# Learning Instance and Task-Aware Dynamic Kernels for Few-Shot Learning

Rongkai Ma[1], Pengfei Fang[1,2,3(✉)], Gil Avraham[4], Yan Zuo[3],
Tianyu Zhu[1], Tom Drummond[5], Mehrtash Harandi[1,3]

[1]Monash University   [2]Australian National University   [3]CSIRO
[4]Amazon Australia   [5]The University of Melbourne
✉: Corresponding author

**Abstract.** Learning and generalizing to novel concepts with few samples (Few-Shot Learning) is still an essential challenge to real-world applications. A principle way of achieving few-shot learning is to realize a model that can rapidly adapt to the context of a given task. Dynamic networks have been shown capable of learning content-adaptive parameters efficiently, making them suitable for few-shot learning. In this paper, we propose to learn the dynamic kernels of a convolution network as a function of the task at hand, enabling faster generalization. To this end, we obtain our dynamic kernels based on the entire task and each sample, and develop a mechanism further conditioning on each individual channel and position independently. This results in dynamic kernels that simultaneously attend to the global information whilst also considering minuscule details available. We empirically show that our model improves performance on few-shot classification and detection tasks, achieving a tangible improvement over several baseline models. This includes state-of-the-art results on four few-shot classification benchmarks: *mini*-ImageNet, *tiered*-ImageNet, CUB and FC100 and competitive results on a few-shot detection dataset: MS COCO-PASCAL-VOC.

## 1 Introduction

Despite the great success of the modern deep neural networks (DNNs), in many cases, the problem of adapting a DNN with only a handful of labeled data is still challenging. Few-Shot Learning (FSL) aims to address the inefficiencies of modern machine learning frameworks by adapting models trained on large databases for novel tasks with limited data [15,34,35,48,62].

Early approaches of FSL learned a fixed embedding function to encode samples into a latent space, where they could be categorized by their semantic relationships [47,48,62]. However, such fixed approaches do not account for category differences [39], which may exist between already learned tasks and novel tasks. Ignoring these discrepancies can severely limit the adaptability of a model as well as its ability to scale in the FSL setting. Although methods that adapt embeddings [33,52,61] attempt to address this issue, they still utilize fixed models which lack full adaptability and are constrained to previously learned tasks. Another group of approaches, such as [15,35], adapt models with a few optimization
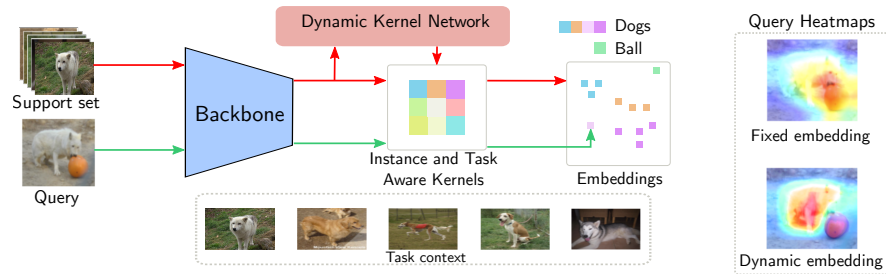
Fig. 1: Typical FSL models learn fixed embeddings, which are not flexible enough to rapidly adapt to novel tasks. Our method instead uses a dynamic kernel network to produce a set of dynamic parameters which are both instance and task-aware

steps. Given the complexity of the loss landscape of a DNN, such methods come short compared to metric-based solutions [56].

Dynamic kernel approaches have been shown to be computationally efficient [64] relative to optimization-based solutions, resulting in the models capable of encoding the novel information at the parameter level [22]. A very recent study [58] showed that dynamic kernel methods are effective when applied to FSL via the learning of input-conditioned filters, enabling the realization of adaptive models; a key limitation of this dynamic kernel method is that only the per-class level information is utilized via class prototypes.

Arguably, when learning new tasks with sparse data, it is vital to fully utilize information on both an instance and task-level for efficiency. As an example, when given the task of differentiating between dog breeds (Fig. 1), both task and instance-level information is required; it is important that we are sensitive enough to distinguish the minuscule details between different dog breeds (instance-level), yet can still focus on the global knowledge required to filter out irrelevant, non-task related objects (task-level).

In this paper, we propose a novel, dynamic approach for FSL tasks. Our method is realized by a dynamic kernel that can encode both instance-level and task-level features via jointly learning a dynamic kernel generator and a context learning module. The dynamic kernel generator adaptively produces an instance kernel by exploiting the information along the spatial and channel dimensions of a feature map in a decoupled manner. It further incorporates information from the frequency domain resulting in a rich set of descriptive representations. The context learning module refines the support features into a task-specific representation which is used to produce the task-specific kernel. The resulting instance and task-specific kernels are fused to obtain a dynamic kernel which is INStance and Task-Aware (INSTA). Our method differentiates from approaches such as FEAT [61] and GLoFa [33] by learning a set of dynamic parameters adaptive to the novel tasks instead of employing fixed models during inference. Furthermore, in contrast to optimization-based methods [15], our approach can

adapt the model parameters without the requirements of backpropagation during the inference stage. We offer the following contributions:

- We propose a novel FSL approach to extract both instance and task-specific information using dynamic kernels.
- We offer the first FSL framework capable of being evaluated on both classification *and* detection tasks.
- Empirically, we offer substantial improvements over several FSL baselines, including optimization-based and metric-based approaches.

## 2   Related Work

The family of few-shot learning literature is broad and diverse. However, those related to this work are mainly the family of optimization-based methods [4,15–17,26,35,41,46,65] and metric-based methods [5,12,45,47,48,52,61,62].

**Optimization based methods.** Optimization-based methods such as MAML [15] or Reptile [35] focus on learning a set of initial model parameters that can generalize to new tasks (or environments) with a small number of optimization steps and samples without severe over-fitting. Furthermore, in most cases, this group of methods present a framework trained with a bi-level optimization setting [17], which provides a feasible solution to adapt the model to the test set from the initialized model.

Our proposed framework is similar to the optimization-based methods [1,2,7,27,32] in the sense that the model parameters are task-adaptive. However, our solution does not require backpropagation during inference to achieve so. Furthermore, our method can be incorporated with optimization-based methods, and empirically we observed that such construction yields performance improvement. This will be demonstrated in § 4.

**Metric-based methods.** In few-shot learning literature, the metric-based methods aim to define a metric to measure the dis/similarity between samples from a few observations [9,24,30,31,42,57,59,63]. ProtoNet [47] achieves this by learning a fixed latent space where the class representations (*i.e.*, prototype), obtained by averaging the features from the same class, are distinctive for different classes. DeepEMD [62] formulates the query-support matching as an optimal transport problem and adopts Earth Mover's distance as the metric. One commonality in the aforementioned methods is the fact that all employ a fixed embedding space when facing novel tasks, which essentially limits their adaptability. On the other hand, many previous methods suggested to adapt the embeddings to the novel tasks [14,19,33,36,44,50,61]. CTM [28] proposes to produce a mask to disregard the uninformative features from the support and query embeddings during inference. MatchingNet [52] uses a memory module to perform sample-wise embedding adaptation and determines the query label by a cosine distance. TADAM [36] proposes to learn a dynamic feature extractor by applying a linear transformation to the embeddings to learn a set of scale and shift parameters. FEAT [61] and GLoFa [33] provide an inspiring way to perform the embedding adaptation using a set function.

Our contribution is complementary to the aforementioned methods. We aim to learn a set of dynamic kernels via exploiting the instance and task-level information according to the task at hand. This results in a more distinctive and descriptive representation, which effectively boosts the performance of these methods directly relying on constructing a metric space.

**Dynamic kernels.** The application of dynamic kernels solutions within the domain of few-shot learning is less explored in the current literature. However, it has been demonstrated useful when labels are abundant [6,8,18,21,49,51,54,64]. Zhou *et al.* [64] have proposed to use decoupled attention over the channel and spatial dimensions. This results in a content-adaptive and efficient method that provides a feasible way to achieve task-adaptive FSL.

To leverage the effectiveness of the dynamic kernel into few-shot learning tasks, [58] propose to learn dynamic filters for the spatial and channels separately via grouping strategy. The resulted kernels are then applied to the query to produce a support-aligned feature. Due to the usage of grouping, the performance might sacrifice for efficiency [64]. Inspired by these methods, our INSTA produces dynamic kernels that are both instance-aware and task (or episodic)-aware while also incorporating valuable frequency patterns; as such, our method produces more informative and representative features.

## 3   Method

In this section, we introduce our proposed INSTA. Note that while we present our method in terms of few-shot classification, the proposed approach is generic and can be seamlessly used to address other few-shot problems, including structured-prediction tasks such as few-shot object detection (see § 4.2 for details).

