

# Robust fine-tuning of zero-shot models

Mitchell Wortsman<sup>\*†</sup>   Gabriel Ilharco<sup>\*†</sup>   Mike Li<sup>‡</sup>   Jong Wook Kim<sup>§</sup>

Hannaneh Hajishirzi<sup>†◦</sup>   Ali Farhadi<sup>\*†</sup>   Hongseok Namkoong<sup>\*‡</sup>   Ludwig Schmidt<sup>†Δ</sup>

## Abstract

Large pre-trained models such as CLIP offer consistent accuracy across a range of data distributions when performing zero-shot inference (i.e., without fine-tuning on a specific dataset). Although existing fine-tuning approaches substantially improve accuracy in-distribution, they also reduce out-of-distribution robustness. We address this tension by introducing a simple and effective method for improving robustness: ensembling the weights of the zero-shot and fine-tuned models. Compared to standard fine-tuning, the resulting weight-space ensembles provide large accuracy improvements out-of-distribution, while matching or improving in-distribution accuracy. On ImageNet and five derived distribution shifts, weight-space ensembles improve out-of-distribution accuracy by 2 to 10 percentage points while increasing in-distribution accuracy by nearly 1 percentage point relative to standard fine-tuning. These improvements come at no additional computational cost during fine-tuning or inference.

## 1 Introduction

A foundational goal of machine learning is to develop models that work reliably across a broad range of data distributions. Recently, large pre-trained models such as CLIP [80] and ALIGN [48] have demonstrated unprecedented robustness to challenging distribution shifts where prior robustness interventions failed to improve performance [92]. While these results point towards pre-training on large, heterogeneous datasets as a promising direction for increasing robustness, an important caveat is that these robustness improvements occur only in the zero-shot setting, i.e., when the model performs inference without fine-tuning on a specific target distribution [3, 8].

In a concrete application, a zero-shot model can be fine-tuned on extra application-specific data, which often yields large performance gains on the target distribution. However, Radford et al. [80] have shown that current fine-tuning techniques carry a hidden cost: the out-of-distribution accuracy of their fine-tuned CLIP models is often lower than that of the original zero-shot models. This leads to a natural question:

*Can zero-shot models be fine-tuned without reducing out-of-distribution accuracy?*

As pre-trained models are becoming a cornerstone of machine learning, techniques for fine-tuning them on downstream applications are increasingly important. Indeed, the question of robustly fine-tuning pre-trained models has recently also been raised as an open problem by several authors [3, 8]. Andreassen et al. [3] explored several fine-tuning approaches but found that none yielded models with improved robustness at high accuracy. Furthermore, Taori et al. [92] demonstrated that no current algorithmic robustness interventions provide consistent gains across the distribution shifts where zero-shot CLIP exhibits robustness on.

In this paper, we introduce a new way of fine-tuning zero-shot models that addresses the aforementioned question and achieves the best of both worlds: increased performance out-of-distribution while maintaining

<sup>\*,\*</sup>These authors contributed equally. Correspondence to {mitchnw, gamaga, schmidt}@cs.washington.edu.

<sup>†</sup>University of Washington <sup>‡</sup>Columbia University <sup>§</sup>OpenAI <sup>◦</sup>Allen Institute for Artificial Intelligence <sup>Δ</sup>Toyota Research

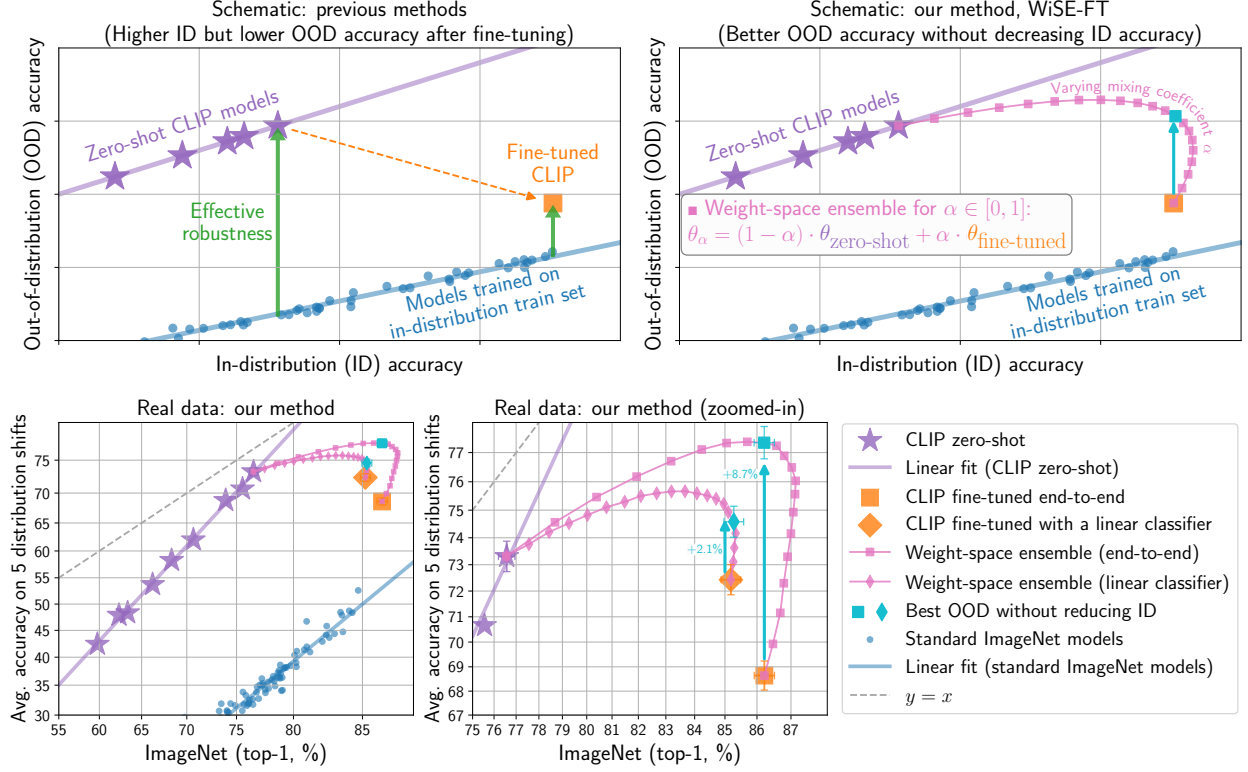


Figure 1: Compared to standard fine-tuning, weight-space ensembles for fine-tuning (WiSE-FT) improve out-of-distribution accuracy without decreasing in-distribution performance. **(Top left)** Zero-shot CLIP models [80] exhibit high effective robustness and moderate in-distribution accuracy, while standard fine-tuning (end-to-end or with a linear classifier) attains higher in-distribution accuracy and less effective robustness. **(Top right)** Our method linearly interpolates between the zero-shot and fine-tuned models with a mixing coefficient  $\alpha$ . Curves are drawn by varying  $\alpha \in [0, 1]$ . **(Bottom)** On five distribution shifts derived from ImageNet [20] (ImageNetV2 [81], ImageNet-R [40], ImageNet Sketch [95], ObjectNet [4], and ImageNet-A [41]), WiSE-FT improves average out-of-distribution accuracy by 8.7 percentage points (pp) when fine-tuning end-to-end (+2.1 pp when fine-tuning a linear classifier) while maintaining accuracy in-distribution.

or even improving in-distribution accuracy relative to standard fine-tuning. Our method (Figure 1) has two steps: first, we fine-tune the zero-shot model on application-specific data. Second, we combine the original zero-shot and fine-tuned models by linearly interpolating between their weights, which is known as weight-space ensembling.

Compared to standard fine-tuning, weight-space ensembles for fine-tuning (WiSE-FT) substantially improve out-of-distribution accuracy without decreasing in-distribution performance. Concretely, on ImageNet and five of the natural distribution shifts studied by Radford et al. [80], WiSE-FT applied to standard end-to-end fine-tuning improves out-of-distribution accuracy by 2 to 15 percentage points (pp) while maintaining the in-distribution accuracy of the fine-tuned model, and improves out-of-distribution accuracy by 1 to 9 pp relative to the zero-shot model. Moreover, WiSE-FT provides large gains when compared to a range of alternative approaches such as regularization, distillation, and ensembling softmax outputs. The robustness comes at no additional computational cost during fine-tuning or inference because the zero-shot and fine-tuned models are ensembled into a single model of the same size.

To understand the robustness gains of WiSE-FT, we empirically analyze ensembling through the lens of

distributional robustness. First, we study WiSE-FT with linear (last layer) fine-tuning as it is amenable to analysis: our procedure reduces to ensembling the outputs of two models, which highlights the complementarity of model predictions as a key property. We illustrate via detailed measurements how the predictions of zero-shot and fine-tuned models are diverse, and that models are more confident on the parts of the test distributions they perform well on.

For end-to-end fine-tuning, we connect our observations to earlier work on the phenomenology of deep learning. Neyshabur et al. [72] found that end-to-end fine-tuning the same model twice yielded two different solutions that were connected via a linear path in weight space along which error remains low, which is also known as linear mode connectivity [29]. The authors therefore concluded that the two solutions are in the same basin of the loss landscape. Interestingly, linear interpolation in weight space succeeds despite non-linearity in the activation functions of the neural networks. Our observations point to a similarly benign loss landscape in the case of WiSE-FT, but the exact shape of the loss landscape and the connection between error in- and out-of-distribution are still an open problem.

The large robustness gains and simplicity of weight-space ensembling offer a promising direction in learning distributionally robust models. While ensembles are a classical technique in machine learning, they are usually studied in a setting where every part of the ensemble is trained on the same data distribution. In contrast, we study ensembling of models for the same classification task, but trained on different data distributions. Our results suggest ensembling models across distributions may offer advantages that have been overlooked by the traditional perspective of measuring performance on a single test distribution.

Beyond the robustness viewpoint, we also study the benefits of weight-space ensembling in the low-data regime. The low-data regime is particularly relevant for fine-tuning in practice because assembling a training set of the size of ImageNet is often prohibitively expensive. When five examples per class are used for fine-tuning, WiSE-FT improves in-distribution performance by 0.3 to 6.1 percentage points compared to the best of the zero-shot and fine-tuned models on a range of seven datasets (ImageNet [20], CIFAR-10, CIFAR-100 [55], Describable Textures [17], Food-101 [9], SUN397 [98], and Stanford Cars [54]).

**Paper outline.** We introduce our experimental setup and the distribution shifts we consider in Section 2. Section 3 then formally describes WiSE-FT. Section 4 contains our main results on ImageNet, along with more experiments in the low-data regime and on additional datasets. Section 5 discusses possible mechanisms behind and further analyzes the observed empirical trends exhibited by WiSE-FT. We review related work in Section 6 and then conclude the paper in Section 7.

## 2 Background and experimental setup

Our main experiments compare the performance of the zero-shot model, fine-tuned model, and their weight-space ensemble. Primarily, we contrast model accuracy on data from two distributions. The first distribution, referred to as in-distribution (ID) and denoted  $\mathcal{D}$ , has both a train set  $\mathcal{S}_{\mathcal{D}}^{\text{tr}} = \{(x_i, y_i)\}_{i=1}^n$  and a disjoint test set consisting of images sampled from  $\mathcal{D}$ . The second distribution, referred to as out-of-distribution (OOD) and denoted  $\mathcal{D}'$ , has only a test set. For a model  $f$  we let  $\text{acc}_{\mathcal{D}}(f)$  and  $\text{acc}_{\mathcal{D}'}(f)$  refer respectively to classification accuracy on the in- and out-of-distribution test sets. We consider  $k$ -way image classification where  $x_i$  is an image with corresponding label  $y_i \in \{1, \dots, k\}$ . The outputs of  $f$  are  $k$  dimensional vectors, where index  $j$  is the non-normalized score assigned to class  $j$ .

**Distribution shifts.** Distribution shifts can be broadly characterized into a) synthetic, e.g.  $\ell_{\infty}$ -adversarial examples or artificial changes in image contrast, brightness, etc. [38, 93, 32, 1]; and b) natural, where samples are not perturbed after acquisition and changes in data distributions arise through naturally occurring variations in lighting, geographic location, crowdsourcing process, image styles, etc. [92, 40, 50, 95, 4, 41]. Our primary focus is natural distribution shifts, which are more likely to occur in the real world.

Specifically, our key results are shown for five distribution shifts illustrated in Figure 2 where  $\mathcal{S}_{\mathcal{D}}^{\text{tr}}$  is ImageNet [20]. We consider:

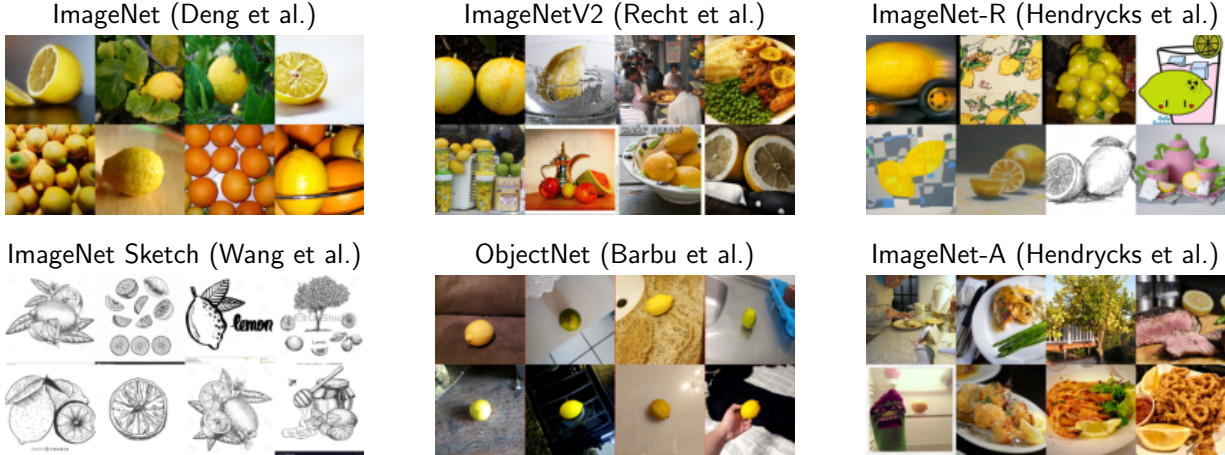


Figure 2: Samples of the class *lemon*, from ImageNet (in-distribution) and the derived out-of-distribution datasets considered in our main experiments [20, 81, 40, 95, 4, 41].

- ImageNet-V2 (IN-V2) [81], a reproduction of the ImageNet test set with distribution shift;
- ImageNet-R (IN-R) [40], renditions (e.g. sculptures, paintings) for 200 ImageNet classes;
- ImageNet Sketch (IN-Sketch) [95], a test set which contains sketches instead of natural images;
- ObjectNet [4], a test set of objects in various scenes with 113 classes overlapping with ImageNet;
- ImageNet-A (IN-A) [41], a test set of natural images misclassified by a ResNet-50 [37] for 200 ImageNet classes.

For consistency, we refer to ImageNet as in-distribution and the remainder as out-of-distribution even for zero-shot models which are never trained in-distribution.

**Effective robustness and scatter plots.** Central to our analysis are effective robustness scatter plots, which illustrate model performance under distribution shift [81, 92]. These scatter plots display in-distribution accuracy on the  $x$ -axis and out-of-distribution accuracy on the  $y$ -axis—a model  $f$  is shown as a point  $(\text{acc}_{\mathcal{D}}(f), \text{acc}_{\mathcal{D}'}(f))$ . Figure 1 exemplifies these scatter plots with both schematics and real data. These plots illustrate *effective robustness* [92], a tool which allows us to compare the robustness of models with different in-distribution accuracies, and we define formally below.

For a number of distribution shifts, accuracy on the standard test set is a reliable predictor of accuracy under distribution shift [92, 68]. In other words, there exists a function  $\beta : [0, 1] \rightarrow [0, 1]$  where  $\text{acc}_{\mathcal{D}'}(f) \approx \beta(\text{acc}_{\mathcal{D}}(f))$  across models  $f$  trained on the in-distribution train set  $\mathcal{S}_{\mathcal{D}}^{\text{tr}}$ . Effective robustness [92] is accuracy beyond this baseline, defined formally as  $\rho(f) = \text{acc}_{\mathcal{D}'}(f) - \beta(\text{acc}_{\mathcal{D}}(f))$ . In these scatter plots, effective robustness is vertical movement above expected out-of-distribution performance (Figure 1, top). Moreover, when we say that a model is robust to distribution shift, we mean that effective robustness is positive. We emphasize that Taori et al. [92] observed that no algorithmic robustness intervention consistently achieves substantial effective robustness across the distribution shifts in Figure 2—the first model to do so was zero-shot CLIP.

When there is a deterministic function relating in- and out-of-distribution accuracy, appropriate axis scaling will enable a linear fit in the corresponding scatter plot. Empirically, when applying logit (or probit) axis scaling, models with zero effective robustness approximately lie on a linear trend [92, 68] (Figure 1).

The scatter plots we display also include a wide range of machine learning models from a comprehensive testbed of evaluations [92, 68], including: models trained on  $\mathcal{S}_{\mathcal{D}}^{\text{tr}}$  (*standard training*); models trained on additional data and fine-tuned using  $\mathcal{S}_{\mathcal{D}}^{\text{tr}}$  (*trained with more data*); and models trained using various *existing*

*robustness interventions*, e.g. special data augmentation [22, 93, 26, 31, 39] or adversarially robust models [18, 83, 85]. Finally, we apply logit axis scaling and display 95% Clopper-Pearson confidence intervals for the accuracies of select points.

**Zero-shot models and CLIP.** Zero-shot CLIP models [80] exhibit effective robustness and lie on a qualitatively different linear trend. These models are pre-trained using a third distribution  $\mathcal{P}$ , e.g., image-caption pairs from the web. Given a set of image-caption pairs  $\{(x_1, s_1), \dots, (x_B, s_B)\}$ , CLIP trains an image-encoder  $g$  and text-encoder  $h$  such that the similarity  $\langle g(x_i), h(s_i) \rangle$  is maximized relative to unaligned pairs.

CLIP-like models perform zero-shot  $k$ -way classification given an image  $x$  and class names  $C = \{c_1, \dots, c_k\}$  by matching  $x$  with potential captions. For instance, using potential captions  $s_i = \text{“a photo of a } \{c_i\} \text{”}$  for each class  $i$ , the zero-shot model predicts the class via  $\arg \max_j \langle g(x), h(s_j) \rangle$ .<sup>1</sup> Equivalently, one can construct  $\mathbf{W}_{\text{zero-shot}} \in \mathbb{R}^{d \times k}$  with columns  $h(s_j)$  and compute outputs  $f(x) = g(x)^\top \mathbf{W}_{\text{zero-shot}}$ . Unless explicitly mentioned, our experiments use the CLIP model ViT-L/14@336px (based on vision transformers [25], additional details provided in Section D.1 of the Appendix).

### 3 Weight-space ensembling for robust fine-tuning

This section describes and motivates our proposed method, WiSE-FT. At a high level, WiSE-FT consists of two steps. First, we fine-tune the zero-shot model on application-specific data. Second, we combine the original zero-shot and fine-tuned models by linearly interpolating between their weights, also referred to as weight-space ensembling.

The zero-shot model excels under distribution shift while standard fine-tuning achieves high accuracy in-distribution. Our motivation is to combine these two models into one that achieves the best of both worlds. Weight-space ensembles are a natural choice as they ensemble models with no extra computational cost. Moreover, previous work has suggested that interpolation in weight space may improve performance when models share part of their optimization trajectory [46, 72].

**Step 1: Standard fine-tuning.** As in Section 2, we let  $\mathcal{S}_D^{\text{tr}}$  denote the dataset used for fine-tuning and  $g$  denote the image encoder used by CLIP [80]. We are now explicit in writing  $g(x, \theta_{\text{enc}})$  where  $x$  is an input image and  $\theta_{\text{enc}}$  are the parameters of the image-encoder  $g$  learned during pre-training.

Standard fine-tuning considers the model  $f(x, \theta) = g(x, \theta_{\text{enc}})^\top \mathbf{W}_{\text{classifier}}$  where  $\mathbf{W}_{\text{classifier}} \in \mathbb{R}^{d \times k}$  is the classification head and  $\theta = [\theta_{\text{enc}}, \mathbf{W}_{\text{classifier}}]$  are the parameters of  $f$ . We then proceed with the optimization problem  $\arg \min_{\theta} \left\{ \sum_{(x_i, y_i) \in \mathcal{S}_D^{\text{tr}}} \ell(f(x_i, \theta), y_i) + \lambda R(\theta) \right\}$  where  $\ell$  is the cross-entropy loss and  $R$  is a regularization term (e.g. weight decay).

