

Diversity-Enhanced and Classification-Aware Prompt Learning for Few-Shot Learning via Stable Diffusion

Anonymous authors

Paper under double-blind review

Abstract

Recent text-to-image generative models have exhibited an impressive ability to generate fairly realistic images from some text prompts. In this work, we explore to leverage off-the-shelf text-to-image generative models to train non-specific downstream few-shot classification model architectures using synthetic dataset to classify real images. Current approaches use hand-crafted or model-generated text prompts of text-to-image generative models to generate desired synthetic images, however, they have limited capability of generating diverse images. Especially, their synthetic datasets have relatively limited relevance to the downstream classification tasks. This makes them fairly hard to guarantee training models from synthetic images are efficient in practice. To address this issue, we propose a method capable of adaptively learning proper text prompts for the off-the-shelf diffusion model to generate diverse and classification-aware synthetic images. Our approach shows consistently improvements in various classification datasets, with results comparable to existing prompt designing methods. We find that replacing data generation strategy of existing zero/few-shot methods with proposed method could consistently improve downstream classification performance across different network architectures, demonstrating its model-agnostic potential for few-shot learning. This makes it possible to train an efficient downstream few-shot learning model from synthetic images generated by proposed method for real problems.

1 Introduction

Recently, deep learning powered by large-scale annotated data has achieved great success in the field of image recognition (He et al., 2016). However, acquiring and curating a large-scale high-quality dataset can be notoriously costly and time-consuming. This is especially challenging for inherently expensive domains, such as medical imaging, remote sensing, etc. Few-shot learning addresses the data issue by training a model using few data from the concerned tasks (Wang et al., 2020; Song et al., 2023; Shu et al., 2018). Generally, few-shot learning models use specialised algorithms and architectures to achieve the objective (Zhang et al., 2023b; Sendera et al., 2023; Zhou et al., 2023a; Baik et al., 2023; Zhang et al., 2023a; Shu et al., 2023c). This limits the variety of model architectures and potential applicability for real-world problems.

To address this limitation, some researches focus on generating images which are then used to train a classification model. It has been early explored based on GANs models (Besnier et al., 2020; Zhang et al., 2021; Jahanian et al., 2021), and become a research hotspot recently due to the rise of text-to-image diffusion models (Ho et al., 2020; Song & Ermon, 2019). Though these methods improve few-shot performance, some studies have found that synthetic datasets often have significant differences from real datasets and simply increasing the number of generated images cannot narrow the gap (Sariyildiz et al., 2023; He et al., 2022).

Fortunately, some researchers (Bansal & Grover, 2023; Shin et al., 2023; Burg et al., 2023a) explore to generate high-quality images for improving the results, and a potential path is to design proper text descriptions (prompts) for text-to-image generation diffusion models to generate desired synthetic images. A direct approach is to construct prompts by formatting class labels according to a template (called vanilla prompt (Sariyildiz et al., 2023; Radford et al., 2021)), such as “a photo of {class}”. To produce more diverse text descriptions, multi-domain prompt (Shipard et al., 2023) additionally provides a list of domains with

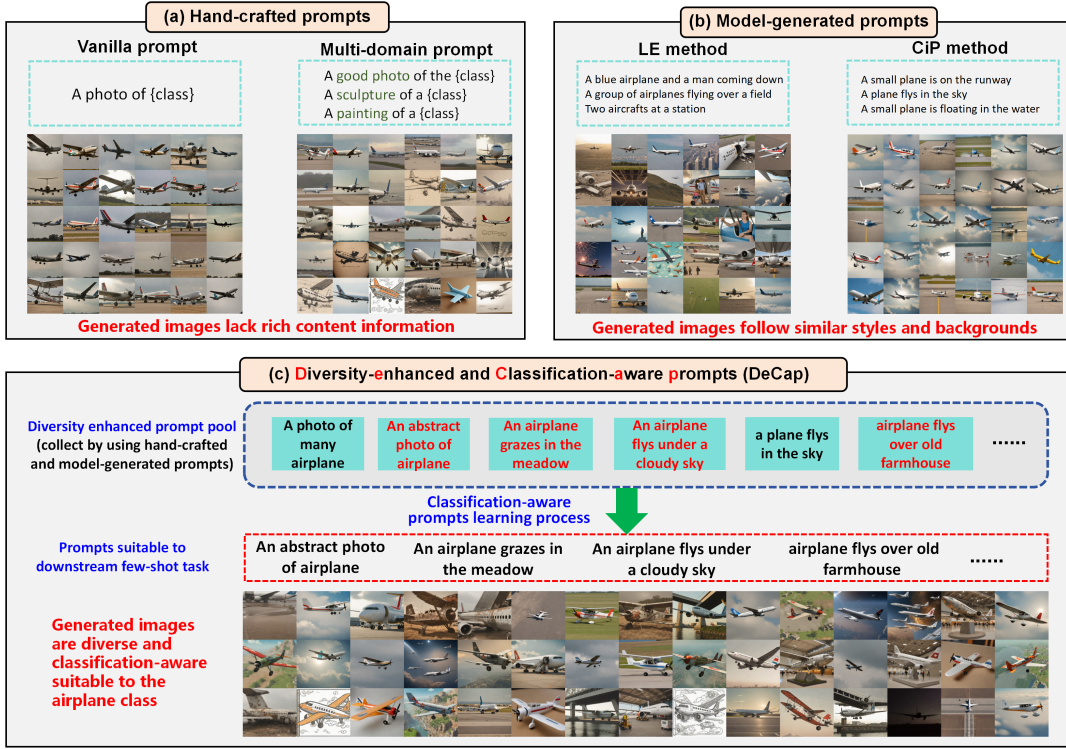


Figure 1: The comparison between existing prompt design methods and proposed DeCap method. Hand-crafted methods usually generate images with different domain information but limited content information. Model-generated methods overcome this shortcoming, while may generate images share similar patterns. DeCap constructs a diversity-enhanced prompt pool that contains potentially all-inclusive prompt information by integrating the advantages of hand-crafted and model-generated methods, and then carry out classification-aware prompt learning process to mine proper prompts suitable to few-shot tasks from the prompt pool. Figure shows the mined prompts for airplane classification.

the prompt, e.g., “a {domain} of a {class}”, to construct a set of prompt templates, in which ‘{domain}’ refers to drawing, painting, sketch, etc. However, these hand-crafted prompts have limited capacity of generating images with rich content information, which usually leads to inferior generalization performance when training downstream models. To improve the content quality of prompts, the language enhancement (LE) method (He et al., 2022) leverages an off-the-shelf word-to-sentence T5 model (Raffel et al., 2020) to automatically expand class names into various sentences with rich content descriptions, containing the class names as language prompts. While this method hardly considers the class-relevant visual information for image classification. The CiP method (Lei et al., 2023) generates high-quality prompts via extracting meaningful captions from real images using the off-the-shelf image captioning models such as BLIP2 (Li et al., 2022), showing a significant improvement in generating informative synthetic images.

Although prompts produced by off-the-shelf foundational models can help generate high-quality images, they still have evident deficiencies in practice. On the one hand, generated prompts tend to share fixed or similar patterns for different images as reported in Wang et al. (2023a), which may limit diversity of synthetic images. For example, as shown in Figure 1, images generated by LE and CiP methods usually follow the similar styles and backgrounds. This limitation, which is even more serious under few-shot setting, may cause subpopulation shift problem (Nagarajan et al., 2020; Yang et al., 2023a), i.e., some subpopulations of synthetic images shift from real-world datasets. On the other hand, existing prompt designing methods have relatively limited relevance to the downstream classification tasks. Generally, the generated text prompts only employ class names or class-relevant visual information, which leads to some noises in generated prompts, e.g., prompts containing noisy labels or additional negative class information

(please also see Figure 3). Therefore, it is relatively hard to guarantee that training models from synthetic images are efficient for downstream classification tasks, which tends to hinder their application effectiveness and reduce their performance stability in real problems.

To alleviate the aforementioned issues, this paper presents a **Diversity-enhanced and Classification-aware prompt (DeCap)** learning strategy to mine proper text prompts for downstream few-shot classification tasks (see Figure 1 for illustration). **Our main idea is to build a diversity prompt pool, and then select suitable prompts from it. In detail, we first combine existing hand-crafted diverse prompt templates and rich content prompt descriptions generated by off-the-shelf foundational models to construct a prompt pool containing potentially all-inclusive diverse prompt information. And then we propose a novel meta-learning approach to select proper prompts tailored for the few-shot learning task from the prompt pool.** The DeCap method involves two nested learning loops: an inner-loop to train a classification model using generated synthetic images, and an outer-loop to search suitable prompts for text-to-image foundational generative models that produce synthetic training images for the inner-level classification model. The given real few-shot images are employed to compute outer-loop meta-objective for helping achieve classification-aware prompt learning. Through iteratively ameliorating both prompts selection and classification model performance, our method is capable of mining proper prompts which are attained specifically suitable to concerned few-shot task.

In summary, we make the following four-fold contributions:

- (1) We propose an approach to automatically learn proper text prompts for text-to-image diffusion models, enabling the generation of diverse, classification-aware synthetic images for few-shot learning tasks in a meta-learning manner.
- (2) We demonstrate that proposed DeCap method could not only automatically discover suitable prompts tailored to the concerned tasks (e.g., Section 4.3.3), but also adaptively filter out noisy and low-quality prompts which are potentially detrimental to classification model (e.g., Fig. 3), by leveraging the proposed classification awareness meta-objective.
- (3) We observe that diverse visual effects and content styles of mined prompts are potentially beneficial for producing high-quality training images tailored to specific tasks (e.g., Table 3 and Appendix D), leading to consistent performance improvements across various datasets, as shown in Table 1. In contrast, existing methods often excel on specific datasets and relatively fixed prompts but struggle with broader adaptability due to their limited prompt flexibility.
- (4) We show that substituting the data generation strategy of existing zero/few-shot methods with our approach could further improve their performance across different algorithms and network architectures as shown in Table 2.

The paper is organized as follows. Section 2 shows the related work. Section 3 presents the proposed method. Section 4 demonstrates experimental results and the conclusion is finally made.

2 Related work

Text-to-Image Diffusion Model. Diffusion model (Ho et al., 2020; Song & Ermon, 2019) has emerged as a research hotspot in the field of image generation recently, due to their impressive generative capabilities. It achieves gradual matching from a Gaussian distribution to an image distribution by reversing the diffusion process. Recently, thanks to large-scale image-text paired datasets (Schuhmann et al., 2022) and the maturity of text-image foundation models such as CLIP, some state-of-the-art text-to-image diffusion models, including DALL-E (Ramesh et al., 2022), GLIDE (Nichol et al., 2022), Imagen (Saharia et al., 2022), and Stable Diffusion (Rombach et al., 2022), can produce a wide variety of highly realistic images, which has greatly propelled research in fields such as art (Wahid et al., 2023), style transfer (Yuan et al., 2022; Zhang et al., 2023e), image controlling (Ruiz et al., 2023; Zhang et al., 2023c; Gal et al., 2022), image editing (Bhattarai et al., 2020; Bodur et al., 2024a;b), data augmentation (Trabucco et al., 2023; Dunlap et al., 2024) etc. In this paper, we explore leveraging off-the-shelf diffusion models to generate high-quality synthetic images for downstream few-shot image recognition.

Synthetic Dataset for Image Recognition. In the early stages, some research (Besnier et al., 2020; Zhang et al., 2021; Jahanian et al., 2021) explored the role of synthetic datasets with GAN models, but their application scenarios are constrained due to the limited capabilities of early GANs. With the emergence of large-scale text-to-image generative models, recent studies have validated the utility of synthetic datasets at a large scale. For example, for classification tasks, Sariyildiz et al. (2023); Bansal & Grover (2023) train synthetic ImageNet datasets from scratch, He et al. (2022); Li et al. (2023) showing that CLIP (Radford et al., 2021) can boost performance from synthetic datasets. Tian et al. (2024) validates the outstanding performance of synthetic dataset using SimCLR and MAE models. In the field of object detection, Karazija et al. (2023a) utilizes the output results of generative model’s cross-attention layers as weak supervision for zero-shot object recognition. Additionally, synthetic datasets are also applied to addressing long-tail problems (Shin et al., 2023).

The data generation strategy could be roughly divided into two categories. One is fine-tuning based method (Azizi et al., 2023; Yuan et al., 2023), which fine-tunes generative models’ parameters using task data. These methods demonstrate strong domain adaptation capabilities on large-scale datasets and can effectively generate samples that conform to the distribution of real dataset. However, it often requires large-scale real datasets. Therefore, the other is prompt designing method to address few-shot learning. They don’t alter the parameters of generative models; instead, it focuses on setting proper prompts for off-the-shelf generative models to generate synthetic datasets. As discussed in Section 1, there exist two methodologies of setting prompts, i.e., hand-crafted and model-generated prompts. While they are not sufficient to generate high-quality images for classification. In this paper, we propose to integrate the advantages of both methodologies to achieve a diversity-enhanced and classification-aware prompt learning strategy. We need to clarify that, different from prompt learning methods (Zhou et al., 2022a;b) specifically designed for multimodal models like CLIP, which directly helps adjust off-the-shell models prediction adapting to the concerned data, we focuses on generating efficient synthetic data for further help train downstream few-shot learning.

Meta Learning. Meta learning (Hospedales et al., 2021; Shu et al., 2023a), also known as learning to learn, focuses on how to quickly adapt and apply previously acquired knowledge when faced with new learning tasks. Meta learning is widely used in few-shot learning (Finn et al., 2017; Shu et al., 2018; Ravi & Larochelle, 2016; Snell et al., 2017), hyperparameter optimization (Franceschi et al., 2018), transfer learning (Jang et al., 2019; Sun et al., 2020), label noise learning (Shu et al., 2019; 2023b; Wu et al., 2021), machine learning automation Xu et al. (2024), etc. For image generation field, meta learning is used to achieve data distillation (Nguyen et al., 2020; Wang et al., 2018; 2023c; Such et al., 2020), data augmentation (Yamaguchi et al., 2024), etc. Different from previous works updating parameters of generative model, we use meta learning to set proper text prompts for diffusion models to generate high-quality images for concerned few-shot learning task.

