# SEGA: Semantic Guided Attention on Visual Prototype for Few-Shot Learning

Fengyuan Yang[1,2], Ruiping Wang[1,2,3], Xilin Chen[1,2]

[1]Key Laboratory of Intelligent Information Processing of Chinese Academy of Sciences (CAS),
Institute of Computing Technology, CAS, Beijing, 100190, China
[2]University of Chinese Academy of Sciences, Beijing, 100049, China
[3]Beijing Academy of Artificial Intelligence, Beijing, 100084, China

`fengyuan.yang@vipl.ict.ac.cn`, {`wangruiping, xlchen`}`@ict.ac.cn`

## Abstract

*Teaching machines to recognize a new category based on few training samples especially only one remains challenging owing to the incomprehensive understanding of the novel category caused by the lack of data. However, human can learn new classes quickly even given few samples since human can tell what discriminative features should be focused on about each category based on both the visual and semantic prior knowledge. To better utilize those prior knowledge, we propose the **SEmantic Guided Attention (SEGA)** mechanism where the semantic knowledge is used to guide the visual perception in a top-down manner about what visual features should be paid attention to when distinguishing a category from the others. As a result, the embedding of the novel class even with few samples can be more discriminative. Concretely, a feature extractor is trained to embed few images of each novel class into a visual prototype with the help of transferring visual prior knowledge from base classes. Then we learn a network that maps semantic knowledge to category-specific attention vectors which will be used to perform feature selection to enhance the visual prototypes. Extensive experiments on miniImageNet, tieredImageNet, CIFAR-FS, and CUB indicate that our semantic guided attention realizes anticipated function and outperforms state-of-the-art results.*

## 1. Introduction

Object recognition has been significantly improved in the past decade with the rapid growth of data scales and the help of deep learning methods [20,47,50]. However, the frequency distribution of visual categories generally presents the form of long-tailed distribution [39] which means it is difficult to collect a sufficient number of samples for categories on the tail part. Even for categories on the head part, a large scale of image annotation is also a heavy and expensive job. Therefore, we study this kind of more realistic task
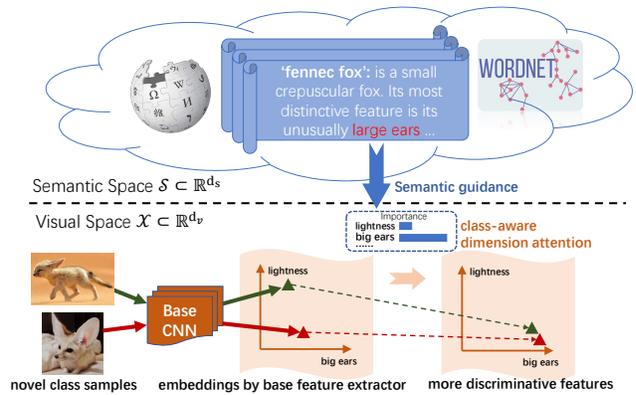


Figure 1: The illustration diagram shows the motivation of ours semantic guided attention. In few-shot learning, the visual cognition of novel class is incomprehensive given few labeled images. However, semantic knowledge can give guidance about what key feature dimensions of this category should be focused on. By applying semantic attention to visual features, we can have more discriminative recognition of the novel class.

named few-shot learning (FSL) which means learning new categories based on few labeled samples.

One of the main challenges of few-shot learning is that when only given few labeled samples, the recognizer cannot get the comprehensive recognition of the novel class. To deal with this challenge, prior knowledge is of vital importance since the reason why we human beings can learn new categories quickly and efficiently is that we have already learned so many base categories before. Therefore in few-shot learning, we usually transfer knowledge from base classes with a large number of labeled images to the target novel classes with a small number of labeled images. Thus most previous works have focused on how to transfer visual prior knowledge efficiently from base classes to novel classes [8, 11, 48, 54]. However, the semantic prior knowledge plays a pivotal role in human learning too. For
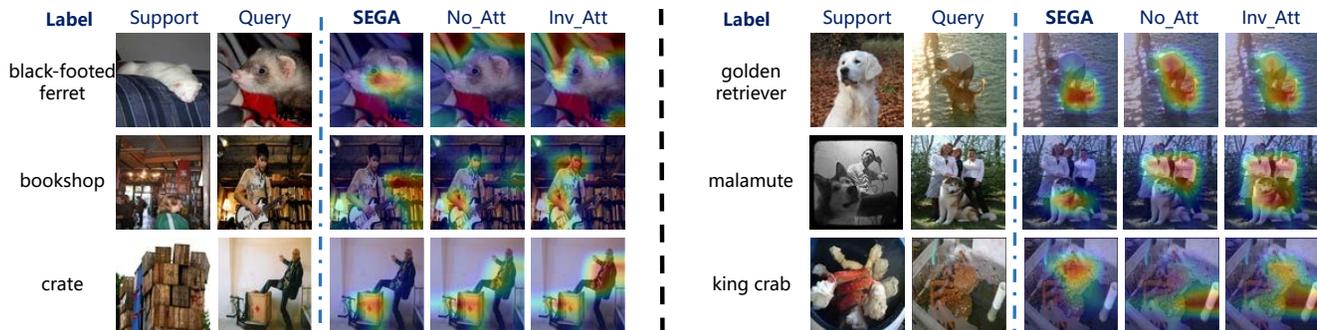
Figure 2: The Grad-CAM visualization testing on miniImageNet's unseen classes under 1-Shot scenario. Column "*Support*" gives the only training sample for each novel class. Column "*Query*" shows the query image of this novel class. The next three columns show the query image's Grad-CAM visualization when applying our semantic guided attention("*SEGA*"), not applying any attention("*No_Att*") and applying the inversed version of our semantic guided attention("*Inv_Att*") respectively. Warmer color with a higher value.

example, the category name along with its semantic explanation is always mentioned when parents teach children to recognize the animals in the atlas, from which children can have more comprehensive recognition about the new category and the relation between this category and categories learned before. More specifically, from the semantic guidance human can get what key features of the category should be focused on and what noisy features should be ignored as illustrated in Figure 1, which could help us to tackle the above challenge of few-shot learning.

