

Supervised Masked Knowledge Distillation for Few-Shot Transformers

Han Lin*, Guangxing Han*[†], Jiawei Ma, Shiyuan Huang, Xudong Lin, Shih-Fu Chang
Columbia University

{hl13199, gh2561, jiawei.m, sh3813, xl2798, sc250}@columbia.edu

Abstract

Vision Transformers (ViTs) emerge to achieve impressive performance on many data-abundant computer vision tasks by capturing long-range dependencies among local features. However, under few-shot learning (FSL) settings on small datasets with only a few labeled data, ViT tends to overfit and suffers from severe performance degradation due to its absence of CNN-like inductive bias. Previous works in FSL avoid such problem either through the help of self-supervised auxiliary losses, or through the dextile uses of label information under supervised settings. But the gap between self-supervised and supervised few-shot Transformers is still unfilled. Inspired by recent advances in self-supervised knowledge distillation and masked image modeling (MIM), we propose a novel **Supervised Masked Knowledge Distillation model (SMKD)** for few-shot Transformers which incorporates label information into self-distillation frameworks. Compared with previous self-supervised methods, we allow intra-class knowledge distillation on both class and patch tokens, and introduce the challenging task of masked patch tokens reconstruction across intra-class images. Experimental results on four few-shot classification benchmark datasets show that our method with simple design outperforms previous methods by a large margin and achieves a new start-of-the-art. Detailed ablation studies confirm the effectiveness of each component of our model. Code for this paper is available here: <https://github.com/HL-hanlin/SMKD>.

1. Introduction

Vision Transformers (ViTs) [19] have emerge as a competitive alternative to Convolutional Neural Networks (CNNs) [35] in recent years, and have achieved impressive performance in many vision tasks including image classification [19, 45, 67, 74], object detection [5, 13, 30–32, 90], and object segmentation [57, 63]. Compared with CNNs, which introduce inductive bias through convolutional kernels with fixed receptive fields [41], the attention layers in

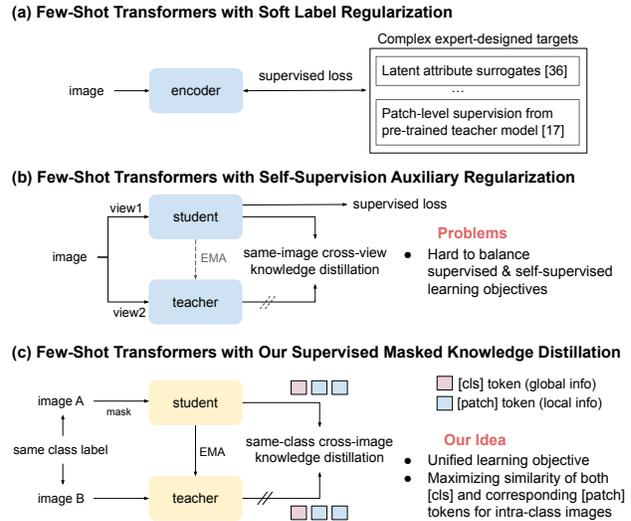


Figure 1. **Comparison of the proposed idea and other existing methods for few-shot Transformers.** Our model mitigates the overfitting issue of few-shot Transformers, by extending the masked knowledge distillation framework into the supervised setting, and enforcing the alignment of [cls] and corresponding [patch] tokens for intra-class images.

ViT allow it to model global token relations to capture long-range token dependencies. However, such flexibility also comes at a cost: ViT is data-hungry and it needs to learn token dependencies purely from data. This property often makes it easy to overfit to datasets with small training set and suffer from severe generalization performance degradation [42, 44]. Therefore, we are motivated to study how to make ViTs generalize well on these small datasets, especially under the few-shot learning (FSL) setting [23, 46, 70] which aims to recognize unseen new instances at test time just from only a few (e.g. one or five) labeled samples from each new categories.

Most of the existing methods mitigate the overfitting issue of few-shot Transformers [17] using various regularizations. For instance, some works utilize label information in a weaker [36], softer [49] way, or use label information efficiently through patch-level supervision [17]. However, these models usually design sophisticated learning targets.

*Equal contribution. [†]Corresponding author.

On the other hand, self-distillation techniques [7, 11], and particularly, the recent masked self-distillation [33, 52, 88], which distills knowledge learned from an uncorrupted image to the knowledge predicted from a masked image, have lead an emerging trend in self-supervised Transformers in various fields [18, 76]. Inspired by such success, recent works in FSL attempt to incorporate self-supervised pretext tasks into the standard supervised learning through auxiliary losses [44, 50, 53], or to adopt a self-supervised pre-training, supervised training two-stage framework to train few-shot Transformers [22, 37]. Compared with traditional supervised methods, self-supervision can learn less biased representations towards base class, which usually leads to better generalization ability for novel classes [47]. However, the two learning objectives of self-supervision and supervision are conflicting and it is hard to balance them during training. Therefore, how to efficiently leverage the strengths of self-supervised learning to alleviate the overfitting issue of supervised training remains a challenge.

In this work, we propose a novel supervised masked knowledge distillation framework (SMKD) for few-shot Transformers, which handles the aforementioned challenge through a natural extension of the self-supervised masked knowledge distillation framework into the supervised setting (shown in Fig. 1). Different from supervised contrastive learning [39] which only utilizes global image features for training, we leverage multi-scale information from the images (both global [cls] token and local [patch] tokens) to formulate the learning objectives, which has been demonstrated to be effective in the recent self-supervised Transformer methods [33, 88]. For global [cls] tokens, we can simply maximize the similarity for intra-class images. However, it is non-trivial and challenging to formulate the learning objectives for local [patch] tokens because we do not have ground-truth patch-level annotations. To address this problem, we propose to estimate the similarity between [patch] tokens across intra-class images using cross-attention, and enforce the alignment of corresponding [patch] tokens. Particularly, *reconstructing masked [patch] tokens across intra-class images increases the difficulty of model learning, thus encouraging learning generalizable few-shot Transformer models by jointly exploiting the holistic knowledge of the image and the similarity of intra-class images.*

As shown in Fig. 2, we compare our model with the existing self-supervised/supervised learning methods. Our model is a natural extension of the supervised contrastive learning method [39] and self-supervised knowledge distillation methods [7, 88]. Thus our model inherits both the advantage of method [39] for effectively leveraging label information, and the advantages of methods [7, 88] for not needing large batch size and negative samples. Meanwhile, the newly-introduced challenging task of masked

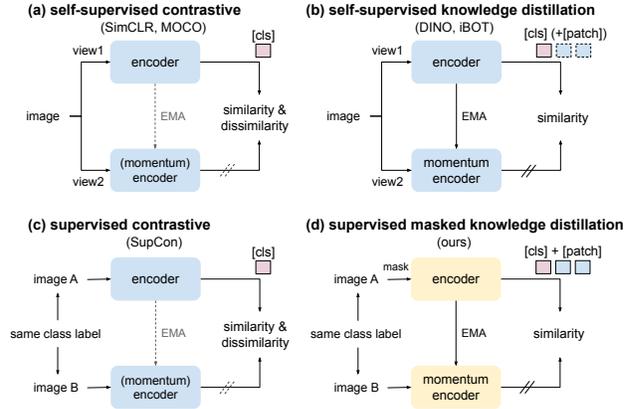


Figure 2. **Comparison of other self-supervised/supervised frameworks.** Our method (d) is a natural extension of (b) and (c), with the newly-introduced challenging task of masked [patch] tokens reconstruction across intra-class images.

[patch] tokens reconstruction across intra-class images makes our method more powerful for learning generalizable few-shot Transformer models.

Compared with contemporary works on few-shot Transformers [17, 36, 37], our framework enjoys several good properties from a practical point of view. (1) Our method does not introduce any additional learnable parameters besides the ViT backbone and projection head, which makes it easy to be combined with other methods [36, 37, 83]. (2) Our method is both effective and training-efficient, with stronger performance and less training time on four few-shot classification benchmarks, compared with [36, 37]. In a nutshell, our main contributions can be summarized as follows:

- We propose a new supervised knowledge distillation framework (SMKD) that incorporates class label information into self-distillation, thus filling the gap between self-supervised knowledge distillation and traditional supervised learning.
- Within the proposed framework, we design two supervised-contrastive losses on both class and patch levels, and introduce the challenging task of masked patch tokens reconstruction across intra-class images.
- Given its simple design, we test our SMKD on four few-shot datasets, and show that it achieves a new SOTA on CIFAR-FS and FC100 by a large margin, as well as competitive performance on *mini*-ImageNet and *tiered*-ImageNet using the simple prototype classification method for few-shot evaluation.

2. Related Work

Few-shot Learning. Few-shot learning aims at fast knowledge transfer from seen base classes to unseen novel classes given only a few labeled samples from each novel class. The meta-learning paradigm, which simulates few-shot tasks episodically to mimic the human-learning process in the

real world, once dominated FSL [26–29, 51, 60, 62, 64, 70]. Surprisingly, recent works [9, 25, 66, 81] have shown that the state-of-the-art meta-learning methods can be outperformed by simple baseline methods using just a distance-based classifier without the complicated design of meta-learning strategies. Therefore, recent methods [21, 36, 49] in FSL start to focus less on meta-learning, and more on learning embeddings with good generalization ability. Our paper follows this trend and proposes a knowledge distillation framework to learn generalizable embeddings.

