

VeCAF: Vision-language Collaborative Active Finetuning with Training Objective Awareness

Rongyu Zhang^{1,2,*}, Zefan Cai^{2,6,*}, Huanrui Yang^{3,*}, Zidong Liu⁴, Denis Gudovskiy⁵, Tomoyuki Okuno⁵, Yohei Nakata⁵, Kurt Keutzer³, Baobao Chang⁶, Yuan Du^{1,†}, Li Du¹, Shanghang Zhang^{2,†}
¹ Nanjing University, ² National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University, ³ University of California, Berkeley, ⁴ Tsinghua University, ⁵ Panasonic, ⁶ School of Software and Microelectronics, Peking University

Abstract

Finetuning a pretrained vision model (PVM) is a common technique for learning downstream vision tasks. However, the conventional finetuning process with randomly sampled data points results in diminished training efficiency. To address this drawback, we propose a novel approach, *Vision-language Collaborative Active Finetuning* (VeCAF). With the emerging availability of labels and natural language annotations of images through web-scale crawling or controlled generation, VeCAF makes use of these information to perform parametric data selection for PVM finetuning. VeCAF incorporates the finetuning objective to select significant data points that effectively guide the PVM towards faster convergence to meet the performance goal. This process is assisted by the inherent semantic richness of the text embedding space which we use to augment image features. Furthermore, the flexibility of text-domain augmentation allows VeCAF to handle out-of-distribution scenarios without external data. Extensive experiments show the leading performance and high computational efficiency of VeCAF that is superior to baselines in both in-distribution and out-of-distribution image classification tasks. On ImageNet, VeCAF uses up to $3.3\times$ less training batches to reach the target performance compared to full finetuning, and achieves an accuracy improvement of 2.7% over the state-of-the-art active finetuning method with the same number of batches.

1. Introduction

Deep learning has made significant progress in the field of computer vision that is typically attributed to the use of large-scale models and datasets [8, 26]. Hence, training such models from scratch is a time-consuming process and demands extensive amount of data. To address this, the pretraining-finetuning [9, 10, 34, 41] paradigm has been recognized as a favorable approach for both vision and language tasks. For vision tasks, a model can be first trained on abundant supervised or unsupervised data and be saved as a pretrained vision model (PVM) [2, 6, 29]. Then, the PVM is finetuned on a labeled dataset for a specific downstream task. By capitalizing on ample pretraining data and conserving valuable training resources during the finetuning stage, this paradigm has archived remarkable adoption in practical applications.

In the real-world deployments, practitioners aim to adapt deep learning models to a certain scenario or tune a model towards a specific performance target with minimal efforts for data selection and with quick training. Training with all available downstream task data can be not only costly but

also can lead to a biased or degraded performance in case of improperly collected data. This motivates the proposal of a data selection framework that can actively select the optimal data subset for finetuning. Previous work on active learning [4, 44] has shown the feasibility of PVM finetuning with only a small subset (e.g., less than 5%) of training data while achieving high performance metrics in downstream tasks. However, this line of work is often limited by the setting of low label availability which hinders its effectiveness to meet the user-specified objectives.

With the growing feasibility to gather large amounts of images with labels and natural language captions in the target domain through web-scale data crawling [36] or controlled generation [25], we find it is practical to explore a novel setting of active finetuning using *annotated* data. Then, we aim to select an optimal subset of training data for finetuning while having faster convergence and/or higher performance metrics as shown in Figure 1. The selection can further be performed in a loop to accommodate the changing model performance during the finetuning process. To this end, we propose to perform an Objective-aware Data Selection (ODS) using a parameterized data selection model. This ODS model reweighs training data distribution

* Equal contributions; † Corresponding authors

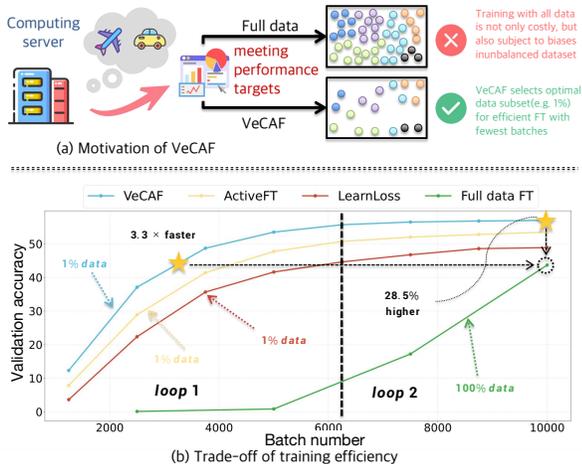


Figure 1. (a) Motivation of VeCAF. We select the optimal subset from a large labeled training set for efficient finetuning (FT) towards a user-specified objective. (b) Training curve comparison on ImageNet-1K validation set. All baselines select 1% of data in each FT loop with the exception of a conventional setup with full-data FT. VeCAF achieves the target accuracy faster with significantly fewer training batches and achieves higher accuracy with the same training cost.

according to the downstream objective and selects a subset that is both diverse and representative to the task.

The pursuit of objective-awareness brings new challenges to the data selection. Intuitively, images with misleading object appearances and complicated backgrounds, as illustrated in Figure 2, often provide more informative supervision. However, the image features extracted by the PVM may not fully capture all the semantic information present in the image. Therefore, PVM image features may miss useful information for the data selection and finetuning process as in previous ActiveFT [44] work. To address this limitation, we propose to leverage semantically rich language embedding spaces of the text encoders (e.g., CLIP [34], mT5 [45], BERT [9] etc.) in our novel Vision-language Collaborative Active Fine-tuning (VeCAF) approach. Specifically, we extract the text embeddings of the captions associated with each image. These captions may be sourced directly from original datasets such as COCO-Caption or alternatively generated by a multimodal Large Language Model (LLM) e.g., BLIP-2 [22]. Then, we propose a Cross-attentive Embedding Augmentation (CEA) to augment the image features extracted by the PVM such that the augmented features can focus more on the rich semantic information of the training samples. Therefore, the CEA facilitates both the active data selection and the finetuning.

Empirically, we demonstrate improved efficiency and performance across different scenarios. VeCAF is evaluated on three image classification datasets including CIFAR-

10 [19], Caltech101 [13], and ImageNet-1K [8] with the pretraining-finetuning paradigm where base models are pre-trained on ImageNet-1K. Results on ImageNet demonstrate that VeCAF can significantly accelerate the PVM convergence speed to the target performance and saves up to $3.3\times$ computational cost when compared to finetuning with all training set. Importantly, VeCAF also addresses out-of-distribution (OOD) scenarios, where we augment image features by target-domain text embeddings derived from the generated image captions. By leveraging the alignment between text and images embeddings, VeCAF increases the likelihood of selecting images that possess the characteristic features of the target domain from the training dataset as shown in Figure 3. In addition, we verify the OOD generalization ability on the corrupted ImageNet-C dataset where VeCAF improves accuracy by over 6% when compared to state-of-the-art active learning methods. Our main contributions are summarized as follows:

- We propose a novel framework, VeCAF, to improve computational efficiency of PVM finetuning using both the training objective and the language-embedded knowledge.
- We propose the Objective-aware Data Selection (ODS), where a parameterized data selection model is optimized for the user-specified objectives and selects a subset that contributes to faster convergence and higher performance metrics.
- We further employ pretrained language encoders with the proposed Cross-attentive Embedding Augmentation (CEA) to enrich semantic information in image features and to provide explicit semantic guidance for data selection and finetuning.

