

---

# GeNIe: Generative Hard Negative Images Through Diffusion

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Data augmentation is crucial in training deep models, preventing them from over-  
2 fitting to limited data. Recent advances in generative AI, e.g., diffusion mod-  
3 els, have enabled more sophisticated augmentation techniques that produce data  
4 resembling natural images. We introduce GeNIe a novel augmentation method  
5 which leverages a latent diffusion model conditioned on a text prompt to combine  
6 two contrasting data points (an image from the source category and a text prompt  
7 from the target category) to generate challenging augmentations. To achieve this,  
8 we adjust the noise level (equivalently, number of diffusion iterations) to ensure  
9 the generated image retains low-level and background features from the source  
10 image while representing the target category, resulting in a *hard negative* sample  
11 for the source category. We further automate and enhance GeNIe by adaptively  
12 adjusting the noise level selection on a per image basis (coined as GeNIe-Ada),  
13 leading to further performance improvements. Our extensive experiments, in both  
14 few-shot and long-tail distribution settings, demonstrate the effectiveness of our  
15 novel augmentation method and its superior performance over the prior art.

## 16 1 Introduction

17 Augmentation has become an integral part of training deep learning models, particularly when faced  
18 with limited training data. For instance, when it comes to image classification with limited number  
19 of samples per class, model generalization ability can be significantly hindered. Simple transfor-  
20 mations like rotation, cropping, and adjustments in brightness artificially diversify the training set,  
21 offering the model a more comprehensive grasp of potential data variations. Hence, augmentation  
22 can serve as a practical strategy to boost the model’s learning capacity, minimizing the risk of overfit-  
23 ting and facilitating effective knowledge transfer from limited labelled data to real-world scenarios.  
24 Various image augmentation methods, encompassing standard transformations, and learning-based  
25 approaches have been proposed [16, 15, 110, 111, 100]. Some augmentation strategies combine two  
26 images possibly from two different categories to generate a new sample image. The simplest ones  
27 in this category are MixUp [111] and CutMix [110] where two images are combined in the pixel  
28 space. However, the resulting augmentations often do not lie within the manifold of natural images  
29 and act as out-of-distribution samples that will not be encountered during testing.

30 Recently, leveraging generative models for data augmentation has gained an upsurge of attention  
31 [100, 83, 63, 35]. These interesting studies, either based on fine-tuning or prompt engineering of  
32 diffusion models, are mostly focused on generating *generic augmentations* without considering the  
33 impact of other classes and incorporating that information into the generative process for a classifi-  
34 cation context. We take a different approach to generate challenging augmentations near the decision  
35 boundaries of a downstream classifier. Inspired by diffusion-based image editing methods [67, 63]  
36 some of which are previously used for data augmentation, we propose to use conditional latent dif-

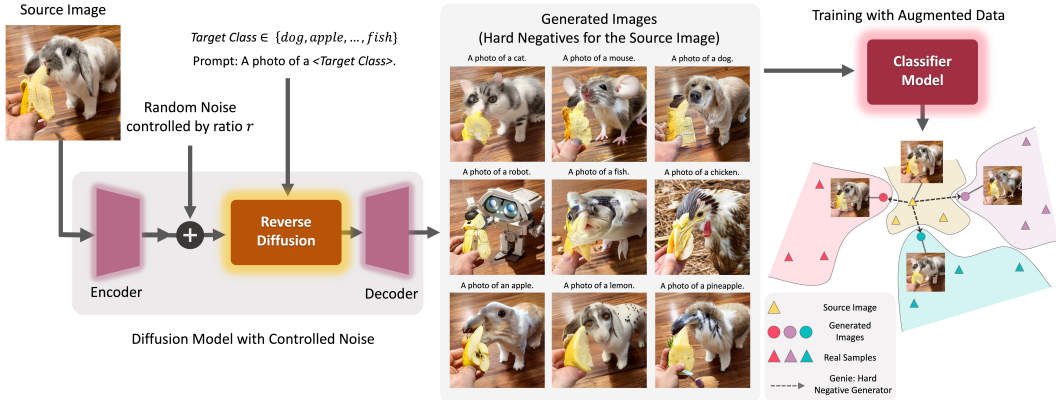


Figure 1: **Generative Hard Negative Images Through Diffusion (GeNIe)**: generates hard negative images that belong to the target category but are similar to the source image from low-level feature and contextual perspectives. GeNIe starts from a source image passing it through a partial noise addition process, and conditioning it on a different target category. By controlling the amount of noise, the reverse latent diffusion process generates images that serve as *hard negatives* for the source category.

fusion models [81] for generating *hard negative* images. Our core idea (coined as GeNIe) is to sample source images from various categories and prompt the diffusion model with a contradictory text corresponding to a different target category. We demonstrate that the choice of noise level (or equivalently number of iterations) for the diffusion process plays a pivotal role in generating images that semantically belong to the target category while retaining low-level features from the source image. We argue that these generated samples serve as *hard negatives* [108, 65] for the source category (or from a dual perspective hard positives for the target category). To further enhance GeNIe, we propose an adaptive noise level selection strategy (dubbed as GeNIe-Ada) enabling it to adjust noise levels automatically per sample.

To establish the impact of GeNIe, we focus on two challenging scenarios: *long-tail* and *few-shot* settings. In real-world applications, data often follows a long-tail distribution, where common scenarios dominate and rare occurrences are underrepresented. For instance, a person jaywalking a highway causes models to struggle with such unusual scenarios. Combating such a bias or lack of sufficient data samples during model training is essential in building robust models for self-driving cars or surveillance systems, to name a few. Same challenge arises in few-shot learning settings where the model has to learn from only a handful of samples. Our extensive quantitative and qualitative experimentation, on a suite of few-shot and long-tail distribution settings, corroborate the effectiveness of the proposed novel augmentation method (GeNIe, GeNIe-Ada) in generating hard negatives, corroborating its significant impact on categories with a limited number of samples. A high-level sketch of GeNIe is illustrated in Fig. 1. Our main contributions are summarized below:

- We introduce GeNIe, a novel yet elegantly simple diffusion-based augmentation method to create challenging augmentations in the manifold of natural images. For the first time, to our best knowledge, GeNIe achieves this by combining two sources of information (a source image, and a contradictory target prompt) through a noise-level adjustment mechanism.
- We further extend GeNIe by automating the noise-level adjustment strategy on a per-sample basis (called GeNIe-Ada), to enable generating hard negative samples in the context of image classification, leading also to further performance enhancement.
- To substantiate the impact of GeNIe, we present a suit of quantitative and qualitative results including extensive experimentation on two challenging tasks: few-shot and long tail distribution settings corroborating that GeNIe (and its extension GeNIe-Ada) significantly improve the downstream classification performance.

## 2 Related Work

**Data Augmentations.** Simple flipping, cropping, colour jittering, and blurring are some forms of image augmentations [91]. These augmentations are commonly adopted in training deep learning models. However, using these data augmentations is not trivial in some domains. For example, using blurring might remove important low-level information from medical images. More advanced

73 approaches, such as MixUp [111] and CutMix [110], mix images and their labels accordingly [37,  
74 59, 47, 17]. However, the resulting augmentations are not natural images anymore, and thus, act  
75 as out-of-distribution samples that will not be seen at test time. Another strand of research tailors  
76 the augmentation strategy through a learning process to fit the training data [23, 16, 15]. Unlike the  
77 above methods, we propose to utilize pre-trained latent diffusion models to generate hard negatives  
78 (in contrast to generic augmentations) through a noise adaptation strategy discussed in Section 3.

79 **Data Augmentation with Generative Models.** Using synthesized images from generative models  
80 to augment training data has been studied before in many domains [30, 86], including domain adap-  
81 tation [41], visual alignment [71], and mitigation of dataset bias [88, 36, 73]. For example, [73]  
82 introduces a methodology aimed at enhancing test set evaluation through augmentation. While pre-  
83 vious methods predominantly relied on GANs [114, 51, 101] as the generative model, more recent  
84 studies promote using diffusion models to augment the data [81, 35, 89, 100, 4, 62, 83, 42, 28, 26, 8].  
85 More specifically, [100, 83, 35, 4] study the effectiveness of text-to-image diffusion models in data  
86 augmentation by diversification of each class with synthetic images. [100] leverages a text-to-image  
87 diffusion model and fine-tunes it on the downstream dataset using textual-inversion [31] to increase  
88 the diversity of existing samples. [83] also utilizes a text-to-image diffusion model, but with a BLIP  
89 [53] model to generate meaningful captions from the existing images. [42] utilizes diffusion models  
90 for augmentation to correct model mistakes. [28] uses CLIP [76] to filter generated images. [26]  
91 utilizes text-based diffusion and a large language model (LLM) to diversify the training data. [8]  
92 uses an LLM to generate text descriptions of failure modes associated with spurious correlations,  
93 which are then used to generate synthetic data through generative models. The challenge here is that  
94 the LLM has little understanding of such failure scenarios and contexts.

95 We take a completely different approach here, without relying on any extra source of information  
96 (e.g., through an LLM). Inspired by image editing approaches such as Boomerang [63] and SDEdit  
97 [67], we propose to adaptively guide a latent diffusion model to generate *hard negatives* images  
98 [65, 108] on a per-sample basis per category. In a nutshell, the aforementioned studies focus on im-  
99 proving the diversity of each class with effective prompts and diffusion models, however, we focus  
100 on generating effective *hard negative* samples for each class by combining two sources of contra-  
101 dicting information (images from the source category and text prompt from the target category).

102 **Language Guided Recognition Models.** Vision-Language foundation models (VLMs) [2, 76, 81,  
103 84, 77, 78] utilize human language to guide the generation of images or to extract features from  
104 images that are aligned with human language. For example, CLIP [76] shows decent zero-shot  
105 performance on many downstream tasks by matching images to their text descriptions. Some recent  
106 works improve the utilization of human language in the prompt [25, 72], and others use a diffusion  
107 model directly as a classifier [49]. Similar to the above, we use a foundation model (Stable Diffusion  
108 1.5 [81]) to improve the downstream task. Concretely, we utilize category names of the downstream  
109 tasks to augment their associate training data with hard negative samples.

110 **Few-Shot Learning.** In Few-shot Learning (FSL), we pre-train a model with abundant data to learn  
111 a rich representation, then fine-tune it on new tasks with only a few available samples. In supervised  
112 FSL [10, 1, 74, 109, 27, 54, 95, 116, 92], pretraining is done on a labeled dataset, whereas in  
113 unsupervised FSL [43, 103, 61, 75, 3, 46, 39, 66, 90] the pre-training has to be conducted on an  
114 unlabeled dataset. We assess the impact of GeNIe on a number of few-shot scenarios and state-of-  
115 the-art baselines by accentuating on its impact on the few-shot inference stage.

### 116 3 Proposed Method: GeNIe

117 Given a source image  $X_S$  from category  $S = \langle \text{source category} \rangle$ , we are interested in generating a  
118 target image  $X_T$  from category  $T = \langle \text{target category} \rangle$ . In doing so, we intend to ensure the low-  
119 level visual features or background context of the source image are preserved, so that we generate  
120 samples that would serve as *hard negatives* for the source image. To this aim, we adopt a conditional  
121 latent diffusion model (such as Stable Diffusion, [81]) conditioned on a text prompt of the following  
122 format “A photo of a  $T = \langle \text{target category} \rangle$ ”.

123 **Key Idea.** GeNIe in its basic form is a simple yet effective augmentation sample generator for  
124 improving a classifier  $f_\theta(\cdot)$  with the following two key aspects: (i) inspired by [63, 67] instead of  
125 adding the full amount of noise  $\sigma_{max}$  and going through all  $N_{max}$  (being typically 50) steps of  
126 denoising, we use less amount of noise ( $r\sigma_{max}$ , with  $r \in (0, 1)$ ) and consequently fewer number  
127 of denoising iterations ( $\lfloor rN_{max} \rfloor$ ); (ii) we prompt the diffusion model with a  $P$  mandating a target



Figure 2: **Effect of noise ratio,  $r$ , in GeNIe:** we employ GeNIe to generate augmentations for the target classes (motorcycle and cat) with varying  $r$ . Smaller  $r$  yields images closely resembling the source semantics, creating an inconsistency with the intended target label. By tracing  $r$  from 0 to 1, augmentations gradually transition from source image characteristics to the target category. However, a distinct shift from the source to the target occurs at a specific  $r$  that may vary for different source images or target categories. For more examples, please refer to Fig. A4.

128 category  $T$  different than the source  $S$ . Hence, we denote the conditional diffusion process as  
 129  $X_r = \text{STDiff}(X_S, P, r)$ . In such a construct, the proximity of the final decoded image  $X_r$  to the  
 130 source image  $X_S$  or the target category defined through the text prompt  $P$  depends on  $r$ . Hence, by  
 131 controlling the amount of noise, we can generate images that blend characteristics of both the text  
 132 prompt  $P$  and the source image  $X_S$ . If we do not provide much of visual details in the text prompt  
 133 (e.g., desired background, etc.), we expect the decoded image  $X_r$  to follow the details of  $X_S$  while  
 134 reflecting the semantics of the text prompt  $P$ . We argue, and demonstrate later, that the newly  
 135 generated samples can serve as *hard negative* examples for the source category  $S$  since they share  
 136 the low-level features of  $X_S$  while representing the semantics of the target category,  $T$ . Notably, the  
 137 source category  $S$  can be randomly sampled or be carefully extracted from the confusion matrix of  
 138  $f_\theta(\cdot)$  based on real training data. The latter might result in even *harder negative* samples being now  
 139 cognizant of model confusions. Finally, we will append our initial dataset with the newly generated  
 140 hard negative samples through GeNIe and (re)train the classifier model.

141 **Enhancing GeNIe: GeNIe-Ada.** One of the remarkable aspects of GeNIe lies in its simple applica-  
 142 tion, requiring only  $X_S$ ,  $P$ , and  $r$ . However, selecting the appropriate value for  $r$  poses a challenge  
 143 as it profoundly influences the outcome. When  $r$  is small, the resulting  $X_r$  tends to closely resemble  
 144  $X_S$ , and conversely, when  $r$  is large (closer to 1), it tends to resemble the semantics of the target  
 145 category. This phenomenon arises because a smaller noise level restricts the capacity of the diffusion  
 146 model to deviate from the semantics of the input  $X_S$ . Thus, a critical question emerges: how can we  
 147 select  $r$  for a particular source image to generate samples that preserve the low-level semantics of  
 148 the source category  $S$  in  $X_S$  while effectively representing the semantics of the target category  $T$ ?  
 149 We propose a method to determine an ideal value for  $r$ .