3 The Proposed DeCap Method

3.1 Preliminary

For a N -classification task, We use $\hat{x}_{ij}^{(k)} = g(\theta_{ij}, \epsilon_k)$ to denote the generated image $\hat{x}_{ij}^{(k)}$ via an off-the-shelf text-to-image foundational models g , where $i \in [N]$, $[N] = \{1, \dots, N\}$ represents the i -th class, $j \in [M]$, $[M] = \{1, \dots, M\}$, where M means how many different prompts for this class, θ_{ij} represents the prompt used to generate this image, $\epsilon_k, k = 1, 2, \dots, l$ represents random gaussian noise, where l means the generation number of each prompt. We denote the synthetic images generated by the prompts θ_{ij} as $X_{ij}^{syn} = \{\hat{x}_{ij}^{(k)}\}_{k=1}^l$. We only study prompt setting for image generation, and we will drop explicit dependence of X_{ij}^{syn} on ϵ_k for brevity in the following, i.e., $X_{ij}^{syn} = g(\theta_{ij})$. Our approach can be directly applied to different diffusion models, and we study the well-known open-sourced Stable Diffusion (SD) model in this work.

Considering a few-shot classification task with real data $D^{real} = \{(x_{ij}, y_{ij}), i = 1, \dots, N, j = 1, \dots, K\}$, where x_{ij}, y_{ij} denote image and its label, and N, K denote the number of classes and real images of each class, respectively. To boost few-shot model performance, it could use SD model to help generate high-quality

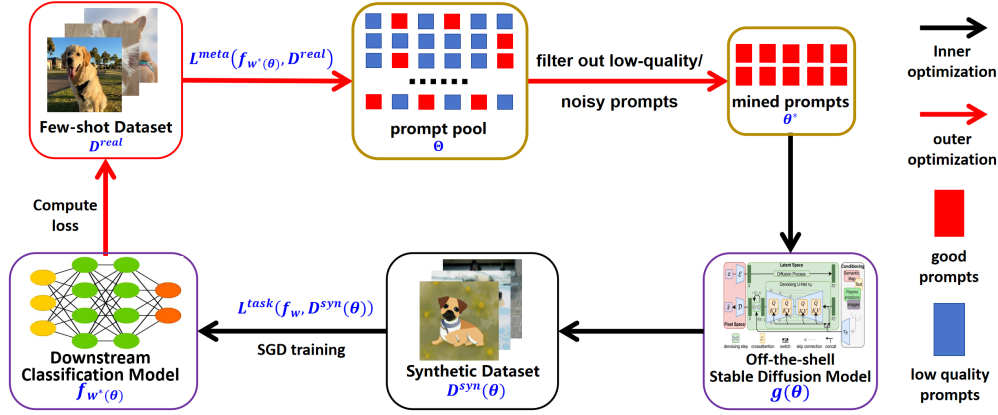


Figure 2: **Overview of the proposed DeCap method.** DeCap training involves two nested training loops. In the inner-loop optimization, we use the selected prompts set θ to generate synthetic dataset and then help train the classification model, while in the outer-loop optimization, we search proper prompts attained specifically suitable to few-shot learning task from pre-constructed prompt pool.

synthetic data for few-shot image recognition tasks. Specifically, the synthetic data could be formulated as

$$X^{syn} = g(\theta), \theta = \{\theta_i, i \in [N]\}, \theta_i = \{\theta_{i1}, \theta_{i2}, \dots, \theta_{iM}\}, \\ X^{syn} = \{X_{ij}^{syn} = g(\theta_{ij}), i \in [N], j \in [M]\}.$$

For simplicity, we denote the synthetic dataset as $D^{syn}(\theta) = \{X^{syn}(\theta), Y\}$, where $Y = \{y_{ij}, i \in [N], j \in [M]\}$. Based on $D^{syn}(\theta)$, we could train a classification network f_w by optimizing the following objective:

$$w^* = \arg \min_{w \in \mathcal{W}} \mathcal{L}^{task}(f_w, D^{syn}(\theta)), \quad (1)$$

where \mathcal{W} denotes parameter space, $\mathcal{L}^{task}(f_w, D^{syn}(\theta)) = \frac{1}{MN} \sum_{i=1}^N \sum_{j=1}^M \mathcal{L}^{task}(f_w(x_{ij}), y_{ij})$, and \mathcal{L}^{task} denote the cross-entropy loss for the few-shot classification task.

As discussed in Section 1, existing prompt design methods may produce a limited diversity of synthetic images, which can degrade the generalization performance of downstream classification models. Especially, we could see that the prompt construction process of existing methods has limited relevance to the downstream classification tasks from Eq.(1), i.e., it overlooks the explicit dependence of w^* on θ . In other words, existing prompt learning methods are potentially classification-agnostic and limited flexibility, which greatly reduces the alignment between synthetic datasets and downstream classification task requirement. To address these two issues, we propose a novel prompt learning strategy called DeCap, which explores to learn proper prompts for generating high-quality images to improve downstream few-shot learning task. We present the method and solving algorithm in Section 3.2 and 3.3, respectively.

3.2 Proposed DeCap Method

The proposed DeCap method firstly constructs a diversity-enhanced prompt pool (Section 3.2.1) by integrating the advantages of hand-crafted and model-generated methods, and then carry out classification-aware prompt learning process (Section 3.2.2) to mine prompts suitable to downstream few-shot classification task.

3.2.1 Diversity-Enhanced Prompt Pool Construction

The prompts constructed by previous methods often focus on specific aspects (see Figure 1). We proposed to integrate the advantages of both hand-crafted and model-generated methods to construct a prompt pool that contains potentially all-inclusive diverse prompt information. To some extent, such prompt pool cannot completely contain expected prompt information, while our experiments demonstrate that it is sufficient to

help generate high-quality images for boost few-shot performance. More comprehensive pool construction strategy is left for future research.

Specifically, we construct the prompt pool Θ containing hand-crafted prompts and model generated prompts, for every class in the dataset. For hand-crafted prompts, we first select some common prompt templates provided by Radford et al. (2021) which contain various domain information. Then we manually add some new prompts into the pool, covering aspects such as color, style, camera angle and so on. Since these prompts describe the object in general terms, we share these prompts for all classes. For model generated prompts, we use BLIP2 model as CiP method (Lei et al., 2023) to produce image captions based on real few-shot data, and utilize T5 model as LE method (He et al., 2022) to generate corresponding class prompts with class labels as information. These prompts describe the object in detail, so different classes will have totally different descriptions. In a nutshell, for each category’s prompt θ_i , it contains two parts: the hand-crafted prompt θ_i^h and the model-generated prompt θ_i^m , i.e., $\theta_i = [\theta_i^h, \theta_i^m]$, where all classes share the same template θ_i^h , while possess private prompt θ_i^m .

After conducting this process, there already exists potentially adequate prompts containing both diverse domain and content information in the prompt pool. We give a simple example about what our prompt pool looks like in Appendix A.2. However, this prompt pool is overly abundant and classification-agnostic, in other words, the pool contains not only proper prompts but also noisy prompts for specific few-shot learning task. An illustration of the necessity of prompt learning process please see Appendix C.1. Therefore, we further propose a classification-aware prompt learning strategy to mine suitable prompts tailored to the concerned few-shot task from the prompt pool in a meta-learning manner to help generate high-quality images.

3.2.2 Classification-Aware Prompt Learning

To establish the connection between prompt setting process and downstream classification model learning, a natural idea is to directly generate images approximating real few-shot image domain, e.g., SDEdit (Meng et al., 2021) and Textual inversion (Gal et al., 2022). However, these methods tend to generate images that are very similar to real few-shot data, which is prone to overfitting, making it hard to generalize to unseen images (Please see Appendix C.3 and Table 4). To overcome this challenge, we explore to utilize real few-shot data at the higher level, i.e., evaluating the classification performance on these real data, rather than focusing on achieving visual similarity to the real data. To achieve this goal, we formulate the prompt learning as the following bi-level optimization objective:

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{L}^{meta}(f_{w^*(\theta)}, D^{real}), \quad (2)$$

$$\text{where } w^*(\theta) = \arg \min_{w \in \mathcal{W}} \mathcal{L}^{task}(f_w, D^{syn}(\theta)), \quad (3)$$

where the inner-level objective (Eq.(3)) is the same as Eq.(1), while the outer-level meta-objective \mathcal{L}^{meta} is computed by evaluating the performance of obtained classification model in Eq.(2) on real few-shot data D^{real} , which utilizes real data to guide the prompt learning at higher level. Note that existing methods predefine the prompts while overlook the explicit dependence of classification model w^* on prompts θ . As a comparison, we explicitly require the performance of classification model to depend on the prompts θ . Specifically, given a prompt set $\theta \in \Theta$, we use these prompts to obtain the synthetic dataset $D^{syn}(\theta)$, and then train the downstream classification model $w^*(\theta)$ on the synthetic dataset. **With such dependency, the classification-aware meta-objective then help to search proper prompts from prompt pool, which aims to improve the classification performance on the real few-shot data.**

3.3 Learning Algorithm of the DeCap Method

To solve above bi-level optimization objective, through iteratively ameliorating both searching prompts at outer-level learning and classification model performance at inner-level learning, our algorithm is capable of mining classification-aware prompts which is attained specifically suitable to downstream few-shot task.

In our implementation, the optimization of $\theta \in \Theta$ is actually a discrete prompt selection problem since Θ represents a discrete prompt pool. Therefore, we use the genetic algorithm (GA) (Katoch et al., 2021)

Algorithm 1 Learning Algorithm of the DeCap Method

Input: Downstream few-shot learning task real dataset D^{real} ; Algorithm iteration number $max\text{-}iter$, population quantity $popsize$; Constructed prompt pool $pool$; off-the-shell text-to-image diffusion model g

Output: Optimal prompt set θ^*

```

1: GA.initial(max-iter,popsize)
2: for  $iter = 1, 2, \dots, max\text{-}iter$  do
3:    $fit = []$ ,  $Pop = []$ 
4:   for  $m = 1, 2, \dots, popsize$  do
5:     # An individual of population
6:      $pop^{(m)} = GA.sample()$ 
7:     # See Algorithm 2 in Appendix A.4
8:      $\theta^{(m)}$ ,  $Y^{syn} = get\_prompt(pop^{(m)}, pool)$ 
9:      $X^{syn}(\theta^{(m)}) = g(\theta^{(m)})$ 
10:    # Eq.(3)
11:     $w^*(\theta^{(m)}) = \arg \min \mathcal{L}^{task}(f_w, D^{syn}(\theta^{(m)}))$ 
12:    # Eq.(2)
13:     $fit_{pop^{(m)}} = \mathcal{L}^{meta}(f_{w^*(\theta^{(m)})}, D^{real})$ 
14:     $fit.append(fit_{pop^{(m)}})$ ,  $Pop.append(pop^{(m)})$ 
15:   end for
16:    $GA.update(fit, Pop)$ 
17: end for
18: return  $GA.best$ 

```

to solve the outer-level optimization objective in Eq.(2). Generally speaking, a genetic algorithm first generates different inputs, then obtains the corresponding value function outputs for these inputs, adjusts the search direction based on the magnitude of the outputs, and eventually completes the optimization process. Therefore, we only need to define the GA’s input and value function for DeCap objective, and then genetic algorithm can be employed to mine proper prompts θ^* from constructed prompt pool Θ in Section 3.2.1.

In our implementation, the input of GA algorithm is defined as a vector of integers. The length of the vector represents the number of prompts selected, and each dimension of the vector corresponds to the index of the selected prompt, with values ranging from 0 to the size of the prompt pool. Under this definition, each input represents a different combination of selected prompts. The value function is defined as the outer-level meta-objective \mathcal{L}^{meta} in Eq.(2). The whole learning algorithm of proposed DeCap method is summarized in Algorithm 1. More details about genetic algorithm please see Appendix A.3. Moreover, we also discuss the selection of optimization methods and the implementation details for DeCap method at Appendix A.1.

4 Experimental Results

4.1 Few-Shot Classification Performance

We compared with existing prompt designing strategies including: (1) vanilla prompt (Radford et al., 2021): using the template “a photo of {class}”. (2) multi-domain prompt: using different text templates from domains provided in Radford et al. (2021). (3) LE (He et al., 2022): using language models for text prompt construction, where the input and output of language models are the class label and a sentence containing the class label, respectively. (4) CiP (Lei et al., 2023): generating captions for real image data using the BLIP2¹ model. We conduct experiments on 12 datasets: CIFAR10, STL-10, Imagenette, Pets, Caltech-101, ImageNet100, EuroSAT, FGVC Aircraft, Country211, DTD, UCF101, Imagenet. Datasets details are introduced in Appendix B.1. For the selection of the classification model, we use the CLIP model, as it has shown powerful classification ability. The training strategy we used strictly follows the settings described

¹<https://huggingface.co/Salesforce/blip2-opt-2.7b>

Table 1: Top-1 classification accuracy on different datasets. **We randomly selected 10 images per class to construct the few-shot datasets.** **Bold scores** represent the best results on each datasets, and the second best scores are marked by underline.