2. Related work

2.1. Active learning

The learning algorithm in active learning is allowed to choose the data from which it learns [38]. There are two main selection criteria: uncertainty [3, 28, 47, 49] and diversity [1, 5, 37, 43]. Uncertainty of the model can aid selection of the most difficult unlabeled data. Early works estimate the uncertainty with various heuristics such as posterior probability [21, 27, 42, 48], entropy [17, 30] and classification margin [39]. Previous works [32, 37] also formulate active learning as an optimization problem. They typically operate in a discrete space that trivially matches the sample distribution of a dataset [11, 15]. However, discrete optimization is harder to solve than the one in continuous space. Also, most previous methods are designed for a from-scratch training without the pretraining stage. Bengar et al. [4] reveals drawbacks of such setting without unsuper-

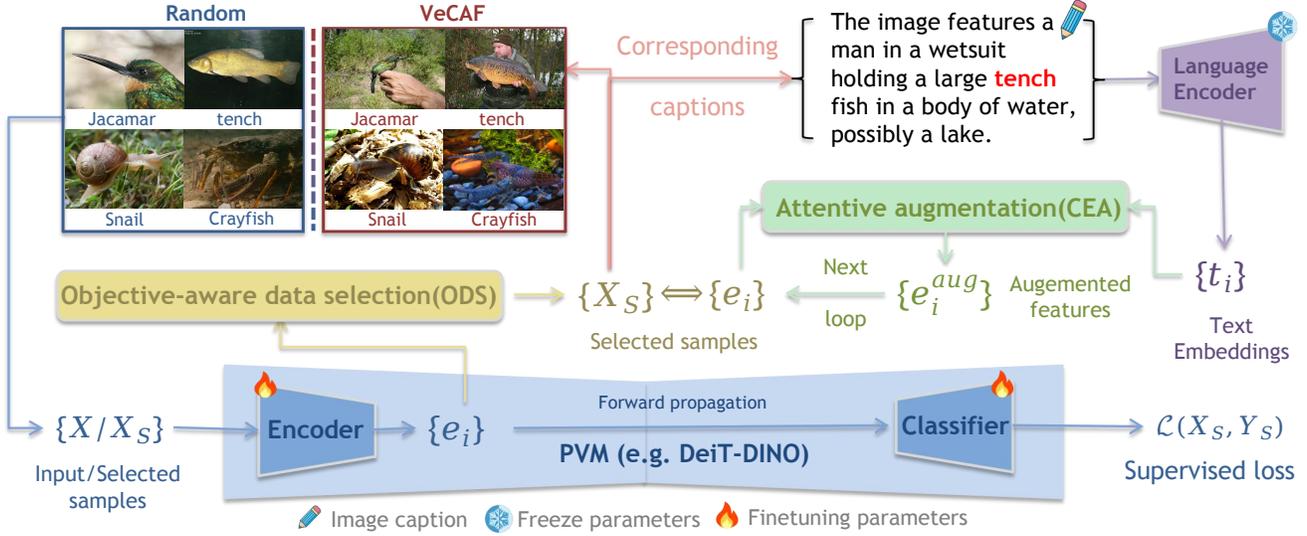


Figure 2. **The overall framework of VeCAF.** In each data selection loop, VeCAF performs an Objective-aware Data Selection (ODS) to select more informative images for finetuning. Cross-attentive Embedding Augmentation (CEA) is performed on the selected images to further enrich the semantic information captured by the image embeddings by incorporating language knowledge of the caption.

vised pretraining. Xie et al. [44] addresses both shortcomings with the proposed continuous-space ActiveFT method that applies selected samples to the finetuning of the pre-trained model in a single pass. VeCAF extends ActiveFT to the practical setting of data selection from a large labeled dataset. VeCAF framework enhances the training efficiency with a training objective-aware data selection and achieves optimal finetuning results with minimal training batches.

2.2. Exploiting language in vision model training

Recent developments in language models, especially vision-language models [22, 23] demonstrate their effectiveness in aligning the embedding space of vision and language to achieve cross-modal generalization. For example, BLIP-2 [22] addresses the modality gap using a lightweight Querying Transformer, while Shikra [7] handles spatial coordinate inputs and outputs in natural language and excels in referential dialogue and general vision-language tasks. Many prior work exploit the connection between the image and text modalities, where they explore the use of language in training better vision models. For example, Ma et al. [31] leverage pretrained language models to design a distribution alignment objective. This objective guides the vision model to learn linguistic representations specific to the task under a semi-supervised setting. Similarly, Fahes et al. [12] utilize CLIP to optimize affine transformations of source domain features. This optimization aligns these features with the target text embeddings while preserving their content and semantics. In our work, we use the language embeddings of image captions to perform text-space augmentations, achieving better sample selection quality in

active learning setting.

3. Method

This section introduces the details of the proposed VeCAF framework. As illustrated in Figure 2, we start by selecting a subset of training data with the PVM image feature space using the proposed ODS method, as introduced in Section 3.1. Then, the selected samples pass through pretrained language model to get semantically-rich text embeddings, which augment the image features with our proposed CEA technique, as formulated in Section 3.2. The augmented image features are used for PVM finetuning as candidates in future rounds of data selection. Finally, we show that VeCAF can overcome challenges of out-of-distribution data in Section 3.3.

3.1. Training objective-aware data selection

Consider a PVM $f(\cdot; w)$ with w weights and a user-defined finetuning objective \mathcal{L} . Given a labeled training set $\mathcal{D} : \{\mathcal{X}, \mathcal{Y}\}$, where \mathcal{X} is the set of image x and \mathcal{Y} is the set of corresponding label y . The goal of our active data selection is to find a subset $\mathcal{S} : \{\mathcal{X}_S, \mathcal{Y}_S\} \subset \mathcal{D}$, such that PVM finetuning with this subset for a fixed number of iterations leads to the largest reduction of the training objective. Formally, we formulate the optimization problem of our objective-aware data selection (ODS) process as

$$\mathcal{S}_{opt} = \arg \min_{\mathcal{S}} \mathbb{E}_{x, y \in \mathcal{D}} [\mathcal{L}(f(x; w - \beta \delta_S), y)], \quad (1)$$

where $\delta_S = \nabla_w \mathbb{E}_{x_s, y_s \in \mathcal{S}} [\mathcal{L}(f(x_s; w), y_s)]$ is the gradient accumulated by finetuning with \mathcal{S} and β is the learning rate.

In a practical setting, we cannot compute the gradient $\nabla_w \mathcal{L}(f(x; w), y)$ for each training example (x, y) before the data selection, as the gradient computation cost almost equals to the cost of full finetuning, which contradicts the purpose of active data selection. In this sense, we make an assumption that a data point with a larger loss contributes more to the convergence speed of the model.

Intuitively, this assumption leads a naive data selection policy of selecting the data points with Top- K training losses. However, previous active learning work [44] has discovered that the diversity of the selected data is also important in order to cover corner cases in the dataset and to avoid overfitting. Therefore, we design the ODS algorithm with the following principles: **1) data point with a larger loss $\mathcal{L}(f(x; w), y)$ shall be selected with a higher probability;** and **2) maintaining the diversity of data points selected in \mathcal{S} .** Analytically, we formulate our ODS objective as

$$\mathcal{S}_{opt} = \arg \min_{\mathcal{S}} D_{KL}(p_{\mathcal{L}}(\mathcal{D}) || p_{\mathcal{S}}(\mathcal{S})) - \lambda R(p_{\mathcal{S}}(\mathcal{S})), \quad (2)$$

where $D_{KL}(\cdot || \cdot)$ is the KL divergence, $R(\cdot)$ is a diversity metric, and λ is a tradeoff factor. $p_{\mathcal{L}}(\mathcal{D})$ is the distribution of the full training set that guides data selection. To follow our first principle, we assign the probability of each data point (x, y) in $p_{\mathcal{L}}(\mathcal{D})$ according to the finetuning objective \mathcal{L} scaled by a Z normalization factor as

$$p_{\mathcal{L}}(x, y) = \mathcal{L}(f(x; w), y) / Z. \quad (3)$$

The distribution of the selected data $p_{\mathcal{S}}(\mathcal{S})$ is determined by the data selection model. As a sanity check, the naive ‘‘Top- K training losses’’ serves as the optimal solution for Equation 2 without considering diversity when $\lambda = 0$.

