# Prompt Augmentation: Improving Image Generation for Rare Entities

William Sun        Wenqing Wei        Nian Lu

*Abstract*—In this study, we show that text-to-image model underperform on prompts with rare entities, and we solve this problem with prompt augmentation by introducing more details with a finetuned language model. We show that prompt augmentation can improve the performance of rare entities of a frozen text-to-image model, especially for prompts with ambiguous entity names or not enough context available. Code released: https://github.com/LongEarW/imagen_rare_eval

*Index Terms*—Text-to-Image, diffusion model, prompt engineering

## I. Introduction

Recent works show that Text-to-Image models could benefit from the scaling of image-text dataset and LLM size. Image synthesis works like Imagen [1] and unCLIP significantly outperform on COCO dataset [2] with zero-shot testing, compared with previous methods that even trained on COCO.

As we have discussed in depth during class, few-shot techniques during test-time has been shown to dramatically improve performance of LLMs in NLP tasks. Since related work such as Imagen has shown that LLM scaling laws greatly affect downstream tasks such as image generation fidelity, we are motivated to try other techniques from NLP literature, namely few-shot prompting. We believe that these techniques will come into play beyond NLP tasks now that task fusion is occurring in the deep learning space.

In this project, we have the following hypotheses:

1) LLM based Text-to-Image models have lower performance on rare objects
2) Prompt engineering by enriching object text descriptions can improve rare object image generation performance

To show these two hypotheses, we will explore how LLM based image synthesis models perform on rare or unseen objects via qualitative and quantitative analysis on image generations. This is opposed to the generic evaluation dataset, COCO, which only includes common entities (80 categories, eg. Person, Toothbrush) [2]. Hence, how will Imagen perform more generally in the transfer learning space? In general, answering this question helps us to better understand how the LLM may be able to encode information about object text for the subsequent diffusion model, because even though the training set for Imagen may not include examples of a particular object, it may arise in the corpus of the LLM training set. In this project, we also develop a method to enrich the prompt to improve the performance.

Thus, inspired by these techniques, we investigate the effect of augmentation on image synthesis tasks with a frozen generation model in this project.

## II. Related Work

**LLM for Text-to-Image** Text encoder is a critical component for text-to-image models. With rising of LLM (large language models), image synthesis works like GLIDE [3] and DALL-E 2 [4] observed performance boosting while using LLM for latent prior. And Imagen discovered that a frozen generic LLM could be used for text encoding for image synthesis task, and the synthesis performance could benefit from scaling of LLMs. Imagen achieved SOTA with T5-XXL [5], the largest LLM trained on C4 dataset at that time. However, it is reasonable to consider if the currently used LLM is "scaled" enough and if any limitation still exists for future improvement. In this study, we compare the model performance between prompts with and without uncommon entities.

**Performance augmentation.** Re-Imagen suggested model performance can be boosted when trained with prompts which are enriched with retrieved information. Similarly, DALL-E 3 [6] came up with utilizing LLM to improve caption quality with more details for training. Inspired by these works, we propose to use LLM for prompt enrichment, and provide the effect on Imagen.

## III. Implementations

### A. Model

As Imagen didn't release the pretrained model, we use DeepFloyd-IF, a public text-to-image with similar implementation of Imagen. Also, we use the first 2 up-sampling diffusion modules to generate images of 256x256 resolution (the last super resolution module for 1024x1024 is skipped because of limitation on GPU resource). This is able to save 30% time during image generation, while the 256x256 resolution image is high enough quality for both quantitative and qualitative analysis.

### B. Dataset

To compare model performance on rare and common entitites, we need a suitable dataset of caption-image pairs. However, DeepFloyd-IF is trained on LAION-A and internal datasets, which we have no access. Hence, we choose a subset of WebQA's entity image and entity description pairs. With similar filtering rules suggested in Re-Imagen, we remove

noisy information like wiki-id and date, and remove entity descriptions shorter than 4 tokens or longer than 18 tokens. After doing this, we have 389,750 entity image and description pairs for evaluation.

