# Recognizing Unseen Objects via Multimodal Intensive Knowledge Graph Propagation

Likang Wu[*]
School of Computer Science and Technology, University of Science and Technology of China & State Key Laboratory of Cognitive Intelligence
Hefei, Anhui, China
wulk@mail.ustc.edu.cn

Zhi Li[*]
Shenzhen International Graduate School, Tsinghua University
Shenzhen, Guangdong, China
zhilizl@sz.tsinghua.edu.cn

Hongke Zhao
College of Management and Economics, Tianjin University
Tianjin, China
hongke@tju.edu.cn

Zhefeng Wang
Huawei Cloud
Hangzhou, Zhejiang, China
wangzhefeng@huawei.com

Qi Liu
School of Computer Science and Technology, University of Science and Technology of China & State Key Laboratory of Cognitive Intelligence
Hefei, Anhui, China
qiliuql@ustc.edu.cn

Baoxing Huai
Huawei Cloud
Hangzhou, Zhejiang, China
huaibaoxing@huawei.com

Nicholas Jing Yuan
Huawei Cloud
Hangzhou, Zhejiang, China
nicholas.jing.yuan@gmail.com

Enhong Chen[†]
School of Computer Science and Technology, University of Science and Technology of China & State Key Laboratory of Cognitive Intelligence
Hefei, Anhui, China
cheneh@ustc.edu.cn

## ABSTRACT

Zero-Shot Learning (ZSL), which aims at automatically recognizing unseen objects, is a promising learning paradigm to understand new real-world knowledge for machines continuously. Recently, the Knowledge Graph (KG) has been proven as an effective scheme for handling the zero-shot task with large-scale and non-attribute data. Prior studies always embed relationships of seen and unseen objects into visual information from existing knowledge graphs to promote the cognitive ability of the unseen data. Actually, real-world knowledge is naturally formed by multimodal facts. Compared with ordinary structural knowledge from a graph perspective, multimodal KG can provide cognitive systems with fine-grained knowledge. For example, the text description and visual content can depict more critical details of a fact than only depending on knowledge triplets. Unfortunately, this multimodal fine-grained knowledge is largely unexploited due to the bottleneck of feature alignment between different modalities. To that end, we propose a multimodal intensive ZSL framework that matches regions of images with corresponding semantic embeddings via a designed dense attention module and self-calibration loss. It makes the semantic transfer process of our ZSL framework learns more differentiated knowledge between entities. Our model also gets rid of the performance limitation of only using rough global features. We conduct extensive experiments and evaluate our model on large-scale real-world data. The experimental results clearly demonstrate the effectiveness of the proposed model in standard zero-shot classification tasks.

## CCS CONCEPTS

• **Information systems** → **Multimedia information systems**; **Data mining**; • **Computing methodologies** → **Knowledge representation and reasoning**.

## KEYWORDS

Knowledge Graph, Multimodal Data, Zero-Shot Learning, Graph Neural Networks

[*]Equal Contributions.
[†]Corresponding Author.

# 1 INTRODUCTION

Zero-Shot Learning (ZSL), which utilizes prior knowledge to recognize unseen objects in the machine learning process [45, 46], becomes an important task to test the machine's cognitive ability. The popular approaches mainly focus on learning the semantic correlation between seen and unseen categories by constructing a unified space from manually annotated attributes [24, 47]. Although these methods have achieved promising performances in ZSL tasks, it is too difficult and time-consuming to annotate elaborate attributes for large-scale data in the real world. Recently, several researchers begin to cope with this challenge by incorporating the Knowledge Graphs (KG) to generate transferable knowledge and promote the cognitive ability of the ZSL models [21, 39].

Although existing KG-based ZSL methods obtain great achievements in the task with large-scale data, they only exploit the structural knowledge to generate the representations of unseen objects. When identifying similar categories, they often fail to pay attention to key distinguishing features, resulting in classification errors. Actually, for this problem, multimodal knowledge (e.g. concept words, text descriptions, referred images, etc.) can help machines better understand the real world in a more fine-grained way and also promote the cognitive ability of ZSL. For example, Figure 1 illustrates a toy example of multimodal knowledge graph in the ZSL task. Now "Horse" and "White Tiger" are seen classes for the machine while "Zebra" is unseen. Then, with structural knowledge, we teach the machine that "Zebra" is a subgenus of "Horse". With text descriptions, we can tell the machine that "Zebra" has black and white stripes same as "White Tiger" to help the machine adaptively capture the fine-grained semantics. Therefore, the machine can automatically recognize "Zebra" as a horse-like animal with black and white stripes after seeing the "Horse" and "White Tiger" images. From this example, we can conclude that multimodal knowledge can help machines promote cognitive ability for the real world furthermore achieve better ZSL performances.

Unfortunately, the exploration of multimodal knowledge in ZSL is largely limited by great challenges, which leads to the lack of mature paradigms for solving ZSL problems via multimodal knowledge graphs. First, it is hard to organize the multimodal knowledge and generate representations for the heterogeneous multimodal data. Second, the semantic associations of different modal knowledge are not easy to capture and construct. For example, the text descriptions have more fine-grained information ("black and white stripes") than concept words as shown in Figure 1. Then, adaptively exploiting fine-grained information to enhance corresponding features in other modalities is significantly difficult, which leads to the dilemma of transferring multimodal knowledge from seen objects to unseen objects. And we know that knowledge transfer is the key to promoting ZSL learning.

To solve the dilemma mentioned above, in this paper, we present a focused study on the zero-shot recognition framework from a novel multimodal knowledge intensive perspective. Specifically, we propose a novel Fine-grained Graph Propagation (FGP) model to deeply exploit the fine-grained multimodal knowledge and transfer it from seen to unseen objects. We first introduce a semantic space to adaptively generate the various fine-grained semantic embeddings from the text descriptions for each concept. Then, a fine-grained
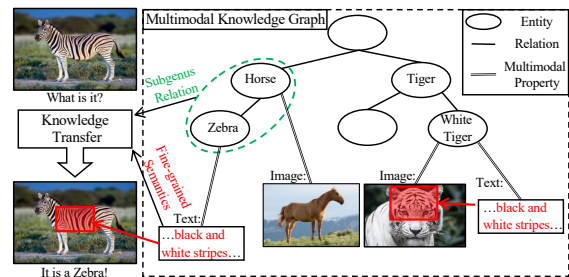


**Figure 1: The knowledge transfer process of typical zero-shot learning with the fine-grained knowledge alignment mechanism on multimodal data.**

visual knowledge learning mechanism is designed to prompt the image classifier to understand more knowledge from both visual and textual modalities, which adaptively maintains the fine-grained features aligned in the visual-semantic space simultaneously. Next, to transfer the multimodal knowledge from seen to unseen objects, we propose a Multi-facet Graph Convolutional Network (Multi-GCN) to learn unified fine-grained class representations for both the seen and unseen objects from the structural knowledge. Finally, we conduct sufficient experiments in real-world large-scale datasets. Experimental results and visualization cases clearly demonstrate the effectiveness of our proposed FGP. In summary, the main contributions could be summarized as follows:

- We present the zero-shot recognition framework from a novel multimodal knowledge intensive perspective.
- We propose a Fine-grained Graph Propagation (FGP) model that deeply exploits the fine-grained multimodal knowledge and transfers it from seen to unseen objects. FGP breaks through the restriction that requires annotation attributes for ensuring visual-semantic alignment.
- We conduct sufficient experiments on real-world large-scale datasets, and the experimental results show the superiority of our model compared with all baselines.

# 2 RELATED WORK

The two most relevant technical points of this paper are zero-shot learning and graph neural networks. We introduce the related research work from these two aspects respectively.

## 2.1 Zero-Shot Learning

We will review zero-shot learning from two perspectives: embedding-based methods and KG-based methods. The motivation of embedding models is to represent all classes by trained vector representations that can be mapped to visual classifiers [4, 6, 11, 12, 51]. [35] trains two unsupervised neural networks for visual and text inputs and then learn a linear mapping between captured image representations and word embedding. After that, [11] proposes a novel framework DeViSE which trains a mapping relation from image to word embeddings via a ConvNet and a transformation layer. Instead of predicting the word embedding directly, [31] designs ConSE to construct the image embedding, which combines the image classification ConvNet and the word embedding model at the same time. In

recent years, many studies implement the attention mechanism to learn fine-grained features to enhance the discriminative power of embedding [19, 20, 47, 48]. [19] propose an interesting feature composition model that captures attribute-based visual features from training instances and combines them to construct fine-grained features for unseen classes. The multimodal model [27] incorporates a joint visual-attribute attention module and a multi-channel explanation module to try explainable zero-shot recognition. Besides, with the help of ontology-based knowledge representation and semantic information, several recent studies explore richer and more competitive prior knowledge to model the inter-class relationship for knowledge transfer of ZSL [13, 14].

