# Webly Supervised Image Classification with Metadata: Automatic Noisy Label Correction via Visual-Semantic Graph

**Jingkang Yang***
SenseTime Research
Dept. of Electrical and Computer
Engineering, Rice University
yangjingkang@sensetime.com

**Weirong Chen***
SenseTime Research
Dept. of Comp. Sci. & Eng., The
Chinese University of Hong Kong
chenweirong@sensetime.com

**Litong Feng**
SenseTime Research
fenglitong@sensetime.com

**Xiaopeng Yan**
SenseTime Research
yanxiaopeng@sensetime.com

**Huabin Zheng**
SenseTime Research
zhenghuabin@sensetime.com

**Wayne Zhang**
SenseTime Research
Qing Yuan Research Institute,
Shanghai Jiao Tong University
wayne.zhang@sensetime.com

## ABSTRACT

Webly supervised learning becomes attractive recently for its efficiency in data expansion without expensive human labeling. However, adopting search queries or hashtags as web labels of images for training brings massive noise that degrades the performance of DNNs. Especially, due to the semantic confusion of query words, the images retrieved by one query may contain tremendous images belonging to other concepts. For example, searching 'tiger cat' on Flickr will return a dominating number of tiger images rather than the cat images. These realistic noisy samples usually have clear visual semantic clusters in the visual space that mislead DNNs from learning accurate semantic labels. To correct real-world noisy labels, expensive human annotations seem indispensable. Fortunately, we find that metadata can provide extra knowledge to discover clean web labels in a labor-free fashion, making it feasible to automatically provide correct semantic guidance among the massive label-noisy web data. In this paper, we propose an automatic label corrector *VSGraph-LC* based on the visual-semantic graph. *VSGraph-LC* starts from anchor selection referring to the semantic similarity between metadata and correct label concepts, and then propagates correct labels from anchors on a visual graph using graph neural network (GNN). Experiments on realistic webly supervised learning datasets Webvision-1000 and NUS-81-Web show the effectiveness and robustness of *VSGraph-LC*. Moreover, *VSGraph-LC* reveals its advantage on the open-set validation set.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; *Supervised learning by classification*; *Neural networks*; *Information extraction.*

## KEYWORDS

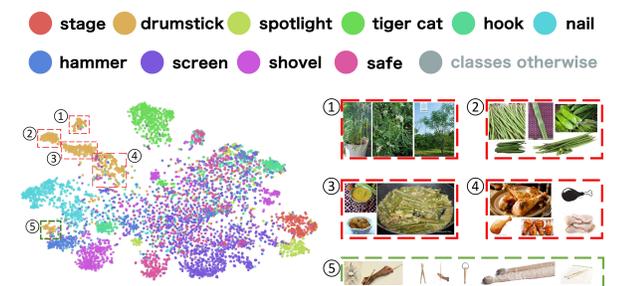Webly Supervised Learning; Metadata; Semantic Label Noise; Visual-Semantic Graph; Graph Neural Networks
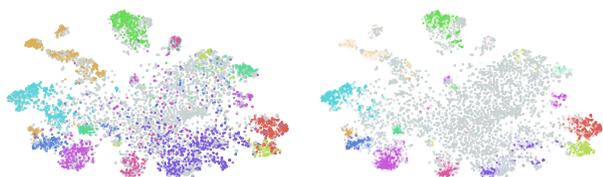
## 1 INTRODUCTION

Deep convolutional neural networks (CNNs) are successful by virtue of large-scale datasets with human annotation [23]. However, human annotation is extremely time-consuming and expensive, which impedes the further expansion of those big datasets [42]. To overcome this limitation, researchers use web crawlers to collect billions of images and annotate them directly using text queries or hashtags [20, 30]. However, due to the ambiguity or polysemy of the query fed into the search engine, label noise is subsequently introduced. Therefore, webly supervised learning, aiming at using huge scalable web-crawled data directly for networks training by suppressing label noise, has attracted great attention recently [1].

Early exploration in this direction relies on human-verified clean subsets. Representative methods using clean subsets include MentorNet [17] and CleanNet [24]. However, with the trend of rapid growth in the size of webly datasets, building clean subsets with manual verification becomes more infeasible, especially when the number of categories exceeds ten thousands [45]. Therefore, recent works prefer models without clean subset dependencies, making webly supervised learning fully automatic. To this end, some works strengthen networks' endurability against label noise using moving average of model predictions [40], loss function modification [11, 31] or regularization [34, 54] . Co-teaching uses two different networks to mutually detect label noise, doubly ensuring the model's denoising ability [13]. Other works identify noisy samples based on some hypotheses including data density [12, 14] or model confidence [50].

(a) Colors reflect web labels. Web label 'Drumstick' shows representative images corresponding to 5 regions of interest. Only Region-5 corresponds to the correct concept
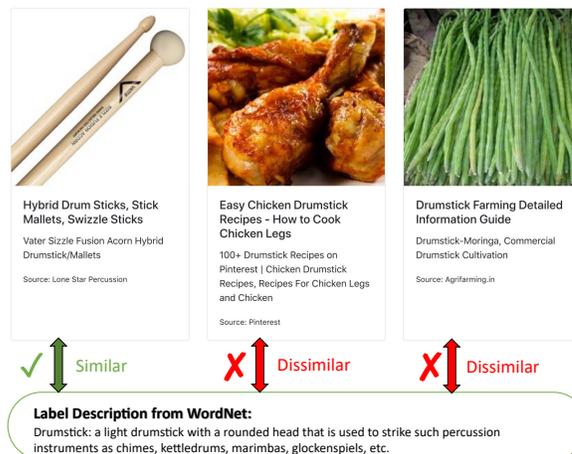


(b) Prediction by Co-teaching          (c) Prediction by VSGraph-LC

**Figure 1: T-SNE visualization [29] of WebVision-pretrained ResNet50 [15] features on 10 selected categories. Three observations are highlighted: (1) CNN models that trained from WebVision can distinguish different semantics within a category, even when semantics mismatch category definition. (2) Severe semantic label noise is a real-world problem, as majority images of class 'drumsticks' deviate the true concept of percussion mallets. (3) Co-teaching fails to correct the majority semantic label noise, but our VSGraph-LC is able to. Node brightness represents prediction confidence**

Although the aforementioned methods effectively enhance the model against label noise, especially for outliers, suppressing *semantic label noise* is critical but untouched. To clarify, semantic label noise exists due to query's polysemy or insufficient semantic resolution. Usually, semantic label noise would be severe in some category, which is composed of a large number of samples reflecting another semantic concept. As webly datasets come from real-world, those off-target semantics are usually out-of-distribution (i.e., deviate from all semantic labels or out of interests in the test sets).