### 3.1   Problem Formulation

In what follows, we will give a brief description of the problem formulation for few-shot classification. The vectors and matrices (or high-dimensional tensors) are denoted by bold lower-case letters (*e.g.*, $\boldsymbol{x}$) and bold upper-case letters (*e.g.*, $\boldsymbol{X}$) throughout this paper. FSL aims to generalize the knowledge acquired by a model from a set of examples $\mathcal{D}_{train} = \{(\boldsymbol{X}_i, y_i) | y_i \in \mathcal{C}_{train}\}$, to novel and unseen tasks $\mathcal{D}_{test} = \{(\boldsymbol{X}_i, y_i) | y_i \in \mathcal{C}_{test}\}, \mathcal{C}_{train} \cap \mathcal{C}_{test} = \emptyset$, in a low data regime.

We follow the meta-learning protocol to formulate FSL with episodic training and testing. Specifically, an episode $\mathcal{E}$ consists of a support set $\mathcal{X}^s = \{(\boldsymbol{X}_{ij}^s, y_i^s) | i = 1, \ldots, N, j = 1, \ldots, K, y_i^s \in \mathcal{C}_{train}\}$, where $\boldsymbol{X}_{ij}^s$ denotes the $j$-th sample in the class $y_i^s$ and the query set $\mathcal{X}^q = \{(\boldsymbol{X}_i^q, y_i^q) | i = 1, \ldots, N\}$, where $\boldsymbol{X}_i^q$ denotes a query example [1] sampled from class $y_i^q$ (the test setup is the same as training but the episodes are sampled from $\mathcal{D}_{test}$). Such a formulation is known as $N$-way $K$-shot, where the goal is to utilize the support samples and their labels to obtain $\Theta^*$, the optimal parameters of the model, such that each query is classified

---

[1] Without losing generality, we use one sample per class as a query for presenting our method. In practice, each episode contains multiple samples per query class.
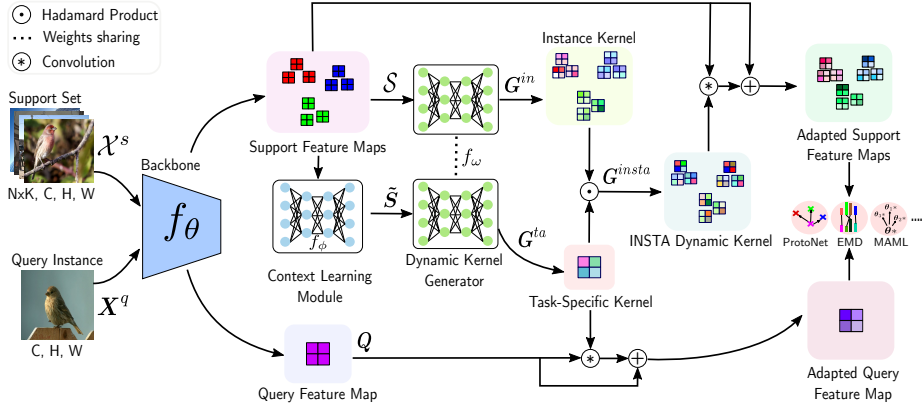
Fig. 2: The framework of our method. Given the support set $\mathcal{X}^s$ and the query sample $\boldsymbol{X}^q$, the backbone network $f_\theta$ first encodes them into a representation space as $\mathcal{S} = \{\boldsymbol{S}_{11}, \ldots, \boldsymbol{S}_{NK}\}$ and $\boldsymbol{Q}$. Then a dynamic generator $f_\omega$ is used to produce an instance kernel for each support sample. Meanwhile, a context learning module $f_\phi$ is used to refine and aggregate the context of the entire support feature maps to produce a task-specific representation, which is then used as the input of $f_\omega$ to obtain the task-specific kernel. Note the instance kernel and task-specific kernel share the parameters across $f_\omega$. Finally, the INSTA dynamic kernel that is both instance-aware and task-aware is applied to the support feature maps, while only the task-specific kernel is applied to the query feature map to obtain the adapted representations

correctly into one of classes in the support set. The object of network training follows that:

$$\Theta^* = \arg\min_{\Theta} \sum_{\boldsymbol{X}_i^q \in \mathcal{X}^q} \mathcal{L}(f_\Theta(\boldsymbol{X}_i^q | \mathcal{X}^s), y_i^q), \tag{1}$$

where $f_\Theta$ represents the entire model parameterized by $\Theta$.

### 3.2 Model Overview

We first provide an overview of our model (see the sketch of the pipeline in Fig. 2). Overall, our framework is composed of three modules, the backbone network $f_\theta$, the dynamic kernel generator $f_\omega$, and the context learning module $f_\phi$. The backbone network first extracts the feature maps from the input images, followed by a two-branch dynamic kernel network to obtain our proposed dynamic kernels. Specifically, the dynamic kernel generator independently refines the features in the support set to produce the instance kernel (*i.e.*, $\boldsymbol{G}^{in}$) in one branch. In another branch, the context learning module $f_\phi$ first produces a task-specific representation by refining and summarizing the features from the entire support set, which is then used as the input of the dynamic kernel generator to produce

the task-specific kernel (*i.e.*, $\boldsymbol{G}^{ta}$) adaptively. The dynamic kernel generator is shared across these two branches to generate the instance kernel and the task-specific kernel. Such a two-branch design enables the network to be aware of both the instance-level feature and global context feature of the support set. The design choice is justified in § 4. Finally, both the instance and task-specific kernels are fused and employed to boost the discrimination of instance features in the support set. The task-specific kernel is applied to the query feature maps for refining the representation without extracting a query instance kernel. Given the adapted support and query features, any post-matching algorithm (*e.g.*, metric-based or optimization-based FSL) can be employed seamlessly to achieve the few-shot classification task.

### 3.3   Dynamic Kernel Generator

In this part, we provide a detailed description of the dynamic kernel generator. The central component of our model is the dynamic kernel generator, which receives a feature map as input and produces a kernel adaptively. An essential problem one may face is that the feature map extracted by modern DNNs usually has a large size (*e.g.*, $640 \times 5 \times 5$ for ResNet-12). As such, we need to identify $c_{out} \times c \times k \times k$ amount of parameters for the dynamic kernel, where the $c_{out}$ is the output channel size (we consider $c = c_{out}$ in our paper), $c$ is the input channel size, and $k$ is the kernel size. This poses a significant issue since the generated dynamic kernel is prone to overfitting given the limited data of FSL. To tackle this problem, we consider designing our dynamic kernel generator in a decomposed manner [64]. As such, we develop a channel kernel network and a spatial kernel network to produce a kernel for each channel and each spatial location independently. This design exploits the information per data-sample to a large extent while reducing the number of parameters of a dynamic kernel, which greatly fits the low data regime of FSL. Fig. 3 illustrates the pipeline of the proposed dynamic kernel generator. In what follows, we detail out the operations of the channel and spatial kernel networks.

**Channel Kernel Network.**   Given a feature map $\boldsymbol{S} \in \mathbb{R}^{c \times h \times w}$, where $c$, $h$, and $w$ denote the size of channel, height and width, a common practice to realize the channel kernel network follows the SE block [20], as done for example in [8]. In the SE block, global average pooling (GAP) is performed over the feature map to encode the global representation, which is fed to the following sub-network for channel-wise scaling. The drawback of this design is that the GAP operator mainly preserves the low frequency components of the feature map, as shown in [37] (averaging is equivalent to low-pass filtering in the frequency domain). Thus, GAP discards important signal patterns in the feature map to a great degree. Clearly, low-frequency components cannot fully characterize the information encoded in a feature map, especially in the low data regime. We will empirically show this in § 4.3. To mitigate this issue, we opt for multi-spectral attention (MSA) to make better use of high-frequency patterns in the feature map. Given a feature map $\boldsymbol{S}$, we equally split the feature map into $n$ smaller tensors along the channel dimension (in our experiments

$n = 16$), as $\{\boldsymbol{S}^0, \boldsymbol{S}^1, \ldots, \boldsymbol{S}^{n-1} | \boldsymbol{S}^i \in \mathbb{R}^{\frac{c}{n} \times h \times w}\}$. Then, each channel of the tensor $\boldsymbol{S}^i$ is processed by a basis function of a 2D-<u>D</u>iscrete <u>C</u>osine <u>T</u>ransform (DCT) following the work of Qin *et al.* [37]. As a result, we obtain a real-valued feature vector in the frequency domain as $\boldsymbol{\tau}^i = \text{DCT}(\boldsymbol{S}^i), \boldsymbol{\tau}^i \in \mathbb{R}^{c/n}$, which is the frequency-encoded vector corresponding to $\boldsymbol{S}^i$. The frequency-encode vector for $\boldsymbol{S}$ is obtained by concatenating $\boldsymbol{\tau}^i$s as:

$$\boldsymbol{\tau} = \text{concat}(\boldsymbol{\tau}^0, \boldsymbol{\tau}^1, \ldots, \boldsymbol{\tau}^{n-1}), \tag{2}$$

where $\boldsymbol{\tau} \in \mathbb{R}^c$. Please refer to the supplementary material for the theoretical aspects of the MSA module. As compared to the global feature obtained by GAP, frequency-encoded feature $\boldsymbol{\tau}$ contains more diverse information patterns, which brings extra discriminative power to our model. This will be empirically discussed in § 4.3.