Our motivation for using semantic guidance is derived from cognitive neurosciences. The reason why a human can perform *"object constancy"*, which is the ability to identify objects across changes in the detailed context including illumination, object pose, and background [53], is that our object recognition system can tell the key discriminative features concerning each category. For the categories with enough samples, the key features can be concluded from a large number of images. But as for the categories with few images, human can also get which key features should be focused on directly from the high-level guidance of semantic information such as the class name. Apart from cognitive neurosciences, we can also find a similar idea in contrastive learning [15,31] which suggests that there is no need to precisely reconstruct the dollar bill but only need to tell key features of the bill to distinguish it from other objects [6]. We cannot rely only on the semantic information, which is compact, to directly reconstruct the visual prototype well enough. However, it makes sense to use the semantic information to get which key features should be paid attention to when performing classification.

Therefore, we propose a more human-like way to utilize label semantic knowledge inspired by the above theories. A feature extractor is trained to transfer the visual prior knowledge from base classes to obtain a feature space. After that, our framework learns a network that maps the class semantic knowledge to the class-specific attention which implies

key dimensions of visual features. Then during testing, the visual feature attention of novel classes can be obtained by taking the names of novel classes as input to the above mapping. Finally, the attention will be applied to the visual prototype to highlight the key features so that a more discriminative representation for the novel class can be obtained.

More specifically, we take the first Grad-CAM [46] result shown in Figure 2 as an example to illustrate our idea. Assuming that our machine has never seen the images of *black-footed ferret* before. If ignoring label semantic information, the model can be easily misguided by background noise and large intra-class variations when learning with only one sample of *black-footed ferret*, thus may take the striped background as the key feature of this category. On the contrary, our method (SEGA) can tell that *black-footed ferret* is a kind of animal based on semantic knowledge, which guides our model to pay attention to key features when dealing with animals.

## 2. Related Works

The mainstream few-shot learning methods utilize the auxiliary task to transfer visual prior knowledge from base classes to few-shot classes [34, 57]. From the perspective of how to design auxiliary tasks, there exist common approaches such as metric-based, hallucination-based, and meta-learning-based. However, the above taxonomy cannot separate current methods well enough. Here we give a more orthogonal taxonomy from the perspective of how to use the support set. The first kind uses the support set to directly generate the classifiers for novel classes, which is related to lazy learning [2]. It contains the metric-based methods [18,30,48,49,54] and the classification weight generation methods [11,17,35,36]. The other kind uses the support set to finetune the end-to-end network, which is related to eager learning. It contains hallucination-based [1,9,14,44] and optimization-based methods [8,37,42,43,56]. Our approach belongs to the lazy learning category.

**Lazy learning methods.** The first kind of lazy learning method is metric-based methods, whose core idea is to learn a metric space in which the samples of the same category are near each other while different categories are far away. This idea can be traced back to NCA [13] and LMNN [59]. Nowadays, there emerge many deep learning metric-based methods such as Matching Networks [54] which learns a nearest neighbor classifier with the help of context information, Prototypical Networks [48] which generates class prototype as the mean of support set samples, Relation Network [49] which uses the neural network to model the distance measurement, *etc*. Another kind of lazy learning method is classification weight generation methods, which generate the class weight for novel classes directly [11, 35, 36]. The cosine classifier is widely used to avoid the norm problem [17, 35]. Our approach lies in this classification weight generation category.

Most recently, there emerge works that use semantic knowledge in few-shot learning [5,25,32,45,60] whose idea comes from a related research area named zero-shot learning [7, 21, 22]. The source of semantic knowledge can be attributes, word embeddings, and even knowledge graphs.

**Semantic knowledge in few-shot learning.** Thanks to the development in natural language processing, we can get label embeddings from the pre-training word embedding models such as GloVe [33]. TriNet [5] takes the class label embedding to hallucinate new samples in semantic feature space. TRAML [24] uses class label embedding to generate adaptive margin loss. In addition, AM3 [60] uses label embedding to generate a semantic prototype which is used to perform a convex combination with the visual prototype to form the final class representation. Furthermore, MultiSem [45] introduces extra semantic like verbal descriptions and more recent work [63] extracts parts/attributes from WordNet [28]. Apart from semantic knowledge from language, correlation knowledge, which can be obtained from knowledge graph (*e.g*. NEIL [29], WordNet [28], *etc*.), can also be helpful in FSL. KTN [32] proposes to construct a GCN in which the node representation comes from label embedding and the edge comes from knowledge graph to transfer knowledge from base class to novel class. KGTN [4] employs a similar idea to construct a GCN whose edge weights are generated by the semantic hierarchy of categories.

As we can see, previous methods can be divided into two paradigms: semantic-dominated (more rely on semantic, *e.g*. TriNet [5], TRAML [24] *etc*.) and multimodal-fusion (fuse semantic and visual equally, *e.g*. AM3 [60], KTN [32] *etc*.). Here in this paper, we propose a new paradigm that is visual-dominated (*i.e*. just use semantic to enhance, *e.g*. our SEGA uses semantic just to generate visual attention instead of reconstruction). The advantages lie in (1) More robust and easier to learn since only need to learn attention while previous two paradigms overuse semantic to re-

construct visual information (*e.g*. AM3 [60] reconstructs semantic prototype directly in visual space). (2) More reasonable since the essence of semantic is invariant (like attention in SEGA), we shouldn't expect it useful in equivariant jobs (like generation/fusion in the previous two paradigms). (3) Customized for FSL since class-specific semantic attention can eliminate background noise and intra-class variations which are inevitably amplified in FSL.

## 3. Approach

We give some fundamental analysis for why semantic knowledge could be helpful in FSL in §3.1; then we establish preliminaries about problem setting in §3.2 and introduce general framwork in §3.3; finally, we focus on our SEGA in §3.4 followed by some discussion in §3.5.

### 3.1. Why Utilize Semantic Knowledge?

The introduction of semantic knowledge is nothing new in ZSL since the ZSL task cannot be completed without semantic knowledge. While in FSL previous works hardly use semantic knowledge. Most recently, there emerge some works that propose to utilize semantic knowledge in FSL [5, 32, 45, 60]. However, there lacks of fundamental analysis about why using semantic knowledge can help in FSL. Thus, we adopt the Canonical Correlation Analysis (CCA) for aligning the visual and semantic features to the same latent space to analyze the correlation of visual space and semantic space. We train a CCA model on the visual features and word embeddings of 64 base classes $\mathcal{D}^{base}$ from miniImageNet, where visual features are obtained by the feature extractor trained on base classes. Then we perform the same CCA model on visual features and word embeddings of 16 validation classes $\mathcal{D}^{val}$ and 20 test classes $\mathcal{D}^{test}$ respectively to calculate the correlation coefficient. There is still a relatively high correlation between visual and semantic space on these novel classes while the correlation coefficient is quite small when using the non-corresponding visual and semantic data to train the CCA (results can be found in the supplementary material). We can draw a conclusion that visual space and semantic space are quite relevant and the alignment between them calculated on base classes can be transferred to novel classes.

