# SEMANTIC-GUIDED LORA PARAMETERS GENERATION

# **Anonymous authors**

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

018

019

021

024

025

026

027 028

029

031

033

034

037

038

040 041

042

043

044

046

047

048

051

052

Paper under double-blind review

# **ABSTRACT**

Generating new Low-Rank Adaptation (LoRA) weights from pre-trained LoRAs has demonstrated strong generalization capabilities across a variety of tasks for efficiently transferring AI models, especially on resource-constrained edges. However, previous studies either merge base LoRAs via weighting coefficients or train a generative model in the closed-world assumption, limiting their efficiency and flexibility in complex edge user cases. This challenge may further increase when there are significant domain shifts between training and deployment. To this end, we propose Semantic-guided LoRA Parameter Generation (SG-LoRA), a tuning-free generative framework to efficiently produce task-specific parameters for unseen tasks in a semantic-to-LoRA pipeline. Concretely, SG-LoRA uses task descriptions as the semantic bridge, measuring their proximity to a set of known expert tasks in a shared embedding space. Based on this semantic guidance, it models the target task's LoRA parameter distribution to generate high-performing parameters for novel tasks. SG-LoRA enables the real-time construction of LoRA models aligned with individual intents by distilling knowledge from prominent LoRA experts and, meanwhile, offering a privacy-preserving solution for personalized model adaptation in a novel zero-shot open-world setting proposed in this work. Extensive experiments on multiple challenging tasks confirm the superior performance and remarkable adaptability of SG-LoRA.

## 1 Introduction

In recent years, deep learning has seen remarkable progress, largely driven by the advent of large-scale pre-trained models (LPMs) Chung et al. (2024); Li et al. (2022); Liu et al. (2023); Rombach et al. (2022); Touvron et al. (2023). Trained on massive and diverse datasets, these models demonstrate exceptional performance across a wide range of downstream tasks Shenaj et al. (2024); Alayrac et al. (2022); Touvron et al. (2023). However, as both model and data scales continue to grow, retraining the entire model becomes increasingly computationally expensive and often infeasible in practice. To mitigate this challenge, parameter-efficient fine-tuning (PEFT) methods have drawn considerable attention Zhang et al. (2023c;b); Ding et al. (2023). Among them, Low-Rank Adaptation (LoRA) Hu et al. (2022) has emerged as a prominent approach. LoRA adapts pre-trained models by introducing a small number of trainable low-rank matrices into existing layers, achieving strong task-specific performance while leaving the original model weights unchanged Sung et al. (2022).

While an increasing number of pre-trained LoRA modules are becoming publicly available, effectively leveraging them in real-world scenarios remains a significant challenge. As shown in Figure. 1(a), we propose the Zero-Shot Open-world Adaption (ZSOA) in this paper, which aims to generate LoRA weights for unseen tasks based on a set of pre-trained LoRAs. ZSOA emphasizes two key aspects: (1) No raw data is available for the unseen task, highlighting the need for rapid adaptability to evolving user intents; and (2) Open-world task coverage, defined by a broad and unconstrained task space in which the unseen tasks may not be directly related to the seen tasks. Compared to traditional LoRAs that need to be fine-tuned on downstream tasks, ZSOA shows data and computation-friendly strength, resulting in more flexibility in practice, particularly in edge environments where data privacy constraints and limited computational resources make large-scale retraining infeasible.

To enable broader applications of LoRA, prior research has explored two main directions, as illustrated in Figure. 1 (b-c), each partially addressing the challenges of ZSOA. The first line of work focuses on merging-based methods, which aim to rapidly construct task-specific models by directly fusing existing LoRA modules at hand Wortsman et al. (2022a); Yadav et al. (2023); Shenaj et al.

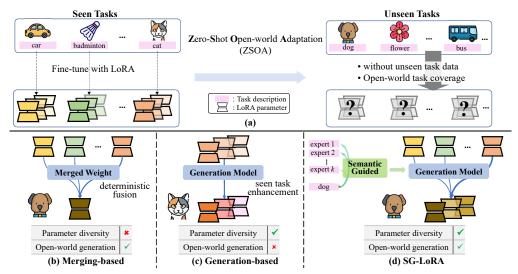


Figure 1: **Motivation of our SG-Lora.** We consider a challenging scenario termed Zero-Shot Open-World Adaptation (ZSOA), where a model is provided rich Loral resources for seen tasks but lacks access to data for unseen tasks during inference, with an unconstrained task space. Conventional Loral adaptation methods are not suitable for ZSOA: merging-based approaches struggle to explore the diversity of Loral parameters, while generation-based methods primarily focus on Loral enhancement for seen task. Our SG-Loral uses task descriptions as semantic guidance to enable conditional Loral generation for unseen tasks in a data-free and open-world manner. (Each color family represents a set of Loral parameters for the same task, for example, brown represents Loral for the 'Dog' task, and yellow represents Loral for the 'Car' task.)

(2024). Although these methods support open-world generation, the generated weights are obtained through deterministic fusion of existing LoRAs, resulting in limited diversity and constraining the model's ability to adapt to flexible or evolving requirements Shenaj et al. (2024); Wu et al. (2024). Moreover, the merging process must be carefully designed, as conflicts may arise when integrating LoRA modules trained on different tasks Zou et al. (2025); Zhao et al. (2024). In parallel, another direction explores generation-based methods, which leverage generative models, such as variational autoencoders (VAEs) Kingma & Welling (2013) or diffusion models Ho et al. (2020), to synthesize new LoRA parameters. By introducing stochasticity, these approaches enable greater diversity in parameter generation and bypass traditional fine-tuning pipelines. However, their success often relies on a closed-world assumption, where training and test tasks are drawn from similar distributions Soro et al. (2024). As a result, these methods are susceptible to task (or domain) shifts scenarios, and struggle to handle open-world tasks.

In this work, we draw inspiration from the human ability to intuitively infer semantic relationships between prior knowledge and new tasks Lake et al. (2017), enabling effective generalization in unfamiliar situations. For example, after learning to recognize cat breeds such as *Birman* and *Egyptian Mau*, a person can identify *British Shorthair* based solely on its textual description by relating it to previously acquired concepts. Motivated by this analogy-driven reasoning process, we propose a Semantic-Guided LoRA Parameter generation framework (SG-LoRA) that adapt LPMs to the ZSOA setting. Specifically, given a set of expert (or base) LoRAs whose weights are trained on seen tasks, we aim to train a semantic-to-LoRA model that takes the semantics of the unseen tasks as inputs and outputs high-performing LoRA weights directly. Notably, the task semantics serve as a bridge between the seen and unseen tasks, guiding our SG-LoRA on how to leverage expert knowledge to generate task-specific LoRA weights. Motivated by prompt engineering Zhou et al. (2022), we adopt the task description to identify the task semantics. These descriptions, typically concise yet semantically rich, are processed by a frozen CLIP text encoder to capture task-level correlations without exposing user-specific data. The task semantics are then modeled as Gaussian distributions according to the relationship between the unseen and seen tasks.