Once we obtain $\boldsymbol{\tau}$, a light-weight network is used to produce the channel kernel values adaptively. This network is realized by a two-layer MLP (or $1 \times 1$ convolution), with architecture as $c \rightarrow \sigma \times c \rightarrow \text{ReLu} \rightarrow k^2 \times c$, where $0 < \sigma < 1$, and $k$ controls the size of receptive field in the dynamic kernel (practically, we set $\sigma = 0.2$ and $k = 3$). Then we reshape the output vector into a $c \times k \times k$-sized tensor, denoted by $\boldsymbol{G}^{ch}$. We can obtain the final channel dynamic kernel $\hat{\boldsymbol{G}}^{ch} \in \mathbb{R}^{c \times h \times w \times k \times k}$ via applying batch normalization (BN) and spatial-wise broadcast to $\boldsymbol{G}^{ch}$.
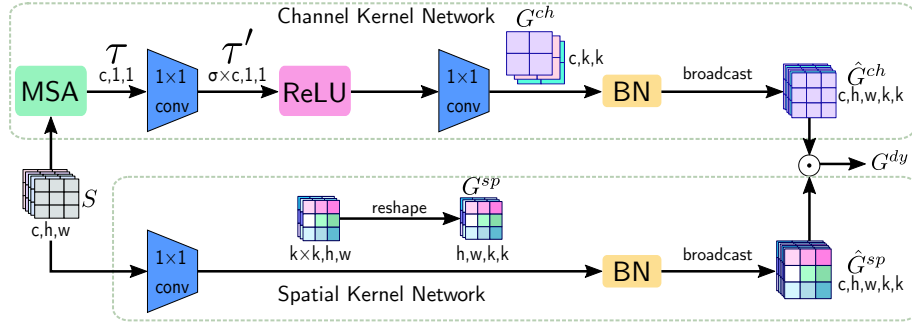


Fig. 3: The architecture of dynamic kernel generator. We adopt a decomposed architecture to produce dynamic kernels for the channels and spatial dimensions independently. This results in a light-weight kernel, which greatly fits the low-data regime of FSL

**Spatial Kernel Network.** Independent to the channel kernel network, our spatial kernel network adaptively produces a convolution kernel $\boldsymbol{G}^{sp}_{a,b} \in \mathbb{R}^{k \times k}$ for each spatial location of a feature map (*i.e.*, $\boldsymbol{S}_{:,a,b}$). This is achieved by using a $1 \times 1$ convolution layer, with the architecture of $(c, k^2, 1, 1)$ and the reshape operation. Therefore, the spatial kernel $\boldsymbol{G}^{sp}$ for all the spatial locations of a

feature map $\boldsymbol{S}$ has a size of $h \times w \times k \times k$. We then obtain the final spatial dynamic kernel $\hat{\boldsymbol{G}}^{sp} \in \mathbb{R}^{c \times h \times w \times k \times k}$ by applying the BN and channel-wise broadcast to $\boldsymbol{G}^{sp}$.

Finally, to unify the channel and spatial dynamic kernels, we apply the Hadamard product between the $\hat{\boldsymbol{G}}^{ch}$ and $\hat{\boldsymbol{G}}^{sp}$ to obtain the dynamic kernel $\boldsymbol{G}^{dy}$ as:

$$\boldsymbol{G}^{dy} = \hat{\boldsymbol{G}}^{sp} \odot \hat{\boldsymbol{G}}^{ch}, \tag{3}$$

where $\boldsymbol{G}^{dy} \in \mathbb{R}^{c \times h \times w \times k \times k}$. Notably, the number of parameters of our dynamic kernel is much less than a normal convolution kernel ($c \times h \times w \times k^2 \ll c_{out} \times c \times k^2$ given $c_{out} = c = 640$ and $h = w = 5$).

### 3.4   Dynamic Kernel

In this part, we discuss the process of obtaining the proposed INSTA dynamic kernel, which consists of the instance kernel and the task-specific kernel.

**Instance Kernel.** We defined the instance kernel as the dynamic kernel extracted from each support feature map. Formally, given the support set images $\mathcal{X}^s = \{\boldsymbol{X}_{11}^s, \ldots, \boldsymbol{X}_{NK}^s | \boldsymbol{X}_{ij}^s \in \mathbb{R}^{C \times H \times W}\}$, where $C$, $H$, and $W$ indicate the channel, height, and width of an image, respectively, the backbone network first extracts a feature map from each image, as $\mathcal{S} = f_\theta(\mathcal{X}^s)$, with $\mathcal{S} = \{\boldsymbol{S}_{11}, \ldots, \boldsymbol{S}_{NK} | \boldsymbol{S}_{ij} \in \mathbb{R}^{c \times h \times w}\}$. The feature map for the query sample can be obtained in a similar fashion, as $\boldsymbol{Q}_i = f_\theta(\boldsymbol{X}_i^q)$. We then use the dynamic kernel generator $f_\omega$ to independently produce a dynamic kernel $\boldsymbol{G}_{ij}^{in} \in \mathbb{R}^{c \times h \times w \times k \times k}$ for each support feature map $\boldsymbol{S}_{ij}$, following the process described in § 3.3.

**Task-Specific Kernel.** Adapting according to the whole context of support set is essential for FSL since it contains essential information of the task [36]. Following this intuition, we propose to learn the task-specific kernel to represent the knowledge of the task encoded in the support set. We achieve this by first using a context learning module $f_\phi$ to produce a fully context-aware representation for the support set by refining and aggregating the intermediate features. To be specific, we use four $1 \times 1$ convolution layers with a summation layer in the middle of the network to aggregate the $N \times K$ support features into one task-specific representation $\tilde{\boldsymbol{S}}$ (Please refer to the supplementary material for the conceptual diagram of $f_\phi$). Formally, this task representation can be obtained by:

$$\tilde{\boldsymbol{S}} = f_\phi(\mathcal{S}), \tag{4}$$

where $\tilde{\boldsymbol{S}} \in \mathbb{R}^{c \times h \times w}$ and the $\mathcal{S}$ denotes the entire support set. Then this task-specific representation $\tilde{\boldsymbol{S}}$ is used as the input to the dynamic kernel generator, following the process described in § 3.3 to adaptively produce our task-specific kernel $\boldsymbol{G}^{ta} \in \mathbb{R}^{c \times h \times w \times k \times k}$.

**INSTA Dynamic Kernel.** Once we have the instance and task-specific kernels, we fuse each instance kernel $\boldsymbol{G}_{ij}^{in}$ with task-specific kernel $\boldsymbol{G}^{ta}$ using the

Hadamard product, such that each dynamic kernel considers both the instance-level and task-level features. Formally, it can be described as:

$$\boldsymbol{G}_{ij}^{insta} = \boldsymbol{G}_{ij}^{in} \odot \boldsymbol{G}^{ta}, \tag{5}$$

where $\boldsymbol{G}_{ij}^{insta}$ is the final fused kernel corresponding to the $ij$-th sample in the support set.

After $\boldsymbol{G}_{ij}^{insta}$ is obtained, we apply it on its corresponding support feature map $\boldsymbol{S}_{ij}$. This is equivalent to first apply the task-specific kernel to extract the features, which is relevant to the task and then apply the instance kernel to further increase the discriminatory power of the resulting feature maps. While we only apply the task-specific kernel to the query feature map. This design choice is justified in § 4.3, where we compare our current design with a variant where we also extract instance kernel from the query feature maps. The dynamic convolution can be implemented using the Hadamard product between the un-folded feature map and the dynamic kernel. In doing so, the unfold operation first samples a $k \times k$ spatial region of the input feature map at each time and then stores the sampled region into extended dimensions, thereby obtaining an unfolded support feature map $\boldsymbol{S}_{ij}^{u} \in \mathbb{R}^{c \times h \times w \times k \times k}$ and an unfolded query feature map $\boldsymbol{Q}_{i}^{u} \in \mathbb{R}^{c \times h \times w \times k \times k}$. Then the adapted support and query feature maps are obtained as:

$$\dot{\boldsymbol{S}}_{ij} = \mathrm{AvgPool2d}(\boldsymbol{S}_{ij}^{u} \odot \boldsymbol{G}_{ij}^{insta}) + \boldsymbol{S}_{ij},$$
$$\dot{\boldsymbol{Q}}_{i} = \mathrm{AvgPool2d}(\boldsymbol{Q}_{i}^{u} \odot \boldsymbol{G}^{ta}) + \boldsymbol{Q}_{i}, \tag{6}$$

where $\dot{\boldsymbol{S}}_{ij} \in \mathbb{R}^{c \times h \times w}$ and $\dot{\boldsymbol{Q}}_{i} \in \mathbb{R}^{c \times h \times w}$. Note that the convolution operation between the dynamic kernel and the feature map in Fig. 2 is equivalent to the average of the Hadamard product between the unfolded feature map and the dynamic kernel over the $k \times k$ dimension (see Fig. 4). In our design, we adopt a residual connection between the original and the updated feature maps to obtain the final adapted features. Having the adapted support and query feature maps, any post-matching algorithms can be adopted to achieve FSL tasks. In § 4, we demonstrate that our algorithm can further boost the performance on various few-shot classification models (*e.g.*, MAML [15], ProtoNet [47] and EMD [62]) and the few-shot detection model [11].