To enable continuous optimization, we optimize Equation 2 using a parameterized data selection model θ_S . For the simplicity, we follow previous work [44] to model the data selection distribution $p_{\mathcal{S}}(\mathcal{S})$ in the lower-dimension image embedding space, where embedding e is produced by the PVM from the input x at a hidden layer. θ_S consists of K centroids in the image embedding space, each selecting the nearest data point. We define the probability of a data point x_i with the corresponding embedding e_i being selected as

$$p_{\mathcal{S}}(x_i) = \exp(\langle e_i, \theta_S^{c_i} \rangle) / \sum_{x_j \in \mathcal{D}} \exp(\langle e_j, \theta_S^{c_j} \rangle), \quad (4)$$

where $\langle \cdot, \cdot \rangle$ denotes the cosine distance, and $\theta_S^{c_i}$ is the closest centroid in θ_S to e_i . We derive the parameterized distribution distance $D(\theta_S) := D_{KL}(p_{\mathcal{L}} || p_{\mathcal{S}})$ for optimization objective Equation 2 using Equation 3 and 4 as

$$\begin{aligned} D(\theta_S) &= \sum_{(x_i, y_i) \in \mathcal{D}} p_{\mathcal{L}}(x_i, y_i) \log \frac{p_{\mathcal{L}}(x_i, y_i)}{p_{\mathcal{S}}(x_i)} \\ &= \mathbb{E}_{p_{\mathcal{L}}} [\log p_{\mathcal{L}}(x_i, y_i)] - \mathbb{E}_{p_{\mathcal{L}}} [\log p_{\mathcal{S}}(x_i)] \\ &= C - \alpha \sum_{(x_i, y_i) \in \mathcal{D}} \mathcal{L}(f(x_i; w), y_i) \langle e_i, \theta_S^{c_i} \rangle, \end{aligned} \quad (5)$$

where C and α are the constants omitted in the formulation.

For the diversity metric $R(\cdot)$, we follow the diversity regularization term proposed in [44] as

$$R(\theta_S) = - \sum_{\theta_S^i} \left[\log \sum_{\theta_S^j, j \neq i} \exp(\langle \theta_S^i, \theta_S^j \rangle) \right]. \quad (6)$$

By substituting Equation 5 and Equation 6 into Equation 2 and by removing constant terms, our final objective in terms of θ_S parameters can be written as

$$\begin{aligned} \theta_S^* &= \arg \min_{\theta_S} - \sum_{(x_i, y_i) \in \mathcal{D}} \mathcal{L}(f(x_i; w), y_i) \langle e_i, \theta_S^{c_i} \rangle \\ &\quad + \lambda \sum_{\theta_S^i} \left[\log \sum_{\theta_S^j, j \neq i} \exp(\langle \theta_S^i, \theta_S^j \rangle) \right]. \end{aligned} \quad (7)$$

The optimization on θ_S is conducted via gradient descent. To further resolve the dependency of the optimized data selection model to its initialization as observed in [44], we consider an independently-initialized data selection model ensemble, denoted as $\{\theta_e\}_{e=1}^E$, in the optimization. We empirically set $E = 5$ to balance the performance-cost tradeoff. After optimizing each data selection model θ_e independently using Equation 7, we can then remove the bias in the initialization by utilizing the mean μ and covariance Σ calculated on θ_e^* . Specifically, given an optimized data selection model θ_1^* , we achieve the final unbiased data selection model as $\theta_S^* = \theta_1^* - \Sigma^{-1}(\theta_1^* - \mu)$.

The optimization in Equation 7 is performed before each finetuning ‘‘loop’’ with the current model weights w . The training data with the closest embedding to each θ_S centroid is selected to form a ‘‘finetuning set’’ with K elements. Then, this set is used to finetune the PVM until the next loop starts. We update ODS weights w after a predetermined number of training batches.

3.2. Cross-attentive embedding augmentation

We empirically find that samples selected by the ODS process tend to consist of multiple objects in both foreground and background, which is a result of coarse image embeddings produced by the PVM. To further improve quality of image embeddings for both data selection and finetuning, we propose Cross-attentive Embedding Augmentation (CEA). Given a training image, CEA leverages

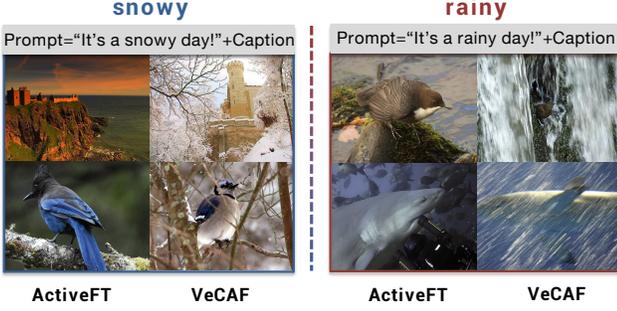


Figure 3. The selected samples of ActiveFT [44] and VeCAF. With the caption augmented: “It is a {snowy/rainy} day!”, VeCAF can select images that correspond to the target domain.

a text encoder to transform the image caption to the corresponding text embedding. Then, the text embedding is used to augment the image embedding with attention-based method.

Given selected sample $\{x_{S_i}\}_{i=1}^K$ in the previous ODS run, the corresponding image caption can be denoted as $\{c_{S_i}\}_{i=1}^K$. Then we feed the caption into a frozen text encoder, e.g. BERT [9], to convert the captions into the text embeddings $\{t_i\}_{i=1}^K$. We use text embeddings with the same dimensions as image embeddings for easier fusion.

Inspired by [16], CEA is conducted as mapping the image embedding e_i towards the corresponding text embedding t_i . To decide the magnitude of the augmentation, we compute a sample-wise attention score $\{\alpha_i\}_{i=1}^K$ with the cosine distance between e_i and t_i as

$$\alpha_i = \text{Softmax}\left(\frac{e_j \cdot t_j}{\|e_j\|_2 \|t_j\|_2}\right) = \frac{\exp\left(\frac{e_i \cdot t_i}{\|e_i\|_2 \|t_i\|_2}\right)}{\sum_{j=1}^K \exp\left(\frac{e_j \cdot t_j}{\|e_j\|_2 \|t_j\|_2}\right)}. \quad (8)$$

The attention score α_i helps to derive the augmented embedding e_i^{aug} using the image embedding e_i and the corresponding text embedding t_i as

$$e_i^{aug} = e_i - \eta \alpha_i (e_i - t_i), \quad (9)$$

where η is the fixed step size. The proposed CEA method enriches the semantic information and improves performance after finetuning as shown in experiments. We further present the pseudo code of VeCAF in Algorithm 1, which specifies the complete procedure of the proposed Vision-language Collaborative Active Finetuning (VeCAF).

