# A data augmentation perspective on diffusion models and retrieval

Max F. Burg[*,1,2], Florian Wenzel[1], Dominik Zietlow[1], Max Horn[1], Osama Makansi[3], Francesco Locatello[1], and Chris Russell[1]

[1]Amazon Web Services, Tübingen, Germany, [2]International Max Planck Research School for Intelligent Systems, University of Tübingen, University of Göttingen, Germany, [3]Amazon, Tübingen, Germany, max.burg@bethgelab.org, {flwenzel, zietld, hornmax, omakans, locatelf, cmruss}@amazon.de

April 21, 2023

## Abstract

Diffusion models excel at generating photorealistic images from text-queries. Naturally, many approaches have been proposed to use these generative abilities to augment training datasets for downstream tasks, such as classification. However, diffusion models are themselves trained on large noisily supervised, but nonetheless, annotated datasets. It is an open question whether the generalization capabilities of diffusion models beyond using the additional data of the pre-training process for augmentation lead to improved downstream performance. We perform a systematic evaluation of existing methods to generate images from diffusion models and study new extensions to assess their benefit for data augmentation. While we find that personalizing diffusion models towards the target data outperforms simpler prompting strategies, we also show that using the training data of the diffusion model alone, via a simple nearest neighbor retrieval procedure, leads to even stronger downstream performance. Overall, our study probes the limitations of diffusion models for data augmentation but also highlights its potential in generating new training data to improve performance on simple downstream vision tasks.

## 1 Introduction

Data augmentation is a key component of training robust and high-performing computer vision models and given its success, it is becoming increasingly sophisticated: From the early simple image transformations (random cropping, flipping, color jittering, and shearing) [6], over augmenting additional training data by combining pairs of images, such as MixUp [44] and CutMix [43], all the way to image augmentations using generative models. Augmentation via image transformations improves robustness towards distortions that resemble the transformation [39] and interpolating augmentations are particularly helpful in situations where diverse training data is scarce [12]. With the success of generative adversarial networks (GANs), generative models finally scaled to high-dimensional domains and allowed the generation of photorealistic images. The idea of using them for data augmentation purposes has been prevalent since their early successes and forms the basis of a new set of data augmentation strategies [47, 13, 2, 8, 7, 23, 20, 26].

Given the emergence of diffusion models (DMs) that outperform GANs in terms of visual quality and diversity [27, 33, 28], using them for data augmentation is a natural next step. Unlike GANs, diffusion models can be easily conditioned using text queries, which allows for more controlled data generation. These models are trained on a large dataset of billions of filtered image text pairs retrieved from the internet, enabling them to generate images of unparalleled variety.

We benchmark a variety of existing augmentation techniques based on diffusion models and propose new extensions and variants. We show that simple text prompts based on class labels suffice for conditioning the DMs to improve the performance of standard classifiers. As images generated by simple prompts match the training distribution of the DM and not the training distribution of the classifiers, we test—inspired by related work on personalized DMs—methods that fine-tune the DM conditioning and optionally the DMs denoising model component. These fine-tuned models outperformed the best prompting strategies for their ability to create even better synthetic data for augmentation (Figure 1).

Surprisingly, we find that simply retrieving nearest neighbors from the DM's training dataset (using the same CLIP-like model [25] used in the DM) with simple prompts is a very strong augmentation strategy. This strategy outperforms all investigated methods based on generating synthetic images (Figure 1), suggesting that the true potential

---

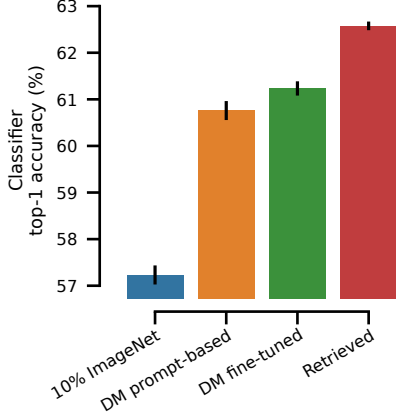[*]Work done during an internship at Amazon Web Services.

Figure 1: *Are generative methods beneficial for data augmentation?* Each bar shows the accuracy (along with the standard error) of the downstream classifier based on the best augmentation methods within each family. While diffusion model based techniques (orange and green) improve over the baseline of only using the original 10% ImageNet data (blue), a simple retrieval method using images from the diffusion model's pre-training dataset directly (red) performs best. This suggests that the generative capabilities of diffusion models for augmentation have not yet been fully leveraged.

of DMs for data augmentation is not yet fully realized, and that there is a large opportunity to improve upon conditioning mechanisms and fine-tuning strategies for DMs.

## 2   Related work

**Data augmentation using generative adversarial networks.**  Data augmentation is widely used when training deep networks. It overcomes some challenges associated with training on small datasets and improving the generalization of the trained models [44]. Manually designed augmentation methods have limited flexibility, and the idea of using ML-generated data for training has attracted attention. Generative adversarial networks (GAN) such as Big-GAN [4] have been used to synthesize images for ImageNet classes [3, 17]. Despite early promising results, the use of GANs to generate synthetic training data has shown limited advantages over traditional data augmentation methods [47]. Diffusion models, on the other hand, might be a better candidate since they are more flexible via general-purpose text-conditioning and exhibit a larger diversity and better image quality [36, 27, 33, 28]. Hence, in this work, we focus on diffusion models.