150 Our intuition suggests that by varying the noise ratio  $r$  from 0 to 1,  $X_r$  will progressively resemble  
 151 category  $S$  in the beginning and category  $T$  towards the end. However, somewhere between 0  
 152 and 1,  $X_r$  will undergo a rapid transition from category  $S$  to  $T$ . This phenomenon is empirically  
 153 observed in our experiments with varying  $r$ , as depicted in Fig. 2. Although the exact reason for this  
 154 rapid change remains uncertain, one possible explanation is that the intermediate points between  
 155 two categories reside far from the natural image manifold, thus, challenging the diffusion model’s  
 156 capability to generate them. Ideally, we should select  $r$  corresponding to just after this rapid semantic  
 157 transition, as at this point,  $X_r$  exhibits the highest similarity to the source image while belonging to  
 158 the target category.

159 We propose to trace the semantic trajectory between  $X_S$  and  $X_T$  through the lens of the classifier  
 160  $f_\theta(\cdot)$ . As shown in Algorithm 1, assuming access to the classifier backbone  $f_\theta(\cdot)$  and at least one  
 161 example  $X_T$  from the target category, we convert both  $X_S$  and  $X_T$  into their respective latent vectors  
 162  $Z_S$  and  $Z_T$  by passing them through  $f_\theta(\cdot)$ . Then, we sample  $M$  values for  $r$  uniformly distributed  
 163  $\in (0, 1)$ , generating their corresponding  $X_r$  and their latent vectors  $Z_r$  for all those  $r$ . Subsequently,  
 164 we calculate  $d_r = \frac{(Z_r - Z_S)^T (Z_T - Z_S)}{\|Z_T - Z_S\|_2}$  as the distance between  $Z_r$  and  $Z_S$  projected onto the vector  
 165 connecting  $Z_S$  and  $Z_T$ . Our hypothesis posits that the rapid semantic transition corresponds to a  
 166 sharp change in this projected distance. Therefore, we sample  $n$  values for  $r$  uniformly distributed

---

**Algorithm 1: GeNIe-Ada**

---

**Require:**  $X_S, X_T, f_\theta(\cdot), \text{STDiff}(\cdot), M$   
 Extract  $Z_S \leftarrow f_\theta(X_S), Z_T \leftarrow f_\theta(X_T)$   
**for**  $m \in [1, M]$  **do**  
    $r \leftarrow \frac{m}{M}, Z_r \leftarrow f_\theta(\text{STDiff}(X, P, r))$   
    $d_m \leftarrow \frac{(Z_r - Z_S)^T (Z_T - Z_S)}{\|Z_T - Z_S\|_2}$   
 $m^* \leftarrow \text{argmax}_m |d_m - d_{m-1}|, \forall m \in [2, M]$   
 $r^* \leftarrow \frac{m^*}{M}$   
**Return:**  $X_{r^*} = \text{STDiff}(X_S, P, r^*)$

---

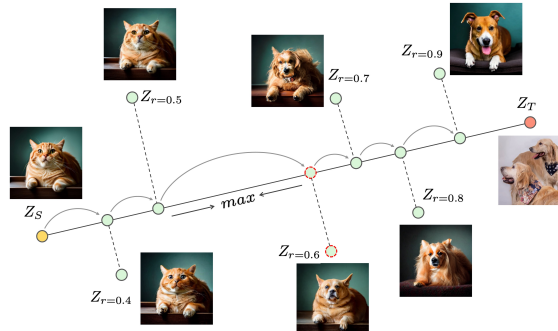


Figure 3: GeNIe-Ada: To choose  $r$  adaptively for each (source image, target category) pair, we propose tracing the semantic trajectory from  $Z_S$  (source image embeddings) to  $Z_T$  (target embeddings) through the lens of the classifier  $f_\theta(\cdot)$  (Algorithm 1). We adaptively select the sample right after the largest semantic shift.

167 between 0 and 1, and analyze the variations in  $d_r$ . We identify the largest gap in  $d_r$  and select the  $r$   
 168 value just after the gap when increasing  $r$ , as detailed in Algorithm 1 and illustrated in Fig. 3.

## 169 4 Experiments

170 Since the impact of augmentation is more pronounced when the training data is limited, we evaluate  
 171 the impact of GeNIe on Few-Shot classification in Section 4.1, Long-Tailed classification in Sec-  
 172 tion 4.2, and fine-grained classification in Section A.2. For GeNIe-Ada in all scenarios, we utilize  
 173 GeNIe to generate augmentations from the noise level set  $\{0.5, 0.6, 0.7, 0.8, 0.9\}$ . The selection of  
 174 the appropriate noise level per source image and target is adaptive, achieved through Algorithm 1.

175 **Baselines.** We use Stable Diffusion 1.5 [81] as our base diffusion model. In all settings,  
 176 we use the same prompt format to generate images for the target class: i.e., “A photo of a  
 177 <target category>”, where we replace the target category with the target category label.  
 178 We generate  $512 \times 512$  images for all methods. For fairness in comparison, we generate the same  
 179 number of new images for each class. We use a single NVIDIA RTX 3090 for image generation.  
 180 We consider 4 diffusion-based baselines and a suite of traditional data augmentation baselines:

181 **Img2Img [63, 67]:** We sample an image from a target class, add noise to its latent representation and  
 182 then pass it along with a prompt for the target category through reverse diffusion. The focus here is  
 183 on a target class for which we generate extra positive samples. Adding large amount of noise leads  
 184 to generating an image less similar to the original image. We use two different noise magnitudes for  
 185 this baseline:  $r = 0.3$  and  $r = 0.7$  and denote them by  $\text{Img2Img}^L$  and  $\text{Img2Img}^H$ , respectively.

186 **Txt2Img [4, 35]:** For this baseline, we omit the forward diffusion process and only use the reverse  
 187 process starting from a text prompt for the target class of interest. This is similar to the base text-  
 188 to-image generation strategy adopted in [81, 35, 89, 4, 62]. Fig. 4 illustrates a set of generated  
 189 augmentation examples for Txt2Img, Img2Img, and GeNIe.

190 **DAFusion [100]:** In this method, an embedding is optimized with a set of images for each class to  
 191 correspond to the classes in the dataset. This approach is introduced in Textual Inversion [32]. We  
 192 optimize an embedding for 5000 iterations for each class in the dataset, followed by augmentation  
 193 similar as the DAFusion method.

194 **Cap2Aug[83]:** It is a recent diffusion-based data augmentation strategy that uses image captions as  
 195 text prompts for an image-to-image diffusion model.

196 **Traditional Data Augmentation:** We consider both weak and strong traditional augmentations.  
 197 More specifically, for weak augmentation we use random resize crop with scaling  $\in [0.2, 1.0]$  and  
 198 horizontal flipping. For strong augmentation, we consider random color jitter, random grayscale,  
 199 and Gaussian blur. For the sake of completeness, we also compare against data augmentations such  
 200 as CutMix [110] and MixUp [111] that combine two images together.

## 201 4.1 Few-shot Classification

202 We assess the impact of GeNIe compared to other augmentations in a number of few-shot classifica-  
 203 tion (FSL) scenarios, where the model has to learn only from the samples contained in the ( $N$ -way,  
 204  $K$ -shot) support set and infer on the query set. Note that this corresponds to an inference-only FSL

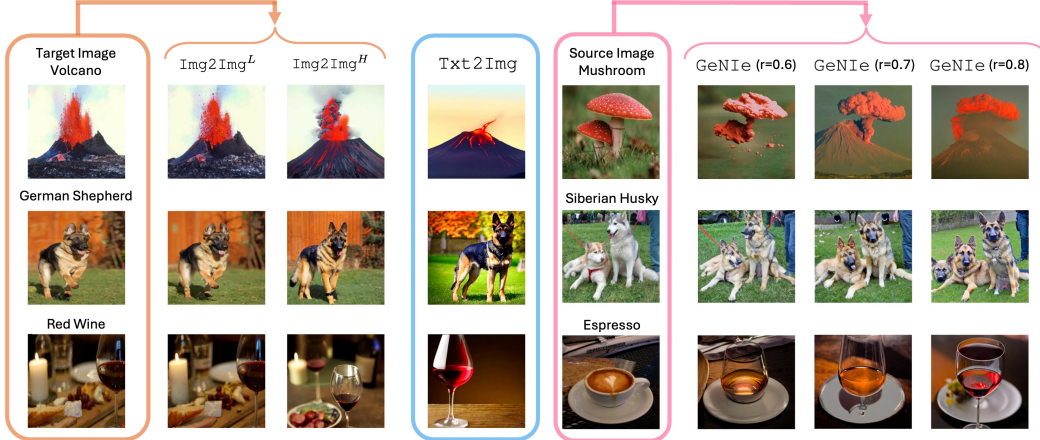


Figure 4: **Visualization of Generative Samples:** We compare GeNIe with two baselines:  $\text{Img2Img}^L$  **augmentation:** both image and text prompt are from the same category. Adding noise does not change the image much, so they are not hard examples.  $\text{Txt2Img}$  **augmentation:** We simply use the text prompt only to generate an image for the desired category (e.g., using a text2image method). Such images may be far from the domain of our task since the generation is not informed by any visual data from our task. **GeNIe augmentation:** We use the target category name in the text prompt only along with the source image.

205 setting where a pretraining stage on an abundant dataset is discarded. The goal is to assess how well  
 206 the model can benefit from the augmentations while keeping the original  $N \times K$  samples intact.

207 **Datasets.** We conduct our few-shot experiments on two most commonly adopted few-shot classi-  
 208 fication datasets: *mini*-Imagenet [79] and *tiered*-Imagenet [80]. *mini*-Imagenet is a subset of Ima-  
 209 geNet [22] for few-shot classification. It contains 100 classes with 600 samples each. We follow  
 210 the predominantly adopted settings of [79, 10] where we split the entire dataset into 64 classes for  
 211 training, 16 for validation and 20 for testing. *tiered*-Imagenet is a larger subset of ImageNet with  
 212 608 classes and a total of 779, 165 images, which are grouped into 34 higher-level nodes in the *Im-*  
 213 *ageNet* human-curated hierarchy. This set of nodes is partitioned into 20, 6, and 8 disjoint sets of  
 214 training, validation, and testing nodes, and the corresponding classes form the respective meta-sets.

215 **Evaluation.** To quantify the impact of different augmentation methods, we evaluate the test-set ac-  
 216 curacies of a state-of-the-art unsupervised few-shot learning method with GeNIe and compare them  
 217 against the accuracies obtained using other augmentation methods. Specifically, we use UniSiam  
 218 [61] pre-trained with ResNet-18, ResNet-34 and ResNet-50 backbones and follow its evaluation  
 219 strategy of fine-tuning a logistic regressor to perform ( $N$ -way,  $K$ -shot) classification on the test sets  
 220 of *mini*- and *tiered*-Imagenet. Following [79], an episode consists of a labeled support-set and an un-  
 221 labelled query-set. The support-set contains  $N$  randomly sampled classes where each class contains  
 222  $K$  samples, whereas the query-set contains  $Q$  randomly sampled unlabeled images per class. We  
 223 conduct our experiments on the two most commonly adopted settings: (5-way, 1-shot) and (5-way,  
 224 5-shot) classification settings. Following the literature, we sample 16-shots per class for the query  
 225 set in both settings. We report the test accuracies along with the 95% confidence interval over 600  
 226 and 1000 episodes for *mini*-ImageNet and *tiered*-ImageNet, respectively.

227 **Implementation Details:** GeNIe generates augmented images for each class using images from all  
 228 other classes as the source image. We use  $r = 0.8$  in our experiments. We generate 4 samples per  
 229 class as augmentations in the 5-way, 1-shot setting and 20 samples per class as augmentations in the  
 230 5-way, 5-shot setting. For the sake of a fair comparison, we ensure that the total number of labelled  
 231 samples in the support set after augmentation remains the same across all different traditional and  
 232 generative augmentation methodologies. Due to the expensive training of embeddings for each class  
 233 in each episode, we only evaluated the DA-Fusion baseline on the first 100 episodes.

234 **Results:** The results on *mini*-Imagenet and *tiered*-Imagenet for both (5-way, 1 and 5-shot) set-  
 235 tings are summarized in Table 1 and Table 2, respectively. Regardless of the choice of back-  
 236 bone, we observe that GeNIe helps consistently improve UniSiam’s performance and outperform  
 237 other supervised and unsupervised few-shot classification methods as well as other diffusion-based  
 238 [100, 63, 82, 35] and classical [110, 111] data augmentation techniques on both datasets, across both  
 239 (5-way, 1 and 5-shot) settings. Our noise adaptive method of selecting optimal augmentations per  
 240 source image (GeNIe-Ada) further improves GeNIe’s performance across all three backbones, both

Table 1: *mini-ImageNet*: We use our augmentations on (5-way, 1-shot) and (5-way, 5-shot) few-shot settings of mini-Imagenet dataset with 3 different backbones (ResNet-18, 34, and 50). We compare with various baselines and show that our augmentations with UniSiam outperform all the baselines including Txt2Img and DAFusion augmentation. The number of generated images per class is 4 for 1-shot and 20 for 5-shot settings.