To handle large non-attribute data with fine-grained features, we concentrate on the popular fashion of using multimodal knowledge graphs [5] (explicit knowledge representations) for this shortage. Researchers have proposed a unified pipeline on how to use knowledge graphs for ZSL object recognition [21, 39]. This idea has been early used in the work [21, 39], where a graph convolutional neural network was trained to predict logistic regression classifiers on top of a pre-trained CNN to distinguish unseen classes. Later, HVEL learned the implicit knowledge and explicit knowledge in the hyperbolic space [26], which better captured the class hierarchy with fewer dimensions. [38] exploits multiple relationships among different categories of KG for zero-shot learning by employing graph convolutional representation and contrastive learning techniques. [49] visually divides a set of images from seen classes into clusters of local image regions according to their visual similarity, and further imposes their class discrimination and semantic relatedness. MGFN [41] develops a multi-granularity fusion network that integrates discriminative information from multiple GCN branches. In other related applications, such as action recognition, DARK [28] trains a factorized model by extracting disentangled feature representations for verbs and nouns and then predicting classification weights using relations in external knowledge graphs. However, there is no effective KG-based method utilizing multimodal features well, since it is challenging to ensure the alignment of information between different modalities without fine-grained annotations. In this paper, we present a novel view to exploit this multimodal knowledge in the KG-based ZSL.

## 2.2 Graph Neural Networks

We just utilize graph neural networks as an efficient tool to handle our knowledge aggregation process in our work, so we only summarize some of the most classic representation learning methods. Graph neural networks (GNNs) [15, 34, 42], especially graph convolutional networks [18], have activated many researchers' motivation in structured data modeling and analysis because of their theory simplicity and model efficiency [2, 7]. They break hard performance bottlenecks in many research areas, such as node classification [23, 44], graph classification [8], and recommendation [33, 52]. The graph spectral theory is first used to derive a graph convolutional layer [18]. Then, the polynomial spectral filters are proposed to greatly reduce the computational cost [8]. [23] proposes the implementation of a linear filter to make further calculation simplification. Along with spectral methods, directly conducting graph convolution in the spatial domain is also investigated by so many

studies [10, 16]. Among them, graph attention network as a representative method has shown the greatest research potential [37], which adaptively specifies weights to the neighbors of a node by attention mechanism [1, 43]. For the GNN-based representation in zero-shot tasks, [40] proposes a GCN-based model named DGPN which follows the principles of locality and compositionality of zero-shot model generalization. [25] optimizes the learned attribute space for ZSL by training a propagation mechanism to refine the semantic attributes according to the neighbors and related classes. And [50] explores a graph construction approach to flexibly create useful category graphs via leveraging diverse correlations between category nodes.

## 3 METHODOLOGY

In this section, we first illustrate our research problem formally and present task-related notations. Along the processing line, we would introduce our Fine-grained Graph Propagation (FGP) framework in detail. Figure 2 illustrates the overall framework of FGP.

### 3.1 Task-Related Notations

We focus on the classic zero-shot learning problem. We use $C, C_{tr}, C_{te}$ to denote the set of all classes, training classes, and test classes, respectively. Note that, here $C_{tr} \cap C_{te} = \emptyset$. Then we define $\mathcal{D}_{tr} = \{(\mathbf{X}_i, c_i)|i = 1, 2, ..., |C_{tr}|\}$ as the training image dataset, where $\mathbf{X}_i$ indicates the $i-$th training image and $c_i \in C_{tr}$ the corresponding class label. Similarly, we get all image set $\mathcal{D}$ and the test image set $\mathcal{D}_{te}$. All classes in $C$ are connected as nodes in a knowledge graph $G$, where each node $i$ is represented by a category name, corresponding text description $T_i$, and images belong to this class. The symmetric adjacency matrix $A \in \mathbb{R}^{C \times C}$ denotes the connections between the classes in $G$, where $A_{i,j} = 1$ means two nodes $i$ and $j$ connected and vice versa. In this setting, zero-shot recognition aims to predict the class labels of a set of test data points to the set of classes $C_{te}$. Unlike traditional classification, the test data points have to be assigned to previously unseen classes.

### 3.2 Semantic Embeddings Generation

Considering the human cognitive process, a person only needs to know two prior conditions if he intends to recognize the unseen object zebra: first, he has seen the horse or related horse photos as well as black and white stripe objects; second, this person has been told the zebra is a horse-like animal with black and white stripes. Here the semantic information of black and white stripes is the key point for humans or machines to identify zebra from various horses, which is also a typical fine-grained feature in this task. Then, how understanding the text description "black and white stripes" is one of the major tasks to simulate the human cognitive process for ZSL. To achieve this goal, in the first step, we design a Semantic Embeddings Generation (SEG) module to extract key semantic representations from text descriptions. It should be noted that there is a pretty necessary requirement to match the key semantics (such as stripe characteristic) with corresponding certain parts in the visual space (images, videos, or reality), which is the so-called (fine-grained) visual-semantic features alignment. In many specific datasets, the semantic alignment has been annotated with well-designed boolean attributes [24] (each sub-attribute is
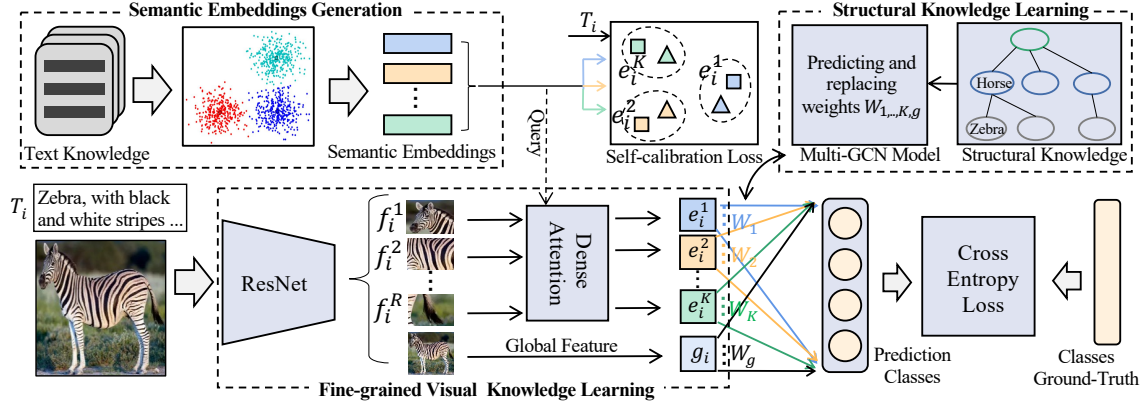
**Figure 2: The overview of our zero-shot recognition framework, i.e., Fine-grained Graph Propagation (FGP). In detail, the two panels on the top indicate the semantic embedding generation process and structural knowledge learning, respectively. The bottom panel denotes our fine-grained knowledge extractor and classifier.**

supposed to pay attention to a distinctive part). However, in the real-world dataset, the annotation for large-scale objects is unaffordable and not realistic. In order to achieve the alignment of the visual-semantic features in large-scale real-world datasets, we propose an unsupervised Semantic Embeddings Generation (SEG) module to extract the key semantic representations from text descriptions in the first step, which would benefit the following fine-grained visual-semantic alignment.

In the SEG module, we propose to utilize the unsupervised clustering algorithm to deal with the unlabeled text corpus that comes from the multimodal KG. As mentioned before, given the description $T_i$ of each concept $i \in C$, we design and adopt several necessary preprocessing steps on the original text descriptions, including:

Step 1. Tokenization, lemmatization, converting uppercase to lowercase, and deleting stop words.