Vision Transformers in FSL. Variants of Transformer [69] have achieved great success in NLP [15, 56], computer vision [19, 67], and multimodal learning [55]. However, the lack of inductive bias makes Transformer infamous for its data-hungry property, which is especially significant in few-shot settings when the amount of data is limited. A line of works focuses on introducing inductive bias back to the Transformer architecture, including methods using pyramid structure [71, 72], shifted windows [45], and explicit convolutional token embeddings [74]. With this said, some recent works [16, 28, 37, 80] still show that few-shot Transformers have the potential of fast adaptation to novel classes. Our work also studies few-shot Transformers, and we show that our method can work well even on the vanilla ViT structure without explicit incorporation of inductive bias.

Self-Supervision for FSL. Self-supervised learning (SSL) has shown great potential in FSL due to its good generalization ability to novel classes. Previous methods incorporate SSL into FSL in various ways. Some works propose to include self-supervised pretext tasks into the standard supervised learning through auxiliary losses [44, 50, 53]. For example, [44] proposed a regularization loss that extracts additional information from images by predicting the geometric distance between patch tokens. [53] designed their loss by discouraging spatially disordered attention maps based on the idea that objects usually occupy connected regions, and [50] derived their auxiliary loss from self-supervision tasks of rotation and exemplar. Some other works [22, 37], instead, adopt a two-stage procedure by pretraining a model via self-supervision before supervised training. [37] takes advantage of self-supervised training with iBOT [88] as a pretext task, and then uses inner loop token importance reweighting for supervised fine-tuning. [22] initializes its model with a pre-trained self-supervised DINO [7], then trained with supervised cross-entropy loss. Our work follows the second branch of work, but the difference is obvious from previous works. Instead of designing complicated training pipeline or leveraging additional learnable modules at inference time, we use supervised training on self-supervised pre-trained model, and focus mainly on bridging the gap between self-supervised and supervised knowledge distillation with minimum extra design.

Another emerging trend in SSL is Masked Image Mod-

eling (MIM) [3, 33, 52, 88], which aims at recovering the patch-level target (e.g. image pixels, patch features) of the masked content in a corrupted input image. In iBOT [88], the class and patch tokens share the same projection head, thus helping the ViT backbone to learn semantically meaningful patch embeddings, which can be used as a supplement to image-level global self-supervision.

3. Our Approach

We present our Supervised Masked Knowledge Distillation (SMKD) framework in this section. Fig. 3 shows an overview of our model. Task definition for few-shot classification is presented in Sec. 3.1. We explain the formulation of SSL knowledge distillation in Sec. 3.2, and how we extend it into our supervised knowledge distillation framework in Sec. 3.3. Training pipeline is given in Sec. 3.4.

3.1. Learning Task Definition

In few-shot classification, we are given two datasets \mathcal{D}_{base} and \mathcal{D}_{novel} with disjoint class labels \mathcal{C}_{base} and \mathcal{C}_{novel} ($\mathcal{C}_{base} \cap \mathcal{C}_{novel} = \emptyset$). The base dataset \mathcal{D}_{base} , which contains abundant labeled samples, is used to train the feature extractor. And the novel dataset \mathcal{D}_{novel} is then used to sample episodes for prototype estimation and few-shot evaluation. In a standard N-way K-shot task, each episode $\mathcal{D}_{epi} = (\mathcal{D}_S, \mathcal{D}_Q)$ covers N classes from \mathcal{C}_{novel} . The support dataset \mathcal{D}_S , which contains K samples from each class, is used for class prototype estimation, and the query dataset \mathcal{D}_Q is then used for evaluation. This task aims at correctly classifying \mathcal{D}_Q into N classes from sampled episodes. And the main focus of this paper is to train a feature extractor with good generalization ability from \mathcal{D}_{base} .

3.2. Preliminary: SSL with Knowledge Distillation

Our work is inspired by the self-supervised knowledge distillation frameworks proposed recently [7]. Specifically, given an input image x uniformly sampled from the training set \mathcal{I} , random data augmentations are applied to generate two augmented views x^1 and x^2 (*we represent image view indices as superscripts and patch indices as subscripts*), which are then fed into the teacher and student networks. The student network, parameterized by θ_s , consists of an encoder with ViT backbone and a projection head with a 3-layer multi-layer perceptron (MLP) followed by l_2 -normalized bottleneck. The ViT backbone first generates a [cls] token, which is then entered into the projection head and outputs a probability distribution $P_s^{[cls]}$ over K dimensions. The teacher network, parameterized by θ_t , is Exponentially Moving Averaged (EMA) updated by the student network θ_s , which distills its knowledge to the student by minimizing the cross-entropy loss over the outputs of the

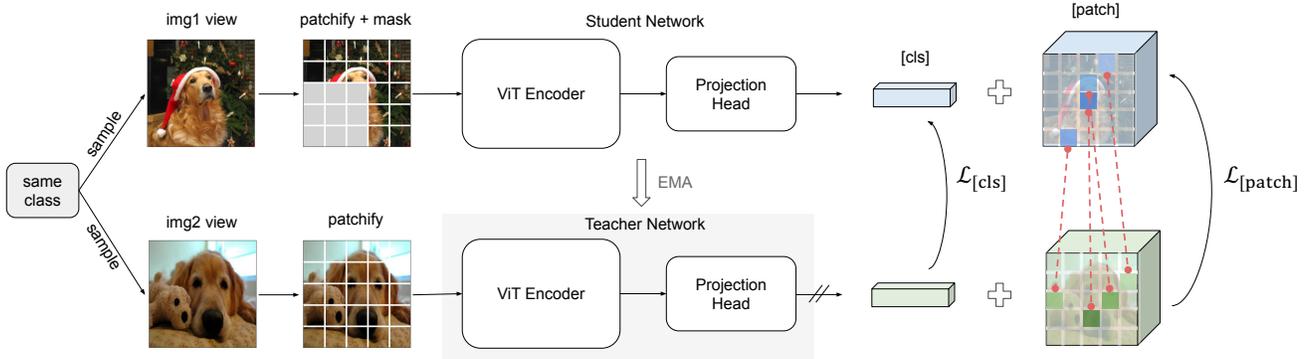


Figure 3. **Overview of our SMKD framework.** Two views are generated from a pair of images sampled (*with replacement*) from the same class. The first view, applied with random blockwise masking, is passed to the student network, while the second unmasked view is passed to the teacher network. Both networks consist of a ViT backbone and a projection head. The parameters of the teacher network are Exponentially Moving Averaged (EMA) updated by the student network. SMKD distills knowledge between intra-class cross-views on both class and patch levels. $\mathcal{L}_{[cls]}$ distills knowledge from $[cls]$ tokens, while $\mathcal{L}_{[patch]}$ distills knowledge from $[patch]$ tokens by finding dense correspondence of matched token pairs (connected by red dashed lines) with highest similarities.

categorical distributions from their projection heads:

$$\mathcal{L}_{[cls]} = \mathcal{H}(P_t^{[cls]}(x^1), P_s^{[cls]}(x^2)) \quad (1)$$

where $\mathcal{H}(x, y) = -x \log y$.

Masked Image Modeling (MIM) [3, 65] can be performed via self-distillation as follows [88]. Given a randomly sampled mask sequence $m \in \{0, 1\}^N$ over an image with N $[patch]$ tokens $x = \{x_i\}_{i=1}^N$, x_i 's with $m_i = 1$ are then replaced by a learnable token embedding $e_{[MASK]}$, which results in a corrupted image $\hat{x} = \{\hat{x}_i\}_{i=1}^N = \{(1 - m_i)x_i + m_i e_{[MASK]}\}_{i=1}^N$. This corrupted image and the original uncorrupted image are fed into the student and teacher networks respectively. The objective of MIM is to recover the masked tokens from the corrupted image, which is equivalent to minimizing the cross-entropy loss between the outputs of the categorical distributions of the student and teacher networks on masked patches:

$$\mathcal{L}_{MIM} = \sum_{i=1}^N m_i \cdot \mathcal{H}(P_t^{[patch]}(x)_i, P_s^{[patch]}(\hat{x})_i) \quad (2)$$

3.3. Supervised Masked Knowledge Distillation

Distill the Class Tokens. Recent self-distillation frameworks [7, 24, 88] distill knowledge on $[cls]$ tokens from cross-view images via Eq.(1). To incorporate label information into such self-supervised frameworks, we further allow knowledge on $[cls]$ tokens to be distilled from intra-class cross-views. This can be achieved by a small extension of the way we sample images. Rather than sampling a single image $x \sim \mathcal{I}$ and generating two views, now we sample two images $x, y \sim \mathcal{I}^c$ (*with replacement*) and generate one view for each of them¹. $\mathcal{I}^c \subseteq \mathcal{I}$ here denotes the set

¹For simplicity of illustration, we show our method with only one augmented view from each image. We sample two views from an image in our implementation. Loss is then averaged over all cross-view pairs.

of images with the same class label c in the training set \mathcal{I} . Specifically, we denote x' and y' as the augmented views generated from image x and y respectively. We apply additional random blockwise masking on x' and denote the resulting corrupted view as \hat{x}' . Then the corrupted view \hat{x}' and uncorrupted view y' are sent to the student and teacher networks separately. Now our supervised-contrastive loss on $[cls]$ tokens becomes:

$$\mathcal{L}_{[cls]} = \mathcal{H}(P_t^{[cls]}(y'), P_s^{[cls]}(\hat{x}')) \quad (3)$$

When x, y are sampled as the same image ($x = y$), our loss performs *masked* self-distillation across two views of the same image similar to Eq.(1). In the other situation when x and y represent different images ($x \neq y$), our loss then minimizes the cross-entropy loss of projected $[cls]$ tokens among all cross-view pairs between images x and y .

