the background, we can see he is standing on a skull with snakes around his feet

Alttext 30 cool bathroom ceiling lights and other lighting ideas for Bathroom ceiling ideas

PaLI a bathroom with a sink , shower and a bench .

RECAP Short A modern bathroom with a rain shower

RECAP Long A very modern bathroom with a combination of sink, shower, and a bed. the lighting in the room is very impressive

Alttext Stetson University Spec Martin Stadium

PaLI a stadium with a tower and a field in the background .

RECAP Short A stadium with a tower in the background

RECAP Long A stadium with a bleacher and a tower. the bleacher is empty and ready to use



Alttext Foyer Window Quilt : Gold shoe girl windows for wednesday

PaLI a window with a stained glass window looking out to a snowy yard .

RECAP Short A view of the snow from a window

RECAP Long A view of the snow covered ground outside the window. the car is passing by the white porch

Alttext white coastal kitchen and dining room

PaLI a dining room with a table and chairs .

RECAP Short A white dining room with a view

RECAP Long A dining room with a round table that is made of marble. there are clear plastic chairs placed around the table. there are some lamps placed on the wall

Alttext Dinosaur Train Cake Images : Party Cakes: Dinosaur Train Cake & Cupcake Display

PaLI a cake and cupcakes on a display with ribbon .

RECAP Short A blue cake and cupcakes

RECAP Long A blue cake with green and orange decorations. the cupcakes are arranged on the second layer of the cake stand. the cake is ready to be served

Figure 13. More examples of captions generated by RECAP variants vs. the original Alttext. Photos taken from LAION.