3.3. Improving out-of-distribution scenarios

In addition to achieving more efficient finetuning on out-of-distribution (OOD) tasks, one of the additional benefits of introducing the text modality is the ability to artificially modify the corresponding image captions as a semantic augmentation. By leveraging the strong language capability

Algorithm 1: Vision-language Collaborative Active Finetuning (VeCAF)

input: Objective-aware data selection $ODS(\cdot;)$, labeled data pool $(\mathcal{X}, \mathcal{Y})$, image caption pool Cap , PVM $f(\cdot; w)$, pretrained language encoder $LM(\cdot)$, data selection loop number L , batch number B for each loop

output: Finetuned vision model $f(\cdot; w_{FT})$

- 1 **for** $loop \in [L]$ **do**
- 2 Obtained the selected sample pool
 $\mathcal{S}_{opt} = ODS((\mathcal{X}, \mathcal{Y}); f(\cdot; w))$;
- 3 Get the corresponding image caption Cap_i for
 $s_i \in \mathcal{S}_{opt}$;
- 4 Transfer the image caption to text embedding as
 $t_i = LM(Cap_i)$;
 /* Cross-attentive embedding augmentation (CEA) */
- 5 CEA attention score computation
 $\alpha_i = \text{Softmax}\left(\frac{e_j \cdot t_j}{\|e_j\|_2 \|t_j\|_2}\right)$;
- 6 Image embedding e_i augmentation
 $e_i^{aug} = e_i - \eta \cdot \alpha_i (e_i - t_i)$;
 /* PVM finetuning with \mathcal{S}_{opt} */
- 7 **for** $batch \in [B]$ **do**
- 8 Sample next batch from \mathcal{S}_{opt} ;
- 9 Calculate the loss with the classifier;
- 10 Optimize $f(\cdot; w_{FT})$ via gradient descent;
- 11 **end**
- 12 $f(\cdot; w) \leftarrow f(\cdot; w_{FT})$
- 13 **end**
- 14 Return the finetuned vision model $f(\cdot; w_{FT})$

provided by the text encoder, we can implicitly alter our requirements for the selected data. This capability enables us to facilitate domain transfer to out-of-distribution (OOD) scenarios with only ID data.

For example, leveraging the CEA technique, we can add descriptive phrases like “It’s a {target_domain} day!” to the image captions. This modification shifts the PVM image embedding towards the “snow” distribution while highlighting the ID sample images that have snowflakes or similar patterns during the data selection process, as demonstrated in Figure 3. By incorporating such textual cues, we can guide the selection process toward images that possess specific characteristics or attributes, so the finetuned model can better generalize in out-of-distribution scenarios.

4. Experiments

In this section, we explain our experiment protocols and conduct experiments on multiple image classification tasks. We demonstrate the superiority of the proposed VeCAF in

Table 1. **Classification accuracy with the fixed training cost.** All methods for each dataset are trained using the fixed number of batches. The percentage value on top reports the ratio of data selected during each loop. Some results marked as N/A (“-”) as explained in Appendix B. Top-1 accuracy with a standard error of 3 repetitions is reported, %.

Method	Loop	CIFAR-10			Caltech101			ImageNet-1K		
		1%	2%	5%	2%	5%	10%	1%	2%	4%
LearnLoss[46]	single-run	85.93 \pm 0.05	91.22 \pm 0.08	93.89 \pm 0.07	46.47 \pm 0.16	43.74 \pm 0.13	65.59 \pm 0.05	49.37 \pm 0.09	57.86 \pm 0.07	63.46 \pm 0.11
	multi-run	87.53 \pm 0.14	92.53 \pm 0.17	94.43 \pm 0.22	47.36 \pm 0.13	44.27 \pm 0.19	66.27 \pm 0.14	49.89 \pm 0.16	58.22 \pm 0.17	64.18 \pm 0.13
TA-VAAL[18]	single-run	85.46 \pm 0.10	92.65 \pm 0.06	94.85 \pm 0.10	59.26 \pm 0.05	58.11 \pm 0.08	66.94 \pm 0.08	-	-	63.86 \pm 0.13
	multi-run	87.74 \pm 0.17	93.77 \pm 0.19	96.01 \pm 0.12	60.57 \pm 0.18	59.25 \pm 0.15	67.32 \pm 0.21	-	-	64.32 \pm 0.16
ALFA-Mix[33]	single-run	86.69 \pm 0.07	92.87 \pm 0.06	95.14 \pm 0.08	59.73 \pm 0.05	58.74 \pm 0.09	67.36 \pm 0.08	-	-	64.03 \pm 0.07
	multi-run	88.14 \pm 0.13	93.26 \pm 0.13	95.75 \pm 0.11	60.74 \pm 0.16	60.47 \pm 0.14	68.25 \pm 0.15	-	-	64.69 \pm 0.19
ActiveFT [44]	single-run	90.91 \pm 0.12	93.80 \pm 0.09	95.39 \pm 0.08	62.86 \pm 0.05	60.55 \pm 0.09	69.34 \pm 0.06	53.96 \pm 0.07	60.33 \pm 0.09	64.72 \pm 0.10
	multi-run	92.79 \pm 0.11	94.17 \pm 0.16	95.92 \pm 0.14	64.43 \pm 0.12	61.97 \pm 0.17	71.42 \pm 0.14	55.67 \pm 0.13	61.86 \pm 0.21	65.18 \pm 0.17
Full Data FT	single-run	93.64 \pm 0.02			62.08 \pm 0.02			57.53 \pm 0.01		
VeCAF(ours)	multi-run	93.57 \pm 0.02	95.27 \pm 0.04	96.24 \pm 0.02	66.33 \pm 0.04	65.15 \pm 0.03	72.21 \pm 0.03	58.31 \pm 0.04	63.76 \pm 0.03	66.57 \pm 0.02

both in-distribution and out-of-distribution scenarios and improve the finetuning efficiency by up to $3.3\times$ and 2.7% on ImageNet-1K. We first introduce the experiment setup in Sec. 4.1. Main results in Sec. 4.2 show the effectiveness and efficiency of the proposed VeCAF compared with other sample selection methods. In Sec. 4.3, we present further analysis from different perspectives.

4.1. Experiment setup

Datasets For model training, we conduct experiments using three image classification datasets: CIFAR-10 [19], class-unbalanced Caltech101 [13], and ImageNet-1K [8]. Details of the dataset can be found in Appendix A. For out-of-distribution evaluation, we evaluate on ImageNet-C [14] which consists of synthetically generated corruptions applied to the ImageNet validation set.

Implementation details In our main experiments, we use the DeiT-B model [40] pretrained with DINO [6] on ImageNet-1K [8] as the PVM for finetuning. Additionally, we present experiments for different PVM architectures and model sizes in Section 4.3 to demonstrate VeCAF generalizability. For all experiments, we resize input images to 224×224 to ensure consistency during both the data selection and the finetuning. In the ODS process, we optimize the data selection model parameters θ using the Adam optimizer with a learning rate of 0.001 until convergence. We follow [40] to finetune the PVM with the selected data.

We adopt the standard protocols outlined in [40] to finetune the DeiT-Base model. For the CIFAR-10 [19] and the Caltech101 [13], we perform supervised finetuning of the pretrained models using the SGD optimizer. We set the learning rate (lr) to $5e-4$, weight decay to $1e-4$, and momen-

tum to 0.9. The total batch number used for training CIFAR-10 and Caltech101 is 750 and 1500, respectively. For the ImageNet-1K [8], the SGD optimizer with the same hyperparameters as CIFAR-10 and Caltech101 is employed. The total number of used batches for training is 125000. These experiments are conducted on two Tesla-A100 GPUs. Each GPU processes a batch size of 256, and we apply cosine learning rate decay on selected subsets of the training data.