	STL-10	CIFAR10	Imagenette	Pets	Caltech-101	Imagenet100	EuroSAT	Aircraft	Country211	DTD	UCF101	Imagenet
<i>without real</i>												
zero-shot	94.26	70.25	97.22	81.85	83.89	70.14	23.11	17.07	13.44	42.39	60.80	58.18
vanilla prompt	95.33	<u>72.37</u>	97.69	82.29	84.74	70.62	31.31	17.04	13.72	49.20	60.83	58.45
multi-domain	94.97	70.66	<u>97.89</u>	83.07	87.56	70.50	30.11	17.85	13.90	49.63	60.83	58.64
LE	94.61	70.33	97.45	83.24	84.03	70.73	29.35	17.73	14.14	52.50	61.06	58.86
CiP	94.92	70.24	97.65	<u>84.04</u>	<u>88.12</u>	<u>70.76</u>	<u>39.91</u>	<u>18.00</u>	<u>14.98</u>	<u>60.69</u>	<u>63.63</u>	<u>59.75</u>
DeCap (ours)	95.91	76.98	97.95	85.36	88.67	71.08	41.94	19.74	15.44	62.39	64.45	60.29
<i>with real</i>												
real-only	94.28	70.33	97.22	81.96	84.51	70.33	24.15	19.41	13.80	42.55	61.09	58.51
vanilla prompt	<u>95.55</u>	<u>76.20</u>	<u>98.00</u>	83.84	89.85	<u>70.87</u>	47.83	18.21	13.77	51.97	60.85	59.12
multi-domain	95.02	74.54	97.92	84.56	90.31	70.62	43.30	18.99	13.90	51.76	60.90	59.17
LE	94.72	71.66	97.49	84.00	84.34	70.46	42.06	20.22	14.37	54.04	61.54	59.85
CiP	95.05	70.51	97.75	<u>85.16</u>	89.86	70.86	<u>49.17</u>	<u>20.31</u>	<u>15.49</u>	62.66	<u>64.66</u>	<u>62.51</u>
DeCap (ours)	95.93	77.19	98.03	85.78	<u>89.87</u>	71.11	50.22	20.64	15.68	63.99	65.21	62.82

in He et al. (2022), where we finetune CLIP with generated data. We use “a photo of {class}” as the text initialization for CLIP tuning for all datasets to eliminate the impact of different initializations on the evaluation of each method. Training and evaluating details are presented in Appendix B.2. Table 1 shows the few-shot classification performance of each method on nine downstream few-shot learning datasets, where “without real” means that we only use synthetic images to train downstream models, while “with real” means that we use both synthetic and real few-shot images. Some ablation studies and analyses on proposed DeCap method please see Appendix C.3.

Using synthetic data to train downstream classification model, all prompt designing methods can improve CLIP zero-shot performance, showing that generating synthetic data is helpful to boost downstream classification model, and proposed DeCap method demonstrates the best classification accuracies across all datasets. Note that existing methods often use relatively fixed-pattern prompts, which exhibits limitations in adapting to diverse classification tasks. As for datasets with simple categories like STL-10, CIFAR-10 and Imagenette, the hand-crafted prompts could achieve superior performance than model-generated prompts, illustrating that the prompts with only class/domain information may be relatively more proper for these tasks; while for datasets with complex categories like Pets, Caltech-101, Imagenet100, EuroSAT, Aircraft and Country211 datasets, the model-generated prompts could achieve better performance than hand-crafted prompts, implying that rich content information is more helpful to address these complex tasks. These results reveal that effective prompts should be set considering the characteristic of concerned tasks.

To address this challenge, proposed DeCap method could adaptively learn proper prompts suitable to the concerned tasks by reconciling class/domain information and rich content information (visualization of mined prompts see Appendix D.2), so as to generate better training images. On the whole, our method achieves an average performance improvement of 1.30% point compared to the best results of existing method on different datasets. Other metrics evaluating few-shot performance including precision, recall and F1-score, are also reported in Appendix C.4, and our method delivers consistent improvements. We also evaluate the adversarial robustness of these methods in Appendix C.5, which further substantiate the high-quality data generation capability of proposed DeCap method.

When using additional real data, CLIP’s performance could be further improved, though the number of real data is relatively smaller than synthetic data. This implies that the quality of real data may be higher than that of synthetic data. All prompt designing methods obtain a further improvement over only using synthetic data. Even so, DeCap method still shows advantages over other methods on the most datasets (only one dataset slightly lower than best method), demonstrating that our approach could genuinely generate

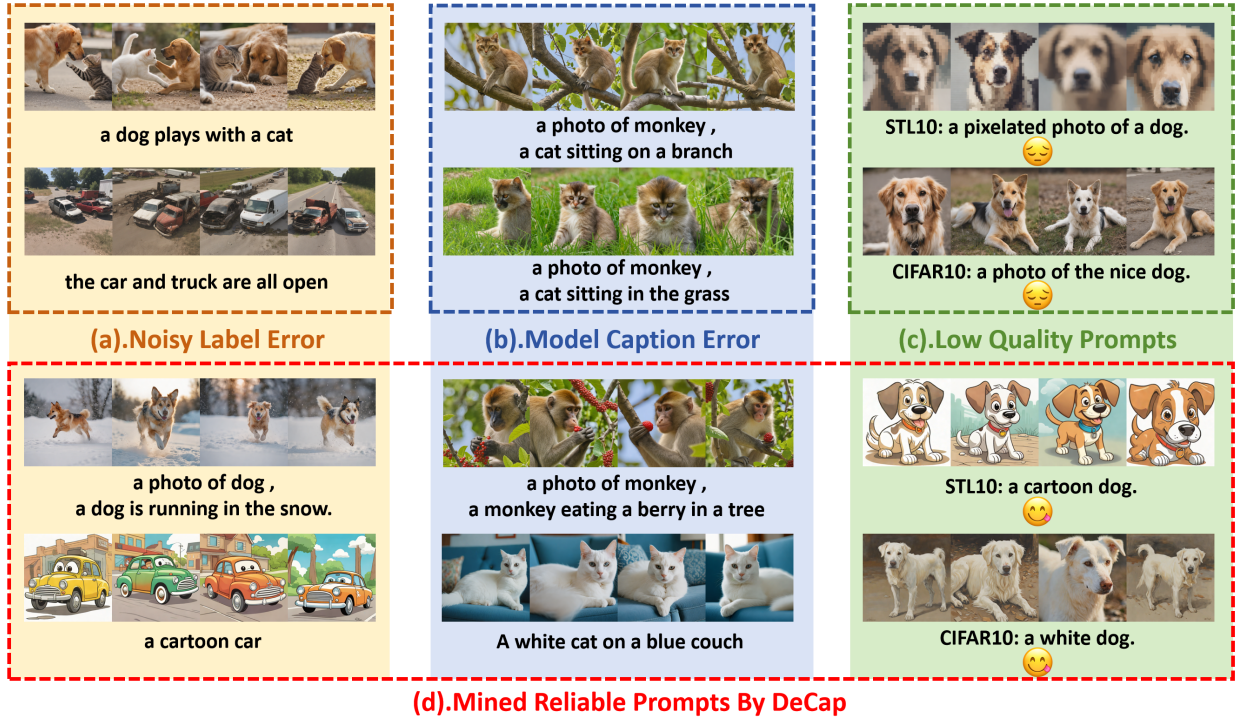


Figure 3: Illustration of (a) noisy label, (b) model caption error or (c) low quality prompts in prompt pool generated by existing hand-crafted and model generated prompt designing methods, and (d) mined reliable prompts by our DeCap method.

high-quality images as a supplement to real few-shot data to further boost the classification performance. These experimental results substantiate the capability of proposed DeCap method in generating high-quality images for boosting few-shot classification performance.

Notice that the CiP method demonstrates strong performance on some tasks by leveraging real data information, while it may be limited effectiveness in utilizing real data via approximating visual similarity. As a comparison, proposed DeCap method utilizes real data at a higher meta-level to facilitate the trained classifier based on synthetic images to perform well on unseen real images. As a result, our method achieves consistently improve over CiP method. To further illustrate this point, we additionally compared two commonly used image augmentation methods including SDEdit (Meng et al., 2021) and Textual inversion (Gal et al., 2022) in the Appendix C.2, which tend to produce similar images based on given real few-shot images. In Appendix.C.7, we demonstrate that even synthetic images with significant difference with real images in terms of visual similarity (see Figure 8), can still bring much improvements to classification performance. These results further show that Decap method exhibits better flexibility and adaptability in prompt setting via utilizing real few-shot images at the meta-level.

4.2 Comparison with SOTA Methods

In Section 4.1, we have shown that DeCap method performs better than other prompt designing methods under the same CLIP classification model. The key goal of DeCap method is to mine proper prompts to generate high-quality data for improving few-shot learning, while it is not confined to specialised algorithms and architectures to complete few-shot learning tasks. To illustrate this, we explore to use synthetic data of DeCap method to evaluate its capacity in improving other zero/few-shot algorithms and architectures.

Table 2: Comparison of DeCap and SOTA methods on different datasets. “Method + DeCap” denotes the performance of replacing original synthetic data strategies of each method with proposed DeCap method without altering any of model architectures for a fair comparison.

	STL-10	CIFAR10	Imagenette	Pets	Caltech-101	Imagenet100	EuroSAT	Aircraft	Country211	DTD	UCF101	Imagenet
FakeIt	52.26	38.45	69.60	29.74	66.20	32.75	48.40	37.70	3.61	31.25	33.30	15.46
With DeCap	60.39	48.80	75.40	55.22	70.51	39.21	51.20	40.60	4.22	41.67	41.07	21.59
SuS-X	95.24	72.77	98.24	79.64	84.57	69.96	33.89	18.30	12.97	46.65	57.18	57.84
With DeCap	95.43	75.89	98.39	80.40	84.89	70.3	37.37	19.83	13.63	56.49	61.8	59.54
CaFo	95.33	85.34	97.66	86.62	94.09	74.64	83.5	26.07	16.20	65.74	72.69	63.41
With DeCap	95.90	86.00	98.06	88.66	94.28	76.28	85.59	32.10	16.88	68.03	75.02	64.29

Specifically, we conducted our experiments on three SOTA algorithms: (1) FakeIt (Sariyildiz et al., 2023): It uses synthetic datasets to train the network on ResNet-50. (2) SuS-X (Udandara et al., 2023): It leverages synthetic datasets as a dynamic support set and extends Tip-Adapter by utilizing the image-text distance. (3) CaFo (Zhang et al., 2023d): It augments few-shot datasets with synthetic data and then combines the predictions of pre-trained CLIP and DINO. In our implementations, we replaced the data generation strategies of these methods with DeCap without altering any of model architectures for a fair comparison, and follow original settings of these methods to train the corresponding classification models. We provide experimental details in Appendix B.3, and more results in Appendix C.6.

Table 2 reports the results. Notice that FakeIt method uses synthetic data to train the ResNet-50 model from scratch, which eliminates the effect of pre-training data for downstream classification tasks. Thus the performance of trained ResNet-50 model could appropriately reflect the quality of synthetic data. The DeCap method achieves a significant improvement of 7.43% points over original data generation strategy of FakeIt, substantiating that our method is classifier-agnostic and capable of generating high-quality data suitable to concerned tasks. Though SuS-X and CaFo methods use pre-trained models, synthetic data of DeCap method could still outperform these methods in the vast majority of datasets. These results demonstrate that synthetic data of our DeCap method are not confined to specialised algorithms and classifier architectures. This implies that our DeCap method is algorithm-agnostic for improving downstream few-shot learning tasks, and hopeful to be readily applied to real-world problems and tasks.

4.3 Why Proposed DeCap Method Perform Better?

In this section, we further present some analysis of DeCap method in two aspects: robustness against noisy or low quality prompts, and data value analysis of synthetic data.

4.3.1 Robustness against noisy or low quality prompts

Existing prompts methods may set prompts for specific classes, which are potentially noisy or low quality for downstream few-shot tasks. E.g., for LE method, it may generate prompts that contain not only the class we want, but also other classes information in the dataset. An illustrated example is presented in Fig.3 (a): for STL10 dataset, when we generate images for “dog”/“car” classes, some images also contain information of “cat”/“truck” classes. Since “cat”/“truck” classes belong to the dataset, these prompts would generate images with noisy labels for the classification of “dog”/“car”. For CiP method, due to the limitations of the BLIP2 model’s capability, it cannot always accurately annotate images, which may result in misidentifications. Although CiP method recognizes this issue and employs a prompt concatenation method like “a photo of {class}, {image caption}” to reduce the influence of noisy captions, we found this may not always work. For example, as shown in Fig.3 (b), when the BLIP2 model mistakenly identifies a monkey as a cat, the defined prompt “a photo of monkey, a cat sitting in a branch” may generate an image that blending features of cat and monkey. The issue of misidentification is particularly prominent in certain tasks, such as CIFAR10, where the low resolution images significantly impact the model’s judgments. This explains why the CiP method performs poorly on CIFAR10 dataset, as presented in Table 1. Moreover,




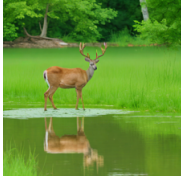

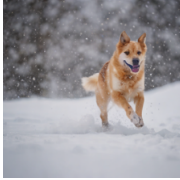



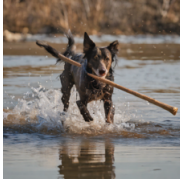

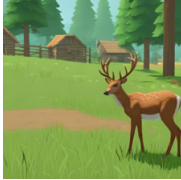
a doodle of the car. 	dog walking on a sunny day 	a photo of the bird. 	deer on a green pond. 
olympic athletes racing cars during racing match. 	a photo of dog , a dog is running in the snow 	a embroidered bird. 	a photo of the clean deer. 
A dynamic car. 	a photo of dog, a dog plays in water with stick 	art of the bird. 	a cartoon deer. 

Table 3: Examples of synthetic images generated by DeCap method for STL-10 dataset.



Figure 4: Illustration of the distributions of learned hand-crafted prompts by DeCap method on (a) CIFAR-10 and (b) STL-10 datasets. Each column represents the same template, and each row indicates which prompts were selected for each class. Black indicates prompt is not selected, orange indicates the prompt is selected once, and white indicates the prompt is selected more than once. For clarity, we removed prompts that were never selected.

hand-crafted prompts often introduce different domain information to construct diverse prompts. Generally, only part of domain information is reliable, while an amount of domain information may be of low quality for the concerned tasks. As shown in Fig.3 (c), though both of prompts could generate images of dog, the improper domain information could hinder the performance of concerned classification models, e.g., the synthetic pixelated images may provide low-quality training data for STL-10 task. In Appendix C.7, we further illustrate influence of prompts with domain information on the synthetic images.

