

# On the Importance of Distractors for Few-Shot Classification

Rajshekhar Das<sup>1</sup>

<sup>1</sup>Carnegie Mellon University

rajshek@andrew.cmu.edu

Yu-Xiong Wang<sup>2</sup>

<sup>2</sup>University of Illinois at Urbana-Champaign

yxw@illinois.edu

José M.F. Moura<sup>1</sup>

moura@andrew.cmu.edu

## Abstract

*Few-shot classification aims at classifying categories of a novel task by learning from just a few (typically, 1 to 5) labelled examples. An effective approach to few-shot classification involves a prior model trained on a large-sample base domain, which is then finetuned over the novel few-shot task to yield generalizable representations. However, task-specific finetuning is prone to overfitting due to the lack of enough training examples. To alleviate this issue, we propose a new finetuning approach based on contrastive learning that reuses unlabelled examples from the base domain in the form of distractors. Unlike the nature of unlabelled data used in prior works, distractors belong to classes that do not overlap with the novel categories. We demonstrate for the first time that inclusion of such distractors can significantly boost few-shot generalization. Our technical novelty includes a stochastic pairing of examples sharing the same category in the few-shot task and a weighting term that controls the relative influence of task-specific negatives and distractors. An important aspect of our finetuning objective is that it is agnostic to distractor labels and hence applicable to various base domain settings. Compared to state-of-the-art approaches, our method shows accuracy gains of up to 12% in cross-domain and up to 5% in unsupervised prior-learning settings. Our code is available at <https://github.com/quantacode/Contrastive-Finetuning.git>*

## 1. Introduction

The ability to learn from very few examples is innate to human intelligence. In contrast, large amounts of labelled examples are required by modern machine learning algorithms to learn a new task. This limits their applicability to domains where data is either expensive to annotate and collect or simply inaccessible due to privacy concerns. To overcome this limitation, few-shot classification has been proposed as a generic framework for learning to classify with very limited supervision [13, 33, 36, 61]. Under this paradigm, most approaches leverage prior knowledge from

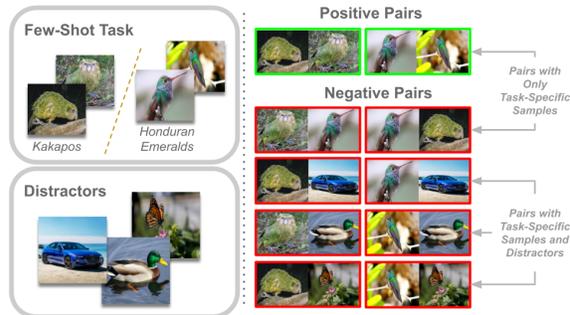


Figure 1: **Classification of Kakapos vs. Honduran Emeralds with just few examples per class and many distractors:** The idea is to leverage unlabelled data in the form of *distractors* that need not be semantically related to the classes in the few-shot task. The hope is that by pairing distractors and task samples as negatives (bottom six red boxes) and encouraging greater dissimilarity between such pairs, image representations of the two classes, Kakapos and Honduran Emeralds, will be pushed farther away. This would ultimately lead to better classification.

a (labelled) *base* domain to solve a novel task by either finetuning-based transfer [10, 66] or meta-learning [13, 17, 50, 59, 61, 64, 73]. In particular, when the base and novel domains are related, the hope is that representations learnt in the base domain can be generalized to novel tasks, thus facilitating *positive* knowledge transfer.

While the above paradigm is effective for tasks that can leverage large datasets like ImageNet [54] as the related base domain, for others, such as rare species classification [72] or medical image classification [74], acquiring necessary prior knowledge can be exceedingly difficult due to the absence of a related base domain with labelled data. To relax such data requirements, recent techniques explore alternative ways such as unsupervised learning [26, 30] or cross-domain learning [1, 17, 46, 68] to obtain representations useful for novel tasks. In the absence of labelled base data, approaches like [26, 29, 30] seek to benefit from self-supervised representation learning over unlabelled data in a related domain. In a more challenging scenario where re-

lated base data is hard to obtain, cross-domain techniques [11, 66, 68] exploit representations learnt in other domains that do not have the same task characteristics as the novel tasks.

Although the issue of learning a good prior representation remains a core focus in few-shot classification, it addresses only a part of the problem. In this work, we investigate the other important aspect, *i.e.*, *effective finetuning specific to the novel task*. Our main motivation comes from recent findings [1, 10, 18] that demonstrate the outperformance of simple finetuning over more sophisticated prior learning techniques such as meta-learning. Despite its effectiveness, we suspect that finetuning might still suffer from overfitting as a consequence of small training set in a few-shot task. To alleviate this situation, we propose to *leverage additional unlabelled data exclusive to the task*. Such datapoints are referred to as *distractors*. For instance, in the case of classifying Honduran Emeralds and Kakapos (rare species of birds), examples of butterflies, cars or ducks can serve as distractors (Fig. 1). By the virtue of its task-exclusivity, distractors can be obtained from various data-abundant domains with categories that could be semantically unrelated to novel task categories. However, in this work, we restrict ourselves to just the base data as a source for distractors. This allows us to efficiently reuse the data under standard settings and directly compare with prior works.

To this end, we pose the imminent question – *Can distractors improve few-shot generalization?* The answer is, somewhat surprisingly, yes. To elucidate how, we propose *ConFT*, a simple finetuning method based on a **contrastive** loss that contrasts pairs of the same class against those from different classes. We show that with a few simple but crucial modifications to the standard contrastive loss, distractors can be incorporated to boost generalization. We hypothesize that in the absence of extensive in-domain supervision for prior experience, distractor-aware finetuning can yield non-trivial gains. Towards the design of the loss function, we adopt an *asymmetric* construction of similarity pairs to ensure that *distractors contribute only through different-class pairs*. Our key insight here is two-fold – 1) generalization in contrastive learning can be influenced by not only same-class but also different-class pairs; 2) construction of different-class pairs is extremely flexible in that it can include samples from task-specific as well as task-exclusive categories. As a test of generality, we study the effect of our finetuning approach in conjunction with two different prior learning setups, namely, cross-domain and unsupervised prior learning. **Our contributions** are as follows.

- We propose contrastive finetuning, *ConFT*, a novel finetuning method for transfer based few-shot classification.

- We show how distractors can be incorporated in a contrastive objective to improve few-shot generalization.
- The proposed method outperforms state-of-the-art approaches by up to 12 points in the cross-domain few-shot learning and up to 5 points in unsupervised prior learning settings.

## 2. Related Work

### 2.1. Few-Shot Classification

Modern algorithms for few-shot classification are predominantly based on meta-learning where the goal is to quickly adapt to novel tasks. These approaches can be broadly classified into three categories: initialization based [13, 40, 41, 51, 55], hallucination based [2, 22, 75], and metric-learning based [4, 33, 59, 61, 64, 73] methods. Despite the growing interest in sophisticated meta-learning techniques, recent works [1, 7, 10, 66] have demonstrated that even simple finetuning based transfer learning [15, 19, 34, 45, 80] can outperform them. Such baselines usually involve cross-entropy training over the base categories followed by finetuning over a disjoint set of novel classes. Following these results, we further the investigation of finetuning for few-shot classification.

**Cross-Domain Few-Shot Classification:** A number of recent works [1, 12, 18, 42, 44, 56, 68, 69] have been proposed to address the cross-domain setup where base and novel classes are not only disjoint but also belong to different domains. Interestingly, [7] demonstrated that in this setup too, finetuning based transfer approaches outperformed popular meta-learning methods by significant margins. Following that, [68] proposed to learn feature-wise transformations via meta-learning to improve few-shot generalization of metric-based approaches. While in standard finetuning, the embedding model is usually frozen to avoid overfitting, recent works like [1, 18] have shown that frozen embeddings can hinder few-shot generalization. In this work, we build upon these developments to propose a more effective finetuning method over the entire embedding model.

In the context of learning from heterogeneous domain, [67] introduced a benchmark for multi-domain few-shot classification. This benchmark has been adopted by some recent works [8, 11, 39, 57]. While multiple base domains can alleviate cross-domain learning, we test our approach on a more challenging setup [68] that only involves a single base domain. Recent works used [68] as a benchmark to evaluate the importance of representation change [42] and spatial contrastive learning [44] in cross-domain few-shot classification. Another related work [69] leveraged unlabelled data from the novel domain in addition to few-shot labelled data to improve the task performance in a similar benchmark [18]. In contrast to [69], we operate under a

limited access to novel domain data, *i.e.* only the few-shot labelled data.

**Unlabelled Data in Few-Shot Classification:** Our use of unlabelled data in the form of distractors is inspired from cognitive neuroscience studies [38] describing the effect of visual distractors on learning and memory. Prior works that use additional unlabelled data for few-shot classification include [5, 16, 37, 52, 63, 76]. Complementary to [5, 16, 63] that exploit unlabelled data via self-supervised objectives in the prior learning phase, we use unlabelled data specifically for task-specific finetuning. Nonetheless, combining both perspectives could yield further benefits and is left for future work.

More related approaches [37, 52] combined heterogeneous unlabelled data, *i.e.*, task-specific data and distractors, in a semi-supervised framework. Our distractor-aware finetuning differs from these works in two important ways: our few-shot classification is strictly inductive in that we do not use unlabelled data specific to the task, and our method leverages distractors instead of treating them as interference that needs to be masked out. The most relevant methods [1, 15], like us, reused the base (or source) domain as a source for additional data. The key difference, however, is that their success relies on effective alignment of the base and novel classes, whereas we benefit from contrasting the two. While the importance of distractor-aware learning has been investigated in the context of object detection [47, 82], their benefit to few-shot generalization has not been studied before.

Recently, [26, 29, 30] have studied few-shot classification in the context of unsupervised prior-learning where the base data is unlabelled. In this work, we evaluate the benefit of contrastive finetuning under this setting and compare it to existing methods.