### 3.2. Problem Formulation

Before testing on novel class, we are given $M$ base classes (denoted as $\mathcal{Y}^b$) for meta-learning purpose. During testing, there are $N$ novel classes (denoted as $\mathcal{Y}^n$) in each few-shot learning task where the base classes and novel classes are disjoint, *i.e.*, $\mathcal{Y}^b \cap \mathcal{Y}^n = \emptyset$. We use the index $\{1, ..., M\}$ to represent the base classes and $\{M + 1, ..., M + N\}$ to represent the novel classes. The base classes dataset (denoted as $\mathcal{D}^{base}$) has plenty of samples per class, while the novel class dataset named sup-
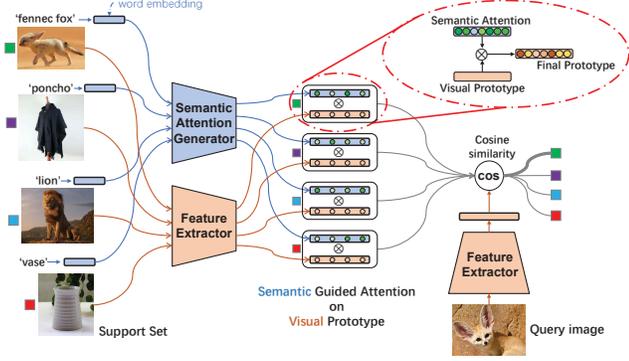
Figure 3: The framework of our proposed Semantic Guided Attention method. We highlight the process of semantic guided attention which is to apply generated semantic attention to visual prototype by Hadamard product. "$\otimes$" denotes the Hadamard product operation and "cos" denotes the cosine classifier.

port set (denoted as $\mathcal{D}^{novel}$) has only $K$ labeled samples per class. The support set contains $N \times K$ labeled images $\mathcal{D}^{novel} = \{(\boldsymbol{x}_i, y_i) \mid \boldsymbol{x}_i \in \mathcal{X}, y_i \in \mathcal{Y}^n\}_{i=1}^{N \times K}$, thus we call it the $N$-Way, $K$-Shot setting. $\mathcal{X} \in \mathbb{R}^{d_v}$ represents the visual space. To derive visual attention from class semantic knowledge, semantic information $\mathcal{S} = \{\boldsymbol{s}^c \in \mathbb{R}^{d_s}\}_{c=1}^{M+N}$ is provided for each class $c \in \mathcal{Y}^b \cup \mathcal{Y}^n$. The goal of FSL is to learn the classifiers for novel classes $f_{fsl} : \mathcal{X} \rightarrow \mathcal{Y}^n$.

## 3.3. Framework

Figure 3 shows the framework of our Semantic Guided Attention (SEGA). It contains three submodules as follows.

**Feature Extractor**. Just like that a human cannot learn a novel category without having seen anything before, the FSL method cannot learn novel class $\mathcal{D}^{novel}$ without the help of base classes $\mathcal{D}^{base}$. The base classes $\mathcal{D}^{base}$ are used to train the *Feature Extractor* in Figure 3. The training procedure follows the classical deep learning paradigm with the common backbone. After training we fix the *Feature Extractor* just like many other works, and now we get the visual space $\mathcal{X} \in \mathbb{R}^{d_v}$. Without doubt, the performance of visual space $\mathcal{X}$ significantly depends on the backbone capacity, the number of classes in $\mathcal{D}^{base}$, and the quality and diversity of samples in $\mathcal{D}^{base}$. So it is worth knowing and keeping all settings the same for experimental study later.

**Cosine Classifier**. Apart from feature extractor, the general deep learning framework consists of a classifier. The standard classifier employs dot-product to calculate classification scores. It is not suitable for FSL setting since the norm of generated weight is not controllable. Therefore during few-shot training we use *Cosine Classifier* proposed by [11, 35] in which the classification score is calculated based on cosine similarity $score_k(\boldsymbol{x}) = \cos \langle \boldsymbol{x}, \boldsymbol{w}_k \rangle = \left\langle \frac{\boldsymbol{x}}{\|\boldsymbol{x}\|}, \frac{\boldsymbol{w}_k}{\|\boldsymbol{w}_k\|} \right\rangle$. With the normalization, the similarity cal-

culation is only based on angle and can generalize well. Besides, the classifier weight is equivalent to class prototype in metric learning when using cosine classifier since the optimization target is mathematically equivalent [35].

**Classification Weight Generator**. The classification weights can be tuned adequately when given plenty of images per class. However, there are not enough samples to tune the classification weight in FSL. As noted before, our method follows the weight generation paradigm which means we will generate the classification weights $\mathcal{W}^{novel} = \{\boldsymbol{w}_c\}_{c=M+1}^{M+N}$ for each novel class. Specifically, our classification weight generator is a semantic guided attention weight generator which will be elaborated later.

**Training Procedure**. To fully mimic human learning process, the training procedure consists of two stages. The first stage is to train the *Feature Extractor* on the base classes $\mathcal{D}^{base}$, which follows the standard classification training paradigm without any semantic knowledge or few-shot concerns. After the first training stage, the *Feature Extractor* will be fixed. The second stage, *i.e.* few-shot training, is of vital importance which contains the training of the *Semantic Guided Attention Weight Generator* and *Cosine Classifier*. Here we adopt the strategy in Dynamic-FSL [11] whose idea is similar with *episodic training* [54]. Compared to *episodic training*, the difference is that our strategy performs the classification task across the whole base classes $\mathcal{Y}^b$ and simulates testing scenario at the same time. More specifically, for each episode, we randomly sample $N$ classes from the base classes $\mathcal{Y}^b$ to act as "novel" classes, then sample $K$ samples from each "novel" class to form a fake $N$-Way $K$-Shot support set. As shown in Figure 3, we can calculate $N$ visual prototypes and enhance them using semantic guided attentions. Thus we get $N$ classification weights which are used to replace the corresponding base classification weights (other weights are also enhanced by their own semantic attentions) in *Cosine Classifier*, and then perform classification and cross-entropy loss calculation.