For example, Figure 1a shows that although web label 'drumsticks' has its correct semantic concept of 'percussion mallets' according to the test set, the majority of training samples actually belong to concepts of 'drumstick trees/vegetable' and 'chicken legs' due to polysemy. These out-of-distribution noisy samples clearly cluster themselves in the visual feature space. As a result, density assumption popularly adopted by previous methods [12, 14] will regard 'chicken legs' and 'drumstick trees/vegetable' samples falsely as clean data. Figure 1b further shows the ineffectiveness of the representative self-training method Co-teaching, where most samples belong to off-target semantics are still predicted as positive.

As webly dataset is crawled from the Internet, text metadata associated with web images has great potential to provide valuable



**Figure 2: Exemplar metadata associated with web images labeled as 'Drumstick'. By comparing metadata with label description, images with *semantic label noise* will be detected**

information, which, however, has been ignored for a long time. In this paper, we aim to take advantage of extra knowledge provided by metadata to suppress label noise, especially semantic label noise. Figure 2 shows that information in metadata and label description can be used to detect semantic label noise. By using off-the-shelf natural language processing (NLP) models [52], we convert the comparison between label description and metadata from human cognition level to text feature space, which ensures a fully automated process to precisely locate samples with correct semantic concepts, without the need for expensive manual annotations.

Hence, we are motivated to build an automatic pipeline for webly supervised learning with metadata. Specifically, we propose a label corrector named VSGraph-LC, which first selects anchor samples for each category through matching label description from Word-Net [33] and metadata of every crawled image using a powerful NLP model XLNet [52]. To help those semantically correct anchors propagate their web labels towards more samples, we leverage a graph neural network (GNN) [18] training on $k$-NN visual feature graph of the entire training set. The corrected labels substitute the former noisy web labels for finetuning our final model.

In summary, our contributions are mainly three-fold:

- We explore two understudied but important factors under webly supervised learning setting: semantic label noise and text metadata.
- A human-labor-free label correcting framework that fully exploits the merits of GNN and CNN is proposed as VSGraph-LC, ignited by anchors automatically selected by metadata.
- The proposed framework is shown effective on NUS-81-Web and WebVision datasets, and reaches the state-of-the-art result on WebVision-1000. VSGraph-LC produces more appealing results if the test set contains out-of-distribution images[1].

---

[1]Known as open-set recognition task, introduced in [22]

## 2 RELATED WORK

### 2.1 Webly Supervised Learning

A formal definition of webly supervised learning is in Section 3.1. Based on the dependency of clean subsets, webly supervised learning methods can be divided into two categories. Methods that use clean subsets to guide the detection of label noise include MentorNet, which learns a dynamic curriculum as a sample-weighting scheme from a human-labeled subset [17]. CleanNet transfers knowledge learned from human-verified samples from a fraction of categories to the entire dataset for sample reweighting [24]. In [48], a probabilistic graphical model is learned from a human-verified subset to depict relationships between images, class labels, and label noise. However, with the trend that webly datasets are exceeding billions of training data with more than ten thousand categories [45], the exponential increase of human annotations seems infeasible.

Therefore, solving webly supervised image classification without any human-verified labels attracts more attention recently. CurriculumNet trains an image classifier with a curriculum arranged in an unsupervised manner [12]. By assuming that samples from high-density regions in visual feature space have more correct labels, a three-stage training strategy is designed to feed model from clean to noisy samples. A similar assumption is adopted in [14], which selects prototypes in high-density regions. All samples get their labels corrected based on the similarity between samples and prototypes. Co-teaching trains two networks simultaneously, letting one be trained on possible clean samples selected by the other. Such cross-update can reduce the self-accumulative error from single model and therefore enhances robustness. However, according to our inspection in Section 1 and Figure 1, previous annotation-free methods are vulnerable to massive semantic label noise. [36] creatively leverages text from a strong pretrained phrase generator to suppress label confusion, but requires a complicated two-stream design with large network architecture. Thus, we are motivated to propose a fully automatic pipeline to solve the semantic label noise problem, aided by metadata crawled with web images.

### 2.2 Metadata and Concept Learning

Online platforms, including search engines and social media, can not only provide abundant web images, but also meaningful metadata. Various computer vision tasks utilize the potential value of metadata, including powerful visual-language models pretraining [28, 38], image retrieval [10, 26], and visual question answering [21, 37]. However, for the problem of webly supervised learning, metadata is unfortunately neglected.

Recent work has also proposed several visual concept discovery approaches based on metadata. In [39], unreliable concepts extracted from metadata are filtered by cross-validation average precision from a simple classifier. The remaining concepts are then clustered as concept vocabulary for downstream tasks. [2] selects visual exemplars using a clustering method with metadata, and manually assigns concepts for cluster exemplars to guide the classifier training. Furthermore, ConceptLearner [55] uses an automatic threshold method for concept allocation to clusters. Still, based on clustering, NEIL [4] establishes a lifelong training system that progressively identifies concepts and expands the dataset for better classifiers. Concept detectors and extra knowledge of concept relationships are also applied in NEIL. Those labor-free clustering methods focus on efficiently defining reasonable concepts from web data, rather than dealing with massive semantic label noise where web labels and target concepts are mismatched. In this paper, we collect precise concept definitions from WordNet [33]. Metadata is utilized to pinpoint reliable images with correct concepts for each category. Our work is orthogonal to concept learning.

### 2.3 Text Embedding

Text embedding is a fundamental topic in NLP and has made considerable progress in the past decade. With the emergence of neural probability language models, word2vec becomes one of the most widely used word embedding methods. This method uses self-supervised representation learning, which assumes that words placed in a similar context have close meanings [32]. Recent works uses powerful Transformer models [41] with carefully-designed pretraining tasks and large-scale corpora, which can even outperform human performance on multiple NLP benchmarks [27, 38, 52]. In our work, we use XLNet to encode the metadata and label description into a vector format.

