# When Does Self-supervision Improve Few-shot Learning?

Jong-Chyi Su[1]      Subhransu Maji[1]      Bharath Hariharan[2]

[1] University of Massachusetts Amherst
[2] Cornell University
{jcsu,smaji}@cs.umass.edu, bharathh@cs.cornell.edu

**Abstract.** We investigate the role of self-supervised learning (SSL) in the context of few-shot learning. Although recent research has shown the benefits of SSL on large unlabeled datasets, its utility on small datasets is relatively unexplored. We find that SSL reduces the relative error rate of few-shot meta-learners by 4%-27%, even when the datasets are small and *only* utilizing images within the datasets. The improvements are greater when the training set is smaller or the task is more challenging. Although the benefits of SSL may increase with larger training sets, we observe that SSL can hurt the performance when the distributions of images used for meta-learning and SSL are different. We conduct a systematic study by varying the degree of domain shift and analyzing the performance of several meta-learners on a multitude of domains. Based on this analysis we present a technique that automatically selects images for SSL from a large, generic pool of unlabeled images for a given dataset that provides further improvements.

## 1   Introduction

Current machine learning algorithms require enormous amounts of training data to learn new tasks. This is an issue for many practical problems across domains such as biology and medicine where labeled data is hard to come by. In contrast, we humans can quickly learn new concepts from limited training data by relying on our past "visual experience". Recent work attempts to emulate this by training a feature representation to classify a training dataset of "base" classes with the hope that the resulting representation generalizes not just to unseen examples of the same classes but also to novel classes, which may have very few training examples (called few-shot learning). However, training for base class classification can force the network to only encode features that are useful for distinguishing between base classes. In the process, it might discard semantic information that is irrelevant for base classes but critical for novel classes. This might be especially true when the base dataset is small or when the class distinctions are challenging.

One way to recover this useful semantic information is to leverage representation learning techniques that do not use class labels, namely, *unsupervised* or *self-supervised learning*. The key idea is to learn about statistical regularities within images, such as the spatial relationship between patches, or its orientation,
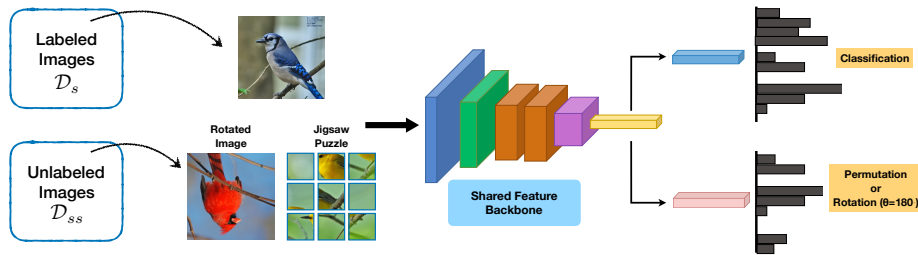
Fig. 1: **Combining supervised and self-supervised losses for few-shot learning.** Self-supervised tasks such as jigsaw puzzle or rotation prediction act as a data-dependent regularizer for the shared feature backbone. Our work investigates how the performance on the *target task domain* ($\mathcal{D}_s$) is impacted by the choice of the *domain used for self-supervision* ($\mathcal{D}_{ss}$).

that might be a cue to semantics. Despite recent advances, these techniques have only been applied to a few domains (*e.g.*, entry-level classes on internet imagery), and under the assumption that large amounts of unlabeled images are available. Their applicability to the general few-shot scenario is unclear. In particular, can these techniques prevent overfitting to base classes and improve performance on novel classes in the few-shot setting? If so, does the benefit generalize across domains and to more challenging tasks? Moreover, can self-supervision boost performance in domains where even unlabeled images are hard to get?

This paper seeks to answer these questions. We show that with *no additional training data*, adding a self-supervised task as an auxiliary task (Fig. 1) improves the performance of existing few-shot techniques on benchmarks across a multitude of domains (Fig. 2), in agreement with conclusions from similar recent work [18]. Intriguingly, we find that the benefits of self-supervision *increase* with the difficulty of the task, for example when training from a smaller base dataset, or with degraded inputs such as low resolution or greyscale images (Fig. 3).

One might surmise that as with traditional SSL, additional unlabeled images might improve performance further. But what unlabeled images should we use for novel problem domains where unlabeled data is not freely available? To answer this, we conduct a series of experiments with additional unlabeled data from different domains. We find that adding more unlabeled images improves performance *only* when the images used for self-supervision are within the *same domain* as the base classes (Fig. 4a); otherwise, they can even *negatively* impact the performance of the few-shot learner (Fig. 4b). Based on this analysis, we present a simple approach that uses a domain classifier to pick similar-domain unlabeled images for self-supervision from a large and generic pool of images (Fig. 5). The resulting method improves over the performance of a model trained with self-supervised learning from images within the dataset (Fig. 6). Taken together, this results in a powerful, general, and practical approach for improving few-shot learning on small datasets in novel domains. Finally, these benefits are also observed in standard classification tasks (Appendix A.3).

## 2 Related Work

**Few-shot learning** Few-shot learning aims to learn representations that generalize well to the novel classes where only a few images are available. To this end, several meta-learning approaches have been proposed that evaluate representations by sampling many few-shot tasks within the domain of a *base* dataset. These include optimization-based meta-learners, such as model-agnostic meta-learner (MAML) [16], gradient unrolling [49], closed-form solvers [4], and convex learners [35]. The second class of methods rely on distance-based classifiers such as matching networks [61] and prototypical networks (ProtoNet) [55]. Another class of methods [19, 47, 48] model the mapping between training data and classifier weights using a feed-forward network.

While the literature is rapidly growing, a recent study by Chen *et al*. [10] has shown that the differences between meta-learners are diminished when deeper networks are used. They develop a strong baseline for few-shot learning and show that the performance of ProtoNet [55] matches or surpasses several recently proposed meta-learners. We build our experiments on top of this work and show that auxiliary self-supervised tasks provide additional benefits across a large array of few-shot benchmarks and across meta-learners.

**Self-supervised learning** Human labels are expensive to collect and hard to scale up. To this end, there has been increasing research interest to investigate learning representations from unlabeled data. In particular, the image itself already contains structural information that can be utilized. One class of methods remove part of the visual data and task the network with predicting what has been removed from the rest in a discriminative manner [34, 46, 59, 68, 69]. Another line of works treat each image (and augmentations of itself) as one class and use contrastive learning as self-supervision [3, 5, 9, 15, 22, 24, 25, 38, 43, 65]. Other self-supervised tasks include predicting rotation [20], relative patch location [13], clusters [7, 8], and number of objects [42], *etc*.

On top of those SSL tasks, combining different tasks can be beneficial [14], in this work we also see its benefit. Asano *et al*. [2] showed that the representations can be learned with only one image and extreme augmentations. We also investigate SSL on a low-data regime, but use SSL as a regularizer instead of a pre-training task. Goyal *et al*. [21] and Kolesnikov *et al*. [31] compared various SSL tasks at scale and concluded that solving jigsaw puzzles and predicting image rotations are among the most effective, motivating the choice of tasks in our experiments. Note that these two works did not include a comparison with contrastive learning approaches.

In addition to pre-training models, SSL can also be used to improve other tasks. For example, Zhai *et al*. [67] showed that self-supervision can be used to improve recognition in a semi-supervised setting and presented results on a partially labeled version of the ImageNet dataset. Carlucci *et al*. [6] used self-supervision to improve domain generalization. In this work we use SSL to improve few-shot learning where the goal is to generalize to novel classes.