## 3.4. Semantic Guided Attention Weight Generator

Two key components in our classification weight generation are *Visual Prototype* and *Semantic Guided Attention*. As shown in Figure 3, the semantic guided attention will be applied to enhance the visual prototype and the result will be set as the final classification weight.

**Visual Prototype.** As noted above, the classifier weight is equivalent to the class prototype in metric learning when using the cosine classifier. Hence, we generate the classification weight based on support sets $\mathcal{D}^{novel}$ just like most other metric-based FSL methods. Following the classical Prototypical Network [48], we get the visual prototype:

$$\mathbf{p}_{avg}^c = \frac{1}{|\mathcal{D}_c^n|} \sum_{(\boldsymbol{x_i}, y_i) \in \mathcal{D}_c^n} \boldsymbol{x_i}, \qquad (1)$$

where $\mathcal{D}_c^n \in \mathcal{D}^{novel}$ is the subset of the support set which contains the samples belonging to class $c$. Without doubt, this prototype generation way is too straightforward which ignores the visual prior knowledge that can be transferred from the base class weights $\mathcal{W}^{base} = \{\boldsymbol{w}_c\}_{c=1}^M$. Therefore, we follow our baseline method [11] to transfer the visual prior from base class weights based on the cosine similarity to enhance the averaging prototype:

$$\mathbf{p}_{att}^c = \frac{1}{|\mathcal{D}_c^n|} \sum_{(\boldsymbol{x_i}, y_i) \in \mathcal{D}_c^n} \sum_{j \in \mathcal{Y}^b} Att\left(\boldsymbol{\phi_q}\boldsymbol{x_i}, \boldsymbol{k_j}\right) \cdot \boldsymbol{w}_j, \quad (2)$$

where $\boldsymbol{\phi_q} \in \mathbb{R}^{d_v \times d_v}$ is a learnable weight matrix and $\left\{\boldsymbol{k_j} \in \mathbb{R}^{d_v}\right\}_{j=1}^M$ is a set of $M$ learnable keys. $\boldsymbol{\phi_q}$ transforms the feature $\boldsymbol{x_i}$ to query vector which will be used to perform attention with $\boldsymbol{k_j}$ by a cosine based attention kernel $Att\left(\cdot, \cdot\right)$. By transferring visual prior knowledge from base classes, we model the final visual prototype as the combination of $\mathbf{p}_{avg}^c$ and $\mathbf{p}_{att}^c$:

$$\mathbf{p}^c = \lambda_1 \times \mathbf{p}_{avg}^c + \lambda_2 \times \mathbf{p}_{att}^c, \quad (3)$$

where $\lambda_1, \lambda_2 \in \mathbb{R}$ are learnable coefficients. The final visual prototype will then be feature-selected by semantic guided attention as introduced next.

**Semantic Guided Attention.** The visual prototype is neither precise nor stable on account of the lack of image samples, thus we propose to use the semantic knowledge to guide the attention on the visual prototype in top-down manner as shown in Figure 3. The semantic knowledge can come from the class labels, attributes, and even knowledge graph. Here we choose the word embeddings of class labels as the semantic knowledge source: $\mathcal{S} = \{\boldsymbol{s}^c \in \mathbb{R}^{d_s}\}_{c=1}^{M+N}$, where $\boldsymbol{s}^c$ is the word embedding of the label of class $c$. $d_s$ is the dimension of the word embedding space. We use an MLP to model the transformation $g : \mathbb{R}^{d_s} \rightarrow \mathbb{R}^{d_v}$ which maps the word embedding in semantic space $\mathbb{R}^{d_s}$ to the visual attention in visual space $\mathcal{X}$: $\boldsymbol{a}_c = g(\boldsymbol{s}^c)$. The last layer of $g$ is a sigmoid function, therefore the visual attention $\boldsymbol{a}_c$ is bounded between [0,1]. Actually, $\boldsymbol{a}_c$ can be understood as a feature selection in visual space $\mathcal{X}$ which selects the vital feature dimensions with respect to class $c$ guided by semantic knowledge $\boldsymbol{s}^c$, and the final classification weight for class $c$ is

$$\boldsymbol{w}_c = \boldsymbol{a}_c \otimes \mathbf{p}^c, \quad (4)$$

where $\otimes$ denotes the Hadamard product (*i.e.*, element-wise product operation). It is worth noting that we use *Cosine Classifier* to perform final classification so that the classification score is calculated by

$$\text{score}_c(\boldsymbol{x}) = t \cdot \cos \langle \boldsymbol{x}, \boldsymbol{w}_c \rangle = t \cdot \cos \langle \boldsymbol{x}, \boldsymbol{a}_c \otimes \mathbf{p}^c \rangle, \quad (5)$$

where $t$ is the temperature coefficient to scale the cosine similarity in order to be better suitable for softmax. As we can see, the visual attention $\boldsymbol{a}_c$ here is also playing a role in transferring knowledge since the more similar the visual attention is, the more similar the classification weight will be. Following the same paradigm we can generate all $N$ novel classification weights $\{\boldsymbol{w}_c\}_{c=M+1}^{M+N}$, then we get the desired few-shot learning classifier $f_{fsl}$.

### 3.5. Discussion

**Difference from the baseline.** Our SEGA is mainly inspired by Dynamic-FSL [11] which is a widely-used framework in FSL proposing the cosine classifier and the classification weight generator. However, [11] doesn't involve semantic information which can play a critical role especially when there is a lack of visual experience. In contrast, our SEGA not only utilizes semantic but also works in a new paradigm that is visual-dominated, while previous methods are either semantic-dominated or multimodal-fusion.

**Further improvements.** As we can see, our approach SEGA is orthogonal to most of the current FSL methods since we introduce semantic knowledge which is another dimension to enhance the visual prototype. That means SEGA can be combined with most unimodal metric-based FSL methods to get the prototype more stable and precise, especially in the case of extremely short of images like in 1-Shot learning scenario. On the other hand, semantic knowledge can come from many other sources. With the development of NLP, the semantic guidance imposed on visual features would be more accurate by exploring more powerful knowledge sources such as visual knowledge base, BERT embedding and so on.

## 4. Experiments

In this section, we evaluate our method on four benchmark datasets and then analyze the effectiveness of it.

### 4.1. Datasets and Settings

**Datasets.** We perform experiments on four widely used FSL benchmarks to verify the effectiveness of the proposed method, *i.e.*, miniImageNet [54], tieredImageNet [40], CIFAR-FS [3], and CUB [55]. miniImageNet and tiered-ImageNet are both derivatives of ImageNet dataset [41], CIFAR-FS is derived from CIFAR-100 dataset [19,52]. The datasets summary can be found in supplementary material.