Such design has two major advantages. (1) It can be implemented efficiently. Instead of intentionally sampling image pairs from the same class, we just need to look through the images in a mini-batch, find image pairs belonging to the same class, and then apply our loss in Eq.(3). (2) Unlike previous works that use either supervised [49, 87] or self-supervised [8, 34] contrastive losses, our method follows the recent trend in SSL works [7, 10, 24, 82] and avoids the need for negative examples.

Distill the Patch Tokens. Beyond the knowledge distillation on the global $[cls]$ tokens, we introduce the challenging task of masked patch tokens reconstruction across intra-class images, to fully exploit the local details of the image for training. Our main intuition here is based on the following hypothesis: for intra-class images, even though their semantic information can be vastly different at the patch level, there should at least exist some patches that share similar semantic meanings.

For each patch k from an input view \mathbf{y}' sent to the teacher network (with its corresponding token embedding defined as \mathbf{f}_k^t), we need to first find the most similar patch k^+ from the masked view $\hat{\mathbf{x}}'$ of the student network (with its corresponding token embedding defined as $\mathbf{f}_{k^+}^s$), and then perform knowledge distillation between these two matched tokens. The token embedding $\mathbf{f}_{k^+}^s$ represents either the k^+ th [patch] token if patch k^+ is unmasked, or the reconstructed patch token if it is masked.

As we do not have any patch-level annotations, we use cosine similarity to find the best-matched patch of k among all the [patch] tokens in the student network:

$$k^+ = \arg \max_{l \in [N]} \frac{\mathbf{f}_k^t \top \mathbf{f}_l^s}{\|\mathbf{f}_k^t\| \|\mathbf{f}_l^s\|} \quad (4)$$

Our patch-level knowledge distillation loss now becomes:

$$\mathcal{L}_{[\text{patch}]} = \sum_{k=1}^N \omega_{k^+} \cdot \mathcal{H}(\mathbf{P}_t^{[\text{patch}]}(\mathbf{y}')_k, \mathbf{P}_s^{[\text{patch}]}(\hat{\mathbf{x}}')_{k^+}) \quad (5)$$

where ω_{k^+} is a scalar representing the weight we give to each loss term. We find that taking it as a constant value is effective enough. The ablation study of more complex designs of ω_{k^+} is given in the appendix.

Our loss shares some similarities with DenseCL [73]. However, the differences are also obvious: (1) Our loss serves as an extension of their self-supervised contrastive loss into a supervised variant. (2) We further incorporate MIM into our design and allow masked patches to be matched, which makes our task harder and leads to more semantically meaningful patch embeddings.

3.4. Training Pipeline

We train our model in two stages: self-supervised pre-training and supervised training. In the first stage, we use the recently proposed MIM framework [88] for self-supervised pre-training. The self-supervised loss is the summation of $\mathcal{L}_{[\text{cls}]}$ and \mathcal{L}_{MIM} in Eq.(1) and Eq.(2) without scaling. In the second stage, we continue training the pre-trained model using our supervised-contrastive losses $\mathcal{L}_{[\text{cls}]}$ and $\mathcal{L}_{[\text{patch}]}$ from Eq.(3) and Eq.(5). We define our training loss as: $\mathcal{L} = \mathcal{L}_{[\text{cls}]} + \lambda \mathcal{L}_{[\text{patch}]}$, where λ controls the relative scale of these two components. A relatively large λ will make our model focus more on localization and less on high-level semantics.

4. Experiment

4.1. Datasets

Our model is evaluated on four widely used and publicly available few-shot classification datasets: *mini*-ImageNet [70], *tiered*-ImageNet [58], CIFAR-FS [4], and FC100

[51]. *mini*-ImageNet and *tiered*-ImageNet are derived from ImageNet [14], while CIFAR-FS and FC100 are derived from CIFAR100 [40]. *mini*-ImageNet contains 100 classes, which are randomly split into 64 base classes for training, 16 classes for few-shot validation, and 20 classes for few-shot evaluation. There are 600 images in each class. *tiered*-ImageNet contains 609 classes with 779165 images in total. The class split for training, few-shot validation and few-shot evaluation are 351, 97, and 160. CIFAR-FS contains 100 classes with class split as 64, 16, and 20. FC100 contains 100 classes with class split as 60, 20, and 20. For both of these two datasets, each class has 600 images with smaller resolutions (32×32) compared with ImageNet.

4.2. Implementation Details

Self-supervised pretraining. We pre-train our Vision Transformer backbone and projection head following the same pipeline in iBOT [88]. Most of the hyper-parameter settings are kept unchanged without tuning. We use a batch size of 640 and a learning rate of 0.0005 decayed with the cosine schedule. *mini*-ImageNet and *tiered*-ImageNet are pre-trained for 1200 epochs, and CIFAR-FS and FC100 are pre-trained for 900 epochs. All models are trained on 8 Nvidia RTX 3090 GPUs. Detailed training parameter settings are included in the appendix.

Supervised knowledge distillation After getting the pre-trained model, we train it with our supervised-contrastive losses. We find that the model with the best performance can converge within 60 epochs in the validation set. We use the same batch size and learning rate as in the pretraining stage. Compared with the first stage, the only extra hyper-parameter is the scaling parameter λ . We set it to make the ratio between $\mathcal{L}_{[\text{patch}]}$ and $\mathcal{L}_{[\text{cls}]}$ roughly around 2. Ablation over different λ is given in the appendix.

Few-shot Evaluation We use the simple prototype classification method (*Prototype*) [51,62] and the linear classifier (*Classifier*) used in S2M2 [50] as our default methods for few-shot evaluation. For each sampled episode data in an N-way K-shot task, *Prototype* first estimates each class prototype using the averaged feature over K support samples. Then a new sampled query image is classified into one of the N classes which has the highest cosine similarity between its feature vector and the class prototype. Instead, *Classifier* trains a linear classifier from the $N \times K$ support samples, which is then used to classify new query samples. More complex evaluation methods (e.g. DeepEMD) are also compatible with our framework. The feature we used for evaluation is the concatenation of the [cls] token with the weighted average [patch] token (*weighted avg pool*). The weights of *weighted avg pool* is the average of the self-attention values of the [cls] token with all heads of the last attention layer. Ablation over different choices of features for evaluation is studied in Sec. 4.4.

Table 1. **Results on mini-ImageNet and tiered-ImageNet.** Top three methods are colored in red, blue and green respectively. A more comprehensive version of this table is shown in the appendix.

Method	Backbone	#Params	miniImageNet,5-way		tieredImageNet,5-way	
			1-shot	5-shot	1-shot	5-shot
DeepEMD [83]	<i>ResNet-12</i>	12.4M	65.91 ± 0.82	82.41 ± 0.56	71.16 ± 0.87	86.03 ± 0.58
IE [59]	<i>ResNet-12</i>	12.4M	67.28 ± 0.80	84.78 ± 0.52	72.21 ± 0.90	87.08 ± 0.58
COSOC [48]	<i>ResNet-12</i>	12.4M	69.28 ± 0.49	85.16 ± 0.42	73.57 ± 0.43	87.57 ± 0.10
Meta-QDA [85]	<i>WRN-28-10</i>	36.5M	67.38 ± 0.55	84.27 ± 0.75	74.29 ± 0.66	89.41 ± 0.77
OM [54]	<i>WRN-28-10</i>	36.5M	66.78 ± 0.30	85.29 ± 0.41	71.54 ± 0.29	87.79 ± 0.46
SUN [17]	<i>ViT</i>	12.5M	67.80 ± 0.45	83.25 ± 0.30	72.99 ± 0.50	86.74 ± 0.33
FewTURE [37]	<i>Swin-Tiny</i>	29.0M	72.40 ± 0.78	86.38 ± 0.49	76.32 ± 0.87	89.96 ± 0.55
HCTransformers [36]	$3 \times ViT-S$	63.0M	74.74 ± 0.17	89.19 ± 0.13	79.67 ± 0.20	91.72 ± 0.11
Ours (Prototype)	<i>ViT-S</i>	21M	74.28 ± 0.18	88.82 ± 0.09	78.83 ± 0.20	91.02 ± 0.12
Ours (Classifier)	<i>ViT-S</i>	21M	74.10 ± 0.17	88.89 ± 0.09	78.81 ± 0.21	91.21 ± 0.11
Ours + HCT [36]	$3 \times ViT-S$	63M	75.32 ± 0.18	89.57 ± 0.09	79.74 ± 0.20	91.68 ± 0.11

Table 2. **Results on CIFAR-FS and FC100.** A more comprehensive version of this table is shown in the appendix.

Method	Backbone	#Params	CIFAR-FS,5-way		FC100,5-way	
			1-shot	5-shot	1-shot	5-shot
BML [89]	<i>ResNet-12</i>	12.4M	73.45 ± 0.47	88.04 ± 0.33	45.00 ± 0.41	63.03 ± 0.41
IE [59]	<i>ResNet-12</i>	12.4M	77.87 ± 0.85	89.74 ± 0.57	47.76 ± 0.77	65.30 ± 0.76
TPMN [75]	<i>ResNet-12</i>	12.4M	75.50 ± 0.90	87.20 ± 0.60	46.93 ± 0.71	63.26 ± 0.74
PSST [12]	<i>WRN-28-10</i>	36.5M	77.02 ± 0.38	88.45 ± 0.35	-	-
Meta-QDA [85]	<i>WRN-28-10</i>	36.5M	75.95 ± 0.59	88.72 ± 0.79	-	-
SUN [17]	<i>ViT</i>	12.5M	78.37 ± 0.46	88.84 ± 0.32	-	-
FewTURE [37]	<i>Swin-Tiny</i>	29.0M	77.76 ± 0.81	88.90 ± 0.59	47.68 ± 0.78	63.81 ± 0.75
HCTransformers [36]	$3 \times ViT-S$	63.0M	78.89 ± 0.18	90.50 ± 0.09	48.27 ± 0.15	66.42 ± 0.16
Ours (Prototype)	<i>ViT-S</i>	21M	80.08 ± 0.18	90.63 ± 0.13	50.38 ± 0.16	68.37 ± 0.16
Ours (Classifier)	<i>ViT-S</i>	21M	79.82 ± 0.18	90.91 ± 0.13	50.28 ± 0.16	68.50 ± 0.16

4.3. Comparison With the State-of-the-arts (SOTAs)

We evaluate our proposed SMKD with the above-mentioned evaluation methods on four few-shot classification datasets. Our goal in this section is to prove the effectiveness of our method given its simple training pipeline and evaluation procedure. This being said, the performance of our method can be further boosted by adopting strategies from contemporary works (see the last row "Ours+HCT [36]" in Table 1 and detailed results in Table 3).