Evaluation protocol We primarily focus on an efficient training setting, where all the algorithms are only allowed to train with the same batch size and the same number of batches. Convergence results with unlimited batches are in Appendix C.3. As a default setting, we perform multi-run data selection before each loop with the fixed training batches for all baselines including VeCAF. Three loops are used in the main experiments, where each loop is defined as 1/3-rd of the total number of training batches.

Baselines We compare VeCAF with three active learning baselines LearnLoss [46], TA-VAAL [18], ALFA-Mix [33], ActiveFT [44] and conventional full data finetuning.

- LearnLoss [46] predicts target losses for unlabeled inputs to identify potential incorrect predictions.
- TA-VAAL [18] considers the data distribution of labeled and unlabeled pools and enhances the learning process by incorporating a ranking loss prediction.
- ALFA-Mix [33] employs interpolations between labeled and unlabeled instances to uncover unrecognized features, which derives an efficient implementation us-

Table 2. **OOD generalization ability.** Models are trained with 1% data subset per loop using the uncorrupted ImageNet. Evaluation results are for the distorted ImageNet-C validation set.

Source	Target (eval.)	Method	Top-1 Acc.
Prompt: It’s a snowy day.			
		Source-only	38.89 \pm 0.20
ImageNet	ImageNet-C	CLIPStyler	40.71 \pm 0.41
		Snowy	ActiveFT
		VeCAF	42.33 \pm 0.03
	Prompt: It’s a foggy day.		
		Source-only	45.55 \pm 0.27
ImageNet-C	Foggy	CLIPStyler	46.25 \pm 0.36
		ActiveFT	42.45 \pm 0.08
		VeCAF	47.71 \pm 0.03

ing a closed-form solution to identify the optimal interpolation that induces changes in predictions.

- ActiveFT [44] focuses on selecting the data subset that exhibits similar distribution to the unlabeled pool, while maintaining the diversity within the selected subset by optimizing a parametric model in the continuous space.
- Full Data FT takes all the data from the training set to finetune the pretrained vision models.

4.2. Main results

In-distribution results *Overall performance:* We present the image classification results in Table 8. The results demonstrate the superior effectiveness of the proposed VeCAF method when compared to other active learning approaches. For a fair comparison, we have also extend the four active learning baselines to a multi-run framework. This means that we apply these methods to select the data finetuning subset at the beginning of each loop using different random seeds. We can observe from the results that traditional active learning methods often encounter challenges within the pretraining-finetuning paradigm, which aligns with our findings. In contrast, VeCAF consistently outperforms other methods across all three datasets, irrespective of the employed sampling ratios. Remarkably, even with low sampling ratios, our method excels in selecting highly representative samples. When compared to full data finetuning with limited training batches, VeCAF achieves higher accuracy even with only 1% of the data being used for finetuning in each epoch. This practical advantage is significant as it allows for supervised finetuning with a smaller number of samples compared to the overall pool size with less steps, thereby reducing training costs. Additionally, it is noteworthy that VeCAF exhibits a more pronounced performance enhancement on more complex datasets, with an increase

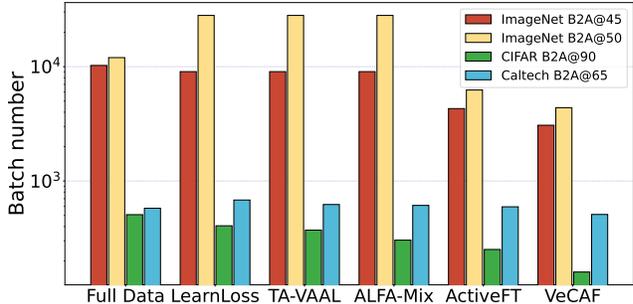


Figure 4. **Comparison of training efficiency.** VeCAF requires significantly fewer training batches to reach the target accuracy (B2A) compared with other baselines and full-data finetuning. Note that the y-axis has an exponential scale.

Table 3. **1-vs.-all accuracy for certain categories on CIFAR-10.** DINO-S [18] is used as PVM with 2% of data in each loop.

Methods	Airplane	Bird	Cat	Deer
ActiveFT[44]	89.50 \pm 0.01	82.70 \pm 0.03	83.00 \pm 0.02	87.50 \pm 0.04
Top-K Loss	82.18 \pm 0.05	81.88 \pm 0.01	75.72 \pm 0.04	84.64 \pm 0.03
VeCAF(ours)	91.24 \pm 0.02	88.80 \pm 0.04	86.41 \pm 0.02	89.12 \pm 0.03

of over 2.7% accuracy on ImageNet-1K compared to 1% on CIFAR-10.

Efficiency enhancement: To provide a deeper understanding of the effectiveness of our proposed method, we have included figures to illustrate the required batch numbers to achieve the target accuracy (B2A) for different approaches on the three datasets. Specifically, we report training batches needed for each method with **B2A@45,50** (i.e., 50% Top-1 accuracy) for *ImageNet*, **B2A@65** for *Caltech101*, and **B2A@90** for *CIFAR-10*, respectively. Figure 4 displays such batch numbers, highlighting the efficiency of VeCAF in comparison to other methods. On ImageNet, VeCAF achieves 3.3 \times acceleration over full data finetuning (3075 v.s. 10250 batches) and outperforms other baselines more significantly as the target accuracy goes higher. Additionally, we plot the training loss as a function of the batch number for different approaches during the finetuning process in Figure 5. This visualization further highlights the convergence speed and performance of VeCAF compared to alternative methods. These figures provide a comprehensive view of the performance and efficiency of VeCAF, emphasizing its effectiveness in terms of training speed and achieving the desired accuracy levels.

Finegrained training objective awareness: The proposed ODS method offers additional flexibility to accommodate finegrained training objective. For verification, we evaluate VeCAF under 1-vs.-all finetuning objective on multiple CIFAR-10 classes. Specifically, given a target class, we set the loss as a binary classification, with target class

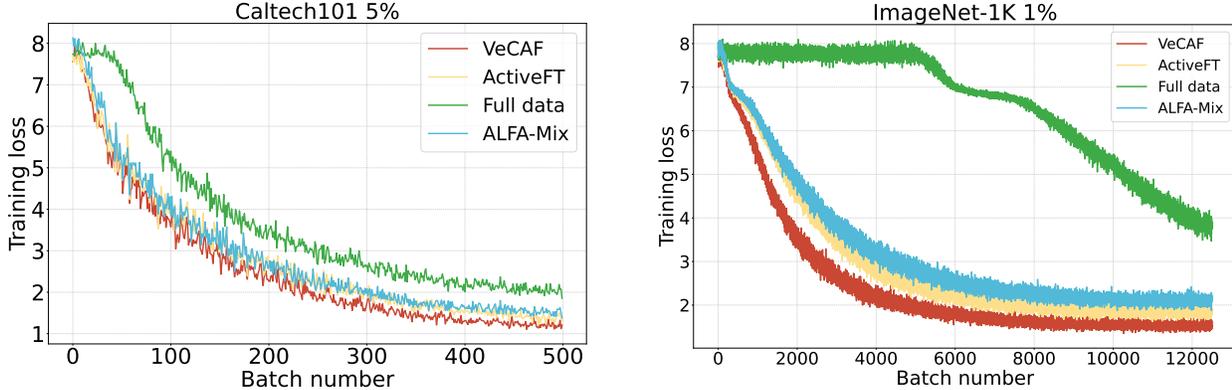


Figure 5. Training loss curve of VeCAF and other baselines including ActiveFT, ALFA-Mix, and Full data FT on Caltech-101 (left) and ImageNet-1K (right) with 5% and 1% data, respectively.

being positive and all others being negative. The modified loss is used in ODS optimization as in Equation 7. A naive objective-aware baseline of selecting samples with the largest loss is also included in the comparison. Table 3 shows that VeCAF accommodates the finegrained objective to better improve the performance of target categories.