### C. Rare Entity Selection

To collect prompts of rare entities, we define rare entities as noun tokens with frequency of one. Namely, we tokenize text with T5-tokenizor, the same tokenizor used in DeepFloyd-IF and Imagen. This is so that our statistics of token frequency matches with the token space of model, and we can take advantage of its subword tokenization and lemmatization features. In addition, we identify the properties of token with spacy, a popular NLP processing toolkit, which tags word property according to text context for better accuracy.

### D. Prompt Enrichment

To automatically enrich prompts, we follow DALL-E-3's prompt augmentation method. Namely, a language model is biased to generate descriptive captions. While DALL-E-3 didn't release their text-to-image generation system, they publish 200 samples of concise caption and descriptive caption pairs. Hence, we come up the idea of finetuning the ChatGPT to generate highly descriptive captions given concise captions.

We use few-shot to finetune ChatGPT with the template below:

> *"messages": [{"role": "system", "content": "You are good at adding details into image description"}, {"role": "user", "content": $concise caption$}, {"role": "assistant", "content": $descriptive caption$}]*

And the finetuned model is used to augment the entity descriptions. Below are is an example of augmented description:

> Raw DESCR: *Bertinoro, lapide ai caduti della prima guerra mondiale*
> Augmented DESCR: *a stone monument stands proudly in the town square of bertinoro, italy, commemorating the fallen soldiers of the first world war. adorned with wreaths of red poppies and surrounded by a manicured garden, the memorial reflects the solemnity and gratitude of the community.*

One note here to consider is whether the finetuned model will still produce a description that is too long/complex for Imagen, especially when compared to the simpler tests that are shown in Imagen. Another possibility is that ChatGPT may produce false information that will then be introduced into the downstream task. Thus, this process may be refined in the future to handle these issues.

## IV. Evaluation

To evaluate the model performance on common and rare entities, and study if augmented prompts could improve performance, we measure the generated images by image quality, image-prompt alignment, and if images are photorealistic. Following the benchmark used in Imagen, we also make human evaluation for human preference.

### A. Quantative Evaluation

We measure the image quality by FID score (1), and alignment by CLIP score (2). For FID, $\mu_1$ and $\sigma_1$ represent the mean and covariance of feature activation of groundtruth images by InceptionV3, while $\mu_2$ and $\sigma_2$ represent those of generated images. The lower FID score indicates higher image quality. For CLIP score, $v_i$ represents the visual CLIP embedding of $i^{th}$ image, and $c_i$ represents the textual CLIP embedding of corresponding caption. The higher CLIP score indicates better alignment between text and image.

$$FID = |\mu_1 - \mu_2| + Tr(\sigma_1 + \sigma_1 - 2\sqrt{\sigma_1 * \sigma_2}) \quad (1)$$

$$CLIP - S = \frac{\sum_i^N \max(\cos(v_i, c_i), 0)}{N} \quad (2)$$

We provide the experiments' result in Table. I. Comparing model performance between baseline prompts and prompts with rare entities, the model significantly underperforms on rare entities with respect to image quality (higher FID), which meet the expectation. And we don't observe significant difference on text-image alignment. While comparing model performance of augmented prompts to prompts with rare entities, we observe augmented prompts bring better text-image alignment (higher CLIP score). However, the image quality further decay (highest FID among all three experiments). We predict that this is because DeepFloyd is not trained with highly descriptive prompts, hence the generated images diverge from the training image sets.