## 2.2. Contrastive Learning

Contrastive learning yields a similarity distribution over data by comparing pairs of different samples [60]. Recently, contrastive learning [20, 21, 58, 62] based methods have emerged as the state of the art for supervised [28, 31, 78] and self-supervised [6, 23, 25, 27, 65, 70, 79] representation learning. While the supervised approaches primarily exploit ground-truth labels to construct same-class pairs, self-supervised techniques leverage domain knowledge in the form of data augmentation to generate such pairs. As a special case, [31] maximized the benefit by integrating both forms of contrastive losses into a single objective. In this work, we use a modified version of the supervised contrastive loss when more than one labelled example is available per category. However, in the extreme case of 1-shot classification, it switches to self-supervised contrastive learning. Recent works such as [11, 44] also explored contrastive learning in the context of few-shot classification.

While they use contrastive objectives at the *prior-learning* stage to learn a general-purpose representation solely on the base domain, our method uses a contrastive objective at *finetuning* to improve the *downstream-task-specific* representation directly on the target domain task with base domain data as distractors. As a design choice, we adopt the contrastive loss over other losses like cross-entropy since it allows us to leverage distractor data that does not belong to the novel categories but improves generalization.

## 3. Our Approach

To achieve the goal of few-shot generalization, our contrastive finetuning method, ConFT, optimizes for two simultaneous objectives. First, it aims to bring task-specific samples that share the same class close to each other; and second, it strives to push apart samples that belong to different classes. This two-fold objective can lead to compact clusters that are well separated amongst each other. In the following sections, we first introduce some notations that we then use to formally describe our approach. An overview of our method is presented in Fig. 2.

### 3.1. Preliminaries

Consider an input space  $\mathcal{X}$  and a categorical label set  $\mathcal{Y} = \{c_1, \dots, c_M\}$  where each of the  $M$  classes is represented via one-hot encoding. A representation space  $\mathcal{R} \subset \mathbb{R}^r$  of the input is defined by the composition of an augmentation function  $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{X}$  and a representation model  $\mathcal{M}_\theta : \mathcal{X} \rightarrow \mathcal{R}$ , parameterized by  $\theta$ . The augmentation function is a composition of standard image transformations such as random cropping, color jittering, horizontal flipping *etc.* Given a small number  $K$ , a few-shot classification task  $\tau$  can be defined as the collection of a support set,  $\tau_{\text{supp}} = \{(x_i, y_i) \mid x_i \in \mathcal{X}, y_i \in \mathcal{Y}, i \in I_{\text{supp}}\}$  with  $K$  examples per class, and a query set,  $\tau_q = \{\tilde{x}_j \mid \tilde{x}_j \in \mathcal{X}, j \in I_q\}$  sampled from the same (but unobserved) classes. Here,  $I_{\text{supp}}$  and  $I_q$  are the collection of indices for the support and query sets, respectively. The few-shot classification goal is to leverage the support set to obtain a classifier for the query samples. In this case, the classifier is constructed over the representation model obtained via *contrastive finetuning* of a prior model,  $\mathcal{M}_{\theta_0}$  over  $\tau_{\text{supp}}$ .

### 3.2. The ConFT Objective

A key component of the ConFT objective is that it includes unlabelled samples, *distractors*, to improve few-shot generalization. Formally, a distractor set,  $S_{\text{dt}} = \{x_i \mid x_i \in \mathcal{X}, i \in I_{\text{dt}}\}$ , drawn from a domain  $D : \mathcal{X} \times \mathcal{Y}_D$  together with the task-specific support set  $\tau_{\text{supp}}$ , constitutes the training data for few-shot learning. Here, the distractor class set  $\mathcal{Y}_D$  is assumed to be *task-exclusive*,  $\mathcal{Y}_D \cap \mathcal{Y} = \emptyset$ . Starting with a support set example  $i$  (a.k.a anchor), we first construct an anchor-negative index set,  $N(i) = \{p \in$

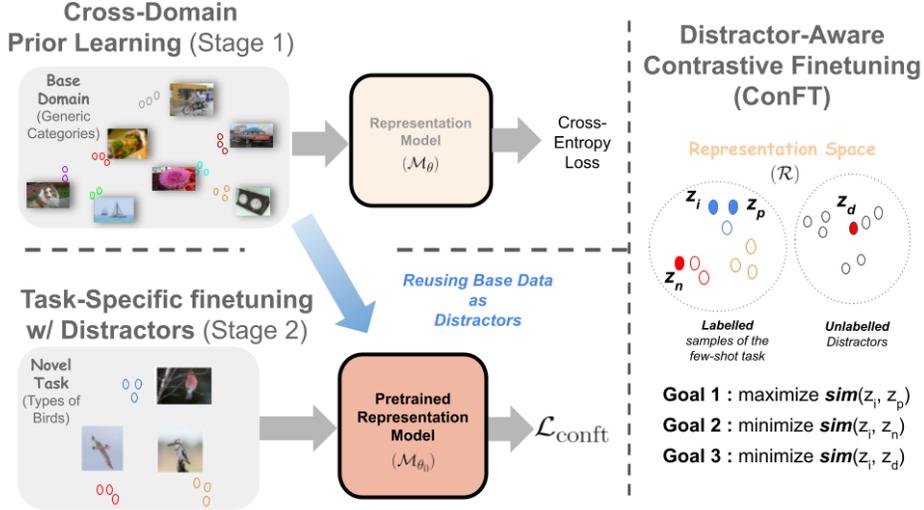


Figure 2: **Contrastive Finetuning in Cross-Domain Few-Shot Learning:** Our contrastive finetuning approach to few-shot classification comprises of two stages: 1) The *prior learning* stage trains a representation model on labelled (under cross-domain settings) base data using a cross-entropy loss; 2) The pretrained representation model is then finetuned over *task-specific* samples as well as *distractors* using a contrastive loss. For each task sample  $z_i$ , the contrastive objective (*right*) maximizes a similarity score,  $\text{sim}$ , over same-class pairs while minimizing it over other pair types. In the absence of enough labelled examples, distractors can improve classification by pushing apart task-specific clusters (here, different classes of birds).

$I_{\text{supp}} \setminus \{y_i \neq y_p\}$ , and an anchor-positive index set  $P(i)$  such that  $y_p = y_i, \forall p \in P(i)$ . Samples indexed by  $N(i)$  are treated as negatives within the task, whereas those indexed by  $I_{\text{dt}}$  act as negatives exclusive to the task. Finally, we define our contrastive loss that uses a  $l_2$ -normalized representation  $z \in \mathbb{R}^r$  as follows

$$\mathcal{L}_{\text{confit}}(\theta) = -\frac{1}{|I_{\text{supp}}|} \sum_{i \in I_{\text{supp}}} \frac{1}{|P(i)|} \sum_{p \in P(i)} \log l_{ip},$$

$$l_{ip} = \frac{\exp\left(\frac{z_i \cdot z_p}{\gamma}\right)}{\exp\left(\frac{z_i \cdot z_p}{\gamma}\right) + \sum_{n \in N(i)} \exp\left(\frac{z_i \cdot z_n}{\gamma}\right) + \sum_{d \in I_{\text{dt}}} \exp\left(\frac{z_i \cdot z_d}{\gamma}\right)}$$
(1)

where  $\gamma$  is a temperature hyper-parameter. The finetuning objective is simply the minimization of  $\mathcal{L}_{\text{confit}}$  to yield optimal parameters  $\theta_\tau$  specific to task  $\tau$ . To classify the query samples, we construct a nearest-mean classifier [12, 48, 60] atop the updated representation  $\mathcal{M}_{\theta_\tau}$ . The class-specific weight vectors are computed as an average over the representations of  $K$  support examples pertaining to that class. The  $j^{\text{th}}$  query sample is then assigned to the class whose weight vector has the largest cosine similarity (and hence, nearest in the Euclidean sense) with the query representation. We use the accuracy of this classifier to compare various baselines in the experiment section.

**Construction of Anchor-Positive Set  $P(i)$ :** To construct an anchor-positive set, we randomly pair task-

samples belonging to the same class with no sample occurring in more than one pair. In each pair, if one is assigned to be the anchor, the other acts as its positive. As an example, in a 5-way 4-shot task, our stochastic pair construction will result in 10 pairs where each of the 5 classes has 2 pairs. In the case when the number of shots is odd, we omit one sample from each class to allow even pairing. The omission is, however, not an issue in the overall scheme of finetuning where multiple steps of gradient descent optimization ensures that eventually each sample gets to participate with equal chance. In the special case 1-shot learning, anchor-positive sets are constructed similar to [6] using augmentation  $\mathcal{A}$ .

### 3.3. Relative Importance of Anchor-Negatives

Given the loss formulation of (1), both task-specific (few-shot) and task-exclusive (distractor) anchor-negatives influence the loss proportionate to their respective batch sizes. While the batch size of task-specific negatives  $N(i)$  is upper bounded by the number of ways  $M$  and the number of shots  $K$ , the batch size of distractors can be made as large as that of the domain itself, *i.e.*,  $|D|$ . In standard contrastive learning paradigms with only task-specific and no task-exclusive training examples, large batch sizes of negatives are known to be beneficial for downstream task performance. However, in our case where both types of negatives exist, naively increasing distractor batch size can be

counterproductive (shown in the supplementary). We suspect that too many distractors might overshadow the effect of task-specific negatives that can be more crucial for generalization. Also, the effect might vary according to the proximity of distractors with respect to task samples in the representation space. Nonetheless, there is a need to balance the undue influence of distractors by adjusting the batch sizes. To avoid an extensive search for an optimal batch size specific to the distractor domain, we propose a domain-agnostic weighting scheme for the anchor-negatives proportional to their batch sizes as follows

$$l_{ip} = \frac{\exp(\frac{z_i \cdot z_p}{\gamma})}{\exp(\frac{z_i \cdot z_p}{\gamma}) + \alpha \sum_{n \in N(i)} \exp(\frac{z_i \cdot z_n}{\gamma}) + (2 - \alpha) \sum_{d \in I_{dt}} \exp(\frac{z_i \cdot z_d}{\gamma})}, \quad (2)$$

$$\alpha = 2 \frac{|I_{dt}|}{|N(i)| + |I_{dt}|}$$

We found that this simple weighting scheme makes the few-shot performance robust to batch size variations and also improves the overall performance (see the supplementary).