ResNet-18					ResNet-34				
Augmentation	Method	Pre-training	1-shot	5-shot	Augmentation	Method	Pre-training	1-shot	5-shot
-	iDeMe-Net [14]	sup.	59.1±0.9	74.6±0.7	Weak	Baseline [10]	sup.	49.8±0.7	73.5±0.7
-	Robust + dist [27]	sup.	63.7±0.6	81.2±0.4	Weak	Baseline++ [10]	sup.	52.7±0.8	76.2±0.6
-	AFHN [54]	sup.	62.4±0.7	78.2±0.6	Weak	SimCLR [9]	unsup.	64.0±0.4	79.8±0.3
Weak	ProtoNet+SSL [94]	sup.+ssl	-	76.6	Weak	SimSiam [12]	unsup.	63.8±0.4	80.4±0.3
Weak	Neg-Cosine [57]	sup.	62.3±0.8	80.9±0.6	Weak	UniSiam+dist [61]	unsup.	<b>65.6±0.4</b>	<b>83.4±0.2</b>
-	Centroid Align[1]	sup.	59.9±0.7	80.4±0.7	Weak	UniSiam [61]	unsup.	64.3±0.8	82.3±0.5
-	Baseline [10]	sup.	59.6±0.8	77.3±0.6	Strong	UniSiam [61]	unsup.	64.5±0.8	82.1±0.6
-	Baseline++ [10]	sup.	59.0±0.8	76.7±0.6	CutMix [110]	UniSiam [61]	unsup.	64.0±0.8	81.7±0.6
Weak	PSST [13]	sup.+ssl	59.5±0.5	77.4±0.5	MixUp [111]	UniSiam [61]	unsup.	63.7±0.8	80.1±0.8
Weak	UMTRA [46]	unsup.	43.1±0.4	53.4±0.3	Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	65.5±0.8	82.9±0.5
Weak	ProtoCLR [66]	unsup.	50.9±0.4	71.6±0.3	Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	70.5±0.8	84.8±0.5
Weak	SimCLR [9]	unsup.	62.6±0.4	79.7±0.3	Txt2Img[4, 35]	UniSiam [61]	unsup.	75.4±0.6	85.5±0.5
Weak	SimSiam [12]	unsup.	62.8±0.4	79.9±0.3	DAFusion [100]	UniSiam [61]	unsup.	64.7±1.9	83.2±1.4
Weak	UniSiam+dist [61]	unsup.	<b>64.1±0.4</b>	<b>82.3±0.3</b>	GeNIe (Ours)	UniSiam [61]	unsup.	<b>77.1±0.6</b>	<b>86.3±0.4</b>
Weak	UniSiam [61]	unsup.	63.1±0.8	81.4±0.5	GeNIe-Ada (Ours)	UniSiam [61]	unsup.	<b>78.5±0.6</b>	<b>86.6±0.4</b>
Strong	UniSiam [61]	unsup.	62.8±0.8	81.2±0.6	<b>ResNet-50</b>				
CutMix [110]	UniSiam [61]	unsup.	62.7±0.8	80.6±0.6	Weak	PDA+Net [11]	unsup.	63.8±0.9	83.1±0.6
MixUp [111]	UniSiam [61]	unsup.	62.1±0.8	80.7±0.6	Weak	Meta-DM [40]	unsup.	66.7±0.4	85.3±0.2
Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	63.9±0.8	82.1±0.5	Weak	UniSiam [61]	unsup.	64.6±0.8	83.4±0.5
Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	69.1±0.7	84.0±0.5	Strong	UniSiam [61]	unsup.	64.8±0.8	83.2±0.5
Txt2Img[4, 35]	UniSiam [61]	unsup.	74.1±0.6	84.6±0.5	CutMix [110]	UniSiam [61]	unsup.	64.3±0.8	83.2±0.5
DAFusion [100]	UniSiam [61]	unsup.	64.3±1.8	82.0±1.4	MixUp [111]	UniSiam [61]	unsup.	63.8±0.8	84.6±0.5
GeNIe (Ours)	UniSiam [61]	unsup.	<b>75.5±0.6</b>	<b>85.4±0.4</b>	Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	66.0±0.8	84.0±0.5
GeNIe-Ada (Ours)	UniSiam [61]	unsup.	<b>76.8±0.6</b>	<b>85.9±0.4</b>	Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	71.1±0.7	85.7±0.5
					Txt2Img[4, 35]	UniSiam [61]	unsup.	76.4±0.6	86.5±0.4
					DAFusion [100]	UniSiam [61]	unsup.	65.7±1.8	83.9±1.2
					GeNIe (Ours)	UniSiam [61]	unsup.	<b>77.3±0.6</b>	<b>87.2±0.4</b>
					GeNIe-Ada (Ours)	UniSiam [61]	unsup.	<b>78.6±0.6</b>	<b>87.9±0.4</b>

241 few-shot settings, and both datasets (*mini* and *tiered-Imagenet*). Few-shot accuracies for ResNet-  
 242 34 computed on *tiered-Imagenet* are reported in Section A.3 of the appendix. Note that employing  
 243 CutMix and MixUp seems to lead to performance degradation compared to weak augmentations,  
 244 probably due to overfitting since these methods can only choose from 4 other classes to mix.

## 245 4.2 Long-Tailed Classification

246 We evaluate our method on long-tailed data, where the number of instances per class is unbalanced,  
 247 with most categories having limited samples (tail). Our goal is to mitigate this bias by augmenting  
 248 the tail of the distribution with generated samples. We evaluate GeNIe using two different backbones  
 249 and methods: the ViT architecture with LViT [107], and ResNet50 with VL-LTR [97].

250 Following LViT [107], we first train an MAE [34] and ViT on the unbalanced dataset without any  
 251 augmentation. Next, we train the Balanced Fine-Tuning stage of LViT by incorporating the aug-  
 252 mentation data generated using GeNIe or other baselines. For ResNet50, we use VL-LTR code to  
 253 fine-tune the CLIP [76] ResNet50 pretrained backbone with generated augmentations by GeNIe.

254 **Dataset:** We perform experiments on ImageNet-LT [60]. It contains 115.8K images from 1,000  
 255 categories. The number of images per class varies from 1280 to 5. Imagenet-LT classes can be  
 256 divided into 3 groups: “Few” with less than 20 images, “Med” with 20 – 100 images, and “Many”  
 257 with more than 100 images. Imagenet-LT uses the same validation set as ImageNet. We augment  
 258 “Few” categories only and limit the number of generated images to 50 samples per class. For GeNIe,  
 259 instead of randomly sampling the source images from other classes, we use a confusion matrix on  
 260 the training data to find the top-4 most confused classes and only consider those classes for random  
 261 sampling of the source image. The source category may be from “Many”, “Med”, or “Few sets”.

262 **Results:** Augmenting training data with GeNIe-Ada improves accuracy on the “Few” set by 11.7%  
 263 and 4.4% compared with LViT only and LViT with Txt2Img augmentation baselines respectively.  
 264 In ResNet50, GeNIe-Ada outperforms Cap2Aug baseline in “Few” categories by 7.6%. The results  
 265 are summarized in Table 3. Please refer to Section A.4 for implementation details.

## 266 4.3 Ablation and Analysis

267 **Semantic Shift from Source to Target Class.** The core motivation behind GeNIe-Ada is that by  
 268 varying the noise ratio  $r$  from 0 to 1, augmented sample  $X_r$  will progressively shift its semantic cat-  
 269 egory from source ( $S$ ) in the beginning to target category ( $T$ ) towards the end. However, somewhere  
 270 between 0 and 1,  $X_r$  will undergo a rapid transition from  $S$  to  $T$ . To demonstrate this hypothesis  
 271 empirically, in Figs. 5 and A5, we visualize pairs of source images and target categories with their re-  
 272 spective GeNIe generated augmentations for different noise ratios  $r$ , along with their corresponding

Table 2: **tiered-ImageNet**: Accuracies (%  $\pm$  std) for 5-way, 1-shot and 5-way, 5-shot classification settings on the test-set. We compare against various SOTA supervised and unsupervised few-shot classification baselines as well as other augmentation methods, with UniSiam [61] pre-trained ResNet-18,50 backbones.

ResNet-18				
Augmentation	Method	Pre-training	1-shot	5-shot
Weak	SimCLR[9]	unsup.	63.4 $\pm$ 0.4	79.2 $\pm$ 0.3
Weak	SimSiam [12]	unsup.	64.1 $\pm$ 0.4	81.4 $\pm$ 0.3
Weak	UniSiam [61]	unsup.	63.1 $\pm$ 0.7	81.0 $\pm$ 0.5
Strong	UniSiam [61]	unsup.	62.8 $\pm$ 0.7	80.9 $\pm$ 0.5
CutMix [110]	UniSiam [61]	unsup.	62.1 $\pm$ 0.7	78.9 $\pm$ 0.6
MixUp [111]	UniSiam [61]	unsup.	62.1 $\pm$ 0.7	78.4 $\pm$ 0.6
Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	63.9 $\pm$ 0.7	81.8 $\pm$ 0.5
Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	68.7 $\pm$ 0.7	83.5 $\pm$ 0.5
Txt2Img[35]	UniSiam [61]	unsup.	72.9 $\pm$ 0.6	84.2 $\pm$ 0.5
DAFusion [100]	UniSiam [61]	unsup.	62.6 $\pm$ 2.1	81.0 $\pm$ 1.5
GeNIe(Ours)	UniSiam [61]	unsup.	<b>73.6<math>\pm</math>0.6</b>	<b>85.0<math>\pm</math>0.4</b>
GeNIe-Ada(Ours)	UniSiam [61]	unsup.	<b>75.1<math>\pm</math>0.6</b>	<b>85.5<math>\pm</math>0.5</b>

ResNet-50				
Augmentation	Method	Pre-training	1-shot	5-shot
Weak	PDA+Net [11]	unsup.	69.0 $\pm$ 0.9	84.2 $\pm$ 0.7
Weak	Meta-DM [40]	unsup.	69.6 $\pm$ 0.4	86.5 $\pm$ 0.3
Weak	UniSiam + dist [61]	unsup.	69.6 $\pm$ 0.4	86.5 $\pm$ 0.3
Weak	UniSiam [61]	unsup.	66.8 $\pm$ 0.7	84.7 $\pm$ 0.5
Strong	UniSiam [61]	unsup.	66.5 $\pm$ 0.7	84.5 $\pm$ 0.5
CutMix [110]	UniSiam [61]	unsup.	66.0 $\pm$ 0.7	83.3 $\pm$ 0.5
MixUp [111]	UniSiam [61]	unsup.	66.1 $\pm$ 0.5	84.1 $\pm$ 0.8
Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	67.8 $\pm$ 0.7	85.3 $\pm$ 0.5
Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	72.4 $\pm$ 0.7	86.7 $\pm$ 0.4
Txt2Img[35]	UniSiam [61]	unsup.	77.1 $\pm$ 0.6	87.3 $\pm$ 0.4
DAFusion [100]	UniSiam [61]	unsup.	66.5 $\pm$ 2.2	84.8 $\pm$ 1.4
GeNIe (Ours)	UniSiam [61]	unsup.	<b>78.0<math>\pm</math>0.6</b>	<b>88.0<math>\pm</math>0.4</b>
GeNIe-Ada (Ours)	UniSiam [61]	unsup.	<b>78.8<math>\pm</math>0.6</b>	<b>88.6<math>\pm</math>0.6</b>

Table 3: **Long-Tailed ImageNet-LT**: We compare different augmentation methods on ImageNet-LT and report Top-1 accuracy for “Few”, “Medium”, and “Many” sets. On the “Few” set and LiVT method, our augmentations improve the accuracy by 11.7 points compared to LiVT original augmentation and 4.4 points compared to Txt2Img. GeNIe-Ada outperforms Cap2Aug baseline in “Few” categories by 7.6%. Refer to Table A4 for a full comparison with prior Long-Tailed methods.

ResNet-50				
Method	Many	Med.	Few	Overall Acc
ResLT [18]	63.3	53.3	40.3	55.1
PaCo [19]	68.2	58.7	41.0	60.0
LWS [44]	62.2	48.6	31.8	51.5
Zero-shot CLIP [76]	60.8	59.3	58.6	59.8
DRO-LT [85]	64.0	49.8	33.1	53.5
VL-LTR [97]	77.8	67.0	50.8	70.1
Cap2Aug [83]	78.5	<b>67.7</b>	51.9	70.9
GeNIe-Ada	<b>79.2</b>	64.6	<b>59.5</b>	<b>71.5</b>

ViT-B				
Method	Many	Med.	Few	Overall Acc
ViT [24]	50.5	23.5	6.9	31.6
MAE [33]	74.7	48.2	19.4	54.5
DeiT [99]	70.4	40.9	12.8	48.4
LiVT [107]	73.6	56.4	41.0	60.9
LiVT + Img2Img <sup>L</sup>	74.3	56.4	34.3	60.5
LiVT + Img2Img <sup>H</sup>	73.8	56.4	45.3	61.6
LiVT + Txt2Img	<b>74.9</b>	55.6	48.3	62.2
LiVT + GeNIe-Ada	74.0	<b>56.9</b>	<b>52.7</b>	<b>63.1</b>

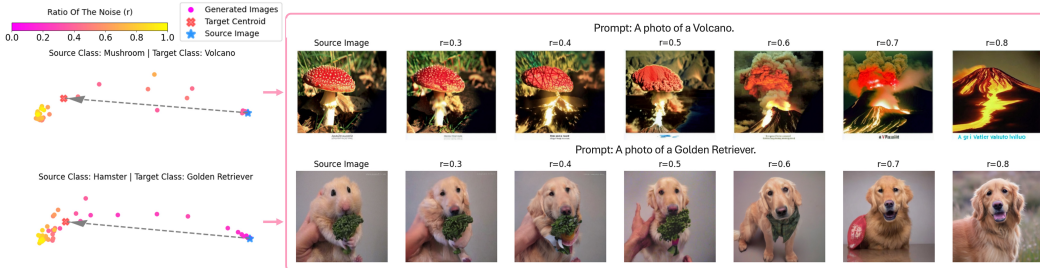


Figure 5: **Embedding visualizations of generative augmentations**: We pass all generative augmentations through DINOv2 ViT-G (serving as an oracle) to extract their corresponding embeddings and visualize them with PCA. As shown, the extent of semantic shifts varies based on both the source image and the target class.

273 PCA-projected embedding scatter plots (on the far left). We extract embeddings for all the images  
 274 using a DINOv2 ViT-G pretrained backbone, which we assume as an oracle model in identifying  
 275 the right category. We observe that as  $r$  increases from 0.3 to 0.8, the images transition to embody  
 276 more of the target category’s semantics while preserving the contextual features of the source image.  
 277 This transition of semantics can also be observed in the embedding plots (on the left) where they  
 278 consistently shift from the proximity of the source image (blue star) to the target class’s centroid  
 279 (red cross) as the noise ratio  $r$  increases. The sparse distribution of points within  $r = [0.4, 0.6]$  for  
 280 the first image and  $r = [0.2, 0.4]$  for the second image aligns with our intuition of a rapid transition  
 281 from category  $S$  to  $T$ , thus empirically affirming our motivation behind GeNIe-Ada.