**Semantic knowledge source.** We use GloVe [33] to be our semantic knowledge source which is a word embedding model trained on the Wikipedia dataset and the dimension of the word embedding is 300. We use it to get word embedding for each category. When the category label contains more than one word (*e.g.* "baseball bat"), most previous methods will generate the embedding by averaging them (*e.g.* $\frac{emb(\text{"baseball"}) + emb(\text{"bat"})}{2}$) which is unreasonable. Besides, most previous methods use a 300-dimension

Table 1: Ablation study of our proposed method on miniImageNet, tieredImageNet, and CIFAR-FS. We report the average classification accuracies (%) on 5000 test episodes of novel categories (with 95% confidence intervals). "Sem." denotes whether to use semantic defined in Equation(4) and "FAKE" means using the non-corresponding label as semantic guidance.

| Sem. | miniImageNet | | tieredImageNet | | CIFAR-FS | |
|------|:---:|:---:|:---:|:---:|:---:|:---:|
| | 5Way 1Shot | 10Way 1Shot | 5Way 1Shot | 10Way 1Shot | 5Way 1Shot | 10Way 1Shot |
| YES | **69.04**±0.26(↑**6**) | **52.71**±0.15(↑**6**) | **72.18**±0.30(↑**4**) | **56.82**±0.21(↑**3**) | **76.24**±0.25(↑**8**) | **61.77**±0.17(↑**8**) |
| NO | 62.81±0.27( − ) | 46.73±0.17( − ) | 68.55±0.31( − ) | 54.01±0.21( − ) | 67.78±0.30( − ) | 53.32±0.21( − ) |
| FAKE | 59.04±0.27(↓4) | 43.58±0.16(↓3) | 64.64±0.31(↓4) | 50.07±0.21(↓4) | 63.27±0.29(↓4) | 48.66±0.19(↓5) |

zero vector instead when there is no annotation found in GloVe's vocabulary. Our method uses the hypernym synset based on WordNet to deal with these problems. For those annotations not found in GloVe, we refer to its hypernym synset in WordNet until there is an annotation found in GloVe. By this way, more accurate semantic information can be explored for each category.

**Implementation details.** All experiments are conducted under in PyTorch framework[1]. For all datasets, we utilize a ResNet-12 as our backbone following most previous works [5, 12, 24, 26, 60]. We also change the number of filters from (64,128,256,512) to (64,160,320,640) same as [23, 38, 51]. To avoid overfitting we follow most prior works [16,26,58,64] to adopt the random crop, color jittering, erasing and Dropblock [10] regularization. The semantic guided attention generator used in all cases is an MLP, with 2 fully connected layers and a dropout layer between them, followed by sigmoid nonlinearity. Other parameters $\lambda_1$, $\lambda_2$, and cosine similarity temperature $t$ are tuned during the training of the generator. We use SGD optimizer with a momentum of 0.9 and weight decay of 5e-4. During the first training stage, we train the *Feature Extractor* for 60 epochs (90 for tieredImageNet), with each epoch consisting of 1000 episodes. As for the second training stage, we train *Semantic Guided Attention Weight Generator* and *Cosine Classifier* for 20 epochs in all cases. We adopt an empirical learning rate scheduler following the practice of [11,23,58]. More details can be found in the supplementary material.

## 4.2. Effectiveness of the Proposed Framework

As shown in Table 1, we conduct ablation studies to verify the effectiveness of the proposed semantic guided attention weight generator. By comparing the performance of with semantic (the first row) and without semantic (the second row which is our baseline [11] under our framework), we can infer that the semantic knowledge can significantly improve performance (*i.e.*, the performance improvement is 6%, 4% and 8% on miniImageNet, tieredImageNet, and CIFAR-FS respectively). It is also worth noting that the correct guidance is rather important since when using fake semantic knowledge (the third row), which means using the irrelevant semantic label to generate class attention, the

performance drops even lower than the result without semantic guidance. Furthermore, we conduct experiments on different kinds of pre-trained word embedding models like Word2Vec [27] and see similar phenomena which can be found in the supplementary material.

## 4.3. Dive Deep into Semantic Attention

**What does semantic attention do?** To better understand how the semantic attention works, we perform the t-SNE visualization. Figure 4(a) and 4(b) shows the change of the prototypes before and after applying the semantic attention under 5-Way 1-Shot scenario. As we can see, before the semantic attention the generated prototypes $\mathbf{p}^c$ are quite unstable, but after applying semantic attention the final prototypes $\mathbf{a}_c \otimes \mathbf{p}^c$ become a lot more stable (*i.e.*, the prototypes of the same class get closer and vice versa). It gives the reason why our model can get significant gain under 1-Shot setting. We also show the results when applying inversed semantic guided attention $(1-\mathbf{a}_c)\otimes\mathbf{p}^c$ in Figure 4(c) which means to ignore dimensions that our model thinks is important while emphasizing the unimportant ones. After that, the prototypes get even more unstable and chaotic which further demonstrates that our SEGA does capture the class-specific discriminative dimensions.

**Where does SEGA pay attention to?** Figure 2 shows the Grad-CAM [46] visualization testing on miniImageNet unseen classes based on our model with ResNet-12 backbone. Noted that the model remains exactly the same for all three results columns, the only difference is the attention vector to be applied. By comparing the results of "SEGA" ($\mathbf{w}_c = \mathbf{a}_c \otimes \mathbf{p}^c$), "No_Att" ($\mathbf{w}_c = \mathbf{p}^c$) and "Inv_Att" ($\mathbf{w}_c = (1 - \mathbf{a}_c) \otimes \mathbf{p}^c$), our SEGA can pay attention to the most crucial class-specific feature instead of the misguiding background noise and large intra-class variations. When only given one sample, if we do not utilize the semantic knowledge from the class label(*e.g.*, *bookshop*), even we human can get confused about what this category exactly is (*e.g.*, when seeing a person in the bookshop). Thus, leveraging semantic knowledge from class labels is of vital importance, from which our model knows it should pay attention to the most crucial feature (*e.g.*, books and bookshelf) instead of the noise (*e.g.*, person) in the image.