However, the focus of most prior works on SSL is to supplant traditional supervised representation learning with unsupervised learning on large unlabeled datasets for downstream tasks. Crucially in almost all prior works, self-supervised representations consistently lag behind fully-supervised ones trained on the same dataset with the same architecture [21, 31]. *In contrast, our work focuses on an important counterexample:* self-supervision can in fact augment standard supervised training for few-shot transfer learning in the low training data regime *without* relying on any external dataset.

The most related work is that of Gidaris *et al.* [18] who also use self-supervision to improve few-shot learning. Although the initial results are similar (Table 3), we further show these benefits on several datasets with harder recognition problems (fine-grained classification) and with deeper models (Section 4.1). Moreover, we present a novel analysis of the impact of the domain of unlabeled data (Section 4.2). Finally, we propose a new and simple approach to automatically select similar-domain unlabeled data for self-supervision (Section 4.3).

**Multi-task learning** Our work is related to multi-task learning, a class of techniques that train on multiple task objectives together to improve each one. Previous works in the computer vision literature have shown moderate benefits by combining tasks such as edge, normal, and saliency estimation for images, or part segmentation and detection for humans [30, 37, 51]. However, there is significant evidence that training on multiple tasks together often hurts performance on individual tasks [30, 37]. Only certain task combinations appear to be mutually beneficial, and sometimes specific architectures are needed. Our key contribution here is showing that self-supervised tasks and few-shot learning are indeed mutually beneficial in this sense.

**Domain selection** On supervised learning, Cui *et al.* [12] used Earth Mover's distance to measure the domain similarity and select the source domain for pre-training. Ngiam *et al.* [39] found more data for pre-training does not always help and proposed to use importance weights to select pre-training data. Task2vec [1] generates a task embedding given a probe network and the target dataset. Such embeddings can help select a better pre-training model from a pool of experts which yields better performance after fine-tuning. Unlike these, we do not assume that the source domain is labeled and rely on self-supervised learning. On self-supervised learning, Goyal *et al.* [21] used two pre-training and target datasets to show the importance of source domain on large-scale self-supervised learning. Unlike this work, we investigate the performance on few-shot learning across a number of domains, as well as investigate methods for domain selection. A concurrent work [62] also investigates the effect of domain shifts on SSL.

## 3   Method

We adopt the commonly used setup for few-shot learning where one is provided with labeled training data for a set of *base* classes $\mathcal{D}_b$ and a much smaller training

set (typically 1-5 examples per class) for *novel* classes $\mathcal{D}_n$. The goal of the few-shot learner is to learn representations on the base classes that lead to good generalization on novel classes. Although in theory the base classes are assumed to have a large number of labeled examples, in practice this number can be quite small for novel or fine-grained domains, *e.g.* less than 5000 images for the birds dataset [63], making it challenging to learn a generalizable representation.

Our framework, as seen in Fig. 1, combines *meta-learning* approaches for few-shot learning with *self-supervised learning*. Denote a labeled training dataset $\mathcal{D}_s$ as $\{(x_i, y_i)\}_{i=1}^n$ consisting of pairs of images $x_i \in \mathcal{X}$ and labels $y_i \in \mathcal{Y}$. A feed-forward convolutional network $f(x)$ maps the input to an embedding space which is then mapped to the label space using a classifier $g$. The overall mapping from the input to the label can be written as $g \circ f(x) : \mathcal{X} \to \mathcal{Y}$. Learning consists of estimating functions $f$ and $g$ that minimize an empirical loss $\ell$ over the training data along with suitable regularization $\mathcal{R}$ over the functions $f$ and $g$. This can be written as:

$$\mathcal{L}_s := \sum_{(x_i,y_i) \in \mathcal{D}_s} \ell\big(g \circ f(x_i), y_i\big) + \mathcal{R}(f, g).$$

A commonly used loss is the cross-entropy loss and a regularizer is the $\ell_2$ norm of the parameters of the functions. In a transfer learning setting $g$ is discarded and relearned on training data for novel classes.

We also consider self-supervised losses $\mathcal{L}_{ss}$ based on labeled data $x \to (\hat{x}, \hat{y})$ that can be derived automatically without any human labeling. Fig. 1 shows two examples: the *jigsaw task* rearranges the input image and uses the index of the permutation as the target label, while the *rotation task* uses the angle of the rotated image as the target label. A separate function $h$ is used to predict these labels from the shared feature backbone $f$ with a self-supervised loss:

$$\mathcal{L}_{ss} := \sum_{x_i \in \mathcal{D}_{ss}} \ell\big(h \circ f(\hat{x}_i), \hat{y}_i\big).$$

Our final loss combines the two: $\mathcal{L} := \mathcal{L}_s + \mathcal{L}_{ss}$ and thus the self-supervised losses act as a data-dependent regularizer for representation learning. The details of these losses are described in the next sections.

Note that the domain of images used for supervised $\mathcal{D}_s$ and self-supervised $\mathcal{D}_{ss}$ losses need not to be identical. In particular, we would like to use larger sets of images for self-supervised learning from related domains. The key questions we ask are: (1) How effective is SSL when $\mathcal{D}_s = \mathcal{D}_{ss}$ especially when we have a small sample of $D_s$? (2) How do the domain shifts between $\mathcal{D}_s$ and $\mathcal{D}_{ss}$ affect generalization performance? and (3) How to select images from a large, generic pool to construct an effective $\mathcal{D}_{ss}$ given a target domain $\mathcal{D}_s$?

### 3.1 Supervised Losses ($\mathcal{L}_s$)

Most of our results are presented using a meta-learner based on prototypical networks [55] that perform episodic training and testing over sampled datasets in stages called meta-training and meta-testing. During meta-training, we randomly

sample $N$ classes from the base set $\mathcal{D}_b$, then we select a support set $\mathcal{S}_b$ with $K$ images per class and another query set $\mathcal{Q}_b$ with $M$ images per class. We call this an $N$-way $K$-shot classification task. The embeddings are trained to predict the labels of the query set $\mathcal{Q}_b$ conditioned on the support set $\mathcal{S}_b$ using a nearest mean (prototype) classifier. The objective is to minimize the prediction loss on the query set. Once training is complete, given the novel dataset $\mathcal{D}_n$, class prototypes are recomputed for classification and query examples are classified based on the distances to the class prototypes.

Prototypical networks are related to distance-based learners such as matching networks [61] or metric-learning based on label similarity [29]. We also present few-shot classification results using a gradient-based meta-learner called MAML [16], and one trained with a standard cross-entropy loss on all the base classes. We also present standard classification results where the test set contains images from the same base categories in Appendix A.3.

### 3.2   Self-supervised Losses ($\mathcal{L}_{ss}$)

We consider two losses motivated by a recent large-scale comparison of the effectiveness of self-supervised learning tasks [21] described below:

- *Jigsaw puzzle task loss.* Here the input image $x$ is tiled into $3{\times}3$ regions and permuted randomly to obtain an input $\hat{x}$. The target label $\hat{y}$ is the index of the permutation. The index (one of 9!) is reduced to one of 35 following the procedure outlined in [41], which grouped the possible permutations based on the hamming distance to control the difficulty of the task.
- *Rotation task loss.* We follow the method of [20] where the input image $x$ is rotated by an angle $\theta \in \{0°, 90°, 180°, 270°\}$ to obtain $\hat{x}$ and the target label $\hat{y}$ is the index of the angle.
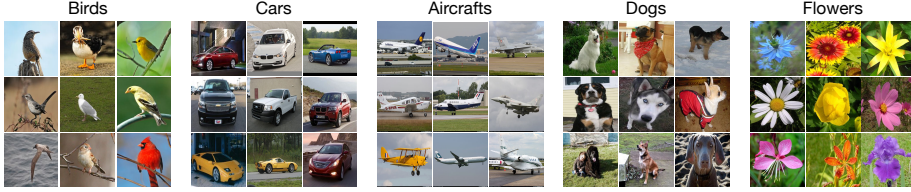
In both cases we use the cross-entropy loss between the target and prediction.