**Data augmentation using diffusion models.**  Recently, diffusion models showed astonishing results for synthesizing images [36, 27, 33, 28]. Numerous approaches have been published that adapt diffusion models to better fit new images and can be used for augmentation. We evaluate methods that employ diffusion models in a guidance-free manner or prompt them without adapting the model. [19] generate variations of a given dataset by first adding noise to the images and then denoising them again, and [34] generate a synthetic ImageNet clone only using the class names of the target dataset. We also evaluate methods for specializing (AKA personalizing) diffusion models into our study. Given a few images of the same object (or concept), [9] learn a joint word embedding (pseudo word) that reflects the subject and can be used to synthesize new variations of it (e.g., in different styles) and [10] recently extended their method to significantly reduce the number of required training steps. [16] follow a similar approach and proposes a method for text-conditioned image editing by fine-tuning the diffusion model and learning a new word embedding that aligns with the input image and the target text. We investigate the usefulness of "personalization" for data augmentation, which was not conclusively addressed in the original papers. Additionally, we propose and evaluate extensions of these methods tailored to data augmentation and provide a thorough evaluation in a unified setting.

Given the fast moving field, there have been multiple concurrently proposed methods for augmentation that we could not include in our study. Some methods employ alternative losses for image generation [29], optimize the features of the embedded images [45] and use alternative prompts to the diffusion model with class descriptions generated by a language model [15]. Other methods focus on fine-tuning the diffusion model learning a unique identifier for the given subject [30], [38] additionally utilize image synthesis, editing [21] and in painting [18, 32], and others focusing on medical data [11, 1].

## 3   Experimental protocol

We augment using a wide-range of generative and retrieval based techniques, evaluating performance on a downstream classification task (Figure 2).

**Dataset.**  To simulate training data in a low-data regime, we sample 10% of the ImageNet Large Scale Visual Recognition Challenge 2012 (ILSVRC2012) [31] training split retaining class imbalance. As the retrieval method did not return sufficient samples for 10 of the 1,000 ImageNet classes, they were excluded from augmentation, model training and evaluation. We additionally sample a disjoint set of images of the same size from the training split for hyperparameter
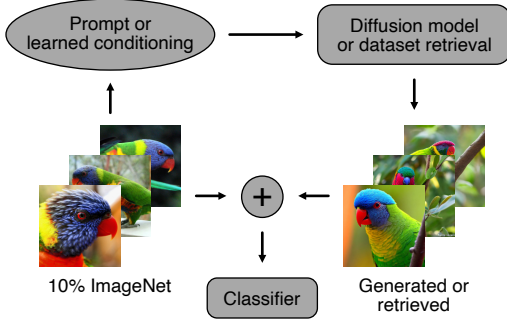
Figure 2: *Experimental protocol.* We generate images by guiding the Stable Diffusion model by text prompts or by learned conditioning, or retrieving images by a nearest neighbor search in the CLIP embedding space of the DM's training data. We then train the downstream classifier on the original 10% ImageNet data augmented by the additional data and evaluate on the original ImageNet validation split.
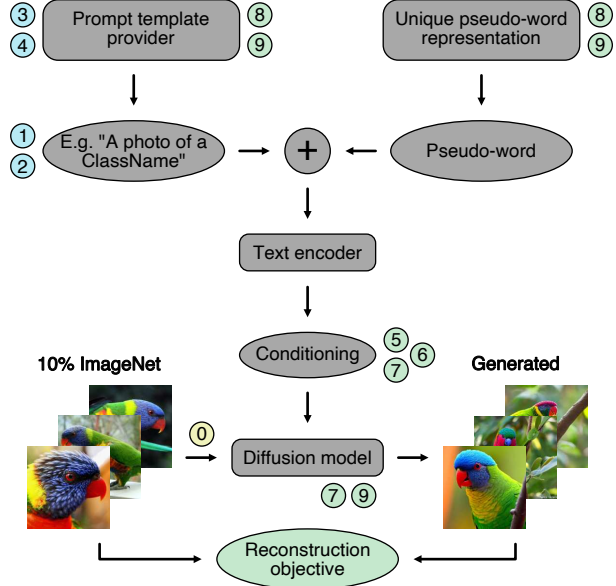


Figure 3: *Taxonomy of diffusion model based augmentation methods.* All considered methods adapt different components of the prompting, conditioning mechanisms, and the fine-tuning of the diffusion model. We reference each method by a circled number (see section 3.1 for details). Some methods edit the prompts while keeping the DM frozen: ①, ② use a single prompt for each class and ③, ④ use multiple prompts from a set of templates. Another family of methods optimize the conditioning vectors for the given images: ⑤, ⑥ only optimize the conditioning vector keeping the DM frozen, while ⑦ also jointly fine-tunes the DM. Instead of optimizing the conditioning vector, ⑧ learns a pseudo-word description of the class using multiple prompts keeping the DM frozen, while ⑨ additionally fine-tunes the DM. ⓪ does not adapt any component of the DM and relies on encoding and decoding to create variations of the given images.

optimization and model selection, which we refer to as the validation split. All trained classification models are evaluated on the original ILSVRC2012 validation split.

**Data augmentation and classifier training.** For each augmentation strategy, we generate 390 samples per class ensuring that the number of images at least tripled per class, as in our subsampled dataset each class contains between 74 and 130 images. We resample the additional data into 5 sets containing the same number of samples per class as our ImageNet subset and train a ResNet-50 classifier [14] on each to derive variance estimates. During training, the samples are further augmented by random resizing and cropping as is typical in ImageNet training. We trained the classifier with a batch size 256 and a learning rate of 0.1 which we divided by 10 when validation accuracy failed to improve for 10 epochs. We stopped training when the validation accuracy did not improve for 20 epochs or after at most 270 epochs and used the highest validation accuracy checkpoint for final scoring. Each model was trained on 8 NVIDIA T4 GPUs using distributed data parallel training.