Step 2. Using the dependency parsing tree in spacy tool[1] to filter short phrases (nouns with modifiers, or modifiers, eg., adj + noun, noun, adj, etc) from sentences, such as "white stripe".

Step 3. Using the pre-trained model Glove [32] to transfer all phrases into embedding vectors, where the embedding $\mathbf{p}_j$ of a phrase is produced by the mean pooling of all word embeddings that belong to it.

For a fair comparison, we obtain word embeddings by the Glove model the same as DGP's approach [21] rather than other more advanced text models. After the preprocessing, we utilize the Kmeans clustering algorithm [29] to cluster all phrase embeddings into $K$ semantic groups. Along this line, we capture the centroids set:

$$\Phi_v = \{\mathbf{v}^{c,k} \in \mathbb{R}^{d_c} | k = 1, 2, ..., K\}, \tag{1}$$

where these centroid $\mathbf{v}^{c,k}$ of obtained $K$ clusters are regarded as our key semantic embeddings, which are expected to have the capacity to concentrate on different and representative fine-grained attributes in the semantic space. From the view of machine understanding, the semantics of these key clusters may be abstract and do not necessarily correspond to the feature attributes concerned

by human experts. We will explore the performance influence of the number of clusters in our experimental section.

## 3.3 Fine-Grained Visual Knowledge Learning

In this section, with the availability of original visual inputs and generated key semantic embeddings, we implement a dense attention module to learn the fine-grained visual knowledge, which ensures the alignment of visual-semantic space at the same time.

The visual attention mechanism is able to generate a feature from the most relevant region of an image to match the key embedding and has been shown to be effective for image classification. Follow the setting of [20], each image in $\mathbf{X}$ is divided into $R$ equal-size grid cells (regions) which are denoted by $\{I^r | r \in R\}$. Here we use a ResNet-50 [17] pretrained on ImageNet [9] to catch the feature vector $\mathbf{f}^r \in \mathbb{R}^{d_f} = f(I_r)$ of region $r$, where $f(\cdot)$ denotes the ResNet-50 network. For a sample, we define the query value of the attention mechanism by $\mathbf{v}^c$ and key values by $\{\mathbf{f}^r | r \in R\}$, then the weight of each query on keys can be calculated as Eq. (2),

$$\alpha_k^r = \frac{\exp(\mathbf{v}^{c,k}\mathbf{w}_\alpha \mathbf{f}^r)}{\exp(\sum_{r'=1}^{R} \mathbf{v}^{c,k}\mathbf{w}_\alpha \mathbf{f}^{r'})}, \tag{2}$$

where $\mathbf{w}_\alpha \in \mathbb{R}^{d_c \times d_f}$ is the parameter to expand the dimension of $\mathbf{v}^{c,k}$ to be equal to the dimension of visual feature $\mathbf{f}^r$. In the weight distribution $\alpha_k = \{\alpha_k^1, \alpha_k^2, ..., \alpha_k^R\}$ of key semantic embedding $\mathbf{v}^{c,k}$, $\alpha_k^r \in [0, 1]$ and $\sum_{r=1}^{R} \alpha_k^r = 1$ means that our model can select different regions according to diverse weight values. Hence we carry out the selection procedure easily and get the output as follows:

$$\mathbf{e_k} = \sum_{r=1}^{R} \alpha_k^r \mathbf{f}^r. \tag{3}$$

Now we obtain the attention-based representation. Each state $\mathbf{e}_k$ contains the enhanced certain region which focuses on different topics. That is to say, our method possesses the ability of extracting fine-grained features from rough information. We aim to further ensure that the key semantic embedding recognizes the required key visual object which corresponds to the description of text data.

---

[1]https://spacy.io

To achieve this goal, we design a self-calibration loss that restrains the distance between the fine-grained feature and corresponding semantic embedding. In detail, we assume that the representations' correlation of focused regions from different images that belong to the same key semantic embedding should be relatively close compared with other matched pair types. For example, if the semantic key concentrates on the "**yellow body**" related topic, the distance to this key embedding from text embedding of class **Lion** and **Orange Cat** should be closer than that from **Bear** or **Panda**, and so on. Along this line, for a class $c$, we first calculate the cosine distances $dis_c \in \mathbb{R}^K$ between its text embedding and our key semantic embeddings. Regarding $dis_c$ as a new label, we just need to approach the similarity of each fine-grained feature $\{\mathbf{e}_k | k \in K\}$ (for convenience, here we omit the sample id $i$ of $\mathbf{e}_k$ obtained from the sample with label $c$) and $dis_c$ constantly for each class $c$.

$$\mathcal{L}_{sc} = \frac{1}{2K} \sum_{k=1}^{K} \| \text{ReLU}(\mathbf{e}_k \mathbf{w}_{sc}) - dis_c{}^k \|^2 . \tag{4}$$

The self-calibration loss $\mathcal{L}_{sc}$ is shown in Eq. (4), where $\mathbf{w}_{sc} \in \mathbb{R}^{d_f \times 1}$ denotes the training mapping parameter. Actually, with the optimization iterations, the similar fine-grained regions in various images would be aligned to the correlative key semantic embedding step by step. In this way, our model gradually understands what kind of visual feature is "yellow body". Having got this far, more than half of the motivation for fine-grained information alignment has been achieved. Note that the optimization of all parameters in this section would be executed in a fine-tuning process.

We will freeze the parameters of ResNet-50 and fine-tune our model on the ImageNet-2012 dataset. Specifically, for $\{\mathbf{e}_k | k \in K\}$, we implement individual fully-connected layer $\mathbf{w}_f^{k\top} \in \mathbb{R}^{d_f \times |C_{tr}|}$ as the classifier for each of them. So the final output is:

$$\mathbf{o} = \mathbf{e}_g \mathbf{w}_g^\top + \mathbf{b_g} + \sum_{k=1}^{K} \mathbf{e}_k \mathbf{w}_f^{k\top} + \mathbf{b}_k, \tag{5}$$

where $\mathbf{e}_g$ and $\mathbf{w}_g$ denote the global feature extracted by ResNet-50 and corresponding global classifier, respectively. The final cross-entropy loss of our zero-shot classification task is defined as follows:

$$\mathcal{L}_{ce} = -\frac{1}{|\mathcal{D}_{tr}|} \sum_{i \in \mathcal{D}_{tr}} \log \frac{\exp(\mathbf{o}_i^{c_i})}{\sum_{c' \in C_{tr}} \exp(\mathbf{o}_i^{c'})}, \tag{6}$$

where $c_i$ denotes the true class label of sample $i$. In the fine-tuning stage, we optimize the loss $\mathcal{L}_{sc}$ for all samples and $\mathcal{L}_{ce}$ simultaneously until they converge to stable status.

## 3.4 Structural Knowledge Learning

To learn the relationships between unseen and seen classes, we draw support from the structural association in the knowledge graph by GCN as well as [21]. It is a semi-supervised message propagation approach that aggregates implicit and explicit correlations of nodes (classes) into the node representations. In detail, we input the complete text embedding of each class as the initial state of GCN to predict the last layer weight (classifier) of CNN (e.g., ResNet) and then replace the original weight with obtained prediction vector. Because the GCN network could also learn the projection weights for unseen classes according to its semi-supervised training style,
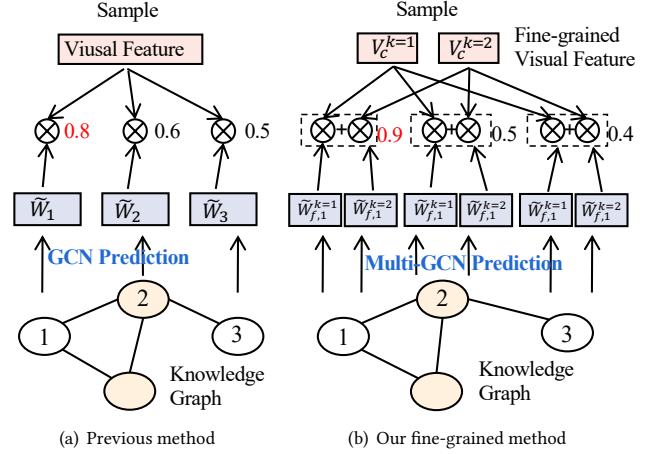


**Figure 3: The difference between previous rough GCN prediction and our fine-grained aligned prediction. On the knowledge graph, the orange nodes denote seen classes with labels and the white nodes denote unseen classes.**

the classification framework obtains the ability to recognize unseen concepts after a fine-tuning of the CNN's parameters (the new classifiers learned by GCN are frozen in the fine-tuning process).