Compared with traditional convolutional backbones, Transformer-based models are still underdeveloped in few-shot classification. Recent methods with Transformer backbones [17, 36, 37] differ a lot in their training pipelines and evaluation procedures. SUN [17] first pre-trains a teacher network with supervised loss, then uses it to generate patch-level supervision for the student network as a supplement to class-level supervision. HCTransformers [36] converts class-level prediction into a latent prototype learning problem, and introduces spectral tokens pooling to merge neighboring tokens with similar semantic meanings adaptively. FewTURE [37], which adopts a similar two-stage training pipeline as ours, shows the effectiveness of inner loop token importance reweighting to avoid supervision collapse.

Despite the success of these works, our SMKD still shows competitive performance given its simple design. Table 1 contains the results of *mini-ImageNet* and *tiered-ImageNet*. With *Prototype* and *Classifier* as the evaluation method, our proposed SMKD outperforms all methods with *ResNet* and *WRN* backbones, and ranked second among methods with Transformer backbones. By adopting the same strategies (patch size of 8 and spectral tokens pooling) as in HCTransformers [36], our method achieves a new SOTA on *mini-ImageNet*. The effect of each of these two strategies is shown in Table 3. Furthermore, as displayed in Table 2, our method performs the best on the two small-resolution datasets (CIFAR-FS and FC100), and outperforms all previous results by great margins (**0.93%** for 1-shot and **0.41%** for 5-shot on CIFAR-FS, and **over 2%** for both 1 and 5-shot on FC100). *In conclusion, our method grows more effective on datasets with smaller resolutions and fewer training images.*

The top competitor to our model is HCTransformer [36]. However, their method is more complex compared with ours in the following two aspects: (1) Their method uses a smaller patch size of 8 rather than 16 (which is used in most other methods as well as ours). As shown in Table

Table 3. **Results for small patch size and spectral tokens pooling.** Results are evaluated on *mini-ImageNet* with *Prototype* and *Classifier* methods, and the best is reported.

Method	Patch 8	Spec. Pool	1-shot	5-shot
HCT [36]	✓		71.27 ± 0.17	84.68 ± 0.10
HCT [36]	✓	✓	74.62 ± 0.20	89.19 ± 0.13
Ours			74.28 ± 0.18	88.82 ± 0.09
Ours	✓		74.00 ± 0.17	89.57 ± 0.09
Ours	✓	✓	75.32 ± 0.18	89.23 ± 0.18

4, this makes their training procedure $4\times$ slower than ours in each epoch. (2) Their method contains three sets of cascaded transformers, each corresponding to a set of teacher-student network with ViT-S backbone. This triples the total number of learnable parameters (63M versus 21M in ours) and incurs much longer inference time.

Table 4. **Training clock time comparison on *mini-ImageNet*.** For HCT, we report the result provided in their paper. All experiments are conducted on 8 Nvidia RTX 3090 GPUs.

	Stage1	Stage2	Stage3
HCT [36]	21.1h (400 epoch)	0.58h (2 epoch)	0.21h (2 epoch)
	Self-Supervised Pre-train	Supervised Training	
Ours	15.0h (1200 epoch)	1.08h (60 epoch)	

For a fair comparison, we compare few-shot classification results of methods with the same ViT-S backbone (21M parameters) on *mini-ImageNet* in Table 5. The results for HCT [36] and FewTURE [37] are copy-pasted from their paper. Our model outperforms them by a large margin when the number of trainable parameters is comparable.

Visualizations of the multi-head self-attention and patch token correspondence are shown in Fig. 4 and Fig. 5.

Table 5. **Comparison results with the same ViT-S backbone on *mini-ImageNet*.** For HCT, we report their result of the first student transformer (training time shown in Stage 1 of Table 4).

Method	Backbone	1-shot	5-shot
HCT [36]	ViT-S	71.27 ± 0.17	84.68 ± 0.10
FewTURE [37]	ViT-S	68.02 ± 0.88	84.51 ± 0.53
Ours	ViT-S	74.28 ± 0.18	88.82 ± 0.09

4.4. Ablation Study

Self-Supervised MIM Pretraining Helps. Our training pipeline contains two stages. A natural question is whether the first self-supervised pretraining stage is necessary. We can see from Table 6 that the classification accuracy of our method drops a lot without pretraining. If we train in a single stage, with the combination of the self-supervised learning losses and supervised contrastive learning losses, the model is still hard to converge without a good initialization (1-shot: 43.71 vs. our 74.28. 5-shot: 60.27 vs. our

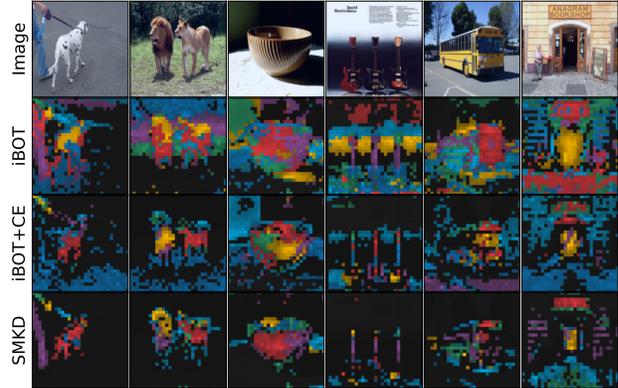


Figure 4. **Visualization of multi-head self-attention maps.** The self-attention of the $[cls]$ tokens with different heads in the last attention layer of ViT are visualized in different colors. iBOT+CE represents the model first pre-trained with iBOT, then trained with CE loss. Our SMKD pays more attention to the foreground objects, especially the most discriminative parts.

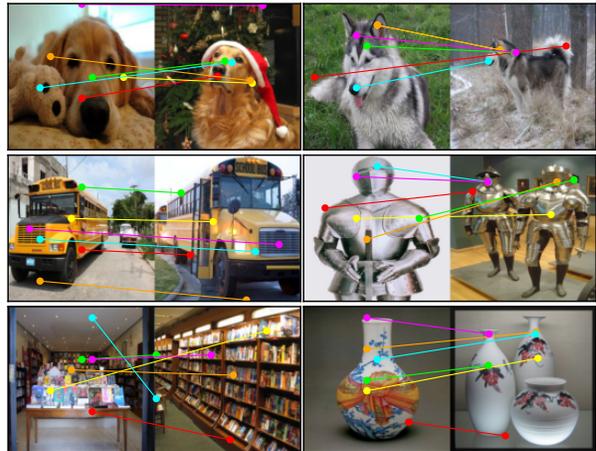


Figure 5. **Visualization of dense correspondence.** We use the patches with the highest self-attention of the $[cls]$ token on each attention head (6 in total) of the last layer of ViT-S as queries. Best matched patches with highest similarities are connected with lines.

88.82 on *mini-ImageNet*). An explanation is that it will be hard for our patch-level knowledge distillation to find matched patches if the ViT backbone is initialized with random weights. What we observe is that our model finds some shortcut solution that prohibits it from proper training. Additionally, we also test the traditional CE loss with and without self-supervised pretraining. The model pre-trained with self-supervision and then trained with supervised CE loss works surprisingly well and even outperforms most of the methods in Table 1 on *mini-ImageNet*. This further supports the usefulness of our two-stage training procedure. For the masking strategy in the first pretraining stage, our ablation study over three masking types (blockwise, random, and no mask) in the appendix shows that blockwise masking balances the best between 1 and 5-shots classification accuracy.

Table 6. **Few-shot evaluation w/wo self-supervised pretraining.** For models with pretraining stage, we pre-train with iBOT for 1200 epochs, then train with CE or our loss for 60 epochs. For models without pretraining stage, we train with supervised losses for 60 epochs directly. Results are evaluated by *Prototype* classification method on *mini-ImageNet*.

Method	Pre-train	1-shot	5-shot
CE		48.98 ± 0.16	67.62 ± 0.14
CE	✓	71.02 ± 0.18	87.25 ± 0.10
Ours		33.14 ± 0.13	45.02 ± 0.14
Ours	✓	74.28 ± 0.18	88.82 ± 0.09

Table 7. **Ablation of different losses in the supervised training stage.** All losses are trained for 60 epochs initialized with the pre-trained iBOT model. The first row represents the pre-trained iBOT without additional supervised training. $\mathcal{L}_{CE} + \mathcal{L}_{[patch]}$ stands for the combination of CE loss and our patch-level loss. Results are evaluated by *Prototype* method on *mini-ImageNet*.