### 2.4 Graph Neural Networks

Graph neural networks (GNNs) become popular in the last five years when combining deep learning and graph theory to learn from graph structured data [47, 56]. Representative models ChebNet [7] and GCN [19] use efficient filtering approaches in graph convolution operators to improve scalability and robustness. Plentiful GNN variants emerge sequentially. The simplest GNN model, simple convolutional network (SGC) [46], claims that while eliminating unnecessary complexity and redundant computations, its simplification can still maintain accuracy.
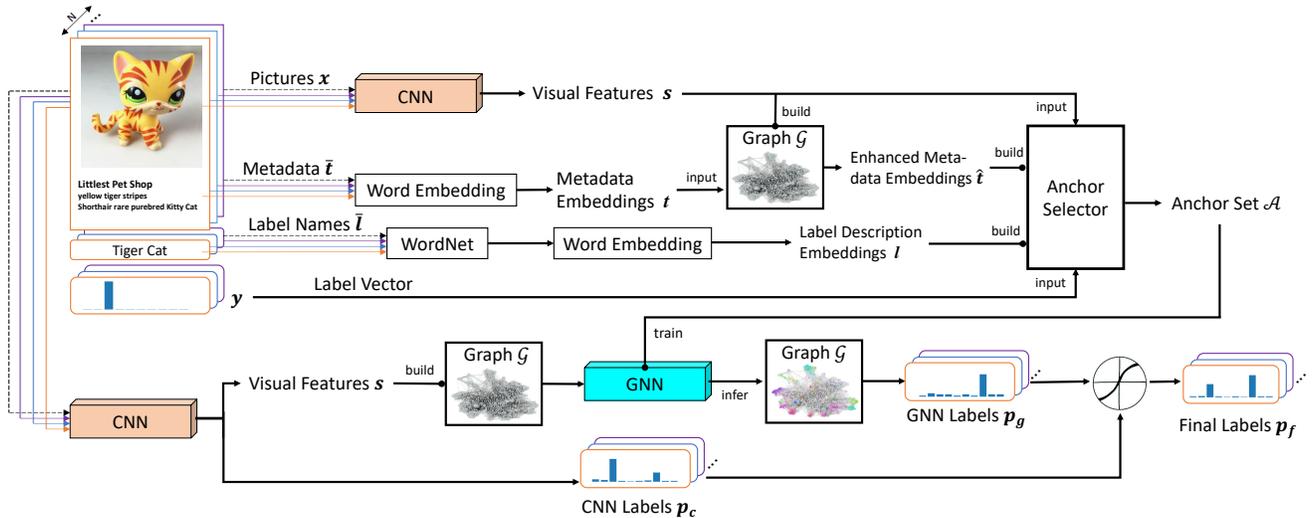
Recent works also attempt to deal with image classification tasks on the graph settings. Some works build a knowledge graph to provide possibility of label coexistence for multi-label classification [5]. Some other works do not foucs on label space, but use the neighborhood on the visual feature space to enhance the classifier [16, 51]. In our work, based on the observation in Figure 1 (even if semantic label noise is severe, features extracted from DNNs can still be clustered based on semantic information), we follow the visual feature graph path and train GNN for label correction.

## 3 PROPOSED METHOD

To correct the massive semantic label noise, our strategy is to use text metadata to provide correct semantic guidance for label correction, on the graph spanned in the training feature space. Since our label correction method is characterized by semantic guidance on the visual graph, we entitle it as a visual-semantic graph-based label corrector, referred as 'VSGraph-LC'. The flow chart of VSGraph-LC is illustrated in Figure 3. The procedures of VSGraph-LC are presented in an algorithm form in the appendix. Details of our approach is explained in the following subsections.

### 3.1 Problem Definition and Notations

Traditional webly supervised learning problems aim to train the CNN model $\mathcal{M}(\theta)$ for the optimal parameter $\hat{\theta}$ from dataset $\mathcal{D} = \{(x_1, y_1^*), \ldots, (x_N, y_N^*)\}$ [25]. Consider the massive label noise, web

**Figure 3: Pipeline of VSGraph-LC. The aim is to provide final labels that are more reliable than web labels for later finetuning. A metadata-based anchor selector is built firstly to provide guidance for GNN training. GNN and CNN labels are the predictions of GNN and the pretrained CNN model, respectively. The final labels take advantages of both CNN and GNN labels**

label $y_i^*$ might not reflect the correct category that $x_i$ belongs to [48]. Thus, our task is to propose a label corrector VSGraph-LC which provides final labels $p_f$ to correct the former noisy labels $y^*$ for finetuning on the pretrained CNN model $\mathcal{M}(\theta_0)$, with the aid of metadata $\bar{t}_i$ crawled together with $x_i$ and $y_i^*$. For notation, we use non-subscript style to represent the matrix that collects all vectors of the entire dataset. For example, $y^*$ represents the matrix collecting all $y_i^*$. We also denote the matrix of label names as $\bar{l}$.

## 3.2 Visual Graph Construction

With a CNN model $\mathcal{M}(\theta_0)$ that pretrained on the entire webly training set $\mathcal{D}$, for every sample $x_i$, we obtain a visual feature $s_i$ that is extracted before fully-connected layer and a CNN label $p_{c_i}$ equivalent to $p(y|x_i, \theta_0)$, the prediction of $\mathcal{M}(\theta_0)$. With visual features $s$, we construct a k-nearest neighbor [9] undirected graph $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$ where node set $\mathcal{V}$ contains every sample in the training set attached with corresponding visual and text features for later use. Information of edge set $\mathcal{E}$ is in weighted adjacency matrix $\mathbf{A}$.

$$A_{ij} = \begin{cases} \cos(s_i, s_j) & , \text{if } v_i \in \mathcal{N}_k(v_j) \text{ or } v_j \in \mathcal{N}_k(v_i) \\ 0 & , \text{otherwise,} \end{cases} \quad (1)$$

where $\mathcal{N}_k(v_i)$ denotes the set of $k$ neighbors of node $v_i$ and $\cos(s_i, s_j)$ calculates the cosine similarity between two features $s_i$ and $s_j$. Till then, the visual graph is built completely.

## 3.3 Text Embedding

In this section, we prepare embeddings for unstructured text metadata and label description. After obtaining the raw metadata $\bar{t}$ from the dataset, we remove all punctuations, digits, and stop words from the raw text, followed by the tokenization, stemming, and lemmatizing. The preprocessed metadata is encoded by an off-the-shelf

document embedding model. Denote the joint process of preprocessing and document embedding as function $E_{doc}$, the metadata embedding $t_i$ can be expressed as

$$t_i = E_{doc}(\bar{t}_i) \quad (2)$$

To encode the label concept, however, we cannot apply $E_{doc}$ directly on the label name $\bar{l}_i$ since the label name itself only contains a few words which may be non-descriptive and semantically confusing. This problem can be solved by involving a semantic knowledge base to obtain a detailed label description for a given concept. In our work, we use WordNet [33], a lexical database that arranges distinct cognitive synonyms (called 'synsets') in a tree structure, to enhance the descriptive power of the label name. Similar to [44], our label description is obtained by extracting the definitions and lemma names of the original web label synset and its adjacent synsets, including hyponyms (subclass of) and member holonyms (part of). The purpose of collecting adjacent synsets is to include potentially related concepts. For instance, the label description of class 'drumstick' is 'drumstick: a stick used for playing a drum' for the original web synset plus 'mallet, hammer: a light drumstick with a rounded head that is used to strike such percussion instruments as chimes, kettledrums, marimbas, glockenspiels, etc.' for its adjacent synsets.