Importantly, simply scaling the number of expert LoRAs does not guarantee performance gains, as they may provide contradictory or irrelevant task knowledge. To resolve this, we design a sparse aggregator that assembles the most semantically relevant expert to integrate rational prior knowledge for a target task. During inference, the system directly generates target LoRA modules aligned with user requirements using only textual task queries. Crucially, the stochastic nature of our trained

generator converts deterministic LoRA construction into probabilistic parameter sampling, enhancing both parameter diversity and dynamic adaptability to evolving user intents. In summary, the main contributions of this work include :

- We introduce Semantic-Guided LoRA Generation, a versatile framework that harnesses semantic task relationships to enable zero-shot open-world adaptation. By conditioning on prior available task knowledge, SG-LoRA can synthesize high-performance LoRA parameters for arbitrary unseen tasks without retraining.
- By seamlessly integrating generated LoRA modules into off-the-shelf LPMs, our method enables fast personalization at inference time. It allows flexible configuration of the expert LoRA repository, supporting task-adaptive LoRAs both within and across datasets, thereby achieving scalable and adaptable model behavior.
- Comprehensive experiments on multiple image-text retrieval benchmarks demonstrate that the proposed method can rapidly generate LoRA parameters achieving performance comparable to traditional LoRA fine-tuning.

## 2 RELATED WORK

## 2.1 LOW-RANK ADAPTATION

Low-Rank Adaptation (LoRA) is a parameter-efficient fine-tuning method to adapt large models to novel tasks by approximating weight updates with low-rank matrices, drastically reducing trainable parameters while maintaining pre-trained knowledge Huang et al. (2023); Zhang et al. (2023a). Specifically, given a pre-trained weight  $W_0$  with input x and hidden state h, LoRA decomposes the weight update  $\Delta W \in \mathbb{R}^{a \times b}$  into two low-dimensional matrices  $B \in \mathbb{R}^{a \times r}$  and  $A \in \mathbb{R}^{r \times b}$ :

$$h = W_0 x + \Delta W x = W_0 x + \gamma B A x, \tag{1}$$

where the rank  $r \ll \min(a,b)$ .  $\gamma$  is a constant scaling hyperparameter that controls the contribution of the LoRA update. Recently, a growing number of high-quality, pre-trained LoRA modules have become publicly available for various architectures Civitai (2024); HuggingFace (2024), including transformers and vision–language models, offering rich resources to accelerate adaptation and deployment on specific downstream tasks.

#### 2.2 Model Merging

Model merging integrates parameter-level knowledge from multiple independently trained networks into a unified model, achieving enhanced capabilities Li et al. (2023); Ilharco et al. (2022); Jang et al. (2024); Wortsman et al. (2022b). The pioneer work Model Soups Wortsman et al. (2022a) establishes weight averaging as a foundational paradigm, showing that averaging fine-tuned models from identical pre-trained bases with varied hyperparameters consistently outperforms individual models. AdapterSoup Chronopoulou et al. (2023) generalizes the Model Soups paradigm to cross-domain adaptation by dynamically averaging domain-specific adapters at test time. This approach preserves the base model's integrity while enhancing out-of-distribution generalization through selective weightspace interpolation of relevant domain knowledge. Recent advances have extended this paradigm to LoRA-based module fusion. For instance, LoraHub Huang et al. (2023) dynamically composes pre-trained LoRA modules by optimizing their weights through few-shot examples from new tasks, leveraging black-box optimization techniques (e.g., CMA-ES) to achieve efficient adaptation without backpropagation. Meanwhile, SemLA Qorbani et al. (2025) introduces a training-free approach by directly comparing test images' visual features with known domain prototypes, using the resulting similarity to efficiently guide adapter retrieval and fusion. However, they either require unknowntask data or involve loading and unloading multiple LoRA adapters for each input, which can be computationally impractical. Moreover, they rely on inflexible, deterministic fusion for unseen tasks.

## 2.3 NEURAL NETWORK PARAMETERS GENERATION

Although generative modeling has advanced considerably, the direct generation of network weights for pre-trained models remains an emerging area of study Knyazev et al. (2021); Zhmoginov et al. (2022);

Knyazev et al. (2023); Wang et al. (2024a). Approaches such as generative hyper-representation learning Ha et al. (2016), neural network diffusion Peebles et al. (2022); Hu et al. (2021); Jin et al. (2024), and kernel density estimation—based methods Soro et al. (2024) have shown promise but remain fundamentally limited to small architectures and unconditional weight generation within fixed distributions. Consequently, these methods struggle to generalize to unseen tasks, constraining their broader applicability. While meta-learning frameworks Nava et al. (2022); Zhang et al. (2024) have enabled powerful joint model generation for visual recognition and few-shot learning, they often neglect the personalization and parameter diversity, restricting the generator's output to classifier heads rather than more flexible and expressive parameter sets, such as LoRA. ICM-LoRAShao et al. (2025) innovatively explores the parameter relationships among tasks through task vectors, but focuses on closed-world, task-specific enhancements of LoRA parameters. As a result, the question of whether one can rapidly generate efficient, user-intent-focused LoRA parameters in open-world settings remains unexplored.

## 3 THE PROPOSED MODEL

In this section, we begin with an overview of essential concepts for understanding semantic-guided LoRA parameter generation, followed by detailed descriptions of our proposed methods.

#### 3.1 PROBLEM DEFINATION AND PRELIMINARY

We define Zero-Shot Open-world Adaptation (ZSOA) as a novel and challenging task setting that requires models to generalize across semantically diverse tasks in open-world scenarios. Unlike conventional zero-shot learning or task transfer—often confined to fixed label spaces or narrow domains—ZSOA focuses on tasks that share a common structural format (e.g., image-text retrieval) but differ substantially in domain, content, or distribution. It emphasizes that queried tasks during inference time are drawn from an undefined and unbounded set, requiring rapid adaptation to novel tasks without access to raw data. This setting reflects realistic deployment scenarios, where task-level generalization must rely solely on prior experience.

In this work, we instantiate ZSOA in the context of fine-grained image-text retrieval, with each task formulated as a retrieval problem over a specific semantic category (e.g., animal species, flower type). Formally, given a set of fine-tuned LoRA module  $\mathcal W$  trained on known tasks  $\mathcal T$ , our goal is to adapt the model to unseen task  $\mathcal T^*$  without accessing any labeled image-text pairs. Let  $f(\mathcal T)$  denote the textual description of task  $\mathcal T$ , we learn a generator G that predicts LoRA parameters for  $\mathcal T^*$  based on semantic descriptions and  $\mathcal W$ :

$$W^* = G(f(\mathcal{T}^*), \mathcal{W}, f(\mathcal{T})) \tag{2}$$

The synthesized LoRA parameter  $\mathcal{W}^*$  is then used to modulate a frozen vision-language backbone, enabling it to perform the image-text retrieval task defined by  $\mathcal{T}^*$ . ZSOA thus extends traditional zero-shot learning by enabling parameter-level generalization , accommodating diverse user-instructed tasks with open-world queries.

# 3.2 LORA PARAMETER DATASET CONSTRUCTION

### 3.2.1 TASK-SPECIFIC LORA TRAINING

The first stage involves constructing a dataset of LoRA parameters. Consider a collection of N distinct tasks  $\mathcal{T} = \{\mathcal{T}_1, \dots, \mathcal{T}_N\}$ , where each  $\mathcal{T}_n$  corresponds to a specific retrieval task (e.g., image-text retrieval for Cat category). In this context, the task is dataset-agnostic—that is, it may originate from the same dataset or from a different one. Given a pre-trained vision-language model (VLM), we train a task-specific LoRA for each  $\mathcal{T}_n$  using corresponding image-text pairs, applying LoRA modules at consistent positions within the VLM. To ensure comparability and reduce variance, all LoRA modules are trained using identical network configurations by default. After training stabilizes, we extract and store the LoRA from the final M epochs, yielding task-specific parameter data:

$$\Delta \mathbf{W}_n = \{\Delta \mathbf{W}_n^m\}_{m=1}^M, \quad \mathbf{d}_n = f(\mathcal{T}_n), \tag{3}$$

where  $\Delta \mathbf{W}_n^m = [\mathbf{B}_n^m, \mathbf{A}_n^m]$  denotes the concatenation of LoRA parameters in layer-wise order for task n from the m-th saved epoch.  $\mathbf{d}_n$  denotes the textual description associated with  $\mathcal{T}_n$ , generated

using the template 'a photo of a <class name>', and encoded by the frozen CLIP text encoder  $f(\cdot)$  to serve as a global semantic representation. Collectively, these form the LoRA parameter dataset:

$$W = \{\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_N\},\tag{4}$$

where each element is optimized for image-text alignment within its respective task.