Loss	1-shot	5-shot
-	60.93 ± 0.17	80.38 ± 0.12
\mathcal{L}_{CE}	71.02 ± 0.18	87.25 ± 0.10
$\mathcal{L}_{[cls]}$	70.21 ± 0.17	87.03 ± 0.10
$\mathcal{L}_{[patch]}$	70.84 ± 0.18	85.90 ± 0.11
$\mathcal{L}_{CE} + \mathcal{L}_{[patch]}$	70.70 ± 0.18	86.77 ± 0.10
$\mathcal{L}_{[cls]} + \mathcal{L}_{[patch]}$	74.28 ± 0.18	88.82 ± 0.09

Class and Patch-Level Distillation Complement Each Other. Our SMKD uses the combination of class-level and patch-level supervised-contrastive losses as the final objective. Here we study several different combinations of supervised losses to confirm our proposed method performs the best. Results in Table 7 show that combining $\mathcal{L}_{[cls]}$ with $\mathcal{L}_{[patch]}$ achieves **13.35%** improvement over baseline for 1-shot and **8.44%** for 5-shot classification, which outperforms each of its components. Sharing the parameters of projection heads for class and patch tokens allows the semantics obtained from their distillation to complement each other, which leads to much better performance. This can also be seen from the comparison between $\mathcal{L}_{CE} + \mathcal{L}_{[patch]}$ and $\mathcal{L}_{[cls]} + \mathcal{L}_{[patch]}$. \mathcal{L}_{CE} alone achieves stronger performance than $\mathcal{L}_{[cls]}$, but combining it with $\mathcal{L}_{[patch]}$ does not provide additional benefits with separate classification heads.

Weighted Average Pooling Boosts Performance. Some recent works [68, 77] find that using average pooling instead of the $[cls]$ token during evaluation can encourage token-level tasks (e.g. localization, segmentation). Inspired by such findings, we evaluate our model with different tokens as well as their combinations. Results from Table 8 show that the performance of the $[cls]$ token can be boosted by **3.02%** for 1-shot and **0.79%** for 5-shot by just concatenating with the weighted average pooling token. This shows

that our proposed method can learn meaningful representations for both $[cls]$ and $[patch]$ tokens. Moreover, weighted average pooling and average pooling tokens share similar information, but the former performs better because it gives less weight to the background.

Table 8. **Ablation of different tokens for few-shot evaluation.** "cls" stands for the $[cls]$ token, "avg pool" stands for using average pooling as token, "weighted avg pool" stands for weighted average $[patch]$ token with weights from the self-attention module from the last block of ViT. The last four rows represent concatenation of these three tokens. All tokens are normalized to unit length. Results are evaluated by *Prototype* classification method.

Embedding	1-shot	5-shot
① cls	71.26 ± 0.17	88.03 ± 0.10
② avg pool	71.65 ± 0.19	86.40 ± 0.11
③ weighted avg pool	71.83 ± 0.19	86.38 ± 0.11
① + ②	73.83 ± 0.18	88.58 ± 0.10
① + ③	74.28 ± 0.18	88.82 ± 0.09
② + ③	71.86 ± 0.19	86.47 ± 0.11
① + ② + ③	73.80 ± 0.18	88.26 ± 0.10

5. Conclusion

In this work, we propose a novel supervised knowledge distillation framework (SMKD) for few-shot Transformers, which extends the self-supervised masked knowledge distillation framework into the supervised setting. With our design of supervised-contrastive losses, we incorporate supervision into both class and patch-level knowledge distillation while still enjoy the benefits of not needing large batch size and negative samples. Evaluation results together with ablation studies demonstrate the superiority of our method given its simple design compared with contemporary works. Our two-stage training is a special case of curriculum learning from the easy samples to the hard ones. We can unify the learning objectives of self-supervised learning and supervised contrastive learning, using a carefully designed curriculum learning strategy for future works. We hope our work can bridge the gap between self-supervised and supervised knowledge distillation, and inspire more works on supervised few-shot learning methods.

6. Acknowledgements

This research is sponsored by Air Force Research Laboratory (AFRL) under agreement number FA8750-19-1-1000. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation therein. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of Air Force Laboratory, DARPA or the U.S. Government.

References

- [1] Arman Afrasiyabi, Jean-François Lalonde, and Christian Gagné. Mixture-based feature space learning for few-shot image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9041–9051, 2021. [14](#)
- [2] Sungyong Baik, Janghoon Choi, Heewon Kim, Dohee Cho, Jaesik Min, and Kyoung Mu Lee. Meta-learning with task-adaptive loss function for few-shot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9465–9474, 2021. [14](#)
- [3] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. [3](#), [4](#)
- [4] Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. *arXiv preprint arXiv:1805.08136*, 2018. [5](#)
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. [1](#)
- [6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. [13](#)
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. [2](#), [3](#), [4](#), [15](#)
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. [4](#)
- [9] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019. [3](#)
- [10] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. [4](#)
- [11] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9640–9649, 2021. [2](#)
- [12] Zhengyu Chen, Jixie Ge, Heshen Zhan, Siteng Huang, and Donglin Wang. Pareto self-supervised training for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13663–13672, 2021. [6](#), [14](#)
- [13] Zhigang Dai, Bolun Cai, Yugeng Lin, and Junying Chen. Up-detr: Unsupervised pre-training for object detection with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1601–1610, 2021. [1](#)
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [5](#)
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [3](#)
- [16] Carl Doersch, Ankush Gupta, and Andrew Zisserman. Crosstransformers: spatially-aware few-shot transfer. *Advances in Neural Information Processing Systems*, 33:21981–21993, 2020. [3](#)
- [17] Bowen Dong, Pan Zhou, Shuicheng Yan, and Wangmeng Zuo. Self-promoted supervision for few-shot transformer. In *European Conference on Computer Vision (ECCV)*, 2022. [1](#), [2](#), [6](#), [13](#), [14](#)
- [18] Xiaoyi Dong, Yinglin Zheng, Jianmin Bao, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, et al. Maskclip: Masked self-distillation advances contrastive language-image pretraining. *arXiv preprint arXiv:2208.12262*, 2022. [2](#)
- [19] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [1](#), [3](#)
- [20] Nanyi Fei, Zhiwu Lu, Tao Xiang, and Songfang Huang. Melr: Meta-learning via modeling episode-level relationships for few-shot learning. In *International Conference on Learning Representations*, 2020. [13](#)
- [21] Yutong Feng, Jianwen Jiang, Mingqian Tang, Rong Jin, and Yue Gao. Rethinking supervised pre-training for better downstream transferring. *arXiv preprint arXiv:2110.06014*, 2021. [3](#)
- [22] Hanan Gani, Muzammal Naseer, and Mohammad Yaqub. How to train vision transformer on small-scale datasets? *arXiv preprint arXiv:2210.07240*, 2022. [2](#), [3](#)
- [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/CVF international conference on computer vision*, pages 8059–8068, 2019. [1](#), [14](#)
- [24] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. [4](#)
- [25] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *European conference on computer vision*, pages 124–141. Springer, 2020. [3](#)
- [26] Guangxing Han, Yicheng He, Shiyuan Huang, Jiawei Ma, and Shih-Fu Chang. Query adaptive few-shot object detection with heterogeneous graph convolutional networks. In

- Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3263–3272, October 2021. [3](#)
- [27] Guangxing Han, Shiyuan Huang, Jiawei Ma, Yicheng He, and Shih-Fu Chang. Meta faster r-cnn: Towards accurate few-shot object detection with attentive feature alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 780–789, 2022. [3](#)
- [28] Guangxing Han, Jiawei Ma, Shiyuan Huang, Long Chen, and Shih-Fu Chang. Few-shot object detection with fully cross-transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5321–5330, June 2022. [3](#)
- [29] Guangxing Han, Jiawei Ma, Shiyuan Huang, Long Chen, Rama Chellappa, and Shih-Fu Chang. Multimodal few-shot object detection with meta-learning based cross-modal prompting. *arXiv preprint arXiv:2204.07841*, 2022. [3](#)
- [30] Guangxing Han, Xuan Zhang, and Chongrong Li. Revisiting faster r-cnn: A deeper look at region proposal network. In *International Conference on Neural Information Processing*, pages 14–24, 2017. [1](#)
- [31] Guangxing Han, Xuan Zhang, and Chongrong Li. Single shot object detection with top-down refinement. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages 3360–3364. IEEE, 2017. [1](#)
- [32] Guangxing Han, Xuan Zhang, and Chongrong Li. Semi-supervised dff: Decoupling detection and feature flow for video object detectors. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1811–1819, 2018. [1](#)
- [33] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. [2, 3](#)
- [34] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. [4](#)
- [35] 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](#)
- [36] Yangji He, Weihang Liang, Dongyang Zhao, Hong-Yu Zhou, Weifeng Ge, Yizhou Yu, and Wenqiang Zhang. Attribute surrogates learning and spectral tokens pooling in transformers for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9119–9129, 2022. [1, 2, 3, 6, 7, 13, 14](#)
- [37] Markus Hiller, Rongkai Ma, Mehrtash Harandi, and Tom Drummond. Rethinking generalization in few-shot classification. *arXiv preprint arXiv:2206.07267*, 2022. [2, 3, 6, 7, 13, 14](#)
- [38] Dahyun Kang, Heeseung Kwon, Juhong Min, and Minsu Cho. Relational embedding for few-shot classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8822–8833, 2021. [14](#)
- [39] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. [2](#)
- [40] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [5](#)
- [41] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. [1](#)
- [42] Seung Hoon Lee, Seunghyun Lee, and Byung Cheol Song. Vision transformer for small-size datasets. *arXiv preprint arXiv:2112.13492*, 2021. [1](#)
- [43] Junjie Li, Zilei Wang, and Xiaoming Hu. Learning intact features by erasing-inpainting for few-shot classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8401–8409, 2021. [13](#)
- [44] Yahui Liu, Enver Sangineto, Wei Bi, Nicu Sebe, Bruno Lepri, and Marco Nadai. Efficient training of visual transformers with small datasets. *Advances in Neural Information Processing Systems*, 34:23818–23830, 2021. [1, 2, 3](#)
- [45] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. [1, 3](#)
- [46] Jiang Lu, Pinghua Gong, Jieping Ye, and Changshui Zhang. Learning from very few samples: A survey. *arXiv preprint arXiv:2009.02653*, 2020. [1](#)
- [47] Yuning Lu, Liangjian Wen, Jianzhuang Liu, Yajing Liu, and Xinmei Tian. Self-supervision can be a good few-shot learner. *arXiv preprint arXiv:2207.09176*, 2022. [2](#)
- [48] Xu Luo, Longhui Wei, Liangjian Wen, Jinrong Yang, Lingxi Xie, Zenglin Xu, and Qi Tian. Rectifying the shortcut learning of background: Shared object concentration for few-shot image recognition. *arXiv preprint arXiv:2107.07746*, 2021. [6, 13](#)
- [49] Jiawei Ma, Hanchen Xie, Guangxing Han, Shih-Fu Chang, Aram Galstyan, and Wael Abd-Almageed. Partner-assisted learning for few-shot image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10573–10582, 2021. [1, 3, 4, 13, 14](#)
- [50] Puneet Mangla, Nupur Kumari, Abhishek Sinha, Mayank Singh, Balaji Krishnamurthy, and Vineeth N Balasubramanian. Charting the right manifold: Manifold mixup for few-shot learning. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2218–2227, 2020. [2, 3, 5](#)
- [51] Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive metric for improved few-shot learning. *Advances in neural information processing systems*, 31, 2018. [3, 5](#)
- [52] Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. A unified view of masked image modeling. *arXiv preprint arXiv:2210.10615*, 2022. [2, 3](#)
- [53] Elia Peruzzo, Enver Sangineto, Yahui Liu, Marco De Nadai, Wei Bi, Bruno Lepri, and Nicu Sebe. Spatial entropy