For each label name $\bar{l}_i \in \{\bar{l}_1, \ldots, \bar{l}_C\}$, where $C$ is total number of categories, we locate the corresponding synset in WordNet and extract its label description, denoted as $\phi(\bar{l}_i)$. Using $E_{doc}$ still, the label description embedding $l_i$ is obtained by

$$l_i = E_{doc}(\phi(\bar{l}_i)). \quad (3)$$

## 3.4 Anchor Selection

From t-SNE visual feature visualization [29] in Figure 1, we observe that even under the same web label category, the CNN model is still able to cluster features according to concepts, but the discriminative

function has a large bias because of semantic label noise. Thus, it is necessary to set anchors for categories to pinpoint the correct concept. We use metadata to achieve this goal.

Firstly, we enhance the metadata embeddings $t$ by applying a graph smoothing function explained in [18]. The enhanced metadata embeddings $\hat{t}$ are calculated as

$$\hat{t} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + w\mathbf{I})\mathbf{D}^{-\frac{1}{2}}t, \tag{4}$$

which aggregates metadata embeddings of neighbors on the visual graph to alleviate noise from metadata. The degree matrix $\mathbf{D}$ is diagonal with its element $D_{ii} = \sum_j A_{ij}$. Self-weight ratio $w$ is a scalar to decide the proportion of origin $t$ in the enhanced $\hat{t}$. $\mathbf{I}$ is an identity matrix.

For sample $i$, $l_{y_i^*}$ denotes the label description embedding of its web label name. Cosine similarity between $\hat{t}_i$ and $l_{y_i^*}$ reflects the possibility that it belongs to its web label's correct concept. We select $m$ anchors from each class to form anchor set $\mathcal{A}$ as Equation 5, igniting the next GNN labeling process. $\tau_{y_i^*}^m$ equals to the $m$-th highest value in web category $y_i^*$, ensuring equal number of anchors are selected in each class.

$$\mathcal{A} = \left\{ (s_i, y_i^*) \,\middle|\, \cos\left(\hat{t}_i, l_{y_i^*}\right) \geq \tau_{y_i^*}^m \right\} \tag{5}$$

## 3.5 Graph Neural Networks Labeling

We conduct GNN training on the graph $\mathcal{G}$. According to occam's razor [3], a basic $L$-layer simple graph convolutional network (SGC) [46] is implemented. For layer $i \in \{1, \ldots, L\}$, the input $h^{(i-1)}$ is transformed into $h^{(i)}$ by

$$h^{(i)} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{A} + w\mathbf{I})\mathbf{D}^{-\frac{1}{2}}h^{(i-1)}\theta_g^{(i)}, \tag{6}$$

where $\theta_g^{(i)}$ is the corresponding trainable parameter. Hyperparameter $w$ follows Equation 4. Although the GNN operation is applied to the entire visual feature set, the loss is only computed on the selected anchor set $\mathcal{A}$. The first layer input $h^{(0)}$ is assigned by visual features $s$, and the final output $h^{(L)} = p(y|s, \theta_g)$, where $\theta_g$ collects all the training parameters $\theta_g^{(i)}$ across layers. The loss function is

$$\mathcal{L}_g = \sum_{(s_i, y_i^*) \in \mathcal{A}} -y_i^* \log\left(h_i^{(L)}\right). \tag{7}$$

After the optimal $\hat{\theta}_g$ is obtained from training, the labeling process starts by applying the GNN inference for the entire visual feature set. The prediction $p(y|s, \hat{\theta}_g)$ is directly utilized as GNN labels $p_g$.

## 3.6 Correct Label Estimation

In the experiments, we observe that GNN labels with high confidence are usually reliable samples with clean background. In contrast, samples with low GNN label confidence are either hard samples or open-set noise, which are difficult or impossible to be classified. To fully exploit these samples beyond the reach of GNN, CNN label $p_c$, the prediction of the pretrained CNN model, is utilized to correct low-confident GNN labels.

Formally, a method for combining GNN and CNN labels is proposed as

$$p_{f_i} = \begin{cases} p_{g_i} & , \max(p_{g_i}) \geq \tau_f \\ \lambda p_{g_i} + (1 - \lambda)p_{c_i} & , \text{otherwise,} \end{cases} \tag{8}$$

where $\lambda$ controls the contribution of CNN labels when GNN labels have lower confidence scores than threshold $\tau_f$.

## 3.7 Summary

In summary, our VSGraph-LC firstly selects anchors based on metadata and label description. Anchors propagate their correct semantic concepts across the visual feature graph. The process of VSGraph-LC finishes when the final label $p_f$ is generated. With the corrected label, we finetune the pretrained CNN model with loss function

$$\mathcal{L}_c = \sum_{(x_i, p_{f_i}) \in \mathcal{D}} -p_{f_i} \log\left(p(y|x_i, \theta)\right). \tag{9}$$

## 4 EXPERIMENTS

In this section, we evaluate our method VSGraph-LC on WebVision-1000 [25] and NUS-WIDE [6]. To accelerate the experiments, we generate a compact Google-500 dataset from massive WebVision-1000 for ablation study and discussion. Our method reaches the state-of-the-art result on WebVision-1000 and proves its robustness and generalization on noisy multi-label dataset NUS-WIDE. Besides, we investigate the progressive training strategy for VSGraph-LC in Section 4.5. We also find our method more powerful on open-set validation set in Section 4.6.

## 4.1 Datasets, Metadata and Configurations

*WebVision-1000.* WebVision-1000 [25] contains $2.4M$ web images crawled from Flickr and Google, with keywords from 1000 class-labels in ILSVRC-2012 [8] (known as ImageNet). The estimated top-1 label accuracy is 48% [12]. As WebVision shares the same 1000 classes with ImageNet, we also use ImageNet validation set along with WebVision-1000's own validation set for evaluation.

For the WebVision dataset, every training sample contains extra attributes crawled along with the image, including the rank in searching result, title, description, tags, source website, etc. Considering both cleanness and descriptiveness, we choose 'title' + 'description' attributes for Google images and 'title' + 'tags' attributes for Flickr images as the metadata.

*Google-500.* Google-500 only keeps images crawled from google websites in WebVision-1000 for their more completed and cleaner metadata than Flickr. Also, we randomly sample one-half categories to alleviate the large consumption of time and GPU resources without losing generalization with a total of 489755 samples. Validation sets of selected categories remain. We mainly use Google-500 for ablation studies. The metadata details refer to WebVision-1000.

*NUS-WIDE.* NUS-WIDE [6] is a real-world web image dataset that contains 269,648 images with the total number of 5018 associated tags crawled from Flickr. Each image contains web tags and human-annotated ground-truth labels for the 81 concepts. Since the dataset does not contain the original web queries, we obtain web labels by extracting labels from images' associated tags among

these 81 labels, i.e. we check whether each of the 81 labels appears in its web tags. It is reported in [6] that on average 50% of the web labels are incorrect and 50% of the ground-truth labels are missing in web labels. All web tags for an image are used as its metadata.