## 3.1 Augmentation methods

The benchmarked augmentation methods can be grouped into four categories: (1) guidance-free diffusion model sampling, (2) simple conditioning techniques with prompts based on the objects' class label, and (3) personalization techniques that fine-tune the diffusion model conditioning and optionally the diffusion model itself to the classifier's data domain. We compare diffusion model approaches to a simple baseline (4) using images retrieved from the dataset

that the diffusion model was trained on.

**Unconditional generation.** Following the procedure of ⓪ BOOMERANG [19], we investigated a guidance-free method that does not require conditioning or updating the diffusion model. Instead, the approach adds noise to individual samples before denoising them.

**Prompt conditioning.** We explore several prompt-based methods of guiding the DM to produce samples for a specific class (Figure 3). ① SIMPLE PROMPT: We condition the model by simple prompts containing the object's class name $n$, prompting the DM with "A photo of $n$." and a version of it, ② SIMPLE PROMPT (NO WS), stripping whites-

3

pace, $w(\cdot)$, from class names, "A photo of $w(n)$." ③ CLIP PROMPTS: We add sampling prompts from the set of CLIP [25] text-encoder templates, e.g. "a photo of many $w(n)$.", "a black and white photo of the $w(n)$.", etc. and ④ SARIYILDIZ ET AL. PROMPTS: a set of templates proposed to create a synthetic ImageNet clone [34].

**Fine-tuning the diffusion model.** We explore various methods for fine-tuning a DM for class-personalized sampling to improve reconstruction of the classification dataset (Figure 3). ⑤ FT CONDITIONING: freezing the DM and optimizing one conditioning per class. ⑥ FT CLUSTER CONDITIONING: optimizing multiple conditionings per class instead of just one. ⑦ inspired by IMAGIC [16], jointly fine-tuning the conditioning and the DM's denoising component. ⑧ TEXTUAL INVERSION [9]: instead of fine-tuning the conditioning, sampling prompt templates and combining them with optimizing a pseudo-word representing the class-concept. ⑨ PSEUDOWORD+DM: combining the previous approach with optimizing the DM's denoising component.

**Laion nearest neighbor retrieval.** As a baseline comparison, we propose using ⑩ RETRIEVAL to select images from the Laion dataset used to train the diffusion model. This method finds nearest neighbor images to the SIMPLE PROMPT (NO WS) class name prompts in the CLIP embedding space.

## 3.2 Implementation and training details

**Diffusion model backbone.** We used the pretrained Stable Diffusion v1.4 network, based on a latent diffusion architecture [28]. We used the provided safety checker to discard generated images if marked as NSFW, replacing them with new samples. As Stable Diffusion was trained on image sizes of 512 px, we kept this resolution for all methods.

**Prompt generation.** For prompt-based sampling methods, we generated prompts based on the ImageNet class names, defined by WordNet [22] synsets representing distinct entities in the WordNet graph. Each synset consists of one or multiple lemmas describing the class, where each lemma can consist of multiple words, e.g., "Tiger shark, Galeocerdo Cuvieri". We link each class via its synset to its class name. If a synset consists of multiple lemmas, we separate them by a comma, resulting in prompts like "A photo of tiger shark.", as we found that providing multiple lemmas led to better performance than using only the first lemma of a synset. Whenever methods inserted the class name into prompt templates, we sampled the templates randomly with replacement. Sariyildiz et al. [34] provided multiple categories of prompt templates (e.g. class name only, class

name with hypernyms, additionally combined with "multiple" and "multiple different" specifications, class name with definition, and class name with hypernym and randomly sampled backgrounds from the places dataset [46]). Here, we sampled for each category the same number of images and randomly across the background templates.

**Fine-tuning.** For all methods that require additional fine-tuning, we trained the model with the default Stable Diffusion optimization objective [28] until the validation loss stagnated or increased – no model was trained for more than 40 epochs. We set hyperparameters in accordance with published works [9, 16, 19] or existing code where available. Scaling to larger batch sizes was implemented by square-root scaling the learning rate to ensure constant gradient variance, and we performed a fine-grained grid-search for the optimal learning rate. For the FT CLUSTER CONDITIONING models, we clustered the training images within a class using k-means on the inception v3 embeddings. The conditioning vectors were initialized by encoding the SIMPLE PROMPT (NO WS) prompts with the DM's text encoder. For textual inversion [9] we fine-tuned the text embedding vector corresponding to the introduced pseudoword $n$ and sample from the provided text templates [9] to optimize the reconstruction objective for our ImageNet subset, freezing all other model components. We initialized the pseudo-word embedding with the final word in the first synset lemma (i.e., for the class "tiger shark" we used "shark"). Where this initialization resulted in multiple initial tokens, we initialized with the mean of the embeddings. For each fine-tuning run we stored checkpoints with the best train and validation loss and those corresponding to 1, 2 and 3 epochs of training. In the case of IMAGIC [16] we found that the validation loss was still decreasing after 40 epochs, however, as the quality of the reconstructed images significantly deteriorated after 8 epochs, we instead stored checkpoints for the first 8 epochs. This is similar to the Imagic training scheme, which optimized the embedding for 100 steps and the U-Net for 1,500 steps on a single image. Although we follow existing procedure as closely as possible, image quality might improve for longer model training runs. We selected the final model checkpoint based on the validation accuracy of a classifier trained on the union of augmented samples and sub-sampled ImageNet.