Different from previous KG-based models, in the weight fitting step, we have to maintain different classifiers for several fine-grained features with different focused topics. It is obviously inappropriate to use a single GCN with the same state to predict the weight parameters of all classifiers. Meanwhile, we still cannot throw away the visual-semantic alignment condition at this stage. To this end, we propose an effective but not complicated Multi-GCN network on the knowledge graph to fit our target parameters. We define each unique GCN of Multi-GCN as:

$$\mathbf{H}_k^{l+1} = \sigma(D^{-1} A \mathbf{H}_k^l \theta_k^l), k = 1, 2, ..., K, g, \tag{7}$$

where the symmetric adjacency matrix $A \in \mathbb{R}^{C \times C}$ denotes the connections between the classes in the knowledge graph, which also includes self-loops. $D_{i,i} = \sum_j A_{i,j}$ is a degree matrix that normalizes rows in $A$ to ensure that the value scale of feature representations is not modified by $A$. $\theta_k^l \in \mathbb{R}^{d_c \times F}$ denotes the trainable weight matrix for layer $l$ with $F$ corresponding to the number of learned filters. $\sigma$ indicates the nonlinear activation function Leaky ReLU. Besides, $\mathbf{H}_k^l$ represents the $l$-th layer's hidden state for classifiers in the $k$-th fine-grained channel. Here for each node $i$, the initial node embedding is defined as follows:

$$\mathbf{H}_{k,i}^0 = \text{MeanPooling}(\mathbf{P}_{k,i}), \tag{8}$$

$$\mathbf{P}_{k,i} = \{\mathbf{p}_j | dis_c(\mathbf{p}_j, \mathbf{v}^{c,k}) \le \forall dis_c(\mathbf{p}_j, \mathbf{v}^{c,\hat{k} \ne k}),$$
$$\hat{k} \le K, \mathbf{p}_j \in T_i\}, \tag{9}$$

where $dis_c(\cdot)$ is also the cosine distance to calculate vectors' similarity, and $\mathbf{P}_{k,i}$ denotes all Glove-based phrase embeddings in the text data $T_i$ that more close to the key semantic embedding $\mathbf{v}^{c,k}$ than others ($\{\mathbf{v}^{c,\hat{k}} | \hat{k} \ne k\}$). In the condition mentioned above, the

---

**Algorithm 1:** Fine-grained Graph Propagation

---

**input** : The set of all classes $C$, training classes $C_{tr}$, and test classes $C_{te}$. The training image dataset $\mathcal{D}_{tr}$, test image set $\mathcal{D}_{te}$. The knowledge graph $G$, the text description $T_i$ of each class $i$

**output**: Prediction values for test dataset

1 Get pre-trained ResNet-50 backbone $\Theta_1$;

2 Randomly initialize attention parameter $\Theta_2$, classification layer $\{\mathbf{w}_f, \mathbf{w}_g\}$, Multi-GCN parameter $\Theta_3$;

3 Preprocess each original text $T_i$ into phrase embeddings;

4 Generate $K$ semantic clusters and centroids $\mathbf{v}^c$ by Kmeans;

5 Freeze ResNet-50 parameter $\Theta_1$ and fine-tune parameters $\{\Theta_2, \mathbf{w}_f, \mathbf{w}_g\}$ with query $\mathbf{v}^c$ on $\mathcal{D}_{tr}$ in the fine-grained visual-semantic alignment process;

6 Fit classifier weights $\{\mathbf{w}_f, \mathbf{w}_g\}$ by the output matrix $\{\widetilde{\mathbf{w}}_f, \widetilde{\mathbf{w}}_g\}$ of Multi-GCN $\Theta_3$ trained with phrase embedding inputs of $C_{tr}$ on $G$;

7 Replace $\{\mathbf{w}_f, \mathbf{w}_g\}$ with $\{\widetilde{\mathbf{w}}_f, \widetilde{\mathbf{w}}_g\}$;

8 Fine-tune $\Theta_1$ on $\mathcal{D}_{tr}$, then output classification values $\widetilde{\mathbf{y}}_{te}$ for all samples in $\mathcal{D}_{te}$;

9 **return** $\widetilde{\mathbf{y}}_{te}$;

---

prediction flow based on the knowledge graph is aligned with the fine-grained level knowledge, which further avoids the information interference of mistaken visual-text (semantic) matchings.

In the loss function $\mathcal{L}_{gcn}$, the Multi-GCN model is trained to predict the classifier weights for the seen classes by optimizing $\mathcal{L}_{gcn}$, where each $\widetilde{\mathbf{w}} \in \mathbb{R}^{C_{tr} \times d_f}$ is the prediction of the GCN for the known classes and corresponds to the $C_{tr}$ rows of the GCN output.

$$\mathcal{L}_{gcn} = \frac{1}{2C_{tr}(K+1)} \Big[ \sum_{i=1}^{C_{tr}} \sum_{j=1}^{d_f} (\mathbf{w}_{g,i,j} - \widetilde{\mathbf{w}}_{g,i,j}) + \\ \sum_{k=1}^{K} \sum_{i=1}^{C_{tr}} \sum_{j=1}^{d_f} (\mathbf{w}_{f,i,j}^k - \widetilde{\mathbf{w}}_{f,i,j}^k) \Big]. \tag{10}$$

To illustrate the GCN-based prediction more clearly, we give the structure comparison in Figure 3, where each number in the node of KG denotes a class id and the red scores denote true predictions. According to this figure, we can find that the classifiers of unseen classes are also able to be learned within the semi-supervised learning process, which makes the zero-shot recognition for unseen classes possible. And our proposed model achieves the individual fine-grained prediction with sub-text-visual alignment to promote the performance. In other words, for each enhanced fine-grained visual feature with a certain semantic topic, we set an individual channel for GCN to implement semi-supervised training.

Finally, for more clear illustration, the whole algorithm flow of our multimodal framework is presented in **Algorithm** 1, where step 3~4 denotes the Semantic Embeddings Generation approach, step 5 describes the training way of Fine-grained Visual Knowledge Learning, and step 6~8 introduces the Structural Knowledge Learning method and the final classification procedure.

# 4 EXPERIMENTS

## 4.1 Datasets

*4.1.1 ImageNet.* We conduct the experiments of our model and baselines on the ImageNet dataset [2]. Following the setting of [11, 21], we adopt the train/test split of them and evaluate the zero-shot classification on the large 21K ImageNet dataset. Since the fine-tuning process is executed on the ImageNet 2012 1K classes, the test datasets consist of three parts: "2-hops", "3-hops", and "All". Hops refer to the distance that classes are away from the ImageNet 2012 1K classes in the ImageNet hierarchy and thus is a measure of how far unseen classes are away from seen classes. "2-hops" includes roughly 1.5K classes within two hops from the seen classes, while "3-hops" contains 7.8K classes. "All" contains close to 21K classes. All these classes are not contained in ImageNet 2012, which was used to pre-train the ResNet-50 model. What's more, we further evaluate the performance when training categories are included as potential labels. Here we called the splits as "2-hops+1K", "3-hops+1K", "All+1K". We will evaluate the Top 1~20 accuracy on these three (sub) datasets.

*4.1.2 AWA2.* We follow the setting of [21] to test the AWA2 dataset [3]. It consists of 50 animal classes, with a total of 37,322 images and an average of 746 per class. The split ensures that there is no overlap between the test classes and the ImageNet 2012 dataset, where 40 classes are used for training and 10 for testing. AWA2 test classes are contained in the 21K ImageNet and several of the training classes (24 out of 40) that are in the proposed split overlap with the ImageNet 2012 dataset. Note that, same as the setting of [21], we also only test the Top-1 accuracy of each model on the zero-shot setting since the relatively small scale of this dataset and the number of unseen classes is not enough to try large indicators.