We consider the two most common variants of fine-tuning: end-to-end, where all values of  $\theta$  are modified, and fine-tuning only a linear classifier, where  $\theta_{\text{enc}}$  is fixed at the value learned during pre-training. While end-to-end fine-tuning offers better accuracy, fine-tuning a linear classifier requires less computational resources allowing a comprehensive suite of experiments [80, 53]. Additional details are provided in Section 4.2 and Appendix D.4.

**Step 2: Weight-space ensembling.** For a *mixing coefficient*  $\alpha \in [0, 1]$ , we consider the *weight-space ensemble* between the zero-shot model with parameters  $\theta_0$  and the model obtained via standard fine-tuning with parameters  $\theta_1$ . The predictions of the corresponding weight-space ensemble **wse** are given by

$$\text{wse}(x, \alpha) = f(x, (1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1) , \quad (1)$$

---

<sup>1</sup>In practice it is helpful to consider a few candidate captions using the class name  $c_i$ , e.g.  $s_i^{(1)} = \text{“a photo of a } \{c_i\} \text{”}$  and  $s_i^{(2)} = \text{“a picture of a } \{c_i\} \text{”}$ , and use their re-normalized average embedding (referred to as prompt ensembling [80]).



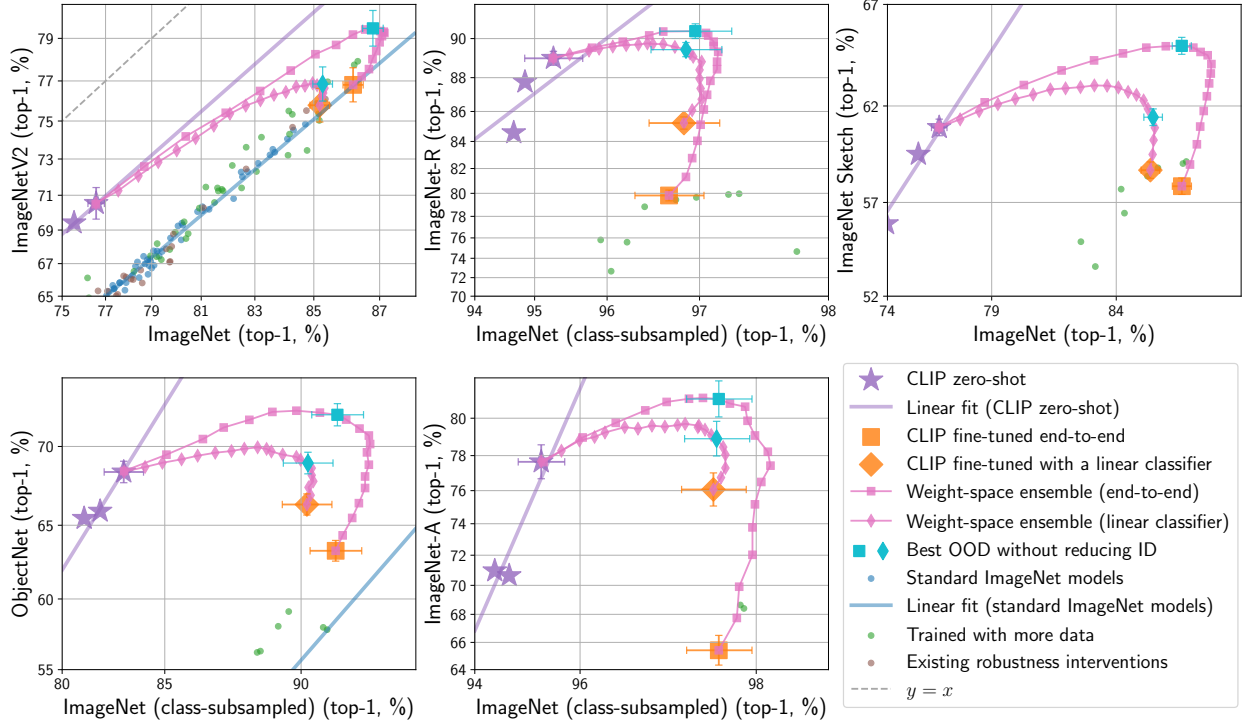


Figure 3: WiSE-FT improves in- and out-of-distribution accuracy across five distribution shifts derived from ImageNet (Figure 2, [20, 81, 40, 95, 4, 41]). Standard ImageNet models, models trained with more data, and existing robustness interventions are from Taori et al. [92]. A zoomed-out version of this figure is available in Appendix F (Figure 14).

i.e., we use the element-wise weighted average of the zero-shot and fine-tuned parameters. When fine-tuning only the linear classifier, weight-space ensembling is equivalent to the traditional output-space ensemble [23]  $(1 - \alpha) \cdot f(x, \theta_0) + \alpha \cdot f(x, \theta_1)$  since Equation 1 decomposes as

$$\begin{aligned} \text{wse}(x, \alpha) &= g(x, \theta_{\text{enc}})^\top ((1 - \alpha) \cdot \mathbf{W}_{\text{zero-shot}} + \alpha \cdot \mathbf{W}_{\text{classifier}}) \\ &= (1 - \alpha) \cdot g(x, \theta_{\text{enc}})^\top \mathbf{W}_{\text{zero-shot}} + \alpha \cdot g(x, \theta_{\text{enc}})^\top \mathbf{W}_{\text{classifier}}. \end{aligned} \quad (2)$$

As neural networks are non-linear with respect to their parameters, ensembling all layers—as we do when end-to-end fine-tuning—typically fails, achieving no better accuracy than a randomly initialized neural network [29]. However, as similarly observed by previous work where part of the optimization trajectory is shared [46, 29, 72], we find that the zero-shot and fine-tuned models are connected by a linear path in weight-space along which accuracy remains high (explored further in Section 5.2).<sup>2</sup>

Remarkably, as we will show in Section 4, WiSE-FT boosts out-of-distribution accuracy relative to the fine-tuned model without decreasing in-distribution performance. These improvements come without any additional computational cost during fine-tuning or inference as a single set of weights are used. We provide pseudocode for WiSE-FT in a PyTorch-like style in Appendix A.

<sup>2</sup>In the context of fine-tuning, previous work has only considered linear paths of high accuracy between two fine-tuned models. As we are using a zero-shot model and  $\mathbf{W}_{\text{classifier}}$  is initialized with  $\mathbf{W}_{\text{zero-shot}}$ , we consider paths between the pre-trained and fine-tuned models.

	IN (ID)	OOD datasets					Avg OOD	Avg ID, OOD
		IN-V2	IN-R	IN-Sketch	ObjectNet*	IN-A		
NS EfficientNet-L2 [99]	88.4	80.2	74.7	47.6	68.5	<b>84.9</b>	71.2	79.8
ViT-G/14 [102]	<b>90.4</b>	<b>83.3</b>	-	-	70.5	-	-	-
Zero-shot ALIGN [48]	76.4	70.1	<b>92.2</b>	-	-	70.1	-	-
CLIP-based models								
Zero-shot [80]	76.2	70.1	88.9	60.2	70.0	77.2	73.3	74.8
Fine-tuned linear classifier [80]	85.4	75.9	84.2	57.4	66.2	75.3	71.8	78.6
Zero-shot (PyTorch)	76.6	70.5	89.0	60.9	69.1	77.7	73.4	75.0
Fine-tuned linear classifier (LC, ours)	85.2	75.8	85.3	58.7	67.2	76.1	72.6	78.9
End-to-end fine-tuned (E2E, ours)	86.2	76.8	79.8	57.9	63.3	65.4	68.6	77.4
Weight-space ensembles (WiSE-FT, ours)								
LC, $\alpha=0.75$	85.1	76.8	88.4	61.9	69.7	78.9	75.1	80.1
LC, $\alpha=0.4$	82.7	75.8	89.7	63.0	70.7	79.6	75.8	79.2
LC, optimal $\alpha$	85.3	76.9	89.8	63.0	70.7	79.7	76.0	80.7
E2E, $\alpha=0.75$	87.0	78.8	86.1	62.5	68.1	75.2	74.1	80.5
E2E, $\alpha=0.4$	86.2	79.2	89.9	<b>65.0</b>	71.9	80.7	77.3	81.8
E2E, optimal $\alpha$	<u>87.1</u>	<u>79.5</u>	<u>90.3</u>	<b>65.0</b>	<b>72.1</b>	<u>81.0</u>	<b>77.6</b>	<b>82.3</b>

Table 1: Compared to standard end-to-end (E2E) fine-tuning, WiSE-FT improves out-of-distribution accuracy by 2 to 15 percentage points (pp), without decreasing in-distribution performance ( $\alpha = 0.4$ ). Alternatively, with a mixing coefficient  $\alpha = 0.75$ , performance on ImageNet improves by 0.9 pp, while out-of-distribution accuracy increases by 2 to 10 pp. The highest overall accuracy for each dataset is in bold, while the highest accuracy among models derived from CLIP are underlined. *Avg OOD* displays the mean performance among the five out-of-distribution datasets, while *Avg ID, OOD* shows the average of ImageNet (ID) and Avg OOD.

## 4 Results

This section presents our key experimental findings. Section 4.1 demonstrates that WiSE-FT boosts accuracy of a fine-tuned CLIP model on five distribution shifts studied by Radford et al. [80], while maintaining or improving accuracy on ImageNet. Section 4.2 then develops and contrasts WiSE-FT with a series of alternatives including regularization, distillation, and ensembling softmax outputs. Next, Section 4.3 demonstrates that when limited labeled data is available, WiSE-FT reliably outperforms both zero-shot and fine-tuned models on the datasets we consider. Finally, Sections 4.4 and 4.5 respectively explore distribution shifts and a variety of model sizes. While many of our experiments consider end-to-end fine-tuning, fine-tuning only a linear classifier requires far less computational resources, allowing a more comprehensive suite of experiments. Accordingly, it is more widely used throughout this section.

### 4.1 Main results: ImageNet and associated distribution shifts

Our main results on ImageNet and five derived distribution shifts are presented in Tables 1-2 and Figure 3. As illustrated in Figure 3, when the mixing coefficient  $\alpha$  varies from 0 to 1 the performance of  $\text{wse}(\cdot, \alpha)$  does not move linearly from the zero-shot model to the fine-tuned model. Accuracy instead detours up and right, improving performance in- and out-of-distribution.

Concretely, compared to standard fine-tuning, WiSE-FT improves out-of-distribution performance by 2 to 15 percentage points (pp), without reducing in-distribution performance ( $\alpha = 0.4$ , end-to-end). Note that fixing  $\alpha = 0.4$  is close to optimal for all out-of-distribution datasets, with only a 0.3 pp difference on average (Table 1). Alternatively, when  $\alpha = 0.75$ , ImageNet performance improves by 0.9 pp while out-of-distribution performance is boosted by 2 to 10 pp compared to standard end-to-end fine-tuning.

Table 2 offers an alternative perspective, demonstrating (a) gains out-of-distribution relative to the fine-tuned

	FT	IN-V2	IN-R	IN- Sketch	ObjectNet	IN-A	Avg
OOD improvement relative to fine-tuned without decreasing ID	LC	1.1	4.2	2.7	2.6	2.9	2.7
	E2E	2.6	10.5	7.1	8.8	15.6	8.9
ID improvement relative to zero-shot without decreasing OOD	LC	8.8	1.7	8.7	7.2	2.4	5.8
	E2E	10.6	1.9	10.6	8.8	2.7	6.9

Table 2: Compared to the fine-tuned model, WiSE-FT improves out-of-distribution accuracy without reducing in-distribution performance. Compared to the zero-shot model, WiSE-FT improves in-distribution accuracy without reducing out-of-distribution performance. All numbers are percentage point improvements, and FT indicates whether end-to-end fine-tuning (E2E) or fine-tuning only the linear classifier (LC). Top and bottom respectively capture vertical movement above the fine-tuned model and horizontal movement to the right of the zero-shot model in associated scatter plots.

	IN-V2	IN-R	IN-Sketch	ObjectNet	IN-A	Avg
Weight-space ensemble	<b>1.1</b>	<b>4.2</b>	2.7	2.6	<b>2.9</b>	<b>2.7</b>
Output-space ensemble	0.8	4.1	<b>3.1</b>	<b>2.7</b>	2.7	<b>2.7</b>
Distillation	0.4	2.7	1.2	1.0	1.7	1.4
Regularize to zero-shot		3.6		1.8		1.1

Table 3: Out-of-distribution accuracy gain (percentage points) without any loss in-distribution relative to the fine-tuned linear classifier for the methods described in Section 4.2.

model without losing in-distribution accuracy and (b) gains in-distribution relative to the zero-shot classifier without losing out-of-distribution accuracy. Although fine-tuning only a linear classifier leads to typically smaller gains compared to end-to-end, accuracy out-of-distribution still improves by 1 to 4 pp without any loss in-distribution.

## 4.2 Comparison to alternative fine-tuning methods

In searching for a method of fine-tuning which preserves robustness, we explore a variety of alternatives. Many exhibit a concave trend in effective robustness plots, although WiSE-FT offers the best results overall. This section describes these methods and their performance is compared in Figure 4. We primarily focus on methods which fine-tune the linear classifier, allowing comprehensive experimentation.

**Random interpolation.** This method uses either the zero-shot or fine-tuned linear classifier depending on a (biased) coin flip. For hyperparameter  $\alpha \in [0, 1]$  outputs are computed as  $(1 - \xi) \cdot f(x, \theta_0) + \xi \cdot f(x, \theta_1)$  where  $\xi$  is a Bernoulli( $\alpha$ ) random variable. For this method and all others with a hyperparameter  $\alpha \in [0, 1]$  we evaluate models for  $\alpha \in \{0, 0.05, 0.1, \dots, 1\}$ .

**Ensembling softmax outputs.** Instead of ensembling in weight space, this method combines softmax probabilities assigned by the zero-shot and fine-tuned linear classifier. Concretely, for hyperparameter  $\alpha \in [0, 1]$  outputs are computed as  $(1 - \alpha) \cdot \text{softmax}(f(x, \theta_0)) + \alpha \cdot \text{softmax}(f(x, \theta_1))$ . This method performs comparably to weight-space ensembling but requires slightly more compute.

**Linear classifier with various regularizers.** We explore fine-tuning linear classifiers with four regularization strategies: no regularization, weight decay, L1 regularization, and label smoothing [70]. Linear-classifiers are trained with mini-batch optimization, using the AdamW optimizer [61, 77] with a cosine-annealing learning rate schedule [60]. This method is significantly faster and less memory-intensive than the L-BFGS

\*Although this table considers ImageNet class names only, ObjectNet provides alternative class names which can improve the performance of zero-shot CLIP by 2.3 percentage points. See Appendix D.3 for details and comparisons.



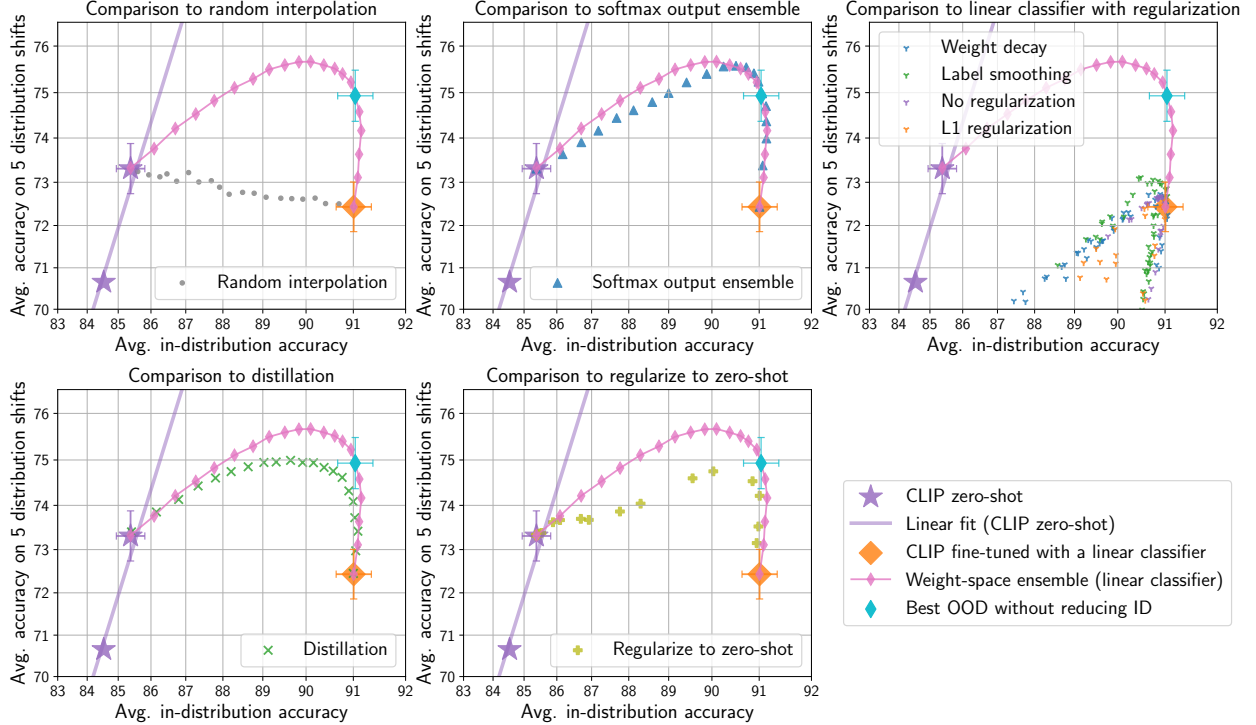


Figure 4: Comparing the relative in- and out-of-distribution performance of weight-space ensembling with the alternatives described in Section 4.2. Many methods follow a concave trend, though weight-space ensembling provides the best performance overall.

implementation used by Radford et al. [80] at ImageNet scale with similar accuracy. Additional details on hyperparameters and more analyses are provided in Section D.2.

**Distillation.** Network distillation [44] trains one network to match the outputs of another. We use this technique to fine-tune while matching the outputs of the zero-shot model, in an attempt to boost out-of-distribution performance. For a hyperparameter  $\alpha \in [0, 1]$  and cross-entropy loss  $\ell$  we fine-tune  $\theta$  according to the minimization objective

$$\sum_{(x_i, y_i) \in \mathcal{S}_D^H} (1 - \alpha) \cdot \ell(f(x_i, \theta), y_i) + \alpha \cdot \ell(f(x_i, \theta), f(x_i, \theta_0)) . \quad (3)$$

**Regularization towards zero-shot.** We train a linear classifier with an additional regularization term which penalizes movement from the zero-shot classifier’s weights. For a hyperparameter  $\lambda \in \{1 \cdot 10^{-8}, 5 \cdot 10^{-8}, 1 \cdot 10^{-7}, \dots, 5 \cdot 10^2\}$  we add the regularization term  $\lambda \|\mathbf{W} - \mathbf{W}_{\text{zero-shot}}\|_F^2$  where  $\mathbf{W}$  is the linear classifier being fine-tuned. In most cases this method performs slightly worse than distillation.

For completeness, Section F of the Appendix displays a version of Figure 4 for each dataset. Table 3 offers an alternative perspective, showing the amount of out-of-distribution accuracy which can be gained without reducing in-distribution performance for each method. Finally, Figure 5 considers baselines for the end-to-end fine-tuned model. We first compare with ensembling in output-space, i.e.,  $(1 - \alpha) \cdot f(x, \theta_0) + \alpha \cdot f(x, \theta_1)$ . Next, we consider the evaluation of checkpoints obtained at every epoch of fine-tuning. We find WiSE-FT outperforms both methods.

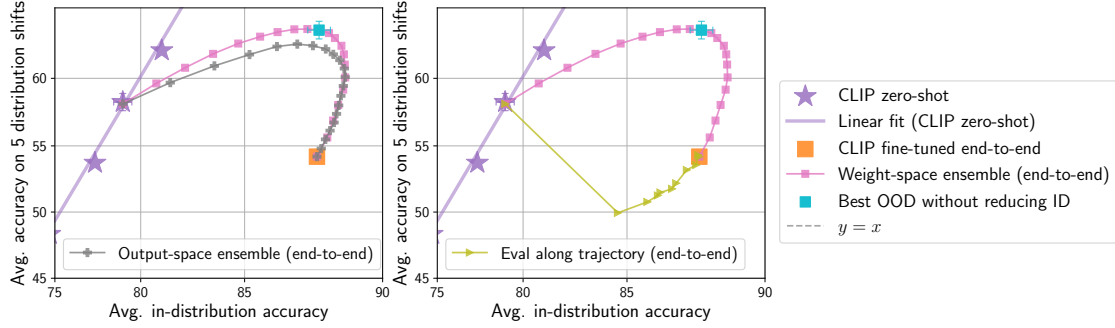


Figure 5: When fine-tuning end-to-end, weight-space ensembles outperform output-space ensembles and checkpoints saved at each epoch of fine-tuning (results shown for the ViT-B/16 CLIP model).