**Nearest neighbor retrieval.** For RETRIEVAL, we used Laion 5b [35], a publicly available dataset of 5 billion image-caption pairs extracted via web-crawler. The data was then filtered by only retaining images where the CLIP image embedding was consistent with the caption embedding. This filtering acts as a weak form of supervision, that retains those images where CLIP is more likely to work.

| Augmentation method | Accuracy (%) |
|---|---|
| **10% ImageNet** | **57.2 ± 0.2** |
| **20% ImageNet** | **70.2 ± 0.3** |
| **Boomerang** [19] ⓪ | **56.3 ± 0.3** |
| Simple prompt (no ws) [proposed] ② | 60.0 ± 0.3 |
| Simple prompt [proposed] ① | 60.1 ± 0.2 |
| Sariyildiz et al. prompts [34] ④ | 60.8 ± 0.2 |
| **CLIP prompts** [proposed] ③ | **60.9 ± 0.2** |
| FT conditioning [proposed] ⑤ | 60.8 ± 0.1 |
| FT cluster conditioning [proposed] ⑥ | 60.9 ± 0.2 |
| Imagic (conditioning & DM) [16] ⑦ | 61.0 ± 0.3 |
| Textual inversion (pseudoword) [9] ⑧ | 61.0 ± 0.4 |
| **Pseudoword+DM** (combining [9, 16]) ⑨ | **61.2 ± 0.2** |
| **Retrieval** [proposed] ⑩ | **62.6 ± 0.1** |

Table 1: *Overall performance of the data augmentation methods.* We report the top-1 accuracy and standard error of the downstream classifier trained on the augmented data. The methods are described in section 3.1 and are grouped into families as described in the beginning of section 4. The circles refer to the taxonomy in Figure 3.

The dataset provides a CLIP embedding nearest neighbors search index for each instance and an official package (clip-retrieval) allows for fast search and retrieval. We used this to retrieve 130 images per class. Images with a Laion aesthetics score of less than 5 were discarded to allow a fair comparison with Stable Diffusion 1.4 which was trained on this subset. For a fair comparison to our generative augmentation methods, we used the same safety checker model to discard images that were marked as NSFW. Due to changes in the availability of the images at the URLs in the dataset and the described filtering steps, it is often necessary to retrieve more than the desired number of images, which we do by increasing the number of nearest neighbors gradually from $1.4 \cdot 130$ to $6 \cdot 1.4 \cdot 130$ when not enough samples were found. To avoid using the same image multiple times, we applied the duplicate detector of the clip-retrieval package, however, this does not detect all near duplicate images and some duplicates are still used.

# 4 Results

Table 1 provides a results summary. In the following, section 4.1 discusses the baselines, section 4.2 the results for the unconditional DM augmentation method, section 4.3 conditioning the DM with text prompts, section 4.4 personalization approaches to diffusion models, and section 4.5 discusses retrieval-based techniques. Each block in Table 1 corresponds to one section.

**Challenges for generative augmentations methods.** Before discussing the results in detail, we highlight four challenges for generative augmentation. *Class ambiguity:* Some class labels in ImageNet have multiple meanings, leading to class ambiguity (see samples in Figure 5A-B). For instance, for some methods, generated images for the class "papillon" show butterflies of the species papilio phorbanta (French: papillon la pâture), however, in ImageNet this class is a dog breed. Moreover, in some cases the label "crane" leads to generated images referring either to the machine or the bird and the label "desktop computer" sometimes leads to generated images showing a laptop on top of desks, instead of a traditional desktop computer. *Diversity:* Another challenge is to generate images of sufficient diversity. Most of the considered augmentation methods generate objects in similar styles and semantic settings. For instance, in ImageNet objects can appear in various semantic styles (e.g., empty and full beer glass in different shapes and contexts), image styles (black and white, bright, dark, etc.), and image qualities (e.g., blurred or JPEG corrupted samples). *Domain shift:* Augmentation methods must be flexible enough to adapt to the distribution of images in ImageNet. *Fidelity:* Finally, generative methods must provide high-quality samples that match the photorealism of ImageNet.

## 4.1 Baselines

To establish the baseline classifier performance , we train on our 10% ImageNet subset (no added data), achieving a top-1 accuracy of $57.2 \pm 0.2$. Since all augmentation methods double the size of the dataset, we consider an upper bound performance by using 20% of original ImageNet (no added data) with a classifier accuracy of $70.2 \pm 0.3$.

## 4.2 Unconditional generation

BOOMERANG [19] sequentially adds Gaussian noise and uses a diffusion model to denoise elements of the target dataset, without altering the diffusion model or prompts. This results in a lower top-1 accuracy of $56.3 \pm 0.3$ compared to the original subset of ImageNet. Figure 4 shows that this method does not induce significant diversity in the dataset and only slightly distorts the original images, leading to an overall decrease in performance. However, since this method is directly applied on the target dataset samples, it does not suffer from domain shift or class ambiguity.

**Summary.** Using diffusion models off the shelf, i.e., without adaptation or conditioning using the target dataset, is not beneficial and can even deteriorate the classifier performance compared to the ImageNet baseline.

A ImageNet  B Boomerang  C Difference

Figure 4: *Example images generated by* BOOMERANG *[19]. The generated images lack diversity and effectively only add noise to the original ones.* **A.** Example images from 10% ImageNet. **B.** The augmentations produced by Boomerang. **C.** There is only a small difference between the original and augmented images.

## 4.3 Prompt-conditioning of the diffusion model

One way to adapt diffusion models to the target dataset is via conditioning mechanisms that guide image generation. We investigate several text prompting strategies.

**Simple prompt conditioning.** We investigate a SIMPLE PROMPT conditioning method that uses the prompt "A photo of $n$.", where $n$ is replaced by the class name, e.g., "A photo of tiger shark, Galeocerdo cuvieri." This improved over the 10% ImageNet baseline to a top-1 classification accuracy of $60.1 \pm 0.2$, but Figure 5B shows clear problems with class ambiguity and a lack of diversity.