*Remark 1.* The proposed dynamic kernels are convolutional filters and hence by nature differ from attention masks. As shown in Fig. 4, the attention mask merely re-weights each element of a feature map [13,20]. In contrast, our proposal in Eq. (6) realizes the dynamic convolution by first performing the element-wise multiplication between the unfolded feature map and the dynamic kernels, whose results are then averaged over the unfolded dimension. This is essentially equivalent to convolution operation (see Fig. 4).

## 4  Experiments

In this section, we first evaluate our method across four standard few-shot classification benchmarks: *mini*-ImageNet [38], *tiered*-ImageNet [39], CUB [53], and
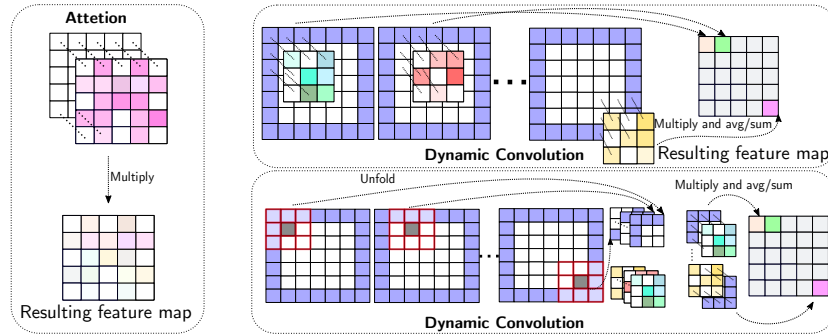
Fig. 4: Schematic comparison of the dynamic convolution and attention mechanism. The 2-D average of the element-wise multiplication between unfolded features and dynamic kernels is equivalent to using adaptive kernels sliding over the original feature map

FC100 [36]. Full details of the implementation are provided in the supplementary material. Furthermore, we evaluate the effectiveness of our framework for the few-shot detection task on the MS COCO and PASCAL VOC datasets [11]. Finally, we provide an ablation study to discuss the effect of each module in our framework. Please refer to supplementary material for a detailed description of each dataset.

### 4.1   Few-Shot Classification

We conduct few-shot classification experiments on three different state-of-the-art models, including MAML [15], ProtoNet [47] and DeepEMD [62] as baselines, and employ the proposed INSTA on top of them. We use the ResNet-10 backbone for the MAML and the ResNet-12 backbone for the other two baselines across all four benchmarks. For a fair comparison, we implement all the baseline models to report the results indicated by "*" across Table 1 – Table 2. Notably, for DeepEMD experiments, we adopt the *open-cv* solver instead of the *qpth* solver originally used in the paper to train the network due to resource capacity, which is indicated by "♣". Moreover, following the the same evaluation protocols in our baseline frameworks [7, 62], we report the mean accuracy with 95% confidence interval (Please refer to supplementary material for more implementation details).

***mini*-ImageNet.** As shown in Table 1, our method improves the performance of all the baseline models by a noticeable margin. It is worthwhile to mention that our INSTA can boost the performance of the baseline model ProtoNet by 4.72% and 3.67%, and achieves 67.01% and 83.13% for 5-way 1-shot and 5-way 5-shot settings, which outperforms many recent published models. Furthermore, we can show an improvement over a strong baseline model, *i.e.* DeepEMD, and achieve

Table 1: Few-shot classification accuracy and 95% confidence interval on *mini*-ImageNet and *tiered*-ImageNet with ResNet backbones

| Model | Backbone | *mini*-ImageNet | | *tiered*-ImageNet | |
|---|---|---|---|---|---|
| | | 5-way 1-shot | 5-way 5-shot | 5-way 1-shot | 5-way 5-shot |
| MetaOptNet [26] | ResNet-12 | $62.64 \pm 0.61$ | $78.63 \pm 0.46$ | $68.23 \pm 0.23$ | $84.03 \pm 0.56$ |
| TADAM [36] | ResNet-12 | $58.50 \pm 0.30$ | $76.70 \pm 0.30$ | - | - |
| FEAT [61] | ResNet-12 | $66.78 \pm 0.20$ | $82.05 \pm 0.14$ | $70.80 \pm 0.23$ | $84.79 \pm 0.16$ |
| CAN [19] | ResNet-12 | $63.85 \pm 0.48$ | $79.44 \pm 0.34$ | $69.89 \pm 0.51$ | $84.23 \pm 0.37$ |
| FRN [57] | ResNet-12 | $66.45 \pm 0.19$ | $82.83 \pm 0.13$ | $72.06 \pm 0.22$ | $86.89 \pm 0.14$ |
| InfoPatch [30] | ResNet-12 | $67.67 \pm 0.45$ | $82.44 \pm 0.31$ | $71.51 \pm 0.52$ | $85.44 \pm 0.35$ |
| GLoFA [33] | ResNet-12 | $66.12 \pm 0.42$ | $81.37 \pm 0.33$ | $69.75 \pm 0.33$ | $83.58 \pm 0.42$ |
| DMF [58] | ResNet-12 | $67.76 \pm 0.46$ | $82.71 \pm 0.31$ | $71.89 \pm 0.52$ | $85.96 \pm 0.35$ |
| MAML* [15] | ResNet-10 | $54.73 \pm 0.87$ | $66.72 \pm 0.81$ | $59.85 \pm 0.97$ | $73.20 \pm 0.81$ |
| **INSTA-MAML*** | ResNet-10 | $\mathbf{56.41 \pm 0.87}$ | $\mathbf{71.56 \pm 0.75}$ | $\mathbf{63.34 \pm 0.92}$ | $\mathbf{78.01 \pm 0.71}$ |
| ProtoNet* [47] | ResNet-12 | $62.29 \pm 0.33$ | $79.46 \pm 0.48$ | $68.25 \pm 0.23$ | $84.01 \pm 0.56$ |
| **INSTA-ProtoNet*** | ResNet-12 | $\mathbf{67.01 \pm 0.30}$ | $\mathbf{83.13 \pm 0.56}$ | $\mathbf{70.65 \pm 0.33}$ | $\mathbf{85.76 \pm 0.59}$ |
| DeepEMD*♣ [62] | ResNet-12 | $67.37 \pm 0.45$ | $83.17 \pm 0.75$ | $73.19 \pm 0.32$ | $86.79 \pm 0.61$ |
| **INSTA-DeepEMD*♣** | ResNet-12 | $\mathbf{68.46 \pm 0.48}$ | $\mathbf{84.21 \pm 0.82}$ | $\mathbf{73.87 \pm 0.31}$ | $\mathbf{88.02 \pm 0.61}$ |

state-of-the-art performance on this dataset. We provide more comparison between our approach and other state-of-the-art methods in our supplementary material for *mini*-ImageNet and *tiered*-ImageNet.

***tiered*-ImageNet.** As shown in Table 1, our model consistently brings the performance gain over baseline models. Among baseline models, our model shows a greater performance improvement over MAML than ProtoNet and DeepEMD. Notably, a significant improvement on 5-way 5-shot can be seen over DeepEMD. We attribute this improvement to our model's context module leveraging the categorical nature of this dataset. The results also show that the DeepEMD exhibits a performance gain with our design and outperforms many other recent models, which achieves the state-of-the-art result on this dataset.

**CUB.** The result on the CUB dataset is shown in Table 2, where the performance of all the three baseline models is improved by integrating the proposed INSTA. For the ProtoNet baseline with ResNet-12 backbone, as an example, 4.29% and 3.72% improvements can be observed for 5-way 1-shot and 5-way 5-shot, respectively. Moreover, among the improved models, our INSTA-DeepEMD achieves the state-of-the-art results, which are 75.26% and 88.12% for 5-way 1-shot and 5-way 5-shot settings, respectively. Please refer to supplementary material for additional results of the INSTA-ProtoNet with ResNet-18 backbone.

**FC100.** Consistent with the observation in the other benchmarks, the results in Table 2 again show that a performance improvement can be achieved on FC100. Furthermore, both INSTA-ProtoNet and INSTA-DeepEMD achieve comparable performance with the recent state-of-the-art models on FC100, which vividly shows the effectiveness of our approach.