## 4. Prior Learning and Distractors

Thus far we have assumed the access to a distractor set,  $S_{dt}$  and a prior model  $\mathcal{M}_\theta$ . In this section, we describe how to obtain them and how distractors boost generalization. Recall that our goal is to achieve few-shot generalization by finetuning a prior model over the few-shot task. However, due to the scarcity of task-specific labelled examples, a reasonably strong prior encoded in the model parameters  $\theta_0$  is crucial for preventing overfitting, especially when using high-capacity models like neural networks. We next describe two different ways of learning such a prior that can serve as a good initialization for subsequent finetuning.

### 4.1. Types of Prior Learning

**Cross-Domain Learning:** In the cross-domain setup, we are provided a labelled dataset,  $D_l = \{(x_i, y_i) | x_i \in \mathcal{X}_{sc}, y_i \in \mathcal{Y}_{sc}\}_{i=1}^{|D_l|}$  drawn from a source domain  $\mathcal{X}_{sc} \times \mathcal{Y}_{sc}$ , such that the categorical label set  $\mathcal{Y}_{sc}$  is disjoint from novel categories  $\mathcal{Y}$ . The key characteristic of this setup is that the distribution of  $M$ -way  $K$ -shot tasks, if constructed out of  $D_l$ , will be significantly different from novel tasks in the target domain. Such distribution shift could arise due to difference in task granularity (*e.g.*, coarse-grained vs. fine-grained) or shift in input distribution or both. In this work, we consider the case where the shift in task granularity is notably more than the input distribution. Towards the goal of learning a reasonably strong prior, we adopt a simple objective that minimizes cross-entropy loss over all categories in  $D_l$ . During finetuning, the distractors are sampled

from  $D_l$ , thus, naturally satisfying the non-overlapping categories assumption with respect to novel tasks.

**Unsupervised Prior Learning:** For unsupervised prior learning, we are given an unlabelled dataset,  $D_u = \{v_i\}_{i=1}^{|D_u|}$  drawn from a source domain  $\mathcal{X}_{su} \times \mathcal{Y}_{su}$ , such that the corresponding labels in  $\mathcal{Y}_{su}$  are unobserved. While there are no explicit assumptions about the task distribution gap in this setting, the strength of the learnt prior is likely to be more reasonable when the distribution gap is small. To learn a suitable prior using  $D_u$ , we use the SimCLR loss [6] as a form of self-supervised objective. Our choice of this objective over others [23] was based on its superior performance found in our preliminary experiments. Priors learnt via self-supervised contrastive objectives on large base datasets (like, ImageNet) have been shown to transfer well to many-shot downstream tasks. In this work we show that such objectives are effective even with smaller base datasets and few-shot downstream tasks. In our experiments, priors learnt in this way already outperform state-of-the-art approaches [26,29,30] that are then further improved by our proposed contrastive finetuning. In this setup, we use  $D_u$  as the source for distractors where the assumption of non-overlapping categories is satisfied with high probability, provided the base dataset is relatively large and encapsulates a wide variety of categorical concepts.

### 4.2. Distractor-Aware Generalization

The most important and perhaps surprising aspect of our method is that distractors, despite being drawn from unrelated (to novel task) categories, can improve generalization. To understand the underlying mechanism, we propose to measure the change in quality of task-specific representation before and after finetuning. Particularly, given a few-shot task with  $M$  classes, we define the subset of query samples,  $I_q^c \subset I_q$  that share the same class<sup>1</sup> and two other quantities – cluster spread  $u_{spread}^q$  and cluster-separation  $u_{sep}^q$  that measure the degree of clustering in the representation space. Specifically,

$$u_{spread}^q(\theta_t) = \frac{1}{M} \sum_{m=1}^M \sum_{\substack{i \in I_q^c \\ j \in I_q^c \setminus \{i\}}} (1 - z_i \cdot z_j), \quad (3)$$

$$u_{sep}^q(\theta_t) = \frac{1}{M} \sum_{m=1}^M \sum_{\substack{i \in I_q^c \\ j \in I_q \setminus I_q^c}} (1 - z_i \cdot z_j), \quad (4)$$

where  $\theta_t$  are the parameters of the representation model after  $t$  finetuning epochs. For each of the above quantities,

<sup>1</sup>Note that the query class labels are considered only for analysis purposes. In practice, they are not observed.

we define the change,  $\delta_*^q(t) = u_*^q(\theta_t) - u_*^q(\theta_0)$ , and relative change,  $\delta_*^{\text{rel},q}(t) = \frac{\delta_*^q(t)}{\kappa(\theta_0)}$  where, the subscript can be *sep* or *spread* and division by a fixed value,  $\kappa(\theta_0)$  ensures scale invariance. Finally, to quantify generalization within a given target domain, we define the average relative change,  $\mathbb{E}_\tau [\delta_*^{\text{rel},q}(t)]$  over a large number of tasks sampled from that domain. The average relative change can also be defined for support examples by simply swapping superscript ‘q’ with ‘s’. Also, in practice, we use  $u_{\text{sep}}^s(\theta_0)$  as the fixed value for  $\kappa(\theta_0)$  irrespective of the superscript or subscript.

### 4.3. A Multitask Variant of ConFT

While our original objective (2) is agnostic to distractor supervision, finetuning in the cross-domain setting can further benefit from distractor labels. To that end, we introduce an auxiliary loss  $\mathcal{L}_{\text{mtce}}$  during finetuning that minimizes the cross-entropy between predicted probabilities and one-hot encodings of the ground-truth label averaged over the base data,  $D_l$ . This leads to a new multitask formulation

$$\mathcal{L}_{\text{mt-conft}} = \mathcal{L}_{\text{conft}} + \lambda \mathcal{L}_{\text{mtce}}, \quad (5)$$

where we fix the relative weighting factor  $\lambda = 1$  in our experiments and use a cosine classifier [7] for  $\mathcal{L}_{\text{mtce}}$ . We found that this simple extension led to significant performance gains in some domains while marginal in others, depending on domain characteristics.

## 5. Experiments

Following sections first introduce some baselines (§5.1) and present our main results for contrastive finetuning in the cross-domain setup (§5.2). Then, §5.3 elucidates the generalization mechanism of ConFT followed by ablations in §5.4. Finally, §5.5 demonstrates the performance of our approach in the unsupervised prior learning setup.

**Datasets and Benchmarks:** We evaluate our proposed finetuning method in a variety of novel domains spanning across two different paradigms for prior learning. For cross-domain evaluations, we adopt the benchmark introduced by [68] that comprises of Cars [35], CUB [77], Places [81], and Plantae [72] as the novel domains and *miniImageNet* [50] as the base domain. Each dataset is split into *train*, *val* and, *test* categories (please refer to the supplementary for details), where tasks sampled from the *test* split are used to evaluate the few-shot performance in respective domains. We use the *val* splits for cross-validating the hyperparameters and the *train* split of *miniImageNet* as our base data. For experiments in unsupervised prior learning, we use the same *train* split of *miniImageNet* to learn a self-supervised representation that is then evaluated for few-shot performance on *miniImageNet-test*. We present additional results on Meta-Dataset [67] in the supplementary.

**Backbone (Representation Model):** Following best practices in cross-domain few-shot learning, we adopt a ResNet10 [24] model for most of our experiments. In the unsupervised learning case, we use a four-layer CNN consistent with existing works except for a reduced filter size from 64 to 20 in the final layer. This modification was found to improve contrastive finetuning performance.

**Optimization and Hyperparameters:** In this work, we evaluate few-shot performance over 5-way 1-shot and 5-way 5-shot tasks with 15 query samples, irrespective of the prior learning setup. For the contrastive finetuning, we use an ADAM [32] optimizer with a suitable learning rate and early-stopping criteria. Our proposed method has a few hyperparameters such as the temperature ( $\gamma$ ), learning rate, early-stopping criteria, and data augmentation ( $\mathcal{A}$ ). However, recent studies [43] have highlighted that excessive hyperparameter tuning on large validation sets can lead to overoptimistic results in limited-labelled data settings like semi-supervised learning. Thus, we keep an extremely small budget for hyperparameter tuning. Among the mentioned hyperparameter, the one with the most number of parameters is the augmentation function  $\mathcal{A}$ . In this work, we do not tune  $\mathcal{A}$  to any specific target domain. Instead, we use a fixed augmentation scheme introduced by [7] for the cross-domain setting and AutoAugment [9] for the unsupervised prior learning case. Please refer to the supplementary for a detailed summary of hyperparameters used in our experiments.

### 5.1. Baseline Comparisons

We begin our evaluations by comparing various baselines for finetuning in Table 1. These include two simple baselines (introduced in [7]) and two strong baselines (introduced in [18]). While the simple baselines freeze the backbones, the others allow finetuning over the entire embedding model. Another key difference is that the simple baselines are evaluated using standard linear evaluation [6, 7], whereas the rest are evaluated using nearest-mean classifiers. We compare the performance of all these baselines to our vanilla and multi-task (MT) versions of ConFT. Following previous works, the learning rates for the simple baselines are kept at 0.01, whereas for others (including ours), we use smaller learning rates (0.005 or 0.0005). We observe that among the baselines, the cosine classifier based baseline, FT-all (CC), outperforms the linear classifier based FT-all (LC). However, both versions of our finetuning approach significantly outperform all baselines across various dataset and shot settings.

