

MBZUAI

Digital.Commons@MBZUAI

Computer Vision Faculty Publications

Scholarly Works

7-16-2021

Unsupervised anomaly instance segmentation for baggage threat recognition

Taimur Hassan

Khalifa University of Science and Technology

Samet Akçay

Durham University

Mohammed Bennamoun

The University of Western Australia

Salman Khan

Mohamed Bin Zayed University of Artificial Intelligence

Naoufel Werghi

Khalifa University of Science and Technology

Follow this and additional works at: <https://dclibrary.mbzuai.ac.ae/cvfp>



Part of the [Artificial Intelligence and Robotics Commons](#)

Preprint: arXiv

- Archived with thanks to arXiv
- Preprint License: CC BY 4.0
- Uploaded 29 March 2022

Recommended Citation

T. Hassan, S. Akçay, M. Bennamoun, S. Khan, and N. Werghi, "Unsupervised anomaly instance segmentation for baggage threat recognition," 2021, arXiv:2107.07333v2

This Article is brought to you for free and open access by the Scholarly Works at Digital.Commons@MBZUAI. It has been accepted for inclusion in Computer Vision Faculty Publications by an authorized administrator of Digital.Commons@MBZUAI. For more information, please contact libraryservices@mbzuai.ac.ae.

Exploring Complementary Strengths of Invariant and Equivariant Representations for Few-Shot Learning

Mamshad Nayeem Rizve[†]

Salman Khan[‡]

Fahad Shahbaz Khan[‡]

Mubarak Shah[†]

[†]Center for Research in Computer Vision, UCF, USA

[‡]Mohamed bin Zayed University of AI, UAE

nayeemrizve@knights.ucf.edu, {salman.khan, fahad.khan}@mbzuai.ac.ae, shah@crcv.ucf.edu

Abstract

In many real-world problems, collecting a large number of labeled samples is infeasible. Few-shot learning (FSL) is the dominant approach to address this issue, where the objective is to quickly adapt to novel categories in presence of a limited number of samples. FSL tasks have been predominantly solved by leveraging the ideas from gradient-based meta-learning and metric learning approaches. However, recent works have demonstrated the significance of powerful feature representations with a simple embedding network that can outperform existing sophisticated FSL algorithms. In this work, we build on this insight and propose a novel training mechanism that simultaneously enforces equivariance and invariance to a general set of geometric transformations. Equivariance or invariance has been employed standalone in the previous works; however, to the best of our knowledge, they have not been used jointly. Simultaneous optimization for both of these contrasting objectives allows the model to jointly learn features that are not only independent of the input transformation but also the features that encode the structure of geometric transformations. These complementary sets of features help generalize well to novel classes with only a few data samples. We achieve additional improvements by incorporating a novel self-supervised distillation objective. Our extensive experimentation shows that even without knowledge distillation our proposed method can outperform current state-of-the-art FSL methods on five popular benchmark datasets. Our codes are available at: <https://github.com/nayeemrizve/invariance-equivariance>.

1. Introduction

In recent years, deep learning methods have made great strides on several challenging problems [29, 72, 28, 6, 7]. This success can be partially attributed to the availability of large-scale labeled datasets [14, 6, 83, 44]. However, acquiring large amounts of labeled data is infeasible in several real-world problems due to practical constraints such as the rarity of an event or the high cost of manual anno-

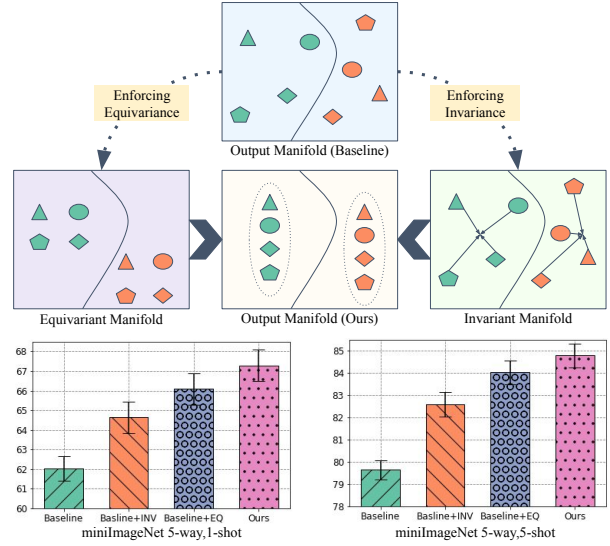


Figure 1. *Approach Overview*: Shapes represent different transformations and colors represent different classes. While the invariant features provide better discrimination, the equivariant features help us learn the internal structure of the data manifold. These complementary representations help us generalize better to new tasks with only a few training samples. By jointly leveraging the strengths of equivariant and invariant features, our method achieves significant improvement over baseline (bottom row).

tation. Few-shot learning (FSL) targets this problem by learning a model on a set of base classes and studies its adaptability to novel classes with only a few samples (typically 1-5) [19, 77, 66, 71]. Remarkably, this setting is different from transfer and self/semi-supervised learning that assumes the availability of pretrained models [64, 81, 36] or large-amounts of unlabeled data [17, 9, 3].

FSL has been predominantly solved using ideas from meta-learning. The two most dominant approaches are optimization-based meta-learning [19, 32, 62] and metric-learning based methods [66, 71, 1]. Both sets of approaches attempt to train a base learner which can be quickly adapted in the presence of a few novel class examples. However, recently it has been shown in [56] that the quick adaptation of the base learner crucially depends on *feature reuse*. Other

recent works [73, 15, 10] have also shown that a baseline feature extractor trained on all the meta-train set can achieve comparable performance to the state-of-the-art meta learning based methods. This brings in an interesting question: *How much further can FSL performance be pushed by simply improving the base feature extractor?*

To answer this question, first, we take a look at the inductive biases in machine learning (ML) algorithms. The optimization of all ML algorithms takes advantage of different inductive biases for hypothesis selection; as the solutions are never unique. The generalization of these algorithms often relies on the effective design of inductive biases, since they encode our priori preference for a particular set of solutions. For instance, regularization methods like ℓ_1/ℓ_2 -penalties [74], dropout [67], or early stopping [53] implicitly impose Occam’s razor in the optimization process by selecting simpler solutions. Likewise, convolutional neural networks (CNN) by design impose translation invariance [2] which makes the internal embeddings translation equivariant. Inspired by this, several methods [12, 20, 16] have attempted to generalize CNNs by imposing *equivariance* to different geometric transformations so that the internal structure of data can be modeled more efficiently. On the other hand, methods like [38] try to be robust against nuisance variations by learning transformation *invariant* features. However, such inductive biases do not provide optimal generalization on FSL tasks and the design of efficient inductive designs for FSL is relatively unexplored.

In this paper, we propose a novel feature learning approach by designing an effective set of inductive biases. We observe that the features required to achieve invariance against input transformations can provide better discrimination, but do not ensure optimal generalization. Similarly, features that focus on transformation discrimination are not optimal for class discrimination but learn equivariant properties that help in learning the data structure leading to better transferability. Therefore, we propose to combine the complementary strengths of both feature types through a multi-task objective that simultaneously seeks to retain both invariant and equivariant features. We argue that learning such generic features encourages the base feature extractor to be more general. We validate this claim by performing extensive experimentation on multiple benchmark datasets. We also conduct thorough ablation studies to demonstrate that enforcing both equivariance and invariance outperforms enforcing either of these objectives alone (see Fig. 1).

Our main contributions are:

- We enforce complimentary *equivariance* and *invariance* to a general set of geometric transformations to model the underlying structure of the data, while remaining discriminative, thereby improving generalization for FSL.
- Instead of extensive architectural changes, we propose

a simple alternative by defining self-supervised tasks as auxiliary supervision. For *equivariance*, we introduce a transformation discrimination task, while an instance discrimination task is developed to learn transformation *invariant* features.

- We demonstrate additional gains with cross-task knowledge distillation that retains the variance properties.

2. Related Works

Few-shot Learning: The FSL approaches generally belong to the meta-learning family, which either learn a generalizable metric space [66, 35, 78, 51] or apply gradient-based updates to obtain a good initialization. In the first class of methods, Siamese networks related a pair of images [35], matching networks applied an LSTM based context encoder to match query and support set images [78], and prototypical networks used the distance between the query and the prototype embedding for class assignment [66]. A task-dependent metric scaling approach to improve FSL was introduced in [51]. The second category use gradient-based meta-learning methods that include using a sequence model (e.g., LSTM) to learn generalizable optimization rules [58], Model-agnostic Meta-Learning (MAML) to find a good initialization that can be quickly adapted to new tasks with minimal supervision [19], and Latent Embedding Optimization (LEO) that applied MAML in the low dimensional space from which high-dimensional parameters can be generated. A few recent efforts, e.g., ProtoMAML [76], combined the complementary strengths of metric-learning and gradient-based meta-learning methods.