Furthermore, Figure 5(b) is a harder task where the query is the CutMix [62] of two novel categories. Even in this
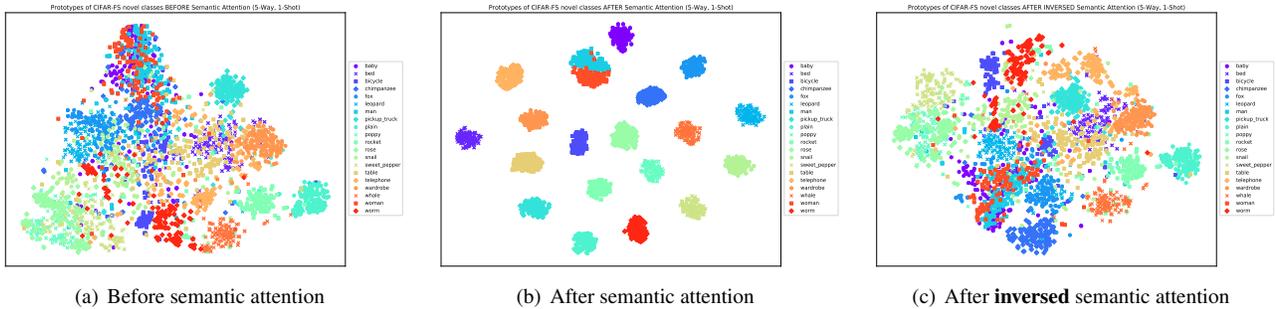
---

[1]The codes are at http://vipl.ict.ac.cn/resources/codes or https://github.com/MartaYang/SEGA

(a) Before semantic attention     (b) After semantic attention     (c) After **inversed** semantic attention

Figure 4: t-SNE visualization of the prototypes in visual space under 5-Way 1-Shot scenario. **(a)** and **(b)** are prototypes before ($\mathbf{p}^c$) and after ($\boldsymbol{a}_c \otimes \mathbf{p}^c$) performing the semantic attention. **(c)** shows the result when applying the inverse attention ($(1 - \boldsymbol{a}_c) \otimes \mathbf{p}^c$). Note that the prototypes are all generated during 600 epochs testing on 20 unseen novel classes on CIFAR-FS (results for miniImageNet can be found in supplementary material) and the point color represents its category. All the prototypes are L2-normalized since we use the cosine classifier.

circumstance, our SEGA can still pay attention to the key part of each corresponding category, which demonstrates the robustness of the attention generated by our model.

More interestingly, since our SEGA can generate class-specific attention, why not apply the intersection of two categories' attentions to get their common ground attributes? Figure 5(a) shows that common attribute "spots" will be highlighted when applying the intersection of attentions of *dalmatian* and *ladybug* (both of them have spots on body), and "long legs" will be highlight when applying *dalmatian* ∩ *saluki* (both of them have long legs). It suggests that our SEGA implicitly establishes a correspondence between semantic knowledge and visual attributes.

**Why does semantic attention work?** Figure 6(a) shows the similarity matrix of the generated attention vectors and their hierarchical clustering tree. It makes sense that visually similar categories cluster together and we can find block-diagonal phenomenon in the similarity matrix (*e.g.*, the attention of baby, man, and woman cluster together and bicycle, rocket, and truck also cluster together). Note that even we do not explicitly exploit WordNet [28] knowledge database, our hierarchical clustering result is quite similar to the ground truth hierarchical structure in WordNet shown in Figure 6(b), which again verifies the effectiveness of SEGA.

### 4.4. Benchmark Comparisons and Evaluations

In this subsection, We make comparison with several popular FSL approaches within the inductive learning framework. Table 2 shows the results on miniImageNet and tieredImageNet dataset. Note that KTN, TriNet, AM3, and our SEGA utilize semantic knowledge while other methods are in unimodal settings. Our method achieves the highest performance especially in 5-Way 1-Shot setting and outperforms the most relevant semantic using method AM3 [60]. Furthermore, we even adapt AM3 to our framework for further fair comparison which can be found in supplementary material. The advantage should be attributed to the more
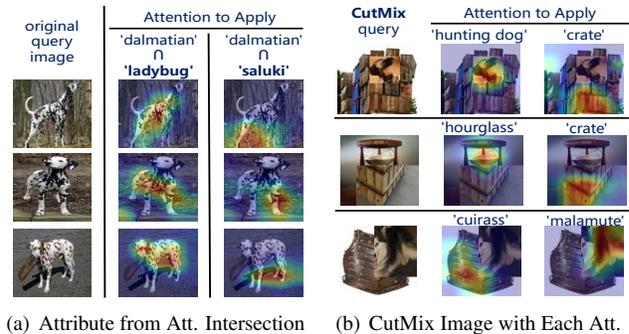


(a) Attribute from Att. Intersection     (b) CutMix Image with Each Att.

Figure 5: **(a)** shows that by applying the intersection of two categories' attentions, the common ground attribute of these two category can be highlighted (*e.g.* "spots" for *dalmatian* ∩ *ladybug* and "long legs" for *dalmatian* ∩ *saluki* ). **(b)** gives the results on the same CutMix query image when applying the attention of each two categories respectively.



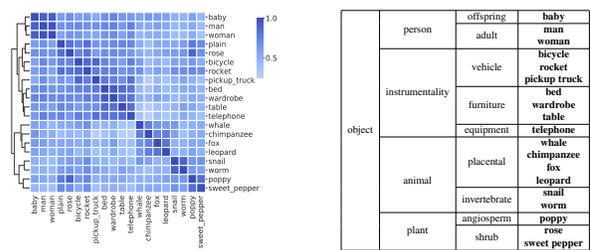(a) Hierarchical-Clustered Heatmap     (b) WordNet Hierarchical Structure

Figure 6: In **(a)**, the attention vector for each class is generated by our attention generation model given novel class names of CIFAR-FS test set. Pearson correlation coefficient and the average-linkage algorithm are used for similarity matrix calculating and hierarchical clustering. **(b)** shows their ground truth hierarchical structure on WordNet, which can be basically matched with the former clustering result.

human-like way to utilize knowledge which is the attention mechanism instead of reconstruction of prototype. Figure 7 shows the performance gain from our SEGA is getting down

Table 2: Comparisons with popular FSL approaches in average classification accuracies (%) on miniImageNet and tieredImageNet. We report the average classification accuracies (%) on 5000 test episodes of novel categories (with 95% confidence intervals). "Sem." denotes whether to leverage semantic knowledge.