Table 2: Few-shot classification accuracy and 95% confidence interval on CUB and FC100 with ResNet backbones

| Model | Backbone | CUB | | FC100 | |
|---|---|---|---|---|---|
| | | 5-way 1-shot | 5-way 5-shot | 5-way 1-shot | 5-way 5-shot |
| RelationNet [48] | ResNet-18 | $67.59 \pm 1.02$ | $82.75 \pm 0.58$ | - | - |
| Chen *et al.* [7] | ResNet-18 | 67.02 | 83.58 | - | - |
| SimpleShot [55] | ResNet-18 | 70.28 | 86.37 | - | - |
| Neg-Margin [29] | ResNet-18 | $72.66 \pm 0.85$ | $89.40 \pm 0.43$ | - | - |
| P-transfer [43] | ResNet-12 | $73.88 \pm 0.87$ | $87.81 \pm 0.48$ | - | - |
| TADAM [36] | ResNet-12 | - | - | $40.10 \pm 0.40$ | $56.10 \pm 0.40$ |
| ConstellationNet [59] | ResNet-12 | - | - | $43.80 \pm 0.20$ | $59.70 \pm 0.20$ |
| MAML* [15] | ResNet-10 | $70.46 \pm 0.97$ | $80.15 \pm 0.73$ | $34.50 \pm 0.69$ | $47.31 \pm 0.68$ |
| **INSTA-MAML*** | ResNet-10 | $\mathbf{73.08 \pm 0.97}$ | $\mathbf{84.26 \pm 0.66}$ | $\mathbf{38.02 \pm 0.70}$ | $\mathbf{51.20 \pm 0.68}$ |
| ProtoNet* [47] | ResNet-12 | $64.76 \pm 0.33$ | $77.99 \pm 0.68$ | $41.54 \pm 0.76$ | $57.08 \pm 0.76$ |
| **INSTA-ProtoNet*** | ResNet-12 | $\mathbf{69.05 \pm 0.32}$ | $\mathbf{81.71 \pm 0.63}$ | $\mathbf{44.12 \pm 0.26}$ | $\mathbf{62.04 \pm 0.75}$ |
| DeepEMD*♣ [62] | ResNet-12 | $74.55 \pm 0.30$ | $87.55 \pm 0.54$ | $45.12 \pm 0.26$ | $61.46 \pm 0.70$ |
| **INSTA-DeepEMD*♣** | ResNet-12 | $\mathbf{75.26 \pm 0.31}$ | $\mathbf{88.12 \pm 0.54}$ | $\mathbf{45.42 \pm 0.26}$ | $\mathbf{62.37 \pm 0.68}$ |

## 4.2   Few-Shot Detection

**Problem Definition.** Given a query and several support images (each support image only contains one object), a few-shot detection task is to output labels and corresponding bounding boxes for all objects in the query image that belong to support categories. Specifically, few-shot detection follows the N-way K-shot setting, where the support set contains N categories, with each category having K samples. In the following, we will discuss the experiment setup of this task. The implementation details on incorporating our INSTA with the baseline model proposed by Fan *et al.* [11] are included in the supplementary material.

**Experiment Setup.** In this experiment, we follow the training and testing settings in [11], where a 2-way 9-shot contrastive training strategy is adopted, and a 20-way 10-shot setting is used for final testing.

**Results.**

Table 3: The few-shot detection results of 6 different average precision (AP) on the test set, including MS COCO and PASCAL VOC datasets with 20-way 10-shot setting

| Model | AP | $AP_{50}$ | $AP_{75}$ | $AP_s$ | $AP_m$ | $AP_l$ |
|---|---|---|---|---|---|---|
| FR [23] | 5.6 | 12.3 | 4.6 | - | - | - |
| Meta [60] | 8.7 | 19.1 | 6.6 | - | - | - |
| Fan *et al.* [11] | 11.1 | 20.4 | 10.6 | 2.8 | 12.3 | 20.7 |
| **Fan *et al.* +INSTA** | **12.5** | **23.6** | **12.1** | **3.3** | **13.2** | **21.4** |

We compare our model against the baseline model and other relevant state-of-the-art few-shot detection models. As shown in Table 3, by incorporating the proposed INSTA on the baseline model, improved performance can be seen across all of the metrics. Specifically, as compared to the baseline model, the $AP_{50}$ is improved by 3.2% and achieves 23.6%, which empirically shows that our dynamic kernels indeed extract more informative and representative features. Moreover, the proposed INSTA improves the performance of detection for objects of all the scales (refer to $AP_s$, $AP_m$, and $AP_l$), which again illustrates the importance of our context module when adapting to different object scales in a given task.

### 4.3 Ablation Study

In the following section, we conduct ablation studies to discuss and verify the effect of each component of our framework, including the instance-kernel, task-specific kernel, MSA *vs*. GAP, and we also evaluate some variants of the proposed INSTA to justify the selection of our current framework. In this study, we use the ProtoNet as the baseline, and the ResNet-12 is adopted as its backbone. Additionally, this study is conducted on the *mini*-ImageNet under the 5-way 5-shot setting. The results are summarized in Table 4.

Table 4: The ablation study of each component in our framework

| ID | Model | Apply to $\boldsymbol{S}_{ij}$ | Apply to $\boldsymbol{Q}$ | 5-way 5-shot |
|---|---|---|---|---|
| (i) | ProtoNet | - | - | $79.46 \pm 0.48$ |
| (ii) | ProtoNet + $\boldsymbol{G}^{ta}$ | $\boldsymbol{G}^{ta}$ | - | $81.56 \pm 0.57$ |
| (iii) | ProtoNet + $\boldsymbol{G}^{ta}$ | $\boldsymbol{G}^{ta}$ | $\boldsymbol{G}^{ta}$ | $82.51 \pm 0.58$ |
| (iv) | ProtoNet + $\boldsymbol{G}^{in}$ | $\boldsymbol{G}^{in}$ | - | $80.81 \pm 0.60$ |
| (v) | ProtoNet + $\boldsymbol{G}^{insta}$ | $\boldsymbol{G}^{ta}, \boldsymbol{G}^{in}$ | - | $81.74 \pm 0.56$ |
| (vi) | INSTA-ProtoNet + GAP | $\boldsymbol{G}^{ta}, \boldsymbol{G}^{in}$ | $\boldsymbol{G}^{ta}$ | $82.07 \pm 0.56$ |
| (vii) | INSTA-ProtoNet + $\boldsymbol{G}_Q^{in}$ | $\boldsymbol{G}^{ta}, \boldsymbol{G}^{in}$ | $\boldsymbol{G}_Q^{in}, \boldsymbol{G}^{ta}$ | $82.24 \pm 0.56$ |
| (viii) | INSTA-ProtoNet w/o sharing $f_\omega$ | $\boldsymbol{G}^{ta}, \boldsymbol{G}^{in}$ | $\boldsymbol{G}^{ta}$ | $81.36 \pm 0.59$ |
| (ix) | INSTA-ProtoNet | $\boldsymbol{G}^{ta}, \boldsymbol{G}^{in}$ | $\boldsymbol{G}^{ta}$ | $\boldsymbol{83.13 \pm 0.56}$ |

**Effectiveness of Task-Specific Kernel.** We study the effect of the task-specific kernel in this experiment (*i.e.*, (ii) in Table 4), where we only apply the task-specific kernel to support samples and leave the query unchanged. By comparing (i) and (ii), we can observe that our task-specific kernel improves the performance of ProtoNet by 2.10%. Moreover, we apply the task-specific kernel on the query in setting (iii), and we can observe a further improvement over the setting (ii), which verifies the effectiveness of our task-specific kernel.
**Effectiveness of Instance Kernel.** In setting (iv), we enable the instance kernel based on ProtoNet (*i.e.*, (i)) and the performance of ProtoNet is improved by 1.35%. Moreover, comparing setting (ii), where the only task-specific kernel

is enabled, to setting (v), where both task-specific and instance kernels are enabled, a performance improvement can also be observed. Both cases illustrate the effectiveness of our instance kernel.

**Effectiveness of query adaptation.** The purpose of this experiment is to study the effect of the adaptation on the query feature. In setting (iii), we enable the task-specific kernel on both support and query features but disable the instance kernel. Compared to setting (ii), where the task-specific kernel is enabled only on the support features, setting (iii) yields a better performance. Furthermore, the comparison between setting (v) and (ix) highlights the importance of adapting the query feature map using the task information for FSL tasks.

**Effectiveness of MSA.** In this study, we show the performance gap between using GAP and MSA (§ 3.3). In setting (vi) we replace the MSA in our final design (ix) with GAP. As the results in Table 4 showed, the MSA indeed helps with learning a more informative and representative dynamic kernel than GAP.

**Instance Kernel for Query.** In our framework, we only use task-specific kernels on the query feature map since the instance kernels obtained from support samples might not be instance level representative for the query sample but provide useful task information. Therefore, we perform extra experiments to verify whether the instance kernel obtained from the query itself can extract better features. As the result in (vii) indicates, the instance kernel extracted from the query feature map does not improve the performance, as compared to our final design.