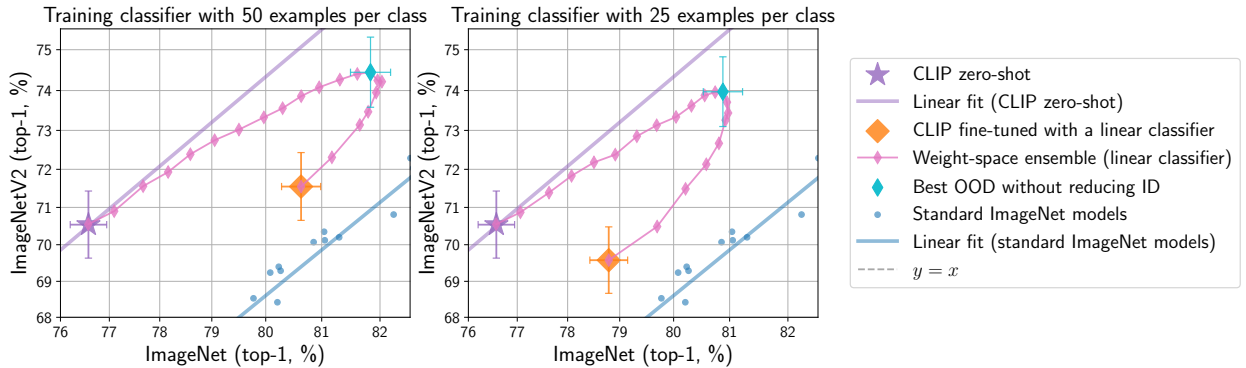


Figure 6: WiSE-FT improves both in- and out-of-distribution accuracy in the low data regime.

### 4.3 In-distribution gains in the low data regime

Our investigation so far has focused primarily on robustness and ImageNet scale. However, in most cases fine-tuning is performed on a much smaller dataset. In this section, we begin by fine-tuning a linear classifier on smaller subsets of ImageNet with 25 and 50 random examples per class. As illustrated in Figure 6, weight-space ensembles provide better accuracy than either the zero-shot model or linear classifier. This accuracy boost holds not only under distribution shift, but also on the standard test set. Recall that when fine-tuning on the full ImageNet training set, performance gains in-distribution were only realised when using the end-to-end variant of WiSE-FT.

Beyond robustness, weight-space ensembles can improve in-distribution accuracy when the training set is small. We examine the performance of weight-space ensembles on ImageNet in addition to a number of datasets considered by Kornblith et al. [53]: CIFAR-10, CIFAR-100 [55], Describable Textures [17], Food-101 [9], SUN397 [98], and Stanford Cars [54] when  $k$  examples per class are used for fine-tuning a linear classifier ( $k = \{1, 5, 10, 25, 50\}$ ). Average results are shown in Figure 7, while Appendix F provides a breakdown for all datasets.

### 4.4 Robustness on additional distribution shifts

This section explores the performance of WiSE-FT on additional natural distribution shifts, when fine-tuning a linear classifier. We consider: (i) ImageNet-Vid-Robust and YTBB-Robust, datasets with distribution shift induced by temporal perturbations in videos [86]; (ii) CIFAR-10.1 [81] and CIFAR-10.2 [62], reproductions of the popular image classification dataset CIFAR-10 [55] with a distribution shift; (iii) WILDS-FMoW, a

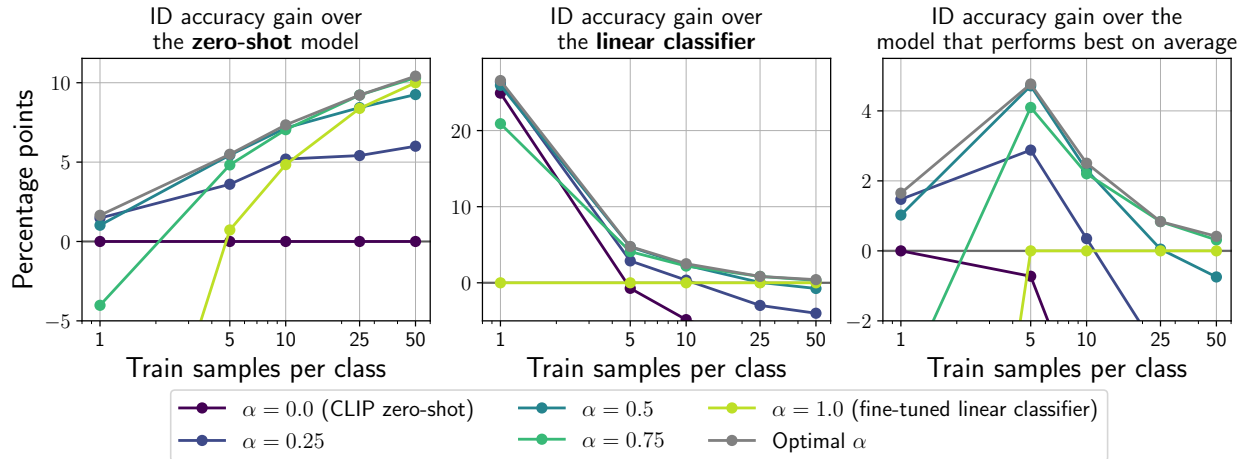


Figure 7: **WiSE-FT can improve in-distribution accuracy over the linear classifier and zero-shot model in the low data regime.** On the  $x$ -axis we consider  $k = \{1, 5, 10, 25, 50\}$  examples per class for fine-tuning. On the  $y$ -axis we display in-distribution accuracy improvements of WiSE-FT averaged over seven datasets [20, 55, 17, 9, 98, 98, 54]. For  $k = 1$ , the zero-shot model outperforms the fine-tuned linear classifier, and ensembles closer to the zero-shot model (small  $\alpha$ ) yield high performance. When more data is available, the reverse is true, and higher values of  $\alpha$  improve in-distribution performance. Appendix F, displays a breakdown for all datasets.

satellite image recognition task where the test set has a geographic and temporal distribution shift [50, 16]; (iv) WILDS-iWildCam, a wildlife recognition task where the test set has a geographic distribution shift [50, 6].

For ImageNet-Vid-Robust and YTBB-Robust we use ImageNet as the in-distribution set for fine-tuning. However, the treatment of these datasets is slightly different from those explored in Section 4.1—the zero-shot classifier is constructed using class names provided by Shankar et al. [86] instead of the ImageNet class names. This technique is referred to by Radford et al. [80] as adapting to class shift, and is only possible for YTBB-Robust (+26.9 pp), ImageNet-Vid-Robust (+8.3 pp) and ObjectNet (+2.3 pp) [80]. For CIFAR-10.1 and CIFAR-10.2 we fine-tune on CIFAR-10. For both WILDS-iWildCam and WILDS-FMoW, we use the in-distribution train set provided by WILDS [50].

Results are illustrated in Figure 8, showcasing success and failure cases of WiSE-FT. On both ImageNet-Vid-Robust and YTBB-Robust, WiSE-FT provides improved out-of distribution performance over the linear classifier without reducing in-distribution accuracy. However, for YTBB-Robust, this solution still performs substantially worse out-of-distribution than the zero-shot model. We anticipate that this discrepancy arises due to the substantial difference out-of-distribution between the zero-shot and fine-tuned model (over 30%). On CIFAR-10.1, WILDS-FMoW, and WILDS-iWildCam, WiSE-FT follows the purple trend (fit to zero-shot CLIP models) across a change in accuracy that is especially notable for the WILDS distribution shifts. In other words, WiSE-FT achieves the effective robustness of a hypothetical better CLIP model. However, only on WILDS-iWildCam is WiSE-FT able to improve accuracy out-of-distribution with zero loss in-distribution. Nevertheless, these results demonstrate that WiSE-FT is a useful technique across diverse natural distribution shifts. On CIFAR-10.1, accuracy out-of-distribution can be improved with minimal loss in-distribution and on WILDS-iWildCam, WiSE-FT outperforms or matches all baselines both in- and out-of-distribution. CIFAR-10.2 is a failure case which remains to be understood, as WiSE-FT provides nearly linear movement between the zero-shot and fine-tuned model.

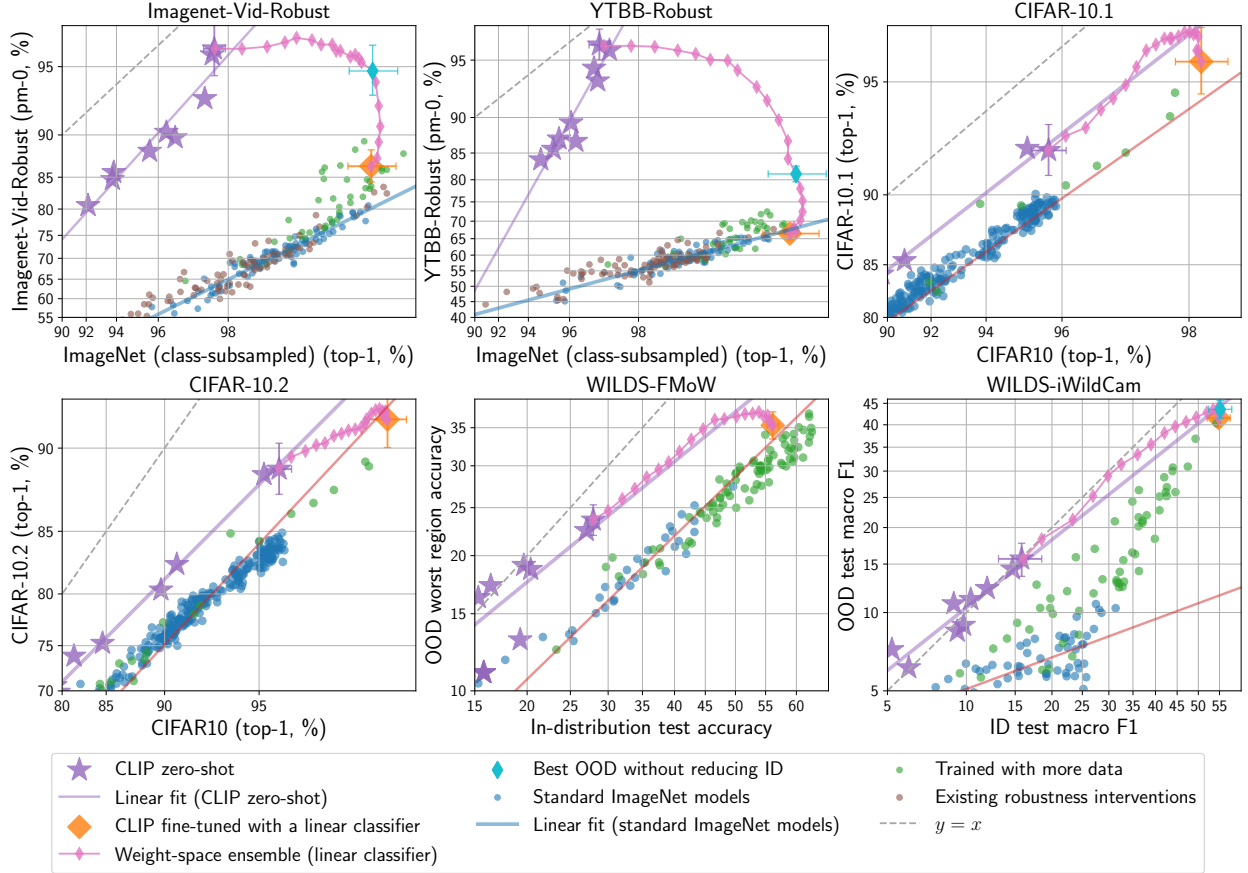


Figure 8: **WiSE-FT provides benefits for additional distribution shifts.** For ImageNet-Vid-Robust, YTBB-Robust [86], and WILDS-iWildCam [50, 6], WiSE-FT improves out-of-distribution accuracy without any loss in-distribution. For CIFAR-10.1 [97] and WILDS-FMoW [50, 16], WiSE-FT climbs the linear trend fit to zero-shot CLIP models.

#### 4.5 Robustness across scales of pre-training compute

The strong correlation between standard test accuracy and accuracy under distribution shift holds from low to high performing models. This offers the opportunity to explore robustness for smaller, easy to run models. Our exploration began with the lowest accuracy CLIP models and similar trends held at scale. Figure 9 shows improved out-of-distribution accuracy with minimal loss in-distribution across orders of magnitude of pre-training compute. Moreover, in Appendix F (Figure 22) we recreate the main results (Figure 3) with a smaller CLIP ViT-B/16 model, finding similar trends.

## 5 Discussion

This section further analyzes the empirical phenomena we observed so far. We begin with the case where only the final linear layer of the network is fine-tuned and predictions from the weight-space ensemble can be factored into the outputs of the zero-shot and fine-tuned model (Equation 2). Next, we connect our observations regarding end-to-end fine-tuning with earlier work on the phenomenology of deep learning [72, 29]. Finally, Section 5.3 investigates whether improvements are limited to fine-tuned CLIP. We find that substantial gains out-of-distribution also occur when ensembling the zero-shot model with an independently

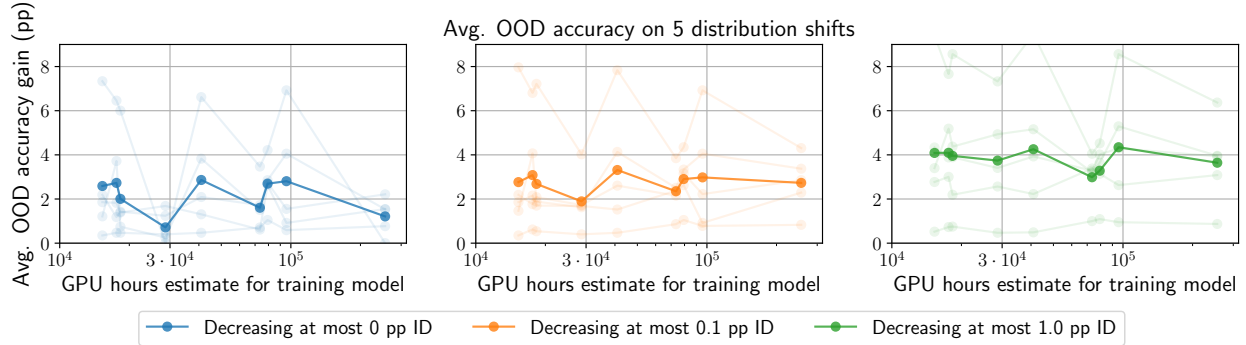


Figure 9: **WiSE-FT provides benefits for all CLIP models.** Accuracy can be improved out-of-distribution relative to the linear classifier with less than  $\epsilon \in \{0, 0.1, 1\}$  percentage points (pp) loss in-distribution across orders of magnitude of training compute. The CLIP model RN50x64 requires the most GPU hours to train.

trained network.

## 5.1 Zero-shot and fine-tuned models are complementary

Fine-tuning only the linear classifier presents a unique opportunity: in this setting weight-space ensembles exactly decompose into output-space ensembles (Equation 2). In other words, the outputs of the weight-space ensemble are a weighted average of the outputs of each model, allowing us to examine the weight-space ensemble by inspecting each component. In this section, we find that the zero-shot and fine-tuned models have diverse predictions, both in- and out-of distribution. Moreover, we show that this diversity is utilized effectively: while the fine-tuned model is more confident in-distribution, the reverse is true out-of-distribution.

**Zero-shot and specialized models are diverse.** An ensemble works best when there is diversity among the predictions of its constituents. If two models make coincident mistakes, so will their ensemble—no benefit will be gained from combining them. In contrast, if two models make uncorrelated errors, it becomes possible that they correct each other’s mistakes. Although diverse predictions do not guarantee ensemble accuracy—two random classifiers are diverse—diversity is commonly studied in conjunction with ensemble performance [57, 56, 71, 7]. Here, we explore several measures of diversity (i) Prediction Diversity, which measures the fraction of examples for which two classifiers disagree but one is correct; (ii) Cohen’s Kappa Complement, the complement of Cohen’s kappa coefficient of agreement between annotators; (iii) Average KL Divergence between softmax probabilities; (iv) Centered Kernel Alignment Complement, the complement of the similarity metric between representations proposed by Kornblith et al. [52] (CKA). More details are provided in Appendix C.

In Figure 10, we show diversity measures between zero-shot models and fine-tuned linear classifiers. Both in- and out-of-distribution, the zero-shot and fine-tuned models are diverse, despite sharing the same backbone. As a point of comparison, we include average diversity measures between two linear classifiers fine-tuned with random splits on half of ImageNet,<sup>3</sup> denoted in orange in Figure 10.

**Models are more confident where they excel.** Model diversity enables high accuracy ensembles if each model is more confident in the domain where they achieve higher accuracy. Typically, the zero-shot model has better out-of-domain performance while the fine-tuned model excels in-distribution. An effective ensemble would leverage each model’s expertise based on which distribution the data is from. Here, we empirically show that this is the case across a number of datasets we consider.

<sup>3</sup>This baseline also compares models with a shared backbone. Two linear classifiers fine-tuned on the same data will converge to the same solution, resulting in negligible diversity and no benefits from ensembling. As a stronger baseline, we use two classifiers fine-tuned on different subsets of ImageNet, with half of the data.

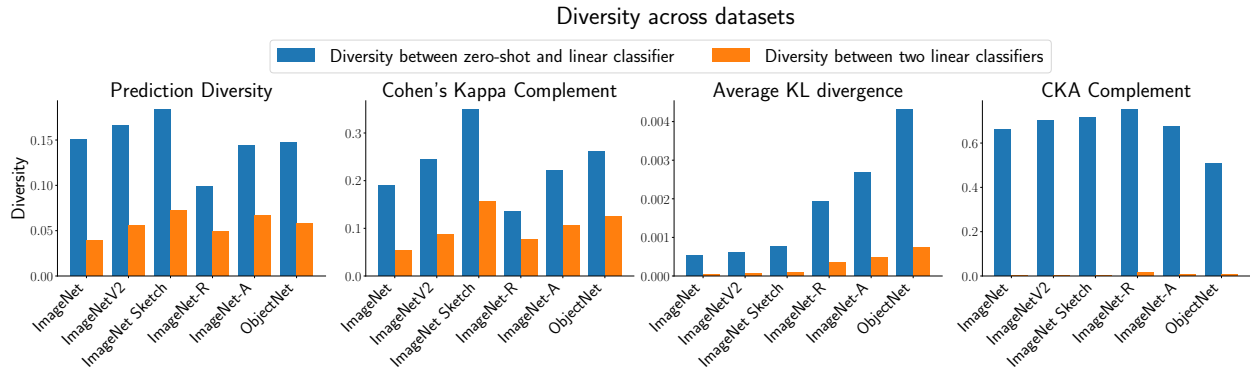


Figure 10: **Zero-shot and fine-tuned models are diverse.** Across numerous datasets and diversity measures, zero-shot and fine-tuned models exhibit high diversity in their predictions. See Appendix C for definitions and details.

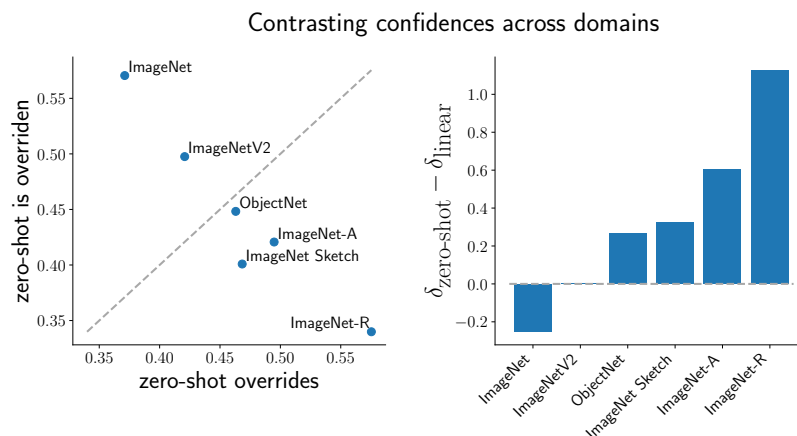


Figure 11: **Zero-shot and fine-tuned models are more confident where they excel.** (Left) In most out-of-distribution datasets, the zero-shot overrides the linear classifier more than it is overridden. The reverse is true for ImageNet (in-distribution). (Right) Similarly, zero-shot models are more confident on out-of-distribution datasets, while the reverse is true in-distribution. The margin  $\delta_f$  measures the average difference between the largest and second largest unnormalized output for classifier  $f$ .