### 3.3   Stochastic Sampling and Training

When the images used for SSL and meta-learning are identical, *i.e.*, $\mathcal{D}_s = \mathcal{D}_{ss}$, the same batch of images are used for computing both losses $\mathcal{L}_s$ and $\mathcal{L}_{ss}$. For experiments investigating the effect of domain shifts described in Section 4.2 and 4.3, where SSL and meta-learner are trained on different domains, *i.e.* $\mathcal{D}_s \neq \mathcal{D}_{ss}$, a separate batch of size of 64 is used for computing $\mathcal{L}_{ss}$. After the two forward passes, one for the supervised task and one for the self-supervised task, the two losses are combined and gradient updates are performed. While other techniques exist [11, 26, 54], simply averaging the two losses performed well.

## 4   Experiments

We first describe the datasets and experimental details. In Section 4.1, we present the results of using SSL to improve few-shot learning on various datasets. In

| Setting | Set | Stats | mini-ImageNet | tiered-ImageNet | Birds | Cars | Aircrafts | Dogs | Flowers |
|---|---|---|---|---|---|---|---|---|---|
| Few-shot transfer | Base | classes | 64 | 351 | 100 | 98 | 50 | 60 | 51 |
| | | images | 38,400 | 448,695 | 5885 | 8162 | 5000 | 10337 | 4129 |
| | Val | classes | 16 | 97 | 50 | 49 | 25 | 30 | 26 |
| | | images | 9,600 | 124,261 | 2950 | 3993 | 2500 | 5128 | 2113 |
| | Novel | classes | 20 | 160 | 50 | 49 | 25 | 30 | 25 |
| | | images | 12,000 | 206,209 | 2953 | 4030 | 2500 | 5115 | 1947 |

Table 1: **Example images and dataset statistics**. For few-shot learning experiments the classes are split into *base*, *val*, and *novel* set. Image representations learned on *base* set are evaluated on the *novel* set while *val* set is used for cross-validation. These datasets vary in the number of classes but are orders of magnitude smaller than ImageNet dataset.

Section 4.2, we show the effect of domain shift between labeled and unlabeled data for SSL. Last, we propose a way to select images from a pool for SSL to further improve the performance of few-shot learning in Section 4.3.

**Datasets and benchmarks** We experiment with datasets across diverse domains: Caltech-UCSD birds [63], Stanford cars [32], FGVC aircrafts [36], Stanford dogs [27], and Oxford flowers [40]. Each dataset contains between 100 and 200 classes with a few thousands of images. We also experiment with the widely-used *mini*-ImageNet [61] and *tiered*-ImageNet [50] benchmarks for few-shot learning. In *mini*-ImageNet, each class has 600 images, wherein *tiered*-ImageNet each class has 732 to 1300 images.

We split classes within a dataset into three disjoint sets: *base, val*, and *novel*. For each class, all the images in the dataset are used in the corresponding set. A model is trained on the base set of categories, validated on the val set, and tested on the novel set of categories given a few examples per class. For birds, we use the same split as [10], where {*base, val, novel*} sets have {100, 50, 50} classes respectively. The same ratio is used for the other four fine-grained datasets. We follow the original splits for *mini*-ImageNet and *tiered*-ImageNet. The statistics of various datasets used in our experiments are shown in Table 1. Notably, fine-grained datasets are significantly smaller.

We also present results on a setting where the base set is "degraded" either by (1) reducing the resolution, (2) removing color, or (3) reducing the number

of training examples. This allows us to study the effectiveness of SSL on even smaller datasets and as a function of the difficulty of the task.

**Meta-learners and feature backbone** We follow the best practices and use the codebase for few-shot learning described in [10]. In particular, we use ProtoNet [55] with a ResNet-18 [23] network as the feature backbone. Their experiments found this to be the best performing. We also present experiments with other meta-learners such as MAML [16] and softmax classifiers in Section 4.1.

**Learning and optimization** We use 5-way (classes) and 5-shot (examples per-class) with 16 query images for training. For experiments using 20% of labeled data, we use 5 query images for training since the minimum number of images per class is 10. The models are trained with ADAM [28] with a learning rate of 0.001 for 60,000 episodes. We report the mean accuracy and 95% confidence interval over 600 test experiments. In each test episode, $N$ classes are selected from the novel set, and for each class 5 support images and 16 query images are selected. We report results for $N = \{5, 20\}$ classes.

**Image sampling and data augmentation** Data augmentation has a significant impact on few-shot learning performance. We follow the data augmentation procedure outlined in [10] which resulted in a strong baseline performance. For label and rotation predictions, images are first resized to 224 pixels for the shorter edge while maintaining the aspect ratio, from which a central crop of 224×224 is obtained. For jigsaw puzzles, we first randomly crop 255×255 region from the original image with random scaling between [0.5, 1.0], then split into 3×3 regions, from which a random crop of size 64×64 is picked. While it might appear that with self-supervision the model effectively sees more images, SSL provides consistent improvements even after extensive data augmentation including cropping, flipping, and color jittering. More experimental details are in Appendix A.5.

**Other experimental results** In Appendix A.3, we show the benefits of using self-supervision for *standard* classification tasks when training the model *from scratch*. We further visualize these models in Appendix A.4 to show that models trained with self-supervision tend to avoid accidental correlation of background features to class labels.

### 4.1   Results on Few-shot Learning

**Self-supervised learning improves few-shot learning** Fig. 2 shows the accuracies of various models on few-shot learning benchmarks. Our ProtoNet baseline matches the results of the *mini*-ImageNet and birds datasets presented in [10] (in their Table A5). Our results show that jigsaw puzzle task improves the ProtoNet baseline on all seven datasets. Specifically, it reduces the *relative error rate* by 4.0%, 8.7%, 19.7%, 8.4%, 4.7%, 15.9%, and 27.8% on *mini*-ImageNet,
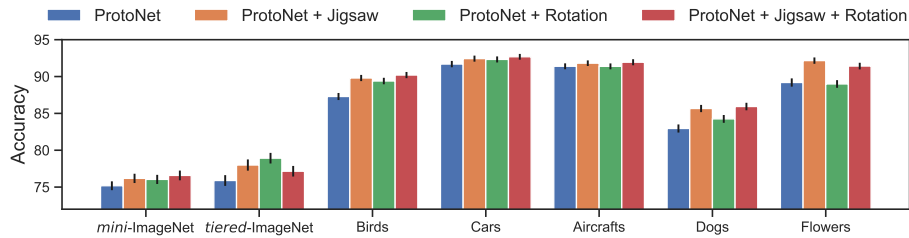
Fig. 2: **Benefits of SSL for few-shot learning tasks.** We show the accuracy of the ProtoNet baseline of using different SSL tasks. The jigsaw task results in an improvement of the 5-way 5-shot classification accuracy across datasets. Combining SSL tasks can be beneficial for some datasets. Here SSL was performed on images within the base classes only. See Appendix A.1 for a tabular version and results for 20-way 5-shot classification.
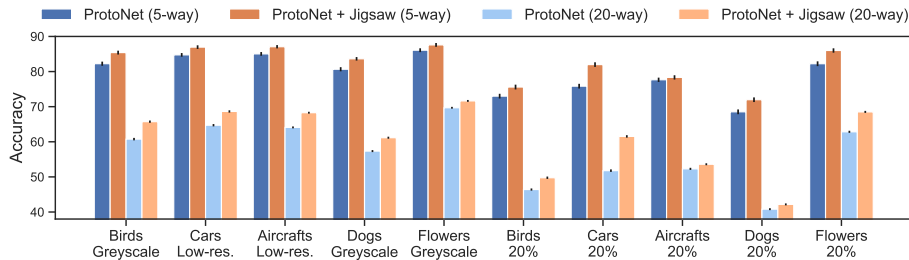


Fig. 3: **Benefits of SSL for *harder* few-shot learning tasks.** We show the accuracy of using the jigsaw puzzle task over ProtoNet baseline on harder versions of the datasets. We see that SSL is effective even on smaller datasets and the relative benefits are higher.