| Models | Backbone | Sem. | miniImageNet | | tieredImageNet | |
|---|---|---|---|---|---|---|
| | | | 5Way-1Shot | 5Way-5Shot | 5Way-1Shot | 5Way-5Shot |
| Matching Networks (NIPS'16) [54] | 4Conv | No | 43.56±0.84 | 55.31±0.73 | - | - |
| MAML (ICML'17) [8] | 4Conv | No | 48.70±1.84 | 63.11±0.92 | 51.67±1.81 | 70.30±1.75 |
| ProtoNet (NIPS'17) [48] | 4Conv | No | 49.42±0.78 | 68.20±0.66 | 53.31±0.89 | 72.69±0.74 |
| Dynamic-FSL (CVPR'18) [11] | 4Conv | No | 56.20±0.86 | 72.81±0.62 | - | - |
| Dynamic-FSL (ours baseline) | ResNet-12 | No | 62.81±0.27 | 78.97±0.18 | 68.55±0.31 | 83.95±0.21 |
| wDAE-GNN (CVPR'19) [12] | WRN-28-10 | No | 61.07±0.15 | 76.75±0.11 | 68.18±0.16 | 83.09±0.12 |
| MetaOptNet (CVPR'19) [23] | ResNet-12 | No | 62.64±0.61 | 78.63±0.46 | 65.99±0.72 | 81.56±0.53 |
| DeepEMD (CVPR'20) [64] | ResNet-12 | No | 65.91±0.82 | **82.41**±0.56 | 71.16±0.87 | **86.03**±0.58 |
| RFS (ECCV'20) [51] | ResNet-12 | No | 64.82±0.60 | 82.14±0.43 | 71.52±0.69 | **86.03**±0.49 |
| Neg-Cosine (ECCV'20) [26] | ResNet-12 | No | 63.85±0.81 | 81.57±0.56 | - | - |
| KTN (ICCV'19) [32] | 4Conv | Yes | 64.42±0.72 | 74.16±0.56 | - | - |
| TriNet (TIP'19) [5] | ResNet-18 | Yes | 58.12±1.37 | 76.92±0.69 | - | - |
| AM3 (NIPS'19) [60] | ResNet-12 | Yes | 65.30±0.49 | 78.10±0.36 | 69.08±0.47 | 82.58±0.31 |
| SEGA (ours) | ResNet-12 | Yes | **69.04**±0.26 | 79.03±0.18 | **72.18**±0.30 | 84.28±0.21 |

Table 3: Results on CUB. Test setting is the same as above.

| Models | CUB | |
|---|---|---|
| | 5Way 1Shot | 5Way 5Shot |
| TriNet (TIP'19) [5] | 69.61±0.46 | 84.10±0.35 |
| MultiSem (CoRR'19) [45] | 76.1±n/a | 82.9±n/a |
| FEAT (CVPR'20) [61] | 68.87±0.22 | 82.90±0.15 |
| DeepEMD (CVPR'20) [64] | 75.65±0.83 | 88.69±0.50 |
| SEGA (ours) | **84.57**±0.22 | **90.85**±0.16 |

Table 4: CIFAR-FS results. Test setting is the same as above.

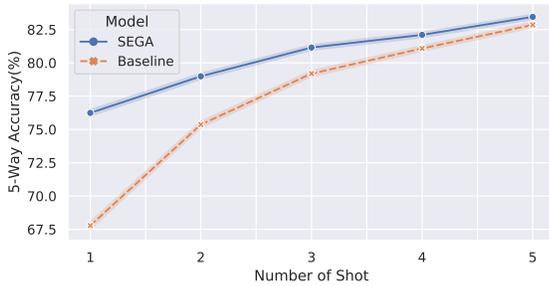| Models | CIFAR-FS | |
|---|---|---|
| | 5Way 1Shot | 5Way 5Shot |
| MAML (ICML'17) [8] | 58.9±1.9 | 71.5±1.0 |
| ProtoNet (NIPS'17) [48] | 55.5±0.7 | 72.0±0.6 |
| MetaOptNet (CVPR'19) [23] | 72.0±0.7 | 84.2±0.5 |
| RFS (ECCV'20) [51] | 73.9±0.8 | **86.9**±0.5 |
| SEGA (ours) | **78.45**±0.24 | 86.00±0.20 |



Figure 7: 5-Way accuracy on CIFAR-FS from 1 to 5 Shot.

when the number of shots goes larger. The reason is that when only given one sample per class the visual prototype is poor and unstable thus the semantic knowledge can help a lot. However, in 5-Shot setting, the visual prototype is getting stable and accurate when given more samples and the gain from semantic information is getting lower (we show the visualization result of 5-Shot in supplementary material which implies that the generated prototypes are already very stable in visual space when given 5 samples). Even though, our performance can still have an advantage over SOTAs in larger shot scenarios when using purer semantic knowledge (e.g. CUB attributes) to guide the attention (as the CUB results shown in Table 3). Results on CIFAR-FS are shown in Table 4, where our method also gets competitive results. We also show our advantage over other SOTAs in computation complexity in the supplementary material.

## 5. Conclusions

In this work, we propose a simple yet effective FSL approach which is accomplished by **SEmantic Guided Attention (SEGA)** on the visual prototype to give top-down guidance on which key features we should focus on. Our proposed approach shows its effectiveness in four popular FSL benchmarks especially when given only one labeled sample for each novel category. Furthermore, we dive deep into how and why our semantic attention works, and further conduct extensive and interesting experiments. Besides, we also analyze the correlation of visual space and semantic space and find out that the alignment calculated on base classes can be transferred and generalized well to novel classes which gives fundamental evidence for the usefulness of the semantic knowledge for FSL task.

# References

[1] Arman Afrasiyabi, Jean-François Lalonde, and Christian Gagné. Associative alignment for few-shot image classification. In *European Conference on Computer Vision (ECCV)*, pages 18–35. Springer, 2020.

[2] David W Aha. *Lazy learning*. Springer Science & Business Media, 2013.

[3] L Bertinetto, J Henriques, PHS Torr, and A Vedaldi. Meta-learning with differentiable closed-form solvers. International Conference on Learning Representations (ICLR), 2019.

[4] Riquan Chen, Tianshui Chen, Xiaolu Hui, Hefeng Wu, Guanbin Li, and Liang Lin. Knowledge graph transfer network for few-shot recognition. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 10575–10582, 2020.