We examine the cases where the models being ensembled disagree. We say the zero-shot model *overrides* the fine-tuned model if their predictions disagree and the zero-shot prediction matches that of the weight-space ensemble. Similarly, if models disagree and the linear classifier prediction matches the ensemble, we say the zero-shot is *overridden*.<sup>4</sup> Figure 11 (left) shows the fraction of samples where the zero-shot model overrides and is overridden by the fine-tuned linear classifier for  $\alpha = 0.5$ . Most out-of-distribution datasets lie below the diagonal—the zero-shot overrides the linear classifier more than it is overridden. The only exception is ImageNetV2, which was collected to closely reproduce ImageNet. Moreover, on ImageNet (in-distribution), the zero-shot is overridden more than it overrides.

Additionally, we are interested in measuring model confidence. Recall that we are ensembling quantities before a softmax is applied, so we avoid criteria that uses probability vectors (e.g. Guo et al. [36]). Instead, we consider the margin  $\delta$  between the largest and second largest output of each classifier. Figure 11 (right) illustrates the difference between  $\delta$  for the zero-shot and linear classifier averaged over samples. We find that

<sup>4</sup>More precisely, if  $\hat{y}_{\text{ens}}$ ,  $\hat{y}_{\text{zs}}$  and  $\hat{y}_{\text{ft}}$  are the classes predicted by the ensemble, zero-shot and fine-tuned model, the zero-shot model overrides the fine-tuned model when  $\hat{y}_{\text{zs}} \neq \hat{y}_{\text{ft}} \wedge \hat{y}_{\text{ens}} = \hat{y}_{\text{zs}}$ , and is overridden when  $\hat{y}_{\text{zs}} \neq \hat{y}_{\text{ft}} \wedge \hat{y}_{\text{ens}} = \hat{y}_{\text{ft}}$ .



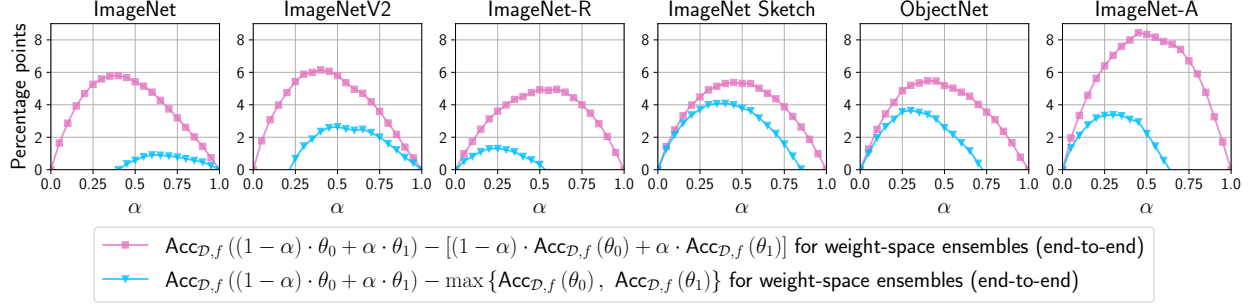


Figure 12: The zero-shot and fine-tuned models exhibit *linear mode connectivity* (Definition 1, [29]) on ImageNet and the main distribution shifts we consider (Observation 1). Moreover, there exists an  $\alpha$  for which weight-space ensembles outperform both the zero-shot and fine-tuned models (Observation 2).

zero-shot models are more confident in their predictions for out-of-distribution datasets, while the reverse is true in-distribution.

## 5.2 An error landscape perspective

Our discussion now focuses on the empirical phenomenon we observe when weight-space ensembling all layers in the network. Specifically, this section formalizes our observations and details related phenomena. Recall that the weight-space ensemble of  $\theta_0$  and  $\theta_1$  is given by  $f(x, (1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1)$  where  $\theta_0$  and  $\theta_1$  are the parameters of the zero-shot and pre-trained model, respectively.

We begin with a natural generalization of linear mode connectivity [29] to the setting where accuracy of the endpoints may differ substantially. For a distribution  $\mathcal{D}$  and model  $f$ , let  $\text{Acc}_{\mathcal{D},f}(\theta)$  denote the accuracy of the model  $f$  evaluated with parameters  $\theta$  on the test set of images sampled from  $\mathcal{D}$ .

**Definition 1:** Parameters  $\theta_0$  and  $\theta_1$  exhibit *linear mode connectivity* with respect to model  $f$  and distribution  $\mathcal{D}$  if, for all  $\alpha \in [0, 1]$ ,

$$\text{Acc}_{\mathcal{D},f}((1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1) \geq [(1 - \alpha) \cdot \text{Acc}_{\mathcal{D},f}(\theta_0) + \alpha \cdot \text{Acc}_{\mathcal{D},f}(\theta_1)]. \quad (4)$$

Note that when  $\text{Acc}_{\mathcal{D},f}(\theta_0) = \text{Acc}_{\mathcal{D},f}(\theta_1)$ , Definition 1 is equivalent to the original definition of Frankle et al. [29].<sup>5</sup>

**Observation 1:** As illustrated in Figure 12, we observe linear mode connectivity on ImageNet and the main distribution shifts we consider (Section 2).

To assist in contextualizing Observation 1, we review related phenomena. Neural networks are nonlinear, and hence weight-space ensembles only achieve good performance in exceptional cases—ensembling two randomly initialized networks in weights-space achieves no better accuracy than a random classifier [29].

Linear mode connectivity has been observed by Frankle et al. [29], Izmailov et al. [46] when part of the training trajectory is shared, and by Neyshabur et al. [72] when two models are fine-tuned with a shared initialization. In particular, the observations of Neyshabur et al. [72] may elucidate why weight-space ensembles attain high accuracy in the setting we consider, as they suggest that fine-tuning remains in a region where solutions are connected by a linear path along which error remains low. However, instead of considering the weight-space ensemble of two fine-tuned models, we consider the weight-space ensemble of the *pre-trained* and fine-tuned models. This is only possible for a pre-trained model capable of zero-shot inference, such as CLIP.

<sup>5</sup>The original definition used  $\geq \frac{1}{2}(\text{Acc}_{\mathcal{D},f}(\theta_0) + \text{Acc}_{\mathcal{D},f}(\theta_1))$  and so Definition 1 is not a strict generalization. Although this original definition does account for some discrepancy between the accuracy of the endpoints, it is less applicable when accuracy differs substantially.

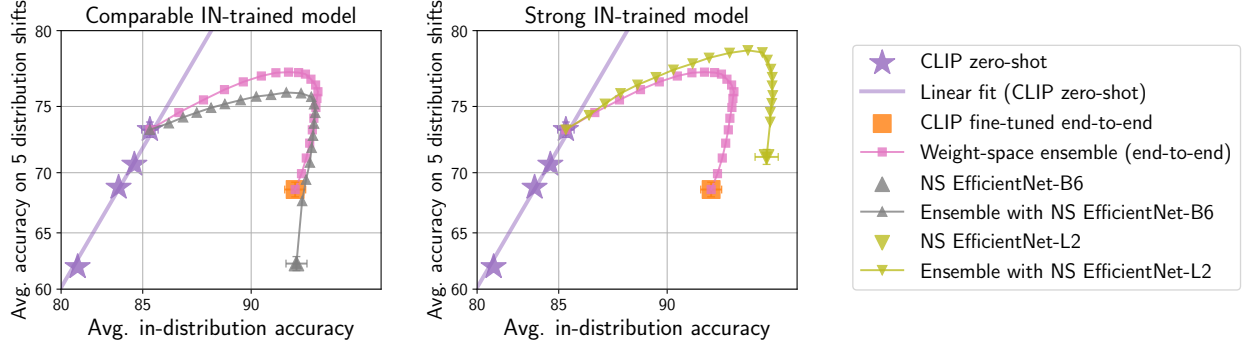


Figure 13: **Ensembling with a zero-shot model improves the out-of-distribution performance of an independently trained model.** (Left) Output-space ensembling with an independently trained model (NoisyStudent EfficientNet-B6) with comparable in-distribution performance to the end-to-end fine-tuned model. (Right) Output-space ensembling with an independently trained model with strong in-distribution performance (NoisyStudent EfficientNet-L2). Results averaged over the five distribution shifts as in Figure 1.

**Observation 2:** As illustrated by Figure 12, on ImageNet and the main distribution shifts we consider, weight-space ensembling (end-to-end) may outperform both the zero-shot and fine-tuned models. Formally, there exists an  $\alpha$  for which

$$\text{Acc}_{\mathcal{D},f}((1-\alpha) \cdot \theta_0 + \alpha \cdot \theta_1) \geq \max\{\text{Acc}_{\mathcal{D},f}(\theta_0), \text{Acc}_{\mathcal{D},f}(\theta_1)\} . \quad (5)$$

We are not the first to observe that when interpolating between models with linear mode connectivity, the accuracy of models in the center may exceed that of either endpoint [46, 72, 97]. Neyshabur et al. [72] speculate that interpolation could produce solutions closer to the true center of a basin. This intuition applied in our setting is schematized in Appendix F (Figure 36). In contrast to Neyshabur et al. [72], we interpolate between models which observe different data.

**Open questions.** Many questions remain: for instance, in this work we follow a linear path from the zero-shot to fine-tuned model in weight space, is there a better method of combining these two models? Moreover, why do we observe that weight-space ensembles outperform output-space ensembles (Figure 5)? Appendix E shows that these two methods are equivalent in certain regimes which may approximate the later stages of training [28], although their connection in the general case is unclear.

### 5.3 Ensembling zero-shot CLIP with independently trained models

So far we have shown that a zero-shot model can be used to improve out-of-distribution performance of the derived fine-tuned model. Here, we investigate whether this improvement is specific to the fine-tuned model. On the contrary, we find that ensembling with robust models can improve out-of-distribution accuracy of *independently trained models*. Note that in the general case where the models being ensembled have different architectures, we are unable to perform weight-space ensembling; instead, we ensemble the outputs of each model. This increases the computational cost of inference, in contrast to the results shown in Section 4.

Concretely, we ensemble zero-shot CLIP with two Noisy Student EfficientNet models [99, 91]: (i) EfficientNet-B6 (Figure 13, left), with in-distribution performance comparable to the end-to-end fine-tuned CLIP model; and (ii) EfficientNet-L2 (Figure 13, right), the strongest model available on PyTorch ImageNet Models [96]. In both cases, we observe substantial improvements from ensembling—13.6 pp and 6.9 pp in average out-of-domain accuracy without reducing in-distribution performance. Details provided in Appendix B.

## 6 Related work

**Robustness.** Understanding how models perform under distribution shift remains an important goal, as real world models may encounter data from new environments [2]. Previous work has studied model behavior under synthetic [38, 93, 64, 32, 27, 1] and natural distribution shift [40, 50, 95, 4, 41]. Interventions used for synthetic shifts do not typically provide robustness to natural distribution shifts [92]. On the contrary, across many distribution shifts, accuracy on the standard test set is a reliable predictor for accuracy under distribution shift [100, 67, 92, 90, 68]. CLIP zero-shot models lie on a separate trend with higher effective robustness compared to models trained in-distribution [80]. Andreassen et al. [3] show that when pre-trained models are adapted to a specific distribution through standard fine-tuning, effective robustness significantly deteriorates at convergence. Exploring better methods for fine-tuning which preserve out-of-distribution performance is the motivation for this work.

**Pre-training and transfer learning.** Pre-training on large amounts of data is a powerful technique for building high-performing machine learning systems [25, 51, 101, 79, 11]. As scale is a strong driver of performance [102, 51, 49, 42, 43], an array of data collection pipelines and pre-training objectives that alleviate the need for explicit human annotation have been proposed. Although traditional supervised learning (e.g., image classification) is an effective pre-training objective [89, 51, 25, 102], more scalable methods that do not require labeled data have attracted significant interest [24, 74, 75, 33, 35, 15, 13, 14, 69, 12]. One increasingly popular class of vision models are those pre-trained with auxiliary language supervision [21, 84, 104, 80].<sup>6</sup> Such data can be gathered from the web with relative ease, enabling the construction of large pre-training datasets [48].

Once models are pre-trained, there are various ways to adapt to downstream tasks. A common choice is fine-tuning a linear classifier, which is computationally cheap [53, 80]. An alternative that requires significantly more computational resources is end-to-end fine-tuning, where all weights of the models can be changed after pre-training. Models trained with language supervision can additionally be used in downstream tasks without the need for labeled data (i.e. zero-shot inference). Zero-shot models trained with language supervision can perform inference by constructing captions associated with each class name and predicting the class which maximizes agreement with a particular image.

**Traditional (output-space) ensembles.** Ensemble methods are a foundational technique in machine learning. Traditional ensemble methods, which we refer to as output-space ensembles, ensemble the predictions (outputs) of many classifiers [23, 5, 10, 58, 30]. Given data  $x$ , the output-space ensemble of classifiers  $f_1, \dots, f_n$  is given by  $F_n(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$ , though non-uniform weighting may also be considered (for boosting methods [30], there is no  $1/n$  factor and each  $f_i$  is a weak learner trained to correct the errors of  $F_{i-1}$ ). Typically, ensemble methods outperform individual classifiers. Moreover, output-space ensembling can provide calibrated uncertainty estimates under distribution shift [58, 76, 88]. In contrast to Lakshminarayanan et al. [58], Ovadia et al. [76], Stickland and Murray [88] we consider the ensemble of two models which have observed different data.

As traditional ensembles require a separate pass through each model to compute the output, output-space ensembles require more computational resources, though previous work has attempted to mitigate this issue (e.g. through distillation into a single model [65, 94]). Compared to an ensemble of 15 models trained on the same dataset, Mustafa et al. [71] find an OOD improvement of 0.8-1.6 pp (on ImageNetV2, ImageNet-R, ObjectNet, and ImageNet-A) by ensembling a similar number of models pre-trained on different datasets. In contrast, we see a larger accuracy improvement of 2-15 pp from only ensembling two models. Moreover, as we are ensembling in weight-space no extra compute is required compared to a single model.

**Weight-space ensembles.** Weight-space ensembles consider a common model  $f$  with different sets of parameters  $\theta_0$  and  $\theta_1$ . Outputs are then computed via  $f(x, (1 - \alpha) \cdot \theta_0 + \alpha \cdot \theta_1)$ , which is also referred to as linearly interpolating in weight-space between  $\theta_0$  and  $\theta_1$  with a coefficient  $\alpha$  [29, 63, 34, 97]. After the weights

---

<sup>6</sup>Previous literature is often inconsistent in the terminology used to describe the use of natural language supervision, ranging from unsupervised to supervised learning.

are combined, weight-space ensembles require no more computational resources than an individual model. When weight-space ensembles achieve high accuracy for  $\alpha \in [0, 1]$ ,  $\theta_0$  and  $\theta_1$  are said to exhibit linear mode connectivity [29]. Izmailov et al. [46] find linear mode connectivity among solutions which lie on the same training trajectory. These solutions may then be averaged in weight space for improved performance. Indeed, averaging the weights along a trajectory is a central method in optimization [82, 78, 73]. For instance Zhang et al. [103] propose optimizing with a set of fast and slow weights. Every  $k$  steps, these two sets of weights are averaged and a new trajectory begins (i.e. the slow weights are fixed while the fast weights may move). Here, we revisit these techniques from a distributional robustness perspective and consider the weight-space ensemble of models which have observed different data.

## 7 Conclusion

Zero-shot models pre-trained on large, heterogeneous datasets offer a promising avenue for building robust machine learning models [79, 48]. On applications where additional data is available, the performance of zero-shot models can be improved by fine-tuning. However, these improvements come at the expense of out-of-distribution robustness. We have presented WiSE-FT, a simple method for fine-tuning zero-shot models that removes the compromise between high accuracy and robustness. Across a number of datasets and models, WiSE-FT matches or improves in-distribution accuracy compared to standard fine-tuning, while substantially improving out-of-distribution performance. Although our investigation is centered around CLIP, we expect that our findings are more broadly applicable to other models and modalities [8, 79, 11]. We view WiSE-FT as a first step towards more sophisticated fine-tuning schemes and anticipate that future work will continue to leverage the robustness of zero-shot models for building more reliable neural networks.

## Acknowledgements

We thank Jesse Dodge, Samir Gadre, Ari Holtzman, Sewon Min, Nam Pho, Sarah Pratt, Alec Radford, and Rohan Taori for helpful discussions and draft feedback.

## References

- [1] Michael A Alcorn, Qi Li, Zhitao Gong, Chengfei Wang, Long Mai, Wei-Shinn Ku, and Anh Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. <https://arxiv.org/abs/1811.11553>.
- [2] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety, 2016. <https://arxiv.org/abs/1606.06565>.
- [3] Anders Andreassen, Yasaman Bahri, Behnam Neyshabur, and Rebecca Roelofs. The evolution of out-of-distribution robustness throughout fine-tuning, 2021. <https://arxiv.org/abs/2106.15831>.
- [4] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. URL <https://proceedings.neurips.cc/paper/2019/file/97af07a14cacba681feacf3012730892-Paper.pdf>.
- [5] Eric Bauer and Ron Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine learning*, 1999. <https://link.springer.com/article/10.1023/A:1007515423169>.
- [6] Sara Beery, Arushi Agarwal, Elijah Cole, and Vighnesh Birodkar. The iwildcam 2021 competition dataset. In *Conference on Computer Vision and Pattern Recognition (CVPR) FGVC8 Workshop*, 2021. <https://arxiv.org/abs/2105.03494>.

- [7] Yijun Bian and Huanhuan Chen. When does diversity help generalization in classification ensembles? *IEEE Transactions on Cybernetics*, 2021. <https://arxiv.org/abs/1910.13631>.
- [8] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models, 2021. <https://arxiv.org/abs/2108.07258>.
- [9] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *European Conference on Computer Vision (ECCV)*, 2014. [https://data.vision.ee.ethz.ch/cvl/datasets\\_extra/food-101/](https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/).
- [10] Leo Breiman. Bagging predictors. *Machine learning*, 1996. <https://link.springer.com/article/10.1007/BF00058655>.
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, et al. Language models are few-shot learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2005.14165>.
- [12] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2006.09882>.
- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning (ICML)*, 2020. <http://proceedings.mlr.press/v119/chen20j.html>.
- [14] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2006.10029>.
- [15] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning, 2020. <https://arxiv.org/abs/2003.04297>.
- [16] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. <https://arxiv.org/abs/1711.07846>.
- [17] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. <https://arxiv.org/abs/1311.3618>.
- [18] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning (ICML)*, 2019. <https://arxiv.org/abs/1902.02918>.
- [19] Alexander D’Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D Hoffman, et al. Underspecification presents challenges for credibility in modern machine learning, 2020. <https://arxiv.org/abs/2011.03395>.
- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009. <https://ieeexplore.ieee.org/document/5206848>.

- [21] Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annotations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. <https://arxiv.org/abs/2006.06666>.
- [22] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout, 2017. <https://arxiv.org/abs/1708.04552>.
- [23] Thomas G Dietterich. Ensemble methods in machine learning. In *International workshop on multiple classifier systems*, 2000. [https://link.springer.com/chapter/10.1007/3-540-45014-9\\_1](https://link.springer.com/chapter/10.1007/3-540-45014-9_1).
- [24] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *International Conference on Computer Vision (ICCV)*, 2015. <https://arxiv.org/abs/1505.05192>.
- [25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021. <https://arxiv.org/abs/2010.11929>.
- [26] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry. Exploring the landscape of spatial robustness. In *International Conference on Machine Learning (ICML)*, 2019. <https://arxiv.org/abs/1712.02779>.
- [27] Kevin Eykholt, Ivan Evtimov, Earlene Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. Robust physical-world attacks on deep learning visual classification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. <https://arxiv.org/abs/1707.08945>.
- [28] Stanislav Fort, Gintare Karolina Dziugaite, Mansheej Paul, Sepideh Kharaghani, Daniel M Roy, and Surya Ganguli. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2010.15110>.
- [29] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning (ICML)*, 2020. <https://arxiv.org/abs/1912.05671>.
- [30] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 1997. <https://www.sciencedirect.com/science/article/pii/S002200009791504X>.
- [31] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations (ICLR)*, 2018. <https://arxiv.org/abs/1811.12231>.
- [32] Robert Geirhos, Carlos R Medina Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. Generalisation in humans and deep neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. <https://arxiv.org/abs/1808.08750>.
- [33] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations (ICLR)*, 2018. <https://arxiv.org/abs/1803.07728>.