282 To further establish this, in Fig. 6, we demonstrate the efficacy of GeNIe in generating hard negatives  
 283 at the decision boundaries of an SVM classifier, which is trained on the labelled support set of  
 284 the few-shot tasks of *mini-Imagenet*, without any augmentations. We then plot source and target  
 285 class probabilities ( $P(Y_S|X_r)$  and  $P(Y_T|X_r)$ , respectively) of the generated augmentation samples  
 286  $X_r$ . For both  $r = 0.6$  and  $0.7$ , there is significant overlap between  $P(Y_S|X_r)$  and  $P(Y_T|X_r)$ ,  
 287 making it difficult for the classifier to decide the correct class. On the right-hand-side, GeNIe-Ada  
 288 automatically selects the best  $r$  resulting in the most overlap between the two distributions, thus  
 289 offering the hardest negative sample among the considered  $r$  values (for more details see A.1).  
 290 Note that a large overlap between distributions is not sufficient to call the generated samples hard  
 291 negatives because they should also belong to the target category. This is, however, confirmed by the  
 292 high Oracle accuracy in Table 4 (elaborated in detail in the following paragraph) which verifies that  
 293 majority of the generated augmentation samples do belong to the target category.



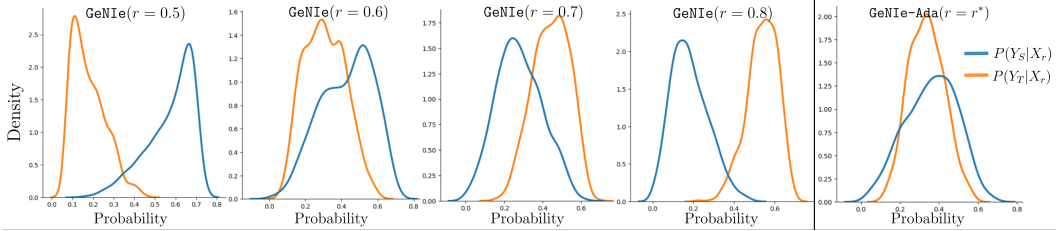


Figure 6: **Why GeNIe augmentations are challenging?** While deciding which class the generated augmentations ( $X_r$ ) belong to is already difficult within  $r = [0.6, 0.7]$  (due to high overlap between  $P(Y_S|X_r)$  and  $P(Y_T|X_r)$ ), GeNIe-Ada selects the best noise threshold ( $r^*$ ) offering the hardest negative sample.

Table 4: **Effect of Noise in GeNIe:** We use the same setting as in Table 1 to study the effect of the amount of noise. As expected (also shown in Fig 5), small noise results in worse accuracy since some generated images may be from the source category rather than the target one. For  $r = 0.5$  only 73% of the generated data is from the target category. This behaviour is also shown in Fig. 2. Notably, reducing the noise level below 0.7 is associated with a decline in oracle accuracy and subsequent degradation in the performance of the final few-shot model. Note that the high oracle accuracy of GeNIe-Ada demonstrates its capability to adaptively select the noise level per source and target, ensuring semantic consistency with the intended target.

Noise	ResNet-18		ResNet-34		ResNet-50		Oracle Acc
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	
GeNIe( $r=0.5$ )	60.42±0.8	74.11±0.6	62.02±0.8	75.80±0.6	63.65±0.9	77.61±0.6	73.4±0.5
GeNIe( $r=0.6$ )	69.66±0.7	80.65±0.5	71.13±0.7	82.21±0.5	72.10±0.7	82.79±0.5	85.8±0.4
GeNIe( $r=0.7$ )	74.50±0.6	83.26±0.5	76.41±0.6	84.44±0.5	77.05±0.6	84.95±0.4	94.5±0.2
GeNIe( $r=0.8$ )	75.45±0.6	85.38±0.4	77.08±0.6	86.28±0.4	77.28±0.6	87.22±0.4	98.2±0.1
GeNIe( $r=0.9$ )	74.96±0.6	85.29±0.4	77.63±0.6	86.17±0.4	77.73±0.6	87.00±0.4	99.3±0.1
GeNIe-Ada	76.79±0.6	85.89±0.4	78.49±0.6	86.55±0.4	78.64±0.6	87.88±0.4	98.9±0.2

294 **Label consistency of the generated samples.** The choice of noise ratio  $r$  is important in producing  
 295 hard negative examples. In Table 4, we present the accuracy of the GeNIe model across various noise  
 296 ratios, alongside the oracle accuracy, which is an ImageNet pre-trained DeiT-Base [98] classifier.  
 297 We observe a decline in the label consistency of generated data (quantified by the performance of  
 298 the oracle model) when decreasing the noise level. Reducing  $r$  also results in a degradation in the  
 299 performance of the final few-shot model (87.2%  $\rightarrow$  77.6%) corroborating that an appropriate choice  
 300 of  $r$  plays a crucial role in our design strategy. We investigate this further in the following paragraph.

301 **Effect of Noise in GeNIe.** We examine the impact of noise on the performance of the few-shot  
 302 model in Table 4. Noise levels  $r \in [0.7, 0.8]$  yield the best performance. Conversely, utilizing noise  
 303 levels below 0.7 diminishes performance due to label inconsistency, as is demonstrated in Table 4  
 304 and Fig 5. As such, determining the appropriate noise level is pivotal for the performance of GeNIe  
 305 to be able to generate challenging hard negatives while maintaining label consistency. An alternative  
 306 approach to finding the optimal noise level involves using GeNIe-Ada to adaptively select the noise  
 307 level for each source image and target class. As demonstrated in Tables 4 and AI, GeNIe-Ada  
 308 achieves performance that is comparable to or surpasses that of GeNIe with fixed noise levels.

## 309 5 Concluding Remarks

310 GeNIe, for the first time to our knowledge, combines contradictory sources of information (a source  
 311 image, and a different target category prompt) through a noise adjustment strategy into a conditional  
 312 latent diffusion model to generate challenging augmentations, which can serve as hard negatives.

313 **Limitation.** The required time to create augmentations through GeNIe is on par with any typical  
 314 diffusion-based competitors [4, 35]; however, this is naturally slower than traditional augmentation  
 315 techniques [110, 111]. This is not a bottleneck in offline augmentation strategies, but can be con-  
 316 sidered a limiting factor in real-time scenarios. Recent studies are already mitigating this through  
 317 advancements in diffusion model efficiency [87, 68, 58]. Another challenge present in any genera-  
 318 tive AI-based augmentation technique is the domain shift between the distribution of training data  
 319 and the downstream context they might be used for augmentation. A possible remedy is to fine-tune  
 320 the diffusion backbone on a rather small dataset from the downstream task.

321 **Broader Impact.** We believe ideas from GeNIe can have a significant impact when it comes to gen-  
 322 erating hard augmentations challenging and thus enhancing downstream tasks beyond classification.  
 323 At the same time, just like any other generative model, GeNIe can also introduce inherent biases  
 324 stemming from the training data used to build its diffusion backbone, which can reflect and amplify  
 325 societal prejudices or inaccuracies. Therefore, it is crucial to carefully mitigate potential biases in  
 326 generative models such as GeNIe to ensure a fair and ethical deployment of deep learning systems.

327 **References**

- 328 [1] Afrasiyabi, A., Lalonde, J.F., Gagné, C.: Associative alignment for few-shot image classifi-  
329 cation. In: ECCV (2019)
- 330 [2] Alayrac, J.B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A.,  
331 Millican, K., Reynolds, M., Ring, R., Rutherford, E., Cabi, S., Han, T., Gong, Z., Saman-  
332 gooei, S., Monteiro, M., Menick, J., Borgeaud, S., Brock, A., Nematzadeh, A., Sharifzadeh,  
333 S., Binkowski, M., Barreira, R., Vinyals, O., Zisserman, A., Simonyan, K.: Flamingo: a  
334 visual language model for few-shot learning (2022)
- 335 [3] Antoniou, A., Storkey, A.: Assume, augment and learn: Unsupervised few-shot meta-  
336 learning via random labels and data augmentation. arxiv:1902.09884 (2019)
- 337 [4] Azizi, S., Kornblith, S., Saharia, C., Norouzi, M., Fleet, D.J.: Synthetic data from diffusion  
338 models improves imagenet classification (2023)
- 339 [5] Bossard, L., Guillaumin, M., Van Gool, L.: Food-101 – mining discriminative components  
340 with random forests. In: European Conference on Computer Vision (2014)
- 341 [6] Cai, J., Wang, Y., Hwang, J.N., et al.: Ace: Ally complementary experts for solving long-  
342 tailed recognition in one-shot. In: ICCV, pp. 112–121 (2021)
- 343 [7] Cao, K., Wei, C., Gaidon, A., Arechiga, N., Ma, T.: Learning imbalanced datasets with label-  
344 distribution-aware margin loss. NeurIPS **32** (2019)
- 345 [8] Chegini, A., Feizi, S.: Identifying and mitigating model failures through few-shot clip-aided  
346 diffusion generation. arXiv preprint arXiv:2312.05464 (2023)
- 347 [9] Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning  
348 of visual representations. In: ICML (2020)
- 349 [10] Chen, W.Y., Liu, Y.C., Kira, Z., Wang, Y.C.F., Huang, J.B.: A closer look at few-shot classi-  
350 fication. In: ICLR (2019)
- 351 [11] Chen, W., Si, C., Wang, W., Wang, L., Wang, Z., Tan, T.: Few-shot learning with part discov-  
352 ery and augmentation from unlabeled images. arXiv preprint arXiv:2105.11874 (2021)
- 353 [12] Chen, X., He, K.: Exploring simple siamese representation learning. In: CVPR (2021)
- 354 [13] Chen, Z., Ge, J., Zhan, H., Huang, S., Wang, D.: Pareto self-supervised training for few-shot  
355 learning. In: CVPR (2021)
- 356 [14] Chen, Z., Fu, Y., Wang, Y.X., Ma, L., Liu, W., Hebert, M.: Image deformation meta-networks  
357 for one-shot learning. In: CVPR (2019)
- 358 [15] Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmen-  
359 tation policies from data (2019)
- 360 [16] Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.V.: Randaugment: Practical automated data aug-  
361 mentation with a reduced search space (2019)
- 362 [17] Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.: Randaugment: Practical automated data aug-  
363 mentation with a reduced search space. In: Larochelle, H., Ranzato, M., Hadsell, R.,  
364 Balcan, M., Lin, H. (eds.) Advances in Neural Information Processing Systems. vol. 33,  
365 pp. 18613–18624. Curran Associates, Inc. (2020), [https://proceedings.neurips.cc/  
366 paper/2020/file/d85b63ef0ccb114d0a3bb7b7d808028f-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/d85b63ef0ccb114d0a3bb7b7d808028f-Paper.pdf)
- 367 [18] Cui, J., Liu, S., Tian, Z., Zhong, Z., Jia, J.: Reslt: Residual learning for long-tailed recogni-  
368 tion. IEEE transactions on pattern analysis and machine intelligence **45**(3), 3695–3706 (2022)
- 369 [19] Cui, J., Zhong, Z., Liu, S., Yu, B., Jia, J.: Parametric contrastive learning. In: Proceedings of  
370 the IEEE/CVF international conference on computer vision. pp. 715–724 (2021)
- 371 [20] Cui, J., Zhong, Z., Liu, S., Yu, B., Jia, J.: Parametric contrastive learning. In: ICCV. pp.  
372 715–724 (2021)
- 373 [21] Cui, Y., Jia, M., Lin, T.Y., Song, Y., Belongie, S.: Class-balanced loss based on effective  
374 number of samples. In: CVPR. pp. 9268–9277 (2019)
- 375 [22] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierar-  
376 chical image database. In: 2009 IEEE conference on computer vision and pattern recognition.  
377 pp. 248–255. Ieee (2009)

- 378 [23] Ding, M., An, B., Xu, Y., Satheesh, A., Huang, F.: SAFLEX: Self-adaptive augmentation via  
379 feature label extrapolation. In: The Twelfth International Conference on Learning Representations  
380 (2024), <https://openreview.net/forum?id=qL6brrBDk2>
- 381 [24] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., De-  
382 hghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is  
383 worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021)
- 384 [25] Dunlap, L., Mohri, C., Zhang, H., Guillory, D., Darrell, T., Gonzalez, J.E., Rohrbach, A.,  
385 Raghunathan, A.: Using language to extend to unseen domains. International Conference on  
386 Learning Representations (ICLR) (2023)
- 387 [26] Dunlap, L., Umino, A., Zhang, H., Yang, J., Gonzalez, J.E., Darrell, T.: Diversify your vision  
388 datasets with automatic diffusion-based augmentation (2023)
- 389 [27] Dvornik, N., Mairal, J., Schmid, C.: Diversity with cooperation: Ensemble methods for few-  
390 shot classification. In: ICCV (2019)
- 391 [28] Feng, C.M., Yu, K., Liu, Y., Khan, S., Zuo, W.: Diverse data augmentation with diffusions  
392 for effective test-time prompt tuning (2023)
- 393 [29] Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep  
394 networks. In: Proceedings of the 34th International Conference on Machine Learning. pp.  
395 1126–1135 (2017)
- 396 [30] Frid-Adar, M., Diamant, I., Klang, E., Amitai, M., Goldberger, J., Greenspan, H.: Gan-  
397 based synthetic medical image augmentation for increased cnn performance in liver lesion  
398 classification. Neurocomputing (2018)
- 399 [31] Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.:  
400 An image is worth one word: Personalizing text-to-image generation using textual inversion.  
401 arXiv preprint arXiv:2208.01618 (2022)
- 402 [32] Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or,  
403 D.: An image is worth one word: Personalizing text-to-image generation using textual  
404 inversion (2022). <https://doi.org/10.48550/ARXIV.2208.01618>, [https://arxiv.org/abs/](https://arxiv.org/abs/2208.01618)  
405 [2208.01618](https://arxiv.org/abs/2208.01618)
- 406 [33] He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.B.: Masked autoencoders are scalable  
407 vision learners. In: CVPR. pp. 15979–15988. IEEE (2022)
- 408 [34] He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable  
409 vision learners (2021)
- 410 [35] He, R., Sun, S., Yu, X., Xue, C., Zhang, W., Torr, P., Bai, S., Qi, X.: Is synthetic data from  
411 generative models ready for image recognition? arXiv preprint arXiv:2210.07574 (2022)
- 412 [36] Hemmat, R.A., Pezeshki, M., Bordes, F., Drozdal, M., Romero-Soriano, A.: Feedback-  
413 guided data synthesis for imbalanced classification (2023)
- 414 [37] Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., Lakshminarayanan, B.: AugMix:  
415 A simple data processing method to improve robustness and uncertainty. Proceedings of the  
416 International Conference on Learning Representations (ICLR) (2020)
- 417 [38] Hong, Y., Zhang, J., Sun, Z., Yan, K.: Safa: Sample-adaptive feature augmentation for long-  
418 tailed image classification. In: ECCV (2022)
- 419 [39] Hsu, K., Levine, S., Finn, C.: Unsupervised learning via meta-learning. In: ICLR (2018)
- 420 [40] Hu, W., Jiang, X., Liu, J., Yang, Y., Tian, H.: Meta-dm: Applications of diffusion models on  
421 few-shot learning (2023)
- 422 [41] Huang, S.W., Lin, C.T., Chen, S.P., an Po-Hao Hsu, Y.Y.W., Lai, S.H.: Auggan: Cross do-  
423 main adaptation with gan-based data augmentation. European Conference on Computer Vision  
424 (2018)
- 425 [42] Jain, S., Lawrence, H., Moitra, A., Madry, A.: Distilling model failures as directions in latent  
426 space. In: ArXiv preprint arXiv:2206.14754 (2022)
- 427 [43] Jang, H., Lee, H., Shin, J.: Unsupervised meta-learning via few-shot pseudo-supervised con-  
428 trastive learning. In: The Eleventh International Conference on Learning Representations  
429 (2022)