Inductive Biases in CNNs: Inductive biases reflect our prior knowledge regarding a particular problem. State of the art CNNs are based on such design choices which range from the convolutional operator (e.g., the weight sharing and translational equivariance) [40], pooling operator (e.g., local neighbourhood relevance) [11], regularization mechanisms (e.g., sparsity with ℓ_1 regularizer) [33], and loss functions (e.g., max-margin boundaries) [27]. Similarly, recurrent architectures and attention mechanisms are biased towards preserving contextual information and being invariant to time translation [2]. A number of approaches have been designed to achieve invariance to nuisances such as natural perturbations [30, 75], viewpoint changes [46], and image transformations [13, 5]. On the other hand, equivariant representations have also been investigated to retain knowledge regarding group actions [12, 54, 63, 42], thereby maintaining meaningful structure in the representations. In this work, we advocate that the representations required to simultaneously achieve invariance and equivariance can be useful for generalization to new tasks with limited data.

Self-supervised Learning for FSL: Our self-supervised loss is inspired by the recent progress in self-supervised

learning (SSL), where proxy tasks are defined to learn transferable representations without adding any manual annotations [57]. The pretext tasks include colorization [39, 82], inpainting [52], relative patch location [17, 50], and amount of rotation applied [24]. Recently, the potential of SSL for FSL was explored in [23, 68]. In [23] a parallel branch with the rotation prediction task to help learn generalizable features was added. Su *et al.* [68] also used rotation and permutation of patches as auxiliary tasks and concluded that SSL is more effective in low-shot regimes and under significant domain shifts. A recent approach employed SimCLR [9] style contrastive learning with augmented pairs to learn improved representations in either unsupervised pretraining [45] or episodic training [18] for FSL.

In contrast to the existing SSL approaches for FSL, we propose to jointly optimize for a complimentary pair of pretext tasks that lead to better generalization. Our novel distillation objective acquires knowledge from the classification as well as proxy task heads and demonstrates further performance improvements. We present our approach next.

3. Our Approach

We first describe the problem setting and the baseline training approach and then present our proposed approach.

3.1. Problem Formulation

Few-shot learning (FSL) operates in two phases, first a model on a set of *base* classes is trained and then at inference a new set of *few-shot* classes are received. We define the base training set as $\mathcal{D}_b = \{(\mathbf{x}, \mathbf{y})\}$, where $\mathbf{x} \in \mathcal{I} \subset \mathbb{R}^{h \times w \times 3}$ is an image, and the one-hot encoded label $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^{N_b}$ can belong to a total of N_b base classes. At inference, a data set of few-shot classes $\mathcal{D}_f = \{(\mathbf{x}, \mathbf{y})\}$ is presented for learning such that the label \mathbf{y} belongs to one of the N_f novel classes, each with a total of K examples (K typically ranges between 1-5). The evaluation setting for few-shot classes is denoted as N_f -way, K -shot. Importantly, the N_b base and N_f few-shot classes belong to totally disjoint sets.

For solving the FSL task, most meta-learning methods [19, 66, 77] have leveraged an episodic training scheme. An episode consists of a small train and test set $(\mathcal{D}_{tr}^i, \mathcal{D}_{ts}^i)$. The examples for the train and test set of an episode are sampled from the same distribution i.e. from the same subset of meta-training classes. Meta-learning methods try to optimize the parameters of the base learner by solving a collection of these episodes. The main motivation is that the evaluation conditions should be emulated in the base training stage. However, following recent works [73, 15, 10], we do not use an episodic training scheme which allows us to train a single generalizable model that can be efficiently used for any-way, any-shot setting without retraining. Specifically,

we train our base learner on the whole base training set \mathcal{D}_b in a supervised manner.

Let's assume our base learner for the FSL task is a neural network, f_Θ , parameterized with parameters Θ . The role of this base learner is to extract good feature embeddings that can generalize for novel classes. The base learner f_Θ can project an input image \mathbf{x} into the embedding space $f_\Theta : \mathbf{x} \rightarrow \mathbf{z}$, such that $\mathbf{z} \in \mathbb{R}^d$. Now, to optimize the parameters of the base learner f_Θ we need a classifier to project the extracted embeddings into the label space. To this end, we introduce a classifier function, f_Φ , with parameters Φ that projects the embeddings \mathbf{z} into the label space \mathcal{Y} i.e., $f_\Phi : \mathbf{z} \rightarrow \tilde{\mathbf{y}}$, such that $\tilde{\mathbf{y}} \in \mathcal{Y}$.

We jointly optimize the parameters of both f_Θ and f_Φ by minimizing cross-entropy loss on the whole base-training set \mathcal{D}_b . The classification loss is given by,

$$\mathcal{L}_{ce} = -\log \frac{\exp(\tilde{\mathbf{y}}_{j:\mathbf{y}_j=1})}{\sum_i \exp(\tilde{\mathbf{y}}_i)}, \text{ s.t., } \mathbf{y} \in \{0, 1\}^{N_b}, \tilde{\mathbf{y}} = f_{\Theta, \Phi}(\mathbf{x}).$$

To regularize the parameters of both of the sub-networks, we add a regularization term. Hence, the learning objective for our baseline training algorithm becomes:

$$\mathcal{L}_{baseline} = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_b} [\mathcal{L}_{ce}(f_{\Theta, \Phi}(\mathbf{x}), \mathbf{y})] + \mathcal{R}(\Theta, \Phi). \quad (1)$$

Here, $\mathcal{R}(\Theta, \Phi)$ is an \mathcal{L}_2 regularization term for the parameters Θ and Φ . Next, we present our inductive objectives.

3.2. Injecting Inductive Biases through SSL

We propose to enforce equivariance and invariance to a general set of geometric transformations \mathcal{T} by simply performing self-supervised learning (SSL). Self-supervision is particularly useful for learning general features without accessing semantic labels. For representation learning, self-supervised methods generally aim for either achieving equivariance to some input transformations or learn to discriminate instances by making the representations invariant. To the best of our knowledge, simultaneous equivariance and invariance to a general set of geometric transformations \mathcal{T} has not been explored in the self-supervised literature. We are the first ones to do so.

The transformation set \mathcal{T} can be obtained from a family of geometric transformations, $\mathcal{D}_{\mathcal{T}}$; $\mathcal{T} \sim \mathcal{D}_{\mathcal{T}}$. Here, $\mathcal{D}_{\mathcal{T}}$ can be interpreted as a family of geometric transformations like Euclidean transformation, Similarity transformation, Affine transformation, and Projective transformation. All of these geometric transformations can be represented with a $\mathbb{R}^{3 \times 3}$ matrix with varying degrees of freedom. However, enforcing equivariance and invariance for a continuous space of geometric transformations, \mathcal{T} , is difficult and may even lead to suboptimal solutions. To overcome this issue, in this work, we quantize the *complete* space of affine transformations. We approximate $\mathcal{D}_{\mathcal{T}}$ by dividing it into M discrete

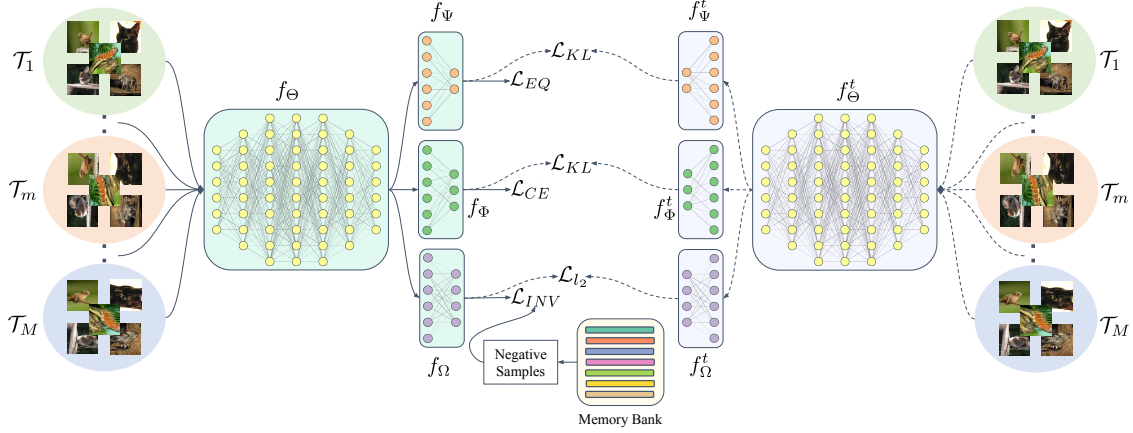


Figure 2. *Network Architecture during Training*: A series of transformed inputs (transformed by applying transformations $\mathcal{T}_1 \dots \mathcal{T}_M$) are provided to a shared feature extractor f_Θ . The resulting embedding is forwarded to three parallel heads f_Ψ , f_Φ and f_Ω that focus on learning equivariant features, discriminative class boundaries, and invariant features, respectively. The resulting output representations are distilled from an old copy of the model (*teacher* model on the right) across multiple-heads to further improve the encoded representations. Notably, a dedicated memory bank of negative samples helps stabilize our invariant contrastive learning.

set of transformations. Here, M can be selected based on the nature of the data and computation budget.