- regularization for vision transformers. *arXiv preprint arXiv:2206.04636*, 2022. 2, 3
- [54] Guodong Qi, Huimin Yu, Zhaohui Lu, and Shuzhao Li. Transductive few-shot classification on the oblique manifold. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8412–8422, 2021. 6, 13
- [55] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3
- [56] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 3
- [57] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12179–12188, 2021. 1
- [58] Mengye Ren, Eleni Triantafyllou, 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
- [59] Mamshad Nayeem Rizve, Salman Khan, Fahad Shahbaz Khan, and Mubarak Shah. Exploring complementary strengths of invariant and equivariant representations for few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10836–10846, June 2021. 6, 13, 14
- [60] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. *arXiv preprint arXiv:1807.05960*, 2018. 3
- [61] 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. 14
- [62] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 3, 5
- [63] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021. 1
- [64] 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. 3
- [65] Hao Tan, Jie Lei, Thomas Wolf, and Mohit Bansal. Vimpac: Video pre-training via masked token prediction and contrastive learning. *arXiv preprint arXiv:2106.11250*, 2021. 4
- [66] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? In *European Conference on Computer Vision*, pages 266–282. Springer, 2020. 3
- [67] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 1, 3
- [68] Hugo Touvron, Matthieu Cord, and Hervé Jégou. Deit iii: Revenge of the vit. *arXiv preprint arXiv:2204.07118*, 2022. 8
- [69] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [70] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016. 1, 3, 5
- [71] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 3
- [72] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. 3
- [73] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3024–3033, 2021. 5
- [74] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22–31, 2021. 1, 3
- [75] Jiamin Wu, Tianzhu Zhang, Yongdong Zhang, and Feng Wu. Task-aware part mining network for few-shot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8433–8442, 2021. 6, 13, 14
- [76] Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for motor control. *arXiv preprint arXiv:2203.06173*, 2022. 2
- [77] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084*, 2021. 8
- [78] Weijian Xu, Yifan Xu, Huaijin Wang, and Zhuowen Tu. Attentional constellation nets for few-shot learning. In *International Conference on Learning Representations*, 2021. 14
- [79] Shuo Yang, Lu Liu, and Min Xu. Free lunch for few-shot learning: Distribution calibration. *arXiv preprint arXiv:2101.06395*, 2021. 13

- [80] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8808–8817, 2020. [3](#), [13](#)
- [81] Nikolaos-Antonios Ypsilantis, Noa Garcia, Guangxing Han, Sarah Ibrahimi, Nanne Van Noord, and Giorgos Tolias. The met dataset: Instance-level recognition for artworks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. [3](#)
- [82] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International Conference on Machine Learning*, pages 12310–12320. PMLR, 2021. [4](#)
- [83] Chi Zhang, Yujun Cai, Guosheng Lin, and Chunhua Shen. Deepemd: Few-shot image classification with differentiable earth mover’s distance and structured classifiers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [2](#), [6](#), [13](#)
- [84] Chi Zhang, Henghui Ding, Guosheng Lin, RuiBo Li, Changhu Wang, and Chunhua Shen. Meta navigator: Search for a good adaptation policy for few-shot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9435–9444, 2021. [13](#), [14](#)
- [85] Xueting Zhang, Debin Meng, Henry Gouk, and Timothy M Hospedales. Shallow bayesian meta learning for real-world few-shot recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 651–660, 2021. [6](#), [13](#), [14](#)
- [86] Jiabao Zhao, Yifan Yang, Xin Lin, Jing Yang, and Liang He. Looking wider for better adaptive representation in few-shot learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 10981–10989, 2021. [13](#)
- [87] Nanxuan Zhao, Zhirong Wu, Rynson WH Lau, and Stephen Lin. What makes instance discrimination good for transfer learning? *arXiv preprint arXiv:2006.06606*, 2020. [4](#)
- [88] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021. [2](#), [3](#), [4](#), [5](#), [13](#)
- [89] Ziqi Zhou, Xi Qiu, Jiangtao Xie, Jianan Wu, and Chi Zhang. Binocular mutual learning for improving few-shot classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8402–8411, 2021. [6](#), [13](#), [14](#)
- [90] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. [1](#)

7. Appendix

7.1. Implementation Details for Model Training

Self-supervised pretraining. We pre-train our Vision Transformer backbone and projection head following the same pipeline in iBOT [88]. Most of the hyper-parameter settings are kept unchanged without tuning. ViT-Small, which has ~ 21 M parameters is used as our default architecture. Our default patch size is set as 16. For the student network, the [cls] token output and [patch] tokens output share the same projection head. This head-sharing strategy is also applied to the teacher network. For both networks, we set the output dimension of projection heads as 8192. We linearly warm up the learning rate for 10 epochs to its base value of $5e-4$, then use cosine schedule to decay it to $1e-5$. Cosine schedule is also used for weight decay from 0.04 to 0.4. Besides, we use the multi-crop strategy [6] with 2 global crops (224×224) and 10 local crops (96×96), with scale range (0.4, 1.0) and (0.05, 0.4) respectively. We found that allowing knowledge distillation between global and local crops from intra-class images harms the performance, which is consistent with [88]. Therefore, local crops here are only used for self-distillation with global crops from the same image. Furthermore, we apply blockwise masking on global crops sent into the student network, with a masking ratio uniformly sampled from [0, 1, 0.5] with probability 0.5, and 0 with probability 0.5. Ablation of different masking strategies is given in Sec. 7.3. Our batch size is set as 640 (batch size per GPU equal to 80). *mini-ImageNet* and *tiered-ImageNet* are pre-trained for 1200 epochs, and CIFAR-FS and FC100 are pre-trained for 900 epochs. All models are trained on 8 Nvidia RTX 3090 GPUs.

Supervised knowledge distillation. After finishing the pretraining stage, we train the model with our supervised-contrastive loss. The best evaluation accuracy on the validation set can usually be achieved within 60 epochs of training. We use the same set of hyper-parameters as the first pretraining stage without further tuning. Ablation of the scaling parameter λ , which controls the relative size of $\mathcal{L}_{[patch]}$ and $\mathcal{L}_{[cls]}$ is given in Sec. 7.3.

7.2. Few-shot Evaluation Results

We present few-shot evaluation results with more methods on the four benchmark datasets here in Table 9 and 10. ViT-based methods are better than the traditional CNN-based methods in general. The ranking of our method remains unchanged.

Table 9. **More comprehensive few-shot evaluation results on mini-ImageNet and tiered-ImageNet.** Top three methods are colored in red, blue and green respectively.