### 5.2. ConFT for Cross-Domain Prior Learning

In this section, we present our main results on cross-domain few-shot learning (see Table 2). We compare our approach with various prior works on the LFT benchmark

Finetuning Method		CUB		Cars		Places		Plantae	
Loss	FT Type	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
Cross-Entropy	fixed-BB (LC) [7]	39.77 ± 0.66	51.33 ± 0.70	33.99 ± 0.64	44.14 ± 0.70	44.53 ± 0.75	55.94 ± 0.69	37.07 ± 0.70	46.58 ± 0.69
	fixed-BB (CC) [7]	43.26 ± 0.76	62.87 ± 0.74	25.33 ± 1.85	50.40 ± 0.74	47.70 ± 0.76	69.48 ± 0.69	40.49 ± 0.77	56.64 ± 0.72
	FT-all (LC) [18]	40.81 ± 0.75	61.82 ± 0.72	34.50 ± 0.67	55.63 ± 0.75	45.91 ± 0.77	68.73 ± 0.73	37.51 ± 0.71	58.33 ± 0.68
	FT-all (CC) [18]	44.30 ± 0.73	67.05 ± 0.69	36.79 ± 0.76	57.65 ± 0.76	49.10 ± 0.78	70.32 ± 0.72	40.31 ± 0.76	61.30 ± 0.75
Contrastive	-	43.42 ± 0.75	62.80 ± 0.76	35.19 ± 0.66	51.41 ± 0.72	49.56 ± 0.80	70.71 ± 0.68	40.39 ± 0.79	55.54 ± 0.69
	ConFT (ours)	45.57 ± 0.76	70.53 ± 0.75	<b>39.11 ± 0.77</b>	61.53 ± 0.75	<b>49.97 ± 0.86</b>	72.09 ± 0.68	<b>43.09 ± 0.78</b>	62.54 ± 0.76
	MT-ConFT (ours)	<b>49.25 ± 0.83</b>	<b>74.45 ± 0.71</b>	37.36 ± 0.69	<b>62.54 ± 0.72</b>	<b>49.94 ± 0.81</b>	<b>72.71 ± 0.69</b>	41.82 ± 0.75	<b>63.01 ± 0.74</b>

Table 1: **Baseline Comparisons.** Results on 1-shot and 5-shot tasks on the LFT benchmark [68]. These results are obtained by averaging over 600 novel tasks, each consisting of 5 classes and 15 queries per class. We also present 95% confidence intervals. The train split of the *mini*ImageNet dataset is used as base data. Here, FT-all denotes the case where the entire embedding model is finetuned. Other abbreviations – BB: Backbone model (ResNet-10), LC: Linear Classifier, CC: Cosine Classifier with a multiplication factor of 10.

Method			1-shot			
Prior Learning	Task Specific Finetuning	Backbone	CUB	Cars	Places	Plantae
AAL [1]	arcmax	ResNet18	47.25 ± 0.76	-	-	-
MN [73]	-	ResNet10	35.89 ± 0.51	30.77 ± 0.47	49.86 ± 0.79	32.70 ± 0.60
MN w/ featTx [68]	-	ResNet10	36.61 ± 0.53	29.82 ± 0.44	51.07 ± 0.68	34.48 ± 0.50
RN [64]	-	ResNet10	42.44 ± 0.77	29.11 ± 0.60	48.64 ± 0.85	33.17 ± 0.64
RN w/ featTx [68]	-	ResNet10	44.07 ± 0.77	28.63 ± 0.59	50.68 ± 0.87	33.14 ± 0.62
GNN [59]	-	ResNet10+	45.69 ± 0.68	31.79 ± 0.51	53.10 ± 0.80	35.60 ± 0.56
GNN w/ featTx [68]	-	ResNet10+	47.47 ± 0.75	31.61 ± 0.53	<b>55.77 ± 0.79</b>	35.95 ± 0.58
MAML [13]	-	Conv4	40.51 ± 0.08	33.57 ± 0.14	-	-
ANIL [49]	-	Conv4	41.12 ± 0.15	34.77 ± 0.31	-	-
BOIL [42]	-	Conv4	44.20 ± 0.15	36.12 ± 0.29	-	-
CE Training	-	ResNet10	43.42 ± 0.75	35.19 ± 0.66	49.56 ± 0.80	40.39 ± 0.79
CE Training	ConFT (ours)	ResNet10	45.57 ± 0.76	<b>39.11 ± 0.77</b>	49.97 ± 0.86	<b>43.09 ± 0.78</b>
CE Training	MT-ConFT (ours)	ResNet10	<b>49.25 ± 0.83</b>	37.36 ± 0.69	49.94 ± 0.81	41.82 ± 0.75

Method			5-shot			
Prior Learning	Task Specific Finetuning	Backbone	CUB	Cars	Places	Plantae
Baseline [7]	-	ResNet18	65.57 ± 0.70	-	-	-
Baseline ++ [7]	-	ResNet18	62.04 ± 0.76	-	-	-
DiversityNCoop [12]	-	ResNet18	66.17 ± 0.55	-	-	-
AAL [1]	arcmax	ResNet18	72.37 ± 0.89	-	-	-
BOIL [42]	-	ResNet12	-	49.71 ± 0.28	-	-
MN [73]	-	ResNet10	51.37 ± 0.77	38.99 ± 0.64	63.16 ± 0.77	46.53 ± 0.68
MN w/ featTx [68]	-	ResNet10	55.23 ± 0.83	41.24 ± 0.65	64.55 ± 0.75	41.69 ± 0.63
RN [64]	-	ResNet10	57.77 ± 0.69	37.33 ± 0.68	63.32 ± 0.76	44.00 ± 0.60
RN w/ featTx [68]	-	ResNet10	59.46 ± 0.71	39.91 ± 0.69	66.28 ± 0.72	45.08 ± 0.59
GNN [59]	-	ResNet10+	62.25 ± 0.65	44.28 ± 0.63	70.84 ± 0.65	52.53 ± 0.59
GNN w/ featTx [68]	-	ResNet10+	66.98 ± 0.68	44.90 ± 0.64	<b>73.94 ± 0.67</b>	53.85 ± 0.62
MAML [13]	-	Conv4	53.09 ± 0.16	44.56 ± 0.21	-	-
ANIL [49]	-	Conv4	55.82 ± 0.21	46.55 ± 0.29	-	-
BOIL [42]	-	Conv4	60.92 ± 0.11	50.64 ± 0.22	-	-
CE Training	-	ResNet10	62.80 ± 0.76	51.41 ± 0.72	70.71 ± 0.68	55.54 ± 0.69
CE Training	ConFT (ours)	ResNet10	70.53 ± 0.75	61.53 ± 0.75	72.09 ± 0.68	62.54 ± 0.76
CE Training	MT-ConFT (ours)	ResNet10	<b>74.45 ± 0.71</b>	<b>62.54 ± 0.72</b>	72.71 ± 0.69	<b>63.01 ± 0.74</b>

Table 2: **Cross-Domain Few-Shot Classification Results.** We present the results with 95% confidence intervals and highlight the best performing methods. The results are an average over 600 tasks. Here, ‘-’ denotes numbers not reported by previous works.

[68]. We observe that overall our proposed approaches, ConFT and MT-ConFT, significantly outperform the best previous results in Cars (by **3 to 12 points**), Plantae (by **7 to 9 points**) and CUB (by **1.7 to 2 points**) domains. We also observe higher gains in the 5-shot setting than the 1-shot case, since more labelled examples can improve few-shot generalization. Further, we find that using the auxiliary loss (MT-ConFT) is more beneficial in the 5-shot case. In fact, it performs worse than ConFT in the 1-shot cases for Cars, Places, and Plantae. Such a degradation could be due to a misalignment between the self-supervised objective (to

which ConFT boils down in the 1-shot case) and the auxiliary cross-entropy loss. In the Places domain, “GNN w/ featTx” yields the best performance, whereas our approach outperforms the rest for the 5-shot case. We suspect that the use of a more sophisticated model in “GNN w/ featTx”, namely, graph neural net [59] built on top of a ResNet-10 model, leads to a better cross-domain generalization when the domain gap is smaller.

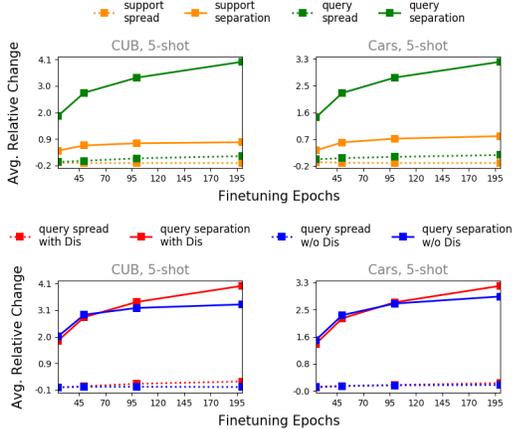


Figure 3: **Understanding Generalization in ConFT.** *Top:* Average relative change in cluster-spread and cluster-separation of support and query samples as a function of finetuning epochs. *Bottom:* Comparing the average relative change in cluster-spread and cluster-separation of only query samples under the presence and absence of distractors. The spread and separation quantities are averaged over 600 tasks for both *top* and *bottom*.

### 5.3. Effect of Distractors on Generalization

In this section, we investigate the central question – *How do distractors improve generalization?* We present two sets of plots in Figure 3 that track the change in cluster-spread and cluster-separation as the finetuning progresses. In the first set, we plot the average relative change,  $\mathbb{E}[\delta_{\text{spread}}^{\text{rel}}(t)]$  and  $\mathbb{E}[\delta_{\text{sep}}^{\text{rel}}(t)]$  (see §4.2) as a function of finetuning epochs,  $t$  for both support and query samples in 2 different settings, namely, CUB (5-shot) and Cars (5-shot). We observe that for support examples (yellow lines), cluster-spread decreases with increasing epochs while the cluster-separation increases. This is indeed what is expected for training datapoints (here, support examples) and serves as a sanity check. For query samples (green plot), on the other hand, both cluster-spread and separation increase with the progress in finetuning epochs. The key observation, however, is that cluster-separation increases to a much greater extent than the cluster-spread, thus improving overall discriminability between classes represented by these clusters. While the increase in cluster-separation hints towards the possible reason behind improved generalization, it is not clear how much of the improvement, if any, is a consequence of incorporating distractors. To delineate the effects of distractors from the contrastive loss itself, we present the second set of plots that compare the average relative change in *query* cluster spread and separation under the presence (red line) and absence (blue line) of distractors for the same data settings. We observe that with increas-

Anchor-Positives		Anchor-Negatives			Accuracy	
SPC	T-Pos	T-Neg	D-Neg	W	CUB, 5-shot	Cars, 5-shot
✓	✓	✓	✓	✓	69.61±0.68	61.01±0.74
✓		✓	✓		70.16±0.70	57.42±0.80
✓		✓	✓		67.44±0.71	59.00±0.73
✓		✓		✓	70.26±0.68	60.58±0.77
✓		✓	✓	✓	70.53±0.75	61.53±0.75

Table 3: **Ablation 1.** Novel task performance with various types of anchor-positives and anchor-negatives. Here, SPC: Stochastic Pair Construction, T: Task, D: Distractor, W: relative weighting.