#### 3.2.2 LORA EXPERT REPOSITORY FORMATION

To construct a reliable expert LoRA space that simulates diverse LoRA resources, we first curate a representative subset of tasks from the full corpus  $\mathcal{W}$ , where each task is associated with its corresponding LoRA parameters and semantic embedding. For each selected task, we compute the mean LoRA parameters  $\mu_e$  over all available adaptations, yielding a distilled representation of its task-specific adaptation pattern. These averaged parameters, along with their associated semantic embeddings, constitute our expert repository:

$$W_{\text{expert}} = \{ (\boldsymbol{\mu}_e, \mathbf{d}_e) \mid e \in \mathcal{E} \}, \quad \boldsymbol{\mu}_e = \frac{1}{M} \Delta \mathbf{W}_e, \tag{5}$$

where  $\mathcal{E}$  represents the selected expert index and  $\mathcal{W}_{expert}$  serves as a compact yet expressive basis for capturing the essential characteristics of each knowledge domain. The remaining LoRA data are then split into training and evaluation sets for subsequent model development and evaluation.

#### 3.3 SEMANTIC-GUIDED LORA PARAMETER GENERATION

Given the collected LoRA repository  $W_{\text{expert}}$ , our SG-LoRA framework generates conditional LoRA parameters under semantic guidance. We first define the task semantics by selecting and combining the most relevant expert LoRAs from the repository using a sparse aggregator. A conditional variational autoencoder (CVAE)Sohn et al. (2015) is then trained to generate target LoRA parameters aligned with the corresponding task semantics. Once trained, the model can generalize to unseen tasks in an open-world query setting without further training.

#### 3.3.1 Construction of task semantics

Intuitively, not all experts contribute equally to an unseen task. Therefore, it is essential to identify and prioritize the most beneficial experts. Fortunately, the CLIP textual encoder is well-suited for this purpose, as it effectively captures rich semantic relationships across tasks. As shown in Eq. 3, we use each task's textual embedding as a global semantic descriptor for its LoRA parameters. For an unseen task  $\mathcal{T}^*$  with textual embedding  $\mathbf{d}^*$ , we compute cosine similarities with all expert embeddings  $\{d_e|e\in\mathcal{E}\}$  and select the top-K experts with the highest similarity scores to form a semantically tailored expert set, indexed by  $\mathcal{I}_{\text{top-}k}$ . Then, we normalize their similarity scores using the softmax function to obtain the fusion coefficients:

function to obtain the fusion coefficients:
$$\alpha_k = \frac{\exp\left(\text{sim}(\mathbf{d}^*, \mathbf{d}_k)/\tau\right)}{\sum_{k' \in \mathcal{I}_{\text{top-}k}} \exp\left(\text{sim}(\mathbf{d}^*, \mathbf{d}_{k'})/\tau\right)}, \quad k \in \mathcal{I}_{\text{top-}K}, \tag{6}$$

where  $\tau > 0$  is a temperature parameter. The semantic vector for task  $\mathcal{T}^*$  is computed as a weighted sum:

$$\boldsymbol{\mu}^* = \sum_{k \in \mathcal{I}_{\text{top-}K}} \alpha_k \cdot \boldsymbol{\mu}_k. \tag{7}$$

The attention strategy in Eq. 7 guides our model to understand the unseen task with the expert LoRAs from the textual perspective, resulting in high-quality semantic representation.

To capture the semantic diversity of the task, we also consider estimating the element-wise variance for task  $\mathcal{T}^*$  under the Law of Total Variance theory:

$$\boldsymbol{\sigma}^{*2} = \sum_{k=1}^{K} \alpha_k \boldsymbol{\sigma}_k^2 + \sum_{k=1}^{K} \alpha_k (\boldsymbol{\mu}_k - \boldsymbol{\mu}^*) \odot (\boldsymbol{\mu}_k - \boldsymbol{\mu}^*), \tag{8}$$

where  $\odot$  denotes element-wise multiplication and  $\sigma_k$  denotes the variance of the k-th expert. This flexible formulation enables the model to better reflect the statistical properties of new tasks, which is crucial for generative modeling. The estimated mean and variance guide the generative process, resulting in more accurate and context-aware outputs, thereby improving generalization and robustness. We leave the derivation of Eq. 8 in Appendix A.2. To simplify, we will use c to represent the task semantics of  $\mathcal{T}^*$  in the following descriptions, e.g.,  $c = \{\mu^*, \sigma^{*2}\}$ .

## 3.3.2 CONDITIONAL LORA PARAMETER GENERATION

We adopt a conditional variational autoencoder framework to generate target LoRA parameters based on the task semantics calculated above. Given a batch of training LoRA tensor  $\boldsymbol{X}$ , the encoder approximates the posterior distribution  $q(\boldsymbol{z}|\boldsymbol{X},c)$  using a multi-layer perceptron (MLP) that takes the  $\boldsymbol{X}$  to be reconstructed and the task semantics c as input. A latent code  $z \sim q(\boldsymbol{z}|\boldsymbol{X},c)$  is sampled and passed to the decoder along with c as condition to reconstruct the original input  $\boldsymbol{X}$ . Unlike traditional VAEs that adopt  $p(\boldsymbol{z}) = \mathcal{N}(0,\mathbf{I})$  as the prior distribution, we here develop a semantic-aware prior for each task  $p(\boldsymbol{z}|c)$ .  $p(\boldsymbol{z}|c)$  is parameterized with stacked MLPs, allowing the model to flexibly represent a task-specific prior distribution based on domain-level statistics.

The model is trained to maximize the evidence lower bound (ELBO), which consists of two terms: the reconstruction and the regularization term:

$$\mathcal{L}_{\text{CVAE}} = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{X},c)} \left[ \|\boldsymbol{X} - \hat{\boldsymbol{X}}\|^2 \right] + \lambda \cdot KL(q(\boldsymbol{z}|\boldsymbol{X},c)\|p(\boldsymbol{z}|c)), \tag{9}$$

where  $\hat{X}$  denotes the reconstructed LoRA parameters,  $KL(\cdot \parallel \cdot)$  is the Kullback-Leibler divergence, and  $\lambda$  controls the relative weight of the KL term. The first term encourages the decoder to reconstruct accurate LoRA parameters, while the second term regularizes the latent space to align the task-specific prior. During inference, a sample z is drawn from the prior distribution p(z|c), and the decoder generates the corresponding custom LoRA parameter.