Unfortunately, these noisy prompts are relatively hard to be filtered using data cleaning strategies such as CLIP filtering (He et al., 2022). To address the issue, proposed DeCap method aims to mine proper prompts suitable to the concerned classification task in a meta-learning manner. As shown in Fig.3 (d), with such higher-level downstream classification-aware outer-loop supervised information, DeCap method could adaptively select effective prompts that help boost downstream classification performance, and discard aforementioned potential noisy prompts that would potentially hurt downstream classification performance.

4.3.2 Data Value Analysis of synthetic data

To better analyze why DeCap method outperforms existing prompt designing methods, we use “leave-one-out” method (Ghorbani & Zou, 2019) to evaluate data valuation, and then select typical high-quality images generated by DeCap method. Specifically, given a dataset S and a measure function V , we use

$\phi_i = V(D \cup \{i\}) - V(D)$ to represent data valuation of the synthetic image i . In our implementation, we use the dataset generated by vanilla prompt method as the benchmark dataset S and classification accuracy as the measure function V . Then we could compute data valuation of synthetic images generated by DeCap method via adding one image at a time. Table 3 visualizes the synthetic images with high data valuations for STL-10 dataset, and more visualizations are shown in Appendix D.1. As it shown, we can see that synthetic images contain various patterns such as image style, background, camera angles, and actions, providing novel, diverse, and meaningful content information for given real few-shot data. This indicates that DeCap method does mine proper diverse and rich content prompts suitable to concerned few-shot learning tasks, naturally leading to its better accuracies than other prompt designing methods.

4.3.3 Visualization of Learned Prompts

Fig. 11 and 12 in Appendix D.2 show that proposed method could adaptively adjust the proportions of model-generated and hand-crafted prompts, so as to reconcile class/domain information and rich content information for different classes. The distributions of learned hand-crafted prompts are shown in Fig.4, we can find that domain information required by each class are distinctly different, which implies that proposed method could adaptively learning suitable prompts for different classes of specific few-shot tasks. Table 15 shows that consistently selected prompts are with high diversity and fine-grained information, including: movement, posture, background, color, quantity, other objects, and so on. More analysis please refer to Appendix D.2. These results substantiate the capability of DeCap method on adaptively mining suitable prompts tailored to the concerned tasks, naturally leading to enhancement of classification performance.

5 Conclusion

We present the DeCap, a novel adaptive prompt learning approach to generate diverse and classification-aware synthetic data for downstream few-shot learning in a meta-learning manner. Proposed DeCap method could mine potential reliable prompts suitable to downstream few-shot learning tasks, demonstrating impressive capabilities in improving downstream classification models for different few-shot learning tasks compared with existing prompt designing methods. We could further boost existing SOTA zero/few-shot learning methods by simply replacing data generation strategy with the proposed method, showing its potential model-agnostic characteristics. Besides, we also provide some intuitive visual interpretation, providing an initial insight into proposed DeCap method. Such an adaptive prompt learning approach based on classification-aware meta-objective is promising to establish the connection between pixel-level data generation and image-level semantic understanding, and hopeful to be employed to other computer vision tasks, like semantic segmentation and object detection, etc.

Impact Statement

In this work, we study the potential high-level computer vision applications of the popular diffusion models. We believe that only using the low-level visual similarity cannot effectively help complete high-level vision tasks. The classification-aware meta-objective provide a novel perspective that establishes the connection between low-level vision and high-level vision, aiming to advance the field of machine learning and computer vision. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J Fleet. Synthetic data from diffusion models improves imagenet classification. *arXiv preprint arXiv:2304.08466*, 2023.

- Sungyong Baik, Myungsub Choi, Janghoon Choi, Heewon Kim, and Kyoung Mu Lee. Learning to learn task-adaptive hyperparameters for few-shot learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets. In *ICLR 2023 Workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*, 2023. URL <https://openreview.net/forum?id=LjGqAFP6rA>.
- Victor Besnier, Himalaya Jain, Andrei Bursuc, Matthieu Cord, and Patrick Pérez. This dataset does not exist: training models from generated images. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2020.
- Binod Bhattarai, Seungryul Baek, Rumeysa Bodur, and Tae-Kyun Kim. Sampling strategies for gan synthetic data. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2303–2307. IEEE, 2020.
- Rumeysa Bodur, Binod Bhattarai, and Tae-Kyun Kim. Prompt augmentation for self-supervised text-guided image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8829–8838, 2024a.
- Rumeysa Bodur, Erhan Gundogdu, Binod Bhattarai, Tae-Kyun Kim, Michael Donoser, and Loris Baz-zani. iedit: Localised text-guided image editing with weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7426–7435, 2024b.
- Max F Burg, Florian Wenzel, Dominik Zietlow, Max Horn, Osama Makansi, Francesco Locatello, and Chris Russell. A data augmentation perspective on diffusion models and retrieval. *arXiv*, 2023a.
- Max F Burg, Florian Wenzel, Dominik Zietlow, Max Horn, Osama Makansi, Francesco Locatello, and Chris Russell. Image retrieval outperforms diffusion models on data augmentation. *Transactions on Machine Learning Research*, 2023b.
- Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, 2011.
- Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 113–123, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Lisa Dunlap, Alyssa Umino, Han Zhang, Jiezhi Yang, Joseph E Gonzalez, and Trevor Darrell. Diversify your vision datasets with automatic diffusion-based augmentation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transactions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazzi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *International conference on machine learning*, pp. 1568–1577. PMLR, 2018.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.

- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016.
- Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In *International conference on machine learning*, pp. 2242–2251. PMLR, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. Is synthetic data from generative models ready for image recognition? *arXiv preprint arXiv:2210.07574*, 2022.
- Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5149–5169, 2021.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Jeremy Howard and Sylvain Gugger. Fastai: a layered api for deep learning. *Information*, 11(2):108, 2020.
- Khawar Islam, Muhammad Zaigham Zaheer, Arif Mahmood, and Karthik Nandakumar. Diffusemix: Label-preserving data augmentation with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27621–27630, 2024.
- Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. Generative models as a data source for multiview representation learning. *arXiv preprint arXiv:2106.05258*, 2021.
- Yunhun Jang, Hankook Lee, Sung Ju Hwang, and Jinwoo Shin. Learning what and where to transfer. In *International conference on machine learning*, pp. 3030–3039. PMLR, 2019.
- Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pp. 709–727. Springer, 2022.
- Laurynas Karazija, Iro Laina, Andrea Vedaldi, and Christian Rupprecht. Diffusion models for zero-shot open-vocabulary segmentation. *arXiv preprint arXiv:2306.09316*, 2023a.
- Laurynas Karazija, Iro Laina, Andrea Vedaldi, and Christian Rupprecht. Diffusion Models for Zero-Shot Open-Vocabulary Segmentation. *arXiv preprint*, 2023b.
- Sourabh Katoch, Sumit Singh Chauhan, and Vijay Kumar. A review on genetic algorithm: past, present, and future. *Multimedia tools and applications*, 80:8091–8126, 2021.
- Hoki Kim. Torchattacks: A pytorch repository for adversarial attacks. *arXiv preprint arXiv:2010.01950*, 2020.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *Artificial intelligence safety and security*, pp. 99–112. Chapman and Hall/CRC, 2018.
- Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. doi: 10.1126/science.aab3050. URL <https://www.science.org/doi/abs/10.1126/science.aab3050>.
- Shiye Lei, Hao Chen, Sen Zhang, Bo Zhao, and Dacheng Tao. Image captions are natural prompts for text-to-image models. *arXiv preprint arXiv:2307.08526*, 2023.
- Bo Li, Haotian Liu, Liangyu Chen, Yong Jae Lee, Chunyuan Li, and Ziwei Liu. Benchmarking and analyzing generative data for visual recognition. *arXiv preprint arXiv:2307.13697*, 2023.
- Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE international conference on computer vision*, pp. 5542–5550, 2017.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Weixin Liang, Girmaw Abebe Tadesse, Daniel Ho, Li Fei-Fei, Matei Zaharia, Ce Zhang, and James Zou. Advances, challenges and opportunities in creating data for trustworthy ai. *Nature Machine Intelligence*, 4(8):669–677, 2022.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013.
- Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021.
- Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*, 2020.
- Timothy Nguyen, Zhouong Chen, and Jaehoon Lee. Dataset meta-learning from kernel ridge-regression. *arXiv preprint arXiv:2011.00050*, 2020.
- Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *ICML*, 2022.
- Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- Ed Pizzi, Sreya Dutta Roy, Sugosh Nagavara Ravindra, Priya Goyal, and Matthijs Douze. A self-supervised descriptor for image copy detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14532–14542, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. 2022.
- Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International conference on learning representations*, 2016.
- Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pp. 5389–5400. PMLR, 2019.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510, 2023.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
- Mert Bulent Sariyildiz, Karteek Alahari, Diane Larlus, and Yannis Kalantidis. Fake it till you make it: Learning transferable representations from synthetic imagenet clones. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. *arXiv preprint arXiv:2311.17042*, 2023.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- Marcin Sendera, Marcin Przewięzlikowski, Konrad Karanowski, Maciej Zięba, Jacek Tabor, and Przemysław Spurek. Hypershot: Few-shot learning by kernel hypernetworks. In *WACV*, 2023.
- Joonghyuk Shin, Minguk Kang, and Jaesik Park. Fill-up: Balancing long-tailed data with generative models. *arXiv preprint arXiv:2306.07200*, 2023.
- Jordan Shipard, Arnold Wiliem, Kien Nguyen Thanh, Wei Xiang, and Clinton Fookes. Diversity is definitely needed: Improving model-agnostic zero-shot classification via stable diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 769–778, June 2023.
- Jun Shu, Zongben Xu, and Deyu Meng. Small sample learning in big data era. *arXiv preprint arXiv:1808.04572*, 2018.
- Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weight-net: Learning an explicit mapping for sample weighting. *Advances in neural information processing systems*, 32, 2019.
- Jun Shu, Deyu Meng, and Zongben Xu. Learning an explicit hyper-parameter prediction function conditioned on tasks. *J. Mach. Learn. Res.*, 24:1–74, 2023a.

- Jun Shu, Xiang Yuan, Deyu Meng, and Zongben Xu. Cmw-net: Learning a class-aware sample weighting mapping for robust deep learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023b.
- Yang Shu, Zhangjie Cao, Jinghan Gao, Jianmin Wang, S Yu Philip, and Mingsheng Long. Omni-training: bridging pre-training and meta-training for few-shot learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023c.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- Kihyuk Sohn, Huiwen Chang, José Lezama, Luisa Polania, Han Zhang, Yuan Hao, Irfan Essa, and Lu Jiang. Visual prompt tuning for generative transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19840–19851, 2023.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Yisheng Song, Ting Wang, Puyu Cai, Subrota K Mondal, and Jyoti Prakash Sahoo. A comprehensive survey of few-shot learning: Evolution, applications, challenges, and opportunities. *ACM Computing Surveys*, 55(13s):1–40, 2023.
- Felipe Petroski Such, Aditya Rawal, Joel Lehman, Kenneth Stanley, and Jeffrey Clune. Generative teaching networks: Accelerating neural architecture search by learning to generate synthetic training data. In *International Conference on Machine Learning*, pp. 9206–9216. PMLR, 2020.
- Qianru Sun, Yaoyao Liu, Zhaozheng Chen, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning through hard tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(3):1443–1456, 2020.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pp. 776–794. Springer, 2020.
- Yonglong Tian, Lijie Fan, Phillip Isola, Huiwen Chang, and Dilip Krishnan. Stablerep: Synthetic images from text-to-image models make strong visual representation learners. *Advances in Neural Information Processing Systems*, 36, 2024.
- Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective data augmentation with diffusion models. *arXiv preprint arXiv:2302.07944*, 2023.
- Vishaal Udandaraao, Ankush Gupta, and Samuel Albanie. Sus-x: Training-free name-only transfer of vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2725–2736, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. *The Caltech-UCSD Birds-200-2011 Dataset*. 7 2011.
- Risqo Wahid, Joel Mero, and Paavo Ritala. Written by chatgpt, illustrated by midjourney: generative ai for content marketing. *Asia Pacific Journal of Marketing and Logistics*, 35(8):1813–1822, 2023.
- Alex Jinpeng Wang, Kevin Qinghong Lin, David Junhao Zhang, Stan Weixian Lei, and Mike Zheng Shou. Too large; data reduction for vision-language pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3147–3157, 2023a.