#### TABLE I
#### PERFORMANCE METRICS

| - | Baseline DESCR | Rare DESCR | Aug DESCR |
|---|---|---|---|
| FID | 95.32 | 125.26 | 132.71 |
| CLIP | 0.2967 | 0.2976 | 0.3188 |

### B. Human Evaluation

For evaluation, we are inspired by the techniques used in Imagen. Specifically, for human evaluation, we try a benchmark (266 instances) similar to DrawBench which is large enough to test the model well, but small enough to perform trials within the team (three human raters) as for the purposes of this project. In Drawbench, raters are ask which set of images is of higher quality and which set of images better represents the text caption. And in our evaluation, the rates are given three questions:

1) Which image has better text-image alignment?
2) Which image is more photorealistic?
3) Which image do you prefer?

Then, to compare the original rare image generation with the enriched generation, we score each pair according to the three prompts with a tuple $(s_1, s_2, s_3)$, where $s_i = 1$ if the original is better, $s_i = 2$ if the enriched is better, or $s_i = 3$ if neither is preferable. We provide the statistics below. Note that the

preference percentages are not 100% in total as some instance is rated as 3 (neither is preferable)

| - | Rare DESCR | Aug DESCR |
|---|---|---|
| Alignment | 29.32% | 57.14% |
| Photorealistic | 26.32% | 62.03% |
| Subjective Preference | 27.44% | 59.40% |

In human evaluation, the enriched generations outperforms the default generation on all metrics (higher preferred rate on alignment, photorealistic and subjective preference), which are partially against with the quantitative results in the previous section. Namely, the quantitative results show augmented prompts lead to lower image quality, while human evaluation suggests augmented prompts have significant advantages on "photorealistic" and "subjective preference". Also, the human evaluation suggests augmented prompts has more significant advantage on alignment.

As to image-to-text alignment, we believe this is because: 1) the enriched prompts has larger length and will be truncated by text encoder of CLIP; 2) CLIP is not trained on highly descriptive captions, some details are omitted. As to subjective preference and photorealistic, we believe the enriched generations contain more details, and give a sense of realistic. In the meanwhile, the enriched generations diverge from the groundtruth images, and leads to higher FID score, which only measures the divergence between two image distributions. Note that the augmented prompts are generated by the fine-tuned ChatGPT, without the supervision from original images.

An interesting note here for Table 2 is that a sizeable portion of the evaluation samples are scored $s_i = 3$. This means that many of the generations, enriched or not, are not preferable, which indicates that the resulting generation is nonsense to the human evaluator. We discuss the reasoning behind this in the next section. We also provide some prompts and generated images that support our assumptions in the next section.

## V. GENERATION ANALYSIS

In the following section, we provide analysis for the results of the human evaluation. Each sample is visualized in the following manner. Each row contains these items in order: image number, original dataset image, original caption, default generation, enriched caption, enriched generation. Next, we provide some examples of how the human evaluation was done.
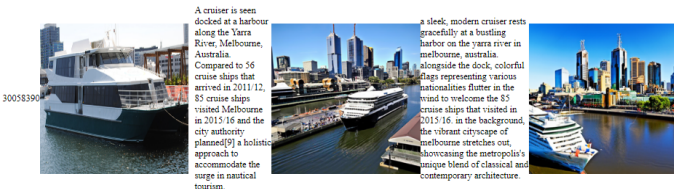


Fig. 1. Human evaluation sample 1.

Figure 1: An example where both generated images are good, but the non-enriched sample is better. It has better alignment because it includes the dock next to the ship, and it is more realistic because the color saturation looks more like a photo. Thus, we give this a rating of (1, 1, 1).
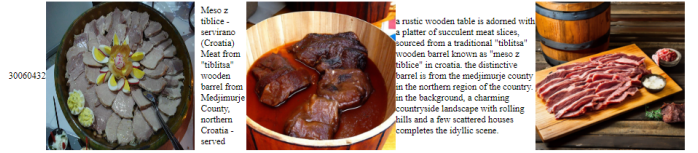


Fig. 2. Human evaluation sample 2.

Figure 2: Another example where both generated images are good, but the enriched sample is better because the sliced meat matches the actual dish, as well as the color palette of the dish itself. Thus, we give this a rating of (2, 2, 2).