## 4 EXPERIMENTS

#### 4.1 EXPERIMENTAL SETTINGS

**Datasets.** The proposed SG-Lora is evaluated on three benchmark datasets. Specifically, we use the widely adopted MS-COCO dataset Lin et al. (2014), a standard benchmark for image-text retrieval, known for its diverse scenes and rich linguistic annotations. To evaluate the model's generalization ability, we further include the OxfordPets dataset Parkhi et al. (2012) and the Flowers102 dataset Nilsback & Zisserman (2008), both of which are originally designed for fine-grained image classification. Moreover, given the inherent ambiguity and limited informativeness of MS-COCO captions and the absence of captions in the other two datasets, we construct synthetic textual descriptions using Qwen2-VL Wang et al. (2024b). For MS-COCO, we use the original training split but regenerate captions for each image, effectively creating a new image–caption dataset, and then divide it into training, validation, and test sets. More details are provided in Appendix A.1.

**Metrics.** The evaluation metric used is Recall@K (R@K), which quantifies the proportion of correct matches appearing in the top-K retrieved candidates. We report R@1, R@5, and R@10 for both image-to-text and text-to-image retrieval scenarios.

Implementation Details. We adopt CLIP ViT-B/16 as our backbone, injecting rank-2 LoRA adapters into the  $W_q$ ,  $W_k$ , and  $W_v$  projection matrices of every Transformer block in the visual encoder. Training is carried out with the Adam optimizer. The CVAE's encoder and prior network each consist of two-layer MLPs with ReLU activations, whereas the decoder is realized as a three-layer MLP with ReLU activations. We set the default values of M, K, and  $\lambda$  in the model to 100, 4, and 1, respectively. All experiments were performed on a single NVIDIA A6000 GPU.

## 4.2 Comparative Methods

To evaluate the effectiveness of the proposed method, we compare it with the following methods: (1) **Zero-Shot CLIP**: The original CLIP model without any adaptation of LoRA modules. (2) **Model Soups**: Consistent with Wortsman et al. (2022a) all LoRA experts in  $W_{\text{expert}}$  are uniformly averaged without considering their relevance to the target task. (3) **AdapterSoup**: The top-K experts from the expert repository are selected based on the semantic similarity vector and with equal weight, assigning each a coefficient of 1/K. This can be seen as an variant of Chronopoulou et al. (2023) revised to our setting. (4) **Top-K** LoRA Weighted: The top-K experts are selected based on the semantic similarity vector, and their weights are computed by applying a softmax function over the similarity scores for adaptive merging. (5) **SG-LoRA**: Our proposed method, which generates

task-specific LoRAs based on semantic proximity. (6) **Oracle:** For each task, LoRA parameters are trained individually on the specific dataset where we evaluate.

# 4.3 Main Results and Discussion

#### 4.3.1 IN-DATASET IMAGE-TEXT RETRIEVAL

We first conducted in-dataset evaluations on MS-COCO and OxfordPets dataset separately, with results shown in Table.1. Several key observations are as follows: 1). Compared to the Zero-Shot CLIP baseline and consistent with previous findings Qorbani et al. (2025), directly merging all experts in the expert repository leads to performance improvements. 2). Selecting semantically relevant experts, like those most related to the current query task from the repository, can further enhance performance. However, naively treating all selected experts equally may result in degraded performance. For instance, in Table 1, AdapterSoup underperforms compared to Top-K LoRA Weighted. This may be because different experts contribute unevenly to the target task. Assigning equal weight ignores these differences and may amplify noise from less relevant experts. By contrast, incorporating semantic weighting coefficients allows the fusion process to account for varying degrees of relevance, leading to more effective integration of expert knowledge and improved retrieval performance. 3). While merging-based approaches still fall short of the performance achieved by directly fine-tuning LoRA on the unseen task, SG-LoRA fully recovers the performance of oracle adapters.

Notably, SG-LORA even outperforms Oracle in certain cases—for example, R@1 for bidirectional retrieval on MS-COCO and R@1 for image-to-text retrieval on OxfordPets. This improvement gains from the efficient compression of expert LoRAs: our trained CVAE integrates the target LoRA using compact yet semantically rich task representations, enabling the generation of target-aligned LoRA parameters by modeling their distribution in the parameter space. Moreover, the Oracle LoRA fine-tuned on unseen tasks sometimes suffers from overfitting, especially when trained on a small set of image-caption pairs. Our SG-LORA helps mitigate this performance drop, likely due to its ability to generalize without relying on target-task data.

Table 1: Model Performance Comparison on MS-COCO and OxfordPets Datasets. The best results are highlighted in bold, and the second-best results are underlined.

Method	MS-COCO						OxfordPets					
		I2T			T2I			I2T			T2I	
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Zero-Shot CLIP Oracle	66.43 72.45	84.31 88.91	89.14 93.41	41.66 53.10	64.63 76.47	73.01 83.97	40.45 55.84	66.27 81.84	77.53 89.13	26.03 40.99	50.66 70.41	62.98 80.39
Model Soups AdapterSoup Top-K LoRA Weighted SG-LoRA	69.37 70.70 71.55 <b>74.31</b>	85.96 86.57 <u>87.54</u> <b>88.78</b>	90.95 91.09 <u>91.69</u> <b>92.50</b>	47.38 48.64 49.85 <b>54.42</b>	69.54 70.51 71.79 <b>75.45</b>	77.97 78.79 <u>79.66</u> <b>82.18</b>	52.54 52.59 53.96 <b>57.15</b>	77.80 78.52 <u>79.41</u> <b>80.40</b>	85.59 86.09 86.53 <b>88.04</b>	33.51 34.05 35.42 37.62	61.77 62.70 64.08 <b>67.16</b>	72.93 73.93 74.99 <b>77.44</b>

## 4.3.2 Cross-Dataset Image-Text Retrieval

Given the flexibility of SG-LORA, we conducted a more challenging cross-dataset evaluation. As shown in Table.2 and 6, SG-LORA consistently outperforms merging-based approaches in these settings. Interestingly, we also observed that models trained on MS-COCO were able to generate LoRA parameters that, in some cases, outperformed those trained directly on the OxfordPets (e.g., a relative improvement of 1.22% in T2I R@1). This may be attributed to the richer expert knowledge available in MS-COCO, which, due to its broader data diversity, enables more extensive exploration of the parameter space during generation, a capability not achievable when generating within the narrower scope of OxfordPets. Another possible reason is that the uniform LoRA training configuration used across datasets (as described in Section 3.2.1) may not be optimal for OxfordPets. Conversely, we find that when the generation model is trained on OxfordPets and applied to MS-COCO, its performance is generally worse than that of models trained on MS-COCO.

Table 2: Cross-Dataset Generalization Performance: Bidirectional Evaluation between MS-COCO and OxfordPets. The best results are highlighted in bold, and the second-best results are underlined.

Method	$MS\text{-}COCO \rightarrow OxfordPets$					OxfordPets $\rightarrow$ MS-COCO						
				T2I			I2T			T2I		
	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10	R@1	R@5	R@10
Zero-Shot CLIP Oracle	40.45 55.84	66.27 81.84	77.53 89.13	26.03 40.99	50.66 70.41	62.98 80.39	66.43 72.45	84.31 88.91	89.14 93.41	41.66 53.10	64.63 76.47	73.01 83.97
Model Soups AdapterSoup Top-K LoRA Weighted SG-LoRA	44.67 45.96 48.13 <b>55.41</b>	70.91 71.83 <u>73.43</u> <b>80.73</b>	80.78 81.42 <u>82.73</u> <b>87.33</b>	30.45 30.88 <u>33.34</u> <b>38.84</b>	56.77 57.32 <u>59.53</u> <b>66.77</b>	68.52 69.08 70.89 <b>76.69</b>	68.58 68.74 68.75 <b>70.81</b>	85.67 85.83 85.77 <b>86.83</b>	90.62 90.63 <u>90.67</u> <b>91.41</b>	44.09 44.19 44.60 46.50	66.55 66.58 66.76 <b>68.73</b>	75.08 <u>75.31</u> 75.25 <b>77.19</b>

#### 4.4 GENERALIZATION IN STANDARD IMAGE-TEXT RETRIEVAL

Considering another complex scenario where the test task may contain image-text pairs from multiple categories, we evaluated retrieval performance on the Flickr30K test set Plummer et al. (2015). Since this dataset has no clear category distinction during retrieval, we randomly selected one caption per image, fed it into the CLIP textual encoder to obtain the textual embedding, and then calculated the mean value as the task description for retrieval. The experimental results are shown in Figure.2, where SG-Lora outperforms Zero-Shot CLIP. Additionally, SG-Lora trained on MS-COCO achieves better results than that trained on OxfordPets. This is because MS-COCO provides more comprehensive expert knowledge, covering a wider range of categories, while OxfordPets primarily focuses on fine-grained distinctions within just two broad categories—cats and dogs. This also indicates that when semantic guidance is more powerful and comprehensive, or more relevant to the downstream task, the generated Lora parameters are also superior.