**Shared Dynamic Kernel Generator.** We hypothesize that sharing the Dynamic Kernel Network for the task-specific kernel and instance kernel encourages learning a more representative instance kernel. To verify this, we further conduct an experiment where two independent dynamic kernel generators are used to produce the task-specific kernel and instance kernel. As the comparison between (ix) and (viii) shows, inferring the instance and task kernels using a shared weight dynamic kernel generator is an essential design choice.

## 5    Conclusion

In this paper, we propose to learn a dynamic embedding function realized by a novel dynamic kernel, which extracts features at both instance-level and task-level while encoding important frequency patterns. Our method improves the performance of several FSL models. This is demonstrated on 4 public few-shot classification datasets, including *mini*-ImageNet, *tiered*-ImageNet, CUB, and FC100 and a few-shot detection dataset, namely MS COO-PASCAL-VOC.

## 6    Implementation Details

### 6.1    Datasets

***mini*-ImageNet.** The *mini*-ImageNet is sampled from ImageNet [10]. This dataset has 100 classes, with each having 600 samples. We follow the standard

protocol [38] to split the dataset into 64 training, 16 validation, and 20 testing classes.

***tiered*-ImageNet.** Similar to *mini*-ImageNet, *tiered*-ImageNet is also a subset of the ImageNet. This dataset consists of 608 classes from 34 categories and is split into 351 classes from 20 categories for training, 97 classes from 6 categories for validation, and 160 classes from 8 categories for testing.

**CUB.** The CUB is a fine-grained dataset, which consists of 11,788 images from 200 different breeds of birds. We follow the standard settings [29], in which the dataset is split into 100/50/50 breeds for training, validation, and testing, respectively.

**FC100.** FC100 dataset is a variant of the standard CIFAR100 dataset [25], which contains images from 100 classes, with each class containing 600 samples. We follow the standard setting [36], where the dataset is split into 60/20/20 classes for training, validation and testing, respectively.

**MS COCO and PASCAL VOC Datasets.** In the few-shot detection task, we follow the protocol used in [11] to construct the dataset, where images from 60 categories of the MS COCO dataset are used for training and images from the rest of 20 common categories between MS COCO and PASCAL VOC datasets are used for testing.

## 6.2    Few-Shot Classification Hyperparameters

**Network and Optimizer.** We use the ResNet-10 backbone [7] for the MAML and the ResNet-12 backbone [61, 62] for the other two baselines across all four benchmarks. Noted additional ResNet-18 back1bone [66] is employed for the ProtoNet experiments on CUB. We fix the size of input images to $84 \times 84$ for ProtoNet and DeepEMD baselines and $224 \times 224$ for MAML baseline (We strictly follow the same pre-processing protocol[2] in the original DeepEMD implementation to implement the metric-based baselines and our model, *i.e.*, DeepEMD, ProtoNet, INSTA-DeepEMD, and INSTA-ProtoNet. For MAML, we strictly follow the pre-processing protocol implemented in [7][3]). We use SGD optimizer for ProtoNet and DeepEMD experiments [61, 62] and AdamW optimizer for MAML experiments [7] across all the datasets. For ProtoNet and DeepEMD baselines, we use L2 regularizer with 0.0005 weight decay factor. In the MAML baseline, the weight dacay factor is 0.01. For ResNet-18 and ResNet-10 backbones, we disable the average pooling and remove the last fully connected (FC) layer to produce the feature maps with size of $512 \times 11 \times 11$ and $512 \times 7 \times 7$, respectively. In the ResNet-12 backbone, the network produces the feature map with a size of $640 \times 5 \times 5$.

**Training.** We follow the good practice in the state-of-the-art models [45, 61, 62], where the network training is split into two stages, *i.e.* pre-training and meta-training stages. During pre-training stage, the backbone network with a FC layer

---

[2] https://github.com/icoz69/DeepEMD.

[3] https://github.com/wyharveychen/CloserLookFewShot.

is trained on all the training classes of the dataset with the standard classification task. We select the network with the highest validation accuracy as the pre-trained backbone network for the meta-training stage. During meta-training stage, we follow the standard episodic training protocol [52] to train the entire model. We set a small learning rate (0.0002) for the backbone and a larger learning rate ($0.0002 \times 25$) for the other modules during meta-training stage. Additionally, we use cosine annealing learning rate scheduler over 200 epochs.
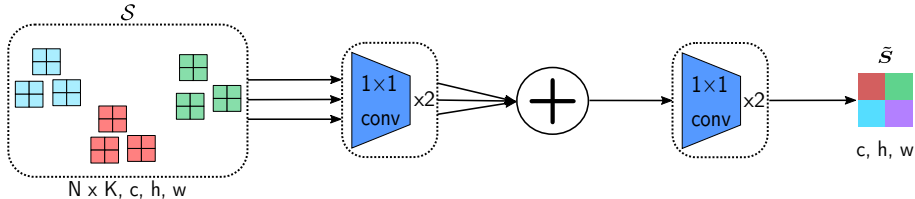


Fig. 5: The conceptual diagram of the context learning module

### 6.3    Few-Shot Detection

For the few-shot detection task, we consider the method proposed by Fan *et al.* [11] as our baseline model, which inherits from the faster-R-CNN [40] framework. Similar to the few-shot classification task, we implement our method to produce dynamic kernels and perform convolution on the feature maps extracted by the backbone network (*i.e.* ResNet-50), which is followed by a region proposal network (RPN), a region of interest pooling (ROI pooling) operation, and the classification and bounding box regression heads, outputting the categories and bounding boxes for the objects to the query image. Notably, For a fair comparison, we do not use any additional data augmentation. Please refer to [11] for more details of the framework and training strategy.

## 7    Additional Experiment Results

### 7.1    *mini*-ImageNet, *tiered*-ImageNet, and CUB

In this part, we provide extra comparison between our model and other state-of-the-art models on *mini*-ImageNet, *tiered*-ImageNet, and CUB. We follow the same evaluation protocol in [62][4] to evaluate the metric-based baseline models and our models (*i.e.*, DeepEMD, ProtoNet, INSTA-DeepEMD, and INSTA-ProtoNet), where 5,000 and 600 episodes are randomly sampled for 1-shot and 5-shot settings. For MAML baseline and INSTA-MAML, we strictly follow the

---

[4] https://github.com/icoz69/DeepEMD/blob/master/eval.py.

same evaluation protocol in [7][5], where 600 episodes are randomly sampled for both 1-shot and 5-shot settings.

Table 5: Few-shot classification accuracy and 95% confidence interval on *mini*-ImageNet and *tiered*-ImageNet with ResNet backbones

| Model | Backbone | *mini*-ImageNet | | *tiered*-ImageNet | |
|---|---|---|---|---|---|
| | | 5-way 1-shot | 5-way 5-shot | 5-way 1-shot | 5-way 5-shot |
| MetaOptNet-SVM [26] | ResNet-12 | $64.09 \pm 0.62$ | $80.00 \pm 0.45$ | - | - |
| Neg-Margin [29] | ResNet-12 | $63.85 \pm 0.81$ | $81.57 \pm 0.56$ | - | - |
| TPN [31] | ResNet-12 | 59.46 | 75.65 | - | - |
| DSN-MR [45] | ResNet-12 | $64.60 \pm 0.72$ | $79.51 \pm 0.50$ | $67.39 \pm 0.82$ | $82.82 \pm 0.56$ |
| $E^3BM$ [32] | ResNet-12 | $63.80 \pm 0.40$ | $80.10 \pm 0.30$ | $71.20 \pm 0.40$ | $85.30 \pm 0.30$ |
| ConstellationNet [59] | ResNet-12 | $64.89 \pm 0.23$ | $79.95 \pm 0.37$ | - | - |
| MELR [14] | ResNet-12 | $67.40 \pm 0.43$ | $83.40 \pm 0.28$ | $72.14 \pm 0.51$ | $87.01 \pm 0.35$ |
| CNL [63] | ResNet-12 | $67.96 \pm 0.98$ | $83.36 \pm 0.51$ | $73.42 \pm 0.95$ | $87.72 \pm 0.75$ |
| MAML* [15] | ResNet-10 | $54.73 \pm 0.87$ | $66.72 \pm 0.81$ | $59.85 \pm 0.97$ | $73.20 \pm 0.81$ |
| **INSTA-MAML*** | ResNet-10 | $\mathbf{56.41 \pm 0.87}$ | $\mathbf{71.56 \pm 0.75}$ | $\mathbf{63.34 \pm 0.92}$ | $\mathbf{78.01 \pm 0.71}$ |
| ProtoNet* [47] | ResNet-12 | $62.29 \pm 0.33$ | $79.46 \pm 0.48$ | $68.25 \pm 0.23$ | $84.01 \pm 0.56$ |
| **INSTA-ProtoNet*** | ResNet-12 | $\mathbf{67.01 \pm 0.30}$ | $\mathbf{83.13 \pm 0.56}$ | $\mathbf{70.65 \pm 0.33}$ | $\mathbf{85.76 \pm 0.59}$ |
| DeepEMD*♣ [62] | ResNet-12 | $67.37 \pm 0.45$ | $83.17 \pm 0.75$ | $73.19 \pm 0.32$ | $86.79 \pm 0.61$ |
| **INSTA-DeepEMD*♣** | ResNet-12 | $\mathbf{68.46 \pm 0.48}$ | $\mathbf{84.21 \pm 0.82}$ | $\mathbf{73.87 \pm 0.31}$ | $\mathbf{88.02 \pm 0.61}$ |