Fig. 3. Human evaluation sample 3.

Figure 3: An example where both generated images are bad. The resulting images are both nonsense and unlike the original image. We find that these occur when the prompt itself is lacking. Thus, we give this a rating of (3, 3, 3).

We admit that this schema for human evaluation can be flawed and biased, as all three of us know which example is the enriched generation during evaluation, but we chose to do this because it was an easier process given the limited time frame. We believe that this is okay for the project because for this project we are only interested in the results as a proof-of-concept. Given the overwhelming preference of augmented generations over the default generations, we can safely move this process to the next step in a real study, where we would choose more human evaluators in a blind voting process. Then, in this case, all we would show is the original caption and the default generation versus the enriched generation.

Generally, because we are inputting difficult prompts into the generator, we find that the image generation can sometimes lead to nonsense. This occurs especially for figure-caption pairs from scientific literature, where the figure and caption build off of one another, so the descriptions in the captions may be significantly lacking. Hence, it is difficult for Imagen to perform well with image alignment without proper descriptions. See Figure 4 for an example where Imagen fails, but proper description provided in the enriched prompt helps significantly. However, in some cases, such as captions following images of diagrams, the subsequent enriched prompt is still not helpful. In this case, a more intelligent approach to prompt enrichment is still needed. Part of the fault here is to do with the original

WebQA dataset. Many of these caption-image pairs that were included in the dataset would not make sense to a human evaluator in the first place, as the caption-image pair of a scientific diagram may be nonsensical. Thus, this could lead to poorer alignment later on during test time.



Fig. 4. Enriched prompt performs much better.

## VI. CONCLUSION AND FUTURE WORK

Through our experimentation, we found that augmented prompts can greatly improve image generation for rare entities, especially for prompts with ambiguous names, lacking in context, or for prompts "unknown" to the model from the text-encoder training set. This makes sense as what we know from few-shot learning in NLP transfers to this task well. Providing more context and logic for the text encoder to work with will greatly enhance the performance in the downstream task of image generation. However, we also observe the limitations of our method:

1) Quantitatively worse image quality with augmented prompt
2) Enriched details are sometimes omitted by the image generator
3) No multi-modal guidance for prompt augmentation

We believe that the first two limitations can be solved with further experimentation and clever prompt engineering techniques (informative prompts within training prompts' distribution). This is because we know how well Imagen can perform in photorealistic image alignment even given complex prompts, so the issue is more tied to correctly representing the enriched prompt correctly. Hence, if time permitted, we plan to work on:

1) Shorten/simplify the generated enriching text within certain length
2) Generate a tree of enriched text and optimize
3) Enrich prompts with supervision from the groundtruth images with the caption ability of GPT
4) Finetune the image generation model with enriched prompts
5) Gather a larger, unbiased human evaluation metric

As discussed above, the first four points are a variety of techniques to try for more clever prompt augmentation experiments. Ablation studies involving these would also be interesting to test out. Finally, as discussed in the previous section, point five would be necessary for a true study.

## 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" arXiv:2205.11487, 2022.
[2] Lin, TY. et al., "Microsoft COCO: Common Objects in Context." In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds) Computer Vision – ECCV 2014. ECCV 2014. Lecture Notes in Computer Science, vol 8693. Springer, Cham. https://doi.org/10.1007/978-3-319-10602-1_48
[3] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Bob McGrew Pamela Mishkin, Ilya Sutskever, and Mark Chen. "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models." In arXiv:2112.10741, 2021.
[4] Aditya Ramesh and Prafulla Dhariwal and Alex Nichol and Casey Chu and Mark Chen "Hierarchical Text-Conditional Image Generation with CLIP Latents" In arXiv:2204.06125, 2022.
[5] 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." JMLR, 21(140), 2020.
[6] James Betker, Gabriel Goh, Li Jing, TimBrooks, Jianfeng Wang, Linjie Li, LongOuyang, JuntangZhuang, JoyceLee, YufeiGuo, WesamManassra, PrafullaDhariwal, CaseyChu, YunxinJiao and Aditya Ramesh "Improving Image Generation with Better Captions" JMLR, 21(140), 2020.https://api.semanticscholar.org/CorpusID:264403242