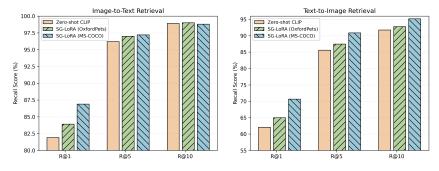


Figure 2: Simulation of task-agnostic evaluation on a general image-text retrieval dataset. We train models on MS-COCO and OxfordPets, respectively, and evaluate them on the Flickr30K test set.

#### 4.5 ABLATION STDUY

**Impact of expert repository configuration.** In the experiment of Table. 2, we observed that the *Cat* expert from the MS-COCO dataset was frequently selected as a semantic condition. To further evaluate the impact of semantically salient experts, we assessed how the inclusion of the *Cat* expert from MS-COCO influences retrieval performance on two unseen cat classes on OxfordPets. As shown in Table. 3, incorporating the *Cat* expert into the expert repository improves performance in most cases, particularly for text-to-image retrieval. This highlights the effectiveness of semantically guided expert selection. This finding also demonstrates the flexibility of our method in constructing task-adaptive expert repositories—particularly when a richer pool of LoRA resources is available.

Table 3: Ablation on expert repository strategy for cross-dataset evaluation. We evaluate how the MS-COCO *Cat* expert affects retrieval on two unseen OxfordPets cat tasks (marked with gray text).

Expert	Egyp	tian Ma	и <b>I2T</b>	Egyp	tian Ma	и <b>Т2</b> І	Expert	P	ersian <b>I</b>	2T	Pe	rsian <b>T</b>	2I
strategy	R@1	R@5	R@10	R@1	R@5	R@10	strategy	R@1	R@5	R@10	R@1	R@5	R@10
							w/o Cat expert						
w/ Cat expert	37.11	63.92	72.16	15.21	35.05	46.91	w/ Cat expert	47.00	79.65	86.00	36.75	64.25	73.75

Table 4: Top-4 expert for Yorkshire Terrier category under mixed-source experts configurations

Expert Task	Source Dataset	Contribution
Scottish Terrier	OxfordPets	0.9221
Dog	MS-COCO	0.0692
Cat	MS-COCO	0.0082
American Bulldog	OxfordPets	0.0005

Given that our model supports open-world expert repository construction, we further conducted a case study where experts from OxfordPets, MS-COCO, and Flowers102 were combined into a mixed-source repository. We also combined training data from both datasets to train the SG-LORA model accordingly. Figure.3 presents a comparison between single-source and mixed-source expert configurations. Additionally, we present the top-4 experts selected by SG-LORA under the mixed-source setup for the unseen *Yorkshire Terrier* category, along with their corresponding weights in Table.4. As shown, the mixed-source experts yield better performance than the single-source experts. The performance even surpasses that of the oracle LoRA model in text-to-image retrieval. These demonstrate the potential of our method in more realistic, real-time application scenarios, where expert repositories are constructed dynamically from heterogeneous data sources.

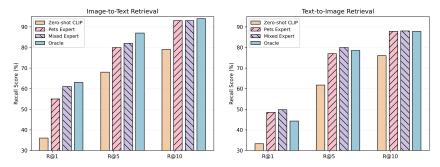


Figure 3: Comparison on expert repository configuration: single-source experts vs. mixed-source experts.

Impact of the number of experts. We conduct ablation studies on the number of experts K used in task semantic construction. As shown in Figure.4, using too few experts leads to insufficient knowledge for generalizing to unseen tasks, while increasing K incorporates more semantic information but may also introduce irrelevant context. Overall, K=4 shows a good balance between expert diversity and relevance, delivering good performance on both datasets.

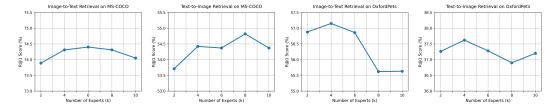


Figure 4: Ablaton study on the number of experts.

### 5 CONCLUSION

In this work, we introduce a novel and challenging task setting—Zero-Shot Open-world Adaptation (ZSOA)—which requires models to generalize across semantically diverse tasks in open-world scenarios. To achieve this, we propose a flexible and efficient approach that dynamically generates task-specific LoRA parameters guided by available LoRA resources. By identifying the most relevant expert knowledge based on semantic similarity and leveraging task semantics in a conditional generative framework, our SG-LoRA models the distribution of unseen parameters in a tuning-free manner. The inherent stochasticity of our generation process further introduces diversity, enhancing adaptability to previously unseen tasks. SG-LoRA is scalable and naturally privacy-preserving, making it well-suited for deployment in sensitive and dynamic real-world environments.

## REPRODUCIBILITY STATEMENT

We provide detailed descriptions of the model, training procedure, and evaluation in the main text. Additional implementation details, hyperparameters, and ablation studies are included in the Appendix A.

## REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Alexandra Chronopoulou, Matthew E Peters, Alexander Fraser, and Jesse Dodge. Adaptersoup: Weight averaging to improve generalization of pretrained language models. *arXiv preprint arXiv:2302.07027*, 2023.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- Civitai. Civitai: The Home of Open-Source Generative AI. https://civitai.com/, 2024. Accessed: November 2024.
- Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. Parameter-efficient fine-tuning of large-scale pre-trained language models. *Nature machine intelligence*, 5(3):220–235, 2023.
- David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. arXiv preprint arXiv:1609.09106, 2016.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Wei Hu, QiHao Zhao, Yangyu Huang, and Fan Zhang. P-diff: Learning classifier with noisy labels based on probability difference distributions. In 2020 25th International Conference on Pattern Recognition (ICPR), pp. 1882–1889. IEEE, 2021.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition. *arXiv preprint arXiv:2307.13269*, 2023.
- HuggingFace. Hugging Face The AI community building the future. https://huggingface.co/, 2024. Accessed: November 2024.
- Gabriel Ilharco, Mitchell Wortsman, Samir Yitzhak Gadre, Shuran Song, Hannaneh Hajishirzi, Simon Kornblith, Ali Farhadi, and Ludwig Schmidt. Patching open-vocabulary models by interpolating weights. *Advances in Neural Information Processing Systems*, 35:29262–29277, 2022.
- Dong-Hwan Jang, Sangdoo Yun, and Dongyoon Han. Model stock: All we need is just a few fine-tuned models. In *European Conference on Computer Vision*, pp. 207–223. Springer, 2024.
- Xiaolong Jin, Kai Wang, Dongwen Tang, Wangbo Zhao, Yukun Zhou, Junshu Tang, and Yang You. Conditional lora parameter generation. *arXiv preprint arXiv:2408.01415*, 2024.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
  - Boris Knyazev, Michal Drozdzal, Graham W Taylor, and Adriana Romero Soriano. Parameter prediction for unseen deep architectures. *Advances in Neural Information Processing Systems*, 34: 29433–29448, 2021.