- Boyu Wang, Jorge A Mendez, Changjian Shui, Fan Zhou, Di Wu, Gezheng Xu, Christian Gagné, and Eric Eaton. Gap minimization for knowledge sharing and transfer. *Journal of Machine Learning Research*, 24 (33):1–57, 2023b.
- Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pp. 10506–10518, 2019.
- Kai Wang, Jianyang Gu, Daquan Zhou, Zheng Zhu, Wei Jiang, and Yang You. Dim: Distilling dataset into generative model. *arXiv preprint arXiv:2303.04707*, 2023c.
- Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM computing surveys*, 53(3):1–34, 2020.
- Zhicai Wang, Longhui Wei, Tan Wang, Heyu Chen, Yanbin Hao, Xiang Wang, Xiangnan He, and Qi Tian. Enhance image classification via inter-class image mixup with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17223–17233, 2024.
- Yichen Wu, Jun Shu, Qi Xie, Qian Zhao, and Deyu Meng. Learning to purify noisy labels via meta soft label corrector. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 10388–10396, 2021.
- Zongben Xu, Jun Shu, and Deyu Meng. Simulating learning methodology (slem): an approach to machine learning automation. *National Science Review*, 11(8):nwae277, 2024.
- Shin’ya Yamaguchi, Daiki Chijiwa, Sekitoshi Kanai, Atsutoshi Kumagai, and Hisashi Kashima. Regularizing neural networks with meta-learning generative models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer look at subpopulation shift. In *ICML*, 2023a.
- Zuhao Yang, Fangneng Zhan, Kunhao Liu, Muyu Xu, and Shijian Lu. Ai-generated images as data source: The dawn of synthetic era. *arXiv preprint arXiv:2310.01830*, 2023b.
- Zhuoran Yu, Chenchen Zhu, Sean Culatana, Raghuraman Krishnamoorthi, Fanyi Xiao, and Yong Jae Lee. Diversify, don’t fine-tune: Scaling up visual recognition training with synthetic images. *arXiv preprint arXiv:2312.02253*, 2023.
- Jianhao Yuan, Francesco Pinto, Adam Davies, and Philip Torr. Not just pretty pictures: Toward interventional data augmentation using text-to-image generators. *arXiv preprint arXiv:2212.11237*, 2022.
- Jianhao Yuan, Jie Zhang, Shuyang Sun, Philip Torr, and Bo Zhao. Real-fake: Effective training data synthesis through distribution matching. *arXiv preprint arXiv:2310.10402*, 2023.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Baoquan Zhang, Xutao Li, Yunming Ye, and Shanshan Feng. Prototype completion for few-shot learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023a.
- Ji Zhang, Lianli Gao, Xu Luo, Hengtao Shen, and Jingkuan Song. Deta: Denoised task adaptation for few-shot learning. In *ICCV*, 2023b.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3836–3847, 2023c.

- Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Yu Qiao, Peng Gao, and Hongsheng Li. Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15211–15222, 2023d.
- Yuxin Zhang, Nisha Huang, Fan Tang, Haibin Huang, Chongyang Ma, Weiming Dong, and Changsheng Xu. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10146–10156, 2023e.
- Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10145–10155, 2021.
- Chenyu Zheng, Guoqiang Wu, and Chongxuan Li. Toward understanding generative data augmentation. *Advances in neural information processing systems*, 36:54046–54060, 2023.
- Fei Zhou, Peng Wang, Lei Zhang, Wei Wei, and Yanning Zhang. Revisiting prototypical network for cross domain few-shot learning. In *CVPR*, 2023a.
- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16816–16825, 2022a.
- Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022b.
- Yongchao Zhou, Hshmat Sahak, and Jimmy Ba. Training on thin air: Improve image classification with generated data. *arXiv preprint arXiv:2305.15316*, 2023b.

A More Details of Proposed DeCap Method

A.1 Discussion About the Optimization Process

We first explain why we don’t use gradient-based methods like Finn et al. (2017); Shu et al. (2023b); Franceschi et al. (2018) for optimization, and then describe our implementation details, which enable fast training.

One may think that in the outer loop (Eq.2), it is better to use gradient-based methods to learn continuous soft prompts like Jia et al. (2022); Zhou et al. (2022b); Sohn et al. (2023). However, the primary challenge lies in the computational cost. The parameters we want to update must be propagated through a large generative model within the computation graph during backpropagation. This results in an overwhelming computational burden even with 80GB of GPU memory, it is difficult to handle the computation for just a few images. More critically, we need to perform gradient optimization over the entire synthetic dataset, which requires backpropagation through thousands of images.

On the contrary, GA algorithm demonstrates clear advantages in our problem setting. On one hand, it discovers prompts with better interpretability as shown in Sec.4.3; on the other hand, it avoids the costly computation of meta-gradients and supports various non-differentiable operations.

In terms of implementation, we can leverage the strengths of the GA algorithm to store the images corresponding to its prompt, thereby avoiding repeated image generation when a prompt is selected multiple times. We also adopt the feature caching strategy proposed by He et al. (2022) to further accelerate training. Additionally, since GA algorithm is suitable for parallel operation, the speed can be further accelerated by parallel computation. These implementation techniques significantly reduce the computational cost of DeCap algorithm: our code can run on a single GPU with 8GB of memory, and the optimization takes, on average, only 0.1 GPU hours per class. For example, on the STL10 dataset (10 classes), we consume about 1 GPU hour, and for ImageNet100 (100 classes), roughly 10 GPU hours are needed.

A.2 Examples of Prompt Pool Construction

In this section, we give a simple example about what our prompt pool looks like.

Let us consider “cat v.s. dog” classification task. Assuming that our hand-crafted prompts are [“a photo of {}”, “a sketch of {}”, “a {} image”] and model-generated prompts are {cat:[“a cat on the grass”, “a cute cat”], dog:[“a barking dog”, “a dog in the room”]}. Then, our prompt pool will be:

```
{cat:[“a photo of {cat}”, “a sketch of {cat}”, “a {cat} image”, “a cat on the grass”, “a cute cat”],
dog:[“a photo of {dog}”, “a sketch of {dog}”, “a {dog} image”, “a barking dog”, “a dog in the room”]}
```

If we randomly select 2 prompts for each class, for example, the 0th and 3th prompts for cat, and 1th and 2th prompts for dog, which represents $pop = [0, 3, 1, 2]$, the selected prompts for generating dataset will be {cat:[“a photo of {cat}”, “a cat on the grass”]; dog:[“a sketch of {dog}”, “a {dog} image”]}.

A.3 GA Algorithm Details

Genetic Algorithm (GA) is an optimization technique inspired by natural selection and genetic processes, widely used for complex combinatorial optimization problem-solving. Its key steps can be summarized as follows:

- Initialization of Population: Randomly generate a set number of individuals (solutions) to form the initial population, with each individual represented by a gene encoding (typically a binary string or real numbers).
- Fitness Evaluation: Assess the fitness of each individual using a fitness function that quantifies their performance based on the problems objectives.

- **Selection:** Select individuals for the next generation based on their fitness values. Common selection methods include roulette wheel selection, tournament selection, and rank selection, where fitter individuals have a higher chance of being chosen.
- **Crossover:** Combine parts of two parent individuals’ genes to produce new offspring. Crossover enhances genetic diversity, with methods like single-point, multi-point, and uniform crossover.
- **Mutation:** Introduce random changes to a portion of an individuals genes with a certain probability, increasing genetic variation and helping to avoid local optima. Mutation can involve flipping gene bits or assigning random values.
- **Population Update:** Merge the offspring with the current population and select suitable individuals based on fitness, often using elitism to retain the best solutions.
- **Termination Condition:** Determine if termination criteria are met, such as reaching a maximum number of iterations, achieving a predefined fitness goal, or when improvements in fitness become negligible.
- **Output Results:** Present the final optimal solution or any satisfactory solutions, along with relevant analysis and validation.

Actually, in Algorithm 1, the *GA.update()* operation means the steps from “selection” to “population update” operation. Our code are based on the scikit-opt library, and we use their default operators. What’s more, unlike traditional meta learning methods(Finn et al., 2017; Shu et al., 2018; Franceschi et al., 2018; Jang et al., 2019; Sun et al., 2020) relying on computing meta gradient to optimize outer-level meta loss, our outer-level optimization does not involve any meta gradient calculation (i.e., derivative-free optimization), and we only execute gradient descent algorithm at the inner-level optimization.

A.4 “Get_Prompt” Method in Algorithm 1

Algorithm 2 shows the “get_prompt” method in Algorithm 1. We denote the number of classes as N , the name of these classes as “class_names”, prompt numbers per class as M .

Algorithm 2 Get_prompt Algorithm

Input: indexes pop , prompt pool $pool$;

Output: prompt set: prompts, labels: Y^{syn} .

```

1: # pop is the index of  $\theta = [\theta_1, \theta_2, \dots, \theta_N]^\top, \theta_i \in \mathbb{R}^M$  in prompt pool
2:  $pop.reshape[N, M]$ 
3: #  $Y^{syn}$  contains every synthetic sample’s label
4:  $prompts = [], Y^{syn} = []$ 
5: for  $i = 1, 2, \dots, N$  do
6:    $class = class\_names[i]$ 
7:    $prompts.append(pool[class][pop[i]])$ 
8:    $Y^{syn}.append(i.repeat[M])$ 
9: end for
10: return prompts,  $Y^{syn}$ 

```

B Implementation Details

B.1 Datasets Details

In this section, we give a brief introduction about datasets we used in Section 4.

CIFAR10(Krizhevsky, 2009): The CIFAR10 dataset contains 10 common classes: airplane, car, bird, cat, dog, deer, frog, horse, ship, truck. Each class contains 6000 color images with 32×32 size. CIFAR10 is widely used in image classification.

STL-10(Coates et al., 2011): The STL-10 dataset contains 10 common classes in real life: airplane, bird, car, cat, deer, dog, horse, monkey, ship, and truck. Although these photos comes from ImageNet, their annotations may be quite different, for example, "dog" class contains various dog breeds.

Imagenette(Howard & Gugger, 2020): Imagenette is a subset of the larger ImageNet dataset, containing 10 easily distinguished classes: tench, English springer, cassette player, chain saw, church, French horn, garbage truck, gas pump, golf ball, parachute. It was created to provide a smaller, more manageable subset for training and testing image classification models.

Pets(Parkhi et al., 2012): The Pets dataset consists of images of 12 different cats breeds and 25 different dogs breeds. It is commonly used for fine-grained classification tasks, where the goal is to classify images into specific subcategories within a broader class.

ImageNet100(Tian et al., 2020): ImageNet100 is a subset of the original ImageNet dataset, containing 100 classes. It serves as a smaller alternative to the full ImageNet dataset for training and evaluating deep learning models for image classification tasks.

Caltech-101(Fei-Fei et al., 2006): The Caltech-101 dataset is a widely used benchmark dataset for object recognition. It contains images of objects belonging to 101 distinct categories, including animals, vehicles, and household items.

EuroSAT(Helber et al., 2019): EuroSAT is a dataset of Sentinel-2 satellite images for land cover classification. It contains 27,000 RGB images across 10 classes, such as agriculture, forest, and water bodies, with a resolution of 64×64 pixels. It is widely used in remote sensing and environmental monitoring tasks.

Aircraft(Maji et al., 2013): The FGVC Aircraft dataset is designed for fine-grained visual classification of aircraft. It includes 10,000 images of 102 different aircraft models, focusing on distinguishing subtle differences between similar models. It is commonly used in fine-grained recognition research.

Country211(Radford et al., 2021): Country211 is a dataset released by OpenAI, designed to assess the geolocation capability of visual representations. It filters the YFCC100m dataset to find 211 countries that have at least 300 photos with GPS coordinates. OpenAI built a balanced dataset with 211 categories, by sampling 200 photos for training and 100 photos for testing, for each country.

B.2 Experimental Settings in Section 4.1

B.2.1 Model Selection

For the pre-trained generative model, we choose the Stable Diffusion XL-Turbo (SDXL-Turbo) model² for its fast generation speed and high quality image generation. This model takes text prompts as input and outputs images at a resolution of 512×512 . During our experiments, we use ResNet-50 as the CLIP image encoder backbone. For classifier tuning (He et al., 2022), different text prompt initializations may cause slight differences in accuracy, but since our method focuses on the dataset quality, we simply use the vanilla template "a photo of {class}" for all the datasets.

B.2.2 Training Setting

Since Stable Diffusion XL-Turbo doesn't use classifier-free guidance, we simply set the guidance scale to 0 and we set inference steps to 2. For inner-level training of classification model, we generated 80 images for each class and trained for 20 epochs using the Adam optimizer with a learning rate from $2e - 3$ to $2e - 5$, equipped with the cosine learning rate schedule. For outer-level training, we set the hyper-parameters of the GA algorithm as follows: popsize of 80, maxiter of 80. Regarding the selection of few-shot datasets,

²<https://huggingface.co/stabilityai/sdxl-turbo>

we randomly selected 10 images per class to construct the few-shot datasets. We do training on 8 NVIDIA A800 GPUs, with pytorch 1.12.1 and Ubuntu 20.04.

B.2.3 Evaluation Settings

The settings of stable diffusion model are the same in Appendix B.2.2. We generated 800 images for each class and fine-tune CLIP for 30 epochs. We use the AdamW optimizer equipped with the cosine learning schedule. After training, we use the fine-tuned CLIP model to do evaluation on real test datasets. All the results are the average over 5 times run, with random seed in 7, 21, 42, 84, 105.

B.3 Experimental Settings in Section 4.2

FakeIt(Sariyildiz et al., 2023): FakeIt use Stable Diffusion V1-4 model and different classifier-free guidance scale, but our generative model are not fit for using classifier-free guidance, so we re-implemented their generation approach under our generative model. Other training settings are the same with original paper, including classification model architecture, training learning rate, data augment strategy and so on.

SuS-X(Udandara et al., 2023): The generative model of SuS-X is Stable Diffusion V1-4. For a better performance comparion, we reimplement SuS-X method with SDXL-Turbo model for higher quality image generation. The prompt strategy and other experimental settings keep the setting in the original paper.

CaFo(Zhang et al., 2023d): Since CaFo utilizes the OpenAI model to generate description for CLIP text initialization, and the original model has been deprecated, we employed the simple template “a photo of {class}” for text initialization across all datasets to ensure fairness. All other experimental settings remain consistent with the original paper. We have to point that CaFo is a few-shot learning method, and we only report the 16-shot result in Table 2 due to space limitation. Other shot results are given in Section C.6.

C More Experimental Analysis of Proposed DeCap Method

C.1 Why Is Adaptive Prompt Learning Necessary?

To validate the necessity of adaptive prompt selection, we implement three prompt selection baseline strategies: (1) randomly selecting the same number of prompts from the prompt pool, denoted by **Random**. (2) using all prompts of the prompt pool, denoted by **All**. (3) randomly substitute half of DeCap prompts to randomly selecting prompts from th prompt pool, denoted by **Half Random**. Table 5 shows the performance comparison on different datasets. All the experimental settings are the same as Appendix B.2.3. We can see that the adaptive prompts selected by DeCap method could significantly improve classification model performance compared to random selection strategy. Besides, although using all prompts in the prompt pool offers more sufficient diversity than subset selection, it suffers from various issues mentioned in Section 4.3.1, which may deteriorates the performance of classification models. This explains that the performance of all prompts is only better than the random selection strategy but not as good as DeCap method. To further illustrate whether or not the most contributive prompts are consistently selected, we randomly substitute half of the prompts learned by DeCap with randomly selecting prompts from prompt pool. It could be seen that when learned prompts are substituted, the performance would be impaired. That is to say, the adaptive prompt learning process is necessary, and the mined prompts indeed contribute much classification performance improvements on the concerned few-shot tasks. These results further substantiate the capability of adaptively mining appropriate prompts for generating high-quality images for downstream few-shot learning tasks.