- [34] Ian J Goodfellow, Oriol Vinyals, and Andrew M Saxe. Qualitatively characterizing neural network optimization problems. In *International Conference on Learning Representations (ICLR)*, 2014. <https://arxiv.org/abs/1412.6544>.
- [35] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *International Conference on Computer Vision (ICCV)*, 2019. <https://arxiv.org/abs/1905.01235>.
- [36] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning (ICML)*, 2017. <https://arxiv.org/abs/1706.04599>.
- [37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. <https://arxiv.org/abs/1512.03385>.
- [38] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *International Conference on Learning Representations (ICLR)*, 2019. <https://arxiv.org/abs/1903.12261>.
- [39] Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. In *International Conference on Learning Representations (ICLR)*, 2020. <https://arxiv.org/abs/1912.02781>.
- [40] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *International Conference on Computer Vision (ICCV)*, 2021. <https://arxiv.org/abs/2006.16241>.
- [41] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. <https://arxiv.org/abs/1907.07174>.
- [42] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling, 2020. <https://arxiv.org/abs/2010.14701>.
- [43] Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for transfer, 2021. <https://arxiv.org/abs/2102.01293>.
- [44] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *Advances in Neural Information Processing Systems (NeurIPS) Deep Learning Workshop*, 2015. <https://arxiv.org/abs/1503.02531>.
- [45] Tin Kam Ho. The random subspace method for constructing decision forests. *IEEE transactions on pattern analysis and machine intelligence*, 1998. <https://ieeexplore.ieee.org/document/709601>.
- [46] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018. <https://arxiv.org/abs/1803.05407>.
- [47] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. <https://arxiv.org/abs/1806.07572>.
- [48] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2102.05918>.

- [49] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. <https://arxiv.org/abs/2001.08361>.
- [50] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2012.07421>.
- [51] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *European Conference on Computer Vision (ECCV)*, 2020. <https://arxiv.org/abs/1912.11370>.
- [52] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International Conference on Machine Learning (ICML)*, 2019. <https://arxiv.org/abs/1905.00414>.
- [53] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. <https://arxiv.org/abs/1805.08974>.
- [54] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *International Conference on Computer Vision (ICCV) Workshops*, 2013. <https://ieeexplore.ieee.org/document/6755945>.
- [55] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images, 2009. <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
- [56] Ludmila I Kuncheva. *Combining pattern classifiers: methods and algorithms*. John Wiley & Sons, 2014.
- [57] Ludmila I Kuncheva and Christopher J Whitaker. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine learning*, 2003. <https://doi.org/10.1023/A:1022859003006>.
- [58] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. <https://arxiv.org/abs/1612.01474>.
- [59] Chunyuan Li, Heerad Farkhor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension of objective landscapes. In *International Conference on Learning Representations (ICLR)*, 2018. <https://arxiv.org/abs/1804.08838>.
- [60] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*, 2016. <https://arxiv.org/abs/1608.03983>.
- [61] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019. <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [62] Shangyun Lu, Bradley Nott, Aaron Olson, Alberto Todeschini, Hossein Vahabi, Yair Carmon, and Ludwig Schmidt. Harder or different? a closer look at distribution shift in dataset reproduction. In *International Conference on Machine Learning (ICML) Workshop on Uncertainty and Robustness in Deep Learning*, 2020. <http://www.gatsby.ucl.ac.uk/~balaji/udl2020/accepted-papers/UDL2020-paper-101.pdf>.

- [63] James Lucas, Juhan Bae, Michael R Zhang, Stanislav Fort, Richard Zemel, and Roger Grosse. Analyzing monotonic linear interpolation in neural network loss landscapes. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2104.11044>.
- [64] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2017. <https://arxiv.org/abs/1706.06083>.
- [65] Andrey Malinin, Bruno Mlodozienec, and Mark Gales. Ensemble distribution distillation. In *International Conference on Learning Representations (ICLR)*, 2019. <https://arxiv.org/abs/1905.00076>.
- [66] Mary L McHugh. Interrater reliability: the kappa statistic. *Biochemia medica*, 2012.
- [67] John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. The effect of natural distribution shift on question answering models. In *International Conference on Machine Learning (ICML)*, 2020. <https://arxiv.org/abs/2004.14444>.
- [68] John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar, Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation between out-of-distribution and in-distribution generalization. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2107.04649>.
- [69] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. <https://arxiv.org/abs/1912.01991>.
- [70] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1906.02629>.
- [71] Basil Mustafa, Carlos Riquelme, Joan Puigcerver, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Deep ensembles for low-data transfer learning, 2020. <https://arxiv.org/abs/2010.06866>.
- [72] Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2008.11687>.
- [73] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms, 2018. <https://arxiv.org/abs/1803.02999>.
- [74] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision (ECCV)*, 2016. <https://arxiv.org/abs/1603.09246>.
- [75] Mehdi Noroozi, Hamed Pirsiavash, and Paolo Favaro. Representation learning by learning to count. In *International Conference on Computer Vision (ICCV)*, 2017. <https://arxiv.org/abs/1708.06734>.
- [76] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1906.02530>.
- [77] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1912.01703>.

- [78] Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization*, 1992. <https://epubs.siam.org/doi/abs/10.1137/0330046?journalCode=sjcodc>.
- [79] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language Models are Unsupervised Multitask Learners, 2019. <https://openai.com/blog/better-language-models/>.
- [80] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2103.00020>.
- [81] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do ImageNet classifiers generalize to ImageNet? In *International Conference on Machine Learning (ICML)*, 2019. <https://arxiv.org/abs/1902.10811>.
- [82] David Ruppert. Efficient estimations from a slowly convergent robbins-monro process, 1988. <https://ecommons.cornell.edu/handle/1813/8664>.
- [83] Hadi Salman, Greg Yang, Jerry Li, Pengchuan Zhang, Huan Zhang, Ilya Razenshteyn, and Sebastian Bubeck. Provably robust deep learning via adversarially trained smoothed classifiers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1906.04584>.
- [84] Mert Bulent Sariyildiz, Julien Perez, and Diane Larlus. Learning visual representations with caption annotations. In *European Conference on Computer Vision (ECCV)*, 2020. <https://arxiv.org/abs/2008.01392>.
- [85] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1904.12843>.
- [86] Vaishaal Shankar, Achal Dave, Rebecca Roelofs, Deva Ramanan, Benjamin Recht, and Ludwig Schmidt. Do image classifiers generalize across time?, 2019. <https://arxiv.org/abs/1906.02168>.
- [87] David B Skalak et al. The sources of increased accuracy for two proposed boosting algorithms. In *American Association for Artificial Intelligence (AAAI), Integrating Multiple Learned Models Workshop*, 1996. <https://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.40.2269&rep=rep1&type=pdf>.
- [88] Asa Cooper Stickland and Iain Murray. Diverse ensembles improve calibration. In *International Conference on Machine Learning (ICML) Workshop on Uncertainty and Robustness in Deep Learning*, 2020. <https://arxiv.org/abs/2007.04206>.
- [89] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *International Conference on Computer Vision (ICCV)*, 2017. <https://arxiv.org/abs/1707.02968>.
- [90] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. <https://arxiv.org/abs/1912.04838>.
- [91] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning (ICML)*, 2019. <https://proceedings.mlr.press/v97/tan19a/tan19a.pdf>.

- [92] Rohan Taori, Achal Dave, Vaishaal Shankar, Nicholas Carlini, Benjamin Recht, and Ludwig Schmidt. Measuring robustness to natural distribution shifts in image classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. <https://arxiv.org/abs/2007.00644>.
- [93] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. Ensemble adversarial training: Attacks and defenses. In *International Conference on Learning Representations (ICLR)*, 2017. <https://arxiv.org/abs/1705.07204>.
- [94] Linh Tran, Bastiaan S Veeling, Kevin Roth, Jakub Swiatkowski, Joshua V Dillon, Jasper Snoek, Stephan Mandt, Tim Salimans, Sebastian Nowozin, and Rodolphe Jenatton. Hydra: Preserving ensemble diversity for model distillation. In *International Conference on Machine Learning (ICML) Workshop on Uncertainty and Robustness in Deep Learning*, 2020. <https://arxiv.org/abs/2001.04694>.
- [95] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1905.13549>.
- [96] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- [97] Mitchell Wortsman, Maxwell C Horton, Carlos Guestrin, Ali Farhadi, and Mohammad Rastegari. Learning neural network subspaces. In *International Conference on Machine Learning (ICML)*, 2021. <https://arxiv.org/abs/2102.10472>.
- [98] Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio Torralba, and Aude Oliva. Sun database: Exploring a large collection of scene categories. *International Journal of Computer Vision*, 2016. <https://link.springer.com/article/10.1007/s11263-014-0748-y>.
- [99] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. <https://arxiv.org/abs/1911.04252>.
- [100] Chhavi Yadav and Léon Bottou. Cold case: The lost mnist digits. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1905.10498>.
- [101] I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. Billion-scale semi-supervised learning for image classification, 2019. <https://arxiv.org/abs/1905.00546>.
- [102] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers, 2021. <https://arxiv.org/abs/2106.04560>.
- [103] Michael R Zhang, James Lucas, Geoffrey Hinton, and Jimmy Ba. Lookahead optimizer: k steps forward, 1 step back. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. <https://arxiv.org/abs/1907.08610>.
- [104] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text, 2020. <https://arxiv.org/abs/2010.00747>.

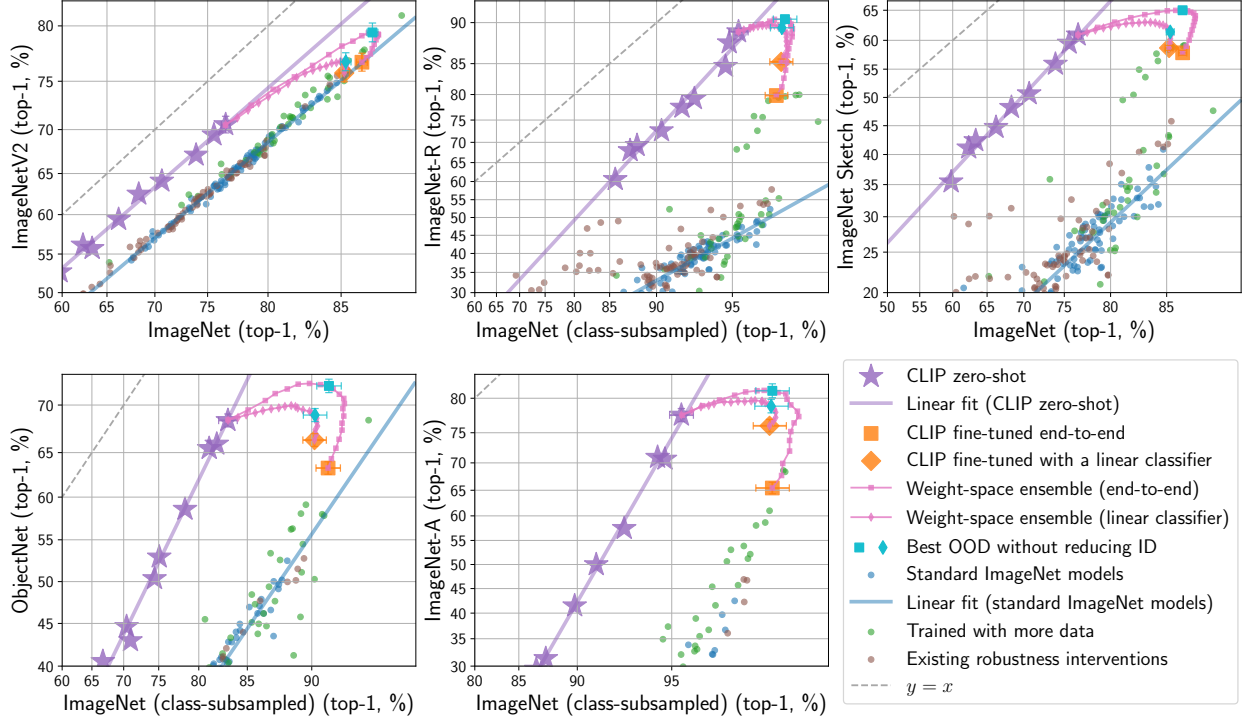


Figure 14: WiSE-FT improves in- and out-of-distribution accuracy across a number of distribution shifts. A zoomed out version of Figure 3.

## A Pseudocode for WiSE-FT

---

### Algorithm 1 Pseudocode for WiSE-FT in a PyTorch-like style

---

```
def wse(model, zeroshot_checkpoint, finetuned_checkpoint, alpha):
    # load state dicts from checkpoints
    theta_0 = torch.load(zeroshot_checkpoint)["state_dict"]
    theta_1 = torch.load(finetuned_checkpoint)["state_dict"]

    # make sure checkpoints are compatible
    assert set(theta_0.keys()) == set(theta_1.keys())

    # interpolate between all weights in the checkpoints
    theta = {
        key: (1-alpha) * theta_0[key] + alpha * theta_1[key]
        for key in theta_0.keys()
    }

    # update the model (in-place) according to the new weights
    model.load_state_dict(theta)

def wise_ft(model, dataset, zeroshot_checkpoint, alpha, hparams):
    # load the zero-shot weights
    theta_0 = torch.load(zeroshot_checkpoint)["state_dict"]
    model.load_state_dict(theta_0)

    # standard fine-tuning
    finetuned_checkpoint = finetune(model, dataset, hparams)

    # perform weight-space ensembling (in-place)
    wse(model, zeroshot_checkpoint, finetuned_checkpoint, alpha)
```

---



		OOD datasets					Avg	Avg
		IN-V2	IN-R	IN-Sketch	ObjectNet	IN-A	OOD	ID,OOD
IN (ID)								
CLIP								
End-to-end fine-tuned	86.2	76.8	79.8	57.9	63.3	65.4	68.6	77.4
WiSE-FT, $\alpha=0.75$	87.0	78.8	86.1	62.5	68.1	75.2	74.1	80.5
WiSE-FT, $\alpha=0.4$	86.2	79.2	89.9	<b>65.0</b>	71.9	80.7	77.3	81.8
WiSE-FT, optimal $\alpha$	87.1	79.5	90.3	<b>65.0</b>	72.1	81.0	77.6	82.3
NS EfficientNet-B6								
No ensemble	86.5	77.7	65.6	47.8	58.3	62.3	62.3	74.4
OSE, $\alpha=0.75$	87.0	78.8	86.4	56.7	66.5	75.9	72.9	80.0
OSE, $\alpha=0.4$	84.3	77.2	89.5	63.8	69.7	79.0	75.8	80.0
OSE, optimal $\alpha$	87.1	79.3	89.7	63.8	69.7	79.3	76.4	81.8
NS EfficientNet-L2								
No ensemble	88.3	80.8	74.6	47.6	69.8	84.7	71.5	79.9
OSE, $\alpha=0.75$	<b>88.6</b>	81.6	88.0	53.4	72.2	<b>87.1</b>	76.5	82.5
OSE, $\alpha=0.4$	85.2	78.5	<b>90.5</b>	63.9	72.6	86.0	78.3	81.8
OSE, optimal $\alpha$	<b>88.6</b>	<b>81.7</b>	<b>90.5</b>	63.9	<b>73.1</b>	<b>87.1</b>	<b>79.3</b>	<b>83.9</b>

Table 4: Accuracy of various independently trained models ensembled with CLIP on ImageNet and derived distribution shifts. OSE denotes output-space ensembling. *Avg OOD* displays the mean performance among the five out-of-distribution datasets, while *Avg ID,OOD* shows the average of ImageNet (ID) and Avg OOD.

## B Independently trained models

In Table 4, we compare performance of ensembles between zero-shot CLIP and two independently trained models, NoisyStudent EfficientNet-B6 and L2 [99, 91], on ImageNet and derived distribution shifts.

## C Diversity measures

Let  $\mathcal{S} = \{(x^{(i)}, y^{(i)}), 1 \leq i \leq N\}$  be a classification set with input data  $x^{(i)}$  and labels  $y^{(i)} \in \{1, \dots, C\}$ , where  $C$  is the number of classes. A classifier  $f$  is a function that maps inputs  $x$  to logits  $f(x) \in \mathbb{R}^C$ , yielding predictions  $\hat{y} = \arg \max_{1 \leq c \leq C} f(x)_c$ . We consider measures of diversity  $\mathcal{M}(f, g, \mathcal{S})$  between two classifiers  $f$  and  $g$  and the dataset  $\mathcal{S}$ . For simplicity,  $\hat{y}_f^{(i)}$  is used to denote the predictions from classifier  $f$  given inputs  $x^{(i)}$  (and similarly for  $g$ ).

**Prediction Diversity (PD).** One of the most intuitive ways to measure diversity between pairs of classifiers is to compute the fraction of samples where they disagree while one is correct [45, 87]. Formally, the prediction diversity PD is defined as:

$$\text{PD}(f, g, \mathcal{S}) = \frac{1}{N} \sum_{1 \leq i \leq N} \mathbb{1} \left[ \left( \hat{y}_f^{(i)} = y^{(i)} \wedge \hat{y}_g^{(i)} \neq y^{(i)} \right) \vee \left( \hat{y}_f^{(i)} \neq y^{(i)} \wedge \hat{y}_g^{(i)} = y^{(i)} \right) \right]. \quad (6)$$

**Cohen’s Kappa Complement (CC).** Cohen’s kappa coefficient is a measure of agreement between two annotators [66]. Here, we use its complement as a diversity measure between two classifiers:

$$\text{CC}(f, g, \mathcal{S}) = 1 - \frac{p_o - p_e}{1 - p_e} = \frac{1 - p_o}{1 - p_e}, \quad (7)$$

where  $p_e$  is the expected agreement between the classifiers and  $p_o$  is the empirical probability of agreement. Formally, if  $n_{f,k}$  is the number of samples where classifier  $f$  predicted label  $k$  (i.e.  $n_{f,k} = \sum_{1 \leq i \leq N} \mathbb{1}[\hat{y}_f^i = k]$ ), then:

$$p_e = \frac{1}{N^2} \sum_{1 \leq c \leq C} n_{f,c} n_{g,c}, \quad p_o = \frac{1}{N} \sum_{1 \leq i \leq N} \mathbb{1}[\hat{y}_f^i = \hat{y}_g^i] \quad (8)$$

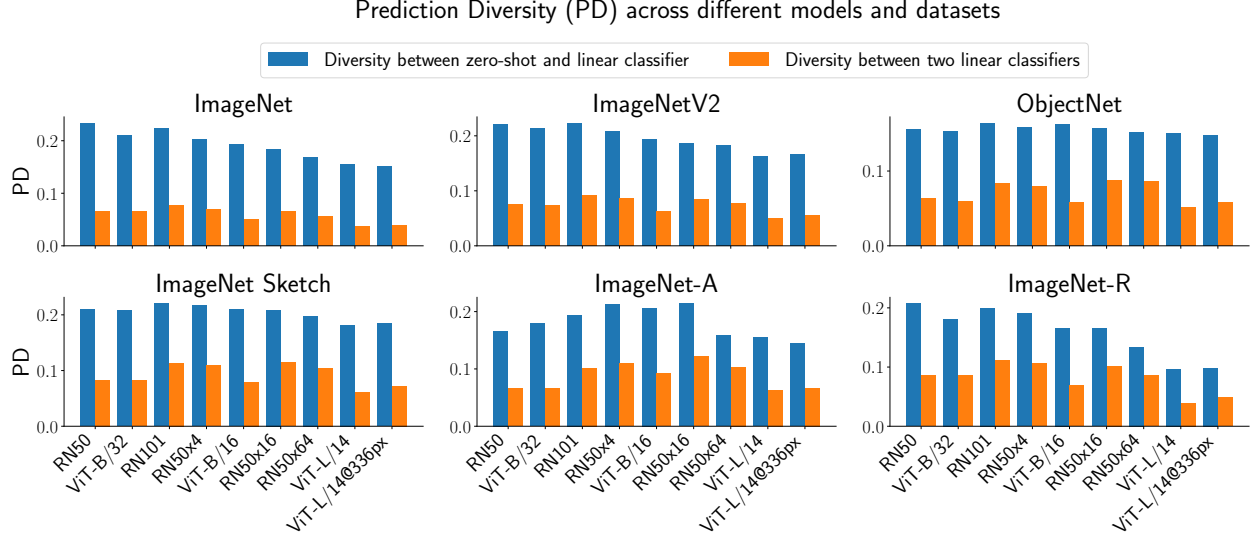


Figure 15: **Prediction Diversity (PD)** for multiple datasets and CLIP models (Equation 6).