Table 6: Few-shot classification accuracy and 95% confidence interval on CUB with ResNet backbones

| Model | Backbone | 5-way 1-shot | 5-way 5-shot |
|---|---|---|---|
| MAML* [15] | ResNet-10 | $70.46 \pm 0.97$ | $80.15 \pm 0.73$ |
| **INSTA-MAML*** | ResNet-10 | $\mathbf{73.08 \pm 0.97}$ | $\mathbf{84.26 \pm 0.66}$ |
| DeepEMD*♣ [62] | ResNet-12 | $74.55 \pm 0.30$ | $87.55 \pm 0.54$ |
| **INSTA-DeepEMD*♣** | ResNet-12 | $\mathbf{75.26 \pm 0.31}$ | $\mathbf{88.12 \pm 0.54}$ |
| ProtoNet* | ResNet-18 | $75.06 \pm 0.30$ | $87.39 \pm 0.48$ |
| **INSTA-ProtoNet*** | ResNet-18 | $\mathbf{77.18 \pm 0.29}$ | $\mathbf{89.54 \pm 0.44}$ |

## 7.2    Meta-Dataset

To verify the effectiveness of our method on the cross-domain few-shot classification problem, we incorporate INSTA into the baseline model simple-CNAPS [3]. In this experiment, we fix the trained baseline model and only fine-tune the

---

[5] https://github.com/wyharveychen/CloserLookFewShot/blob/master/test.py.

Table 7: Few-shot classification results on Meta-dataset with ResNet-18 backbone

| Dataset | Simple-CNAPS | INSTA-Simple-CNAPS |
|---|---|---|
| ILSVRC | $55.5 \pm 1.1$ | $\mathbf{58.5 \pm 1.1}$ |
| Omniglot | $91.0 \pm 0.6$ | $\mathbf{91.9 \pm 0.6}$ |
| Aircraft | $81.2 \pm 0.7$ | $\mathbf{82.4 \pm 0.8}$ |
| Birds | $74.3 \pm 0.9$ | $\mathbf{75.7 \pm 0.8}$ |
| Textures | $66.9 \pm 0.8$ | $\mathbf{67.8 \pm 0.8}$ |
| Quick Draw | $76.7 \pm 0.8$ | $\mathbf{76.8 \pm 0.8}$ |
| Fungi | $47.5 \pm 1.0$ | $\mathbf{49.2 \pm 1.1}$ |
| VGG Flowers | $\mathbf{90.5 \pm 0.6}$ | $90.4 \pm 0.6$ |
| Traffic Signs | $72.0 \pm 0.7$ | $\mathbf{74.1 \pm 0.7}$ |
| MSCOCO | $47.3 \pm 1.1$ | $\mathbf{53.9 \pm 1.1}$ |

modules to generate the INSTA dynamic kernels (*i.e.*, dynamic kernel generator and context learning module). We follow the same implementation of simple-CNAPS[6] to conduct this experiment (*e.g.*, $8 \times 10^{-3}$ as learning rate, Adam as the optimizer, *etc.*). As the results in Table 7 suggested, INSTA improves the baseline over almost all the datasets, which again shows the effectiveness of our proposed INSTA dynamic kernels.

### 7.3   Ablation Study

In this part, we provide extra ablation studies on the effect of the residual connection and the spatial size of our dynamic kernel.

**Residual Connection**. In this experiment, we study the effect of the residual connection in our framework. In setting (i) of the Table 8, we disable the residual connection between the adapted and original features. The result suggests that the residual connection is an essential design choice for our framework.

**Kernel Size**. We provide an extra study on the effect of the spatial size of our dynamic kernel. Given that the feature map extracted by ResNet-12 has spatial size $5 \times 5$, the dynamic kernel size is constrained smaller or equal to $5 \times 5$. Therefore, in this study, we compare the results when the dynamic kernel size $k = 5 \times 5$ to our final design ($k = 3 \times 3$).

## 8   Multi-Spectral Attention

In this section, we provide more details for using the 2D-DCT to obtain the frequency-encoded vector. We first introduce the basis function of the 2D-DCT.

---

[6] https://github.com/peymanbateni/simple-cnaps/tree/master/simple-cnaps-src.

Table 8: The extra ablation study for the effect of the residual connection and spatial size of our dynamic kernel

| ID | Model | 5-way 5-shot |
|---|---|---|
| (i) | INSTA-ProtoNet w/o residual | 80.03 |
| (ii) | INSTA-ProtoNet w/ $5 \times 5$ $\boldsymbol{G}^{dy}$ | 82.43 |
| (iii) | INSTA-ProtoNet | 83.13 |

The basis $B_{u,v}^{a,b}$ of the 2D-DCT is given by:

$$B_{u,v}^{a,b} = \cos(\frac{\pi u}{h}(a + \frac{1}{2}))\cos(\frac{\pi v}{w}(b + \frac{1}{2})), \tag{7}$$

where $u, v$ are the frequency components of a basis. Then the frequency-encoded vector of a 3D-tensor $\boldsymbol{S}^i \in \mathbb{R}^{\frac{c}{n} \times h \times w}$ can be obtained by:

$$\boldsymbol{\tau}^i = \sum_{a=0}^{h-1} \sum_{b=0}^{w-1} \boldsymbol{S}_{:,a,b} B_{u_i,v_i}^{a,b}$$

$$s.t. \ i \in \{0, 1, \ldots, n-1\},$$

where $\boldsymbol{\tau}^i \in \mathbb{R}^{\frac{c}{n}}$ is the frequency-encoded vector, $h$ and $w$ are the height and width of the input signal. We pick the lowest 16 frequency components for our basis function according to [37]. Finally, we concatenate all the frequency-encoded vectors as:

$$\boldsymbol{\tau} = \text{concat}(\boldsymbol{\tau}^0, \boldsymbol{\tau}^1, \ldots, \boldsymbol{\tau}^{n-1}), \tag{8}$$

where $\boldsymbol{\tau} \in \mathbb{R}^c$.

# References

1. Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M.W., Pfau, D., Schaul, T., Shillingford, B., De Freitas, N.: Learning to learn by gradient descent by gradient descent. In: Advances in neural information processing systems. pp. 3981–3989 (2016) 3

2. Antoniou, A., Edwards, H., Storkey, A.: How to train your maml. In: International Conference on Learning Representations (2019) 3

3. Bateni, P., Goyal, R., Masrani, V., Wood, F., Sigal, L.: Improved few-shot visual classification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) 17

4. Bertinetto, L., Henriques, J.F., Torr, P., Vedaldi, A.: Meta-learning with differentiable closed-form solvers. In: International Conference on Learning Representations (2018) 3

5. Bertinetto, L., Henriques, J.F., Valmadre, J., Torr, P.H., Vedaldi, A.: Learning feedforward one-shot learners. In: Proceedings of the 30th International Conference on Neural Information Processing Systems. pp. 523–531 (2016) 3

6.  Bolukbasi, T., Wang, J., Dekel, O., Saligrama, V.: Adaptive neural networks for efficient inference. In: International Conference on Machine Learning. pp. 527–536. PMLR (2017) 4

7.  Chen, W.Y., Liu, Y.C., Kira, Z., Wang, Y.C.F., Huang, J.B.: A closer look at few-shot classification. arXiv preprint arXiv:1904.04232 (2019) 3, 10, 12, 15, 17

8.  Chen, Y., Dai, X., Liu, M., Chen, D., Yuan, L., Liu, Z.: Dynamic convolution: Attention over convolution kernels. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11030–11039 (2020) 4, 6

9.  Choi, J., Krishnamurthy, J., Kembhavi, A., Farhadi, A.: Structured set matching networks for one-shot part labeling. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3627–3636 (2018) 3

10. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) 14

11. Fan, Q., Zhuo, W., Tang, C.K., Tai, Y.W.: Few-shot object detection with attention-rpn and multi-relation detector. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4013–4022 (2020) 9, 10, 12, 15, 16

12. Fang, P., Harandi, M., Petersson, L.: Kernel methods in hyperbolic spaces. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10665–10674 (2021) 3

13. Fang, P., Zhou, J., Roy, S.K., Ji, P., Petersson, L., Harandi, M.: Attention in attention networks for person retrieval. IEEE Transactions on Pattern Analysis and Machine Intelligence pp. 1–1 (2021) 9

14. Fei, N., Lu, Z., Xiang, T., Huang, S.: Melr: Meta-learning via modeling episode-level relationships for few-shot learning. In: International Conference on Learning Representations (2020) 3, 17

15. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. arXiv preprint arXiv:1703.03400 (2017) 1, 2, 3, 9, 10, 11, 12, 17