C.2 Compared with Data Augmentation Strategies Using Diffusion Model

Image augmentation methods based on diffusion models usually edit the given real images to generate novel synthetic images for training the classifier. In other words, these methods attempt to generate novel images based on visual similarity. Generally, prompt-based data generation methods mainly focus on how to generate images when data are scarce or unavailable for the concerned task, while data augmentation

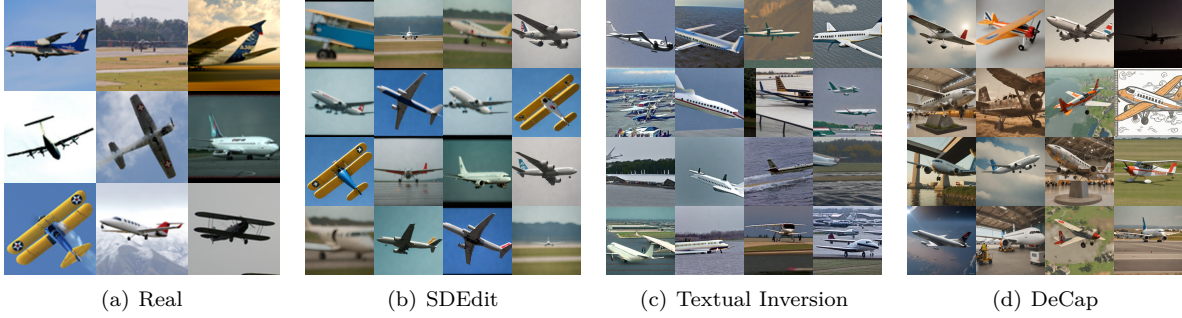


Figure 5: Examples of synthetic images by SDEdit, Textual Inversion, and DeCap methods. Compared with the other two methods, DeCap shows greater diversity and visual differences from real images, while DeCap method outperforms the other two methods in classification performance.

Table 4: Comparison of augmentation-based data generation methods and proposed DeCap method.

Methods	SDEdit	Textual inversion	DeCap
STL10	94.74	94.86	95.90
Caltech-101	85.82	84.76	88.67

methods focus on how to further boost the performance when data are relatively enough. When the real images are sparse, the synthetic images of data augmentation methods may lack diversity, which are prone to overfitting. In Table 4, we compared our method with two commonly used image augmentation methods on STL10 and Caltech-101 datasets: 1) SDEdit (Meng et al., 2021), using real image as the initial point of diffusion inverse process. 2) Textual inversion (Gal et al., 2022), learning the most similar prompts related to the given images. The results show that augmentation-based data generation methods achieve sub-optimal performance, while proposed DeCap method achieve significant performance improvements. This implies that compared with visual similarity guidance, our classification-aware meta-objective is promising to set appropriate text prompts to boost the performance of concerned few-shot tasks. Figure 5 shows some typical synthetic images generated by SDEdit, Textual Inversion, and DeCap methods. As it shown, proposed DeCap method could generate images with diversity and visual differences from real images, even some unseen “new” images. Please also see Appendix C.7, our DeCap method could generate images that are visually distinct from real images while still maintain similar classification performance.

C.3 Ablation Study

We first conducted ablation experiments on two important parameters of our method: the number of prompts selected per class and the iteration count of the GA algorithm. The experimental results of the number of prompts selected per class are reported in Table 6. We find that fewer prompts may lead to low dataset diversity, possibly hindering model performance, while more prompts may increase optimization difficulty, making it hard to find the optimal solution. We suggest to set the number of prompts selected per class as 20. Table 7 shows the performance of different iterations numbers of GA algorithm. We observed that performance of classification model converges around 80 generations. In our all experiments, we suggest to set the iterations number of GA algorithm as 80. **Moreover, we also report results of replacing the generative model from SDXL-Turbo to Playground-v2 in Tab.9. We can see that our method still work when generative models are changed.**

We further explore the result on extremely low-shot cases (e.g., 1-shot). We provide the results for the 1-shot scenario on the STL10 dataset in Table 8. In extreme cases where real samples are severely lacking, our method may underperform compared to existing approaches due to the inability to fully represent the task. Under such conditions, we recommend leveraging the inherent knowledge of the generative model. We believe this result aligns with common sense: when there is a large amount of data, we should make full use of real data. However, when data is extremely scarce, relying on the knowledge inherent in large models

Table 5: Comparion of different random selection strategies with adaptive prompt selection strategy by DeCap method.

	STL10	CIFAR10	Imagenette	Pets	Caltech-101	Imagenet100	EuroSAT	Aircraft	Country211
Random	94.74	71.94	97.25	82.26	85.76	70.48	36.76	17.88	14.56
All	95.19	74.41	97.53	82.45	84.45	70.66	37.26	18.57	14.52
Half Random	95.15	75.58	97.86	84.93	86.87	70.66	38.44	18.27	14.92
DeCap (Ours)	95.91	76.98	97.95	85.36	88.67	71.08	41.94	19.74	15.44

Table 6: Ablation study on the number of selecting prompts per class.

5	10	20	40
95.73	95.82	95.93	95.81

Table 7: Ablation study on the iteration number of GA algorithm.

20it	40it	60it	80it
95.73	95.87	95.93	95.93

Table 8: Ablation study on extremely low shot real images case (1-shot).

1shot	10shot	vanilla	CiP
95.15	95.93	95.33	94.92

may be a better option because the given data cannot fully represent the task. And we will explore to better address the issue in the future work.

Table 9: The result of replacing the generative model from SDXL-Turbo to Playground-v2.

	vanilla	multi	LE	CiP	DeCap
DTD	50.53	52.29	54.09	63.13	63.78
UCF101	60.89	60.91	61.35	64.26	64.84

C.4 Experimental Results on Different Evaluation Metrics

In this section, we give results of other evaluation metrics including precision (Table 10), recall (Table 11) and F1-score (Table 12), which are commonly used in few-shot learning, to further explore the robustness and generalization ability of DeCap method. Some brief introduction about these metrics are given as follows:

- **Precision:** Precision measures the accuracy of positive predictions. It is defined as:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

Precision answers the question: *Of all the instances predicted as positive, how many are actually positive?*

- **Recall:** Recall, also known as sensitivity or true positive rate, measures the ability of the model to correctly identify positive instances. It is defined as:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$$

Recall answers the question: *Of all the actual positive instances, how many were correctly predicted?*

- **F1 Score:** The F1-score is the harmonic mean of precision and recall, providing a single metric that balances both. It is defined as:

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F1-score is particularly useful when the class distribution is uneven or when precision and recall are equally important.

The results demonstrate that our method performs well on these metrics, indicating that it not only achieves high accuracy but also excels in identifying positive samples and is more cautious when dealing with them. It more comprehensively illustrates the robustness and generalization of our method.

Table 10: **Precision** results of different methods among all datasets.

	vanilla	multi	LE	CiP	DeCap
STL10	95.35	94.90	94.17	95.16	95.63
CIFAR10	76.61	76.55	77.32	77.23	77.40
Im-10	97.27	97.30	97.27	97.30	97.34
Pets	82.58	84.76	84.17	84.52	85.70
Caltech-101	84.42	84.57	85.27	85.39	85.43
Im-100	69.82	71.27	69.28	73.03	71.90
EuroSAT	43.01	38.45	50.50	47.32	49.36
Aircraft	18.38	18.85	20.85	18.67	20.85
Country211	17.20	17.26	18.10	17.12	17.77

Table 11: **Recall** results of different methods among all datasets.

	vanilla	multi	LE	CiP	DeCap
STL10	95.31	94.78	94.04	94.72	95.57
CIFAR10	72.39	69.63	68.24	68.63	76.99
Imagenette	97.25	97.28	97.25	97.28	97.33
Pets	81.69	82.05	82.13	83.12	84.52
Caltech-101	85.21	86.02	85.37	86.07	86.54
Imagenet100	68.32	69.98	67.00	69.56	70.94
EuroSAT	31.07	29.62	28.31	40.62	42.22
Aircraft	17.03	17.84	17.72	17.98	19.71
Country211	13.72	13.90	14.14	14.98	15.44

C.5 Adversarial Robustness

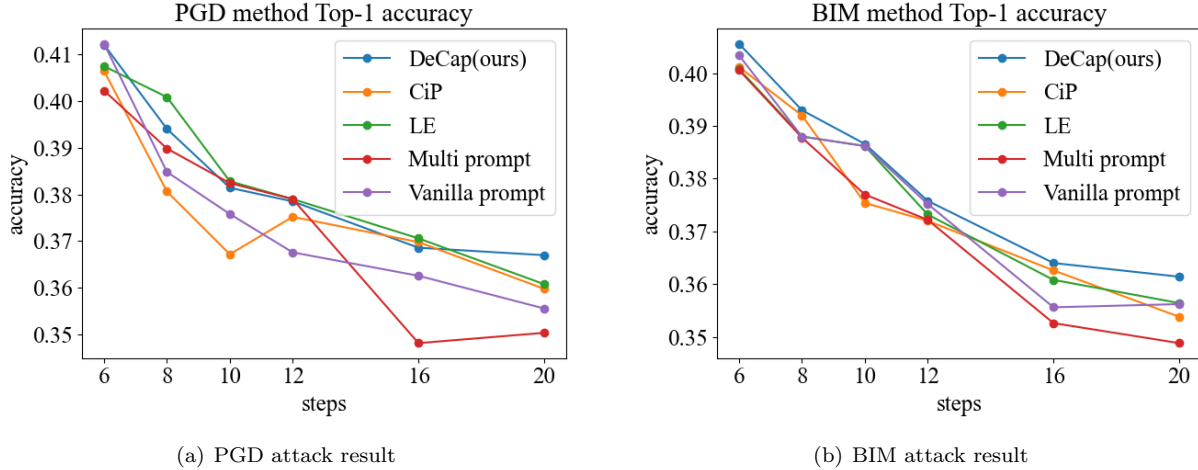


Figure 6: **Adversarial robustness of classification models trained with synthetic images using different prompts designing methods.** We report results on ImageNet100 validation set under two adversarial attack methods. The horizontal axis represents the number of steps taken in the attack, and the vertical axis represents the accuracy of the trained classification models on the validation set after the attack.

Adversarial learning aims to evaluate model robustness by adding small perturbations to the input data, causing the model to make false predictions but making little difference to human observers. We use two common attack methods: BIM (Basic Iterative Method) attack (Kurakin et al., 2018) and PGD (Projected

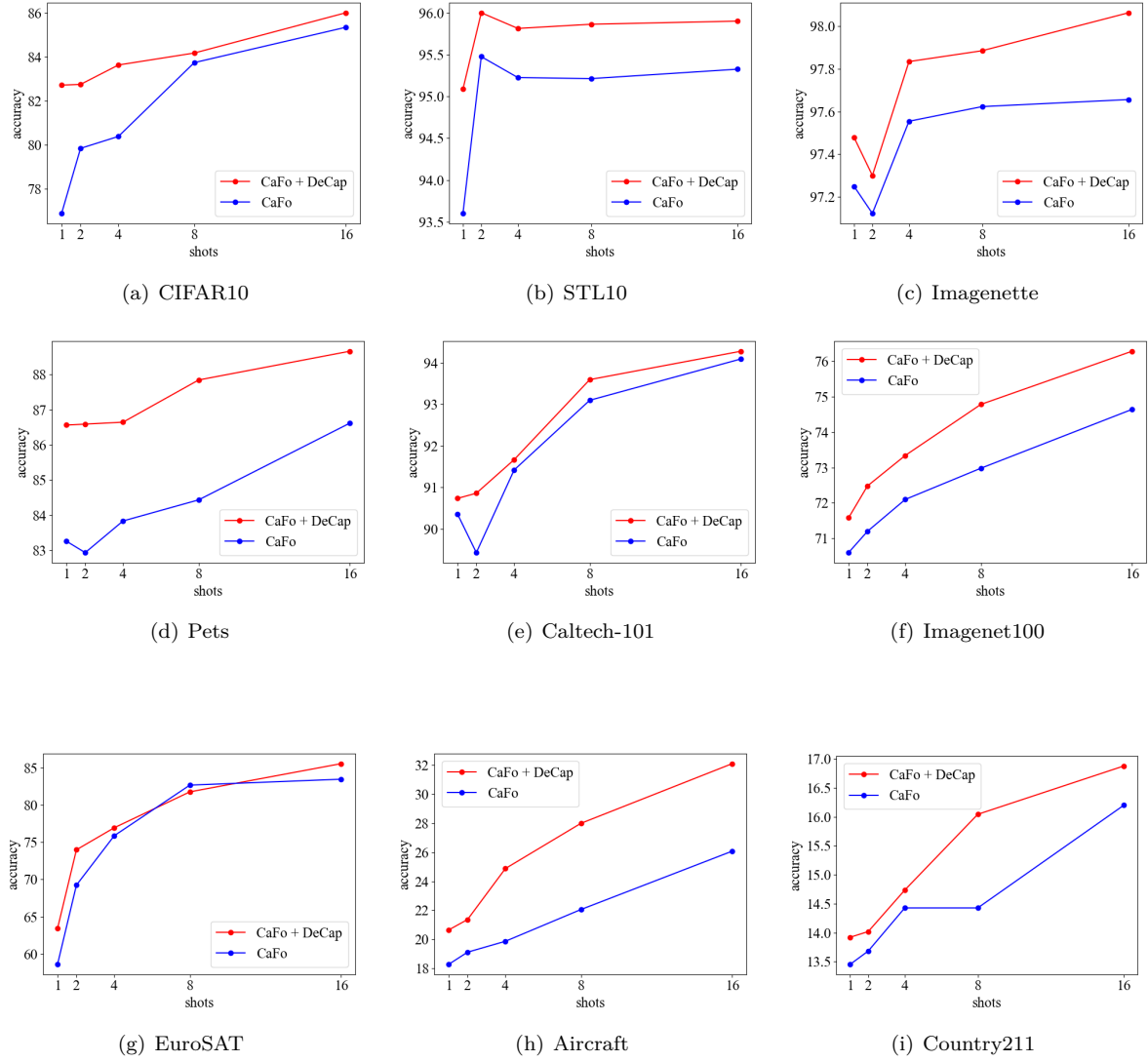


Figure 7: Classification accuracies on different shots of each class for CaFo and CaFo + DeCap Methods

Table 12: **F1-score** results of different methods among all datasets.

	vanilla	multi	LE	CiP	DeCap
STL10	95.29	94.71	93.99	94.74	95.58
CIFAR10	71.89	69.04	67.30	69.19	76.89
Imagenette	97.23	97.26	97.23	97.26	97.31
Pets	81.28	81.96	82.09	83.23	84.67
Caltech-101	82.04	82.87	82.12	83.45	83.73
Imagenet100	67.06	68.99	65.52	69.42	70.25
EuroSAT	26.36	24.43	24.38	36.36	39.34
Aircraft	15.10	15.98	16.00	15.95	17.81
Country211	13.18	13.35	13.62	14.36	14.93

Table 13: More experimental results of SuS-X and CaFo methods. This table presents the classification accuracies of other data generation algorithms on these two zero/few-shot methods. DeCap still achieves leading performance in the vast majority of datasets.