**KL Divergence (KL).** The Kullback-Leibler divergence measures how different a probability distribution is from another. Let  $p_f^{(i)} = \text{softmax}(f(x^{(i)}))$  for a classifier  $f$ , and let  $p_{f,c}^{(i)}$  be the probability assigned to class  $c$ . We consider the average KL-divergence over all samples as a diversity measure:

$$\text{KL}(f, g, \mathcal{S}) = \frac{1}{N} \sum_{1 \leq i \leq N} \sum_{1 \leq c \leq C} p_{f,c}^{(i)} \log \left( \frac{p_{f,c}^{(i)}}{p_{g,c}^{(i)}} \right). \quad (9)$$

**Centered Kernel Alignment Complement (CKAC).** CKA is a similarity measure that compares two different sets of high-dimensional representations [52]. It is commonly used for comparing representations of two neural networks, or determining correspondences between two hidden layers of the same network. CKA measures the agreement between two matrices containing the pair-wise similarities of all samples in a dataset, where each matrix is constructed according to the representations of a model. More formally, let  $S \in \mathbb{R}^{N \times d}$  denote the  $d$ -dimensional features for all samples in a dataset  $\mathcal{S}$ , pre-processed to center the columns. For two models  $f$  and  $g$  yielding similarity matrices  $S_f$  and  $S_g$ , CKA is defined as:

$$\text{CKA}(f, g, \mathcal{S}) = \frac{\|S_g^\top S_f\|_F^2}{\|S_f^\top S_f\|_F \|S_g^\top S_g\|_F}, \quad (10)$$

where  $\|S\|_F$  denotes the Frobenius norm of the matrix  $S$ . Larger CKA values indicate larger similarities between the representations of the two models, and thus, smaller diversity. We define the diversity measure CKAC as:

$$\text{CKAC} = 1 - \text{CKA}. \quad (11)$$

Note that CKAC is computationally expensive to compute for large datasets. For this reason, in our experiments with distributions larger than 10,000 samples, we randomly sample 10,000 to compute this measure.

**Diversity across different architectures** We extend Figure 10 to show results for all combinations of diversity measures, datasets, and CLIP models. Similarly to before, the baselines compares models with the same encoder, with two linear classifiers trained on different subsets of ImageNet with half of the data. Results are shown in Figures 15-18.

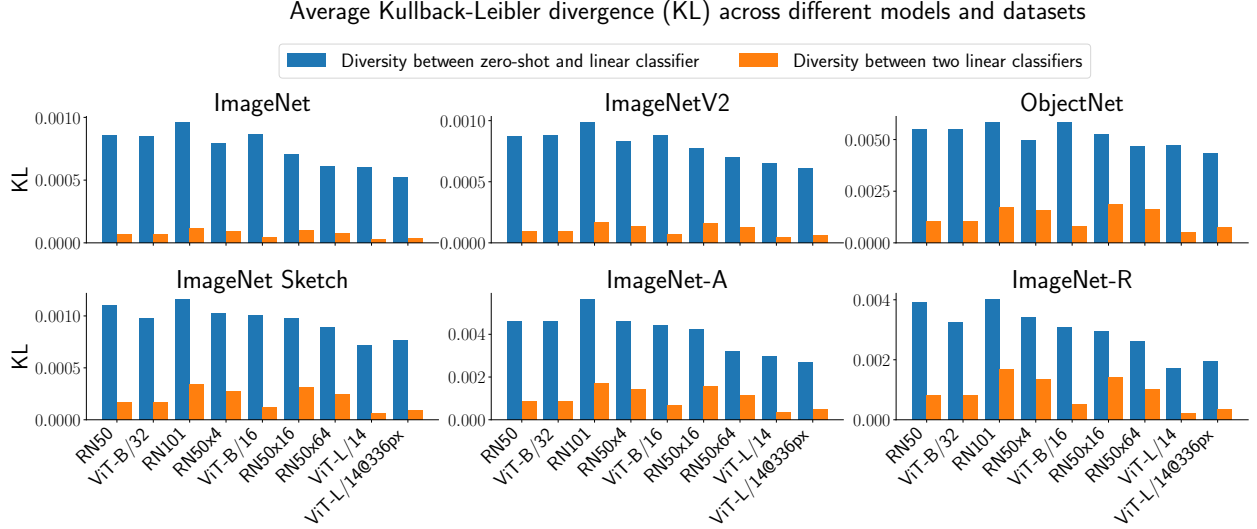


Figure 16: Cohen’s Kappa Complement (CC) for multiple datasets and CLIP models (Equation 7).

## D Experimental details

### D.1 CLIP Zero-shot

This section extends Section 2 with more details on inference with the CLIP zero-shot model. First, in all settings we use the CLIP model ViT-L/14@336px except for: (i) on CIFAR-10 and CIFAR-100 we find that ViT-L/14 performs slightly better than ViT-L/14@336px. (ii) In Figures 5 and 22 where we use ViT-B/16. Second, CLIP learns a temperature parameter which is factored into the learned weight matrix  $\mathbf{W}_{\text{zero-shot}}$  described in Section 2. Finally, to construct  $\mathbf{W}_{\text{zero-shot}}$  we ensemble the 80 prompts provided by CLIP at <https://github.com/openai/CLIP>. However, we manually engineer prompts for five datasets: WILDS-FMoW, WILDS-iWildCam, Stanford Cars, Describable Textures and Food-101. Exact prompts will be released along with code.

### D.2 Fine-tuning a linear classifier

This section extends the description of linear classifier training from Section 4.2 with details on hyperparameters and additional analyses. In each of the four regularization strategies—no regularization, weight decay, L1 regularization, and label smoothing—we run 64 hyperparameter configurations. For each trial, mini-batch size is drawn uniformly from  $\{64, 128, 256\}$  and learning rate is set to  $10^{-\beta}$  with  $\beta$  chosen uniformly at random from the range  $[0, 4]$ . Hyperparameters for each regularization strategy are as follows: (i) The weight decay coefficient is set to  $10^{-\lambda}$  where  $\lambda$  is chosen uniformly at random from  $[0, 4]$  for each trial; (ii) The L1 regularization coefficient is set to  $10^{-\lambda}$  where  $\lambda$  is chosen uniformly at random from  $[4, 8]$  for each trial; (iii) The label smoothing [70] coefficient  $\lambda$  is chosen uniformly at random from  $[0, 0.25]$  for each trial. The linear classifier used for ensembling attains the best performance in-distribution. The hyperparameters from this trial are then used in the distillation and regularization experiments described in Section 4.2. In the low-data regime (Section 4.3), this process is repeated for each  $k$  and dataset.

Figure 19 demonstrates that various regularization strategies largely move along the same parabolic trend—even linear classifiers trained without explicit regularization. We expect that for these classifiers the role of explicit regularization is replaced by the regularization implicit in stochastic mini-batch optimization. Figure 20 demonstrates that increasing learning rate moves monotonically along this trend.

When training linear classifiers with  $k$  images per class as in Section 4.3 the number of epochs is scaled

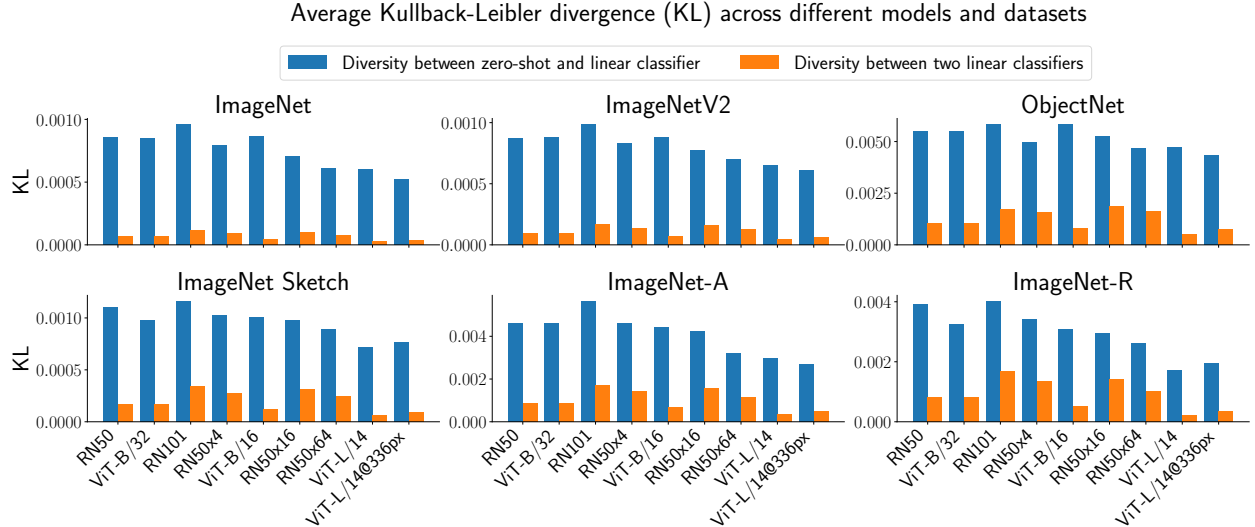


Figure 17: **Average KL Divergence (KL)** for multiple datasets and CLIP models (Equation 9).

approximately inversely proportional to the amount of data removed (e.g., with half the data we train for twice as many epochs so the number of iterations is consistent). To choose the number of epochs we use default PyTorch AdamW hyperparameters (learning rate 0.001, weight decay 0.01) and double the number of epochs until performance saturates.

### D.3 ObjectNet

The zero-shot models in Table 2 use the ImageNet class names instead of the ObjectNet class names. However, this *adaptation to class shift* improves performance by 2.3% [80]. Out of the 5 datasets used for the majority of the experiments in Section 3, ObjectNet is the only dataset for which this is possible. In Figure 21 we compare weight-space ensembles with and without adaptation to class shift.

### D.4 End-to-end fine-tuning

Fine-tuning end-to-end uses the AdamW optimizer [61, 77] with a batch size of 256,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ , weight decay of 0.1 and a cosine-annealing learning rate schedule [60] with 500 warm-up steps. We use a learning rate of  $2 \times 10^{-5}$ , gradient clipping at global norm 1 and train for a total of 10 epochs. We use the same data augmentations as [80], randomly cropping a square from resized images with the largest dimension being 336 pixels for ViT-L/14@336px and 224 for the remaining models.

## E When do weight-space approximate output-space ensembles?

In practice we observe a difference between weight-space and output-space ensembling. However, it is worth noting that these two methods of ensembling are not as different as they initially appear. In certain regimes a weight-space ensemble approximates the corresponding output-space ensemble—for instance, when training is well approximated by a linear expansion, referred to as the NTK regime [47]. Fort et al. [28] find that a linear expansion becomes more accurate in the later phase of neural network training, a phase which closely resembles fine-tuning.

Consider the set  $\Theta = \{(1 - \alpha)\theta_0 + \alpha\theta_1 : \alpha \in [0, 1]\}$  consisting of all  $\theta$  which lie on the linear path between  $\theta_0$  and  $\theta_1$ .

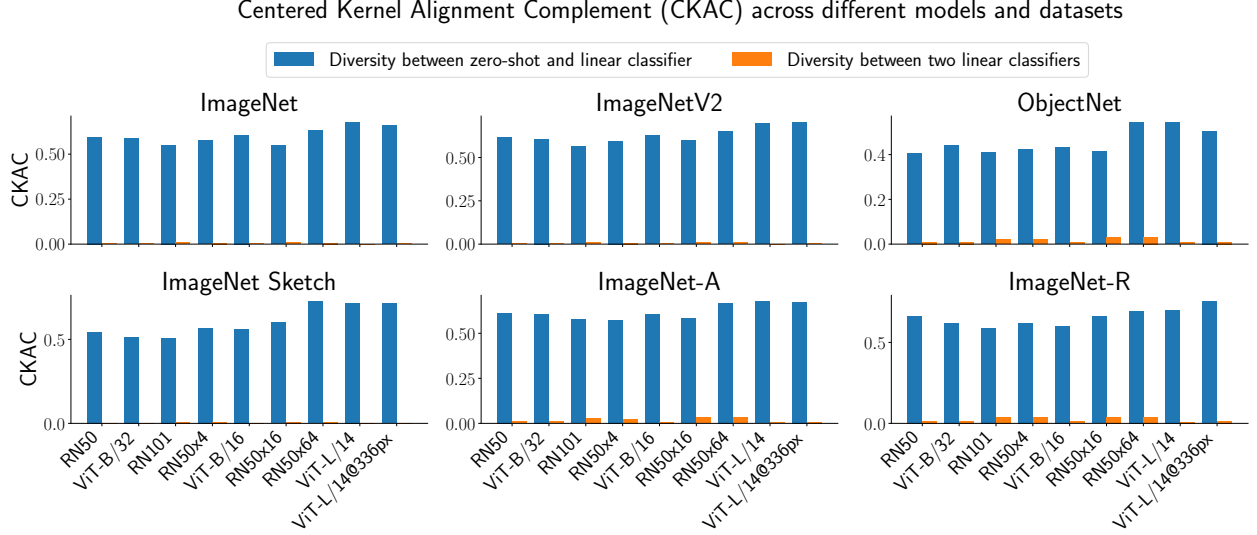


Figure 18: **Central Kernel Alignment Complement (CKAC)** for multiple datasets and CLIP models (Equation 11).

**Proposition 1.** When  $f(\theta) = f(\theta_0) + \nabla f(\theta_0)^\top (\theta - \theta_0)$  for all  $\theta \in \Theta$ , the weight- and output-space ensemble of  $\theta_0$  and  $\theta_1$  are equivalent.

*Proof.* We may begin with the weight-space ensemble and retrieve the output-space ensemble

$$f((1 - \alpha)\theta_0 + \alpha\theta_1) \tag{12}$$

$$= f(\theta_0) + \nabla f(\theta_0)^\top ((1 - \alpha)\theta_0 + \alpha\theta_1 - \theta_0) \tag{13}$$

$$= f(\theta_0) + \alpha \nabla f(\theta_0)^\top (\theta_1 - \theta_0) \tag{14}$$

$$= f(\theta_0) + \alpha \nabla f(\theta_0)^\top (\theta_1 - \theta_0) + \alpha f(\theta_0) - \alpha f(\theta_0) \tag{15}$$

$$= (1 - \alpha)f(\theta_0) + \alpha \left( f(\theta_0) + \nabla f(\theta_0)^\top (\theta_1 - \theta_0) \right) \tag{16}$$

$$= (1 - \alpha)f(\theta_0) + \alpha f(\theta_1) \tag{17}$$

where the first and final line follow by the linearity assumption.  $\square$

## F Additional Figures

This section provides supplemental figures:

- Figure 14 provides a zoomed out version of Figure 3.
- Figure 22 recreates our main result (Figure 3) with the smaller CLIP ViT-B/16 model.
- We compare weight-space ensembling to a series of alternatives as in Section 4.2 and Figure 4. However, instead of displaying average in- and out-of-distribution we show the comparison separately for each dataset (Figures 23, 24, 25, 26, and 27).
- Figures 28 and 29 show a breakdown of in-distribution gains in the low-data regime for the seven datasets averaged in Figure 7.
- We show samples from the datasets studied in our main experiments in Figures 30-35, along with the predictions of the zero-shot, linear classifier, and weight-space ensemble.
- Figure 36 provides a schematic for the average error landscape.

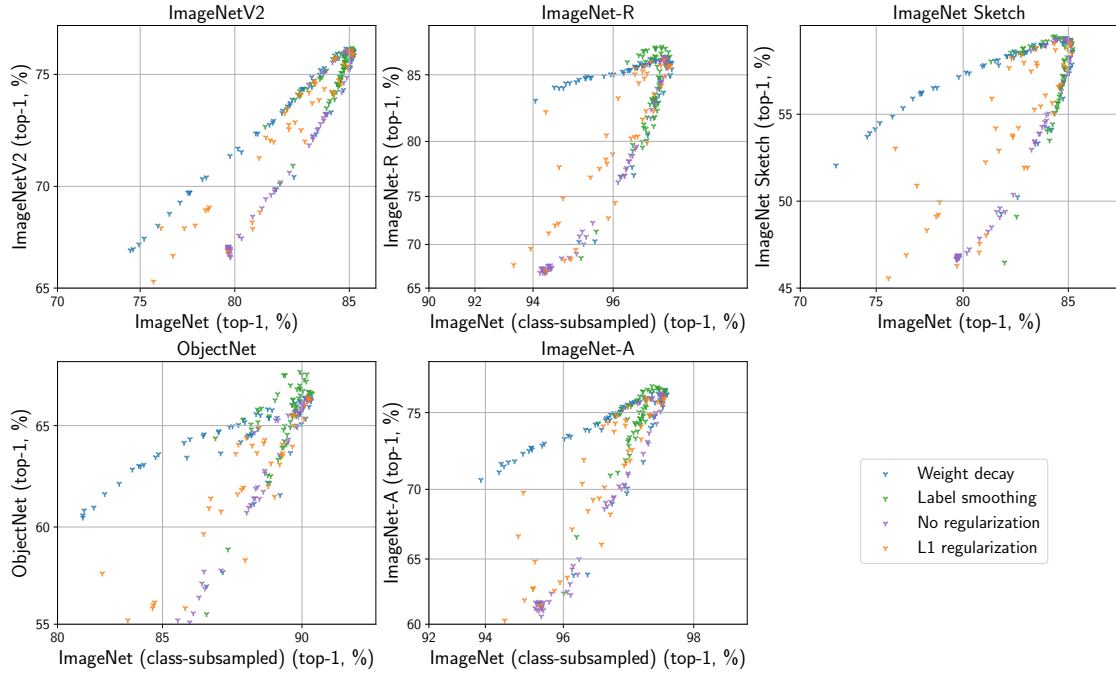


Figure 19: **Various regularizers trace similar trends when fine-tuning a linear classifier.** Comparing the relative in- and out-of-distribution performance of fine-tuning a linear classifier with the various methods of regularization discussed in Section D.2.



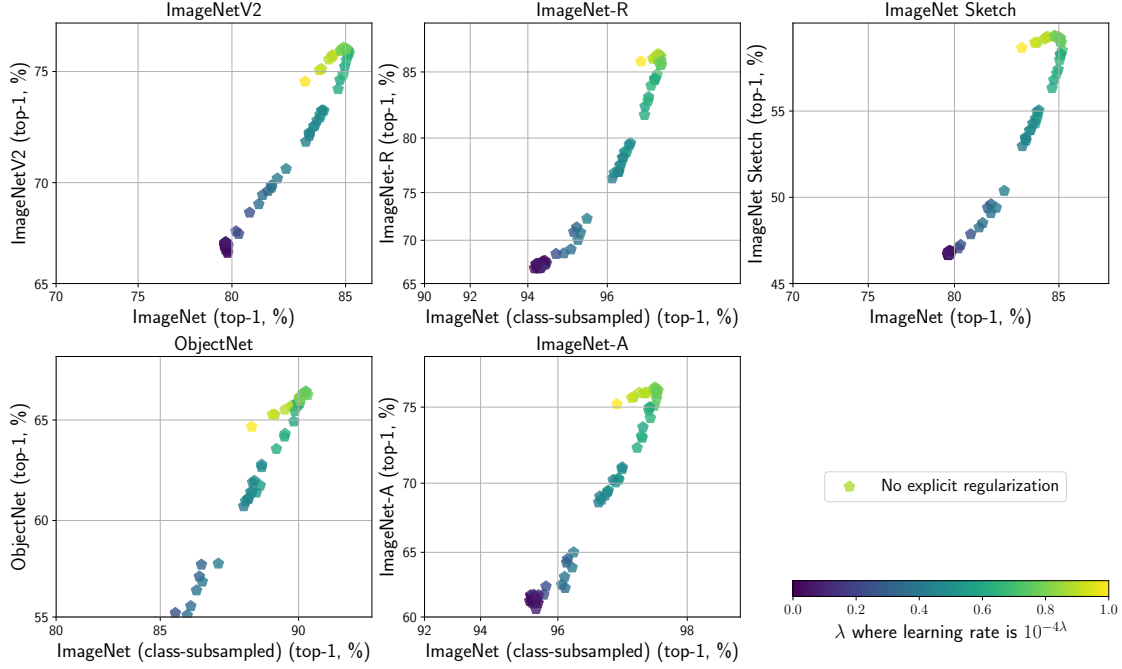


Figure 20: Comparing the relative in- and out-of-distribution performance of fine-tuning a linear classifier with various learning rates and no explicit regularization. As discussed in Section D.2, batch size is chosen randomly from  $\{64, 128, 256\}$  for each experiment. As learning rate increases the linear classifiers follow a parabolic trend similar to the trend followed by explicit regularization (see Figure 19).