*Configuration details.* ResNet50 [15] is selected as our CNN model in all experiments. For all experiments, we set batch size as 256 and mini-batch size as 32 trained on 8 GPUs, except in WebVision-1000 we set batch size as 1024 on 32 GPUs. We use the standard SGD with the momentum of 0.9 and weight decay of $10^{-4}$. A warm-start linearly reaches the initial learning rate in the first 10 epochs. The remained epochs are ruled by a cosine learning rate scheduler. A simple class reweighting is performed to deal with class imbalance. Google-500/WebVision-1000 requires training from scratch with 120/150 epochs and an initial learning rate of 0.1/0.4. In the finetuning stage, initial learning rate is cut half from origin one without warm-start. Training on NUS-WIDE requires an ImageNet pretrained model with 100 epochs and learning rate 0.002. For the GNN model, we take the 1-layer SGC model with learning rate of 0.1, 5000 epochs and Adam optimizer with weight decay of $10^{-6}$.

## 4.2 Google-500

In this section, we experiment VSGraph-LC on the Google-500 dataset. We first evaluate the anchor selector and scrutinize the quality of selected anchors in Section 4.2.1. We then confirm that the final label calculation policy in Section 3.6 can achieve a better result than only using any one of the component labels. Finally, we evaluate the performance of VSGraph-LC under various hyperparameters, showing its robustness.

*4.2.1 Anchor Selector.* With metadata embeddings $t$ and label description embeddings $l$ well prepared, anchor selector will be built according to Section 3.4. Firstly we choose $k = 5$ for $k$-NN graph $\mathcal{G}$ and self weight $w = 0$ for text feature aggregation, which means the enhanced text feature only depends on neighbours' embeddings. The selection of hyperparameter will be explained in Section 4.2.3. Figure 4 shows anchors from 3 classes with lowest classification accuracy according to pretrained model $\mathcal{M}(\theta_0)$. Top 10 anchors selected by different methods are shown. The visualization shows that $\mathcal{M}(\theta_0)$ predicts high confidence on samples with an incorrect semantic concept regarding class 'drumstick' and 'tiger cat', which reflects that $\mathcal{M}(\theta_0)$ is misled by massive semantic label noise from these classes. Anchors selected by metadata without graph enhancement can also make mistakes. Fortunately, when the graph-enhanced text features are introduced, those mistakes made by isolated metadata are largely mitigated by insurance from samples' neighborhood. In the following experiments, we keep $m = 10$ for the Google-500 anchor selector and expect the selected anchors to provide reliable guidance for the graph neural network.

*4.2.2 GNN Training, Final Labeling and Results.* With graph $\mathcal{G}$ constructed, simple graph convolutional network (SGC) [46] is utilized for training followed by Section 3.5. Table 1 shows the deficient performance if we only use the GNN label as the final label for finetune. As explained in Section 3.6, with GNN, anchors will only propagate their labels to nearby clean and easy samples on the graph, whereas for those hard samples or open-set noise, GNN prefers to give them very low prediction scores, acting like

**Table 1: Ablation study on Google-500 dataset**

| Method | WebVision | | ImageNet | |
|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 |
| Pretrained model | 66.96 | 82.68 | 61.54 | 78.89 |
| Co-teaching | 67.61 | 84.04 | 62.18 | 80.98 |
| Finetune by $p_g$ only | 64.79 | 81.22 | 60.39 | 78.76 |
| Finetune by $p_c$ only | 67.83 | 83.93 | 62.68 | 80.68 |
| Finetune by $p_f$ | **68.14** | **84.46** | **63.16** | **81.45** |

dropping numerous data that harm the data-driven DNNs. Thus, we set $\lambda = 0.5$ for Equation 8, therefore the model finetuned by the final label will have a large improvement than the pretrained model $\mathcal{M}(\theta_0)$ and the model finetuned by only CNN or GNN labels. The results also have advantages over Co-teaching.

*4.2.3 Hyperparameters.* In this section, we try several values for critical hyperparameters of $k$, $w$ for k-NN graph, and $\lambda$, $\tau_f$ for final label calculation in Equation 8. Mean values of WebVision/ImageNet top-1/top-5 are used as average accuracy to indicate the overall performance. Figure 5a shows the effectiveness of using k-NN, where $k = 5$ ensures an optimal result. Figure 5b shows $\lambda = 0.5$ has constant merit comparing to $\lambda = 0$. $\tau_f = 0.7$ can assist to achieve optimal average accuracy. However, although the advantages of selected hyperparameters, the accuracy difference does not vary much for $\tau_f \in [0.5, 0.9]$, $k \in [3, 8]$ and any $w$, showing the robustness of proposed VSGraph-LC to hyperparameters.

## 4.3 WebVision-1000

Table 2 reports experimental results on WebVision-1000. VSGraph-LC gains a large improvement on WebVision top-1 by more than 1.2%. The advantage is even larger on ImageNet, especially on ImageNet top-5, showing a good generalization ability. In comparison, finetuning by CNN labels only obtains weaker improvements.

Notice that our method also outperforms the state-of-the-art methods on WebVision and ImageNet validation sets. Except CurriculumNet, MentorNet and CleanNet both adopt extra human-verified datasets to train a guidance network first, while MentorNet chooses a backbone of InceptionResNetV2 which is stronger than our ResNet50. With these disadvantages, however, our method can still obtain a better performance compared to the above methods. In addition, Multimodal uses ImageNet data for training visual embedding and a query-image pairs dataset for training phrased generation, with stronger InceptionV3 being the backbone. Our VSGraph-LC can still exceed Multimodel in WebVision top-1 accuracy.

## 4.4 NUS-WIDE

NUS-WIDE [6] is another real-world web image dataset. Different from WebVision, NUS-WIDE is designed for multi-label image classification. Previous works [43, 57] only use the dataset with ground-truth labels for standard multi-label learning, while we are interested in the real-world label noise setting in multi-label learning. In this case, we train our model with web labels and evaluate

**Anchors by Model Confidence** | **Anchors by Metadata w/o Graph Enhancement** | **Anchors by Metadata w/ Graph Enhancement (k=5)**