	STL10	CIFAR10	Imagenette	Pets	Caltech-101	Imagenet100	EuroSAT	Aircraft	Country211
SuS-X									
Multi	95.29	72.28	98.24	78.99	84.89	70.08	34.43	18.45	13.09
LE	95.39	73.62	98.37	79.20	84.51	69.98	27.41	18.33	13.01
CiP	95.21	73.19	98.34	79.77	84.76	70.14	33.02	19.47	13.63
DeCap	95.43	75.89	98.39	80.40	84.89	70.30	37.37	19.83	13.63
CaFo									
Multi	95.48	85.76	97.88	86.92	94.27	75.58	85.19	31.71	16.65
LE	95.51	86.00	97.78	88.06	94.03	75.00	85.63	26.37	15.27
CiP	95.33	84.05	97.88	87.38	94.27	75.92	85.64	27.18	14.83
DeCap	95.90	86.00	98.06	88.66	94.28	76.28	85.59	32.10	16.88

Gradient Descent) attack (Madry et al., 2018). The BIM employs an iterative gradient ascent approach, where at each step, BIM perturbs the image along the gradient direction predicted by the model. It can be written as $x_{i+1} = x_i + \epsilon \nabla_{x_i} J_\theta(x_i, y)$, where x_0 denotes the original image, y denotes its label, and $\nabla_{x_i} J$ means the gradient of loss function w.r.t. x_i . PGD further projects the adversarial examples into an ϵ -ball around the original image.

We use classification model weights obtained from Section 4.1 and implement adversarial attack on ImageNet100 validation dataset. We use torchattacks Kim (2020) library to conduct this experiment. We select attack step size ϵ as $1/255$ for these two methods, and Fig.6 reports the attack result on different attack steps. We found that model-generated prompts, due to their rich content details, have a slight advantage in adversarial robustness compared to hand-crafted prompts. Moreover, since DeCap integrates the strengths of hand-crafted and model-generated prompts methods, it consistently performs well in terms of resilience against adversarial attacks.

C.6 More Experimental Results of Section 4.2

In Section 4.2, due to space limitations, we only present the compared results of the original algorithms and DeCap method. In Table 13, we additionally provide the compared results of other prompt design methods from Table 1 on the SuS-X and CaFo algorithms. The experimental results are aligned with conclusions in Section 4.2.

Fig.7 shows more experimental results on different shots of each class for CaFo and CaFo + DeCap Methods. The experimental results are also aligned with conclusions in Section 4.2.



Figure 8: (Left) Examples of real images, synthetic images by hand-crafted prompts and DeCap methods on PACS Sketch dataset. (Right) Performance comparison between prior hand-crafted prompts and DeCap methods. Although DeCap appears visually distinct from real data, its classification performance can approach that of training on the full dataset and significantly outperforms the prior prompts “a sketch figure of {}”.

C.7 Are prompts with domain information enough for classification?

To illustrate this point, we conducted experiments on the Sketch subclass of the PACS dataset Li et al. (2017), and the experiment settings are the same in Section B.2. For this dataset, all images follow the same style. As shown in Fig.8, although using the prior information “a sketch figure of ” yields better results than random domain information, while synthetic images by DeCap method displays different visual effects. Specifically, the prompts DeCap method discovered not only include the “sketch” category but also encompassed a wide range of domain information. Despite these images being visually distinct from real images, their classification performance significantly surpassed that of using only the “sketch” prompts, even approaching the performance achieved by training on the entire real dataset. This result demonstrates that diverse prompts help boost classification performance, and DeCap method genuinely focuses on mining proper prompts for improving classification performance rather than visual effects. Even so, these selected prompts are hard to be interpreted by human vision. We hope that a rational theoretical insight could characterize such phenomenon in the future study.

D Visualization of Synthetic Images and Learned Prompts

D.1 Visualization of Synthetic Images

Fig.9 and Fig.10 show more visualized examples of synthetic images on Pets and Imagenet100 datasets generated by DeCap method. As it shown, proposed DeCap method could generate diverse and multi-domain images, and thus obtain better classification performance.

D.2 More Visualization of Learned Prompts

We will demonstrate that DeCap method can adaptively learn proper and dataset-specific prompts that are suitable for concerned tasks from the following three aspects.

Firstly, the ratios of the number of model-generated and hand-crafted prompts for each class are varying, as shown in Fig. 11 and 12. This reflects that our method could adaptively adjust the proportions that reconcile class/domain information and rich content information for different classes.

Secondly, though the hand-crafted prompts follow the same templates, we can see that different classes may learn relatively different prompts in Fig.4. This further reveals our method could adaptively learn classification-aware prompts for each class, so as to achieve better performance on downstream tasks. Moreover, we additionally give some examples about the consistently selected model-generated prompts during the optimization process to further highlight the significance of integrating fine-grained prompt descriptions. As we can see in Table 15, the consistently selected prompts show high diversity and fine-grained informa-

tion, including: movement, posture, background, color, quantity, other objects, and so on. This pictures significantly help to provide classification-benefit features.

Lastly, Table 14 shows that though STL-10 and CIFAR-10 datasets have some same categories, the learned prompts by our method could be almostly different. This demonstrated that our DeCap method could learning proper prompts suitable to concerned few-shot datasets. For instance, we can see that learned prompts for the STL-10 dataset are realistic, while learned prompts for CIFAR-10 dataset are of low-resolution imagery. Notice that these prompts are well aligned with prior knowledge of these datasets.

Moreover, we additionally display the complete set of prompt pool of the “airplane” class in STL-10 dataset in Table 16, to offer a more intuitive understanding for the characteristic of our method stated above. And we further give visualizations that demonstrate the prompt selection process over the course of optimization, including image examples and the evolution of prompts, please see Fig.13.

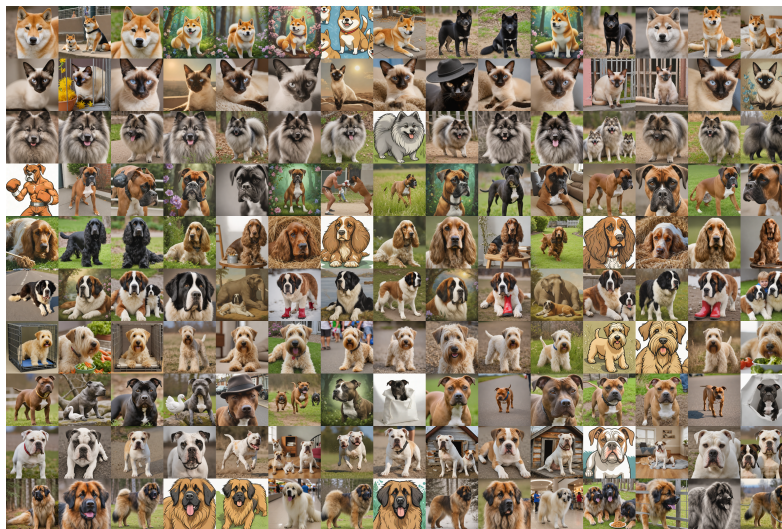


Figure 9: Examples of generated images on Pets dataset by DeCap method.

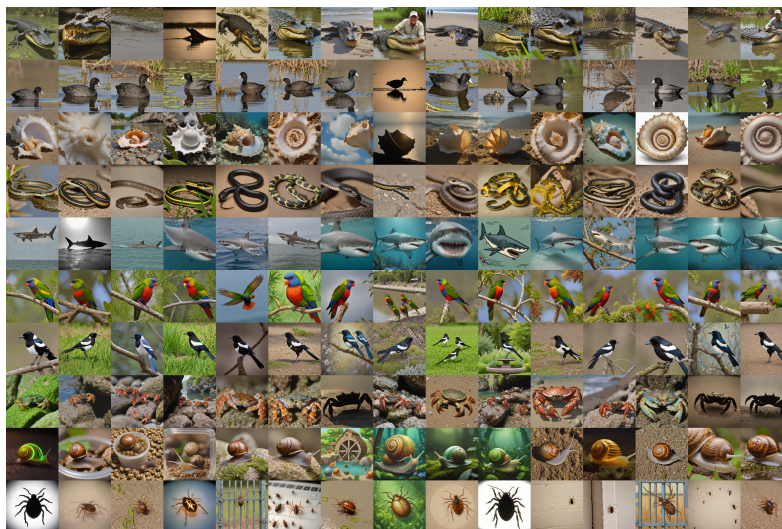


Figure 10: Examples of generated images on Imagenet100 dataset by DeCap method.

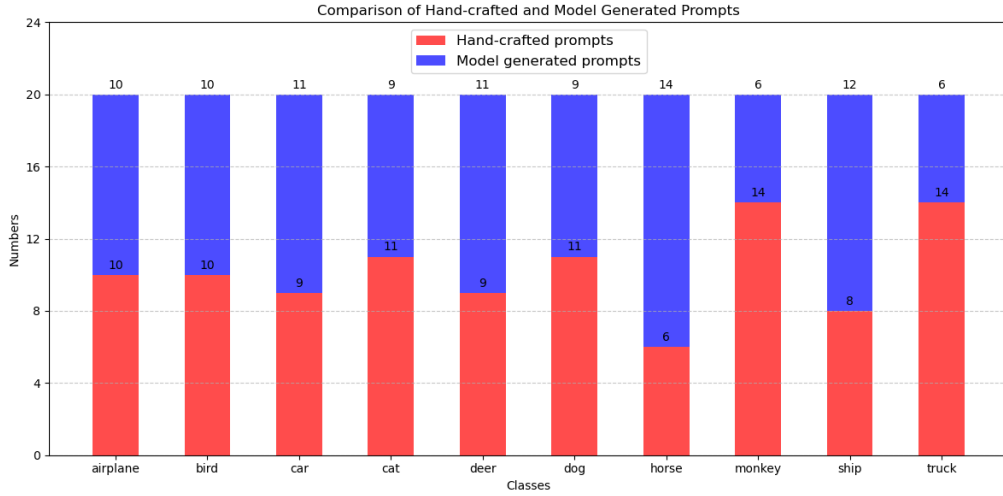


Figure 11: Illustrations of the number of hand-crafted prompts vs the number of model-generated prompts mined by DeCap method on STL-10 dataset.