For training, we generate M transformed copies of an input image \mathbf{x} by applying all M transformations. Then we combine all of these transformed images together into a single tensor, $\mathbf{x}_{all} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{M-1}\}$. Here, \mathbf{x}_i is the input image \mathbf{x} transformed through i^{th} transformation, \mathcal{T}_i (the subscript of \mathbf{x}_i is dropped in the subsequent discussion for clarity). We send this composite input to the network and optimize for both equivariance and invariance. The training is performed in a multi-task fashion. In addition to the classification head, which is needed for the baseline supervised training, two other heads are added on top of the base learner, as shown in Figure 2. One of these heads is used for enforcing equivariance, and the other is used for enforcing invariance. This multi-task training scheme ensures that the base learner retains both transformation equivariant and invariant features in the output embedding. We explain each component of our inductive loss below.

3.2.1 Enforcing Equivariance

As discussed above, equivariant features help us encode the inherent structure of data that improves generalization of features to new tasks. To enforce equivariance for the set \mathcal{T} comprising of M quantized transformations, we introduce an MLP f_Ψ with parameters Ψ . The role of f_Ψ is to project the output embeddings from the base learner \mathbf{z} into an equivariant space i.e., $f_\Psi : \mathbf{z} \rightarrow \tilde{\mathbf{u}}$, where $\tilde{\mathbf{u}} \in \mathcal{U} \subset \mathbb{R}^M$.

In order to train the network, we create proxy labels without any manual supervision. For a specific transformation, a M dimensional one-hot encoded vector $\mathbf{u} \in \{0, 1\}^M$ (such that $\sum_i \mathbf{u}_i = 1$) is used to represent the label for f_Ψ . Once proxy labels are assigned, training is performed in a

supervised manner with the cross-entropy loss, as follows:

$$\mathcal{L}_{eq} = -\log \frac{\exp(\tilde{\mathbf{u}}_{j:\mathbf{u}_j=1})}{\sum_i \exp(\tilde{\mathbf{u}}_i)}, \text{ s.t., } \tilde{\mathbf{u}} = f_{\Theta, \Psi}(\mathbf{x}). \quad (2)$$

This supervised training with proxy labels in the equivariant space \mathcal{U} ensures that the output embedding \mathbf{z} retains transformation equivariant features.

3.2.2 Enforcing Invariance

While equivariant representations are important to encode the structure in data, they may not be optimal for class discrimination. This is because the transformations we consider are nuisance variations that do not change the image class, therefore a good feature extractor should also encode representations that are independent of these input variations. To enforce invariance to the set \mathcal{T} consisting of M quantized transformations, we introduce another MLP f_Ω with parameters Ω . The role of f_Ω is to project the output embeddings from the base learner \mathbf{z} into an invariant space i.e., $f_\Omega : \mathbf{z} \rightarrow \mathbf{v}$ where $\mathbf{v} \in \mathcal{V} \subset \mathbb{R}^D$ and D is the dimension of the invariant embedding.

To optimize for invariance we leverage a contrastive loss [26] for instance discrimination. We enforce invariance by maximizing the similarity between an embedding \mathbf{v}^m corresponding to a transformed image (after undergoing m^{th} transformation \mathcal{T}_m), and the reference embedding \mathbf{v}^0 (embedding from the original image without applying any transformation \mathcal{T}). Importantly, we note that selecting negatives within a batch is not sufficient to obtain discriminant representations [79, 48]. We employ a memory bank in our contrastive loss to sample more *negative samples* without arbitrarily increasing the batch size. Further, the memory bank allows a stable convergence behavior [48]. Our learning objective is as follows:

$$\mathcal{L}_{in} = -\frac{1}{M} \sum_{m=0}^{M-1} \log(h(\mathbf{v}^r, \mathbf{v}^m)) \begin{cases} m \neq 0 \rightarrow \mathbf{v}^r = \mathbf{v}^0 \\ m = 0 \rightarrow \mathbf{v}^r = \tilde{\mathbf{v}}^0 \end{cases} \quad (3)$$

where, m denotes the transformation index, $\tilde{\mathbf{v}}^0$ represents a previous copy of the reference \mathbf{v}^0 held in the memory and the function $h(\cdot)$ is defined as,

$$h(\mathbf{v}^r, \mathbf{v}^m) = \frac{\exp\left(\frac{s(\mathbf{v}^r, \mathbf{v}^m)}{\tau}\right)}{\exp\left(\frac{s(\mathbf{v}^r, \mathbf{v}^m)}{\tau}\right) + \sum_{\mathbf{v}' \in \mathcal{D}_n} \exp\left(\frac{s(\mathbf{v}', \mathbf{v}^m)}{\tau}\right)}.$$

Here, $s(\cdot)$ is a similarity function, τ is the temperature, and \mathcal{D}_n is the set of *negative samples* drawn from the memory bank for a particular minibatch. Note that we also maximize the similarity between the reference embedding \mathbf{v}^0 and its past representation $\tilde{\mathbf{v}}^0$ which helps stabilize the learning.

3.2.3 Multi-head Distillation

Once the invariant and equivariant representations are learned by our model, we use self-distillation to train a new model using outputs from the previous model as anchor points (Fig. 2). Note that in typical knowledge distillation [31], information is exchanged from a larger model (teacher) to a smaller one (student) by matching their softened outputs. In contrast, the outputs from the same models are matched in the self-distillation [21] where the smooth predictions encode inter-label dependencies, thereby helping the model to learn better representations.

In our case, a simple knowledge distillation by pairing the logits [73] would not ensure the transfer of invariant and equivariant representations learned by the previous model version. Therefore, we extend the idea of logit-based knowledge distillation and employ it to our invariant and equivariant embedding embeddings. Specifically, in parallel to minimizing the Kullback Leibler (KL) divergence for the soft output of supervised classifier head f_Φ , we also minimize the KL divergence between the outputs of the equivariant head f_Ψ . Since the output of our invariant head f_Ω is not a probability distribution, we minimize a \mathcal{L}_2 loss for distilling the knowledge at this head. The overall learning objective for knowledge distillation is as follows:

$$\mathcal{L}_{kd} = \text{KL}(f_{\Theta, \Phi}^t(\mathbf{x}), f_{\Theta, \Phi}(\mathbf{x})) + \text{KL}(f_{\Theta, \Psi}^t(\mathbf{x}), f_{\Theta, \Psi}(\mathbf{x})) + \mathcal{L}_2(f_{\Theta, \Omega}^t(\mathbf{x}), f_{\Theta, \Omega}(\mathbf{x})). \quad (4)$$

Here, $f_{(\cdot, \cdot)}^t$ and $f_{(\cdot, \cdot)}$ are the teacher and student networks for distillation, respectively.