**Tackling class ambiguity by white space removal.** While generally mitigating class ambiguity is a hard problem, we focus on the ambiguity introduced by class names composed of multiple words. To this end, we investigated the variant SIMPLE PROMPT (NO WS) (no white space) that removed the white space from the class name used by SIMPLE PROMPT, e.g., "A photo of desktopcomputer.". While class ambiguity reduced for some classes (e.g., desktop computer; Figure 5C) it was ineffective for others (e.g., papillon) and overall resulted in slightly degraded performance $60.0 \pm 0.3$ compared to keeping the white space. We explore more advanced methods in the following sections.

**Improving sample diversity by diverse prompt templates.** To increase the diversity of the samples, we made use of multiple prompt templates. We use CLIP PROMPTS which randomly selects one of the text templates provided by CLIP [25] (see Figure 5D for examples and [25] for a full list), increasing classification accuracy to $60.9 \pm 0.2$. This method performed best among all prompt-based techniques. Interestingly, the overall performance increased even though some prompts lead to synthetic images that did

not match the style of ImageNet samples (e.g., "a cartoon photo of the papillon.") or images with texture-like contents instead of objects (e.g., papillon, image to the bottom right in Figure 5D). Surprisingly, these slightly more elaborate templates just adding few more words compared to SIMPLE PROMPT (NO WS) (e.g., "a bad photo of a papillon.") improved class ambiguity in some cases (e.g., papillon; Figure 5D). A similar method using multiple prompt templates was introduced by Sariyildiz et al. [34] (see section 3.2 for details) and performs slightly worse in our augmentation setting ($60.8 \pm 0.2$ classification accuracy).

**Summary.** Augmenting a dataset with images sampled by prompt-based conditioning techniques improves the downstream classifier performance, but various challenges remain. ImageNet is more than a decade old, and some images included in the dataset are older and of different style when compared to the images created by the recently trained Stable Diffusion model (e.g., desktop computer; Figure 5A-D). In other cases, generated images do not match the domain of ImageNet samples, for instance because they are too artificial, sometimes even like a computer rendering (e.g., motor scooter; Figure 5A-D), their texture does not match (e.g., papillon; Figure 5A,D) or the prompt does not match the desired style (e.g., cartoon; Figure 5A,D). This tendency may have been amplified by Stable Diffusion v1.4 being trained on a subset of Laion, which was filtered to contain only aesthetic images.[1]

## 4.4 Personalizing the diffusion model

In the previous section, we found that adapting the diffusion models by only editing the prompt offers limited improvement when used for augmentation. This section explores a more advanced set of methods that additionally fine-tune parts of the diffusion model to "personalize" it to the 10% ImageNet training images by optimizing the DM reconstruction objective (c.f. Figure 3). Originally, these methods were proposed in the context of personalizing the model to a specific object or concept, i.e., a single or only a few images. In the following, we explore extensions for augmenting larger target datasets.

**Fine-tuning the conditioning vectors.** The general diffusion model architecture has multiple components that can be fine-tuned (see Figure 3). We start by fine-tuning one conditioning vector per class to optimize the reconstruction objective for our ImageNet subsample, while keeping the other model weights fixed. This method, dubbed FT CONDITIONING (fine-tune conditioning), achieves an augmentation accuracy of $60.8 \pm 0.1$ and is on par with the best

---

[1] https://github.com/CompVis/stable-diffusion

Figure 5: *Example images obtained from the investigated augmentation methods.* **A.** 10% ImageNet original images. **B.-D.** Images generated by our prompt-based sampling techniques: (B) SIMPLE PROMPT, (C) SIMPLE PROMPT (NO WS) and (D) CLIP PROMPTS, respectively, where in (D) we show images to the prompts "a bad photo of a $w(n)$", "a black and white photo of the $w(n)$", "a cartoon $w(n)$", and "a photo of many $w(n)$". **E.-H.** Images generated by the fine-tuned diffusion models, specifically (E) FT CONDITIONING, (F) TEXTUAL INVERSION, (G) IMAGIC, and (H) PSEUDOWORD+DM. **I.** Examples of RETRIEVAL from the diffusion model's training set. Best viewed when zoomed in.

prompt editing method CLIP PROMPTS. Interestingly, while this method achieves good augmentation performance, the generated images do not look photorealistic, and suffer from noise and missing backgrounds (Figure 5E).

**Fine-tuning clusters of conditioning vectors.** Using a single conditioning vector for all images from one class might be insufficient to capture the full class variability. To see if this is the case, we explore FT CLUSTER CONDITION-ING. We generate $k$ clusters of images per class, before fitting a conditioning vector to each individual cluster (see section 3.2 for details). We find that using $k = 5$ clusters slightly improves the accuracy by $0.1$ percentage points over using $k = 1$. For larger $k = 10$ and $k = 15$, performance decreased ($60.8 \pm 0.2$ and $60.6 \pm 0.3$ accuracy). This might be due to the smaller number of images per cluster, making the fine-tuning more prone to over-fitting.

**Textual inversion.** To reduce the number of unrealistic images, we explored a variant of TEXTUAL INVERSION [9]. This method learns a pseudoword $n$ representing the class

concept combined with a randomly drawn textual description of the generated image style (e.g., "a photo of a $n$", "a rendering of a $n$", etc.; see [9] for a full list). While this method improves the photorealism of the generated samples (Figure 5F), it only slightly improves the augmentation accuracy over the simple fine-tuning of conditioning vectors (FT CONDITIONING) by $0.1$ percentage points (Table 1).