Method	Backbone	#Params	miniImageNet,5-way		tieredImageNet,5-way	
			1-shot	5-shot	1-shot	5-shot
DeepEMD [83]	<i>ResNet-12</i>	12.4M	65.91 \pm 0.82	82.41 \pm 0.56	71.16 \pm 0.87	86.03 \pm 0.58
IE [59]	<i>ResNet-12</i>	12.4M	67.28 \pm 0.80	84.78 \pm 0.52	72.21 \pm 0.90	87.08 \pm 0.58
BML [89]	<i>ResNet-12</i>	12.4M	67.04 \pm 0.63	83.63 \pm 0.29	68.99 \pm 0.50	85.49 \pm 0.34
PAL [49]	<i>ResNet-12</i>	12.4M	69.37 \pm 0.64	84.40 \pm 0.44	72.25 \pm 0.72	86.95 \pm 0.47
TPMN [75]	<i>ResNet-12</i>	12.4M	67.64 \pm 0.63	83.44 \pm 0.43	72.24 \pm 0.70	86.55 \pm 0.63
MN+MC [84]	<i>ResNet-12</i>	12.4M	67.14 \pm 0.80	83.82 \pm 0.51	74.58 \pm 0.88	86.73 \pm 0.61
DC [79]	<i>ResNet-12</i>	12.4M	68.57 \pm 0.55	82.88 \pm 0.42	78.19 \pm 0.25	89.90 \pm 0.41
MELR [20]	<i>ResNet-12</i>	12.4M	67.40 \pm 0.43	83.40 \pm 0.28	72.14 \pm 0.51	87.01 \pm 0.35
COSOC [48]	<i>ResNet-12</i>	12.4M	69.28 \pm 0.49	85.16 \pm 0.42	73.57 \pm 0.43	87.57 \pm 0.10
CSEI [43]	<i>ResNet-12</i>	12.4M	68.94 \pm 0.28	85.07 \pm 0.50	73.76 \pm 0.32	87.83 \pm 0.59
CNL [86]	<i>ResNet-12</i>	12.4M	67.96 \pm 0.98	83.36 \pm 0.51	73.42 \pm 0.95	87.72 \pm 0.75
FEAT [80]	<i>WRN-28-10</i>	36.5M	65.10 \pm 0.20	81.11 \pm 0.14	70.41 \pm 0.23	84.38 \pm 0.16
Meta-QDA [85]	<i>WRN-28-10</i>	36.5M	67.38 \pm 0.55	84.27 \pm 0.75	74.29 \pm 0.66	89.41 \pm 0.77
OM [54]	<i>WRN-28-10</i>	36.5M	66.78 \pm 0.30	85.29 \pm 0.41	71.54 \pm 0.29	87.79 \pm 0.46
SUN [17]	<i>ViT</i>	12.5M	67.80 \pm 0.45	83.25 \pm 0.30	72.99 \pm 0.50	86.74 \pm 0.33
FewTURE [37]	<i>ViT-S</i>	21.0M	68.02 \pm 0.88	84.51 \pm 0.53	72.96 \pm 0.92	86.43 \pm 0.67
FewTURE [37]	<i>Swin-Tiny</i>	29.0M	72.40 \pm 0.78	86.38 \pm 0.49	76.32 \pm 0.87	89.96 \pm 0.55
HCT (Prototype) [36]	$3 \times$ <i>ViT-S</i>	63.0M	74.74 \pm 0.17	85.66 \pm 0.10	79.67 \pm 0.20	89.27 \pm 0.13
HCT (Classifier) [36]	$3 \times$ <i>ViT-S</i>	63.0M	74.62 \pm 0.20	89.19 \pm 0.13	79.57 \pm 0.20	91.72 \pm 0.11
Ours (Prototype)	<i>ViT-S</i>	21.0M	74.28 \pm 0.18	88.82 \pm 0.09	78.83 \pm 0.20	91.02 \pm 0.12
Ours (Classifier)	<i>ViT-S</i>	21.0M	74.10 \pm 0.17	88.89 \pm 0.09	78.81 \pm 0.21	91.21 \pm 0.11
Ours + HCT [36]	$3 \times$ <i>ViT-S</i>	63.0M	75.32 \pm 0.18	89.57 \pm 0.09	79.74 \pm 0.20	91.68 \pm 0.11

Table 10. **More comprehensive few-shot evaluation results on CIFAR-FS and FC100.** Top three methods are colored in red, blue and green respectively.

Method	Backbone	#Params	CIFAR-FS,5-way		FC100,5-way	
			1-shot	5-shot	1-shot	5-shot
DSN-MR [61]	<i>ResNet-12</i>	12.4M	75.60 ± 0.90	86.20 ± 0.60	-	-
BML [89]	<i>ResNet-12</i>	12.4M	73.45 ± 0.47	88.04 ± 0.33	45.00 ± 0.41	63.03 ± 0.41
IE [59]	<i>ResNet-12</i>	12.4M	77.87 ± 0.85	89.74 ± 0.57	47.76 ± 0.77	65.30 ± 0.76
PAL [49]	<i>ResNet-12</i>	12.4M	77.10 ± 0.70	88.00 ± 0.50	47.20 ± 0.60	64.00 ± 0.60
TPMN [75]	<i>ResNet-12</i>	12.4M	75.50 ± 0.90	87.20 ± 0.60	46.93 ± 0.71	63.26 ± 0.74
MN+MC [84]	<i>ResNet-12</i>	12.4M	74.63 ± 0.91	86.45 ± 0.59	46.40 ± 0.81	61.33 ± 0.71
RENet [38]	<i>ResNet-12</i>	12.4M	74.51 ± 0.46	86.60 ± 0.32	-	-
ConstellationNet [78]	<i>ResNet-12</i>	12.4M	75.40 ± 0.20	86.80 ± 0.20	43.80 ± 0.20	59.70 ± 0.20
ALFA+MeTAL [2]	<i>ResNet-12</i>	12.4M	-	-	44.54 ± 0.50	58.44 ± 0.42
MixtFSL [1]	<i>ResNet-12</i>	12.4M	-	-	41.50 ± 0.67	58.39 ± 0.62
CC+rot [23]	<i>WRN-28-10</i>	36.5M	73.62 ± 0.31	86.05 ± 0.22	-	-
PSST [12]	<i>WRN-28-10</i>	36.5M	77.02 ± 0.38	88.45 ± 0.35	-	-
Meta-QDA [85]	<i>WRN-28-10</i>	36.5M	75.95 ± 0.59	88.72 ± 0.79	-	-
SUN [17]	<i>ViT</i>	12.5M	78.37 ± 0.46	88.84 ± 0.32	-	-
FewTURE [37]	<i>ViT-S</i>	21.0M	76.10 ± 0.88	86.14 ± 0.64	46.20 ± 0.79	63.14 ± 0.73
FewTURE [37]	<i>Swin-Tiny</i>	29.0M	77.76 ± 0.81	88.90 ± 0.59	47.68 ± 0.78	63.81 ± 0.75
HCT (Prototype) [36]	$3 \times$ <i>ViT-S</i>	63.0M	78.89 ± 0.18	87.73 ± 0.11	48.27 ± 0.15	61.49 ± 0.15
HCT (Classifier) [36]	$3 \times$ <i>ViT-S</i>	63.0M	78.88 ± 0.18	90.50 ± 0.09	48.15 ± 0.16	66.42 ± 0.16
Ours (Prototype)	<i>ViT-S</i>	21M	80.08 ± 0.18	90.63 ± 0.13	50.38 ± 0.16	68.37 ± 0.16
Ours (Classifier)	<i>ViT-S</i>	21M	79.82 ± 0.18	90.91 ± 0.13	50.28 ± 0.16	68.50 ± 0.16

7.3. Additional Ablation Studies

Why $\mathcal{L}_{[cls]} + \mathcal{L}_{MIM}$ in stage 1? Our insight is that the $[cls]$ tokens in global loss have better high-level semantics, but often disregard the rich local structures. While the MIM loss \mathcal{L}_{MIM} constructed from $[patch]$ tokens can remedy this problem, increase task difficulty, and work as strong data augmentations. In Table 11, we can find that using both losses in stage 1 gives the best results.

Table 11. Ablation of SSL tasks in stage 1 on *mini-ImageNet*.

Stage1				Stage2: $\mathcal{L}_{[cls]} + \mathcal{L}_{[patch]}$	
$\mathcal{L}_{[cls]}$	\mathcal{L}_{MIM}	1-shot	5-shot	1-shot	5-shot
✓		58.55	78.90	72.93	88.07
	✓	27.66	33.82	37.03	50.95
✓	✓	60.93	80.38	74.28	88.82

Masking Strategies. We use blockwise masking as our default in the main text. In Table 12, we test random mask and no mask while keeping all other hyper-parameters unchanged. "Block Mask \rightarrow No Mask" represents self-supervised pretraining with blockwise masking, and supervised training with no mask. Using either a random mask or block mask can boost the classification accuracy in the first self-supervised pretraining stage, but their advantage over no mask decreases in the second supervised training stage. We choose blockwise masking as our default strategy since it balances 1 and 5-shot classification accuracy the best.

Scaling Parameter λ . This parameter controls the relative importance of class-level and patch-level losses in our final loss: $\mathcal{L} = \mathcal{L}_{[cls]} + \lambda \mathcal{L}_{[patch]}$. A relatively large value of λ will put more focus on localization and less on high-level semantics. Here in Table 13, we test different λ values by keeping the base of $\mathcal{L}_{[cls]}$ to 1 and scale $\mathcal{L}_{[patch]}$. As we can see, the λ parameter influences 1-shot classification accuracy more than 5-shot. We choose $\lambda = 0.25$ as our default (which makes the ratio of $\mathcal{L}_{[patch]}/\mathcal{L}_{[cls]}$ roughly around 2) since it has best 1-shot performance and competitive 5-shot accuracy.

Table 12. Ablation over different masking strategy in self-supervised pretraining stage.

Masking Strategy	Self-supervised Pre-train		Supervised Training	
	1-shot	5-shot	1-shot	5-shot
No Mask	59.15 \pm 0.17	79.23 \pm 0.12	73.94 \pm 0.17	88.93 \pm 0.09
Block Mask	60.93 \pm 0.17	80.38 \pm 0.12	74.01 \pm 0.17	88.89 \pm 0.09
Random Mask	60.94 \pm 0.18	79.62 \pm 0.13	74.07 \pm 0.18	88.66 \pm 0.09
Block Mask \rightarrow No Mask	60.93 \pm 0.17	80.38 \pm 0.12	73.44 \pm 0.17	88.87 \pm 0.09

Table 13. The influence of different ratio between $\mathcal{L}_{[cls]}$ and $\mathcal{L}_{[patch]}$.