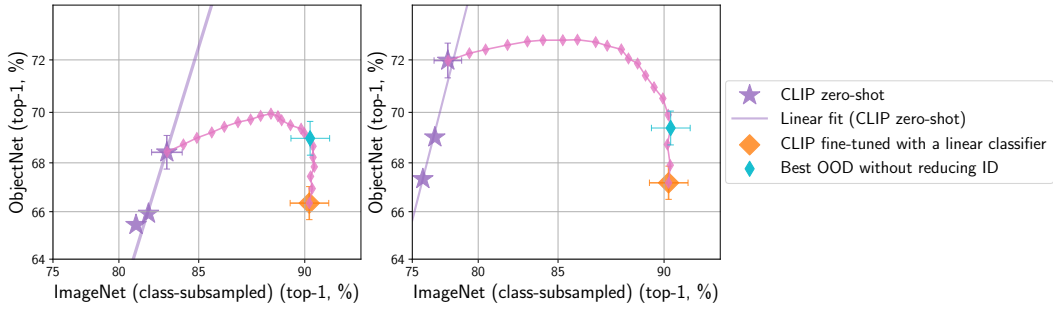


Figure 21: Effective robustness scatter plots for ObjectNet, with and without adapting to class shift. **Left:** Using ImageNet class names to construct the zero-shot classifier. **Right:** Using ObjectNet class names to construct the zero-shot classifier.

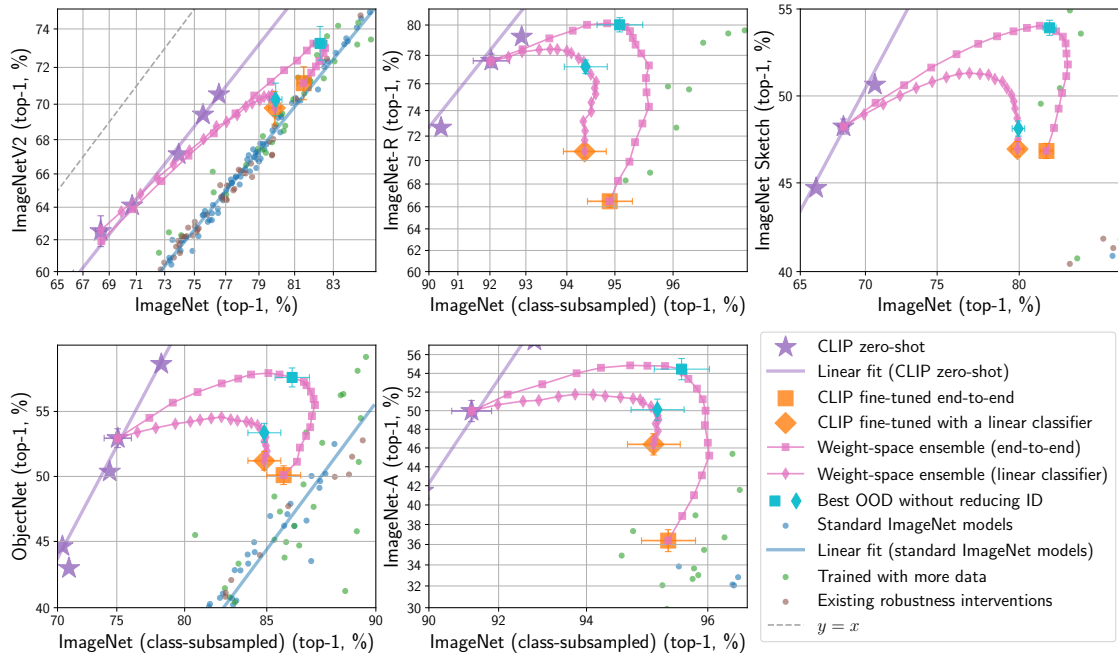


Figure 22: WiSE-FT improves in- and out-of-distribution accuracy across a number of distribution shifts. A version of Figure 14 with a smaller CLIP ViT-B/16 model.

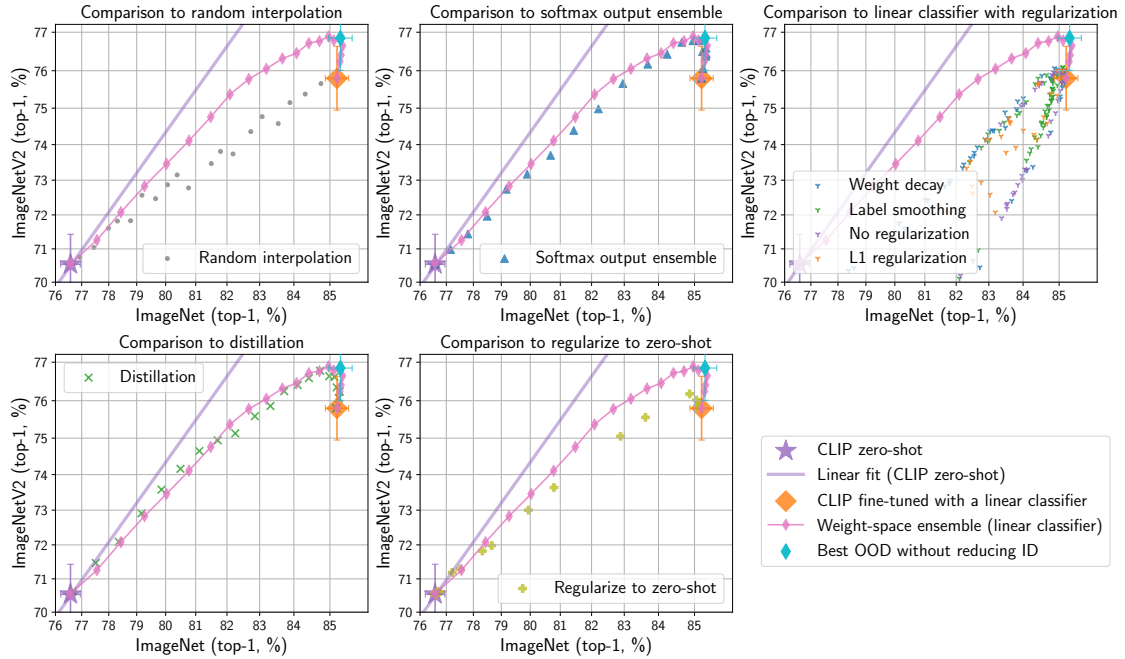


Figure 23: Comparing the relative in- and out-of-distribution performance of WiSE-FT with the alternatives described in Section 4.2 for ImageNetV2.

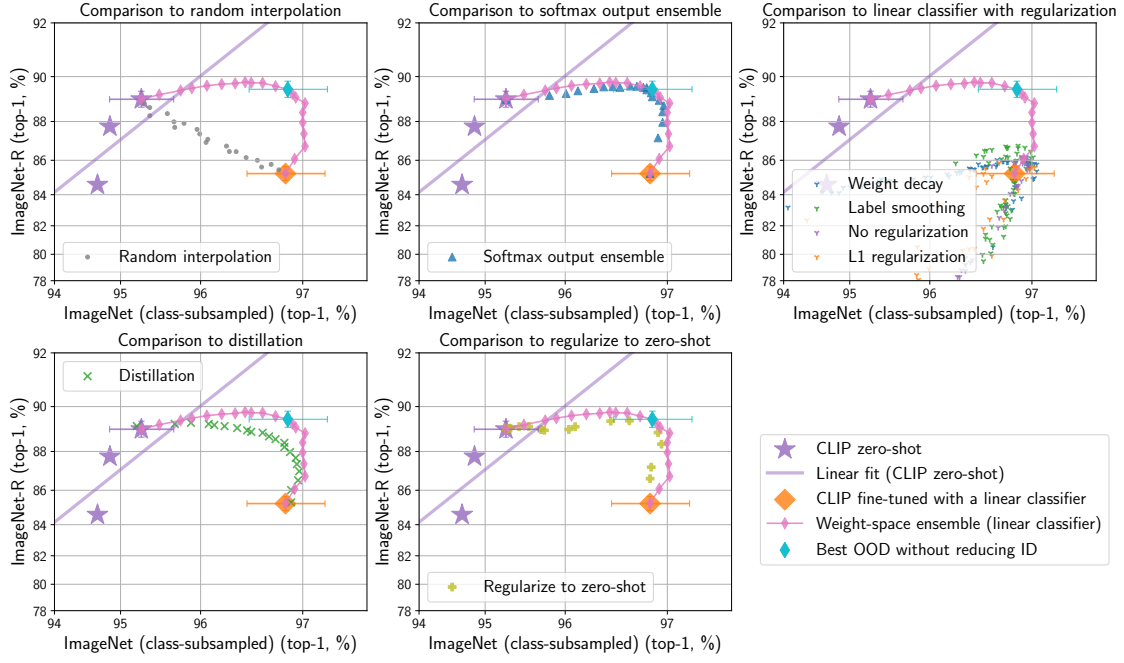


Figure 24: Comparing the relative in- and out-of-distribution performance of WiSE-FT with the alternatives described in Section 4.2 for ImageNet-R.

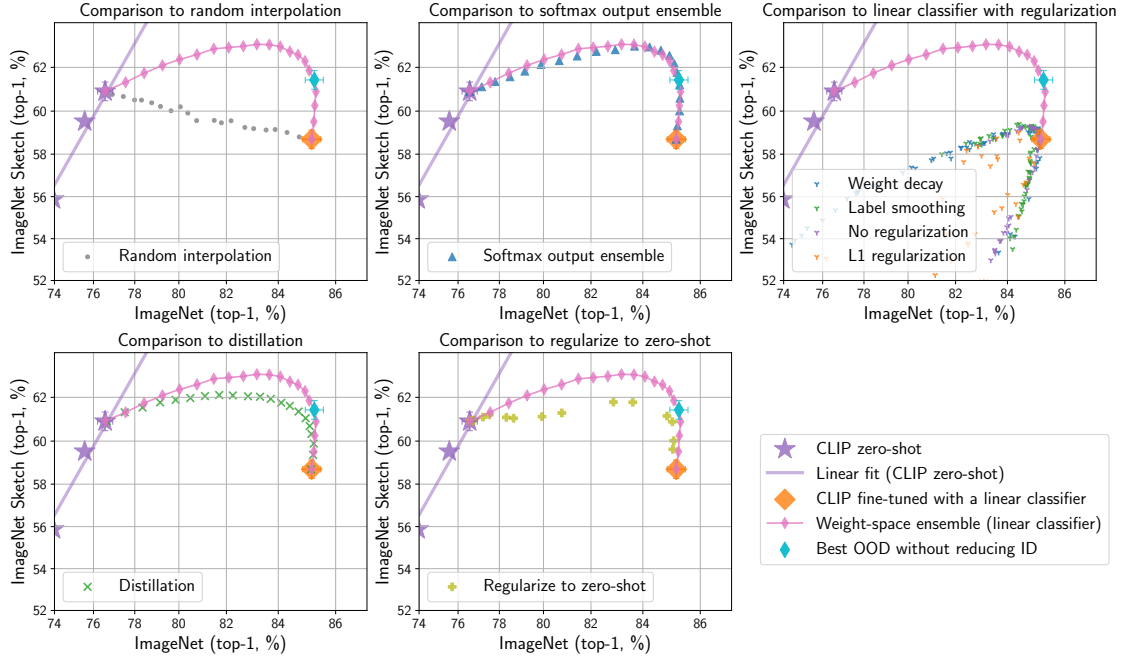


Figure 25: Comparing the relative in- and out-of-distribution performance of WiSE-FT with the alternatives described in Section 4.2 for ImageNet Sketch.

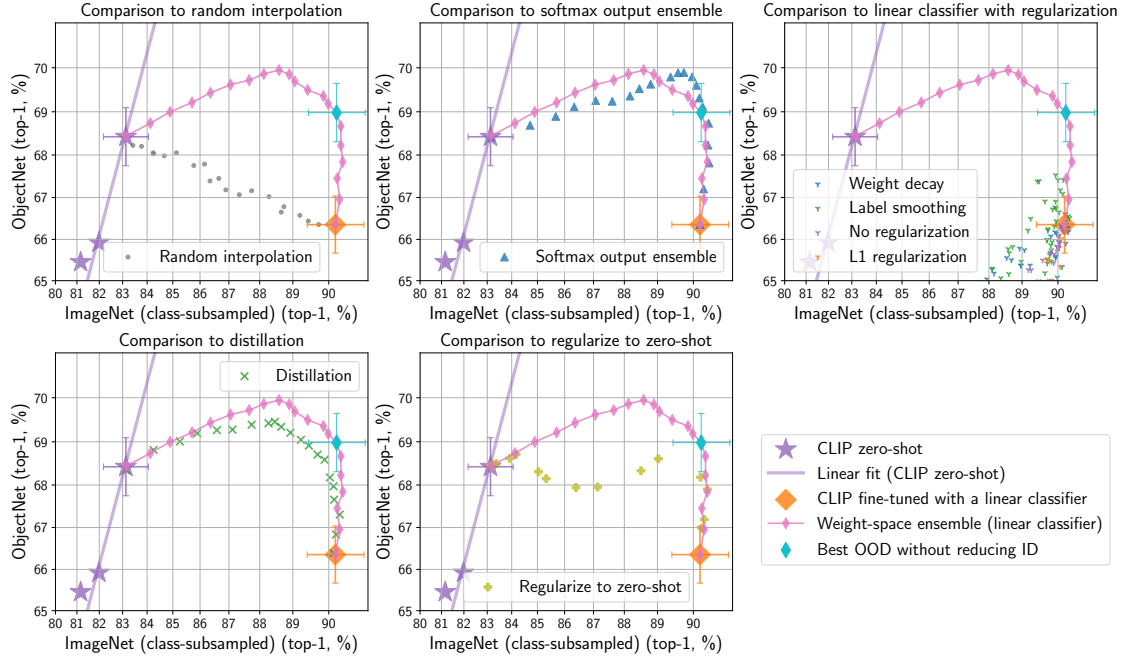


Figure 26: Comparing the relative in- and out-of-distribution performance of WiSE-FT with the alternatives described in Section 4.2 for ObjectNet.

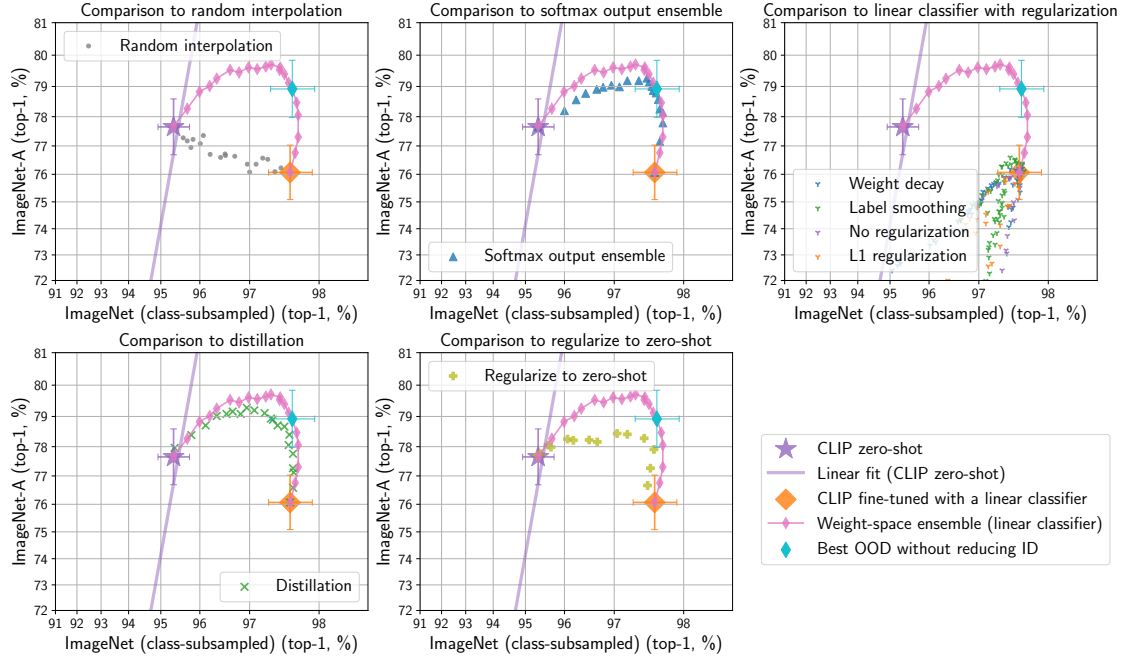


Figure 27: Comparing the relative in- and out-of-distribution performance of WiSE-FT with the alternatives described in Section 4.2 for ImageNet-A.

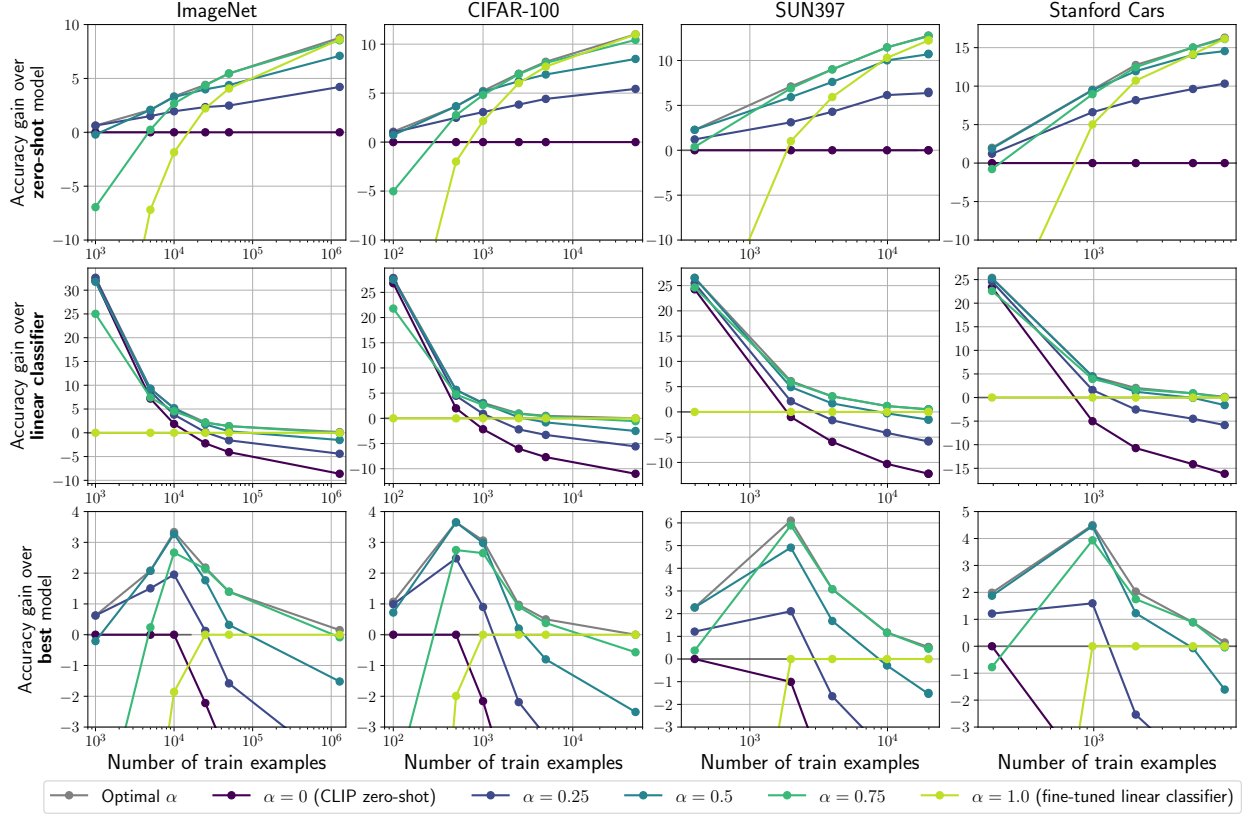


Figure 28: WiSE-FT improves in-distribution accuracy over the linear classifier and zero-shot model in the low data regime. On the  $x$ -axis we consider  $k = \{1, 5, 10, 25, 50\}$  examples per class and the full training set. On the  $y$ -axis we consider the in-distribution accuracy improvement of WiSE-FT over the **(top)** zero-shot model, **(middle)** fine-tuned linear classifier, and **(bottom)** best of the zero-shot and fine-tuned linear classifier.