[5] Zitian Chen, Yanwei Fu, Yinda Zhang, Yu-Gang Jiang, Xiangyang Xue, and Leonid Sigal. Multi-level semantic feature augmentation for one-shot learning. *IEEE Transactions on Image Processing (TIP)*, 28(9):4594–4605, 2019.

[6] Robert Epstein. The empty brain. *Aeon, May*, 18:2016, 2016.

[7] Ali Farhadi, Ian Endres, Derek Hoiem, and David Forsyth. Describing objects by their attributes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1778–1785. IEEE, 2009.

[8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning (ICML)*, pages 1126–1135. PMLR, 2017.

[9] Hang Gao, Zheng Shou, Alireza Zareian, Hanwang Zhang, and Shih-Fu Chang. Low-shot learning via covariance-preserving adversarial augmentation networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 975–985, 2018.

[10] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. *Advances in Neural Information Processing Systems (NeurIPS)*, 31:10727–10737, 2018.

[11] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4367–4375, 2018.

[12] Spyros Gidaris and Nikos Komodakis. Generating classification weights with gnn denoising autoencoders for few-shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21–30, 2019.

[13] Jacob Goldberger, Geoffrey E Hinton, Sam Roweis, and Russ R Salakhutdinov. Neighbourhood components analysis. *Advances in Neural Information Processing Systems (NeurIPS)*, 17:513–520, 2004.

[14] Jiechao Guan, Zhiwu Lu, Tao Xiang, Aoxue Li, An Zhao, and Ji-Rong Wen. Zero and few shot learning with semantic feature synthesis and competitive learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.

[15] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9729–9738, 2020.

[16] Ruibing Hou, Hong Chang, Bingpeng Ma, Shiguang Shan, and Xilin Chen. Cross attention network for few-shot classification. *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

[17] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 831–839, 2019.

[18] Gregory Koch, Richard Zemel, and Ruslan Salakhutdinov. Siamese neural networks for one-shot image recognition. In *International Conference on Machine Learning workshop (ICMLW)*, volume 2. Lille, 2015.

[19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[20] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

[21] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Learning to detect unseen object classes by between-class attribute transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 951–958. IEEE, 2009.

[22] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 36(3):453–465, 2013.

[23] Kwonjoon Lee, Subhransu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10657–10665, 2019.

[24] Aoxue Li, Weiran Huang, Xu Lan, Jiashi Feng, Zhenguo Li, and Liwei Wang. Boosting few-shot learning with adaptive margin loss. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12576–12584, 2020.

[25] Aoxue Li, Tiange Luo, Zhiwu Lu, Tao Xiang, and Liwei Wang. Large-scale few-shot learning: Knowledge transfer with class hierarchy. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7212–7220, 2019.

[26] Bin Liu, Yue Cao, Yutong Lin, Qi Li, Zheng Zhang, Mingsheng Long, and Han Hu. Negative margin matters: Understanding margin in few-shot classification. In *European Conference on Computer Vision (ECCV)*, pages 438–455. Springer, 2020.

[27] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

[28] George A Miller. *WordNet: An electronic lexical database*. MIT press, 1998.

[29] Tom Mitchell, William Cohen, Estevam Hruschka, Partha Talukdar, Bishan Yang, Justin Betteridge, Andrew Carlson, Bhavana Dalvi, Matt Gardner, Bryan Kisiel, et al. Never-ending learning. *Communications of the ACM*, 61(5):103–115, 2018.

[30] Maximilian Nickel and Douwe Kiela. Poincar\'e embeddings for learning hierarchical representations. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[31] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[32] Zhimao Peng, Zechao Li, Junge Zhang, Yan Li, Guo-Jun Qi, and Jinhui Tang. Few-shot image recognition with knowledge transfer. In *IEEE International Conference on Computer Vision (ICCV)*, pages 441–449, 2019.

[33] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.

[34] Guo-Jun Qi and Jiebo Luo. Small data challenges in big data era: A survey of recent progress on unsupervised and semi-supervised methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.

[35] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5822–5830, 2018.

[36] Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting parameters from activations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7229–7238, 2018.

[37] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. 2016.

[38] Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Few-shot learning with embedded class models and shot-free meta training. In *IEEE International Conference on Computer Vision (ICCV)*, pages 331–339, 2019.

[39] William J Reed. The pareto, zipf and other power laws. *Economics letters*, 74(1):15–19, 2001.

[40] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. In *International Conference on Learning Representations (ICLR)*, 2018.

[41] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.

[42] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *International Conference on Learning Representations (ICLR)*, 2018.

[43] Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.

[44] Edgar Schonfeld, Sayna Ebrahimi, Samarth Sinha, Trevor Darrell, and Zeynep Akata. Generalized zero-and few-shot learning via aligned variational autoencoders. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8247–8255, 2019.

[45] Eli Schwartz, Leonid Karlinsky, Rogerio Feris, Raja Giryes, and Alex M Bronstein. Baby steps towards few-shot learning with multiple semantics. *arXiv preprint arXiv:1906.01905*, 2019.

[46] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017.

[47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

[48] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4077–4087, 2017.

[49] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1199–1208, 2018.

[50] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1–9, 2015.

[51] 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 (ECCV)*, pages 266–282. Springer, 2020.

[52] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 30(11):1958–1970, 2008.

[53] Shimon Ullman. *High-level vision: Object recognition and visual cognition*, volume 2. MIT press Cambridge, MA, 1996.

[54] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3630–3638, 2016.

[55] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.

[56] Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E Gonzalez. Tafe-net: Task-aware feature embeddings for low shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1831–1840, 2019.

[57] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)*, 53(3):1–34, 2020.

[58] Zeyuan Wang, Yifan Zhao, Jia Li, and Yonghong Tian. Co-operative bi-path metric for few-shot learning. In *ACM Inter-*

*national Conference on Multimedia (ACMMM)*, pages 1524–1532, 2020.

[59] Kilian Q Weinberger and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10(2), 2009.

[60] Chen Xing, Negar Rostamzadeh, Boris Oreshkin, and Pedro O O Pinheiro. Adaptive cross-modal few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 32:4847–4857, 2019.

[61] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8808–8817, 2020.

[62] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *IEEE International Conference on Computer Vision (ICCV)*, pages 6023–6032, 2019.

[63] Baoquan Zhang, Xutao Li, Yunming Ye, Zhichao Huang, and Lisai Zhang. Prototype completion with primitive knowledge for few-shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3754–3762, 2021.

[64] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12203–12213, 2020.