# Augmented Better Instances

30250205 — Aloe boiteani

a vibrant green aloe boiteani plant sits on a rustic wooden windowsill, basking in the warm sunlight. its thick, fleshy leaves are patterned with light green spots and lined with small, soft spikes. behind the window, a cozy room with white walls and potted plants exudes a calm and inviting ambiance.

30262289

In the Zero series, Zero was redesigned to look both sleeker and more human.

in the zero series, the iconic character zero has been given a stunning redesign. he appears sleeker and more humanoid, with a streamlined body and sharp, angular features. clad in a new, cutting-edge black and red armor, zero gazes intensely ahead, brimming with determination and power. the futuristic setting around him enhances the sense of technological advancement, creating an awe-inspiring image.

30267877

A chiropractor performs an adjustment on a patient.

a chiropractor in a white lab coat gently places their hands on a relaxed patient, who is lying face down on a padded table. the patient's expression reflects a sense of relief as the chiropractor applies a precise adjustment to their spine. sunlight streams in through a large window, casting a warm glow in the serene chiropractic office.

30298993

Intellectual activities such as playing chess or regular social interaction have been linked to a reduced risk of Alzheimer's disease in epidemiological studies, although no causal relationship has been found.

in a cozy living room filled with warm sunlight, an elderly couple engages in a spirited game of chess. their gaze focused on the intricately carved pieces, their hands move with precision and determination. a bookshelf lined with books on various intellectual topics, including psychology and neurology, serves as a backdrop. the room exudes a sense of intellectual curiosity and a love for engaging in stimulating activities.

30335172

Viola adunca 6719 Viola Aduncahookedspur Violet, Early Blue Violet; auth. Sm.

a vibrant cluster of viola aduncahookedspur violets, or early blue violets, rests upon a mossy rock in a peaceful forest clearing. sunlight filters through the canopy above, casting a gentle glow on the flowers and their surroundings. the scene is filled with the delicate beauty of the wildflowers, their petals ranging from purple to blue, while nearby ferns and shrubs create a lush backdrop.
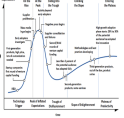
30033254

Claude Monet - Meules, milieu du jour

in monet's famous painting "meules, milieu du jour," the scene is set in the middle of a sunlit day in the french countryside. a series of haystacks is scattered across the field, casting long shadows on the gently undulating ground. vibrant hues of yellow, orange, and green paint the landscape, while an old farmhouse stands in the distance, creating a peaceful and idyllic rural scene.

# Raw Better Instances

30299645 — Traditional leben

a traditional leben, or arabic tent, sits in the golden sands of the desert at sunset. adorned with colorful geometric patterns, the tent's open flaps reveal soft cushions and ornate carpets inside. a gentle breeze stirs the surrounding dunes as a camel peacefully grazes nearby.

30308759

A sleeping Arctic fox with its fluffy tail wrapped around itself and over its face

a serene arctic fox lies in a snowy field, its fluffy white tail wrapped protectively around its body and tucked over its face. the fox's closed eyes and relaxed expression reveal a sense of tranquility, while the surrounding landscape glistens under the soft light of the arctic sun.

30342699

A visual representation of capital gain with coins, as the essential nature of capital gain is accrual of capital.

a sprawling field of vibrant green grass, with a single large, ancient-looking tree in the center, provides the backdrop for the concept of capital gain. stacks of shiny gold coins rest at the base of the tree, while smaller, sprouting saplings representing future growth surround them. the bright sun above casts a warm glow, emphasizing the wealth that can accumulate over time through capital gain.