**(a) Selected Anchors for Class 'Drumstick'**

**(b) Selected Anchors for Class 'Spotlight'**

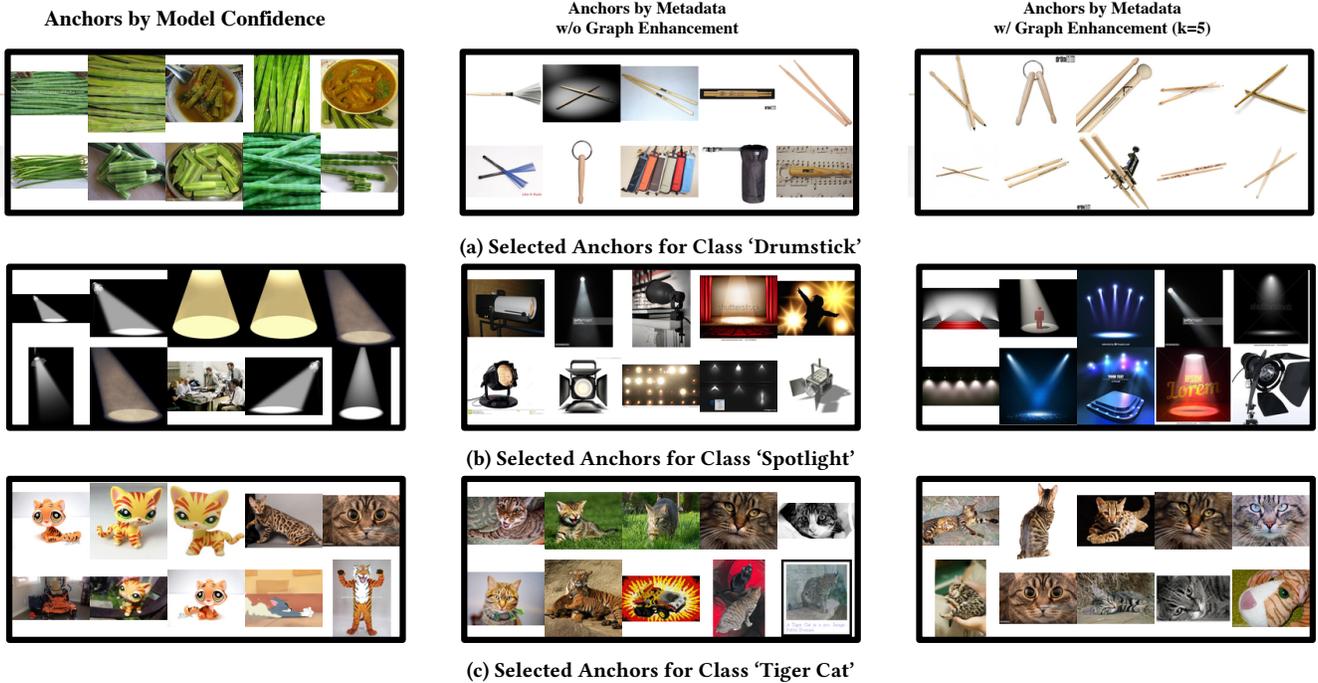**(c) Selected Anchors for Class 'Tiger Cat'**

**Figure 4: Exemplar anchors for three most noisy classes in Google-500. Different columns indicate anchors selected by different methods. Using metadata with graph enhancement has perceptible advantages compared to other methods**
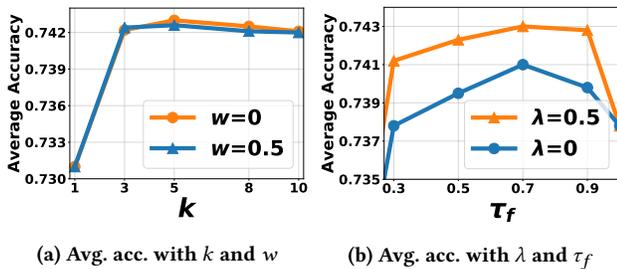


**(a) Avg. acc. with $k$ and $w$**     **(b) Avg. acc. with $\lambda$ and $\tau_f$**

**Figure 5: Comparison between different hyperparameters**

**Table 2: The state-of-the-art results on WebVision-1000**

| Method | Backbone | WebVision | | ImageNet | |
|---|---|---|---|---|---|
| | | Top-1 | Top-5 | Top-1 | Top-5 |
| MentorNet [17] | InceptionResNetV2 | 72.60 | 88.90 | 64.20 | 84.80 |
| CleanNet [24] | ResNet50 | 70.31 | 87.77 | 63.42 | 84.59 |
| CurriculumNet [12] | InceptionV2 | 72.10 | 89.20 | 64.80 | 84.90 |
| Multimodal [36] | InceptionV3 | 73.15 | 89.73 | - | - |
| Pretrained model | ResNet50 | 74.25 | 89.84 | 68.28 | 86.23 |
| Finetune by $p_c$ only | ResNet50 | 75.15 | 89.93 | 69.07 | 86.76 |
| Finetune by $p_f$ | ResNet50 | **75.48** | **90.15** | **69.42** | **87.29** |

it on ground-truth labels. To distinguish from previous multi-label classifier learned without noise, we denote their and our setting

**Table 3: Results on NUS-81-Web with noisy web labels for training. $K = 3$ is used for calculating C-F1 and O-F1**

| Method | C-F1 | O-F1 | mAP |
|---|---|---|---|
| Pretrained model | 37.51 | 39.59 | 43.94 |
| Finetune by $p_c$ only | 37.62 | 39.15 | 43.99 |
| Finetune by $p_f$ | **38.58** | **40.16** | **44.83** |

as NUS-81 and NUS-81-Web, respectively. The only difference between NUS-81 and NUS-81-Web is that the NUS-81-Web training set uses web labels rather than annotated ground-truth labels. The train-test split policy adopts the official one. For experiments, we remove samples without any label within the 81 label set. Moreover, the label descriptions are obtained by identifying the most relevant synset on WordNet for each label.

*4.4.1 Evaluation Metrics.* For multi-label classification, we compare the overall F1-measure (O-F1), per-class (also known as macro-averaged) F1-measure (C-F1), and mean average precision (mAP) to evaluate the performance. Following [43], we use top $K = 3$ highest confidence labels for each image as the prediction and compare with the ground-truth labels.

*4.4.2 Anchor Selector and GNN Training.* With the high noisy ratio and small dataset size for NUS-81-Web, the model trained directly from scratch using web labels cannot provide strong visual features

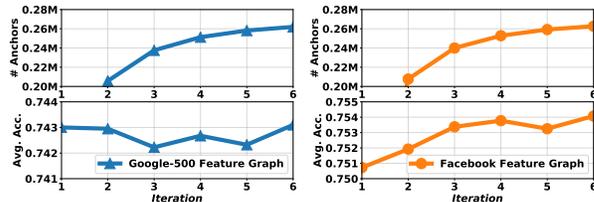**Table 4: Results on Google-500 open-set problem**

| Method | WebVision | | ImageNet | |
|---|---|---|---|---|
| | C-P/C-R | C-F1 | C-P/C-R | C-F1 |
| Pretrained model | 40.44/67.23 | 50.50 | 37.75/60.36 | 46.45 |
| Co-teaching | 40.38/69.21 | 51.00 | 38.27/63.58 | 47.78 |
| Finetune by $p_c$ | 40.34/68.21 | 50.70 | 38.54/62.10 | 47.56 |
| Finetune by $p_g$ | 59.64/55.63 | **57.56** | 56.74/48.89 | 52.52 |
| Finetune by $p_f$ | 52.47/62.95 | 57.24 | 49.66/56.36 | **52.80** |



(a) Google-500 Feature Graph  (b) Off-the-Shelf Feature Graph

**Figure 6: Exploration on progressive training, which only works with graph built by high-quality visual features**

for building $k$-NN graph. Therefore, we finetune the ImageNet-pretrained model for our pretrained model $\mathcal{M}(\theta_0)$ instead. Regarding hyperparameters, we set $k = 10$ for $k$-NN graph, self-weight $w = 0$ for text enhancement and $m = 50$ for anchor selector.