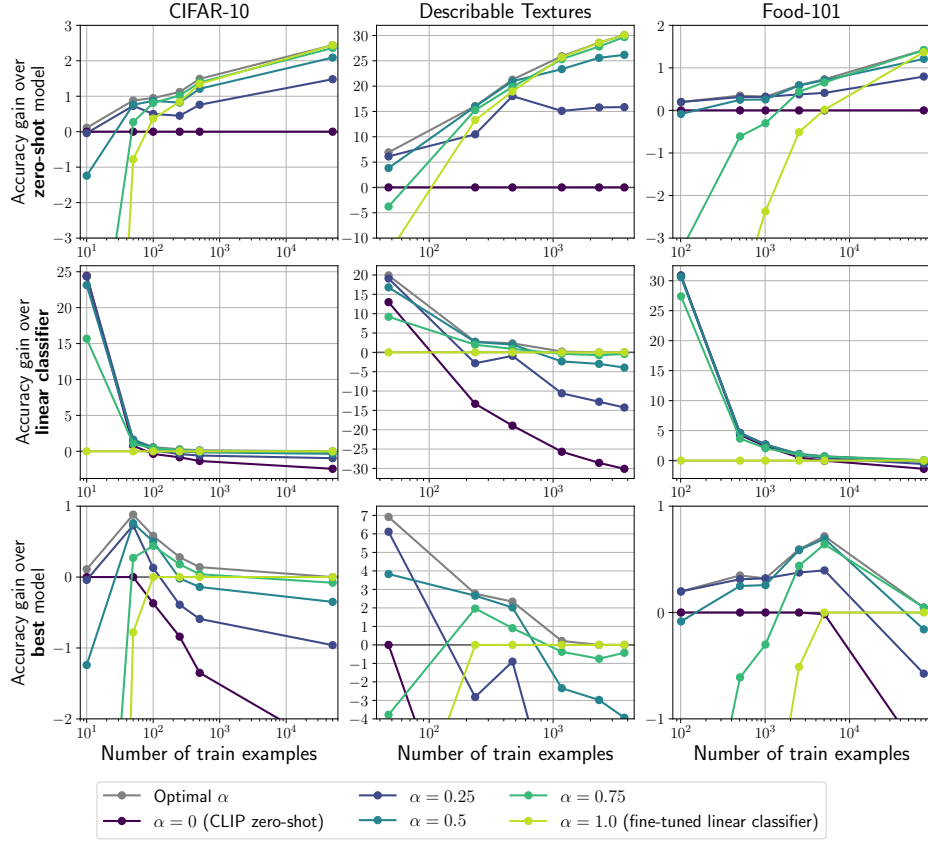


Figure 29: WiSE-FT improves in-distribution accuracy over the linear classifier and zero-shot model in the low data regime. On the  $x$ -axis we consider  $k = \{1, 5, 10, 25, 50\}$  examples per class and the full training set. On the  $y$ -axis we consider the in-distribution accuracy improvement of WiSE-FT over the **(top)** zero-shot model, **(middle)** fine-tuned linear classifier, and **(bottom)** best of the zero-shot and fine-tuned linear classifier.



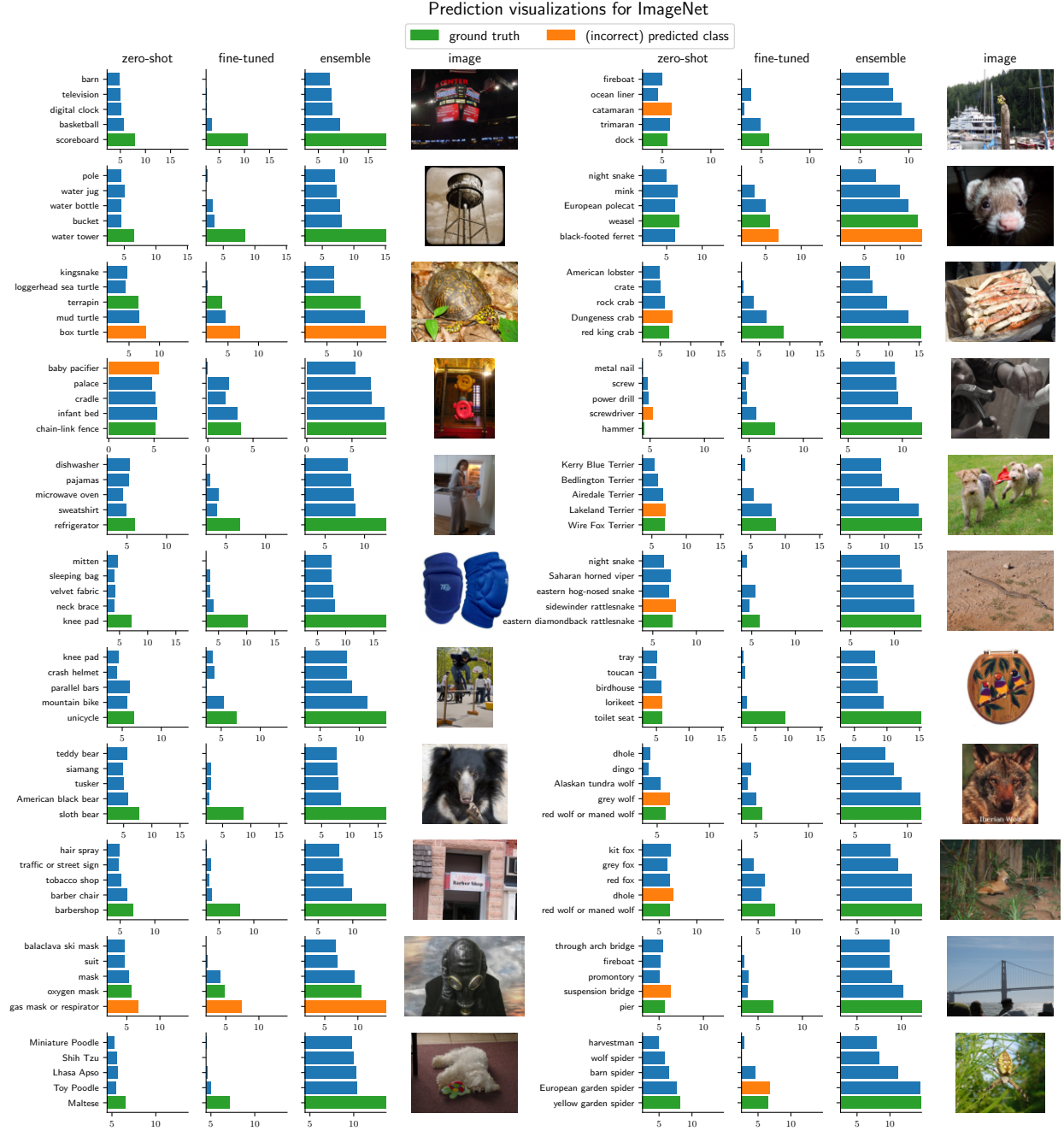


Figure 30: Visualization of predictions on **ImageNet** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.

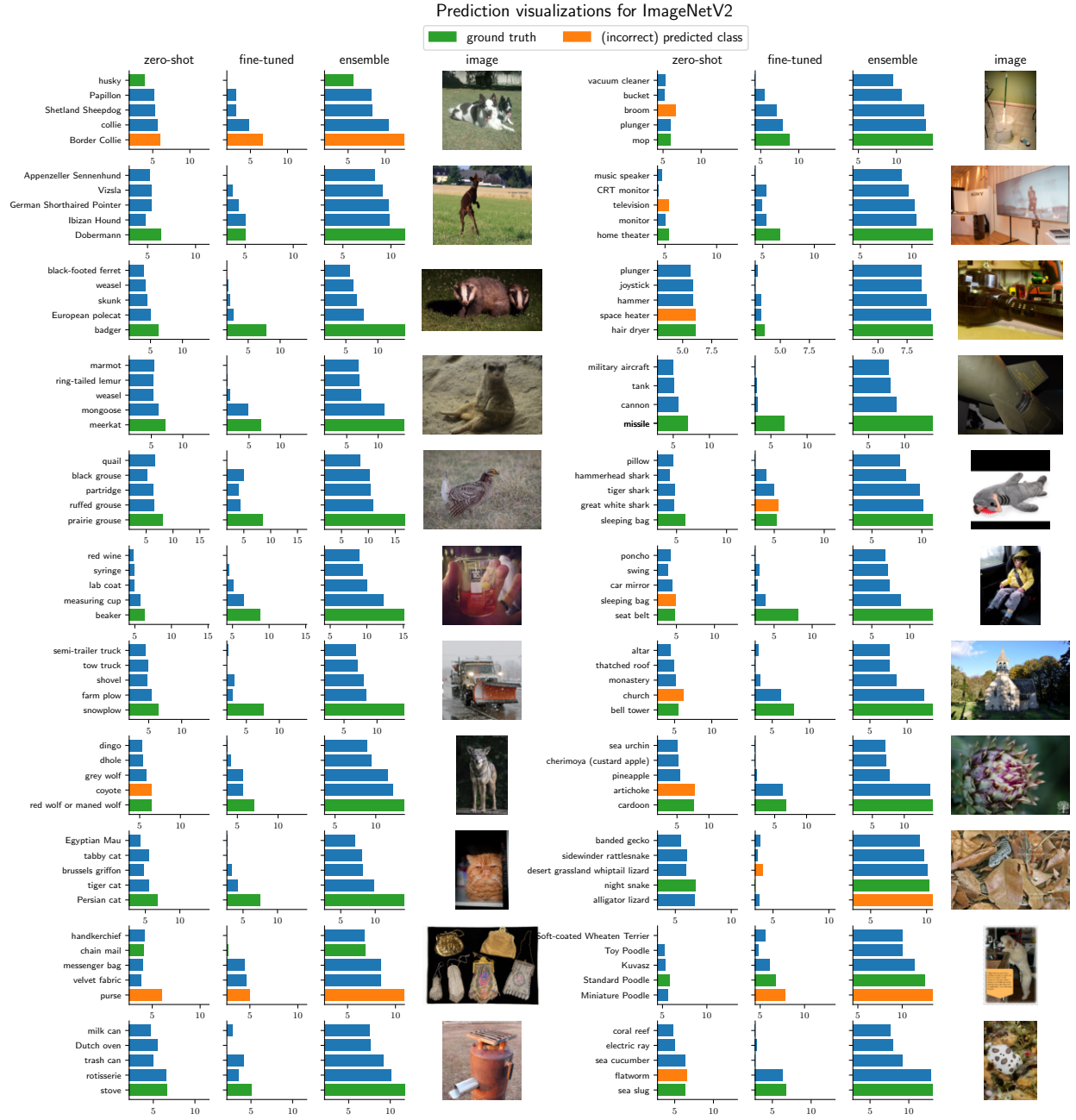


Figure 31: Visualization of predictions on **ImageNet V2** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.

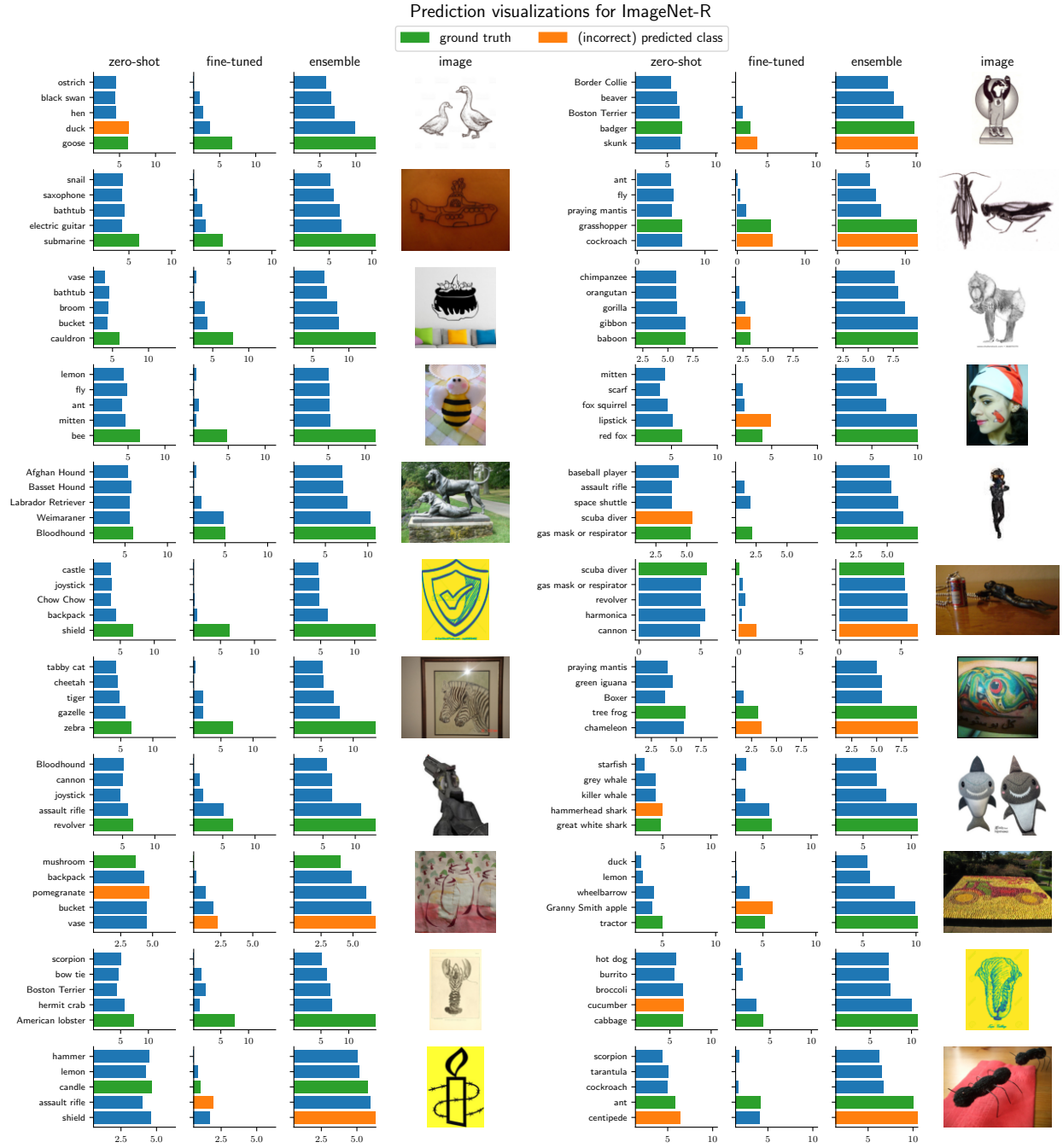


Figure 32: Visualization of predictions on **ImageNet-R** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.

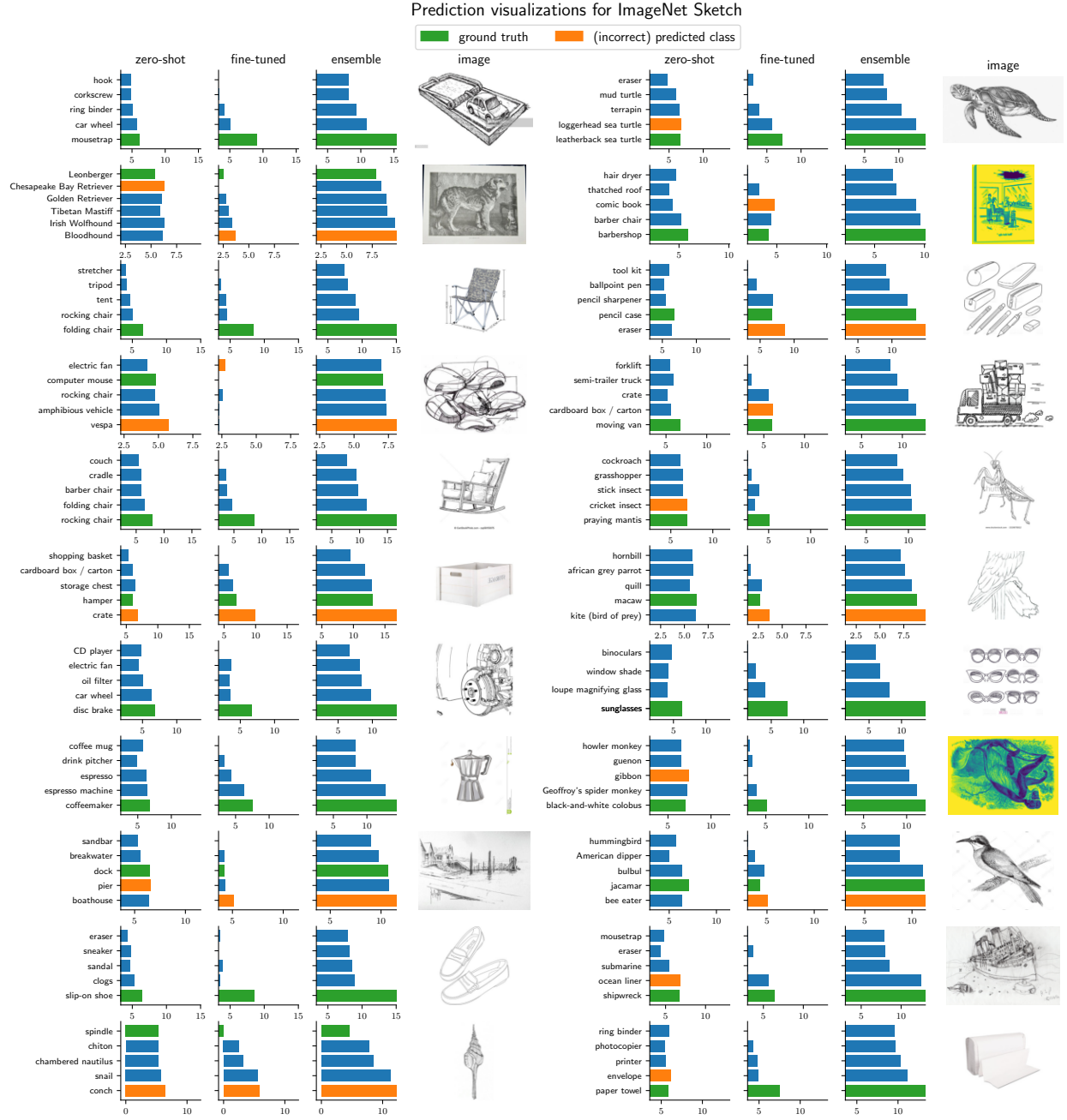


Figure 33: Visualization of predictions on **ImageNet Sketch** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.



Figure 34: Visualization of predictions on **ObjectNet** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.



Figure 35: Visualization of predictions on **ImageNet-A** for the zero-shot, fine-tuned linear classifier and weight-space ensemble with  $\alpha = 0.75$  for random samples (left) and random samples where one of the classifiers is correct (right). For the zero-shot and linear classifiers, we display logits scaled by  $(1 - \alpha)$  and  $\alpha$ , respectively.



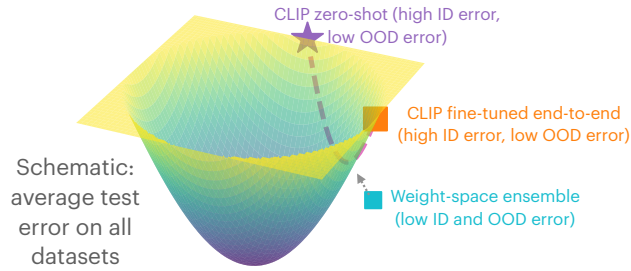


Figure 36: Schematic of the average error landscape. Li et al. [59] observe that solution spaces for a given task are high dimensional, while D’Amour et al. [19], Wortsman et al. [97] observe that movement within the solution space can change model performance on other data distributions. According to these observations, it is possible that the model finds a solution that performs well on the downstream task during fine-tuning without leaving a region of low error on the original task. Moreover, interpolating between the two solutions may travel closer to a true minimum [72, 46].