Out-of-distribution results We also conduct experiments using ImageNet-C test set to assess its superiority under out-of-distribution (OOD) scenarios, as presented in Table 2. Specifically, we consider the domain adaptation scenarios of clear→snowy and clear→foggy.

To provide a comprehensive evaluation, we compare VeCAF with the state-of-the-art baseline, CLIPstyler [20]. Both VeCAF and CLIPstyler leverage the pretrained CLIP model and offer generalization capabilities for source images. However, it is essential to note that VeCAF uses active learning concept to select source-domain data points for finetuning, while CLIPstyler specializes in style transfer. Furthermore, we compare the results obtained in the source-only setting, where models are finetuned exclusively on the source data, with full training set and ActiveFT [44] baselines. 1% of the training data is used in each finetuning loop for ActiveFT and VeCAF.

In these experiments, we finetune the classifier using uncorrupted (source) ImageNet data points selected by ActiveFT and VeCAF, while CLIPstyler finetunes on stylized images. We can see from Table 2 that VeCAF consistently outperforms the three baselines by up to 4.2%. CLIPstyler achieves better performance than other baselines, but multiple artifacts in the stylized images limit its performance. ActiveFT underperforms in the OOD setting with only 36.36% and 42.45% accuracy, respectively. This is a result of overfitting as ActiveFT selects data only based on the source domain distribution. On the other hand, VeCAF is superior in generalization ability, by leveraging a domain-specific text embedding augmentation on selection and finetuning enabled by the proposed CEA.

Table 4. Top-1 accuracy of VeCAF with different types and sizes of PVM backbones, and choices of pretrained language encoders. ImageNet-1K results are reported with 1% data.

PVM \ LM	CLIP	BERT-L	mT5-L	GPT2-L
	DeiT-S	53.51±0.02	53.64±0.03	53.82±0.03
DeiT-B	58.31±0.04	58.63±0.04	58.71±0.03	58.87±0.01
SwinT-S	53.66±0.03	53.71±0.03	53.89±0.02	54.01±0.02
SwinT-B	56.76±0.01	56.87±0.03	56.98±0.03	57.13±0.02
XciT-M	58.48±0.04	58.70±0.03	58.77±0.03	58.95±0.03
XciT-L	61.13±0.01	61.37±0.01	61.56±0.03	61.78±0.03

4.3. Results analysis

Generality of VeCAF The proposed VeCAF framework can be used to finetune various pretrained vision models (PVMs) with the help of different language encoders (LM). We apply it to DeiT-B [40] pretrained using the DINO framework [18] and several other PVMs such as DeiT-S, Swin-Transformer-S/B [29], and XciT-M/L [2]. Furthermore, we use other LMs including BERT-L [9], mT5-L [45], and GPT2-L [35] for augmentation by text embeddings. Table 4 reports results on ImageNet-1K using different pairs of PVMs and LMs. VeCAF demonstrates its versatility and capability to adapt to different PVM model architectures, and the flexibility to be used with different text encoder models per user preferences. These results prove VeCAF to be a general active data selection method, that can incorporate the text embedding information from various configurations of LMs to improve the efficiency.

Embedding visualization In Table 6, we present the image embedding visualization of sample in the CIFAR-10 training set using UMAP dimension reduction. With each

Table 5. **Ablation study for the proposed techniques in VeCAF.** Data selection ratio is set to 1% for CIFAR-10, 5% for Caltech101, and 1% for ImageNet-1K in each loop.

ODS	CEA	CIFAR-10	Caltech-101	ImageNet-1K
-	-	92.79 \pm 0.12	60.55 \pm 0.10	55.67 \pm 0.13
✓	-	93.27 \pm 0.03	64.11 \pm 0.04	57.73 \pm 0.04
-	✓	93.15 \pm 0.05	63.04 \pm 0.06	56.13 \pm 0.05
✓	✓	93.57\pm0.02	65.15\pm0.03	58.31\pm0.04

Table 6. FLOPs and time for VeCAF, ActiveFT, and Full-FT

Method	Training batches	Total FLOPs (G, all batches)				Wallclock
		CEA	ODS/DS	FT	ALL	
VeCAF	3000	9.53 \times 10 ³	1.5 \times 10 ²	2.11 \times 10 ⁵	2.21 \times 10 ⁵	2.4h
Full-FT	10000	-	-	7.01 \times 10 ⁵	7.01 \times 10 ⁵	3.1h

method selecting 1% of data from the training set, black dots represent the samples selected by both ActiveFT and VeCAF, red stars denote samples only selected by VeCAF, and blue stars denote samples only selected by ActiveFT. While maintaining the diversity of data selection, samples chosen by VeCAF appear to be closer to the boundaries compared to those selected by ActiveFT. This confirms that the proposed ODS helps to select important samples around the decision boundaries. This is a result of our selection strategy that incorporates training objective and, therefore, helps the PVM to learn the subtle differences between categories more efficiently in the finetuning.

Efficiency analysis We analyze the overhead of VLM in each data selection loop in Table 6. We estimate the backward computations triple the forward pass following the PyTorch report [24]. The CLIP-ViT-L model we use requires about 11 \times the FLOPs of DeiT, but only needs to be inferred once in the loop. This leads to the FLOPs overhead of CEA to be merely 4.5% of the finetuning cost. The resulting FLOPs reduction ratio is therefore similar to the batch number reduction ratio. Note that this analysis is consistent across datasets as all data are resized to 224 \times 224.

Ablation study We first verify the importance of proposed techniques, the Objective-aware Data Selection (ODS) and Cross-attentive Embedding Augmentation (CEA) in Table 5. Specifically, ODS can be disabled by removing the \mathcal{L} term in Equation 7. ODS significantly improves model performance, especially with limited data, enhancing classification accuracy and expediting the finetuning process. CEA further improves classification accuracy by integrating rich semantic information from text embeddings, enhancing model generalization and capturing under-

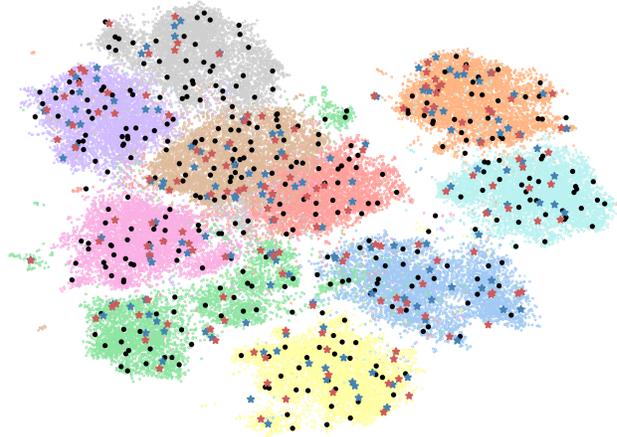


Figure 6. **UMAP visualization of training image embeddings.** Each background color represents one class. For selected samples, ★ suggests being selected by VeCAF only, ★ suggests by ActiveFT [44] only, and ● suggests by both.