*4.4.3 Results.* Experimental results on NUS-81-Web are shown in Table 3. Our proposed method outperforms baseline model and the model finetuned by only CNN labels for all three metrics. For the mAP score, we achieve 0.9% improvement using VSGraph-LC compared to the pretrained model. Our method can also increase both C-F1 and O-F1 by 1.1% and 0.6%, respectively. Unlike the performance on WebVision, finetuning with CNN labels brings no significant improvement on C-F1 and mAP, and even a 0.4% drop on O-F1. This demonstrates that under a high noisy-level setting like NUS-81-Web, labels corrected by our VSGraph-LC method are more reliable than the CNN labels.

## 4.5 Discussion on Progressive Training

In this section, we discuss the performance of progressive GNN training. The procedure is as follows: After the standard GNN training is completed as Section 3.5, samples with confidence over threshold $\tau_f$ will be selected as anchors, labeled by GNN labels, for GNN training in the next round. Progressive GNN training utilizes the static graph that built by the pretrained CNN features. We suppose that with increasing iterations, the performance of VSGraph-LC will keep growing, boosted by more anchors. However, Figure 6a shows that more iterations cannot guarantee a better performance under Google-500 settings in Section 4.2, even though the number of anchors increases according to top Figure 6a. We assume that the failure attributes to the quality of graph structure $\mathcal{G}$ built by purely Google-500 base model features. Unable to build the connection between easy and hard samples of the identical category, the existing graph $\mathcal{G}$ only constrains GNN labels within easy samples, thus disables progressive training. To prove this, we build the visual graph using strong features extracted by an off-the-shelf ResNeXt-101 model provided by [49], which is trained on more than 940 million images in semi-weakly supervised learning fashion. We evaluate progressive training again by ONLY changing the graph structure, keeping the original 5000 anchors and all node features and scores the same. In this case, initial anchors can reach hard and more informative samples, accumulated as new anchors for progressive training. During the iteration, $\mathcal{G}$ remains unchanged as before. The results show a stable improvement iteration by iteration, indicating the importance of visual graph structure for progressive training of VSGraph-LC.

## 4.6 Discussion on Open-Set Recognition (OSR)

Due to the complexity of real-world scenarios, DNN is ideally resistant to open-set images, *i.e.*, for test images which do not belong to any category in the validation set, CNN models should produce low confidence for them. To demonstrate our model's superiority on open-set recognition task, we use the final model trained from the Google-500 training set and evaluate it on the entire WebVision/ImageNet validation dataset. To be specific, since the WebVision/ImageNet validation sets share the same 1000 classes, we set these 1000 classes except Google-500's classes as open-set classes (i.e. 500 for classification and other 500 as open set). OSR expects unconfident predictions for open-set images. If CNN models predict an image with confidence beyond a threshold (0.2 in our paper), we consider it is classified into the predicted category, followed by [35, 53]. Thus, those open-set images might be falsely classified into Google-500's classes. We take per-class precision (C-R), recall (C-R), and F1-measure (C-F1) as metrics in Table 4, where VSGraph-LC has tremendous improvements than the pretrained model and Co-teaching, showing our model fits the real-world open-set tasks well.

We also find that finetuning by $p_c$ does not work in OSR as it tends to give every sample high confidence even for open-set images. However, $p_g$ from GNN only gives high confidence to samples near anchors, thus can reject open-set samples and outperforms $p_c$ on OSR. But GNN might miss some hard positive samples, harming performance in the close-set setting. We therefore introduce $p_f$ to combine $p_c$ and $p_g$, whose OSR ability, in some cases, might be weaker than only using $p_g$ alone because of $p_c$.

## 5 CONCLUSIONS

In this paper, we focus on webly supervised learning task and highlight two understudied but critical factors: semantic label noise and text metadata. Based on our extensive exploration on them, we gain insights that CNN model that pretrained from entire webly dataset is able to provide a visual feature space where similar semantic images cluster themselves. With efficient usage of metadata, an effective and automatic label corrector VSGraph-LC is proposed.

# REFERENCES

[1] Görkem Algan and Ilkay Ulusoy. 2019. Image Classification with Deep Learning in the Presence of Noisy Labels: A Survey. *arXiv preprint arXiv:1912.05170* (2019).

[2] Tamara L Berg and David A Forsyth. 2006. Animals on the web. In *CVPR*, Vol. 2. IEEE, 1463–1470.

[3] Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. 1987. Occam's razor. *Inform. Process. Lett.* 24, 6 (1987), 377–380.

[4] Xinlei Chen, Abhinav Shrivastava, and Abhinav Gupta. 2013. Neil: Extracting visual knowledge from web data. In *ICCV*. 1409–1416.

[5] Zhao-Min Chen, Xiu-Shen Wei, Peng Wang, and Yanwen Guo. 2019. Multi-label image recognition with graph convolutional networks. In *CVPR*. 5177–5186.

[6] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. 2009. NUS-WIDE: a real-world web image database from National University of Singapore. In *ACM CIVR*. 1–9.

[7] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *NIPS*. 3844–3852.

[8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*. IEEE, 248–255.

[9] Sahibsingh A Dudani. 1976. The distance-weighted k-nearest-neighbor rule. *TSMC* 4 (1976), 325–327.

[10] Daniela Espinoza-Molina and Mihai Datcu. 2013. Earth-observation image retrieval based on content, semantics, and metadata. *TGRS* 51, 11 (2013), 5145–5159.

[11] Aritra Ghosh, Himanshu Kumar, and PS Sastry. 2017. Robust loss functions under label noise for deep neural networks. In *AAAI*.

[12] Sheng Guo, Weilin Huang, Haozhi Zhang, Chenfan Zhuang, Dengke Dong, Matthew R Scott, and Dinglong Huang. 2018. Curriculumnet: Weakly supervised learning from large-scale web images. In *ECCV*. 135–150.

[13] Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi Sugiyama. 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NIPS*. 8527–8537.

[14] Jiangfan Han, Ping Luo, and Xiaogang Wang. 2019. Deep self-learning from noisy labels. In *ICCV*. 5138–5147.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–778.

[16] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. 2019. Label propagation for deep semi-supervised learning. In *CVPR*. 5070–5079.

[17] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. 2018. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*. 2304–2313.

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

[19] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.