## 4.2 Baselines

We compare our model to some representative approaches which are suitable for non-attribute datasets. Specifically, **DeViSE** [11] leverages textual data to learn semantic relationships between different labels, and explicitly maps images into a rich semantic embedding space. **ConSE** [31] projects visual features into the semantic space as a convex combination of several semantic embeddings of the $T$ nearest seen classes which are weighted by the probabilities that the image belongs to seen classes. **EXEM** [4] regards the cluster centers of visual feature vectors as the target semantic representations and leverages structural relations well on the cluster. **SYNC** [3] aims to align the semantic space with the visual space via phantom object classes to connect seen and unseen classes. **GCNZ** [39] uses both semantic embeddings and categorical relationships to predict the final classifiers. GCNZ inputs semantic embedding for each concept node (representing visual category) in the knowledge graph. **DGP** [21] is a KG-based graph embedding model which explicitly exploits the hierarchical structure of the KG to perform ZSL by propagating knowledge through the dense connectivity structure. **GPCL** [38] exploits multiple relationships among different categories of KG for zero-shot learning by employing graph convolutional representation and contrastive learning
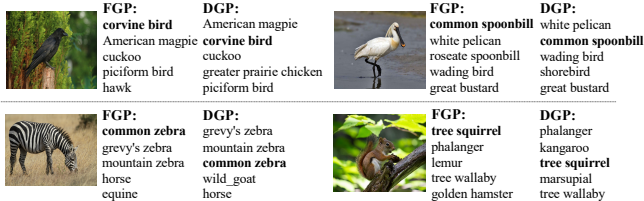
---

[2] https://image-net.org/download.php
[3] http://cvml.ist.ac.at/AwA2/

**Table 1: Top-k accuracy results for ZSL setting(only testing on unseen classes) on three subsets based on the ImageNet dataset.**

| Method type | Model | Hit@k (%) on 2-hops | | | | | Hit@k (%) on 3-hops | | | | | Hit@k (%) on All | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 5 | 10 | 20 | 1 | 2 | 5 | 10 | 20 | 1 | 2 | 5 | 10 | 20 |
| Embedding-based method | DeViSE | 6.7 | 11.2 | 19.4 | 28.1 | 38.3 | 2.1 | 3.5 | 6.3 | 9.5 | 14.1 | 1.0 | 1.8 | 3.0 | 4.6 | 7.1 |
| | ConSE | 8.3 | 12.9 | 21.8 | 30.9 | 41.7 | 2.6 | 4.1 | 7.3 | 11.1 | 16.4 | 1.3 | 2.1 | 3.8 | 5.8 | 8.7 |
| | SYNC | 10.5 | 17.7 | 28.6 | 40.1 | 52.0 | 2.9 | 4.9 | 9.2 | 14.2 | 20.9 | 1.4 | 2.4 | 4.5 | 7.1 | 10.9 |
| | EXEM | 12.5 | 19.5 | 32.3 | 43.7 | 55.2 | 3.6 | 5.9 | 10.7 | 16.1 | 23.1 | 1.8 | 2.9 | 5.3 | 8.2 | 12.2 |
| Knowledge Graph-based method | GCNZ | 19.8 | 33.3 | 53.2 | 65.4 | 74.6 | 4.1 | 7.5 | 14.2 | 20.2 | 27.7 | 1.8 | 3.3 | 6.3 | 9.1 | 12.7 |
| | DGP | 26.6 | 40.7 | 60.3 | 72.3 | 81.3 | 6.3 | 10.7 | 19.3 | 27.7 | 37.7 | 3.0 | 5.0 | 9.3 | 13.9 | 19.8 |
| | GPCL | **28.4** | **43.0** | **62.6** | **74.5** | **82.9** | 7.0 | **11.7** | 20.7 | 29.2 | 39.0 | 3.3 | 5.6 | 10.1 | **14.7** | 20.5 |
| | HVEL | 13.3 | 20.8 | 39.2 | 52.7 | 62.4 | 6.5 | 10.6 | 18.8 | 25.8 | 35.2 | **3.7** | **5.9** | 10.3 | 13.0 | 16.4 |
| | MGFN | 26.9 | 40.9 | 60.3 | 72.7 | 81.4 | 6.3 | 10.8 | 19.6 | 27.9 | 37.9 | 3.2 | 5.7 | 9.6 | 14.1 | 19.9 |
| | FGP (ours) | 26.4 | 41.2 | 61.0 | 72.8 | 81.7 | **7.3** | 11.0 | **21.3** | **29.8** | **39.4** | **3.7** | 5.5 | **10.5** | **14.7** | **20.9** |

**Table 2: Top-k accuracy results with General ZSL setting(testing instances includes seen and unseen classes simultaneously) on three subsets based on the ImageNet dataset. ‡ denotes the result from [39]. Here we also follow [21, 26] do not use SYNC and EXEM due to the algorithm characteristics.**

| Method type | Model | Hit@k (%) on 2-hops | | | | | Hit@k (%) on 3-hops | | | | | Hit@k (%) on All | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 5 | 10 | 20 | 1 | 2 | 5 | 10 | 20 | 1 | 2 | 5 | 10 | 20 |
| Embedding-based method | DeViSE | 0.8 | 2.7 | 7.9 | 14.2 | 22.7 | 0.5 | 1.4 | 3.4 | 5.9 | 9.7 | 0.3 | 0.8 | 1.9 | 3.2 | 5.3 |
| | ConSE | 0.3 | 6.2 | 17.0 | 24.9 | 33.5 | 0.2 | 2.2 | 5.9 | 9.7 | 14.3 | 0.2 | 1.2 | 3.0 | 5.0 | 7.5 |
| | ConSE‡ | 0.1 | 11.2 | 24.3 | 29.1 | 32.7 | 0.2 | 3.2 | 7.3 | 10.0 | 12.2 | 0.1 | 1.5 | 3.5 | 4.9 | 6.2 |
| Knowledge Graph-based method | GCNZ | 9.7 | 20.4 | 42.6 | 57.0 | 68.2 | 2.2 | 5.1 | 11.9 | 18.0 | 25.6 | 1.0 | 2.3 | 5.3 | 8.1 | 11.7 |
| | DGP | 10.3 | 26.4 | 50.3 | 65.2 | 76.0 | 2.9 | 7.1 | 16.1 | 24.9 | 35.1 | 1.4 | 3.4 | 7.9 | 12.6 | 18.7 |
| | GPCL | 7.0 | 26.8 | 52.5 | **67.5** | **77.9** | 2.0 | 7.1 | 17.3 | 26.2 | **36.5** | 1.0 | 3.4 | 8.5 | 13.2 | **19.3** |
| | HVEL | 6.4 | 11.9 | 27.2 | 35.3 | 45.2 | 3.6 | **8.7** | 15.3 | 20.5 | 29.1 | 2.2 | 4.6 | 9.2 | 12.7 | 15.5 |
| | MGFN | **13.5** | 28.1 | 51.1 | 65.4 | 76.1 | 3.7 | 7.3 | 16.5 | 25.3 | 35.9 | 1.8 | 3.5 | 8.2 | 12.8 | 18.9 |
| | FGP (ours) | 12.0 | **28.8** | **53.0** | 65.4 | 75.8 | **4.0** | 8.7 | **17.9** | **26.6** | 36.2 | **2.5** | **4.8** | **9.4** | **13.5** | 18.4 |



**Figure 4: The qualitative comparison cases of top-5 classification results on some unseen categories under the ZSL setting. The correct category is marked as bold in each case.**

techniques. **HVEL** [26] further learns hierarchical-aware image embeddings in hyperbolic space for zero-shot learning, which is capable of preserving the hierarchical structure of semantic classes. **MGFN** [41] develops a multi-granularity fusion network that integrates discriminative information from multiple GCN branches.