3.2.4 Overall Objective

Finally, we obtain the resultant loss for injecting the desired inductive biases by combining both equivariant \mathcal{L}_{eq} , invari-

ant \mathcal{L}_{in} , and multi-head distillation \mathcal{L}_{kd} losses:

$$\mathcal{L}_{inductive} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_b, \mathbf{v}' \sim \mathcal{D}_n} \left[\mathcal{L}_{eq}(f_{\Theta, \Psi}(\mathbf{x}), \mathbf{u}) + \mathcal{L}_{in}(f_{\Theta, \Omega}(\mathbf{x}), \mathbf{v}') + \mathcal{L}_{kd}(f_{\Theta, \Phi}^t(\mathbf{x}), f_{\Theta, \Psi}^t(\mathbf{x}), f_{\Theta, \Omega}^t(\mathbf{x})) \right].$$

The overall loss is simply a combination of inductive and baseline objectives,

$$\mathcal{L} = \mathcal{L}_{baseline} + \mathcal{L}_{inductive}. \quad (5)$$

3.3. Few-Shot Evaluation

For evaluation, we test our base learner f_Θ by sampling FSL tasks from a held-out test set comprising of images from novel classes. Each FSL task contains a support set and a corresponding query set $\{D_{supp}, D_{query}\}$; both contain images from the same subset of test classes. Using f_Θ , we obtain embeddings for the images of both D_{supp} and D_{query} . Following [73], we train a simple logistic regression classifier based on the image embeddings and the corresponding labels from the D_{supp} . We use that linear classifier to infer the labels of the query embeddings.

4. Experimental Evaluation

Datasets: We evaluate our method on five popular benchmark FSL datasets. Two of these datasets are subset of the CIFAR100 dataset: CIFAR-FS [4] and FC100 [51]. Another two are derivatives of the ImageNet [14] dataset: miniImageNet [77] and tieredImageNet [61]. The CIFAR-FS dataset is constructed by randomly splitting the 100 classes of the CIFAR-100 dataset into 64, 16, and 20 train, validation, and test splits. FC100 dataset makes the FSL task more challenging by making the splits more diverse; the FC100 train, validation, and test splits contain 60, 20, and 20 classes. Following [59], we use 64, 16, and 20 classes of the miniImageNet dataset for training, validation, and testing. The tieredImageNet dataset contains 608 ImageNet classes that are grouped into 34 high-level categories, and we use 20/351, 6/97, and 8/160 categories/classes for training, validation, and testing. We also evaluate our method on the newly proposed Meta-Dataset [76], which contains 10 diverse datasets to make the FSL task more challenging and closer to realistic classification problems.

Implementation Details: Following [73, 47, 51, 41], we use a ResNet-12 network as our base learner to conduct experiments on CIFAR-FS, FC100, miniImageNet, tieredImageNet datasets. Following [73, 41], we also apply Drop-block [22] regularizer to our Resnet-12 base learner. For Meta-Dataset experiments we use a Resnet-18 [29] network as our base learner to be consistent with [73]. We instantiate both of our equivariant and invariant embedding learners (f_Ψ, f_Ω) with an MLP consisting of a single hidden layer. The classifier, f_Φ , is instantiated with a single linear layer.

Methods	Backbone	1-Shot	5-Shot
MAML[19]	32-32-32-32	58.90 \pm 1.90	71.50 \pm 1.00
Proto-Net [†] [66]	64-64-64-64	55.50 \pm 0.70	72.00 \pm 0.60
Relation Net[71]	64-96-128-256	55.00 \pm 1.00	69.30 \pm 0.80
R2D2[4]	96-192-384-512	65.30 \pm 0.20	79.40 \pm 0.10
Shot-Free[60]	ResNet-12	69.20	84.70
TEWAM[55]	ResNet-12	70.40	81.30
Proto-Net [†] [66]	ResNet-12	72.20 \pm 0.70	83.50 \pm 0.50
MetaOptNet[41]	ResNet-12	72.60 \pm 0.70	84.30 \pm 0.50
Boosting[23]	WRN-28-10	73.60 \pm 0.30	86.00 \pm 0.20
Fine-tuning[15]	WRN-28-10	76.58 \pm 0.68	85.79 \pm 0.50
DSN-MR[65]	ResNet-12	75.60 \pm 0.90	86.20 \pm 0.60
MABAS[34]	ResNet-12	73.51 \pm 0.92	85.49 \pm 0.68
RFS-Simple[73]	ResNet-12	71.50 \pm 0.80	86.00 \pm 0.50
RFS-Distill[73]	ResNet-12	73.90 \pm 0.80	86.90 \pm 0.50
Ours	ResNet-12	76.83 \pm 0.82	89.26 \pm 0.58
Ours-Distill	ResNet-12	77.87 \pm 0.85	89.74 \pm 0.57

Table 1. Average 5-way few-shot classification accuracy with 95% confidence intervals on **CIFAR-FS** dataset; [†]trained on train and validation sets. Top two results are shown in red and blue.

Methods	Backbone	1-Shot	5-Shot
Proto-Net [†] [66]	64-64-64-64	35.30 \pm 0.60	48.60 \pm 0.60
Proto-Net [†] [66]	ResNet-12	37.50 \pm 0.60	52.50 \pm 0.60
TADAM[51]	ResNet-12	40.10 \pm 0.40	56.10 \pm 0.40
MetaOptNet[41]	ResNet-12	41.10 \pm 0.60	55.50 \pm 0.60
MTL[69]	ResNet-12	45.10 \pm 1.80	57.60 \pm 0.90
Fine-tuning[15]	WRN-28-10	43.16 \pm 0.59	57.57 \pm 0.55
MABAS[34]	ResNet-12	42.31 \pm 0.75	57.56 \pm 0.78
RFS-Simple[73]	ResNet-12	42.60 \pm 0.70	59.10 \pm 0.60
RFS-Distill[73]	ResNet-12	44.60 \pm 0.70	60.90 \pm 0.60
Ours	ResNet-12	47.38 \pm 0.79	64.43 \pm 0.77
Ours-Distill	ResNet-12	47.76 \pm 0.77	65.30 \pm 0.76

Table 2. Average 5-way few-shot classification accuracy with 95% confidence intervals on **FC100** dataset; [†]trained on train and validation sets. Top two results are shown in red and blue.

We use SGD optimizer with a momentum of 0.9 in all experiments. For CIFAR-FS, FC100, miniImageNet, tiered-ImageNet datasets we set the initial learning rate to 0.05 and use a weight decay of $5e-4$. For experiments on CIFAR-FS, FC100, miniImageNet datasets, we train for 65 epochs; the learning rate is decayed by a factor of 10 after the first 60 epochs. We train for 60 epochs for experiments on the tieredImageNet dataset; the learning rate is decayed by a factor of 10 for 3 times after the first 30 epochs. For Meta-Dataset experiments, we set the initial learning rate to 0.1 and use a weight decay of $1e-4$. We train our method for 90 epochs and decay the learning rate by a factor of 10 every 30 epochs. We use a batch size of 64 in all of our experiments except on Meta-Dataset where the batch size is set to 256 following [73]. For Meta-dataset experiments, we use standard data augmentation which includes random horizontal flip and random resized crop. For all the other dataset experiments we use random crop, color jittering and random hor-

Methods	Backbone	1-Shot	5-Shot
MAML[19]	32-32-32-32	48.70 \pm 1.84	63.11 \pm 0.92
Matching Net [77]	64-64-64-64	43.56 \pm 0.84	55.31 \pm 0.73
Proto-Net [†] [66]	64-64-64-64	49.42 \pm 0.78	68.20 \pm 0.66
Relation Net[71]	64-96-128-256	50.44 \pm 0.82	65.32 \pm 0.70
R2D2[4]	96-192-384-512	51.20 \pm 0.60	68.80 \pm 0.10
SNAIL[47]	ResNet-12	55.71 \pm 0.99	68.88 \pm 0.92
AdaResNet[49]	ResNet-12	56.88 \pm 0.62	71.94 \pm 0.57
TADAM[51]	ResNet-12	58.50 \pm 0.30	76.70 \pm 0.30
Shot-Free[60]	ResNet-12	59.04	77.64
TEWAM[55]	ResNet-12	60.07	75.90
MTL[69]	ResNet-12	61.20 \pm 1.80	75.50 \pm 0.80
MetaOptNet[41]	ResNet-12	62.64 \pm 0.61	78.63 \pm 0.46
Boosting[23]	WRN-28-10	63.77 \pm 0.45	80.70 \pm 0.33
Fine-tuning[15]	WRN-28-10	57.73 \pm 0.62	78.17 \pm 0.49
LEO-trainval [†] [62]	WRN-28-10	61.76 \pm 0.08	77.59 \pm 0.12
Deep DTN[8]	ResNet-12	63.45 \pm 0.86	77.91 \pm 0.62
AFHN[43]	ResNet-18	62.38 \pm 0.72	78.16 \pm 0.56
AWGIM[25]	WRN-28-10	63.12 \pm 0.08	78.40 \pm 0.11
DSN-MR[65]	ResNet-12	64.60 \pm 0.72	79.51 \pm 0.50
MABAS[34]	ResNet-12	65.08 \pm 0.86	82.70 \pm 0.54
RFS-Simple[73]	ResNet-12	62.02 \pm 0.63	79.64 \pm 0.44
RFS-Distill[73]	ResNet-12	64.82 \pm 0.60	82.14 \pm 0.43
Ours	ResNet-12	66.82 \pm 0.80	84.35 \pm 0.51
Ours-Distill	ResNet-12	67.28 \pm 0.80	84.78 \pm 0.52

Table 3. Average 5-way few-shot classification accuracy with 95% confidence intervals on **miniImageNet** dataset; [†]trained on train and validation sets. Top two results are shown in red and blue.

izontal flip for data augmentation following [73, 41]. Consistent with [73], we use a temperature coefficient of 4.0 for our knowledge distillation experiments. For all datasets, we perform one stage of distillation. We sample 600 FSL tasks to report our scores on all datasets except Meta-Dataset.