- 430 [44] Kang, B., Xie, S., Rohrbach, M., Yan, Z., Gordo, A., Feng, J., Kalantidis, Y.: Decoupling rep-  
431 resentation and classifier for long-tailed recognition. arXiv preprint arXiv:1910.09217 (2019)
- 432 [45] Kang, B., Xie, S., Rohrbach, M., Yan, Z., Gordo, A., Feng, J., Kalantidis, Y.: Decoupling  
433 representation and classifier for long-tailed recognition. In: ICLR (2020)
- 434 [46] Khodadadeh, S., Boloni, L., Shah, M.: Unsupervised meta-learning for few-shot image clas-  
435 sification. In: NeurIPS (2019)
- 436 [47] Kim, J.H., Choo, W., Song, H.O.: Puzzle mix: Exploiting saliency and local statistics for  
437 optimal mixup. In: International Conference on Machine Learning. pp. 5275–5285. PMLR  
438 (2020)
- 439 [48] Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3D object representations for fine-grained catego-  
440 rization. In: Workshop on 3D Representation and Recognition. Sydney, Australia (2013)
- 441 [49] Li, A.C., Prabhudesai, M., Duggal, S., Brown, E., Pathak, D.: Your diffusion model is secretly  
442 a zero-shot classifier (2023)
- 443 [50] Li, B., Han, Z., Li, H., Fu, H., Zhang, C.: Trustworthy long-tailed classification. In: CVPR.  
444 pp. 6970–6979 (2022)
- 445 [51] Li, D., Ling, H., Kim, S.W., Kreis, K., Barriuso, A., Fidler, S., Torralba, A.: Bigdatasetgan:  
446 Synthesizing imagenet with pixel-wise annotations (2022)
- 447 [52] Li, J., Tan, Z., Wan, J., Lei, Z., Guo, G.: Nested collaborative learning for long-tailed visual  
448 recognition. In: CVPR. pp. 6949–6958 (2022)
- 449 [53] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified  
450 vision-language understanding and generation (2022)
- 451 [54] Li, K., Zhang, Y., Li, K., Fu, Y.: Adversarial feature hallucination networks for few-shot  
452 learning. In: CVPR (2020)
- 453 [55] Li, M., Cheung, Y.m., Lu, Y., et al.: Long-tailed visual recognition via gaussian clouded logit  
454 adjustment. In: CVPR. pp. 6929–6938 (2022)
- 455 [56] Li, T., Cao, P., Yuan, Y., Fan, L., Yang, Y., Feris, R.S., Indyk, P., Katabi, D.: Targeted  
456 supervised contrastive learning for long-tailed recognition. In: CVPR. pp. 6918–6928 (2022)
- 457 [57] Liu, B., Cao, Y., Lin, Y., Li, Q., Zhang, Z., Long, M., Hu, H.: Negative margin matters:  
458 Understanding margin in few-shot classification. In: ECCV (2020)
- 459 [58] Liu, X., Zhang, X., Ma, J., Peng, J., et al.: InstafLOW: One step is enough for high-quality  
460 diffusion-based text-to-image generation. In: The Twelfth International Conference on Learn-  
461 ing Representations (2023)
- 462 [59] Liu, Z., Li, S., Wu, D., Liu, Z., Chen, Z., Wu, L., Li, S.Z.: Automix: Unveiling the power of  
463 mixup for stronger classifiers. In: Computer Vision—ECCV 2022: 17th European Conference,  
464 Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIV. pp. 441–458. Springer (2022)
- 465 [60] Liu, Z., Miao, Z., Zhan, X., Wang, J., Gong, B., Yu, S.X.: Large-scale long-tailed recognition  
466 in an open world. In: CVPR (2019)
- 467 [61] Lu, Y., Wen, L., Liu, J., Liu, Y., Tian, X.: Self-supervision can be a good few-shot learner. In:  
468 European Conference on Computer Vision. pp. 740–758. Springer (2022)
- 469 [62] Luo, X.J., Wang, S., Wu, Z., Sakaridis, C., Cheng, Y., Fan, D.P., Gool, L.V.: Camdiff: Cam-  
470 ouflage image augmentation via diffusion model (2023)
- 471 [63] Luzi, L., Siahkoohi, A., Mayer, P.M., Casco-Rodriguez, J., Baraniuk, R.: Boomerang: Local  
472 sampling on image manifolds using diffusion models (2022)
- 473 [64] Maji, S., Rahtu, E., Kannala, J., Blaschko, M.B., Vedaldi, A.: Fine-grained visual classifica-  
474 tion of aircraft. arXiv preprint arXiv:1306.5151 (2013)
- 475 [65] Mao, J., Xiao, T., Jiang, Y., Cao, Z.: What can help pedestrian detection? (2017)
- 476 [66] Medina, C., Devos, A., Grossglauser, M.: Self-supervised prototypical transfer learning for  
477 few-shot classification. In: ICMLW (2020)
- 478 [67] Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Guided image  
479 synthesis and editing with stochastic differential equations. arXiv preprint arXiv:2108.01073  
480 (2021)

- 481 [68] Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S., Ho, J., Salimans, T.: On distilla-  
482 tion of guided diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer  
483 Vision and Pattern Recognition. pp. 14297–14306 (2023)
- 484 [69] Menon, A.K., Jayasumana, S., Rawat, A.S., Jain, H., Veit, A., Kumar, S.: Long-tail learning  
485 via logit adjustment. In: ICLR (2021)
- 486 [70] Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P.,  
487 Haziza, D., Massa, F., El-Nouby, A., Assran, M., Ballas, N., Galuba, W., Howes, R., Huang,  
488 P.Y., Li, S.W., Misra, I., Rabbat, M., Sharma, V., Synnaeve, G., Xu, H., Jegou, H., Mairal,  
489 J., Labatut, P., Joulin, A., Bojanowski, P.: DINOv2: Learning robust visual features without  
490 supervision (2023)
- 491 [71] Peebles, W., Zhu, J.Y., Zhang, R., Torralba, A., Efros, A., Shechtman, E.: Gan-supervised  
492 dense visual alignment. In: CVPR (2022)
- 493 [72] Petryk, S., Dunlap, L., Nasser, K., Gonzalez, J., Darrell, T., Rohrbach, A.: On guid-  
494 ing visual attention with language specification. In: Conference on Computer Vision and  
495 Pattern Recognition (CVPR) (2022). <https://doi.org/10.48550/ARXIV.2202.08926>, <https://arxiv.org/abs/2202.08926>
- 497 [73] Prabhu, V., Yenamandra, S., Chattopadhyay, P., Hoffman, J.: Lance: Stress-testing visual  
498 models by generating language-guided counterfactual images. *Advances in Neural Informa-*  
499 *tion Processing Systems* **36** (2024)
- 500 [74] Qiao, S., Liu, C., Shen, W., Yuille, A.: Few-shot image recognition by predicting parameters  
501 from activations. In: CVPR (2018)
- 502 [75] Qin, T., Li, W., Shi, Y., Yang, G.: Unsupervised few-shot learning via distribution shift-based  
503 augmentation. *arxiv:2004.05805* (2020)
- 504 [76] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell,  
505 A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models  
506 from natural language supervision. In: ICML (2021)
- 507 [77] Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image  
508 generation with clip latents. *arXiv preprint arXiv:2204.06125* **1**(2), 3 (2022)
- 509 [78] Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I.:  
510 Zero-shot text-to-image generation. In: ICML (2021)
- 511 [79] Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2017)
- 512 [80] Ren, M., Ravi, S., Triantafillou, E., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle,  
513 H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. In: International  
514 Conference on Learning Representations (2018), [https://openreview.net/forum?id=](https://openreview.net/forum?id=HJcSzz-CZ)  
515 [HJcSzz-CZ](https://openreview.net/forum?id=HJcSzz-CZ)
- 516 [81] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image syn-  
517 thesis with latent diffusion models. In: CVPR (2022)
- 518 [82] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image syn-  
519 thesis with latent diffusion models (2021)
- 520 [83] Roy, A., Shah, A., Shah, K., Roy, A., Chellappa, R.: Cap2aug: Caption guided image to  
521 image data augmentation (2023)
- 522 [84] Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gon-  
523 tijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion  
524 models with deep language understanding. *Advances in Neural Information Processing Sys-*  
525 *tems* **35**, 36479–36494 (2022)
- 526 [85] Samuel, D., Chechik, G.: Distributional robustness loss for long-tail learning. In: ICCV  
527 (2021)
- 528 [86] Sankaranarayanan, S., Balaji, Y., Castillo, C.D., Chellappa, R.: Generate to adapt: Aligning  
529 domains using generative adversarial networks. *Conference on Computer Vision and Pattern*  
530 *Recognition (CVPR)* (2018)
- 531 [87] Sauer, A., Lorenz, D., Blattmann, A., Rombach, R.: Adversarial diffusion distillation. *arXiv*  
532 *preprint arXiv:2311.17042* (2023)

- 533 [88] Sharmanska, V., Hendricks, L.A., Darrell, T., Quadrianto, N.: Contrastive examples for ad-  
534 dressing the tyranny of the majority. CoRR **abs/2004.06524** (2020), [https://arxiv.org/  
535 abs/2004.06524](https://arxiv.org/abs/2004.06524)
- 536 [89] Shipard, J., Wiliem, A., Thanh, K.N., Xiang, W., Fookes, C.: Boosting zero-shot classification  
537 with synthetic data diversity via stable diffusion. arXiv preprint arXiv:2302.03298 (2023)
- 538 [90] Shirekar, O.K., Singh, A., Jamali-Rad, H.: Self-attention message passing for contrastive  
539 few-shot learning. In: Proceedings of the IEEE/CVF Winter Conference on Applications of  
540 Computer Vision (WACV). pp. 5426–5436 (January 2023)
- 541 [91] Shorten, C., Khoshgoftaar, T.M.: A survey on image data augmentation for deep learning.  
542 Journal of big data **6**(1), 1–48 (2019)
- 543 [92] Singh, A.R., Jamali-Rad, H.: Transductive decoupled variational inference for few-shot  
544 classification. Transactions on Machine Learning Research (2023), [https://openreview.  
545 net/forum?id=bomdTc9HyL](https://openreview.net/forum?id=bomdTc9HyL)
- 546 [93] Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: Advances  
547 in Neural Information Processing Systems (2017)
- 548 [94] Su, J.C., Maji, S., Hariharan, B.: When does self-supervision improve few-shot learning? In:  
549 ECCV (2020)
- 550 [95] Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H., Hospedales, T.M.: Learning to compare:  
551 Relation network for few-shot learning. In: CVPR (2018)
- 552 [96] Tang, K., Huang, J., Zhang, H.: Long-tailed classification by keeping the good and removing  
553 the bad momentum causal effect. NeurIPS **33**, 1513–1524 (2020)
- 554 [97] Tian, C., Wang, W., Zhu, X., Dai, J., Qiao, Y.: Vi-ltr: Learning class-wise visual-linguistic  
555 representation for long-tailed visual recognition. In: ECCV 2022 (2022)
- 556 [98] Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-  
557 efficient image transformers and distillation through attention (2021)
- 558 [99] Touvron, H., Cord, M., Jégou, H.: Deit iii: Revenge of the vit. In: ECCV (2022)
- 559 [100] Trabucco, B., Doherty, K., Gurinas, M.A., Salakhutdinov, R.: Effective data augmentation  
560 with diffusion models. In: The Twelfth International Conference on Learning Representations  
561 (2024), <https://openreview.net/forum?id=ZwzUA9zeAg>
- 562 [101] Tritrong, N., Rewatbowornwong, P., Suwajanakorn, S.: Repurposing gans for one-shot se-  
563 mantic part segmentation. In: IEEE Conference on Computer Vision and Pattern Recognition  
564 (CVPR) (2021)
- 565 [102] Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011  
566 dataset (2011)
- 567 [103] Wang, H., Deng, Z.H.: Contrastive prototypical network with wasserstein confidence penalty.  
568 In: European Conference on Computer Vision. pp. 665–682. Springer (2022)
- 569 [104] Wang, H., Fu, S., He, X., Fang, H., Liu, Z., Hu, H.: Towards calibrated hyper-sphere repre-  
570 sentation via distribution overlap coefficient for long-tailed learning. In: ECCV (2022)
- 571 [105] Wang, X., Lian, L., Miao, Z., Liu, Z., Yu, S.X.: Long-tailed recognition by routing diverse  
572 distribution-aware experts. In: ICLR. OpenReview.net (2021)
- 573 [106] Xu, Y., Li, Y.L., Li, J., Lu, C.: Constructing balance from imbalance for long-tailed image  
574 recognition. In: ECCV. pp. 38–56. Springer (2022)
- 575 [107] Xu, Z., Liu, R., Yang, S., Chai, Z., Yuan, C.: Learning imbalanced data with vision trans-  
576 formers (2023)
- 577 [108] Xuan, H., Stylianou, A., Liu, X., Pless, R.: Hard negative examples are hard, but useful  
578 (2021)
- 579 [109] Ye, H.J., Hu, H., Zhan, D.C., Sha, F.: Few-shot learning via embedding adaptation with  
580 set-to-set functions. In: CVPR (2020)
- 581 [110] Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to  
582 train strong classifiers with localizable features. In: ICCV. pp. 6023–6032 (2019)