**Fine-tuning the denoising.** We now explore fine-tuning the DMs denoising module jointly with the conditioning vector. This idea stems from IMAGIC [16] and results in an augmentation accuracy of $61.0 \pm 0.3$ (examples in Figure 5G), which is on par with TEXTUAL INVERSION.

Finally, we combine the best-performing methods: jointly optimizing a pseudo-word per class (TEXTUAL INVERSION [9]) and the DMs denoising module (IMAGIC [16]). We denote this method by PSEU-DOWORD+DM and it has an accuracy of $61.2 \pm 0.2$, outperforming all other DM based techniques investigated. The generated images also show improved photorealism over TEXTUAL INVERSION and IMAGIC (Figure 5H).

**Summary.** Personalizing diffusion models improves in matching the domain of the 10% ImageNet images better (e.g., generated images for "desktop computer" now look similarly old as in ImageNet), improve upon prompt-based techniques to reduce class-ambiguity (e.g. papillon and desktop computer showed ambiguity in Figure 5B but not in panels E-H), show good sample variety and the best-performing PSEUDOWORD+DM method enhances photorealism and reduces the artistic style of generated images. This model performs best across all generative augmentation techniques investigated, and suggests that we can leverage personalization techniques to combine the DM's knowledge of billions of annotated images with learning the domain distribution of the data we want to augment.

## 4.5 Retrieving from the pre-training data

We showed that diffusion models are helpful for creating augmentation data, however, it is unclear how much value the DM's generative capabilities add compared to simply using their pre-training data directly for augmentation. To answer this question, we propose a simple RETRIEVAL method fetching images from the DM's pre-training data that are semantically closest to the SIMPLE PROMPT (NO WS) prompts (see section 3.2 for details). Augmenting with this data outperformed all investigated diffusion model based approaches at $62.6 \pm 0.1$ top-1 accuracy, while being computationally less demanding. As retrieved images are real-world images, they improved over diffusion model generated images in terms of photorealism, and retrieved images show good variety and detail. However, retrieval also suffers from the mentioned class ambiguity ("papillon" and "mailbag, postbag" in Figure 5I), reflecting a mismatch of concepts from CLIP latent space to ImageNet classes as we measured semantic similarity as distance of CLIP embeddings.

**When does augmentation help?** To better understand the failure cases of the augmentation methods, we check for each class in ImageNet if the augmentation method is beneficial. To this end, for each class, we compute the performance improvement of the downstream classifier using the augmentation method compared to only using the original ImageNet samples. Figure 6 shows the distribution of improvements for the RETRIEVAL augmentation method and the best DM-based method PSEUDOWORD+DM. Both methods improve the performance for most of the classes, however, there are still many classes where the performance is decreased up to 10%. The distribution of improvements for RETRIEVAL is similarly shaped as for PSEUDOWORD+DM, but significantly pushed to the right. Systematically investigating these failure cases might be a fruitful avenue to further improve generative augmentation.
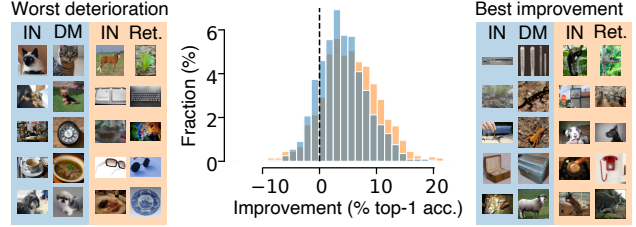


Figure 6: *Augmentation improvement for each class.* The distribution of improvements for each class is shown for the best-performing DM-based method PSEUDOWORD+DM (blue) and RETRIEVAL (orange). For each class the improvement is computed by comparing the downstream classification performance using augmentations compared to only using the original 10% ImageNet samples. Images to the left and right show examples of classes with worst deterioration and greatest improvement.

**Summary.** Augmenting with retrieved images outperforms all DM based approaches, indicating that DM generated images still have significant shortcomings even after personalization. However, retrieval still underperforms the 20% ImageNet upper bound performance.

# 5 Discussion

Diffusion models have shown their effectiveness in many application areas, and using them for data augmentation is an intriguing research direction.

We have evaluated multiple methods to prompt and personalize diffusion models on their usefulness for data augmentation, showing that none of them beat the simple baseline of retrieving images from the diffusion model's pre-training dataset. Why is the retrieval baseline so hard to beat? One reason might be that retrieval potentially accesses more information, as the pre-training dataset is usually much larger than the weights of the generative models trained on it. Although it has been argued that diffusion models might partially compress the training data [37, 5], it is still unclear if the generative model captured all relevant information. However, diffusion models could possibly improve upon the so-far superior retrieval by generating a large number of additional data and more diverse and compositionally novel images, for instance by generating out-of-domain samples (e.g., "a photograph of an astronaut riding a horse") [27, 33, 28]. Furthermore, diffusion models allow, in principle, for more controlled adaptation than retrieval methods. We showed that personalization methods are a good step in this direction, however, they typically focus only on creating variants of specific given images. To unlock the true potential of diffusion models for data aug-

mentation, new methods that capture the target dataset manifold as a whole, are needed.

**Limitations and future work.** While using a subset of ImageNet is an established benchmark for generative data augmentation methods [19, 29, 15], a limitation of our work is that we only consider this setting. It would be interesting to extend our analysis to other datasets (e.g., medical images) and to explore the investigated methods for out-of-distribution generalization [24, 40]. Diffusion models are a very active research area, and new methods and applications are published on a daily basis. Hence, our experiments could only capture a subset of possible methods and newly proposed extensions of them (in total, eleven methods). We compared all methods with the same augmentation budget, however, diffusion models can generate arbitrarily many samples. Thus, it would be interesting to explore their scaling behavior for higher augmentation ratios. We have shown that simple retrieval is a very strong competitor for data augmentation. Further improvements could be introduced by diversifying the retrieval set [41, 42] or retrieving images using linear combinations of inputs [47].