For our geometric transformations, we sample from a complete space of similarity transformation and use four rotation transformations: $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$, two scaling transformations: $\{0.67, 1.0\}$ and three aspect ration transformations: $\{0.67, 1.0, 1.33\}$. These geometric transformations can be applied in any combination. For all of our experiments, we set the total number of applied transformations to 16. Additional details and experiments with more geometric transformations are included in the supplementary materials. For the contrastive loss, we use a memory bank that stores 64-dimensional embedding of instances; we sample 6400 negative samples from the memory bank for each mini-batch and set the value of τ to 1.0.

4.1. Results

We present our results on five popular benchmark FSL datasets in Table 1-5 which demonstrates that even without multi-head distillation our proposed method consistently outperforms the current state-of-the-art (SOTA) FSL methods on both 5-way 1-shot and 5-way 5-shot tasks. By virtue of our novel representation learning approach which re-

Methods	Backbone	1-Shot	5-Shot
MAML[19]	32-32-32-32	51.67 \pm 1.81	70.30 \pm 1.75
Proto-Net [†] [66]	64-64-64-64	53.31 \pm 0.89	72.69 \pm 0.74
Relation Net[71]	64-96-128-256	54.48 \pm 0.93	71.32 \pm 0.78
Shot-Free[60]	ResNet-12	63.52	82.59
MetaOptNet[41]	ResNet-12	65.99 \pm 0.72	81.56 \pm 0.53
Boosting[23]	WRN-28-10	70.53 \pm 0.51	84.98 \pm 0.36
Fine-tuning[15]	WRN-28-10	66.58 \pm 0.70	85.55 \pm 0.48
LEO-trainval [†] [62]	WRN-28-10	66.33 \pm 0.05	81.44 \pm 0.09
AWGIM[25]	WRN-28-10	67.69 \pm 0.11	82.82 \pm 0.13
DSN-MR[65]	ResNet-12	67.39 \pm 0.82	82.85 \pm 0.56
RFS-Simple[73]	ResNet-12	69.74 \pm 0.72	84.41 \pm 0.55
RFS-Distill[73]	ResNet-12	71.52 \pm 0.69	86.03 \pm 0.49
Ours	ResNet-12	71.87 \pm 0.89	86.82 \pm 0.58
Ours-Distill	ResNet-12	72.21 \pm 0.90	87.08 \pm 0.58

Table 4. Average 5-way few-shot classification accuracy with 95% confidence intervals on **tieredImageNet** dataset; [†]trained on train and validation sets. Top two results are shown in red and blue.

tains both the transformation invariant and equivariant features in the learned embeddings, our proposed method improves over the baseline RFS-Simple [73] method across all datasets by 2-5% for both 1-shot and 5-shot tasks. To be more specific, our method outperforms the current best results on CIFAR-FS dataset (Table 1) by 1.3% in the 1-shot task whereas for the 5-shot task it improves the score by 2.8%. However, unlike [15], which achieves the current best results on the CIFAR-FS 1-shot task, we do not perform any transductive fine-tuning. For FC100 dataset (Table 2) we observe an even bigger improvement; 2.7% and 4.4% for 1 and 5-shot, respectively. We see similar trends in mini-ImageNet and tieredImageNet (Table 3,4) where we consistently improve over the current SOTA methods by 0.7-2.2%.

For the Meta-Dataset [76], we train our model on the ILSVRC train split and test on 10 diverse datasets. Our results in Table 5 demonstrate that our method outperforms the fo-Proto-MAML [76] across all 10 datasets. Even without multi-head distillation, we outperform both simple and distilled version of the RFS method on 6 out of 10 datasets. Overall, we perform favorably well against the RFS, achieving a new SOTA result on the challenging Meta-Dataset.

4.2. Ablations

To study the contribution of different components of our method we do a thorough ablation study on three benchmark FSL datasets: miniImageNet, CIFAR-FS, and FC100 (Table 6). On these three datasets, our baseline supervised training achieves 62.02%, 71.50%, and 42.60% average accuracy on 5-way 1-shot task respectively; which is the same as RFS-Simple [73]. By enforcing invariance we obtain 2.62%, 2%, and 3.5% improvements respectively. Likewise, enforcing equivariance gives 4.07%, 4.87%, and 4.13% improvements over the baseline respectively. On the

Dataset	fo-Proto-MAML	RFS		Ours	Ours-Distill
		Simple	Distill		
ILSVRC	49.53	60.14	61.48	60.64	61.36
Omniglot	63.37	64.92	64.31	65.55	65.53
Aircraft	55.95	63.12	62.32	65.65	66.58
Birds	68.66	77.69	79.47	77.84	78.23
Textures	66.49	78.59	79.28	81.07	80.42
Quick Draw	51.52	62.48	60.83	57.91	59.02
Fungi	39.96	47.12	48.53	49.26	49.50
VGG Flower	87.15	91.60	91.00	92.06	92.66
Traffic Signs	48.83	77.51	76.33	78.92	79.92
MSCOCO	43.74	57.00	59.28	55.07	55.68
Mean Accuracy	57.52	68.02	68.28	68.40	68.89

Table 5. Results on Meta-Dataset. Average accuracy (%) is reported with variable number of ways and shots, following the setup in [76]. 1000 tasks are sampled for evaluation. Top two results are shown in red and blue.

other hand, we get even bigger improvements by simultaneously optimizing for both equivariance and invariance; achieving 4.8%, 5.33%, and 4.78% improvements on top of the baseline supervised training. Besides, joint training gives 1.3%-3.3% improvement over only invariance training and 0.5%-0.7% improvement in comparison to only equivariance training. We observe similar trends for 5-way 5-shot task. This consistent improvement across all datasets for both tasks empirically validates our claim that joint optimization for both equivariance and invariance is beneficial for FSL tasks. Our ablation study also shows that the multi-head distillation improves the performance over the standard logit-level supervised distillation across all datasets.

Effect of the number of Transformations: To investigate the effect of the total number of applied transformations, we perform an ablation study on the CIFAR-FS validation set by varying the number of transformations, M . We present the results in Table 7, which demonstrates that initially, the performance of our method improves with the increasing M . However, the performance starts to saturate beyond a particular point. We hypothesize that the performance for an increasing number of transformations decreases since discriminating a higher number of transformations is more difficult and the model spends more representation capability for solving this harder task. A similar trend is observed in [24], where increasing the number of recognizable rotations does not lead to better performance. Based on Table 7 results, we set the value of M to 16 for all of our experiments and do not *tune* the M value from dataset to dataset.

4.3. Analysis

We do a t-SNE visualization of the output embeddings from f_{Θ} for the test images of miniImageNet to demonstrate the effectiveness of our method (see Fig. 3). We observe that the base learner trained in a supervised manner can retain good class discrimination even for unseen test classes. However, as evident in Fig. 3, the class boundaries are not precise and compact. Enforcing invariance on top of the

Method	miniImageNet, 5-Way		CIFAR-FS, 5-Way		FC100, 5-Way	
	1-Shot	5-Shot	1-Shot	5-Shot	1-Shot	5-Shot
Baseline Training	62.02 \pm 0.63	79.64 \pm 0.44	71.50 \pm 0.80	86.00 \pm 0.50	42.60 \pm 0.70	59.10 \pm 0.60
Ours with only Invariance	64.64 \pm 0.80	82.59 \pm 0.54	73.50 \pm 0.86	87.55 \pm 0.61	46.10 \pm 0.78	63.18 \pm 0.76
Ours with only Equivariance	66.09 \pm 0.80	84.03 \pm 0.53	76.37 \pm 0.83	89.08 \pm 0.58	46.73 \pm 0.79	64.09 \pm 0.75
Ours with Equi and Invar (W/O KD)	66.82 \pm 0.80	84.35 \pm 0.51	76.83 \pm 0.82	89.26 \pm 0.58	47.38 \pm 0.79	64.43 \pm 0.77
Ours with Supervised KD	66.95 \pm 0.78	84.39 \pm 0.52	76.92 \pm 0.85	89.34 \pm 0.57	47.70 \pm 0.81	65.09 \pm 0.76
Ours Full	67.28 \pm 0.80	84.78 \pm 0.52	77.87 \pm 0.85	89.74 \pm 0.57	47.76 \pm 0.77	65.30 \pm 0.76

Table 6. Ablation study on **miniImageNet**, **CIFAR-FS**, and **FC100** datasets.