- 583 [111] Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk mini-  
584 mization. In: ICLR (2018)
- 585 [112] Zhang, S., Li, Z., Yan, S., He, X., Sun, J.: Distribution alignment: A unified framework for  
586 long-tail visual recognition. In: CVPR. pp. 2361–2370 (2021)
- 587 [113] Zhang, Y., Hooi, B., Hong, L., Feng, J.: Test-agnostic long-tailed recognition by test-time  
588 aggregating diverse experts with self-supervision. arXiv preprint arXiv:2107.09249 (2021)
- 589 [114] Zhang, Y., Ling, H., Gao, J., Yin, K., Lafleche, J.F., Barriuso, A., Torralba, A., Fidler, S.:  
590 Datasetgan: Efficient labeled data factory with minimal human effort. In: CVPR (2021)
- 591 [115] Zhong, Z., Cui, J., Liu, S., Jia, J.: Improving calibration for long-tailed recognition. In:  
592 CVPR. pp. 16489–16498. Computer Vision Foundation / IEEE (2021)
- 593 [116] Zhou, Z., Qiu, X., Xie, J., Wu, J., Zhang, C.: Binocular mutual learning for improving few-  
594 shot classification. In: ICCV (2021)
- 595 [117] Zhu, J., Wang, Z., Chen, J., Chen, Y.P.P., Jiang, Y.G.: Balanced contrastive learning for long-  
596 tailed visual recognition. In: CVPR. pp. 6908–6917 (2022)

597 **A Appendix**

598 **A.1 Analyzing GeNIe, GeNIe-Ada’s Class-Probabilities**

599 The core aim of GeNIe and GeNIe-Ada is to address the failure modes of a classifier  
 600 by generating *challenging* samples located near the decision boundary of each class pair,  
 601 which facilitates the learning process in effectively enhancing the decision boundary between  
 602 classes. As summarized in Table 4 and illustrated in Fig. 5, we have empirically corroborated that GeNIe and GeNIe-Ada can respectively produce samples  $X_r, X_{r^*}$  that are negative with respect to the source image  $X_S$ , while semantically belonging to the class  $T$ . To

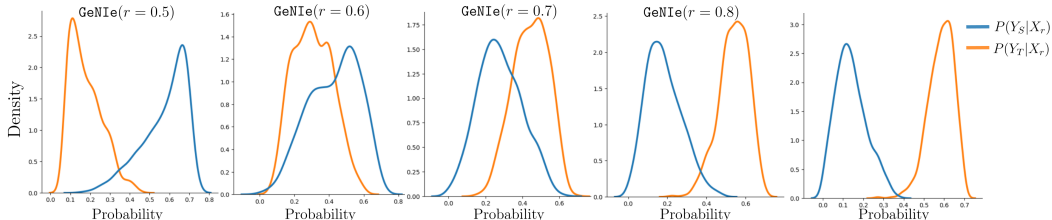


Figure A1:  $P(Y_S|X_r)$  and  $P(Y_T|X_r)$  for  $r \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$ . On average, the classifier confidently predicts the source class more than the target class for  $X_r$  for  $r = 0.5$ , and vice-versa for  $r = 0.8, 0.9$ . However, for  $r = 0.6, 0.7$ , the classifier struggles to classify  $X_r$ , indicating that the augmented samples are located closer to the decision boundary.

604 further analyze the effectiveness of GeNIe and GeNIe-Ada, we compare the source class-  
 605 probabilities  $P(Y_S|X_r)$  and target-class probabilities  $P(Y_T|X_r)$  of augmented samples  $X_r$ .  
 606 To compute these class probabilities, we first fit an SVM classifier  
 607 (as followed in UniSiam [61]) only on the labelled support set em-  
 608 beddings of each episode in the *mini*Imagenet test dataset. Then,  
 609 we perform inference using each episode’s SVM classifier on its re-  
 610 spective  $X_r$ ’s and extract its class probabilities of belonging to its  
 611 source class  $S$  and target class  $T$ . These per augmentation-sample  
 612 source and target class probabilities are then averaged for each  
 613 episode for each  $r \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$  in the case of GeNIe  
 614 and for the optimal  $r = r^*$  per sample in the case of GeNIe-Ada,  
 615 plotted as density plots in Fig. A1, Fig. A2, respectively. Fig. A1  
 616 illustrates that  $P(Y_S|X_r)$  and  $P(Y_T|X_r)$  have significant overlap  
 617 in the case of  $r \in \{0.6, 0.7\}$  indicating class-confusion for  $X_r$ .

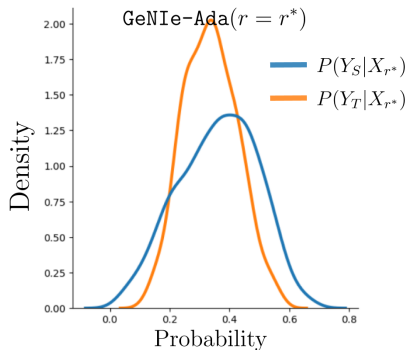


Figure A2: Significant overlap between  $P(Y_S|X_{r^*})$  and  $P(Y_T|X_{r^*})$  indicates high class-confusion for augmented samples generated by GeNIe-Ada.

619 Furthermore, Fig. A2 illustrates that when using the optimal  $r = r^*$   
 620 found by GeNIe-Ada per sample,  $P(Y_S|X_r)$  and  $P(Y_T|X_r)$  signif-  
 621 icantly overlap around probability scores of 0.2 – 0.45, indicating  
 622 class confusion for GeNIe-Ada augmentations. This corroborates  
 623 with our analysis in Section 4.3, Table 4 and additionally empirically  
 624 proves that the augmented samples generated by GeNIe for  
 625  $r \in \{0.6, 0.7\}$  and GeNIe-Ada for  $r = r^*$  are actually located near  
 626 the decision boundary of each class pair.

627 **A.2 Fine-grained Few-shot Classification**

628 To further investigate the impact of the proposed method, we compare GeNIe with other text-based  
 629 data augmentation techniques across four distinct fine-grained datasets in a 20-way, 1-shot classifica-  
 630 tion setting. We employ the pre-trained DINOv2 ViT-G [70] backbone as a feature extractor to  
 631 derive features from training images. Subsequently, an SVM classifier is trained on these features,  
 632 and we report the Top-1 accuracy of the model on the test set.

633 **Datasets:** We assess our method on several datasets: Food101 [5] with 101 classes of various foods,  
 634 CUB200 [102] with 200 bird species classes, Cars196 [48] with 196 car model classes, and FGVC-  
 635 Aircraft [64] with 41 aircraft manufacturer classes. We provide detailed information around fine-  
 636 grained datasets in Table A2. The reported metric is the average Top-1 accuracy over 100 episodes.



637 Each episode involves sampling 20 classes and 1-shot from the training set, with the final model  
 638 evaluated on the respective test set.

639 **Implementation Details:** We enhance the basic prompt by incorporating the superclass name for  
 640 the fine-grained dataset: “A photo of a <target class>, a type of <superclass>”. For instance,  
 641 in the *food* dataset and the *burger* class, our prompt reads: “A photo of a *burger*, a type of *food*.” No  
 642 additional augmentation is used for generative methods in this context. We generate 19 samples for  
 643 both cases of our method and also the baseline with weak augmentation.

644 **Results:** Table A1 summarizes the results. GeNIe helps outperform all other baselines and aug-  
 645 mentations, including Txt2Img, by margins upto 0.5% on CUB200 [102], 6.6% on Cars196 [48],  
 646 0.1% on Food101 [5] and 5.3% on FGVC-Aircraft [64]. Notably, GeNIe exhibits great effectiveness  
 647 in more challenging datasets, outperforming the baseline with traditional augmentation by about  
 648 38% for the Cars dataset and by roughly 17% for the Aircraft dataset. It can be observed here that  
 649 GeNIe-Ada performs on-par with GeNIe with a fixed noise level, eliminating the necessity for noise  
 650 level search in GeNIe.

Table A1: **Few-shot Learning on Fine-grained dataset:** We utilize an SVM classifier trained atop the DI-NOV2 ViT-G pretrained backbone, reporting Top-1 accuracy for the test set of each dataset. The baseline is an SVM trained on the same backbone using weak augmentation. Across all datasets, GeNIe surpasses this baseline.

Method	Birds	Cars	Foods	Aircraft
	CUB200 [102]	Cars196 [48]	Food101 [5]	Aircraft [64]
Baseline	90.3	49.8	82.9	29.2
Img2Img <sup>L</sup> [63]	90.7	50.4	87.4	31.0
Img2Img <sup>H</sup> [63]	91.3	56.4	91.7	34.7
Txt2Img[35]	92.0	81.3	93.0	41.7
GeNIe (r=0.5)	92.0	84.6	91.5	39.8
GeNIe (r=0.6)	92.2	87.1	92.5	45.0
GeNIe (r=0.7)	92.5	<b>87.9</b>	92.9	<b>47.0</b>
GeNIe (r=0.8)	92.5	87.7	<b>93.1</b>	46.5
GeNIe (r=0.9)	92.4	87.1	<b>93.1</b>	45.7
GeNIe-Ada	<b>92.6</b>	<b>87.9</b>	<b>93.1</b>	46.9

Table A2: Train and test split details of the fine-grained datasets. We use the provided train set for few-shot task generation, and the provided test sets for our evaluation. For the Aircraft dataset we use manufacturer hierarchy.

Dataset	Classes	Train samples	Test samples
CUB200 [102]	200	5994	5794
Food101 [5]	101	75750	25250
Cars [48]	196	8144	8041
Aircraft [64]	41	6,667	3333

651 **A.3 Few-shot Classification with ResNet-34 on *tiered*Imagenet**

652 We follow the same evaluation protocol here as mentioned in section 4.1. As summarized in Table  
 653 A3, GeNIe and GeNIe-Ada outperform all other classical and generative data augmentation  
 654 techniques.

655 **A.4 Additional details of Long-Tail experiments**

656 We present a comprehensive version of Table 3 to benchmark the performance with different back-  
 657 bone architectures (e.g., ResNet50) and to compare against previous long-tail baselines; this is de-  
 658 tailed in Table A4.

659 **Implementation Details of LViT:** We download the pre-trained ViT-B of LViT [107] and finetune  
 660 it with Bal-BCE loss proposed therein on the augmented dataset. Training takes 2 hours on four  
 661 NVIDIA RTX 3090 GPUs. We use the same hyperparameters as in [107] for finetuning: 100 epochs,  
 662  $lr = 0.008$ , batch size of 1024, CutMix and MixUp for the data augmentation.

663 **Implementation Details of VL-LTR:** We use the official code of VL-LTR [97] for our experiments.  
 664 We use a pre-trained CLIP ResNet-50 backbone. We followed the hyperparameters reported in VL-

Table A3: *tiered-ImageNet*: Accuracies ( $\% \pm \text{std}$ ) for 5-way, 1-shot and 5-way, 5-shot classification settings on the test-set. We compare against various SOTA supervised and unsupervised few-shot classification baselines as well as other augmentation methods, with UniSiam [61] pre-trained ResNet-34 backbone.

ResNet-34				
Augmentation	Method	Pre-training	1-shot	5-shot
Weak	MAML + dist [29]	sup.	51.7±1.8	70.3±1.7
Weak	ProtoNet [93]	sup.	52.0±1.2	72.1±1.5
Weak	UniSiam + dist [61]	unsup.	68.7±0.4	85.7±0.3
Weak	UniSiam [61]	unsup.	65.0±0.7	82.5±0.5
Strong	UniSiam [61]	unsup.	64.8±0.7	82.4±0.5
CutMix [110]	UniSiam [61]	unsup.	63.8±0.7	80.3±0.6
MixUp [111]	UniSiam [61]	unsup.	64.1±0.7	80.0±0.6
Img2Img <sup>L</sup> [63]	UniSiam [61]	unsup.	66.1±0.7	83.1±0.5
Img2Img <sup>H</sup> [63]	UniSiam [61]	unsup.	70.4±0.7	84.7±0.5
Txt2Img[35]	UniSiam [61]	unsup.	75.0±0.6	85.4±0.4
DAFusion [100]	UniSiam [61]	unsup.	64.1±2.1	82.8±1.4
GeNIe (Ours)	UniSiam [61]	unsup.	<b>75.7±0.6</b>	<b>86.0±0.4</b>
GeNIe-Ada (Ours)	UniSiam [61]	unsup.	<b>76.9±0.6</b>	<b>86.3±0.2</b>

665 LTR [97]. We augment only “Few” category and train the backbone with the VL-LTR [97] method.  
 666 Training takes 4 hours on 8 NVIDIA RTX 3090 GPUs.

## 667 A.5 More Visualizations

668 Additional qualitative results resembling the style presented in Fig. 4 are presented in Fig. A3, and  
 669 more visuals akin to Fig. 2 can be found in Fig. A4. Moreover, we also present more visualization  
 670 similar to the style in Fig. 5 in Fig. A5.

Table A4: **Long-Tailed ImageNet-LT**: We compare different augmentation methods on ImageNet-LT and report Top-1 accuracy for “Few”, “Medium”, and “Many” sets. † indicates results with ResNeXt50. \*: indicates training with 384 resolution so is not directly comparable with other methods with 224 resolution. On the “Few” set and LiVT method, our augmentations improve the accuracy by 11.7 points compared to LiVT original augmentation and 4.4 points compared to Txt2Img.