Table 7. Ablation study of the number of data selection loops on the CIFAR-10 dataset with 5% data selection.

# Loops	2	3	4	5
CIFAR	95.89 \pm 0.02	96.01 \pm 0.02	96.14 \pm 0.03	96.24\pm0.03

lying semantics. We further study the impact of the number of loops using CIFAR-10 dataset in Table 7. The total number of training batches is the same in all experiments with equal division by the loop count. It is clear that the model performance is enhanced with the increase in the number of loops. However, performing additional loops of data selection leads to overhead in the total training time. We set number of loop to three throughout our experiments to balance the trade-off. Detailed analysis on the training time is provided in Appendix C.2.

5. Conclusions and future work

In this paper, we improved the efficiency of finetuning a PVM towards a user-specified performance target with a novel active data selection framework, VeCAF. VeCAF finds a subset of training data that leads to faster convergence with an objective-aware data selection model, and additionally utilizes the text-domain knowledge of pretrained VLM to augment image embeddings. Through extensive experiments, we demonstrated the superior performance of the proposed approach when compared to baseline active data selection methods and the finetuning with all data. In future work, we aim to further unleash the potential of text-domain augmentation by improving certain finegrained performance metrics in vision domain tasks, and by extending active data selection to active data generation for additional performance and efficiency improvements.

References

- [1] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [2] Alaaeldin Ali, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. XcIT: Cross-covariance image transformers. *Advances in Neural Information Processing Systems*, 2021.
- [3] William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9368–9377, 2018.
- [4] Javad Zolfaghari Bengar, Joost van de Weijer, Bartłomiej Twardowski, and Bogdan Raducanu. Reducing label effort: Self-supervised meets active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [5] Zefan Cai, Baobao Chang, and Wenjuan Han. Human-in-the-loop through chain-of-thought. *arXiv preprint arXiv:2306.07932*, 2023.
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [7] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal LLM’s referential dialogue magic. *arXiv:2306.15195*, 2023.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [11] Brochu Eric, Nando Freitas, and Abhijeet Ghosh. Active preference learning with discrete choice data. *Advances in Neural Information Processing Systems*, 2007.
- [12] Mohammad Fahes, Tuan-Hung Vu, Andrei Bursuc, Patrick Pérez, and Raoul de Charette. PODA: Prompt-driven zero-shot domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [13] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2004.
- [14] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [15] Aditi Jha, Zoe C Ashwood, and Jonathan W Pillow. Bayesian active learning for discrete latent variable models. *arXiv:2202.13426*, 2022.
- [16] Shaoxiong Ji, Shirui Pan, Guodong Long, Xue Li, Jing Jiang, and Zi Huang. Learning private neural language modeling with attentive aggregation. In *International joint conference on neural networks (IJCNN)*, 2019.
- [17] Ajay J Joshi, Fatih Porikli, and Nikolaos Papanikolopoulos. Multi-class active learning for image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [18] Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [19] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *MS thesis, University of Toronto*, 2009.
- [20] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [21] David D. Lewis and Jason Catlett. Heterogenous uncertainty sampling for supervised learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 1994.
- [22] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv:2301.12597*, 2023.
- [23] Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, et al. M³IT: A large-scale dataset towards multi-modal multilingual instruction tuning. *arXiv:2306.04387*, 2023.
- [24] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, et al. Pytorch distributed: Experiences on accelerating data parallel training. *arXiv preprint arXiv:2006.15704*, 2020.
- [25] Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Open-vocabulary object segmentation with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7667–7676, 2023.
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollar, and Larry Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014.
- [27] Jiaming Liu, Rongyu Zhang, Xiaowei Chi, Xiaoqi Li, Ming Lu, Yandong Guo, and Shanghang Zhang. Multi-latent space alignments for unsupervised domain adaptation in multi-view 3d object detection. *arXiv preprint arXiv:2211.17126*, 2022.
- [28] Zhuoming Liu, Hao Ding, Huaping Zhong, Weijia Li, Jifeng Dai, and Conghui He. Influence selection for active learning. In *Proceedings of the IEEE/CVF Conference on Computer*

- Vision and Pattern Recognition (CVPR)*, 2021.
- [29] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [30] Wenjie Luo, Alex Schwing, and Raquel Urtasun. Latent structured active learning. *Advances in Neural Information Processing Systems*, 2013.
- [31] Wenxuan Ma, Shuang Li, JinMing Zhang, Chi Harold Liu, Jingxuan Kang, Yulin Wang, and Gao Huang. Borrowing knowledge from pre-trained language model: A new data-efficient visual learning paradigm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [32] Rafid Mahmood, Sanja Fidler, and Marc T Law. Low-budget active learning via wasserstein distance: An integer programming approach. In *International Conference on Learning Representations (ICLR)*, 2022.
- [33] Amin Parvaneh, Ehsan Abbasnejad, Damien Teney, Gholamreza Reza Haffari, Anton Van Den Hengel, and Javen Qin-feng Shi. Active learning by feature mixing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [34] 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 *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [35] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [36] 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.
- [37] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [38] Burr Settles. Active learning literature survey. Technical report, University of Wisconsin–Madison, 2010.
- [39] Simon Tong and Daphne Koller. Support vector machine active learning with applications to text classification. *JMLR*, 2001.
- [40] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [41] 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:2302.13971*, 2023.
- [42] Keze Wang, Dongyu Zhang, Ya Li, Ruimao Zhang, and Liang Lin. Cost-effective active learning for deep image classification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2016.
- [43] Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*, 2023.
- [44] Yichen Xie, Han Lu, Junchi Yan, Xiaokang Yang, Masayoshi Tomizuka, and Wei Zhan. Active finetuning: Exploiting annotation budget in the pretraining-finetuning paradigm. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [45] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. In *North American Chapter of the Association for Computational Linguistics*, 2020.
- [46] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [47] Tianning Yuan, Fang Wan, Mengying Fu, Jianzhuang Liu, Songcen Xu, Xiangyang Ji, and Qixiang Ye. Multiple instance active learning for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [48] Rongyu Zhang, Xiaowei Chi, Guiliang Liu, Wenyi Zhang, Yuan Du, and Fangxin Wang. Unimodal training-multimodal prediction: Cross-modal federated learning with hierarchical aggregation. *arXiv preprint arXiv:2303.15486*, 2023.
- [49] Rongyu Zhang, Yulin Luo, Jiaming Liu, Huanrui Yang, Zhen Dong, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, Yuan Du, et al. Efficient deweather mixture-of-experts with uncertainty-aware feature-wise linear modulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 16812–16820, 2024.

Appendix

In the supplementary material, we provide additional information for the main paper. We start by providing the datasets information in Appendix A and additional details on our reproduction of previous active learning baselines in Appendix B, and explains the missing results indicated by “-” in the experiment section of the paper. Additional insights into loss convergence on the CIFAR dataset are meticulously documented in Appendix C.1. We delve into the efficiency of our approach in Appendix C.2, providing a compelling narrative on the method’s expediency. Lastly, an extensive evaluation of the model’s accuracy, achieved with an unrestricted count of training batches, is detailed in Appendix C.3, solidifying the robustness of our experimental findings.

Besides this document, we also include the source code of VeCAF in the supplementary material. Please check the README file for details.