M	Description	1-Shot	5-Shot
3	Aspect-Ratio	65.13 \pm 0.93	81.22 \pm 0.66
4	Rotation	66.56 \pm 0.92	82.64 \pm 0.64
8	Rotation, Scale	67.46 \pm 0.92	82.80 \pm 0.64
12	Aspect-Ratio, Rotation	68.04 \pm 0.93	83.48 \pm 0.64
16	Aspect-Ratio, Rotation, Scale	68.20 \pm 0.92	83.63 \pm 0.62
20	Aspect-Ratio, Rotation, Scale	68.07 \pm 0.90	83.53 \pm 0.61

Table 7. Ablation Study on **CIFAR-FS** validation set with different values of M . We choose $M = 16$ for all the experiments.

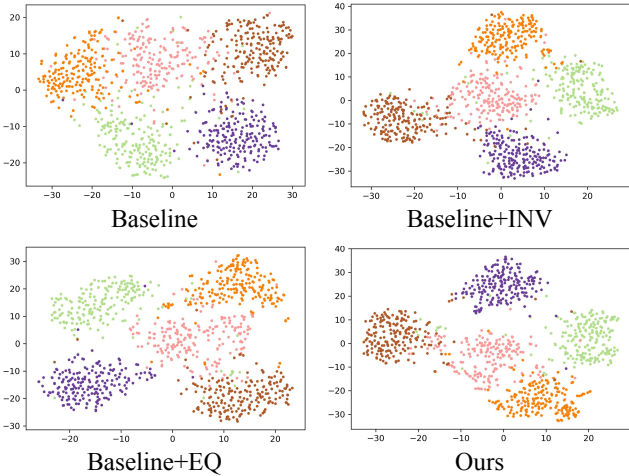


Figure 3. t-SNE visualization of features for 1000 randomly sampled images from 5 randomly selected test classes of **miniImageNet** dataset. In our case, the learned embeddings provide better discrimination for unseen test classes.

base learner leads to more compact class boundaries; however, the sample embeddings of different classes are still relatively closer to one another. On the other hand, enforcing equivariance leads to class representations that are well spread out since it retains the transformation equivariant information in the embedding space. Finally, our proposed method takes advantage of both of these complementary properties and generates embeddings that lead to more compact clusters and discriminative class boundaries.

4.4. Alternate Self-Supervision Losses

In Table 8, to further analyze the performance improvement of our method, we conduct a set of experiments where commonly used self-supervised objectives like solving jig-

Method	1-Shot	5-Shot
Baseline Training	62.02 \pm 0.63	79.64 \pm 0.44
Baseline + Jigsaw Puzzle [50]	63.98 \pm 0.79	81.08 \pm 0.55
Baseline + Location Pred [70]	64.39 \pm 0.81	81.75 \pm 0.54
Baseline + Context Pred [17]	64.72 \pm 0.79	81.83 \pm 0.54
Baseline + Rotation [24]	65.25 \pm 0.80	82.85 \pm 0.54
Ours (W/O KD)	66.82 \pm 0.80	84.35 \pm 0.51

Table 8. FSL with different SSL objectives on **miniImageNet** dataset.

saw puzzles [50], patch location prediction [70], context prediction [17], rotation classification [24] are added on top of the base learner as an auxiliary task. We found that our proposed method which aims to learn representations that retain both transformation invariant and equivariant information outperforms all of these SSL tasks by a good margin. Besides, we have noticed that the patch-based SSL tasks [50, 70, 17] generally underperform in comparison to SSL tasks that rely on changing the global statistics of the image while maintaining the local statistics; this conclusion is in line with the experimental results from [23].

5. Conclusion

In this work, we explored a set of inductive biases that help us learn highly discriminative and transferable representations for FSL. Specifically, we showed that simultaneously learning equivariant and invariant representations to a set of generic transformations results in retaining a complementary set of features that work well for novel classes. We also designed a novel multi-head knowledge distillation objective which delivers additional gains. We conducted extensive ablation to empirically validate our claim that joint optimization for invariance and equivariance leads to more generic and transferable features. We obtained new state-of-the-art results on four popular benchmark FSL datasets as well as on the newly proposed challenging Meta-Dataset.

Acknowledgements This research is based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA R&D Contract No. D17PC00345. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the offi-

cial policies or endorsements, either expressed or implied, of the ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.

A. Supplementary Materials Overview

In the supplementary materials we include the following: additional details about the applied geometric transformations (Section B), additional results with the transformations sampled from the complete space of affine transformations (Section C), ablation study on the coefficient of inductive loss (Section D), ablation study on the temperature of knowledge distillation (Section E), effect of successive self knowledge distillation (Section F), and effect of enforcing invariance and equivariance for supervised classification (Section G).

B. Geometric Transformations

For our geometric transformations, we sample from a complete space of similarity transformation and use four rotation transformations: $\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$, two scaling transformations: $\{0.67, 1.0\}$ and three aspect ratio transformations: $\{0.67, 1.0, 1.33\}$. Different combinations of these transformations lead to different values of M (total number of applied transformations). An ablation study on the value of M is included in section 4.2 of the main paper. In Table 9 we include the complete description of different values of M that we use in our experiments.

C. Additional Results with Affine Transformations

We perform a set of experiments where the objective is to sample geometric transformation from the complete space of affine transformations. To this end, we quantize the affine transformation space according to Table 10. This leads to 972 distinct geometric transformations. Since it's not feasible to apply all the 972 transformations on an input image x to obtain the input tensor $x_{all} = \{x_0, x_1, \dots, x_{971}\}$, we randomly sample 10 geometric transformations from the set of 972 transformations. We apply these randomly sampled 10 geometric transformations on an input image x and generate the input tensor x_{all} . The results of these experiments are presented in Table 11. From Table 11 it's evident that training with either invariance or equivariance improves over the baseline training for both 1 and 5 shot tasks (2.5-3.7% improvement). Joint optimization for both invariance and equivariance provides additional improvement of $\sim 1\%$. Even though the experiments with geometric transformations sampled from the complete affine transformation space do not improve over the training with $M = 16$ (description of $M = 16$ is available in Table 9), the experiments demonstrate consistent improvement when

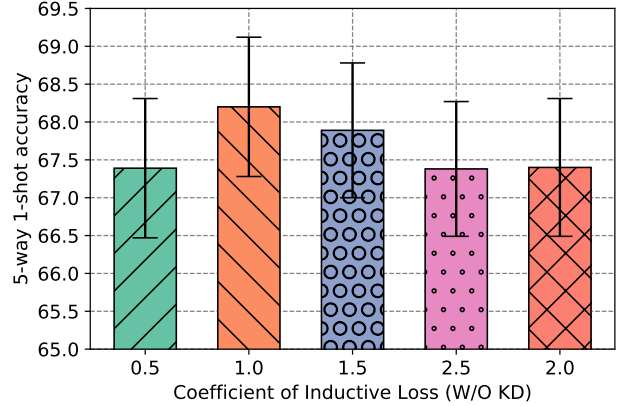


Figure 4. Ablation study on **CIFAR-FS** validation set with different coefficients of the inductive loss (W/O KD); the reported score is average 5-way 1-shot classification accuracy with 95% confidence intervals.

both invariance and equivariance are enforced simultaneously. This provides additional support for our claim that enforcing both invariance and equivariance is beneficial for learning good general representations for solving challenging FSL tasks.

D. Ablation Study for Coefficient of Inductive Loss

We conduct an ablation study to measure the effect of different values of the coefficient of inductive loss (without multi-head distillation) on the CIFAR-FS [4] validation set; the results of 5-way 1-shot FSL tasks are presented in fig. 4. From fig. 4 it is evident that the proposed method is fairly robust to the different values of the coefficient of the inductive loss. However, the best performance is obtained when we set the loss coefficient to 1.0. Based on this ablation study, we use a loss coefficient of 1.0 for the inductive loss in all of our experiments.