## 4.3 Implementation Details

We use a ResNet-50 model that has been pre-trained on the ImageNet 2012 dataset. We extract a feature map at the last convolutional layer whose size is $7 \times 7 \times 2048$ and treat it as a set of

**Table 3: Top-1 accuracy (%) results for zero-shot recognition on the AWA2 dataset.**

| Method type | Model | ZSL | GZSL |
|---|---|---|---|
| Embedding-based method | ConSE | 44.5 | 5.8 |
| | DeViSE | 59.7 | 16.0 |
| | SYNC | 46.6 | - |
| | EXEM | 55.6 | - |
| Knowledge Graph-based method | GCNZ | 70.7 | 37.3 |
| | DGP | 77.3 | 40.2 |
| | GPCL | 77.8 | 41.4 |
| | HVEL | 77.0 | 39.8 |
| | MGFN | 74.2 | 38.3 |
| | FGP (ours) | **79.1** | **43.3** |

features from $R = 7 \times 7$ regions. The GloVe model trained on the Wikipedia dataset is utilized to obtain the feature representation of concepts' descriptions in the KG (we use the descriptions from WordNet [30] in this paper). Every single GCN in the Multi-GCN model consists of layers with hidden dimensions of 2048 and the final output dimension corresponds to the number of weights in the last layer of the ResNet-50 architecture, 2049 for weights and bias. In addition, we regularize the outputs into similar ranges by

L2-Normalization. Similarly, we also normalize the ground truth weights produced by CNN. To avoid over-fitting, we implement Dropout [36] with a dropout rate of 0.5 in each layer of GCN. Our Multi-GCN model is trained for 3000 epochs with a learning rate of 0.001 and weight decay of 0.0005 by the Adam optimizer [22]. The negative slope of leaky ReLU is 0.2. The proposed FGP model is performed on an NVIDIA Tesla V100 GPU. Fine-tuning is done for 20 epochs using SGD with a learning rate of 0.0001 and momentum of 0.9. What's more, we explore the influence of $K$, the number of key semantic embeddings, and $L$, the number of GCN layers in the ablation studies. The dimension $d_c = 300$ of input vectors of Kmeans is the same as the output dimension of Glove.

## 4.4 Performance Comparison

*4.4.1 Quantitative Performance Comparison.* We report the quantitative results of ZSL testing and General ZSL testing on ImageNet in Table 1 and Table 2, respectively. To give a more intuitive comparison, some qualitative comparison results are shown in the next subsection. The results of baselines are taken from [21, 26, 38]. Following the setting of them, we utilize the Top-k accuracy to evaluate the model as well, and the EXEM model and the SYNC model are only used on the ZSL task due to their characteristics. Overall, our proposed model FGP achieves obviously better results than other comparing methods on most of the testing tasks. Compared with previous powerful KG-based models DGP, GPCL, MGFN, and HVEL, with the help of fine-grained knowledge, our model reflects a more comprehensive performance on different scale settings. It is easy to find that HVEL shows a serious bias on the "All" test data, which outperforms all baselines including DGP on this experimental part but is significantly weaker than DGP, GCNZ on both ZSL and GZSL evaluations when recognizing the "2-hops" dataset. Similarly, the GPCL model performs well on the "2-hops" ZSL task but is not enough good as our model on other settings, especially the GZSL tasks. On the contrary, although the bias tendency of DGP is not so critical as the same as former, this model still performs not good on the biggest dataset "All".

Going deeper into the results produced by FGP, it beats all rivals in the "2-hops+1K" and "All + 1K" tasks and reaches the highest Top-1 scores 4.0% and 2.5%, which provide 8.1% and 13.6% relative improvement over the previous state-of-the-art. Besides, on the "3-hops" and "All" datasets of the ZSL setting, FGP still shows the best performance. Specifically, FGP surpasses all baselines on most indicators with a distinct margin, where HVEL is slightly better than our model on Top-2 accuracy but not ideal on other indicators. To summarize, the single-modality KG model learns the semantic transfer by relationships, which gradually loses the recognition ability as the number of hops increases. And the fine-grained features captured by FGP (multimodal KG) alleviate this shortage.

On the other hand, the further exploration results on the AWA2 dataset are also shown in Table 3. Here we follow the setting of DGP [21] to test the Top-1 accuracy of each model since the relatively small scale of this dataset and the number of unseen classes is not enough to try large indicators. Obviously, our proposed model gets the best result of 79.1% and 43.3% accuracy on the ZSL and GZSL settings, respectively.
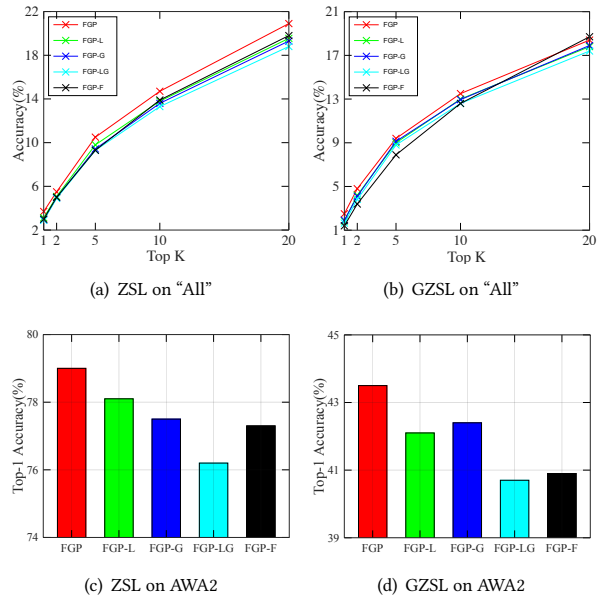


(a) ZSL on "All"    (b) GZSL on "All"

(c) ZSL on AWA2    (d) GZSL on AWA2

**Figure 5: The module analysis of ablation studies of our proposed model FGP for both the ZSL and GZSL tasks on the "All" test set and the AWA2 set. -L, -G, -LG, -F indicates FGP w/o the self-calibration loss $\mathcal{L}_{sc}$, w/o the visual global features $e_g$, w/o $\mathcal{L}_{sc}$ and $e_g$, and only uses visual global features.**

*4.4.2 Qualitative Comparison.* To give a more intuitive comparison, some qualitative comparison results produced by models are shown in Figure 4. The samples come from the "2-hops" set and their corresponding predicting tags are produced by DGP and FGP under the ZSL setting (only testing on unseen classes), since the DGP model inspires our motivation. Intuitively, from the top-5 classification results of each testing image, we observe that DGP and FGP generally provide coherent top-5 results due to the relationships of classes in the knowledge graph. According to Figure 4, with the help of more fine-grained visual features and textual features, our model obtains better performance in many samples on zero-shot classification. For example, although "common zebra" and "grevy's zebra" are very similar in body shape, they can easily be distinguished by the shape of their ears, where the former has pointed ears and the latter has round ears.

## 4.5 Ablation Studies

*4.5.1 Module Analysis.* To further verify the effectiveness of each modified module in our model, we design an ablation experiment carefully. In detail, the experiment includes FGP (our proposed complete model), FGP-L (FGP w/o the self-calibration loss $\mathcal{L}_{sc}$), FGP-G (FGP w/o the visual global features $e_g$), FGP-LG (FGP w/o $\mathcal{L}_{sc}$ and the visual global features $e_g$), and FGP-F (the model only uses visual global features, it can be viewed as DGP). We run these models on the ImageNet 2012 dataset and test their ZSL and GZSL performances both on the "All" set and AWA2 set, relevant experimental results are all recorded in Figure 5. Figure 5(a)-(b) offer the Top-1~20 accuracy on "All". Figure 5(c)-(d) offer the Top-1 accuracy

**Table 4: Results of the fine-tuning experiments on the "All" dataset. (-f1), (-f2), and (-f1f2) indicate the FGP model without the first fine-tuning, the second fine-tuning and both fine-tunings, respectively. We also present the results of DGP and DGP(-f) (without the fine-tuning) for comparison.**

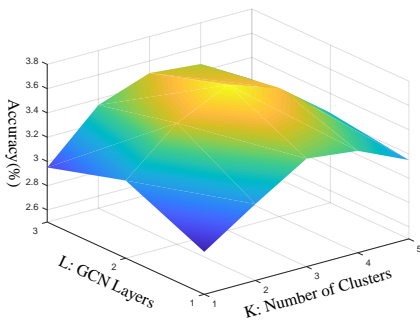| Test set | Model | Hit@k (%) | | | | |
|---|---|---|---|---|---|---|
| | | 1 | 2 | 5 | 10 | 20 |
| **All** | DGP (-f) | 2.6 | 4.7 | 9.0 | 13.5 | 19.1 |
| | DGP | 3.0 | 5.0 | 9.3 | 13.9 | 19.8 |
| | FGP (-f1) | 2.6 | 4.5 | 8.7 | 13.0 | 18.4 |
| | FGP (-f2) | 3.2 | 5.2 | 9.6 | 13.7 | 19.6 |
| | FGP (-f1f2) | 2.3 | 4.3 | 8.7 | 13.2 | 17.7 |
| | FGP | **3.7** | **5.5** | **10.5** | **14.7** | **20.9** |



**Figure 6: The visualization results of our hyperparameter tune-up experiment for the number of GCN layers and semantic clusters on the "All" dataset.**

on AWA2. We can find that the complete model FGP achieves the best results in both two tasks, but there are no distinct performance margins among other sub-models. The worst performer is FGP-LG. Due to the lack of corresponding modules, the effects of FGP-L, FGP-G, and FGP-F are obviously worse, which shows that our self-calibration loss, global visual features, and fine-grained features are all necessary and useful in the framework.