*tiered*-ImageNet, birds, cars, aircrafts, dogs, and flowers datasets respectively. Predicting rotations also improves the ProtoNet baseline on most of the datasets, except for aircrafts and flowers. We speculate this is because most flower images are symmetrical, and airplanes are usually horizontal, making the rotation task too hard or too trivial respectively to benefit the main task. In addition, combining these two SSL tasks can be beneficial sometimes. A tabular version and the results of 20-way classification are included in Appendix A.1.

**Gains are larger for harder tasks** Fig. 3 shows the performance on the degraded version of the same datasets (first five groups). For cars and aircrafts we use low-resolution images where the images are down-sampled by a factor of *four* and up-sampled back to 224×224 with bilinear interpolation. For natural categories we discard color. Low-resolution images are considerably harder to classify for man-made categories while color information is most useful for

| Loss | Birds | Cars | Aircrafts | Dogs | Flowers |
|------|-------|------|-----------|------|---------|
|  | 5-way 5-shot | | | | |
| Softmax | 81.5±0.5 | 87.7±0.5 | 89.2±0.4 | 77.6±0.6 | 91.0±0.5 |
| Softmax + Jigsaw | 83.9±0.5 | 90.6±0.5 | 89.6±0.4 | 77.8±0.6 | 91.1±0.5 |
| MAML | 81.2±0.7 | 86.9±0.6 | 88.8±0.5 | 77.3±0.7 | 79.0±0.9 |
| MAML + Jigsaw | 81.1±0.7 | 89.0±0.5 | 89.1±0.5 | 77.3±0.7 | 82.6±0.7 |
| ProtoNet | 87.3±0.5 | 91.7±0.4 | 91.4±0.4 | 83.0±0.6 | 89.2±0.6 |
| ProtoNet + Jigsaw | **89.8±0.4** | **92.4±0.4** | **91.8±0.4** | **85.7±0.5** | **92.2±0.4** |

Table 2: **Performance on few-shot learning using different meta-learners.** Using jigsaw puzzle loss improves different meta-learners on most of the datasets. ProtoNet with jigsaw loss performs the best on all five datasets.

natural categories [56]. On birds and dogs datasets, the improvements using self-supervision (3.2% and 2.9% on 5-way 5-shot) are higher compared to color images (2.5% and 2.7%), similarly on the cars and aircrafts datasets with low-resolution images (2.2% and 2.1% vs. 0.7% and 0.4%). We also conduct an experiment where only 20% of the images in the base categories are used for both SSL and meta-learning (last five groups in Fig. 3). This results in a much smaller training set than standard few-shot benchmarks: 20% of the birds dataset amounts to only roughly 3% of the popular *mini*-ImageNet dataset. We find larger benefits from SSL in this setting. For example, the gain from the jigsaw puzzle loss for 5-way 5-shot car classification increases from 0.7% (original dataset) to 7.0% (20% training data).

**Improvements generalize to other meta-learners** We combine SSL with other meta-learners and find the combination to be effective. In particular, we use MAML [16] and a standard feature extractor trained with cross-entropy loss (softmax) as in [10]. Table 2 compares meta-learners based on a ResNet-18 network trained with and without *jigsaw puzzle loss*. We observe that the average 5-way 5-shot accuracies across five fine-grained datasets for softmax, MAML, and ProtoNet improve from 85.5%, 82.6%, and 88.5% to 86.6%, 83.8%, and 90.4% respectively when combined with the jigsaw puzzle task. Self-supervision improves performance across different meta-learners and different datasets; however, ProtoNet trained with self-supervision is the best model across all datasets.

**Self-supervision alone is not enough** SSL alone significantly lags behind supervised learning in our experiments. For example, a ResNet-18 trained with SSL alone achieve 32.9% (w/ jigsaw) and 33.7% (w/ rotation) 5-way 5-shot accuracy averaged across five fine-grained datasets. While this is better than a random initialization (29.5%), it is dramatically worse than one trained with a simple cross-entropy loss (85.5%) on the labels (details in Table 4 in Appendix A.1). Surprisingly, we also found that initialization with SSL followed by meta-learning did
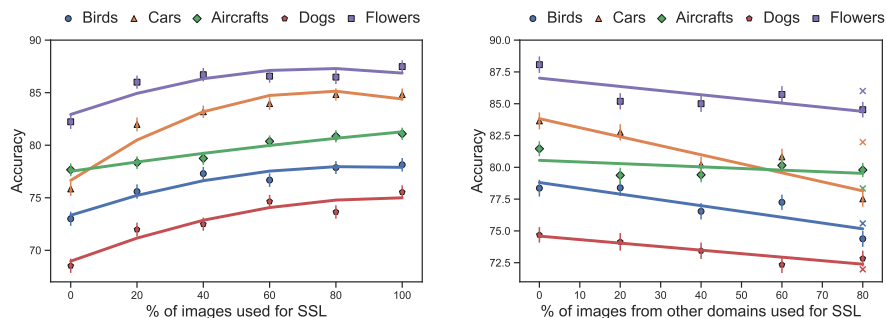
| Model | mage Size | Backbone | SSL | Accuracy (%) |
|---|---|---|---|---|
| MAML [16] | | Conv4-64 | - | 63.1 |
| ProtoNet [55] | | Conv4-64 | - | 68.2 |
| RelationNet [57] | | Conv4-64 | - | 65.3 |
| LwoF [19] | 84×84 | Conv4-64 | - | 72.8 |
| PFA [48]* | | WRN-28-10 | - | 73.7 |
| TADAM [44] | | ResNet-12 | - | 76.7 |
| LEO [53]* | | WRN-28-10 | - | 77.6 |
| MetaOptNet-SVM [35]† | | ResNet-12 | - | 78.6 |
| Chen *et al.* [10] | 84×84 | Conv4-64 | - | 64.2 |
| (ProtoNet) | 224×224 | ResNet-18 | - | 73.7 |
| | | Conv4-64 | - | 70.0 |
| | | | Rotation | 71.7 |
| Gidaris *et al.* [18] | 84×84 | Conv4-512 | - | 71.6 |
| (ProtoNet) | | | Rotation | 74.0 |
| | | WRN-28-10 | - | 68.7 |
| | | | Rotation | 72.1 |
| | | | - | 75.2 |
| **Ours** | 224×224 | ResNet-18 | Rotation | 76.0 |
| (ProtoNet) | | | Jigsaw | 76.2 |
| | | | Rot.+Jig. | 76.6 |

Table 3: **Comparison with prior works on *mini*-ImageNet.** 5-shot 5-way classification accuracies on 600 test episodes are reported. The implementation details including image size, backbone model, and training are different in each paper. *validation classes are used for training. †dropblock [17], label smoothing, and weight decay are used.

*not* yield improvements over meta-learning starting from random initialization, supporting the view that SSL acts as a feature regularizer.

**Few-shot learning as an evaluation for self-supervised tasks** The few-shot classification task provides a way of evaluating the effectiveness of self-supervised tasks. For example, on 5-way 5-shot aircrafts classification, training with only jigsaw and rotation task gives 38.8% and 29.5% respectively, suggesting that rotation is not an effective self-supervised task for airplanes. We speculate that it might be because the task is too easy as airplanes are usually horizontal.