30389072

Natural French manicure Fing'rs Natural French Spitzen

a hand with perfectly shaped nails adorned with a natural french manicure featuring the fing'rs natural french spitzen. the translucent white tips of the nails blend seamlessly with the soft pink base, giving an elegant and sophisticated look to the hand. sunshine streams in from the side, highlighting the manicure's natural beauty.

30386729

General hype cycle for technology

a grand, bustling market square buzzing with the energy of innovation and progress. towering holographic advertisements display the latest technological marvels, while street performers showcase futuristic gadgets. people from all walks of life excitedly discuss and test the newest inventions, creating a vibrant atmosphere of anticipation and curiosity.
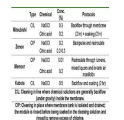
30011208

Under the groundsheet - geograph.org.uk - 1527254 Under the groundsheet I was amazed to find this little Common Newt (Lissotriton vulgaris) under the groundsheet of my awning when I was packing things away.

a serene garden scene is captured in the photo, with a large green awning shading a cozy nook. underneath the groundsheet, colorful flowers bloom, and a surprised individual gently lifts the corner of the fabric to reveal a tiny common newt. the sun shines warmly, casting a golden light on the discovery, as nearby plants and trees create a peaceful atmosphere.

# Both not preferred Instances

30214656

A ritmo de marimbas en el baile de negras, Masaya Nicaragua tomada por Maynor Valenzuela

in the vibrant streets of masaya, nicaragua, a group of locals dance to the rhythm of marimbas under the warm afternoon sun. women in traditional black dresses, adorned with colorful embroidered details, gracefully move their feet while their partners twirl them around. the scene is captured by maynor valenzuela, a local photographer, who emphasizes the energy and joy of the moment in his photograph.

30314325 — Dentarene sarcina.

a heavily pregnant dentarene, a mythical creature with the body of a lion and the head of an eagle, rests peacefully in a sunlit forest clearing. its wings are folded protectively around its round belly, and colorful wildflowers bloom at its feet. a sense of anticipation fills the air as other forest animals gather to offer support and celebrate the imminent arrival of the dentarene's young.

30049024

Intensive chemical cleaning protocols for four MBR suppliers (the exact protocol for chemical cleaning can vary from a plant to another)

four MBR suppliers engage in an intensive chemical cleaning process, each utilizing a different protocol specific to their plant. the industrial facility setting is bustling with activity as workers in protective gear meticulously follow their respective procedures. large stainless steel tanks, pipelines, and a control panel dominate the foreground, showcasing the intricate nature of the operation. in the background, brightly lit research and development labs hint at ongoing innovation in the field.

30079120

Homeopathic remedies – ineffective for treating cancer

a serene and brightly lit homeopathic clinic, with shelves lined with small bottles of water and diluted substances. in the center, a compassionate practitioner explains to a patient that homeopathic remedies are ineffective for treating cancer, while offering alternative complementary therapies to support their well-being.

30087099

Each of the Spanish parties had its recommendation to voters.

a vibrant plaza in a spanish village serves as the backdrop for a lively political rally. supporters of different parties gather, with each wearing the distinctive colors and symbols of their chosen group. flags and banners flutter in the wind, while musicians play traditional tunes, adding to the energetic atmosphere.

30140634

Origin[x]. Giant Papillon [UK]See also:Checkered Giant [US]Miniature Papillon11–25 lb(5.0–11.3 kg)ShortErect[Colored butterfly, eye circles, cheek spots, ear base, saddle, and rump spots; all on a base of white.] "All [BRC] recognised colours are admissible."NoYesEUGermany

a giant papillon, originating from the uk, is depicted in a lush meadow. its coat, consisting of a base of white, is adorned with intricate patterns of colored butterflies, eye circles, cheek spots, ear base spots, and saddle and rump spots. the majestic butterfly dog stands tall with its short, erect ears accentuating its alert expression. a few other miniature papillons playfully surround the giant papillon, basking in the picturesque english countryside.