16. Flennerhag, S., Rusu, A.A., Pascanu, R., Visin, F., Yin, H., Hadsell, R.: Meta-learning with warped gradient descent. arXiv preprint arXiv:1909.00025 (2019) 3

17. Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., Pontil, M.: Bilevel programming for hyperparameter optimization and meta-learning. arXiv preprint arXiv:1806.04910 (2018) 3

18. Ha, D., Dai, A., Le, Q.V.: Hypernetworks. arXiv preprint arXiv:1609.09106 (2016) 4

19. Hou, R., Chang, H., Ma, B., Shan, S., Chen, X.: Cross attention network for few-shot classification. arXiv preprint arXiv:1910.07677 (2019) 3, 11

20. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7132–7141 (2018) 6, 9

21. Huang, G., Chen, D., Li, T., Wu, F., van der Maaten, L., Weinberger, K.: Multi-scale dense networks for resource efficient image classification. In: International Conference on Learning Representations (2018) 4

22. Jia, X., De Brabandere, B., Tuytelaars, T., Gool, L.V.: Dynamic filter networks. Advances in neural information processing systems **29**, 667–675 (2016) 2

23. Kang, B., Liu, Z., Wang, X., Yu, F., Feng, J., Darrell, T.: Few-shot object detection via feature reweighting. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8420–8429 (2019) 12

24. Koch, G., Zemel, R., Salakhutdinov, R., et al.: Siamese neural networks for one-shot image recognition. In: ICML deep learning workshop. vol. 2. Lille (2015) 3
25. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009) 15
26. Lee, K., Maji, S., Ravichandran, A., Soatto, S.: Meta-learning with differentiable convex optimization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 10657–10665 (2019) 3, 11, 17
27. Lee, Y., Choi, S.: Gradient-based meta-learning with learned layerwise metric and subspace. In: International Conference on Machine Learning. pp. 2927–2936. PMLR (2018) 3
28. Li, H., Eigen, D., Dodge, S., Zeiler, M., Wang, X.: Finding task-relevant features for few-shot learning by category traversal. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1–10 (2019) 3
29. Liu, B., Cao, Y., Lin, Y., Li, Q., Zhang, Z., Long, M., Hu, H.: Negative margin matters: Understanding margin in few-shot classification. arXiv preprint arXiv:2003.12060 (2020) 12, 15, 17
30. Liu, C., Fu, Y., Xu, C., Yang, S., Li, J., Wang, C., Zhang, L.: Learning a few-shot embedding model with contrastive learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 8635–8643 (2021) 3, 11
31. Liu, Y., Lee, J., Park, M., Kim, S., Yang, E., Hwang, S.J., Yang, Y.: Learning to propagate labels: Transductive propagation network for few-shot learning. arXiv preprint arXiv:1805.10002 (2018) 3, 17
32. Liu, Y., Schiele, B., Sun, Q.: An ensemble of epoch-wise empirical bayes for few-shot learning. In: European Conference on Computer Vision. pp. 404–421. Springer (2020) 3, 17
33. Lu, S., Ye, H.J., Zhan, D.C.: Tailoring embedding function to heterogeneous few-shot tasks by global and local feature adaptors. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 8776–8783 (2021) 1, 2, 3, 11
34. Ma, R., Fang, P., Drummond, T., Harandi, M.: Adaptive poincaré point to set distance for few-shot classification. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 1926–1934 (2022) 1
35. Nichol, A., Achiam, J., Schulman, J.: On first-order meta-learning algorithms. arXiv preprint arXiv:1803.02999 (2018) 1, 3
36. Oreshkin, B.N., Rodriguez, P., Lacoste, A.: Tadam: task dependent adaptive metric for improved few-shot learning. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems. pp. 719–729 (2018) 3, 8, 10, 11, 12, 15
37. Qin, Z., Zhang, P., Wu, F., Li, X.: Fcanet: Frequency channel attention networks. arXiv preprint arXiv:2012.11879 (2020) 6, 7, 19
38. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2017) 9, 15
39. Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J.B., Larochelle, H., Zemel, R.S.: Meta-learning for semi-supervised few-shot classification. arXiv preprint arXiv:1803.00676 (2018) 1, 9
40. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems 28, 91–99 (2015) 16
41. Rusu, A.A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., Hadsell, R.: Meta-learning with latent embedding optimization. arXiv preprint arXiv:1807.05960 (2018) 3

42. Satorras, V.G., Estrach, J.B.: Few-shot learning with graph neural networks. In: International Conference on Learning Representations (2018) 3
43. Shen, Z., Liu, Z., Qin, J., Savvides, M., Cheng, K.T.: Partial is better than all: Revisiting fine-tuning strategy for few-shot learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 9594–9602 (2021) 12
44. Shyam, P., Gupta, S., Dukkipati, A.: Attentive recurrent comparators. In: International Conference on Machine Learning. pp. 3173–3181. PMLR (2017) 3
45. Simon, C., Koniusz, P., Nock, R., Harandi, M.: Adaptive subspaces for few-shot learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4136–4145 (2020) 3, 15, 17
46. Simon, C., Koniusz, P., Nock, R., Harandi, M.: On modulating the gradient for meta-learning. In: European Conference on Computer Vision. pp. 556–572. Springer (2020) 3
47. Snell, J., Swersky, K., Zemel, R.: Prototypical networks for few-shot learning. In: Advances in neural information processing systems. pp. 4077–4087 (2017) 1, 3, 9, 10, 11, 12, 17
48. Sung, F., Yang, Y., Zhang, L., Xiang, T., Torr, P.H., Hospedales, T.M.: Learning to compare: Relation network for few-shot learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1199–1208 (2018) 1, 3, 12
49. Teerapittayanon, S., McDanel, B., Kung, H.T.: Branchynet: Fast inference via early exiting from deep neural networks. In: 2016 23rd International Conference on Pattern Recognition (ICPR). pp. 2464–2469. IEEE (2016) 4
50. Triantafillou, E., Zemel, R., Urtasun, R.: Few-shot learning through an information retrieval lens. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. pp. 2252–2262 (2017) 3
51. Veit, A., Belongie, S.: Convolutional networks with adaptive inference graphs. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 3–18 (2018) 4
52. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: Advances in neural information processing systems. pp. 3630–3638 (2016) 1, 3, 16
53. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset (2011) 9
54. Wang, X., Yu, F., Dou, Z.Y., Darrell, T., Gonzalez, J.E.: Skipnet: Learning dynamic routing in convolutional networks. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 409–424 (2018) 4
55. Wang, Y., Chao, W.L., Weinberger, K.Q., van der Maaten, L.: Simpleshot: Revisiting nearest-neighbor classification for few-shot learning. arXiv preprint arXiv:1911.04623 (2019) 12
56. Wang, Y., Yao, Q., Kwok, J.T., Ni, L.M.: Generalizing from a few examples: A survey on few-shot learning. ACM Computing Surveys (CSUR) 53(3), 1–34 (2020) 2
57. Wertheimer, D., Tang, L., Hariharan, B.: Few-shot classification with feature map reconstruction networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8012–8021 (2021) 3, 11
58. Xu, C., Fu, Y., Liu, C., Wang, C., Li, J., Huang, F., Zhang, L., Xue, X.: Learning dynamic alignment via meta-filter for few-shot learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5182–5191 (2021) 2, 4, 11

59. Xu, W., Wang, H., Tu, Z., et al.: Attentional constellation nets for few-shot learning. In: International Conference on Learning Representations (2020) 3, 12, 17
60. Yan, X., Chen, Z., Xu, A., Wang, X., Liang, X., Lin, L.: Meta r-cnn: Towards general solver for instance-level low-shot learning. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9577–9586 (2019) 12
61. Ye, H.J., Hu, H., Zhan, D.C., Sha, F.: Few-shot learning via embedding adaptation with set-to-set functions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8808–8817 (2020) 1, 2, 3, 11, 15
62. Zhang, C., Cai, Y., Lin, G., Shen, C.: Deepemd: Few-shot image classification with differentiable earth mover's distance and structured classifiers. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12203–12213 (2020) 1, 3, 9, 10, 11, 12, 15, 16, 17
63. Zhao, J., Yang, Y., Lin, X., Yang, J., He, L.: Looking wider for better adaptive representation in few-shot learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 10981–10989 (2021) 3, 17
64. Zhou, J., Jampani, V., Pi, Z., Liu, Q., Yang, M.H.: Decoupled dynamic filter networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6647–6656 (2021) 2, 4, 6
65. Zhu, T., Ma, R., Harandi, M., Drummond, T.: Learning online for unified segmentation and tracking models. In: 2021 International Joint Conference on Neural Networks (IJCNN). pp. 1–8 (2021) 3
66. Ziko, I., Dolz, J., Granger, E., Ayed, I.B.: Laplacian regularized few-shot learning. In: International Conference on Machine Learning. pp. 11660–11670. PMLR (2020) 15