**Comparison with prior works** Our results also echo those of [18] who find that the rotation task improves on *mini*- and *tiered*-ImageNet. In addition we show the improvement still holds when using deeper networks, higher resolution images, and in fine-grained domains. We provide a comparison with other few-shot learning methods in Table 3.

(a) Effect of number of images on SSL.      (b) Effect of domain shift on SSL.

Fig. 4: **Effect of size and domain of SSL on 5-way 5-shot classification accuracy.** **(a)** More unlabeled data from the same domain for SSL improves the performance of the meta-learner. **(b)** Replacing a fraction (x-axis) of the images with those from other domains makes SSL less effective.

### 4.2    Analyzing the Effect of Domain Shift for Self-supervision

Scaling SSL to massive unlabeled datasets that are readily available for some domains is a promising avenue for improvement. *However, do more unlabeled data always help for a task in hand?* This question hasn't been sufficiently addressed in the literature as most prior works study the effectiveness of SSL on a curated set of images, such as ImageNet, and their transferability to a handful of tasks. We conduct a series of experiments to characterize the effect of size and distribution $\mathcal{D}_{ss}$ of images used for SSL in the context of few-shot learning on domain $\mathcal{D}_s$.

First, we investigate if SSL on unlabeled data from the same domain improves the meta-learner. We use 20% of the images in the base categories for meta-learning identical to the setting in Fig. 3. The labels of the remaining 80% data are withheld and only the images are used for SSL. We systematically vary the number of images used by SSL from 20% to 100%. The results are presented in Fig. 4a. The accuracy improves with the size of the unlabeled set with diminishing returns. Note that 0% corresponds to no SSL and 20% corresponds to using only the labeled images for SSL ($\mathcal{D}_s = \mathcal{D}_{ss}$).

Fig. 4b shows an experiment where a fraction of the unlabeled images are replaced with images from other four datasets. For example, 20% along the x-axis for birds indicate that 20% of the images in the base set are replaced by images drawn uniformly at random from other datasets. Since the numbers of images used for SSL is identical, the x-axis from left to right represents increasing amounts of domain shifts between $\mathcal{D}_s$ and $\mathcal{D}_{ss}$. We observe that the effectiveness of SSL decreases as the fraction of out-of-domain images increases. Importantly, training with SSL on the available 20% within domain images (shown as crosses) is often (on 3 out of 5 datasets) better than increasing the set of images by five times to include out of domain images.
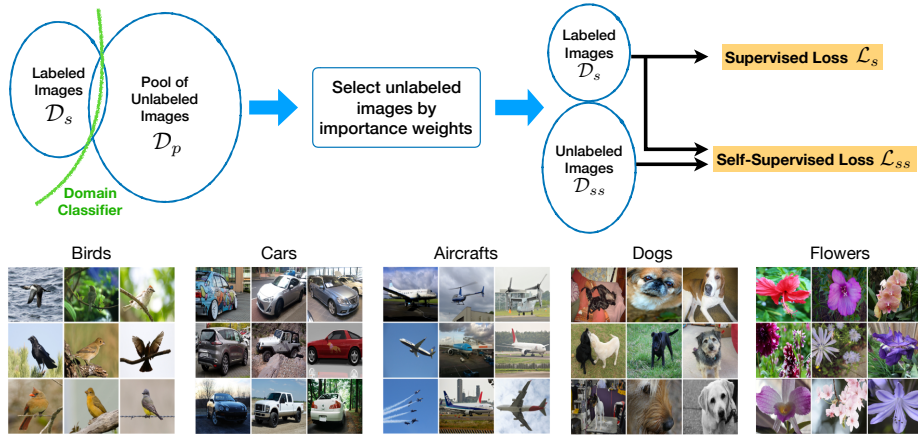
Fig. 5: **Overview of domain selection for self-supervision. Top:** We first train a domain classifier using $\mathcal{D}_s$ and (a subset of) $\mathcal{D}_p$, then select images using the predictions from the domain classifier for self-supervision. **Bottom:** Selected images of each dataset using importance weights.
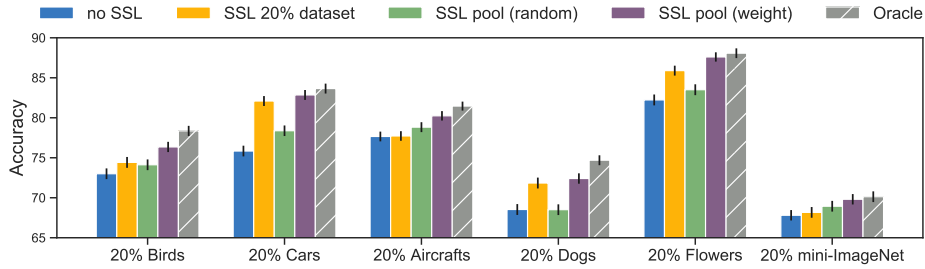


Fig. 6: **Effectiveness of selected images for SSL.** With random selection, the extra unlabeled data often hurts the performance, while those sampled using the *importance weights* improve performance on all five datasets. A tabular version is shown in Appendix A.2.

### 4.3 Selecting Images for Self-supervision

Based on the above analysis we propose a simple method to select images for SSL from a large, generic pool of unlabeled images in a dataset dependent manner. We use a "domain weighted" model to select the top images based on a domain classifier, in our case a binary logistic regression model trained with images from the source domain $\mathcal{D}_s$ as the positive class and images from the pool $\mathcal{D}_p$ as the negative class based on ResNet-101 image features. The top images are selected according to the ratio $p(x \in \mathcal{D}_s)/p(x \in \mathcal{D}_p)$. Note that these *importance weights* account for the domain shift. Fig. 5 shows an overview of the selection process.

We evaluate this approach using a pool of images $\mathcal{D}_p$ consisting of (1) the training images of the "bounding box" subset of Open Images V5 [33] which has 1,743,042 images from 600 classes, and (2) iNaturalist 2018 dataset [60] which has 461,939 images from 8162 species. For each dataset, we use 20% of the labeled images as $\mathcal{D}_s$. The rest 80% of the data are only used as the "oracle" where the unlabeled data are drawn from the exact same distribution as $\mathcal{D}_s$. We show some of the selected images for self-supervision $\mathcal{D}_{ss}$ in Fig. 5.

Fig. 6 shows the results of ProtoNet trained on 20% labeled examples with jigsaw puzzle as self-supervision. To have a fair comparison, for methods of selecting images from the pool, we select the same number (80% of the original labeled dataset size) of images as $\mathcal{D}_{ss}$. We report the mean accuracy of five runs. "SSL with 20% dataset" denotes a baseline of only using $\mathcal{D}_s$ for self-supervision ($\mathcal{D}_s = \mathcal{D}_{ss}$), which is our reference "lower bound". SSL pool "(random)" and "(weight)" denote two approaches of selecting images for self-supervision. The former selects images uniformly at random, which is detrimental for cars, dogs, and flowers. The pool selected according to the *importance weights* provides significant improvements over "no SSL", "SSL with 20% dataset", and "random selection" baselines on all five datasets. The oracle is trained with the remaining 80% of the original dataset as $\mathcal{D}_{ss}$, which is a reference "upper bound".

## 5    Conclusion

Self-supervision improves the performance on few-shot learning tasks across a range of different domains. Surprisingly, we found that self-supervision is more beneficial for more challenging problems, especially when the number of images used for self-supervision is small, orders of magnitude smaller than previously reported results. This has a practical benefit that the images within small datasets can be used for self-supervision without relying on a large-scale external dataset. We have also shown that additional unlabeled images can improve performance only if they are from the *same or similar* domains. Finally, for domains where unlabeled data is limited, we present a novel, simple approach to automatically identify such similar-domain images from a larger pool.