- Boris Knyazev, Doha Hwang, and Simon Lacoste-Julien. Can we scale transformers to predict parameters of diverse imagenet models? In *International Conference on Machine Learning*, pp. 17243–17259. PMLR, 2023.
  - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
    - Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behavioral and brain sciences*, 40:e253, 2017.
    - Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
    - Weishi Li, Yong Peng, Miao Zhang, Liang Ding, Han Hu, and Li Shen. Deep model fusion: A survey. arXiv preprint arXiv:2309.15698, 2023.
    - Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pp. 740–755. Springer, 2014.
    - Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
    - Elvis Nava, Seijin Kobayashi, Yifei Yin, Robert K Katzschmann, and Benjamin F Grewe. Metalearning via classifier (-free) diffusion guidance. *arXiv preprint arXiv:2210.08942*, 2022.
    - Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In 2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing, pp. 722–729. IEEE, 2008.
    - 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.
    - William Peebles, Ilija Radosavovic, Tim Brooks, Alexei A Efros, and Jitendra Malik. Learning to learn with generative models of neural network checkpoints. *arXiv preprint arXiv:2209.12892*, 2022.
    - Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015.
    - Reza Qorbani, Gianluca Villani, Theodoros Panagiotakopoulos, Marc Botet Colomer, Linus Härenstam-Nielsen, Mattia Segu, Pier Luigi Dovesi, Jussi Karlgren, Daniel Cremers, Federico Tombari, et al. Semantic library adaptation: Lora retrieval and fusion for open-vocabulary semantic segmentation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 9804–9815, 2025.
    - 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.
    - Yihua Shao, Minxi Yan, Yang Liu, Siyu Chen, Wenjie Chen, Xinwei Long, Ziyang Yan, Lei Li, Chenyu Zhang, Nicu Sebe, et al. In-context meta lora generation. *arXiv preprint arXiv:2501.17635*, 2025.
  - Donald Shenaj, Ondrej Bohdal, Mete Ozay, Pietro Zanuttigh, and Umberto Michieli. Lora. rar: Learning to merge loras via hypernetworks for subject-style conditioned image generation. *arXiv* preprint arXiv:2412.05148, 2024.
    - Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems*, 28, 2015.

- Bedionita Soro, Bruno Andreis, Hayeon Lee, Wonyong Jeong, Song Chong, Frank Hutter, and Sung Ju Hwang. Diffusion-based neural network weights generation. *arXiv preprint arXiv:2402.18153*, 2024.
  - Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Vl-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5227–5237, 2022.
  - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
  - Kai Wang, Dongwen Tang, Boya Zeng, Yida Yin, Zhaopan Xu, Yukun Zhou, Zelin Zang, Trevor Darrell, Zhuang Liu, and Yang You. Neural network diffusion. arXiv preprint arXiv:2402.13144, 2024a.
  - Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024b.
  - Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International conference on machine learning*, pp. 23965–23998. PMLR, 2022a.
  - Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 7959–7971, 2022b.
  - Xun Wu, Shaohan Huang, and Furu Wei. Mixture of loRA experts. In *The Twelfth International Conference on Learning Representations*, 2024.
  - Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36:7093–7115, 2023.
  - Baoquan Zhang, Chuyao Luo, Demin Yu, Xutao Li, Huiwei Lin, Yunming Ye, and Bowen Zhang. Metadiff: Meta-learning with conditional diffusion for few-shot learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 16687–16695, 2024.
  - Jinghan Zhang, Junteng Liu, Junxian He, et al. Composing parameter-efficient modules with arithmetic operation. *Advances in Neural Information Processing Systems*, 36:12589–12610, 2023a.
  - Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adalora: Adaptive budget allocation for parameter-efficient finetuning. *arXiv preprint arXiv:2303.10512*, 2023b.
  - Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023c.
  - Ziyu Zhao, Tao Shen, Didi Zhu, Zexi Li, Jing Su, Xuwu Wang, Kun Kuang, and Fei Wu. Merging loras like playing lego: Pushing the modularity of lora to extremes through rank-wise clustering. arXiv preprint arXiv:2409.16167, 2024.
  - Andrey Zhmoginov, Mark Sandler, and Maksym Vladymyrov. Hypertransformer: Model generation for supervised and semi-supervised few-shot learning. In *International Conference on Machine Learning*, pp. 27075–27098. PMLR, 2022.
  - Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision (IJCV)*, 2022.
  - Xiandong Zou, Mingzhu Shen, Christos-Savvas Bouganis, and Yiren Zhao. Cached multi-lora composition for multi-concept image generation. *arXiv preprint arXiv:2502.04923*, 2025.

## A APPENDIX

#### A.1 EXPERIMENTAL DETAILS

#### A.1.1 IMAGE-TEXT RETRIEVAL DATASET

In this work, we generate fine-grained captions for the image-text retrieval task using Qwen2-VL. As illustrated in Figure.5, our in-context learning approach provides the Multimodal Large Language Model (MLLM) with a small number of demonstration examples, enabling it to generate detailed captions for target images. Using this method, we have produced four diverse and highly relevant captions for each image in the entire OxfordPets dataset and the Flowers102 dataset, a subset of the MS-COCO dataset. These fine-grained descriptions serve as a valuable resource for downstream tasks such as fine-grained image-text retrieval and facilitate further research in ZSOA.

# You are an image description assistant. Please analyze the given image carefully and provide a detailed description of its contents [Demonstration 1] Image: <image\_path1> Input: A photo of pug Output: The pug has a fawn coat, black mask, and wrinkled face with a short-muzzled expression. It wears a green collar with a tag, sits on a rock, and has a solemn gaze. [Demonstration 2] Image: <image path2> Input: A photo of birman. Output: The Birman cat has a sleek cream and brown coat, striking blue eyes, and a dark brown nose. It wears a blue collar and is perched on a scratching post. [Target Query] Image: <image\_path> Input: A photo of {class name} Output: fine-grained caption for the target image

Figure 5: Examples of caption generation using Qwen2-VL.

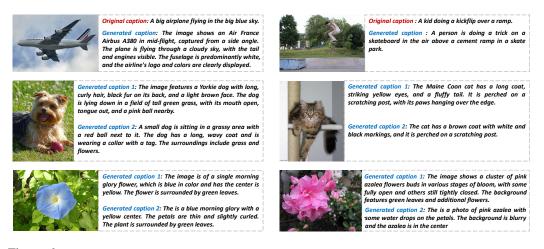


Figure 6: Illustration of generated captions: results on the MS-COCO (top), OxfordPets (middle), and Flowers102 (bottom) datasets.

As shown in Figure.6, we present examples of image-text pairs from three datasets. Compared to the original captions in the COCO dataset (labeled as *Original Caption* in the top), our generated captions more accurately reflect the content of the corresponding images. Thus, these datasets enable us to assess the model's robustness both in in-domain retrieval and in generalizing to novel domains and compositional scenarios.