*4.5.2 Fine-tuning Analysis.* There are 2 times of fine-tuning for model training in our proposed scheme, which can be easily distinguished by reviewing **Algorithm** 1. The first fine-tuning in step 5 is used for optimizing the dense attention network and the second fine-tuning in step 8 can correct deviations between the feature extractor (including dense attention) and the replaced classifier weights learned by Multi-GCN. We observe that fine-tuning consistently improves performance for both models in all our experiments. Ablation studies that highlight the impact of fine-tuning for the "All" scenario can be found in Table 4. FGP (-f1f2) is viewed as the accuracy that is achieved after training the FGP model without both fine-tunings. FGP (-f1) and FGP (-f2) is used to denote the results for FGP without the first fine-tuning and the second fine-tuning, respectively. We further report the accuracy achieved by the DGP model without fine-tuning (DGP(-f)) since DGP is the base approach of FGP. We observe that the effect of FGP (-f1) and FGP (-f1f2) is much lower than FGP (-f2) and the complete model FGP, that is because the parameters of dense attention module cannot

be optimized within their training phrase. Overall, fine-tuning the zero-shot learning model consistently leads to improved results.

## 4.6 Hyperparameter Tune-up

In our model, there are two major hyperparameters: GCN layers $L$ and the number of semantic clusters $K$, which affect the neighbor hops of nodes in KG and the fine-grained level respectively. We intend to explore corresponding influences bring back by these two hyperparameters in this subsection. As shown in Figure 6, the Top-1 accuracy of the FGP model is plotted in a 3D chart with different hyperparameters on the "All" dataset, where the value ranges of $L$ and $K$ are $[1, 2, 3]$ and $[1, 2, 3, 4, 5]$, respectively. In this figure, the warmer the tone, the higher the accuracy. Clearly, our model ascends to the peak when $L = 2$ and $K = 3$. The too-small value of $K$ leads to more overlapping of semantic space, and too-large $K$ would bring too many model parameters that result in over-fitting. Cater to the structural characteristics of the dataset, $K = 3$ makes the model achieve relatively better performances than other situations, which forms the ridge of our result distribution. So the optimal number of semantic clusters is 3 in this dataset. However, GCN layers don't show an obvious effect with the regular pattern on the performance of FGP, except for the lowest score when $K = 1$. It shows that the idea of only aggregating features of entities with long paths (more neighbors) in ordinary KG has great limitations in our zero-shot learning framework.

## 5 CONCLUSION

In this paper, we proposed a novel Fine-grained Graph Propagation (FGP) model based on the multimodal KG, which enhanced the semantic transfer process at a fine-grained knowledge level for the zero-shot recognition task. More specifically, we first generated key semantic embeddings which matched with relevant regions of input images in the fine-grained level visual-semantic alignment process. Then a Multi-GCN was used to learn the classifiers that mapped into all seen and unseen classes. Along this line, after a slight fine-tuning, our model achieved the goal of embedding fine-grained features in the multimodal KG-based zero-shot learning framework. We conducted extensive experiments on large real-world datasets and the results clearly demonstrated the effectiveness of FGP compared with several state-of-the-art methods. Hence our learning paradigm can be viewed as an effective booster to achieve the breakthrough of related applications which face the dilemma of sparse manual labels. The further consideration of additional knowledge graph resources such as Wikipedia for settings where these are available for a subset of nodes is a future direction.

## 6 ACKNOWLEDGMENTS

## REFERENCES

[1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. *ICLR* (2015).

[2] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Van-dergheynst. 2017. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine* 34, 4 (2017), 18–42.

[3] Soravit Changpinyo, Wei-Lun Chao, Boqing Gong, and Fei Sha. 2016. Synthesized classifiers for zero-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 5327–5336.

[4] Soravit Changpinyo, Wei-Lun Chao, and Fei Sha. 2017. Predicting visual ex-emplars of unseen classes for zero-shot learning. In *Proceedings of the IEEE international conference on computer vision.* 3476–3485.

[5] Liyi Chen, Zhi Li, Tong Xu, Han Wu, Zhefeng Wang, Nicholas Jing Yuan, and Enhong Chen. 2022. Multi-modal Siamese Network for Entity Alignment. In *Proc. of KDD.*

[6] Long Chen, Hanwang Zhang, Jun Xiao, Wei Liu, and Shih-Fu Chang. 2018. Zero-shot visual recognition using semantics-preserving adversarial embedding net-works. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* 1043–1052.

[7] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 257–266.

[8] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolu-tional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems.* 3844–3852.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition.* Ieee, 248–255.

[10] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, and Alán Aspuru-Guzik. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems.* 2224–2232.

[11] Andrea Frome, Greg Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. 2013. Devise: A deep visual-semantic embedding model. (2013).

[12] Yanwei Fu and Leonid Sigal. 2016. Semi-supervised vocabulary-informed learning. In *Proceedings of CVPR.* 5337–5346.

[13] Yuxia Geng, Jiaoyan Chen, Zhuo Chen, Jeff Z Pan, Zhiquan Ye, Zonggang Yuan, Yantao Jia, and Huajun Chen. 2021. OntoZSL: Ontology-enhanced Zero-shot Learning. In *Proceedings of the Web Conference 2021.* 3325–3336.

[14] Yuxia Geng, Jiaoyan Chen, Wen Zhang, Yajing Xu, Zhuo Chen, Jeff Z Pan, Yufeng Huang, Feiyu Xiong, and Huajun Chen. 2022. Disentangled Ontology Embedding for Zero-shot Learning. *arXiv preprint arXiv:2206.03739* (2022).

[15] Marco Gori, Gabriele Monfardini, and Franco Scarselli. 2005. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, Vol. 2. IEEE, 729–734.

[16] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NIPS.* 1024–1034.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 770–778.

[18] Mikael Henaff, Joan Bruna, and Yann LeCun. 2015. Deep convolutional networks on graph-structured data. *arXiv preprint arXiv:1506.05163* (2015).

[19] Dat Huynh and Ehsan Elhamifar. 2020. Compositional zero-shot learning via fine-grained dense feature composition. *Advances in Neural Information Processing Systems* 33 (2020), 19849–19860.

[20] Dat Huynh and Ehsan Elhamifar. 2020. Fine-grained generalized zero-shot learning via dense attribute-based attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 4483–4493.

[21] Michael Kampffmeyer, Yinbo Chen, Xiaodan Liang, Hao Wang, Yujia Zhang, and Eric P Xing. 2019. Rethinking knowledge graph propagation for zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 11487–11496.

[22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic opti-mization. *arXiv preprint arXiv:1412.6980* (2014).

[23] Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. *ICLR* (2017).

[24] Lu Liu, Tianyi Zhou, Guodong Long, Jing Jiang, and Chengqi Zhang. 2020. At-tribute Propagation Network for Graph Zero-shot Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 4868–4875.

[25] Lu Liu, Tianyi Zhou, Guodong Long, Jing Jiang, and Chengqi Zhang. 2020. At-tribute Propagation Network for Graph Zero-Shot Learning. *AAAI* (2020).

[26] Shaoteng Liu, Jingjing Chen, Liangming Pan, Chong-Wah Ngo, Tat-Seng Chua, and Yu-Gang Jiang. 2020. Hyperbolic visual embedding learning for zero-shot recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 9273–9281.

[27] Yu Liu and Tinne Tuytelaars. 2020. A Deep Multi-Modal Explanation Model for Zero-Shot Learning. *IEEE Transactions on Image Processing* (2020).