Distractor Domain Size				Accuracy	
512	1024	2048	38400	CUB 5-shot	Cars 5-shot
✓				70.34 ± 0.71	61.16 ± 0.76
	✓			70.28 ± 0.70	61.11 ± 0.77
		✓		70.92 ± 0.69	61.31 ± 0.76
			✓	70.53 ± 0.75	61.53 ± 0.75

Table 4: **Ablation 2.** Novel task performance with varying sizes of the distractor domain, *i.e.*, *miniImageNet-train*. Note that this is *different* from distractor batch size  $|S_{\text{dt}}|$ .

ing finetuning epochs the gap between cluster-separation and cluster-spread widens to a larger extent in the presence of distractors than in their absence. This leads to our conclusion that **distractors help generalization by increasing task-specific cluster separation.**

### 5.4. Ablations

In this section, we introduce a few important ablations that help deconstruct the ConFT and MT-ConFT objectives. In Table 3, we compare our stochastic anchor-positive construction with naive inclusion of all positives for every anchor. We also, ablate the contribution of each type of anchor-negatives: task-specific and distractors and compare that to relative weighting of the two. In Table 4, we studied the importance of distractor domain size and found that novel task performance is fairly robust to the size of the distractor domain. This is particularly encouraging since we need not store the entire base data during finetuning.

### 5.5. ConFT for Unsupervised Prior Learning

In Table 5, we demonstrate the generality of contrastive finetuning by evaluating on the unsupervised prior learning benchmark *miniImageNet*. The key distinction from cross-domain settings is that we do not have labelled base data to learn from. So, we leverage self-supervised contrastive learning [6] on the unlabelled base data and show that it outperforms state of the art by 1 to 2 points. Finetuning the resultant representation with our ConFT objective further improves the accuracy by 2to 4 points. This is particularly significant, as the results come very close to supervised baselines that serve as performance upper bound in this set-

Method			5-way <i>mini</i> ImageNet	
Prior Learning	Finetuning	BB	1-shot	5-shot
Sup. MML	MFT	Conv4	46.81 $\pm$ 0.77	62.13 $\pm$ 0.72
Sup. PN	-	Conv4	46.56 $\pm$ 0.76	62.29 $\pm$ 0.71
-	RandInit	Conv4	27.59 $\pm$ 0.59	38.48 $\pm$ 0.66
BG-MML [26]	MFT	Conv4	36.24 $\pm$ 0.74	51.28 $\pm$ 0.6
BG-PN [26]	-	Conv4	36.62 $\pm$ 0.70	50.16 $\pm$ 0.7
DC-MML [26]	MFT	Conv4	39.90 $\pm$ 0.74	53.97 $\pm$ 0.70
DC-PN [26]	-	Conv4	39.18 $\pm$ 0.71	53.36 $\pm$ 0.70
U-MML [29]	MFT	Conv4	39.93	50.73
LG-MML [30]	MFT	Conv4	40.19 $\pm$ 0.58	54.56 $\pm$ 0.55
LG-PN [30]	-	Conv4	40.05 $\pm$ 0.60	52.53 $\pm$ 0.51
SimCLR [6]	-	Conv4	41.54 $\pm$ 0.61	56.57 $\pm$ 0.59
SimCLR	ConFT	Conv4	<b>43.45<math>\pm</math>0.60</b>	<b>60.02<math>\pm</math>0.57</b>

Table 5: **Unsupervised Prior Learning.** The results are averaged over 1000 novel tasks and are presented with 95% confidence intervals. Here, MFT refers to meta-style finetuning [13]. MML: Maml, PN: ProtoNet, DC: DeepCluster, U: Umtra, LG: Lasium-Gan, BG: BiGAN, Sup.: Supervised, BB: Backbone.

ting [30].

## 6. Conclusion

We introduce a novel contrastive finetuning approach to few-shot classification. Specifically, our method leverages distractors to improve generalization by encouraging cluster separation of the novel task samples. We show that our method leads to significant performance gains in both cross-domain and unsupervised prior learning setups.

## References

- [1] Arman Afrasiyabi, Jean-François Lalonde, and Christian Gagné. Associative alignment for few-shot image classification. In *ECCV*, 2020. 1, 2, 3, 7
- [2] Antreas Antoniou, Amos J. Storkey, and Harrison A Edwards. Data augmentation generative adversarial networks. *arXiv*, 2017. 2
- [3] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *NeurIPS*. 2007. 12
- [4] Jane Bromley, James Bentz, Leon Bottou, Isabelle Guyon, Yann Lecun, Cliff Moore, Eduard Sackinger, and Rookpak Shah. Signature verification using a "siamese" time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7:25, 1993. 2
- [5] Guy Bukchin, Eli Schwartz, Kate Saenko, Ori Shahar, R. Feris, Raja Giryes, and Leonid Karlinsky. Fine-grained angular contrastive learning with coarse labels. *arXiv*, 2020. 3
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv*, 2020. 3, 4, 5, 6, 8, 9, 13
- [7] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *ICLR*, 2019. 2, 6, 7, 15
- [8] Yinbo Chen, Xiaolong Wang, Zhuang Liu, Huijuan Xu, and Trevor Darrell. A new meta-baseline for few-shot learning. *arXiv*, 2020. 2
- [9] Ekin D. Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation policies from data. *arXiv*, 2018. 6
- [10] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *ICLR*, 2020. 1, 2
- [11] Carl Doersch, Ankush Gupta, and Andrew Zisserman. Crosstransformers: Spatially-aware few-shot transfer. In *NeurIPS*, 2020. 2, 3, 17
- [12] Nikita Dvornik, Cordelia Schmid, and Julien Mairal. Diversity with cooperation: Ensemble methods for few-shot classification. In *ICCV*, 2019. 2, 4, 7
- [13] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 1, 2, 7, 9
- [14] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *CVPR*, 2017. 12
- [15] Weifeng Ge and Yizhou Yu. Borrowing treasures from the wealthy: Deep transfer learning through selective joint finetuning. In *CVPR*, 2017. 2, 3
- [16] Spyros Gidaris, Andrei Bursuc, Nikos Komodakis, Patrick Pérez, and Matthieu Cord. Boosting few-shot visual learning with self-supervision. In *CVPR*, 2019. 3
- [17] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. *CVPR*, 2018. 1
- [18] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *ECCV*, 2020. 2, 6, 7
- [19] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: Transfer learning through adaptive fine-tuning. *arXiv*, 2018. 2
- [20] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS*, 2010. 3
- [21] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, 2006. 3
- [22] Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *ICCV*, 2017. 2
- [23] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 3, 5
- [24] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6, 13, 15
- [25] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *ICLR*, 2019. 3

- [26] Kyle Hsu, Sergey Levine, and Chelsea Finn. Unsupervised learning via meta-learning. In *ICLR*, 2019. 1, 3, 5, 9, 12
- [27] Olivier J. Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, S. M. Ali Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv*, 2020. 3
- [28] Konstantinos Kamnitsas, Daniel C. Castro, Loïc Le Folgoc, Ian Walker, Ryutaro Tanno, Daniel Rueckert, Ben Glocker, Antonio Criminisi, and Aditya V. Nori. Semi-supervised learning via compact latent space clustering. In *ICML*, 2018. 3
- [29] Siavash Khodadadeh, Ladislau Boloni, and Mubarak Shah. Unsupervised meta-learning for few-shot image classification. In *NeurIPS*, 2019. 1, 3, 5, 9
- [30] Siavash Khodadadeh, Sharare Zehtabian, Saeed Vahidian, Weijia Wang, Bill Lin, and Ladislau Boloni. Unsupervised meta-learning through latent-space interpolation in generative models. In *ICLR*, 2021. 1, 3, 5, 9
- [31] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *NeurIPS*, 2020. 3, 14
- [32] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, 2015. 6
- [33] Gregory R. Koch. Siamese neural networks for one-shot image recognition. In *ICML Deep Learning Workshop. Vol. 2.*, 2015. 1, 2
- [34] Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do better imagenet models transfer better? *CVPR*, 2019. 2
- [35] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV Workshops*, 2013. 6
- [36] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350:1332 – 1338, 2015. 1
- [37] Xinzhe Li, Qianru Sun, Yaoyao Liu, Qin Zhou, Shibao Zheng, Tat-Seng Chua, and Bernt Schiele. Learning to self-train for semi-supervised few-shot classification. In *NeurIPS*, 2019. 3
- [38] Chen Liang, Xiao Yang, Neisarg Dave, Drew Wham, Bart Pursel, and C. Lee Giles. Distractor generation for multiple choice questions using learning to rank. In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, 2018. 3
- [39] Orchid Majumder, Avinash Ravichandran, Subhansu Maji, M. Polito, Rahul Bhotika, and Stefano Soatto. Revisiting contrastive learning for few-shot classification. *arXiv*, 2021. 2
- [40] Munkhdalai, Tsendsuren, Yu, and Hong. Meta networks. In *ICML*, 2017. 2
- [41] Alex Nichol and John Schulman. Reptile: A scalable meta-learning algorithm. *arXiv*, 2018. 2
- [42] Jaehoon Oh, Hyungjun Yoo, ChangHwan Kim, and Se-Young Yun. {BOIL}: Towards representation change for few-shot learning. In *ICLR*, 2021. 2, 7
- [43] Avital Oliver, Augustus Odena, Colin Raffel, Ekin D Cubuk, and Ian J Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. *arXiv*, 2018. 6
- [44] Yassine Ouali, Céline Hudelot, and Myriam Tami. Spatial contrastive learning for few-shot classification. *arXiv*, 2020. 2, 3, 15, 16
- [45] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 2010. 2
- [46] Cheng Perng Phoo and Bharath Hariharan. Self-training for few-shot transfer across extreme task differences. In *ICLR*, 2021. 1
- [47] Horst Possegger, Thomas Mauthner, and Horst Bischof. In defense of color-based model-free tracking. In *CVPR*, June 2015. 3
- [48] Hang Qi, David Lowe, and Matthew Brown. Low-shot learning with imprinted weights. In *CVPR*, 2018. 4
- [49] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. In *ICLR*, 2020. 7
- [50] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 1, 6
- [51] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 2
- [52] Mengye Ren, Sachin Ravi, Eleni Triantafillou, Jake Snell, Kevin Swersky, Josh B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classification. In *ICLR*, 2018. 3
- [53] Joshua David Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with hard negative samples. In *ICLR*, 2021. 15
- [54] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 1
- [55] Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *ICLR*, 2019. 2
- [56] Jeongun Ryu, Jaewoong Shin, H. B. Lee, and Sung Ju Hwang. Metaperturb: Transferable regularizer for heterogeneous tasks and architectures. *arXiv*, 2020. 2
- [57] Tonmoy Saikia, T. Brox, and C. Schmid. Optimized generic feature learning for few-shot classification across domains. *arXiv*, 2020. 2
- [58] Ruslan Salakhutdinov and Geoff Hinton. Learning a nonlinear embedding by preserving class neighbourhood structure. In *AISTATS*, 2007. 3
- [59] Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *ICLR*, 2018. 1, 2, 7
- [60] Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *ICML*, 2019. 3, 4
- [61] Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017. 1, 2, 15, 16