Future work could investigate if using other self-supervised tasks can also improve few-shot learning, in particular constrastive learning approaches [3,22,24, 38,58]. Future work could also investigate how and when self-supervision improves generalization across self-supervised and supervised tasks empirically [1,66].

### Acknowledgement

# References

1. Achille, A., Lam, M., Tewari, R., Ravichandran, A., Maji, S., Fowlkes, C., Soatto, S., Perona, P.: Task2Vec: Task embedding for meta-learning. In: ICCV (2019) 4, 14
2. Asano, Y.M., Rupprecht, C., Vedaldi, A.: A critical analysis of self-supervision, or what we can learn from a single image. In: ICLR (2020) 3
3. Bachman, P., Hjelm, R.D., Buchwalter, W.: Learning representations by maximizing mutual information across views. arXiv preprint arXiv:1906.00910 (2019) 3, 14
4. Bertinetto, L., Henriques, J.F., Torr, P.H., Vedaldi, A.: Meta-learning with differentiable closed-form solvers. In: ICLR (2019) 3
5. Bojanowski, P., Joulin, A.: Unsupervised learning by predicting noise. In: ICML (2017) 3
6. Carlucci, F.M., D'Innocente, A., Bucci, S., Caputo, B., Tommasi, T.: Domain generalization by solving jigsaw puzzles. In: CVPR (2019) 3
7. Caron, M., Bojanowski, P., Joulin, A., Douze, M.: Deep clustering for unsupervised learning of visual features. In: ECCV (2018) 3
8. Caron, M., Bojanowski, P., Mairal, J., Joulin, A.: Unsupervised pre-training of image features on non-curated data. In: ICCV (2019) 3
9. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: ICML (2020) 3
10. Chen, W.Y., Liu, Y.C., Kira, Z., Wang, Y.C., Huang, J.B.: A closer look at few-shot classification. In: ICLR (2019) 3, 7, 8, 10, 11, 21
11. Chen, Z., Badrinarayanan, V., Lee, C.Y., Rabinovich, A.: Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In: ICML (2018) 6
12. Cui, Y., Song, Y., Sun, C., Howard, A., Belongie, S.: Large scale fine-grained categorization and domain-specific transfer learning. In: CVPR (2018) 4
13. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: ICCV (2015) 3
14. Doersch, C., Zisserman, A.: Multi-task self-supervised visual learning. In: ICCV (2017) 3
15. Dosovitskiy, A., Springenberg, J.T., Riedmiller, M., Brox, T.: Discriminative unsupervised feature learning with convolutional neural networks. In: NeurIPS (2014) 3
16. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML (2017) 3, 6, 8, 10, 11
17. Ghiasi, G., Lin, T.Y., Le, Q.V.: Dropblock: A regularization method for convolutional networks. In: NeurIPS (2018) 11
18. Gidaris, S., Bursuc, A., Komodakis, N., Pérez, P., Cord, M.: Boosting few-shot visual learning with self-supervision. In: ICCV (2019) 2, 4, 11
19. Gidaris, S., Komodakis, N.: Dynamic few-shot visual learning without forgetting. In: CVPR (2018) 3, 11
20. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. In: ICLR (2018) 3, 6
21. Goyal, P., Mahajan, D., Gupta, A., Misra, I.: Scaling and benchmarking self-supervised visual representation learning. In: ICCV (2019) 3, 4, 6
22. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: CVPR (2020) 3, 14
23. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) 8

24. Hénaff, O.J., Razavi, A., Doersch, C., Eslami, S., Oord, A.v.d.: Data-efficient image recognition with contrastive predictive coding. arXiv preprint arXiv:1905.09272 (2019) 3, 14

25. Hjelm, R.D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., Bengio, Y.: Learning deep representations by mutual information estimation and maximization. In: ICLR (2019) 3

26. Kendall, A., Gal, Y., Cipolla, R.: Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In: CVPR (2018) 6

27. Khosla, A., Jayadevaprakash, N., Yao, B., Fei-Fei, L.: Novel dataset for fine-grained image categorization. In: First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2011) 7

28. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015) 8

29. Koch, G., Zemel, R., Salakhutdinov, R.: Siamese neural networks for one-shot image recognition. In: ICML deep learning workshop. vol. 2 (2015) 6

30. Kokkinos, I.: Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In: CVPR (2017) 4

31. Kolesnikov, A., Zhai, X., Beyer, L.: Revisiting self-supervised visual representation learning. In: CVPR (2019) 3, 4

32. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3D object representations for fine-grained categorization. In: 4th International IEEE Workshop on 3D Representation and Recognition (3DRR). Sydney, Australia (2013) 7

33. Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J., Krasin, I., Pont-Tuset, J., Kamali, S., Popov, S., Malloci, M., Duerig, T., Ferrari, V.: The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. arXiv:1811.00982 (2018) 14

34. Larsson, G., Maire, M., Shakhnarovich, G.: Learning representations for automatic colorization. In: ECCV (2016) 3

35. Lee, K., Maji, S., Ravichandran, A., Soatto, S.: Meta-learning with differentiable convex optimization. In: CVPR (2019) 3, 11

36. Maji, S., Rahtu, E., Kannala, J., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151 (2013) 7

37. Maninis, K.K., Radosavovic, I., Kokkinos, I.: Attentive single-tasking of multiple tasks. In: CVPR (2019) 4

38. Misra, I., van der Maaten, L.: Self-supervised learning of pretext-invariant representations. In: CVPR (2020) 3, 14

39. Ngiam, J., Peng, D., Vasudevan, V., Kornblith, S., Le, Q.V., Pang, R.: Domain adaptive transfer learning with specialist models. arXiv preprint arXiv:1811.07056 (2018) 4

40. Nilsback, M.E., Zisserman, A.: A visual vocabulary for flower classification. In: CVPR (2006) 7

41. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: ECCV (2016) 6, 23

42. Noroozi, M., Pirsiavash, H., Favaro, P.: Representation learning by learning to count. In: ICCV (2017) 3

43. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018) 3

44. Oreshkin, B., López, P.R., Lacoste, A.: Tadam: Task dependent adaptive metric for improved few-shot learning. In: NeurIPS (2018) 11

45. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: PyTorch: An imperative style, high-performance deep learning library. In: NeurIPS (2019) 21
46. Pathak, D., Krähenbühl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: CVPR (2016) 3
47. Qi, H., Brown, M., Lowe, D.G.: Low-shot learning with imprinted weights. In: CVPR (2018) 3
48. Qiao, S., Liu, C., Shen, W., Yuille, A.L.: Few-shot image recognition by predicting parameters from activations. In: CVPR (2018) 3, 11
49. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2017) 3
50. Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle, H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. In: ICLR (2018) 7
51. Ren, Z., Lee, Y.J.: Cross-domain self-supervised multi-task feature learning using synthetic imagery. In: CVPR (2018) 4
52. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet large scale visual recognition challenge. International Journal of Computer Vision (IJCV) (2015) 22
53. Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. arXiv preprint arXiv:1807.05960 (2018) 11
54. Sener, O., Koltun, V.: Multi-task learning as multi-objective optimization. In: NeurIPS (2018) 6
55. Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: NeurIPS (2017) 3, 5, 8, 11
56. Su, J.C., Maji, S.: Adapting models to signal degradation using distillation. In: BMVC (2017) 10
57. Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H., Hospedales, T.M.: Learning to compare: Relation network for few-shot learning. In: CVPR (2018) 11
58. Tian, Y., Krishnan, D., Isola, P.: Contrastive multiview coding. In: ECCV (2020) 14
59. Trinh, T.H., Luong, M.T., Le, Q.V.: Selfie: Self-supervised pretraining for image embedding. arXiv preprint arXiv:1906.02940 (2019) 3
60. Van Horn, G., Mac Aodha, O., Song, Y., Cui, Y., Sun, C., Shepard, A., Adam, H., Perona, P., Belongie, S.: The iNaturalist species classification and detection dataset. In: CVPR (2018) 14
61. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: NeurIPS (2016) 3, 6, 7
62. Wallace, B., Hariharan, B.: Extending and analyzing self-supervised learning across domains. In: ECCV (2020) 4
63. Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., Perona, P.: Caltech-UCSD Birds 200. Tech. Rep. CNS-TR-2010-001, California Institute of Technology (2010) 5, 7
64. Wertheimer, D., Hariharan, B.: Few-shot learning with localization in realistic settings. In: CVPR (2019) 21
65. Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. In: CVPR (2018) 3