Table 5: Selected Expert Tasks for Each Datasets

MS-COCO	OxfordPets	Flowers102
Airplane	Abyssinian	Sweet pea
Truck	American bulldog	Tiger lily
Traffic Light	American Pit Bull Terrier	Monkshood
Cat	Birman	King Protea
Horse	Bombay	Corn Poppy
Giraffe	British Shorthair	Daffodil
Handbag	German Shorthaired	Sunflower
Snowboard	Havanese	Osteospermum
Wine Glass	Keeshond	Anthurium
Banana	Leonberger	Hibiscus
Hot Dog	Newfoundland	Desert-Rose
Laptop	Scottish Terrier	Mallow

### A.1.2 LORA EXPERT REPOSITORY DETAILS

We present the expert tasks included in each dataset's expert repository in Table. 5, as set in the main manuscript.

#### A.2 AGGREGATION FUNCTION FOR SEMATIC PRIOR

We begin with the classical identity in probability theory known as the Law of Total Variance: for a random vector  $\Theta$  and a discrete conditioning variable C, the following identity holds:

$$Var(\Theta) = \mathbb{E}_C \left[ Var(\Theta \mid C) \right] + Var_C \left( \mathbb{E}[\Theta \mid C] \right), \tag{10}$$

**Proof Sketch:** Using the identity  $Var(\Theta) = \mathbb{E}[\Theta^2] - (\mathbb{E}[\Theta])^2$  and the Law of Iterated Expectations, we write:

$$\mathbb{E}[\Theta^2] = \mathbb{E}_C \left[ \mathbb{E}[\Theta^2 \mid C] \right], \quad \mathbb{E}[\Theta] = \mathbb{E}_C \left[ \mathbb{E}[\Theta \mid C] \right], \tag{11}$$

From the definition of conditional variance,  $Var(\Theta \mid C) = \mathbb{E}[\Theta^2 \mid C] - (\mathbb{E}[\Theta \mid C])^2$ , we substitute and obtain:

$$\operatorname{Var}(\Theta) = \mathbb{E}_{C} \left[ \operatorname{Var}(\Theta \mid C) \right] + \mathbb{E}_{C} \left[ \left( \mathbb{E}[\Theta \mid C] \right)^{2} \right] - \left( \mathbb{E}_{C} \left[ \mathbb{E}[\Theta \mid C] \right] \right)^{2}, \tag{12}$$

where the last two terms equal  $\operatorname{Var}_C(\mathbb{E}[\Theta \mid C])$ .

**Task Semantic via Aggregation.** Suppose we have K expert tasks, each providing a LoRA parameter set with mean  $\mu_i$  and variance  $\sigma_i^2$ . For an unseen task  $\mathcal{T}^*$  with expert weight vector  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_k]^T$ , where  $\alpha_i \geq 0$  and  $\sum_i \alpha_i = 1$ , we construct the prior as a weighted combination. The prior mean is:

$$\mu^* = \sum_{i=1}^k \alpha_i \mu_i,\tag{13}$$

Applying the Law of Total Variance, the element-wise prior variance for  $\mathcal{T}^*$  is:

$$\sigma^{2^*} = \sum_{k=1}^{K} \alpha_k \sigma_k^2 + \sum_{k=1}^{K} \alpha_k (\mu_k - \mu^*) \odot (\mu_k - \mu^*), \tag{14}$$

where ⊙ denotes element-wise multiplication. In Eq.14, the first term, which represents the expectation of the conditional variances (computed as a weighted average of the per-task variances), quantifies the uncertainty within each expert task. The second term, which is the variance of the conditional means, quantifies the dispersion or uncertainty among different expert tasks. Eq.13 and Eq.14 together define the task semantic, providing a more informative prior than modeling the parameters as a standard Gaussian, as it ensures that both intra-task variation and inter-task uncertainty are captured when modeling the latent prior for each unseen task.

#### A.3 ADDITIONAL RESULTS ON CROSS-DATASET IMAGE-TEXT RETRIEVAL

As shown in Fig. 6, we present cross-dataset evaluation with the model trained on MS-COCO and tested on Flowers102. It can be observed that when there exists a noticeable gap between the expert task and the inference task, although the performance generally surpasses Zero-Shot CLIP, both AdapterSoup and Top-K LoRA Weighted fail to significantly outperform Model Soup. In contrast, the semantics-guided SG-Lora demonstrates a strong ability in generating high-performance Lora parameters.

Table 6: Cross-dataset evaluation with the model trained on MS-COCO and tested on Flowers 102. The best results are highlighted in bold, and the second-best results are underlined.

Method	Į I	2T Metr	ics	<b>T2I Metrics</b>			
11201100	R@1	R@5	R@10	R@1	R@5	R@10	
Zero-Shot CLIP	22.30	53.57	69.41	16.33	46.79	68.49	
Oracle	26.23	59.17	77.55	21.57	57.56	80.13	
Model Soups	23.50	54.84	73.85	17.23	50.96	73.48	
AdapterSoup	24.21	54.26	72.83	17.76	50.93	<u>73.83</u>	
Top-K LoRA Weighted	23.98	52.83	71.91	17.63	50.69	73.74	
SG-LoRA	<b>26.83</b>	<b>56.63</b>	<b>74.16</b>	20.52	<b>53.71</b>	<b>76.69</b>	

#### A.4 ADDITIONAL RESULTS ON CLASSIFICATION TASK

In addition to the image-text retrieval task, we also explored the performance of SG-Lora on classification tasks. Specifically, we selected 20 superclasses from CIFAR-100 Krizhevsky et al. (2009) as 20 distinct tasks, with each task corresponding to a 5-class classification. The superclasses were chosen as defined in the official CIFAR-100 hierarchy. We selected 8 of these tasks as expert Lora and performed inference on 6 unseen tasks, with the results presented in Table.7. We observed that, compared to Zero-Shot CLIP, both Model Soups and the selection of expert parameters to construct Lora improved classification performance on unseen tasks. Furthermore, SG-Lora achieved the best performance, indicating that our method is also applicable to classification tasks.

However, we found that SG-Lora still has a significant performance gap compared to fine-tuning Lora directly on unseen tasks using training data. We hypothesize that this is because, for classification tasks, the independence (or orthogonality, where decision boundaries between tasks may be orthogonal) between tasks is much more pronounced. In contrast, for image-text retrieval tasks, which operate at a finer granularity, there exist stronger inter-task correlations—this inherent property facilitates feasible Lora parameter transfer. As shown in Figure.8, we visualized the average Lora parameters for the 20 tasks and found that, although semantically similar tasks are closer in parameter space, their distances remain relatively sparse compared to image-text retrieval tasks in Figure.7. We plan to conduct further exploration on this issue in the future.

Additionally, we conducted an ablation study on textual description for the classification task, as shown in Table.8. It can be observed that, for the current task setting, more detailed textual descriptions, which account for the specific categories within each task, can better capture the semantic relationships between tasks, thus leading to improved performance.

Table 7: Model Performance of image classification on CIFAR-100 superclass.