- [62] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NeurIPS*, 2016. 3
- [63] Jong-Chyi Su, Subhransu Maji, and Bharath Hariharan. When does self-supervision improve few-shot learning? In *ECCV*, 2020. 3, 15, 16
- [64] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip H.S. Torr, and Timothy M. Hospedales. Learning to compare: Relation network for few-shot learning. In *CVPR*, 2018. 1, 2, 7
- [65] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv*, 2019. 3
- [66] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: A good embedding is all you need? *arXiv*, 2020. 1, 2, 15
- [67] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. Meta-dataset: A dataset of datasets for learning to learn from few examples. In *ICLR*, 2020. 2, 6
- [68] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. Cross-domain few-shot classification via learned feature-wise transformation. In *ICLR*, 2020. 1, 2, 6, 7, 12, 13
- [69] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. Cross-domain few-shot classification via learned feature-wise transformation. In *ICLR*, 2020. 2
- [70] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv*, 2018. 3
- [71] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 9, 2008. 12
- [72] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2016. 1, 6
- [73] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NIPS*. 2016. 1, 2, 7
- [74] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M. Summers. Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1
- [75] Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan. Low-shot learning from imaginary data. In *CVPR*, 2018. 2
- [76] Yu-Xiong Wang and Martial Hebert. Learning from small sample sets by combining unsupervised meta-training with cnns. In *NeurIPS*, 2016. 3
- [77] Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010. 6
- [78] Zhirong Wu, Alexei A Efros, and Stella Yu. Improving generalization via scalable neighborhood component analysis. In *ECCV*, 2018. 3
- [79] Zhirong Wu, Yuanjun Xiong, X Yu Stella, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 3
- [80] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *NeurIPS*, 2014. 2
- [81] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, , and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 40(6):1452–1464, 2018. 6
- [82] Zheng Zhu, Qiang Wang, Li Bo, Wei Wu, Junjie Yan, and Weiming Hu. Distractor-aware siamese networks for visual object tracking. In *ECCV*, 2018. 3

## A. Overview of ConFT

Algorithm 1 provides an overview of our distractor-aware contrastive finetuning approach ConFT.

## B. Additional Experimental Details

### B.1. Data Domains

Problem Setup (Prior Learning)					
Domain	Dataset	# categories per split			
		train	val	test	
Cross-Domain	Base	<i>miniImageNet</i>	64	16	20
	Novel	CUB	100	50	50
	Novel	Cars	98	49	49
	Novel	Places	183	91	91
Unsupervised	Base and Novel	<i>miniImageNet</i>	64	16	20

Table 6: **Dataset statistics for both cross-domain [68] and unsupervised prior learning settings [26].** Each dataset is split into *train*, *val*, and *test* categories. For the cross-domain setup, the *train* split of *miniImageNet* is always used as the base domain whereas the *test* splits of other datasets are used as the novel domain on which few-shot evaluation is performed. For the unsupervised prior learning setup, *train* split of *miniImageNet* is stripped off its labels to emulate an unlabelled base domain, whereas the *test* split is used as the novel domain. In both setups, *val* splits are used to cross-validate hyperparameters specific to the associated novel domain.

In the main paper, we evaluated our finetuning method on various datasets that serve as base or novel domains in cross-domain as well as unsupervised prior learning settings. Here, in Table 6, we summarize the statistics of these datasets along with their specific use as base or novel domain. Additionally, in Table 7, we visualize these domains, both qualitatively and quantitatively, to provide a reference to their relative proximity in the representation space. This proximity provides a rough estimate of how related two domains are and consequently, the degree of knowledge transfer across domains for cross-domain few-shot classification.

For the qualitative visualization in Table 7, we use t-SNE [71] to embed features of randomly sampled datapoints from each domain onto a 2-dimensional space. These features are obtained from the pretrained ResNet10 model (see §B.2 for training details) and are used for our cross-domain experiments. For quantitative visualization, we compute *Proxy A-distance* [3, 14], or PAD, between the base domain (here, *miniImageNet*) and a novel domain as a measure of their closeness in the representation space. To compute PAD, we train a binary classifier over the same ResNet10

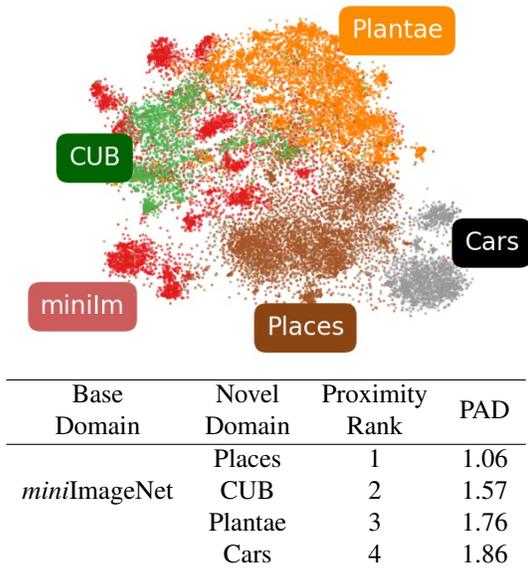


Table 7: **Qualitative and quantitative visualization of the base and novel domains in the cross-domain benchmark [68].** We use t-SNE to visualize the base and novel domains in our cross-domain benchmark. The domain names are presented in boxes with colors that match the corresponding domains in the scatter plot. Here, “miniIm” refers to the *miniImageNet* domain. We also compute the *Proxy A-distance* (PAD) [3, 14] between the base domain and a novel domain as a measure of their relatedness in the representation space. Smaller the PAD value, closer is the novel domain to the base and hence, more related. The PAD values are also used to rank the novel domains according to their proximity to the base domain with the closest domain ranked the highest.

model used for t-SNE but with frozen embedding weights. The classifier distinguishes between randomly drawn samples of the base and novel domains. Denoting  $\epsilon$  as the generalization error of this classifier, the PAD  $\in [0, 2]$  is calculated as  $2(1 - 2\epsilon)$ . Thus, a lower PAD value implies higher generalization error which, in turn, signifies that the base and novel domains are too similar to be distinguished well enough. Finally, the PAD values are used to rank each novel domain, such that the highest rank is assigned to the one closest to the base domain *i.e.*, *miniImageNet*. These ranks correlate well with the t-SNE visualization as well. For instance, CUB and Places, which are ranked higher than Cars and Plantae, are also closer to *miniImageNet* in the t-SNE plot.

---

**Algorithm 1** Distractor-Aware Contrastive Finetuning

---

**Input:** Distractor Dataset ( $D$ ), Prior Model ( $\mathcal{M}_{\theta_0}$ ), few-shot task ( $\tau$ ), Number of Finetuning Epochs ( $J_{\text{ft}}$ ), Augmentation Function ( $\mathcal{A}$ ), Temperature Coefficient ( $\gamma$ ), Learning Rate ( $\eta$ )

**Output:** Finetuned Model Parameters ( $\theta_\tau$ )

- 1: shuffle  $D$
  - 2: **for**  $j \leftarrow 1$  to  $J_{\text{ft}}$  **do**
  - 3: From  $D$ , randomly sample a fixed size batch  $S_{\text{dt}}$  without replacement
  - 4: Using  $\mathcal{A}$  augment each support sample  $x_i, \forall i \in I_{\text{supp}}$
  - 5: For each augmented support sample, define i) anchor-positive index set  $P(i)$ ; ii) anchor-negative index set  $N(i)$  specific to  $\tau$ ; and iii) distractor index set  $I_{\text{dt}}$
  - 6: For all samples, compute  $z_i = \frac{h_i}{\|h_i\|_2}$ , where  $h = \mathcal{M}_\theta(\mathcal{A}(x_i)), \forall i \in I_{\text{supp}}$  and  $h = \mathcal{M}_\theta(x_i), \forall i \in I_{\text{dt}}$
  - 7: Evaluate  $\mathcal{L}_{\text{confit}}(\theta)$  using the quantities computed in previous steps
  - 8: Update model parameters  $\theta \leftarrow \theta - \eta \nabla \mathcal{L}_{\text{confit}}(\theta)$
  - 9: **if**  $j = |D|$  **then**
  - 10: shuffle  $D$
  - 11: **end if**
  - 12: **end for**
- 

## B.2. Prior Learning

As described in the main paper, we use a ResNet10 model [24] as our prior embedding for cross-domain few-shot classification. To avoid specialized hyperparameter tuning while training the prior model, we simply use the pretrained weights<sup>2</sup> made available by [68]. This model was originally trained on all 64 categories of the *miniImageNet train* split.