[28] Zhekun Luo, Shalini Ghosh, Devin Guillory, Keizo Kato, Trevor Darrell, and Huijuan Xu. 2022. Disentangled Action Recognition with Knowledge Bases. In

[29] James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, Vol. 1. Oakland, CA, USA, 281–297.

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

[31] Mohammad Norouzi, Tomas Mikolov, Samy Bengio, Yoram Singer, Jonathon Shlens, Andrea Frome, Greg S Corrado, and Jeffrey Dean. 2014. Zero-shot learning by convex combination of semantic embeddings. *ICLR* (2014).

[32] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP).* 1532–1543.

[33] Zhaopeng Qiu, Yunfan Hu, and Xian Wu. 2022. Graph Neural News Recommen-dation with User Existing and Potential Interest Modeling. *ACM Trans. Knowl. Discov. Data* 16, 5 (2022), 96:1–96:17.

[34] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE Transactions on Neural Networks* 20, 1 (2008), 61–80.

[35] Richard Socher, Milind Ganjoo, Hamsa Sridhar, Osbert Bastani, Christopher D Manning, and Andrew Y Ng. 2014. Zero-shot learning through cross-modal transfer. *NIPS* (2014).

[36] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15, 1 (2014), 1929–1958.

[37] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph attention networks. *ICLR* (2018).

[38] Jin Wang and Bo Jiang. 2021. Zero-Shot Learning via Contrastive Learning on Dual Knowledge Graphs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 885–892.

[39] Xiaolong Wang, Yufei Ye, and Abhinav Gupta. 2018. Zero-shot recognition via semantic embeddings and knowledge graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 6857–6866.

[40] Zheng Wang, Jialong Wang, Yuchen Guo, and Zhiguo Gong. 2021. Zero-shot Node Classification with Decomposed Graph Prototype Network. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining.* 1769–1779.

[41] Jiwei Wei, Haotian Sun, Yang Yang, Xing Xu, Jingjing Li, and Heng Tao Shen. 2022. Semantic guided knowledge graph for large-scale zero-shot learning. *Journal of Visual Communication and Image Representation* 88 (2022), 103629.

[42] Likang Wu, Zhi Li, Hongke Zhao, Qi Liu, Jun Wang, Mengdi Zhang, and Enhong Chen. 2021. Learning the Implicit Semantic Representation on Graph-Structured Data. *arXiv preprint arXiv:2101.06471* (2021).

[43] Likang Wu, Zhi Li, Hongke Zhao, Zhen Pan, Qi Liu, and Enhong Chen. 2020. Es-timating Early Fundraising Performance of Innovations via Graph-Based Market Environment Model.. In *AAAI.* 6396–6403.

[44] Likang Wu, Hongke Zhao, Zhi Li, Zhenya Huang, Qi Liu, and Enhong Chen. 2023. Learning the Explainable Semantic Relations via Unified Graph Topic-Disentangled Neural Networks. *ACM Transactions on Knowledge Discovery from Data* (2023).

[45] Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, et al. 2023. A Survey on Large Language Models for Recommendation. *arXiv preprint arXiv:2305.19860* (2023).

[46] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. 2018. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence* 41, 9 (2018).

[47] Guo-Sen Xie, Li Liu, Xiaobo Jin, Fan Zhu, Zheng Zhang, Jie Qin, Yazhou Yao, and Ling Shao. 2019. Attentive region embedding network for zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 9384–9393.

[48] Guo-Sen Xie, Li Liu, Fan Zhu, Fang Zhao, Zheng Zhang, Yazhou Yao, Jie Qin, and Ling Shao. 2020. Region graph embedding network for zero-shot learning. In *European conference on computer vision.* Springer, 562–580.

[49] Wenjia Xu, Yongqin Xian, Jiuniu Wang, Bernt Schiele, and Zeynep Akata. 2022. VGSE: Visually-Grounded Semantic Embeddings for Zero-Shot Learning. In *Proceedings of the CVPR.* 9316–9325.

[50] Caixia Yan, Qinghua Zheng, Xiaojun Chang, Minnan Luo, Chung-Hsing Yeh, and Alexander G. Hauptman. 2020. Semantics-Preserving Graph Propagation for Zero-Shot Object Detection. *IEEE Transactions on Image Processing* (2020).

[51] Chunjie Zhang, Chao Liang, and Yao Zhao. 2022. Exemplar-Based, Semantic Guided Zero-Shot Visual Recognition. *IEEE Transactions on Image Processing* 31 (2022), 3056–3065.

[52] Zhi Zheng, Chao Wang, Tong Xu, Dazhong Shen, Penggang Qin, Xiangyu Zhao, Baoxing Huai, Xian Wu, and Enhong Chen. 2023. Interaction-aware drug package recommendation via policy gradient. *ACM Transactions on Information Systems* 41, 1 (2023), 1–32.

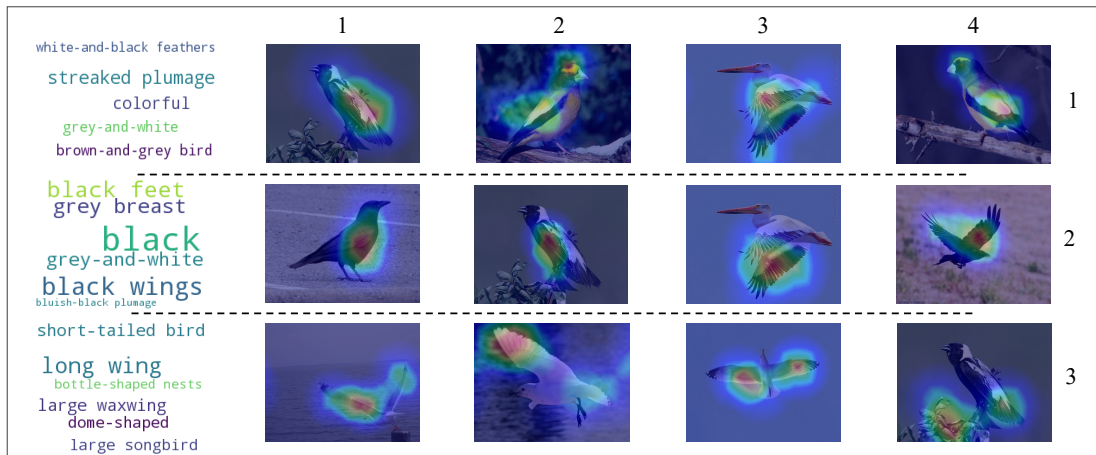*Proceedings of the 2022 NAACL.* 559–572.

**Figure 7: The visualization results of attention maps with corresponding semantic embeddings on the bird subset. The most representative words of each semantic cluster are presented in the form of word cloud. The attention distribution of each image is clearly displayed through the heat map.**

## A  APPENDICES

### A.1  Fine-Grained Features Visualization

We intend to observe the focused topics of key semantic embeddings and the status of fine-grained level visual-semantic alignment. To present a more intuitive and concise qualitative result, we conduct this experiment on the collected bird subset from ImageNet which has relatively unified and distinct strong-visibility features (e.g. body, wings, feet, etc.). We show the semantic phrase clusters of our Kmeans module (here generating three clusters in our setting which are represented by the word clouds in Figure 7, the font size of a phrase denotes its centrality weight in a cluster) and corresponding cases (the images with attention heatmaps). According to Figure 7, there is a major focused topic in each semantic word cloud, and the attention heatmaps of their cases are able to match

them to a considerable extent. For instance, keeping up with respective semantic contents, the images in the first row can strengthen the streaked and mixed-color parts, while the second row concentrates on the solid dark regions. We can find that the major focused semantics of the word clouds of the first two rows are "streaked, colorful" and the dark tone like "black, grey" too. Compared with the formers, the semantic topic of the third row does not seem very clear and obvious, so its case performances are not as good as the first two. The $case_{3,4}$ assigns some larger weights in unimportant areas, but $case_{3,1}$, $case_{3,2}$, and $case_{3,3}$ are all able to find the regions in original images that correspond to the important "long/large wing" semantic information. To sum up, the visualization shows that our model can realize the motivation of fine-grained feature alignment in the KG-based zero-shot learning, which helps the cognitive system distinguish similar but locally different objects more precisely such as $case_{1,1}$ and $case_{2,1}$ in Figure 7.