Method	Accuracy
Zero-Shot CLIP	72.30
Oracle	91.43
Model Soups	75.63
AdapterSoup	72.60
Top-K LoRA Weighted	72.70
SG-LoRA	77.50

Table 8: Ablation study of textual description on CIFAR-100 superclass Classification

Description	Accuracy (%)
A photo of a <superclass name=""></superclass>	75.77
This is a classification task for recognizing <i><super< i=""></super<></i>	class name>,
which includes class_1,, class_5	77.50

## A.5 MORE ABLATIONS AND ANALYSIS

Table 9: Ablation study on modalities of semantic prior condition

Condition	Met	trics	Dataset		
	I2T R@1	T2I R@1	2 444500		
Visual	73.16	52.70	MS-COCO		
Textual	74.31	54.42	MS-COCO		
Visual	86.30	70.12	Flickr30K		
Textual	86.90	70.66	FIICKIOUK		
Condition	Met	trics	Dataset		
	Accı	ıracy			
Visual	73	.83	CIFAR-100		
Textual	77	.50	CIFAR-100		

Ablation on modalities of semantic priors. The construction of semantic priors serves as the foundation for our SG-Loral In Table.9, we compare the performance of semantic conditions across different modalities, where the conditional task description from each modality directly influences the selection of experts for unseen tasks by affecting  $\alpha_k$  in Eq.6. The visual condition is obtained by averaging the visual embeddings of training set images within each task dataset using a frozen CLIP visual encoder. Experimental results show that the textual condition better captures the semantic relationships between tasks. This could be attributed to two factors. Firstly, the high degree of condensation in textual semantics might play a role. Secondly, the disparities between training and test set images (or the presence of noisy images) within a task could result in inaccuracies in the visual prior condition.

**Impact of expert repository configuration.** Consistent with Table 3, we further evaluated the impact of incorporating the *Dog* expert from the MS-COCO dataset on retrieval performance for two unseen dog tasks in the OxfordPets dataset. As demonstrated in Table 10, including the *Dog* expert in the expert repository consistently improves the performance.

**Visulization of LoRA parameters.** To further investigate the parameter diversity of SG-Lora, we conducted evaluations on the unseen 'Zebra' task from the MS-COCO dataset at different training stages and visualized the generated Loras using t-SNE. As shown in Figure 9, we observe that the distribution of Loras generated by SG-Lora gradually aligns with that of Oracle Loras (directly trained in image-caption pairs), while still preserving diversity rather than extensively overlapping with the Oracle. This indicate that, by injected stochasticity, our method effectively explores the high-performance Loras in the parameter space.

Additionally, by examining the bottom subfigure, we observe that in parameter space, both Adapter-Soup (Tok-K LoRA Merging) and Tok-K LoRA Weighted lie closer to the mean of the Oracle LoRA compared to Model Soup. This is because the latter treats all experts equally, whereas the former two provide more informative semantic guidance, allowing the LoRA parameters to be better tailored to the current unseen task.

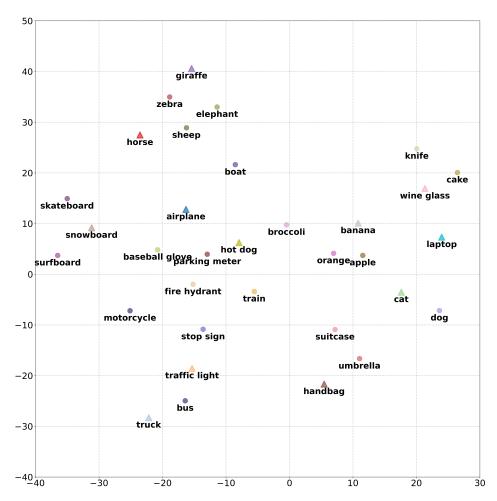


Figure 7: t-SNE visualization of the averaged LoRA parameters on MS-COCO dataset for image-text retrieval task. Triangular markers indicate expert LoRAs. Semantically similar LoRA parameters tend to cluster closely together.

Table 10: Ablation study on expert repository strategy for cross-dataset evaluation. We assess the impact of the *Dog* expert from MS-COCO dataset on the retrieval performance in two unseen dog tasks from OxfordPets dataset (marked with gray text).

Expert	I	2T Metri	ics	T2I Metrics						
strategy	R@1	R@5	R@10	R@1	R@5	R@10				
Pug										
w/o <i>Dog</i> expert w <i>Dog</i> expert	45.00 46.00	66.00 67.00	78.00 78.00	32.00 32.00	51.00 51.75	61.25 62.75				
	Сһіһиаһиа									
w/o <i>Dog</i> expert w <i>Dog</i> expert	64.00 66.00	89.00 90.00	93.00 95.00	52.50 55.25	81.00 81.25	88.25 89.75				

# A.6 LIMITATIONS AND FUTURE WORK

Although our method achieves strong performance, there remain several directions for future exploration. First, the current expert repository assumes that all LoRA experts share a same structure, such as having the same rank. In practice, however, publicly available expert LoRAs may be heterogeneous.

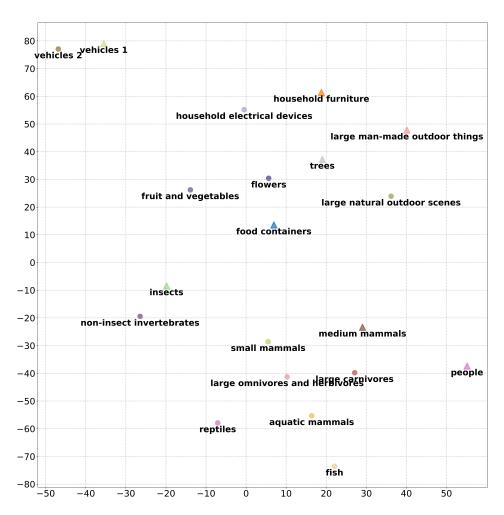


Figure 8: t-SNE visualization of the averaged LoRA parameters on CIFAR-100 for classification task. Triangular markers indicate expert LoRAs. Although semantically similar LoRA parameters tend to cluster closely together, the distribution appears sparser compared to the MS-COCO dataset due to stronger task independence.

Extending the conditional LoRA generation process to account for such heterogeneity could enable the integration of a broader range of expert knowledge. Moreover, our current approach employs the standard CLIP textual template as the task description. Leveraging task descriptions derived from LLM-based reasoning may provide more accurate and semantically aligned guidance.

## A.7 THE USE OF LARGE LANGUAGE MODELS

In this work, Large Language Models (LLMs) were used exclusively for language polishing and spelling correction.

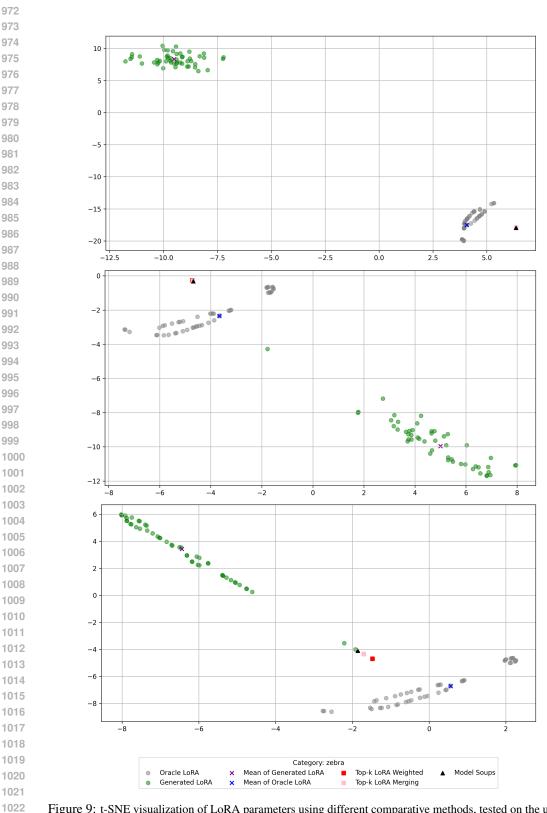


Figure 9: t-SNE visualization of LoRA parameters using different comparative methods, tested on the unseen 'Zebra' category at different training stages. For SG-LORA and Oracle LoRA, we randomly sampled 50 samples each. The subfigures from top to bottom represent increasing CVAE training epochs.