For the unsupervised prior learning, we train a modified four-layer convolution neural network (CNN), using the recently proposed self-supervised contrastive learning objective [6]. As proposed in [6], we use a 128-dimensional linear projection head on top of the CNN for better generalizability of learnt representations. We train the model with a batch size of 512, temperature coefficient 0.1, and the same augmentation scheme introduced in [6]. Further, we use ADAM optimizer with initial learning rate of  $1e - 3$ , and a weight decay of  $1e - 5$ .

## B.3. Hyperparameter Details

Our proposed contrastive finetuning involves a few hyperparameters such as temperature, learning rate, early-stopping criteria, distractor batch size, and data augmentation scheme. For early-stopping criteria, we set a predetermined range of epochs up to which the pretrained embedding model is finetuned. Here, one finetuning epoch refers to one pass through all the samples of the few-shot task (exclusive of distractors). The range of these epochs along with other hyperparameters are summarized in Table 8 and Table 9. Additionally, we also show the final hyperparameter values used for finetuning in the cross-domain and unsuper-

<sup>2</sup><https://github.com/hytseng0509/CrossDomainFewShot>

vised prior learning settings (the corresponding experiments were reported in the main paper).

## C. Additional Ablations

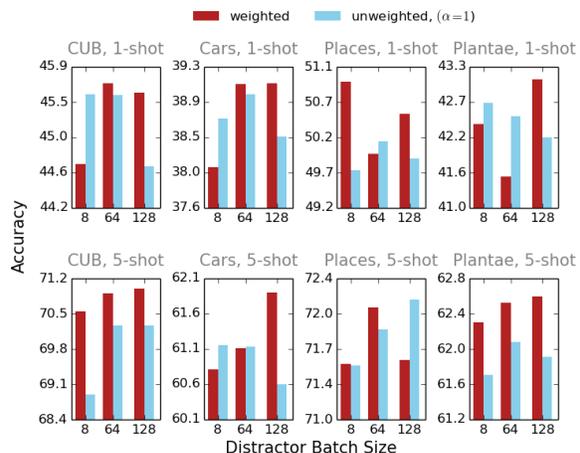


Figure 4: **Comparing the effect of distractor batch size,  $|S_{\text{dt}}|$ , on the weighted and unweighted versions of  $\mathcal{L}_{\text{confit}}$ .** The red and blue bars represent weighted and unweighted versions of  $\mathcal{L}_{\text{confit}}$ , respectively, where  $\alpha$  represents the parameter used to relatively weigh task-specific and task-exclusive (distractors) anchor-negative terms. For each novel domain and shot setting per domain, we compare the performance of two versions in terms of the classification accuracy of unseen samples given a novel task at various distractor batch sizes. These accuracies, as in all other cross-domain experiments, are averaged over 600 randomly chosen novel tasks.

Hyperparameter	Range	CUB		Cars		Places		Plantae	
		1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
learning rate	{5e-4, 5e-3}	5e-3	5e-3	5e-3	5e-3	5e-4	5e-4	5e-3	5e-3
temperature, $\gamma$	{0.05, 0.1, 0.5}	0.1	0.1	0.05	0.05	0.1	0.05	0.1	0.1
distractor batch size, $ S_{dt} $	{64, 128}	64	128	128	128	64	64	128	128
early stopping epoch	{50, 100, 200, 300, 400}	100	100	400	300	200	50	100	100

Table 8: **Hyperparameter details for ConFT with cross-domain prior learning.** This table summarizes the range of various hyperparameters used for finetuning. Additionally, we report the cross-validated values used for the cross-domain prior learning setup. The input image resolution used in this setup is  $224 \times 224$ .

Hyperparameter	Range	<i>miniImageNet</i>	
		1-shot	5-shot
learning rate	{5e-4, 5e-3}	5e-4	5e-4
temperature, $\gamma$	{0.05, 0.1, 0.5}	0.05	0.05
distractor batch size, $ S_{dt} $	{64, 128}	64	64
early-stopping epoch	{50, 100, 200, 300, 400, 500}	400	400

Table 9: **Hyperparameter details for ConFT with unsupervised prior learning.** This table summarizes the range of various hyperparameters used for finetuning. Additionally, we report the cross-validated values used for the unsupervised prior learning setup. The input image resolution used in this setup is  $84 \times 84$ .

Similarity-Pair Construction	Cars	
	1-shot	5-shot
Standard	$37.09 \pm 0.76$	$60.72 \pm 0.74$
Asymmetric (ours)	<b><math>39.11 \pm 0.77</math></b>	<b><math>61.53 \pm 0.75</math></b>

Table 10: **Comparing our proposed asymmetric construction of similarity pairs against standard construction.** Results are shown for both 1-shot and 5-shot tasks sampled from the Cars domain with *miniImageNet* as the base domain. These results are averaged over 600 random novel tasks and are reported with ( $\pm$ ) 95% confidence intervals. Despite using more supervision in the form of distractor labels, the standard pair construction under-performs our (distractor) label-agnostic asymmetric pair construction.

In this section we elucidate the importance of two modifications introduced to the standard contrastive loss, namely, asymmetric construction of similarity pairs and relative weighting of anchor-negative terms.

### C.1. Asymmetric Construction of Similarity Pairs

Our proposed finetuning approach is a general contrastive learning framework for incorporating additional *unlabelled* data in the form of distractors. While construction of positive distractor pairs (that share the same class) is difficult in the absence of distractor labels, constructing anchor-negatives, with anchors being task-specific samples, is much easier following the non-overlapping assumption of task and distractor categories. This results in an asym-

metric construction of similarity pairs where distractors, unlike task-specific samples, can meaningfully participate only as anchor-negatives. In fact, this asymmetry is critical in the unsupervised prior learning setup, where distractors are sampled from an unlabelled base domain. In the case of cross-domain prior learning, however, we have a labelled base data as a source for distractors. To motivate our asymmetric pair construction in this case, we compare it to a standard construction that allows distractors to additionally participate as anchor-positives. To form such an anchor-positive, a distractor is paired with another distractor sharing the same class. Here, anchor-negatives with respect to a distractor include all the datapoints that do not share the class with it. This includes samples from both the novel few-shot task and other distractors. Overall, the resulting form of the contrastive loss can be viewed as applying supervised contrastive objective [31] (without augmentation-based positives) to the union set of task samples and distractors within a training batch. In Table 10, we evaluate these two types of pair constructions on the cross-domain setting, *miniImageNet*  $\rightarrow$  Cars. Interestingly enough, our formulation of the contrastive loss with asymmetric pair construction yields superior performance despite using less supervision than the supervised contrastive loss.

### C.2. Importance of Weighted Negatives

Another important component of our loss is the relative weighting parameter  $\alpha$  that balances the effect of task-specific and task-exclusive (distractor based) anchor-negative terms. To validate the utility of such a weighting scheme, we compare the weighted version of  $\mathcal{L}_{\text{conft}}$  to its unweighted version *i.e.*,  $\alpha = 1$ . Following the results for various novel domains and shot settings in Figure 4, we make the following observations. The weighted loss (*red* bars) performance improves with larger distractor batch sizes in most cases (5 out of 8). The improvement is more pronounced for domains like Cars and Plantae that are farther away from the base dataset - *miniImageNet* (see Table 7). For closer domains like CUB or Places, we sometimes notice a sweet spot at batch size = 64. In contrast, the unweighted version (*blue* bars) experiences a performance drop with increasing batch sizes, when the novel

domains are farther from the base domain. In other cases, the trends are inconclusive. The most important observation, however, comes from comparing the two versions of the loss. Specifically, the weighted version not only outperforms the unweighted loss at higher batch sizes but also results in the best performance in almost every setting. The only exception is Places, 5-shot where the unweighted loss yields the best performance. A possible explanation is as follows: due to the similarity of Places (novel domain) and *mini*ImageNet (base domain) in the embedding space (see Table 7), distractor samples from Places may serve as hard negatives that are important for effective contrastive learning [53]. Thus, down-weighting their contribution at higher batch sizes would degrade the final performance.

### C.3. Data Augmentation

Augmentation		CUB	Cars
Task Samples	Distractors	5-shot	5-shot
-	-	69.90 ± 0.75	58.64 ± 0.88
✓	-	70.53 ± 0.75	61.53 ± 0.75

Table 11: In this ablation we compare the few-shot performance when a prior embedding is finetuned (using ConFT) with or without augmentation to task-specific samples. Note that, we never use augmentation for distractors in our experiments.

Yet another important component of our contrastive finetuning objective is the data augmentation function  $\mathcal{A}$ . To avoid extensive tuning of large hyperparameter space associated with  $\mathcal{A}$ , we adopt a fixed augmentation strategy introduced in [7]. In Table 11, we show the benefit of using this strategy to augment samples specific to the novel task. Following preliminary investigations, we found that augmenting distractors did not make much difference. Hence, we never apply data augmentation to distractors in our experiments.

### C.4. Loss Type

In Table 12, we compare contrastive and cross-entropy finetuning in conjunction with the auxiliary cross-entropy objective (MT). While the two objectives yield similar performance for the CUB case, contrastive finetuning outperforms cross-entropy loss based finetuning in Cars. These results show that the contrastive loss could be a better choice for few-shot classification.