**Conclusion.** We have shown that a simple retrieval baseline can outperform a wide range of diffusion model based augmentations. However, given the fast rate of progress in this field it is not possible to definitively say that retrieval can not be beaten, and we believe that diffusion models have the potential to improve over this baseline. Nonetheless, the strength of retrieval's performance makes it clear that future works using diffusion models for augmentation should also compare against this baseline. We hope that our paper provides ground for researchers to benchmark generative augmentation methods and assess their benefit by comparing them with retrieval baselines. In the longer-term, we hope that new methods can be developed which combine retrieval and generation, leading to greater improvements in the diversity and quality of augmented images.

# 6 Acknowledgements

# References

[1] Mohamed Akrout, Bálint Gyepesi, Péter Holló, Adrienn Poór, Blága Kincső, Stephen Solis, Katrina Cirone, Jeremy Kawahara, Dekker Slade, Latif Abid, Máté Kovács, and István Fazekas. Diffusion-based Data Augmentation for Skin Disease Classification: Impact Across Original Medical Datasets to Fully Synthetic Images, Jan. 2023. arXiv:2301.04802 [cs, eess].

[2] Antreas Antoniou, Amos Storkey, and Harrison Edwards. Data augmentation generative adversarial networks. *arXiv preprint arXiv:1711.04340*, 2017.

[3] Victor Besnier, Himalaya Jain, Andrei Bursuc, Matthieu Cord, and Patrick Pérez. This dataset does not exist: training models from generated images. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2020.

[4] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.

[5] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting Training Data from Diffusion Models, Jan. 2023. arXiv:2301.13188 [cs].

[6] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020.

[7] Cristóbal Esteban, Stephanie L Hyland, and Gunnar Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.

[8] Maayan Frid-Adar, Idit Diamant, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit Greenspan. Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification. *Neurocomputing*, 321:321–331, 2018.

[9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion, Aug. 2022. arXiv:2208.01618 [cs].

[10] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. Designing an Encoder for Fast Personalization of Text-to-Image Models, Feb. 2023. arXiv:2302.12228 [cs].

[11] Sahra Ghalebikesabi, Leonard Berrada, Sven Gowal, Ira Ktena, Robert Stanforth, Jamie Hayes, Soham De, Samuel L. Smith, Olivia Wiles, and Borja Balle. Differentially Private Diffusion Models Generate Useful Synthetic Images, Feb. 2023. arXiv:2302.13861 [cs, stat].

[12] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2918–2928, 2021.

[13] Partha Ghosh, Dominik Zietlow, Michael J Black, Larry S Davis, and Xiaochen Hu. Invgan: invertible gans. In *Pattern Recognition: 44th DAGM German Conference, DAGM GCPR 2022, Konstanz, Germany, September 27–30, 2022, Proceedings*, pages 3–19. Springer, 2022.

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[15] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. Is synthetic data from generative models ready for image recognition?, Oct. 2022. arXiv:2210.07574 [cs].

[16] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-Based Real Image Editing with Diffusion Models, Oct. 2022. arXiv:2210.09276 [cs].

[17] Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Bigdatasetgan: Synthesizing imagenet with pixel-wise annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21330–21340, 2022.

[18] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. RePaint: Inpainting using Denoising Diffusion Probabilistic Models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11451–11461, June 2022. ISSN: 2575-7075.

[19] Lorenzo Luzi, Ali Siahkoohi, Paul M. Mayer, Josue Casco-Rodriguez, and Richard Baraniuk. Boomerang: Local sampling on image manifolds using diffusion models, Oct. 2022. arXiv:2210.12100 [cs, stat].

[20] Giovanni Mariani, Florian Scheidegger, Roxana Istrate, Costas Bekas, and Cristiano Malossi. Bagan: Data augmentation with balancing gan. *arXiv preprint arXiv:1803.09655*, 2018.

[21] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021.

[22] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.

[23] Saman Motamed, Patrik Rogalla, and Farzad Khalvati. Data augmentation using generative adversarial networks (gans) for gan-based detection of pneumonia and covid-19 in chest x-ray images. *Informatics in Medicine Unlocked*, 27:100779, 2021.

[24] Jielin Qiu, Yi Zhu, Xingjian Shi, Florian Wenzel, Zhiqiang Tang, Ding Zhao, Bo Li, and Mu Li. Are multimodal models robust to image and text perturbations? *arXiv preprint arXiv:2212.08044*, 2022.

[25] 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. *arXiv:2103.00020 [cs]*, Feb. 2021. arXiv: 2103.00020.

[26] Vikram V Ramaswamy, Sunnie SY Kim, and Olga Russakovsky. Fair attribute classification through latent space de-biasing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9301–9310, 2021.

[27] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical Text-Conditional Image Generation with CLIP Latents, Apr. 2022. arXiv:2204.06125 [cs].

[28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.

[29] Aniket Roy, Anshul Shah, Ketul Shah, Anirban Roy, and Rama Chellappa. DiffAlign : Few-shot learning using diffusion based synthesis and alignment, Dec. 2022. arXiv:2212.05404 [cs].

[30] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation, Aug. 2022. arXiv:2208.12242 [cs].

[31] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

[32] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-Image Diffusion Models. In *ACM SIGGRAPH 2022 Conference Proceedings*, SIGGRAPH '22, pages 1–10, New York, NY, USA, July 2022. Association for Computing Machinery.