ResNet-50				
Method	Many	Med.	Few	Overall Acc
CE [21]	64.0	33.8	5.8	41.6
LDAM [7]	60.4	46.9	30.7	49.8
c-RT [45]	61.8	46.2	27.3	49.6
$\tau$ -Norm [45]	59.1	46.9	30.7	49.4
Causal [96]	62.7	48.8	31.6	51.8
Logit Adj. [69]	61.1	47.5	27.6	50.1
RIDE(4E)† [105]	68.3	53.5	35.9	56.8
MiSLAS [115]	62.9	50.7	34.3	52.7
DisAlign [112]	61.3	52.2	31.4	52.9
ACE† [6]	71.7	54.6	23.5	56.6
PaCo† [20]	68.0	56.4	37.2	58.2
TADE† [113]	66.5	<b>57.0</b>	43.5	58.8
TSC [56]	63.5	49.7	30.4	52.4
GCL [55]	63.0	52.7	37.1	54.5
TLC [50]	68.9	55.7	40.8	55.1
BCL† [117]	67.6	54.6	36.6	57.2
NCL [52]	67.3	55.4	39.0	57.7
SAFA [38]	63.8	49.9	33.4	53.1
DOC [104]	65.1	52.8	34.2	55.0
DLSA [106]	67.8	54.5	38.8	57.5
ResLT [18]	63.3	53.3	40.3	55.1
PaCo [19]	68.2	58.7	41.0	60.0
LWS [44]	62.2	48.6	31.8	51.5
Zero-shot CLIP [76]	60.8	59.3	58.6	59.8
DRO-LT [85]	64.0	49.8	33.1	53.5
VL-LTR [97]	77.8	67.0	50.8	70.1
Cap2Aug [83]	78.5	<b>67.7</b>	51.9	70.9
GeNIe-Ada	<b>79.2</b>	64.6	<b>59.5</b>	<b>71.5</b>
ViT-B				
LiVT* [107]	76.4	59.7	42.7	63.8
ViT [24]	50.5	23.5	6.9	31.6
MAE [33]	74.7	48.2	19.4	54.5
DeiT [99]	70.4	40.9	12.8	48.4
LiVT [107]	73.6	56.4	41.0	60.9
LiVT + Img2Img <sup>L</sup>	74.3	56.4	34.3	60.5
LiVT + Img2Img <sup>H</sup>	73.8	56.4	45.3	61.6
LiVT + Txt2Img	<b>74.9</b>	55.6	48.3	62.2
LiVT + GeNIe (r=0.8)	74.5	56.7	50.9	62.8
LiVT + GeNIe-Ada	74.0	<b>56.9</b>	<b>52.7</b>	<b>63.1</b>

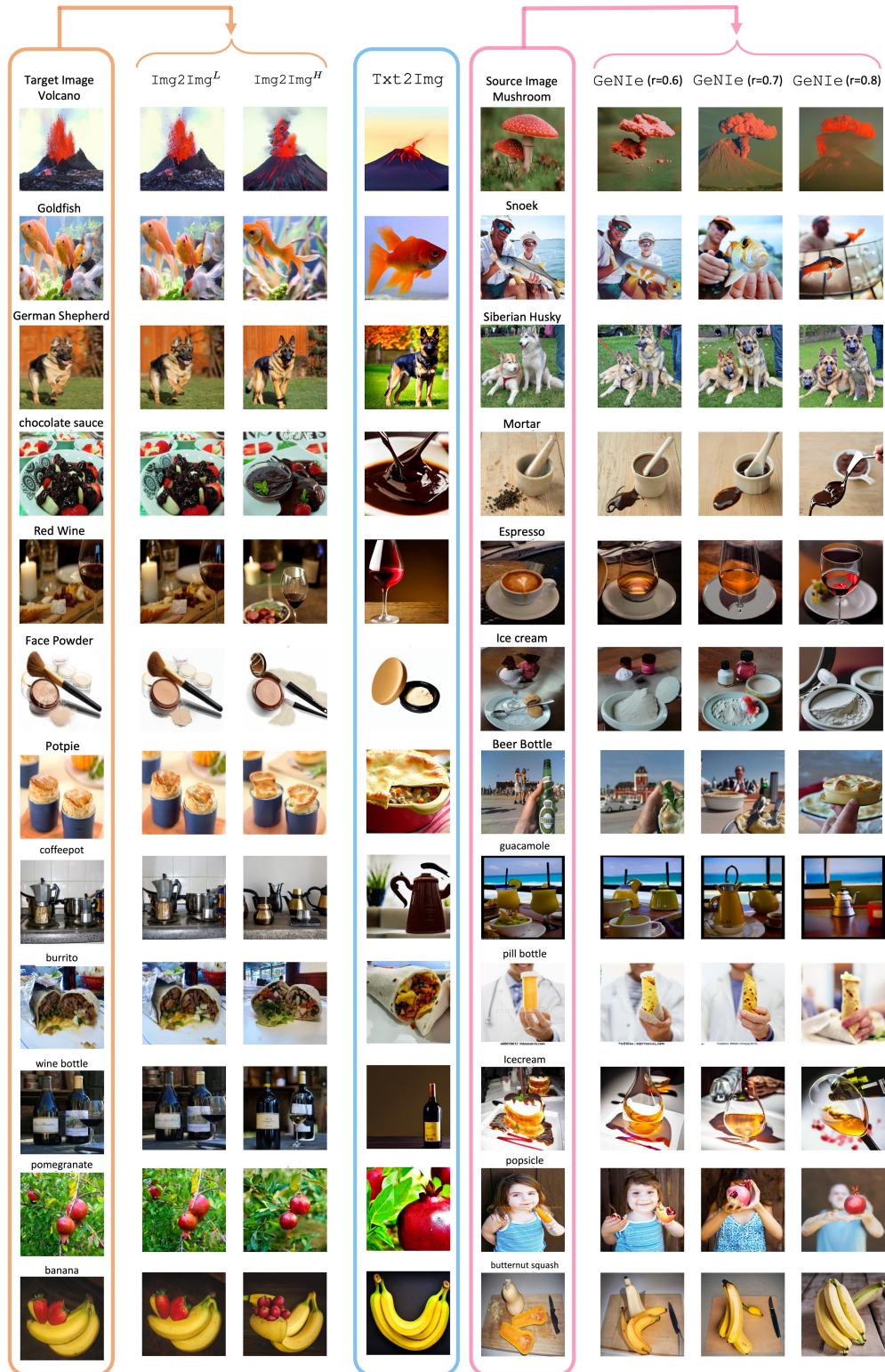


Figure A3: **Visualization of Generative Samples:** More visualization akin to Fig. 4. We compare GeNIe with two baselines:  $\text{Img2Img}^L$  **augmentation** uses both image and text prompt from the same category, resulting in less challenging examples.  $\text{Txt2Img}$  **augmentation** generates images based solely on a text prompt, potentially deviating from the task's visual domain. GeNIe **augmentation** incorporates the target category name in the text prompt along with the source image, producing desired images with an optimal amount of noise, and balancing the impact of the source image and text prompt.

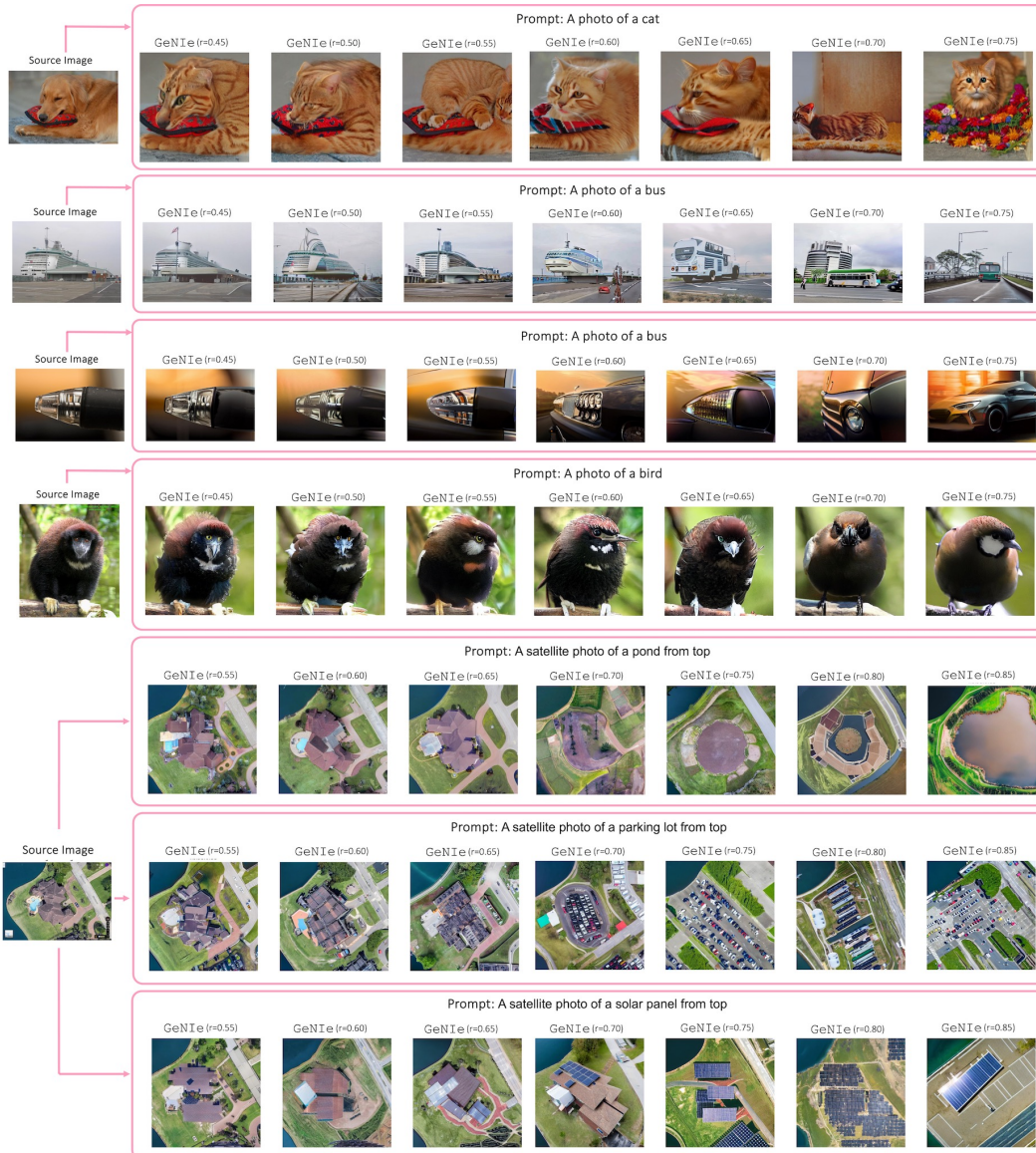


Figure A4: **Effect of noise in GeNIe:** Akin to Fig. 2, we use GeNIe to create augmentations with varying noise levels. As is illustrated in the examples above, a reduced amount of noise leads to images closely mirroring the semantics of the source images, causing a misalignment with the intended target label.

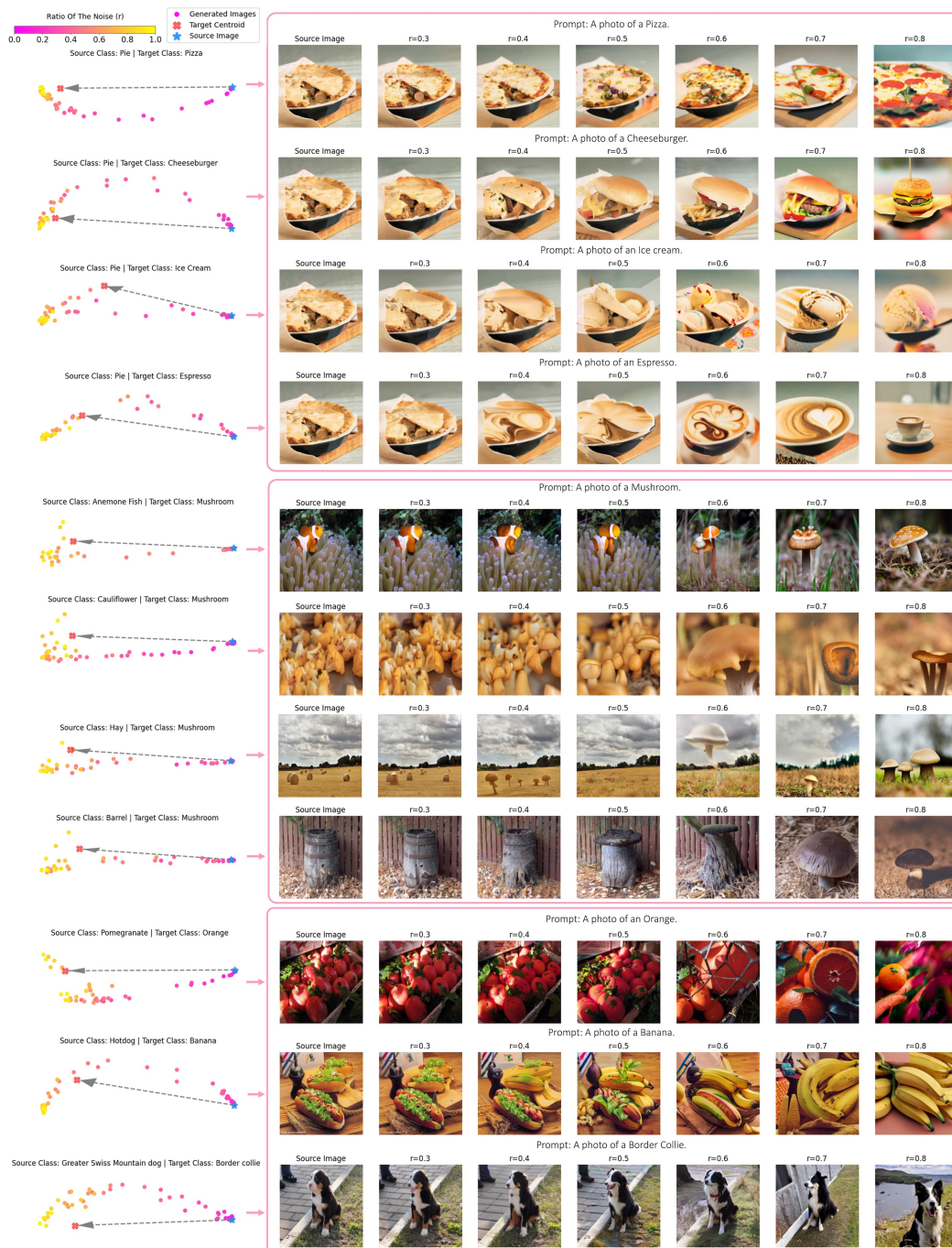


Figure A5: **Effect of noise in GeNIe:** Similar to Fig. 5, we pass all the generated augmentations through the DinoV2 ViT-G model, which acts as our oracle model, to obtain their associated embeddings. Subsequently, we employ PCA for visualization purposes. The visualization reveals that the magnitude of semantic transformations is contingent upon both the source image and the specified target category.