## D. ConFT as a General Finetuning Approach

In Table 13, we validate the complementary effect of our finetuning approach to a variety of prior learning schemes. Specifically, we compare our simple cross-entropy objective with ProtoNet [61] and ProtoNet with auxiliary self-

Prior Learning	Method	CUB	Cars
	Task Specific Finetuning	5-shot	5-shot
CE Training	MT-ceFT ( $\beta = 1$ )	71.35 ± 0.70	58.97 ± 0.76
CE Training	MT-ceFT ( $\beta = 10$ )	74.32 ± 0.69	60.01 ± 0.74
CE Training	MT-ConFT ( $\beta = 1$ )	71.65 ± 0.74	61.25 ± 0.70
CE Training	MT-ConFT ( $\beta = 10$ )	74.45 ± 0.71	62.54 ± 0.72

Table 12: **Ablation.** Cross-entropy/contrastive finetuning with a multi-task (MT) cross entropy objective. Here, all cross entropy objectives are based on cosine classifier with a multiplying factor,  $\beta$

supervision [63]. Both of these approaches are based on meta-learning, and were originally proposed for in-domain few-shot classification where base and novel tasks follow the same distribution. Nevertheless, the embeddings thus learnt are readily applicable to cross-domain tasks as well. For the auxiliary self-supervision, we use image rotation as our pretext task. While previous work [63] has demonstrated the improvement in in-domain few-shot generalization resulting from rotation based self-supervision, we found that the improvement is marginal in our cross-domain setting (see ProtoNet without finetuning vs. ProtoNet + Rot. without finetuning in Table 13), except for when the novel domain is Cars. To obtain these results, we use the official implementation<sup>3</sup> of [63] with the same hyperparameters (such as loss weighting term) but different backbone. As our pretrained embedding, we trained a ProtoNet model (with auxiliary self-supervision) based on ResNet10 [24] architecture. Our main observation from Table 13 is as follows: while better prior learning objectives such as those with auxiliary self-supervision can improve few-shot classification in the novel domains, finetuning with ConFT consistently leads to large improvements over the prior embeddings.

## E. Additional Comparison with Prior Work

In Table 14, we report additional comparison with a concurrent work SCL [44] that introduces attention-based spatial contrastive objective in the prior-learning phase. For a fair comparison to SCL, we adopt the same backbone based on the ResNet12 architecture which was originally proposed in [66]. While the spatial contrastive objective benefits from larger image resolution ( $224 \times 224$ ), we found it significantly increases the time for finetuning in our case, especially given the larger backbone. So, in this case, we conduct our experiments with a smaller resolution of  $84 \times 84$ . Despite the drop in resolution, our finetuning based approach over simple cross-entropy prior learning outperforms the more sophisticated SCL by significant margins in CUB (7 points) and Cars (13 points). While we attain similar performance in the case of Plantae, we underperform in Places domain. This gap can be understood as a

<sup>3</sup>[https://github.com/cvl-umass/fsl\\_ssl](https://github.com/cvl-umass/fsl_ssl)

Method			5-shot			
Prior Learning	Task Specific Finetuning	Backbone	CUB	Cars	Places	Plantae
ProtoNet [61]	-	ResNet10	58.80 ± 0.77	44.07 ± 0.69	71.03 ± 0.72	51.33 ± 0.72
ProtoNet [61]	ConFT (ours)	ResNet10	<b>66.63 ± 0.69</b>	<b>59.27 ± 0.73</b>	<b>72.05 ± 0.71</b>	<b>58.83 ± 0.76</b>
ProtoNet + Rot. [63]	-	ResNet10	58.68 ± 0.75	46.48 ± 0.71	71.20 ± 0.75	51.93 ± 0.67
ProtoNet + Rot. [63]	ConFT (ours)	ResNet10	<b>66.75 ± 0.71</b>	<b>61.67 ± 0.75</b>	<b>73.91 ± 0.70</b>	<b>60.38 ± 0.75</b>
CE Training	-	ResNet10	62.80 ± 0.76	51.41 ± 0.72	70.71 ± 0.68	55.54 ± 0.69
CE Training	ConFT (ours)	ResNet10	<b>70.53 ± 0.75</b>	<b>61.53 ± 0.75</b>	<b>72.09 ± 0.68</b>	<b>62.54 ± 0.76</b>

Table 13: **Combining ConFT with different pretraining schemes for cross-domain prior learning.** We present the results for 5-way 5-shot tasks averaged over 600 such tasks with ( $\pm$ ) 95% confidence intervals. The highlighted numbers demonstrate that ConFT consistently improves the few-shot performance of prior embeddings across data domains.

Method			1-shot			
Prior Learning	Task Specific Finetuning	Backbone	CUB	Cars	Places	Plantae
SCL [44]	-	ResNet12	50.09 ± 0.7	34.93 ± 0.6	<b>60.32 ± 0.8</b>	40.23 ± 0.6
CE Training	-	ResNet12	50.00 ± 0.77	34.88 ± 0.64	55.62 ± 0.91	38.47 ± 0.72
CE Training	ConFT (ours)	ResNet12	<b>52.01 ± 0.82</b>	<b>39.54 ± 0.68</b>	56.66 ± 0.85	<b>40.90 ± 0.73</b>

Method			5-shot			
Prior Learning	Task Specific Finetuning	Backbone	CUB	Cars	Places	Plantae
SCL [44]	-	ResNet12	68.81 ± 0.6	52.22 ± 0.7	<b>76.51 ± 0.6</b>	<b>59.91 ± 0.6</b>
CE Training	-	ResNet12	69.75 ± 0.73	49.92 ± 0.74	73.79 ± 0.67	54.66 ± 0.77
CE Training	ConFT (ours)	ResNet12	<b>76.49 ± 0.63</b>	<b>64.87 ± 0.70</b>	74.22 ± 0.71	<b>59.23 ± 0.77</b>

Table 14: **Additional Prior Work Comparison.** SCL introduces a novel attention-based spatial contrastive objective for prior learning. While we employ a much simpler cross-entropy objective for prior learning (see CE training *without* ConFT), finetuning the prior embedding with ConFT outperforms SCL significantly in two (CUB and Cars) out of four domains. Our approach yields competitive results for Plantae as well. Further, due to the complementary nature of finetuning, the best performance might be achieved by combining SCL with our ConFT.

consequence of a stronger SCL based prior embedding for *mini*ImageNet and greater similarity of the *mini*ImageNet domain to Places as opposed to other novel domains (see Table 7). Nonetheless, our finetuning is complimentary to SCL, and hence we suspect that the best performance could be achieved by combining it with our ConFT.

## F. Meta-Dataset Results

In this section, we present the results of our ConFT approach on Meta-Dataset (see Table 15). Here, we use an off-the-shelf ResNet18 model<sup>4</sup> pretrained on ImageNet-train-split of Meta-Dataset using just cross-entropy objective. In order to maintain consistency with pretraining, our finetuning operates at a small image resolution of  $84 \times 84$ . In this experiments, we keep most of the hyperparameters fixed across all datasets. In particular, we use a temperature of 0.1, a distractor batch size of 128, and a learning rate of  $5e - 5$ . The early stopping epoch is cross-validated using the meta-validation splits of respective datasets. We observe that our approach outperforms the state of the art in 7 out of 10 datasets and sometimes by a significant margin.

This is despite the fact that our input resolution is much smaller compared to  $224 \times 224$  in the state of the art and our approach does *not* benefit from a transductive setting. Finally, our results reinforce the superiority of simple finetuning over more complex meta-learning frameworks (*e.g.* cross-attention based) even when the domain gap is large.

<sup>4</sup><https://github.com/peymanbateni/simple-cnaps>

Method	Target Datasets				
	ILSVRC	Omni	Aircraft	Birds	DTD
PN [11]	41.87 $\pm$ 0.89	61.33 $\pm$ 1.13	39.40 $\pm$ 0.78	65.57 $\pm$ 0.73	59.06 $\pm$ 0.60
CTX [11]	51.70 $\pm$ 0.90	84.24 $\pm$ 0.79	62.29 $\pm$ 0.73	<b>79.38<math>\pm</math>0.54</b>	<b>65.86<math>\pm</math>0.58</b>
CTX+SC [11]	51.29 $\pm$ 0.89	86.14 $\pm$ 0.74	<b>69.74<math>\pm</math>0.67</b>	74.85 $\pm$ 0.62	63.84 $\pm$ 0.62
CTX+SC+Aug [11]	52.56 $\pm$ 0.86	87.53 $\pm$ 0.61	64.28 $\pm$ 0.71	73.27 $\pm$ 0.63	64.72 $\pm$ 0.63
ConFT (ours)	<b>72.07<math>\pm</math>0.71</b>	<b>98.22<math>\pm</math>0.17</b>	68.44 $\pm$ 0.70	74.93 $\pm$ 0.67	63.11 $\pm$ 0.70

Method	Target Dataset				
	QDraw	Fungi	Flower	Sign	COCO
PN [11]	47.86 $\pm$ 0.80	41.64 $\pm$ 1.02	83.88 $\pm$ 0.48	44.84 $\pm$ 0.88	41.14 $\pm$ 0.82
CTX [11]	63.36 $\pm$ 0.73	49.43 $\pm$ 0.98	92.74 $\pm$ 0.29	68.31 $\pm$ 0.71	48.63 $\pm$ 0.79
CTX+SC [11]	64.11 $\pm$ 0.67	48.87 $\pm$ 0.91	93.00 $\pm$ 0.30	70.62 $\pm$ 0.68	48.45 $\pm$ 0.83
CTX+SC+Aug [11]	66.90 $\pm$ 0.66	48.22 $\pm$ 0.94	93.23 $\pm$ 0.28	78.45 $\pm$ 0.60	56.61 $\pm$ 0.78
ConFT (ours)	<b>80.02<math>\pm</math>0.6</b>	<b>50.16<math>\pm</math>0.80</b>	<b>94.52<math>\pm</math>0.29</b>	<b>88.22<math>\pm</math>0.59</b>	<b>70.73<math>\pm</math>0.79</b>

Table 15: **Meta-Dataset Results (5-shot)**. Cross-domain results of our distractor-aware contrastive finetuning (ConFT) on transfer from ImageNet-only are presented here. The accuracies are averaged over 600 evaluation tasks with 95% confidence intervals. PN: Prototypical Net, SC: SimCLR Episodes.