A. Datasets

CIFAR-10 [19] consists of 60,000 images with a resolution of 32×32 pixels, divided into 10 categories. The training set contains 50,000 images while the test set contains 10,000 images. The Caltech101 [13] dataset consists of images from 101 object categories with 40 to 800 images per category. Most classes contain around 50 images, and the image resolution is approximately 300×200 pixels. ImageNet-1K [8] is a larger dataset with 1,331,167 images belonging to 1,000 classes. The training set consists of 1,281,167 images while the validation set contains 50,000 images. All the training sets from these datasets are considered candidate pools for selection. We also leverage the ImageNet-C [14] as an OOD test set to evaluate our VeCAF under out-of-distribution scenarios. ImageNet-C is an openly available dataset that consists of algorithmically generated corruptions, including blur and noise, applied to the ImageNet test set. It serves as a valuable resource for evaluating the robustness and generalization capabilities of computer vision models.

B. Active learning reproducing details

In our study, we incorporate three well-established active learning algorithms, namely LearnLoss [46], TA-VAAL [18], and ALFA-Mix [33], within the pretraining-finetuning paradigm for image classification tasks. Specifically, we evaluate these algorithms on three widely used datasets, namely CIFAR-10, Caltech101, and ImageNet. To ensure a systematic and consistent evaluation, all three algorithms employ a batch-selection strategy for sample acquisition during the active learning process.

In Table 1 of our paper, the presence of “-” is attributed to the nature of traditional active learning methods, which re-

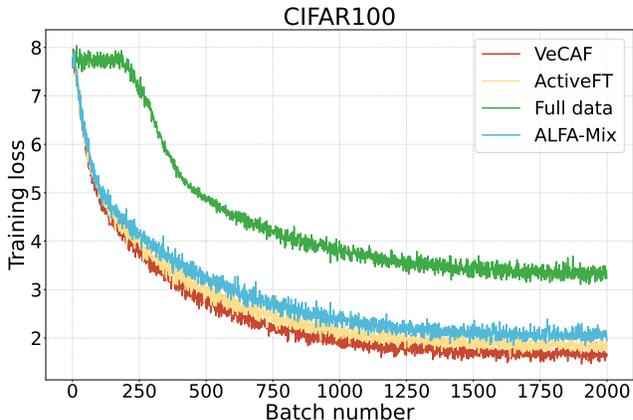


Figure 7. Training curve of VeCAF and baselines including ActiveFT, ALFA-Mix, and Full data FT on 5% CIFAR-100.

quire a small initial set randomly sampled at the beginning of the process. It is important to note that the performance of these active learning algorithms on this initial set is comparable to random sampling. Therefore, to avoid redundancy, we have omitted reporting duplicate results for these random initial sets. For instance, in the case of CIFAR-10, we exclude reporting results for the sampling ratio of 0.5%, and for Caltech101, results for the sampling ratio of 1% are not reported. Moreover, for the ImageNet-1K dataset, the reporting results of sampling ratios of 1% and 2% are omitted as the size of the initial set is 2.5% according to the reported setting of these previous papers. By excluding duplicate results for random initial sets and smaller sampling ratios, we aim to present clear and concise information in Table 1, focusing on the most relevant and informative performance metrics for the active learning methods applied in our study.

C. Additional experiment results

C.1. Loss convergence analysis

Given that the loss convergence on CIFAR-10 data is sufficiently rapid for each baseline and thus may not effectively highlight the advantages of our proposed VeCAF, we instead present the loss convergence of CIFAR-100, which possesses analogous domain characteristics to CIFAR-10, as illustrated in Figure 7. Mirroring the depiction in Figure 4 of the primary text, Figure 7 exhibits a comparable trend in loss convergence, where our proposed VeCAF not only converges more swiftly but also to a lower loss value in comparison to the baselines. This demonstration underscores the enhanced convergence efficiency of VeCAF, both in terms of speed and performance.

Table 8. Convergence accuracy with unlimited number of batches.

Method	Loop	CIFAR-10			Caltech101			ImageNet-1K		
		1%	2%	5%	2%	5%	10%	1%	2%	4%
Full Data FT	single-run	99.31 \pm 0.01			88.24 \pm 0.02			82.76 \pm 0.01		
LearnLoss[46]	single-run	90.07 \pm 0.02	93.67 \pm 0.03	95.99 \pm 0.02	62.88 \pm 0.02	73.09 \pm 0.03	83.04 \pm 0.04	52.97 \pm 0.03	60.14 \pm 0.03	61.93 \pm 0.03
	multi-run	90.25 \pm 0.03	94.21 \pm 0.03	96.32 \pm 0.04	63.74 \pm 0.02	73.25 \pm 0.03	83.31 \pm 0.03	53.66 \pm 0.05	60.49 \pm 0.03	62.32 \pm 0.04
ActiveFT [44]	single-run	92.31 \pm 0.02	95.46 \pm 0.02	98.18 \pm 0.04	73.69 \pm 0.02	81.33 \pm 0.03	86.78 \pm 0.02	56.87 \pm 0.04	63.19 \pm 0.03	66.01 \pm 0.03
	multi-run	92.95 \pm 0.01	95.87 \pm 0.03	98.54 \pm 0.02	74.22 \pm 0.02	81.88 \pm 0.03	87.04 \pm 0.02	57.11 \pm 0.03	63.46 \pm 0.03	66.21 \pm 0.02
VeCAF(ours)	multi-run	93.87\pm0.02	96.47\pm0.01	98.97\pm0.01	75.36\pm0.01	83.62\pm0.02	87.72\pm0.01	59.41\pm0.02	65.64\pm0.01	68.52\pm0.02

Table 9. Running time to select various percentages of samples from the Caltech101 training set for each data selection loop.

Sel. ratio	ALFA-Mix	LearnLoss	ActiveFT	VeCAF
2%	6m45s	1m42s	12.02s	16.38s
10%	52m31s	23m17s	13.36s	18.87s

C.2. Time complexity of data selection

Efficiency is a crucial aspect of the VeCAF, and it is desirable for it to operate in a time-efficient manner, so as to reduce the overhead of data selection in each training loop. In our study, we evaluate the time required to select various proportions of training samples from the Caltech101 dataset with selection ratio 2% and 10%, as shown in 9. Here we consider the image caption of each training sample readily available as they can be generated offline for only once, while the time for performing ODS, text embedding generation, and CEA are included in the reported VeCAF time. Traditional active learning algorithms like LearnLoss [46] and ALFA-Mix [33] requires multiple trial model updates to gradually adjust the selected data, where these trial updates constitute the majority of the time in the data selection process, making them significantly inefficient. In contrast, ActiveFT and our proposed VeCAF method perform sample selection in a single pass at the beginning of each data selection loop, eliminating the need for performing trial updates on the model. This results in significant time savings compared to traditional approaches. The slight time increase of VeCAF over ActiveFT is caused by the text embedding CEA process. Meanwhile, considering the > 150s model training time in each loop, this 4-5s (3%) time overhead is negligible, and also justifiable considering the benefits offered by VeCAF.

C.3. Accuracy under unlimited training batches

This work concentrates on an efficient training paradigm, and accordingly, we have presented most of our experimental outcomes in the primary manuscript using a limited

number of batches. This approach inherently benefits methodologies that enable quicker convergence. To thoroughly assess VeCAF’s convergence efficacy, we have lifted the constraints on the number of training batches in this section and conducted a comparative analysis of VeCAF’s final convergence metrics against established active learning frameworks. The results, delineated in 8, demonstrate that VeCAF not only achieves expedited convergence but also surpasses previous active learning strategies in terms of final performance metrics. Remarkably, VeCAF attains performance on par with comprehensive finetuning by utilizing merely 5% of the data for CIFAR-10 and 10% for Caltech101. These findings suggest that VeCAF is capable of significantly enhancing both computational and data efficiency throughout the PVM finetuning procedure.