E. Ablation Study for Knowledge Distillation Temperature

To analyse the effect of knowledge distillation temperature (for Kullback Leibler (KL) divergence losses) we conduct an ablation study on the validation set of CIFAR-FS [4] dataset. From fig. 5 we can observe that the proposed method with multi-head distillation objective is not very sensitive to the temperature coefficient of knowledge distillation. The proposed method achieves similar performance on the CIFAR-FS validation set when the value of distillation temperature is set to 4.0 and 5.0. Based on this ablation study and to be consistent with [73], we set the value of the coefficient of knowledge distillation temperature to 4.0 in all of our experiments.

M	Description
3	$\text{AR}:\{0.67, 1.0, 1.33\}$
4	$\text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$
8	$\text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\} \times \text{S}:\{0.67, 1.0\}$
12	$\text{AR}:\{0.67, 1.0, 1.33\} \times \text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$
16	$(\text{AR}:\{0.67, 1.0, 1.33\} \times \text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}) \cup (\text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\} \times \text{S}:\{0.67\})$
20	$(\text{AR}:\{0.67, 1.0, 1.33\} \times \text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}) \cup (\text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\} \times \text{S}:\{0.67\} \times \text{AR}:\{0.67, 1.33\})$
24	$\text{AR}:\{0.67, 1.0, 1.33\} \times \text{ROT}:\{0^\circ, 90^\circ, 180^\circ, 270^\circ\} \times \text{S}:\{0.67, 1.0\}$

Table 9. Complete description of different values of M based on different combination of aspect ratio (AR), rotation (ROT), and scaling (S) transformations.

Transformation	Quantized Values
Rotation	$\{0^\circ, 90^\circ, 180^\circ, 270^\circ\}$
Translation(X)	$\{-0.2, 0.0, 0.2\}$
Translation(Y)	$\{-0.2, 0.0, 0.2\}$
Scale	$\{0.67, 1.0, 1.33\}$
Aspect-Ratio	$\{0.67, 1.0, 1.33\}$
Shear	$\{-20^\circ, 0^\circ, 20^\circ\}$

Table 10. Quantization of the space of Affine transformations.

Method	1-Shot	5-Shot
Baseline Training	62.02 ± 0.63	79.64 ± 0.44
Ours with only Invar (affine)	65.55 ± 0.81	82.17 ± 0.52
Ours with only Equi (affine)	65.70 ± 0.79	82.47 ± 0.53
Ours with Equi and Invar (affine)	66.82 ± 0.79	82.96 ± 0.53
Ours with Equi and Invar ($M=16$)	66.82 ± 0.80	84.35 ± 0.51

Table 11. Average 5-way few-shot classification accuracy with 95% confidence intervals on **miniImageNet** dataset; trained with different geometric transformations.

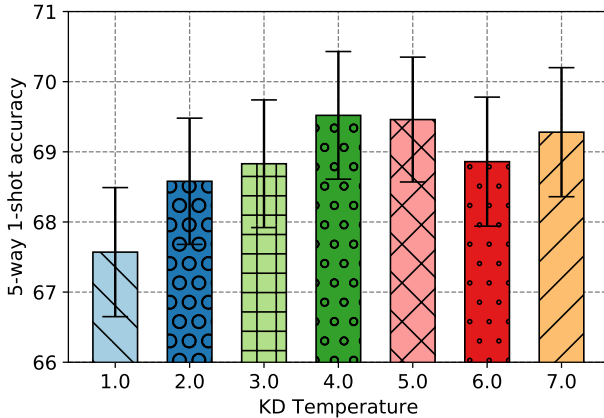


Figure 5. Ablation study on **CIFAR-FS** validation set with different values of knowledge distillation temperature; the reported score is average 5-way 1-shot classification accuracy with 95% confidence intervals.

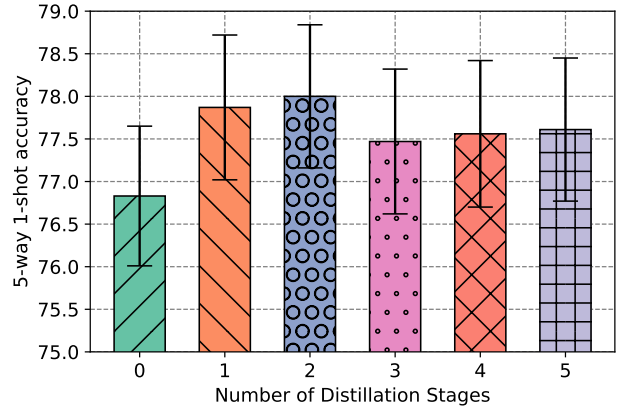


Figure 6. Evaluation of different knowledge distillation stages on **CIFAR-FS** dataset; the reported score is average 5-way 1-shot classification accuracy with 95% confidence intervals.

F. Effect of Successive Distillation

In all of our experiments, we use only one stage of multi-head knowledge distillation. To further investigate the effect of knowledge distillation we perform multiple stages of self knowledge distillation on **CIFAR-FS** [4] dataset. The results are presented in fig. 6. Here, the 0^{th} distillation stage is the base learner trained with only the supervised baseline loss ($\mathcal{L}_{baseline}$), equivariant loss (\mathcal{L}_{eq}), and invariant loss (\mathcal{L}_{in}). From fig. 6, we observe that the performance in the FSL task improves for the first 2 stages of distillation, after that the performance saturates. Besides, the improvement from stage 1 to stage 2 is minimal ($\sim 0.1\%$). Therefore, to make the proposed method more computationally efficient we perform only one stage of distillation in all of our experiments.

G. Invariance and Equivariance for Supervised Classification

To demonstrate the effectiveness of complementary strengths of invariant and equivariant representations we conduct fully supervised classification experiments on

Method	Error Rate (%)
Supervised Baseline	18.78
Ours with only Invariance	18.56
Ours with only Equivariance	16.95
Ours with Equi and Invar (W/O KD)	16.84

Table 12. Results with invariance and equivariance for supervised classification on **CIFAR-100** dataset.

benchmark CIFAR-100 dataset [37]. For these experiments, we use the standard Wide-Resnet-28-10 [80] architecture as the backbone. For training, we use an SGD optimizer with an initial learning rate of 0.1. We set the momentum to 0.9 and use a weight decay of $5e-4$. For all the experiments, the training is performed for 200 epochs where the learning rate is decayed by a factor of 5 at epochs 60, 120, and 160. We use a batch size of 128 for all the experiments as well as a dropout rate of 0.3. The training augmentations include standard data augmentations: random crop and random horizontal flip. For enforcing invariance and equivariance, we set the value of M to 12 for computational efficiency; description of $M = 12$ is available in Table 9. We do not perform knowledge distillation for these experiments. The results of these experiments are presented in Table 12.

From Table 12, we can notice that enforcing invariance provides little improvement (0.2%) over the supervised baseline. This is expected since the train and test data is coming from the same distribution and same set of classes; making the class boundaries compact (for seen classes) doesn’t provide that much additional benefit. However, in the case of FSL we observe that enforcing invariance over baseline provides 2.62%, 2%, and 3.5% improvement for miniImageNet [77], CIFAR-FS [4], and FC100 [51] datasets respectively (section 4.2 of main text). On the other hand, enforcing equivariance for supervised classification provides better improvement (1.8%) since it helps the model to better learn the structure of data. Even though enforcing equivariance provides noticeable improvement for supervised classification, in the case of FSL we obtain a much bigger improvement of 4.07%, 4.87%, and 4.13% for miniImageNet [77], CIFAR-FS [4], and FC100 [51] datasets respectively (section 4.2 of main text). Finally, joint optimization for both invariance and equivariance achieves the best performance and provides minimal but consistent improvement of 0.1% over enforcing only equivariance. However, joint optimization provides a much larger improvement on FSL tasks (see section 4.2 of the main text). From these experiments, we conclude that, although enforcing both invariance and equivariance is beneficial for supervised classification, injecting these inductive biases is more crucial for FSL tasks since the inductive inference for FSL tasks is more challenging (inference on unseen/novel classes).