[33] 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, May 2022. arXiv:2205.11487 [cs].

[34] Mert Bulent Sariyildiz, Karteek Alahari, Diane Larlus, and Yannis Kalantidis. Fake it till you make it: Learning(s) from a synthetic ImageNet clone, Dec. 2022. arXiv:2212.08420 [cs].

[35] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.

[36] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.

[37] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models, Dec. 2022. arXiv:2212.03860 [cs].

[38] Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective Data Augmentation With Diffusion Models, Feb. 2023. arXiv:2302.07944 [cs].

[39] Florian Wenzel, Andrea Dittadi, Peter Vincent Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, et al. Assaying out-of-distribution generalization in transfer learning. *arXiv preprint arXiv:2207.09239*, 2022.

[40] Florian Wenzel, Andrea Dittadi, Peter Vincent Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, Bernhard Schölkopf, and Francesco Locatello. Assaying Out-Of-Distribution Generalization in Transfer Learning, Oct. 2022. arXiv:2207.09239 [cs, stat].

[41] Florian Wenzel, Jasper Snoek, Dustin Tran, and Rodolphe Jenatton. Hyperparameter ensembles for robustness and uncertainty quantification. In *Advances in Neural Information Processing Systems*, 2020.

[42] Yisong Yue and Carlos Guestrin. Linear submodular bandits and their application to diversified retrieval. In *Advances in Neural Information Processing Systems*, 2011.

[43] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.

[44] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.

[45] Yifan Zhang, Daquan Zhou, Bryan Hooi, Kai Wang, and Jiashi Feng. Expanding Small-Scale Datasets with Guided Imagination, Nov. 2022. arXiv:2211.13976 [cs].

[46] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.

[47] Dominik Zietlow, Michael Lohaus, Guha Balakrishnan, Matthäus Kleindessner, Francesco Locatello, Bernhard Schölkopf, and Chris Russell. Leveling down in computer vision: Pareto inefficiencies in fair deep classifiers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10410–10421, 2022.
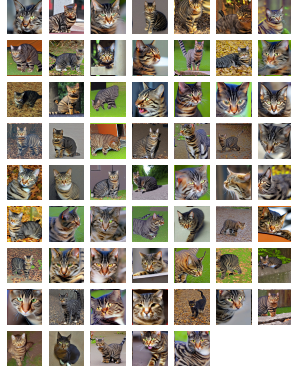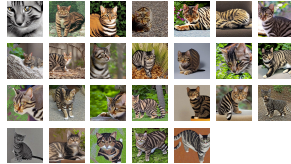
# Supplementary material

**A** Tiger cat

ImageNet          FT cluster conditioning

Cluster 1



Cluster 2



Cluster 3



Cluster 4



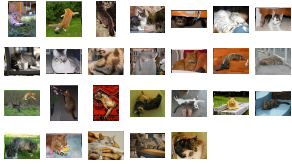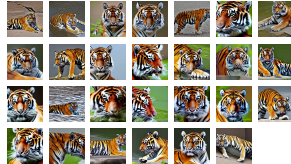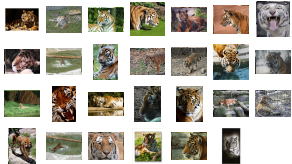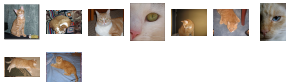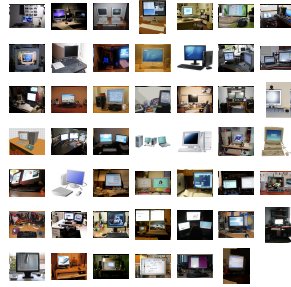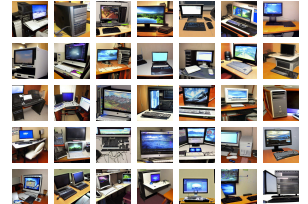Cluster 5



**B** Desktop computer

ImageNet          FT cluster conditioning
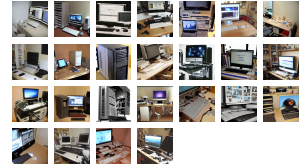
Cluster 1



Cluster 2



Cluster 3



Cluster 4



Cluster 5



Figure 7: FT CLUSTER CONDITIONING *with $k = 5$ clusters compared to ImageNet.* Semantically similar ImageNet images are clustered together and one conditioning is learned for each cluster to reconstruct the training images (see section 3 for details). We exclude images resembling human faces to preserve data privacy. **A.** Examples for the class "tiger cat" which is ambiguous in ImageNet itself (left column). **B.** Examples for the class "desktop computer". Best viewed when zoomed in.

Figure 8: CLIP PROMPTS *examples for each CLIP text template.*

a photo of a papillon

a photo of one papillon

a rendering of a papillon

a close-up photo of the papillon

a cropped photo of the papillon

a rendition of the papillon

the photo of a papillon

a photo of the clean papillon

a photo of a clean papillon

a rendition of a papillon

a photo of a dirty papillon

a photo of a nice papillon

a dark photo of the papillon

a good photo of a papillon

a photo of my papillon

a photo of the nice papillon

a photo of the cool papillon

a photo of the small papillon

a close-up photo of a papillon

a photo of the weird papillon

a bright photo of the papillon

a photo of the large papillon

a cropped photo of a papillon

a photo of a cool papillon

a photo of the papillon

a photo of a small papillon

a good photo of the papillon

Figure 9: *examples for each* TEXTUAL INVERSION *text template.*

Epoch 1

Epoch 4 (used)

Epoch 10

Epoch 20

Epoch 30

Epoch 40

Figure 10: *Examples of* IMAGIC *optimization for various epochs.*