66. Zamir, A.R., Sax, A., Shen, W., Guibas, L.J., Malik, J., Savarese, S.: Taskonomy: Disentangling task transfer learning. In: CVPR. pp. 3712–3722 (2018) 14
67. Zhai, X., Oliver, A., Kolesnikov, A., Beyer, L.: S4L: Self-supervised semi-supervised learning. In: ICCV (2019) 3
68. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: ECCV (2016) 3
69. Zhang, R., Isola, P., Efros, A.A.: Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In: CVPR (2017) 3

# A   Appendix

In Appendix A.1 and Appendix A.2, we provide all the numbers of the figures in Section 4.1 and Section 4.3 separately. We show that SSL can also improve traditional fine-grained classification in Appendix A.3 and its model visualization in Appendix A.4. Last, we describe the implementation details in Appendix A.5.

## A.1   Results on Few-shot Learning

Table 4 shows the performance of ProtoNet with different self-supervision on seven datasets. We also test the accuracy of the model on novel classes when trained *only* with self-supervision on the base set of images. Compared to the randomly initialized model ("None" rows), training the network to predict rotations gives around 2% to 21% improvements on all datasets, while solving jigsaw puzzles only improves on aircrafts and flowers. However, these numbers are significantly worse than learning with supervised labels on the base set, in line with the current literature.

Table 5 shows the performance of ProtoNet with *jigsaw puzzle loss* on harder benchmarks. The results on the degraded version of the datasets are shown in the top part, and the bottom part shows the results of using only 20% of the images in the base categories. The gains using SSL are higher in this setting.

## A.2   Results on Selecting Images for SSL

Table 6 shows the performance of selecting images for self-supervision, a tabular version of Figure 5 in Section 4.3. "Pool (random)" uniformly samples images proportional to the size of each dataset, while the "pool (weight)" one tends to pick more images from related domains.

## A.3   Results on Standard Fine-grained Classification

Here we present results on standard fine-grained classification tasks. Different from few-shot transfer learning, all the classes are seen in the training set and the test set contains novel images from the same classes. We use the standard training and test splits provided in the datasets. We investigate if SSL can improve the training of deep networks (*e.g.* ResNet-18 network) when *trained from scratch* (*i.e.* with random initialization) using images and labels in the training set only.

| Loss | mini-ImageNet | tiered-ImageNet | Birds | Cars | Aircrafts | Dogs | Flowers |
|---|---|---|---|---|---|---|---|
| | | | 5-way 5-shot | | | | |
| ProtoNet(PN) | 75.2±0.6 | 75.9±0.7 | 87.3±0.5 | 91.7±0.4 | 91.4±0.4 | 83.0±0.6 | 89.2±0.6 |
| PN+Jigsaw | 76.2±0.6 | 78.0±0.7 | 89.8±0.4 | 92.4±0.4 | 91.8±0.4 | 85.7±0.5 | **92.2±0.4** |
| *Rel. err. red.* | *4.0%* | *8.7%* | *19.7%* | *8.4%* | *4.7%* | *15.9%* | *27.8%* |
| Jigsaw | 25.6±0.5 | 24.9±0.4 | 25.7±0.5 | 25.3±0.5 | 38.8±0.6 | 24.3±0.5 | 50.5±0.7 |
| PN+Rotation | 76.0±0.6 | **78.9±0.7** | 89.4±0.4 | 92.3±0.4 | 91.4±0.4 | 84.3±0.5 | 89.0±0.5 |
| Rotation | 51.4±0.7 | 50.7±0.8 | 33.1±0.6 | 29.4±0.5 | 29.5±0.5 | 27.3±0.5 | 49.4±0.7 |
| PN+Jig.+Rot. | **76.6±0.7** | 77.2±0.7 | **90.2±0.4** | **92.7±0.4** | **91.9±0.4** | **85.9±0.5** | 91.4±0.5 |
| None | 31.0±0.5 | 28.9±0.5 | 26.7±0.5 | 25.2±0.5 | 28.1±0.5 | 25.3±0.5 | 42.3±0.8 |
| | | | 20-way 5-shot | | | | |
| ProtoNet(PN) | 46.6±0.3 | 49.7±0.4 | 69.3±0.3 | 78.7±0.3 | 78.6±0.3 | 61.6±0.3 | 75.4±0.3 |
| PN+Jigsaw | 47.8±0.3 | **52.4±0.4** | 73.7±0.3 | 79.1±0.3 | **79.1±0.2** | 65.4±0.3 | **79.2±0.3** |
| Jigsaw | 9.2±0.2 | 7.5±0.1 | 8.1±0.1 | 7.1±0.1 | 15.4±0.2 | 7.1±0.1 | 25.7±0.2 |
| PN+Rotation | 48.2±0.3 | **52.4±0.4** | 72.9±0.3 | **80.0±0.3** | 78.4±0.2 | 63.4±0.3 | 73.9±0.3 |
| Rotation | 27.4±0.2 | 25.7±0.3 | 12.9±0.2 | 9.3±0.2 | 9.8±0.2 | 8.8±0.1 | 26.3±0.2 |
| PN+Jig.+Rot. | **49.0±0.3** | 51.2±0.4 | **75.0±0.3** | 79.8±0.3 | 79.0±0.2 | **66.2±0.3** | 78.6±0.3 |
| None | 10.8±0.1 | 11.0±0.2 | 9.3±0.2 | 7.5±0.1 | 8.9±0.1 | 7.8±0.1 | 22.6±0.2 |

Table 4: **Performance on few-shot learning tasks.** The mean accuracy (%) and the 95% confidence interval of 600 randomly chosen test experiments are reported for various combinations of loss functions. The top part shows the accuracy on 5-way 5-shot classification tasks, while the bottom part shows the same on 20-way 5-shot. Adding self-supervised losses to the ProtoNet loss improves the performance on all seven datasets on 5-way classification results. On 20-way classification, the improvements are even larger. The last row indicates results with a randomly initialized network. The top part of this table corresponds to Figure 2 in Section 4.1.

The accuracy of using various loss functions are shown in Table 7. Training with self-supervision improves performance across datasets. On birds, cars, and dogs, predicting rotation gives 4.1%, 3.1%, and 3.0% improvements, while on aircrafts and flowers, the *jigsaw puzzle loss* yields 0.9% and 3.6% improvements.