References

- [1] Kelsey Allen, Evan Shelhamer, Hanul Shin, and Joshua Tenenbaum. Infinite mixture prototypes for few-shot learning. volume 97 of *Proceedings of Machine Learning Research*, pages 232–241, Long Beach, California, USA, 09–15 Jun 2019. PMLR. 1
- [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018. 2
- [3] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A Raffel. Mixmatch: A holistic approach to semi-supervised learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 5049–5059. Curran Associates, Inc., 2019. 1
- [4] Luca Bertinetto, Joao F. Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *International Conference on Learning Representations*, 2019. 5, 6, 9, 10, 11
- [5] Alexander Buslaev, Vladimir I Iglovikov, Eugene Khvedchenya, Alex Parinov, Mikhail Druzhinin, and Alexandr A Kalinin. Albumentations: fast and flexible image augmentations. *Information*, 11(2):125, 2020. 2
- [6] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 1
- [7] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 1
- [8] Mengting Chen, Yuxin Fang, Xinggang Wang, Heng Luo, Yifeng Geng, Xinyu Zhang, C. Huang, W. Liu, and Bo Wang. Diversity transfer network for few-shot learning. In *AAAI*, 2020. 6
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020. 1, 3
- [10] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *International Conference on Learning Representations*, 2019. 2, 3
- [11] N. Cohen and A. Shashua. Inductive bias of deep convolutional networks through pooling geometry. *International Conference on Learning Representations (ICLR)*, 2017. 2
- [12] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999, 2016. 2
- [13] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation

- policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 2
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009. 1, 5
- [15] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations*, 2020. 2, 3, 6, 7
- [16] Sander Dieleman, Jeffrey De Fauw, and Koray Kavukcuoglu. Exploiting cyclic symmetry in convolutional neural networks. *arXiv preprint arXiv:1602.02660*, 2016. 2
- [17] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1422–1430, 2015. 1, 3, 8
- [18] Carl Doersch, Ankush Gupta, and Andrew Zisserman. Crosstransformers: spatially-aware few-shot transfer. *Advances in Neural Information Processing Systems*, 2020. 3
- [19] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. 1, 2, 3, 6, 7
- [20] Marc Finzi, Samuel Stanton, Pavel Izmailov, and Andrew Gordon Wilson. Generalizing convolutional neural networks for equivariance to lie groups on arbitrary continuous data. *arXiv preprint arXiv:2002.12880*, 2020. 2
- [21] Tommaso Furlanello, Zachary Chase Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born-again neural networks. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1602–1611. PMLR, 2018. 5
- [22] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. In *Advances in Neural Information Processing Systems*, pages 10727–10737, 2018. 5
- [23] Spyros Gidaris, Andrei Bursuc, Nikos Komodakis, Patrick Pérez, and Matthieu Cord. Boosting few-shot visual learning with self-supervision. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8059–8068, 2019. 3, 6, 7, 8
- [24] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *International Conference on Learning Representations*, 2018. 3, 7, 8
- [25] Yiluan Guo and Ngai-Man Cheung. Attentive weights generation for few shot learning via information maximization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13499–13508, 2020. 6, 7
- [26] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006. 4
- [27] Munawar Hayat, Salman Khan, Syed Waqas Zamir, Jianbing Shen, and Ling Shao. Gaussian affinity for max-margin class imbalanced learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6469–6479, 2019. 2
- [28] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 1
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 5
- [30] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. 2
- [31] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. 5
- [32] Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11719–11727, 2019. 1
- [33] Salman Khan, Hossein Rahmani, Syed Afaq Ali Shah, and Mohammed Bannamoun. A guide to convolutional neural networks for computer vision. *Synthesis Lectures on Computer Vision*, 8(1):1–207, 2018. 2
- [34] Jaekyeom Kim, Hyoungseok Kim, and Gunhee Kim. Model-agnostic boundary-adversarial sampling for test-time generalization in few-shot learning. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 599–617, Cham, 2020. Springer International Publishing. 6
- [35] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2. Lille, 2015. 2
- [36] Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better? In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2661–2671, 2019. 1
- [37] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research). 11
- [38] Dmitry Laptev, Nikolay Savinov, Joachim M Buhmann, and Marc Pollefeys. Ti-pooling: transformation-invariant pooling for feature learning in convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 289–297, 2016. 2
- [39] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European conference on computer vision*, pages 577–593. Springer, 2016. 3
- [40] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361(10):1995, 1995. 2
- [41] Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex op-

- timization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10657–10665, 2019. 5, 6, 7
- [42] Jan Eric Lenssen, Matthias Fey, and Pascal Libuschewski. Group equivariant capsule networks. In *Advances in Neural Information Processing Systems*, pages 8844–8853, 2018. 2
- [43] Kai Li, Yulun Zhang, Kunpeng Li, and Yun Fu. Adversarial feature hallucination networks for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13470–13479, 2020. 6
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2014. 1
- [45] Carlos Medina, Arnout Devos, and Matthias Grossglauser. Self-supervised prototypical transfer learning for few-shot classification. *arXiv preprint arXiv:2006.11325*, 2020. 3
- [46] Michael Milford, Chunhua Shen, Stephanie Lowry, Niko Suenderhauf, Sareh Shirazi, Guosheng Lin, Fayao Liu, Edward Pepperell, Cesar Lerma, Ben Upcroft, et al. Sequence searching with deep-learned depth for condition-and viewpoint-invariant route-based place recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 18–25, 2015. 2
- [47] Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *International Conference on Learning Representations*, 2018. 5, 6
- [48] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020. 4
- [49] Tsendsuren Munkhdalai, Kingdi Yuan, Soroush Mehri, and Adam Trischler. Rapid adaptation with conditionally shifted neurons. volume 80 of *Proceedings of Machine Learning Research*, pages 3664–3673, Stockholmmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. 6
- [50] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016. 3, 8
- [51] Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. In *Advances in Neural Information Processing Systems*, pages 721–731, 2018. 2, 5, 6, 11
- [52] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 3
- [53] Lutz Prechelt. Early stopping-but when? In *Neural Networks: Tricks of the trade*, pages 55–69. Springer, 1998. 2
- [54] Guo-Jun Qi, Liheng Zhang, Feng Lin, and Xiao Wang. Learning generalized transformation equivariant representations via autoencoding transformations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [55] Limeng Qiao, Yemin Shi, Jia Li, Yaowei Wang, Tiejun Huang, and Yonghong Tian. Transductive episodic-wise adaptive metric for few-shot learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3603–3612, 2019. 6
- [56] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. In *International Conference on Learning Representations*, 2020. 1
- [57] Jathushan Rajasegaran, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Mubarak Shah. Self-supervised knowledge distillation for few-shot learning. <https://arxiv.org/abs/2006.09785>, 2020. 3
- [58] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2017. 2
- [59] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 5
- [60] Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Few-shot learning with embedded class models and shot-free meta training. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 331–339, 2019. 6, 7
- [61] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676*, 2018. 5
- [62] Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *International Conference on Learning Representations*, 2019. 1, 6, 7
- [63] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *Advances in neural information processing systems*, pages 3856–3866, 2017. 2
- [64] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 806–813, 2014. 1
- [65] Christian Simon, Piotr Koniusz, Richard Nock, and Mehrtash Harandi. Adaptive subspaces for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4136–4145, 2020. 6, 7
- [66] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017. 1, 2, 3, 6, 7
- [67] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. 2
- [68] Jong-Chyi Su, Subhransu Maji, and Bharath Hariharan. When does self-supervision improve few-shot learning? *European Conference on Computer Vision*, 2020. 3
- [69] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 403–412, 2019. 6

- [70] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A Efros. Unsupervised domain adaptation through self-supervision. *arXiv preprint arXiv:1909.11825*, 2019. [8](#)
- [71] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. Learning to compare: Relation network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1199–1208, 2018. [1](#), [6](#), [7](#)
- [72] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. [1](#)
- [73] 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 preprint arXiv:2003.11539*, 2020. [2](#), [3](#), [5](#), [6](#), [7](#), [9](#)
- [74] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996. [2](#)
- [75] Florian Tramèr and Dan Boneh. Adversarial training and robustness for multiple perturbations. In *Advances in Neural Information Processing Systems*, pages 5866–5876, 2019. [2](#)
- [76] 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 *International Conference on Learning Representations*, 2020. [2](#), [5](#), [7](#)
- [77] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3630–3638. Curran Associates, Inc., 2016. [1](#), [3](#), [5](#), [6](#), [11](#)
- [78] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016. [2](#)
- [79] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018. [4](#)
- [80] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard C. Wilson and William A. P. Smith, editors, *Proceedings of the British Machine Vision Conference (BMVC)*, pages 87.1–87.12. BMVA Press, September 2016. [11](#)
- [81] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722, 2018. [1](#)
- [82] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. [3](#)
- [83] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017. [1](#)