## 671 **NeurIPS Paper Checklist**

### 672 **1. Claims**

673 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s  
674 contributions and scope?

675 Answer: [Yes]

676 Justification: We demonstrate the effectiveness of our augmentation method through empirical  
677 comparison with four different generative augmentation baselines across two scenarios: few-shot  
678 and long-tail classification. Additionally, we perform analytical experiments on our augmented  
679 samples to illustrate their nature as hard negatives.

680 Guidelines:

- 681 • The answer NA means that the abstract and introduction do not include the claims made in the  
682 paper.
- 683 • The abstract and/or introduction should clearly state the claims made, including the contributions  
684 made in the paper and important assumptions and limitations. A No or NA answer to this question  
685 will not be perceived well by the reviewers.
- 686 • The claims made should match theoretical and experimental results, and reflect how much the  
687 results can be expected to generalize to other settings.
- 688 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not  
689 attained by the paper.

### 690 **2. Limitations**

691 Question: Does the paper discuss the limitations of the work performed by the authors?

692 Answer: [Yes]

693 Justification: We discuss about the limitations of our method in Sec 5

694 Guidelines:

- 695 • The answer NA means that the paper has no limitation while the answer No means that the paper  
696 has limitations, but those are not discussed in the paper.
- 697 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 698 • The paper should point out any strong assumptions and how robust the results are to violations of  
699 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,  
700 asymptotic approximations only holding locally). The authors should reflect on how these as-  
701 sumptions might be violated in practice and what the implications would be.
- 702 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested  
703 on a few datasets or with a few runs. In general, empirical results often depend on implicit  
704 assumptions, which should be articulated.
- 705 • The authors should reflect on the factors that influence the performance of the approach. For  
706 example, a facial recognition algorithm may perform poorly when image resolution is low or  
707 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide  
708 closed captions for online lectures because it fails to handle technical jargon.
- 709 • The authors should discuss the computational efficiency of the proposed algorithms and how they  
710 scale with dataset size.
- 711 • If applicable, the authors should discuss possible limitations of their approach to address prob-  
712 lems of privacy and fairness.
- 713 • While the authors might fear that complete honesty about limitations might be used by reviewers  
714 as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t  
715 acknowledged in the paper. The authors should use their best judgment and recognize that indi-  
716 vidual actions in favor of transparency play an important role in developing norms that preserve  
717 the integrity of the community. Reviewers will be specifically instructed to not penalize honesty  
718 concerning limitations.

### 719 **3. Theory Assumptions and Proofs**

720 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
721 complete (and correct) proof?

722 Answer: [NA]

723 Justification: We do not have theoretical results.

724 Guidelines:

- 725 • The answer NA means that the paper does not include theoretical results.
- 726 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 727 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 728 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
- 729 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
- 730 intuition.
- 731 • Inversely, any informal proof provided in the core of the paper should be complemented by formal
- 732 proofs provided in appendix or supplemental material.
- 733 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 734 4. Experimental Result Reproducibility

735 Question: Does the paper fully disclose all the information needed to reproduce the main exper-  
736 imental results of the paper to the extent that it affects the main claims and/or conclusions of the  
737 paper (regardless of whether the code and data are provided or not)?

738 Answer: [Yes]

739 Justification: We provide implementation details in each experimental section. Additionally, we  
740 include the code as supplementary material and plan to release it publicly.

741 Guidelines:

- 742 • The answer NA means that the paper does not include experiments.
- 743 • If the paper includes experiments, a No answer to this question will not be perceived well by the
- 744 reviewers: Making the paper reproducible is important, regardless of whether the code and data
- 745 are provided or not.
- 746 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
- 747 their results reproducible or verifiable.
- 748 • Depending on the contribution, reproducibility can be accomplished in various ways. For exam-
- 749 ple, if the contribution is a novel architecture, describing the architecture fully might suffice, or if
- 750 the contribution is a specific model and empirical evaluation, it may be necessary to either make
- 751 it possible for others to replicate the model with the same dataset, or provide access to the model.
- 752 In general, releasing code and data is often one good way to accomplish this, but reproducibility
- 753 can also be provided via detailed instructions for how to replicate the results, access to a hosted
- 754 model (e.g., in the case of a large language model), releasing of a model checkpoint, or other
- 755 means that are appropriate to the research performed.
- 756 • While NeurIPS does not require releasing code, the conference does require all submissions to
- 757 provide some reasonable avenue for reproducibility, which may depend on the nature of the con-
- 758 tribution. For example
- 759 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce
- 760 that algorithm.
- 761 (b) If the contribution is primarily a new model architecture, the paper should describe the architec-
- 762 ture clearly and fully.
- 763 (c) If the contribution is a new model (e.g., a large language model), then there should either be a
- 764 way to access this model for reproducing the results or a way to reproduce the model (e.g., with
- 765 an open-source dataset or instructions for how to construct the dataset).
- 766 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are wel-
- 767 come to describe the particular way they provide for reproducibility. In the case of closed-source
- 768 models, it may be that access to the model is limited in some way (e.g., to registered users), but
- 769 it should be possible for other researchers to have some path to reproducing or verifying the
- 770 results.

#### 771 5. Open access to data and code

772 Question: Does the paper provide open access to the data and code, with sufficient instructions to  
773 faithfully reproduce the main experimental results, as described in supplemental material?

774 Answer: [Yes]

775 Justification: We provide implementation details in each experimental section. Additionally, we  
776 include the code as supplementary material and plan to release it publicly.



777 Guidelines:

- 778 • The answer NA means that paper does not include experiments requiring code.
- 779 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 780
- 781 • While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 782
- 783
- 784 • The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 785
- 786
- 787 • The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 788
- 789 • The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- 790
- 791
- 792 • At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 793
- 794 • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
- 795

#### 796 6. Experimental Setting/Details

797 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

799 Answer: [Yes]

800 Justification: We provide implementation details and dataset details in each experimental section.

801 Guidelines:

- 802 • The answer NA means that the paper does not include experiments.
- 803 • The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- 804
- 805 • The full details can be provided either with the code, in appendix, or as supplemental material.

#### 806 7. Experiment Statistical Significance

807 Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

809 Answer: [Yes]

810 Justification: We repeat few-shot training for 600 episodes on mini-ImageNet and 1000 episodes on tiered-ImageNet, reporting the mean and variance for each method.

812 Guidelines:

- 813 • The answer NA means that the paper does not include experiments.
- 814 • The authors should answer “Yes” if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- 815
- 816
- 817 • The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- 818
- 819
- 820 • The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- 821
- 822 • The assumptions made should be given (e.g., Normally distributed errors).
- 823 • It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- 824 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- 825
- 826
- 827 • For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 828

- 829 • If error bars are reported in tables or plots, The authors should explain in the text how they were  
830 calculated and reference the corresponding figures or tables in the text.

## 831 8. Experiments Compute Resources

832 Question: For each experiment, does the paper provide sufficient information on the computer  
833 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
834 ments?

835 Answer: [Yes]

836 Justification: We provide implementation and dataset details in each experimental section. Ad-  
837 ditionally, we elaborate on the required resources, including GPUs and training hours, for each  
838 experiment.

839 Guidelines:

- 840 • The answer NA means that the paper does not include experiments.
- 841 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud  
842 provider, including relevant memory and storage.
- 843 • The paper should provide the amount of compute required for each of the individual experimental  
844 runs as well as estimate the total compute.
- 845 • The paper should disclose whether the full research project required more compute than the ex-  
846 periments reported in the paper (e.g., preliminary or failed experiments that didn't make it into  
847 the paper).

## 848 9. Code Of Ethics

849 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS  
850 Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

851 Answer: [Yes]

852 Justification: We reviewed the NeurIPS Code of Ethics.

853 Guidelines:

- 854 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 855 • If the authors answer No, they should explain the special circumstances that require a deviation  
856 from the Code of Ethics.
- 857 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due  
858 to laws or regulations in their jurisdiction).

## 859 10. Broader Impacts

860 Question: Does the paper discuss both potential positive societal impacts and negative societal  
861 impacts of the work performed?

862 Answer: [Yes]

863 Justification: We discuss about broader impact in Conclusion.

864 Guidelines:

- 865 • The answer NA means that there is no societal impact of the work performed.
- 866 • If the authors answer NA or No, they should explain why their work has no societal impact or  
867 why the paper does not address societal impact.
- 868 • Examples of negative societal impacts include potential malicious or unintended uses (e.g., dis-  
869 information, generating fake profiles, surveillance), fairness considerations (e.g., deployment of  
870 technologies that could make decisions that unfairly impact specific groups), privacy considera-  
871 tions, and security considerations.
- 872 • The conference expects that many papers will be foundational research and not tied to particular  
873 applications, let alone deployments. However, if there is a direct path to any negative applications,  
874 the authors should point it out. For example, it is legitimate to point out that an improvement in  
875 the quality of generative models could be used to generate deepfakes for disinformation. On the  
876 other hand, it is not needed to point out that a generic algorithm for optimizing neural networks  
877 could enable people to train models that generate Deepfakes faster.

- 878 • The authors should consider possible harms that could arise when the technology is being used  
879 as intended and functioning correctly, harms that could arise when the technology is being used  
880 as intended but gives incorrect results, and harms following from (intentional or unintentional)  
881 misuse of the technology.
- 882 • If there are negative societal impacts, the authors could also discuss possible mitigation strategies  
883 (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for moni-  
884 toring misuse, mechanisms to monitor how a system learns from feedback over time, improving  
885 the efficiency and accessibility of ML).

### 886 1. Safeguards

887 Question: Does the paper describe safeguards that have been put in place for responsible release of  
888 data or models that have a high risk for misuse (e.g., pretrained language models, image generators,  
889 or scraped datasets)?

890 Answer: [NA]

891 Justification: We believe our work does not have such risks.

892 Guidelines:

- 893 • The answer NA means that the paper poses no such risks.
- 894 • Released models that have a high risk for misuse or dual-use should be released with necessary  
895 safeguards to allow for controlled use of the model, for example by requiring that users adhere to  
896 usage guidelines or restrictions to access the model or implementing safety filters.
- 897 • Datasets that have been scraped from the Internet could pose safety risks. The authors should  
898 describe how they avoided releasing unsafe images.
- 899 • We recognize that providing effective safeguards is challenging, and many papers do not require  
900 this, but we encourage authors to take this into account and make a best faith effort.

### 901 2. Licenses for existing assets

902 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,  
903 properly credited and are the license and terms of use explicitly mentioned and properly respected?

904 Answer: [Yes]

905 Justification: We cited all datasets and code used in our paper.

906 Guidelines:

- 907 • The answer NA means that the paper does not use existing assets.
- 908 • The authors should cite the original paper that produced the code package or dataset.
- 909 • The authors should state which version of the asset is used and, if possible, include a URL.
- 910 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 911 • For scraped data from a particular source (e.g., website), the copyright and terms of service of  
912 that source should be provided.
- 913 • If assets are released, the license, copyright information, and terms of use in the package should  
914 be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for  
915 some datasets. Their licensing guide can help determine the license of a dataset.
- 916 • For existing datasets that are re-packaged, both the original license and the license of the derived  
917 asset (if it has changed) should be provided.
- 918 • If this information is not available online, the authors are encouraged to reach out to the asset's  
919 creators.

### 920 3. New Assets

921 Question: Are new assets introduced in the paper well documented and is the documentation pro-  
922 vided alongside the assets?

923 Answer: [NA]

924 Justification: We do not release new assets.

925 Guidelines:

- 926 • The answer NA means that the paper does not release new assets.
- 927 • Researchers should communicate the details of the dataset/code/model as part of their submis-  
928 sions via structured templates. This includes details about training, license, limitations, etc.

- 929 • The paper should discuss whether and how consent was obtained from people whose asset is  
930 used.
- 931 • At submission time, remember to anonymize your assets (if applicable). You can either create an  
932 anonymized URL or include an anonymized zip file.

933 **4. Crowdsourcing and Research with Human Subjects**

934 Question: For crowdsourcing experiments and research with human subjects, does the paper in-  
935 clude the full text of instructions given to participants and screenshots, if applicable, as well as  
936 details about compensation (if any)?

937 Answer: [NA]

938 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

939 Guidelines:

- 940 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
941 subjects.
- 942 • Including this information in the supplemental material is fine, but if the main contribution of the  
943 paper involves human subjects, then as much detail as possible should be included in the main  
944 paper.
- 945 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other  
946 labor should be paid at least the minimum wage in the country of the data collector.

947 **5. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Sub-**  
948 **jects**

949 Question: Does the paper describe potential risks incurred by study participants, whether such  
950 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or  
951 an equivalent approval/review based on the requirements of your country or institution) were ob-  
952 tained?

953 Answer: [NA]

954 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

955 Guidelines:

- 956 • The answer NA means that the paper does not involve crowdsourcing nor research with human  
957 subjects.
- 958 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be  
959 required for any human subjects research. If you obtained IRB approval, you should clearly state  
960 this in the paper.
- 961 • We recognize that the procedures for this may vary significantly between institutions and loca-  
962 tions, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their  
963 institution.
- 964 • For initial submissions, do not include any information that would break anonymity (if applica-  
965 ble), such as the institution conducting the review.