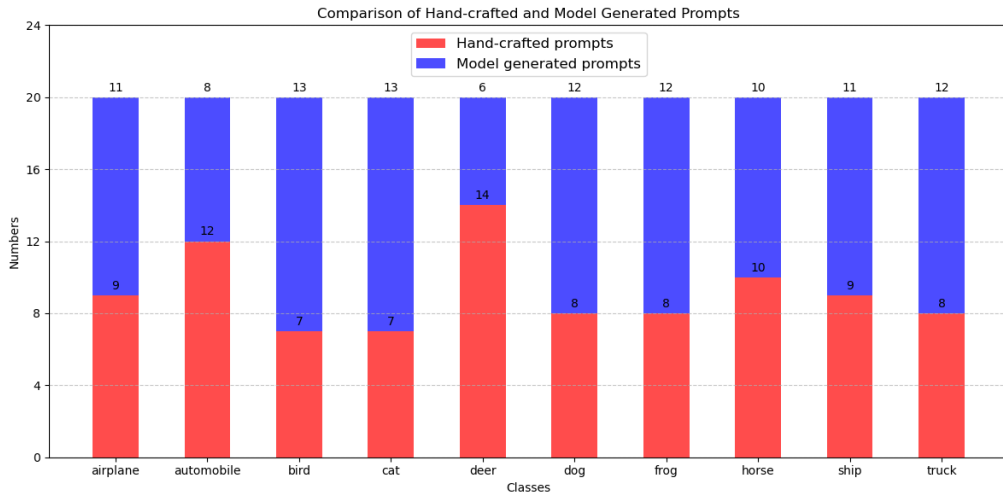
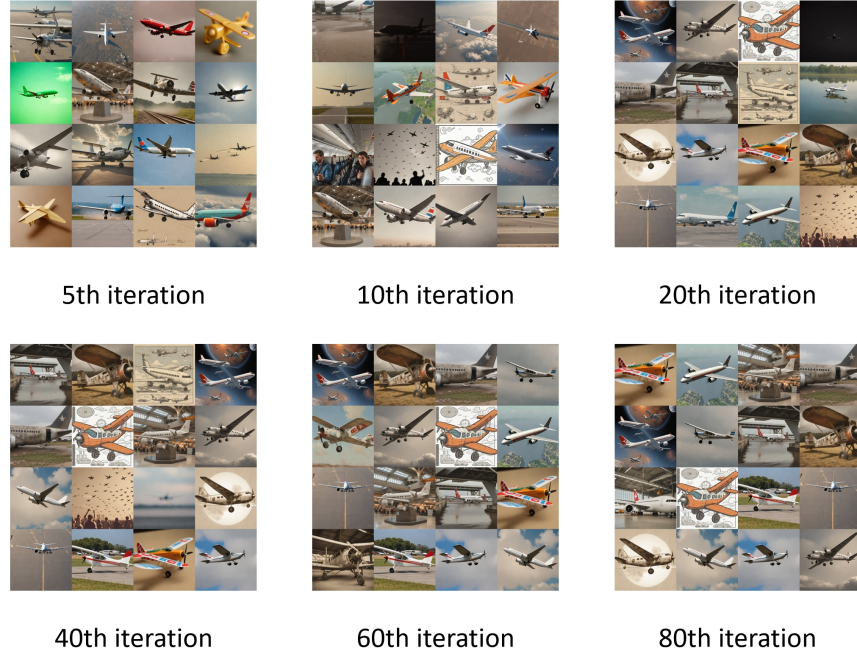


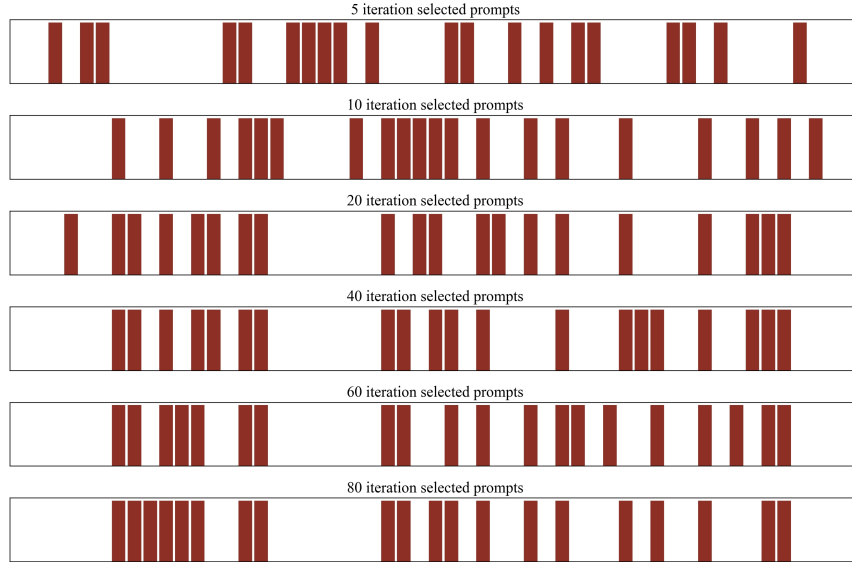
Figure 12: Illustrations of the number of hand-crafted prompts vs the number of model-generated prompts mined by DeCap method on CIFAR10 dataset.

Table 14: Illustration of mined prompts for “deer” class on different datasets. DeCap method selects completely different prompts for the same class across different datasets, demonstrating its ability to adaptively learn the prompts suited to each specific dataset.

STL-10	A deer is grazing the woods. deer are grazing under a tree a photo of the clean deer. A silhouette of deer. a deer and young man roam around during a december game a photo of a deer. a deer in a video game. a toy deer a deer on a pond the cartoon deer. the hornets and deer are on a ridge a deer. deer resting with the grazing padou atop old farmhouse An ink painting of a deer deer on a green pond. the toy deer. a brown bear eats the deer A glossy deer. a photo of a large deer. a group of deer on prairie are seen grazing in their natural habitat
CIFAR10	a photo of deer, a wild deer in the wild a deer on a farm A soft-focus deer. art of a deer. deer and their prey on the northern slopes a photo of deer, a deer standing in the snow with a sky background a pixelated photo of the deer. several deer grazing in the desert fox and a deer on the grounds of a city a rendering of a deer. a photo of deer, a group of deers standing in a field A silhouette of deer. a photo of deer, a deer is standing in the grass a deer is grazing an ancient inscription. a photo of deer, a herd of deer in the desert deer and the munro. A pair of deer on a trail. a hunt deer on a desert land deer and the munro. a pixelated photo of the deer.



(a) Image examples during different optimization iterations.



(b) The evolution of prompts during different optimization iterations.

Figure 13: Image examples and the evolution of prompts during different optimization iterations. For clarify, Fig.(b) shows only the selected prompts and omits the rest.

Table 15: Examples of model-generated prompts which are consistently selected during optimization process.






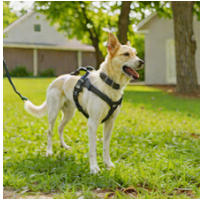



<p>a bird sitting on a branch.</p> 	<p>cars that have to make an effort to turn off.</p> 	<p>A black cat is in a room where the window is down.</p> 
<p>A truck with lots of people on it.</p> 	<p>A deer is grazing the woods.</p> 	<p>A dog is standing in its yard with a harness on it.</p> 
<p>A white horse in the open barn.</p> 	<p>dogs inside a home on a summer.</p> 	<p>a semi truck driving down a rural road.</p> 

Table 16: The prompt pool of “airplane” class in STL-10 dataset.
Mined prompts by DeCap method are highlighted in **bold**.

<i>Model-generated prompts</i>			
a photo of airplane , a small plane is parked on the runway	a photo of airplane , a large passenger jet flying through a blue sky	a photo of airplane , a small plane is on the runway	a photo of airplane , a yellow airplane flying through a blue sky
a photo of airplane , a large passenger jet flying through the sky	a photo of airplane , a small plane flying in the sky	a photo of airplane , a small plane is floating in the water	a photo of airplane , a large passenger jet sitting on a runway
a photo of airplane , a small plane flying in the sky	a photo of airplane , a small blue airplane is taking off from the runway	a photo of airplane , a plane is parked on the tarmac	a photo of airplane , two small planes are sitting on the water
a photo of airplane , a plane flying in the sky	a photo of airplane , a small plane flying over a mountain range	a photo of airplane , a plane is on the runway	a photo of airplane , a small plane flying through the air
a photo of airplane , a large white plane	a photo of airplane , two planes flying in the sky	a photo of airplane , a plane flying in the sky	a photo of airplane , a small plane flying over a city
a photo of airplane , a plane flying in the sky	a photo of airplane , a small plane flying through the air	a photo of airplane , a plane flying in the sky	a photo of airplane , a plane is parked on the tarmac
a photo of airplane , a plane flying in the sky	a photo of airplane , a plane flying in the sky	a photo of airplane , a small plane is parked on the water	a photo of airplane , a plane flying in the sky
a photo of airplane , a small plane sitting on a snowy field	a photo of airplane , a small plane flying in the sky	An airplane that has been seen flying over another airplane.	A plane is in a parking lot.
airplane that you bought a few years ago	A small airplane is flying over a highway at a time.	An airplane that is parked in an airport	The plane has an engine, a seat, a console, a charger, and
An aircraft is in the flight over a lake.	The airplanes are all parked inside the parking lot.	A plane is in the air.	plane of a small aircraft.
A red and white airplane with a green and green color scheme.	An airplane parked on the runway near a pier.	An airplane that has just broken ground behind it.	A plane parked next to one of the airplanes above it’s engine.
Airplanes in space that are not as big as usual.	An airplane parked along a highway.	A small airplane that’s flying at low speeds under a cloudy sky.	airplanes need people to work hard at the zoo
An airplane parked next to a bridge.	Some airplanes flying over people.	The airplane is in a green sky with blue skies.	The airplane with the lights is about to be docked.
An airplane with three engines and a propeller.	an airplane with a window	An airplane parked on top of a hill next to it	airplane on the tracks.
An airplane on an airplane track	A plane with tires on it flying away from it.	An airplane is parked on a runway at a airport.	These airplanes are in a wing.
a commercial airplane traveling in july.	airplane in flight... a photo and video	Two air planes all flying in a row.	A modern airplane is arriving in the air.

An airplane in the middle of nowhere with its doors lowered.	airplane is parked in a parking lot	An airplane that is coming in to land.	passengers in an airplane in the rain
airplanes flying at a rate of 2 to 3 mph on a sunday	An airplane on a runway next to a small green field.	an airplane on an airport runway	a classic red blue airplane is shown in the cockpit with bright colors as well.
planes in a dry pit	airplane and other objects in the air	an airplane makes an outgoing landing on the ground	An old airplane is coming down the track.
A man attempting to board a commercial airplane.	Small airplanes with wing lights attached to them.	A small airplane with the tail mounted up.	An airplane in a flight path with some passengers nearby.
airplane on an old building	An aircraft goes up through a window dripping with smoke and debris.	airplanes that have been converted to jet engines	airplanes cruising in the bay.
A boy is running with an airplane that is on the runway.	A blue airplane has its wings shut.	An airplane is about to land in a parking lot and be delivered.	two airplanes parked at the airport
A white airplane on the runway with blue ice.	jet airplane is ready for a test	An airplane is sitting on a ground with all three engines on the ground.	There's one airplane in the cockpit which is parked by another airplane.
An airplane is in the air.	A family is on a small airplane at a hotel.	aircraft carrier and an airplane together with some gulls.	airplane inside of the airplane
The airplane is looking down.	An airplane is shown flying on a runway.	small bodied airplane on a plane	An airplane parked next to fireworks on the sky.
An airplane that appears to be on the runway.	An airplane that is very close to the ground in an airport.	Various aircraft and airplanes are getting ready for flight.	an airplane is seen arriving on a runway
Two aircrafts in a white airplane at a station.	The airplane landed.	an airplane that is making a flying flight	airplane on the runway at the airport
this airplane was able to take off with just a small amount of effort to get the	An airplane that is in a flying position.	An airplane making its way between jets.	airplane sitting in air
a large old plane sits off the fuel tank	aircraft carrier and its crew arriving in an airplane	airplane on the runway	A family of airplanes are in a building.
plane flies around city	an airplane about to land in a desert	An electric airplane in the sky.	an airplane that is making it's way around the tarmac
airplanes on the runway	the crew of airplane on board	The airplane is in the air.	jet airplane wing during maintenance
A blue and white airplane with white wing panels.	A commercial airplane flying under the radar.	The airplane has been damaged by the winds.	An airplane flying near a tarmac.

Hand-crafted prompts

a good photo of the airplane.	a photo of many airplane.	a sculpture of a airplane.	a photo of the hard to see airplane.
a low resolution photo of the airplane.	a rendering of a airplane.	graffiti of a airplane.	a bad photo of the airplane.
a cropped photo of the airplane.	a tattoo of a airplane.	the embroidered airplane.	a photo of a hard to see airplane.
a bright photo of a airplane.	a photo of a clean airplane.	a photo of a dirty airplane.	a dark photo of the airplane.
a drawing of a airplane.	a photo of my airplane.	the plastic airplane.	a photo of the cool airplane.
a close-up photo of a airplane.	a black and white photo of the airplane.	a painting of the airplane.	a painting of a airplane.
a pixelated photo of the airplane.	a sculpture of the airplane.	a bright photo of the airplane.	a cropped photo of a airplane.
a plastic airplane.	a photo of the dirty airplane.	a jpeg corrupted photo of a airplane.	a blurry photo of the airplane.
a photo of the airplane.	a bad photo of a airplane.	a rendering of the airplane.	a airplane in a video game.
a photo of one airplane.	a doodle of a airplane.	a close-up photo of the airplane.	a photo of a airplane.
the origami airplane.	the airplane in a video game.	a sketch of a airplane.	a doodle of the airplane.
a airplane.	a origami airplane.	a low resolution photo of a airplane.	the toy airplane.
a rendition of the airplane.	a photo of the clean airplane.	a photo of a large airplane.	a rendition of a airplane.
a photo of a nice airplane.	a photo of a weird airplane.	a blurry photo of a airplane.	a cartoon airplane.
art of a airplane.	a sketch of the airplane.	a embroidered airplane.	a pixelated photo of a airplane.
a jpeg corrupted photo of the airplane.	a good photo of a airplane.	a photo of the nice airplane.	a photo of the small airplane.
a photo of the weird airplane.	the cartoon airplane.	art of the airplane.	a drawing of the airplane.
a photo of the large airplane.	a black and white photo of a airplane.	a dark photo of a airplane.	graffiti of the airplane.
a toy airplane.	a photo of a cool airplane.	a photo of a small airplane.	a tattoo of the airplane.
a digital style airplane	a colorful airplane	a modern style airplane	an abstract photo of airplane
a cartoon style airplane	a virtual style airplane	An ink painting of a airplane	a toy airplane
A model airplane.	a red airplane	a blue airplane	a yellow airplane
a black airplane	a white airplane	An old airplane.	A futuristic airplane.
A minimalist airplane.	A detailed illustration of airplane.	A close-up of airplane.	A shadowy figure of airplane.
A silhouette of airplane.	A bright and vibrant airplane.	An abstract concept of airplane.	A vintage style airplane.
A neon-lit airplane.	A monochrome airplane.	A watercolor painting of airplane.	A sketch of airplane.
A digital art of airplane.	A handcrafted airplane.	An aerial view of airplane.	A side profile of airplane.

A textured airplane.	A glossy airplane.	A matte airplane.	A glowing airplane.
A rustic airplane.	A weathered airplane.	A sparkling airplane.	A serene airplane.
A chaotic airplane.	A whimsical airplane.	A dynamic airplane.	A frozen moment of airplane.
A soft-focus airplane.	A high-contrast airplane.	A sepia-toned airplane.	A saturated airplane.
An isolated airplane.	A mirrored airplane.	A panoramic view of airplane.	An enchanted airplane.