[20] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. 2017. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from https://github. com/openimages* 2, 3 (2017), 18.

[21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV* 123, 1 (2017), 32–73.

[22] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Tom Duerig, et al. 2018. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982* (2018).

[23] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *nature* 521, 7553 (2015), 436–444.

[24] Kuang-Huei Lee, Xiaodong He, Lei Zhang, and Linjun Yang. 2018. Cleannet: Transfer learning for scalable image classifier training with label noise. In *CVPR*. 5447–5456.

[25] Wen Li, Limin Wang, Wei Li, Eirikur Agustsson, and Luc Van Gool. 2017. Webvision database: Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862* (2017).

[26] Zechao Li and Jinhui Tang. 2015. Weakly supervised deep metric learning for community-contributed image retrieval. *IEEE Transactions on Multimedia* 17, 11 (2015), 1989–1999.

[27] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).

[28] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*. 13–23.

[29] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* 9, Nov (2008), 2579–2605.

[30] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. 2018. Exploring the limits of weakly supervised pretraining. In *ECCV*. 181–196.

[31] Naresh Manwani and PS Sastry. 2013. Noise tolerance under risk minimization. *IEEE Transactions on Cybernetics* 43, 3 (2013), 1146–1151.

[32] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*. 3111–3119.

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

[34] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *TPAMI* 41, 8 (2018), 1979–1993.

[35] Pramuditha Perera, Ramesh Nallapati, and Bing Xiang. 2019. Ocgan: One-class novelty detection using gans with constrained latent representations. In *CVPR*. 2898–2906.

[36] Manan Shah, Krishnamurthy Viswanathan, Chun-Ta Lu, Ariel Fuxman, Zhen Li, Aleksei Timofeev, Chao Jia, and Chen Sun. 2019. Inferring Context from Pixels for Multimodal Image Classification. In *ACM CIKM*. ACM, 189–198.

[37] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*. 2556–2565.

[38] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2020. VL-BERT: Pre-training of Generic Visual-Linguistic Representations. In *ICLR*.

[39] Chen Sun, Chuang Gan, and Ram Nevatia. 2015. Automatic concept discovery from parallel text and visual corpora. In *ICCV*. 2596–2604.

[40] Daiki Tanaka, Daiki Ikami, Toshihiko Yamasaki, and Kiyoharu Aizawa. 2018. Joint optimization framework for learning with noisy labels. In *CVPR*. 5552–5560.

[41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*. 5998–6008.

[42] M Mitchell Waldrop. 2019. News Feature: What are the limits of deep learning? *PNAS* 116, 4 (2019), 1074–1077.

[43] Jiang Wang, Yi Yang, Junhua Mao, Zhiheng Huang, Chang Huang, and Wei Xu. 2016. Cnn-rnn: A unified framework for multi-label image classification. In *CVPR*. 2285–2294.

[44] Tingting Wei, Yonghe Lu, Huiyou Chang, Qiang Zhou, and Xianyu Bao. 2015. A semantic approach for text clustering using WordNet and lexical chains. *Expert Systems with Applications* 42, 4 (2015), 2264–2275.

[45] Baoyuan Wu, Weidong Chen, Yanbo Fan, Yong Zhang, Jinlong Hou, Jie Liu, and Tong Zhang. 2019. Tencent ml-images: A large-scale multi-label image database for visual representation learning. *IEEE Access* 7 (2019), 172683–172693.

[46] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying Graph Convolutional Networks. In *ICML*. PMLR, 6861–6871.

[47] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. 2020. A comprehensive survey on graph neural networks. *TNNLS* (2020).

[48] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. 2015. Learning from massive noisy labeled data for image classification. In *CVPR*. 2691–2699.

[49] I Zeki Yalniz, Hervé Jégou, Kan Chen, Manohar Paluri, and Dhruv Mahajan. 2019. Billion-scale semi-supervised learning for image classification. *arXiv preprint arXiv:1905.00546* (2019).

[50] J. Yang, L. Feng, W. Chen, X. Yan, H. Zheng, P. Luo, and W. Zhang. 2020. Webly Supervised Image Classification with Self-Contained Confidence. In *ECCV*.

[51] Lei Yang, Xiaohang Zhan, Dapeng Chen, Junjie Yan, Chen Change Loy, and Dahua Lin. 2019. Learning to cluster faces on an affinity graph. In *CVPR*. 2298–2306.

[52] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*. 5753–5763.

[53] Ryota Yoshihashi, Wen Shao, Rei Kawakami, Shaodi You, Makoto Iida, and Takeshi Naemura. 2019. Classification-reconstruction learning for open-set recognition. In *CVPR*. 4016–4025.

[54] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. Mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412* (2017).

[55] Bolei Zhou, Vignesh Jagadeesh, and Robinson Piramuthu. 2015. Conceptlearner: Discovering visual concepts from weakly labeled image collections. In *CVPR*. 1492–1500.

[56] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2018. Graph neural networks: A review of methods and applications. *arXiv preprint arXiv:1812.08434* (2018).

[57] Feng Zhu, Hongsheng Li, Wanli Ouyang, Nenghai Yu, and Xiaogang Wang. 2017. Learning spatial regularization with image-level supervisions for multi-label image classification. In *CVPR*. 5513–5522.

# A ALGORITHM

---

**Algorithm 1** VSGraph-LC (depicted as Figure 3)

---

**Input**: Dataset $\mathcal{D} = \{x, y^*\}$, raw metadata $\bar{t}$ for every sample, raw label names $\bar{l}$, number of anchors per class $m$.

**Output**: Final corrected labels $p_f$ for finetuning $\mathcal{M}(\theta_0)$.

▷ *Prepare Word Embeddings*
  1: Obtain word embeddings $t$ from raw metadata $\bar{t}$ by Eq.2.
  2: Obtain word embeddings $l$ from label names $\bar{l}$ by Eq.3.
▷ *Obtain CNN Labels*
  3: Obtain pretrained CNN model $\mathcal{M}(\theta_0)$ from the entire $\{x, y^*\}$.
  4: Obtain visual features $s$ and CNN labels $p_c$ from $\mathcal{M}(\theta_0)$.
▷ *Obtain GNN Labels*
  5: Build graph $\mathcal{G}$ using $s$, with adjacency matrix $\mathbf{A}$ by Eq.1.
  6: Obtain graph-enhanced text features $\hat{t}$ by Eq.4.
  7: Generate anchor selector by Eq.5 with $m$, obtain anchor set $\mathcal{A}$.
  8: Train GNN on $\mathcal{G}$ supervised by $\mathcal{A}$ for optimal $\hat{\theta}_g$.
  9: Serve GNN prediction $p(y|s, \hat{\theta}_g)$ as GNN labels $p_g$.
▷ *Obtain Final Labels for Finetuning*
 10: Obtain final pseudo labels $p_f$ by Eq.8.

---