$\mathcal{L}_{[patch]}/\mathcal{L}_{[cls]}$	1-shot	5-shot
≈ 4 ($\lambda = 0.9$)	73.45 \pm 0.17	88.89 \pm 0.09
≈ 2 ($\lambda = 0.45$)	74.28 \pm 0.18	88.82 \pm 0.09
≈ 1 ($\lambda = 0.2$)	74.01 \pm 0.17	88.89 \pm 0.09
≈ 0.5 ($\lambda = 0.1$)	73.11 \pm 0.17	88.44 \pm 0.09

Weighting Parameter ω_{k+} in $\mathcal{L}_{[patch]}$. This parameter in Eq.(5) gives weights to each component of our patch-level contrastive loss $\mathcal{L}_{[patch]}$. We set $\omega_{k+} = 1/N$ in our main text due to its simplicity. In Table 14, we compare it (*Simple Avg*) with another variant (*Self-Attention Weighted Avg*), which uses the averaged self-attention weights over attention heads of the [cls] token with all [patch] tokens in the last attention layer of teacher network to aggregate pairwise patch matching losses. As found in [7], the self-attention of ViTs is good at capturing foreground regions. So we use it here as a way to highlight foreground objects and to attenuate irrelevant background information. Our default simple average outperforms this variant on both 1 and 5-shot classification accuracies. One explanation is as follows. If the foreground objects of two intra-class images differ a lot, then *Self-Attention Weighted Avg* tends to minimize $\mathcal{L}_{[patch]}$ by decreasing the weights of the losses associated with the patches covering these foreground objects, which makes our model deviate from optimum.

Table 14. Different weighting schemes of patch-level supervised-contrastive loss.

Weighting Scheme	1-shot	5-shot
Simple Avg	74.28 \pm 0.18	88.82 \pm 0.09
Self-Attention Weighted Avg	74.11 \pm 0.18	88.52 \pm 0.10

Comparison with smaller backbones: To make ViT-S (~ 21 M parameters) comparable with ResNet-12 (~ 12 M parameters), we trim by half either its embedding dimension (d_{embed}) or the number of attention heads (#heads). From Table 15, trimming #heads by half only results in little drop in accuracy, which still outperforms the best method with ResNet-12 backbones. The training speed also increases by 10% with fewer #heads. Given this result, our comparison now becomes complete: our method outperforms both shallow (ResNet-12) and deep (WRN-28-10) CNN-based backbones, as well as ViTs with the same (see Table 5) or more (see Table 3 & Table 4) parameters.

Table 15. ViT-S with similar size as ResNet-12 on *mini*-ImageNet

Backbone	d_{embed}	#heads	#param	Stage1		Stage2	
				1-shot	5-shot	1-shot	5-shot
ViT-S	192	6	11M	60.70	79.56	71.14	87.12
ViT-S	384	3	11M	62.12	81.27	72.70	87.90
ViT-S	384	6	21M	60.93	80.38	74.28	88.82
ResNet12	-	-	12M	-	-	69.37	85.16
WRN-28-10	-	-	36M	-	-	67.38	85.29

7.4. Visualizations

We visualize more self-attention maps and dense correspondence in Fig. 6 and Fig. 7.

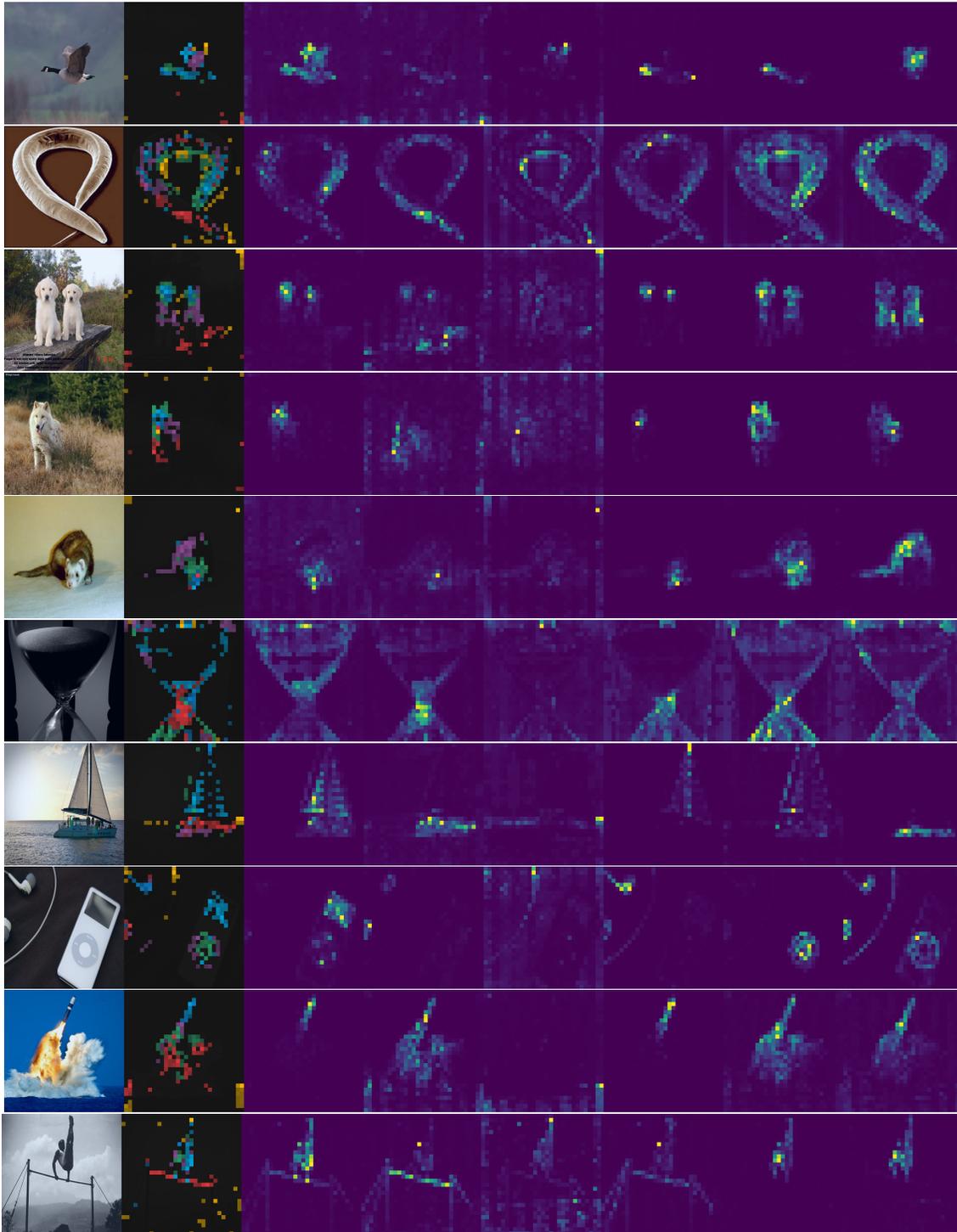


Figure 6. **Visualization of multi-head self-attention maps.** The self-attention of the `[cls]` tokens with different heads in the last attention layer of ViT are visualized in different colors in the second column. The last six columns visualize each attention head. Images are from the test set of *mini-ImageNet*.

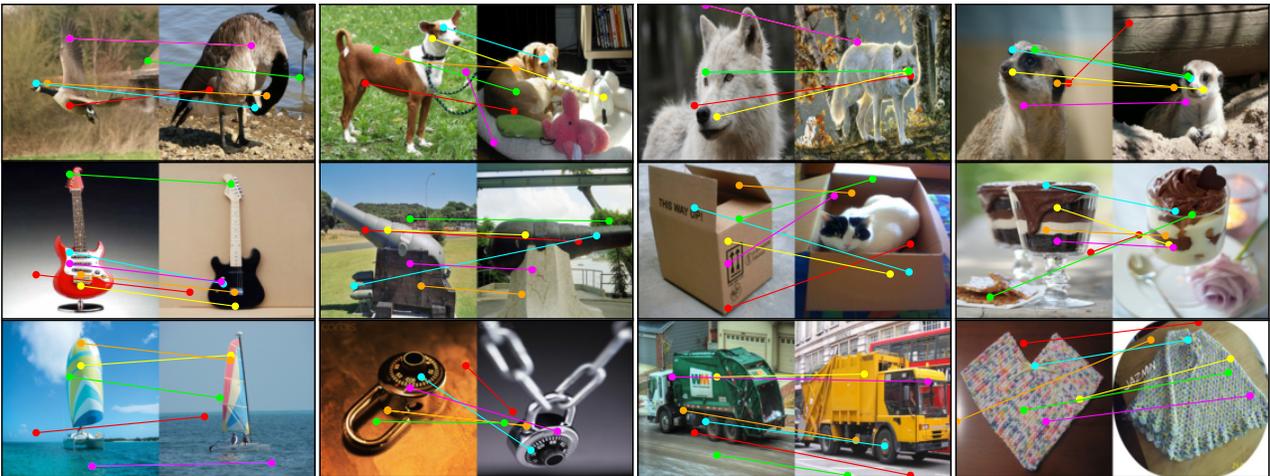


Figure 7. **Visualization of dense correspondence.** We use the patches with the highest self-attention of the $[cls]$ token on each attention head (6 in total) of the last layer of ViT-S as queries. Best-matched patches with the highest similarities are connected with lines. Images are from the validation and testing set of *mini-ImageNet*.