## A.4   Visualization of Learned Models

To understand why the representation generalizes, we visualize what pixels contribute the most to the correct classification for various models. In particular, for each image and model, we compute the gradient of the logits (predictions before softmax) for the correct class with respect to the input image. The magnitude of the gradient at each pixel is a proxy for its importance and is

| Loss | Birds Greyscale | Cars Low-resolution | Aircrafts Low-resolution | Dogs Greyscale | Flowers Greyscale |
|---|---|---|---|---|---|
| | 5-way 5-shot | | | | |
| ProtoNet | 82.2±0.6 | 84.8±0.5 | 85.0±0.5 | 80.7±0.6 | 86.1±0.6 |
| ProtoNet + Jigsaw | 85.4±0.6 | 87.0±0.5 | 87.1±0.5 | 83.6±0.5 | 87.6±0.5 |
| | 20-way 5-shot | | | | |
| ProtoNet | 60.8±0.4 | 64.7±0.3 | 64.1±0.3 | 57.4±0.3 | 69.7±0.3 |
| ProtoNet + Jigsaw | 65.7±0.3 | 68.6±0.3 | 68.3±0.3 | 61.2±0.3 | 71.6±0.3 |
| Loss | 20% Birds | 20% Cars | 20% Aircrafts | 20% Dogs | 20% Flowers |
| | 5-way 5-shot | | | | |
| ProtoNet | 73.0±0.7 | 75.8±0.7 | 77.7±0.6 | 68.5±0.7 | 82.2±0.7 |
| ProtoNet + Jigsaw | 75.4±0.7 | 82.8±0.6 | 78.4±0.6 | 69.1±0.7 | 86.0±0.6 |
| | 20-way 5-shot | | | | |
| ProtoNet | 46.4±0.3 | 51.8±0.4 | 52.3±0.3 | 40.8±0.3 | 62.8±0.3 |
| ProtoNet + Jigsaw | 49.8±0.3 | 61.5±0.4 | 53.6±0.3 | 42.2±0.3 | 68.5±0.3 |

Table 5: **Performance on *harder* few-shot learning tasks.** Accuracies are reported on novel set for 5-way 5-shot and 20-way 5-shot classification with degraded inputs, and with a subset (20%) of the images in the base set. The loss of color or resolution, and the smaller training set size make the tasks more challenging as seen by the drop in the performance of the ProtoNet baseline. However the improvements of using the *jigsaw puzzle loss* are higher in comparison to the results presented in Table 4.

| Method | 20% Birds | 20% Cars | 20% Aircrafts | 20% Dogs | 20% Flowers | 20% *mini-ImageNet* |
|---|---|---|---|---|---|---|
| No SSL | 73.0±0.7 | 75.8±0.7 | 77.7±0.6 | 68.5±0.7 | 82.2±0.7 | 67.81±0.65 |
| SSL 20% dataset | 74.4±0.7 | 82.1±0.6 | 77.7±0.6 | 71.8±0.7 | 85.9±0.6 | 68.47±0.66 |
| SSL Pool (random) | 74.1±0.7 | 78.4±0.7 | 78.8±0.6 | 68.5±0.7 | 83.5±0.7 | 68.94±0.68 |
| SSL Pool (weight) | **76.4±0.6** | **82.9±0.6** | **80.2±0.6** | **72.4±0.7** | **87.6±0.6** | **69.81±0.65** |
| *SSL 100% (oracle)* | *78.4±0.6* | *83.7±0.6* | *81.5±0.6* | *74.7±0.6* | *88.1±0.6* | *70.13±0.67* |

Table 6: **Performance on selecting images for self-supervision.** Adding more unlabeled images selected randomly from a pool often hurts the performance. Selecting similar images by importance weights improves on all five datasets.

visualized as "saliency maps". Figure 7 shows these maps for various images and models trained with and without self-supervision on the standard classification task. It appears that the self-supervised models tend to focus more on the

| Loss | Birds | Cars | Aircrafts | Dogs | Flowers |
|---|---|---|---|---|---|
| Softmax | 47.0 | 72.6 | 69.9 | 51.4 | 72.8 |
| Softmax + Jigsaw | 49.2 | 73.2 | **70.8** | 53.5 | **76.4** |
| Softmax + Rotation | **51.1** | **75.7** | 70.0 | **54.4** | 73.5 |

Table 7: **Performance on standard fine-grained classification tasks.** Per-image accuracy (%) on the test set are reported. Using self-supervision improves the accuracy of a ResNet-18 network trained *from scratch* over the baseline of supervised training with cross-entropy (softmax) loss on all five datasets.

foreground regions, as seen by the amount of bright pixels within the bounding box. One hypothesis is that self-supervised tasks force the model to rely less on background features, which might be accidentally correlated to the class labels. For fine-grained recognition, localization indeed improves performance when training from few examples (see [64] for a contemporary evaluation of the role of localization for few-shot learning).

### A.5   Experimental Details

**Optimization details on few-shot learning**  During training, especially for the jigsaw puzzle task, we found it to be beneficial to *not* track the running mean and variance for the batch normalization layer, and instead estimate them for each batch independently. We hypothesize that this is because the inputs contain both full-sized images and small patches, which might have different statistics. At test time we do the same. We found the accuracy goes up as the batch size increases but saturates at a size of 64.

When training with supervised and self-supervised loss, a trade-off term $\lambda$ between the losses can be used, thus the total loss is $\mathcal{L} = (1 - \lambda)\mathcal{L}_s + \lambda\mathcal{L}_{ss}$. We find that simply use $\lambda = 0.5$ works the best, except for training on *mini*- and *tiered*-ImageNet with jigsaw loss, where we set $\lambda = 0.3$. We suspect that this is because the variation of the image size and the categories are higher, making the self-supervision harder to train with limited data. When both jigsaw and rotation losses are used, we set $\lambda = 0.5$ and the two self-supervised losses are averaged for $\mathcal{L}_{ss}$.

For training meta-learners, we use 16 query images per class for each training episode. When only 20% of labeled data are used, 5 query images per class are used. For MAML, we use 10 query images and the approximation method for backpropagation as proposed in [10] to reduce the GPU memory usage. When training with self-supervised loss, it is added when computing the loss in the outer loop. We use PyTorch [45] for our experiments.

**Optimization details on domain classifier**  For the domain classifier, we first obtain features from the penultimate-layer (2048 dimensional) from a ResNet-101

Fig. 7: **Saliency maps for various images and models.** For each image we visualize the magnitude of the gradient with respect to the correct class for models trained with various loss functions. The magnitudes are scaled to the same range for easier visualization. The models trained with self-supervision often have lower energy on the background regions when there is clutter. We highlight a few examples with blue borders and the bounding-box of the object for each image is shown in red.

model pre-trained on ImageNet [52]. We then train a binary logistic regression model with weight decay using LBFGS for 1000 iterations. The images from the labeled dataset are the positive class and from the pool of unlabeled data are the negative class. A subset of negative images are selected uniformly at random with 10 times the size of positive images. A loss for the positive class is scaled by the inverse of its frequency to account for the significantly larger number of negative examples.

**Optimization details on standard classification** For standard classification (Appendix A.3) we train a ResNet-18 network *from scratch*. All the models are trained with ADAM optimizer with a learning rate of 0.001 for 600 epochs with a batch size of 16. We track the running statistics for the batch normalization layer for the softmax baselines following the conventional setting, *i.e.* w/o self-supervised loss, but do not track these statistics when training with self-supervision.

**Architectures for self-supervised tasks** For jigsaw puzzle task, we follow the architecture of [41] where it was first proposed. The ResNet18 results in a 512-dimensional feature for each input, and we add a fully-connected (`fc`) layer with 512-units on top. The nine patches give nine 512-dimensional feature vectors, which are concatenated. This is followed by a `fc` layer, projecting the feature vector from 4608 to 4096 dimensions, and a `fc` layer with 35-dimensional outputs corresponding to the 35 permutations for the jigsaw task.

For rotation prediction task, the 512-dimensional output of ResNet-18 is passed through three `fc` layers with {512, 128, 128, 4} units. The predictions correspond to the four rotation angles. Between each `fc` layer, a `ReLU` activation and a dropout layer with a dropout probability of 0.5 are added.